

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

```
*****
34820 Fri Jan 3 22:11:52 2014
new/usr/src/uts/common/cpr/cpr_main.c
patch cpu-pause-func-deglobalize
*****
```

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/systm.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/nodctl.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 curthreadpfn;  
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(hrttime_t 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      */
771
772 #if defined(__x86)
773     /*
774      * Initialize our syscall handlers
775      */
776     init_cpu_syscall(CPU);
777
778 #endif
779
780     i_cpr_pre_resume_cpus();
781
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      * Restore this cpu to use the regular cpu_pause(), so that
794      * online and offline will work correctly
795      */
796     cpu_pause_func = NULL;
797
798     /*
799      * clear the affinity set in cpr_suspend_cpus()
800      */
801     affinity_clear();
802
803     /*
804      * offline all the cpus that were brought online during suspend
805      */
806     cpr_restore_offline();
807
808     mutex_exit(&cpu_lock);
809
810 void
811 cpr_unpause_cpus(void)
812 {
813     /*
814      * Now restore the system back to what it was before we suspended
815      */
816
817     PMD(PMD_SX, ("cpr_unpause_cpus: restoring system\n"));
818
819     mutex_enter(&cpu_lock);
820
821     /*
822      * Restore this cpu to use the regular cpu_pause(), so that
823      * online and offline will work correctly
824      */
825     cpu_pause_func = NULL;
826
827     /*
828      * Restart the paused cpus
829      */
830     start_cpus();
831
832     /*
833      * clear the affinity set in cpr_suspend_cpus()
834      */
835     affinity_clear();
836
837     /*
838      * offline all the cpus that were brought online during suspend
839      */
840     cpr_restore_offline();
841
842     mutex_exit(&cpu_lock);
843
844     /*
845      * Bring the system back up from a checkpoint, at this point
846      * the VM has been minimally restored by boot, the following
847      * are executed sequentially:
848      */

```

```

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 pwron_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     rb_dump:
955     /*
956      * IMPORTANT: SENSITIVE RESUME SEQUENCE
957      *
958      * DO NOT ADD ANY INITIALIZATION STEP BEFORE THIS POINT!!
959      */
960     rb_nodump:
961     /*
962      * If we did suspend to RAM, we didn't generate a dump
963      */
964     PMD(PMD_SX, ("cpr_resume: CPR DMA callback\n"))
965     (void) callb_execute_class(CB_CL_CPR_DMA, CB_CODE_CPR_RESUME);
966     if (cpr_suspend_succeeded) {
967         PMD(PMD_SX, ("cpr_resume: CPR RPC callback\n"))
968         (void) callb_execute_class(CB_CL_CPR_RPC, CB_CODE_CPR_RESUME);

```

```

970         }
971
972         prom_resume_prepost();
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 #endif
1001
1002     /*
1003      * let the tmp callout catch up.
1004      */
1005     PMD(PMD_SX, ("cpr_resume: CPR_CALLBACK callback\n"))
1006     (void) callb_execute_class(CB_CL_CPR_CALLBACK, CB_CODE_CPR_RESUME);
1007
1008     i_cpr_enable_intr();
1009
1010     mutex_enter(&cpu_lock);
1011     PMD(PMD_SX, ("cpr_resume: cyclic resume\n"))
1012     cyclic_resume();
1013     mutex_exit(&cpu_lock);
1014
1015     PMD(PMD_SX, ("cpr_resume: handle xc\n"))
1016     i_cpr_handle_xc(0); /* turn it off to allow xc assertion */
1017
1018     PMD(PMD_SX, ("cpr_resume: CPR_POST_KERNEL callback\n"))
1019     (void) callb_execute_class(CB_CL_CPR_POST_KERNEL, CB_CODE_CPR_RESUME);
1020
1021     /*
1022      * statistics gathering
1023      */
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(&pwrn_tv);
1033
1034         ctp = &cpr_term.tm_shutdown;
1035         if (sleepytype == 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", &pwrn_tv);
1040
1041             str = " prom time";
1042             CPR_STAT_EVENT_START_TMZ(str, &pwrn_tv);
1043             ctp = &cpr_term.tm_cprboot_start;
1044             CPR_STAT_EVENT_END_TMZ(str, ctp);
1045
1046             str = " read statefile";
1047             CPR_STAT_EVENT_START_TMZ(str, ctp);
1048             ctp = &cpr_term.tm_cprboot_end;
1049             CPR_STAT_EVENT_END_TMZ(str, ctp);
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             mutex_exit(&cpu_lock);
1085         }
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
1102             (sleeptype == CPR_TORAM &&
1103             CPR->c_substate == C_ST_NODUMP) {
1104                 if (cpr_test_point == FORCE_SUSPEND_TO_RAM) {
1105                     PMD(PMD_SX, ("cpr_resume: resume device "
1106                         "warn\n"))
1107                     cpr_err(CE_WARN, str);
1108                 } else {
1109                     PMD(PMD_SX, ("cpr_resume: resume device "
1110                         "panic\n"))
1111                     cpr_err(CE_PANIC, str);
1112             } else {
1113                 PMD(PMD_SX, ("cpr_resume: resume device warn\n"))
1114                 cpr_err(CE_WARN, str);
1115             }
1116         }
1117
1118     CPR_STAT_EVENT_END(" start drivers");
1119     CPR_DEBUG(CPR_DEBUG1, "done\n");
1120
1121 #if defined(__x86)
1122     /*
1123      * If cpr_resume_uniprocs is set, then unpause all the processors
1124      * that were paused before resuming the drivers
1125      */
1126     if (cpr_resume_uniprocs) {
1127         mutex_enter(&cpu_lock);
1128         start_cpus();
1129         mutex_exit(&cpu_lock);
1130     }
1131 #endif
1132
1133 /*
1134  * If we had disabled modunloading in this cpr resume cycle (i.e. we
1135  * resumed from a state earlier than C_ST_SUSPEND_DEVICES), re-enable
1136  * modunloading now.
1137 */
1138 if (CPR->c_substate != C_ST_SUSPEND_DEVICES) {
1139     PMD(PMD_SX, ("cpr_resume: modload enable\n"))
1140     modunload_enable();
1141 }
1142
1143 /*
1144  * Hooks needed by lock manager prior to resuming.
1145  * Refer to code for more comments.
1146 */
1147 PMD(PMD_SX, ("cpr_resume: lock mgr\n"))
1148 cpr_lock_mgr(lm_cprresume);
1149
1150 #if defined(__sparc)
1151 /*
1152  * This is a partial (half) resume during cpr suspend, we
1153  * haven't yet given up on the suspend. On return from here,
1154  * cpr_suspend() will try to reallocate and retry the suspend.
1155 */
1156 if (CPR->c_substate == C_ST_DUMP_NOSPC) {
1157     return (0);
1158 }
1159
1160 if (sleeptype == CPR_TODISK) {
1161     rb_statef_alloc();
1162     cpr_statef_close();
1163
1164 rb_disable_ufs_logging:
1165     /*
1166      * if ufs logging was disabled, re-enable
1167

```

```

1167
1168             */
1169         }
1170     #endif
1171
1172 rb_pm_reattach_noinvol:
1173     /*
1174      * When pm_reattach_noinvol() succeeds, modunload_thread will
1175      * remain disabled until after cpr suspend passes the
1176      * C_ST_STOP_KERNEL_THREADS state. If any failure happens before
1177      * cpr suspend reaches this state, we'll need to enable modunload
1178      * thread during rollback.
1179 */
1180 if (CPR->c_substate == C_ST_DISABLE_UFS_LOGGING ||
1181     CPR->c_substate == C_ST_STATEF_ALLOC ||
1182     CPR->c_substate == C_ST_SUSPEND_DEVICES ||
1183     CPR->c_substate == C_ST_STOP_KERNEL_THREADS) {
1184     PMD(PMD_SX, ("cpr_resume: reattach noinvol fini\n"))
1185     pm_reattach_noinvol_fini();
1186 }
1187
1188 PMD(PMD_SX, ("cpr_resume: CPR POST USER callback\n"))
1189 (void) callb_execute_class(CB_CL_CPR_POST_USER, CB_CODE_CPR_RESUME);
1190 PMD(PMD_SX, ("cpr_resume: CPR PROMPRINTF callback\n"))
1191 (void) callb_execute_class(CB_CL_CPR_PROMPRINTF, CB_CODE_CPR_RESUME);
1192
1193 PMD(PMD_SX, ("cpr_resume: restore direct levels\n"))
1194 pm_restore_direct_levels();
1195
1196 rb_stop_user_threads:
1197     CPR_DEBUG(CPR_DEBUG1, "starting user threads...");
1198     PMD(PMD_SX, ("cpr_resume: starting user threads\n"))
1199     cpr_start_user_threads();
1200     CPR_DEBUG(CPR_DEBUG1, "done\n");
1201     /*
1202      * Ask Xorg to resume the frame buffer, and wait for it to happen
1203      */
1204     mutex_enter(&srn_clone_lock);
1205     if (srn_signal) {
1206         PMD(PMD_SX, ("cpr_suspend: (*srn_signal)(..., "
1207                         "SRN_NORMAL_RESUME)\n"))
1208         srn_inuse = 1; /* because (*srn_signal) cv_waits */
1209         (*srn_signal)(SRN_TYPE_APM, SRN_NORMAL_RESUME);
1210         srn_inuse = 0;
1211     } else {
1212         PMD(PMD_SX, ("cpr_suspend: srn_signal NULL\n"))
1213     }
1214     mutex_exit(&srn_clone_lock);
1215
1216 #if defined(__sparc)
1217 rb_mp_offline:
1218     if (cpr_mp_online())
1219         cpr_err(CE_WARN, "Failed to online all the processors.");
1220 #endif
1221
1222 rb_others:
1223     PMD(PMD_SX, ("cpr_resume: dep thread\n"))
1224     pm_dispatch_to_dep_thread(PM_DEP_WK_CPR_RESUME, NULL, NULL,
1225                               PM_DEP_WAIT, NULL, 0);
1226
1227 PMD(PMD_SX, ("cpr_resume: CPR PM callback\n"))
1228 (void) callb_execute_class(CB_CL_CPR_PM, CB_CODE_CPR_RESUME);
1229
1230 if (cpr_suspend_succeeded) {
1231     cpr_stat_record_events();
1232 }


```

```
1234 #if defined(__sparc)
1235     if (sleeptype == CPR_TODISK && !cpr_reusable_mode)
1236         cpr_clear_definfo();
1237 #endif

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");

1245     CPR_STAT_EVENT_END("Resume Total");

1247     CPR_STAT_EVENT_START_TMZ("WHOLE CYCLE", &wholecycle_tv);
1248     CPR_STAT_EVENT_END("WHOLE CYCLE");

1250     if (cpr_debug & CPR_DEBUG1)
1251         cmn_err(CE_CONT, "\nThe system is back where you left!\n");

1253     CPR_STAT_EVENT_START("POST CPR DELAY");

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);

1260     CPR_STAT_EVENT_PRINT();
1261 #endif /* CPR_STAT */
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

52333 Fri Jan 3 22:11:52 2014
new/usr/src/uts/common/disp/cmt.c
patch cpu-pause-func-deglobalize

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;  
342  
343     ASSERT(MUTEX_HELD(&cpu_lock));  
344  
345     parent = pg->cmt_parent;  
346     if (parent == NULL) {  
347         /*  
348         * Nothing to do  
349         */  
350         return;  
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     }  
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     }  
398     children = pg->cmt_children;  
399     pg->cmt_children = parent->cmt_children;  
400     parent->cmt_children = children;  
402     nchildren = pg->cmt_nchildren;  
403     pg->cmt_nchildren = parent->cmt_nchildren;  
404     parent->cmt_nchildren = nchildren;  
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;  
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 */  
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;  
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         }  
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_itr_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
1538              group_expand(pg->cmt_siblings,
1539                          GROUP_SIZE(pg->cmt_siblings) + cap_needed);
1540
1541          /*
1542           * 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         /*
1592          * Because we are iterating over the same group that we are
1593          * operating on (removing the children), first add all of PG's
1594          * children to the parent's children set, and once we are done
1595          * iterating, empty PG's children set.
1596          */
1597         if (pg->cmt_children != NULL) {
1598             children = pg->cmt_children;
1599
1600             group_iter_init(&child_iter);
1601             while ((child = group_iterate(children, &child_iter)) != 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     group_empty(pg->cmt_children);
1615 }
1616
1617 /*
1618  * Reset the callbacks to the defaults
1619  */
1620 pg_callback_set_defaults((pg_t *)pg);
1621
1622 /*
1623  * Update all the CPU lineages in each of PG's CPUs
1624  */
1625 PG_CPU_ITR_INIT(pg, cpu_iter);
1626 while ((cpu = pg_cpu_next(&cpu_iter)) != NULL) {
1627     pg_cmt_t           *cpu_pg;
1628     group_iter_t        liter; /* Iterator for the lineage */
1629     cpu_pg_t            *cpd;   /* CPU's PG data */
1630
1631     /*
1632      * The CPU's lineage is under construction still
1633      * references the bootstrap CPU PG data structure.
1634      */
1635     if (pg_cpu_is_bootstrapped(cpu))
1636         cpd = pgdata;
1637     else
1638         cpd = cpu->cpu_pg;
1639
1640     /*
1641      * Iterate over the CPU's PGs updating the children
1642      * of the PG being promoted, since they have a new
1643      * parent and siblings set.
1644      */
1645     group_iter_init(&liter);
1646     while ((cpu_pg = group_iterate(&cpd->pgs,
1647                                   &liter)) != NULL) {
1648         if (cpu_pg->cmt_parent == pg) {
1649             cpu_pg->cmt_parent = pg->cmt_parent;
1650             cpu_pg->cmt_siblings = pg->cmt_siblings;
1651         }
1652     }
1653
1654     /*
1655      * Update the CPU's lineages
1656      */
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
1673 {

```

```
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 Fri Jan 3 22:11:52 2014
new/usr/src/uts/common/disp/cpupart.c

patch cpu-pause-func-deglobalize

_____ 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 lgrpid;
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             /*
```

1

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

```
377                                         * There are still threads bound to the partition
378                                         */
379                                         cpu_state_change_notify(cp->cpu_id, CPU_CPUPART_IN);
380                                         return (EBUSY);
381 }
382 }

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 }

393 /*
394  * Stop further threads weak binding to this cpu.
395  */
396 cpu_inmotion = cp;
397 membrar_enter();

399 /*
400  * Notify the Processor Groups subsystem that the CPU
401  * will be moving cpu partitions. This is done before
402  * CPUs are paused to provide an opportunity for any
403  * needed memory allocations.
404  */
405 pg_cpupart_out(cp, oldpp);
406 pg_cpupart_in(cp, newpp);

408 again:
409     if (move_threads) {
410         int loop_count;
411         /*
412          * Check for threads strong or weak bound to this CPU.
413          */
414         for (loop_count = 0; disp_bound_threads(cp, 0); loop_count++) {
415             if (loop_count >= 5) {
416                 cpu_state_change_notify(cp->cpu_id,
417                                         CPU_CPUPART_IN);
418                 pg_cpupart_out(cp, newpp);
419                 pg_cpupart_in(cp, oldpp);
420                 cpu_inmotion = NULL;
421                 return (EBUSY); /* some threads still bound */
422             }
423             delay(1);
424         }
425     }

427 /*
428  * Before we actually start changing data structures, notify
429  * the cyclic subsystem that we want to move this CPU out of its
430  * partition.
431  */
432 if (!cyclic_move_out(cp)) {
433     /*
434      * This CPU must be the last CPU in a processor set with
435      * a bound cyclic.
436      */
437     cpu_state_change_notify(cp->cpu_id, CPU_CPUPART_IN);
438     pg_cpupart_out(cp, newpp);
439     pg_cpupart_in(cp, oldpp);
440     cpu_inmotion = NULL;
441     return (EBUSY);
442 }
```

2

```

444 pause_cpus(cp, NULL);
445 pause_cpus(cp);
446
447 if (move_threads) {
448     /*
449      * The thread on cpu before the pause thread may have read
450      * cpu_inmotion before we raised the barrier above. Check
451      * again.
452     */
453     if (disp_bound_threads(cp, 1)) {
454         start_cpus();
455         goto again;
456     }
457 }
458
459 /*
460  * Now that CPUs are paused, let the PG subsystem perform
461  * any necessary data structure updates.
462  */
463 pg_cpupart_move(cp, oldpp, newpp);
464
465 /* save this cpu's lgroup -- it'll be the same in the new partition */
466 lgrp_id = cp->cpu_lpl->lpl_lgrp_id;
467
468 cpu_lpl = cp->cpu_lpl;
469 /*
470  * let the lgroup framework know cp has left the partition
471  */
472 lgrp_config(LGRP_CONFIG_CPU_PART_DEL, (uintptr_t)cp, lgrp_id);
473
474 /* move out of old partition */
475 oldpp->cp_ncpus--;
476 if (oldpp->cp_ncpus > 0) {
477
478     ncp = cp->cpu_prev_part->cpu_next_part = cp->cpu_next_part;
479     cp->cpu_next_part->cpu_prev_part = cp->cpu_prev_part;
480     if (oldpp->cp_cpulist == cp) {
481         oldpp->cp_cpulist = ncp;
482     }
483 } else {
484     ncp = oldpp->cp_cpulist = NULL;
485     cp_num_parts_nonempty--;
486     ASSERT(cp_num_parts_nonempty != 0);
487 }
488 oldpp->cp_gen++;
489
490 /* move into new partition */
491 newlist = newpp->cp_cpulist;
492 if (newlist == NULL) {
493     newpp->cp_cpulist = cp->cpu_next_part = cp->cpu_prev_part = cp;
494     cp_num_parts_nonempty++;
495     ASSERT(cp_num_parts_nonempty != 0);
496 } else {
497     cp->cpu_next_part = newlist;
498     cp->cpu_prev_part = newlist->cpu_prev_part;
499     newlist->cpu_prev_part->cpu_next_part = cp;
500     newlist->cpu_prev_part = cp;
501 }
502 cp->cpu_part = newpp;
503 newpp->cp_ncpus++;
504 newpp->cp_gen++;
505
506 ASSERT(bitset_is_null(&newpp->cp_haltset));
507 ASSERT(bitset_is_null(&oldpp->cp_haltset));

```

```

509
510     /*
511      * let the lgroup framework know cp has entered the partition
512      */
513      lgrp_config(LGRP_CONFIG_CPUPART_ADD, (uintptr_t)cp, lgrpid);
514
515     /*
516      * If necessary, move threads off processor.
517      */
518      if (move_threads) {
519          ASSERT(ncp != NULL);
520
521          /*
522           * Walk thru the active process list to look for
523           * threads that need to have a new home lgroup,
524           * or the last CPU they run on is the same CPU
525           * being moved out of the partition.
526          */
527
528          for (p = p_ractive; p != NULL; p = p->p_next) {
529
530              t = p->p_tlist;
531
532              if (t == NULL)
533                  continue;
534
535              lgrp_diff_lpl = 0;
536
537              do {
538
539                  ASSERT(t->t_lpl != NULL);
540
541                  /*
542                   * Update the count of how many threads are
543                   * in this CPU's lgroup but have a different lpl
544                  */
545
546                  if (t->t_lpl != cpu_lpl &&
547                      t->t_lpl->lpl_lgrpid == lgrpid)
548                      lgrp_diff_lpl++;
549
550                  /*
551                   * If the lgroup that t is assigned to no
552                   * longer has any CPUs in t's partition,
553                   * we'll have to choose a new lgroup for t.
554                  */
555
556                  if (!LGRP_CPUS_IN_PART(t->t_lpl->lpl_lgrpid,
557                                         t->t_cpupart)) {
558                      lgrp_move_thread(t,
559                                      lgrp_choose(t, t->t_cpupart), 0);
560                  }
561
562                  /*
563                   * make sure lpl points to our own partition
564                   */
565                  ASSERT(t->t_lpl >= t->t_cpupart->cp_lgrploads &&
566                         (t->t_lpl < t->t_cpupart->cp_lgrploads +
567                          t->t_cpupart->cp_nlgrploads));
568
569                  ASSERT(t->t_lpl->lpl_ncpu > 0);
570
571                  /*
572                   * Update CPU last ran on if it was this CPU */
573                  if (t->t_cpu == cp && t->t_cpupart == oldpp &&
574                      t->t_bound_cpu != cp) {
575                      t->t_cpu = disp_lowpri_cpu(ncp,
576                                              t->t_lpl, t->t_pri, NULL);
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025
1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079
1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133
1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187
1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241
1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295
1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349
1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403
1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457
1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511
1512
1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528
1529
1530
1531
1532
1533
1534
1535
1536
1537
1538
1539
1540
1541
1542
1543
1544
1545
1546
1547
1548
1549
1550
1551
1552
1553
1554
1555
1556
1557
1558
1559
1560
1561
1562
1563
1564
1565
1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619
1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673
1674
1675
1676
1677
1678
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727
1728
1729
1730
1731
1732
1733
1734
1735
1736
1737
1738
1739
1740
1741
1742
1743
1744
1745
1746
1747
1748
1749
1750
1751
1752
1753
1754
1755
1756
1757
1758
1759
1760
1761
1762
1763
1764
1765
1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1779
1780
1781
1782
1783
1784
1785
1786
1787
1788
1789
1789
1790
1791
1792
1793
1794
1795
1796
1797
1798
1799
1799
1800
1801
1802
1803
1804
1805
1806
1807
1808
1809
1809
1810
1811
1812
1813
1814
1815
1816
1817
1818
1819
1819
1820
1821
1822
1823
1824
1825
1826
1827
1828
1829
1829
1830
1831
1832
1833
1834
1835
1836
1837
1838
1839
1839
1840
1841
1842
1843
1844
1845
1846
1847
1848
1849
1849
1850
1851
1852
1853
1854
1855
1856
1857
1858
1859
1859
1860
1861
1862
1863
1864
1865
1866
1867
1868
1869
1869
1870
1871
1872
1873
1874
1875
1876
1877
1878
1879
1879
1880
1881
1882
1883
1884
1885
1886
1887
1888
1889
1889
1890
1891
1892
1893
1894
1895
1896
1897
1898
1899
1899
1900
1901
1902
1903
1904
1905
1906
1907
1908
1909
1909
1910
1911
1912
1913
1914
1915
1916
1917
1918
1919
1919
1920
1921
1922
1923
1924
1925
1926
1927
1928
1929
1929
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1939
1940
1941
1942
1943
1944
1945
1946
1947
1948
1949
1949
1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1989
1990
1991
1992
1993
1994
1995
1996
1997
1998
1999
1999
2000
2001
2002
2003
2004
2005
2006
2007
2008
2009
2009
2010
2011
2012
2013
2014
2015
2016
2017
2018
2019
2019
2020
2021
2022
2023
2024
2025
2026
2027
2028
2029
2029
2030
2031
2032
2033
2034
2035
2036
2037
2038
2039
2039
2040
2041
2042
2043
2044
2045
2046
2047
2048
2049
2049
2050
2051
2052
2053
2054
2055
2056
2057
2058
2059
2059
2060
2061
2062
2063
2064
2065
2066
2067
2068
2069
2069
2070
2071
2072
2073
2074
2075
2076
2077
2078
2079
2079
2080
2081
2082
2083
2084
2085
2086
2087
2088
2089
2089
2090
2091
2092
2093
2094
2095
2096
2097
2098
2099
2099
2100
2101
2102
2103
2104
2105
2106
2107
2108
2109
2109
2110
2111
2112
2113
2114
2115
2116
2117
2118
2119
2119
2120
2121
2122
2123
2124
2125
2126
2127
2128
2129
2129
2130
2131
2132
2133
2134
2135
2136
2137
2138
2139
2139
2140
2141
2142
2143
2144
2145
2146
2147
2148
2149
2149
2150
2151
2152
2153
2154
2155
2156
2157
2158
2159
2159
2160
2161
2162
2163
2164
2165
2166
2167
2168
2169
2169
2170
2171
2172
2173
2174
2175
2176
2177
2178
2179
2179
2180
2181
2182
2183
2184
2185
2186
2187
2188
2189
2189
2190
2191
2192
2193
2194
2195
2196
2197
2198
2198
2199
2200
2201
2202
2203
2204
2205
2206
2207
2208
2209
2209
2210
2211
2212
2213
2214
2215
2216
2217
2218
2219
2219
2220
2221
2222
2223
2224
2225
2226
2227
2228
2229
2229
2230
2231
2232
2233
2234
2235
2236
2237
2238
2239
2239
2240
2241
2242
2243
2244
2245
2246
2247
2248
2249
2249
2250
2251
2252
2253
2254
2255
2256
2257
2258
2259
2259
2260
2261
2262
2263
2264
2265
2266
2267
2268
2269
2269
2270
2271
2272
2273
2274
2275
2276
2277
2278
2279
2279
2280
2281
2282
2283
2284
2285
2286
2287
2288
2289
2289
2290
2291
2292
2293
2294
2295
2296
2297
2298
2298
2299
2299
2300
2301
2302
2303
2304
2305
2306
2307
2308
2309
2309
2310
2311
2312
2313
2314
2315
2316
2317
2318
2319
2319
2320
2321
2322
2323
2324
2325
2326
2327
2328
2329
2329
2330
2331
2332
2333
2334
2335
2336
2337
2338
2339
2339
2340
2341
2342
2343
2344
2345
2346
2347
2348
2349
2349
2350
2351
2352
2353
2354
2355
2356
2357
2358
2359
2359
2360
2361
2362
2363
2364
2365
2366
2367
2368
2369
2369
2370
2371
2372
2373
2374
2375
2376
2377
2378
2379
2379
2380
2381
2382
2383
2384
2385
2386
2387
2388
2389
2389
2390
2391
2392
2393
2394
2395
2396
2397
2398
2398
2399
2399
2400
2401
2402
2403
2404
2405
2406
2407
2408
2409
2409
2410
2411
2412
2413
2414
2415
2416
2417
2418
2419
2419
2420
2421
2422
2423
2424
2425
2426
2427
2428
2429
2429
2430
2431
2432
2433
2434
2435
2436
2437
2438
2439
2439
2440
2441
2442
2443
2444
2445
2446
2447
2448
2449
2449
2450
2451
2452
2453
2454
2455
2456
2457
2458
2459
2459
2460
2461
2462
2463
2464
2465
2466
2467
2468
2469
2469
2470
2471
2472
2473
2474
2475
2476
2477
2478
2479
2479
2480
2481
2482
2483
2484
2485
2486
2487
2488
2489
2489
2490
2491
2492
2493
2494
2495
2496
2497
2497
2498
2499
2499
2500
2501
2502
2503
2504
2505
2506
2507
2508
2509
2509
2510
2511
2512
2513
2514
2515
2516
2517
2518
2519
2519
2520
2521
2522
2523
2524
2525
2526
2527
2528
2529
2529
2530
2531
2532
2533
2534
2535
2536
2537
2538
2539
2539
2540
2541
2542
2543
2544
2545
2546
2547
2548
2549
2549
2550
2551
2552
2553
2554
2555
2556
2557
2558
2559
2559
2560
2561
2562
2563
2564
2565
2566
2567
2568
2569
2569
2570
2571
2572
2573
2574
2575
2576
2577
2578
2579
2579
2580
2581
2582
2583
2584
2585
2586
2587
2588
2589
2589
2590
2591
2592
2593
2594
2595
2596
2597
2598
2598
2599
2599
2600
2601
2602
2603
2604
2605
2606
2607
2608
2609
2609
2610
2611
2612
2613
2614
2615
2616
2617
2618
2619
2619
2620
2621
2622
2623
2624
2625
2626
2627
2628
2629
2629
2630
2631
2632
2633
2634
2635
2636
2637
2638
2639
2639
2640
2641
2642
2643
2644
2645
2646
2647
2648
2649
2649
2650
2651
2652
2653
2654
2655
2656
2657
2658
2659
2659
2660
2661
2662
2663
2664
2665
2666
2667
2668
2669
2669
2670
2671
2672
2673
2674
2675
2676
2677
2678
2679
2679
2680
2681
2682
2683
2684
2685
2686
2687
2688
2689
2689
2690
2691
2692
2693
2694
2695
2696
2697
2698
2698
2699
2699
2700
2701
2702
2703
2704
2705
2706
2707
2708
2709
2709
2710
2711
2712
2713
2714
2715
2716
2717
2718
2719
2719
2720
2721
2722
2723
2724
2725
2726
2727
2728
2729
2729
2730
2731
2732
2733
2734
2735
2736
2737
2738
2739
2739
2740
2741
2742
2743
2744
2745
2746
2747
2748
2749
2749
2750
2751
2752
2753
2754
2755
2756
2757
2758
2759
2759
2760
2761
2762
2763
2764
2765
2766
2767
2768
2769
2769
2770
2771
2772
2773
2774
2775
2776
2777
2778
2779
2779
2780
2781
2782
2783
2784
2785
2786
2787
2788
2789
2789
2790
2791
2792
2793
2794
2795
2796
2797
2798
2798
2799
2799
2800
2801
2802
2803
2804
2805
2806
2807
2808
2809
2809
2810
2811
2812
2813
2814
2815
2816
2817
2818
2819
2819
2820
2821
2822
2823
2824
2825
2826
2827
2828
2829
2829
2830
2831
2832
2833
2834
2835
2836
2837
2838
2839
2839
2840
2841
2842
2843
2844
2845
2846
2847
2848
2849
2849
2850
2851
2852
2853
2854
2855
2856
2857
2858
2859
2859
2860
2861

```

```

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         klgpset_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;

595 do {
596     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_CPU_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 == old_pp &&
624        t->t_bound_cpu != cp) {
625         t->t_bound_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_runrun = 1;
642         cp->cpu_kprunrun = 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 }
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());
822
823     pp = kmem_zalloc(sizeof (cpupart_t), KM_SLEEP);
824     pp->cp_nlgrploads = lgrp_plat_max_lgrps();
825     pp->cp_lgrploads = kmem_zalloc(sizeof (lpl_t) * pp->cp_nlgrploads,
826                                     KM_SLEEP);
827
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     klgpset_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);
856
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);
862
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);
869     pause_cpus(NULL);
870     pp->cp_next = cp_list_head;
871     pp->cp_prev = cp_list_head->cp_prev;
872     cp_list_head->cp_prev->cp_next = pp;
873     cp_list_head->cp_prev = pp;
874     start_cpus();
875     mutex_exit(&cpu_lock);
876
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;
958
959     ASSERT(pool_lock_held());
960     mutex_enter(&cpu_lock);
961
962     pp = cpupart_find(psid);
963     if (pp == NULL || pp == &cp_default) {
964         mutex_exit(&cpu_lock);
965         return (EINVAL);
966     }
967
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     }
976
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     }
984
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);
994
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);
1008
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);
1015     pause_cpus(NULL);
1016     pp->cp_prev->cp_next = pp->cp_next;
1017     pp->cp_next->cp_prev = pp->cp_prev;
1018     if (cp_list_head == pp)
1019         cp_list_head = pp->cp_next;
1020     start_cpus();
1021
1022     if (cp_id_next > pp->cp_id)
1023         cp_id_next = pp->cp_id;
1024
1025     if (pp->cp_kstat)
1026         kstat_delete(pp->cp_kstat);
1027
1028     cp_numparts--;
1029
1030     disp_kp_free(&pp->cp_kp_queue);
1031
1032     cpupart_lpl_teardown(pp);
1033
1034     kmem_free(pp, sizeof (cpupart_t));
1035     mutex_exit(&cpu_lock);
1036
1037 }


---


unchanged_portion_omitted

```

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

68046 Fri Jan 3 22:11:52 2014
new/usr/src/uts/common/disp/disp.c
patch setfrontbackdq
patch cpu-pause-func-deglobalize

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;
323 ASSERT(MUTEX_HELD(&cpu_lock));
325 disp_mem = kmem_zalloc(NCPU *
326 sizeof (struct disp_queue_info), KM_SLEEP);
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;
343 pause_cpus(NULL, NULL);
344 pause_cpus(NULL);
345 for (i = 0; i < num; i++)
346 disp_dq_assign(&disp_mem[i], numpris);
start_cpus();
348 /*
349 * I must free all of the memory after starting the cpus because
350 * I can not risk sleeping in kmem_free while the cpus are stopped.
351 */
352 for (i = 0; i < num; i++)
353 disp_dq_free(&disp_mem[i]);
355 kmem_free(disp_mem, NCPU * sizeof (struct disp_queue_info));
356 }
unchanged_portion_omitted
1151 /*
1152 * setbackdq() keeps runqs balanced such that the difference in length
1153 * between the chosen rung and the next one is no more than RUNQ_MAX_DIFF.
1154 * For threads with priorities below RUNQ_MATCH_PRI levels, the runq's lengths
1155 * must match. When per-thread TS_RUNQMATCH flag is set, setbackdq() will
1156 * try to keep runqs perfectly balanced regardless of the thread priority.
1157 */
1158 #define RUNQ_MATCH_PRI 16 /* pri below which queue lengths must match */
1159 #define RUNQ_MAX_DIFF 2 /* maximum rung length difference */

1

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

1160 #define RUNQ_LEN(cp, pri) ((cp)->cpu_disp->disp_q[pri].dq_srunct)
1162 /*
1163 * Macro that evaluates to true if it is likely that the thread has cache
1164 * warmth. This is based on the amount of time that has elapsed since the
1165 * thread last ran. If that amount of time is less than "rechoose_interval"
1166 * ticks, then we decide that the thread has enough cache warmth to warrant
1167 * some affinity for t->t_cpu.
1168 */
1169 #define THREAD_HAS_CACHE_WARMTH(thread) \
1170 ((thread == curthread) || \
1171 ((ddi_get_lbolt() - thread->t_disp_time) <= rechoose_interval))
1173 #endif /* ! codereview */
1174 /*
1175 * Put the specified thread on the front/back of the dispatcher queue
1176 * corresponding to its current priority.
1177 * Put the specified thread on the back of the dispatcher
1178 * queue corresponding to its current priority.
1179 * Called with the thread in transition, onproc or stopped state and locked
1180 * (transition implies locked) and at high spl. Returns with the thread in
1181 * TS_RUN state and still locked.
1182 * Called with the thread in transition, onproc or stopped state
1183 * and locked (transition implies locked) and at high spl.
1184 * Returns with the thread in TS_RUN state and still locked.
1185 static void
1186 setfrontbackdq(kthread_t *tp, boolean_t front)
1187 void
1188 setbackdq(kthread_t *tp)
1189 {
1190 dispq_t *dq;
1191 disp_t *dp;
1192 cpu_t *cp;
1193 pri_t tpri;
1194 boolean_t bound;
1195 int bound;
1196 boolean_t self;
1197 ASSERT(THREAD_LOCK_HELD(tp));
1198 ASSERT((tp->t_schedflag & TS_ALLSTART) == 0);
1199 ASSERT(!thread_on_queue(tp)); /* make sure tp isn't on a rung */
1200 /*
1201 * If thread is "swapped" or on the swap queue don't
1202 * queue it, but wake sched.
1203 */
1204 if ((tp->t_schedflag & (TS_LOAD | TS_ON_SWAPQ)) != TS_LOAD) {
1205 disp_swapped_setrun(tp);
1206 return;
1207 }
1208 self = (tp == curthread);
1209 bound = (tp->t_bound_cpu || tp->t_weakbound_cpu);
1210 if (tp->t_bound_cpu || tp->t_weakbound_cpu)
1211 bound = 1;
1212 else
1213 bound = 0;
1214 tpri = DISP_PRIO(tp);
1215 if (ncpus == 1)
1216 cp = tp->t_cpu;
1217 else if (!bound) {
1218 if (tpri >= kpqpri) {

2

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

3

```

1213             setkpdq(tp, front ? SETKP_FRONT : SETKP_BACK);
1214             setkpdq(tp, SETKP_BACK);
1215             return;
1216         }
1217     cp = tp->t_cpu;
1218     if (!front) {
1219 #endif /* ! codereview */
1220     /*
1221      * We'll generally let this thread continue to run where
1222      * it last ran...but will consider migration if:
1223      * - We thread probably doesn't have much cache warmth.
1224      * - The CPU where it last ran is the target of an offline
1225      *   request.
1226      * - The thread last ran outside its home lgroup.
1227     */
1228     if ((!THREAD_HAS_CACHE_WARMTH(tp)) || (cp == cpu_inmotion))
1229     if ((tp->t_cpu == cpu_inmotion) {
1230         cp = disp_lowpri_cpu(tp->t_cpu, tp->t_lpl, tpri,
1231     } else if (!LGRP_CONTAINS_CPU(tp->t_lpl->lpl_lgrp, tp->t_lpl,
1232         cp = disp_lowpri_cpu(tp->t_cpu, tp->t_lpl, tpri,
1233         self ? tp->t_cpu : NULL);
1234     }
1235     } else {
1236         cp = tp->t_cpu;
1237     }
1238     if (tp->t_cpupart == cp->cpu_part) {
1239         if (front) {
1240             /*
1241              * We'll generally let this thread continue to run
1242              * where it last ran, but will consider migration
1243              * - The thread last ran outside its home lgroup
1244              * - The CPU where it last ran is the target of an
1245              *   offline request (a thread_nomigrate() on the
1246              *   motion CPU relies on this when forcing a priority
1247              *   - The thread isn't the highest priority thread
1248              *   it last ran, and it is considered not likely
1249              *   have significant cache warmth.
1250             */
1251         if (!LGRP_CONTAINS_CPU(tp->t_lpl->lpl_lgrp, cp)
1252             (cp == cpu_inmotion)) {
1253                 cp = disp_lowpri_cpu(tp->t_cpu, tp->t_lpl,
1254                 self ? cp : NULL);
1255             } else if ((tpri < cp->cpu_disp->disp_maxrunpri)
1256             (!THREAD_HAS_CACHE_WARMTH(tp))) {
1257                 cp = disp_lowpri_cpu(tp->t_cpu, tp->t_lpl,
1258                 NULL);
1259             }
1260         } else {
1261 #endif /* ! codereview */
1262             int qlen;
1263             /*
1264              * Perform any CMT load balancing
1265              */
1266             cp = cmt_balance(tp, cp);
1267             /*
1268              * Balance across the run queues
1269              */
1270             qlen = RUNQ_LEN(cp, tpri);
1271             if (tpri >= RUNQ_MATCH_PRI &&
1272

```

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

