```
*******************************************************
   95160 Tue May  7 07:35:42 2019
new/usr/src/uts/common/os/cpu.c
10923 thread_affinity_set(CPU_CURRENT) can skip cpu_lock
Reviewed by: Jerry Jelinek <jerry.jelinek@joyent.com>
Reviewed by: John Levon <john.levon@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) 1991, 2010, Oracle and/or its affiliates. All rights reserved.
  23  * Copyright (c) 2012 by Delphix. All rights reserved.
  24  * Copyright 2018 Joyent, Inc.
  25  */

  27 /*
  28  * Architecture-independent CPU control functions.
  29  */

  31 #include <sys/types.h>
  32 #include <sys/param.h>
  33 #include <sys/var.h>
  34 #include <sys/thread.h>
  35 #include <sys/cpuvar.h>
  36 #include <sys/cpu_event.h>
  37 #include <sys/kstat.h>
  38 #include <sys/uadmin.h>
  39 #include <sys/systm.h>
  40 #include <sys/errno.h>
  41 #include <sys/cmn_err.h>
  42 #include <sys/procset.h>
  43 #include <sys/processor.h>
  44 #include <sys/debug.h>
  45 #include <sys/cpupart.h>
  46 #include <sys/lgrp.h>
  47 #include <sys/pset.h>
  48 #include <sys/pghw.h>
  49 #include <sys/kmem.h>
  50 #include <sys/kmem_impl.h>      /* to set per-cpu kmem_cache offset */
  51 #include <sys/atomic.h>
  52 #include <sys/callb.h>
  53 #include <sys/vtrace.h>
  54 #include <sys/cyclic.h>
  55 #include <sys/bitmap.h>
  56 #include <sys/nvpair.h>
  57 #include <sys/pool_pset.h>
  58 #include <sys/msacct.h>
  59 #include <sys/time.h>
```

```
  60 #include <sys/archsystm.h>
  61 #include <sys/sdt.h>
  62 #if defined(__x86) || defined(__amd64)
  63 #include <sys/x86_archext.h>
  64 #endif
  65 #include <sys/callo.h>

  67 extern int      mp_cpu_start(cpu_t *);
  68 extern int      mp_cpu_stop(cpu_t *);
  69 extern int      mp_cpu_poweron(cpu_t *);
  70 extern int      mp_cpu_poweroff(cpu_t *);
  71 extern int      mp_cpu_configure(int);
  72 extern int      mp_cpu_unconfigure(int);
  73 extern void     mp_cpu_faulted_enter(cpu_t *);
  74 extern void     mp_cpu_faulted_exit(cpu_t *);

  76 extern int cmp_cpu_to_chip(processorid_t cpuid);
  77 #ifdef __sparcv9
  78 extern char *cpu_fru_fmri(cpu_t *cp);
  79 #endif

  81 static void cpu_add_active_internal(cpu_t *cp);
  82 static void cpu_remove_active(cpu_t *cp);
  83 static void cpu_info_kstat_create(cpu_t *cp);
  84 static void cpu_info_kstat_destroy(cpu_t *cp);
  85 static void cpu_stats_kstat_create(cpu_t *cp);
  86 static void cpu_stats_kstat_destroy(cpu_t *cp);

  88 static int cpu_sys_stats_ks_update(kstat_t *ksp, int rw);
  89 static int cpu_vm_stats_ks_update(kstat_t *ksp, int rw);
  90 static int cpu_stat_ks_update(kstat_t *ksp, int rw);
  91 static int cpu_state_change_hooks(int, cpu_setup_t, cpu_setup_t);

  93 /*
  94  * cpu_lock protects ncpus, ncpus_online, cpu_flag, cpu_list, cpu_active,
  95  * max_cpu_seqid_ever, and dispatch queue reallocations.  The lock ordering with
  96  * respect to related locks is:
  97  *
  98  *      cpu_lock --> thread_free_lock  --->  p_lock  --->  thread_lock()
  99  *
 100  * Warning:  Certain sections of code do not use the cpu_lock when
 101  * traversing the cpu_list (e.g. mutex_vector_enter(), clock()).  Since
 102  * all cpus are paused during modifications to this list, a solution
 103  * to protect the list is too either disable kernel preemption while
 104  * walking the list, *or* recheck the cpu_next pointer at each
 105  * iteration in the loop.  Note that in no cases can any cached
 106  * copies of the cpu pointers be kept as they may become invalid.
 107  */
 108 kmutex_t        cpu_lock;
 109 cpu_t           *cpu_list;              /* list of all CPUs */
 110 cpu_t           *clock_cpu_list;        /* used by clock to walk CPUs */
 111 cpu_t           *cpu_active;            /* list of active CPUs */
 112 static cpuset_t cpu_available;          /* set of available CPUs */
 113 cpuset_t        cpu_seqid_inuse;        /* which cpu_seqids are in use */

 115 cpu_t           **cpu_seq;              /* ptrs to CPUs, indexed by seq_id */

 117 /*
 118  * max_ncpus keeps the max cpus the system can have. Initially
 119  * it's NCPU, but since most archs scan the devtree for cpus
 120  * fairly early on during boot, the real max can be known before
 121  * ncpus is set (useful for early NCPU based allocations).
 122  */
 123 int max_ncpus = NCPU;
 124 /*
 125  * platforms that set max_ncpus to maxiumum number of cpus that can be
```

```
126  * dynamically added will set boot_max_ncpus to the number of cpus found
127  * at device tree scan time during boot.
