

new/usr/src/cmd/mdb/Makefile.kmdb.files

1

```
*****
2524 Fri Apr 6 17:24:52 2018
new/usr/src/cmd/mdb/Makefile.kmdb.files
9210 remove KMDB branch debugging support
9211 ::crregs could do with cr2/cr3 support
9209 ::ttrace should be able to filter by thread
Reviewed by: Patrick Mooney <patrick.mooney@joyent.com>
Reviewed by: Yuri Pankov <yuripv@yuripv.net>
*****
1 #
2 # CDDL HEADER START
3 #
4 # The contents of this file are subject to the terms of the
5 # Common Development and Distribution License (the "License").
6 # You may not use this file except in compliance with the License.
7 #
8 # You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9 # or http://www.opensolaris.org/os/licensing.
10 # See the License for the specific language governing permissions
11 # and limitations under the License.
12 #
13 # When distributing Covered Code, include this CDDL HEADER in each
14 # file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 # If applicable, add the following below this CDDL HEADER, with the
16 # fields enclosed by brackets "[]" replaced with your own identifying
17 # information: Portions Copyright [yyyy] [name of copyright owner]
18 #
19 # CDDL HEADER END
20 #
21 #
22 # Copyright 2009 Sun Microsystems, Inc. All rights reserved.
23 # Use is subject to license terms.
24 #
25 #
26 #
27 # Copyright (c) 2012 by Delphix. All rights reserved.
28 # Copyright (c) 2018 Joyent, Inc. All rights reserved.
29 # Copyright (c) 2012 Joyent, Inc. All rights reserved.
30 #
31 KMDBSRCS += \
32     ffs.c \
33     kaif_start.c \
34     mdb.c \
35     mdb_addrvec.c \
36     mdb_argvec.c \
37     mdb_callb.c \
38     mdb_cmdbuf.c \
39     mdb_cmds.c \
40     kvm_cpu.c \
41     kmdb_conf.c \
42     kmdb_context.c \
43     kmdb_create.c \
44     mdb_ctf.c \
45     kmdb_ctf_open.c \
46     mdb_debug.c \
47     kmdb_demangle.c \
48     mdb_disasm.c \
49     kmdb_dl.c \
50     kmdb_dpi.c \
51     mdb_dump.c \
52     mdb_err.c \
53     mdb_evset.c \
54     kmdb_fault.c \
55     kmdb_fdio.c \
56     mdb_fmt.c \
```

new/usr/src/cmd/mdb/Makefile.kmdb.files

2

```
56     mdb_frame.c \
57     mdb_gelf.c \
58     mdb_help.c \
59     mdb_io.c \
60     kmdb_kdi.c \
61     kmdb_kvm.c \
62     mdb_logio.c \
63     mdb_list.c \
64     mdb_macalias.c \
65     kmdb_main.c \
66     mdb_modapi.c \
67     mdb_module.c \
68     kmdb_module.c \
69     kmdb_module_load.c \
70     mdb_nm.c \
71     mdb_nv.c \
72     mdb_pipeio.c \
73     mdb_print.c \
74     kmdb_promio.c \
75     kmdb_promif.c \
76     mdb_set.c \
77     kmdb_shell.c \
78     mdb_signal.c \
79     mdb_string.c \
80     mdb_strio.c \
81     kmdb_stubs.c \
82     mdb_tab.c \
83     mdb_target.c \
84     kmdb_terminfo.c \
85     mdb_termio.c \
86     mdb_typedef.c \
87     mdb_umem.c \
88     kmdb_umemglue.c \
89     mdb_value.c \
90     mdb_vcb.c \
91     mdb_wcb.c \
92     mdb_whatisc.c \
93     kmdb_wr.c
94
95 KMDBML +=
96
97 KMDBOBS = $(KMDBSRCS:%.c=%.o) $(KMDBML:%.s=%.o)
98
99 PROMSRCS +=
100
101 PROMOBS = $(PROMSRCS:%.c=%.o)
102
103 KCTLSRCS += \
104     kctl_auxv.c \
105     kctl_dmod.c \
106     kctl_err.c \
107     kctl_main.c \
108     kctl_mod.c \
109     kctl_string.c \
110     kctl_wr.c
111
112 KCTLML +=
113
114 KCTLOBS = $(KCTLSRCS:%.c=%.o) $(KCTLML:%.s=%.o)
115
116 SRCS += $(KMDBSRCS) $(PROMSRCS)
117 MLSRCS += $(KMDBML)
118 OBS = $(SRCS:%.c=%.o) $(KMDBML:%.s=%.o)
119
120 ALLOBS = $(OBS) $(KCTLOBS)
121 ALLLINTFILES = $(ALLOBS:%.o=%.ln)
```

new/usr/src/cmd/mdb/Makefile.kmdb.files

3

```
123 # files that need KMDB_VERSION defined
124 VERSFILES = \
125     kmdb_conf.c \
126     kctl_main.c
127 VERSOBS = ${VERSFILES:%.c=%.o}
```

new/usr/src/cmd/mdb/common/kmdb/kmdb_dpi_impl.h

1

```
*****
3496 Fri Apr 6 17:24:52 2018
new/usr/src/cmd/mdb/common/kmdb/kmdb_dpi_impl.h
9210 remove KMDB branch debugging support
9211 ::crregs could do with cr2/cr3 support
9209 ::ttrace should be able to filter by thread
Reviewed by: Patrick Mooney <patrick.mooney@joyent.com>
Reviewed by: Yuri Pankov <yuripv@yuripv.net>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22  * Copyright 2007 Sun Microsystems, Inc. All rights reserved.
23  * Use is subject to license terms.
24  *
25  * Copyright 2018 Joyent, Inc.
26  */

28 #ifndef _KMDB_DPI_IMPL_H
29 #define _KMDB_DPI_IMPL_H

29 #pragma ident "%Z%M% %I% %E% SMI"

31 #include <setjmp.h>
32 #ifdef __sparc
33 #include <sys/regset.h>
34 #endif /* __sparc */
35 #include <sys/types.h>

37 #include <kmdb/kmdb_auxv.h>
38 #include <kmdb/kmdb_dpi.h>
39 #include <mdb/mdb_kreg.h>
40 #include <mdb/mdb_target.h>

42 #ifdef __cplusplus
43 extern "C" {
44 #endif

46 extern jmp_buf *kmdb_dpi_fault_pcb;

48 /*
49  * The routines used by the kmdb side of the DPI to access the saved state
50  * of the current kernel instance, and to control that instance. A populated
51  * version of this vector is provided by the DPI backend used to control the
52  * machine. General use of the kmdb DPI is not via direct invocation of the
53  * functions in this ops vector, but rather flows through the convenience
54  * wrappers in kmdb_dpi.c.
55  */
```

new/usr/src/cmd/mdb/common/kmdb/kmdb_dpi_impl.h

2

```
56 struct dpi_ops {
57     int (*dpo_init)(kmdb_auxv_t *);

59     void (*dpo_debugger_activate)(kdi_debugvec_t **, uint_t);
60     void (*dpo_debugger_deactivate)(void);

62     void (*dpo_enter_mon)(void);

64     void (*dpo_modchg_register)(void (*)(struct modctl *, int));
65     void (*dpo_modchg_cancel)(void);

67     int (*dpo_get_cpu_state)(int);
68     int (*dpo_get_master_cpuid)(void);

70     const mdb_tgt_gregset_t *(*dpo_get_gregs)(int);
71     int (*dpo_get_register)(const char *, kreg_t *);
72     int (*dpo_set_register)(const char *, kreg_t);
73 #ifdef __sparc
74     int (*dpo_get_rwin)(int, int, struct rwindow *);
75     int (*dpo_get_nwin)(int);
76 #endif

78     int (*dpo_brkpt_arm)(uintptr_t, mdb_instr_t *);
79     int (*dpo_brkpt_disarm)(uintptr_t, mdb_instr_t);

81     int (*dpo_wapt_validate)(kmdb_wapt_t *);
82     int (*dpo_wapt_reserve)(kmdb_wapt_t *);
83     void (*dpo_wapt_release)(kmdb_wapt_t *);
84     void (*dpo_wapt_arm)(kmdb_wapt_t *);
85     void (*dpo_wapt_disarm)(kmdb_wapt_t *);
86     int (*dpo_wapt_match)(kmdb_wapt_t *);

88     int (*dpo_step)(void);
89 #if defined(__i386) || defined(__amd64)
90     void (*dpo_step_branch)(void);
91 #endif

90     uintptr_t (*dpo_call)(uintptr_t, uint_t, const uintptr_t *);

92     void (*dpo_dump_crumbs)(uintptr_t, int);

97 #if defined(__i386) || defined(__amd64)
98     void (*dpo_msr_add)(const kdi_msr_t *);
99     uint64_t (*dpo_msr_get)(int, uint_t);
100 #endif

94 #ifdef __sparc
95     void (*dpo_kernpanic)(int);
96 #endif
97 };

_____unchanged_portion_omitted_____
```

```

*****
63016 Fri Apr 6 17:24:53 2018
new/usr/src/cmd/mdb/common/kmdb/kmdb_kvm.c
9210 remove KMDB branch debugging support
9211 ::crregs could do with cr2/cr3 support
9209 ::ttrace should be able to filter by thread
Reviewed by: Patrick Mooney <patrick.mooney@joyent.com>
Reviewed by: Yuri Pankov <yuripv@yuripv.net>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22  * Copyright (c) 2004, 2010, Oracle and/or its affiliates. All rights reserved.
23  * Copyright (c) 2013 by Delphix. All rights reserved.
24  *
25  * Copyright 2018 Joyent, Inc.
26  */

28 #include <kmdb/kmdb_kvm.h>
29 #include <kmdb/kvm.h>
30 #include <kmdb/kmdb_kdi.h>
31 #include <kmdb/kmdb_promif.h>
32 #include <kmdb/kmdb_module.h>
33 #include <kmdb/kmdb_asmutil.h>
34 #include <mdb/mdb_types.h>
35 #include <mdb/mdb_conf.h>
36 #include <mdb/mdb_err.h>
37 #include <mdb/mdb_modapi.h>
38 #include <mdb/mdb_target_impl.h>
39 #include <mdb/mdb_debug.h>
40 #include <mdb/mdb_string.h>
41 #include <mdb/mdb_ctf.h>
42 #include <mdb/mdb_kreg_impl.h>
43 #include <mdb/mdb_ks.h>
44 #include <mdb/mdb.h>

46 #include <strings.h>
47 #include <dlfcn.h>
48 #include <sys/isa_defs.h>
49 #include <sys/kobj.h>
50 #include <sys/kobj_impl.h>
51 #include <sys/bitmap.h>
52 #include <vm/as.h>

54 static const char KMT_RTLD_NAME[] = "krtld";
55 static const char KMT_MODULE[] = "mdb_ks";
56 static const char KMT_CTFPARENT[] = "genunix";

```

```

58 static mdb_list_t kmt_defbpb_list; /* List of current deferred bp's */
59 static int kmt_defbpb_lock; /* For list, running kernel holds */
60 static uint_t kmt_defbpb_modchg_isload; /* Whether mod change is load/unload */
61 static struct modctl *kmt_defbpb_modchg_modctl; /* modctl for defbpb checking */
62 static uint_t kmt_defbpb_num; /* Number of referenced def'd bp's */
63 static int kmt_defbpb_bpspec; /* vespec for def'd bp activation bp */

65 static const mdb_se_ops_t kmt_brkpt_ops;
66 static const mdb_se_ops_t kmt_wapt_ops;

68 static void kmt_sync(mdb_tgt_t *);

70 typedef struct kmt_symarg {
71     mdb_tgt_sym_f *sym_cb; /* Caller's callback function */
72     void *sym_data; /* Callback function argument */
73     uint_t sym_type; /* Symbol type/binding filter */
74     mdb_syminfo_t sym_info; /* Symbol id and table id */
75     const char *sym_obj; /* Containing object */
76 } kmt_symarg_t;
    unchanged portion omitted

550 /*ARGSUSED*/
551 static int
552 kmt_status_dcmd(uintptr_t addr, uint_t flags, int argc, const mdb_arg_t *argv)
553 {
554     kmt_data_t *kmt = mdb.m_target->t_data;
555     struct utsname uts;
556     char uuid[37];
557     kreg_t tt;

558     if (mdb_tgt_readsym(mdb.m_target, MDB_TGT_AS_VIRT, &uts, sizeof (uts),
559         "unix", "utsname") != sizeof (uts)) {
560         warn("failed to read 'utsname' struct from kernel\n");
561         bzero(&uts, sizeof (uts));
562         (void) strcpy(uts.nodename, "unknown machine");
563     }

565     mdb_printf("debugging live kernel (%d-bit) on %s\n",
566         (int)(sizeof (void *) * NBBY),
567         (uts.nodename == '\0' ? "(not set)" : uts.nodename));
568     mdb_printf("operating system: %s %s (%s)\n",
569         uts.release, uts.version, uts.machine);

571     if (mdb_tgt_readsym(mdb.m_target, MDB_TGT_AS_VIRT, uuid, sizeof (uuid),
572         "genunix", "dump_osimage_uuid") != sizeof (uuid)) {
573         warn("failed to read 'dump_osimage_uuid' string from kernel\n");
574         (void) strcpy(uuid, "(error)");
575     } else if (*uuid == '\0') {
576         (void) strcpy(uuid, "(not set)");
577     } else if (uuid[36] != '\0') {
578         (void) strcpy(uuid, "(invalid)");
579     }
580     mdb_printf("image uuid: %s\n", uuid);

581     if (kmt->kmt_cpu != NULL) {
582         mdb_printf("CPU-specific support: %s\n",
583             kmt_cpu_name(kmt->kmt_cpu));
584     }

582     mdb_printf("DTrace state: %s\n", (kmdb_kdi_dtrace_get_state() ==
583         KDI_DTSTATE_DTRACE_ACTIVE ? "active (debugger breakpoints cannot
584         "be armed)" : "inactive"));

586     (void) kmdb_dpi_get_register("tt", &tt);
587     mdb_printf("stopped on: %s\n", kmt_trapname(tt));

```

```

589     (void) kmt_dmod_status("pending dmod loads:", KMDB_MC_STATE_LOADING);
590     (void) kmt_dmod_status("pending dmod unloads:",
591         KMDB_MC_STATE_UNLOADING);

593     return (DCMD_OK);
594 }
    unchanged portion omitted

2373 static void
2374 kmt_destroy(mdb_tgt_t *t)
2375 {
2376     kmt_data_t *kmt = t->t_data;
2377     kmt_module_t *km, *pkm;

2379     mdb_nv_destroy(&kmt->kmt_modules);
2380     for (km = mdb_list_prev(&kmt->kmt_modlist); km != NULL; km = pkm) {
2381         pkm = mdb_list_prev(km);
2382         mdb_free(km, sizeof (kmt_module_t));
2383     }

2385     if (!kmt_defbp_lock)
2386         kmt_defbp_destroy_all();

2388     if (kmt->kmt_trapmap != NULL)
2389         mdb_free(kmt->kmt_trapmap, BT_SIZEOFMAP(kmt->kmt_trapmax));

2395     if (kmt->kmt_cpu != NULL)
2396         kmt_cpu_destroy(kmt->kmt_cpu);

2391     if (kmt != NULL)
2392         mdb_free(kmt, sizeof (kmt_data_t));
2393 }

2395 static const mdb_tgt_ops_t kmt_ops = {
2396     kmt_setflags,          /* t_setflags */
2397     (int (*)()) mdb_tgt_notsup, /* t_setcontext */
2398     kmt_activate,         /* t_activate */
2399     (void (*)()) mdb_tgt_nop, /* t_deactivate */
2400     kmt_periodic,        /* t_periodic */
2401     kmt_destroy,         /* t_destroy */
2402     kmt_name,            /* t_name */
2403     (const char (*)()) mdb_conf_isa, /* t_isa */
2404     kmt_platform,       /* t_platform */
2405     kmt_uname,          /* t_uname */
2406     kmt_dmodel,         /* t_dmodel */
2407     (ssize_t (*)()) mdb_tgt_notsup, /* t_aread */
2408     (ssize_t (*)()) mdb_tgt_notsup, /* t_awrite */
2409     kmt_read,           /* t_vread */
2410     kmt_write,          /* t_vwrite */
2411     kmt_pread,          /* t_pread */
2412     kmt_pwrite,         /* t_pwrite */
2413     kmt_read,           /* t_fread */
2414     kmt_write,          /* t_fwrite */
2415     kmt_ioread,         /* t_ioread */
2416     kmt_iowrite,        /* t_iowrite */
2417     kmt_vtop,           /* t_vtop */
2418     kmt_lookup_by_name, /* t_lookup_by_name */
2419     kmt_lookup_by_addr, /* t_lookup_by_addr */
2420     kmt_symbol_iter,    /* t_symbol_iter */
2421     kmt_mapping_iter,   /* t_mapping_iter */
2422     kmt_object_iter,    /* t_object_iter */
2423     kmt_addr_to_map,    /* t_addr_to_map */
2424     kmt_name_to_map,    /* t_name_to_map */
2425     kmt_addr_to_ctf,    /* t_addr_to_ctf */
2426     kmt_name_to_ctf,    /* t_name_to_ctf */
2427     kmt_status,         /* t_status */

```

```

2428     (int (*)()) mdb_tgt_notsup, /* t_run */
2429     kmt_step,                 /* t_step */
2430     kmt_step_out,            /* t_step_out */
2438     kmt_step_branch,        /* t_step_branch */
2431     kmt_next,                /* t_next */
2432     kmt_continue,           /* t_cont */
2433     (int (*)()) mdb_tgt_notsup, /* t_signal */
2434     kmt_add_vbrkpt,         /* t_add_vbrkpt */
2435     kmt_add_sbrkpt,         /* t_add_sbrkpt */
2436     kmt_add_pwapt,          /* t_add_pwapt */
2437     kmt_add_vwapt,          /* t_add_vwapt */
2438     kmt_add_iowapt,         /* t_add_iowapt */
2439     (int (*)()) mdb_tgt_null, /* t_add_sysenter */
2440     (int (*)()) mdb_tgt_null, /* t_add_sysexit */
2441     (int (*)()) mdb_tgt_null, /* t_add_signal */
2442     kmt_add_trap,           /* t_add_fault */
2443     kmt_getareg,            /* t_getareg */
2444     kmt_putareg,           /* t_putareg */
2445     (int (*)()) mdb_tgt_nop, /* XXX t_stack_iter */
2446     (int (*)()) mdb_tgt_notsup /* t_auxv */
2447 };

2449 /*
2450  * Called immediately upon resumption of the system after a step or continue.
2451  * Allows us to synchronize kmt's view of the world with reality.
2452  */
2453 /*ARGSUSED*/
2454 static void
2455 kmt_sync(mdb_tgt_t *t)
2456 {
2457     kmt_data_t *kmt = t->t_data;
2458     int symavail;

2460     mdb_dprintf(MDB_DBG_KMOD, "synchronizing with kernel\n");

2462     symavail = kmt->kmt_symavail;
2463     kmt->kmt_symavail = FALSE;

2465     /*
2466      * Resync our view of the world if the modules have changed, or if we
2467      * didn't have any symbols coming into this function. The latter will
2468      * only happen on startup.
2469      */
2470     if (kmdb_kdi_mods_changed() || !symavail)
2471         kmt_modlist_update(t);

2473     /*
2474      * It would be nice if we could run this less frequently, perhaps
2475      * after a dvec-initiated trigger.
2476      */
2477     kmdb_module_sync();

2479     kmt->kmt_symavail = TRUE;

2481     mdb_dprintf(MDB_DBG_KMOD, "synchronization complete\n");

2483     kmt_defbp_prune();

2485     if (kmt_defbp_num > 0 && kmt_defbp_bpspec == 0 &&
2486         kmdb_kdi_dtrace_get_state() != KDI_DTSTATE_DTRACE_ACTIVE) {
2487         /*
2488          * Deferred breakpoints were created while DTrace was active,
2489          * and consequently the deferred breakpoint enabling mechanism
2490          * wasn't activated. Activate it now, and then try to activate
2491          * the deferred breakpoints. We do this so that we can catch
2492          * the ones which may apply to modules that have been loaded

```

```
2493         * while they were waiting for DTrace to deactivate.
2494         */
2495         (void) kmt_defbp_activate(t);
2496         (void) mdb_tgt_sespec_activate_all(t);
2497     }
2498
2499     if (kmt->kmt_cpu_retry && ((kmt->kmt_cpu = kmt_cpu_create(t)) !=
2500         NULL || errno != EAGAIN))
2501         kmt->kmt_cpu_retry = FALSE;
2502
2503     (void) mdb_tgt_status(t, &t->t_status);
2504 }
2505
2506 /*
2507  * This routine executes while the kernel is running.
2508  */
2509 /*ARGSUSED*/
2510 int
2511 kmdb_kvm_create(mdb_tgt_t *t, int argc, const char *argv[])
2512 {
2513     kmt_data_t *kmt;
2514
2515     if (argc != 0)
2516         return (set_errno(EINVAL));
2517
2518     kmt = mdb_zalloc(sizeof (kmt_data_t), UM_SLEEP);
2519     t->t_data = kmt;
2520     t->t_ops = &kmt_ops;
2521     t->t_flags |= MDB_TGT_F_RDWR; /* kmdb is always r/w */
2522
2523     (void) mdb_nv_insert(&mdb.m_nv, "cpuid", &kmt_cpuid_disc, 0,
2524         MDB_NV_PERSIST | MDB_NV_RDONLY);
2525
2526     (void) mdb_nv_create(&kmt->kmt_modules, UM_SLEEP);
2527
2528     kmt_init_isadep(t);
2529
2530     kmt->kmt_symavail = FALSE;
2531     kmt->kmt_cpu_retry = TRUE;
2532
2533     bzero(&kmt_defbp_list, sizeof (mdb_list_t));
2534
2535     return (0);
2536
2537 create_err:
2538     kmt_destroy(t);
2539
2540     return (-1);
2541 }
2542
2543 _____unchanged_portion_omitted_____
```

```

*****
4411 Fri Apr 6 17:24:53 2018
new/usr/src/cmd/mdb/common/kmdb/kvm.h
9210 remove KMDB branch debugging support
9211 ::crregs could do with cr2/cr3 support
9209 ::ttrace should be able to filter by thread
Reviewed by: Patrick Mooney <patrick.mooney@joyent.com>
Reviewed by: Yuri Pankov <yuripv@yuripv.net>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22  * Copyright 2007 Sun Microsystems, Inc. All rights reserved.
23  * Use is subject to license terms.
24  *
25  * Copyright 2018 Joyent, Inc.
26  */

28 #ifndef _KVM_H
29 #define _KVM_H

29 #pragma ident "%Z%M% %I% %E% SMI"

31 /*
32  * The kmdb target
33  */

35 #include <mdb/mdb_modapi.h>
36 #include <mdb/mdb_target.h>
37 #include <kmdb/kmdb_dpi.h>
38 #include <kmdb/kvm_isadep.h>
39 #include <kmdb/kvm_cpu.h>

40 #include <sys/kobj.h>

42 #ifdef __cplusplus
43 extern "C" {
44 #endif

46 #define KM_F_PRIMARY 1

48 #define KMT_TRAP_NOTENUM -1 /* Glob for unnamed traps */
49 #define KMT_TRAP_ALL -2 /* Glob for all traps */

51 typedef struct kmt_module {
52     mdb_list_t km_list; /* List forward/back pointers */
53     char *km_name; /* Module name */
54     char km_seen;

```

```

55     GElf_Ehdr km_ehdr;

57     mdb_gelf_symtab_t *km_symtab;
58     Shdr km_symtab_hdr;
59     Shdr km_strtab_hdr;
60     const void *km_symtab_va;
61     const void *km_strtab_va;

63     uintptr_t km_text_va;
64     size_t km_text_size;
65     uintptr_t km_data_va;
66     size_t km_data_size;
67     uintptr_t km_bss_va;
68     size_t km_bss_size;
69     const void *km_ctf_va;
70     size_t km_ctf_size;

72     ctf_file_t *km_ctfp;
73     struct modctl km_modctl;
74     struct module km_module;
75     int km_flags;
76 } kmt_module_t;

78 typedef struct kmt_data {
79     const mdb_tgt_regdesc_t *kmt_rds; /* Register description table */
80     mdb_nv_t kmt_modules; /* Hash table of modules */
81     mdb_list_t kmt_modlist; /* List of mods in load order */
82     const char *kmt_rtld_name; /* Module containing krtld */
83     caddr_t kmt_writemap; /* Used to map PAs for writes */
84     size_t kmt_writemapsz; /* Size of same */
85     mdb_map_t kmt_map; /* Persistent map for callers */
86     ulong_t *kmt_trapmap;
87     size_t kmt_trapmax;
89     kmt_cpu_t *kmt_cpu; /* CPU-specific plugin */
90     int kmt_cpu_retry; /* Try CPU detect again? */
88     int kmt_symavail; /* Symbol resolution allowed */
89     uint_t kmt_narmedbrkpts; /* Number of armed brkpts */
90 #if defined(__i386) || defined(__amd64)
91     struct {
92         GElf_Sym _kmt_cmhint;
93         GElf_Sym _kmt_cmnttrap;
94         GElf_Sym _kmt_sysenter;
95         GElf_Sym _kmt_brand_sysenter;
96 #if defined(__amd64)
97         GElf_Sym _kmt_syscall;
98         GElf_Sym _kmt_brand_syscall;
99 #endif
100     } kmt_intrsyms;
101 #endif
102 } kmt_data_t;

_____unchanged_portion_omitted_____

127 extern void kmt_printregs(const mdb_tgt_gregset_t *gregs);

129 extern const char *kmt_def_dismode(void);

131 extern void kmt_init_isadep(mdb_tgt_t *);
132 extern void kmt_startup_isadep(mdb_tgt_t *);

134 extern ssize_t kmt_write(mdb_tgt_t *, const void *, size_t, uintptr_t);
135 extern ssize_t kmt_pwrite(mdb_tgt_t *, const void *, size_t, physaddr_t);
136 extern ssize_t kmt_rw(mdb_tgt_t *, void *, size_t, uint64_t,
137     ssize_t (*)(void *, size_t, uint64_t));
138 extern ssize_t kmt_writer(void *, size_t, uint64_t);
139 extern ssize_t kmt_ioread(mdb_tgt_t *, void *, size_t, uintptr_t);
140 extern ssize_t kmt_iowrite(mdb_tgt_t *, const void *, size_t, uintptr_t);

```

```
142 extern int kmt_step_out(mdb_tgt_t *, uintptr_t *);
146 extern int kmt_step_branch(mdb_tgt_t *);
143 extern int kmt_next(mdb_tgt_t *, uintptr_t *);

145 extern int kmt_stack(uintptr_t, uint_t, int, const mdb_arg_t *);
146 extern int kmt_stackv(uintptr_t, uint_t, int, const mdb_arg_t *);
147 extern int kmt_stackr(uintptr_t, uint_t, int, const mdb_arg_t *);
148 extern int kmt_cpustack(uintptr_t, uint_t, int, const mdb_arg_t *, int, int);

150 extern const char *kmt_trapname(int);

152 #ifdef __cplusplus
153 }
_____unchanged_portion_omitted_____
```



```

*****
84154 Fri Apr 6 17:24:53 2018
new/usr/src/cmd/mdb/common/mdb/mdb_cmds.c
9210 remove KMDB branch debugging support
9211 ::crregs could do with cr2/cr3 support
9209 ::ttrace should be able to filter by thread
Reviewed by: Patrick Mooney <patrick.mooney@joyent.com>
Reviewed by: Yuri Pankov <yuripv@yuripv.net>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */

22 /*
23  * Copyright 2009 Sun Microsystems, Inc. All rights reserved.
24  * Use is subject to license terms.
25 */

27 /*
28  * Copyright (c) 2012 by Delphix. All rights reserved.
29  * Copyright (c) 2018 Joyent, Inc. All rights reserved.
30  * Copyright (c) 2015 Joyent, Inc. All rights reserved.
31  * Copyright (c) 2013 Josef 'Jeff' Sipek <jeffpc@josefsipek.net>
32  * Copyright (c) 2015, 2017 by Delphix. All rights reserved.
33  */

34 #include <sys/elf.h>
35 #include <sys/elf_SPARC.h>

37 #include <libproc.h>
38 #include <stdlib.h>
39 #include <string.h>
40 #include <fcntl.h>
41 #include <errno.h>
42 #include <alloca.h>
43 #include <libctf.h>
44 #include <ctype.h>

46 #include <mdb/mdb_string.h>
47 #include <mdb/mdb_argvec.h>
48 #include <mdb/mdb_nv.h>
49 #include <mdb/mdb_fmt.h>
50 #include <mdb/mdb_target.h>
51 #include <mdb/mdb_err.h>
52 #include <mdb/mdb_debug.h>
53 #include <mdb/mdb_conf.h>
54 #include <mdb/mdb_module.h>
55 #include <mdb/mdb_modapi.h>
56 #include <mdb/mdb_stdlib.h>

```

```

57 #include <mdb/mdb_lex.h>
58 #include <mdb/mdb_io_impl.h>
59 #include <mdb/mdb_help.h>
60 #include <mdb/mdb_disasm.h>
61 #include <mdb/mdb_frame.h>
62 #include <mdb/mdb_evset.h>
63 #include <mdb/mdb_print.h>
64 #include <mdb/mdb_nm.h>
65 #include <mdb/mdb_set.h>
66 #include <mdb/mdb_demangle.h>
67 #include <mdb/mdb_ctf.h>
68 #include <mdb/mdb_whatish>
69 #include <mdb/mdb_whatish_impl.h>
70 #include <mdb/mdb_macalias.h>
71 #include <mdb/mdb_tab.h>
72 #include <mdb/mdb_typedef.h>
73 #ifdef _KMDB
74 #include <kmdb/kmdb_kdi.h>
75 #endif
76 #include <mdb/mdb.h>

78 #ifdef __sparc
79 #define SETHI_MASK      0xc1c00000
80 #define SETHI_VALUE    0x01000000

82 #define IS_SETHI(machcode)      (((machcode) & SETHI_MASK) == SETHI_VALUE)

84 #define OP(machcode)      ((machcode) >> 30)
85 #define OP3(machcode)    (((machcode) >> 19) & 0x3f)
86 #define RD(machcode)    (((machcode) >> 25) & 0x1f)
87 #define RS1(machcode)   (((machcode) >> 14) & 0x1f)
88 #define I(machcode)     (((machcode) >> 13) & 0x01)

90 #define IMM13(machcode)  ((machcode) & 0x1fff)
91 #define IMM22(machcode) ((machcode) & 0x3fffff)

93 #define OP_ARITH_MEM_MASK 0x2
94 #define OP_ARITH         0x2
95 #define OP_MEM           0x3

97 #define OP3_CC_MASK      0x10
98 #define OP3_COMPLEX_MASK 0x20

100 #define OP3_ADD          0x00
101 #define OP3_OR           0x02
102 #define OP3_XOR          0x03

104 #ifndef R_07
105 #define R_07             0xf
106 #endif
107 #endif /* __sparc */

109 static mdb_tgt_addr_t
110 write_uint8(mdb_tgt_as_t as, mdb_tgt_addr_t addr, uint64_t ull, uint_t rdback)
111 {
112     uint8_t o, n = (uint8_t)ull;

114     if (rdback && mdb_tgt_aread(mdb.m_target, as, &o, sizeof (o),
115         addr) == -1)
116         return (addr);

118     if (mdb_tgt_awrite(mdb.m_target, as, &n, sizeof (n), addr) == -1)
119         return (addr);

121     if (rdback) {
122         if (mdb_tgt_aread(mdb.m_target, as, &n, sizeof (n), addr) == -1)

```

```
123         return (addr);
124
125         mdb_iob_printf(mdb.m_out, "%-#*lla%16T%-#8x=%8T0x%x\n",
126             mdb_iob_getmargin(mdb.m_out), addr, o, n);
127     }
128
129     return (addr + sizeof (n));
130 }
    _____unchanged_portion_omitted_____
2734 static int
2735 cmd_step(uintptr_t addr, uint_t flags, int argc, const mdb_arg_t *argv)
2736 {
2737     int (*func)(mdb_tgt_t *, mdb_tgt_status_t *) = &mdb_tgt_step;
2738     const char *name = "single-step";
2739
2740     if (argc > 0 && argv->a_type == MDB_TYPE_STRING) {
2741         if (strcmp(argv->a_un.a_str, "out") == 0) {
2742             func = &mdb_tgt_step_out;
2743             name = "step (out)";
2744             argv++;
2745             argc--;
2746         } else if (strcmp(argv->a_un.a_str, "branch") == 0) {
2747             func = &mdb_tgt_step_branch;
2748             name = "step (branch)";
2749             argv++;
2750             argc--;
2746         } else if (strcmp(argv->a_un.a_str, "over") == 0) {
2747             func = &mdb_tgt_next;
2748             name = "step (over)";
2749             argv++;
2750             argc--;
2751         }
2752     }
2753
2754     return (cmd_cont_common(addr, flags, argc, argv, func, name));
2755 }
    _____unchanged_portion_omitted_____
```

```

*****
30681 Fri Apr 6 17:24:54 2018
new/usr/src/cmd/mdb/common/mdb/mdb_kproc.c
9210 remove KMDB branch debugging support
9211 ::crregs could do with cr2/cr3 support
9209 ::ttrace should be able to filter by thread
Reviewed by: Patrick Mooney <patrick.mooney@joyent.com>
Reviewed by: Yuri Pankov <yuripv@yuripv.net>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22  * Copyright 2008 Sun Microsystems, Inc. All rights reserved.
23  * Use is subject to license terms.
24  *
25  * Copyright 2018 Joyent, Inc.
26  */

26 #pragma ident "%Z%M% %I% %E% SMI"

28 /*
29  * Kernel Process View Target
30  *
31  * The kproc target is activated when the user is debugging a kernel using the
32  * kvm target and executes a ::context dcmd to change the debugger view to one
33  * of the running processes. The kvm target's t_setcontext operation will
34  * create and activate a kproc target in response to this call. The kproc
35  * target itself is built upon the kvm target's libkvm cookie and the ability
36  * to read information from the kernel itself and the ability to read the
37  * address space of a particular user process with kvm_aread(). It also relies
38  * on a special set of functions provided by the kvm target's mdb_ks support
39  * module in order to bootstrap: specifically, given the initial proc pointer,
40  * mdb_ks provides functions to return the set of address space mappings, the
41  * address space pointer itself, the aux vector saved in the u-area,
42  * and the process data model. The kproc target maintains a list of address
43  * space mappings (kp_map_t) and load objects (kp_file_t), and for each load
44  * object will attempt to read the corresponding dynamic symbol table. In
45  * order to bootstrap, the target uses the AT_BASE and AT_ENTRY aux vector
46  * elements to locate the dynamic linker and executable mappings. With these
47  * mappings in place, we initialize a librtld_db agent on the target (see
48  * mdb_pservice.c for how this is done), and then process each load object
49  * found in the link-map chain. In order to simplify the construction of
50  * symbol tables for each load object, we would like make use of our existing
51  * library of GElf processing code. Since the MDB GElf code uses mdb_io
52  * objects to read in an ELF file, we simply define a new type of mdb_io object
53  * where each read operation is translated into a call to kproc's t_vread
54  * function to read from the range of the address space defined by the mapping
55  * as if it were a file.

```

```

56 */
58 #include <sys/types.h>
59 #include <sys/proc.h>
60 #include <sys/auxv.h>
62 #include <strings.h>
63 #include <limits.h>
64 #include <rtld_db.h>
65 #include <procfcs.h>
66 #include <dlfcn.h>
67 #include <kvm.h>
69 #include <mdb/mdb_target_impl.h>
70 #include <mdb/mdb_debug.h>
71 #include <mdb/mdb_string.h>
72 #include <mdb/mdb_err.h>
73 #include <mdb/mdb_ks.h>
74 #include <mdb/mdb_gelf.h>
75 #include <mdb/mdb_io_impl.h>
76 #include <mdb/mdb.h>
78 typedef struct kp_symarg {
79     mdb_tgt_sym_f *sym_cb;           /* Caller's callback function */
80     void *sym_data;                 /* Callback function argument */
81     uint_t sym_type;                /* Symbol type/binding filter */
82     uintptr_t sym_adjust;           /* Symbol value adjustment */
83     mdb_syminfo_t sym_info;         /* Symbol id and table id */
84     const char *sym_obj;            /* Containing object */
85 } kp_symarg_t;
86 #define unchanged_portion_omitted

886 static const mdb_tgt_ops_t kproc_ops = {
887     (int (*)()) mdb_tgt_notsup,      /* t_setflags */
888     kp_setcontext,                  /* t_setcontext */
889     kp_activate,                    /* t_activate */
890     kp_deactivate,                  /* t_deactivate */
891     (void (*)()) mdb_tgt_nop,       /* t_periodic */
892     kp_destroy,                     /* t_destroy */
893     kp_name,                         /* t_name */
894     kp_isa,                          /* t_isa */
895     kp_platform,                    /* t_platform */
896     kp_uname,                        /* t_uname */
897     kp_dmodel,                      /* t_dmodel */
898     (ssize_t (*)()) mdb_tgt_notsup, /* t_aread */
899     (ssize_t (*)()) mdb_tgt_notsup, /* t_awrite */
900     kp_vread,                       /* t_vread */
901     kp_vwrite,                      /* t_vwrite */
902     (ssize_t (*)()) mdb_tgt_notsup, /* t_pread */
903     (ssize_t (*)()) mdb_tgt_notsup, /* t_pwrite */
904     (ssize_t (*)()) mdb_tgt_notsup, /* t_fread */
905     (ssize_t (*)()) mdb_tgt_notsup, /* t_fwrite */
906     (ssize_t (*)()) mdb_tgt_notsup, /* t_ioread */
907     (ssize_t (*)()) mdb_tgt_notsup, /* t_iowrite */
908     kp_vtop,                        /* t_vtop */
909     kp_lookup_by_name,               /* t_lookup_by_name */
910     kp_lookup_by_addr,              /* t_lookup_by_addr */
911     kp_symbol_iter,                 /* t_symbol_iter */
912     kp_mapping_iter,                /* t_mapping_iter */
913     kp_object_iter,                 /* t_object_iter */
914     kp_addr_to_map,                  /* t_addr_to_map */
915     kp_name_to_map,                  /* t_name_to_map */
916     (struct ctf_file *(*())()) mdb_tgt_null, /* t_addr_to_ctf */
917     (struct ctf_file *(*())()) mdb_tgt_null, /* t_name_to_ctf */
918     kp_status,                       /* t_status */
919     (int (*)()) mdb_tgt_notsup,      /* t_run */

```

```
920     (int (*)()) mdb_tgt_notsup,      /* t_step */
921     (int (*)()) mdb_tgt_notsup,      /* t_step_out */
922     (int (*)()) mdb_tgt_notsup,      /* t_step_branch */
922     (int (*)()) mdb_tgt_notsup,      /* t_next */
923     (int (*)()) mdb_tgt_notsup,      /* t_cont */
924     (int (*)()) mdb_tgt_notsup,      /* t_signal */
925     (int (*)()) mdb_tgt_null,         /* t_add_sbrkpt */
926     (int (*)()) mdb_tgt_null,         /* t_add_vbrkpt */
927     (int (*)()) mdb_tgt_null,         /* t_add_pwapt */
928     (int (*)()) mdb_tgt_null,         /* t_add_vwapt */
929     (int (*)()) mdb_tgt_null,         /* t_add_iowapt */
930     (int (*)()) mdb_tgt_null,         /* t_add_sysenter */
931     (int (*)()) mdb_tgt_null,         /* t_add_sysexit */
932     (int (*)()) mdb_tgt_null,         /* t_add_signal */
933     (int (*)()) mdb_tgt_null,         /* t_add_fault */
934     (int (*)()) mdb_tgt_notsup,       /* t_getareg XXX */
935     (int (*)()) mdb_tgt_notsup,       /* t_putareg XXX */
936     (int (*)()) mdb_tgt_notsup,       /* t_stack_iter XXX */
937     kp_auxv                            /* t_auxv */
938 };
_____unchanged_portion_omitted_____
```

new/usr/src/cmd/mdb/common/mdb/mdb_proc.c

1

```
*****
148244 Fri Apr 6 17:24:54 2018
new/usr/src/cmd/mdb/common/mdb/mdb_proc.c
9210 remove KMDB branch debugging support
9211 ::crregs could do with cr2/cr3 support
9209 ::ttrace should be able to filter by thread
Reviewed by: Patrick Mooney <patrick.mooney@joyent.com>
Reviewed by: Yuri Pankov <yuripv@yuripv.net>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */

22 /*
23  * Copyright 2010 Sun Microsystems, Inc. All rights reserved.
24  * Use is subject to license terms.
25 */
26 /*
27  * Copyright 2018 Joyent, Inc.
27  * Copyright 2015 Joyent, Inc.
28  * Copyright (c) 2014 by Delphix. All rights reserved.
29 */

31 /*
32  * User Process Target
33  *
34  * The user process target is invoked when the -u or -p command-line options
35  * are used, or when an ELF executable file or ELF core file is specified on
36  * the command-line. This target is also selected by default when no target
37  * options are present. In this case, it defaults the executable name to
38  * "a.out". If no process or core file is currently attached, the target
39  * functions as a kind of virtual /dev/zero (in accordance with adb(1)
40  * semantics); reads from the virtual address space return zeroes and writes
41  * fail silently. The proc target itself is designed as a wrapper around the
42  * services provided by libproc.so: t->t_pshandle is set to the struct
43  * ps_prochandle pointer returned as a handle by libproc. The target also
44  * opens the executable file itself using the MDB GElf services, for
45  * interpreting the .symtab and .dynsym if no libproc handle has been
46  * initialized, and for handling i/o to and from the object file. Currently,
47  * the only ISA-dependent portions of the proc target are the $r and ::fpregs
48  * dcmds, the callbacks for t_next() and t_step_out(), and the list of named
49  * registers; these are linked in from the proc_isadep.c file for each ISA and
50  * called from the common code in this file.
51  *
52  * The user process target implements complete user process control using the
53  * facilities provided by libproc.so. The MDB execution control model and
54  * an overview of software event management is described in mdb_target.c. The
55  * proc target implements breakpoints by replacing the instruction of interest
56  * with a trap instruction, and then restoring the original instruction to step
```

new/usr/src/cmd/mdb/common/mdb/mdb_proc.c

2

```
57  * over the breakpoint. The idea of replacing program text with instructions
58  * that transfer control to the debugger dates back as far as 1951 [1]. When
59  * the target stops, we replace each breakpoint with the original instruction
60  * as part of the disarm operation. This means that no special processing is
61  * required for t_vread() because the instrumented instructions will never be
62  * seen by the debugger once the target stops. Some debuggers have improved
63  * start/stop performance by leaving breakpoint traps in place and then
64  * handling a read from a breakpoint address as a special case. Although this
65  * improves efficiency for a source-level debugger, it runs somewhat contrary
66  * to the philosophy of the low-level debugger. Since we remove the
67  * instructions, users can apply other external debugging tools to the process
68  * once it has stopped (e.g. the proc(1) tools) and not be misled by MDB
69  * instrumentation. The tracing of faults, signals, system calls, and
70  * watchpoints and general process inspection is implemented directly using
71  * the mechanisms provided by /proc, as described originally in [2] and [3].
72  *
73  * References
74  *
75  * [1] S. Gill, "The Diagnosis Of Mistakes In Programmes on the EDSAC",
76  * Proceedings of the Royal Society Series A Mathematical and Physical
77  * Sciences, Cambridge University Press, 206(1087), May 1951, pp. 538-554.
78  *
79  * [2] T.J. Killian, "Processes as Files", Proceedings of the USENIX Association
80  * Summer Conference, Salt Lake City, June 1984, pp. 203-207.
81  *
82  * [3] Roger Faulkner and Ron Gomes, "The Process File System and Process
83  * Model in UNIX System V", Proceedings of the USENIX Association
84  * Winter Conference, Dallas, January 1991, pp. 243-252.
85  */

87 #include <mdb/mdb_proc.h>
88 #include <mdb/mdb_disasm.h>
89 #include <mdb/mdb_signal.h>
90 #include <mdb/mdb_string.h>
91 #include <mdb/mdb_module.h>
92 #include <mdb/mdb_debug.h>
93 #include <mdb/mdb_conf.h>
94 #include <mdb/mdb_err.h>
95 #include <mdb/mdb_types.h>
96 #include <mdb/mdb.h>

98 #include <sys/utsname.h>
99 #include <sys/wait.h>
100 #include <sys/stat.h>
101 #include <termio.h>
102 #include <signal.h>
103 #include <stdio_ext.h>
104 #include <stdlib.h>
105 #include <string.h>

107 #define PC_FAKE -1UL /* illegal pc value unequal 0 */
108 #define PANIC_BUFSIZE 1024

110 static const char PT_EXEC_PATH[] = "a.out"; /* Default executable */
111 static const char PT_CORE_PATH[] = "core"; /* Default core file */

113 static const pt_ptl_ops_t proc_lwp_ops;
114 static const pt_ptl_ops_t proc_tdb_ops;
115 static const mdb_se_ops_t proc_brkpt_ops;
116 static const mdb_se_ops_t proc_wapt_ops;

118 static int pt_setrun(mdb_tgt_t *, mdb_tgt_status_t *, int);
119 static void pt_activate_common(mdb_tgt_t *);
120 static mdb_tgt-vespec_f pt_ignore_sig;
121 static mdb_tgt-se_f pt_fork;
122 static mdb_tgt-se_f pt_exec;
```

```

124 static int pt_lookup_by_name_thr(mdb_tgt_t *, const char *,
125     const char *, GElf_Sym *, mdb_syminfo_t *, mdb_tgt_tid_t);
126 static int tblbase(mdb_tgt_t *, mdb_tgt_tid_t, Lmid_t, const char *,
127     psaddr_t *);
129 /*
130  * When debugging postmortem, we don't resolve names as we may very well not
131  * be on a system on which those names resolve.
132  */
133 #define PT_LIBPROC_RESOLVE(P) \
134     (!(mdb.m_flags & MDB_FL_LMRAW) && Pstate(P) != PS_DEAD)
136 /*
137  * The Perror_printf() function interposes on the default, empty libproc
138  * definition. It will be called to report additional information on complex
139  * errors, such as a corrupt core file. We just pass the args to vwarn.
140  */
141 /*ARGSUSED*/
142 void
143 Perror_printf(struct ps_prochandle *P, const char *format, ...)
144 {
145     va_list alist;
147     va_start(alist, format);
148     vwarn(format, alist);
149     va_end(alist);
150 }

```

unchanged portion omitted

```

4662 static const mdb_tgt_ops_t proc_ops = {
4663     pt_setflags, /* t_setflags */
4664     (int (*)()) mdb_tgt_notsup, /* t_setcontext */
4665     pt_activate, /* t_activate */
4666     pt_deactivate, /* t_deactivate */
4667     pt_periodic, /* t_periodic */
4668     pt_destroy, /* t_destroy */
4669     pt_name, /* t_name */
4670     (const char (*)()) mdb_conf_isa, /* t_isa */
4671     pt_platform, /* t_platform */
4672     pt_uname, /* t_uname */
4673     pt_dmodel, /* t_dmodel */
4674     (ssize_t (*)()) mdb_tgt_notsup, /* t_aread */
4675     (ssize_t (*)()) mdb_tgt_notsup, /* t_awrite */
4676     pt_vread, /* t_vread */
4677     pt_vwrite, /* t_vwrite */
4678     (ssize_t (*)()) mdb_tgt_notsup, /* t_pread */
4679     (ssize_t (*)()) mdb_tgt_notsup, /* t_pwrite */
4680     pt_fread, /* t_fread */
4681     pt_fwrite, /* t_fwrite */
4682     (ssize_t (*)()) mdb_tgt_notsup, /* t_iread */
4683     (ssize_t (*)()) mdb_tgt_notsup, /* t_iowrite */
4684     (int (*)()) mdb_tgt_notsup, /* t_vtop */
4685     pt_lookup_by_name, /* t_lookup_by_name */
4686     pt_lookup_by_addr, /* t_lookup_by_addr */
4687     pt_symbol_iter, /* t_symbol_iter */
4688     pt_mapping_iter, /* t_mapping_iter */
4689     pt_object_iter, /* t_object_iter */
4690     pt_addr_to_map, /* t_addr_to_map */
4691     pt_name_to_map, /* t_name_to_map */
4692     pt_addr_to_ctf, /* t_addr_to_ctf */
4693     pt_name_to_ctf, /* t_name_to_ctf */
4694     pt_status, /* t_status */
4695     pt_run, /* t_run */
4696     pt_step, /* t_step */

```

```

4697     pt_step_out, /* t_step_out */
4698     (int (*)()) mdb_tgt_notsup, /* t_step_branch */
4698     pt_next, /* t_next */
4699     pt_continue, /* t_cont */
4700     pt_signal, /* t_signal */
4701     pt_add_vbrkpt, /* t_add_vbrkpt */
4702     pt_add_sbrkpt, /* t_add_sbrkpt */
4703     (int (*)()) mdb_tgt_null, /* t_add_pwapt */
4704     pt_add_vwapt, /* t_add_vwapt */
4705     (int (*)()) mdb_tgt_null, /* t_add_iowapt */
4706     pt_add_sysenter, /* t_add_sysenter */
4707     pt_add_sysexit, /* t_add_sysexit */
4708     pt_add_signal, /* t_add_signal */
4709     pt_add_fault, /* t_add_fault */
4710     pt_getareg, /* t_getareg */
4711     pt_putareg, /* t_putareg */
4712     pt_stack_iter, /* t_stack_iter */
4713     pt_auxv, /* t_auxv */
4714 };

```

unchanged portion omitted

```

*****
11832 Fri Apr 6 17:24:54 2018
new/usr/src/cmd/mdb/common/mdb/mdb_rawfile.c
9210 remove KMDB branch debugging support
9211 ::crregs could do with cr2/cr3 support
9209 ::ttrace should be able to filter by thread
Reviewed by: Patrick Mooney <patrick.mooney@joyent.com>
Reviewed by: Yuri Pankov <yuripv@yuripv.net>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22  * Copyright 2008 Sun Microsystems, Inc. All rights reserved.
23  * Use is subject to license terms.
24  *
25  * Copyright 2018 Joyent, Inc.
26  */

26 #pragma ident "%Z%M% %I% %E% SMI"

28 /*
29  * Raw File Target
30  *
31  * The raw file target is invoked whenever a file of unrecognizable type is
32  * specified on the command line, or when raw file examination is forced using
33  * the -f option. If one file is specified, that file will be opened as the
34  * "object" file. If two files are specified, the second one will be opened
35  * as the "core" file. Each file is opened using the fdio backend, which
36  * internally supports both byte-oriented i/o and block-oriented i/o as needed.
37  */

39 #include <mdb/mdb_modapi.h>
40 #include <mdb/mdb_target_impl.h>
41 #include <mdb/mdb_io_impl.h>
42 #include <mdb/mdb_conf.h>
43 #include <mdb/mdb_err.h>
44 #include <mdb/mdb.h>

46 #include <sys/dtrace.h>
47 #include <fcntl.h>

49 typedef struct rf_data {
50     mdb_io_t *r_object_fio;
51     mdb_io_t *r_core_fio;
52 } rf_data_t;
    unchanged_portion_omitted_

358 static const mdb_tgt_ops_t rawfile_ops = {

```

```

359     rf_setflags, /* t_setflags */
360     (int (*)()) mdb_tgt_notsup, /* t_setcontext */
361     rf_activate, /* t_activate */
362     rf_deactivate, /* t_deactivate */
363     (void (*)()) mdb_tgt_nop, /* t_periodic */
364     rf_destroy, /* t_destroy */
365     rf_name, /* t_name */
366     (const char (*)()) mdb_conf_isa, /* t_isa */
367     (const char (*)()) mdb_conf_platform, /* t_platform */
368     (int (*)()) mdb_tgt_notsup, /* t_uname */
369     (int (*)()) mdb_tgt_notsup, /* t_dmodel */
370     rf_aread, /* t_aread */
371     rf_awrite, /* t_awrite */
372     rf_vread, /* t_vread */
373     rf_vwrite, /* t_vwrite */
374     rf_pread, /* t_pread */
375     rf_pwrite, /* t_pwrite */
376     rf_fread, /* t_fread */
377     rf_fwrite, /* t_fwrite */
378     (ssize_t (*)()) mdb_tgt_notsup, /* t_ioread */
379     (ssize_t (*)()) mdb_tgt_notsup, /* t_iowrite */
380     (int (*)()) mdb_tgt_notsup, /* t_vtop */
381     (int (*)()) mdb_tgt_notsup, /* t_lookup_by_name */
382     (int (*)()) mdb_tgt_notsup, /* t_lookup_by_addr */
383     (int (*)()) mdb_tgt_notsup, /* t_symbol_iter */
384     rf_mapping_iter, /* t_mapping_iter */
385     rf_mapping_iter, /* t_object_iter */
386     (const mdb_map_t (*)()) mdb_tgt_null, /* t_addr_to_map */
387     (const mdb_map_t (*)()) mdb_tgt_null, /* t_name_to_map */
388     (struct ctf_file (*)()) mdb_tgt_null, /* t_addr_to_ctf */
389     (struct ctf_file (*)()) mdb_tgt_null, /* t_name_to_ctf */
390     rf_status, /* t_status */
391     (int (*)()) mdb_tgt_notsup, /* t_run */
392     (int (*)()) mdb_tgt_notsup, /* t_step */
393     (int (*)()) mdb_tgt_notsup, /* t_step_out */
394     (int (*)()) mdb_tgt_notsup, /* t_step_branch */
394     (int (*)()) mdb_tgt_notsup, /* t_next */
395     (int (*)()) mdb_tgt_notsup, /* t_cont */
396     (int (*)()) mdb_tgt_notsup, /* t_signal */
397     (int (*)()) mdb_tgt_null, /* t_add_vbrkpt */
398     (int (*)()) mdb_tgt_null, /* t_add_sbrkpt */
399     (int (*)()) mdb_tgt_null, /* t_add_pwapt */
400     (int (*)()) mdb_tgt_null, /* t_add_vwapt */
401     (int (*)()) mdb_tgt_null, /* t_add_iowapt */
402     (int (*)()) mdb_tgt_null, /* t_add_sysenter */
403     (int (*)()) mdb_tgt_null, /* t_add_sysexit */
404     (int (*)()) mdb_tgt_null, /* t_add_signal */
405     (int (*)()) mdb_tgt_null, /* t_add_fault */
406     (int (*)()) mdb_tgt_notsup, /* t_getareg */
407     (int (*)()) mdb_tgt_notsup, /* t_putareg */
408     (int (*)()) mdb_tgt_notsup, /* t_stack_iter */
409     (int (*)()) mdb_tgt_notsup, /* t_auxv */
410 };
    unchanged_portion_omitted_

```

```

*****
64895 Fri Apr 6 17:24:54 2018
new/usr/src/cmd/mdb/common/mdb/mdb_target.c
9210 remove KMDB branch debugging support
9211 ::crregs could do with cr2/cr3 support
9209 ::ttrace should be able to filter by thread
Reviewed by: Patrick Mooney <patrick.mooney@joyent.com>
Reviewed by: Yuri Pankov <yuripv@yuripv.net>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22  * Copyright 2009 Sun Microsystems, Inc. All rights reserved.
23  * Use is subject to license terms.
24  *
25  * Copyright 2018 Joyent, Inc.
26  */
27
28 /*
29  * MDB Target Layer
30  *
31  * The *target* is the program being inspected by the debugger. The MDB target
32  * layer provides a set of functions that insulate common debugger code,
33  * including the MDB Module API, from the implementation details of how the
34  * debugger accesses information from a given target. Each target exports a
35  * standard set of properties, including one or more address spaces, one or
36  * more symbol tables, a set of load objects, and a set of threads that can be
37  * examined using the interfaces in <mdb/mdb_target.h>. This technique has
38  * been employed successfully in other debuggers, including [1], primarily
39  * to improve portability, although the term "target" often refers to the
40  * encapsulation of architectural or operating system-specific details. The
41  * target abstraction is useful for MDB because it allows us to easily extend
42  * the debugger to examine a variety of different program forms. Primarily,
43  * the target functions validate input arguments and then call an appropriate
44  * function in the target ops vector, defined in <mdb/mdb_target_impl.h>.
45  * However, this interface layer provides a very high level of flexibility for
46  * separating the debugger interface from instrumentation details. Experience
47  * has shown this kind of design can facilitate separating out debugger
48  * instrumentation into an external agent [2] and enable the development of
49  * advanced instrumentation frameworks [3]. We want MDB to be an ideal
50  * extensible framework for the development of such applications.
51  *
52  * Aside from a set of wrapper functions, the target layer also provides event
53  * management for targets that represent live executing programs. Our model of
54  * events is also extensible, and is based upon work in [3] and [4]. We define
55  * a *software event* as a state transition in the target program (for example,
56  * the transition of the program counter to a location of interest) that is
57  * observed by the debugger or its agent. A *software event specifier* is a

```

```

58 * description of a class of software events that is used by the debugger to
59 * instrument the target so that the corresponding software events can be
60 * observed. In MDB, software event specifiers are represented by the
61 * mdb_sespec_t structure, defined in <mdb/mdb_target_impl.h>. As the user,
62 * the internal debugger code, and MDB modules may all wish to observe software
63 * events and receive appropriate notification and callbacks, we do not expose
64 * software event specifiers directly as part of the user interface. Instead,
65 * clients of the target layer request that events be observed by creating
66 * new *virtual event specifiers*. Each virtual specifier is named by a unique
67 * non-zero integer (the VID), and is represented by a mdb_vespec_t structure.
68 * One or more virtual specifiers are then associated with each underlying
69 * software event specifier. This design enforces the constraint that the
70 * target must only insert one set of instrumentation, regardless of how many
71 * times the target layer was asked to trace a given event. For example, if
72 * multiple clients request a breakpoint at a particular address, the virtual
73 * specifiers will map to the same sespec, ensuring that only one breakpoint
74 * trap instruction is actually planted at the given target address. When no
75 * virtual specifiers refer to an sespec, it is no longer needed and can be
76 * removed, along with the corresponding instrumentation.
77 *
78 * The following state transition diagram illustrates the life cycle of a
79 * software event specifier and example transitions:
80 *
81 *
82 *
83 *
84 *
85 *
86 *
87 *
88 *
89 *
90 *
91 *
92 *
93 *
94 * The MDB execution control model is based upon the synchronous debugging
95 * model exported by Solaris proc(4). A target program is set running or the
96 * debugger is attached to a running target. On ISTOP (stop on event of
97 * interest), one target thread is selected as the representative. The
98 * algorithm for selecting the representative is target-specific, but we assume
99 * that if an observed software event has occurred, the target will select the
100 * thread that triggered the state transition of interest. The other threads
101 * are stopped in sympathy with the representative as soon as possible. Prior
102 * to continuing the target, we plant our instrumentation, transitioning event
103 * specifiers from the ACTIVE to the ARMED state, and then back again when the
104 * target stops. We then query each active event specifier to learn which ones
105 * are matched, and then invoke the callbacks associated with their vespecs.
106 * If an OS error occurs while attempting to arm or disarm a specifier, the
107 * specifier is transitioned to the ERROR state; we will attempt to arm it
108 * again at the next continue. If no target process is under our control or
109 * if an event is not currently applicable (e.g. a deferred breakpoint on an
110 * object that is not yet loaded), it remains in the IDLE state. The target
111 * implementation should intercept object load events and then transition the
112 * specifier to the ACTIVE state when the corresponding object is loaded.
113 *
114 * To simplify the debugger implementation and allow targets to easily provide
115 * new types of observable events, most of the event specifier management is
116 * done by the target layer. Each software event specifier provides an ops
117 * vector of subroutines that the target layer can call to perform the
118 * various state transitions described above. The target maintains two lists
119 * of mdb_sespec_t's: the t_idle list (IDLE state) and the t_active list
120 * (ACTIVE, ARMED, and ERROR states). Each mdb_sespec_t maintains a list of
121 * associated mdb_vespec_t's. If an sespec is IDLE or ERROR, its se_errno
122 * field will have an errno value specifying the reason for its inactivity.
123 * The vespec stores the client's callback function and private data, and the

```

```

+-----+ delete +-----+ cont/ +-----+
(| ( DEAD )) <----- ( ACTIVE ) <----- ( ARMED )
+-----+
^ load/unload ^ failure/
delete | object / \ reset | failure
^ v
+-----+ +-----+
+---- ( IDLE ) ( ERR ) <----+
+-----+ +-----+
| |
+-----+

```



```

124 * arguments used to construct the sespec. All objects are reference counted
125 * so we can destroy an object when it is no longer needed. The mdb_sespec_t
126 * invariants for the respective states are as follows:
127 *
128 * IDLE: on t_idle list, se_data == NULL, se_errno != 0, se_ctor not called
129 * ACTIVE: on t_active list, se_data valid, se_errno == 0, se_ctor called
130 * ARMED: on t_active list, se_data valid, se_errno == 0, se_ctor called
131 * ERROR: on t_active list, se_data valid, se_errno != 0, se_ctor called
132 *
133 * Additional commentary on specific state transitions and issues involving
134 * event management can be found below near the target layer functions.
135 *
136 * References
137 *
138 * [1] John Gilmore, "Working in GDB", Technical Report, Cygnus Support,
139 *     1.84 edition, 1994.
140 *
141 * [2] David R. Hanson and Mukund Raghavachari, "A Machine-Independent
142 *     Debugger", Software--Practice and Experience, 26(11), 1277-1299(1996).
143 *
144 * [3] Michael W. Shapiro, "RDB: A System for Incremental Replay Debugging",
145 *     Technical Report CS-97-12, Department of Computer Science,
146 *     Brown University.
147 *
148 * [4] Daniel B. Price, "New Techniques for Replay Debugging", Technical
149 *     Report CS-98-05, Department of Computer Science, Brown University.
150 */

```

```

152 #include <mdb/mdb_target_impl.h>
153 #include <mdb/mdb_debug.h>
154 #include <mdb/mdb_modapi.h>
155 #include <mdb/mdb_err.h>
156 #include <mdb/mdb_callb.h>
157 #include <mdb/mdb_gelf.h>
158 #include <mdb/mdb_io_impl.h>
159 #include <mdb/mdb_string.h>
160 #include <mdb/mdb_signal.h>
161 #include <mdb/mdb_frame.h>
162 #include <mdb/mdb.h>

```

```

164 #include <sys/stat.h>
165 #include <sys/param.h>
166 #include <sys/signal.h>
167 #include <strings.h>
168 #include <stdlib.h>
169 #include <errno.h>

```

```

171 /*
172 * Define convenience macros for referencing the set of vespec flag bits that
173 * are preserved by the target implementation, and the set of bits that
174 * determine automatic ve_hits == ve_limit behavior.
175 */
176 #define T_IMPL_BITS \
177     (MDB_TGT_SPEC_INTERNAL | MDB_TGT_SPEC_SILENT | MDB_TGT_SPEC_MATCHED | \
178     MDB_TGT_SPEC_DELETED)

```

```

180 #define T_AUTO_BITS \
181     (MDB_TGT_SPEC_AUTOSTOP | MDB_TGT_SPEC_AUTODEL | MDB_TGT_SPEC_AUTODIS)

```

```

183 /*
184 * Define convenience macro for referencing target flag pending continue bits.
185 */

```

```

186 #define T_CONT_BITS \
187     (MDB_TGT_F_STEP | MDB_TGT_F_STEP_OUT | MDB_TGT_F_NEXT | MDB_TGT_F_CONT)
188 #define T_F_STEP \
189     (MDB_TGT_F_STEP | MDB_TGT_F_STEP_OUT | MDB_TGT_F_STEP_BRANCH | \
190     MDB_TGT_F_NEXT | MDB_TGT_F_CONT)

```

```

189 mdb_tgt_t *
190 mdb_tgt_create(mdb_tgt_ctor_f *ctor, int flags, int argc, const char *argv[])
191 {
192     mdb_module_t *mp;
193     mdb_tgt_t *t;
194
195     if (flags & ~MDB_TGT_F_ALL) {
196         (void) set_errno(EINVAL);
197         return (NULL);
198     }
199
200     t = mdb_zalloc(sizeof (mdb_tgt_t), UM_SLEEP);
201     mdb_list_append(&mdb.m_tgtlist, t);
202
203     t->t_module = &mdb.m_rmod;
204     t->t_matched = T_SE_END;
205     t->t_flags = flags;
206     t->t_vepos = 1;
207     t->t_veneg = 1;
208
209     for (mp = mdb.m_mhead; mp != NULL; mp = mp->mod_next) {
210         if (ctor == mp->mod_tgt_ctor) {
211             t->t_module = mp;
212             break;
213         }
214     }
215
216     if (ctor(t, argc, argv) != 0) {
217         mdb_list_delete(&mdb.m_tgtlist, t);
218         mdb_free(t, sizeof (mdb_tgt_t));
219         return (NULL);
220     }
221
222     mdb_dprintf(MDB_DBG_TGT, "t_create %s (%p)\n",
223               t->t_module->mod_name, (void *)t);
224
225     (void) t->t_ops->t_status(t, &t->t_status);
226     return (t);
227 }

```

unchanged_portion_omitted

```

1047 /*
1048 * This function provides the low-level target continue algorithm. We proceed
1049 * in three phases: (1) we arm the active sespecs, except the specs matched at
1050 * the time we last stopped, (2) we call se_cont() on any matched sespecs to
1051 * step over these event transitions, and then arm the corresponding sespecs,
1052 * and (3) we call the appropriate low-level continue routine. Once the
1053 * target stops again, we determine which sespecs were matched, and invoke the
1054 * appropriate vespec callbacks and perform other vespec maintenance.
1055 */
1056 static int
1057 tgt_continue(mdb_tgt_t *t, mdb_tgt_status_t *tsp,
1058             int (*t_cont)(mdb_tgt_t *, mdb_tgt_status_t *))
1059 {
1060     mdb_var_t *hitv = mdb_nv_lookup(&mdb.m_nv, "hits");
1061     uintptr_t pc = t->t_status.st_pc;
1062     int error = 0;
1063
1064     mdb_sespec_t *sep, *nsep, *matched;
1065     mdb-vespec_t *vsep, *nvsep;
1066     uintptr_t addr;
1067
1068     uint_t cbits = 0; /* union of pending continue bits */
1069     uint_t ncont = 0; /* # of callbacks that requested cont */
1070     uint_t n = 0; /* # of callbacks */

```

```

1072 /*
1073  * If the target is undead, dead, or lost, we no longer allow continue.
1074  * This effectively forces the user to use ::kill or ::run after death.
1075  */
1076 if (t->t_status.st_state == MDB_TGT_UNDEAD)
1077     return (set_errno(EMDB_TGTZOMB));
1078 if (t->t_status.st_state == MDB_TGT_DEAD)
1079     return (set_errno(EMDB_TGTCORE));
1080 if (t->t_status.st_state == MDB_TGT_LOST)
1081     return (set_errno(EMDB_TGTLOST));
1082
1083 /*
1084  * If any of single-step, step-over, or step-out is pending, it takes
1085  * precedence over an explicit or pending continue, because these are
1086  * all different specialized forms of continue.
1087  */
1088 if (t->t_flags & MDB_TGT_F_STEP)
1089     t_cont = t->t_ops->t_step;
1090 else if (t->t_flags & MDB_TGT_F_NEXT)
1091     t_cont = t->t_ops->t_step;
1092 else if (t->t_flags & MDB_TGT_F_STEP_BRANCH)
1093     t_cont = t->t_ops->t_cont;
1094 else if (t->t_flags & MDB_TGT_F_STEP_OUT)
1095     t_cont = t->t_ops->t_cont;
1096
1097 /*
1098  * To handle step-over, we ask the target to find the address past the
1099  * next control transfer instruction. If an address is found, we plant
1100  * a temporary breakpoint there and continue; otherwise just step.
1101  */
1102 if ((t->t_flags & MDB_TGT_F_NEXT) && !(t->t_flags & MDB_TGT_F_STEP)) {
1103     if (t->t_ops->t_next(t, &addr) == -1 || mdb_tgt_add_vbrkpt(t,
1104         addr, MDB_TGT_SPEC_HIDDEN | MDB_TGT_SPEC_TEMPORARY,
1105         no_se_f, NULL) == 0) {
1106         mdb_dprintf(MDB_DBG_TGT, "next falling back to step: "
1107             "%s\n", mdb_strerror(errno));
1108     } else
1109         t_cont = t->t_ops->t_cont;
1110 }
1111
1112 /*
1113  * To handle step-out, we ask the target to find the return address of
1114  * the current frame, plant a temporary breakpoint there, and continue.
1115  */
1116 if (t->t_flags & MDB_TGT_F_STEP_OUT) {
1117     if (t->t_ops->t_step_out(t, &addr) == -1)
1118         return (-1); /* errno is set for us */
1119
1120     if (mdb_tgt_add_vbrkpt(t, addr, MDB_TGT_SPEC_HIDDEN |
1121         MDB_TGT_SPEC_TEMPORARY, no_se_f, NULL) == 0)
1122         return (-1); /* errno is set for us */
1123 }
1124
1125 /*
1126  * To handle step-branch, we ask the target to enable it for the coming
1127  * continue. Step-branch is incompatible with step, so don't enable it
1128  * if we're going to be stepping.
1129  */
1130 if (t->t_flags & MDB_TGT_F_STEP_BRANCH && t_cont == t->t_ops->t_cont) {
1131     if (t->t_ops->t_step_branch(t) == -1)
1132         return (-1); /* errno is set for us */
1133 }
1134
1135 (void) mdb_signal_block(SIGHUP);
1136 (void) mdb_signal_block(SIGTERM);

```

```

1125     mdb_intr_disable();
1126
1127     t->t_flags &= ~T_CONT_BITS;
1128     t->t_flags |= MDB_TGT_F_BUSY;
1129     mdb_tgt_sespec_arm_all(t);
1130
1131     ASSERT(t->t_matched != NULL);
1132     matched = t->t_matched;
1133     t->t_matched = T_SE_END;
1134
1135     if (mdb.m_term != NULL)
1136         IOP_SUSPEND(mdb.m_term);
1137
1138     /*
1139      * Iterate over the matched sespec list, performing autostop processing
1140      * and clearing the matched bit for each associated vespec. We then
1141      * invoke each sespec's se_cont callback in order to continue past
1142      * the corresponding event. If the matched list has more than one
1143      * sespec, we assume that the se_cont callbacks are non-interfering.
1144      */
1145     for (sep = matched; sep != T_SE_END; sep = sep->se_matched) {
1146         for (vep = mdb_list_next(&sep->se_elist; vep != NULL; ) {
1147             if ((vep->ve_flags & MDB_TGT_SPEC_AUTOSTOP) &&
1148                 (vep->ve_limit && vep->ve_hits == vep->ve_limit))
1149                 vep->ve_hits = 0;
1150
1151                 vep->ve_flags &= ~MDB_TGT_SPEC_MATCHED;
1152                 vep = mdb_list_next(vep);
1153         }
1154
1155         if (sep->se_ops->se_cont(t, sep, &t->t_status) == -1) {
1156             error = errno ? errno : -1;
1157             tgt_disarm_sespecs(t);
1158             break;
1159         }
1160
1161         if (!(t->t_status.st_flags & MDB_TGT_ISTOP)) {
1162             tgt_disarm_sespecs(t);
1163             if (t->t_status.st_state == MDB_TGT_UNDEAD)
1164                 mdb_tgt_sespec_idle_all(t, EMDB_TGTZOMB, TRUE);
1165             else if (t->t_status.st_state == MDB_TGT_LOST)
1166                 mdb_tgt_sespec_idle_all(t, EMDB_TGTLOST, TRUE);
1167             break;
1168         }
1169     }
1170
1171     /*
1172      * Clear the se_matched field for each matched sespec, and drop the
1173      * reference count since the sespec is no longer on the matched list.
1174      */
1175     for (sep = matched; sep != T_SE_END; sep = nsep) {
1176         nsep = sep->se_matched;
1177         sep->se_matched = NULL;
1178         mdb_tgt_sespec_rele(t, sep);
1179     }
1180
1181     /*
1182      * If the matched list was non-empty, see if we hit another event while
1183      * performing se_cont() processing. If so, don't bother continuing any
1184      * further. If not, arm the sespecs on the old matched list by calling
1185      * mdb_tgt_sespec_arm_all() again and then continue by calling t_cont.
1186      */
1187     if (matched != T_SE_END) {
1188         if (error != 0 || !(t->t_status.st_flags & MDB_TGT_ISTOP))
1189             goto out; /* abort now if se_cont() failed */

```

```

1191         if ((t->t_matched = tgt_match_sespecs(t, FALSE)) != T_SE_END) {
1192             tgt_disarm_sespecs(t);
1193             goto out;
1194         }
1196         mdb_tgt_sespec_arm_all(t);
1197     }
1199     if (t_cont != t->t_ops->t_step || pc == t->t_status.st_pc) {
1200         if (t_cont(t, &t->t_status) != 0)
1201             error = errno ? errno : -1;
1202     }
1204     tgt_disarm_sespecs(t);
1206     if (t->t_flags & MDB_TGT_F_UNLOAD)
1207         longjmp(mdb.m_frame->f_pcb, MDB_ERR_QUIT);
1209     if (t->t_status.st_state == MDB_TGT_UNDEAD)
1210         mdb_tgt_sespec_idle_all(t, EMDB_TGTZOMB, TRUE);
1211     else if (t->t_status.st_state == MDB_TGT_LOST)
1212         mdb_tgt_sespec_idle_all(t, EMDB_TGTLOST, TRUE);
1213     else if (t->t_status.st_flags & MDB_TGT_ISTOP)
1214         t->t_matched = tgt_match_sespecs(t, TRUE);
1215 out:
1216     if (mdb.m_term != NULL)
1217         IOP_RESUME(mdb.m_term);
1219     (void) mdb_signal_unblock(SIGTERM);
1220     (void) mdb_signal_unblock(SIGHUP);
1221     mdb_intr_enable();
1223     for (sep = t->t_matched; sep != T_SE_END; sep = sep->se_matched) {
1224         /*
1225          * When we invoke a ve_callback, it may in turn request that the
1226          * target continue immediately after callback processing is
1227          * complete. We only allow this to occur if *all* callbacks
1228          * agree to continue. To implement this behavior, we keep a
1229          * count (ncont) of such requests, and only apply the cumulative
1230          * continue bits (cbits) to the target if ncont is equal to the
1231          * total number of callbacks that are invoked (n).
1232          */
1233         for (vep = mdb_list_next(&sep->se_elist);
1234              vep != NULL; vep = nvep, n++) {
1235             /*
1236              * Place an extra hold on the current vespec and pick
1237              * up the next pointer before invoking the callback: we
1238              * must be prepared for the vespec to be deleted or
1239              * moved to a different list by the callback.
1240              */
1241             mdb_tgt-vespec_hold(t, vep);
1242             nvep = mdb_list_next(vep);
1244             vep->ve_flags |= MDB_TGT_SPEC_MATCHED;
1245             vep->ve_hits++;
1247             mdb_nv_set_value(mdb.m_dot, t->t_status.st_pc);
1248             mdb_nv_set_value(hitv, vep->ve_hits);
1250             ASSERT((t->t_flags & T_CONT_BITS) == 0);
1251             vep->ve_callback(t, vep->ve_id, vep->ve_data);
1253             ncont += (t->t_flags & T_CONT_BITS) != 0;
1254             cbits |= (t->t_flags & T_CONT_BITS);
1255             t->t_flags &= ~T_CONT_BITS;

```

```

1257         if (vep->ve_limit && vep->ve_hits == vep->ve_limit) {
1258             if (vep->ve_flags & MDB_TGT_SPEC_AUTODEL)
1259                 (void) mdb_tgt-vespec_delete(t,
1260                     vep->ve_id);
1261             else if (vep->ve_flags & MDB_TGT_SPEC_AUTODIS)
1262                 (void) mdb_tgt-vespec_disable(t,
1263                     vep->ve_id);
1264         }
1266         if (vep->ve_limit && vep->ve_hits < vep->ve_limit) {
1267             if (vep->ve_flags & MDB_TGT_SPEC_AUTOSTOP)
1268                 (void) mdb_tgt_continue(t, NULL);
1269         }
1271         mdb_tgt-vespec_rele(t, vep);
1272     }
1273 }
1275     if (t->t_matched != T_SE_END && ncont == n)
1276         t->t_flags |= cbits; /* apply continues (see above) */
1278     mdb_tgt_sespec_prune_all(t);
1280     t->t_status.st_flags &= ~MDB_TGT_BUSY;
1281     t->t_flags &= ~MDB_TGT_F_BUSY;
1283     if (tsp != NULL)
1284         bcopy(&t->t_status, tsp, sizeof (mdb_tgt_status_t));
1286     if (error != 0)
1287         return (set_errno(error));
1289     return (0);
1290 }
1291 unchanged portion omitted
1398 int
1399 mdb_tgt_step_branch(mdb_tgt_t *t, mdb_tgt_status_t *tsp)
1400 {
1401     t->t_flags |= MDB_TGT_F_STEP_BRANCH; /* set flag even if tgt not busy */
1402     return (tgt_request_continue(t, tsp, 0, t->t_ops->t_cont));
1403 }
1387 int
1388 mdb_tgt_next(mdb_tgt_t *t, mdb_tgt_status_t *tsp)
1389 {
1390     t->t_flags |= MDB_TGT_F_NEXT; /* set flag even if tgt not busy */
1391     return (tgt_request_continue(t, tsp, 0, t->t_ops->t_step));
1392 }
1393 unchanged portion omitted

```

new/usr/src/cmd/mdb/common/mdb/mdb_target.h

1

```
*****
22252 Fri Apr 6 17:24:55 2018
new/usr/src/cmd/mdb/common/mdb/mdb_target.h
9210 remove KMDB branch debugging support
9211 ::crregs could do with cr2/cr3 support
9209 ::ttrace should be able to filter by thread
Reviewed by: Patrick Mooney <patrick.mooney@joyent.com>
Reviewed by: Yuri Pankov <yuripv@yuripv.net>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License, Version 1.0 only
6  * (the "License"). You may not use this file except in compliance
7  * with the License.
8  *
9  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
10 * or http://www.opensolaris.org/os/licensing.
11 * See the License for the specific language governing permissions
12 * and limitations under the License.
13 *
14 * When distributing Covered Code, include this CDDL HEADER in each
15 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
16 * If applicable, add the following below this CDDL HEADER, with the
17 * fields enclosed by brackets "[]" replaced with your own identifying
18 * information: Portions Copyright [yyyy] [name of copyright owner]
19 *
20 * CDDL HEADER END
21 */
22 /*
23 * Copyright 2005 Sun Microsystems, Inc. All rights reserved.
24 * Use is subject to license terms.
25 *
26 * Copyright 2018 Joyent, Inc.
27 */

29 #ifndef _MDB_TARGET_H
30 #define _MDB_TARGET_H

32 #include <sys/utsname.h>
33 #include <sys/types.h>
34 #include <gelf.h>

36 #ifdef __cplusplus
37 extern "C" {
38 #endif

40 /*
41  * Forward declaration of the target structure: the target itself is defined in
42  * mdb_tgt_impl.h and is opaque with respect to callers of this interface.
43  */

45 struct mdb_tgt;
46 struct mdb_arg;
47 struct ctf_file;

49 typedef struct mdb_tgt mdb_tgt_t;

51 extern void mdb_create_builtin_tgts(void);
52 extern void mdb_create_loadable_disasms(void);

54 /*
55  * Target Constructors
56  *
57  * These functions are used to create a complete debugger target. The
```

new/usr/src/cmd/mdb/common/mdb/mdb_target.h

2

```
58 * constructor is passed as an argument to mdb_tgt_create().
59 */

61 extern int mdb_value_tgt_create(mdb_tgt_t *, int, const char *[]);
62 #ifndef _KMDB
63 extern int mdb_kvm_tgt_create(mdb_tgt_t *, int, const char *[]);
64 extern int mdb_proc_tgt_create(mdb_tgt_t *, int, const char *[]);
65 extern int mdb_kproc_tgt_create(mdb_tgt_t *, int, const char *[]);
66 extern int mdb_rawfile_tgt_create(mdb_tgt_t *, int, const char *[]);
67 #else
68 extern int kmdb_kvm_create(mdb_tgt_t *, int, const char *[]);
69 #endif

71 /*
72  * Targets are created by calling mdb_tgt_create() with an optional set of
73  * target flags, an argument list, and a target constructor (see above):
74  */

76 #define MDB_TGT_F_RDWR          0x0001 /* Open for writing (else read-only) */
77 #define MDB_TGT_F_ALLOWIO      0x0002 /* Allow I/O mem access (live only) */
78 #define MDB_TGT_F_FORCE        0x0004 /* Force open (even if non-exclusive) */
79 #define MDB_TGT_F_PRELOAD      0x0008 /* Preload all symbol tables */
80 #define MDB_TGT_F_NOLOAD       0x0010 /* Do not do load-object processing */
81 #define MDB_TGT_F_NOSTOP       0x0020 /* Do not stop target on attach */
82 #define MDB_TGT_F_STEP         0x0040 /* Single-step is pending */
83 #define MDB_TGT_F_STEP_OUT     0x0080 /* Step-out is pending */
84 #define MDB_TGT_F_NEXT         0x0100 /* Step-over is pending */
85 #define MDB_TGT_F_CONT         0x0200 /* Continue is pending */
86 #define MDB_TGT_F_BUSY         0x0400 /* Target is busy executing */
87 #define MDB_TGT_F_ASIO         0x0800 /* Use t_ared and t_awrite for i/o */
88 #define MDB_TGT_F_UNLOAD       0x1000 /* Unload has been requested */
89 #define MDB_TGT_F_ALL          0x1fff /* Mask of all valid flags */
90 #define MDB_TGT_F_STEP_BRANCH  0x0100 /* Step-branch is pending */
91 #define MDB_TGT_F_NEXT         0x0200 /* Step-over is pending */
92 #define MDB_TGT_F_CONT         0x0400 /* Continue is pending */
93 #define MDB_TGT_F_BUSY         0x0800 /* Target is busy executing */
94 #define MDB_TGT_F_ASIO         0x1000 /* Use t_ared and t_awrite for i/o */
95 #define MDB_TGT_F_UNLOAD       0x2000 /* Unload has been requested */
96 #define MDB_TGT_F_ALL          0x3fff /* Mask of all valid flags */

91 typedef int mdb_tgt_ctor_f(mdb_tgt_t *, int, const char *[]);

93 extern mdb_tgt_t *mdb_tgt_create(mdb_tgt_ctor_f *, int, int, const char *[]);
94 extern void mdb_tgt_destroy(mdb_tgt_t *);

96 extern int mdb_tgt_getflags(mdb_tgt_t *);
97 extern int mdb_tgt_setflags(mdb_tgt_t *, int);
98 extern int mdb_tgt_setcontext(mdb_tgt_t *, void *);

100 /*
101  * Targets are activated and de-activated by the debugger framework. An
102  * activation occurs after construction when the target becomes the current
103  * target in the debugger. A target is de-activated prior to its destructor
104  * being called by mdb_tgt_destroy, or when another target is activated.
105  * These callbacks are suitable for loading support modules and other tasks.
106  */
107 extern void mdb_tgt_activate(mdb_tgt_t *);

109 /*
110  * Prior to issuing a new command prompt, the debugger framework calls the
111  * target's periodic callback to allow it to load new modules or perform
112  * other background tasks.
113  */
114 extern void mdb_tgt_periodic(mdb_tgt_t *);

116 /*
```



```
386 #define MDB_TGT_SPEC_AUTODIS 0x0080 /* Disable when match limit reached */
387 #define MDB_TGT_SPEC_AUTOSTOP 0x0100 /* Stop when match limit reached */
388 #define MDB_TGT_SPEC_STICKY 0x0200 /* Do not delete as part of :z */

390 #define MDB_TGT_SPEC_HIDDEN (MDB_TGT_SPEC_INTERNAL | MDB_TGT_SPEC_SILENT)

392 typedef struct mdb_tgt_spec_desc {
393     int spec_id; /* Event specifier id (VID) */
394     uint_t spec_flags; /* Flags (see above) */
395     uint_t spec_hits; /* Count of number of times matched */
396     uint_t spec_limit; /* Limit on number of times matched */
397     int spec_state; /* State (see above) */
398     int spec_errno; /* Last error code (if IDLE or ERROR) */
399     uintptr_t spec_base; /* Start of affected memory region */
400     size_t spec_size; /* Size of affected memory region */
401     void *spec_data; /* Callback private data */
402 } mdb_tgt_spec_desc_t;
unchanged portion omitted
```

new/usr/src/cmd/mdb/common/mdb/mdb_target_impl.h

1

```
*****
14888 Fri Apr 6 17:24:55 2018
new/usr/src/cmd/mdb/common/mdb/mdb_target_impl.h
9210 remove KMDB branch debugging support
9211 ::crregs could do with cr2/cr3 support
9209 ::ttrace should be able to filter by thread
Reviewed by: Patrick Mooney <patrick.mooney@joyent.com>
Reviewed by: Yuri Pankov <yuripv@yuripv.net>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22  * Copyright 2006 Sun Microsystems, Inc. All rights reserved.
23  * Use is subject to license terms.
24 */
25 /*
26  * Copyright (c) 2018, Joyent, Inc. All rights reserved.
26  * Copyright (c) 2012, Joyent, Inc. All rights reserved.
27 */

29 #ifndef _MDB_TARGET_IMPL_H
30 #define _MDB_TARGET_IMPL_H

32 #include <mdb/mdb_target.h>
33 #include <mdb/mdb_module.h>
34 #include <mdb/mdb_list.h>
35 #include <mdb/mdb_gelf.h>
36 #include <sys/auxv.h>

38 #ifdef __cplusplus
39 extern "C" {
40 #endif

42 #ifdef _MDB

44 /*
45  * Target Operations
46  *
47  * This ops vector implements the set of primitives which can be used by the
48  * debugger to interact with the target, and encompasses most of the calls
49  * found in <mdb/mdb_target.h>. The remainder of the target interface is
50  * implemented by common code that invokes these primitives or manipulates
51  * the common target structures directly.
52 */

54 typedef struct mdb_tgt_ops {
55     int (*t_setflags)(mdb_tgt_t *, int);
56     int (*t_setcontext)(mdb_tgt_t *, void *);
```

new/usr/src/cmd/mdb/common/mdb/mdb_target_impl.h

2

```
58     void (*t_activate)(mdb_tgt_t *);
59     void (*t_deactivate)(mdb_tgt_t *);
60     void (*t_periodic)(mdb_tgt_t *);
61     void (*t_destroy)(mdb_tgt_t *);

63     const char *(*t_name)(mdb_tgt_t *);
64     const char *(*t_isa)(mdb_tgt_t *);
65     const char *(*t_platform)(mdb_tgt_t *);
66     int (*t_uname)(mdb_tgt_t *, struct utsname *);
67     int (*t_dmodel)(mdb_tgt_t *);

69     ssize_t (*t_aread)(mdb_tgt_t *,
70         mdb_tgt_as_t, void *, size_t, mdb_tgt_addr_t);

72     ssize_t (*t_awrite)(mdb_tgt_t *,
73         mdb_tgt_as_t, const void *, size_t, mdb_tgt_addr_t);

75     ssize_t (*t_vread)(mdb_tgt_t *, void *, size_t, uintptr_t);
76     ssize_t (*t_vwrite)(mdb_tgt_t *, const void *, size_t, uintptr_t);
77     ssize_t (*t_pread)(mdb_tgt_t *, void *, size_t, physaddr_t);
78     ssize_t (*t_pwrite)(mdb_tgt_t *, const void *, size_t, physaddr_t);
79     ssize_t (*t_fread)(mdb_tgt_t *, void *, size_t, uintptr_t);
80     ssize_t (*t_fwrite)(mdb_tgt_t *, const void *, size_t, uintptr_t);
81     ssize_t (*t_ioread)(mdb_tgt_t *, void *, size_t, uintptr_t);
82     ssize_t (*t_iowrite)(mdb_tgt_t *, const void *, size_t, uintptr_t);

84     int (*t_vtop)(mdb_tgt_t *, mdb_tgt_as_t, uintptr_t, physaddr_t);

86     int (*t_lookup_by_name)(mdb_tgt_t *,
87         const char *, const char *, GElf_Sym *, mdb_syminfo_t);

89     int (*t_lookup_by_addr)(mdb_tgt_t *,
90         uintptr_t, uint_t, char *, size_t, GElf_Sym *, mdb_syminfo_t);

92     int (*t_symbol_iter)(mdb_tgt_t *,
93         const char *, uint_t, uint_t, mdb_tgt_sym_f *, void *);

95     int (*t_mapping_iter)(mdb_tgt_t *, mdb_tgt_map_f *, void *);
96     int (*t_object_iter)(mdb_tgt_t *, mdb_tgt_map_f *, void *);

98     const mdb_map_t *(*t_addr_to_map)(mdb_tgt_t *, uintptr_t);
99     const mdb_map_t *(*t_name_to_map)(mdb_tgt_t *, const char *);
100    struct ctf_file *(*t_addr_to_ctf)(mdb_tgt_t *, uintptr_t);
101    struct ctf_file *(*t_name_to_ctf)(mdb_tgt_t *, const char *);

103    int (*t_status)(mdb_tgt_t *, mdb_tgt_status_t *);
104    int (*t_run)(mdb_tgt_t *, int, const struct mdb_arg *);
105    int (*t_step)(mdb_tgt_t *, mdb_tgt_status_t *);
106    int (*t_step_out)(mdb_tgt_t *, uintptr_t *);
107    int (*t_step_branch)(mdb_tgt_t *);
107    int (*t_next)(mdb_tgt_t *, uintptr_t *);
108    int (*t_cont)(mdb_tgt_t *, mdb_tgt_status_t *);
109    int (*t_signal)(mdb_tgt_t *, int);

111    int (*t_add_vbrkpt)(mdb_tgt_t *, uintptr_t,
112        int, mdb_tgt_se_f *, void *);
113    int (*t_add_sbrkpt)(mdb_tgt_t *, const char *,
114        int, mdb_tgt_se_f *, void *);

116    int (*t_add_pwapt)(mdb_tgt_t *, physaddr_t, size_t, uint_t,
117        int, mdb_tgt_se_f *, void *);
118    int (*t_add_vwapt)(mdb_tgt_t *, uintptr_t, size_t, uint_t,
119        int, mdb_tgt_se_f *, void *);
120    int (*t_add_iowapt)(mdb_tgt_t *, uintptr_t, size_t, uint_t,
121        int, mdb_tgt_se_f *, void *);
```

```
123     int (*t_add_sysenter)(mdb_tgt_t *, int, int, mdb_tgt_se_f *, void *);
124     int (*t_add_sysexit)(mdb_tgt_t *, int, int, mdb_tgt_se_f *, void *);
125     int (*t_add_signal)(mdb_tgt_t *, int, int, mdb_tgt_se_f *, void *);
126     int (*t_add_fault)(mdb_tgt_t *, int, int, mdb_tgt_se_f *, void *);

128     int (*t_getareg)(mdb_tgt_t *, mdb_tgt_tid_t, const char *,
129                     mdb_tgt_reg_t *);
130     int (*t_putareg)(mdb_tgt_t *, mdb_tgt_tid_t, const char *,
131                     mdb_tgt_reg_t);

133     int (*t_stack_iter)(mdb_tgt_t *, const mdb_tgt_gregset_t *,
134                         mdb_tgt_stack_f *, void *);

136     int (*t_auxv)(mdb_tgt_t *, const auxv_t **auxvp);
137 } mdb_tgt_ops_t;
_____ unchanged portion omitted
```



```

*****
6175 Fri Apr 6 17:24:55 2018
new/usr/src/cmd/mdb/common/mdb/mdb_value.c
9210 remove KMDB branch debugging support
9211 ::crregs could do with cr2/cr3 support
9209 ::ttrace should be able to filter by thread
Reviewed by: Patrick Mooney <patrick.mooney@joyent.com>
Reviewed by: Yuri Pankov <yuripv@yuripv.net>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22  * Copyright 2008 Sun Microsystems, Inc. All rights reserved.
23  * Use is subject to license terms.
24  *
25  * Copyright 2018 Joyent, Inc.
26  */
27
28 #pragma ident "%Z%M% %I% %E% SMI"
29
30 /*
31  * Immediate Value Target
32  *
33  * The immediate value target is used when the '=' verb is used to
34  * format an immediate value, or with ::print -i. The target is
35  * initialized with a specific value, and then simply copies bytes from
36  * this integer in its read routine. Two notes:
37  *
38  * (1) the address parameter of value_read is treated as an offset into
39  * the immediate value.
40  *
41  * (2) on big-endian systems, we need to be careful about the place we
42  * copy data from. If the caller specified a typesize in the argv array
43  * we use that for offsetting, otherwise we use the read size.
44  * If the user didn't specify the typesize, then 'addr' is ignored,
45  * and all reads are at an offset of 0 into the immediate value. This
46  * covers both the usage of ::print -i, and the semantics of adb
47  * commands like "0x1234=X", which should produce 0x1234 as a result;
48  * the adb model is for it to act like a cast down to the smaller
49  * integer type; this is handled as mentioned.
50  */
51
52 #include <mdb/mdb_target_impl.h>
53 #include <mdb/mdb_types.h>
54 #include <mdb/mdb_conf.h>
55 #include <mdb/mdb_err.h>
56
57 #include <sys/isa_defs.h>

```

```

56 #include <strings.h>
57
58 void mdb_value_tgt_destroy(mdb_tgt_t *);
59
60 typedef struct mdb_value_data {
61     uintmax_t mvd_data;
62     size_t mvd_typesize;
63 } mdb_value_data_t;
64
65 unchanged_portion_omitted
66
67 static const mdb_tgt_ops_t value_ops = {
68     (int (*)()) mdb_tgt_notsup, /* t_setflags */
69     (int (*)()) mdb_tgt_notsup, /* t_setcontext */
70     (void (*)()) mdb_tgt_nop, /* t_activate */
71     (void (*)()) mdb_tgt_nop, /* t_deactivate */
72     (void (*)()) mdb_tgt_nop, /* t_periodic */
73     mdb_value_tgt_destroy, /* t_destroy */
74     (const char *(*()) mdb_tgt_null, /* t_name */
75     (const char *(*()) mdb_conf_isa, /* t_isa */
76     (const char *(*()) mdb_conf_platform, /* t_platform */
77     (int (*)()) mdb_tgt_notsup, /* t_uname */
78     (int (*)()) mdb_tgt_notsup, /* t_dmodel */
79     (ssize_t (*)()) mdb_tgt_notsup, /* t_aread */
80     (ssize_t (*)()) mdb_tgt_notsup, /* t_awrite */
81     value_read, /* t_vread */
82     value_write, /* t_vwrite */
83     (ssize_t (*)()) mdb_tgt_notsup, /* t_pread */
84     (ssize_t (*)()) mdb_tgt_notsup, /* t_pwrite */
85     value_read, /* t_fread */
86     value_write, /* t_fwrite */
87     value_read, /* t_loread */
88     value_write, /* t_iowrite */
89     (int (*)()) mdb_tgt_notsup, /* t_vtop */
90     (int (*)()) mdb_tgt_notsup, /* t_lookup_by_name */
91     (int (*)()) mdb_tgt_notsup, /* t_lookup_by_addr */
92     (int (*)()) mdb_tgt_notsup, /* t_symbol_iter */
93     (int (*)()) mdb_tgt_notsup, /* t_mapping_iter */
94     (int (*)()) mdb_tgt_notsup, /* t_object_iter */
95     (const mdb_map_t *(*()) mdb_tgt_null, /* t_addr_to_map */
96     (const mdb_map_t *(*()) mdb_tgt_null, /* t_name_to_map */
97     (struct ctf_file *(*()) mdb_tgt_null, /* t_addr_to_ctf */
98     (struct ctf_file *(*()) mdb_tgt_null, /* t_name_to_ctf */
99     (int (*)()) mdb_tgt_notsup, /* t_status */
100    (int (*)()) mdb_tgt_notsup, /* t_run */
101    (int (*)()) mdb_tgt_notsup, /* t_step */
102    (int (*)()) mdb_tgt_notsup, /* t_step_out */
103    (int (*)()) mdb_tgt_notsup, /* t_step_branch */
104    (int (*)()) mdb_tgt_notsup, /* t_next */
105    (int (*)()) mdb_tgt_notsup, /* t_cont */
106    (int (*)()) mdb_tgt_notsup, /* t_signal */
107    (int (*)()) mdb_tgt_null, /* t_add_vbrkpt */
108    (int (*)()) mdb_tgt_null, /* t_add_sbrkpt */
109    (int (*)()) mdb_tgt_null, /* t_add_pwapt */
110    (int (*)()) mdb_tgt_null, /* t_add_vwapt */
111    (int (*)()) mdb_tgt_null, /* t_add_iowapt */
112    (int (*)()) mdb_tgt_null, /* t_add_sysenter */
113    (int (*)()) mdb_tgt_null, /* t_add_sysexit */
114    (int (*)()) mdb_tgt_null, /* t_add_signal */
115    (int (*)()) mdb_tgt_null, /* t_add_fault */
116    (int (*)()) mdb_tgt_notsup, /* t_getareg */
117    (int (*)()) mdb_tgt_notsup, /* t_putareg */
118    (int (*)()) mdb_tgt_nop, /* t_stack_iter */
119    (int (*)()) mdb_tgt_notsup, /* t_auxv */
120 };
121
122 unchanged_portion_omitted

```

```

*****
23287 Fri Apr 6 17:24:55 2018
new/usr/src/cmd/mdb/i86pc/modules/unix/i86mmu.c
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22 * Copyright 2009 Sun Microsystems, Inc. All rights reserved.
23 * Use is subject to license terms.
24 *
25 * Copyright 2018 Joyent, Inc.
26 */

28 /*
29 * This part of the file contains the mdb support for dcmds:
30 *      ::memseg_list
31 * and walkers for:
32 *      memseg - a memseg list walker for ::memseg_list
33 *
34 */

36 #include <sys/types.h>
37 #include <sys/machparam.h>
38 #include <sys/controlregs.h>
39 #include <sys/mach_mmu.h>
40 #ifdef __xpv
41 #include <sys/hypervisor.h>
42 #endif
43 #include <vm/as.h>

45 #include <mdb/mdb_modapi.h>
46 #include <mdb/mdb_target.h>

48 #include <vm/page.h>
49 #include <vm/hat_i86.h>

51 #define VA_SIGN_BIT (1UL << 47)
52 #define VA_SIGN_EXTEND(va) (((va) ^ VA_SIGN_BIT) - VA_SIGN_BIT)

54 struct pfn2pp {
55     pfn_t pfn;
56     page_t *pp;
57 };
    unchanged_portion_omitted

```

```

401 /*
402  * Print a PTE in more human friendly way. The PTE is assumed to be in
403  * a level 0 page table, unless -l specifies another level.
404  *
405  * The PTE value can be specified as the -p option, since on a 32 bit kernel
406  * with PAE running it's larger than a uintptr_t.
407  */
408 static int
409 do_pte_dcmd(int level, uint64_t pte)
410 {
411     static char *attr[] = {
412         "wrbck", "wrthru", "uncached", "uncached",
413         "wrbck", "wrthru", "wrcombine", "uncached"};
414     int pat_index = 0;
415     pfn_t mfn;
416
417     mdb_printf("pte=0x%llr: ", pte);
418     mdb_printf("pte=%llr: ", pte);
419     if (PTE_GET(pte, mmu.pt_nx))
420         mdb_printf("noexec ");
421
422     mfn = pte2mfn(pte, level);
423     mdb_printf("%s=0x%llr ", is_xpv ? "mfn" : "pfn", mfn);
424
425     if (PTE_GET(pte, mmu.pt_nx))
426         mdb_printf("noexec ");
427
428     if (PTE_GET(pte, PT_NOCONSIST))
429         mdb_printf("noconsist ");
430
431     if (PTE_GET(pte, PT_NOSYNC))
432         mdb_printf("nosync ");
433
434     if (PTE_GET(pte, mmu.pt_global))
435         mdb_printf("global ");
436
437     if (level > 0 && PTE_GET(pte, PT_PAGESIZE))
438         mdb_printf("largepage ");
439
440     if (level > 0 && PTE_GET(pte, PT_MOD))
441         mdb_printf("mod ");
442
443     if (level > 0 && PTE_GET(pte, PT_REF))
444         mdb_printf("ref ");
445
446     if (PTE_GET(pte, PT_USER))
447         mdb_printf("user ");
448
449     if (PTE_GET(pte, PT_WRITABLE))
450         mdb_printf("write ");
451
452     /*
453      * Report non-standard cacheability
454      */
455     pat_index = 0;
456     if (level > 0) {
457         if (PTE_GET(pte, PT_PAGESIZE) && PTE_GET(pte, PT_PAT_LARGE))
458             pat_index += 4;
459     } else {
460         if (PTE_GET(pte, PT_PAT_4K))
461             pat_index += 4;
462     }
463
464     if (PTE_GET(pte, PT_NOCACHE))
465         pat_index += 2;

```

```

462     if (PTE_GET(pte, PT_WRITETHRU))
463         pat_index += 1;

465     if (pat_index != 0)
466         mdb_printf("%s", attr[pat_index]);

468     if (PTE_GET(pte, PT_VALID) == 0)
469         mdb_printf(" !VALID ");

471     mdb_printf("\n");
472     return (DCMD_OK);
473 }

475 /*
476  * Print a PTE in more human friendly way. The PTE is assumed to be in
477  * a level 0 page table, unless -l specifies another level.
478  *
479  * The PTE value can be specified as the -p option, since on a 32 bit kernel
480  * with PAE running it's larger than a uintptr_t.
481  */
482 /*ARGSUSED*/
483 int
484 pte_dcmd(uintptr_t addr, uint_t flags, int argc, const mdb_arg_t *argv)
485 {
486     uint64_t level = 0;
487     int level = 0;
488     uint64_t pte = 0;
489     char *level_str = NULL;
490     char *pte_str = NULL;

485     init_mmu();

487     if (mmu.num_level == 0)
488         return (DCMD_ERR);

490     if ((flags & DCMD_ADDRSPEC) == 0)
491     if (mdb_getopts(argc, argv,
492         'p', MDB_OPT_STR, &pte_str,
493         'l', MDB_OPT_STR, &level_str) != argc)
494         return (DCMD_USAGE);

493     if (mdb_getopts(argc, argv,
494         'l', MDB_OPT_UINT64, &level) != argc)
495     /*
496      * parse the PTE to decode, if it's 0, we don't do anything
497      */
498     if (pte_str != NULL) {
499         pte = mdb_strtoull(pte_str);
500     } else {
501         if ((flags & DCMD_ADDRSPEC) == 0)
502             return (DCMD_USAGE);
503         pte = addr;
504     }
505     if (pte == 0)
506         return (DCMD_OK);

497     if (level > mmu.max_level) {
498         mdb_warn("invalid level %lu\n", level);
499     /*
500      * parse the level if supplied
501      */
502     if (level_str != NULL) {
503         level = mdb_strtoull(level_str);
504         if (level < 0 || level > mmu.max_level)
505             return (DCMD_ERR);
506     }
507 }

```

```

502     if (addr == 0)
503         return (DCMD_OK);

505     return (do_pte_dcmd((int)level, addr));
506     return (do_pte_dcmd(level, pte));
507 }
508 }
509 }
510 }
511 }
512 }
513 }
514 }
515 }
516 }
517 }
518 }
519 }
520 }
521 }
522 }
523 }
524 }
525 }
526 }
527 }
528 }
529 }
530 }
531 }
532 }
533 }
534 }
535 }
536 }
537 }
538 }
539 }
540 }
541 }
542 }
543 }
544 }
545 }
546 }
547 }
548 }
549 }
550 }
551 }
552 }
553 }
554 }
555 }
556 }
557 }
558 }
559 }
560 }
561 }
562 }
563 }
564 }
565 }
566 }
567 }
568 }
569 }
570 }
571 }
572 }
573 }
574 }
575 }
576 }
577 }
578 }
579 }
580 }
581 }
582 }
583 }
584 }
585 }
586 }
587 }
588 }
589 }
590 }
591 }
592 }
593 }
594 }
595 }
596 }
597 }
598 }
599 }
600 }
601 }
602 }
603 }
604 }
605 }
606 }
607 }
608 }
609 }
610 }
611 }
612 }
613 }
614 }
615 }
616 }
617 }
618 }
619 }
620 }
621 }
622 }
623 }
624 }
625 }
626 }
627 }
628 }
629 }
630 }
631 }
632 }
633 }
634 }
635 }
636 }
637 }
638 }
639 }
640 }
641 }
642 }
643 }
644 }
645 }
646 }
647 }
648 }
649 }
650 }
651 }
652 }
653 }
654 }
655 }
656 }
657 }
658 }
659 }
660 }
661 }
662 }
663 }
664 }
665 }
666 }
667 }
668 }
669 }
670 }
671 }
672 }
673 }
674 }
675 }
676 }
677 }
678 }
679 }
680 }
681 }
682 }
683 }
684 }
685 }
686 }
687 }
688 }
689 }
690 }
691 }
692 }
693 }
694 }
695 }
696 }
697 }
698 }
699 }
700 }
701 }
702 }
703 }
704 }
705 }
706 }
707 }
708 }
709 }
710 }
711 }
712 }
713 }
714 }
715 }
716 }
717 }
718 }
719 }
720 }
721 }
722 }
723 }
724 }
725 }
726 }
727 }
728 }
729 }
730 }
731 }
732 }
733 }
734 }
735 }
736 }
737 }
738 }
739 }
740 }
741 }
742 }
743 }
744 }
745 }
746 }
747 }
748 }
749 }
750 }
751 }
752 }
753 }
754 }
755 }
756 }
757 }
758 }
759 }
760 }
761 }
762 }
763 }
764 }
765 }
766 }
767 }
768 }
769 }
770 }
771 }
772 }
773 }
774 }
775 }
776 }
777 }
778 }
779 }
780 }
781 }
782 }
783 }
784 }
785 }
786 }
787 }
788 }
789 }
790 }
791 }
792 }
793 }
794 }
795 }
796 }
797 }
798 }
799 }
800 }
801 }
802 }
803 }
804 }
805 }
806 }
807 }
808 }
809 }
810 }
811 }
812 }
813 }
814 }
815 }
816 }
817 }
818 }
819 }
820 }
821 }
822 }
823 }
824 }
825 }
826 }
827 }
828 }
829 }
830 }
831 }
832 }
833 }
834 }
835 }
836 }
837 }
838 }
839 }
840 }
841 }
842 }
843 }
844 }
845 }
846 }
847 }
848 }
849 }
850 }
851 }
852 }
853 }
854 }
855 }
856 }
857 }
858 }
859 }
860 }
861 }
862 }
863 }
864 }
865 }
866 }
867 }
868 }
869 }
870 }
871 }
872 }
873 }
874 }
875 }
876 }
877 }
878 }
879 }
880 }
881 }
882 }
883 }
884 }
885 }
886 }
887 }
888 }
889 }
890 }
891 }
892 }
893 }
894 }
895 }
896 }
897 }
898 }
899 }
900 }
901 }
902 }
903 }
904 }
905 }
906 }
907 }
908 }
909 }
910 }
911 }
912 }
913 }
914 }
915 }
916 }
917 }
918 }
919 }
920 }
921 }
922 }
923 }
924 }
925 }
926 }
927 }
928 }
929 }
930 }
931 }
932 }
933 }
934 }
935 }
936 }
937 }
938 }
939 }
940 }
941 }
942 }
943 }
944 }
945 }
946 }
947 }
948 }
949 }
950 }
951 }
952 }
953 }
954 }
955 }
956 }
957 }
958 }
959 }
960 }
961 }
962 }
963 }
964 }
965 }
966 }
967 }
968 }
969 }
970 }
971 }
972 }
973 }
974 }
975 }
976 }
977 }
978 }
979 }
980 }
981 }
982 }
983 }
984 }
985 }
986 }
987 }
988 }
989 }
990 }
991 }
992 }
993 }
994 }
995 }
996 }
997 }
998 }
999 }
1000 }

```

```

561     }
562
563     /*
564     * read the hat and its hash table
565     */
566     if (mdb_vread(&hat, sizeof (hat), (uintptr_t)hatp) == -1) {
567         mdb_warn("Couldn't read struct hat\n");
568         return (DCMD_ERR);
569     }
570
571     /*
572     * read the htable hashtable
573     */
574     for (level = 0; level <= mmu.max_level; ++level) {
575         if (level == TOP_LEVEL(&hat))
576             base = 0;
577         else
578             base = addr & mmu.level_mask[level + 1];
579
580         for (h = 0; h < hat.hat_num_hash; ++h) {
581             if (mdb_vread(&ht, sizeof (htable_t *),
582                 (uintptr_t)(hat.hat_ht_hash + h)) == -1) {
583                 mdb_warn("Couldn't read htable\n");
584                 return (DCMD_ERR);
585             }
586             for (; ht != NULL; ht = htable.ht_next) {
587                 if (mdb_vread(&htable, sizeof (htable_t),
588                     (uintptr_t)ht) == -1) {
589                     mdb_warn("Couldn't read htable\n");
590                     return (DCMD_ERR);
591                 }
592
593                 if (htable.ht_vaddr != base ||
594                     htable.ht_level != level)
595                     continue;
596
597                 pte = get_pte(&hat, &htable, addr);
598
599                 if (print_level) {
600                     mdb_printf("\tlevel=%d htable=0x%p "
601                         "pte=0x%llr\n", level, ht, pte);
602                     mdb_printf("\tlevel=%d htable=%p "
603                         "pte=%llr\n", level, ht, pte);
604                 }
605                 if (!PTE_ISVALID(pte)) {
606                     mdb_printf("Address %p is unmapped.\n",
607                         addr);
608                     return (DCMD_ERR);
609                 }
610                 if (found)
611                     continue;
612
613                 if (PTE_IS_LGPG(pte, level))
614                     paddr = mdb_ma_to_pa(pte &
615                         PT_PADDR_LGPG);
616                 else
617                     paddr = mdb_ma_to_pa(pte & PT_PADDR);
618                 paddr += addr & mmu.level_offset[level];
619                 if (pap != NULL)
620                     *pap = paddr;
621                 if (mfnp != NULL)
622                     *mfnp = pte2mfn(pte, level);
623                 found = 1;
624             }

```

```

625     }
626 }
627
628 done:
629     if (!found)
630         return (DCMD_ERR);
631     return (DCMD_OK);
632 }
633
634 unchanged_portion_omitted
635
636
637
638
639 /*
640 * Report all hat's that either use PFN as a page table or that map the page.
641 */
642 static int
643 do_report_maps(pfn_t pfn)
644 {
645     struct hat *hatp;
646     struct hat hat;
647     htable_t *ht;
648     htable_t htable;
649     uintptr_t base;
650     int h;
651     int level;
652     int entry;
653     x86pte_t pte;
654     x86pte_t buf;
655     x86pte32_t *pte32 = (x86pte32_t *)&buf;
656     physaddr_t paddr;
657     size_t len;
658
659     /*
660     * The hats are kept in a list with khat at the head.
661     */
662     for (hatp = khat; hatp != NULL; hatp = hat.hat_next) {
663         /*
664         * read the hat and its hash table
665         */
666         if (mdb_vread(&hat, sizeof (hat), (uintptr_t)hatp) == -1) {
667             mdb_warn("Couldn't read struct hat\n");
668             return (DCMD_ERR);
669         }
670
671         /*
672         * read the htable hashtable
673         */
674         paddr = 0;
675         for (h = 0; h < hat.hat_num_hash; ++h) {
676             if (mdb_vread(&ht, sizeof (htable_t *),
677                 (uintptr_t)(hat.hat_ht_hash + h)) == -1) {
678                 mdb_warn("Couldn't read htable\n");
679                 return (DCMD_ERR);
680             }
681             for (; ht != NULL; ht = htable.ht_next) {
682                 if (mdb_vread(&htable, sizeof (htable_t),
683                     (uintptr_t)ht) == -1) {
684                     mdb_warn("Couldn't read htable\n");
685                     return (DCMD_ERR);
686                 }
687
688                 /*
689                 * only report kernel addresses once
690                 */
691                 if (hatp != khat &&
692                     htable.ht_vaddr >= kernelbase)
693                     continue;

```

```

743         /*
744         * Is the PFN a pagetable itself?
745         */
746         if (htable.ht_pfn == pfn) {
747             mdb_printf("Pagetable for "
748                 "hat=%p htable=%p\n", hatp, ht);
749             continue;
750         }
751
752         /*
753         * otherwise, examine page mappings
754         */
755         level = htable.ht_level;
756         if (level > mmu.max_page_level)
757             continue;
758         paddr = mmu_ptob((physaddr_t)htable.ht_pfn);
759         for (entry = 0;
760             entry < HTABLE_NUM_PTES(&htable);
761             ++entry) {
762
763             base = htable.ht_vaddr + entry *
764                 mmu.level_size[level];
765
766             /*
767             * only report kernel addresses once
768             */
769             if (hatp != khat &&
770                 base >= kernelbase)
771                 continue;
772
773             len = mdb_pread(&pte, mmu.pte_size,
774                 len = mdb_pread(&buf, mmu.pte_size,
775                 paddr + entry * mmu.pte_size);
776             if (len != mmu.pte_size)
777                 return (DCMD_ERR);
778             if (mmu.pte_size == sizeof (x86pte_t))
779                 pte = buf;
780             else
781                 pte = *pte32;
782
783             if ((pte & PT_VALID) == 0)
784                 continue;
785             if (level == 0 || !(pte & PT_PAGESIZE))
786                 pte &= PT_PADDR;
787             else
788                 pte &= PT_PADDR_LGPG;
789             if (mmu_btop(mdb_ma_to_pa(pte)) != pfn)
790                 continue;
791             mdb_printf("hat=%p maps addr=%p\n",
792                 hatp, (caddr_t)base);
793         }
794     }
795 }
796
797 done:
798     return (DCMD_OK);
799 }
800
801 unchanged portion omitted
802
803 static int
804 do_ptable_dcmd(pfn_t pfn, uint64_t level)
805 do_ptable_dcmd(pfn_t pfn)
806 {
807     struct hat *hatp;
808     struct hat hat;

```

```

831     htable_t *ht;
832     htable_t htable;
833     uintptr_t base;
834     int h;
835     int level;
836     int entry;
837     uintptr_t pagesize;
838     x86pte_t pte;
839     x86pte_t buf;
840     x86pte32_t *pte32 = (x86pte32_t *)&buf;
841     physaddr_t paddr;
842     size_t len;
843
844     /*
845     * The hats are kept in a list with khat at the head.
846     */
847     for (hatp = khat; hatp != NULL; hatp = hat.hat_next) {
848         /*
849         * read the hat and its hash table
850         */
851         if (mdb_vread(&hat, sizeof (hat), (uintptr_t)hatp) == -1) {
852             mdb_warn("Couldn't read struct hat\n");
853             return (DCMD_ERR);
854         }
855
856         /*
857         * read the htable hashtable
858         */
859         paddr = 0;
860         for (h = 0; h < hat.hat_num_hash; ++h) {
861             if (mdb_vread(&ht, sizeof (htable_t *),
862                 (uintptr_t)(hat.hat_ht_hash + h)) == -1) {
863                 mdb_warn("Couldn't read htable\n");
864                 return (DCMD_ERR);
865             }
866             for (; ht != NULL; ht = htable.ht_next) {
867                 if (mdb_vread(&htable, sizeof (htable_t),
868                     (uintptr_t)ht) == -1) {
869                     mdb_warn("Couldn't read htable\n");
870                     return (DCMD_ERR);
871                 }
872             }
873
874             /*
875             * Is this the PFN for this htable
876             */
877             if (htable.ht_pfn == pfn)
878                 goto found_it;
879         }
880     }
881
882 found_it:
883     if (htable.ht_pfn == pfn) {
884         mdb_printf("htable=%p\n", ht);
885         if (level == (uint64_t)-1) {
886             level = htable.ht_level;
887         } else if (htable.ht_level != level) {
888             mdb_warn("htable has level %d but forcing level %lu\n",
889                 htable.ht_level, level);
890         }
891         base = htable.ht_vaddr;
892         pagesize = mmu.level_size[level];
893     } else {
894         if (level == (uint64_t)-1)
895             level = 0;
896         mdb_warn("couldn't find matching htable, using level=%lu, "

```

```

895     "base address=0x0\n", level);
896     mdb_printf("Unknown pagetable - assuming level/addr 0");
897     level = 0; /* assume level == 0 for PFN */
898     base = 0;
899     pagesize = mmu.level_size[level];
900     pagesize = MMU_PAGESIZE;
901 }
902
903 paddr = mmu_ptob((physaddr_t)pfn);
904 for (entry = 0; entry < mmu.ptes_per_table; ++entry) {
905     len = mdb_pread(&buf, mmu.pte_size,
906         paddr + entry * mmu.pte_size);
907     if (len != mmu.pte_size)
908         return (DCMD_ERR);
909     if (mmu.pte_size == sizeof (x86pte_t))
910         pte = buf;
911     else
912         pte = *pte32;
913
914     if (pte == 0)
915         continue;
916
917     mdb_printf("[%3d] va=0x%p ", entry,
918         VA_SIGN_EXTEND(base + entry * pagesize));
919     mdb_printf("[%3d] va=%p ", entry, base + entry * pagesize);
920     do_pte_dcmd(level, pte);
921 }
922
923 done:
924     return (DCMD_OK);
925 }
926
927 /*
928  * Dump the page table at the given PFN
929  */
930 /*ARGSUSED*/
931 int
932 ptable_dcmd(uintptr_t addr, uint_t flags, int argc, const mdb_arg_t *argv)
933 {
934     pfn_t pfn;
935     uint_t mflag = 0;
936     uint64_t level = (uint64_t)-1;
937
938     init_mmu();
939
940     if (mmu.num_level == 0)
941         return (DCMD_ERR);
942
943     if ((flags & DCMD_ADDRSPEC) == 0)
944         return (DCMD_USAGE);
945
946     if (mdb_getopts(argc, argv,
947         'm', MDB_OPT_SETBITS, TRUE, &mflag,
948         'l', MDB_OPT_UINT64, &level, NULL) != argc)
949         'm', MDB_OPT_SETBITS, TRUE, &mflag, NULL) != argc)
950         return (DCMD_USAGE);
951
952     if (level != (uint64_t)-1 && level > mmu.max_level) {
953         mdb_warn("invalid level %lu\n", level);
954         return (DCMD_ERR);
955     }
956
957     pfn = (pfn_t)addr;
958     if (mflag)
959         pfn = mdb_mfn_to_pfn(pfn);

```

```

953     return (do_ptable_dcmd(pfn, level));
954     return (do_ptable_dcmd(pfn));
955 }
956
957 _unchanged_portion_omitted_
958
959 /*
960  * Dump the htables for the given hat
961  */
962 /*ARGSUSED*/
963 int
964 htables_dcmd(uintptr_t addr, uint_t flags, int argc, const mdb_arg_t *argv)
965 {
966     hat_t *hat;
967
968     init_mmu();
969
970     if (mmu.num_level == 0)
971         return (DCMD_ERR);
972
973     if ((flags & DCMD_ADDRSPEC) == 0)
974         return (DCMD_USAGE);
975
976     hat = (hat_t *)addr;
977
978     return (do_htables_dcmd(hat));
979 }
980
981 static uintptr_t
982 entry2va(size_t *entries)
983 {
984     uintptr_t va = 0;
985
986     for (level_t l = mmu.max_level; l >= 0; l--)
987         va += entries[l] << mmu.level_shift[l];
988
989     return (VA_SIGN_EXTEND(va));
990 }
991
992 static void
993 ptmap_report(size_t *entries, uintptr_t start,
994     boolean_t user, boolean_t writable, boolean_t wflag)
995 {
996     uint64_t curva = entry2va(entries);
997
998     mdb_printf("mapped %s,%s range of %lu bytes: %a-%a\n",
999         user ? "user" : "kernel", writable ? "writable" : "read-only",
1000     curva - start, start, curva - 1);
1001     if (wflag && start >= kernelbase)
1002         (void) mdb_call_dcmd("whatis", start, DCMD_ADDRSPEC, 0, NULL);
1003 }
1004
1005 int
1006 ptmap_dcmd(uintptr_t addr, uint_t flags, int argc, const mdb_arg_t *argv)
1007 {
1008     physaddr_t paddr;
1009     size_t entry;
1010     uintptr_t start;
1011     boolean_t writable;
1012     boolean_t user;
1013     boolean_t wflag;
1014     level_t curlevel;
1015
1016     if ((flags & DCMD_ADDRSPEC) == 0)
1017         return (DCMD_USAGE);
1018
1019     if (mdb_getopts(argc, argv,

```

```

1054     'w', MDB_OPT_SETBITS, TRUE, &wflag, NULL) != argc)
1055         return (DCMD_USAGE);

1057     init_mmu();

1059     if (mmu.num_level == 0)
1060         return (DCMD_ERR);

1062     curlevel = mmu.max_level;

1064     paddr[curlevel] = addr & MMU_PAGEMASK;

1066     for (;;) {
1067         physaddr_t pte_addr;
1068         x86pte_t pte;

1070         pte_addr = paddr[curlevel] +
1071             (entry[curlevel] << mmu.pte_size_shift);

1073         if (mdb_pread(&pte, sizeof (pte), pte_addr) != sizeof (pte)) {
1074             mdb_warn("couldn't read pte at %p", pte_addr);
1075             return (DCMD_ERR);
1076         }

1078         if (PTE_GET(pte, PT_VALID) == 0) {
1079             if (start != (uintptr_t)-1) {
1080                 ptmap_report(entry, start,
1081                     user, writable, wflag);
1082                 start = (uintptr_t)-1;
1083             }
1084         } else if (curlevel == 0 || PTE_GET(pte, PT_PAGESIZE)) {
1085             if (start == (uintptr_t)-1) {
1086                 start = entry2va(entry);
1087                 user = PTE_GET(pte, PT_USER);
1088                 writable = PTE_GET(pte, PT_WRITABLE);
1089             } else if (user != PTE_GET(pte, PT_USER) ||
1090                 writable != PTE_GET(pte, PT_WRITABLE)) {
1091                 ptmap_report(entry, start,
1092                     user, writable, wflag);
1093                 start = entry2va(entry);
1094                 user = PTE_GET(pte, PT_USER);
1095                 writable = PTE_GET(pte, PT_WRITABLE);
1096             }
1097         } else {
1098             /* Descend a level. */
1099             physaddr_t pa = mmu_ptob(pte2mfn(pte, curlevel));
1100             paddr[--curlevel] = pa;
1101             entry[curlevel] = 0;
1102             continue;
1103         }

1105         while (++entry[curlevel] == mmu.ptes_per_table) {
1106             /* Ascend back up. */
1107             entry[curlevel] = 0;
1108             if (curlevel == mmu.max_level) {
1109                 if (start != (uintptr_t)-1) {
1110                     ptmap_report(entry, start,
1111                         user, writable, wflag);
1112                 }
1113                 goto out;
1114             }

1116             curlevel++;
1117         }
1118     }

```

```

1120 out:
1121     return (DCMD_OK);
1122 }
_____unchanged_portion_omitted_

```

new/usr/src/cmd/mdb/i86pc/modules/unix/i86mmu.h

1

```
*****
2075 Fri Apr 6 17:24:55 2018
new/usr/src/cmd/mdb/i86pc/modules/unix/i86mmu.h
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22 * Copyright 2009 Sun Microsystems, Inc. All rights reserved.
23 * Use is subject to license terms.
24 *
25 * Copyright 2018 Joyent, Inc.
26 */

28 #ifndef _I86MMU_H
29 #define _I86MMU_H

31 #ifdef __cplusplus
32 extern "C" {
33 #endif

35 extern int pte_dcmd(uintptr_t addr, uint_t flags, int argc,
36     const mdb_arg_t *argv);

38 extern int report_maps_dcmd(uintptr_t addr, uint_t flags, int argc,
39     const mdb_arg_t *argv);

41 extern int htables_dcmd(uintptr_t addr, uint_t flags, int argc,
42     const mdb_arg_t *argv);

44 extern int ptable_dcmd(uintptr_t addr, uint_t flags, int argc,
45     const mdb_arg_t *argv);
47 extern int ptmap_dcmd(uintptr_t addr, uint_t flags, int argc,
48     const mdb_arg_t *argv);

50 extern int va2pfn_dcmd(uintptr_t addr, uint_t flags, int argc,
51     const mdb_arg_t *argv);

53 extern int mfntopfndcmd(uintptr_t addr, uint_t flags, int argc,
54     const mdb_arg_t *argv);

56 extern int pfntomfndcmd(uintptr_t addr, uint_t flags, int argc,
57     const mdb_arg_t *argv);

59 extern int memseg_list(uintptr_t addr, uint_t flags, int argc,
```

new/usr/src/cmd/mdb/i86pc/modules/unix/i86mmu.h

2

```
60     const mdb_arg_t *argv);

62 extern int memseg_walk_init(mdb_walk_state_t *);
63 extern int memseg_walk_step(mdb_walk_state_t *);
64 extern void memseg_walk_fini(mdb_walk_state_t *);

66 extern void free_mmu(void);

68 #ifdef __cplusplus
69 }
    unchanged_portion_omitted
```


new/usr/src/cmd/mdb/i86pc/modules/unix/unix.c

1

```
*****
25628 Fri Apr 6 17:24:56 2018
new/usr/src/cmd/mdb/i86pc/modules/unix/unix.c
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>
9210 remove KMDB branch debugging support
9211 ::crregs could do with cr2/cr3 support
9209 ::ttrace should be able to filter by thread
Reviewed by: Patrick Mooney <patrick.mooney@joyent.com>
Reviewed by: Yuri Pankov <yuripv@yuripv.net>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22 * Copyright (c) 1999, 2010, Oracle and/or its affiliates. All rights reserved.
23 * Copyright 2018 Joyent, Inc.
23 * Copyright 2015 Joyent, Inc.
24 */

26 #include <mdb/mdb_modapi.h>
27 #include <mdb/mdb_ctf.h>
28 #include <sys/cpuvar.h>
29 #include <sys/system.h>
30 #include <sys/traptrace.h>
31 #include <sys/x_call.h>
32 #include <sys/xc_levels.h>
33 #include <sys/avintr.h>
34 #include <sys/system.h>
35 #include <sys/trap.h>
36 #include <sys/mutex.h>
37 #include <sys/mutex_impl.h>
38 #include "i86mmu.h"
39 #include "unix_sup.h"
40 #include <sys/apix.h>
41 #include <sys/x86_archext.h>
42 #include <sys/bitmap.h>
43 #include <sys/controlregs.h>

45 #define TT_HDLR_WIDTH 17

48 /* apix only */
49 static apix_impl_t *d_apix[NCPU];
50 static int use_apix = 0;

52 static int
53 ttrace_ttr_size_check(void)
```

new/usr/src/cmd/mdb/i86pc/modules/unix/unix.c

2

```
54 {
55     mdb_ctf_id_t ttrtid;
56     ssize_t ttr_size;

58     if (mdb_ctf_lookup_by_name("trap_trace_rec_t", &ttrtid) != 0 ||
59         mdb_ctf_type_resolve(ttrtid, &ttrtid) != 0) {
60         mdb_warn("failed to determine size of trap_trace_rec_t; "
61             "non-TRAPTRACE kernel?\n");
62         return (0);
63     }

65     if ((ttr_size = mdb_ctf_type_size(ttrtid)) !=
66         sizeof (trap_trace_rec_t)) {
67         /*
68          * On Intel machines, this will happen when TTR_STACK_DEPTH
69          * is changed. This code could be smarter, and could
70          * dynamically adapt to different depths, but not until a
71          * need for such adaptation is demonstrated.
72          */
73         mdb_warn("size of trap_trace_rec_t (%d bytes) doesn't "
74             "match expected %d\n", ttr_size, sizeof (trap_trace_rec_t));
75         return (0);
76     }

78     return (1);
79 }
_____ unchanged_portion_omitted_

409 typedef struct ttrace_dcmd {
410     processorid_t ttd_cpu;
411     uint_t ttd_extended;
412     uintptr_t ttd_kthread;
413     trap_trace_ctl_t ttd_ttc[NCPU];
414 } ttrace_dcmd_t;

416 #if defined(__amd64)

418 #define DUMP(reg) #reg, regs->r_##reg
419 #define THREEREGS " %3s: %16lx %3s: %16lx %3s: %16lx\n"

421 static void
422 ttrace_dumpregs(trap_trace_rec_t *rec)
423 {
424     struct regs *regs = &rec->ttr_regs;

426     mdb_printf(THREEREGS, DUMP(rdi), DUMP(rsi), DUMP(rdx));
427     mdb_printf(THREEREGS, DUMP(rcx), DUMP(r8), DUMP(r9));
428     mdb_printf(THREEREGS, DUMP(rax), DUMP(rbx), DUMP(rbp));
429     mdb_printf(THREEREGS, DUMP(r10), DUMP(r11), DUMP(r12));
430     mdb_printf(THREEREGS, DUMP(r13), DUMP(r14), DUMP(r15));
431     mdb_printf(THREEREGS, DUMP(ds), DUMP(es), DUMP(fs));
432     mdb_printf(THREEREGS, DUMP(gs), "trp", regs->r_trapno, DUMP(err));
433     mdb_printf(THREEREGS, DUMP(rip), DUMP(cs), DUMP(rfl));
434     mdb_printf(THREEREGS, DUMP(rsp), DUMP(ss), "cr2", rec->ttr_cr2);
435     mdb_printf(" %3s: %16lx %3s: %16lx\n",
436         "fsb", regs->r_fsbbase,
437         "gsb", regs->r_gsbbase);
438     mdb_printf("\n");
439 }
_____ unchanged_portion_omitted_

461 #endif /* __amd64 */

463 int
464 ttrace_walk(uintptr_t addr, trap_trace_rec_t *rec, ttrace_dcmd_t *dcmd)
465 {
```

```

466     struct regs *regs = &rec->ttr_regs;
467     processorid_t cpu = -1, i;

469     for (i = 0; i < NCPU; i++) {
470         if (addr >= dcmd->ttd_ttc[i].ttc_first &&
471             addr < dcmd->ttd_ttc[i].ttc_limit) {
472             cpu = i;
473             break;
474         }
475     }

477     if (cpu == -1) {
478         mdb_warn("couldn't find %p in any trap trace ctl\n", addr);
479         return (WALK_ERR);
480     }

482     if (dcmd->ttd_cpu != -1 && cpu != dcmd->ttd_cpu)
483         return (WALK_NEXT);

485     if (dcmd->ttd_kthread != 0 &&
486         dcmd->ttd_kthread != rec->ttr_curthread)
487         return (WALK_NEXT);

489     mdb_printf("%3d %15llx ", cpu, rec->ttr_stamp);

491     for (i = 0; ttrace_hdlr[i].t_hdlr != NULL; i++) {
492         if (rec->ttr_marker != ttrace_hdlr[i].t_marker)
493             continue;
494         mdb_printf("%4s ", ttrace_hdlr[i].t_name);
495         if (ttrace_hdlr[i].t_hdlr(rec) == -1)
496             return (WALK_ERR);
497     }

499     mdb_printf(" %a\n", regs->r_pc);

501     if (dcmd->ttd_extended == FALSE)
502         return (WALK_NEXT);

504     if (rec->ttr_marker == TT_INTERRUPT)
505         ttrace_intr_detail(rec);
506     else
507         ttrace_dumpregs(rec);

509     if (rec->ttr_sdepth > 0) {
510         for (i = 0; i < rec->ttr_sdepth; i++) {
511             if (i >= TTR_STACK_DEPTH) {
512                 mdb_printf("%17s*** invalid ttr_sdepth (is %d, "
513                     "should be <= %d)\n", " ", rec->ttr_sdepth,
514                     TTR_STACK_DEPTH);
515                 break;
516             }

518             mdb_printf("%17s %a()\n", " ", rec->ttr_stack[i]);
519         }
520         mdb_printf("\n");
521     }

523     return (WALK_NEXT);
524 }

526 int
527 ttrace(uintptr_t addr, uint_t flags, int argc, const mdb_arg_t *argv)
528 {
529     ttrace_dcmd_t dcmd;
530     trap_trace_ctl_t *ttc = dcmd.ttd_ttc;
531     trap_trace_rec_t rec;

```

```

532     size_t ttc_size = sizeof (trap_trace_ctl_t) * NCPU;

534     if (!ttrace_ttr_size_check())
535         return (WALK_ERR);

537     bzero(&dcmd, sizeof (dcmd));
538     dcmd.ttd_cpu = -1;
539     dcmd.ttd_extended = FALSE;

541     if (mdb_readsym(ttc, ttc_size, "trap_trace_ctl") == -1) {
542         mdb_warn("symbol 'trap_trace_ctl' not found; "
543             "non-TRAPTRACE kernel?\n");
544         return (DCMD_ERR);
545     }

547     if (mdb_getopts(argc, argv,
548         'x', MDB_OPT_SETBITS, TRUE, &dcmd.ttd_extended,
549         't', MDB_OPT_UINTPTR, &dcmd.ttd_kthread, NULL) != argc)
550         return (DCMD_USAGE);

552     if (DCMD_HDRSPEC(flags)) {
553         mdb_printf("%3s %15s %4s %2s %-*s\n", "CPU",
554             "TIMESTAMP", "TYPE", "Vec", TT_HDLR_WIDTH, "HANDLER",
555             "EIP");
556     }

558     if (flags & DCMD_ADDRSPEC) {
559         if (addr >= NCPU) {
560             if (mdb_vread(&rec, sizeof (rec), addr) == -1) {
561                 mdb_warn("couldn't read trap trace record "
562                     "at %p", addr);
563                 return (DCMD_ERR);
564             }

566             if (ttrace_walk(addr, &rec, &dcmd) == WALK_ERR)
567                 return (DCMD_ERR);

569             return (DCMD_OK);
570         }
571         dcmd.ttd_cpu = addr;
572     }

574     if (mdb_readvar(&use_apix, "apix_enable") == -1) {
575         mdb_warn("failed to read apix_enable");
576         use_apix = 0;
577     }

579     if (use_apix) {
580         if (mdb_readvar(&d_apixs, "apixs") == -1) {
581             mdb_warn("\nfailed to read apixs.");
582             return (DCMD_ERR);
583         }
584         /* change to apix ttrace interrupt handler */
585         ttrace_hdlr[4].t_hdlr = ttrace_apix_interrupt;
586     }

588     if (mdb_walk("ttrace", (mdb_walk_cb_t)ttrace_walk, &dcmd) == -1) {
589         mdb_warn("couldn't walk 'ttrace'");
590         return (DCMD_ERR);
591     }

593     return (DCMD_OK);
594 }

```

unchanged portion omitted

```

752 static void
753 ptable_help(void)
754 {
755     mdb_printf(
756         "Given a PFN holding a page table, print its contents, and\n"
757         "the address of the corresponding htable structure.\n"
758         "\n"
759         "-m Interpret the PFN as an MFN (machine frame number)\n"
760         "-l force page table level (3 is top)\n");
761     "-m Interpret the PFN as an MFN (machine frame number)\n");
762 }

763 static void
764 ptmap_help(void)
765 {
766     mdb_printf(
767         "Report all mappings represented by the page table hierarchy\n"
768         "rooted at the given cr3 value / physical address.\n"
769         "\n"
770         "-w run ::what is on mapping start addresses\n");
771 }

772 /*
773 * NSEC_SHIFT is replicated here (it is not defined in a header file),
774 * but for amusement, the reader is directed to the comment that explains
775 * the rationale for this particular value on x86. Spoiler: the value is
776 * selected to accommodate 60 MHz Pentiums! (And a confession: if the voice
777 * in that comment sounds too familiar, it's because your author also wrote
778 * that code -- some fifteen years prior to this writing in 2011...)
779 */
780 #define NSEC_SHIFT 5
781 #define NSEC_SHIFT 5

782 /*ARGSUSED*/
783 static int
784 scalehrtime_cmd(uintptr_t addr, uint_t flags, int argc, const mdb_arg_t *argv)
785 {
786     uint32_t nsec_scale;
787     hrttime_t tsc = addr, hrt;
788     unsigned int *tscp = (unsigned int *)&tsc;
789     uintptr_t scalehrtimef;
790     uint64_t scale;
791     GElf_Sym sym;

792     if (!(flags & DCMD_ADDRSPEC)) {
793         if (argc != 1)
794             return (DCMD_USAGE);

795         switch (argv[0].a_type) {
796             case MDB_TYPE_STRING:
797                 tsc = mdb_strtoull(argv[0].a_un.a_str);
798                 break;
799             case MDB_TYPE_IMMEDIATE:
800                 tsc = argv[0].a_un.a_val;
801                 break;
802             default:
803                 return (DCMD_USAGE);
804         }
805     }

806     if (mdb_readsym(&scalehrtimef,
807         sizeof (scalehrtimef), "scalehrtimef") == -1) {
808         mdb_warn("couldn't read 'scalehrtimef'");
809         return (DCMD_ERR);
810     }

811     if (mdb_lookup_by_name("tsc_scalehrtime", &sym) == -1) {

```

```

812         mdb_warn("couldn't find 'tsc_scalehrtime'");
813         return (DCMD_ERR);
814     }

815     if (sym.st_value != scalehrtimef) {
816         mdb_warn("::scalehrtime requires that scalehrtimef "
817             "be set to tsc_scalehrtime\n");
818         return (DCMD_ERR);
819     }

820     if (mdb_readsym(&nsec_scale, sizeof (nsec_scale), "nsec_scale") == -1) {
821         mdb_warn("couldn't read 'nsec_scale'");
822         return (DCMD_ERR);
823     }

824     scale = (uint64_t)nsec_scale;
825     hrt = ((uint64_t)tscp[1] * scale) << NSEC_SHIFT;
826     hrt += ((uint64_t)tscp[0] * scale) >> (32 - NSEC_SHIFT);

827     mdb_printf("0x%llx\n", hrt);

828     return (DCMD_OK);
829 }

830 unchanged portion omitted

831 #ifdef _KMDB
832 /* ARGSUSED */
833 static int
834 crregs_dcnd(uintptr_t addr, uint_t flags, int argc, const mdb_arg_t *argv)
835 {
836     ulong_t cr0, cr2, cr3, cr4;
837     ulong_t cr0, cr4;
838     static const mdb_bitmask_t cr0_flag_bits[] = {
839         { "PE", CR0_PE, CR0_PE },
840         { "MP", CR0_MP, CR0_MP },
841         { "EM", CR0_EM, CR0_EM },
842         { "TS", CR0_TS, CR0_TS },
843         { "ET", CR0_ET, CR0_ET },
844         { "NE", CR0_NE, CR0_NE },
845         { "WP", CR0_WP, CR0_WP },
846         { "AM", CR0_AM, CR0_AM },
847         { "NW", CR0_NW, CR0_NW },
848         { "CD", CR0_CD, CR0_CD },
849         { "PG", CR0_PG, CR0_PG },
850         { NULL, 0, 0 }
851     };

852     static const mdb_bitmask_t cr3_flag_bits[] = {
853         { "PCD", CR3_PCD, CR3_PCD },
854         { "PWT", CR3_PWT, CR3_PWT },
855         { NULL, 0, 0 }
856     };

857     static const mdb_bitmask_t cr4_flag_bits[] = {
858         { "VME", CR4_VME, CR4_VME },
859         { "PVI", CR4_PVI, CR4_PVI },
860         { "TSD", CR4_TSD, CR4_TSD },
861         { "DE", CR4_DE, CR4_DE },
862         { "PSE", CR4_PSE, CR4_PSE },
863         { "PAE", CR4_PAE, CR4_PAE },
864         { "MCE", CR4_MCE, CR4_MCE },
865         { "PGE", CR4_PGE, CR4_PGE },
866         { "PCE", CR4_PCE, CR4_PCE },
867         { "OSFXSR", CR4_OSFXSR, CR4_OSFXSR },
868         { "OSXMMEXCPT", CR4_OSXMMEXCPT, CR4_OSXMMEXCPT }
869     };

```

```

943     { "VMXE",          CR4_VMXE,      CR4_VMXE },
944     { "SMXE",          CR4_SMXE,      CR4_SMXE },
945     { "PCIDE",         CR4_PCIDE,     CR4_PCIDE },
946     { "OSXSAVE",      CR4_OSXSAVE,   CR4_OSXSAVE },
947     { "SMEP",          CR4_SMEP,      CR4_SMEP },
948     { "SMAP",          CR4_SMAP,      CR4_SMAP },
949     { NULL,            0,              0 }
950 };

952     cr0 = kmdb_unix_getcr0();
953     cr2 = kmdb_unix_getcr2();
954     cr3 = kmdb_unix_getcr3();
955     cr4 = kmdb_unix_getcr4();
956     mdb_printf("%%cr0 = 0x%lx <b>\n", cr0, cr0, cr0_flag_bits);
957     mdb_printf("%%cr2 = 0x%lx <a>\n", cr2, cr2);

959     if ((cr4 & CR4_PCIDE)) {
960         mdb_printf("%%cr3 = 0x%lx <pfn:0x%lx pcid:%lu>\n", cr3,
961             cr3 >> MMU_PAGESHIFT, cr3 & MMU_PAGEOFFSET);
962     } else {
963         mdb_printf("%%cr3 = 0x%lx <pfn:0x%lx flags:%b>\n", cr3,
964             cr3 >> MMU_PAGESHIFT, cr3, cr3_flag_bits);
965     }

967     mdb_printf("%%cr4 = 0x%lx <b>\n", cr4, cr4, cr4_flag_bits);

927     mdb_printf("%%cr0 = 0x%08x <b>\n", cr0, cr0, cr0_flag_bits);
928     mdb_printf("%%cr4 = 0x%08x <b>\n", cr4, cr4, cr4_flag_bits);
969     return (DCMD_OK);
970 }
971 #endif

973 static const mdb_dcmd_t dcmds[] = {
974     { "gate_desc", ":", "dump a gate descriptor", gate_desc },
975     { "idt", ":[-v]", "dump an IDT", idt },
976     { "ttrace", "[-x] [-t kthread]", "dump trap trace buffers", ttrace },
977     { "ttrace", "[-x]", "dump trap trace buffers", ttrace },
978     { "vatopfn", ":[-a as]", "translate address to physical page",
979         va2pfn_dcmd },
980     { "report_maps", ":[-m]",
981         "Given PFN, report mappings / page table usage",
982         report_maps_dcmd, report_maps_help },
983     { "htables", "", "Given hat_t *, lists all its htable_t * values",
984         htables_dcmd, htables_help },
985     { "ptable", ":[-lm]", "Given PFN, dump contents of a page table",
986         "ptable", ":[-m]", "Given PFN, dump contents of a page table",
987         ptable_dcmd, ptable_help },
988     { "ptmap", ":", "Given a cr3 value, dump all mappings",
989         ptmap_dcmd, ptmap_help },
990     { "pte", ":[-l N]", "print human readable page table entry",
991         "pte", ":[-p XXXXX] [-l N]", "print human readable page table entry",
992         pte_dcmd },
993     { "pfntomfn", ":", "convert physical page to hypervisor machine page",
994         pfntomfn_dcmd },
995     { "mfntopf", ":", "convert hypervisor machine page to physical page",
996         mfntopf_dcmd },
997     { "memseg_list", ":", "show memseg list", memseg_list },
998     { "scalehrtime", ":",
999         "scale an unscaled high-res time", scalehrtime_cmd },
1000     { "x86_featureset", NULL, "dump the x86_featureset vector",
1001         x86_featureset_dcmd },
1002     #ifdef _KMDB
1003     { "crregs", NULL, "dump control registers", crregs_dcmd },
1004     #endif
1005     { NULL }
1006 };

```

unchanged_portion_omitted

new/usr/src/cmd/mdb/i86pc/modules/unix/unix_sup.h

1

807 Fri Apr 6 17:24:56 2018

new/usr/src/cmd/mdb/i86pc/modules/unix/unix_sup.h

9210 remove KMDB branch debugging support

9211 ::crregs could do with cr2/cr3 support

9209 ::ttrace should be able to filter by thread

Reviewed by: Patrick Mooney <patrick.mooney@joyent.com>

Reviewed by: Yuri Pankov <yuripv@yuripv.net>

```
1 /*
2  * This file and its contents are supplied under the terms of the
3  * Common Development and Distribution License ("CDDL"), version 1.0.
4  * You may only use this file in accordance with the terms of version
5  * 1.0 of the CDDL.
6  *
7  * A full copy of the text of the CDDL should have accompanied this
8  * source. A copy of the CDDL is also available via the Internet at
9  * http://www.illumos.org/license/CDDL.
10 */
```

```
12 /*
13  * Copyright 2018 Joyent, Inc.
14  * Copyright 2015 Joyent, Inc.
15 */
```

```
16 #ifndef _UNIX_SUP_H
17 #define _UNIX_SUP_H
```

```
19 /*
20  * Support routines for unix.
21 */
```

```
23 #include <sys/types.h>
```

```
25 #ifdef __cplusplus
26 extern "C" {
27 #endif
```

```
29 extern ulong_t kmdb_unix_getcr0(void);
30 extern ulong_t kmdb_unix_getcr2(void);
31 extern ulong_t kmdb_unix_getcr3(void);
32 extern ulong_t kmdb_unix_getcr4(void);
```

```
34 #ifdef __cplusplus
35 }
```

_____ unchanged portion omitted

new/usr/src/cmd/mdb/i86pc/modules/unix/unix_sup.s

1

1500 Fri Apr 6 17:24:56 2018
new/usr/src/cmd/mdb/i86pc/modules/unix/unix_sup.s
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>
9210 remove KMDB branch debugging support
9211 ::crregs could do with cr2/cr3 support
9209 ::ttrace should be able to filter by thread
Reviewed by: Patrick Mooney <patrick.mooney@joyent.com>
Reviewed by: Yuri Pankov <yuripv@yuripv.net>

```
1 /*
2  * This file and its contents are supplied under the terms of the
3  * Common Development and Distribution License ("CDDL"), version 1.0.
4  * You may only use this file in accordance with the terms of version
5  * 1.0 of the CDDL.
6  *
7  * A full copy of the text of the CDDL should have accompanied this
8  * source. A copy of the CDDL is also available via the Internet at
9  * http://www.illumos.org/license/CDDL.
10 */
```

```
12 /*
13  * Copyright 2018 Joyent, Inc.
14  * Copyright 2015 Joyent, Inc.
15 */
```

```
16 #if !defined(__lint)
17     .file "unix_sup.s"
18 #endif /* __lint */
```

```
20 /*
21  * Support routines for the unix kmdb module
22 */
```

```
24 #include <sys/asm_linkage.h>
```

```
26 #if defined(__lint)
```

```
28 #include <sys/types.h>
```

```
30 ulong_t
31 kmdb_unix_getcr0(void)
32 { return (0); }
```

```
34 ulong_t
35 kmdb_unix_getcr3(void)
36 { return (0); }
```

```
38 ulong_t
39 kmdb_unix_getcr4(void)
40 { return (0); }
```

```
42 #else /* __lint */
```

```
44 #if defined(__amd64)
45     ENTRY(kmdb_unix_getcr0)
46     movq %cr0, %rax
47     ret
48     SET_SIZE(kmdb_unix_getcr0)
```

```
50     ENTRY(kmdb_unix_getcr2)
51     movq %cr2, %rax
52     ret
53     SET_SIZE(kmdb_unix_getcr2)
```

new/usr/src/cmd/mdb/i86pc/modules/unix/unix_sup.s

2

```
55     ENTRY(kmdb_unix_getcr3)
56     movq %cr3, %rax
57     ret
58     SET_SIZE(kmdb_unix_getcr3)
```

```
60     ENTRY(kmdb_unix_getcr4)
61     movq %cr4, %rax
62     ret
63     SET_SIZE(kmdb_unix_getcr4)
_____ unchanged_portion_omitted_
```

```
71     ENTRY(kmdb_unix_getcr2)
72     movl %cr2, %eax
73     ret
74     SET_SIZE(kmdb_unix_getcr2)

76     ENTRY(kmdb_unix_getcr3)
77     movl %cr3, %eax
78     ret
79     SET_SIZE(kmdb_unix_getcr3)
```

```
81     ENTRY(kmdb_unix_getcr4)
82     movl %cr4, %eax
83     ret
84     SET_SIZE(kmdb_unix_getcr4)
_____ unchanged_portion_omitted_
```

new/usr/src/cmd/mdb/intel/Makefile.kmdb

1

1842 Fri Apr 6 17:24:56 2018
new/usr/src/cmd/mdb/intel/Makefile.kmdb
9210 remove KMDB branch debugging support
9211 ::crregs could do with cr2/cr3 support
9209 ::ttrace should be able to filter by thread
Reviewed by: Patrick Mooney <patrick.mooney@joyent.com>
Reviewed by: Yuri Pankov <yuripv@yuripv.net>

```
1 #
2 # CDDL HEADER START
3 #
4 # The contents of this file are subject to the terms of the
5 # Common Development and Distribution License (the "License").
6 # You may not use this file except in compliance with the License.
7 #
8 # You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9 # or http://www.opensolaris.org/os/licensing.
10 # See the License for the specific language governing permissions
11 # and limitations under the License.
12 #
13 # When distributing Covered Code, include this CDDL HEADER in each
14 # file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 # If applicable, add the following below this CDDL HEADER, with the
16 # fields enclosed by brackets "[]" replaced with your own identifying
17 # information: Portions Copyright [yyyy] [name of copyright owner]
18 #
19 # CDDL HEADER END
20 #
21 #
22 # Copyright 2009 Sun Microsystems, Inc. All rights reserved.
23 # Use is subject to license terms.
24 #
25 # Copyright 2018 Joyent, Inc.
26 #
```

```
28 PROMSRCS += \
29     prom_env.c \
30     prom_getchar.c \
31     prom_init.c \
32     prom_printf.c \
33     prom_putchar.c
```

```
35 KMDBSRCS += \
36     kaif.c \
37     kmdb_dpi_isadep.c \
38     kmdb_fault_isadep.c \
39     kmdb_kdi_isadep.c \
40     kmdb_promif_isadep.c \
41     kvm_cpu_amd.c \
42     kvm_cpu_p4.c \
43     kvm_isadep.c
```

```
43 KMDBML += \
44     kmdb_asmutil.s \
45     kmdb_setcontext.s
```

```
47 KCTLSRCS += \
48     kctl_isadep.c
```

```
50 CTXOFFUSERS = \
51     kmdb_setcontext.o
```

```
53 CPPFLAGS += -DDIS_TEXT
```

```
55 $(CTXOFFUSERS) $(CTXOFFUSERS:%.o=%.ln): kmdb_context_off.h
```

new/usr/src/cmd/mdb/intel/Makefile.kmdb

2

```
57 kaif_activate.o kaif_activate.ln := CPPFLAGS += -D_MACHDEP -D_KMEMUSER
```

```
59 STANDBLIBS += \
60     ../libstandctf/libstandctf.so \
61     $(SRC)/lib/libumem/$(MACHDIR)/libstandumem.so \
62     ../libstand/libstand.a
```

```
64 KMDBLIBS = $(STANDBLIBS) ../mdb_ks/kmod/mdb_ks
```

```
66 MAPFILE_SOURCES = \
67     $(MAPFILE_SOURCES_COMMON) \
68     ../../kmdb/kmdb_dpi_isadep.h \
69     $(MAPFILE_SOURCES_$(MACH))
```

```
71 %.o: ../../../../uts/intel/promif/%.c
72     $(COMPILE.c) $<
73     $(CTFCONVERT_O)
```

```
75 %.ln: ../../../../uts/intel/promif/%.c
76     $(LINT.c) -c $<
```

new/usr/src/cmd/mdb/intel/ia32/Makefile.kmdb

1

1053 Fri Apr 6 17:24:56 2018

new/usr/src/cmd/mdb/intel/ia32/Makefile.kmdb

9210 remove KMDB branch debugging support

9211 ::crregs could do with cr2/cr3 support

9209 ::ttrace should be able to filter by thread

Reviewed by: Patrick Mooney <patrick.mooney@joyent.com>

Reviewed by: Yuri Pankov <yuripv@yuripv.net>

```
1 #
2 # CDDL HEADER START
3 #
4 # The contents of this file are subject to the terms of the
5 # Common Development and Distribution License (the "License").
6 # You may not use this file except in compliance with the License.
7 #
8 # You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9 # or http://www.opensolaris.org/os/licensing.
10 # See the License for the specific language governing permissions
11 # and limitations under the License.
12 #
13 # When distributing Covered Code, include this CDDL HEADER in each
14 # file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 # If applicable, add the following below this CDDL HEADER, with the
16 # fields enclosed by brackets "[]" replaced with your own identifying
17 # information: Portions Copyright [yyyy] [name of copyright owner]
18 #
19 # CDDL HEADER END
20 #
21 #
22 # Copyright 2007 Sun Microsystems, Inc. All rights reserved.
23 # Use is subject to license terms.
24 #
25 # Copyright 2018 Joyent, Inc.
26 #ident "%Z%M% %I% %E% SMI"
```

```
27 KMDBML += \
28     kaif_invoke.s \
29     kmdb_start.s
```

```
31 KMDBSRCS += \
32     mdb_ia32util.c \
33     kmdb_makecontext.c
34     kmdb_makecontext.c \
35     kvm_cpu_p6.c
```

```
35 SACPPFLAGS = -D__$(MACH)
```


new/usr/src/cmd/mdb/intel/kmdb/kaif.c

1

19998 Fri Apr 6 17:24:56 2018

new/usr/src/cmd/mdb/intel/kmdb/kaif.c

8956 Implement KPTI

Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>

Reviewed by: Robert Mustacchi <rm@joyent.com>

9210 remove KMDB branch debugging support

9211 ::crregs could do with cr2/cr3 support

9209 ::ttrace should be able to filter by thread

Reviewed by: Patrick Mooney <patrick.mooney@joyent.com>

Reviewed by: Yuri Pankov <yuripv@yuripv.net>

```
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22 * Copyright 2007 Sun Microsystems, Inc. All rights reserved.
23 * Use is subject to license terms.
24 *
25 * Copyright 2018 Joyent, Inc.
26 */

26 #pragma ident "%Z%M% %I% %E% SMI"

28 /*
29  * The debugger/"PROM" interface layer
30  *
31  * It makes more sense on SPARC. In reality, these interfaces deal with three
32  * things: setting break/watchpoints, stepping, and interfacing with the KDI to
33  * set up kmdb's IDT handlers.
34  */

36 #include <kmdb/kmdb_dpi_impl.h>
37 #include <kmdb/kmdb_kdi.h>
38 #include <kmdb/kmdb_umemglue.h>
39 #include <kmdb/kaif.h>
40 #include <kmdb/kmdb_io.h>
41 #include <kmdb/kaif_start.h>
42 #include <mdb/mdb_err.h>
43 #include <mdb/mdb_debug.h>
44 #include <mdb/mdb_isautil.h>
45 #include <mdb/mdb_io_impl.h>
46 #include <mdb/mdb_kreg_impl.h>
47 #include <mdb/mdb.h>

49 #include <sys/types.h>
50 #include <sys/bitmap.h>
51 #include <sys/termios.h>
52 #include <sys/kdi_impl.h>
```

new/usr/src/cmd/mdb/intel/kmdb/kaif.c

2

53 #include <sys/sysmacros.h>

```
55 /*
56  * This is the area containing the saved state when we enter
57  * via kmdb's IDT entries.
58  */
59 kdi_cpucsave_t    *kaif_cpucsave;
60 int               kaif_ncpusave;
61 kdi_drreg_t       kaif_drreg;

63 uint32_t          kaif_waptmap;

65 int               kaif_trap_switch;

67 void (*kaif_modchg_cb)(struct modctl *, int);

69 enum {
70     M_SYSRET      = 0x07, /* after M_ESC */
71     M_ESC         = 0x0f,
72     M_SYSEXIT     = 0x35, /* after M_ESC */
73     M_REX_LO      = 0x40, /* first REX prefix */
74     M_REX_HI      = 0x4f, /* last REX prefix */
75     M_PUSHF       = 0x9c, /* pushfl and pushfq */
76     M_POPF        = 0x9d, /* popfl and popfq */
77     M_INT3        = 0xcc,
78     M_INTX        = 0xcd,
79     M_INT0        = 0xce,
80     M_IRET        = 0xcf,
81     M_CLI         = 0xfa,
82     M_STI         = 0xfb
83 };
    unchanged_portion_omitted

260 /*
261  * Refuse to single-step or break within any stub that loads a user %cr3 value.
262  * As the KDI traps are not careful to restore such a %cr3, this can all go
263  * wrong, both spectacularly and subtly.
264  */
265 static boolean_t
266 kaif_toxic_text(uintptr_t addr)
267 {
268     static GElf_Sym toxic_syms[2] = { 0, };
269     size_t i;

271     if (toxic_syms[0].st_name == NULL) {
272         if (mdb_tgt_lookup_by_name(mdb.m_target, MDB_TGT_OBJ_EXEC,
273             "tr_iret_user", &toxic_syms[0], NULL) != 0)
274             warn("couldn't find tr_iret_user\n");
275         if (mdb_tgt_lookup_by_name(mdb.m_target, MDB_TGT_OBJ_EXEC,
276             "tr_mmu_flush_user_range", &toxic_syms[1], NULL) != 0)
277             warn("couldn't find tr_mmu_flush_user_range\n");
278     }

280     for (i = 0; i < ARRAY_SIZE(toxic_syms); i++) {
281         if (addr >= toxic_syms[i].st_value &&
282             addr - toxic_syms[i].st_value < toxic_syms[i].st_size)
283             return (B_TRUE);
284     }

286     return (B_FALSE);
287 }

289 static int
290 kaif_brkpt_arm(uintptr_t addr, mdb_instr_t *instrp)
291 {
292     mdb_instr_t bkpt = KAIF_BREAKPOINT_INSTR;
```

```

294     if (kaif_toxic_text(addr)) {
295         warn("%a cannot be a breakpoint target\n", addr);
296         return (set_errno(EMDB_TGTNOTSUP));
297     }
299     if (mdb_tgt_vread(mdb.m_target, instrp, sizeof (mdb_instr_t), addr) !=
300         sizeof (mdb_instr_t))
301         return (-1); /* errno is set for us */
303     if (mdb_tgt_vwrite(mdb.m_target, &bkpt, sizeof (mdb_instr_t), addr) !=
304         sizeof (mdb_instr_t))
305         return (-1); /* errno is set for us */
307     return (0);
308 }
    unchanged portion omitted
472 static int
473 kaif_step(void)
474 {
475     kreg_t pc, fl, oldfl, newfl, sp;
476     mdb_tgt_addr_t npc;
477     mdb_instr_t instr;
478     int emulated = 0, rchk = 0;
479     size_t pcoff = 0;
481     (void) kmdb_dpi_get_register("pc", &pc);
483     if (kaif_toxic_text(pc)) {
484         warn("%a cannot be stepped\n", pc);
485         return (set_errno(EMDB_TGTNOTSUP));
486     }
488     if ((npc = mdb_dis_nextins(mdb.m_disasm, mdb.m_target,
489         MDB_TGT_AS_VIRT, pc)) == pc) {
490         warn("failed to decode instruction at %a for step\n", pc);
491         return (set_errno(EINVAL));
492     }
494     /*
495     * Stepping behavior depends on the type of instruction.  It does not
496     * depend on the presence of a REX prefix, as the action we take for a
497     * given instruction doesn't currently vary for 32-bit instructions
498     * versus their 64-bit counterparts.
499     */
500     do {
501         if (mdb_tgt_vread(mdb.m_target, &instr, sizeof (mdb_instr_t),
502             pc + pcoff) != sizeof (mdb_instr_t)) {
503             warn("failed to read at %p for step",
504                 (void *) (pc + pcoff));
505             return (-1);
506         }
507     } while (pcoff++, (instr >= M_REX_LO && instr <= M_REX_HI && !rchk++));
509     switch (instr) {
510     case M_IRET:
511         warn("iret cannot be stepped\n");
512         return (set_errno(EMDB_TGTNOTSUP));
514     case M_INT3:
515     case M_INTX:
516     case M_INT0:
517         warn("int cannot be stepped\n");
518         return (set_errno(EMDB_TGTNOTSUP));

```

```

520     case M_ESC:
521         if (mdb_tgt_vread(mdb.m_target, &instr, sizeof (mdb_instr_t),
522             pc + pcoff) != sizeof (mdb_instr_t)) {
523             warn("failed to read at %p for step",
524                 (void *) (pc + pcoff));
525             return (-1);
526         }
528     switch (instr) {
529     case M_SYSRET:
530         warn("sysret cannot be stepped\n");
531         return (set_errno(EMDB_TGTNOTSUP));
532     case M_SYSEXIT:
533         warn("sysexit cannot be stepped\n");
534         return (set_errno(EMDB_TGTNOTSUP));
535     }
536     break;
538     /*
539     * Some instructions need to be emulated.  We need to prevent direct
540     * manipulations of EFLAGS, so we'll emulate cli, sti, pushfl and
541     * popfl also receive special handling, as they manipulate both EFLAGS
542     * and %esp.
543     */
544     case M_CLLI:
545         (void) kmdb_dpi_get_register(FLAGS_REG_NAME, &fl);
546         fl &= ~KREG_EFLAGS_IF_MASK;
547         (void) kmdb_dpi_set_register(FLAGS_REG_NAME, fl);
549         emulated = 1;
550         break;
552     case M_STI:
553         (void) kmdb_dpi_get_register(FLAGS_REG_NAME, &fl);
554         fl |= (1 << KREG_EFLAGS_IF_SHIFT);
555         (void) kmdb_dpi_set_register(FLAGS_REG_NAME, fl);
557         emulated = 1;
558         break;
560     case M_POPLF:
561         /*
562         * popfl will restore a pushed EFLAGS from the stack, and could
563         * in so doing cause IF to be turned on, if only for a brief
564         * period.  To avoid this, we'll secretly replace the stack's
565         * EFLAGS with our decaffeinated brand.  We'll then manually
566         * load our EFLAGS copy with the real verion after the step.
567         */
568         (void) kmdb_dpi_get_register("sp", &sp);
569         (void) kmdb_dpi_get_register(FLAGS_REG_NAME, &fl);
571         if (mdb_tgt_vread(mdb.m_target, &newfl, sizeof (kreg_t),
572             sp) != sizeof (kreg_t)) {
573             warn("failed to read " FLAGS_REG_NAME
574                 " at %p for popfl step\n", (void *) sp);
575             return (set_errno(EMDB_TGTNOTSUP)); /* XXX ? */
576         }
578         fl = (fl & ~KREG_EFLAGS_IF_MASK) | KREG_EFLAGS_TF_MASK;
580         if (mdb_tgt_vwrite(mdb.m_target, &fl, sizeof (kreg_t),
581             sp) != sizeof (kreg_t)) {
582             warn("failed to update " FLAGS_REG_NAME
583                 " at %p for popfl step\n", (void *) sp);
584             return (set_errno(EMDB_TGTNOTSUP)); /* XXX ? */
585         }