```

1274                                     !(tp->t_schedflag & TS_RUNQMATCH))
1275                                     qlen == RUNQ_MAX_DIFF;
1276                                     if (qlen > 0) {
1277                                         cpu_t *newcp;
1278
1279                                         if (tp->t_lpl->lpl_lgrp_id == LGRP_ROOTID)
1280                                             newcp = cp->cpu_next_part;
1281                                         } else if ((newcp = cp->cpu_next_lpl) ==
1282                                             newcp = cp->cpu_next_part;
1283                                         }
1284
1285                                         if (RUNQ_LEN(newcp, tpri) < qlen) {
1286                                             DTRACE_PROBE3(runq_balance,
1287                                               kthread_t *, tp,
1288                                               cpu_t *, cp, cpu_t *, newcp)
1289                                         cp = newcp;
1290                                         }
1291                                     }
1292 #endif /* ! codereview */
1293     } else {
1294         /*
1295          * Migrate to a cpu in the new partition.
1296          */
1297         cp = disp_lowpri_cpu(tp->t_cpupart->cp_cpulist,
1298                             tp->t_lpl, tp->t_pri, NULL);
1299     }
1300 }
1301
1302 #endif /* ! codereview */
1303     ASSERT((cp->cpu_flags & CPU QUIESCED) == 0);
1304 } else {
1305     /*
1306      * It is possible that t_weakbound_cpu != t_bound_cpu (for
1307      * a short time until weak binding that existed when the
1308      * strong binding was established has dropped) so we must
1309      * favour weak binding over strong.
1310      */
1311     cp = tp->t_weakbound_cpu ?
1312         tp->t_weakbound_cpu : tp->t_bound_cpu;
1313 }
1314
1315 #endif /* ! codereview */
1316 /*
1317  * A thread that is ONPROC may be temporarily placed on the run queue
1318  * but then chosen to run again by disp. If the thread we're placing on
1319  * the queue is in TS_ONPROC state, don't set its t_waitrq until a
1320  * replacement process is actually scheduled in swtch(). In this
1321  * situation, curthread is the only thread that could be in the ONPROC
1322  * state.
1323 */
1324 if ((!self) && (tp->t_waitrq == 0)) {
1325     hrtimer_t curtime;
1326
1327     curtime = gethrtime_unscaled();
1328     (void) cpu_update_pct(tp, curtime);
1329     tp->t_waitrq = curtime;
1330 } else {
1331     (void) cpu_update_pct(tp, gethrtime_unscaled());
1332 }
1333
1334 dp = cp->cpu_disp;
1335 disp_lock_enter_high(&dp->disp_lock);
1336
1337 DTRACE_SCHED3(enqueue, kthread_t *, tp, disp_t *, dp, int, front);
1338 if (front) {
1339     TRACE 2(TR FAC DISP, TR FRONTO, "frontrq:pri %d tid %p", tpri,

```

```

1340             tp);
1341     } else {
1329         DTRACE_SCHED3(enqueue, kthread_t *, tp, disp_t *, dp, int, 0);
1342         TRACE_3(TR_FAC_DISP, TR_BACKQ, "setbackdq:pri %d cpu %p tid %p",
1343                 tpri, cp, tp);
1344     }
1345 #endif /* ! codereview */

1347 #ifndef NPROBE
1348     /* Kernel probe */
1349     if (tnf_tracing_active)
1350         tnf_thread_queue(tp, cp, tpri);
1351 #endif /* NPROBE */

1353     ASSERT(tpri >= 0 && tpri < dp->disp_npri);

1355     THREAD_RUN(tp, &dp->disp_lock);           /* set t_state to TS_RUN */
1356     tp->t_disp_queue = dp;
1357     tp->t_link = NULL;

1359     dq = &dp->disp_q[tpri];
1360     dp->disp_nrunnable++;
1361     if (!bound)
1362         dp->disp_steam = 0;
1363     membar_enter();

1365     if (dq->dq_sruncont++ != 0) {
1366         if (front) {
1367             ASSERT(dq->dq_last != NULL);
1368             tp->t_link = dq->dq_first;
1369             dq->dq_first = tp;
1370         } else {
1371 #endif /* ! codereview */
1372             ASSERT(dq->dq_first != NULL);
1373             dq->dq_last->t_link = tp;
1374             dq->dq_last = tp;
1375         }
1376 #endif /* ! codereview */
1377     } else {
1378         ASSERT(dq->dq_first == NULL);
1379         ASSERT(dq->dq_last == NULL);
1380         dq->dq_first = dq->dq_last = tp;
1381         BT_SET(dp->disp_qactmap, tpri);
1382         if (tpri > dp->disp_maxrunpri) {
1383             dp->disp_maxrunpri = tpri;
1384             membar_enter();
1385             cpu_resched(cp, tpri);
1386         }
1387     }

1389     if (!bound && tpri > dp->disp_max_unbound_pri) {
1390         if (self && dp->disp_max_unbound_pri == -1 && cp == CPU) {
1391             /*
1392             * If there are no other unbound threads on the
1393             * run queue, don't allow other CPUs to steal
1394             * this thread while we are in the middle of a
1395             * context switch. We may just switch to it
1396             * again right away. CPU_DISP_DONTSTEAL is cleared
1397             * in swtch and swtch_to.
1398             */
1399             cp->cpu_disp_flags |= CPU_DISP_DONTSTEAL;
1400         }
1401         dp->disp_max_unbound_pri = tpri;
1402     }
1404 #endif /* ! codereview */

```

```

1405         (*disp_enq_thread)(cp, bound);
1406     }

1408 /*
1409  * Put the specified thread on the back of the dispatcher
1410  * queue corresponding to its current priority.
1411  *
1412  * Called with the thread in transition, onproc or stopped state
1413  * and locked (transition implies locked) and at high spl.
1414  * Returns with the thread in TS_RUN state and still locked.
1415  */
1416 void
1417 setbackdq(kthread_t *tp)
1418 {
1419     setfrontbackdq(tp, B_FALSE);
1420 }

1422 /*
1423 #endif /* ! codereview */
1424  * Put the specified thread on the front of the dispatcher
1425  * queue corresponding to its current priority.
1426  *
1427  * Called with the thread in transition, onproc or stopped state
1428  * and locked (transition implies locked) and at high spl.
1429  * Returns with the thread in TS_RUN state and still locked.
1430  */
1431 void
1432 setfrontdq(kthread_t *tp)
1433 {
1434     setfrontbackdq(tp, B_TRUE);
1435     disp_t          *dp;
1436     dispq_t         *dq;
1437     cpu_t          *cp;
1438     pri_t          *tpri;
1439     int             bound;

1440     ASSERT(THREAD_LOCK_HELD(tp));
1441     ASSERT((tp->t_schedflag & TS_ALLSTART) == 0);
1442     ASSERT(!thread_on_queue(tp)); /* make sure tp isn't on a rung */

1443     /*
1444      * If thread is "swapped" or on the swap queue don't
1445      * queue it, but wake sched.
1446      */
1447     if ((tp->t_schedflag & (TS_LOAD | TS_ON_SWAPQ)) != TS_LOAD) {
1448         disp_swapped_setrun(tp);
1449         return;
1450     }

1451     if (tp->t_bound_cpu || tp->t_weakbound_cpu)
1452         bound = 1;
1453     else
1454         bound = 0;

1455     tpri = DISP_PRIO(tp);
1456     if (ncpus == 1)
1457         cp = tp->t_cpu;
1458     else if (!bound) {
1459         if (tpri >= kpqpri) {
1460             setkpqdq(tp, SETKP_FRONT);
1461             return;
1462         }
1463         cp = tp->t_cpu;
1464         if (tp->t_cupart == cp->cpu_part) {
1465             /*
1466              * We'll generally let this thread continue to run
1467            */
1468         }
1469     }

```

```

1268             * where it last ran, but will consider migration if:
1269             * - The thread last ran outside its home lgroup.
1270             * - The CPU where it last ran is the target of an
1271               offline request (a thread_nomigrate() on the in
1272               motion CPU relies on this when forcing a preempt).
1273             * - The thread isn't the highest priority thread where
1274               it last ran, and it is considered not likely to
1275               have significant cache warmth.
1276         */
1277     if ((!LGRP_CONTAINS_CPU(tp->t_lpl->lpl_lgrp, cp)) ||
1278         (cp == cpu_imotion)) {
1279         cp = disp_lowpri_cpu(tp->t_cpu, tp->t_lpl, tpri,
1280                             (tp == curthread) ? cp : NULL);
1281     } else if ((tpri < cp->cpu_disp->disp_maxrunpri) &&
1282                 (!THREAD_HAS_CACHE_WARMTH(tp))) {
1283         cp = disp_lowpri_cpu(tp->t_cpu, tp->t_lpl, tpri,
1284                             NULL);
1285     }
1286   } else {
1287     /*
1288      * Migrate to a cpu in the new partition.
1289      */
1290     cp = disp_lowpri_cpu(tp->t_cpupart->cp_cpulist,
1291                          tp->t_lpl, tp->t_pri, NULL);
1292   }
1293   ASSERT((cp->cpu_flags & CPU QUIESCED) == 0);
1294 } else {
1295   /*
1296      * It is possible that t_weakbound_cpu != t_bound_cpu (for
1297      * a short time until weak binding that existed when the
1298      * strong binding was established has dropped) so we must
1299      * favour weak binding over strong.
1300   */
1301   cp = tp->t_weakbound_cpu ?
1302       tp->t_weakbound_cpu : tp->t_bound_cpu;
1303 }
1304
1305 /*
1306  * A thread that is ONPROC may be temporarily placed on the run queue
1307  * but then chosen to run again by disp. If the thread we're placing on
1308  * the queue is in TS_ONPROC state, don't set its t_waitrq until a
1309  * replacement process is actually scheduled in swtch(). In this
1310  * situation, curthread is the only thread that could be in the ONPROC
1311  * state.
1312 */
1313 if ((tp != curthread) && (tp->t_waitrq == 0)) {
1314   hrtimer_t curtime;
1315
1316   curtime = gethrtime_unscaled();
1317   (void) cpu_update_pct(tp, curtime);
1318   tp->t_waitrq = curtime;
1319 } else {
1320   (void) cpu_update_pct(tp, gethrtime_unscaled());
1321 }
1322
1323 dp = cp->cpu_disp;
1324 disp_lock_enter_high(&dp->disp_lock);
1325
1326 TRACE_2(TR_FAC_DISP, TR_FRONTQ, "frontq:pri %d tid %p", tpri, tp);
1327 DTTRACE_SCHED3(enqueue, kthread_t *, tp, disp_t *, dp, int, 1);
1328
1329 #ifndef NPROBE
1330   /* Kernel probe */
1331   if (tnf_tracing_active)
1332     tnf_thread_queue(tp, cp, tpri);
1333 #endif /* NPROBE */

```

```

1335   ASSERT(tpri >= 0 && tpri < dp->disp_npri);
1336
1337   THREAD_RUN(tp, &dp->disp_lock); /* set TS_RUN state and lock */
1338   tp->t_disp_queue = dp;
1339
1340   dq = &dp->disp_q[tpri];
1341   dp->disp_nrunnable++;
1342   if (!bound)
1343     dp->disp_stal = 0;
1344   membar_enter();
1345
1346   if (dq->dq_srunct++ != 0) {
1347     ASSERT(dq->dq_last != NULL);
1348     tp->t_link = dq->dq_first;
1349     dq->dq_first = tp;
1350   } else {
1351     ASSERT(dq->dq_last == NULL);
1352     ASSERT(dq->dq_first == NULL);
1353     tp->t_link = NULL;
1354     dq->dq_first = dq->dq_last = tp;
1355     BT_SET(dp->disp_qactmap, tpri);
1356     if (tpri > dp->disp_maxrunpri) {
1357       dp->disp_maxrunpri = tpri;
1358       membar_enter();
1359       cpu_resched(cp, tpri);
1360     }
1361   }
1362
1363   if (!bound && tpri > dp->disp_max_unbound_pri) {
1364     if (tp == curthread && dp->disp_max_unbound_pri == -1 &&
1365         cp == CPU) {
1366       /*
1367        * If there are no other unbound threads on the
1368        * run queue, don't allow other CPUs to steal
1369        * this thread while we are in the middle of a
1370        * context switch. We may just switch to it
1371        * again right away. CPU_DISP_DONTSTEAL is cleared
1372        * in swtch and swtch_to.
1373       */
1374     }
1375     cp->cpu_disp_flags |= CPU_DISP_DONTSTEAL;
1376   }
1377   dp->disp_max_unbound_pri = tpri;
1378   (*disp_eng_thread)(cp, bound);
1379
1380   unchanged_portion_omitted

```

new/usr/src/uts/common/os/cpu.c

```
*****
94655 Fri Jan 3 22:11:53 2014
new/usr/src/uts/common/os/cpu.c
patch cpu-pause-func-deglobalize
*****
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
30 #include <sys/types.h>
31 #include <sys/param.h>
32 #include <sys/var.h>
33 #include <sys/thread.h>
34 #include <sys/cpuvar.h>
35 #include <sys/cpu_event.h>
36 #include <sys/kstat.h>
37 #include <sys/uadmin.h>
38 #include <sys/sysctl.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> /* to set per-cpu kmem_cache offset */
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)
```

1

new/usr/src/uts/common/os/cpu.c

```
62 #include <sys/x86_archext.h>
63 #endif
64 #include <sys/callo.h>
65
66 extern int      mp_cpu_start(cpu_t *);
67 extern int      mp_cpu_stop(cpu_t *);
68 extern int      mp_cpu_poweron(cpu_t *);
69 extern int      mp_cpu_poweroff(cpu_t *);
70 extern int      mp_cpu_configure(int);
71 extern int      mp_cpu_unconfigure(int);
72 extern void     mp_cpu_faulted_enter(cpu_t *);
73 extern void     mp_cpu_faulted_exit(cpu_t *);
74
75 extern int      cmp_cpu_to_chip(processorid_t cpuid);
76 #ifdef __sparcv9
77 extern char     *cpu_fru_fmr1(cpu_t *cp);
78#endif
79
80 static void     cpu_add_active_internal(cpu_t *cp);
81 static void     cpu_remove_active(cpu_t *cp);
82 static void     cpu_info_kstat_create(cpu_t *cp);
83 static void     cpu_info_kstat_destroy(cpu_t *cp);
84 static void     cpu_stats_kstat_create(cpu_t *cp);
85 static void     cpu_stats_kstat_destroy(cpu_t *cp);
86
87 static int      cpu_sys_stats_ks_update(kstat_t *ksp, int rw);
88 static int      cpu_vm_stats_ks_update(kstat_t *ksp, int rw);
89 static int      cpu_stat_ks_update(kstat_t *ksp, int rw);
90 static int      cpu_state_change_hooks(int, cpu_setup_t, cpu_setup_t);
91
92 /*
93 * cpu_lock protects ncpus, ncpus_online, cpu_flag, cpu_list, cpu_active,
94 * max_cpu_seqid_ever, and dispatch queue reallocations. The lock ordering with
95 * respect to related locks is:
96 *
97 *     cpu_lock --> thread_free_lock ---> p_lock ---> thread_lock()
98 *
99 * Warning: Certain sections of code do not use the cpu_lock when
100 * traversing the cpu_list (e.g. mutex_vector_enter(), clock()). Since
101 * all cpus are paused during modifications to this list, a solution
102 * to protect the list is too either disable kernel preemption while
103 * walking the list, *or* recheck the cpu_next pointer at each
104 * iteration in the loop. Note that in no cases can any cached
105 * copies of the cpu pointers be kept as they may become invalid.
106 */
107 kmutex_t        cpu_lock;
108 cpu_t          *cpu_list;           /* list of all CPUs */
109 cpu_t          *clock_cpu_list;    /* used by clock to walk CPUs */
110 cpu_t          *cpu_active;        /* list of active CPUs */
111 static cpuset_t cpu_available;    /* set of available CPUs */
112 cpuset_t        cpu_seqid_inuse;   /* which cpu_seqids are in use */
113
114 cpu_t          **cpu_seq;         /* ptrs to CPUs, indexed by seq_id */
115
116 /*
117 * max_ncpus keeps the max cpus the system can have. Initially
118 * it's NCPU, but since most archs scan the devtree for cpus
119 * fairly early on during boot, the real max can be known before
120 * ncpus is set (useful for early NCPU based allocations).
121 */
122 int max_ncpus = NCPU;
123 /*
124 * platforms that set max_ncpus to maximum number of cpus that can be
125 * dynamically added will set boot_max_ncpus to the number of cpus found
126 * at device tree scan time during boot.
127 */
```

2

```

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 outch;
214     kstat_named_t msg;
215     kstat_named_t sema;
216     kstat_named_t namei;
217     kstat_named_t ufsiget;
218     kstat_named_t ufsdirblk;
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];

```

[new/usr/src/uts/common/os/cpu.c](#)

5

```

769 long lindex = index;
770
771 ASSERT((curthread->t_bound_cpu != NULL) || (*safe == PAUSE_DIE));
772
773 while (*safe != PAUSE_DIE) {
774     *safe = PAUSE_READY;
775     membar_enter(); /* make sure stores are flushed */
776     sema_v(&cpi->cp_sem); /* signal requesting thread */
777
778     /*
779      * Wait here until all pause threads are running. That
780      * indicates that it's safe to do the spl. Until
781      * cpu_pause_info.cp_go is set, we don't want to spl
782      * because that might block clock interrupts needed
783      * to preempt threads on other CPUs.
784      */
785     while (cpi->cp_go == 0)
786         ;
787
788     /*
789      * Even though we are at the highest disp prio, we need
790      * to block out all interrupts below LOCK_LEVEL so that
791      * an intr doesn't come in, wake up a thread, and call
792      * setbackdq/setfrontdq.
793      */
794     s = splhigh();
795
796     /*
797      * if cp_func has been set then call it using index as the
798      * argument, currently only used by cpr_suspend_cpus().
799      * This function is used as the code to execute on the
800      * "paused" cpu's when a machine comes out of a sleep state
801      * and CPU's were powered off. (could also be used for
802      * hotplugging CPU's).
803      * if cpu_pause_func() has been set then call it using
804      * index as the argument, currently only used by
805      * cpr_suspend_cpus(). This function is used as the
806      * code to execute on the "paused" cpu's when a machine
807      * comes out of a sleep state and CPU's were powered off.
808      * (could also be used for hotplugging CPU's).
809      */
810     if (cpi->cp_func != NULL)
811         (*cpi->cp_func)((void *)lindex);
812     if (cpu_pause_func != NULL)
813         (*cpu_pause_func)((void *)lindex);
814
815     mach_cpu_pause(safe);
816
817     splx(s);
818     /*
819      * Waiting is at an end. Switch out of cpu_pause
820      * loop and resume useful work.
821      */
822     swtch();
823 }
824
825 mutex_enter(&pause_free_mutex);
826 *safe = PAUSE_DEAD;
827 cv_broadcast(&pause_free_cv);
828 mutex_exit(&pause_free_mutex);
829 }

```

unchanged_portion_omitted

```

974 /*
975  * Pause all of the CPUs except the one we are on by creating a high
976  * priority thread bound to those CPUs.
977 */

```

new/usr/src/uts/common/os/cpu.c

```

978 * Note that one must be extremely careful regarding code
979 * executed while CPUs are paused. Since a CPU may be paused
980 * while a thread scheduling on that CPU is holding an adaptive
981 * lock, code executed with CPUs paused must not acquire adaptive
982 * (or low-level spin) locks. Also, such code must not block,
983 * since the thread that is supposed to initiate the wakeup may
984 * never run.
985 *
986 * With a few exceptions, the restrictions on code executed with CPUs
987 * paused match those for code executed at high-level interrupt
988 * context.
989 */
990 void
991 pause_cpus(cpu_t *off_cp, void *(*func)(void *))
992 {
993     processorid_t    cpu_id;
994     int             i;
995     struct _cpu_pause_info *cpi = &cpu_pause_info;

997     ASSERT(MUTEX_HELD(&cpu_lock));
998     ASSERT(cpi->cp_paused == NULL);
999     cpi->cp_count = 0;
1000    cpi->cp_go = 0;
1001    for (i = 0; i < NCPU; i++)
1002        safe_list[i] = PAUSE_IDLE;
1003    kpreempt_disable();

1005    cpi->cp_func = func;

1007 #endif /* ! codereview */
1008 /*
1009  * If running on the cpu that is going offline, get off it.
1010  * This is so that it won't be necessary to rechoose a CPU
1011  * when done.
1012  */
1013 if (CPU == off_cp)
1014     cpu_id = off_cp->cpu_next_part->cpu_id;
1015 else
1016     cpu_id = CPU->cpu_id;
1017 affinity_set(cpu_id);

1019 /*
1020  * Start the pause threads and record how many were started
1021  */
1022 cpi->cp_count = cpu_pause_start(cpu_id);

1024 /*
1025  * Now wait for all CPUs to be running the pause thread.
1026  */
1027 while (cpi->cp_count > 0) {
1028     /*
1029      * Spin reading the count without grabbing the disp
1030      * lock to make sure we don't prevent the pause
1031      * threads from getting the lock.
1032      */
1033     while (sema_held(&cpi->cp_sem))
1034         ;
1035     if (sema_try(&cpi->cp_sem))
1036         --cpi->cp_count;
1037 }
1038 cpi->cp_go = 1;                                /* all have reached cpu_pause */

1040 /*
1041  * Now wait for all CPUs to spl. (Transition from PAUSE_READY
1042  * to PAUSE_WAIT.)

```

```

1043     */
1044     for (i = 0; i < NCPU; 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));
1069
1070     if (cpun >= NCPU || 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;
1084
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_cputart_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         return (0);
1276     }

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_OFFLINE 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      */
1443     ASSERT(curthread->t_cpu != cp);
1444     for (p = p->practive; p != NULL; p = p->p_next) {
1446         t = p->p_tlist;
1448         if (t == NULL)
1449             continue;
1451         lgrp_diff_lpl = 0;
1453         do {
1454             ASSERT(t->t_lpl != NULL);
1455             /*
1456              * Taking last CPU in lpl offline
1457              * Rehome thread if it is in this lpl
1458              * Otherwise, update the count of how many
1459              * threads are in this CPU's lgroup but have
1460              * a different lpl.
1461             */
1463             if (cpu_lpl->lpl_ncpu == 0) {
1464                 if (t->t_lpl == cpu_lpl)
1465                     lgrp_move_thread(t,
1466                                     lgrp_choose(t,
1467                                     t->t_cpupart), 0);
1468                 else if (t->t_lpl->lpl_lgrpid ==
1469                           cpu_lpl->lpl_lgrpid)
1470                     lgrp_diff_lpl++;
1471             }
1472             ASSERT(t->t_lpl->lpl_ncpu > 0);
1474             /*
1475              * Update CPU last ran on if it was this CPU
1476              */
1477             if (t->t_cpu == cp && t->t_bound_cpu != cp)
1478                 t->t_cpu = disp_lowpri_cpu(ncp,
1479                                             t->t_lpl, t->t_pri, NULL);
1480             ASSERT(t->t_cpu != cp || t->t_bound_cpu == cp ||
1481                   t->t_weakbound_cpu == cp);
1483             t = t->t_forw;
1484         } while (t != p->p_tlist);
1486         /*
1487          * Didn't find any threads in the same lgroup as this
1488          * CPU with a different lpl, so remove the lgroup from
1489          * the process lgroup bitmask.
1490          */
1492         if (lgrp_diff_lpl == 0)
1493             klgrpset_del(p->p_lgrpset, cpu_lpl->lpl_lgrpid);
1494     }
1496     /*
1497      * Walk thread list looking for threads that need to be
1498      * rehomed, since there are some threads that are not in
1499      * their process's p_tlist.
1500      */
1502     t = curthread;

```

```

1503         do {
1504             ASSERT(t != NULL && t->t_lpl != NULL);
1506             /*
1507              * Rehome threads with same lpl as this CPU when this
1508              * is the last CPU in the lpl.
1509              */
1511             if ((cpu_lpl->lpl_ncpu == 0) && (t->t_lpl == cpu_lpl))
1512                 lgrp_move_thread(t,
1513                                 lgrp_choose(t, t->t_cpupart), 1);
1515             ASSERT(t->t_lpl->lpl_ncpu > 0);
1517             /*
1518              * Update CPU last ran on if it was this CPU
1519              */
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;
1529         } while (t != curthread);
1530         ASSERT((cp->cpu_flags & (CPUFAULTED | CPU_SPARE)) == 0);
1531         cp->cpu_flags |= CPUOFFLINE;
1532         disp_cpu_inactive(cp);
1533         if (!no_quiesce)
1534             cp->cpu_flags |= CPUQUIESCED;
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     }
1544     out:
1545     cpu_inmotion = NULL;
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);
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);
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);

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     }

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;

1739     ASSERT(MUTEX_HELD(&cpu_lock));
1740     ASSERT(cpu_list != NULL); /* list started in cpu_list_init */
1742     lgrp_config(LGRP_CONFIG_CPU_ADD, (uintptr_t)cp, 0);

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'ed 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     for (seqid = 0; CPU_IN_SET(cpu_seqid_inuse, seqid); seqid++)
1772         continue;
1773     CPUSet_ADD(cpu_seqid_inuse, seqid);
1774     cp->cpu_seqid = seqid;

1776     if (seqid > max_cpu_seqid_ever)
1777         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;

1786     /*
1787      * allocate a pause thread for this CPU.
1788      */
1789     cpu_pause_alloc(cp);

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;

1803     init_cpu_mstate(cp, CMS_SYSTEM);

1805     pool_pset_mod = gethrtime();
1806 }

1808 /*
1809  * Do the opposite of cpu_add_unit().
1810 */
1811 void
1812 cpu_del_unit(int cpuid)
1813 {
1814     struct cpu    *cp, *cpnext;

1816     ASSERT(MUTEX_HELD(&cpu_lock));
1817     cp = cpu[cpuid];
1818     ASSERT(cp != NULL);

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);

1825     /*
1826      * Tear down the CPU's physical ID cache, and update any
1827      * processor groups
1828      */
1829     pg_cpu_fini(cp, NULL);
1830     pghw_physid_destroy(cp);

1832     /*
1833      * Destroy kstat stuff.
1834      */
1835     cpu_info_kstat_destroy(cp);
1836     term_cpu_mstate(cp);
1837     /*
1838      * Free up pause thread.
1839      */
1840     cpu_pause_free(cp);
1841     CPUSet_DEL(cpu_available, cp->cpu_id);
1842     cpu[cp->cpu_id] = NULL;
1843     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
1862     cpnext = cp->cpu_next;
1863     cp->cpu_prev->cpu_next = cp->cpu_next;
1864     cp->cpu_next->cpu_prev = cp->cpu_prev;
1865     if (cp == cpu_list)
1866         cpu_list = cpnext;
1867
1868     /*
1869      * Signals that the cpu has been deleted (see above).
1870      */
1871     cp->cpu_next = NULL;
1872     cp->cpu_prev = NULL;
1873
1874     start_cpus();
1875
1876     CPUSER_SET_DEL(cpu_seqid_inuse, cp->cpu_seqid);
1877     ncpus--;
1878     lgrp_config(LGRP_CONFIG_CPU_DEL, (uintptr_t)cp, 0);
1879
1880 } unchanged_portion_omitted
1881
1882 /*
1883  * Add a CPU to the list of active CPUs.
1884  * This is called from machine-dependent layers when a new CPU is started.
1885  */
1886 void
1887 cpu_add_active(cpu_t *cp)
1888 {
1889     pg_cpupart_in(cp, cp->cpu_part);
1890
1891     pause_cpus(NULL, NULL);
1892     pause_cpus(NULL);
1893     cpu_add_active_internal(cp);
1894     start_cpus();
1895
1896     cpu_stats_kstat_create(cp);
1897     cpu_create_intrstat(cp);
1898     lgrp_kstat_create(cp);
1899     cpu_state_change_notify(cp->cpu_id, CPU_INIT);
1900 } unchanged_portion_omitted

```

new/usr/src/uts/common/os/cpu_event.c

```
*****
30588 Fri Jan 3 22:11:53 2014
new/usr/src/uts/common/os/cpu_event.c
patch cpu-pause-func-deglobalize
*****
_____ 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;
```

1

new/usr/src/uts/common/os/cpu_event.c

```
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     static void
445     cpu_idle_remove_callback(cpu_idle_cb_impl_t *cip)
446     {
447         int i, found = 0;
448         int unlock = 0, unpause = 0;
449         cpu_idle_cb_state_t *sp;
450
451         ASSERT(MUTEX_HELD(&cpu_idle_cb_lock));
452
453         /* Try to acquire cpu_lock if not held yet. */
454         if (!MUTEX_HELD(&cpu_lock)) {
455             mutex_enter(&cpu_lock);
456             unlock = 1;
457         }
458
459         /*
460          * Pause all other CPUs.
461          * It's guaranteed that no other threads will access cpu_idle_cb_array
462          * after pause_cpus().
463          */
464         if (!cpus_paused()) {
465             pause_cpus(NULL, NULL);
466             pause_cpus(NULL);
467             unpause = 1;
468         }
469
470         /* Remove cip from array. */
471         for (i = 0; i < cpu_idle_cb_curr; i++) {
472             if (found == 0) {
473                 if (cpu_idle_cb_array[i].impl == cip) {
474                     found = 1;
475                 }
476             } else {
477                 cpu_idle_cb_array[i - 1] = cpu_idle_cb_array[i];
478             }
479         }
480         ASSERT(found != 0);
481         cpu_idle_cb_curr--;
482
483         /*
484          * Reset property ready flag for all CPUs if no registered callback
485          * left because cpu_idle_enter/exit will stop updating property if
486          * there's no callback registered.
487          */
488         if (cpu_idle_cb_curr == 0) {
489             for (sp = cpu_idle_cb_state, i = 0; i < max_ncpus; i++, sp++) {
490                 sp->v.ready = B_FALSE;
491             }
492         }
493     }
```

2

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

unchanged portion omitted

```
*****
21527 Fri Jan 3 22:11:53 2014
new/usr/src/uts/common/os/cpu_pm.c
patch cpu-pause-func-deglobalize
*****
_____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;
178
179     mutex_enter(&cpu_lock);
180     if (new_policy == cpupm_policy) {
181         mutex_exit(&cpu_lock);
182         return (result);
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         cpupm_policy = CPUPM_POLICY_DISABLED;
193         start_cpus();
194
195         result = cmt_pad_disable(PGHW_POW_ACTIVE);
196
197         /*
198          * Once PAD has been enabled, it should always be possible
199          * to disable it.
200          */
201         ASSERT(result == 0);
202
203         /*
204          * Bring all the active power domains to the maximum
205          * performance state.
206          */
207         cpupm_state_change_global(CPUPM_DTYPE_ACTIVE,
208             CPUPM_STATE_MAX_PERF);
209
210         break;
211     case CPUPM_POLICY_ELASTIC:
212
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         cpupm_policy = CPUPM_POLICY_ELASTIC;
233         start_cpus();
234
235         break;
236     default:
237         cmn_err(CE_WARN, "Attempt to set unknown CPUPM policy %d\n",
238                 new_policy);
239         ASSERT(0);
240         break;
241     }
242     mutex_exit(&cpu_lock);
243
244     return (result);
245 }
_____unchanged_portion_omitted_____

```

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

```

new/usr/src/uts/common/os/lgrp.c

```
*****  
119440 Fri Jan 3 22:11:53 2014  
new/usr/src/uts/common/os/lgrp.c  
patch cpu-pause-func-deglobalize  
*****  
_____ 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 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     klgrpset_t     changed;  
1243     int             count;  
1244     int             i;  
1245     lgrp_t          *my_lgrp;  
1246     lgrp_id_t       lgrp_id;  
1247     mnodeset_t      mnodeset_mask = ((mnodeset_t)1 << mnode);  
1248     boolean_t        drop_lock = B_FALSE;  
1249     boolean_t        need_synch = B_FALSE;  
  
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     }  
  
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;  
  
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_nmnodes. This restores the inconsistency  
1274     * in topology, temporarily introduced by lgrp_mem_fini().  
1275     */  
1276     if (! (is_copy_rename && (lgrp_root->lgrp_mnodes == mnodeset_mask)) &&  
1277         lgrp_root->lgrp_mnodes & mnodeset_mask) {  
1278         if (drop_lock)  
1279             mutex_exit(&cpu_lock);  
1280         return;  
1281     }  
  
1283     /*
```

1

new/usr/src/uts/common/os/lgrp.c

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

2

```

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, lgrpid);
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 * LOCKING for this routine is a bit tricky. Usually it is called without
1373 * cpu_lock and it must grab cpu_lock here to prevent racing with other
1374 * callers. During DR of the board containing the caged memory it may be called
1375 * with cpu_lock already held and CPUs paused.
1376 *
1377 * If the deletion is part of the DR copy-rename and the deleted mnode is the
1378 * only one present in the lgrp_root->lgrp_mnodes, all the topology is updated,
1379 * but lgrp_root->lgrp_nmnodes is left intact. Later, lgrp_mem_init() will insert
1380 * the same mnode back into the topology. See lgrp_mem_rename() and
1381 * lgrp_mem_init() for additional details.
1382 */
1383 void
1384 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    lgrpid;
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      * In the special case of DR copy-rename with the only mnode being
1453      * removed, the lgrp_mnodes for the root is always non-zero, but we
1454      * still need to update the lgroup topology.
1455      */
1456     if ((my_lgrp->lgrp_nmnodes > 0) &&
1457         !(is_copy_rename && (my_lgrp == lgrp_root) &&
1458           (my_lgrp->lgrp_mnodes == mnodes_mask))) {
1459         if (drop_lock)
1460             mutex_exit(&cpu_lock);
1461         return;
1462     }
1463
1464     /*
1465      * This lgroup does not contain any memory now
1466      */
1467     klgrpset_clear(my_lgrp->lgrp_set[LGRP_RSRC_MEM]);
1468
1469     /*
1470      * Remove this lgroup from lgroup topology if it does not contain any
1471      * resources now
1472      */
1473     lgrpid = my_lgrp->lgrp_id;
1474     count = 0;
1475     klgrpset_clear(changed);
1476     if (lgrp_rsets_empty(my_lgrp->lgrp_set)) {
1477
1478         /*
1479          * Delete lgroup when no more resources
1480          */
1481         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                 lgrpid))
1498             continue;
1499
1500         klgrpset_del(lgrp->lgrp_set[LGRP_RSRC_MEM], lgrpid);
1501     }
1502     if (drop_lock)
1503         mutex_exit(&cpu_lock);
1504 }
unchanged_portion_omitted_
```

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

```
1
```

```
*****  
36945 Fri Jan 3 22:11:54 2014  
new/usr/src/uts/common/os/lgrp_topo.c  
patch cpu-pause-func-deglobalize  
*****  
_____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;  
  
1466     count = 0;  
1467     if (*changed)  
1468         klgpset_clear(*changed);  
  
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     }  
  
1482     mutex_enter(&cpu_lock);  
1483     pause_cpus(NULL, NULL);  
1483     pause_cpus(NULL);  
  
1485     /*  
1486      * Look for any leaf lgroup without its latency set, finish adding it  
1487      * to the lgroup topology assuming that it exists and has the root  
1488      * lgroup as its parent, and update the memory nodes of all lgroups  
1489      * that have it as a memory resource.  
1490      */  
1491     for (i = 0; i < lgrp_count; i++) {  
1492         lgrp = lgrps[i];  
  
1494     /*  
1495      * Skip non-existent and non-leaf lgroups and any lgroup  
1496      * with its latency set already  
1497      */  
1498     if (lgrp == NULL || lgrp->lgrp_id == LGRP_NONE ||  
1499         lgrp->lgrp_childcnt != 0 || lgrp->lgrp_latency != 0)  
1500         continue;  
  
1502 #ifdef DEBUG  
1503     if (lgrp_topo_debug > 1)  
1504         prom_printf("\nlgrp_topo_update: updating lineage "
```

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

```
2
```

```
1505                                         "of lgrp %d at 0x%p\n", lgrp->lgrp_id,  
1506                                         (void *)lgrp);  
1507 }  
1508 #endif /* DEBUG */  
  
1510         count += lgrp_leaf_add(lgrp, lgrps, lgrp_count, &changes);  
1511         if (*changed)  
1512             klgpset_or(*changed, changes);  
1514         if (!klgpset_isempty(changes))  
1515             (void) lgrp_mnode_update(changes, NULL);  
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 }  
  
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     }  
1530     start_cpus();  
1531     mutex_exit(&cpu_lock);  
1533 }  
1534 }  
_____unchanged_portion_omitted_____
```

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

```
1
```

```
*****  
82758 Fri Jan 3 22:11:54 2014  
new/usr/src/uts/common/os/mem_config.c  
patch cpu-pause-func-deglobalize  
*****  
_____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;  
  
3304     mutex_enter(&cpu_lock);  
3305     pause_cpus(NULL, NULL);  
3305  
3307     cp = cpu_list;  
3308     do {  
3309         vc = cp->cpu_vm_data;  
3310         vc->vc_pnum_memseg = NULL;  
3311         vc->vc_pnext_memseg = NULL;  
  
3313     } while ((cp = cp->cpu_next) != cpu_list);  
  
3315     start_cpus();  
3316     mutex_exit(&cpu_lock);  
3317 }  
_____unchanged_portion_omitted_____
```

```
*****
30244 Fri Jan 3 22:11:54 2014
new/usr/src/uts/common/sys/cpuvar.h
patch cpu-pause-func-deglobalize
*****
_____ 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 \
666 void cpu_pause_init(void); \
667 cpu_t *cpu_get(processorid_t cpun); /* get the CPU struct associated */
```

```
668 int cpu_online(cpu_t *cp); /* take cpu online */ \
669 int cpu_offline(cpu_t *cp, int flags); /* take cpu offline */ \
670 int cpu_spare(cpu_t *cp, int flags); /* take cpu to spare */ \
671 int cpu_faulted(cpu_t *cp, int flags); /* take cpu to faulted */ \
672 int cpu_poweron(cpu_t *cp); /* take powered-off cpu to offline */ \
673 int cpu_poweroff(cpu_t *cp); /* take offline cpu to powered-off */ \
674 \
675 cpu_t *cpu_intr_next(cpu_t *cp); /* get next online CPU taking intrs */ \
676 int cpu_intr_count(cpu_t *cp); /* count # of CPUs handling intrs */ \
677 int cpu_intr_on(cpu_t *cp); /* CPU taking I/O interrupts? */ \
678 void cpu_intr_enable(cpu_t *cp); /* enable I/O interrupts */ \
679 void cpu_intr_disable(cpu_t *cp); /* disable I/O interrupts */ \
680 void cpu_intr_alloc(cpu_t *cp, int n); /* allocate interrupt threads */ \
681 \
682 /* \
683  * Routines for checking CPU states. \
684 */ \
685 /* \
686  * cpu_is_online(cpu_t *); /* check if CPU is online */ \
687  * cpu_is_nointr(cpu_t *); /* check if CPU can service intrs */ \
688  * cpu_is_active(cpu_t *); /* check if CPU can run threads */ \
689  * cpu_is_offline(cpu_t *); /* check if CPU is offline */ \
690  * 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, hrtime_t); \
727 extern void get_cpu_mstate(struct cpu *, hrtime_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 extern void /* */
```

```

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 Fri Jan 3 22:11:54 2014
new/usr/src/uts/i86pc/i86hvm/io/xpv/xpv_support.c
patch cpu-pause-func-deglobalize
*****
_____ 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;

442     cmn_err(CE_NOTE, "Domain suspending for save/migrate");
444     SUSPEND_DEBUG("xen_suspend_domain\n");

446     /*
447      * We only want to suspend the PV devices, since the emulated devices
448      * are suspended by saving the emulated device state. The PV devices
449      * are all children of the xpvd nexus device. So we search the
450      * device tree for the xpvd node to use as the root of the tree to
451      * be suspended.
452     */
453     if (xpvd_dip == NULL)
454         ddi_walk_devs(ddi_root_node(), check_xpvd, NULL);

456     /*
457      * suspend interrupts and devices
458     */
459     if (xpvd_dip != NULL)
460         (void) xen_suspend_devices(ddi_get_child(xpvd_dip));
461     else
462         cmn_err(CE_WARN, "No PV devices found to suspend");
463     SUSPEND_DEBUG("xenbus_suspend\n");
464     xenbus_suspend();

466     mutex_enter(&cpu_lock);

468     /*
469      * Suspend on vcpu 0
470     */
471     thread_affinity_set(curthread, 0);
472     kpreempt_disable();

474     if (ncpus > 1)
475         pause_cpus(NULL, NULL);
476         pause_cpus(NULL);
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      * From here on in, we can't take locks.
484     */

486     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);

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);

507     SUSPEND_DEBUG("gnttab_resume\n");
508     gnttab_resume();

510     SUSPEND_DEBUG("ec_resume\n");
511     ec_resume();

513     intr_restore(flags);

515     if (ncpus > 1)
516         start_cpus();

518     mutex_exit(&ec_lock);
519     mutex_exit(&cpu_lock);

521     /*
522      * Now we can take locks again.
523     */
525     rtcsync();

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);

533     thread_affinity_clear(curthread);
534     kpreempt_enable();

536     SUSPEND_DEBUG("finished xen_suspend_domain\n");

538 }
539 }_____ unchanged_portion_omitted_
```