128  */
129 int boot_max_ncpus = -1;
130 int boot_ncpus = -1;
131 /*
132  * Maximum possible CPU id.  This can never be >= NCPU since NCPU is
133  * used to size arrays that are indexed by CPU id.
134  */
135 processorid_t max_cpuid = NCPU - 1;

137 /*
138  * Maximum cpu_seqid was given. This number can only grow and never shrink. It
139  * can be used to optimize NCPU loops to avoid going through CPUs which were
140  * never on-line.
141  */
142 processorid_t max_cpu_seqid_ever = 0;

144 int ncpus = 1;
145 int ncpus_online = 1;

147 /*
148  * CPU that we're trying to offline.  Protected by cpu_lock.
149  */
150 cpu_t *cpu_inmotion;

152 /*
153  * Can be raised to suppress further weakbinding, which are instead
154  * satisfied by disabling preemption.  Must be raised/lowered under cpu_lock,
155  * while individual thread weakbinding synchronization is done under thread
156  * lock.
157  */
158 int weakbindingbarrier;

160 /*
161  * Variables used in pause_cpus().
162  */
163 static volatile char safe_list[NCPU];

165 static struct _cpu_pause_info {
166         int             cp_spl;        /* spl saved in pause_cpus() */
167         volatile int    cp_go;         /* Go signal sent after all ready */
168         int             cp_count;      /* # of CPUs to pause */
169         ksema_t         cp_sem;        /* synch pause_cpus & cpu_pause */
170         kthread_id_t    cp_paused;
171         void            *(*cp_func)(void *);
172 } cpu_pause_info;
```
_____*unchanged_portion_omitted_*

```
388 /*
389  * Set affinity for a specified CPU.
390  *
391  * Specifying a cpu_id of CPU_CURRENT, allowed _only_ when setting affinity for
392  * curthread, will set affinity to the CPU on which the thread is currently
393  * running.  For other cpu_id values, the caller must ensure that the
394  * referenced CPU remains valid, which can be done by holding cpu_lock across
395  * this call.
396  *
397  * CPU affinity is guaranteed after return of thread_affinity_set().  If a
398  * caller setting affinity to CPU_CURRENT requires that its thread not migrate
399  * CPUs prior to a successful return, it should take extra precautions (such as
400  * their own call to kpreempt_disable) to ensure that safety.
401  *
402  * A CPU affinity reference count is maintained by thread_affinity_set and
403  * thread_affinity_clear (incrementing and decrementing it, respectively),
404  * maintaining CPU affinity while the count is non-zero, and allowing regions
```

```
405  * of code which require affinity to be nested.
389  * A reference count is incremented and the affinity is held until the
390  * reference count is decremented to zero by thread_affinity_clear().
391  * This is so regions of code requiring affinity can be nested.
392  * Caller needs to ensure that cpu_id remains valid, which can be
393  * done by holding cpu_lock across this call, unless the caller
394  * specifies CPU_CURRENT in which case the cpu_lock will be acquired
395  * by thread_affinity_set and CPU->cpu_id will be the target CPU.
406  */
407 void
408 thread_affinity_set(kthread_id_t t, int cpu_id)
409 {
410         cpu_t *cp;
401         int             c;

412         ASSERT(!(t == curthread && t->t_weakbound_cpu != NULL));

414         if (cpu_id == CPU_CURRENT) {
415                 VERIFY3P(t, ==, curthread);
416                 kpreempt_disable();
417                 cp = CPU;
418         } else {
405         if ((c = cpu_id) == CPU_CURRENT) {
406                 mutex_enter(&cpu_lock);
407                 cpu_id = CPU->cpu_id;
408         }
419                 /*
420                  * We should be asserting that cpu_lock is held here, but
421                  * the NCA code doesn't acquire it.  The following assert
422                  * should be uncommented when the NCA code is fixed.
423                  *
424                  * ASSERT(MUTEX_HELD(&cpu_lock));
425                  */
426                 VERIFY((cpu_id >= 0) && (cpu_id < NCPU));
416         ASSERT((cpu_id >= 0) && (cpu_id < NCPU));
427                 cp = cpu[cpu_id];

429                 /* user must provide a good cpu_id */
430                 VERIFY(cp != NULL);
431         }

418         ASSERT(cp != NULL);                 /* user must provide a good cpu_id */
433         /*
434          * If there is already a hard affinity requested, and this affinity
435          * conflicts with that, panic.
436          */
437         thread_lock(t);
438         if (t->t_affinitycnt > 0 && t->t_bound_cpu != cp) {
439                 panic("affinity_set: setting %p but already bound to %p",
440                     (void *)cp, (void *)t->t_bound_cpu);
441         }
442         t->t_affinitycnt++;
443         t->t_bound_cpu = cp;

445         /*
446          * Make sure we're running on the right CPU.
447          */
448         if (cp != t->t_cpu || t != curthread) {
449                 ASSERT(cpu_id != CPU_CURRENT);
450                 force_thread_migrate(t);         /* drops thread lock */
451         } else {
452                 thread_unlock(t);
453         }

455         if (cpu_id == CPU_CURRENT) {
456                 kpreempt_enable();
```

```
457              }
440              if (c == CPU_CURRENT)
441                      mutex_exit(&cpu_lock);
458 }
_____unchanged_portion_omitted_
```