```

```

586         break;
587     }

589     if (emulated) {
590         (void) kmdb_dpi_set_register("pc", npc);
591         return (0);
592     }

594     /* Do the step with IF off, and TF (step) on */
595     (void) kmdb_dpi_get_register(FLAGS_REG_NAME, &oldfl);
596     (void) kmdb_dpi_set_register(FLAGS_REG_NAME,
597         ((oldfl | (1 << KREG_EFLAGS_TF_SHIFT)) & ~KREG_EFLAGS_IF_MASK));

599     kmdb_dpi_resume_master(); /* ... there and back again ... */

601     /* EFLAGS has now changed, and may require tuning */

603     switch (instr) {
604     case M_POPLF:
605         /*
606          * Use the EFLAGS we grabbed before the pop - see the pre-step
607          * M_POPLF comment.
608          */
609         (void) kmdb_dpi_set_register(FLAGS_REG_NAME, newfl);
610         return (0);

612     case M_PUSHF:
613         /*
614          * We pushed our modified EFLAGS (with IF and TF turned off)
615          * onto the stack. Replace the pushed version with our
616          * unmodified one.
617          */
618         (void) kmdb_dpi_get_register("sp", &sp);

620         if (mdb_tgt_vwrite(mdb.m_target, &oldfl, sizeof (kreg_t),
621             sp) != sizeof (kreg_t)) {
622             warn("failed to update pushed " FLAGS_REG_NAME
623                 " at %p after pushfl step\n", (void *)sp);
624             return (set_errno(EMDB_TGTNOTSUP)); /* XXX ? */
625         }

627         /* Go back to using the EFLAGS we were using before the step */
628         (void) kmdb_dpi_set_register(FLAGS_REG_NAME, oldfl);
629         return (0);

631     default:
632         /*
633          * The stepped instruction may have altered EFLAGS. We only
634          * really care about the value of IF, and we know the stepped
635          * instruction didn't alter it, so we can simply copy the
636          * pre-step value. We'll also need to turn TF back off.
637          */
638         (void) kmdb_dpi_get_register(FLAGS_REG_NAME, &fl);
639         (void) kmdb_dpi_set_register(FLAGS_REG_NAME,
640             ((fl & ~(KREG_EFLAGS_TF_MASK|KREG_EFLAGS_IF_MASK)) |
641             (oldfl & KREG_EFLAGS_IF_MASK)));
642         return (0);
643     }
644 }

606 /*
607  * The target has already configured the chip for branch step, leaving us to
608  * actually make the machine go. Due to a number of issues involving
609  * the potential alteration of system state via instructions like sti, cli,
610  * pushfl, and popfl, we're going to treat this like a normal system resume.
611  * All CPUs will be released, on the kernel's IDT. Our primary concern is

```

```

612  * the alteration/storage of our TF'd EFLAGS via pushfl and popfl. There's no
613  * real workaround - we don't have opcode breakpoints - so the best we can do is
614  * to ensure that the world won't end if someone does bad things to EFLAGS.
615  *
616  * Two things can happen:
617  * 1. EFLAGS.TF may be cleared, either maliciously or via a popfl from saved
618  * state. The CPU will continue execution beyond the branch, and will not
619  * reenter the debugger unless brought/sent in by other means.
620  * 2. Someone may pushlf the TF'd EFLAGS, and may stash a copy of it somewhere.
621  * When the saved version is popfl'd back into place, the debugger will be
622  * re-entered on a single-step trap.
623  */
624 static void
625 kaif_step_branch(void)
626 {
627     kreg_t fl;

629     (void) kmdb_dpi_get_register(FLAGS_REG_NAME, &fl);
630     (void) kmdb_dpi_set_register(FLAGS_REG_NAME,
631         (fl | (1 << KREG_EFLAGS_TF_SHIFT)));

633     kmdb_dpi_resume_master();

635     (void) kmdb_dpi_set_register(FLAGS_REG_NAME, fl);
636 }

646 /*ARGSUSED*/
647 static uintptr_t
648 kaif_call(uintptr_t funcva, uint_t argc, const uintptr_t argv[])
649 {
650     return (kaif_invoke(funcva, argc, argv));
651 }
652
653 _____unchanged_portion_omitted_____

727 static void
728 kaif_msr_add(const kdi_msr_t *msrs)
729 {
730     kdi_msr_t *save;
731     size_t nr_msrs = 0;
732     size_t i;

734     while (msrs[nr_msrs].msr_num != 0)
735         nr_msrs++;
736     /* we want to copy the terminating kdi_msr_t too */
737     nr_msrs++;

739     save = mdb_zalloc(sizeof (kdi_msr_t) * nr_msrs * kaif_ncpusave,
740         UM_SLEEP);

742     for (i = 0; i < kaif_ncpusave; i++)
743         bcopy(msrs, &save[nr_msrs * i], sizeof (kdi_msr_t) * nr_msrs);

745     kmdb_kdi_set_debug_msrs(save);
746 }

748 static uint64_t
749 kaif_msr_get(int cpuid, uint_t num)
750 {
751     kdi_cpuseave_t *save;
752     kdi_msr_t *msr;
753     int i;

755     if ((save = kaif_cpuid2save(cpuid)) == NULL)
756         return (-1); /* errno is set for us */

758     msr = save->krs_msr;

```

```
760     for (i = 0; msr[i].msr_num != 0; i++) {
761         if (msr[i].msr_num == num && (msr[i].msr_type & KDI_MSR_READ))
762             return (msr[i].kdi_msr_val);
763     }
764
765     return (0);
766 }
```

```
735 void
736 kaif_trap_set_debugger(void)
737 {
738     kmdb_kdi_idt_switch(NULL);
739 }
```

unchanged_portion_omitted

```
833 dpi_ops_t kmdb_dpi_ops = {
834     kaif_init,
835     kaif_activate,
836     kmdb_kdi_deactivate,
837     kaif_enter_mon,
838     kaif_modchg_register,
839     kaif_modchg_cancel,
840     kaif_get_cpu_state,
841     kaif_get_master_cpuid,
842     kaif_get_gregs,
843     kaif_get_register,
844     kaif_set_register,
845     kaif_brkpt_arm,
846     kaif_brkpt_disarm,
847     kaif_wapt_validate,
848     kaif_wapt_reserve,
849     kaif_wapt_release,
850     kaif_wapt_arm,
851     kaif_wapt_disarm,
852     kaif_wapt_match,
853     kaif_step,
854     kaif_step_branch,
855     kaif_call,
856     kaif_dump_crumbs,
857     kaif_msr_add,
858     kaif_msr_get,
859 };
```

unchanged_portion_omitted

```

*****
3622 Fri Apr 6 17:24:57 2018
new/usr/src/cmd/mdb/intel/kmdb/kmdb_dpi_isadep.c
9210 remove KMDB branch debugging support
9211 ::crregs could do with cr2/cr3 support
9209 ::ttrace should be able to filter by thread
Reviewed by: Patrick Mooney <patrick.mooney@joyent.com>
Reviewed by: Yuri Pankov <yuripv@yuripv.net>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22  * Copyright 2007 Sun Microsystems, Inc. All rights reserved.
23  * Use is subject to license terms.
24  *
25  * Copyright 2018 Joyent, Inc.
26  */
27
28 #pragma ident "%Z%M% %I% %E% SMI"
29
30 /*
31  * Intel-specific portions of the DPI
32  */
33
34 #include <sys/types.h>
35 #include <sys/trap.h>
36
37 #include <kmdb/kmdb_dpi_impl.h>
38 #include <kmdb/kmdb_fault.h>
39 #include <kmdb/kmdb_kdi.h>
40 #include <mdb/mdb_err.h>
41 #include <mdb/mdb_debug.h>
42 #include <mdb/mdb_kreg.h>
43 #include <mdb/mdb.h>
44
45 void
46 kmdb_dpi_handle_fault(kreg_t trapno, kreg_t pc, kreg_t sp, int cpuid)
47 {
48     kmdb_kdi_system_claim();
49
50     mdb_dprintf(MDB_DBG_DPI, "\ndpi_handle_fault: trapno %u, pc 0x%0?p, "
51               "sp 0x%0?p\n", (int)trapno, pc, sp);
52
53     switch (trapno) {
54     case T_GPFLT:
55         errno = EACCES;
56     default:
57         errno = EMDB_NOMAP;
58     }
59 }

```

```

56     }
57
58     if (kmdb_dpi_fault_pcb != NULL) {
59         longjmp(*kmdb_dpi_fault_pcb, 1);
60         /*NOTREACHED*/
61     }
62
63     /* Debugger fault */
64     kmdb_fault(trapno, pc, sp, cpuid);
65 }
66
67 _____ unchanged_portion_omitted _____
68
69 void
70 kmdb_dpi_reboot(void)
71 {
72     /*
73      * We're going to skip all of the niceties we employ in resume_common,
74      * as we don't plan to ever return.
75      */
76     longjmp(kmdb_dpi_entry_pcb, KMDB_DPI_CMD_REBOOT);
77 }
78
79 void
80 kmdb_dpi_msr_add(const kdi_msr_t *msrs)
81 {
82     mdb.m_dpi->dpo_msr_add(msrs);
83 }
84
85 uint64_t
86 kmdb_dpi_msr_get(uint_t msr)
87 {
88     return (mdb.m_dpi->dpo_msr_get(DPI_MASTER_CPUID, msr));
89 }
90
91 uint64_t
92 kmdb_dpi_msr_get_by_cpu(int cpuid, uint_t msr)
93 {
94     return (mdb.m_dpi->dpo_msr_get(cpuid, msr));
95 }
96
97 _____ unchanged_portion_omitted _____

```

new/usr/src/cmd/mdb/intel/kmdb/kmdb_dpi_isadep.h

1

1342 Fri Apr 6 17:24:57 2018

new/usr/src/cmd/mdb/intel/kmdb/kmdb_dpi_isadep.h

9210 remove KMDB branch debugging support

9211 ::crregs could do with cr2/cr3 support

9209 ::ttrace should be able to filter by thread

Reviewed by: Patrick Mooney <patrick.mooney@joyent.com>

Reviewed by: Yuri Pankov <yuripv@yuripv.net>

```
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22  * Copyright 2007 Sun Microsystems, Inc. All rights reserved.
23  * Use is subject to license terms.
24  *
25  * Copyright 2018 Joyent, Inc.
26  */
```

```
28 #ifndef _KMDB_DPI_ISADEP_H
29 #define _KMDB_DPI_ISADEP_H
```

```
29 #pragma ident "%Z%M% %I% %E% SMI"
```

```
31 #ifndef _ASM
32 #include <mdb/mdb_isautil.h>
33 #include <sys/kdi_machimpl.h>
34 #endif
```

```
36 #ifdef __cplusplus
37 extern "C" {
38 #endif
```

```
40 #ifndef _ASM
```

```
42 extern void kmdb_dpi_handle_fault(kreg_t, kreg_t, kreg_t, int);
```

```
44 extern void kmdb_dpi_reboot(void) __NORETURN;
```

```
46 extern void kmdb_dpi_msr_add(const kdi_msr_t *);
47 extern uint64_t kmdb_dpi_msr_get(uint_t);
48 extern uint64_t kmdb_dpi_msr_get_by_cpu(int, uint_t);
```

```
46 #endif /* _ASM */
```

```
48 #ifdef __cplusplus
49 }
```

unchanged_portion_omitted

new/usr/src/cmd/mdb/intel/kmdb/kmdb_kdi_isadep.c

1

2397 Fri Apr 6 17:24:57 2018

new/usr/src/cmd/mdb/intel/kmdb/kmdb_kdi_isadep.c

9210 remove KMDB branch debugging support

9211 ::crregs could do with cr2/cr3 support

9209 ::ttrace should be able to filter by thread

Reviewed by: Patrick Mooney <patrick.mooney@joyent.com>

Reviewed by: Yuri Pankov <yuripv@yuripv.net>

```
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22  * Copyright 2007 Sun Microsystems, Inc. All rights reserved.
23  * Use is subject to license terms.
24  *
25  * Copyright 2018 Joyent, Inc.
26  */
```

```
26 #pragma ident "%Z%M% %I% %E% SMI"
```

```
28 #include <sys/types.h>
29 #include <sys/kdi_impl.h>
30 #include <sys/segments.h>
31 #include <sys/cpuvar.h>
```

```
33 #include <mdb/mdb_debug.h>
34 #include <mdb/mdb_err.h>
35 #include <mdb/mdb_umem.h>
36 #include <kmdb/kmdb_dpi.h>
37 #include <mdb/mdb.h>
```

```
39 /*ARGSUSED*/
40 void
41 kmdb_kdi_stop_slaves(int my_cpuid, int doxc)
42 {
43     /* Stop other CPUs if there are CPUs to stop */
44     mdb.m_kdi->mkdi_stop_slaves(my_cpuid, doxc);
45 }
```

_____unchanged_portion_omitted_____

```
106 void
107 kmdb_kdi_set_debug_msrs(kdi_msr_t *msrs)
108 {
109     mdb.m_kdi->mkdi_set_debug_msrs(msrs);
110 }
```

```
106 void
107 kmdb_kdi_memrange_add(caddr_t base, size_t len)
```

new/usr/src/cmd/mdb/intel/kmdb/kmdb_kdi_isadep.c

2

```
108 {
109     mdb.m_kdi->mkdi_memrange_add(base, len);
110 }
```

_____unchanged_portion_omitted_____

new/usr/src/cmd/mdb/intel/kmdb/kmdb_kdi_isadep.h

1

1627 Fri Apr 6 17:24:57 2018

new/usr/src/cmd/mdb/intel/kmdb/kmdb_kdi_isadep.h

9210 remove KMDB branch debugging support

9211 ::crregs could do with cr2/cr3 support

9209 ::ttrace should be able to filter by thread

Reviewed by: Patrick Mooney <patrick.mooney@joyent.com>

Reviewed by: Yuri Pankov <yuripv@yuripv.net>

```
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22  * Copyright 2007 Sun Microsystems, Inc. All rights reserved.
23  * Use is subject to license terms.
24  *
25  * Copyright 2018 Joyent, Inc.
26  */
```

```
28 #ifndef _KMDB_KDI_ISADEP_H
```

```
29 #define _KMDB_KDI_ISADEP_H
```

```
29 #pragma ident "%Z%M% %I% %E% SMI"
```

```
31 #include <sys/types.h>
```

```
32 #include <sys/kdi_machimpl.h>
```

```
34 #include <mdb/mdb_target.h>
```

```
36 #ifdef __cplusplus
```

```
37 extern "C" {
```

```
38 #endif
```

```
40 struct gate_desc;
```

```
42 extern void kmdb_kdi_activate(kdi_main_t, kdi_cpusave_t *, int);
```

```
43 extern void kmdb_kdi_deactivate(void);
```

```
45 extern void kmdb_kdi_idt_switch(kdi_cpusave_t *);
```

```
47 extern void kmdb_kdi_update_drreg(kdi_drreg_t *);
```

```
48 extern void kmdb_kdi_set_debug_msrs(kdi_msr_t *);
```

```
49 extern uintptr_t kmdb_kdi_get_userlimit(void);
```

```
51 extern int kmdb_kdi_get_cpuinfo(uint_t *, uint_t *, uint_t *);
```

```
53 extern void kmdb_kdi_memrange_add(caddr_t, size_t);
```

new/usr/src/cmd/mdb/intel/kmdb/kmdb_kdi_isadep.h

2

```
55 extern void kmdb_kdi_reboot(void);
```

```
57 #ifdef __cplusplus
```

```
58 }
```

```
_____ unchanged_portion_omitted_
```


new/usr/src/cmd/mdb/intel/kmdb/kvm_isadep.c

1

```
*****
14222 Fri Apr 6 17:24:57 2018
new/usr/src/cmd/mdb/intel/kmdb/kvm_isadep.c
9210 remove KMDB branch debugging support
9211 ::crregs could do with cr2/cr3 support
9209 ::ttrace should be able to filter by thread
Reviewed by: Patrick Mooney <patrick.mooney@joyent.com>
Reviewed by: Yuri Pankov <yuripv@yuripv.net>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22  * Copyright 2007 Sun Microsystems, Inc. All rights reserved.
23  * Use is subject to license terms.
24  *
25  * Copyright 2018 Joyent, Inc.
26  */
27
28 #pragma ident "%Z%%M% %I% %E% SMI"
29
30 /*
31  * isa-dependent portions of the kmdb target
32  */
33
34 #include <kmdb/kvm.h>
35 #include <kmdb/kvm_cpu.h>
36 #include <kmdb/kmdb_kdi.h>
37 #include <kmdb/kmdb_asmutil.h>
38 #include <mdb/mdb_debug.h>
39 #include <mdb/mdb_err.h>
40 #include <mdb/mdb_list.h>
41 #include <mdb/mdb_target_impl.h>
42 #include <mdb/mdb_isautil.h>
43 #include <mdb/mdb_kreg_impl.h>
44 #include <mdb/mdb.h>
45
46 #include <sys/types.h>
47 #include <sys/frame.h>
48 #include <sys/trap.h>
49 #include <sys/bitmap.h>
50 #include <sys/pci_impl.h>
51
52 /* Higher than the highest trap number for which we have a defined specifier */
53 #define KMT_MAXTRAPNO 0x20
54
55 #define IOPORTLIMIT 0xffff /* XXX find a new home for this */
56
57 const char *
```

new/usr/src/cmd/mdb/intel/kmdb/kvm_isadep.c

2

```
55 kmt_def_dismode(void)
56 {
57 #ifdef __amd64
58     return ("amd64");
59 #else
60     return ("ia32");
61 #endif
62 }
63
64 unchanged_portion_omitted_
65
66 104 int
67 105 kmt_step_branch(mdb_tgt_t *t)
68 106 {
69 107     kmt_data_t *kmt = t->t_data;
70
71 109     return (kmt_cpu_step_branch(t, kmt->kmt_cpu));
72 110 }
73
74 103 /*
75 104  * Return the address of the next instruction following a call, or return -1
76 105  * and set errno to EAGAIN if the target should just single-step.
77 106  */
78 107 int
79 108 kmt_next(mdb_tgt_t *t, uintptr_t *p)
80 109 {
81 110     kreg_t pc;
82 111     mdb_instr_t instr;
83
84 113     (void) kmdb_dpi_get_register("pc", &pc);
85
86 115     if (mdb_tgt_vread(t, &instr, sizeof (mdb_instr_t), pc) !=
87 116         sizeof (mdb_instr_t))
88 117         return (-1); /* errno is set for us */
89
90 119     return (mdb_isa_next(t, p, pc, instr));
91 120 }
92
93 unchanged_portion_omitted_
94
95 359 int
96 360 kmt_msr_validate(const kdi_msr_t *msr)
97 361 {
98 362     uint64_t val;
99
100 364     for (/* */; msr->msr_num != 0; msr++) {
101 365         if (kmt_rwmsr(msr->msr_num, &val, rdmsr) < 0)
102 366             return (0);
103 367     }
104
105 369     return (1);
106 370 }
107
108 350 /*ARGSUSED*/
109 351 ssize_t
110 352 kmt_write(mdb_tgt_t *t, const void *buf, size_t nbytes, uintptr_t addr)
111 353 {
112 354     if (!(t->t_flags & MDB_TGT_F_ALLOWIO) &&
113 355         (nbytes = kmdb_kdi_range_is_nontoxic(addr, nbytes, 1)) == 0)
114 356         return (set_errno(EMDB_NOMAP));
115
116 358     /*
117 359     * No writes to user space are allowed. If we were to allow it, we'd
118 360     * be in the unfortunate situation where kmdb could place a breakpoint
119 361     * on a userspace executable page; this dirty page would end up being
120 362     * flushed back to disk, incurring sadness when it's next executed.
121 363     * Besides, we can't allow trapping in from userspace anyway.
122 364     */
```

new/usr/src/cmd/mdb/intel/kmdb/kvm_isadep.c

3

```
365     if (addr < kmdb_kdi_get_userlimit())
366         return (set_errno(EMDB_TGTNOTSUP));
368     return (kmt_rw(t, (void *)buf, nbytes, addr, kmt_writer));
369 }
_____unchanged_portion_omitted_____
```

new/usr/src/cmd/mdb/intel/kmdb/kvm_isadep.h

1

1670 Fri Apr 6 17:24:58 2018

new/usr/src/cmd/mdb/intel/kmdb/kvm_isadep.h

9210 remove KMDB branch debugging support

9211 ::crregs could do with cr2/cr3 support

9209 ::ttrace should be able to filter by thread

Reviewed by: Patrick Mooney <patrick.mooney@joyent.com>

Reviewed by: Yuri Pankov <yuripv@yuripv.net>

```
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22  * Copyright 2007 Sun Microsystems, Inc. All rights reserved.
23  * Use is subject to license terms.
24  *
25  * Copyright 2018 Joyent, Inc.
26  */

28 #ifndef _KVM_ISADEP_H
29 #define _KVM_ISADEP_H

31 #pragma ident "%Z%%M% %I% %E% SMI"

32 extern "C" {
33 #endif

35 extern uintptr_t kmt_invoke(uintptr_t, uint_t, const uintptr_t *);

37 extern void kmt_in(void *, size_t, uintptr_t);
38 extern void kmt_out(void *, size_t, uintptr_t);

40 extern int kmt_in_dcnd(uintptr_t, uint_t, int, const mdb_arg_t *);
41 extern int kmt_out_dcnd(uintptr_t, uint_t, int, const mdb_arg_t *);

43 extern int kmt_rdmsr(uintptr_t, uint_t, int, const mdb_arg_t *);
44 extern int kmt_wrmsr(uintptr_t, uint_t, int, const mdb_arg_t *);

46 extern int kmt_rdpcicfg(uintptr_t, uint_t, int, const mdb_arg_t *);
47 extern int kmt_wrpcicfg(uintptr_t, uint_t, int, const mdb_arg_t *);
48 extern int kmt_msr_validate(const kdi_msr_t *);

49 #ifdef __cplusplus
50 }
_____unchanged_portion_omitted_
```

```

*****
9127 Fri Apr 6 17:24:58 2018
new/usr/src/cmd/mdb/intel/mdb/kvm_amd64dep.c
9210 remove KMDB branch debugging support
9211 ::crregs could do with cr2/cr3 support
9209 ::ttrace should be able to filter by thread
Reviewed by: Patrick Mooney <patrick.mooney@joyent.com>
Reviewed by: Yuri Pankov <yuripv@yuripv.net>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22  * Copyright 2008 Sun Microsystems, Inc. All rights reserved.
23  * Use is subject to license terms.
24  *
25  * Copyright 2018 Joyent, Inc.
26  */

26 #pragma ident "%Z%M %I %E% SMI"

28 /*
29  * Libkvm Kernel Target Intel 64-bit component
30  *
31  * This file provides the ISA-dependent portion of the libkvm kernel target.
32  * For more details on the implementation refer to mdb_kvm.c.
33  */

35 #include <sys/types.h>
36 #include <sys/reg.h>
37 #include <sys/frame.h>
38 #include <sys/stack.h>
39 #include <sys/sysmacros.h>
40 #include <sys/panic.h>
41 #include <sys/privregs.h>
42 #include <strings.h>

44 #include <mdb/mdb_target_impl.h>
45 #include <mdb/mdb_disasm.h>
46 #include <mdb/mdb_modapi.h>
47 #include <mdb/mdb_conf.h>
48 #include <mdb/mdb_kreg_impl.h>
49 #include <mdb/mdb_amd64util.h>
50 #include <mdb/kvm_isadep.h>
51 #include <mdb/mdb_kvm.h>
52 #include <mdb/mdb_err.h>
53 #include <mdb/mdb_debug.h>
54 #include <mdb/mdb.h>

```

```

56 /*ARGUSED*/
57 int
58 kt_regs(uintptr_t addr, uint_t flags, int argc, const mdb_arg_t *argv)
59 {
60     mdb_amd64_printregs((const mdb_tgt_gregset_t *)addr);
61     return (DCMD_OK);
62 }

unchanged_portion_omitted

105 const mdb_tgt_ops_t kt_amd64_ops = {
106     kt_setflags,           /* t_setflags */
107     kt_setcontext,       /* t_setcontext */
108     kt_activate,         /* t_activate */
109     kt_deactivate,       /* t_deactivate */
110     (void (*)(void)) mdb_tgt_nop, /* t_periodic */
111     kt_destroy,          /* t_destroy */
112     kt_name,             /* t_name */
113     (const char (*)(void)) mdb_conf_isa, /* t_isa */
114     kt_platform,         /* t_platform */
115     kt_uname,            /* t_uname */
116     kt_dmodel,           /* t_dmodel */
117     kt_aread,            /* t_aread */
118     kt_awrite,           /* t_awrite */
119     kt_vread,            /* t_vread */
120     kt_vwrite,           /* t_vwrite */
121     kt_pread,            /* t_pread */
122     kt_pwrite,           /* t_pwrite */
123     kt_fread,            /* t_fread */
124     kt_fwrite,           /* t_fwrite */
125     (ssize_t (*)(void)) mdb_tgt_notsup, /* t_iorread */
126     (ssize_t (*)(void)) mdb_tgt_notsup, /* t_iowrite */
127     kt_vtop,             /* t_vtop */
128     kt_lookup_by_name,   /* t_lookup_by_name */
129     kt_lookup_by_addr,   /* t_lookup_by_addr */
130     kt_symbol_iter,      /* t_symbol_iter */
131     kt_mapping_iter,     /* t_mapping_iter */
132     kt_object_iter,      /* t_object_iter */
133     kt_addr_to_map,      /* t_addr_to_map */
134     kt_name_to_map,      /* t_name_to_map */
135     kt_addr_to_ctf,      /* t_addr_to_ctf */
136     kt_name_to_ctf,      /* t_name_to_ctf */
137     kt_status,           /* t_status */
138     (int (*)(void)) mdb_tgt_notsup, /* t_run */
139     (int (*)(void)) mdb_tgt_notsup, /* t_step */
140     (int (*)(void)) mdb_tgt_notsup, /* t_step_out */
141     (int (*)(void)) mdb_tgt_notsup, /* t_step_branch */
142     (int (*)(void)) mdb_tgt_notsup, /* t_next */
143     (int (*)(void)) mdb_tgt_notsup, /* t_cont */
144     (int (*)(void)) mdb_tgt_null, /* t_signal */
145     (int (*)(void)) mdb_tgt_null, /* t_add_vbrkpt */
146     (int (*)(void)) mdb_tgt_null, /* t_add_sbrkpt */
147     (int (*)(void)) mdb_tgt_null, /* t_add_pwapt */
148     (int (*)(void)) mdb_tgt_null, /* t_add_vwapt */
149     (int (*)(void)) mdb_tgt_null, /* t_add_iowapt */
150     (int (*)(void)) mdb_tgt_null, /* t_add_sysenter */
151     (int (*)(void)) mdb_tgt_null, /* t_add_sysexit */
152     (int (*)(void)) mdb_tgt_null, /* t_add_signal */
153     (int (*)(void)) mdb_tgt_null, /* t_add_fault */
154     kt_getareg,          /* t_getareg */
155     kt_putareg,          /* t_putareg */
156     mdb_amd64_kvm_stack_iter, /* t_stack_iter */
157     (int (*)(void)) mdb_tgt_notsup /* t_auxv */
};

unchanged_portion_omitted

```

```

*****
9471 Fri Apr 6 17:24:58 2018
new/usr/src/cmd/mdb/intel/mdb/kvm_ia32dep.c
9210 remove KMDB branch debugging support
9211 ::crregs could do with cr2/cr3 support
9209 ::ttrace should be able to filter by thread
Reviewed by: Patrick Mooney <patrick.mooney@joyent.com>
Reviewed by: Yuri Pankov <yuripv@yuripv.net>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22  * Copyright 2008 Sun Microsystems, Inc. All rights reserved.
23  * Use is subject to license terms.
24  *
25  * Copyright 2018 Joyent, Inc.
26  */

26 #pragma ident "%Z%M %I %E% SMI"

28 /*
29  * Libkvm Kernel Target Intel 32-bit component
30  *
31  * This file provides the ISA-dependent portion of the libkvm kernel target.
32  * For more details on the implementation refer to mdb_kvm.c.
33  */

35 #include <sys/types.h>
36 #include <sys/regset.h>
37 #include <sys/frame.h>
38 #include <sys/stack.h>
39 #include <sys/sysmacros.h>
40 #include <sys/panic.h>
41 #include <strings.h>

43 #include <mdb/mdb_target_impl.h>
44 #include <mdb/mdb_disasm.h>
45 #include <mdb/mdb_modapi.h>
46 #include <mdb/mdb_conf.h>
47 #include <mdb/mdb_kreg_impl.h>
48 #include <mdb/mdb_ia32util.h>
49 #include <mdb/kvm_isadep.h>
50 #include <mdb/mdb_kvm.h>
51 #include <mdb/mdb_err.h>
52 #include <mdb/mdb_debug.h>
53 #include <mdb/mdb.h>

```

```

56 /*ARGUSED*/
57 int
58 kt_regs(uintptr_t addr, uint_t flags, int argc, const mdb_arg_t *argv)
59 {
60     mdb_ia32_printregs((const mdb_tgt_gregset_t *)addr);
61     return (DCMD_OK);
62 }

unchanged_portion_omitted

105 const mdb_tgt_ops_t kt_ia32_ops = {
106     kt_setflags,           /* t_setflags */
107     kt_setcontext,       /* t_setcontext */
108     kt_activate,         /* t_activate */
109     kt_deactivate,       /* t_deactivate */
110     (void (*)( )) mdb_tgt_nop, /* t_periodic */
111     kt_destroy,          /* t_destroy */
112     kt_name,             /* t_name */
113     (const char (*)( )) mdb_conf_isa, /* t_isa */
114     kt_platform,         /* t_platform */
115     kt_uname,            /* t_uname */
116     kt_dmodel,           /* t_dmodel */
117     kt_aread,            /* t_aread */
118     kt_awrite,           /* t_awrite */
119     kt_vread,            /* t_vread */
120     kt_vwrite,           /* t_vwrite */
121     kt_pread,            /* t_pread */
122     kt_pwrite,           /* t_pwrite */
123     kt_fread,            /* t_fread */
124     kt_fwrite,           /* t_fwrite */
125     (ssize_t (*)( )) mdb_tgt_notsup, /* t_ioroad */
126     (ssize_t (*)( )) mdb_tgt_notsup, /* t_iowrite */
127     kt_vtop,             /* t_vtop */
128     kt_lookup_by_name,   /* t_lookup_by_name */
129     kt_lookup_by_addr,   /* t_lookup_by_addr */
130     kt_symbol_iter,      /* t_symbol_iter */
131     kt_mapping_iter,     /* t_mapping_iter */
132     kt_object_iter,      /* t_object_iter */
133     kt_addr_to_map,       /* t_addr_to_map */
134     kt_name_to_map,       /* t_name_to_map */
135     kt_addr_to_ctf,       /* t_addr_to_ctf */
136     kt_name_to_ctf,       /* t_name_to_ctf */
137     kt_status,           /* t_status */
138     (int (*)( )) mdb_tgt_notsup, /* t_run */
139     (int (*)( )) mdb_tgt_notsup, /* t_step */
140     (int (*)( )) mdb_tgt_notsup, /* t_step_out */
141     (int (*)( )) mdb_tgt_notsup, /* t_step_branch */
142     (int (*)( )) mdb_tgt_notsup, /* t_next */
143     (int (*)( )) mdb_tgt_notsup, /* t_cont */
144     (int (*)( )) mdb_tgt_notsup, /* t_signal */
145     (int (*)( )) mdb_tgt_null, /* t_add_vbrkpt */
146     (int (*)( )) mdb_tgt_null, /* t_add_sbrkpt */
147     (int (*)( )) mdb_tgt_null, /* t_add_pwapt */
148     (int (*)( )) mdb_tgt_null, /* t_add_vwapt */
149     (int (*)( )) mdb_tgt_null, /* t_add_iowapt */
150     (int (*)( )) mdb_tgt_null, /* t_add_sysenter */
151     (int (*)( )) mdb_tgt_null, /* t_add_sysexit */
152     (int (*)( )) mdb_tgt_null, /* t_add_signal */
153     (int (*)( )) mdb_tgt_null, /* t_add_fault */
154     kt_getareg,          /* t_getareg */
155     kt_putareg,          /* t_putareg */
156     mdb_ia32_kvm_stack_iter, /* t_stack_iter */
157     (int (*)( )) mdb_tgt_notsup, /* t_auxv */
158 };

unchanged_portion_omitted

```

```

*****
16759 Fri Apr 6 17:24:58 2018
new/usr/src/cmd/mdb/intel/mdb/mdb_amd64util.c
9441 kmdb should stash %cr3 in kdiregs
Reviewed by: John Levon <john.levon@joyent.com>
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */

22 /*
23  * Copyright 2007 Sun Microsystems, Inc. All rights reserved.
24  * Use is subject to license terms.
25 */
26 /*
27  * Copyright (c) 2018, Joyent, Inc. All rights reserved.
27  * Copyright (c) 2012, Joyent, Inc. All rights reserved.
28  * Copyright 2014 Nexenta Systems, Inc. All rights reserved.
29  * Copyright (c) 2013 by Delphix. All rights reserved.
30 */

32 #include <sys/types.h>
33 #include <sys/reg.h>
34 #include <sys/privregs.h>
35 #include <sys/stack.h>
36 #include <sys/frame.h>

38 #include <mdb/mdb_target_impl.h>
39 #include <mdb/mdb_kreg_impl.h>
40 #include <mdb/mdb_debug.h>
41 #include <mdb/mdb_modapi.h>
42 #include <mdb/mdb_amd64util.h>
43 #include <mdb/mdb_ctf.h>
44 #include <mdb/mdb_err.h>
45 #include <mdb/mdb.h>

47 #include <saveargs.h>

49 /*
50  * This array is used by the getareg and putareg entry points, and also by our
51  * register variable discipline.
52 */

54 const mdb_tgt_regdesc_t mdb_amd64_kregs[] = {
55     { "savfp", KREG_SAVFP, MDB_TGT_R_EXPORT },
56     { "savpc", KREG_SAVPC, MDB_TGT_R_EXPORT },

```

```

57     { "rdi", KREG_RDI, MDB_TGT_R_EXPORT },
58     { "edi", KREG_RDI, MDB_TGT_R_EXPORT }, MDB_TGT_R_32 },
59     { "di", KREG_RDI, MDB_TGT_R_EXPORT }, MDB_TGT_R_16 },
60     { "dil", KREG_RDI, MDB_TGT_R_EXPORT }, MDB_TGT_R_8L },
61     { "rsi", KREG_RSI, MDB_TGT_R_EXPORT },
62     { "esi", KREG_RSI, MDB_TGT_R_EXPORT }, MDB_TGT_R_32 },
63     { "si", KREG_RSI, MDB_TGT_R_EXPORT }, MDB_TGT_R_16 },
64     { "sil", KREG_RSI, MDB_TGT_R_EXPORT }, MDB_TGT_R_8L },
65     { "rdx", KREG_RDX, MDB_TGT_R_EXPORT },
66     { "edx", KREG_RDX, MDB_TGT_R_EXPORT }, MDB_TGT_R_32 },
67     { "dx", KREG_RDX, MDB_TGT_R_EXPORT }, MDB_TGT_R_16 },
68     { "dh", KREG_RDX, MDB_TGT_R_EXPORT }, MDB_TGT_R_8H },
69     { "dl", KREG_RDX, MDB_TGT_R_EXPORT }, MDB_TGT_R_8L },
70     { "rcx", KREG_RCX, MDB_TGT_R_EXPORT },
71     { "ecx", KREG_RCX, MDB_TGT_R_EXPORT }, MDB_TGT_R_32 },
72     { "cx", KREG_RCX, MDB_TGT_R_EXPORT }, MDB_TGT_R_16 },
73     { "ch", KREG_RCX, MDB_TGT_R_EXPORT }, MDB_TGT_R_8H },
74     { "cl", KREG_RCX, MDB_TGT_R_EXPORT }, MDB_TGT_R_8L },
75     { "r8", KREG_R8, MDB_TGT_R_EXPORT },
76     { "r8d", KREG_R8, MDB_TGT_R_EXPORT }, MDB_TGT_R_32 },
77     { "r8w", KREG_R8, MDB_TGT_R_EXPORT }, MDB_TGT_R_16 },
78     { "r8l", KREG_R8, MDB_TGT_R_EXPORT }, MDB_TGT_R_8L },
79     { "r9", KREG_R9, MDB_TGT_R_EXPORT },
80     { "r9d", KREG_R8, MDB_TGT_R_EXPORT }, MDB_TGT_R_32 },
81     { "r9w", KREG_R8, MDB_TGT_R_EXPORT }, MDB_TGT_R_16 },
82     { "r9l", KREG_R8, MDB_TGT_R_EXPORT }, MDB_TGT_R_8L },
83     { "rax", KREG_RAX, MDB_TGT_R_EXPORT },
84     { "eax", KREG_RAX, MDB_TGT_R_EXPORT }, MDB_TGT_R_32 },
85     { "ax", KREG_RAX, MDB_TGT_R_EXPORT }, MDB_TGT_R_16 },
86     { "ah", KREG_RAX, MDB_TGT_R_EXPORT }, MDB_TGT_R_8H },
87     { "al", KREG_RAX, MDB_TGT_R_EXPORT }, MDB_TGT_R_8L },
88     { "rbx", KREG_RBX, MDB_TGT_R_EXPORT },
89     { "ebx", KREG_RBX, MDB_TGT_R_EXPORT }, MDB_TGT_R_32 },
90     { "bx", KREG_RBX, MDB_TGT_R_EXPORT }, MDB_TGT_R_16 },
91     { "bh", KREG_RBX, MDB_TGT_R_EXPORT }, MDB_TGT_R_8H },
92     { "bl", KREG_RBX, MDB_TGT_R_EXPORT }, MDB_TGT_R_8L },
93     { "rbp", KREG_RBP, MDB_TGT_R_EXPORT },
94     { "ebp", KREG_RBP, MDB_TGT_R_EXPORT }, MDB_TGT_R_32 },
95     { "bp", KREG_RBP, MDB_TGT_R_EXPORT }, MDB_TGT_R_16 },
96     { "bpl", KREG_RBP, MDB_TGT_R_EXPORT }, MDB_TGT_R_8L },
97     { "r10", KREG_R10, MDB_TGT_R_EXPORT },
98     { "r10d", KREG_R10, MDB_TGT_R_EXPORT }, MDB_TGT_R_32 },
99     { "r10w", KREG_R10, MDB_TGT_R_EXPORT }, MDB_TGT_R_16 },
100    { "r10l", KREG_R10, MDB_TGT_R_EXPORT }, MDB_TGT_R_8L },
101    { "r11", KREG_R11, MDB_TGT_R_EXPORT },
102    { "r11d", KREG_R11, MDB_TGT_R_EXPORT }, MDB_TGT_R_32 },
103    { "r11w", KREG_R11, MDB_TGT_R_EXPORT }, MDB_TGT_R_16 },
104    { "r11l", KREG_R11, MDB_TGT_R_EXPORT }, MDB_TGT_R_8L },
105    { "r12", KREG_R12, MDB_TGT_R_EXPORT },
106    { "r12d", KREG_R12, MDB_TGT_R_EXPORT }, MDB_TGT_R_32 },
107    { "r12w", KREG_R12, MDB_TGT_R_EXPORT }, MDB_TGT_R_16 },
108    { "r12l", KREG_R12, MDB_TGT_R_EXPORT }, MDB_TGT_R_8L },
109    { "r13", KREG_R13, MDB_TGT_R_EXPORT },
110    { "r13d", KREG_R13, MDB_TGT_R_EXPORT }, MDB_TGT_R_32 },
111    { "r13w", KREG_R13, MDB_TGT_R_EXPORT }, MDB_TGT_R_16 },
112    { "r13l", KREG_R13, MDB_TGT_R_EXPORT }, MDB_TGT_R_8L },
113    { "r14", KREG_R14, MDB_TGT_R_EXPORT },
114    { "r14d", KREG_R14, MDB_TGT_R_EXPORT }, MDB_TGT_R_32 },
115    { "r14w", KREG_R14, MDB_TGT_R_EXPORT }, MDB_TGT_R_16 },
116    { "r14l", KREG_R14, MDB_TGT_R_EXPORT }, MDB_TGT_R_8L },
117    { "r15", KREG_R15, MDB_TGT_R_EXPORT },
118    { "r15d", KREG_R15, MDB_TGT_R_EXPORT }, MDB_TGT_R_32 },
119    { "r15w", KREG_R15, MDB_TGT_R_EXPORT }, MDB_TGT_R_16 },
120    { "r15l", KREG_R15, MDB_TGT_R_EXPORT }, MDB_TGT_R_8L },
121    { "ds", KREG_DS, MDB_TGT_R_EXPORT },
122    { "es", KREG_ES, MDB_TGT_R_EXPORT },

```

```

123     { "fs", KREG_FS, MDB_TGT_R_EXPORT },
124     { "gs", KREG_GS, MDB_TGT_R_EXPORT },
125     { "trapno", KREG_TRAPNO, MDB_TGT_R_EXPORT | MDB_TGT_R_PRIV },
126     { "err", KREG_ERR, MDB_TGT_R_EXPORT | MDB_TGT_R_PRIV },
127     { "rip", KREG_RIP, MDB_TGT_R_EXPORT },
128     { "cs", KREG_CS, MDB_TGT_R_EXPORT },
129     { "rflags", KREG_RFLAGS, MDB_TGT_R_EXPORT },
130     { "eflags", KREG_EFLAGS, MDB_TGT_R_EXPORT | MDB_TGT_R_32 },
131     { "rsp", KREG_RSP, MDB_TGT_R_EXPORT },
132     { "esp", KREG_RSP, MDB_TGT_R_EXPORT | MDB_TGT_R_32 },
133     { "sp", KREG_RSP, MDB_TGT_R_EXPORT | MDB_TGT_R_16 },
134     { "spl", KREG_RSP, MDB_TGT_R_EXPORT | MDB_TGT_R_8L },
135     { "ss", KREG_SS, MDB_TGT_R_EXPORT },
136     { "gsbase", KREG_GSBASE, MDB_TGT_R_EXPORT },
137     { "kgsbase", KREG_KGSBASE, MDB_TGT_R_EXPORT },
138     { "cr2", KREG_CR2, MDB_TGT_R_EXPORT },
139     { "cr3", KREG_CR3, MDB_TGT_R_EXPORT },
140     { NULL, 0, 0 }
141 };

143 void
144 mdb_amd64_printregs(const mdb_tgt_gregset_t *gregs)
145 {
146     const kreg_t *kregs = &gregs->kregs[0];
147     kreg_t rflags = kregs[KREG_RFLAGS];

149 #define GETREG2(x) ((uintptr_t)kregs[(x)], ((uintptr_t)kregs[(x)])

151     mdb_printf("%rax = 0x%0?p %15A %r9 = 0x%0?p %A\n",
152         GETREG2(KREG_RAX), GETREG2(KREG_R9));
153     mdb_printf("%rbx = 0x%0?p %15A %r10 = 0x%0?p %A\n",
154         GETREG2(KREG_RBX), GETREG2(KREG_R10));
155     mdb_printf("%rcx = 0x%0?p %15A %r11 = 0x%0?p %A\n",
156         GETREG2(KREG_RCX), GETREG2(KREG_R11));
157     mdb_printf("%rdx = 0x%0?p %15A %r12 = 0x%0?p %A\n",
158         GETREG2(KREG_RDX), GETREG2(KREG_R12));
159     mdb_printf("%rsi = 0x%0?p %15A %r13 = 0x%0?p %A\n",
160         GETREG2(KREG_RSI), GETREG2(KREG_R13));
161     mdb_printf("%rdi = 0x%0?p %15A %r14 = 0x%0?p %A\n",
162         GETREG2(KREG_RDI), GETREG2(KREG_R14));
163     mdb_printf("%r8 = 0x%0?p %15A %r15 = 0x%0?p %A\n\n",
164         GETREG2(KREG_R8), GETREG2(KREG_R15));

166     mdb_printf("%rip = 0x%0?p %A\n", GETREG2(KREG_RIP));
167     mdb_printf("%rbp = 0x%0?p\n", kregs[KREG_RBP]);
168     mdb_printf("%rsp = 0x%0?p\n", kregs[KREG_RSP]);

170     mdb_printf("%rflags = 0x%08x\n", rflags);

172     mdb_printf(" id=%u vip=%u vif=%u ac=%u vm=%u rf=%u nt=%u iopl=0x%x\n",
173         (rflags & KREG_EFLAGS_ID_MASK) >> KREG_EFLAGS_ID_SHIFT,
174         (rflags & KREG_EFLAGS_VIP_MASK) >> KREG_EFLAGS_VIP_SHIFT,
175         (rflags & KREG_EFLAGS_VIF_MASK) >> KREG_EFLAGS_VIF_SHIFT,
176         (rflags & KREG_EFLAGS_AC_MASK) >> KREG_EFLAGS_AC_SHIFT,
177         (rflags & KREG_EFLAGS_VM_MASK) >> KREG_EFLAGS_VM_SHIFT,
178         (rflags & KREG_EFLAGS_RF_MASK) >> KREG_EFLAGS_RF_SHIFT,
179         (rflags & KREG_EFLAGS_NT_MASK) >> KREG_EFLAGS_NT_SHIFT,
180         (rflags & KREG_EFLAGS_IOPL_MASK) >> KREG_EFLAGS_IOPL_SHIFT);

182     mdb_printf(" status=<%s,%s,%s,%s,%s,%s,%s,%s,%s>\n\n",
183         (rflags & KREG_EFLAGS_OF_MASK) ? "OF" : "of",
184         (rflags & KREG_EFLAGS_DF_MASK) ? "DF" : "df",
185         (rflags & KREG_EFLAGS_IF_MASK) ? "IF" : "if",
186         (rflags & KREG_EFLAGS_TF_MASK) ? "TF" : "tf",
187         (rflags & KREG_EFLAGS_SF_MASK) ? "SF" : "sf",
188         (rflags & KREG_EFLAGS_ZF_MASK) ? "ZF" : "zf",

```

```

189         (rflags & KREG_EFLAGS_AF_MASK) ? "AF" : "af",
190         (rflags & KREG_EFLAGS_PF_MASK) ? "PF" : "pf",
191         (rflags & KREG_EFLAGS_CF_MASK) ? "CF" : "cf");

193     mdb_printf("%cs = 0x%04x\t%ds = 0x%04x\t"
194         "%es = 0x%04x\t%fs = 0x%04x\n", kregs[KREG_CS], kregs[KREG_DS],
195         kregs[KREG_ES], kregs[KREG_FS] & 0xffff);
196     mdb_printf("%gs = 0x%04x\t%gsbase = 0x%lx\t%kgsbase = 0x%lx\n",
197         kregs[KREG_GS] & 0xffff, kregs[KREG_GSBASE], kregs[KREG_KGSBASE]);
198     mdb_printf("%trapno = 0x%x\t%err = 0x%x\t%cr2 = 0x%lx\t"
199         "%cr3 = 0x%lx\n", kregs[KREG_TRAPNO], kregs[KREG_ERR],
200         kregs[KREG_CR2], kregs[KREG_CR3]);
201     mdb_printf("%24s%cs = 0x%04x\t%ds = 0x%04x\t%es = 0x%04x\n",
202         " ", kregs[KREG_CS], kregs[KREG_DS], kregs[KREG_ES]);

192     mdb_printf("%trapno = 0x%x\t%fs = 0x%04x\t%gs = 0x%04x\n",
193         kregs[KREG_TRAPNO], (kregs[KREG_FS] & 0xffff),
194         (kregs[KREG_GS] & 0xffff));
195     mdb_printf(" %err = 0x%x\n", kregs[KREG_ERR]);
201 }
    _____unchanged_portion_omitted_____

```

new/usr/src/cmd/mdb/intel/mdb/mdb_kreg.h

1

```
*****
5339 Fri Apr 6 17:24:58 2018
new/usr/src/cmd/mdb/intel/mdb/mdb_kreg.h
9441 kmdb should stash %cr3 in kdiregs
Reviewed by: John Levon <john.levon@joyent.com>
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22  * Copyright 2007 Sun Microsystems, Inc. All rights reserved.
23  * Use is subject to license terms.
24  *
25  * Copyright 2018 Joyent, Inc.
26  */

28 #ifndef _MDB_KREG_H
29 #define _MDB_KREG_H

29 #pragma ident "%Z%M% %I% %E% SMI"

31 #include <sys/kdi_regs.h>
32 #ifndef _ASM
33 #include <sys/types.h>
34 #endif

36 #ifdef __cplusplus
37 extern "C" {
38 #endif

40 #ifndef _ASM
41 #ifndef __amd64
42 typedef uint64_t kreg_t;
43 #else /* __amd64 */
44 typedef uint32_t kreg_t;
45 #endif /* __amd64 */
46 #endif /* !_ASM */

48 #define KREG_NGREG      KDIREG_NGREG

50 /*
51  * The order of these registers corresponds to a slightly altered struct regs,
52  * in the order kmdb entry pushes onto the stack.
53  */

55 #ifdef __amd64
```

new/usr/src/cmd/mdb/intel/mdb/mdb_kreg.h

2

```
57 #define KREG_SAVFP      KDIREG_SAVFP
58 #define KREG_SAVPC      KDIREG_SAVPC
59 #define KREG_RDI        KDIREG_RDI
60 #define KREG_RSI        KDIREG_RSI
61 #define KREG_RDX        KDIREG_RDX
62 #define KREG_RCX        KDIREG_RCX
63 #define KREG_R8         KDIREG_R8
64 #define KREG_R9         KDIREG_R9
65 #define KREG_RAX        KDIREG_RAX
66 #define KREG_RBX        KDIREG_RBX
67 #define KREG_RBP        KDIREG_RBP
68 #define KREG_R10        KDIREG_R10
69 #define KREG_R11        KDIREG_R11
70 #define KREG_R12        KDIREG_R12
71 #define KREG_R13        KDIREG_R13
72 #define KREG_R14        KDIREG_R14
73 #define KREG_R15        KDIREG_R15
74 #define KREG_DS         KDIREG_DS
75 #define KREG_ES         KDIREG_ES
76 #define KREG_FS         KDIREG_FS
77 #define KREG_GS         KDIREG_GS
78 #define KREG_GSBASE     KDIREG_GSBASE
79 #define KREG_KGSBASE    KDIREG_KGSBASE
80 #define KREG_TRAPNO     KDIREG_TRAPNO
81 #define KREG_ERR        KDIREG_ERR
82 #define KREG_CR2        KDIREG_CR2
83 #define KREG_CR3        KDIREG_CR3
84 #define KREG_RIP        KDIREG_RIP
85 #define KREG_CS         KDIREG_CS
86 #define KREG_RFLAGS     KDIREG_RFLAGS
87 #define KREG_RSP        KDIREG_RSP
88 #define KREG_SS         KDIREG_SS

90 #define KREG_PC         KREG_RIP
91 #define KREG_SP         KREG_RSP
92 #define KREG_FP         KREG_RBP

94 #else /* __amd64 */

96 #define KREG_SAVFP      KDIREG_SAVFP
97 #define KREG_SAVPC      KDIREG_SAVPC
98 #define KREG_SS         KDIREG_SS
99 #define KREG_GS         KDIREG_GS
100 #define KREG_FS         KDIREG_FS
101 #define KREG_ES         KDIREG_ES
102 #define KREG_DS         KDIREG_DS
103 #define KREG_EDI        KDIREG_EDI
104 #define KREG_ESI        KDIREG_ESI
105 #define KREG_EBP        KDIREG_EBP
106 #define KREG_ESP        KDIREG_ESP
107 #define KREG_EBX        KDIREG_EBX
108 #define KREG_EDX        KDIREG_EDX
109 #define KREG_ECX        KDIREG_ECX
110 #define KREG_EAX        KDIREG_EAX
111 #define KREG_TRAPNO     KDIREG_TRAPNO
112 #define KREG_ERR        KDIREG_ERR
113 #define KREG_EIP        KDIREG_EIP
114 #define KREG_CS         KDIREG_CS
115 #define KREG_EFLAGS     KDIREG_EFLAGS
116 #define KREG_UESP       KDIREG_UESP

118 #define KREG_PC         KREG_EIP
119 #define KREG_SP         KREG_ESP
120 #define KREG_FP         KREG_EBP
```



```

122 #endif /* __amd64 */
124 #define KREG_EFLAGS_ID_MASK      0x00200000
125 #define KREG_EFLAGS_ID_SHIFT    21
127 #define KREG_EFLAGS_VIP_MASK     0x00100000
128 #define KREG_EFLAGS_VIP_SHIFT   20
130 #define KREG_EFLAGS_VIF_MASK     0x00080000
131 #define KREG_EFLAGS_VIF_SHIFT   19
133 #define KREG_EFLAGS_AC_MASK     0x00040000
134 #define KREG_EFLAGS_AC_SHIFT    18
136 #define KREG_EFLAGS_VM_MASK     0x00020000
137 #define KREG_EFLAGS_VM_SHIFT    17
139 #define KREG_EFLAGS_RF_MASK     0x00010000
140 #define KREG_EFLAGS_RF_SHIFT    16
142 #define KREG_EFLAGS_NT_MASK     0x00004000
143 #define KREG_EFLAGS_NT_SHIFT    14
145 #define KREG_EFLAGS_IOPL_MASK   0x00003000
146 #define KREG_EFLAGS_IOPL_SHIFT  12
148 #define KREG_EFLAGS_OF_MASK     0x00000800
149 #define KREG_EFLAGS_OF_SHIFT    11
151 #define KREG_EFLAGS_DF_MASK     0x00000400
152 #define KREG_EFLAGS_DF_SHIFT    10
154 #define KREG_EFLAGS_IF_MASK     0x00000200
155 #define KREG_EFLAGS_IF_SHIFT    9
157 #define KREG_EFLAGS_TF_MASK     0x00000100
158 #define KREG_EFLAGS_TF_SHIFT    8
160 #define KREG_EFLAGS_SF_MASK     0x00000080
161 #define KREG_EFLAGS_SF_SHIFT    7
163 #define KREG_EFLAGS_ZF_MASK     0x00000040
164 #define KREG_EFLAGS_ZF_SHIFT    6
166 #define KREG_EFLAGS_AF_MASK     0x00000010
167 #define KREG_EFLAGS_AF_SHIFT    4
169 #define KREG_EFLAGS_PF_MASK     0x00000004
170 #define KREG_EFLAGS_PF_SHIFT    2
172 #define KREG_EFLAGS_CF_MASK     0x00000001
173 #define KREG_EFLAGS_CF_SHIFT    0
175 /* %dr7 */
176 #define KREG_DRCTL_WP_BASESHIFT  16
177 #define KREG_DRCTL_WP_INCRSHIFT  4
178 #define KREG_DRCTL_WP_LENSHIFT   2
179 #define KREG_DRCTL_WP_LENRRWMASK 0xf
181 #define KREG_DRCTL_WP_EXEC       0
182 #define KREG_DRCTL_WP_WONLY     1
183 #define KREG_DRCTL_WP_IORW      2
184 #define KREG_DRCTL_WP_RW        3
186 #define KREG_DRCTL_WP_SHIFT(n) \
187     (KREG_DRCTL_WP_BASESHIFT + KREG_DRCTL_WP_INCRSHIFT * (n))

```

```

188 #define KREG_DRCTL_WP_MASK(n) \
189     (KREG_DRCTL_WP_LENRRWMASK << KREG_DRCTL_WP_SHIFT(n))
190 #define KREG_DRCTL_WP_LENRRW(n, len, rw) \
191     (((len) << KREG_DRCTL_WP_LENSHIFT) | (rw) << KREG_DRCTL_WP_SHIFT(n))
193 #define KREG_DRCTL_WPEN_INCRSHIFT 2
194 #define KREG_DRCTL_WPEN_MASK(n) \
195     (3 << (KREG_DRCTL_WPEN_INCRSHIFT * (n)))
196 #define KREG_DRCTL_WPEN(n)      KREG_DRCTL_WPEN_MASK(n)
198 /* %dr6 */
199 #define KREG_DRSTAT_BT_MASK     0x00008000
200 #define KREG_DRSTAT_BS_MASK    0x00004000
201 #define KREG_DRSTAT_BD_MASK    0x00002000
203 #define KREG_DRSTAT_WP_MASK(n) (1 << (n))
205 #ifdef __cplusplus
206 }

```

unchanged portion omitted

```

*****
13668 Fri Apr 6 17:24:58 2018
new/usr/src/cmd/mdb/sparc/kmdb/kvm_isadep.c
9210 remove KMDB branch debugging support
9211 ::crregs could do with cr2/cr3 support
9209 ::ttrace should be able to filter by thread
Reviewed by: Patrick Mooney <patrick.mooney@joyent.com>
Reviewed by: Yuri Pankov <yuripv@yuripv.net>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License, Version 1.0 only
6  * (the "License"). You may not use this file except in compliance
7  * with the License.
8  *
9  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
10 * or http://www.opensolaris.org/os/licensing.
11 * See the License for the specific language governing permissions
12 * and limitations under the License.
13 *
14 * When distributing Covered Code, include this CDDL HEADER in each
15 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
16 * If applicable, add the following below this CDDL HEADER, with the
17 * fields enclosed by brackets "[]" replaced with your own identifying
18 * information: Portions Copyright [yyyy] [name of copyright owner]
19 *
20 * CDDL HEADER END
21 */
22 /*
23 * Copyright 2006 Sun Microsystems, Inc. All rights reserved.
24 * Use is subject to license terms.
25 *
26 * Copyright 2018 Joyent, Inc.
27 */

27 #pragma ident      "%Z%M% %I%      %E% SMI"

29 /*
30 * isa-dependent portions of the kmdb target
31 */

33 #include <mdb/mdb_kreg_impl.h>
34 #include <mdb/mdb_debug.h>
35 #include <mdb/mdb_modapi.h>
36 #include <mdb/mdb_v9util.h>
37 #include <mdb/mdb_target_impl.h>
38 #include <mdb/mdb_err.h>
39 #include <mdb/mdb_umem.h>
40 #include <kmdb/kmdb_kdi.h>
41 #include <kmdb/kmdb_dpi.h>
42 #include <kmdb/kmdb_promif.h>
43 #include <kmdb/kmdb_asmutil.h>
44 #include <kmdb/kvm.h>
45 #include <mdb/mdb.h>

47 #include <sys/types.h>
48 #include <sys/stack.h>
49 #include <sys/regset.h>
50 #include <sys/sysmacros.h>
51 #include <sys/bitmap.h>
52 #include <sys/machtrap.h>
53 #include <sys/trap.h>

55 /* Higher than the highest trap number for which we have a specific specifier */

```

```

56 #define KMT_MAXTRAPNO    0x1fff

58 #define OP(x)              ((x) >> 30)
59 #define OP3(x)            (((x) >> 19) & 0x3f)
60 #define RD(x)            (((x) >> 25) & 0x1f)
61 #define RS1(x)           (((x) >> 14) & 0x1f)
62 #define RS2(x)           ((x) & 0x1f)

64 #define OP_ARITH          0x2

66 #define OP3_OR            0x02
67 #define OP3_SAVE         0x3c
68 #define OP3_RESTORE      0x3d

70 static int
71 kmt_stack_iter(mdb_tgt_t *t, const mdb_tgt_gregset_t *gsp,
72               mdb_tgt_stack_f *func, void *arg, int cpuid)
73 {
74     const mdb_tgt_gregset_t *grp;
75     mdb_tgt_gregset_t gregs;
76     kreg_t *kregs = &gregs.kregs[0];
77     long nwin, stopwin, canrestore, wp, i, sp;
78     long argv[6];

80     /*
81      * If gsp isn't null, we were asked to dump a trace from a
82      * specific location. The normal iterator can handle that.
83      */
84     if (gsp != NULL) {
85         if (cpuid != DPI_MASTER_CPUID)
86             warn("register set provided - ignoring cpu argument\n");
87         return (mdb_kvm_v9stack_iter(t, gsp, func, arg));
88     }

90     if (kmdb_dpi_get_cpu_state(cpuid) < 0) {
91         warn("failed to iterate through stack for cpu %u", cpuid);
92         return (DCMD_ERR);
93     }

95     /*
96      * We're being asked to dump the trace for the current CPU.
97      * To do that, we need to iterate first through the saved
98      * register windors. If there's more to the trace than that,
99      * we'll hand off to the normal iterator.
100     */
101     if ((grp = kmdb_dpi_get_gregs(cpuid)) == NULL) {
102         warn("failed to retrieve registers for cpu %d", cpuid);
103         return (DCMD_ERR);
104     }

106     bcopy(grp, &gregs, sizeof (mdb_tgt_gregset_t));

108     wp = kregs[KREG_CWP];
109     canrestore = kregs[KREG_CANRESTORE];
110     nwin = kmdb_dpi_get_nwin(cpuid);
111     stopwin = ((wp + nwin) - canrestore - 1) % nwin;

113     mdb_dprintf(MDB_DBG_KMOD, "dumping cwp = %lu, canrestore = %lu, "
114               "stopwin = %lu\n", wp, canrestore, stopwin);

116     for (;;) {
117         struct rwindow rwin;

119         for (i = 0; i < 6; i++)
120             argv[i] = kregs[KREG_IO + i];

```

```

122         if (kregs[KREG_PC] != 0 &&
123             func(arg, kregs[KREG_PC], 6, argv, &gregs) != 0)
124             return (0);

126         kregs[KREG_PC] = kregs[KREG_I7];
127         kregs[KREG_NPC] = kregs[KREG_PC] + 4;

129         if ((sp = kregs[KREG_FP] + STACK_BIAS) == STACK_BIAS || sp == 0)
130             return (0); /* Stop if we're at the end of stack */

132         if (sp & (STACK_ALIGN - 1))
133             return (set_errno(EMDB_STKALIGN));

135         wp = (wp + nwin - 1) % nwin;

137         if (wp == stopwin)
138             break;

140         bcopy(&kregs[KREG_I0], &kregs[KREG_O0], 8 * sizeof (kreg_t));

142         if (kmdb_dpi_get_rwin(cpuid, wp, &rwin) < 0) {
143             warn("unable to get registers from window %ld\n", wp);
144             return (-1);
145         }

147         for (i = 0; i < 8; i++)
148             kregs[KREG_L0 + i] = (uintptr_t)rwin.rw_local[i];
149         for (i = 0; i < 8; i++)
150             kregs[KREG_IO + i] = (uintptr_t)rwin.rw_in[i];
151     }

153     mdb_dprintf(MDB_DBG_KMOD, "dumping wp %ld and beyond normally\n", wp);

155     /*
156     * hack - if we null out pc here, iterator won't print the frame
157     * that corresponds to the current set of registers. That's what we
158     * want because we just printed them above.
159     */
160     kregs[KREG_PC] = 0;
161     return (mdb_kvm_v9stack_iter(t, &gregs, func, arg));
162 }

```

unchanged portion omitted

```

380 /*ARGSUSED*/
381 int
382 kmt_step_branch(mdb_tgt_t *t)
383 {
384     return (set_errno(EMDB_TGTHWNOTSUP));
385 }

```

```

380 static const char *
381 regno2name(int idx)
382 {
383     const mdb_tgt_regdesc_t *rd;

385     for (rd = mdb_sparcv9_kregs; rd->rd_name != NULL; rd++) {
386         if (idx == rd->rd_num)
387             return (rd->rd_name);
388     }

390     ASSERT(rd->rd_name != NULL);

392     return ("unknown");
393 }

```

unchanged portion omitted

```

*****
15729 Fri Apr 6 17:24:59 2018
new/usr/src/cmd/mdb/sparc/mdb/kvm_v7dep.c
9210 remove KMDB branch debugging support
9211 ::crregs could do with cr2/cr3 support
9209 ::ttrace should be able to filter by thread
Reviewed by: Patrick Mooney <patrick.mooney@joyent.com>
Reviewed by: Yuri Pankov <yuripv@yuripv.net>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22  * Copyright 2008 Sun Microsystems, Inc. All rights reserved.
23  * Use is subject to license terms.
24  *
25  * Copyright 2018 Joyent, Inc.
26  */

26 #pragma ident "%Z%M% %I% %E% SMI"

28 /*
29  * Libkvm Kernel Target SPARC v7 component
30  *
31  * This file provides the ISA-dependent portion of the libkvm kernel target.
32  * For more details on the implementation refer to mdb_kvm.c.
33  */

35 #include <sys/types.h>
36 #include <sys/machtypes.h>
37 #include <sys/regset.h>
38 #include <sys/frame.h>
39 #include <sys/stack.h>
40 #include <sys/sysmacros.h>
41 #include <sys/panic.h>
42 #include <strings.h>

44 #include <mdb/mdb_target_impl.h>
45 #include <mdb/mdb_disasm.h>
46 #include <mdb/mdb_modapi.h>
47 #include <mdb/mdb_conf.h>
48 #include <mdb/mdb_kreg.h>
49 #include <mdb/mdb_kvm.h>
50 #include <mdb/mdb_err.h>
51 #include <mdb/mdb_debug.h>
52 #include <mdb/mdb.h>

54 /*
55  * The mdb_tgt_gregset type is opaque to callers of the target interface

```

```

56 * and to our own target common code. We now can define it explicitly.
57 */
58 struct mdb_tgt_gregset {
59     kreg_t kregs[KREG_NGREG];
60 };
_____ unchanged_portion_omitted _____

371 const mdb_tgt_ops_t kt_sparcv7_ops = {
372     kt_setflags,                /* t_setflags */
373     kt_setcontext,             /* t_setcontext */
374     kt_activate,               /* t_activate */
375     kt_deactivate,             /* t_deactivate */
376     (void (*)()) mdb_tgt_nop,  /* t_periodic */
377     kt_destroy,                 /* t_destroy */
378     kt_name,                    /* t_name */
379     (const char (*)()) mdb_conf_isa, /* t_isa */
380     kt_platform,                /* t_platform */
381     kt_uname,                   /* t_uname */
382     kt_dmodel,                  /* t_dmodel */
383     kt_aread,                   /* t_aread */
384     kt_awrite,                  /* t_awrite */
385     kt_vread,                   /* t_vread */
386     kt_vwrite,                  /* t_vwrite */
387     kt_pread,                   /* t_pread */
388     kt_pwrite,                  /* t_pwrite */
389     kt_fread,                   /* t_fread */
390     kt_fwrite,                  /* t_fwrite */
391     (ssize_t (*)()) mdb_tgt_notsup, /* t_ioread */
392     (ssize_t (*)()) mdb_tgt_notsup, /* t_iowrite */
393     kt_vtop,                     /* t_vtop */
394     kt_lookup_by_name,           /* t_lookup_by_name */
395     kt_lookup_by_addr,          /* t_lookup_by_addr */
396     kt_symbol_iter,             /* t_symbol_iter */
397     kt_mapping_iter,            /* t_mapping_iter */
398     kt_object_iter,             /* t_object_iter */
399     kt_addr_to_map,              /* t_addr_to_map */
400     kt_name_to_map,             /* t_name_to_map */
401     kt_addr_to_ctf,             /* t_addr_to_ctf */
402     kt_name_to_ctf,             /* t_name_to_ctf */
403     kt_status,                  /* t_status */
404     (int (*)()) mdb_tgt_notsup, /* t_run */
405     (int (*)()) mdb_tgt_notsup, /* t_step */
406     (int (*)()) mdb_tgt_notsup, /* t_step_out */
407     (int (*)()) mdb_tgt_notsup, /* t_step_branch */
408     (int (*)()) mdb_tgt_notsup, /* t_next */
409     (int (*)()) mdb_tgt_notsup, /* t_cont */
410     (int (*)()) mdb_tgt_notsup, /* t_signal */
411     (int (*)()) mdb_tgt_null,   /* t_add_vbrkpt */
412     (int (*)()) mdb_tgt_null,   /* t_add_sbrkpt */
413     (int (*)()) mdb_tgt_null,   /* t_add_pwapt */
414     (int (*)()) mdb_tgt_null,   /* t_add_vwapt */
415     (int (*)()) mdb_tgt_null,   /* t_add_iowapt */
416     (int (*)()) mdb_tgt_null,   /* t_add_sysenter */
417     (int (*)()) mdb_tgt_null,   /* t_add_sysexit */
418     (int (*)()) mdb_tgt_null,   /* t_add_signal */
419     (int (*)()) mdb_tgt_null,   /* t_add_fault */
420     kt_getareg,                 /* t_getareg */
421     kt_putareg,                 /* t_putareg */
422     kt_stack_iter,              /* t_stack_iter */
423     (int (*)()) mdb_tgt_notsup, /* t_auxv */
_____ unchanged_portion_omitted _____

```

```

*****
10883 Fri Apr 6 17:24:59 2018
new/usr/src/cmd/mdb/sparc/mdb/kvm_v9dep.c
9210 remove KMDB branch debugging support
9211 ::crregs could do with cr2/cr3 support
9209 ::ttrace should be able to filter by thread
Reviewed by: Patrick Mooney <patrick.mooney@joyent.com>
Reviewed by: Yuri Pankov <yuripv@yuripv.net>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22  * Copyright 2008 Sun Microsystems, Inc. All rights reserved.
23  * Use is subject to license terms.
24  *
25  * Copyright 2018 Joyent, Inc.
26  */

26 #pragma ident "%Z%%M% %I% %E% SMI"

28 /*
29  * Libkvm Kernel Target SPARC v9 component
30  *
31  * This file provides the ISA-dependent portion of the libkvm kernel target.
32  * For more details on the implementation refer to mdb_kvm.c. The SPARC v9
33  * ISA code is actually compiled into *both* the sparcv7 and sparcv9 MDB
34  * binaries because we need to deal with the sparcv9 CPU registers when
35  * debugging a 32-bit crash dump from a kernel running on a sparcv9 CPU.
36  */

38 #ifndef __sparcv9cpu
39 #define __sparcv9cpu
40 #endif

42 #include <sys/types.h>
43 #include <sys/machtypes.h>
44 #include <sys/regset.h>
45 #include <sys/frame.h>
46 #include <sys/stack.h>
47 #include <sys/sysmacros.h>
48 #include <sys/panic.h>
49 #include <strings.h>

51 #include <mdb/mdb_target_impl.h>
52 #include <mdb/mdb_disasm.h>
53 #include <mdb/mdb_modapi.h>
54 #include <mdb/mdb_conf.h>
55 #include <mdb/mdb_kreg_impl.h>

```

```

56 #include <mdb/mdb_v9util.h>
57 #include <mdb/mdb_kvm.h>
58 #include <mdb/mdb_err.h>
59 #include <mdb/mdb_debug.h>
60 #include <mdb/mdb.h>

62 #ifndef STACK_BIAS
63 #define STACK_BIAS 0
64 #endif

66 static int
67 kt_getareg(mdb_tgt_t *t, mdb_tgt_tid_t tid,
68            const char *rname, mdb_tgt_reg_t *rp)
69 {
70     const mdb_tgt_regdesc_t *rdp;
71     kt_data_t *kt = t->t_data;

73     if (tid != kt->k_tid)
74         return (set_errno(EMDB_NOREGS));

76     for (rdp = kt->k_rds; rdp->rd_name != NULL; rdp++) {
77         if (strcmp(rname, rdp->rd_name) == 0) {
78             *rp = kt->k_regs->kregs[rdp->rd_num];
79             return (0);
80         }
81     }

83     return (set_errno(EMDB_BADREG));
84 }

unchanged_portion_omitted

184 const mdb_tgt_ops_t kt_sparcv9_ops = {
185     kt_setflags, /* t_setflags */
186     kt_setcontext, /* t_setcontext */
187     kt_activate, /* t_activate */
188     kt_deactivate, /* t_deactivate */
189     (void (*)()) mdb_tgt_nop, /* t_periodic */
190     kt_destroy, /* t_destroy */
191     kt_name, /* t_name */
192     (const char (*)()) mdb_conf_isa, /* t_isa */
193     kt_platform, /* t_platform */
194     kt_uname, /* t_uname */
195     kt_dmodel, /* t_dmodel */
196     kt_aread, /* t_aread */
197     kt_awrite, /* t_awrite */
198     kt_vread, /* t_vread */
199     kt_vwrite, /* t_vwrite */
200     kt_pread, /* t_pread */
201     kt_pwrite, /* t_pwrite */
202     kt_fread, /* t_fread */
203     kt_fwrite, /* t_fwrite */
204     (ssize_t (*)()) mdb_tgt_notsup, /* t_ioread */
205     (ssize_t (*)()) mdb_tgt_notsup, /* t_iowrite */
206     kt_vtop, /* t_vtop */
207     kt_lookup_by_name, /* t_lookup_by_name */
208     kt_lookup_by_addr, /* t_lookup_by_addr */
209     kt_symbol_iter, /* t_symbol_iter */
210     kt_mapping_iter, /* t_mapping_iter */
211     kt_object_iter, /* t_object_iter */
212     kt_addr_to_map, /* t_addr_to_map */
213     kt_name_to_map, /* t_name_to_map */
214     kt_addr_to_ctf, /* t_addr_to_ctf */
215     kt_name_to_ctf, /* t_name_to_ctf */
216     kt_status, /* t_status */
217     (int (*)()) mdb_tgt_notsup, /* t_run */
218     (int (*)()) mdb_tgt_notsup, /* t_step */

```

```
219     (int (*)()) mdb_tgt_notsup,      /* t_step_out */
220     (int (*)()) mdb_tgt_notsup,      /* t_step_branch */
220     (int (*)()) mdb_tgt_notsup,      /* t_next */
221     (int (*)()) mdb_tgt_notsup,      /* t_cont */
222     (int (*)()) mdb_tgt_notsup,      /* t_signal */
223     (int (*)()) mdb_tgt_null,         /* t_add_vbrkpt */
224     (int (*)()) mdb_tgt_null,         /* t_add_sbrkpt */
225     (int (*)()) mdb_tgt_null,         /* t_add_pwapt */
226     (int (*)()) mdb_tgt_null,         /* t_add_iowapt */
227     (int (*)()) mdb_tgt_null,         /* t_add_vwapt */
228     (int (*)()) mdb_tgt_null,         /* t_add_sysenter */
229     (int (*)()) mdb_tgt_null,         /* t_add_sysexit */
230     (int (*)()) mdb_tgt_null,         /* t_add_signal */
231     (int (*)()) mdb_tgt_null,         /* t_add_fault */
232     kt_getareg,                       /* t_getareg */
233     kt_putareg,                       /* t_putareg */
234     mdb_kvm_v9stack_iter,             /* t_stack_iter */
235     (int (*)()) mdb_tgt_notsup,      /* t_auxv */
236 };
    unchanged_portion_omitted
```

```

*****
12018 Fri Apr 6 17:24:59 2018
new/usr/src/uts/common/sys/sysmacros.h
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*      Copyright (c) 1984, 1986, 1987, 1988, 1989 AT&T */
22 /*      All Rights Reserved      */

25 /*
26 * Copyright 2008 Sun Microsystems, Inc.  All rights reserved.
27 * Use is subject to license terms.
28 *
29 * Copyright 2013 Nexenta Systems, Inc.  All rights reserved.
30 *
31 * Copyright 2018 Joyent Inc.
32 */

34 #ifndef _SYS_SYSMACROS_H
35 #define _SYS_SYSMACROS_H

37 #include <sys/param.h>
38 #include <sys/stddef.h>

40 #ifdef __cplusplus
41 extern "C" {
42 #endif

44 /*
45 * Some macros for units conversion
46 */
47 /*
48 * Disk blocks (sectors) and bytes.
49 */
50 #define dtob(DD)      ((DD) << DEV_BSHIFT)
51 #define btod(BB)      (((BB) + DEV_BSIZE - 1) >> DEV_BSHIFT)
52 #define btodt(BB)     ((BB) >> DEV_BSHIFT)
53 #define lbtod(BB)     (((offset_t)(BB) + DEV_BSIZE - 1) >> DEV_BSHIFT)

55 /* common macros */
56 #ifndef MIN
57 #define MIN(a, b)      ((a) < (b) ? (a) : (b))
58 #endif
59 #ifndef MAX

```

```

60 #define MAX(a, b)      ((a) < (b) ? (b) : (a))
61 #endif
62 #ifndef ABS
63 #define ABS(a)         ((a) < 0 ? -(a) : (a))
64 #endif
65 #ifndef SIGNOF
66 #define SIGNOF(a)      ((a) < 0 ? -1 : (a) > 0)
67 #endif

69 #ifdef _KERNEL

71 /*
72 * Convert a single byte to/from binary-coded decimal (BCD).
73 */
74 extern unsigned char byte_to_bcd[256];
75 extern unsigned char bcd_to_byte[256];

77 #define BYTE_TO_BCD(x)  byte_to_bcd[(x) & 0xff]
78 #define BCD_TO_BYTE(x) bcd_to_byte[(x) & 0xff]

80 #endif /* _KERNEL */

82 /*
83 * WARNING: The device number macros defined here should not be used by device
84 * drivers or user software. Device drivers should use the device functions
85 * defined in the DDI/DKI interface (see also ddi.h). Application software
86 * should make use of the library routines available in makedev(3). A set of
87 * new device macros are provided to operate on the expanded device number
88 * format supported in SVR4. Macro versions of the DDI device functions are
89 * provided for use by kernel proper routines only. Macro routines bmajor(),
90 * major(), minor(), emajor(), eminor(), and makedev() will be removed or
91 * their definitions changed at the next major release following SVR4.
92 */

94 #define O_BITSMAJOR    7      /* # of SVR3 major device bits */
95 #define O_BITSMINOR    8      /* # of SVR3 minor device bits */
96 #define O_MAXMAJ       0x7f   /* SVR3 max major value */
97 #define O_MAXMIN       0xff   /* SVR3 max minor value */

100 #define L_BITSMAJOR32  14     /* # of SVR4 major device bits */
101 #define L_BITSMINOR32 18     /* # of SVR4 minor device bits */
102 #define L_MAXMAJ32     0x3fff /* SVR4 max major value */
103 #define L_MAXMIN32     0x3fff /* MAX minor for 3b2 software drivers. */
104 /* For 3b2 hardware devices the minor is */
105 /* restricted to 256 (0-255) */

107 #ifdef _LP64
108 #define L_BITSMAJOR    32     /* # of major device bits in 64-bit Solaris */
109 #define L_BITSMINOR    32     /* # of minor device bits in 64-bit Solaris */
110 #define L_MAXMAJ       0xfffffffful /* max major value */
111 #define L_MAXMIN       0xfffffffful /* max minor value */
112 #else
113 #define L_BITSMAJOR32  L_BITSMAJOR32
114 #define L_BITSMINOR32 L_BITSMINOR32
115 #define L_MAXMAJ32     L_MAXMAJ32
116 #define L_MAXMIN32     L_MAXMIN32
117 #endif

119 #ifdef _KERNEL

121 /* major part of a device internal to the kernel */

123 #define major(x)        (major_t)((((unsigned)(x)) >> O_BITSMINOR) & O_MAXMAJ)
124 #define bmajor(x)       (major_t)((((unsigned)(x)) >> O_BITSMINOR) & O_MAXMAJ)

```

```

126 /* get internal major part of expanded device number */
128 #define getmajor(x)      (major_t)((((dev_t)(x)) >> L_BITSMINOR) & L_MAXMAJ)
130 /* minor part of a device internal to the kernel */
132 #define minor(x)        (minor_t)((x) & O_MAXMIN)
134 /* get internal minor part of expanded device number */
136 #define getminor(x)     (minor_t)((x) & L_MAXMIN)
138 #else /* _KERNEL */
140 /* major part of a device external from the kernel (same as emajor below) */
142 #define major(x)        (major_t)((((unsigned)(x)) >> O_BITSMINOR) & O_MAXMAJ)
144 /* minor part of a device external from the kernel (same as eminor below) */
146 #define minor(x)        (minor_t)((x) & O_MAXMIN)
148 #endif /* _KERNEL */
150 /* create old device number */
152 #define makedev(x, y) (unsigned short)((((x) << O_BITSMINOR) | ((y) & O_MAXMIN))
154 /* make an new device number */
156 #define makedevice(x, y) (dev_t)((((dev_t)(x) << L_BITSMINOR) | ((y) & L_MAXMIN))
159 /*
160 * emajor() allows kernel/driver code to print external major numbers
161 * eminor() allows kernel/driver code to print external minor numbers
162 */
164 #define emajor(x) \
165     (major_t)((((unsigned int)(x) >> O_BITSMINOR) > O_MAXMAJ) ? \
166     NODEV : (((unsigned int)(x) >> O_BITSMINOR) & O_MAXMAJ)
168 #define eminor(x) \
169     (minor_t)((x) & O_MAXMIN)
171 /*
172 * get external major and minor device
173 * components from expanded device number
174 */
175 #define getemajor(x)      (major_t)((((dev_t)(x) >> L_BITSMINOR) > L_MAXMAJ) ? \
176     NODEV : (((dev_t)(x) >> L_BITSMINOR) & L_MAXMAJ))
177 #define geteminor(x)     (minor_t)((x) & L_MAXMIN)
179 /*
180 * These are versions of the kernel routines for compressing and
181 * expanding long device numbers that don't return errors.
182 */
183 #if (L_BITSMINOR == L_BITSMINOR) && (L_BITSMINOR32 == L_BITSMINOR32)
185 #define DEVCML(x)      (x)
186 #define DEVEXPL(x)    (x)
188 #else
190 #define DEVCML(x)      \
191     (dev32_t)((((x) >> L_BITSMINOR) > L_MAXMAJ32) || \

```

```

192     ((x) & L_MAXMIN) > L_MAXMIN32) ? NODEV32 : \
193     (((x) >> L_BITSMINOR) << L_BITSMINOR32) | ((x) & L_MAXMIN32)))
195 #define DEVEXPL(x)      \
196     (((x) == NODEV32) ? NODEV : \
197     makedevice(((x) >> L_BITSMINOR32) & L_MAXMAJ32, (x) & L_MAXMIN32))
199 #endif /* L_BITSMINOR32 ... */
201 /* convert to old (SVR3.2) dev format */
203 #define cmpdev(x) \
204     (o_dev_t)((((x) >> L_BITSMINOR) > O_MAXMAJ) || \
205     ((x) & L_MAXMIN) > O_MAXMIN) ? NODEV : \
206     (((x) >> L_BITSMINOR) << O_BITSMINOR) | ((x) & O_MAXMIN))
208 /* convert to new (SVR4) dev format */
210 #define expdev(x) \
211     (dev_t)((((dev_t)((x) >> O_BITSMINOR) & O_MAXMAJ) << L_BITSMINOR) | \
212     ((x) & O_MAXMIN))
214 /*
215 * Macro for checking power of 2 address alignment.
216 */
217 #define IS_P2ALIGNED(v, a) (((uintptr_t)(v)) & ((uintptr_t)(a) - 1)) == 0)
219 /*
220 * Macros for counting and rounding.
221 */
222 #define howmany(x, y)      (((x)+(y)-1)/(y))
223 #define roundup(x, y)    (((x)+(y)-1)/(y))*y)
225 /*
226 * Macro to determine if value is a power of 2
227 */
228 #define ISP2(x)           (((x) & ((x) - 1)) == 0)
230 /*
231 * Macros for various sorts of alignment and rounding. The "align" must
232 * be a power of 2. Often times it is a block, sector, or page.
233 */
235 /*
236 * return x rounded down to an align boundary
237 * eg, P2ALIGN(1200, 1024) == 1024 (1*align)
238 * eg, P2ALIGN(1024, 1024) == 1024 (1*align)
239 * eg, P2ALIGN(0x1234, 0x100) == 0x1200 (0x12*align)
240 * eg, P2ALIGN(0x5600, 0x100) == 0x5600 (0x56*align)
241 */
242 #define P2ALIGN(x, align)      ((x) & ~(align - 1))
244 /*
245 * return x % (mod) align
246 * eg, P2PHASE(0x1234, 0x100) == 0x34 (x-0x12*align)
247 * eg, P2PHASE(0x5600, 0x100) == 0x00 (x-0x56*align)
248 */
249 #define P2PHASE(x, align)     ((x) & ((align) - 1))
251 /*
252 * return how much space is left in this block (but if it's perfectly
253 * aligned, return 0).
254 * eg, P2NPHASE(0x1234, 0x100) == 0xcc (0x13*align-x)
255 * eg, P2NPHASE(0x5600, 0x100) == 0x00 (0x56*align-x)
256 */
257 #define P2NPHASE(x, align)    (-(x) & ((align) - 1))

```



```

259 /*
260 * return x rounded up to an align boundary
261 * eg, P2ROUNDUP(0x1234, 0x100) == 0x1300 (0x13*align)
262 * eg, P2ROUNDUP(0x5600, 0x100) == 0x5600 (0x56*align)
263 */
264 #define P2ROUNDUP(x, align)          (-(~(x) & ~(align)))

266 /*
267 * return the ending address of the block that x is in
268 * eg, P2END(0x1234, 0x100) == 0x12ff (0x13*align - 1)
269 * eg, P2END(0x5600, 0x100) == 0x56ff (0x57*align - 1)
270 */
271 #define P2END(x, align)              (-(~(x) & ~(align)))

273 /*
274 * return x rounded up to the next phase (offset) within align.
275 * phase should be < align.
276 * eg, P2PHASEUP(0x1234, 0x100, 0x10) == 0x1310 (0x13*align + phase)
277 * eg, P2PHASEUP(0x5600, 0x100, 0x10) == 0x5610 (0x56*align + phase)
278 */
279 #define P2PHASEUP(x, align, phase)   ((phase) - (((phase) - (x)) & ~(align)))

281 /*
282 * return TRUE if adding len to off would cause it to cross an align
283 * boundary.
284 * eg, P2BOUNDARY(0x1234, 0xe0, 0x100) == TRUE (0x1234 + 0xe0 == 0x1314)
285 * eg, P2BOUNDARY(0x1234, 0x50, 0x100) == FALSE (0x1234 + 0x50 == 0x1284)
286 */
287 #define P2BOUNDARY(off, len, align) \
288     (((off) ^ ((off) + (len) - 1)) > (align) - 1)

290 /*
291 * Return TRUE if they have the same highest bit set.
292 * eg, P2SAMEHIGHBIT(0x1234, 0x1001) == TRUE (the high bit is 0x1000)
293 * eg, P2SAMEHIGHBIT(0x1234, 0x3010) == FALSE (high bit of 0x3010 is 0x2000)
294 */
295 #define P2SAMEHIGHBIT(x, y)         (((x) ^ (y)) < ((x) & (y)))

297 /*
298 * Typed version of the P2* macros. These macros should be used to ensure
299 * that the result is correctly calculated based on the data type of (x),
300 * which is passed in as the last argument, regardless of the data
301 * type of the alignment. For example, if (x) is of type uint64_t,
302 * and we want to round it up to a page boundary using "PAGESIZE" as
303 * the alignment, we can do either
304 *     P2ROUNDUP(x, (uint64_t)PAGESIZE)
305 * or
306 *     P2ROUNDUP_TYPED(x, PAGESIZE, uint64_t)
307 */
308 #define P2ALIGN_TYPED(x, align, type) \
309     ((type)(x) & ~(type)(align))
310 #define P2PHASE_TYPED(x, align, type) \
311     ((type)(x) & ((type)(align) - 1))
312 #define P2NPHASE_TYPED(x, align, type) \
313     (~(type)(x) & ((type)(align) - 1))
314 #define P2ROUNDUP_TYPED(x, align, type) \
315     (-(~(type)(x) & ~(type)(align)))
316 #define P2END_TYPED(x, align, type) \
317     (-(~(type)(x) & ~(type)(align)))
318 #define P2PHASEUP_TYPED(x, align, phase, type) \
319     ((type)(phase) - (((type)(phase) - (type)(x)) & ~(type)(align)))
320 #define P2CROSS_TYPED(x, y, align, type) \
321     (((type)(x) ^ (type)(y)) > (type)(align) - 1)
322 #define P2SAMEHIGHBIT_TYPED(x, y, type) \
323     (((type)(x) ^ (type)(y)) < ((type)(x) & (type)(y)))

```

```

325 /*
326 * Macros to atomically increment/decrement a variable. mutex and var
327 * must be pointers.
328 */
329 #define INCR_COUNT(var, mutex) mutex_enter(mutex), (*(var))++, mutex_exit(mutex)
330 #define DECR_COUNT(var, mutex) mutex_enter(mutex), (*(var))--, mutex_exit(mutex)

332 /*
333 * Macros to declare bitfields - the order in the parameter list is
334 * Low to High - that is, declare bit 0 first. We only support 8-bit bitfields
335 * because if a field crosses a byte boundary it's not likely to be meaningful
336 * without reassembly in its nonnative endianness.
337 */
338 #if defined(_BIT_FIELDS_LTOH)
339 #define DECL_BITFIELD2(_a, _b) \
340     uint8_t _a, _b \
341 #define DECL_BITFIELD3(_a, _b, _c) \
342     uint8_t _a, _b, _c \
343 #define DECL_BITFIELD4(_a, _b, _c, _d) \
344     uint8_t _a, _b, _c, _d \
345 #define DECL_BITFIELD5(_a, _b, _c, _d, _e) \
346     uint8_t _a, _b, _c, _d, _e \
347 #define DECL_BITFIELD6(_a, _b, _c, _d, _e, _f) \
348     uint8_t _a, _b, _c, _d, _e, _f \
349 #define DECL_BITFIELD7(_a, _b, _c, _d, _e, _f, _g) \
350     uint8_t _a, _b, _c, _d, _e, _f, _g \
351 #define DECL_BITFIELD8(_a, _b, _c, _d, _e, _f, _g, _h) \
352     uint8_t _a, _b, _c, _d, _e, _f, _g, _h
353 #elif defined(_BIT_FIELDS_HTOH)
354 #define DECL_BITFIELD2(_a, _b) \
355     uint8_t _b, _a \
356 #define DECL_BITFIELD3(_a, _b, _c) \
357     uint8_t _c, _b, _a \
358 #define DECL_BITFIELD4(_a, _b, _c, _d) \
359     uint8_t _d, _c, _b, _a \
360 #define DECL_BITFIELD5(_a, _b, _c, _d, _e) \
361     uint8_t _e, _d, _c, _b, _a \
362 #define DECL_BITFIELD6(_a, _b, _c, _d, _e, _f) \
363     uint8_t _f, _e, _d, _c, _b, _a \
364 #define DECL_BITFIELD7(_a, _b, _c, _d, _e, _f, _g) \
365     uint8_t _g, _f, _e, _d, _c, _b, _a \
366 #define DECL_BITFIELD8(_a, _b, _c, _d, _e, _f, _g, _h) \
367     uint8_t _h, _g, _f, _e, _d, _c, _b, _a
368 #else
369 #error One of _BIT_FIELDS_LTOH or _BIT_FIELDS_HTOH must be defined
370 #endif /* _BIT_FIELDS_LTOH */

372 #if !defined(ARRAY_SIZE)
370 /* avoid any possibility of clashing with <stddef.h> version */
371 #if (defined(_KERNEL) || defined(_FAKE_KERNEL)) && !defined(_KMUSER)

373 #define ARRAY_SIZE(x) (sizeof (x) / sizeof (x[0]))
374 #endif

375 #endif /* _KERNEL, !_KMUSER */

376 #ifdef __cplusplus
377 }
    unchanged_portion_omitted

```

```

*****
6325 Fri Apr 6 17:24:59 2018
new/usr/src/uts/i86pc/Makefile.files
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>
*****
1 #
2 # CDDL HEADER START
3 #
4 # The contents of this file are subject to the terms of the
5 # Common Development and Distribution License (the "License").
6 # You may not use this file except in compliance with the License.
7 #
8 # You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9 # or http://www.opensolaris.org/os/licensing.
10 # See the License for the specific language governing permissions
11 # and limitations under the License.
12 #
13 # When distributing Covered Code, include this CDDL HEADER in each
14 # file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 # If applicable, add the following below this CDDL HEADER, with the
16 # fields enclosed by brackets "[]" replaced with your own identifying
17 # information: Portions Copyright [yyyy] [name of copyright owner]
18 #
19 # CDDL HEADER END
20 #
21 #
22 #
23 # Copyright (c) 1992, 2010, Oracle and/or its affiliates. All rights reserved.
24 #
25 # Copyright (c) 2010, Intel Corporation.
26 # Copyright 2018 Joyent, Inc.
27 # Copyright 2017 Joyent, Inc.
28 #
29 # This Makefile defines file modules in the directory uts/i86pc
30 # and its children. These are the source files which are i86pc
31 # "implementation architecture" dependent.
32 #
33 #
34 # object lists
35 #
36 CORE_OBJS += \
37 acpi_stubs.o \
38 biosdisk.o \
39 bios_call.o \
40 cbe.o \
41 cmi.o \
42 cmi_hw.o \
43 cms.o \
44 comm_page.o \
45 confunix.o \
46 cpu_idle.o \
47 cpuid.o \
48 cpuid_subr.o \
49 cpupm.o \
50 cpupm_mach.o \
51 cpupm_amd.o \
52 cpupm_intel.o \
53 cpupm_throttle.o \
54 cpu_acpi.o \
55 dis_tables.o \
56 ddi_impl.o \
57 dtrace_subr.o \
58 dvma.o \

```

```

59 fpu_subr.o \
60 fakebop.o \
61 fastboot.o \
62 fb_swtdch.o \
63 graphics.o \
64 hardclk.o \
65 hat_i86.o \
66 hat_kdi.o \
67 hment.o \
68 hold_page.o \
69 hrtimers.o \
70 htable.o \
71 hypercall.o \
72 hypersubr.o \
73 i86_mmu.o \
74 ibft.o \
75 instr_size.o \
76 intr.o \
77 kboot_mmu.o \
78 kdi_subr.o \
79 kdi_idt.o \
80 kdi_idthdl.o \
81 kdi_asm.o \
82 lgrpplat.o \
83 mach_kdi.o \
84 mach_sysconfig.o \
85 machdep.o \
86 mem_config.o \
87 mem_config_stubs.o \
88 mem_config_arch.o \
89 memlist_new.o \
90 memnode.o \
91 microcode.o \
92 microfind.o \
93 mlsetup.o \
94 mp_call.o \
95 mp_implfuncs.o \
96 mp_machdep.o \
97 mp_pc.o \
98 mp_startup.o \
99 memscrub.o \
100 mpcore.o \
101 notes.o \
102 pci_bios.o \
103 pci_cfgacc.o \
104 pci_cfgacc_x86.o \
105 pci_cfgspace.o \
106 pci_mech1.o \
107 pci_mech1_amd.o \
108 pci_mech2.o \
109 pci_neptune.o \
110 pci_orion.o \
111 pmem.o \
112 ppage.o \
113 pwrnow.o \
114 speedstep.o \
115 ssp.o \
116 startup.o \
117 timestamp.o \
118 todpc_subr.o \
119 trap.o \
120 turbo.o \
121 vm_machdep.o \
122 xpv_platform.o \
x_call.o \

```

```

124 #
125 #     Add the SMBIOS subsystem object files directly to the list of objects
126 #     built into unix itself; this is all common code except for smb_dev.c.
127 #
128 CORE_OBJS += $(SMBIOS_OBJS)

130 #
131 # These get compiled twice:
132 # - once in the dboot (direct boot) identity mapped code
133 # - once for use during early startup in unix
134 #
135 BOOT_DRIVER_OBJS = \
136     boot_console.o \
137     boot_keyboard.o \
138     boot_keyboard_table.o \
139     boot_vga.o \
140     boot_mmu.o \
141     dboot_multiboot2.o \
142     $(FONT_OBJS)

144 CORE_OBJS += $(BOOT_DRIVER_OBJS)

146 #
147 #     locore.o is special. It must be the first file relocated so that it
148 #     it is relocated just where its name implies.
149 #
150 SPECIAL_OBJS_32 += \
151     locore.o \
152     fast_trap_asm.o \
153     interrupt.o \
154     syscall_asm.o

156 SPECIAL_OBJS_64 += \
157     locore.o \
158     fast_trap_asm.o \
159     interrupt.o \
160     syscall_asm_amd64.o \
161     kpti_trampoline.o \
161     syscall_asm_amd64.o

163 SPECIAL_OBJS += $(SPECIAL_OBJS_$(CLASS))

165 #
166 # Objects that get compiled into the identity mapped PT_LOAD section of unix
167 # to handle the earliest part of booting.
168 #
169 DBOOT_OBJS_32 =

171 DBOOT_OBJS_64 += dboot_elfload.o

173 DBOOT_OBJS += \
174     dboot_asm.o \
175     dboot_grub.o \
176     dboot_printf.o \
177     dboot_startkern.o \
178     memcpy.o \
179     memset.o \
180     muldiv.o \
181     shal.o \
182     string.o \
183     $(BOOT_DRIVER_OBJS) \
184     $(DBOOT_OBJS_$(CLASS))

186 #
187 #     driver and misc modules
188 #

```

```

189 GFX_PRIVATE_OBJS += gfx_private.o gfxp_pci.o gfxp_segmap.o \
190     gfxp_devmap.o gfxp_vgext.o gfxp_vm.o vgasubr.o
191 FIPE_OBJS += fipec_drv.o fipec_pm.o
192 IOAT_OBJS += ioat.o ioat_rs.o ioat_ioctl.o ioat_chan.o
193 ISANEXUS_OBJS += isa.o dma_engine.o i8237A.o
194 PCIE_MISC_OBJS += pcie_acpi.o pciehp_acpi.o pcie_x86.o
195 PCI_E_NEXUS_OBJS += npe.o npe_misc.o
196 PCI_E_NEXUS_OBJS += pci_common.o pci_kstats.o pci_tools.o
197 PCINEXUS_OBJS += pci.o pci_common.o pci_kstats.o pci_tools.o
198 PCPLUSMP_OBJS += apic.o apic_regops.o psm_common.o apic_introp.o \
199     mp_platform_common.o mp_platform_misc.o \
200     hpet_acpi.o apic_common.o apic_timer.o
201 APIX_OBJS += apix.o apic_regops.o psm_common.o apix_intr.o apix_utils.o \
202     apix_irm.o mp_platform_common.o hpet_acpi.o apic_common.o \
203     apic_timer.o apix_regops.o

206 ACPI_DRV_OBJS += acpi_drv.o acpi_video.o
207 ACPINEX_OBJS += acpinex_drv.o acpinex_event.o

209 CPUDRV_OBJS += \
210     cpudrv.o \
211     cpudrv_mach.o

213 PPM_OBJS += ppm_subr.o ppm.o ppm_plat.o

215 ACPIPPM_OBJS += acpippm.o acpisleep.o
216 ACPIDEV_OBJS += acpidev_drv.o \
217     acpidev_scope.o acpidev_device.o \
218     acpidev_container.o \
219     acpidev_cpu.o \
220     acpidev_dr.o \
221     acpidev_memory.o \
222     acpidev_pci.o \
223     acpidev_resource.o \
224     acpidev_util.o

226 DRMACH ACPI_OBJS += drmach_acpi.o dr_util.o drmach_err.o

228 DR_OBJS += dr.o dr_cpu.o dr_err.o dr_io.o dr_mem_acpi.o dr_quiesce.o dr_util.o

230 ROOTNEX_OBJS += rootnex.o immu.o immu_dmar.o immu_dvma.o \
231     immu_intrmap.o immu_qinv.o immu_regs.o

233 TZMON_OBJS += tzmon.o
234 UPVC_OBJS += upvc.o psm_common.o
235 XSVC_OBJS += xsvc.o
236 AMD_IOMMU_OBJS += amd_iommu.o amd_iommu_impl.o amd_iommu_acpi.o \
237     amd_iommu_cmd.o amd_iommu_log.o amd_iommu_page_tables.o

239 #
240 #     Build up defines and paths.
241 #
242 ALL_DEFS += -Di86pc
243 INC_PATH += -I$(UTSBASE)/i86pc -I$(SRC)/common
244 INC_PATH += -I$(UTSBASE)/i86pvp -I$(UTSBASE)/common/xen

246 #
247 # Since the assym files are derived, the dependencies must be explicit for
248 # all files including this file. (This is only actually required in the
249 # instance when the .nse_depinfo file does not exist.) It may seem that
250 # the lint targets should also have a similar dependency, but they don't
251 # since only C headers are included when #defined(__lint) is true.
252 #

254 ASSYM_DEPS += \

```

```
255     copy.o           \
256     desctbls_asm.o  \
257     ddi_i86_asm.o   \
258     exception.o     \
259     fast_trap_asm.o \
260     float.o          \
261     i86_subr.o       \
262     interrupt.o     \
263     lock_prim.o     \
264     locore.o        \
265     mpcore.o        \
266     sseblk.o        \
267     swtch.o         \
268     syscall_asm.o   \
269     syscall_asm_amd64.o \
270     kpti_trampoline.o \
271     cpr_wakecode.o

273 CPR_IMPL_OBJS = cpr_impl.o    cpr_wakecode.o

275 $(KDI_ASSYM_DEPS:%=$(OBJSDIR)/%):      $(DSFDIR)/$(OBJSDIR)/kdi_assym.h

276 ASSYM_DEPS += kdi_asm.o
```

new/usr/src/uts/i86pc/io/gfx_private/gfxp_vm.c

1

```
*****
9011 Fri Apr 6 17:24:59 2018
new/usr/src/uts/i86pc/io/gfx_private/gfxp_vm.c
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */

22 /*
23  * Copyright 2007 Sun Microsystems, Inc. All rights reserved.
24  * Use is subject to license terms.
25  *
26  * Copyright 2018 Joyent, Inc.
27  */

29 #include <sys/debug.h>
30 #include <sys/types.h>
31 #include <sys/param.h>
32 #include <sys/time.h>
33 #include <sys/buf.h>
34 #include <sys/errno.h>
35 #include <sys/system.h>
36 #include <sys/conf.h>
37 #include <sys/signal.h>
38 #include <sys/file.h>
39 #include <sys/uio.h>
40 #include <sys/ioctl.h>
41 #include <sys/map.h>
42 #include <sys/proc.h>
43 #include <sys/user.h>
44 #include <sys/mman.h>
45 #include <sys/cred.h>
46 #include <sys/open.h>
47 #include <sys/stat.h>
48 #include <sys/utsname.h>
49 #include <sys/kmem.h>
50 #include <sys/cmn_err.h>
51 #include <sys/vnode.h>
52 #include <vm/page.h>
53 #include <vm/as.h>
54 #include <vm/hat.h>
55 #include <vm/seg.h>
56 #include <vm/seg_kmem.h>
57 #include <vm/hat_i86.h>
58 #include <sys/vmstmm.h>
59 #include <sys/ddi.h>
```

new/usr/src/uts/i86pc/io/gfx_private/gfxp_vm.c

2

```
60 #include <sys/devops.h>
61 #include <sys/sunddi.h>
62 #include <sys/ddi_impldefs.h>
63 #include <sys/fs/snnode.h>
64 #include <sys/pci.h>
65 #include <sys/modctl.h>
66 #include <sys/uio.h>
67 #include <sys/visual_io.h>
68 #include <sys/fbio.h>
69 #include <sys/ddidmareq.h>
70 #include <sys/tnf_probe.h>
71 #include <sys/kstat.h>
72 #include <sys/callb.h>
73 #include <sys/promif.h>
74 #include <sys/atomic.h>
75 #include <sys/gfx_private.h>

77 #ifdef __xpv
78 #include <sys/hypervisor.h>
79 #endif

81 /*
82  * Create a kva mapping for a pa (start..start+size) with
83  * the specified cache attributes (mode).
84  */
85 gfxp_kva_t
86 gfxp_map_kernel_space(uint64_t start, size_t size, uint32_t mode)
87 {
88     uint_t poffset;
89     uint64_t base;
90     pgcnt_t npages;
91     caddr_t cvaddr;
92     int hat_flags;
93     uint_t hat_attr;
94     pfn_t pfn;

96     if (size == 0)
97         return (0);

99 #ifdef __xpv
100 /*
101  * The hypervisor doesn't allow r/w mappings to some pages, such as
102  * page tables, gdt, etc. Detect %cr3 to notify users of this interface.
103  */
104 if (start == mmu_ptob(mmu_btop(getcr3_pa())))
105 if (start == mmu_ptob(mmu_btop(getcr3())))
106     return (0);
107 #endif

108 if (mode == GFXP_MEMORY_CACHED)
109     hat_attr = HAT_STORECACHING_OK;
110 else if (mode == GFXP_MEMORY_WRITECOMBINED)
111     hat_attr = HAT_MERGING_OK | HAT_PLAT_NOCACHE;
112 else /* GFXP_MEMORY_UNCACHED */
113     hat_attr = HAT_STRICTORDER | HAT_PLAT_NOCACHE;
114 hat_flags = HAT_LOAD_LOCK;
115 poffset = start & PAGEOFFSET;
116 base = start - poffset;
117 npages = btopr(size + poffset);
118 cvaddr = vmem_alloc(heap_arena, ptob(npages), VM_NOSLEEP);
119 if (cvaddr == NULL)
120     return (NULL);

122 #ifdef __xpv
123 ASSERT(DOMAIN_IS_INITDOMAIN(xen_info));
124 pfn = xen_assign_pfn(mmu_btop(base));
```

```
125 #else
126     pfn = btop(base);
127 #endif

129     hat_devload(kas.a_hat, cvaddr, ptob(npages), pfn,
130                PROT_READ|PROT_WRITE|hat_attr, hat_flags);
131     return (cvaddr + pgoffset);
132 }
    unchanged_portion_omitted

300 /*
301  * Like gfxp_map_kernel_space, but
302  * just the hat_devload part.
303  */
304 void
305 gfxp_load_kernel_space(uint64_t start, size_t size,
306                        uint32_t mode, caddr_t cvaddr)
307 {
308     uint_t pgoffset;
309     uint64_t base;
310     pgcnt_t npages;
311     int hat_flags;
312     uint_t hat_attr;
313     pfn_t pfn;

315     if (size == 0)
316         return;

318 #ifdef __xpv
319     /*
320      * The hypervisor doesn't allow r/w mappings to some pages, such as
321      * page tables, gdt, etc. Detect %cr3 to notify users of this interface.
322      */
323     if (start == mmu_ptob(mmu_btop(getcr3_pa())))
324         if (start == mmu_ptob(mmu_btop(getcr3())))
325             return;
326 #endif

327     if (mode == GFXP_MEMORY_CACHED)
328         hat_attr = HAT_STORECACHING_OK;
329     else if (mode == GFXP_MEMORY_WRITECOMBINED)
330         hat_attr = HAT_MERGING_OK | HAT_PLAT_NOCACHE;
331     else
332         /* GFXP_MEMORY_UNCACHED */
333         hat_attr = HAT_STRICTORDER | HAT_PLAT_NOCACHE;

335     hat_flags = HAT_LOAD_LOCK;

337     pgoffset = start & PAGEOFFSET;
338     base = start - pgoffset;
339     npages = btopr(size + pgoffset);

341 #ifdef __xpv
342     ASSERT(DOMAIN_IS_INITDOMAIN(xen_info));
343     pfn = xen_assign_pfn(mmu_btop(base));
344 #else
345     pfn = btop(base);
346 #endif

348     hat_devload(kas.a_hat, cvaddr, ptob(npages), pfn,
349                PROT_READ|PROT_WRITE|hat_attr, hat_flags);
350 }
    unchanged_portion_omitted
```

new/usr/src/uts/i86pc/ml/fb_sw_tch_src.s

1

```
*****
9444 Fri Apr 6 17:25:00 2018
new/usr/src/uts/i86pc/ml/fb_sw_tch_src.s
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */

22 /*
23  * Copyright 2008 Sun Microsystems, Inc. All rights reserved.
24  * Use is subject to license terms.
25  * Copyright 2018 Joyent, Inc.
26 */

29 #if defined(__lint)

31 int fb_sw_tch_silence_lint = 0;

33 #else

35 #include <sys/asm_linkage.h>
36 #include <sys/segments.h>
37 #include <sys/controlregs.h>
38 #include <sys/machparam.h>
39 #include <sys/multiboot.h>
40 #include <sys/fastboot.h>
41 #include "assym.h"

43 /*
44  * This code is to switch from 64-bit or 32-bit to protected mode.
45  */

47 /*
48  * For debugging with LEDs
49  */
50 #define FB_OUTB_ASM(val)      \
51     movb    val, %al;        \
52     outb   $0x80;

55 #define DISABLE_PAGING      \
56     movl   %cr4, %eax;      \
57     btrl  $17, %eax;        /* clear PCIDE bit */
58     movl  %eax, %cr4;      \
59     movl  %cr0, %eax;      \
```

new/usr/src/uts/i86pc/ml/fb_sw_tch_src.s

2

```
60     btrl  $31, %eax;        /* clear PG bit */
61     movl  %eax, %cr0;

63 /*
64  * This macro contains common code for 64/32-bit versions of copy_sections().
65  * On entry:
66  *   fbf points to the fboot_file_t
67  *   snum contains the number of sections
68  * Registers that would be clobbered:
69  *   fbs, snum, %eax, %ecx, %edi, %esi.
70  * NOTE: fb_dest_pa is supposed to be in the first 1GB,
71  * therefore it is safe to use 32-bit register to hold it's value
72  * even for 64-bit code.
73  */

75 #define COPY_SECT(fbf, fbs, snum)      \
76     lea   FB_SECTIONS(fbf), fbs;      \
77     xorl  %eax, %eax;                  \
78 1:    movl  FB_DEST_PA(fbf), %esi;      \
79     addl  FB_SEC_OFFSET(fbs), %esi;    \
80     movl  FB_SEC_PADDR(fbs), %edi;     \
81     movl  FB_SEC_SIZE(fbs), %ecx;     \
82     rep                                     \
83     movsb;                                \
84     /* Zero BSS */                       \
85     movl  FB_SEC_BSS_SIZE(fbs), %ecx;  \
86     rep                                     \
87     stosb;                                \
88     add   $FB_SECTIONS_INCR, fbs;      \
89     dec  snum;                            \
90     jnz  1b;

93     .globl _start
94 _start:

96     /* Disable interrupts */
97     cli

99 #if defined(__amd64)
100 /* Switch to a low memory stack */
101     movq  $_start, %rsp
102     addq  $FASTBOOT_STACK_OFFSET, %rsp

104 /*
105  * Copy from old stack to new stack
106  * If the content before fi_valid gets bigger than 0x200 bytes,
107  * the reserved stack size above will need to be changed.
108  */
109     movq  %rdi, %rsi; /* source from old stack */
110     movq  %rsp, %rdi; /* destination on the new stack */
111     movq  $FI_VALID, %rcx /* size to copy */
112     rep                                     \
113     smovb

115 #elif defined(__i386)
116     movl  0x4(%esp), %esi /* address of fastboot info struct */

118 /* Switch to a low memory stack */
119     movl  $_start, %esp
120     addl  $FASTBOOT_STACK_OFFSET, %esp

122 /* Copy struct to stack */
123     movl  %esp, %edi; /* destination on the new stack */
124     movl  $FI_VALID, %ecx /* size to copy */
125     rep
```

```

126         smovb
128 #endif
130 #if defined(__amd64)
132         xorl    %eax, %eax
133         xorl    %edx, %edx
135         movl    $MSR_AMD_FSBASE, %ecx
136         wrmsr
138         movl    $MSR_AMD_GSBASE, %ecx
139         wrmsr
141         movl    $MSR_AMD_KGSBASE, %ecx
142         wrmsr
144 #endif
145         /*
146          * zero out all the registers to make sure they're 16 bit clean
147          */
148 #if defined(__amd64)
149         xorq    %r8, %r8
150         xorq    %r9, %r9
151         xorq    %r10, %r10
152         xorq    %r11, %r11
153         xorq    %r12, %r12
154         xorq    %r13, %r13
155         xorq    %r14, %r14
156         xorq    %r15, %r15
157 #endif
158         xorl    %eax, %eax
159         xorl    %ebx, %ebx
160         xorl    %ecx, %ecx
161         xorl    %edx, %edx
162         xorl    %ebp, %ebp
164 #if defined(__amd64)
165         /*
166          * Load our own GDT
167          */
168         lgdt    gdt_info
169 #endif
170         /*
171          * Load our own IDT
172          */
173         lidt    idt_info
175 #if defined(__amd64)
176         /*
177          * Invalidate all TLB entries.
178          * Load temporary pagetables to copy kernel and boot-archive
179          */
180         movq    %cr4, %rax
181         andq    $_BITNOT(CR4_PGE), %rax
182         movq    %rax, %cr4
183         movq    FI_PAGETABLE_PA(%rsp), %rax
184         movq    %rax, %cr3
186         leaq    FI_FILES(%rsp), %rbx    /* offset to the files */
188         /* copy unix to final destination */
189         movq    FI_LAST_TABLE_PA(%rsp), %rsi    /* page table PA */
190         leaq    _MUL(FASTBOOT_UNIX, FI_FILES_INCR)(%rbx), %rdi
191         call    map_copy

```

```

193         /* copy boot archive to final destination */
194         movq    FI_LAST_TABLE_PA(%rsp), %rsi    /* page table PA */
195         leaq    _MUL(FASTBOOT_BOOTARCHIVE, FI_FILES_INCR)(%rbx), %rdi
196         call    map_copy
198         /* Copy sections if there are any */
199         leaq    _MUL(FASTBOOT_UNIX, FI_FILES_INCR)(%rbx), %rdi
200         movl    FB_SECTCNT(%rdi), %esi
201         cmpl    $0, %esi
202         je     lf
203         call    copy_sections
204 1:
205         /*
206          * Shut down 64 bit mode. First get into compatibility mode.
207          */
208         movq    %rsp, %rax
209         pushq   $B32DATA_SEL
210         pushq   %rax
211         pushf
212         pushq   $B32CODE_SEL
213         pushq   $1f
214         iretq
216         .code32
217 1:
218         movl    $B32DATA_SEL, %eax
219         movw    %ax, %ss
220         movw    %ax, %ds
221         movw    %ax, %es
222         movw    %ax, %fs
223         movw    %ax, %gs
225         /*
226          * Disable long mode by:
227          * - shutting down paging (bit 31 of cr0). This will flush the
228          *   TLBs.
229          * - turning off PCID in cr4
230          * - disabling LME (long mode enable) in EFER (extended feature reg)
231          */
232 #endif
233         DISABLE_PAGING    /* clobbers %eax */
235 #if defined(__amd64)
236         ljmp    $B32CODE_SEL, $1f
237 1:
238 #endif
240         /*
241          * Clear PGE, PAE and PSE flags as dboot expects them to be
242          * cleared.
243          */
244         movl    %cr4, %eax
245         andl    $_BITNOT(CR4_PGE | CR4_PAE | CR4_PSE), %eax
246         movl    %eax, %cr4
248 #if defined(__amd64)
249         movl    $MSR_AMD_EFER, %ecx    /* Extended Feature Enable */
250         rdmsr
251         btcl    $8, %eax    /* bit 8 Long Mode Enable bit */
252         wrmsr
254 #elif defined(__i386)
255         /*
256          * If fi_has_pae is set, re-enable paging with PAE.
257          */

```



```

258     leal   FI_FILES(%esp), %ebx    /* offset to the files */
259     movl   FI_HAS_PAE(%esp), %edi  /* need to enable paging or not */
260     cmpl   $0, %edi
261     je     paging_on              /* no need to enable paging */

263     movl   FI_LAST_TABLE_PA(%esp), %esi /* page table PA */

265     /*
266     * Turn on PAE
267     */
268     movl   %cr4, %eax
269     orl   $CR4_PAE, %eax
270     movl   %eax, %cr4

272     /*
273     * Load top pagetable base address into cr3
274     */
275     movl   FI_PAGETABLE_PA(%esp), %eax
276     movl   %eax, %cr3

278     movl   %cr0, %eax
279     orl   $_CONST(CR0_PG | CR0_WP | CR0_AM), %eax
280     andl  $_BITNOT(CR0_NW | CR0_CD), %eax
281     movl   %eax, %cr0
282     jmp   paging_on
283 paging_on:

285     /* copy unix to final destination */
286     leal   _MUL(FASTBOOT_UNIX, FI_FILES_INCR)(%ebx), %edx
287     call  map_copy

289     /* copy boot archive to final destination */
290     leal   _MUL(FASTBOOT_BOOTARCHIVE, FI_FILES_INCR)(%ebx), %edx
291     call  map_copy

293     /* Disable paging one more time */
294     DISABLE_PAGING

296     /* Copy sections if there are any */
297     leal   _MUL(FASTBOOT_UNIX, FI_FILES_INCR)(%ebx), %edx
298     movl   FB_SECTCNT(%edx), %eax
299     cmpl   $0, %eax
300     je     1f
301     call  copy_sections
302 1:

304     /* Whatever flags we turn on we need to turn off */
305     movl   %cr4, %eax
306     andl  $_BITNOT(CR4_PAE), %eax
307     movl   %eax, %cr4
308 #endif /* __i386 */

310 dboot_jump:
311     /* Jump to dboot */
312     movl   $DBOOT_ENTRY_ADDRESS, %edi
313     movl   FI_NEW_MBI_PA(%esp), %ebx
314     movl   $MIBOOTLOADER_MAGIC, %eax
315     jmp   *%edi

317 #if defined(__amd64)

319     .code64
320     ENTRY_NP(copy_sections)
321     /*
322     * On entry
323     *     %rdi points to the fboot_file_t

```

```

324     *     %rsi contains number of sections
325     */
326     movq   %rdi, %rdx
327     movq   %rsi, %r9

329     COPY_SECT(%rdx, %r8, %r9)
330     ret
331     SET_SIZE(copy_sections)
_____unchanged_portion_omitted_____

```

new/usr/src/uts/i86pc/ml/genassym.c

1

```
*****
5723 Fri Apr 6 17:25:00 2018
new/usr/src/uts/i86pc/ml/genassym.c
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22 * Copyright (c) 1992, 2010, Oracle and/or its affiliates. All rights reserved.
23 *
24 * Copyright 2018 Joyent, Inc.
25 */

27 #ifndef _GENASSYM
28 #define _GENASSYM
29 #endif

31 #define exit    kern_exit

33 #include <sys/types.h>
34 #include <sys/param.h>
35 #include <sys/system.h>
36 #include <sys/elf_notes.h>
37 #include <sys/thread.h>
38 #include <sys/rwlock.h>
39 #include <sys/proc.h>
40 #include <sys/cpuvar.h>
41 #include <sys/clock.h>
42 #include <sys/trap.h>
43 #include <sys/modctl.h>
44 #include <sys/traptrace.h>
45 #include <vm/seg.h>
46 #include <sys/avintr.h>
47 #include <sys/pic.h>
48 #include <sys/pit.h>
49 #include <sys/fp.h>
50 #include <sys/disp.h>
51 #include <sys/archsystem.h>
52 #include <sys/x86_archext.h>
53 #include <sys/sunddi.h>
54 #include <sys/mach_mmu.h>

56 #if defined(__xpv)
57 #include <sys/hypervisor.h>
58 #endif
```

new/usr/src/uts/i86pc/ml/genassym.c

2

```
60 #undef  exit    /* unhide exit, see comment above */
61 extern void exit(int);

63 /*
64  * Proactively discourage anyone from referring to structures or
65  * member offsets in this program.
66  */
67 #define struct    struct...
68 #define OFFSET    OFFSET...

70 int
71 main(int argc, char *argv[])
72 {
73     printf("#define\tT_AST 0x%x\n", T_AST);

74     printf("#define\tLOCK_LEVEL 0x%x\n", LOCK_LEVEL);
75     printf("#define\tCLOCK_LEVEL 0x%x\n", CLOCK_LEVEL);
76     printf("#define\tDISP_LEVEL 0x%x\n", DISP_LEVEL);
77     printf("#define\tPIL_MAX 0x%x\n", PIL_MAX);
78     printf("#define\tHIGH_LEVELS 0x%x\n", HIGH_LEVELS);
79     printf("#define\tCPU_INTR_ACTV_HIGH_LEVEL_MASK 0x%x\n",
80            CPU_INTR_ACTV_HIGH_LEVEL_MASK);

81     printf("#define\tPIC_NSEOI 0x%x\n", PIC_NSEOI);
82     printf("#define\tPIC_SEOI_LVL7 0x%x\n", PIC_SEOI_LVL7);

84     printf("#define\tNANOSEC 0x%llx\n", NANOSEC);
85     printf("#define\tADJ_SHIFT 0x%x\n", ADJ_SHIFT);

87     printf("#define\tSSLEEP 0x%x\n", SSLEEP);
88     printf("#define\tSRUN 0x%x\n", SRUN);
89     printf("#define\tSONPROC 0x%x\n", SONPROC);

91     printf("#define\tT_INTR_THREAD 0x%x\n", T_INTR_THREAD);
92     printf("#define\tFREE_THREAD 0x%x\n", TS_FREE);
93     printf("#define\tTS_FREE 0x%x\n", TS_FREE);
94     printf("#define\tTS_ZOMB 0x%x\n", TS_ZOMB);
95     printf("#define\tTP_MSACCT 0x%x\n", TP_MSACCT);
96     printf("#define\tTP_WATCHPT 0x%x\n", TP_WATCHPT);
97     printf("#define\tONPROC_THREAD 0x%x\n", TS_ONPROC);

99     printf("#define\tS_READ 0x%x\n", (int)S_READ);
100    printf("#define\tS_WRITE 0x%x\n", (int)S_WRITE);
101    printf("#define\tS_EXEC 0x%x\n", (int)S_EXEC);
102    printf("#define\tS_OTHER 0x%x\n", (int)S_OTHER);

104    printf("#define\tNORMALRETURN 0x%x\n", (int)NORMALRETURN);
105    printf("#define\tLWP_USER 0x%x\n", LWP_USER);
106    printf("#define\tLWP_SYS 0x%x\n", LWP_SYS);
107    printf("#define\tLMS_USER 0x%x\n", LMS_USER);
108    printf("#define\tLMS_SYSTEM 0x%x\n", LMS_SYSTEM);

110    printf("#define\tSSE_MXCSR_EFLAGS 0x%x\n", SSE_MXCSR_EFLAGS);

112    printf("#define\tFP_487 0x%x\n", FP_487);
113    printf("#define\tFP_486 0x%x\n", FP_486);
114    printf("#define\tFPU_CW_INIT 0x%x\n", FPU_CW_INIT);
115    printf("#define\tFPU_EN 0x%x\n", FPU_EN);
116    printf("#define\tFPU_VALID 0x%x\n", FPU_VALID);

118    printf("#define\tFP_NO 0x%x\n", FP_NO);
119    printf("#define\tFP_SW 0x%x\n", FP_SW);
120    printf("#define\tFP_HW 0x%x\n", FP_HW);
121    printf("#define\tFP_287 0x%x\n", FP_287);
122    printf("#define\tFP_387 0x%x\n", FP_387);
123    printf("#define\t__FP_SSE 0x%x\n", __FP_SSE);
```

```

125     printf("#define\tFP_FNSAVE 0x%x\n", FP_FNSAVE);
126     printf("#define\tFP_FXSAVE 0x%x\n", FP_FXSAVE);
127     printf("#define\tFP_XSAVE 0x%x\n", FP_XSAVE);
128
129     printf("#define\tAV_INT_SPURIOUS 0x%x\n", AV_INT_SPURIOUS);
130
131     printf("#define\tCPU_READY 0x%x\n", CPU_READY);
132     printf("#define\tCPU QUIESCED 0x%x\n", CPU QUIESCED);
133
134     printf("#define\tMCMD_PORT 0x%x\n", MCMD_PORT);
135     printf("#define\tSCMD_PORT 0x%x\n", SCMD_PORT);
136     printf("#define\tMIMR_PORT 0x%x\n", MIMR_PORT);
137     printf("#define\tSIMR_PORT 0x%x\n", SIMR_PORT);
138
139     printf("#define\tDMP_NOSYNC 0x%x\n", DMP_NOSYNC);
140
141     printf("#define\tRW_WRITER\t0x%x\n", RW_WRITER);
142     printf("#define\tRW_READER\t0x%x\n", RW_READER);
143
144     printf("#define\tNSYSCALL 0x%x\n", NSYSCALL);
145
146     printf("#define\tSE_32RVAL1 0x%x\n", SE_32RVAL1);
147     printf("#define\tSE_32RVAL2 0x%x\n", SE_32RVAL2);
148     printf("#define\tSE_64RVAL 0x%x\n", SE_64RVAL);
149
150     printf("#define\tMAXSYSARGS 0x%x\n", MAXSYSARGS);
151
152     /* Hack value just to allow clock to be kicked */
153     printf("#define\tNSEC_PER_CLOCK_TICK 0x%llx\n", NANOSEC / 100);
154
155     printf("#define\tNSEC_PER_COUNTER_TICK 0x%llx\n", NANOSEC / PIT_HZ);
156
157     printf("#define\tPITCTRO_PORT 0x%x\n", PITCTRO_PORT);
158     printf("#define\tPITCTL_PORT 0x%x\n", PITCTL_PORT);
159     printf("#define\tPIT_COUNTDOWN 0x%x\n",
160           PIT_CO | PIT_LOADMODE | PIT_NDIVMODE);
161
162     printf("#define\tNBPW 0x%x\n", NBPW);
163
164     printf("#define\tDDI_ACCATTR_IO_SPACE 0x%x\n", DDI_ACCATTR_IO_SPACE);
165     printf("#define\tDDI_ACCATTR_DIRECT 0x%x\n", DDI_ACCATTR_DIRECT);
166     printf("#define\tDDI_ACCATTR_CPU_VADDR 0x%x\n", DDI_ACCATTR_CPU_VADDR);
167     printf("#define\tDDI_DEV_AUTOINCR 0x%x\n", DDI_DEV_AUTOINCR);
168
169     printf("#define\tMMU_STD_PAGESIZE 0x%x\n", (uint_t)MMU_STD_PAGESIZE);
170     printf("#define\tMMU_STD_PAGEMASK 0x%x\n", (uint_t)MMU_STD_PAGEMASK);
171     printf("#define\tFOUR_MEG 0x%x\n", (uint_t)FOUR_MEG);
172
173     printf("#define\tTRAPTR_NENT 0x%x\n", TRAPTR_NENT);
174
175     printf("#define\tCPU_DTRACE_NOFAULT 0x%x\n", CPU_DTRACE_NOFAULT);
176     printf("#define\tCPU_DTRACE_BADADDR 0x%x\n", CPU_DTRACE_BADADDR);
177     printf("#define\tCPU_DTRACE_DIVZERO 0x%x\n", CPU_DTRACE_DIVZERO);
178     printf("#define\tCPU_DTRACE_ILLOP 0x%x\n", CPU_DTRACE_ILLOP);
179
180     printf("#define\tMODS_NOUNLOAD 0x%x\n", MODS_NOUNLOAD);
181     printf("#define\tMODS_WEAK 0x%x\n", MODS_WEAK);
182     printf("#define\tMODS_INSTALLED 0x%x\n", MODS_INSTALLED);
183
184     printf("#define\tKPREEMPT_SYNC 0x%x\n", KPREEMPT_SYNC);
185
186 #if defined(__xpv)
187     printf("#define\tSHUTDOWN_reboot 0x%x\n", SHUTDOWN_reboot);
188     printf("#define\tSCHEDOP_block 0x%x\n", SCHEDOP_block);
189     printf("#define\tVGCF_IN_KERNEL 0x%x\n", VGCF_IN_KERNEL);

```

```

171 #endif
172     return (0);
173 }
_____unchanged_portion_omitted_

```

```

*****
23312 Fri Apr 6 17:25:00 2018
new/usr/src/uts/i86pc/ml/kpti_trampoline.s
9441 kmdb should stash %cr3 in kdiregs
Reviewed by: John Levon <john.levon@joyent.com>
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>
*****
1 /*
2  * This file and its contents are supplied under the terms of the
3  * Common Development and Distribution License ("CDDL"), version 1.0.
4  * You may only use this file in accordance with the terms of version
5  * 1.0 of the CDDL.
6  *
7  * A full copy of the text of the CDDL should have accompanied this
8  * source. A copy of the CDDL is also available via the Internet at
9  * http://www.illumos.org/license/CDDL.
10 */
11 /*
12  * Copyright 2018 Joyent, Inc.
13 */
14
15 /*
16  * This file contains the trampolines that are used by KPTI in order to be
17  * able to take interrupts/trap/etc while on the "user" page table.
18  *
19  * We don't map the full kernel text into the user page table: instead we
20  * map this one small section of trampolines (which compiles to ~13 pages).
21  * These trampolines are set in the IDT always (so they will run no matter
22  * whether we're on the kernel or user page table), and their primary job is to
23  * pivot us to the kernel %cr3 and %rsp without ruining everything.
24  *
25  * All of these interrupts use the amd64 IST feature when we have KPTI enabled,
26  * meaning that they will execute with their %rsp set to a known location, even
27  * if we take them in the kernel.
28  *
29  * Over in desctbls.c (for cpu0) and mp_pc.c (other cpus) we set up the IST
30  * stack to point at &cpu->cpu_m.mcpu_kpti.kf_tr_rsp. You can see the mcpu_kpti
31  * (a struct kpti_frame) defined in machcpuvar.h. This struct is set up to be
32  * page-aligned, and we map the page it's on into both page tables. Using a
33  * struct attached to the cpu_t also means that we can use %rsp-relative
34  * addressing to find anything on the cpu_t, so we don't have to touch %gs or
35  * %GSBASE at all on incoming interrupt trampolines (which can get pretty hairy).
36  *
37  * This little struct is where the CPU will push the actual interrupt frame.
38  * Then, in the trampoline, we change %cr3, then figure out our destination
39  * stack pointer and "pivot" to it (set %rsp and re-push the CPU's interrupt
40  * frame). Then we jump to the regular ISR in the kernel text and carry on as
41  * normal.
42  *
43  * We leave the original frame and any spilled regs behind in the kpti_frame
44  * lazily until we want to return to userland. Then, we clear any spilled
45  * regs from it, and overwrite the rest with our iret frame. When switching
46  * this cpu to a different process (in hat_switch), we bzero the whole region to
47  * make sure nothing can leak between processes.
48  *
49  * When we're returning back to the original place we took the interrupt later
50  * (especially if it was in userland), we have to jmp back to the "return
51  * trampolines" here, since when we set %cr3 back to the user value, we need to
52  * be executing from code here in these shared pages and not the main kernel
53  * text again. Even though it should be fine to iret directly from kernel text
54  * when returning to kernel code, we make things jmp to a trampoline here just
55  * for consistency.
56  *
57  * Note that with IST, it's very important that we always must have pivoted

```

```

58  * away from the IST stack before we could possibly take any other interrupt
59  * on the same IST (unless it's an end-of-the-world fault and we don't care
60  * about coming back from it ever).
61  *
62  * This is particularly relevant to the dbgtrap/brktrap trampolines, as they
63  * regularly have to happen from within trampoline code (e.g. in the sysenter
64  * single-step case) and then return to the world normally. As a result, these
65  * two are IST'd to their own kpti_frame right above the normal one (in the same
66  * page), so they don't clobber their parent interrupt.
67  *
68  * To aid with debugging, we also IST the page fault (#PF/pftrap), general
69  * protection fault (#GP/gptrap) and stack fault (#SS/stktrap) interrupts to
70  * their own separate kpti_frame. This ensures that if we take one of these
71  * due to a bug in trampoline code, we preserve the original trampoline
72  * state that caused the trap.
73  *
74  * NMI, MCE and dblfault interrupts also are taken on their own dedicated IST
75  * stacks, since they can interrupt another ISR at any time. These stacks are
76  * full-sized, however, and not a little kpti_frame struct. We only set %cr3 in
77  * their trampolines (and do it unconditionally), and don't bother pivoting
78  * away. We're either going into the panic() path, or we're going to return
79  * straight away without rescheduling, so it's fine to not be on our real
80  * kthread stack (and some of the state we want to go find it with might be
81  * corrupt!)
82  *
83  * Finally, for these "special" interrupts (NMI/MCE/double fault) we use a
84  * special %cr3 value we stash here in the text (kpti_safe_cr3). We set this to
85  * point at the PML4 for kas early in boot and never touch it again. Hopefully
86  * it survives whatever corruption brings down the rest of the kernel!
87  *
88  * Syscalls are different to interrupts (at least in the SYSENTER/SYSCALL64
89  * cases) in that they do not push an interrupt frame (and also have some other
90  * effects). In the syscall trampolines, we assume that we can only be taking
91  * the call from userland and use SWAPGS and an unconditional overwrite of %cr3.
92  * We do not do any stack pivoting for syscalls (and we leave SYSENTER's
93  * existing %rsp pivot untouched) -- instead we spill registers into
94  * %gs:CPU_KPTI_* as we need to.
95  *
96  * Note that the normal %cr3 values do not cause invalidations with PCIDE - see
97  * hat_switch().
98  */
99
100 /*
101  * The macros here mostly line up with what's in kdi_idthdl.s, too, so if you
102  * fix bugs here check to see if they should be fixed there as well.
103  */
104
105 #include <sys/asm_linkage.h>
106 #include <sys/asm_misc.h>
107 #include <sys/regset.h>
108 #include <sys/privregs.h>
109 #include <sys/psw.h>
110 #include <sys/machbrand.h>
111 #include <sys/param.h>
112
113 #if defined(__lint)
114
115 #include <sys/types.h>
116 #include <sys/thread.h>
117 #include <sys/system.h>
118
119 #else /* __lint */
120
121 #include <sys/segments.h>
122 #include <sys/pcb.h>
123 #include <sys/trap.h>

```

```

124 #include <sys/ftrace.h>
125 #include <sys/traptrace.h>
126 #include <sys/clock.h>
127 #include <sys/model.h>
128 #include <sys/panic.h>

130 #if defined(__xpv)
131 #include <sys/hypervisor.h>
132 #endif

134 #include "assym.h"

136 .data
137     DGDEF3(kpti_enable, 8, 8)
138     .fill 1, 8, 1

140 .section ".text";
141 .align MMU_PAGESIZE

143 .global kpti_trap_start
144 kpti_trap_start:
145     nop

147 /* This will be set by mlsetup, and then double-checked later */
148 .global kpti_safe_cr3
149 kpti_safe_cr3:
150     .quad 0
151     SET_SIZE(kpti_safe_cr3)

153 /* startup_kmem() will overwrite this */
154 .global kpti_kbase
155 kpti_kbase:
156     .quad KERNELBASE
157     SET_SIZE(kpti_kbase)

159 #define SET_KERNEL_CR3(spillreg) \
160     mov    %cr3, spillreg; \
161     mov    spillreg, %gs:CPU_KPTI_TR_CR3; \
162     mov    %gs:CPU_KPTI_KCR3, spillreg; \
163     cmp    $0, spillreg; \
164     je     2f; \
165     mov    spillreg, %cr3; \
166 2:

168 #if DEBUG
169 #define SET_USER_CR3(spillreg) \
170     mov    %cr3, spillreg; \
171     mov    spillreg, %gs:CPU_KPTI_TR_CR3; \
172     mov    %gs:CPU_KPTI_UCR3, spillreg; \
173     mov    spillreg, %cr3
174 #else
175 #define SET_USER_CR3(spillreg) \
176     mov    %gs:CPU_KPTI_UCR3, spillreg; \
177     mov    spillreg, %cr3
178 #endif

180 #define PIVOT_KPTI_STK(spillreg) \
181     mov    %rsp, spillreg; \
182     mov    %gs:CPU_KPTI_RET_RSP, %rsp; \
183     pushq T_FRAMERET_SS(spillreg); \
184     pushq T_FRAMERET_RSP(spillreg); \
185     pushq T_FRAMERET_RFLAGS(spillreg); \
186     pushq T_FRAMERET_CS(spillreg); \
187     pushq T_FRAMERET_RIP(spillreg)

```

```

190 #define INTERRUPT_TRAMPOLINE_P(errpush) \
191     pushq %r13; \
192     pushq %r14; \
193     subq  $KPTI_R14, %rsp; \
194     /* Save current %cr3. */ \
195     mov    %cr3, %r14; \
196     mov    %r14, KPTI_TR_CR3(%rsp); \
197     \
198     cmpw  $KCS_SEL, KPTI_CS(%rsp); \
199     je     3f; \
200 1: \
201     /* Change to the "kernel" %cr3 */ \
202     mov    KPTI_KCR3(%rsp), %r14; \
203     cmp    $0, %r14; \
204     je     2f; \
205     mov    %r14, %cr3; \
206 2: \
207     /* Get our cpu_t in %r13 */ \
208     mov    %rsp, %r13; \
209     and    $(~(MMU_PAGESIZE - 1)), %r13; \
210     subq  $CPU_KPTI_START, %r13; \
211     /* Use top of the kthread stk */ \
212     mov    CPU_THREAD(%r13), %r14; \
213     mov    T_STACK(%r14), %r14; \
214     addq  $REGSIZE+MINFRAME, %r14; \
215     jmp   4f; \
216 3: \
217     /* Check the %rsp in the frame. */ \
218     /* Is it above kernel base? */ \
219     mov    kpti_kbase, %r14; \
220     cmp    %r14, KPTI_RSP(%rsp); \
221     jb     1b; \
222     /* Use the %rsp from the trap frame */ \
223     mov    KPTI_RSP(%rsp), %r14; \
224     and    $(~0xf), %r14; \
225 4: \
226     mov    %rsp, %r13; \
227     /* %r14 contains our destination stk */ \
228     mov    %r14, %rsp; \
229     pushq KPTI_SS(%r13); \
230     pushq KPTI_RSP(%r13); \
231     pushq KPTI_RFLAGS(%r13); \
232     pushq KPTI_CS(%r13); \
233     pushq KPTI_RIP(%r13); \
234     errpush; \
235     mov    KPTI_R14(%r13), %r14; \
236     mov    KPTI_R13(%r13), %r13

238 #define INTERRUPT_TRAMPOLINE_NOERR \
239     INTERRUPT_TRAMPOLINE_P(/**/)

241 #define INTERRUPT_TRAMPOLINE \
242     INTERRUPT_TRAMPOLINE_P(pushq KPTI_ERR(%r13))

244 /*
245  * This is used for all interrupts that can plausibly be taken inside another
246  * interrupt and are using a kpti_frame stack (so #BP, #DB, #GP, #PF, #SS).
247  *
248  * We check for whether we took the interrupt while in another trampoline, in
249  * which case we need to use the kthread stack.
250  */
251 #define DBG_INTERRUPT_TRAMPOLINE_P(errpush) \
252     pushq %r13; \
253     pushq %r14; \
254     subq  $KPTI_R14, %rsp; \
255     /* Check for clobbering */ \

```

```

256     cmp     $0, KPTI_FLAG(%rsp); \
257     je      1f; \
258     /* Don't worry, this totally works */ \
259     int     $8; \
260 1: \
261     movq    $1, KPTI_FLAG(%rsp); \
262     /* Save current %cr3. */ \
263     mov     %cr3, %r14; \
264     mov     %r14, KPTI_TR_CR3(%rsp); \
265 \
266     cmpw    $KCS_SEL, KPTI_CS(%rsp); \
267     je      4f; \
268 2: \
269     /* Change to the "kernel" %cr3 */ \
270     mov     KPTI_KCR3(%rsp), %r14; \
271     cmp     $0, %r14; \
272     je      3f; \
273     mov     %r14, %cr3; \
274 3: \
275     /* Get our cpu_t in %r13 */ \
276     mov     %rsp, %r13; \
277     and     $(~(MMU_PAGESIZE - 1)), %r13; \
278     subq    $CPU_KPTI_START, %r13; \
279     /* Use top of the kthread stk */ \
280     mov     CPU_THREAD(%r13), %r14; \
281     mov     T_STACK(%r14), %r14; \
282     addq    $REGSIZE+MINFRAME, %r14; \
283     jmp     6f; \
284 4: \
285     /* Check the %rsp in the frame. */ \
286     /* Is it above kernel base? */ \
287     /* If not, treat as user. */ \
288     mov     kpti_kbase, %r14; \
289     cmp     %r14, KPTI_RSP(%rsp); \
290     jb      2b; \
291     /* Is it within the kpti_frame page? */ \
292     /* If it is, treat as user interrupt */ \
293     mov     %rsp, %r13; \
294     and     $(~(MMU_PAGESIZE - 1)), %r13; \
295     mov     KPTI_RSP(%rsp), %r14; \
296     and     $(~(MMU_PAGESIZE - 1)), %r14; \
297     cmp     %r13, %r14; \
298     je      2b; \
299     /* Were we in trampoline code? */ \
300     leaq   kpti_trampoline_start, %r14; \
301     cmp     %r14, KPTI_RIP(%rsp); \
302     jb      5f; \
303     leaq   kpti_trampoline_end, %r14; \
304     cmp     %r14, KPTI_RIP(%rsp); \
305     ja      5f; \
306     /* If we were, change %cr3: we might */ \
307     /* have interrupted before it did. */ \
308     mov     KPTI_KCR3(%rsp), %r14; \
309     mov     %r14, %cr3; \
310 5: \
311     /* Use the %rsp from the trap frame */ \
312     mov     KPTI_RSP(%rsp), %r14; \
313     and     $(~0xf), %r14; \
314 6: \
315     mov     %rsp, %r13; \
316     /* %r14 contains our destination stk */ \
317     mov     %r14, %rsp; \
318     pushq   KPTI_SS(%r13); \
319     pushq   KPTI_RSP(%r13); \
320     pushq   KPTI_RFLAGS(%r13); \
321     pushq   KPTI_CS(%r13); \

```

```

322     pushq   KPTI_RIP(%r13); \
323     errpush; \
324     mov     KPTI_R14(%r13), %r14; \
325     movq    $0, KPTI_FLAG(%r13); \
326     mov     KPTI_R13(%r13), %r13 \
\
328 #define DBG_INTERRUPT_TRAMPOLINE_NOERR \
329     DBG_INTERRUPT_TRAMPOLINE_P(/**/) \
\
331 #define DBG_INTERRUPT_TRAMPOLINE \
332     DBG_INTERRUPT_TRAMPOLINE_P(pushq KPTI_ERR(%r13)) \
\
334     /* \
335     * These labels (_start and _end) are used by trap.c to determine if \
336     * we took an interrupt like an NMI during the return process. \
337     */ \
338 .global tr_sysc_ret_start \
339 tr_sysc_ret_start: \
\
341     /* \
342     * Syscall return trampolines. \
343     * \
344     * These are expected to be called on the kernel %gs. tr_sysret[ql] are \
345     * called after %rsp is changed back to the user value, so we have no \
346     * stack to work with. tr_sysexit has a kernel stack (but has to \
347     * preserve rflags, soooo). \
348     */ \
349     ENTRY_NP(tr_sysretq) \
350     cmpq    $1, kpti_enable \
351     jne     1f \
\
353     mov     %r13, %gs:CPU_KPTI_R13 \
354     SET_USER_CR3(%r13) \
355     mov     %gs:CPU_KPTI_R13, %r13 \
356     /* Zero these to make sure they didn't leak from a kernel trap */ \
357     movq    $0, %gs:CPU_KPTI_R13 \
358     movq    $0, %gs:CPU_KPTI_R14 \
359 1: \
360     swapgs \
361     sysretq \
362     SET_SIZE(tr_sysretq) \
\
364     ENTRY_NP(tr_sysretl) \
365     cmpq    $1, kpti_enable \
366     jne     1f \
\
368     mov     %r13, %gs:CPU_KPTI_R13 \
369     SET_USER_CR3(%r13) \
370     mov     %gs:CPU_KPTI_R13, %r13 \
371     /* Zero these to make sure they didn't leak from a kernel trap */ \
372     movq    $0, %gs:CPU_KPTI_R13 \
373     movq    $0, %gs:CPU_KPTI_R14 \
374 1: \
375     SWAPGS \
376     SYSRETL \
377     SET_SIZE(tr_sysretl) \
\
379     ENTRY_NP(tr_sysexit) \
380     /* \
381     * Note: we want to preserve RFLAGS across this branch, since sysexit \
382     * (unlike sysret above) does not restore RFLAGS for us. \
383     * \
384     * We still have the real kernel stack (sysexit does restore that), so \
385     * we can use pushfq/popfq. \
386     */ \
387     pushfq \

```

```

389     cmpq    $1, kpti_enable
390     jne     lf

392     /* Have to pop it back off now before we change %cr3! */
393     popfq
394     mov     %r13, %gs:CPU_KPTI_R13
395     SET_USER_CR3(%r13)
396     mov     %gs:CPU_KPTI_R13, %r13
397     /* Zero these to make sure they didn't leak from a kernel trap */
398     movq    $0, %gs:CPU_KPTI_R13
399     movq    $0, %gs:CPU_KPTI_R14
400     jmp     2f
401 1:
402     popfq
403 2:
404     swapgs
405     sti
406     sysexit
407     SET_SIZE(tr_sysexit)

409 .global tr_sysc_ret_end
410 tr_sysc_ret_end:

412     /*
413     * Syscall entry trampolines.
414     */

416 #if DEBUG
417 #define MK_SYSCALL_TRAMPOLINE(isr) \
418     ENTRY_NP(tr_**/isr); \
419     swapgs; \
420     mov     %r13, %gs:CPU_KPTI_R13; \
421     mov     %cr3, %r13; \
422     mov     %r13, %gs:CPU_KPTI_TR_CR3; \
423     mov     %gs:CPU_KPTI_KCR3, %r13; \
424     mov     %r13, %cr3; \
425     mov     %gs:CPU_KPTI_R13, %r13; \
426     swapgs; \
427     jmp     isr; \
428     SET_SIZE(tr_**/isr)
429 #else
430 #define MK_SYSCALL_TRAMPOLINE(isr) \
431     ENTRY_NP(tr_**/isr); \
432     swapgs; \
433     mov     %r13, %gs:CPU_KPTI_R13; \
434     mov     %gs:CPU_KPTI_KCR3, %r13; \
435     mov     %r13, %cr3; \
436     mov     %gs:CPU_KPTI_R13, %r13; \
437     swapgs; \
438     jmp     isr; \
439     SET_SIZE(tr_**/isr)
440 #endif

442     MK_SYSCALL_TRAMPOLINE(sys_syscall)
443     MK_SYSCALL_TRAMPOLINE(sys_syscall32)
444     MK_SYSCALL_TRAMPOLINE(brand_sys_syscall)
445     MK_SYSCALL_TRAMPOLINE(brand_sys_syscall32)

447     /*
448     * SYSENTER is special. The CPU is really not very helpful when it
449     * comes to preserving and restoring state with it, and as a result
450     * we have to do all of it by hand. So, since we want to preserve
451     * RFLAGS, we have to be very careful in these trampolines to not
452     * clobber any bits in it. That means no cmpqs or branches!
453     */

```

```

454     ENTRY_NP(tr_sys_sysenter)
455     swapgs
456     mov     %r13, %gs:CPU_KPTI_R13
457 #if DEBUG
458     mov     %cr3, %r13
459     mov     %r13, %gs:CPU_KPTI_TR_CR3
460 #endif
461     mov     %gs:CPU_KPTI_KCR3, %r13
462     mov     %r13, %cr3
463     mov     %gs:CPU_KPTI_R13, %r13
464     jmp     _sys_sysenter_post_swapgs
465     SET_SIZE(tr_sys_sysenter)

467     ENTRY_NP(tr_brand_sys_sysenter)
468     swapgs
469     mov     %r13, %gs:CPU_KPTI_R13
470 #if DEBUG
471     mov     %cr3, %r13
472     mov     %r13, %gs:CPU_KPTI_TR_CR3
473 #endif
474     mov     %gs:CPU_KPTI_KCR3, %r13
475     mov     %r13, %cr3
476     mov     %gs:CPU_KPTI_R13, %r13
477     jmp     _brand_sys_sysenter_post_swapgs
478     SET_SIZE(tr_brand_sys_sysenter)

480 #define MK_SYSCALL_INT_TRAMPOLINE(isr) \
481     ENTRY_NP(tr_**/isr); \
482     swapgs; \
483     mov     %r13, %gs:CPU_KPTI_R13; \
484     SET_KERNEL_CR3(%r13); \
485     mov     %gs:CPU_THREAD, %r13; \
486     mov     T_STACK(%r13), %r13; \
487     addq    $REGSIZE+MINFRAME, %r13; \
488     mov     %r13, %rsp; \
489     pushq   %gs:CPU_KPTI_SS; \
490     pushq   %gs:CPU_KPTI_RSP; \
491     pushq   %gs:CPU_KPTI_RFLAGS; \
492     pushq   %gs:CPU_KPTI_CS; \
493     pushq   %gs:CPU_KPTI_RIP; \
494     mov     %gs:CPU_KPTI_R13, %r13; \
495     SWAPGS; \
496     jmp     isr; \
497     SET_SIZE(tr_**/isr)

499     MK_SYSCALL_INT_TRAMPOLINE(brand_sys_syscall_int)
500     MK_SYSCALL_INT_TRAMPOLINE(sys_syscall_int)

502     /*
503     * Interrupt/trap return trampolines
504     */

506 .global tr_intr_ret_start
507 tr_intr_ret_start:

509     ENTRY_NP(tr_iret_auto)
510     cmpq    $1, kpti_enable
511     jne     tr_iret_kernel
512     cmpw    $KCS_SEL, T_FRAMERET_CS(%rsp)
513     je     tr_iret_kernel
514     jmp     tr_iret_user
515     SET_SIZE(tr_iret_auto)

517     ENTRY_NP(tr_iret_kernel)
518     /*
519     * Yes, this does nothing extra. But this way we know if we see iret

```

```

520     * elsewhere, then we've failed to properly consider trampolines there.
521     */
522     iretq
523     SET_SIZE(tr_iret_kernel)

525     ENTRY_NP(tr_iret_user)
526     cmpq   $1, kpti_enable
527     jne    lf

529     swapgs
530     mov    %r13, %gs:CPU_KPTI_R13
531     PIVOT_KPTI_STK(%r13)
532     SET_USER_CR3(%r13)
533     mov    %gs:CPU_KPTI_R13, %r13
534     /* Zero these to make sure they didn't leak from a kernel trap */
535     movq   $0, %gs:CPU_KPTI_R13
536     movq   $0, %gs:CPU_KPTI_R14
537     swapgs
538 1:
539     iretq
540     SET_SIZE(tr_iret_user)

542     /*
543     * This special return trampoline is for KDI's use only (with kmdb).
544     *
545     * KDI/kmdb do not use swapgs -- they directly write the GSBASE MSR
546     * instead. This trampoline runs after GSBASE has already been changed
547     * back to the userland value (so we can't use %gs).
548     *
549     * Instead, the caller gives us a pointer to the kpti_dbg frame in %r13.
550     * The KPTI_R13 member in the kpti_dbg has already been set to what the
551     * real %r13 should be before we IRET.
552     *
553     * Additionally, KDI keeps a copy of the incoming %cr3 value when it
554     * took an interrupt, and has put that back in the kpti_dbg area for us
555     * to use, so we don't do any sniffing of %cs here. This is important
556     * so that debugging code that changes %cr3 is possible.
557     */
558     ENTRY_NP(tr_iret_kdi)
559     movq   %r14, KPTI_R14(%r13) /* %r14 has to be preserved by us */

561     movq   %rsp, %r14 /* original %rsp is pointing at IRET frame */
562     leaq   KPTI_TOP(%r13), %rsp
563     pushq  T_FRAMERET_SS(%r14)
564     pushq  T_FRAMERET_RSP(%r14)
565     pushq  T_FRAMERET_RFLAGS(%r14)
566     pushq  T_FRAMERET_CS(%r14)
567     pushq  T_FRAMERET_RIP(%r14)

569     movq   KPTI_TR_CR3(%r13), %r14
570     movq   %r14, %cr3

572     movq   KPTI_R14(%r13), %r14
573     movq   KPTI_R13(%r13), %r13 /* preserved by our caller */

575     iretq
576     SET_SIZE(tr_iret_kdi)

578 .global tr_intr_ret_end
579 tr_intr_ret_end:

581     /*
582     * Interrupt/trap entry trampolines
583     */

585     /* CPU pushed an error code, and ISR wants one */

```

```

586 #define MK_INTR_TRAMPOLINE(isr) \
587     ENTRY_NP(tr_**/isr); \
588     INTERRUPT_TRAMPOLINE; \
589     jmp    isr; \
590     SET_SIZE(tr_**/isr)

592     /* CPU didn't push an error code, and ISR doesn't want one */
593 #define MK_INTR_TRAMPOLINE_NOERR(isr) \
594     ENTRY_NP(tr_**/isr); \
595     push   $0; \
596     INTERRUPT_TRAMPOLINE_NOERR; \
597     jmp    isr; \
598     SET_SIZE(tr_**/isr)

600     /* CPU pushed an error code, and ISR wants one */
601 #define MK_DBG_INTR_TRAMPOLINE(isr) \
602     ENTRY_NP(tr_**/isr); \
603     DBG_INTERRUPT_TRAMPOLINE; \
604     jmp    isr; \
605     SET_SIZE(tr_**/isr)

607     /* CPU didn't push an error code, and ISR doesn't want one */
608 #define MK_DBG_INTR_TRAMPOLINE_NOERR(isr) \
609     ENTRY_NP(tr_**/isr); \
610     push   $0; \
611     DBG_INTERRUPT_TRAMPOLINE_NOERR; \
612     jmp    isr; \
613     SET_SIZE(tr_**/isr)

616     MK_INTR_TRAMPOLINE_NOERR(div0trap)
617     MK_DBG_INTR_TRAMPOLINE_NOERR(dbgtrap)
618     MK_DBG_INTR_TRAMPOLINE_NOERR(brktrap)
619     MK_INTR_TRAMPOLINE_NOERR(ovflotrap)
620     MK_INTR_TRAMPOLINE_NOERR(boundstrap)
621     MK_INTR_TRAMPOLINE_NOERR(invoptrap)
622     MK_INTR_TRAMPOLINE_NOERR(ndptrap)
623     MK_INTR_TRAMPOLINE(invtsstrap)
624     MK_INTR_TRAMPOLINE(segnptrap)
625     MK_DBG_INTR_TRAMPOLINE(stktrap)
626     MK_DBG_INTR_TRAMPOLINE(gptrap)
627     MK_DBG_INTR_TRAMPOLINE(pftrap)
628     MK_INTR_TRAMPOLINE_NOERR(resvtrap)
629     MK_INTR_TRAMPOLINE_NOERR(ndperr)
630     MK_INTR_TRAMPOLINE(achktrap)
631     MK_INTR_TRAMPOLINE_NOERR(xmtrap)
632     MK_INTR_TRAMPOLINE_NOERR(ivaltrap)
633     MK_INTR_TRAMPOLINE_NOERR(fasttrap)
634     MK_INTR_TRAMPOLINE_NOERR(dtrace_ret)

636     /*
637     * These are special because they can interrupt other traps, and
638     * each other. We don't need to pivot their stacks, because they have
639     * dedicated IST stack space, but we need to change %cr3.
640     */
641     ENTRY_NP(tr_nmiint)
642     pushq  %r13
643     mov    kpti_safe_cr3, %r13
644     mov    %r13, %cr3
645     popq  %r13
646     jmp    nmiint
647     SET_SIZE(tr_nmiint)

649 #if !defined(__xpv)
650     ENTRY_NP(tr_syserrtrap)
651     /*

```



```

652 * If we got here we should always have a zero error code pushed.
653 * The INT $0x8 instr doesn't seem to push one, though, which we use
654 * as an emergency panic in the other trampolines. So adjust things
655 * here.
656 */
657 cmpq $0, (%rsp)
658 je lf
659 pushq $0
660 1:
661 pushq %r13
662 mov kpti_safe_cr3, %r13
663 mov %r13, %cr3
664 popq %r13
665 jmp syserrtrap
666 SET_SIZE(tr_syserrtrap)
667 #endif

669 ENTRY_NP(tr_mcetrap)
670 pushq %r13
671 mov kpti_safe_cr3, %r13
672 mov %r13, %cr3
673 popq %r13
674 jmp mcetrap
675 SET_SIZE(tr_mcetrap)

677 /*
678 * Interrupts start at 32
679 */
680 #define MKIVCT(n) \
681 ENTRY_NP(tr_ivct/**/n) \
682 push $0; \
683 INTERRUPT_TRAMPOLINE; \
684 push $n - 0x20; \
685 jmp cmnint; \
686 SET_SIZE(tr_ivct/**/n)

688 MKIVCT(32); MKIVCT(33); MKIVCT(34); MKIVCT(35);
689 MKIVCT(36); MKIVCT(37); MKIVCT(38); MKIVCT(39);
690 MKIVCT(40); MKIVCT(41); MKIVCT(42); MKIVCT(43);
691 MKIVCT(44); MKIVCT(45); MKIVCT(46); MKIVCT(47);
692 MKIVCT(48); MKIVCT(49); MKIVCT(50); MKIVCT(51);
693 MKIVCT(52); MKIVCT(53); MKIVCT(54); MKIVCT(55);
694 MKIVCT(56); MKIVCT(57); MKIVCT(58); MKIVCT(59);
695 MKIVCT(60); MKIVCT(61); MKIVCT(62); MKIVCT(63);
696 MKIVCT(64); MKIVCT(65); MKIVCT(66); MKIVCT(67);
697 MKIVCT(68); MKIVCT(69); MKIVCT(70); MKIVCT(71);
698 MKIVCT(72); MKIVCT(73); MKIVCT(74); MKIVCT(75);
699 MKIVCT(76); MKIVCT(77); MKIVCT(78); MKIVCT(79);
700 MKIVCT(80); MKIVCT(81); MKIVCT(82); MKIVCT(83);
701 MKIVCT(84); MKIVCT(85); MKIVCT(86); MKIVCT(87);
702 MKIVCT(88); MKIVCT(89); MKIVCT(90); MKIVCT(91);
703 MKIVCT(92); MKIVCT(93); MKIVCT(94); MKIVCT(95);
704 MKIVCT(96); MKIVCT(97); MKIVCT(98); MKIVCT(99);
705 MKIVCT(100); MKIVCT(101); MKIVCT(102); MKIVCT(103);
706 MKIVCT(104); MKIVCT(105); MKIVCT(106); MKIVCT(107);
707 MKIVCT(108); MKIVCT(109); MKIVCT(110); MKIVCT(111);
708 MKIVCT(112); MKIVCT(113); MKIVCT(114); MKIVCT(115);
709 MKIVCT(116); MKIVCT(117); MKIVCT(118); MKIVCT(119);
710 MKIVCT(120); MKIVCT(121); MKIVCT(122); MKIVCT(123);
711 MKIVCT(124); MKIVCT(125); MKIVCT(126); MKIVCT(127);
712 MKIVCT(128); MKIVCT(129); MKIVCT(130); MKIVCT(131);
713 MKIVCT(132); MKIVCT(133); MKIVCT(134); MKIVCT(135);
714 MKIVCT(136); MKIVCT(137); MKIVCT(138); MKIVCT(139);
715 MKIVCT(140); MKIVCT(141); MKIVCT(142); MKIVCT(143);
716 MKIVCT(144); MKIVCT(145); MKIVCT(146); MKIVCT(147);
717 MKIVCT(148); MKIVCT(149); MKIVCT(150); MKIVCT(151);

```

```

718 MKIVCT(152); MKIVCT(153); MKIVCT(154); MKIVCT(155);
719 MKIVCT(156); MKIVCT(157); MKIVCT(158); MKIVCT(159);
720 MKIVCT(160); MKIVCT(161); MKIVCT(162); MKIVCT(163);
721 MKIVCT(164); MKIVCT(165); MKIVCT(166); MKIVCT(167);
722 MKIVCT(168); MKIVCT(169); MKIVCT(170); MKIVCT(171);
723 MKIVCT(172); MKIVCT(173); MKIVCT(174); MKIVCT(175);
724 MKIVCT(176); MKIVCT(177); MKIVCT(178); MKIVCT(179);
725 MKIVCT(180); MKIVCT(181); MKIVCT(182); MKIVCT(183);
726 MKIVCT(184); MKIVCT(185); MKIVCT(186); MKIVCT(187);
727 MKIVCT(188); MKIVCT(189); MKIVCT(190); MKIVCT(191);
728 MKIVCT(192); MKIVCT(193); MKIVCT(194); MKIVCT(195);
729 MKIVCT(196); MKIVCT(197); MKIVCT(198); MKIVCT(199);
730 MKIVCT(200); MKIVCT(201); MKIVCT(202); MKIVCT(203);
731 MKIVCT(204); MKIVCT(205); MKIVCT(206); MKIVCT(207);
732 MKIVCT(208); MKIVCT(209); MKIVCT(210); MKIVCT(211);
733 MKIVCT(212); MKIVCT(213); MKIVCT(214); MKIVCT(215);
734 MKIVCT(216); MKIVCT(217); MKIVCT(218); MKIVCT(219);
735 MKIVCT(220); MKIVCT(221); MKIVCT(222); MKIVCT(223);
736 MKIVCT(224); MKIVCT(225); MKIVCT(226); MKIVCT(227);
737 MKIVCT(228); MKIVCT(229); MKIVCT(230); MKIVCT(231);
738 MKIVCT(232); MKIVCT(233); MKIVCT(234); MKIVCT(235);
739 MKIVCT(236); MKIVCT(237); MKIVCT(238); MKIVCT(239);
740 MKIVCT(240); MKIVCT(241); MKIVCT(242); MKIVCT(243);
741 MKIVCT(244); MKIVCT(245); MKIVCT(246); MKIVCT(247);
742 MKIVCT(248); MKIVCT(249); MKIVCT(250); MKIVCT(251);
743 MKIVCT(252); MKIVCT(253); MKIVCT(254); MKIVCT(255);

745 /*
746 * We're PCIDE, but we don't have INVPCID. The only way to invalidate a
747 * PCID other than the current one, then, is to load its cr3 then
748 * invlpg. But loading kE_user_cr3 means we can longer access our
749 * caller's text mapping (or indeed, its stack). So this little helper
750 * has to live within our trampoline text region.
751 *
752 * Called as tr_mmu_flush_user_range(addr, len, pgsz, cr3)
753 */
754 ENTRY_NP(tr_mmu_flush_user_range)
755 push %rbx
756 /* When we read cr3, it never has the NOINVL bit set. */
757 mov %cr3, %rax
758 movq $CR3_NOINVL_BIT, %rbx
759 orq %rbx, %rax

761 mov %rcx, %cr3
762 add %rdi, %rsi
763 .align ASM_ENTRY_ALIGN
764 1:
765 invlpg (%rdi)
766 add %rdx, %rdi
767 cmp %rsi, %rdi
768 jb 1b
769 mov %rax, %cr3
770 pop %rbx
771 retq
772 SET_SIZE(tr_mmu_flush_user_range)

774 .align MMU_PAGESIZE
775 .global kpti_trampoline_end
776 kpti_trampoline_end:
777 nop

779 #endif /* __lint */