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

```
*****
23453 Fri Jan 3 22:11:54 2014
new/usr/src/uts/i86pc/io/dr/dr_quiesce.c
patch cpu-pause-func-deglobalize
*****
_____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     force = dr_cmd_flags(handle) & SBD_FLAG_FORCE;
796
797     prom_printf("\nDR: suspending user threads...\n");
798     srh->sr_suspend_state = DR_SRSTATE_USER;
799     if (((rc = dr_stop_user_threads(srh)) != DDI_SUCCESS) &&
800         dr_check_user_stop_result) {
801         dr_resume(srh);
802         return (rc);
803     }
804
805     if (!force) {
806         struct dr_ref drc = {0};
807
808         prom_printf("\nDR: checking devices...\n");
809         dev_errs_idx = 0;
810
811         drc.arr = dev_errs;
812         drc.idx = &dev_errs_idx;
813         drc.len = DR_MAX_ERR_INT;
814
815         /*
816          * Since the root node can never go away, it
817          * doesn't have to be held.
818          */
819         ddi_walk_devs(ddi_root_node(), dr_check_unsafe_major, &drc);
820         if (dev_errs_idx) {
821             handle->h_err = drerr_int(ESBD_UNSAFE, dev_errs,
822                                       dev_errs_idx, 1);
823             dr_resume(srh);
824             return (DDI_FAILURE);
825         }
826         PR_QR("done\n");
827     } else {
828         prom_printf("\nDR: dr_suspend invoked with force flag\n");
829     }
830
831 #ifndef SKIP_SYNC
832     /*
833      * This sync swap out all user pages
834      */
835     vfs_sync(SYNC_ALL);
836 #endif
837
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         return (rc);
865     }
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

new/usr/src/uts/i86pc/io/pcplusmp/apic.c

34720 Fri Jan 3 22:11:55 2014

new/usr/src/uts/i86pc/io/pcplusmp/apic.c

patch apic-simplify

patch remove-apic_cr8pri

patch spacing-fix

patch apic-task-reg-write-dup

 unchanged_portion_omitted

```
137       /*  
138        * The ipl of an ISR at vector X is apic_vectortoipl[X>>4]  
139        * NOTE that this is vector as passed into intr_enter which is  
140        * programmed vector - 0x20 (APIC_BASE_VECT)  
141        */
```

```
143 uchar_t apic_ipltopri[MAXIPL + 1];     /* unix ipl to apic pri */  
144        /* The taskpri to be programmed into apic to mask given ipl */
```

```
146 #if defined(__amd64)  
147 uchar_t apic_cr8pri[MAXIPL + 1];     /* unix ipl to cr8 pri */  
148 #endif
```

```
146 /*  
147  * Correlation of the hardware vector to the IPL in use, initialized  
148  * from apic_vectortoipl[] in apic_init(). The final IPLs may not correlate  
149  * to the IPLs in apic_vectortoipl on some systems that share interrupt lines  
150  * connected to errata-stricken IOAPICs  
151  */  
152 uchar_t apic_ipls[APIC_AVAIL_VECTOR];
```

```
154 /*  
155  * Patchable global variables.  
156 */
```

```
157 int      apic_enable_hwsoftint = 0;     /* 0 - disable, 1 - enable    */  
158 int      apic_enable_bind_log = 1;     /* 1 - display interrupt binding log */
```

```
160 /*  
161  * Local static data  
162 */
```

```
163 static struct psm_ops apic_ops = {  
164     apic_probe,  
  
165     apic_init,  
166     apic_picinit,  
167     apic_intr_enter,  
168     apic_intr_exit,  
169     apic_setspl,  
170     apic_addspl,  
171     apic_delspl,  
172     apic_disable_intr,  
173     apic_enable_intr,  
174     (int (*)(int))NULL,                    /* psm_softlvl_to_irq */  
175     (void (*)(int))NULL,                    /* psm_set_softintr */  
  
176     apic_set_idlecpu,  
177     apic_unset_idlecpu,  
  
178     apic_clkinit,  
179     apic_getclkirq,
```

```
180     (void (*)(void))NULL,                  /* psm_hrtimeinit */  
181     apic_gethrtime,  
  
182     apic_get_next_processorid,  
183     apic_cpu_start,  
184     apic_post_cpu_start,  
185     apic_shutdown,
```

1

new/usr/src/uts/i86pc/io/pcplusmp/apic.c

```
190        apic_get_ipivect,  
191        apic_send_ipi,  
  
193        (int (*)(dev_info_t *, int))NULL,     /* psm_translate_irq */  
194        (void (*)(int, char *))NULL,            /* psm_notify_error */  
195        (void (*)(int))NULL,                    /* psm_notify_func */  
196        apic_timer_reprogram,  
197        apic_timer_enable,  
198        apic_timer_disable,  
199        apic_post_cyclic_setup,  
200        apic_preshutdown,  
201        apic_intr_ops,  
202        apic_state,  
203        apic_cpu_ops,  
204 };
```

 unchanged_portion_omitted

```
274 void  
275 apic_init(void)  
276 {  
277     int i;  
278     int j = 1;  
  
280     psm_get_ioapicid = apic_get_ioapicid;  
281     psm_get_localapicid = apic_get_localapicid;  
282     psm_xlate_vector_by_irq = apic_xlate_vector_by_irq;  
  
284     apic_ipltopri[0] = APIC_VECTOR_PER_IPL; /* leave 0 for idle */  
285     for (i = 0; i < (APIC_AVAIL_VECTOR / APIC_VECTOR_PER_IPL); i++) {  
286         if ((i < ((APIC_AVAIL_VECTOR / APIC_VECTOR_PER_IPL) - 1)) &&  
287             (apic_vectortoipl[i + 1] == apic_vectortoipl[i]))  
288             /* get to highest vector at the same ipl */  
289             continue;  
290         for (j <= apic_vectortoipl[i]; j++) {  
291             apic_ipltopri[j] = (i << APIC_IPL_SHIFT) +  
292                            APIC_BASE_VECT;  
293         }  
294     }  
295     for (j < MAXIPL + 1; j++)  
296         /* fill up any empty ipltopri slots */  
297         apic_ipltopri[j] = (i << APIC_IPL_SHIFT) + APIC_BASE_VECT;  
298     apic_init_common();  
299 #ifndef __amd64  
300 #if defined(__amd64)  
301     /*  
302         * Make cpu-specific interrupt info point to cr8pri vector  
303         */  
304     for (i = 0; i <= MAXIPL; i++)  
305         apic_cr8pri[i] = apic_ipltopri[i] >> APIC_IPL_SHIFT;  
306     CPU->cpu_pri_data = apic_cr8pri;  
307 #else  
308     if (cpuid_have_cr8access(CPU))  
309         apic_have_32bit_cr8 = 1;  
310 #endif /* !__amd64 */  
311 #endif /* __amd64 */  
312 }
```

 unchanged_portion_omitted

528 #endif /* DEBUG */

```
530 /*  
531  * platform_intr_enter  
532 */  
533  *     Called at the beginning of the interrupt service routine to  
534  *     mask all level equal to and below the interrupt priority  
535  *     of the interrupting vector. An EOI should be given to  
536  *     the interrupt controller to enable other HW interrupts.
```

2

```

537 *
538 *      Return -1 for spurious interrupts
539 *
540 */
541 /*ARGSUSED*/
542 static int
543 apic_intr_enter(int ipl, int *vectorp)
544 {
545     uchar_t vector;
546     int nipl;
547     int irq;
548     ulong_t iflag;
549     apic_cpus_info_t *cpu_infop;

551 /*
552  * The real vector delivered is (*vectorp + 0x20), but our caller
553  * subtracts 0x20 from the vector before passing it to us.
554  * (That's why APIC_BASE_VECT is 0x20.)
555  */
556     vector = (uchar_t)*vectorp;

558 /* if interrupted by the clock, increment apic_nsec_since_boot */
559     if (vector == apic_clkvect) {
560         if (!apic_oneshot) {
561             /* NOTE: this is not MT aware */
562             apic_hrtime_stamp++;
563             apic_nsec_since_boot += apic_nsec_per_intr;
564             apic_hrtime_stamp++;
565             last_count_read = apic_hertz_count;
566             apic_redistribute_compute();
567         }

569         /* We will avoid all the book keeping overhead for clock */
570         nipl = apic_ipls[vector];
572
573         *vectorp = apic_vector_to_irq[vector + APIC_BASE_VECT];

574         apic_reg_ops->apic_write_task_reg(apic_ipltopri[nipl]);
575         apic_reg_ops->apic_send_eoi(0);
576         if (apic_mode == LOCAL_APIC) {
577 #if defined(__amd64)
578             setcr8((ulong_t)(apic_ipltopri[nipl] >>
579                         APIC_IPL_SHIFT));
580 #else
581             if (apic_have_32bit_cr8)
582                 setcr8((ulong_t)(apic_ipltopri[nipl] >>
583                             APIC_IPL_SHIFT));
584             else
585                 LOCAL_APIC_WRITE_REG(APIC_TASK_REG,
586                                         (uint32_t)apic_ipltopri[nipl]);
587 #endif
588         } else {
589             LOCAL_APIC_WRITE_REG(APIC_EOI_REG, 0);
590         }
591     }
592
593     X2APIC_WRITE(APIC_TASK_REG, apic_ipltopri[nipl]);
594     X2APIC_WRITE(APIC_EOI_REG, 0);
595 }

596     return (nipl);
597 }

598     cpu_infop = &apic_cpus[psm_get_cpu_id()];

599     if (vector == (APIC_SPUR_INTR - APIC_BASE_VECT)) {
600         cpu_infop->aci_spur_cnt++;
601         return (APIC_INT_SPURIOUS);
602     }

```

```

587 /* Check if the vector we got is really what we need */
588 if (apic_revector_pending) {
589     /*
590      * Disable interrupts for the duration of
591      * the vector translation to prevent a self-race for
592      * the apic_revector_lock. This cannot be done
593      * in apic_xlate_vector because it is recursive and
594      * we want the vector translation to be atomic with
595      * respect to other (higher-priority) interrupts.
596     */
597     iflag = intr_clear();
598     vector = apic_xlate_vector(vector + APIC_BASE_VECT) -
599             APIC_BASE_VECT;
600     intr_restore(iflag);
601 }
602
603 nipl = apic_ipls[vector];
604 *vector = irq = apic_vector_to_irq(vector + APIC_BASE_VECT);
605
606 apic_reg_ops->apic_write_task_reg(apic_ipltopri[nipl]);
607 if (apic_mode == LOCAL_APIC) {
608 #if defined(__amd64)
609     setcr8((ulong_t)(apic_ipltopri[nipl] >> APIC_IPL_SHIFT));
610 #else
611     if (apic_have_32bit_cr8)
612         setcr8((ulong_t)(apic_ipltopri[nipl] >>
613                         APIC_IPL_SHIFT));
614     else
615         LOCAL_APIC_WRITE_REG(APIC_TASK_REG,
616                               (uint32_t)apic_ipltopri[nipl]);
617 #endif
618 } else {
619     X2APIC_WRITE(APIC_TASK_REG, apic_ipltopri[nipl]);
620 }
621
622 cpu_infop->aci_current[nipl] = (uchar_t)irq;
623 cpu_infop->aci_curipl = (uchar_t)nipl;
624 cpu_infop->aci_ISR_in_progress |= 1 << nipl;
625
626 /*
627  * apic_level_intr could have been assimilated into the irq struct.
628  * but, having it as a character array is more efficient in terms of
629  * cache usage. So, we leave it as is.
630  */
631 if (!apic_level_intr[irq]) {
632     apic_reg_ops->apic_send_eoi(0);
633     if (apic_mode == LOCAL_APIC) {
634         LOCAL_APIC_WRITE_REG(APIC_EOI_REG, 0);
635     } else {
636         X2APIC_WRITE(APIC_EOI_REG, 0);
637     }
638 }
639
640 #ifdef DEBUG
641     APIC_DEBUG_BUF_PUT(vector);
642     APIC_DEBUG_BUF_PUT(irq);
643     APIC_DEBUG_BUF_PUT(nipl);
644     APIC_DEBUG_BUF_PUT(psm_get_cpu_id());
645     if ((apic_stretch_interrupts) && (apic_stretch_ISR & (1 << nipl)))
646         drv_usecwait(apic_stretch_interrupts);
647
648     if (apic_break_on_cpu == psm_get_cpu_id())
649         apic_break();
650 #endif /* DEBUG */
651     return (nipl);
652 }
```

```
633 }
_____unchanged_portion_omitted_____  
  
649 /*
650  * Any changes made to this function must also change X2APIC
651  * version of intr_exit.
652 */
653 void
654 apic_intr_exit(int prev_ipl, int irq)
655 {
656     apic_cpus_info_t *cpu_infop;  
  
658     local_apic_write_task_reg(apic_ipktopri[prev_ipl]);
700 #if defined(__amd64)
701     setcr8((ulong_t)apic_cr8pri[prev_ipl]);
702 #else
703     if (apic_have_32bit_cr8)
704         setcr8((ulong_t)(apic_ipktopri[prev_ipl] >> APIC_IPL_SHIFT));
705     else
706         apicadr[APIC_TASK_REG] = apic_ipktopri[prev_ipl];
707 #endif  
  
660     APIC_INTR_EXIT();
661 }  
_____unchanged_portion_omitted_____  
  
685 /*
686  * Mask all interrupts below or equal to the given IPL.
687  * Any changes made to this function must also change X2APIC
688  * version of setspl.
689 */
690 static void
691 apic_setspl(int ipl)
692 {
693     local_apic_write_task_reg(apic_ipktopri[ipl]);
742 #if defined(__amd64)
743     setcr8((ulong_t)apic_cr8pri[ipl]);
744 #else
745     if (apic_have_32bit_cr8)
746         setcr8((ulong_t)(apic_ipktopri[ipl] >> APIC_IPL_SHIFT));
747     else
748         apicadr[APIC_TASK_REG] = apic_ipktopri[ipl];
749 #endif  
  
695     /* interrupts at ipl above this cannot be in progress */
696     apic_cpus[psm_get_cpu_id()].aci_ISR_in_progress &= (2 << ipl) - 1;
697     /*
698      * this is a patch fix for the ALR QSMP P5 machine, so that interrupts
699      * have enough time to come in before the priority is raised again
700      * during the idle() loop.
701      */
702     if (apic_setspl_delay)
703         (void) apic_reg_ops->apic_get_pri();
704 }  
_____unchanged_portion_omitted_____  

```

```
new/usr/src/uts/i86pc/io/pcplusmp/apic_regops.c
```

```
*****
9758 Fri Jan 3 22:11:55 2014
new/usr/src/uts/i86pc/io/pcplusmp/apic_regops.c
patch use-x2apic-feature
patch apic-simplify
patch apic-task-reg-write-dup
*****
_____unchanged_portion_omitted_____
86 int apic_have_32bit_cr8 = 0;
88 /* The default ops is local APIC (Memory Mapped IO) */
89 apic_reg_ops_t *apic_reg_ops = &local_apic_regs_ops;
91 /*
92  * APIC register ops related data structures and functions.
93 */
94 void apic_send_EOI();
95 void apic_send_directed_EOI(uint32_t irq);
97 #define X2APIC_CPUID_BIT      21
97 #define X2APIC_ENABLE_BIT     10
99 /*
100 * Local APIC Implementation
101 */
102 static uint64_t
103 local_apic_read(uint32_t reg)
104 {
105     return ((uint32_t)apicadr[reg]);
106 }
108 static void
109 local_apic_write(uint32_t reg, uint64_t value)
110 {
111     LOCAL_APIC_WRITE_REG(reg, (uint32_t)value);
112     apicadr[reg] = (uint32_t)value;
112 }
_____unchanged_portion_omitted_____
126 static void
127 local_apic_write_task_reg(uint64_t value)
128 {
129 #if defined(__amd64)
130     setcr8((ulong_t)(value >> APIC_IPL_SHIFT));
131 #else
132     if (apic_have_32bit_cr8)
133         setcr8((ulong_t)(value >> APIC_IPL_SHIFT));
134     else
135         LOCAL_APIC_WRITE_REG(APIC_TASK_REG, (uint32_t)value);
136     apicadr[APIC_TASK_REG] = (uint32_t)value;
136 #endif
137 }
_____unchanged_portion_omitted_____
232 int
233 apic_detect_x2apic(void)
234 {
236     struct cpuid_regs cp;
235     if (x2apic_enable == 0)
236         return (0);
238     return is_x86_feature(x86_featureset, X86FSET_X2APIC);
241     cp.cp_ecx = 1;
242     (void) __cpuid_insn(&cp);
```

```
1
```

```
new/usr/src/uts/i86pc/io/pcplusmp/apic_regops.c
```

```
244     return ((cp.cp_ecx & (0x1 << X2APIC_CPUID_BIT)) ? 1 : 0);
239 }
_____unchanged_portion_omitted_____

```

```
2
```

new/usr/src/uts/i86pc/io/ppm/acpisleep.c

```
*****
4291 Fri Jan 3 22:11:55 2014
new/usr/src/uts/i86pc/io/ppm/acpisleep.c
patch cpu-pause-func-deglobalize
*****
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/archsystem.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/machsystem.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 */
```

1

new/usr/src/uts/i86pc/io/ppm/acpisleep.c

```
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        }
110
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) {
```

2

```
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 Fri Jan 3 22:11:55 2014  
new/usr/src/uts/i86pc/os/cpr_impl.c  
patch cpu-pause-func-deglobalize  
*****  
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/cpuid.c
```

```
*****  
122611 Fri Jan 3 22:11:55 2014  
new/usr/src/uts/i86pc/os/cpuid.c  
patch x2apic-x86fset  
patch remove-unused-vars  
*****  
1 /*  
2 * CDDL HEADER START  
3 *  
4 * The contents of this file are subject to the terms of the  
5 * Common Development and Distribution License (the "License").  
6 * You may not use this file except in compliance with the License.  
7 *  
8 * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE  
9 * or http://www.opensolaris.org/os/licensing.  
10 * See the License for the specific language governing permissions  
11 * and limitations under the License.  
12 *  
13 * When distributing Covered Code, include this CDDL HEADER in each  
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.  
15 * If applicable, add the following below this CDDL HEADER, with the  
16 * fields enclosed by brackets "[]" replaced with your own identifying  
17 * information: Portions Copyright [yyyy] [name of copyright owner]  
18 *  
19 * CDDL HEADER END  
20 */  
21 /*  
22 * Copyright (c) 2004, 2010, Oracle and/or its affiliates. All rights reserved.  
23 * Copyright (c) 2011 by Delphix. All rights reserved.  
24 * Copyright 2013 Nexenta Systems, Inc. All rights reserved.  
25 * Copyright 2014 Josef "Jeff" Sipek <jeffpc@josefsipek.net>  
26 #endif /* ! codereview */  
27 */  
28 /*  
29 * Copyright (c) 2010, Intel Corporation.  
30 * All rights reserved.  
31 */  
32 /*  
33 * Portions Copyright 2009 Advanced Micro Devices, Inc.  
34 */  
35 /*  
36 * Copyright (c) 2012, Joyent, Inc. All rights reserved.  
37 */  
38 /*  
39 * Various routines to handle identification  
40 * and classification of x86 processors.  
41 */  
  
43 #include <sys/types.h>  
44 #include <sys/archsystm.h>  
45 #include <sys/x86_archext.h>  
46 #include <sys/kmem.h>  
47 #include <sys/sysm.h>  
48 #include <sys/cmn_err.h>  
49 #include <sys/sunddi.h>  
50 #include <sys/sunndi.h>  
51 #include <sys/cpuvar.h>  
52 #include <sys/processor.h>  
53 #include <sys/sysmacros.h>  
54 #include <sys/pg.h>  
55 #include <sys/fp.h>  
56 #include <sys/controlregs.h>  
57 #include <sys/bitmap.h>  
58 #include <sys/auxv_386.h>  
59 #include <sys/memnode.h>  
60 #include <sys/pci_cfgspace.h>
```

1

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

```
62 #ifdef __xpv  
63 #include <sys/hypervisor.h>  
64 #else  
65 #include <sys/ontrap.h>  
66 #endif  
  
68 /*  
69 * Pass 0 of cpuid feature analysis happens in locore. It contains special code  
70 * to recognize Cyrix processors that are not cpuid-compliant, and to deal with  
71 * them accordingly. For most modern processors, feature detection occurs here  
72 * in pass 1.  
73 */  
74 /* Pass 1 of cpuid feature analysis happens just at the beginning of mlsetup()  
75 * for the boot CPU and does the basic analysis that the early kernel needs.  
76 * x86_featureset is set based on the return value of cpuid_pass1() of the boot  
77 * CPU.  
78 */  
79 /* Pass 1 includes:  
80 *  
81 * o Determining vendor/model/family/stepping and setting x86_type and  
82 * x86_vendor accordingly.  
83 * o Processing the feature flags returned by the cpuid instruction while  
84 * applying any workarounds or tricks for the specific processor.  
85 * o Mapping the feature flags into Solaris feature bits (X86_*).  
86 * o Processing extended feature flags if supported by the processor,  
87 * again while applying specific processor knowledge.  
88 * o Determining the CMT characteristics of the system.  
89 */  
90 /* Pass 1 is done on non-boot CPUs during their initialization and the results  
91 * are used only as a meager attempt at ensuring that all processors within the  
92 * system support the same features.  
93 */  
94 /* Pass 2 of cpuid feature analysis happens just at the beginning  
95 * of startup(). It just copies in and corrects the remainder  
96 * of the cpuid data we depend on: standard cpuid functions that we didn't  
97 * need for pass1 feature analysis, and extended cpuid functions beyond the  
98 * simple feature processing done in pass1.  
99 */  
100 /* Pass 3 of cpuid analysis is invoked after basic kernel services; in  
101 * particular kernel memory allocation has been made available. It creates a  
102 * readable brand string based on the data collected in the first two passes.  
103 */  
104 /* Pass 4 of cpuid analysis is invoked after post_startup() when all  
105 * the support infrastructure for various hardware features has been  
106 * initialized. It determines which processor features will be reported  
107 * to userland via the aux vector.  
108 */  
109 /* All passes are executed on all CPUs, but only the boot CPU determines what  
110 * features the kernel will use.  
111 */  
112 /* Much of the worst junk in this file is for the support of processors  
113 * that didn't really implement the cpuid instruction properly.  
114 */  
115 /* NOTE: The accessor functions (cpuid_get*) are aware of, and ASSERT upon,  
116 * the pass numbers. Accordingly, changes to the pass code may require changes  
117 * to the accessor code.  
118 */  
  
120 uint_t x86_vendor = X86_VENDOR_IntelClone;  
121 uint_t x86_type = X86_TYPE_OTHER;  
122 uint_t x86_clflush_size = 0;  
  
124 uint_t pentiumpro_bug4046376;  
25 uint_t pentiumpro_bug4064495;
```

2

```

126 uchar_t x86_featureset[BT_SIZEOFMAP(NUM_X86_FEATURES)];
```

```

128 static char *x86_feature_names[NUM_X86_FEATURES] = {
129     "lgpg",
130     "tsc",
131     "msr",
132     "mtrr",
133     "pge",
134     "de",
135     "cmov",
136     "mmx",
137     "mca",
138     "pae",
139     "cv8",
140     "pat",
141     "sep",
142     "sse",
143     "sse2",
144     "htt",
145     "asysc",
146     "nx",
147     "sse3",
148     "cx16",
149     "cmp",
150     "tscp",
151     "mwait",
152     "sse4a",
153     "cpuid",
154     "ssse3",
155     "sse4_1",
156     "sse4_2",
157     "lgpg",
158     "clfsch",
159     "64",
160     "aes",
161     "pclmulqdq",
162     "xsave",
163     "avx",
164     "vmx",
165     "svm",
166     "topoext",
167     "f16c",
168     "rdrand",
169     "x2apic",
69     "rdrand"
170 };


---


unchanged portion omitted

```

119 uint_t enable486;

219 static size_t xsave_state_size = 0;
220 uint64_t xsave_bv_all = (XFEATURE_LEGACY_FP | XFEATURE_SSE);
221 boolean_t xsave_force_disable = B_FALSE;

223 /*
224 * This is set to platform type we are running on.
225 */
226 static int platform_type = -1;

228 #if !defined(__xpv)
229 /*
230 * Variable to patch if hypervisor platform detection needs to be
231 * disabled (e.g. platform_type will always be HW_NATIVE if this is 0).
232 */
233 int enable_platform_detection = 1;
234 #endif

```


```

```

236 /*
237 * monitor/mwait info.
238 *
239 * size_actual and buf_actual are the real address and size allocated to get
240 * proper mwait_buf alignment. buf_actual and size_actual should be passed
241 * to kmem_free(). Currently kmem_alloc() and mwait happen to both use
242 * processor cache-line alignment, but this is not guaranteed in the future.
243 */
244 struct mwait_info {
245     size_t          mon_min;      /* min size to avoid missed wakeups */
246     size_t          mon_max;      /* size to avoid false wakeups */
247     size_t          size_actual;   /* size actually allocated */
248     void           *buf_actual;   /* memory actually allocated */
249     uint32_t        support;      /* processor support of monitor/mwait */
250 };
unchanged portion omitted
936 void
937 cpuid_pass1(cpu_t *cpu, uchar_t *featureset)
938 {
939     uint32_t mask_ecx, mask_edx;
940     struct cpuid_info *cpi;
941     struct cpuid_regs *cp;
942     int xcpuid;
943 #if !defined(__xpv)
944     extern int idle_cpu_prefer_mwait;
945 #endif
946
947     /*
948     * Space statically allocated for BSP, ensure pointer is set
949     */
950     if (cpu->cpu_id == 0) {
951         if (cpu->cpu_m.mcpi == NULL)
952             cpu->cpu_m.mcpi = &cpuid_info0;
953     }
954
955     add_x86_feature(featureset, X86FSET_CPUID);
956
957     cpi = cpu->cpu_m.mcpi;
958     ASSERT(cpi != NULL);
959     cp = &cpi->cpi_std[0];
960     cp->cp_eax = 0;
961     cpi->cpi_maxeax = __cpuid_insn(cp);
962     {
963         uint32_t *iptr = (uint32_t *)cpi->cpi_vendorstr;
964         *iptr++ = cp->cp_ebx;
965         *iptr++ = cp->cp_edx;
966         *iptr++ = cp->cp_ecx;
967         *(char *)&cpi->cpi_vendorstr[12] = '\0';
968     }
969
970     cpi->cpi_vendor = __cpuid_vendorstr_to_vendorcode(cpi->cpi_vendorstr);
971     x86_vendor = cpi->cpi_vendor; /* for compatibility */
972
973     /*
974     * Limit the range in case of weird hardware
975     */
976     if (cpi->cpi_maxeax > CPI_MAXEAX_MAX)
977         cpi->cpi_maxeax = CPI_MAXEAX_MAX;
978     if (cpi->cpi_maxeax < 1)
979         goto pass1_done;
980
981     cp = &cpi->cpi_std[1];
982     cp->cp_eax = 1;
983     (void) __cpuid_insn(cp);

```

[new/usr/src/uts/i86pc/os/cpuid.c](#)

5

```

985      /*
986       * Extract identifying constants for easy access.
987       */
988     cpi->cpi_model = CPI_MODEL(cpi);
989     cpi->cpi_family = CPI_FAMILY(cpi);

991     if (cpi->cpi_family == 0xf)
992         cpi->cpi_family += CPI_FAMILY_XTD(cpi);

994     /*
995      * Beware: AMD uses "extended model" iff base *FAMILY* == 0xf.
996      * Intel, and presumably everyone else, uses model == 0xf, as
997      * one would expect (max value means possible overflow). Sigh.
998      */
999

1000    switch (cpi->cpi_vendor) {
1001      case X86_VENDOR_Intel:
1002          if (IS_EXTENDED_MODEL_INTEL(cpi))
1003              cpi->cpi_model += CPI_MODEL_XTD(cpi) << 4;
1004          break;
1005      case X86_VENDOR_AMD:
1006          if (CPI_FAMILY(cpi) == 0xf)
1007              cpi->cpi_model += CPI_MODEL_XTD(cpi) << 4;
1008          break;
1009      default:
1010          if (cpi->cpi_model == 0xf)
1011              cpi->cpi_model += CPI_MODEL_XTD(cpi) << 4;
1012          break;
1013    }

1015    cpi->cpi_step = CPI_STEP(cpi);
1016    cpi->cpi_brandid = CPI_BRANDID(cpi);

1018    /*
1019     * *default* assumptions:
1020     * - believe %edx feature word
1021     * - ignore %ecx feature word
1022     * - 32-bit virtual and physical addressing
1023     */
1024    mask_edx = 0xffffffff;
1025    mask_ecx = 0;

1027    cpi->cpi_pabits = cpi->cpi_vabits = 32;

1029    switch (cpi->cpi_vendor) {
1030      case X86_VENDOR_Intel:
1031          if (cpi->cpi_family == 5)
1032              x86_type = X86_TYPE_P5;
1033          else if (IS_LEGACY_P6(cpi)) {
1034              x86_type = X86_TYPE_P6;
1035              pentiumpro_bug4046376 = 1;
1036              pentiumpro_bug4064495 = 1;
1037              /*
1038               * Clear the SEP bit when it was set erroneously
1039               */
1040              if (cpi->cpi_model < 3 && cpi->cpi_step < 3)
1041                  cp->cp_edx &= ~CPUID_INTC_EDX_SEP;
1042          } else if (IS_NEW_F6(cpi) || cpi->cpi_family == 0xf) {
1043              x86_type = X86_TYPE_P4;
1044              /*
1045               * We don't currently depend on any of the %ecx
1046               * features until Prescott, so we'll only check
1047               * this from P4 onwards. We might want to revisit
1048               * that idea later.
1049               */

```

```
new/usr/src/uts/i86pc/os/cpuid.c
1049             mask_ecx = 0xffffffff;
1050         } else if (cpi->cpi_family > 0xf)
1051             mask_ecx = 0xffffffff;
1052         /*
1053          * We don't support MONITOR/MWAIT if leaf 5 is not available
1054          * to obtain the monitor linesize.
1055          */
1056         if (cpi->cpi_maxeax < 5)
1057             mask_ecx &= ~CPUID_INTC_ECX_MON;
1058         break;
1059     case X86_VENDOR_IntelClone:
1060     default:
1061         break;
1062     case X86_VENDOR_AMD:
1063 #if defined(OPTERON_ERRATUM_108)
1064         if (cpi->cpi_family == 0xf && cpi->cpi_model == 0xe) {
1065             cp->cp_eax = (0x0f0 & cp->cp_eax) | 0xc0;
1066             cpi->cpi_model = 0xc;
1067         } else
1068 #endif
1069         if (cpi->cpi_family == 5) {
1070             /*
1071              * AMD K5 and K6
1072              *
1073              * These CPUs have an incomplete implementation
1074              * of MCA/MCE which we mask away.
1075              */
1076         mask_edx &= ~(CPUID_INTC_EDX_MCE | CPUID_INTC_EDX_MCA);

1077         /*
1078          * Model 0 uses the wrong (APIC) bit
1079          * to indicate PGE. Fix it here.
1080          */
1081         if (cpi->cpi_model == 0) {
1082             if (cp->cp_edx & 0x200) {
1083                 cp->cp_edx &= ~0x200;
1084                 cp->cp_edx |= CPUID_INTC_EDX_PGE;
1085             }
1086         }
1087

1088         /*
1089          * Early models had problems w/ MMX; disable.
1090          */
1091         if (cpi->cpi_model < 6)
1092             mask_edx &= ~CPUID_INTC_EDX_MMX;
1093     }
1094

1095         /*
1096          * For newer families, SSE3 and CX16, at least, are valid;
1097          * enable all
1098          */
1099         if (cpi->cpi_family >= 0xf)
1100             mask_ecx = 0xffffffff;
1101         /*
1102          * We don't support MONITOR/MWAIT if leaf 5 is not available
1103          * to obtain the monitor linesize.
1104          */
1105         if (cpi->cpi_maxeax < 5)
1106             mask_ecx &= ~CPUID_INTC_ECX_MON;

1107 #if !defined(__xpv)
1108     /*
1109      * Do not use MONITOR/MWAIT to halt in the idle loop on any AMD
1110      * processors. AMD does not intend MWAIT to be used in the cpu
1111      * idle loop on current and future processors. 10h and future
1112      * AMD processors use more power in MWAIT than HLT.
1113
1114
```

```

1115             * Pre-family-10h Opterons do not have the MWAIT instruction.
1116             */
1117         idle_cpu_prefer_mwait = 0;
1118 #endif

1120         break;
1121     case X86_VENDOR_TM:
1122     /*
1123         * workaround the NT workaround in CMS 4.1
1124         */
1125     if (cpi->cpi_family == 5 && cpi->cpi_model == 4 &&
1126         (cpi->cpi_step == 2 || cpi->cpi_step == 3))
1127         cp->cp_edx |= CPUID_INTC_EDX_CX8;
1128     break;
1129     case X86_VENDOR_Centaur:
1130     /*
1131         * workaround the NT workarounds again
1132         */
1133     if (cpi->cpi_family == 6)
1134         cp->cp_edx |= CPUID_INTC_EDX_CX8;
1135     break;
1136     case X86_VENDOR_Cyrix:
1137     /*
1138         * We rely heavily on the probing in locore
1139         * to actually figure out what parts, if any,
1140         * of the Cyrix cpuid instruction to believe.
1141         */
1142     switch (x86_type) {
1143     case X86_TYPE_CYRIX_486:
1144         mask_edx = 0;
1145         break;
1146     case X86_TYPE_CYRIX_6x86:
1147         mask_edx = 0;
1148         break;
1149     case X86_TYPE_CYRIX_6x86L:
1150         mask_edx =
1151             CPUID_INTC_EDX_DE |
1152             CPUID_INTC_EDX_CX8;
1153         break;
1154     case X86_TYPE_CYRIX_6x86MX:
1155         mask_edx =
1156             CPUID_INTC_EDX_DE |
1157             CPUID_INTC_EDX_MSR |
1158             CPUID_INTC_EDX_CX8 |
1159             CPUID_INTC_EDX_PGE |
1160             CPUID_INTC_EDX_CMOV |
1161             CPUID_INTC_EDX_MMX;
1162         break;
1163     case X86_TYPE_CYRIX_GXm:
1164         mask_edx =
1165             CPUID_INTC_EDX_MSR |
1166             CPUID_INTC_EDX_CX8 |
1167             CPUID_INTC_EDX_CMOV |
1168             CPUID_INTC_EDX_MMX;
1169         break;
1170     case X86_TYPE_CYRIX_MediasGX:
1171         break;
1172     case X86_TYPE_CYRIX_MII:
1173     case X86_TYPE_VIA_CYRIX_III:
1174         mask_edx =
1175             CPUID_INTC_EDX_DE |
1176             CPUID_INTC_EDX_TSC |
1177             CPUID_INTC_EDX_MSR |
1178             CPUID_INTC_EDX_CX8 |
1179             CPUID_INTC_EDX_PGE |
1180             CPUID_INTC_EDX_CMOV |

```

```

1181                                         CPUID_INTC_EDX_MMX;
1182                                         break;
1183                                     default:
1184                                         break;
1185                                     }
1186                                 }
1187                             }

1188 #if defined(__xpv)
1189 /*
1190     * Do not support MONITOR/MWAIT under a hypervisor
1191     */
1192     mask_ecx &= ~CPUID_INTC_ECX_MON;
1193 /*
1194     * Do not support XSAVE under a hypervisor for now
1195     */
1196     xsave_force_disable = B_TRUE;
1197

1198 #endif /* __xpv */

1199 if (xsave_force_disable) {
1200     mask_ecx &= ~CPUID_INTC_ECX_XSAVE;
1201     mask_ecx &= ~CPUID_INTC_ECX_AVX;
1202     mask_ecx &= ~CPUID_INTC_ECX_F16C;
1203 }
1204

1205 /*
1206     * Now we've figured out the masks that determine
1207     * which bits we choose to believe, apply the masks
1208     * to the feature words, then map the kernel's view
1209     * of these feature words into its feature word.
1210     */
1211 cp->cp_edx &= mask_edx;
1212 cp->cp_ecx &= mask_ecx;

1213 /*
1214     * apply any platform restrictions (we don't call this
1215     * immediately after __cpuid_insn here, because we need the
1216     * workarounds applied above first)
1217     */
1218 platform_cpuid_mangle(cpi->cpi_vendor, 1, cp);

1219 /*
1220     * fold in overrides from the "eeprom" mechanism
1221     */
1222 cp->cp_edx |= cpuid_feature_edx_include;
1223 cp->cp_edx &= ~cpuid_feature_edx_exclude;

1224 cp->cp_ecx |= cpuid_feature_ecx_include;
1225 cp->cp_ecx &= ~cpuid_feature_ecx_exclude;

1226 if (cp->cp_edx & CPUID_INTC_EDX_PSE) {
1227     add_x86_feature(featureset, X86FSET_LARGEPAGE);
1228 }
1229 if (cp->cp_edx & CPUID_INTC_EDX_TSC) {
1230     add_x86_feature(featureset, X86FSET_TSC);
1231 }
1232 if (cp->cp_edx & CPUID_INTC_EDX_MSR) {
1233     add_x86_feature(featureset, X86FSET_MSR);
1234 }
1235 if (cp->cp_edx & CPUID_INTC_EDX_MTRR) {
1236     add_x86_feature(featureset, X86FSET_MTRR);
1237 }
1238 if (cp->cp_edx & CPUID_INTC_EDX_PGE) {
1239     add_x86_feature(featureset, X86FSET_PGE);
1240 }

```

