

new/usr/src/uts/common/cpr/cpr_main.c

```
*****
34820 Thu May 1 19:50:17 2014
new/usr/src/uts/common/cpr/cpr_main.c
XXXX pass in cpu_pause_func via pause_cpus
*****  
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 */  
  
26 /*
27 * This module contains the guts of checkpoint-resume mechanism.
28 * All code in this module is platform independent.
29 */  
  
31 #include <sys/types.h>
32 #include <sys/errno.h>
33 #include <sys/callb.h>
34 #include <sys/processor.h>
35 #include <sys/machsysm.h>
36 #include <sys/clock.h>
37 #include <sys/vfs.h>
38 #include <sys/kmem.h>
39 #include <nfs/lm.h>
40 #include <sys/sysm.h>
41 #include <sys/cpr.h>
42 #include <sys/bootconf.h>
43 #include <sys/cyclic.h>
44 #include <sys/filio.h>
45 #include <sys/fs/ufs_filio.h>
46 #include <sys/epm.h>
47 #include <sys/modctl.h>
48 #include <sys/reboot.h>
49 #include <sys/kdi.h>
50 #include <sys/promif.h>
51 #include <sys/srn.h>
52 #include <sys/cpr_impl.h>  
  
54 #define PPM(dip) ((dev_info_t *)DEVI(dip)->devi_pm_ppm)  
  
56 extern struct cpr_terminator cpr_term;  
  
58 extern int cpr_alloc_statefile(int);
59 extern void cpr_start_kernel_threads(void);
60 extern void cpr_abbreviate_devpath(char *, char *);
61 extern void cpr_convert_promtime(cpr_time_t *);
```

1

new/usr/src/uts/common/cpr/cpr_main.c

```
62 extern void cpr_send_notice(void);
63 extern void cpr_set_bitmap_size(void);
64 extern void cpr_stat_init();
65 extern void cpr_statef_close(void);
66 extern void flush_windows(void);
67 extern void (*srn_signal)(int, int);
68 extern void init_cpu_syscall(struct cpu *);
69 extern void i_cpr_pre_resume_cpus();
70 extern void i_cpr_post_resume_cpus();
71 extern int cpr_is_ufs(struct vfs *);  
  
73 extern int pm_powering_down;
74 extern kmutex_t srn_clone_lock;
75 extern int srn_inuse;  
  
77 static int cpr_suspend(int);
78 static int cpr_resume(int);
79 static void cpr_suspend_init(int);
80 #if defined(__x86)
81 static int cpr_suspend_cpus(void);
82 static void cpr_resume_cpus(void);
83 #endif
84 static int cpr_all_online(void);
85 static void cpr_restore_offline(void);  
  
87 cpr_time_t wholecycle_tv;
88 int cpr_suspend_succeeded;
89 pfnt curthreaddpfn;
90 int curthreadremapped;  
  
92 extern cpuset_t cpu_ready_set;
93 extern void *(cpu_pause_func)(void *);  
  
94 extern processorid_t i_cpr_bootcpuid(void);
95 extern cpu_t *i_cpr_bootcpu(void);
96 extern void tsc_adjust_delta(hrtimedelta tdelta);
97 extern void tsc_resume(void);
98 extern int tsc_resume_in_cyclic;  
  
100 /*
101 * Set this variable to 1, to have device drivers resume in an
102 * uniprocessor environment. This is to allow drivers that assume
103 * that they resume on a UP machine to continue to work. Should be
104 * deprecated once the broken drivers are fixed
105 */
106 int cpr_resume_uniproc = 0;  
  
108 /*
109 * save or restore abort_enable; this prevents a drop
110 * to kadb or prom during cpr_resume_devices() when
111 * there is no kbd present; see abort_sequence_enter()
112 */
113 static void
114 cpr_sae(int stash)
115 {
116     static int saved_ae = -1;  
  
118     if (stash) {
119         saved_ae = abort_enable;
120         abort_enable = 0;
121     } else if (saved_ae != -1) {
122         abort_enable = saved_ae;
123         saved_ae = -1;
124     }
125 }  
_____unchanged_portion_omitted_____
```

2

```

384 int
385 cpr_suspend_cpus(void)
386 {
387     int ret = 0;
388     extern void *i_cpr_save_context(void *arg);
389
390     mutex_enter(&cpu_lock);
391
392     /*
393      * the machine could not have booted without a bootcpu
394      */
395     ASSERT(i_cpr_bootcpu() != NULL);
396
397     /*
398      * bring all the offline cpus online
399      */
400     if ((ret = cpr_all_online()) == 0) {
401         mutex_exit(&cpu_lock);
402         return (ret);
403     }
404
405     /*
406      * Set the affinity to be the boot processor
407      * This is cleared in either cpr_resume_cpus() or cpr_unpause_cpus()
408      */
409     affinity_set(i_cpr_bootcpuid());
410
411     ASSERT(CPU->cpu_id == 0);
412
413     PMD(PMD_SX, ("curthread running on bootcpu\n"));
414
415     /*
416      * pause all other running CPUs and save the CPU state at the sametime
417      */
418     pause_cpus(NULL, i_cpr_save_context);
419     cpu_pause_func = i_cpr_save_context;
420     pause_cpus(NULL);
421
422     mutex_exit(&cpu_lock);
423 }


---


unchanged_portion_omitted_

```

764 void
765 cpr_resume_cpus(void)
766 {
767 /*
768 * this is a cut down version of start_other_cpus()
769 * just do the initialization to wake the other cpus
770 */
772 #if defined(__x86)
773 /*
774 * Initialize our syscall handlers
775 */
776 init_cpu_syscall(CPU);
778 #endif
780 i_cpr_pre_resume_cpus();
782 /*
783 * Restart the paused cpus
784 */

```

785     mutex_enter(&cpu_lock);
786     start_cpus();
787     mutex_exit(&cpu_lock);
788
789     i_cpr_post_resume_cpus();
790
791     mutex_enter(&cpu_lock);
792
793     /*
794      * Restore this cpu to use the regular cpu_pause(), so that
795      * online and offline will work correctly
796      */
797     cpu_pause_func = NULL;
798
799     /*
800      * clear the affinity set in cpr_suspend_cpus()
801      */
802     affinity_clear();
803
804     /*
805      * offline all the cpus that were brought online during suspend
806      */
807     cpr_restore_offline();
808
809     mutex_exit(&cpu_lock);
810 }
811
812 void
813 cpr_unpause_cpus(void)
814 {
815     /*
816      * Now restore the system back to what it was before we suspended
817      */
818     PMD(PMD_SX, ("cpr_unpause_cpus: restoring system\n"));
819
820     mutex_enter(&cpu_lock);
821
822     /*
823      * Restore this cpu to use the regular cpu_pause(), so that
824      * online and offline will work correctly
825      */
826     cpu_pause_func = NULL;
827
828     /*
829      * Restart the paused cpus
830      */
831     start_cpus();
832
833     /*
834      * clear the affinity set in cpr_suspend_cpus()
835      */
836     affinity_clear();
837
838     /*
839      * offline all the cpus that were brought online during suspend
840      */
841     cpr_restore_offline();
842
843     /*
844      * Bring the system back up from a checkpoint, at this point
845      * the VM has been minimally restored by boot, the following
846      * are executed sequentially:
847      */

```

```

838 *      - machdep setup and enable interrupts (mp startup if it's mp)
839 *      - resume all devices
840 *      - restart daemons
841 *      - put all threads back on run queue
842 */
843 static int
844 cpr_resume(int sleeptype)
845 {
846     cpr_time_t pwrone_tv, *ctp;
847     char *str;
848     int rc = 0;
849
850     /*
851      * The following switch is used to resume the system
852      * that was suspended to a different level.
853      */
854     CPR_DEBUG(CPR_DEBUG1, "\nEntering cpr_resume...\n");
855     PMD(PMD_SX, ("cpr_resume %x\n", sleeptype));
856
857     /*
858      * Note:
859      *
860      * The rollback labels rb_xyz do not represent the cpr resume
861      * state when event 'xyz' has happened. Instead they represent
862      * the state during cpr suspend when event 'xyz' was being
863      * entered (and where cpr suspend failed). The actual call that
864      * failed may also need to be partially rolled back, since they
865      * aren't atomic in most cases. In other words, rb_xyz means
866      * "roll back all cpr suspend events that happened before 'xyz',
867      * and the one that caused the failure, if necessary."
868      */
869     switch (CPR->c_substate) {
870 #if defined(__sparc)
871     case C_ST_DUMP:
872         /*
873          * This is most likely a full-fledged cpr_resume after
874          * a complete and successful cpr suspend. Just roll back
875          * everything.
876          */
877         ASSERT(sleeptype == CPR_TODISK);
878         break;
879
880     case C_ST_REUSEABLE:
881     case C_ST_DUMP_NOSPC:
882     case C_ST_SETPROPS_0:
883     case C_ST_SETPROPS_1:
884         /*
885          * C_ST_REUSEABLE and C_ST_DUMP_NOSPC are the only two
886          * special switch cases here. The other two do not have
887          * any state change during cpr_suspend() that needs to
888          * be rolled back. But these are exit points from
889          * cpr_suspend, so theoretically (or in the future), it
890          * is possible that a need for roll back of a state
891          * change arises between these exit points.
892          */
893         ASSERT(sleeptype == CPR_TODISK);
894         goto rb_dump;
895 #endif
896
897     case C_ST_NODUMP:
898         PMD(PMD_SX, ("cpr_resume: NODUMP\n"));
899         goto rb_nodump;
900
901     case C_ST_STOP_KERNEL_THREADS:
902         PMD(PMD_SX, ("cpr_resume: STOP_KERNEL_THREADS\n"));
903         goto rb_stop_kernel_threads;

```

```

905     case C_ST_SUSPEND_DEVICES:
906         PMD(PMD_SX, ("cpr_resume: SUSPEND_DEVICES\n"));
907         goto rb_suspend_devices;
908
909 #if defined(__sparc)
910     case C_ST_STATEF_ALLOC:
911         ASSERT(sleeptype == CPR_TODISK);
912         goto rb_statef_alloc;
913
914     case C_ST_DISABLE_UFS_LOGGING:
915         ASSERT(sleeptype == CPR_TODISK);
916         goto rb_disable_ufs_logging;
917 #endif
918
919     case C_ST_PM_REATTACH_NOINVOL:
920         PMD(PMD_SX, ("cpr_resume: REATTACH_NOINVOL\n"));
921         goto rb_pm_reattach_noinvol;
922
923     case C_ST_STOP_USER_THREADS:
924         PMD(PMD_SX, ("cpr_resume: STOP_USER_THREADS\n"));
925         goto rb_stop_user_threads;
926
927 #if defined(__sparc)
928     case C_ST_MP_OFFLINE:
929         PMD(PMD_SX, ("cpr_resume: MP_OFFLINE\n"));
930         goto rb_mp_offline;
931 #endif
932
933 #if defined(__x86)
934     case C_ST_MP_PAUSED:
935         PMD(PMD_SX, ("cpr_resume: MP_PAUSED\n"));
936         goto rb_mp_paused;
937 #endif
938
939     default:
940         PMD(PMD_SX, ("cpr_resume: others\n"));
941         goto rb_others;
942     }
943
944 rb_all:
945     /*
946      * perform platform-dependent initialization
947      */
948     if (cpr_suspend_succeeded)
949         i_cpr_machdep_setup();
950
951     /*
952      * system did not really go down if we jump here
953      */
954
955 rb_dump:
956     /*
957      * IMPORTANT: SENSITIVE RESUME SEQUENCE
958      *
959      * DO NOT ADD ANY INITIALIZATION STEP BEFORE THIS POINT!!
960      */
961 rb_nodump:
962     /*
963      * If we did suspend to RAM, we didn't generate a dump
964      */
965     PMD(PMD_SX, ("cpr_resume: CPR DMA callback\n"));
966     (void) callb_execute_class(CB_CL_CPR_DMA, CB_CODE_CPR_RESUME);
967     if (cpr_suspend_succeeded) {
968         PMD(PMD_SX, ("cpr_resume: CPR RPC callback\n"));
969         (void) callb_execute_class(CB_CL_CPR_RPC, CB_CODE_CPR_RESUME);

```

```

970         }
971
972         prom_resume_repost();
973 #if !defined(__sparc)
974         /*
975          * Need to sync the software clock with the hardware clock.
976          * On Sparc, this occurs in the sparc-specific cbe. However
977          * on x86 this needs to be handled _before_ we bring other cpu's
978          * back online. So we call a resume function in timestamp.c
979         */
980         if (tsc_resume_in_cyclic == 0)
981             tsc_resume();
982
983 #endif
984
985 #if defined(__sparc)
986     if (cpr_suspend_succeeded && (boothowto & RB_DEBUG))
987         kdi_dvec_cpr_restart();
988 #endif
989
990 #if defined(__x86)
991 rb_mp_paused:
992     PT(PT_RMPO);
993     PMD(PMD_SX, ("resume aux cpus\n"))
994
995     if (cpr_suspend_succeeded) {
996         cpr_resume_cpus();
997     } else {
998         cpr_unpause_cpus();
999     }
1000 }
1001 #endif
1002
1003 /*
1004  * let the tmp callout catch up.
1005 */
1006 PMD(PMD_SX, ("cpr_resume: CPR CALLOUT callback\n"))
1007 (void) callb_execute_class(CB_CL_CPR_CALLOUT, CB_CODE_CPR_RESUME);
1008
1009 i_cpr_enable_intr();
1010
1011 mutex_enter(&cpu_lock);
1012 PMD(PMD_SX, ("cpr_resume: cyclic resume\n"))
1013 cyclic_resume();
1014 mutex_exit(&cpu_lock);
1015
1016 PMD(PMD_SX, ("cpr_resume: handle xc\n"))
1017 i_cpr_handle_xc(); /* turn it off to allow xc assertion */
1018
1019 PMD(PMD_SX, ("cpr_resume: CPR POST KERNEL callback\n"))
1020 (void) callb_execute_class(CB_CL_CPR_POST_KERNEL, CB_CODE_CPR_RESUME);
1021
1022 /*
1023  * statistics gathering
1024 */
1025 if (cpr_suspend_succeeded) {
1026     /*
1027      * Prevent false alarm in tod_validate() due to tod
1028      * value change between suspend and resume
1029      */
1030     cpr_tod_status_set(TOD_CPR_RESUME_DONE);
1031
1032     cpr_convert_promtime(&pwrone_tv);
1033
1034     ctp = &cpr_term.tm_shutdown;
1035     if (sleeptype == CPR_TODISK)

```

```

1036             CPR_STAT_EVENT_END_TMZ(" write statefile", ctp);
1037             CPR_STAT_EVENT_END_TMZ("Suspend Total", ctp);
1038
1039             CPR_STAT_EVENT_START_TMZ("Resume Total", &pwrone_tv);
1040             str = " prom time";
1041             CPR_STAT_EVENT_START_TMZ(str, &pwrone_tv);
1042             ctp = &cpr_term.tm_cprboot_start;
1043             CPR_STAT_EVENT_END_TMZ(str, ctp);
1044
1045             str = " read statefile";
1046             CPR_STAT_EVENT_START_TMZ(str, ctp);
1047             ctp = &cpr_term.tm_cprboot_end;
1048             CPR_STAT_EVENT_END_TMZ(str, ctp);
1049
1050         }
1051
1052 rb_stop_kernel_threads:
1053     /*
1054      * Put all threads back to where they belong; get the kernel
1055      * daemons straightened up too. Note that the callback table
1056      * locked during cpr_stop_kernel_threads() is released only
1057      * in cpr_start_kernel_threads(). Ensure modunloading is
1058      * disabled before starting kernel threads, we don't want
1059      * modunload thread to start changing device tree underneath.
1060      */
1061     PMD(PMD_SX, ("cpr_resume: modunload disable\n"))
1062     modunload_disable();
1063     PMD(PMD_SX, ("cpr_resume: start kernel threads\n"))
1064     cpr_start_kernel_threads();
1065
1066 rb_suspend_devices:
1067     CPR_DEBUG(CPR_DEBUG1, "resuming devices...");
1068     CPR_STAT_EVENT_START(" start drivers");
1069
1070     PMD(PMD_SX,
1071         ("cpr_resume: rb_suspend_devices: cpr_resume_uniproc = %d\n",
1072          cpr_resume_uniproc));
1073
1074 #if defined(__x86)
1075     /*
1076      * If cpr_resume_uniproc is set, then pause all the other cpus
1077      * apart from the current cpu, so that broken drivers that think
1078      * that they are on a uniprocessor machine will resume
1079      */
1080     if (cpr_resume_uniproc) {
1081         mutex_enter(&cpu_lock);
1082         pause_cpus(NULL, NULL);
1083         pause_cpus(NULL);
1084     }
1085 #endif
1086
1087     /*
1088      * The policy here is to continue resume everything we can if we did
1089      * not successfully finish suspend; and panic if we are coming back
1090      * from a fully suspended system.
1091      */
1092     PMD(PMD_SX, ("cpr_resume: resume devices\n"))
1093     rc = cpr_resume_devices(ddi_root_node(), 0);
1094
1095     cpr_sae(0);
1096
1097     str = "Failed to resume one or more devices.";
1098
1099     if (rc) {
1100         if (CPR->c_substate == C_ST_DUMP ||
```

```

1101     (sleeptype == CPR_TORAM &&
1102      CPR->c_substate == C_ST_NODUMP)) {
1103         if (cpr_test_point == FORCE_SUSPEND_TO_RAM) {
1104             PMD(PMD_SX, ("cpr_resume: resume device "
1105                         "warn\n"))
1106             cpr_err(CE_WARN, str);
1107         } else {
1108             PMD(PMD_SX, ("cpr_resume: resume device "
1109                         "panic\n"))
1110             cpr_err(CE_PANIC, str);
1111         }
1112     } else {
1113         PMD(PMD_SX, ("cpr_resume: resume device warn\n"))
1114         cpr_err(CE_WARN, str);
1115     }
1116 }

1117 CPR_STAT_EVENT_END(" start drivers");
1118 CPR_DEBUG(CPR_DEBUG1, "done\n");

1119 #if defined(__x86)
1120 /*
1121  * If cpr_resume_uniproc is set, then unpause all the processors
1122  * that were paused before resuming the drivers
1123  */
1124 if (cpr_resume_uniproc) {
1125     mutex_enter(&cpu_lock);
1126     start_cpus();
1127     mutex_exit(&cpu_lock);
1128 }
1129 #endiff

1130 /*
1131  * If we had disabled modunloading in this cpr resume cycle (i.e. we
1132  * resumed from a state earlier than C_ST_SUSPEND_DEVICES), re-enable
1133  * modunloading now.
1134  */
1135 if (CPR->c_substate != C_ST_SUSPEND_DEVICES) {
1136     PMD(PMD_SX, ("cpr_resume: modload enable\n"))
1137     modunload_enable();
1138 }
1139 /*
1140  * Hooks needed by lock manager prior to resuming.
1141  * Refer to code for more comments.
1142  */
1143 PMD(PMD_SX, ("cpr_resume: lock mgr\n"))
1144 cpr_lock_mgr(lm_cprresume);

1145 #if defined(__sparc)
1146 /*
1147  * This is a partial (half) resume during cpr suspend, we
1148  * haven't yet given up on the suspend. On return from here,
1149  * cpr_suspend() will try to reallocate and retry the suspend.
1150  */
1151 if (CPR->c_substate == C_ST_DUMP_NOSPC) {
1152     return (0);
1153 }
1154 */

1155 if (sleeptype == CPR_TODISK) {
1156     rb_statef_alloc:
1157         cpr_statef_close();
1158 }

1159 if (sleeptype == CPR_TORAM) {
1160     rb_disable_ufs_logging:
1161         /* if ufs logging was disabled, re-enable
1162

```

```

1163         */
1164         (void) cpr_ufs_logging(1);
1165     }
1166 #endiff

1167 rb_pm_reattach_noinvol:
1168 /*
1169  * When pm_reattach_noinvol() succeeds, modunload_thread will
1170  * remain disabled until after cpr suspend passes the
1171  * C_ST_STOP_KERNEL_THREADS state. If any failure happens before
1172  * cpr suspend reaches this state, we'll need to enable modunload
1173  * thread during rollback.
1174
1175 if (CPR->c_substate == C_ST_DISABLE_UFS_LOGGING ||
1176     CPR->c_substate == C_ST_STATEF_ALLOC ||
1177     CPR->c_substate == C_ST_SUSPEND_DEVICES ||
1178     CPR->c_substate == C_ST_STOP_KERNEL_THREADS) {
1179     PMD(PMD_SX, ("cpr_resume: reattach noinvol fini\n"))
1180     pm_reattach_noinvol_fini();
1181 }
1182
1183 PMD(PMD_SX, ("cpr_resume: CPR POST USER callback\n"))
1184 (void) callb_execute_class(CB_CL_CPR_POST_USER, CB_CODE_CPR_RESUME);
1185 PMD(PMD_SX, ("cpr_resume: CPR PROMPRINTF callback\n"))
1186 (void) callb_execute_class(CB_CL_CPR_PROMPRINTF, CB_CODE_CPR_RESUME);

1187 PMD(PMD_SX, ("cpr_resume: restore direct levels\n"))
1188 pm_restore_direct_levels();

1189 rb_stop_user_threads:
1190 CPR_DEBUG(CPR_DEBUG1, "starting user threads...");
1191 PMD(PMD_SX, ("cpr_resume: starting user threads\n"))
1192 cpr_start_user_threads();
1193 CPR_DEBUG(CPR_DEBUG1, "done\n");
1194 /*
1195  * Ask Xorg to resume the frame buffer, and wait for it to happen
1196  */
1197 mutex_enter(&srn_clone_lock);
1198 if (srn_signal) {
1199     PMD(PMD_SX, ("cpr_suspend: (*srn_signal)(..., "
1200                 "SRN_NORMAL_RESUME)\n"))
1201     srn_inuse = 1; /* because (*srn_signal) cv_waits */
1202     (*srn_signal)(SRN_TYPE_APM, SRN_NORMAL_RESUME);
1203     srn_inuse = 0;
1204 } else {
1205     PMD(PMD_SX, ("cpr_suspend: srn_signal NULL\n"))
1206 }
1207 mutex_exit(&srn_clone_lock);

1208 #if defined(__sparc)
1209 rb_mp_offline:
1210     if (cpr_mp_online())
1211         cpr_err(CE_WARN, "Failed to online all the processors.");
1212 #endiff

1213 rb_others:
1214 PMD(PMD_SX, ("cpr_resume: dep_thread\n"))
1215 pm_dispatch_to_dep_thread(PM_DEP_WK_CPR_RESUME, NULL, NULL,
1216                           PM_DEP_WAIT, NULL, 0);
1217
1218 PMD(PMD_SX, ("cpr_resume: CPR PM callback\n"))
1219 (void) callb_execute_class(CB_CL_CPR_PM, CB_CODE_CPR_RESUME);

1220 if (cpr_suspend_succeeded) {
1221     cpr_stat_record_events();
1222 }