```

```

*****
52082 Fri Apr 6 17:25:00 2018
new/usr/src/uts/i86pc/ml/locore.s
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */

22 /*
23  * Copyright (c) 1992, 2010, Oracle and/or its affiliates. All rights reserved.
24 */
25 /*
26  * Copyright (c) 2018 Joyent, Inc.
27  * Copyright (c) 2016, Joyent, Inc. All rights reserved.
28 */

29 /*      Copyright (c) 1990, 1991 UNIX System Laboratories, Inc. */
30 /*      Copyright (c) 1984, 1986, 1987, 1988, 1989, 1990 AT&T */
31 /*      All Rights Reserved */

33 /*      Copyright (c) 1987, 1988 Microsoft Corporation */
34 /*      All Rights Reserved */

37 #include <sys/asm_linkage.h>
38 #include <sys/asm_misc.h>
39 #include <sys/regset.h>
40 #include <sys/privregs.h>
41 #include <sys/psw.h>
42 #include <sys/reboot.h>
43 #include <sys/x86_archext.h>
44 #include <sys/machparam.h>

46 #if defined(__lint)

48 #include <sys/types.h>
49 #include <sys/thread.h>
50 #include <sys/system.h>
51 #include <sys/lgrp.h>
52 #include <sys/regset.h>
53 #include <sys/link.h>
54 #include <sys/bootconf.h>
55 #include <sys/bootsvcs.h>

57 #else /* __lint */

```

```

59 #include <sys/segments.h>
60 #include <sys/pcb.h>
61 #include <sys/trap.h>
62 #include <sys/ftrace.h>
63 #include <sys/traptrace.h>
64 #include <sys/clock.h>
65 #include <sys/cmn_err.h>
66 #include <sys/pit.h>
67 #include <sys/panic.h>

69 #if defined(__xpv)
70 #include <sys/hypervisor.h>
71 #endif

73 #include "assym.h"

75 /*
76  * Our assumptions:
77  *   - We are running in protected-paged mode.
78  *   - Interrupts are disabled.
79  *   - The GDT and IDT are the callers; we need our copies.
80  *   - The kernel's text, initialized data and bss are mapped.
81  *
82  * Our actions:
83  *   - Save arguments
84  *   - Initialize our stack pointer to the thread 0 stack (t0stack)
85  *   - and leave room for a phony "struct regs".
86  *   - Our GDT and IDT need to get munged.
87  *   - Since we are using the boot's GDT descriptors, we need
88  *   - to copy them into our GDT before we switch to ours.
89  *   - We start using our GDT by loading correct values in the
90  *   - selector registers (cs=KCS_SEL, ds=es=ss=KDS_SEL, fs=KFS_SEL,
91  *   - gs=KGS_SEL).
92  *   - The default LDT entry for syscall is set.
93  *   - We load the default LDT into the hardware LDT register.
94  *   - We load the default TSS into the hardware task register.
95  *   - Check for cpu type, i.e. 486 vs. P5 vs. P6 etc.
96  *   - mlsetup(%esp) gets called.
97  *   - We change our appearance to look like the real thread 0.
98  *   - (NOTE: making ourselves to be a real thread may be a noop)
99  *   - main() gets called. (NOTE: main() never returns).
100 *
101 * NOW, the real code!
102 */
103 /*
104  * The very first thing in the kernel's text segment must be a jump
105  * to the os/fakebop.c startup code.
106 */
107 .text
108 jmp    _start

110 /*
111  * Globals:
112 */
113 .globl _locore_start
114 .globl mlsetup
115 .globl main
116 .globl panic
117 .globl t0stack
118 .globl t0
119 .globl sysp
120 .globl edata

122 /*
123  * call back into boot - sysp (bootsvcs.h) and bootops (bootconf.h)
124 */

```

```

125     .globl  bootops
126     .globl  bootopsp

128     /*
129     * NOTE: t0stack should be the first thing in the data section so that
130     * if it ever overflows, it will fault on the last kernel text page.
131     */
132     .data
133     .comm   t0stack, DEFAULTSTKSZ, 32
134     .comm   t0, 4094, 32

136 #endif /* __lint */

139 #if defined(__amd64)

141 #if defined(__lint)

143 /* ARGSUSED */
144 void
145 _locore_start(struct boot_syscalls *sysp, ulong_t rsi, struct bootops *bop)
146 {}

148 #else /* __lint */

150     /*
151     * kobj_init() vectors us back to here with (note) a slightly different
152     * set of arguments than _start is given (see lint prototypes above).
153     *
154     * XXX Make this less vile, please.
155     */
156     ENTRY_NP(_locore_start)

158     /*
159     * %rdi = boot services (should die someday)
160     * %rdx = bootops
161     * end
162     */

164     leaq   edata(%rip), %rbp      /* reference edata for ksyms */
165     movq   $0, (%rbp)           /* limit stack back trace */

167     /*
168     * Initialize our stack pointer to the thread 0 stack (t0stack)
169     * and leave room for a "struct regs" for lwp0. Note that the
170     * stack doesn't actually align to a 16-byte boundary until just
171     * before we call mlsetup because we want to use %rsp to point at
172     * our regs structure.
173     */
174     leaq   t0stack(%rip), %rsp
175     addq   $_CONST(DEFAULTSTKSZ - REGSIZE), %rsp
176 #if (REGSIZE & 15) == 0
177     subq   $8, %rsp
178 #endif
179     /*
180     * Save call back for special x86 boot services vector
181     */
182     movq   %rdi, sysp(%rip)

184     movq   %rdx, bootops(%rip)   /* save bootops */
185     movq   $bootops, bootopsp(%rip)

187     /*
188     * Save arguments and flags, if only for debugging ..
189     */
190     movq   %rdi, REGOFF_RDI(%rsp)

```

```

191     movq   %rsi, REGOFF_RSI(%rsp)
192     movq   %rdx, REGOFF_RDX(%rsp)
193     movq   %rcx, REGOFF_RCX(%rsp)
194     movq   %r8, REGOFF_R8(%rsp)
195     movq   %r9, REGOFF_R9(%rsp)
196     pushf
197     popq   %r11
198     movq   %r11, REGOFF_RFL(%rsp)

200 #if !defined(__xpv)
201     /*
202     * Enable write protect and alignment check faults.
203     */
204     movq   %cr0, %rax
205     orq   $_CONST(CR0_WP|CR0_AM), %rax
206     andq   $_BITNOT(CR0_WT|CR0_CE), %rax
207     movq   %rax, %cr0
208 #endif /* __xpv */

210     /*
211     * (We just assert this works by virtue of being here)
212     */
213     bts   $X86FSET_CPUID, x86_featureset(%rip)

215     /*
216     * mlsetup() gets called with a struct regs as argument, while
217     * main takes no args and should never return.
218     */
219     xorl   %ebp, %ebp
220     movq   %rsp, %rdi
221     pushq  %rbp
222     /* (stack pointer now aligned on 16-byte boundary right here) */
223     movq   %rsp, %rbp
224     call   mlsetup
225     call   main
226     /* NOTREACHED */
227     leaq   __return_from_main(%rip), %rdi
228     xorl   %eax, %eax
229     call   panic
230     SET_SIZE(_locore_start)
    _____unchanged_portion_omitted_____

1075 #endif /* __lint */
1076 #endif /* !__amd64 */

1079 /*
1080 * For stack layout, see privregs.h
1081 * When cmntrap gets called, the error code and trap number have been pushed.
1082 * When cmntrap_pushed gets called, the entire struct regs has been pushed.
1083 */

1085 #if defined(__lint)

1087 /* ARGSUSED */
1088 void
1089 cmntrap()
1090 {}

1092 #else /* __lint */

1094     .globl  trap                /* C handler called below */

1096 #if defined(__amd64)

1098     ENTRY_NP2(cmntrap, _cmntrap)

```

```

1100     INTR_PUSH
1102     ALTENTRY(cmntrap_pushed)
1104     movq   %rsp, %rbp
1106     /*
1107     * - if this is a #pf i.e. T_PGFLT, %r15 is live
1108     *   and contains the faulting address i.e. a copy of %cr2
1109     *
1110     * - if this is a #db i.e. T_SGLSTP, %r15 is live
1111     *   and contains the value of %db6
1112     */
1114     TRACE_PTR(%rdi, %rbx, %ebx, %rcx, $TT_TRAP) /* Uses labels 8 and 9 */
1115     TRACE_REGS(%rdi, %rsp, %rbx, %rcx) /* Uses label 9 */
1116     TRACE_STAMP(%rdi) /* Clobbers %eax, %edx, uses 9 */
1118     /*
1119     * We must first check if DTrace has set its NOFAULT bit. This
1120     * regrettably must happen before the trap stack is recorded, because
1121     * this requires a call to getpcstack() and may induce recursion if an
1122     * fbt::getpcstack: enabling is inducing the bad load.
1123     */
1124     movl   %gs:CPU_ID, %eax
1125     shlq  $CPU_CORE_SHIFT, %rax
1126     leaq  cpu_core(%rip), %r8
1127     addq  %r8, %rax
1128     movw  CPUC_DTRACE_FLAGS(%rax), %cx
1129     testw $CPU_DTRACE_NOFAULT, %cx
1130     jnz   .dtrace_induced
1132     TRACE_STACK(%rdi)
1134     movq  %rbp, %rdi
1135     movq  %r15, %rsi
1136     movl  %gs:CPU_ID, %edx
1138     /*
1139     * We know that this isn't a DTrace non-faulting load; we can now safely
1140     * reenale interrupts. (In the case of pagefaults, we enter through an
1141     * interrupt gate.)
1142     */
1143     ENABLE_INTR_FLAGS
1145     call  trap /* trap(rp, addr, cpuid) handles all traps */
1146     jmp   _sys_rtt
1148 .dtrace_induced:
1149     cmpw  $KCS_SEL, REGOFF_CS(%rbp) /* test CS for user-mode trap */
1150     jne   3f /* if from user, panic */
1152     cmpl  $T_PGFLT, REGOFF_TRAPNO(%rbp)
1153     je    1f
1155     cmpl  $T_GPFLT, REGOFF_TRAPNO(%rbp)
1156     je    0f
1158     cmpl  $T_ILLINST, REGOFF_TRAPNO(%rbp)
1159     je    0f
1161     cmpl  $T_ZERODIV, REGOFF_TRAPNO(%rbp)
1162     jne   4f /* if not PF/GP/UD/DE, panic */
1164     orw  $CPU_DTRACE_DIVZERO, %cx

```

```

1165     movw  %cx, CPUC_DTRACE_FLAGS(%rax)
1166     jmp   2f
1168     /*
1169     * If we've taken a GPF, we don't (unfortunately) have the address that
1170     * induced the fault. So instead of setting the fault to BADADDR,
1171     * we'll set the fault to ILL0P.
1172     */
1173 0:
1174     orw  $CPU_DTRACE_ILL0P, %cx
1175     movw %cx, CPUC_DTRACE_FLAGS(%rax)
1176     jmp  2f
1177 1:
1178     orw  $CPU_DTRACE_BADADDR, %cx
1179     movw %cx, CPUC_DTRACE_FLAGS(%rax) /* set fault to bad addr */
1180     movq %r15, CPUC_DTRACE_ILLOVAL(%rax) /* fault addr is illegal value */
1181     /*
1182     2:
1183     movq  REGOFF_RIP(%rbp), %rdi
1184     movq  %rdi, %r12
1185     call  dtrace_instr_size
1186     addq  %rax, %r12
1187     movq  %r12, REGOFF_RIP(%rbp)
1188     INTR_POP
1189     jmp   tr_iret_auto
1189     IRET
1190     /*NOTREACHED*/
1191 3:
1192     leaq  dtrace_badflags(%rip), %rdi
1193     xorl  %eax, %eax
1194     call  panic
1195 4:
1196     leaq  dtrace_badtrap(%rip), %rdi
1197     xorl  %eax, %eax
1198     call  panic
1199     SET_SIZE(cmntrap)
1200     unchanged portion omitted
1201
1202 #endif /* __i386 */
1203
1204 #endif /* __lint */
1205
1206 /*
1207 * Return from _sys_trap routine.
1208 */
1209
1210 #if defined(__lint)
1211 void
1212 lwp_rtt_initial(void)
1213 {}
1214
1215 void
1216 lwp_rtt(void)
1217 {}
1218
1219 void
1220 _sys_rtt(void)
1221 {}
1222 #else /* __lint */
1223 #if defined(__amd64)
1224     ENTRY_NP(lwp_rtt_initial)
1225     movq  %gs:CPU_THREAD, %r15

```

```

1512      movq   T_STACK(%r15), %rsp      /* switch to the thread stack */
1513      movq   %rsp, %rbp
1514      call   __dtrace_probe__proc_start
1515      jmp    _lwp_rtt

1517      ENTRY_NP(lwp_rtt)

1519      /*
1520      * r14 lwp
1521      * rdx lwp->lwp_procp
1522      * r15 curthread
1523      */

1525      movq   %gs:CPU_THREAD, %r15
1526      movq   T_STACK(%r15), %rsp      /* switch to the thread stack */
1527      movq   %rsp, %rbp
1528  _lwp_rtt:
1529      call   __dtrace_probe__proc_lwp__start
1530      movq   %gs:CPU_LWP, %r14
1531      movq   LWP_PROCP(%r14), %rdx

1533      /*
1534      * XX64 Is the stack misaligned correctly at this point?
1535      *      If not, we need to do a push before calling anything ..
1536      */

1538  #if defined(DEBUG)
1539      /*
1540      * If we were to run lwp_savectx at this point -without-
1541      * pcb_update being set to 1, we'd end up sampling the hardware
1542      * state left by the previous running lwp, rather than setting
1543      * the values requested by the lwp creator.  Bad.
1544      */
1545      testb  $0x1, PCB_RUPDATE(%r14)
1546      jne   lf
1547      leaq  _no_pending_updates(%rip), %rdi
1548      movl  $_LINE__, %esi
1549      movq  %r14, %rdx
1550      xorl  %eax, %eax
1551      call  panic
1552  _no_pending_updates:
1553      .string "locore.s:%d lwp_rtt(lwp %p) but pcb_rupdate != 1"
1554  1:
1555  #endif

1557      /*
1558      * If agent lwp, clear %fs and %gs
1559      */
1560      cmpq  %r15, P_AGENTTP(%rdx)
1561      jne  lf
1562      xorl  %ecx, %ecx
1563      movq  %rcx, REGOFF_FS(%rsp)
1564      movq  %rcx, REGOFF_GS(%rsp)
1565      movw  %cx, LWP_PCB_FS(%r14)
1566      movw  %cx, LWP_PCB_GS(%r14)
1567  1:
1568      call  dtrace_systrace_rtt
1569      movq  REGOFF_RDX(%rsp), %rsi
1570      movq  REGOFF_RAX(%rsp), %rdi
1571      call  post_syscall      /* post_syscall(rvall, rval2) */

1573      /*
1574      * set up to take fault on first use of fp
1575      */
1576      STTS(%rdi)

```

```

1578      /*
1579      * XXX - may want a fast path that avoids sys_rtt_common in the
1580      * most common case.
1581      */
1582      ALTENTRY(_sys_rtt)
1583      CLI(%rax)      /* disable interrupts */
1584      ALTENTRY(_sys_rtt_ints_disabled)
1585      movq  %rsp, %rdi      /* pass rp to sys_rtt_common */
1586      call  sys_rtt_common  /* do common sys_rtt tasks */
1587      testq %rax, %rax      /* returning to userland? */
1588      jz    sr_sup

1590      /*
1591      * Return to user
1592      */
1593      ASSERT_UPCALL_MASK_IS_SET
1594      cmpw  $UCS_SEL, REGOFF_CS(%rsp) /* test for native (64-bit) lwp? */
1595      je    sys_rtt_syscall

1597      /*
1598      * Return to 32-bit userland
1599      */
1600      ALTENTRY(sys_rtt_syscall32)
1601      USER32_POP
1602      jmp   tr_iret_user
1602      IRET
1603      /*NOTREACHED*/

1605      ALTENTRY(sys_rtt_syscall)
1606      /*
1607      * Return to 64-bit userland
1608      */
1609      USER_POP
1610      ALTENTRY(nopop_sys_rtt_syscall)
1611      jmp   tr_iret_user
1611      IRET
1612      /*NOTREACHED*/
1613      SET_SIZE(nopop_sys_rtt_syscall)

1615      /*
1616      * Return to supervisor
1617      * NOTE: to make the check in trap() that tests if we are executing
1618      * segment register fixup/restore code work properly, sr_sup MUST be
1619      * after _sys_rtt .
1620      */
1621      ALTENTRY(sr_sup)
1622      /*
1623      * Restore regs before doing iretq to kernel mode
1624      */
1625      INTR_POP
1626      jmp   tr_iret_kernel
1626      IRET
1627      .globl _sys_rtt_end
1628  _sys_rtt_end:
1629      /*NOTREACHED*/
1630      SET_SIZE(sr_sup)
1630      _____unchanged_portion_omitted_____

```

```

*****
15320 Fri Apr 6 17:25:00 2018
new/usr/src/uts/i86pc/ml/mpcore.s
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22 * Copyright (c) 1992, 2010, Oracle and/or its affiliates. All rights reserved.
23 */
24 /*
25 * Copyright (c) 2010, Intel Corporation.
26 * All rights reserved.
27 *
28 * Copyright 2018 Joyent, Inc.
29 */
30
31 #include <sys/asm_linkage.h>
32 #include <sys/asm_misc.h>
33 #include <sys/regset.h>
34 #include <sys/privregs.h>
35 #include <sys/x86_archext.h>
36
37 #if !defined(__lint)
38 #include <sys/segments.h>
39 #include "assym.h"
40 #endif
41
42 /*
43  * Our assumptions:
44  *   - We are running in real mode.
45  *   - Interrupts are disabled.
46  *   - Selectors are equal (cs == ds == ss) for all real mode code
47  *   - The GDT, IDT, ktss and page directory has been built for us
48  *
49  * Our actions:
50  * Start CPU:
51  *   - We start using our GDT by loading correct values in the
52  *   selector registers (cs=KCS_SEL, ds=es=ss=KDS_SEL, fs=KFS_SEL,
53  *   gs=KGS_SEL).
54  *   - We change over to using our IDT.
55  *   - We load the default LDT into the hardware LDT register.
56  *   - We load the default TSS into the hardware task register.
57  *   - call mp_startup(void) indirectly through the T_PC
58  * Stop CPU:
59  *   - Put CPU into halted state with interrupts disabled

```

```

60 *
61 */
62
63 #if defined(__lint)
64
65 void
66 real_mode_start_cpu(void)
67 {}
68
69 void
70 real_mode_stop_cpu_stage1(void)
71 {}
72
73 void
74 real_mode_stop_cpu_stage2(void)
75 {}
76
77 #else /* __lint */
78
79 #if defined(__amd64)
80
81     ENTRY_NP(real_mode_start_cpu)
82
83     /*
84      * NOTE: The GNU assembler automatically does the right thing to
85      * generate data size operand prefixes based on the code size
86      * generation mode (e.g. .code16, .code32, .code64) and as such
87      * prefixes need not be used on instructions EXCEPT in the case
88      * of address prefixes for code for which the reference is not
89      * automatically of the default operand size.
90      */
91     .code16
92     cli
93     movw     %cs, %ax
94     movw     %ax, %ds /* load cs into ds */
95     movw     %ax, %ss /* and into ss */
96
97     /*
98      * Helps in debugging by giving us the fault address.
99      *
100     * Remember to patch a hlt (0xf4) at cmntrap to get a good stack.
101     */
102     movl     $0xffc, %esp
103     movl     %cr0, %eax
104
105     /*
106     * Enable protected-mode, write protect, and alignment mask
107     */
108     orl     $(CR0_PE|CR0_WP|CR0_AM), %eax
109     movl     %eax, %cr0
110
111     /*
112     * Do a jmp immediately after writing to cr0 when enabling protected
113     * mode to clear the real mode prefetch queue (per Intel's docs)
114     */
115     jmp     pestart
116
117 pestart:
118     /*
119     * 16-bit protected mode is now active, so prepare to turn on long
120     * mode.
121     *
122     * Note that we currently assume that if we're attempting to run a
123     * kernel compiled with (__amd64) #defined, the target CPU has long
124     * mode support.
125     */

```

```

127 #if 0
128 /*
129  * If there's a chance this might not be true, the following test should
130  * be done, with the no_long_mode branch then doing something
131  * appropriate:
132  */
134 movl    $0x80000000, %eax    /* get largest extended CPUID */
135 cpuid
136 cmpl    $0x80000000, %eax    /* check if > 0x80000000 */
137 jbe     no_long_mode        /* nope, no long mode */
138 movl    $0x80000001, %eax
139 cpuid
140 btl     $29, %edx            /* check for long mode */
141 jnc     no_long_mode        /* long mode not supported */
142 #endif
144 /*
145  * Add any initial cr4 bits
146  */
147 movl    %cr4, %eax
148 addr32 orl    CR4OFF, %eax
150 /*
151  * Enable PAE mode (CR4.PAE)
152  */
153 orl     $CR4_PAE, %eax
154 movl    %eax, %cr4
156 /*
157  * Point cr3 to the 64-bit long mode page tables.
158  *
159  * Note that these MUST exist in 32-bit space, as we don't have
160  * a way to load %cr3 with a 64-bit base address for the page tables
161  * until the CPU is actually executing in 64-bit long mode.
162  */
163 addr32 movl    CR3OFF, %eax
164 movl    %eax, %cr3
166 /*
167  * Set long mode enable in EFER (EFER.LME = 1)
168  */
169 movl    $MSR_AMD_EFER, %ecx
170 rdmsr
171 orl     $AMD_EFER_LME, %eax
172 wrmsr
174 /*
175  * Finally, turn on paging (CR0.PG = 1) to activate long mode.
176  */
177 movl    %cr0, %eax
178 orl     $CR0_PG, %eax
179 movl    %eax, %cr0
181 /*
182  * The instruction after enabling paging in CR0 MUST be a branch.
183  */
184 jmp     long_mode_active
186 long_mode_active:
187 /*
188  * Long mode is now active but since we're still running with the
189  * original 16-bit CS we're actually in 16-bit compatability mode.
190  *
191  * We have to load an intermediate GDT and IDT here that we know are

```

```

192  * in 32-bit space before we can use the kernel's GDT and IDT, which
193  * may be in the 64-bit address space, and since we're in compatability
194  * mode, we only have access to 16 and 32-bit instructions at the
195  * moment.
196  */
197 addr32 lgdtl    TEMPGDTOFF    /* load temporary GDT */
198 addr32 lidt    TEMPIDTOFF    /* load temporary IDT */
200 /*
201  * Do a far transfer to 64-bit mode. Set the CS selector to a 64-bit
202  * long mode selector (CS.L=1) in the temporary 32-bit GDT and jump
203  * to the real mode platter address of long_mode_64 as until the 64-bit
204  * CS is in place we don't have access to 64-bit instructions and thus
205  * can't reference a 64-bit %rip.
206  */
207 pushl    $TEMP_CS64_SEL
208 addr32 pushl    LM64OFF
209 lretl
211 .globl    long_mode_64
212 long_mode_64:
213 .code64
214 /*
215  * We are now running in long mode with a 64-bit CS (EFER.LMA=1,
216  * CS.L=1) so we now have access to 64-bit instructions.
217  *
218  * First, set the 64-bit GDT base.
219  */
220 .globl    rm_platter_pa
221 movl    rm_platter_pa, %eax
222 lgdtq    GDTROFF(%rax)    /* load 64-bit GDT */
224 /*
225  * Save the CPU number in %r11; get the value here since it's saved in
226  * the real mode platter.
227  */
228 movl    CPUNOFF(%rax), %r11
230 /*
231  * Add rm_platter_pa to %rsp to point it to the same location as seen
232  * from 64-bit mode.
233  */
234 addq    %rax, %rsp
236 /*
237  * Now do an lretq to load CS with the appropriate selector for the
238  * kernel's 64-bit GDT and to start executing 64-bit setup code at the
239  * virtual address where boot originally loaded this code rather than
240  * the copy in the real mode platter's rm_code array as we've been
241  * doing so far.
242  */
243 pushq    $KCS_SEL
244 pushq    $kernel_cs_code
245 lretq
246 .globl    real_mode_start_cpu_end
247 real_mode_start_cpu_end:
248 nop
250 kernel_cs_code:
251 /*
252  * Complete the balance of the setup we need to before executing
253  * 64-bit kernel code (namely init rsp, TSS, LGDT, FS and GS).
254  */
255 .globl    rm_platter_va
256 movq    rm_platter_va, %rax
257 lidt    IDTROFF(%rax)

```

```

259     movw    $KDS_SEL, %ax
260     movw    %ax, %ds
261     movw    %ax, %es
262     movw    %ax, %ss

264     movw    $KTSS_SEL, %ax      /* setup kernel TSS */
265     ltr     %ax

267     xorw    %ax, %ax           /* clear LDTR */
268     lldt   %ax

270     /*
271     * Set GS to the address of the per-cpu structure as contained in
272     * cpu[cpu_number].
273     *
274     * Unfortunately there's no way to set the 64-bit gsbase with a mov,
275     * so we have to stuff the low 32 bits in %eax and the high 32 bits in
276     * %edx, then call wrmsr.
277     */
278     leaq   cpu(%rip), %rdi
279     movl   (%rdi, %r11, 8), %eax
280     movl   4(%rdi, %r11, 8), %edx
281     movl   $MSR_AMD_GSBASE, %ecx
282     wrmsr

284     /*
285     * Init FS and KernelGSBase.
286     *
287     * Based on code in mlsetup(), set them both to 8G (which shouldn't be
288     * valid until some 64-bit processes run); this will then cause an
289     * exception in any code that tries to index off them before they are
290     * properly setup.
291     */
292     xorl   %eax, %eax          /* low 32 bits = 0 */
293     movl   $2, %edx           /* high 32 bits = 2 */
294     movl   $MSR_AMD_FSBASE, %ecx
295     wrmsr

297     movl   $MSR_AMD_KGSBASE, %ecx
298     wrmsr

300     /*
301     * Init %rsp to the exception stack set in tss_istl and create a legal
302     * AMD64 ABI stack frame
303     */
304     movq   %gs:CPU_TSS, %rax
305     movq   TSS_IST1(%rax), %rsp
306     pushq  $0                /* null return address */
307     pushq  $0                /* null frame pointer terminates stack trace */
308     movq   %rsp, %rbp        /* stack aligned on 16-byte boundary */

310     movq   %cr0, %rax
311     andq   $~(CR0_TS|CR0_EM), %rax /* clear emulate math chip bit */
312     orq   $(CR0_MP|CR0_NE), %rax
313     movq   %rax, %cr0       /* set machine status word */

315     /*
316     * Before going any further, enable usage of page table NX bit if
317     * that's how our page tables are set up.
318     */
319     bt     $X86FSET_NX, x86_featureset(%rip)
320     jnc   lf
321     movl   $MSR_AMD_EFER, %ecx
322     rdmsr
323     orl   $AMD_EFER_NXE, %eax

```

```

324     wrmsr
325 1:

327     /*
328     * Complete the rest of the setup and call mp_startup().
329     */
330     movq   %gs:CPU_THREAD, %rax /* get thread ptr */
331     call   *T_PC(%rax)         /* call mp_startup_boot */
329     call   *T_PC(%rax)         /* call mp_startup */
332     /* not reached */
333     int    $20                /* whoops, returned somehow! */

335     SET_SIZE(real_mode_start_cpu)

337 #elif defined(__i386)

339     ENTRY_NP(real_mode_start_cpu)

341 #if !defined(__GNUC_AS__)

343     cli
344     D16 movw    %cs, %eax
345     movw    %eax, %ds          /* load cs into ds */
346     movw    %eax, %ss          /* and into ss */

348     /*
349     * Helps in debugging by giving us the fault address.
350     *
351     * Remember to patch a hlt (0xf4) at cmntrap to get a good stack.
352     */
353     D16 movl   $0xffc, %esp

355     D16 A16 lgdt   %cs:GDTR0FF
356     D16 A16 lidt   %cs:IDTR0FF
357     D16 A16 movl   %cs:CR4OFF, %eax /* set up CR4, if desired */
358     D16 andl   %eax, %eax
359     D16 A16 je    no_cr4

361     D16 movl   %eax, %ecx
362     D16 movl   %cr4, %eax
363     D16 orl   %ecx, %eax
364     D16 movl   %eax, %cr4

365 no_cr4:
366     D16 A16 movl   %cs:CR3OFF, %eax
367     A16 movl   %eax, %cr3
368     movl   %cr0, %eax

370     /*
371     * Enable protected-mode, paging, write protect, and alignment mask
372     */
373     D16 orl   $[CR0_PG|CR0_PE|CR0_WP|CR0_AM], %eax
374     movl   %eax, %cr0
375     jmp    pestart

377 pestart:
378     D16 pushl   $KCS_SEL
379     D16 pushl   $kernel_cs_code
380     D16 lret
381     .globl real_mode_start_cpu_end
382 real_mode_start_cpu_end:
383     nop

385     .globl kernel_cs_code
386 kernel_cs_code:
387     /*
388     * At this point we are with kernel's cs and proper eip.

```



```

389      *
390      * We will be executing not from the copy in real mode platter,
391      * but from the original code where boot loaded us.
392      *
393      * By this time GDT and IDT are loaded as is cr3.
394      */
395      movw    $KFS_SEL,%eax
396      movw    %eax,%fs
397      movw    $KGS_SEL,%eax
398      movw    %eax,%gs
399      movw    $KDS_SEL,%eax
400      movw    %eax,%ds
401      movw    %eax,%es
402      movl    %gs:CPU_TSS,%esi
403      movw    %eax,%ss
404      movl    TSS_ESP0(%esi),%esp
405      movw    $KTSS_SEL,%ax
406      ltr    %ax
407      xorw    %ax,%ax          /* clear LDTR */
408      lldt   %ax
409      movl    %cr0,%edx
410      andl   $-1![(CR0_TS|CR0_EM)],%edx /* clear emulate math chip bit */
411      orl    $(CR0_MP|CR0_NE),%edx
412      movl   %edx,%cr0        /* set machine status word */

414     /*
415     * Before going any further, enable usage of page table NX bit if
416     * that's how our page tables are set up.
417     */
418     bt     $X86FSET_NX, x86_featureset
419     jnc   lf
420     movl  %cr4,%ecx
421     andl  $CR4_PAE,%ecx
422     jz   lf
423     movl  $MSR_AMD_EFER,%ecx
424     rdmsr
425     orl   $AMD_EFER_NXE,%eax
426     wrmsr
427 1:
428     movl  %gs:CPU_THREAD,%eax /* get thread ptr */
429     call  *T_PC(%eax)        /* call mp_startup */
430     /* not reached */
431     int   $20                /* whoops, returned somehow! */

433 #else

435     cli
436     mov   %cs,%ax
437     mov   %eax,%ds          /* load cs into ds */
438     mov   %eax,%ss          /* and into ss */

440     /*
441     * Helps in debugging by giving us the fault address.
442     *
443     * Remember to patch a hlt (0xf4) at cmntrap to get a good stack.
444     */
445     D16 mov    $0xffc,%esp

447     D16 A16 lgdtl  %cs:GDTR0FF
448     D16 A16 lidt1  %cs:IDTR0FF
449     D16 A16 mov    %cs:CR4OFF,%eax /* set up CR4, if desired */
450     D16 and    %eax,%eax
451     D16 A16 je    no_cr4

453     D16 mov    %eax,%ecx
454     D16 mov    %cr4,%eax

```

```

455     D16 or    %ecx,%eax
456     D16 mov    %eax,%cr4
457 no_cr4:
458     D16 A16 mov    %cs:CR3OFF,%eax
459     A16 mov    %eax,%cr3
460     mov    %cr0,%eax

462     /*
463     * Enable protected-mode, paging, write protect, and alignment mask
464     */
465     D16 or    $(CR0_PG|CR0_PE|CR0_WP|CR0_AM),%eax
466     mov    %eax,%cr0
467     jmp    pestart

469 pestart:
470     D16 pushl   $KCS_SEL
471     D16 pushl   $kernel_cs_code
472     D16 lret
473     .globl real_mode_start_cpu_end
474 real_mode_start_cpu_end:
475     nop
476     .globl kernel_cs_code
477 kernel_cs_code:
478     /*
479     * At this point we are with kernel's cs and proper eip.
480     *
481     * We will be executing not from the copy in real mode platter,
482     * but from the original code where boot loaded us.
483     *
484     * By this time GDT and IDT are loaded as is cr3.
485     */
486     mov    $KFS_SEL,%ax
487     mov    %eax,%fs
488     mov    $KGS_SEL,%ax
489     mov    %eax,%gs
490     mov    $KDS_SEL,%ax
491     mov    %eax,%ds
492     mov    %eax,%es
493     mov    %gs:CPU_TSS,%esi
494     mov    %eax,%ss
495     mov    TSS_ESP0(%esi),%esp
496     mov    $(KTSS_SEL),%ax
497     ltr    %ax
498     xorw   %ax,%ax          /* clear LDTR */
499     lldt   %ax
500     mov    %cr0,%edx
501     and   $~(CR0_TS|CR0_EM),%edx /* clear emulate math chip bit */
502     or    $(CR0_MP|CR0_NE),%edx
503     mov    %edx,%cr0        /* set machine status word */

505     /*
506     * Before going any farther, enable usage of page table NX bit if
507     * that's how our page tables are set up. (PCIDE is enabled later on).
508     * that's how our page tables are set up.
509     */
510     bt     $X86FSET_NX, x86_featureset
511     jnc   lf
512     movl  %cr4,%ecx
513     andl  $CR4_PAE,%ecx
514     jz   lf
515     movl  $MSR_AMD_EFER,%ecx
516     rdmsr
517     orl   $AMD_EFER_NXE,%eax
518 wrmsr
519 1:
520     mov    %gs:CPU_THREAD,%eax /* get thread ptr */

```

new/usr/src/uts/i86pc/ml/mpcore.s

9

```
520     call    *T_PC(%eax)          /* call mp_startup */
521     /* not reached */
522     int     $20                  /* whoops, returned somehow! */
523 #endif

525     SET_SIZE(real_mode_start_cpu)
unchanged_portion_omitted
```

new/usr/src/uts/i86pc/ml/offsets.in

1

8676 Fri Apr 6 17:25:00 2018
new/usr/src/uts/i86pc/ml/offsets.in
9441 kmdb should stash %cr3 in kdiregs
Reviewed by: John Levon <john.levon@joyent.com>
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>

```
1 \
2 \ Copyright (c) 2004, 2010, Oracle and/or its affiliates. All rights reserved.
3 \ Copyright 2012 Garrett D'Amore <garrett@damore.org>. All rights reserved.
4 \ Copyright 2018 Joyent, Inc.
4 \ Copyright 2016 Joyent, Inc.
5 \
6 \ CDDL HEADER START
7 \
8 \ The contents of this file are subject to the terms of the
9 \ Common Development and Distribution License (the "License").
10 \ You may not use this file except in compliance with the License.
11 \
12 \ You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
13 \ or http://www.opensolaris.org/os/licensing.
14 \ See the License for the specific language governing permissions
15 \ and limitations under the License.
16 \
17 \ When distributing Covered Code, include this CDDL HEADER in each
18 \ file and include the License file at usr/src/OPENSOLARIS.LICENSE.
19 \ If applicable, add the following below this CDDL HEADER, with the
20 \ fields enclosed by brackets "[]" replaced with your own identifying
21 \ information: Portions Copyright [yyyy] [name of copyright owner]
22 \
23 \ CDDL HEADER END
24 \
```

```
27 \
28 \ offsets.in: input file to produce assym.h using the ctfstabs program
29 \
```

```
31 #ifndef _GENASSYM
32 #define _GENASSYM
33 #endif
```

```
35 #define SIZES 1
```

```
37 #include <sys/types.h>
38 #include <sys/bootsvcs.h>
39 #include <sys/system.h>
40 #include <sys/sysinfo.h>
41 #include <sys/user.h>
42 #include <sys/thread.h>
43 #include <sys/proc.h>
44 #include <sys/cpuvar.h>
45 #include <sys/tss.h>
46 #include <sys/privregs.h>
47 #include <sys/segments.h>
48 #include <sys/devops.h>
49 #include <sys/ddi_impldefs.h>
50 #include <vm/as.h>
51 #include <sys/avintr.h>
52 #include <sys/pic.h>
53 #include <sys/rm_platter.h>
54 #include <sys/stream.h>
55 #include <sys/strsubr.h>
56 #include <sys/sunddi.h>
```

new/usr/src/uts/i86pc/ml/offsets.in

2

```
57 #include <sys/modctl.h>
58 #include <sys/traptrace.h>
59 #include <sys/onttrap.h>
60 #include <sys/lgrp.h>
61 #include <sys/dtrace.h>
62 #include <sys/brand.h>
63 #include <sys/fastboot.h>
64 #include <sys/cpr_wakecode.h>
65 #include <sys/comm_page.h>
```

```
67 proc          PROC_SIZE
68     p_link
69     p_next
70     p_child
71     p_sibling
72     p_sig
73     p_flag
74     p_tlist
75     p_as
76     p_lockp
77     p_user
78     p_ldt
79     p_ldt_desc
80     p_model
81     p_pctx
82     p_agenttp
83     p_zone
84     p_brand
85     p_brand_data
```

```
87 _kthread      THREAD_SIZE
88     t_pcb          T_LABEL
89     t_lock
90     t_lockstat
91     t_lockp
92     t_lock_flush
93     t_kpri_req
94     t_oldspl
95     t_pri
96     t_pil
97     t_lwp
98     t_procp
99     t_link
100    t_state
101    t_mstate
102    t_preempt_lk
103    t_stk          T_STACK
104    t_swap
105    t_lwpchan.lc_wchan T_WCHAN
106    t_flag        T_FLAGS
107    t_ctx
108    t_lofault
109    t_onfault
110    t_onttrap
111    t_cpu
112    t_lpl
113    t_bound_cpu
114    t_intr
115    t_forw
116    t_back
117    t_sig
118    t_tid
119    t_pre_sys
120    t_preempt
121    t_proc_flag
122    t_startpc
```

```

123     t_sysnum
124     t_intr_start
125     _tu._ts._t_astflag      T_ASTFLAG
126     _tu._ts._t_post_sys    T_POST_SYS
127     _tu._t_post_sys_ast    T_POST_SYS_AST
128     t_copyops
129 #ifdef  __amd64
130     t_useracc
131 #endif

133 ctxop
134     save_op                CTXOP_SAVE

136 as
137     a_hat

139 user    USIZEBYTES
140     u_comm
141     u_signal

143 _label_t
144     val        LABEL_VAL

146 #define LABEL_PC LABEL_VAL
147 #define LABEL_SP _CONST(LABEL_VAL + LABEL_VAL_INCR)
148 #define T_PC _CONST(T_LABEL + LABEL_PC)
149 #define T_SP _CONST(T_LABEL + LABEL_SP)

151 _klwp
152     lwp_thread
153     lwp_procp
154     lwp_brand
155     lwp_eosys
156     lwp_regs
157     lwp_arg
158     lwp_ap
159     lwp_cursig
160     lwp_state
161     lwp_mstate.ms_acct      LWP_MS_ACCT
162     lwp_mstate.ms_prev     LWP_MS_PREV
163     lwp_mstate.ms_start    LWP_MS_START
164     lwp_mstate.ms_state_start LWP_MS_STATE_START
165     lwp_pcb
166     lwp_ru.sysc            LWP_RU_SYSC

168 #define LWP_ACCT_USER _CONST(LWP_MS_ACCT + _MUL(LMS_USER, LWP_MS_ACCT_
169 #define LWP_ACCT_SYSTEM _CONST(LWP_MS_ACCT + _MUL(LMS_SYSTEM, LWP_MS_ACC

171 fpu_ctx
172     fpu_regs                FPU_CTX_FPU_REGS
173     fpu_flags                FPU_CTX_FPU_FLAGS
174     fpu_xsave_mask          FPU_CTX_FPU_XSAVE_MASK

176 fxsave_state    FXSAVE_STATE_SIZE
177     fx_fsw        FXSAVE_STATE_FSW
178     fx_mxcsr_mask FXSAVE_STATE_MXCSR_MASK

181 autovec        AUTOVECSIZE
182     av_vector
183     av_intarg1
184     av_intarg2
185     av_ticksp
186     av_link
187     av_prilevel
188     av_dip

```

```

190 av_head
191     avh_link
192     avh_hi_pri
193     avh_lo_pri

195 cpu
196     cpu_id
197     cpu_flags
198     cpu_self
199     cpu_thread
200     cpu_thread_lock
201     cpu_kprunrun
202     cpu_lwp
203     cpu_fpowner
204     cpu_idle_thread
205     cpu_intr_thread
206     cpu_intr_actv
207     cpu_base_spl
208     cpu_intr_stack
209     cpu_stats.sys.cpumigrate      CPU_STATS_SYS_CPUMIGRATE
210     cpu_stats.sys.intr            CPU_STATS_SYS_INTR
211     cpu_stats.sys.intrblk         CPU_STATS_SYS_INTRBLK
212     cpu_stats.sys.syscall         CPU_STATS_SYS_SYSCALL
213     cpu_profile_pc
214     cpu_profile_upc
215     cpu_profile_pil
216     cpu_ftrace.ftd_state          CPU_FTRACE_STATE
217     cpu_mstate
218     cpu_intracct

220 #define CPU_INTR_ACTV_REF _CONST(CPU_INTR_ACTV + 2)

222 cpu
223     cpu_m.pil_high_start          CPU_PIL_HIGH_START
224     cpu_m.intrstat                CPU_INTRSTAT
225     cpu_m.mcpu_current_hat        CPU_CURRENT_HAT
226     cpu_m.mcpu_gdt                CPU_GDT
227     cpu_m.mcpu_idt                CPU_IDT
228     cpu_m.mcpu_tss                CPU_TSS
229     cpu_m.mcpu_softinfo           CPU_SOFTINFO
230     cpu_m.mcpu_pri                CPU_PRI
231 #if defined(__xpv)
232     cpu_m.mcpu_vcpu_info           CPU_VCPU_INFO
233 #endif

235 cpu
236     cpu_m.mcpu_kpti.kf_kernel_cr3 CPU_KPTI_KCR3
237     cpu_m.mcpu_kpti.kf_user_cr3   CPU_KPTI_UCR3
238     cpu_m.mcpu_kpti.kf_tr_rsp     CPU_KPTI_TR_RSP
239     cpu_m.mcpu_kpti.kf_tr_cr3     CPU_KPTI_TR_CR3
240     cpu_m.mcpu_kpti.kf_r13        CPU_KPTI_R13
241     cpu_m.mcpu_kpti.kf_r14        CPU_KPTI_R14
242     cpu_m.mcpu_kpti.kf_tr_ret_rsp CPU_KPTI_RET_RSP

244     cpu_m.mcpu_kpti.kf_ss          CPU_KPTI_SS
245     cpu_m.mcpu_kpti.kf_rsp         CPU_KPTI_RSP
246     cpu_m.mcpu_kpti.kf_rflags     CPU_KPTI_RFLAGS
247     cpu_m.mcpu_kpti.kf_cs         CPU_KPTI_CS
248     cpu_m.mcpu_kpti.kf_rip        CPU_KPTI_RIP
249     cpu_m.mcpu_kpti.kf_err        CPU_KPTI_ERR

251     cpu_m.mcpu_pad2                CPU_KPTI_START
252     cpu_m.mcpu_pad3                CPU_KPTI_END

254     cpu_m.mcpu_kpti_dbg            CPU_KPTI_DBG

```

```

256 kpti_frame
257     kf_r14      KPTI_R14
258     kf_r13      KPTI_R13
259     kf_err      KPTI_ERR
260     kf_rip      KPTI_RIP
261     kf_cs       KPTI_CS
262     kf_rflags   KPTI_RFLAGS
263     kf_rsp      KPTI_RSP
264     kf_ss       KPTI_SS

266     kf_tr_rsp   KPTI_TOP

268     kf_kernel_cr3 KPTI_KCR3
269     kf_user_cr3  KPTI_UCR3
270     kf_tr_ret_rsp KPTI_RET_RSP
271     kf_tr_cr3   KPTI_TR_CR3

273     kf_tr_flag  KPTI_FLAG

275 standard_pic
276     c_curmask
277     c_iplmask

279 ddi_dma_impl
280     dmai_rflags
281     dmai_rdip

283 dev_info
284     devi_ops          DEVI_DEV_OPS
285     devi_bus_ctl
286     devi_bus_dma_ctl
287     devi_bus_dma_allochdl
288     devi_bus_dma_freehdl
289     devi_bus_dma_bindhdl
290     devi_bus_dma_unbindhdl
291     devi_bus_dma_flush
292     devi_bus_dma_win

294 dev_ops
295     devo_bus_ops      DEVI_BUS_OPS

297 bus_ops
298     bus_ctl          OPS_CTL
299     bus_dma_map      OPS_MAP
300     bus_dma_ctl     OPS_MCTL
301     bus_dma_allochdl OPS_ALLOCHDL
302     bus_dma_freehdl OPS_FREEHDL
303     bus_dma_bindhdl OPS_BINDHDL
304     bus_dma_unbindhdl OPS_UNBINDHDL
305     bus_dma_flush   OPS_FLUSH
306     bus_dma_win     OPS_WIN

308 sysent  SYSENT_SIZE  SYSENT_SIZE_SHIFT
309     sy_callc
310     sy_flags
311     sy_narg

313 stdata
314     sd_lock

316 queue
317     q_flag
318     q_next
319     q_stream
320     q_syncq

```

```

321     q_qinfo

323 qinit
324     qi_putp

326 syncq
327     sq_flags
328     sq_count
329     sq_lock
330     sq_wait

332 rm_platter
333     rm_idt_lim      IDTROFF
334     rm_gdt_lim      GDTROFF
335     rm_pdbcr       CR3OFF
336     rm_cpu         CPUNOFF
337     rm_cr4         CR4OFF
338     rm_cpu_halt_code CPUHALTCODEOFF
339     rm_cpu_halted  CPUHALTEDOFF

341 ddi_acc_impl
342     ahi_acc_attr   ACC_ATTR
343     ahi_get8      ACC_GETB
344     ahi_get16     ACC_GETW
345     ahi_get32     ACC_GETL
346     ahi_get64     ACC_GETLL
347     ahi_put8      ACC_PUTB
348     ahi_put16     ACC_PUTW
349     ahi_put32     ACC_PUTL
350     ahi_put64     ACC_PUTLL
351     ahi_rep_get8  ACC_REP_GETB
352     ahi_rep_get16 ACC_REP_GETW
353     ahi_rep_get32 ACC_REP_GETL
354     ahi_rep_get64 ACC_REP_GETLL
355     ahi_rep_put8  ACC_REP_PUTB
356     ahi_rep_put16 ACC_REP_PUTW
357     ahi_rep_put32 ACC_REP_PUTL
358     ahi_rep_put64 ACC_REP_PUTLL

360 on_trap_data
361     ot_prot
362     ot_trap
363     ot_trampoline
364     ot_jmpbuf
365     ot_prev
366     ot_handle
367     ot_padi

369 trap_trace_ctl_t  __TRAPTR_SIZE TRAPTR_SIZE_SHIFT
370     ttc_next      TRAPTR_NEXT
371     ttc_first     TRAPTR_FIRST
372     ttc_limit     TRAPTR_LIMIT

374 trap_trace_rec_t  TRAP_ENT_SIZE
375     ttr_cr2
376     ttr_info.idt_entry.vector  TTR_VECTOR
377     ttr_info.idt_entry.ipl     TTR_IPL
378     ttr_info.idt_entry.spl     TTR_SPL
379     ttr_info.idt_entry.pri     TTR_PRI
380     ttr_info.gate_entry.sysnum TTR_SYSNUM
381     ttr_marker
382     ttr_stamp
383     ttr_curthread
384     ttr_sdepth
385     ttr_stack

```

```

387 lgrp_ld
388     lpl_lgrp_id

390 dtrace_id_t      DTRACE_IDSIZE

392 cpu_core         CPU_CORE_SIZE   CPU_CORE_SHIFT
393     cpuc_dtrace_flags
394     cpuc_dtrace_illval

396 timespec        TIMESPEC_SIZE

398 gate_desc        GATE_DESC_SIZE

400 desctbr_t        DESC_TBR_SIZE
401     dtr_limit
402     dtr_base

404 mod_stub_info    MODS_SIZE
405     mods_func_addr   MODS_INSTFCN
406     mods_errfcn     MODS_RETFCN
407     mods_flag

409 \#define          TRAP_TSIZE          _MUL(TRAP_ENT_SIZE, TRAPTR_NENT)

411 copyops
412     cp_copyin
413     cp_xcopyin
414     cp_copyout
415     cp_xcopyout
416     cp_copyinstr
417     cp_copyoutstr
418     cp_fuword8
419     cp_fuword16
420     cp_fuword32
421     cp_fuword64
422     cp_suword8
423     cp_suword16
424     cp_suword32
425     cp_suword64
426     cp_physio

428 brand
429     b_machops

431 brand_proc_data_t
432     spd_handler

434 fastboot_file_t
435     fb_va
436     fb_pte_list_va
437     fb_pte_list_pa
438     fb_dest_pa
439     fb_size
440     fb_next_pa
441     fb_sections
442     fb_sectcnt

444 fastboot_section_t
445     fb_sec_offset
446     fb_sec_paddr
447     fb_sec_size
448     fb_sec_bss_size

450 fastboot_info_t
451     fi_files
452     fi_has_pae

```

```

453     fi_pagetable_va
454     fi_pagetable_pa
455     fi_last_table_pa
456     fi_new_mbi_pa
457     fi_valid

459 zone
460     zone_brand_data

462 wc_cpu   WC_CPU_SIZE
463     wc_retaddr
464     wc_virtaddr
465     wc_cr0
466     wc_cr3
467     wc_cr4
468     wc_cr8
469     wc_fs
470     wc_fsbase
471     wc_gs
472     wc_gsbase
473     wc_kgsbase
474     wc_r8
475     wc_r9
476     wc_r10
477     wc_r11
478     wc_r12
479     wc_r13
480     wc_r14
481     wc_r15
482     wc_rax
483     wc_rbp
484     wc_rbx
485     wc_rcx
486     wc_rdi
487     wc_rdx
488     wc_rsi
489     wc_rsp
490     wc_gdt_limit   WC_GDT
491     wc_gdt_base
492     wc_idt_limit   WC_IDT
493     wc_idt_base
494     wc_tr
495     wc_ldt
496     wc_eflags
497     wc_ebx
498     wc_edi
499     wc_esi
500     wc_ebp
501     wc_esp
502     wc_esp
503     wc_ss
504     wc_cs
505     wc_ds
506     wc_es
507     wc_cpu_id
508     wc_saved_stack

510 wc_wakecode
511     wc_cpu

513 comm_page_s   COMM_PAGE_S_SIZE

```

new/usr/src/uts/i86pc/ml/syscall_asm_amd64.s

1

```
*****
38851 Fri Apr 6 17:25:01 2018
new/usr/src/uts/i86pc/ml/syscall_asm_amd64.s
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22 * Copyright (c) 2004, 2010, Oracle and/or its affiliates. All rights reserved.
23 * Copyright 2018 Joyent, Inc.
23 * Copyright 2015 Joyent, Inc.
24 * Copyright (c) 2016 by Delphix. All rights reserved.
25 */

27 #include <sys/asm_linkage.h>
28 #include <sys/asm_misc.h>
29 #include <sys/regset.h>
30 #include <sys/privregs.h>
31 #include <sys/psw.h>
32 #include <sys/machbrand.h>

34 #if defined(__lint)

36 #include <sys/types.h>
37 #include <sys/thread.h>
38 #include <sys/system.h>

40 #else /* __lint */

42 #include <sys/segments.h>
43 #include <sys/pcb.h>
44 #include <sys/trap.h>
45 #include <sys/ftrace.h>
46 #include <sys/traptrace.h>
47 #include <sys/clock.h>
48 #include <sys/model.h>
49 #include <sys/panic.h>

51 #if defined(__xpv)
52 #include <sys/hypervisor.h>
53 #endif

55 #include "assym.h"

57 #endif /* __lint */
```

new/usr/src/uts/i86pc/ml/syscall_asm_amd64.s

2

```
59 /*
60 * We implement five flavours of system call entry points
61 *
62 * - syscall/sysretq (amd64 generic)
63 * - syscall/sysretl (i386 plus SYSC bit)
64 * - sysenter/sysexit (i386 plus SEP bit)
65 * - int/iret (i386 generic)
66 * - lcall/iret (i386 generic)
67 *
68 * The current libc included in Solaris uses int/iret as the base unoptimized
69 * kernel entry method. Older libc implementations and legacy binaries may use
70 * the lcall call gate, so it must continue to be supported.
71 *
72 * System calls that use an lcall call gate are processed in trap() via a
73 * segment-not-present trap, i.e. lcalls are extremely slow(!).
74 *
75 * The basic pattern used in the 32-bit SYSC handler at this point in time is
76 * to have the bare minimum of assembler, and get to the C handlers as
77 * quickly as possible.
78 *
79 * The 64-bit handler is much closer to the sparcv9 handler; that's
80 * because of passing arguments in registers. The 32-bit world still
81 * passes arguments on the stack -- that makes that handler substantially
82 * more complex.
83 *
84 * The two handlers share a few code fragments which are broken
85 * out into preprocessor macros below.
86 *
87 * XX64 come back and speed all this up later. The 32-bit stuff looks
88 * especially easy to speed up the argument copying part ..
89 *
90 *
91 * Notes about segment register usage (c.f. the 32-bit kernel)
92 *
93 * In the 32-bit kernel, segment registers are dutifully saved and
94 * restored on all mode transitions because the kernel uses them directly.
95 * When the processor is running in 64-bit mode, segment registers are
96 * largely ignored.
97 *
98 * %cs and %ss
99 * controlled by the hardware mechanisms that make mode transitions
100 *
101 * The remaining segment registers have to either be pointing at a valid
102 * descriptor i.e. with the 'present' bit set, or they can NULL descriptors
103 *
104 * %ds and %es
105 * always ignored
106 *
107 * %fs and %gs
108 * fsbase and gsbase are used to control the place they really point at.
109 * The kernel only depends on %gs, and controls its own gsbase via swapgs
110 *
111 * Note that loading segment registers is still costly because the GDT
112 * lookup still happens (this is because the hardware can't know that we're
113 * not setting up these segment registers for a 32-bit program). Thus we
114 * avoid doing this in the syscall path, and defer them to lwp context switch
115 * handlers, so the register values remain virtualized to the lwp.
116 */

118 #if defined(SYSCALLTRACE)
119 #define ORL_SYSCALLTRACE(r32) \
120 orl syscalltrace(%rip), r32
121 #else
122 #define ORL_SYSCALLTRACE(r32)
123 #endif
```

```

125 /*
126 * In the 32-bit kernel, we do absolutely nothing before getting into the
127 * brand callback checks. In 64-bit land, we do swagps and then come here.
128 * We assume that the %rsp- and %r15-stashing fields in the CPU structure
129 * are still unused.
130 *
131 * Check if a brand_mach_ops callback is defined for the specified callback_id
132 * type. If so invoke it with the kernel's %gs value loaded and the following
133 * data on the stack:
134 *
135 * stack: -----
136 *      32 | callback pointer
137 *      24 | user (or interrupt) stack pointer
138 *      16 | lwp pointer
139 *      v  8 | userland return address
140 *      0   | callback wrapper return addr
141 * -----
142 *
143 * Since we're pushing the userland return address onto the kernel stack
144 * we need to get that address without accessing the user's stack (since we
145 * can't trust that data). There are different ways to get the userland
146 * return address depending on how the syscall trap was made:
147 *
148 * a) For sys_syscall and sys_syscall32 the return address is in %rcx.
149 * b) For sys_sysenter the return address is in %rdx.
150 * c) For sys_int80 and sys_syscall_int (int91), upon entry into the macro,
151 * the stack pointer points at the state saved when we took the interrupt:
152 *
153 *      -----
154 *      | user's %ss
155 *      | user's %esp
156 *      | EFLAGS register
157 *      | user's %cs
158 *      | user's %eip
159 *      -----
160 *
161 * The 2nd parameter to the BRAND_CALLBACK macro is either the
162 * BRAND_URET_FROM_REG or BRAND_URET_FROM_INTR_STACK macro. These macros are
163 * used to generate the proper code to get the userland return address for
164 * each syscall entry point.
165 *
166 * The interface to the brand callbacks on the 64-bit kernel assumes %r15
167 * is available as a scratch register within the callback. If the callback
168 * returns within the kernel then this macro will restore %r15. If the
169 * callback is going to return directly to userland then it should restore
170 * %r15 before returning to userland.
171 */
172 #define BRAND_URET_FROM_REG(rip_reg) \
173     pushq   rip_reg                /* push the return address */
174
175 /*
176 * The interrupt stack pointer we saved on entry to the BRAND_CALLBACK macro
177 * is currently pointing at the user return address (%eip).
178 */
179 #define BRAND_URET_FROM_INTR_STACK() \
180     movq    %gs:CPU_RTMP_RSP, %r15 /* grab the intr. stack pointer */ ; \
181     pushq   (%r15)                /* push the return address */
182
183 #define BRAND_CALLBACK(callback_id, push_userland_ret) \
184     movq    %rsp, %gs:CPU_RTMP_RSP /* save the stack pointer */ ; \
185     movq    %r15, %gs:CPU_RTMP_R15 /* save %r15 */ ; \
186     movq    %gs:CPU_THREAD, %r15 /* load the thread pointer */ ; \
187     movq    T_STACK(%r15), %rsp /* switch to the kernel stack */ ; \
188     subq    $16, %rsp /* save space for 2 pointers */ ; \
189     pushq   %r14 /* save %r14 */ ; \
190     movq    %gs:CPU_RTMP_RSP, %r14 /* restore %r14 */ ; \
191     movq    %r14, 8(%rsp) /* stash the user stack pointer */ ; \

```

```

191     popq    %r14 /* restore %r14 */ ; \
192     movq    T_LWP(%r15), %r15 /* load the lwp pointer */ ; \
193     pushq   %r15 /* push the lwp pointer */ ; \
194     movq    LWP_PROCP(%r15), %r15 /* load the proc pointer */ ; \
195     movq    P_BRAND(%r15), %r15 /* load the brand pointer */ ; \
196     movq    B_MACHOPS(%r15), %r15 /* load the machops pointer */ ; \
197     movq    _CONST(_MUL(callback_id, CPTRSIZE))(%r15), %r15 /* load callback_id */ ; \
198     cmpq    $0, %r15 /* if callback_id == 0 */ ; \
199     je      lf /* if callback_id == 0, jump to lf */ ; \
200     movq    %r15, 16(%rsp) /* save the callback pointer */ ; \
201     push_userland_ret /* push the return address */ ; \
202     call    *24(%rsp) /* call callback */ ; \
203 1:     movq    %gs:CPU_RTMP_R15, %r15 /* restore %r15 */ ; \
204     movq    %gs:CPU_RTMP_RSP, %rsp /* restore the stack pointer */ ; \
205
206 #define MSTATE_TRANSITION(from, to) \
207     movl    $from, %edi ; \
208     movl    $to, %esi ; \
209     call    syscall_mstate
210
211 /*
212 * Check to see if a simple (direct) return is possible i.e.
213 *
214 * if (t->t_post_sys_ast | syscalltrace |
215 *      lwp->lwp_pcb.pcb_rupdate == 1)
216 *     do full version ;
217 *
218 * Preconditions:
219 * - t is curthread
220 * Postconditions:
221 * - condition code NE is set if post-sys is too complex
222 * - rtmp is zeroed if it isn't (we rely on this!)
223 * - ltmp is smashed
224 */
225 #define CHECK_POSTSYS_NE(t, ltmp, rtmp) \
226     movq    T_LWP(t), ltmp ; \
227     movzbl  PCB_RUPDATE(ltmp), rtmp ; \
228     ORL_SYSCALLTRACE(rtmp) ; \
229     orl    T_POST_SYS_AST(t), rtmp ; \
230     cmpl   $0, rtmp
231
232 /*
233 * Fix up the lwp, thread, and eflags for a successful return
234 *
235 * Preconditions:
236 * - zwreg contains zero
237 */
238 #define SIMPLE_SYSCALL_POSTSYS(t, lwp, zwreg) \
239     movb    $LWP_USER, LWP_STATE(lwp) ; \
240     movw    zwreg, T_SYSNUM(t) ; \
241     andb    $_CONST(0xffff - PS_C), REGOFF_RFL(%rsp)
242
243 /*
244 * ASSERT(lwptoregs(lwp) == rp);
245 *
246 * This may seem obvious, but very odd things happen if this
247 * assertion is false
248 *
249 * Preconditions:
250 * (%rsp is ready for normal call sequence)
251 * Postconditions (if assertion is true):
252 * %r11 is smashed
253 *
254 * ASSERT(rp->r_cs == descnum)
255 *
256 * The code selector is written into the regs structure when the

```



```

257 * lwp stack is created. We use this ASSERT to validate that
258 * the regs structure really matches how we came in.
259 *
260 * Preconditions:
261 *   (%rsp is ready for normal call sequence)
262 * Postconditions (if assertion is true):
263 *   -none-
264 *
265 * ASSERT(lwp->lwp_pcb.pcb_rupdate == 0);
266 *
267 * If this is false, it meant that we returned to userland without
268 * updating the segment registers as we were supposed to.
269 *
270 * Note that we must ensure no interrupts or other traps intervene
271 * between entering privileged mode and performing the assertion,
272 * otherwise we may perform a context switch on the thread, which
273 * will end up setting pcb_rupdate to 1 again.
274 */
275 #if defined(DEBUG)

277 #if !defined(__lint)

279 __lwptoregs_msg:
280     .string "syscall_asm_amd64.s:%d lwptoregs(%p) [%p] != rp [%p]"

282 __codesel_msg:
283     .string "syscall_asm_amd64.s:%d rp->r_cs [%ld] != %ld"

285 __no_rupdate_msg:
286     .string "syscall_asm_amd64.s:%d lwp %p, pcb_rupdate != 0"

288 #endif /* !__lint */

290 #define ASSERT_LWPTOREGS(lwp, rp) \
291     movq   LWP_REGS(lwp), %r11; \
292     cmpq   rp, %r11; \
293     je     7f; \
294     leaq  __lwptoregs_msg(%rip), %rdi; \
295     movl  $_LINE__, %esi; \
296     movq  lwp, %rdx; \
297     movq  %r11, %rcx; \
298     movq  rp, %r8; \
299     xorl  %eax, %eax; \
300     call  panic; \
301 7:

303 #define ASSERT_NO_RUPDATE_PENDING(lwp) \
304     testb  $0x1, PCB_RUPDATE(lwp); \
305     je     8f; \
306     movq  lwp, %rdx; \
307     leaq  __no_rupdate_msg(%rip), %rdi; \
308     movl  $_LINE__, %esi; \
309     xorl  %eax, %eax; \
310     call  panic; \
311 8:

313 #else
314 #define ASSERT_LWPTOREGS(lwp, rp)
315 #define ASSERT_NO_RUPDATE_PENDING(lwp)
316 #endif

318 /*
319 * Do the traptrace thing and restore any registers we used
320 * in situ. Assumes that %rsp is pointing at the base of
321 * the struct regs, obviously ..
322 */

```

```

323 #ifdef TRAPTRACE
324 #define SYSCALL_TRAPTRACE(ttype) \
325     TRACE_PTR(%rdi, %rbx, %ebx, %rcx, ttype); \
326     TRACE_REGS(%rdi, %rsp, %rbx, %rcx); \
327     TRACE_STAMP(%rdi); /* rdtsc clobbers %eax, %edx */ \
328     movq   REGOFF_RAX(%rsp), %rax; \
329     movq   REGOFF_RBX(%rsp), %rbx; \
330     movq   REGOFF_RCX(%rsp), %rcx; \
331     movq   REGOFF_RDX(%rsp), %rdx; \
332     movl   %eax, TTR_SYSNUM(%rdi); \
333     movq   REGOFF_RDI(%rsp), %rdi

335 #define SYSCALL_TRAPTRACE32(ttype) \
336     SYSCALL_TRAPTRACE(ttype); \
337     /* paranoia: clean the top 32-bits of the registers */ \
338     orl   %eax, %eax; \
339     orl   %ebx, %ebx; \
340     orl   %ecx, %ecx; \
341     orl   %edx, %edx; \
342     orl   %edi, %edi

343 #else /* TRAPTRACE */
344 #define SYSCALL_TRAPTRACE(ttype)
345 #define SYSCALL_TRAPTRACE32(ttype)
346 #endif /* TRAPTRACE */

348 /*
349 * The 64-bit libc syscall wrapper does this:
350 *
351 * fn(<args>)
352 * {
353 *     movq   %rcx, %r10    -- because syscall smashes %rcx
354 *     movl   $CODE, %eax
355 *     syscall
356 *     <error processing>
357 * }

359 * Thus when we come into the kernel:
360 *
361 * %rdi, %rsi, %rdx, %r10, %r8, %r9 contain first six args
362 * %rax is the syscall number
363 * %r12-%r15 contain caller state
364 *
365 * The syscall instruction arranges that:
366 *
367 * %rcx contains the return %rip
368 * %r11 contains bottom 32-bits of %rflags
369 * %rflags is masked (as determined by the SFMASK msr)
370 * %cs is set to UCS_SEL (as determined by the STAR msr)
371 * %ss is set to UDS_SEL (as determined by the STAR msr)
372 * %rip is set to sys_syscall (as determined by the LSTAR msr)
373 *
374 * Or in other words, we have no registers available at all.
375 * Only swapgs can save us!
376 *
377 * Under the hypervisor, the swapgs has happened already. However, the
378 * state of the world is very different from that we're familiar with.
379 *
380 * In particular, we have a stack structure like that for interrupt
381 * gates, except that the %cs and %ss registers are modified for reasons
382 * that are not entirely clear. Critically, the %rcx/%r11 values do
383 * *not* reflect the usage of those registers under a 'real' syscall[1];
384 * the stack, therefore, looks like this:
385 *
386 * 0x0(rsp)    potentially junk %rcx
387 * 0x8(rsp)    potentially junk %r11
388 * 0x10(rsp)   user %rip

```

```

389 *      0x18(%rsp)      modified %cs
390 *      0x20(%rsp)      user %rflags
391 *      0x28(%rsp)      user %rsp
392 *      0x30(%rsp)      modified %ss
393 *
394 *
395 * and before continuing on, we must load the %rip into %rcx and the
396 * %rflags into %r11.
397 *
398 * [1] They used to, and we relied on it, but this was broken in 3.1.1.
399 * Sigh.
400 */
401 #if defined(__xpv)
402 #define XPV_SYSCALL_PROD \
403     movq    0x10(%rsp), %rcx; \
404     movq    0x20(%rsp), %r11; \
405     movq    0x28(%rsp), %rsp
406 #else
407 #define XPV_SYSCALL_PROD /* nothing */
408 #endif
409
410 #if defined(__lint)
411
412 /*ARGSUSED*/
413 void
414 sys_syscall()
415 {}
416
417 void
418 _allsyscalls()
419 {}
420
421 size_t _allsyscalls_size;
422
423 #else /* __lint */
424
425     ENTRY_NP2(brand_sys_syscall, _allsyscalls)
426     SWAPGS /* kernel gsbase */
427     XPV_SYSCALL_PROD
428     BRAND_CALLBACK(BRAND_CB_SYSCALL, BRAND_URET_FROM_REG(%rcx))
429     jmp     noprod_sys_syscall
430
431     ALTENTRY(sys_syscall)
432     SWAPGS /* kernel gsbase */
433     XPV_SYSCALL_PROD
434
435 noprod_sys_syscall:
436     movq    %r15, %gs:CPU_RTMP_R15
437     movq    %rsp, %gs:CPU_RTMP_RSP
438
439     movq    %gs:CPU_THREAD, %r15
440     movq    T_STACK(%r15), %rsp /* switch from user to kernel stack */
441
442     ASSERT_UPCALL_MASK_IS_SET
443
444     movl    $UCS_SEL, REGOFF_CS(%rsp)
445     movq    %rcx, REGOFF_RIP(%rsp) /* syscall: %rip -> %rcx */
446     movq    %r11, REGOFF_RFL(%rsp) /* syscall: %rfl -> %r11d */
447     movl    $UDS_SEL, REGOFF_SS(%rsp)
448
449     movl    %eax, %eax /* wrapper: sysc# -> %eax */
450     movq    %rdi, REGOFF_RDI(%rsp)
451     movq    %rsi, REGOFF_RSI(%rsp)
452     movq    %rdx, REGOFF_RDX(%rsp)
453     movq    %r10, REGOFF_RCX(%rsp) /* wrapper: %rcx -> %r10 */
454     movq    %r10, %rcx /* arg[3] for direct calls */

```

```

456     movq    %r8, REGOFF_R8(%rsp)
457     movq    %r9, REGOFF_R9(%rsp)
458     movq    %rax, REGOFF_RAX(%rsp)
459     movq    %rbx, REGOFF_RBX(%rsp)
460
461     movq    %rbp, REGOFF_RBP(%rsp)
462     movq    %r10, REGOFF_R10(%rsp)
463     movq    %gs:CPU_RTMP_RSP, %r11
464     movq    %r11, REGOFF_RSP(%rsp)
465     movq    %r12, REGOFF_R12(%rsp)
466
467     movq    %r13, REGOFF_R13(%rsp)
468     movq    %r14, REGOFF_R14(%rsp)
469     movq    %gs:CPU_RTMP_R15, %r10
470     movq    %r10, REGOFF_R15(%rsp)
471     movq    $0, REGOFF_SAVFP(%rsp)
472     movq    $0, REGOFF_SAVPC(%rsp)
473
474     /*
475     * Copy these registers here in case we end up stopped with
476     * someone (like, say, /proc) messing with our register state.
477     * We don't -restore- them unless we have to in update_sregs.
478     *
479     * Since userland -can't- change fsbase or gsbase directly,
480     * and capturing them involves two serializing instructions,
481     * we don't bother to capture them here.
482     */
483     xorl    %ebx, %ebx
484     movw    %ds, %bx
485     movq    %rbx, REGOFF_DS(%rsp)
486     movw    %es, %bx
487     movq    %rbx, REGOFF_ES(%rsp)
488     movw    %fs, %bx
489     movq    %rbx, REGOFF_FS(%rsp)
490     movw    %gs, %bx
491     movq    %rbx, REGOFF_GS(%rsp)
492
493     /*
494     * If we're trying to use TRAPTRACE though, I take that back: we're
495     * probably debugging some problem in the SWAPGS logic and want to know
496     * what the incoming gsbase was.
497     *
498     * Since we already did SWAPGS, record the KGSBASE.
499     */
500     #if defined(DEBUG) && defined(TRAPTRACE) && !defined(__xpv)
501     movl    $MSR_AMD_KGSBASE, %ecx
502     rdmsr
503     movl    %eax, REGOFF_GSBASE(%rsp)
504     movl    %edx, REGOFF_GSBASE+4(%rsp)
505     #endif
506
507     /*
508     * Machine state saved in the regs structure on the stack
509     * First six args in %rdi, %rsi, %rdx, %rcx, %r8, %r9
510     * %eax is the syscall number
511     * %rsp is the thread's stack, %r15 is curthread
512     * REG_RSP(%rsp) is the user's stack
513     */
514
515     SYSCALL_TRAPTRACE($TT_SYSC64)
516
517     movq    %rsp, %rbp
518
519     movq    T_LWP(%r15), %r14
520     ASSERT_NO_RUPDATE_PENDING(%r14)

```

```

521     ENABLE_INTR_FLAGS

523     MSTATE_TRANSITION(LMS_USER, LMS_SYSTEM)
524     movl   REGOFF_RAX(%rsp), %eax /* (%rax damaged by mstate call) */

526     ASSERT_LWPTOREGS(%r14, %rsp)

528     movb   $LWP_SYS, LWP_STATE(%r14)
529     incq   LWP_RU_SYSC(%r14)
530     movb   $NORMALRETURN, LWP_EOSYS(%r14)

532     incq   %gs:CPU_STATS_SYS_SYSCALL

534     movw   %ax, T_SYSNUM(%r15)
535     movzbl T_PRE_SYS(%r15), %ebx
536     ORL_SYSCALLTRACE(%ebx)
537     testl  %ebx, %ebx
538     jne    _syscall_pre

540 _syscall_invoke:
541     movq   REGOFF_RDI(%rbp), %rdi
542     movq   REGOFF_RSI(%rbp), %rsi
543     movq   REGOFF_RDX(%rbp), %rdx
544     movq   REGOFF_RCX(%rbp), %rcx
545     movq   REGOFF_R8(%rbp), %r8
546     movq   REGOFF_R9(%rbp), %r9

548     cmpl  $NSYSCALL, %eax
549     jae   _syscall_ill
550     shll  $SYSENT_SIZE_SHIFT, %eax
551     leaq  sysent(%rax), %rbx

553     call  *SY_CALLC(%rbx)

555     movq  %rax, %r12
556     movq  %rdx, %r13

558     /*
559     * If the handler returns two ints, then we need to split the
560     * 64-bit return value into two 32-bit values.
561     */
562     testw  $SE_32RVAL2, SY_FLAGS(%rbx)
563     je     5f
564     movq   %r12, %r13
565     shrq   $32, %r13 /* upper 32-bits into %edx */
566     movl   %r12d, %r12d /* lower 32-bits into %eax */
567 5:
568     /*
569     * Optimistically assume that there's no post-syscall
570     * work to do. (This is to avoid having to call syscall_mstate())
571     * with interrupts disabled)
572     */
573     MSTATE_TRANSITION(LMS_SYSTEM, LMS_USER)

575     /*
576     * We must protect ourselves from being descheduled here;
577     * If we were, and we ended up on another cpu, or another
578     * lwp got in ahead of us, it could change the segment
579     * registers without us noticing before we return to userland.
580     */
581     CLI(%r14)
582     CHECK_POSTSYS_NE(%r15, %r14, %ebx)
583     jne    _syscall_post

585     /*
586     * We need to protect ourselves against non-canonical return values

```

```

587     * because Intel doesn't check for them on sysret (AMD does). Canonical
588     * addresses on current amd64 processors only use 48-bits for VAs; an
589     * address is canonical if all upper bits (47-63) are identical. If we
590     * find a non-canonical %rip, we opt to go through the full
591     * _syscall_post path which takes us into an iretq which is not
592     * susceptible to the same problems sysret is.
593     *
594     * We're checking for a canonical address by first doing an arithmetic
595     * shift. This will fill in the remaining bits with the value of bit 63.
596     * If the address were canonical, the register would now have either all
597     * zeroes or all ones in it. Therefore we add one (inducing overflow)
598     * and compare against 1. A canonical address will either be zero or one
599     * at this point, hence the use of ja.
600     *
601     * At this point, r12 and r13 have the return value so we can't use
602     * those registers.
603     */
604     movq   REGOFF_RIP(%rsp), %rcx
605     sarq   $47, %rcx
606     incq   %rcx
607     cmpq   $1, %rcx
608     ja     _syscall_post

611     SIMPLE_SYSCALL_POSTSYS(%r15, %r14, %bx)

613     movq   %r12, REGOFF_RAX(%rsp)
614     movq   %r13, REGOFF_RDX(%rsp)

616     /*
617     * To get back to userland, we need the return %rip in %rcx and
618     * the return %rfl in %r1d. The sysretq instruction also arranges
619     * to fix up %cs and %ss; everything else is our responsibility.
620     */
621     movq   REGOFF_RDI(%rsp), %rdi
622     movq   REGOFF_RSI(%rsp), %rsi
623     movq   REGOFF_RDX(%rsp), %rdx
624     /* %rcx used to restore %rip value */

626     movq   REGOFF_R8(%rsp), %r8
627     movq   REGOFF_R9(%rsp), %r9
628     movq   REGOFF_RAX(%rsp), %rax
629     movq   REGOFF_RBX(%rsp), %rbx

631     movq   REGOFF_RBP(%rsp), %rbp
632     movq   REGOFF_R10(%rsp), %r10
633     /* %r11 used to restore %rfl value */
634     movq   REGOFF_R12(%rsp), %r12

636     movq   REGOFF_R13(%rsp), %r13
637     movq   REGOFF_R14(%rsp), %r14
638     movq   REGOFF_R15(%rsp), %r15

640     movq   REGOFF_RIP(%rsp), %rcx
641     movl   REGOFF_RFL(%rsp), %r1d

643 #if defined(__xpv)
644     addq   $REGOFF_RIP, %rsp
645 #else
646     movq   REGOFF_RSP(%rsp), %rsp
647 #endif

649     /*
650     * There can be no instructions between the ALTENTRY below and
651     * SYSRET or we could end up breaking brand support. See label usage
652     * in snl_brand_syscall_callback for an example.

```

```

653     */
654     ASSERT_UPCALL_MASK_IS_SET
655 #if defined(__xpv)
656     SYSRETQ
657     ALTENTRY(nopop_sys_syscall_swaggs_sysretq)
658
659     /*
660     * We can only get here after executing a brand syscall
661     * interposition callback handler and simply need to
662     * "sysretq" back to userland. On the hypervisor this
663     * involves the iret hypercall which requires us to construct
664     * just enough of the stack needed for the hypercall.
665     * (rip, cs, rflags, rsp, ss).
666     */
667     movq    %rsp, %gs:CPU_RTMP_RSP          /* save user's rsp */
668     movq    %gs:CPU_THREAD, %r11
669     movq    T_STACK(%r11), %rsp
670
671     movq    %rcx, REGOFF_RIP(%rsp)
672     movl    $UCS_SEL, REGOFF_CS(%rsp)
673     movq    %gs:CPU_RTMP_RSP, %r11
674     movq    %r11, REGOFF_RSP(%rsp)
675     pushfq
676     popq   %r11                          /* hypercall enables ints */
677     movq    %r11, REGOFF_RFL(%rsp)
678     movl    $UDS_SEL, REGOFF_SS(%rsp)
679     addq    $REGOFF_RIP, %rsp
680     /*
681     * XXPV: see comment in SYSRETQ definition for future optimization
682     *       we could take.
683     */
684     ASSERT_UPCALL_MASK_IS_SET
685     SYSRETQ
686 #else
687     ALTENTRY(nopop_sys_syscall_swaggs_sysretq)
688     jmp     tr_sysretq
689     SWAPGS                                /* user gsbase */
690     SYSRETQ
691 #endif
692
693     /*NOTREACHED*/
694     SET_SIZE(nopop_sys_syscall_swaggs_sysretq)
695     unchanged_portion_omitted
696 #endif /* __lint */
697
698 #if defined(__lint)
699
700 /*ARGSUSED*/
701 void
702 sys_syscall32()
703 {}
704
705 #else /* __lint */
706
707     ENTRY_NP(brand_sys_syscall32)
708     SWAPGS                                /* kernel gsbase */
709     XPV_TRAP_POP
710     BRAND_CALLBACK(BRAND_CB_SYSCALL32, BRAND_URET_FROM_REG(%rcx))
711     jmp     nopop_sys_syscall32
712
713     ALTENTRY(sys_syscall32)
714     SWAPGS                                /* kernel gsbase */
715     XPV_TRAP_POP
716
717 nopop_sys_syscall32:
718     movl    %esp, %r10d

```

```

719     movq    %gs:CPU_THREAD, %r15
720     movq    T_STACK(%r15), %rsp
721     movl    %eax, %eax
722
723     movl    $U32CS_SEL, REGOFF_CS(%rsp)
724     movl    %ecx, REGOFF_RIP(%rsp)
725     movq    %r11, REGOFF_RFL(%rsp)      /* syscall: %rip -> %rcx */
726     movq    %r10, REGOFF_RSP(%rsp)      /* syscall: %rfl -> %r10 */
727     movl    $UDS_SEL, REGOFF_SS(%rsp)
728
729 _syscall32_save:
730     movl    %edi, REGOFF_RDI(%rsp)
731     movl    %esi, REGOFF_RSI(%rsp)
732     movl    %ebp, REGOFF_RBP(%rsp)
733     movl    %ebx, REGOFF_RBX(%rsp)
734     movl    %edx, REGOFF_RDX(%rsp)
735     movl    %ecx, REGOFF_RCX(%rsp)
736     movl    %eax, REGOFF_RAX(%rsp)      /* wrapper: sysc# -> %eax */
737     movq    $0, REGOFF_SAVFP(%rsp)
738     movq    $0, REGOFF_SAVPC(%rsp)
739
740     /*
741     * Copy these registers here in case we end up stopped with
742     * someone (like, say, /proc) messing with our register state.
743     * We don't -restore- them unless we have to in update_sregs.
744     */
745     * Since userland -can't- change fsbase or gsbase directly,
746     * we don't bother to capture them here.
747     */
748     xorl    %ebx, %ebx
749     movw    %ds, %bx
750     movq    %rbx, REGOFF_DS(%rsp)
751     movw    %es, %bx
752     movq    %rbx, REGOFF_ES(%rsp)
753     movw    %fs, %bx
754     movq    %rbx, REGOFF_FS(%rsp)
755     movw    %gs, %bx
756     movq    %rbx, REGOFF_GS(%rsp)
757
758     /*
759     * If we're trying to use TRAPTRACE though, I take that back: we're
760     * probably debugging some problem in the SWAPGS logic and want to know
761     * what the incoming gsbase was.
762     */
763     * Since we already did SWAPGS, record the KGSBASE.
764     */
765 #if defined(DEBUG) && defined(TRAPTRACE) && !defined(__xpv)
766     movl    $MSR_AMD_KGSBASE, %ecx
767     rdmsr
768     movl    %eax, REGOFF_GSBASE(%rsp)
769     movl    %edx, REGOFF_GSBASE+4(%rsp)
770 #endif
771
772     /*
773     * Application state saved in the regs structure on the stack
774     * %eax is the syscall number
775     * %rsp is the thread's stack, %r15 is curthread
776     * REG_RSP(%rsp) is the user's stack
777     */
778
779     SYSCALL_TRAPTRACE32($TT_SYSC)
780
781     movq    %rsp, %rbp
782
783     movq    T_LWP(%r15), %r14
784     ASSERT_NO_RUPDATE_PENDING(%r14)

```

```

816     ENABLE_INTR_FLAGS

818     MSTATE_TRANSITION(LMS_USER, LMS_SYSTEM)
819     movl    REGOFF_RAX(%rsp), %eax /* (%rax damaged by mstate call) */

821     ASSERT_LWPTOREGS(%r14, %rsp)

823     incq    %gs:CPU_STATS_SYS_SYSCALL

825     /*
826     * Make some space for MAXSYSARGS (currently 8) 32-bit args placed
827     * into 64-bit (long) arg slots, maintaining 16 byte alignment. Or
828     * more succinctly:
829     *
830     * SA(MAXSYSARGS * sizeof (long)) == 64
831     */
832 #define SYS_DROP        64          /* drop for args */
833     subq    $SYS_DROP, %rsp
834     movb    $LWP_SYS, LWP_STATE(%r14)
835     movq    %r15, %rdi
836     movq    %rsp, %rsi
837     call   syscall_entry

839     /*
840     * Fetch the arguments copied onto the kernel stack and put
841     * them in the right registers to invoke a C-style syscall handler.
842     * %rax contains the handler address.
843     *
844     * Ideas for making all this go faster of course include simply
845     * forcibly fetching 6 arguments from the user stack under lofault
846     * protection, reverting to copyin_args only when watchpoints
847     * are in effect.
848     *
849     * (If we do this, make sure that exec and libthread leave
850     * enough space at the top of the stack to ensure that we'll
851     * never do a fetch from an invalid page.)
852     *
853     * Lots of ideas here, but they won't really help with bringup B-)
854     * Correctness can't wait, performance can wait a little longer ..
855     */

857     movq    %rax, %rbx
858     movl    0(%rsp), %edi
859     movl    8(%rsp), %esi
860     movl    0x10(%rsp), %edx
861     movl    0x18(%rsp), %ecx
862     movl    0x20(%rsp), %r8d
863     movl    0x28(%rsp), %r9d

865     call   *SY_CALLC(%rbx)

867     movq    %rbp, %rsp /* pop the args */

869     /*
870     * amd64 syscall handlers -always- return a 64-bit value in %rax.
871     * On the 32-bit kernel, they always return that value in %eax:%edx
872     * as required by the 32-bit ABI.
873     *
874     * Simulate the same behaviour by unconditionally splitting the
875     * return value in the same way.
876     */
877     movq    %rax, %r13
878     shrq    $32, %r13 /* upper 32-bits into %edx */
879     movl    %eax, %r12d /* lower 32-bits into %eax */

```

```

881     /*
882     * Optimistically assume that there's no post-syscall
883     * work to do. (This is to avoid having to call syscall_mstate())
884     * with interrupts disabled)
885     */
886     MSTATE_TRANSITION(LMS_SYSTEM, LMS_USER)

888     /*
889     * We must protect ourselves from being descheduled here;
890     * If we were, and we ended up on another cpu, or another
891     * lwp got in ahead of us, it could change the segment
892     * registers without us noticing before we return to userland.
893     */
894     CLI(%r14)
895     CHECK_POSTSYS_NE(%r15, %r14, %ebx)
896     jne     __full_syscall_postsys32
897     SIMPLE_SYSCALL_POSTSYS(%r15, %r14, %bx)

899     /*
900     * To get back to userland, we need to put the return %rip in %rcx and
901     * the return %rfl in %r11d. The sysret instruction also arranges
902     * to fix up %cs and %ss; everything else is our responsibility.
903     */

905     movl    %r12d, %eax /* %eax: rval1 */
906     movl    REGOFF_RBX(%rsp), %ebx
907     /* %ecx used for return pointer */
908     movl    %r13d, %edx /* %edx: rval2 */
909     movl    REGOFF_RBP(%rsp), %ebp
910     movl    REGOFF_RSI(%rsp), %esi
911     movl    REGOFF_RDI(%rsp), %edi

913     movl    REGOFF_RFL(%rsp), %r11d /* %r11 -> eflags */
914     movl    REGOFF_RIP(%rsp), %ecx /* %ecx -> %eip */
915     movl    REGOFF_RSP(%rsp), %esp

917     ASSERT_UPCALL_MASK_IS_SET
918     ALTENTRY(nopop_sys_syscall32_swapgs_sysretl)
919     jmp     tr_sysretl
920     SWAPGS /* user gsbases */
921     SYSRETL
922     SET_SIZE(nopop_sys_syscall32_swapgs_sysretl)

```

unchanged portion omitted

```

939 #endif /* __lint */

941 /*
942 * System call handler via the sysenter instruction
943 * Used only for 32-bit system calls on the 64-bit kernel.
944 *
945 * The caller in userland has arranged that:
946 *
947 * - %eax contains the syscall number
948 * - %ecx contains the user %esp
949 * - %edx contains the return %eip
950 * - the user stack contains the args to the syscall
951 *
952 * Hardware and (privileged) initialization code have arranged that by
953 * the time the sysenter instructions completes:
954 *
955 * - %rip is pointing to sys_sysenter (below).
956 * - %cs and %ss are set to kernel text and stack (data) selectors.
957 * - %rsp is pointing at the lwp's stack
958 * - interrupts have been disabled.
959 *
960 * Note that we are unable to return both "rvals" to userland with

```

```

961 * this call, as %edx is used by the sysexit instruction.
962 *
963 * One final complication in this routine is its interaction with
964 * single-stepping in a debugger. For most of the system call mechanisms, the
965 * CPU automatically clears the single-step flag before we enter the kernel.
966 * The sysenter mechanism does not clear the flag, so a user single-stepping
967 * through a libc routine may suddenly find themselves single-stepping through the
968 * kernel. To detect this, kmdb and trap() both compare the trap %pc to the
969 * [brand_]sys_enter addresses on each single-step trap. If it finds that we
970 * have single-stepped to a sysenter entry point, it explicitly clears the flag
971 * and executes the sys_sysenter routine.
972 *
973 * single-stepping in a debugger. For most of the system call mechanisms,
974 * the CPU automatically clears the single-step flag before we enter the
975 * kernel. The sysenter mechanism does not clear the flag, so a user
976 * single-stepping through a libc routine may suddenly find themselves
977 * single-stepping through the kernel. To detect this, kmdb compares the
978 * trap %pc to the [brand_]sys_enter addresses on each single-step trap.
979 * If it finds that we have single-stepped to a sysenter entry point, it
980 * explicitly clears the flag and executes the sys_sysenter routine.
981 */
982 #if defined(__lint)
983 void
984 sys_sysenter()
985 {}
986 #else /* __lint */
987 #define ENTRY_NP(brand_sys_sysenter) \
988     SWAPGS /* kernel gsbase */ \
989     ALTENTRY(_brand_sys_sysenter_post_swapgs)
990 #define BRAND_CALLBACK(BRAND_CB_SYSENTER, BRAND_URET_FROM_REG(%rdx)) \
991     /* \
992     * Jump over sys_sysenter to allow single-stepping as described \
993     * above. \
994     */ \
995     jmp _sys_sysenter_post_swapgs
996 #endif
1000 ALTENTRY(sys_sysenter)
1001 SWAPGS /* kernel gsbase */
1002 ALTENTRY(_sys_sysenter_post_swapgs)
1003 #define ALTENTRY(_sys_sysenter_post_swapgs) \
1004     movq %gs:CPU_THREAD, %r15
1005 #endif
1006 movl $U32CS_SEL, REGOFF_CS(%rsp)
1007 movl %ecx, REGOFF_RSP(%rsp) /* wrapper: %esp -> %ecx */
1008 movl %edx, REGOFF_RIP(%rsp) /* wrapper: %eip -> %edx */
1009 /*

```

```

1010 * NOTE: none of the instructions that run before we get here should
1011 * clobber bits in (R)FLAGS! This includes the kpti trampoline.
1012 */
1013 pushfq
1014 popq %r10
1015 movl $UDS_SEL, REGOFF_SS(%rsp)
1016
1017 /*
1018 * Set the interrupt flag before storing the flags to the
1019 * flags image on the stack so we can return to user with
1020 * interrupts enabled if we return via sys_rtt_syscall32
1021 */
1022 orq $SPS_IE, %r10
1023 movq %r10, REGOFF_RFL(%rsp)
1024
1025 movl %edi, REGOFF_RDI(%rsp)
1026 movl %esi, REGOFF_RSI(%rsp)
1027 movl %ebp, REGOFF_RBP(%rsp)
1028 movl %ebx, REGOFF_RBX(%rsp)
1029 movl %edx, REGOFF_RDX(%rsp)
1030 movl %ecx, REGOFF_RCX(%rsp)
1031 movl %eax, REGOFF_RAX(%rsp) /* wrapper: sysc# -> %eax */
1032 movq $0, REGOFF_SAVFP(%rsp)
1033 movq $0, REGOFF_SAVPC(%rsp)
1034
1035 /*
1036 * Copy these registers here in case we end up stopped with
1037 * someone (like, say, /proc) messing with our register state.
1038 * We don't -restore- them unless we have to in update_sregs.
1039 */
1040 * Since userland -can't- change fsbase or gsbase directly,
1041 * we don't bother to capture them here.
1042 */
1043 xorl %ebx, %ebx
1044 movw %ds, %bx
1045 movq %rbx, REGOFF_DS(%rsp)
1046 movw %es, %bx
1047 movq %rbx, REGOFF_ES(%rsp)
1048 movw %fs, %bx
1049 movq %rbx, REGOFF_FS(%rsp)
1050 movw %gs, %bx
1051 movq %rbx, REGOFF_GS(%rsp)
1052
1053 /*
1054 * If we're trying to use TRAPTRACE though, I take that back: we're
1055 * probably debugging some problem in the SWAPGS logic and want to know
1056 * what the incoming gsbase was.
1057 */
1058 * Since we already did SWAPGS, record the KGSBASE.
1059 */
1060 #if defined(DEBUG) && defined(TRAPTRACE) && !defined(__xpv)
1061 movl $MSR_AMD_KGSBASE, %ecx
1062 rdmsr
1063 movl %eax, REGOFF_GSBASE(%rsp)
1064 movl %edx, REGOFF_GSBASE+4(%rsp)
1065 #endif
1066
1067 /*
1068 * Application state saved in the regs structure on the stack
1069 * %eax is the syscall number
1070 * %rsp is the thread's stack, %r15 is curthread
1071 * REG_RSP(%rsp) is the user's stack
1072 */
1073
1074 SYSCALL_TRAPTRACE($TT_SYSENTER)

```

```

1076     movq    %rsp, %rbp
1078     movq    T_LWP(%r15), %r14
1079     ASSERT_NO_RUPDATE_PENDING(%r14)
1081     ENABLE_INTR_FLAGS
1083     /*
1084     * Catch 64-bit process trying to issue sysenter instruction
1085     * on Nocona based systems.
1086     */
1087     movq    LWP_PROCP(%r14), %rax
1088     cmpq    $DATAMODEL_ILP32, P_MODEL(%rax)
1089     je      7f
1091     /*
1092     * For a non-32-bit process, simulate a #ud, since that's what
1093     * native hardware does. The traptrace entry (above) will
1094     * let you know what really happened.
1095     */
1096     movq    $T_ILLINST, REGOFF_TRAPNO(%rsp)
1097     movq    REGOFF_CS(%rsp), %rdi
1098     movq    %rdi, REGOFF_ERR(%rsp)
1099     movq    %rsp, %rdi
1100     movq    REGOFF_RIP(%rsp), %rsi
1101     movl    %gs:CPU_ID, %edx
1102     call   trap
1103     jmp    _sys_rtt
1104 7:
1106     MSTATE_TRANSITION(LMS_USER, LMS_SYSTEM)
1107     movl    REGOFF_RAX(%rsp), %eax /* (%rax damaged by mstate calls) */
1109     ASSERT_LWPTOREGS(%r14, %rsp)
1111     incq    %gs:CPU_STATS_SYS_SYSCALL
1113     /*
1114     * Make some space for MAXSYSARGS (currently 8) 32-bit args
1115     * placed into 64-bit (long) arg slots, plus one 64-bit
1116     * (long) arg count, maintaining 16 byte alignment.
1117     */
1118     subq    $SYS_DROP, %rsp
1119     movb    $LWP_SYS, LWP_STATE(%r14)
1120     movq    %r15, %rdi
1121     movq    %rsp, %rsi
1122     call   syscall_entry
1124     /*
1125     * Fetch the arguments copied onto the kernel stack and put
1126     * them in the right registers to invoke a C-style syscall handler.
1127     * %rax contains the handler address.
1128     */
1129     movq    %rax, %rbx
1130     movl    0(%rsp), %edi
1131     movl    8(%rsp), %esi
1132     movl    0x10(%rsp), %edx
1133     movl    0x18(%rsp), %ecx
1134     movl    0x20(%rsp), %r8d
1135     movl    0x28(%rsp), %r9d
1137     call   *SY_CALLC(%rbx)
1139     movq    %rbp, %rsp /* pop the args */
1141     /*

```

```

1142     * amd64 syscall handlers -always- return a 64-bit value in %rax.
1143     * On the 32-bit kernel, the always return that value in %eax:%edx
1144     * as required by the 32-bit ABI.
1145     *
1146     * Simulate the same behaviour by unconditionally splitting the
1147     * return value in the same way.
1148     */
1149     movq    %rax, %r13
1150     shrq    $32, %r13 /* upper 32-bits into %edx */
1151     movl    %eax, %r12d /* lower 32-bits into %eax */
1153     /*
1154     * Optimistically assume that there's no post-syscall
1155     * work to do. (This is to avoid having to call syscall_mstate())
1156     * with interrupts disabled)
1157     */
1158     MSTATE_TRANSITION(LMS_SYSTEM, LMS_USER)
1160     /*
1161     * We must protect ourselves from being descheduled here;
1162     * If we were, and we ended up on another cpu, or another
1163     * lwp got int ahead of us, it could change the segment
1164     * registers without us noticing before we return to userland.
1165     *
1166     * This cli is undone in the tr_sysexit trampoline code.
1167     */
1168     cli
1169     CHECK_POSTSYS_NE(%r15, %r14, %ebx)
1170     jne    _full_syscall_postsys32
1171     SIMPLE_SYSCALL_POSTSYS(%r15, %r14, %bx)
1173     /*
1174     * To get back to userland, load up the 32-bit registers and
1175     * sysexit back where we came from.
1176     */
1178     /*
1179     * Interrupts will be turned on by the 'sti' executed just before
1180     * sysexit. The following ensures that restoring the user's rflags
1181     * doesn't enable interrupts too soon.
1182     */
1183     andq    $_BITNOT(PS_IE), REGOFF_RFL(%rsp)
1185     /*
1186     * (There's no point in loading up %edx because the sysexit
1187     * mechanism smashes it.)
1188     */
1189     movl    %r12d, %eax
1190     movl    REGOFF_RBX(%rsp), %ebx
1191     movl    REGOFF_RBP(%rsp), %ebp
1192     movl    REGOFF_RSI(%rsp), %esi
1193     movl    REGOFF_RDI(%rsp), %edi
1195     movl    REGOFF_RIP(%rsp), %edx /* sysexit: %edx -> %eip */
1196     pushq  REGOFF_RFL(%rsp)
1197     popfq
1198     movl    REGOFF_RSP(%rsp), %ecx /* sysexit: %ecx -> %esp */
1199     ALTENTRY(sys_sysenter_swapgs_sysexit)
1200     jmp    tr_sysexit
1154     swapgs
1155     sti
1156     sysexit
1201     SET_SIZE(sys_sysenter_swapgs_sysexit)
1206 #endif /* __lint */

```

```

1208 /*
1209 * This is the destination of the "int $T_SYSCALLINT" interrupt gate, used by
1210 * the generic i386 libc to do system calls. We do a small amount of setup
1211 * before jumping into the existing sys_syscall32 path.
1212 */
1213 #if defined(__lint)

1215 /*ARGSUSED*/
1216 void
1217 sys_syscall_int()
1218 {}

1220 #else /* __lint */

1222     ENTRY_NP(brand_sys_syscall_int)
1223     SWAPGS /* kernel gsbase */
1224     XPV_TRAP_POP
1225     call smap_enable
1226     BRAND_CALLBACK(BRAND_CB_INT91, BRAND_URET_FROM_INTR_STACK())
1227     jmp nopop_syscall_int

1229     ALTENTRY(sys_syscall_int)
1230     SWAPGS /* kernel gsbase */
1231     XPV_TRAP_POP
1232     call smap_enable

1234 nopop_syscall_int:
1235     movq %gs:CPU_THREAD, %r15
1236     movq T_STACK(%r15), %rsp
1237     movl %eax, %eax
1238     /*
1239     * Set t_post_sys on this thread to force ourselves out via the slow
1240     * path. It might be possible at some later date to optimize this out
1241     * and use a faster return mechanism.
1242     */
1243     movb $1, T_POST_SYS(%r15)
1244     CLEAN_CS
1245     jmp _syscall32_save
1246     /*
1247     * There should be no instructions between this label and SWAPGS/IRET
1248     * or we could end up breaking branded zone support. See the usage of
1249     * this label in lx_brand_int80_callback and snl_brand_int91_callback
1250     * for examples.
1251     *
1252     * We want to swapgs to maintain the invariant that all entries into
1253     * tr_iret_user are done on the user gsbase.
1254     */
1255     ALTENTRY(sys_sysint_swapgs_iret)
1256     SWAPGS
1257     jmp tr_iret_user
1258     SWAPGS /* user gsbase */
1259     IRET
1258     /*NOTREACHED*/
1259     SET_SIZE(sys_sysint_swapgs_iret)

```

unchanged portion omitted

new/usr/src/uts/i86pc/os/cpuid.c

1

```
*****
137103 Fri Apr 6 17:25:01 2018
new/usr/src/uts/i86pc/os/cpuid.c
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22 * Copyright (c) 2004, 2010, Oracle and/or its affiliates. All rights reserved.
23 * Copyright (c) 2011, 2016 by Delphix. All rights reserved.
24 * Copyright 2013 Nexenta Systems, Inc. All rights reserved.
25 * Copyright 2014 Josef "Jeff" Sipek <jeffpc@josefsipek.net>
26 */
27 /*
28 * Copyright (c) 2010, Intel Corporation.
29 * All rights reserved.
30 */
31 /*
32 * Portions Copyright 2009 Advanced Micro Devices, Inc.
33 */
34 /*
35 * Copyright 2018 Joyent, Inc.
36 * Copyright 2017 Joyent, Inc.
37 */
38 * Various routines to handle identification
39 * and classification of x86 processors.
40 */

42 #include <sys/types.h>
43 #include <sys/archsystem.h>
44 #include <sys/x86_archext.h>
45 #include <sys/kmem.h>
46 #include <sys/system.h>
47 #include <sys/cmn_err.h>
48 #include <sys/sunddi.h>
49 #include <sys/sunndi.h>
50 #include <sys/cpuvar.h>
51 #include <sys/processor.h>
52 #include <sys/sysmacros.h>
53 #include <sys/pg.h>
54 #include <sys/fp.h>
55 #include <sys/controlregs.h>
56 #include <sys/bitmap.h>
57 #include <sys/auxv_386.h>
58 #include <sys/memnode.h>
```

new/usr/src/uts/i86pc/os/cpuid.c

2

```
59 #include <sys/pci_cfgspace.h>
60 #include <sys/comm_page.h>
61 #include <sys/mach_mmu.h>
62 #include <sys/tsc.h>

64 #ifdef __xpv
65 #include <sys/hypervisor.h>
66 #else
67 #include <sys/ontrap.h>
68 #endif

70 /*
71 * Pass 0 of cpuid feature analysis happens in locore. It contains special code
72 * to recognize Cyrix processors that are not cpuid-compliant, and to deal with
73 * them accordingly. For most modern processors, feature detection occurs here
74 * in pass 1.
75 *
76 * Pass 1 of cpuid feature analysis happens just at the beginning of mlsetup()
77 * for the boot CPU and does the basic analysis that the early kernel needs.
78 * x86_featureset is set based on the return value of cpuid_pass1() of the boot
79 * CPU.
80 *
81 * Pass 1 includes:
82 *
83 *   o Determining vendor/model/family/stepping and setting x86_type and
84 *     x86_vendor accordingly.
85 *   o Processing the feature flags returned by the cpuid instruction while
86 *     applying any workarounds or tricks for the specific processor.
87 *   o Mapping the feature flags into illumos feature bits (X86_*).
88 *   o Mapping the feature flags into Solaris feature bits (X86_*).
89 *   o Processing extended feature flags if supported by the processor,
90 *     again while applying specific processor knowledge.
91 *   o Determining the CMT characteristics of the system.
92 *
93 * Pass 1 is done on non-boot CPUs during their initialization and the results
94 * are used only as a meager attempt at ensuring that all processors within the
95 * system support the same features.
96 *
97 * Pass 2 of cpuid feature analysis happens just at the beginning
98 * of startup(). It just copies in and corrects the remainder
99 * of the cpuid data we depend on: standard cpuid functions that we didn't
100 * need for pass1 feature analysis, and extended cpuid functions beyond the
101 * simple feature processing done in pass1.
102 *
103 * Pass 3 of cpuid analysis is invoked after basic kernel services; in
104 * particular kernel memory allocation has been made available. It creates a
105 * readable brand string based on the data collected in the first two passes.
106 *
107 * Pass 4 of cpuid analysis is invoked after post_startup() when all
108 * the support infrastructure for various hardware features has been
109 * initialized. It determines which processor features will be reported
110 * to userland via the aux vector.
111 *
112 * All passes are executed on all CPUs, but only the boot CPU determines what
113 * features the kernel will use.
114 *
115 * Much of the worst junk in this file is for the support of processors
116 * that didn't really implement the cpuid instruction properly.
117 *
118 * NOTE: The accessor functions (cpuid_get*) are aware of, and ASSERT upon,
119 * the pass numbers. Accordingly, changes to the pass code may require changes
120 * to the accessor code.
121 */

122 uint_t x86_vendor = X86_VENDOR_IntelClone;
123 uint_t x86_type = X86_TYPE_OTHER;
```

```

124 uint_t x86_clflush_size = 0;

126 #if defined(__xpv)
127 int x86_use_pcid = 0;
128 int x86_use_invpcid = 0;
129 #else
130 int x86_use_pcid = -1;
131 int x86_use_invpcid = -1;
132 #endif

134 uint_t pentiumpro_bug4046376;

136 uchar_t x86_featureset[BT_SIZEOFMAP(NUM_X86_FEATURES)];

138 static char *x86_feature_names[NUM_X86_FEATURES] = {
139     "lgpg",
140     "tsc",
141     "msr",
142     "mtrr",
143     "pge",
144     "de",
145     "cmov",
146     "mmx",
147     "mca",
148     "pae",
149     "cv8",
150     "pat",
151     "sep",
152     "sse",
153     "sse2",
154     "htt",
155     "asysc",
156     "nx",
157     "sse3",
158     "cx16",
159     "cmp",
160     "tscp",
161     "mwait",
162     "sse4a",
163     "cpuid",
164     "ssse3",
165     "sse4_1",
166     "sse4_2",
167     "lgpg",
168     "clflush",
169     "64",
170     "aes",
171     "pclmulqdq",
172     "xsave",
173     "avx",
174     "vmx",
175     "svm",
176     "topoext",
177     "f16c",
178     "rdrand",
179     "x2apic",
180     "avx2",
181     "bmi1",
182     "bmi2",
183     "fma",
184     "smep",
185     "smap",
186     "adx",
187     "rdseed",
188     "mpx",
189     "avx512f",

```

```

190     "avx512dq",
191     "avx512pf",
192     "avx512er",
193     "avx512cd",
194     "avx512bw",
195     "avx512vl",
196     "avx512fma",
197     "avx512vbmi",
198     "avx512_vpopcntdq",
199     "avx512_4vnniw",
200     "avx512_4fmaps",
201     "xsaveopt",
202     "xsavec",
203     "xsaves",
204     "sha",
205     "umip",
206     "pku",
207     "ospke",
208     "pcid",
209     "invpcid",
210 };
    unchanged portion omitted

997 void
998 cpuid_pass1(cpu_t *cpu, uchar_t *featureset)
999 {
1000     uint32_t mask_ecx, mask_edx;
1001     struct cpuid_info *cpi;
1002     struct cpuid_regs *cp;
1003     int xcpuid;
1004     #if !defined(__xpv)
1005     extern int idle_cpu_prefer_mwait;
1006     #endif

1008     /*
1009     * Space statically allocated for BSP, ensure pointer is set
1010     */
1011     if (cpu->cpu_id == 0) {
1012         if (cpu->cpu_m.mcpu_cpi == NULL)
1013             cpu->cpu_m.mcpu_cpi = &cpuid_info0;
1014     }

1016     add_x86_feature(featureset, X86FSET_CPUID);

1018     cpi = cpu->cpu_m.mcpu_cpi;
1019     ASSERT(cpi != NULL);
1020     cp = &cpi->cpi_std[0];
1021     cp->cp_eax = 0;
1022     cpi->cpi_maxeax = __cpuid_insn(cp);
1023     {
1024         uint32_t *iptr = (uint32_t *)cpi->cpi_vendorstr;
1025         *iptr++ = cp->cp_ebx;
1026         *iptr++ = cp->cp_edx;
1027         *iptr++ = cp->cp_ecx;
1028         *(char *)&cpi->cpi_vendorstr[12] = '\0';
1029     }

1031     cpi->cpi_vendor = _cpuid_vendorstr_to_vendorcode(cpi->cpi_vendorstr);
1032     x86_vendor = cpi->cpi_vendor; /* for compatibility */

1034     /*
1035     * Limit the range in case of weird hardware
1036     */
1037     if (cpi->cpi_maxeax > CPI_MAXEAX_MAX)
1038         cpi->cpi_maxeax = CPI_MAXEAX_MAX;
1039     if (cpi->cpi_maxeax < 1)

```

```

1040         goto pass1_done;
1042     cp = &cpi->cpi_std[1];
1043     cp->cp_eax = 1;
1044     (void) __cpuid_insn(cp);
1046     /*
1047     * Extract identifying constants for easy access.
1048     */
1049     cpi->cpi_model = CPI_MODEL(cpi);
1050     cpi->cpi_family = CPI_FAMILY(cpi);
1052     if (cpi->cpi_family == 0xf)
1053         cpi->cpi_family += CPI_FAMILY_XTD(cpi);
1055     /*
1056     * Beware: AMD uses "extended model" iff base *FAMILY* == 0xf.
1057     * Intel, and presumably everyone else, uses model == 0xf, as
1058     * one would expect (max value means possible overflow). Sigh.
1059     */
1061     switch (cpi->cpi_vendor) {
1062     case X86_VENDOR_Intel:
1063         if (IS_EXTENDED_MODEL_INTEL(cpi))
1064             cpi->cpi_model += CPI_MODEL_XTD(cpi) << 4;
1065         break;
1066     case X86_VENDOR_AMD:
1067         if (CPI_FAMILY(cpi) == 0xf)
1068             cpi->cpi_model += CPI_MODEL_XTD(cpi) << 4;
1069         break;
1070     default:
1071         if (cpi->cpi_model == 0xf)
1072             cpi->cpi_model += CPI_MODEL_XTD(cpi) << 4;
1073         break;
1074     }
1076     cpi->cpi_step = CPI_STEP(cpi);
1077     cpi->cpi_brandid = CPI_BRANDID(cpi);
1079     /*
1080     * *default* assumptions:
1081     * - believe %edx feature word
1082     * - ignore %ecx feature word
1083     * - 32-bit virtual and physical addressing
1084     */
1085     mask_edx = 0xffffffff;
1086     mask_ecx = 0;
1088     cpi->cpi_pabits = cpi->cpi_vabits = 32;
1090     switch (cpi->cpi_vendor) {
1091     case X86_VENDOR_Intel:
1092         if (cpi->cpi_family == 5)
1093             x86_type = X86_TYPE_P5;
1094         else if (IS_LEGACY_P6(cpi)) {
1095             x86_type = X86_TYPE_P6;
1096             pentiumpro_bug4046376 = 1;
1097             /*
1098             * Clear the SEP bit when it was set erroneously
1099             */
1100             if (cpi->cpi_model < 3 && cpi->cpi_step < 3)
1101                 cp->cp_edx &= ~CPIUID_INTC_EDX_SEP;
1102         } else if (IS_NEW_F6(cpi) || cpi->cpi_family == 0xf) {
1103             x86_type = X86_TYPE_P4;
1104             /*
1105             * We don't currently depend on any of the %ecx

```

```

1106         * features until Prescott, so we'll only check
1107         * this from P4 onwards. We might want to revisit
1108         * that idea later.
1109         */
1110         mask_ecx = 0xffffffff;
1111     } else if (cpi->cpi_family > 0xf)
1112         mask_ecx = 0xffffffff;
1113     /*
1114     * We don't support MONITOR/MWAIT if leaf 5 is not available
1115     * to obtain the monitor linesize.
1116     */
1117     if (cpi->cpi_maxeax < 5)
1118         mask_ecx &= ~CPIUID_INTC_ECX_MON;
1119     break;
1120     case X86_VENDOR_IntelClone:
1121     default:
1122         break;
1123     case X86_VENDOR_AMD:
1124     #if defined(OPTERON_ERRATUM_108)
1125         if (cpi->cpi_family == 0xf && cpi->cpi_model == 0xe) {
1126             cp->cp_eax = (0xf0f & cp->cp_eax) | 0xc0;
1127             cpi->cpi_model = 0xc;
1128         } else
1129     #endif
1130         if (cpi->cpi_family == 5) {
1131             /*
1132             * AMD K5 and K6
1133             *
1134             * These CPUs have an incomplete implementation
1135             * of MCA/MCE which we mask away.
1136             */
1137             mask_edx &= ~(CPIUID_INTC_EDX_MCE | CPIUID_INTC_EDX_MCA);
1139             /*
1140             * Model 0 uses the wrong (APIC) bit
1141             * to indicate PGE. Fix it here.
1142             */
1143             if (cpi->cpi_model == 0) {
1144                 if (cp->cp_edx & 0x200) {
1145                     cp->cp_edx &= ~0x200;
1146                     cp->cp_edx |= CPIUID_INTC_EDX_PGE;
1147                 }
1148             }
1150             /*
1151             * Early models had problems w/ MMX; disable.
1152             */
1153             if (cpi->cpi_model < 6)
1154                 mask_edx &= ~CPIUID_INTC_EDX_MMX;
1155         }
1157     /*
1158     * For newer families, SSE3 and CX16, at least, are valid;
1159     * enable all
1160     */
1161     if (cpi->cpi_family >= 0xf)
1162         mask_ecx = 0xffffffff;
1163     /*
1164     * We don't support MONITOR/MWAIT if leaf 5 is not available
1165     * to obtain the monitor linesize.
1166     */
1167     if (cpi->cpi_maxeax < 5)
1168         mask_ecx &= ~CPIUID_INTC_ECX_MON;
1170     #if !defined(__xpv)
1171     /*

```

```

1172     * Do not use MONITOR/MWAIT to halt in the idle loop on any AMD
1173     * processors. AMD does not intend MWAIT to be used in the cpu
1174     * idle loop on current and future processors. 10h and future
1175     * AMD processors use more power in MWAIT than HLT.
1176     * Pre-family-10h Opterons do not have the MWAIT instruction.
1177     */
1178     idle_cpu_prefer_mwait = 0;
1179 #endif

1181     break;
1182 case X86_VENDOR_TM:
1183     /*
1184     * workaround the NT workaround in CMS 4.1
1185     */
1186     if (cpi->cpi_family == 5 && cpi->cpi_model == 4 &&
1187         (cpi->cpi_step == 2 || cpi->cpi_step == 3))
1188         cp->cp_edx |= CPUID_INTC_EDX_CX8;
1189     break;
1190 case X86_VENDOR_Centaur:
1191     /*
1192     * workaround the NT workarounds again
1193     */
1194     if (cpi->cpi_family == 6)
1195         cp->cp_edx |= CPUID_INTC_EDX_CX8;
1196     break;
1197 case X86_VENDOR_Cyrix:
1198     /*
1199     * We rely heavily on the probing in locore
1200     * to actually figure out what parts, if any,
1201     * of the Cyrix cpuid instruction to believe.
1202     */
1203     switch (x86_type) {
1204     case X86_TYPE_CYRIX_486:
1205         mask_edx = 0;
1206         break;
1207     case X86_TYPE_CYRIX_6x86:
1208         mask_edx = 0;
1209         break;
1210     case X86_TYPE_CYRIX_6x86L:
1211         mask_edx =
1212             CPUID_INTC_EDX_DE |
1213             CPUID_INTC_EDX_CX8;
1214         break;
1215     case X86_TYPE_CYRIX_6x86MX:
1216         mask_edx =
1217             CPUID_INTC_EDX_DE |
1218             CPUID_INTC_EDX_MSR |
1219             CPUID_INTC_EDX_CX8 |
1220             CPUID_INTC_EDX_PGE |
1221             CPUID_INTC_EDX_CMOV |
1222             CPUID_INTC_EDX_MMX;
1223         break;
1224     case X86_TYPE_CYRIX_GXm:
1225         mask_edx =
1226             CPUID_INTC_EDX_MSR |
1227             CPUID_INTC_EDX_CX8 |
1228             CPUID_INTC_EDX_CMOV |
1229             CPUID_INTC_EDX_MMX;
1230         break;
1231     case X86_TYPE_CYRIX_MediaGX:
1232         break;
1233     case X86_TYPE_CYRIX_MII:
1234     case X86_TYPE_VIA_CYRIX_III:
1235         mask_edx =
1236             CPUID_INTC_EDX_DE |
1237             CPUID_INTC_EDX_TSC |

```

```

1238             CPUID_INTC_EDX_MSR |
1239             CPUID_INTC_EDX_CX8 |
1240             CPUID_INTC_EDX_PGE |
1241             CPUID_INTC_EDX_CMOV |
1242             CPUID_INTC_EDX_MMX;
1243         break;
1244     default:
1245         break;
1246     }
1247     break;
1248 }

1250 #if defined(__xpv)
1251 /*
1252 * Do not support MONITOR/MWAIT under a hypervisor
1253 */
1254 mask_ecx &= ~CPUID_INTC_ECX_MON;
1255 /*
1256 * Do not support XSAVE under a hypervisor for now
1257 */
1258 xsave_force_disable = B_TRUE;

1260 #endif /* __xpv */

1262 if (xsave_force_disable) {
1263     mask_ecx &= ~CPUID_INTC_ECX_XSAVE;
1264     mask_ecx &= ~CPUID_INTC_ECX_AVX;
1265     mask_ecx &= ~CPUID_INTC_ECX_F16C;
1266     mask_ecx &= ~CPUID_INTC_ECX_FMA;
1267 }

1269 /*
1270 * Now we've figured out the masks that determine
1271 * which bits we choose to believe, apply the masks
1272 * to the feature words, then map the kernel's view
1273 * of these feature words into its feature word.
1274 */
1275 cp->cp_edx &= mask_edx;
1276 cp->cp_ecx &= mask_ecx;

1278 /*
1279 * apply any platform restrictions (we don't call this
1280 * immediately after __cpuid_insn here, because we need the
1281 * workarounds applied above first)
1282 */
1283 platform_cpuid_mangle(cpi->cpi_vendor, 1, cp);

1285 /*
1286 * In addition to ecx and edx, Intel is storing a bunch of instruction
1287 * set extensions in leaf 7's ebx, ecx, and edx.
1288 */
1289 if (cpi->cpi_vendor == X86_VENDOR_Intel && cpi->cpi_maxeax >= 7) {
1290     struct cpuid_regs *ecp;
1291     ecp = &cpi->cpi_std[7];
1292     ecp->cp_eax = 7;
1293     ecp->cp_ecx = 0;
1294     (void) __cpuid_insn(ecp);
1295     /*
1296     * If XSAVE has been disabled, just ignore all of the
1297     * extended-save-area dependent flags here.
1298     */
1299     if (xsave_force_disable) {
1300         ecp->cp_ebx &= ~CPUID_INTC_EBX_7_0_BMI1;
1301         ecp->cp_ebx &= ~CPUID_INTC_EBX_7_0_BMI2;
1302         ecp->cp_ebx &= ~CPUID_INTC_EBX_7_0_AVX2;
1303         ecp->cp_ebx &= ~CPUID_INTC_EBX_7_0_MPX;

```

```

1304         ecp->cp_ebx &= ~CPUID_INTC_EBX_7_0_ALL_AVX512;
1305         ecp->cp_ecx &= ~CPUID_INTC_ECX_7_0_ALL_AVX512;
1306         ecp->cp_edx &= ~CPUID_INTC_EDX_7_0_ALL_AVX512;
1307     }
1309     if (ecp->cp_ebx & CPUID_INTC_EBX_7_0_SMEP)
1310         add_x86_feature(featureset, X86FSET_SMEP);
1312     if (ecp->cp_ebx & CPUID_INTC_EBX_7_0_INVPCID) {
1313         add_x86_feature(featureset, X86FSET_INVPCID);
1314     }
1316     /*
1317     * We check disable_smmap here in addition to in_startup_smmap()
1318     * to ensure CPUs that aren't the boot CPU don't accidentally
1319     * include it in the feature set and thus generate a mismatched
1320     * x86 feature set across CPUs. Note that at this time we only
1321     * enable SMAP for the 64-bit kernel.
1322     */
1323     #if defined(__amd64)
1324     if (ecp->cp_ebx & CPUID_INTC_EBX_7_0_SMAP &&
1325         disable_smmap == 0)
1326         add_x86_feature(featureset, X86FSET_SMAP);
1327     #endif
1328     if (ecp->cp_ebx & CPUID_INTC_EBX_7_0_MMX)
1329         add_x86_feature(featureset, X86FSET_MMX);
1331     if (ecp->cp_ebx & CPUID_INTC_EBX_7_0_RDSEED)
1332         add_x86_feature(featureset, X86FSET_RDSEED);
1334     if (ecp->cp_ebx & CPUID_INTC_EBX_7_0_ADX)
1335         add_x86_feature(featureset, X86FSET_ADX);
1336 }
1338 /*
1339 * fold in overrides from the "eeprom" mechanism
1340 */
1341 cp->cp_edx |= cpuid_feature_edx_include;
1342 cp->cp_edx &= ~cpuid_feature_edx_exclude;
1344 cp->cp_ecx |= cpuid_feature_ecx_include;
1345 cp->cp_ecx &= ~cpuid_feature_ecx_exclude;
1347 if (cp->cp_edx & CPUID_INTC_EDX_PSE) {
1348     add_x86_feature(featureset, X86FSET_LARGEPAGE);
1349 }
1350 if (cp->cp_edx & CPUID_INTC_EDX_TSC) {
1351     add_x86_feature(featureset, X86FSET_TSC);
1352 }
1353 if (cp->cp_edx & CPUID_INTC_EDX_MSR) {
1354     add_x86_feature(featureset, X86FSET_MSR);
1355 }
1356 if (cp->cp_edx & CPUID_INTC_EDX_MTRR) {
1357     add_x86_feature(featureset, X86FSET_MTRR);
1358 }
1359 if (cp->cp_edx & CPUID_INTC_EDX_PGE) {
1360     add_x86_feature(featureset, X86FSET_PGE);
1361 }
1362 if (cp->cp_edx & CPUID_INTC_EDX_CMOV) {
1363     add_x86_feature(featureset, X86FSET_CMOV);
1364 }
1365 if (cp->cp_edx & CPUID_INTC_EDX_MMX) {
1366     add_x86_feature(featureset, X86FSET_MMX);
1367 }
1368 if ((cp->cp_edx & CPUID_INTC_EDX_MCE) != 0 &&
1369     (cp->cp_edx & CPUID_INTC_EDX_MCA) != 0) {

```

```

1370         add_x86_feature(featureset, X86FSET_MCA);
1371     }
1372     if (cp->cp_edx & CPUID_INTC_EDX_PAE) {
1373         add_x86_feature(featureset, X86FSET_PAE);
1374     }
1375     if (cp->cp_edx & CPUID_INTC_EDX_CX8) {
1376         add_x86_feature(featureset, X86FSET_CX8);
1377     }
1378     if (cp->cp_ecx & CPUID_INTC_ECX_CX16) {
1379         add_x86_feature(featureset, X86FSET_CX16);
1380     }
1381     if (cp->cp_edx & CPUID_INTC_EDX_PAT) {
1382         add_x86_feature(featureset, X86FSET_PAT);
1383     }
1384     if (cp->cp_edx & CPUID_INTC_EDX_SEP) {
1385         add_x86_feature(featureset, X86FSET_SEP);
1386     }
1387     if (cp->cp_edx & CPUID_INTC_EDX_FXSR) {
1388         /*
1389         * In our implementation, fxsave/fxrstor
1390         * are prerequisites before we'll even
1391         * try and do SSE things.
1392         */
1393         if (cp->cp_edx & CPUID_INTC_EDX_SSE) {
1394             add_x86_feature(featureset, X86FSET_SSE);
1395         }
1396         if (cp->cp_edx & CPUID_INTC_EDX_SSE2) {
1397             add_x86_feature(featureset, X86FSET_SSE2);
1398         }
1399         if (cp->cp_ecx & CPUID_INTC_ECX_SSE3) {
1400             add_x86_feature(featureset, X86FSET_SSE3);
1401         }
1402         if (cp->cp_ecx & CPUID_INTC_ECX_SSSE3) {
1403             add_x86_feature(featureset, X86FSET_SSSE3);
1404         }
1405         if (cp->cp_ecx & CPUID_INTC_ECX_SSE4_1) {
1406             add_x86_feature(featureset, X86FSET_SSE4_1);
1407         }
1408         if (cp->cp_ecx & CPUID_INTC_ECX_SSE4_2) {
1409             add_x86_feature(featureset, X86FSET_SSE4_2);
1410         }
1411         if (cp->cp_ecx & CPUID_INTC_ECX_AES) {
1412             add_x86_feature(featureset, X86FSET_AES);
1413         }
1414         if (cp->cp_ecx & CPUID_INTC_ECX_PCLMULQDQ) {
1415             add_x86_feature(featureset, X86FSET_PCLMULQDQ);
1416         }
1417     }
1418     if (cpi->cpi_std[7].cp_ebx & CPUID_INTC_EBX_7_0_SHA)
1419         add_x86_feature(featureset, X86FSET_SHA);
1421     if (cpi->cpi_std[7].cp_ecx & CPUID_INTC_ECX_7_0_UMIP)
1422         add_x86_feature(featureset, X86FSET_UMIP);
1423     if (cpi->cpi_std[7].cp_ecx & CPUID_INTC_ECX_7_0_PKU)
1424         add_x86_feature(featureset, X86FSET_PKU);
1425     if (cpi->cpi_std[7].cp_ecx & CPUID_INTC_ECX_7_0_OSPKE)
1426         add_x86_feature(featureset, X86FSET_OSPKE);
1428     if (cp->cp_ecx & CPUID_INTC_ECX_XSAVE) {
1429         add_x86_feature(featureset, X86FSET_XSAVE);
1430     }
1431     /* We only test AVX & AVX512 when there is XSAVE */
1433     if (cp->cp_ecx & CPUID_INTC_ECX_AVX) {
1434         add_x86_feature(featureset,
1435             X86FSET_AVX);

```

```

1437      /*
1438      * Intel says we can't check these without also
1439      * checking AVX.
1440      */
1441      if (cp->cp_ecx & CPUID_INTC_ECX_F16C)
1442          add_x86_feature(featureset,
1443                          X86FSET_F16C);
1444
1445      if (cp->cp_ecx & CPUID_INTC_ECX_FMA)
1446          add_x86_feature(featureset,
1447                          X86FSET_FMA);
1448
1449      if (cpi->cpi_std[7].cp_ebx &
1450          CPUID_INTC_EBX_7_0_BMI1)
1451          add_x86_feature(featureset,
1452                          X86FSET_BMI1);
1453
1454      if (cpi->cpi_std[7].cp_ebx &
1455          CPUID_INTC_EBX_7_0_BMI2)
1456          add_x86_feature(featureset,
1457                          X86FSET_BMI2);
1458
1459      if (cpi->cpi_std[7].cp_ebx &
1460          CPUID_INTC_EBX_7_0_AVX2)
1461          add_x86_feature(featureset,
1462                          X86FSET_AVX2);
1463    }
1464
1465    if (cpi->cpi_std[7].cp_ebx &
1466        CPUID_INTC_EBX_7_0_AVX512F) {
1467        add_x86_feature(featureset, X86FSET_AVX512F);
1468
1469        if (cpi->cpi_std[7].cp_ebx &
1470            CPUID_INTC_EBX_7_0_AVX512DQ)
1471            add_x86_feature(featureset,
1472                            X86FSET_AVX512DQ);
1473        if (cpi->cpi_std[7].cp_ebx &
1474            CPUID_INTC_EBX_7_0_AVX512IFMA)
1475            add_x86_feature(featureset,
1476                            X86FSET_AVX512FMA);
1477        if (cpi->cpi_std[7].cp_ebx &
1478            CPUID_INTC_EBX_7_0_AVX512PF)
1479            add_x86_feature(featureset,
1480                            X86FSET_AVX512PF);
1481        if (cpi->cpi_std[7].cp_ebx &
1482            CPUID_INTC_EBX_7_0_AVX512ER)
1483            add_x86_feature(featureset,
1484                            X86FSET_AVX512ER);
1485        if (cpi->cpi_std[7].cp_ebx &
1486            CPUID_INTC_EBX_7_0_AVX512CD)
1487            add_x86_feature(featureset,
1488                            X86FSET_AVX512CD);
1489        if (cpi->cpi_std[7].cp_ebx &
1490            CPUID_INTC_EBX_7_0_AVX512BW)
1491            add_x86_feature(featureset,
1492                            X86FSET_AVX512BW);
1493        if (cpi->cpi_std[7].cp_ebx &
1494            CPUID_INTC_EBX_7_0_AVX512VL)
1495            add_x86_feature(featureset,
1496                            X86FSET_AVX512VL);
1497
1498        if (cpi->cpi_std[7].cp_ecx &
1499            CPUID_INTC_ECX_7_0_AVX512VBMI)
1500            add_x86_feature(featureset,
1501                            X86FSET_AVX512VBMI);

```

```

1502      if (cpi->cpi_std[7].cp_ecx &
1503          CPUID_INTC_ECX_7_0_AVX512VPOPCDQ)
1504          add_x86_feature(featureset,
1505                          X86FSET_AVX512VPOPCDQ);
1506
1507      if (cpi->cpi_std[7].cp_edx &
1508          CPUID_INTC_EDX_7_0_AVX512NNIW)
1509          add_x86_feature(featureset,
1510                          X86FSET_AVX512NNIW);
1511      if (cpi->cpi_std[7].cp_edx &
1512          CPUID_INTC_EDX_7_0_AVX512FMAPS)
1513          add_x86_feature(featureset,
1514                          X86FSET_AVX512FMAPS);
1515    }
1516  }
1517 }
1518
1519 if (cpi->cpi_vendor == X86_VENDOR_Intel) {
1520     if (cp->cp_ecx & CPUID_INTC_ECX_PCID) {
1521         add_x86_feature(featureset, X86FSET_PCID);
1522     }
1523 }
1524
1525 if (cp->cp_ecx & CPUID_INTC_ECX_X2APIC) {
1526     add_x86_feature(featureset, X86FSET_X2APIC);
1527 }
1528 if (cp->cp_edx & CPUID_INTC_EDX_DE) {
1529     add_x86_feature(featureset, X86FSET_DE);
1530 }
1531 #if !defined(__xpv)
1532 if (cp->cp_ecx & CPUID_INTC_ECX_MON) {
1533
1534     /*
1535     * We require the CLFLUSH instruction for erratum workaround
1536     * to use MONITOR/MWAIT.
1537     */
1538     if (cp->cp_edx & CPUID_INTC_EDX_CLFSH) {
1539         cpi->cpi_mwait.support |= MWAIT_SUPPORT;
1540         add_x86_feature(featureset, X86FSET_MWAIT);
1541     } else {
1542         extern int idle_cpu_assert_cflush_monitor;
1543
1544         /*
1545         * All processors we are aware of which have
1546         * MONITOR/MWAIT also have CLFLUSH.
1547         */
1548         if (idle_cpu_assert_cflush_monitor) {
1549             ASSERT((cp->cp_ecx & CPUID_INTC_ECX_MON) &&
1550                 (cp->cp_edx & CPUID_INTC_EDX_CLFSH));
1551         }
1552     }
1553 }
1554 #endif /* __xpv */
1555
1556 if (cp->cp_ecx & CPUID_INTC_ECX_VMX) {
1557     add_x86_feature(featureset, X86FSET_VMX);
1558 }
1559
1560 if (cp->cp_ecx & CPUID_INTC_ECX_RDRAND)
1561     add_x86_feature(featureset, X86FSET_RDRAND);
1562
1563 /*
1564 * Only need it first time, rest of the cpus would follow suit.
1565 * we only capture this for the bootcpu.
1566 */
1567 if (cp->cp_edx & CPUID_INTC_EDX_CLFSH) {

```

```

1568         add_x86_feature(featureset, X86FSET_CLFSH);
1569         x86_clflush_size = (BITX(cp->cp_ebx, 15, 8) * 8);
1570     }
1571     if (is_x86_feature(featureset, X86FSET_PAE))
1572         cpi->cpi_pabits = 36;

1574     /*
1575     * Hyperthreading configuration is slightly tricky on Intel
1576     * and pure clones, and even trickier on AMD.
1577     *
1578     * (AMD chose to set the HTT bit on their CMP processors,
1579     * even though they're not actually hyperthreaded. Thus it
1580     * takes a bit more work to figure out what's really going
1581     * on ... see the handling of the CMP_LGCY bit below)
1582     */
1583     if (cp->cp_edx & CPUID_INTC_EDX_HTT) {
1584         cpi->cpi_ncpu_per_chip = CPI_CPU_COUNT(cpi);
1585         if (cpi->cpi_ncpu_per_chip > 1)
1586             add_x86_feature(featureset, X86FSET_HTT);
1587     } else {
1588         cpi->cpi_ncpu_per_chip = 1;
1589     }

1591     if (cpi->cpi_vendor == X86_VENDOR_Intel && cpi->cpi_maxeax >= 0xD &&
1592         !xsave_force_disable) {
1593         struct cpuid_regs r, *ecp;

1595         ecp = &r;
1596         ecp->cp_eax = 0xD;
1597         ecp->cp_ecx = 1;
1598         ecp->cp_edx = ecp->cp_ebx = 0;
1599         (void) __cpuid_insn(ecp);

1601         if (ecp->cp_eax & CPUID_INTC_EAX_D_1_XSAVEOPT)
1602             add_x86_feature(featureset, X86FSET_XSAVEOPT);
1603         if (ecp->cp_eax & CPUID_INTC_EAX_D_1_XSAVEC)
1604             add_x86_feature(featureset, X86FSET_XSAVEC);
1605         if (ecp->cp_eax & CPUID_INTC_EAX_D_1_XSAVES)
1606             add_x86_feature(featureset, X86FSET_XSAVES);
1607     }

1609     /*
1610     * Work on the "extended" feature information, doing
1611     * some basic initialization for cpuid_pass2()
1612     */
1613     xcpuid = 0;
1614     switch (cpi->cpi_vendor) {
1615     case X86_VENDOR_Intel:
1616         /*
1617         * On KVM we know we will have proper support for extended
1618         * cpuid.
1619         */
1620         if ((IS_NEW_F6(cpi) || cpi->cpi_family >= 0xf ||
1621             (get_hwenv() == HW_KVM && cpi->cpi_family == 6 &&
1622              (cpi->cpi_model == 6 || cpi->cpi_model == 2)))
1623             xcpuid++;
1624         break;
1625     case X86_VENDOR_AMD:
1626         if (cpi->cpi_family > 5 ||
1627             (cpi->cpi_family == 5 && cpi->cpi_model >= 1))
1628             xcpuid++;
1629         break;
1630     case X86_VENDOR_Cyrix:
1631         /*
1632         * Only these Cyrix CPUs are -known- to support
1633         * extended cpuid operations.

```

```

1634         /*
1635         if (x86_type == X86_TYPE_VIA_CYRIX_III ||
1636             x86_type == X86_TYPE_CYRIX_GXm)
1637             xcpuid++;
1638         break;
1639     case X86_VENDOR_Centaur:
1640     case X86_VENDOR_TM:
1641     default:
1642         xcpuid++;
1643         break;
1644     }

1646     if (xcpuid) {
1647         cp = &cpi->cpi_extd[0];
1648         cp->cp_eax = 0x80000000;
1649         cpi->cpi_xmaxeax = __cpuid_insn(cp);
1650     }

1652     if (cpi->cpi_xmaxeax & 0x80000000) {

1654         if (cpi->cpi_xmaxeax > CPI_XMAXEAX_MAX)
1655             cpi->cpi_xmaxeax = CPI_XMAXEAX_MAX;

1657         switch (cpi->cpi_vendor) {
1658         case X86_VENDOR_Intel:
1659         case X86_VENDOR_AMD:
1660             if (cpi->cpi_xmaxeax < 0x80000001)
1661                 break;
1662             cp = &cpi->cpi_extd[1];
1663             cp->cp_eax = 0x80000001;
1664             (void) __cpuid_insn(cp);

1666             if (cpi->cpi_vendor == X86_VENDOR_AMD &&
1667                 cpi->cpi_family == 5 &&
1668                 cpi->cpi_model == 6 &&
1669                 cpi->cpi_step == 6) {
1670                 /*
1671                 * K6 model 6 uses bit 10 to indicate SYSC
1672                 * Later models use bit 11. Fix it here.
1673                 */
1674                 if (cp->cp_edx & 0x400) {
1675                     cp->cp_edx &= ~0x400;
1676                     cp->cp_edx |= CPUID_AMD_EDX_SYSC;
1677                 }
1678             }

1680             platform_cpuid_mangle(cpi->cpi_vendor, 0x80000001, cp);

1682             /*
1683             * Compute the additions to the kernel's feature word.
1684             */
1685             if (cp->cp_edx & CPUID_AMD_EDX_NX) {
1686                 add_x86_feature(featureset, X86FSET_NX);
1687             }

1689             /*
1690             * Regardless whether or not we boot 64-bit,
1691             * we should have a way to identify whether
1692             * the CPU is capable of running 64-bit.
1693             */
1694             if (cp->cp_edx & CPUID_AMD_EDX_LM) {
1695                 add_x86_feature(featureset, X86FSET_64);
1696             }

1698 #if defined(__amd64)
1699             /* 1 GB large page - enable only for 64 bit kernel */

```

```

1700     if (cp->cp_edx & CPUID_AMD_EDX_1GPG) {
1701         add_x86_feature(featureset, X86FSET_1GPG);
1702     }
1703 #endif

1705     if ((cpi->cpi_vendor == X86_VENDOR_AMD) &&
1706         (cpi->cpi_std[1].cp_edx & CPUID_INTC_EDX_FXSR) &&
1707         (cp->cp_ecx & CPUID_AMD_ECX_SSE4A)) {
1708         add_x86_feature(featureset, X86FSET_SSE4A);
1709     }

1711     /*
1712     * If both the HTT and CMP_LGCY bits are set,
1713     * then we're not actually HyperThreaded.  Read
1714     * "AMD CPUID Specification" for more details.
1715     */
1716     if (cpi->cpi_vendor == X86_VENDOR_AMD &&
1717         is_x86_feature(featureset, X86FSET_HTT) &&
1718         (cp->cp_ecx & CPUID_AMD_ECX_CMP_LGCY)) {
1719         remove_x86_feature(featureset, X86FSET_HTT);
1720         add_x86_feature(featureset, X86FSET_CMP);
1721     }
1722 #if defined(__amd64)
1723     /*
1724     * It's really tricky to support syscall/sysret in
1725     * the i386 kernel; we rely on sysenter/sysexit
1726     * instead.  In the amd64 kernel, things are -way-
1727     * better.
1728     */
1729     if (cp->cp_edx & CPUID_AMD_EDX_SYSC) {
1730         add_x86_feature(featureset, X86FSET_ASYSC);
1731     }

1733     /*
1734     * While we're thinking about system calls, note
1735     * that AMD processors don't support sysenter
1736     * in long mode at all, so don't try to program them.
1737     */
1738     if (x86_vendor == X86_VENDOR_AMD) {
1739         remove_x86_feature(featureset, X86FSET_SEP);
1740     }
1741 #endif

1742     if (cp->cp_edx & CPUID_AMD_EDX_TSCP) {
1743         add_x86_feature(featureset, X86FSET_TSCP);
1744     }

1746     if (cp->cp_ecx & CPUID_AMD_ECX_SVM) {
1747         add_x86_feature(featureset, X86FSET_SVM);
1748     }

1750     if (cp->cp_ecx & CPUID_AMD_ECX_TOPOEXT) {
1751         add_x86_feature(featureset, X86FSET_TOPOEXT);
1752     }
1753     break;
1754 default:
1755     break;
1756 }

1758 /*
1759 * Get CPUID data about processor cores and hyperthreads.
1760 */
1761 switch (cpi->cpi_vendor) {
1762 case X86_VENDOR_Intel:
1763     if (cpi->cpi_maxeax >= 4) {
1764         cp = &cpi->cpi_std[4];
1765         cp->cp_eax = 4;

```

```

1766         cp->cp_ecx = 0;
1767         (void) __cpuid_insn(cp);
1768         platform_cpuid_mangle(cpi->cpi_vendor, 4, cp);
1769     }
1770     /*FALLTHROUGH*/
1771 case X86_VENDOR_AMD:
1772     if (cpi->cpi_xmaxeax < 0x80000008)
1773         break;
1774     cp = &cpi->cpi_extd[8];
1775     cp->cp_eax = 0x80000008;
1776     (void) __cpuid_insn(cp);
1777     platform_cpuid_mangle(cpi->cpi_vendor, 0x80000008, cp);

1779     /*
1780     * Virtual and physical address limits from
1781     * cpuid override previously guessed values.
1782     */
1783     cpi->cpi_pabits = BITX(cp->cp_eax, 7, 0);
1784     cpi->cpi_vabits = BITX(cp->cp_eax, 15, 8);
1785     break;
1786 default:
1787     break;
1788 }

1790 /*
1791 * Derive the number of cores per chip
1792 */
1793 switch (cpi->cpi_vendor) {
1794 case X86_VENDOR_Intel:
1795     if (cpi->cpi_maxeax < 4) {
1796         cpi->cpi_ncore_per_chip = 1;
1797         break;
1798     } else {
1799         cpi->cpi_ncore_per_chip =
1800             BITX((cpi->cpi_std[4].cp_eax, 31, 26) + 1;
1801     }
1802     break;
1803 case X86_VENDOR_AMD:
1804     if (cpi->cpi_xmaxeax < 0x80000008) {
1805         cpi->cpi_ncore_per_chip = 1;
1806         break;
1807     } else {
1808         /*
1809         * On family 0xf cpuid fn 2 ECX[7:0] "NC" is
1810         * 1 less than the number of physical cores on
1811         * the chip.  In family 0x10 this value can
1812         * be affected by "downcoring" - it reflects
1813         * 1 less than the number of cores actually
1814         * enabled on this node.
1815         */
1816         cpi->cpi_ncore_per_chip =
1817             BITX((cpi->cpi_extd[8].cp_ecx, 7, 0) + 1;
1818     }
1819     break;
1820 default:
1821     cpi->cpi_ncore_per_chip = 1;
1822     break;
1823 }

1825 /*
1826 * Get CPUID data about TSC Invariance in Deep C-State.
1827 */
1828 switch (cpi->cpi_vendor) {
1829 case X86_VENDOR_Intel:
1830     if (cpi->cpi_maxeax >= 7) {
1831         cp = &cpi->cpi_extd[7];

```



```

1832         cp->cp_eax = 0x80000007;
1833         cp->cp_ecx = 0;
1834         (void) __cpuid_insn(cp);
1835     }
1836     }
1837     default:
1838         break;
1839     }
1840 } else {
1841     cpi->cpi_ncore_per_chip = 1;
1842 }

1844 /*
1845  * If more than one core, then this processor is CMP.
1846  */
1847 if (cpi->cpi_ncore_per_chip > 1) {
1848     add_x86_feature(featureset, X86FSET_CMP);
1849 }

1851 /*
1852  * If the number of cores is the same as the number
1853  * of CPUs, then we cannot have HyperThreading.
1854  */
1855 if (cpi->cpi_ncpu_per_chip == cpi->cpi_ncore_per_chip) {
1856     remove_x86_feature(featureset, X86FSET_HTT);
1857 }

1859 cpi->cpi_apicid = CPI_APIC_ID(cpi);
1860 cpi->cpi_procnodes_per_pkg = 1;
1861 cpi->cpi_cores_per_compunit = 1;
1862 if (is_x86_feature(featureset, X86FSET_HTT) == B_FALSE &&
1863     is_x86_feature(featureset, X86FSET_CMP) == B_FALSE) {
1864     /*
1865      * Single-core single-threaded processors.
1866      */
1867     cpi->cpi_chipid = -1;
1868     cpi->cpi_clogid = 0;
1869     cpi->cpi_coreid = cpu->cpu_id;
1870     cpi->cpi_pkgcoreid = 0;
1871     if (cpi->cpi_vendor == X86_VENDOR_AMD)
1872         cpi->cpi_procnodeid = BITX(cpi->cpi_apicid, 3, 0);
1873     else
1874         cpi->cpi_procnodeid = cpi->cpi_chipid;
1875 } else if (cpi->cpi_ncpu_per_chip > 1) {
1876     if (cpi->cpi_vendor == X86_VENDOR_Intel)
1877         cpuid_intel_getids(cpu, featureset);
1878     else if (cpi->cpi_vendor == X86_VENDOR_AMD)
1879         cpuid_amd_getids(cpu);
1880     else {
1881         /*
1882          * All other processors are currently
1883          * assumed to have single cores.
1884          */
1885         cpi->cpi_coreid = cpi->cpi_chipid;
1886         cpi->cpi_pkgcoreid = 0;
1887         cpi->cpi_procnodeid = cpi->cpi_chipid;
1888         cpi->cpi_compunitid = cpi->cpi_chipid;
1889     }
1890 }

1892 /*
1893  * Synthesize chip "revision" and socket type
1894  */
1895 cpi->cpi_chiprev = _cpuid_chiprev(cpi->cpi_vendor, cpi->cpi_family,
1896     cpi->cpi_model, cpi->cpi_step);
1897 cpi->cpi_chiprevstr = _cpuid_chiprevstr(cpi->cpi_vendor,

```

```

1898         cpi->cpi_family, cpi->cpi_model, cpi->cpi_step);
1899     cpi->cpi_socket = _cpuid_skt(cpi->cpi_vendor, cpi->cpi_family,
1900         cpi->cpi_model, cpi->cpi_step);

1902     /*
1903     * While we're here, check for the AMD "Error Pointer Zero/Restore"
1904     * feature. This can be used to setup the FP save handlers
1905     * appropriately.
1906     */
1907     if (cpi->cpi_vendor == X86_VENDOR_AMD) {
1908         if (cpi->cpi_xmaxeax >= 0x80000008 &&
1909             cpi->cpi_extd[8].cp_ebx & CPUID_AMD_EBX_ERR_PTR_ZERO) {
1910             /* Special handling for AMD FP not necessary. */
1911             cpi->cpi_fp_aml_save = 0;
1912         } else {
1913             cpi->cpi_fp_aml_save = 1;
1914         }
1915     }

1917 pass1_done:
1918     cpi->cpi_pass = 1;
1919 }

    unchanged portion omitted

5017 void
5018 enable_pcid(void)
5019 {
5020     if (x86_use_pcid == -1)
5021         x86_use_pcid = is_x86_feature(x86_featureset, X86FSET_PCID);

5023     if (x86_use_invpcid == -1) {
5024         x86_use_invpcid = is_x86_feature(x86_featureset,
5025             X86FSET_INVPCID);
5026     }

5028     if (!x86_use_pcid)
5029         return;

5031     /*
5032     * Intel say that on setting PCIDE, it immediately starts using the PCID
5033     * bits; better make sure there's nothing there.
5034     */
5035     ASSERT((getcr3() & MMU_PAGEOFFSET) == PCID_NONE);

5037     setcr4(getcr4() | CR4_PCIDE);
5038 }

5040 /*
5041  * Setup necessary registers to enable XSAVE feature on this processor.
5042  * This function needs to be called early enough, so that no xsave/xrstor
5043  * ops will execute on the processor before the MSR's are properly set up.
5044  *
5045  * Current implementation has the following assumption:
5046  * - cpuid_pass1() is done, so that X86 features are known.
5047  * - fpu_probe() is done, so that fp_save_mech is chosen.
5048  */
5049 void
5050 xsave_setup_msr(cpu_t *cpu)
5051 {
5052     ASSERT(fp_save_mech == FP_XSAVE);
5053     ASSERT(is_x86_feature(x86_featureset, X86FSET_XSAVE));

5055     /* Enable OSXSAVE in CR4. */
5056     setcr4(getcr4() | CR4_OSXSAVE);
5057     /*
5058     * Update SW copy of ECX, so that /dev/cpu/self/cpuid will report

```

new/usr/src/uts/i86pc/os/cpuid.c

19

```
5059     * correct value.  
5060     */  
5061     cpu->cpu_m.mcpu_cpi->cpi_std[1].cp_ecx |= CPUID_INTC_ECX_OSXSAVE;  
5062     setup_xfem();  
5063 }  
_____unchanged_portion_omitted_____
```

```

*****
66820 Fri Apr 6 17:25:01 2018
new/usr/src/uts/i86pc/os/fakebop.c
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */

22 /*
23  * Copyright 2010 Sun Microsystems, Inc. All rights reserved.
24  * Use is subject to license terms.
25  *
26  * Copyright (c) 2010, Intel Corporation.
27  * All rights reserved.
28  *
29  * Copyright 2018 Joyent, Inc. All rights reserved.
30  * Copyright 2013 Joyent, Inc. All rights reserved.
31  */

32 /*
33  * This file contains the functionality that mimics the boot operations
34  * on SPARC systems or the old boot.bin/multiboot programs on x86 systems.
35  * The x86 kernel now does everything on its own.
36  */

38 #include <sys/types.h>
39 #include <sys/bootconf.h>
40 #include <sys/bootsvcs.h>
41 #include <sys/bootinfo.h>
42 #include <sys/multiboot.h>
43 #include <sys/multiboot2.h>
44 #include <sys/multiboot2_impl.h>
45 #include <sys/bootvfs.h>
46 #include <sys/bootprops.h>
47 #include <sys/varargs.h>
48 #include <sys/param.h>
49 #include <sys/machparam.h>
50 #include <sys/machsystem.h>
51 #include <sys/archsystem.h>
52 #include <sys/boot_console.h>
53 #include <sys/cmn_err.h>
54 #include <sys/system.h>
55 #include <sys/promif.h>
56 #include <sys/archsystem.h>
57 #include <sys/x86_archext.h>
58 #include <sys/kobj.h>

```

```

59 #include <sys/privregs.h>
60 #include <sys/sysmacros.h>
61 #include <sys/ctype.h>
62 #include <sys/fastboot.h>
63 #ifdef __xpv
64 #include <sys/hypervisor.h>
65 #include <net/if.h>
66 #endif
67 #include <vm/kboot_mmu.h>
68 #include <vm/hat_pte.h>
69 #include <sys/kobj.h>
70 #include <sys/kobj_lex.h>
71 #include <sys/pci_cfgspace_impl.h>
72 #include <sys/fastboot_impl.h>
73 #include <sys/acpi/acconfig.h>
74 #include <sys/acpi/acpi.h>

76 static int have_console = 0; /* set once primitive console is initialized */
77 static char *boot_args = "";

79 /*
80  * Debugging macros
81  */
82 static uint_t kbm_debug = 0;
83 #define DBG_MSG(s) { if (kbm_debug) bop_printf(NULL, "%s", s); }
84 #define DBG(x) { if (kbm_debug) \
85     bop_printf(NULL, "%s is %" PRIx64 "\n", #x, (uint64_t)(x)); \
86 }

88 #define PUT_STRING(s) { \
89     char *cp; \
90     for (cp = (s); *cp; ++cp) \
91         bcons_putchar(*cp); \
92 }

94 bootops_t bootop; /* simple bootops we'll pass on to kernel */
95 struct bsys_mem bm;

97 /*
98  * Boot info from "glue" code in low memory. xbootp is used by:
99  * do_bop_phys_alloc(), do_bsys_alloc() and boot_prop_finish().
100 */
101 static struct xboot_info *xbootp;
102 static uintptr_t next_virt; /* next available virtual address */
103 static paddr_t next_phys; /* next available physical address from dboot */
104 static paddr_t high_phys = -(paddr_t)1; /* last used physical address */

106 /*
107  * buffer for vsnprintf for console I/O
108  */
109 #define BUFFERSIZE 512
110 static char buffer[BUFFERSIZE];

112 /*
113  * stuff to store/report/manipulate boot property settings.
114  */
115 typedef struct bootprop {
116     struct bootprop *bp_next;
117     char *bp_name;
118     uint_t bp_vlen;
119     char *bp_value;
120 } bootprop_t;
121 unchanged portion omitted
836 typedef int (*bios_func_t)(int, bios_regs_t *);

838 /*ARGSUSED*/

```

```
839 static void
840 do_bsys_doint(bootops_t *bop, int intnum, struct bop_regs *rp)
841 {
842     #if defined(__xpv)
843         prom_panic("unsupported call to BOP_DOINT()\n");
844     #else /* __xpv */
845         static int firsttime = 1;
846         bios_func_t bios_func = (bios_func_t)(void *) (uintptr_t)0x5000;
847         bios_regs_t br;
848
849         /*
850          * We're about to disable paging; we shouldn't be PCID enabled.
851          */
852         if (getcr4() & CR4_PCIDE)
853             prom_panic("do_bsys_doint() with PCID enabled\n");
854
855         /*
856          * The first time we do this, we have to copy the pre-packaged
857          * low memory bios call code image into place.
858          */
859         if (firsttime) {
860             extern char bios_image[];
861             extern uint32_t bios_size;
862
863             bcopy(bios_image, (void *)bios_func, bios_size);
864             firsttime = 0;
865         }
866
867         br.ax = rp->eax.word.ax;
868         br.bx = rp->ebx.word.bx;
869         br.cx = rp->ecx.word.cx;
870         br.dx = rp->edx.word.dx;
871         br.bp = rp->ebp.word.bp;
872         br.si = rp->esi.word.si;
873         br.di = rp->edi.word.di;
874         br.ds = rp->ds;
875         br.es = rp->es;
876
877         DBG_MSG("Doing BIOS call...");
878         DBG(br.ax);
879         DBG(br.bx);
880         DBG(br.dx);
881         rp->eflags = bios_func(intnum, &br);
882         DBG_MSG("done\n");
883
884         rp->eax.word.ax = br.ax;
885         rp->ebx.word.bx = br.bx;
886         rp->ecx.word.cx = br.cx;
887         rp->edx.word.dx = br.dx;
888         rp->ebp.word.bp = br.bp;
889         rp->esi.word.si = br.si;
890         rp->edi.word.di = br.di;
891         rp->ds = br.ds;
892         rp->es = br.es;
893     #endif /* __xpv */
894 }
_____unchanged_portion_omitted_____
```

```

*****
57338 Fri Apr 6 17:25:02 2018
new/usr/src/uts/i86pc/os/intr.c
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21
22 /*
23  * Copyright (c) 2004, 2010, Oracle and/or its affiliates. All rights reserved.
24  * Copyright (c) 2018 Joyent, Inc. All rights reserved.
25  * Copyright (c) 2012, Joyent, Inc. All rights reserved.
26 */
27
28 /*
29  * To understand the present state of interrupt handling on i86pc, we must
30  * first consider the history of interrupt controllers and our way of handling
31  * interrupts.
32  *
33  * History of Interrupt Controllers on i86pc
34  * -----
35  *
36  * Intel 8259 and 8259A
37  *
38  * The first interrupt controller that attained widespread use on i86pc was
39  * the Intel 8259(A) Programmable Interrupt Controller that first saw use with
40  * the 8086. It took up to 8 interrupt sources and combined them into one
41  * output wire. Up to 8 8259s could be slaved together providing up to 64 IRQs.
42  * With the switch to the 8259A, level mode interrupts became possible. For a
43  * long time on i86pc the 8259A was the only way to handle interrupts and it
44  * had its own set of quirks. The 8259A and its corresponding interval timer
45  * the 8254 are programmed using outb and inb instructions.
46  *
47  * Intel Advanced Programmable Interrupt Controller (APIC)
48  *
49  * Starting around the time of the introduction of the P6 family
50  * microarchitecture (i686) Intel introduced a new interrupt controller.
51  * Instead of having the series of slaved 8259A devices, Intel opted to outfit
52  * each processor with a Local APIC (lapic) and to outfit the system with at
53  * least one, but potentially more, I/O APICs (ioapic). The lapics and ioapics
54  * initially communicated over a dedicated bus, but this has since been
55  * replaced. Each physical core and even hyperthread currently contains its
56  * own local apic, which is not shared. There are a few exceptions for
57  * hyperthreads, but that does not usually concern us.
58  *
59  * Instead of talking directly to 8259 for status, sending End Of Interrupt

```

```

59 * (EOI), etc. a microprocessor now communicates directly to the lapic. This
60 * also allows for each microprocessor to be able to have independent controls.
61 * The programming method is different from the 8259. Consumers map the lapic
62 * registers into uncacheable memory to read and manipulate the state.
63 *
64 * The number of addressable interrupt vectors was increased to 256. However
65 * vectors 0-31 are reserved for the processor exception handling, leaving the
66 * remaining vectors for general use. In addition to hardware generated
67 * interrupts, the lapic provides a way for generating inter-processor
68 * interrupts (IPI) which are the basis for CPU cross calls and CPU pokes.
69 *
70 * AMD ended up implementing the Intel APIC architecture in lieu of their work
71 * with Cyrix.
72 *
73 * Intel x2apic
74 *
75 * The x2apic is an extension to the lapic which started showing up around the
76 * same time as the Sandy Bridge chipsets. It provides a new programming mode
77 * as well as new features. The goal of the x2apic is to solve a few problems
78 * with the previous generation of lapic and the x2apic is backwards compatible
79 * with the previous programming and model. The only downsides to using the
80 * backwards compatibility is that you are not able to take advantage of the new
81 * x2apic features.
82 *
83 * o The APIC ID is increased from an 8-bit value to a 32-bit value. This
84 * increases the maximum number of addressable physical processors beyond
85 * 256. This new ID is assembled in a similar manner as the information that
86 * is obtainable by the extended cpuid topology leaves.
87 *
88 * o A new means of generating IPIs was introduced.
89 *
90 * o Instead of memory mapping the registers, the x2apic only allows for
91 * programming it through a series of wrmsrs. This has important semantic
92 * side effects. Recall that the registers were previously all mapped to
93 * uncacheable memory which meant that all operations to the local apic were
94 * serializing instructions. With the switch to using wrmsrs this has been
95 * relaxed and these operations can no longer be assumed to be serializing
96 * instructions.
97 *
98 * Note for the rest of this we are only going to concern ourselves with the
99 * apic and x2apic which practically all of i86pc has been using now for
100 * quite some time.
101 *
102 * Interrupt Priority Levels
103 * -----
104 *
105 * On i86pc systems there are a total of fifteen interrupt priority levels
106 * (ipls) which range from 1-15. Level 0 is for normal processing and
107 * non-interrupt processing. To manipulate these values the family of spl
108 * functions (which date back to UNIX on the PDP-11) are used. Specifically,
109 * splr() to raise the priority level and splx() to lower it. One should not
110 * generally call setspl() directly.
111 *
112 * Both i86pc and the supported SPARC platforms honor the same conventions for
113 * the meaning behind these IPLs. The most important IPL is the platform's
114 * LOCK_LEVEL (0xa on i86pc). If a thread is above LOCK_LEVEL it _must_ not
115 * sleep on any synchronization object. The only allowed synchronization
116 * primitive is a mutex that has been specifically initialized to be a spin
117 * lock (see mutex_init(9F)). Another important level is DISP_LEVEL (0xb on
118 * i86pc). You must be at DISP_LEVEL if you want to control the dispatcher.
119 * The XC_HI_PIL is the highest level (0xf) and is used during cross-calls.
120 *
121 * Each interrupt that is registered in the system fires at a specific IPL.
122 * Generally most interrupts fire below LOCK_LEVEL.
123 *
124 * PSM Drivers

```

```

125 * -----
126 *
127 * We currently have three sets of PSM (platform specific module) drivers
128 * available. uppc, pcplusmp, and apix. uppc (uni-processor PC) is the original
129 * driver that interacts with the 8259A and 8254. In general, it is not used
130 * anymore given the prevalence of the apic.
131 *
132 * The system prefers to use the apix driver over the pcplusmp driver. The apix
133 * driver requires HW support for an x2apic. If there is no x2apic HW, apix
134 * will not be used. In general we prefer using the apix driver over the
135 * pcplusmp driver because it gives us much more flexibility with respect to
136 * interrupts. In the apix driver each local apic has its own independent set
137 * of interrupts, whereas the pcplusmp driver only has a single global set of
138 * interrupts. This is why pcplusmp only supports a finite number of interrupts
139 * per IPL -- generally 16, often less. The apix driver supports using either
140 * the x2apic or the local apic programming modes. The programming mode does not
141 * change the number of interrupts available, just the number of processors
142 * that we can address. For the apix driver, the x2apic mode is enabled if the
143 * system supports interrupt re-mapping, otherwise the module manages the
144 * x2apic in local mode.
145 *
146 * When there is no x2apic present, we default back to the pcplusmp PSM driver.
147 * In general, this is not problematic unless you have more than 256
148 * processors in the machine or you do not have enough interrupts available.
149 *
150 * Controlling Interrupt Generation on i86pc
151 * -----
152 *
153 * There are two different ways to manipulate which interrupts will be
154 * generated on i86pc. Each offers different degrees of control.
155 *
156 * The first is through the flags register (eflags and rflags on i386 and amd64
157 * respectively). The IF bit determines whether or not interrupts are enabled
158 * or disabled. This is manipulated in one of several ways. The most common way
159 * is through the cli and sti instructions. These clear the IF flag and set it,
160 * respectively, for the current processor. The other common way is through the
161 * use of the intr_clear and intr_restore functions.
162 *
163 * Assuming interrupts are not blocked by the IF flag, then the second form is
164 * through the Processor-Priority Register (PPR). The PPR is used to determine
165 * whether or not a pending interrupt should be delivered. If the ipl of the
166 * new interrupt is higher than the current value in the PPR, then the lpic
167 * will either deliver it immediately (if interrupts are not in progress) or it
168 * will deliver it once the current interrupt processing has issued an EOI. The
169 * highest unmasked interrupt will be the one delivered.
170 *
171 * The PPR register is based upon the max of the following two registers in the
172 * lpic, the TPR register (also known as CR8 on amd64) that can be used to
173 * mask interrupt levels, and the current vector. Because the pcplusmp module
174 * always sets TPR appropriately early in the do_interrupt path, we can usually
175 * just think that the PPR is the TPR. The pcplusmp module also issues an EOI
176 * once it has set the TPR, so higher priority interrupts can come in while
177 * we're servicing a lower priority interrupt.
178 *
179 * Handling Interrupts
180 * -----
181 *
182 * Interrupts can be broken down into three categories based on priority and
183 * source:
184 *
185 *   o High level interrupts
186 *   o Low level hardware interrupts
187 *   o Low level software interrupts
188 *
189 * High Level Interrupts
190 *

