

```

*****
77297 Tue Sep 3 20:26:48 2013
new/usr/src/cmd/mdb/common/modules/zfs/zfs.c
4101 metaslab_debug should allow for fine-grained control
4102 space_maps should store more information about themselves
4103 space_map object blocksize should be increased
4104 ::spa_space no longer works
4105 removing a mirrored log device results in a leaked object
4106 asynchronously load metaslab
Reviewed by: Matthew Ahrens <mahrens@delphix.com>
Reviewed by: Adam Leventhal <ahl@delphix.com>
Reviewed by: Sebastien Roy <seb@delphix.com>
*****
_____unchanged_portion_omitted_____

1458 typedef struct mdb_space_map_phys_t {
1459     uint64_t smp_alloc;
1460 } mdb_space_map_phys_t;

1462 typedef struct mdb_space_map {
1463     uint64_t sm_size;
1464     uint64_t sm_alloc;
1465     uintptr_t sm_phys;
1466 } mdb_space_map_t;

1468 typedef struct mdb_range_tree {
1469     uint64_t rt_space;
1470 } mdb_range_tree_t;

1472 typedef struct mdb metaslab {
1473     uintptr_t ms_alloctree[TXG_SIZE];
1474     uintptr_t ms_freetree[TXG_SIZE];
1475     uintptr_t ms_tree;
1476     uintptr_t ms_sm;
1477     space_map_t ms_allocmap[TXG_SIZE];
1478     space_map_t ms_freemap[TXG_SIZE];
1479     space_map_t ms_map;
1480     space_map_obj_t ms_smo;
1481     space_map_obj_t ms_smo_syncing;
1482 } mdb metaslab_t;

1484 typedef struct space_data {
1485     uint64_t ms_alloctree[TXG_SIZE];
1486     uint64_t ms_freetree[TXG_SIZE];
1487     uint64_t ms_tree;
1488     uint64_t ms_allocmap[TXG_SIZE];
1489     uint64_t ms_freemap[TXG_SIZE];
1490     uint64_t ms_map;
1491     uint64_t avail;
1492     uint64_t nowavail;
1493 } space_data_t;

1494 /* ARGSUSED */
1495 static int
1496 space_cb(uintptr_t addr, const void *unknown, void *arg)
1497 {
1498     space_data_t *sd = arg;
1499     mdb metaslab_t ms;
1500     mdb_range_tree_t rt;
1501     mdb_space_map_t sm;
1502     mdb_space_map_phys_t smp = { 0 };
1503     int i;

1504     if (mdb_ctf_vread(&ms, "metaslab_t", "mdb metaslab_t",
1505         addr, 0) == -1)
1506         if (GETMEMB(addr, "metaslab", ms_allocmap, ms.ms_allocmap) ||

```

```

1482     GETMEMB(addr, "metaslab", ms_freemap, ms.ms_freemap) ||
1483     GETMEMB(addr, "metaslab", ms_map, ms.ms_map) ||
1484     GETMEMB(addr, "metaslab", ms_smo, ms.ms_smo) ||
1485     GETMEMB(addr, "metaslab", ms_smo_syncing, ms.ms_smo_syncing)) {
1500         return (WALK_ERR);

1502     for (i = 0; i < TXG_SIZE; i++) {

1504         if (mdb_ctf_vread(&rt, "range_tree_t",
1505             "mdb_range_tree_t", ms.ms_alloctree[i], 0) == -1)
1506             sd->ms_alloctree[i] += rt.rt_space;

1508         if (mdb_ctf_vread(&rt, "range_tree_t",
1509             "mdb_range_tree_t", ms.ms_freetree[i], 0) == -1)
1510             sd->ms_freetree[i] += rt.rt_space;
1511     }

1513     if (mdb_ctf_vread(&rt, "range_tree_t",
1514         "mdb_range_tree_t", ms.ms_tree, 0) == -1 ||
1515         mdb_ctf_vread(&sm, "space_map_t",
1516             "mdb_space_map_t", ms.ms_sm, 0) == -1)
