

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

```
*****
68010 Sat Nov 15 21:14:32 2014
new/usr/src/uts/common/disp/disp.c
patch setfrontbackdq
*****
unchanged_portion_omitted_
1151 /*
1152 * setbackdq() keeps rungs 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 rung's lengths
1155 * must match. When per-thread TS_RUNQMATCH flag is set, setbackdq() will
1156 * try to keep rungs 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 */
1160 #define RUNQ_LEN(cp, pri) ((cp)->cpu_disp->disp_q[pri].dq_sruncnt)

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 *
1180 * Called with the thread in transition, onproc or stopped state and locked
1181 * (transition implies locked) and at high spl. Returns with the thread in
1182 * TS_RUN state and still locked.
1183 * Called with the thread in transition, onproc or stopped state
1184 * and locked (transition implies locked) and at high spl.
1185 * Returns with the thread in TS_RUN state and still locked.
1186 */
1187 static void
1188 setfrontbackdq(kthread_t *tp, boolean_t front)
1189 void
1190 setbackdq(kthread_t *tp)
1191 {
1192     dispq_t        *dq;
1193     disp_t          *dp;
1194     cpu_t           *cp;
1195     pri_t           tpri;
1196     boolean_t       bound;
1197     int              bound;
1198     boolean_t       self;
1199
1200     ASSERT(THREAD_LOCK_HELD(tp));
1201     ASSERT((tp->t_schedflag & TS_ALLSTART) == 0);
1202     ASSERT(!thread_on_queue(tp)); /* make sure tp isn't on a rung */
1203
1204     /*
1205      * If thread is "swapped" or on the swap queue don't
1206      * queue it, but wake sched.
1207      */
1208     if ((tp->t_schedflag & (TS_LOAD | TS_ON_SWAPQ)) != TS_LOAD) {
1209         disp_swapped_setrun(tp);
```

1

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

```
1202             return;
1203         }
1204         self = (tp == curthread);
1205         bound = (tp->t_bound_cpu || tp->t_weakbound_cpu);
1206         if (tp->t_bound_cpu || tp->t_weakbound_cpu)
1207             bound = 1;
1208         else
1209             bound = 0;
1210
1211         tpri = DISP_PRIO(tp);
1212         if (ncpus == 1)
1213             cp = tp->t_cpu;
1214         else if (!bound) {
1215             if (tpri >= kpqpri) {
1216                 setkpqdq(tp, front ? SETKP_FRONT : SETKP_BACK);
1217                 setkpqdq(tp, SETKP_BACK);
1218             }
1219             cp = tp->t_cpu;
1220         }
1221     #endif /* ! codereview */
1222     /*
1223      * We'll generally let this thread continue to run where
1224      * it last ran...but will consider migration if:
1225      * - We thread probably doesn't have much cache warmth.
1226      * - The CPU where it last ran is the target of an offli
1227      *   request.
1228      * - The thread last ran outside its home lgroup.
1229      */
1230     if ((!THREAD_HAS_CACHE_WARMTH(tp)) ||
1231         (cp == cpu_inmotion))
1232         cp = disp_lowpri_cpu(cp, tp->t_lpl, tpri, NULL);
1233     else if (!LGRP_CONTAINS_CPU(tp->t_lpl->lpl_lgrp, cp))
1234         cp = disp_lowpri_cpu(cp, tp->t_lpl, tpri,
1235                             self ? cp : NULL);
1236
1237     if ((!THREAD_HAS_CACHE_WARMTH(tp)) ||
1238         (tp->t_cpu == cpu_inmotion))
1239         cp = disp_lowpri_cpu(tp->t_cpu, tp->t_lpl, tpri, NULL);
1240     else if (!LGRP_CONTAINS_CPU(tp->t_lpl->lpl_lgrp, tp->t_cpu))
1241         cp = disp_lowpri_cpu(tp->t_cpu, tp->t_lpl, tpri,
1242                             self ? tp->t_cpu : NULL);
1243     else {
1244         cp = tp->t_cpu;
1245     }
1246
1247     if (tp->t_cpupart == cp->cpu_part) {
1248         if (front) {
1249             /*
1250              * We'll generally let this thread continue to r
1251              * where it last ran, but will consider migratio
1252              * - The thread last ran outside its home lgroup
1253              * - The CPU where it last ran is the target of
1254              *   offline request (a thread_nomigrate() on th
1255              *   motion CPU relies on this when forcing a pr
1256              * - The thread isn't the highest priority threa
1257              *   it last ran, and it is considered not likel
1258              *   have significant cache warmth.
1259              */
1260             if ((!LGRP_CONTAINS_CPU(tp->t_lpl->lpl_lgrp, cp)
1261                  (cp == cpu_inmotion)))
1262                 cp = disp_lowpri_cpu(cp, tp->t_lpl, tpri)
```