```

```

191 * High level interrupts encompasses both hardware-sourced and software-sourced
192 * interrupts. Examples of high level hardware interrupts include the serial
193 * console. High level software-sourced interrupts are still delivered through
194 * the local apic through IPIS. This is primarily cross calls.
195 *
196 * When a high level interrupt comes in, we will raise the SPL and then pin the
197 * current lwp to the processor. We will use its lwp, but our own interrupt
198 * stack and process the high level interrupt in-situ. These handlers are
199 * designed to be very short in nature and cannot go to sleep, only block on a
200 * spin lock. If the interrupt has a lot of work to do, it must generate a
201 * low-priority software interrupt that will be processed later.
202 *
203 * Low level hardware interrupts
204 *
205 * Low level hardware interrupts start off like their high-level cousins. The
206 * current CPU contains a number of kernel threads (kthread_t) that can be used
207 * to process low level interrupts. These are shared between both low level
208 * hardware and software interrupts. Note that while we run with our
209 * kthread_t, we borrow the pinned threads lwp_t until such a time as we hit a
210 * synchronization object. If we hit one and need to sleep, then the scheduler
211 * will instead create the rest of what we need.
212 *
213 * Low level software interrupts
214 *
215 * Low level software interrupts are handled in a similar way as hardware
216 * interrupts, but the notification vector is different. Each CPU has a bitmask
217 * of pending software interrupts. We can notify a CPU to process software
218 * interrupts through a specific trap vector as well as through several
219 * checks that are performed throughout the code. These checks will look at
220 * processing software interrupts as we lower our spl.
221 *
222 * We attempt to process the highest pending software interrupt that we can
223 * which is greater than our current IPL. If none currently exist, then we move
224 * on. We process a software interrupt in a similar fashion to a hardware
225 * interrupt.
226 *
227 * Traditional Interrupt Flow
228 * -----
229 *
230 * The following diagram tracks the flow of the traditional uppc and pcplusmp
231 * interrupt handlers. The apix driver has its own version of do_interrupt().
232 * We come into the interrupt handler with all interrupts masked by the IF
233 * flag. This is because we set up the handler using an interrupt-gate, which
234 * is defined architecturally to have cleared the IF flag for us.
235 *
236 * +-----+ +-----+ +-----+
237 * | _interrupt() |--->| do_interrupt() |--->| *setlvl() |
238 * +-----+ +-----+ +-----+
239 *
240 *
241 * low-level HW int | softint
242 *
243 * +-----+ +-----+
244 * | intr_thread_ | | hi-level int
245 * | prolog() | | +-----+
246 * +-----+ +-----+ | hilevel_
247 * | | | | | intr_
248 * | | | | | prolog() |-----+
249 * +-----+ +-----+ | +-----+
250 * | switch_sp_ | | | On intr
251 * | and_call() | | | Stack
252 * +-----+ +-----+ | | +-----+
253 * | | | | | | | | | switch_sp_
254 * | | | | | | | | | and_call()
255 * +-----+ +-----+ | +-----+
256 * | dispatch_ | | | dispatch_
257 * | | | | | | | | | hilevel()

```

```

257 * | hardint() |
258 * +-----+
259 * | v |
260 * +-----+ +-----+ +-----+ +-----+
261 * | sti |->| av_dispatch_autovect |->| cli |-----+
262 * +-----+ +-----+ +-----+
263 *
264 *
265 * | v |
266 * +-----+
267 * | for each |
268 * | handler |
269 * | *intr() |
270 * +-----+
271 *
272 * | v |
273 * +-----+
274 * | intr_thread_ |
275 * | epilog() |
276 * +-----+
277 *
278 * | v |
279 * +-----+
280 * | low-level |
281 * +-----+
282 *
283 * | v |
284 * +-----+
285 * | *setlvlx() |
286 * +-----+
287 *
288 * | v |
289 * +-----+
290 * | return |<----| softint pending? |----->| dosoftint() |<----+
291 * | no |-----+ | yes |-----+
292 *
293 * | softint pil too low |
294 * +-----+
295 *
296 * | v |
297 * +-----+
298 * | dispatch_ |
299 * | softint() |<----| switch_sp_ |<----| *setspl() |
300 * | and_call() | +-----+
301 *
302 * | v |
303 * +-----+
304 * | sti |->| av_dispatch_autovect |->| cli |->| dosoftint_ |
305 * | | +-----+ | epilog() |
306 * +-----+
307 *
308 * | v |
309 * +-----+
310 * | interrupt |
311 * | thread |
312 * | blocked |
313 * +-----+
314 *
315 * | v |
316 * +-----+
317 * | set_base_spl() |->| *setlvlx() |->| splhigh() |->| sti() |->| swtch() |
318 * +-----+ +-----+ +-----+ +-----+
319 *
320 * Calls made on Interrupt Stacks and Epilogue routines
321 *
322 * We use the switch_sp_and_call() assembly routine to switch our sp to the
323 * interrupt stacks and then call the appropriate dispatch function. In the
324 * case of interrupts which may block, softints and hardints, we always ensure
325 * that we are still on the interrupt thread when we call the epilog routine.
326 * This is not just important, it's necessary. If the interrupt thread blocked,
327 * we won't return from our switch_sp_and_call() function and instead we'll go
328 * through and set ourselves up to swtch() directly.
329 *
330 * New Interrupt Flow
331 * -----

```

```

323 *
324 * The apix module has its own interrupt path. This is done for various
325 * reasons. The first is that rather than having global interrupt vectors, we
326 * now have per-cpu vectors.
327 *
328 * The other substantial change is that the apix design does not use the TPR to
329 * mask interrupts below the current level. In fact, except for one special
330 * case, it does not use the TPR at all. Instead, it only uses the IF flag
331 * (cli/sti) to either block all interrupts or allow any interrupts to come in.
332 * The design is such that when interrupts are allowed to come in, if we are
333 * currently servicing a higher priority interrupt, the new interrupt is treated
334 * as pending and serviced later. Specifically, in the pcpusmp module's
335 * apix_intr_enter() the code masks interrupts at or below the current
336 * IPL using the TPR before sending EOI, whereas the apix module's
337 * apix_intr_enter() simply sends EOI.
338 *
339 * The one special case where the apix code uses the TPR is when it calls
340 * through the apix_reg_ops function pointer apix_write_task_reg in
341 * apix_init_intr() to initially mask all levels and then finally to enable all
342 * levels.
343 *
344 * Recall that we come into the interrupt handler with all interrupts masked
345 * by the IF flag. This is because we set up the handler using an
346 * interrupt-gate which is defined architecturally to have cleared the IF flag
347 * for us.
348 *
349 * +-----+ +-----+
350 * | _interrupt() |-----| apix_do_interrupt() |
351 * +-----+ +-----+
352 *
353 * | v |
354 * +-----+
355 * | hard int? |-----+ | softint? |
356 * | | +-----+ | (but no low-level looping) |
357 * | *setlvl() | +-----+
358 * +-----+
359 * | v |
360 * +-----+
361 * | check IPL |
362 * +-----+
363 * | v |
364 * +-----+
365 * | apix_add_ |
366 * | pending_ |
367 * | hardint() |<----| low-level int |
368 * +-----+ +-----+ | hi-level int |
369 * | | +-----+ | | +-----+
370 * | | | apix_intr_ | | | apix_hilevel_ |
371 * | | | thread_prolog() | | | intr_prolog() |
372 * | | +-----+ | | +-----+
373 * | | | |
374 * | | | switch_sp_ | | | switch_sp_ |
375 * | | | and_call() | | | and_call() |
376 * | | +-----+ | | +-----+
377 * | | | |
378 * | | | apix_dispatch_ | | | apix_dispatch_ |
379 * | | | lowlevel() | | | hilevel() |
380 * | | +-----+ | | +-----+
381 * | | | |
382 * | | | +-----+ | | | +-----+
383 * | | | |apix_dispatch_by_vector() | | | |
384 * | | | +-----+ | | | +-----+
385 * | | | |XC_HI_PIL| | | |
386 * | | | +-----+ | | | +-----+
387 * | | | |sti| |*intr()| |cli| | |
388 * | | | +-----+ | | | +-----+
389 * | | | | v | | | |
390 * | | | | | | | |
391 * | | | | | | | |
392 * | | | | | | | |
393 * | | | | | | | |
394 * | | | | | | | |
395 * | | | | | | | |
396 * | | | | | | | |
397 * | | | | | | | |
398 * | | | | | | | |
399 * | | | | | | | |
400 * | | | | | | | |
401 * | | | | | | | |
402 * | | | | | | | |
403 * | | | | | | | |
404 * | | | | | | | |
405 * | | | | | | | |
406 * | | | | | | | |
407 * | | | | | | | |
408 * | | | | | | | |
409 * | | | | | | | |
410 * | | | | | | | |
411 * | | | | | | | |
412 * | | | | | | | |
413 * | | | | | | | |
414 * | | | | | | | |
415 * | | | | | | | |
416 * | | | | | | | |
417 * | | | | | | | |
418 * | | | | | | | |
419 * | | | | | | | |
420 * | | | | | | | |
421 * | | | | | | | |
422 * | | | | | | | |
423 * | | | | | | | |
424 * | | | | | | | |
425 * | | | | | | | |
426 * | | | | | | | |
427 * | | | | | | | |
428 * | | | | | | | |
429 * | | | | | | | |
430 * | | | | | | | |
431 * | | | | | | | |
432 * | | | | | | | |
433 * | | | | | | | |
434 * | | | | | | | |
435 * | | | | | | | |
436 * | | | | | | | |
437 * | | | | | | | |
438 * | | | | | | | |
439 * | | | | | | | |
440 * | | | | | | | |
441 * | | | | | | | |
442 * | | | | | | | |
443 * | | | | | | | |
444 * | | | | | | | |
445 * | | | | | | | |
446 * | | | | | | | |
447 * | | | | | | | |
448 * | | | | | | | |
449 * | | | | | | | |
450 * | | | | | | | |
451 * | | | | | | | |
452 * | | | | | | | |
453 * | | | | | | | |
454 * | | | | | | | |
455 * | | | | | | | |
456 * | | | | | | | |
457 * | | | | | | | |
458 * | | | | | | | |
459 * | | | | | | | |
460 * | | | | | | | |
461 * | | | | | | | |
462 * | | | | | | | |
463 * | | | | | | | |
464 * | | | | | | | |
465 * | | | | | | | |
466 * | | | | | | | |
467 * | | | | | | | |
468 * | | | | | | | |
469 * | | | | | | | |
470 * | | | | | | | |
471 * | | | | | | | |
472 * | | | | | | | |
473 * | | | | | | | |
474 * | | | | | | | |
475 * | | | | | | | |
476 * | | | | | | | |
477 * | | | | | | | |
478 * | | | | | | | |
479 * | | | | | | | |
480 * | | | | | | | |
481 * | | | | | | | |
482 * | | | | | | | |
483 * | | | | | | | |
484 * | | | | | | | |
485 * | | | | | | | |
486 * | | | | | | | |
487 * | | | | | | | |
488 * | | | | | | | |
489 * | | | | | | | |
490 * | | | | | | | |
491 * | | | | | | | |
492 * | | | | | | | |
493 * | | | | | | | |
494 * | | | | | | | |
495 * | | | | | | | |
496 * | | | | | | | |
497 * | | | | | | | |
498 * | | | | | | | |
499 * | | | | | | | |
500 * | | | | | | | |
501 * | | | | | | | |
502 * | | | | | | | |
503 * | | | | | | | |
504 * | | | | | | | |
505 * | | | | | | | |
506 * | | | | | | | |
507 * | | | | | | | |
508 * | | | | | | | |
509 * | | | | | | | |
510 * | | | | | | | |
511 * | | | | | | | |
512 * | | | | | | | |
513 * | | | | | | | |
514 * | | | | | | | |
515 * | | | | | | | |
516 * | | | | | | | |
517 * | | | | | | | |
518 * | | | | | | | |
519 * | | | | | | | |
520 * | | | | | | | |
521 * | | | | | | | |
522 * | | | | | | | |
523 * | | | | | | | |
524 * | | | | | | | |
525 * | | | | | | | |
526 * | | | | | | | |
527 * | | | | | | | |
528 * | | | | | | | |
529 * | | | | | | | |
530 * | | | | | | | |
531 * | | | | | | | |
532 * | | | | | | | |
533 * | | | | | | | |
534 * | | | | | | | |
535 * | | | | | | | |
536 * | | | | | | | |
537 * | | | | | | | |
538 * | | | | | | | |
539 * | | | | | | | |
540 * | | | | | | | |
541 * | | | | | | | |
542 * | | | | | | | |
543 * | | | | | | | |
544 * | | | | | | | |
545 * | | | | | | | |
546 * | | | | | | | |
547 * | | | | | | | |
548 * | | | | | | | |
549 * | | | | | | | |
550 * | | | | | | | |
551 * | | | | | | | |
552 * | | | | | | | |
553 * | | | | | | | |
554 * | | | | | | | |
555 * | | | | | | | |
556 * | | | | | | | |
557 * | | | | | | | |
558 * | | | | | | | |
559 * | | | | | | | |
560 * | | | | | | | |
561 * | | | | | | | |
562 * | | | | | | | |
563 * | | | | | | | |
564 * | | | | | | | |
565 * | | | | | | | |
566 * | | | | | | | |
567 * | | | | | | | |
568 * | | | | | | | |
569 * | | | | | | | |
570 * | | | | | | | |
571 * | | | | | | | |
572 * | | | | | | | |
573 * | | | | | | | |
574 * | | | | | | | |
575 * | | | | | | | |
576 * | | | | | | | |
577 * | | | | | | | |
578 * | | | | | | | |
579 * | | | | | | | |
580 * | | | | | | | |
581 * | | | | | | | |
582 * | | | | | | | |
583 * | | | | | | | |
584 * | | | | | | | |
585 * | | | | | | | |
586 * | | | | | | | |
587 * | | | | | | | |
588 * | | | | | | | |
589 * | | | | | | | |
590 * | | | | | | | |
591 * | | | | | | | |
592 * | | | | | | | |
593 * | | | | | | | |
594 * | | | | | | | |
595 * | | | | | | | |
596 * | | | | | | | |
597 * | | | | | | | |
598 * | | | | | | | |
599 * | | | | | | | |
600 * | | | | | | | |
601 * | | | | | | | |
602 * | | | | | | | |
603 * | | | | | | | |
604 * | | | | | | | |
605 * | | | | | | | |
606 * | | | | | | | |
607 * | | | | | | | |
608 * | | | | | | | |
609 * | | | | | | | |
610 * | | | | | | | |
611 * | | | | | | | |
612 * | | | | | | | |
613 * | | | | | | | |
614 * | | | | | | | |
615 * | | | | | | | |
616 * | | | | | | | |
617 * | | | | | | | |
618 * | | | | | | | |
619 * | | | | | | | |
620 * | | | | | | | |
621 * | | | | | | | |
622 * | | | | | | | |
623 * | | | | | | | |
624 * | | | | | | | |
625 * | | | | | | | |
626 * | | | | | | | |
627 * | | | | | | | |
628 * | | | | | | | |
629 * | | | | | | | |
630 * | | | | | | | |
631 * | | | | | | | |
632 * | | | | | | | |
633 * | | | | | | | |
634 * | | | | | | | |
635 * | | | | | | | |
636 * | | | | | | | |
637 * | | | | | | | |
638 * | | | | | | | |
639 * | | | | | | | |
640 * | | | | | | | |
641 * | | | | | | | |
642 * | | | | | | | |
643 * | | | | | | | |
644 * | | | | | | | |
645 * | | | | | | | |
646 * | | | | | | | |
647 * | | | | | | | |
648 * | | | | | | | |
649 * | | | | | | | |
650 * | | | | | | | |
651 * | | | | | | | |
652 * | | | | | | | |
653 * | | | | | | | |
654 * | | | | | | | |
655 * | | | | | | | |
656 * | | | | | | | |
657 * | | | | | | | |
658 * | | | | | | | |
659 * | | | | | | | |
660 * | | | | | | | |
661 * | | | | | | | |
662 * | | | | | | | |
663 * | | | | | | | |
664 * | | | | | | | |
665 * | | | | | | | |
666 * | | | | | | | |
667 * | | | | | | | |
668 * | | | | | | | |
669 * | | | | | | | |
670 * | | | | | | | |
671 * | | | | | | | |
672 * | | | | | | | |
673 * | | | | | | | |
674 * | | | | | | | |
675 * | | | | | | | |
676 * | | | | | | | |
677 * | | | | | | | |
678 * | | | | | | | |
679 * | | | | | | | |
680 * | | | | | | | |
681 * | | | | | | | |
682 * | | | | | | | |
683 * | | | | | | | |
684 * | | | | | | | |
685 * | | | | | | | |
686 * | | | | | | | |
687 * | | | | | | | |
688 * | | | | | | | |
689 * | | | | | | | |
690 * | | | | | | | |
691 * | | | | | | | |
692 * | | | | | | | |
693 * | | | | | | | |
694 * | | | | | | | |
695 * | | | | | | | |
696 * | | | | | | | |
697 * | | | | | | | |
698 * | | | | | | | |
699 * | | | | | | | |
700 * | | | | | | | |
701 * | | | | | | | |
702 * | | | | | | | |
703 * | | | | | | | |
704 * | | | | | | | |
705 * | | | | | | | |
706 * | | | | | | | |
707 * | | | | | | | |
708 * | | | | | | | |
709 * | | | | | | | |
710 * | | | | | | | |
711 * | | | | | | | |
712 * | | | | | | | |
713 * | | | | | | | |
714 * | | | | | | | |
715 * | | | | | | | |
716 * | | | | | | | |
717 * | | | | | | | |
718 * | | | | | | | |
719 * | | | | | | | |
720 * | | | | | | | |
721 * | | | | | | | |
722 * | | | | | | | |
723 * | | | | | | | |
724 * | | | | | | | |
725 * | | | | | | | |
726 * | | | | | | | |
727 * | | | | | | | |
728 * | | | | | | | |
729 * | | | | | | | |
730 * | | | | | | | |
731 * | | | | | | | |
732 * | | | | | | | |
733 * | | | | | | | |
734 * | | | | | | | |
735 * | | | | | | | |
736 * | | | | | | | |
737 * | | | | | | | |
738 * | | | | | | | |
739 * | | | | | | | |
740 * | | | | | | | |
741 * | | | | | | | |
742 * | | | | | | | |
743 * | | | | | | | |
744 * | | | | | | | |
745 * | | | | | | | |
746 * | | | | | | | |
747 * | | | | | | | |
748 * | | | | | | | |
749 * | | | | | | | |
750 * | | | | | | | |
751 * | | | | | | | |
752 * | | | | | | | |
753 * | | | | | | | |
754 * | | | | | | | |
755 * | | | | | | | |
756 * | | | | | | | |
757 * | | | | | | | |
758 * | | | | | | | |
759 * | | | | | | | |
760 * | | | | | | | |
761 * | | | | | | | |
762 * | | | | | | | |
763 * | | | | | | | |
764 * | | | | | | | |
765 * | | | | | | | |
766 * | | | | | | | |
767 * | | | | | | | |
768 * | | | | | | | |
769 * | | | | | | | |
770 * | | | | | | | |
771 * | | | | | | | |
772 * | | | | | | | |
773 * | | | | | | | |
774 * | | | | | | | |
775 * | | | | | | | |
776 * | | | | | | | |
777 * | | | | | | | |
778 * | | | | | | | |
779 * | | | | | | | |
780 * | | | | | | | |
781 * | | | | | | | |
782 * | | | | | | | |
783 * | | | | | | | |
784 * | | | | | | | |
785 * | | | | | | | |
786 * | | | | | | | |
787 * | | | | | | | |
788 * | | | | | | | |
789 * | | | | | | | |
790 * | | | | | | | |
791 * | | | | | | | |
792 * | | | | | | | |
793 * | | | | | | | |
794 * | | | | | | | |
795 * | | | | | | | |
796 * | | | | | | | |
797 * | | | | | | | |
798 * | | | | | | | |
799 * | | | | | | | |
800 * | | | | | | | |
801 * | | | | | | | |
802 * | | | | | | | |
803 * | | | | | | | |
804 * | | | | | | | |
805 * | | | | | | | |
806 * | | | | | | | |
807 * | | | | | | | |
808 * | | | | | | | |
809 * | | | | | | | |
810 * | | | | | | | |
811 * | | | | | | | |
812 * | | | | | | | |
813 * | | | | | | | |
814 * | | | | | | | |
815 * | | | | | | | |
816 * | | | | | | | |
817 * | | | | | | | |
818 * | | | | | | | |
819 * | | | | | | | |
820 * | | | | | | | |
821 * | | | | | | | |
822 * | | | | | | | |
823 * | | | | | | | |
824 * | | | | | | | |
825 * | | | | | | | |
826 * | | | | | | | |
827 * | | | | | | | |
828 * | | | | | | | |
829 * | | | | | | | |
830 * | | | | | | | |
831 * | | | | | | | |
832 * | | | | | | | |
833 * | | | | | | | |
834 * | | | | | | | |
835 * | | | | | | | |
836 * | | | | | | | |
837 * | | | | | | | |
838 * | | | | | | | |
839 * | | | | | | | |
840 * | | | | | | | |
841 * | | | | | | | |
842 * | | | | | | | |
843 * | | | | | | | |
844 * | | | | | | | |
845 * | | | | | | | |
846 * | | | | | | | |
847 * | | | | | | | |
848 * | | | | | | | |
849 * | | | | | | | |
850 * | | | | | | | |
851 * | | | | | | | |
852 * | | | | | | | |
853 * | | | | | | | |
854 * | | | | | | | |
855 * | | | | | | | |
856 * | | | | | | | |
857 * | | | | | | | |
858 * | | | | | | | |
859 * | | | | | | | |
860 * | | | | | | | |
861 * | | | | | | | |
862 * | | | | | | | |
863 * | | | | | | | |
864 * | | | | | | | |
865 * | | | | | | | |
866 * | | | | | | | |
867 * | | | | | | | |
868 * | | | | | | | |
869 * | | | | | | | |
870 * | | | | | | | |
871 * | | | | | | | |
872 * | | | | | | | |
873 * | | | | | | | |
874 * | | | | | | | |
875 * | | | | | | | |
876 * | | | | | | | |
877 * | | | | | | | |
878 * | | | | | | | |
879 * | | | | | | | |
880 * | | | | | | | |
881 * | | | | | | | |
882 * | | | | | | | |
883 * | | | | | | | |
884 * | | | | | | | |
885 * | | | | | | | |
886 * | | | | | | | |
887 * | | | | | | | |
888 * | | | | | | | |
889 * | | | | | | | |
890 * | | | | | | | |
891 * | | | | | | | |
892 * | | | | | | | |
893 * | | | | | | | |
894 * | | | | | | | |
895 * | | | | | | | |
896 * | | | | | | | |
897 * | | | | | | | |
898 * | | | | | | | |
899 * | | | | | | | |
900 * | | | | | | | |
901 * | | | | | | | |
902 * | | | | | | | |
903 * | | | | | | | |
904 * | | | | | | | |
905 * | | | | | | | |
906 * | | | | | | | |
907 * | | | | | | | |
908 * | | | | | | | |
909 * | | | | | | | |
910 * | | | | | | | |
911 * | | | | | | | |
912 * | | | | | | | |
913 * | | | | | | | |
914 * | | | | | | | |
915 * | | | | | | | |
916 * | | | | | | | |
917 * | | | | | | | |
918 * | | | | | | | |
919 * | | | | | | | |
920 * | | | | | | | |
921 * | | | | | | | |
922 * | | | | | | | |
923 * | | | | | | | |
924 * | | | | | | | |
925 * | | | | | | | |
926 * | | | | | | | |
927 * | | | | | | | |
928 * | | | | | | | |
929 * | | | | | | | |
930 * | | | | | | | |
931 * | | | | | | | |
932 * | | | | | | | |
933 * | | | | | | | |
934 * | | | | | | | |
935 * | | | | | | | |
936 * | | | | | | | |
937 * | | | | | | | |
938 * | | | | | | | |
939 * | | | | | | | |
940 * | | | | | | | |
941 * | | | | | | | |
942 * | | | | | | | |
943 * | | | | | | | |
944 * | | | | | | | |
945 * | | | | | | | |
946 * | | | | | | | |
947 * | | | | | | | |
948 * | | | | | | | |
949 * | | | | | | | |
950 * | | | | | | | |
951 * | | | | | | | |
952 * | | | | | | | |
953 * | | | | | | | |
954 * | | | | | | | |
955 * | | | | | | | |
956 * | | | | | | | |
957 * | | | | | | | |
958 * | | | | | | | |
959 * | | | | | | | |
960 * | | | | | | | |
961 * | | | | | | | |
962 * | | | | | | | |
963 * | | | | | | | |
964 * | | | | | | | |
965 * | | | | | | | |
966 * | | | | | | | |
967 * | | | | | | | |
968 * | | | | | | | |
969 * | | | | | | | |
970 * | | | | | | | |
971 * | | | | | | | |
972 * | | | | | | | |
973 * | | | | | | | |
974 * | | | | | | | |
975 * | | | | | | | |
976 * | | | | | | | |
977 * | | | | | | | |
978 * | | | | | | | |
979 * | | | | | | | |
980 * | | | | | | | |
981 * | | | | | | | |
982 * | | | | | | | |
983 * | | | | | | | |
984 * | | | | | | | |
985 * | | | | | | | |
986 * | | | | | | | |
987 * | | | | | | | |
988 * | | | | | | | |
989 * | | | | | | | |
990 * | | | | | | | |
991 * | | | | | | | |
992 * | | | | | | | |
993 * | | | | | | | |
994 * | | | | | | | |
995 * | | | | | | | |
996 * | | | | | | | |
997 * | | | | | | | |
998 * | | | | | | | |
999 * | | | | | | | |
1000 * | | | | | | | |

```



```

521 */
522 void
523 set_base_spl(void)
524 {
525     struct cpu *cpu = CPU;
526     uint16_t active = (uint16_t)cpu->cpu_intr_actv;

528     cpu->cpu_base_spl = active == 0 ? 0 : bsrw_insn(active);
529 }
_____ unchanged_portion_omitted _____

1429 /*
1430 * Common tasks always done by _sys_rtt, called with interrupts disabled.
1431 * Returns 1 if returning to userland, 0 if returning to system mode.
1432 */
1433 int
1434 sys_rtt_common(struct regs *rp)
1435 {
1436     kthread_t *tp;
1437     extern void mutex_exit_critical_start();
1438     extern long mutex_exit_critical_size;
1439     extern void mutex_owner_running_critical_start();
1440     extern long mutex_owner_running_critical_size;

1442 loop:

1444     /*
1445     * Check if returning to user
1446     */
1447     tp = CPU->cpu_thread;
1448     if (USERMODE(rp->r_cs)) {
1449         /*
1450         * Check if AST pending.
1451         */
1452         if (tp->t_astflag) {
1453             /*
1454             * Let trap() handle the AST
1455             */
1456             sti();
1457             rp->r_trapno = T_AST;
1458             trap(rp, (caddr_t)0, CPU->cpu_id);
1459             cli();
1460             goto loop;
1461         }

1463 #if defined(__amd64)
1464     /*
1465     * We are done if segment registers do not need updating.
1466     */
1467     if (tp->t_lwp->lwp_pcb.pcb_rupdate == 0)
1468         return (1);

1470     if (update_sregs(rp, tp->t_lwp)) {
1471         /*
1472         * 1 or more of the selectors is bad.
1473         * Deliver a SIGSEGV.
1474         */
1475         proc_t *p = ttoproc(tp);

1477         sti();
1478         mutex_enter(&p->p_lock);
1479         tp->t_lwp->lwp_cursig = SIGSEGV;
1480         mutex_exit(&p->p_lock);
1481         psig();
1482         tp->t_sig_check = 1;

```

```

1483         cli();
1484     }
1485     tp->t_lwp->lwp_pcb.pcb_rupdate = 0;

1487 #endif /* __amd64 */
1488     return (1);
1489 }

1491 #if !defined(__xpv)
1492 /*
1493  * Assert that we're not trying to return into the syscall return
1494  * trampolines. Things will go baaaaad if we try to do that.
1495  *
1496  * Note that none of these run with interrupts on, so this should
1497  * never happen (even in the sysexit case the STI doesn't take effect
1498  * until after sysexit finishes).
1499  */
1500     extern void tr_sysc_ret_start();
1501     extern void tr_sysc_ret_end();
1502     ASSERT(!(rp->r_pc >= (uintptr_t)tr_sysc_ret_start &&
1503            rp->r_pc <= (uintptr_t)tr_sysc_ret_end));
1504 #endif

1506 /*
1507  * Here if we are returning to supervisor mode.
1508  * Check for a kernel preemption request.
1509  */
1510     if (CPU->cpu_kprunrun && (rp->r_ps & PS_IE)) {

1512         /*
1513         * Do nothing if already in kpreempt
1514         */
1515         if (!tp->t_preempt_lk) {
1516             tp->t_preempt_lk = 1;
1517             sti();
1518             kpreempt(1); /* asynchronous kpreempt call */
1519             cli();
1520             tp->t_preempt_lk = 0;
1521         }
1522     }

1524     /*
1525     * If we interrupted the mutex_exit() critical region we must
1526     * reset the PC back to the beginning to prevent missed wakeups
1527     * See the comments in mutex_exit() for details.
1528     */
1529     if ((uintptr_t)rp->r_pc - (uintptr_t)mutex_exit_critical_start <
1530        mutex_exit_critical_size) {
1531         rp->r_pc = (greg_t)mutex_exit_critical_start;
1532     }

1534     /*
1535     * If we interrupted the mutex_owner_running() critical region we
1536     * must reset the PC back to the beginning to prevent dereferencing
1537     * of a freed thread pointer. See the comments in mutex_owner_running
1538     * for details.
1539     */
1540     if ((uintptr_t)rp->r_pc -
1541        (uintptr_t)mutex_owner_running_critical_start <
1542        mutex_owner_running_critical_size) {
1543         rp->r_pc = (greg_t)mutex_owner_running_critical_start;
1544     }

1546     return (0);
1547 }
_____ unchanged_portion_omitted _____

```

new/usr/src/uts/i86pc/os/mach_kdi.c

1

```
*****
4133 Fri Apr 6 17:25:02 2018
new/usr/src/uts/i86pc/os/mach_kdi.c
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22 * Copyright 2007 Sun Microsystems, Inc. All rights reserved.
23 * Use is subject to license terms.
24 *
25 * Copyright 2018 Joyent, Inc.
26 */

26 #pragma ident      "%Z%M% %I%      %E% SMI"

28 /*
29 * Kernel/Debugger Interface (KDI) routines.  Called during debugger under
30 * various system states (boot, while running, while the debugger has control).
31 * Functions intended for use while the debugger has control may not grab any
32 * locks or perform any functions that assume the availability of other system
33 * services.
34 */

36 #include <sys/system.h>
37 #include <sys/x86_archext.h>
38 #include <sys/kdi_impl.h>
39 #include <sys/smp_impldefs.h>
40 #include <sys/psm_types.h>
41 #include <sys/segments.h>
42 #include <sys/archsystem.h>
43 #include <sys/controlregs.h>
44 #include <sys/trap.h>
45 #include <sys/kobj.h>
46 #include <sys/kobj_impl.h>
47 #include <sys/mach_mmu.h>

49 void
50 kdi_idt_write(gate_desc_t *gate, uint_t vec)
51 {
52     gate_desc_t *idt = CPU->cpu_m.mcpu_idt;

54     /*
55      * See kdi_idtr_set().
56      */
57     if (idt == NULL) {
```

new/usr/src/uts/i86pc/os/mach_kdi.c

2

```
58         desctbr_t idtr;
59         rd_idtr(&idtr);
60         idt = (gate_desc_t *)idtr.dtr_base;
61     }

63     idt[vec] = *gate;
64 }

unchanged_portion_omitted_

116 void
117 kdi_flush_caches(void)
118 {
119     reload_cr3();
120 }

116 extern void kdi_slave_entry(void);

118 void
119 kdi_stop_slaves(int cpu, int doxc)
120 {
121     if (doxc)
122         kdi_xc_others(cpu, kdi_slave_entry);
123 }

unchanged_portion_omitted_
```

```

*****
14906 Fri Apr 6 17:25:02 2018
new/usr/src/uts/i86pc/os/mlsetup.c
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22 * Copyright (c) 2012 Gary Mills
23 *
24 * Copyright (c) 1993, 2010, Oracle and/or its affiliates. All rights reserved.
25 * Copyright (c) 2011 by Delphix. All rights reserved.
26 * Copyright 2018 Joyent, Inc.
27 * Copyright 2016 Joyent, Inc.
28 */
29 * Copyright (c) 2010, Intel Corporation.
30 * All rights reserved.
31 */

33 #include <sys/types.h>
34 #include <sys/sysmacros.h>
35 #include <sys/disp.h>
36 #include <sys/promif.h>
37 #include <sys/clock.h>
38 #include <sys/cpuvar.h>
39 #include <sys/stack.h>
40 #include <vm/as.h>
41 #include <vm/hat.h>
42 #include <sys/reboot.h>
43 #include <sys/avintr.h>
44 #include <sys/vtrace.h>
45 #include <sys/proc.h>
46 #include <sys/thread.h>
47 #include <sys/cpupart.h>
48 #include <sys/pset.h>
49 #include <sys/copyops.h>
50 #include <sys/pg.h>
51 #include <sys/disp.h>
52 #include <sys/debug.h>
53 #include <sys/sunddi.h>
54 #include <sys/x86_archext.h>
55 #include <sys/privregs.h>
56 #include <sys/machsystem.h>
57 #include <sys/ontrap.h>
58 #include <sys/bootconf.h>

```

```

59 #include <sys/boot_console.h>
60 #include <sys/kdi_machimpl.h>
61 #include <sys/archsystem.h>
62 #include <sys/promif.h>
63 #include <sys/pci_cfgspace.h>
64 #include <sys/bootvfs.h>
65 #include <sys/tsc.h>
66 #ifdef __xpv
67 #include <sys/hypervisor.h>
68 #else
69 #include <sys/xpv_support.h>
70 #endif

72 /*
73  * some globals for patching the result of cpuid
74  * to solve problems w/ creative cpu vendors
75 */

77 extern uint32_t cpuid_feature_ecx_include;
78 extern uint32_t cpuid_feature_ecx_exclude;
79 extern uint32_t cpuid_feature_edx_include;
80 extern uint32_t cpuid_feature_edx_exclude;

82 /*
83  * Set console mode
84 */
85 static void
86 set_console_mode(uint8_t val)
87 {
88     struct bop_regs rp = {0};

90     rp.eax.byte.ah = 0x0;
91     rp.eax.byte.al = val;
92     rp.ebx.word.bx = 0x0;

94     BOP_DOINT(bootops, 0x10, &rp);
95 }

98 /*
99  * Setup routine called right before main(). Interposing this function
100 * before main() allows us to call it in a machine-independent fashion.
101 */
102 void
103 mlsetup(struct regs *rp)
104 {
105     u_longlong_t prop_value;
106     extern struct classfuncs sys_classfuncs;
107     extern disp_t cpu0_disp;
108     extern char t0stack[];
109     extern int post_fastreboot;
110     extern uint64_t plat_dr_options;

112     ASSERT_STACK_ALIGNED();

114     /*
115      * initialize cpu_self
116      */
117     cpu[0]->cpu_self = cpu[0];

119 #if defined(__xpv)
120     /*
121      * Point at the hypervisor's virtual cpu structure
122      */
123     cpu[0]->cpu_m.mcpu_vcpu_info = &HYPERVISOR_shared_info->vcpu_info[0];
124 #endif

```

```

126  /*
127  * check if we've got special bits to clear or set
128  * when checking cpu features
129  */

131  if (bootprop_getval("cpuid_feature_ecx_include", &prop_value) != 0)
132      cpuid_feature_ecx_include = 0;
133  else
134      cpuid_feature_ecx_include = (uint32_t)prop_value;

136  if (bootprop_getval("cpuid_feature_ecx_exclude", &prop_value) != 0)
137      cpuid_feature_ecx_exclude = 0;
138  else
139      cpuid_feature_ecx_exclude = (uint32_t)prop_value;

141  if (bootprop_getval("cpuid_feature_edx_include", &prop_value) != 0)
142      cpuid_feature_edx_include = 0;
143  else
144      cpuid_feature_edx_include = (uint32_t)prop_value;

146  if (bootprop_getval("cpuid_feature_edx_exclude", &prop_value) != 0)
147      cpuid_feature_edx_exclude = 0;
148  else
149      cpuid_feature_edx_exclude = (uint32_t)prop_value;

151  #if !defined(__xpv)
152  /*
153  * Check to see if KPTI has been explicitly enabled or disabled.
154  * We have to check this before init_desctbls().
155  */
156  if (bootprop_getval("kpti", &prop_value) == 0) {
157      kpti_enable = (uint64_t)(prop_value == 1);
158      prom_printf("unix: forcing kpti to %s due to boot argument\n",
159                 (kpti_enable == 1) ? "ON" : "OFF");
160  } else {
161      kpti_enable = 1;
162  }

164  if (bootprop_getval("pcid", &prop_value) == 0 && prop_value == 0) {
165      prom_printf("unix: forcing pcid to OFF due to boot argument\n");
166      x86_use_pcid = 0;
167  } else if (kpti_enable != 1) {
168      x86_use_pcid = 0;
169  }
170  #endif

172  /*
173  * Initialize idt0, gdt0, ldt0_default, ktss0 and dftss.
174  */
175  init_desctbls();

177  /*
178  * lgrp_init() and possibly cpuid_pass1() need PCI config
179  * space access
180  */

181  #if defined(__xpv)
182      if (DOMAIN_IS_INITDOMAIN(xen_info))
183          pci_cfgspace_init();
184  #else
185      pci_cfgspace_init();
186  /*
187  * Initialize the platform type from CPU 0 to ensure that
188  * determine_platform() is only ever called once.
189  */
190  determine_platform();

```

```

191  #endif

193  /*
194  * The first lightweight pass (pass0) through the cpuid data
195  * was done in locore before mlsetup was called. Do the next
196  * pass in C code.
197  *
198  * The x86_featureset is initialized here based on the capabilities
199  * of the boot CPU. Note that if we choose to support CPUs that have
200  * different feature sets (at which point we would almost certainly
201  * want to set the feature bits to correspond to the feature
202  * minimum) this value may be altered.
203  */
204  cpuid_pass1(cpu[0], x86_featureset);

206  #if !defined(__xpv)
207      if ((get_hwenv() & HW_XEN_HVM) != 0)
208          xen_hvm_init();

210  /*
211  * Before we do anything with the TSCs, we need to work around
212  * Intel erratum BT81. On some CPUs, warm reset does not
213  * clear the TSC. If we are on such a CPU, we will clear TSC ourselves
214  * here. Other CPUs will clear it when we boot them later, and the
215  * resulting skew will be handled by tsc_sync_master()/_slave();
216  * note that such skew already exists and has to be handled anyway.
217  *
218  * We do this only on metal. This same problem can occur with a
219  * hypervisor that does not happen to virtualise a TSC that starts from
220  * zero, regardless of CPU type; however, we do not expect hypervisors
221  * that do not virtualise TSC that way to handle writes to TSC
222  * correctly, either.
223  */
224  if (get_hwenv() == HW_NATIVE &&
225      cpuid_getvendor(CPU) == X86_VENDOR_Intel &&
226      cpuid_getfamily(CPU) == 6 &&
227      (cpuid_getmodel(CPU) == 0x2d || cpuid_getmodel(CPU) == 0x3e) &&
228      is_x86_feature(x86_featureset, X86FSET_TSC)) {
229      (void) wrmsr(REG_TSC, 0UL);
230  }

232  /*
233  * Patch the tsc_read routine with appropriate set of instructions,
234  * depending on the processor family and architecture, to read the
235  * time-stamp counter while ensuring no out-of-order execution.
236  * Patch it while the kernel text is still writable.
237  *
238  * Note: tsc_read is not patched for intel processors whose family
239  * is >6 and for amd whose family >f (in case they don't support rdtscp
240  * instruction, unlikely). By default tsc_read will use cpuid for
241  * serialization in such cases. The following code needs to be
242  * revisited if intel processors of family >= f retains the
243  * instruction serialization nature of mfence instruction.
244  * Note: tsc_read is not patched for x86 processors which do
245  * not support "mfence". By default tsc_read will use cpuid for
246  * serialization in such cases.
247  *
248  * The Xen hypervisor does not correctly report whether rdtscp is
249  * supported or not, so we must assume that it is not.
250  */
251  if ((get_hwenv() & HW_XEN_HVM) == 0 &&
252      is_x86_feature(x86_featureset, X86FSET_TSCP))
253      patch_tsc_read(TSC_TSCP);
254  else if (cpuid_getvendor(CPU) == X86_VENDOR_AMD &&
255      cpuid_getfamily(CPU) <= 0xf &&
256      is_x86_feature(x86_featureset, X86FSET_SSE2))

```

```

257     patch_tsc_read(TSC_RDTSC_MFENCE);
258     else if (cpuid_getvendor(CPU) == X86_VENDOR_Intel &&
259            cpuid_getfamily(CPU) <= 6 &&
260            is_x86_feature(x86_featureset, X86FSET_SSE2))
261         patch_tsc_read(TSC_RDTSC_LFENCE);
263 #endif /* !__xpv */

265 #if defined(__i386) && !defined(__xpv)
266     /*
267      * Some i386 processors do not implement the rdtsc instruction,
268      * or at least they do not implement it correctly. Patch them to
269      * return 0.
270      */
271     if (!is_x86_feature(x86_featureset, X86FSET_TSC))
272         patch_tsc_read(TSC_NONE);
273 #endif /* __i386 && !__xpv */

275 #if defined(__amd64) && !defined(__xpv)
276     patch_memops(cpuid_getvendor(CPU));
277 #endif /* __amd64 && !__xpv */

279 #if !defined(__xpv)
280     /* XXPV what, if anything, should be dorked with here under xen? */

282     /*
283      * While we're thinking about the TSC, let's set up %cr4 so that
284      * userland can issue rdtsc, and initialize the TSC_AUX value
285      * (the cpuid) for the rdtscp instruction on appropriately
286      * capable hardware.
287      */
288     if (is_x86_feature(x86_featureset, X86FSET_TSC))
289         setcr4(getcr4() & ~CR4_TSD);

291     if (is_x86_feature(x86_featureset, X86FSET_TSCP))
292         (void) wrmsr(MSR_AMD_TSCAUX, 0);

294     /*
295      * Let's get the other %cr4 stuff while we're here. Note, we defer
296      * enabling CR4_SMAP until startup_end(); however, that's importantly
297      * before we start other CPUs. That ensures that it will be synced out
298      * to other CPUs.
299      */
300     if (is_x86_feature(x86_featureset, X86FSET_DE))
301         setcr4(getcr4() | CR4_DE);

303     if (is_x86_feature(x86_featureset, X86FSET_SMEP))
304         setcr4(getcr4() | CR4_SMEP);
305 #endif /* __xpv */

307     /*
308      * initialize t0
309      */
310     t0.t_stk = (caddr_t)rp - MINFRAME;
311     t0.t_stkbase = t0stack;
312     t0.t_pri = maxclsypri - 3;
313     t0.t_schedflag = TS_LOAD | TS_DONT_SWAP;
314     t0.t_procp = &p0;
315     t0.t_plockp = &p0lock.pl_lock;
316     t0.t_lwp = &lwp0;
317     t0.t_forw = &t0;
318     t0.t_back = &t0;
319     t0.t_next = &t0;
320     t0.t_prev = &t0;
321     t0.t_cpu = cpu[0];
322     t0.t_disp_queue = &cpu0_disp;

```

```

323     t0.t_bind_cpu = PBIND_NONE;
324     t0.t_bind_pset = PS_NONE;
325     t0.t_bindflag = (uchar_t)default_binding_mode;
326     t0.t_cpupart = &cp_default;
327     t0.t_clfuncs = &sys_classfuncs.thread;
328     t0.t_copyops = NULL;
329     THREAD_ONPROC(&t0, CPU);

331     lwp0.lwp_thread = &t0;
332     lwp0.lwp_regs = (void *)rp;
333     lwp0.lwp_procp = &p0;
334     t0.t_tid = p0.p_lwpcnt = p0.p_lwprcnt = p0.p_lwpid = 1;

336     p0.p_exec = NULL;
337     p0.p_stat = SRUN;
338     p0.p_flag = SSYS;
339     p0.p_tlist = &t0;
340     p0.p_stksize = 2*PAGESIZE;
341     p0.p_stkpageszc = 0;
342     p0.p_as = &kas;
343     p0.p_lockp = &p0lock;
344     p0.p_brkpageszc = 0;
345     p0.p_tl_lgrpid = LGRP_NONE;
346     p0.p_tr_lgrpid = LGRP_NONE;
347     psecflags_default(&p0.p_secflags);

349     sigorset(&p0.p_ignore, &ignoredefault);

351     CPU->cpu_thread = &t0;
352     bzero(&cpu0_disp, sizeof (disp_t));
353     CPU->cpu_disp = &cpu0_disp;
354     CPU->cpu_disp->disp_cpu = CPU;
355     CPU->cpu_dispthread = &t0;
356     CPU->cpu_idle_thread = &t0;
357     CPU->cpu_flags = CPU_READY | CPU_RUNNING | CPU_EXISTS | CPU_ENABLE;
358     CPU->cpu_dispatch_pri = t0.t_pri;

360     CPU->cpu_id = 0;

362     CPU->cpu_pri = 12; /* initial PIL for the boot CPU */

364     /*
365      * The kernel doesn't use LDTs unless a process explicitly requests one.
366      */
367     p0.p_ldt_desc = null_sdesc;

369     /*
370      * Initialize thread/cpu microstate accounting
371      */
372     init_mstate(&t0, LMS_SYSTEM);
373     init_cpu_mstate(CPU, CMS_SYSTEM);

375     /*
376      * Initialize lists of available and active CPUs.
377      */
378     cpu_list_init(CPU);

380     pg_cpu_bootstrap(CPU);

382     /*
383      * Now that we have taken over the GDT, IDT and have initialized
384      * active CPU list it's time to inform kmdb if present.
385      */
386     if (boothowto & RB_DEBUG)
387         kdi_idt_sync();

```

```

389     if (BOP_GETPROPLEN(bootops, "efi-systab") < 0) {
390         /*
391          * In BIOS system, explicitly set console to text mode (0x3)
392          * if this is a boot post Fast Reboot, and the console is set
393          * to CONS_SCREEN_TEXT.
394          */
395         if (post_fastreboot &&
396             boot_console_type(NULL) == CONS_SCREEN_TEXT) {
397             set_console_mode(0x3);
398         }
399     }

401     /*
402     * If requested (boot -d) drop into kmdb.
403     *
404     * This must be done after cpu_list_init() on the 64-bit kernel
405     * since taking a trap requires that we re-compute gsbase based
406     * on the cpu list.
407     */
408     if (boothowto & RB_DEBUGENTER)
409         kmdb_enter();

411     cpu_vm_data_init(CPU);

413     rp->r_fp = 0; /* terminate kernel stack traces! */

415     prom_init("kernel", (void *)NULL);

417     /* User-set option overrides firmware value. */
418     if (bootprop_getval(PLAT_DR_OPTIONS_NAME, &prop_value) == 0) {
419         plat_dr_options = (uint64_t)prop_value;
420     }
421 #if defined(__xpv)
422     /* No support of DR operations on xpv */
423     plat_dr_options = 0;
424 #else /* __xpv */
425     /* Flag PLAT_DR_FEATURE_ENABLED should only be set by DR driver. */
426     plat_dr_options &= ~PLAT_DR_FEATURE_ENABLED;
427 #ifndef __amd64
428     /* Only enable CPU/memory DR on 64 bits kernel. */
429     plat_dr_options &= ~PLAT_DR_FEATURE_MEMORY;
430     plat_dr_options &= ~PLAT_DR_FEATURE_CPU;
431 #endif /* __amd64 */
432 #endif /* __xpv */

434     /*
435     * Get value of "plat_dr_physmax" boot option.
436     * It overrides values calculated from MSCT or SRAT table.
437     */
438     if (bootprop_getval(PLAT_DR_PHYSMAX_NAME, &prop_value) == 0) {
439         plat_dr_physmax = ((uint64_t)prop_value) >> PAGESHIFT;
440     }

442     /* Get value of boot_ncpus. */
443     if (bootprop_getval(BOOT_NCPUS_NAME, &prop_value) != 0) {
444         boot_ncpus = NCPU;
445     } else {
446         boot_ncpus = (int)prop_value;
447         if (boot_ncpus <= 0 || boot_ncpus > NCPU)
448             boot_ncpus = NCPU;
449     }

451     /*
452     * Set max_ncpus and boot_max_ncpus to boot_ncpus if platform doesn't
453     * support CPU DR operations.
454     */

```

```

455     if (plat_dr_support_cpu() == 0) {
456         max_ncpus = boot_max_ncpus = boot_ncpus;
457     } else {
458         if (bootprop_getval(PLAT_MAX_NCPUS_NAME, &prop_value) != 0) {
459             max_ncpus = NCPU;
460         } else {
461             max_ncpus = (int)prop_value;
462             if (max_ncpus <= 0 || max_ncpus > NCPU) {
463                 max_ncpus = NCPU;
464             }
465             if (boot_ncpus > max_ncpus) {
466                 boot_ncpus = max_ncpus;
467             }
468         }

470         if (bootprop_getval(BOOT_MAX_NCPUS_NAME, &prop_value) != 0) {
471             boot_max_ncpus = boot_ncpus;
472         } else {
473             boot_max_ncpus = (int)prop_value;
474             if (boot_max_ncpus <= 0 || boot_max_ncpus > NCPU) {
475                 boot_max_ncpus = boot_ncpus;
476             } else if (boot_max_ncpus > max_ncpus) {
477                 boot_max_ncpus = max_ncpus;
478             }
479         }
480     }

482     /*
483     * Initialize the lgrp framework
484     */
485     lgrp_init(LGRP_INIT_STAGE1);

487     if (boothowto & RB_HALT) {
488         prom_printf("unix: kernel halted by -h flag\n");
489         prom_enter_mon();
490     }

492     ASSERT_STACK_ALIGNED();

494     /*
495     * Fill out cpu_ucode_info. Update microcode if necessary.
496     */
497     ucode_check(CPU);

499     if (workaround_errata(CPU) != 0)
500         panic("critical workaround(s) missing for boot cpu");
501 }

```

unchanged_portion_omitted

```

*****
18011 Fri Apr 6 17:25:02 2018
new/usr/src/uts/i86pc/os/mp_pc.c
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22 * Copyright (c) 2007, 2010, Oracle and/or its affiliates. All rights reserved.
23 */
24 /*
25 * Copyright (c) 2010, Intel Corporation.
26 * All rights reserved.
27 */
28 /*
29 * Copyright 2018 Joyent, Inc
30 * Copyright 2011 Joyent, Inc. All rights reserved.
31 */
32 /*
33 * Welcome to the world of the "real mode platter".
34 * See also startup.c, mpcore.s and apic.c for related routines.
35 */
36
37 #include <sys/types.h>
38 #include <sys/system.h>
39 #include <sys/cpuvar.h>
40 #include <sys/cpu_module.h>
41 #include <sys/kmem.h>
42 #include <sys/archsystem.h>
43 #include <sys/machsystem.h>
44 #include <sys/controlregs.h>
45 #include <sys/x86_archext.h>
46 #include <sys/smp_impldefs.h>
47 #include <sys/sysmacros.h>
48 #include <sys/mach_mmu.h>
49 #include <sys/promif.h>
50 #include <sys/cpu.h>
51 #include <sys/cpu_event.h>
52 #include <sys/sunndi.h>
53 #include <sys/fs/dv_node.h>
54 #include <vm/hat_i86.h>
55 #include <vm/as.h>
56
57 extern cpuset_t cpu_ready_set;

```

```

59 extern int mp_start_cpu_common(cpu_t *cp, boolean_t boot);
60 extern void real_mode_start_cpu(void);
61 extern void real_mode_start_cpu_end(void);
62 extern void real_mode_stop_cpu_stage1(void);
63 extern void real_mode_stop_cpu_stage1_end(void);
64 extern void real_mode_stop_cpu_stage2(void);
65 extern void real_mode_stop_cpu_stage2_end(void);
66
67 void rmp_gdt_init(rm_platter_t *);
68
69 /*
70 * Fill up the real mode platter to make it easy for real mode code to
71 * kick it off. This area should really be one passed by boot to kernel
72 * and guaranteed to be below 1MB and aligned to 16 bytes. Should also
73 * have identical physical and virtual address in paged mode.
74 */
75 static ushort_t *warm_reset_vector = NULL;
76
77 int
78 mach_cpucontext_init(void)
79 {
80     ushort_t *vec;
81     ulong_t addr;
82     struct rm_platter *rm = (struct rm_platter *)rm_platter_va;
83
84     if (!(vec = (ushort_t *)psm_map_phys(WARM_RESET_VECTOR,
85     sizeof (vec), PROT_READ | PROT_WRITE)))
86         return (-1);
87
88     /*
89     * setup secondary cpu bios boot up vector
90     * Write page offset to 0x467 and page frame number to 0x469.
91     */
92     addr = (ulong_t)((caddr_t)rm->rm_code - (caddr_t)rm) + rm_platter_pa;
93     vec[0] = (ushort_t)(addr & PAGEOFFSET);
94     vec[1] = (ushort_t)((addr & (0xffff & PAGEMASK)) >> 4);
95     warm_reset_vector = vec;
96
97     /* Map real mode platter into kas so kernel can access it. */
98     hat_devload(kas.a_hat,
99     (caddr_t)(uintptr_t)rm_platter_pa, MMU_PAGESIZE,
100     btop(rm_platter_pa), PROT_READ | PROT_WRITE | PROT_EXEC,
101     HAT_LOAD_NOCONSIST);
102
103     /* Copy CPU startup code to rm_platter if it's still during boot. */
104     if (!plat_dr_enabled()) {
105         ASSERT((size_t)real_mode_start_cpu_end -
106             (size_t)real_mode_start_cpu <= RM_PLATTER_CODE_SIZE);
107         bcopy((caddr_t)real_mode_start_cpu, (caddr_t)rm->rm_code,
108             (size_t)real_mode_start_cpu_end -
109             (size_t)real_mode_start_cpu);
110     }
111
112     return (0);
113 }
114
115 _____ unchanged_portion_omitted _____
116
125 #if defined(__amd64)
126 extern void *long_mode_64(void);
127 #endif /* __amd64 */
128
129 /*ARGSUSED*/
130 void
131 rmp_gdt_init(rm_platter_t *rm)
132 {

```

```

134 #if defined(__amd64)
135     /* Use the kas address space for the CPU startup thread. */
136     if (mmu_ptob(kas.a_hat->hat_htable->ht_pfn) > 0xffffffffUL) {
137         if (MAKECR3(kas.a_hat->hat_htable->ht_pfn) > 0xffffffffUL)
138             panic("Cannot initialize CPUs; kernel's 64-bit page tables\n"
139                 "located above 4G in physical memory (@ 0x%lx)",
140                 mmu_ptob(kas.a_hat->hat_htable->ht_pfn));
141     }
142     MAKECR3(kas.a_hat->hat_htable->ht_pfn);
143
144     /*
145     * Setup pseudo-descriptors for temporary GDT and IDT for use ONLY
146     * by code in real_mode_start_cpu():
147     *
148     * GDT[0]: NULL selector
149     * GDT[1]: 64-bit CS: Long = 1, Present = 1, bits 12, 11 = 1
150     *
151     * Clear the IDT as interrupts will be off and a limit of 0 will cause
152     * the CPU to triple fault and reset on an NMI, seemingly as reasonable
153     * a course of action as any other, though it may cause the entire
154     * platform to reset in some cases...
155     */
156     rm->rm_temp_gdt[0] = 0;
157     rm->rm_temp_gdt[TEMPGDT_KCODE64] = 0x20980000000000ULL;
158
159     rm->rm_temp_gdt_lim = (ushort_t)(sizeof (rm->rm_temp_gdt) - 1);
160     rm->rm_temp_gdt_base = rm_platter_pa +
161         (uint32_t)offsetof(rm_platter_t, rm_temp_gdt);
162     rm->rm_temp_idt_lim = 0;
163     rm->rm_temp_idt_base = 0;
164
165     /*
166     * Since the CPU needs to jump to protected mode using an identity
167     * mapped address, we need to calculate it here.
168     */
169     rm->rm_longmode64_addr = rm_platter_pa +
170         (uint32_t)((uintptr_t)long_mode_64 -
171         (uintptr_t)real_mode_start_cpu);
172 #endif /* __amd64 */
173 }
174
175 static void *
176 mach_cpucontext_alloc_tables(struct cpu *cp)
177 {
178     tss_t *ntss;
179     struct cpu_tables *ct;
180     size_t ctsize;
181
182     /*
183     * Allocate space for stack, tss, gdt and idt. We round the size
184     * allotted for cpu_tables up, so that the TSS is on a unique page.
185     * This is more efficient when running in virtual machines.
186     */
187     ctsize = P2ROUNDUP(sizeof (*ct), PAGESIZE);
188     ct = kmem_zalloc(ctsize, KM_SLEEP);
189     ct = kmem_zalloc(P2ROUNDUP(sizeof (*ct), PAGESIZE), KM_SLEEP);
190     if ((uintptr_t)ct & PAGEOFFSET)
191         panic("mach_cpucontext_alloc_tables: cpu%d misaligned tables",
192             cp->cpu_id);
193
194     ntss = cp->cpu_tss = &ct->ct_tss;
195
196 #if defined(__amd64)
197     uintptr_t va;
198     size_t len;

```

```

197     /*
198     * #DF (double fault).
199     */
200     ntss->tss_ist1 = (uintptr_t)&ct->ct_stack1[sizeof (ct->ct_stack1)];
201     ntss->tss_ist1 = (uint64_t)&ct->ct_stack[sizeof (ct->ct_stack)];
202
203     /*
204     * #NM (non-maskable interrupt)
205     */
206     ntss->tss_ist2 = (uintptr_t)&ct->ct_stack2[sizeof (ct->ct_stack2)];
207
208     /*
209     * #MC (machine check exception / hardware error)
210     */
211     ntss->tss_ist3 = (uintptr_t)&ct->ct_stack3[sizeof (ct->ct_stack3)];
212
213     /*
214     * #DB, #BP debug interrupts and KDI/kmdb
215     */
216     ntss->tss_ist4 = (uintptr_t)&cp->cpu_m.mcpu_kpti_dbg.kf_tr_rsp;
217
218     if (kpti_enable == 1) {
219         /*
220         * #GP, #PF, #SS fault interrupts
221         */
222         ntss->tss_ist5 = (uintptr_t)&cp->cpu_m.mcpu_kptiflt.kf_tr_rsp;
223
224         /*
225         * Used by all other interrupts
226         */
227         ntss->tss_ist6 = (uint64_t)&cp->cpu_m.mcpu_kpti.kf_tr_rsp;
228
229         /*
230         * On AMD64 we need to make sure that all of the pages of the
231         * struct cpu_tables are punched through onto the user CPU for
232         * kpti.
233         *
234         * The final page will always be the TSS, so treat that
235         * separately.
236         */
237         for (va = (uintptr_t)ct, len = ctsize - MMU_PAGESIZE;
238             len >= MMU_PAGESIZE;
239             len -= MMU_PAGESIZE, va += MMU_PAGESIZE) {
240             /* The doublefault stack must be RW */
241             hati_cpu_punchin(cp, va, PROT_READ | PROT_WRITE);
242         }
243         ASSERT3U((uintptr_t)ntss, ==, va);
244         hati_cpu_punchin(cp, (uintptr_t)ntss, PROT_READ);
245     }
246 #elif defined(__i386)
247
248     ntss->tss_esp0 = ntss->tss_esp1 = ntss->tss_esp2 = ntss->tss_esp =
249         (uint32_t)&ct->ct_stack1[sizeof (ct->ct_stack1)];
250         (uint32_t)&ct->ct_stack[sizeof (ct->ct_stack)];
251
252     ntss->tss_ss0 = ntss->tss_ss1 = ntss->tss_ss2 = ntss->tss_ss = KDS_SEL;
253
254     ntss->tss_eip = (uint32_t)cp->cpu_thread->t_pc;
255
256     ntss->tss_cs = KCS_SEL;
257     ntss->tss_ds = ntss->tss_es = KDS_SEL;
258     ntss->tss_fs = KFS_SEL;
259     ntss->tss_gs = KGS_SEL;
260 #endif /* __i386 */

```



```

262     /*
263     * Set I/O bit map offset equal to size of TSS segment limit
264     * for no I/O permission map. This will cause all user I/O
265     * instructions to generate #gp fault.
266     */
267     ntss->tss_bitmapbase = sizeof (*ntss);

269     /*
270     * Setup kernel tss.
271     */
272     set_syssegd((system_desc_t *)&cp->cpu_gdt[GDT_KTSS], cp->cpu_tss,
273               sizeof (*cp->cpu_tss) - 1, SDT_SYSTSS, SEL_KPL);

275     return (ct);
276 }

278 void *
279 mach_cpucontext_xalloc(struct cpu *cp, int optype)
280 {
281     size_t len;
282     struct cpu_tables *ct;
283     rm_platter_t *rm = (rm_platter_t *)rm_platter_va;
284     static int cpu_halt_code_ready;

286     if (optype == MACH_CPUCONTEXT_OP_STOP) {
287         ASSERT(plat_dr_enabled());

289         /*
290         * The WARM_RESET_VECTOR has a limitation that the physical
291         * address written to it must be page-aligned. To work around
292         * this limitation, the CPU stop code has been splitted into
293         * two stages.
294         * The stage 2 code, which implements the real logic to halt
295         * CPUs, is copied to the rm_cpu_halt_code field in the real
296         * mode platter. The stage 1 code, which simply jumps to the
297         * stage 2 code in the rm_cpu_halt_code field, is copied to
298         * rm_code field in the real mode platter and it may be
299         * overwritten after the CPU has been stopped.
300         */
301         if (!cpu_halt_code_ready) {
302             /*
303             * The rm_cpu_halt_code field in the real mode platter
304             * is used by the CPU stop code only. So only copy the
305             * CPU stop stage 2 code into the rm_cpu_halt_code
306             * field on the first call.
307             */
308             len = (size_t)real_mode_stop_cpu_stage2_end -
309                  (size_t)real_mode_stop_cpu_stage2;
310             ASSERT(len <= RM_PLATTER_CPU_HALT_CODE_SIZE);
311             bcopy((caddr_t)real_mode_stop_cpu_stage2,
312                  (caddr_t)rm->rm_cpu_halt_code, len);
313             cpu_halt_code_ready = 1;
314         }

316         /*
317         * The rm_code field in the real mode platter is shared by
318         * the CPU start, CPU stop, CPR and fast reboot code. So copy
319         * the CPU stop stage 1 code into the rm_code field every time.
320         */
321         len = (size_t)real_mode_stop_cpu_stage1_end -
322              (size_t)real_mode_stop_cpu_stage1;
323         ASSERT(len <= RM_PLATTER_CODE_SIZE);
324         bcopy((caddr_t)real_mode_stop_cpu_stage1,
325              (caddr_t)rm->rm_code, len);
326         rm->rm_cpu_halted = 0;

```

```

328         return (cp->cpu_m.mcpu_mach_ctx_ptr);
329     } else if (optype != MACH_CPUCONTEXT_OP_START) {
330         return (NULL);
331     }

333     /*
334     * Only need to allocate tables when starting CPU.
335     * Tables allocated when starting CPU will be reused when stopping CPU.
336     */
337     ct = mach_cpucontext_alloc_tables(cp);
338     if (ct == NULL) {
339         return (NULL);
340     }

342     /* Copy CPU startup code to rm_platter for CPU hot-add operations. */
343     if (plat_dr_enabled()) {
344         bcopy((caddr_t)real_mode_start_cpu, (caddr_t)rm->rm_code,
345              (size_t)real_mode_start_cpu_end -
346              (size_t)real_mode_start_cpu);
347     }

349     /*
350     * Now copy all that we've set up onto the real mode platter
351     * for the real mode code to digest as part of starting the cpu.
352     */
353     rm->rm_idt_base = cp->cpu_idt;
354     rm->rm_idt_lim = sizeof (*cp->cpu_idt) * NIDT - 1;
355     rm->rm_gdt_base = cp->cpu_gdt;
356     rm->rm_gdt_lim = sizeof (*cp->cpu_gdt) * NGDT - 1;

358     /*
359     * CPU needs to access kernel address space after powering on.
360     * When hot-adding CPU at runtime, directly use top level page table
361     * of kas other than the return value of getcr3(). getcr3() returns
362     * current process's top level page table, which may be different from
363     * the one of kas.
364     */
365     rm->rm_pdbr = MAKECR3(kas.a_hat->hat_htable->ht_pfn, PCID_NONE);
366     rm->rm_pdbr = MAKECR3(kas.a_hat->hat_htable->ht_pfn);
367     rm->rm_cpu = cp->cpu_id;

369     /*
370     * We need to mask off any bits set on our boot CPU that can't apply
371     * while the subject CPU is initializing. If appropriate, they are
372     * enabled later on.
373     * For hot-adding CPU at runtime, Machine Check and Performance Counter
374     * should be disabled. They will be enabled on demand after CPU powers
375     * on successfully
376     */
377     rm->rm_cr4 = getcr4();
378     rm->rm_cr4 &= ~(CR4_MCE | CR4_PCE | CR4_PCIDE);
379     rm->rm_cr4 &= ~(CR4_MCE | CR4_PCE);

381     rmp_gdt_init(rm);

383     return (ct);
384 }

```

unchanged portion omitted

new/usr/src/uts/i86pc/os/mp_startup.c

1

```
*****
51721 Fri Apr 6 17:25:02 2018
new/usr/src/uts/i86pc/os/mp_startup.c
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */

22 /*
23  * Copyright (c) 1992, 2010, Oracle and/or its affiliates. All rights reserved.
24 */
25 /*
26  * Copyright (c) 2010, Intel Corporation.
27  * All rights reserved.
28 */
29 /*
30  * Copyright 2018 Joyent, Inc.
31  * Copyright 2016 Joyent, Inc.
32  * Copyright 2013 Nexenta Systems, Inc. All rights reserved.
33 */

34 #include <sys/types.h>
35 #include <sys/thread.h>
36 #include <sys/cpuvar.h>
37 #include <sys/cpu.h>
38 #include <sys/t_lock.h>
39 #include <sys/param.h>
40 #include <sys/proc.h>
41 #include <sys/disp.h>
42 #include <sys/class.h>
43 #include <sys/cmn_err.h>
44 #include <sys/debug.h>
45 #include <sys/note.h>
46 #include <sys/asm_linkage.h>
47 #include <sys/x_call.h>
48 #include <sys/system.h>
49 #include <sys/var.h>
50 #include <sys/vtrace.h>
51 #include <vm/hat.h>
52 #include <vm/as.h>
53 #include <vm/seg_kmem.h>
54 #include <vm/seg_kp.h>
55 #include <sys/segments.h>
56 #include <sys/kmem.h>
57 #include <sys/stack.h>
58 #include <sys/smp_impldefs.h>
```

new/usr/src/uts/i86pc/os/mp_startup.c

2

```
59 #include <sys/x86_archext.h>
60 #include <sys/machsystem.h>
61 #include <sys/traptrace.h>
62 #include <sys/clock.h>
63 #include <sys/cpc_impl.h>
64 #include <sys/pg.h>
65 #include <sys/cmt.h>
66 #include <sys/dtrace.h>
67 #include <sys/archsystem.h>
68 #include <sys/fp.h>
69 #include <sys/reboot.h>
70 #include <sys/kdi_machimpl.h>
71 #include <vm/hat_i86.h>
72 #include <vm/vm_dep.h>
73 #include <sys/memnode.h>
74 #include <sys/pci_cfgspace.h>
75 #include <sys/mach_mmu.h>
76 #include <sys/sysmacros.h>
77 #if defined(__xpv)
78 #include <sys/hypervisor.h>
79 #endif
80 #include <sys/cpu_module.h>
81 #include <sys/ontrap.h>

82 struct cpu cpus[1] __aligned(MMU_PAGESIZE);
83 struct cpu *cpu[NCPU] = {&cpus[0]};
84 struct cpu *cpu_free_list;
85 struct cpu *cpu_core[NCPU];
86 struct cpu cpus[1]; /* CPU data */
87 struct cpu *cpu[NCPU] = {&cpus[0]}; /* pointers to all CPUs */
88 struct cpu *cpu_free_list; /* list for released CPUs */
89 struct cpu cpu_core[NCPU]; /* cpu_core structures */

90 #define cpu_next_free cpu_prev

91 /*
92  * Useful for disabling MP bring-up on a MP capable system.
93 */
94 int use_mp = 1;

95 /*
96  * to be set by a PSM to indicate what cpus
97  * are sitting around on the system.
98 */
99 cpuset_t mp_cpus;

101 /*
102  * This variable is used by the hat layer to decide whether or not
103  * critical sections are needed to prevent race conditions. For sun4m,
104  * this variable is set once enough MP initialization has been done in
105  * order to allow cross calls.
106 */
107 int flushes_require_xcalls;

109 cpuset_t cpu_ready_set; /* initialized in startup() */

111 static void mp_startup_boot(void);
112 static void mp_startup_hotplug(void);

114 static void cpu_sep_enable(void);
115 static void cpu_sep_disable(void);
116 static void cpu_asysc_enable(void);
117 static void cpu_asysc_disable(void);

119 /*
120  * Init CPU info - get CPU type info for processor_info system call.
```

```

121 */
122 void
123 init_cpu_info(struct cpu *cp)
124 {
125     processor_info_t *pi = &cp->cpu_type_info;
126
127     /*
128      * Get clock-frequency property for the CPU.
129      */
130     pi->pi_clock = cpu_freq;
131
132     /*
133      * Current frequency in Hz.
134      */
135     cp->cpu_curr_clock = cpu_freq_hz;
136
137     /*
138      * Supported frequencies.
139      */
140     if (cp->cpu_supp_freqs == NULL) {
141         cpu_set_supp_freqs(cp, NULL);
142     }
143
144     (void) strcpy(pi->pi_processor_type, "i386");
145     if (fpu_exists)
146         (void) strcpy(pi->pi_fputypes, "i387 compatible");
147
148     cp->cpu_idstr = kmem_zalloc(CPU_IDSTRLEN, KM_SLEEP);
149     cp->cpu_brandstr = kmem_zalloc(CPU_IDSTRLEN, KM_SLEEP);
150
151     /*
152      * If called for the BSP, cp is equal to current CPU.
153      * For non-BSPs, cpuid info of cp is not ready yet, so use cpuid info
154      * of current CPU as default values for cpu_idstr and cpu_brandstr.
155      * They will be corrected in mp_startup_common() after cpuid_pass1()
156      * has been invoked on target CPU.
157      */
158     (void) cpuid_getidstr(CPU, cp->cpu_idstr, CPU_IDSTRLEN);
159     (void) cpuid_getbrandstr(CPU, cp->cpu_brandstr, CPU_IDSTRLEN);
160 }
161
162 /*
163  * Configure syscall support on this CPU.
164  */
165 /*ARGSUSED*/
166 void
167 init_cpu_syscall(struct cpu *cp)
168 {
169     kpreempt_disable();
170
171 #if defined(__amd64)
172     if (is_x86_feature(x86_featureset, X86FSET_MSR) &&
173         is_x86_feature(x86_featureset, X86FSET_ASYSC)) {
174         uint64_t flags;
175
176 #if !defined(__xpv)
177 #if !defined(__lint)
178         /*
179          * The syscall instruction imposes a certain ordering on
180          * segment selectors, so we double-check that ordering
181          * here.
182          */
183         CTASSERT(KDS_SEL == KCS_SEL + 8);
184         CTASSERT(UDS_SEL == U32CS_SEL + 8);
185         CTASSERT(UCS_SEL == U32CS_SEL + 16);
186         ASSERT(KDS_SEL == KCS_SEL + 8);

```

```

183         ASSERT(UDS_SEL == U32CS_SEL + 8);
184         ASSERT(UCS_SEL == U32CS_SEL + 16);
185     }
186 #endif
187
188     /*
189      * Turn syscall/sysret extensions on.
190      */
191     cpu_asysc_enable();
192
193     /*
194      * Program the magic registers ..
195      */
196     wrmsr(MSR_AMD_STAR,
197           ((uint64_t)(U32CS_SEL << 16 | KCS_SEL)) << 32);
198     if (kpti_enable == 1) {
199         wrmsr(MSR_AMD_LSTAR,
200             (uint64_t)(uintptr_t)tr_sys_syscall);
201         wrmsr(MSR_AMD_CSTAR,
202             (uint64_t)(uintptr_t)tr_sys_syscall32);
203     } else {
204         wrmsr(MSR_AMD_LSTAR,
205             (uint64_t)(uintptr_t)sys_syscall);
206         wrmsr(MSR_AMD_CSTAR,
207             (uint64_t)(uintptr_t)sys_syscall32);
208     }
209     wrmsr(MSR_AMD_LSTAR, (uint64_t)(uintptr_t)sys_syscall);
210     wrmsr(MSR_AMD_CSTAR, (uint64_t)(uintptr_t)sys_syscall32);
211
212     /*
213      * This list of flags is masked off the incoming
214      * %rfl when we enter the kernel.
215      */
216     flags = PS_IE | PS_T;
217     if (is_x86_feature(x86_featureset, X86FSET_SMAP) == B_TRUE)
218         flags |= PS_ACHK;
219     wrmsr(MSR_AMD_SFMASK, flags);
220 #endif
221
222     /*
223      * On 32-bit kernels, we use sysenter/sysexit because it's too
224      * hard to use syscall/sysret, and it is more portable anyway.
225      *
226      * On 64-bit kernels on Nocona machines, the 32-bit syscall
227      * variant isn't available to 32-bit applications, but sysenter is.
228      */
229     if (is_x86_feature(x86_featureset, X86FSET_MSR) &&
230         is_x86_feature(x86_featureset, X86FSET_SEP)) {
231
232 #if !defined(__xpv)
233 #if !defined(__lint)
234         /*
235          * The sysenter instruction imposes a certain ordering on
236          * segment selectors, so we double-check that ordering
237          * here. See "sysenter" in Intel document 245471-012, "IA-32
238          * Intel Architecture Software Developer's Manual Volume 2:
239          * Instruction Set Reference"
240          */
241         CTASSERT(KDS_SEL == KCS_SEL + 8);
242         ASSERT(KDS_SEL == KCS_SEL + 8);
243
244         CTASSERT(U32CS_SEL == ((KCS_SEL + 16) | 3));
245         CTASSERT(UDS_SEL == U32CS_SEL + 8);
246         ASSERT32(UCS_SEL == ((KCS_SEL + 16) | 3));
247         ASSERT32(UDS_SEL == UCS_SEL + 8);

```

```

233     ASSERT64(U32CS_SEL == ((KCS_SEL + 16) | 3));
234     ASSERT64(UDS_SEL == U32CS_SEL + 8);
237 #endif

239     cpu_sep_enable();

241     /*
242      * resume() sets this value to the base of the threads stack
243      * via a context handler.
244      */
245     wrmsr(MSR_INTC_SEP_ESP, 0);

247     if (kpti_enable == 1) {
248         wrmsr(MSR_INTC_SEP_EIP,
249             (uint64_t)(uintptr_t)tr_sys_sysenter);
250     } else {
251         wrmsr(MSR_INTC_SEP_EIP,
252             (uint64_t)(uintptr_t)sys_sysenter);
253     }
254     wrmsr(MSR_INTC_SEP_EIP, (uint64_t)(uintptr_t)sys_sysenter);
255 }

256     kpreempt_enable();
257 }

```

unchanged portion omitted

```

275 #endif /* !defined(__xpv) */

277 /*
278  * Multiprocessor initialization.
279  * Allocate and initialize the cpu structure, TRAPTRACE buffer, and the
280  * startup and idle threads for the specified CPU.
281  * Parameter boot is true for boot time operations and is false for CPU
282  * DR operations.
283  */
284 static struct cpu *
285 mp_cpu_configure_common(int cpun, boolean_t boot)
286 {
287     struct cpu *cp;
288     kthread_id_t tp;
289     caddr_t sp;
290     proc_t *procp;
291 #if !defined(__xpv)
292     extern int idle_cpu_prefer_mwait;
293     extern void cpu_idle_mwait();
294 #endif
295     extern void idle();
296     extern void cpu_idle();
297
299 #ifdef TRAPTRACE
300     trap_trace_ctl_t *ttc = &trap_trace_ctl[cpun];
301 #endif

303     ASSERT(MUTEX_HELD(&cpu_lock));
304     ASSERT(cpun < NCPU && cpu[cpun] == NULL);

306     if (cpu_free_list == NULL) {
307         cp = kmem_zalloc(sizeof (*cp), KM_SLEEP);
308     } else {
309         cp = cpu_free_list;
310         cpu_free_list = cp->cpu_next_free;
311     }

313     cp->cpu_m.mcpu_istamp = cpun << 16;

315     /* Create per CPU specific threads in the process p0. */

```

```

316     procp = &p0;

318     /*
319      * Initialize the dispatcher first.
320      */
321     disp_cpu_init(cp);

323     cpu_vm_data_init(cp);

325     /*
326      * Allocate and initialize the startup thread for this CPU.
327      * Interrupt and process switch stacks get allocated later
328      * when the CPU starts running.
329      */
330     tp = thread_create(NULL, 0, NULL, NULL, 0, procp,
331         TS_STOPPED, maxclsypri);

333     /*
334      * Set state to TS_ONPROC since this thread will start running
335      * as soon as the CPU comes online.
336      *
337      * All the other fields of the thread structure are setup by
338      * thread_create().
339      */
340     THREAD_ONPROC(tp, cp);
341     tp->t_preempt = 1;
342     tp->t_bound_cpu = cp;
343     tp->t_affinitycnt = 1;
344     tp->t_cpu = cp;
345     tp->t_disp_queue = cp->cpu_disp;

347     /*
348      * Setup thread to start in mp_startup_common.
349      */
350     sp = tp->t_stk;
351     tp->t_sp = (uintptr_t)(sp - MINFRAME);
352 #if defined(__amd64)
353     tp->t_sp -= STACK_ENTRY_ALIGN; /* fake a call */
354 #endif

355     /*
356      * Setup thread start entry point for boot or hotplug.
357      */
358     if (boot) {
359         tp->t_pc = (uintptr_t)mp_startup_boot;
360     } else {
361         tp->t_pc = (uintptr_t)mp_startup_hotplug;
362     }

364     cp->cpu_id = cpun;
365     cp->cpu_self = cp;
366     cp->cpu_thread = tp;
367     cp->cpu_lwp = NULL;
368     cp->cpu_dispthread = tp;
369     cp->cpu_dispatch_pri = DISP_PRI0(tp);

371     /*
372      * cpu_base_spl must be set explicitly here to prevent any blocking
373      * operations in mp_startup_common from causing the spl of the cpu
374      * to drop to 0 (allowing device interrupts before we're ready) in
375      * resume().
376      * cpu_base_spl MUST remain at LOCK_LEVEL until the cpu is CPU_READY.
377      * As an extra bit of security on DEBUG kernels, this is enforced with
378      * an assertion in mp_startup_common() -- before cpu_base_spl is set
379      * to its proper value.
380      */
381     cp->cpu_base_spl = ipltospl(LOCK_LEVEL);

```

```

383      /*
384       * Now, initialize per-CPU idle thread for this CPU.
385       */
386      tp = thread_create(NULL, PAGESIZE, idle, NULL, 0, procp, TS_ONPROC, -1);

388      cp->cpu_idle_thread = tp;

390      tp->t_preempt = 1;
391      tp->t_bound_cpu = cp;
392      tp->t_affinitycnt = 1;
393      tp->t_cpu = cp;
394      tp->t_disp_queue = cp->cpu_disp;

396      /*
397       * Bootstrap the CPU's PG data
398       */
399      pg_cpu_bootstrap(cp);

401      /*
402       * Perform CPC initialization on the new CPU.
403       */
404      kcpc_hw_init(cp);

406      /*
407       * Allocate virtual addresses for cpu_caddr1 and cpu_caddr2
408       * for each CPU.
409       */
410      setup_vaddr_for_ppcopy(cp);

412      /*
413       * Allocate page for new GDT and initialize from current GDT.
414       */
415      #if !defined(__lint)
416          ASSERT((sizeof (*cp->cpu_gdt) * NGDT) <= PAGESIZE);
417      #endif
418      cp->cpu_gdt = kmem_zalloc(PAGESIZE, KM_SLEEP);
419      bcopy(CPU->cpu_gdt, cp->cpu_gdt, (sizeof (*cp->cpu_gdt) * NGDT));

421      #if defined(__i386)
422          /*
423           * setup kernel %gs.
424           */
425          set_usegd(&cp->cpu_gdt[GDT_GS], cp, sizeof (struct cpu) -1, SDT_MEMRWA,
426                  SEL_KPL, 0, 1);
427      #endif

429      /*
430       * Allocate pages for the CPU LDT.
431       * If we have more than one node, each cpu gets a copy of IDT
432       * local to its node. If this is a Pentium box, we use cpu 0's
433       * IDT. cpu 0's IDT has been made read-only to workaround the
434       * cmpxchgl register bug
435       */
432      cp->cpu_m.mcpu_ldt = kmem_zalloc(LDT_CPU_SIZE, KM_SLEEP);
433      cp->cpu_m.mcpu_ldt_len = 0;

435      /*
436       * Allocate a per-CPU IDT and initialize the new IDT to the currently
437       * running CPU.
438       */
439      #if !defined(__lint)
440          if (system_hardware.hd_nodes && x86_type != X86_TYPE_P5) {
441              ASSERT((sizeof (*CPU->cpu_idt) * NIDT) <= PAGESIZE);
442              cp->cpu_idt = kmem_alloc(PAGESIZE, KM_SLEEP);

```

```

430          cp->cpu_idt = kmem_zalloc(PAGESIZE, KM_SLEEP);
431          bcopy(CPU->cpu_idt, cp->cpu_idt, PAGESIZE);
432      } else {
433          cp->cpu_idt = CPU->cpu_idt;
434      }

445      /*
446       * alloc space for cpuid info
447       */
448      cpuid_alloc_space(cp);
449      #if !defined(__xpv)
450          if (is_x86_feature(x86_featureset, X86FSET_MWAIT) &&
451              idle_cpu_prefer_mwait) {
452              cp->cpu_m.mcpu_mwait = cpuid_mwait_alloc(cp);
453              cp->cpu_m.mcpu_idle_cpu = cpu_idle_mwait;
454          } else
455      #endif
456          cp->cpu_m.mcpu_idle_cpu = cpu_idle;

458      init_cpu_info(cp);

460      #if !defined(__xpv)
461          init_cpu_id_gdt(cp);
462      #endif

464      /*
465       * alloc space for ucode_info
466       */
467      ucode_alloc_space(cp);
468      xc_init_cpu(cp);
469      hat_cpu_online(cp);

471      #ifdef TRAPTRACE
472          /*
473           * If this is a TRAPTRACE kernel, allocate TRAPTRACE buffers
474           */
475          ttc->ttc_first = (uintptr_t)kmem_zalloc(trap_trace_bufsize, KM_SLEEP);
476          ttc->ttc_next = ttc->ttc_first;
477          ttc->ttc_limit = ttc->ttc_first + trap_trace_bufsize;
478      #endif

480      /*
481       * Record that we have another CPU.
482       */
483      /*
484       * Initialize the interrupt threads for this CPU
485       */
486      cpu_intr_alloc(cp, NINTR_THREADS);

488      cp->cpu_flags = CPU_OFFLINE | CPU_QUIESCED | CPU_POWEROFF;
489      cpu_set_state(cp);

491      /*
492       * Add CPU to list of available CPUs. It'll be on the active list
493       * after mp_startup_common().
494       */
495      cpu_add_unit(cp);

497      return (cp);
498  }

500  /*
501  * Undo what was done in mp_cpu_configure_common
502  */
503  static void
504  mp_cpu_unconfigure_common(struct cpu *cp, int error)

```

```

505 {
506     ASSERT(MUTEX_HELD(&cpu_lock));

508     /*
509      * Remove the CPU from the list of available CPUs.
510      */
511     cpu_del_unit(cp->cpu_id);

513     if (error == ETIMEDOUT) {
514         /*
515          * The cpu was started, but never *seemed* to run any
516          * code in the kernel; it's probably off spinning in its
517          * own private world, though with potential references to
518          * our kmem-allocated IDTs and GDTs (for example).
519          *
520          * Worse still, it may actually wake up some time later,
521          * so rather than guess what it might or might not do, we
522          * leave the fundamental data structures intact.
523          */
524         cp->cpu_flags = 0;
525         return;
526     }

528     /*
529      * At this point, the only threads bound to this CPU should
530      * special per-cpu threads: it's idle thread, it's pause threads,
531      * and it's interrupt threads. Clean these up.
532      */
533     cpu_destroy_bound_threads(cp);
534     cp->cpu_idle_thread = NULL;

536     /*
537      * Free the interrupt stack.
538      */
539     segkp_release(segkp,
540         cp->cpu_intr_stack - (INTR_STACK_SIZE - SA(MINFRAME)));
541     cp->cpu_intr_stack = NULL;

543 #ifdef TRAPTRACE
544     /*
545      * Discard the trap trace buffer
546      */
547     {
548         trap_trace_ctl_t *ttc = &trap_trace_ctl[cp->cpu_id];

550         kmem_free((void *)ttc->ttc_first, trap_trace_bufsize);
551         ttc->ttc_first = NULL;
552     }
553 #endif

555     hat_cpu_offline(cp);

557     ucode_free_space(cp);

559     /* Free CPU ID string and brand string. */
560     if (cp->cpu_idstr) {
561         kmem_free(cp->cpu_idstr, CPU_IDSTRLEN);
562         cp->cpu_idstr = NULL;
563     }
564     if (cp->cpu_brandstr) {
565         kmem_free(cp->cpu_brandstr, CPU_IDSTRLEN);
566         cp->cpu_brandstr = NULL;
567     }

569 #if !defined(__xpv)
570     if (cp->cpu_m.mcpu_mwait != NULL) {

```

```

571         cpuid_mwait_free(cp);
572         cp->cpu_m.mcpu_mwait = NULL;
573     }
574 #endif
575     cpuid_free_space(cp);

577     if (cp->cpu_idt != CPU->cpu_idt)
578         kmem_free(cp->cpu_idt, PAGESIZE);
579     cp->cpu_idt = NULL;

581     kmem_free(cp->cpu_m.mcpu_ldt, LDT_CPU_SIZE);
582     cp->cpu_m.mcpu_ldt = NULL;
583     cp->cpu_m.mcpu_ldt_len = 0;

585     kmem_free(cp->cpu_gdt, PAGESIZE);
586     cp->cpu_gdt = NULL;

588     if (cp->cpu_supp_freqs != NULL) {
589         size_t len = strlen(cp->cpu_supp_freqs) + 1;
590         kmem_free(cp->cpu_supp_freqs, len);
591         cp->cpu_supp_freqs = NULL;
592     }

594     teardown_vaddr_for_ppcopy(cp);

596     kcpc_hw_fini(cp);

598     cp->cpu_dispthread = NULL;
599     cp->cpu_thread = NULL; /* discarded by cpu_destroy_bound_threads() */

601     cpu_vm_data_destroy(cp);

603     xc_fini_cpu(cp);
604     disp_cpu_fini(cp);

606     ASSERT(cp != CPU0);
607     bzero(cp, sizeof (*cp));
608     cp->cpu_next_free = cpu_free_list;
609     cpu_free_list = cp;
610 }

    unchanged_portion_omitted

1663 /*
1664  * Startup function for 'other' CPUs (besides boot cpu).
1665  * Called from real_mode_start.
1666  *
1667  * WARNING: until CPU_READY is set, mp_startup_common and routines called by
1668  * mp_startup_common should not call routines (e.g. kmem_free) that could call
1669  * hat_unload which requires CPU_READY to be set.
1670  */
1671 static void
1672 mp_startup_common(boolean_t boot)
1673 {
1674     cpu_t *cp = CPU;
1675     uchar_t new_x86_featureset[BT_SIZEOFMAP(NUM_X86_FEATURES)];
1676     extern void cpu_event_init_cpu(cpu_t *);

1678     /*
1679      * We need to get TSC on this proc synced (i.e., any delta
1680      * from cpu0 accounted for) as soon as we can, because many
1681      * many things use gethrtime/pc_gethrestime, including
1682      * interrupts, cmn_err, etc. Before we can do that, we want to
1683      * clear TSC if we're on a buggy Sandy/Ivy Bridge CPU, so do that
1684      * right away.
1685      */
1686     bzero(new_x86_featureset, BT_SIZEOFMAP(NUM_X86_FEATURES));

```

```

1687     cpuid_pass1(cp, new_x86_featureset);
1689
1689 if (boot && get_hwenv() == HW_NATIVE &&
1690     cpuid_getvendor(CPU) == X86_VENDOR_Intel &&
1691     cpuid_getfamily(CPU) == 6 &&
1692     (cpuid_getmodel(CPU) == 0x2d || cpuid_getmodel(CPU) == 0x3e) &&
1693     is_x86_feature(new_x86_featureset, X86FSET_TSC)) {
1694     (void) wrmsr(REG_TSC, 0UL);
1695 }
1697
1697 /* Let the control CPU continue into tsc_sync_master() */
1698 mp_startup_signal(&procset_slave, cp->cpu_id);
1700 #ifndef __xpv
1701 if (tsc_gethrtime_enable)
1702     tsc_sync_slave();
1703 #endif
1705
1705 /*
1706  * Once this was done from assembly, but it's safer here; if
1707  * it blocks, we need to be able to swtch() to and from, and
1708  * since we get here by calling t_pc, we need to do that call
1709  * before swtch() overwrites it.
1710  */
1711 (void) (*ap_mlsetup)();
1713 #ifndef __xpv
1714 /*
1715  * Program this cpu's PAT
1716  */
1717 pat_sync();
1718 #endif
1720
1720 /*
1721  * Set up TSC_AUX to contain the cpuid for this processor
1722  * for the rdtscl instruction.
1723  */
1724 if (is_x86_feature(x86_featureset, X86FSET_TSCP))
1725     (void) wrmsr(MSR_AMD_TSCAUX, cp->cpu_id);
1727
1727 /*
1728  * Initialize this CPU's syscall handlers
1729  */
1730 init_cpu_syscall(cp);
1732
1732 /*
1733  * Enable interrupts with spl set to LOCK_LEVEL. LOCK_LEVEL is the
1734  * highest level at which a routine is permitted to block on
1735  * an adaptive mutex (allows for cpu poke interrupt in case
1736  * the cpu is blocked on a mutex and halts). Setting LOCK_LEVEL blocks
1737  * device interrupts that may end up in the hat layer issuing cross
1738  * calls before CPU_READY is set.
1739  */
1740 splx(ipלטospl(LOCK_LEVEL));
1741 sti();
1743
1743 /*
1744  * Do a sanity check to make sure this new CPU is a sane thing
1745  * to add to the collection of processors running this system.
1746  *
1747  * XXX Clearly this needs to get more sophisticated, if x86
1748  * systems start to get built out of heterogeneous CPUs; as is
1749  * likely to happen once the number of processors in a configuration
1750  * gets large enough.
1751  */
1752 if (compare_x86_featureset(x86_featureset, new_x86_featureset) ==

```

```

1753     B_FALSE) {
1754     cmn_err(CE_CONT, "cpu%d: featureset\n", cp->cpu_id);
1755     print_x86_featureset(new_x86_featureset);
1756     cmn_err(CE_WARN, "cpu%d feature mismatch", cp->cpu_id);
1757 }
1759
1759 /*
1760  * There exists a small subset of systems which expose differing
1761  * MWAIT/MONITOR support between CPUs. If MWAIT support is absent from
1762  * the boot CPU, but is found on a later CPU, the system continues to
1763  * operate as if no MWAIT support is available.
1764  *
1765  * The reverse case, where MWAIT is available on the boot CPU but not
1766  * on a subsequently initialized CPU, is not presently allowed and will
1767  * result in a panic.
1768  */
1769 if (is_x86_feature(x86_featureset, X86FSET_MWAIT) !=
1770     is_x86_feature(new_x86_featureset, X86FSET_MWAIT)) {
1771     if (!is_x86_feature(x86_featureset, X86FSET_MWAIT)) {
1772         remove_x86_feature(new_x86_featureset, X86FSET_MWAIT);
1773     } else {
1774         panic("unsupported mixed cpu mwait support detected");
1775     }
1776 }
1778
1778 /*
1779  * We could be more sophisticated here, and just mark the CPU
1780  * as "faulted" but at this point we'll opt for the easier
1781  * answer of dying horribly. Provided the boot cpu is ok,
1782  * the system can be recovered by booting with use_mp set to zero.
1783  */
1784 if (workaround_errata(cp) != 0)
1785     panic("critical workaround(s) missing for cpu%d", cp->cpu_id);
1787
1787 /*
1788  * We can touch cpu_flags here without acquiring the cpu_lock here
1789  * because the cpu_lock is held by the control CPU which is running
1790  * mp_start_cpu_common().
1791  * Need to clear CPU QUIESCED flag before calling any function which
1792  * may cause thread context switching, such as kmem_alloc() etc.
1793  * The idle thread checks for CPU QUIESCED flag and loops for ever if
1794  * it's set. So the startup thread may have no chance to switch back
1795  * again if it's switched away with CPU QUIESCED set.
1796  */
1797 cp->cpu_flags &= ~(CPU_POWEROFF | CPU QUIESCED);
1799
1799 enable_pcid();
1801
1801 /*
1802  * Setup this processor for XSAVE.
1803  */
1804 if (fp_save_mech == FP_XSAVE) {
1805     xsave_setup_msr(cp);
1806 }
1808
1808 cpuid_pass2(cp);
1809 cpuid_pass3(cp);
1810 cpuid_pass4(cp, NULL);
1812
1812 /*
1813  * Correct cpu_idstr and cpu_brandstr on target CPU after
1814  * cpuid_pass1() is done.
1815  */
1816 (void) cpuid_getidstr(cp, cp->cpu_idstr, CPU_IDSTRLEN);
1817 (void) cpuid_getbrandstr(cp, cp->cpu_brandstr, CPU_IDSTRLEN);

```

```

1819     cp->cpu_flags |= CPU_RUNNING | CPU_READY | CPU_EXISTS;
1821     post_startup_cpu_fixups();
1823     cpu_event_init_cpu(cp);
1825     /*
1826     * Enable preemption here so that contention for any locks acquired
1827     * later in mp_startup_common may be preempted if the thread owning
1828     * those locks is continuously executing on other CPUs (for example,
1829     * this CPU must be preemptible to allow other CPUs to pause it during
1830     * their startup phases). It's safe to enable preemption here because
1831     * the CPU state is pretty-much fully constructed.
1832     */
1833     curthread->t_preempt = 0;
1835     /* The base spl should still be at LOCK LEVEL here */
1836     ASSERT(cp->cpu_base_spl == ipltospl(LOCK_LEVEL));
1837     set_base_spl(); /* Restore the spl to its proper value */
1839     pghw_physid_create(cp);
1840     /*
1841     * Delegate initialization tasks, which need to access the cpu_lock,
1842     * to mp_start_cpu_common() because we can't acquire the cpu_lock here
1843     * during CPU DR operations.
1844     */
1845     mp_startup_signal(&procset_slave, cp->cpu_id);
1846     mp_startup_wait(&procset_master, cp->cpu_id);
1847     pg_cmt_cpu_startup(cp);
1849     if (boot) {
1850         mutex_enter(&cpu_lock);
1851         cp->cpu_flags &= ~CPU_OFFLINE;
1852         cpu_enable_intr(cp);
1853         cpu_add_active(cp);
1854         mutex_exit(&cpu_lock);
1855     }
1857     /* Enable interrupts */
1858     (void) spl0();
1860     /*
1861     * Fill out cpu_ucode_info. Update microcode if necessary.
1862     */
1863     ucode_check(cp);
1865 #ifndef __xpv
1866     {
1867         /*
1868         * Set up the CPU module for this CPU. This can't be done
1869         * before this CPU is made CPU_READY, because we may (in
1870         * heterogeneous systems) need to go load another CPU module.
1871         * The act of attempting to load a module may trigger a
1872         * cross-call, which will ASSERT unless this cpu is CPU_READY.
1873         */
1874         cmi_hdl_t hdl;
1876         if ((hdl = cmi_init(CMI_HDL_NATIVE, cmi_ntv_hwchipid(CPU),
1877             cmi_ntv_hwcoid(CPU), cmi_ntv_hwstrandid(CPU))) != NULL) {
1878             if (is_x86_feature(x86_featureset, X86FSET_MCA))
1879                 cmi_mca_init(hdl);
1880             cp->cpu_m.mcpu_cmi_hdl = hdl;
1881         }
1882     }
1883 #endif /* __xpv */

```

```

1885     if (boothowto & RB_DEBUG)
1886         kdi_cpu_init();
1888     /*
1889     * Setting the bit in cpu_ready_set must be the last operation in
1890     * processor initialization; the boot CPU will continue to boot once
1891     * it sees this bit set for all active CPUs.
1892     */
1893     CPUSET_ATOMIC_ADD(cpu_ready_set, cp->cpu_id);
1895     (void) mach_cpu_create_device_node(cp, NULL);
1897     cmn_err(CE_CONT, "?cpu%d: %s\n", cp->cpu_id, cp->cpu_idstr);
1898     cmn_err(CE_CONT, "?cpu%d: %s\n", cp->cpu_id, cp->cpu_brandstr);
1899     cmn_err(CE_CONT, "?cpu%d initialization complete - online\n",
1900         cp->cpu_id);
1902     /*
1903     * Now we are done with the startup thread, so free it up.
1904     */
1905     thread_exit();
1906     panic("mp_startup: cannot return");
1907     /*NOTREACHED*/
1908 }

```

unchanged portion omitted


```

*****
88262 Fri Apr 6 17:25:03 2018
new/usr/src/uts/i86pc/os/startup.c
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */

22 /*
23  * Copyright (c) 1993, 2010, Oracle and/or its affiliates. All rights reserved.
24  * Copyright 2012 DEY Storage Systems, Inc. All rights reserved.
25  * Copyright 2017 Nexenta Systems, Inc.
26  * Copyright (c) 2018 Joyent, Inc.
26  * Copyright 2015 Joyent, Inc.
27  * Copyright (c) 2015 by Delphix. All rights reserved.
28  */
29 /*
30  * Copyright (c) 2010, Intel Corporation.
31  * All rights reserved.
32  */

34 #include <sys/types.h>
35 #include <sys/t_lock.h>
36 #include <sys/param.h>
37 #include <sys/sysmacros.h>
38 #include <sys/signal.h>
39 #include <sys/system.h>
40 #include <sys/user.h>
41 #include <sys/mman.h>
42 #include <sys/vm.h>
43 #include <sys/conf.h>
44 #include <sys/avintr.h>
45 #include <sys/autoconf.h>
46 #include <sys/disp.h>
47 #include <sys/class.h>
48 #include <sys/bitmap.h>

50 #include <sys/privregs.h>

52 #include <sys/proc.h>
53 #include <sys/buf.h>
54 #include <sys/kmem.h>
55 #include <sys/mem.h>
56 #include <sys/kstat.h>

58 #include <sys/reboot.h>

```

```

60 #include <sys/cred.h>
61 #include <sys/vnode.h>
62 #include <sys/file.h>

64 #include <sys/procfs.h>

66 #include <sys/vfs.h>
67 #include <sys/cmn_err.h>
68 #include <sys/utsname.h>
69 #include <sys/debug.h>
70 #include <sys/kdi.h>

72 #include <sys/dumphdr.h>
73 #include <sys/bootconf.h>
74 #include <sys/memlist_plat.h>
75 #include <sys/varargs.h>
76 #include <sys/promif.h>
77 #include <sys/modctl.h>

79 #include <sys/sunddi.h>
80 #include <sys/sunndi.h>
81 #include <sys/ndi_impldefs.h>
82 #include <sys/ddidmareq.h>
83 #include <sys/psw.h>
84 #include <sys/regset.h>
85 #include <sys/clock.h>
86 #include <sys/pte.h>
87 #include <sys/tss.h>
88 #include <sys/stack.h>
89 #include <sys/trap.h>
90 #include <sys/fp.h>
91 #include <vm/kboot_mmu.h>
92 #include <vm/anon.h>
93 #include <vm/as.h>
94 #include <vm/page.h>
95 #include <vm/seg.h>
96 #include <vm/seg_dev.h>
97 #include <vm/seg_kmem.h>
98 #include <vm/seg_kpm.h>
99 #include <vm/seg_map.h>
100 #include <vm/seg_vn.h>
101 #include <vm/seg_kp.h>
102 #include <sys/memnode.h>
103 #include <vm/vm_dep.h>
104 #include <sys/thread.h>
105 #include <sys/sysconf.h>
106 #include <sys/vm_machparam.h>
107 #include <sys/archsystem.h>
108 #include <sys/machsystem.h>
109 #include <vm/hat.h>
110 #include <vm/hat_i86.h>
111 #include <sys/pmem.h>
112 #include <sys/smp_impldefs.h>
113 #include <sys/x86_archext.h>
114 #include <sys/cpuvar.h>
115 #include <sys/segments.h>
116 #include <sys/clconf.h>
117 #include <sys/kobj.h>
118 #include <sys/kobj_lex.h>
119 #include <sys/cpc_impl.h>
120 #include <sys/cpu_module.h>
121 #include <sys/smbios.h>
122 #include <sys/debug_info.h>
123 #include <sys/bootinfo.h>
124 #include <sys/ddi_periodic.h>

```

```

125 #include <sys/systeminfo.h>
126 #include <sys/multiboot.h>
127 #include <sys/ramdisk.h>

129 #ifdef __xpv

131 #include <sys/hypervisor.h>
132 #include <sys/xen_mmu.h>
133 #include <sys/evtchn_impl.h>
134 #include <sys/gnttab.h>
135 #include <sys/xpv_panic.h>
136 #include <xen/sys/xenbus_comms.h>
137 #include <xen/public/physdev.h>

139 extern void xen_late_startup(void);

141 struct xen_evt_data cpu0_evt_data;

143 #else /* __xpv */
144 #include <sys/memlist_impl.h>

146 extern void mem_config_init(void);
147 #endif /* __xpv */

149 extern void progressbar_init(void);
150 extern void brand_init(void);
151 extern void pcf_init(void);
152 extern void pg_init(void);
153 extern void ssp_init(void);

155 extern int size_pse_array(pgcnt_t, int);

157 #if defined(_SOFT_HOSTID)

159 #include <sys/rtc.h>

161 static int32_t set_soft_hostid(void);
162 static char hostid_file[] = "/etc/hostid";

164 #endif

166 void *gfx_devinfo_list;

168 #if defined(__amd64) && !defined(__xpv)
169 extern void immu_startup(void);
170 #endif

172 /*
173  * XXX make declaration below "static" when drivers no longer use this
174  * interface.
175  */
176 extern caddr_t p0_va; /* Virtual address for accessing physical page 0 */

178 /*
179  * segkp
180  */
181 extern int segkp_fromheap;

183 static void kvm_init(void);
184 static void startup_init(void);
185 static void startup_memlist(void);
186 static void startup_kmem(void);
187 static void startup_modules(void);
188 static void startup_vm(void);
189 static void startup_end(void);
190 static void layout_kernel_va(void);

```

```

192 /*
193  * Declare these as initialized data so we can patch them.
194  */
195 #ifdef __i386

197 /*
198  * Due to virtual address space limitations running in 32 bit mode, restrict
199  * the amount of physical memory configured to a max of PHYSMEM pages (16g).
200  */
201 * If the physical max memory size of 64g were allowed to be configured, the
202 * size of user virtual address space will be less than 1g. A limited user
203 * address space greatly reduces the range of applications that can run.
204 *
205 * If more physical memory than PHYSMEM is required, users should preferably
206 * run in 64 bit mode which has far looser virtual address space limitations.
207 *
208 * If 64 bit mode is not available (as in IA32) and/or more physical memory
209 * than PHYSMEM is required in 32 bit mode, physmem can be set to the desired
210 * value or to 0 (to configure all available memory) via eeprom(1M). kernelbase
211 * should also be carefully tuned to balance out the need of the user
212 * application while minimizing the risk of kernel heap exhaustion due to
213 * kernelbase being set too high.
214  */
215 #define PHYSMEM 0x400000

217 #else /* __amd64 */

219 /*
220  * For now we can handle memory with physical addresses up to about
221  * 64 Terabytes. This keeps the kernel above the VA hole, leaving roughly
222  * half the VA space for seg_kpm. When systems get bigger than 64TB this
223  * code will need revisiting. There is an implicit assumption that there
224  * are no *huge* holes in the physical address space too.
225  */
226 #define TERABYTE (1ul << 40)
227 #define PHYSMEM_MAX64 mmu_btop(64 * TERABYTE)
228 #define PHYSMEM PHYSMEM_MAX64
229 #define AMD64_VA_HOLE_END 0xFFFF80000000000ul

231 #endif /* __amd64 */

233 pgcnt_t physmem = PHYSMEM;
234 pgcnt_t obp_pages; /* Memory used by PROM for its text and data */

236 char *kobj_file_buf;
237 int kobj_file_bufsize; /* set in /etc/system */

239 /* Global variables for MP support. Used in mp_startup */
240 caddr_t rm_platter_va = 0;
241 uint32_t rm_platter_pa;

243 int auto_lpg_disable = 1;

245 /*
246  * Some CPUs have holes in the middle of the 64-bit virtual address range.
247  */
248 uintptr_t hole_start, hole_end;

250 /*
251  * kpm mapping window
252  */
253 caddr_t kpm_vbase;
254 size_t kpm_size;
255 static int kpm_desired;
256 #ifdef __amd64

```

```

257 static uintptr_t segkpm_base = (uintptr_t)SEGKPM_BASE;
258 #endif

260 /*
261 * Configuration parameters set at boot time.
262 */

264 caddr_t econtig;          /* end of first block of contiguous kernel */

266 struct bootops          *bootops = 0; /* passed in from boot */
267 struct bootops          **bootopsp;
268 struct boot_syscalls    *sysp;       /* passed in from boot */

270 char bootblock_fstype[16];

272 char kern_bootargs[OBP_MAXPATHLEN];
273 char kern_bootfile[OBP_MAXPATHLEN];

275 /*
276 * ZFS zio segment. This allows us to exclude large portions of ZFS data that
277 * gets cached in kmem caches on the heap. If this is set to zero, we allocate
278 * zio buffers from their own segment, otherwise they are allocated from the
279 * heap. The optimization of allocating zio buffers from their own segment is
280 * only valid on 64-bit kernels.
281 */
282 #if defined(__amd64)
283 int segzio_fromheap = 0;
284 #else
285 int segzio_fromheap = 1;
286 #endif

288 /*
289 * Give folks an escape hatch for disabling SMAP via kmdb. Doesn't work
290 * post-boot.
291 */
292 int disable_smap = 0;

294 /*
295 * new memory fragmentations are possible in startup() due to BOP_ALLOCS. this
296 * depends on number of BOP_ALLOC calls made and requested size, memory size
297 * combination and whether boot.bin memory needs to be freed.
298 */
299 #define POSS_NEW_FRAGMENTS    12

301 /*
302 * VM data structures
303 */
304 long page_hashsz;        /* Size of page hash table (power of two) */
305 unsigned int page_hashsz_shift; /* log2(page_hashsz) */
306 struct page *pp_base;    /* Base of initial system page struct array */
307 struct page **page_hash; /* Page hash table */
308 pad_mutex_t *pse_mutex; /* Locks protecting pp->p_selock */
309 size_t pse_table_size;   /* Number of mutexes in pse_mutex[] */
310 int pse_shift;          /* log2(pse_table_size) */
311 struct seg ktextseg;     /* Segment used for kernel executable image */
312 struct seg kalloc;      /* Segment used for "valloc" mapping */
313 struct seg kpsseg;      /* Segment used for pageable kernel virt mem */
314 struct seg kmapseg;     /* Segment used for generic kernel mappings */
315 struct seg kdebugseg;   /* Segment used for the kernel debugger */

317 struct seg *segkmap = &kmapseg; /* Kernel generic mapping segment */
318 static struct seg *segmap = &kmapseg; /* easier to use name for in here */

320 struct seg *segkp = &kpsseg; /* Pageable kernel virtual memory segment */

322 #if defined(__amd64)

```

```

323 struct seg kvseg_core; /* Segment used for the core heap */
324 struct seg kpmseg;     /* Segment used for physical mapping */
325 struct seg *segkpm = &kpmseg; /* 64bit kernel physical mapping segment */
326 #else
327 struct seg *segkpm = NULL; /* Unused on IA32 */
328 #endif

330 caddr_t segkp_base;    /* Base address of segkp */
331 caddr_t segzio_base;  /* Base address of segzio */
332 #if defined(__amd64)
333 pgcnt_t segkpsize = btop(SEGKPDEFESIZE); /* size of segkp segment in pages */
334 #else
335 pgcnt_t segkpsize = 0;
336 #endif
337 pgcnt_t segziosize = 0; /* size of zio segment in pages */

339 /*
340 * A static DR page_t VA map is reserved that can map the page structures
341 * for a domain's entire RA space. The pages that back this space are
342 * dynamically allocated and need not be physically contiguous. The DR
343 * map size is derived from KPM size.
344 * This mechanism isn't used by x86 yet, so just stubs here.
345 */
346 int ppvm_enable = 0; /* Static virtual map for page structs */
347 page_t *ppvm_base = NULL; /* Base of page struct map */
348 pgcnt_t ppvm_size = 0; /* Size of page struct map */

350 /*
351 * VA range available to the debugger
352 */
353 const caddr_t kdi_segdebugbase = (const caddr_t)SEGDEBUGBASE;
354 const size_t kdi_segdebugsize = SEGDEBUGSIZE;

356 struct memseg *memseg_base;
357 struct vnode unused_pages_vp;

359 #define FOURGB    0x100000000LL

361 struct memlist *memlist;

363 caddr_t s_text; /* start of kernel text segment */
364 caddr_t e_text; /* end of kernel text segment */
365 caddr_t s_data; /* start of kernel data segment */
366 caddr_t e_data; /* end of kernel data segment */
367 caddr_t modtext; /* start of loadable module text reserved */
368 caddr_t e_modtext; /* end of loadable module text reserved */
369 caddr_t moddata; /* start of loadable module data reserved */
370 caddr_t e_moddata; /* end of loadable module data reserved */

372 struct memlist *phys_install; /* Total installed physical memory */
373 struct memlist *phys_avail; /* Total available physical memory */
374 struct memlist *bios_rsvd; /* Bios reserved memory */

376 /*
377 * kphysm_init returns the number of pages that were processed
378 */
379 static pgcnt_t kphysm_init(page_t *, pgcnt_t);

381 #define IO_PROP_SIZE    64 /* device property size */

383 /*
384 * a couple useful roundup macros
385 */
386 #define ROUND_UP_PAGE(x) \
387 ((uintptr_t)P2ROUNDUP((uintptr_t)(x), (uintptr_t)MMU_PAGESIZE))
388 #define ROUND_UP_LPAGE(x) \

```

```

389      ((uintptr_t)P2ROUNDUP((uintptr_t)(x), mmu.level_size[1]))
390 #define ROUND_UP_4MEG(x) \
391      ((uintptr_t)P2ROUNDUP((uintptr_t)(x), (uintptr_t)FOUR_MEG))
392 #define ROUND_UP_TOPLLEVEL(x) \
393      ((uintptr_t)P2ROUNDUP((uintptr_t)(x), mmu.level_size[mmu.max_level]))

395 /*
396 *      32-bit Kernel's Virtual memory layout.
397 *      +-----+
398 *      |-----|
399 *      | 0xFFC00000 |-----| ARGSBASE
400 *      | debugger   |-----|
401 *      | 0xFF800000 |-----| SEGDEBUGBASE
402 *      | Kernel Data |-----|
403 *      | 0xFE000000 |-----|
404 *      | Kernel Text |-----|
405 *      | 0xFE800000 |-----| - KERNEL_TEXT (0xFB400000 on Xen)
406 *      | GDT         |-----| - GDT page (GDT_VA)
407 *      | debug info  |-----| - debug info (DEBUG_INFO_VA)
408 *      |
409 *      | page_t structures
410 *      | memsegs, memlists,
411 *      | page hash, etc.
412 *      |-----| - ekernelheap, valloc_base (floating)
413 *      | (segkp is just an arena in the heap)
414 *      |
415 *      | kvseg
416 *      |
417 *      |-----| - kernelheap (floating)
418 *      | Segkmap
419 *      |-----| - segmap_start (floating)
420 *      | 0xC3002000 |-----|
421 *      | Red Zone
422 *      | 0xC3000000 |-----| - kernelbase / userlimit (floating)
423 *      |
424 *      | Shared objects
425 *      |
426 *      |-----|
427 *      | user data
428 *      |-----|
429 *      | user text
430 *      | 0x08048000 |-----|
431 *      | user stack
432 *      |-----|
433 *      | invalid
434 *      | 0x00000000 |-----|
435 *      |
436 *      |
437 *      |
438 *      |
439 *      |
440 *      | 0xFFFFFFFF.FFC00000 |-----| - ARGSBASE
441 *      | debugger (?)
442 *      | 0xFFFFFFFF.FF800000 |-----| - SEGDEBUGBASE
443 *      | unused
444 *      |-----|
445 *      | Kernel Data
446 *      |-----|
447 *      | Kernel Text
448 *      |-----| - KERNEL_TEXT
449 *      |-----| - debug info (DEBUG_INFO_VA)
450 *      |-----| - GDT page (GDT_VA)
451 *      |-----| - IDT page (IDT_VA)
452 *      |-----| - LDT pages (LDT_VA)
453 *      |-----| - debug info (DEBUG_INFO_VA)

```

64-bit Kernel's Virtual memory layout. (assuming 64 bit app)

```

454 *      Core heap
455 *      0xFFFFFFFF.C0000000 |-----| (used for loadable modules)
456 *      Kernel heap
457 *      |-----| - core_base / ekernelheap
458 *      0xFFFFFXXX.XXX00000 |-----| - kernelheap (floating)
459 *      segmap
460 *      0xFFFFFXXX.XXX00000 |-----| - segmap_start (floating)
461 *      device mappings
462 *      0xFFFFFXXX.XXX00000 |-----| - toxic_addr (floating)
463 *      segzio
464 *      0xFFFFFXXX.XXX00000 |-----| - segzio_base (floating)
465 *      segkp
466 *      --- |-----| - segkp_base (floating)
467 *      page_t structures
468 *      memsegs, memlists,
469 *      page hash, etc.
470 *      0xFFFFF00.00000000 |-----| - valloc_base (lower if >256GB)
471 *      segkpm
472 *      0xFFFFFE00.00000000 |-----|
473 *      Red Zone
474 *      0xFFFFFD80.00000000 |-----| - KERNELBASE (lower if >256GB)
475 *      User stack
476 *      |-----| - User space memory
477 *      shared objects, etc
478 *      |-----| (grows downwards)
479 *      |
480 *      0xFFFFF8000.00000000 |-----|
481 *      VA Hole / unused
482 *      |-----|
483 *      0x00008000.00000000 |-----|
484 *      |
485 *      |
486 *      |
487 *      |-----|
488 *      user heap
489 *      |-----| (grows upwards)
490 *      user data
491 *      |-----|
492 *      user text
493 *      0x00000000.04000000 |-----|
494 *      invalid
495 *      0x00000000.00000000 |-----|
496 *      |
497 *      A 32 bit app on the 64 bit kernel sees the same layout as on the 32 bit
498 *      kernel, except that userlimit is raised to 0xfe000000
499 *      |
500 *      Floating values:
501 *      |
502 *      valloc_base: start of the kernel's memory management/tracking data
503 *      structures. This region contains page_t structures for
504 *      physical memory, memsegs, memlists, and the page hash.
505 *      |
506 *      core_base: start of the kernel's "core" heap area on 64-bit systems.
507 *      This area is intended to be used for global data as well as for module
508 *      text/data that does not fit into the nucleus pages. The core heap is
509 *      restricted to a 2GB range, allowing every address within it to be
510 *      accessed using rip-relative addressing
511 *      |
512 *      ekernelheap: end of kernelheap and start of segmap.
513 *      |
514 *      kernelheap: start of kernel heap. On 32-bit systems, this starts right
515 *      above a red zone that separates the user's address space from the
516 *      kernel's. On 64-bit systems, it sits above segkp and segkpm.
517 *      |
518 *      segmap_start: start of segmap. The length of segmap can be modified
519 *      through eeprom. The default length is 16MB on 32-bit systems and 64MB

```

```

520 * on 64-bit systems.
521 *
522 * kernelbase: On a 32-bit kernel the default value of 0xd4000000 will be
523 * decreased by 2X the size required for page_t. This allows the kernel
524 * heap to grow in size with physical memory. With sizeof(page_t) == 80
525 * bytes, the following shows the values of kernelbase and kernel heap
526 * sizes for different memory configurations (assuming default segmap and
527 * segkp sizes).
528 *
529 *      mem      size for      kernelbase      kernel heap
530 *      size      page_t's      size              size
531 *      ----      -
532 *      1gb      0x01400000      0xd1800000      684MB
533 *      2gb      0x02800000      0xcf000000      704MB
534 *      4gb      0x05000000      0xca000000      744MB
535 *      6gb      0x07800000      0xc5000000      784MB
536 *      8gb      0x0a000000      0xc0000000      824MB
537 *      16gb     0x14000000      0xac000000      984MB
538 *      32gb     0x28000000      0x84000000      1304MB
539 *      64gb     0x50000000      0x34000000      1944MB (*)
540 *
541 * kernelbase is less than the abi minimum of 0xc0000000 for memory
542 * configurations above 8gb.
543 *
544 * (*) support for memory configurations above 32gb will require manual tuning
545 * of kernelbase to balance out the need of user applications.
546 */

548 /* real-time-clock initialization parameters */
549 extern time_t process_rtc_config_file(void);

551 uintptr_t      kernelbase;
552 uintptr_t      postbootkernelbase; /* not set till boot loader is gone */
553 uintptr_t      eprom_kernelbase;
554 size_t         segmapsize;
555 uintptr_t      segmap_start;
556 int           segmapfreelists;
557 pgcnt_t       npages;
558 pgcnt_t       orig_npages;
559 size_t         core_size; /* size of "core" heap */
560 uintptr_t      core_base; /* base address of "core" heap */

562 /*
563 * List of bootstrap pages. We mark these as allocated in startup.
564 * release_bootstrap() will free them when we're completely done with
565 * the bootstrap.
566 */
567 static page_t *bootpages;

569 /*
570 * boot time pages that have a vnode from the ramdisk will keep that forever.
571 */
572 static page_t *rd_pages;

574 /*
575 * Lower 64K
576 */
577 static page_t *lower_pages = NULL;
578 static int lower_pages_count = 0;

580 struct system_hardware system_hardware;

582 /*
583 * Enable some debugging messages concerning memory usage...
584 */
585 static void

```

```

586 print_memlist(char *title, struct memlist *mp)
587 {
588     prom_printf("MEMLIST: %s:\n", title);
589     while (mp != NULL) {
590         prom_printf("\tAddress 0x%" PRIx64 " , size 0x%" PRIx64 "\n",
591             mp->ml_address, mp->ml_size);
592         mp = mp->ml_next;
593     }
594 }

_____ unchanged_portion_omitted _____

922 static void
923 kpm_init()
924 {
925     struct segkpm_crargs b;

927     /*
928     * These variables were all designed for sfmmu in which segkpm is
929     * mapped using a single pagesize - either 8KB or 4MB. On x86, we
930     * might use 2+ page sizes on a single machine, so none of these
931     * variables have a single correct value. They are set up as if we
932     * always use a 4KB pagesize, which should do no harm. In the long
933     * run, we should get rid of KPM's assumption that only a single
934     * pagesize is used.
935     */
936     kpm_pgshft = MMU_PAGESHIFT;
937     kpm_pgsz = MMU_PAGESIZE;
938     kpm_pgoff = MMU_PAGEOFFSET;
939     kpm_p2pshft = 0;
940     kpm_npgs = 1;
941     ASSERT(((uintptr_t)kpm_vbase & (kpm_pgsz - 1)) == 0);

943     PRM_POINT("about to create segkpm");
944     rw_enter(&kas.a_lock, RW_WRITER);

946     if (seg_attach(&kas, kpm_vbase, kpm_size, segkpm) < 0)
947         panic("cannot attach segkpm");

949     b.prot = PROT_READ | PROT_WRITE;
950     b.nvcolors = 1;

952     if (segkpm_create(segkpm, (caddr_t)&b) != 0)
953         panic("segkpm_create segkpm");

955     rw_exit(&kas.a_lock);

957     kpm_enable = 1;

959     /*
960     * As the KPM was disabled while setting up the system, go back and fix
961     * CPU zero's access to its user page table. This is a bit gross, but
962     * we have a chicken and egg problem otherwise.
963     */
964     ASSERT(CPU->cpu_hat_info->hci_user_l3ptes == NULL);
965     CPU->cpu_hat_info->hci_user_l3ptes =
966         (x86pte_t *)hat_kpm_mapin_pfn(CPU->cpu_hat_info->hci_user_l3pfn);
967 }

_____ unchanged_portion_omitted _____

1423 /*
1424 * Layout the kernel's part of address space and initialize kmem allocator.
1425 */
1426 static void
1427 startup_kmem(void)
1428 {
1429     extern void page_set_colorequiv_arr(void);

```

```

1430 #if !defined(__xpv)
1431     extern uint64_t kpti_kbase;
1432 #endif

1434     PRM_POINT("startup_kmem() starting...");

1436 #if defined(__amd64)
1437     if (eprom_kernelbase && eprom_kernelbase != KERNELBASE)
1438         cmn_err(CE_NOTE, "!kernelbase cannot be changed on 64-bit "
1439             "systems.");
1440     kernelbase = segkpm_base - KERNEL_REDZONE_SIZE;
1441     core_base = (uintptr_t)COREHEAP_BASE;
1442     core_size = (size_t)MISC_VA_BASE - COREHEAP_BASE;
1443 #else
1444     /* __i386 */
1445     /* We configure kernelbase based on:
1446     *
1447     * 1. user specified kernelbase via eeprom command. Value cannot exceed
1448     *    KERNELBASE_MAX. we large page align eprom_kernelbase
1449     *
1450     * 2. Default to KERNELBASE and adjust to 2X less the size for page_t.
1451     *    On large memory systems we must lower kernelbase to allow
1452     *    enough room for page_t's for all of memory.
1453     *
1454     * The value set here, might be changed a little later.
1455     */
1456     if (eprom_kernelbase) {
1457         kernelbase = eprom_kernelbase & mmu.level_mask[1];
1458         if (kernelbase > KERNELBASE_MAX)
1459             kernelbase = KERNELBASE_MAX;
1460     } else {
1461         kernelbase = (uintptr_t)KERNELBASE;
1462         kernelbase -= ROUND_UP_4MEG(2 * valloc_sz);
1463     }
1464     ASSERT((kernelbase & mmu.level_offset[1]) == 0);
1465     core_base = valloc_base;
1466     core_size = 0;
1467 #endif
1468 /* __i386 */

1469     PRM_DEBUG(core_base);
1470     PRM_DEBUG(core_size);
1471     PRM_DEBUG(kernelbase);

1473 #if defined(__i386)
1474     segkp_fromheap = 1;
1475 #endif
1476 /* __i386 */

1477     ekernelheap = (char *)core_base;
1478     PRM_DEBUG(ekernelheap);

1480     /*
1481     * Now that we know the real value of kernelbase,
1482     * update variables that were initialized with a value of
1483     * KERNELBASE (in common/conf/param.c).
1484     *
1485     * XXX The problem with this sort of hackery is that the
1486     * compiler just may feel like putting the const declarations
1487     * (in param.c) into the .text section. Perhaps they should
1488     * just be declared as variables there?
1489     */

1491     *(uintptr_t *)&_kernelbase = kernelbase;
1492     *(uintptr_t *)&_userlimit = kernelbase;
1493 #if defined(__amd64)
1494     *(uintptr_t *)&_userlimit -= KERNELBASE - USERLIMIT;
1495 #if !defined(__xpv)

```

```

1496     kpti_kbase = kernelbase;
1497 #endif
1498 #else
1499     *(uintptr_t *)&_userlimit32 = _userlimit;
1500 #endif
1501     PRM_DEBUG(_kernelbase);
1502     PRM_DEBUG(_userlimit);
1503     PRM_DEBUG(_userlimit32);

1505     /* We have to re-do this now that we've modified _userlimit. */
1506     mmu_calc_user_slots();

1508     layout_kernel_va();

1510 #if defined(__i386)
1511     /*
1512     * If segmap is too large we can push the bottom of the kernel heap
1513     * higher than the base. Or worse, it could exceed the top of the
1514     * VA space entirely, causing it to wrap around.
1515     */
1516     if (kernelheap >= ekernelheap || (uintptr_t)kernelheap < kernelbase)
1517         panic("too little address space available for kernelheap,"
1518             " use eeprom for lower kernelbase or smaller segmapsize");
1519 #endif
1520 /* __i386 */

1521     /*
1522     * Initialize the kernel heap. Note 3rd argument must be > 1st.
1523     */
1524     kernelheap_init(kernelheap, ekernelheap,
1525         kernelheap + MMU_PAGESIZE,
1526         (void *)core_base, (void *) (core_base + core_size));

1528 #if defined(__xpv)
1529     /*
1530     * Link pending events struct into cpu struct
1531     */
1532     CPU->cpu_m.mcpu_evt_pend = &cpu0_evt_data;
1533 #endif
1534     /*
1535     * Initialize kernel memory allocator.
1536     */
1537     kmem_init();

1539     /*
1540     * Factor in colerequiv to check additional 'equivalent' bins
1541     */
1542     page_set_colorequiv_arr();

1544     /*
1545     * print this out early so that we know what's going on
1546     */
1547     print_x86_featureset(x86_featureset);

1549     /*
1550     * Initialize bp_mapin().
1551     */
1552     bp_init(MMU_PAGESIZE, HAT_STORECACHING_OK);

1554     /*
1555     * orig_npages is non-zero if physmem has been configured for less
1556     * than the available memory.
1557     */
1558     if (orig_npages) {
1559         cmn_err(CE_WARN, "!%slimiting physmem to 0x%lx of 0x%lx pages",
1560             (npages == PHYSMEM ? "Due to virtual address space " : ""),
1561             npages, orig_npages);

```

```

1562     }
1563 #if defined(__i386)
1564     if (eprom_kernelbase && (eprom_kernelbase != kernelbase))
1565         cmn_err(CE_WARN, "kernelbase value, User specified 0x%lx, "
1566              "System using 0x%lx",
1567              (uintptr_t)eprom_kernelbase, (uintptr_t)kernelbase);
1568 #endif

1570 #ifdef KERNELBASE_ABI_MIN
1571 if (kernelbase < (uintptr_t)KERNELBASE_ABI_MIN) {
1572     cmn_err(CE_NOTE, "!kernelbase set to 0x%lx, system is not "
1573          "i386 ABI compliant.", (uintptr_t)kernelbase);
1574 }
1575 #endif

1577 #ifndef __xpv
1578 if (plat_dr_support_memory()) {
1579     mem_config_init();
1580 }
1581 #else /* __xpv */
1582 /*
1583  * Some of the xen start information has to be relocated up
1584  * into the kernel's permanent address space.
1585  */
1586 PRM_POINT("calling xen_relocate_start_info()");
1587 xen_relocate_start_info();
1588 PRM_POINT("xen_relocate_start_info() done");

1590 /*
1591  * (Update the vcpu pointer in our cpu structure to point into
1592  * the relocated shared info.)
1593  */
1594 CPU->cpu_m.mcpu_vcpu_info =
1595     &HYPERVISOR_shared_info->vcpu_info[CPU->cpu_id];
1596 #endif /* __xpv */

1598     PRM_POINT("startup_kmem() done");
1599 }
_____ unchanged_portion_omitted_____

2013 /*
2014  * Finish initializing the VM system, now that we are no longer
2015  * relying on the boot time memory allocators.
2016  */
2017 static void
2018 startup_vm(void)
2019 {
2020     struct segmap_crargs a;

2022     extern int use_brk_lpg, use_stk_lpg;

2024     PRM_POINT("startup_vm() starting...");

2026     /*
2027      * Initialize the hat layer.
2028      */
2029     hat_init();

2031     /*
2032      * Do final allocations of HAT data structures that need to
2033      * be allocated before quiescing the boot loader.
2034      */
2035     PRM_POINT("Calling hat_kern_alloc()...");
2036     hat_kern_alloc((caddr_t)segmap_start, segmapsize, ekernelheap);
2037     PRM_POINT("hat_kern_alloc() done");

```

```

2039 #ifndef __xpv
2040     /*
2041      * Setup Page Attribute Table
2042      */
2043     pat_sync();
2044 #endif

2046     /*
2047      * The next two loops are done in distinct steps in order
2048      * to be sure that any page that is doubly mapped (both above
2049      * KERNEL_TEXT and below kernelbase) is dealt with correctly.
2050      * Note this may never happen, but it might someday.
2051      */
2052     bootpages = NULL;
2053     PRM_POINT("Protecting boot pages");

2055     /*
2056      * Protect any pages mapped above KERNEL_TEXT that somehow have
2057      * page_t's. This can only happen if something weird allocated
2058      * in this range (like kadb/kmdb).
2059      */
2060     protect_boot_range(KERNEL_TEXT, (uintptr_t)-1, 0);

2062     /*
2063      * Before we can take over memory allocation/mapping from the boot
2064      * loader we must remove from our free page lists any boot allocated
2065      * pages that stay mapped until release_bootstrap().
2066      */
2067     protect_boot_range(0, kernelbase, 1);

2070     /*
2071      * Switch to running on regular HAT (not boot_mmu)
2072      */
2073     PRM_POINT("Calling hat_kern_setup()...");
2074     hat_kern_setup();

2076     /*
2077      * It is no longer safe to call BOP_ALLOC(), so make sure we don't.
2078      */
2079     bop_no_more_mem();

2081     PRM_POINT("hat_kern_setup() done");

2083     hat_cpu_online(CPU);

2085     /*
2086      * Initialize VM system
2087      */
2088     PRM_POINT("Calling kvm_init()...");
2089     kvm_init();
2090     PRM_POINT("kvm_init() done");

2092     /*
2093      * Tell kmdb that the VM system is now working
2094      */
2095     if (boothowto & RB_DEBUG)
2096         kdi_dvec_vmready();

2098 #if defined(__xpv)
2099     /*
2100      * Populate the I/O pool on domain 0
2101      */
2102     if (DOMAIN_IS_INITDOMAIN(xen_info)) {
2103         extern long populate_io_pool(void);
2104         long init_io_pool_cnt;

```

```

2106         PRM_POINT("Populating reserve I/O page pool");
2107         init_io_pool_cnt = populate_io_pool();
2108         PRM_DEBUG(init_io_pool_cnt);
2109     }
2110 #endif
2111     /*
2112     * Mangle the brand string etc.
2113     */
2114     cpuid_pass3(CPU);

2116 #if defined(__amd64)

2118     /*
2119     * Create the device arena for toxic (to dtrace/kmdb) mappings.
2120     */
2121     device_arena = vmem_create("device", (void *)toxic_addr,
2122         toxic_size, MMU_PAGESIZE, NULL, NULL, NULL, 0, VM_SLEEP);

2124 #else /* __i386 */

2126     /*
2127     * allocate the bit map that tracks toxic pages
2128     */
2129     toxic_bit_map_len = btop((ulong_t)(valloc_base - kernelbase));
2130     PRM_DEBUG(toxic_bit_map_len);
2131     toxic_bit_map =
2132         kmem_zalloc(BT_SIZEOFMAP(toxic_bit_map_len), KM_NOSLEEP);
2133     ASSERT(toxic_bit_map != NULL);
2134     PRM_DEBUG(toxic_bit_map);

2136 #endif /* __i386 */

2139     /*
2140     * Now that we've got more VA, as well as the ability to allocate from
2141     * it, tell the debugger.
2142     */
2143     if (boothowto & RB_DEBUG)
2144         kdi_dvec_memavail();

2124     /*
2125     * The following code installs a special page fault handler (#pf)
2126     * to work around a pentium bug.
2127     */
2128 #if !defined(__amd64) && !defined(__xpv)
2129     if (x86_type == X86_TYPE_P5) {
2130         desc_t idtr;
2131         gate_desc_t *newidt;

2133         if ((newidt = kmem_zalloc(MMU_PAGESIZE, KM_NOSLEEP)) == NULL)
2134             panic("failed to install pentium_pftap");

2136         bcopy(idt0, newidt, NIDT * sizeof (*idt0));
2137         set_gatesegd(&newidt[T_PGFLT], &pentium_pftap,
2138             KCS_SEL, SDT_SYSIGT, TRP_KPL, 0);

2140         (void) as_setprot(&kas, (caddr_t)newidt, MMU_PAGESIZE,
2141             PROT_READ | PROT_EXEC);

2143         CPU->cpu_idt = newidt;
2144         idtr.dtr_base = (uintptr_t)CPU->cpu_idt;
2145         idtr.dtr_limit = (NIDT * sizeof (*idt0)) - 1;
2146         wr_idtr(&idtr);
2147     }
2148 #endif /* !__amd64 */

```

```

2146 #if !defined(__xpv)
2147     /*
2148     * Map page pfn=0 for drivers, such as kd, that need to pick up
2149     * parameters left there by controllers/BIOS.
2150     */
2151     PRM_POINT("setup up p0_va");
2152     p0_va = i86devmap(0, 1, PROT_READ);
2153     PRM_DEBUG(p0_va);
2154 #endif

2156     cmn_err(CE_CONT, "?mem = %luK (0x%lx)\n",
2157         physinstalled << (MMU_PAGESHIFT - 10), ptob(physinstalled));

2159     /*
2160     * disable automatic large pages for small memory systems or
2161     * when the disable flag is set.
2162     * Do not yet consider page sizes larger than 2m/4m.
2163     */
2164     if (!auto_lpg_disable && mmu.max_page_level > 0) {
2165         max_uheap_lpsize = LEVEL_SIZE(1);
2166         max_ustack_lpsize = LEVEL_SIZE(1);
2167         max_privmap_lpsize = LEVEL_SIZE(1);
2168         max_uidata_lpsize = LEVEL_SIZE(1);
2169         max_utext_lpsize = LEVEL_SIZE(1);
2170         max_shm_lpsize = LEVEL_SIZE(1);
2171     }
2172     if (physmem < privm_lpg_min_physmem || mmu.max_page_level == 0 ||
2173         auto_lpg_disable) {
2174         use_brk_lpg = 0;
2175         use_stk_lpg = 0;
2176     }
2177     mcntl0_lpsize = LEVEL_SIZE(mmu.umax_page_level);

2180     PRM_POINT("Calling hat_init_finish()...");
2181     hat_init_finish();
2182     PRM_POINT("hat_init_finish() done");

2184     /*
2185     * Initialize the segkp segment type.
2186     */
2187     rw_enter(&kas.a_lock, RW_WRITER);
2188     PRM_POINT("Attaching segkp");
2189     if (segkp_fromheap) {
2190         segkp->s_as = &kas;
2191     } else if (seg_attach(&kas, (caddr_t)segkp_base, mmu_ptob(segkpsize),
2192         segkp) < 0) {
2193         panic("startup: cannot attach segkp");
2194         /*NOTREACHED*/
2195     }
2196     PRM_POINT("Doing segkp_create()");
2197     if (segkp_create(segkp) != 0) {
2198         panic("startup: segkp_create failed");
2199         /*NOTREACHED*/
2200     }
2201     PRM_DEBUG(segkp);
2202     rw_exit(&kas.a_lock);

2204     /*
2205     * kpm segment
2206     */
2207     segmap_kpm = 0;
2208     if (kpm_desired)
2209         if (kpm_desired) {
2210             kpm_init();

```



```

2214         kpm_enable = 1;
2215     }

2211     /*
2212     * Now create segmap segment.
2213     */
2214     rw_enter(&kas.a_lock, RW_WRITER);
2215     if (seg_attach(&kas, (caddr_t)segmap_start, segmapsize, segmap) < 0) {
2216         panic("cannot attach segmap");
2217         /*NOTREACHED*/
2218     }
2219     PRM_DEBUG(segmap);

2221     a.prot = PROT_READ | PROT_WRITE;
2222     a.shmsize = 0;
2223     a.nfreelist = segmapfreelists;

2225     if (segmap_create(segmap, (caddr_t)&a) != 0)
2226         panic("segmap_create segmap");
2227     rw_exit(&kas.a_lock);

2229     setup_vaddr_for_ppcopy(CPU);

2231     segdev_init();
2232 #if defined(__xpv)
2233     if (DOMAIN_IS_INITDOMAIN(xen_info))
2234 #endif
2235         pmem_init();

2237     PRM_POINT("startup_vm() done");
2238 }
    unchanged portion omitted

2250 static void
2251 startup_end(void)
2252 {
2253     int i;
2254     extern void setx86isalist(void);
2255     extern void cpu_event_init(void);

2257     PRM_POINT("startup_end() starting...");

2259     /*
2260     * Perform tasks that get done after most of the VM
2261     * initialization has been done but before the clock
2262     * and other devices get started.
2263     */
2264     kern_setup1();

2266     /*
2267     * Perform CPC initialization for this CPU.
2268     */
2269     kcpc_hw_init(CPU);

2271     /*
2272     * Initialize cpu event framework.
2273     */
2274     cpu_event_init();

2276 #if defined(OPTERON_WORKAROUND_6323525)
2277     if (opteron_workaround_6323525)
2278         patch_workaround_6323525();
2279 #endif

2280     /*
2281     * If needed, load TOD module now so that ddi_get_time(9F) etc. work
2282     * (For now, "needed" is defined as set tod_module_name in /etc/system)

```

```

2283     /*
2284     if (tod_module_name != NULL) {
2285         PRM_POINT("load_tod_module()");
2286         load_tod_module(tod_module_name);
2287     }

2289 #if defined(__xpv)
2290     /*
2291     * Forceload interposing TOD module for the hypervisor.
2292     */
2293     PRM_POINT("load_tod_module()");
2294     load_tod_module("xpvtod");
2295 #endif

2297     /*
2298     * Configure the system.
2299     */
2300     PRM_POINT("Calling configure()...");
2301     configure(); /* set up devices */
2302     PRM_POINT("configure() done");

2304     /*
2305     * We can now setup for XSAVE because fpu_probe is done in configure().
2306     */
2307     if (fp_save_mech == FP_XSAVE) {
2308         xsave_setup_msr(CPU);
2309     }

2311     /*
2312     * Set the isa_list string to the defined instruction sets we
2313     * support.
2314     */
2315     setx86isalist();
2316     cpu_intr_alloc(CPU, NINTR_THREADS);
2317     psm_install();

2319     /*
2320     * We're done with bootops. We don't unmap the bootstrap yet because
2321     * we're still using bootsvcs.
2322     */
2323     PRM_POINT("NULLing out bootops");
2324     *bootopsp = (struct bootops *)NULL;
2325     bootops = (struct bootops *)NULL;

2327 #if defined(__xpv)
2328     ec_init_debug_irq();
2329     xs_domu_init();
2330 #endif

2332 #if !defined(__xpv)
2333 #if defined(__amd64) && !defined(__xpv)
2334     /*
2335     * Intel IOMMU has been setup/initialized in ddi_impl.c
2336     * Start it up now.
2337     */
2337     immu_startup();

2339     /*
2340     * Now that we're no longer going to drop into real mode for a BIOS call
2341     * via bootops, we can enable PCID (which requires CR0.PG).
2342     */
2343     enable_pcid();
2344 #endif

2346     PRM_POINT("Enabling interrupts");
2347     (*picinitf)();

```

```
2348     sti();
2349 #if defined(__xpv)
2350     ASSERT(CPU->cpu_m.mcpu_vcpu_info->evtchn_upcall_mask == 0);
2351     xen_late_startup();
2352 #endif

2354     (void) add_avsoftintr((void *)&softlevell_hdl, 1, softlevell,
2355                          "softlevell", NULL, NULL); /* XXX to be moved later */

2357     /*
2358     * Register software interrupt handlers for ddi_periodic_add(9F).
2359     * Software interrupts up to the level 10 are supported.
2360     */
2361     for (i = DDI_IPL_1; i <= DDI_IPL_10; i++) {
2362         (void) add_avsoftintr((void *)&softlevel_hdl[i-1], i,
2363                              (avfunc)ddi_periodic_softintr, "ddi_periodic",
2364                              (caddr_t)(uintptr_t)i, NULL);
2365     }

2367 #if !defined(__xpv)
2368     if (modload("drv", "amd_iommu") < 0) {
2369         PRM_POINT("No AMD IOMMU present\n");
2370     } else if (ddi_hold_installed_driver(ddi_name_to_major(
2371         "amd_iommu")) == NULL) {
2372         prom_printf("ERROR: failed to attach AMD IOMMU\n");
2373     }
2374 #endif
2375     post_startup_cpu_fixups();

2377     PRM_POINT("startup_end() done");
2378 }
unchanged_portion_omitted
```

new/usr/src/uts/i86pc/os/trap.c

1

```
*****
64156 Fri Apr 6 17:25:03 2018
new/usr/src/uts/i86pc/os/trap.c
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21
22 /*
23  * Copyright (c) 1992, 2010, Oracle and/or its affiliates. All rights reserved.
24  */
25
26 /*      Copyright (c) 1990, 1991 UNIX System Laboratories, Inc. */
27 /*      Copyright (c) 1984, 1986, 1987, 1988, 1989, 1990 AT&T */
28 /*              All Rights Reserved */
29 /*              */
30 /*      Copyright (c) 1987, 1988 Microsoft Corporation */
31 /*              All Rights Reserved */
32 /*              */
33
34 /*
35  * Copyright 2018 Joyent, Inc.
36  * Copyright 2017 Joyent, Inc.
37  */
38 #include <sys/types.h>
39 #include <sys/sysmacros.h>
40 #include <sys/param.h>
41 #include <sys/signal.h>
42 #include <sys/unistd.h>
43 #include <sys/user.h>
44 #include <sys/proc.h>
45 #include <sys/disp.h>
46 #include <sys/class.h>
47 #include <sys/core.h>
48 #include <sys/syscall.h>
49 #include <sys/cpuvar.h>
50 #include <sys/vm.h>
51 #include <sys/sysinfo.h>
52 #include <sys/fault.h>
53 #include <sys/stack.h>
54 #include <sys/psw.h>
55 #include <sys/regset.h>
56 #include <sys/fp.h>
57 #include <sys/trap.h>
58 #include <sys/kmem.h>
```

new/usr/src/uts/i86pc/os/trap.c

2

```
59 #include <sys/vtrace.h>
60 #include <sys/cmn_err.h>
61 #include <sys/prsystem.h>
62 #include <sys/mutex_impl.h>
63 #include <sys/machsystem.h>
64 #include <sys/archsystem.h>
65 #include <sys/sdt.h>
66 #include <sys/avintr.h>
67 #include <sys/kobj.h>
68
69 #include <vm/hat.h>
70
71 #include <vm/seg_kmem.h>
72 #include <vm/as.h>
73 #include <vm/seg.h>
74 #include <vm/hat_pte.h>
75 #include <vm/hat_i86.h>
76
77 #include <sys/procfs.h>
78
79 #include <sys/reboot.h>
80 #include <sys/debug.h>
81 #include <sys/debugreg.h>
82 #include <sys/modctl.h>
83 #include <sys/aio_impl.h>
84 #include <sys/tnf.h>
85 #include <sys/tnf_probe.h>
86 #include <sys/cred.h>
87 #include <sys/mman.h>
88 #include <sys/x86_archext.h>
89 #include <sys/copyops.h>
90 #include <c2/audit.h>
91 #include <sys/ftrace.h>
92 #include <sys/panic.h>
93 #include <sys/traptrace.h>
94 #include <sys/ontrap.h>
95 #include <sys/cpc_impl.h>
96 #include <sys/bootconf.h>
97 #include <sys/bootinfo.h>
98 #include <sys/promif.h>
99 #include <sys/mach_mmu.h>
100 #if defined(__xpv)
101 #include <sys/hypervisor.h>
102 #endif
103 #include <sys/contract/process_impl.h>
104
105 #define USER      0x10000      /* user-mode flag added to trap type */
106
107 static const char *trap_type_mnemonic[] = {
108     "de",    "db",    "2",    "bp",
109     "of",    "br",    "ud",    "nm",
110     "df",    "9",    "ts",    "np",
111     "ss",    "gp",    "pf",    "15",
112     "mf",    "ac",    "mc",    "xf"
113 };
114
115 unchanged portion omitted
116
117 #endif /* OPTERON_ERRATUM_91 */
118
119 /*
120  * Called from the trap handler when a processor trap occurs.
121  *
122  * Note: All user-level traps that might call stop() must exit
123  * trap() by 'goto out' or by falling through.
124  * Note Also: trap() is usually called with interrupts enabled, (PS_IE == 1)
125  * however, there are paths that arrive here with PS_IE == 0 so special care
126  */
```

```

462 * must be taken in those cases.
463 */
464 void
465 trap(struct regs *rp, caddr_t addr, processorid_t cpuid)
466 {
467     kthread_t *ct = curthread;
468     enum seg_rw rw;
469     unsigned type;
470     proc_t *p = ttproc(ct);
471     klwp_t *lwp = ttolwp(ct);
472     uintptr_t lofault;
473     label_t *onfault;
474     faultcode_t pagefault(), res, errcode;
475     enum fault_type fault_type;
476     k_siginfo_t siginfo;
477     uint_t fault = 0;
478     int mstate;
479     int sicode = 0;
480     int watchcode;
481     int watchpage;
482     caddr_t vaddr;
483     int singlestep_twiddle;
483     size_t sz;
484     int ta;
485 #ifdef __amd64
486     uchar_t instr;
487 #endif
488
489     ASSERT_STACK_ALIGNED();
490
491     type = rp->r_trapno;
492     CPU_STATS_ADDQ(CPU, sys, trap, 1);
493     ASSERT(ct->t_schedflag & TS_DONT_SWAP);
494
495     if (type == T_PGFLT) {
496
497         errcode = rp->r_err;
498         if (errcode & PF_ERR_WRITE)
499             rw = S_WRITE;
500         else if ((caddr_t)rp->r_pc == addr ||
501                (mmu.pt_nx != 0 && (errcode & PF_ERR_EXEC)))
502             rw = S_EXEC;
503         else
504             rw = S_READ;
505
506 #if defined(__i386)
507     /*
508      * Pentium Pro work-around
509      */
510     if ((errcode & PF_ERR_PROT) && pentiumpro_bug4046376) {
511         uint_t attr;
512         uint_t priv_violation;
513         uint_t access_violation;
514
515         if (hat_getattr(addr < (caddr_t)kernelbase ?
516                        curproc->p_as->a_hat : kas.a_hat, addr, &attr)
517             == -1) {
518             errcode &= ~PF_ERR_PROT;
519         } else {
520             priv_violation = (errcode & PF_ERR_USER) &&
521                             !(attr & PROT_USER);
522             access_violation = (errcode & PF_ERR_WRITE) &&
523                               !(attr & PROT_WRITE);
524             if (!priv_violation && !access_violation)
525                 goto cleanup;
526         }
527     }
528 #endif
529
530     }
531
532     }
533
534     }
535
536     }
537
538     }
539
540     }
541
542     }
543
544     }
545
546     }
547
548     }
549
550     }
551
552     }
553
554     }
555
556     }
557
558     }
559
560     }
561
562     }
563
564     }
565
566     }
567
568     }
569
570     }
571
572     }
573
574     }
575
576     }
577
578     }
579
580     }
581
582     }
583
584     }
585
586     }
587
588     }
589
590     }
591
592     }
593
594     }
595
596     }
597
598     }
599
600     }
601
602     }
603
604     }
605
606     }
607
608     }
609
610     }
611
612     }
613
614     }
615
616     }
617
618     }
619
620     }
621
622     }
623
624     }
625
626     }
627
628     }
629
630     }
631
632     }
633
634     }
635
636     }
637
638     }
639
640     }
641
642     }
643
644     }
645
646     }
647
648     }
649
650     }
651
652     }
653
654     }
655
656     }
657
658     }
659
660     }
661
662     }
663
664     }
665
666     }
667
668     }
669
670     }
671
672     }
673
674     }
675
676     }
677
678     }
679
680     }
681
682     }
683
684     }
685
686     }
687
688     }
689
690     }
691
692     }
693
694     }
695
696     }
697
698     }
699
700     }
701
702     }
703
704     }
705
706     }
707
708     }
709
710     }
711
712     }
713
714     }
715
716     }
717
718     }
719
720     }
721
722     }
723
724     }
725
726     }
727
728     }
729
730     }
731
732     }
733
734     }
735
736     }
737
738     }
739
740     }
741
742     }
743
744     }
745
746     }
747
748     }
749
750     }
751
752     }
753
754     }
755
756     }
757
758     }
759
760     }
761
762     }
763
764     }
765
766     }
767
768     }
769
770     }
771
772     }
773
774     }
775
776     }
777
778     }
779
780     }
781
782     }
783
784     }
785
786     }
787
788     }
789
790     }
791
792     }
793
794     }
795
796     }
797
798     }
799
800     }
801
802     }
803
804     }
805
806     }
807
808     }
809
810     }
811
812     }
813
814     }
815
816     }
817
818     }
819
820     }
821
822     }
823
824     }
825
826     }
827
828     }
829
830     }
831
832     }
833
834     }
835
836     }
837
838     }
839
840     }
841
842     }
843
844     }
845
846     }
847
848     }
849
850     }
851
852     }
853
854     }
855
856     }
857
858     }
859
860     }
861
862     }
863
864     }
865
866     }
867
868     }
869
870     }
871
872     }
873
874     }
875
876     }
877
878     }
879
880     }
881
882     }
883
884     }
885
886     }
887
888     }
889
890     }
891
892     }
893
894     }
895
896     }
897
898     }
899
900     }
901
902     }
903
904     }
905
906     }
907
908     }
909
910     }
911
912     }
913
914     }
915
916     }
917
918     }
919
920     }
921
922     }
923
924     }
925
926     }
927
928     }
929
930     }
931
932     }
933
934     }
935
936     }
937
938     }
939
940     }
941
942     }
943
944     }
945
946     }
947
948     }
949
950     }
951
952     }
953
954     }
955
956     }
957
958     }
959
960     }
961
962     }
963
964     }
965
966     }
967
968     }
969
970     }
971
972     }
973
974     }
975
976     }
977
978     }
979
980     }
981
982     }
983
984     }
985
986     }
987
988     }
989
990     }
991
992     }
993
994     }
995
996     }
997
998     }
999
1000    }

```

```

527     }
528 #endif /* __i386 */
529
530     } else if (type == T_SGLSTP && lwp != NULL)
531         lwp->lwp_pcb.pcb_drstat = (uintptr_t)addr;
532
533     if (tdebug)
534         showregs(type, rp, addr);
535
536     if (USERMODE(rp->r_cs)) {
537         /*
538          * Set up the current cred to use during this trap. u_cred
539          * no longer exists. t_cred is used instead.
540          * The current process credential applies to the thread for
541          * the entire trap. If trapping from the kernel, this
542          * should already be set up.
543          */
544         if (ct->t_cred != p->p_cred) {
545             cred_t *oldcred = ct->t_cred;
546             /*
547              * DTrace accesses t_cred in probe context. t_cred
548              * must always be either NULL, or point to a valid,
549              * allocated cred structure.
550              */
551             ct->t_cred = crgetcred();
552             crfree(oldcred);
553         }
554         ASSERT(lwp != NULL);
555         type |= USER;
556         ASSERT(lwptoregs(lwp) == rp);
557         lwp->lwp_state = LWP_SYS;
558
559         switch (type) {
560         case T_PGFLT + USER:
561             if ((caddr_t)rp->r_pc == addr)
562                 mstate = LMS_TFAULT;
563             else
564                 mstate = LMS_DFAULT;
565             break;
566         default:
567             mstate = LMS_TRAP;
568             break;
569         }
570         /* Kernel probe */
571         TNF_PROBE_1(thread_state, "thread", /* CSTYLED */,
572                   tnf_microstate, state, mstate);
573         mstate = new_mstate(ct, mstate);
574
575         bzero(&siginfo, sizeof (siginfo));
576     }
577
578     switch (type) {
579     case T_PGFLT + USER:
580     case T_SGLSTP:
581     case T_SGLSTP + USER:
582     case T_BPTFLT + USER:
583         break;
584
585     default:
586         FTRACE_2("trap(): type=0x%lx, regs=0x%lx",
587                (ulong_t)type, (ulong_t)rp);
588         break;
589     }
590
591     switch (type) {
592     case T_SIMDFPE:
593
594     }
595
596     }
597
598     }
599
600     }
601
602     }
603
604     }
605
606     }
607
608     }
609
610     }
611
612     }
613
614     }
615
616     }
617
618     }
619
620     }
621
622     }
623
624     }
625
626     }
627
628     }
629
630     }
631
632     }
633
634     }
635
636     }
637
638     }
639
640     }
641
642     }
643
644     }
645
646     }
647
648     }
649
650     }
651
652     }
653
654     }
655
656     }
657
658     }
659
660     }
661
662     }
663
664     }
665
666     }
667
668     }
669
670     }
671
672     }
673
674     }
675
676     }
677
678     }
679
680     }
681
682     }
683
684     }
685
686     }
687
688     }
689
690     }
691
692     }
693
694     }
695
696     }
697
698     }
699
700     }
701
702     }
703
704     }
705
706     }
707
708     }
709
710     }
711
712     }
713
714     }
715
716     }
717
718     }
719
720     }
721
722     }
723
724     }
725
726     }
727
728     }
729
730     }
731
732     }
733
734     }
735
736     }
737
738     }
739
740     }
741
742     }
743
744     }
745
746     }
747
748     }
749
750     }
751
752     }
753
754     }
755
756     }
757
758     }
759
760     }
761
762     }
763
764     }
765
766     }
767
768     }
769
770     }
771
772     }
773
774     }
775
776     }
777
778     }
779
780     }
781
782     }
783
784     }
785
786     }
787
788     }
789
790     }
791
792     }
793
794     }
795
796     }
797
798     }
799
800     }
801
802     }
803
804     }
805
806     }
807
808     }
809
810     }
811
812     }
813
814     }
815
816     }
817
818     }
819
820     }
821
822     }
823
824     }
825
826     }
827
828     }
829
830     }
831
832     }
833
834     }
835
836     }
837
838     }
839
840     }
841
842     }
843
844     }
845
846     }
847
848     }
849
850     }
851
852     }
853
854     }
855
856     }
857
858     }
859
860     }
861
862     }
863
864     }
865
866     }
867
868     }
869
870     }
871
872     }
873
874     }
875
876     }
877
878     }
879
880     }
881
882     }
883
884     }
885
886     }
887
888     }
889
890     }
891
892     }
893
894     }
895
896     }
897
898     }
899
900     }
901
902     }
903
904     }
905
906     }
907
908     }
909
910     }
911
912     }
913
914     }
915
916     }
917
918     }
919
920     }
921
922     }
923
924     }
925
926     }
927
928     }
929
930     }
931
932     }
933
934     }
935
936     }
937
938     }
939
940     }
941
942     }
943
944     }
945
946     }
947
948     }
949
950     }
951
952     }
953
954     }
955
956     }
957
958     }
959
960     }
961
962     }
963
964     }
965
966     }
967
968     }
969
970     }
971
972     }
973
974     }
975
976     }
977
978     }
979
980     }
981
982     }
983
984     }
985
986     }
987
988     }
989
990     }
991
992     }
993
994     }
995
996     }
997
998     }
999
1000    }

```

```

593     /* Make sure we enable interrupts before die()ing */
594     sti(); /* The SIMD exception comes in via cmnintrap */
595     /*FALLTHROUGH*/
596 default:
597     if (type & USER) {
598         if (tudebug)
599             showregs(type, rp, (caddr_t)0);
600         printf("trap: Unknown trap type %d in user mode\n",
601             type & ~USER);
602         siginfo.si_signo = SIGILL;
603         siginfo.si_code = ILL_ILLLTRP;
604         siginfo.si_addr = (caddr_t)rp->r_pc;
605         siginfo.si_trapno = type & ~USER;
606         fault = FLTILL;
607         break;
608     } else {
609         (void) die(type, rp, addr, cpuid);
610         /*NOTREACHED*/
611     }
612
613 case T_PGFLT: /* system page fault */
614     /*
615     * If we're under on_trap() protection (see <sys/ontrap.h>),
616     * set ot_trap and bounce back to the on_trap() call site
617     * via the installed trampoline.
618     */
619     if ((ct->t_ontrap != NULL) &&
620         (ct->t_ontrap->ot_prot & OT_DATA_ACCESS)) {
621         ct->t_ontrap->ot_trap |= OT_DATA_ACCESS;
622         rp->r_pc = ct->t_ontrap->ot_trampoline;
623         goto cleanup;
624     }
625
626     /*
627     * If we have an Instruction fault in kernel mode, then that
628     * means we've tried to execute a user page (SMEP) or both of
629     * PAE and NXE are enabled. In either case, given that it's a
630     * kernel fault, we should panic immediately and not try to make
631     * any more forward progress. This indicates a bug in the
632     * kernel, which if execution continued, could be exploited to
633     * wreak havoc on the system.
634     */
635     if (errcode & PF_ERR_EXEC) {
636         (void) die(type, rp, addr, cpuid);
637     }
638
639     /*
640     * We need to check if SMAP is in play. If SMAP is in play, then
641     * any access to a user page will show up as a protection
642     * violation. To see if SMAP is enabled we first check if it's a
643     * user address and whether we have the feature flag set. If we
644     * do and the interrupted registers do not allow for user
645     * accesses (PS_ACHK is not enabled), then we need to die
646     * immediately.
647     */
648     if (addr < (caddr_t)kernelbase &&
649         is_x86_feature(x86_featureset, X86FSET_SMAP) == B_TRUE &&
650         (rp->r_ps & PS_ACHK) == 0) {
651         (void) die(type, rp, addr, cpuid);
652     }
653
654     /*
655     * See if we can handle as pagefault. Save lofault and onfault
656     * across this. Here we assume that an address less than
657     * KERNELBASE is a user fault. We can do this as copy.s
658     * routines verify that the starting address is less than

```

```

659     * KERNELBASE before starting and because we know that we
660     * always have KERNELBASE mapped as invalid to serve as a
661     * "barrier".
662     */
663     lofault = ct->t_lofault;
664     onfault = ct->t_onfault;
665     ct->t_lofault = 0;
666
667     mstate = new_mstate(ct, LMS_KFAULT);
668
669     if (addr < (caddr_t)kernelbase) {
670         res = pagefault(addr,
671             (errcode & PF_ERR_PROT)? F_PROT: F_INVAL, rw, 0);
672         if (res == FC_NOMAP &&
673             addr < p->p_usrstack &&
674             grow(addr))
675             res = 0;
676     } else {
677         res = pagefault(addr,
678             (errcode & PF_ERR_PROT)? F_PROT: F_INVAL, rw, 1);
679     }
680     (void) new_mstate(ct, mstate);
681
682     /*
683     * Restore lofault and onfault. If we resolved the fault, exit.
684     * If we didn't and lofault wasn't set, die.
685     */
686     ct->t_lofault = lofault;
687     ct->t_onfault = onfault;
688     if (res == 0)
689         goto cleanup;
690
691 #if defined(OPTERON_ERRATUM_93) && defined(LP64)
692     if (lofault == 0 && opteron_erratum_93) {
693         /*
694         * Workaround for Opteron Erratum 93. On return from
695         * a System Management Interrupt at a HLT instruction
696         * the %rip might be truncated to a 32 bit value.
697         * BIOS is supposed to fix this, but some don't.
698         * If this occurs we simply restore the high order bits.
699         * The HLT instruction is 1 byte of 0xf4.
700         */
701         uintptr_t rip = rp->r_pc;
702
703         if ((rip & 0xfffffffful) == rip) {
704             rip |= 0xfffffffful << 32;
705             if (hat_getpfnum(kas.a_hat, (caddr_t)rip) !=
706                 PFN_INVALID &&
707                 (*(uchar_t *)rip == 0xf4 ||
708                 *(uchar_t *) (rip - 1) == 0xf4)) {
709                 rp->r_pc = rip;
710                 goto cleanup;
711             }
712         }
713     }
714 #endif /* OPTERON_ERRATUM_93 && LP64 */
715
716 #ifdef OPTERON_ERRATUM_91
717     if (lofault == 0 && opteron_erratum_91) {
718         /*
719         * Workaround for Opteron Erratum 91. Prefetches may
720         * generate a page fault (they're not supposed to do
721         * that!). If this occurs we simply return back to the
722         * instruction.
723         */
724         caddr_t pc = (caddr_t)rp->r_pc;

```

```

726     /*
727     * If the faulting PC is not mapped, this is a
728     * legitimate kernel page fault that must result in a
729     * panic. If the faulting PC is mapped, it could contain
730     * a prefetch instruction. Check for that here.
731     */
732     if (hat_getpfnum(kas.a_hat, pc) != PFN_INVALID) {
733         if (cmp_to_prefetch((uchar_t *)pc)) {
734             #ifdef DEBUG
735                 cmn_err(CE_WARN, "Opteron erratum 91 "
736                     "occurred: kernel prefetch"
737                     " at %p generated a page fault!",
738                     (void *)rp->r_pc);
739             #endif /* DEBUG */
740             goto cleanup;
741         }
742     }
743     (void) die(type, rp, addr, cpuid);
744 }
745 #endif /* OPTERON_ERRATUM_91 */

747     if (lofault == 0)
748         (void) die(type, rp, addr, cpuid);

750     /*
751     * Cannot resolve fault. Return to lofault.
752     */
753     if (lodebug) {
754         showregs(type, rp, addr);
755         traceregs(rp);
756     }
757     if (FC_CODE(res) == FC_OBJERR)
758         res = FC_ERRNO(res);
759     else
760         res = EFAULT;
761     rp->r_r0 = res;
762     rp->r_pc = ct->t_lofault;
763     goto cleanup;

765     case T_PGFLT + USER: /* user page fault */
766         if (faultdebug) {
767             char *fault_str;

769             switch (rw) {
770                 case S_READ:
771                     fault_str = "read";
772                     break;
773                 case S_WRITE:
774                     fault_str = "write";
775                     break;
776                 case S_EXEC:
777                     fault_str = "exec";
778                     break;
779                 default:
780                     fault_str = "";
781                     break;
782             }
783             printf("user %s fault: addr=0x%lx errcode=0x%x\n",
784                 fault_str, (uintptr_t)addr, errcode);
785         }

787     #if defined(OPTERON_ERRATUM_100) && defined(LP64)
788         /*
789         * Workaround for AMD erratum 100
790         */

```

```

791     * A 32-bit process may receive a page fault on a non
792     * 32-bit address by mistake. The range of the faulting
793     * address will be
794     *
795     * 0xffffffff80000000 .. 0xffffffffffffffff or
796     * 0x0000000100000000 .. 0x000000017fffffff
797     *
798     * The fault is always due to an instruction fetch, however
799     * the value of r_pc should be correct (in 32 bit range),
800     * so we ignore the page fault on the bogus address.
801     */
802     if (p->p_model == DATAMODEL_ILP32 &&
803         (0xffffffff80000000 <= (uintptr_t)addr ||
804         (0x100000000 <= (uintptr_t)addr &&
805         (uintptr_t)addr <= 0x17fffffff))) {
806         if (!opteron_erratum_100)
807             panic("unexpected erratum #100");
808         if (rp->r_pc <= 0xffffffff)
809             goto out;
810     }
811 #endif /* OPTERON_ERRATUM_100 && LP64 */

813     ASSERT(!(curthread->t_flag & T_WATCHPT));
814     watchpage = (pr_watch_active(p) && pr_is_watchpage(addr, rw));
815 #ifdef __i386
816     /*
817     * In 32-bit mode, the lcall (system call) instruction fetches
818     * one word from the stack, at the stack pointer, because of the
819     * way the call gate is constructed. This is a bogus
820     * read and should not be counted as a read watchpoint.
821     * We work around the problem here by testing to see if
822     * this situation applies and, if so, simply jumping to
823     * the code in locore.s that fields the system call trap.
824     * The registers on the stack are already set up properly
825     * due to the match between the call gate sequence and the
826     * trap gate sequence. We just have to adjust the pc.
827     */
828     if (watchpage && addr == (caddr_t)rp->r_sp &&
829         rw == S_READ && instr_is_lcall_syscall((caddr_t)rp->r_pc)) {
830         extern void watch_syscall(void);

832         rp->r_pc += LCALLSIZE;
833         watch_syscall(); /* never returns */
834         /* NOTREACHED */
835     }
836 #endif /* __i386 */
837     vaddr = addr;
838     if (!watchpage || (sz = instr_size(rp, &vaddr, rw)) <= 0)
839         fault_type = (errcode & PF_ERR_PROT)? F_PROT: F_INVALID;
840     else if ((watchcode = pr_is_watchpoint(&vaddr, &ta,
841         sz, NULL, rw)) != 0) {
842         if (ta) {
843             do_watch_step(vaddr, sz, rw,
844                 watchcode, rp->r_pc);
845             fault_type = F_INVALID;
846         } else {
847             bzero(&siginfo, sizeof(siginfo));
848             siginfo.si_signo = SIGTRAP;
849             siginfo.si_code = watchcode;
850             siginfo.si_addr = vaddr;
851             siginfo.si_trapafter = 0;
852             siginfo.si_pc = (caddr_t)rp->r_pc;
853             fault = FLTWATCH;
854             break;
855         }
856     } else {

```

```

857     /* XXX pr_watch_emul() never succeeds (for now) */
858     if (rw != S_EXEC && pr_watch_emul(rp, vaddr, rw))
859         goto out;
860     do_watch_step(vaddr, sz, rw, 0, 0);
861     fault_type = F_INVAL;
862 }

864 res = pagefault(addr, fault_type, rw, 0);

866 /*
867  * If pagefault() succeeded, ok.
868  * Otherwise attempt to grow the stack.
869  */
870 if (res == 0 ||
871     (res == FC_NOMAP &&
872     addr < p->p_usrstack &&
873     grow(addr))) {
874     lwp->lwp_lastfault = FLTPAGE;
875     lwp->lwp_lastfaddr = addr;
876     if (prismember(&p->p_fltmask, FLTPAGE)) {
877         bzero(&siginfo, sizeof (siginfo));
878         siginfo.si_addr = addr;
879         (void) stop_on_fault(FLTPAGE, &siginfo);
880     }
881     goto out;
882 } else if (res == FC_PROT && addr < p->p_usrstack &&
883     (mmu.pt_nx != 0 && (errcode & PF_ERR_EXEC))) {
884     report_stack_exec(p, addr);
885 }

887 #ifdef OPTERON_ERRATUM_91
888 /*
889  * Workaround for Opteron Erratum 91. Prefetches may generate a
890  * page fault (they're not supposed to do that!). If this
891  * occurs we simply return back to the instruction.
892  *
893  * We rely on copyin to properly fault in the page with r_pc.
894  */
895 if (opteron_erratum_91 &&
896     addr != (caddr_t)rp->r_pc &&
897     instr_is_prefetch((caddr_t)rp->r_pc)) {
898 #ifdef DEBUG
899     cmn_err(CE_WARN, "Opteron erratum 91 occurred: "
900         "prefetch at %p in pid %d generated a trap!",
901         (void *)rp->r_pc, p->p_pid);
902 #endif /* DEBUG */
903     goto out;
904 }
905 #endif /* OPTERON_ERRATUM_91 */

907 if (tudebug)
908     showregs(type, rp, addr);
909 /*
910  * In the case where both pagefault and grow fail,
911  * set the code to the value provided by pagefault.
912  * We map all errors returned from pagefault() to SIGSEGV.
913  */
914 bzero(&siginfo, sizeof (siginfo));
915 siginfo.si_addr = addr;
916 switch (FC_CODE(res)) {
917 case FC_HWERR:
918 case FC_NOSUPPORT:
919     siginfo.si_signo = SIGBUS;
920     siginfo.si_code = BUS_ADRERR;
921     fault = FLTACCESS;
922     break;

```

```

923 case FC_ALIGN:
924     siginfo.si_signo = SIGBUS;
925     siginfo.si_code = BUS_ADRALN;
926     fault = FLTACCESS;
927     break;
928 case FC_OBJERR:
929     if ((siginfo.si_errno = FC_ERRNO(res)) != EINTR) {
930         siginfo.si_signo = SIGBUS;
931         siginfo.si_code = BUS_OBJERR;
932         fault = FLTACCESS;
933     }
934     break;
935 default: /* FC_NOMAP or FC_PROT */
936     siginfo.si_signo = SIGSEGV;
937     siginfo.si_code =
938         (res == FC_NOMAP)? SEGV_MAPERR : SEGV_ACCERR;
939     fault = FLTBOUNDS;
940     break;
941 }
942 break;

944 case T_ILLINST + USER: /* invalid opcode fault */
945 /*
946  * If the syscall instruction is disabled due to LDT usage, a
947  * user program that attempts to execute it will trigger a #ud
948  * trap. Check for that case here. If this occurs on a CPU which
949  * doesn't even support syscall, the result of all of this will
950  * be to emulate that particular instruction.
951  */
952 if (p->p_ldt != NULL &&
953     ldt_rewrite_syscall(rp, p, X86FSET_ASYS))
954     goto out;

956 #ifdef __amd64
957 /*
958  * Emulate the LAHF and SAHF instructions if needed.
959  * See the instr_is_lsahf function for details.
960  */
961 if (p->p_model == DATAMODEL_LP64 &&
962     instr_is_lsahf((caddr_t)rp->r_pc, &instr)) {
963     emulate_lsahf(rp, instr);
964     goto out;
965 }
966 #endif

968 /*FALLTHROUGH*/

970 if (tudebug)
971     showregs(type, rp, (caddr_t)0);
972 siginfo.si_signo = SIGILL;
973 siginfo.si_code = ILL_ILLOPC;
974 siginfo.si_addr = (caddr_t)rp->r_pc;
975 fault = FLTILL;
976 break;

978 case T_ZERODIV + USER: /* integer divide by zero */
979     if (tudebug && tudebugfpe)
980         showregs(type, rp, (caddr_t)0);
981     siginfo.si_signo = SIGFPE;
982     siginfo.si_code = FPE_INTDIV;
983     siginfo.si_addr = (caddr_t)rp->r_pc;
984     fault = FLTIZDIV;
985     break;

987 case T_OVFLW + USER: /* integer overflow */
988     if (tudebug && tudebugfpe)

```

```

989     showregs(type, rp, (caddr_t)0);
990     siginfo.si_signo = SIGFPE;
991     siginfo.si_code = FPE_INTOVF;
992     siginfo.si_addr = (caddr_t)rp->r_pc;
993     fault = FLTIOVF;
994     break;

996 case T_NOEXTFLT + USER: /* math coprocessor not available */
997     if (tudebug && tudebugfpe)
998         showregs(type, rp, addr);
999     if (fpnoextflt(rp)) {
1000         siginfo.si_signo = SIGILL;
1001         siginfo.si_code = ILL_ILLOPC;
1002         siginfo.si_addr = (caddr_t)rp->r_pc;
1003         fault = FLTILL;
1004     }
1005     break;

1007 case T_EXTOVRFALT: /* extension overrun fault */
1008     /* check if we took a kernel trap on behalf of user */
1009     {
1010         extern void ndptrap_frstor(void);
1011         if (rp->r_pc != (uintptr_t)ndptrap_frstor) {
1012             sti(); /* T_EXTOVRFALT comes in via cmninttrap */
1013             (void) die(type, rp, addr, cpuid);
1014         }
1015         type |= USER;
1016     }
1017     /*FALLTHROUGH*/
1018 case T_EXTOVRFALT + USER: /* extension overrun fault */
1019     if (tudebug && tudebugfpe)
1020         showregs(type, rp, addr);
1021     if (fpextovrfalt(rp)) {
1022         siginfo.si_signo = SIGSEGV;
1023         siginfo.si_code = SEGV_MAPERR;
1024         siginfo.si_addr = (caddr_t)rp->r_pc;
1025         fault = FLTBOUNDS;
1026     }
1027     break;

1029 case T_EXTERRFLT: /* x87 floating point exception pending */
1030     /* check if we took a kernel trap on behalf of user */
1031     {
1032         extern void ndptrap_frstor(void);
1033         if (rp->r_pc != (uintptr_t)ndptrap_frstor) {
1034             sti(); /* T_EXTERRFLT comes in via cmninttrap */
1035             (void) die(type, rp, addr, cpuid);
1036         }
1037         type |= USER;
1038     }
1039     /*FALLTHROUGH*/

1041 case T_EXTERRFLT + USER: /* x87 floating point exception pending */
1042     if (tudebug && tudebugfpe)
1043         showregs(type, rp, addr);
1044     if (sicode = fpexterrflt(rp)) {
1045         siginfo.si_signo = SIGFPE;
1046         siginfo.si_code = sicode;
1047         siginfo.si_addr = (caddr_t)rp->r_pc;
1048         fault = FLTFPE;
1049     }
1050     break;

1052 case T_SIMDFPE + USER: /* SSE and SSE2 exceptions */
1053     if (tudebug && tudebugsse)
1054         showregs(type, rp, addr);

```

```

1055     if (!is_x86_feature(x86_featureset, X86FSET_SSE) &&
1056         !is_x86_feature(x86_featureset, X86FSET_SSE2)) {
1057         /*
1058          * There are rumours that some user instructions
1059          * on older CPUs can cause this trap to occur; in
1060          * which case send a SIGILL instead of a SIGFPE.
1061          */
1062         siginfo.si_signo = SIGILL;
1063         siginfo.si_code = ILL_ILLTRP;
1064         siginfo.si_addr = (caddr_t)rp->r_pc;
1065         siginfo.si_trapno = type & ~USER;
1066         fault = FLTILL;
1067     } else if ((sicode = fpsimderrflt(rp)) != 0) {
1068         siginfo.si_signo = SIGFPE;
1069         siginfo.si_code = sicode;
1070         siginfo.si_addr = (caddr_t)rp->r_pc;
1071         fault = FLTFPE;
1072     }

1074     sti(); /* The SIMD exception comes in via cmninttrap */
1075     break;

1077 case T_BPTFLT: /* breakpoint trap */
1078     /*
1079      * Kernel breakpoint traps should only happen when kmdb is
1080      * active, and even then, it'll have interposed on the IDT, so
1081      * control won't get here. If it does, we've hit a breakpoint
1082      * without the debugger, which is very strange, and very
1083      * fatal.
1084      */
1085     if (tudebug && tudebugbpt)
1086         showregs(type, rp, (caddr_t)0);

1088     (void) die(type, rp, addr, cpuid);
1089     break;

1091 case T_SGLSTP: /* single step/hw breakpoint exception */

1093 #if !defined(__xpv)
1094     /* Now evaluate how we got here */
1095     if (lwp != NULL && (lwp->lwp_pcb.pcb_drstat & DR_SINGLESTEP)) {
1096         /*
1097          * We'd never normally get here, as kmdb handles its own single
1098          * step traps. There is one nasty exception though, as
1099          * described in more detail in sys_sysenter(). Note that
1100          * checking for all four locations covers both the KPTI and the
1101          * non-KPTI cases correctly: the former will never be found at
1102          * (brand_)sys_sysenter, and vice versa.
1103          * i386 single-steps even through lcalls which
1104          * change the privilege level. So we take a trap at
1105          * the first instruction in privileged mode.
1106          *
1107          * Set a flag to indicate that upon completion of
1108          * the system call, deal with the single-step trap.
1109          *
1110          * The same thing happens for sysenter, too.
1111          */
1112         if (lwp != NULL && (lwp->lwp_pcb.pcb_drstat & DR_SINGLESTEP)) {
1113             if (rp->r_pc == (greg_t)brand_sys_sysenter ||
1114                 rp->r_pc == (greg_t)sys_sysenter ||
1115                 rp->r_pc == (greg_t)tr_brand_sys_sysenter ||
1116                 rp->r_pc == (greg_t)tr_sys_sysenter) {
1117
1118                 rp->r_pc += 0x3; /* sizeof (swaps) */
1119
1120                 singlestep_twiddle = 0;

```



```

1107         if (rp->r_pc == (uintptr_t)sys_sysenter ||
1108             rp->r_pc == (uintptr_t)brand_sys_sysenter) {
1109             singlestep_twiddle = 1;
1110 #if defined(__amd64)
1111             /*
1112              * Since we are already on the kernel's
1113              * %gs, on 64-bit systems the sysenter case
1114              * needs to adjust the pc to avoid
1115              * executing the swapgs instruction at the
1116              * top of the handler.
1117              */
1118             if (rp->r_pc == (uintptr_t)sys_sysenter)
1119                 rp->r_pc = (uintptr_t)
1120                     _sys_sysenter_post_swapgs;
1121             else
1122                 rp->r_pc = (uintptr_t)
1123                     _brand_sys_sysenter_post_swapgs;
1124 #endif
1125         }
1126 #if defined(__i386)
1127         else if (rp->r_pc == (uintptr_t)sys_call ||
1128             rp->r_pc == (uintptr_t)brand_sys_call) {
1129             singlestep_twiddle = 1;
1130         }
1131 #endif
1132         else {
1133             /* not on sysenter/syscall; uregs available */
1134             if (tudebug && tudebugbpt)
1135                 showregs(type, rp, (caddr_t)0);
1136         }
1137         if (singlestep_twiddle) {
1138             rp->r_ps &= ~PS_T; /* turn off trace */
1139             lwp->lwp_pcb.pcb_flags |= DEBUG_PENDING;
1140             ct->t_post_sys = 1;
1141             aston(curthread);
1142             goto cleanup;
1143         } else {
1144             if (tudebug && tudebugbpt)
1145                 showregs(type, rp, (caddr_t)0);
1146         }
1147     }
1148 #endif /* !__xpv */
1149
1150 /* XXX - needs review on debugger interface? */
1151 if (boothowto & RB_DEBUG)
1152     debug_enter((char *)NULL);
1153 else
1154     (void) die(type, rp, addr, cpuid);
1155 break;
1156
1157 case T_NMIFLT: /* NMI interrupt */
1158     printf("Unexpected NMI in system mode\n");
1159     goto cleanup;
1160
1161 case T_NMIFLT + USER: /* NMI interrupt */
1162     printf("Unexpected NMI in user mode\n");
1163     break;
1164
1165 case T_GPFLT: /* general protection violation */
1166     /*
1167      * Any #GP that occurs during an on_trap .. no_trap bracket
1168      * with OT_DATA_ACCESS or OT_SEGMENT_ACCESS protection,
1169      * or in a on_fault .. no_fault bracket, is forgiven
1170      * and we trampoline. This protection is given regardless
1171      * of whether we are 32/64 bit etc - if a distinction is
1172      * required then define new on_trap protection types.