```

1247     if (cp->cp_edx & CPUID_INTC_EDX_CMOV) {
1248         add_x86_feature(featureset, X86FSET_CMOV);
1249     }
1250     if (cp->cp_edx & CPUID_INTC_EDX_MMX) {
1251         add_x86_feature(featureset, X86FSET_MMX);
1252     }
1253     if ((cp->cp_edx & CPUID_INTC_EDX_MCE) != 0 &&
1254         (cp->cp_edx & CPUID_INTC_EDX_MCA) != 0) {
1255         add_x86_feature(featureset, X86FSET_MCA);
1256     }
1257     if (cp->cp_edx & CPUID_INTC_EDX_PAE) {
1258         add_x86_feature(featureset, X86FSET_PAE);
1259     }
1260     if (cp->cp_edx & CPUID_INTC_EDX_CX8) {
1261         add_x86_feature(featureset, X86FSET_CX8);
1262     }
1263     if (cp->cp_ecx & CPUID_INTC_ECX_CX16) {
1264         add_x86_feature(featureset, X86FSET_CX16);
1265     }
1266     if (cp->cp_edx & CPUID_INTC_EDX_PAT) {
1267         add_x86_feature(featureset, X86FSET_PAT);
1268     }
1269     if (cp->cp_edx & CPUID_INTC_EDX_SEP) {
1270         add_x86_feature(featureset, X86FSET_SEP);
1271     }
1272     if (cp->cp_edx & CPUID_INTC_EDX_FXSR) {
1273         /*
1274          * In our implementation, fxsave/fxrstor
1275          * are prerequisites before we'll even
1276          * try and do SSE things.
1277         */
1278     if (cp->cp_edx & CPUID_INTC_EDX_SSE) {
1279         add_x86_feature(featureset, X86FSET_SSE);
1280     }
1281     if (cp->cp_edx & CPUID_INTC_EDX_SSE2) {
1282         add_x86_feature(featureset, X86FSET_SSE2);
1283     }
1284     if (cp->cp_ecx & CPUID_INTC_ECX_SSE3) {
1285         add_x86_feature(featureset, X86FSET_SSE3);
1286     }
1287     if (cp->cp_ecx & CPUID_INTC_ECX_SSSE3) {
1288         add_x86_feature(featureset, X86FSET_SSSE3);
1289     }
1290     if (cp->cp_ecx & CPUID_INTC_ECX_SSE4_1) {
1291         add_x86_feature(featureset, X86FSET_SSE4_1);
1292     }
1293     if (cp->cp_ecx & CPUID_INTC_ECX_SSE4_2) {
1294         add_x86_feature(featureset, X86FSET_SSE4_2);
1295     }
1296     if (cp->cp_ecx & CPUID_INTC_ECX_AES) {
1297         add_x86_feature(featureset, X86FSET_AES);
1298     }
1299     if (cp->cp_ecx & CPUID_INTC_ECX_PCLMULQDQ) {
1300         add_x86_feature(featureset, X86FSET_PCLMULQDQ);
1301     }
1302
1303     if (cp->cp_ecx & CPUID_INTC_ECX_XSAVE) {
1304         add_x86_feature(featureset, X86FSET_XSAVE);
1305
1306         /* We only test AVX when there is XSAVE */
1307         if (cp->cp_ecx & CPUID_INTC_ECX_AVX) {
1308             add_x86_feature(featureset,
1309                             X86FSET_AVX);
1310
1311         if (cp->cp_ecx & CPUID_INTC_ECX_F16C)
1312             add_x86_feature(featureset,

```

```

1313                                         X86FSET_F16C);
1314
1315     }
1316 }
1317     if (cp->cp_ecx & CPUID_INTC_ECX_X2APIC) {
1318         add_x86_feature(featureset, X86FSET_X2APIC);
1319     }
1320 #endif /* ! codereview */
1321     if (cp->cp_edx & CPUID_INTC_EDX_DE) {
1322         add_x86_feature(featureset, X86FSET_DE);
1323     }
1324 #if !defined(__xpv)
1325     if (cp->cp_ecx & CPUID_INTC_ECX_MON) {
1326
1327         /*
1328          * We require the CLFLUSH instruction for erratum workaround
1329          * to use MONITOR/MWAIT.
1330         */
1331         if (cp->cp_edx & CPUID_INTC_EDX_CLFSH) {
1332             cpi->cpi_mwait.support |= MWAIT_SUPPORT;
1333             add_x86_feature(featureset, X86FSET_MWAIT);
1334         } else {
1335             extern int idle_cpu_assert_cflush_monitor;
1336
1337             /*
1338              * All processors we are aware of which have
1339              * MONITOR/MWAIT also have CLFLUSH.
1340             */
1341             if (idle_cpu_assert_cflush_monitor) {
1342                 ASSERT((cp->cp_ecx & CPUID_INTC_ECX_MON) &&
1343                        (cp->cp_edx & CPUID_INTC_EDX_CLFSH));
1344             }
1345         }
1346     }
1347 #endif /* __xpv */
1348
1349     if (cp->cp_ecx & CPUID_INTC_ECX_VMX) {
1350         add_x86_feature(featureset, X86FSET_VMX);
1351     }
1352
1353     if (cp->cp_ecx & CPUID_INTC_ECX_RDRAND)
1354         add_x86_feature(featureset, X86FSET_RDRAND);
1355
1356     /*
1357      * Only need it first time, rest of the cpus would follow suit.
1358      * we only capture this for the bootcpu.
1359     */
1360     if (cp->cp_edx & CPUID_INTC_EDX_CLFSH) {
1361         add_x86_feature(featureset, X86FSET_CLFSH);
1362         x86_clflush_size = (BITX(cp->cp_ebx, 15, 8) * 8);
1363     }
1364     if (is_x86_feature(featureset, X86FSET_PAE))
1365         cpi->cpi_pabits = 36;
1366
1367     /*
1368      * Hyperthreading configuration is slightly tricky on Intel
1369      * and pure clones, and even trickier on AMD.
1370      *
1371      * (AMD chose to set the HTT bit on their CMP processors,
1372      * even though they're not actually hyperthreaded. Thus it
1373      * takes a bit more work to figure out what's really going
1374      * on ... see the handling of the CMP_LGCV bit below)
1375     */
1376     if (cp->cp_edx & CPUID_INTC_EDX_HTT) {
1377         cpi->cpi_ncpu_per_chip = CPI_CPU_COUNT(cpi);
1378         if (cpi->cpi_ncpu_per_chip > 1)

```

```

1379         add_x86_feature(featureset, X86FSET_HTT);
1380     } else {
1381         cpi->cpi_ncpu_per_chip = 1;
1382     }
1384     /*
1385      * Work on the "extended" feature information, doing
1386      * some basic initialization for cpuid_pass2()
1387      */
1388     xcpuid = 0;
1389     switch (cpi->cpi_vendor) {
1390     case X86_VENDOR_Intel:
1391         if (IS_NEW_F6(cpi) || cpi->cpi_family >= 0xf)
1392             xcpuid++;
1393         break;
1394     case X86_VENDOR_AMD:
1395         if (cpi->cpi_family > 5 ||
1396             (cpi->cpi_family == 5 && cpi->cpi_model >= 1))
1397             xcpuid++;
1398         break;
1399     case X86_VENDOR_Cyrix:
1400         /*
1401          * Only these Cyrix CPUs are -known- to support
1402          * extended cpuid operations.
1403          */
1404         if (x86_type == X86_TYPE_VIA_CYRIX_III ||
1405             x86_type == X86_TYPE_CYRIX_GXm)
1406             xcpuid++;
1407         break;
1408     case X86_VENDOR_Centaur:
1409     case X86_VENDOR_TM:
1410     default:
1411         xcpuid++;
1412         break;
1413     }
1415     if (xcpuid) {
1416         cp = &cpi->cpi_extd[0];
1417         cp->cp_eax = 0x80000000;
1418         cpi->cpi_xmaxeax = __cpuid_insn(cp);
1419     }
1421     if (cpi->cpi_xmaxeax & 0x80000000) {
1422         if (cpi->cpi_xmaxeax > CPI_XMAXEAX_MAX)
1423             cpi->cpi_xmaxeax = CPI_XMAXEAX_MAX;
1424
1425         switch (cpi->cpi_vendor) {
1426         case X86_VENDOR_Intel:
1427         case X86_VENDOR_AMD:
1428             if (cpi->cpi_xmaxeax < 0x80000001)
1429                 break;
1430             cp = &cpi->cpi_extd[1];
1431             cp->cp_eax = 0x80000001;
1432             (void) __cpuid_insn(cp);
1433
1434             if (cpi->cpi_vendor == X86_VENDOR_AMD &&
1435                 cpi->cpi_family == 5 &&
1436                 cpi->cpi_model == 6 &&
1437                 cpi->cpi_step == 6) {
1438                 /*
1439                  * K6 model 6 uses bit 10 to indicate SYSC
1440                  * Later models use bit 11. Fix it here.
1441                  */
1442                 if (cp->cp_edx & 0x400) {
1443                     cp->cp_edx &= ~0x400;
1444

```

```

1445             cp->cp_edx |= CPUID_AMD_EDX_SYSC;
1446         }
1447     }
1449
1450     platform_cpuid_mangle(cpi->cpi_vendor, 0x80000001, cp);
1451     /*
1452      * Compute the additions to the kernel's feature word.
1453      */
1454     if (cp->cp_edx & CPUID_AMD_EDX_NX) {
1455         add_x86_feature(featureset, X86FSET_NX);
1456     }
1457
1458     /*
1459      * Regardless whether or not we boot 64-bit,
1460      * we should have a way to identify whether
1461      * the CPU is capable of running 64-bit.
1462      */
1463     if (cp->cp_edx & CPUID_AMD_EDX_LM) {
1464         add_x86_feature(featureset, X86FSET_64);
1465     }
1466
1467 #if defined(__amd64)
1468     /*
1469      * 1 GB large page - enable only for 64 bit kernel */
1470     if (cp->cp_edx & CPUID_AMD_EDX_1GPG) {
1471         add_x86_feature(featureset, X86FSET_1GPG);
1472     }
1473
1474     if ((cpi->cpi_vendor == X86_VENDOR_AMD) &&
1475         (cpi->cpi_std[1].cp_edx & CPUID_INTC_EDX_FXSR) &&
1476         (cp->cp_ecx & CPUID_AMD_ECX_SSE4A)) {
1477         add_x86_feature(featureset, X86FSET_SSE4A);
1478     }
1479
1480     /*
1481      * If both the HTT and CMP_LGCY bits are set,
1482      * then we're not actually HyperThreaded. Read
1483      * "AMD CPUID Specification" for more details.
1484      */
1485     if (cpi->cpi_vendor == X86_VENDOR_AMD &&
1486         is_x86_feature(featureset, X86FSET_HTT) &&
1487         (cp->cp_ecx & CPUID_AMD_ECX_CMP_LGCY)) {
1488         remove_x86_feature(featureset, X86FSET_HTT);
1489         add_x86_feature(featureset, X86FSET_CMP);
1490     }
1491 #if defined(__amd64)
1492
1493     /*
1494      * It's really tricky to support syscall/sysret in
1495      * the i386 kernel; we rely on sysenter/sysexit
1496      * instead. In the amd64 kernel, things are -way-
1497      * better.
1498      */
1499     if (cp->cp_edx & CPUID_AMD_EDX_SYSC) {
1500         add_x86_feature(featureset, X86FSET_ASYSC);
1501     }
1502
1503     /*
1504      * While we're thinking about system calls, note
1505      * that AMD processors don't support sysenter
1506      * in long mode at all, so don't try to program them.
1507      */
1508     if (x86_vendor == X86_VENDOR_AMD) {
1509         remove_x86_feature(featureset, X86FSET_SEP);
1510     }
1511 #endif

```

```

1577
1578
1579
1580
1581
1582
1583
1584
1585     /* On family 0xf cpuid fn 2 ECX[7:0] "NC" is
1586     * 1 less than the number of physical cores on
1587     * the chip. In family 0x10 this value can
1588     * be affected by "downcoreing" - it reflects
1589     * 1 less than the number of cores actually
1590     * enabled on this node.
1591     */
1592     cpi->cpi_ncore_per_chip =
1593         BITX((cpi)->cpi_extd[8].cp_ecx, 7, 0) + 1;
1594 }
1595 break;
1596 default:
1597     cpi->cpi_ncore_per_chip = 1;
1598     break;
1599 }
1600
1601 /*
1602  * Get CPUID data about TSC Invariance in Deep C-State.
1603  */
1604 switch (cpi->cpi_vendor) {
1605 case X86_VENDOR_Intel:
1606     if (cpi->cpi_maxeax >= 7) {
1607         cp = &cpi->cpi_extd[7];
1608         cp->cp_eax = 0x80000007;
1609         cp->cp_ecx = 0;
1610         (void) __cpuid_insn(cp);
1611     }
1612     break;
1613 default:
1614     break;
1615 }
1616 } else {
1617     cpi->cpi_ncore_per_chip = 1;
1618 }
1619
1620 /*
1621  * If more than one core, then this processor is CMP.
1622  */
1623 if (cpi->cpi_ncore_per_chip > 1) {
1624     add_x86_feature(featureset, X86FSET_CMP);
1625 }
1626
1627 /*
1628  * If the number of cores is the same as the number
1629  * of CPUs, then we cannot have HyperThreading.
1630  */
1631 if (cpi->cpi_ncpu_per_chip == cpi->cpi_ncore_per_chip) {
1632     remove_x86_feature(featureset, X86FSET_HTT);
1633 }
1634
1635 cpi->cpi_apicid = CPI_APIC_ID(cpi);
1636 cpi->cpi_procnodes_per_pkg = 1;
1637 cpi->cpi_cores_per_computunit = 1;
1638 if (is_x86_feature(featureset, X86FSET_HTT) == B_FALSE &&
1639     is_x86_feature(featureset, X86FSET_CMP) == B_FALSE) {
1640     /*
1641      * Single-core single-threaded processors.
1642      */
1643     cpi->cpi_chipid = -1;
1644     cpi->cpi_clogid = 0;
1645     cpi->cpi_coreid = cpu->cpu_id;
1646     cpi->cpi_pkgcoreid = 0;
1647     if (cpi->cpi_vendor == X86_VENDOR_AMD)
1648         cpi->cpi_procnodeid = BITX(cpi->cpi_apicid, 3, 0);
1649     else

```

```

1643         cpi->cpi_procnodeid = cpi->cpi_chipid;
1644     } else if (cpi->cpi_ncpu_per_chip > 1) {
1645         if (cpi->cpi_vendor == X86_VENDOR_Intel)
1646             cpuid_intel_getids(cpu, featureset);
1647         else if (cpi->cpi_vendor == X86_VENDOR_AMD)
1648             cpuid_amd_getids(cpu);
1649         else {
1650             /*
1651             * All other processors are currently
1652             * assumed to have single cores.
1653             */
1654             cpi->cpi_coreid = cpi->cpi_chipid;
1655             cpi->cpi_pkgcoreid = 0;
1656             cpi->cpi_procnodeid = cpi->cpi_chipid;
1657             cpi->cpi_computunitid = cpi->cpi_chipid;
1658         }
1659     }
1660     /*
1661      * Synthesize chip "revision" and socket type
1662      */
1663     cpi->cpi_chiprev = _cpuid_chiprev(cpi->cpi_vendor, cpi->cpi_family,
1664                                         cpi->cpi_model, cpi->cpi_step);
1665     cpi->cpi_chiprevstr = _cpuid_chiprevstr(cpi->cpi_vendor,
1666                                              cpi->cpi_family, cpi->cpi_model, cpi->cpi_step);
1667     cpi->cpi_socket = _cpuid_skt(cpi->cpi_vendor, cpi->cpi_family,
1668                                   cpi->cpi_model, cpi->cpi_step);
1669
1671 pass1_done:
1672     cpi->cpi_pass = 1;
1673 }
1675 /*
1676  * Make copies of the cpuid table entries we depend on, in
1677  * part for ease of parsing now, in part so that we have only
1678  * one place to correct any of it, in part for ease of
1679  * later export to userland, and in part so we can look at
1680  * this stuff in a crash dump.
1681 */
1683 /*ARGSUSED*/
1684 void
1685 cpuid_pass2(cpu_t *cpu)
1686 {
1687     uint_t n, nmax;
1688     int i;
1689     struct cpuid_regs *cp;
1690     uint8_t *dp;
1691     uint32_t *iptr;
1692     struct cpuid_info *cpi = cpu->cpu_m.mcpu_cpi;
1694     ASSERT(cpi->cpi_pass == 1);
1696     if (cpi->cpi_maxeax < 1)
1697         goto pass2_done;
1699     if ((nmax = cpi->cpi_maxeax + 1) > NMAX_CPI_STD)
1700         nmax = NMAX_CPI_STD;
1701     /*
1702      * (We already handled n == 0 and n == 1 in pass 1)
1703      */
1704     for (n = 2, cp = &cpi->cpi_std[2]; n < nmax; n++, cp++) {
1705         cp->cp_eax = n;
1707         /*
1708          * CPUID function 4 expects %ecx to be initialized

```

```

1709     * with an index which indicates which cache to return
1710     * information about. The OS is expected to call function 4
1711     * with %ecx set to 0, 1, 2, ... until it returns with
1712     * EAX[4:0] set to 0, which indicates there are no more
1713     * caches.
1714     *
1715     * Here, populate cpi_std[4] with the information returned by
1716     * function 4 when %ecx == 0, and do the rest in cpuid_pass3()
1717     * when dynamic memory allocation becomes available.
1718     *
1719     * Note: we need to explicitly initialize %ecx here, since
1720     * function 4 may have been previously invoked.
1721     */
1722     if (n == 4)
1723         cp->cp_ecx = 0;
1725     (void) __cpuid_insn(cp);
1726     platform_cpuid_mangle(cpi->cpi_vendor, n, cp);
1727     switch (n) {
1728     case 2:
1729         /*
1730          * "the lower 8 bits of the %eax register
1731          * contain a value that identifies the number
1732          * of times the cpuid [instruction] has to be
1733          * executed to obtain a complete image of the
1734          * processor's caching systems."
1735          *
1736          * How *do* they make this stuff up?
1737          */
1738         cpi->cpi_ncache = sizeof (*cp) *
1739                         BITX(cp->cp_eax, 7, 0);
1740         if (cpi->cpi_ncache == 0)
1741             break;
1742         cpi->cpi_ncache--;           /* skip count byte */
1744         /*
1745          * Well, for now, rather than attempt to implement
1746          * this slightly dubious algorithm, we just look
1747          * at the first 15 ..
1748          */
1749         if (cpi->cpi_ncache > (sizeof (*cp) - 1))
1750             cpi->cpi_ncache = sizeof (*cp) - 1;
1752         dp = cpi->cpi_cacheinfo;
1753         if (BITX(cp->cp_eax, 31, 31) == 0) {
1754             uint8_t *p = (void *)&cp->cp_eax;
1755             for (i = 1; i < 4; i++)
1756                 if (p[i] != 0)
1757                     *dp++ = p[i];
1759         }
1760         if (BITX(cp->cp_ebx, 31, 31) == 0) {
1761             uint8_t *p = (void *)&cp->cp_ebx;
1762             for (i = 0; i < 4; i++)
1763                 if (p[i] != 0)
1764                     *dp++ = p[i];
1765         }
1766         if (BITX(cp->cp_ecx, 31, 31) == 0) {
1767             uint8_t *p = (void *)&cp->cp_ecx;
1768             for (i = 0; i < 4; i++)
1769                 if (p[i] != 0)
1770                     *dp++ = p[i];
1771         }
1772         if (BITX(cp->cp_edx, 31, 31) == 0) {
1773             uint8_t *p = (void *)&cp->cp_edx;
1774             for (i = 0; i < 4; i++)
1775                 if (p[i] != 0)

```

```

1775             *dp++ = p[i];
1776         }
1777         break;
1778
1779     case 3: /* Processor serial number, if PSN supported */
1780         break;
1781
1782     case 4: /* Deterministic cache parameters */
1783         break;
1784
1785     case 5: /* Monitor/Mwait parameters */
1786     {
1787         size_t mwait_size;
1788
1789         /*
1790          * check cpi_mwait.support which was set in cpuid_pss1
1791          */
1792         if (!(cpi->cpi_mwait.support & MWAIT_SUPPORT))
1793             break;
1794
1795         /*
1796          * Protect ourself from insane mwait line size.
1797          * Workaround for incomplete hardware emulator(s).
1798          */
1799         mwait_size = (size_t)MWAIT_SIZE_MAX(cpi);
1800         if (mwait_size < sizeof(uint32_t) ||
1801             !ISP2(mwait_size)) {
1802 #if DEBUG
1803             cmn_err(CE_NOTE, "Cannot handle cpu %d mwait "
1804                     "size %ld", cpu->cpu_id, (long)mwait_size);
1805 #endif
1806             break;
1807         }
1808
1809         cpi->cpi_mwait.mon_min = (size_t)MWAIT_SIZE_MIN(cpi);
1810         cpi->cpi_mwait.mon_max = mwait_size;
1811         if (MWAIT_EXTENSION(cpi)) {
1812             cpi->cpi_mwait.support |= MWAIT_EXTENSIONS;
1813             if (MWAIT_INT_ENABLE(cpi))
1814                 cpi->cpi_mwait.support |=
1815                         MWAIT_ECX_INT_ENABLE;
1816         }
1817         break;
1818     }
1819     default:
1820         break;
1821     }
1822 }
1823
1824 if (cpi->cpi_maxeax >= 0xB && cpi->cpi_vendor == X86_VENDOR_Intel) {
1825     struct cpuid_regs regs;
1826
1827     cp = &regs;
1828     cp->cp_eax = 0xB;
1829     cp->cp_edx = cp->cp_ebx = cp->cp_ecx = 0;
1830
1831     (void) __cpuid_insn(cp);
1832
1833     /*
1834      * Check CPUID.EAX=0BH, ECX=0H:EBX is non-zero, which
1835      * indicates that the extended topology enumeration leaf is
1836      * available.
1837      */
1838     if (cp->cp_ebx) {
1839         uint32_t x2apic_id;
1840         uint_t coreid_shift = 0;

```

```

1841         uint_t ncpu_per_core = 1;
1842         uint_t chipid_shift = 0;
1843         uint_t ncpu_per_chip = 1;
1844         uint_t i;
1845         uint_t level;
1846
1847         for (i = 0; i < CPI_FNB_ECX_MAX; i++) {
1848             cp->cp_eax = 0xB;
1849             cp->cp_ecx = i;
1850
1851             (void) __cpuid_insn(cp);
1852             level = CPI_CPU_LEVEL_TYPE(cp);
1853
1854             if (level == 1) {
1855                 x2apic_id = cp->cp_edx;
1856                 coreid_shift = BITX(cp->cp_eax, 4, 0);
1857                 ncpu_per_core = BITX(cp->cp_ebx, 15, 0);
1858             } else if (level == 2) {
1859                 x2apic_id = cp->cp_edx;
1860                 chipid_shift = BITX(cp->cp_eax, 4, 0);
1861                 ncpu_per_chip = BITX(cp->cp_ebx, 15, 0);
1862             }
1863         }
1864
1865         cpi->cpi_apicid = x2apic_id;
1866         cpi->cpi_ncpu_per_chip = ncpu_per_chip;
1867         cpi->cpi_ncore_per_chip = ncpu_per_chip /
1868             ncpu_per_core;
1869         cpi->cpi_chipid = x2apic_id >> chipid_shift;
1870         cpi->cpi_clogid = x2apic_id & ((1 << chipid_shift) - 1);
1871         cpi->cpi_coreid = x2apic_id >> coreid_shift;
1872         cpi->cpi_pkgcoreid = cpi->cpi_clogid >> coreid_shift;
1873     }
1874
1875     /* Make cp NULL so that we don't stumble on others */
1876     cp = NULL;
1877 }
1878
1879 /*
1880  * XSAVE enumeration
1881  */
1882 if (cpi->cpi_maxeax >= 0xD) {
1883     struct cpuid_regs regs;
1884     boolean_t cpuid_d_valid = B_TRUE;
1885
1886     cp = &regs;
1887     cp->cp_eax = 0xD;
1888     cp->cp_edx = cp->cp_ebx = cp->cp_ecx = 0;
1889
1890     (void) __cpuid_insn(cp);
1891
1892     /*
1893      * Sanity checks for debug
1894      */
1895     if ((cp->cp_eax & XFEATURE_LEGACY_FP) == 0 ||
1896         (cp->cp_eax & XFEATURE_SSE) == 0) {
1897         cpuid_d_valid = B_FALSE;
1898     }
1899
1900     cpi->cpi_xsav_xsav_hw_features_low = cp->cp_eax;
1901     cpi->cpi_xsav_xsav_hw_features_high = cp->cp_edx;
1902     cpi->cpi_xsav_xsav_max_size = cp->cp_ecx;
1903
1904     /*
1905      * If the hw supports AVX, get the size and offset in the save
1906      * area for the ymm state.
1907     */

```

```

1907 */
1908 if (cpi->cpi_xsave.xsav_hw_features_low & XFEATURE_AVX) {
1909     cp->cp_eax = 0xD;
1910     cp->cp_ecx = 2;
1911     cp->cp_edx = cp->cp_ebx = 0;
1912
1913     (void) __cpuid_insn(cp);
1914
1915     if (cp->cp_ebx != CPUID_LEAFD_2_YMM_OFFSET ||
1916         cp->cp_eax != CPUID_LEAFD_2_YMM_SIZE) {
1917             cpuid_d_valid = B_FALSE;
1918         }
1919
1920     cpi->cpi_xsave.ymm_size = cp->cp_eax;
1921     cpi->cpi_xsave.ymm_offset = cp->cp_ebx;
1922 }
1923
1924 if (is_x86_feature(x86_featureset, X86FSET_XSAVE)) {
1925     xsave_state_size = 0;
1926 } else if (cpuid_d_valid) {
1927     xsave_state_size = cpi->cpi_xsave.xsav_max_size;
1928 } else {
1929     /* Broken CPUID 0xD, probably in HVM */
1930     cmn_err(CE_WARN, "cpu%d: CPUID.0xD returns invalid \"%s\" value: hw_low = %d, hw_high = %d, xsave_size = %d\""
1931             ", ymm_size = %d, ymm_offset = %d\\n",
1932             cp->cpu_id, cpi->cpi_xsave.xsav_hw_features_low,
1933             cpi->cpi_xsave.xsav_hw_features_high,
1934             (int)cpi->cpi_xsave.xsav_max_size,
1935             (int)cpi->cpi_xsave.ymm_size,
1936             (int)cpi->cpi_xsave.ymm_offset);
1937
1938     if (xsave_state_size != 0) {
1939         /*
1940         * This must be a non-boot CPU. We cannot
1941         * continue, because boot cpu has already
1942         * enabled XSAVE.
1943         */
1944         ASSERT(cpu->cpu_id != 0);
1945         cmn_err(CE_PANIC, "cpu%d: we have already \"%s\""
1946                 "enabled XSAVE on boot cpu, cannot "
1947                 "continue.", cpu->cpu_id);
1948     } else {
1949         /*
1950         * Must be from boot CPU, OK to disable XSAVE.
1951         */
1952         ASSERT(cpu->cpu_id == 0);
1953         remove_x86_feature(x86_featureset,
1954                             X86FSET_XSAVE);
1955         remove_x86_feature(x86_featureset, X86FSET_AVX);
1956         CPI_FEATURES_ECX(cpi) &= ~CPUID_INTC_ECX_XSAVE;
1957         CPI_FEATURES_ECX(cpi) &= ~CPUID_INTC_ECX_AVX;
1958         CPI_FEATURES_ECX(cpi) &= ~CPUID_INTC_ECX_F16C;
1959         xsave_force_disable = B_TRUE;
1960     }
1961 }
1962
1963 }

1964 if ((cpi->cpi_xmaxeax & 0x80000000) == 0)
1965     goto pass2_done;

1966 if ((nmax = cpi->cpi_xmaxeax - 0x80000000 + 1) > NMAX_CPI_EXTD)
1967     nmax = NMAX_CPI_EXTD;
1968
1969 /*
1970 * Copy the extended properties, fixing them as we go.

```

```

1973 * (We already handled n == 0 and n == 1 in pass 1)
1974 */
1975 iptr = (void *)cpi->cpi_brandstr;
1976 for (n = 2, cp = &cpi->cpi_extd[2]; n < nmax; cp++, n++) {
1977     cp->cp_eax = 0x80000000 + n;
1978     (void) __cpuid_insn(cp);
1979     platform_cpuid_mangle(cpi->cpi_vendor, 0x80000000 + n, cp);
1980     switch (n) {
1981         case 2:
1982         case 3:
1983         case 4:
1984             /*
1985                 * Extract the brand string
1986             */
1987             *iptr++ = cp->cp_eax;
1988             *iptr++ = cp->cp_ebx;
1989             *iptr++ = cp->cp_ecx;
1990             *iptr++ = cp->cp_edx;
1991             break;
1992     case 5:
1993         switch (cpi->cpi_vendor) {
1994             case X86_VENDOR_AMD:
1995                 /*
1996                     * The Athlon and Duron were the first
1997                     * parts to report the sizes of the
1998                     * TLB for large pages. Before then,
1999                     * we don't trust the data.
2000                 */
2001                 if (cpi->cpi_family < 6 ||
2002                     (cpi->cpi_family == 6 &&
2003                         cpi->cpi_model < 1))
2004                     cp->cp_eax = 0;
2005                 break;
2006         default:
2007             break;
2008     }
2009     break;
2010 }
2011 case 6:
2012     switch (cpi->cpi_vendor) {
2013         case X86_VENDOR_AMD:
2014             /*
2015                 * The Athlon and Duron were the first
2016                 * AMD parts with L2 TLB's.
2017                 * Before then, don't trust the data.
2018             */
2019             if (cpi->cpi_family < 6 ||
2020                 cpi->cpi_family == 6 &&
2021                     cpi->cpi_model < 1)
2022                 cp->cp_eax = cp->cp_ebx = 0;
2023             /*
2024                 * AMD Duron rev A0 reports L2
2025                 * cache size incorrectly as 1K
2026                 * when it is really 64K
2027             */
2028             if (cpi->cpi_family == 6 &&
2029                 cpi->cpi_model == 3 &&
2030                     cpi->cpi_step == 0) {
2031                 cp->cp_ecx &= 0xfffff;
2032                 cp->cp_ecx |= 0x400000;
2033             }
2034             break;
2035     case X86_VENDOR_Cyrix: /* VIA C3 */
2036             /*
2037                 * VIA C3 processors are a bit messed
2038                 * up w.r.t. encoding cache sizes in %ecx
2039             */

```

```

2039         if (cpi->cpi_family != 6)
2040             break;
2041         /*
2042          * model 7 and 8 were incorrectly encoded
2043          *
2044          * xxx is model 8 really broken?
2045          */
2046         if (cpi->cpi_model == 7 || cpi->cpi_model == 8)
2047             cp->cp_ecx =
2048                 BITX(cp->cp_ecx, 31, 24) << 16 |
2049                 BITX(cp->cp_ecx, 23, 16) << 12 |
2050                 BITX(cp->cp_ecx, 15, 8) << 8 |
2051                 BITX(cp->cp_ecx, 7, 0);
2052
2053         /*
2054          * model 9 stepping 1 has wrong associativity
2055          */
2056         if (cpi->cpi_model == 9 && cpi->cpi_step == 1)
2057             cp->cp_ecx |= 8 << 12;
2058         break;
2059     case X86_VENDOR_Intel:
2060     /*
2061        * Extended L2 Cache features function.
2062        * First appeared on Prescott.
2063        */
2064     default:
2065         break;
2066     }
2067     break;
2068 default:
2069     break;
2070 }
2071 }

2072 pass2_done:
2073     cpi->cpi_pass = 2;
2074 }
2075 }

2076 static const char *
2077 intel_cpubrand(const struct cpuid_info *cpi)
2078 {
2079     int i;
2080
2081     if (!is_x86_feature(x86_featureset, X86FSET_CPUID) ||
2082         cpi->cpi_maxeax < 1 || cpi->cpi_family < 5)
2083         return ("i486");
2084
2085     switch (cpi->cpi_family) {
2086     case 5:
2087         return ("Intel Pentium(r)");
2088     case 6:
2089         switch (cpi->cpi_model) {
2090             uint_t celeron, xeon;
2091             const struct cpuid_regs *cp;
2092
2093             case 0:
2094             case 1:
2095             case 2:
2096                 return ("Intel Pentium(r) Pro");
2097             case 3:
2098             case 4:
2099                 return ("Intel Pentium(r) II");
2100             case 6:
2101                 return ("Intel Celeron(r)");
2102             case 5:
2103             case 7:
2104                 celeron = xeon = 0;

```

```

2105         cp = &cpi->cpi_std[2]; /* cache info */
2106
2107         for (i = 1; i < 4; i++) {
2108             uint_t tmp;
2109
2110             tmp = (cp->cp_eax >> (8 * i)) & 0xff;
2111             if (tmp == 0x40)
2112                 celeron++;
2113             if (tmp >= 0x44 && tmp <= 0x45)
2114                 xeon++;
2115         }
2116
2117         for (i = 0; i < 2; i++) {
2118             uint_t tmp;
2119
2120             tmp = (cp->cp_ebx >> (8 * i)) & 0xff;
2121             if (tmp == 0x40)
2122                 celeron++;
2123             else if (tmp >= 0x44 && tmp <= 0x45)
2124                 xeon++;
2125         }
2126
2127         for (i = 0; i < 4; i++) {
2128             uint_t tmp;
2129
2130             tmp = (cp->cp_ecx >> (8 * i)) & 0xff;
2131             if (tmp == 0x40)
2132                 celeron++;
2133             else if (tmp >= 0x44 && tmp <= 0x45)
2134                 xeon++;
2135         }
2136
2137         for (i = 0; i < 4; i++) {
2138             uint_t tmp;
2139
2140             tmp = (cp->cp_edx >> (8 * i)) & 0xff;
2141             if (tmp == 0x40)
2142                 celeron++;
2143             else if (tmp >= 0x44 && tmp <= 0x45)
2144                 xeon++;
2145         }
2146
2147         if (celeron)
2148             return ("Intel Celeron(r)");
2149         if (xeon)
2150             return (cpi->cpi_model == 5 ?
2151                     "Intel Pentium(r) II Xeon(tm)" :
2152                     "Intel Pentium(r) III Xeon(tm)");
2153         return (cpi->cpi_model == 5 ?
2154                     "Intel Pentium(r) II or Pentium(r) II Xeon(tm)" :
2155                     "Intel Pentium(r) III or Pentium(r) III Xeon(tm)");
2156         default:
2157             break;
2158     }
2159     default:
2160         break;
2161     }
2162
2163     /* BrandID is present if the field is nonzero */
2164     if (cpi->cpi_brandid != 0) {
2165         static const struct {
2166             uint_t bt_bid;
2167             const char *bt_str;
2168         } brand_tbl[] = {
2169             { 0x1, "Intel(r) Celeron(r)" },
2170             { 0x2, "Intel(r) Pentium(r) III" },

```

```

2171 { 0x3, "Intel(r) Pentium(r) III Xeon(tm)" },
2172 { 0x4, "Intel(r) Pentium(r) III" },
2173 { 0x6, "Mobile Intel(r) Pentium(r) III" },
2174 { 0x7, "Mobile Intel(r) Celeron(r)" },
2175 { 0x8, "Intel(r) Pentium(r) 4" },
2176 { 0x9, "Intel(r) Pentium(r) 4" },
2177 { 0xa, "Intel(r) Celeron(r)" },
2178 { 0xb, "Intel(r) Xeon(tm)" },
2179 { 0xc, "Intel(r) Xeon(tm) MP" },
2180 { 0xe, "Mobile Intel(r) Pentium(r) 4" },
2181 { 0xf, "Mobile Intel(r) Celeron(r)" },
2182 { 0x11, "Mobile Genuine Intel(r)" },
2183 { 0x12, "Intel(r) Celeron(r) M" },
2184 { 0x13, "Mobile Intel(r) Celeron(r)" },
2185 { 0x14, "Intel(r) Celeron(r)" },
2186 { 0x15, "Mobile Genuine Intel(r)" },
2187 { 0x16, "Intel(r) Pentium(r) M" },
2188 { 0x17, "Mobile Intel(r) Celeron(r)" }
2189 };
2190 uint_t btblmax = sizeof(brand_tbl) / sizeof(brand_tbl[0]);
2191 uint_t sgn;
2192
2193 sgn = (cpi->cpi_family << 8) |
2194     (cpi->cpi_model << 4) | cpi->cpi_step;
2195
2196 for (i = 0; i < btblmax; i++)
2197     if (brand_tbl[i].bt_bid == cpi->cpi_brandid)
2198         break;
2199
2200 if (i < btblmax) {
2201     if (sgn == 0x6b1 && cpi->cpi_brandid == 3)
2202         return ("Intel(r) Celeron(r)");
2203     if (sgn < 0xf13 && cpi->cpi_brandid == 0xb)
2204         return ("Intel(r) Xeon(tm) MP");
2205     if (sgn < 0xf13 && cpi->cpi_brandid == 0xe)
2206         return ("Intel(r) Xeon(tm)");
2207 }
2208
2209 return (NULL);
2210
2211 }
2212
2213 static const char *
2214 amd_cpubrand(const struct cpuid_info *cpi)
2215 {
2216     if (!is_x86_feature(x86_featureset, X86FSET_CPUID) ||
2217         cpi->cpi_maxeax < 1 || cpi->cpi_family < 5)
2218         return ("i486 compatible");
2219
2220     switch (cpi->cpi_family) {
2221         case 5:
2222             switch (cpi->cpi_model) {
2223                 case 0:
2224                 case 1:
2225                 case 2:
2226                 case 3:
2227                 case 4:
2228                 case 5:
2229                     return ("AMD-K5(r)");
2230                 case 6:
2231                 case 7:
2232                     return ("AMD-K6(r)");
2233                 case 8:
2234                     return ("AMD-K6(r)-2");
2235                 case 9:
2236                     return ("AMD-K6(r)-III");

```

```

2237     default:
2238         return ("AMD (family 5)");
2239     }
2240
2241     case 6:
2242         switch (cpi->cpi_model) {
2243             case 1:
2244                 return ("AMD-K7(tm)");
2245             case 0:
2246             case 2:
2247             case 4:
2248                 return ("AMD Athlon(tm)");
2249             case 3:
2250             case 7:
2251                 return ("AMD Duron(tm)");
2252             case 6:
2253             case 8:
2254             case 10:
2255                 /*
2256                  * Use the L2 cache size to distinguish
2257                  */
2258             case 11:
2259                 return ((cpi->cpi_extd[6].cp_ecx >> 16) >= 256 ?
2260                         "AMD Athlon(tm)" : "AMD Duron(tm)");
2261             default:
2262                 return ("AMD (family 6)");
2263             }
2264         default:
2265             break;
2266     }
2267
2268     if (cpi->cpi_family == 0xf && cpi->cpi_model == 5 &&
2269         cpi->cpi_brandid != 0) {
2270         switch (BITX(cpi->cpi_brandid, 7, 5)) {
2271             case 3:
2272                 return ("AMD Opteron(tm) UP 1xx");
2273             case 4:
2274                 return ("AMD Opteron(tm) DP 2xx");
2275             case 5:
2276                 return ("AMD Opteron(tm) MP 8xx");
2277             default:
2278                 return ("AMD Opteron(tm)");
2279         }
2280     }
2281
2282     return (NULL);
2283 }
2284
2285 static const char *
2286 cyrix_cpubrand(struct cpuid_info *cpi, uint_t type)
2287 {
2288     if (!is_x86_feature(x86_featureset, X86FSET_CPUID) ||
2289         cpi->cpi_maxeax < 1 || cpi->cpi_family < 5 |||
2290         type == X86_TYPE_CYRIX_486)
2291         return ("i486 compatible");
2292
2293     switch (type) {
2294         case X86_TYPE_CYRIX_6x86:
2295             return ("Cyrix 6x86");
2296         case X86_TYPE_CYRIX_6x86L:
2297             return ("Cyrix 6x86L");
2298         case X86_TYPE_CYRIX_6x86MX:
2299             return ("Cyrix 6x86MX");
2300         case X86_TYPE_CYRIX_GXm:
2301             return ("Cyrix GXm");
2302         case X86_TYPE_CYRIX_MediaGX:
2303             return ("Cyrix MediaGX");
2304         case X86_TYPE_CYRIX_MII:
2305             return ("Cyrix MII");
2306     }

```

