

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

1

```
*****
74605 Fri Apr 11 14:23:38 2014
new/usr/src/uts/common/os/clock.c
4748 use an enum for tod_faulted_global
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*      Copyright (c) 1984, 1986, 1987, 1988, 1989 AT&T */
22 /*      All Rights Reserved */

24 /*
25  * Copyright (c) 1988, 2010, Oracle and/or its affiliates. All rights reserved.
26  * Copyright (c) 2013, Joyent, Inc. All rights reserved.
27  */

29 #include <sys/param.h>
30 #include <sys/t_lock.h>
31 #include <sys/types.h>
32 #include <sys/tuneable.h>
33 #include <sys/sysmacros.h>
34 #include <sys/system.h>
35 #include <sys/cpuvar.h>
36 #include <sys/lgrp.h>
37 #include <sys/user.h>
38 #include <sys/proc.h>
39 #include <sys/callo.h>
40 #include <sys/kmem.h>
41 #include <sys/var.h>
42 #include <sys/cmn_err.h>
43 #include <sys/swap.h>
44 #include <sys/vmsystem.h>
45 #include <sys/class.h>
46 #include <sys/time.h>
47 #include <sys/debug.h>
48 #include <sys/vtrace.h>
49 #include <sys/spl.h>
50 #include <sys/atomic.h>
51 #include <sys/dumphdr.h>
52 #include <sys/archsystem.h>
53 #include <sys/fs/swapnode.h>
54 #include <sys/panic.h>
55 #include <sys/disp.h>
56 #include <sys/msacct.h>
57 #include <sys/mem_cage.h>

59 #include <vm/page.h>
60 #include <vm/anon.h>
61 #include <vm/rm.h>
```