```

```

1144         *
1145         * On amd64, we can get a #gp from referencing addresses
1146         * in the virtual address hole e.g. from a copyin or in
1147         * update_sregs while updating user segment registers.
1148         *
1149         * On the 32-bit hypervisor we could also generate one in
1150         * mfn_to_pfn by reaching around or into where the hypervisor
1151         * lives which is protected by segmentation.
1152         */
1153
1154     /*
1155      * If we're under on_trap() protection (see <sys/ontrap.h>),
1156      * set ot_trap and trampoline back to the on_trap() call site
1157      * for OT_DATA_ACCESS or OT_SEGMENT_ACCESS.
1158      */
1159     if (ct->t_ontrap != NULL) {
1160         int ttype = ct->t_ontrap->ot_prot &
1161             (OT_DATA_ACCESS | OT_SEGMENT_ACCESS);
1162
1163         if (ttype != 0) {
1164             ct->t_ontrap->ot_trap |= ttype;
1165             if (tudebug)
1166                 showregs(type, rp, (caddr_t)0);
1167             rp->r_pc = ct->t_ontrap->ot_trampoline;
1168             goto cleanup;
1169         }
1170     }
1171
1172     /*
1173      * If we're under lofault protection (copyin etc.),
1174      * longjmp back to lofault with an EFAULT.
1175      */
1176     if (ct->t_lofault) {
1177         /*
1178          * Fault is not resolvable, so just return to lofault
1179          */
1180         if (lodebug) {
1181             showregs(type, rp, addr);
1182             traceregs(rp);
1183         }
1184         rp->r_r0 = EFAULT;
1185         rp->r_pc = ct->t_lofault;
1186         goto cleanup;
1187     }
1188
1189     /*
1190      * We fall through to the next case, which repeats
1191      * the OT_SEGMENT_ACCESS check which we've already
1192      * done, so we'll always fall through to the
1193      * T_STKFLT case.
1194      */
1195     /*FALLTHROUGH*/
1196     case T_SEGFLT: /* segment not present fault */
1197         /*
1198          * One example of this is #NP in update_sregs while
1199          * attempting to update a user segment register
1200          * that points to a descriptor that is marked not
1201          * present.
1202          */
1203         if (ct->t_ontrap != NULL &&
1204             ct->t_ontrap->ot_prot & OT_SEGMENT_ACCESS) {
1205             ct->t_ontrap->ot_trap |= OT_SEGMENT_ACCESS;
1206             if (tudebug)
1207                 showregs(type, rp, (caddr_t)0);
1208             rp->r_pc = ct->t_ontrap->ot_trampoline;
1209             goto cleanup;

```

```

1210     }
1211     /*FALLTHROUGH*/
1212     case T_STKFLT: /* stack fault */
1213     case T_TSSFLT: /* invalid TSS fault */
1214         if (tudebug)
1215             showregs(type, rp, (caddr_t)0);
1216         if (kern_gpfault(rp))
1217             (void) die(type, rp, addr, cpuid);
1218         goto cleanup;
1219
1220 /*
1221  * ONLY 32-bit PROCESSES can USE a PRIVATE LDT! 64-bit apps
1222  * should have no need for them, so we put a stop to it here.
1223  *
1224  * So: not-present fault is ONLY valid for 32-bit processes with
1225  * a private LDT trying to do a system call. Emulate it.
1226  *
1227  * #gp fault is ONLY valid for 32-bit processes also, which DO NOT
1228  * have a private LDT, and are trying to do a system call. Emulate it.
1229  */
1230
1231     case T_SEGFLT + USER: /* segment not present fault */
1232     case T_GPFLT + USER: /* general protection violation */
1233 #ifdef _SYSCALL32_IMPL
1234         if (p->p_model != DATAMODEL_NATIVE) {
1235 #endif /* _SYSCALL32_IMPL */
1236             if (instr_is_lcall_syscall((caddr_t)rp->r_pc)) {
1237                 if (type == T_SEGFLT + USER)
1238                     ASSERT(p->p_ldt != NULL);
1239
1240             if ((p->p_ldt == NULL && type == T_GPFLT + USER) ||
1241                 type == T_SEGFLT + USER) {
1242
1243                 /*
1244                  * The user attempted a system call via the obsolete
1245                  * call gate mechanism. Because the process doesn't have
1246                  * an LDT (i.e. the ldtr contains 0), a #gp results.
1247                  * Emulate the syscall here, just as we do above for a
1248                  * #np trap.
1249                  */
1250
1251                 /*
1252                  * Since this is a not-present trap, rp->r_pc points to
1253                  * the trapping lcall instruction. We need to bump it
1254                  * to the next insn so the app can continue on.
1255                  */
1256                 rp->r_pc += LCALLSIZE;
1257                 lwp->lwp_regs = rp;
1258
1259                 /*
1260                  * Normally the microstate of the LWP is forced back to
1261                  * LMS_USER by the syscall handlers. Emulate that
1262                  * behavior here.
1263                  */
1264                 mstate = LMS_USER;
1265
1266                 dosyscall();
1267                 goto out;
1268             }
1269         }
1270 #ifdef _SYSCALL32_IMPL
1271     }
1272 #endif /* _SYSCALL32_IMPL */
1273     /*
1274     * If the current process is using a private LDT and the
1275     * trapping instruction is sysenter, the sysenter instruction

```

```

1276     * has been disabled on the CPU because it destroys segment
1277     * registers. If this is the case, rewrite the instruction to
1278     * be a safe system call and retry it. If this occurs on a CPU
1279     * which doesn't even support sysenter, the result of all of
1280     * this will be to emulate that particular instruction.
1281     */
1282     if (p->p_ldt != NULL &&
1283         ldt_rewrite_syscall(rp, p, X86FSET_SEP))
1284         goto out;
1285
1286     /*FALLTHROUGH*/
1287
1288     case T_BOUNDFLT + USER: /* bound fault */
1289     case T_STKFLT + USER: /* stack fault */
1290     case T_TSSFLT + USER: /* invalid TSS fault */
1291         if (tudebug)
1292             showregs(type, rp, (caddr_t)0);
1293         siginfo.si_signo = SIGSEGV;
1294         siginfo.si_code = SEGV_MAPERR;
1295         siginfo.si_addr = (caddr_t)rp->r_pc;
1296         fault = FLTBOUNDS;
1297         break;
1298
1299     case T_ALIGNMENT + USER: /* user alignment error (486) */
1300         if (tudebug)
1301             showregs(type, rp, (caddr_t)0);
1302         bzero(&siginfo, sizeof(siginfo));
1303         siginfo.si_signo = SIGBUS;
1304         siginfo.si_code = BUS_ADRALN;
1305         siginfo.si_addr = (caddr_t)rp->r_pc;
1306         fault = FLTACCESS;
1307         break;
1308
1309     case T_SGLSTP + USER: /* single step/hw breakpoint exception */
1310         if (tudebug && tudebugbpt)
1311             showregs(type, rp, (caddr_t)0);
1312
1313         /* Was it single-stepping? */
1314         if (lwp->lwp_pcb.pcb_drstat & DR_SINGLESTEP) {
1315             pcb_t *pcb = &lwp->lwp_pcb;
1316
1317             rp->r_ps &= ~PS_T;
1318             /*
1319              * If both NORMAL_STEP and WATCH_STEP are in effect,
1320              * give precedence to WATCH_STEP. If neither is set,
1321              * user must have set the PS_T bit in %efl; treat this
1322              * as NORMAL_STEP.
1323              */
1324             if ((fault = undo_watch_step(&siginfo)) == 0 &&
1325                 ((pcb->pcb_flags & NORMAL_STEP) ||
1326                  !(pcb->pcb_flags & WATCH_STEP))) {
1327                 siginfo.si_signo = SIGTRAP;
1328                 siginfo.si_code = TRAP_TRACE;
1329                 siginfo.si_addr = (caddr_t)rp->r_pc;
1330                 fault = FLTRTRACE;
1331             }
1332             pcb->pcb_flags &= ~(NORMAL_STEP|WATCH_STEP);
1333         }
1334         break;
1335
1336     case T_BPTFLT + USER: /* breakpoint trap */
1337         if (tudebug && tudebugbpt)
1338             showregs(type, rp, (caddr_t)0);
1339         /*
1340         * int 3 (the breakpoint instruction) leaves the pc referring
1341         * to the address one byte after the breakpointed address.

```

```

1342     * If the P_PR_BPTADJ flag has been set via /proc, We adjust
1343     * it back so it refers to the breakpointed address.
1344     */
1345     if (p->p_proc_flag & P_PR_BPTADJ)
1346         rp->r_pc--;
1347     siginfo.si_signo = SIGTRAP;
1348     siginfo.si_code = TRAP_BRKPT;
1349     siginfo.si_addr = (caddr_t)rp->r_pc;
1350     fault = FLTBPT;
1351     break;

1353 case T_AST:
1354     /*
1355     * This occurs only after the cs register has been made to
1356     * look like a kernel selector, either through debugging or
1357     * possibly by functions like setcontext(). The thread is
1358     * about to cause a general protection fault at common_iret()
1359     * in locore. We let that happen immediately instead of
1360     * doing the T_AST processing.
1361     */
1362     goto cleanup;

1364 case T_AST + USER: /* profiling, resched, h/w error pseudo trap */
1365     if (lwp->lwp_pcb.pcb_flags & ASYNC_HWERR) {
1366         proc_t *p = ttproc(curthread);
1367         extern void print_msg_hwerr(ctid_t ct_id, proc_t *p);

1369         lwp->lwp_pcb.pcb_flags &= ~ASYNC_HWERR;
1370         print_msg_hwerr(p->p_ct_process->conp_contract.ct_id,
1371             p);
1372         contract_process_hwerr(p->p_ct_process, p);
1373         siginfo.si_signo = SIGKILL;
1374         siginfo.si_code = SI_NOINFO;
1375     } else if (lwp->lwp_pcb.pcb_flags & CPC_OVERFLOW) {
1376         lwp->lwp_pcb.pcb_flags &= ~CPC_OVERFLOW;
1377         if (kcpc_overflow_ast()) {
1378             /*
1379             * Signal performance counter overflow
1380             */
1381             if (tudebug)
1382                 showregs(type, rp, (caddr_t)0);
1383             bzero(&siginfo, sizeof(siginfo));
1384             siginfo.si_signo = SIGEMT;
1385             siginfo.si_code = EMT_CPCOVF;
1386             siginfo.si_addr = (caddr_t)rp->r_pc;
1387             fault = FLTCPCOVF;
1388         }
1389     }

1391     break;
1392 }

1394 /*
1395  * We can't get here from a system trap
1396  */
1397 ASSERT(type & USER);

1399 if (fault) {
1400     /* We took a fault so abort single step. */
1401     lwp->lwp_pcb.pcb_flags &= ~(NORMAL_STEP|WATCH_STEP);
1402     /*
1403     * Remember the fault and fault address
1404     * for real-time (SIGPROF) profiling.
1405     */
1406     lwp->lwp_lastfault = fault;
1407     lwp->lwp_lastfaddr = siginfo.si_addr;

```

```

1409     DTRACE_PROC2(fault, int, fault, ksiginfo_t *, &siginfo);

1411     /*
1412     * If a debugger has declared this fault to be an
1413     * event of interest, stop the lwp. Otherwise just
1414     * deliver the associated signal.
1415     */
1416     if (siginfo.si_signo != SIGKILL &&
1417         prismember(&p->p_fltmask, fault) &&
1418         stop_on_fault(fault, &siginfo) == 0)
1419         siginfo.si_signo = 0;
1420     }

1422     if (siginfo.si_signo)
1423         trapsig(&siginfo, (fault != FLTFPE && fault != FLTCPCOVF));

1425     if (lwp->lwp_oweupc)
1426         profil_tick(rp->r_pc);

1428     if (ct->t_astflag | ct->t_sig_check) {
1429         /*
1430         * Turn off the AST flag before checking all the conditions that
1431         * may have caused an AST. This flag is on whenever a signal or
1432         * unusual condition should be handled after the next trap or
1433         * syscall.
1434         */
1435         astoff(ct);
1436         /*
1437         * If a single-step trap occurred on a syscall (see above)
1438         * recognize it now. Do this before checking for signals
1439         * because deferred_singlestep_trap() may generate a SIGTRAP to
1440         * the LWP or may otherwise mark the LWP to call issig(FORREAL).
1441         */
1442         if (lwp->lwp_pcb.pcb_flags & DEBUG_PENDING)
1443             deferred_singlestep_trap((caddr_t)rp->r_pc);

1445         ct->t_sig_check = 0;

1447         /*
1448         * As in other code paths that check against TP_CHANGEBIND,
1449         * we perform the check first without p_lock held -- only
1450         * acquiring p_lock in the unlikely event that it is indeed
1451         * set. This is safe because we are doing this after the
1452         * astoff(); if we are racing another thread setting
1453         * TP_CHANGEBIND on us, we will pick it up on a subsequent
1454         * lap through.
1455         */
1456         if (curthread->t_proc_flag & TP_CHANGEBIND) {
1457             mutex_enter(&p->p_lock);
1458             if (curthread->t_proc_flag & TP_CHANGEBIND) {
1459                 timer_lwpbind();
1460                 curthread->t_proc_flag &= ~TP_CHANGEBIND;
1461             }
1462             mutex_exit(&p->p_lock);
1463         }

1465         /*
1466         * for kaio requests that are on the per-process poll queue,
1467         * aio->aio_pollq, they're AIO_POLL bit is set, the kernel
1468         * should copyout their result_t to user memory. by copying
1469         * out the result_t, the user can poll on memory waiting
1470         * for the kaio request to complete.
1471         */
1472         if (p->p_aio)
1473             aio_cleanup(0);

```

```

1474     /*
1475     * If this LWP was asked to hold, call holdlwp(), which will
1476     * stop. holdlwps() sets this up and calls pokelwps() which
1477     * sets the AST flag.
1478     *
1479     * Also check TP_EXITLWP, since this is used by fresh new LWPs
1480     * through lwp_rtt(). That flag is set if the lwp_create(2)
1481     * syscall failed after creating the LWP.
1482     */
1483     if (ISHOLD(p))
1484         holdlwp();

1486     /*
1487     * All code that sets signals and makes ISSIG evaluate true must
1488     * set t_astflag afterwards.
1489     */
1490     if (ISSIG_PENDING(ct, lwp, p)) {
1491         if (issig(FORREAL))
1492             psig();
1493         ct->t_sig_check = 1;
1494     }

1496     if (ct->t_rprof != NULL) {
1497         realsigprof(0, 0, 0);
1498         ct->t_sig_check = 1;
1499     }

1501     /*
1502     * /proc can't enable/disable the trace bit itself
1503     * because that could race with the call gate used by
1504     * system calls via "lcall". If that happened, an
1505     * invalid EFLAGS would result. prstep()/prnostep()
1506     * therefore schedule an AST for the purpose.
1507     */
1508     if (lwp->lwp_pcb.pcb_flags & REQUEST_STEP) {
1509         lwp->lwp_pcb.pcb_flags &= ~REQUEST_STEP;
1510         rp->r_ps |= PS_T;
1511     }
1512     if (lwp->lwp_pcb.pcb_flags & REQUEST_NOSTEP) {
1513         lwp->lwp_pcb.pcb_flags &= ~REQUEST_NOSTEP;
1514         rp->r_ps &= ~PS_T;
1515     }
1516 }

1518 out: /* We can't get here from a system trap */
1519     ASSERT(type & USER);

1521     if (ISHOLD(p))
1522         holdlwp();

1524     /*
1525     * Set state to LWP_USER here so preempt won't give us a kernel
1526     * priority if it occurs after this point. Call CL_TRAPRET() to
1527     * restore the user-level priority.
1528     *
1529     * It is important that no locks (other than spinlocks) be entered
1530     * after this point before returning to user mode (unless lwp_state
1531     * is set back to LWP_SYS).
1532     */
1533     lwp->lwp_state = LWP_USER;

1535     if (ct->t_trapret) {
1536         ct->t_trapret = 0;
1537         thread_lock(ct);
1538         CL_TRAPRET(ct);
1539         thread_unlock(ct);

```

```

1540     }
1541     if (CPU->cpu_runrun || curthread->t_schedflag & TS_ANYWAITQ)
1542         preempt();
1543     prunstop();
1544     (void) new_mstate(ct, mstate);

1546     /* Kernel probe */
1547     TNF_PROBE_1(thread_state, "thread", /* CSTYLEL */,
1548         tnf_microstate, state, LMS_USER);

1550     return;

1552 cleanup: /* system traps end up here */
1553     ASSERT(!(type & USER));
1554 }
_____unchanged_portion_omitted_____

1668 /*
1669 * Print out debugging info.
1670 */
1671 static void
1672 showregs(uint_t type, struct regs *rp, caddr_t addr)
1673 {
1674     int s;

1676     s = spl7();
1677     type &= ~USER;
1678     if (PTOU(curproc)->u_comm[0])
1679         printf("%s: ", PTOU(curproc)->u_comm);
1680     if (type < TRAP_TYPES)
1681         printf("#%s %s\n", trap_type_mnemonic[type], trap_type[type]);
1682     else
1683         switch (type) {
1684             case T_SYSCALL:
1685                 printf("Syscall Trap:\n");
1686                 break;
1687             case T_AST:
1688                 printf("AST\n");
1689                 break;
1690             default:
1691                 printf("Bad Trap = %d\n", type);
1692                 break;
1693         }
1694     if (type == T_PGFLT) {
1695         printf("Bad %s fault at addr=0x%lx\n",
1696             USERMODE(rp->r_cs) ? "user": "kernel", (uintptr_t)addr);
1697     } else if (addr) {
1698         printf("addr=0x%lx\n", (uintptr_t)addr);
1699     }

1701     printf("pid=%d, pc=0x%lx, sp=0x%lx, eflags=0x%lx\n",
1702         ttoproc(curthread) && ttoproc(curthread)->p_pidp,
1703         ttoproc(curthread)->p_pid : 0, rp->r_pc, rp->r_sp, rp->r_ps);

1705 #if defined(__lint)
1706     /*
1707     * this clause can be deleted when lint bug 4870403 is fixed
1708     * (lint thinks that bit 32 is illegal in a %b format string)
1709     */
1710     printf("cr0: %x cr4: %b\n",
1711         (uint_t)getcr0(), (uint_t)getcr4(), FMT_CR4);
1712 #else
1713     printf("cr0: %b cr4: %b\n",
1714         (uint_t)getcr0(), FMT_CR0, (uint_t)getcr4(), FMT_CR4);
1715 #endif /* __lint */

```

```
1717     printf("cr2: %lx ", getcr2());
1741     printf("cr2: %lx", getcr2());
1718 #if !defined(__xpv)
1719     printf("cr3: %lx ", getcr3());
1743     printf("cr3: %lx", getcr3());
1720 #if defined(__amd64)
1721     printf("cr8: %lx\n", getcr8());
1722 #endif
1723 #endif
1724     printf("\n");

1726     dumpregs(rp);
1727     splx(s);
1728 }
    _____
    unchanged_portion_omitted_

1822 #endif /* __i386 */

1824 /*
1825  * Test to see if the instruction is part of _sys_rtt (or the KPTI trampolines
1826  * which are used by _sys_rtt).
1827  * Test to see if the instruction is part of _sys_rtt.
1828  * Again on the hypervisor if we try to IRET to user land with a bad code
1829  * or stack selector we will get vectored through xen_failsafe_callback.
1830  * In which case we assume we got here via _sys_rtt since we only allow
1831  * IRET to user land to take place in _sys_rtt.
1832  */
1833 static int
1834 instr_is_sys_rtt(caddr_t pc)
1835 {
1836     extern void _sys_rtt(), _sys_rtt_end();

1838 #if !defined(__xpv)
1839     extern void tr_sysc_ret_start(), tr_sysc_ret_end();
1840     extern void tr_intr_ret_start(), tr_intr_ret_end();

1842     if ((uintptr_t)pc >= (uintptr_t)tr_sysc_ret_start &&
1843         (uintptr_t)pc <= (uintptr_t)tr_sysc_ret_end)
1844         return (1);

1846     if ((uintptr_t)pc >= (uintptr_t)tr_intr_ret_start &&
1847         (uintptr_t)pc <= (uintptr_t)tr_intr_ret_end)
1848         return (1);
1849 #endif

1851     if ((uintptr_t)pc < (uintptr_t)_sys_rtt ||
1852         (uintptr_t)pc > (uintptr_t)_sys_rtt_end)
1853         return (0);

1855     return (1);
1856 }
    _____
    unchanged_portion_omitted_
```

new/usr/src/uts/i86pc/sys/mach_mmu.h

1

```
*****
5476 Fri Apr 6 17:25:03 2018
new/usr/src/uts/i86pc/sys/mach_mmu.h
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22 * Copyright 2007 Sun Microsystems, Inc. All rights reserved.
23 * Use is subject to license terms.
24 *
25 * Copyright 2018 Joyent, Inc.
26 */

28 #ifndef _SYS_MACH_MMU_H
29 #define _SYS_MACH_MMU_H

31 #ifdef __cplusplus
32 extern "C" {
33 #endif

35 #ifndef _ASM

37 #include <sys/types.h>
38 #include <sys/system.h>

40 /*
41  * Platform-dependent MMU routines and types.
42  *
43  * WARNING: this header file is used by both dboot and i86pc, so don't go using
44  * normal kernel headers.
45  */

47 #define TWO_MEG      (2 * 1024 * 1024)

49 /*
50  * This is:
51  *   The kernel nucleus pagesizes, ie: bi->bi_kseg_size
52  *   The grub 64 bit file load address (see multiboot header in dboot_grub.s)
53  *   The grub 32 bit and hypervisor physical load addresses of
54  *   the kernel text/data (see Mapfile.unix)
55  */
56 #define FOUR_MEG      (4 * 1024 * 1024)
```

new/usr/src/uts/i86pc/sys/mach_mmu.h

2

```
58 #define ONE_GIG      (1024 * 1024 * 1024)
59 #define FOUR_GIG     ((uint64_t)4 * ONE_GIG)

61 #define MMU_STD_PAGESIZE      4096
62 #ifdef __amd64
63 #define MMU_STD_PAGEMASK     0xFFFFFFFFFFFFFFFFULL
64 #else
65 #define MMU_STD_PAGEMASK     0xFFFFFFFF000ULL
66 #endif

68 /*
69  * Defines for the bits in X86 and AMD64 Page Tables
70  *
71  * Notes:
72  *
73  * Largepages and PAT bits:
74  *
75  * bit 7 at level 0 is the PAT bit
76  * bit 7 above level 0 is the Pagesize bit (set for large page)
77  * bit 12 (when a large page) is the PAT bit
78  *
79  * In Solaris the PAT/PWT/PCD values are set up so that:
80  *
81  * PAT & PWT -> Write Protected
82  * PAT & PCD -> Write Combining
83  * PAT by itself (PWT == 0 && PCD == 0) yields uncacheable (same as PCD == 1)
84  *
85  *
86  * Permission bits:
87  *
88  * - PT_USER must be set in all levels for user pages
89  * - PT_WRITE must be set in all levels for user writable pages
90  * - PT_NX applies if set at any level
91  *
92  * For these, we use the "allow" settings in all tables above level 0 and only
93  * ever disable things in PTEs.
94  *
95  * The use of PT_GLOBAL and PT_NX depend on being enabled in processor
96  * control registers. Hence, we use a variable to reference these bit
97  * masks. During hat_kern_setup() if the feature isn't enabled we
98  * clear out the variables.
99  */

100 #define PT_VALID      (0x001) /* a valid translation is present */
101 #define PT_WRITABLE  (0x002) /* the page is writable */
102 #define PT_USER      (0x004) /* the page is accessible by user mode */
103 #define PT_WRITETHRU (0x008) /* write back caching is disabled (non-PAT) */
104 #define PT_NOCACHE   (0x010) /* page is not cacheable (non-PAT) */
105 #define PT_REF       (0x020) /* page was referenced */
106 #define PT_MOD       (0x040) /* page was modified */
107 #define PT_PAGESIZE  (0x080) /* above level 0, indicates a large page */
108 #define PT_PAT_4K    (0x080) /* at level 0, used for write combining */
109 #define PT_GLOBAL    (0x100) /* the mapping is global */
110 #define PT_SOFTWARE  (0xe00) /* software bits */

112 #define PT_PAT_LARGE  (0x1000) /* PAT bit for large pages */

114 #define PT_PTPBITS    (PT_VALID | PT_USER | PT_WRITABLE | PT_REF)
115 #define PT_FLAGBITS   (0xffff) /* for masking off flag bits */

117 /*
118  * The software bits are used by the HAT to track attributes.
119  * Note that the attributes are inclusive as the values increase.
120  *
121  * PT_NOSYNC - The PT_REF/PT_MOD bits are not sync'd to page_t.
122  * The hat will install them as always set.
123  */
```

```
124 * PT_NOCONSIST - There is no hment entry for this mapping.
125 *
126 * PT_FOREIGN - used for the hypervisor, check via
127 *             (pte & PT_SOFTWARE) >= PT_FOREIGN
128 *             as it might set 0x800 for foreign grant table mappings.
129 */
130 #define PT_NOSYNC      (0x200) /* PTE was created with HAT_NOSYNC */
131 #define PT_NOCONSIST  (0x400) /* PTE was created with HAT_LOAD_NOCONSIST */
132 #define PT_FOREIGN    (0x600) /* MFN mapped on the hypervisor has no PFN */

134 #ifndef _BOOT
136 extern ulong_t getcr3(void);
137 extern void setcr3(ulong_t);

139 #define getcr3_pa() (getcr3() & MMU_PAGEMASK)
140 #define getpcid() ((getcr4() & CR4_PCIDE) ? \
141                  (getcr3() & MMU_PAGEOFFSET) : PCID_NONE)

143 extern void mmu_invlpg(caddr_t);

145 #endif

147 #ifdef __xpv
148 #include <sys/xen_mmu.h>
149 #else
150 #include <sys/pc_mmu.h>
151 #endif

153 /*
154 * The software extraction for a single Page Table Entry will always
155 * be a 64 bit unsigned int. If running a non-PAE hat, the page table
156 * access routines know to extend/shorten it to 32 bits.
157 */
158 typedef uint64_t x86pte_t;
159 typedef uint32_t x86pte32_t;

161 x86pte_t get_pteval(paddr_t, uint_t);
162 void set_pteval(paddr_t, uint_t, uint_t, x86pte_t);
163 paddr_t make_ptable(x86pte_t *, uint_t);
164 x86pte_t *find_pte(uint64_t, paddr_t *, uint_t, uint_t);
165 x86pte_t *map_pte(paddr_t, uint_t);

154 #ifndef _BOOT
155 ulong_t getcr3();
156 #endif

167 extern uint_t *shift_amt;
168 extern uint_t ptes_per_table;
169 extern paddr_t top_page_table;
170 extern uint_t top_level;
171 extern uint_t pte_size;
172 extern uint_t shift_amt_nopae[];
173 extern uint_t shift_amt_pae[];
174 extern uint32_t lpagesize;

176 #ifdef __cplusplus
177 }
_____ unchanged_portion_omitted
```

```

*****
7583 Fri Apr 6 17:25:03 2018
new/usr/src/uts/i86pc/sys/machcpuvar.h
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22 * Copyright 2009 Sun Microsystems, Inc. All rights reserved.
23 * Use is subject to license terms.
24 */
25 /*
26 * Copyright 2018 Joyent, Inc.
27 * Copyright 2011 Joyent, Inc. All rights reserved.
28 */

29 #ifndef _SYS_MACHCPUVAR_H
30 #define _SYS_MACHCPUVAR_H

31
32 #ifdef __cplusplus
33 extern "C" {
34 #endif

35
36 #include <sys/inttypes.h>
37 #include <sys/x_call.h>
38 #include <sys/tss.h>
39 #include <sys/segments.h>
40 #include <sys/rm_platter.h>
41 #include <sys/avintr.h>
42 #include <sys/pte.h>
43 #include <sys/stddef.h>
44 #include <sys/debug.h>
45 #include <sys/cpuvar.h>

46
47 #ifndef _ASM
48 /*
49 * On a virtualized platform a virtual cpu may not be actually
50 * on a physical cpu, especially in situations where a configuration has
51 * more vcpus than pcpus. This function tells us (if it's able) if the
52 * specified vcpu is currently running on a pcpu. Note if it is not
53 * known or not able to determine, it will return the unknown state.
54 */
55 #define VCPU_STATE_UNKNOWN 0
56 #define VCPU_ON_PCPU 1
57 #define VCPU_NOT_ON_PCPU 2

```

```

59 extern int vcpu_on_pcpu(processorid_t);

61 /*
62 * Machine specific fields of the cpu struct
63 * defined in common/sys/cpuvar.h.
64 *
65 * Note: This is kinda kludgy but seems to be the best
66 * of our alternatives.
67 */

69 struct cpuid_info;
70 struct cpu_ucode_info;
71 struct cmi_hdl;

73 /*
74 * A note about the hypervisor affinity bits: a one bit in the affinity mask
75 * means the corresponding event channel is allowed to be serviced
76 * by this cpu.
77 */
78 struct xen_evt_data {
79     ulong_t    pending_sel[PIL_MAX + 1]; /* event array selectors */
80     ulong_t    pending_evts[PIL_MAX + 1][sizeof(ulong_t) * 8];
81     ulong_t    evt_affinity[sizeof(ulong_t) * 8]; /* service on cpu */
82 };

84 struct kpti_frame {
85     uint64_t    kf_lower_redzone;

87     /* Stashed value of %cr3 when we entered the trampoline. */
88     greg_t      kf_tr_cr3;

90     /*
91     * We use %r13-r14 as scratch registers in the trampoline code,
92     * so stash those here "below" the rest of the stack so they can be
93     * pushed/popped if needed.
94     */
95     greg_t      kf_r14;
96     greg_t      kf_r13;

98     /*
99     * Part of this struct is used as the HW stack frame when taking an
100    * interrupt on the user page table. The CPU is going to push a bunch
101    * of regs onto the stack pointer set in the TSS/IDT (which we set to
102    * &kf_rsp here).
103    *
104    * This is only a temporary holding area for them (we'll move them over
105    * to the real interrupt stack once we've set %cr3).
106    *
107    * Note that these must be cleared during a process switch on this cpu.
108    */
109     greg_t      kf_err; /* Bottom of initial hw stack frame */
110     greg_t      kf_rip;
111     greg_t      kf_cs;
112     greg_t      kf_rflags;
113     greg_t      kf_rsp;
114     greg_t      kf_ss;

116     greg_t      kf_tr_rsp; /* Top of HW stack frame */
117     /* We also write this with the %rsp value on tramp entry */

119     /* Written to 0x1 when this kpti_frame is in use. */
120     uint64_t    kf_tr_flag;

122     uint64_t    kf_middle_redzone;

124     /*

```



```

125  * The things we need to write to %cr3 to change between page tables.
126  * These live "above" the HW stack.
127  */
128  greg_t      kf_kernel_cr3;
129  greg_t      kf_user_cr3;
130  greg_t      kf_tr_ret_rsp;

132  uint64_t    kf_unused;          /* For 16-byte align */

134  uint64_t    kf_upper_redzone;
135  };

137  /*
138  * This first value, MACHCPU_SIZE is the size of all the members in the cpu_t
139  * AND struct machcpu, before we get to the mcpu_pad and the kpti area.
140  * The KPTI is used to contain per-CPU data that is visible in both sets of
141  * page-tables, and hence must be page-aligned and page-sized. See
142  * hat_pcp_setup().
143  *
144  * There is a CTASSERT in os/intr.c that checks these numbers.
145  */
146  #define MACHCPU_SIZE      (572 + 1584)
147  #define MACHCPU_PAD      (MMU_PAGESIZE - MACHCPU_SIZE)
148  #define MACHCPU_PAD2     (MMU_PAGESIZE - 16 - 3 * sizeof (struct kpti_frame))

150  struct machcpu {
151  /*
152  * x_call fields - used for interprocessor cross calls
153  */
154  struct xc_msg *xc_msgbox;
155  struct xc_msg *xc_free;
156  xc_data_t xc_data;
157  uint32_t xc_wait_cnt;
158  volatile uint32_t xc_work_cnt;

160  int mcpu_nodeid;          /* node-id */
161  int mcpu_pri;            /* CPU priority */

163  struct hat *mcpu_current_hat; /* cpu's current hat */

165  struct hat_cpu_info *mcpu_hat_info;

167  volatile ulong_t mcpu_tlb_info;

169  /* i86 hardware table addresses that cannot be shared */

171  user_desc_t *mcpu_gdt;   /* GDT */
172  gate_desc_t *mcpu_idt;  /* current IDT */

174  tss_t *mcpu_tss;        /* TSS */
175  void *mcpu_ldt;
176  size_t mcpu_ldt_len;

178  kmutex_t mcpu_ppaddr_mutex;
179  caddr_t mcpu_caddr1;    /* per cpu CADDR1 */
180  caddr_t mcpu_caddr2;    /* per cpu CADDR2 */
181  uint64_t mcpu_caddr1pte;
182  uint64_t mcpu_caddr2pte;

184  struct softint mcpu_softinfo;
185  uint64_t pil_high_start[HIGH_LEVELS];
186  uint64_t intrstat[PIL_MAX + 1][2];

188  struct cpuid_info *mcpu_cpi;

190  #if defined(__amd64)

```

```

191  greg_t mcpu_rtmp_rsp;    /* syscall: temporary %rsp stash */
192  greg_t mcpu_rtmp_r15;   /* syscall: temporary %r15 stash */
193  #endif

195  struct vcpu_info *mcpu_vcpu_info;
196  uint64_t mcpu_gdtpa;    /* hypervisor: GDT physical address */

198  uint16_t mcpu_intr_pending; /* hypervisor: pending intrpt levels */
199  uint16_t mcpu_ec_mbox;     /* hypervisor: evtchn_dev mailbox */
200  struct xen_evt_data *mcpu_evt_pending; /* hypervisor: pending events */

202  volatile uint32_t *mcpu_mwait; /* MONITOR/MWAIT buffer */
203  void (*mcpu_idle_cpu)(void); /* idle function */
204  uint16_t mcpu_idle_type;     /* CPU next idle type */
205  uint16_t max_cstates;       /* supported max cstates */

207  struct cpu_ucose_info *mcpu_ucose_info;

209  void *mcpu_pm_mach_state;
210  struct cmi_hdl *mcpu_cmi_hdl;
211  void *mcpu_mach_ctx_ptr;

213  /*
214  * A stamp that is unique per processor and changes
215  * whenever an interrupt happens. Useful for detecting
216  * if a section of code gets interrupted.
217  * The high order 16 bits will hold the cpu->cpu_id.
218  * The low order bits will be incremented on every interrupt.
219  */
220  volatile uint32_t mcpu_istamp;

222  char mcpu_pad[MACHCPU_PAD];

224  /* This is the start of the page */
225  char mcpu_pad2[MACHCPU_PAD2];
226  struct kpti_frame mcpu_kpti;
227  struct kpti_frame mcpu_kpti_flt;
228  struct kpti_frame mcpu_kpti_dbg;
229  char mcpu_pad3[16];
230  };

232  #define NINTR_THREADS (LOCK_LEVEL-1) /* number of interrupt threads */
233  #define MWAIT_HALTED (1) /* mcpu_mwait set when halting */
234  #define MWAIT_RUNNING (0) /* mcpu_mwait set to wakeup */
235  #define MWAIT_WAKEUP_IPI (2) /* need IPI to wakeup */
236  #define MWAIT_WAKEUP(cpu) (*(cpu->cpu_m.mcpu_mwait) = MWAIT_RUNNING)

238  #endif /* _ASM */

240  /* Please DON'T add any more of this namespace-poisoning sewage here */

242  #define cpu_nodeid cpu_m.mcpu_nodeid
243  #define cpu_pri cpu_m.mcpu_pri
244  #define cpu_current_hat cpu_m.mcpu_current_hat
245  #define cpu_hat_info cpu_m.mcpu_hat_info
246  #define cpu_ppaddr_mutex cpu_m.mcpu_ppaddr_mutex
247  #define cpu_gdt cpu_m.mcpu_gdt
248  #define cpu_idt cpu_m.mcpu_idt
249  #define cpu_tss cpu_m.mcpu_tss
250  #define cpu_ldt cpu_m.mcpu_ldt
251  #define cpu_caddr1 cpu_m.mcpu_caddr1
252  #define cpu_caddr2 cpu_m.mcpu_caddr2
253  #define cpu_softinfo cpu_m.mcpu_softinfo
254  #define cpu_caddr1pte cpu_m.mcpu_caddr1pte
255  #define cpu_caddr2pte cpu_m.mcpu_caddr2pte

```

new/usr/src/uts/i86pc/sys/machcpuvar.h

5

```
256 #ifdef __cplusplus
257 }
```

_____unchanged_portion_omitted_

new/usr/src/uts/i86pc/sys/machparam.h

1

```
*****
9789 Fri Apr 6 17:25:04 2018
new/usr/src/uts/i86pc/sys/machparam.h
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22 * Copyright (c) 1992, 2010, Oracle and/or its affiliates. All rights reserved.
23 * Copyright (c) 2015 by Delphix. All rights reserved.
24 * Copyright 2018 Joyent, Inc.
25 * Copyright 2016 Joyent, Inc.
26 */
27 /*      Copyright (c) 1988 AT&T */
28 /*      All Rights Reserved      */

31 #ifndef _SYS_MACHPARAM_H
32 #define _SYS_MACHPARAM_H

34 #if !defined(_ASM)
35 #include <sys/types.h>

37 #if defined(__xpv)
38 #include <sys/xpv_impl.h>
39 #endif

41 #endif

43 #ifdef __cplusplus
44 extern "C" {
45 #endif

47 #ifndef _ASM
48 #define ADDRESS_C(c)    c ## ul
49 #else /* _ASM */
50 #define ADDRESS_C(c)    (c)
51 #endif /* _ASM */

53 /*
54  * Machine dependent parameters and limits.
55  */

57 #if defined(__amd64)
58 /*
```

new/usr/src/uts/i86pc/sys/machparam.h

2

```
59  * If NCPU grows beyond 256, sizing for the x86 comm page will require
60  * adjustment.
61  */
62 #define NCPU      256
63 #define NCPU_LOG2      8
64 #elif defined(__i386)
65 #define NCPU      32
66 #define NCPU_LOG2      5
67 #endif

69 /* NCPU_P2 is NCPU rounded to a power of 2 */
70 #define NCPU_P2 (1 << NCPU_LOG2)

72 /*
73  * The value defined below could grow to 16. hat structure and
74  * page_t have room for 16 nodes.
75  */
76 #define MAXNODES      4
77 #define NUMA_NODEMASK      0x0f

79 /*
80  * Define the FPU symbol if we could run on a machine with an external
81  * FPU (i.e. not integrated with the normal machine state like the vax).
82  *
83  * The fpu is defined in the architecture manual, and the kernel hides
84  * its absence if it is not present, that's pretty integrated, no?
85  */

87 /* supported page sizes */
88 #define MMU_PAGE_SIZES      3

90 /*
91  * MMU_PAGES* describes the physical page size used by the mapping hardware.
92  * PAGES* describes the logical page size used by the system.
93  */

95 #define MMU_PAGESIZE      0x1000          /* 4096 bytes */
96 #define MMU_PAGESHIFT      12           /* log2(MMU_PAGESIZE) */

98 #if !defined(_ASM)
99 #define MMU_PAGEOFFSET      (MMU_PAGESIZE-1) /* Mask of address bits in page */
100 #else /* _ASM */
101 #define MMU_PAGEOFFSET      _CONST(MMU_PAGESIZE-1) /* assembler lameness */
102 #endif /* _ASM */

104 #define MMU_PAGEMASK      (~MMU_PAGEOFFSET)

106 #define PAGESIZE      0x1000          /* All of the above, for logical */
107 #define PAGESHIFT      12
108 #define PAGEOFFSET      (PAGESIZE - 1)
109 #define PAGEMASK      (~PAGEOFFSET)

111 /*
112  * DATA_ALIGN is used to define the alignment of the Unix data segment.
113  */
114 #define DATA_ALIGN      PAGESIZE

116 /*
117  * DEFAULT_KERNEL_THREAD stack size (in pages).
118  */
119 #if defined(__amd64)
120 #define DEFAULTSTKSZ_NPGS      5
121 #elif defined(__i386)
122 #define DEFAULTSTKSZ_NPGS      3
123 #endif
```

```

125 #if !defined(_ASM)
126 #define DEFAULTSTKSZ (DEFAULTSTKSZ_NPGS * PAGESIZE)
127 #else /* !_ASM */
128 #define DEFAULTSTKSZ _MUL(DEFAULTSTKSZ_NPGS, PAGESIZE) /* as(1) lameness */
129 #endif /* !_ASM */

131 /*
132 * KERNELBASE is the virtual address at which the kernel segments start in
133 * all contexts.
134 *
135 * KERNELBASE is not fixed. The value of KERNELBASE can change with
136 * installed memory or on 32 bit systems the eprom variable 'eprom_kernelbase'.
137 *
138 * common/conf/param.c requires a compile time defined value for KERNELBASE.
139 * This value is save in the variable _kernelbase. _kernelbase may then be
140 * modified with to a different value in i86pc/os/startup.c.
141 *
142 * Most code should be using kernelbase, which resolves to a reference to
143 * _kernelbase.
144 */
145 #define KERNEL_TEXT_amd64 UINT64_C(0xfffffffffb800000)

147 #ifdef __i386

149 #define KERNEL_TEXT_i386 ADDRESS_C(0xfe800000)

151 /*
152 * We don't use HYPERVISOR_VIRT_START, as we need both the PAE and non-PAE
153 * versions in our code. We always compile based on the lower PAE address.
154 */
155 #define KERNEL_TEXT_i386_xpv \
156 (HYPERVISOR_VIRT_START_PAE - 3 * ADDRESS_C(0x400000))

158 #endif /* __i386 */

160 #if defined(__amd64)

162 #define KERNELBASE ADDRESS_C(0xfffffd8000000000)

164 /*
165 * Size of the unmapped "red zone" at the very bottom of the kernel's
166 * address space. Corresponds to 1 slot in the toplevel pagetable.
167 */
168 #define KERNEL_REDZONE_SIZE ((uintptr_t)1 << 39)

170 /*
171 * Base of 'core' heap area, which is used for kernel and module text/data
172 * that must be within a 2GB range to allow for rip-relative addressing.
173 */
174 #define COREHEAP_BASE ADDRESS_C(0xffffffffc0000000)

176 /*
177 * Beginning of the segkpm window. A lower value than this is used if
178 * physical addresses exceed 1TB. See i86pc/os/startup.c
179 */
180 #define SEGKPM_BASE ADDRESS_C(0xfffffe0000000000)

182 /*
183 * This is valloc_base, above seg_kpm, but below everything else.
184 * A lower value than this may be used if SEGKPM_BASE is adjusted.
185 * See i86pc/os/startup.c
186 */
187 #define VALLOC_BASE ADDRESS_C(0xfffff00000000000)

189 /*
190 * default and boundary sizes for segkp

```

```

191 */
192 #define SEGKPDEFESIZE (2L * 1024L * 1024L * 1024L) /* 2G */
193 #define SEGKPMAXSIZE (8L * 1024L * 1024L * 1024L) /* 8G */
194 #define SEGKPMINSIZE (200L * 1024 * 1024L) /* 200M */

196 /*
197 * minimum size for segzio
198 */
199 #define SEGZIOMINSIZE (400L * 1024 * 1024L) /* 400M */

201 /*
202 * During intial boot we limit heap to the top 4Gig.
203 */
204 #define BOOT_KERNELHEAP_BASE ADDRESS_C(0xffffffff00000000)

206 /*
207 * VMWare works best if we don't use the top 64Meg of memory for amd64.
208 * Set KERNEL_TEXT to top_o_memory - 64Meg - 8 Meg for 8Meg of nucleus pages.
209 */
210 #define PROMSTART ADDRESS_C(0xffc00000)
211 #define KERNEL_TEXT KERNEL_TEXT_amd64

213 /*
214 * Virtual address range available to the debugger
215 */
216 #define SEGDEBUGBASE ADDRESS_C(0xfffffffff8000000)
217 #define SEGDEBUGSIZE ADDRESS_C(0x4000000)

219 /*
220 * Define upper limit on user address space
221 */
222 * In amd64, the upper limit on a 64-bit user address space is 1 large page
223 * (2MB) below kernelbase. The upper limit for a 32-bit user address space
224 * is 1 small page (4KB) below the top of the 32-bit range. The 64-bit
225 * limit give dtrace the red zone it needs below kernelbase. The 32-bit
226 * limit gives us a small red zone to detect address-space overruns in a
227 * user program.
228 *
229 * On the hypervisor, we limit the user to memory below the VA hole.
230 * Subtract 1 large page for a red zone.
231 */
232 #if defined(__xpv)
233 #define USERLIMIT ADDRESS_C(0x00007ffffe000000)
234 #else
235 #define USERLIMIT ADDRESS_C(0xfffffd7fffe00000)
236 #endif

238 #ifdef bug_5074717_is_fixed
239 #define USERLIMIT32 ADDRESS_C(0xfffff000)
240 #else
241 #define USERLIMIT32 ADDRESS_C(0xfefff000)
242 #endif

244 #elif defined(__i386)

246 #ifdef DEBUG
247 #define KERNELBASE ADDRESS_C(0xc8000000)
248 #else
249 #define KERNELBASE ADDRESS_C(0xd4000000)
250 #endif

252 #define KERNELBASE_MAX ADDRESS_C(0xe0000000)

254 /*
255 * The i386 ABI requires that the user address space be at least 3Gb
256 * in size. KERNELBASE_ABI_MIN is used as the default KERNELBASE for

```

```

257 * physical memory configurations > 4gb.
258 */
259 #define KERNELBASE_ABI_MIN      ADDRESS_C(0xc0000000)

261 /*
262 * Size of the unmapped "red zone" at the very bottom of the kernel's
263 * address space. Since segmap start immediately above the red zone, this
264 * needs to be MAXBSIZE aligned.
265 */
266 #define KERNEL_REDZONE_SIZE     MAXBSIZE

268 /*
269 * This is the last 4MB of the 4G address space. Some psm modules
270 * need this region of virtual address space mapped 1-1
271 * The top 64MB of the address space is reserved for the hypervisor.
272 */
273 #define PROMSTART               ADDRESS_C(0xf0000000)
274 #ifndef __xpv
275 #define KERNEL_TEXT              KERNEL_TEXT_i386_xpv
276 #else
277 #define KERNEL_TEXT              KERNEL_TEXT_i386
278 #endif

280 /*
281 * Virtual address range available to the debugger
282 * We place it just above the kernel text (4M) and kernel data (4M).
283 */
284 #define SEGDEBUGBASE            (KERNEL_TEXT + ADDRESS_C(0x800000))
285 #define SEGDEBUGSIZE           ADDRESS_C(0x400000)

287 /*
288 * Define upper limit on user address space
289 */
290 #define USERLIMIT              KERNELBASE
291 #define USERLIMIT32            USERLIMIT

293 #endif /* __i386 */

295 /*
296 * Reserve pages just below KERNEL_TEXT for the GDT, IDT, LDT, TSS and debug
297 * info.
298 * Reserve pages just below KERNEL_TEXT for the GDT, IDT, TSS and debug info.
299 * For now, DEBUG_INFO_VA must be first in this list for "xm" initiated dumps
300 * of solaris domUs to be usable with mdb. Relying on a fixed VA is not viable
301 * long term, but it's the best we've got for now.
302 */
303 #if !defined(_ASM)
304 #define DEBUG_INFO_VA          (KERNEL_TEXT - MMU_PAGESIZE)
305 #define GDT_VA                 (DEBUG_INFO_VA - MMU_PAGESIZE)
306 #define IDT_VA                 (GDT_VA - MMU_PAGESIZE)
307 #define LDT_VA                 (IDT_VA - (16 * MMU_PAGESIZE))
308 #define KTSS_VA               (LDT_VA - MMU_PAGESIZE)
309 #define DFTSS_VA              (IDT_VA - MMU_PAGESIZE)
310 #define MISC_VA_BASE          (DFTSS_VA)
311 #define MISC_VA_SIZE          (KERNEL_TEXT - MISC_VA_BASE)
312 #endif /* !_ASM */

314 #if !defined(_ASM) && !defined(_KMDB)
315 extern uintptr_t kernelbase, segmap_start, segmapsize;
316 #endif

318 /*
319 * ARGSBASE is the base virtual address of the range which
320 * the kernel uses to map the arguments for exec.

```

```

321 */
322 #define ARGSBASE                PROMSTART

324 /*
325 * reserve space for modules
326 */
327 #define MODTEXT (1024 * 1024 * 2)
328 #define MODDATA (1024 * 300)

330 /*
331 * The heap has a region allocated from it of HEAPTEXT_SIZE bytes specifically
332 * for module text.
333 */
334 #define HEAPTEXT_SIZE          (64 * 1024 * 1024) /* bytes */

336 /*
337 * Size of a kernel threads stack. It must be a whole number of pages
338 * since the segment it comes from will only allocate space in pages.
339 */
340 #define T_STACKSZ              2*PAGESIZE

342 /*
343 * Size of a cpu startup thread stack. (It must be a whole number of pages
344 * since the containing segment only allocates space in pages.)
345 */

347 #define STARTUP_STKSZ          3*PAGESIZE

349 /*
350 * Bus types
351 */
352 #define BTISA                   1
353 #define BTEISA                  2
354 #define BTMCA                   3

356 #ifdef __cplusplus
357 }

```

unchanged_portion_omitted

new/usr/src/uts/i86pc/sys/machprivregs.h

1

```
*****
3763 Fri Apr 6 17:25:04 2018
new/usr/src/uts/i86pc/sys/machprivregs.h
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */

22 /*
23  * Copyright 2007 Sun Microsystems, Inc. All rights reserved.
24  * Use is subject to license terms.
25  *
26  * Copyright 2018 Joyent, Inc.
27  */

29 #ifndef _SYS_MACHPRIVREGS_H
30 #define _SYS_MACHPRIVREGS_H

32 /*
33  * Platform dependent instruction sequences for manipulating
34  * privileged state
35  */

37 #ifdef __cplusplus
38 extern "C" {
39 #endif

41 #define ASSERT_UPCALL_MASK_IS_SET          /* empty */

43 /*
44  * CLI and STI
45  */

47 #define CLI(r)                             \
48     cli

50 #define STI                                 \
51     sti

53 /*
54  * Used to re-enable interrupts in the body of exception handlers
55  */

57 #if defined(__amd64)
```

new/usr/src/uts/i86pc/sys/machprivregs.h

2

```
59 #define ENABLE_INTR_FLAGS                 \
60     pushq  $F_ON;                         \
61     popfq

63 #elif defined(__i386)

65 #define ENABLE_INTR_FLAGS                 \
66     pushl  $F_ON;                         \
67     popfl

69 #endif /* __i386 */

71 /*
72  * IRET and SWAPGS
73  */
74 #if defined(__amd64)

76 #define IRET        iretq
77 #define SYSRETQ    sysretq
78 #define SYSRETL    sysretl
79 #define SWAPGS      swapgs
80 #define XPV_TRAP_POP /* empty */
81 #define XPV_TRAP_PUSH /* empty */

83 #elif defined(__i386)

85 #endif /* __i386 */

87 #define XPV_TRAP_POP /* empty */
88 #define XPV_TRAP_PUSH /* empty */
89 #define CLEAN_CS    /* empty */

92 /*
93  * Macros for saving the original segment registers and restoring them
94  * for fast traps.
95  */
96 #if defined(__amd64)

98 /*
99  * Smaller versions of INTR_PUSH and INTR_POP for fast traps.
100 * The following registers have been pushed onto the stack by
101 * hardware at this point:
102 *
103 *     greg_t  r_rip;
104 *     greg_t  r_cs;
105 *     greg_t  r_rfl;
106 *     greg_t  r_rsp;
107 *     greg_t  r_ss;
108 *
109 * This handler is executed both by 32-bit and 64-bit applications.
110 * 64-bit applications allow us to treat the set (%rdi, %rsi, %rdx,
111 * %rcx, %r8, %r9, %r10, %r11, %rax) as volatile across function calls.
112 * However, 32-bit applications only expect (%eax, %edx, %ecx) to be volatile
113 * across a function call -- in particular, %esi and %edi MUST be saved!
114 *
115 * We could do this differently by making a FAST_INTR_PUSH32 for 32-bit
116 * programs, and FAST_INTR_PUSH for 64-bit programs, but it doesn't seem
117 * particularly worth it.
118  */
119 #define FAST_INTR_PUSH                     \
120     INTGATE_INIT_KERNEL_FLAGS;           \
121     subq  $REGOFF_RIP, %rsp;
```

```
122     movq    %rsi, REGOFF_RSI(%rsp); \  
123     movq    %rdi, REGOFF_RDI(%rsp); \  
124     swapgs  
  
126 #define FAST_INTR_POP          \  
127     swapgs; \  
128     movq    REGOFF_RSI(%rsp), %rsi; \  
129     movq    REGOFF_RDI(%rsp), %rdi; \  
130     addq    $REGOFF_RIP, %rsp  
  
132 #define FAST_INTR_RETURN      jmp tr_iret_user  
132 #define FAST_INTR_RETURN      iretq  
  
134 #elif defined(__i386)  
  
136 #define FAST_INTR_PUSH        \  
137     cld; \  
138     __SEGREGS_PUSH           \  
139     __SEGREGS_LOAD_KERNEL  
  
141 #define FAST_INTR_POP          \  
142     __SEGREGS_POP  
  
144 #define FAST_INTR_RETURN      iret  
  
146 #endif /* __i386 */  
  
148 /*  
149  * Handling the CR0.TS bit for floating point handling.  
150  *  
151  * When the TS bit is *set*, attempts to touch the floating  
152  * point hardware will result in a #nm trap.  
153  */  
154 #if defined(__amd64)  
  
156 #define STTS(rtmp)             \  
157     movq    %cr0, rtmp; \  
158     orq    $CR0_TS, rtmp; \  
159     movq    rtmp, %cr0  
  
161 #elif defined(__i386)  
  
163 #define STTS(rtmp)             \  
164     movl   %cr0, rtmp; \  
165     orl   $CR0_TS, rtmp; \  
166     movl   rtmp, %cr0  
  
168 #endif /* __i386 */  
  
170 #define CLTS                    \  
171     clts  
  
173 #ifdef __cplusplus  
174 }  
  
unchanged_portion_omitted
```

new/usr/src/uts/i86pc/sys/pc_mmu.h

1

```
*****
1833 Fri Apr 6 17:25:04 2018
new/usr/src/uts/i86pc/sys/pc_mmu.h
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22 * Copyright 2007 Sun Microsystems, Inc. All rights reserved.
23 * Use is subject to license terms.
24 *
25 * Copyright 2018 Joyent, Inc.
26 */

28 #ifndef _SYS_PC_MMU_H
29 #define _SYS_PC_MMU_H

31 #ifdef __cplusplus
32 extern "C" {
33 #endif

35 /*
36  * Platform-dependent MMU routines and types for real x86 hardware.
37  *
38  * WARNING: this header file is used by both dboot and i86pc, so don't go using
39  * normal kernel headers.
40  */

42 #define IN_HYPERVISOR_VA(va) (__lintzero)

44 void reload_cr3(void);

46 #define pa_to_ma(pa) (pa)
47 #define ma_to_pa(ma) (ma)
48 #define pfn_to_mfn(pfn) (pfn)
49 #define mfn_to_pfn(mfn) (mfn)

51 #ifndef _BOOT

53 extern uint64_t kpti_safe_cr3;
53 void mmu_tlbflush_entry(caddr_t);
54 void setcr3(ulong_t);

55 #define INVPCID_ADDR (0)
```

new/usr/src/uts/i86pc/sys/pc_mmu.h

2

```
56 #define INVPCID_ID (1)
57 #define INVPCID_ALL_GLOBAL (2)
58 #define INVPCID_ALL_NONGLOBAL (3)

60 extern void invpcid_insn(uint64_t, uint64_t, uintptr_t);
61 extern void tr_mmu_flush_user_range(uint64_t, size_t, size_t, uint64_t);

63 #if defined(__GNUC__)
64 #include <asm/mmu.h>
65 #endif

67 #endif /* !_BOOT */

69 #ifdef __cplusplus
70 }
_____unchanged_portion_omitted_
```



```

*****
4107 Fri Apr 6 17:25:04 2018
new/usr/src/uts/i86pc/sys/rm_platter.h
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22 * Copyright (c) 1992, 2010, Oracle and/or its affiliates. All rights reserved.
23 */
24 /*
25 * Copyright (c) 2010, Intel Corporation.
26 * All rights reserved.
27 */
28 /*
29 * Copyright 2018 Joyent, Inc.
30 * Copyright 2011 Joyent, Inc. All rights reserved.
31 */

32 #ifndef _SYS_RM_PLATTER_H
33 #define _SYS_RM_PLATTER_H

34
35 #include <sys/types.h>
36 #include <sys/tss.h>
37 #include <sys/segments.h>

38
39 #ifdef __cplusplus
40 extern "C" {
41 #endif

42
43 #define RM_PLATTER_CODE_SIZE      0x400
44 #define RM_PLATTER_CPU_HALT_CODE_SIZE  0x100

45
46 typedef struct rm_platter {
47     char    rm_code[RM_PLATTER_CODE_SIZE];
48     char    rm_cpu_halt_code[RM_PLATTER_CPU_HALT_CODE_SIZE];
49 #if defined(__amd64)
50     /*
51      * The compiler will want to 64-bit align the 64-bit rm_gdt_base
52      * pointer, so we need to add an extra four bytes of padding here to
53      * make sure rm_gdt_lim and rm_gdt_base will align to create a proper
54      * ten byte GDT pseudo-descriptor.
55      */
56     uint32_t    rm_gdt_pad;
57 #endif /* __amd64 */
58     ushort_t    rm_debug;

```

```

59     ushort_t    rm_gdt_lim;    /* stuff for lgdt */
60     user_desc_t    *rm_gdt_base;
61 #if defined(__amd64)
62     /*
63      * The compiler will want to 64-bit align the 64-bit rm_idt_base
64      * pointer, so we need to add an extra four bytes of padding here to
65      * make sure rm_idt_lim and rm_idt_base will align to create a proper
66      * ten byte IDT pseudo-descriptor.
67      */
68     uint32_t    rm_idt_pad;
69 #endif /* __amd64 */
70     ushort_t    rm_cpu_halted; /* non-zero if CPU has been halted */
71     ushort_t    rm_idt_lim;    /* stuff for lidt */
72     gate_desc_t    *rm_idt_base;
73     uint_t    rm_pdbrr;    /* cr3 value */
74     uint_t    rm_cpu;    /* easy way to know which CPU we are */
75     uint_t    rm_filler3;
76     uint_t    rm_cr4;    /* cr4 value on cpu0 */
77 #if defined(__amd64)
78     /*
79      * Temporary GDT for the brief transition from real mode to protected
80      * mode before a CPU continues on into long mode.
81      *
82      * Putting it here assures it will be located in identity mapped memory
83      * (va == pa, 1:1).
84      *
85      * rm_temp_gdt is sized to hold only a null descriptor in slot zero
86      * and a 64-bit code descriptor in slot one.
87      *
88      * rm_temp_[g]idt_lim and rm_temp_[g]idt_base are the pseudo-descriptors
89      * for the temporary GDT and IDT, respectively.
90      */
91     uint64_t    rm_temp_gdt[2];
92     ushort_t    rm_temp_gdt_desc_pad; /* filler to align GDT desc */
93     ushort_t    rm_temp_gdt_lim;
94     uint32_t    rm_temp_gdt_base;
95     ushort_t    rm_temp_idt_desc_pad; /* filler to align IDT desc */
96     ushort_t    rm_temp_idt_lim;
97     uint32_t    rm_temp_idt_base;

98
99     /*
100    * The code executing in the rm_platter needs the offset into the
101    * platter at which the 64-bit code starts, so have mp_startup
102    * calculate it and store it here.
103    */
104     uint32_t    rm_longmode64_addr;
105 #endif /* __amd64 */
106 } rm_platter_t;

107
108 /*
109 * cpu tables put within a single structure two of the tables which need to be
110 * allocated when a CPU starts up.
111 *
112 * Note: the tss should be 16 byte aligned for best performance on amd64
113 * Since DEFAULTSTKSZ is a multiple of PAGESIZE tss will be aligned.
114 */
115 struct cpu_tables {
116     /* IST stacks */
117     char    ct_stack1[DEFAULTSTKSZ]; /* dblfault */
118     #if !defined(__xpv)
119     char    ct_stack2[DEFAULTSTKSZ]; /* nmi */
120     char    ct_stack3[DEFAULTSTKSZ]; /* mce */
121 #endif
122     char    ct_stack[DEFAULTSTKSZ];
123     tss_t    ct_tss;
124 };

```

unchanged_portion_omitted

```

*****
131966 Fri Apr 6 17:25:04 2018
new/usr/src/uts/i86pc/vm/hat_i86.c
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>
9208 hati_demap_func should take pagesize into account
Reviewed by: Robert Mustacchi <rm@joyent.com>
Reviewed by: Tim Kordas <tim.kordas@joyent.com>
Reviewed by: Yuri Pankov <yuripv@yuripv.net>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[ ]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22  * Copyright (c) 1992, 2010, Oracle and/or its affiliates. All rights reserved.
23 */
24 /*
25  * Copyright (c) 2010, Intel Corporation.
26  * All rights reserved.
27 */
28 /*
29  * Copyright 2011 Nexenta Systems, Inc. All rights reserved.
30  * Copyright 2018 Joyent, Inc. All rights reserved.
31  * Copyright (c) 2014, 2015 by Delphix. All rights reserved.
32 */
33
34 /*
35  * VM - Hardware Address Translation management for i386 and amd64
36  *
37  * Implementation of the interfaces described in <common/vm/hat.h>
38  *
39  * Nearly all the details of how the hardware is managed should not be
40  * visible outside this layer except for misc. machine specific functions
41  * that work in conjunction with this code.
42  *
43  * Routines used only inside of i86pc/vm start with hati_ for HAT Internal.
44  */
45
46 /*
47  * amd64 HAT Design
48  *
49  * -----
50  * Background
51  * -----
52  *
53  * On x86, the address space is shared between a user process and the kernel.
54  * This is different from SPARC. Conventionally, the kernel lives at the top of
55  * the address space and the user process gets to enjoy the rest of it. If you

```

```

56 * look at the image of the address map in uts/i86pc/os/startup.c, you'll get a
57 * rough sense of how the address space is laid out and used.
58 *
59 * Every unique address space is represented by an instance of a HAT structure
60 * called a 'hat_t'. In addition to a hat_t structure for each process, there is
61 * also one that is used for the kernel (kas.a_hat), and each CPU ultimately
62 * also has a HAT.
63 *
64 * Each HAT contains a pointer to its root page table. This root page table is
65 * what we call an L3 page table in illumos and Intel calls the PML4. It is the
66 * physical address of the L3 table that we place in the %cr3 register which the
67 * processor uses.
68 *
69 * Each of the many layers of the page table is represented by a structure
70 * called an htable_t. The htable_t manages a set of 512 8-byte entries. The
71 * number of entries in a given page table is constant across all different
72 * level page tables. Note, this is only true on amd64. This has not always been
73 * the case on x86.
74 *
75 * Each entry in a page table, generally referred to as a PTE, may refer to
76 * another page table or a memory location, depending on the level of the page
77 * table and the use of large pages. Importantly, the top-level L3 page table
78 * (PML4) only supports linking to further page tables. This is also true on
79 * systems which support a 5th level page table (which we do not currently
80 * support).
81 *
82 * Historically, on x86, when a process was running on CPU, the root of the page
83 * table was inserted into %cr3 on each CPU on which it was currently running.
84 * When processes would switch (by calling hat_switch()), then the value in %cr3
85 * on that CPU would change to that of the new HAT. While this behavior is still
86 * maintained in the xpv kernel, this is not what is done today.
87 *
88 * -----
89  * Per-CPU Page Tables
90  * -----
91 *
92 * Throughout the system the 64-bit kernel has a notion of what it calls a
93 * per-CPU page table or PCP. The notion of a per-CPU page table was originally
94 * introduced as part of the original work to support x86 PAE. On the 64-bit
95 * kernel, it was originally used for 32-bit processes running on the 64-bit
96 * kernel. The rationale behind this was that each 32-bit process could have all
97 * of its memory represented in a single L2 page table as each L2 page table
98 * entry represents 1 GbE of memory.
99 *
100 * Following on from this, the idea was that given that all of the L3 page table
101 * entries for 32-bit processes are basically going to be identical with the
102 * exception of the first entry in the page table, why not share those page
103 * table entries. This gave rise to the idea of a per-CPU page table.
104 *
105 * The way this works is that we have a member in the machcpu_t called the
106 * mcpu_hat_info. That structure contains two different 4k pages: one that
107 * represents the L3 page table and one that represents an L2 page table. When
108 * the CPU starts up, the L3 page table entries are copied in from the kernel's
109 * page table. The L3 kernel entries do not change throughout the lifetime of
110 * the kernel. The kernel portion of these L3 pages for each CPU have the same
111 * records, meaning that they point to the same L2 page tables and thus see a
112 * consistent view of the world.
113 *
114 * When a 32-bit process is loaded into this world, we copy the 32-bit process's
115 * four top-level page table entries into the CPU's L2 page table and then set
116 * the CPU's first L3 page table entry to point to the CPU's L2 page.
117 * Specifically, in hat_pcp_update(), we're copying from the process's
118 * HAT_COPIED_32 HAT into the page tables specific to this CPU.
119 *
120 * As part of the implementation of kernel page table isolation, this was also
121 * extended to 64-bit processes. When a 64-bit process runs, we'll copy their L3

```

```

122 * PTEs across into the current CPU's L3 page table. (As we can't do the
123 * first-L3-entry trick for 64-bit processes, ->hci_pcp_l2ptes is unused in this
124 * case.)
125 *
126 * The use of per-CPU page tables has a lot of implementation ramifications. A
127 * HAT that runs a user process will be flagged with the HAT_COPIED flag to
128 * indicate that it is using the per-CPU page table functionality. In tandem
129 * with the HAT, the top-level htable_t will be flagged with the HTABLE_COPIED
130 * flag. If the HAT represents a 32-bit process, then we will also set the
131 * HAT_COPIED_32 flag on that hat_t.
132 *
133 * These two flags work together. The top-level htable_t when using per-CPU page
134 * tables is 'virtual'. We never allocate a ptable for this htable_t (i.e.
135 * ht->ht_pfn is PFN_INVALID). Instead, when we need to modify a PTE in an
136 * HTABLE_COPIED ptable, x86pte_access_pagetable() will redirect any accesses to
137 * ht_hat->hat_copied_ptes.
138 *
139 * Of course, such a modification won't actually modify the HAT_PCP page tables
140 * that were copied from the HAT_COPIED htable. When we change the top level
141 * page table entries (L2 PTEs for a 32-bit process and L3 PTEs for a 64-bit
142 * process), we need to make sure to trigger hat_pcp_update() on all CPUs that
143 * are currently tied to this HAT (including the current CPU).
144 *
145 * To do this, PCP piggy-backs on TLB invalidation, specifically via the
146 * hat_tlb_inval() path from link_ptp() and unlink_ptp().
147 *
148 * (Importantly, in all such cases, when this is in operation, the top-level
149 * entry should not be able to refer to an actual page table entry that can be
150 * changed and consolidated into a large page. If large page consolidation is
151 * required here, then there will be much that needs to be reconsidered.)
152 *
153 * -----
154 * Kernel Page Table Isolation and the Per-CPU HAT
155 * -----
156 *
157 * All Intel CPUs that support speculative execution and paging are subject to a
158 * series of bugs that have been termed 'Meltdown'. These exploits allow a user
159 * process to read kernel memory through cache side channels and speculative
160 * execution. To mitigate this on vulnerable CPUs, we need to use a technique
161 * called kernel page table isolation. What this requires is that we have two
162 * different page table roots. When executing in kernel mode, we will use a %cr3
163 * value that has both the user and kernel pages. However when executing in user
164 * mode, we will need to have a %cr3 that has all of the user pages; however,
165 * only a subset of the kernel pages required to operate.
166 *
167 * These kernel pages that we need mapped are:
168 *
169 *   o Kernel Text that allows us to switch between the cr3 values.
170 *   o The current global descriptor table (GDT)
171 *   o The current interrupt descriptor table (IDT)
172 *   o The current task switching state (TSS)
173 *   o The current local descriptor table (LDT)
174 *   o Stacks and scratch space used by the interrupt handlers
175 *
176 * For more information on the stack switching techniques, construction of the
177 * trampolines, and more, please see i86pc/ml/kpti_trampoline.s. The most
178 * important part of these mappings are the following two constraints:
179 *
180 *   o The mappings are all per-CPU (except for read-only text)
181 *   o The mappings are static. They are all established before the CPU is
182 *     started (with the exception of the boot CPU).
183 *
184 * To facilitate the kernel page table isolation we employ our per-CPU
185 * page tables discussed in the previous section and add the notion of a per-CPU
186 * HAT. Fundamentally we have a second page table root. There is both a kernel
187 * page table (hci_pcp_l3ptes), and a user L3 page table (hci_user_l3ptes).

```

```

188 * Both will have the user page table entries copied into them, the same way
189 * that we discussed in the section 'Per-CPU Page Tables'.
190 *
191 * The complex part of this is how do we construct the set of kernel mappings
192 * that should be present when running with the user page table. To answer that,
193 * we add the notion of a per-CPU HAT. This HAT functions like a normal HAT,
194 * except that it's not really associated with an address space the same way
195 * that other HATs are.
196 *
197 * This HAT lives off of the 'struct hat_cpu_info' which is a member of the
198 * machcpu in the member hci_user_hat. We use this per-CPU HAT to create the set
199 * of kernel mappings that should be present on this CPU. The kernel mappings
200 * are added to the per-CPU HAT through the function hati_cpu_punchin(). Once a
201 * mapping has been punched in, it may not be punched out. The reason that we
202 * opt to leverage a HAT structure is that it knows how to allocate and manage
203 * all of the lower level page tables as required.
204 *
205 * Because all of the mappings are present at the beginning of time for this CPU
206 * and none of the mappings are in the kernel pageable segment, we don't have to
207 * worry about faulting on these HAT structures and thus the notion of the
208 * current HAT that we're using is always the appropriate HAT for the process
209 * (usually a user HAT or the kernel's HAT).
210 *
211 * A further constraint we place on the system with these per-CPU HATs is that
212 * they are not subject to htable_steal(). Because each CPU will have a rather
213 * fixed number of page tables, the same way that we don't steal from the
214 * kernel's HAT, it was determined that we should not steal from this HAT due to
215 * the complications involved and somewhat criminal nature of htable_steal().
216 *
217 * The per-CPU HAT is initialized in hat_pcp_setup() which is called as part of
218 * onlining the CPU, but before the CPU is actually started. The per-CPU HAT is
219 * removed in hat_pcp_takedown() which is called when a CPU is being offlined to
220 * be removed from the system (which is different from what psradm usually
221 * does).
222 *
223 * Finally, once the CPU has been onlined, the set of mappings in the per-CPU
224 * HAT must not change. The HAT related functions that we call are not meant to
225 * be called when we're switching between processes. For example, it is quite
226 * possible that if they were, they would try to grab an htable_mutex which
227 * another thread might have. One needs to treat hat_switch() as though they
228 * were above LOCK_LEVEL and therefore _must not_ block under any circumstance.
229 */
230
231 #include <sys/machparam.h>
232 #include <sys/machsystem.h>
233 #include <sys/mman.h>
234 #include <sys/types.h>
235 #include <sys/system.h>
236 #include <sys/cpuvar.h>
237 #include <sys/thread.h>
238 #include <sys/proc.h>
239 #include <sys/cpu.h>
240 #include <sys/kmem.h>
241 #include <sys/disp.h>
242 #include <sys/shm.h>
243 #include <sys/sysmacros.h>
244 #include <sys/machparam.h>
245 #include <sys/vmem.h>
246 #include <sys/vmsystem.h>
247 #include <sys/promif.h>
248 #include <sys/var.h>
249 #include <sys/x86_archext.h>
250 #include <sys/atomic.h>
251 #include <sys/bitmap.h>
252 #include <sys/controlregs.h>
253 #include <sys/bootconf.h>

```

```

254 #include <sys/bootsvcs.h>
255 #include <sys/bootinfo.h>
256 #include <sys/archsystem.h>

258 #include <vm/seg_kmem.h>
259 #include <vm/hat_i86.h>
260 #include <vm/as.h>
261 #include <vm/seg.h>
262 #include <vm/page.h>
263 #include <vm/seg_kp.h>
264 #include <vm/seg_kpm.h>
265 #include <vm/vm_dep.h>
266 #ifdef __xpv
267 #include <sys/hypervisor.h>
268 #endif
269 #include <vm/kboot_mmu.h>
270 #include <vm/seg_spt.h>

272 #include <sys/cmn_err.h>

274 /*
275  * Basic parameters for hat operation.
276  */
277 struct hat_mmu_info mmu;

279 /*
280  * The page that is the kernel's top level pagetable.
281  */
282 * For 32 bit PAE support on i86pc, the kernel hat will use the 1st 4 entries
283 * on this 4K page for its top level page table. The remaining groups of
284 * 4 entries are used for per processor copies of user PCP pagetables for
285 * 4 entries are used for per processor copies of user VLP pagetables for
286 * running threads. See hat_switch() and reload_pae32() for details.
287 *
288 * pcp_page[0..3] - level==2 PTEs for kernel HAT
289 * pcp_page[4..7] - level==2 PTEs for user thread on cpu 0
290 * pcp_page[8..11] - level==2 PTE for user thread on cpu 1
291 * vlp_page[0..3] - level==2 PTEs for kernel HAT
292 * vlp_page[4..7] - level==2 PTEs for user thread on cpu 0
293 * vlp_page[8..11] - level==2 PTE for user thread on cpu 1
294 * etc...
295 *
296 * On the 64-bit kernel, this is the normal root of the page table and there is
297 * nothing special about it when used for other CPUs.
298 */
299 static x86pte_t *pcp_page;
300 static x86pte_t *vlp_page;

302 /*
303  * forward declaration of internal utility routines
304  */
305 static x86pte_t hati_update_pte(htable_t *ht, uint_t entry, x86pte_t expected,
306 x86pte_t new);

308 /*
309  * The kernel address space exists in all non-HAT_COPIED HATs. To implement this
310  * the kernel reserves a fixed number of entries in the topmost level(s) of page
311  * tables. The values are setup during startup and then copied to every user hat
312  * created by hat_alloc(). This means that kernelbase must be:
313  * The kernel address space exists in all HATs. To implement this the
314  * kernel reserves a fixed number of entries in the topmost level(s) of page
315  * tables. The values are setup during startup and then copied to every user
316  * hat created by hat_alloc(). This means that kernelbase must be:
317  *
318  * 4Meg aligned for 32 bit kernels
319  * 512Gig aligned for x86_64 64 bit kernel

```

```

311 *
312 * The hat_kernel_range_ts describe what needs to be copied from kernel hat
313 * to each user hat.
314 */
315 typedef struct hat_kernel_range {
316     level_t      hkr_level;
317     uintptr_t    hkr_start_va;
318     uintptr_t    hkr_end_va;      /* zero means to end of memory */
319 } hat_kernel_range_t;
320 #define NUM_KERNEL_RANGE 2
321 static hat_kernel_range_t kernel_ranges[NUM_KERNEL_RANGE];
322 static int num_kernel_ranges;

324 uint_t use_boot_reserve = 1; /* cleared after early boot process */
325 uint_t can_steal_post_boot = 0; /* set late in boot to enable stealing */

327 /*
328  * enable_lgpg: controls lg page support for user applications.
329  * By default, lg pages are exported to user applications. enable_lgpg can
330  * be set to 0 to not export.
331  */
332 int     enable_lgpg = 1;

334 /*
335  * AMD shanghai processors provide better management of lgb ptes in its tlb.
336  * By default, lg page support will be disabled for pre-shanghai AMD
337  * processors that don't have optimal tlb support for the lg page size.
338  * chk_optimal_lgtlb can be set to 0 to force lg page support on sub-optimal
339  * processors.
340  */
341 int     chk_optimal_lgtlb = 1;

344 #ifdef DEBUG
345 uint_t maplgcnt;
346 #endif

349 /*
350  * A cpuset for all cpus. This is used for kernel address cross calls, since
351  * the kernel addresses apply to all cpus.
352  */
353 cpuset_t khata_cpuset;

355 /*
356  * management stuff for hat structures
357  */
358 kmutex_t      hat_list_lock;
359 kcondvar_t    hat_list_cv;
360 kmem_cache_t  *hat_cache;
361 kmem_cache_t  *hat_hash_cache;
362 kmem_cache_t  *hat32_hash_cache;
363 kmem_cache_t  *vlp_hash_cache;

364 /*
365  * Simple statistics
366  */
367 struct hatstats hatstat;

369 /*
370  * Some earlier hypervisor versions do not emulate cmpxchg of PTEs
371  * correctly. For such hypervisors we must set PT_USER for kernel
372  * entries ourselves (normally the emulation would set PT_USER for
373  * kernel entries and PT_USER|PT_GLOBAL for user entries). pt_kern is
374  * thus set appropriately. Note that dboot/kbm is OK, as only the full
375  * HAT uses cmpxchg() and the other paths (hypercall etc.) were never

```

```

376 * incorrect.
377 */
378 int pt_kern;

191 /*
192 * useful stuff for atomic access/clearing/setting REF/MOD/RO bits in page_t's.
193 */
194 extern void atomic_orb(uchar_t *addr, uchar_t val);
195 extern void atomic_andb(uchar_t *addr, uchar_t val);

380 #ifndef __xpv
381 extern pfn_t memseg_get_start(struct memseg *);
382 #endif

384 #define PP_GETRM(pp, rmmask)    (pp->p_nrm & rmmask)
385 #define PP_ISMOD(pp)          PP_GETRM(pp, P_MOD)
386 #define PP_ISREF(pp)          PP_GETRM(pp, P_REF)
387 #define PP_ISRO(pp)           PP_GETRM(pp, P_RO)

389 #define PP_SETRM(pp, rm)       atomic_orb(&(pp->p_nrm), rm)
390 #define PP_SETRM(pp, P_MOD)    PP_SETRM(pp, P_MOD)
391 #define PP_SETREF(pp)          PP_SETRM(pp, P_REF)
392 #define PP_SETRO(pp)           PP_SETRM(pp, P_RO)

394 #define PP_CLRRM(pp, rm)       atomic_andb(&(pp->p_nrm), ~(rm))
395 #define PP_CLRRM(pp, P_MOD)    PP_CLRRM(pp, P_MOD)
396 #define PP_CLRRM(pp, P_REF)    PP_CLRRM(pp, P_REF)
397 #define PP_CLRRM(pp, P_RO)     PP_CLRRM(pp, P_RO)
398 #define PP_CLRRM(pp, P_MOD | P_REF | P_RO) PP_CLRRM(pp, P_MOD | P_REF | P_RO)

400 /*
401 * kmem cache constructor for struct hat
402 */
403 /*ARGSUSED*/
404 static int
405 hati_constructor(void *buf, void *handle, int kmflags)
406 {
407     hat_t *hat = buf;

409     mutex_init(&hat->hat_mutex, NULL, MUTEX_DEFAULT, NULL);
410     bzero(hat->hat_pages_mapped,
411           sizeof (pgcnt_t) * (mmu.max_page_level + 1));
412     hat->hat_ism_pgcnt = 0;
413     hat->hat_stats = 0;
414     hat->hat_flags = 0;
415     CPuset_ZERO(hat->hat_cpus);
416     hat->hat_htable = NULL;
417     hat->hat_ht_hash = NULL;
418     return (0);
419 }

421 /*
422 * Put it at the start of the global list of all hats (used by stealing)
423 */
424 * kas.a_hat is not in the list but is instead used to find the
425 * first and last items in the list.
426 *
427 * - kas.a_hat->hat_next points to the start of the user hats.
428 *   The list ends where hat->hat_next == NULL
429 *
430 * - kas.a_hat->hat_prev points to the last of the user hats.
431 *   The list begins where hat->hat_prev == NULL
432 */
433 static void
434 hat_list_append(hat_t *hat)
435 {

```

```

436     mutex_enter(&hat_list_lock);
437     hat->hat_prev = NULL;
438     hat->hat_next = kas.a_hat->hat_next;
439     if (hat->hat_next)
440         hat->hat_next->hat_prev = hat;
441     else
442         kas.a_hat->hat_prev = hat;
443     kas.a_hat->hat_next = hat;
444     mutex_exit(&hat_list_lock);
445 }

447 /*
448 * Allocate a hat structure for as. We also create the top level
449 * htable and initialize it to contain the kernel hat entries.
450 */
451 hat_t *
452 hat_alloc(struct as *as)
453 {
454     hat_t *hat;
455     htable_t *ht; /* top level htable */
456     uint_t use_copied;
457     uint_t use_vlp;
458     r;
459     hat_kernel_range_t *rp;
460     uintptr_t va;
461     uintptr_t eva;
462     uint_t start;
463     uint_t cnt;
464     htable_t *src;
465     boolean_t use_hat32_cache;

466     /*
467     * Once we start creating user process HATs we can enable
468     * the htable_steal() code.
469     */
470     if (can_steal_post_boot == 0)
471         can_steal_post_boot = 1;

473     ASSERT(AS_WRITE_HELD(as));
474     hat = kmem_cache_alloc(hat_cache, KM_SLEEP);
475     hat->hat_as = as;
476     mutex_init(&hat->hat_mutex, NULL, MUTEX_DEFAULT, NULL);
477     ASSERT(hat->hat_flags == 0);

479 #if defined(__xpv)
480     /*
481     * No PCP stuff on the hypervisor due to the 64-bit split top level
482     * No VLP stuff on the hypervisor due to the 64-bit split top level
483     * page tables. On 32-bit it's not needed as the hypervisor takes
484     * care of copying the top level PTEs to a below 4Gig page.
485     */
486     use_copied = 0;
487     use_hat32_cache = B_FALSE;
488     hat->hat_max_level = mmu.max_level;
489     hat->hat_num_copied = 0;
490     hat->hat_flags = 0;
491     use_vlp = 0;
492 #else /* __xpv */

493     /*
494     * All processes use HAT_COPIED on the 64-bit kernel if KPTI is
495     * turned on.
496     */
497     if (ttoproc(curthread)->p_model == DATAMODEL_ILP32) {
498         use_copied = 1;
499         hat->hat_max_level = mmu.max_level32;

```

```

499     hat->hat_num_copied = mmu.num_copied_ents32;
500     use_hat32_cache = B_TRUE;
501     hat->hat_flags |= HAT_COPIED_32;
502     HATSTAT_INC(hs_hat_copied32);
503 } else if (kpti_enable == 1) {
504     use_copied = 1;
505     hat->hat_max_level = mmu.max_level;
506     hat->hat_num_copied = mmu.num_copied_ents;
507     use_hat32_cache = B_FALSE;
508     HATSTAT_INC(hs_hat_copied64);
509 } else {
510     use_copied = 0;
511     use_hat32_cache = B_FALSE;
512     hat->hat_max_level = mmu.max_level;
513     hat->hat_num_copied = 0;
514     hat->hat_flags = 0;
515     HATSTAT_INC(hs_hat_normal64);
516 }
517 /* 32 bit processes uses a VLP style hat when running with PAE */
518 #if defined(__amd64)
519     use_vlp = (tproc(curthread)->p_model == DATAMODEL_ILP32);
520 #elif defined(__i386)
521     use_vlp = mmu.pae_hat;
522 #endif
523 #endif /* __xpv */
524 if (use_copied) {
525     hat->hat_flags |= HAT_COPIED;
526     bzero(hat->hat_copied_ptes, sizeof(hat->hat_copied_ptes));
527 if (use_vlp) {
528     hat->hat_flags = HAT_VLP;
529     bzero(hat->hat_vlp_ptes, VLP_SIZE);
530 }
531 }
532 /*
533  * Allocate the htable hash. For 32-bit PCP processes we use the
534  * hat32_hash_cache. However, for 64-bit PCP processes we do not as the
535  * number of entries that they have to handle is closer to
536  * hat_hash_cache in count (though there will be more wastage when we
537  * have more DRAM in the system and thus push down the user address
538  * range).
539  * Allocate the htable hash
540  */
541 if (use_hat32_cache) {
542     hat->hat_num_hash = mmu.hat32_hash_cnt;
543     hat->hat_ht_hash = kmem_cache_alloc(hat32_hash_cache, KM_SLEEP);
544 if ((hat->hat_flags & HAT_VLP) {
545     hat->hat_num_hash = mmu.vlp_hash_cnt;
546     hat->hat_ht_hash = kmem_cache_alloc(vlp_hash_cache, KM_SLEEP);
547 } else {
548     hat->hat_num_hash = mmu.hash_cnt;
549     hat->hat_ht_hash = kmem_cache_alloc(hat_hash_cache, KM_SLEEP);
550 }
551 bzero(hat->hat_ht_hash, hat->hat_num_hash * sizeof(htable_t *));
552 }
553 /*
554  * Initialize Kernel HAT entries at the top of the top level page
555  * tables for the new hat.
556  */
557 hat->hat_htable = NULL;
558 hat->hat_ht_cached = NULL;
559 XPV_DISALLOW_MIGRATE();
560 ht = htable_create(hat, (uintptr_t)0, TOP_LEVEL(hat), NULL);
561 hat->hat_htable = ht;
562 #if defined(__amd64)
563 if (hat->hat_flags & HAT_COPIED)

```

```

312     if (hat->hat_flags & HAT_VLP)
313         goto init_done;
314 #endif
315
316 for (r = 0; r < num_kernel_ranges; ++r) {
317     rp = &kernel_ranges[r];
318     for (va = rp->hkr_start_va; va != rp->hkr_end_va;
319          va += cnt * LEVEL_SIZE(rp->hkr_level)) {
320
321         if (rp->hkr_level == TOP_LEVEL(hat))
322             ht = hat->hat_htable;
323         else
324             ht = htable_create(hat, va, rp->hkr_level,
325                                NULL);
326
327         start = htable_va2entry(va, ht);
328         cnt = HTABLE_NUM_PTES(ht) - start;
329         eva = va +
330              ((uintptr_t)cnt << LEVEL_SHIFT(rp->hkr_level));
331         if (rp->hkr_end_va != 0 &&
332             (eva > rp->hkr_end_va || eva == 0))
333             cnt = htable_va2entry(rp->hkr_end_va, ht) -
334                  start;
335 #if defined(__i386) && !defined(__xpv)
336         if (ht->ht_flags & HTABLE_COPIED) {
337             bcopy(&pcp_page[start],
338                  &hat->hat_copied_ptes[start],
339                  cnt * sizeof(x86pte_t));
340             continue;
341         }
342 #endif
343         src = htable_lookup(kas.a_hat, va, rp->hkr_level);
344         ASSERT(src != NULL);
345         x86pte_copy(src, ht, start, cnt);
346         htable_release(src);
347     }
348 }
349
350 init_done:
351 #if defined(__xpv)
352     /*
353      * Pin top level page tables after initializing them
354      */
355     xen_pin(hat->hat_htable->ht_pfn, mmu.max_level);
356 #if defined(__amd64)
357     xen_pin(hat->hat_user_ptable, mmu.max_level);
358 #endif
359 #endif
360 #endif
361     XPV_ALLOW_MIGRATE();
362
363     hat_list_append(hat);
364
365     return (hat);
366 }
367
368 #if !defined(__xpv)
369 /*
370  * Cons up a HAT for a CPU. This represents the user mappings. This will have
371  * various kernel pages punched into it manually. Importantly, this hat is
372  * ineligible for stealing. We really don't want to deal with this ever
373  * faulting and figuring out that this is happening, much like we don't with

```

```

614 * kas.
615 */
616 static hat_t *
617 hat_cpu_alloc(cpu_t *cpu)
618 {
619     hat_t *hat;
620     htable_t *ht;

622     hat = kmem_cache_alloc(hat_cache, KM_SLEEP);
623     hat->hat_as = NULL;
624     mutex_init(&hat->hat_mutex, NULL, MUTEX_DEFAULT, NULL);
625     hat->hat_max_level = mmu.max_level;
626     hat->hat_num_copied = 0;
627     hat->hat_flags = HAT_PCP;

629     hat->hat_num_hash = mmu.hash_cnt;
630     hat->hat_ht_hash = kmem_cache_alloc(hat_hash_cache, KM_SLEEP);
631     bzero(hat->hat_ht_hash, hat->hat_num_hash * sizeof(htable_t *));

633     hat->hat_next = hat->hat_prev = NULL;

635     /*
636     * Because this HAT will only ever be used by the current CPU, we'll go
637     * ahead and set the CPUSET up to only point to the CPU in question.
638     * Put it at the start of the global list of all hats (used by stealing)
639     *
640     * kas.a_hat is not in the list but is instead used to find the
641     * first and last items in the list.
642     *
643     * - kas.a_hat->hat_next points to the start of the user hats.
644     *   The list ends where hat->hat_next == NULL
645     *
646     * - kas.a_hat->hat_prev points to the last of the user hats.
647     *   The list begins where hat->hat_prev == NULL
648     */
649     CPUSET_ADD(hat->hat_cpus, cpu->cpu_id);
650     mutex_enter(&hat_list_lock);
651     hat->hat_prev = NULL;
652     hat->hat_next = kas.a_hat->hat_next;
653     if (hat->hat_next)
654         hat->hat_next->hat_prev = hat;
655     else
656         kas.a_hat->hat_prev = hat;
657     kas.a_hat->hat_next = hat;
658     mutex_exit(&hat_list_lock);

660     hat->hat_htable = NULL;
661     hat->hat_ht_cached = NULL;
662     ht = htable_create(hat, (uintptr_t)0, TOP_LEVEL(hat), NULL);
663     hat->hat_htable = ht;

665     hat_list_append(hat);

667     return (hat);
668 }
669 #endif /* !__xpv */

671 /*
672 * process has finished executing but as has not been cleaned up yet.
673 */
674 /* ARGSUSED */
675 void
676 hat_free_start(hat_t *hat)
677 {
678     ASSERT(AS_WRITE_HELD(hat->hat_as));

```

```

661     /*
662     * If the hat is currently a stealing victim, wait for the stealing
663     * to finish. Once we mark it as HAT_FREEING, htable_steal()
664     * won't look at its pagetables anymore.
665     */
666     mutex_enter(&hat_list_lock);
667     while (hat->hat_flags & HAT_VICTIM)
668         cv_wait(&hat_list_cv, &hat_list_lock);
669     hat->hat_flags |= HAT_FREEING;
670     mutex_exit(&hat_list_lock);
671 }

673 /*
674 * An address space is being destroyed, so we destroy the associated hat.
675 */
676 void
677 hat_free_end(hat_t *hat)
678 {
679     kmem_cache_t *cache;

681     ASSERT(hat->hat_flags & HAT_FREEING);

683     /*
684     * must not be running on the given hat
685     */
686     ASSERT(CPU->cpu_current_hat != hat);

688     /*
689     * Remove it from the list of HATs
690     */
691     mutex_enter(&hat_list_lock);
692     if (hat->hat_prev)
693         hat->hat_prev->hat_next = hat->hat_next;
694     else
695         kas.a_hat->hat_next = hat->hat_next;
696     if (hat->hat_next)
697         hat->hat_next->hat_prev = hat->hat_prev;
698     else
699         kas.a_hat->hat_prev = hat->hat_prev;
700     mutex_exit(&hat_list_lock);
701     hat->hat_next = hat->hat_prev = NULL;

703 #if defined(__xpv)
704     /*
705     * On the hypervisor, unpin top level page table(s)
706     */
707     VERIFY3U(hat->hat_flags & HAT_PCP, ==, 0);
708     xen_unpin(hat->hat_htable->ht_pfn);
709 #if defined(__amd64)
710     xen_unpin(hat->hat_user_ptable);
711 #endif
712 #endif

714     /*
715     * Make a pass through the htables freeing them all up.
716     */
717     htable_purge_hat(hat);

719     /*
720     * Decide which kmem cache the hash table came from, then free it.
721     */
722     if (hat->hat_flags & HAT_COPIED) {
723 #if defined(__amd64)
724         if (hat->hat_flags & HAT_COPIED_32) {
725             cache = hat32_hash_cache;
726         } else {

```

```

458     if (hat->hat_flags & HAT_VLP)
459         cache = vlp_hash_cache;
460     else
461         cache = hat_hash_cache;
462     }
463 #endif
464     cache = hat32_hash_cache;
465 #endif
466 } else {
467     cache = hat_hash_cache;
468 }
469 kmem_cache_free(cache, hat->hat_ht_hash);
470 hat->hat_ht_hash = NULL;
471
472 hat->hat_flags = 0;
473 hat->hat_max_level = 0;
474 hat->hat_num_copied = 0;
475 kmem_cache_free(hat_cache, hat);
476 }
477
478 unchanged portion omitted
479
480 /*
481 * Determine the number of slots that are in used in the top-most level page
482 * table for user memory. This is based on _userlimit. In effect this is similar
483 * to htable_va2entry, but without the convenience of having an htable.
484 */
485 void
486 mmu_calc_user_slots(void)
487 {
488     uint_t ent, nptes;
489     uintptr_t shift;
490
491     nptes = mmu.top_level_count;
492     shift = _userlimit >> mmu.level_shift[mmu.max_level];
493     ent = shift & (nptes - 1);
494
495     /*
496      * Ent tells us the slot that the page for _userlimit would fit in. We
497      * need to add one to this to cover the total number of entries.
498      */
499     mmu.top_level_uslots = ent + 1;
500
501     /*
502      * When running 32-bit compatability processes on a 64-bit kernel, we
503      * will only need to use one slot.
504      */
505     mmu.top_level_uslots32 = 1;
506
507     /*
508      * Record the number of PCP page table entries that we'll need to copy
509      * around. For 64-bit processes this is the number of user slots. For
510      * 32-bit proceses, this is 4 1 GiB pages.
511      */
512     mmu.num_copied_ents = mmu.top_level_uslots;
513     mmu.num_copied_ents32 = 4;
514 }
515
516 /*
517 * Initialize hat data structures based on processor MMU information.
518 */
519 void
520 mmu_init(void)
521 {
522     uint_t max_htables;
523     uint_t pa_bits;
524     uint_t va_bits;

```

```

839     int i;
840
841     /*
842      * If CPU enabled the page table global bit, use it for the kernel
843      * This is bit 7 in CR4 (PGE - Page Global Enable).
844      */
845     if (is_x86_feature(x86_featureset, X86FSET_PGE) &&
846         (getcr4() & CR4_PGE) != 0)
847         mmu.pt_global = PT_GLOBAL;
848
849 #if !defined(__xpv)
850     /*
851      * The 64-bit x86 kernel has split user/kernel page tables. As such we
852      * cannot have the global bit set. The simplest way for us to deal with
853      * this is to just say that pt_global is zero, so the global bit isn't
854      * present.
855      */
856     if (kpti_enable == 1)
857         mmu.pt_global = 0;
858 #endif
859
860     /*
861      * Detect NX and PAE usage.
862      */
863     mmu.pae_hat = kbm_pae_support;
864     if (kbm_nx_support)
865         mmu.pt_nx = PT_NX;
866     else
867         mmu.pt_nx = 0;
868
869     /*
870      * Use CPU info to set various MMU parameters
871      */
872     cpuid_get_addrsize(CPU, &pa_bits, &va_bits);
873
874     if (va_bits < sizeof(void *) * NBBY) {
875         mmu.hole_start = (1ul << (va_bits - 1));
876         mmu.hole_end = 0ul - mmu.hole_start - 1;
877     } else {
878         mmu.hole_end = 0;
879         mmu.hole_start = mmu.hole_end - 1;
880     }
881 #if defined(OPTERON_ERRATUM_121)
882     /*
883      * If erratum 121 has already been detected at this time, hole_start
884      * contains the value to be subtracted from mmu.hole_start.
885      */
886     ASSERT(hole_start == 0 || opteron_erratum_121 != 0);
887     hole_start = mmu.hole_start - hole_start;
888 #else
889     hole_start = mmu.hole_start;
890 #endif
891     hole_end = mmu.hole_end;
892
893     mmu.highest_pfn = mmu_btop((1ull << pa_bits) - 1);
894     if (mmu.pae_hat == 0 && pa_bits > 32)
895         mmu.highest_pfn = PFN_4G - 1;
896
897     if (mmu.pae_hat) {
898         mmu.pte_size = 8; /* 8 byte PTEs */
899         mmu.pte_size_shift = 3;
900     } else {
901         mmu.pte_size = 4; /* 4 byte PTEs */
902         mmu.pte_size_shift = 2;
903     }

```



```

905     if (mmu.pae_hat && !is_x86_feature(x86_featureset, X86FSET_PAE))
906         panic("Processor does not support PAE");

908     if (!is_x86_feature(x86_featureset, X86FSET_CX8))
909         panic("Processor does not support cmpxchg8b instruction");

911 #if defined(__amd64)

913     mmu.num_level = 4;
914     mmu.max_level = 3;
915     mmu.ptes_per_table = 512;
916     mmu.top_level_count = 512;

918     /*
919      * 32-bit processes only use 1 GB ptes.
920      */
921     mmu.max_level32 = 2;

923     mmu.level_shift[0] = 12;
924     mmu.level_shift[1] = 21;
925     mmu.level_shift[2] = 30;
926     mmu.level_shift[3] = 39;

928 #elif defined(__i386)

930     if (mmu.pae_hat) {
931         mmu.num_level = 3;
932         mmu.max_level = 2;
933         mmu.ptes_per_table = 512;
934         mmu.top_level_count = 4;

936         mmu.level_shift[0] = 12;
937         mmu.level_shift[1] = 21;
938         mmu.level_shift[2] = 30;

940     } else {
941         mmu.num_level = 2;
942         mmu.max_level = 1;
943         mmu.ptes_per_table = 1024;
944         mmu.top_level_count = 1024;

946         mmu.level_shift[0] = 12;
947         mmu.level_shift[1] = 22;
948     }

950 #endif /* __i386 */

952     for (i = 0; i < mmu.num_level; ++i) {
953         mmu.level_size[i] = 1UL << mmu.level_shift[i];
954         mmu.level_offset[i] = mmu.level_size[i] - 1;
955         mmu.level_mask[i] = ~mmu.level_offset[i];
956     }

958     set_max_page_level();
959     mmu_calc_user_slots();

961     mmu_page_sizes = mmu.max_page_level + 1;
962     mmu_exported_page_sizes = mmu.umax_page_level + 1;

964     /* restrict legacy applications from using pagesizes lg and above */
965     mmu_legacy_page_sizes =
966         (mmu_exported_page_sizes > 2) ? 2 : mmu_exported_page_sizes;

969     for (i = 0; i <= mmu.max_page_level; ++i) {
970         mmu.pte_bits[i] = PT_VALID | pt_kern;

```

```

971         if (i > 0)
972             mmu.pte_bits[i] |= PT_PAGESIZE;
973     }

975     /*
976      * NOTE Legacy 32 bit PAE mode only has the P_VALID bit at top level.
977      */
978     for (i = 1; i < mmu.num_level; ++i)
979         mmu.ptp_bits[i] = PT_PTPBITS;

981 #if defined(__i386)
982     mmu.ptp_bits[2] = PT_VALID;
983 #endif

985     /*
986      * Compute how many hash table entries to have per process for htables.
987      * We start with 1 page's worth of entries.
988      * If physical memory is small, reduce the amount need to cover it.
989      */
990     max_htables = physmax / mmu.ptes_per_table;
991     mmu.hash_cnt = MMU_PAGESIZE / sizeof (htable_t *);
992     while (mmu.hash_cnt > 16 && mmu.hash_cnt >= max_htables)
993         mmu.hash_cnt >>= 1;
994     mmu.hat32_hash_cnt = mmu.hash_cnt;
995     mmu.vlp_hash_cnt = mmu.hash_cnt;
996 #endif

997 #if defined(__amd64)
998     /*
999      * If running in 64 bits and physical memory is large,
1000      * increase the size of the cache to cover all of memory for
1001      * a 64 bit process.
1002      */
1003     #define HASH_MAX_LENGTH 4
1004     while (mmu.hash_cnt * HASH_MAX_LENGTH < max_htables)
1005         mmu.hash_cnt <<= 1;
1006 #endif
1007 }

1010 /*
1011  * initialize hat data structures
1012  */
1013 void
1014 hat_init()
1015 {
1016     #if defined(__i386)
1017         /*
1018          * _userlimit must be aligned correctly
1019          */
1020         if ((*_userlimit & LEVEL_MASK(1)) != *_userlimit) {
1021             prom_printf("hat_init(): _userlimit=%p, not aligned at %p\n",
1022                 (void *)_userlimit, (void *)LEVEL_SIZE(1));
1023             halt("hat_init(): Unable to continue");
1024         }
1025     #endif

1027     cv_init(&hat_list_cv, NULL, CV_DEFAULT, NULL);

1029     /*
1030      * initialize kmem caches
1031      */
1032     htable_init();
1033     hment_init();

1035     hat_cache = kmem_cache_create("hat_t",

```

```

1036         sizeof (hat_t), 0, hati_constructor, NULL, NULL,
1037         NULL, 0, 0);

1039     hat_hash_cache = kmem_cache_create("HatHash",
1040     mmu.hash_cnt * sizeof (htable_t *), 0, NULL, NULL, NULL,
1041     NULL, 0, 0);

1043     /*
1044     * 32-bit PCP hats can use a smaller hash table size on large memory
1045     * machines
1046     * VLP hats can use a smaller hash table size on large memroy machines
1047     */
1048     if (mmu.hash_cnt == mmu.hat32_hash_cnt) {
1049         hat32_hash_cache = hat_hash_cache;
1050     } else if (mmu.hash_cnt == mmu.vlp_hash_cnt) {
1051         vlp_hash_cache = hat_hash_cache;
1052     } else {
1053         hat32_hash_cache = kmem_cache_create("Hat32Hash",
1054         mmu.hat32_hash_cnt * sizeof (htable_t *), 0, NULL, NULL,
1055         NULL, NULL, 0, 0);
1056         vlp_hash_cache = kmem_cache_create("HatVlpHash",
1057         mmu.vlp_hash_cnt * sizeof (htable_t *), 0, NULL, NULL, NULL,
1058         NULL, 0, 0);
1059     }

1060     /*
1061     * Set up the kernel's hat
1062     */
1063     AS_LOCK_ENTER(&kas, RW_WRITER);
1064     kas.a_hat = kmem_cache_alloc(hat_cache, KM_NOSLEEP);
1065     mutex_init(&kas.a_hat->hat_mutex, NULL, MUTEX_DEFAULT, NULL);
1066     kas.a_hat->hat_as = &kas;
1067     kas.a_hat->hat_flags = 0;
1068     AS_LOCK_EXIT(&kas);

1069     CPuset_ZERO(khat_cpuset);
1070     CPuset_ADD(khat_cpuset, CPU->cpu_id);

1071     /*
1072     * The kernel HAT doesn't use PCP regardless of architectures.
1073     */
1074     ASSERT3U(mmu.max_level, >, 0);
1075     kas.a_hat->hat_max_level = mmu.max_level;
1076     kas.a_hat->hat_num_copied = 0;

1077     /*
1078     * The kernel hat's next pointer serves as the head of the hat list .
1079     * The kernel hat's prev pointer tracks the last hat on the list for
1080     * htable_steal() to use.
1081     */
1082     kas.a_hat->hat_next = NULL;
1083     kas.a_hat->hat_prev = NULL;

1084     /*
1085     * Allocate an htable hash bucket for the kernel
1086     * XX64 - tune for 64 bit procs
1087     */
1088     kas.a_hat->hat_num_hash = mmu.hash_cnt;
1089     kas.a_hat->hat_ht_hash = kmem_cache_alloc(hat_hash_cache, KM_NOSLEEP);
1090     bzero(kas.a_hat->hat_ht_hash, mmu.hash_cnt * sizeof (htable_t *));

1091     /*
1092     * zero out the top level and cached htable pointers
1093     */
1094     kas.a_hat->hat_ht_cached = NULL;
1095     kas.a_hat->hat_htable = NULL;

```

```

1097     /*
1098     * Pre-allocate hrm_hashtab before enabling the collection of
1099     * refmod statistics. Allocating on the fly would mean us
1100     * running the risk of suffering recursive mutex enters or
1101     * deadlocks.
1102     */
1103     hrm_hashtab = kmem_zalloc(HRM_HASHSIZE * sizeof (struct hrmstat *),
1104     KM_SLEEP);
1105 }

1106 extern void kpti_trampoline_start();
1107 extern void kpti_trampoline_end();

1108 extern void kdi_isr_start();
1109 extern void kdi_isr_end();

1110 extern gate_desc_t kdi_idt[NIDT];

1111 /*
1112 * Prepare per-CPU pagetables for all processes on the 64 bit kernel.
1113 * Prepare CPU specific pagetables for VLP processes on 64 bit kernels.
1114 * Each CPU has a set of 2 pagetables that are reused for any 32 bit
1115 * process it runs. They are the top level pagetable, hci_pcp_l3ptes, and
1116 * the next to top level table for the bottom 512 Gig, hci_pcp_l2ptes.
1117 * process it runs. They are the top level pagetable, hci_vlp_l3ptes, and
1118 * the next to top level table for the bottom 512 Gig, hci_vlp_l2ptes.
1119 */
1120 /* ARGSUSED */
1121 static void
1122 hat_pcp_setup(struct cpu *cpu)
1123 {
1124     hat_vlp_setup(struct cpu *cpu)
1125 {
1126     #if !defined(__xpv)
1127     #if defined(__amd64) && !defined(__xpv)
1128     struct hat_cpu_info *hci = cpu->cpu_hat_info;
1129     uintptr_t va;
1130     size_t len;
1131     pfn_t pfn;
1132     /*
1133     * allocate the level==2 page table for the bottom most
1134     * 512Gig of address space (this is where 32 bit apps live)
1135     */
1136     ASSERT(hci != NULL);
1137     hci->hci_pcp_l2ptes = kmem_zalloc(MMU_PAGESIZE, KM_SLEEP);
1138     hci->hci_vlp_l2ptes = kmem_zalloc(MMU_PAGESIZE, KM_SLEEP);
1139     /*
1140     * Allocate a top level pagetable and copy the kernel's
1141     * entries into it. Then link in hci_pcp_l2ptes in the 1st entry.
1142     * entries into it. Then link in hci_vlp_l2ptes in the 1st entry.
1143     */
1144     hci->hci_pcp_l3ptes = kmem_zalloc(MMU_PAGESIZE, KM_SLEEP);
1145     hci->hci_pcp_l3pfn =
1146     hat_getpfnnum(kas.a_hat, (caddr_t)hci->hci_pcp_l3ptes);
1147     ASSERT3U(hci->hci_pcp_l3pfn, !=, PFN_INVALID);
1148     bcopy(pcp_page, hci->hci_pcp_l3ptes, MMU_PAGESIZE);
1149     hci->hci_vlp_l3ptes = kmem_zalloc(MMU_PAGESIZE, KM_SLEEP);
1150     hci->hci_vlp_pfn =
1151     hat_getpfnnum(kas.a_hat, (caddr_t)hci->hci_vlp_l3ptes);
1152     ASSERT(hci->hci_vlp_pfn != PFN_INVALID);
1153     bcopy(vlp_page, hci->hci_vlp_l3ptes, MMU_PAGESIZE);

```

```

1149 hci->hci_pcp_l2pfn =
1150     hat_getpfnum(kas.a_hat, (caddr_t)hci->hci_pcp_l2ptes);
1151 ASSERT3U(hci->hci_pcp_l2pfn, !=, PFN_INVALID);

1153 /*
1154  * Now go through and allocate the user version of these structures.
1155  * Unlike with the kernel version, we allocate a hat to represent the
1156  * top-level page table as that will make it much simpler when we need
1157  * to patch through user entries.
1158  */
1159 hci->hci_user_hat = hat_cpu_alloc(cpu);
1160 hci->hci_user_l3pfn = hci->hci_user_hat->hat_htable->ht_pfn;
1161 ASSERT3U(hci->hci_user_l3pfn, !=, PFN_INVALID);
1162 hci->hci_user_l3ptes =
1163     (x86pte_t *)hat_kpm_mapin_pfn(hci->hci_user_l3pfn);

1165 /* Skip the rest of this if KPTI is switched off at boot. */
1166 if (kpti_enable != 1)
1167     return;

1169 /*
1170  * OK, now that we have this we need to go through and punch the normal
1171  * holes in the CPU's hat for this. At this point we'll punch in the
1172  * following:
1173  *
1174  *   o GDT
1175  *   o IDT
1176  *   o LDT
1177  *   o Trampoline Code
1178  *   o machcpu KPTI page
1179  *   o kmdb ISR code page (just trampolines)
1180  *
1181  * If this is cpu0, then we also can initialize the following because
1182  * they'll have already been allocated.
1183  *
1184  *   o TSS for CPU 0
1185  *   o Double Fault for CPU 0
1186  *
1187  * The following items have yet to be allocated and have not been
1188  * punched in yet. They will be punched in later:
1189  *
1190  *   o TSS (mach_cpucontext_alloc_tables())
1191  *   o Double Fault Stack (mach_cpucontext_alloc_tables())
1192  */
1193 hati_cpu_punchin(cpu, (uintptr_t)cpu->cpu_gdt, PROT_READ);
1194 hati_cpu_punchin(cpu, (uintptr_t)cpu->cpu_idt, PROT_READ);

1196 /*
1197  * As the KDI IDT is only active during kmdb sessions (including single
1198  * stepping), typically we don't actually need this punched in (we
1199  * consider the routines that switch to the user cr3 to be toxic). But
1200  * if we ever accidentally end up on the user cr3 while on this IDT,
1201  * we'd prefer not to triple fault.
1202  */
1203 hati_cpu_punchin(cpu, (uintptr_t)&kdi_idt, PROT_READ);

1205 CTASSERT(((uintptr_t)&kpti_trampoline_start % MMU_PAGESIZE) == 0);
1206 CTASSERT(((uintptr_t)&kpti_trampoline_end % MMU_PAGESIZE) == 0);
1207 for (va = (uintptr_t)&kpti_trampoline_start;
1208      va < (uintptr_t)&kpti_trampoline_end; va += MMU_PAGESIZE) {
1209     hati_cpu_punchin(cpu, va, PROT_READ | PROT_EXEC);
1210 }

1212 VERIFY3U(((uintptr_t)cpu->cpu_m.mcpu_ldt) % MMU_PAGESIZE, ==, 0);
1213 for (va = (uintptr_t)cpu->cpu_m.mcpu_ldt, len = LDT_CPU_SIZE;
1214      len >= MMU_PAGESIZE; va += MMU_PAGESIZE, len -= MMU_PAGESIZE) {

```

```

1215     hati_cpu_punchin(cpu, va, PROT_READ);
1216 }

1218 /* mcpu_pad2 is the start of the page containing the kpti_frames. */
1219 hati_cpu_punchin(cpu, (uintptr_t)&cpu->cpu_m.mcpu_pad2[0],
1220     PROT_READ | PROT_WRITE);

1222 if (cpu == &cpus[0]) {
1223     /*
1224     * CPU0 uses a global for its double fault stack to deal with
1225     * the chicken and egg problem. We need to punch it into its
1226     * user HAT.
1227     */
1228     extern char dblfault_stack0[];

1230     hati_cpu_punchin(cpu, (uintptr_t)cpu->cpu_m.mcpu_tss,
1231         PROT_READ);

1233     for (va = (uintptr_t)dblfault_stack0,
1234          len = DEFAULTSTKSZ; len >= MMU_PAGESIZE;
1235          va += MMU_PAGESIZE, len -= MMU_PAGESIZE) {
1236         hati_cpu_punchin(cpu, va, PROT_READ | PROT_WRITE);
1237     }
1238 }

1240 CTASSERT(((uintptr_t)&kdi_isr_start % MMU_PAGESIZE) == 0);
1241 CTASSERT(((uintptr_t)&kdi_isr_end % MMU_PAGESIZE) == 0);
1242 for (va = (uintptr_t)&kdi_isr_start;
1243      va < (uintptr_t)&kdi_isr_end; va += MMU_PAGESIZE) {
1244     hati_cpu_punchin(cpu, va, PROT_READ | PROT_EXEC);
1245 }
1246 #endif /* !_xpv */
1247     pfn = hat_getpfnum(kas.a_hat, (caddr_t)hci->hci_vlp_l2ptes);
1248     ASSERT(pfn != PFN_INVALID);
1249     hci->hci_vlp_l3ptes[0] = MAKEPTP(pfn, 2);
1250 #endif /* !_amd64 && !_xpv */
1251 }

1249 /*ARGSUSED*/
1250 static void
1251 hat_pcp_takedown(cpu_t *cpu)
1252 {
1253     #if !defined(__xpv)
1254     #if defined(__amd64) && !defined(__xpv)
1255         struct hat_cpu_info *hci;

1256         if ((hci = cpu->cpu_hat_info) == NULL)
1257             return;
1258         if (hci->hci_pcp_l2ptes != NULL)
1259             kmem_free(hci->hci_pcp_l2ptes, MMU_PAGESIZE);
1260         if (hci->hci_pcp_l3ptes != NULL)
1261             kmem_free(hci->hci_pcp_l3ptes, MMU_PAGESIZE);
1262         if (hci->hci_user_hat != NULL) {
1263             hat_free_start(hci->hci_user_hat);
1264             hat_free_end(hci->hci_user_hat);
1265         }
1266         if (hci->hci_vlp_l2ptes)
1267             kmem_free(hci->hci_vlp_l2ptes, MMU_PAGESIZE);
1268         if (hci->hci_vlp_l3ptes)
1269             kmem_free(hci->hci_vlp_l3ptes, MMU_PAGESIZE);
1270     #endif
1271     #endif
1272 }

1276 /*

```

```

1277 * Finish filling in the kernel hat.
1278 * Pre fill in all top level kernel page table entries for the kernel's
1279 * part of the address range. From this point on we can't use any new
1280 * kernel large pages if they need PTE's at max_level
1281 *
1282 * create the kmap mappings.
1283 */
1284 void
1285 hat_init_finish(void)
1286 {
1287     size_t     size;
1288     uint_t     r = 0;
1289     uintptr_t  va;
1290     hat_kernel_range_t *rp;

1293 /*
1294  * We are now effectively running on the kernel hat.
1295  * Clearing use_boot_reserve shuts off using the pre-allocated boot
1296  * reserve for all HAT allocations. From here on, the reserves are
1297  * only used when avoiding recursion in kmem_alloc().
1298  */
1299     use_boot_reserve = 0;
1300     htable_adjust_reserve();

1302 /*
1303  * User HATS are initialized with copies of all kernel mappings in
1304  * higher level page tables. Ensure that those entries exist.
1305  */
1306 #if defined(__amd64)

1308     NEXT_HKR(r, 3, kernelbase, 0);
1309 #if defined(__xpv)
1310     NEXT_HKR(r, 3, HYPERVISOR_VIRT_START, HYPERVISOR_VIRT_END);
1311 #endif

1313 #elif defined(__i386)

1315 #if !defined(__xpv)
1316     if (mmu.pae_hat) {
1317         va = kernelbase;
1318         if ((va & LEVEL_MASK(2)) != va) {
1319             va = P2ROUNDUP(va, LEVEL_SIZE(2));
1320             NEXT_HKR(r, 1, kernelbase, va);
1321         }
1322         if (va != 0)
1323             NEXT_HKR(r, 2, va, 0);
1324     } else
1325 #endif /* __xpv */
1326     NEXT_HKR(r, 1, kernelbase, 0);

1328 #endif /* __i386 */

1330     num_kernel_ranges = r;

1332 /*
1333  * Create all the kernel pagetables that will have entries
1334  * shared to user HATs.
1335  */
1336     for (r = 0; r < num_kernel_ranges; ++r) {
1337         rp = &kernel_ranges[r];
1338         for (va = rp->hkr_start_va; va != rp->hkr_end_va;
1339              va += LEVEL_SIZE(rp->hkr_level)) {
1340             htable_t *ht;

1342             if (IN_HYPERVISOR_VA(va))

```

```

1343             continue;

1345             /* can/must skip if a page mapping already exists */
1346             if (rp->hkr_level <= mmu.max_page_level &&
1347                 (ht = htable_getpage(kas.a_hat, va, NULL)) !=
1348                 NULL) {
1349                 htable_release(ht);
1350                 continue;
1351             }

1353             (void) htable_create(kas.a_hat, va, rp->hkr_level - 1,
1354                                 NULL);
1355         }
1356     }

1358 /*
1359  * 32 bit PAE metal kernels use only 4 of the 512 entries in the
1360  * page holding the top level pagetable. We use the remainder for
1361  * the "per CPU" page tables for PCP processes.
1362  * the "per CPU" page tables for VLP processes.
1363  * Map the top level kernel pagetable into the kernel to make
1364  * it easy to use bcopy access these tables.
1365  * PAE is required for the 64-bit kernel which uses this as well to
1366  * perform the per-CPU pagetables. See the big theory statement.
1367  */
1368     if (mmu.pae_hat) {
1369         pcp_page = vmem_alloc(heap_arena, MMU_PAGESIZE, VM_SLEEP);
1370         hat_devload(kas.a_hat, (caddr_t)pcp_page, MMU_PAGESIZE,
1371                    vlp_page = vmem_alloc(heap_arena, MMU_PAGESIZE, VM_SLEEP);
1372                    hat_devload(kas.a_hat, (caddr_t)vlp_page, MMU_PAGESIZE,
1373                               kas.a_hat->hat_htable->ht_pfn,
1374                               PROT_WRITE |
1375                               PROT_READ | HAT_NOSYNC | HAT_UNORDERED_OK,
1376                               HAT_LOAD | HAT_LOAD_NOCONSIST);
1377     }
1378     hat_pcp_setup(CPU);
1379     hat_vlp_setup(CPU);

1380 /*
1381  * Create kmap (cached mappings of kernel PTEs)
1382  * for 32 bit we map from segmap_start .. ekernelheap
1383  * for 64 bit we map from segmap_start .. segmap_start + segmapsize;
1384  */
1385 #if defined(__i386)
1386     size = (uintptr_t)ekernelheap - segmap_start;
1387 #elif defined(__amd64)
1388     size = segmapsize;
1389 #endif
1390     hat_kmap_init((uintptr_t)segmap_start, size);

1392 #if !defined(__xpv)
1393     ASSERT3U(kas.a_hat->hat_htable->ht_pfn, !=, PFN_INVALID);
1394     ASSERT3U(kpti_safe_cr3, ==,
1395             MAKECR3(kas.a_hat->hat_htable->ht_pfn, PCID_KERNEL));
1396 #endif
1397 }

1399 /*
1400  * On 32 bit PAE mode, PTE's are 64 bits, but ordinary atomic memory references
1401  * are 32 bit, so for safety we must use atomic_cas_64() to install these.
1402  */
1403 #ifdef __i386
1404     static void

```

```

1405 reload_pae32(hat_t *hat, cpu_t *cpu)
1406 {
1407     x86pte_t *src;
1408     x86pte_t *dest;
1409     x86pte_t pte;
1410     int i;
1411
1412     /*
1413      * Load the 4 entries of the level 2 page table into this
1414      * cpu's range of the pcp page and point cr3 at them.
1415      * cpu's range of the vlp page and point cr3 at them.
1416      */
1417     ASSERT(mmu.pae_hat);
1418     src = hat->hat_copied_ptes;
1419     dest = pcp_page + (cpu->cpu_id + 1) * MAX_COPIED_PTES;
1420     for (i = 0; i < MAX_COPIED_PTES; ++i) {
1421         src = hat->hat_vlp_ptes;
1422         dest = vlp_page + (cpu->cpu_id + 1) * VLP_NUM_PTES;
1423         for (i = 0; i < VLP_NUM_PTES; ++i) {
1424             for (;;) {
1425                 pte = dest[i];
1426                 if (pte == src[i])
1427                     break;
1428                 if (atomic_cas_64(dest + i, pte, src[i]) != src[i])
1429                     break;
1430             }
1431         }
1432     }
1433 #endif
1434
1435 /*
1436 * Update the PCP data on the CPU cpu to the one on the hat. If this is a 32-bit
1437 * process, then we must update the L2 pages and then the L3. If this is a
1438 * 64-bit process then we must update the L3 entries.
1439 */
1440 static void
1441 hat_pcp_update(cpu_t *cpu, const hat_t *hat)
1442 {
1443     ASSERT3U(hat->hat_flags & HAT_COPIED, !=, 0);
1444
1445     if ((hat->hat_flags & HAT_COPIED_32) != 0) {
1446         const x86pte_t *l2src;
1447         x86pte_t *l2dst, *l3ptes, *l3uptes;
1448         /*
1449          * This is a 32-bit process. To set this up, we need to do the
1450          * following:
1451          *
1452          * - Copy the 4 L2 PTEs into the dedicated L2 table
1453          * - Zero the user L3 PTEs in the user and kernel page table
1454          * - Set the first L3 PTE to point to the CPU L2 table
1455          */
1456         l2src = hat->hat_copied_ptes;
1457         l2dst = cpu->cpu_hat_info->hci_pcp_l2ptes;
1458         l3ptes = cpu->cpu_hat_info->hci_pcp_l3ptes;
1459         l3uptes = cpu->cpu_hat_info->hci_user_l3ptes;
1460
1461         l2dst[0] = l2src[0];
1462         l2dst[1] = l2src[1];
1463         l2dst[2] = l2src[2];
1464         l2dst[3] = l2src[3];
1465
1466         /*
1467          * Make sure to use the mmu to get the number of slots. The
1468          * number of PLP entries that this has will always be less as
1469          * it's a 32-bit process.
1470          */

```

```

1467         bzero(l3ptes, sizeof (x86pte_t) * mmu.top_level_uslots);
1468         l3ptes[0] = MAKEPTP(cpu->cpu_hat_info->hci_pcp_l2pfn, 2);
1469         bzero(l3uptes, sizeof (x86pte_t) * mmu.top_level_uslots);
1470         l3uptes[0] = MAKEPTP(cpu->cpu_hat_info->hci_pcp_l2pfn, 2);
1471     } else {
1472         /*
1473          * This is a 64-bit process. To set this up, we need to do the
1474          * following:
1475          *
1476          * - Zero the 4 L2 PTEs in the CPU structure for safety
1477          * - Copy over the new user L3 PTEs into the kernel page table
1478          * - Copy over the new user L3 PTEs into the user page table
1479          */
1480         ASSERT3S(kpti_enable, ==, 1);
1481         bzero(cpu->cpu_hat_info->hci_pcp_l2ptes, sizeof (x86pte_t) * 4);
1482         bcopy(hat->hat_copied_ptes, cpu->cpu_hat_info->hci_pcp_l3ptes,
1483             sizeof (x86pte_t) * mmu.top_level_uslots);
1484         bcopy(hat->hat_copied_ptes, cpu->cpu_hat_info->hci_user_l3ptes,
1485             sizeof (x86pte_t) * mmu.top_level_uslots);
1486     }
1487 }
1488
1489 static void
1490 reset_kpti(struct kpti_frame *fr, uint64_t kcr3, uint64_t ucr3)
1491 {
1492     ASSERT3U(fr->kf_tr_flag, ==, 0);
1493 #if DEBUG
1494     if (fr->kf_kernel_cr3 != 0) {
1495         ASSERT3U(fr->kf_lower_redzone, ==, 0xdeadbeefdeadbeef);
1496         ASSERT3U(fr->kf_middle_redzone, ==, 0xdeadbeefdeadbeef);
1497         ASSERT3U(fr->kf_upper_redzone, ==, 0xdeadbeefdeadbeef);
1498     }
1499 #endif
1500
1501     bzero(fr, offsetof(struct kpti_frame, kf_kernel_cr3));
1502     bzero(&fr->kf_unused, sizeof (struct kpti_frame) -
1503         offsetof(struct kpti_frame, kf_unused));
1504
1505     fr->kf_kernel_cr3 = kcr3;
1506     fr->kf_user_cr3 = ucr3;
1507     fr->kf_tr_ret_rsp = (uintptr_t)&fr->kf_tr_rsp;
1508
1509     fr->kf_lower_redzone = 0xdeadbeefdeadbeef;
1510     fr->kf_middle_redzone = 0xdeadbeefdeadbeef;
1511     fr->kf_upper_redzone = 0xdeadbeefdeadbeef;
1512 }
1513
1514 #ifdef __xpv
1515 static void
1516 hat_switch_xen(hat_t *hat)
1517 {
1518     struct mmuext_op t[2];
1519     uint_t retcnt;
1520     uint_t opcnt = 1;
1521     uint64_t newcr3;
1522
1523     ASSERT(!(hat->hat_flags & HAT_COPIED));
1524     ASSERT(!(getcr4() & CR4_PCIDE));
1525
1526     newcr3 = MAKECR3((uint64_t)hat->hat_htable->ht_pfn, PCID_NONE);
1527
1528     t[0].cmd = MMUEXT_NEW_BASEPTR;
1529     t[0].arg1.mfn = mmu_btop(pa_to_ma(newcr3));
1530
1531     /*
1532      * There's an interesting problem here, as to what to actually specify

```

```

1533     * when switching to the kernel hat. For now we'll reuse the kernel hat
1534     * again.
1535     */
1536     t[1].cmd = MMUEXT_NEW_USER_BASEPTR;
1537     if (hat == kas.a_hat)
1538         t[1].arg1.mfn = mmu_btop(pa_to_ma(newcr3));
1539     else
1540         t[1].arg1.mfn = pfn_to_mfn(hat->hat_user_ptable);
1541     ++opcnt;

1543     if (HYPERVISOR_mmuext_op(t, opcnt, &retcnt, DOMID_SELF) < 0)
1544         panic("HYPERVISOR_mmu_update() failed");
1545     ASSERT(retcnt == opcnt);
1546 }
1547 #endif /* __xpv */

1549 /*
1550  * Switch to a new active hat, maintaining bit masks to track active CPUs.
1551  *
1552  * With KPTI, all our HATs except kas should be using PCP. Thus, to switch
1553  * HATs, we need to copy over the new user PTEs, then set our trampoline context
1554  * as appropriate.
1555  *
1556  * If lacking PCID, we then load our new cr3, which will flush the TLB: we may
1557  * have established userspace TLB entries via kernel accesses, and these are no
1558  * longer valid. We have to do this eagerly, as we just deleted this CPU from
1559  * ->hat_cpus, so would no longer see any TLB shootdowns.
1560  *
1561  * With PCID enabled, things get a little more complicated. We would like to
1562  * keep TLB context around when entering and exiting the kernel, and to do this,
1563  * we partition the TLB into two different spaces:
1564  *
1565  * PCID_KERNEL is defined as zero, and used both by kas and all other address
1566  * spaces while in the kernel (post-trampoline).
1567  *
1568  * PCID_USER is used while in userspace. Therefore, userspace cannot use any
1569  * lingering PCID_KERNEL entries to kernel addresses it should not be able to
1570  * read.
1571  *
1572  * The trampoline cr3s are set not to invalidate on a mov to %cr3. This means if
1573  * we take a journey through the kernel without switching HATs, we have some
1574  * hope of keeping our TLB state around.
1575  *
1576  * On a hat switch, rather than deal with any necessary flushes on the way out
1577  * of the trampolines, we do them upfront here. If we're switching from kas, we
1578  * shouldn't need any invalidation.
1579  *
1580  * Otherwise, we can have stale userspace entries for both PCID_USER (what
1581  * happened before we move onto the kcr3) and PCID_KERNEL (any subsequent
1582  * userspace accesses such as ddi_copyin()). Since setcr3() won't do these
1583  * flushes on its own in PCIDE, we'll do a non-flushing load and then
1584  * invalidate everything.
1585  * On the 32-bit PAE hypervisor, %cr3 is a 64-bit value, on metal it
1586  * remains a 32-bit value.
1587  */
1588 void
1589 hat_switch(hat_t *hat)
1590 {
1591     uint64_t newcr3;
1592     cpu_t *cpu = CPU;
1593     hat_t *old = cpu->cpu_current_hat;

1594     /*
1595      * set up this information first, so we don't miss any cross calls
1596      */
1597     if (old != NULL) {

```

```

1596         if (old == hat)
1597             return;
1598         if (old != kas.a_hat)
1599             CPUSET_ATOMIC_DEL(old->hat_cpus, cpu->cpu_id);
1600     }

1602     /*
1603      * Add this CPU to the active set for this HAT.
1604      */
1605     if (hat != kas.a_hat) {
1606         CPUSET_ATOMIC_ADD(hat->hat_cpus, cpu->cpu_id);
1607     }
1608     cpu->cpu_current_hat = hat;

1610 #if defined(__xpv)
1611     hat_switch_xen(hat);
1612 #else
1613     struct hat_cpu_info *info = cpu->cpu_m.mcpu_hat_info;
1614     uint64_t pcide = getcr4() & CR4_PCIDE;
1615     uint64_t kcr3, ucr3;
1616     pfn_t tl_kpfn;
1617     ulong_t flag;
1618     /*
1619      * now go ahead and load cr3
1620      */
1621     if (hat->hat_flags & HAT_VLP) {
1622 #if defined(__amd64)
1623         x86pte_t *vlpptep = cpu->cpu_hat_info->hci_vlp_l2ptes;
1624 #endif

1625     EQUIV(kpti_enable, !mmu.pt_global);

1626     if (hat->hat_flags & HAT_COPIED) {
1627         hat_pcp_update(cpu, hat);
1628         tl_kpfn = info->hci_pcp_l3pfn;
1629         VLP_COPY(hat->hat_vlp_ptes, vlpptep);
1630         newcr3 = MAKECR3(cpu->cpu_hat_info->hci_vlp_pfn);
1631 #elif defined(__i386)
1632         reload_pae32(hat, cpu);
1633         newcr3 = MAKECR3(kas.a_hat->hat_htable->ht_pfn) +
1634             (cpu->cpu_id + 1) * VLP_SIZE;
1635 #endif
1636     } else {
1637         IMPLY(kpti_enable, hat == kas.a_hat);
1638         tl_kpfn = hat->hat_htable->ht_pfn;
1639         newcr3 = MAKECR3((uint64_t)hat->hat_htable->ht_pfn);
1640     }
1641 #ifdef __xpv
1642     {
1643         struct mmuext_op t[2];
1644         uint_t retcnt;
1645         uint_t opcnt = 1;

1646         if (pcide) {
1647             ASSERT(kpti_enable);

1648             kcr3 = MAKECR3(tl_kpfn, PCID_KERNEL) | CR3_NOINVL_BIT;
1649             ucr3 = MAKECR3(info->hci_user_l3pfn, PCID_USER) |
1650                 CR3_NOINVL_BIT;

1651             setcr3(kcr3);
1652             if (old != kas.a_hat)
1653                 mmu_flush_tlb(FLUSH_TLB_ALL, NULL);
1654         } else {
1655             kcr3 = MAKECR3(tl_kpfn, PCID_NONE);
1656             ucr3 = kpti_enable ?
1657                 MAKECR3(info->hci_user_l3pfn, PCID_NONE) :

```

```

1643         0;
1644     }
1645     setcr3(kcr3);
1646 }

1032     t[0].cmd = MMUEXT_NEW_BASEPTR;
1033     t[0].arg1.mfn = mmu_btop(pa_to_ma(newcr3));
1034 #if defined(__amd64)
1648     /*
1649     * We will already be taking shutdowns for our new HAT, and as KPTI
1650     * invpcid emulation needs to use kf_user_cr3, make sure we don't get
1651     * any cross calls while we're inconsistent. Note that it's harmless to
1652     * have a *stale* kf_user_cr3 (we just did a FLUSH_TLB_ALL), but a
1653     * *zero* kf_user_cr3 is not going to go very well.
1654     * There's an interesting problem here, as to what to
1655     * actually specify when switching to the kernel hat.
1656     * For now we'll reuse the kernel hat again.
1657     */
1658     if (pcide)
1659     {
1660         flag = intr_clear();
1661         t[1].cmd = MMUEXT_NEW_USER_BASEPTR;
1662         if (hat == kas.a_hat)
1663             t[1].arg1.mfn = mmu_btop(pa_to_ma(newcr3));
1664         else
1665             t[1].arg1.mfn = pfn_to_mfn(hat->hat_user_ptable);
1666         ++opcnt;
1667     }
1668 #endif /* __amd64 */
1669     if (HYPERVISOR_mmuext_op(t, opcnt, &retcnt, DOMID_SELF) < 0)
1670         panic("HYPERVISOR_mmu_update() failed");
1671     ASSERT(retcnt == opcnt);

1672     reset_kpti(&cpu->cpu_m.mcpu_kpti, kcr3, ucr3);
1673     reset_kpti(&cpu->cpu_m.mcpu_kptiflt, kcr3, ucr3);
1674     reset_kpti(&cpu->cpu_m.mcpu_kpti_dbg, kcr3, ucr3);

1675     if (pcide)
1676         intr_restore(flag);

1677 #endif /* !__xpv */

1678 }
1679 #else
1680     setcr3(newcr3);
1681 #endif
1682 ASSERT(cpu == CPU);
1683 }
1684 #endif /* unchanged_portion_omitted */

1894 /*
1895 * This the set of PTE bits for PFN, permissions and caching
1896 * that are allowed to change on a HAT_LOAD_REMAP
1897 */
1898 #define PT_REMAP_BITS \
1899     (PT_PADDR | PT_NX | PT_WRITABLE | PT_WRITETHRU | \
1900     PT_NOCACHE | PT_PAT_4K | PT_PAT_LARGE | PT_IGNORE | PT_REF | PT_MOD)

1901 #define REMAPASSERT(EX) if (!(EX)) panic("hati_pte_map: " #EX)
1902 /*
1903 * Do the low-level work to get a mapping entered into a HAT's pagetables
1904 * and in the mapping list of the associated page_t.
1905 */
1906 static int
1907 hati_pte_map(
1908     htable_t     *ht,
1909     uint_t       entry,
1910     page_t       *pp,

```

```

1912     x86pte_t     pte,
1913     int          flags,
1914     void         *pte_ptr)
1915 {
1916     hat_t        *hat = ht->ht_hat;
1917     x86pte_t     old_pte;
1918     level_t      l = ht->ht_level;
1919     hment_t      *hm;
1920     uint_t       is_consist;
1921     uint_t       is_locked;
1922     int          rv = 0;

1923     /*
1924     * Is this a consistent (ie. need mapping list lock) mapping?
1925     */
1926     is_consist = (pp != NULL && (flags & HAT_LOAD_NOCONSIST) == 0);

1927     /*
1928     * Track locked mapping count in the htable. Do this first,
1929     * as we track locking even if there already is a mapping present.
1930     */
1931     is_locked = (flags & HAT_LOAD_LOCK) != 0 && hat != kas.a_hat;
1932     if (is_locked)
1933         HTABLE_LOCK_INC(ht);

1934     /*
1935     * Acquire the page's mapping list lock and get an hment to use.
1936     * Note that hment_prepare() might return NULL.
1937     */
1938     if (is_consist) {
1939         x86_hm_enter(pp);
1940         hm = hment_prepare(ht, entry, pp);
1941     }

1942     /*
1943     * Set the new pte, retrieving the old one at the same time.
1944     */
1945     old_pte = x86pte_set(ht, entry, pte, pte_ptr);

1946     /*
1947     * Did we get a large page / page table collision?
1948     */
1949     if (old_pte == LPAGE_ERROR) {
1950         if (is_locked)
1951             HTABLE_LOCK_DEC(ht);
1952         rv = -1;
1953         goto done;
1954     }

1955     /*
1956     * If the mapping didn't change there is nothing more to do.
1957     */
1958     if (PTE_EQUIV(pte, old_pte))
1959         goto done;

1960     /*
1961     * Install a new mapping in the page's mapping list
1962     */
1963     if (!PTE_ISVALID(old_pte)) {
1964         if (is_consist) {
1965             hment_assign(ht, entry, pp, hm);
1966             x86_hm_exit(pp);
1967         } else {
1968             ASSERT(flags & HAT_LOAD_NOCONSIST);
1969         }
1970     }
1971 #if defined(__amd64)

```

```

1978     if (ht->ht_flags & HTABLE_COPIED) {
1366     if (ht->ht_flags & HTABLE_VLP) {
1979         cpu_t *cpu = CPU;
1980         hat_pcp_update(cpu, hat);
1368         x86pte_t *vlpptep = cpu->cpu_hat_info->hci_vlp_l2ptes;
1369         VLP_COPY(hat->hat_vlp_ptes, vlpptep);
1981     }
1982 #endif
1983     HTABLE_INC(ht->ht_valid_cnt);
1984     PGCNT_INC(hat, 1);
1985     return (rv);
1986 }

1988 /*
1989  * Remap's are more complicated:
1990  * - HAT_LOAD_REMAP must be specified if changing the pfn.
1991  * - We also require that NOCONSIST be specified.
1992  * - Otherwise only permission or caching bits may change.
1993  */
1994 if (!PTE_ISPAGE(old_pte, 1))
1995     panic("non-null/page mapping pte=" FMT_PTE, old_pte);

1997 if (PTE2PFN(old_pte, 1) != PTE2PFN(pte, 1)) {
1998     REMAPASSERT(flags & HAT_LOAD_REMAP);
1999     REMAPASSERT(flags & HAT_LOAD_NOCONSIST);
2000     REMAPASSERT(PTE_GET(old_pte, PT_SOFTWARE) >= PT_NOCONSIST);
2001     REMAPASSERT(pf_is_memory(PTE2PFN(old_pte, 1)) ==
2002                 pf_is_memory(PTE2PFN(pte, 1)));
2003     REMAPASSERT(!is_consist);
2004 }

2006 /*
2007  * We only let remaps change the certain bits in the PTE.
2008  */
2009 if (PTE_GET(old_pte, ~PT_REMAP_BITS) != PTE_GET(pte, ~PT_REMAP_BITS))
2010     panic("remap bits changed: old_pte="FMT_PTE", pte="FMT_PTE"\n",
2011           old_pte, pte);

2013 /*
2014  * We don't create any mapping list entries on a remap, so release
2015  * any allocated hment after we drop the mapping list lock.
2016  */
2017 done:
2018 if (is_consist) {
2019     x86_hm_exit(pp);
2020     if (hm != NULL)
2021         hment_free(hm);
2022 }
2023 return (rv);
2024 }

2026 /*
2027  * Internal routine to load a single page table entry. This only fails if
2028  * we attempt to overwrite a page table link with a large page.
2029  */
2030 static int
2031 hati_load_common(
2032     hat_t      *hat,
2033     uintptr_t  va,
2034     page_t     *pp,
2035     uint_t     attr,
2036     uint_t     flags,
2037     level_t    level,
2038     pfn_t      pfn)
2039 {
2040     htable_t    *ht;

```

```

2041     uint_t      entry;
2042     x86pte_t    pte;
2043     int         rv = 0;

2045     /*
2046     * The number 16 is arbitrary and here to catch a recursion problem
2047     * early before we blow out the kernel stack.
2048     */
2049     ++curthread->t_hatdepth;
2050     ASSERT(curthread->t_hatdepth < 16);

2052     ASSERT(hat == kas.a_hat || (hat->hat_flags & HAT_PCP) != 0 ||
2053           AS_LOCK_HELD(hat->hat_as));
1441     ASSERT(hat == kas.a_hat || AS_LOCK_HELD(hat->hat_as));

2055     if (flags & HAT_LOAD_SHARE)
2056         hat->hat_flags |= HAT_SHARED;

2058     /*
2059     * Find the page table that maps this page if it already exists.
2060     */
2061     ht = htable_lookup(hat, va, level);

2063     /*
2064     * We must have HAT_LOAD_NOCONSIST if page_t is NULL.
2065     */
2066     if (pp == NULL)
2067         flags |= HAT_LOAD_NOCONSIST;

2069     if (ht == NULL) {
2070         ht = htable_create(hat, va, level, NULL);
2071         ASSERT(ht != NULL);
2072     }
2073     /*
2074     * htable_va2entry checks this condition as well, but it won't include
2075     * much useful info in the panic. So we do it in advance here to include
2076     * all the context.
2077     */
2078     if (ht->ht_vaddr > va || va > HTABLE_LAST_PAGE(ht)) {
2079         panic("hati_load_common: bad htable: va=%p, last page=%p, "
2080               "ht->ht_vaddr=%p, ht->ht_level=%d", (void *)va,
2081               (void *)HTABLE_LAST_PAGE(ht), (void *)ht->ht_vaddr,
2082               (int)ht->ht_level);
2083     }
2084     entry = htable_va2entry(va, ht);

2086     /*
2087     * a bunch of paranoid error checking
2088     */
2089     ASSERT(ht->ht_busy > 0);
1467     if (ht->ht_vaddr > va || va > HTABLE_LAST_PAGE(ht))
1468         panic("hati_load_common: bad htable %p, va %p",
1469               (void *)ht, (void *)va);
2090     ASSERT(ht->ht_level == level);

2092     /*
2093     * construct the new PTE
2094     */
2095     if (hat == kas.a_hat)
2096         attr &= ~PROT_USER;
2097     pte = hati_mkpte(pfn, attr, level, flags);
2098     if (hat == kas.a_hat && va >= kernelbase)
2099         PTE_SET(pte, mmu.pt_global);

2101     /*
2102     * establish the mapping

```



```

2103      */
2104      rv = hati_pte_map(ht, entry, pp, pte, flags, NULL);

2106      /*
2107      * release the htable and any reserves
2108      */
2109      htable_release(ht);
2110      --curthread->t_hatdepth;
2111      return (rv);
2112 }

    unchanged_portion_omitted

2537 #if !defined(__xpv)
2538 /*
2539 * Cross call service routine to demap a range of virtual
2540 * pages on the current CPU or flush all mappings in TLB.
2541 * Cross call service routine to demap a virtual page on
2542 * the current CPU or flush all mappings in TLB.
2543 */
2544 static int
2545 hati_demap_func(xc_arg_t a1, xc_arg_t a2, xc_arg_t a3)
2546 {
2547     _NOTE(ARGUNUSED(a3));
2548     hat_t      *hat = (hat_t *)a1;
2549     tlb_range_t *range = (tlb_range_t *)a2;
2550     caddr_t     addr = (caddr_t)a2;
2551     size_t      len = (size_t)a3;

2552     /*
2553     * If the target hat isn't the kernel and this CPU isn't operating
2554     * in the target hat, we can ignore the cross call.
2555     */
2556     if (hat != kas.a_hat && hat != CPU->cpu_current_hat)
2557         return (0);

2558     if (range->tr_va != DEMAP_ALL_ADDR) {
2559         mmu_flush_tlb(FLUSH_TLB_RANGE, range);
2560     }
2561     /*
2562     * For a normal address, we flush a range of contiguous mappings
2563     */
2564     if ((uintptr_t)addr != DEMAP_ALL_ADDR) {
2565         for (size_t i = 0; i < len; i += MMU_PAGESIZE)
2566             mmu_tlbflush_entry(addr + i);
2567     }
2568     return (0);
2569 }

2570 /*
2571 * We are flushing all of userspace.
2572 * Otherwise we reload cr3 to effect a complete TLB flush.
2573 */
2574 /*
2575 * When using PCP, we first need to update this CPU's idea of the PCP
2576 * PTEs.
2577 * A reload of cr3 on a VLP process also means we must also recopy in
2578 * the pte values from the struct hat
2579 */
2580 if (hat->hat_flags & HAT_COPIED) {
2581     if (hat->hat_flags & HAT_VLP) {
2582 #if defined(__amd64)
2583         hat_pcp_update(CPU, hat);
2584         x86pte_t *vlpptes = CPU->cpu_hat_info->hci_vlp_l2ptes;

2585         VLP_COPY(hat->hat_vlp_ptes, vlpptes);
2586 #endif
2587     }
2588     #elif defined(__i386)
2589         reload_pae32(hat, CPU);
2590     #endif
2591 }