1517         return (WALK_ERR);
1518     sd->ms_allocmap[0] += ms.ms_allocmap[0].sm_space;
1519     sd->ms_allocmap[1] += ms.ms_allocmap[1].sm_space;
1520     sd->ms_allocmap[2] += ms.ms_allocmap[2].sm_space;
1521     sd->ms_allocmap[3] += ms.ms_allocmap[3].sm_space;
1522     sd->ms_freemap[0] += ms.ms_freemap[0].sm_space;
1523     sd->ms_freemap[1] += ms.ms_freemap[1].sm_space;
1524     sd->ms_freemap[2] += ms.ms_freemap[2].sm_space;
1525     sd->ms_freemap[3] += ms.ms_freemap[3].sm_space;
1526     sd->ms_map += ms.ms_map.sm_space;
1527     sd->avail += ms.ms_map.sm_size - ms.ms_smo.smo_alloc;
1528     sd->nowavail += ms.ms_map.sm_size - ms.ms_smo_syncing.smo_alloc;

1529     if (sm.sm_phys != NULL) {
1530         (void) mdb_ctf_vread(&smp, "space_map_phys_t",
1531             "mdb_space_map_phys_t", sm.sm_phys, 0);
1532     }

1533     sd->ms_tree += rt.rt_space;
1534     sd->avail += sm.sm_size - sm.sm_alloc;
1535     sd->nowavail += sm.sm_size - smp.smp_alloc;

1536     return (WALK_NEXT);
1537 }

1538 /*
1539  * ::spa_space [-b]
1540  *
1541  * Given a spa_t, print out it's on-disk space usage and in-core
1542  * estimates of future usage. If -b is given, print space in bytes.
1543  * Otherwise print in megabytes.
1544  */
1545 /* ARGSUSED */
1546 static int
1547 spa_space(uintptr_t addr, uint_t flags, int argc, const mdb_arg_t *argv)
1548 {
1549     mdb spa_t spa;
1550     uintptr_t dp_root_dir;
1551     mdb_dsl_dir_t dd;
1552     mdb_dsl_dir_phys_t dsp;
1553     uint64_t children;
1554     uintptr_t childaddr;
1555     space_data_t sd;
1556     int shift = 20;
1557     char *suffix = "M";

```

```

1551     int bytes = B_FALSE;

1553     if (mdb_getopts(argc, argv, 'b', MDB_OPT_SETBITS, TRUE, &bytes, NULL) !=
1554         argc)
1555         return (DCMD_USAGE);
1556     if (!(flags & DCMD_ADDRSPEC))
1557         return (DCMD_USAGE);

1559     if (bytes) {
1560         shift = 0;
1561         suffix = "";
1562     }

1564     if (GETMEMB(addr, "spa", spa_dsl_pool, spa.spa_dsl_pool) ||
1565         GETMEMB(addr, "spa", spa_root_vdev, spa.spa_root_vdev) ||
1566         GETMEMB(spa.spa_root_vdev, "vdev", vdev_children, children) ||
1567         GETMEMB(spa.spa_root_vdev, "vdev", vdev_child, childaddr) ||
1568         GETMEMB(spa.spa_dsl_pool, "dsl_pool",
1569                 dp_root_dir, dp_root_dir) ||
1570         GETMEMB(dp_root_dir, "dsl_dir", dd_phys, dd.dd_phys) ||
1571         GETMEMB(dp_root_dir, "dsl_dir",
1572                 dd_space_towrite, dd.dd_space_towrite) ||
1573         GETMEMB(dd.dd_phys, "dsl_dir_phys",
1574                 dd_used_bytes, dsp.dd_used_bytes) ||
1575         GETMEMB(dd.dd_phys, "dsl_dir_phys",
1576                 dd_compressed_bytes, dsp.dd_compressed_bytes) ||
1577         GETMEMB(dd.dd_phys, "dsl_dir_phys",
1578                 dd_uncompressed_bytes, dsp.dd_uncompressed_bytes)) {
1579         return (DCMD_ERR);
1580     }