```

2303     return ("Cyrix M2");
2304 case X86_TYPE_VIA_CYRIX_III:
2305     return ("VIA Cyrix M3");
2306 default:
2307     /*
2308      * Have another wild guess ...
2309      */
2310     if (cpi->cpi_family == 4 && cpi->cpi_model == 9)
2311         return ("Cyrix 5x86");
2312     else if (cpi->cpi_family == 5) {
2313         switch (cpi->cpi_model) {
2314             case 2:
2315                 return ("Cyrix 6x86"); /* Cyrix M1 */
2316             case 4:
2317                 return ("Cyrix MediaGX");
2318             default:
2319                 break;
2320         }
2321     } else if (cpi->cpi_family == 6) {
2322         switch (cpi->cpi_model) {
2323             case 0:
2324                 return ("Cyrix 6x86MX"); /* Cyrix M2? */
2325             case 5:
2326             case 6:
2327             case 7:
2328             case 8:
2329             case 9:
2330                 return ("VIA C3");
2331             default:
2332                 break;
2333         }
2334     }
2335     break;
2336 }
2337 return (NULL);
2338 }

2340 /*
2341 * This only gets called in the case that the CPU extended
2342 * feature brand string (0x80000002, 0x80000003, 0x80000004)
2343 * aren't available, or contain null bytes for some reason.
2344 */
2345 static void
2346 fabricate_brandstr(struct cpuid_info *cpi)
2347 {
2348     const char *brand = NULL;

2349     switch (cpi->cpi_vendor) {
2350 case X86_VENDOR_Intel:
2351         brand = intel_cpubrand(cpi);
2352         break;
2353 case X86_VENDOR_AMD:
2354         brand = amd_cpubrand(cpi);
2355         break;
2356 case X86_VENDOR_Cyrix:
2357         brand = cyrix_cpubrand(cpi, x86_type);
2358         break;
2359 case X86_VENDOR_NexGen:
2360         if (cpi->cpi_family == 5 && cpi->cpi_model == 0)
2361             brand = "NexGen Nx586";
2362         break;
2363 case X86_VENDOR_Centaur:
2364         if (cpi->cpi_family == 5)
2365             switch (cpi->cpi_model) {
2366                 case 4:
2367                     brand = "Centaur C6";

```

```

2369     break;
2370 case 8:
2371     brand = "Centaur C2";
2372     break;
2373 case 9:
2374     brand = "Centaur C3";
2375     break;
2376 default:
2377     break;
2378 }
2379 break;
2380 case X86_VENDOR_Rise:
2381     if (cpi->cpi_family == 5 &&
2382         (cpi->cpi_model == 0 || cpi->cpi_model == 2))
2383         brand = "Rise mP6";
2384     break;
2385 case X86_VENDOR_SiS:
2386     if (cpi->cpi_family == 5 && cpi->cpi_model == 0)
2387         brand = "SiS 55x";
2388     break;
2389 case X86_VENDOR_TM:
2390     if (cpi->cpi_family == 5 && cpi->cpi_model == 4)
2391         brand = "Transmeta Crusoe TM3x00 or TM5x00";
2392     break;
2393 case X86_VENDOR_NSC:
2394 case X86_VENDOR_UMC:
2395 default:
2396     break;
2397 }
2398 if (brand) {
2399     (void) strcpy((char *)cpi->cpi_brandstr, brand);
2400     return;
2401 }

2402 /*
2403 * If all else fails ...
2404 */
2405 (void) sprintf(cpi->cpi_brandstr, sizeof (cpi->cpi_brandstr),
2406 "%s.%d.%d", cpi->cpi_vendorstr, cpi->cpi_family,
2407 cpi->cpi_model, cpi->cpi_step);
2408 }

2409 */

2410 /*
2411 * This routine is called just after kernel memory allocation
2412 * becomes available on cpu0, and as part of mp_startup() on
2413 * the other cpus.
2414 *
2415 * Fixup the brand string, and collect any information from cpuid
2416 * that requires dynamically allocated storage to represent.
2417 */
2418 /*ARGSUSED*/
2419 void
2420 cpuid_pass3(cpu_t *cpu)
2421 {
2422     int i, max, shft, level, size;
2423     struct cpuid_regs regs;
2424     struct cpuid_regs *cp;
2425     struct cpuid_info *cpi = cpu->cpu_m.mcpu_cpi;
2426
2427     ASSERT(cpi->cpi_pass == 2);

2428 */

2429 /*
2430 * Function 4: Deterministic cache parameters
2431 *
2432 * Take this opportunity to detect the number of threads
2433 * sharing the last level cache, and construct a corresponding
2434

```

```

2435     * cache id. The respective cpuid_info members are initialized
2436     * to the default case of "no last level cache sharing".
2437     */
2438     cpi->cpi_ncpu_shr_last_cache = 1;
2439     cpi->cpi_last_lvl_cacheid = cpi->cpu_id;

2441 if (cpi->cpi_maxeax >= 4 && cpi->cpi_vendor == X86_VENDOR_Intel) {
2442
2443     /*
2444      * Find the # of elements (size) returned by fn 4, and along
2445      * the way detect last level cache sharing details.
2446      */
2447     bzero(&regs, sizeof (regs));
2448     cp = &regs;
2449     for (i = 0, max = 0; i < CPI_FN4_ECX_MAX; i++) {
2450         cp->cp_eax = 4;
2451         cp->cp_ecx = i;
2452
2453         (void) __cpuid_insn(cp);
2454
2455         if (CPI_CACHE_TYPE(cp) == 0)
2456             break;
2457         level = CPI_CACHE_LVL(cp);
2458         if (level > max) {
2459             max = level;
2460             cpi->cpi_ncpu_shr_last_cache =
2461                 CPI_NTHR_SHR_CACHE(cp) + 1;
2462         }
2463     }
2464     cpi->cpi_std_4_size = size = i;
2465
2466     /*
2467      * Allocate the cpi_std_4 array. The first element
2468      * references the regs for fn 4, %ecx == 0, which
2469      * cpuid_pass2() stashed in cpi->cpi_std[4].
2470      */
2471     if (size > 0) {
2472         cpi->cpi_std_4 =
2473             kmem_alloc(size * sizeof (cp), KM_SLEEP);
2474         cpi->cpi_std_4[0] = &cpi->cpi_std[4];
2475
2476         /*
2477          * Allocate storage to hold the additional regs
2478          * for function 4, %ecx == 1 .. cpi_std_4_size.
2479          *
2480          * The regs for fn 4, %ecx == 0 has already
2481          * been allocated as indicated above.
2482          */
2483         for (i = 1; i < size; i++) {
2484             cp = cpi->cpi_std_4[i] =
2485                 kmem_zalloc(sizeof (regs), KM_SLEEP);
2486             cp->cp_eax = 4;
2487             cp->cp_ecx = i;
2488
2489             (void) __cpuid_insn(cp);
2490         }
2491     }
2492
2493     /*
2494      * Determine the number of bits needed to represent
2495      * the number of CPUs sharing the last level cache.
2496      *
2497      * Shift off that number of bits from the APIC id to
2498      * derive the cache id.
2499      */
2500     shft = 0;
2501     for (i = 1; i < cpi->cpi_ncpu_shr_last_cache; i <= 1)

```

```

2501                                         shft++;
2502                                         cpi->cpi_last_lvl_cacheid = cpi->cpi_apicid >> shft;
2503 }
2504
2505 /*
2506  * Now fixup the brand string
2507 */
2508 if ((cpi->cpi_xmaxeax & 0x80000000) == 0) {
2509     fabricate_brandstr(cpi);
2510 } else {
2511
2512     /*
2513      * If we successfully extracted a brand string from the cpuid
2514      * instruction, clean it up by removing leading spaces and
2515      * similar junk.
2516      */
2517     if (cpi->cpi_brandstr[0]) {
2518         size_t maxlen = sizeof (cpi->cpi_brandstr);
2519         char *src, *dst;
2520
2521         dst = src = (char *)cpi->cpi_brandstr;
2522         src[maxlen - 1] = '\0';
2523
2524         /*
2525          * strip leading spaces
2526          */
2527         while (*src == ' ')
2528             src++;
2529
2530         /*
2531          * Remove any 'Genuine' or "Authentic" prefixes
2532          */
2533         if (strncmp(src, "Genuine ", 8) == 0)
2534             src += 8;
2535         if (strncmp(src, "Authentic ", 10) == 0)
2536             src += 10;
2537
2538         /*
2539          * Now do an in-place copy.
2540          * Map (R) to (r) and (TM) to (tm).
2541          * The era of teletypes is long gone, and there's
2542          * -really- no need to shout.
2543          */
2544         while (*src != '\0') {
2545             if (src[0] == '(') {
2546                 if (strncmp(src + 1, "R)", 2) == 0) {
2547                     (void) strcpy(dst, "(r)", 3);
2548                     src += 3;
2549                     dst += 3;
2550                     continue;
2551                 }
2552                 if (strncmp(src + 1, "TM)", 3) == 0) {
2553                     (void) strcpy(dst, "(tm)", 4);
2554                     src += 4;
2555                     dst += 4;
2556                     continue;
2557                 }
2558             }
2559             *dst++ = *src++;
2560         }
2561         *dst = '\0';
2562
2563         /*
2564          * Finally, remove any trailing spaces
2565          */
2566         while (--dst > cpi->cpi_brandstr)
2567             if (*dst == ' ')
2568                 *dst = '\0';

```

```

2567             else
2568                 break;
2569         } else
2570             fabricate_brandstr(cpi);
2571     }
2572     cpi->cpi_pass = 3;
2573 }

2575 /*
2576  * This routine is called out of bind_hwcap() much later in the life
2577  * of the kernel (post_startup()). The job of this routine is to resolve
2578  * the hardware feature support and kernel support for those features into
2579  * what we're actually going to tell applications via the aux vector.
2580 */
2581 void
2582 cpuid_pass4(cpu_t *cpu, uint_t *hwcap_out)
2583 {
2584     struct cpuid_info *cpi;
2585     uint_t hwcap_flags = 0, hwcap_flags_2 = 0;

2587     if (cpu == NULL)
2588         cpu = CPU;
2589     cpi = cpu->cpu_m.mcpu_cpi;

2591     ASSERT(cpi->cpi_pass == 3);

2593     if (cpi->cpi_maxeax >= 1) {
2594         uint32_t *edx = &cpi->cpi_support[STD_EDX_FEATURES];
2595         uint32_t *ecx = &cpi->cpi_support[STD_ECX_FEATURES];

2596         *edx = CPI_FEATURES_EDX(cpi);
2597         *ecx = CPI_FEATURES_ECX(cpi);

2598         /*
2599          * [these require explicit kernel support]
2600          */
2601         if (!is_x86_feature(x86_featureset, X86FSET_SEP))
2602             *edx &= ~CPUID_INTC_EDX_SEP;

2603         if (!is_x86_feature(x86_featureset, X86FSET_SSE))
2604             *edx &= ~(CPUID_INTC_FXSR | CPUID_INTC_EDX_SSE);
2605         if (!is_x86_feature(x86_featureset, X86FSET_SSE2))
2606             *edx &= ~CPUID_INTC_EDX_SSE2;

2607         if (!is_x86_feature(x86_featureset, X86FSET_HTT))
2608             *edx &= ~CPUID_INTC_EDX_HTT;

2609         if (!is_x86_feature(x86_featureset, X86FSET_SSE3))
2610             *ecx &= ~CPUID_INTC_ECX_SSE3;

2611         if (!is_x86_feature(x86_featureset, X86FSET_SSSE3))
2612             *ecx &= ~CPUID_INTC_ECX_SSSE3;
2613         if (!is_x86_feature(x86_featureset, X86FSET_SSE4_1))
2614             *ecx &= ~CPUID_INTC_ECX_SSE4_1;
2615         if (!is_x86_feature(x86_featureset, X86FSET_SSE4_2))
2616             *ecx &= ~CPUID_INTC_ECX_SSE4_2;
2617         if (!is_x86_feature(x86_featureset, X86FSET_AES))
2618             *ecx &= ~CPUID_INTC_ECX_AES;
2619         if (!is_x86_feature(x86_featureset, X86FSET_PCLMULQDQ))
2620             *ecx &= ~CPUID_INTC_ECX_PCLMULQDQ;
2621         if (!is_x86_feature(x86_featureset, X86FSET_XSAVE |
2622             CPUID_INTC_ECX_OSXSAVE));
2623         if (!is_x86_feature(x86_featureset, X86FSET_AVX))
2624             *ecx &= ~CPUID_INTC_ECX_AVX;
2625         if (!is_x86_feature(x86_featureset, X86FSET_F16C))
2626

```

```

2633             *ecx &= ~CPUID_INTC_ECX_F16C;

2635         /*
2636          * [no explicit support required beyond x87 fp context]
2637          */
2638         if (!fpu_exists)
2639             *edx &= ~(CPUID_INTC_EDX_FPU | CPUID_INTC_EDX_MMX);

2641         /*
2642          * Now map the supported feature vector to things that we
2643          * think userland will care about.
2644          */
2645         if (*edx & CPUID_INTC_EDX_SEP)
2646             hwcap_flags |= AV_386_SEP;
2647         if (*edx & CPUID_INTC_EDX_SSE)
2648             hwcap_flags |= AV_386_FXSR | AV_386_SSE;
2649         if (*edx & CPUID_INTC_EDX_SSE2)
2650             hwcap_flags |= AV_386_SSE2;
2651         if (*ecx & CPUID_INTC_ECX_SSE3)
2652             hwcap_flags |= AV_386_SSE3;
2653         if (*ecx & CPUID_INTC_ECX_SSSE3)
2654             hwcap_flags |= AV_386_SSSE3;
2655         if (*ecx & CPUID_INTC_ECX_SSE4_1)
2656             hwcap_flags |= AV_386_SSE4_1;
2657         if (*ecx & CPUID_INTC_ECX_SSE4_2)
2658             hwcap_flags |= AV_386_SSE4_2;
2659         if (*ecx & CPUID_INTC_ECX_MOVBE)
2660             hwcap_flags |= AV_386_MOVBE;
2661         if (*ecx & CPUID_INTC_ECX_AES)
2662             hwcap_flags |= AV_386_AES;
2663         if (*ecx & CPUID_INTC_ECX_PCLMULQDQ)
2664             hwcap_flags |= AV_386_PCLMULQDQ;
2665         if ((*ecx & CPUID_INTC_ECX_XSAVE) &&
2666             (*ecx & CPUID_INTC_ECX_OSXSAVE)) {
2667             hwcap_flags |= AV_386_XSAVE;

2669             if (*ecx & CPUID_INTC_ECX_AVX) {
2670                 hwcap_flags |= AV_386_AVX;
2671                 if (*ecx & CPUID_INTC_ECX_F16C)
2672                     hwcap_flags_2 |= AV_386_2_F16C;
2673             }
2674         }
2675         if (*ecx & CPUID_INTC_ECX_VMX)
2676             hwcap_flags |= AV_386_VMX;
2677         if (*ecx & CPUID_INTC_ECX_POPCNT)
2678             hwcap_flags |= AV_386_POPCNT;
2679         if (*edx & CPUID_EDX_FPU)
2680             hwcap_flags |= AV_386_FPU;
2681         if (*edx & CPUID_EDX_MMX)
2682             hwcap_flags |= AV_386_MMX;

2684         if (*edx & CPUID_INTC_EDX_TSC)
2685             hwcap_flags |= AV_386_TSC;
2686         if (*edx & CPUID_INTC_EDX_CX8)
2687             hwcap_flags |= AV_386_CX8;
2688         if (*edx & CPUID_INTC_EDX_CMOV)
2689             hwcap_flags |= AV_386_CMOV;
2690         if (*ecx & CPUID_INTC_ECX_CX16)
2691             hwcap_flags |= AV_386_CX16;

2693         if (*ecx & CPUID_INTC_ECX_RDRAND)
2694             hwcap_flags_2 |= AV_386_2_RDRAND;
2695     }

2697     if (cpi->cpi_xmaxeax < 0x80000001)
2698         goto pass4_done;

```

```

2700     switch (cpi->cpi_vendor) {
2701         struct cpuid_regs cp;
2702         uint32_t *edx, *ecx;
2703
2704     case X86_VENDOR_Intel:
2705         /*
2706          * Seems like Intel duplicated what we necessary
2707          * here to make the initial crop of 64-bit OS's work.
2708          * Hopefully, those are the only "extended" bits
2709          * they'll add.
2710         */
2711     /*FALLTHROUGH*/
2712
2713     case X86_VENDOR_AMD:
2714         edx = &cpi->cpi_support[AMD_EDX_FEATURES];
2715         ecx = &cpi->cpi_support[AMD_ECX_FEATURES];
2716
2717         *edx = CPI_FEATURES_XTD_EDX(cpi);
2718         *ecx = CPI_FEATURES_XTD_ECX(cpi);
2719
2720         /*
2721          * [these features require explicit kernel support]
2722         */
2723         switch (cpi->cpi_vendor) {
2724             case X86_VENDOR_Intel:
2725                 if (!is_x86_feature(x86_featureset, X86FSET_TSCP))
2726                     *edx &= ~CPUID_AMD_EDX_TSCP;
2727                 break;
2728
2729             case X86_VENDOR_AMD:
2730                 if (!is_x86_feature(x86_featureset, X86FSET_TSCP))
2731                     *edx &= ~CPUID_AMD_EDX_TSCP;
2732                 if (!is_x86_feature(x86_featureset, X86FSET_SSE4A))
2733                     *ecx &= ~CPUID_AMD_ECX_SSE4A;
2734                 break;
2735
2736             default:
2737                 break;
2738         }
2739
2740         /*
2741          * [no explicit support required beyond
2742          * x87 fp context and exception handlers]
2743         */
2744         if (!fpu_exists)
2745             *edx &= ~(CPUID_AMD_EDX_MMXamd |
2746                      CPUID_AMD_EDX_3DNow | CPUID_AMD_EDX_3DNowx);
2747
2748         if (!is_x86_feature(x86_featureset, X86FSET_NX))
2749             *edx &= ~CPUID_AMD_EDX_NX;
2750 #if !defined(__amd64)
2751         *edx &= ~CPUID_AMD_EDX_LM;
2752 #endif
2753
2754         /*
2755          * Now map the supported feature vector to
2756          * things that we think userland will care about.
2757         */
2758 #if defined(__amd64)
2759         if (*edx & CPUID_AMD_EDX_SYSC)
2760             hwcap_flags |= AV_386_AMD_SYSC;
2761
2762         if (*edx & CPUID_AMD_EDX_MMXamd)
2763             hwcap_flags |= AV_386_AMD_MMX;
2764         if (*edx & CPUID_AMD_EDX_3DNow)
2765             hwcap_flags |= AV_386_AMD_3DNow;

```

```

2765         if (*edx & CPUID_AMD_EDX_3DNowx)
2766             hwcap_flags |= AV_386_AMD_3DNowx;
2767         if (*ecx & CPUID_AMD_ECX_SVM)
2768             hwcap_flags |= AV_386_AMD_SVM;
2769
2770         switch (cpi->cpi_vendor) {
2771             case X86_VENDOR_AMD:
2772                 if (*edx & CPUID_AMD_EDX_TSCP)
2773                     hwcap_flags |= AV_386_TSCP;
2774                 if (*ecx & CPUID_AMD_ECX_AHF64)
2775                     hwcap_flags |= AV_386_AHF;
2776                 if (*ecx & CPUID_AMD_ECX_SSE4A)
2777                     hwcap_flags |= AV_386_AMD_SSE4A;
2778                 if (*ecx & CPUID_AMD_ECX_LZCNT)
2779                     hwcap_flags |= AV_386_AMD_LZCNT;
2780                 break;
2781
2782             case X86_VENDOR_Intel:
2783                 if (*edx & CPUID_AMD_EDX_TSCP)
2784                     hwcap_flags |= AV_386_TSCP;
2785
2786                 /*
2787                  * Arrgh.
2788                  * Intel uses a different bit in the same word.
2789                 */
2790                 if (*ecx & CPUID_INTC_ECX_AHF64)
2791                     hwcap_flags |= AV_386_AHF;
2792                 break;
2793
2794             default:
2795                 break;
2796         }
2797         break;
2798
2799     case X86_VENDOR_TM:
2800         cp.cp_eax = 0x80860001;
2801         (void) __cpuid_insn(&cp);
2802         cpi->cpi_support[TM_EDX_FEATURES] = cp.cp_edx;
2803
2804     default:
2805         break;
2806     }
2807
2808 pass4_done:
2809     cpi->cpi_pass = 4;
2810     if (hwcap_out != NULL) {
2811         hwcap_out[0] = hwcap_flags;
2812         hwcap_out[1] = hwcap_flags_2;
2813     }
2814 }
2815
2816 /*
2817  * Simulate the cpuid instruction using the data we previously
2818  * captured about this CPU. We try our best to return the truth
2819  * about the hardware, independently of kernel support.
2820  */
2821
2822 uint32_t
2823 cpuid_insn(cpu_t *cpu, struct cpuid_regs *cp)
2824 {
2825     struct cpuid_info *cpi;
2826     struct cpuid_regs *xcp;
2827
2828     if (cpu == NULL)
2829         cpu = CPU;
2830     cpi = cpu->cpu_m.mcpu_cpi;

```

```

2832     ASSERT(cpuid_checkpass(cpu, 3));
2834
2835     /* CPUID data is cached in two separate places: cpi_std for standard
2836     * CPUID functions, and cpi_extd for extended CPUID functions.
2837     */
2838     if (cp->cp_eax <= cpi->cpi_maxeax && cp->cp_eax < NMAX_CPI_STD)
2839         xcp = &cpi->cpi_std[cp->cp_eax];
2840     else if (cp->cp_eax >= 0x80000000 && cp->cp_eax <= cpi->cpi_xmaxeax &&
2841             cp->cp_eax < 0x80000000 + NMAX_CPI_EXTD)
2842         xcp = &cpi->cpi_extd[cp->cp_eax - 0x80000000];
2843     else
2844         /*
2845          * The caller is asking for data from an input parameter which
2846          * the kernel has not cached. In this case we go fetch from
2847          * the hardware and return the data directly to the user.
2848          */
2849     return (_cpuid_insn(cp));
2851
2852     cp->cp_eax = xcp->cp_eax;
2853     cp->cp_ebx = xcp->cp_ebx;
2854     cp->cp_ecx = xcp->cp_ecx;
2855     cp->cp_edx = xcp->cp_edx;
2856     return (cp->cp_eax);
2857 }
2858 int
2859 cpuid_checkpass(cpu_t *cpu, int pass)
2860 {
2861     return (cpu != NULL && cpu->cpu_m.mcpu_cpi != NULL &&
2862             cpu->cpu_m.mcpu_cpi->cpi_pass >= pass);
2863 }
2864 int
2865 cpuid_getbrandstr(cpu_t *cpu, char *s, size_t n)
2866 {
2867     ASSERT(cpuid_checkpass(cpu, 3));
2868
2869     return (snprintf(s, n, "%s", cpu->cpu_m.mcpu_cpi->cpi_brandstr));
2870 }
2871
2872 int
2873 cpuid_is_cmt(cpu_t *cpu)
2874 {
2875     if (cpu == NULL)
2876         cpu = CPU;
2877
2878     ASSERT(cpuid_checkpass(cpu, 1));
2879
2880     return (cpu->cpu_m.mcpu_cpi->cpi_chipid >= 0);
2881 }
2882
2883 /*
2884  * AMD and Intel both implement the 64-bit variant of the syscall
2885  * instruction (syscall1q), so if there's -any- support for syscall,
2886  * cpuid currently says "yes, we support this".
2887  *
2888  * However, Intel decided to -not- implement the 32-bit variant of the
2889  * syscall instruction, so we provide a predicate to allow our caller
2890  * to test that subtlety here.
2891  *
2892  *
2893  * XXPV Currently, 32-bit syscall instructions don't work via the hypervisor,
2894  * even in the case where the hardware would in fact support it.
2895  */
2896 /*ARGSUSED*/

```

```

2897 int
2898 cpuid_syscall32_insn(cpu_t *cpu)
2899 {
2900     ASSERT(cpuid_checkpass((cpu == NULL ? CPU : cpu), 1));
2901
2902 #if !defined(__xpv)
2903     if (cpu == NULL)
2904         cpu = CPU;
2905
2906     /*CSTYLED*/
2907     {
2908         struct cpuid_info *cpi = cpu->cpu_m.mcpu_cpi;
2909
2910         if (cpi->cpi_vendor == X86_VENDOR_AMD &&
2911             cpi->cpi_xmaxeax >= 0x80000001 &&
2912             (CPI_FEATURES_XTD_EDX(cpi) & CPUID_AMD_EDX_SYSC))
2913             return (1);
2914     }
2915 #endif
2916     return (0);
2917 }
2918
2919 int
2920 cpuid_getidstr(cpu_t *cpu, char *s, size_t n)
2921 {
2922     struct cpuid_info *cpi = cpu->cpu_m.mcpu_cpi;
2923
2924     static const char fmt[] =
2925         "x86 (%s %X family %d model %d step %d clock %d MHz)";
2926     static const char fmt_ht[] =
2927         "x86 (chipid 0x%X %s %X family %d model %d step %d clock %d MHz)";
2928
2929     ASSERT(cpuid_checkpass(cpu, 1));
2930
2931     if (cpuid_is_cmt(cpu))
2932         return (snprintf(s, n, fmt_ht, cpi->cpi_chipid,
2933                         cpi->cpi_vendorstr, cpi->cpi_std[1].cp_eax,
2934                         cpi->cpi_family, cpi->cpi_model,
2935                         cpi->cpi_step, cpu->cpu_type_info.pi_clock));
2936     return (snprintf(s, n, fmt,
2937                     cpi->cpi_vendorstr, cpi->cpi_std[1].cp_eax,
2938                     cpi->cpi_family, cpi->cpi_model,
2939                     cpi->cpi_step, cpu->cpu_type_info.pi_clock));
2940 }
2941
2942 const char *
2943 cpuid_getvendorstr(cpu_t *cpu)
2944 {
2945     ASSERT(cpuid_checkpass(cpu, 1));
2946     return ((const char *)cpu->cpu_m.mcpu_cpi->cpi_vendorstr);
2947 }
2948
2949 uint_t
2950 cpuid_getvendor(cpu_t *cpu)
2951 {
2952     ASSERT(cpuid_checkpass(cpu, 1));
2953     return (cpu->cpu_m.mcpu_cpi->cpi_vendor);
2954 }
2955
2956 uint_t
2957 cpuid_getfamily(cpu_t *cpu)
2958 {
2959     ASSERT(cpuid_checkpass(cpu, 1));
2960     return (cpu->cpu_m.mcpu_cpi->cpi_family);
2961 }

```

```

2963 uint_t
2964 cpuid_getmodel(cpu_t *cpu)
2965 {
2966     ASSERT(cpuid_checkpass(cpu, 1));
2967     return (cpu->cpu_m.mcpu_cpi->cpi_model);
2968 }

2970 uint_t
2971 cpuid_get_ncpu_per_chip(cpu_t *cpu)
2972 {
2973     ASSERT(cpuid_checkpass(cpu, 1));
2974     return (cpu->cpu_m.mcpu_cpi->cpi_ncpu_per_chip);
2975 }

2977 uint_t
2978 cpuid_get_ncore_per_chip(cpu_t *cpu)
2979 {
2980     ASSERT(cpuid_checkpass(cpu, 1));
2981     return (cpu->cpu_m.mcpu_cpi->cpi_ncore_per_chip);
2982 }

2984 uint_t
2985 cpuid_get_ncpu_sharing_last_cache(cpu_t *cpu)
2986 {
2987     ASSERT(cpuid_checkpass(cpu, 2));
2988     return (cpu->cpu_m.mcpu_cpi->cpi_ncpu_shr_last_cache);
2989 }

2991 id_t
2992 cpuid_get_last_lvl_cacheid(cpu_t *cpu)
2993 {
2994     ASSERT(cpuid_checkpass(cpu, 2));
2995     return (cpu->cpu_m.mcpu_cpi->cpi_last_lvl_cacheid);
2996 }

2998 uint_t
2999 cpuid_getstep(cpu_t *cpu)
3000 {
3001     ASSERT(cpuid_checkpass(cpu, 1));
3002     return (cpu->cpu_m.mcpu_cpi->cpi_step);
3003 }

3005 uint_t
3006 cpuid_getsig(struct cpu *cpu)
3007 {
3008     ASSERT(cpuid_checkpass(cpu, 1));
3009     return (cpu->cpu_m.mcpu_cpi->cpi_std[1].cp_eax);
3010 }

3012 uint32_t
3013 cpuid_getchiprev(struct cpu *cpu)
3014 {
3015     ASSERT(cpuid_checkpass(cpu, 1));
3016     return (cpu->cpu_m.mcpu_cpi->cpi_chiprev);
3017 }

3019 const char *
3020 cpuid_getchiprevstr(struct cpu *cpu)
3021 {
3022     ASSERT(cpuid_checkpass(cpu, 1));
3023     return (cpu->cpu_m.mcpu_cpi->cpi_chiprevstr);
3024 }

3026 uint32_t
3027 cpuid_getsockettype(struct cpu *cpu)
3028 {

```

```

3029     ASSERT(cpuid_checkpass(cpu, 1));
3030     return (cpu->cpu_m.mcpu_cpi->cpi_socket);
3031 }

3033 const char *
3034 cpuid_getsocketstr(cpu_t *cpu)
3035 {
3036     static const char *socketstr = NULL;
3037     struct cpuid_info *cpi;
3038
3039     ASSERT(cpuid_checkpass(cpu, 1));
3040     cpi = cpu->cpu_m.mcpu_cpi;
3041
3042     /* Assume that socket types are the same across the system */
3043     if (socketstr == NULL)
3044         socketstr = _cpuid_sktstr(cpi->cpi_vendor, cpi->cpi_family,
3045                               cpi->cpi_model, cpi->cpi_step);
3046
3047     return (socketstr);
3048 }
3049

3051 int
3052 cpuid_get_chipid(cpu_t *cpu)
3053 {
3054     ASSERT(cpuid_checkpass(cpu, 1));
3055
3056     if (cpuid_is_cmt(cpu))
3057         return (cpu->cpu_m.mcpu_cpi->cpi_chipid);
3058     return (cpu->cpu_id);
3059 }

3061 id_t
3062 cpuid_get_coreid(cpu_t *cpu)
3063 {
3064     ASSERT(cpuid_checkpass(cpu, 1));
3065     return (cpu->cpu_m.mcpu_cpi->cpi_coreid);
3066 }

3068 int
3069 cpuid_get_pkgcoreid(cpu_t *cpu)
3070 {
3071     ASSERT(cpuid_checkpass(cpu, 1));
3072     return (cpu->cpu_m.mcpu_cpi->cpi_pkgcoreid);
3073 }

3075 int
3076 cpuid_get_clogid(cpu_t *cpu)
3077 {
3078     ASSERT(cpuid_checkpass(cpu, 1));
3079     return (cpu->cpu_m.mcpu_cpi->cpi_clogid);
3080 }

3082 int
3083 cpuid_get_cacheid(cpu_t *cpu)
3084 {
3085     ASSERT(cpuid_checkpass(cpu, 1));
3086     return (cpu->cpu_m.mcpu_cpi->cpi_last_lvl_cacheid);
3087 }

3089 uint_t
3090 cpuid_get_procnodeid(cpu_t *cpu)
3091 {
3092     ASSERT(cpuid_checkpass(cpu, 1));
3093     return (cpu->cpu_m.mcpu_cpi->cpi_procnodeid);
3094 }

```

```

3096 uint_t
3097 cpuid_get_procnodes_per_pkg(cpu_t *cpu)
3098 {
3099     ASSERT(cpuid_checkpass(cpu, 1));
3100     return (cpu->cpu_m.mcpu_cpi->cpi_procnodes_per_pkg);
3101 }

3103 uint_t
3104 cpuid_get_computunitid(cpu_t *cpu)
3105 {
3106     ASSERT(cpuid_checkpass(cpu, 1));
3107     return (cpu->cpu_m.mcpu_cpi->cpi_computunitid);
3108 }

3110 uint_t
3111 cpuid_get_cores_per_computunit(cpu_t *cpu)
3112 {
3113     ASSERT(cpuid_checkpass(cpu, 1));
3114     return (cpu->cpu_m.mcpu_cpi->cpi_cores_per_computunit);
3115 }

3117 /*ARGSUSED*/
3118 int
3119 cpuid_have_cr8access(cpu_t *cpu)
3120 {
3121 #if defined(__amd64)
3122     return (1);
3123 #else
3124     struct cpuid_info *cpi;
3125
3126     ASSERT(cpu != NULL);
3127     cpi = cpu->cpu_m.mcpu_cpi;
3128     if (cpi->cpi_vendor == X86_VENDOR_AMD && cpi->cpi_maxeax >= 1 &&
3129         (CPI_FEATURES_XTD_ECX(cpi) & CPUID_AMD_ECX_CR8D) != 0)
3130         return (1);
3131     return (0);
3132 #endif
3133 }

3135 uint32_t
3136 cpuid_get_apicid(cpu_t *cpu)
3137 {
3138     ASSERT(cpuid_checkpass(cpu, 1));
3139     if (cpu->cpu_m.mcpu_cpi->cpi_maxeax < 1) {
3140         return (UINT32_MAX);
3141     } else {
3142         return (cpu->cpu_m.mcpu_cpi->cpi_apicid);
3143     }
3144 }

3146 void
3147 cpuid_get_addrsize(cpu_t *cpu, uint_t *pabits, uint_t *vabits)
3148 {
3149     struct cpuid_info *cpi;
3150
3151     if (cpu == NULL)
3152         cpu = CPU;
3153     cpi = cpu->cpu_m.mcpu_cpi;
3154
3155     ASSERT(cpuid_checkpass(cpu, 1));
3156
3157     if (pabits)
3158         *pabits = cpi->cpi_pabits;
3159     if (vabits)
3160         *vabits = cpi->cpi_vabits;

```

```

3161 }

3163 /*
3164  * Returns the number of data TLB entries for a corresponding
3165  * pagesize. If it can't be computed, or isn't known, the
3166  * routine returns zero. If you ask about an architecturally
3167  * impossible pagesize, the routine will panic (so that the
3168  * host implementor knows that things are inconsistent.)
3169 */
3170 uint_t
3171 cpuid_get_dtib_nent(cpu_t *cpu, size_t pagesize)
3172 {
3173     struct cpuid_info *cpi;
3174     uint_t dtib_nent = 0;
3175
3176     if (cpu == NULL)
3177         cpu = CPU;
3178     cpi = cpu->cpu_m.mcpu_cpi;
3179
3180     ASSERT(cpuid_checkpass(cpu, 1));
3181
3182     /*
3183      * Check the L2 TLB info
3184      */
3185     if (cpi->cpi_xmaxeax >= 0x80000006) {
3186         struct cpuid_regs *cp = &cpi->cpi_extd[6];
3187
3188         switch (pagesize) {
3189             case 4 * 1024:
3190                 /*
3191                  * All zero in the top 16 bits of the register
3192                  * indicates a unified TLB. Size is in low 16 bits.
3193                  */
3194                 if ((cp->cp_ebx & 0xffff0000) == 0)
3195                     dtib_nent = cp->cp_ebx & 0x0000ffff;
3196                 else
3197                     dtib_nent = BITX(cp->cp_ebx, 27, 16);
3198                 break;
3199
3200             case 2 * 1024 * 1024:
3201                 if ((cp->cp_eax & 0xffff0000) == 0)
3202                     dtib_nent = cp->cp_eax & 0x0000ffff;
3203                 else
3204                     dtib_nent = BITX(cp->cp_eax, 27, 16);
3205                 break;
3206
3207             default:
3208                 panic("unknown L2 pagesize");
3209                 /*NOTREACHED*/
3210         }
3211     }
3212
3213     if (dtib_nent != 0)
3214         return (dtib_nent);
3215
3216     /*
3217      * No L2 TLB support for this size, try L1.
3218      */
3219     if (cpi->cpi_xmaxeax >= 0x80000005) {
3220         struct cpuid_regs *cp = &cpi->cpi_extd[5];
3221
3222         switch (pagesize) {
3223             case 4 * 1024:
3224                 dtib_nent = BITX(cp->cp_ebx, 23, 16);
3225                 break;
3226

```

```

3227         case 2 * 1024 * 1024:
3228             dtlb_nent = BITX(cp->cp_eax, 23, 16);
3229             break;
3230         default:
3231             panic("unknown L1 d-TLB pagesize");
3232             /*NOTREACHED*/
3233     }
3234 }
3235
3236     return (dtlb_nent);
3237 }
3238 */
3239 /* Return 0 if the erratum is not present or not applicable, positive
3240 * if it is, and negative if the status of the erratum is unknown.
3241 *
3242 * See "Revision Guide for AMD Athlon(tm) 64 and AMD Opteron(tm)
3243 * Processors" #25759, Rev 3.57, August 2005
3244 */
3245
3246 int
3247 cpuid_opteron_erratum(cpu_t *cpu, uint_t erratum)
3248 {
3249     struct cpuid_info *cpi = cpu->cpu_m.mcpu_cpi;
3250     uint_t eax;
3251
3252     /*
3253     * Bail out if this CPU isn't an AMD CPU, or if it's
3254     * a legacy (32-bit) AMD CPU.
3255     */
3256     if (cpi->cpi_vendor != X86_VENDOR_AMD ||
3257         cpi->cpi_family == 4 || cpi->cpi_family == 5 ||
3258         cpi->cpi_family == 6)
3259
3260         return (0);
3261
3262     eax = cpi->cpi_std[1].cp_eax;
3263
3264 #define SH_B0(eax)      (eax == 0xf40 || eax == 0xf50)
3265 #define SH_B3(eax)      (eax == 0xf51)
3266 #define B(eax)          (SH_B0(eax) || SH_B3(eax))
3267
3268 #define SH_C0(eax)      (eax == 0xf48 || eax == 0xf58)
3269
3270 #define SH(CG)(eax)    (eax == 0xf4a || eax == 0xf5a || eax == 0xf7a)
3271 #define DH(CG)(eax)    (eax == 0xfc0 || eax == 0xfe0 || eax == 0xff0)
3272 #define CH(CG)(eax)    (eax == 0xf82 || eax == 0xfb2)
3273 #define CG(eax)        (SH(CG)(eax) || DH(CG)(eax) || CH(CG)(eax))
3274
3275 #define SH_D0(eax)      (eax == 0x10f40 || eax == 0x10f50 || eax == 0x10f70)
3276 #define DH_D0(eax)      (eax == 0x10fc0 || eax == 0x10ff0)
3277 #define CH_D0(eax)      (eax == 0x10f80 || eax == 0x10fb0)
3278 #define D0(eax)         (SH_D0(eax) || DH_D0(eax) || CH_D0(eax))
3279
3280 #define SH_E0(eax)      (eax == 0x20f50 || eax == 0x20f40 || eax == 0x20f70)
3281 #define JH_E1(eax)      (eax == 0x20f10) /* JH8_E0 had 0x20f30 */
3282 #define DH_E3(eax)      (eax == 0x20fc0 || eax == 0x20ff0)
3283 #define SH_E4(eax)      (eax == 0x20f51 || eax == 0x20f71)
3284 #define BH_E4(eax)      (eax == 0x20fb1)
3285 #define SH_E5(eax)      (eax == 0x20f42)
3286 #define DH_E6(eax)      (eax == 0x20ff2 || eax == 0x20fc2)
3287 #define JH_E6(eax)      (eax == 0x20f12 || eax == 0x20f32)
3288 #define EX(eax)         (SH_E0(eax) || JH_E1(eax) || DH_E3(eax) || \
3289                         SH_E4(eax) || BH_E4(eax) || SH_E5(eax) || \
3290                         DH_E6(eax) || JH_E6(eax))
3291
3292 #define DR_AX(eax)      (eax == 0x100f00 || eax == 0x100f01 || eax == 0x100f02)

```

