

new/usr/src/uts/i86pc/vm/vm_machdep.c

100184 Thu Apr 25 02:44:40 2019

new/usr/src/uts/i86pc/vm/vm_machdep.c

10806 mnode_range_setup() makes assumptions about mnodes

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

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

Reviewed by: Toomas Soome <tsoome@me.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 2019, Joyent, Inc.
27 * Copyright 2018 Joyent, Inc.
28 */

30 /* Copyright (c) 1984, 1986, 1987, 1988, 1989 AT&T */
31 /*      All Rights Reserved */
32 /*
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/sysctl.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/rman.h>
55 #include <sys/vnode.h>
56 #include <sys/cred.h>
57 #include <sys/exec.h>
```

1

new/usr/src/uts/i86pc/vm/vm_machdep.c

```
58 #include <sys/exechdr.h>
59 #include <sys/debug.h>
60 #include <sys/vmsystm.h>
61 #include <sys/swap.h>
62 #include <sys/dumpfd.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_archtev.h>
81 #include <sys/elf_386.h>
82 #include <sys/cmn_err.h>
83 #include <sys/archsystm.h>
84 #include <sys/machsystm.h>
85 #include <sys/secflags.h>

87 #include <sys/vtrace.h>
88 #include <sys/ddidmreq.h>
89 #include <sys/promif.h>
90 #include <sys/memnode.h>
91 #include <sys/stack.h>
92 #include <util/qsort.h>
93 #include <sys/taskq.h>

95 #ifdef __xpv
96
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 config_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_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_config_pfclist(uint_t);

122 /*
123 * percentage of phys mem to hold in the i/o pool
```

2

```

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;
129
130 #endif /* __xpv */
131
132 uint_t vac_colors = 1;
133
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;
142
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_pfnnhi;
150     int      mnr_mnode;
151     int      mnr_memrange;           /* index into memranges[] */
152     int      mnr_next;             /* next lower PA mnoderange */
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_mtsc_pgcnt;
163     } *mnr_mts;
164 #endif
165 } mnoderange_t;
166
167 unchanged_portion omitted
168 pfn_t *memranges = &arch_memranges[0];
169 int nranges = NUM_MEM_RANGES;
170
171 /*
172 * This combines mem_node_config and memranges into one data
173 * structure to be used for page list management.
174 */
175 static mnoderange_t *mnoderanges;
176 static int mnoderangecnt;
177 static int mtype4g;
178 static int mtype16m;
179 static int mtypetop;
180 mnoderange_t *mnoderanges;
181 int mnoderangecnt;
182 int mtype4g;
183 int mtype16m;
184 int mtypetop; /* index of highest pfn'ed mnoderange */
185
186 /*
187 * 4g memory management variables for systems with more than 4g of memory:
188 */
189 * physical memory below 4g is required for 32bit dma devices and, currently,
190 * for kmem memory. On systems with more than 4g of memory, the pool of memory

```