```

```

2573     }

2574     mmu_flush_tlb(FLUSH_TLB_NONGLOBAL, NULL);
2575     reload_cr3();
2576     return (0);
2577 }

2578 #define TLBIDLE_CPU_HALTED    (0x1UL)
2579 #define TLBIDLE_INVALID_ALL  (0x2UL)
2580 /*
2581 * Flush all TLB entries, including global (ie. kernel) ones.
2582 */
2583 static void
2584 flush_all_tlb_entries(void)
2585 {
2586     ulong_t cr4 = getcr4();

2587     if (cr4 & CR4_PGE) {
2588         setcr4(cr4 & ~(ulong_t)CR4_PGE);
2589         setcr4(cr4);

2590         /*
2591          * 32 bit PAE also needs to always reload_cr3()
2592          */
2593         if (mmu.max_level == 2)
2594             reload_cr3();
2595     } else {
2596         reload_cr3();
2597     }
2598 }

2599 #define TLB_CPU_HALTED    (01ul)
2600 #define TLB_INVALID_ALL  (02ul)
2601 #define CAS_TLB_INFO(cpu, old, new) \
2602     atomic_cas_ulong((ulong_t *)&(cpu)->cpu_m.mcpu_tlb_info, (old), (new))

2603 /*
2604 * Record that a CPU is going idle
2605 */
2606 void
2607 tlb_going_idle(void)
2608 {
2609     atomic_or_ulong((ulong_t *)&CPU->cpu_m.mcpu_tlb_info,
2610         TLBIDLE_CPU_HALTED);
2611     atomic_or_ulong((ulong_t *)&CPU->cpu_m.mcpu_tlb_info, TLB_CPU_HALTED);
2612 }

2613 /*
2614 * Service a delayed TLB flush if coming out of being idle.
2615 * It will be called from cpu idle notification with interrupt disabled.
2616 */
2617 void
2618 tlb_service(void)
2619 {
2620     ulong_t tlb_info;
2621     ulong_t found;

2622     /*
2623     * We only have to do something if coming out of being idle.
2624     */
2625     tlb_info = CPU->cpu_m.mcpu_tlb_info;
2626     if (tlb_info & TLBIDLE_CPU_HALTED) {
2627         if (tlb_info & TLB_CPU_HALTED) {
2628             ASSERT(CPU->cpu_current_hat == kas.a_hat);
2629         }
2630     }
2631 }

```

```

2612         * Atomic clear and fetch of old state.
2613         */
2614         while ((found = CAS_TLB_INFO(CPU, tlb_info, 0)) != tlb_info) {
2615             ASSERT(found & TLBIDLE_CPU_HALTED);
2622             ASSERT(found & TLB_CPU_HALTED);
2616             tlb_info = found;
2617             SMT_PAUSE();
2618         }
2619         if (tlb_info & TLBIDLE_INVAL_ALL)
2620             mmu_flush_tlb(FLUSH_TLB_ALL, NULL);
2626         if (tlb_info & TLB_INVAL_ALL)
2627             flush_all_tlb_entries();
2621     }
2622 }
2623 #endif /* !__xpv */

2625 /*
2626  * Internal routine to do cross calls to invalidate a range of pages on
2627  * all CPUs using a given hat.
2628  */
2629 void
2630 hat_tlb_inval_range(hat_t *hat, tlb_range_t *in_range)
2637 hat_tlb_inval_range(hat_t *hat, uintptr_t va, size_t len)
2631 {
2632     extern int     flushes_require_xcalls; /* from mp_startup.c */
2633     cpuset_t       justme;
2634     cpuset_t       cpus_to_shootdown;
2635     tlb_range_t    range = *in_range;
2636 #ifndef __xpv
2637     cpuset_t       check_cpus;
2638     cpu_t          *cpup;
2639     int            c;
2640 #endif

2642     /*
2643      * If the hat is being destroyed, there are no more users, so
2644      * demap need not do anything.
2645      */
2646     if (hat->hat_flags & HAT_FREEING)
2647         return;

2649     /*
2650      * If demapping from a shared pagetable, we best demap the
2651      * entire set of user TLBs, since we don't know what addresses
2652      * these were shared at.
2653      */
2654     if (hat->hat_flags & HAT_SHARED) {
2655         hat = kas.a_hat;
2656         range.tr_va = DEMAP_ALL_ADDR;
2662         va = DEMAP_ALL_ADDR;
2657     }

2659     /*
2660      * if not running with multiple CPUs, don't use cross calls
2661      */
2662     if (panicstr || !flushes_require_xcalls) {
2663 #ifndef __xpv
2664         if (range.tr_va == DEMAP_ALL_ADDR) {
2670             if (va == DEMAP_ALL_ADDR) {
2665                 xen_flush_tlb();
2666             } else {
2667                 for (size_t i = 0; i < TLB_RANGE_LEN(&range);
2668                     i += MMU_PAGESIZE) {
2669                     xen_flush_va((caddr_t)(range.tr_va + i));
2673                 for (size_t i = 0; i < len; i += MMU_PAGESIZE)
2674                     xen_flush_va((caddr_t)(va + i));

```

```

2670         }
2671     }
2672 #else
2673     (void) hati_demap_func((xc_arg_t)hat, (xc_arg_t)&range, 0);
2673     (void) hati_demap_func((xc_arg_t)hat,
2674                             (xc_arg_t)va, (xc_arg_t)len);
2674 #endif
2675     return;
2676 }

2679     /*
2680      * Determine CPUs to shutdown. Kernel changes always do all CPUs.
2681      * Otherwise it's just CPUs currently executing in this hat.
2682      */
2683     kpreempt_disable();
2684     CPUSSET_ONLY(justme, CPU->cpu_id);
2685     if (hat == kas.a_hat)
2686         cpus_to_shootdown = khat_cpuset;
2687     else
2688         cpus_to_shootdown = hat->hat_cpus;

2690 #ifndef __xpv
2691     /*
2692      * If any CPUs in the set are idle, just request a delayed flush
2693      * and avoid waking them up.
2694      */
2695     check_cpus = cpus_to_shootdown;
2696     for (c = 0; c < NCPU && !CPUSSET_ISNULL(check_cpus); ++c) {
2697         along_t tlb_info;

2699         if (!CPU_IN_SET(check_cpus, c))
2700             continue;
2701         CPUSSET_DEL(check_cpus, c);
2702         cpup = cpu[c];
2703         if (cpup == NULL)
2704             continue;

2706         tlb_info = cpup->cpu_m.mcpu_tlb_info;
2707         while (tlb_info == TLBIDLE_CPU_HALTED) {
2708             (void) CAS_TLB_INFO(cpup, TLBIDLE_CPU_HALTED,
2709                                 TLBIDLE_CPU_HALTED | TLBIDLE_INVAL_ALL);
2712             while (tlb_info == TLB_CPU_HALTED) {
2713                 (void) CAS_TLB_INFO(cpup, TLB_CPU_HALTED,
2714                                     TLB_CPU_HALTED | TLB_INVAL_ALL);
2715                 SMT_PAUSE();
2716                 tlb_info = cpup->cpu_m.mcpu_tlb_info;
2717             }
2718             if (tlb_info == (TLBIDLE_CPU_HALTED | TLBIDLE_INVAL_ALL)) {
2719                 if (tlb_info == (TLB_CPU_HALTED | TLB_INVAL_ALL)) {
2720                     HATSTAT_INC(hs_tlb_inval_delayed);
2721                     CPUSSET_DEL(cpus_to_shootdown, c);
2722                 }
2723             }
2724         }
2725     }
2726 #endif

2720     if (CPUSSET_ISNULL(cpus_to_shootdown) ||
2721         CPUSSET_ISEQUAL(cpus_to_shootdown, justme)) {

2723 #ifndef __xpv
2724         if (range.tr_va == DEMAP_ALL_ADDR) {
2725             if (va == DEMAP_ALL_ADDR) {
2726                 xen_flush_tlb();
2727             } else {
2728                 for (size_t i = 0; i < TLB_RANGE_LEN(&range);

```

```

2729         xen_flush_va((caddr_t)(range.tr_va + i));
2132         for (size_t i = 0; i < len; i += MMU_PAGESIZE)
2133             xen_flush_va((caddr_t)(va + i));
2730     }
2731 }
2732 #else
2733     (void) hati_demap_func((xc_arg_t)hat, (xc_arg_t)&range, 0);
2136     (void) hati_demap_func((xc_arg_t)hat,
2137         (xc_arg_t)va, (xc_arg_t)len);
2734 #endif

2736     } else {

2738         CPUSET_ADD(cpus_to_shutdown, CPU->cpu_id);
2739 #ifdef __xpv
2740         if (range.tr_va == DEMAP_ALL_ADDR) {
2144             if (va == DEMAP_ALL_ADDR) {
2741                 xen_gflush_tlb(cpus_to_shutdown);
2742             } else {
2743                 for (size_t i = 0; i < TLB_RANGE_LEN(&range);
2744                     i += MMU_PAGESIZE) {
2745                     xen_gflush_va((caddr_t)(range.tr_va + i),
2147                         for (size_t i = 0; i < len; i += MMU_PAGESIZE) {
2148                             xen_gflush_va((caddr_t)(va + i),
2746                                 cpus_to_shutdown);
2747                         }
2748                 }
2749 #else
2750                 xc_call((xc_arg_t)hat, (xc_arg_t)&range, 0,
2153                 xc_call((xc_arg_t)hat, (xc_arg_t)va, (xc_arg_t)len,
2751                     CPUSET2BV(cpus_to_shutdown), hati_demap_func);
2752 #endif

2754     }
2755     kpreempt_enable();
2756 }

2758 void
2759 hat_tlb_inval(hat_t *hat, uintptr_t va)
2760 {
2761     /*
2762      * Create range for a single page.
2763      */
2764     tlb_range_t range;
2765     range.tr_va = va;
2766     range.tr_cnt = 1; /* one page */
2767     range.tr_level = MIN_PAGE_LEVEL; /* pages are MMU_PAGESIZE */

2769     hat_tlb_inval_range(hat, &range);
2164     hat_tlb_inval_range(hat, va, MMU_PAGESIZE);
2770 }

```

unchanged portion omitted

```

2935 /*
2331  * Do the callbacks for ranges being unloaded.
2332  */
2333 typedef struct range_info {
2334     uintptr_t     rng_va;
2335     ulong_t      rng_cnt;
2336     level_t      rng_level;
2337 } range_info_t;

2339 /*
2936  * Invalidate the TLB, and perform the callback to the upper level VM system,
2937  * for the specified ranges of contiguous pages.
2938  */

```

```

2939 static void
2940 handle_ranges(hat_t *hat, hat_callback_t *cb, uint_t cnt, tlb_range_t *range)
2344 handle_ranges(hat_t *hat, hat_callback_t *cb, uint_t cnt, range_info_t *range)
2941 {
2942     while (cnt > 0) {
2347         size_t len;

2943         --cnt;
2944         hat_tlb_inval_range(hat, &range[cnt]);
2350         len = range[cnt].rng_cnt << LEVEL_SHIFT(range[cnt].rng_level);
2351         hat_tlb_inval_range(hat, (uintptr_t)range[cnt].rng_va, len);

2946         if (cb != NULL) {
2947             cb->hcb_start_addr = (caddr_t)range[cnt].tr_va;
2354             cb->hcb_start_addr = (caddr_t)range[cnt].rng_va;
2948             cb->hcb_end_addr = cb->hcb_start_addr;
2949             cb->hcb_end_addr += range[cnt].tr_cnt <<
2950                 LEVEL_SHIFT(range[cnt].tr_level);
2356             cb->hcb_end_addr += len;
2951             cb->hcb_function(cb);
2952         }
2953     }
2954 }

2956 /*
2957  * Unload a given range of addresses (has optional callback)
2958  *
2959  * Flags:
2960  * define      HAT_UNLOAD                0x00
2961  * define      HAT_UNLOAD_NOSYNC        0x02
2962  * define      HAT_UNLOAD_UNLOCK       0x04
2963  * define      HAT_UNLOAD_OTHER        0x08 - not used
2964  * define      HAT_UNLOAD_UNMAP       0x10 - same as HAT_UNLOAD
2965  */
2966 #define MAX_UNLOAD_CNT (8)
2967 void
2968 hat_unload_callback(
2969     hat_t      *hat,
2970     caddr_t    addr,
2971     size_t     len,
2972     uint_t     flags,
2973     hat_callback_t *cb)
2974 {
2975     uintptr_t   vaddr = (uintptr_t)addr;
2976     uintptr_t   eaddr = vaddr + len;
2977     httable_t   *ht = NULL;
2978     uint_t      entry;
2979     uintptr_t   contig_va = (uintptr_t)-1L;
2980     tlb_range_t r[MAX_UNLOAD_CNT];
2386     range_info_t r[MAX_UNLOAD_CNT];
2981     uint_t      r_cnt = 0;
2982     x86pte_t    old_pte;

2984     XPV_DISALLOW_MIGRATE();
2985     ASSERT(hat == kas.a_hat || eaddr <= _userlimit);
2986     ASSERT(IS_PAGEALIGNED(vaddr));
2987     ASSERT(IS_PAGEALIGNED(eaddr));

2989     /*
2990     * Special case a single page being unloaded for speed. This happens
2991     * quite frequently, COW faults after a fork() for example.
2992     */
2993     if (cb == NULL && len == MMU_PAGESIZE) {
2994         ht = httable_getpte(hat, vaddr, &entry, &old_pte, 0);
2995         if (ht != NULL) {
2996             if (PTE_ISVALID(old_pte)) {

```

```

2997         hat_pte_unmap(ht, entry, flags, old_pte,
2998                     NULL, B_TRUE);
2999     }
3000     htable_release(ht);
3001 }
3002 XPV_ALLOW_MIGRATE();
3003 return;
3004 }

3006 while (vaddr < eaddr) {
3007     old_pte = htable_walk(hat, &ht, &vaddr, eaddr);
3008     if (ht == NULL)
3009         break;

3011     ASSERT(!IN_VA_HOLE(vaddr));

3013     if (vaddr < (uintptr_t)addr)
3014         panic("hat_unload_callback(): unmap inside large page");

3016     /*
3017      * We'll do the call backs for contiguous ranges
3018      */
3019     if (vaddr != contig_va ||
3020         (r_cnt > 0 && r[r_cnt - 1].tr_level != ht->ht_level)) {
3021         (r_cnt > 0 && r[r_cnt - 1].rng_level != ht->ht_level)) {
3022             if (r_cnt == MAX_UNLOAD_CNT) {
3023                 handle_ranges(hat, cb, r_cnt, r);
3024                 r_cnt = 0;
3025             }
3026             r[r_cnt].tr_va = vaddr;
3027             r[r_cnt].tr_cnt = 0;
3028             r[r_cnt].tr_level = ht->ht_level;
3029             r[r_cnt].rng_va = vaddr;
3030             r[r_cnt].rng_cnt = 0;
3031             r[r_cnt].rng_level = ht->ht_level;
3032             ++r_cnt;
3033         }

3034     /*
3035      * Unload one mapping (for a single page) from the page tables.
3036      * Note that we do not remove the mapping from the TLB yet,
3037      * as indicated by the tlb=FALSE argument to hat_pte_unmap().
3038      * handle_ranges() will clear the TLB entries with one call to
3039      * hat_tlb_inval_range() per contiguous range. This is
3040      * safe because the page can not be reused until the
3041      * callback is made (or we return).
3042      */
3043     entry = htable_va2entry(vaddr, ht);
3044     hat_pte_unmap(ht, entry, flags, old_pte, NULL, B_FALSE);
3045     ASSERT(ht->ht_level <= mmu.max_page_level);
3046     vaddr += LEVEL_SIZE(ht->ht_level);
3047     contig_va = vaddr;
3048     ++r[r_cnt - 1].tr_cnt;
3049     ++r[r_cnt - 1].rng_cnt;
3050 }
3051 if (ht)
3052     htable_release(ht);

3054 /*
3055  * handle last range for callbacks
3056  */
3057 if (r_cnt > 0)
3058     handle_ranges(hat, cb, r_cnt, r);
3059 XPV_ALLOW_MIGRATE();
3060 }

```

```

3058 /*
3059  * Invalidate a virtual address translation on a slave CPU during
3060  * panic() dumps.
3061  */
3062 void
3063 hat_flush_range(hat_t *hat, caddr_t va, size_t size)
3064 {
3065     ssize_t sz;
3066     caddr_t endva = va + size;

3068     while (va < endva) {
3069         sz = hat_getpagesize(hat, va);
3070         if (sz < 0) {
3071             #ifdef __xpv
3072                 xen_flush_tlb();
3073             #else
3074                 mmu_flush_tlb(FLUSH_TLB_ALL, NULL);
3075                 flush_all_tlb_entries();
3076             #endif
3077             break;
3078         }
3079         #ifdef __xpv
3080             xen_flush_va(va);
3081         #else
3082             mmu_flush_tlb_kpage((uintptr_t)va);
3083             mmu_tlbflush_entry(va);
3084         #endif
3085         va += sz;
3086     }
3087     unchanged_portion_omitted

3089 /*
3090  * hat_unshare() is similar to hat_unload_callback(), but
3091  * we have to look for empty shared pagetables. Note that
3092  * hat_unshare() is always invoked against an entire segment.
3093  */
3094 /*ARGSUSED*/
3095 void
3096 hat_unshare(hat_t *hat, caddr_t addr, size_t len, uint_t ismszc)
3097 {
3098     uint64_t    vaddr = (uintptr_t)addr;
3099     uintptr_t    eaddr = vaddr + len;
3100     htable_t    *ht = NULL;
3101     uint_t      need_demaps = 0;
3102     int         flags = HAT_UNLOAD_UNMAP;
3103     level_t     l;

3105     ASSERT(hat != kas.a_hat);
3106     ASSERT(eaddr <= _userlimit);
3107     ASSERT(IS_PAGEALIGNED(vaddr));
3108     ASSERT(IS_PAGEALIGNED(eaddr));
3109     XPV_DISALLOW_MIGRATE();

3111     /*
3112      * First go through and remove any shared pagetables.
3113      *
3114      * Note that it's ok to delay the TLB shutdown till the entire range is
3115      * finished, because if hat_pageunload() were to unload a shared
3116      * pagetable page, its hat_tlb_inval() will do a global TLB invalidate.
3117      */
3118     l = mmu.max_page_level;
3119     if (l == mmu.max_level)
3120         --l;
3121     for (; l >= 0; --l) {

```

```

3722     for (vaddr = (uintptr_t)addr; vaddr < eaddr;
3723          vaddr = (vaddr & LEVEL_MASK(1 + 1)) + LEVEL_SIZE(1 + 1)) {
3724         ASSERT(!IN_VA_HOLE(vaddr));
3725         /*
3726          * find a pagetable that maps the current address
3727          */
3728         ht = htable_lookup(hat, vaddr, 1);
3729         if (ht == NULL)
3730             continue;
3731         if (ht->ht_flags & HTABLE_SHARED_PFN) {
3732             /*
3733              * clear page count, set valid_cnt to 0,
3734              * let htable_release() finish the job
3735              */
3736             hat->hat_ism_pgcnt -= ht->ht_valid_cnt <<
3737                 (LEVEL_SHIFT(ht->ht_level) - MMU_PAGESHIFT);
3738             ht->ht_valid_cnt = 0;
3739             need_demaps = 1;
3740         }
3741         htable_release(ht);
3742     }
3743 }

3745 /*
3746  * flush the TLBs - since we're probably dealing with MANY mappings
3747  * we just do a full invalidation.
3748  * we do just one CR3 reload.
3749  */
3749 if (!(hat->hat_flags & HAT_FREEING) && need_demaps)
3750     hat_tlb_inval(hat, DEMAP_ALL_ADDR);

3752 /*
3753  * Now go back and clean up any unaligned mappings that
3754  * couldn't share pagetables.
3755  */
3756 if (!is_it_dism(hat, addr))
3757     flags |= HAT_UNLOAD_UNLOCK;
3758 hat_unload(hat, addr, len, flags);
3759 XPV_ALLOW_MIGRATE();
3760 }
unchanged_portion_omitted

4502 /*
4503  * Release a CPU private mapping for the given address.
4504  * We decrement the htable valid count so it might be destroyed.
4505  */
4506 /*ARGSUSED1*/
4507 void
4508 hat_memppte_release(caddr_t addr, hat_memppte_t pte_pa)
4509 {
4510     htable_t      *ht;

4512     XPV_DISALLOW_MIGRATE();
4513     /*
4514      * invalidate any left over mapping and decrement the htable valid count
4515      */
4516 #ifdef __xpvm
4517     if (HYPERVISOR_update_va_mapping((uintptr_t)addr, 0,
4518         UVMF_INVLPG | UVMF_LOCAL))
4519         panic("HYPERVISOR_update_va_mapping() failed");
4520 #else
4521     {
4522         x86pte_t *ptepr;

4524         ptepr = x86pte_mapin(mmu_btop(pte_pa),
4525             (pte_pa & MMU_PAGEOFFSET) >> mmu.pte_size_shift, NULL);

```

```

4526         if (mmu.pae_hat)
4527             *ptepr = 0;
4528         else
4529             *(x86pte32_t *)ptepr = 0;
4530         mmu_flush_tlb_kpage((uintptr_t)addr);
4531         mmu_tlbflush_entry(addr);
4532         x86pte_mapout();
4533     }
4534 #endif

4535     ht = htable_getpte(kas.a_hat, ALIGN2PAGE(addr), NULL, NULL, 0);
4536     if (ht == NULL)
4537         panic("hat_memppte_release(): invalid address");
4538     ASSERT(ht->ht_level == 0);
4539     HTABLE_DEC(ht->ht_valid_cnt);
4540     htable_release(ht);
4541     XPV_ALLOW_MIGRATE();
4542 }

4544 /*
4545  * Apply a temporary CPU private mapping to a page. We flush the TLB only
4546  * on this CPU, so this ought to have been called with preemption disabled.
4547  */
4548 void
4549 hat_memppte_remap(
4550     pfn_t          pfn,
4551     caddr_t        addr,
4552     hat_memppte_t pte_pa,
4553     uint_t         attr,
4554     uint_t         flags)
4555 {
4556     uintptr_t      va = (uintptr_t)addr;
4557     x86pte_t       pte;

4559     /*
4560      * Remap the given PTE to the new page's PFN. Invalidate only
4561      * on this CPU.
4562      */
4563 #ifdef DEBUG
4564     htable_t      *ht;
4565     uint_t         entry;

4567     ASSERT(IS_PAGEALIGNED(va));
4568     ASSERT(!IN_VA_HOLE(va));
4569     ht = htable_getpte(kas.a_hat, va, &entry, NULL, 0);
4570     ASSERT(ht != NULL);
4571     ASSERT(ht->ht_level == 0);
4572     ASSERT(ht->ht_valid_cnt > 0);
4573     ASSERT(ht->ht_pfn == mmu_btop(pte_pa));
4574     htable_release(ht);
4575 #endif
4576     XPV_DISALLOW_MIGRATE();
4577     pte = hati_mkpte(pfn, attr, 0, flags);
4578 #ifdef __xpvm
4579     if (HYPERVISOR_update_va_mapping(va, pte, UVMF_INVLPG | UVMF_LOCAL))
4580         panic("HYPERVISOR_update_va_mapping() failed");
4581 #else
4582     {
4583         x86pte_t *ptepr;

4585         ptepr = x86pte_mapin(mmu_btop(pte_pa),
4586             (pte_pa & MMU_PAGEOFFSET) >> mmu.pte_size_shift, NULL);
4587         if (mmu.pae_hat)
4588             *(x86pte_t *)ptepr = pte;
4589         else
4590             *(x86pte32_t *)ptepr = (x86pte32_t)pte;

```

```

4591         mmu_flush_tlb_kpage((uintptr_t)addr);
4592         mmu_tlbFlush_entry(addr);
4593         x86pte_mapout();
4594     }
4595     #endif
4596     XPV_ALLOW_MIGRATE();
4597 }
4598
4599     unchanged_portion_omitted
4600
4601 /*
4602  * HAT part of cpu initialization.
4603 */
4604 void
4605 hat_cpu_online(struct cpu *cpup)
4606 {
4607     if (cpup != CPU) {
4608         x86pte_cpu_init(cpup);
4609         hat_pcp_setup(cpup);
4610         hat_vlp_setup(cpup);
4611     }
4612     CPuset_ATOMIC_ADD(khat_cpuset, cpup->cpu_id);
4613 }
4614
4615 /*
4616  * HAT part of cpu deletion.
4617  * (currently, we only call this after the cpu is safely passivated.)
4618 */
4619 void
4620 hat_cpu_offline(struct cpu *cpup)
4621 {
4622     ASSERT(cpup != CPU);
4623
4624     CPuset_ATOMIC_DEL(khat_cpuset, cpup->cpu_id);
4625     hat_pcp_tear_down(cpup);
4626     hat_vlp_tear_down(cpup);
4627     x86pte_cpu_fini(cpup);
4628 }
4629
4630     unchanged_portion_omitted
4631
4632 #endif /* __xpv */
4633
4634 /*
4635  * Helper function to punch in a mapping that we need with the specified
4636  * attributes.
4637 */
4638 void
4639 hati_cpu_punchin(cpu_t *cpu, uintptr_t va, uint_t attrs)
4640 {
4641     int ret;
4642     pfn_t pfn;
4643     hat_t *cpu_hat = cpu->cpu_hat_info->hci_user_hat;
4644
4645     ASSERT3S(kpti_enable, ==, 1);
4646     ASSERT3P(cpu_hat, !=, NULL);
4647     ASSERT3U(cpu_hat->hat_flags & HAT_PCP, ==, HAT_PCP);
4648     ASSERT3U(va & MMU_PAGEOFFSET, ==, 0);
4649
4650     pfn = hat_getpfnum(kas.a_hat, (caddr_t)va);
4651     VERIFY3U(pfn, !=, PFN_INVALID);
4652
4653     /*
4654      * We purposefully don't try to find the page_t. This means that this
4655      * will be marked PT_NOCONSIST; however, given that this is pretty much
4656      * a static mapping that we're using we should be relatively OK.
4657      */
4658     attrs |= HAT_STORECACHING_OK;
4659     ret = hati_load_common(cpu_hat, va, NULL, attrs, 0, 0, pfn);

```

```

5114         VERIFY3S(ret, ==, 0);
5115     }

```

new/usr/src/uts/i86pc/vm/hat_i86.h

1

```
*****
9504 Fri Apr 6 17:25:05 2018
new/usr/src/uts/i86pc/vm/hat_i86.h
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22 * Copyright 2008 Sun Microsystems, Inc. All rights reserved.
23 * Use is subject to license terms.
24 */
25 /*
26 * Copyright (c) 2014 by Delphix. All rights reserved.
27 * Copyright 2018 Joyent, Inc.
28 */
29
30 #ifndef _VM_HAT_I86_H
31 #define _VM_HAT_I86_H
32
33
34 #ifdef __cplusplus
35 extern "C" {
36 #endif
37
38 /*
39  * VM - Hardware Address Translation management.
40  *
41  * This file describes the contents of the x86_64 HAT data structures.
42  */
43 #include <sys/types.h>
44 #include <sys/t_lock.h>
45 #include <sys/cpuvar.h>
46 #include <sys/x_call.h>
47 #include <vm/seg.h>
48 #include <vm/page.h>
49 #include <sys/vmparam.h>
50 #include <sys/vm_machparam.h>
51 #include <sys/promif.h>
52 #include <vm/hat_pte.h>
53 #include <vm/htable.h>
54 #include <vm/hment.h>
55
56 /*
57  * The essential data types involved:
58  *
59  * htable_t - There is one of these for each page table and it is used
```

new/usr/src/uts/i86pc/vm/hat_i86.h

2

```
60 *
61 * by the HAT to manage the page table.
62 * hment_t - Links together multiple PTEs to a single page.
63 */
64
65 /*
66  * Maximum number of per-CPU pagetable entries that we'll need to cache in the
67  * HAT. See the big theory statement in uts/i86pc/vm/hat_i86.c for more
68  * information.
69  * VLP processes have a 32 bit address range, so their top level is 2 and
70  * with only 4 PTEs in that table.
71 */
72 #if defined(__xpv)
73 * The Xen hypervisor does not use per-CPU pagetables (PCP). Define a single
74 * struct member for it at least to make life easier and not make the member
75 * conditional.
76 #define MAX_COPIED_PTES 1
77 #else
78 /*
79  * The 64-bit kernel may have up to 512 PTEs present in it for a given process.
80  */
81 #define MAX_COPIED_PTES 512
82 #endif /* __xpv */
83 #define VLP_LEVEL (2)
84 #define VLP_NUM_PTES (4)
85 #define VLP_SIZE (VLP_NUM_PTES * sizeof(x86pte_t))
86 #define TOP_LEVEL(h) (((h)->hat_flags & HAT_VLP) ? VLP_LEVEL : mmu.max_level)
87 #define VLP_COPY(fromptep, toptep) { \
88     toptep[0] = fromptep[0]; \
89     toptep[1] = fromptep[1]; \
90     toptep[2] = fromptep[2]; \
91     toptep[3] = fromptep[3]; \
92 }
93
94 #define TOP_LEVEL(h) (((h)->hat_max_level))
95
96 /*
97  * The hat struct exists for each address space.
98  */
99 struct hat {
100     kmutex_t hat_mutex;
101     struct as *hat_as;
102     uint_t hat_stats;
103     pgcnt_t hat_pages_mapped[MAX_PAGE_LEVEL + 1];
104     pgcnt_t hat_ism_pgcnt;
105     cpuset_t hat_cpus;
106     uint16_t hat_flags;
107     uint8_t hat_max_level; /* top level of this HAT */
108     uint_t hat_num_copied; /* Actual num of hat_copied_ptes[] */
109     htable_t *hat_htable; /* top level htable */
110     struct hat *hat_next;
111     struct hat *hat_prev;
112     uint_t hat_num_hash; /* number of htable hash buckets */
113     htable_t **hat_ht_hash; /* htable hash buckets */
114     htable_t *hat_ht_cached; /* cached free htables */
115     x86pte_t hat_copied_ptes[MAX_COPIED_PTES];
116     x86pte_t hat_vlp_ptes[VLP_NUM_PTES];
117 #if defined(__amd64) && defined(__xpv)
118     pfn_t hat_user_ptable; /* alt top ptable for user mode */
119 #endif
120 };
121 typedef struct hat hat_t;
122
123 #define PGCNT_INC(hat, level) \
```

```

113     atomic_inc_ulong(&(hat)->hat_pages_mapped[level]);
114 #define PGCNT_DEC(hat, level) \
115     atomic_dec_ulong(&(hat)->hat_pages_mapped[level]);

117 /*
118  * Flags for the hat_flags field. For more information, please see the big
119  * theory statement on the HAT design in uts/i86pc/vm/hat_i86.c.
120  * Flags for the hat_flags field
121  * HAT_FREEING - set when HAT is being destroyed - mostly used to detect that
122  * demap()s can be avoided.
123  *
124  * HAT_COPIED - Indicates this HAT is a source for per-cpu page tables: see the
125  * big comment in hat_i86.c for a description.
126  * HAT_VLP - indicates a 32 bit process has a virtual address range less than
127  * the hardware's physical address range. (VLP->Virtual Less-than Physical)
128  * Note - never used on the hypervisor.
129  * HAT_COPIED_32 - HAT_COPIED, but for an ILP32 process.
130  *
131  * HAT_VICTIM - This is set while a hat is being examined for page table
132  * stealing and prevents it from being freed.
133  *
134  * HAT_SHARED - The hat has exported it's page tables via hat_share()
135  *
136  * HAT_PINNED - On the hypervisor, indicates the top page table has been pinned.
137  *
138  * HAT_PCP - Used for the per-cpu user page table (i.e. associated with a CPU,
139  * not a process).
140  */
141 #define HAT_FREEING      (0x0001)
142 #define HAT_VICTIM      (0x0002)
143 #define HAT_SHARED      (0x0004)
144 #define HAT_PINNED      (0x0008)
145 #define HAT_COPIED      (0x0010)
146 #define HAT_COPIED_32  (0x0020)
147 #define HAT_PCP         (0x0040)
148 #define HAT_VLP         (0x0002)
149 #define HAT_VICTIM      (0x0004)
150 #define HAT_SHARED      (0x0008)
151 #define HAT_PINNED      (0x0010)

152 /*
153  * Additional platform attribute for hat_devload() to force no caching.
154  */
155 #define HAT_PLAT_NOCACHE      (0x100000)

156 /*
157  * Simple statistics for the HAT. These are just counters that are
158  * atomically incremented. They can be reset directly from the kernel
159  * debugger.
160  */
161 struct hatstats {
162     ulong_t hs_reap_attempts;
163     ulong_t hs_reaped;
164     ulong_t hs_steals;
165     ulong_t hs_ptable_allocs;
166     ulong_t hs_ptable_frees;
167     ulong_t hs_htable_rgets; /* allocs from reserve */
168     ulong_t hs_htable_rputs; /* putbacks to reserve */
169     ulong_t hs_htable_shared; /* number of htables shared */
170     ulong_t hs_htable_unshared; /* number of htables unshared */
171     ulong_t hs_hm_alloc;
172     ulong_t hs_hm_free;
173     ulong_t hs_hm_put_reserve;
174     ulong_t hs_hm_get_reserve;

```

```

171     ulong_t hs_hm_steals;
172     ulong_t hs_hm_steal_exam;
173     ulong_t hs_tlb_inval_delayed;
174     ulong_t hs_hat_copied64;
175     ulong_t hs_hat_copied32;
176     ulong_t hs_hat_normal64;
177 };
178 extern struct hatstats hatstat;
179 #ifdef DEBUG
180 #define HATSTAT_INC(x)  (++hatstat.x)
181 #else
182 #define HATSTAT_INC(x)  (0)
183 #endif

184 #if defined(_KERNEL)

185 /*
186  * Useful macro to align hat_XXX() address arguments to a page boundary
187  */
188 #define ALIGN2PAGE(a)  (((uintptr_t)(a) & MMU_PAGEMASK)
189 #define IS_PAGEALIGNED(a)  (((uintptr_t)(a) & MMU_PAGEOFFSET) == 0)

190 extern uint_t khat_running; /* set at end of hat_kern_setup() */
191 extern cpuset_t khat_cpuset; /* cpuset for kernel address demap Xcalls */
192 extern kmutex_t hat_list_lock;
193 extern kcondvar_t hat_list_cv;

200 /*
201  * Interfaces to setup a cpu private mapping (ie. preemption disabled).
202  * The attr and flags arguments are the same as for hat_devload().
203  * setup() must be called once, then any number of calls to remap(),
204  * followed by a final call to release()
205  *
206  * Used by ppcopy(), page_zero(), the memscrubber, and the kernel debugger.
207  */
208 typedef paddr_t hat_mempte_t; /* phys addr of PTE */
209 extern hat_mempte_t hat_mempte_setup(caddr_t addr);
210 extern void hat_mempte_remap(pfn_t, caddr_t, hat_mempte_t,
211     uint_t attr, uint_t flags);
212 extern void hat_mempte_release(caddr_t addr, hat_mempte_t);

213 /*
214  * Interfaces to manage which thread has access to htable and hment reserves.
215  * The USE_HAT_RESERVES macro should always be recomputed in full. Its value
216  * (due to curthread) can change after any call into kmem/vmem.
217  */
218 #define USE_HAT_RESERVES() \
219     (use_boot_reserve || curthread->t_hatdepth > 1 || \
220     panicstr != NULL || vmem_is_populator())

221 /*
222  * initialization stuff needed by by startup, mp_startup...
223  */
224 extern void hat_cpu_online(struct cpu *);
225 extern void hat_cpu_offline(struct cpu *);
226 extern void setup_vaddr_for_ppcopy(struct cpu *);
227 extern void teardown_vaddr_for_ppcopy(struct cpu *);
228 extern void clear_boot_mappings(uintptr_t, uintptr_t);

229 /*
230  * magic value to indicate that all TLB entries should be demapped.
231  */

```



```

237 #define DEMAP_ALL_ADDR (~(uintptr_t)0)

239 /*
240 * not in any include file???
241 */
242 extern void halt(char *fmt);

244 /*
245 * x86 specific routines for use online in setup or i86pc/vm files
246 */
247 extern void hat_kern_alloc(caddr_t segmap_base, size_t segmap_size,
248     caddr_t ekernelheap);
249 extern void hat_kern_setup(void);
250 extern void hat_tlb_inval(struct hat *hat, uintptr_t va);
251 extern void hat_pte_unmap(htable_t *ht, uint_t entry, uint_t flags,
252     x86pte_t old_pte, void *pte_ptr, boolean_t tlb);
253 extern void hat_init_finish(void);
254 extern caddr_t hat_kpm_pfn2va(pfn_t pfn);
255 extern pfn_t hat_kpm_va2pfn(caddr_t);
256 extern page_t *hat_kpm_vaddr2page(caddr_t);
257 extern uintptr_t hat_kernelbase(uintptr_t);
258 extern void hat_kmap_init(uintptr_t base, size_t len);

259 extern hment_t *hati_page_unmap(page_t *pp, htable_t *ht, uint_t entry);

261 extern void mmu_calc_user_slots(void);
262 extern void hat_tlb_inval(struct hat *hat, uintptr_t va);
263 extern void hat_switch(struct hat *hat);
264 #if !defined(__xpv)
265 /*
266 * routines to deal with delayed TLB invalidations for idle CPUs
267 */
268 extern void tlb_going_idle(void);
269 extern void tlb_service(void);
270 #endif

271 #define TLB_RANGE_LEN(r) ((r)->tr_cnt << LEVEL_SHIFT((r)->tr_level))

272 /*
273 * A range of virtual pages for purposes of demapping.
274 * Hat switch function invoked to load a new context into %cr3
275 */
276 typedef struct tlb_range {
277     uintptr_t tr_va; /* address of page */
278     ulong_t tr_cnt; /* number of pages in range */
279     int8_t tr_level; /* page table level */
280 } tlb_range_t;
281 extern void hat_switch(struct hat *hat);

282 #if defined(__xpv)

283 #define XPV_DISALLOW_MIGRATE() xen_block_migrate()
284 #define XPV_ALLOW_MIGRATE() xen_allow_migrate()

285 #define mmu_flush_tlb_page(va) mmu_invlpg((caddr_t)va)
286 #define mmu_flush_tlb_kpage(va) mmu_invlpg((caddr_t)va)

287 #ifdef __xpv
288 /*
289 * Interfaces to use around code that maps/unmaps grant table references.
290 */
291 extern void hat_prepare_mapping(hat_t *, caddr_t, uint64_t *);
292 extern void hat_release_mapping(hat_t *, caddr_t);

293 #define XPV_DISALLOW_MIGRATE() xen_block_migrate()
294 #define XPV_ALLOW_MIGRATE() xen_allow_migrate()

```

```

290 #else

291 #define XPV_DISALLOW_MIGRATE() /* nothing */
292 #define XPV_ALLOW_MIGRATE() /* nothing */

293 #define pfn_is_foreign(pfn) __lintzero

294 typedef enum flush_tlb_type {
295     FLUSH_TLB_ALL = 1,
296     FLUSH_TLB_NONGLOBAL = 2,
297     FLUSH_TLB_RANGE = 3,
298 } flush_tlb_type_t;
299 #endif

300 extern void mmu_flush_tlb(flush_tlb_type_t, tlb_range_t *);
301 extern void mmu_flush_tlb_kpage(uintptr_t);
302 extern void mmu_flush_tlb_page(uintptr_t);

303 extern void hati_cpu_punchin(cpu_t *cpu, uintptr_t va, uint_t attrs);

304 /*
305 * routines to deal with delayed TLB invalidations for idle CPUs
306 */
307 extern void tlb_going_idle(void);
308 extern void tlb_service(void);

309 #endif /* !__xpv */

310 #endif /* _KERNEL */

311 #ifdef __cplusplus
312 }
313 #endif

```

unchanged portion omitted

new/usr/src/uts/i86pc/vm/hat_kdi.c

1

```
*****
8256 Fri Apr 6 17:25:05 2018
new/usr/src/uts/i86pc/vm/hat_kdi.c
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */

22 /*
23  * Copyright 2007 Sun Microsystems, Inc. All rights reserved.
24  * Use is subject to license terms.
25  *
26  * Copyright 2018 Joyent, Inc.
27  */

29 /*
30  * HAT interfaces used by the kernel debugger to interact with the VM system.
31  * These interfaces are invoked when the world is stopped. As such, no blocking
32  * operations may be performed.
33  */

35 #include <sys/cpuvar.h>
36 #include <sys/kdi_impl.h>
37 #include <sys/errno.h>
38 #include <sys/system.h>
39 #include <sys/sysmacros.h>
40 #include <sys/mman.h>
41 #include <sys/bootconf.h>
42 #include <sys/cmn_err.h>
43 #include <vm/seg_kmem.h>
44 #include <vm/hat_i86.h>
45 #if defined(__xpv)
46 #include <sys/hypervisor.h>
47 #endif
48 #include <sys/bootinfo.h>
49 #include <vm/kboot_mmu.h>
50 #include <sys/machsystem.h>

52 /*
53  * The debugger needs direct access to the PTE of one page table entry
54  * in order to implement vtop and physical read/writes
55  */
56 static uintptr_t hat_kdi_page = 0; /* vaddr for phsical page accesses */
57 static uint_t use_kbm = 1;
58 uint_t hat_kdi_use_pae; /* if 0, use x86pte32_t for pte type */
```

new/usr/src/uts/i86pc/vm/hat_kdi.c

2

```
60 #if !defined(__xpv)
61 static x86pte_t *hat_kdi_pte = NULL; /* vaddr of pte for hat_kdi_page */
62 #endif

64 /*
65  * Get the address for remapping physical pages during boot
66  */
67 void
68 hat_boot_kdi_init(void)
69 {
70     hat_kdi_page = (uintptr_t)kvm_push(0); /* first call gets address... */
71 }

unchanged_portion_omitted

143 #else
144 #define kdi_mtop(m) (m)
145 #define kdi_ptom(p) (p)
146 #endif

148 /*ARGSUSED*/
149 int
150 kdi_vtop(uintptr_t va, uint64_t *pap)
151 {
152     uintptr_t vaddr = va;
153     size_t len;
154     pfn_t pfn;
155     uint_t prot;
156     int level;
157     x86pte_t pte;
158     int index;

160     /*
161      * if the mmu struct isn't relevant yet, we need to probe
162      * the boot loader's pagetables.
163      */
164     if (!khat_running) {
165         if (kvm_probe(&vaddr, &len, &pfn, &prot) == 0)
166             return (ENOENT);
167         if (vaddr > va)
168             return (ENOENT);
169         if (vaddr < va)
170             pfn += mmu_btop(va - vaddr);
171         *pap = pfn_to_pa(pfn) + (vaddr & MMU_PAGEOFFSET);
172         return (0);
173     }

175     /*
176      * We can't go through normal hat routines, so we'll use
177      * kdi_pread() to walk the page tables
178      */
179     #if defined(__xpv)
180     *pap = pfn_to_pa(CPU->cpu_current_hat->hat_htable->ht_pfn);
181     #else
182     *pap = getcr3_pa();
183     *pap = getcr3() & MMU_PAGEMASK;
184     #endif

185     for (level = mmu.max_level; ; --level) {
186         index = (va >> LEVEL_SHIFT(level)) & (mmu.ptes_per_table - 1);
187         *pap += index << mmu.pte_size_shift;
188         pte = 0;
189         if (kdi_pread((caddr_t)&pte, mmu.pte_size, *pap, &len) != 0)
190             return (ENOENT);
191         if (pte == 0)
192             return (ENOENT);
193         if (level > 0 && level <= mmu.max_page_level &&
194             (pte & PT_PAGESIZE)) {
```

```

194         *pap = kdi_mtop(pte & PT_PADDR_LGPG);
195         break;
196     } else {
197         *pap = kdi_mtop(pte & PT_PADDR);
198         if (level == 0)
199             break;
200     }
201 }
202 *pap += va & LEVEL_OFFSET(level);
203 return (0);
204 }

206 static int
207 kdi_prw(caddr_t buf, size_t nbytes, uint64_t pa, size_t *ncopiedp, int doread)
208 {
209     size_t ncopied = 0;
210     off_t  pgoff;
211     size_t sz;
212     caddr_t va;
213     caddr_t from;
214     caddr_t to;
215     x86pte_t pte;

217     /*
218      * if this is called before any initialization - fail
219      */
220     if (hat_kdi_page == 0)
221         return (EAGAIN);

222     while (nbytes > 0) {
223         /*
224          * figure out the addresses and construct a minimal PTE
225          */
226         pgoff = pa & MMU_PAGEOFFSET;
227         sz = MIN(nbytes, MMU_PAGESIZE - pgoff);
228         va = (caddr_t)hat_kdi_page + pgoff;
229         pte = kdi_ptom(mmu_ptob(mmu_btop(pa))) | PT_VALID;
230         if (doread) {
231             if (doread) {
232                 from = va;
233                 to = buf;
234             } else {
235                 PTE_SET(pte, PT_WRITABLE);
236                 from = buf;
237                 to = va;
238             }
239         }

240         /*
241          * map the physical page
242          */
243         if (use_kbm)
244             (void) kbm_push(pa);
245 #if defined(__xpv)
246         else
247             (void) HYPERVISOR_update_va_mapping(
248                 (uintptr_t)va, pte, UVMF_INVLPG);
249 #else
250         else if (hat_kdi_use_pae)
251             *hat_kdi_pte = pte;
252         else
253             *(x86pte32_t *)hat_kdi_pte = pte;
254         mmu_flush_tlb_kpage(hat_kdi_page);
255 #endif
256     }
257     bcopy(from, to, sz);

```

```

259     /*
260      * erase the mapping
261      */
262     if (use_kbm)
263         kbm_pop();
264 #if defined(__xpv)
265     else
266         (void) HYPERVISOR_update_va_mapping(
267             (uintptr_t)va, 0, UVMF_INVLPG);
268 #else
269     else if (hat_kdi_use_pae)
270         *hat_kdi_pte = 0;
271     else
272         *(x86pte32_t *)hat_kdi_pte = 0;
273     mmu_flush_tlb_kpage(hat_kdi_page);
274 #endif
275     mmu_tlbflush_entry((caddr_t)hat_kdi_page);
276     buf += sz;
277     pa += sz;
278     nbytes -= sz;
279     ncopied += sz;
280 }

282     if (ncopied == 0)
283         return (ENOENT);

285     *ncopiedp = ncopied;
286     return (0);
287 }

    unchanged_portion_omitted

301 #if !defined(__xpv)
302 /*
303  * This gets used for flushing the TLB on all the slaves just prior to doing a
304  * kdi_prw(). It's unclear why this was originally done, since kdi_prw() itself
305  * will flush any lingering hat_kdi_page mappings, but let's presume it was a
306  * good idea.
307  */
308 void
309 kdi_flush_caches(void)
310 {
311     mmu_flush_tlb(FLUSH_TLB_ALL, NULL);
312 }
313 #endif

315 /*
316  * Return the number of bytes, relative to the beginning of a given range, that
317  * are non-toxic (can be read from and written to with relative impunity).
318  */
319 /*ARGSUSED*/
320 size_t
321 kdi_range_is_nontoxic(uintptr_t va, size_t sz, int write)
322 {
323     #if defined(__amd64)
324         extern uintptr_t toxic_addr;
325         extern size_t toxic_size;
326     #endif

327     /*
328      * Check 64 bit toxic range.
329      */
330     if (toxic_addr != 0 &&
331         va + sz >= toxic_addr &&
332         va < toxic_addr + toxic_size)
333         return (va < toxic_addr ? toxic_addr - va : 0);

```

```
335     /*
336     * avoid any Virtual Address hole
337     */
338     if (va + sz >= hole_start && va < hole_end)
339         return (va < hole_start ? hole_start - va : 0);
341     return (sz);
343 #elif defined(__i386)
344     extern void *device_arena_contains(void *, size_t, size_t *);
345     uintptr_t v;
347     v = (uintptr_t)device_arena_contains((void *)va, sz, NULL);
348     if (v == 0)
349         return (sz);
350     else if (v <= va)
351         return (0);
352     else
353         return (v - va);
355 #endif /* __i386 */
356 }
unchanged_portion_omitted
```

new/usr/src/uts/i86pc/vm/hat_pte.h

1

```
*****
9324 Fri Apr 6 17:25:05 2018
new/usr/src/uts/i86pc/vm/hat_pte.h
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>
9208 hati_demap_func should take pagesize into account
Reviewed by: Robert Mustacchi <rm@joyent.com>
Reviewed by: Tim Kordas <tim.kordas@joyent.com>
Reviewed by: Yuri Pankov <yuripv@yuripv.net>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[ ]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22  * Copyright 2007 Sun Microsystems, Inc. All rights reserved.
23  * Use is subject to license terms.
24  * Copyright 2018 Joyent, Inc.
25  */

27 #ifndef _VM_HAT_PTE_H
28 #define _VM_HAT_PTE_H

30 #ifdef __cplusplus
31 extern "C" {
32 #endif

34 #include <sys/types.h>
35 #include <sys/mach_mmu.h>

37 /*
38  * macros to get/set/clear the PTE fields
39  */
40 #define PTE_SET(p, f) ((p) |= (f))
41 #define PTE_CLR(p, f) ((p) &= ~(x86pte_t)(f))
42 #define PTE_GET(p, f) ((p) & (f))

44 /*
45  * Handy macro to check if a pagetable entry or pointer is valid
46  */
47 #define PTE_ISVALID(p) PTE_GET(p, PT_VALID)

49 /*
50  * Does a PTE map a large page.
51  */
52 #define PTE_IS_LGPG(p, l) ((l) > 0 && PTE_GET((p), PT_PAGESIZE))

54 /*
55  * does this PTE represent a page (not a pointer to another page table)?
```

new/usr/src/uts/i86pc/vm/hat_pte.h

2

```
56 */
57 #define PTE_ISPAGE(p, l) \
58 (PTE_ISVALID(p) && ((l) == 0 || PTE_GET(p, PT_PAGESIZE)))

60 /*
61  * Handy macro to check if 2 PTE's are the same - ignores REF/MOD bits.
62  * On the 64 bit hypervisor we also have to ignore the high order
63  * software bits and the global/user bit which are set/cleared
64  * capriciously (by the hypervisor!)
65  */
66 #if defined(__amd64) && defined(__xpv)
67 #define PT_IGNORE ((0x7fful << 52) | PT_GLOBAL | PT_USER)
68 #else
69 #define PT_IGNORE (0)
70 #endif
71 #define PTE_EQUIV(a, b) (((a) | (PT_IGNORE | PT_REF | PT_MOD)) == \
72 ((b) | (PT_IGNORE | PT_REF | PT_MOD)))

74 /*
75  * Shorthand for converting a PTE to it's pfn.
76  */
77 #define PTE2MFN(p, l) \
78 mmu_btop(PTE_GET((p), PTE_IS_LGPG((p), (l)) ? PT_PADDR_LGPG : PT_PADDR))
79 #ifdef __xpv
80 #define PTE2PFN(p, l) pte2pfn(p, l)
81 #else
82 #define PTE2PFN(p, l) PTE2MFN(p, l)
83 #endif

85 #define PT_NX (0x8000000000000000ull)
86 #define PT_PADDR (0x000fffffffffff000ull)
87 #define PT_PADDR_LGPG (0x000ffffffffffe000ull) /* phys addr for large pages */

89 /*
90  * Macros to create a PTP or PTE from the pfn and level
91  */
92 #ifdef __xpv

94 /*
95  * we use the highest order bit in physical address pfns to mark foreign mfns
96  */
97 #ifdef _LP64
98 #define PFN_IS_FOREIGN_MFN (1ul << 51)
99 #else
100 #define PFN_IS_FOREIGN_MFN (1ul << 31)
101 #endif

103 #define MAKEPTP(pfn, l) \
104 (pa_to_ma(pfn_to_pa(pfn)) | mmu.ptp_bits[(l) + 1])
105 #define MAKEPTE(pfn, l) \
106 ((pfn & PFN_IS_FOREIGN_MFN) ? \
107 ((pfn_to_pa(pfn) & ~PFN_IS_FOREIGN_MFN) | mmu.pte_bits[l]) | \
108 PT_FOREIGN | PT_REF | PT_MOD) : \
109 (pa_to_ma(pfn_to_pa(pfn)) | mmu.pte_bits[l])
110 #else
111 #define MAKEPTP(pfn, l) \
112 (pfn_to_pa(pfn) | mmu.ptp_bits[(l) + 1])
113 #define MAKEPTE(pfn, l) \
114 (pfn_to_pa(pfn) | mmu.pte_bits[l])
115 #endif

117 /*
118  * The idea of "level" refers to the level where the page table is used in the
119  * the hardware address translation steps. The level values correspond to the
120  * following names of tables used in AMD/Intel architecture documents:
121  */
```

```

122 *      AMD/INTEL name      Level #
123 *      -----
124 *      Page Map Level 4      3
125 *      Page Directory Pointer 2
126 *      Page Directory        1
127 *      Page Table            0
128 *
129 * The numbering scheme is such that the values of 0 and 1 can correspond to
130 * the pagesize codes used for MPSS support. For now the Maximum level at
131 * which you can have a large page is a constant, that may change in
132 * future processors.
133 *
134 * The type of "level_t" is signed so that it can be used like:
135 *      level_t l;
136 *      ...
137 *      while (--l >= 0)
138 *      ...
139 */
140 #define MAX_NUM_LEVEL      4
141 #define MAX_PAGE_LEVEL     2
142 #define MIN_PAGE_LEVEL     0
143 typedef int8_t level_t;
144 #define LEVEL_SHIFT(l)    (mmu.level_shift[l])
145 #define LEVEL_SIZE(l)    (mmu.level_size[l])
146 #define LEVEL_OFFSET(l)  (mmu.level_offset[l])
147 #define LEVEL_MASK(l)    (mmu.level_mask[l])
148
149 /*
150 * Macros to:
151 * Check for a PFN above 4Gig and 64Gig for 32 bit PAE support
152 */
153 #define PFN_4G              (4ull * (1024 * 1024 * 1024 / MMU_PAGESIZE))
154 #define PFN_64G            (64ull * (1024 * 1024 * 1024 / MMU_PAGESIZE))
155 #define PFN_ABOVE4G(pfn) ((pfn) >= PFN_4G)
156 #define PFN_ABOVE64G(pfn) ((pfn) >= PFN_64G)
157
158 /*
159 * The CR3 register holds the physical address of the top level page table,
160 * along with the current PCID if any.
161 * The CR3 register holds the physical address of the top level page table.
162 */
163 #define MAKECR3(pfn, pcid)  (mmu_ptob(pfn) | pcid)
164 #define MAKECR3(pfn)       mmu_ptob(pfn)
165
166 /*
167 * HAT/MMU parameters that depend on kernel mode and/or processor type
168 */
169 struct htable;
170 struct hat_mmu_info {
171     x86pte_t pt_nx;          /* either 0 or PT_NX */
172     x86pte_t pt_global;     /* either 0 or PT_GLOBAL */
173
174     pfn_t highest_pfn;
175
176     uint_t num_level;        /* number of page table levels in use */
177     uint_t max_level;        /* just num_level - 1 */
178     uint_t max_page_level;   /* maximum level at which we can map a page */
179     uint_t umax_page_level;  /* max user page map level */
180     uint_t ptes_per_table;   /* # of entries in lower level page tables */
181     uint_t top_level_count;  /* # of entries in top-level page table */
182     uint_t top_level_uslots; /* # of user slots in top-level page table */
183     uint_t num_copied_ents;  /* # of PCP-copied PTEs to create */
184     /* 32-bit versions of values */
185     uint_t top_level_uslots32;
186     uint_t max_level32;
187     uint_t num_copied_ents32;

```

```

176     uint_t top_level_count; /* # of entries in top most level page table */
177
178     uint_t hash_cnt;        /* cnt of entries in htable_hash_cache */
179     uint_t hat32_hash_cnt; /* cnt of entries in 32-bit htable_hash_cache */
180     uint_t vlp_hash_cnt;    /* cnt of entries in vlp htable_hash_cache */
181
182     uint_t pae_hat;         /* either 0 or 1 */
183
184     uintptr_t hole_start;   /* start of VA hole (or -1 if none) */
185     uintptr_t hole_end;     /* end of VA hole (or 0 if none) */
186
187     struct htable **kmap_htables; /* htables for segmap + 32 bit heap */
188     x86pte_t *kmap_ptes;     /* mapping of pagetables that map kmap */
189     uintptr_t kmap_addr;     /* start addr of kmap */
190     uintptr_t kmap_eaddr;    /* end addr of kmap */
191
192     uint_t pte_size;        /* either 4 or 8 */
193     uint_t pte_size_shift; /* either 2 or 3 */
194     x86pte_t ptp_bits[MAX_NUM_LEVEL]; /* bits set for interior PTP */
195     x86pte_t pte_bits[MAX_NUM_LEVEL]; /* bits set for leaf PTE */
196
197     /*
198     * A range of VA used to window pages in the i86pc/vm code.
199     * See PWIN_XXX macros.
200     */
201     caddr_t pwin_base;
202     caddr_t pwin_pte_va;
203     paddr_t pwin_pte_pa;
204
205     /*
206     * The following tables are equivalent to PAGEXXXXXX at different levels
207     * in the page table hierarchy.
208     */
209     uint_t level_shift[MAX_NUM_LEVEL]; /* PAGESHIFT for given level */
210     uintptr_t level_size[MAX_NUM_LEVEL]; /* PAGESIZE for given level */
211     uintptr_t level_offset[MAX_NUM_LEVEL]; /* PAGEOFFSET for given level */
212     uintptr_t level_mask[MAX_NUM_LEVEL]; /* PAGEMASK for given level */
213 };
214
215 _____unchanged_portion_omitted_____

```

new/usr/src/uts/i86pc/vm/htable.c

1

```
*****
60560 Fri Apr 6 17:25:05 2018
new/usr/src/uts/i86pc/vm/htable.c
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */

22 /*
23  * Copyright (c) 2004, 2010, Oracle and/or its affiliates. All rights reserved.
24  * Copyright (c) 2014 by Delphix. All rights reserved.
25  * Copyright 2018 Joyent, Inc.
26  * Copyright 2015 Joyent, Inc.
27 */

28 #include <sys/types.h>
29 #include <sys/sysmacros.h>
30 #include <sys/kmem.h>
31 #include <sys/atomic.h>
32 #include <sys/bitmap.h>
33 #include <sys/machparam.h>
34 #include <sys/machsystem.h>
35 #include <sys/mman.h>
36 #include <sys/system.h>
37 #include <sys/cpuvar.h>
38 #include <sys/thread.h>
39 #include <sys/proc.h>
40 #include <sys/cpu.h>
41 #include <sys/kmem.h>
42 #include <sys/disp.h>
43 #include <sys/vmem.h>
44 #include <sys/vmsystem.h>
45 #include <sys/promif.h>
46 #include <sys/var.h>
47 #include <sys/x86_archext.h>
48 #include <sys/archsystem.h>
49 #include <sys/bootconf.h>
50 #include <sys/dumphdr.h>
51 #include <vm/seg_kmem.h>
52 #include <vm/seg_kpm.h>
53 #include <vm/hat.h>
54 #include <vm/hat_i86.h>
55 #include <sys/cmn_err.h>
56 #include <sys/panic.h>

58 #ifdef __xpv
```

new/usr/src/uts/i86pc/vm/htable.c

2

```
59 #include <sys/hypervisor.h>
60 #include <sys/xpv_panic.h>
61 #endif

63 #include <sys/bootinfo.h>
64 #include <vm/kboot_mmu.h>

66 static void x86pte_zero(htable_t *dest, uint_t entry, uint_t count);

68 kmem_cache_t *htable_cache;

70 /*
71  * The variable htable_reserve_amount, rather than HTABLE_RESERVE_AMOUNT,
72  * is used in order to facilitate testing of the htable_steal() code.
73  * By resetting htable_reserve_amount to a lower value, we can force
74  * stealing to occur. The reserve amount is a guess to get us through boot.
75 */
76 #define HTABLE_RESERVE_AMOUNT (200)
77 uint_t htable_reserve_amount = HTABLE_RESERVE_AMOUNT;
78 kmutex_t htable_reserve_mutex;
79 uint_t htable_reserve_cnt;
80 htable_t *htable_reserve_pool;

82 /*
83  * Used to hand test htable_steal().
84 */
85 #ifdef DEBUG
86 ulong_t force_steal = 0;
87 ulong_t ptable_cnt = 0;
88 #endif

90 /*
91  * This variable is so that we can tune this via /etc/system
92  * Any value works, but a power of two <= mmu.ptes_per_table is best.
93 */
94 uint_t htable_steal_passes = 8;

96 /*
97  * mutex stuff for access to htable hash
98 */
99 #define NUM_HTABLE_MUTEX 128
100 kmutex_t htable_mutex[NUM_HTABLE_MUTEX];
101 #define HTABLE_MUTEX_HASH(h) ((h) & (NUM_HTABLE_MUTEX - 1))

103 #define HTABLE_ENTER(h) mutex_enter(&htable_mutex[HTABLE_MUTEX_HASH(h)]);
104 #define HTABLE_EXIT(h) mutex_exit(&htable_mutex[HTABLE_MUTEX_HASH(h)]);

106 /*
107  * forward declarations
108 */
109 static void link_ptp(htable_t *higher, htable_t *new, uintptr_t vaddr);
110 static void unlink_ptp(htable_t *higher, htable_t *old, uintptr_t vaddr);
111 static void htable_free(htable_t *ht);
112 static x86pte_t *x86pte_access_pagetable(htable_t *ht, uint_t index);
113 static void x86pte_release_pagetable(htable_t *ht);
114 static x86pte_t x86pte_cas(htable_t *ht, uint_t entry, x86pte_t old,
115 x86pte_t new);

117 /*
118  * A counter to track if we are stealing or reaping htables. When non-zero
119  * htable_free() will directly free htables (either to the reserve or kmem)
120  * instead of putting them in a hat's htable cache.
121 */
122 uint32_t htable_dont_cache = 0;

124 /*
```

```

125 * Track the number of active pagetables, so we can know how many to reap
126 */
127 static uint32_t active_ptables = 0;

129 #ifdef __xpv
130 /*
131  * Deal with hypervisor complications.
132  */
133 void
134 xen_flush_va(caddr_t va)
135 {
136     struct mmuext_op t;
137     uint_t count;

139     if (IN_XPV_PANIC()) {
140         mmu_flush_tlb_page((uintptr_t)va);
140         mmu_tlbflush_entry((caddr_t)va);
141     } else {
142         t.cmd = MMUEXT_INVLPG_LOCAL;
143         t.arg1.linear_addr = (uintptr_t)va;
144         if (HYPERVISOR_mmuext_op(&t, 1, &count, DOMID_SELF) < 0)
145             panic("HYPERVISOR_mmuext_op() failed");
146         ASSERT(count == 1);
147     }
148 }

150 void
151 xen_gflush_va(caddr_t va, cpuset_t cpus)
152 {
153     struct mmuext_op t;
154     uint_t count;

156     if (IN_XPV_PANIC()) {
157         mmu_flush_tlb_page((uintptr_t)va);
157         mmu_tlbflush_entry((caddr_t)va);
158         return;
159     }

161     t.cmd = MMUEXT_INVLPG_MULTI;
162     t.arg1.linear_addr = (uintptr_t)va;
163     /*LINTED: constant in conditional context*/
164     set_xen_guest_handle(t.arg2.vcpumask, &cpus);
165     if (HYPERVISOR_mmuext_op(&t, 1, &count, DOMID_SELF) < 0)
166         panic("HYPERVISOR_mmuext_op() failed");
167     ASSERT(count == 1);
168 }

```

unchanged portion omitted

```

563 /*
564 * This routine steals htables from user processes.  Called by htable_reap
565 * (reap=TRUE) or htable_alloc (reap=FALSE).
566 */
567 static htable_t *
568 htable_steal(uint_t cnt, boolean_t reap)
569 {
570     hat_t      *hat = kas.a_hat;      /* list starts with khat */
571     htable_t    *list = NULL;
572     htable_t    *ht;
573     uint_t      stolen = 0;
574     uint_t      pass, passes;
575     uint_t      threshold;

577     /*
578      * Limit htable_steal_passes to something reasonable
579      */
580     if (htable_steal_passes == 0)

```

```

581         htable_steal_passes = 1;
582     if (htable_steal_passes > mmu.ptes_per_table)
583         htable_steal_passes = mmu.ptes_per_table;

585     /*
586      * If we're stealing merely as part of kmem reaping (versus stealing
587      * to assure forward progress), we don't want to actually steal any
588      * active htables.  (Stealing active htables merely to give memory
589      * back to the system can inadvertently kick off an htable crime wave
590      * as active processes repeatedly steal htables from one another,
591      * plummeting the system into a kind of HAT lawlessness that can
592      * become so violent as to impede the one thing that can end it: the
593      * freeing of memory via ARC reclaim and other means.)  So if we're
594      * reaping, we limit ourselves to the first pass that steals cached
595      * htables that aren't in use -- which gives memory back, but averts
596      * the entire breakdown of social order.
597      */
598     passes = reap ? 0 : htable_steal_passes;

600     /*
601      * Loop through all user hats.  The 1st pass takes cached htables that
602      * aren't in use.  The later passes steal by removing mappings, too.
603      */
604     atomic_inc_32(&htable_dont_cache);
605     for (pass = 0; pass <= passes && stolen < cnt; ++pass) {
606         threshold = pass * mmu.ptes_per_table / htable_steal_passes;

608         mutex_enter(&hat_list_lock);

610         /* skip the first hat (kernel) */
611         hat = kas.a_hat->hat_next;
612         for (;;) {
613             /*
614              * Skip any hat that is already being stolen from.
615              *
616              * We skip SHARED hats, as these are dummy
617              * hats that host ISM shared page tables.
618              *
619              * We also skip if HAT_FREEING because hat_pte_unmap()
620              * won't zero out the PTE's.  That would lead to hitting
621              * stale PTEs either here or under hat_unload() when we
622              * steal and unload the same page table in competing
623              * threads.
624              *
625              * We skip HATs that belong to CPUs, to make our lives
626              * simpler.
627              */
628             while (hat != NULL && (hat->hat_flags &
629                 (HAT_VICTIM | HAT_SHARED | HAT_FREEING |
630                 HAT_PCP)) != 0) {
631                 while (hat != NULL &&
632                     (hat->hat_flags &
633                     (HAT_VICTIM | HAT_SHARED | HAT_FREEING)) != 0)
634                     hat = hat->hat_next;
635             }

637             if (hat == NULL)
638                 break;

640             /*
641              * Mark the HAT as a stealing victim so that it is
642              * not freed from under us, e.g. in as_free()
643              */
644             hat->hat_flags |= HAT_VICTIM;
645             mutex_exit(&hat_list_lock);

```



```

644      /*
645       * Take any htables from the hat's cached "free" list.
646       */
647      hat_enter(hat);
648      while ((ht = hat->hat_ht_cached) != NULL &&
649             stolen < cnt) {
650          hat->hat_ht_cached = ht->ht_next;
651          ht->ht_next = list;
652          list = ht;
653          ++stolen;
654      }
655      hat_exit(hat);

657      /*
658       * Don't steal active htables on first pass.
659       */
660      if (pass != 0 && (stolen < cnt))
661          htable_steal_active(hat, cnt, threshold,
662                             &stolen, &list);

664      /*
665       * do synchronous teardown for the reap case so that
666       * we can forget hat; at this time, hat is
667       * guaranteed to be around because HAT_VICTIM is set
668       * (see htable_free() for similar code)
669       */
670      for (ht = list; (ht) && (reap); ht = ht->ht_next) {
671          if (ht->ht_hat == NULL)
672              continue;
673          ASSERT(ht->ht_hat == hat);
674 #if defined(__xpv) && defined(__amd64)
675          ASSERT(!(ht->ht_flags & HTABLE_COPIED));
676          if (ht->ht_level == mmu.max_level) {
677              if (!(ht->ht_flags & HTABLE_VLP) &&
678                  ht->ht_level == mmu.max_level) {
679                  ptable_free(hat->hat_user_ptable);
680                  hat->hat_user_ptable = PFN_INVALID;
681              }
682 #endif
683          /*
684           * forget the hat
685           */
686          ht->ht_hat = NULL;
687      }

689      mutex_enter(&hat_list_lock);

691      /*
692       * Are we finished?
693       */
694      if (stolen == cnt) {
695          /*
696           * Try to spread the pain of stealing,
697           * move victim HAT to the end of the HAT list.
698           */
699          if (pass >= 1 && cnt == 1 &&
700              kas.a_hat->hat_prev != hat)
701              move_victim(hat);
702          /*
703           * We are finished
704           */
705      }

707      /*
708       * Clear the victim flag, hat can go away now (once
709       * the lock is dropped)

```

```

708      /*
709       * if (hat->hat_flags & HAT_VICTIM) {
710           ASSERT(hat != kas.a_hat);
711           hat->hat_flags &= ~HAT_VICTIM;
712           cv_broadcast(&hat_list_cv);
713       }

715       /* move on to the next hat */
716       hat = hat->hat_next;
717   }

719       mutex_exit(&hat_list_lock);

721   }
722   ASSERT(!MUTEX_HELD(&hat_list_lock));

724   atomic_dec_32(&htable_dont_cache);
725   return (list);
726 }

unchanged_portion_omitted

727 /*
728  * Allocate an htable, stealing one or using the reserve if necessary
729  */
730 static htable_t *
731 htable_alloc(
732     hat_t          *hat,
733     uintptr_t      vaddr,
734     level_t        level,
735     htable_t       *shared)
736 {
737     htable_t      *ht = NULL;
738     uint_t         is_copied;
739     uint_t         is_vlp;
740     uint_t         is_bare = 0;
741     uint_t         need_to_zero = 1;
742     int            kmflags = (can_steal_post_boot ? KM_NOSLEEP : KM_SLEEP);

743     if (level < 0 || level > TOP_LEVEL(hat))
744         panic("htable_alloc(): level %d out of range\n", level);

745     is_copied = (hat->hat_flags & HAT_COPIED) &&
746                level == hat->hat_max_level;
747     if (is_copied || shared != NULL)
748         is_vlp = (hat->hat_flags & HAT_VLP) && level == VLP_LEVEL;
749     if (is_vlp || shared != NULL)
750         is_bare = 1;

751     /*
752      * First reuse a cached htable from the hat_ht_cached field, this
753      * avoids unnecessary trips through kmem/page allocators.
754      */
755     if (hat->hat_ht_cached != NULL && !is_bare) {
756         hat_enter(hat);
757         ht = hat->hat_ht_cached;
758         if (ht != NULL) {
759             hat->hat_ht_cached = ht->ht_next;
760             need_to_zero = 0;
761             /* XX64 ASSERT() they're all zero somehow */
762             ASSERT(ht->ht_pfn != PFN_INVALID);
763         }
764         hat_exit(hat);
765     }

767     if (ht == NULL) {
768         /*

```

```

817     * Allocate an htable, possibly refilling the reserves.
818     */
819     if (USE_HAT_RESERVES()) {
820         ht = htable_get_reserve();
821     } else {
822         /*
823          * Donate successful htable allocations to the reserve.
824          */
825         for (;;) {
826             ht = kmem_cache_alloc(htable_cache, kmflags);
827             if (ht == NULL)
828                 break;
829             ht->ht_pfn = PFN_INVALID;
830             if (USE_HAT_RESERVES() ||
831                 htable_reserve_cnt >= htable_reserve_amount)
832                 break;
833             htable_put_reserve(ht);
834         }
835     }

837     /*
838     * allocate a page for the hardware page table if needed
839     */
840     if (ht != NULL && !is_bare) {
841         ht->ht_hat = hat;
842         ht->ht_pfn = ptable_alloc((uintptr_t)ht);
843         if (ht->ht_pfn == PFN_INVALID) {
844             if (USE_HAT_RESERVES())
845                 htable_put_reserve(ht);
846             else
847                 kmem_cache_free(htable_cache, ht);
848             ht = NULL;
849         }
850     }
851 }

853     /*
854     * If allocations failed, kick off a kmem_reap() and resort to
855     * htable steal(). We may spin here if the system is very low on
856     * memory. If the kernel itself has consumed all memory and kmem_reap()
857     * can't free up anything, then we'll really get stuck here.
858     * That should only happen in a system where the administrator has
859     * misconfigured VM parameters via /etc/system.
860     */
861     while (ht == NULL && can_steal_post_boot) {
862         kmem_reap();
863         ht = htable_steal(1, B_FALSE);
864         HATSTAT_INC(hs_steals);

866         /*
867         * If we stole for a bare htable, release the pagetable page.
868         */
869         if (ht != NULL) {
870             if (is_bare) {
871                 ptable_free(ht->ht_pfn);
872                 ht->ht_pfn = PFN_INVALID;
873             }
874             /*
875             * make stolen page table writable again in kpm
876             */
877             } else if (kpm_vbase && xen_kpm_page(ht->ht_pfn,
878                 PT_VALID | PT_WRITABLE) < 0) {
879                 panic("failure making kpm r/w pfn=0x%x",
880                     ht->ht_pfn);
881             }
882         }

```

```

883     }
884 }

886     /*
887     * All attempts to allocate or steal failed. This should only happen
888     * if we run out of memory during boot, due perhaps to a huge
889     * boot_archive. At this point there's no way to continue.
890     */
891     if (ht == NULL)
892         panic("htable_alloc(): couldn't steal\n");

894 #if defined(__amd64) && defined(__xpv)
895     /*
896     * Under the 64-bit hypervisor, we have 2 top level page tables.
897     * If this allocation fails, we'll resort to stealing.
898     * We use the stolen page indirectly, by freeing the
899     * stolen htable first.
900     */
901     if (level == mmu.max_level) {
902         for (;;) {
903             htable_t *stolen;

905             hat->hat_user_ptable = ptable_alloc((uintptr_t)ht + 1);
906             if (hat->hat_user_ptable != PFN_INVALID)
907                 break;
908             stolen = htable_steal(1, B_FALSE);
909             if (stolen == NULL)
910                 panic("2nd steal ptable failed\n");
911             htable_free(stolen);
912         }
913         block_zero_no_xmm(kpm_vbase + pfn_to_pa(hat->hat_user_ptable),
914             MMU_PAGESIZE);
915     }
916 #endif

918     /*
919     * Shared page tables have all entries locked and entries may not
920     * be added or deleted.
921     */
922     ht->ht_flags = 0;
923     if (shared != NULL) {
924         ASSERT(shared->ht_valid_cnt > 0);
925         ht->ht_flags |= HTABLE_SHARED_PFN;
926         ht->ht_pfn = shared->ht_pfn;
927         ht->ht_lock_cnt = 0;
928         ht->ht_valid_cnt = 0;           /* updated in hat_share() */
929         ht->ht_shares = shared;
930         need_to_zero = 0;
931     } else {
932         ht->ht_shares = NULL;
933         ht->ht_lock_cnt = 0;
934         ht->ht_valid_cnt = 0;
935     }

937     /*
938     * setup flags, etc. for copied page tables.
939     * setup flags, etc. for VLP htables
940     */
941     if (is_copied) {
942         ht->ht_flags |= HTABLE_COPIED;
943     }
944     if (is_vlp) {
945         ht->ht_flags |= HTABLE_VLP;
946         ASSERT(ht->ht_pfn == PFN_INVALID);
947         need_to_zero = 0;
948     }

```



```

1131     HTABLE_DEC(higher->ht_valid_cnt);
1132 }

1134 /*
1135 * Link an entry for a new table at vaddr and level into the existing table
1136 * one level higher. We are always holding the HASH_ENTER() when doing this.
1137 */
1138 static void
1139 link_ptp(htable_t *higher, htable_t *new, uintptr_t vaddr)
1140 {
1141     uint_t      entry = htable_va2entry(vaddr, higher);
1142     x86pte_t    newptp = MAKEPTP(new->ht_pfn, new->ht_level);
1143     x86pte_t    found;

1145     ASSERT(higher->ht_busy > 0);

1147     ASSERT(new->ht_level != mmu.max_level);

1149     HTABLE_INC(higher->ht_valid_cnt);

1151     found = x86pte_cas(higher, entry, 0, newptp);
1152     if ((found & ~PT_REF) != 0)
1153         panic("HAT: ptp not 0, found=" FMT_PTE, found);

1155     /*
1156     * When a top level PTE changes for a copied htable, we must trigger a
1157     * hat_pcp_update() on all HAT CPUs.
1158     *
1159     * When any top level VLP page table entry changes, we must issue
1160     * a reload of cr3 on all processors using it.
1161     * We also need to do this for the kernel hat on PAE 32 bit kernel.
1162     */
1163     if (
1164     #ifdef __i386
1165         (higher->ht_hat == kas.a_hat &&
1166          higher->ht_level == higher->ht_hat->hat_max_level) ||
1167         (higher->ht_hat == kas.a_hat && higher->ht_level == VLP_LEVEL) ||
1168     #endif
1169         (higher->ht_flags & HTABLE_COPIED)
1170         (higher->ht_flags & HTABLE_VLP))
1171         hat_tlb_inval(higher->ht_hat, DEMAP_ALL_ADDR);
1172 }
1173
1174     unchanged portion omitted

1285 /*
1286 * Find the htable for the pagetable at the given level for the given address.
1287 * If found acquires a hold that eventually needs to be htable_release()d
1288 */
1289 htable_t *
1290 htable_lookup(hat_t *hat, uintptr_t vaddr, level_t level)
1291 {
1292     uintptr_t    base;
1293     uint_t      hashval;
1294     htable_t    *ht = NULL;

1296     ASSERT(level >= 0);
1297     ASSERT(level <= TOP_LEVEL(hat));

1299     if (level == TOP_LEVEL(hat)) {
1300     #if defined(__amd64)
1301         /*
1302         * 32 bit address spaces on 64 bit kernels need to check
1303         * for overflow of the 32 bit address space
1304         */
1305         if ((hat->ht_flags & HAT_COPIED_32) &&
1306             vaddr >= ((uint64_t)1 << 32))

```

```

1298     if ((hat->ht_flags & HAT_VLP) && vaddr >= ((uint64_t)1 << 32))
1299         return (NULL);
1300 #endif
1301     base = 0;
1302 } else {
1303     base = vaddr & LEVEL_MASK(level + 1);
1304 }

1314     hashval = HTABLE_HASH(hat, base, level);
1315     HTABLE_ENTER(hashval);
1316     for (ht = hat->hat_ht_hash[hashval]; ht; ht = ht->ht_next) {
1317         if (ht->ht_hat == hat &&
1318             ht->ht_vaddr == base &&
1319             ht->ht_level == level)
1320             break;
1321     }
1322     if (ht)
1323         ++ht->ht_busy;

1325     HTABLE_EXIT(hashval);
1326     return (ht);
1327 }
1328
1329     unchanged portion omitted
1330 #endif /* __i386 */

1345 /*
1346 * Disable preemption and establish a mapping to the pagetable with the
1347 * given pfn. This is optimized for there case where it's the same
1348 * pfn as we last used referenced from this CPU.
1349 */
1350 static x86pte_t *
1351 x86pte_access_pagetable(htable_t *ht, uint_t index)
1352 {
1353     /*
1354     * HTABLE_COPIED pagetables are contained in the hat_t
1355     * VLP pagetables are contained in the hat_t
1356     */
1357     if (ht->ht_flags & HTABLE_COPIED) {
1358         ASSERT3U(index, <, ht->ht_hat->hat_num_copied);
1359         return (PT_INDEX_PTR(ht->ht_hat->hat_copied_ptes, index));
1360     }
1361     if (ht->ht_flags & HTABLE_VLP)
1362         return (PT_INDEX_PTR(ht->ht_hat->hat_vlp_ptes, index));
1363     return (x86pte_mapin(ht->ht_pfn, index, ht));
1364 }

1365 /*
1366 * map the given pfn into the page table window.
1367 */
1368 /* ARGSUSED */
1369 x86pte_t *
1370 x86pte_mapin(pfn_t pfn, uint_t index, htable_t *ht)
1371 {
1372     x86pte_t *pteptr;
1373     x86pte_t pte = 0;
1374     x86pte_t newpte;
1375     int x;

1376     ASSERT(pfn != PFN_INVALID);

1377     if (!khat_running) {
1378         caddr_t va = kbm_remap_window(pfn_to_pa(pfn), 1);
1379         return (PT_INDEX_PTR(va, index));
1380     }

1381     /*

```

```

1983      * If kpm is available, use it.
1984      */
1985      if (kpm_vbase)
1986          return (PT_INDEX_PTR(hat_kpm_pfn2va(pfn), index));

1988      /*
1989      * Disable preemption and grab the CPU's hci_mutex
1990      */
1991      kpreempt_disable();

1993      ASSERT(CPU->cpu_hat_info != NULL);
1994      ASSERT(!(getcr4() & CR4_PCIDE));

1996      mutex_enter(&CPU->cpu_hat_info->hci_mutex);
1997      x = PWIN_TABLE(CPU->cpu_id);
1998      pteptr = (x86pte_t *)PWIN_PTE_VA(x);
1999 #ifdef __xpv
2000      if (mmu.pae_hat)
2001          pte = *pteptr;
2002      else
2003          pte = *(x86pte32_t *)pteptr;
2004 #endif

2006      newpte = MAKEPTE(pfn, 0) | mmu.pt_global | mmu.pt_nx;

2008      /*
2009      * For hardware we can use a writable mapping.
2010      */
2011 #ifdef __xpv
2012      if (IN_XPV_PANIC())
2013 #endif
2014          newpte |= PT_WRITABLE;

2016      if (!PTE_EQUIV(newpte, pte)) {

2018 #ifdef __xpv
2019          if (!IN_XPV_PANIC()) {
2020              xen_map(newpte, PWIN_VA(x));
2021          } else
2022 #endif
2023          {
2024              XPV_ALLOW_PAGETABLE_UPDATES();
2025              if (mmu.pae_hat)
2026                  *pteptr = newpte;
2027              else
2028                  *(x86pte32_t *)pteptr = newpte;
2029              XPV_DISALLOW_PAGETABLE_UPDATES();
2030              mmu_flush_tlb_kpage((uintptr_t)PWIN_VA(x));
2031              mmu_tlbflush_entry((caddr_t)(PWIN_VA(x)));
2032          }
2033      return (PT_INDEX_PTR(PWIN_VA(x), index));
2034 }

2036 /*
2037 * Release access to a page table.
2038 */
2039 static void
2040 x86pte_release_pagetable(htable_t *ht)
2041 {
2042     if (ht->ht_flags & HTABLE_COPIED)
2043     /*
2044      * nothing to do for VLP htables
2045      */
2046     if (ht->ht_flags & HTABLE_VLP)
2047         return;

```

```

2045     x86pte_mapout();
2046 }
    _____unchanged_portion_omitted_____

2093 /*
2094 * Atomic unconditional set of a page table entry, it returns the previous
2095 * value. For pre-existing mappings if the PFN changes, then we don't care
2096 * about the old pte's REF / MOD bits. If the PFN remains the same, we leave
2097 * the MOD/REF bits unchanged.
2098 *
2099 * If asked to overwrite a link to a lower page table with a large page
2100 * mapping, this routine returns the special value of LPAGE_ERROR. This
2101 * allows the upper HAT layers to retry with a smaller mapping size.
2102 */
2103 x86pte_t
2104 x86pte_set(htable_t *ht, uint_t entry, x86pte_t new, void *ptr)
2105 {
2106     x86pte_t      old;
2107     x86pte_t      prev;
2108     x86pte_t      *ptep;
2109     level_t       l = ht->ht_level;
2110     x86pte_t      pfn_mask = (l != 0) ? PT_PADDR_LGPG : PT_PADDR;
2111     x86pte_t      n;
2112     uintptr_t     addr = htable_e2va(ht, entry);
2113     hat_t         *hat = ht->ht_hat;

2115     ASSERT(new != 0); /* don't use to invalidate a PTE, see x86pte_update */
2116     ASSERT(!(ht->ht_flags & HTABLE_SHARED_PFN));
2117     if (ptr == NULL)
2118         ptep = x86pte_access_pagetable(ht, entry);
2119     else
2120         ptep = ptr;

2122     /*
2123     * Install the new PTE. If remapping the same PFN, then
2124     * copy existing REF/MOD bits to new mapping.
2125     */
2126     do {
2127         prev = GET_PTE(ptep);
2128         n = new;
2129         if ((PTE_ISVALID(n) && (prev & pfn_mask) == (new & pfn_mask))
2130             n |= prev & (PT_REF | PT_MOD);

2132         /*
2133         * Another thread may have installed this mapping already,
2134         * flush the local TLB and be done.
2135         */
2136         if (prev == n) {
2137             old = new;
2138 #ifdef __xpv
2139             if (!IN_XPV_PANIC())
2140                 xen_flush_va((caddr_t)addr);
2141             else
2142 #endif
2143                 mmu_flush_tlb_page(addr);
2144                 mmu_tlbflush_entry((caddr_t)addr);
2145             goto done;
2147         /*
2148         * Detect if we have a collision of installing a large
2149         * page mapping where there already is a lower page table.
2150         */
2151         if (l > 0 && (prev & PT_VALID) && !(prev & PT_PAGESIZE)) {
2152             old = LPAGE_ERROR;

```

```

2153         goto done;
2154     }

2156     XPV_ALLOW_PAGETABLE_UPDATES();
2157     old = CAS_PTE(pte, prev, n);
2158     XPV_DISALLOW_PAGETABLE_UPDATES();
2159 } while (old != prev);

2161 /*
2162  * Do a TLB demap if needed, ie. the old pte was valid.
2163  *
2164  * Note that a stale TLB writeback to the PTE here either can't happen
2165  * or doesn't matter. The PFN can only change for NOSYNC|NOCONSIST
2166  * mappings, but they were created with REF and MOD already set, so
2167  * no stale writeback will happen.
2168  *
2169  * Segmap is the only place where remaps happen on the same pfn and for
2170  * that we want to preserve the stale REF/MOD bits.
2171  */
2172 if (old & PT_REF)
2173     hat_tlb_inval(hat, addr);

2175 done:
2176 if (ptr == NULL)
2177     x86pte_release_pagetable(ht);
2178 return (old);
2179 }

2181 /*
2182  * Atomic compare and swap of a page table entry. No TLB invalidates are done.
2183  * This is used for links between pagetables of different levels.
2184  * Note we always create these links with dirty/access set, so they should
2185  * never change.
2186  */
2187 x86pte_t
2188 x86pte_cas(htable_t *ht, uint_t entry, x86pte_t old, x86pte_t new)
2189 {
2190     x86pte_t     pte;
2191     x86pte_t     *pteptr;
2192 #ifdef __xpv
2193     /*
2194      * We can't use writable pagetables for upper level tables, so fake it.
2195      */
2196     mmu_update_t t[2];
2197     int cnt = 1;
2198     int count;
2199     maddr_t ma;

2201     if (!IN_XPV_PANIC()) {
2202         ASSERT(!(ht->ht_flags & HTABLE_COPIED));
2203         ASSERT(!(ht->ht_flags & HTABLE_VLP)); /* no VLP yet */
2204         ma = pa_to_ma(PT_INDEX_PHYSADDR(pfn_to_pa(ht->ht_pfn), entry));
2205         t[0].ptr = ma | MMU_NORMAL_PT_UPDATE;
2206         t[0].val = new;

2207 #if defined(__amd64)
2208         /*
2209          * On the 64-bit hypervisor we need to maintain the user mode
2210          * top page table too.
2211          */
2212         if (ht->ht_level == mmu.max_level && ht->ht_hat != kas.a_hat) {
2213             ma = pa_to_ma(PT_INDEX_PHYSADDR(pfn_to_pa(
2214                 ht->ht_hat->hat_user_ptable), entry));
2215             t[1].ptr = ma | MMU_NORMAL_PT_UPDATE;
2216             t[1].val = new;
2217             ++cnt;

```

```

2218     }
2219 #endif /* __amd64 */

2221     if (HYPERVISOR_mmu_update(t, cnt, &count, DOMID_SELF))
2222         panic("HYPERVISOR_mmu_update() failed");
2223     ASSERT(count == cnt);
2224     return (old);
2225 }
2226 #endif

2227 pte = x86pte_access_pagetable(ht, entry);
2228 XPV_ALLOW_PAGETABLE_UPDATES();
2229 pte = CAS_PTE(pte, old, new);
2230 XPV_DISALLOW_PAGETABLE_UPDATES();
2231 x86pte_release_pagetable(ht);
2232 return (pte);
2233 }

    _____ unchanged_portion_omitted _____

2355 #ifndef __xpv
2356 /*
2357  * Copy page tables - this is just a little more complicated than the
2358  * previous routines. Note that it's also not atomic! It also is never
2359  * used for HTABLE_COPIED pagetables.
2360  * used for VLP pagetables.
2361  */
2362 void
2363 x86pte_copy(htable_t *src, htable_t *dest, uint_t entry, uint_t count)
2364 {
2365     caddr_t src_va;
2366     caddr_t dst_va;
2367     size_t size;
2368     x86pte_t *pteptr;
2369     x86pte_t pte;

2370     ASSERT(khat_running);
2371     ASSERT(!(dest->ht_flags & HTABLE_COPIED));
2372     ASSERT(!(src->ht_flags & HTABLE_COPIED));
2373     ASSERT(!(dest->ht_flags & HTABLE_VLP));
2374     ASSERT(!(src->ht_flags & HTABLE_VLP));
2375     ASSERT(!(src->ht_flags & HTABLE_SHARED_PFN));
2376     ASSERT(!(dest->ht_flags & HTABLE_SHARED_PFN));

2377     /*
2378      * Acquire access to the CPU pagetable windows for the dest and source.
2379      */
2380     dst_va = (caddr_t)x86pte_access_pagetable(dest, entry);
2381     if (kpm_vbase) {
2382         src_va = (caddr_t)
2383             PT_INDEX_PTR(hat_kpm_pfn2va(src->ht_pfn), entry);
2384     } else {
2385         uint_t x = PWIN_SRC(CPU->cpu_id);
2386         ASSERT(!(getcr4() & CR4_PCIDE));

2387         /*
2388          * Finish defining the src pagetable mapping
2389          */
2390         src_va = (caddr_t)PT_INDEX_PTR(PWIN_VA(x), entry);
2391         pte = MAKEPTE(src->ht_pfn, 0) | mmu.pt_global | mmu.pt_nx;
2392         pteptr = (x86pte_t *)PWIN_PTE_VA(x);
2393         if (mmu.pae_hat)
2394             *pteptr = pte;
2395         else
2396             *(x86pte32_t *)pteptr = pte;
2397         mmu_flush_tlb_kpage((uintptr_t)PWIN_VA(x));
2398         mmu_tlbflush_entry((caddr_t)(PWIN_VA(x)));

```



```

*****
11423 Fri Apr 6 17:25:06 2018
new/usr/src/uts/i86pc/vm/htable.h
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22 * Copyright 2007 Sun Microsystems, Inc. All rights reserved.
23 * Use is subject to license terms.
24 */
25 /*
26 * Copyright (c) 2014 by Delphix. All rights reserved.
27 * Copyright 2018 Joyent, Inc.
28 */

30 #ifndef _VM_HTABLE_H
31 #define _VM_HTABLE_H

33 #ifdef __cplusplus
34 extern "C" {
35 #endif

37 #if defined(__GNUC__) && defined(_ASM_INLINES) && defined(_KERNEL)
38 #include <asm/htable.h>
39 #endif

41 extern void atomic_andb(uint8_t *addr, uint8_t value);
42 extern void atomic_orb(uint8_t *addr, uint8_t value);
43 extern void atomic_incl6(uint16_t *addr);
44 extern void atomic_decl6(uint16_t *addr);
44 extern void mmu_tlbflush_entry(caddr_t addr);

46 /*
47  * Each hardware page table has an htable_t describing it.
48  *
49  * We use a reference counter mechanism to detect when we can free an htable.
50  * In the implementation the reference count is split into 2 separate counters:
51  *
52  *   ht_busy is a traditional reference count of uses of the htable pointer
53  *
54  *   ht_valid_cnt is a count of how references are implied by valid PTE/PTP
55  *   entries in the pagetable
56  *
57  * ht_busy is only incremented by htable_lookup() or htable_create()
58  * while holding the appropriate hash_table mutex. While installing a new

```

```

59  * valid PTE or PTP, in order to increment ht_valid_cnt a thread must have
60  * done an htable_lookup() or htable_create() but not the htable_release yet.
61  *
62  * htable_release(), while holding the mutex, can know that if
63  * busy == 1 and valid_cnt == 0, the htable can be free'd.
64  *
65  * The fields have been ordered to make htable_lookup() fast. Hence,
66  * ht_hat, ht_vaddr, ht_level and ht_next need to be clustered together.
67  */
68 struct htable {
69     struct htable *ht_next;      /* forward link for hash table */
70     struct hat    *ht_hat;      /* hat this mapping comes from */
71     uintptr_t    ht_vaddr;      /* virt addr at start of this table */
72     int8_t       ht_level;      /* page table level: 0=4K, 1=2M, ... */
73     uint8_t      ht_flags;      /* see below */
74     int16_t      ht_busy;       /* implements locking protocol */
75     int16_t      ht_valid_cnt;  /* # of valid entries in this table */
76     uint32_t     ht_lock_cnt;   /* # of locked entries in this table */
77     /* never used for kernel hat */
78     pfn_t        ht_pfn;        /* pfn of page of the pagetable */
79     struct htable *ht_prev;     /* backward link for hash table */
80     struct htable *ht_parent;   /* htable that points to this htable */
81     struct htable *ht_shares;   /* for HTABLE_SHARED_PFN only */
82 };
83 typedef struct htable htable_t;

85 /*
86  * Flags values for htable ht_flags field:
87  *
88  * HTABLE_COPIED - This is the top level htable of a HAT being used with per-CPU
89  * pagetables.
90  * HTABLE_VLP - this is the top level htable of a VLP HAT.
91  *
92  * HTABLE_SHARED_PFN - this htable had its PFN assigned from sharing another
93  * htable. Used by hat_share() for ISM.
94 #define HTABLE_COPIED (0x01)
95 #define HTABLE_VLP (0x01)
96 #define HTABLE_SHARED_PFN (0x02)

97 /*
98  * The htable hash table hashing function. The 28 is so that high
99  * order bits are include in the hash index to skew the wrap
100 * around of addresses. Even though the hash buckets are stored per
101 * hat we include the value of hat pointer in the hash function so
102 * that the secondary hash for the htable mutex winds up begin different in
103 * every address space.
104 */
105 #define HTABLE_HASH(hat, va, lvl) \
106     (((va) >> LEVEL_SHIFT(1)) + ((va) >> 28) + (lvl) + \
107     ((uintptr_t)(hat) >> 4) & ((hat)->hat_num_hash - 1))

109 /*
110  * Each CPU gets a unique hat_cpu_info structure in cpu_hat_info. For more
111  * information on its use and members, see uts/i86pc/vm/hat_i86.c.
112  * Each CPU gets a unique hat_cpu_info structure in cpu_hat_info.
113 */
113 struct hat_cpu_info {
114     kmutex_t hci_mutex;      /* mutex to ensure sequential usage */
115 #if defined(__amd64)
116     pfn_t    hci_pcp_l3pfn;  /* pfn of hci_pcp_l3ptes */
117     pfn_t    hci_pcp_l2pfn;  /* pfn of hci_pcp_l2ptes */
118     x86pte_t *hci_pcp_l3ptes; /* PCP Level==3 pagetable (top) */
119     x86pte_t *hci_pcp_l2ptes; /* PCP Level==2 pagetable */
120     struct hat *hci_user_hat; /* CPU specific HAT */
121     pfn_t    hci_user_l3pfn; /* pfn of hci_user_l3ptes */

```



```

122     x86pte_t *hci_user_l3ptes;    /* PCP User L3 pagetable */
123     pfn_t    hci_vlp_pfn;        /* pfn of hci_vlp_l3ptes */
124     x86pte_t *hci_vlp_l3ptes;    /* VLP Level==3 pagetable (top) */
125     x86pte_t *hci_vlp_l2ptes;    /* VLP Level==2 pagetable */
126 #endif /* __amd64 */
127 };

128 /*
129 * Compute the last page aligned VA mapped by an htable.
130 * Given a va and a level, compute the virtual address of the start of the
131 * next page at that level.
132 *
133 * XX64 - The check for the VA hole needs to be better generalized.
134 */
135 #if defined(__amd64)
136 #define HTABLE_NUM_PTES(ht)      (((ht)->ht_flags & HTABLE_COPIED) ? \
137     ((ht)->ht_level == mmu.max_level) ? 512 : 4) : 512)
138 #define HTABLE_NUM_PTES(ht)     (((ht)->ht_flags & HTABLE_VLP) ? 4 : 512)
139 #define HTABLE_LAST_PAGE(ht)    \
140     ((ht)->ht_level == mmu.max_level ? ((uintptr_t)0UL - MMU_PAGESIZE) : \
141     ((ht)->ht_vaddr - MMU_PAGESIZE + \
142     ((uintptr_t)HTABLE_NUM_PTES(ht) << LEVEL_SHIFT((ht)->ht_level)))
143 #define NEXT_ENTRY_VA(va, l)    \
144     ((va & LEVEL_MASK(l)) + LEVEL_SIZE(l) == mmu.hole_start ? \
145     mmu.hole_end : (va & LEVEL_MASK(l)) + LEVEL_SIZE(l))
146 #elif defined(__i386)
147 #define HTABLE_NUM_PTES(ht)     \
148     (!mmu.pae_hat ? 1024 : ((ht)->ht_level == 2 ? 4 : 512))
149 #define HTABLE_LAST_PAGE(ht)    \
150     ((ht)->ht_vaddr - MMU_PAGESIZE + \
151     ((uintptr_t)HTABLE_NUM_PTES(ht) << LEVEL_SHIFT((ht)->ht_level)))
152 #define NEXT_ENTRY_VA(va, l)    \
153     ((va & LEVEL_MASK(l)) + LEVEL_SIZE(l))
154 #endif

155 #endif

156 #if defined(_KERNEL)
157 /*
158 * initialization function called from hat_init()
159 */
160 extern void htable_init(void);

161 /*
162 * Functions to lookup, or "lookup and create", the htable corresponding
163 * to the virtual address "vaddr" in the "hat" at the given "level" of
164 * page tables. htable_lookup() may return NULL if no such entry exists.
165 *
166 * On return the given htable is marked busy (a shared lock) - this prevents
167 * the htable from being stolen or freed) until htable_release() is called.
168 *
169 * If kalloc_flag is set on an htable_create() we can't call kmem allocation
170 * routines for this htable, since it's for the kernel hat itself.
171 *
172 * htable_acquire() is used when an htable pointer has been extracted from
173 * an hment and we need to get a reference to the htable.
174 */
175 extern htable_t *htable_lookup(struct hat *hat, uintptr_t vaddr, level_t level);
176 extern htable_t *htable_create(struct hat *hat, uintptr_t vaddr, level_t level,
177     htable_t *shared);

```

```

178 extern void htable_acquire(htable_t *);
179 extern void htable_release(htable_t *ht);
180 extern void htable_destroy(htable_t *ht);

181 /*
182 * Code to free all remaining htables for a hat. Called after the hat is no
183 * longer in use by any thread.
184 */
185 extern void htable_purge_hat(struct hat *hat);

186 /*
187 * Find the htable, page table entry index, and PTE of the given virtual
188 * address. If not found returns NULL. When found, returns the htable_t *,
189 * sets entry, and has a hold on the htable.
190 */
191 extern htable_t *htable_getpte(struct hat *, uintptr_t, uint_t *, x86pte_t *,
192     level_t);

193 /*
194 * Similar to hat_getpte(), except that this only succeeds if a valid
195 * page mapping is present.
196 */
197 extern htable_t *htable_getpage(struct hat *hat, uintptr_t va, uint_t *entry);

198 /*
199 * Called to allocate initial/additional htables for reserve.
200 */
201 extern void htable_initial_reserve(uint_t);
202 extern void htable_reserve(uint_t);

203 /*
204 * Used to readjust the htable reserve after the reserve list has been used.
205 * Also called after boot to release left over boot reserves.
206 */
207 extern void htable_adjust_reserve(void);

208 /*
209 * return number of bytes mapped by all the htables in a given hat
210 */
211 extern size_t htable_mapped(struct hat *);

212 /*
213 * Attach initial pagetables as htables
214 */
215 extern void htable_attach(struct hat *, uintptr_t, level_t, struct htable *,
216     pfn_t);

217 /*
218 * Routine to find the next populated htable at or above a given virtual
219 * address. Can specify an upper limit, or HTABLE_WALK_TO_END to indicate
220 * that it should search the entire address space. Similar to
221 * hat_getpte(), but used for walking through address ranges. It can be
222 * used like this:
223 *
224 *     va = ...
225 *     ht = NULL;
226 *     while (va < end_va) {
227 *         pte = htable_walk(hat, &ht, &va, end_va);
228 *         if (!pte)
229 *             break;
230 *
231 *         ... code to operate on page at va ...
232 *
233 *         va += LEVEL_SIZE(ht->ht_level);

```

```

250 *      }
251 *      if (ht)
252 *          htable_release(ht);
253 *
254 */
255 extern x86pte_t htable_walk(struct hat *hat, htable_t **ht, uintptr_t *va,
256                          uintptr_t eaddr);

258 #define HTABLE_WALK_TO_END ((uintptr_t)-1)

260 /*
261 * Utilities convert between virtual addresses and page table entry indices.
262 */
263 extern uint_t htable_va2entry(uintptr_t va, htable_t *ht);
264 extern uintptr_t htable_e2va(htable_t *ht, uint_t entry);

266 /*
267 * Interfaces that provide access to page table entries via the htable.
268 *
269 * Note that all accesses except x86pte_copy() and x86pte_zero() are atomic.
270 */
271 extern void      x86pte_cpu_init(cpu_t *);
272 extern void      x86pte_cpu_fini(cpu_t *);

274 extern x86pte_t x86pte_get(htable_t *, uint_t entry);

276 /*
277 * x86pte_set returns LPAGE_ERROR if it's asked to overwrite a page table
278 * link with a large page mapping.
279 */
280 #define LPAGE_ERROR (~(x86pte_t)1)
281 extern x86pte_t x86pte_set(htable_t *, uint_t entry, x86pte_t new, void *);

283 extern x86pte_t x86pte_inval(htable_t *ht, uint_t entry,
284                             x86pte_t old, x86pte_t *ptr, boolean_t tlb);

286 extern x86pte_t x86pte_update(htable_t *ht, uint_t entry,
287                              x86pte_t old, x86pte_t new);

289 extern void      x86pte_copy(htable_t *src, htable_t *dest, uint_t entry,
290                              uint_t cnt);

292 /*
293 * access to a pagetable knowing only the pfn
294 */
295 extern x86pte_t *x86pte_mapin(pfn_t, uint_t, htable_t *);
296 extern void x86pte_mapout(void);

298 /*
299 * these are actually inlines for "lock; incw", "lock; decw", etc. instructions.
300 */
301 #define HTABLE_INC(x)      atomic_inc16((uint16_t *)&x)
302 #define HTABLE_DEC(x)      atomic_dec16((uint16_t *)&x)
303 #define HTABLE_LOCK_INC(ht) atomic_inc_32(&(ht)->ht_lock_cnt)
304 #define HTABLE_LOCK_DEC(ht) atomic_dec_32(&(ht)->ht_lock_cnt)

306 #ifdef __xpv
307 extern void xen_flush_va(caddr_t va);
308 extern void xen_gflush_va(caddr_t va, cpuset_t);
309 extern void xen_flush_tlb(void);
310 extern void xen_gflush_tlb(cpuset_t);
311 extern void xen_pin(pfn_t, level_t);
312 extern void xen_unpin(pfn_t);
313 extern int xen_kpm_page(pfn_t, uint_t);
315 /*

```

```

316 * The hypervisor maps all page tables into our address space read-only.
317 * Under normal circumstances, the hypervisor then handles all updates to
318 * the page tables underneath the covers for us. However, when we are
319 * trying to dump core after a hypervisor panic, the hypervisor is no
320 * longer available to do these updates. To work around the protection
321 * problem, we simply disable write-protect checking for the duration of a
322 * pagetable update operation.
323 */
324 #define XPV_ALLOW_PAGETABLE_UPDATES() \
325 { \
326     if (IN_XPV_PANIC()) \
327         setcr0((getcr0() & ~CR0_WP) & 0xffffffff); \
328 }
329 #define XPV_DISALLOW_PAGETABLE_UPDATES() \
330 { \
331     if (IN_XPV_PANIC() > 0) \
332         setcr0((getcr0() | CR0_WP) & 0xffffffff); \
333 }

335 #else /* __xpv */

337 #define XPV_ALLOW_PAGETABLE_UPDATES()
338 #define XPV_DISALLOW_PAGETABLE_UPDATES()

340 #endif

342 #endif /* _KERNEL */

345 #ifdef __cplusplus
346 }
_____ unchanged portion omitted

```

new/usr/src/uts/i86pc/vm/i86_mmu.c

1

```
*****
15116 Fri Apr 6 17:25:06 2018
new/usr/src/uts/i86pc/vm/i86_mmu.c
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22 * Copyright 2010 Sun Microsystems, Inc. All rights reserved.
23 * Use is subject to license terms.
24 *
25 * Copyright 2018 Joyent, Inc.
26 */

28 #include <sys/t_lock.h>
29 #include <sys/memlist.h>
30 #include <sys/cpuvar.h>
31 #include <sys/vmem.h>
32 #include <sys/mman.h>
33 #include <sys/vm.h>
34 #include <sys/kmem.h>
35 #include <sys/cmn_err.h>
36 #include <sys/debug.h>
37 #include <sys/vm_machparam.h>
38 #include <sys/tss.h>
39 #include <sys/vnode.h>
40 #include <vm/hat.h>
41 #include <vm/anon.h>
42 #include <vm/as.h>
43 #include <vm/page.h>
44 #include <vm/seg.h>
45 #include <vm/seg_kmem.h>
46 #include <vm/seg_map.h>
47 #include <vm/hat_i86.h>
48 #include <sys/promif.h>
49 #include <sys/x86_archext.h>
50 #include <sys/system.h>
51 #include <sys/archsystem.h>
52 #include <sys/sunddi.h>
53 #include <sys/ddidmareq.h>
54 #include <sys/controlregs.h>
55 #include <sys/reboot.h>
56 #include <sys/kdi.h>
57 #include <sys/bootconf.h>
58 #include <sys/bootsvcs.h>
59 #include <sys/bootinfo.h>
```

new/usr/src/uts/i86pc/vm/i86_mmu.c

2

```
60 #include <vm/kboot_mmu.h>

62 #ifdef __xpv
63 #include <sys/hypervisor.h>
64 #endif

66 #define ON_USER_HAT(cpu) \
67     ((cpu)->cpu_m.mcpu_current_hat != NULL && \
68     (cpu)->cpu_m.mcpu_current_hat != kas.a_hat)
64 caddr_t
65 i86devmap(pfn_t pf, pgcnt_t pgcnt, uint_t prot)
66 {
67     caddr_t addr;
68     caddr_t addr1;
69     page_t *pp;

71     addr1 = addr = vmem_alloc(heap_arena, mmu_ptob(pgcnt), VM_SLEEP);

73     for (; pgcnt != 0; addr += MMU_PAGESIZE, ++pf, --pgcnt) {
74         pp = page_numtopp_nolock(pf);
75         if (pp == NULL) {
76             hat_devload(kas.a_hat, addr, MMU_PAGESIZE, pf,
77                 prot | HAT_NOSYNC, HAT_LOAD_LOCK);
78         } else {
79             hat_memload(kas.a_hat, addr, pp,
80                 prot | HAT_NOSYNC, HAT_LOAD_LOCK);
81         }
82     }

84     return (addr1);
85 }

70 /*
88 * This routine is like page_numtopp, but accepts only free pages, which
89 * it allocates (unfrees) and returns with the exclusive lock held.
90 * It is used by machdep.c/dma_init() to find contiguous free pages.
91 *
92 * XXX this and some others should probably be in vm_machdep.c
93 */
94 page_t *
95 page_numtopp_alloc(pfn_t pfn)
96 {
97     page_t *pp;

99 retry:
100     pp = page_numtopp_nolock(pfn);
101     if (pp == NULL) {
102         return (NULL);
103     }

105     if (!page_trylock(pp, SE_EXCL)) {
106         return (NULL);
107     }

109     if (page_pptonum(pp) != pfn) {
110         page_unlock(pp);
111         goto retry;
112     }

114     if (!PP_ISFREE(pp)) {
115         page_unlock(pp);
116         return (NULL);
117     }
118     if (pp->p_szc) {
119         page_demote_free_pages(pp);
120         page_unlock(pp);
```

```

121         goto retry;
122     }
123
124     /* If associated with a vnode, destroy mappings */
125
126     if (pp->p_vnode) {
127
128         page_destroy_free(pp);
129
130         if (!page_lock(pp, SE_EXCL, (kmutex_t *)NULL, P_NO_RECLAIM)) {
131             return (NULL);
132         }
133
134         if (page_pptonum(pp) != pfnnum) {
135             page_unlock(pp);
136             goto retry;
137         }
138     }
139
140     if (!PP_ISFREE(pp)) {
141         page_unlock(pp);
142         return (NULL);
143     }
144
145     if (!page_reclaim(pp, (kmutex_t *)NULL))
146         return (NULL);
147
148     return (pp);
149 }
150
151 /*
152  * Flag is not set early in boot. Once it is set we are no longer
153  * using boot's page tables.
154  */
155 uint_t khat_running = 0;
156
157 /*
158  * This procedure is callable only while the boot loader is in charge of the
159  * MMU. It assumes that PA == VA for page table pointers. It doesn't live in
160  * kboot_mmu.c since it's used from common code.
161  */
162 pfn_t
163 va_to_pfn(void *vaddr)
164 {
165     uintptr_t    des_va = ALIGN2PAGE(vaddr);
166     uintptr_t    va = des_va;
167     size_t       len;
168     uint_t       prot;
169     pfn_t        pfn;
170
171     if (khat_running)
172         panic("va_to_pfn(): called too late\n");
173
174     if (kvm_probe(&va, &len, &pfn, &prot) == 0)
175         return (PFN_INVALID);
176     if (va > des_va)
177         return (PFN_INVALID);
178     if (va < des_va)
179         pfn += mmu_btop(des_va - va);
180     return (pfn);
181 }
182
183 unchanged_portion_omitted
184 #endif
185
186 /*
187  * Routine to pre-allocate data structures for hat_kern_setup(). It computes

```

```

223  * how many pagetables it needs by walking the boot loader's page tables.
224  */
225 /*ARGSUSED*/
226 void
227 hat_kern_alloc(
228     caddr_t segmap_base,
229     size_t segmap_size,
230     caddr_t ekernelheap)
231 {
232     uintptr_t    last_va = (uintptr_t)-1;    /* catch 1st time */
233     uintptr_t    va = 0;
234     size_t       size;
235     pfn_t        pfn;
236     uint_t       prot;
237     uint_t       table_cnt = 1;
238     uint_t       mapping_cnt;
239     level_t      start_level;
240     level_t      l;
241     struct memlist *pmem;
242     level_t      lpagel = mmu.max_page_level;
243     uint64_t      paddr;
244     int64_t       psize;
245     int           nwindows;
246
247     if (kpm_size > 0) {
248         /*
249          * Create the kpm page tables. When running on the
250          * hypervisor these are made read/only at first.
251          * Later we'll add write permission where possible.
252          */
253         for (pmem = phys_install; pmem; pmem = pmem->ml_next) {
254             paddr = pmem->ml_address;
255             psize = pmem->ml_size;
256             while (psize >= MMU_PAGESIZE) {
257                 /* find the largest page size */
258                 for (l = lpagel; l > 0; l--) {
259                     if ((paddr & LEVEL_OFFSET(l)) == 0 &&
260                         psize > LEVEL_SIZE(l))
261                         break;
262                 }
263             }
264 #if defined(__xpv)
265             /*
266              * Create read/only mappings to avoid
267              * conflicting with pagetable usage
268              */
269             xen_kpm_create(paddr, l);
270 #else
271             kvm_map((uintptr_t)kpm_vbase + paddr, paddr,
272                 l, l);
273 #endif
274             paddr += LEVEL_SIZE(l);
275             psize -= LEVEL_SIZE(l);
276         }
277     }
278
279     /*
280      * If this machine doesn't have a kpm segment, we need to allocate
281      * a small number of 'windows' which can be used to map pagetables.
282      */
283     nwindows = (kpm_size == 0) ? 2 * NCPU : 0;
284
285 #if defined(__xpv)
286     /*
287      * On a hypervisor, these windows are also used by the xpv_panic

```

```

289     * code, where we need one window for each level of the pagetable
290     * hierarchy.
291     */
292     nwindows = MAX(nwindows, mmu.max_level);
293 #endif

295     if (nwindows != 0) {
296         /*
297          * Create the page windows and 1 page of VA in
298          * which we map the PTEs of those windows.
299          */
300         mmu.pwin_base = vmem_xalloc(heap_arena, nwindows * MMU_PAGESIZE,
301             LEVEL_SIZE(1), 0, 0, NULL, NULL, VM_SLEEP);
302         ASSERT(nwindows <= MMU_PAGESIZE / mmu.pte_size);
303         mmu.pwin_pte_va = vmem_xalloc(heap_arena, MMU_PAGESIZE,
304             MMU_PAGESIZE, 0, 0, NULL, NULL, VM_SLEEP);

306         /*
307          * Find/Create the page table window mappings.
308          */
309         paddr = 0;
310         (void) find_pte((uintptr_t)mmu.pwin_base, &paddr, 0, 0);
311         ASSERT(paddr != 0);
312         ASSERT((paddr & MMU_PAGEOFFSET) == 0);
313         mmu.pwin_pte_pa = paddr;
314 #ifdef __xpv
315         (void) find_pte((uintptr_t)mmu.pwin_pte_va, NULL, 0, 0);
316         kbm_read_only((uintptr_t)mmu.pwin_pte_va, mmu.pwin_pte_pa);
317 #else
318         kbm_map((uintptr_t)mmu.pwin_pte_va, mmu.pwin_pte_pa, 0, 1);
319 #endif
320     }

322     /*
323     * Walk the boot loader's page tables and figure out
324     * how many tables and page mappings there will be.
325     */
326     while (kbm_probe(&va, &size, &pf, &prot) != 0) {
327         /*
328          * At each level, if the last_va falls into a new htable,
329          * increment table_cnt. We can stop at the 1st level where
330          * they are in the same htable.
331          */
332         start_level = 0;
333         while (start_level <= mmu.max_page_level) {
334             if (size == LEVEL_SIZE(start_level))
335                 break;
336             start_level++;
337         }

339         for (l = start_level; l < mmu.max_level; ++l) {
340             if (va >> LEVEL_SHIFT(l + 1) ==
341                 last_va >> LEVEL_SHIFT(l + 1))
342                 break;
343             ++table_cnt;
344         }
345         last_va = va;
346         l = (start_level == 0) ? 1 : start_level;
347         va = (va & LEVEL_MASK(l)) + LEVEL_SIZE(l);
348     }

350     /*
351     * Besides the boot loader mappings, we're going to fill in
352     * the entire top level page table for the kernel. Make sure there's
353     * enough reserve for that too.
354     */

```

```

355         table_cnt += mmu.top_level_count - ((kernelbase >>
356             LEVEL_SHIFT(mmu.max_level)) & (mmu.top_level_count - 1));

439 #if defined(__i386)
440 /*
441  * The 32 bit PAE hat allocates tables one level below the top when
442  * kernelbase isn't 1 Gig aligned. We'll just be sloppy and allocate
443  * a bunch more to the reserve. Any unused will be returned later.
444  * Note we've already counted these mappings, just not the extra
445  * pagetables.
446  */
447     if (mmu.pae_hat != 0 && (kernelbase & LEVEL_OFFSET(mmu.max_level)) != 0)
448         table_cnt += mmu.ptes_per_table -
449             ((kernelbase & LEVEL_OFFSET(mmu.max_level)) >>
450                 LEVEL_SHIFT(mmu.max_level - 1));
451 #endif

453     /*
454     * Add 1/4 more into table_cnt for extra slop. The unused
455     * slop is freed back when we htable_adjust_reserve() later.
456     */
457     table_cnt += table_cnt >> 2;

464     /*
465     * We only need mapping entries (hments) for shared pages.
466     * This should be far, far fewer than the total possible,
467     * We'll allocate enough for 1/16 of all possible PTEs.
468     */
469     mapping_cnt = (table_cnt * mmu.ptes_per_table) >> 4;

471     /*
472     * Now create the initial htable/hment reserves
473     */
474     htable_initial_reserve(table_cnt);
475     hment_reserve(mapping_cnt);
476     x86pte_cpu_init(CPU);
477 }

480 /*
481 * This routine handles the work of creating the kernel's initial mappings
482 * by deciphering the mappings in the page tables created by the boot program.
483 *
484 * We maintain large page mappings, but only to a level 1 pagesize.
485 * The boot loader can only add new mappings once this function starts.
486 * In particular it can not change the pagesize used for any existing
487 * mappings or this code breaks!
488 */

490 void
491 hat_kern_setup(void)
492 {
493     /*
494     * Attach htables to the existing pagetables
495     */
496     /* BEGIN CSTYLED */
497     htable_attach(kas.a_hat, 0, mmu.max_level, NULL,
498 #ifdef __xpv
499         mmu_btop(xen_info->pt_base - ONE_GIG));
500 #else
501         mmu_btop(getcr3_pa());
502 #endif
503     /* END CSTYLED */
504 }

505 #if defined(__xpv)

```

```

500 #if defined(__i386) && !defined(__xpv)
501     CPU->cpu_tss->tss_cr3 = dftss0->tss_cr3 = getcr3();
502 #endif /* __i386 */

504 #if defined(__xpv) && defined(__amd64)
505 /*
506  * Try to make the kpm mappings r/w. Failures here are OK, as
507  * it's probably just a pagetable
508  */
509     xen_kpm_finish_init();
510 #endif

513 /*
514  * The kernel HAT is now officially open for business.
515  */
516     khat_running = 1;

518     CPuset_ATOMIC_ADD(kas.a_hat->hat_cpus, CPU->cpu_id);
519     CPU->cpu_current_hat = kas.a_hat;
520 }

522 #ifndef __xpv

524 /*
525  * Note that the INVPCID_ALL* variants can be used even in the !PCIDE case, but
526  * INVPCID_ADDR isn't.
527  */
528 static void
529 invpcid(uint64_t type, uint64_t pcid, uintptr_t addr)
530 {
531     ulong_t flag;
532     uint64_t cr4;

534     if (x86_use_invpcid == 1) {
535         ASSERT(is_x86_feature(x86_featureset, X86FSET_INVPCID));
536         invpcid_insn(type, pcid, addr);
537         return;
538     }

539     switch (type) {
540     case INVPCID_ALL_GLOBAL:
541         flag = intr_clear();
542         cr4 = getcr4();
543         setcr4(cr4 & ~(ulong_t)CR4_PGE);
544         setcr4(cr4 | CR4_PGE);
545         intr_restore(flag);
546         break;

547     case INVPCID_ALL_NONGLOBAL:
548         if (!(getcr4() & CR4_PCIDE)) {
549             reload_cr3();
550         } else {
551             flag = intr_clear();
552             cr4 = getcr4();
553             setcr4(cr4 & ~(ulong_t)CR4_PGE);
554             setcr4(cr4 | CR4_PGE);
555             intr_restore(flag);
556         }
557         break;

558     case INVPCID_ADDR:
559         if (pcid == PCID_USER) {
560             flag = intr_clear();
561             ASSERT(addr < kernelbase);
562             ASSERT(ON_USER_HAT(CPU));
563             ASSERT(CPU->cpu_m.mcpu_kpti.kf_user_cr3 != 0);

```

```

564         tr_mmu_flush_user_range(addr, MMU_PAGESIZE,
565                                 MMU_PAGESIZE, CPU->cpu_m.mcpu_kpti.kf_user_cr3);
566         intr_restore(flag);
567     } else {
568         mmu_invlpg((caddr_t)addr);
569     }
570     break;

572 default:
573     panic("unsupported invpcid(%lu)", type);
574     break;
575 }

577 /*
578  * Flush one kernel mapping.
579  * We want to assert on kernel space here mainly for reasoning about the PCIDE
580  * case: namely, this flush should never need to flush a non-current PCID
581  * mapping. This presumes we never have reason to flush the kernel regions
582  * available to PCID_USER (the trampolines and so on). It also relies on
583  * PCID_KERNEL == PCID_NONE.
584  */
585 void
586 mmu_flush_tlb_kpage(uintptr_t va)
587 {
588     ASSERT(va >= kernelbase);
589     ASSERT(getpcid() == PCID_KERNEL);
590     mmu_invlpg((caddr_t)va);
591 }

593 /*
594  * Flush one mapping: local CPU version of hat_tlb_inval().
595  * If this is a userspace address in the PCIDE case, we need two invalidations,
596  * one for any potentially stale PCID_USER mapping, as well as any established
597  * while in the kernel.
598  */
599 void
600 mmu_flush_tlb_page(uintptr_t va)
601 {
602     ASSERT(getpcid() == PCID_KERNEL);

604     if (va >= kernelbase) {
605         mmu_flush_tlb_kpage(va);
606         return;
607     }

608     if (!(getcr4() & CR4_PCIDE)) {
609         mmu_invlpg((caddr_t)va);
610         return;
611     }

612     /*
613      * Yes, kas will need to flush below kernelspace, at least during boot.
614      * But there's no PCID_USER context.
615      */
616     if (ON_USER_HAT(CPU))
617         invpcid(INVPCID_ADDR, PCID_USER, va);
618     invpcid(INVPCID_ADDR, PCID_KERNEL, va);
619 }

621 static void
622 mmu_flush_tlb_range(uintptr_t addr, size_t len, size_t pgsz)
623 {
624     EQUIV(addr < kernelbase, (addr + len - 1) < kernelbase);

```

```

533     ASSERT(len > 0);
534     ASSERT(pgsz != 0);

536     if (!(getcr4() & CR4_PCIDE) || x86_use_invpcid == 1) {
537         for (uintptr_t va = addr; va < (addr + len); va += pgsz)
538             mmu_flush_tlb_page(va);
539         return;
540     }

542     /*
543     * As an emulated invpcid() in the PCIDE case requires jumping
544     * cr3s, we batch the invalidations. We should only need to flush the
545     * user range if we're on a user-space HAT.
546     */
547     if (addr < kernelbase && ON_USER_HAT(CPU)) {
548         ulong_t flag = intr_clear();
549         ASSERT(CPU->cpu_m.mcpu_kpti.kf_user_cr3 != 0);
550         tr_mmu_flush_user_range(addr, len, pgsz,
551             CPU->cpu_m.mcpu_kpti.kf_user_cr3);
552         intr_restore(flag);
553     }

555     for (uintptr_t va = addr; va < (addr + len); va += pgsz)
556         mmu_invpg((caddr_t)va);
557 }

559 /*
560 * MMU TLB (and PT cache) flushing on this CPU.
561 *
562 * FLUSH_TLB_ALL: invalidate everything, all PCIDs, all PT_GLOBAL.
563 * FLUSH_TLB_NONGLOBAL: invalidate all PCIDs, excluding PT_GLOBAL.
564 * FLUSH_TLB_RANGE: invalidate the given range, including PCID_USER
565 * mappings as appropriate. If using invpcid, PT_GLOBAL mappings are not
566 * invalidated.
567 */
568 void
569 mmu_flush_tlb(flush_tlb_type_t type, tlb_range_t *range)
570 {
571     ASSERT(getpcid() == PCID_KERNEL);

573     switch (type) {
574     case FLUSH_TLB_ALL:
575         ASSERT(range == NULL);
576         invpcid(INVPCID_ALL_GLOBAL, 0, 0);
577         break;

579     case FLUSH_TLB_NONGLOBAL:
580         ASSERT(range == NULL);
581         invpcid(INVPCID_ALL_NONGLOBAL, 0, 0);
582         break;

584     case FLUSH_TLB_RANGE: {
585         mmu_flush_tlb_range(range->tr_va, TLB_RANGE_LEN(range),
586             LEVEL_SIZE(range->tr_level));
587         break;
588     }

590     default:
591         panic("invalid call mmu_flush_tlb(%d)", type);
592         break;
593     }
594 }

596 #endif /* ! __xpv */

```

```

*****
11762 Fri Apr 6 17:25:06 2018
new/usr/src/uts/i86pc/vm/kboot_mmu.c
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */

22 /*
23  * Copyright 2008 Sun Microsystems, Inc. All rights reserved.
24  * Use is subject to license terms.
25  *
26  * Copyright 2018 Joyent, Inc.
27  */

29 #include <sys/types.h>
30 #include <sys/system.h>
31 #include <sys/archsystem.h>
32 #include <sys/debug.h>
33 #include <sys/bootconf.h>
34 #include <sys/bootsvcs.h>
35 #include <sys/bootinfo.h>
36 #include <sys/mman.h>
37 #include <sys/cmn_err.h>
38 #include <sys/param.h>
39 #include <sys/machparam.h>
40 #include <sys/machsystem.h>
41 #include <sys/promif.h>
42 #include <sys/kobj.h>
43 #ifdef __xpv
44 #include <sys/hypervisor.h>
45 #endif
46 #include <vm/kboot_mmu.h>
47 #include <vm/hat_pte.h>
48 #include <vm/hat_i86.h>
49 #include <vm/seg_kmem.h>

51 #if 0
52 /*
53  * Joe's debug printing
54  */
55 #define DBG(x) \
56     bop_printf(NULL, "kboot_mmu.c: %s is %" PRIx64 "\n", #x, (uint64_t)(x));
57 #else
58 #define DBG(x) /* naught */
59 #endif

```

```

61 /*
62  * Page table and memory stuff.
63  */
64 static caddr_t window;
65 static caddr_t pte_to_window;

67 /*
68  * this are needed by mmu_init()
69  */
70 int kbm_nx_support = 0; /* NX bit in PTEs is in use */
71 int kbm_pae_support = 0; /* PAE is 64 bit Page table entries */
72 int kbm_pge_support = 0; /* PGE is Page table global bit enabled */
73 int kbm_largepage_support = 0;
74 uint_t kbm_nucleus_size = 0;

76 #define BOOT_SHIFT(1) (shift_amt[1])
77 #define BOOT_SZ(1) ((size_t)1 << BOOT_SHIFT(1))
78 #define BOOT_OFFSET(1) (BOOT_SZ(1) - 1)
79 #define BOOT_MASK(1) (~BOOT_OFFSET(1))

81 /*
82  * Initialize memory management parameters for boot time page table management
83  */
84 void
85 kbm_init(struct xboot_info *bi)
86 {
87     /*
88      * configure mmu information
89      */
90     kbm_nucleus_size = (uintptr_t)bi->bi_kseg_size;
91     kbm_largepage_support = bi->bi_use_largepage;
92     kbm_nx_support = bi->bi_use_nx;
93     kbm_pae_support = bi->bi_use_pae;
94     kbm_pge_support = bi->bi_use_pge;
95     window = bi->bi_pt_window;
96     DBG(window);
97     pte_to_window = bi->bi_pte_to_pt_window;
98     DBG(pte_to_window);
99     if (kbm_pae_support) {
100         shift_amt = shift_amt_pae;
101         ptes_per_table = 512;
102         pte_size = 8;
103         lpagesize = TWO_MEG;
104 #ifdef __amd64
105         top_level = 3;
106 #else
107         top_level = 2;
108 #endif
109     } else {
110         shift_amt = shift_amt_nopae;
111         ptes_per_table = 1024;
112         pte_size = 4;
113         lpagesize = FOUR_MEG;
114         top_level = 1;
115     }

117 #ifdef __xpv
118     xen_info = bi->bi_xen_start_info;
119     mfn_list = (mfn_t *)xen_info->mfn_list;
120     DBG(mfn_list);
121     mfn_count = xen_info->nr_pages;
122     DBG(mfn_count);
123 #endif
124     top_page_table = bi->bi_top_page_table;
125     DBG(top_page_table);

```



```

126 }

128 /*
129  * Change the addressible page table window to point at a given page
130  */
131 /*ARGUSED*/
132 void *
133 kbm_remap_window(paddr_t physaddr, int writeable)
134 {
135     x86pte_t pt_bits = PT_NOCONSIST | PT_VALID | PT_WRITABLE;

137     DBG(physaddr);

139 #ifdef __xpv
140     if (!writeable)
141         pt_bits &= ~PT_WRITABLE;
142     if (HYPERVISOR_update_va_mapping((uintptr_t)window,
143         pa_to_ma(physaddr) | pt_bits, UVMF_INVLPG | UVMF_LOCAL) < 0)
144         bop_panic("HYPERVISOR_update_va_mapping() failed");
145 #else
146     if (kbm_pae_support)
147         *((x86pte_t *)pte_to_window) = physaddr | pt_bits;
148     else
149         *((x86pte32_t *)pte_to_window) = physaddr | pt_bits;
150     mmu_invlpg(window);
151     mmu_tlbflush_entry(window);
152 #endif
153     DBG(window);
154     return (window);
155 }

156 /*
157  * Add a mapping for the physical page at the given virtual address.
158  */
159 void
160 kbm_map(uintptr_t va, paddr_t pa, uint_t level, uint_t is_kernel)
161 {
162     x86pte_t *ptep;
163     paddr_t pte_physaddr;
164     x86pte_t pteval;

166     if (khat_running)
167         panic("kbm_map() called too late");

169     pteval = pa_to_ma(pa) | PT_NOCONSIST | PT_VALID | PT_WRITABLE;
170     if (level >= 1)
171         pteval |= PT_PAGESIZE;
172     if (kbm_pge_support && is_kernel)
173         pteval |= PT_GLOBAL;

175 #ifdef __xpv
176     /*
177      * try update_va_mapping first - fails if page table is missing.
178      */
179     if (HYPERVISOR_update_va_mapping(va, pteval,
180         UVMF_INVLPG | UVMF_LOCAL) == 0)
181         return;
182 #endif

184     /*
185      * Find the pte that will map this address. This creates any
186      * missing intermediate level page tables.
187      */
188     ptep = find_pte(va, &pte_physaddr, level, 0);
189     if (ptep == NULL)
190         bop_panic("kbm_map: find_pte returned NULL");

```

```

192 #ifdef __xpv
193     if (HYPERVISOR_update_va_mapping(va, pteval, UVMF_INVLPG | UVMF_LOCAL))
194         bop_panic("HYPERVISOR_update_va_mapping() failed");
195 #else
196     if (kbm_pae_support)
197         *ptep = pteval;
198     else
199         *((x86pte32_t *)ptep) = pteval;
200     mmu_invlpg((caddr_t)va);
198     mmu_tlbflush_entry((caddr_t)va);
201 #endif
202 }
    unchanged_portion_omitted

329 /*
330  * Destroy a boot loader page table 4K mapping.
331  */
332 void
333 kbm_unmap(uintptr_t va)
334 {
335     if (khat_running)
336         panic("kbm_unmap() called too late");
337     else {
338 #ifdef __xpv
339         (void) HYPERVISOR_update_va_mapping(va, 0,
340             UVMF_INVLPG | UVMF_LOCAL);
341 #else
342         x86pte_t *ptep;
343         level_t level = 0;
344         uint_t probe_only = 1;

346         ptep = find_pte(va, NULL, level, probe_only);
347         if (ptep == NULL)
348             return;

350         if (kbm_pae_support)
351             *ptep = 0;
352         else
353             *((x86pte32_t *)ptep) = 0;
354         mmu_invlpg((caddr_t)va);
352         mmu_tlbflush_entry((caddr_t)va);
355 #endif
356     }
357 }

360 /*
361  * Change a boot loader page table 4K mapping.
362  * Returns the pfn of the old mapping.
363  */
364 pfn_t
365 kbm_remap(uintptr_t va, pfn_t pfn)
366 {
367     x86pte_t *ptep;
368     level_t level = 0;
369     uint_t probe_only = 1;
370     x86pte_t pte_val = pa_to_ma(pfn_to_pa(pfn)) | PT_WRITABLE |
371         PT_NOCONSIST | PT_VALID;
372     x86pte_t old_pte;

374     if (khat_running)
375         panic("kbm_remap() called too late");
376     ptep = find_pte(va, NULL, level, probe_only);
377     if (ptep == NULL)

```

```

378         bop_panic("kvm_remap: find_pte returned NULL");
380     if (kvm_pae_support)
381         old_pte = *ptep;
382     else
383         old_pte = *((x86pte32_t *)ptep);
385 #ifdef __xpv
386     if (HYPERVISOR_update_va_mapping(va, pte_val, UVMF_INVLPG | UVMF_LOCAL))
387         bop_panic("HYPERVISOR_update_va_mapping() failed");
388 #else
389     if (kvm_pae_support)
390         *((x86pte_t *)ptep) = pte_val;
391     else
392         *((x86pte32_t *)ptep) = pte_val;
393     mmu_invlpg((caddr_t)va);
394     mmu_tlbflush_entry((caddr_t)va);
395 #endif
396     if (!(old_pte & PT_VALID) || ma_to_pa(old_pte) == -1)
397         return (PFN_INVALID);
398     return (mmu_btop(ma_to_pa(old_pte)));
399 }
402 /*
403  * Change a boot loader page table 4K mapping to read only.
404  */
405 void
406 kvm_read_only(uintptr_t va, paddr_t pa)
407 {
408     x86pte_t pte_val = pa_to_ma(pa) |
409         PT_NOCONSIST | PT_REF | PT_MOD | PT_VALID;
411 #ifdef __xpv
412     if (HYPERVISOR_update_va_mapping(va, pte_val, UVMF_INVLPG | UVMF_LOCAL))
413         bop_panic("HYPERVISOR_update_va_mapping() failed");
414 #else
415     x86pte_t *ptep;
416     level_t level = 0;
418     ptep = find_pte(va, NULL, level, 0);
419     if (ptep == NULL)
420         bop_panic("kvm_read_only: find_pte returned NULL");
422     if (kvm_pae_support)
423         *ptep = pte_val;
424     else
425         *((x86pte32_t *)ptep) = pte_val;
426     mmu_invlpg((caddr_t)va);
427     mmu_tlbflush_entry((caddr_t)va);
428 #endif
429 }
430
431 #ifdef __xpv
432 void
433 kvm_pop(void)
434 {
435     #ifdef __xpv
436     if (HYPERVISOR_update_va_mapping((uintptr_t)window, save_pte,
437         UVMF_INVLPG | UVMF_LOCAL) < 0)
438         bop_panic("HYPERVISOR_update_va_mapping() failed");
439     #else
440     if (kvm_pae_support)
441         *((x86pte_t *)pte_to_window) = save_pte;
442     else

```

```

463         *((x86pte32_t *)pte_to_window) = save_pte;
464     mmu_invlpg(window);
465     mmu_tlbflush_entry(window);
466 #endif
467 }
468
469 #ifdef __xpv
470 void
471 kvm_pop(void)
472 {
473     #ifdef __xpv
474     if (HYPERVISOR_update_va_mapping((uintptr_t)window, save_pte,
475         UVMF_INVLPG | UVMF_LOCAL) < 0)
476         bop_panic("HYPERVISOR_update_va_mapping() failed");
477     #else
478     if (kvm_pae_support)
479         *((x86pte_t *)pte_to_window) = save_pte;
480     else

```

```

*****
99871 Fri Apr 6 17:25:06 2018
new/usr/src/uts/i86pc/vm/vm_machdep.c
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22 * Copyright (c) 1992, 2010, Oracle and/or its affiliates. All rights reserved.
23 */
24 /*
25 * Copyright (c) 2010, Intel Corporation.
26 * All rights reserved.
27 * Copyright 2018 Joyent, Inc.
27 * Copyright 2016 Joyent, Inc.
28 */

30 /* Copyright (c) 1984, 1986, 1987, 1988, 1989 AT&T */
31 /* All Rights Reserved */

33 /*
34 * Portions of this source code were derived from Berkeley 4.3 BSD
35 * under license from the Regents of the University of California.
36 */

38 /*
39 * UNIX machine dependent virtual memory support.
40 */

42 #include <sys/types.h>
43 #include <sys/param.h>
44 #include <sys/system.h>
45 #include <sys/user.h>
46 #include <sys/proc.h>
47 #include <sys/kmem.h>
48 #include <sys/vmem.h>
49 #include <sys/buf.h>
50 #include <sys/cpuvar.h>
51 #include <sys/lgrp.h>
52 #include <sys/disp.h>
53 #include <sys/vm.h>
54 #include <sys/mman.h>
55 #include <sys/vnode.h>
56 #include <sys/cred.h>
57 #include <sys/exec.h>
58 #include <sys/exechdr.h>

```

```

59 #include <sys/debug.h>
60 #include <sys/vmsystem.h>
61 #include <sys/swap.h>
62 #include <sys/dumphdr.h>
63 #include <sys/random.h>

65 #include <vm/hat.h>
66 #include <vm/as.h>
67 #include <vm/seg.h>
68 #include <vm/seg_kp.h>
69 #include <vm/seg_vn.h>
70 #include <vm/page.h>
71 #include <vm/seg_kmem.h>
72 #include <vm/seg_kpm.h>
73 #include <vm/vm_dep.h>

75 #include <sys/cpu.h>
76 #include <sys/vm_machparam.h>
77 #include <sys/memlist.h>
78 #include <sys/bootconf.h> /* XXX the memlist stuff belongs in memlist_plat.h */
79 #include <vm/hat_i86.h>
80 #include <sys/x86_archext.h>
81 #include <sys/elf_386.h>
82 #include <sys/cmn_err.h>
83 #include <sys/archsystem.h>
84 #include <sys/machsystem.h>
85 #include <sys/secflags.h>

87 #include <sys/vtrace.h>
88 #include <sys/ddidmareq.h>
89 #include <sys/promif.h>
90 #include <sys/memnode.h>
91 #include <sys/stack.h>
92 #include <util/qsorth.h>
93 #include <sys/taskq.h>

95 #ifdef __xpv

97 #include <sys/hypervisor.h>
98 #include <sys/xen_mmu.h>
99 #include <sys/balloon_impl.h>

101 /*
102 * domain 0 pages usable for DMA are kept pre-allocated and kept in
103 * distinct lists, ordered by increasing mfn.
104 */
105 static kmutex_t io_pool_lock;
106 static kmutex_t contig_list_lock;
107 static page_t *io_pool_4g; /* pool for 32 bit dma limited devices */
108 static page_t *io_pool_16m; /* pool for 24 bit dma limited legacy devices */
109 static long io_pool_cnt;
110 static long io_pool_cnt_max = 0;
111 #define DEFAULT_IO_POOL_MIN 128
112 static long io_pool_cnt_min = DEFAULT_IO_POOL_MIN;
113 static long io_pool_cnt_lowater = 0;
114 static long io_pool_shrink_attempts; /* how many times did we try to shrink */
115 static long io_pool_shrinks; /* how many times did we really shrink */
116 static long io_pool_grows; /* how many times did we grow */
117 static mfn_t start_mfn = 1;
118 static caddr_t io_pool_kva; /* use to alloc pages when needed */

120 static int create_contig_pfnlist(uint_t);

122 /*
123 * percentage of phys mem to hold in the i/o pool
124 */

```

```

125 #define DEFAULT_IO_POOL_PCT 2
126 static long io_pool_physmem_pct = DEFAULT_IO_POOL_PCT;
127 static void page_io_pool_sub(page_t **, page_t *, page_t *);
128 int ioalloc_dbg = 0;

130 #endif /* __xpv */

132 uint_t vac_colors = 1;

134 int largepagesupport = 0;
135 extern uint_t page_create_new;
136 extern uint_t page_create_exists;
137 extern uint_t page_create_putbacks;
138 /*
139  * Allow users to disable the kernel's use of SSE.
140  */
141 extern int use_sse_pagecopy, use_sse_pagezero;

143 /*
144  * combined memory ranges from mnode and memranges[] to manage single
145  * mnode/mtype dimension in the page lists.
146  */
147 typedef struct {
148     pfn_t     mnr_pfnlo;
149     pfn_t     mnr_pfnhi;
150     int       mnr_mnode;
151     int       mnr_memrange; /* index into memranges[] */
152     int       mnr_next; /* next lower PA mnode range */
153     int       mnr_exists;
154     /* maintain page list stats */
155     pgcnt_t  mnr_mt_clpgcnt; /* cache list cnt */
156     pgcnt_t  mnr_mt_flpgcnt[MMU_PAGE_SIZES]; /* free list cnt per szc */
157     pgcnt_t  mnr_mt_totcnt; /* sum of cache and free lists */
158 #ifdef DEBUG
159     struct mnr_mts { /* mnode/mtype szc stats */
160         pgcnt_t  mnr_mts_pgcnt;
161         int       mnr_mts_colors;
162         pgcnt_t  *mnr_mts_pgcnt;
163     } *mnr_mts;
164 #endif
165 } mnode_range_t;
166 #ifndef unchanged_portion_omitted
290 #endif

292 uint_t mmu_page_sizes;

294 /* How many page sizes the users can see */
295 uint_t mmu_exported_page_sizes;

297 /* page sizes that legacy applications can see */
298 uint_t mmu_legacy_page_sizes;

300 /*
301  * Number of pages in 1 GB. Don't enable automatic large pages if we have
302  * fewer than this many pages.
303  */
304 pgcnt_t shm_lpg_min_physmem = 1 << (30 - MMU_PAGESHIFT);
305 pgcnt_t privm_lpg_min_physmem = 1 << (30 - MMU_PAGESHIFT);

307 /*
308  * Maximum and default segment size tunables for user private
309  * and shared anon memory, and user text and initialized data.
310  * These can be patched via /etc/system to allow large pages
311  * to be used for mapping application private and shared anon memory.
312  */
313 size_t mcntl0_lpsize = MMU_PAGESIZE;

```

```

314 size_t max_uheap_lpsize = MMU_PAGESIZE;
315 size_t default_uheap_lpsize = MMU_PAGESIZE;
316 size_t max_ustack_lpsize = MMU_PAGESIZE;
317 size_t default_ustack_lpsize = MMU_PAGESIZE;
318 size_t max_privmap_lpsize = MMU_PAGESIZE;
319 size_t max_uidata_lpsize = MMU_PAGESIZE;
320 size_t max_utext_lpsize = MMU_PAGESIZE;
321 size_t max_shm_lpsize = MMU_PAGESIZE;

324 /*
325  * initialized by page_coloring_init().
326  */
327 uint_t page_colors;
328 uint_t page_colors_mask;
329 uint_t page_coloring_shift;
330 int cpu_page_colors;
331 static uint_t l2_colors;

333 /*
334  * Page freelists and cachelists are dynamically allocated once mnode_rangecnt
335  * and page_colors are calculated from the l2 cache n-way set size. Within a
336  * mnode range, the page freelist and cachelist are hashed into bins based on
337  * color. This makes it easier to search for a page within a specific memory
338  * range.
339  */
340 #define PAGE_COLORS_MIN 16

342 page_t ***page_freelists;
343 page_t ***page_cachelists;

346 /*
347  * Used by page layer to know about page sizes
348  */
349 hw_page_size_t hw_page_array[MAX_NUM_LEVEL + 1];

351 kmutex_t *fpc_mutex[NPC_MUTEX];
352 kmutex_t *cpc_mutex[NPC_MUTEX];

354 /* Lock to protect mnode_ranges array for memory DR operations. */
355 static kmutex_t mnode_range_lock;

357 /*
358  * Only let one thread at a time try to coalesce large pages, to
359  * prevent them from working against each other.
360  */
361 static kmutex_t contig_lock;
362 #define CONTIG_LOCK() mutex_enter(&contig_lock);
363 #define CONTIG_UNLOCK() mutex_exit(&contig_lock);

365 #define PFN_16M (mmu_btop((uint64_t)0x1000000))

367 caddr_t
368 i86devmap(pfn_t pf, pgcnt_t pgcnt, uint_t prot)
369 {
370     caddr_t addr;
371     caddr_t addr1;
372     page_t *pp;

374     addr1 = addr = vmem_alloc(heap_arena, mmu_ptob(pgcnt), VM_SLEEP);

376     for (; pgcnt != 0; addr += MMU_PAGESIZE, ++pf, --pgcnt) {
377         pp = page_numtopp_nolock(pf);
378         if (pp == NULL) {
379             hat_devload(kas.a_hat, addr, MMU_PAGESIZE, pf,

```

```

380         prot | HAT_NOSYNC, HAT_LOAD_LOCK);
381     } else {
382         hat_memload(kas.a_hat, addr, pp,
383         prot | HAT_NOSYNC, HAT_LOAD_LOCK);
384     }
385 }
387     return (addr1);
388 }

390 /*
391  * This routine is like page_numtopp, but accepts only free pages, which
392  * it allocates (unfrees) and returns with the exclusive lock held.
393  * It is used by machdep.c/dma_init() to find contiguous free pages.
394  */
395 page_t *
396 page_numtopp_alloc(pfn_t pfnnum)
397 {
398     page_t *pp;

400 retry:
401     pp = page_numtopp_nolock(pfnnum);
402     if (pp == NULL) {
403         return (NULL);
404     }

406     if (!page_trylock(pp, SE_EXCL)) {
407         return (NULL);
408     }

410     if (page_pptonum(pp) != pfnnum) {
411         page_unlock(pp);
412         goto retry;
413     }

415     if (!PP_ISFREE(pp)) {
416         page_unlock(pp);
417         return (NULL);
418     }
419     if (pp->p_szc) {
420         page_demote_free_pages(pp);
421         page_unlock(pp);
422         goto retry;
423     }

425     /* If associated with a vnode, destroy mappings */

427     if (pp->p_vnode) {

429         page_destroy_free(pp);

431         if (!page_lock(pp, SE_EXCL, (kmutex_t *)NULL, P_NO_RECLAIM)) {
432             return (NULL);
433         }

435         if (page_pptonum(pp) != pfnnum) {
436             page_unlock(pp);
437             goto retry;
438         }
439     }

441     if (!PP_ISFREE(pp)) {
442         page_unlock(pp);
443         return (NULL);
444     }

```

```

446     if (!page_reclaim(pp, (kmutex_t *)NULL))
447         return (NULL);

449     return (pp);
450 }

452 /*
453  * Return the optimum page size for a given mapping
454  */
455 /*ARGSUSED*/
456 size_t
457 map_pgsz(int matype, struct proc *p, caddr_t addr, size_t len, int memcntl)
458 {
459     level_t l = 0;
460     size_t pgsz = MMU_PAGESIZE;
461     size_t max_lpsize;
462     uint_t mszc;

464     ASSERT(matype != MAPPGSZ_VA);

466     if (matype != MAPPGSZ_ISM && phymem < privm_lpg_min_phymem) {
467         return (MMU_PAGESIZE);
468     }

470     switch (matype) {
471     case MAPPGSZ_HEAP:
472     case MAPPGSZ_STK:
473         max_lpsize = memcntl ? mcntl0_lpsize : (matype ==
474         MAPPGSZ_HEAP ? max_uheap_lpsize : max_ustack_lpsize);
475         if (max_lpsize == MMU_PAGESIZE) {
476             return (MMU_PAGESIZE);
477         }
478         if (len == 0) {
479             len = (matype == MAPPGSZ_HEAP ? p->p_brkbase +
480             p->p_brksize - p->p_bssbase : p->p_stksize;
481         }
482         len = (matype == MAPPGSZ_HEAP ? MAX(len,
483         default_uheap_lpsize) : MAX(len, default_ustack_lpsize));

485     /*
486      * use the pages size that best fits len
487      */
488     for (l = mmu.umax_page_level; l > 0; --l) {
489         if (LEVEL_SIZE(l) > max_lpsize || len < LEVEL_SIZE(l)) {
490             continue;
491         } else {
492             pgsz = LEVEL_SIZE(l);
493         }
494         break;
495     }

497     mszc = (matype == MAPPGSZ_HEAP ? p->p_brkpageszc :
498     p->p_stkpageszc);
499     if (addr == 0 && (pgsz < hw_page_array[mszc].hp_size)) {
500         pgsz = hw_page_array[mszc].hp_size;
501     }
502     return (pgsz);

504     case MAPPGSZ_ISM:
505         for (l = mmu.umax_page_level; l > 0; --l) {
506             if (len >= LEVEL_SIZE(l))
507                 return (LEVEL_SIZE(l));
508         }
509         return (LEVEL_SIZE(0));
510     }
511     return (pgsz);

```

new/usr/src/uts/i86pc/vm/vm_machdep.c

7

512 }

unchanged_portion_omitted

```

*****
5573 Fri Apr 6 17:25:07 2018
new/usr/src/uts/i86xpv/Makefile.files
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>
*****
1 #
2 # CDDL HEADER START
3 #
4 # The contents of this file are subject to the terms of the
5 # Common Development and Distribution License (the "License").
6 # You may not use this file except in compliance with the License.
7 #
8 # You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9 # or http://www.opensolaris.org/os/licensing.
10 # See the License for the specific language governing permissions
11 # and limitations under the License.
12 #
13 # When distributing Covered Code, include this CDDL HEADER in each
14 # file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 # If applicable, add the following below this CDDL HEADER, with the
16 # fields enclosed by brackets "[]" replaced with your own identifying
17 # information: Portions Copyright [yyyy] [name of copyright owner]
18 #
19 # CDDL HEADER END
20 #
21 #
22 #
23 # Copyright (c) 2007, 2010, Oracle and/or its affiliates. All rights reserved.
24 #
25 # Copyright 2018 Joyent, Inc.
26 #
27 #
28 # This Makefile defines file modules in the directory uts/i86xpv
29 # and its children. These are the source files which are i86xpv
30 # "implementation architecture" dependent.
31 #
32 #
33 #
34 # object lists
35 #
36 CORE_OBJJS += \
37     acpi_stubs.o \
38     balloon.o \
39     biosdisk.o \
40     cbe.o \
41     cmi.o \
42     cmi_hw.o \
43     cms.o \
44     confunix.o \
45     cpuid.o \
46     cpuid_subr.o \
47     cpupm.o \
48     cpupm_mach.o \
49     dis_tables.o \
50     ddi_impl.o \
51     dtrace_subr.o \
52     dvma.o \
53     fakebop.o \
54     fpu_subr.o \
55     fastboot.o \
56     fb_swtd.o \
57     graphics.o \
58     hardclk.o \
59     hat_i86.o \

```

```

60     hat_kdi.o \
61     hment.o \
62     hold_page.o \
63     hrtimers.o \
64     htable.o \
65     i86_mmu.o \
66     ibft.o \
67     instr_size.o \
68     intr.o \
69     kboot_mmu.o \
70     kdi_subr.o \
71     kdi_idt.o \
72     kdi_idthdl.o \
73     kdi_asm.o \
74     lgrplat.o \
75     mach_kdi.o \
76     mach_sysconfig.o \
77     machdep.o \
78     mem_config_stubs.o \
79     memnode.o \
80     microcode.o \
81     mlsetup.o \
82     mp_call.o \
83     mp_implfuncs.o \
84     mp_machdep.o \
85     mp_startup.o \
86     memscrub.o \
87     notes.o \
88     pci_bios.o \
89     pci_cfgacc.o \
90     pci_cfgacc_x86.o \
91     pci_cfgspace.o \
92     pci_mech1.o \
93     pci_mech2.o \
94     pci_neptune.o \
95     pci_orion.o \
96     pmem.o \
97     ppage.o \
98     startup.o \
99     ssp.o \
100     xpv_timestamp.o \
101     todpc_subr.o \
102     trap.o \
103     vm_machdep.o \
104     x_call.o \
105 #
106 # Add the SMBIOS subsystem object files directly to the list of objects
107 # built into unix itself; this is all common code except for smb_dev.c.
108 #
109 CORE_OBJJS += $(SMBIOS_OBJJS)
110 #
111 #
112 # These get compiled twice:
113 # - once in the dboot (direct boot) identity mapped code
114 # - once for use during early startup in unix
115 #
116 BOOT_DRIVER_OBJJS = \
117     boot_console.o \
118     boot_keyboard.o \
119     boot_keyboard_table.o \
120     boot_mmu.o \
121     boot_vga.o \
122     boot_xconsole.o \
123     dboot_multiboot2.o \
124     $(FONT_OBJJS)

```

```

126 CORE_OBJS += $(BOOT_DRIVER_OBJS)

128 #
129 # Extra XEN files separated out for now.
130 #
131 CORE_OBJS += \
132     cpr_driver.o \
133     evtchn.o \
134     gnttab.o \
135     hypercall.o \
136     hyperevent.o \
137     hypersubr.o \
138     mp_xen.o \
139     panic_asm.o \
140     xenguest.o \
141     xenbus_client.o \
142     xenbus_comms.o \
143     xenbus_probe.o \
144     xenbus_xs.o \
145     xen_machdep.o \
146     xen_mmu.o \
147     xpv_panic.o \
148     xvdi.o

150 #
151 #     locore.o is special. It must be the first file relocated so that it
152 #     it is relocated just where its name implies.
153 #
154 SPECIAL_OBJS_32 += \
155     locore.o \
156     fast_trap_asm.o \
157     interrupt.o \
158     syscall_asm.o

160 SPECIAL_OBJS_64 += \
161     locore.o \
162     fast_trap_asm.o \
163     interrupt.o \
164     syscall_asm_amd64.o \
165     kpti_trampoline.o \
166     syscall_asm_amd64.o

167 SPECIAL_OBJS += $(SPECIAL_OBJS_$(CLASS))

169 #
170 # object files used to boot into full kernel
171 #
172 DBOOT_OBJS_32 = muldiv.o

174 DBOOT_OBJS_64 =

176 DBOOT_OBJS += \
177     dboot_asm.o \
178     dboot_printf.o \
179     dboot_startkern.o \
180     dboot_xen.o \
181     hypercall.o \
182     hypersubr.o \
183     memcpy.o \
184     memset.o \
185     string.o \
186     $(BOOT_DRIVER_OBJS) \
187     $(DBOOT_OBJS_$(CLASS))

189 #

```

```

190 #             driver & misc modules
191 #
192 BALLOON_OBJS += balloon_drv.o
193 DOMCAPS_OBJS += domcaps.o
194 EVTCHN_OBJS += evtchn_dev.o
195 GFX_PRIVATE_OBJS += gfx_private.o gfxp_pci.o gfxp_segmap.o \
196     gfxp_devmap.o gfxp_vgtext.o gfxp_vm.o vgasubr.o
197 IOAT_OBJS += ioat.o ioat_rs.o ioat_ioctl.o ioat_chan.o
198 ISANEXUS_OBJS += isa.o dma_engine.o i8237A.o
199 PCI_E_NEXUS_OBJS += npe.o npe_misc.o
200 PCI_E_NEXUS_OBJS += pci_common.o pci_kstats.o pci_tools.o
201 PCINEXUS_OBJS += pci.o pci_common.o pci_kstats.o pci_tools.o
202 PRIVCMD_OBJS += seg_mf.o privcmd.o privcmd_hcall.o
203 ROOTNEX_OBJS += rootnex.o
204 XPVTOD_OBJS += xpvtod.o
205 XPV_AUTOCONFIG_OBJS += xpv_autoconfig.o
206 XPV_PSM_OBJS += xpv_psm.o mp_platform_common.o mp_platform_xpv.o \
207     apic_regops.o psm_common.o xpv_intr.o
208 XPV_UPPC_OBJS += xpv_uppc.o psm_common.o
209 XENBUS_OBJS += xenbus_dev.o
210 XENCONS_OBJS += xencons.o
211 XPVD_OBJS += xpvd.o
212 XPVTAP_OBJS += xpvtap.o blk_common.o seg_mf.o
213 XNB_OBJS += xnb.o
214 XNBE_OBJS += xnbe.o
215 XNBO_OBJS += xnbo.o
216 XNBU_OBJS += xnbu.o
217 XNF_OBJS += xnf.o
218 XSVC_OBJS += xsvc.o
219 XDF_OBJS += xdf.o
220 XDB_OBJS += xdb.o
221 XDT_OBJS += xdt.o

223 #
224 #             Build up defines and paths.
225 #
226 INC_PATH += -I$(UTSBASE)/i86xpv -I$(UTSBASE)/i86pc -I$(SRC)/common \
227     -I$(UTSBASE)/common/xen

229 #
230 # Since the assym files are derived, the dependencies must be explicit for
231 # all files including this file. (This is only actually required in the
232 # instance when the .nse_depinfo file does not exist.) It may seem that
233 # the lint targets should also have a similar dependency, but they don't
234 # since only C headers are included when #defined(__lint) is true.
235 #

237 ASSYM_DEPS += \
238     copy.o \
239     desctbls_asm.o \
240     ddi_i86_asm.o \
241     exception.o \
242     fast_trap_asm.o \
243     float.o \
244     hyperevent.o \
245     i86_subr.o \
246     kdi_asm.o \
247     interrupt.o \
248     lock_prim.o \
249     locore.o \
250     panic_asm.o \
251     sseblk.o \
252     swtch.o \
253     syscall_asm.o \
254     syscall_asm_amd64.o

```


new/usr/src/uts/i86xpv/Makefile.files

5

```
256 $(KDI_ASSYM_DEPS:%=$(OBJSDIR)/%): $(DSFDIR)/$(OBJSDIR)/kdi_assym.h
```

```
256 ASSYM_DEPS += kdi_asm.o
```

new/usr/src/uts/i86xpv/os/xpv_panic.c

1

```
*****
27893 Fri Apr 6 17:25:07 2018
new/usr/src/uts/i86xpv/os/xpv_panic.c
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22 * Copyright (c) 2012 Gary Mills
23 * Copyright 2016 PALO, Richard.
24 *
25 * Copyright (c) 2010, Oracle and/or its affiliates. All rights reserved.
26 *
27 * Copyright 2018 Joyent, Inc.
28 */

30 #include <sys/types.h>
31 #include <sys/clock.h>
32 #include <sys/psm.h>
33 #include <sys/archsystem.h>
34 #include <sys/machsystem.h>
35 #include <sys/compress.h>
36 #include <sys/modctl.h>
37 #include <sys/trap.h>
38 #include <sys/panic.h>
39 #include <sys/regset.h>
40 #include <sys/frame.h>
41 #include <sys/kobj.h>
42 #include <sys/apic.h>
43 #include <sys/apic_timer.h>
44 #include <sys/dumphdr.h>
45 #include <sys/mem.h>
46 #include <sys/x86_archext.h>
47 #include <sys/xpv_panic.h>
48 #include <sys/boot_console.h>
49 #include <sys/bootsvcs.h>
50 #include <sys/consdev.h>
51 #include <vm/hat_pte.h>
52 #include <vm/hat_i86.h>

54 /* XXX: need to add a PAE version too, if we ever support both PAE and non */
55 #if defined(__i386)
56 #define XPV_FILENAME "/boot/xen-syms"
57 #else
58 #define XPV_FILENAME "/boot/amd64/xen-syms"
59 #endif
```

new/usr/src/uts/i86xpv/os/xpv_panic.c

2

```
60 #define XPV_MODNAME "xpv"

62 int xpv_panicking = 0;

64 struct module *xpv_module;
65 struct modctl *xpv_modctl;

67 #define ALIGN(x, a) ((a) == 0 ? (uintptr_t)(x) : \
68 ((uintptr_t)(x) + (uintptr_t)(a) - 1) & ~(uintptr_t)(a) - 1))

70 /* Pointer to the xpv_panic_info structure handed to us by Xen. */
71 static struct panic_info *xpv_panic_info = NULL;

73 /* Timer support */
74 #define NSEC_SHIFT 5
75 #define T_XPV_TIMER 0xd1
76 #define XPV_TIMER_INTERVAL 1000 /* 1000 microseconds */
77 static uint32_t *xpv_apicadr = NULL;
78 static uint_t nsec_scale;

80 /* IDT support */
81 #pragma align 16(xpv_panic_idt)
82 static gate_desc_t xpv_panic_idt[NIDT]; /* interrupt descriptor table */

84 /* Xen pagetables mapped into our HAT's ptable windows */
85 static pfn_t ptable_pfn[MAX_NUM_LEVEL];

87 /* Number of MMU_PAGESIZE pages we're adding to the Solaris dump */
88 static int xpv_dump_pages;

90 /*
91 * There are up to two large swathes of RAM that we don't want to include
92 * in the dump: those that comprise the Xen version of segkpm. On 32-bit
93 * systems there is no such region of memory. On 64-bit systems, there
94 * should be just a single contiguous region that corresponds to all of
95 * physical memory. The tricky bit is that Xen's heap sometimes lives in
96 * the middle of their segkpm, and is mapped using only kpm-like addresses.
97 * In that case, we need to skip the swathes before and after Xen's heap.
98 */
99 uintptr_t kpm1_low = 0;
100 uintptr_t kpm1_high = 0;
101 uintptr_t kpm2_low = 0;
102 uintptr_t kpm2_high = 0;

104 /*
105 * Some commonly used values that we don't want to recompute over and over.
106 */
107 static int xpv_panic_nptes[MAX_NUM_LEVEL];
108 static ulong_t xpv_panic_cr3;
109 static uintptr_t xpv_end;

111 static void xpv_panic_console_print(const char *fmt, ...);
112 static void (*xpv_panic_printf)(const char *, ...) = xpv_panic_console_print;

114 #define CONSOLE_BUF_SIZE 256
115 static char console_buffer[CONSOLE_BUF_SIZE];
116 static boolean_t use_polledio;

118 /*
119 * Pointers to machine check panic info (if any).
120 */
121 xpv_mca_panic_data_t *xpv_mca_panic_data = NULL;

123 static void
124 xpv_panic_putc(int m)
125 {
```

```
126     struct cons_polledio *c = cons_polledio;
128     /* This really shouldn't happen */
129     if (boot_console_type(NULL) == CONS_HYPERVISOR)
130         return;
132     if (use_polledio == B_TRUE)
133         c->cons_polledio_putchar(c->cons_polledio_argument, m);
134     else
135         bcons_putchar(m);
136 }
unchanged_portion_omitted

160 static void
161 xpv_panic_map(int level, pfn_t pfn)
162 {
163     x86pte_t pte, *pteptr;
165     /*
166      * The provided pfn represents a level 'level' page table. Map it
167      * into the 'level' slot in the list of page table windows.
168      */
169     pteptr = (x86pte_t *)PWIN_PTE_VA(level);
170     pte = pfn_to_pa(pfn) | PT_VALID;
172     XPV_ALLOW_PAGETABLE_UPDATES();
173     if (mmu.pae_hat)
174         *pteptr = pte;
175     else
176         *(x86pte32_t *)pteptr = pte;
177     XPV_DISALLOW_PAGETABLE_UPDATES();
179     mmu_flush_tlb_page((uintptr_t)PWIN_VA(level));
177     mmu_tlbflush_entry(PWIN_VA(level));
180 }
unchanged_portion_omitted
```

```

*****
13843 Fri Apr 6 17:25:08 2018
new/usr/src/uts/intel/Makefile.rules
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>
*****
1 #
2 # CDDL HEADER START
3 #
4 # The contents of this file are subject to the terms of the
5 # Common Development and Distribution License (the "License").
6 # You may not use this file except in compliance with the License.
7 #
8 # You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9 # or http://www.opensolaris.org/os/licensing.
10 # See the License for the specific language governing permissions
11 # and limitations under the License.
12 #
13 # When distributing Covered Code, include this CDDL HEADER in each
14 # file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 # If applicable, add the following below this CDDL HEADER, with the
16 # fields enclosed by brackets "[]" replaced with your own identifying
17 # information: Portions Copyright [yyyy] [name of copyright owner]
18 #
19 # CDDL HEADER END
20 #
21 #
22 # Copyright 2009 Sun Microsystems, Inc. All rights reserved.
23 # Use is subject to license terms.
24 # Copyright 2018 Joyent, Inc. All rights reserved.
24 # Copyright 2012 Joyent, Inc. All rights reserved.
25 # Copyright 2017 Nexenta Systems, Inc.
26 #

28 #
29 # This Makefile defines all file modules and build rules for the
30 # directory uts/intel and its children. These are the source files which
31 # are specific to the intel processor.
32 #
33 # The following two-level ordering must be maintained in this file.
34 # Lines are sorted first in order of decreasing specificity based on
35 # the first directory component. That is, sun4u rules come before
36 # sparc rules come before common rules.
37 #
38 # Lines whose initial directory components are equal are sorted
39 # alphabetically by the remaining components.

41 #
42 # Need a way to distinguish between the ia32 and amd64 subdirs.
43 #
44 SUBARCH_DIR_32 = ia32
45 SUBARCH_DIR_64 = amd64
46 SUBARCH_DIR = $(SUBARCH_DIR_$(CLASS))

48 #
49 # Section 1a: C object build rules
50 #
51 $(OBJDIR)/%.o: $(SRC)/common/fs/%.c
52 $(COMPILE.c) -o $@ $<
53 $(CTFCONVERT_O)

55 $(OBJDIR)/%.o: $(UTSBASE)/common/io/power/%.c
56 $(COMPILE.c) -o $@ $<
57 $(CTFCONVERT_O)