```

3293 #define DR_B0(eax)      (eax == 0x100f20)
3294 #define DR_B1(eax)      (eax == 0x100f21)
3295 #define DR_BA(eax)      (eax == 0x100f2a)
3296 #define DR_B2(eax)      (eax == 0x100f22)
3297 #define DR_B3(eax)      (eax == 0x100f23)
3298 #define RB_C0(eax)      (eax == 0x100f40)
3299
3300 switch (erratum) {
3301     case 1:
3302         return (cpi->cpi_family < 0x10);
3303         /* what does the asterisk mean? */
3304         return (B(eax) || SH_C0(eax) || CG(eax));
3305     case 51:
3306         return (B(eax));
3307     case 57:
3308         return (cpi->cpi_family <= 0x11);
3309     case 58:
3310         return (B(eax));
3311     case 60:
3312         return (cpi->cpi_family <= 0x11);
3313     case 61:
3314     case 62:
3315     case 63:
3316     case 64:
3317     case 65:
3318     case 66:
3319     case 68:
3320     case 69:
3321     case 70:
3322     case 71:
3323         return (B(eax));
3324     case 72:
3325         return (SH_B0(eax));
3326     case 74:
3327         return (B(eax));
3328     case 75:
3329         return (cpi->cpi_family < 0x10);
3330     case 76:
3331         return (B(eax));
3332     case 77:
3333         return (cpi->cpi_family <= 0x11);
3334     case 78:
3335         return (B(eax) || SH_C0(eax));
3336     case 79:
3337         return (B(eax) || SH_C0(eax) || CG(eax) || D0(eax) || EX(eax));
3338     case 80:
3339     case 81:
3340     case 82:
3341         return (B(eax));
3342     case 83:
3343         return (B(eax) || SH_C0(eax) || CG(eax));
3344     case 85:
3345         return (cpi->cpi_family < 0x10);
3346     case 86:
3347         return (SH_C0(eax) || CG(eax));
3348     case 88:
3349     #if !defined(__amd64)
3350         return (0);
3351     #else
3352         return (B(eax) || SH_C0(eax));
3353     #endif
3354     case 89:
3355         return (cpi->cpi_family < 0x10);
3356     case 90:
3357         return (B(eax) || SH_C0(eax) || CG(eax));
3358     case 91:

```

```

3359     case 92:
3360         return (B(eax) || SH_C0(eax));
3361     case 93:
3362         return (SH_C0(eax));
3363     case 94:
3364         return (B(eax) || SH_C0(eax) || CG(eax));
3365     case 95:
3366 #if !defined(__amd64)
3367         return (0);
3368 #else
3369         return (B(eax) || SH_C0(eax));
3370 #endif
3371     case 96:
3372         return (B(eax) || SH_C0(eax) || CG(eax));
3373     case 97:
3374     case 98:
3375         return (SH_C0(eax) || CG(eax));
3376     case 99:
3377         return (B(eax) || SH_C0(eax) || CG(eax) || D0(eax));
3378     case 100:
3379         return (B(eax) || SH_C0(eax));
3380     case 101:
3381     case 103:
3382         return (B(eax) || SH_C0(eax) || CG(eax) || D0(eax));
3383     case 104:
3384         return (SH_C0(eax) || CG(eax) || D0(eax));
3385     case 105:
3386     case 106:
3387     case 107:
3388         return (B(eax) || SH_C0(eax) || CG(eax) || D0(eax));
3389     case 108:
3390         return (DH(CG(eax)));
3391     case 109:
3392         return (SH_C0(eax) || CG(eax) || D0(eax));
3393     case 110:
3394         return (D0(eax) || EX(eax));
3395     case 111:
3396         return (CG(eax));
3397     case 112:
3398         return (B(eax) || SH_C0(eax) || CG(eax) || D0(eax) || EX(eax));
3399     case 113:
3400         return (eax == 0x20fc0);
3401     case 114:
3402         return (SH_E0(eax) || JH_E1(eax) || DH_E3(eax));
3403     case 115:
3404         return (SH_E0(eax) || JH_E1(eax));
3405     case 116:
3406         return (SH_E0(eax) || JH_E1(eax) || DH_E3(eax));
3407     case 117:
3408         return (B(eax) || SH_C0(eax) || CG(eax) || D0(eax));
3409     case 118:
3410         return (SH_E0(eax) || JH_E1(eax) || SH_E4(eax) || BH_E4(eax) ||
3411                 JH_E6(eax));
3412     case 121:
3413         return (B(eax) || SH_C0(eax) || CG(eax) || D0(eax) || EX(eax));
3414     case 122:
3415         return (cpi->cpi_family < 0x10 || cpi->cpi_family == 0x11);
3416     case 123:
3417         return (JH_E1(eax) || BH_E4(eax) || JH_E6(eax));
3418     case 131:
3419         return (cpi->cpi_family < 0x10);
3420     case 6336786:
3421         /*
3422          * Test for AdvPowerMgmtInfo.TscPStateInvariant
3423          * if this is a K8 family or newer processor
3424          */

```

```

3425             if (CPI_FAMILY(cpi) == 0xf) {
3426                 struct cpuid_regs regs;
3427                 regs.cp_eax = 0x80000007;
3428                 (void) __cpuid_insn(&regs);
3429                 return (!(regs.cp_edx & 0x100));
3430             }
3431             return (0);
3432         case 6323525:
3433             return (((((eax > 12) & 0xff00) + (eax & 0xf00)) |
3434                     ((eax > 4) & 0xf) | ((eax > 12) & 0xf00)) < 0xf40);
3435
3436         case 6671130:
3437             /*
3438              * check for processors (pre-Shanghai) that do not provide
3439              * optimal management of lgb ptes in its tlb.
3440             */
3441             return (cpi->cpi_family == 0x10 && cpi->cpi_model < 4);
3442
3443         case 298:
3444             return (DR_AX(eax) || DR_B0(eax) || DR_B1(eax) || DR_BA(eax) ||
3445                     DR_B2(eax) || RB_C0(eax));
3446
3447         case 721:
3448 #if defined(__amd64)
3449             return (cpi->cpi_family == 0x10 || cpi->cpi_family == 0x12);
3450 #else
3451             return (0);
3452 #endif
3453
3454         default:
3455             return (-1);
3456     }
3457
3458 }
3459
3460 /*
3461  * Determine if specified erratum is present via OSVW (OS Visible Workaround).
3462  * Return 1 if erratum is present, 0 if not present and -1 if indeterminate.
3463  */
3464 int
3465 osvw_opteron_erratum(cpu_t *cpu, uint_t erratum)
3466 {
3467     struct cpuid_info      *cpi;
3468     uint_t                  oswwid;
3469     static int               oswffeature = -1;
3470     uint64_t                osvwlength;
3471
3472     cpi = cpu->cpu_m.mcpu_cpi;
3473
3474     /* confirm OSVW supported */
3475     if (oswffeature == -1) {
3476         oswffeature = cpi->cpi_extd[1].cp_ecx & CPUID_AMD_ECX_OSVW;
3477     } else {
3478         /* assert that oswv feature setting is consistent on all cpus */
3479         ASSERT(oswffeature ==
3480               (cpi->cpi_extd[1].cp_ecx & CPUID_AMD_ECX_OSVW));
3481     }
3482     if (!oswffeature)
3483         return (-1);
3484
3485     osvwlength = rdmsr(MSR_AMD_OSVW_ID_LEN) & OSVW_ID_LEN_MASK;
3486
3487     switch (erratum) {
3488     case 298: /* oswwid is 0 */
3489         oswwid = 0;
3490     }

```

```

3491     if (osvwlength <= (uint64_t)osvwid) {
3492         /* oswid 0 is unknown */
3493         return (-1);
3494     }
3495
3496     /*
3497      * Check the OSVW STATUS MSR to determine the state
3498      * of the erratum where:
3499      * 0 - fixed by HW
3500      * 1 - BIOS has applied the workaround when BIOS
3501      * workaround is available. (Or for other errata,
3502      * OS workaround is required.)
3503      * For a value of 1, caller will confirm that the
3504      * erratum 298 workaround has indeed been applied by BIOS.
3505
3506      * A 1 may be set in cpus that have a HW fix
3507      * in a mixed cpu system. Regarding erratum 298:
3508      * In a multiprocessor platform, the workaround above
3509      * should be applied to all processors regardless of
3510      * silicon revision when an affected processor is
3511      * present.
3512  */
3514
3515     return (rdmsr(MSR_AMD_OSVW_STATUS +
3516                 (osvwid / OSVW_ID_CNT_PER_MSR)) &
3517                 (1ULL << (osvwid % OSVW_ID_CNT_PER_MSR)));
3518
3519     default:
3520         return (-1);
3521 }
3522 static const char assoc_str[] = "associativity";
3523 static const char line_str[] = "line-size";
3524 static const char size_str[] = "size";
3525
3526 static void
3527 add_cache_prop(dev_info_t *devi, const char *label, const char *type,
3528                 uint32_t val)
3529 {
3530     char buf[128];
3531
3532     /*
3533      * ndi_prop_update_int() is used because it is desirable for
3534      * DDI_PROP_HW_DEF and DDI_PROP_DONTSLLEEP to be set.
3535      */
3536
3537     if (snprintf(buf, sizeof (buf), "%s-%s", label, type) < sizeof (buf))
3538         (void) ndi_prop_update_int(DDI_DEV_T_NONE, devi, buf, val);
3539 }
3540
3541 /*
3542  * Intel-style cache/tlb description
3543  *
3544  * Standard cpuid level 2 gives a randomly ordered
3545  * selection of tags that index into a table that describes
3546  * cache and tlb properties.
3547 */
3548
3549 static const char l1_icache_str[] = "l1-icache";
3550 static const char l1_dcache_str[] = "l1-dcache";
3551 static const char l2_cache_str[] = "l2-cache";
3552 static const char l3_cache_str[] = "l3-cache";
3553 static const char itlb4k_str[] = "itlb-4K";
3554 static const char dtlb4k_str[] = "dtlb-4K";
3555 static const char itlb2M_str[] = "itlb-2M";
3556 static const char itlb4M_str[] = "itlb-4M";

```

```

3557 static const char dtlb4M_str[] = "dtlb-4M";
3558 static const char dtlb24_str[] = "dtlb0-2M-4M";
3559 static const char itlb424_str[] = "itlb-4K-2M-4M";
3560 static const char itlb24_str[] = "itlb-2M-4M";
3561 static const char dtlb44_str[] = "dtlb-4K-4M";
3562 static const char s11_dcache_str[] = "sectored-l1-dcache";
3563 static const char s12_cache_str[] = "sectored-l2-cache";
3564 static const char itrace_str[] = "itrace-cache";
3565 static const char s13_cache_str[] = "sectored-l3-cache";
3566 static const char sh_12_tlb4k_str[] = "shared-l2-tlb-4k";
3567
3568 static const struct cachetab {
3569     uint8_t          ct_code;
3570     uint8_t          ct_assoc;
3571     uint16_t         ct_line_size;
3572     size_t           ct_size;
3573     const char       *ct_label;
3574 } intel_ctab[] = {
3575     /*
3576      * maintain descending order!
3577      *
3578      * Codes ignored - Reason
3579      * -----
3580      * 40H - intel_cpuid_4_cache_info() disambiguates 12/13 cache
3581      * f0H/f1H - Currently we do not interpret prefetch size by design
3582      */
3583     { 0xe4, 16, 64, 8*1024*1024, l3_cache_str },
3584     { 0xe3, 16, 64, 4*1024*1024, l3_cache_str },
3585     { 0xe2, 16, 64, 2*1024*1024, l3_cache_str },
3586     { 0xde, 12, 64, 6*1024*1024, l3_cache_str },
3587     { 0xdd, 12, 64, 3*1024*1024, l3_cache_str },
3588     { 0xdc, 12, 64, ((1*1024*1024)+(512*1024)), l3_cache_str },
3589     { 0xd8, 8, 64, 4*1024*1024, l3_cache_str },
3590     { 0xd7, 8, 64, 2*1024*1024, l3_cache_str },
3591     { 0xd6, 8, 64, 1*1024*1024, l3_cache_str },
3592     { 0xd2, 4, 64, 2*1024*1024, l3_cache_str },
3593     { 0xd1, 4, 64, 1*1024*1024, l3_cache_str },
3594     { 0xd0, 4, 64, 512*1024, l3_cache_str },
3595     { 0xca, 4, 0, 512, sh_12_tlb4k_str },
3596     { 0xc0, 4, 0, 8, dtlb44_str },
3597     { 0xba, 4, 0, 64, dtlb4k_str },
3598     { 0xb4, 4, 0, 256, dtlb4k_str },
3599     { 0xb3, 4, 0, 128, dtlb4k_str },
3600     { 0xb2, 4, 0, 64, itlb4k_str },
3601     { 0xb0, 4, 0, 128, itlb4k_str },
3602     { 0x87, 8, 64, 1024*1024, l2_cache_str },
3603     { 0x86, 4, 64, 512*1024, l2_cache_str },
3604     { 0x85, 8, 32, 2*1024*1024, l2_cache_str },
3605     { 0x84, 8, 32, 1024*1024, l2_cache_str },
3606     { 0x83, 8, 32, 512*1024, l2_cache_str },
3607     { 0x82, 8, 32, 256*1024, l2_cache_str },
3608     { 0x80, 8, 64, 512*1024, l2_cache_str },
3609     { 0x7f, 2, 64, 512*1024, l2_cache_str },
3610     { 0x7d, 8, 64, 2*1024*1024, s12_cache_str },
3611     { 0x7c, 8, 64, 1024*1024, s12_cache_str },
3612     { 0x7b, 8, 64, 512*1024, s12_cache_str },
3613     { 0x7a, 8, 64, 256*1024, s12_cache_str },
3614     { 0x79, 8, 64, 128*1024, s12_cache_str },
3615     { 0x78, 8, 64, 1024*1024, l2_cache_str },
3616     { 0x73, 8, 0, 64*1024, itrace_str },
3617     { 0x72, 8, 0, 32*1024, itrace_str },
3618     { 0x71, 8, 0, 16*1024, itrace_str },
3619     { 0x70, 8, 0, 12*1024, itrace_str },
3620     { 0x68, 4, 64, 32*1024, s11_dcache_str },
3621     { 0x67, 4, 64, 16*1024, s11_dcache_str },
3622     { 0x66, 4, 64, 8*1024, s11_dcache_str },

```

```

3623     { 0x60, 8, 64, 16*1024, s11_dcache_str},
3624     { 0x5d, 0, 0, 256, dtlb44_str},
3625     { 0x5c, 0, 0, 128, dtlb44_str},
3626     { 0x5b, 0, 0, 64, dtlb44_str},
3627     { 0x5a, 4, 0, 32, dtlb24_str},
3628     { 0x59, 0, 0, 16, dtlb4k_str},
3629     { 0x57, 4, 0, 16, dtlb4k_str},
3630     { 0x56, 4, 0, 16, dtlb4M_str},
3631     { 0x55, 0, 0, 7, itlb24_str},
3632     { 0x52, 0, 0, 256, itlb424_str},
3633     { 0x51, 0, 0, 128, itlb424_str},
3634     { 0x50, 0, 0, 64, itlb424_str},
3635     { 0x4f, 0, 0, 32, itlb4k_str},
3636     { 0x4e, 24, 64, 6*1024*1024, l2_cache_str},
3637     { 0x4d, 16, 64, 16*1024*1024, l3_cache_str},
3638     { 0x4c, 12, 64, 12*1024*1024, l3_cache_str},
3639     { 0x4b, 16, 64, 8*1024*1024, l3_cache_str},
3640     { 0x4a, 12, 64, 6*1024*1024, l3_cache_str},
3641     { 0x49, 16, 64, 4*1024*1024, l3_cache_str},
3642     { 0x48, 12, 64, 3*1024*1024, l2_cache_str},
3643     { 0x47, 8, 64, 8*1024*1024, l3_cache_str},
3644     { 0x46, 4, 64, 4*1024*1024, l3_cache_str},
3645     { 0x45, 4, 32, 2*1024*1024, l2_cache_str},
3646     { 0x44, 4, 32, 1024*1024, l2_cache_str},
3647     { 0x43, 4, 32, 512*1024, l2_cache_str},
3648     { 0x42, 4, 32, 256*1024, l2_cache_str},
3649     { 0x41, 4, 32, 128*1024, l2_cache_str},
3650     { 0x3e, 4, 64, 512*1024, s12_cache_str},
3651     { 0x3d, 6, 64, 384*1024, s12_cache_str},
3652     { 0x3c, 4, 64, 256*1024, s12_cache_str},
3653     { 0x3b, 2, 64, 128*1024, s12_cache_str},
3654     { 0x3a, 6, 64, 192*1024, s12_cache_str},
3655     { 0x39, 4, 64, 128*1024, s12_cache_str},
3656     { 0x30, 8, 64, 32*1024, l1_icache_str},
3657     { 0x2c, 8, 64, 32*1024, l1_dcache_str},
3658     { 0x29, 8, 64, 4096*1024, s13_cache_str},
3659     { 0x25, 8, 64, 2048*1024, s13_cache_str},
3660     { 0x23, 8, 64, 1024*1024, s13_cache_str},
3661     { 0x22, 4, 64, 512*1024, s13_cache_str},
3662     { 0xe, 6, 64, 24*1024, l1_dcache_str},
3663     { 0xd, 4, 32, 16*1024, l1_dcache_str},
3664     { 0xc, 4, 32, 16*1024, l1_dcache_str},
3665     { 0xb, 4, 0, 4, itlb4M_str},
3666     { 0xa, 2, 32, 8*1024, l1_dcache_str},
3667     { 0x8, 4, 32, 16*1024, l1_icache_str},
3668     { 0x6, 4, 32, 8*1024, l1_icache_str},
3669     { 0x5, 4, 0, 32, dtlb4M_str},
3670     { 0x4, 4, 0, 8, dtlb4M_str},
3671     { 0x3, 4, 0, 64, dtlb4k_str},
3672     { 0x2, 4, 0, 2, itlb4M_str},
3673     { 0x1, 4, 0, 32, itlb4k_str},
3674     { 0 }
3675 };

3677 static const struct cachetab cyrix_ctab[] = {
3678     { 0x70, 4, 0, 32, "tlb-4K" },
3679     { 0x80, 4, 16, 16*1024, "l1-cache" },
3680     { 0 }
3681 };

3683 /*
3684 * Search a cache table for a matching entry
3685 */
3686 static const struct cachetab *
3687 find_cachent(const struct cachetab *ct, uint_t code)
3688 {

```

```

3689     if (code != 0) {
3690         for (; ct->ct_code != 0; ct++)
3691             if (ct->ct_code <= code)
3692                 break;
3693         if (ct->ct_code == code)
3694             return (ct);
3695     }
3696     return (NULL);
3697 }

3698 /*
3699 * Populate cachetab entry with L2 or L3 cache-information using
3700 * cpuid function 4. This function is called from intel_walk_cacheinfo()
3701 * when descriptor 0x49 is encountered. It returns 0 if no such cache
3702 * information is found.
3703 */
3704 static int
3705 intel_cpuid_4_cache_info(struct cachetab *ct, struct cpuid_info *cpi)
3706 {
3707     uint32_t level, i;
3708     int ret = 0;

3709     for (i = 0; i < cpi->cpi_std_4_size; i++) {
3710         level = CPI_CACHE_LVL(cpi->cpi_std_4[i]);

3711         if (level == 2 || level == 3) {
3712             ct->ct_assoc = CPI_CACHE_WAYS(cpi->cpi_std_4[i]) + 1;
3713             ct->ct_line_size =
3714                 CPI_CACHE_COH_LN_SZ(cpi->cpi_std_4[i]) + 1;
3715             ct->ct_size = ct->ct_assoc *
3716                 (CPI_CACHE_PARTS(cpi->cpi_std_4[i]) + 1) *
3717                 ct->ct_line_size *
3718                 (cpi->cpi_std_4[i]->cp_ecx + 1);

3719             if (level == 2) {
3720                 ct->ct_label = l2_cache_str;
3721             } else if (level == 3) {
3722                 ct->ct_label = l3_cache_str;
3723             }
3724             ret = 1;
3725         }
3726     }
3727 }

3728     return (ret);
3729 }

3730 }

3731 }

3732 }

3733 }

3734 }

3735 /*
3736 * Walk the cacheinfo descriptor, applying 'func' to every valid element
3737 * The walk is terminated if the walker returns non-zero.
3738 */
3739 static void
3740 intel_walk_cacheinfo(struct cpuid_info *cpi,
3741     void *arg, int (*func)(void *, const struct cachetab *))
3742 {
3743     const struct cachetab *ct;
3744     struct cachetab des_49_ct, des_b1_ct;
3745     uint8_t *dp;
3746     int i;

3747     if ((dp = cpi->cpi_cacheinfo) == NULL)
3748         return;
3749     for (i = 0; i < cpi->cpi_ncache; i++, dp++) {
3750         /*
3751         * For overloaded descriptor 0x49 we use cpuid function 4
3752         * if supported by the current processor, to create
3753         * cache information.
3754     }

```

```

3755     * For overloaded descriptor 0xb1 we use X86_PAE flag
3756     * to disambiguate the cache information.
3757     */
3758     if (*dp == 0x49 && cpi->cpi_maxeax >= 0x4 &&
3759         intel_cpuid_4_cache_info(&des_49_ct, cpi) == 1) {
3760         ct = &des_49_ct;
3761     } else if (*dp == 0xb1) {
3762         des_b1_ct.ct_code = 0xb1;
3763         des_b1_ct.ct_assoc = 4;
3764         des_b1_ct.ct_line_size = 0;
3765         if (is_x86_feature(x86_featureset, X86FSET_PAE)) {
3766             des_b1_ct.ct_size = 8;
3767             des_b1_ct.ct_label = itlb2M_str;
3768         } else {
3769             des_b1_ct.ct_size = 4;
3770             des_b1_ct.ct_label = itlb4M_str;
3771         }
3772         ct = &des_b1_ct;
3773     } else {
3774         if ((ct = find_cacheent(intel_ctab, *dp)) == NULL) {
3775             continue;
3776         }
3777     }
3778     if (func(arg, ct) != 0) {
3779         break;
3780     }
3781 }
3782 }

3783 }

3784 /*

3785 * (Like the Intel one, except for Cyrix CPUs)
3786 */
3787 static void
3788 cyrix_walk_cacheinfo(struct cpuid_info *cpi,
3789     void *arg, int (*func)(void *, const struct cachetab *))
3790 {
3791     const struct cachetab *ct;
3792     uint8_t *dp;
3793     int i;

3794     if ((dp = cpi->cpi_cacheinfo) == NULL)
3795         return;
3796     for (i = 0; i < cpi->cpi_ncache; i++, dp++) {
3797         /*
3798             * Search Cyrix-specific descriptor table first ..
3799             */
3800         if ((ct = find_cacheent(cyrix_ctab, *dp)) != NULL) {
3801             if (func(arg, ct) != 0)
3802                 break;
3803             continue;
3804         }
3805         /*
3806             * .. else fall back to the Intel one
3807             */
3808         if ((ct = find_cacheent(intel_ctab, *dp)) != NULL) {
3809             if (func(arg, ct) != 0)
3810                 break;
3811             continue;
3812         }
3813     }
3814 }

3815 }

3816 }

3817 */

3818 /* A cacheinfo walker that adds associativity, line-size, and size properties
3819 * to the devinfo node it is passed as an argument.
3820 */

```

```

3821 */
3822 static int
3823 add_cacheent_props(void *arg, const struct cachetab *ct)
3824 {
3825     dev_info_t *devi = arg;
3826
3827     add_cache_prop(devi, ct->ct_label, assoc_str, ct->ct_assoc);
3828     if (ct->ct_line_size != 0)
3829         add_cache_prop(devi, ct->ct_label, line_str,
3830                         ct->ct_line_size);
3831     add_cache_prop(devi, ct->ct_label, size_str, ct->ct_size);
3832     return (0);
3833 }

3834 static const char fully_assoc[] = "fully-associative?";

3835 /*
3836  * AMD style cache/tlb description
3837  *
3838  * Extended functions 5 and 6 directly describe properties of
3839  * tlbs and various cache levels.
3840  */
3841 static void
3842 add_amd_assoc(dev_info_t *devi, const char *label, uint_t assoc)
3843 {
3844     switch (assoc) {
3845     case 0: /* reserved; ignore */
3846         break;
3847     default:
3848         add_cache_prop(devi, label, assoc_str, assoc);
3849         break;
3850     case 0xff:
3851         add_cache_prop(devi, label, fully_assoc, 1);
3852         break;
3853     }
3854 }

3855 static void
3856 add_amd_tlb(dev_info_t *devi, const char *label, uint_t assoc, uint_t size)
3857 {
3858     if (size == 0)
3859         return;
3860     add_cache_prop(devi, label, size_str, size);
3861     add_amd_assoc(devi, label, assoc);
3862 }

3863 static void
3864 add_amd_cache(dev_info_t *devi, const char *label,
3865                 uint_t size, uint_t assoc, uint_t lines_per_tag, uint_t line_size)
3866 {
3867     if (size == 0 || line_size == 0)
3868         return;
3869     add_amd_assoc(devi, label, assoc);
3870     /*
3871         * Most AMD parts have a sectored cache. Multiple cache lines are
3872         * associated with each tag. A sector consists of all cache lines
3873         * associated with a tag. For example, the AMD K6-III has a sector
3874         * size of 2 cache lines per tag.
3875         */
3876     if (lines_per_tag != 0)
3877         add_cache_prop(devi, label, "lines-per-tag", lines_per_tag);
3878     add_cache_prop(devi, label, line_str, line_size);
3879     add_cache_prop(devi, label, size_str, size * 1024);
3880 }

3881
3882
3883
3884
3885

```

```

3887 static void
3888 add_amd_l2_assoc(dev_info_t *devi, const char *label, uint_t assoc)
3889 {
3890     switch (assoc) {
3891         case 0: /* off */
3892             break;
3893         case 1:
3894         case 2:
3895         case 4:
3896             add_cache_prop(devi, label, assoc_str, assoc);
3897             break;
3898         case 6:
3899             add_cache_prop(devi, label, assoc_str, 8);
3900             break;
3901         case 8:
3902             add_cache_prop(devi, label, assoc_str, 16);
3903             break;
3904         case 0xf:
3905             add_cache_prop(devi, label, fully_assoc, 1);
3906             break;
3907         default: /* reserved; ignore */
3908             break;
3909     }
3910 }

3912 static void
3913 add_amd_l2_tlb(dev_info_t *devi, const char *label, uint_t assoc, uint_t size)
3914 {
3915     if (size == 0 || assoc == 0)
3916         return;
3917     add_amd_l2_assoc(devi, label, assoc);
3918     add_cache_prop(devi, label, size_str, size);
3919 }

3921 static void
3922 add_amd_l2_cache(dev_info_t *devi, const char *label,
3923                     uint_t size, uint_t assoc, uint_t lines_per_tag, uint_t line_size)
3924 {
3925     if (size == 0 || assoc == 0 || line_size == 0)
3926         return;
3927     add_amd_l2_assoc(devi, label, assoc);
3928     if (lines_per_tag != 0)
3929         add_cache_prop(devi, label, "lines-per-tag", lines_per_tag);
3930     add_cache_prop(devi, label, line_str, line_size);
3931     add_cache_prop(devi, label, size_str, size * 1024);
3932 }

3934 static void
3935 amd_cache_info(struct cpuid_info *cpi, dev_info_t *devi)
3936 {
3937     struct cpuid_regs *cp;

3938     if (cpi->cpi_xmaxeax < 0x80000005)
3939         return;
3940     cp = &cpi->cpi_extd[5];

3941     /*
3942      * 4M/2M L1 TLB configuration
3943      *
3944      * We report the size for 2M pages because AMD uses two
3945      * TLB entries for one 4M page.
3946      */
3947     add_amd_tlb(devi, "dtlb-2M",
3948                 BITX(cp->cp_eax, 31, 24), BITX(cp->cp_eax, 23, 16));
3949     add_amd_tlb(devi, "itlb-2M",
3950                 BITX(cp->cp_eax, 15, 8), BITX(cp->cp_eax, 7, 0));
3951
3952

```

```

3954     /*
3955      * 4K L1 TLB configuration
3956      */
3957
3958     switch (cpi->cpi_vendor) {
3959         uint_t nentries;
3960     case X86_VENDOR_TM:
3961         if (cpi->cpi_family >= 5) {
3962             /*
3963              * Crusoe processors have 256 TLB entries, but
3964              * cpuid data format constrains them to only
3965              * reporting 255 of them.
3966              */
3967             if ((nentries = BITX(cp->cp_ebx, 23, 16)) == 255)
3968                 nentries = 256;
3969             /*
3970              * Crusoe processors also have a unified TLB
3971              */
3972             add_amd_tlb(devi, "tlb-4K", BITX(cp->cp_ebx, 31, 24),
3973                         nentries);
3974             break;
3975         }
3976         /*FALLTHROUGH*/
3977     default:
3978         add_amd_tlb(devi, itlb4k_str,
3979                     BITX(cp->cp_ebx, 31, 24), BITX(cp->cp_ebx, 23, 16));
3980         add_amd_tlb(devi, dtlb4k_str,
3981                     BITX(cp->cp_ebx, 15, 8), BITX(cp->cp_ebx, 7, 0));
3982         break;
3983     }
3984
3985     /*
3986      * data L1 cache configuration
3987      */
3988
3989     add_amd_cache(devi, ll_dcache_str,
3990                   BITX(cp->cp_ecx, 31, 24), BITX(cp->cp_ecx, 23, 16),
3991                   BITX(cp->cp_ecx, 15, 8), BITX(cp->cp_ecx, 7, 0));
3992
3993     /*
3994      * code L1 cache configuration
3995      */
3996
3997     add_amd_cache(devi, ll_icache_str,
3998                   BITX(cp->cp_edx, 31, 24), BITX(cp->cp_edx, 23, 16),
3999                   BITX(cp->cp_edx, 15, 8), BITX(cp->cp_edx, 7, 0));
4000
4001     if (cpi->cpi_xmaxeax < 0x80000006)
4002         return;
4003     cp = &cpi->cpi_extd[6];
4004
4005     /* Check for a unified L2 TLB for large pages */
4006
4007     if (BITX(cp->cp_eax, 31, 16) == 0)
4008         add_amd_l2_tlb(devi, "12-tlb-2M",
4009                         BITX(cp->cp_eax, 15, 12), BITX(cp->cp_eax, 11, 0));
4010     else {
4011         add_amd_l2_tlb(devi, "12-dtlb-2M",
4012                         BITX(cp->cp_eax, 31, 28), BITX(cp->cp_eax, 27, 16));
4013         add_amd_l2_tlb(devi, "12-itlb-2M",
4014                         BITX(cp->cp_eax, 15, 12), BITX(cp->cp_eax, 11, 0));
4015     }
4016
4017     /* Check for a unified L2 TLB for 4K pages */

```

```

4019     if (BITX(cp->cp_ebx, 31, 16) == 0) {
4020         add_amd_12_tlb(devi, "12-tlb-4K",
4021                         BITX(cp->cp_eax, 15, 12), BITX(cp->cp_eax, 11, 0));
4022     } else {
4023         add_amd_12_tlb(devi, "12-dtlb-4K",
4024                         BITX(cp->cp_eax, 31, 28), BITX(cp->cp_eax, 27, 16));
4025         add_amd_12_tlb(devi, "12-itlb-4K",
4026                         BITX(cp->cp_eax, 15, 12), BITX(cp->cp_eax, 11, 0));
4027     }
4028
4029     add_amd_12_cache(devi, 12_cache_str,
4030                         BITX(cp->cp_ecx, 31, 16), BITX(cp->cp_ecx, 15, 12),
4031                         BITX(cp->cp_ecx, 11, 8), BITX(cp->cp_ecx, 7, 0));
4032 }
4034 /*
4035 * There are two basic ways that the x86 world describes it cache
4036 * and tlb architecture - Intel's way and AMD's way.
4037 *
4038 * Return which flavor of cache architecture we should use
4039 */
4040 static int
4041 x86_which_cacheinfo(struct cpuid_info *cpi)
4042 {
4043     switch (cpi->cpi_vendor) {
4044     case X86_VENDOR_Intel:
4045         if (cpi->cpi_maxeax >= 2)
4046             return (X86_VENDOR_Intel);
4047         break;
4048     case X86_VENDOR_AMD:
4049         /*
4050          * The K5 model 1 was the first part from AMD that reported
4051          * cache sizes via extended cpuid functions.
4052         */
4053         if (cpi->cpi_family > 5 ||
4054             (cpi->cpi_family == 5 && cpi->cpi_model >= 1))
4055             return (X86_VENDOR_AMD);
4056         break;
4057     case X86_VENDOR_TM:
4058         if (cpi->cpi_family >= 5)
4059             return (X86_VENDOR_AMD);
4060     /*FALLTHROUGH*/
4061     default:
4062         /*
4063          * If they have extended CPU data for 0x80000005
4064          * then we assume they have AMD-format cache
4065          * information.
4066          *
4067          * If not, and the vendor happens to be Cyrix,
4068          * then try our-Cyrix specific handler.
4069          *
4070          * If we're not Cyrix, then assume we're using Intel's
4071          * table-driven format instead.
4072         */
4073         if (cpi->cpi_xmaxeax >= 0x80000005)
4074             return (X86_VENDOR_AMD);
4075         else if (cpi->cpi_vendor == X86_VENDOR_Cyrix)
4076             return (X86_VENDOR_Cyrix);
4077         else if (cpi->cpi_maxeax >= 2)
4078             return (X86_VENDOR_Intel);
4079         break;
4080     }
4081     return (-1);
4082 }
4084 void

```

```

4085 cpuid_set_cpu_properties(void *dip, processorid_t cpu_id,
4086     struct cpuid_info *cpi)
4087 {
4088     dev_info_t *cpu_devi;
4089     int create;
4090
4091     cpu_devi = (dev_info_t *)dip;
4092
4093     /* device_type */
4094     (void) ndi_prop_update_string(DDI_DEV_T_NONE, cpu_devi,
4095                                 "device_type", "cpu");
4096
4097     /* reg */
4098     (void) ndi_prop_update_int(DDI_DEV_T_NONE, cpu_devi,
4099                               "reg", cpu_id);
4100
4101     /* cpu-mhz, and clock-frequency */
4102     if (cpu_freq > 0) {
4103         long long mul;
4104
4105         (void) ndi_prop_update_int(DDI_DEV_T_NONE, cpu_devi,
4106                                   "cpu-mhz", cpu_freq);
4107         if ((mul = cpu_freq * 1000000LL) <= INT_MAX)
4108             (void) ndi_prop_update_int(DDI_DEV_T_NONE, cpu_devi,
4109                                       "clock-frequency", (int)mul);
4110     }
4111
4112     if (!is_x86_feature(x86_featureset, X86FSET_CPUID)) {
4113         return;
4114     }
4115
4116     /* vendor-id */
4117     (void) ndi_prop_update_string(DDI_DEV_T_NONE, cpu_devi,
4118                                 "vendor-id", cpi->cpi_vendorstr);
4119
4120     if (cpi->cpi_maxeax == 0) {
4121         return;
4122     }
4123
4124     /*
4125      * family, model, and step
4126      */
4127     (void) ndi_prop_update_int(DDI_DEV_T_NONE, cpu_devi,
4128                               "family", CPI_FAMILY(cpi));
4129     (void) ndi_prop_update_int(DDI_DEV_T_NONE, cpu_devi,
4130                               "cpu-model", CPI_MODEL(cpi));
4131     (void) ndi_prop_update_int(DDI_DEV_T_NONE, cpu_devi,
4132                               "stepping-id", CPI_STEP(cpi));
4133
4134     /* type */
4135     switch (cpi->cpi_vendor) {
4136     case X86_VENDOR_Intel:
4137         create = 1;
4138         break;
4139     default:
4140         create = 0;
4141         break;
4142     }
4143     if (create)
4144         (void) ndi_prop_update_int(DDI_DEV_T_NONE, cpu_devi,
4145                                   "type", CPI_TYPE(cpi));
4146
4147     /* ext-family */
4148     switch (cpi->cpi_vendor) {
4149     case X86_VENDOR_Intel:
4150     case X86_VENDOR_AMD:

```

```

4151         create = cpi->cpi_family >= 0xf;
4152         break;
4153     default:
4154         create = 0;
4155         break;
4156     }
4157     if (create)
4158         (void) ndi_prop_update_int(DDI_DEV_T_NONE, cpu_devi,
4159             "ext-family", CPI_FAMILY_XTD(cpi));
4160
4161     /* ext-model */
4162     switch (cpi->cpi_vendor) {
4163     case X86_VENDOR_Intel:
4164         create = IS_EXTENDED_MODEL_INTEL(cpi);
4165         break;
4166     case X86_VENDOR_AMD:
4167         create = CPI_FAMILY(cpi) == 0xf;
4168         break;
4169     default:
4170         create = 0;
4171         break;
4172     }
4173     if (create)
4174         (void) ndi_prop_update_int(DDI_DEV_T_NONE, cpu_devi,
4175             "ext-model", CPI_MODEL_XTD(cpi));
4176
4177     /* generation */
4178     switch (cpi->cpi_vendor) {
4179     case X86_VENDOR_AMD:
4180         /*
4181          * AMD K5 model 1 was the first part to support this
4182          */
4183         create = cpi->cpi_xmaxeax >= 0x80000001;
4184         break;
4185     default:
4186         create = 0;
4187         break;
4188     }
4189     if (create)
4190         (void) ndi_prop_update_int(DDI_DEV_T_NONE, cpu_devi,
4191             "generation", BITX((cpi)->cpi_extd[1].cp_eax, 11, 8));
4192
4193     /* brand-id */
4194     switch (cpi->cpi_vendor) {
4195     case X86_VENDOR_Intel:
4196         /*
4197          * brand id first appeared on Pentium III Xeon model 8,
4198          * and Celeron model 8 processors and Opteron
4199          */
4200         create = cpi->cpi_family > 6 ||
4201             (cpi->cpi_family == 6 && cpi->cpi_model >= 8);
4202         break;
4203     case X86_VENDOR_AMD:
4204         create = cpi->cpi_family >= 0xf;
4205         break;
4206     default:
4207         create = 0;
4208         break;
4209     }
4210     if (create && cpi->cpi_brandid != 0) {
4211         (void) ndi_prop_update_int(DDI_DEV_T_NONE, cpu_devi,
4212             "brand-id", cpi->cpi_brandid);
4213     }
4214
4215     /* chunks, and apic-id */
4216     switch (cpi->cpi_vendor) {

```