```

```
1234 #if defined(__sparc)
1235     if (sleepetype == CPR_TODISK && !cpr_reusable_mode)
1236         cpr_clear_definfo();
1237 #endif
1238
1239     i_cpr_free_cpus();
1240     CPR_DEBUG(CPR_DEBUG1, "Sending SIGTHAW...");
1241     PMD(PMD_SX, ("cpr_resume: SIGTHAW\n"))
1242     cpr_signal_user(SIGTHAW);
1243     CPR_DEBUG(CPR_DEBUG1, "done\n");
1244
1245     CPR_STAT_EVENT_END("Resume Total");
1246
1247     CPR_STAT_EVENT_START_TMZ("WHOLE CYCLE", &wholecycle_tv);
1248     CPR_STAT_EVENT_END("WHOLE CYCLE");
1249
1250     if (cpr_debug & CPR_DEBUG1)
1251         cmn_err(CE_CONT, "\nThe system is back where you left!\n");
1252
1253     CPR_STAT_EVENT_START("POST CPR DELAY");
1254
1255 #ifdef CPR_STAT
1256     ctp = &cpr_term.tm_shutdown;
1257     CPR_STAT_EVENT_START_TMZ("PWROFF TIME", ctp);
1258     CPR_STAT_EVENT_END_TMZ("PWROFF TIME", &pwrone_tv);
1259
1260     CPR_STAT_EVENT_PRINT();
1261 #endif /* CPR_STAT */
1262
1263     PMD(PMD_SX, ("cpr_resume returns %x\n", rc))
1264     return (rc);
1265 }
```

unchanged portion omitted

new/usr/src/uts/common/disp/cmt.c

```
*****
5233 Thu May 1 19:50:17 2014
new/usr/src/uts/common/disp/cmt.c
XXXX pass in cpu_pause_func via pause_cpus
*****
_____unchanged_portion_omitted_____
```

```
320 /*
321 * Promote PG above it's current parent.
322 * This is only legal if PG has an equal or greater number of CPUs than its
323 * parent.
324 *
325 * This routine operates on the CPU specific processor group data (for the CPUs
326 * in the PG being promoted), and may be invoked from a context where one CPU's
327 * PG data is under construction. In this case the argument "pgdata", if not
328 * NULL, is a reference to the CPU's under-construction PG data.
329 */
330 static void
331 cmt_hier_promote(pg_cmt_t *pg, cpu_pg_t *pgdata)
332 {
333     pg_cmt_t      *parent;
334     group_t        *children;
335     cpu_t          *cpu;
336     group_iter_t   iter;
337     pg_cpu_itr_t  cpu_iter;
338     int             r;
339     int             err;
340     int             nchildren;
341
342     ASSERT(MUTEX_HELD(&cpu_lock));
343
344     parent = pg->cmt_parent;
345     if (parent == NULL) {
346         /*
347         * Nothing to do
348         */
349         return;
350     }
351
352     ASSERT(PG_NUM_CPUS((pg_t *)pg) >= PG_NUM_CPUS((pg_t *)parent));
353
354     /*
355     * We're changing around the hierarchy, which is actively traversed
356     * by the dispatcher. Pause CPUS to ensure exclusivity.
357     */
358     pause_cpus(NULL, NULL);
359     pause_cpus(NULL);
360
361     /*
362     * If necessary, update the parent's sibling set, replacing parent
363     * with PG.
364     */
365     if (parent->cmt_siblings) {
366         if (group_remove(parent->cmt_siblings, parent, GRP_NORESIZE)
367             != -1) {
368             r = group_add(parent->cmt_siblings, pg, GRP_NORESIZE);
369             ASSERT(r != -1);
370         }
371
372         /*
373         * If the parent is at the top of the hierarchy, replace it's entry
374         * in the root lgroup's group of top level PGs.
375         */
376         if (parent->cmt_parent == NULL &&
377             parent->cmt_siblings != &cmt_root->cl_pgs) {
```

1

```
new/usr/src/uts/common/disp/cmt.c
*****
378     if (group_remove(&cmt_root->cl_pgs, parent, GRP_NORESIZE)
379         != -1) {
380         r = group_add(&cmt_root->cl_pgs, pg, GRP_NORESIZE);
381         ASSERT(r != -1);
382     }
383 }
384
385 /*
386 * We assume (and therefore assert) that the PG being promoted is an
387 * only child of it's parent. Update the parent's children set
388 * replacing PG's entry with the parent (since the parent is becoming
389 * the child). Then have PG and the parent swap children sets and
390 * children counts.
391 */
392 ASSERT(GROUP_SIZE(parent->cmt_children) <= 1);
393 if (group_remove(parent->cmt_children, pg, GRP_NORESIZE) != -1) {
394     r = group_add(parent->cmt_children, parent, GRP_NORESIZE);
395     ASSERT(r != -1);
396 }
397
398 children = pg->cmt_children;
399 pg->cmt_children = parent->cmt_children;
400 parent->cmt_children = children;
401
402 nchildren = pg->cmt_nchildren;
403 pg->cmt_nchildren = parent->cmt_nchildren;
404 parent->cmt_nchildren = nchildren;
405
406 /*
407 * Update the sibling references for PG and it's parent
408 */
409 pg->cmt_siblings = parent->cmt_siblings;
410 parent->cmt_siblings = pg->cmt_children;
411
412 /*
413 * Update any cached lineages in the per CPU pg data.
414 */
415 PG_CPU_ITR_INIT(pg, cpu_iter);
416 while ((cpu = pg_cpu_next(&cpu_iter)) != NULL) {
417     int           idx;
418     int           sz;
419     pg_cmt_t     *cpu_pg;
420     cpu_pg_t     *pgd; /* CPU's PG data */
421
422     /*
423     * The CPU's whose lineage is under construction still
424     * references the bootstrap CPU PG data structure.
425     */
426     if (pg_cpu_is_bootstrapped(cpu))
427         pgd = pgdata;
428     else
429         pgd = cpu->cpu_pg;
430
431     /*
432     * Iterate over the CPU's PGs updating the children
433     * of the PG being promoted, since they have a new parent.
434     */
435     group_iter_init(&iter);
436     while ((cpu_pg = group_iterate(&pgd->cmt_pgs, &iter)) != NULL) {
437         if (cpu_pg->cmt_parent == pg) {
438             cpu_pg->cmt_parent = parent;
439         }
440     }
441
442     /*
443     * Update the CMT load balancing lineage
```

2

```

444         */
445         if ((idx = group_find(&pgd->cmt_pgs, (void *)pg)) == -1) {
446             /*
447              * Unless this is the CPU who's lineage is being
448              * constructed, the PG being promoted should be
449              * in the lineage.
450             */
451             ASSERT(pg_cpu_is_bootstrapped(cpu));
452             continue;
453         }
454
455         ASSERT(idx > 0);
456         ASSERT(GROUP_ACCESS(&pgd->cmt_pgs, idx - 1) == parent);
457
458         /*
459          * Have the child and the parent swap places in the CPU's
460          * lineage
461         */
462         group_remove_at(&pgd->cmt_pgs, idx);
463         group_remove_at(&pgd->cmt_pgs, idx - 1);
464         err = group_add_at(&pgd->cmt_pgs, parent, idx);
465         ASSERT(err == 0);
466         err = group_add_at(&pgd->cmt_pgs, pg, idx - 1);
467         ASSERT(err == 0);
468
469         /*
470          * Ensure cmt_lineage references CPU's leaf PG.
471          * Since cmt_pgs is top-down ordered, the bottom is the last
472          * element.
473         */
474         if ((sz = GROUP_SIZE(&pgd->cmt_pgs)) > 0)
475             pgd->cmt_lineage = GROUP_ACCESS(&pgd->cmt_pgs, sz - 1);
476     }
477
478     /*
479      * Update the parent references for PG and it's parent
480      */
481     pg->cmt_parent = parent->cmt_parent;
482     parent->cmt_parent = pg;
483
484     start_cpus();
485 }


---


unchanged_portion_omitted
1455 /*
1456  * Prune PG, and all other instances of PG's hardware sharing relationship
1457  * from the CMT PG hierarchy.
1458  *
1459  * This routine operates on the CPU specific processor group data (for the CPUs
1460  * in the PG being pruned), and may be invoked from a context where one CPU's
1461  * PG data is under construction. In this case the argument "pgdata", if not
1462  * NULL, is a reference to the CPU's under-construction PG data.
1463  */
1464 static int
1465 pg_cmt_prune(pg_cmt_t *pg_bad, pg_cmt_t **lineage, int *sz, cpu_pg_t *pgdata)
1466 {
1467     group_t        *hwset, *children;
1468     int            i, j, r, size = *sz;
1469     group_iter_t   hw_iter, child_iter;
1470     pg_cpu_iter_t cpu_iter;
1471     pg_cmt_t      *pg, *child;
1472     cpu_t          *cpu;
1473     int            cap_needed;
1474     pghw_type_t    hw;
1475
1476     ASSERT(MUTEX_HELD(&cpu_lock));

```

```

1478     /*
1479      * Inform pghw layer that this PG is pruned.
1480      */
1481     pghw_cmt_fini((pghw_t *)pg_bad);
1482
1483     hw = ((pghw_t *)pg_bad)->pghw_hw;
1484
1485     if (hw == PGHW_POW_ACTIVE) {
1486         cmn_err(CE_NOTE, "!Active CPUPM domain groups look suspect. "
1487                 "Event Based CPUPM Unavailable");
1488     } else if (hw == PGHW_POW_IDLE) {
1489         cmn_err(CE_NOTE, "!Idle CPUPM domain groups look suspect. "
1490                 "Dispatcher assisted CPUPM disabled.");
1491     }
1492
1493     /*
1494      * Find and eliminate the PG from the lineage.
1495      */
1496     for (i = 0; i < size; i++) {
1497         if (lineage[i] == pg_bad) {
1498             for (j = i; j < size - 1; j++)
1499                 lineage[j] = lineage[j + 1];
1500             *sz = size - 1;
1501             break;
1502         }
1503     }
1504
1505     /*
1506      * We'll prune all instances of the hardware sharing relationship
1507      * represented by pg. But before we do that (and pause CPUs) we need
1508      * to ensure the hierarchy's groups are properly sized.
1509      */
1510     hwset = pghw_set_lookup(hw);
1511
1512     /*
1513      * Blacklist the hardware so future processor groups of this type won't
1514      * participate in CMT thread placement.
1515      *
1516      * XXX
1517      * For heterogeneous system configurations, this might be overkill.
1518      * We may only need to blacklist the illegal PGs, and other instances
1519      * of this hardware sharing relationship may be ok.
1520      */
1521     cmt_hw_blacklisted[hw] = 1;
1522
1523     /*
1524      * For each of the PGs being pruned, ensure sufficient capacity in
1525      * the siblings set for the PG's children
1526      */
1527     group_iter_init(&hw_iter);
1528     while ((pg = group_iterate(hwset, &hw_iter)) != NULL) {
1529         /*
1530          * PG is being pruned, but if it is bringing up more than
1531          * one child, ask for more capacity in the siblings group.
1532          */
1533         cap_needed = 0;
1534         if (pg->cmt_children &&
1535             GROUP_SIZE(pg->cmt_children) > 1) {
1536             cap_needed = GROUP_SIZE(pg->cmt_children) - 1;
1537             group_expand(pg->cmt_siblings,
1538                         GROUP_SIZE(pg->cmt_siblings) + cap_needed);
1539
1540             /*
1541              * If this is a top level group, also ensure the

```

```

1543             * capacity in the root lgrp level CMT grouping.
1544             */
1545         if (pg->cmt_parent == NULL &&
1546             pg->cmt_siblings != &cmt_root->cl_pgs) {
1547             group_expand(&cmt_root->cl_pgs,
1548                         GROUP_SIZE(&cmt_root->cl_pgs) + cap_needed);
1549             cmt_root->cl_npgs += cap_needed;
1550         }
1551     }
1552 }
1553
1554 /*
1555  * We're operating on the PG hierarchy. Pause CPUs to ensure
1556  * exclusivity with respect to the dispatcher.
1557 */
1558 pause_cpus(NULL, NULL);
1559 pause_cpus(NULL);
1560
1561 /*
1562  * Prune all PG instances of the hardware sharing relationship
1563  * represented by pg.
1564 */
1565 group_iter_init(&hw_iter);
1566 while ((pg = group_iterate(hwset, &hw_iter)) != NULL) {
1567
1568     /*
1569      * Remove PG from it's group of siblings, if it's there.
1570      */
1571     if (pg->cmt_siblings) {
1572         (void) group_remove(pg->cmt_siblings, pg, GRP_NORESIZE);
1573     }
1574     if (pg->cmt_parent == NULL &&
1575         pg->cmt_siblings != &cmt_root->cl_pgs) {
1576         (void) group_remove(&cmt_root->cl_pgs, pg,
1577                             GRP_NORESIZE);
1578     }
1579
1580     /*
1581      * Indicate that no CMT policy will be implemented across
1582      * this PG.
1583      */
1584     pg->cmt_policy = CMT_NO_POLICY;
1585
1586     /*
1587      * Move PG's children from it's children set to it's parent's
1588      * children set. Note that the parent's children set, and PG's
1589      * siblings set are the same thing.
1590      *
1591      * Because we are iterating over the same group that we are
1592      * operating on (removing the children), first add all of PG's
1593      * children to the parent's children set, and once we are done
1594      * iterating, empty PG's children set.
1595      */
1596     if (pg->cmt_children != NULL) {
1597         children = pg->cmt_children;
1598
1599         group_iter_init(&child_iter);
1600         while ((child = group_iterate(children, &child_iter))
1601                != NULL) {
1602             if (pg->cmt_siblings != NULL) {
1603                 r = group_add(pg->cmt_siblings, child,
1604                               GRP_NORESIZE);
1605                 ASSERT(r == 0);
1606
1607             if (pg->cmt_parent == NULL &&
1608                 pg->cmt_siblings !=

```

```

1608                 &cmt_root->cl_pgs) {
1609                 r = group_add(&cmt_root->cl_pgs,
1610                               child, GRP_NORESIZE);
1611                 ASSERT(r == 0);
1612             }
1613         }
1614     }
1615     group_empty(pg->cmt_children);
1616 }
1617
1618 /*
1619  * Reset the callbacks to the defaults
1620  */
1621 pg_callback_set_defaults((pg_t *)pg);
1622
1623 /*
1624  * Update all the CPU lineages in each of PG's CPUs
1625  */
1626 PG_CPU_ITR_INIT(pg, cpu_iter);
1627 while ((cpu = pg_cpu_next(&cpu_iter)) != NULL) {
1628     pg_cmt_t          *cpu_pg;
1629     group_iter_t       liter; /* Iterator for the lineage */
1630     cpu_pg_t          *cpd;   /* CPU's PG data */
1631
1632     /*
1633      * The CPU's lineage is under construction still
1634      * references the bootstrap CPU PG data structure.
1635      */
1636     if (pg_cpu_is_bootstrapped(cpu))
1637         cpd = pgdata;
1638     else
1639         cpd = cpu->cpu_pg;
1640
1641     /*
1642      * Iterate over the CPU's PGs updating the children
1643      * of the PG being promoted, since they have a new
1644      * parent and siblings set.
1645      */
1646     group_iter_init(&liter);
1647     while ((cpu_pg = group_iterate(&cpd->pgs,
1648                                   &liter)) != NULL) {
1649         if (cpu_pg->cmt_parent == pg) {
1650             cpu_pg->cmt_parent = pg->cmt_parent;
1651             cpu_pg->cmt_siblings = pg->cmt_siblings;
1652         }
1653     }
1654
1655     /*
1656      * Update the CPU's lineages
1657      *
1658      * Remove the PG from the CPU's group used for CMT
1659      * scheduling.
1660      */
1661     (void) group_remove(&cpd->cmt_pgs, pg, GRP_NORESIZE);
1662 }
1663 start_cpus();
1664 return (0);
1665 }
1666
1667 /*
1668  * Disable CMT scheduling
1669  */
1670 static void
1671 pg_cmt_disable(void)
1672 {

```

```
1674     cpu_t          *cpu;
1676     ASSERT(MUTEX_HELD(&cpu_lock));
1678     pause_cpus(NULL, NULL);
1679     pause_cpus(NULL);
1680     cpu = cpu_list;
1681     do {
1682         if (cpu->cpu_pg)
1683             group_empty(&cpu->cpu_pg->cmt_pgs);
1684     } while ((cpu = cpu->cpu_next) != cpu_list);
1685     cmt_sched_disabled = 1;
1686     start_cpus();
1687     cmn_err(CE_NOTE, "!CMT thread placement optimizations unavailable");
1688 }


---

unchanged portion omitted
```

```
new/usr/src/uts/common/disp/cpupart.c
```

```
*****
30551 Thu May 1 19:50:17 2014
new/usr/src/uts/common/disp/cpupart.c
XXXX pass in cpu_pause_func via pause_cpus
*****
_____unchanged_portion_omitted_____
```

```
319 static int
320 cpupart_move_cpu(cpu_t *cp, cpupart_t *newpp, int forced)
321 {
322     cpupart_t *oldpp;
323     cpu_t *ncp, *newlist;
324     kthread_t *t;
325     int move_threads = 1;
326     lgrp_id_t lgrp_id;
327     proc_t *p;
328     int lgrp_diff_lpl;
329     lpl_t *cpu_lpl;
330     int ret;
331     boolean_t unbind_all_threads = (forced != 0);

333     ASSERT(MUTEX_HELD(&cpu_lock));
334     ASSERT(newpp != NULL);

336     oldpp = cp->cpu_part;
337     ASSERT(oldpp != NULL);
338     ASSERT(oldpp->cp_ncpus > 0);

340     if (newpp == oldpp) {
341         /*
342          * Don't need to do anything.
343         */
344         return (0);
345     }

347     cpu_state_change_notify(cp->cpu_id, CPU_CPUPART_OUT);

349     if (!disp_bound_partition(cp, 0)) {
350         /*
351          * Don't need to move threads if there are no threads in
352          * the partition. Note that threads can't enter the
353          * partition while we're holding cpu_lock.
354         */
355         move_threads = 0;
356     } else if (oldpp->cp_ncpus == 1) {
357         /*
358          * The last CPU is removed from a partition which has threads
359          * running in it. Some of these threads may be bound to this
360          * CPU.
361          *
362          * Attempt to unbind threads from the CPU and from the processor
363          * set. Note that no threads should be bound to this CPU since
364          * cpupart_move_threads will refuse to move bound threads to
365          * other CPUs.
366         */
367         (void) cpu_unbind(oldpp->cp_cpulist->cpu_id, B_FALSE);
368         (void) cpupart_unbind_threads(oldpp, B_FALSE);

370         if (!disp_bound_partition(cp, 0)) {
371             /*
372              * No bound threads in this partition any more
373              */
374             move_threads = 0;
375         } else {
376             /*
377             */
378         }
379     }
380     if (move_threads) {
381         /*
382          * If forced flag is set unbind any threads from this CPU.
383          * Otherwise unbind soft-bound threads only.
384          */
385         if ((ret = cpu_unbind(cp->cpu_id, unbind_all_threads)) != 0) {
386             cpu_state_change_notify(cp->cpu_id, CPU_CPUPART_IN);
387             return (ret);
388         }
389         /*
390          * Stop further threads weak binding to this cpu.
391          */
392         cpu_inmotion = cp;
393         membar_enter();
394         /*
395          * Notify the Processor Groups subsystem that the CPU
396          * will be moving cpu partitions. This is done before
397          * CPUs are paused to provide an opportunity for any
398          * needed memory allocations.
399          */
400         pg_cpupart_out(cp, oldpp);
401         pg_cpupart_in(cp, newpp);
402     }
403     again:
404     if (move_threads) {
405         int loop_count;
406         /*
407          * Check for threads strong or weak bound to this CPU.
408          */
409         for (loop_count = 0; disp_bound_threads(cp, 0); loop_count++) {
410             if (loop_count >= 5) {
411                 cpu_state_change_notify(cp->cpu_id,
412                                         CPU_CPUPART_IN);
413                 pg_cpupart_out(cp, newpp);
414                 pg_cpupart_in(cp, oldpp);
415                 cpu_inmotion = NULL;
416                 return (EBUSY); /* some threads still bound */
417             }
418             delay(1);
419         }
420     }
421     /*
422      * Before we actually start changing data structures, notify
423      * the cyclic subsystem that we want to move this CPU out of its
424      * partition.
425      */
426     if (!cyclic_move_out(cp)) {
427         /*
428          * This CPU must be the last CPU in a processor set with
429          * a bound cyclic.
430          */
431         cpu_state_change_notify(cp->cpu_id, CPU_CPUPART_IN);
432         pg_cpupart_out(cp, newpp);
433         pg_cpupart_in(cp, oldpp);
434         cpu_inmotion = NULL;
435         return (EBUSY);
436     }
437 }
```

```
1
```

```
new/usr/src/uts/common/disp/cpupart.c
*****
377     /*
378      * There are still threads bound to the partition
379      */
380     cpu_state_change_notify(cp->cpu_id, CPU_CPUPART_IN);
381     return (EBUSY);
382 }
383 */
384 /*
385  * If forced flag is set unbind any threads from this CPU.
386  * Otherwise unbind soft-bound threads only.
387  */
388 if ((ret = cpu_unbind(cp->cpu_id, unbind_all_threads)) != 0) {
389     cpu_state_change_notify(cp->cpu_id, CPU_CPUPART_IN);
390     return (ret);
391 }
392 /*
393  * Stop further threads weak binding to this cpu.
394  */
395 cpu_inmotion = cp;
396 membar_enter();
397 /*
398  * Notify the Processor Groups subsystem that the CPU
399  * will be moving cpu partitions. This is done before
400  * CPUs are paused to provide an opportunity for any
401  * needed memory allocations.
402  */
403 pg_cpupart_out(cp, oldpp);
404 pg_cpupart_in(cp, newpp);
405
406 again:
407 if (move_threads) {
408     int loop_count;
409     /*
410      * Check for threads strong or weak bound to this CPU.
411      */
412     for (loop_count = 0; disp_bound_threads(cp, 0); loop_count++) {
413         if (loop_count >= 5) {
414             cpu_state_change_notify(cp->cpu_id,
415                                     CPU_CPUPART_IN);
416             pg_cpupart_out(cp, newpp);
417             pg_cpupart_in(cp, oldpp);
418             cpu_inmotion = NULL;
419             return (EBUSY); /* some threads still bound */
420         }
421         delay(1);
422     }
423 }
424 /*
425  * Before we actually start changing data structures, notify
426  * the cyclic subsystem that we want to move this CPU out of its
427  * partition.
428  */
429 if (!cyclic_move_out(cp)) {
430     /*
431      * This CPU must be the last CPU in a processor set with
432      * a bound cyclic.
433      */
434     cpu_state_change_notify(cp->cpu_id, CPU_CPUPART_IN);
435     pg_cpupart_out(cp, newpp);
436     pg_cpupart_in(cp, oldpp);
437     cpu_inmotion = NULL;
438     return (EBUSY);
439 }
```

```
2
```



```

574         }
575         t = t->t_forw;
576     } while (t != p->p_tlist);

578     /*
579      * Didn't find any threads in the same lgroup as this
580      * CPU with a different lpl, so remove the lgroup from
581      * the process lgroup bitmask.
582      */
583
584     if (lgrp_diff_lpl)
585         klgrpset_del(p->p_lgrpset, lgrp_id);
586 }

588     /*
589      * Walk thread list looking for threads that need to be
590      * rehomed, since there are some threads that are not in
591      * their process's p_tlist.
592      */
593
594     t = curthread;

596 do {
597     ASSERT(t != NULL && t->t_lpl != NULL);

599     /*
600      * If the lgroup that t is assigned to no
601      * longer has any CPUs in t's partition,
602      * we'll have to choose a new lgroup for t.
603      * Also, choose best lgroup for home when
604      * thread has specified lgroup affinities,
605      * since there may be an lgroup with more
606      * affinity available after moving CPUs
607      * around.
608      */
609     if (!LGRP_CPUS_IN_PART(t->t_lpl->lpl_lgrp_id,
610                           t->t_cpupart) || t->t_lgrp_affinity) {
611         lgrp_move_thread(t,
612                           lgrp_choose(t, t->t_cpupart), 1);
613     }

615     /* make sure lpl points to our own partition */
616     ASSERT((t->t_lpl >= t->t_cpupart->cp_lgrploads) &&
617            (t->t_lpl < t->t_cpupart->cp_lgrploads +
618             t->t_cpupart->cp_nlgrploads));
619
620     ASSERT(t->t_lpl->lpl_ncpu > 0);

622     /* Update CPU last ran on if it was this CPU */
623     if (t->t_cpu == cp && t->t_cpupart == oldpp &&
624        t->t_bound_cpu != cp) {
625         t->t_cpu = disp_lowpri_cpu(ncp, t->t_lpl,
626                                    t->t_pri, NULL);
627     }

629         t = t->t_next;
630     } while (t != curthread);

632     /*
633      * Clear off the CPU's run queue, and the kp queue if the
634      * partition is now empty.
635      */
636     disp_cpu_inactive(cp);

638     /*
639      * Make cp switch to a thread from the new partition.