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

2

```
62 #include <sys/cyclic.h>
63 #include <sys/cputpart.h>
64 #include <sys/rctl.h>
65 #include <sys/task.h>
66 #include <sys/sdt.h>
67 #include <sys/ddi_periodic.h>
68 #include <sys/random.h>
69 #include <sys/modctl.h>

71 /*
72  * for NTP support
73  */
74 #include <sys/timex.h>
75 #include <sys/inttypes.h>

77 #include <sys/sunddi.h>
78 #include <sys/clock_impl.h>

80 /*
81  * clock() is called straight from the clock cyclic; see clock_init().
82  *
83  * Functions:
84  *     reprime clock
85  *     maintain date
86  *     jab the scheduler
87  */

89 extern kcondvar_t      fsflush_cv;
90 extern sysinfo_t      sysinfo;
91 extern vminfo_t      vminfo;
92 extern int      idleswtch; /* flag set while idle in pswtch() */
93 extern hrtime_t volatile devinfo_freeze;

95 /*
96  * high-precision avenrun values. These are needed to make the
97  * regular avenrun values accurate.
98  */
99 static uint64_t hp_avenrun[3];
100 int      avenrun[3]; /* FSCALED average run queue lengths */
101 time_t time; /* time in seconds since 1970 - for compatibility only */

103 static struct loadavg_s loadavg;
104 /*
105  * Phase/frequency-lock loop (PLL/FLL) definitions
106  *
107  * The following variables are read and set by the ntp_adjtime() system
108  * call.
109  *
110  * time_state shows the state of the system clock, with values defined
111  * in the timex.h header file.
112  *
113  * time_status shows the status of the system clock, with bits defined
114  * in the timex.h header file.
115  *
116  * time_offset is used by the PLL/FLL to adjust the system time in small
117  * increments.
118  *
119  * time_constant determines the bandwidth or "stiffness" of the PLL.
120  *
121  * time_tolerance determines maximum frequency error or tolerance of the
122  * CPU clock oscillator and is a property of the architecture; however,
123  * in principle it could change as result of the presence of external
124  * discipline signals, for instance.
125  *
126  * time_precision is usually equal to the kernel tick variable; however,
127  * in cases where a precision clock counter or external clock is
```

```

128 * available, the resolution can be much less than this and depend on
129 * whether the external clock is working or not.
130 *
131 * time_maxerror is initialized by a ntp_adjtime() call and increased by
132 * the kernel once each second to reflect the maximum error bound
133 * growth.
134 *
135 * time_esterror is set and read by the ntp_adjtime() call, but
136 * otherwise not used by the kernel.
137 */
138 int32_t time_state = TIME_OK; /* clock state */
139 int32_t time_status = STA_UNSYNC; /* clock status bits */
140 int32_t time_offset = 0; /* time offset (us) */
141 int32_t time_constant = 0; /* pll time constant */
142 int32_t time_tolerance = MAXFREQ; /* frequency tolerance (scaled ppm) */
143 int32_t time_precision = 1; /* clock precision (us) */
144 int32_t time_maxerror = MAXPHASE; /* maximum error (us) */
145 int32_t time_esterror = MAXPHASE; /* estimated error (us) */

147 /*
148 * The following variables establish the state of the PLL/FLL and the
149 * residual time and frequency offset of the local clock. The scale
150 * factors are defined in the timex.h header file.
151 *
152 * time_phase and time_freq are the phase increment and the frequency
153 * increment, respectively, of the kernel time variable.
154 *
155 * time_freq is set via ntp_adjtime() from a value stored in a file when
156 * the synchronization daemon is first started. Its value is retrieved
157 * via ntp_adjtime() and written to the file about once per hour by the
158 * daemon.
159 *
160 * time_adj is the adjustment added to the value of tick at each timer
161 * interrupt and is recomputed from time_phase and time_freq at each
162 * seconds rollover.
163 *
164 * time_reftime is the second's portion of the system time at the last
165 * call to ntp_adjtime(). It is used to adjust the time_freq variable
166 * and to increase the time_maxerror as the time since last update
167 * increases.
168 */
169 int32_t time_phase = 0; /* phase offset (scaled us) */
170 int32_t time_freq = 0; /* frequency offset (scaled ppm) */
171 int32_t time_adj = 0; /* tick adjust (scaled 1 / hz) */
172 int32_t time_reftime = 0; /* time at last adjustment (s) */

174 /*
175 * The scale factors of the following variables are defined in the
176 * timex.h header file.
177 *
178 * pps_time contains the time at each calibration interval, as read by
179 * microtime(). pps_count counts the seconds of the calibration
180 * interval, the duration of which is nominally pps_shift in powers of
181 * two.
182 *
183 * pps_offset is the time offset produced by the time median filter
184 * pps_tf[], while pps_jitter is the dispersion (jitter) measured by
185 * this filter.
186 *
187 * pps_freq is the frequency offset produced by the frequency median
188 * filter pps_ff[], while pps_stabil is the dispersion (wander) measured
189 * by this filter.
190 *
191 * pps_usec is latched from a high resolution counter or external clock
192 * at pps_time. Here we want the hardware counter contents only, not the
193 * contents plus the time_tv.usec as usual.