2

```

1254             self ? cp : NULL);
1255         } else if ((tpri < cp->cpu_disp->disp_maxrunpri)
1256             (!THREAD_HAS_CACHE_WARMTH(tp))) {
1257             cp = disp_lowpri_cpu(cp, tp->t_lpl, tpri
1258                                 NULL);
1259         }
1260     } else {
1261 #endif /* ! codereview */
1262         int qlen;
1263
1264         /*
1265          * Perform any CMT load balancing
1266          */
1267         cp = cmt_balance(tp, cp);
1268
1269         /*
1270          * Balance across the run queues
1271          */
1272         qlen = RUNQ_LEN(cp, tpri);
1273         if (tpri >= RUNQ_MATCH_PRI &&
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     }
1293 #endif /* ! codereview */
1294     } else {
1295         /*
1296          * Migrate to a cpu in the new partition.
1297          */
1298         cp = disp_lowpri_cpu(tp->t_cpupart->cp_cpulist,
1299                             tp->t_lpl, tp->t_pri, NULL);
1300     }
1301 #endif /* ! codereview */
1302     ASSERT((cp->cpu_flags & CPU QUIESCED) == 0);
1303     } else {
1304         /*
1305          * It is possible that t_weakbound_cpu != t_bound_cpu (for
1306          * a short time until weak binding that existed when the
1307          * strong binding was established has dropped) so we must
1308          * favour weak binding over strong.
1309          */
1310         cp = tp->t_weakbound_cpu ?
1311             tp->t_weakbound_cpu : tp->t_bound_cpu;
1312     }
1313
1314 #endif /* ! codereview */
1315     /*
1316      * A thread that is ONPROC may be temporarily placed on the run queue
1317      * but then chosen to run again by disp. If the thread we're placing on
1318      * the queue is in TS_ONPROC state, don't set its t_waitrq until a
1319

```

```

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             hrttime_t curtime;
1326
1327             curtime = gethrttime_unscaled();
1328             (void) cpu_update_pct(tp, curtime);
1329             tp->t_waitrq = curtime;
1330         } else {
1331             (void) cpu_update_pct(tp, gethrttime_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_FRONTQ, "frontq:pri %d tid %p", tpri,
1340                     tp);
1341         } else {
1342             DTRACE_SCHED3(enqueue, kthread_t *, tp, disp_t *, dp, int, 0);
1343             TRACE_3(TR_FAC_DISP, TR_BACKQ, "setbackdq:pri %d cpu %p tid %p",
1344                     tpri, cp, tp);
1345         }
1346     }
1347 #endif /* ! codereview */
1348
1349 #ifndef NPROBE
1350     /* Kernel probe */
1351     if (tnf_tracing_active)
1352         tnf_thread_queue(tp, cp, tpri);
1353 #endif /* NPROBE */
1354
1355     ASSERT(tpri >= 0 && tpri < dp->disp_npri);
1356
1357     THREAD_RUN(tp, &dp->disp_lock); /* set t_state to TS_RUN */
1358     tp->t_disp_queue = dp;
1359     tp->t_link = NULL;
1360
1361     dq = &dp->disp_q[tpri];
1362     dp->disp_nrunnable++;
1363     if (!bound)
1364         dp->disp_stal = 0;
1365     membar_enter();
1366
1367     if (dq->dq_srunct++ != 0) {
1368         if (front) {
1369             ASSERT(dq->dq_last != NULL);
1370             tp->t_link = dq->dq_first;
1371             dq->dq_first = tp;
1372         } else {
1373             ASSERT(dq->dq_first != NULL);
1374             dq->dq_last->t_link = tp;
1375             dq->dq_last = tp;
1376         }
1377     }
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     }
1385     membar_enter();
1386