```

```

640         */
641         cp->cpu_rnrun = 1;
642         cp->cpu_kprnrun = 1;
643     }

645     cpu_inmotion = NULL;
646     start_cpus();

648     /*
649      * Let anyone interested know that cpu has been added to the set.
650      */
651     cpu_state_change_notify(cp->cpu_id, CPU_CPUPART_IN);

653     /*
654      * Now let the cyclic subsystem know that it can reshuffle cyclics
655      * bound to the new processor set.
656      */
657     cyclic_move_in(cp);

659     return (0);
660 }



---


unchanged_portion_omitted

812 /*
813  * Create a new partition. On MP systems, this also allocates a
814  * kpreempt disp queue for that partition.
815  */
816 int
817 cpupart_create(psetid_t *psid)
818 {
819     cpupart_t      *pp;
820
821     ASSERT(pool_lock_held());

823     pp = kmalloc(sizeof (cpupart_t), KM_SLEEP);
824     pp->cp_nlgrploads = lgrp_plat_max_lgrps();
825     pp->cp_lgrploads = kmalloc(sizeof (lpl_t) * pp->cp_nlgrploads,
826                                KM_SLEEP);

828     mutex_enter(&cpu_lock);
829     if (cp_numparts == cp_max_numparts) {
830         mutex_exit(&cpu_lock);
831         kmem_free(pp->cp_lgrploads, sizeof (lpl_t) * pp->cp_nlgrploads);
832         pp->cp_lgrploads = NULL;
833         kmem_free(pp, sizeof (cpupart_t));
834         return (ENOMEM);
835     }
836     cp_numparts++;
837     /* find the next free partition ID */
838     while (cpupart_find(CPTOPS(cp_id_next)) != NULL)
839         cp_id_next++;
840     pp->cp_id = cp_id_next++;
841     pp->cp_ncpus = 0;
842     pp->cp_cpulist = NULL;
843     pp->cp_attr = 0;
844     klgrpset_clear(pp->cp_lgrpset);
845     pp->cp_kp_queue.disp_maxrunpri = -1;
846     pp->cp_kp_queue.disp_max_unbound_pri = -1;
847     pp->cp_kp_queue.disp_cpu = NULL;
848     pp->cp_gen = 0;
849     DISP_LOCK_INIT(&pp->cp_kp_queue.disp_lock);
850     *psid = CPTOPS(pp->cp_id);
851     disp_kp_alloc(&pp->cp_kp_queue, v.v_nglobpris);
852     cpupart_kstat_create(pp);
853     cpupart_lpl_initialize(pp);

```

```

855     bitset_init(&pp->cp_cmt_pgs);
857     /*
858      * Initialize and size the partition's bitset of halted CPUs.
859      */
860     bitset_init_fanout(&pp->cp_haltset, cp_haltset_fanout);
861     bitset_resize(&pp->cp_haltset, max_ncpus);

863     /*
864      * Pause all CPUs while changing the partition list, to make sure
865      * the clock thread (which traverses the list without holding
866      * cpu_lock) isn't running.
867      */
868     pause_cpus(NULL, NULL);
868     pause_cpus(NULL);
869     pp->cp_next = cp_list_head;
870     pp->cp_prev = cp_list_head->cp_prev;
871     cp_list_head->cp_prev->cp_next = pp;
872     cp_list_head->cp_prev = pp;
873     start_cpus();
874     mutex_exit(&cpu_lock);

876     return (0);
877 }



---



unchanged_portion_omitted



949 */
950 * Destroy a partition.
951 */
952 int
953 cpupart_destroy(psetid_t psid)
954 {
955     cpu_t    *cp, *first_cp;
956     cpupart_t *pp, *newpp;
957     int       err = 0;

959     ASSERT(pool_lock_held());
960     mutex_enter(&cpu_lock);

962     pp = cpupart_find(psid);
963     if (pp == NULL || pp == &cp_default) {
964         mutex_exit(&cpu_lock);
965         return (EINVAL);
966     }

968     /*
969      * Unbind all the threads currently bound to the partition.
970      */
971     err = cpupart_unbind_threads(pp, B_TRUE);
972     if (err) {
973         mutex_exit(&cpu_lock);
974         return (err);
975     }

977     newpp = &cp_default;
978     while ((cp = pp->cp_cpulist) != NULL) {
979         if (err = cpupart_move_cpu(cp, newpp, 0)) {
980             mutex_exit(&cpu_lock);
981             return (err);
982         }
983     }

985     ASSERT(bitset_is_null(&pp->cp_cmt_pgs));
986     ASSERT(bitset_is_null(&pp->cp_haltset));

```

```

988     /*
989      * Teardown the partition's group of active CMT PGs and halted
990      * CPUs now that they have all left.
991      */
992     bitset_fini(&pp->cp_cmt_pgs);
993     bitset_fini(&pp->cp_haltset);

995     /*
996      * Reset the pointers in any offline processors so they won't
997      * try to rejoin the destroyed partition when they're turned
998      * online.
999      */
1000    first_cp = cp = CPU;
1001    do {
1002        if (cp->cpu_part == pp) {
1003            ASSERT(cp->cpu_flags & CPU_OFFLINE);
1004            cp->cpu_part = newpp;
1005        }
1006        cp = cp->cpu_next;
1007    } while (cp != first_cp);

1009    /*
1010     * Pause all CPUs while changing the partition list, to make sure
1011     * the clock thread (which traverses the list without holding
1012     * cpu_lock) isn't running.
1013     */
1014     pause_cpus(NULL, NULL);
1014     pause_cpus(NULL);
1015     pp->cp_prev->cp_next = pp->cp_next;
1016     pp->cp_next->cp_prev = pp->cp_prev;
1017     if (cp_list_head == pp)
1018         cp_list_head = pp->cp_next;
1019     start_cpus();

1021     if (cp_id_next > pp->cp_id)
1022         cp_id_next = pp->cp_id;

1024     if (pp->cp_kstat)
1025         kstat_delete(pp->cp_kstat);

1027     cp_numparts--;
1029     disp_kp_free(&pp->cp_kp_queue);
1031     cpupart_lpl_teardown(pp);

1033     kmem_free(pp, sizeof (cpupart_t));
1034     mutex_exit(&cpu_lock);

1036     return (err);
1037 }



---



unchanged_portion_omitted


```

```
*****  
 70612 Thu May 1 19:50:18 2014  
new/usr/src/uts/common/disp/disp.c  
XXXX pass in cpu_pause_func via pause_cpus  
*****  
unchanged_portion_omitted  
  
312 /*  
313  * For each CPU, allocate new dispatch queues  
314  * with the stated number of priorities.  
315 */  
316 static void  
317 cpu_dispqalloc(int numpris)  
318 {  
319     cpu_t    *cpup;  
320     struct disp_queue_info  *disp_mem;  
321     int i, num;  
322  
323     ASSERT(MUTEX_HELD(&cpu_lock));  
324  
325     disp_mem = kmem_zalloc(NCPU *  
326                           sizeof (struct disp_queue_info), KM_SLEEP);  
327  
328     /*  
329      * This routine must allocate all of the memory before stopping  
330      * the cpus because it must not sleep in kmem_alloc while the  
331      * CPUs are stopped. Locks they hold will not be freed until they  
332      * are restarted.  
333      */  
334     i = 0;  
335     cpup = cpu_list;  
336     do {  
337         disp_dq_alloc(&disp_mem[i], numpris, cpup->cpu_disp);  
338         i++;  
339         cpup = cpup->cpu_next;  
340     } while (cpup != cpu_list);  
341     num = i;  
342  
343     pause_cpus(NULL, NULL);  
344     pause_cpus(NULL);  
345     for (i = 0; i < num; i++)  
346         disp_dq_assign(&disp_mem[i], numpris);  
347     start_cpus();  
348  
349     /*  
350      * I must free all of the memory after starting the cpus because  
351      * I can not risk sleeping in kmem_free while the cpus are stopped.  
352      */  
353     for (i = 0; i < num; i++)  
354         disp_dq_free(&disp_mem[i]);  
355  
356 }  
unchanged_portion_omitted
```

```
*****
94655 Thu May 1 19:50:18 2014
new/usr/src/uts/common/os/cpu.c
XXXX pass in cpu_pause_func via pause_cpus
*****
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 */

25 /*
26 * Architecture-independent CPU control functions.
27 */
28 */

29 #include <sys/types.h>
30 #include <sys/param.h>
31 #include <sys/var.h>
32 #include <sys/thread.h>
33 #include <sys/cpuvar.h>
34 #include <sys/cpu_event.h>
35 #include <sys/kstat.h>
36 #include <sys/uadmin.h>
37 #include <sys/uadmin.h>
38 #include <sys/sysm.h>
39 #include <sys/errno.h>
40 #include <sys/cmn_err.h>
41 #include <sys/procset.h>
42 #include <sys/processor.h>
43 #include <sys/debug.h>
44 #include <sys/cpupart.h>
45 #include <sys/lgrp.h>
46 #include <sys/pset.h>
47 #include <sys/pghw.h>
48 #include <sys/kmem.h>
49 #include <sys/kmem_impl.h>
50 #include <sys/atomic.h>
51 #include <sys/callb.h>
52 #include <sys/vtrace.h>
53 #include <sys/cyclic.h>
54 #include <sys/bitmap.h>
55 #include <sys/nvpair.h>
56 #include <sys/pool_pset.h>
57 #include <sys/msacct.h>
58 #include <sys/time.h>
59 #include <sys/archsysm.h>
60 #include <sys/sdt.h>
61 #if defined(__x86) || defined(__amd64)
```

/* to set per-cpu kmem_cache offset */

```
62 #include <sys/x86_archext.h>
63 #endif
64 #include <sys/callo.h>

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

73 extern int cmp_cpu_to_chip(processorid_t cpuid);
74 #ifdef __sparcv9
75 extern char *cpu_fru_fmri(cpu_t *cp);
76#endif

77 static void cpu_add_active_internal(cpu_t *cp);
78 static void cpu_remove_active(cpu_t *cp);
79 static void cpu_info_kstat_create(cpu_t *cp);
80 static void cpu_info_kstat_destroy(cpu_t *cp);
81 static void cpu_stats_kstat_create(cpu_t *cp);
82 static void cpu_stats_kstat_destroy(cpu_t *cp);

83 static int cpu_sys_stats_ks_update(kstat_t *ksp, int rw);
84 static int cpu_vm_stats_ks_update(kstat_t *ksp, int rw);
85 static int cpu_stat_ks_update(kstat_t *ksp, int rw);
86 static int cpu_state_change_hooks(int, cpu_setup_t, cpu_setup_t);

87 static int
88 /* cpu_lock protects ncpus, ncpus_online, cpu_flag, cpu_list, cpu_active,
89 * max_cpu_seqid_ever, and dispatch queue reallocations. The lock ordering with
90 * respect to related locks is:
91 *
92 *      cpu_lock --> thread_free_lock ---> p_lock ---> thread_lock()
93 *
94 * Warning: Certain sections of code do not use the cpu_lock when
95 * traversing the cpu_list (e.g. mutex_vector_enter(), clock()). Since
96 * all cpus are paused during modifications to this list, a solution
97 * to protect the list is too either disable kernel preemption while
98 * walking the list, *or* recheck the cpu_next pointer at each
99 * iteration in the loop. Note that in no cases can any cached
100 * copies of the cpu pointers be kept as they may become invalid.
101 */
102 kmutex_t cpu_lock;
103 cpu_t *cpu_list; /* list of all CPUs */
104 cpu_t *clock_cpu_list; /* used by clock to walk CPUs */
105 cpu_t *cpu_active; /* list of active CPUs */
106 static cpuset_t cpu_available; /* set of available CPUs */
107 cpuset_t cpu_seqid_inuse; /* which cpu_seqids are in use */

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

109 /* max_ncpus keeps the max cpus the system can have. Initially
110 * it's NCPU, but since most archs scan the devtree for cpus
111 * fairly early on during boot, the real max can be known before
112 * ncpus is set (useful for early NCPU based allocations).
113 */
114 int max_ncpus = NCPU;
115 /* platforms that set max_ncpus to maximum number of cpus that can be
116 * dynamically added will set boot_max_ncpus to the number of cpus found
117 * at device tree scan time during boot.
118 */
119 /*
```

```

128 int boot_max_ncpus = -1;
129 int boot_ncpus = -1;
130 /*
131  * Maximum possible CPU id. This can never be >= NCPU since NCPU is
132  * used to size arrays that are indexed by CPU id.
133  */
134 processorid_t max_cpuid = NCPU - 1;

136 /*
137  * Maximum cpu_seqid was given. This number can only grow and never shrink. It
138  * can be used to optimize NCPU loops to avoid going through CPUs which were
139  * never on-line.
140 */
141 processorid_t max_cpu_seqid_ever = 0;

143 int ncpus = 1;
144 int ncpus_online = 1;

146 /*
147  * CPU that we're trying to offline. Protected by cpu_lock.
148 */
149 cpu_t *cpu_inmotion;

151 /*
152  * Can be raised to suppress further weakbinding, which are instead
153  * satisfied by disabling preemption. Must be raised/lowered under cpu_lock,
154  * while individual thread weakbinding synchronization is done under thread
155  * lock.
156 */
157 int weakbindingbarrier;

159 /*
160  * Variables used in pause_cpus().
161 */
162 static volatile char safe_list[NCPU];

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

174 static kmutex_t pause_free_mutex;
175 static kcondvar_t pause_free_cv;

176 void *(*cpu_pause_func)(void *) = NULL;

178 static struct cpu_sys_stats_ks_data {
179     kstat_named_t cpu_ticks_idle;
180     kstat_named_t cpu_ticks_user;
181     kstat_named_t cpu_ticks_kernel;
182     kstat_named_t cpu_ticks_wait;
183     kstat_named_t cpu_nsec_idle;
184     kstat_named_t cpu_nsec_user;
185     kstat_named_t cpu_nsec_kernel;
186     kstat_named_t cpu_nsec_dtrace;
187     kstat_named_t cpu_nsec_intr;
188     kstat_named_t cpu_load_intr;
189     kstat_named_t wait_ticks_io;
190     kstat_named_t dtrace_probes;
191     kstat_named_t bread;

```

```

192     kstat_named_t bwrite;
193     kstat_named_t lread;
194     kstat_named_t lwrite;
195     kstat_named_t phread;
196     kstat_named_t phwrite;
197     kstat_named_t pswitch;
198     kstat_named_t trap;
199     kstat_named_t intr;
200     kstat_named_t syscall;
201     kstat_named_t sysread;
202     kstat_named_t syswrite;
203     kstat_named_t sysfork;
204     kstat_named_t sysvfork;
205     kstat_named_t sysexec;
206     kstat_named_t readch;
207     kstat_named_t writech;
208     kstat_named_t rcvint;
209     kstat_named_t xmtint;
210     kstat_named_t mdmint;
211     kstat_named_t rawch;
212     kstat_named_t canch;
213     kstat_named_t outhch;
214     kstat_named_t msg;
215     kstat_named_t sema;
216     kstat_named_t namei;
217     kstat_named_t ufsiget;
218     kstat_named_t ufsirblk;
219     kstat_named_t ufsipage;
220     kstat_named_t ufsinopage;
221     kstat_named_t procovf;
222     kstat_named_t intrthread;
223     kstat_named_t intrblk;
224     kstat_named_t intrunpin;
225     kstat_named_t idlethread;
226     kstat_named_t inv_swch;
227     kstat_named_t nthreads;
228     kstat_named_t cpumigrate;
229     kstat_named_t xcalls;
230     kstat_named_t mutex_adenters;
231     kstat_named_t rw_rdfails;
232     kstat_named_t rw_wrfails;
233     kstat_named_t modload;
234     kstat_named_t modunload;
235     kstat_named_t bawrite;
236     kstat_named_t iowait;
237 } cpu_sys_stats_ks_data_template = {
238     unchanged_portion_omitted_
239 };

251 /*
252  * This routine is called to place the CPUs in a safe place so that
253  * one of them can be taken off line or placed on line. What we are
254  * trying to do here is prevent a thread from traversing the list
255  * of active CPUs while we are changing it or from getting placed on
256  * the run queue of a CPU that has just gone off line. We do this by
257  * creating a thread with the highest possible prio for each CPU and
258  * having it call this routine. The advantage of this method is that
259  * we can eliminate all checks for CPU_ACTIVE in the disp routines.
260  * This makes disp faster at the expense of making p_online() slower
261  * which is a good trade off.
262 */
263 static void
264 cpu_pause(int index)
265 {
266     int s;
267     struct _cpu_pause_info *cpi = &cpu_pause_info;
268     volatile char *safe = &safe_list[index];

```

```

769     long    lindex = index;
771
772     ASSERT((curthread->t_bound_cpu != NULL) || (*safe == PAUSE_DIE));
773
774     while (*safe != PAUSE_DIE) {
775         *safe = PAUSE_READY;
776         membar_enter(); /* make sure stores are flushed */
777         sema_v(&cpi->cp_sem); /* signal requesting thread */
778
779         /*
780          * Wait here until all pause threads are running. That
781          * indicates that it's safe to do the spl. Until
782          * cpu_pause_info.cp_go is set, we don't want to spl
783          * because that might block clock interrupts needed
784          * to preempt threads on other CPUs.
785        */
786     while (cpi->cp_go == 0)
787         ;
788
789     /*
790      * Even though we are at the highest disp prio, we need
791      * to block out all interrupts below LOCK_LEVEL so that
792      * an intr doesn't come in, wake up a thread, and call
793      * setbackdq/setfrontdq.
794    */
795     s = splhigh();
796
797     /*
798      * if cp_func has been set then call it using index as the
799      * argument, currently only used by cpr_suspend_cpus().
800      * This function is used as the code to execute on the
801      * "paused" cpu's when a machine comes out of a sleep state
802      * and CPU's were powered off. (could also be used for
803      * hotplugging CPU's).
804      * if cpu_pause_func() has been set then call it using
805      * index as the argument, currently only used by
806      * cpr_suspend_cpus(). This function is used as the
807      * code to execute on the "paused" cpu's when a machine
808      * comes out of a sleep state and CPU's were powered off.
809      * (could also be used for hotplugging CPU's).
810    */
811
812     if (cpi->cp_func != NULL)
813         (*cpi->cp_func)((void *)lindex);
814     if (cpu_pause_func != NULL)
815         (*cpu_pause_func)((void *)lindex);
816
817     mach_cpu_pause(safe);
818
819     mutex_enter(&pause_free_mutex);
820     *safe = PAUSE_DEAD;
821     cv_broadcast(&pause_free_cv);
822     mutex_exit(&pause_free_mutex);
823 }
824
825 unchanged_portion_omitted
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 * Pause all of the CPUs except the one we are on by creating a high
876 * priority thread bound to those CPUs.
877 */

```

```

978     /*
979      * Note that one must be extremely careful regarding code
980      * executed while CPUs are paused. Since a CPU may be paused
981      * while a thread scheduling on that CPU is holding an adaptive
982      * lock, code executed with CPUs paused must not acquire adaptive
983      * (or low-level spin) locks. Also, such code must not block,
984      * since the thread that is supposed to initiate the wakeup may
985      * never run.
986
987      * With a few exceptions, the restrictions on code executed with CPUs
988      * paused match those for code executed at high-level interrupt
989      * context.
990    */
991    void
992    pause_cpus(cpu_t *off_cp, void *(*func)(void *));
993 {
994     processorid_t    cpu_id;
995     int             i;
996     struct _cpu_pause_info *cpi = &cpu_pause_info;
997
998     ASSERT(MUTEX_HELD(&cpu_lock));
999     ASSERT(cpi->cp_paused == NULL);
1000    cpi->cp_count = 0;
1001    cpi->cp_go = 0;
1002    for (i = 0; i < NCPU; i++)
1003        safe_list[i] = PAUSE_IDLE;
1004    kpreempt_disable();
1005
1006    cpi->cp_func = func;
1007 #endif /* ! codereview */
1008
1009    /*
1010      * If running on the cpu that is going offline, get off it.
1011      * This is so that it won't be necessary to rechoose a CPU
1012      * when done.
1013    */
1014    if (CPU == off_cp)
1015        cpu_id = off_cp->cpu_next_part->cpu_id;
1016    else
1017        cpu_id = CPU->cpu_id;
1018    affinity_set(cpu_id);
1019
1020    /*
1021      * Start the pause threads and record how many were started
1022    */
1023    cpi->cp_count = cpu_pause_start(cpu_id);
1024
1025    /*
1026      * Now wait for all CPUs to be running the pause thread.
1027    */
1028    while (cpi->cp_count > 0) {
1029        /*
1030          * Spin reading the count without grabbing the disp
1031          * lock to make sure we don't prevent the pause
1032          * threads from getting the lock.
1033        */
1034        while (sema_held(&cpi->cp_sem))
1035            ;
1036        if (sema_tryr(&cpi->cp_sem))
1037            --cpi->cp_count;
1038    }
1039    cpi->cp_go = 1; /* all have reached cpu_pause */
1040
1041    /*
1042      * Now wait for all CPUs to spl. (Transition from PAUSE_READY
1043      * to PAUSE_WAIT.)
1044    */

```

```

1043     /*
1044      for (i = 0; i < NCPUs; i++) {
1045          while (safe_list[i] != PAUSE_WAIT)
1046              ;
1047      }
1048      cpi->cp_spl = splhigh();           /* block dispatcher on this CPU */
1049      cpi->cp_paused = curthread;
1050  }

1052  /*
1053   * Check whether the current thread has CPUs paused
1054  */
1055 int
1056 cpus_paused(void)
1057 {
1058     if (cpu_pause_info.cp_paused != NULL) {
1059         ASSERT(cpu_pause_info.cp_paused == curthread);
1060         return (1);
1061     }
1062     return (0);
1063 }

1065 static cpu_t *
1066 cpu_get_all(processorid_t cpun)
1067 {
1068     ASSERT(MUTEX_HELD(&cpu_lock));

1070     if (cpun >= NCPUs || cpun < 0 || !CPU_IN_SET(cpu_available, cpun))
1071         return (NULL);
1072     return (cpu[cpun]);
1073 }

1075 /*
1076  * Check whether cpun is a valid processor id and whether it should be
1077  * visible from the current zone. If it is, return a pointer to the
1078  * associated CPU structure.
1079  */
1080 cpu_t *
1081 cpu_get(processorid_t cpun)
1082 {
1083     cpu_t *c;

1085     ASSERT(MUTEX_HELD(&cpu_lock));
1086     c = cpu_get_all(cpun);
1087     if (c != NULL && !INGLOBALZONE(curproc) && pool_pset_enabled() &&
1088         zone_pset_get(curproc->p_zone) != cpupart_query_cpu(c))
1089         return (NULL);
1090     return (c);
1091 }

1093 /*
1094  * The following functions should be used to check CPU states in the kernel.
1095  * They should be invoked with cpu_lock held. Kernel subsystems interested
1096  * in CPU states should *not* use cpu_get_state() and various P_ONLINE/etc
1097  * states. Those are for user-land (and system call) use only.
1098  */

1100 /*
1101  * Determine whether the CPU is online and handling interrupts.
1102  */
1103 int
1104 cpu_is_online(cpu_t *cpu)
1105 {
1106     ASSERT(MUTEX_HELD(&cpu_lock));
1107     return (cpu_flagged_online(cpu->cpu_flags));
1108 }

```

```

1110 /*
1111  * Determine whether the CPU is offline (this includes spare and faulted).
1112 */
1113 int
1114 cpu_is_offline(cpu_t *cpu)
1115 {
1116     ASSERT(MUTEX_HELD(&cpu_lock));
1117     return (cpu_flagged_offline(cpu->cpu_flags));
1118 }

1120 /*
1121  * Determine whether the CPU is powered off.
1122 */
1123 int
1124 cpu_is_poweredoff(cpu_t *cpu)
1125 {
1126     ASSERT(MUTEX_HELD(&cpu_lock));
1127     return (cpu_flagged_poweredoff(cpu->cpu_flags));
1128 }

1130 /*
1131  * Determine whether the CPU is handling interrupts.
1132 */
1133 int
1134 cpu_is_nointr(cpu_t *cpu)
1135 {
1136     ASSERT(MUTEX_HELD(&cpu_lock));
1137     return (cpu_flagged_nointr(cpu->cpu_flags));
1138 }

1140 /*
1141  * Determine whether the CPU is active (scheduling threads).
1142 */
1143 int
1144 cpu_is_active(cpu_t *cpu)
1145 {
1146     ASSERT(MUTEX_HELD(&cpu_lock));
1147     return (cpu_flagged_active(cpu->cpu_flags));
1148 }

1150 /*
1151  * Same as above, but these require cpu_flags instead of cpu_t pointers.
1152 */
1153 int
1154 cpu_flagged_online(cpu_flag_t cpu_flags)
1155 {
1156     return (cpu_flagged_active(cpu_flags) &&
1157             (cpu_flags & CPU_ENABLE));
1158 }

1160 int
1161 cpu_flagged_offline(cpu_flag_t cpu_flags)
1162 {
1163     return (((cpu_flags & CPU_POWEROFF) == 0) &&
1164             ((cpu_flags & (CPU_READY | CPU_OFFLINE)) != CPU_READY));
1165 }

1167 int
1168 cpu_flagged_poweredoff(cpu_flag_t cpu_flags)
1169 {
1170     return ((cpu_flags & CPU_POWEROFF) == CPU_POWEROFF);
1171 }

1173 int
1174 cpu_flagged_nointr(cpu_flag_t cpu_flags)