```

216 * below 4g can be depleted without any paging activity given that there is
217 * likely to be sufficient memory above 4g.
218 *
219 * physmax4g is set true if the largest pfn is over 4g. The rest of the
220 * 4g memory management code is enabled only when physmax4g is true.
221 *
222 * maxmem4g is the count of the maximum number of pages on the page lists
223 * with physical addresses below 4g. It can be a lot less than 4g given that
224 * BIOS may reserve large chunks of space below 4g for hot plug pci devices,
225 * agg aperture etc.
226 *
227 * freemem4g maintains the count of the number of available pages on the
228 * page lists with physical addresses below 4g.
229 *
230 * DESFREE4G specifies the desired amount of below 4g memory. It defaults to
231 * 6% (desfree4gshift = 4) of maxmem4g.
232 *
233 * RESTRICT4G_ALLOC returns true if freemem4g falls below DESFREE4G
234 * and the amount of physical memory above 4g is greater than freemem4g.
235 * In this case, page_get_* routines will restrict below 4g allocations
236 * for requests that don't specifically require it.
237 */
238
239 #define DESFREE4G      (maxmem4g >> desfree4gshift)
240
241 #define RESTRICT4G_ALLOC \
242     (physmax4g && (freemem4g < DESFREE4G) && ((freemem4g << 1) < freemem))
243
244 static pgcnt_t maxmem4g;
245 static pgcnt_t freemem4g;
246 static int     physmax4g;
247 static int     desfree4gshift = 4; /* maxmem4g shift to derive DESFREE4G */
248
249 /*
250 * 16m memory management:
251 *
252 * reserve some amount of physical memory below 16m for legacy devices.
253 *
254 * RESTRICT16M_ALLOC returns true if there are sufficient free pages above
255 * 16m or if the 16m pool drops below DESFREE16M.
256 *
257 * In this case, general page allocations via page_get_{free,cache}list
258 * routines will be restricted from allocating from the 16m pool. Allocations
259 * that require specific pfn ranges (page_get_anystart) and PG_PANIC allocations
260 * are not restricted.
261 */
262
263 #define FREEMEM16M      MTYPE_FREEMEM(mtype16m)
264 #define DESFREE16M      desfree16m
265 #define RESTRICT16M_ALLOC(freemem, pgcnt, flags) \
266     (mtype16m != -1 && (freemem != 0) && ((flags & PG_PANIC) == 0) && \
267     ((freemem != 0) && ((flags & PG_PANIC) == 0) && \
268     ((freemem >= (FREEMEM16M)) || \
269     (FREEMEM16M < (DESFREE16M + pgcnt))))
270
271 static pgcnt_t desfree16m = 0x380;
272
273 /*
274 * This can be patched via /etc/system to allow old non-PAE aware device
275 * drivers to use kmem_alloc'd memory on 32 bit systems with > 4Gig RAM.
276 */
277 #ifdef VM_STATS
278 struct {
279     ulong_t pga_alloc;

```

new/usr/src/uts/i86pc/vm/vm_machdep.c

5

```

281         ulong_t pga_notfullrange;
282         ulong_t pga_nulldmaattr;
283         ulong_t pga_allocok;
284         ulong_t pga_allocfailed;
285         ulong_t pgma_alloc;
286         ulong_t pgma_allocok;
287         ulong_t pgma_allocfailed;
288         ulong_t pgma_allocempty;
289 } pga_vmstats;
_____unchanged_portion_omitted

1392 static int
1393 mnoderange_cmp(const void *v1, const void *v2)
1394 {
1395     const mnoderange_t *m1 = v1;
1396     const mnoderange_t *m2 = v2;
1397
1398     if (m1->mnr_pfnlo < m2->mnr_pfnlo)
1399         return (-1);
1400     return (m1->mnr_pfnlo > m2->mnr_pfnlo);
1401 }

1392 /*
1393  * mnoderange_range_setup() initializes mnoderanges.
1394  */
1403 void
1404 mnoderange_range_setup(mnoderange_t *mnoderanges)
1405 {
1406     mnoderange_t *mp;
1407     size_t nr_ranges;
1408     size_t mnode;
1409     mnoderange_t *mp = mnoderanges;
1410     int mnode, mri;
1411     int mindex = 0; /* current index */
1412     int i, j;
1413     pfn_t hipfn;
1414     int last, hi;

1410     for (mnode = 0, nr_ranges = 0, mp = mnoderanges;
1411          mnode < max_mem_nodes; mnode++) {
1412         size_t mri = nranges - 1;

1415         for (mnode = 0; mnode < max_mem_nodes;
1416              if (mem_node_config[mnode].exists)
1417                  continue;

1409             mri = nranges - 1;

1417             while (MEMRANGEHI(mri) < mem_no
1418                 mri--;

1420
1421             while (mri >= 0 && mem_node_con
1422                 MEMRANGELO(mri)) {
1423                     mp->mnr_pfnlo = MAX(MEM
1424                     mnoderanges->mnr_pfnlo,
1425                     mem_node_config[mno
1426                     mp->mnr_pfghi = MIN(MEM
1427                     mnoderanges->mnr_pfghi,
1428                     mem_node_config[mno
1429                     mp->mnr_mnode = mnode;
1430                     mp->mnr_memrange = mri;
1431                     mp->mnr_next = -1;
1432                     mp->mnr_exists = 1;
1433                     mp++;
1434                     nr_ranges++;
1435                     mnoderanges->mnr_mnode

```

new/usr/src/uts/i86pc/vm/vm_machdep.c

```

1421 mmoderanges->mnr_memrange = mri;
1422 mmoderanges->mnr_exists = 1;
1423 mmoderanges++;
1424 mindex++;
1425     if (mem_node_config[mnode].physmax > MEMRANGEHI(mri))
1426         mri--;
1427     else
1428         break;
1429 }
1430 }
1431 /*  

1432 * mmoderangecnt can be larger than nr_ranges when memory DR is  

1433 * supposedly supported.  