```

```

194 *
195 * pps_valid counts the number of seconds since the last PPS update. It
196 * is used as a watchdog timer to disable the PPS discipline should the
197 * PPS signal be lost.
198 *
199 * pps_glitch counts the number of seconds since the beginning of an
200 * offset burst more than tick/2 from current nominal offset. It is used
201 * mainly to suppress error bursts due to priority conflicts between the
202 * PPS interrupt and timer interrupt.
203 *
204 * pps_intcnt counts the calibration intervals for use in the interval-
205 * adaptation algorithm. It's just too complicated for words.
206 */
207 struct timeval pps_time; /* kernel time at last interval */
208 int32_t pps_tf[] = {0, 0, 0}; /* pps time offset median filter (us) */
209 int32_t pps_offset = 0; /* pps time offset (us) */
210 int32_t pps_jitter = MAXTIME; /* time dispersion (jitter) (us) */
211 int32_t pps_ff[] = {0, 0, 0}; /* pps frequency offset median filter */
212 int32_t pps_freq = 0; /* frequency offset (scaled ppm) */
213 int32_t pps_stabil = MAXFREQ; /* frequency dispersion (scaled ppm) */
214 int32_t pps_usec = 0; /* microsec counter at last interval */
215 int32_t pps_valid = PPS_VALID; /* pps signal watchdog counter */
216 int32_t pps_glitch = 0; /* pps signal glitch counter */
217 int32_t pps_count = 0; /* calibration interval counter (s) */
218 int32_t pps_shift = PPS_SHIFT; /* interval duration (s) (shift) */
219 int32_t pps_intcnt = 0; /* intervals at current duration */

221 /*
222 * PPS signal quality monitors
223 *
224 * pps_jitcnt counts the seconds that have been discarded because the
225 * jitter measured by the time median filter exceeds the limit MAXTIME
226 * (100 us).
227 *
228 * pps_calcnt counts the frequency calibration intervals, which are
229 * variable from 4 s to 256 s.
230 *
231 * pps_errcnt counts the calibration intervals which have been discarded
232 * because the wander exceeds the limit MAXFREQ (100 ppm) or where the
233 * calibration interval jitter exceeds two ticks.
234 *
235 * pps_stbcnt counts the calibration intervals that have been discarded
236 * because the frequency wander exceeds the limit MAXFREQ / 4 (25 us).
237 */
238 int32_t pps_jitcnt = 0; /* jitter limit exceeded */
239 int32_t pps_calcnt = 0; /* calibration intervals */
240 int32_t pps_errcnt = 0; /* calibration errors */
241 int32_t pps_stbcnt = 0; /* stability limit exceeded */

243 kcondvar_t lbolt_cv;

245 /*
246 * Hybrid lbolt implementation:
247 *
248 * The service historically provided by the lbolt and lbolt64 variables has
249 * been replaced by the ddi_get_lbolt() and ddi_get_lbolt64() routines, and the
250 * original symbols removed from the system. The once clock driven variables are
251 * now implemented in an event driven fashion, backed by gethrtime() coarsed to
252 * the appropriate clock resolution. The default event driven implementation is
253 * complemented by a cyclic driven one, active only during periods of intense
254 * activity around the DDI lbolt routines, when a lbolt specific cyclic is
255 * reprogramed to fire at a clock tick interval to serve consumers of lbolt who
256 * rely on the original low cost of consulting a memory position.
257 *
258 * The implementation uses the number of calls to these routines and the
259 * frequency of these to determine when to transition from event to cyclic