```

```

1175 {
1176     return (cpu_flagged_active(cpu_flags) &&
1177            (cpu_flags & CPU_ENABLE) == 0);
1178 }

1180 int
1181 cpu_flagged_active(cpu_flag_t cpu_flags)
1182 {
1183     return (((cpu_flags & (CPU_POWEROFF | CPU_FAULTED | CPU_SPARE)) == 0) &&
1184            ((cpu_flags & (CPU_READY | CPU_OFFLINE)) == CPU_READY));
1185 }

1187 /*
1188  * Bring the indicated CPU online.
1189  */
1190 int
1191 cpu_online(cpu_t *cp)
1192 {
1193     int error = 0;

1195 /*
1196  * Handle on-line request.
1197  * This code must put the new CPU on the active list before
1198  * starting it because it will not be paused, and will start
1199  * using the active list immediately. The real start occurs
1200  * when the CPU QUIESCED flag is turned off.
1201  */
1203     ASSERT(MUTEX_HELD(&cpu_lock));

1205 /*
1206  * Put all the cpus into a known safe place.
1207  * No mutexes can be entered while CPUs are paused.
1208  */
1209     error = mp_cpu_start(cp); /* arch-dep hook */
1210     if (error == 0) {
1211         pg_cputpart_in(cp, cp->cpu_part);
1212         pause_cpus(NULL, NULL);
1213         pause_cpus(NULL);
1214         cpu_add_active_internal(cp);
1215         if (cp->cpu_flags & CPU_FAULTED) {
1216             cp->cpu_flags &= ~CPU_FAULTED;
1217             mp_cpu_faulted_exit(cp);
1218         }
1219         cp->cpu_flags &= ~(CPU QUIESCED | CPU_OFFLINE | CPU_FROZEN |
1220                           CPU_SPARE);
1221         CPU_NEW_GENERATION(cp);
1222         start_cpus();
1223         cpu_stats_kstat_create(cp);
1224         cpu_create_intrstat(cp);
1225         lgrp_kstat_create(cp);
1226         cpu_state_change_notify(cp->cpu_id, CPU_ON);
1227         cpu_intr_enable(cp); /* arch-dep hook */
1228         cpu_state_change_notify(cp->cpu_id, CPU_INTR_ON);
1229         cpu_set_state(cp);
1230         cyclic_online(cp);
1231         /*
1232          * This has to be called only after cyclic_online(). This
1233          * function uses cyclics.
1234         */
1235         callout_cpu_online(cp);
1236         poke_cpu(cp->cpu_id);
1237     }

1238     return (error);
1239 }

```

```

1241 /*
1242  * Take the indicated CPU offline.
1243  */
1244 int
1245 cpu_offline(cpu_t *cp, int flags)
1246 {
1247     cpupart_t *pp;
1248     int error = 0;
1249     cpu_t *ncp;
1250     int intr_enable;
1251     int cyclic_off = 0;
1252     int callout_off = 0;
1253     int loop_count;
1254     int no_quiesce = 0;
1255     int (*bound_func)(struct cpu *, int);
1256     kthread_t *t;
1257     lpl_t *cpu_lpl;
1258     proc_t *p;
1259     int lgrp_diff_lpl;
1260     boolean_t unbind_all_threads = (flags & CPU_FORCED) != 0;

1262     ASSERT(MUTEX_HELD(&cpu_lock));

1264 /*
1265  * If we're going from faulted or spare to offline, just
1266  * clear these flags and update CPU state.
1267  */
1268     if (cp->cpu_flags & (CPU_FAULTED | CPU_SPARE)) {
1269         if (cp->cpu_flags & CPU_FAULTED) {
1270             cp->cpu_flags &= ~CPU_FAULTED;
1271             mp_cpu_faulted_exit(cp);
1272         }
1273         cp->cpu_flags &= ~CPU_SPARE;
1274         cpu_set_state(cp);
1275     }
1276     return (0);

1278 /*
1279  * Handle off-line request.
1280  */
1281     pp = cp->cpu_part;
1282     /*
1283      * Don't offline last online CPU in partition
1284      */
1285     if (ncpus_online <= 1 || pp->cp_ncpus <= 1 || cpu_intr_count(cp) < 2)
1286         return (EBUSY);
1287     /*
1288      * Unbind all soft-bound threads bound to our CPU and hard bound threads
1289      * if we were asked to.
1290      */
1291     error = cpu_unbind(cp->cpu_id, unbind_all_threads);
1292     if (error != 0)
1293         return (error);
1294     /*
1295      * We shouldn't be bound to this CPU ourselves.
1296      */
1297     if (curthread->t_bound_cpu == cp)
1298         return (EBUSY);

1300 /*
1301  * Tell interested parties that this CPU is going offline.
1302  */
1303     CPU_NEW_GENERATION(cp);
1304     cpu_state_change_notify(cp->cpu_id, CPU_OFF);

```

```

1306     /*
1307      * Tell the PG subsystem that the CPU is leaving the partition
1308      */
1309      pg_cpupart_out(cp, pp);

1311     /*
1312      * Take the CPU out of interrupt participation so we won't find
1313      * bound kernel threads. If the architecture cannot completely
1314      * shut off interrupts on the CPU, don't quiesce it, but don't
1315      * run anything but interrupt thread... this is indicated by
1316      * the CPU_OFLINE flag being on but the CPU QUIESCE flag being
1317      * off.
1318      */
1319      intr_enable = cp->cpu_flags & CPU_ENABLE;
1320      if (intr_enable)
1321          no_quiesce = cpu_intr_disable(cp);

1323     /*
1324      * Record that we are aiming to offline this cpu. This acts as
1325      * a barrier to further weak binding requests in thread_nomigrate
1326      * and also causes cpu_choose, disp_lowpri_cpu and setfrontdq to
1327      * lean away from this cpu. Further strong bindings are already
1328      * avoided since we hold cpu_lock. Since threads that are set
1329      * runnable around now and others coming off the target cpu are
1330      * directed away from the target, existing strong and weak bindings
1331      * (especially the latter) to the target cpu stand maximum chance of
1332      * being able to unbind during the short delay loop below (if other
1333      * unbound threads compete they may not see cpu in time to unbind
1334      * even if they would do so immediately.
1335      */
1336      cpu_inmotion = cp;
1337      membar_enter();

1339     /*
1340      * Check for kernel threads (strong or weak) bound to that CPU.
1341      * Strongly bound threads may not unbind, and we'll have to return
1342      * EBUSY. Weakly bound threads should always disappear - we've
1343      * stopped more weak binding with cpu_inmotion and existing
1344      * bindings will drain imminently (they may not block). Nonetheless
1345      * we will wait for a fixed period for all bound threads to disappear.
1346      * Inactive interrupt threads are OK (they'll be in TS_FREE
1347      * state). If test finds some bound threads, wait a few ticks
1348      * to give short-lived threads (such as interrupts) chance to
1349      * complete. Note that if no_quiesce is set, i.e. this cpu
1350      * is required to service interrupts, then we take the route
1351      * that permits interrupt threads to be active (or bypassed).
1352      */
1353      bound_func = no_quiesce ? disp_bound_threads : disp_bound_anythreads;

1355 again: for (loop_count = 0; (*bound_func)(cp, 0); loop_count++) {
1356     if (loop_count >= 5) {
1357         error = EBUSY; /* some threads still bound */
1358         break;
1359     }

1361     /*
1362      * If some threads were assigned, give them
1363      * a chance to complete or move.
1364      *
1365      * This assumes that the clock_thread is not bound
1366      * to any CPU, because the clock_thread is needed to
1367      * do the delay(hz/100).
1368      *
1369      * Note: we still hold the cpu_lock while waiting for
1370      * the next clock tick. This is OK since it isn't
1371      * needed for anything else except processor_bind(2),

```

```

1372             * and system initialization. If we drop the lock,
1373             * we would risk another p_online disabling the last
1374             * processor.
1375             */
1376             delay(hz/100);
1377         }

1379         if (error == 0 && callout_off == 0) {
1380             callout_cpu_offline(cp);
1381             callout_off = 1;
1382         }

1384         if (error == 0 && cyclic_off == 0) {
1385             if (!cyclic_offline(cp)) {
1386                 /*
1387                  * We must have bound cyclics...
1388                  */
1389                 error = EBUSY;
1390                 goto out;
1391             }
1392             cyclic_off = 1;
1393         }

1395         /*
1396          * Call mp_cpu_stop() to perform any special operations
1397          * needed for this machine architecture to offline a CPU.
1398          */
1399         if (error == 0)
1400             error = mp_cpu_stop(cp); /* arch-dep hook */

1402         /*
1403          * If that all worked, take the CPU offline and decrement
1404          * ncpus_online.
1405          */
1406         if (error == 0) {
1407             /*
1408              * Put all the cpus into a known safe place.
1409              * No mutexes can be entered while CPUs are paused.
1410              */
1411             pause_cpus(cp, NULL);
1412             pause_cpus(cp);
1413             /*
1414              * Repeat the operation, if necessary, to make sure that
1415              * all outstanding low-level interrupts run to completion
1416              * before we set the CPU QUIESCED flag. It's also possible
1417              * that a thread has weak bound to the cpu despite our raising
1418              * cpu_inmotion above since it may have loaded that
1419              * value before the barrier became visible (this would have
1420              * to be the thread that was on the target cpu at the time
1421              * we raised the barrier).
1422              */
1423             if ((!no_quiesce && cp->cpu_intr_actv != 0) ||
1424                 (*bound_func)(cp, 1)) {
1425                 start_cpus();
1426                 (void) mp_cpu_start(cp);
1427                 goto again;
1428             }
1429             ncp = cp->cpu_next_part;
1430             cpu_lpl = cp->cpu_lpl;
1431             ASSERT(cpu_lpl != NULL);

1432             /*
1433              * Remove the CPU from the list of active CPUs.
1434              */
1435             cpu_remove_active(cp);

```

```

1437      /*
1438       * Walk the active process list and look for threads
1439       * whose home lgroup needs to be updated, or
1440       * the last CPU they run on is the one being offline now.
1441       */
1442
1443     ASSERT(curthread->t_cpu != cp);
1444     for (p = pactive; p != NULL; p = p->p_next) {
1445
1446         t = p->p_tlist;
1447
1448         if (t == NULL)
1449             continue;
1450
1451         lgrp_diff_lpl = 0;
1452
1453         do {
1454             ASSERT(t->t_lpl != NULL);
1455
1456             /*
1457              * Taking last CPU in lpl offline
1458              * Rehome thread if it is in this lpl
1459              * Otherwise, update the count of how many
1460              * threads are in this CPU's lgroup but have
1461              * a different lpl.
1462             */
1463
1464             if (cpu_lpl->lpl_ncpu == 0) {
1465                 if (t->t_lpl == cpu_lpl)
1466                     lgrp_move_thread(t,
1467                                     lgrp_choose(t,
1468                                     t->t_cputpart), 0);
1469                 else if (t->t_lpl->lpl_lgrpid ==
1470                          cpu_lpl->lpl_lgrpid)
1471                     lgrp_diff_lpl++;
1472             }
1473             ASSERT(t->t_lpl->lpl_ncpu > 0);
1474
1475             /*
1476              * Update CPU last ran on if it was this CPU
1477              */
1478             if (t->t_cpu == cp && t->t_bound_cpu != cp)
1479                 t->t_cpu = disp_lowpri_cpu(ncp,
1480                                             t->t_lpl, t->t_pri, NULL);
1481             ASSERT(t->t_cpu != cp || t->t_bound_cpu == cp ||
1482                   t->t_weakbound_cpu == cp);
1483
1484             t = t->t_forw;
1485         } while (t != p->p_tlist);
1486
1487         /*
1488          * Didn't find any threads in the same lgroup as this
1489          * CPU with a different lpl, so remove the lgroup from
1490          * the process lgroup bitmask.
1491         */
1492
1493         if (lgrp_diff_lpl == 0)
1494             klgpset_del(p->p_lgrpset, cpu_lpl->lpl_lgrpid);
1495     }
1496
1497     /*
1498      * Walk thread list looking for threads that need to be
1499      * rehomed, since there are some threads that are not in
1500      * their process's p_tlist.
1501     */
1502
1503     t = curthread;

```

```

1503     do {
1504         ASSERT(t != NULL && t->t_lpl != NULL);
1505
1506         /*
1507          * Rehome threads with same lpl as this CPU when this
1508          * is the last CPU in the lpl.
1509         */
1510
1511         if ((cpu_lpl->lpl_ncpu == 0) && (t->t_lpl == cpu_lpl))
1512             lgrp_move_thread(t,
1513                             lgrp_choose(t, t->t_cputpart), 1);
1514
1515         ASSERT(t->t_lpl->lpl_ncpu > 0);
1516
1517         /*
1518          * Update CPU last ran on if it was this CPU
1519         */
1520
1521         if (t->t_cpu == cp && t->t_bound_cpu != cp) {
1522             t->t_cpu = disp_lowpri_cpu(ncp,
1523                                         t->t_lpl, t->t_pri, NULL);
1524         }
1525         ASSERT(t->t_cpu != cp || t->t_bound_cpu == cp ||
1526               t->t_weakbound_cpu == cp);
1527         t = t->t_next;
1528
1529         } while (t != curthread);
1530         ASSERT((cp->cpu_flags & (CPU_FAULTED | CPU_SPARE)) == 0);
1531         cp->cpu_flags |= CPU_OFFLINE;
1532         disp_cpu_inactive(cp);
1533         if (!no_quiesce)
1534             cp->cpu_flags |= CPU QUIESCED;
1535         ncpus_online--;
1536         cpu_set_state(cp);
1537         cpu_inmotion = NULL;
1538         start_cpus();
1539         cpu_stats_kstat_destroy(cp);
1540         cpu_delete_intrstat(cp);
1541         lgrp_kstat_destroy(cp);
1542     }
1543
1544     out:
1545     cpu_inmotion = NULL;
1546
1547     /*
1548      * If we failed, re-enable interrupts.
1549      * Do this even if cpu_intr_disable returned an error, because
1550      * it may have partially disabled interrupts.
1551     */
1552     if (error && intr_enable)
1553         cpu_intr_enable(cp);
1554
1555     /*
1556      * If we failed, but managed to offline the cyclic subsystem on this
1557      * CPU, bring it back online.
1558     */
1559     if (error && cyclic_off)
1560         cyclic_online(cp);
1561
1562     /*
1563      * If we failed, but managed to offline callouts on this CPU,
1564      * bring it back online.
1565     */
1566     if (error && callout_off)
1567         callout_cpu_online(cp);

```

```

1569      /*
1570      * If we failed, tell the PG subsystem that the CPU is back
1571      */
1572      pg_cpupart_in(cp, pp);
1573
1574      /*
1575      * If we failed, we need to notify everyone that this CPU is back on.
1576      */
1577      if (error != 0) {
1578          CPU_NEW_GENERATION(cp);
1579          cpu_state_change_notify(cp->cpu_id, CPU_ON);
1580          cpu_state_change_notify(cp->cpu_id, CPU_INTR_ON);
1581      }
1582
1583      return (error);
1584 }


---

unchanged_portion_omitted
1731 /*
1732 * Insert a CPU into the list of available CPUs.
1733 */
1734 void
1735 cpu_add_unit(cpu_t *cp)
1736 {
1737     int seqid;
1738
1739     ASSERT(MUTEX_HELD(&cpu_lock));
1740     ASSERT(cpu_list != NULL); /* list started in cpu_list_init */
1741
1742     lgrp_config(LGRP_CONFIG_CPU_ADD, (uintptr_t)cp, 0);
1743
1744     /*
1745     * Note: most users of the cpu_list will grab the
1746     * cpu_lock to insure that it isn't modified. However,
1747     * certain users can't or won't do that. To allow this
1748     * we pause the other cpus. Users who walk the list
1749     * without cpu_lock, must disable kernel preemption
1750     * to insure that the list isn't modified underneath
1751     * them. Also, any cached pointers to cpu structures
1752     * must be revalidated by checking to see if the
1753     * cpu_next pointer points to itself. This check must
1754     * be done with the cpu_lock held or kernel preemption
1755     * disabled. This check relies upon the fact that
1756     * old cpu structures are not free'd or cleared after
1757     * then are removed from the cpu_list.
1758
1759     * Note that the clock code walks the cpu list dereferencing
1760     * the cpu_part pointer, so we need to initialize it before
1761     * adding the cpu to the list.
1762
1763     cp->cpu_part = &cp_default;
1764     (void) pause_cpus(NULL, NULL);
1765     (void) pause_cpus(NULL);
1766     cp->cpu_next = cpu_list;
1767     cp->cpu_prev = cpu_list->cpu_prev;
1768     cpu_list->cpu_prev->cpu_next = cp;
1769     cpu_list->cpu_prev = cp;
1770     start_cpus();
1771
1772     for (seqid = 0; CPU_IN_SET(cpu_seqid_inuse, seqid); seqid++)
1773         continue;
1774     CPUSet_ADD(cpu_seqid_inuse, seqid);
1775     cp->cpu_seqid = seqid;
1776
1777     if (seqid > max_cpu_seqid_ever)
1778         max_cpu_seqid_ever = seqid;

```

```

1779     ASSERT(ncpus < max_ncpus);
1780     ncpus++;
1781     cp->cpu_cache_offset = KMEM_CPU_CACHE_OFFSET(cp->cpu_seqid);
1782     cpu[cp->cpu_id] = cp;
1783     CPUSet_ADD(cpu_available, cp->cpu_id);
1784     cpu_seq[cp->cpu_seqid] = cp;
1785
1786     /*
1787     * allocate a pause thread for this CPU.
1788     */
1789     cpu_pause_alloc(cp);
1790
1791     /*
1792     * So that new CPUs won't have NULL prev_onln and next_onln pointers,
1793     * link them into a list of just that CPU.
1794     * This is so that disp_lowpri_cpu will work for thread_create in
1795     * pause_cpus() when called from the startup thread in a new CPU.
1796     */
1797     cp->cpu_next_onln = cp;
1798     cp->cpu_prev_onln = cp;
1799     cpu_info_kstat_create(cp);
1800     cp->cpu_next_part = cp;
1801     cp->cpu_prev_part = cp;
1802
1803     init_cpu_mstate(cp, CMS_SYSTEM);
1804
1805     pool_pset_mod = gethrtime();
1806 }
1807
1808 /*
1809 * Do the opposite of cpu_add_unit().
1810 */
1811 void
1812 cpu_del_unit(int cpuid)
1813 {
1814     struct cpu     *cp, *cpnext;
1815
1816     ASSERT(MUTEX_HELD(&cpu_lock));
1817     cp = cpu[cpuid];
1818     ASSERT(cp != NULL);
1819
1820     ASSERT(cp->cpu_next_onln == cp);
1821     ASSERT(cp->cpu_prev_onln == cp);
1822     ASSERT(cp->cpu_next_part == cp);
1823     ASSERT(cp->cpu_prev_part == cp);
1824
1825     /*
1826     * Tear down the CPU's physical ID cache, and update any
1827     * processor groups
1828     */
1829     pg_cpu_fini(cp, NULL);
1830     pg_hw_physid_destroy(cp);
1831
1832     /*
1833     * Destroy kstat stuff.
1834     */
1835     cpu_info_kstat_destroy(cp);
1836     term_cpu_mstate(cp);
1837
1838     /*
1839     * Free up pause thread.
1840     */
1841     cpu_pause_free(cp);
1842     CPUSet_DEL(cpu_available, cp->cpu_id);
1843     cpu[cp->cpu_id] = NULL;
1844     cpu_seq[cp->cpu_seqid] = NULL;

```

```

1845  /*
1846   * The clock thread and mutex_vector_enter cannot hold the
1847   * cpu_lock while traversing the cpu list, therefore we pause
1848   * all other threads by pausing the other cpus. These, and any
1849   * other routines holding cpu pointers while possibly sleeping
1850   * must be sure to call kpreempt_disable before processing the
1851   * list and be sure to check that the cpu has not been deleted
1852   * after any sleeps (check cp->cpu_next != NULL). We guarantee
1853   * to keep the deleted cpu structure around.
1854   *
1855   * Note that this MUST be done AFTER cpu_available
1856   * has been updated so that we don't waste time
1857   * trying to pause the cpu we're trying to delete.
1858   */
1859   (void) pause_cpus(NULL, NULL);
1860   (void) pause_cpus(NULL);

1861   cpnext = cp->cpu_next;
1862   cp->cpu_prev->cpu_next = cp->cpu_next;
1863   cp->cpu_next->cpu_prev = cp->cpu_prev;
1864   if (cp == cpu_list)
1865     cpu_list = cpnext;

1866   /*
1867    * Signals that the cpu has been deleted (see above).
1868    */
1869   cp->cpu_next = NULL;
1870   cp->cpu_prev = NULL;

1873   start_cpus();

1875   CPUSET_DEL(cpu_seqid_inuse, cp->cpu_seqid);
1876   ncpus--;
1877   lgrp_config(LGRP_CONFIG_CPU_DEL, (uintptr_t)cp, 0);

1879   pool_pset_mod = gethrtime();
1880 }

_____unchanged_portion_omitted_____

1922 /*
1923  * Add a CPU to the list of active CPUs.
1924  * This is called from machine-dependent layers when a new CPU is started.
1925  */
1926 void
1927 cpu_add_active(cpu_t *cp)
1928 {
1929   pg_cpupart_in(cp, cp->cpu_part);

1931   pause_cpus(NULL, NULL);
1932   pause_cpus(NULL);
1933   cpu_add_active_internal(cp);
1934   start_cpus();

1935   cpu_stats_kstat_create(cp);
1936   cpu_create_intrstat(cp);
1937   lgrp_kstat_create(cp);
1938   cpu_state_change_notify(cp->cpu_id, CPU_INIT);
1939 }

_____unchanged_portion_omitted_____