```

4217         /*
4218          * first available on Pentium IV and Opteron (K8)
4219          */
4220
4221     case X86_VENDOR_Intel:
4222         create = IS_NEW_F6(cpi) || cpi->cpi_family >= 0xf;
4223         break;
4224     case X86_VENDOR_AMD:
4225         create = cpi->cpi_family >= 0xf;
4226         break;
4227     default:
4228         create = 0;
4229         break;
4230     }
4231     if (create) {
4232         (void) ndi_prop_update_int(DDI_DEV_T_NONE, cpu_devi,
4233             "chunks", CPI_CHUNKS(cpi));
4234         (void) ndi_prop_update_int(DDI_DEV_T_NONE, cpu_devi,
4235             "apic-id", cpi->cpi_apicid);
4236         if (cpi->cpi_chipid >= 0) {
4237             (void) ndi_prop_update_int(DDI_DEV_T_NONE, cpu_devi,
4238                 "chip#", cpi->cpi_chipid);
4239             (void) ndi_prop_update_int(DDI_DEV_T_NONE, cpu_devi,
4240                 "clog#", cpi->cpi_clogid);
4241         }
4242
4243     /* cpuid-features */
4244     (void) ndi_prop_update_int(DDI_DEV_T_NONE, cpu_devi,
4245             "cpuid-features", CPI_FEATURES_EDX(cpi));
4246
4247     /* cpuid-features-ecx */
4248     switch (cpi->cpi_vendor) {
4249     case X86_VENDOR_Intel:
4250         create = IS_NEW_F6(cpi) || cpi->cpi_family >= 0xf;
4251         break;
4252     case X86_VENDOR_AMD:
4253         create = cpi->cpi_family >= 0xf;
4254         break;
4255     default:
4256         create = 0;
4257         break;
4258     }
4259     if (create)
4260         (void) ndi_prop_update_int(DDI_DEV_T_NONE, cpu_devi,
4261             "cpuid-features-ecx", CPI_FEATURES_ECX(cpi));
4262
4263     /* ext-cpuid-features */
4264     switch (cpi->cpi_vendor) {
4265     case X86_VENDOR_Intel:
4266     case X86_VENDOR_AMD:
4267     case X86_VENDOR_Cyrix:
4268     case X86_VENDOR_TM:
4269     case X86_VENDOR_Centaur:
4270         create = cpi->cpi_xmaxeax >= 0x80000001;
4271         break;
4272     default:
4273         create = 0;
4274         break;
4275     }
4276     if (create) {
4277         (void) ndi_prop_update_int(DDI_DEV_T_NONE, cpu_devi,
4278             "ext-cpuid-features", CPI_FEATURES_XTD_EDX(cpi));
4279         (void) ndi_prop_update_int(DDI_DEV_T_NONE, cpu_devi,
4280             "ext-cpuid-features-ecx", CPI_FEATURES_XTD_ECX(cpi));
4281     }
4282

```

```

4284 /*
4285  * Brand String first appeared in Intel Pentium IV, AMD K5
4286  * model 1, and Cyrix GXm. On earlier models we try and
4287  * simulate something similar .. so this string should always
4288  * same -something- about the processor, however lame.
4289 */
4290 (void) ndi_prop_update_string(DDI_DEV_T_NONE, cpu_devi,
4291     "brand-string", cpi->cpi_brandstr);
4292
4293 /*
4294  * Finally, cache and tlb information
4295 */
4296 switch (x86_which_cacheinfo(cpi)) {
4297 case X86_VENDOR_Intel:
4298     intel_walk_cacheinfo(cpi, cpu_devi, add_cacheent_props);
4299     break;
4300 case X86_VENDOR_Cyrix:
4301     cyrix_walk_cacheinfo(cpi, cpu_devi, add_cacheent_props);
4302     break;
4303 case X86_VENDOR_AMD:
4304     amd_cache_info(cpi, cpu_devi);
4305     break;
4306 default:
4307     break;
4308 }
4309
4310 struct l2info {
4311     int *l2i_csz;
4312     int *l2i_lsz;
4313     int *l2i_assoc;
4314     int l2i_ret;
4315 };
4316
4317 /*
4318  * A cacheinfo walker that fetches the size, line-size and associativity
4319  * of the L2 cache
4320 */
4321 static int
4322 intel_l2cinfo(void *arg, const struct cachetab *ct)
4323 {
4324     struct l2info *l2i = arg;
4325     int *ip;
4326
4327     if (ct->ct_label != l2_cache_str &&
4328         ct->ct_label != sl2_cache_str)
4329         return (0); /* not an L2 -- keep walking */
4330
4331     if ((ip = l2i->l2i_csz) != NULL)
4332         *ip = ct->ct_size;
4333     if ((ip = l2i->l2i_lsz) != NULL)
4334         *ip = ct->ct_line_size;
4335     if ((ip = l2i->l2i_assoc) != NULL)
4336         *ip = ct->ct_assoc;
4337     l2i->l2i_ret = ct->ct_size;
4338     return (1); /* was an L2 -- terminate walk */
4339 }
4340
4341 /*
4342  * AMD L2/L3 Cache and TLB Associativity Field Definition:
4343 */
4344 *
4345 * Unlike the associativity for the L1 cache and tlb where the 8 bit
4346 * value is the associativity, the associativity for the L2 cache and
4347 * tlb is encoded in the following table. The 4 bit L2 value serves as
4348 * an index into the amd_afd[] array to determine the associativity.

```

```

4349 *      -1 is undefined. 0 is fully associative.
4350 */
4351 static int amd_afd[] =
4352     {-1, 1, 2, -1, 4, -1, 8, -1, 16, -1, 32, 48, 64, 96, 128, 0};
4353
4354 static void
4355 amd_l2cacheinfo(struct cpuid_info *cpi, struct l2info *l2i)
4356 {
4357     struct cpuid_regs *cp;
4358     uint_t size, assoc;
4359     int i;
4360     int *ip;
4361
4362     if (cpi->cpi_xmaxeax < 0x80000006)
4363         return;
4364     cp = &cpi->cpi_extd[6];
4365
4366     if ((i = BITX(cp->cp_ecx, 15, 12)) != 0 &&
4367         (size = BITX(cp->cp_ecx, 31, 16)) != 0) {
4368         uint_t cachesz = size * 1024;
4369         assoc = amd_afd[i];
4370
4371         ASSERT(assoc != -1);
4372
4373         if ((ip = l2i->l2i_csz) != NULL)
4374             *ip = cachesz;
4375         if ((ip = l2i->l2i_lsz) != NULL)
4376             *ip = BITX(cp->cp_ecx, 7, 0);
4377         if ((ip = l2i->l2i_assoc) != NULL)
4378             *ip = assoc;
4379         l2i->l2i_ret = cachesz;
4380     }
4381
4382 }
4383
4384 int
4385 getl2cacheinfo(cpu_t *cpu, int *csz, int *lsz, int *assoc)
4386 {
4387     struct cpuid_info *cpi = cpu->cpu_m.mcpu_cpi;
4388     struct l2info __l2info, *l2i = &__l2info;
4389
4390     l2i->l2i_csz = csz;
4391     l2i->l2i_lsz = lsz;
4392     l2i->l2i_assoc = assoc;
4393     l2i->l2i_ret = -1;
4394
4395     switch (x86_which_cacheinfo(cpi)) {
4396     case X86_VENDOR_Intel:
4397         intel_walk_cacheinfo(cpi, l2i, intel_l2cinfo);
4398         break;
4399     case X86_VENDOR_Cyrix:
4400         cyrix_walk_cacheinfo(cpi, l2i, intel_l2cinfo);
4401         break;
4402     case X86_VENDOR_AMD:
4403         amd_l2cacheinfo(cpi, l2i);
4404         break;
4405     default:
4406         break;
4407     }
4408     return (l2i->l2i_ret);
4409 }
4410
4411 #if !defined(__xpv)
4412 uint32_t *
4413 cpuid_mwait_alloc(cpu_t *cpu)

```

```

4415 {
4416     uint32_t      *ret;
4417     size_t        mwait_size;
4418
4419     ASSERT(cpuid_checkpass(CPU, 2));
4420
4421     mwait_size = CPU->cpu_m.mcpu_cpi->cpi_mwait.mon_max;
4422     if (mwait_size == 0)
4423         return (NULL);
4424
4425     /*
4426      * kmem_alloc() returns cache line size aligned data for mwait_size
4427      * allocations. mwait_size is currently cache line sized. Neither
4428      * of these implementation details are guaranteed to be true in the
4429      * future.
4430
4431      * First try allocating mwait_size as kmem_alloc() currently returns
4432      * correctly aligned memory. If kmem_alloc() does not return
4433      * mwait_size aligned memory, then use mwait_size ROUNDUP.
4434
4435      * Set cpi_mwait.buf_actual and cpi_mwait.size_actual in case we
4436      * decide to free this memory.
4437
4438     ret = kmem_zalloc(mwait_size, KM_SLEEP);
4439     if (ret == (uint32_t *)P2ROUNDUP((uintptr_t)ret, mwait_size)) {
4440         cpu->cpu_m.mcpu_cpi->cpi_mwait.buf_actual = ret;
4441         cpu->cpu_m.mcpu_cpi->cpi_mwait.size_actual = mwait_size;
4442         *ret = MWAIT_RUNNING;
4443         return (ret);
4444     } else {
4445         kmem_free(ret, mwait_size);
4446         ret = kmem_zalloc(mwait_size * 2, KM_SLEEP);
4447         cpu->cpu_m.mcpu_cpi->cpi_mwait.buf_actual = ret;
4448         cpu->cpu_m.mcpu_cpi->cpi_mwait.size_actual = mwait_size * 2;
4449         ret = (uint32_t *)P2ROUNDUP((uintptr_t)ret, mwait_size);
4450         *ret = MWAIT_RUNNING;
4451         return (ret);
4452     }
4453 }
4454
4455 void
4456 cpuid_mwait_free(cpu_t *cpu)
4457 {
4458     if (cpu->cpu_m.mcpu_cpi == NULL) {
4459         return;
4460     }
4461
4462     if (cpu->cpu_m.mcpu_cpi->cpi_mwait.buf_actual != NULL &
4463         cpu->cpu_m.mcpu_cpi->cpi_mwait.size_actual > 0) {
4464         kmem_free(cpu->cpu_m.mcpu_cpi->cpi_mwait.buf_actual,
4465                   cpu->cpu_m.mcpu_cpi->cpi_mwait.size_actual);
4466     }
4467
4468     cpu->cpu_m.mcpu_cpi->cpi_mwait.buf_actual = NULL;
4469     cpu->cpu_m.mcpu_cpi->cpi_mwait.size_actual = 0;
4470 }
4471
4472 void
4473 patch_tsc_read(int flag)
4474 {
4475     size_t cnt;
4476
4477     switch (flag) {
4478     case X86_NO_TSC:
4479         cnt = &no_rdtsc_end - &no_rdtsc_start;
4480         (void) memcpy((void *)tsc_read, (void *)&no_rdtsc_start, cnt);

```

```

4481             break;
4482     case X86_HAVE_TSCK:
4483         cnt = &tsc_end - &tsc_start;
4484         (void) memcpy((void *)tsc_read, (void *)&tsc_start, cnt);
4485         break;
4486     case X86_TSC_MFENCE:
4487         cnt = &tsc_mfence_end - &tsc_mfence_start;
4488         (void) memcpy((void *)tsc_read,
4489                      (void *)&tsc_mfence_start, cnt);
4490         break;
4491     case X86_TSC_LFENCE:
4492         cnt = &tsc_lfence_end - &tsc_lfence_start;
4493         (void) memcpy((void *)tsc_read,
4494                      (void *)&tsc_lfence_start, cnt);
4495         break;
4496     default:
4497         break;
4498     }
4499 }
4500
4501 int
4502 cpuid_deep_cstates_supported(void)
4503 {
4504     struct cpuid_info *cpi;
4505     struct cpuid_regs regs;
4506
4507     ASSERT(cpuid_checkpass(CPU, 1));
4508
4509     cpi = CPU->cpu_m.mcpu_cpi;
4510
4511     if (!is_x86_feature(x86_featureset, X86FSET_CPUID))
4512         return (0);
4513
4514     switch (cpi->cpi_vendor) {
4515     case X86_VENDOR_Intel:
4516         if (cpi->cpi_xmaxeax < 0x80000007)
4517             return (0);
4518
4519         /*
4520          * TSC run at a constant rate in all ACPI C-states?
4521          */
4522         regs.cp_eax = 0x80000007;
4523         (void) __cpuid_insn(&regs);
4524         return (regs.cp_edx & CPUID_TSC_CSTATE_INVARIANCE);
4525
4526     default:
4527         return (0);
4528     }
4529 }
4530
4531 #endif /* !__xpv */
4532
4533 void
4534 post_startup_cpu_fixups(void)
4535 {
4536 #ifndef __xpv
4537     /*
4538      * Some AMD processors support C1E state. Entering this state will
4539      * cause the local APIC timer to stop, which we can't deal with at
4540      * this time.
4541      */
4542     if (cpuid_getvendor(CPU) == X86_VENDOR_AMD) {
4543         on_trap_data_t otd;
4544         uint64_t reg;
4545
4546         if (!on_trap(&otd, OT_DATA_ACCESS)) {

```

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

59

```

4547         reg = rdmsr(MSR_AMD_INT_PENDING_CMP_HALT);
4548         /* Disable C1E state if it is enabled by BIOS */
4549         if ((reg >> AMD_ACTONCMPHALT_SHIFT) &
4550             AMD_ACTONCMPHALT_MASK) {
4551             reg &= ~(AMD_ACTONCMPHALT_MASK <<
4552                     AMD_ACTONCMPHALT_SHIFT);
4553             wrmsr(MSR_AMD_INT_PENDING_CMP_HALT, reg);
4554         }
4555     }
4556     no_trap();
4557 }
4558 #endif /* ! __xpv */
4559 }

4561 /*
4562  * Setup necessary registers to enable XSAVE feature on this processor.
4563  * This function needs to be called early enough, so that no xsave/xrstor
4564  * ops will execute on the processor before the MSRs are properly set up.
4565  *
4566  * Current implementation has the following assumption:
4567  * - cpuid_pass1() is done, so that X86 features are known.
4568  * - fpu_probe() is done, so that fp_save_mech is chosen.
4569  */
4570 void
4571 xsave_setup_msr(cpu_t *cpu)
4572 {
4573     ASSERT(fp_save_mech == FP_XSAVE);
4574     ASSERT(is_x86_feature(x86_featureset, X86FSET_XSAVE));

4576     /* Enable OSXSAVE in CR4. */
4577     setcr4(getcr4() | CR4 OSXSAVE);
4578     /*
4579      * Update SW copy of ECX, so that /dev/cpu/self/cpuid will report
4580      * correct value.
4581      */
4582     cpu->cpu_m.mcpu_cpi->cpi_std[1].cp_ecx |= CPUID_INTC_ECX OSXSAVE;
4583     setup_xfem();
4584 }

4586 /*
4587  * Starting with the Westmere processor the local
4588  * APIC timer will continue running in all C-states,
4589  * including the deepest C-states.
4590  */
4591 int
4592 cpuid_arat_supported(void)
4593 {
4594     struct cpuid_info *cpi;
4595     struct cpuid_regs regs;

4597     ASSERT(cpuid_checkpass(CPU, 1));
4598     ASSERT(is_x86_feature(x86_featureset, X86FSET_CPUID));

4600     cpi = CPU->cpu_m.mcpu_cpi;

4602     switch (cpi->cpi_vendor) {
4603     case X86_VENDOR_Intel:
4604         /*
4605          * Always-running Local APIC Timer is
4606          * indicated by CPUID.6.EAX[2].
4607          */
4608         if (cpi->cpi_maxeax >= 6) {
4609             regs.cp_eax = 6;
4610             (void) cpuid_insn(NULL, &regs);
4611             return (regs.cp_eax & CPUID_CSTATE_ARAT);
4612         } else {

```

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

```

4613         }
4614     default:
4615         return (0);
4616     }
4617 }
4618 }

4620 /*
4621 * Check support for Intel ENERGY_PERF_BIAS feature
4622 */
4623 int
4624 cpuid_iepb_supported(struct cpu *cp)
4625 {
4626     struct cpuid_info *cpi = cp->cpu_m.mcpu_cpi;
4627     struct cpuid_regs regs;

4628     ASSERT(cpuid_checkpass(cp, 1));

4629     if (!!(is_x86_feature(x86_featureset, X86FSET_CPUID)) ||
4630         !!(is_x86_feature(x86_featureset, X86FSET_MSR))) {
4631         return (0);
4632     }

4633     /*
4634      * Intel ENERGY_PERF_BIAS MSR is indicated by
4635      * capability bit CPUID.6.ECX.3
4636      */
4637     if ((cpi->cpi_vendor != X86_VENDOR_Intel) || (cpi->cpi_maxeax < 6))
4638         return (0);

4639     regs.cp_eax = 0x6;
4640     (void) cpuid_insn(NULL, &regs);
4641     return (regs.cp_ecx & CPUID_EPB_SUPPORT);
4642 }

4643 /*
4644 * Check support for TSC deadline timer
4645 */
4646 *
4647 * TSC deadline timer provides a superior software programming
4648 * model over local APIC timer that eliminates "time drifts".
4649 * Instead of specifying a relative time, software specifies an
4650 * absolute time as the target at which the processor should
4651 * generate a timer event.
4652 */
4653 int
4654 cpuid_deadline_tsc_supported(void)
4655 {
4656     struct cpuid_info *cpi = CPU->cpu_m.mcpu_cpi;
4657     struct cpuid_regs regs;

4658     ASSERT(cpuid_checkpass(CPU, 1));
4659     ASSERT(is_x86_feature(x86_featureset, X86FSET_CPUID));

4660     switch (cpi->cpi_vendor) {
4661     case X86_VENDOR_Intel:
4662         if (cpi->cpi_maxeax >= 1) {
4663             regs.cp_eax = 1;
4664             (void) cpuid_insn(NULL, &regs);
4665             return (regs.cp_ecx & CPUID_DEADLINE_TSC);
4666         } else {
4667             return (0);
4668         }
4669     default:
4670         return (0);
4671     }
4672 }
4673 }
4674 }
4675 }
4676 }
4677 }
4678 }

```

```

4680 #if defined(__amd64) && !defined(__xpv)
4681 /*
4682 * Patch in versions of bcopy for high performance Intel Nhm processors
4683 * and later...
4684 */
4685 void
4686 patch_memops(uint_t vendor)
4687 {
4688     size_t cnt, i;
4689     caddr_t to, from;
4690
4691     if ((vendor == X86_VENDOR_Intel) &&
4692         is_x86_feature(x86_featureset, X86FSET_SSE4_2)) {
4693         cnt = &bcopy_patch_end - &bcopy_patch_start;
4694         to = &bcopy_ck_size;
4695         from = &bcopy_patch_start;
4696         for (i = 0; i < cnt; i++) {
4697             *to++ = *from++;
4698         }
4699     }
4700 }
4701 #endif /* __amd64 && !__xpv */

4703 /*
4704 * This function finds the number of bits to represent the number of cores per
4705 * chip and the number of strands per core for the Intel platforms.
4706 * It re-uses the x2APIC cpuid code of the cpuid_pass2().
4707 */
4708 void
4709 cpuid_get_ext_topo(uint_t vendor, uint_t *core_nborts, uint_t *strand_nborts)
4710 {
4711     struct cpuid_regs regs;
4712     struct cpuid_regs *cp = &regs;
4713
4714     if (vendor != X86_VENDOR_Intel) {
4715         return;
4716     }

4718     /* if the cpuid level is 0xB, extended topo is available. */
4719     cp->cp_eax = 0;
4720     if (__cpuid_insn(cp) >= 0xB) {

4722         cp->cp_eax = 0xB;
4723         cp->cp_edx = cp->cp_ebx = cp->cp_ecx = 0;
4724         (void) __cpuid_insn(cp);

4726         /*
4727          * Check CPUID.EAX=0BH, ECX=0H:EBX is non-zero, which
4728          * indicates that the extended topology enumeration leaf is
4729          * available.
4730          */
4731         if (cp->cp_ebx) {
4732             uint_t coreid_shift = 0;
4733             uint_t chipid_shift = 0;
4734             uint_t i;
4735             uint_t level;

4737             for (i = 0; i < CPI_FNB_ECX_MAX; i++) {
4738                 cp->cp_eax = 0xB;
4739                 cp->cp_ecx = i;

4741                 (void) __cpuid_insn(cp);
4742                 level = CPI_CPU_LEVEL_TYPE(cp);

4744             if (level == 1) {

```

```

4745
4746
4747
4748
4749
4750
4751
4752
4753
4754
4755
4756
4757
4758
4759

4761
4762
4763
4764
4765
4766
4767 }

/*
 * Thread level processor topology
 * Number of bits shift right APIC ID
 * to get the coreid.
 */
coreid_shift = BITX(cp->cp_eax, 4, 0);
} else if (level == 2) {
/*
 * Core level processor topology
 * Number of bits shift right APIC ID
 * to get the chipid.
 */
chipid_shift = BITX(cp->cp_eax, 4, 0);
}

if (coreid_shift > 0 && chipid_shift > coreid_shift) {
    *strand_nborts = coreid_shift;
    *core_nborts = chipid_shift - coreid_shift;
}
}
```

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

34387 Fri Jan 3 22:11:56 2014

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

patch cpu-pause-func-deglobalize

```
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/_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/sysm.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/ddidmareq.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/swapp.h>  
97 #include <sys/thread.h>  
98 #include <sys/sysconf.h>  
99 #include <sys/vm_machparam.h>  
100 #include <sys/archsystm.h>  
101 #include <sys/machsystm.h>  
102 #include <sys/machclock.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";
211     if (fcn == AD_FASTREBOOT && !newkernel.fi_valid)
212         fcn = AD_BOOT;
214     if (!panicstr) {
215         kpreempt_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     }
227     if (force_shutdown_method != AD_UNKNOWN)
228         fcn = force_shutdown_method;
230     /*
231      * XXX - rconsrv 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     rconsrv = NULL;
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");
246     if (IN_XPV_PANIC())
247         reset();
249     /*
250      * We can't bring up the console from above lock level, so do it now
251      */
252     pm_cfb_check_and_powerup();
254     /* make sure there are no more changes to the device tree */
255     devtree_freeze();
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();

265 #if defined(__xpv)
266     /*
267      * XXXPV Should probably think some more about how we deal
268      * with panicing before it's really safe to panic.
269      * On hypervisors, we reboot very quickly.. Perhaps panic
270      * should only attempt to recover by rebooting if,
271      * say, we were able to mount the root filesystem,
272      * or if we successfully launched init(1m).
273     */
274     if (panicstr && proc_init == NULL)
275         (void) HYPERVISOR_shutdown(SHUTDOWN_poweroff);
276 #endif
277     /*
278      * stop other cpus and raise our priority. since there is only
279      * one active cpu after this, and our priority will be too high
280      * for us to be preempted, we're essentially single threaded
281      * from here on out.
282     */
283     (void) spl6();
284     if (!panicstr) {
285         mutex_enter(&cpu_lock);
286         pause_cpus(NULL, NULL);
287         pause_cpus(NULL);
288         mutex_exit(&cpu_lock);
289     }

290     /*
291      * If the system is panicking, the preloaded kernel is valid, and
292      * fastreboot_onpanic has been set, and the system has been up for
293      * longer than fastreboot_onpanic_uptime (default to 10 minutes),
294      * choose Fast Reboot.
295     */
296     if (fcn == AD_BOOT && panicstr && newkernel.fi_valid &&
297         fastreboot_onpanic &&
298         (panic_lbolt - lboot_at_boot) > fastreboot_onpanic_uptime) {
299         fcn = AD_FASTREBOOT;
300     }

302     /*
303      * Try to quiesce devices.
304     */
305     if (is_first_quiesce) {
306         /*
307          * Clear is_first_quiesce before calling quiesce_devices()
308          * so that if quiesce_devices() causes panics, it will not
309          * be invoked again.
310         */
311         is_first_quiesce = 0;

313         quiesce_active = 1;
314         quiesce_devices(ddi_root_node(), &reset_status);
315         if (reset_status == -1) {
316             if (fcn == AD_FASTREBOOT && !force_fastreboot) {
317                 prom_printf("Driver(s) not capable of fast "
318                             "reboot.\n");
319                 prom_printf(fallback_str);
320                 fastreboot_capable = 0;
321                 fcn = AD_BOOT;
322             } else if (fcn != AD_FASTREBOOT)
323                 fastreboot_capable = 0;
324         }

```

```

325             quiesce_active = 0;
326         }

328         /*
329          * Try to reset devices. reset_leaves() should only be called
330          * a) when there are no other threads that could be accessing devices,
331          * and
332          * b) on a system that's not capable of fast reboot (fastreboot_capable
333          * being 0), or on a system where quiesce_devices() failed to
334          * complete (quiesce_active being 1).
335         */
336         if (is_first_reset && (!fastreboot_capable || quiesce_active)) {
337             /*
338              * Clear is_first_reset before calling reset_devices()
339              * so that if reset_devices() causes panics, it will not
340              * be invoked again.
341             */
342             is_first_reset = 0;
343             reset_leaves();
344         }

346         /* Verify newkernel checksum */
347         if (fastreboot_capable && fcn == AD_FASTREBOOT &&
348             fastboot_cksum_verify(newkernel) != 0) {
349             fastreboot_capable = 0;
350             prom_printf("Fast reboot: checksum failed for the new "
351                         "kernel.\n");
352             prom_printf(fallback_str);
353         }

355         (void) spl8();

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

371 } unchanged_portion_omitted

```

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

17010 Fri Jan 3 22:11:56 2014

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

patch cpu-pause-func-deglobalize

```
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

36 #include <sys/types.h>
37 #include <sys/sysm.h>
38 #include <sys/cpuvar.h>
39 #include <sys/cpu_module.h>
40 #include <sys/kmem.h>
41 #include <sys/archsysm.h>
42 #include <sys/machsysm.h>
43 #include <sys/controlregs.h>
44 #include <sys/x86_archext.h>
45 #include <sys/smp_impldefs.h>
46 #include <sys/sysmacros.h>
47 #include <sys/mach_mmu.h>
48 #include <sys/promif.h>
49 #include <sys/cpu.h>
50 #include <sys/cpu_event.h>
51 #include <sys/sunudi.h>
52 #include <sys/fs/dv_node.h>
53 #include <vm/hat_i86.h>
54 #include <vm/as.h>

55 extern cpuset_t cpu_ready_set;

56 extern int mp_start_cpu_common(cpu_t *cp, boolean_t boot);
57 extern void real_mode_start_cpu(void);
58 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 /* Fill up the real mode platter to make it easy for real mode code to
69 * kick it off. This area should really be one passed by boot to kernel
70 * and guaranteed to be below 1MB and aligned to 16 bytes. Should also
71 * have identical physical and virtual address in paged mode.
72 */
73 static ushort_t *warm_reset_vector = NULL;

74 int
75 mach_cpucontext_init(void)
76 {
77     ushort_t *vec;
78     ulong_t addr;
79     struct rm_platter *rm = (struct rm_platter *)rmpl_platter_va;
80
81     if (!(vec = (ushort_t *)psm_map_phys(WARM_RESET_VECTOR,
82                                         sizeof(vec), PROT_READ | PROT_WRITE)))
83         return (-1);

84     /*
85      * setup secondary cpu bios boot up vector
86      * Write page offset to 0x467 and page frame number to 0x469.
87      */
88     addr = (ulong_t)((caddr_t)rm->rm_code - (caddr_t)rm) + rmplatter_pa;
89     vec[0] = (ushort_t)(addr & PAGEOFFSET);
90     vec[1] = (ushort_t)((addr & 0xffff & PAGEMASK) >> 4);
91     warm_reset_vector = vec;

92     /* Map real mode platter into kas so kernel can access it. */
93     hat_devload(kas.a_hat,
94                 (caddr_t)(uintptr_t)rmplatter_pa, MMU_PAGESIZE,
95                 btop(rmplatter_pa), PROT_READ | PROT_WRITE | PROT_EXEC,
96                 HAT_LOAD_NOCONSIST);

97     /* Copy CPU startup code to rm_platter if it's still during boot. */
98     if (!plat_dr_enabled()) {
99         ASSERT((size_t)real_mode_start_cpu_end -
100                (size_t)real_mode_start_cpu <= RM_PLATTER_CODE_SIZE);
101         bcopy((caddr_t)real_mode_start_cpu, (caddr_t)rm->rm_code,
102               (size_t)real_mode_start_cpu_end -
103               (size_t)real_mode_start_cpu);
104     }
105
106     return (0);
107 }
```

unchanged portion omitted

2

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

```
1
```

```
*****  
18803 Fri Jan 3 22:11:56 2014  
new/usr/src/uts/i86pc/os/x_call.c
```

```
patch cpu-pause-func-deglobalize
```

```
*****  
_____ 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 {  
    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 }  
299 }  
_____ unchanged_portion_omitted _____
```

```
new/usr/src/uts/i86xpv/os/mp_xen.c
```

```
1
```

```
*****
25091 Fri Jan 3 22:11:56 2014
new/usr/src/uts/i86xpv/os/mp_xen.c
patch cpu-pause-func-deglobalize
*****
_____unchanged_portion_omitted_____
```

```
579 void
580 mp_enter_barrier(void)
581 {
582     hrtimer_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);
590     pause_cpus(NULL);
591
592     while (!done) {
593         done = 1;
594         poke_allowed = 0;
595
596         if (xpv_gethrtime() - last_poke_time > POKE_TIMEOUT) {
597             last_poke_time = xpv_gethrtime();
598             poke_allowed = 1;
599         }
600
601         for (i = 0; i < NCPUs; i++) {
602             cpu_t *cp = cpu_get(i);
603
604             if (cp == NULL || cp == CPU)
605                 continue;
606
607             switch (cpu_phase[i]) {
608                 case CPU_PHASE_NONE:
609                     cpu_phase[i] = CPU_PHASE_WAIT_SAFE;
610                     poke_cpu(i);
611                     done = 0;
612                     break;
613
614                 case CPU_PHASE_WAIT_SAFE:
615                     if (poke_allowed)
616                         poke_cpu(i);
617                     done = 0;
618                     break;
619
620                 case CPU_PHASE_SAFE:
621                 case CPU_PHASE_POWERED_OFF:
622                     break;
623             }
624
625             SMT_PAUSE();
626         }
627     }
628
629     _____unchanged_portion_omitted_____
```

```
new/usr/src/uts/intel/sys/x86_archext.h
```

```
1
```

```
*****
28993 Fri Jan 3 22:11:56 2014
new/usr/src/uts/intel/sys/x86_archext.h
patch x2apic-x86fset
patch remove-unused-vars
*****
1 /*
2 * CDDL HEADER START
3 *
4 * The contents of this file are subject to the terms of the
5 * Common Development and Distribution License (the "License").
6 * You may not use this file except in compliance with the License.
7 *
8 * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9 * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22 * Copyright (c) 1995, 2010, Oracle and/or its affiliates. All rights reserved.
23 * Copyright (c) 2011 by Delphix. All rights reserved.
24 * Copyright 2012 Nexenta Systems, Inc. All rights reserved.
25 */
26 /*
27 * Copyright (c) 2010, Intel Corporation.
28 * All rights reserved.
29 */
30 /*
31 * Copyright (c) 2012, Joyent, Inc. All rights reserved.
32 * Copyright 2012 Jens Elkner <jel+illumos@cs.uni-magdeburg.de>
33 * Copyright 2012 Hans Rosenfeld <rosenfeld@grumpf.hope-2000.org>
34 */

36 #ifndef _SYS_X86_ARCHEXT_H
37 #define _SYS_X86_ARCHEXT_H

39 #if !defined(_ASM)
40 #include <sys/regset.h>
41 #include <sys/processor.h>
42 #include <vm/seg_enum.h>
43 #include <vm/page.h>
44 #endif /* _ASM */

46 #ifdef __cplusplus
47 extern "C" {
48 #endif

50 /*
51 * cpuid instruction feature flags in %edx (standard function 1)
52 */

54 #define CPUID_INTC_EDX_FPU 0x00000001 /* x87 fpu present */
55 #define CPUID_INTC_EDX_VME 0x00000002 /* virtual-8086 extension */
56 #define CPUID_INTC_EDX_DE 0x00000004 /* debugging extensions */
57 #define CPUID_INTC_EDX_PSE 0x00000008 /* page size extension */
58 #define CPUID_INTC_EDX_TSC 0x00000010 /* time stamp counter */
59 #define CPUID_INTC_EDX_MSR 0x00000020 /* rdmsr and wrmsr */
60 #define CPUID_INTC_EDX_PAE 0x00000040 /* physical addr extension */

*****
```

```
new/usr/src/uts/intel/sys/x86_archext.h
```

```
2
```

```
61 #define CPUID_INTC_EDX_MCE 0x00000080 /* machine check exception */
62 #define CPUID_INTC_EDX_CX8 0x00000100 /* cmpxchg8b instruction */
63 #define CPUID_INTC_EDX_APIC 0x00000200 /* local APIC */
64 /* 0x400 - reserved */
65 #define CPUID_INTC_EDX_SEP 0x00000800 /* sysenter and sysexit */
66 #define CPUID_INTC_EDX_MTRR 0x00001000 /* memory type range reg */
67 #define CPUID_INTC_EDX_PGE 0x00002000 /* page global enable */
68 #define CPUID_INTC_EDX_MCA 0x00004000 /* machine check arch */
69 #define CPUID_INTC_EDX_CMOV 0x00008000 /* conditional move insns */
70 #define CPUID_INTC_EDX_PAT 0x00010000 /* page attribute table */
71 #define CPUID_INTC_EDX_PSE36 0x00020000 /* 36-bit pagesize extension */
72 #define CPUID_INTC_EDX_PSN 0x00040000 /* processor serial number */
73 #define CPUID_INTC_EDX_CLFSH 0x00080000 /* clflush instruction */
74 /* 0x10000 - reserved */
75 #define CPUID_INTC_EDX_DS 0x00200000 /* debug store exists */
76 #define CPUID_INTC_EDX_ACPI 0x00400000 /* monitoring + clock ctrl */
77 #define CPUID_INTC_EDX_MMX 0x00800000 /* MMX instructions */
78 #define CPUID_INTC_EDX_FXSR 0x01000000 /* fxsave and fxrstor */
79 #define CPUID_INTC_EDX_SSE 0x02000000 /* streaming SIMD extensions */
80 #define CPUID_INTC_EDX_SSE2 0x04000000 /* SSE extensions */
81 #define CPUID_INTC_EDX_SS 0x08000000 /* self-snoop */
82 #define CPUID_INTC_EDX_HTT 0x10000000 /* Hyper Thread Technology */
83 #define CPUID_INTC_EDX_TM 0x20000000 /* thermal monitoring */
84 #define CPUID_INTC_EDX_IA64 0x40000000 /* Itanium emulating IA32 */
85 #define CPUID_INTC_EDX_PBE 0x80000000 /* Pending Break Enable */

87 #define FMT_CPUID_INTC_EDX \
88 " \\"20" \
89 " \\"40pbe\37ia64\36tm\35htt\34ss\33sse2\32sse\31fxsr" \
90 " \\"30mmx\27acpi\26ds\24clfsh\23psn\22pse36\21pat" \
91 " \\"20cmov\17mca\16pge\15mtrr\14sep\12apic\11cx8" \
92 " \\"10mce\7pae\6msr\5tsc\4pse\3de\2vme\1fpu" \
93 \
94 /* \
95 * cpuid instruction feature flags in %ecx (standard function 1)
96 */

98 #define CPUID_INTC_ECX_SSE3 0x00000001 /* Yet more SSE extensions */
99 #define CPUID_INTC_ECX_PCLMULQDQ 0x00000002 /* PCLMULQDQ insn */
100 /* 0x00000004 - reserved */
101 #define CPUID_INTC_ECX_MON 0x00000008 /* MONITOR/MWAIT */
102 #define CPUID_INTC_ECX_DSCP 0x00000010 /* CPL-qualified debug store */
103 #define CPUID_INTC_ECX_VMX 0x00000020 /* Hardware VM extensions */
104 #define CPUID_INTC_ECX_SMX 0x00000040 /* Secure mode extensions */
105 #define CPUID_INTC_ECX_EST 0x00000080 /* enhanced SpeedStep */
106 #define CPUID_INTC_ECX_TM2 0x00000100 /* thermal monitoring */
107 #define CPUID_INTC_ECX_SSSE3 0x00000200 /* Supplemental SSE3 insns */
108 #define CPUID_INTC_ECX_CID 0x00000400 /* L1 context ID */
109 /* 0x00000800 - reserved */
110 /* 0x00001000 - reserved */
111 #define CPUID_INTC_ECX_CX16 0x00002000 /* cmpxchg16 */
112 #define CPUID_INTC_ECX_ETPRD 0x00004000 /* extended task pri messages */
113 /* 0x00008000 - reserved */
114 /* 0x00010000 - reserved */
115 /* 0x00020000 - reserved */
116 #define CPUID_INTC_ECX_DCA 0x00040000 /* direct cache access */
117 #define CPUID_INTC_ECX_SSE4_1 0x00080000 /* SSE4.1 insns */
118 #define CPUID_INTC_ECX_SSE4_2 0x00100000 /* SSE4.2 insns */
119 #define CPUID_INTC_ECX_X2APIC 0x00200000 /* x2APIC */

120 #endif /* ! codereview */
121 #define CPUID_INTC_ECX_MOVBE 0x00400000 /* MOVBE insn */
122 #define CPUID_INTC_ECX_POPCNT 0x00800000 /* POPCNT insn */
123 #define CPUID_INTC_ECX_AES 0x02000000 /* AES insns */
124 #define CPUID_INTC_ECX_XSAVE 0x04000000 /* XSAVE/XRESTOR insns */
125 #define CPUID_INTC_ECX_OSXSAVE 0x08000000 /* OS supports XSAVE insns */
126 #define CPUID_INTC_ECX_AVX 0x10000000 /* AVX supported */
```

```

127 #define CPUID_INTC_ECX_F16C 0x20000000 /* F16C supported */
128 #define CPUID_INTC_ECX_RDRAND 0x40000000 /* RDRAND supported */
129 #define CPUID_INTC_ECX_HV 0x80000000 /* Hypervisor */