```

```

59 $(OBJDIR)/%.o: $(SRC)/common/util/i386/%.s
60 $(COMPILE.s) -o $@ $<

62 $(OBJDIR)/%.o: $(UTSBASE)/intel/brand/sn1/%.s
63 $(COMPILE.s) -o $@ $<

65 $(OBJDIR)/%.o: $(UTSBASE)/intel/brand/solaris10/%.s
66 $(COMPILE.s) -o $@ $<

68 $(OBJDIR)/%.o: $(UTSBASE)/intel/dtrace/%.c
69 $(COMPILE.c) -o $@ $<
70 $(CTFCONVERT_O)

72 $(OBJDIR)/%.o: $(UTSBASE)/intel/dtrace/%.s
73 $(COMPILE.s) -o $@ $<

75 $(OBJDIR)/%.o: $(UTSBASE)/intel/fs/proc/%.c
76 $(COMPILE.c) -o $@ $<
77 $(CTFCONVERT_O)

79 $(OBJDIR)/%.o: $(UTSBASE)/intel/ia32/ml/%.s
80 $(COMPILE.s) -o $@ $<

82 $(OBJDIR)/%.o: $(UTSBASE)/intel/ia32/os/%.c
83 $(COMPILE.c) -o $@ $<
84 $(CTFCONVERT_O)

86 $(OBJDIR)/%.o: $(UTSBASE)/intel/ia32/promif/%.c
87 $(COMPILE.c) -o $@ $<
88 $(CTFCONVERT_O)

90 $(OBJDIR)/%.o: $(UTSBASE)/intel/ia32/syscall/%.c
91 $(COMPILE.c) -o $@ $<
92 $(CTFCONVERT_O)

94 $(OBJDIR)/%.o: $(UTSBASE)/intel/io/%.c
95 $(COMPILE.c) -o $@ $<
96 $(CTFCONVERT_O)

98 $(OBJDIR)/%.o: $(UTSBASE)/intel/io/acpica/%.c
99 $(COMPILE.c) -o $@ $<
100 $(CTFCONVERT_O)

102 $(OBJDIR)/%.o: $(UTSBASE)/intel/io/acpica/%.s
103 $(COMPILE.s) -o $@ $<

105 $(OBJDIR)/%.o: $(UTSBASE)/intel/io/acpica/debugger/%.c
106 $(COMPILE.c) -o $@ $<
107 $(CTFCONVERT_O)

109 $(OBJDIR)/%.o: $(UTSBASE)/intel/io/acpica/events/%.c
110 $(COMPILE.c) -o $@ $<
111 $(CTFCONVERT_O)

113 $(OBJDIR)/%.o: $(UTSBASE)/intel/io/acpica/hardware/%.c
114 $(COMPILE.c) -o $@ $<
115 $(CTFCONVERT_O)

117 $(OBJDIR)/%.o: $(UTSBASE)/intel/io/acpica/dispatcher/%.c
118 $(COMPILE.c) -o $@ $<
119 $(CTFCONVERT_O)

121 $(OBJDIR)/%.o: $(UTSBASE)/intel/io/acpica/executer/%.c
122 $(COMPILE.c) -o $@ $<
123 $(CTFCONVERT_O)

```

new/usr/src/uts/intel/Makefile.rules

3

```

125 $(OBJSDIR)/%.o: $(UTSBASE)/intel/io/acpica/parser/%.c
126 $(COMPILE.c) -o $@ $<
127 $(CTFCONVERT_O)

129 $(OBJSDIR)/%.o: $(UTSBASE)/intel/io/acpica/namespace/%.c
130 $(COMPILE.c) -o $@ $<
131 $(CTFCONVERT_O)

133 $(OBJSDIR)/%.o: $(UTSBASE)/intel/io/acpica/resources/%.c
134 $(COMPILE.c) -o $@ $<
135 $(CTFCONVERT_O)

137 $(OBJSDIR)/%.o: $(UTSBASE)/intel/io/acpica/tables/%.c
138 $(COMPILE.c) -o $@ $<
139 $(CTFCONVERT_O)

141 $(OBJSDIR)/%.o: $(UTSBASE)/intel/io/acpica/utilities/%.c
142 $(COMPILE.c) -o $@ $<
143 $(CTFCONVERT_O)

145 $(OBJSDIR)/%.o: $(UTSBASE)/intel/io/acpica/disassembler/%.c
146 $(COMPILE.c) -o $@ $<
147 $(CTFCONVERT_O)

149 $(OBJSDIR)/%.o: $(UTSBASE)/intel/io/agpgart/%.c
150 $(COMPILE.c) -o $@ $<
151 $(CTFCONVERT_O)

153 $(OBJSDIR)/%.o: $(UTSBASE)/intel/io/agpmaster/%.c
154 $(COMPILE.c) -o $@ $<
155 $(CTFCONVERT_O)

157 $(OBJSDIR)/%.o: $(UTSBASE)/intel/io/amd8111s/%.c
158 $(COMPILE.c) -o $@ $<
159 $(CTFCONVERT_O)

161 $(OBJSDIR)/%.o: $(UTSBASE)/intel/io/amr/%.c
162 $(COMPILE.c) -o $@ $<
163 $(CTFCONVERT_O)

165 $(OBJSDIR)/%.o: $(UTSBASE)/intel/io/drm/%.c
166 $(COMPILE.c) -o $@ $<
167 $(CTFCONVERT_O)

169 $(OBJSDIR)/%.o: $(UTSBASE)/intel/io/hotplug/pcicfg/%.c
170 $(COMPILE.c) -o $@ $<
171 $(CTFCONVERT_O)

173 $(OBJSDIR)/%.o: $(UTSBASE)/intel/io/ipmi/%.c
174 $(COMPILE.c) -o $@ $<
175 $(CTFCONVERT_O)

177 $(OBJSDIR)/%.o: $(UTSBASE)/intel/io/intel_nb5000/%.c
178 $(COMPILE.c) -o $@ $<
179 $(CTFCONVERT_O)

181 $(OBJSDIR)/%.o: $(UTSBASE)/intel/io/intel_nhm/%.c
182 $(COMPILE.c) -o $@ $<
183 $(CTFCONVERT_O)

185 $(OBJSDIR)/%.o: $(SRC)/common/mc/mc-amd/%.c
186 $(COMPILE.c) -o $@ $<
187 $(CTFCONVERT_O)

189 $(OBJSDIR)/%.o: $(UTSBASE)/intel/io/mc-amd/%.c
190 $(COMPILE.c) -o $@ $<

```

new/usr/src/uts/intel/Makefile.rules

4

```

191 $(CTFCONVERT_O)

193 $(OBJSDIR)/%.o: $(UTSBASE)/intel/io/pci/%.c
194 $(COMPILE.c) -o $@ $<
195 $(CTFCONVERT_O)

197 $(OBJSDIR)/%.o: $(UTSBASE)/intel/io/pciex/%.c
198 $(COMPILE.c) -o $@ $<
199 $(CTFCONVERT_O)

201 $(OBJSDIR)/%.o: $(UTSBASE)/intel/io/dktp/controller/ata/%.c
202 $(COMPILE.c) -o $@ $<
203 $(CTFCONVERT_O)

205 $(OBJSDIR)/%.o: $(UTSBASE)/intel/io/dktp/dcdev/%.c
206 $(COMPILE.c) -o $@ $<
207 $(CTFCONVERT_O)

209 $(OBJSDIR)/%.o: $(UTSBASE)/intel/io/dktp/disk/%.c
210 $(COMPILE.c) -o $@ $<
211 $(CTFCONVERT_O)

213 $(OBJSDIR)/%.o: $(UTSBASE)/intel/io/dktp/drvobj/%.c
214 $(COMPILE.c) -o $@ $<
215 $(CTFCONVERT_O)

217 $(OBJSDIR)/%.o: $(UTSBASE)/intel/io/dktp/hba/ghd/%.c
218 $(COMPILE.c) -o $@ $<
219 $(CTFCONVERT_O)

221 $(OBJSDIR)/%.o: $(UTSBASE)/intel/io/dnet/%.c
222 $(COMPILE.c) -o $@ $<
223 $(CTFCONVERT_O)

225 $(OBJSDIR)/%.o: $(UTSBASE)/intel/io/scsi/adapters/arcmsr/%.c
226 $(COMPILE.c) -o $@ $<
227 $(CTFCONVERT_O)

229 $(OBJSDIR)/%.o: $(UTSBASE)/intel/io/scsi/targets/%.c
230 $(COMPILE.c) -o $@ $<
231 $(CTFCONVERT_O)

233 $(OBJSDIR)/%.o: $(UTSBASE)/intel/io/vgatext/%.c
234 $(COMPILE.c) -o $@ $<
235 $(CTFCONVERT_O)

237 $(OBJSDIR)/%.o: $(UTSBASE)/intel/io/vmxnet3s/%.c
238 $(COMPILE.c) -o $@ $<
239 $(CTFCONVERT_O)

241 $(OBJSDIR)/%.o: $(UTSBASE)/intel/io/scsi/adapters/pvscsi/%.c
242 $(COMPILE.c) -o $@ $<
243 $(CTFCONVERT_O)

245 $(OBJSDIR)/%.o: $(UTSBASE)/intel/nskern/%.s
246 $(COMPILE.s) -o $@ $<

248 $(OBJSDIR)/%.o: $(UTSBASE)/intel/os/%.c
249 $(COMPILE.c) -o $@ $<
250 $(CTFCONVERT_O)

252 $(OBJSDIR)/%.o: $(UTSBASE)/intel/pcbe/%.c
253 $(COMPILE.c) -o $@ $<
254 $(CTFCONVERT_O)

256 $(OBJSDIR)/%.o: $(UTSBASE)/intel/promif/%.c

```

```

257 $(COMPILE.c) -o $@ $<
258 $(CTFCONVERT_O)

260 $(OBJDIR)/%.o: $(UTSBASE)/intel/syscall/%.c
261 $(COMPILE.c) -o $@ $<
262 $(CTFCONVERT_O)

264 $(OBJDIR)/%.o: $(UTSBASE)/common/os/%.c
265 $(COMPILE.c) -o $@ $<
266 $(CTFCONVERT_O)

268 $(OBJDIR)/%.o: $(UTSBASE)/intel/kdi/%.c
269 $(COMPILE.c) -o $@ $<
270 $(CTFCONVERT_O)

272 $(OBJDIR)/%.o: $(UTSBASE)/intel/kdi/%.s
273 $(COMPILE.s) -o $@ $<

275 $(OBJDIR)/%.o: $(UTSBASE)/intel/kdi/$(SUBARCH_DIR)/%.s
276 $(COMPILE.s) -o $@ $<

275 $(OBJDIR)/%.o: $(UTSBASE)/intel/zfs/%.c
276 $(COMPILE.c) -o $@ $<
277 $(CTFCONVERT_O)

279 #
280 # krtld compiled into unix
281 #

283 KRTLD_INC_PATH = -I$(UTSBASE)/common/krtld -I$(UTSBASE)/intel/sys
284 KRTLD_INC_PATH += -I$(UTSBASE)/intel/$(SUBARCH_DIR)/krtld

286 KRTLD_CPPFLAGS_32 = -DELFTARGET_386
287 KRTLD_CPPFLAGS_64 = -DELFTARGET_AMD64 -DMODDIR_SUFFIX=\"amd64\"
288 KRTLD_CPPFLAGS = $(KRTLD_CPPFLAGS_$(CLASS)) -DKRTLD

290 $(OBJDIR)/%.o: $(UTSBASE)/common/krtld/%.c
291 $(COMPILE.c) $(KRTLD_INC_PATH) $(KRTLD_CPPFLAGS) -o $@ $<
292 $(CTFCONVERT_O)

294 $(OBJDIR)/%.o: $(UTSBASE)/intel/$(SUBARCH_DIR)/krtld/%.c
295 $(COMPILE.c) $(KRTLD_INC_PATH) $(KRTLD_CPPFLAGS) -o $@ $<
296 $(CTFCONVERT_O)

298 #
299 # _DBOOT indicates that krtld is called from a dboot ELF section
300 #
301 $(OBJDIR)/kobj.o := CPPFLAGS += -D_DBOOT

303 $(OBJDIR)/%.o: $(UTSBASE)/intel/$(SUBARCH_DIR)/krtld/%.s
304 $(COMPILE.s) $(KRTLD_INC_PATH) $(KRTLD_CPPFLAGS) -o $@ $<
305 $(CTFCONVERT_O)

307 $(OBJDIR)/%.o: $(SRC)/common/util/$(SUBARCH_DIR)/%.c
308 $(COMPILE.c) $(KRTLD_INC_PATH) $(KRTLD_CPPFLAGS) -o $@ $<
309 $(CTFCONVERT_O)

312 #
313 # Section lb: Lint 'object' build rules.
314 #
315 $(LINTSDIR)/%.ln: $(SRC)/common/fs/%.c
316 @$(LHEAD) $(LINT.c) $< $(LTAIL))

318 $(LINTSDIR)/%.ln: $(SRC)/common/util/i386/%.s
319 @$(LHEAD) $(LINT.s) $< $(LTAIL))

```

```

321 $(LINTSDIR)/%.ln: $(UTSBASE)/intel/brand/sn1/%.s
322 @$(LHEAD) $(LINT.s) $< $(LTAIL))

324 $(LINTSDIR)/%.ln: $(UTSBASE)/intel/brand/solaris10/%.s
325 @$(LHEAD) $(LINT.s) $< $(LTAIL))

327 $(LINTSDIR)/%.ln: $(UTSBASE)/intel/dtrace/%.c
328 @$(LHEAD) $(LINT.c) $< $(LTAIL))

330 $(LINTSDIR)/%.ln: $(UTSBASE)/intel/dtrace/%.s
331 @$(LHEAD) $(LINT.s) $< $(LTAIL))

333 $(LINTSDIR)/%.ln: $(UTSBASE)/intel/zfs/%.c
334 @$(LHEAD) $(LINT.c) $< $(LTAIL))

336 $(LINTSDIR)/%.ln: $(UTSBASE)/intel/fs/proc/%.c
337 @$(LHEAD) $(LINT.c) $< $(LTAIL))

339 $(LINTSDIR)/%.ln: $(UTSBASE)/intel/ia32/ml/%.s
340 @$(LHEAD) $(LINT.s) $< $(LTAIL))

342 $(LINTSDIR)/%.ln: $(UTSBASE)/intel/ia32/os/%.c
343 @$(LHEAD) $(LINT.c) $< $(LTAIL))

345 $(LINTSDIR)/%.ln: $(UTSBASE)/intel/ia32/promif/%.c
346 @$(LHEAD) $(LINT.c) $< $(LTAIL))

348 $(LINTSDIR)/%.ln: $(UTSBASE)/intel/ia32/syscall/%.c
349 @$(LHEAD) $(LINT.c) $< $(LTAIL))

351 $(LINTSDIR)/%.ln: $(UTSBASE)/intel/io/%.c
352 @$(LHEAD) $(LINT.c) $< $(LTAIL))

354 $(LINTSDIR)/%.ln: $(UTSBASE)/intel/io/acpica/%.s
355 @$(LHEAD) $(LINT.s) $< $(LTAIL))

357 $(LINTSDIR)/%.ln: $(UTSBASE)/intel/io/acpica/%.c
358 @$(LHEAD) $(LINT.c) $< $(LTAIL))

360 $(LINTSDIR)/%.ln: $(UTSBASE)/intel/io/acpica/debugger/%.c
361 @$(LHEAD) $(LINT.c) $< $(LTAIL))

363 $(LINTSDIR)/%.ln: $(UTSBASE)/intel/io/acpica/events/%.c
364 @$(LHEAD) $(LINT.c) $< $(LTAIL))

366 $(LINTSDIR)/%.ln: $(UTSBASE)/intel/io/acpica/hardware/%.c
367 @$(LHEAD) $(LINT.c) $< $(LTAIL))

369 $(LINTSDIR)/%.ln: $(UTSBASE)/intel/io/acpica/dispatcher/%.c
370 @$(LHEAD) $(LINT.c) $< $(LTAIL))

372 $(LINTSDIR)/%.ln: $(UTSBASE)/intel/io/acpica/executer/%.c
373 @$(LHEAD) $(LINT.c) $< $(LTAIL))

375 $(LINTSDIR)/%.ln: $(UTSBASE)/intel/io/acpica/parser/%.c
376 @$(LHEAD) $(LINT.c) $< $(LTAIL))

378 $(LINTSDIR)/%.ln: $(UTSBASE)/intel/io/acpica/namespace/%.c
379 @$(LHEAD) $(LINT.c) $< $(LTAIL))

381 $(LINTSDIR)/%.ln: $(UTSBASE)/intel/io/acpica/resources/%.c
382 @$(LHEAD) $(LINT.c) $< $(LTAIL))

384 $(LINTSDIR)/%.ln: $(UTSBASE)/intel/io/acpica/tables/%.c
385 @$(LHEAD) $(LINT.c) $< $(LTAIL))

```

```

387 $(LINTS_DIR)/%.ln: $(UTSBASE)/intel/io/acpica/utilities/%.c
388     @$(LHEAD) $(LINT.c) $< $(LTAIL))

390 $(LINTS_DIR)/%.ln: $(UTSBASE)/intel/io/acpica/disassembler/%.c
391     @$(LHEAD) $(LINT.c) $< $(LTAIL))

393 $(LINTS_DIR)/%.ln: $(UTSBASE)/intel/io/agpgart/%.c
394     @$(LHEAD) $(LINT.c) $< $(LTAIL))

396 $(LINTS_DIR)/%.ln: $(UTSBASE)/intel/io/agpmaster/%.c
397     @$(LHEAD) $(LINT.c) $< $(LTAIL))

399 $(LINTS_DIR)/%.ln: $(UTSBASE)/intel/io/amd8111s/%.c
400     @$(LHEAD) $(LINT.c) $< $(LTAIL))

402 $(LINTS_DIR)/%.ln: $(UTSBASE)/intel/io/amr/%.c
403     @$(LHEAD) $(LINT.c) $< $(LTAIL))

405 $(LINTS_DIR)/%.ln: $(UTSBASE)/intel/io/scsi/adapters/arcmsr/%.c
406     @$(LHEAD) $(LINT.c) $< $(LTAIL))

408 $(LINTS_DIR)/%.ln: $(UTSBASE)/intel/io/drm/%.c
409     @$(LHEAD) $(LINT.c) $< $(LTAIL))

411 $(LINTS_DIR)/%.ln: $(UTSBASE)/intel/io/hotplug/pcicfg/%.c
412     @$(LHEAD) $(LINT.c) $< $(LTAIL))

414 $(LINTS_DIR)/%.ln: $(UTSBASE)/intel/io/hotplug/pci/%.c
415     @$(LHEAD) $(LINT.c) $< $(LTAIL))

417 $(LINTS_DIR)/%.ln: $(UTSBASE)/intel/io/intel_nb5000/%.c
418     @$(LHEAD) $(LINT.c) $< $(LTAIL))

420 $(LINTS_DIR)/%.ln: $(UTSBASE)/intel/io/intel_nhm/%.c
421     @$(LHEAD) $(LINT.c) $< $(LTAIL))

423 $(LINTS_DIR)/%.ln: $(UTSBASE)/intel/io/ipmi/%.c
424     @$(LHEAD) $(LINT.c) $< $(LTAIL))

426 $(LINTS_DIR)/%.ln: $(SRC)/common/mc/mc-amd/%.c
427     @$(LHEAD) $(LINT.c) $< $(LTAIL))

429 $(LINTS_DIR)/%.ln: $(UTSBASE)/intel/io/mc-amd/%.c
430     @$(LHEAD) $(LINT.c) $< $(LTAIL))

432 $(LINTS_DIR)/%.ln: $(UTSBASE)/intel/io/pci/%.c
433     @$(LHEAD) $(LINT.c) $< $(LTAIL))

435 $(LINTS_DIR)/%.ln: $(UTSBASE)/intel/io/pciex/%.c
436     @$(LHEAD) $(LINT.c) $< $(LTAIL))

438 $(LINTS_DIR)/%.ln: $(UTSBASE)/intel/io/dktp/controller/ata/%.c
439     @$(LHEAD) $(LINT.c) $< $(LTAIL))

441 $(LINTS_DIR)/%.ln: $(UTSBASE)/intel/io/dktp/dcdev/%.c
442     @$(LHEAD) $(LINT.c) $< $(LTAIL))

444 $(LINTS_DIR)/%.ln: $(UTSBASE)/intel/io/dktp/disk/%.c
445     @$(LHEAD) $(LINT.c) $< $(LTAIL))

447 $(LINTS_DIR)/%.ln: $(UTSBASE)/intel/io/dktp/drvoobj/%.c
448     @$(LHEAD) $(LINT.c) $< $(LTAIL))

450 $(LINTS_DIR)/%.ln: $(UTSBASE)/intel/io/dktp/hba/ghd/%.c
451     @$(LHEAD) $(LINT.c) $< $(LTAIL))

```

```

453 $(LINTS_DIR)/%.ln: $(UTSBASE)/intel/io/dnet/%.c
454     @$(LHEAD) $(LINT.c) $< $(LTAIL))

456 $(LINTS_DIR)/%.ln: $(UTSBASE)/intel/io/scsi/targets/%.c
457     @$(LHEAD) $(LINT.c) $< $(LTAIL))

459 $(LINTS_DIR)/%.ln: $(UTSBASE)/intel/io/vgatext/%.c
460     @$(LHEAD) $(LINT.c) $< $(LTAIL))

462 $(LINTS_DIR)/%.ln: $(UTSBASE)/intel/io/vmxnet3s/%.c
463     @$(LHEAD) $(LINT.c) $< $(LTAIL))

465 $(LINTS_DIR)/%.ln: $(UTSBASE)/intel/io/scsi/adapters/pvscsi/%.c
466     @$(LHEAD) $(LINT.c) $< $(LTAIL))

468 $(LINTS_DIR)/%.ln: $(UTSBASE)/intel/os/%.c
469     @$(LHEAD) $(LINT.c) $< $(LTAIL))

471 $(LINTS_DIR)/%.ln: $(UTSBASE)/intel/pcbe/%.c
472     @$(LHEAD) $(LINT.c) $< $(LTAIL))

474 $(LINTS_DIR)/%.ln: $(UTSBASE)/intel/promif/%.c
475     @$(LHEAD) $(LINT.c) $< $(LTAIL))

477 $(LINTS_DIR)/%.ln: $(UTSBASE)/intel/syscall/%.c
478     @$(LHEAD) $(LINT.c) $< $(LTAIL))

480 $(LINTS_DIR)/%.ln: $(UTSBASE)/common/os/%.c
481     @$(LHEAD) $(LINT.c) $< $(LTAIL))

483 $(LINTS_DIR)/%.ln: $(UTSBASE)/intel/kdi/%.c
484     @$(LHEAD) $(LINT.c) $< $(LTAIL))

486 $(LINTS_DIR)/%.ln: $(UTSBASE)/intel/kdi/%.s
487     @$(LHEAD) $(LINT.s) $< $(LTAIL))

492 $(LINTS_DIR)/%.ln: $(UTSBASE)/intel/kdi/$(SUBARCH_DIR)/%.s
493     @$(LHEAD) $(LINT.s) $< $(LTAIL))

489 $(LINTS_DIR)/%.ln: $(UTSBASE)/intel/nskern/%.s
490     @$(LHEAD) $(LINT.s) $< $(LTAIL))

492 #
493 # krtld lints
494 #
495 $(LINTS_DIR)/%.ln: $(UTSBASE)/common/krtld/%.c
496     @$(LHEAD) $(LINT.c) $(KRTLDC_PATH) $(KRTLDCPPFLAGS) $< $(LTAIL))

498 $(LINTS_DIR)/%.ln: $(UTSBASE)/intel/$(SUBARCH_DIR)/krtld/%.c
499     @$(LHEAD) $(LINT.c) $(KRTLDC_PATH) $(KRTLDCPPFLAGS) $< $(LTAIL))

501 $(LINTS_DIR)/%.ln: $(UTSBASE)/intel/$(SUBARCH_DIR)/krtld/%.s
502     @$(LHEAD) $(LINT.s) $(KRTLDC_PATH) $(KRTLDCPPFLAGS) $< $(LTAIL))

504 $(LINTS_DIR)/%.ln: $(SRC)/common/util/$(SUBARCH_DIR)/%.c
505     @$(LHEAD) $(LINT.c) $(KRTLDC_PATH) $(KRTLDCPPFLAGS) $< $(LTAIL))

507 $(OBSJ_DIR)/kobj.ln := CPPFLAGS += -D_DBOOT

```

new/usr/src/uts/intel/amd64/sys/kdi_regs.h

1

2035 Fri Apr 6 17:25:08 2018
new/usr/src/uts/intel/amd64/sys/kdi_regs.h
9441 kmdb should stash %cr3 in kdiregs
Reviewed by: John Levon <john.levon@joyent.com>
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>
9210 remove KMDB branch debugging support
9211 ::crregs could do with cr2/cr3 support
9209 ::ttrace should be able to filter by thread
Reviewed by: Patrick Mooney <patrick.mooney@joyent.com>
Reviewed by: Yuri Pankov <yuripv@yuripv.net>

```
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21
22 /*
23 * Copyright 2007 Sun Microsystems, Inc. All rights reserved.
24 * Use is subject to license terms.
25 *
26 * Copyright 2018 Joyent, Inc.
27 */
```

```
29 #ifndef _AMD64_SYS_KDI_REGS_H
30 #define _AMD64_SYS_KDI_REGS_H
```

```
30 #pragma ident "%Z%M% %I% %E% SMI"
```

```
32 #ifdef __cplusplus
33 extern "C" {
34 #endif
```

```
36 #define KDIREG_NGREG 31
```

```
36 /*
37  * A modified version of struct regs layout.
38 */
```

```
40 #define KDIREG_SAVFP 0
41 #define KDIREG_SAVPC 1
42 #define KDIREG_RDI 2
43 #define KDIREG_RSI 3
44 #define KDIREG_RDX 4
45 #define KDIREG_RCX 5
46 #define KDIREG_R8 6
47 #define KDIREG_R9 7
48 #define KDIREG_RAX 8
```

new/usr/src/uts/intel/amd64/sys/kdi_regs.h

2

```
49 #define KDIREG_RBX 9
50 #define KDIREG_RBP 10
51 #define KDIREG_R10 11
52 #define KDIREG_R11 12
53 #define KDIREG_R12 13
54 #define KDIREG_R13 14
55 #define KDIREG_R14 15
56 #define KDIREG_R15 16
57 #define KDIREG_FSBASE 17
58 #define KDIREG_GSBASE 18
59 #define KDIREG_KGSBASE 19
60 #define KDIREG_CR2 20
61 #define KDIREG_CR3 21
62 #define KDIREG_DS 22
63 #define KDIREG_ES 23
64 #define KDIREG_FS 24
65 #define KDIREG_GS 25
66 #define KDIREG_TRAPNO 26
67 #define KDIREG_ERR 27
68 #define KDIREG_RIP 28
69 #define KDIREG_CS 29
70 #define KDIREG_RFLAGS 30
71 #define KDIREG_RSP 31
72 #define KDIREG_SS 32
62 #define KDIREG_DS 20
63 #define KDIREG_ES 21
64 #define KDIREG_FS 22
65 #define KDIREG_GS 23
66 #define KDIREG_TRAPNO 24
67 #define KDIREG_ERR 25
68 #define KDIREG_RIP 26
69 #define KDIREG_CS 27
70 #define KDIREG_RFLAGS 28
71 #define KDIREG_RSP 29
72 #define KDIREG_SS 30
```

```
74 #define KDIREG_NGREG (KDIREG_SS + 1)
```

```
76 #define KDIREG_PC KDIREG_RIP
77 #define KDIREG_SP KDIREG_RSP
78 #define KDIREG_FP KDIREG_RBP
```

```
78 #ifndef _ASM
```

```
80 /* Patch point for MSR clearing. */
81 #define KDI_MSR_PATCH \
82     nop; nop; nop; nop; \
83     nop; nop; nop; nop; \
84     nop; nop; nop; nop; \
85     nop; nop; nop; nop; \
86     nop
```

```
88 #endif /* _ASM */
```

```
90 #define KDI_MSR_PATCHOFF 8 /* bytes of code before patch point */
91 #define KDI_MSR_PATCHSZ 17 /* bytes in KDI_MSR_PATCH, above */
```

```
80 #ifdef __cplusplus
81 }
```

```
unchanged_portion_omitted
```


new/usr/src/uts/intel/asm/htable.h

1

```
*****
2038 Fri Apr 6 17:25:09 2018
new/usr/src/uts/intel/asm/htable.h
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License, Version 1.0 only
6  * (the "License"). You may not use this file except in compliance
7  * with the License.
8  *
9  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
10 * or http://www.opensolaris.org/os/licensing.
11 * See the License for the specific language governing permissions
12 * and limitations under the License.
13 *
14 * When distributing Covered Code, include this CDDL HEADER in each
15 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
16 * If applicable, add the following below this CDDL HEADER, with the
17 * fields enclosed by brackets "[]" replaced with your own identifying
18 * information: Portions Copyright [yyyy] [name of copyright owner]
19 *
20 * CDDL HEADER END
21 */
22 /*
23 * Copyright 2004 Sun Microsystems, Inc. All rights reserved.
24 * Use is subject to license terms.
25 *
26 * Copyright 2018 Joyent, Inc.
27 */

29 #ifndef _ASM_HTABLE_H
30 #define _ASM_HTABLE_H

32 #include <sys/ccompile.h>
33 #include <sys/types.h>

35 #ifdef __cplusplus
36 extern "C" {
37 #endif

39 #if !defined(__lint) && defined(__GNUC__)

39 #if defined(__i386) || defined(__amd64)

41 /*
42  * This set of atomic operations are designed primarily
43  * for some ia32 hat layer operations.
44  */

46 extern __GNU_INLINE void
47 atomic_orb(uint8_t *addr, uint8_t value)
48 {
49     __asm__ __volatile__(
50         "lock; orb %%dl,%0"
51         : "=m" (*addr)
52         : "d" (value), "m" (*addr)
53         : "cc");
54 }
55
56 extern __GNU_INLINE void
```

new/usr/src/uts/intel/asm/htable.h

2

```
87 mmu_tlbflush_entry(caddr_t addr)
88 {
89     __asm__ __volatile__(
90         "invlpg %0"
91         : "=m" (*addr)
92         : "m" (*addr));
93 }

95 #endif /* __i386 || __amd64 */

86 #endif /* !__lint && __GNUC__ */

88 #ifdef __cplusplus
89 }
90
91 _____unchanged_portion_omitted_____
```

```

*****
2393 Fri Apr 6 17:25:10 2018
new/usr/src/uts/intel/asm/mmu.h
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22 * Copyright 2007 Sun Microsystems, Inc. All rights reserved.
23 * Use is subject to license terms.
24 *
25 * Copyright 2018 Joyent, Inc.
26 */

28 #ifndef _ASM_MMU_H
29 #define _ASM_MMU_H

31 #include <sys/ccompile.h>
32 #include <sys/types.h>

34 #ifdef __cplusplus
35 extern "C" {
36 #endif

38 #if defined(__GNUC__)
39 #if defined(__GNUC__) && !defined(__xpv)

40 #if !defined(__xpv)
41 #if defined(__amd64)

42 extern __GNU_INLINE ulong_t
43 getcr3(void)
44 {
45     uint64_t value;

47     __asm__ __volatile__(
48         "movq %%cr3, %0"
49         : "=r" (value));
50     return (value);
51 }
52
53 _____unchanged_portion_omitted_____

60 extern __GNU_INLINE void
61 reload_cr3(void)
62 {
63     setcr3(getcr3());

```

```

64 }

66 #elif defined(__i386)

62 extern __GNU_INLINE ulong_t
63 getcr4(void)
69 getcr3(void)
64 {
65     uint64_t value;
66     uint32_t value;

67     __asm__ __volatile__(
68         "movq %%cr4, %0"
69         : "=r" (value));
70     return (value);
71 }

73 extern __GNU_INLINE void
74 setcr4(ulong_t value)
80 setcr3(ulong_t value)
75 {
76     __asm__ __volatile__(
77         "movq %0, %%cr4"
78         : /* no output */
79         : "r" (value));
80 }

81 _____unchanged_portion_omitted_____

88 /*
89  * We clobber memory: we're not writing anything, but we don't want to
90  * potentially get re-ordered beyond the TLB flush.
91  */
92 extern __GNU_INLINE void
93 invpcid_insn(uint64_t type, uint64_t pcid, uintptr_t addr)
94 {
95     uint64_t pcid_desc[2] = { pcid, addr };
96     __asm__ __volatile__(
97         "invpcid %0, %1"
98         : /* no output */
99         : "m" (*pcid_desc), "r" (type)
100         : "memory");
101 }
94 #endif

103 #endif /* !__xpv */
96 #endif /* __GNUC__ && !__xpv */

105 extern __GNU_INLINE void
106 mmu_invlpg(caddr_t addr)
107 {
108     __asm__ __volatile__(
109         "invlpg %0"
110         : "=m" (*addr)
111         : "m" (*addr));
112 }

114 #endif /* __GNUC__ */

116 #ifdef __cplusplus
117 }

118 _____unchanged_portion_omitted_____

```

```

*****
34750 Fri Apr 6 17:25:10 2018
new/usr/src/uts/intel/ia32/ml/exception.s
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>
*****
1 /*
2  * Copyright (c) 2004, 2010, Oracle and/or its affiliates. All rights reserved.
3  * Copyright (c) 2013, 2014 by Delphix. All rights reserved.
4  * Copyright (c) 2018 Joyent, Inc.
4  * Copyright (c) 2017 Joyent, Inc.
5  */

7 /*
8  * Copyright (c) 1989, 1990 William F. Jolitz.
9  * Copyright (c) 1990 The Regents of the University of California.
10 * All rights reserved.
11 *
12 * Redistribution and use in source and binary forms, with or without
13 * modification, are permitted provided that the following conditions
14 * are met:
15 * 1. Redistributions of source code must retain the above copyright
16 * notice, this list of conditions and the following disclaimer.
17 * 2. Redistributions in binary form must reproduce the above copyright
18 * notice, this list of conditions and the following disclaimer in the
19 * documentation and/or other materials provided with the distribution.
20 * 3. All advertising materials mentioning features or use of this software
21 * must display the following acknowledgement:
22 *   This product includes software developed by the University of
23 *   California, Berkeley and its contributors.
24 * 4. Neither the name of the University nor the names of its contributors
25 * may be used to endorse or promote products derived from this software
26 * without specific prior written permission.
27 *
28 * THIS SOFTWARE IS PROVIDED BY THE REGENTS AND CONTRIBUTORS ``AS IS'' AND
29 * ANY EXPRESS OR IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED TO, THE
30 * IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS FOR A PARTICULAR PURPOSE
31 * ARE DISCLAIMED. IN NO EVENT SHALL THE REGENTS OR CONTRIBUTORS BE LIABLE
32 * FOR ANY DIRECT, INDIRECT, INCIDENTAL, SPECIAL, EXEMPLARY, OR CONSEQUENTIAL
33 * DAMAGES (INCLUDING, BUT NOT LIMITED TO, PROCUREMENT OF SUBSTITUTE GOODS
34 * OR SERVICES; LOSS OF USE, DATA, OR PROFITS; OR BUSINESS INTERRUPTION)
35 * HOWEVER CAUSED AND ON ANY THEORY OF LIABILITY, WHETHER IN CONTRACT, STRICT
36 * LIABILITY, OR TORT (INCLUDING NEGLIGENCE OR OTHERWISE) ARISING IN ANY WAY
37 * OUT OF THE USE OF THIS SOFTWARE, EVEN IF ADVISED OF THE POSSIBILITY OF
38 * SUCH DAMAGE.
39 *
40 * $FreeBSD: src/sys/amd64/amd64/exception.S,v 1.113 2003/10/15 02:04:52 peter E
41 */

43 #include <sys/asm_linkage.h>
44 #include <sys/asm_misc.h>
45 #include <sys/trap.h>
46 #include <sys/psw.h>
47 #include <sys/regset.h>
48 #include <sys/privregs.h>
49 #include <sys/dtrace.h>
50 #include <sys/x86_archext.h>
51 #include <sys/traptrace.h>
52 #include <sys/machparam.h>

54 /*
55  * only one routine in this file is interesting to lint
56  */

58 #if defined(__lint)

```

```

60 void
61 ndptrap_frstor(void)
62 {}

64 #else

66 #include "assym.h"

68 /*
69  * push $0 on stack for traps that do not
70  * generate an error code. This is so the rest
71  * of the kernel can expect a consistent stack
72  * from from any exception.
73  *
74  * Note that for all exceptions for amd64
75  * %r11 and %rcx are on the stack. Just pop
76  * them back into their appropriate registers and let
77  * it get saved as is running native.
78  */

80 #if defined(__xpv) && defined(__amd64)

82 #define NPTRAP_NOERR(trapno) \
83     pushq $0; \
84     pushq $trapno

86 #define TRAP_NOERR(trapno) \
87     XPV_TRAP_POP; \
88     NPTRAP_NOERR(trapno)

90 /*
91  * error code already pushed by hw
92  * onto stack.
93  */
94 #define TRAP_ERR(trapno) \
95     XPV_TRAP_POP; \
96     pushq $trapno

98 #else /* __xpv && __amd64 */

100 #define TRAP_NOERR(trapno) \
101     push $0; \
102     push $trapno

104 #define NPTRAP_NOERR(trapno) TRAP_NOERR(trapno)

106 /*
107  * error code already pushed by hw
108  * onto stack.
109  */
110 #define TRAP_ERR(trapno) \
111     push $trapno

113 #endif /* __xpv && __amd64 */

115 /*
116  * These are the stacks used on cpu0 for taking double faults,
117  * NMIs and MCEs (the latter two only on amd64 where we have IST).
118  *
119  * We define them here instead of in a C file so that we can page-align
120  * them (gcc won't do that in a .c file).
121  */
122 .data
123 DGDEF3(dblfault_stack0, DEFAULTSTKSZ, MMU_PAGESIZE)
124 .fill DEFAULTSTKSZ, 1, 0

```

```

125 DGDEF3(nmi_stack0, DEFAULTSTKSZ, MMU_PAGESIZE)
126 .fill  DEFAULTSTKSZ, 1, 0
127 DGDEF3(mce_stack0, DEFAULTSTKSZ, MMU_PAGESIZE)
128 .fill  DEFAULTSTKSZ, 1, 0

130 /*
131  * #DE
132  */
133 ENTRY_NP(div0trap)
134 TRAP_NOERR(T_ZERODIV) /* $0 */
135 jmp  cmntrap
136 SET_SIZE(div0trap)

138 /*
139  * #DB
140  *
141  * Fetch %dr6 and clear it, handing off the value to the
142  * cmntrap code in %r15/%esi
143  */
144 ENTRY_NP(dbgtrap)
145 TRAP_NOERR(T_SGLSTP) /* $1 */

147 #if defined(__amd64)
148 #if !defined(__xpv) /* no sysenter support yet */
149 /*
150  * If we get here as a result of single-stepping a sysenter
151  * instruction, we suddenly find ourselves taking a #db
152  * in kernel mode -before- we've swappgs'ed. So before we can
153  * take the trap, we do the swappgs here, and fix the return
154  * %rip in trap() so that we return immediately after the
155  * swappgs in the sysenter handler to avoid doing the swappgs again.
156  *
157  * Nobody said that the design of sysenter was particularly
158  * elegant, did they?
159  */

161 pushq  %r11

163 /*
164  * At this point the stack looks like this:
165  *
166  * (high address)      r_ss
167  *                      r_rsp
168  *                      r_rfl
169  *                      r_cs
170  *                      r_rip      <-- %rsp + 24
171  *                      r_err      <-- %rsp + 16
172  *                      r_trapno   <-- %rsp + 8
173  * (low address)      %r11      <-- %rsp
174  */
175 leaq  sys_sysenter(%rip), %r11
176 cmpq  %r11, 24(%rsp) /* Compare to saved r_rip on the stack */
177 je    1f
178 leaq  brand_sys_sysenter(%rip), %r11
179 cmpq  %r11, 24(%rsp) /* Compare to saved r_rip on the stack */
180 je    1f
181 leaq  tr_sys_sysenter(%rip), %r11
182 cmpq  %r11, 24(%rsp)
183 je    1f
184 leaq  tr_brand_sys_sysenter(%rip), %r11
185 cmpq  %r11, 24(%rsp)
186 jne  2f
187 1:  SWAPGS
188 2:  popq  %r11
189 #endif /* !__xpv */

```

```

191     INTR_PUSH
192 #if defined(__xpv)
193     movl  $6, %edi
194     call  kdi_dreg_get
195     movq  %rax, %r15 /* %db6 -> %r15 */
196     movl  $6, %edi
197     movl  $0, %esi
198     call  kdi_dreg_set /* 0 -> %db6 */
199 #else
200     movq  %db6, %r15
201     xorl  %eax, %eax
202     movq  %rax, %db6
203 #endif

205 #elif defined(__i386)

207     INTR_PUSH
208 #if defined(__xpv)
209     pushl  $6
210     call  kdi_dreg_get
211     addl  $4, %esp
212     movl  %eax, %esi /* %dr6 -> %esi */
213     pushl  $0
214     pushl  $6
215     call  kdi_dreg_set /* 0 -> %dr6 */
216     addl  $8, %esp
217 #else
218     movl  %db6, %esi
219     xorl  %eax, %eax
220     movl  %eax, %db6
221 #endif
222 #endif /* __i386 */

224     jmp  cmntrap_pushed
225     SET_SIZE(dbgtrap)

227 #if defined(__amd64)
228 #if !defined(__xpv)

230 /*
231  * Macro to set the gsbases to the address of the struct cpu
232  * for this processor. If we came from userland, set kgsbase else
233  * set gsbases. We find the proper cpu struct by looping through
234  * the cpu structs for all processors till we find a match for the gdt
235  * of the trapping processor. The stack is expected to be pointing at
236  * the standard regs pushed by hardware on a trap (plus error code and trapno).
237  *
238  * It's ok for us to clobber gsbases here (and possibly end up with both gsbases
239  * and kgsbases set to the same value) because we're not going back the normal
240  * way out of here (via IRET). Where we're going, we don't need no user %gs.
241  */
242 #define SET_CPU_GSBASE
243     subq  $REGOFF_TRAPNO, %rsp; /* save regs */
244     movq  %rax, REGOFF_RAX(%rsp);
245     movq  %rbx, REGOFF_RBX(%rsp);
246     movq  %rcx, REGOFF_RCX(%rsp);
247     movq  %rdx, REGOFF_RDX(%rsp);
248     movq  %rbp, REGOFF_RBP(%rsp);
249     movq  %rsp, %rbp;
250     subq  $16, %rsp; /* space for gdt */
251     sgdt  6(%rsp);
252     movq  8(%rsp), %rcx; /* %rcx has gdt to match */
253     xorl  %ebx, %ebx; /* loop index */
254     leaq  cpu(%rip), %rdx; /* cpu pointer array */
255 1:
256     movq  (%rdx, %rbx, LONGSIZE), %rax; /* get cpu[i] */

```

```

257     cmpq    $0x0, %rax;          /* cpu[i] == NULL ? */
258     je      2f;                  /* yes, continue */
259     cmpq    %rcx, CPU_GDT(%rax); /* gdt == cpu[i]->cpu_gdt ? */
260     je      3f;                  /* yes, go set gsbases */
261 2:
262     incl    %ebx;                /* i++ */
263     cmpl    $NCPU, %ebx;         /* i < NCPU ? */
264     jb     1b;                  /* yes, loop */
265 /* XXX BIG trouble if we fall thru here. We didn't find a gdt match */
266 3:
267     movl    $MSR_AMD_KGSBASE, %ecx;
268     cmpw    $KCS_SEL, REGOFF_CS(%rbp); /* trap from kernel? */
269     jne    4f;                  /* no, go set KGSBASE */
270     movl    $MSR_AMD_GSBASE, %ecx; /* yes, set GSBASE */
271     mfence;                      /* OPTERON_ERRATUM_88 */
272 4:
273     movq    %rax, %rdx;          /* write base register */
274     shrq    $32, %rdx;
275     wrmsr;
276     movq    REGOFF_RDX(%rbp), %rdx; /* restore regs */
277     movq    REGOFF_RCX(%rbp), %rcx;
278     movq    REGOFF_RBX(%rbp), %rbx;
279     movq    REGOFF_RAX(%rbp), %rax;
280     movq    %rbp, %rsp;
281     movq    REGOFF_RBP(%rsp), %rbp;
282     addq    $REGOFF_TRAPNO, %rsp /* pop stack */

284 #else /* __xpv */

286 #define SET_CPU_GSBASE /* noop on the hypervisor */

288 #endif /* __xpv */
289 #endif /* __amd64 */

292 #if defined(__amd64)

294     /*
295     * #NMI
296     *
297     * XXPV: See 6532669.
298     */
299     ENTRY_NP(nmiint)
300     TRAP_NOERR(T_NMIFLT) /* $2 */

302     SET_CPU_GSBASE

304     /*
305     * Save all registers and setup segment registers
306     * with kernel selectors.
307     */
308     INTR_PUSH
309     INTGATE_INIT_KERNEL_FLAGS

311     TRACE_PTR(%r12, %rax, %eax, %rdx, $TT_TRAP)
312     TRACE_REGS(%r12, %rsp, %rax, %rbx)
313     TRACE_STAMP(%r12)

315     movq    %rsp, %rbp

317     movq    %rbp, %rdi
318     call   av_dispatch_nmivect

320     INTR_POP
321     jmp    tr_iret_auto
297     IRET

```

```

322     /*NOTREACHED*/
323     SET_SIZE(nmiint)
_____ unchanged_portion_omitted

402 #if defined(__amd64)

404     ENTRY_NP(invoptrap)

406     XPV_TRAP_POP

408     cmpw    $KCS_SEL, 8(%rsp)
409     jne     ud_user

411 #if defined(__xpv)
412     movb    $0, 12(%rsp) /* clear saved upcall_mask from %cs */
413 #endif
414     push    $0 /* error code -- zero for #UD */
415 ud_kernel:
416     push    $0xdddd /* a dummy trap number */
417     INTR_PUSH
418     movq    REGOFF_RIP(%rsp), %rdi
419     movq    REGOFF_RSP(%rsp), %rsi
420     movq    REGOFF_RAX(%rsp), %rdx
421     pushq   (%rsi)
422     movq    %rsp, %rsi
423     subq    $8, %rsp
424     call   dtrace_invop
425     ALTERNATIVE(dtrace_invop_callsite)
426     addq    $16, %rsp
427     cmpl    $DTRACE_INVOP_PUSHL_EBP, %eax
428     je     ud_push
429     cmpl    $DTRACE_INVOP_LEAVE, %eax
430     je     ud_leave
431     cmpl    $DTRACE_INVOP_NOP, %eax
432     je     ud_nop
433     cmpl    $DTRACE_INVOP_RET, %eax
434     je     ud_ret
435     jmp    ud_trap

437 ud_push:
438     /*
439     * We must emulate a "pushq %rbp". To do this, we pull the stack
440     * down 8 bytes, and then store the base pointer.
441     */
442     INTR_POP
443     subq    $16, %rsp /* make room for %rbp */
444     pushq   %rax /* push temp */
445     movq    24(%rsp), %rax /* load calling RIP */
446     addq    $1, %rax /* increment over trapping instr */
447     movq    %rax, 8(%rsp) /* store calling RIP */
448     movq    32(%rsp), %rax /* load calling CS */
449     movq    %rax, 16(%rsp) /* store calling CS */
450     movq    40(%rsp), %rax /* load calling RFLAGS */
451     movq    %rax, 24(%rsp) /* store calling RFLAGS */
452     movq    48(%rsp), %rax /* load calling RSP */
453     subq    $8, %rax /* make room for %rbp */
454     movq    %rax, 32(%rsp) /* store calling RSP */
455     movq    56(%rsp), %rax /* load calling SS */
456     movq    %rax, 40(%rsp) /* store calling SS */
457     movq    32(%rsp), %rax /* reload calling RSP */
458     movq    %rbp, (%rax) /* store %rbp there */
459     popq    %rax /* pop off temp */
460     jmp    tr_iret_kernel /* return from interrupt */
436     IRET /* return from interrupt */
461     /*NOTREACHED*/

```

```

463 ud_leave:
464 /*
465  * We must emulate a "leave", which is the same as a "movq %rbp, %rsp"
466  * followed by a "popq %rbp". This is quite a bit simpler on amd64
467  * than it is on i386 -- we can exploit the fact that the %rsp is
468  * explicitly saved to effect the pop without having to reshuffle
469  * the other data pushed for the trap.
470  */
471 INTR_POP
472 pushq  %rax                /* push temp */
473 movq   8(%rsp), %rax       /* load calling RIP */
474 addq   $1, %rax           /* increment over trapping instr */
475 movq   %rax, 8(%rsp)       /* store calling RIP */
476 movq   (%rbp), %rax        /* get new %rbp */
477 addq   $8, %rbp           /* adjust new %rsp */
478 movq   %rbp, 32(%rsp)     /* store new %rsp */
479 movq   %rax, %rbp         /* set new %rbp */
480 popq   %rax               /* pop off temp */
481 jmp   tr_iret_kernel    /* return from interrupt */
457 IRET                          /* return from interrupt */
482 /*NOTREACHED*/

484 ud_nop:
485 /*
486  * We must emulate a "nop". This is obviously not hard: we need only
487  * advance the %rip by one.
488  */
489 INTR_POP
490 incq   (%rsp)
491 jmp   tr_iret_kernel
467 IRET
492 /*NOTREACHED*/

494 ud_ret:
495 INTR_POP
496 pushq  %rax                /* push temp */
497 movq   32(%rsp), %rax       /* load %rsp */
498 movq   (%rax), %rax         /* load calling RIP */
499 movq   %rax, 8(%rsp)       /* store calling RIP */
500 addq   $8, 32(%rsp)        /* adjust new %rsp */
501 popq   %rax               /* pop off temp */
502 jmp   tr_iret_kernel    /* return from interrupt */
478 IRET                          /* return from interrupt */
503 /*NOTREACHED*/

505 ud_trap:
506 /*
507  * We're going to let the kernel handle this as a normal #UD. If,
508  * however, we came through #BP and are spoofing #UD (in this case,
509  * the stored error value will be non-zero), we need to de-spoof
510  * the trap by incrementing %rip and pushing T_BPTFLT.
511  */
512 cmpq   $0, REGOFF_ERR(%rsp)
513 je     ud_ud
514 incq   REGOFF_RIP(%rsp)
515 addq   $REGOFF_RIP, %rsp
516 NPTRAP_NOERR(T_BPTFLT) /* $3 */
517 jmp   cmntrap

519 ud_ud:
520 addq   $REGOFF_RIP, %rsp
521 ud_user:
522 NPTRAP_NOERR(T_ILLINST)
523 jmp   cmntrap
524 SET_SIZE(invoptrap)
unchanged portion omitted

```

```

718 #else /* __xpv */

720 ENTRY_NP(ndptrap)
721 /*
722  * We want to do this quickly as every lwp using fp will take this
723  * after a context switch -- we do the frequent path in ndptrap_frstor
724  * below; for all other cases, we let the trap code handle it
725  */
726 pushq  %rax
727 pushq  %rbx
728 cmpw   $KCS_SEL, 24(%rsp) /* did we come from kernel mode? */
729 jne    lf
730 LOADCPU(%rax) /* if yes, don't swapgs */
731 jmp   2f
732 1:
733 SWAPGS /* if from user, need swapgs */
734 LOADCPU(%rax)
735 SWAPGS
736 2:
737 /*
738  * Xrstor needs to use edx as part of its flag.
739  * NOTE: have to push rdx after "cmpw ...24(%rsp)", otherwise rsp+$24
740  * will not point to CS.
741  */
742 pushq  %rdx
743 cmpl   $0, fpu_exists(%rip)
744 je     .handle_in_trap /* let trap handle no fp case */
745 movq   CPU_THREAD(%rax), %rbx /* %rbx = curthread */
746 movl   $FPU_EN, %eax
747 movq   T_LWP(%rbx), %rbx /* %rbx = lwp */
748 testq  %rbx, %rbx
749 jz     .handle_in_trap /* should not happen? */
750 #if LWP_PCB_FPU != 0
751 addq   $LWP_PCB_FPU, %rbx /* &lwp->lwp_pcb.pcb_fpu */
752 #endif
753 testl  %eax, PCB_FPU_FLAGS(%rbx)
754 jz     .handle_in_trap /* must be the first fault */
755 clts
756 andl   $_BITNOT(FPU_VALID), PCB_FPU_FLAGS(%rbx)
757 #if FPU_CTX_FPU_REGS != 0
758 addq   $FPU_CTX_FPU_REGS, %rbx
759 #endif

761 movl   FPU_CTX_FPU_XSAVE_MASK(%rbx), %eax /* for xrstor */
762 movl   FPU_CTX_FPU_XSAVE_MASK+4(%rbx), %edx /* for xrstor */

764 /*
765  * the label below is used in trap.c to detect FP faults in
766  * kernel due to user fault.
767  */
768 ALTENTRY(ndptrap_frstor)
769 movq   (%rbx), %rbx /* fpu_regs.kfpu_u.kfpu_XX pointer */
770 .globl _patch_xrstorq_rbx
771 _patch_xrstorq_rbx:
772 fxrstorq (%rbx)
773 popq   %rdx
774 popq   %rbx
775 popq   %rax
776 jmp   tr_iret_auto
752 IRET
777 /*NOTREACHED*/

779 .handle_in_trap:
780 popq   %rdx
781 popq   %rbx

```

```

782     popq     %rax
783     TRAP_NOERR(T_NOEXTFLT) /* $7 */
784     jmp     cmnintrap
785     SET_SIZE(ndptrap_frstor)
_____
unchanged_portion_omitted_

1151 #endif /* !__amd64 */

1153     ENTRY_NP(resvtrap)
1154     TRAP_NOERR(T_RESVTRAP) /* (reserved) */
1155     TRAP_NOERR(15) /* (reserved) */
1156     jmp     cmntrap
1157     SET_SIZE(resvtrap)
_____
unchanged_portion_omitted_

1233     ENTRY_NP(ivaltrap)
1234     TRAP_NOERR(T_INVALIDTRAP) /* very invalid */
1235     TRAP_NOERR(30) /* very invalid */
1236     jmp     cmntrap
1237     SET_SIZE(ivaltrap)

1214     ENTRY_NP(ivalint)
1215     TRAP_NOERR(31) /* even more so */
1216     jmp     cmnint
1217     SET_SIZE(ivalint)

1238     .globl fasttable

1240 #if defined(__amd64)

1242     ENTRY_NP(fasttrap)
1243     cmpl     $T_LASTFAST, %eax
1244     ja      lf
1245     orl     %eax, %eax /* (zero extend top 32-bits) */
1246     leaq   fasttable(%rip), %r11
1247     leaq   (%r11, %rax, CLONGSIZE), %r11
1248     jmp     *(%r11)
1249 1:
1250     /*
1251     * Fast syscall number was illegal. Make it look
1252     * as if the INT failed. Modify %rip to point before the
1253     * INT, push the expected error code and fake a GP fault.
1254     *
1255     * XXX Why make the error code be offset into idt + 1?
1256     * Instead we should push a real (soft?) error code
1257     * on the stack and #gp handler could know about fasttraps?
1258     */
1259     XPV_TRAP_POP

1261     subq    $2, (%rsp) /* XXX int insn 2-bytes */
1262     pushq   $_CONST(_MUL(T_FASTTRAP, GATE_DESC_SIZE) + 2)

1264 #if defined(__xpv)
1265     pushq   %r11
1266     pushq   %rcx
1267 #endif
1268     jmp     gptrap
1269     SET_SIZE(fasttrap)
_____
unchanged_portion_omitted_

1299 #if defined(__amd64)

1301     /*
1302     * RFLAGS 24 bytes up the stack from %rsp.
1303     * XXX a constant would be nicer.
1304     */

```

```

1305     ENTRY_NP(fast_null)
1306     XPV_TRAP_POP
1307     orq     $SPS_C, 24(%rsp) /* set carry bit in user flags */
1308     jmp     tr_iret_auto
1289     IRET
1309     /*NOTREACHED*/
1310     SET_SIZE(fast_null)
_____
unchanged_portion_omitted_

```

```

*****
81879 Fri Apr 6 17:25:10 2018
new/usr/src/uts/intel/ia32/ml/i86_subr.s
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */

22 /*
23 * Copyright (c) 1992, 2010, Oracle and/or its affiliates. All rights reserved.
24 * Copyright 2014 Nexenta Systems, Inc. All rights reserved.
25 * Copyright (c) 2014 by Delphix. All rights reserved.
26 * Copyright 2018 Joyent, Inc.
26 * Copyright 2016 Joyent, Inc.
27 */

29 /*
30 * Copyright (c) 1990, 1991 UNIX System Laboratories, Inc.
31 * Copyright (c) 1984, 1986, 1987, 1988, 1989, 1990 AT&T
32 * All Rights Reserved
33 */

35 /*
36 * Copyright (c) 2009, Intel Corporation.
37 * All rights reserved.
38 */

40 /*
41 * General assembly language routines.
42 * It is the intent of this file to contain routines that are
43 * independent of the specific kernel architecture, and those that are
44 * common across kernel architectures.
45 * As architectures diverge, and implementations of specific
46 * architecture-dependent routines change, the routines should be moved
47 * from this file into the respective ../arch -k/subr.s file.
48 */

50 #include <sys/asm_linkage.h>
51 #include <sys/asm_misc.h>
52 #include <sys/panic.h>
53 #include <sys/ontrap.h>
54 #include <sys/regset.h>
55 #include <sys/privregs.h>
56 #include <sys/reboot.h>
57 #include <sys/psw.h>
58 #include <sys/x86_archext.h>

```

```

60 #if defined(__lint)
61 #include <sys/types.h>
62 #include <sys/system.h>
63 #include <sys/thread.h>
64 #include <sys/archsystem.h>
65 #include <sys/byteorder.h>
66 #include <sys/dtrace.h>
67 #include <sys/ftrace.h>
68 #else /* __lint */
69 #include "assym.h"
70 #endif /* __lint */
71 #include <sys/dditypes.h>

73 /*
74 * on_fault()
75 *
76 * Catch lofault faults. Like setjmp except it returns one
77 * if code following causes uncorrectable fault. Turned off
78 * by calling no_fault(). Note that while under on_fault(),
79 * SMAP is disabled. For more information see
80 * uts/intel/ia32/ml/copy.s.
81 */

83 #if defined(__lint)

85 /* ARGSUSED */
86 int
87 on_fault(label_t *ljb)
88 { return (0); }

90 void
91 no_fault(void)
92 {}

94 #else /* __lint */

96 #if defined(__amd64)

98     ENTRY(on_fault)
99     movq    %gs:CPU_THREAD, %rsi
100     leaq   catch_fault(%rip), %rdx
101     movq   %rdi, T_ONFAULT(%rsi)          /* jumpbuf in t_onfault */
102     movq   %rdx, T_LOFAULT(%rsi)        /* catch_fault in t_lofault */
103     call   smap_disable                 /* allow user accesses */
104     jmp    setjmp                       /* let setjmp do the rest */

106 catch_fault:
107     movq   %gs:CPU_THREAD, %rsi
108     movq   T_ONFAULT(%rsi), %rdi        /* address of save area */
109     xorl   %eax, %eax
110     movq   %rax, T_ONFAULT(%rsi)       /* turn off onfault */
111     movq   %rax, T_LOFAULT(%rsi)      /* turn off lofault */
112     call   smap_enable                 /* disallow user accesses */
113     jmp    longjmp                     /* let longjmp do the rest */
114     SET_SIZE(on_fault)

_____unchanged_portion_omitted_____

428 #endif /* __i386 */
429 #endif /* __lint */

431 /*
432 * Invalidate a single page table entry in the TLB
433 */

435 #if defined(__lint)

```



```

437 /* ARGSUSED */
438 void
439 mmu_invlpg(caddr_t m)
439 mmu_tlbflush_entry(caddr_t m)
440 {}

442 #else /* __lint */

444     ENTRY(mmu_invlpg)
444 #if defined(__amd64)

446     ENTRY(mmu_tlbflush_entry)
445     invlpg (%rdi)
446     ret
447     SET_SIZE(mmu_invlpg)
449     SET_SIZE(mmu_tlbflush_entry)

451 #elif defined(__i386)

453     ENTRY(mmu_tlbflush_entry)
454     movl 4(%esp), %eax
455     invlpg (%eax)
456     ret
457     SET_SIZE(mmu_tlbflush_entry)

459 #endif /* __i386 */
449 #endif /* __lint */

452 /*
453  * Get/Set the value of various control registers
454  */

456 #if defined(__lint)

458 ulong_t
459 getcr0(void)
460 { return (0); }

462 /* ARGSUSED */
463 void
464 setcr0(ulong_t value)
465 {}

467 ulong_t
468 getcr2(void)
469 { return (0); }

471 ulong_t
472 getcr3(void)
473 { return (0); }

475 #if !defined(__xpv)
476 /* ARGSUSED */
477 void
478 setcr3(ulong_t val)
479 {}

481 void
482 reload_cr3(void)
483 {}
484 #endif

486 ulong_t
487 getcr4(void)

```

```

488 { return (0); }

490 /* ARGSUSED */
491 void
492 setcr4(ulong_t val)
493 {}

495 #if defined(__amd64)

497 ulong_t
498 getcr8(void)
499 { return (0); }

501 /* ARGSUSED */
502 void
503 setcr8(ulong_t val)
504 {}

506 #endif /* __amd64 */

508 #else /* __lint */

510 #if defined(__amd64)

512     ENTRY(getcr0)
513     movq %cr0, %rax
514     ret
515     SET_SIZE(getcr0)
515     unchanged_portion_omitted

```

```

*****
25007 Fri Apr 6 17:25:11 2018
new/usr/src/uts/intel/ia32/ml/swtch.s
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22 * Copyright 2007 Sun Microsystems, Inc. All rights reserved.
23 * Use is subject to license terms.
24 */

26 /*
27 * Copyright (c) 2018 Joyent, Inc.
28 * Copyright (c) 2013, Joyent, Inc. All rights reserved.
29 */

30 /*
31 * Process switching routines.
32 */

34 #if defined(__lint)
35 #include <sys/thread.h>
36 #include <sys/system.h>
37 #include <sys/time.h>
38 #else /* __lint */
39 #include "assym.h"
40 #endif /* __lint */

42 #include <sys/asm_linkage.h>
43 #include <sys/asm_misc.h>
44 #include <sys/regset.h>
45 #include <sys/privregs.h>
46 #include <sys/stack.h>
47 #include <sys/segments.h>
48 #include <sys/psw.h>

50 /*
51 * resume(thread_id_t t);
52 *
53 * a thread can only run on one processor at a time. there
54 * exists a window on MPs where the current thread on one
55 * processor is capable of being dispatched by another processor.
56 * some overlap between outgoing and incoming threads can happen
57 * when they are the same thread. in this case where the threads
58 * are the same, resume() on one processor will spin on the incoming

```

```

59 * thread until resume() on the other processor has finished with
60 * the outgoing thread.
61 *
62 * The MMU context changes when the resuming thread resides in a different
63 * process. Kernel threads are known by resume to reside in process 0.
64 * The MMU context, therefore, only changes when resuming a thread in
65 * a process different from curproc.
66 *
67 * resume_from_intr() is called when the thread being resumed was not
68 * passivated by resume (e.g. was interrupted). This means that the
69 * resume lock is already held and that a restore context is not needed.
70 * Also, the MMU context is not changed on the resume in this case.
71 *
72 * resume_from_zombie() is the same as resume except the calling thread
73 * is a zombie and must be put on the deathrow list after the CPU is
74 * off the stack.
75 */

77 #if !defined(__lint)

79 #if LWP_PCB_FPU != 0
80 #error LWP_PCB_FPU MUST be defined as 0 for code in swtch.s to work
81 #endif /* LWP_PCB_FPU != 0 */

83 #endif /* !__lint */

85 #if defined(__amd64)

87 /*
88 * Save non-volatile regs other than %rsp (%rbx, %rbp, and %r12 - %r15)
89 *
90 * The stack frame must be created before the save of %rsp so that tracebacks
91 * of swtch()ed-out processes show the process as having last called swtch().
92 */
93 #define SAVE_REGS(thread_t, retaddr) \
94     movq    %rbp, T_RBP(thread_t); \
95     movq    %rbx, T_RBX(thread_t); \
96     movq    %r12, T_R12(thread_t); \
97     movq    %r13, T_R13(thread_t); \
98     movq    %r14, T_R14(thread_t); \
99     movq    %r15, T_R15(thread_t); \
100    pushq   %rbp; \
101    movq    %rsp, %rbp; \
102    movq    %rsp, T_SP(thread_t); \
103    movq    retaddr, T_PC(thread_t); \
104    movq    %rdi, %r12; \
105    call    __dtrace_probe__sched_off__cpu

107 /*
108 * Restore non-volatile regs other than %rsp (%rbx, %rbp, and %r12 - %r15)
109 *
110 * We load up %rsp from the label_t as part of the context switch, so
111 * we don't repeat that here.
112 *
113 * We don't do a 'leave,' because reloading %rsp/%rbp from the label_t
114 * already has the effect of putting the stack back the way it was when
115 * we came in.
116 */
117 #define RESTORE_REGS(scratch_reg) \
118     movq    %gs:CPU_THREAD, scratch_reg; \
119     movq    T_RBP(scratch_reg), %rbp; \
120     movq    T_RBX(scratch_reg), %rbx; \
121     movq    T_R12(scratch_reg), %r12; \
122     movq    T_R13(scratch_reg), %r13; \
123     movq    T_R14(scratch_reg), %r14; \
124     movq    T_R15(scratch_reg), %r15

```

```

126 /*
127 * Get pointer to a thread's hat structure
128 */
129 #define GET_THREAD_HATP(hatp, thread_t, scratch_reg) \
130     movq    T_PROCP(thread_t), hatp; \
131     movq    P_AS(hatp), scratch_reg; \
132     movq    A_HAT(scratch_reg), hatp

134 #define TSC_READ() \
135     call    tsc_read; \
136     movq    %rax, %r14;

138 /*
139 * If we are resuming an interrupt thread, store a timestamp in the thread
140 * structure. If an interrupt occurs between tsc_read() and its subsequent
141 * store, the timestamp will be stale by the time it is stored. We can detect
142 * this by doing a compare-and-swap on the thread's timestamp, since any
143 * interrupt occurring in this window will put a new timestamp in the thread's
144 * t_intr_start field.
145 */
146 #define STORE_INTR_START(thread_t) \
147     testw   $T_INTR_THREAD, T_FLAGS(thread_t); \
148     jz     1f; \
149 0: \
150     TSC_READ(); \
151     movq    T_INTR_START(thread_t), %rax; \
152     cmpxchgq %r14, T_INTR_START(thread_t); \
153     jnz    0b; \
154 1:

156 #elif defined (__i386)

158 /*
159 * Save non-volatile registers (%ebp, %esi, %edi and %ebx)
160 *
161 * The stack frame must be created before the save of %esp so that tracebacks
162 * of swtch()ed-out processes show the process as having last called swtch().
163 */
164 #define SAVE_REGS(thread_t, retaddr) \
165     movl   %ebp, T_EBP(thread_t); \
166     movl   %ebx, T_EBX(thread_t); \
167     movl   %esi, T_ESI(thread_t); \
168     movl   %edi, T_EDI(thread_t); \
169     pushl  %ebp; \
170     movl   %esp, %ebp; \
171     movl   %esp, T_SP(thread_t); \
172     movl   retaddr, T_PC(thread_t); \
173     movl   8(%ebp), %edi; \
174     pushl  %edi; \
175     call   __dtrace_probe__sched_off__cpu; \
176     addl   $CLONGSIZE, %esp

178 /*
179 * Restore non-volatile registers (%ebp, %esi, %edi and %ebx)
180 *
181 * We don't do a 'leave,' because reloading %rsp/%rbp from the label_t
182 * already has the effect of putting the stack back the way it was when
183 * we came in.
184 */
185 #define RESTORE_REGS(scratch_reg) \
186     movl   %gs:CPU_THREAD, scratch_reg; \
187     movl   T_EBP(scratch_reg), %ebp; \
188     movl   T_EBX(scratch_reg), %ebx; \
189     movl   T_ESI(scratch_reg), %esi; \
190     movl   T_EDI(scratch_reg), %edi

```

```

192 /*
193 * Get pointer to a thread's hat structure
194 */
195 #define GET_THREAD_HATP(hatp, thread_t, scratch_reg) \
196     movl   T_PROCP(thread_t), hatp; \
197     movl   P_AS(hatp), scratch_reg; \
198     movl   A_HAT(scratch_reg), hatp

200 /*
201 * If we are resuming an interrupt thread, store a timestamp in the thread
202 * structure. If an interrupt occurs between tsc_read() and its subsequent
203 * store, the timestamp will be stale by the time it is stored. We can detect
204 * this by doing a compare-and-swap on the thread's timestamp, since any
205 * interrupt occurring in this window will put a new timestamp in the thread's
206 * t_intr_start field.
207 */
208 #define STORE_INTR_START(thread_t) \
209     testw   $T_INTR_THREAD, T_FLAGS(thread_t); \
210     jz     1f; \
211     pushl  %ecx; \
212 0: \
213     pushl  T_INTR_START(thread_t); \
214     pushl  T_INTR_START+4(thread_t); \
215     call   tsc_read; \
216     movl   %eax, %ebx; \
217     movl   %edx, %ecx; \
218     popl   %edx; \
219     popl   %eax; \
220     cmpxchg8b T_INTR_START(thread_t); \
221     jnz    0b; \
222     popl   %ecx; \
223 1:

225 #endif /* __amd64 */

227 #if defined(__lint)

229 /* ARGSUSED */
230 void
231 resume(kthread_t *t)
232 {}

234 #else /* __lint */

236 #if defined(__amd64)

238     .global kpti_enable

240     ENTRY(resume)
241     movq   %gs:CPU_THREAD, %rax
242     leaq   resume_return(%rip), %r11

244     /*
245     * Deal with SMAP here. A thread may be switched out at any point while
246     * it is executing. The thread could be under on_fault() or it could be
247     * pre-empted while performing a copy interruption. If this happens and
248     * we're not in the context of an interrupt which happens to handle
249     * saving and restoring rflags correctly, we may lose our SMAP related
250     * state.
251     *
252     * To handle this, as part of being switched out, we first save whether
253     * or not userland access is allowed ($PS_ACHK in rflags) and store that
254     * in t_useracc on the kthread_t and unconditionally enable SMAP to
255     * protect the system.
256     */

```

```

257     * Later, when the thread finishes resuming, we potentially disable smap
258     * if PS_ACHK was present in rflags. See uts/intel/ia32/ml/copy.s for
259     * more information on rflags and SMAP.
260     */
261     pushfq
262     popq   %rsi
263     andq   $PS_ACHK, %rsi
264     movq   %rsi, T_USERACC(%rax)
265     call   smap_enable

267     /*
268     * Save non-volatile registers, and set return address for current
269     * thread to resume_return.
270     *
271     * %r12 = t (new thread) when done
272     */
273     SAVE_REGS(%rax, %r11)

276     LOADCPU(%r15)
277     movq   CPU_THREAD(%r15), %r13          /* %r15 = CPU */
                                           /* %r13 = curthread */

279     /*
280     * Call savectx if thread has installed context ops.
281     *
282     * Note that if we have floating point context, the save op
283     * (either fpsave_begin or fpxsave_begin) will issue the
284     * async save instruction (fnsave or fxsave respectively)
285     * that we fwait for below.
286     */
287     cmpq   $0, T_CTX(%r13)                /* should current thread savectx? */
288     je     .nosavectx                    /* skip call when zero */

290     movq   %r13, %rdi                    /* arg = thread pointer */
291     call   savectx                      /* call ctx ops */
292 .nosavectx:

294     /*
295     * Call saveptcx if process has installed context ops.
296     */
297     movq   T_PROCP(%r13), %r14           /* %r14 = proc */
298     cmpq   $0, P_PCTX(%r14)             /* should current thread savectx? */
299     je     .nosaveptcx                  /* skip call when zero */

301     movq   %r14, %rdi                    /* arg = proc pointer */
302     call   saveptcx                      /* call ctx ops */
303 .nosaveptcx:

305     /*
306     * Temporarily switch to the idle thread's stack
307     */
308     movq   CPU_IDLE_THREAD(%r15), %rax   /* idle thread pointer */

310     /*
311     * Set the idle thread as the current thread
312     */
313     movq   T_SP(%rax), %rsp              /* It is safe to set rsp */
314     movq   %rax, CPU_THREAD(%r15)

316     /*
317     * Switch in the hat context for the new thread
318     *
319     */
320     GET_THREAD_HATP(%rdi, %r12, %r11)
321     call   hat_switch

```

```

323     /*
324     * Clear and unlock previous thread's t_lock
325     * to allow it to be dispatched by another processor.
326     */
327     movb   $0, T_LOCK(%r13)

329     /*
330     * IMPORTANT: Registers at this point must be:
331     *           %r12 = new thread
332     *
333     * Here we are in the idle thread, have dropped the old thread.
334     */
335     ALTENTRY(_resume_from_idle)
336     /*
337     * spin until dispatched thread's mutex has
338     * been unlocked. this mutex is unlocked when
339     * it becomes safe for the thread to run.
340     */
341 .lock_thread_mutex:
342     lock
343     btsl   $0, T_LOCK(%r12)             /* attempt to lock new thread's mutex */
344     jnc    .thread_mutex_locked        /* got it */

346 .spin_thread_mutex:
347     pause
348     cmpb   $0, T_LOCK(%r12)             /* check mutex status */
349     jz     .lock_thread_mutex          /* clear, retry lock */
350     jmp    .spin_thread_mutex          /* still locked, spin... */

352 .thread_mutex_locked:
353     /*
354     * Fix CPU structure to indicate new running thread.
355     * Set pointer in new thread to the CPU structure.
356     */
357     LOADCPU(%r13)                      /* load current CPU pointer */
358     cmpq   %r13, T_CPU(%r12)
359     je     .setup_cpu

361     /* cp->cpu_stats.sys.cpumigrate++ */
362     incq   CPU_STATS_SYS_CPUMIGRATE(%r13)
363     movq   %r13, T_CPU(%r12)           /* set new thread's CPU pointer */

365 .setup_cpu:
366     /*
367     * Setup rsp0 (kernel stack) in TSS to curthread's saved regs
368     * structure. If this thread doesn't have a regs structure above
369     * the stack -- that is, if lwp_stk_init() was never called for the
370     * thread -- this will set rsp0 to the wrong value, but it's harmless
371     * as it's a kernel thread, and it won't actually attempt to implicitly
372     * use the rsp0 via a privilege change.
373     * Setup rsp0 (kernel stack) in TSS to curthread's stack.
374     * (Note: Since we don't have saved 'regs' structure for all
375     * the threads we can't easily determine if we need to
376     * change rsp0. So, we simply change the rsp0 to bottom
377     * of the thread stack and it will work for all cases.)
378     */
379     * Note that when we have KPTI enabled on amd64, we never use this
380     * value at all (since all the interrupts have an IST set).
381     * XX64 - Is this correct?
382     */
383     movq   CPU_TSS(%r13), %r14
384     #if !defined(__xpv)
385     cmpq   $1, kpti_enable
386     jne    1f
387     leaq   CPU_KPTI_TR_RSP(%r13), %rax
388     jmp    2f

```

```

383 1:
384     movq    T_STACK(%r12), %rax
385     addq    $REGSIZE+MINFRAME, %rax /* to the bottom of thread stack */
386 2:
376 #if !defined(__xpv)
387     movq    %rax, TSS_RSP0(%r14)
388 #else
389     movq    T_STACK(%r12), %rax
390     addq    $REGSIZE+MINFRAME, %rax /* to the bottom of thread stack */
391     movl    $KDS_SEL, %edi
392     movq    %rax, %rsi
393     call   HYPERVISOR_stack_switch
394 #endif /* __xpv */

396     movq    %r12, CPU_THREAD(%r13) /* set CPU's thread pointer */
397     mfence /* synchronize with mutex_exit() */
398     xorl    %ebp, %ebp /* make $<threadlist behave better */
399     movq    T_LWP(%r12), %rax /* set associated lwp to */
400     movq    %rax, CPU_LWP(%r13) /* CPU's lwp ptr */

402     movq    T_SP(%r12), %rsp /* switch to outgoing thread's stack */
403     movq    T_PC(%r12), %r13 /* saved return addr */

405     /*
406     * Call restorectx if context ops have been installed.
407     */
408     cmpq    $0, T_CTX(%r12) /* should resumed thread restorectx? */
409     jz     .norestorectx /* skip call when zero */
410     movq    %r12, %rdi /* arg = thread pointer */
411     call   restorectx /* call ctx ops */
412 .norestorectx:

414     /*
415     * Call restorepctx if context ops have been installed for the proc.
416     */
417     movq    T_PROCP(%r12), %rcx
418     cmpq    $0, P_PCTX(%rcx)
419     jz     .norestorepctx
420     movq    %rcx, %rdi
421     call   restorepctx
422 .norestorepctx:

424     STORE_INTR_START(%r12)

426     /*
427     * If we came into swtch with the ability to access userland pages, go
428     * ahead and restore that fact by disabling SMAP. Clear the indicator
429     * flag out of paranoia.
430     */
431     movq    T_USERACC(%r12), %rax /* should we disable smap? */
432     cmpq    $0, %rax /* skip call when zero */
433     jz     .nosmap
434     xorq    %rax, %rax
435     movq    %rax, T_USERACC(%r12)
436     call   smap_disable
437 .nosmap:

439     /*
440     * Restore non-volatile registers, then have spl0 return to the
441     * resuming thread's PC after first setting the priority as low as
442     * possible and blocking all interrupt threads that may be active.
443     */
444     movq    %r13, %rax /* save return address */
445     RESTORE_REGS(%r11)
446     pushq   %rax /* push return address for spl0() */
447     call   __dtrace_probe__sched_on__cpu

```

```

448     jmp     spl0

450 resume_return:
451     /*
452     * Remove stack frame created in SAVE_REGS()
453     */
454     addq    $CLONGSIZE, %rsp
455     ret
456     SET_SIZE(_resume_from_idle)
_____ unchanged_portion_omitted_

```

```

*****
38018 Fri Apr 6 17:25:11 2018
new/usr/src/uts/intel/ia32/os/desctbls.c
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */

22 /*
23  * Copyright (c) 2004, 2010, Oracle and/or its affiliates. All rights reserved.
24 */

26 /*
27  * Copyright 2018 Joyent, Inc. All rights reserved.
28  * Copyright 2011 Joyent, Inc. All rights reserved.
29 */

30 /*
31  * Copyright (c) 1992 Terrence R. Lambert.
32  * Copyright (c) 1990 The Regents of the University of California.
33  * All rights reserved.
34  *
35  * This code is derived from software contributed to Berkeley by
36  * William Jolitz.
37  *
38  * Redistribution and use in source and binary forms, with or without
39  * modification, are permitted provided that the following conditions
40  * are met:
41  * 1. Redistributions of source code must retain the above copyright
42  * notice, this list of conditions and the following disclaimer.
43  * 2. Redistributions in binary form must reproduce the above copyright
44  * notice, this list of conditions and the following disclaimer in the
45  * documentation and/or other materials provided with the distribution.
46  * 3. All advertising materials mentioning features or use of this software
47  * must display the following acknowledgement:
48  * This product includes software developed by the University of
49  * California, Berkeley and its contributors.
50  * 4. Neither the name of the University nor the names of its contributors
51  * may be used to endorse or promote products derived from this software
52  * without specific prior written permission.
53  *
54  * THIS SOFTWARE IS PROVIDED BY THE REGENTS AND CONTRIBUTORS ``AS IS'' AND
55  * ANY EXPRESS OR IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED TO, THE
56  * IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS FOR A PARTICULAR PURPOSE
57  * ARE DISCLAIMED. IN NO EVENT SHALL THE REGENTS OR CONTRIBUTORS BE LIABLE
58  * FOR ANY DIRECT, INDIRECT, INCIDENTAL, SPECIAL, EXEMPLARY, OR CONSEQUENTIAL

```

```

59  * DAMAGES (INCLUDING, BUT NOT LIMITED TO, PROCUREMENT OF SUBSTITUTE GOODS
60  * OR SERVICES; LOSS OF USE, DATA, OR PROFITS; OR BUSINESS INTERRUPTION)
61  * HOWEVER CAUSED AND ON ANY THEORY OF LIABILITY, WHETHER IN CONTRACT, STRICT
62  * LIABILITY, OR TORT (INCLUDING NEGLIGENCE OR OTHERWISE) ARISING IN ANY WAY
63  * OUT OF THE USE OF THIS SOFTWARE, EVEN IF ADVISED OF THE POSSIBILITY OF
64  * SUCH DAMAGE.
65  *
66  * from: @(#)machdep.c 7.4 (Berkeley) 6/3/91
67 */

69 #include <sys/types.h>
70 #include <sys/sysmacros.h>
71 #include <sys/tss.h>
72 #include <sys/segments.h>
73 #include <sys/trap.h>
74 #include <sys/cpuvar.h>
75 #include <sys/bootconf.h>
76 #include <sys/x86_archext.h>
77 #include <sys/controlregs.h>
78 #include <sys/archsystem.h>
79 #include <sys/machsystem.h>
80 #include <sys/kobj.h>
81 #include <sys/cmn_err.h>
82 #include <sys/reboot.h>
83 #include <sys/kdi.h>
84 #include <sys/mach_mmu.h>
85 #include <sys/system.h>
86 #include <sys/note.h>

88 #ifdef __xpv
89 #include <sys/hypervisor.h>
90 #include <vm/as.h>
91 #endif

93 #include <sys/promif.h>
94 #include <sys/bootinfo.h>
95 #include <vm/kboot_mmu.h>
96 #include <vm/hat_pte.h>

98 /*
99  * cpu0 and default tables and structures.
100 */
101 user_desc_t *gdt0;
102 #if !defined(__xpv)
103 desctbr_t gdt0_default_r;
104 #endif

106 gate_desc_t *idt0; /* interrupt descriptor table */
107 #if defined(__i386)
108 desctbr_t idt0_default_r; /* describes idt0 in IDTR format */
109 #endif

111 tss_t *ktss0; /* kernel task state structure */

113 #if defined(__i386)
114 tss_t *dftss0; /* #DF double-fault exception */
115 #endif /* __i386 */

117 user_desc_t zero_udesc; /* base zero user desc native procs */
118 user_desc_t null_udesc; /* null user descriptor */
119 system_desc_t null_sdesc; /* null system descriptor */

121 #if defined(__amd64)
122 user_desc_t zero_u32desc; /* 32-bit compatibility procs */
123 #endif /* __amd64 */

```

```

125 #if defined(__amd64)
126 user_desc_t    ucs_on;
127 user_desc_t    ucs_off;
128 user_desc_t    ucs32_on;
129 user_desc_t    ucs32_off;
130 #endif /* __amd64 */

132 /*
133  * If the size of this is changed, you must update hat_pcp_setup() and the
134  * definitions in exception.s
135  */
136 extern char dblfault_stack0[DEFAULTSTKSZ];
137 extern char nmi_stack0[DEFAULTSTKSZ];
138 extern char mce_stack0[DEFAULTSTKSZ];
131 #pragma align 16(dblfault_stack0)
132 char          dblfault_stack0[DEFAULTSTKSZ];

140 extern void    fast_null(void);
141 extern hrtime_t get_hrtime(void);
142 extern hrtime_t gethrvtime(void);
143 extern hrtime_t get_hrestime(void);
144 extern uint64_t getlgrp(void);

146 void (*(fasttable[]))(void) = {
147     fast_null,                /* T_FNULL routine */
148     fast_null,                /* T_FGETFP routine (initially null) */
149     fast_null,                /* T_FSETFP routine (initially null) */
150     (void (*)())get_hrtime,   /* T_GETHRTIME */
151     (void (*)())gethrvtime,   /* T_GETHRVTIME */
152     (void (*)())get_hrestime, /* T_GETHRESTIME */
153     (void (*)())getlgrp,      /* T_GETLGRP */
154 };
155
156 _____unchanged_portion_omitted_____

315 #endif /* __i386 */

317 /*
318  * Install gate segment descriptor for interrupt, trap, call and task gates.
319  *
320  * For 64 bit native if we have KPTI enabled, we use the IST stack mechanism on
321  * all interrupts. We have different ISTs for each class of exceptions that are
322  * most likely to occur while handling an existing exception; while many of
323  * these are just going to panic, it's nice not to trample on the existing
324  * exception state for debugging purposes.
325  *
326  * Normal interrupts are all redirected unconditionally to the KPTI trampoline
327  * stack space. This unifies the trampoline handling between user and kernel
328  * space (and avoids the need to touch %gs).
329  *
330  * The KDI IDT *all* uses the DBG IST: consider single stepping tr_pfttrap, when
331  * we do a read from KMDB that cause another #PF. Without its own IST, this
332  * would stomp on the kernel's mcpu_kptiflt frame.
333  */
334 uint_t
335 idt_vector_to_ist(uint_t vector)

315 #if defined(__amd64)

317 /*ARGUSED*/
318 void
319 set_gatesegd(gate_desc_t *dp, void (*func)(void), selector_t sel,
320             uint_t type, uint_t dpl, uint_t vector)
336 {
337 #if defined(__xpv)
338     _NOTE(ARGUNUSED(vector));
339     return (IST_NONE);

```

```

340 #else
341     switch (vector) {
342         /* These should always use IST even without KPTI enabled. */
343         case T_DBLFLT:
344             return (IST_DF);
345         case T_NMIFLT:
346             return (IST_NMI);
347         case T_MCE:
348             return (IST_MCE);
322         dp->sgd_loffset = (uintptr_t)func;
323         dp->sgd_hioffset = (uintptr_t)func >> 16;
324         dp->sgd_hi64offset = (uintptr_t)func >> (16 + 16);

350         case T_BPTFLT:
351         case T_SGLSTP:
352             if (kpti_enable == 1) {
353                 return (IST_DBG);
354             }
355             return (IST_NONE);
356         case T_STKFLT:
357         case T_GPFLT:
358         case T_PGFLT:
359             if (kpti_enable == 1) {
360                 return (IST_NESTABLE);
361             }
362             return (IST_NONE);
363         default:
364             if (kpti_enable == 1) {
365                 return (IST_DEFAULT);
366             }
367             return (IST_NONE);
368     }
326     dp->sgd_selector = (uint16_t)sel;

328     /*
329     * For 64 bit native we use the IST stack mechanism
330     * for double faults. All other traps use the CPL = 0
331     * (tss_rsp0) stack.
332     */
333 #if !defined(__xpv)
334     if (vector == T_DBLFLT)
335         dp->sgd_ist = 1;
336     else
337         dp->sgd_ist = 0;
338 #endif

340     dp->sgd_type = type;
341     dp->sgd_dpl = dpl;
342     dp->sgd_p = 1;
370 }

345 #elif defined(__i386)

347 /*ARGUSED*/
372 void
373 set_gatesegd(gate_desc_t *dp, void (*func)(void), selector_t sel,
374             uint_t type, uint_t dpl, uint_t ist)
350     uint_t type, uint_t dpl, uint_t unused)
375 {
376     dp->sgd_loffset = (uintptr_t)func;
377     dp->sgd_hioffset = (uintptr_t)func >> 16;
378     dp->sgd_hi64offset = (uintptr_t)func >> (16 + 16);

379     dp->sgd_selector = (uint16_t)sel;
380     dp->sgd_ist = ist;
356     dp->sgd_stkcpy = 0; /* always zero bytes */

```

```

381     dp->sgd_type = type;
382     dp->sgd_dpl = dpl;
383     dp->sgd_p = 1;
384 }

362 #endif /* __i386 */

386 /*
387  * Updates a single user descriptor in the the GDT of the current cpu.
388  * Caller is responsible for preventing cpu migration.
389  */

391 void
392 gdt_update_usegd(uint_t sidx, user_desc_t *udp)
393 {
394 #if defined(__xpv)
395     uint64_t dpa = CPU->cpu_m.mcpu_gdtpa + sizeof (*udp) * sidx;
396
397     if (HYPERVISOR_update_descriptor(pa_to_ma(dpa), *(uint64_t *)udp))
398         panic("gdt_update_usegd: HYPERVISOR_update_descriptor");
399
401 #else /* __xpv */
402     CPU->cpu_gdt[sidx] = *udp;
403
404 #endif /* __xpv */
405 #endif /* __xpv */
406 }
407
408 unchanged_portion_omitted
409
410
422 #endif /* __xpv */
423 #endif /* __i386 */

425 /*
426  * Build kernel IDT.
427  *
428  * Note that for amd64 we pretty much require every gate to be an interrupt
429  * gate which blocks interrupts atomically on entry; that's because of our
430  * dependency on using 'swaps' every time we come into the kernel to find
431  * the cpu structure. If we get interrupted just before doing that, %cs could
432  * be in kernel mode (so that the trap prolog doesn't do a swaps), but
433  * %gsbase is really still pointing at something in userland. Bad things will
434  * ensue. We also use interrupt gates for i386 as well even though this is not
435  * required for some traps.
436  *
437  * Perhaps they should have invented a trap gate that does an atomic swaps?
438  */
439 static void
440 init_idt_common(gate_desc_t *idt)
441 {
442     set_gatesegd(&idt[T_ZERODIV],
443                 (kpti_enable == 1) ? &tr_div0trap : &div0trap,
444                 KCS_SEL, SDT_SYSIGT, TRP_KPL, idt_vector_to_ist(T_ZERODIV));
445     set_gatesegd(&idt[T_SGLSTP],
446                 (kpti_enable == 1) ? &tr_dbgtrap : &dbgtrap,
447                 KCS_SEL, SDT_SYSIGT, TRP_KPL, idt_vector_to_ist(T_SGLSTP));
448     set_gatesegd(&idt[T_NMIFLT],
449                 (kpti_enable == 1) ? &tr_nmiint : &nmiint,
450                 KCS_SEL, SDT_SYSIGT, TRP_KPL, idt_vector_to_ist(T_NMIFLT));
451     set_gatesegd(&idt[T_BPTFLT],
452                 (kpti_enable == 1) ? &tr_brktrap : &brktrap,
453                 KCS_SEL, SDT_SYSIGT, TRP_UPL, idt_vector_to_ist(T_BPTFLT));
454     set_gatesegd(&idt[T_OVFLW],
455                 (kpti_enable == 1) ? &tr_ovftrap : &ovftrap,
456                 KCS_SEL, SDT_SYSIGT, TRP_UPL, idt_vector_to_ist(T_OVFLW));
457     set_gatesegd(&idt[T_BOUNDFLT],

```

```

958     (kpti_enable == 1) ? &tr_boundstrap : &boundstrap,
959     KCS_SEL, SDT_SYSIGT, TRP_KPL, idt_vector_to_ist(T_BOUNDFLT));
960     set_gatesegd(&idt[T_ILLINST],
961                 (kpti_enable == 1) ? &tr_invoptrap : &invoptrap,
962                 KCS_SEL, SDT_SYSIGT, TRP_KPL, idt_vector_to_ist(T_ILLINST));
963     set_gatesegd(&idt[T_NOEXTFLT],
964                 (kpti_enable == 1) ? &tr_ndptrap : &ndptrap,
965                 KCS_SEL, SDT_SYSIGT, TRP_KPL, idt_vector_to_ist(T_NOEXTFLT));
920     set_gatesegd(&idt[T_ZERODIV], &div0trap, KCS_SEL, SDT_SYSIGT, TRP_KPL,
921                 0);
922     set_gatesegd(&idt[T_SGLSTP], &dbgtrap, KCS_SEL, SDT_SYSIGT, TRP_KPL,
923                 0);
924     set_gatesegd(&idt[T_NMIFLT], &nmiint, KCS_SEL, SDT_SYSIGT, TRP_KPL,
925                 0);
926     set_gatesegd(&idt[T_BPTFLT], &brktrap, KCS_SEL, SDT_SYSIGT, TRP_UPL,
927                 0);
928     set_gatesegd(&idt[T_OVFLW], &ovftrap, KCS_SEL, SDT_SYSIGT, TRP_UPL,
929                 0);
930     set_gatesegd(&idt[T_BOUNDFLT], &boundstrap, KCS_SEL, SDT_SYSIGT,
931                 TRP_KPL, 0);
932     set_gatesegd(&idt[T_ILLINST], &invoptrap, KCS_SEL, SDT_SYSIGT, TRP_KPL,
933                 0);
934     set_gatesegd(&idt[T_NOEXTFLT], &ndptrap, KCS_SEL, SDT_SYSIGT, TRP_KPL,
935                 0);

967     /*
968     * double fault handler.
969     *
970     * Note that on the hypervisor a guest does not receive #df faults.
971     * Instead a failsafe event is injected into the guest if its selectors
972     * and/or stack is in a broken state. See xen_failsafe_callback.
973     */
974 #if !defined(__xpv)
975     set_gatesegd(&idt[T_DBLFLT],
976                 (kpti_enable == 1) ? &tr_syserrtrap : &syserrtrap,
977                 KCS_SEL, SDT_SYSIGT, TRP_KPL, idt_vector_to_ist(T_DBLFLT));
945 #if defined(__amd64)
947     set_gatesegd(&idt[T_DBLFLT], &syserrtrap, KCS_SEL, SDT_SYSIGT, TRP_KPL,
948                 T_DBLFLT);

950 #elif defined(__i386)

952     /*
953     * task gate required.
954     */
955     set_gatesegd(&idt[T_DBLFLT], NULL, DFTSS_SEL, SDT_SYSTASKGT, TRP_KPL,
956                 0);

958 #endif /* __i386 */
978 #endif /* !__xpv */

980     /*
981     * T_EXTOVRFILT coprocessor-segment-overflow not supported.
982     */
983     set_gatesegd(&idt[T_TSSFLT],
984                 (kpti_enable == 1) ? &tr_invtstrap : &invtstrap,
985                 KCS_SEL, SDT_SYSIGT, TRP_KPL, idt_vector_to_ist(T_TSSFLT));
986     set_gatesegd(&idt[T_SEGFLT],
987                 (kpti_enable == 1) ? &tr_segnttrap : &segnttrap,
988                 KCS_SEL, SDT_SYSIGT, TRP_KPL, idt_vector_to_ist(T_SEGFLT));
989     set_gatesegd(&idt[T_STKFLT],
990                 (kpti_enable == 1) ? &tr_stktrap : &stktrap,
991                 KCS_SEL, SDT_SYSIGT, TRP_KPL, idt_vector_to_ist(T_STKFLT));
992     set_gatesegd(&idt[T_GPFLT],
993                 (kpti_enable == 1) ? &tr_gptrap : &gptrap,

```



```

994     KCS_SEL, SDT_SYSGT, TRP_KPL, idt_vector_to_ist(T_GPFLT));
995     set_gatesegd(&idt[T_PGFLT],
996     (kpti_enable == 1) ? &tr_pftrap : &pftrap,
997     KCS_SEL, SDT_SYSGT, TRP_KPL, idt_vector_to_ist(T_PGFLT));
998     set_gatesegd(&idt[T_EXTERRFLT],
999     (kpti_enable == 1) ? &tr_ndperr : &ndperr,
1000     KCS_SEL, SDT_SYSGT, TRP_KPL, idt_vector_to_ist(T_EXTERRFLT));
1001     set_gatesegd(&idt[T_ALIGNMENT],
1002     (kpti_enable == 1) ? &tr_achktrap : &achktrap,
1003     KCS_SEL, SDT_SYSGT, TRP_KPL, idt_vector_to_ist(T_ALIGNMENT));
1004     set_gatesegd(&idt[T_MCE],
1005     (kpti_enable == 1) ? &tr_mcetrap : &mcetrap,
1006     KCS_SEL, SDT_SYSGT, TRP_KPL, idt_vector_to_ist(T_MCE));
1007     set_gatesegd(&idt[T_SIMDFPE],
1008     (kpti_enable == 1) ? &tr_xmtrap : &xmtrap,
1009     KCS_SEL, SDT_SYSGT, TRP_KPL, idt_vector_to_ist(T_SIMDFPE));

965     set_gatesegd(&idt[T_TSSFLT], &invtstrap, KCS_SEL, SDT_SYSGT, TRP_KPL,
966     0);
967     set_gatesegd(&idt[T_SEGFLT], &segnptrap, KCS_SEL, SDT_SYSGT, TRP_KPL,
968     0);
969     set_gatesegd(&idt[T_STKFLT], &stktrap, KCS_SEL, SDT_SYSGT, TRP_KPL, 0);
970     set_gatesegd(&idt[T_GPFLT], &gptrap, KCS_SEL, SDT_SYSGT, TRP_KPL, 0);
971     set_gatesegd(&idt[T_PGFLT], &pftrap, KCS_SEL, SDT_SYSGT, TRP_KPL, 0);
972     set_gatesegd(&idt[T_EXTERRFLT], &ndperr, KCS_SEL, SDT_SYSGT, TRP_KPL,
973     0);
974     set_gatesegd(&idt[T_ALIGNMENT], &achktrap, KCS_SEL, SDT_SYSGT,
975     TRP_KPL, 0);
976     set_gatesegd(&idt[T_MCE], &mcetrap, KCS_SEL, SDT_SYSGT, TRP_KPL, 0);
977     set_gatesegd(&idt[T_SIMDFPE], &xmtrap, KCS_SEL, SDT_SYSGT, TRP_KPL, 0);

1011     /*
1012     * install fast trap handler at 210.
1013     */
1014     set_gatesegd(&idt[T_FASTTRAP],
1015     (kpti_enable == 1) ? &tr_fasttrap : &fasttrap,
1016     KCS_SEL, SDT_SYSGT, TRP_UPL, idt_vector_to_ist(T_FASTTRAP));
982     set_gatesegd(&idt[T_FASTTRAP], &fasttrap, KCS_SEL, SDT_SYSGT, TRP_UPL,
983     0);

1018     /*
1019     * System call handler.
1020     */
1021     set_gatesegd(&idt[T_SYSCALLINT],
1022     (kpti_enable == 1) ? &tr_sys_syscall_int : &sys_syscall_int,
1023     KCS_SEL, SDT_SYSGT, TRP_UPL, idt_vector_to_ist(T_SYSCALLINT));
988 #if defined(__amd64)
989     set_gatesegd(&idt[T_SYSCALLINT], &sys_syscall_int, KCS_SEL, SDT_SYSGT,
990     TRP_UPL, 0);

992 #elif defined(__i386)
993     set_gatesegd(&idt[T_SYSCALLINT], &sys_call, KCS_SEL, SDT_SYSGT,
994     TRP_UPL, 0);
995 #endif /* __i386 */

1025     /*
1026     * Install the DTrace interrupt handler for the pid provider.
1027     */
1028     set_gatesegd(&idt[T_DTRACE_RET],
1029     (kpti_enable == 1) ? &tr_dtrace_ret : &dtrace_ret,
1030     KCS_SEL, SDT_SYSGT, TRP_UPL, idt_vector_to_ist(T_DTRACE_RET));
1000     set_gatesegd(&idt[T_DTRACE_RET], &dtrace_ret, KCS_SEL,
1001     SDT_SYSGT, TRP_UPL, 0);

1032     /*
1033     * Prepare interposing descriptor for the syscall handler

```

```

1034     * and cache copy of the default descriptor.
1035     */
1036     brand_tbl[0].ih_inum = T_SYSCALLINT;
1037     brand_tbl[0].ih_default_desc = idt0[T_SYSCALLINT];

1039     set_gatesegd(&(brand_tbl[0].ih_interp_desc),
1040     (kpti_enable == 1) ? &tr_brand_sys_syscall_int :
1041     &brand_sys_syscall_int, KCS_SEL, SDT_SYSGT, TRP_UPL,
1042     idt_vector_to_ist(T_SYSCALLINT));
1010 #if defined(__amd64)
1011     set_gatesegd(&(brand_tbl[0].ih_interp_desc), &brand_sys_syscall_int,
1012     KCS_SEL, SDT_SYSGT, TRP_UPL, 0);
1013 #elif defined(__i386)
1014     set_gatesegd(&(brand_tbl[0].ih_interp_desc), &brand_sys_call,
1015     KCS_SEL, SDT_SYSGT, TRP_UPL, 0);
1016 #endif /* __i386 */

1044     brand_tbl[1].ih_inum = 0;
1045 }
    unchanged_portion_omitted

1055 #else /* __xpv */

1057 static void
1058 init_idt(gate_desc_t *idt)
1059 {
1060     char    ivctname[80];
1061     void    (*ivctptr)(void);
1062     int     i;

1064     /*
1065     * Initialize entire table with 'reserved' trap and then overwrite
1066     * specific entries. T_EXTTOVRFLT (9) is unsupported and reserved
1067     * since it can only be generated on a 386 processor. 15 is also
1068     * unsupported and reserved.
1069     */
1070 #if !defined(__xpv)
1071     for (i = 0; i < NIDT; i++) {
1072         set_gatesegd(&idt[i],
1073         (kpti_enable == 1) ? &tr_resvtrap : &resvtrap,
1074         KCS_SEL, SDT_SYSGT, TRP_KPL,
1075         idt_vector_to_ist(T_RESVTRAP));
1076     }
1077 #else
1078     for (i = 0; i < NIDT; i++) {
1079         for (i = 0; i < NIDT; i++)
1080             set_gatesegd(&idt[i], &resvtrap, KCS_SEL, SDT_SYSGT, TRP_KPL,
1081             IST_NONE);
1082 #endif
1046     };

1084     /*
1085     * 20-31 reserved
1086     */
1087 #if !defined(__xpv)
1088     for (i = 20; i < 32; i++) {
1089         set_gatesegd(&idt[i],
1090         (kpti_enable == 1) ? &tr_invaltrap : &invaltrap,
1091         KCS_SEL, SDT_SYSGT, TRP_KPL,
1092         idt_vector_to_ist(T_INVALTRAP));
1093     }
1094 #else
1095     for (i = 20; i < 32; i++) {
1051         for (i = 20; i < 32; i++)
1051             set_gatesegd(&idt[i], &invaltrap, KCS_SEL, SDT_SYSGT, TRP_KPL,
1096

```

```

1097             IST_NONE;
1098         }
1099     #endif
1053         0);

1101     /*
1102     * interrupts 32 - 255
1103     */
1104     for (i = 32; i < 256; i++) {
1105 #if !defined(__xpv)
1106         (void) snprintf(ivctname, sizeof (ivctname),
1107             (kpti_enable == 1) ? "tr_ivct%d" : "ivct%d", i);
1108 #else
1109         (void) snprintf(ivctname, sizeof (ivctname), "ivct%d", i);
1110 #endif
1111         ivctptr = (void (*)(void))kobj_getsymvalue(ivctname, 0);
1112         if (ivctptr == NULL)
1113             panic("kobj_getsymvalue(%s) failed", ivctname);

1115         set_gatesegd(&idt[i], ivctptr, KCS_SEL, SDT_SYSIGT, TRP_KPL,
1116             idt_vector_to_ist(i));
1064         set_gatesegd(&idt[i], ivctptr, KCS_SEL, SDT_SYSIGT, TRP_KPL, 0);
1117     }

1119     /*
1120     * Now install the common ones. Note that it will overlay some
1121     * entries installed above like T_SYSCALLINT, T_FASTTRAP etc.
1122     */
1123     init_idt_common(idt);
1124 }

    unchanged_portion_omitted

1144 #if !defined(__xpv)
1093 #if defined(__amd64)

1146 static void
1147 init_tss(void)
1148 {
1149     extern struct cpu cpus[];
1098     /*
1099     * tss_rsp0 is dynamically filled in by resume() on each context switch.
1100     * All exceptions but #DF will run on the thread stack.
1101     * Set up the double fault stack here.
1102     */
1103     ktss0->tss_ist1 =
1104         (uint64_t)&dblfault_stack0[sizeof (dblfault_stack0)];

1151     /*
1152     * tss_rsp0 is dynamically filled in by resume() (in swtch.s) on each
1153     * context switch but it'll be overwritten with this same value anyway.
1154     * Set I/O bit map offset equal to size of TSS segment limit
1155     * for no I/O permission map. This will force all user I/O
1156     * instructions to generate #gp fault.
1157     */
1155     if (kpti_enable == 1) {
1156         ktss0->tss_rsp0 = (uint64_t)&cpus->cpu_m.mcpu_kpti.kf_tr_rsp;
1157     }
1111     ktss0->tss_bitmapbase = sizeof (*ktss0);

1159     /* Set up the IST stacks for double fault, NMI, MCE. */
1160     ktss0->tss_ist1 = (uintptr_t)&dblfault_stack0[sizeof (dblfault_stack0)];
1161     ktss0->tss_ist2 = (uintptr_t)&nmi_stack0[sizeof (nmi_stack0)];
1162     ktss0->tss_ist3 = (uintptr_t)&mce_stack0[sizeof (mce_stack0)];
1113     /*
1114     * Point %tr to descriptor for ktss0 in gdt.
1115     */

```

```

1116         wr_tsr(KTSS_SEL);
1117     }

1119 #elif defined(__i386)

1121 static void
1122 init_tss(void)
1123 {
1164     /*
1165     * This IST stack is used for #DB,#BP (debug) interrupts (when KPTI is
1166     * enabled), and also for KDI (always).
1167     * ktss0->tss_esp dynamically filled in by resume() on each
1168     * context switch.
1169     */
1168     ktss0->tss_ist4 = (uint64_t)&cpus->cpu_m.mcpu_kpti_dbg.kf_tr_rsp;
1128     ktss0->tss_ss0 = KDS_SEL;
1129     ktss0->tss_eip = (uint32_t)_start;
1130     ktss0->tss_ds = ktss0->tss_es = ktss0->tss_ss = KDS_SEL;
1131     ktss0->tss_cs = KCS_SEL;
1132     ktss0->tss_fs = KFS_SEL;
1133     ktss0->tss_gs = KGS_SEL;
1134     ktss0->tss_ldt = ULDT_SEL;

1170     if (kpti_enable == 1) {
1171         /* This IST stack is used for #GP,#PF,#SS (fault) interrupts. */
1172         ktss0->tss_ist5 =
1173             (uint64_t)&cpus->cpu_m.mcpu_kptiflt.kf_tr_rsp;
1136     /*
1137     * Initialize double fault tss.
1138     */
1139     dftss0->tss_esp0 = (uint32_t)&dblfault_stack0[sizeof (dblfault_stack0)];
1140     dftss0->tss_ss0 = KDS_SEL;

1175         /* This IST stack is used for all other intrs (for KPTI). */
1176         ktss0->tss_ist6 = (uint64_t)&cpus->cpu_m.mcpu_kpti.kf_tr_rsp;
1177     }
1142     /*
1143     * tss_cr3 will get initialized in hat_kern_setup() once our page
1144     * tables have been setup.
1145     */
1146     dftss0->tss_eip = (uint32_t)syserrtrp;
1147     dftss0->tss_esp = (uint32_t)&dblfault_stack0[sizeof (dblfault_stack0)];
1148     dftss0->tss_cs = KCS_SEL;
1149     dftss0->tss_ds = KDS_SEL;
1150     dftss0->tss_es = KDS_SEL;
1151     dftss0->tss_ss = KDS_SEL;
1152     dftss0->tss_fs = KFS_SEL;
1153     dftss0->tss_gs = KGS_SEL;

1179     /*
1180     * Set I/O bit map offset equal to size of TSS segment limit
1181     * for no I/O permission map. This will force all user I/O
1182     * instructions to generate #gp fault.
1183     */
1184     ktss0->tss_bitmapbase = sizeof (*ktss0);

1186     /*
1187     * Point %tr to descriptor for ktss0 in gdt.
1188     */
1189     wr_tsr(KTSS_SEL);
1190 }

1168 #endif /* __i386 */
1192 #endif /* !__xpv */

1194 #if defined(__xpv)

```

```

1196 void
1197 init_desctbls(void)
1198 {
1199     uint_t vec;
1200     user_desc_t *gdt;
1201
1202     /*
1203      * Setup and install our GDT.
1204      */
1205     gdt = init_gdt();
1206
1207     /*
1208      * Store static pa of gdt to speed up pa_to_ma() translations
1209      * on lwp context switches.
1210      */
1211     ASSERT(IS_P2ALIGNED((uintptr_t)gdt, PAGESIZE));
1212     CPU->cpu_gdt = gdt;
1213     CPU->cpu_m.mcpu_gdtpa = pfn_to_pa(va_to_pfn(gdt));
1214
1215     /*
1216      * Setup and install our IDT.
1217      */
1218     #if !defined(__lint)
1219         ASSERT(NIDT * sizeof (*idt0) <= PAGESIZE);
1220     #endif
1221     idt0 = (gate_desc_t *)BOP_ALLOC(bootops, (caddr_t)IDT_VA,
1222                                     PAGESIZE, PAGESIZE);
1223     bzero(idt0, PAGESIZE);
1224     init_idt(idt0);
1225     for (vec = 0; vec < NIDT; vec++)
1226         xen_idt_write(&idt0[vec], vec);
1227
1228     CPU->cpu_idt = idt0;
1229
1230     /*
1231      * set default kernel stack
1232      */
1233     xen_stack_switch(KDS_SEL,
1234                     (ulong_t)&dblfault_stack0[sizeof (dblfault_stack0)]);
1235
1236     xen_init_callbacks();
1237
1238     init_ldt();
1239 }
1240
1241 #else /* __xpv */
1242
1243 void
1244 init_desctbls(void)
1245 {
1246     user_desc_t *gdt;
1247     desctbr_t idtr;
1248
1249     /*
1250      * Allocate IDT and TSS structures on unique pages for better
1251      * performance in virtual machines.
1252      */
1253     #if !defined(__lint)
1254         ASSERT(NIDT * sizeof (*idt0) <= PAGESIZE);
1255     #endif
1256     idt0 = (gate_desc_t *)BOP_ALLOC(bootops, (caddr_t)IDT_VA,
1257                                     PAGESIZE, PAGESIZE);
1258     bzero(idt0, PAGESIZE);
1259     #if !defined(__lint)
1260         ASSERT(sizeof (*ktss0) <= PAGESIZE);

```

```

1261 #endif
1262     ktss0 = (tss_t *)BOP_ALLOC(bootops, (caddr_t)KTSS_VA,
1263                               PAGESIZE, PAGESIZE);
1264     bzero(ktss0, PAGESIZE);
1265
1266     #if defined(__i386)
1267     #if !defined(__lint)
1268         ASSERT(sizeof (*dftss0) <= PAGESIZE);
1269     #endif
1270     dftss0 = (tss_t *)BOP_ALLOC(bootops, (caddr_t)DFTSS_VA,
1271                                 PAGESIZE, PAGESIZE);
1272     bzero(dftss0, PAGESIZE);
1273 #endif
1274
1275     /*
1276      * Setup and install our GDT.
1277      */
1278     gdt = init_gdt();
1279     ASSERT(IS_P2ALIGNED((uintptr_t)gdt, PAGESIZE));
1280     CPU->cpu_gdt = gdt;
1281
1282     /*
1283      * Initialize this CPU's LDT.
1284      */
1285     CPU->cpu_m.mcpu_ldt = BOP_ALLOC(bootops, (caddr_t)LDT_VA,
1286                                   LDT_CPU_SIZE, PAGESIZE);
1287     bzero(CPU->cpu_m.mcpu_ldt, LDT_CPU_SIZE);
1288     CPU->cpu_m.mcpu_ldt_len = 0;
1289
1290     /*
1291      * Setup and install our IDT.
1292      */
1293     init_idt(idt0);
1294
1295     idtr.dtr_base = (uintptr_t)idt0;
1296     idtr.dtr_limit = (NIDT * sizeof (*idt0)) - 1;
1297     wr_idtr(&idtr);
1298     CPU->cpu_idt = idt0;
1299
1300     #if defined(__i386)
1301     /*
1302      * We maintain a description of idt0 in convenient IDTR format
1303      * for #pf's on some older pentium processors. See pentium_pftap().
1304      */
1305     idt0_default_r = idtr;
1306     #endif /* __i386 */
1307
1308     init_tss();
1309     CPU->cpu_tss = ktss0;
1310     init_ldt();
1311
1312     /* Stash this so that the NMI,MCE,#DF and KDI handlers can use it. */
1313     kpti_safe_cr3 = (uint64_t)getcr3();
1314 }
1315
1316 unchanged_portion_omitted
1317
1318 void
1319 brand_interpositioning_enable(void)
1320 {
1321     gate_desc_t *idt = CPU->cpu_idt;
1322     int i;

```

```

1350     ASSERT(curthread->t_preempt != 0 || getpil() >= DISP_LEVEL);

1352     for (i = 0; brand_tbl[i].ih_inum; i++) {
1353         idt[brand_tbl[i].ih_inum] = brand_tbl[i].ih_interp_desc;
1354 #if defined(__xpv)
1355         xen_idt_write(&idt[brand_tbl[i].ih_inum],
1356                     brand_tbl[i].ih_inum);
1357 #endif
1358     }

1360 #if defined(__amd64)
1361 #if defined(__xpv)

1363     /*
1364     * Currently the hypervisor only supports 64-bit syscalls via
1365     * syscall instruction. The 32-bit syscalls are handled by
1366     * interrupt gate above.
1367     */
1368     xen_set_callback(brand_sys_syscall, CALLBACKTYPE_syscall,
1369                     CALLBACKF_mask_events);

1371 #else

1373     if (is_x86_feature(x86_featureset, X86FSET_ASYSC)) {
1374         if (kpti_enable == 1) {
1375             wrmsr(MSR_AMD_LSTAR, (uintptr_t)tr_brand_sys_syscall);
1376             wrmsr(MSR_AMD_CSTAR, (uintptr_t)tr_brand_sys_syscall32);
1377         } else {
1378             wrmsr(MSR_AMD_LSTAR, (uintptr_t)brand_sys_syscall);
1379             wrmsr(MSR_AMD_CSTAR, (uintptr_t)brand_sys_syscall32);
1380         }
1381     }

1383 #endif
1384 #endif /* __amd64 */

1386     if (is_x86_feature(x86_featureset, X86FSET_SEP)) {
1387         if (kpti_enable == 1) {
1388             wrmsr(MSR_INTC_SEP_EIP,
1389                 (uintptr_t)tr_brand_sys_sysenter);
1390         } else {
1391             if (is_x86_feature(x86_featureset, X86FSET_SEP))
1392                 wrmsr(MSR_INTC_SEP_EIP, (uintptr_t)brand_sys_sysenter);
1393         }
1394     }

1396 /*
1397  * Disable interpositioning on the system call path by rewriting the
1398  * sys{call|enter} MSRs and the syscall-related entries in the IDT to use
1399  * the standard entry points, which bypass the interpositioning hooks.
1400  */
1401 void
1402 brand_interpositioning_disable(void)
1403 {
1404     gate_desc_t     *idt = CPU->cpu_idt;
1405     int i;

1407     ASSERT(curthread->t_preempt != 0 || getpil() >= DISP_LEVEL);

1409     for (i = 0; brand_tbl[i].ih_inum; i++) {
1410         idt[brand_tbl[i].ih_inum] = brand_tbl[i].ih_default_desc;
1411 #if defined(__xpv)
1412         xen_idt_write(&idt[brand_tbl[i].ih_inum],
1413                     brand_tbl[i].ih_inum);

```

```

1414 #endif
1415     }

1417 #if defined(__amd64)
1418 #if defined(__xpv)

1420     /*
1421     * See comment above in brand_interpositioning_enable.
1422     */
1423     xen_set_callback(sys_syscall, CALLBACKTYPE_syscall,
1424                     CALLBACKF_mask_events);

1426 #else

1428     if (is_x86_feature(x86_featureset, X86FSET_ASYSC)) {
1429         if (kpti_enable == 1) {
1430             wrmsr(MSR_AMD_LSTAR, (uintptr_t)tr_sys_syscall);
1431             wrmsr(MSR_AMD_CSTAR, (uintptr_t)tr_sys_syscall32);
1432         } else {
1433             wrmsr(MSR_AMD_LSTAR, (uintptr_t)sys_syscall);
1434             wrmsr(MSR_AMD_CSTAR, (uintptr_t)sys_syscall32);
1435         }
1436     }

1438 #endif
1439 #endif /* __amd64 */

1441     if (is_x86_feature(x86_featureset, X86FSET_SEP)) {
1442         if (kpti_enable == 1) {
1443             wrmsr(MSR_INTC_SEP_EIP, (uintptr_t)tr_sys_sysenter);
1444         } else {
1445             if (is_x86_feature(x86_featureset, X86FSET_SEP))
1446                 wrmsr(MSR_INTC_SEP_EIP, (uintptr_t)sys_sysenter);
1447         }
1448     }

```

unchanged_portion_omitted

```

*****
21718 Fri Apr 6 17:25:12 2018
new/usr/src/uts/intel/ia32/os/sysi86.c
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22 * Copyright (c) 1992, 2010, Oracle and/or its affiliates. All rights reserved.
23 * Copyright 2018 Joyent, Inc.
24 */

26 /*      Copyright (c) 1990, 1991 UNIX System Laboratories, Inc. */
27 /*      Copyright (c) 1984, 1986, 1987, 1988, 1989, 1990 AT&T      */
28 /*      All Rights Reserved      */

30 /*      Copyright (c) 1987, 1988 Microsoft Corporation      */
31 /*      All Rights Reserved      */

33 #include <sys/param.h>
34 #include <sys/types.h>
35 #include <sys/sysmacros.h>
36 #include <sys/system.h>
37 #include <sys/signal.h>
38 #include <sys/errno.h>
39 #include <sys/fault.h>
40 #include <sys/syscall.h>
41 #include <sys/cpuvar.h>
42 #include <sys/sysi86.h>
43 #include <sys/psw.h>
44 #include <sys/cred.h>
45 #include <sys/policy.h>
46 #include <sys/thread.h>
47 #include <sys/debug.h>
48 #include <sys/onttrap.h>
49 #include <sys/privregs.h>
50 #include <sys/x86_archext.h>
51 #include <sys/vmem.h>
52 #include <sys/kmem.h>
53 #include <sys/mman.h>
54 #include <sys/archsystem.h>
55 #include <vm/hat.h>
56 #include <vm/as.h>
57 #include <vm/seg.h>
58 #include <vm/seg_kmem.h>
59 #include <vm/faultcode.h>

```

```

60 #include <sys/fp.h>
61 #include <sys/cmn_err.h>
62 #include <sys/segments.h>
63 #include <sys/clock.h>
64 #include <vm/hat_i86.h>
65 #if defined(__xpv)
66 #include <sys/hypervisor.h>
67 #include <sys/note.h>
68 #endif

70 static void ldt_alloc(proc_t *, uint_t);
71 static void ldt_free(proc_t *);
72 static void ldt_dup(proc_t *, proc_t *);
73 static void ldt_grow(proc_t *, uint_t);

75 /*
76  * sysi86 System Call
77  */

79 /* ARGSUSED */
80 int
81 sysi86(short cmd, uintptr_t arg1, uintptr_t arg2, uintptr_t arg3)
82 {
83     struct ssd ssd;
84     int error = 0;
85     int c;
86     proc_t *pp = curproc;

88     switch (cmd) {

90         /*
91          * The SI86V86 subsystem call of the SYSI86 system call
92          * supports only one subcode -- V86SC_IOPL.
93          */
94         case SI86V86:
95             if (arg1 == V86SC_IOPL) {
96                 struct regs *rp = lwptoregs(ttolwp(curthread));
97                 greg_t oldpl = rp->r_ps & PS_IOPL;
98                 greg_t newpl = arg2 & PS_IOPL;

100                 /*
101                  * Must be privileged to run this system call
102                  * if giving more io privilege.
103                  */
104                 if (newpl > oldpl && (error =
105                     secpolicy_sys_config(CRED(), B_FALSE)) != 0)
106                     return (set_errno(error));
107             #if defined(__xpv)
108                 kpreempt_disable();
109                 installctx(curthread, NULL, xen_disable_user_iopl,
110                     xen_enable_user_iopl, NULL, NULL,
111                     xen_disable_user_iopl, NULL);
112                 xen_enable_user_iopl();
113                 kpreempt_enable();
114             #else
115                 rp->r_ps ^= oldpl ^ newpl;
116             #endif

117             } else
118                 error = EINVAL;
119             break;

121         /*
122          * Set a segment descriptor
123          */
124         case SI86DSCR:
125             /*

```

```

126     * There are considerable problems here manipulating
127     * resources shared by many running lwps.  Get everyone
128     * into a safe state before changing the LDT.
129     */
130     if (curthread != pp->p_agenttp && !holdlwps(SHOLDFORK1)) {
131         error = EINTR;
132         break;
133     }
134
135     if (get_udatamodel() == DATAMODEL_LP64) {
136         error = EINVAL;
137         break;
138     }
139
140     if (copyin((caddr_t)arg1, &ssd, sizeof (ssd)) < 0) {
141         error = EFAULT;
142         break;
143     }
144
145     error = setdscr(&ssd);
146
147     mutex_enter(&pp->p_lock);
148     if (curthread != pp->p_agenttp)
149         continuelwps(pp);
150     mutex_exit(&pp->p_lock);
151     break;
152
153     case SI86FPHW:
154         c = fp_kind & 0xff;
155         if (suword32((void *)arg1, c) == -1)
156             error = EFAULT;
157         break;
158
159     case SI86FPSTART:
160         /*
161          * arg1 is the address of _fp_hw
162          * arg2 is the desired x87 FCW value
163          * arg3 is the desired SSE MXCSR value
164          * a return value of one means SSE hardware, else none.
165          */
166         c = fp_kind & 0xff;
167         if (suword32((void *)arg1, c) == -1) {
168             error = EFAULT;
169             break;
170         }
171         fpsetcw((uint16_t)arg2, (uint32_t)arg3);
172         return ((fp_kind & __FP_SSE) ? 1 : 0);
173
174     /* real time clock management commands */
175
176     case WTODC:
177         if ((error = secpolicy_settime(CRED())) == 0) {
178             timestruc_t ts;
179             mutex_enter(&tod_lock);
180             gethrestime(&ts);
181             tod_set(ts);
182             mutex_exit(&tod_lock);
183         }
184         break;
185
186     /* Give some timezone playing room */
187     #define ONEWEEK (7 * 24 * 60 * 60)
188
189     case SGMTL:
190         /*
191          * Called from 32 bit land, negative values

```

```

192     * are not sign extended, so we do that here
193     * by casting it to an int and back.  We also
194     * clamp the value to within reason and detect
195     * when a 64 bit call overflows an int.
196     */
197     if ((error = secpolicy_settime(CRED())) == 0) {
198         int newlag = (int)arg1;
199
200     #ifdef _SYSCALL32_IMPL
201         if (get_udatamodel() == DATAMODEL_NATIVE &&
202             (long)newlag != (long)arg1) {
203             error = EOVERFLOW;
204         } else
205     #endif
206         if (newlag >= -ONEWEEK && newlag <= ONEWEEK)
207             sgmtl(newlag);
208         else
209             error = EOVERFLOW;
210     }
211     break;
212
213     case GGMTL:
214         if (get_udatamodel() == DATAMODEL_NATIVE) {
215             if (sulword((void *)arg1, ggmtl()) == -1)
216                 error = EFAULT;
217     #ifdef _SYSCALL32_IMPL
218             } else {
219                 time_t gmtl;
220
221                 if ((gmtl = ggmtl()) > INT32_MAX) {
222                     /*
223                      * Since gmt_lag can at most be
224                      * +/- 12 hours, something is
225                      * *seriously* messed up here.
226                      */
227                     error = EOVERFLOW;
228                 } else if (suword32((void *)arg1, (int32_t)gmtl) == -1)
229                     error = EFAULT;
230     #endif
231             }
232         break;
233
234     case RTCSYNC:
235         if ((error = secpolicy_settime(CRED())) == 0)
236             rtcsync();
237         break;
238
239     /* END OF real time clock management commands */
240
241     default:
242         error = EINVAL;
243         break;
244     }
245     return (error == 0 ? 0 : set_errno(error));
246 }
247
248 unchanged_portion_omitted
249
250 #endif /* __i386 */
251
252 /*
253  * Load LDT register with the current process's LDT.
254  */
255 static void
256 ldt_load(void)
257 {
258     #if defined(__xpv)

```

```

348     xen_set_ldt(get_ssd_base(&curproc->p_ldt_desc),
349               curproc->p_ldtlimit + 1);
350 #else
351     size_t len;
352     system_desc_t desc;
353
354     /*
355      * Before we can use the LDT on this CPU, we must install the LDT in the
356      * user mapping table.
357      */
358     len = (curproc->p_ldtlimit + 1) * sizeof (user_desc_t);
359     bcopy(curproc->p_ldt, CPU->cpu_m.mcpu_ldt, len);
360     CPU->cpu_m.mcpu_ldt_len = len;
361     set_syssegd(&desc, CPU->cpu_m.mcpu_ldt, len - 1, SDT_SYSLDT, SEL_KPL);
362     *((system_desc_t *)&CPU->cpu_gdt[GDT_LDT]) = desc;
363
364     *((system_desc_t *)&CPU->cpu_gdt[GDT_LDT]) = curproc->p_ldt_desc;
365     wr_ldtr(ULDT_SEL);
366 #endif
367 }
368
369 /*
370 * Store a NULL selector in the LDTR. All subsequent illegal references to
371 * the LDT will result in a #gp.
372 */
373 void
374 ldt_unload(void)
375 {
376     #if defined(__xpv)
377     xen_set_ldt(NULL, 0);
378     #else
379     *((system_desc_t *)&CPU->cpu_gdt[GDT_LDT]) = null_sdesc;
380     wr_ldtr(0);
381
382     bzero(CPU->cpu_m.mcpu_ldt, CPU->cpu_m.mcpu_ldt_len);
383     CPU->cpu_m.mcpu_ldt_len = 0;
384 #endif
385 }
386
387 unchanged_portion_omitted
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739

```

```

741     ldt = kmem_zalloc(ldtsz, KM_SLEEP);
742     ASSERT(IS_P2ALIGNED(ldt, PAGE_SIZE));
743
744 #if defined(__xpv)
745     if (xen_ldt_setprot(ldt, ldtsz, PROT_READ))
746         panic("ldt_alloc:xen_ldt_setprot(PROT_READ) failed");
747 #endif
748
749     pp->p_ldt = ldt;
750     pp->p_ldtlimit = nsels - 1;
751     set_syssegd(&pp->p_ldt_desc, ldt, ldtsz - 1, SDT_SYSLDT, SEL_KPL);
752
753     if (pp == curproc) {
754         kpreempt_disable();
755         ldt_load();
756         kpreempt_enable();
757     }
758 }
759
760 unchanged_portion_omitted
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000

```

```
885         panic("ldt_grow:xen_ldt_setprot(PROT_READ) failed");
886 #endif
888     pp->p_ldt = nldt;
889     pp->p_ldtlimit = nsels - 1;
891     /*
892     * write new ldt segment descriptor.
893     */
894     set_syssegd(&pp->p_ldt_desc, nldt, nldtsz - 1, SDT_SYSLDT, SEL_KPL);
896     /*
897     * load the new ldt.
898     */
899     kpreempt_disable();
900     ldt_load();
901     kpreempt_enable();
903     kmem_free(oldt, oldtsz);
904 }
_____unchanged_portion_omitted_____
```


new/usr/src/uts/intel/ia32/sys/kdi_regs.h

1

```
*****
1850 Fri Apr 6 17:25:12 2018
new/usr/src/uts/intel/ia32/sys/kdi_regs.h
9210 remove KMDB branch debugging support
9211 ::crregs could do with cr2/cr3 support
9209 ::ttrace should be able to filter by thread
Reviewed by: Patrick Mooney <patrick.mooney@joyent.com>
Reviewed by: Yuri Pankov <yuripv@yuripv.net>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */

22 /*
23  * Copyright 2007 Sun Microsystems, Inc. All rights reserved.
24  * Use is subject to license terms.
25  *
26  * Copyright 2018 Joyent, Inc.
27  */

29 #ifndef _IA32_SYS_KDI_REGS_H
30 #define _IA32_SYS_KDI_REGS_H

30 #pragma ident "%Z%M% %I% %E% SMI"

32 #ifdef __cplusplus
33 extern "C" {
34 #endif

36 #define KDIREG_NGREG 21

38 /*
39  * %ss appears in a different place than a typical struct regs, since the
40  * machine won't save %ss on a trap entry from the same privilege level.
41  */

43 #define KDIREG_SAVFP 0
44 #define KDIREG_SAVPC 1
45 #define KDIREG_SS 2
46 #define KDIREG_GS 3
47 #define KDIREG_FS 4
48 #define KDIREG_ES 5
49 #define KDIREG_DS 6
50 #define KDIREG EDI 7
51 #define KDIREG_ESI 8
52 #define KDIREG_EBP 9
53 #define KDIREG_ESP 10
54 #define KDIREG_EBX 11
55 #define KDIREG_EDX 12
```

new/usr/src/uts/intel/ia32/sys/kdi_regs.h

2

```
56 #define KDIREG_ECX 13
57 #define KDIREG_EAX 14
58 #define KDIREG_TRAPNO 15
59 #define KDIREG_ERR 16
60 #define KDIREG_EIP 17
61 #define KDIREG_CS 18
62 #define KDIREG_EFLAGS 19
63 #define KDIREG_UESP 20

65 #define KDIREG_PC KDIREG_EIP
66 #define KDIREG_SP KDIREG_ESP
67 #define KDIREG_FP KDIREG_EBP

69 #ifdef _ASM

71 /* Patch point for MSR clearing. */
72 #define KDI_MSR_PATCH \
73     nop; nop; nop; nop; \
74     nop; nop; nop; nop; \
75     nop; nop; nop; nop; \
76     nop

78 #endif /* _ASM */

80 #define KDI_MSR_PATCHOFF 8 /* bytes of code before patch point */
81 #define KDI_MSR_PATCHSZ 13 /* bytes in KDI_MSR_PATCH, above */

69 #ifdef __cplusplus
70 }
_____unchanged_portion_omitted_____
```

new/usr/src/uts/intel/ia32/sys/trap.h

1

```
*****
4210 Fri Apr 6 17:25:12 2018
new/usr/src/uts/intel/ia32/sys/trap.h
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*      Copyright (c) 1990, 1991 UNIX System Laboratories, Inc. */
22 /*      Copyright (c) 1984, 1986, 1987, 1988, 1989, 1990 AT&T */
23 /*      All Rights Reserved */
24
25 /*
26  * Copyright 2010 Sun Microsystems, Inc. All rights reserved.
27  * Use is subject to license terms.
28  *
29  * Copyright 2018 Joyent, Inc.
30  */
31
32 #ifndef __IA32_SYS_TRAP_H
33 #define __IA32_SYS_TRAP_H
34
35 #ifndef __cplusplus
36 extern "C" {
37 #endif
38
39 /*
40  * Trap type values
41  */
42
43 #define T_ZERODIV      0x0    /* #de divide by 0 error */
44 #define T_SGLSTP      0x1    /* #db single step */
45 #define T_NMIFLT      0x2    /* #nm NMI */
46 #define T_BPTFLT      0x3    /* #bp breakpoint fault, INT3 insn */
47 #define T_OVFLW      0x4    /* #of INTO overflow fault */
48 #define T_BOUNDFLT    0x5    /* #br BOUND insn fault */
49 #define T_ILLNST      0x6    /* #ud invalid opcode fault */
50 #define T_NOEXTFLT    0x7    /* #nm device not available: x87 */
51 #define T_DBLFLT      0x8    /* #df double fault */
52 #define T_EXTOVRFLT   0x9    /* [not generated: 386 only] */
53 #define T_TSSFLT      0xa    /* #ts invalid TSS fault */
54 #define T_SEGFLT      0xb    /* #np segment not present fault */
55 #define T_STKFLT      0xc    /* #ss stack fault */
56 #define T_GPFLT      0xd    /* #gp general protection fault */
57 #define T_PGFLT      0xe    /* #pf page fault */
58 #define T_RESVTRAP    0xf    /* reserved */
59 #define T_EXTERRFLT   0x10   /* #mf x87 FPU error fault */
```

new/usr/src/uts/intel/ia32/sys/trap.h

2

```
60 #define T_ALIGNMENT  0x11   /* #ac alignment check error */
61 #define T_MCE         0x12   /* #mc machine check exception */
62 #define T_SIMDFPE     0x13   /* #xm SSE/SSE exception */
63 #define T_DBGENTR     0x14   /* #db debugger entry */
64 #define T_INVALIDTRAP 0x1e   /* invalid */
65 #define T_ENDPERR     0x21   /* emulated extension error flt */
66 #define T_ENOEXTFLT   0x20   /* emulated ext not present */
67 #define T_FASTTRAP    0xd2   /* fast system call */
68 #define T_SYSCALLINT  0x91   /* general system call */
69 #define T_DTRACE_RET  0x92   /* DTrace pid return */
70 #define T_INTR80      0x80   /* int80 handler for linux emulation */
71 #define T_SOFTINT     0x50fd  /* pseudo softint trap type */
72
73 /*
74  * Pseudo traps.
75  */
76 #define T_INTERRUPT   0x100
77 #define T_FAULT       0x200
78 #define T_AST         0x400
79 #define T_SYSCALL     0x180
80
81
82 /*
83  * Values of error code on stack in case of page fault
84  */
85
86 #define PF_ERR_MASK    0x01   /* Mask for error bit */
87 #define PF_ERR_PAGE    0x00   /* page not present */
88 #define PF_ERR_PROT    0x01   /* protection error */
89 #define PF_ERR_WRITE   0x02   /* fault caused by write (else read) */
90 #define PF_ERR_USER    0x04   /* processor was in user mode */
91                               /* (else supervisor) */
92 #define PF_ERR_EXEC    0x10   /* attempt to execute a No eXec page (AMD) */
93                               /* or kernel tried to execute a user page */
94                               /* (Intel SMEP) */
95
96 /*
97  * Definitions for fast system call subfunctions
98  */
99 #define T_FNULL        0      /* Null trap for testing */
100 #define T_FGETFP       1      /* Get emulated FP context */
101 #define T_FSETFP       2      /* Set emulated FP context */
102 #define T_GETHRTIME    3      /* Get high resolution time */
103 #define T_GETHRVTIME   4      /* Get high resolution virtual time */
104 #define T_GETHRESTIME  5      /* Get high resolution time */
105 #define T_GETLGRP      6      /* Get home lgrp */
106
107 #define T_LASTFAST     6      /* Last valid subfunction */
108
109 /*
110  * Offsets for an interrupt/trap frame.
111  */
112 #define T_FRAME_ERR    0
113 #define T_FRAME_RIP    8
114 #define T_FRAME_CS     16
115 #define T_FRAME_RFLAGS 24
116 #define T_FRAME_RSP    32
117 #define T_FRAME_SS     40
118
119 #define T_FRAMERET_RIP  0
120 #define T_FRAMERET_CS   8
121 #define T_FRAMERET_RFLAGS 16
122 #define T_FRAMERET_RSP 24
123 #define T_FRAMERET_SS  32
124
125 #ifndef __cplusplus
```

new/usr/src/uts/intel/ia32/sys/trap.h

3

126 }

unchanged_portion_omitted

```

*****
17745 Fri Apr 6 17:25:12 2018
new/usr/src/uts/intel/kdi/kdi_asm.s
9441 kmdb should stash %cr3 in kdiregs
Reviewed by: John Levon <john.levon@joyent.com>
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */

22 /*
23  * Copyright 2007 Sun Microsystems, Inc. All rights reserved.
24  * Use is subject to license terms.
25  *
26  * Copyright 2018 Joyent, Inc.
27  */

27 #pragma ident      "%Z%M% %I%      %E% SMI"

29 /*
30  * Debugger entry and exit for both master and slave CPUs. kdi_idthdl.s contains
31  * the IDT stubs that drop into here (mainly via kdi_cmmint).
32  * Debugger entry for both master and slave CPUs
33  */

34 #if defined(__lint)
35 #include <sys/types.h>
36 #else
37 #endif

38 #include <sys/segments.h>
39 #include <sys/asm_linkage.h>
40 #include <sys/controlregs.h>
41 #include <sys/x86_archext.h>
42 #include <sys/privregs.h>
43 #include <sys/machprivregs.h>
44 #include <sys/kdi_regs.h>
45 #include <sys/psw.h>
46 #include <sys/uadmin.h>
47 #ifdef __xpv
48 #include <sys/hypervisor.h>
49 #endif

50 #ifdef _ASM
51 #include <kdi_assym.h>

```

```

51 #include <assym.h>

53 /* clobbers %rdx, %rcx, returns addr in %rax, CPU ID in %rbx */
54 #define GET_CPUSAVE_ADDR \
55     movzbq  %gs:CPU_ID, %rbx;          \
56     movq    %rbx, %rax;                \
57     movq    $KRS_SIZE, %rcx;           \
58     mulq    %rcx;                       \
59     movq    $kdi_cpusave, %rdx;        \
60     /*CSTYLED*/                          \
61     addq    (%rdx), %rax

63 /*
64  * Save copies of the IDT and GDT descriptors. Note that we only save the IDT
65  * and GDT if the IDT isn't ours, as we may be legitimately re-entering the
66  * debugger through the trap handler. We don't want to clobber the saved IDT
67  * in the process, as we'd end up resuming the world on our IDT.
68  */
69 #define SAVE_IDTGDT \
70     movq    %gs:CPU_IDT, %r11;          \
71     leaq   kdi_idt(%rip), %rsi;         \
72     cmpq   %rsi, %r11;                  \
73     je     lf;                           \
74     movq   %r11, KRS_IDT(%rax);         \
75     movq   %gs:CPU_GDT, %r11;           \
76     movq   %r11, KRS_GDT(%rax);        \
77 lf:

79 #ifdef __xpv

81 /*
82  * Already on kernel gsbase via the hypervisor.
83  */
84 #define SAVE_GSBASE(reg) /* nothing */
85 #define RESTORE_GSBASE(reg) /* nothing */

87 #else

89 #define SAVE_GSBASE(base) \
90     movl   $MSR_AMD_GSBASE, %ecx;       \
91     rdmsr;                               \
92     shlq   $32, %rdx;                    \
93     orq    %rax, %rdx;                   \
94     movq   %rdx, REG_OFF(KDIREG_GSBASE)(base); \
95     movl   $MSR_AMD_KGSBASE, %ecx;       \
96     rdmsr;                               \
97     shlq   $32, %rdx;                    \
98     orq    %rax, %rdx;                   \
99     movq   %rdx, REG_OFF(KDIREG_KGSBASE)(base) \
100     movq   %rdx, REG_OFF(KDIREG_GSBASE)(base)

101 /*
102  * We shouldn't have stomped on KGSBASE, so don't try to restore it.
103  */
104 #define RESTORE_GSBASE(base) \
105     movq   REG_OFF(KDIREG_GSBASE)(base), %rdx; \
106     movq   %rdx, %rax;                   \
107     shrq   $32, %rdx;                     \
108     movl   $MSR_AMD_GSBASE, %ecx;         \
109     wrmsr

111 #endif /* __xpv */

113 /*
114  * %ss, %rsp, %rflags, %cs, %rip, %err, %trapno are already on the stack.
115  * %ss, %rsp, %rflags, %cs, %rip, %err, %trapno are already on the stack. Note

```

```

106 * that on the hypervisor, we skip the save/restore of GSBASE: it's slow, and
107 * unnecessary.
115 */
116 #define KDI_SAVE_REGS(base) \
117     movq   %rdi, REG_OFF(KDIREG_RDI)(base); \
118     movq   %rsi, REG_OFF(KDIREG_RSI)(base); \
119     movq   %rdx, REG_OFF(KDIREG_RDX)(base); \
120     movq   %rcx, REG_OFF(KDIREG_RCX)(base); \
121     movq   %r8, REG_OFF(KDIREG_R8)(base); \
122     movq   %r9, REG_OFF(KDIREG_R9)(base); \
123     movq   %rax, REG_OFF(KDIREG_RAX)(base); \
124     movq   %rbx, REG_OFF(KDIREG_RBX)(base); \
125     movq   %rbp, REG_OFF(KDIREG_RBP)(base); \
126     movq   %r10, REG_OFF(KDIREG_R10)(base); \
127     movq   %r11, REG_OFF(KDIREG_R11)(base); \
128     movq   %r12, REG_OFF(KDIREG_R12)(base); \
129     movq   %r13, REG_OFF(KDIREG_R13)(base); \
130     movq   %r14, REG_OFF(KDIREG_R14)(base); \
131     movq   %r15, REG_OFF(KDIREG_R15)(base); \
132     movq   %rbp, REG_OFF(KDIREG_SAVFP)(base); \
133     movq   REG_OFF(KDIREG_RIP)(base), %rax; \
134     movq   %rax, REG_OFF(KDIREG_SAVPC)(base); \
135     movq   %cr2, %rax; \
136     movq   %rax, REG_OFF(KDIREG_CR2)(base); \
137     clrq   %rax; \
138     movw   %ds, %ax; \
139     movq   %rax, REG_OFF(KDIREG_DS)(base); \
140     movw   %es, %ax; \
141     movq   %rax, REG_OFF(KDIREG_ES)(base); \
142     movw   %fs, %ax; \
143     movq   %rax, REG_OFF(KDIREG_FS)(base); \
144     movw   %gs, %ax; \
145     movq   %rax, REG_OFF(KDIREG_GS)(base); \
146     SAVE_GSBASE(base)

148 #define KDI_RESTORE_REGS(base) \
149     movq   base, %rdi; \
150     RESTORE_GSBASE(%rdi); \
151     movq   REG_OFF(KDIREG_ES)(%rdi), %rax; \
152     movw   %ax, %es; \
153     movq   REG_OFF(KDIREG_DS)(%rdi), %rax; \
154     movw   %ax, %ds; \
155     movq   REG_OFF(KDIREG_CR2)(base), %rax; \
156     movq   %rax, %cr2; \
157     movq   REG_OFF(KDIREG_R15)(%rdi), %r15; \
158     movq   REG_OFF(KDIREG_R14)(%rdi), %r14; \
159     movq   REG_OFF(KDIREG_R13)(%rdi), %r13; \
160     movq   REG_OFF(KDIREG_R12)(%rdi), %r12; \
161     movq   REG_OFF(KDIREG_R11)(%rdi), %r11; \
162     movq   REG_OFF(KDIREG_R10)(%rdi), %r10; \
163     movq   REG_OFF(KDIREG_RBP)(%rdi), %rbp; \
164     movq   REG_OFF(KDIREG_RBX)(%rdi), %rbx; \
165     movq   REG_OFF(KDIREG_RAX)(%rdi), %rax; \
166     movq   REG_OFF(KDIREG_R9)(%rdi), %r9; \
167     movq   REG_OFF(KDIREG_R8)(%rdi), %r8; \
168     movq   REG_OFF(KDIREG_RCX)(%rdi), %rcx; \
169     movq   REG_OFF(KDIREG_RDX)(%rdi), %rdx; \
170     movq   REG_OFF(KDIREG_RSI)(%rdi), %rsi; \
171     movq   REG_OFF(KDIREG_RDI)(%rdi), %rdi

173 /*
174 * Given the address of the current CPU's cpusave area in %rax, the following
175 * macro restores the debugging state to said CPU. Restored state includes
176 * the debug registers from the global %dr variables.
177 * the debug registers from the global %dr variables, and debugging MSRs from
178 * the CPU save area. This code would be in a separate routine, but for the

```

```

167 * fact that some of the MSRs are jump-sensitive. As such, we need to minimize
168 * the number of jumps taken subsequent to the update of said MSRs. We can
169 * remove one jump (the ret) by using a macro instead of a function for the
170 * debugging state restoration code.
171 */
172 * Takes the cpusave area in %rdi as a parameter.
173 * Takes the cpusave area in %rdi as a parameter, clobbers %rax-%rdx
174 */
175 #define KDI_RESTORE_DEBUGGING_STATE \
176     pushq  %rdi; \
177     leaq   kdi_dreg(%rip), %r15; \
178     movl   $7, %edi; \
179     movq   DR_CTL(%r15), %rsi; \
180     call   kdi_dreg_set; \
181     movl   $6, %edi; \
182     movq   $KDIREG_DRSTAT_RESERVED, %rsi; \
183     call   kdi_dreg_set; \
184     movl   $0, %edi; \
185     movq   DRADDR_OFF(0)(%r15), %rsi; \
186     call   kdi_dreg_set; \
187     movl   $1, %edi; \
188     movq   DRADDR_OFF(1)(%r15), %rsi; \
189     call   kdi_dreg_set; \
190     movl   $2, %edi; \
191     movq   DRADDR_OFF(2)(%r15), %rsi; \
192     call   kdi_dreg_set; \
193     movl   $3, %edi; \
194     movq   DRADDR_OFF(3)(%r15), %rsi; \
195     call   kdi_dreg_set; \
196     popq   %rdi; \
197     popq   %rdi; \
198     /* \
199     * Write any requested MSRs. \
200     */ \
201     movq   KRS_MSR(%rdi), %rbx; \
202     cmpq   $0, %rbx; \
203     je     3f; \
204     1: \
205     movl   MSR_NUM(%rbx), %ecx; \
206     cmpl   $0, %ecx; \
207     je     3f; \
208     movl   MSR_TYPE(%rbx), %edx; \
209     cmpl   $KDI_MSR_WRITE, %edx; \
210     jne   2f; \
211     movq   MSR_VALP(%rbx), %rdx; \
212     movl   0(%rdx), %eax; \
213     movl   4(%rdx), %edx; \
214     wrmsr; \
215     2: \
216     addq   $MSR_SIZE, %rbx; \
217     jmp   1b; \
218     3: \
219     /* \
220     * We must not branch after re-enabling LBR. If \
221     * kdi_wsr_wrexite_msr is set, it contains the number \
222     * of the MSR that controls LBR. kdi_wsr_wrexite_valp \
223     * contains the value that is to be written to enable \
224     * LBR. \
225     */ \
226     leaq   kdi_msr_wrexite_msr(%rip), %rcx; \
227     movl   (%rcx), %ecx;

```

```

231     cmpl     $0, %ecx;
232     je      1f;
233
234     leaq    kdi_msr_wrexite_valp(%rip), %rdx;
235     movq    (%rdx), %rdx;
236     movl    0(%rdx), %eax;
237     movl    4(%rdx), %edx;
238
239     wrmsr;
240 1:

```

```

205 /*
206  * Each cpusave buffer has an area set aside for a ring buffer of breadcrumbs.
207  * The following macros manage the buffer.
208  */

```

```

210 /* Advance the ring buffer */
211 #define ADVANCE_CRUMB_POINTER(cpusave, tmp1, tmp2) \
212     movq    KRS_CURCRUMBIDX(cpusave), tmp1; \
213     cmpq    $[KDI_NCRUMBS - 1], tmp1; \
214     jge     1f;
215 /* Advance the pointer and index */
216     addq    $1, tmp1;
217     movq    tmp1, KRS_CURCRUMBIDX(cpusave); \
218     movq    KRS_CURCRUMB(cpusave), tmp1; \
219     addq    $KRM_SIZE, tmp1; \
220     jmp     2f;
221 1: /* Reset the pointer and index */
222     movq    $0, KRS_CURCRUMBIDX(cpusave); \
223     leaq    KRS_CRUMBS(cpusave), tmp1; \
224 2:     movq    tmp1, KRS_CURCRUMB(cpusave); \
225 /* Clear the new crumb */
226     movq    $KDI_NCRUMBS, tmp2; \
227 3:     movq    $0, -4(tmp1, tmp2, 4); \
228     decq    tmp2; \
229     jnz     3b

```

```

231 /* Set a value in the current breadcrumb buffer */
232 #define ADD_CRUMB(cpusave, offset, value, tmp) \
233     movq    KRS_CURCRUMB(cpusave), tmp; \
234     movq    value, offset(tmp)

```

```

273 #endif /* _ASM */

```

```

275 #if defined(__lint)
276 void
277 kdi_cmrint(void)
278 {
279 }
280 #else /* __lint */

```

```

236 /* XXX implement me */
237 ENTRY_NP(kdi_nmiint)
238     clrq    %rcx
239     movq    (%rcx), %rcx
240     SET_SIZE(kdi_nmiint)

```

```

242 /*
243  * The main entry point for master CPUs. It also serves as the trap
244  * handler for all traps and interrupts taken during single-step.
245  */
246 ENTRY_NP(kdi_cmrint)
247 ALTENTRY(kdi_master_entry)

```

```

249     pushq   %rax
250     CLI(%rax)

```

```

251     popq    %rax

```

```

253 /* Save current register state */
254     subq    $REG_OFF(KDIREG_TRAPNO), %rsp
255     KDI_SAVE_REGS(%rsp)

```

```

257 #ifdef __xpv
258 /*
259  * Clear saved_upcall_mask in unused byte of cs slot on stack.
260  * It can only confuse things.
261  */
262     movb    $0, REG_OFF(KDIREG_CS)+4(%rsp)
263 #endif

```

```

265 #if !defined(__xpv)
266 /*
267  * Switch to the kernel's GSBASE. Neither GSBASE nor the ill-named
268  * KGSBASE can be trusted, as the kernel may or may not have already
269  * done a swappgs. All is not lost, as the kernel can divine the correct
270  * value for us. Note that the previous GSBASE is saved in the
271  * KDI_SAVE_REGS macro to prevent a usermode process's GSBASE from being
272  * blown away. On the hypervisor, we don't need to do this, since it's
273  * ensured we're on our requested kernel GSBASE already.
274  */
275     subq    $10, %rsp
276     sgdt    (%rsp)
277     movq    2(%rsp), %rdi /* gdt base now in %rdi */
278     addq    $10, %rsp
279     call   kdi_gdt2gsbase /* returns kernel's GSBASE in %rax */

```

```

281     movq    %rax, %rdx
282     shrq    $32, %rdx
283     movl    $MSR_AMD_GSBASE, %ecx
284     wrmsr

```

```

286 /*
287  * In the trampoline we stashed the incoming %cr3. Copy this into
288  * the kdiregs for restoration and later use.
289  */
290     mov     %gs:(CPU_KPTI_DBG+KPTI_TR_CR3), %rdx
291     mov     %rdx, REG_OFF(KDIREG_CR3)(%rsp)
292 /*
293  * Switch to the kernel's %cr3. From the early interrupt handler
294  * until now we've been running on the "paranoid" %cr3 (that of kas
295  * from early in boot).
296  *
297  * If we took the interrupt from somewhere already on the kas/paranoid
298  * %cr3 though, don't change it (this could happen if kcr3 is corrupt
299  * and we took a gptrap earlier from this very code).
300  */
301     cmpq    %rdx, kpti_safe_cr3
302     je      .no_kcr3
303     mov     %gs:CPU_KPTI_KCR3, %rdx
304     cmpq    $0, %rdx
305     je      .no_kcr3
306     mov     %rdx, %cr3
307 .no_kcr3:

```

```

309 #endif /* __xpv */

```

```

311     GET_CPUSAVE_ADDR /* %rax = cpusave, %rbx = CPU ID */

```

```

313     ADVANCE_CRUMB_POINTER(%rax, %rcx, %rdx)

```

```

315     ADD_CRUMB(%rax, KRM_CPU_STATE, $KDI_CPU_STATE_MASTER, %rdx)

```

```

317     movq    REG_OFF(KDIREG_RIP)(%rsp), %rcx
318     ADD_CRUMB(%rax, KRM_PC, %rcx, %rdx)
319     ADD_CRUMB(%rax, KRM_SP, %rsp, %rdx)
320     movq    REG_OFF(KDIREG_TRAPNO)(%rsp), %rcx
321     ADD_CRUMB(%rax, KRM_TRAPNO, %rcx, %rdx)

323     movq    %rsp, %rbp
324     pushq   %rax

326     /*
327     * Were we in the debugger when we took the trap (i.e. was %esp in one
328     * of the debugger's memory ranges)?
329     */
330     leaq    kdi_memranges, %rcx
331     movl    kdi_nmemranges, %edx
332 1:
333     cmpq    MR_BASE(%rcx), %rsp
334 1:     cmpq    MR_BASE(%rcx), %rsp
335     jl     2f          /* below this range -- try the next one */
336     cmpq    MR_LIM(%rcx), %rsp
337     jg     2f          /* above this range -- try the next one */
338     jmp    3f          /* matched within this range */

339 2:
340     decl    %edx
341     decl    %edx
342     jz     kdi_save_common_state /* %rsp not within debugger memory */
343     addq   $MR_SIZE, %rcx
344     jmp    1b

345 3:
346     /*
347     * The master is still set. That should only happen if we hit a trap
348     * while running in the debugger. Note that it may be an intentional
349     * fault. kmdb_dpi_handle_fault will sort it all out.
350     */

351     movq    REG_OFF(KDIREG_TRAPNO)(%rbp), %rdi
352     movq    REG_OFF(KDIREG_RIP)(%rbp), %rsi
353     movq    REG_OFF(KDIREG_RSP)(%rbp), %rdx
354     movq    %rbx, %rcx /* cpuid */

356     call   kdi_dvec_handle_fault

358     /*
359     * If we're here, we ran into a debugger problem, and the user
360     * elected to solve it by having the debugger debug itself. The
361     * state we're about to save is that of the debugger when it took
362     * the fault.
363     */

365     jmp    kdi_save_common_state

367     SET_SIZE(kdi_master_entry)
368     unchanged_portion_omitted

390 #endif /* __lint */

370 /*
371 * The cross-call handler for slave CPUs.
372 *
373 * The debugger is single-threaded, so only one CPU, called the master, may be
374 * running it at any given time. The other CPUs, known as slaves, spin in a
375 * busy loop until there's something for them to do. This is the entry point
376 * for the slaves - they'll be sent here in response to a cross-call sent by the
377 * master.
378 */

```

```

402 #if defined(__lint)
403 char kdi_slave_entry_patch;

405 void
406 kdi_slave_entry(void)
407 {
408 }
409 #else /* __lint */
410     .globl kdi_slave_entry_patch;

380     ENTRY_NP(kdi_slave_entry)

414     /* kdi_msr_add_clrentry knows where this is */
415 kdi_slave_entry_patch:
416     KDI_MSR_PATCH;

382     /*
383     * Cross calls are implemented as function calls, so our stack currently
384     * looks like one you'd get from a zero-argument function call. That
385     * is, there's the return %rip at %rsp, and that's about it. We need
386     * to make it look like an interrupt stack. When we first save, we'll
387     * reverse the saved %ss and %rip, which we'll fix back up when we've
388     * freed up some general-purpose registers. We'll also need to fix up
389     * the saved %rsp.
390     */

392     pushq   %rsp          /* pushed value off by 8 */
393     pushfq
394     CLI(%rax)
395     pushq   $KCS_SEL
396     clrq   %rax
397     movw   %ss, %ax
398     pushq  %rax          /* rip should be here */
399     pushq  $-1          /* phony trap error code */
400     pushq  $-1          /* phony trap number */

402     subq   $REG_OFF(KDIREG_TRAPNO), %rsp
403     KDI_SAVE_REGS(%rsp)

405     movq   %cr3, %rax
406     movq   %rax, REG_OFF(KDIREG_CR3)(%rsp)

408     movq   REG_OFF(KDIREG_SS)(%rsp), %rax
409     xchgq  REG_OFF(KDIREG_RIP)(%rsp), %rax
410     movq   %rax, REG_OFF(KDIREG_SS)(%rsp)

412     movq   REG_OFF(KDIREG_RSP)(%rsp), %rax
413     addq   $8, %rax
414     movq   %rax, REG_OFF(KDIREG_RSP)(%rsp)

416     /*
417     * We've saved all of the general-purpose registers, and have a stack
418     * that is irettable (after we strip down to the error code)
419     */

421     GET_CPUSAVE_ADDR /* %rax = cpusave, %rbx = CPU ID */

423     ADVANCE_CRUMB_POINTER(%rax, %rcx, %rdx)

425     ADD_CRUMB(%rax, KRM_CPU_STATE, $KDI_CPU_STATE_SLAVE, %rdx)

427     movq   REG_OFF(KDIREG_RIP)(%rsp), %rcx
428     ADD_CRUMB(%rax, KRM_PC, %rcx, %rdx)

430     pushq  %rax

```

```

431      jmp      kdi_save_common_state
433      SET_SIZE(kdi_slave_entry)

468 #endif /* __lint */

435 /*
436 * The state of the world:
437 *
438 * The stack has a complete set of saved registers and segment
439 * selectors, arranged in the kdi_regs.h order. It also has a pointer
440 * to our cpusave area.
441 *
442 * We need to save, into the cpusave area, a pointer to these saved
443 * registers. First we check whether we should jump straight back to
444 * the kernel. If not, we save a few more registers, ready the
445 * machine for debugger entry, and enter the debugger.
446 */

483 #if !defined(__lint)

448      ENTRY_NP(kdi_save_common_state)

450      popq    %rdi          /* the cpusave area */
451      movq    %rsp, KRS_GREGS(%rdi) /* save ptr to current saved regs */

453      pushq   %rdi
454      call    kdi_trap_pass
455      cmpq    $1, %rax
456      je     kdi_pass_to_kernel
457      popq    %rax /* cpusave in %rax */

459      SAVE_IDTGDGT

461 #if !defined(__xpv)
462      /* Save off %cr0, and clear write protect */
463      movq    %cr0, %rcx
464      movq    %rcx, KRS_CR0(%rax)
465      andq    $_BITNOT(CR0_WP), %rcx
466      movq    %rcx, %cr0
467 #endif

469      /* Save the debug registers and disable any active watchpoints */

471      movq    %rax, %r15          /* save cpusave area ptr */
472      movl    $7, %edi
473      call    kdi_dreg_get
474      movq    %rax, KRS_DRCTL(%r15)

476      andq    $_BITNOT(KDIREG_DRCTL_WPALLEN_MASK), %rax
477      movq    %rax, %rsi
478      movl    $7, %edi
479      call    kdi_dreg_set

481      movl    $6, %edi
482      call    kdi_dreg_get
483      movq    %rax, KRS_DRSTAT(%r15)

485      movl    $0, %edi
486      call    kdi_dreg_get
487      movq    %rax, KRS_DROFF(0)(%r15)

489      movl    $1, %edi
490      call    kdi_dreg_get
491      movq    %rax, KRS_DROFF(1)(%r15)

```

```

493      movl    $2, %edi
494      call    kdi_dreg_get
495      movq    %rax, KRS_DROFF(2)(%r15)

497      movl    $3, %edi
498      call    kdi_dreg_get
499      movq    %rax, KRS_DROFF(3)(%r15)

501      movq    %r15, %rax          /* restore cpu save area to rax */

540      /*
541      * Save any requested MSRs.
542      */
543      movq    KRS_MSR(%rax), %rcx
544      cmpq    $0, %rcx
545      je     no_msr

547      pushq   %rax          /* rdmsr clobbers %eax */
548      movq    %rcx, %rbx

550 1:
551      movl    MSR_NUM(%rbx), %ecx
552      cmpl    $0, %ecx
553      je     msr_done

555      movl    MSR_TYPE(%rbx), %edx
556      cmpl    $KDI_MSR_READ, %edx
557      jne    msr_next

559      rdmsr          /* addr in %ecx, value into %edx:%eax */
560      movl    %eax, MSR_VAL(%rbx)
561      movl    %edx, _CONST(MSR_VAL + 4)(%rbx)

563 msr_next:
564      addq    $MSR_SIZE, %rbx
565      jmp     1b

567 msr_done:
568      popq    %rax

570 no_msr:
503      clrq    %rbp          /* stack traces should end here */

505      pushq   %rax
506      movq    %rax, %rdi          /* cpusave */

508      call    kdi_debugger_entry

510      /* Pass cpusave to kdi_resume */
511      popq    %rdi

513      jmp     kdi_resume

515      SET_SIZE(kdi_save_common_state)

585 #endif /* !__lint */

517 /*
518 * Resume the world. The code that calls kdi_resume has already
519 * decided whether or not to restore the IDT.
520 */
591 #if defined(__lint)
592 void
593 kdi_resume(void)
594 {
595 }