```

```
*****
30588 Thu May 1 19:50:18 2014
new/usr/src/uts/common/os/cpu_event.c
XXXX pass in cpu_pause_func via pause_cpus
*****
_____unchanged_portion_omitted_____
368 static void
369 cpu_idle_insert_callback(cpu_idle_cb_impl_t *cip)
370 {
371     int unlock = 0, unpause = 0;
372     int i, cnt_new = 0, cnt_old = 0;
373     char *buf_new = NULL, *buf_old = NULL;
374
375     ASSERT(MUTEX_HELD(&cpu_idle_cb_lock));
376
377     /*
378      * Expand array if it's full.
379      * Memory must be allocated out of pause/start_cpus() scope because
380      * kmem_zalloc() can't be called with KM_SLEEP flag within that scope.
381      */
382     if (cpu_idle_cb_curr == cpu_idle_cb_max) {
383         cnt_new = cpu_idle_cb_max + CPU_IDLE_ARRAY_CAPACITY_INC;
384         buf_new = (char *)kmem_zalloc(cnt_new *
385                                         sizeof (cpu_idle_cb_item_t), KM_SLEEP);
386     }
387
388     /* Try to acquire cpu_lock if not held yet. */
389     if (!MUTEX_HELD(&cpu_lock)) {
390         mutex_enter(&cpu_lock);
391         unlock = 1;
392     }
393
394     /*
395      * Pause all other CPUs (and let them run pause thread).
396      * It's guaranteed that no other threads will access cpu_idle_cb_array
397      * after pause_cpus().
398      */
399     if (!cpus_paused()) {
400         pause_cpus(NULL, NULL);
401         pause_cpus(NULL);
402         unpause = 1;
403     }
404
405     /* Copy content to new buffer if needed. */
406     if (buf_new != NULL) {
407         buf_old = (char *)cpu_idle_cb_array;
408         cnt_old = cpu_idle_cb_max;
409         if (buf_old != NULL) {
410             ASSERT(cnt_old != 0);
411             bcopy(cpu_idle_cb_array, buf_new,
412                   sizeof (cpu_idle_cb_item_t) * cnt_old);
413         }
414         cpu_idle_cb_array = (cpu_idle_cb_item_t *)buf_new;
415         cpu_idle_cb_max = cnt_new;
416     }
417
418     /* Insert into array according to priority. */
419     ASSERT(cpu_idle_cb_curr < cpu_idle_cb_max);
420     for (i = cpu_idle_cb_curr; i > 0; i--) {
421         if (cpu_idle_cb_array[i - 1].impl->priority >= cip->priority) {
422             break;
423         }
424         cpu_idle_cb_array[i] = cpu_idle_cb_array[i - 1];
425     }
426     cpu_idle_cb_array[i].arg = cip->argument;
427     cpu_idle_cb_array[i].enter = cip->callback->idle_enter;
428 }
```

```
426     cpu_idle_cb_array[i].exit = cip->callback->idle_exit;
427     cpu_idle_cb_array[i].impl = cip;
428     cpu_idle_cb_curr++;
429
430     /* Resume other CPUs from paused state if needed. */
431     if (unpause) {
432         start_cpus();
433     }
434     if (unlock) {
435         mutex_exit(&cpu_lock);
436     }
437
438     /* Free old resource if needed. */
439     if (buf_old != NULL) {
440         ASSERT(cnt_old != 0);
441         kmem_free(buf_old, cnt_old * sizeof (cpu_idle_cb_item_t));
442     }
443 }
444
445 static void
446 cpu_idle_remove_callback(cpu_idle_cb_impl_t *cip)
447 {
448     int i, found = 0;
449     int unlock = 0, unpause = 0;
450     cpu_idle_cb_state_t *sp;
451
452     ASSERT(MUTEX_HELD(&cpu_idle_cb_lock));
453
454     /* Try to acquire cpu_lock if not held yet. */
455     if (!MUTEX_HELD(&cpu_lock)) {
456         mutex_enter(&cpu_lock);
457         unlock = 1;
458     }
459
460     /*
461      * Pause all other CPUs.
462      * It's guaranteed that no other threads will access cpu_idle_cb_array
463      * after pause_cpus().
464      */
465     if (!cpus_paused()) {
466         pause_cpus(NULL, NULL);
467         pause_cpus(NULL);
468         unpause = 1;
469     }
470
471     /* Remove cip from array. */
472     for (i = 0; i < cpu_idle_cb_curr; i++) {
473         if (found == 0) {
474             if (cpu_idle_cb_array[i].impl == cip) {
475                 found = 1;
476             }
477         } else {
478             cpu_idle_cb_array[i - 1] = cpu_idle_cb_array[i];
479         }
480     }
481     ASSERT(found != 0);
482     cpu_idle_cb_curr--;
483
484     /*
485      * Reset property ready flag for all CPUs if no registered callback
486      * left because cpu_idle_enter/exit will stop updating property if
487      * there's no callback registered.
488      */
489     if (cpu_idle_cb_curr == 0) {
490         for (sp = cpu_idle_cb_state, i = 0; i < max_ncpus; i++, sp++) {
491             sp->v.ready = B_FALSE;
492         }
493     }
494 }
```

```
491     }
493     /* Resume other CPUs from paused state if needed. */
494     if (unpause) {
495         start_cpus();
496     }
497     if (unlock) {
498         mutex_exit(&cpu_lock);
499     }
500 }
```

unchanged_portion_omitted

```
*****
21527 Thu May 1 19:50:18 2014
new/usr/src/uts/common/os/cpu_pm.c
XXXX pass in cpu_pause_func via pause_cpus
*****
_____unchanged_portion_omitted_____
172 int
173 cpupm_set_policy(cpupm_policy_t new_policy)
174 {
175     static int      gov_init = 0;
176     int             result = 0;
177
178     mutex_enter(&cpu_lock);
179     if (new_policy == cpupm_policy) {
180         mutex_exit(&cpu_lock);
181         return (result);
182     }
183
184     /*
185      * Pausing CPUs causes a high priority thread to be scheduled
186      * on all other CPUs (besides the current one). This locks out
187      * other CPUs from making CPUPM state transitions.
188      */
189     switch (new_policy) {
190     case CPUPM_POLICY_DISABLED:
191         pause_cpus(NULL, NULL);
192         pause_cpus(NULL);
193         cpupm_policy = CPUPM_POLICY_DISABLED;
194         start_cpus();
195
196         result = cmt_pad_disable(PGHW_POW_ACTIVE);
197
198         /*
199          * Once PAD has been enabled, it should always be possible
200          * to disable it.
201          */
202         ASSERT(result == 0);
203
204         /*
205          * Bring all the active power domains to the maximum
206          * performance state.
207          */
208         cpupm_state_change_global(CPUPM_DTYPE_ACTIVE,
209             CPUPM_STATE_MAX_PERF);
210
211         break;
212     case CPUPM_POLICY_ELASTIC:
213         result = cmt_pad_enable(PGHW_POW_ACTIVE);
214         if (result < 0) {
215             /*
216              * Failed to enable PAD across the active power
217              * domains, which may well be because none were
218              * enumerated.
219              */
220             break;
221         }
222
223         /*
224          * Initialize the governor parameters the first time through.
225          */
226         if (gov_init == 0) {
227             cpupm_governor_initialize();
228             gov_init = 1;
229         }
230
231         pause_cpus(NULL, NULL);
232         pause_cpus(NULL);
233         cpupm_policy = CPUPM_POLICY_ELASTIC;
234         start_cpus();
235
236         break;
237     default:
238         cmn_err(CE_WARN, "Attempt to set unknown CPUPM policy %d\n",
239             new_policy);
240         ASSERT(0);
241         break;
242     }
243     mutex_exit(&cpu_lock);
244
245 }_____unchanged_portion_omitted_____

```

```
231     pause_cpus(NULL, NULL);
232     pause_cpus(NULL);
233     cpupm_policy = CPUPM_POLICY_ELASTIC;
234     start_cpus();
235
236     break;
237     default:
238         cmn_err(CE_WARN, "Attempt to set unknown CPUPM policy %d\n",
239             new_policy);
240         ASSERT(0);
241         break;
242     }
243     mutex_exit(&cpu_lock);
244
245 }_____unchanged_portion_omitted_____

```

```
*****
119440 Thu May 1 19:50:19 2014
new/usr/src/uts/common/os/lgrp.c
XXXX pass in cpu_pause_func via pause_cpus
*****
_____unchanged_portion_omitted_____
1225 /*
1226 * Called to indicate that the lgrp with platform handle "hand" now
1227 * contains the memory identified by "mnode".
1228 *
1229 * LOCKING for this routine is a bit tricky. Usually it is called without
1230 * cpu_lock and it must must grab cpu_lock here to prevent racing with other
1231 * callers. During DR of the board containing the caged memory it may be called
1232 * with cpu_lock already held and CPUs paused.
1233 *
1234 * If the insertion is part of the DR copy-rename and the inserted mnode (and
1235 * only this mnode) is already present in the lgrp_root->lgrp_mnodes set, we are
1236 * dealing with the special case of DR copy-rename described in
1237 * lgrp_mem_rename().
1238 */
1239 void
1240 lgrp_mem_init(int mnode, lgrp_handle_t hand, boolean_t is_copy_rename)
1241 {
1242     klgpset_t    changed;
1243     int          count;
1244     int          i;
1245     lgrp_t       *my_lgrp;
1246     lgrp_id_t    lgrp_id;
1247     mnodeset_t   mnodes_mask = ((mnodeset_t)1 << mnode);
1248     boolean_t    drop_lock = B_FALSE;
1249     boolean_t    need_synch = B_FALSE;
1250
1251     /*
1252     * Grab CPU lock (if we haven't already)
1253     */
1254     if (!MUTEX_HELD(&cpu_lock)) {
1255         mutex_enter(&cpu_lock);
1256         drop_lock = B_TRUE;
1257     }
1258
1259     /*
1260     * This routine may be called from a context where we already
1261     * hold cpu_lock, and have already paused cpus.
1262     */
1263     if (!cpus_paused())
1264         need_synch = B_TRUE;
1265
1266     /*
1267     * Check if this mnode is already configured and return immediately if
1268     * it is.
1269     *
1270     * NOTE: in special case of copy-rename of the only remaining mnode,
1271     * lgrp_mem_fini() refuses to remove the last mnode from the root, so we
1272     * recognize this case and continue as usual, but skip the update to
1273     * the lgrp_mnodes and the lgrp_mmnodes. This restores the inconsistency
1274     * in topology, temporarily introduced by lgrp_mem_fini().
1275     */
1276     if (! (is_copy_rename && (lgrp_root->lgrp_mnodes == mnodes_mask)) &&
1277         lgrp_root->lgrp_mnodes & mnodes_mask) {
1278         if (drop_lock)
1279             mutex_exit(&cpu_lock);
1280         return;
1281     }
1282
1283     /*

```

```
1284     * Update lgroup topology with new memory resources, keeping track of
1285     * which lgroups change
1286     */
1287     count = 0;
1288     klgpset_clear(changed);
1289     my_lgrp = lgrp_hand_to_lgrp(hand);
1290     if (my_lgrp == NULL) {
1291         /* new lgrp */
1292         my_lgrp = lgrp_create();
1293         lgrp_id = my_lgrp->lgrp_id;
1294         my_lgrp->lgrp_plathand = hand;
1295         my_lgrp->lgrp_latency = lgrp_plat_latency(hand, hand);
1296         klgpset_add(my_lgrp->lgrp_leaves, lgrp_id);
1297         klgpset_add(my_lgrp->lgrp_set[LGRP_RSRC_MEM], lgrp_id);
1298
1299         if (need_synch)
1300             pause_cpus(NULL, NULL);
1301             pause_cpus(NULL);
1302         count = lgrp_leaf_add(my_lgrp, lgrp_table, lgrp_alloc_max + 1,
1303                               &changed);
1304         if (need_synch)
1305             start_cpus();
1306     } else if (my_lgrp->lgrp_latency == 0 && lgrp_plat_latency(hand, hand)
1307               > 0) {
1308         /*
1309         * Leaf lgroup was created, but latency wasn't available
1310         * then. So, set latency for it and fill in rest of lgroup
1311         * topology now that we know how far it is from other leaf
1312         * lgroups.
1313         */
1314         klgpset_clear(changed);
1315         lgrp_id = my_lgrp->lgrp_id;
1316         if (!klgpset_ismember(my_lgrp->lgrp_set[LGRP_RSRC_MEM],
1317                               lgrp_id))
1318             klgpset_add(my_lgrp->lgrp_set[LGRP_RSRC_MEM], lgrp_id);
1319         if (need_synch)
1320             pause_cpus(NULL, NULL);
1321             pause_cpus(NULL);
1322         count = lgrp_leaf_add(my_lgrp, lgrp_table, lgrp_alloc_max + 1,
1323                               &changed);
1324         if (need_synch)
1325             start_cpus();
1326     } else if (!klgpset_ismember(my_lgrp->lgrp_set[LGRP_RSRC_MEM],
1327                                 my_lgrp->lgrp_id)) {
1328         /*
1329         * Add new lgroup memory resource to existing lgroup
1330         */
1331         lgrp_id = my_lgrp->lgrp_id;
1332         klgpset_add(my_lgrp->lgrp_set[LGRP_RSRC_MEM], lgrp_id);
1333         klgpset_add(changed, lgrp_id);
1334         count++;
1335         for (i = 0; i <= lgrp_alloc_max; i++) {
1336             lgrp_t      *lgrp;
1337             lgrp = lgrp_table[i];
1338             if (!LGRP_EXISTS(lgrp) ||
1339                 !lgrp_rsets_member(lgrp->lgrp_set, lgrp_id))
1340                 continue;
1341             klgpset_add(lgrp->lgrp_set[LGRP_RSRC_MEM], lgrp_id);
1342             klgpset_add(changed, lgrp->lgrp_id);
1343             count++;
1344         }
1345     }
1346
1347     /*

```

```

1348     * Add memory node to lgroup and remove lgroup from ones that need
1349     * to be updated
1350     */
1351     if (!(my_lgrp->lgrp_mnodes & mnodes_mask)) {
1352         my_lgrp->lgrp_mnodes |= mnodes_mask;
1353         my_lgrp->lgrp_nmnodes++;
1354     }
1355     klgrpset_del(changed, lgrp_id);
1356
1357     /*
1358     * Update memory node information for all lgroups that changed and
1359     * contain new memory node as a resource
1360     */
1361     if (count)
1362         (void) lgrp_mnode_update(changed, NULL);
1363
1364     if (drop_lock)
1365         mutex_exit(&cpu_lock);
1366 }
1367
1368 /**
1369 * Called to indicate that the lgroup associated with the platform
1370 * handle "hand" no longer contains given memory node
1371 */
1372
1373 * LOCKING for this routine is a bit tricky. Usually it is called without
1374 * cpu_lock and it must grab cpu_lock here to prevent racing with other
1375 * callers. During DR of the board containing the caged memory it may be called
1376 * with cpu_lock already held and CPUs paused.
1377
1378 * If the deletion is part of the DR copy-rename and the deleted mnode is the
1379 * only one present in the lgrp_root->lgrp_mnodes, all the topology is updated,
1380 * but lgrp_root->lgrp_mnodes is left intact. Later, lgrp_mem_init() will insert
1381 * the same mnode back into the topology. See lgrp_mem_rename() and
1382 * lgrp_mem_init() for additional details.
1383 */
1384 void lgrp_mem_fini(int mnode, lgrp_handle_t hand, boolean_t is_copy_rename)
1385 {
1386     klgrpset_t     changed;
1387     int            count;
1388     int            i;
1389     lgrp_t          *my_lgrp;
1390     lgrp_id_t       lgrp_id;
1391     mnodeset_t      mnodes_mask;
1392     boolean_t       drop_lock = B_FALSE;
1393     boolean_t       need_synch = B_FALSE;
1394
1395     /*
1396     * Grab CPU lock (if we haven't already)
1397     */
1398     if (!MUTEX_HELD(&cpu_lock)) {
1399         mutex_enter(&cpu_lock);
1400         drop_lock = B_TRUE;
1401     }
1402
1403     /*
1404     * This routine may be called from a context where we already
1405     * hold cpu_lock and have already paused cpus.
1406     */
1407     if (!cpus_paused())
1408         need_synch = B_TRUE;
1409
1410     my_lgrp = lgrp_hand_to_lgrp(hand);
1411
1412     /*
1413     * The lgrp *must* be pre-existing

```

```

1414     */
1415     ASSERT(my_lgrp != NULL);
1416
1417     /*
1418     * Delete memory node from lgroups which contain it
1419     */
1420     mnodes_mask = ((mnodeset_t)1 << mnode);
1421     for (i = 0; i <= lgrp_alloc_max; i++) {
1422         lgrp_t *lgrp = lgrp_table[i];
1423
1424         /*
1425         * Skip any non-existent lgroups and any lgroups that don't
1426         * contain leaf lgroup of memory as a memory resource
1427         */
1428         if (!LGRP_EXISTS(lgrp) ||
1429             !(lgrp->lgrp_mnodes & mnodes_mask))
1430             continue;
1431
1432         /*
1433         * Avoid removing the last mnode from the root in the DR
1434         * copy-rename case. See lgrp_mem_rename() for details.
1435         */
1436         if (is_copy_rename &&
1437             (lgrp == lgrp_root) && (lgrp->lgrp_mnodes == mnodes_mask))
1438             continue;
1439
1440         /*
1441         * Remove memory node from lgroup.
1442         */
1443         lgrp->lgrp_mnodes &= ~mnodes_mask;
1444         lgrp->lgrp_nmnodes--;
1445         ASSERT(lgrp->lgrp_nmnodes >= 0);
1446     }
1447     ASSERT(lgrp_root->lgrp_nmnodes > 0);
1448
1449     /*
1450     * Don't need to update lgroup topology if this lgroup still has memory.
1451     */
1452     /*
1453     * In the special case of DR copy-rename with the only mnode being
1454     * removed, the lgrp_mnodes for the root is always non-zero, but we
1455     * still need to update the lgroup topology.
1456     */
1457     if ((my_lgrp->lgrp_nmnodes > 0) &&
1458         !(is_copy_rename && (my_lgrp == lgrp_root) &&
1459           (my_lgrp->lgrp_mnodes == mnodes_mask))) {
1460         if (drop_lock)
1461             mutex_exit(&cpu_lock);
1462         return;
1463     }
1464
1465     /*
1466     * This lgroup does not contain any memory now
1467     */
1468     klgrpset_clear(my_lgrp->lgrp_set[LGRP_RSRC_MEM]);
1469
1470     /*
1471     * Remove this lgroup from lgroup topology if it does not contain any
1472     * resources now
1473     */
1474     lgrp_id = my_lgrp->lgrp_id;
1475     count = 0;
1476     klgrpset_clear(changed);
1477     if (lgrp_rsets_empty(my_lgrp->lgrp_set)) {
1478
1479         /*
1480         * Delete lgroup when no more resources
1481         */
1482         if (need_synch)

```

```
1480         pause_cpus(NULL, NULL);
1480         pause_cpus(NULL);
1481     count = lgrp_leaf_delete(my_lgrp, lgrp_table,
1482         lgrp_alloc_max + 1, &changed);
1483     ASSERT(count > 0);
1484     if (need_synch)
1485         start_cpus();
1486 } else {
1487     /*
1488      * Remove lgroup from memory resources of any lgroups that
1489      * contain it as such
1490      */
1491     for (i = 0; i <= lgrp_alloc_max; i++) {
1492         lgrp_t             *lgrp;
1493
1494         lgrp = lgrp_table[i];
1495         if (!LGRP_EXISTS(lgrp) ||
1496             !klgrpset_ismember(lgrp->lgrp_set[LGRP_RSRC_MEM],
1497                 lgrp_id))
1498             continue;
1499
1500         klgrpset_del(lgrp->lgrp_set[LGRP_RSRC_MEM], lgrp_id);
1501     }
1502     if (drop_lock)
1503         mutex_exit(&cpu_lock);
1504 }
```

unchanged portion omitted

```
new/usr/src/uts/common/os/lgrp_topo.c
```

```
*****  
36945 Thu May 1 19:50:19 2014  
new/usr/src/uts/common/os/lgrp_topo.c  
XXXX pass in cpu_pause_func via pause_cpus  
*****
```

```
_____unchanged_portion_omitted_____
```

```
1448 /*  
1449 * Update lgroup topology for any leaves that don't have their latency set  
1450 *  
1451 * This may happen on some machines when the lgroup platform support doesn't  
1452 * know the latencies between nodes soon enough to provide it when the  
1453 * resources are being added. If the lgroup platform code needs to probe  
1454 * memory to determine the latencies between nodes, it must wait until the  
1455 * CPUs become active so at least one CPU in each node can probe memory in  
1456 * each node.  
1457 */  
1458 int  
1459 lgrp_topo_update(lgrp_t **lgrps, int lgrp_count, klgpset_t *changed)  
1460 {  
1461     klgpset_t     changes;  
1462     int           count;  
1463     int           i;  
1464     lgrp_t        *lgrp;  
1465  
1466     count = 0;  
1467     if (changed)  
1468         klgpset_clear(*changed);  
1469  
1470     /*  
1471      * For UMA machines, make sure that root lgroup contains all  
1472      * resources. The root lgrp should also name itself as its own leaf  
1473      */  
1474     if (nlgrps == 1) {  
1475         for (i = 0; i < LGRP_RSRC_COUNT; i++)  
1476             klgpset_add(lgrp_root->lgrp_set[i],  
1477                         lgrp_root->lgrp_id);  
1478         klgpset_add(lgrp_root->lgrp_leaves, lgrp_root->lgrp_id);  
1479         return (0);  
1480     }  
1481  
1482     mutex_enter(&cpu_lock);  
1483     pause_cpus(NULL, NULL);  
1484     pause_cpus(NULL);  
1485  
1486     /*  
1487      * Look for any leaf lgroup without its latency set, finish adding it  
1488      * to the lgroup topology assuming that it exists and has the root  
1489      * lgroup as its parent, and update the memory nodes of all lgroups  
1490      * that have it as a memory resource.  
1491      */  
1492     for (i = 0; i < lgrp_count; i++) {  
1493         lgrp = lgrps[i];  
1494  
1495         /*  
1496          * Skip non-existent and non-leaf lgroups and any lgroup  
1497          * with its latency set already  
1498          */  
1499         if (lgrp == NULL || lgrp->lgrp_id == LGRP_NONE ||  
1500             lgrp->lgrp_childcnt != 0 || lgrp->lgrp_latency != 0)  
1501             continue;  
1502 #ifdef DEBUG  
1503     if (lgrp_topo_debug > 1) {  
1504         prom_printf("\nlgrp_topo_update: updating lineage "
```

```
1
```

```
new/usr/src/uts/common/os/lgrp_topo.c
```

```
1505                                         "of lgrp %d at 0x%p\n", lgrp->lgrp_id,  
1506                                         (void *)lgrp);  
1507 }  
1508 #endif /* DEBUG */  
1509  
1510         count += lgrp_leaf_add(lgrp, lgrps, lgrp_count, &changes);  
1511         if (changed)  
1512             klgpset_or(*changed, changes);  
1513  
1514         if (!klgpset_isempty(changes))  
1515             (void) lgrp_mnode_update(changes, NULL);  
1516  
1517 #ifdef DEBUG  
1518     if (lgrp_topo_debug > 1 && changed)  
1519         prom_printf("lgrp_topo_update: changed %d lgrps: "  
1520                     "0x%llx\n",  
1521                     count, (u_longlong_t)*changed);  
1522 #endif /* DEBUG */  
1523 }  
1524  
1525     if (lgrp_topo_levels < LGRP_TOPO_LEVELS && lgrp_topo_levels == 2) {  
1526         count += lgrp_topo_flatten(2, lgrps, lgrp_count, changed);  
1527         (void) lpl_topo_flatten(2);  
1528     }  
1529  
1530     start_cpus();  
1531     mutex_exit(&cpu_lock);  
1532  
1533 }  
1534 }  
_____unchanged_portion_omitted_____
```

```
2
```

```
*****
82758 Thu May 1 19:50:19 2014
new/usr/src/uts/common/os/mem_config.c
XXXX pass in cpu_pause_func via pause_cpus
*****
_____unchanged_portion_omitted_____
3295 /*
3296  * Invalidate memseg pointers in cpu private vm data caches.
3297 */
3298 static void
3299 memseg_cpu_vm_flush()
3300 {
3301     cpu_t *cp;
3302     vm_cpu_data_t *vc;
3303
3304     mutex_enter(&cpu_lock);
3305     pause_cpus(NULL, NULL);
3306     pause_cpus(NULL);
3307
3308     cp = cpu_list;
3309     do {
3310         vc = cp->cpu_vm_data;
3311         vc->vc_pnum_memseg = NULL;
3312         vc->vc_pnext_memseg = NULL;
3313     } while ((cp = cp->cpu_next) != cpu_list);
3314
3315     start_cpus();
3316     mutex_exit(&cpu_lock);
3317 }
_____unchanged_portion_omitted_____
```

```
*****
30244 Thu May 1 19:50:19 2014
new/usr/src/uts/common/sys/cpuvar.h
XXXX pass in cpu_pause_func via pause_cpus
*****
_____unchanged_portion_omitted_____
611 #define CPU_STATS(cp, stat) \
612     ((cp)->cpu_stats.stat) \
613 \
614 /* \
615  * Increment CPU generation value. \
616  * This macro should be called whenever CPU goes on-line or off-line. \
617  * Updates to cpu_generation should be protected by cpu_lock. \
618 */ \
619 #define CPU_NEW_GENERATION(cp) ((cp)->cpu_generation++) \
620 \
621 #endif /* _KERNEL || _KMEMUSER */ \
622 \
623 /* \
624  * CPU support routines. \
625 */ \
626 #if defined(_KERNEL) && defined(__STDC__) /* not for genassym.c */ \
627 \
628 struct zone; \
629 \
630 void cpu_list_init(cpu_t *); \
631 void cpu_add_unit(cpu_t *); \
632 void cpu_del_unit(int cpuid); \
633 void cpu_add_active(cpu_t *); \
634 void cpu_kstat_init(cpu_t *); \
635 void cpu_visibility_add(cpu_t *, struct zone *); \
636 void cpu_visibility_remove(cpu_t *, struct zone *); \
637 void cpu_visibility_configure(cpu_t *, struct zone *); \
638 void cpu_visibility_unconfigure(cpu_t *, struct zone *); \
639 void cpu_visibility_online(cpu_t *, struct zone *); \
640 void cpu_visibility_offline(cpu_t *, struct zone *); \
641 void cpu_create_intrstat(cpu_t *); \
642 void cpu_delete_intrstat(cpu_t *); \
643 int cpu_kstat_intrstat_update(kstat_t *, int); \
644 void cpu_intr_swth_enter(kthread_t *); \
645 void cpu_intr_swth_exit(kthread_t *); \
646 \
647 void mbox_lock_init(void); /* initialize cross-call locks */ \
648 void mbox_init(int cpun); /* initialize cross-calls */ \
649 void poke_cpu(int cpun); /* interrupt another CPU (to preempt) */ \
650 \
651 /* \
652  * values for safe_list. Pause state that CPUs are in. \
653 */ \
654 #define PAUSE_IDLE 0 /* normal state */ \
655 #define PAUSE_READY 1 /* paused thread ready to spl */ \
656 #define PAUSE_WAIT 2 /* paused thread is spl-ed high */ \
657 #define PAUSE_DIE 3 /* tell pause thread to leave */ \
658 #define PAUSE_DEAD 4 /* pause thread has left */ \
659 \
660 void mach_cpu_pause(volatile char *); \
661 \
662 void pause_cpus(cpu_t *off_cp, void *(*func)(void *)); \
663 void pause_cpus(cpu_t *off_cp); \
664 void start_cpus(void); \
665 int cpus_paused(void); \
666 \
667 void cpu_pause_init(void); \
668 cpu_t *cpu_get(processorid_t cpun); /* get the CPU struct associated */
```

```
669 int cpu_online(cpu_t *cp); /* take cpu online */ \
670 int cpu_offline(cpu_t *cp, int flags); /* take cpu offline */ \
671 int cpu_spare(cpu_t *cp, int flags); /* take cpu to spare */ \
672 int cpu_faulted(cpu_t *cp, int flags); /* take cpu to faulted */ \
673 int cpu_poweron(cpu_t *cp); /* take powered-off cpu to offline */ \
674 int cpu_poweroff(cpu_t *cp); /* take offline cpu to powered-off */ \
675 \
676 cpu_t *cpu_intr_next(cpu_t *cp); /* get next online CPU taking intrs */ \
677 int cpu_intr_count(cpu_t *cp); /* count # of CPUs handling intrs */ \
678 int cpu_intr_on(cpu_t *cp); /* CPU taking I/O interrupts? */ \
679 void cpu_intr_enable(cpu_t *cp); /* enable I/O interrupts */ \
680 int cpu_intr_disable(cpu_t *cp); /* disable I/O interrupts */ \
681 void cpu_intr_alloc(cpu_t *cp, int n); /* allocate interrupt threads */ \
682 \
683 /* \
684  * Routines for checking CPU states. \
685 */ \
686 int cpu_is_online(cpu_t *); /* check if CPU is online */ \
687 int cpu_is_nointr(cpu_t *); /* check if CPU can service intrs */ \
688 int cpu_is_active(cpu_t *); /* check if CPU can run threads */ \
689 int cpu_is_offline(cpu_t *); /* check if CPU is offline */ \
690 int cpu_isPoweredoff(cpu_t *); /* check if CPU is powered off */ \
691 \
692 int cpu_flagged_online(cpu_flag_t); /* flags show CPU is online */ \
693 int cpu_flagged_nointr(cpu_flag_t); /* flags show CPU not handling intrs */ \
694 int cpu_flagged_active(cpu_flag_t); /* flags show CPU scheduling threads */ \
695 int cpu_flagged_offline(cpu_flag_t); /* flags show CPU is offline */ \
696 int cpu_flaggedPoweredoff(cpu_flag_t); /* flags show CPU is powered off */ \
697 \
698 /* \
699  * The processor_info(2) state of a CPU is a simplified representation suitable \
700  * for use by an application program. Kernel subsystems should utilize the \
701  * internal per-CPU state as given by the cpu_flags member of the cpu structure, \
702  * as this information may include platform- or architecture-specific state \
703  * critical to a subsystem's disposition of a particular CPU. \
704 */ \
705 void cpu_set_state(cpu_t *); /* record/timestamp current state */ \
706 int cpu_get_state(cpu_t *); /* get current cpu state */ \
707 const char *cpu_get_state_str(cpu_t *); /* get current cpu state as string */ \
708 \
709 void cpu_set_curr_clock(uint64_t); /* indicate the current CPU's freq */ \
710 void cpu_set_supp_freqs(cpu_t *, const char *); /* set the CPU supported */ \
711 /* frequencies */ \
712 \
713 int cpu_configure(int); \
714 int cpu_unconfigure(int); \
715 void cpu_destroy_bound_threads(cpu_t *cp); \
716 \
717 extern int cpu_bind_thread(kthread_t *tp, processorid_t bind, \
718 processorid_t *obind, int *error); \
719 extern int cpu_unbind(processorid_t cpu_id, boolean_t force); \
720 extern void thread_affinity_set(kthread_t *t, int cpu_id); \
721 extern void thread_affinity_clear(kthread_t *t); \
722 extern void affinity_set(int cpu_id); \
723 extern void affinity_clear(void); \
724 extern void init_cpu_mstate(struct cpu *, int); \
725 extern void term_cpu_mstate(struct cpu *); \
726 extern void new_cpu_mstate(int, hrtimer_t); \
727 extern void get_cpu_mstate(struct cpu *, hrtimer_t *); \
728 extern void thread_nomigrate(void); \
729 extern void thread_allownmigrate(void); \
730 extern void weakbinding_stop(void); \
731 extern void weakbinding_start(void); \
732 \
733 /*
```

```

735 * The following routines affect the CPUs participation in interrupt processing,
736 * if that is applicable on the architecture. This only affects interrupts
737 * which aren't directed at the processor (not cross calls).
738 *
739 * cpu_disable_intr returns non-zero if interrupts were previously enabled.
740 */
741 int     cpu_disable_intr(struct cpu *cp); /* stop issuing interrupts to cpu */
742 void    cpu_enable_intr(struct cpu *cp); /* start issuing interrupts to cpu */

744 /*
745 * The mutex cpu_lock protects cpu_flags for all CPUs, as well as the ncpus
746 * and ncpus_online counts.
747 */
748 extern kmutex_t cpu_lock;      /* lock protecting CPU data */

750 /*
751 * CPU state change events
752 *
753 * Various subsystems need to know when CPUs change their state. They get this
754 * information by registering CPU state change callbacks using
755 * register_cpu_setup_func(). Whenever any CPU changes its state, the callback
756 * function is called. The callback function is passed three arguments:
757 *
758 *   Event, described by cpu_setup_t
759 *   CPU ID
760 *   Transparent pointer passed when registering the callback
761 *
762 * The callback function is called with cpu_lock held. The return value from the
763 * callback function is usually ignored, except for CPU_CONFIG and CPU_UNCONFIG
764 * events. For these two events, non-zero return value indicates a failure and
765 * prevents successful completion of the operation.
766 *
767 * New events may be added in the future. Callback functions should ignore any
768 * events that they do not understand.
769 *
770 * The following events provide notification callbacks:
771 *
772 * CPU_INIT    A new CPU is started and added to the list of active CPUs
773 *               This event is only used during boot
774 *
775 * CPU_CONFIG   A newly inserted CPU is prepared for starting running code
776 *               This event is called by DR code
777 *
778 * CPU_UNCONFIG CPU has been powered off and needs cleanup
779 *               This event is called by DR code
780 *
781 * CPU_ON       CPU is enabled but does not run anything yet
782 *
783 * CPU_INTR_ON CPU is enabled and has interrupts enabled
784 *
785 * CPU_OFF      CPU is going offline but can still run threads
786 *
787 * CPU_CPUPART_OUT    CPU is going to move out of its partition
788 *
789 * CPU_CPUPART_IN     CPU is going to move to a new partition
790 *
791 * CPU_SETUP      CPU is set up during boot and can run threads
792 */
793 typedef enum {
794     CPU_INIT,
795     CPU_CONFIG,
796     CPU_UNCONFIG,
797     CPU_ON,
798     CPU_OFF,
799     CPU_CPUPART_IN,
800     CPU_CPUPART_OUT,