```

```

260 * driven and vice-versa. These values are kept on a per CPU basis for
261 * scalability reasons and to prevent CPUs from constantly invalidating a single
262 * cache line when modifying a global variable. The transition from event to
263 * cyclic mode happens once the thresholds are crossed, and activity on any CPU
264 * can cause such transition.
265 *
266 * The lbolt_hybrid function pointer is called by ddi_get_lbolt() and
267 * ddi_get_lbolt64(), and will point to lbolt_event_driven() or
268 * lbolt_cyclic_driven() according to the current mode. When the thresholds
269 * are exceeded, lbolt_event_driven() will reprogram the lbolt cyclic to
270 * fire at a nsec_per_tick interval and increment an internal variable at
271 * each firing. lbolt_hybrid will then point to lbolt_cyclic_driven(), which
272 * will simply return the value of such variable. lbolt_cyclic() will attempt
273 * to shut itself off at each threshold interval (sampling period for calls
274 * to the DDI lbolt routines), and return to the event driven mode, but will
275 * be prevented from doing so if lbolt_cyclic_driven() is being heavily used.
276 *
277 * lbolt_bootstrap is used during boot to serve lbolt consumers who don't wait
278 * for the cyclic subsystem to be initialized.
279 *
280 */
281 int64_t lbolt_bootstrap(void);
282 int64_t lbolt_event_driven(void);
283 int64_t lbolt_cyclic_driven(void);
284 int64_t (*lbolt_hybrid)(void) = lbolt_bootstrap;
285 uint_t lbolt_ev_to_cyclic(caddr_t, caddr_t);

287 /*
288 * lbolt's cyclic, installed by clock_init().
289 */
290 static void lbolt_cyclic(void);

292 /*
293 * Tunable to keep lbolt in cyclic driven mode. This will prevent the system
294 * from switching back to event driven, once it reaches cyclic mode.
295 */
296 static boolean_t lbolt_cyc_only = B_FALSE;

298 /*
299 * Cache aligned, per CPU structure with lbolt usage statistics.
300 */
301 static lbolt_cpu_t *lb_cpu;

303 /*
304 * Single, cache aligned, structure with all the information required by
305 * the lbolt implementation.
306 */
307 lbolt_info_t *lb_info;

310 int one_sec = 1; /* turned on once every second */
311 static int fsflushcnt; /* counter for t_fsflushr */
312 int dosynctodr = 1; /* patchable; enable/disable sync to TOD chip */
313 int tod_needsync = 0; /* need to sync tod chip with software time */
314 static int tod_broken = 0; /* clock chip doesn't work */
315 time_t boot_time = 0; /* Boot time in seconds since 1970 */
316 cyclic_id_t clock_cyclic; /* clock()'s cyclic_id */
317 cyclic_id_t deadman_cyclic; /* deadman()'s cyclic_id */

319 extern void clock_tick_schedule(int);

321 static int lgrp_ticks; /* counter to schedule lgrp load calcs */

323 /*
324 * for tod fault detection
325 */

```

```

326 #define TOD_REF_FREQ ((longlong_t)(NANOSEC))
327 #define TOD_STALL_THRESHOLD (TOD_REF_FREQ * 3 / 2)
328 #define TOD_JUMP_THRESHOLD (TOD_REF_FREQ / 2)
329 #define TOD_FILTER_N 4
330 #define TOD_FILTER_SETTLE (4 * TOD_FILTER_N)
331 static enum tod_fault_type tod_faulted = TOD_NOFAULT;
331 static int tod_faulted = TOD_NOFAULT;

333 static int tod_status_flag = 0; /* used by tod_validate() */

335 static hrtime_t prev_set_tick = 0; /* gethrtime() prior to tod_set() */
336 static time_t prev_set_tod = 0; /* tv_sec value passed to tod_set() */

338 /* patchable via /etc/system */
339 int tod_validate_enable = 1;

341 /* Diagnose/Limit messages about delay(9F) called from interrupt context */
342 int delay_from_interrupt_diagnose = 0;
343 volatile uint32_t delay_from_interrupt_msg = 20;

345 /*
346 * On non-SPARC systems, TOD validation must be deferred until gethrtime
347 * returns non-zero values (after mach_clkinit's execution).
348 * On SPARC systems, it must be deferred until after hrtime_base
349 * and hres_last_tick are set (in the first invocation of hres_tick).
350 * Since in both cases the prerequisites occur before the invocation of
351 * tod_get() in clock(), the deferment is lifted there.
352 */
353 static boolean_t tod_validate_deferred = B_TRUE;

355 /*
356 * tod_fault_table[] must be aligned with
357 * enum tod_fault_type in systm.h
358 */
359 static char *tod_fault_table[] = {
360     "Reversed", /* TOD_REVERSED */
361     "Stalled", /* TOD_STALLED */
362     "Jumped", /* TOD_JUMPED */
363     "Changed in Clock Rate", /* TOD_RATECHANGED */
364     "Is Read-Only" /* TOD_RDONLY */
365     /*
366      * no strings needed for TOD_NOFAULT
367      */
368 };
_____unchanged_portion_omitted_____

```