```



```

596 #else /* __lint */
521 /* cpusave in %rdi */
522 ENTRY_NP(kdi_resume)
524 /*
525  * Send this CPU back into the world
526  */
527 #if !defined(__xpv)
528 movq   KRS_CR0(%rdi), %rdx
529 movq   %rdx, %cr0
530 #endif
532 KDI_RESTORE_DEBUGGING_STATE
534 movq   KRS_GREGS(%rdi), %rsp
536 #if !defined(__xpv)
537 /*
538  * If we're going back via tr_iret_kdi, then we want to copy the
539  * final %cr3 we're going to back into the kpti_dbg area now.
540  *
541  * Since the trampoline needs to find the kpti_dbg too, we enter it
542  * with %r13 set to point at that. The real %r13 (to restore before
543  * the iret) we stash in the kpti_dbg itself.
544  */
545 movq   %gs:CPU_SELF, %r13 /* can't leaq %gs:*, use self-ptr */
546 addq   $CPU_KPTI_DBG, %r13
548 movq   REG_OFF(KDIREG_R13)(%rsp), %rdx
549 movq   %rdx, KPTI_R13(%r13)
551 movq   REG_OFF(KDIREG_CR3)(%rsp), %rdx
552 movq   %rdx, KPTI_TR_CR3(%r13)
554 /* The trampoline will undo this later. */
555 movq   %r13, REG_OFF(KDIREG_R13)(%rsp)
556 #endif
558 KDI_RESTORE_REGS(%rsp)
559 addq   $REG_OFF(KDIREG_RIP), %rsp /* Discard state, trapno, err */
560 /*
561  * The common trampoline code will restore %cr3 to the right value
562  * for either kernel or userland.
563  */
564 #if !defined(__xpv)
565 jmp    tr_iret_kdi
566 #else
567 IRET
568 #endif
569 /*NOTREACHED*/
570 SET_SIZE(kdi_resume)
618 #endif /* __lint */
620 #if !defined(__lint)
572 ENTRY_NP(kdi_pass_to_kernel)
574 popq   %rdi /* cpusave */
576 movq   $KDI_CPU_STATE_NONE, KRS_CPU_STATE(%rdi)
578 /*
579  * Find the trap and vector off the right kernel handler. The trap
580  * handler will expect the stack to be in trap order, with %rip being

```

```

581  * the last entry, so we'll need to restore all our regs. On i86xpv
582  * we'll need to compensate for XPV_TRAP_POP.
583  *
584  * We're hard-coding the three cases where KMDB has installed permanent
585  * handlers, since after we KDI_RESTORE_REGS(), we don't have registers
586  * to work with; we can't use a global since other CPUs can easily pass
587  * through here at the same time.
588  *
589  * Note that we handle T_DBGENTR since userspace might have tried it.
590  */
591 movq   KRS_GREGS(%rdi), %rsp
592 movq   REG_OFF(KDIREG_TRAPNO)(%rsp), %rdi
593 cmpq   $T_SGLSTP, %rdi
594 je     1f
595 cmpq   $T_BPTFLT, %rdi
596 je     2f
597 cmpq   $T_DBGENTR, %rdi
598 je     3f
599 /*
600  * Hmm, unknown handler. Somebody forgot to update this when they
601  * added a new trap interposition... try to drop back into kmdb.
602  */
603 int    $T_DBGENTR
605 #define CALL_TRAP_HANDLER(name) \
606 KDI_RESTORE_REGS(%rsp); \
607 /* Discard state, trapno, err */ \
608 addq   $REG_OFF(KDIREG_RIP), %rsp; \
609 XPV_TRAP_PUSH; \
610 jmp    %cs:name
612 1:
613 CALL_TRAP_HANDLER(dbgtrap)
614 /*NOTREACHED*/
615 2:
616 CALL_TRAP_HANDLER(brktrap)
617 /*NOTREACHED*/
618 3:
619 CALL_TRAP_HANDLER(inaltrap)
620 /*NOTREACHED*/
622 SET_SIZE(kdi_pass_to_kernel)
unchanged portion omitted
692 #endif /* !__lint */
694 #if defined(__lint)
695 /*ARGSUSED*/
696 void
697 kdi_cpu_debug_init(kdi_cpusave_t *save)
698 {
699 }
700 #else /* __lint */
642 ENTRY_NP(kdi_cpu_debug_init)
643 pushq  %rbp
644 movq   %rsp, %rbp
646 pushq  %rbx /* macro will clobber %rbx */
647 KDI_RESTORE_DEBUGGING_STATE
648 popq   %rbx
650 leave
651 ret
652 SET_SIZE(kdi_cpu_debug_init)

```

```
654 #define GETDREG(name, r) \
655     ENTRY_NP(name); \
656     movq    r, %rax; \
657     ret; \
658     SET_SIZE(name)

660 #define SETDREG(name, r) \
661     ENTRY_NP(name); \
662     movq    %rdi, r; \
663     ret; \
664     SET_SIZE(name)

666     GETDREG(kdi_getdr0, %dr0)
667     GETDREG(kdi_getdr1, %dr1)
668     GETDREG(kdi_getdr2, %dr2)
669     GETDREG(kdi_getdr3, %dr3)
670     GETDREG(kdi_getdr6, %dr6)
671     GETDREG(kdi_getdr7, %dr7)

673     SETDREG(kdi_setdr0, %dr0)
674     SETDREG(kdi_setdr1, %dr1)
675     SETDREG(kdi_setdr2, %dr2)
676     SETDREG(kdi_setdr3, %dr3)
677     SETDREG(kdi_setdr6, %dr6)
678     SETDREG(kdi_setdr7, %dr7)

713     SET_SIZE(kdi_cpu_debug_init)
680 #endif /* !__lint */
```

new/usr/src/uts/intel/kdi/kdi_idt.c

1

```
*****
12124 Fri Apr 6 17:25:12 2018
new/usr/src/uts/intel/kdi/kdi_idt.c
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>
9210 remove KMDB branch debugging support
9211 ::crregs could do with cr2/cr3 support
9209 ::ttrace should be able to filter by thread
Reviewed by: Patrick Mooney <patrick.mooney@joyent.com>
Reviewed by: Yuri Pankov <yuripv@yuripv.net>
9207 kdi_idt: Cast GATESEG_GETOFFSET through uintptr_t
Reviewed by: Yuri Pankov <yuripv@yuripv.net>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22  * Copyright 2009 Sun Microsystems, Inc. All rights reserved.
23  * Use is subject to license terms.
24  *
25  * Copyright 2018 Joyent, Inc.
26  */

28 /*
29  * Management of KMDB's IDT, which is installed upon KMDB activation.
30  *
31  * Debugger activation has two flavors, which cover the cases where KMDB is
32  * loaded at boot, and when it is loaded after boot. In brief, in both cases,
33  * the KDI needs to interpose upon several handlers in the IDT. When
34  * mod-loaded KMDB is deactivated, we undo the IDT interposition, restoring the
35  * handlers to what they were before we started.
36  *
37  * We also take over the entirety of IDT (except the double-fault handler) on
38  * the active CPU when we're in kmdb so we can handle things like page faults
39  * sensibly.
40  *
41  * Boot-loaded KMDB
42  *
43  * When we're first activated, we're running on boot's IDT. We need to be able
44  * to function in this world, so we'll install our handlers into boot's IDT.
45  * This is a little complicated: we're using the fake cpu_t set up by
46  * boot_kdi_tmpinit(), so we can't access cpu_idt directly. Instead,
47  * kdi_idt_write() notices that cpu_idt is NULL, and works around this problem.
48  *
49  * Later, when we're about to switch to the kernel's IDT, it'll call us via
50  * kdi_idt_sync(), allowing us to add our handlers to the new IDT. While
51  * boot-loaded KMDB can't be unloaded, we still need to save the descriptors we
52  * replace so we can pass traps back to the kernel as necessary.
```

new/usr/src/uts/intel/kdi/kdi_idt.c

2

```
53 *
54 * The last phase of boot-loaded KMDB activation occurs at non-boot CPU
55 * startup. We will be called on each non-boot CPU, thus allowing us to set up
56 * any watchpoints that may have been configured on the boot CPU and interpose
57 * on the given CPU's IDT. We don't save the interposed descriptors in this
58 * case -- see kdi_cpu_init() for details.
59 *
60 * Mod-loaded KMDB
61 *
62 * This style of activation is much simpler, as the CPUs are already running,
63 * and are using their own copy of the kernel's IDT. We simply interpose upon
64 * each CPU's IDT. We save the handlers we replace, both for deactivation and
65 * for passing traps back to the kernel. Note that for the hypervisors'
66 * benefit, we need to xcall to the other CPUs to do this, since we need to
67 * actively set the trap entries in its virtual IDT from that vcpu's context
68 * rather than just modifying the IDT table from the CPU running kdi_activate().
69 */

71 #include <sys/types.h>
72 #include <sys/segments.h>
73 #include <sys/trap.h>
74 #include <sys/cpuvar.h>
75 #include <sys/reboot.h>
76 #include <sys/sunddi.h>
77 #include <sys/archsystem.h>
78 #include <sys/kdi_impl.h>
79 #include <sys/x_call.h>
80 #include <ia32/sys/psw.h>
81 #include <vm/hat_i86.h>

83 #define KDI_GATE_NVECS 3

85 #define KDI_IDT_NOSAVE 0
86 #define KDI_IDT_SAVE 1

88 #define KDI_IDT_DTYPE_KERNEL 0
89 #define KDI_IDT_DTYPE_BOOT 1

91 kdi_cpusave_t *kdi_cpusave;
92 int kdi_ncpusave;

94 static kdi_main_t kdi_kmdb_main;

96 kdi_drreg_t kdi_drreg;

98 #ifndef __amd64
99 /* Used to track the current set of valid kernel selectors. */
100 uint32_t kdi_cs;
101 uint32_t kdi_ds;
102 uint32_t kdi_fs;
103 uint32_t kdi_gs;
104 #endif

103 uint_t kdi_msr_wrexite_msr;
104 uint64_t *kdi_msr_wrexite_valp;

106 uintptr_t kdi_kernel_handler;

108 int kdi_trap_switch;

110 #define KDI_MEMRANGES_MAX 2

112 kdi_memrange_t kdi_memranges[KDI_MEMRANGES_MAX];
113 int kdi_nmemranges;

115 typedef void idt_hdlr_f(void);
```

```

117 extern idt_hdr_f kdi_trap0, kdi_trap1, kdi_int2, kdi_trap3, kdi_trap4;
118 extern idt_hdr_f kdi_trap5, kdi_trap6, kdi_trap7, kdi_trap9;
119 extern idt_hdr_f kdi_traperr10, kdi_traperr11, kdi_traperr12;
120 extern idt_hdr_f kdi_traperr13, kdi_traperr14, kdi_trap16, kdi_traperr17;
120 extern idt_hdr_f kdi_traperr13, kdi_traperr14, kdi_trap16, kdi_trap17;
121 extern idt_hdr_f kdi_trap18, kdi_trap19, kdi_trap20, kdi_ivct32;
122 extern idt_hdr_f kdi_invaltrap;
123 extern size_t kdi_ivct_size;
124 extern char kdi_slave_entry_patch;

125 typedef struct kdi_gate_spec {
126     uint_t kgs_vec;
127     uint_t kgs_dpl;
128 } kdi_gate_spec_t;
_____
139 static gate_desc_t kdi_kgates[KDI_GATE_NVECS];

141 extern gate_desc_t kdi_idt[NIDT];
142 gate_desc_t kdi_idt[NIDT];

143 struct idt_description {
144     uint_t id_low;
145     uint_t id_high;
146     idt_hdr_f *id_basehldr;
147     size_t *id_incrp;
148 } idt_description[] = {
149     { T_ZERODIV, 0, kdi_trap0, NULL },
150     { T_SGLSTP, 0, kdi_trap1, NULL },
151     { T_NMIFLT, 0, kdi_int2, NULL },
152     { T_BPTFLT, 0, kdi_trap3, NULL },
153     { T_OVFLW, 0, kdi_trap4, NULL },
154     { T_BOUNDFLT, 0, kdi_trap5, NULL },
155     { T_ILLINST, 0, kdi_trap6, NULL },
156     { T_NOEXTFLT, 0, kdi_trap7, NULL },
157 #if !defined(__xpv)
158     { T_DBLFLT, 0, syserrtrap, NULL },
159 #endif
160     { T_EXTTOVRFLT, 0, kdi_trap9, NULL },
161     { T_TSSFLT, 0, kdi_traperr10, NULL },
162     { T_SEGFLT, 0, kdi_traperr11, NULL },
163     { T_STKFLT, 0, kdi_traperr12, NULL },
164     { T_GPFLT, 0, kdi_traperr13, NULL },
165     { T_PGFLT, 0, kdi_traperr14, NULL },
166     { 15, 0, kdi_invaltrap, NULL },
167     { T_EXTERRFLT, 0, kdi_trap16, NULL },
168     { T_ALIGNMENT, 0, kdi_traperr17, NULL },
169     { T_ALIGNMENT, 0, kdi_trap17, NULL },
170     { T_MCE, 0, kdi_trap18, NULL },
171     { T_SIMDFPE, 0, kdi_trap19, NULL },
172     { T_DBGENTR, 0, kdi_trap20, NULL },
173     { 21, 31, kdi_invaltrap, NULL },
174     { 32, 255, kdi_ivct32, &kdi_ivct_size },
175 };

177 void
178 kdi_idt_init(selector_t sel)
179 {
180     struct idt_description *id;
181     int i;

183     for (id = idt_description; id->id_basehldr != NULL; id++) {
184         uint_t high = id->id_high != 0 ? id->id_high : id->id_low;
185         size_t incr = id->id_incrp != NULL ? *id->id_incrp : 0;

```

```

187 #if !defined(__xpv)
188     if (kpti_enable && sel == KCS_SEL && id->id_low == T_DBLFLT)
189         id->id_basehldr = tr_syserrtrap;
190 #endif

192     for (i = id->id_low; i <= high; i++) {
193         caddr_t hldr = (caddr_t)id->id_basehldr +
194             incr * (i - id->id_low);
195         set_gatesegd(&kdi_idt[i], (void (*)())hldr, sel,
196             SDT_SIGSIGT, TRP_KPL, IST_DBG);
197         SDT_SIGSIGT, TRP_KPL, i);
198     }
199 }

197 /*
198  * Patch caller-provided code into the debugger's IDT handlers. This code is
199  * used to save MSRs that must be saved before the first branch. All handlers
200  * are essentially the same, and end with a branch to kdi_cmhint. To save the
201  * MSR, we need to patch in before the branch. The handlers have the following
202  * structure: KDI_MSR_PATCHOFF bytes of code, KDI_MSR_PATCHSZ bytes of
203  * patchable space, followed by more code.
204  */
205 void
206 kdi_idt_patch(caddr_t code, size_t sz)
207 {
208     int i;

210     ASSERT(sz <= KDI_MSR_PATCHSZ);

212     for (i = 0; i < sizeof(kdi_idt) / sizeof(struct gate_desc); i++) {
213         gate_desc_t *gd;
214         uchar_t *patch;

216         if (i == T_DBLFLT)
217             continue; /* uses kernel's handler */

219         gd = &kdi_idt[i];
220         patch = (uchar_t *)GATESEG_GETOFFSET(gd) + KDI_MSR_PATCHOFF;

222         /*
223          * We can't ASSERT that there's a nop here, because this may be
224          * a debugger restart. In that case, we're copying the new
225          * patch point over the old one.
226          */
227         /* FIXME: dtrace fbt ... */
228         bcopy(code, patch, sz);

230         /* Fill the rest with nops to be sure */
231         while (sz < KDI_MSR_PATCHSZ)
232             patch[sz++] = 0x90; /* nop */
233     }
234 }

201 static void
202 kdi_idt_gates_install(selector_t sel, int saveold)
203 {
204     gate_desc_t gates[KDI_GATE_NVECS];
205     int i;

207     bzero(gates, sizeof(*gates));

209     for (i = 0; i < KDI_GATE_NVECS; i++) {
210         const kdi_gate_spec_t *gs = &kdi_gate_specs[i];
211         uintptr_t func = GATESEG_GETOFFSET(&kdi_idt[gs->kgs_vec]);

```

```

212         set_gatesegd(&gates[i], (void (*)())func, sel, SDT_SYSIGT,
213                    gs->kgs_dpl, IST_DBG);
214         gs->kgs_dpl, gs->kgs_vec);
215     }
216     for (i = 0; i < KDI_GATE_NVECS; i++) {
217         uint_t vec = kdi_gate_specs[i].kgs_vec;
218
219         if (saveold)
220             kdi_kgates[i] = CPU->cpu_m.mcpu_idt[vec];
221
222         kdi_idt_write(&gates[i], vec);
223     }
224 }
225
226 unchanged_portion_omitted
227
228 /*
229 * On some processors, we'll need to clear a certain MSR before proceeding into
230 * the debugger. Complicating matters, this MSR must be cleared before we take
231 * any branches. We have patch points in every trap handler, which will cover
232 * all entry paths for master CPUs. We also have a patch point in the slave
233 * entry code.
234 */
235 static void
236 kdi_msr_add_clreentry(uint_t msr)
237 {
238     #ifdef __amd64
239         uchar_t code[] = {
240             0x51, 0x50, 0x52,          /* pushq %rcx, %rax, %rdx */
241             0xb9, 0x00, 0x00, 0x00, 0x00, /* movl $MSRNUM, %ecx */
242             0x31, 0xc0,              /* clr %eax */
243             0x31, 0xd2,              /* clr %edx */
244             0x0f, 0x30,              /* wrmsr */
245             0x5a, 0x58, 0x59         /* popq %rdx, %rax, %rcx */
246         };
247         uchar_t *patch = &code[4];
248     #else
249         uchar_t code[] = {
250             0x60,                    /* pushal */
251             0xb9, 0x00, 0x00, 0x00, 0x00, /* movl $MSRNUM, %ecx */
252             0x31, 0xc0,              /* clr %eax */
253             0x31, 0xd2,              /* clr %edx */
254             0x0f, 0x30,              /* wrmsr */
255             0x61                     /* popal */
256         };
257         uchar_t *patch = &code[2];
258     #endif
259
260     bcopy(&msr, patch, sizeof (uint32_t));
261
262     kdi_idt_patch((caddr_t)code, sizeof (code));
263
264     bcopy(code, &kdi_slave_entry_patch, sizeof (code));
265 }
266
267 static void
268 kdi_msr_add_wrexist(uint_t msr, uint64_t *valp)
269 {
270     kdi_msr_wrexist_msr = msr;
271     kdi_msr_wrexist_valp = valp;
272 }
273
274 void
275 kdi_set_debug_msrs(kdi_msr_t *msrs)
276 {
277     int nmsrs, i;

```

```

322     ASSERT(kdi_cpusave[0].krs_msr == NULL);
323
324     /* Look in CPU0's MSRs for any special MSRs. */
325     for (nmsrs = 0; msrs[nmsrs].msr_num != 0; nmsrs++) {
326         switch (msrs[nmsrs].msr_type) {
327             case KDI_MSR_CLEARENTRY:
328                 kdi_msr_add_clreentry(msrs[nmsrs].msr_num);
329                 break;
330
331             case KDI_MSR_WRITEDELAY:
332                 kdi_msr_add_wrexist(msrs[nmsrs].msr_num,
333                                     msrs[nmsrs].kdi_msr_valp);
334                 break;
335         }
336     }
337     nmsrs++;
338
339     for (i = 0; i < kdi_ncpusave; i++)
340         kdi_cpusave[i].krs_msr = &msrs[nmsrs * i];
341 }
342
343 void
344 kdi_update_drreg(kdi_drreg_t *drreg)
345 {
346     kdi_drreg = *drreg;
347 }
348
349 unchanged_portion_omitted
350
351 /*
352 * Activation for CPUs other than the boot CPU, called from that CPU's
353 * mp_startup(). We saved the kernel's descriptors when we initialized the
354 * boot CPU, so we don't want to do it again. Saving the handlers from this
355 * CPU's IDT would actually be dangerous with the CPU initialization method in
356 * use at the time of this writing. With that method, the startup code creates
357 * the IDTs for slave CPUs by copying the one used by the boot CPU, which has
358 * already been interposed upon by KMDB. Were we to interpose again, we'd
359 * replace the kernel's descriptors with our own in the save area. By not
360 * saving, but still overwriting, we'll work in the current world, and in any
361 * future world where the IDT is generated from scratch.
362 */
363 void
364 kdi_cpu_init(void)
365 {
366     kdi_idt_gates_install(KCS_SEL, KDI_IDT_NOSAVE);
367     /* Load the debug registers. */
368     /* Load the debug registers and MSRs */
369     kdi_cpu_debug_init(&kdi_cpusave[CPU->cpu_id]);
370 }
371
372 unchanged_portion_omitted
373
374 void
375 kdi_activate(kdi_main_t main, kdi_cpusave_t *cpusave, uint_t ncpusave)
376 {
377     int i;
378     cpuset_t cpuset;
379
380     CPuset_ALL(cpuset);
381
382     kdi_cpusave = cpusave;
383     kdi_ncpusave = ncpusave;
384
385     kdi_kmdb_main = main;
386
387     for (i = 0; i < kdi_ncpusave; i++) {

```

```

318         kdi_cpusave[i].krs_cpu_id = i;

320         kdi_cpusave[i].krs_curcrumb =
321             &kdi_cpusave[i].krs_crumbs[KDI_NCRUMBS - 1];
322         kdi_cpusave[i].krs_curcrumbidx = KDI_NCRUMBS - 1;
323     }

325     if (boothowto & RB_KMDB)
326         kdi_idt_init(KMDBCODE_SEL);
327     else
328         kdi_idt_init(KCS_SEL);

330     /* The initial selector set. Updated by the debugger-entry code */
331 #ifndef __amd64
332     kdi_cs = B32CODE_SEL;
333     kdi_ds = kdi_fs = kdi_gs = B32DATA_SEL;
334 #endif

336     kdi_memranges[0].mr_base = kdi_segdebugbase;
337     kdi_memranges[0].mr_lim = kdi_segdebugbase + kdi_segdebugsize - 1;
338     kdi_nmemranges = 1;

340     kdi_drreg.dr_ctl = KDIREG_DRCTL_RESERVED;
341     kdi_drreg.dr_stat = KDIREG_DRSTAT_RESERVED;

451     kdi_msr_wrexite_msr = 0;
452     kdi_msr_wrexite_valp = NULL;

343     if (boothowto & RB_KMDB) {
344         kdi_idt_gates_install(KMDBCODE_SEL, KDI_IDT_NOSAVE);
345     } else {
346         xc_call(0, 0, 0, CPUSET2BV(cpuset),
347             (xc_func_t)kdi_cpu_activate);
348     }
349 }

unchanged_portion_omitted_

368 /*
369 * We receive all breakpoints and single step traps. Some of them,
370 * including those from userland and those induced by DTrace providers,
371 * are intended for the kernel, and must be processed there. We adopt
372 * this ours-until-proven-otherwise position due to the painful
373 * consequences of sending the kernel an unexpected breakpoint or
374 * single step. Unless someone can prove to us that the kernel is
375 * prepared to handle the trap, we'll assume there's a problem and will
376 * give the user a chance to debug it.
377 */
378 int
379 kdi_trap_pass(kdi_cpusave_t *cpusave)
380 {
381     greg_t tt = cpusave->krs_gregs[KDIREG_TRAPNO];
382     greg_t pc = cpusave->krs_gregs[KDIREG_PC];
383     greg_t cs = cpusave->krs_gregs[KDIREG_CS];

385     if (USERMODE(cs))
386         return (1);

388     if (tt != T_BPTFLT && tt != T_SGLSTP)
389         return (0);

391     if (tt == T_BPTFLT && kdi_dtrace_get_state() ==
392         KDI_DTSTATE_DTRACE_ACTIVE)
393         return (1);

395     /*
396     * See the comments in the kernel's T_SGLSTP handler for why we need to

```

```

397         * do this.
398         */
399 #if !defined(__xpv)
400     if (tt == T_SGLSTP &&
401         (pc == (greg_t)sys_sysenter || pc == (greg_t)brand_sys_sysenter ||
402          pc == (greg_t)tr_sys_sysenter ||
403          pc == (greg_t)tr_brand_sys_sysenter)) {
404 #else
405     if (tt == T_SGLSTP &&
406         (pc == (greg_t)sys_sysenter || pc == (greg_t)brand_sys_sysenter)) {
407 #endif
408         (pc == (greg_t)sys_sysenter || pc == (greg_t)brand_sys_sysenter))
409         return (1);

411     return (0);
412 }

unchanged_portion_omitted_

```

new/usr/src/uts/intel/kdi/kdi_idthdl.s

1

8885 Fri Apr 6 17:25:13 2018

new/usr/src/uts/intel/kdi/kdi_idthdl.s

8956 Implement KPTI

Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>

Reviewed by: Robert Mustacchi <rm@joyent.com>

9210 remove KMDB branch debugging support

9211 ::crregs could do with cr2/cr3 support

9209 ::ttrace should be able to filter by thread

Reviewed by: Patrick Mooney <patrick.mooney@joyent.com>

Reviewed by: Yuri Pankov <yuripv@yuripv.net>

```
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22 * Copyright 2007 Sun Microsystems, Inc. All rights reserved.
23 * Use is subject to license terms.
24 *
25 * Copyright 2018 Joyent, Inc.
26 */

26 #pragma ident      "%Z%M% %I%      %E% SMI"

28 /*
29  * Companion to kdi_asm.s - the implementation of the trap and interrupt
30  * handlers. For the most part, these handlers do the same thing - they
31  * push a trap number onto the stack, followed by a jump to kdi_cmint.
32  * Each trap and interrupt has its own handler because each one pushes a
33  * different number.
34  */

36 #if defined(__lint)
37 #include <sys/types.h>
38 #else

40 #include <sys/asm_linkage.h>
41 #include <sys/asm_misc.h>
42 #include <sys/machprivregs.h>
43 #include <sys/privregs.h>
44 #include <sys/kdi_regs.h>
45 #include <sys/trap.h>
46 #include <sys/param.h>

48 #include <kdi_assym.h>
49 #include <assym.h>
39 /* Nothing in this file is of interest to lint. */
40 #if !defined(__lint)
```

new/usr/src/uts/intel/kdi/kdi_idthdl.s

2

```
51 /*
52  * The default ASM_ENTRY_ALIGN (16) wastes far too much space.
53  * The default ASM_ENTRY_ALIGN (16) wastes far too much space. Pay no
54  * attention to the fleet of nop's we're adding to each handler.
55 */
56 #undef ASM_ENTRY_ALIGN
57 #define ASM_ENTRY_ALIGN 8

58 /*
59  * Generic trap and interrupt handlers.
60  * We need the .align in ENTRY_NP (defined to be ASM_ENTRY_ALIGN) to match our
61  * manual .align (KDI_MSR_PATCHOFF) in order to ensure that the space reserved
62  * at the beginning of the handler for code is exactly KDI_MSR_PATCHOFF bytes
63  * long. Note that the #error below isn't supported by the preprocessor invoked
64  * by as(1), and won't stop the build, but it'll emit a noticeable error message
65  * which won't escape the filters.
66 */
67 #if ASM_ENTRY_ALIGN != KDI_MSR_PATCHOFF
68 #error "ASM_ENTRY_ALIGN != KDI_MSR_PATCHOFF"
69 this won't assemble
70 #endif

71 #if defined(__xpv)
72 /*
73  * kdi_idt_patch will, on certain processors, replace the patch points below
74  * with MSR-clearing code. kdi_id_patch has intimate knowledge of the size of
75  * the nop hole, as well as the structure of the handlers. Do not change
76  * anything here without also changing kdi_idt_patch.
77 */

78 #define INTERRUPT_TRAMPOLINE
79 /*
80  * Generic trap and interrupt handlers.
81 */

82 #else
83 #if defined(__xpv) && defined(__amd64)

84 /*
85  * If we're !xpv, then we will need to support KPTI (kernel page table
86  * isolation), where we have separate page tables for user and kernel modes.
87  * There's more detail about this in kpti_trampoline.s and hat_i86.c
88  * The hypervisor places r11 and rcx on the stack.
89 */

90 #define INTERRUPT_TRAMPOLINE \
91     pushq   %r13; \
92     pushq   %r14; \
93     subq    $KPTI_R14, %rsp; \
94     /* Check for clobbering */ \
95     cmp     $0, KPTI_FLAG(%rsp); \
96     je      1f; \
97     /* Don't worry, this totally works */ \
98     int     $8; \
99     1: \
100    movq    $1, KPTI_FLAG(%rsp); \
101    /* Save current %cr3. */ \
102    mov     %cr3, %r14; \
103    mov     %r14, KPTI_TR_CR3(%rsp); \
104    /* Switch to paranoid %cr3. */ \
105    mov     kpti_safe_cr3, %r14; \
106    mov     %r14, %cr3; \
107    \
108    cmpw    $KCS_SEL, KPTI_CS(%rsp); \
109    je      3f; \
```

```

93 2:
94     /* Get our cpu_t in %r13 */
95     mov     %rsp, %r13;
96     and     $(~(MMU_PAGESIZE - 1)), %r13;
97     subq   $CPU_KPTI_START, %r13;
98     /* Use top of the kthread stk */
99     mov     CPU_THREAD(%r13), %r14;
100    mov     T_STACK(%r14), %r14;
101    addq   $REGSIZE+MINFRAME, %r14;
102    jmp     5f;
103 3:
104    /* Check the %rsp in the frame. */
105    /* Is it above kernel base? */
106    mov     kpti_kbase, %r14;
107    cmp     %r14, KPTI_RSP(%rsp);
108    jb     2b;
109    /* Is it within the kpti_frame page? */
110    mov     %rsp, %r13;
111    and     $(~(MMU_PAGESIZE - 1)), %r13;
112    mov     KPTI_RSP(%rsp), %r14;
113    and     $(~(MMU_PAGESIZE - 1)), %r14;
114    cmp     %r13, %r14;
115    je     2b;
116    /* Use the %rsp from the trap frame. */
117    /* We already did %cr3. */
118    mov     KPTI_RSP(%rsp), %r14;
119    and     $(~0xf), %r14;
120 5:
121    mov     %rsp, %r13;
122    /* %r14 contains our destination stk */
123    mov     %r14, %rsp;
124    pushq   KPTI_SS(%r13);
125    pushq   KPTI_RSP(%r13);
126    pushq   KPTI_RFLAGS(%r13);
127    pushq   KPTI_CS(%r13);
128    pushq   KPTI_RIP(%r13);
129    pushq   KPTI_ERR(%r13);
130    mov     KPTI_R14(%r13), %r14;
131    movq   $0, KPTI_FLAG(%r13);
132    mov     KPTI_R13(%r13), %r13
79 #define TRAP_NOERR(trapno) \
80     popq   %rcx;
81     popq   %r11;
82     pushq  $trapno
134 #endif /* !_xpv */
84 #define TRAP_ERR(trapno) \
85     popq   %rcx;
86     popq   %r11;
87     pushq  $0;
88     pushq  $trapno
90 #else
92 #define TRAP_NOERR(trapno) \
93     push  $trapno
95 #define TRAP_ERR(trapno) \
96     push  $0;
97     push  $trapno
99 #endif /* __xpv && __amd64 */
137 #define MKIVCT(n) \
138     ENTRY_NP(kdi_ivct/**/n/**/);

```

```

139     XPV_TRAP_POP;
140     push  $0; /* err */
141     INTERRUPT_TRAMPOLINE;
142     push  $n;
143     TRAP_ERR(n);
144     .align KDI_MSR_PATCHOFF;
145     KDI_MSR_PATCH;
146     jmp   kdi_cmint;
147     SET_SIZE(kdi_ivct/**/n/**/)
148 #define MKTRAPHDLR(n) \
149     ENTRY_NP(kdi_trap/**/n);
150     XPV_TRAP_POP;
151     push  $0; /* err */
152     INTERRUPT_TRAMPOLINE;
153     push  $n;
154     TRAP_ERR(n);
155     .align KDI_MSR_PATCHOFF;
156     KDI_MSR_PATCH;
157     jmp   kdi_cmint;
158     SET_SIZE(kdi_trap/**/n/**/)
159 #define MKTRAPERRHDLR(n) \
160     ENTRY_NP(kdi_traperr/**/n);
161     XPV_TRAP_POP;
162     INTERRUPT_TRAMPOLINE;
163     push  $n;
164     TRAP_NOERR(n);
165     .align KDI_MSR_PATCHOFF;
166     KDI_MSR_PATCH;
167     jmp   kdi_cmint;
168     SET_SIZE(kdi_traperr/**/n)
169 #if !defined(__xpv)
170 #define MKNMIHDLR \
171     ENTRY_NP(kdi_int2);
172     push  $0;
173     push  $2;
174     pushq %r13;
175     mov   kpti_safe_cr3, %r13;
176     mov   %r13, %cr3;
177     popq  %r13;
178     TRAP_NOERR(2);
179     .align KDI_MSR_PATCHOFF;
180     KDI_MSR_PATCH;
181     jmp   kdi_nmiint;
182     SET_SIZE(kdi_int2)
183 #define MKMCEHDLR \
184     ENTRY_NP(kdi_trap18);
185     push  $0;
186     push  $18;
187     pushq %r13;
188     mov   kpti_safe_cr3, %r13;
189     mov   %r13, %cr3;
190     popq  %r13;
191     jmp   kdi_cmint;
192     SET_SIZE(kdi_trap18)
193 #else
194 #define MKNMIHDLR \
195     ENTRY_NP(kdi_int2);
196     push  $0;
197     push  $2;
198     jmp   kdi_nmiint;
199     SET_SIZE(kdi_int2)

```



```

193 #define MKMCEHDLR \
194     ENTRY_NP(kdi_trap18);          \
195     push    $0;                    \
196     push    $18;                   \
197     jmp     kdi_cmnint;            \
198     SET_SIZE(kdi_trap18)
199 #endif

201 /*
202 * The only way we should reach here is by an explicit "int 0x.." which is
203 * defined not to push an error code.
204 */
205 #define MKINVALHDLR \
206     ENTRY_NP(kdi_invaltrap);      \
207     XPV_TRAP_POP;                 \
208     push    $0; /* err */         \
209     INTERRUPT_TRAMPOLINE;         \
210     push    $255;                  \
211     TRAP_NOERR(255);              \
212     .align  KDI_MSR_PATCHOFF;     \
213     KDI_MSR_PATCH;               \
214     jmp     kdi_cmnint;            \
215     SET_SIZE(kdi_invaltrap)

214     .data
215     DGDEF3(kdi_idt, 16 * NIDT, MMU_PAGESIZE)
216     .fill  MMU_PAGESIZE, 1, 0

218 #if !defined(__xpv)
219 .section ".text"
220 .align MMU_PAGESIZE
221 .global kdi_isr_start
222 kdi_isr_start:
223     nop

225 .global kpti_safe_cr3
226 .global kpti_kbase
227 #endif

229 /*
230 * The handlers themselves
231 */

233     MKINVALHDLR
234     MKTRAPHDLR(0)
235     MKTRAPHDLR(1)
236     MKNMIHDLR/*2*/
237     MKTRAPHDLR(3)
238     MKTRAPHDLR(4)
239     MKTRAPHDLR(5)
240     MKTRAPHDLR(6)
241     MKTRAPHDLR(7)
242     MKTRAPHDLR(9)
243     MKTRAPHDLR(15)
244     MKTRAPHDLR(16)
245     MKMCEHDLR/*18*/
246     MKTRAPHDLR(17)
247     MKTRAPHDLR(18)
248     MKTRAPHDLR(19)
249     MKTRAPHDLR(20)

249     MKTRAPERRHDLR(8)
250     MKTRAPERRHDLR(10)
251     MKTRAPERRHDLR(11)
252     MKTRAPERRHDLR(12)
253     MKTRAPERRHDLR(13)

```

```

254     MKTRAPERRHDLR(14)
255     MKTRAPERRHDLR(17)

257     .globl  kdi_ivct_size
258 kdi_ivct_size:
259     .NWORD [kdi_ivct33-kdi_ivct32]

261     /* 10 billion and one interrupt handlers */
262 kdi_ivct_base:
263     MKIVCT(32);    MKIVCT(33);    MKIVCT(34);    MKIVCT(35);
264     MKIVCT(36);    MKIVCT(37);    MKIVCT(38);    MKIVCT(39);
265     MKIVCT(40);    MKIVCT(41);    MKIVCT(42);    MKIVCT(43);
266     MKIVCT(44);    MKIVCT(45);    MKIVCT(46);    MKIVCT(47);
267     MKIVCT(48);    MKIVCT(49);    MKIVCT(50);    MKIVCT(51);
268     MKIVCT(52);    MKIVCT(53);    MKIVCT(54);    MKIVCT(55);
269     MKIVCT(56);    MKIVCT(57);    MKIVCT(58);    MKIVCT(59);
270     MKIVCT(60);    MKIVCT(61);    MKIVCT(62);    MKIVCT(63);
271     MKIVCT(64);    MKIVCT(65);    MKIVCT(66);    MKIVCT(67);
272     MKIVCT(68);    MKIVCT(69);    MKIVCT(70);    MKIVCT(71);
273     MKIVCT(72);    MKIVCT(73);    MKIVCT(74);    MKIVCT(75);
274     MKIVCT(76);    MKIVCT(77);    MKIVCT(78);    MKIVCT(79);
275     MKIVCT(80);    MKIVCT(81);    MKIVCT(82);    MKIVCT(83);
276     MKIVCT(84);    MKIVCT(85);    MKIVCT(86);    MKIVCT(87);
277     MKIVCT(88);    MKIVCT(89);    MKIVCT(90);    MKIVCT(91);
278     MKIVCT(92);    MKIVCT(93);    MKIVCT(94);    MKIVCT(95);
279     MKIVCT(96);    MKIVCT(97);    MKIVCT(98);    MKIVCT(99);
280     MKIVCT(100);   MKIVCT(101);   MKIVCT(102);   MKIVCT(103);
281     MKIVCT(104);   MKIVCT(105);   MKIVCT(106);   MKIVCT(107);
282     MKIVCT(108);   MKIVCT(109);   MKIVCT(110);   MKIVCT(111);
283     MKIVCT(112);   MKIVCT(113);   MKIVCT(114);   MKIVCT(115);
284     MKIVCT(116);   MKIVCT(117);   MKIVCT(118);   MKIVCT(119);
285     MKIVCT(120);   MKIVCT(121);   MKIVCT(122);   MKIVCT(123);
286     MKIVCT(124);   MKIVCT(125);   MKIVCT(126);   MKIVCT(127);
287     MKIVCT(128);   MKIVCT(129);   MKIVCT(130);   MKIVCT(131);
288     MKIVCT(132);   MKIVCT(133);   MKIVCT(134);   MKIVCT(135);
289     MKIVCT(136);   MKIVCT(137);   MKIVCT(138);   MKIVCT(139);
290     MKIVCT(140);   MKIVCT(141);   MKIVCT(142);   MKIVCT(143);
291     MKIVCT(144);   MKIVCT(145);   MKIVCT(146);   MKIVCT(147);
292     MKIVCT(148);   MKIVCT(149);   MKIVCT(150);   MKIVCT(151);
293     MKIVCT(152);   MKIVCT(153);   MKIVCT(154);   MKIVCT(155);
294     MKIVCT(156);   MKIVCT(157);   MKIVCT(158);   MKIVCT(159);
295     MKIVCT(160);   MKIVCT(161);   MKIVCT(162);   MKIVCT(163);
296     MKIVCT(164);   MKIVCT(165);   MKIVCT(166);   MKIVCT(167);
297     MKIVCT(168);   MKIVCT(169);   MKIVCT(170);   MKIVCT(171);
298     MKIVCT(172);   MKIVCT(173);   MKIVCT(174);   MKIVCT(175);
299     MKIVCT(176);   MKIVCT(177);   MKIVCT(178);   MKIVCT(179);
300     MKIVCT(180);   MKIVCT(181);   MKIVCT(182);   MKIVCT(183);
301     MKIVCT(184);   MKIVCT(185);   MKIVCT(186);   MKIVCT(187);
302     MKIVCT(188);   MKIVCT(189);   MKIVCT(190);   MKIVCT(191);
303     MKIVCT(192);   MKIVCT(193);   MKIVCT(194);   MKIVCT(195);
304     MKIVCT(196);   MKIVCT(197);   MKIVCT(198);   MKIVCT(199);
305     MKIVCT(200);   MKIVCT(201);   MKIVCT(202);   MKIVCT(203);
306     MKIVCT(204);   MKIVCT(205);   MKIVCT(206);   MKIVCT(207);
307     MKIVCT(208);   MKIVCT(209);   MKIVCT(210);   MKIVCT(211);
308     MKIVCT(212);   MKIVCT(213);   MKIVCT(214);   MKIVCT(215);
309     MKIVCT(216);   MKIVCT(217);   MKIVCT(218);   MKIVCT(219);
310     MKIVCT(220);   MKIVCT(221);   MKIVCT(222);   MKIVCT(223);
311     MKIVCT(224);   MKIVCT(225);   MKIVCT(226);   MKIVCT(227);
312     MKIVCT(228);   MKIVCT(229);   MKIVCT(230);   MKIVCT(231);
313     MKIVCT(232);   MKIVCT(233);   MKIVCT(234);   MKIVCT(235);
314     MKIVCT(236);   MKIVCT(237);   MKIVCT(238);   MKIVCT(239);
315     MKIVCT(240);   MKIVCT(241);   MKIVCT(242);   MKIVCT(243);
316     MKIVCT(244);   MKIVCT(245);   MKIVCT(246);   MKIVCT(247);
317     MKIVCT(248);   MKIVCT(249);   MKIVCT(250);   MKIVCT(251);
318     MKIVCT(252);   MKIVCT(253);   MKIVCT(254);   MKIVCT(255);

```

```
320 #if !defined(__xpv)
321 .section ".text"
322 .align MMU_PAGESIZE
323 .global kdi_isr_end
324 kdi_isr_end:
325     nop
326 #endif

328 #endif /* !__lint */
```

new/usr/src/uts/intel/kdi/kdi_offsets.in

1

1740 Fri Apr 6 17:25:13 2018

new/usr/src/uts/intel/kdi/kdi_offsets.in

8956 Implement KPTI

Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>

Reviewed by: Robert Mustacchi <rm@joyent.com>

9210 remove KMDB branch debugging support

9211 ::crregs could do with cr2/cr3 support

9209 ::ttrace should be able to filter by thread

Reviewed by: Patrick Mooney <patrick.mooney@joyent.com>

Reviewed by: Yuri Pankov <yuripv@yuripv.net>

```
1 \
2 \ Copyright 2007 Sun Microsystems, Inc. All rights reserved.
3 \ Use is subject to license terms.
4 \
5 \ Copyright 2018 Joyent, Inc.
6 \
7 \ CDDL HEADER START
8 \
9 \ The contents of this file are subject to the terms of the
10 \ Common Development and Distribution License (the "License").
11 \ You may not use this file except in compliance with the License.
12 \
13 \ You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
14 \ or http://www.opensolaris.org/os/licensing.
15 \ See the License for the specific language governing permissions
16 \ and limitations under the License.
17 \
18 \ When distributing Covered Code, include this CDDL HEADER in each
19 \ file and include the License file at usr/src/OPENSOLARIS.LICENSE.
20 \ If applicable, add the following below this CDDL HEADER, with the
21 \ fields enclosed by brackets "[]" replaced with your own identifying
22 \ information: Portions Copyright [yyyy] [name of copyright owner]
23 \
24 \ CDDL HEADER END
25 \
26 \ ident "%Z%M% %I% %E% SMI"
27 \
28 \ Keep in sync with kdi_state.h
29 \
```

```
29 #include <sys/cpuvar.h>
30 #include <sys/kdi_impl.h>
```

```
32 kdi_memrange_t MR_SIZE
33 mr_base
34 mr_lim

36 kdi_crumb_t KRM_SIZE
37 krm_cpu_state
38 krm_pc
39 krm_sp
40 krm_trapno
41 krm_flag

43 kdi_drreg_t
44 dr_ctl
45 dr_stat
46 dr_addr

50 kdi_msr_t MSR_SIZE
51 msr_num
52 msr_type
```

new/usr/src/uts/intel/kdi/kdi_offsets.in

2

```
53 _u._msr_valp MSR_VALP
54 _u._msr_val MSR_VAL

48 kdi_cpusave_t KRS_SIZE
49 krs_gregs
50 krs_dr
51 krs_dr.dr_ctl KRS_DRCTL
52 krs_dr.dr_stat KRS_DRSTAT
53 krs_gdt
54 krs_idt
55 krs_cr0
64 krs_msr
56 krs_cpu_state
57 krs_curcrumbidx
58 krs_curcrumb
59 krs_crumbs

70 cpu
71 cpu_id

61 greg_t KREG_SIZE

75 #if defined(__amd64)
63 #define REG_SHIFT 3
77 #else
78 #define REG_SHIFT 2
79 #endif

65 #define DRADDR_IDX(num) _CONST(_MUL(num, DR_ADDR_INCR))
66 #define DRADDR_OFF(num) _CONST(DRADDR_IDX(num) + DR_ADDR)
67 #define KRS_DROFF(num) _CONST(DRADDR_OFF(num) + KRS_DR)
68 #define REG_OFF(reg) _CONST(_CONST(reg) << REG_SHIFT)
69 #define KDIREG_OFF(reg) _CONST(_MUL(KREG_SIZE, reg) + KRS_GREGS)
```

```

*****
4230 Fri Apr 6 17:25:13 2018
new/usr/src/uts/intel/os/arch_kdi.c
9210 remove KMDB branch debugging support
9211 ::crregs could do with cr2/cr3 support
9209 ::ttrace should be able to filter by thread
Reviewed by: Patrick Mooney <patrick.mooney@joyent.com>
Reviewed by: Yuri Pankov <yuripv@yuripv.net>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22  * Copyright 2009 Sun Microsystems, Inc. All rights reserved.
23  * Use is subject to license terms.
24  *
25  * Copyright 2018 Joyent, Inc.
26  */

28 /*
29  * Kernel/Debugger Interface (KDI) routines. Called during debugger under
30  * various system states (boot, while running, while the debugger has control).
31  * Functions intended for use while the debugger has control may not grab any
32  * locks or perform any functions that assume the availability of other system
33  * services.
34  */

36 #include <sys/system.h>
37 #include <sys/x86_archext.h>
38 #include <sys/kdi_impl.h>
39 #include <sys/smp_impldefs.h>
40 #include <sys/psm_types.h>
41 #include <sys/segments.h>
42 #include <sys/archsystem.h>
43 #include <sys/controlregs.h>
44 #include <sys/trap.h>
45 #include <sys/kobj.h>
46 #include <sys/kobj_impl.h>
47 #include <sys/clock_impl.h>

49 static void
50 kdi_system_claim(void)
51 {
52     lbolt_debug_entry();

54     psm_notifyf(PSM_DEBUG_ENTER);
55 }
_____unchanged_portion_omitted_____

```

```

146 /*
147  * On Intel, most of these are shared between i86*, so this is really an
148  * arch_kdi_init().
149  */
150 void
151 mach_kdi_init(kdi_t *kdi)
152 {
153     kdi->kdi_plat_call = kdi_plat_call;
154     kdi->kdi_kmdb_enter = kmdb_enter;
155     kdi->mkdi_activate = kdi_activate;
156     kdi->mkdi_deactivate = kdi_deactivate;
157     kdi->mkdi_idt_switch = kdi_idt_switch;
158     kdi->mkdi_update_drreg = kdi_update_drreg;
159     kdi->mkdi_set_debug_msrs = kdi_set_debug_msrs;
160     kdi->mkdi_get_userlimit = kdi_get_userlimit;
161     kdi->mkdi_get_cpuserinfo = kdi_get_cpuserinfo;
162     kdi->mkdi_stop_slaves = kdi_stop_slaves;
163     kdi->mkdi_start_slaves = kdi_start_slaves;
164     kdi->mkdi_slave_wait = kdi_slave_wait;
165     kdi->mkdi_memrange_add = kdi_memrange_add;
166     kdi->mkdi_reboot = kdi_reboot;
167 }
_____unchanged_portion_omitted_____

```

new/usr/src/uts/intel/sys/archsystem.h

1

```
*****
6372 Fri Apr 6 17:25:13 2018
new/usr/src/uts/intel/sys/archsystem.h
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */

22 /*
23  * Copyright (c) 1993, 2010, Oracle and/or its affiliates. All rights reserved.
24  * Copyright 2018 Joyent, Inc.
25  * Copyright 2017 Joyent, Inc.
26 */

27 #ifndef _SYS_ARCHSYSTEM_H
28 #define _SYS_ARCHSYSTEM_H

30 /*
31  * A selection of ISA-dependent interfaces
32 */

34 #include <vm/seg_enum.h>
35 #include <vm/page.h>

37 #ifdef __cplusplus
38 extern "C" {
39 #endif

41 #ifdef _KERNEL

43 extern greg_t getfp(void);
44 extern int getpil(void);

46 extern ulong_t getcr0(void);
47 extern void setcr0(ulong_t);
48 extern ulong_t getcr8(void);
49 extern void setcr8(ulong_t);
50 extern ulong_t getcr2(void);
51 extern void clflush_insn(caddr_t addr);
52 extern void mfence_insn(void);

54 #if defined(__i386)
55 extern uint16_t getgs(void);
56 extern void setgs(uint16_t);

58 extern void patch_sse(void);
```

new/usr/src/uts/intel/sys/archsystem.h

2

```
59 extern void patch_sse2(void);
60 #endif

62 extern void patch_xsave(void);
63 extern kmem_cache_t *fpsave_cache;

65 extern void cli(void);
66 extern void sti(void);

68 extern void tenmicrosec(void);

70 extern void restore_int_flag(ulong_t);
71 extern void intr_restore(ulong_t);
72 extern ulong_t clear_int_flag(void);
73 extern ulong_t intr_clear(void);
74 extern ulong_t getflags(void);
75 extern int interrupts_enabled(void);

77 extern void int3(void);
78 extern void int18(void);
79 extern void int20(void);
80 extern void int_cmci(void);

82 #if defined(__amd64)
83 extern void sys_syscall(), tr_sys_syscall();
84 extern void sys_syscall32(), tr_sys_syscall32();
85 extern void sys_syscall();
86 extern void sys_syscall32();
87 extern void sys_lcall32();
88 extern void sys_syscall_int();
89 extern void tr_sys_syscall_int();
90 extern void brand_sys_syscall(), tr_brand_sys_syscall();
91 extern void brand_sys_syscall32(), tr_brand_sys_syscall32();
92 extern void brand_sys_syscall_int(), tr_brand_sys_syscall_int();
93 extern void brand_sys_syscall();
94 extern void brand_sys_syscall32();
95 extern void brand_sys_syscall_int();
96 extern int update_sregs();
97 extern void reset_sregs();
98 #elif defined(__i386)
99 extern void sys_call();
100 extern void tr_sys_call();
101 extern void brand_sys_call();
102 #endif
103 extern void sys_sysenter();
104 extern void tr_sys_sysenter();
105 extern void _sys_sysenter_post_swapgs();
106 extern void brand_sys_sysenter();
107 extern void tr_brand_sys_sysenter();
108 extern void _brand_sys_sysenter_post_swapgs();

109 extern void dosyscall(void);

111 extern void bind_hwcap(void);

113 extern uint16_t inw(int port);
114 extern uint32_t inl(int port);
115 extern void outw(int port, uint16_t value);
116 extern void outl(int port, uint32_t value);

118 extern void pc_reset(void) __NORETURN;
119 extern void efi_reset(void) __NORETURN;
120 extern void reset(void) __NORETURN;
121 extern int goany(void);

123 extern void setgregs(klwp_t *, gregset_t);
```

```

120 extern void getgregs(klwp_t *, gregset_t);
121 extern void setfpregs(klwp_t *, fpregset_t *);
122 extern void getfpregs(klwp_t *, fpregset_t *);

124 #if defined(_SYSCALL32_IMPL)
125 extern void getgregs32(klwp_t *, gregset32_t);
126 extern void setfpregs32(klwp_t *, fpregset32_t *);
127 extern void getfpregs32(klwp_t *, fpregset32_t *);
128 #endif

130 struct fpu_ctx;

132 extern void fp_free(struct fpu_ctx *, int);
133 extern void fp_save(struct fpu_ctx *);
134 extern void fp_restore(struct fpu_ctx *);

136 extern int fpu_pentium_fdivbug;

138 extern void sep_save(void *);
139 extern void sep_restore(void *);

141 extern void brand_interpositioning_enable(void);
142 extern void brand_interpositioning_disable(void);

144 struct regs;

146 extern int instr_size(struct regs *, caddr_t *, enum seg_rw);

148 extern int enable_cbcpc; /* patchable in /etc/system */

150 extern uint_t cpu_hwcap_flags;
151 extern uint_t cpu_freq;
152 extern uint64_t cpu_freq_hz;

154 extern int use_sse_pagecopy;
155 extern int use_sse_pagezero;
156 extern int use_sse_copy;

158 extern caddr_t i86devmap(pfn_t, pgcnt_t, uint_t);
159 extern page_t *page_numtopp_alloc(pfn_t pfnnum);

161 extern void hwblkclr(void *, size_t);
162 extern void hwblkpagecopy(const void *, void *);
163 extern void page_copy_no_xmm(void *dst, void *src);
164 extern void block_zero_no_xmm(void *dst, int len);
165 #define BLOCKZEROALIGN (4 * sizeof (void *))

167 extern void (*kcpc_hw_enable_cpc_intr)(void);

169 extern void init_desctbbs(void);

171 extern user_desc_t *cpu_get_gdt(void);

173 extern void switch_sp_and_call(void *, void (*)(uint_t, uint_t), uint_t,
174     uint_t);
175 extern hrtime_t (*gethrtimef)(void);
176 extern hrtime_t (*gethrtimeunscaledf)(void);
177 extern void (*scalehrtimef)(hrtime_t *);
178 extern uint64_t (*unscalehrtimef)(hrtime_t);
179 extern void (*gethrestimef)(timestruc_t *);

181 extern void av_dispatch_softvect(uint_t);
182 extern void av_dispatch_autovect(uint_t);
183 extern uint_t atomic_btr32(uint32_t *, uint_t);
184 extern uint_t bsrw_insn(uint16_t);
185 extern int sys_rtt_common(struct regs *);

```

```

186 extern void fakesoftint(void);

188 extern void *plat_traceback(void *);

190 /*
191 * The following two macros are the four byte instruction sequence of stac, ret
192 * and clac, ret. These are used in startup_smap() as a part of properly setting
193 * up the valid instructions. For more information on SMAP, see
194 * uts/intel/ia32/ml/copy.s.
195 */
196 #define SMAP_CLAC_INSTR 0xc3ca010f
197 #define SMAP_STAC_INSTR 0xc3cb010f
198 extern void smap_disable(void);
199 extern void smap_enable(void);

201 #if defined(__xpv)
202 extern void xen_init_callbacks(void);
203 extern void xen_set_callback(void (*)(void), uint_t, uint_t);
204 extern void xen_printf(const char *, ...);
205 #define cpr_dprintf xen_printf
206 extern int xpv_panicking;
207 #define IN_XPV_PANIC() (xpv_panicking > 0)
208 #else
209 extern void setup_mca(void);
210 extern void pat_sync(void);
211 extern void patch_tsc_read(int);
212 #if defined(__amd64) && !defined(__xpv)
213 extern void patch_memops(uint_t);
214 #endif /* defined(__amd64) && !defined(__xpv) */
215 extern void setup_xfem(void);
216 #define cpr_dprintf prom_printf
217 #define IN_XPV_PANIC() (__lintzero)
218 #endif

220 #endif /* _KERNEL */

222 #if defined(_KERNEL) || defined(_BOOT)
223 extern uint8_t inb(int port);
224 extern void outb(int port, uint8_t value);
225 #endif

227 #ifdef __cplusplus
228 }

```

_____ unchanged portion omitted

new/usr/src/uts/intel/sys/controlregs.h

1

```
*****
8003 Fri Apr 6 17:25:13 2018
new/usr/src/uts/intel/sys/controlregs.h
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>
9210 remove KMDB branch debugging support
9211 ::crregs could do with cr2/cr3 support
9209 ::ttrace should be able to filter by thread
Reviewed by: Patrick Mooney <patrick.mooney@joyent.com>
Reviewed by: Yuri Pankov <yuripv@yuripv.net>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22 * Copyright (c) 2004, 2010, Oracle and/or its affiliates. All rights reserved.
23 * Copyright 2018, Joyent, Inc.
24 */

26 #ifndef _SYS_CONTROLREGS_H
27 #define _SYS_CONTROLREGS_H

29 #ifndef _ASM
30 #include <sys/types.h>
31 #endif

33 #ifdef __cplusplus
34 extern "C" {
35 #endif

37 /*
38 * This file describes the x86 architecture control registers which
39 * are part of the privileged architecture.
40 *
41 * Many of these definitions are shared between IA-32-style and
42 * AMD64-style processors.
43 */

45 /* CR0 Register */

47 #define CR0_PG 0x80000000 /* paging enabled */
48 #define CR0_CD 0x40000000 /* cache disable */
49 #define CR0_NW 0x20000000 /* not writethrough */
50 #define CR0_AM 0x00040000 /* alignment mask */
51 #define CR0_WP 0x00010000 /* write protect */
52 #define CR0_NE 0x00000020 /* numeric error */
53 #define CR0_ET 0x00000010 /* extension type */
```

new/usr/src/uts/intel/sys/controlregs.h

2

```
54 #define CR0_TS 0x00000008 /* task switch */
55 #define CR0_EM 0x00000004 /* emulation */
56 #define CR0_MP 0x00000002 /* monitor coprocessor */
57 #define CR0_PE 0x00000001 /* protection enabled */

59 /* XX64 eliminate these compatibility defines */

61 #define CR0_CE CR0_CD
62 #define CR0_WT CR0_NW

64 #define FMT_CR0 \
65     "\20\40pg\37cd\36nw\35am\21wp\6ne\5et\4ts\3em\2mp\1pe"

67 /*
68 * Set the FPU-related control bits to explain to the processor that
69 * we're managing FPU state:
70 * - set monitor coprocessor (allow TS bit to control FPU)
71 * - set numeric exception (disable IGNNE# mechanism)
72 * - set task switch (#nm on first fp instruction)
73 * - clear emulate math bit (cause we're not emulating!)
74 */
75 #define CR0_ENABLE_FPU_FLAGS(cr) \
76     (((cr) | CR0_MP | CR0_NE | CR0_TS) & ~(uint32_t)~CR0_EM)

78 /*
79 * Set the FPU-related control bits to explain to the processor that
80 * we're -not- managing FPU state:
81 * - set emulate (all fp instructions cause #nm)
82 * - clear monitor coprocessor (so fwait/wait doesn't #nm)
83 */
84 #define CR0_DISABLE_FPU_FLAGS(cr) \
85     (((cr) | CR0_EM) & ~(uint32_t)~CR0_MP)

87 /* CR3 Register */

89 #define CR3_PCD 0x00000010 /* cache disable */
90 #define CR3_PWT 0x00000008 /* write through */
91 #if defined(_ASM)
92 #define CR3_NOINVL_BIT 0x8000000000000000
93 #else
94 #define CR3_NOINVL_BIT 0x8000000000000000ULL /* no invalidation */
95 #endif
96 #define PCID_NONE 0x000 /* generic PCID */
97 #define PCID_KERNEL 0x000 /* kernel's PCID */
98 #define PCID_USER 0x001 /* user-space PCID */

92 #define FMT_CR3 "\20\5pcd\4pwt"

100 /* CR4 Register */

102 #define CR4_VME 0x0001 /* virtual-8086 mode extensions */
103 #define CR4_PVI 0x0002 /* protected-mode virtual interrupts */
104 #define CR4_TSD 0x0004 /* time stamp disable */
105 #define CR4_DE 0x0008 /* debugging extensions */
106 #define CR4_PSE 0x0010 /* page size extensions */
107 #define CR4_PAE 0x0020 /* physical address extension */
108 #define CR4_MCE 0x0040 /* machine check enable */
109 #define CR4_PGE 0x0080 /* page global enable */
110 #define CR4_PCE 0x0100 /* perf-monitoring counter enable */
111 #define CR4_OSFCSR 0x0200 /* OS fxsave/fxrstor support */
112 #define CR4_OSXMMEXCPT 0x0400 /* OS unmasked exception support */
113 /* 0x0800 reserved */
114 /* 0x1000 reserved */
115 #define CR4_VMXE 0x2000
116 #define CR4_SMXE 0x4000
117 #define CR4_PCIDE 0x20000 /* PCID enable */
```

```

118 #define CR4_OSXSAVE      0x40000      /* OS xsave/xrestore support */
119 #define CR4_SMEP        0x100000     /* NX for user pages in kernel */
120 #define CR4_SMAP        0x200000     /* kernel can't access user pages */

122 #define FMT_CR4          \
123     "\20\26smap\25smep\23osxsave\22pcide" \
124     "\20\26smap\25smep\23osxsave" \
125     "\17smxe\16vmxe\13xme\12fxsr\11pce\10pge" \
126     "\7mce\6pae\5pse\4de\3tsd\2pvi\lvme"

127 /*
128 * Enable the SSE-related control bits to explain to the processor that
129 * we're managing XMM state and exceptions
130 */
131 #define CR4_ENABLE_SSE_FLAGS(cr) \
132     ((cr) | CR4_OSFXSR | CR4_OSXMMEXCPT)

134 /*
135 * Disable the SSE-related control bits to explain to the processor
136 * that we're NOT managing XMM state
137 */
138 #define CR4_DISABLE_SSE_FLAGS(cr) \
139     ((cr) & ~(uint32_t)(CR4_OSFXSR | CR4_OSXMMEXCPT))

141 /* Intel's SYSENTER configuration registers */

143 #define MSR_INTC_SEP_CS 0x174          /* kernel code selector MSR */
144 #define MSR_INTC_SEP_ESP 0x175        /* kernel esp MSR */
145 #define MSR_INTC_SEP_EIP 0x176        /* kernel eip MSR */

147 /* Intel's microcode registers */
148 #define MSR_INTC_UCODE_WRITE 0x79     /* microcode write */
149 #define MSR_INTC_UCODE_REV 0x8b      /* microcode revision */
150 #define INTC_UCODE_REV_SHIFT 32      /* Bits 63:32 */

152 /* Intel's platform identification */
153 #define MSR_INTC_PLATFORM_ID 0x17     /* Bit 52:50 */
154 #define INTC_PLATFORM_ID_SHIFT 50    /* Bit 52:50 */
155 #define INTC_PLATFORM_ID_MASK 0x7

157 /* AMD's EFER register */

159 #define MSR_AMD_EFER 0xc0000080      /* extended feature enable MSR */

161 #define AMD_EFER_FFXSR 0x4000        /* fast fxsave/fxrstor */
162 #define AMD_EFER_SVME 0x1000         /* svm enable */
163 #define AMD_EFER_NXE 0x0800          /* no-execute enable */
164 #define AMD_EFER_LMA 0x0400          /* long mode active (read-only) */
165 #define AMD_EFER_LME 0x0100          /* long mode enable */
166 #define AMD_EFER_SCE 0x0001          /* system call extensions */

168 #define FMT_AMD_EFER \
169     "\20\17ffxsr\15svme\14nxe\13lma\11lme\1sce"

171 /* AMD's SYSCFG register */

173 #define MSR_AMD_SYSCFG 0xc0000010     /* system configuration MSR */

175 #define AMD_SYSCFG_TOM2 0x200000      /* MtrrTom2En */
176 #define AMD_SYSCFG_MVDM 0x100000     /* MtrrVarDramEn */
177 #define AMD_SYSCFG_MFDM 0x080000     /* MtrrFixDramModEn */
178 #define AMD_SYSCFG_MFDE 0x040000     /* MtrrFixDramEn */

180 #define FMT_AMD_SYSCFG \
181     "\20\26tom2\25mvdM\24mfDM\23mfde"

```

```

183 /* AMD's syscall/sysret MSRs */

185 #define MSR_AMD_STAR 0xc0000081      /* %cs:%ss:%cs:%ss:%eip for syscall */
186 #define MSR_AMD_LSTAR 0xc0000082    /* target %rip of 64-bit syscall */
187 #define MSR_AMD_CSTAR 0xc0000083    /* target %rip of 32-bit syscall */
188 #define MSR_AMD_SFMASK 0xc0000084   /* syscall flag mask */

190 /* AMD's FS.base and GS.base MSRs */

192 #define MSR_AMD_FSBASE 0xc0000100    /* 64-bit base address for %fs */
193 #define MSR_AMD_GSBASE 0xc0000101    /* 64-bit base address for %gs */
194 #define MSR_AMD_KGSBASE 0xc0000102   /* swapgs swaps this with gsbases */
195 #define MSR_AMD_TSCAUX 0xc0000103    /* %ecx value on rdtscp insn */

197 /* AMD's configuration MSRs, weakly documented in the revision guide */

199 #define MSR_AMD_DC_CFG 0xc0011022

201 #define AMD_DC_CFG_DIS_CNV_WC_SSO (UINT64_C(1) << 3)
202 #define AMD_DC_CFG_DIS_SMC_CHK_BUF (UINT64_C(1) << 10)

204 /* AMD's HWCR MSR */

206 #define MSR_AMD_HWCR 0xc0010015

208 #define AMD_HWCR_TLBCACHEDIS (UINT64_C(1) << 3)
209 #define AMD_HWCR_FFDIS 0x000400     /* disable TLB Flush Filter */
210 #define AMD_HWCR_MCI_STATUS_WREN 0x40000 /* enable write of MCI_STATUS */

212 /* AMD's NorthBridge Config MSR, SHOULD ONLY BE WRITTEN TO BY BIOS */

214 #define MSR_AMD_NB_CFG 0xc001001f

216 #define AMD_NB_CFG_SRQ_HEARTBEAT (UINT64_C(1) << 20)
217 #define AMD_NB_CFG_SRQ_SPR (UINT64_C(1) << 32)

219 #define MSR_AMD_BU_CFG 0xc0011023

221 #define AMD_BU_CFG_E298 (UINT64_C(1) << 1)

223 #define MSR_AMD_DE_CFG 0xc0011029

225 #define AMD_DE_CFG_E721 (UINT64_C(1))

227 /* AMD's OSVW MSRs */
228 #define MSR_AMD_OSVW_ID_LEN 0xc0010140
229 #define MSR_AMD_OSVW_STATUS 0xc0010141

232 #define OSVW_ID_LEN_MASK 0xffffFULL
233 #define OSVW_ID_CNT_PER_MSR 64

235 /*
236 * Enable PCI Extended Configuration Space (ECS) on Greyhound
237 */
238 #define AMD_GH_NB_CFG_EN_ECS (UINT64_C(1) << 46)

240 /* AMD microcode patch loader */
241 #define MSR_AMD_PATCHLEVEL 0x8b
242 #define MSR_AMD_PATCHLOADER 0xc0010020

244 #ifdef __cplusplus
245 }
246 #endif

248 #endif /* !_SYS_CONTROLREGS_H */

```


new/usr/src/uts/intel/sys/kdi_machimpl.h

1

```
*****
3809 Fri Apr 6 17:25:13 2018
new/usr/src/uts/intel/sys/kdi_machimpl.h
9210 remove KMDB branch debugging support
9211 ::crregs could do with cr2/cr3 support
9209 ::ttrace should be able to filter by thread
Reviewed by: Patrick Mooney <patrick.mooney@joyent.com>
Reviewed by: Yuri Pankov <yuripv@yuripv.net>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22  * Copyright 2007 Sun Microsystems, Inc. All rights reserved.
23  * Use is subject to license terms.
24 *
25 * Copyright 2018 Joyent, Inc.
26 */

28 #ifndef _SYS_KDI_MACHIMPL_H
29 #define _SYS_KDI_MACHIMPL_H

29 #pragma ident "%Z%M% %I% %E% SMI"

31 /*
32  * The Kernel/Debugger interface. The operations provided by the kdi_t,
33  * defined below, comprise the Debugger -> Kernel portion of the interface,
34  * and are to be used only when the system has been stopped.
35 */

37 #include <sys/modctl.h>
38 #include <sys/types.h>
39 #include <sys/cpuvar.h>
40 #include <sys/kdi_regs.h>

42 #ifdef __cplusplus
43 extern "C" {
44 #endif

46 typedef void (*kdi_main_t)(kdi_cpusave_t *);

48 typedef struct kdi_memrange {
49     caddr_t mr_base;
50     caddr_t mr_lim;
51 } kdi_memrange_t;

53 #define KDI_MEMRANGES_MAX 2

55 typedef struct kdi_mach {
```

new/usr/src/uts/intel/sys/kdi_machimpl.h

2

```
56 void (*mkdi_activate)(kdi_main_t, kdi_cpusave_t *, uint_t);
57 void (*mkdi_deactivate)(void);

59 void (*mkdi_idt_switch)(kdi_cpusave_t *);

61 void (*mkdi_update_drreg)(kdi_drreg_t *);
62 void (*mkdi_set_debug_msrs)(kdi_msr_t *);

63 uintptr_t (*mkdi_get_userlimit)(void);

65 int (*mkdi_get_cpuinfo)(uint_t *, uint_t *, uint_t *);

67 void (*mkdi_stop_slaves)(int, int);

69 void (*mkdi_start_slaves)(void);

71 void (*mkdi_slave_wait)(void);

73 void (*mkdi_memrange_add)(caddr_t, size_t);

75 void (*mkdi_reboot)(void);
76 } kdi_mach_t;

78 #define mkdi_activate kdi_mach.mkdi_activate
79 #define mkdi_deactivate kdi_mach.mkdi_deactivate
80 #define mkdi_idt_switch kdi_mach.mkdi_idt_switch
81 #define mkdi_update_drreg kdi_mach.mkdi_update_drreg
83 #define mkdi_set_debug_msrs kdi_mach.mkdi_set_debug_msrs
82 #define mkdi_get_userlimit kdi_mach.mkdi_get_userlimit
83 #define mkdi_get_cpuinfo kdi_mach.mkdi_get_cpuinfo
84 #define mkdi_stop_slaves kdi_mach.mkdi_stop_slaves
85 #define mkdi_start_slaves kdi_mach.mkdi_start_slaves
86 #define mkdi_slave_wait kdi_mach.mkdi_slave_wait
87 #define mkdi_memrange_add kdi_mach.mkdi_memrange_add
88 #define mkdi_reboot kdi_mach.mkdi_reboot

90 extern void hat_kdi_init(void);

92 extern ulong_t kdi_getdr0(void), kdi_getdr1(void), kdi_getdr2(void);
93 extern ulong_t kdi_getdr3(void), kdi_getdr6(void), kdi_getdr7(void);
94 extern void kdi_setdr0(ulong_t), kdi_setdr1(ulong_t), kdi_setdr2(ulong_t);
95 extern void kdi_setdr3(ulong_t), kdi_setdr6(ulong_t), kdi_setdr7(ulong_t);
96 extern ulong_t kdi_dreg_get(int);
97 extern void kdi_dreg_set(int, ulong_t);
98 extern void kdi_update_drreg(kdi_drreg_t *);
101 extern void kdi_set_debug_msrs(kdi_msr_t *);
99 extern void kdi_cpu_debug_init(kdi_cpusave_t *);

101 extern void kdi_cpu_init(void);
102 extern void kdi_xc_others(int, void (*)(void));
103 extern void kdi_start_slaves(void);
104 extern void kdi_slave_wait(void);

106 extern void kdi_idtr_set(gate_desc_t *, size_t);
107 extern void kdi_idt_write(struct gate_desc *, uint_t);
108 extern void kdi_idt_sync(void);
109 extern void kdi_idt_switch(kdi_cpusave_t *);
110 #ifdef __xpv
111 extern void kdi_idtr_write(desctbr_t *);
112 #else
113 #define kdi_idtr_write(idtr) wr_idtr(idtr)
114 #endif

116 extern void kdi_activate(kdi_main_t, kdi_cpusave_t *, uint_t);
117 extern void kdi_deactivate(void);
118 extern void kdi_stop_slaves(int, int);
```

new/usr/src/uts/intel/sys/kdi_machimpl.h

3

```
119 extern void kdi_memrange_add(caddr_t, size_t);  
120 extern void kdi_reboot(void);
```

```
122 #ifdef __cplusplus  
123 }
```

_____unchanged_portion_omitted_____

new/usr/src/uts/intel/sys/kdi_regs.h

1

```
*****
3021 Fri Apr 6 17:25:14 2018
new/usr/src/uts/intel/sys/kdi_regs.h
9210 remove KMDB branch debugging support
9211 ::crregs could do with cr2/cr3 support
9209 ::ttrace should be able to filter by thread
Reviewed by: Patrick Mooney <patrick.mooney@joyent.com>
Reviewed by: Yuri Pankov <yuripv@yuripv.net>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22  * Copyright 2007 Sun Microsystems, Inc. All rights reserved.
23  * Use is subject to license terms.
24  *
25  * Copyright 2018 Joyent, Inc.
26  */

28 #ifndef _SYS_KDI_REGS_H
29 #define _SYS_KDI_REGS_H

29 #pragma ident "%Z%M% %I% %E% SMI"

31 #ifndef _ASM
32 #include <sys/types.h>
33 #include <sys/segments.h>
34 #include <sys/regset.h>
35 #include <sys/privregs.h>
36 #endif

38 #if defined(__amd64)
39 #include <amd64/sys/kdi_regs.h>
40 #elif defined(__i386)
41 #include <ia32/sys/kdi_regs.h>
42 #endif

44 #ifdef __cplusplus
45 extern "C" {
46 #endif

48 #define KDI_NCRUMBS 5

50 #define KDI_CPU_STATE_NONE 0
51 #define KDI_CPU_STATE_MASTER 1
52 #define KDI_CPU_STATE_SLAVE 2

54 #define KDIREG_DRCTL_WPALLEN_MASK 0x000000ff
55 #define KDIREG_DRSTAT_RESERVED 0xffff0fff
```

new/usr/src/uts/intel/sys/kdi_regs.h

2

```
56 #define KDIREG_DRCTL_RESERVED 0x00000700

58 #define KDI_MSR_READ 0x1 /* read during entry (unlimited) */
59 #define KDI_MSR_WRITE 0x2 /* write during exit (unlimited) */
60 #define KDI_MSR_WRITEDELAY 0x4 /* write after last branch (<= 1) */
61 #define KDI_MSR_CLEARENTRY 0x3 /* clear before 1st branch (<= 1) */

58 #ifndef _ASM

60 /*
61  * We maintain a ring buffer of bread crumbs for debugging purposes. The
62  * current buffer pointer is advanced along the ring with each intercepted
63  * trap (debugger entry, invalid memory access, fault during step, etc).
64  */
65 typedef struct kdi_crumb {
66     greg_t krm_cpu_state; /* This CPU's state at last entry */
67     greg_t krm_pc; /* Instruction pointer at trap */
68     greg_t krm_sp; /* Stack pointer at trap */
69     greg_t krm_trapno; /* The last trap number */
70     greg_t krm_flag; /* KAIF_CRUMB_F_* */
71 } kdi_crumb_t;
_____ unchanged_portion_omitted _____

89 typedef struct kdi_msr {
90     uint_t msr_num;
91     uint_t msr_type;
92     union {
93         uint64_t *_msr_valp;
94         uint64_t _msr_val;
95     } _u;
96 } kdi_msr_t;

98 #define kdi_msr_val _u._msr_val
99 #define kdi_msr_valp _u._msr_valp

84 /*
85  * Data structure used to hold all of the state for a given CPU.
86  */
87 typedef struct kdi_cpusave {
88     greg_t *krs_gregs; /* saved registers */
90     kdi_drreg_t krs_dr; /* saved debug registers */
92     user_desc_t *krs_gdt; /* GDT address */
93     gate_desc_t *krs_idt; /* IDT address */
95     greg_t krs_cr0; /* saved %cr0 */

114     kdi_msr_t *krs_msr; /* ptr to MSR save area */

97     uint_t krs_cpu_state; /* KDI_CPU_STATE_* mstr/slv */
98     uint_t krs_cpu_flushed; /* Have caches been flushed? */
99     uint_t krs_cpu_id; /* this CPU's ID */

101     /* Bread crumb ring buffer */
102     ulong_t krs_curcrumbidx; /* Current krs_crums idx */
103     kdi_crumb_t *krs_curcrumb; /* Pointer to current crumb */
104     kdi_crumb_t krs_crums[KDI_NCRUMBS]; /* Crumbs */
105 } kdi_cpusave_t;
_____ unchanged_portion_omitted _____
```

new/usr/src/uts/intel/sys/segments.h

1

```
*****
24903 Fri Apr 6 17:25:14 2018
new/usr/src/uts/intel/sys/segments.h
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>
*****
1 /*
2  * Copyright (c) 2004, 2010, Oracle and/or its affiliates. All rights reserved.
3  */
4 /*
5  * Copyright 2018 Joyent, Inc.
6  * Copyright 2016 Joyent, Inc.
7  */
8 #ifndef _SYS_SEGMENTS_H
9 #define _SYS_SEGMENTS_H
10
11 #ifdef __cplusplus
12 extern "C" {
13 #endif
14
15 /*
16  * Copyright (c) 1989, 1990 William F. Jolitz
17  * Copyright (c) 1990 The Regents of the University of California.
18  * All rights reserved.
19  *
20  * This code is derived from software contributed to Berkeley by
21  * William Jolitz.
22  *
23  * Redistribution and use in source and binary forms, with or without
24  * modification, are permitted provided that the following conditions
25  * are met:
26  * 1. Redistributions of source code must retain the above copyright
27  * notice, this list of conditions and the following disclaimer.
28  * 2. Redistributions in binary form must reproduce the above copyright
29  * notice, this list of conditions and the following disclaimer in the
30  * documentation and/or other materials provided with the distribution.
31  * 3. All advertising materials mentioning features or use of this software
32  * must display the following acknowledgement:
33  * This product includes software developed by the University of
34  * California, Berkeley and its contributors.
35  * 4. Neither the name of the University nor the names of its contributors
36  * may be used to endorse or promote products derived from this software
37  * without specific prior written permission.
38  *
39  * THIS SOFTWARE IS PROVIDED BY THE REGENTS AND CONTRIBUTORS ``AS IS'' AND
40  * ANY EXPRESS OR IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED TO, THE
41  * IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS FOR A PARTICULAR PURPOSE
42  * ARE DISCLAIMED. IN NO EVENT SHALL THE REGENTS OR CONTRIBUTORS BE LIABLE
43  * FOR ANY DIRECT, INDIRECT, INCIDENTAL, SPECIAL, EXEMPLARY, OR CONSEQUENTIAL
44  * DAMAGES (INCLUDING, BUT NOT LIMITED TO, PROCUREMENT OF SUBSTITUTE GOODS
45  * OR SERVICES; LOSS OF USE, DATA, OR PROFITS; OR BUSINESS INTERRUPTION)
46  * HOWEVER CAUSED AND ON ANY THEORY OF LIABILITY, WHETHER IN CONTRACT, STRICT
47  * LIABILITY, OR TORT (INCLUDING NEGLIGENCE OR OTHERWISE) ARISING IN ANY WAY
48  * OUT OF THE USE OF THIS SOFTWARE, EVEN IF ADVISED OF THE POSSIBILITY OF
49  * SUCH DAMAGE.
50  *
51  * from: @(#)segments.h 7.1 (Berkeley) 5/9/91
52  * $FreeBSD: src/sys/i386/include/segments.h,v 1.34 2003/09/10 01:07:04
53  * jhb Exp $
54  *
55  * 386 Segmentation Data Structures and definitions
56  * William F. Jolitz (william@ernie.berkeley.edu) 6/20/1989
57  */
```

new/usr/src/uts/intel/sys/segments.h

2

```
59 #include <sys/tss.h>
60
61 /*
62  * Selector register format
63  * CS, DS, ES, FS, GS, SS
64  *
65  * 15          3 2 1 0
66  * +-----+-----+-----+
67  * |          SI          |TI|RPL|
68  * +-----+-----+-----+
69  *
70  * SI = selector index
71  * TI = table indicator (0 = GDT, 1 = LDT)
72  * RPL = requestor privilege level
73  */
74 #if !defined(_ASM) || defined(_GNU_C_AS__)
75 #define IDTSEL(s) ((s) << 3) /* index to selector */
76 #define SEL_GDT(s, r) (IDTSEL(s) | r) /* global sel */
77 #else
78 #define IDTSEL(s) [s << 3]
79 #define SEL_GDT(s, r) [IDTSEL(s) | r]
80 #endif
81
82 #define SELTOIDX(s) ((s) >> 3) /* selector to index */
83
84 /*
85  * SEL_(KPL,UPL,XPL) is the RPL or DPL value for code and data selectors
86  * and their descriptors respectively.
87  *
88  * TRP_(KPL,UPL,XPL) is used to indicate the DPL for system gates only.
89  *
90  * This distinction is important to support para-virt guests on the
91  * 64-bit hypervisor. Both guest kernel and user run in ring 3 and the
92  * hypervisor runs in ring 0. When the kernel creates its trap and
93  * interrupt gates it needs a way to prevent any arbitrary int $n
94  * instruction from entering a gate that is not expected. The hypervisor
95  * allows these gates to have a DPL from 1 to 3. By creating gates
96  * with a DPL below user (ring 3) the int $n will generate a #gp fault
97  * which the hypervisor catches and forwards to the guest.
98  */
99 #if defined(__xpv)
100
101 #if defined(__amd64)
102 #define SEL_XPL 0 /* hypervisor privilege level */
103 #define SEL_KPL 3 /* both kernel and user in ring 3 */
104 #define TRP_KPL 1 /* system gate priv (user blocked) */
105 #define TRP_XPL 0 /* system gate priv (hypervisor) */
106
107 #define IST_DBG 0
108 #elif defined(__i386)
109 #define SEL_XPL 0 /* hypervisor privilege level */
110 #define SEL_KPL 1 /* kernel privilege level */
111 #define TRP_KPL SEL_KPL /* system gate priv (user blocked) */
112 #endif
113 #endif /* __i386 */
114
115 #define TRP_XPL 0 /* system gate priv (hypervisor) */
116
117 #else /* __xpv */
118 #define SEL_KPL 0 /* kernel privilege level on metal */
119 #define TRP_KPL SEL_KPL /* system gate priv (user blocked) */
120 #endif

```

```

114 #define IST_DF      1
115 #define IST_NMI     2
116 #define IST_MCE     3
117 #define IST_DBG     4
118 #define IST_NESTABLE 5
119 #define IST_DEFAULT 6

121 #endif /* __xpv */

123 #define IST_NONE     0

125 #define SEL_UPL      3 /* user priority level */
126 #define TRP_UPL     3 /* system gate priv (user allowed) */
127 #define SEL_TI_LDT  4 /* local descriptor table */
128 #define SEL_LDT(s)  (IDXTOSEL(s) | SEL_TI_LDT | SEL_UPL) /* local sel */
129 #define CPL_MASK    3 /* RPL mask for selector */
130 #define SELISLDT(s) (((s) & SEL_TI_LDT) == SEL_TI_LDT)
131 #define SELISUPL(s) (((s) & CPL_MASK) == SEL_UPL)

133 #ifndef _ASM

135 typedef uint16_t selector_t; /* selector register */

137 /*
138 * Hardware descriptor table register format for GDT and IDT.
139 */
140 #if defined(__amd64)

142 #pragma pack(2)
143 typedef struct {
144     uint16_t dtr_limit; /* table limit */
145     uint64_t dtr_base; /* table base address */
146 } desc_t;
147 #endif
148 #endif
149 #endif
150 #endif
151 #endif
152 #endif
153 #endif
154 #endif
155 #endif
156 #endif
157 #endif
158 #endif
159 #endif
160 #endif
161 #endif
162 #endif
163 #endif
164 #endif
165 #endif
166 #endif
167 #endif
168 #endif
169 #endif
170 #endif
171 #endif
172 #endif
173 #endif
174 #endif
175 #endif
176 #endif
177 #endif
178 #endif
179 #endif
180 #endif
181 #endif
182 #endif
183 #endif
184 #endif
185 #endif
186 #endif
187 #endif
188 #endif
189 #endif
190 #endif
191 #endif
192 #endif
193 #endif
194 #endif
195 #endif
196 #endif
197 #endif
198 #endif
199 #endif
200 #endif
201 #endif
202 #endif
203 #endif
204 #endif
205 #endif
206 #endif
207 #endif
208 #endif
209 #endif
210 #endif
211 #endif
212 #endif
213 #endif
214 #endif
215 #endif
216 #endif
217 #endif
218 #endif
219 #endif
220 #endif
221 #endif
222 #endif
223 #endif
224 #endif
225 #endif
226 #endif
227 #endif
228 #endif
229 #endif
230 #endif
231 #endif
232 #endif
233 #endif
234 #endif
235 #endif
236 #endif
237 #endif
238 #endif
239 #endif
240 #endif
241 #endif
242 #endif
243 #endif
244 #endif
245 #endif
246 #endif
247 #endif
248 #endif
249 #endif
250 #endif
251 #endif
252 #endif
253 #endif
254 #endif
255 #endif
256 #endif
257 #endif
258 #endif
259 #endif
260 #endif
261 #endif
262 #endif
263 #endif
264 #endif
265 #endif
266 #endif
267 #endif
268 #endif
269 #endif
270 #endif
271 #endif
272 #endif
273 #endif
274 #endif
275 #endif
276 #endif
277 #endif
278 #endif
279 #endif
280 #endif
281 #endif
282 #endif
283 #endif
284 #endif
285 #endif
286 #endif
287 #endif
288 #endif
289 #endif
290 #endif
291 #endif
292 #endif
293 #endif
294 #endif
295 #endif
296 #endif
297 #endif
298 #endif
299 #endif
300 #endif
301 #endif
302 #endif
303 #endif
304 #endif
305 #endif
306 #endif
307 #endif
308 #endif
309 #endif
310 #endif
311 #endif
312 #endif
313 #endif
314 #endif
315 #endif
316 #endif
317 #endif
318 #endif
319 #endif
320 #endif
321 #endif
322 #endif
323 #endif
324 #endif
325 #endif
326 #endif
327 #endif
328 #endif
329 #endif
330 #endif
331 #endif
332 #endif
333 #endif
334 #endif
335 #endif
336 #endif
337 #endif
338 #endif
339 #endif
340 #endif
341 #endif
342 #endif
343 #endif
344 #endif
345 #endif
346 #endif
347 #endif
348 #endif
349 #endif
350 #endif
351 #endif
352 #endif
353 #endif
354 #endif
355 #endif
356 #endif
357 #endif
358 #endif
359 #endif
360 #endif
361 #endif
362 #endif
363 #endif
364 #endif
365 #endif
366 #endif
367 #endif
368 #endif
369 #endif
370 #endif
371 #endif
372 #endif
373 #endif
374 #endif
375 #endif
376 #endif
377 #endif
378 #endif
379 #endif
380 #endif
381 #endif
382 #endif
383 #endif
384 #endif
385 #endif
386 #endif
387 #endif
388 #endif
389 #endif
390 #endif
391 #endif
392 #endif
393 #endif
394 #endif
395 #endif
396 #endif
397 #endif
398 #endif
399 #endif
400 #endif
401 #endif
402 #endif
403 #endif
404 #endif
405 #endif
406 #endif
407 #endif
408 #endif
409 #endif
410 #endif
411 #endif
412 #endif
413 #endif
414 #endif
415 #endif
416 #endif
417 #endif
418 #endif
419 #endif
420 #endif
421 #endif
422 #endif
423 #endif
424 #endif
425 #endif
426 #endif
427 #endif
428 #endif
429 #endif
430 #endif
431 #endif
432 #endif
433 #endif
434 #endif
435 #endif
436 #endif
437 #endif
438 #endif
439 #endif
440 #endif
441 #endif
442 #endif
443 #endif
444 #endif
445 #endif
446 #endif
447 #endif
448 #endif
449 #endif
450 #endif
451 #endif
452 #endif
453 #endif
454 #endif
455 #endif
456 #endif
457 #endif
458 #endif
459 #endif
460 #endif
461 #endif
462 #endif
463 #endif
464 #endif
465 #endif
466 #endif
467 #endif
468 #endif
469 #endif
470 #endif
471 #endif
472 #endif
473 #endif
474 #endif
475 #endif
476 #endif
477 #endif
478 #endif
479 #endif
480 #endif
481 #endif
482 #endif
483 #endif
484 #endif
485 #endif
486 #endif
487 #endif
488 #endif
489 #endif
490 #endif
491 #endif
492 #endif
493 #endif
494 #endif
495 #endif
496 #endif
497 #endif
498 #endif
499 #endif
500 #endif
501 #endif
502 #endif
503 #endif
504 #endif
505 #endif
506 #endif
507 #endif
508 #endif
509 #endif
510 #endif
511 #endif
512 #endif
513 #endif
514 #endif
515 #endif
516 #endif
517 #endif
518 #endif
519 #endif
520 #endif
521 #endif
522 #endif
523 #endif
524 #endif
525 #endif
526 #endif
527 #endif
528 #endif
529 #endif
530 #endif
531 #endif
532 #endif
533 #endif
534 #endif
535 #endif
536 #endif
537 #endif
538 #endif
539 #endif
540 #endif
541 #endif
542 #endif
543 #endif
544 #endif
545 #endif
546 #endif
547 #endif
548 #endif
549 #endif
550 #endif
551 #endif
552 #endif
553 #endif
554 #endif
555 #endif
556 #endif
557 #endif
558 #endif
559 #endif
560 #endif
561 #endif
562 #endif
563 #endif
564 #endif
565 #endif
566 #endif
567 #endif
568 #endif
569 #endif
570 #endif
571 #endif
572 #endif
573 #endif
574 #endif
575 #endif
576 #endif
577 #endif
578 #endif
579 #endif
580 #endif
581 #endif
582 #endif
583 #endif
584 #endif
585 #endif
586 #endif
587 #endif
588 #endif
589 #endif
590 #endif
591 #endif
592 #endif
593 #endif
594 #endif
595 #endif
596 #endif
597 #endif
598 #endif
599 #endif
600 #endif
601 #endif
602 #endif
603 #endif
604 #endif
605 #endif
606 #endif
607 #endif
608 #endif
609 #endif
610 #endif
611 #endif
612 #endif
613 #endif
614 #endif
615 #endif
616 #endif
617 #endif
618 #endif
619 #endif
620 #endif
621 #endif
622 #endif
623 #endif
624 #endif
625 #endif
626 #endif
627 #endif
628 #endif
629 #endif
630 #endif
631 #endif
632 #endif
633 #endif
634 #endif
635 #endif
636 #endif
637 #endif
638 #endif
639 #endif
640 #endif
641 #endif
642 #endif
643 #endif
644 #endif
645 #endif
646 #endif
647 #endif
648 #endif
649 #endif
650 #endif
651 #endif
652 #endif
653 #endif
654 #endif
655 #endif
656 #endif
657 #endif
658 #endif
659 #endif
660 #endif
661 #endif
662 #endif
663 #endif
664 #endif
665 #endif
666 #endif
667 #endif
668 #endif
669 #endif
670 #endif
671 #endif
672 #endif
673 #endif
674 #endif
675 #endif
676 #endif
677 #endif
678 #endif
679 #endif
680 #endif
681 #endif
682 #endif
683 #endif
684 #endif
685 #endif
686 #endif
687 #endif
688 #endif
689 #endif
690 #endif
691 #endif
692 #endif
693 #endif
694 #endif
695 #endif
696 #endif
697 #endif
698 #endif
699 #endif
700 #endif
701 #endif
702 #endif
703 #endif
704 #endif
705 #endif
706 #endif
707 #endif
708 #endif
709 #endif
710 #endif
711 #endif
712 #endif
713 #endif
714 #endif
715 #endif
716 #endif
717 #endif
718 #endif
719 #endif
720 #endif
721 #endif
722 #endif
723 #endif
724 #endif
725 #endif
726 #endif
727 #endif
728 #endif
729 #endif
730 #endif
731 #endif
732 #endif
733 #endif
734 #endif
735 #endif
736 #endif
737 #endif
738 #endif
739 #endif
740 #endif
741 #endif
742 #endif
743 #endif
744 #endif
745 #endif
746 #endif
747 #endif
748 #endif
749 #endif
750 #endif
751 #endif
752 #endif
753 #endif
754 #endif
755 #endif
756 #endif
757 #endif
758 #endif
759 #endif
760 #endif
761 #endif
762 #endif
763 #endif
764 #endif
765 #endif
766 #endif
767 #endif
768 #endif
769 #endif
770 #endif
771 #endif
772 #endif
773 #endif
774 #endif
775 #endif
776 #endif
777 #endif
778 #endif
779 #endif
780 #endif
781 #endif
782 #endif
783 #endif
784 #endif
785 #endif
786 #endif
787 #endif
788 #endif
789 #endif
790 #endif
791 #endif
792 #endif
793 #endif
794 #endif
795 #endif
796 #endif
797 #endif
798 #endif
799 #endif
800 #endif
801 #endif
802 #endif
803 #endif
804 #endif
805 #endif
806 #endif
807 #endif
808 #endif
809 #endif
810 #endif
811 #endif
812 #endif
813 #endif
814 #endif
815 #endif
816 #endif
817 #endif
818 #endif
819 #endif
820 #endif
821 #endif
822 #endif
823 #endif
824 #endif
825 #endif
826 #endif
827 #endif
828 #endif
829 #endif
830 #endif
831 #endif
832 #endif
833 #endif
834 #endif
835 #endif
836 #endif
837 #endif
838 #endif
839 #endif
840 #endif
841 #endif
842 #endif
843 #endif
844 #endif
845 #endif
846 #endif
847 #endif
848 #endif
849 #endif
850 #endif
851 #endif
852 #endif
853 #endif
854 #endif
855 #endif
856 #endif
857 #endif
858 #endif
859 #endif
860 #endif
861 #endif
862 #endif
863 #endif
864 #endif
865 #endif
866 #endif
867 #endif
868 #endif
869 #endif
870 #endif
871 #endif
872 #endif
873 #endif
874 #endif
875 #endif
876 #endif
877 #endif
878 #endif
879 #endif
880 #endif
881 #endif
882 #endif
883 #endif
884 #endif
885 #endif
886 #endif
887 #endif
888 #endif
889 #endif
890 #endif
891 #endif
892 #endif
893 #endif
894 #endif
895 #endif
896 #endif
897 #endif
898 #endif
899 #endif
900 #endif
901 #endif
902 #endif
903 #endif
904 #endif
905 #endif
906 #endif
907 #endif
908 #endif
909 #endif
910 #endif
911 #endif
912 #endif
913 #endif
914 #endif
915 #endif
916 #endif
917 #endif
918 #endif
919 #endif
920 #endif
921 #endif
922 #endif
923 #endif
924 #endif
925 #endif
926 #endif
927 #endif
928 #endif
929 #endif
930 #endif
931 #endif
932 #endif
933 #endif
934 #endif
935 #endif
936 #endif
937 #endif
938 #endif
939 #endif
940 #endif
941 #endif
942 #endif
943 #endif
944 #endif
945 #endif
946 #endif
947 #endif
948 #endif
949 #endif
950 #endif
951 #endif
952 #endif
953 #endif
954 #endif
955 #endif
956 #endif
957 #endif
958 #endif
959 #endif
960 #endif
961 #endif
962 #endif
963 #endif
964 #endif
965 #endif
966 #endif
967 #endif
968 #endif
969 #endif
970 #endif
971 #endif
972 #endif
973 #endif
974 #endif
975 #endif
976 #endif
977 #endif
978 #endif
979 #endif
980 #endif
981 #endif
982 #endif
983 #endif
984 #endif
985 #endif
986 #endif
987 #endif
988 #endif
989 #endif
990 #endif
991 #endif
992 #endif
993 #endif
994 #endif
995 #endif
996 #endif
997 #endif
998 #endif
999 #endif
1000 #endif

```

```

416 extern int ldt_update_seg(user_desc_t *, user_desc_t *);

418 #if defined(__xpv)

420 extern int xen_idt_to_trap_info(uint_t, gate_desc_t *, void *);
421 extern void xen_idt_write(gate_desc_t *, uint_t);

423 #endif /* __xen */

425 void init_boot_gdt(user_desc_t *);

427 #endif /* _ASM */

429 /*
430 * Common segment parameter definitions for granularity, default
431 * operand size and operaton mode.
432 */
433 #define SDP_BYTES 0 /* segment limit scaled to bytes */
434 #define SDP_PAGES 1 /* segment limit scaled to pages */
435 #define SDP_OP32 1 /* code and data default operand = 32 bits */
436 #define SDP_LONG 1 /* long mode code segment (64 bits) */
437 #define SDP_SHORT 0 /* compat/legacy code segment (32 bits) */
438 /*
439 * System segments and gate types.
440 *
441 * In long mode i386 32-bit ldt, tss, call, interrupt and trap gate
442 * types are redefined into 64-bit equivalents.
443 */
444 #define SDT_SYSNUL 0 /* system null */
445 #define SDT_SYS286TSS 1 /* system 286 TSS available */
446 #define SDT_SYSLDT 2 /* system local descriptor table */
447 #define SDT_SYS286BSY 3 /* system 286 TSS busy */
448 #define SDT_SYS286CGT 4 /* system 286 call gate */
449 #define SDT_SYSTASKGT 5 /* system task gate */
450 #define SDT_SYS286IGT 6 /* system 286 interrupt gate */
451 #define SDT_SYS286TGT 7 /* system 286 trap gate */
452 #define SDT_SYSNUL2 8 /* system null again */
453 #define SDT_SYSTSS 9 /* system TSS available */
454 #define SDT_SYSNUL3 10 /* system null again */
455 #define SDT_SYSTSSBSY 11 /* system TSS busy */
456 #define SDT_SYSCGT 12 /* system call gate */
457 #define SDT_SYSNUL4 13 /* system null again */
458 #define SDT_SYSIGT 14 /* system interrupt gate */
459 #define SDT_SYSTGT 15 /* system trap gate */

461 /*
462 * Memory segment types.
463 *
464 * While in long mode expand-down, writable and accessed type field
465 * attributes are ignored. Only the conforming bit is loaded by hardware
466 * for long mode code segment descriptors.
467 */
468 #define SDT_MEMRO 16 /* read only */
469 #define SDT_MEMROA 17 /* read only accessed */
470 #define SDT_MEMRW 18 /* read write */
471 #define SDT_MEMRWA 19 /* read write accessed */
472 #define SDT_MEMROD 20 /* read only expand dwn limit */
473 #define SDT_MEMRODA 21 /* read only expand dwn limit accessed */
474 #define SDT_MEMRWD 22 /* read write expand dwn limit */
475 #define SDT_MEMRWDA 23 /* read write expand dwn limit accessed */
476 #define SDT_MEME 24 /* execute only */
477 #define SDT_MEMEA 25 /* execute only accessed */
478 #define SDT_MEMER 26 /* execute read */
479 #define SDT_MEMERA 27 /* execute read accessed */
480 #define SDT_MEMEC 28 /* execute only conforming */

```

new/usr/src/uts/intel/sys/segments.h

```

481 #define SDT_MEMEAC 29 /* execute only accessed conforming */
482 #define SDT_MEMERC 30 /* execute read conforming */
483 #define SDT_MEMERAC 31 /* execute read accessed conforming */

485 /*
486 * Entries in the Interrupt Descriptor Table (IDT)
487 */
488 #define IDT_DE 0 /* #DE: Divide Error */
489 #define IDT_DB 1 /* #DB: Debug */
490 #define IDT_NMI 2 /* Nonmaskable External Interrupt */
491 #define IDT_BP 3 /* #BP: Breakpoint */
492 #define IDT_OF 4 /* #OF: Overflow */
493 #define IDT_BR 5 /* #BR: Bound Range Exceeded */
494 #define IDT_UD 6 /* #UD: Undefined/Invalid Opcode */
495 #define IDT_NM 7 /* #NM: No Math Coprocessor */
496 #define IDT_DF 8 /* #DF: Double Fault */
497 #define IDT_FPUGP 9 /* Coprocessor Segment Overrun */
498 #define IDT_TS 10 /* #TS: Invalid TSS */
499 #define IDT_NP 11 /* #NP: Segment Not Present */
500 #define IDT_SS 12 /* #SS: Stack Segment Fault */
501 #define IDT_GP 13 /* #GP: General Protection Fault */
502 #define IDT_PF 14 /* #PF: Page Fault */
503 #define IDT_MF 16 /* #MF: FPU Floating-Point Error */
504 #define IDT_AC 17 /* #AC: Alignment Check */
505 #define IDT_MC 18 /* #MC: Machine Check */
506 #define IDT_XF 19 /* #XF: SIMD Floating-Point Exception */
507 #define NIDT 256 /* size in entries of IDT */

509 /*
510 * Entries in the Global Descriptor Table (GDT)
511 *
512 * We make sure to space the system descriptors (LDT's, TSS')
513 * such that they are double gdt slot aligned. This is because
514 * in long mode system segment descriptors expand to 128 bits.
515 *
516 * GDT_LWPFS and GDT_LWPGS must be the same for both 32 and 64-bit
517 * kernels. See setup_context in libc. 64-bit processes must set
518 * %fs or %gs to null selector to use 64-bit fsbase or gsbase
519 * respectively.
520 */
521 #define GDT_NULL 0 /* null */
522 #define GDT_B32DATA 1 /* dboot 32 bit data descriptor */
523 #define GDT_B32CODE 2 /* dboot 32 bit code descriptor */
524 #define GDT_B16CODE 3 /* bios call 16 bit code descriptor */
525 #define GDT_B16DATA 4 /* bios call 16 bit data descriptor */
526 #define GDT_B64CODE 5 /* dboot 64 bit code descriptor */
527 #define GDT_BGSTMP 7 /* kmdb descriptor only used early in boot */
528 #define GDT_CPUID 16 /* store numeric id of current CPU */

530 #if defined(__amd64)

532 #define GDT_KCODE 6 /* kernel code seg %cs */
533 #define GDT_KDATA 7 /* kernel data seg %ds */
534 #define GDT_U32CODE 8 /* 32-bit process on 64-bit kernel %cs */
535 #define GDT_UDATA 9 /* user data seg %ds (32 and 64 bit) */
536 #define GDT_UCODE 10 /* native user code seg %cs */
537 #define GDT_LDT 12 /* (12-13) LDT for current process */
538 #define GDT_KTSS 14 /* (14-15) kernel tss */
539 #define GDT_FS GDT_NULL /* kernel %fs segment selector */
540 #define GDT_GS GDT_NULL /* kernel %gs segment selector */
541 #define GDT_LWPFS 55 /* lwp private %fs segment selector (32-bit) */
542 #define GDT_LWPGS 56 /* lwp private %gs segment selector (32-bit) */
543 #define GDT_BRANDMIN 57 /* first entry in GDT for brand usage */
544 #define GDT_BRANDMAX 61 /* last entry in GDT for brand usage */
545 #define NGDT 62 /* number of entries in GDT */

```

5

new/usr/src/uts/intel/sys/segments.h

```

547 /*
548 * This selector is only used in the temporary GDT used to bring additional
549 * CPUs from 16-bit real mode into long mode in real_mode_start().
550 */
551 #define TEMP_GDT_KCODE64 1 /* 64-bit code selector */

553 #elif defined(__i386)

555 #define GDT_LDT 40 /* LDT for current process */
556 #define GDT_KTSS 42 /* kernel tss */
557 #define GDT_KCODE 43 /* kernel code seg %cs */
558 #define GDT_KDATA 44 /* kernel data seg %ds */
559 #define GDT_UCODE 45 /* native user code seg %cs */
560 #define GDT_UDATA 46 /* user data seg %ds (32 and 64 bit) */
561 #define GDT_DBFLT 47 /* double fault #DF selector */
562 #define GDT_FS 53 /* kernel %fs segment selector */
563 #define GDT_GS 54 /* kernel %gs segment selector */
564 #define GDT_LWPFS 55 /* lwp private %fs segment selector */
565 #define GDT_LWPGS 56 /* lwp private %gs segment selector */
566 #define GDT_BRANDMIN 57 /* first entry in GDT for brand usage */
567 #define GDT_BRANDMAX 61 /* last entry in GDT for brand usage */
568 #if !defined(__xpv)
569 #define NGDT 90 /* number of entries in GDT */
570 #else
571 #define NGDT 512 /* single 4K page for the hypervisor */
572 #endif

574 #endif /* __i386 */

576 /*
577 * Convenient selector definitions.
578 */

580 /*
581 * XXPV 64 bit Xen only allows the guest %cs/%ss be the private ones it
582 * provides, not the ones we create for ourselves. See FLAT_RING3_CS64 in
583 * public/arch-x86_64.h
584 */
585 * 64-bit Xen runs paravirtual guests in ring 3 but emulates them running in
586 * ring 0 by clearing CPL in %cs value pushed on guest exception stacks.
587 * Therefore we will have KCS_SEL value indicate ring 0 and use that everywhere
588 * in the kernel. But in the few files where we initialize segment registers or
589 * create and update descriptors we will explicitly OR in SEL_KPL (ring 3) for
590 * kernel %cs. See desctbls.c for an example.
591 */

593 #if defined(__xpv) && defined(__amd64)
594 #define KCS_SEL 0xe030 /* FLAT_RING3_CS64 & 0xFFFF0 */
595 #define KDS_SEL 0xe02b /* FLAT_RING3_SS64 */
596 #else
597 #define KCS_SEL SEL_GDT(GDT_KCODE, SEL_KPL)
598 #define KDS_SEL SEL_GDT(GDT_KDATA, SEL_KPL)
599 #endif

601 #define UCS_SEL SEL_GDT(GDT_UCODE, SEL_UPL)
602 #if defined(__amd64)
603 #define TEMP_CS64_SEL SEL_GDT(TEMP_GDT_KCODE64, SEL_KPL)
604 #define U32CS_SEL SEL_GDT(GDT_U32CODE, SEL_UPL)
605 #endif

607 #define UDS_SEL SEL_GDT(GDT_UDATA, SEL_UPL)
608 #define ULDT_SEL SEL_GDT(GDT_LDT, SEL_KPL)
609 #define KTSS_SEL SEL_GDT(GDT_KTSS, SEL_KPL)
610 #define DFTSS_SEL SEL_GDT(GDT_DBFLT, SEL_KPL)
611 #define KFS_SEL 0
612 #define KGS_SEL SEL_GDT(GDT_GS, SEL_KPL)

```

6

```

613 #define LWPPFS_SEL      SEL_GDT(GDT_LWPPFS, SEL_UPL)
614 #define LWPGS_SEL      SEL_GDT(GDT_LWPGS, SEL_UPL)
615 #define BRANDMIN_SEL    SEL_GDT(GDT_BRANDMIN, SEL_UPL)
616 #define BRANDMAX_SEL    SEL_GDT(GDT_BRANDMAX, SEL_UPL)

618 #define B64CODE_SEL     SEL_GDT(GDT_B64CODE, SEL_KPL)
619 #define B32CODE_SEL     SEL_GDT(GDT_B32CODE, SEL_KPL)
620 #define B32DATA_SEL     SEL_GDT(GDT_B32DATA, SEL_KPL)
621 #define B16CODE_SEL     SEL_GDT(GDT_B16CODE, SEL_KPL)
622 #define B16DATA_SEL     SEL_GDT(GDT_B16DATA, SEL_KPL)

624 /*
625  * Temporary %gs descriptor used by kmdb with -d option. Only lives
626  * in boot's GDT and is not copied into kernel's GDT from boot.
627  */
628 #define KMDBGS_SEL      SEL_GDT(GDT_BGSTMP, SEL_KPL)

630 /*
631  * Selector used for kdi_idt when kmdb has taken over the IDT.
632  */
633 #if defined(__amd64)
634 #define KMDBCODE_SEL    B64CODE_SEL
635 #else
636 #define KMDBCODE_SEL    B32CODE_SEL
637 #endif

639 /*
640  * Entries in default Local Descriptor Table (LDT) for every process.
641  */
642 #define LDT_SYSCALL      0      /* call gate for libc.a (obsolete) */
643 #define LDT_SIGCALL      1      /* EOL me, call gate for static sigreturn */
644 #define LDT_RESVD1       2      /* old user %cs */
645 #define LDT_RESVD2       3      /* old user %ds */
646 #define LDT_ALTSYSCALL  4      /* alternate call gate for system calls */
647 #define LDT_ALTSIGCALL  5      /* EOL me, alternate call gate for sigreturn */
648 #define LDT_UDBASE       6      /* user descriptor base index */
649 #define MINNLDT         512    /* Current min solaris ldt size (1 4K page) */
650 #define MAXNLDT         8192   /* max solaris ldt size (16 4K pages) */

652 #ifndef _KERNEL
653 #define LDT_CPU_SIZE     (16 * 4096) /* Size of kernel per-CPU allocation */
654 #endif

656 #ifndef _ASM

658 extern gate_desc_t      *idt0;
659 extern desc_t          idt0_default_reg;
660 extern user_desc_t     *gdt0;

662 extern user_desc_t     zero_udesc;
663 extern user_desc_t     null_udesc;
664 extern system_desc_t   null_sdesc;

666 #if defined(__amd64)
667 extern user_desc_t     zero_u32desc;
668 #endif
669 #if defined(__amd64)
670 extern user_desc_t     ucs_on;
671 extern user_desc_t     ucs_off;
672 extern user_desc_t     ucs32_on;
673 extern user_desc_t     ucs32_off;
674 #endif /* __amd64 */

676 extern tss_t *ktss0;

678 #if defined(__i386)

```

```

679 extern tss_t *dftss0;
680 #endif /* __i386 */

682 extern void div0trap(), dbgtrap(), nmiint(), brktrap(), ovflotrap();
683 extern void boundstrap(), invoptrap(), ndptrap();
684 #if !defined(__xpv)
685 extern void syserrtrap();
686 #endif
687 extern void invaltrap(), invtsstrap(), segnptrap(), stktrap();
688 extern void gptrap(), pfttrap(), ndperr();
689 extern void overrun(), resvtrap();
690 extern void _start(), cmnint();
691 extern void achktrap(), mcetrap();
692 extern void xmttrap();
693 extern void fasttrap();
694 extern void dtrace_ret();

696 /* KPTI trampolines */
697 extern void tr_invaltrap();
698 extern void tr_div0trap(), tr_dbgtrap(), tr_nmiint(), tr_brktrap();
699 extern void tr_ovflotrap(), tr_boundstrap(), tr_invoptrap(), tr_ndptrap();
700 #if !defined(__xpv)
701 extern void tr_syserrtrap();
702 #endif
703 extern void tr_invaltrap(), tr_invtsstrap(), tr_segnptrap(), tr_stktrap();
704 extern void tr_gptrap(), tr_pfttrap(), tr_ndperr();
705 extern void tr_overrun(), tr_resvtrap();
706 extern void tr_achktrap(), tr_mcetrap();
707 extern void tr_xmttrap();
708 extern void tr_fasttrap();
709 extern void tr_dtrace_ret();

711 #if !defined(__amd64)
712 extern void pentium_pfttrap();
713 #endif

715 extern uint64_t kpti_enable;

717 #endif /* _ASM */

719 #ifdef __cplusplus
720 }
    unchanged_portion_omitted

```

new/usr/src/uts/intel/sys/x86_archext.h

1

```
*****
32559 Fri Apr 6 17:25:14 2018
new/usr/src/uts/intel/sys/x86_archext.h
8956 Implement KPTI
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: Robert Mustacchi <rm@joyent.com>
9215 update CPUID defines
Reviewed by: Yuri Pankov <yuripv@yuripv.net>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22  * Copyright (c) 1995, 2010, Oracle and/or its affiliates. All rights reserved.
23  * Copyright (c) 2011 by Delphix. All rights reserved.
24 */
25 /*
26  * Copyright (c) 2010, Intel Corporation.
27  * All rights reserved.
28 */
29 /*
30  * Copyright 2018 Joyent, Inc.
30  * Copyright 2017 Joyent, Inc.
31  * Copyright 2012 Jens Elkner <jel+illumos@cs.uni-magdeburg.de>
32  * Copyright 2012 Hans Rosenfeld <rosenfeld@grumpf.hope-2000.org>
33  * Copyright 2014 Josef 'Jeff' Sipek <jeffpc@josefsipek.net>
34  * Copyright 2018 Nexenta Systems, Inc.
35 */
37 #ifndef _SYS_X86_ARCHEXT_H
38 #define _SYS_X86_ARCHEXT_H
39
40 #if !defined(_ASM)
41 #include <sys/regset.h>
42 #include <sys/processor.h>
43 #include <vm/seg_enum.h>
44 #include <vm/page.h>
45 #endif /* _ASM */
46
47 #ifdef __cplusplus
48 extern "C" {
49 #endif
50
51 /*
52  * cpuid instruction feature flags in %edx (standard function 1)
53 */
54
55 #define CPUID_INTC_EDX_FPU 0x00000001 /* x87 fpu present */
56 #define CPUID_INTC_EDX_VME 0x00000002 /* virtual-8086 extension */
```

new/usr/src/uts/intel/sys/x86_archext.h

2

```
57 #define CPUID_INTC_EDX_DE 0x00000004 /* debugging extensions */
58 #define CPUID_INTC_EDX_PSE 0x00000008 /* page size extension */
59 #define CPUID_INTC_EDX_TSC 0x00000010 /* time stamp counter */
60 #define CPUID_INTC_EDX_MSR 0x00000020 /* rdmsr and wrmsr */
61 #define CPUID_INTC_EDX_PAE 0x00000040 /* physical addr extension */
62 #define CPUID_INTC_EDX_MCE 0x00000080 /* machine check exception */
63 #define CPUID_INTC_EDX_CX8 0x00000100 /* cmpxchg8b instruction */
64 #define CPUID_INTC_EDX_APIC 0x00000200 /* local APIC */
65 /* 0x400 - reserved */
66 #define CPUID_INTC_EDX_SEP 0x00000800 /* sysenter and sysexit */
67 #define CPUID_INTC_EDX_MTRR 0x00001000 /* memory type range reg */
68 #define CPUID_INTC_EDX_PGE 0x00002000 /* page global enable */
69 #define CPUID_INTC_EDX_MCA 0x00004000 /* machine check arch */
70 #define CPUID_INTC_EDX_CMOV 0x00008000 /* conditional move insns */
71 #define CPUID_INTC_EDX_PAT 0x00010000 /* page attribute table */
72 #define CPUID_INTC_EDX_PSE36 0x00020000 /* 36-bit pagesize extension */
73 #define CPUID_INTC_EDX_PSN 0x00040000 /* processor serial number */
74 #define CPUID_INTC_EDX_CLFSH 0x00080000 /* clflush instruction */
75 /* 0x10000 - reserved */
76 #define CPUID_INTC_EDX_DS 0x00200000 /* debug store exists */
77 #define CPUID_INTC_EDX_ACPI 0x00400000 /* monitoring + clock ctrl */
78 #define CPUID_INTC_EDX_MMX 0x00800000 /* MMX instructions */
79 #define CPUID_INTC_EDX_FXSR 0x01000000 /* fxsave and fxrstor */
80 #define CPUID_INTC_EDX_SSE 0x02000000 /* streaming SIMD extensions */
81 #define CPUID_INTC_EDX_SSE2 0x04000000 /* SSE extensions */
82 #define CPUID_INTC_EDX_SS 0x08000000 /* self-snoop */
83 #define CPUID_INTC_EDX_HTT 0x10000000 /* Hyper Thread Technology */
84 #define CPUID_INTC_EDX_TM 0x20000000 /* thermal monitoring */
85 #define CPUID_INTC_EDX_IA64 0x40000000 /* Itanium emulating IA32 */
86 #define CPUID_INTC_EDX_PBE 0x80000000 /* Pending Break Enable */
87
88 /*
89  * cpuid instruction feature flags in %ecx (standard function 1)
90 */
91
92 #define CPUID_INTC_ECX_SSE3 0x00000001 /* Yet more SSE extensions */
93 #define CPUID_INTC_ECX_PCLMULQDQ 0x00000002 /* PCLMULQDQ insn */
94 #define CPUID_INTC_ECX_DTES64 0x00000004 /* 64-bit DS area */
95 /* 0x00000004 - reserved */
96 #define CPUID_INTC_ECX_MON 0x00000008 /* MONITOR/MWAIT */
97 #define CPUID_INTC_ECX_DSCPL 0x00000010 /* CPL-qualified debug store */
98 #define CPUID_INTC_ECX_VMX 0x00000020 /* Hardware VM extensions */
99 #define CPUID_INTC_ECX_SMX 0x00000040 /* Secure mode extensions */
100 #define CPUID_INTC_ECX_EST 0x00000080 /* enhanced SpeedStep */
101 #define CPUID_INTC_ECX_TM2 0x00000100 /* thermal monitoring */
102 #define CPUID_INTC_ECX_SSSE3 0x00000200 /* Supplemental SSE3 insns */
103 #define CPUID_INTC_ECX_CID 0x00000400 /* L1 context ID */
104 /* 0x00000800 - reserved */
105 #define CPUID_INTC_ECX_FMA 0x00001000 /* Fused Multiply Add */
106 #define CPUID_INTC_ECX_CX16 0x00002000 /* cmpxchg16 */
107 #define CPUID_INTC_ECX_ETPRD 0x00004000 /* extended task pri messages */
108 #define CPUID_INTC_ECX_PDCM 0x00008000 /* Perf/Debug Capability MSR */
109 /* 0x00008000 - reserved */
110 /* 0x00010000 - reserved */
111 #define CPUID_INTC_ECX_PCID 0x00020000 /* process-context ids */
112 /* 0x00020000 - reserved */
113 #define CPUID_INTC_ECX_DCA 0x00040000 /* direct cache access */
114 #define CPUID_INTC_ECX_SSE4_1 0x00080000 /* SSE4.1 insns */
115 #define CPUID_INTC_ECX_SSE4_2 0x00100000 /* SSE4.2 insns */
116 #define CPUID_INTC_ECX_X2APIC 0x00200000 /* x2APIC */
117 #define CPUID_INTC_ECX_MOVBE 0x00400000 /* MOVBE insn */
118 #define CPUID_INTC_ECX_POPCNT 0x00800000 /* POPCNT insn */
119 #define CPUID_INTC_ECX_TSCDL 0x01000000 /* Deadline TSC */
120 #define CPUID_INTC_ECX_AES 0x02000000 /* AES insns */
121 #define CPUID_INTC_ECX_XSAVE 0x04000000 /* XSAVE/XRESTOR insns */
122 #define CPUID_INTC_ECX_OSXSAVE 0x08000000 /* OS supports XSAVE insns */
```



```

120 #define CPUID_INTC_ECX_AVX      0x10000000    /* AVX supported */
121 #define CPUID_INTC_ECX_FL16C    0x20000000    /* FL16C supported */
122 #define CPUID_INTC_ECX_RDRAND   0x40000000    /* RDRAND supported */
123 #define CPUID_INTC_ECX_HV      0x80000000    /* Hypervisor */

125 /*
126 * cpuid instruction feature flags in %edx (extended function 0x80000001)
127 */

129 #define CPUID_AMD_EDX_FPU      0x00000001    /* x87 fpu present */
130 #define CPUID_AMD_EDX_VME      0x00000002    /* virtual-8086 extension */
131 #define CPUID_AMD_EDX_DE       0x00000004    /* debugging extensions */
132 #define CPUID_AMD_EDX_PSE      0x00000008    /* page size extensions */
133 #define CPUID_AMD_EDX_TSC      0x00000010    /* time stamp counter */
134 #define CPUID_AMD_EDX_MSR      0x00000020    /* rdmsr and wrmsr */
135 #define CPUID_AMD_EDX_PAE      0x00000040    /* physical addr extension */
136 #define CPUID_AMD_EDX_MCE      0x00000080    /* machine check exception */
137 #define CPUID_AMD_EDX_CX8      0x00000100    /* cmpxchg8b instruction */
138 #define CPUID_AMD_EDX_APIC     0x00000200    /* local APIC */
139                                     /* 0x00000400 - sysc on K6m6 */
140 #define CPUID_AMD_EDX_SYSC     0x00000800    /* AMD: syscall and sysret */
141 #define CPUID_AMD_EDX_MTRR     0x00001000    /* memory type and range reg */
142 #define CPUID_AMD_EDX_PGE      0x00002000    /* page global enable */
143 #define CPUID_AMD_EDX_MCA      0x00004000    /* machine check arch */
144 #define CPUID_AMD_EDX_CMOV     0x00008000    /* conditional move insns */
145 #define CPUID_AMD_EDX_PAT      0x00010000    /* K7: page attribute table */
146 #define CPUID_AMD_EDX_FCMOV    0x00010000    /* FCMOVcc etc. */
147 #define CPUID_AMD_EDX_PSE36    0x00020000    /* 36-bit pagesize extension */
148                                     /* 0x00040000 - reserved */
149                                     /* 0x00080000 - reserved */
150 #define CPUID_AMD_EDX_NX       0x00100000    /* AMD: no-execute page prot */
151                                     /* 0x00200000 - reserved */
152 #define CPUID_AMD_EDX_MMXamd   0x00400000    /* AMD: MMX extensions */
153 #define CPUID_AMD_EDX_MMX      0x00800000    /* MMX instructions */
154 #define CPUID_AMD_EDX_FXSR     0x01000000    /* fxsave and fxrstor */
155 #define CPUID_AMD_EDX_FFXSR    0x02000000    /* fast fxsave/fxrstor */
156 #define CPUID_AMD_EDX_LGPG     0x04000000    /* 1GB page */
157 #define CPUID_AMD_EDX_TSCP     0x08000000    /* rdtscp instruction */
158                                     /* 0x10000000 - reserved */
159 #define CPUID_AMD_EDX_LM       0x20000000    /* AMD: long mode */
160 #define CPUID_AMD_EDX_3DNowx   0x40000000    /* AMD: extensions to 3DNow! */
161 #define CPUID_AMD_EDX_3DNow    0x80000000    /* AMD: 3DNow! instructions */

163 #define CPUID_AMD_ECX_AHF64    0x00000001    /* LAHF and SAHF in long mode */
164 #define CPUID_AMD_ECX_CMP_LGCV 0x00000002    /* AMD: multicore chip */
165 #define CPUID_AMD_ECX_SVM      0x00000004    /* AMD: secure VM */
166 #define CPUID_AMD_ECX_EAS      0x00000008    /* extended apic space */
167 #define CPUID_AMD_ECX_CR8D     0x00000010    /* AMD: 32-bit mov %cr8 */
168 #define CPUID_AMD_ECX_LZCNT    0x00000020    /* AMD: LZCNT insn */
169 #define CPUID_AMD_ECX_SSE4A    0x00000040    /* AMD: SSE4A insns */
170 #define CPUID_AMD_ECX_MAS      0x00000080    /* AMD: MisAlignSse mmode */
171 #define CPUID_AMD_ECX_3DNP     0x00000100    /* AMD: 3DNowPrefetch */
172 #define CPUID_AMD_ECX_OSVW     0x00000200    /* AMD: OSVW */
173 #define CPUID_AMD_ECX_IBS      0x00000400    /* AMD: IBS */
174 #define CPUID_AMD_ECX_SSE5     0x00000800    /* AMD: Extended AVX */
175 #define CPUID_AMD_ECX_SSE5     0x00000800    /* AMD: SSE5 */
176 #define CPUID_AMD_ECX_SKINIT   0x00001000    /* AMD: SKINIT */
177 #define CPUID_AMD_ECX_WDT      0x00002000    /* AMD: WDT */
178                                     /* 0x00004000 - reserved */
179 #define CPUID_AMD_ECX_LWP      0x00008000    /* AMD: Lightweight profiling */
180 #define CPUID_AMD_ECX_FMA4     0x00010000    /* AMD: 4-operand FMA support */
181                                     /* 0x00020000 - reserved */
182                                     /* 0x00040000 - reserved */
183 #define CPUID_AMD_ECX_NIDMSR   0x00080000    /* AMD: Node ID MSR */
184 #define CPUID_AMD_ECX_TBM      0x00100000    /* AMD: trailing bit manips. */

```

```

185 #define CPUID_AMD_ECX_TOPOEXT  0x00400000    /* AMD: Topology Extensions */

187 /*
188 * AMD uses %ebx for some of their features (extended function 0x80000008).
189 */
190 #define CPUID_AMD_EBX_ERR_PTR_ZERO 0x00000004 /* AMD: FP Err. Ptr. Zero */

192 /*
193 * Intel now seems to have claimed part of the "extended" function
194 * space that we previously for non-Intel implementors to use.
195 * More excitingly still, they've claimed bit 20 to mean LAHF/SAHF
196 * is available in long mode i.e. what AMD indicate using bit 0.
197 * On the other hand, everything else is labelled as reserved.
198 */
199 #define CPUID_INTC_ECX_AHF64    0x00100000    /* LAHF and SAHF in long mode */

201 /*
202 * Intel also uses cpuid leaf 7 to have additional instructions and features.
203 * Like some other leaves, but unlike the current ones we care about, it
204 * requires us to specify both a leaf in %eax and a sub-leaf in %ecx. To deal
205 * with the potential use of additional sub-leaves in the future, we now
206 * specifically label the EBX features with their leaf and sub-leaf.
207 */
208 #define CPUID_INTC_EBX_7_0_BMI1 0x00000008    /* BMI1 instrs */
209 #define CPUID_INTC_EBX_7_0_HLE  0x00000010    /* HLE */
210 #define CPUID_INTC_EBX_7_0_AVX2 0x00000020    /* AVX2 supported */
211 #define CPUID_INTC_EBX_7_0_SMEP 0x00000080    /* SMEP in CR4 */
212 #define CPUID_INTC_EBX_7_0_BMI2 0x00000100    /* BMI2 instrs */
213 #define CPUID_INTC_EBX_7_0_INVPCID 0x00000400 /* invpcid instr */
214 #define CPUID_INTC_EBX_7_0_MPX  0x00004000    /* Mem. Prot. Ext. */
215 #define CPUID_INTC_EBX_7_0_AVX512F 0x00010000 /* AVX512 foundation */
216 #define CPUID_INTC_EBX_7_0_AVX512DQ 0x00020000 /* AVX512DQ */
217 #define CPUID_INTC_EBX_7_0_RDSEED 0x00040000 /* RDSEED instr */
218 #define CPUID_INTC_EBX_7_0_ADX   0x00080000    /* ADX instrs */
219 #define CPUID_INTC_EBX_7_0_SMAP  0x01000000    /* SMAP in CR 4 */
220 #define CPUID_INTC_EBX_7_0_AVX512IFMA 0x00200000 /* AVX512IFMA */
221 #define CPUID_INTC_EBX_7_0_CLWB  0x01000000    /* CLWB */
222 #define CPUID_INTC_EBX_7_0_AVX512PF 0x04000000 /* AVX512PF */
223 #define CPUID_INTC_EBX_7_0_AVX512ER 0x08000000 /* AVX512ER */
224 #define CPUID_INTC_EBX_7_0_AVX512CD 0x10000000 /* AVX512CD */
225 #define CPUID_INTC_EBX_7_0_SHA   0x20000000    /* SHA extensions */
226 #define CPUID_INTC_EBX_7_0_AVX512BW 0x40000000 /* AVX512BW */
227 #define CPUID_INTC_EBX_7_0_AVX512VL 0x80000000 /* AVX512VL */

229 #define CPUID_INTC_EBX_7_0_ALL_AVX512 \
230 (CPUID_INTC_EBX_7_0_AVX512F | CPUID_INTC_EBX_7_0_AVX512DQ | \
231 CPUID_INTC_EBX_7_0_AVX512IFMA | CPUID_INTC_EBX_7_0_AVX512PF | \
232 CPUID_INTC_EBX_7_0_AVX512ER | CPUID_INTC_EBX_7_0_AVX512CD | \
233 CPUID_INTC_EBX_7_0_AVX512BW | CPUID_INTC_EBX_7_0_AVX512VL)

235 #define CPUID_INTC_ECX_7_0_AVX512VBMI 0x00000002 /* AVX512VBMI */
236 #define CPUID_INTC_ECX_7_0_UMIP      0x00000004 /* UMIP */
237 #define CPUID_INTC_ECX_7_0_PKU       0x00000008 /* umode prot. keys */
238 #define CPUID_INTC_ECX_7_0_OSPKE     0x00000010 /* OSPKE */
239 #define CPUID_INTC_ECX_7_0_AVX512VPOPCDQ 0x00004000 /* AVX512 VPOPCNTDQ */

241 #define CPUID_INTC_ECX_7_0_ALL_AVX512 \
242 (CPUID_INTC_ECX_7_0_AVX512VBMI | CPUID_INTC_ECX_7_0_AVX512VPOPCDQ)

244 #define CPUID_INTC_EDX_7_0_AVX5124NNIW 0x00000004 /* AVX512 4NNIW */
245 #define CPUID_INTC_EDX_7_0_AVX5124FMAPS 0x00000008 /* AVX512 4FMAPS */

247 #define CPUID_INTC_EDX_7_0_ALL_AVX512 \
248 (CPUID_INTC_EDX_7_0_AVX5124NNIW | CPUID_INTC_EDX_7_0_AVX5124FMAPS)

250 /*

```

```

251 * Intel also uses cpuid leaf 0xd to report additional instructions and features
252 * when the sub-leaf in %ecx == 1. We label these using the same convention as
253 * with leaf 7.
254 */
255 #define CPUID_INTC_EAX_D_1_XSAVEOPT 0x00000001 /* xsaveopt inst. */
256 #define CPUID_INTC_EAX_D_1_XSAVEC 0x00000002 /* xsavec inst. */
257 #define CPUID_INTC_EAX_D_1_XSAVES 0x00000008 /* xsaves inst. */

259 #define REG_PAT 0x277
260 #define REG_TSC 0x10 /* timestamp counter */
261 #define REG_APIC_BASE_MSR 0x1b
262 #define REG_X2APIC_BASE_MSR 0x800 /* The MSR address offset of x2APIC */

264 #if !defined(__xpv)
265 /*
266 * AMD ClE
267 */
268 #define MSR_AMD_INT_PENDING_CMP_HALT 0xC0010055
269 #define AMD_ACTONCMPHALT_SHIFT 27
270 #define AMD_ACTONCMPHALT_MASK 3
271 #endif

273 #define MSR_DEBUGCTL 0x1d9
275 #define DEBUGCTL_LBR 0x01
276 #define DEBUGCTL_BTF 0x02

278 /* Intel P6, AMD */
279 #define MSR_LBR_FROM 0x1db
280 #define MSR_LBR_TO 0x1dc
281 #define MSR_LEX_FROM 0x1dd
282 #define MSR_LEX_TO 0x1de

284 /* Intel P4 (pre-Prescott, non P4 M) */
285 #define MSR_P4_LBSTK_TOS 0x1da
286 #define MSR_P4_LBSTK_0 0x1db
287 #define MSR_P4_LBSTK_1 0x1dc
288 #define MSR_P4_LBSTK_2 0x1dd
289 #define MSR_P4_LBSTK_3 0x1de

291 /* Intel Pentium M */
292 #define MSR_P6M_LBSTK_TOS 0x1c9
293 #define MSR_P6M_LBSTK_0 0x040
294 #define MSR_P6M_LBSTK_1 0x041
295 #define MSR_P6M_LBSTK_2 0x042
296 #define MSR_P6M_LBSTK_3 0x043
297 #define MSR_P6M_LBSTK_4 0x044
298 #define MSR_P6M_LBSTK_5 0x045
299 #define MSR_P6M_LBSTK_6 0x046
300 #define MSR_P6M_LBSTK_7 0x047

302 /* Intel P4 (Prescott) */
303 #define MSR_PR4_LBSTK_TOS 0x1da
304 #define MSR_PR4_LBSTK_FROM_0 0x680
305 #define MSR_PR4_LBSTK_FROM_1 0x681
306 #define MSR_PR4_LBSTK_FROM_2 0x682
307 #define MSR_PR4_LBSTK_FROM_3 0x683
308 #define MSR_PR4_LBSTK_FROM_4 0x684
309 #define MSR_PR4_LBSTK_FROM_5 0x685
310 #define MSR_PR4_LBSTK_FROM_6 0x686
311 #define MSR_PR4_LBSTK_FROM_7 0x687
312 #define MSR_PR4_LBSTK_FROM_8 0x688
313 #define MSR_PR4_LBSTK_FROM_9 0x689
314 #define MSR_PR4_LBSTK_FROM_10 0x68a
315 #define MSR_PR4_LBSTK_FROM_11 0x68b
316 #define MSR_PR4_LBSTK_FROM_12 0x68c

```

```

317 #define MSR_PR4_LBSTK_FROM_13 0x68d
318 #define MSR_PR4_LBSTK_FROM_14 0x68e
319 #define MSR_PR4_LBSTK_FROM_15 0x68f
320 #define MSR_PR4_LBSTK_TO_0 0x6c0
321 #define MSR_PR4_LBSTK_TO_1 0x6c1
322 #define MSR_PR4_LBSTK_TO_2 0x6c2
323 #define MSR_PR4_LBSTK_TO_3 0x6c3
324 #define MSR_PR4_LBSTK_TO_4 0x6c4
325 #define MSR_PR4_LBSTK_TO_5 0x6c5
326 #define MSR_PR4_LBSTK_TO_6 0x6c6
327 #define MSR_PR4_LBSTK_TO_7 0x6c7
328 #define MSR_PR4_LBSTK_TO_8 0x6c8
329 #define MSR_PR4_LBSTK_TO_9 0x6c9
330 #define MSR_PR4_LBSTK_TO_10 0x6ca
331 #define MSR_PR4_LBSTK_TO_11 0x6cb
332 #define MSR_PR4_LBSTK_TO_12 0x6cc
333 #define MSR_PR4_LBSTK_TO_13 0x6cd
334 #define MSR_PR4_LBSTK_TO_14 0x6ce
335 #define MSR_PR4_LBSTK_TO_15 0x6cf

337 #define MCI_CTL_VALUE 0xffffffff

339 #define MTRR_TYPE_UC 0
340 #define MTRR_TYPE_WC 1
341 #define MTRR_TYPE_WT 4
342 #define MTRR_TYPE_WP 5
343 #define MTRR_TYPE_WB 6
344 #define MTRR_TYPE_UC_ 7

346 /*
347 * For Solaris we set up the page attribute table in the following way:
348 * PAT0 Write-Back
349 * PAT1 Write-Through
350 * PAT2 Uncacheable-
351 * PAT3 Uncacheable
352 * PAT4 Write-Back
353 * PAT5 Write-Through
354 * PAT6 Write-Combine
355 * PAT7 Uncacheable
356 * The only difference from h/w default is entry 6.
357 */
358 #define PAT_DEFAULT_ATTRIBUTE \
359 ((uint64_t)MTRR_TYPE_WB | \
360 ((uint64_t)MTRR_TYPE_WT << 8) | \
361 ((uint64_t)MTRR_TYPE_UC_ << 16) | \
362 ((uint64_t)MTRR_TYPE_UC << 24) | \
363 ((uint64_t)MTRR_TYPE_WB << 32) | \
364 ((uint64_t)MTRR_TYPE_WT << 40) | \
365 ((uint64_t)MTRR_TYPE_WC << 48) | \
366 ((uint64_t)MTRR_TYPE_UC << 56))

368 #define X86FSET_LARGE PAGE 0
369 #define X86FSET_TSC 1
370 #define X86FSET_MSR 2
371 #define X86FSET_MTRR 3
372 #define X86FSET_PGE 4
373 #define X86FSET_DE 5
374 #define X86FSET_CMOV 6
375 #define X86FSET_MMX 7
376 #define X86FSET_MCA 8
377 #define X86FSET_PAE 9
378 #define X86FSET_CX8 10
379 #define X86FSET_PAT 11
380 #define X86FSET_SEP 12
381 #define X86FSET_SSE 13
382 #define X86FSET_SSE2 14

```

```

383 #define X86FSET_HTT          15
384 #define X86FSET_ASYSC       16
385 #define X86FSET_NX          17
386 #define X86FSET_SSE3       18
387 #define X86FSET_CX16       19
388 #define X86FSET_CMP         20
389 #define X86FSET_TSCP        21
390 #define X86FSET_MWAIT       22
391 #define X86FSET_SSE4A       23
392 #define X86FSET_CPUID       24
393 #define X86FSET_SSSE3       25
394 #define X86FSET_SSE4_1      26
395 #define X86FSET_SSE4_2      27
396 #define X86FSET_1GPG        28
397 #define X86FSET_CLFSH       29
398 #define X86FSET_64          30
399 #define X86FSET_AES          31
400 #define X86FSET_PCLMULQDQ   32
401 #define X86FSET_XSAVE       33
402 #define X86FSET_AVX         34
403 #define X86FSET_VMX         35
404 #define X86FSET_SVM         36
405 #define X86FSET_TOPOEXT     37
406 #define X86FSET_F16C        38
407 #define X86FSET_RDRAND      39
408 #define X86FSET_X2APIC     40
409 #define X86FSET_AVX2        41
410 #define X86FSET_BMI1       42
411 #define X86FSET_BMI2       43
412 #define X86FSET_FMA         44
413 #define X86FSET_SMEP        45
414 #define X86FSET_SMAP        46
415 #define X86FSET_ADX         47
416 #define X86FSET_RDSEED     48
417 #define X86FSET_MPX        49
418 #define X86FSET_AVX512F     50
419 #define X86FSET_AVX512DQ    51
420 #define X86FSET_AVX512PF    52
421 #define X86FSET_AVX512ER    53
422 #define X86FSET_AVX512CD    54
423 #define X86FSET_AVX512BW    55
424 #define X86FSET_AVX512VL    56
425 #define X86FSET_AVX512FMA   57
426 #define X86FSET_AVX512VBMI  58
427 #define X86FSET_AVX512VPOPCDQ 59
428 #define X86FSET_AVX512NNIW  60
429 #define X86FSET_AVX512FMAPS 61
430 #define X86FSET_XSAVEOPT    62
431 #define X86FSET_XSAVEC      63
432 #define X86FSET_XSAVES      64
433 #define X86FSET_SHA         65
434 #define X86FSET_UMIP        66
435 #define X86FSET_PKU         67
436 #define X86FSET_OSPKE       68
437 #define X86FSET_PCID        69
438 #define X86FSET_INVPCID     70

440 /*
441  * Intel Deep C-State invariant TSC in leaf 0x80000007.
442  */
443 #define CPUID_TSC_CSTATE_INVARIANCE    (0x100)

445 /*
446  * Intel Deep C-state always-running local APIC timer
447  */
448 #define CPUID_CSTATE_ARAT              (0x4)

```

```

450 /*
451  * Intel ENERGY_PERF_BIAS MSR indicated by feature bit CPUID.6.ECX[3].
452  */
453 #define CPUID_EPB_SUPPORT              (1 << 3)

455 /*
456  * Intel TSC deadline timer
457  */
458 #define CPUID_DEADLINE_TSC            (1 << 24)

460 /*
461  * x86_type is a legacy concept; this is supplanted
462  * for most purposes by x86_featureset; modern CPUs
463  * should be X86_TYPE_OTHER
464  */
465 #define X86_TYPE_OTHER                 0
466 #define X86_TYPE_486                   1
467 #define X86_TYPE_P5                     2
468 #define X86_TYPE_P6                     3
469 #define X86_TYPE_CYRIX_486              4
470 #define X86_TYPE_CYRIX_6x86L           5
471 #define X86_TYPE_CYRIX_6x86            6
472 #define X86_TYPE_CYRIX_GXm             7
473 #define X86_TYPE_CYRIX_6x86MX          8
474 #define X86_TYPE_CYRIX_MediaGX         9
475 #define X86_TYPE_CYRIX_MII             10
476 #define X86_TYPE_VIA_CYRIX_III         11
477 #define X86_TYPE_P4                     12

479 /*
480  * x86_vendor allows us to select between
481  * implementation features and helps guide
482  * the interpretation of the cpuid instruction.
483  */
484 #define X86_VENDOR_Intel                 0
485 #define X86_VENDORSTR_Intel              "GenuineIntel"

487 #define X86_VENDOR_IntelClone            1

489 #define X86_VENDOR_AMD                    2
490 #define X86_VENDORSTR_AMD                "AuthenticAMD"

492 #define X86_VENDOR_Cyrix                  3
493 #define X86_VENDORSTR_CYRIX              "CyrixInstead"

495 #define X86_VENDOR_UMC                    4
496 #define X86_VENDORSTR_UMC                "UMC UMC UMC "

498 #define X86_VENDOR_NexGen                 5
499 #define X86_VENDORSTR_NexGen              "NexGenDriven"

501 #define X86_VENDOR_Centaur                 6
502 #define X86_VENDORSTR_Centaur              "CentaurHauls"

504 #define X86_VENDOR_Rise                    7
505 #define X86_VENDORSTR_Rise                 "RiseRiseRise"

507 #define X86_VENDOR_SiS                     8
508 #define X86_VENDORSTR_SiS                  "SiS SiS SiS "

510 #define X86_VENDOR_TM                      9
511 #define X86_VENDORSTR_TM                   "GenuineTMx86"

513 #define X86_VENDOR_NSC                     10
514 #define X86_VENDORSTR_NSC                  "Geode by NSC"

```

```

516 /*
517 * Vendor string max len + \0
518 */
519 #define X86_VENDOR_STRLEN      13

521 /*
522 * Some vendor/family/model/stepping ranges are commonly grouped under
523 * a single identifying banner by the vendor. The following encode
524 * that "revision" in a uint32_t with the 8 most significant bits
525 * identifying the vendor with X86_VENDOR_*, the next 8 identifying the
526 * family, and the remaining 16 typically forming a bitmask of revisions
527 * within that family with more significant bits indicating "later" revisions.
528 */

530 #define X86_CHIPREV_VENDOR_MASK      0xff000000u
531 #define X86_CHIPREV_VENDOR_SHIFT    24
532 #define X86_CHIPREV_FAMILY_MASK     0x00ff0000u
533 #define X86_CHIPREV_FAMILY_SHIFT    16
534 #define X86_CHIPREV_REV_MASK        0x0000ffffu

536 #define X86_CHIPREV_VENDOR(x) \
537     (((x) & X86_CHIPREV_VENDOR_MASK) >> X86_CHIPREV_VENDOR_SHIFT)
538 #define X86_CHIPREV_FAMILY(x) \
539     (((x) & X86_CHIPREV_FAMILY_MASK) >> X86_CHIPREV_FAMILY_SHIFT)
540 #define X86_CHIPREV_REV(x) \
541     ((x) & X86_CHIPREV_REV_MASK)

543 /* True if x matches in vendor and family and if x matches the given rev mask */
544 #define X86_CHIPREV_MATCH(x, mask) \
545     (_X86_CHIPREV_VENDOR(x) == _X86_CHIPREV_VENDOR(mask) && \
546     _X86_CHIPREV_FAMILY(x) == _X86_CHIPREV_FAMILY(mask) && \
547     ((_X86_CHIPREV_REV(x) & X86_CHIPREV_REV(mask)) != 0))

549 /* True if x matches in vendor and family, and rev is at least minx */
550 #define X86_CHIPREV_ATLEAST(x, minx) \
551     (_X86_CHIPREV_VENDOR(x) == _X86_CHIPREV_VENDOR(minx) && \
552     _X86_CHIPREV_FAMILY(x) == _X86_CHIPREV_FAMILY(minx) && \
553     _X86_CHIPREV_REV(x) >= _X86_CHIPREV_REV(minx))

555 #define X86_CHIPREV_MKREV(vendor, family, rev) \
556     (((uint32_t)(vendor) << X86_CHIPREV_VENDOR_SHIFT | \
557     (family) << X86_CHIPREV_FAMILY_SHIFT | (rev))

559 /* True if x matches in vendor, and family is at least minx */
560 #define X86_CHIPFAM_ATLEAST(x, minx) \
561     (_X86_CHIPREV_VENDOR(x) == _X86_CHIPREV_VENDOR(minx) && \
562     _X86_CHIPREV_FAMILY(x) >= _X86_CHIPREV_FAMILY(minx))

564 /* Revision default */
565 #define X86_CHIPREV_UNKNOWN      0x0

567 /*
568 * Definitions for AMD Family 0xf. Minor revisions C0 and CG are
569 * sufficiently different that we will distinguish them; in all other
570 * case we will identify the major revision.
571 */
572 #define X86_CHIPREV_AMD_F_REV_B _X86_CHIPREV_MKREV(X86_VENDOR_AMD, 0xf, 0x0001)
573 #define X86_CHIPREV_AMD_F_REV_C0 _X86_CHIPREV_MKREV(X86_VENDOR_AMD, 0xf, 0x0002)
574 #define X86_CHIPREV_AMD_F_REV_CG _X86_CHIPREV_MKREV(X86_VENDOR_AMD, 0xf, 0x0004)
575 #define X86_CHIPREV_AMD_F_REV_D _X86_CHIPREV_MKREV(X86_VENDOR_AMD, 0xf, 0x0008)
576 #define X86_CHIPREV_AMD_F_REV_E _X86_CHIPREV_MKREV(X86_VENDOR_AMD, 0xf, 0x0010)
577 #define X86_CHIPREV_AMD_F_REV_F _X86_CHIPREV_MKREV(X86_VENDOR_AMD, 0xf, 0x0020)
578 #define X86_CHIPREV_AMD_F_REV_G _X86_CHIPREV_MKREV(X86_VENDOR_AMD, 0xf, 0x0040)

580 /*

```

```

581 * Definitions for AMD Family 0x10. Rev A was Engineering Samples only.
582 */
583 #define X86_CHIPREV_AMD_10_REV_A \
584     _X86_CHIPREV_MKREV(X86_VENDOR_AMD, 0x10, 0x0001)
585 #define X86_CHIPREV_AMD_10_REV_B \
586     _X86_CHIPREV_MKREV(X86_VENDOR_AMD, 0x10, 0x0002)
587 #define X86_CHIPREV_AMD_10_REV_C2 \
588     _X86_CHIPREV_MKREV(X86_VENDOR_AMD, 0x10, 0x0004)
589 #define X86_CHIPREV_AMD_10_REV_C3 \
590     _X86_CHIPREV_MKREV(X86_VENDOR_AMD, 0x10, 0x0008)
591 #define X86_CHIPREV_AMD_10_REV_D0 \
592     _X86_CHIPREV_MKREV(X86_VENDOR_AMD, 0x10, 0x0010)
593 #define X86_CHIPREV_AMD_10_REV_D1 \
594     _X86_CHIPREV_MKREV(X86_VENDOR_AMD, 0x10, 0x0020)
595 #define X86_CHIPREV_AMD_10_REV_E \
596     _X86_CHIPREV_MKREV(X86_VENDOR_AMD, 0x10, 0x0040)

598 /*
599 * Definitions for AMD Family 0x11.
600 */
601 #define X86_CHIPREV_AMD_11_REV_B \
602     _X86_CHIPREV_MKREV(X86_VENDOR_AMD, 0x11, 0x0002)

604 /*
605 * Definitions for AMD Family 0x12.
606 */
607 #define X86_CHIPREV_AMD_12_REV_B \
608     _X86_CHIPREV_MKREV(X86_VENDOR_AMD, 0x12, 0x0002)

610 /*
611 * Definitions for AMD Family 0x14.
612 */
613 #define X86_CHIPREV_AMD_14_REV_B \
614     _X86_CHIPREV_MKREV(X86_VENDOR_AMD, 0x14, 0x0002)
615 #define X86_CHIPREV_AMD_14_REV_C \
616     _X86_CHIPREV_MKREV(X86_VENDOR_AMD, 0x14, 0x0004)

618 /*
619 * Definitions for AMD Family 0x15
620 */
621 #define X86_CHIPREV_AMD_15OR_REV_B2 \
622     _X86_CHIPREV_MKREV(X86_VENDOR_AMD, 0x15, 0x0001)

624 #define X86_CHIPREV_AMD_15TN_REV_A1 \
625     _X86_CHIPREV_MKREV(X86_VENDOR_AMD, 0x15, 0x0002)

627 /*
628 * Various socket/package types, extended as the need to distinguish
629 * a new type arises. The top 8 byte identifies the vendor and the
630 * remaining 24 bits describe 24 socket types.
631 */

633 #define X86_SOCKET_VENDOR_SHIFT      24
634 #define X86_SOCKET_VENDOR(x)        ((x) >> X86_SOCKET_VENDOR_SHIFT)
635 #define X86_SOCKET_TYPE_MASK        0x00ffffff
636 #define X86_SOCKET_TYPE(x)          ((x) & X86_SOCKET_TYPE_MASK)

638 #define X86_SOCKET_MKVAL(vendor, bitval) \
639     (((uint32_t)(vendor) << X86_SOCKET_VENDOR_SHIFT | (bitval))

641 #define X86_SOCKET_MATCH(s, mask) \
642     (_X86_SOCKET_VENDOR(s) == _X86_SOCKET_VENDOR(mask) && \
643     (_X86_SOCKET_TYPE(s) & X86_SOCKET_TYPE(mask)) != 0)

645 #define X86_SOCKET_UNKNOWN 0x0
646 /*

```

```

647      * AMD socket types
648      */
649 #define X86_SOCKET_754      _X86_SOCKET_MKVAL(X86_VENDOR_AMD, 0x000001)
650 #define X86_SOCKET_939      _X86_SOCKET_MKVAL(X86_VENDOR_AMD, 0x000002)
651 #define X86_SOCKET_940      _X86_SOCKET_MKVAL(X86_VENDOR_AMD, 0x000004)
652 #define X86_SOCKET_S1g1     _X86_SOCKET_MKVAL(X86_VENDOR_AMD, 0x000008)
653 #define X86_SOCKET_AM2      _X86_SOCKET_MKVAL(X86_VENDOR_AMD, 0x000010)
654 #define X86_SOCKET_F1207    _X86_SOCKET_MKVAL(X86_VENDOR_AMD, 0x000020)
655 #define X86_SOCKET_S1g2     _X86_SOCKET_MKVAL(X86_VENDOR_AMD, 0x000040)
656 #define X86_SOCKET_S1g3     _X86_SOCKET_MKVAL(X86_VENDOR_AMD, 0x000080)
657 #define X86_SOCKET_AM       _X86_SOCKET_MKVAL(X86_VENDOR_AMD, 0x000100)
658 #define X86_SOCKET_AM2R2    _X86_SOCKET_MKVAL(X86_VENDOR_AMD, 0x000200)
659 #define X86_SOCKET_AM3      _X86_SOCKET_MKVAL(X86_VENDOR_AMD, 0x000400)
660 #define X86_SOCKET_G34      _X86_SOCKET_MKVAL(X86_VENDOR_AMD, 0x000800)
661 #define X86_SOCKET_AS2      _X86_SOCKET_MKVAL(X86_VENDOR_AMD, 0x001000)
662 #define X86_SOCKET_C32      _X86_SOCKET_MKVAL(X86_VENDOR_AMD, 0x002000)
663 #define X86_SOCKET_S1g4     _X86_SOCKET_MKVAL(X86_VENDOR_AMD, 0x004000)
664 #define X86_SOCKET_FT1      _X86_SOCKET_MKVAL(X86_VENDOR_AMD, 0x008000)
665 #define X86_SOCKET_FM1      _X86_SOCKET_MKVAL(X86_VENDOR_AMD, 0x010000)
666 #define X86_SOCKET_FS1      _X86_SOCKET_MKVAL(X86_VENDOR_AMD, 0x020000)
667 #define X86_SOCKET_AM3R2    _X86_SOCKET_MKVAL(X86_VENDOR_AMD, 0x040000)
668 #define X86_SOCKET_FP2      _X86_SOCKET_MKVAL(X86_VENDOR_AMD, 0x080000)
669 #define X86_SOCKET_FS1R2    _X86_SOCKET_MKVAL(X86_VENDOR_AMD, 0x100000)
670 #define X86_SOCKET_FM2      _X86_SOCKET_MKVAL(X86_VENDOR_AMD, 0x200000)

672 /*
673  * xgetbv/xsetbv support
674  * See section 13.3 in vol. 1 of the Intel developers manual.
675  */

677 #define XFEATURE_ENABLED_MASK    0x0
678 /*
679  * XFEATURE_ENABLED_MASK values (eax)
680  * See setup_xfem().
681  */
682 #define XFEATURE_LEGACY_FP      0x1
683 #define XFEATURE_SSE            0x2
684 #define XFEATURE_AVX            0x4
685 #define XFEATURE_MPX            0x18 /* 2 bits, both 0 or 1 */
686 #define XFEATURE_AVX512        0xe0 /* 3 bits, all 0 or 1 */
687 /* bit 8 unused */
688 #define XFEATURE_PKRU           0x200
689 #define XFEATURE_FP_ALL \
690 (XFEATURE_LEGACY_FP | XFEATURE_SSE | XFEATURE_AVX | XFEATURE_MPX | \
691  XFEATURE_AVX512 | XFEATURE_PKRU)

693 #if !defined(_ASM)

695 #if defined(_KERNEL) || defined(_KMEMUSER)

697 #define NUM_X86_FEATURES      71
698 extern uchar_t x86_featureset[];

700 extern void free_x86_featureset(void *featureset);
701 extern boolean_t is_x86_feature(void *featureset, uint_t feature);
702 extern void add_x86_feature(void *featureset, uint_t feature);
703 extern void remove_x86_feature(void *featureset, uint_t feature);
704 extern boolean_t compare_x86_featureset(void *setA, void *setB);
705 extern void print_x86_featureset(void *featureset);

708 extern uint_t x86_type;
709 extern uint_t x86_vendor;
710 extern uint_t x86_clflush_size;

```

```

712 extern uint_t pentiumpro_bug4046376;

714 extern const char CyrixInstead[];

716 #endif

718 #if defined(_KERNEL)

720 /*
721  * This structure is used to pass arguments and get return values back
722  * from the CPUID instruction in __cpuid_insn() routine.
723  */
724 struct cpuid_regs {
725     uint32_t      cp_eax;
726     uint32_t      cp_ebx;
727     uint32_t      cp_ecx;
728     uint32_t      cp_edx;
729 };

731 extern int x86_use_pcid;
732 extern int x86_use_invpcid;

734 /*
735  * Utility functions to get/set extended control registers (XCR)
736  * Initial use is to get/set the contents of the XFEATURE_ENABLED_MASK.
737  */
738 extern uint64_t get_xcr(uint_t);
739 extern void set_xcr(uint_t, uint64_t);

741 extern uint64_t rdmsr(uint_t);
742 extern void wrmsr(uint_t, const uint64_t);
743 extern uint64_t xrdmsr(uint_t);
744 extern void xwrmsr(uint_t, const uint64_t);
745 extern int checked_rdmsr(uint_t, uint64_t *);
746 extern int checked_wrmsr(uint_t, uint64_t);

748 extern void invalidate_cache(void);
749 extern ulong_t getcr4(void);
750 extern void setcr4(ulong_t);

752 extern void mtrr_sync(void);

754 extern void cpu_fast_syscall_enable(void *);
755 extern void cpu_fast_syscall_disable(void *);

757 struct cpu;

759 extern int cpuid_checkpass(struct cpu *, int);
760 extern uint32_t cpuid_insn(struct cpu *, struct cpuid_regs *);
761 extern uint32_t __cpuid_insn(struct cpuid_regs *);
762 extern int cpuid_getbrandstr(struct cpu *, char *, size_t);
763 extern int cpuid_getidstr(struct cpu *, char *, size_t);
764 extern const char *cpuid_getvendorstr(struct cpu *);
765 extern uint_t cpuid_getvendor(struct cpu *);
766 extern uint_t cpuid_getfamily(struct cpu *);
767 extern uint_t cpuid_getmodel(struct cpu *);
768 extern uint_t cpuid_getstep(struct cpu *);
769 extern uint_t cpuid_getsig(struct cpu *);
770 extern uint_t cpuid_get_ncpu_per_chip(struct cpu *);
771 extern uint_t cpuid_get_ncore_per_chip(struct cpu *);
772 extern uint_t cpuid_get_ncpu_sharing_last_cache(struct cpu *);
773 extern id_t cpuid_get_last_lvl_cacheid(struct cpu *);
774 extern int cpuid_get_chipid(struct cpu *);
775 extern id_t cpuid_get_coreid(struct cpu *);
776 extern int cpuid_get_pkgcoreid(struct cpu *);
777 extern int cpuid_get_clogid(struct cpu *);

```

```

778 extern int cpuid_get_cacheid(struct cpu *);
779 extern uint32_t cpuid_get_apicid(struct cpu *);
780 extern uint_t cpuid_get_procnodid(struct cpu *cpu);
781 extern uint_t cpuid_get_procnodes_per_pkg(struct cpu *cpu);
782 extern uint_t cpuid_get_compunitid(struct cpu *cpu);
783 extern uint_t cpuid_get_cores_per_compunit(struct cpu *cpu);
784 extern size_t cpuid_get_xsave_size();
785 extern boolean_t cpuid_need_fp_excp_handling();
786 extern int cpuid_is_cmt(struct cpu *);
787 extern int cpuid_syscall32_insn(struct cpu *);
788 extern int getl2cacheinfo(struct cpu *, int *, int *, int *);

790 extern uint32_t cpuid_getchiprev(struct cpu *);
791 extern const char *cpuid_getchiprevstr(struct cpu *);
792 extern uint32_t cpuid_getsockettype(struct cpu *);
793 extern const char *cpuid_getsocketstr(struct cpu *);

795 extern int cpuid_have_cr8access(struct cpu *);

797 extern int cpuid_opteron_erratum(struct cpu *, uint_t);

799 struct cpuid_info;

801 extern void setx86isalist(void);
802 extern void cpuid_alloc_space(struct cpu *);
803 extern void cpuid_free_space(struct cpu *);
804 extern void cpuid_pass1(struct cpu *, uchar_t *);
805 extern void cpuid_pass2(struct cpu *);
806 extern void cpuid_pass3(struct cpu *);
807 extern void cpuid_pass4(struct cpu *, uint_t *);
808 extern void cpuid_set_cpu_properties(void *, processorid_t,
809     struct cpuid_info *);

811 extern void cpuid_get_addrsize(struct cpu *, uint_t *, uint_t *);
812 extern uint_t cpuid_get_dtlb_nent(struct cpu *, size_t);

814 #if !defined(__xpv)
815 extern uint32_t *cpuid_mwait_alloc(struct cpu *);
816 extern void cpuid_mwait_free(struct cpu *);
817 extern int cpuid_deep_cstates_supported(void);
818 extern int cpuid_arat_supported(void);
819 extern int cpuid_iepb_supported(struct cpu *);
820 extern int cpuid_deadline_tsc_supported(void);
821 extern void vmware_port(int, uint32_t *);
822 #endif

824 struct cpu_ucode_info;

826 extern void ucode_alloc_space(struct cpu *);
827 extern void ucode_free_space(struct cpu *);
828 extern void ucode_check(struct cpu *);
829 extern void ucode_cleanup();

831 #if !defined(__xpv)
832 extern char _tsc_mfence_start;
833 extern char _tsc_mfence_end;
834 extern char _tscp_start;
835 extern char _tscp_end;
836 extern char _no_rdtsc_start;
837 extern char _no_rdtsc_end;
838 extern char _tsc_lfence_start;
839 extern char _tsc_lfence_end;
840 #endif

842 #if !defined(__xpv)
843 extern char bcopy_patch_start;

```

```

844 extern char bcopy_patch_end;
845 extern char bcopy_ck_size;
846 #endif

848 extern void post_startup_cpu_fixups(void);

850 extern uint_t workaround_errata(struct cpu *);

852 #if defined(OPTERON_ERRATUM_93)
853 extern int opteron_erratum_93;
854 #endif

856 #if defined(OPTERON_ERRATUM_91)
857 extern int opteron_erratum_91;
858 #endif

860 #if defined(OPTERON_ERRATUM_100)
861 extern int opteron_erratum_100;
862 #endif

864 #if defined(OPTERON_ERRATUM_121)
865 extern int opteron_erratum_121;
866 #endif

868 #if defined(OPTERON_WORKAROUND_6323525)
869 extern int opteron_workaround_6323525;
870 extern void patch_workaround_6323525(void);
871 #endif

873 #if !defined(__xpv)
874 extern void determine_platform(void);
875 #endif
876 extern int get_hwenv(void);
877 extern int is_controldom(void);

879 extern void enable_pcid(void);

881 extern void xsave_setup_msr(struct cpu *);

883 /*
884  * Hypervisor signatures
885  */
886 #define HVSIG_XEN_HVM "XenVMMXenVMM"
887 #define HVSIG_VMWARE "VMwareVMware"
888 #define HVSIG_KVM "KVMKVMKVM"
889 #define HVSIG_MICROSOFT "Microsoft Hv"

891 /*
892  * Defined hardware environments
893  */
894 #define HW_NATIVE (1 << 0) /* Running on bare metal */
895 #define HW_XEN_PV (1 << 1) /* Running on Xen PVM */

897 #define HW_XEN_HVM (1 << 2) /* Running on Xen HVM */
898 #define HW_VMWARE (1 << 3) /* Running on VMware hypervisor */
899 #define HW_KVM (1 << 4) /* Running on KVM hypervisor */
900 #define HW_MICROSOFT (1 << 5) /* Running on Microsoft hypervisor */

902 #define HW_VIRTUAL (HW_XEN_HVM | HW_VMWARE | HW_KVM | HW_MICROSOFT)

904 #endif /* _KERNEL */

906 #endif /* !_ASM */

908 /*
909  * VMware hypervisor related defines

```

new/usr/src/uts/intel/sys/x86_archext.h

15

```
910 */
911 #define VMWARE_HVMAGIC      0x564d5868
912 #define VMWARE_HVPORT      0x5658
913 #define VMWARE_HVCMD_GETVERSION 0x0a
914 #define VMWARE_HVCMD_GETTSCFREQ 0x2d

916 #ifdef __cplusplus
917 }
_____unchanged_portion_omitted_
```