```

```

801     CPU_SETUP,
802     CPU_INTR_ON
803 } cpu_setup_t;


---


unchanged portion omitted

```

```
*****
18246 Thu May 1 19:50:20 2014
new/usr/src/uts/i86pc/i86hvm/io/xpv/xpv_support.c
XXXX pass in cpu_pause_func via pause_cpus
*****
_____unchanged_portion_omitted_____
429 /*
430  * Top level routine to direct suspend/resume of a domain.
431  */
432 void
433 xen_suspend_domain(void)
434 {
435     extern void rtcsync(void);
436     extern void ec_resume(void);
437     extern kmutex_t ec_lock;
438     struct xen_add_to_physmap xatp;
439     ulong_t flags;
440     int err;
441
442     cmn_err(CE_NOTE, "Domain suspending for save/migrate");
443     SUSPEND_DEBUG("xen_suspend_domain\n");
444
445     /*
446      * We only want to suspend the PV devices, since the emulated devices
447      * are suspended by saving the emulated device state. The PV devices
448      * are all children of the xpvd nexus device. So we search the
449      * device tree for the xpvd node to use as the root of the tree to
450      * be suspended.
451     */
452     if (xpvd_dip == NULL)
453         ddi_walk_devs(ddi_root_node(), check_xpvd, NULL);
454
455     /*
456      * suspend interrupts and devices
457     */
458     if (xpvd_dip != NULL)
459         (void) xen_suspend_devices(ddi_get_child(xpvd_dip));
460     else
461         cmn_err(CE_WARN, "No PV devices found to suspend");
462     SUSPEND_DEBUG("xenbus_suspend\n");
463     xenbus_suspend();
464
465     mutex_enter(&cpu_lock);
466
467     /*
468      * Suspend on vcpu 0
469     */
470     thread_affinity_set(curthread, 0);
471     kpreempt_disable();
472
473     if (ncpus > 1)
474         pause_cpus(NULL, NULL);
475     pause_cpus(NULL);
476
477     /*
478      * We can grab the ec_lock as it's a spinlock with a high SPL. Hence
479      * any holder would have dropped it to get through pause_cpus().
480     */
481     mutex_enter(&ec_lock);
482
483     /*
484      * From here on in, we can't take locks.
485     */
486
487     flags = intr_clear();
```

```
488     SUSPEND_DEBUG("HYPERVISOR_suspend\n");
489     /*
490      * At this point we suspend and sometime later resume.
491      * Note that this call may return with an indication of a cancelled
492      * for now no matter what the return we do a full resume of all
493      * suspended drivers, etc.
494     */
495     (void) HYPERVISOR_shutdown(SHUTDOWN_suspend);
496
497     /*
498      * Point HYPERVISOR_shared_info to the proper place.
499     */
500     xatp.domid = DOMID_SELF;
501     xatp.idx = 0;
502     xatp.space = XENMAPSPACE_shared_info;
503     xatp.gpfn = xen_shared_info_frame;
504     if ((err = HYPERVISOR_memory_op(XENMEM_add_to_physmap, &xatp)) != 0)
505         panic("Could not set shared_info page. error: %d", err);
506
507     SUSPEND_DEBUG("gnttab_resume\n");
508     gnttab_resume();
509
510     SUSPEND_DEBUG("ec_resume\n");
511     ec_resume();
512
513     intr_restore(flags);
514
515     if (ncpus > 1)
516         start_cpus();
517
518     mutex_exit(&ec_lock);
519     mutex_exit(&cpu_lock);
520
521     /*
522      * Now we can take locks again.
523     */
524
525     rtcsync();
526
527     SUSPEND_DEBUG("xenbus_resume\n");
528     xenbus_resume();
529     SUSPEND_DEBUG("xen_resume_devices\n");
530     if (xpvd_dip != NULL)
531         (void) xen_resume_devices(ddi_get_child(xpvd_dip), 0);
532
533     thread_affinity_clear(curthread);
534     kpreempt_enable();
535
536     SUSPEND_DEBUG("finished xen_suspend_domain\n");
537
538 }
539
540 _____unchanged_portion_omitted_____
```

new/usr/src/uts/i86pc/io/dr/dr_quiesce.c

23453 Thu May 1 19:50:20 2014
new/usr/src/uts/i86pc/io/dr/dr_quiesce.c
XXXX pass in cpu_pause_func via pause_cpus

_____unchanged_portion_omitted_____

```
785 int
786 dr_suspend(dr_sr_handle_t *srh)
787 {
788     dr_handle_t      *handle;
789     int               force;
790     int               dev_errs_idx;
791     uint64_t          dev_errs[DR_MAX_ERR_INT];
792     int               rc = DDI_SUCCESS;
793
794     handle = srh->sr_dr_handlep;
795
796     force = dr_cmd_flags(handle) & SBD_FLAG_FORCE;
797
798     prom_printf("\nDR: suspending user threads...\n");
799     srh->sr_suspend_state = DR_SRSTATE_USER;
800     if (((rc = dr_stop_user_threads(srh)) != DDI_SUCCESS) &&
801         dr_check_user_stop_result) {
802         dr_resume(srh);
803         return (rc);
804     }
805
806     if (!force) {
807         struct dr_ref drc = {0};
808
809         prom_printf("\nDR: checking devices...\n");
810         dev_errs_idx = 0;
811
812         drc.arr = dev_errs;
813         drc.idx = &dev_errs_idx;
814         drc.len = DR_MAX_ERR_INT;
815
816         /*
817          * Since the root node can never go away, it
818          * doesn't have to be held.
819          */
820         ddi_walk_devs(ddi_root_node(), dr_check_unsafe_major, &drc);
821         if (dev_errs_idx) {
822             handle->h_err = drerr_int(ESBD_UNSAFE, dev_errs,
823                                     dev_errs_idx, 1);
824             dr_resume(srh);
825             return (DDI_FAILURE);
826         }
827         PR_QR("done\n");
828     } else {
829         prom_printf("\nDR: dr_suspend invoked with force flag\n");
830     }
831
832 #ifndef SKIP_SYNC
833 /*
834  * This sync swap out all user pages
835  */
836 vfs_sync(SYNC_ALL);
837#endif
838 /*
839  * special treatment for lock manager
840  */
841 lm_cprsuspend();
```

1

new/usr/src/uts/i86pc/io/dr/dr_quiesce.c

```
844 #ifndef SKIP_SYNC
845     /*
846      * sync the file system in case we never make it back
847      */
848     sync();
849#endif
850
851     /*
852      * now suspend drivers
853      */
854     prom_printf("DR: suspending drivers...\n");
855     srh->sr_suspend_state = DR_SRSTATE_DRIVER;
856     srh->sr_err_idx = 0;
857     /* No parent to hold busy */
858     if ((rc = dr_suspend_devices(ddi_root_node(), srh)) != DDI_SUCCESS) {
859         if (srh->sr_err_idx && srh->sr_dr_handlep) {
860             (srh->sr_dr_handlep)->h_err = drerr_int(ESBD_SUSPEND,
861                                         srh->sr_err_ints, srh->sr_err_idx, 1);
862         }
863         dr_resume(srh);
864     }
865     return (rc);
866
867     drmach_suspend_last();
868
869     /*
870      * finally, grab all cpus
871      */
872     srh->sr_suspend_state = DR_SRSTATE_FULL;
873
874     mutex_enter(&cpu_lock);
875     pause_cpus(NULL, NULL);
876     pause_cpus(NULL);
877     dr_stop_intr();
878
879 }  
_____unchanged_portion_omitted_____
```

2

```
*****
4291 Thu May 1 19:50:20 2014
new/usr/src/uts/i86pc/io/ppm/acpisleep.c
XXXX pass in cpu_pause_func via pause_cpus
*****
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 #include <sys/types.h>
28 #include <sys/smp_impldefs.h>
29 #include <sys/promif.h>
31 #include <sys/kmem.h>
32 #include <sys/archsysm.h>
33 #include <sys/cpuvar.h>
34 #include <sys/pte.h>
35 #include <vm/seg_kmem.h>
36 #include <sys/epm.h>
37 #include <sys/cpr.h>
38 #include <sys/machsysm.h>
39 #include <sys/clock.h>
41 #include <sys/cpr_wakecode.h>
42 #include <sys/acpi/acpi.h>
44 #ifdef OLDPMCODE
45 #include "acpi.h"
46 #endif
48 #include <sys/x86_archext.h>
49 #include <sys/reboot.h>
50 #include <sys/cpu_module.h>
51 #include <sys/kdi.h>
53 /*
54 * S3 stuff
55 */
57 int acpi_rtc_wake = 0x0; /* wake in N seconds */
59 #if 0 /* debug */
60 static uint8_t branchbuf[64 * 1024]; /* for the HDT branch trace stuff */
61 #endif /* debug */

```

```
63 extern int boothowto;
65 #define BOOTCPU 0 /* cpu 0 is always the boot cpu */
67 extern void kernel_wc_code(void);
68 extern tod_ops_t *tod_ops;
69 extern int flushes_require_xcalls;
70 extern int tsc_gethrtime_enable;
72 extern cpuset_t cpu_ready_set;
73 extern void (*cpu_pause_func)(void *);
75 /*
76 * This is what we've all been waiting for!
77 */
78 int
79 acpi_enter_sleepstate(s3a_t *s3ap)
80 {
81     ACPI_PHYSICAL_ADDRESS wakephys = s3ap->s3a_wakephys;
82     caddr_t wakevirt = rm_platter_va;
83     /*LINTED*/
84     wakecode_t *wp = (wakecode_t *)wakevirt;
85     uint_t Sx = s3ap->s3a_state;
87     PT(PT_SWV);
88     /* Set waking vector */
89     if (AcpiSetFirmwareWakingVector(wakephys) != AE_OK) {
90         PT(PT_SWV_FAIL);
91         PMD(PMD_SX, ("Can't SetFirmwareWakingVector(%lx)\n",
92                     (long)wakephys));
93         goto insomnia;
94     }
96     PT(PT_EWE);
97     /* Enable wake events */
98     if (AcpiEnableEvent(ACPI_EVENT_POWER_BUTTON, 0) != AE_OK) {
99         PT(PT_EWE_FAIL);
100        PMD(PMD_SX, ("Can't EnableEvent(POWER_BUTTON)\n"));
101    }
102    if (acpi_rtc_wake > 0) {
103        /* clear the RTC bit first */
104        (void) AcpiWriteBitRegister(ACPI_BITREG_RT_CLOCK_STATUS, 1);
105        PT(PT_RTCW);
106        if (AcpiEnableEvent(ACPI_EVENT_RTC, 0) != AE_OK) {
107            PT(PT_RTCW_FAIL);
108            PMD(PMD_SX, ("Can't EnableEvent(RTC)\n"));
109        }
111    /*
112     * Set RTC to wake us in a wee while.
113     */
114    mutex_enter(&tod_lock);
115    PT(PT_TOD);
116    TODOP_SETWAKE(tod_ops, acpi_rtc_wake);
117    mutex_exit(&tod_lock);
118 }
120 /*
121 * Prepare for sleep ... could've done this earlier?
122 */
123 PT(PT_SXP);
124 PMD(PMD_SX, ("Calling AcpiEnterSleepStatePrep(%d) ...\n", Sx))
125 if (AcpiEnterSleepStatePrep(Sx) != AE_OK) {
```

```
126             PMD(PMD_SX, ("... failed\n"))
127         goto insomnia;
128     }
129
130     switch (s3ap->s3a_test_point) {
131     case DEVICE_SUSPEND_TO_RAM:
132     case FORCE_SUSPEND_TO_RAM:
133     case LOOP_BACK_PASS:
134         return (0);
135     case LOOP_BACK_FAIL:
136         return (1);
137     default:
138         ASSERT(s3ap->s3a_test_point == LOOP_BACK_NONE);
139     }
140
141     /*
142      * Tell the hardware to sleep.
143      */
144     PT(PT_SXE);
145     PMD(PMD_SX, ("Calling AcpiEnterSleepState(%d) ...\\n", Sx))
146     if (AcpiEnterSleepState(Sx) != AE_OK) {
147         PT(PT_SXE_FAIL);
148         PMD(PMD_SX, ("... failed!\\n"))
149     }
150
151 insomnia:
152     PT(PT_INSOM);
153     /* cleanup is done in the caller */
154     return (1);
155 }
```

unchanged portion omitted

```
*****  
27105 Thu May 1 19:50:20 2014  
new/usr/src/uts/i86pc/os/cpr_impl.c  
XXXX pass in cpu_pause_func via pause_cpus  
*****  
unchanged_portion_omitted  
724 /*  
725 * Stop all other cpu's before halting or rebooting. We pause the cpu's  
726 * instead of sending a cross call.  
727 * Stolen from sun4/os/mp_states.c  
728 */  
730 static int cpu_are_paused; /* sic */  
732 void  
733 i_cpr_stop_other_cpus(void)  
734 {  
735     mutex_enter(&cpu_lock);  
736     if (cpu_are_paused) {  
737         mutex_exit(&cpu_lock);  
738         return;  
739     }  
740     pause_cpus(NULL, NULL);  
740     pause_cpus(NULL);  
741     cpu_are_paused = 1;  
743     mutex_exit(&cpu_lock);  
744 }  
unchanged_portion_omitted
```

new/usr/src/uts/i86pc/os/machdep.c

```
*****
34387 Thu May 1 19:50:20 2014
new/usr/src/uts/i86pc/os/machdep.c
XXXX pass in cpu_pause_func via pause_cpus
*****
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 */

30 #include <sys/types.h>
31 #include <sys/t_lock.h>
32 #include <sys/param.h>
33 #include <sys/segments.h>
34 #include <sys/sysmacros.h>
35 #include <sys/signal.h>
36 #include <sys/systm.h>
37 #include <sys/user.h>
38 #include <sys/mman.h>
39 #include <sys/vm.h>

41 #include <sys/disp.h>
42 #include <sys/class.h>

44 #include <sys/proc.h>
45 #include <sys/buf.h>
46 #include <sys/kmem.h>

48 #include <sys/reboot.h>
49 #include <sys/uadmin.h>
50 #include <sys/callb.h>

52 #include <sys/cred.h>
53 #include <sys/vnode.h>
54 #include <sys/file.h>

56 #include <sys/procfs.h>
57 #include <sys/acct.h>

59 #include <sys/vfs.h>
60 #include <sys/dnlc.h>
61 #include <sys/var.h>
```

1

new/usr/src/uts/i86pc/os/machdep.c

```
62 #include <sys/cmn_err.h>
63 #include <sys/utsname.h>
64 #include <sys/debug.h>
66 #include <sys/dumphdr.h>
67 #include <sys/bootconf.h>
68 #include <sys/varargs.h>
69 #include <sys/promif.h>
70 #include <sys/modctl.h>
72 #include <sys/consdev.h>
73 #include <sys/frame.h>
75 #include <sys/sunddi.h>
76 #include <sys/ddidmreq.h>
77 #include <sys/psw.h>
78 #include <sys/regset.h>
79 #include <sys/privregs.h>
80 #include <sys/clock.h>
81 #include <sys/tss.h>
82 #include <sys/cpu.h>
83 #include <sys/stack.h>
84 #include <sys/trap.h>
85 #include <sys/pic.h>
86 #include <vm/hat.h>
87 #include <vm/anon.h>
88 #include <vm/as.h>
89 #include <vm/page.h>
90 #include <vm/seg.h>
91 #include <vm/seg_kmem.h>
92 #include <vm/seg_map.h>
93 #include <vm/seg_vn.h>
94 #include <vm/seg_kp.h>
95 #include <vm/hat_i86.h>
96 #include <sys/swap.h>
97 #include <sys/thread.h>
98 #include <sys/sysconf.h>
99 #include <sys/vm_machparam.h>
100 #include <sys/archsysm.h>
101 #include <sys/machsysm.h>
102 #include <sys/machlock.h>
103 #include <sys/x_call.h>
104 #include <sys/instance.h>

106 #include <sys/time.h>
107 #include <sys/smp_impldefs.h>
108 #include <sys/psm_types.h>
109 #include <sys/atomic.h>
110 #include <sys/panic.h>
111 #include <sys/cpuvar.h>
112 #include <sys/dtrace.h>
113 #include <sys/bl.h>
114 #include <sys/nvpair.h>
115 #include <sys/x86_archext.h>
116 #include <sys/pool_pset.h>
117 #include <sys/autoconf.h>
118 #include <sys/mem.h>
119 #include <sys/dumphdr.h>
120 #include <sys/compress.h>
121 #include <sys/cpu_module.h>
122 #if defined(__xpv)
123 #include <sys/hypervisor.h>
124 #include <sys/xpv_panic.h>
125 #endif
127 #include <sys/fastboot.h>
```