131 #define FMT_CPUID_INTC_ECX \
132 "20" \
133 "\37rdrand\36f16c\35avx\34osxsav\33xsave" \
134 "\32aes" \
135 "\30popcnt\27movbe\26x2apic\25sse4.2\24sse4.1\23dca" \
119 "\30popcnt\27movbe\25sse4.2\24sse4.1\23dca" \
136 "\20\17etprd\16cx16\13cid\12ssse3\11tm2" \
137 "\10est\7smx\6vmx\5dscpl\4mon\2pcmulqdq\1sse3" \
138 /*

140 * cpuid instruction feature flags in %edx (extended function 0x80000001)
141 */

143 #define CPUID_AMD_EDX_FPU 0x00000001 /* x87 fpu present */
144 #define CPUID_AMD_EDX_VME 0x00000002 /* virtual-8086 extension */
145 #define CPUID_AMD_EDX_DE 0x00000004 /* debugging extensions */
146 #define CPUID_AMD_EDX_PSE 0x00000008 /* page size extensions */
147 #define CPUID_AMD_EDX_TSC 0x00000010 /* time stamp counter */
148 #define CPUID_AMD_EDX_MSR 0x00000020 /* rdmrs and wrmsr */
149 #define CPUID_AMD_EDX_PAE 0x00000040 /* physical addr extension */
150 #define CPUID_AMD_EDX_MCE 0x00000080 /* machine check exception */
151 #define CPUID_AMD_EDX_CX8 0x00000100 /* cmpxchg8b instruction */
152 #define CPUID_AMD_EDX_APIC 0x00000200 /* local APIC */
153 /* 0x00000400 - sysc on K6m6 */
154 #define CPUID_AMD_EDX_SYSC 0x00000800 /* AMD: syscall and sysret */
155 #define CPUID_AMD_EDX_MTRR 0x00001000 /* memory type and range reg */
156 #define CPUID_AMD_EDX_PGE 0x00002000 /* page global enable */
157 #define CPUID_AMD_EDX_MCA 0x00004000 /* machine check arch */
158 #define CPUID_AMD_EDX_CMov 0x00008000 /* conditional move insns */
159 #define CPUID_AMD_EDX_PAT 0x00010000 /* K7: page attribute table */
160 #define CPUID_AMD_EDX_FCMOV 0x00010000 /* FCMOVcc etc. */
161 #define CPUID_AMD_EDX_PSE36 0x00020000 /* 36-bit pagesize extension */
162 /* 0x00040000 - reserved */
163 /* 0x00080000 - reserved */
164 #define CPUID_AMD_EDX_NX 0x00100000 /* AMD: no-execute page prot */
165 /* 0x00200000 - reserved */
166 #define CPUID_AMD_EDX_MMXamd 0x00400000 /* MMX extensions */
167 #define CPUID_AMD_EDX_MMX 0x00800000 /* MMX instructions */
168 #define CPUID_AMD_EDX_FXSR 0x01000000 /* fxsave and fxrstor */
169 #define CPUID_AMD_EDX_FFXSR 0x02000000 /* fast fxsave/fxrstor */
170 #define CPUID_AMD_EDX_1GPG 0x04000000 /* 1GB page */
171 #define CPUID_AMD_EDX_TSCP 0x08000000 /* rdtscp instruction */
172 /* 0x10000000 - reserved */
173 #define CPUID_AMD_EDX_LM 0x20000000 /* AMD: long mode */
174 #define CPUID_AMD_EDX_3DNowx 0x40000000 /* AMD: extensions to 3DNow! */
175 #define CPUID_AMD_EDX_3DNow 0x80000000 /* AMD: 3DNow! instructions */

177 #define FMT_CPUID_AMD_EDX \
178 "20" \
179 "\40a3d\37a3d+\36lm\34tscp\32ffxsr\31fxsr" \
180 "\30mmx\27mmxext\25nx\2pse\21pat" \
181 "\20cmov\17mca\16pge\15mtr\14syscall\12apic\11cx8" \
182 "\10mce\7pae\6msr\5tsc\4pse\3de\2vme\1fp" \
183 /*

184 #define CPUID_AMD_ECX_AHF64 0x00000001 /* LAHF and SAHF in long mode */
185 #define CPUID_AMD_ECX_CMP_LGCY 0x00000002 /* AMD: multicore chip */
186 #define CPUID_AMD_ECX_SVM 0x00000004 /* AMD: secure VM */
187 #define CPUID_AMD_ECX_EAS 0x00000008 /* extended apic space */
188 #define CPUID_AMD_ECX_CR8D 0x00000010 /* AMD: 32-bit mov %cr8 */
189 #define CPUID_AMD_ECX_LZCNT 0x00000020 /* AMD: LZCNT insn */
190 #define CPUID_AMD_ECX_SSE4A 0x00000040 /* AMD: SSE4A insns */
191 #define CPUID_AMD_ECX_MAS 0x00000080 /* AMD: MisAlignSse mnode */

```

```

192 #define CPUID_AMD_ECX_3DNP 0x00000100 /* AMD: 3DNowPrefetch */
193 #define CPUID_AMD_ECX_OSVW 0x00000200 /* AMD: OSWV */
194 #define CPUID_AMD_ECX_IBS 0x00000400 /* AMD: IBS */
195 #define CPUID_AMD_ECX_SSE5 0x00000800 /* AMD: SSE5 */
196 #define CPUID_AMD_ECX_SKINIT 0x00001000 /* AMD: SKINIT */
197 #define CPUID_AMD_ECX_WDT 0x00002000 /* AMD: WDT */
198 #define CPUID_AMD_ECX_TOPOEXT 0x00400000 /* AMD: Topology Extensions */

200 #define FMT_CPUID_AMD_ECX \
201 "20" \
202 "\22topoext" \
203 "\14wdt\13skinit\12sse5\11ibs\10osvw\93dnp\8mas" \
204 "\7sse4a\6lzcnt\5cr8d\3svm\2lcmplgcy\lahf64" \
205 /*

207 * Intel now seems to have claimed part of the "extended" function
208 * space that we previously reserved for non-Intel implementors to use.
209 * More excitingly still, they've claimed bit 20 to mean LAHF/SAHF
210 * is available in long mode i.e. what AMD indicate using bit 0.
211 * On the other hand, everything else is labelled as reserved.
212 */
213 #define CPUID_INTC_ECX_AHF64 0x00100000 /* LAHF and SAHF in long mode */

216 #define P5_MCHADDR 0x0
217 #define P5_CESR 0x11
218 #define P5_CTR0 0x12
219 #define P5_CTR1 0x13

221 #define K5_MCHADDR 0x0
222 #define K5_MCHTYPE 0x01
223 #define K5_TSC 0x10
224 #define K5_TR12 0x12

226 #define REG_PAT 0x277

228 #define REG_MCO_CTL 0x400
229 #define REG_MCS_MISC 0x417
230 #define REG_PERFCTR0 0xc1
231 #define REG_PERFCTR1 0xc2

233 #define REG_PERFEVNT0 0x186
234 #define REG_PERFEVNT1 0x187

236 #define REG_TSC 0x10 /* timestamp counter */
237 #define REG_APIC_BASE_MSR 0x1b
238 #define REG_X2APIC_BASE_MSR 0x800 /* The MSR address offset of x2APIC */

240 #if !defined(__xpv)
241 /*
242 * AMD C1E
243 */
244 #define MSR_AMD_INT_PENDING_CMP_HALT 0xC0010055
245 #define AMD_ACTONCMPHALT_SHIFT 27
246 #define AMD_ACTONCMPHALT_MASK 3
247#endif

249 #define MSR_DEBUGCTL 0x1d9

251 #define DEBUGCTL_LBR 0x01
252 #define DEBUGCTL_BTF 0x02

254 /* Intel P6, AMD */
255 #define MSR_LBR_FROM 0x1db
256 #define MSR_LBR_TO 0x1dc
257 #define MSR_LEX_FROM 0x1dd

```

```

258 #define MSR_LEX_TO          0x1de
260 /* Intel P4 (pre-Prescott, non P4 M) */
261 #define MSR_P4_LBSTK_TOS      0x1da
262 #define MSR_P4_LBSTK_0         0x1db
263 #define MSR_P4_LBSTK_1         0x1dc
264 #define MSR_P4_LBSTK_2         0x1dd
265 #define MSR_P4_LBSTK_3         0x1de
267 /* Intel Pentium M */
268 #define MSR_P6M_LBSTK_TOS      0x1c9
269 #define MSR_P6M_LBSTK_0         0x040
270 #define MSR_P6M_LBSTK_1         0x041
271 #define MSR_P6M_LBSTK_2         0x042
272 #define MSR_P6M_LBSTK_3         0x043
273 #define MSR_P6M_LBSTK_4         0x044
274 #define MSR_P6M_LBSTK_5         0x045
275 #define MSR_P6M_LBSTK_6         0x046
276 #define MSR_P6M_LBSTK_7         0x047
278 /* Intel P4 (Prescott) */
279 #define MSR_PRP4_LBSTK_TOS      0x1da
280 #define MSR_PRP4_LBSTK_FROM_0    0x680
281 #define MSR_PRP4_LBSTK_FROM_1    0x681
282 #define MSR_PRP4_LBSTK_FROM_2    0x682
283 #define MSR_PRP4_LBSTK_FROM_3    0x683
284 #define MSR_PRP4_LBSTK_FROM_4    0x684
285 #define MSR_PRP4_LBSTK_FROM_5    0x685
286 #define MSR_PRP4_LBSTK_FROM_6    0x686
287 #define MSR_PRP4_LBSTK_FROM_7    0x687
288 #define MSR_PRP4_LBSTK_FROM_8    0x688
289 #define MSR_PRP4_LBSTK_FROM_9    0x689
290 #define MSR_PRP4_LBSTK_FROM_10   0x68a
291 #define MSR_PRP4_LBSTK_FROM_11   0x68b
292 #define MSR_PRP4_LBSTK_FROM_12   0x68c
293 #define MSR_PRP4_LBSTK_FROM_13   0x68d
294 #define MSR_PRP4_LBSTK_FROM_14   0x68e
295 #define MSR_PRP4_LBSTK_FROM_15   0x68f
296 #define MSR_PRP4_LBSTK_TO_0       0x6c0
297 #define MSR_PRP4_LBSTK_TO_1       0x6c1
298 #define MSR_PRP4_LBSTK_TO_2       0x6c2
299 #define MSR_PRP4_LBSTK_TO_3       0x6c3
300 #define MSR_PRP4_LBSTK_TO_4       0x6c4
301 #define MSR_PRP4_LBSTK_TO_5       0x6c5
302 #define MSR_PRP4_LBSTK_TO_6       0x6c6
303 #define MSR_PRP4_LBSTK_TO_7       0x6c7
304 #define MSR_PRP4_LBSTK_TO_8       0x6c8
305 #define MSR_PRP4_LBSTK_TO_9       0x6c9
306 #define MSR_PRP4_LBSTK_TO_10      0x6ca
307 #define MSR_PRP4_LBSTK_TO_11      0x6cb
308 #define MSR_PRP4_LBSTK_TO_12      0x6cc
309 #define MSR_PRP4_LBSTK_TO_13      0x6cd
310 #define MSR_PRP4_LBSTK_TO_14      0x6ce
311 #define MSR_PRP4_LBSTK_TO_15      0x6cf
313 #define MCI_CTL_VALUE           0xffffffff
315 #define MTRR_TYPE_UC             0
316 #define MTRR_TYPE_WC             1
317 #define MTRR_TYPE_WT             4
318 #define MTRR_TYPE_WP             5
319 #define MTRR_TYPE_WB             6
320 #define MTRR_TYPE_UC_             7
322 /*
323 * For Solaris we set up the page attribute table in the following way:

```

```

324 * PAT0 Write-Back
325 * PAT1 Write-Through
326 * PAT2 Uncacheable-
327 * PAT3 Uncacheable
328 * PAT4 Write-Back
329 * PAT5 Write-Through
330 * PAT6 Write-Combine
331 * PAT7 Uncacheable
332 * The only difference from h/w default is entry 6.
333 */
334 #define PAT_DEFAULT_ATTRIBUTE \
335 (((uint64_t)MTRR_TYPE_WB | \
336 ((uint64_t)MTRR_TYPE_WT << 8) | \
337 ((uint64_t)MTRR_TYPE_UC_ << 16) | \
338 ((uint64_t)MTRR_TYPE_UC << 24) | \
339 ((uint64_t)MTRR_TYPE_WB << 32) | \
340 ((uint64_t)MTRR_TYPE_WT << 40) | \
341 ((uint64_t)MTRR_TYPE_WC << 48) | \
342 ((uint64_t)MTRR_TYPE_UC << 56)) \
344 #define X86FSET_LARGEPAGE      0
345 #define X86FSET_TSC            1
346 #define X86FSET_MSR            2
347 #define X86FSET_MTRR           3
348 #define X86FSET_PGE            4
349 #define X86FSET_DE             5
350 #define X86FSET_CMOV           6
351 #define X86FSET_MMX            7
352 #define X86FSET_MCA            8
353 #define X86FSET_PAE            9
354 #define X86FSET_CX8             10
355 #define X86FSET_PAT            11
356 #define X86FSET_SEP             12
357 #define X86FSET_SSE             13
358 #define X86FSET_SSE2            14
359 #define X86FSET_HTT             15
360 #define X86FSET_ASYSC           16
361 #define X86FSET_NX              17
362 #define X86FSET_SSE3            18
363 #define X86FSET_CX16            19
364 #define X86FSET_CMP             20
365 #define X86FSET_TSCP            21
366 #define X86FSET_MWAIT           22
367 #define X86FSET_SSE4A           23
368 #define X86FSET_CPUID           24
369 #define X86FSET_SSSE3           25
370 #define X86FSET_SSE4_1           26
371 #define X86FSET_SSE4_2           27
372 #define X86FSET_1GPG             28
373 #define X86FSET_CLFSH           29
374 #define X86FSET_64              30
375 #define X86FSET_AES              31
376 #define X86FSET_PCLMULQDQ        32
377 #define X86FSET_XSAVE            33
378 #define X86FSET_AVX              34
379 #define X86FSET_VMX              35
380 #define X86FSET_SVM              36
381 #define X86FSET_TOPOEXT          37
382 #define X86FSET_F16C             38
383 #define X86FSET_RDRAND           39
384 #define X86FSET_X2APIC           40
385 #endif /* ! codereview */

387 /*
388 * flags to patch tsc_read routine.
389 */

```

```

390 #define X86_NO_TSC          0x0
391 #define X86_HAVE_TSCP        0x1
392 #define X86_TSC_MFENCE       0x2
393 #define X86_TSC_LFENCE       0x4

395 /*
396  * Intel Deep C-State invariant TSC in leaf 0x80000007.
397 */
398 #define CPUID_TSC_CSTATE_INVARIAENCE    (0x100)

400 /*
401  * Intel Deep C-state always-running local APIC timer
402 */
403 #define CPUID_CSTATE_ARAT      (0x4)

405 /*
406  * Intel ENERGY_PERF_BIAS MSR indicated by feature bit CPUID.6.ECX[3].
407 */
408 #define CPUID_EPB_SUPPORT     (1 << 3)

410 /*
411  * Intel TSC deadline timer
412 */
413 #define CPUID_DEADLINE_TSC    (1 << 24)

415 /*
416  * x86_type is a legacy concept; this is supplanted
417  * for most purposes by x86_featureset; modern CPUs
418  * should be X86_TYPE_OTHER
419 */
420 #define X86_TYPE_OTHER        0
421 #define X86_TYPE_486           1
422 #define X86_TYPE_P5            2
423 #define X86_TYPE_P6            3
424 #define X86_TYPE_CYRIX_486     4
425 #define X86_TYPE_CYRIX_6x86L   5
426 #define X86_TYPE_CYRIX_6x86    6
427 #define X86_TYPE_CYRIX_GXm     7
428 #define X86_TYPE_CYRIX_6x86MX   8
429 #define X86_TYPE_CYRIX_MediaGX 9
430 #define X86_TYPE_CYRIX_MII     10
431 #define X86_TYPE_VIA_CYRIX_III 11
432 #define X86_TYPE_P4            12

434 /*
435  * x86_vendor allows us to select between
436  * implementation features and helps guide
437  * the interpretation of the cpuid instruction.
438 */
439 #define X86_VENDOR_Intel       0
440 #define X86_VENDORSTR_Intel    "GenuineIntel"

442 #define X86_VENDOR_IntelClone 1

444 #define X86_VENDOR_AMD         2
445 #define X86_VENDORSTR_AMD     "AuthenticAMD"

447 #define X86_VENDOR_Cyrix        3
448 #define X86_VENDORSTR_CYRIX   "CyrixInstead"

450 #define X86_VENDOR_UMC          4
451 #define X86_VENDORSTR_UMC     "UMC UMC UMC "

453 #define X86_VENDOR_NexGen      5
454 #define X86_VENDORSTR_NexGen   "NexGenDriven"

```

```

456 #define X86_VENDOR_Centaur      6
457 #define X86_VENDORSTR_Centaur  "CentaurHauls"

459 #define X86_VENDOR_Rise         7
460 #define X86_VENDORSTR_Rise    "RiseRiseRise"

462 #define X86_VENDOR_SIS          8
463 #define X86_VENDORSTR_SIS     "SIS SIS SIS"

465 #define X86_VENDOR_TM          9
466 #define X86_VENDORSTR_TM     "GenuineTMx86"

468 #define X86_VENDOR_NSC          10
469 #define X86_VENDORSTR_NSC    "Geode by NSC"

471 /*
472  * Vendor string max len + \0
473 */
474 #define X86_VENDOR_STRLEN       13

476 /*
477  * Some vendor/family/model/stepping ranges are commonly grouped under
478  * a single identifying banner by the vendor. The following encode
479  * that "revision" in a uint32_t with the 8 most significant bits
480  * identifying the vendor with X86_VENDOR_*, the next 8 identifying the
481  * family, and the remaining 16 typically forming a bitmask of revisions
482  * within that family with more significant bits indicating "later" revisions.
483 */

485 #define _X86_CHIPREV_VENDOR_MASK 0xff000000u
486 #define _X86_CHIPREV_VENDOR_SHIFT 24
487 #define _X86_CHIPREV_FAMILY_MASK 0x0ff0000u
488 #define _X86_CHIPREV_FAMILY_SHIFT 16
489 #define _X86_CHIPREV_REV_MASK 0x0000ffffu

491 #define _X86_CHIPREV_VENDOR(x) \
492 (((x) & _X86_CHIPREV_VENDOR_MASK) >> _X86_CHIPREV_VENDOR_SHIFT)
493 #define _X86_CHIPREV_FAMILY(x) \
494 (((x) & _X86_CHIPREV_FAMILY_MASK) >> _X86_CHIPREV_FAMILY_SHIFT)
495 #define _X86_CHIPREV_REV(x) \
496 (((x) & _X86_CHIPREV_REV_MASK))

498 /* True if x matches in vendor and family and if x matches the given rev mask */
499 #define X86_CHIPREV_MATCH(x, mask) \
500  (_X86_CHIPREV_VENDOR(x) == _X86_CHIPREV_VENDOR(mask) && \
501  _X86_CHIPREV_FAMILY(x) == _X86_CHIPREV_FAMILY(mask) && \
502  (( _X86_CHIPREV_REV(x) & _X86_CHIPREV_REV(mask)) != 0))

504 /* True if x matches in vendor and family, and rev is at least minx */
505 #define X86_CHIPREV_ATLEAST(x, minx) \
506  (_X86_CHIPREV_VENDOR(x) == _X86_CHIPREV_VENDOR(minx) && \
507  _X86_CHIPREV_FAMILY(x) == _X86_CHIPREV_FAMILY(minx) && \
508  _X86_CHIPREV_REV(x) >= _X86_CHIPREV_REV(minx))

510 #define _X86_CHIPREV_MKREV(vendor, family, rev) \
511  ((uint32_t)(vendor) << _X86_CHIPREV_VENDOR_SHIFT | \
512  (family) << _X86_CHIPREV_FAMILY_SHIFT | (rev))

514 /* True if x matches in vendor, and family is at least minx */
515 #define X86_CHIPFAM_ATLEAST(x, minx) \
516  (_X86_CHIPREV_VENDOR(x) == _X86_CHIPREV_VENDOR(minx) && \
517  _X86_CHIPREV_FAMILY(x) >= _X86_CHIPREV_FAMILY(minx))

519 /* Revision default */
520 #define X86_CHIPREV_UNKNOWN     0x0

```

```

522 /*
523  * Definitions for AMD Family 0xf. Minor revisions C0 and CG are
524  * sufficiently different that we will distinguish them; in all other
525  * case we will identify the major revision.
526 */
527 #define X86_CHIPREV_AMD_F_REV_B _X86_CHIPREV_MKREV(X86_VENDOR_AMD, 0xf, 0x0001)
528 #define X86_CHIPREV_AMD_F_REV_C0 _X86_CHIPREV_MKREV(X86_VENDOR_AMD, 0xf, 0x0002)
529 #define X86_CHIPREV_AMD_F_REV(CG _X86_CHIPREV_MKREV(X86_VENDOR_AMD, 0xf, 0x0004)
530 #define X86_CHIPREV_AMD_F_REV_D _X86_CHIPREV_MKREV(X86_VENDOR_AMD, 0xf, 0x0008)
531 #define X86_CHIPREV_AMD_F_REV_E _X86_CHIPREV_MKREV(X86_VENDOR_AMD, 0xf, 0x0010)
532 #define X86_CHIPREV_AMD_F_REV_F _X86_CHIPREV_MKREV(X86_VENDOR_AMD, 0xf, 0x0020)
533 #define X86_CHIPREV_AMD_F_REV_G _X86_CHIPREV_MKREV(X86_VENDOR_AMD, 0xf, 0x0040)

535 /*
536  * Definitions for AMD Family 0x10. Rev A was Engineering Samples only.
537 */
538 #define X86_CHIPREV_AMD_10_REV_A \
539     _X86_CHIPREV_MKREV(X86_VENDOR_AMD, 0x10, 0x0001)
540 #define X86_CHIPREV_AMD_10_REV_B \
541     _X86_CHIPREV_MKREV(X86_VENDOR_AMD, 0x10, 0x0002)
542 #define X86_CHIPREV_AMD_10_REV_C2 \
543     _X86_CHIPREV_MKREV(X86_VENDOR_AMD, 0x10, 0x0004)
544 #define X86_CHIPREV_AMD_10_REV_C3 \
545     _X86_CHIPREV_MKREV(X86_VENDOR_AMD, 0x10, 0x0008)
546 #define X86_CHIPREV_AMD_10_REV_D0 \
547     _X86_CHIPREV_MKREV(X86_VENDOR_AMD, 0x10, 0x0010)
548 #define X86_CHIPREV_AMD_10_REV_D1 \
549     _X86_CHIPREV_MKREV(X86_VENDOR_AMD, 0x10, 0x0020)
550 #define X86_CHIPREV_AMD_10_REV_E \
551     _X86_CHIPREV_MKREV(X86_VENDOR_AMD, 0x10, 0x0040)

553 /*
554  * Definitions for AMD Family 0x11.
555 */
556 #define X86_CHIPREV_AMD_11_REV_B \
557     _X86_CHIPREV_MKREV(X86_VENDOR_AMD, 0x11, 0x0002)
559 /*
560  * Definitions for AMD Family 0x12.
561 */
562 #define X86_CHIPREV_AMD_12_REV_B \
563     _X86_CHIPREV_MKREV(X86_VENDOR_AMD, 0x12, 0x0002)

565 /*
566  * Definitions for AMD Family 0x14.
567 */
568 #define X86_CHIPREV_AMD_14_REV_B \
569     _X86_CHIPREV_MKREV(X86_VENDOR_AMD, 0x14, 0x0002)
570 #define X86_CHIPREV_AMD_14_REV_C \
571     _X86_CHIPREV_MKREV(X86_VENDOR_AMD, 0x14, 0x0004)

573 /*
574  * Definitions for AMD Family 0x15
575 */
576 #define X86_CHIPREV_AMD_15OR_REV_B2 \
577     _X86_CHIPREV_MKREV(X86_VENDOR_AMD, 0x15, 0x0001)

579 #define X86_CHIPREV_AMD_15TN_REV_A1 \
580     _X86_CHIPREV_MKREV(X86_VENDOR_AMD, 0x15, 0x0002)

582 /*
583  * Various socket/package types, extended as the need to distinguish
584  * a new type arises. The top 8 byte identifies the vendor and the
585  * remaining 24 bits describe 24 socket types.
586 */

```

```

588 #define _X86_SOCKET_VENDOR_SHIFT 24
589 #define _X86_SOCKET_VENDOR(x) ((x) >> _X86_SOCKET_VENDOR_SHIFT)
590 #define _X86_SOCKET_TYPE_MASK 0x00ffff
591 #define _X86_SOCKET_TYPE(x) ((x) & _X86_SOCKET_TYPE_MASK)

593 #define _X86_SOCKET_MKVAL(vendor, bitval) \
594     ((uint32_t)(vendor) << _X86_SOCKET_VENDOR_SHIFT | (bitval))

596 #define X86_SOCKET_MATCH(s, mask) \
597     (_X86_SOCKET_VENDOR(s) == _X86_SOCKET_VENDOR(mask) && \
598     (_X86_SOCKET_TYPE(s) & _X86_SOCKET_TYPE(mask)) != 0)

600 #define X86_SOCKET_UNKNOWN 0x0
601 /* *
602  * AMD socket types
603 */
604 #define X86_SOCKET_754 _X86_SOCKET_MKVAL(X86_VENDOR_AMD, 0x000001)
605 #define X86_SOCKET_939 _X86_SOCKET_MKVAL(X86_VENDOR_AMD, 0x000002)
606 #define X86_SOCKET_940 _X86_SOCKET_MKVAL(X86_VENDOR_AMD, 0x000004)
607 #define X86_SOCKET_S1g1 _X86_SOCKET_MKVAL(X86_VENDOR_AMD, 0x000008)
608 #define X86_SOCKET_AM2 _X86_SOCKET_MKVAL(X86_VENDOR_AMD, 0x000010)
609 #define X86_SOCKET_F1207 _X86_SOCKET_MKVAL(X86_VENDOR_AMD, 0x000020)
610 #define X86_SOCKET_S1g2 _X86_SOCKET_MKVAL(X86_VENDOR_AMD, 0x000040)
611 #define X86_SOCKET_S1g3 _X86_SOCKET_MKVAL(X86_VENDOR_AMD, 0x000080)
612 #define X86_SOCKET_AM _X86_SOCKET_MKVAL(X86_VENDOR_AMD, 0x000100)
613 #define X86_SOCKET_AM2R2 _X86_SOCKET_MKVAL(X86_VENDOR_AMD, 0x000200)
614 #define X86_SOCKET_AM3 _X86_SOCKET_MKVAL(X86_VENDOR_AMD, 0x000400)
615 #define X86_SOCKET_G34 _X86_SOCKET_MKVAL(X86_VENDOR_AMD, 0x000800)
616 #define X86_SOCKET_ASB2 _X86_SOCKET_MKVAL(X86_VENDOR_AMD, 0x001000)
617 #define X86_SOCKET_C32 _X86_SOCKET_MKVAL(X86_VENDOR_AMD, 0x002000)
618 #define X86_SOCKET_S1g4 _X86_SOCKET_MKVAL(X86_VENDOR_AMD, 0x004000)
619 #define X86_SOCKET_FT1 _X86_SOCKET_MKVAL(X86_VENDOR_AMD, 0x008000)
620 #define X86_SOCKET_FM1 _X86_SOCKET_MKVAL(X86_VENDOR_AMD, 0x010000)
621 #define X86_SOCKET_FS1 _X86_SOCKET_MKVAL(X86_VENDOR_AMD, 0x020000)
622 #define X86_SOCKET_AM3R2 _X86_SOCKET_MKVAL(X86_VENDOR_AMD, 0x040000)
623 #define X86_SOCKET_PP2 _X86_SOCKET_MKVAL(X86_VENDOR_AMD, 0x080000)
624 #define X86_SOCKET_FS1R2 _X86_SOCKET_MKVAL(X86_VENDOR_AMD, 0x100000)
625 #define X86_SOCKET_FM2 _X86_SOCKET_MKVAL(X86_VENDOR_AMD, 0x200000)

627 /*
628  * xgetbv/xsetbv support
629 */
631 #define XFEATURE_ENABLED_MASK 0x0
632 /*
633  * XFEATURE_ENABLED_MASK values (eax)
634 */
635 #define XFEATURE_LEGACY_FP 0x1
636 #define XFEATURE_SSE 0x2
637 #define XFEATURE_AVX 0x4
638 #define XFEATURE_MAX XFEATURE_AVX
639 #define XFEATURE_FP_ALL \
640     (XFEATURE_LEGACY_FP|XFEATURE_SSE|XFEATURE_AVX)

642 #if !defined(_ASM)
644 #if defined(_KERNEL) || defined(_KMEMUSER)
646 #define NUM_X86_FEATURES 41
648 #define NUM_X86_FEATURES 40
647 extern uchar_t x86_featureset[];
649 extern void free_x86_featureset(void *featureset);
650 extern boolean_t is_x86_feature(void *featureset, uint_t feature);
651 extern void add_x86_feature(void *featureset, uint_t feature);
652 extern void remove_x86_feature(void *featureset, uint_t feature);

```

```
653 extern boolean_t compare_x86_featureset(void *setA, void *setB);
654 extern void print_x86_featureset(void *featureset);

657 extern uint_t x86_type;
658 extern uint_t x86_vendor;
659 extern uint_t x86_clflush_size;

661 extern uint_t pentiumpro_bug4046376;
384 extern uint_t pentiumpro_bug4064495;

386 extern uint_t enable486;

663 extern const char CyrixInstead[];

665 #endif

667 #if defined(_KERNEL)

669 /*
670  * This structure is used to pass arguments and get return values back
671  * from the CPUID instruction in __cpuid_insn() routine.
672 */
673 struct cpuid_regs {
674     uint32_t cp_eax;
675     uint32_t cp_ebx;
676     uint32_t cp_ecx;
677     uint32_t cp_edx;
678 };


---

unchanged portion omitted
```

```
new/usr/src/uts/sun4/os/mp_states.c
```

```
1
```

```
*****  
6667 Fri Jan 3 22:11:57 2014  
new/usr/src/uts/sun4/os/mp_states.c  
patch cpu-pause-func-deglobalize  
*****  
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     }  
199     if (ncpus > 1)  
200         intr_redist_all_cpus_shutdown();  
201  
202     pause_cpus(NULL, NULL);  
203     pause_cpus(NULL);  
204     cpu_are_paused = 1;  
205     mutex_exit(&cpu_lock);  
206 }
```

unchanged_portion_omitted

```
*****  
16840 Fri Jan 3 22:11:57 2014  
new/usr/src/uts/sun4/os/prom_subr.c  
patch cpu-pause-func-deglobalize  
*****  
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);  
413     pause_cpus(NULL);  
  
415     /* If some other cpu is entering or is in the prom, spin */  
416     while (prom_cpu || mutex_owner(&prom_mutex)) {  
  
418         start_cpus();  
419         mutex_enter(&prom_mutex);  
  
421         /* Wait for other cpu to exit prom */  
422         while (prom_cpu)  
423             cv_wait(&prom_cv, &prom_mutex);  
  
425         mutex_exit(&prom_mutex);  
426         pause_cpus(NULL, NULL);  
426         pause_cpus(NULL);  
427     }  
  
429     /* At this point all cpus are paused and none are in the prom */  
430 }
```

unchanged_portion_omitted

new/usr/src/uts/sun4u/io/mem_cache.c

```
*****
24533 Fri Jan 3 22:11:57 2014
new/usr/src/uts/sun4u/io/mem_cache.c
patch cpu-pause-func-deglobalize
*****
_____ 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_ECSTATE_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 =
```

1

new/usr/src/uts/sun4u/io/mem_cache.c

```
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++;
/*
    * 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_ECSTATE_MASK;
if (state != PN_ECSTATE_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(
```

2

new/usr/src/uts/sun4u/io/mem_cache.c

3

```

678                                     tag_addr, pattern);
679
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         }
699     /*
700      * We bind ourself to a CPU and send cross trap to
701      * ourselves. On return from xt_one we can rely on the
702      * data in tag_data being filled in. Normally one would
703      * do a xt_sync to make sure that the CPU has completed
704      * the cross trap call xt_one.
705     */
706     xt_one(cache_info->cpu_id,
707            (xfcfunc_t *) (get_l3_tag_t11),
708            tag_addr, (uint64_t) (&tag_data));
709     state = tag_data & CH_ECSTATE_MASK;
710     if (state != PN_ECSTATE_NA) {
711         cmn_err(CE_WARN,
712                 "L3 RETIRE failed for index 0x%x\n",
713                 cache_info->index);
714     retire_failures++;
715     affinity_clear();
716     return (EIO);
717 }
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                (xfcfunc_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);

```

```

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      * We bind ourself to a CPU and send cross trap to
824      * ourself. On return from xt_one we can rely on the
825      * data in clop being filled in. Normally one would
826      * do a xt_sync to make sure that the CPU has completed
827      * the cross trap call xt_one.
828      */
829     xt_one(cache_info->cpu_id,
830           (xcfunc_t *) (get_ecache_dtags_t11),
831           afar, (uint64_t) (&clop));
832     mutex_enter(&cpu_lock);
833     (void) start_cpus();
834     mutex_exit(&cpu_lock);
835     affinity_clear();
836     switch (cache_info->cache) {
837         case L2_CACHE_TAG:
838             for (i = 0; i < PN_CACHE_NWAYS; i++) {
839                 Lxcache_tag_data[i] =
840                     clop.clo_data.chd_l2_data
841 #ifdef DEBUG
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
849                     break;
850         case L3_CACHE_TAG:
851             for (i = 0; i < PN_CACHE_NWAYS; i++) {
852                 Lxcache_tag_data[i] =
853                     clop.clo_data.chd_ec_data
854 #ifdef DEBUG
855                     [i].ec_tag;
856
857                     last_error_injected_bit =
858                         last_l3tag_error_injected_bit;
859                     last_error_injected_way =
860                         last_l3tag_error_injected_way;
861
862                     break;
863     default:
864         return (ENOTSUP);
865     } /* end if switch(cache) */
866
867     if ((cmd == MEM_CACHE_READ_ERROR_INJECTED_TAGS) &&
868         (inject_anonymous_tag_error == 0) &&
869         (last_error_injected_way >= 0) &&
870         (last_error_injected_way <= 3)) {
871         pattern = ((uint64_t)1 <<
872                     last_error_injected_bit);
873
874         /*
875          * If error bit is ECC we need to make sure
876          * ECC on all all WAYS are corrupted.
877          */

```

```

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     } else
893         Lxcache_tag_data
894             [last_error_injected_way] ^=
895                 pattern;
896
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
907     default:
908         return (ENOTSUP);
909     } /* end if switch(cmd) */
910
911     return (ret_val);
912
913 }
```

unchanged_portion_omitted

new/usr/src/uts/sun4u/ngdr/io/dr_quiesce.c

25005 Fri Jan 3 22:11:57 2014
new/usr/src/uts/sun4u/ngdr/io/dr_quiesce.c
patch cpu-pause-func-deglobalize

_____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     handle = srh->sr_dr_handlep;
833     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);

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     }

849     if (!force) {
850         struct dr_ref drc = {0};

852         prom_printf("\nDR: checking devices...\n");
853         dev_errs_idx = 0;

855         drc.arr = dev_errs;
856         drc.idx = &dev_errs_idx;
857         drc.len = DR_MAX_ERR_INT;

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     }

875 #ifndef SKIP_SYNC
876     /*
877      * This sync swap out all user pages
878      */
879     vfs_sync(SYNC_ALL);
880 #endif
```

1

new/usr/src/uts/sun4u/ngdr/io/dr_quiesce.c

```
882     /*
883      * special treatment for lock manager
884      */
885     lm_cprssuspend();

887 #ifndef SKIP_SYNC
888     /*
889      * sync the file system in case we never make it back
890      */
891     sync();
892 #endif

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     /* No parent to hold busy */
901     if ((rc = dr_suspend_devices(ddi_root_node(), srh)) != DDI_SUCCESS) {
902         if (srh->sr_err_idx && srh->sr_dr_handlep) {
903             (srh->sr_dr_handlep)->h_err = drerr_int(ESBD_SUSPEND,
904                                         srh->sr_err_ints, srh->sr_err_idx, 1);
905         }
906         dr_resume(srh);
907         return (rc);
908     }

910     drmach_suspend_last();

912     /*
913      * finally, grab all cpus
914      */
915     srh->sr_suspend_state = DR_SRSTATE_FULL;

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     }

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);

937     mutex_enter(&cpu_lock);
938     pause_cpus(NULL, NULL);
939     pause_cpus(NULL);
940     dr_stop_intr();

941     return (rc);
942 }
```

_____unchanged_portion_omitted_____

2

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

```
*****
50490 Fri Jan 3 22:11:57 2014
new/usr/src/uts/sun4u/os/cpr_impl.c
patch cpu-pause-func-deglobalize
*****
_____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 kctx = 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     }
274
275     i_cpr_xcall(i_cpr_clear_entries);
276 } else
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 }
```

```
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/sbdp_quiesce.c
```

```
*****  
19584 Fri Jan 3 22:11:57 2014  
new/usr/src/uts/sun4u/serengeti/io/sbdp_quiesce.c  
patch cpu-pause-func-deglobalize  
*****  
unchanged_portion_omitted
```

```
764 int  
765 sbdp_suspend(sbdp_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("\nsbdp_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/sbdp_quiesce.c  
*****  
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
```

```
2
```

```
new/usr/src/uts/sun4v/os/mpo.c
```

```
1
```

```
*****
54782 Fri Jan 3 22:11:58 2014
new/usr/src/uts/sun4v/os/mpo.c
patch cpu-pause-func-deglobalize
*****
_____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
219 static void
220 mpo_wr_lock()
221 {
222     mutex_enter(&cpu_lock);
223     pause_cpus(NULL, NULL);
224     pause_cpus(NULL);
225 }
_____unchanged_portion_omitted_____

```

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

```

*****21873 Fri Jan 3 22:11:58 2014*****
new/usr/src/uts/sun4v/os/suspend.c
patch cpu-pause-func-deglobalize
*****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             pg_cpu_inactive(cp);
399         pgps[id] = cp->cpu_pg;
400         pg_cpu_bootstrap(cp);
401     }
402     start_cpus();
403
404     /*
405      * pg_cpu_fini* and pg_cpu_init* must be called while CPUs are
406      * not paused. Use two separate loops here so that we do not

```

1

```

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 < NCPUs; 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 < NCPUs; 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         pause_cpus(NULL);
435         for (id = 0; id < NCPUs; id++) {
436             if ((cp = cpu_get(id)) == NULL)
437                 continue;
438             cp->cpu_pg = pgps[id];
439             if ((cp->cpu_flags & CPU_OFFLINE) == 0)
440                 pg_cpu_active(cp);
441         }
442         start_cpus();
443
444         mutex_exit(&cpu_lock);
445     }
446     (void) md_fini_handle(mdp);
447 }

_____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;

602     ASSERT(suspend_supported());
603     DBG("suspend: %s", __func__);
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     pause_cpus(NULL);
620     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: %ld\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