```

```

1385         cpu_resched(cp, tpri);
1386     }
1387 }
1388 if (!bound && tpri > dp->disp_max_unbound_pri) {
1389     if (self && dp->disp_max_unbound_pri == -1 && cp == CPU) {
1390         /*
1391          * If there are no other unbound threads on the
1392          * run queue, don't allow other CPUs to steal
1393          * this thread while we are in the middle of a
1394          * context switch. We may just switch to it
1395          * again right away. CPU_DISP_DONTSTEAL is cleared
1396          * in swtch and swtch_to.
1397          */
1398         cp->cpu_disp_flags |= CPU_DISP_DONTSTEAL;
1399     }
1400     dp->disp_max_unbound_pri = tpri;
1401 }
1402 }

1403 #endif /* ! codereview */
1404 (*disp_enq_thread)(cp, bound);
1405 }
1406 }

1407 /*
1408  * Put the specified thread on the back of the dispatcher
1409  * queue corresponding to its current priority.
1410  *
1411  * Called with the thread in transition, onproc or stopped state
1412  * and locked (transition implies locked) and at high spl.
1413  * Returns with the thread in TS_RUN state and still locked.
1414  */
1415 void
1416 setbackdq(kthread_t *tp)
1417 {
1418     setfrontbackdq(tp, B_FALSE);
1419 }
1420 }

1421 /*
1422 #endif /* ! codereview */
1423  * Put the specified thread on the front of the dispatcher
1424  * queue corresponding to its current priority.
1425  *
1426  * Called with the thread in transition, onproc or stopped state
1427  * and locked (transition implies locked) and at high spl.
1428  * Returns with the thread in TS_RUN state and still locked.
1429  */
1430 void
1431 setfrontdq(kthread_t *tp)
1432 {
1433     setfrontbackdq(tp, B_TRUE);
1434     disp_t      *dp;
1435     dispq_t     *dq;
1436     cpu_t       *cp;
1437     pri_t       tpri;
1438     int          bound;
1439
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     /*
1445      * If thread is "swapped" or on the swap queue don't
1446      * queue it, but wake sched.
1447      */
1448     if ((tp->t_schedflag & (TS_LOAD | TS_ON_SWAPQ)) != TS_LOAD) {
1449         disp_swapped_setrun(tp);
1450     }

```

```

1248         return;
1249     }
1250     if (tp->t_bound_cpu || tp->t_weakbound_cpu)
1251         bound = 1;
1252     else
1253         bound = 0;
1254
1255     tpri = DISP_PRIO(tp);
1256     if (ncpus == 1)
1257         cp = tp->t_cpu;
1258     else if (!bound) {
1259         if (tpri >= kpqpri) {
1260             setkpqdq(tp, SETKP_FRONT);
1261             return;
1262         }
1263         cp = tp->t_cpu;
1264     if (tp->t_cpupart == cp->cpu_part) {
1265         /*
1266          * We'll generally let this thread continue to run
1267          * where it last ran, but will consider migration if:
1268          * - The thread last ran outside it's home lgroup.
1269          * - The CPU where it last ran is the target of an
1270          * offline request (a thread_nomigrate() on the in
1271          * motion CPU relies on this when forcing a preempt).
1272          * - The thread isn't the highest priority thread where
1273          * it last ran, and it is considered not likely to
1274          * have significant cache warmth.
1275          */
1276     if ((!LGRP_CONTAINS_CPU(tp->t_lpl->lpl_lgrp, cp)) ||
1277         (cp == cpu_inmotion)) {
1278         cp = disp_lowpri_cpu(tp->t_cpu, tp->t_lpl, tpri,
1279                             (tp == curthread) ? cp : NULL);
1280     } else if ((tpri < cp->cpu_disp->disp_maxrunpri) &&
1281                (!THREAD_HAS_CACHE_WARMTH(tp))) {
1282         cp = disp_lowpri_cpu(tp->t_cpu, tp->t_lpl, tpri,
1283                             NULL);
1284     }
1285     } else {
1286         /*
1287          * Migrate to a cpu in the new partition.
1288          */
1289         cp = disp_lowpri_cpu(tp->t_cpupart->cp_cpulist,
1290                             tp->t_lpl, tp->t_pri, NULL);
1291     }
1292     ASSERT((cp->cpu_flags & CPU QUIESCED) == 0);
1293 }
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             hrtime_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         DTRACE_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 */
1334
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_stole = 0;
1344     membar_enter();
1345
1346     if (dq->dq_sruncont++ != 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_gactmap, 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             cp->cpu_disp_flags |= CPU_DISP_DONTSTEAL;
1375         }
1376         dp->disp_max_unbound_pri = tpri;
1377     }
1378     (*disp_eng_thread)(cp, bound);
1379 }

```

unchanged portion omitted