2

```

128 #include <sys/machelf.h>
129 #include <sys/kobj.h>
130 #include <sys/multiboot.h>

132 #ifdef _TRAPTRACE
133 #include <sys/traptrace.h>
134 #endif /* _TRAPTRACE */

136 #include <c2/audit.h>
137 #include <sys/clock_impl.h>

139 extern void audit_enterprom(int);
140 extern void audit_exitprom(int);

142 /*
143 * Tunable to enable apix PSM; if set to 0, pcplusmp PSM will be used.
144 */
145 int apix_enable = 1;

147 int apic_nvidia_io_max = 0; /* no. of NVIDIA i/o apics */

149 /*
150 * Occassionally the kernel knows better whether to power-off or reboot.
151 */
152 int force_shutdown_method = AD_UNKNOWN;

154 /*
155 * The panicbuf array is used to record messages and state:
156 */
157 char panicbuf[PANICBUFSIZE];

159 /*
160 * Flags to control Dynamic Reconfiguration features.
161 */
162 uint64_t plat_dr_options;

164 /*
165 * Maximum physical address for memory DR operations.
166 */
167 uint64_t plat_dr_physmax;

169 /*
170 * maxphys - used during physio
171 * klustsize - used for klustering by swapfs and specfs
172 */
173 int maxphys = 56 * 1024; /* XXX See vm_subr.c - max b_count in physio */
174 int klustsize = 56 * 1024;

176 caddr_t p0_va; /* Virtual address for accessing physical page 0 */

178 /*
179 * defined here, though unused on x86,
180 * to make kstat_fr.c happy.
181 */
182 int vac;

184 void debug_enter(char *);

186 extern void pm_cfb_check_and_powerup(void);
187 extern void pm_cfb_rele(void);

189 extern fastboot_info_t newkernel;

191 /*
192 * Machine dependent code to reboot.
193 * "mdep" is interpreted as a character pointer; if non-null, it is a pointer

```

```

194 * to a string to be used as the argument string when rebooting.
195 *
196 * "invoke_cb" is a boolean. It is set to true when mdboot() can safely
197 * invoke CB_CL_MDBOOT callbacks before shutting the system down, i.e. when
198 * we are in a normal shutdown sequence (interrupts are not blocked, the
199 * system is not panic'ing or being suspended).
200 */
201 /*ARGSUSED*/
202 void
203 mdboot(int cmd, int fcn, char *mdep, boolean_t invoke_cb)
204 {
205     processorid_t bootcpuid = 0;
206     static int is_first_quiesce = 1;
207     static int is_first_reset = 1;
208     int reset_status = 0;
209     static char fallback_str[] = "Falling back to regular reboot.\n";
210
211     if (fcn == AD_FASTREBOOT && !newkernel.fi_valid)
212         fcn = AD_BOOT;
213
214     if (!panicstr) {
215         kp preempt_disable();
216         if (fcn == AD_FASTREBOOT) {
217             mutex_enter(&cpu_lock);
218             if (CPU_ACTIVE(cpu_get(bootcpuid))) {
219                 affinity_set(bootcpuid);
220             }
221             mutex_exit(&cpu_lock);
222         } else {
223             affinity_set(CPU_CURRENT);
224         }
225     }
226
227     if (force_shutdown_method != AD_UNKNOWN)
228         fcn = force_shutdown_method;
229
230     /*
231      * XXX - rconsvp is set to NULL to ensure that output messages
232      * are sent to the underlying "hardware" device using the
233      * monitor's printf routine since we are in the process of
234      * either rebooting or halting the machine.
235      */
236     rconsvp = NULL;
237
238     /*
239      * Print the reboot message now, before pausing other cpus.
240      * There is a race condition in the printing support that
241      * can deadlock multiprocessor machines.
242      */
243     if (!(fcn == AD_HALT || fcn == AD_POWEROFF))
244         prom_printf("rebooting...\n");
245
246     if (IN_XPV_PANIC())
247         reset();
248
249     /*
250      * We can't bring up the console from above lock level, so do it now
251      */
252     pm_cfb_check_and_powerup();
253
254     /* make sure there are no more changes to the device tree */
255     devtree_freeze();
256
257     if (invoke_cb)
258         (void) callb_execute_class(CB_CL_MDBOOT, NULL);

```

```

260     /*
261      * Clear any unresolved UEs from memory.
262      */
263     page_retire_mdboot();

264 #if defined(__xpv)
265     /*
266      * XXXPV Should probably think some more about how we deal
267      * with panicing before it's really safe to panic.
268      * On hypervisors, we reboot very quickly.. Perhaps panic
269      * should only attempt to recover by rebooting if,
270      * say, we were able to mount the root filesystem,
271      * or if we successfully launched init(1m).
272      */
273     if (panicstr && proc_init == NULL)
274         (void) HYPERVISOR_shutdown(SHUTDOWN_poweroff);
275 #endif
276     /*
277      * stop other cpus and raise our priority. since there is only
278      * one active cpu after this, and our priority will be too high
279      * for us to be preempted, we're essentially single threaded
280      * from here on out.
281      */
282     (void) spl6();
283     if (!panicstr) {
284         mutex_enter(&cpu_lock);
285         pause_cpus(NULL, NULL);
286         pause_cpus(NULL);
287         mutex_exit(&cpu_lock);
288     }

289     /*
290      * If the system is panicking, the preloaded kernel is valid, and
291      * fastreboot_onpanic has been set, and the system has been up for
292      * longer than fastreboot_onpanic_uptime (default to 10 minutes),
293      * choose Fast Reboot.
294      */
295     if (fcn == AD_BOOT && panicstr && newkernel.fi_valid &&
296         fastreboot_onpanic &&
297         (panic_lbolt - lbolt_at_boot) > fastreboot_onpanic_uptime) {
298         fcn = AD_FASTREBOOT;
299     }

300     /*
301      * Try to quiesce devices.
302      */
303     if (is_first_quiesce) {
304         /*
305          * Clear is_first_quiesce before calling quiesce_devices()
306          * so that if quiesce_devices() causes panics, it will not
307          * be invoked again.
308          */
309         is_first_quiesce = 0;

310         quiesce_active = 1;
311         quiesce_devices(ddi_root_node(), &reset_status);
312         if (reset_status == -1) {
313             if (fcn == AD_FASTREBOOT && !force_fastreboot) {
314                 prom_printf("Driver(s) not capable of fast "
315                             "reboot.\n");
316                 prom_printf(fallback_str);
317                 fastreboot_capable = 0;
318                 fcn = AD_BOOT;
319             } else if (fcn != AD_FASTREBOOT)
320                 fastreboot_capable = 0;
321         }
322     }

```

```

325             quiesce_active = 0;
326         }

327         /*
328          * Try to reset devices. reset_leaves() should only be called
329          * a) when there are no other threads that could be accessing devices,
330          * and
331          * b) on a system that's not capable of fast reboot (fastreboot_capable
332          * being 0), or on a system where quiesce_devices() failed to
333          * complete (quiesce_active being 1).
334          */
335         if (is_first_reset && (!fastreboot_capable || quiesce_active)) {
336             /*
337              * Clear is_first_reset before calling reset_devices()
338              * so that if reset_devices() causes panics, it will not
339              * be invoked again.
340              */
341             is_first_reset = 0;
342             reset_leaves();
343         }

344         /* Verify newkernel checksum */
345         if (fastreboot_capable && fcn == AD_FASTREBOOT &&
346             fastboot_cksum_verify(&newkernel) != 0) {
347             fastreboot_capable = 0;
348             prom_printf("Fast reboot: checksum failed for the new "
349                         "kernel.\n");
350             prom_printf(fallback_str);
351         }

352         /*
353          * (void) spl8();
354          */
355         if (fastreboot_capable && fcn == AD_FASTREBOOT) {
356             /*
357              * psm_shutdown is called within fast_reboot()
358              */
359             fast_reboot();
360         } else {
361             (*psm_shutdownf)(cmd, fcn);
362             if (fcn == AD_HALT || fcn == AD_POWEROFF)
363                 halt((char *)NULL);
364             else
365                 prom_reboot("");
366         }
367     }
368     /*NOTREACHED*/
369 }
370
371 } unchanged_portion_omitted

```

```
new/usr/src/uts/i86pc/os/mp_pc.c
```

```
*****  
 17010 Thu May 1 19:50:20 2014  
new/usr/src/uts/i86pc/os/mp_pc.c  
XXXX pass in cpu_pause_func via pause_cpus  
*****
```

```
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 2011 Joyent, Inc. All rights reserved.  
30 */  
31 /*  
32 * Welcome to the world of the "real mode platter".  
33 * See also startup.c, mpcore.s and apic.c for related routines.  
34 */  
35 #include <sys/types.h>  
36 #include <sys/sysm.h>  
37 #include <sys/cpuvar.h>  
38 #include <sys/cpu_module.h>  
39 #include <sys/kmem.h>  
40 #include <sys/archsysm.h>  
41 #include <sys/machsysm.h>  
42 #include <sys/controlregs.h>  
43 #include <sys/x86_archext.h>  
44 #include <sys/smp_impldefs.h>  
45 #include <sys/sysmacros.h>  
46 #include <sys/mach_mmu.h>  
47 #include <sys/promif.h>  
48 #include <sys/cpu.h>  
49 #include <sys/cpu_event.h>  
50 #include <sys/sunndi.h>  
51 #include <sys/fs/dv_node.h>  
52 #include <vm/hat_i86.h>  
53 #include <vm/as.h>  
54  
55 extern cpuset_t cpu_ready_set;  
56  
57 extern int mp_start_cpu_common(cpu_t *cp, boolean_t boot);  
58 extern void real_mode_start_cpu(void);  
59 extern void real_mode_start_cpu_end(void);
```

```
1
```

```
new/usr/src/uts/i86pc/os/mp_pc.c
```

```
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 extern void (*cpu_pause_func)(void *);  
67 void rmpl_gdt_init(rmpl_platter_t *);  
68 /*  
69  * Fill up the real mode platter to make it easy for real mode code to  
70  * kick it off. This area should really be one passed by boot to kernel  
71  * and guaranteed to be below 1MB and aligned to 16 bytes. Should also  
72  * have identical physical and virtual address in paged mode.  
73  */  
74 static ushort_t *warm_reset_vector = NULL;  
75  
76 int  
77 mach_cpucontext_init(void)  
78 {  
79     ushort_t *vec;  
80     ulong_t addr;  
81     struct rmplatter *rm = (struct rmplatter *)rmplatter_va;  
82  
83     if (!(vec = (ushort_t *)psm_map_phys(WARM_RESET_VECTOR,  
84         sizeof(vec), PROT_READ | PROT_WRITE)))  
85         return (-1);  
86  
87     /*  
88      * setup secondary cpu bios boot up vector  
89      * Write page offset to 0x467 and page frame number to 0x469.  
90      */  
91     addr = (ulong_t)((caddr_t)rm->rm_code - (caddr_t)rm) + rmplatter_pa;  
92     vec[0] = (ushort_t)(addr & PAGEOFFSET);  
93     vec[1] = (ushort_t)((addr & (0xffff & PAGEMASK)) >> 4);  
94     warm_reset_vector = vec;  
95  
96     /* Map real mode platter into kas so kernel can access it. */  
97     hat_devload(kas.a_hat,  
98             (caddr_t)(uintptr_t)rmplatter_pa, MMU_PAGESIZE,  
99             btop(rmplatter_pa), PROT_READ | PROT_WRITE | PROT_EXEC,  
100            HAT_LOAD_NOCONSIST);  
101  
102     /* Copy CPU startup code to rmplatter if it's still during boot. */  
103     if (!plat_dr_enabled()) {  
104         ASSERT((size_t)real_mode_start_cpu_end -  
105                 (size_t)real_mode_start_cpu <= RM_PLATTER_CODE_SIZE);  
106         bcopy((caddr_t)real_mode_start_cpu, (caddr_t)rm->rm_code,  
107               (size_t)real_mode_start_cpu_end -  
108               (size_t)real_mode_start_cpu);  
109     }  
110  
111     return (0);  
112 }  
113 }  
_____ unchanged_portion_omitted _____
```

```
2
```

```
*****  
18803 Thu May 1 19:50:21 2014  
new/usr/src/uts/i86pc/os/x_call.c  
XXXX pass in cpu_pause_func via pause_cpus  
*****  
unchanged portion omitted  
266 #define XC_FLUSH_MAX_WAITS 1000  
268 /* Flush inflight message buffers. */  
269 int  
270 xc_flush_cpu(struct cpu *cpup)  
271 {  
272     int i;  
274     ASSERT((cpup->cpu_flags & CPU_READY) == 0);  
276     /*  
277      * Pause all working CPUs, which ensures that there's no CPU in  
278      * function xc_common().  
279      * This is used to work around a race condition window in xc_common()  
280      * between checking CPU_READY flag and increasing working item count.  
281      */  
282     pause_cpus(cpup, NULL);  
282     pause_cpus(cpup);  
283     start_cpus();  
285     for (i = 0; i < XC_FLUSH_MAX_WAITS; i++) {  
286         if (cpup->cpu_m.xc_work_cnt == 0) {  
287             break;  
288         }  
289         DELAY(1);  
290     }  
291     for (; i < XC_FLUSH_MAX_WAITS; i++) {  
292         if (!BT_TEST(xc_priority_set, cpup->cpu_id)) {  
293             break;  
294         }  
295         DELAY(1);  
296     }  
298     return (i >= XC_FLUSH_MAX_WAITS ? ETIME : 0);  
299 }  
unchanged portion omitted
```

```
*****
2591 Thu May 1 19:50:21 2014
new/usr/src/uts/i86xpv/os/mp_xen.c
XXXX pass in cpu_pause_func via pause_cpus
*****
_____unchanged_portion_omitted_____
579 void
580 mp_enter_barrier(void)
581 {
582     hrtime_t last_poke_time = 0;
583     int poke_allowed = 0;
584     int done = 0;
585     int i;
586
587     ASSERT(MUTEX_HELD(&cpu_lock));
588
589     pause_cpus(NULL, NULL);
589     pause_cpus(NULL);
590
591     while (!done) {
592         done = 1;
593         poke_allowed = 0;
594
595         if (xpv_gethrtime() - last_poke_time > POKE_TIMEOUT) {
596             last_poke_time = xpv_gethrtime();
597             poke_allowed = 1;
598         }
599
600         for (i = 0; i < NCPU; i++) {
601             cpu_t *cp = cpu_get(i);
602
603             if (cp == NULL || cp == CPU)
604                 continue;
605
606             switch (cpu_phase[i]) {
607                 case CPU_PHASE_NONE:
608                     cpu_phase[i] = CPU_PHASE_WAIT_SAFE;
609                     poke_cpu(i);
610                     done = 0;
611                     break;
612
613                 case CPU_PHASE_WAIT_SAFE:
614                     if (poke_allowed)
615                         poke_cpu(i);
616                     done = 0;
617                     break;
618
619                 case CPU_PHASE_SAFE:
620                 case CPU_PHASE_POWERED_OFF:
621                     break;
622             }
623
624             SMT_PAUSE();
625         }
626     }
627 }
_____unchanged_portion_omitted_____
```

```
*****  
6667 Thu May 1 19:50:21 2014  
new/usr/src/uts/sun4/os/mp_states.c  
XXXX pass in cpu_pause_func via pause_cpus  
*****  
unchanged_portion_omitted  
187 /*  
188  * Stop all other cpu's before halting or rebooting. We pause the cpu's  
189  * instead of sending a cross call.  
190 */  
191 void  
192 stop_other_cpus(void)  
193 {  
194     mutex_enter(&cpu_lock);  
195     if (cpu_are_paused) {  
196         mutex_exit(&cpu_lock);  
197         return;  
198     }  
200     if (ncpus > 1)  
201         intr_redist_all_cpus_shutdown();  
203     pause_cpus(NULL, NULL);  
203     pause_cpus(NULL);  
204     cpu_are_paused = 1;  
206     mutex_exit(&cpu_lock);  
207 }  
unchanged_portion_omitted
```

```
*****  
16840 Thu May 1 19:50:21 2014  
new/usr/src/uts/sun4/os/prom_subr.c  
XXXX pass in cpu_pause_func via pause_cpus  
*****  
unchanged_portion_omitted  
406 /*  
407 * This routine is a special form of pause_cpus(). It ensures that  
408 * prom functions are callable while the cpus are paused.  
409 */  
410 void  
411 promsafe_pause_cpus(void)  
412 {  
413     pause_cpus(NULL, NULL);  
414     pause_cpus(NULL);  
415     /* If some other cpu is entering or is in the prom, spin */  
416     while (prom_cpu || mutex_owner(&prom_mutex)) {  
417         start_cpus();  
418         mutex_enter(&prom_mutex);  
419  
420         /* Wait for other cpu to exit prom */  
421         while (prom_cpu)  
422             cv_wait(&prom_cv, &prom_mutex);  
423  
424         mutex_exit(&prom_mutex);  
425         pause_cpus(NULL, NULL);  
426         pause_cpus(NULL);  
427     }  
428     /* At this point all cpus are paused and none are in the prom */  
429 }  
unchanged_portion_omitted
```

```
*****
24533 Thu May 1 19:50:21 2014
new/usr/src/uts/sun4u/io/mem_cache.c
XXXX pass in cpu_pause_func via pause_cpus
*****
_____unchanged_portion_omitted_____
553 static int
554 mem_cache_ioctl_ops(int cmd, int mode, cache_info_t *cache_info)
555 {
556     int ret_val = 0;
557     uint64_t afar, tag_addr;
558     ch_cpu_logout_t clop;
559     uint64_t Lxcache_tag_data[PN_CACHE_NWAYS];
560     int i, retire_retry_count;
561     cpu_t *cpu;
562     uint64_t tag_data;
563     uint8_t state;
564
565     if (cache_info->way >= PN_CACHE_NWAYS)
566         return (EINVAL);
567     switch (cache_info->cache) {
568         case L2_CACHE_TAG:
569         case L2_CACHE_DATA:
570             if (cache_info->index >=
571                 (PN_L2_SET_SIZE/PN_L2_LINESIZE))
572                 return (EINVAL);
573             break;
574         case L3_CACHE_TAG:
575         case L3_CACHE_DATA:
576             if (cache_info->index >=
577                 (PN_L3_SET_SIZE/PN_L3_LINESIZE))
578                 return (EINVAL);
579             break;
580     default:
581         return (ENOTSUP);
582     }
583     /*
584      * Check if we have a valid cpu ID and that
585      * CPU is ONLINE.
586      */
587     mutex_enter(&cpu_lock);
588     cpu = cpu_get(cache_info->cpu_id);
589     if ((cpu == NULL) || (!cpu_is_online(cpu))) {
590         mutex_exit(&cpu_lock);
591         return (EINVAL);
592     }
593     mutex_exit(&cpu_lock);
594     pattern = 0; /* default value of TAG PA when cacheline is retired. */
595     switch (cmd) {
596         case MEM_CACHE_RETIRE:
597             tag_addr = get_tag_addr(cache_info);
598             pattern |= PN_ESTATE_NA;
599             retire_retry_count = 0;
600             affinity_set(cache_info->cpu_id);
601             switch (cache_info->cache) {
602                 case L2_CACHE_DATA:
603                 case L2_CACHE_TAG:
604                     if (((cache_info->bit & MSB_BIT_MASK) ==
605                         MSB_BIT_MASK)
606                         pattern |= PN_L2TAG_PA_MASK;
607             retry_l2_retire:
608                 if (tag_addr_collides(tag_addr,
609                     cache_info->cache,
610                     retire_l2_start, retire_l2_end))
611                     ret_val =
```

```
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
    retire_l2_alternate(
        tag_addr, pattern);
    else
        ret_val = retire_l2(tag_addr,
            pattern);
    if (ret_val == 1) {
        /*
         * cacheline was in retired
         * STATE already.
         * so return success.
         */
        ret_val = 0;
    }
    if (ret_val < 0) {
        cmn_err(CE_WARN,
            "retire_l2() failed. index = 0x%x way %d. Retrying...\n",
            cache_info->index,
            cache_info->way);
        if (retire_retry_count >= 2) {
            retire_failures++;
            affinity_clear();
            return (EIO);
        }
        retire_retry_count++;
        goto retry_l2_retire;
    }
    if (ret_val == 2)
        l2_flush_retries_done++;
/*
 * We bind ourselves to a CPU and send cross trap to
 * ourselves. On return from xt_one we can rely on the
 * data in tag_data being filled in. Normally one would
 * do a xt_sync to make sure that the CPU has completed
 * the cross trap call xt_one.
*/
xt_one(cache_info->cpu_id,
    (xfcfunc_t *)(get_l2_tag_t11),
    tag_addr, (uint64_t)(&tag_data));
state = tag_data & CH_ESTATE_MASK;
if (state != PN_ESTATE_NA) {
    retire_failures++;
    print_l2_tag(tag_addr,
        tag_data);
    cmn_err(CE_WARN,
        "L2 RETIRE:failed for index 0x%x way %d. Retrying...\n",
        cache_info->index,
        cache_info->way);
    if (retire_retry_count >= 2) {
        retire_failures++;
        affinity_clear();
        return (EIO);
    }
    retire_retry_count++;
    goto retry_l2_retire;
}
break;
case L3_CACHE_TAG:
case L3_CACHE_DATA:
    if (((cache_info->bit & MSB_BIT_MASK) ==
        MSB_BIT_MASK)
        pattern |= PN_L3TAG_PA_MASK;
    if (tag_addr_collides(tag_addr,
        cache_info->cache,
        retire_l3_start, retire_l3_end))
        ret_val =
            retire_l3_alternate(
```

```

678
679         tag_addr, pattern);
680     else
681         ret_val = retire_l3(tag_addr,
682                             pattern);
683     if (ret_val == 1) {
684         /*
685          * cacheline was in retired
686          * STATE already.
687          * so return success.
688         */
689     ret_val = 0;
690     }
691     if (ret_val < 0) {
692         cmn_err(CE_WARN,
693                 "retire_l3() failed. ret_val = %d index = 0x%x\n",
694                 ret_val,
695                 cache_info->index);
696     retire_failures++;
697     affinity_clear();
698     return (EIO);
699   }
700   /*
701    * We bind ourself to a CPU and send cross trap to
702    * ourselves. On return from xt_one we can rely on the
703    * data in tag_data being filled in. Normally one would
704    * do a xt_sync to make sure that the CPU has completed
705    * the cross trap call xt_one.
706   */
707   xt_one(cache_info->cpu_id,
708         (xcfunc_t *) (get_l3_tag_t11),
709         tag_addr, (uint64_t) (&tag_data));
710   state = tag_data & CH_ECSTATE_MASK;
711   if (state != PN_ECSTATE_NA) {
712     cmn_err(CE_WARN,
713             "L3 RETIRE failed for index 0x%x\n",
714             cache_info->index);
715     retire_failures++;
716     affinity_clear();
717     return (EIO);
718   }
719   break;
720 }
721 affinity_clear();
722 break;
723 case MEM_CACHE_UNRETIRE:
724   tag_addr = get_tag_addr(cache_info);
725   pattern = PN_ECSTATE_INV;
726   affinity_set(cache_info->cpu_id);
727   switch (cache_info->cache) {
728     case L2_CACHE_DATA:
729     case L2_CACHE_TAG:
730   /*
731    * We bind ourself to a CPU and send cross trap to
732    * ourselves. On return from xt_one we can rely on the
733    * data in tag_data being filled in. Normally one would
734    * do a xt_sync to make sure that the CPU has completed
735    * the cross trap call xt_one.
736   */
737   xt_one(cache_info->cpu_id,
738         (xcfunc_t *) (get_l2_tag_t11),
739         tag_addr, (uint64_t) (&tag_data));
740   state = tag_data & CH_ECSTATE_MASK;
741   if (state != PN_ECSTATE_NA) {
742     affinity_clear();
743     return (EINVAL);
744   }
745   if (tag_addr_collides(tag_addr,
746                         cache_info->cache,
747                         unretire_l2_start, unretire_l2_end))
748     ret_val =
749       unretire_l2_alternate(
750       tag_addr, pattern);
751   else
752     ret_val =
753       unretire_l2(tag_addr,
754                  pattern);
755   if (ret_val != 0) {
756     cmn_err(CE_WARN,
757             "unretire_l2() failed. ret_val = %d index = 0x%x\n",
758             ret_val,
759             cache_info->index);
760   retire_failures++;
761   affinity_clear();
762   return (EIO);
763 }
764 break;
765 case L3_CACHE_TAG:
766 case L3_CACHE_DATA:
767 /*
768  * We bind ourself to a CPU and send cross trap to
769  * ourselves. On return from xt_one we can rely on the
770  * data in tag_data being filled in. Normally one would
771  * do a xt_sync to make sure that the CPU has completed
772  * the cross trap call xt_one.
773 */
774 xt_one(cache_info->cpu_id,
775       (xcfunc_t *) (get_l3_tag_t11),
776       tag_addr, (uint64_t) (&tag_data));
777 state = tag_data & CH_ECSTATE_MASK;
778 if (state != PN_ECSTATE_NA) {
779   affinity_clear();
780   return (EINVAL);
781 }
782 if (tag_addr_collides(tag_addr,
783                       cache_info->cache,
784                       unretire_l3_start, unretire_l3_end))
785   ret_val =
786     unretire_l3_alternate(
787     tag_addr, pattern);
788 else
789   ret_val =
790     unretire_l3(tag_addr,
791                 pattern);
792 if (ret_val != 0) {
793   cmn_err(CE_WARN,
794           "unretire_l3() failed. ret_val = %d index = 0x%x\n",
795           ret_val,
796           cache_info->index);
797   affinity_clear();
798   return (EIO);
799 }
800 break;
801 affinity_clear();
802 break;
803 case MEM_CACHE_ISRETIRED:
804 case MEM_CACHE_STATE:
805   return (ENOTSUP);
806 case MEM_CACHE_READ_TAGS:
807   break;
808 #ifdef DEBUG
809 case MEM_CACHE_READ_ERROR_INJECTED_TAGS:
810   break;
811 
```

```

812   }
813   if (tag_addr_collides(tag_addr,
814                         cache_info->cache,
815                         unretire_l2_start, unretire_l2_end))
816     ret_val =
817       unretire_l2_alternate(
818       tag_addr, pattern);
819   else
820     ret_val =
821       unretire_l2(tag_addr,
822                  pattern);
823   if (ret_val != 0) {
824     cmn_err(CE_WARN,
825             "unretire_l2() failed. ret_val = %d index = 0x%x\n",
826             ret_val,
827             cache_info->index);
828   retire_failures++;
829   affinity_clear();
830   return (EIO);
831 }
832 break;
833 
```

```

810 #endif
811
812     /*
813      * Read tag and data for all the ways at a given afar
814      */
815     afar = (uint64_t)(cache_info->index
816                       << PN_CACHE_LINE_SHIFT);
817     mutex_enter(&cpu_lock);
818     affinity_set(cache_info->cpu_id);
819     (void) pause_cpus(NULL, NULL);
820     (void) pause_cpus(NULL);
821     mutex_exit(&cpu_lock);
822
823     /*
824      * We bind ourself to a CPU and send cross trap to
825      * ourselves. On return from xt_one we can rely on the
826      * data in clop being filled in. Normally one would
827      * do a xt_sync to make sure that the CPU has completed
828      * the cross trap call xt_one.
829      */
830     xt_one(cache_info->cpu_id,
831            (xfunc_t *) (get_ecache_dtags_t11),
832            afar, (uint64_t) (&clop));
833     mutex_enter(&cpu_lock);
834     (void) start_cpus();
835     mutex_exit(&cpu_lock);
836     affinity_clear();
837     switch (cache_info->cache) {
838         case L2_CACHE_TAG:
839             for (i = 0; i < PN_CACHE_NWAYS; i++) {
840                 Lxcache_tag_data[i] =
841                     clop.clo_data.chd_l2_data
842                     [i].ec_tag;
843             }
844             last_error_injected_bit =
845                 last_l2tag_error_injected_bit;
846             last_error_injected_way =
847                 last_l2tag_error_injected_way;
848             break;
849         case L3_CACHE_TAG:
850             for (i = 0; i < PN_CACHE_NWAYS; i++) {
851                 Lxcache_tag_data[i] =
852                     clop.clo_data.chd_ec_data
853                     [i].ec_tag;
854             }
855             last_error_injected_bit =
856                 last_l3tag_error_injected_bit;
857             last_error_injected_way =
858                 last_l3tag_error_injected_way;
859     #endif
860
861         default:
862             return (ENOTSUP);
863     } /* end if switch(cache) */
864
865     if ((cmd == MEM_CACHE_READ_ERROR_INJECTED_TAGS) &&
866         (inject_anonymous_tag_error == 0) &&
867         (last_error_injected_way >= 0) &&
868         (last_error_injected_way <= 3)) {
869         pattern = ((uint64_t)1 <-
870                     last_error_injected_bit);
871
872         /*
873          * If error bit is ECC we need to make sure
874          * ECC on all all WAYS are corrupted.
875         */