1582     mdb_printf("dd_space_towrite = %llu%s %llu%s %llu%s %llu%s\n",
1583               dd.dd_space_towrite[0] >> shift, suffix,
1584               dd.dd_space_towrite[1] >> shift, suffix,
1585               dd.dd_space_towrite[2] >> shift, suffix,
1586               dd.dd_space_towrite[3] >> shift, suffix);

1588     mdb_printf("dd_phys.dd_used_bytes = %llu%s\n",
1589               dsp.dd_used_bytes >> shift, suffix);
1590     mdb_printf("dd_phys.dd_compressed_bytes = %llu%s\n",
1591               dsp.dd_compressed_bytes >> shift, suffix);
1592     mdb_printf("dd_phys.dd_uncompressed_bytes = %llu%s\n",
1593               dsp.dd_uncompressed_bytes >> shift, suffix);

1595     bzero(&sd, sizeof (sd));
1596     if (mdb_pwalk("metaslab", space_cb, &sd, addr) != 0) {
1597         mdb_warn("can't walk metaslabs");
1598         return (DCMD_ERR);
1599     }

1601     mdb_printf("ms_allocmap = %llu%s %llu%s %llu%s %llu%s\n",
1602               sd.ms_allocmap[0] >> shift, suffix,
1603               sd.ms_allocmap[1] >> shift, suffix,
1604               sd.ms_allocmap[2] >> shift, suffix,
1605               sd.ms_allocmap[3] >> shift, suffix);
1606     mdb_printf("ms_freemap = %llu%s %llu%s %llu%s %llu%s\n",
1607               sd.ms_freemap[0] >> shift, suffix,
1608               sd.ms_freemap[1] >> shift, suffix,
1609               sd.ms_freemap[2] >> shift, suffix,
1610               sd.ms_freemap[3] >> shift, suffix);
1611     mdb_printf("ms_tree = %llu%s\n", sd.ms_tree >> shift, suffix);
1612     sd.ms_freemap[0] >> shift, suffix,

```

```

1581         sd.ms_freemap[1] >> shift, suffix,
1582         sd.ms_freemap[2] >> shift, suffix,
1583         sd.ms_freemap[3] >> shift, suffix);
1584     mdb_printf("ms_map = %llu%s\n", sd.ms_map >> shift, suffix);
1612     mdb_printf("last synced avail = %llu%s\n", sd.avail >> shift, suffix);
1613     mdb_printf("current syncing avail = %llu%s\n",
1614               sd.nowavail >> shift, suffix);

1616     return (DCMD_OK);
1617 }

```

unchanged_portion_omitted

```

*****
90879 Tue Sep 3 20:26:50 2013
new/usr/src/cmd/zdb/zdb.c
4101 metaslab_debug should allow for fine-grained control
4102 space_maps should store more information about themselves
4103 space_map object blocksize should be increased
4104 ::spa_space no longer works
4105 removing a mirrored log device results in a leaked object
4106 asynchronously load metaslab
Reviewed by: Matthew Ahrens <mahrens@delphix.com>
Reviewed by: Adam Leventhal <ahl@delphix.com>
Reviewed by: Sebastien Roy <seb@delphix.com>
*****
unchanged_portion_omitted_

245 const char histo_stars[] = "*****";
246 const int histo_width = sizeof (histo_stars) - 1;

248 static void
249 dump_histogram(const uint64_t *histo, int size, int offset)
249 dump_histogram(const uint64_t *histo, int size)
250 {
251     int i;
252     int minidx = size - 1;
253     int maxidx = 0;
254     uint64_t max = 0;

256     for (i = 0; i < size; i++) {
257         if (histo[i] > max)