1434 * For now do a simple sort of the mmoderanges array to fill in  

1435 * the mnr_next fields. Since mindex is expected to be relatively  

1436 * small, using a simple O(N^2) algorithm.  

1437 */
1438 VERIFY3U(nr_ranges, <=, mmoderangecnt);
1439 qsort(mmoderanges, nr_ranges, sizeof (mmoderange_t), mmoderange_cmp);
1440 /*  

1441 * If some intrepid soul takes the axe to the memory DR code, we can  

1442 * remove ->mnr_next altogether, as we just sorted by ->mnr_pfnlo order.  

1443 *  

1444 * The VERIFY3U() above can be "==" then too.  

1445 */
1446 for (size_t i = 1; i < nr_ranges; i++)
1447     mmoderanges[i].mnr_next = i - 1;
1448
1449 mtypetop = nr_ranges - 1;
1450 mtype16m = pfn_2_mtype(PFN_16MEG - 1); /* Can be -1 ... */
1451 if (physmax4g)
1452     mtype4g = pfn_2_mtype(0xfffff);
1453 for (i = 0; i < mindex; i++) {
1454     if (mp[i].mnr_pfnlo == 0)           /* find lowest */
1455         break;
1456 }
1457 ASSERT(i < mindex);
1458 last = i;
1459 mtype16m = last;
1460 mp[last].mnr_next = -1;
1461 for (i = 0; i < mindex - 1; i++) {
1462     hipfn = (pfn_t)(-1);
1463     hi = -1;
1464     /* find next highest mnode range */
1465     for (j = 0; j < mindex; j++) {
1466         if (mp[j].mnr_pfnlo > mp[last].mnr_pfnlo &&
1467             mp[j].mnr_pfnlo < hipfn) {
1468             hipfn = mp[j].mnr_pfnlo;
1469             hi = j;
1470         }
1471     }
1472     mp[hi].mnr_next = last;
1473     last = hi;
1474 }
1475 mtypetop = last;
1476 }
1477
1478 unchanged portion omitted
1479 /*
1480 * Called once at startup to configure page_coloring data structures and
1481 * does the 1st page_free()/page_freelist_add().
1482 */
1483 void

```

```
1963 page_coloring_setup(caddr_t pcmemaddr)
1964 {
1965     int      i;
1966     int      j;
1967     int      k;
1968     caddr_t  addr;
1969     int      colors;

1971     /*
1972     * do page coloring setup
1973     */
1974     addr = pcmemaddr;

1976     mnoderanges = (mnoderange_t *)addr;
1977     addr += (mnoderangecnt * sizeof (mnoderange_t));

1979     mnode_range_setup(mnoderanges);

1981     if (physmax4g)
1982         mtype4g = pfn_2_mtype(0xfffffff);

1981     for (k = 0; k < NPC_MUTEX; k++) {
1982         fpc_mutex[k] = (kmutex_t *)addr;
1983         addr += (max_mem_nodes * sizeof (kmutex_t));
1984     }
1985     for (k = 0; k < NPC_MUTEX; k++) {
1986         cpc_mutex[k] = (kmutex_t *)addr;
1987         addr += (max_mem_nodes * sizeof (kmutex_t));
1988     }
1989     page_freelists = (page_t ****)addr;
1990     addr += (mnoderangecnt * sizeof (page_t ****));

1992     page_cachelists = (page_t ***)addr;
1993     addr += (mnoderangecnt * sizeof (page_t **));

1995     for (i = 0; i < mnoderangecnt; i++) {
1996         page_freelists[i] = (page_t ***)addr;
1997         addr += (mmu_page_sizes * sizeof (page_t **));

1999         for (j = 0; j < mmu_page_sizes; j++) {
2000             colors = page_get_pagecolors(j);
2001             page_freelists[i][j] = (page_t **)addr;
2002             addr += (colors * sizeof (page_t *));
2003         }
2004         page_cachelists[i] = (page_t **)addr;
2005         addr += (page_colors * sizeof (page_t *));
2006     }
2007 }
```

unchanged portion omitted