```

```

875
876
877
878
879
880
881
882
883
884
885 #endif
886
887     if ((last_error_injected_bit >= 6) &&
888         (last_error_injected_bit <= 14)) {
889         for (i = 0; i < PN_CACHE_NWAYS; i++) {
890             Lxcache_tag_data[i] ^
891                 pattern;
892         }
893     } else
894         Lxcache_tag_data
895             [last_error_injected_way] ^
896                 pattern;
897
898
899     if (ddi_copyout((caddr_t)Lxcache_tag_data,
900                     (caddr_t)cache_info->datap,
901                     sizeof (Lxcache_tag_data), mode)
902                     != DDI_SUCCESS) {
903             return (EFAULT);
904         }
905     break; /* end of READ_TAGS */
906     default:
907         return (ENOTSUP);
908     } /* end if switch(cmd) */
909     return (ret_val);
910 }
```

unchanged_portion_omitted

```
*****
25005 Thu May 1 19:50:22 2014
new/usr/src/uts/sun4u/ngdr/io/dr_quiesce.c
XXXX pass in cpu_pause_func via pause_cpus
*****
_____unchanged_portion_omitted_____
822 int
823 dr_suspend(dr_sr_handle_t *srh)
824 {
825     dr_handle_t    *handle;
826     int            force;
827     int            dev_errs_idx;
828     uint64_t       dev_errs[DR_MAX_ERR_INT];
829     int            rc = DDI_SUCCESS;
831
832     handle = srh->sr_dr_handlep;
833
834     force = dr_cmd_flags(handle) & SBD_FLAG_FORCE;
835
836     /* update the signature block
837     */
838     CPU_SIGNATURE(OS_SIG, SIGST QUIESCE_INPROGRESS, SIGSUBST_NULL,
839                   CPU->cpu_id);
840
841     prom_printf("\nDR: suspending user threads...\n");
842     srh->sr_suspend_state = DR_SRSTATE_USER;
843     if (((rc = dr_stop_user_threads(srh)) != DDI_SUCCESS) &&
844         dr_check_user_stop_result) {
845         dr_resume(srh);
846         return (rc);
847     }
848
849     if (!force) {
850         struct dr_ref drc = {0};
851
852         prom_printf("\nDR: checking devices...\n");
853         dev_errs_idx = 0;
854
855         drc.arr = dev_errs;
856         drc.idx = &dev_errs_idx;
857         drc.len = DR_MAX_ERR_INT;
858
859         /*
860         * Since the root node can never go away, it
861         * doesn't have to be held.
862         */
863         ddi_walk_devs(ddi_root_node(), dr_check_unsafe_major, &drc);
864         if (dev_errs_idx) {
865             handle->h_err = drerr_int(ESBD_UNSAFE, dev_errs,
866                                       dev_errs_idx, 1);
867             dr_resume(srh);
868             return (DDI_FAILURE);
869         }
870         PR_QR("done\n");
871     } else {
872         prom_printf("\nDR: dr_suspend invoked with force flag\n");
873     }
874
875 #ifndef SKIP_SYNC
876     /*
877      * This sync swap out all user pages
878      */
879     vfs_sync(SYNC_ALL);
880 #endif

```

```
882     /*
883      * special treatment for lock manager
884      */
885     lm_cprssuspend();
886
887 #ifndef SKIP_SYNC
888     /*
889      * sync the file system in case we never make it back
890      */
891     sync();
892 #endif
893
894     /*
895      * now suspend drivers
896      */
897     prom_printf("DR: suspending drivers...\n");
898     srh->sr_suspend_state = DR_SRSTATE_DRIVER;
899     srh->sr_err_idx = 0;
900
901     /* No parent to hold busy */
902     if ((rc = dr_suspend_devices(ddi_root_node(), srh)) != DDI_SUCCESS) {
903         if (srh->sr_err_idx && srh->sr_dr_handlep) {
904             (srh->sr_dr_handlep)->h_err = drerr_int(ESBD_SUSPEND,
905                                         srh->sr_err_ints, srh->sr_err_idx, 1);
906         }
907         dr_resume(srh);
908     }
909
910     drmach_suspend_last();
911
912     /*
913      * finally, grab all cpus
914      */
915     srh->sr_suspend_state = DR_SRSTATE_FULL;
916
917     /*
918      * if watchdog was activated, disable it
919      */
920     if (watchdog_activated) {
921         mutex_enter(&tod_lock);
922         tod_ops.tod_clear_watchdog_timer();
923         mutex_exit(&tod_lock);
924         srh->sr_flags |= SR_FLAG_WATCHDOG;
925     } else {
926         srh->sr_flags &= ~(SR_FLAG_WATCHDOG);
927     }
928
929     /*
930      * Update the signature block.
931      * This must be done before cpus are paused, since on Starcat the
932      * cpu signature update aquires an adaptive mutex in the iosram driver.
933      * Blocking with cpus paused can lead to deadlock.
934      */
935     CPU_SIGNATURE(OS_SIG, SIGST QUIESCED, SIGSUBST_NULL, CPU->cpu_id);
936
937     mutex_enter(&cpu_lock);
938     pause_cpus(NULL, NULL);
939     pause_cpus(NULL);
940     dr_stop_intr();
941
942 }
943
944 _____unchanged_portion_omitted_____

```

```
new/usr/src/uts/sun4u/os/cpr_impl.c
```

```
*****
50490 Thu May 1 19:50:22 2014
new/usr/src/uts/sun4u/os/cpr_impl.c
XXXX pass in cpu_pause_func via pause_cpus
*****
_____unchanged_portion_omitted_____
```

```
215 /*  
216  * launch slave cpus into kernel text, pause them,  
217  * and restore the original prom pages  
218 */  
219 void  
220 i_cpr_mp_setup(void)  
221 {  
222     extern void restart_other_cpu(int);  
223     cpu_t *cp;  
224  
225     uint64_t kcontextreg;  
226  
227     /*  
228      * Do not allow setting page size codes in MMU primary context  
229      * register while using cif wrapper. This is needed to work  
230      * around OBP incorrect handling of this MMU register.  
231      */  
232     kcontextreg = 0;  
233  
234     /*  
235      * reset cpu_ready_set so x_calls work properly  
236      */  
237     CPUSET_ZERO(cpu_ready_set);  
238     CPUSET_ADD(cpu_ready_set, getprocessorid());  
239  
240     /*  
241      * setup cif to use the cookie from the new/tmp prom  
242      * and setup tmp handling for calling prom services.  
243      */  
244     i_cpr_cif_setup(CIF_SPLICE);  
245  
246     /*  
247      * at this point, only the nucleus and a few cpr pages are  
248      * mapped in. once we switch to the kernel trap table,  
249      * we can access the rest of kernel space.  
250      */  
251     prom_set_traptable(&trap_table);  
252  
253     if (ncpus > 1) {  
254         sfmmu_init_tsbs();  
255  
256         mutex_enter(&cpu_lock);  
257         /*  
258          * All of the slave cpus are not ready at this time,  
259          * yet the cpu structures have various cpu_flags set;  
260          * clear cpu_flags and mutex_ready.  
261          * Since we are coming up from a CPU suspend, the slave cpus  
262          * are frozen.  
263          */  
264         for (cp = CPU->cpu_next; cp != CPU; cp = cp->cpu_next) {  
265             cp->cpu_flags = CPU_FROZEN;  
266             cp->cpu_m.mutex_ready = 0;  
267         }  
268  
269         for (cp = CPU->cpu_next; cp != CPU; cp = cp->cpu_next)  
270             restart_other_cpu(cp->cpu_id);  
271  
272         pause_cpus(NULL, NULL);  
273     }
```

```
1
```

```
new/usr/src/uts/sun4u/os/cpr_impl.c
```

```
272     pause_cpus(NULL);  
273     mutex_exit(&cpu_lock);  
274  
275     } else  
276         i_cpr_xcall(i_cpr_clear_entries);  
277     i_cpr_clear_entries(0, 0);  
278  
279     /*  
280      * now unlink the cif wrapper; WARNING: do not call any  
281      * prom_xxx() routines until after prom pages are restored.  
282      */  
283     i_cpr_cif_setup(CIF_UNLINK);  
284  
285     (void) i_cpr_prom_pages(CPR_PROM_RESTORE);  
286  
287     /* allow setting page size codes in MMU primary context register */  
288     kcontextreg = kctx;  
289 }  
_____unchanged_portion_omitted_____
```

```
2
```

```

new/usr/src/uts/sun4u/serengeti/io/sbdb_quiesce.c
*****
19584 Thu May 1 19:50:22 2014
new/usr/src/uts/sun4u/serengeti/io/sbdb_quiesce.c
XXXX pass in cpu_pause_func via pause_cpus
*****
_____unchanged_portion_omitted_____
764 int
765 sbdp_suspend(sbdb_sr_handle_t *srh)
766 {
767     int force;
768     int rc = DDI_SUCCESS;
770     force = (srh && (srh->sr_flags & SBDP_IOCTL_FLAG_FORCE));
772     /*
773      * if no force flag, check for unsafe drivers
774      */
775     if (force) {
776         SBDP_DBG_QR("\nsbdb_suspend invoked with force flag");
777     }
779     /*
780      * update the signature block
781      */
782     CPU_SIGNATURE(OS_SIG, SIGST QUIESCE_INPROGRESS, SIGSUBST_NULL,
783                   CPU->cpu_id);
785     /*
786      * first, stop all user threads
787      */
788     SBDP_DBG_QR("SBDP: suspending user threads...\n");
789     SR_SET_STATE(srh, SBDP_SRSTATE_USER);
790     if ((rc = sbdp_stop_user_threads(srh)) != DDI_SUCCESS) &&
791         sbdp_check_user_stop_result) {
792         sbdp_resume(srh);
793         return (rc);
794     }
796 #ifndef SKIP_SYNC
797     /*
798      * This sync swap out all user pages
799      */
800     vfs_sync(SYNC_ALL);
801 #endif
803     /*
804      * special treatment for lock manager
805      */
806     lm_cprsuspend();
808 #ifndef SKIP_SYNC
809     /*
810      * sync the file system in case we never make it back
811      */
812     sync();
814 #endif
815     /*
816      * now suspend drivers
817      */
818     SBDP_DBG_QR("SBDP: suspending drivers...\n");
819     SR_SET_STATE(srh, SBDP_SRSTATE_DRIVER);
821     /*
822      * Root node doesn't have to be held in any way.

```

```

1 new/usr/src/uts/sun4u/serengeti/io/sbdb_quiesce.c
2
823     */
824     if ((rc = sbdp_suspend_devices(ddi_root_node(), srh)) != DDI_SUCCESS) {
825         sbdp_resume(srh);
826         return (rc);
827     }
829     /*
830      * finally, grab all cpus
831      */
832     SR_SET_STATE(srh, SBDP_SRSTATE_FULL);
834     /*
835      * if watchdog was activated, disable it
836      */
837     if (watchdog_activated) {
838         mutex_enter(&tod_lock);
839         saved_watchdog_seconds = tod_ops.tod_clear_watchdog_timer();
840         mutex_exit(&tod_lock);
841         SR_SET_FLAG(srh, SR_FLAG_WATCHDOG);
842     } else {
843         SR_CLEAR_FLAG(srh, SR_FLAG_WATCHDOG);
844     }
846     mutex_enter(&cpu_lock);
847     pause_cpus(NULL, NULL);
847     pause_cpus(NULL);
848     sbdp_stop_intr();
850     /*
851      * update the signature block
852      */
853     CPU_SIGNATURE(OS_SIG, SIGST QUIESCED, SIGSUBST_NULL, CPU->cpu_id);
855 }
856 _____unchanged_portion_omitted_____

```

```
*****
54782 Thu May 1 19:50:22 2014
new/usr/src/uts/sun4v/os/mpo.c
XXXX pass in cpu_pause_func via pause_cpus
*****
_____unchanged_portion_omitted_____
210 /*
211 * The MPO locks are to protect the MPO metadata while that
212 * information is updated as a result of a memory DR operation.
213 * The read lock must be acquired to read the metadata and the
214 * write locks must be acquired to update it.
215 */
216 #define mpo_rd_lock    kpreempt_disable
217 #define mpo_rd_unlock   kpreempt_enable
218 static void
219 mpo_wr_lock()
220 {
221     mutex_enter(&cpu_lock);
222     pause_cpus(NULL, NULL);
223     pause_cpus(NULL);
224     mutex_exit(&cpu_lock);
225 }
_____unchanged_portion_omitted_____
```

new/usr/src/uts/sun4v/os/suspend.c

```
*****
21873 Thu May 1 19:50:22 2014
new/usr/src/uts/sun4v/os/suspend.c
XXXX pass in cpu_pause_func via pause_cpus
*****
_____unchanged_portion_omitted_____
348 /*
349  * Obtain an updated MD from the hypervisor and update cpunodes, CPU HW
350  * sharing data structures, and processor groups.
351 */
352 static void
353 update_cpu_mappings(void)
354 {
355     md_t          *mdp;
356     processorid_t id;
357     cpu_t          *cp;
358     cpu_pg_t      *pgps[NCPU];
359
360     if ((mdp = md_get_handle()) == NULL) {
361         DBG("suspend: md_get_handle failed");
362         return;
363     }
364
365     DBG("suspend: updating CPU mappings");
366
367     mutex_enter(&cpu_lock);
368
369     setup_chip_mappings(mdp);
370     setup_exec_unit_mappings(mdp);
371     for (id = 0; id < NCPU; id++) {
372         if ((cp = cpu_get(id)) == NULL)
373             continue;
374         cpu_map_exec_units(cp);
375     }
376
377     /*
378      * Re-calculate processor groups.
379      *
380      * First tear down all PG information before adding any new PG
381      * information derived from the MD we just downloaded. We must
382      * call pg_cpu_inactive and pg_cpu_active with CPUs paused and
383      * we want to minimize the number of times pause_cpus is called.
384      * Inactivating all CPUs would leave PGs without any active CPUs,
385      * so while CPUs are paused, call pg_cpu_inactive and swap in the
386      * bootstrap PG structure saving the original PG structure to be
387      * fini'd afterwards. This prevents the dispatcher from encountering
388      * PGs in which all CPUs are inactive. Offline CPUs are already
389      * inactive in their PGs and shouldn't be reactivated, so we must
390      * not call pg_cpu_inactive or pg_cpu_active for those CPUs.
391      */
392     pause_cpus(NULL, NULL);
393     pause_cpus(NULL);
394     for (id = 0; id < NCPU; id++) {
395         if ((cp = cpu_get(id)) == NULL)
396             continue;
397         if ((cp->cpu_flags & CPU_OFFLINE) == 0)
398             pgps[id] = cp->cpu_pg;
399         pg_cpu_bootstrap(cp);
400     }
401     start_cpus();
402
403     /*
404      * pg_cpu_fini* and pg_cpu_init* must be called while CPUs are
405      * not paused. Use two separate loops here so that we do not
```

1

2

```
new/usr/src/uts/sun4v/os/suspend.c
406     * initialize PG data for CPUs until all the old PG data structures
407     * are torn down.
408     */
409     for (id = 0; id < NCPU; id++) {
410         if ((cp = cpu_get(id)) == NULL)
411             continue;
412         pg_cpu_fini(cp, pgps[id]);
413         mpo_cpu_remove(id);
414     }
415
416     /*
417      * Initialize PG data for each CPU, but leave the bootstrapped
418      * PG structure in place to avoid running with any PGs containing
419      * nothing but inactive CPUs.
420      */
421     for (id = 0; id < NCPU; id++) {
422         if ((cp = cpu_get(id)) == NULL)
423             continue;
424         mpo_cpu_add(mdp, id);
425         pgps[id] = pg_cpu_init(cp, B_TRUE);
426     }
427
428     /*
429      * Now that PG data has been initialized for all CPUs in the
430      * system, replace the bootstrapped PG structure with the
431      * initialized PG structure and call pg_cpu_active for each CPU.
432      */
433     pause_cpus(NULL, NULL);
434     for (id = 0; id < NCPU; id++) {
435         if ((cp = cpu_get(id)) == NULL)
436             continue;
437         cp->cpu_pg = pgps[id];
438         if ((cp->cpu_flags & CPU_OFFLINE) == 0)
439             pg_cpu_active(cp);
440     }
441     start_cpus();
442
443     mutex_exit(&cpu_lock);
444
445     (void) md_fini_handle(mdp);
446 }
_____unchanged_portion_omitted_____
585 /*
586  * Suspends the OS by pausing CPUs and calling into the HV to initiate
587  * the suspend. When the HV routine hv_guest_suspend returns, the system
588  * will be resumed. Must be called after a successful call to suspend_pre.
589  * suspend_post must be called after suspend_start, whether or not
590  * suspend_start returns an error.
591 */
592 /*ARGSUSED*/
593 int
594 suspend_start(char *error_reason, size_t max_reason_len)
595 {
596     uint64_t      source_tick;
597     uint64_t      source_stick;
598     uint64_t      rv;
599     timestruc_t   source_tod;
600     int           spl;
601
602     ASSERT(suspend_supported());
603     DBG("suspend: %s", __func__);
604
605     sfmmu_ctxdoms_lock();
```

```

607     mutex_enter(&cpu_lock);
609     /* Suspend the watchdog */
610     watchdog_suspend();
612     /* Record the TOD */
613     mutex_enter(&tod_lock);
614     source_tod = tod_get();
615     mutex_exit(&tod_lock);
617     /* Pause all other CPUs */
618     pause_cpus(NULL, NULL);
619     DBG_PROM("suspend: CPUs paused\n");
621     /* Suspend cyclics */
622     cyclic_suspend();
623     DBG_PROM("suspend: cyclics suspended\n");
625     /* Disable interrupts */
626     spl = spl8();
627     DBG_PROM("suspend: spl8()\n");
629     source_tick = gettick_counter();
630     source_stick = gettick();
631     DBG_PROM("suspend: source_tick: 0x%lx\n", source_tick);
632     DBG_PROM("suspend: source_stick: 0x%lx\n", source_stick);

634     /*
635      * Call into the HV to initiate the suspend. hv_guest_suspend()
636      * returns after the guest has been resumed or if the suspend
637      * operation failed or was cancelled. After a successful suspend,
638      * the %tick and %stick registers may have changed by an amount
639      * that is not proportional to the amount of time that has passed.
640      * They may have jumped forwards or backwards. Some variation is
641      * allowed and accounted for using suspend_tick_stick_max_delta,
642      * but otherwise this jump must be uniform across all CPUs and we
643      * operate under the assumption that it is (maintaining two global
644      * offset variables--one for %tick and one for %stick.)
645      */
646     DBG_PROM("suspend: suspending... \n");
647     rv = hv_guest_suspend();
648     if (rv != 0) {
649         splx(spl);
650         cyclic_resume();
651         start_cpus();
652         watchdog_resume();
653         mutex_exit(&cpu_lock);
654         sfmmu_ctxdoms_unlock();
655         DBG("suspend: failed, rv: %d\n", rv);
656         return (rv);
657     }

659     suspend_count++;

661     /* Update the global tick and stick offsets and the preserved TOD */
662     set_tick_offsets(source_tick, source_stick, &source_tod);

664     /* Ensure new offsets are globally visible before resuming CPUs */
665     membar_sync();

667     /* Enable interrupts */
668     splx(spl);

670     /* Set the {%tick,%stick}.NPT bits on all CPUs */
671     if (enable_user_tick_stick_emulation) {

```

```

672             xc_all((xfcfunc_t *)enable_tick_stick_npt, NULL, NULL);
673             xt_sync(cpu_ready_set);
674             ASSERT(gettick_npt() != 0);
675             ASSERT(getstick_npt() != 0);
676         }

678         /* If emulation is enabled, but not currently active, enable it */
679         if (enable_user_tick_stick_emulation && !tick_stick_emulation_active) {
680             tick_stick_emulation_active = B_TRUE;
681         }

683         sfmmu_ctxdoms_remove();

685         /* Resume cyclics, unpause CPUs */
686         cyclic_resume();
687         start_cpus();

689         /* Set the TOD */
690         mutex_enter(&tod_lock);
691         tod_set(source_tod);
692         mutex_exit(&tod_lock);

694         /* Re-enable the watchdog */
695         watchdog_resume();

697         mutex_exit(&cpu_lock);

699         /* Download the latest MD */
700         if ((rv = mach_descrip_update()) != 0)
701             cmn_err(CE_PANIC, "suspend: mach_descrip_update failed: %ld",
702                     rv);

704         sfmmu_ctxdoms_update();
705         sfmmu_ctxdoms_unlock();

707         /* Get new MD, update CPU mappings/relationships */
708         if (suspend_update_cpu_mappings)
709             update_cpu_mappings();

711         DBG("suspend: target tick: 0x%lx", gettick_counter());
712         DBG("suspend: target stick: 0x%llx", gettick());
713         DBG("suspend: user %tick/%stick emulation is %d",
714             tick_stick_emulation_active);
715         DBG("suspend: finished");

717     }
718 }
```

unchanged_portion_omitted