258             max = histo[i];
259         if (histo[i] > 0 && i > maxidx)
260             maxidx = i;
261         if (histo[i] > 0 && i < minidx)
262             minidx = i;
263     }

265     if (max < histo_width)
266         max = histo_width;

268     for (i = minidx; i <= maxidx; i++) {
269         (void) printf("\t\t\t%3u: %6llu %s\n",
270             i + offset, (u_longlong_t)histo[i],
271             i, (u_longlong_t)histo[i],
272             &histo_stars[(max - histo[i]) * histo_width / max]);
273     }

275 static void
276 dump_zap_stats(objset_t *os, uint64_t object)
277 {
278     int error;
279     zap_stats_t zs;

281     error = zap_get_stats(os, object, &zs);
282     if (error)
283         return;

285     if (zs.zs_ptrtbl_len == 0) {
286         ASSERT(zs.zs_num_blocks == 1);
287         (void) printf("\tmicrozap: %llu bytes, %llu entries\n",
288             (u_longlong_t)zs.zs_blocksize,
289             (u_longlong_t)zs.zs_num_entries);
290         return;
291     }

293     (void) printf("\tFat ZAP stats:\n");

```

```

295     (void) printf("\t\tPointer table:\n");
296     (void) printf("\t\t\t%llu elements\n",
297         (u_longlong_t)zs.zs_ptrtbl_len);
298     (void) printf("\t\t\tz_tz_blk: %llu\n",
299         (u_longlong_t)zs.zs_ptrtbl_zt_blk);
300     (void) printf("\t\t\tz_tz_numblks: %llu\n",
301         (u_longlong_t)zs.zs_ptrtbl_zt_numblks);
302     (void) printf("\t\t\tz_tz_shift: %llu\n",
303         (u_longlong_t)zs.zs_ptrtbl_zt_shift);
304     (void) printf("\t\t\tz_tz_blks_copied: %llu\n",
305         (u_longlong_t)zs.zs_ptrtbl_blks_copied);
306     (void) printf("\t\t\tz_tz_nextblk: %llu\n",
307         (u_longlong_t)zs.zs_ptrtbl_nextblk);

309     (void) printf("\t\tZAP entries: %llu\n",
310         (u_longlong_t)zs.zs_num_entries);
311     (void) printf("\t\tLeaf blocks: %llu\n",
312         (u_longlong_t)zs.zs_num_leafs);
313     (void) printf("\t\tTotal blocks: %llu\n",
314         (u_longlong_t)zs.zs_num_blocks);
315     (void) printf("\t\tzap_block_type: 0x%llx\n",
316         (u_longlong_t)zs.zs_block_type);
317     (void) printf("\t\tzap_magic: 0x%llx\n",
318         (u_longlong_t)zs.zs_magic);
319     (void) printf("\t\tzap_salt: 0x%llx\n",
320         (u_longlong_t)zs.zs_salt);

322     (void) printf("\t\tLeafs with 2^n pointers:\n");
323     dump_histogram(zs.zs_leafs_with_2n_pointers, ZAP_HISTOGRAM_SIZE, 0);
323     dump_histogram(zs.zs_leafs_with_2n_pointers, ZAP_HISTOGRAM_SIZE);

325     (void) printf("\t\tBlocks with n*5 entries:\n");
326     dump_histogram(zs.zs_blocks_with_n5_entries, ZAP_HISTOGRAM_SIZE, 0);
326     dump_histogram(zs.zs_blocks_with_n5_entries, ZAP_HISTOGRAM_SIZE);

328     (void) printf("\t\tBlocks n/10 full:\n");
329     dump_histogram(zs.zs_blocks_n_tenths_full, ZAP_HISTOGRAM_SIZE, 0);
329     dump_histogram(zs.zs_blocks_n_tenths_full, ZAP_HISTOGRAM_SIZE);

331     (void) printf("\t\tEntries with n chunks:\n");
332     dump_histogram(zs.zs_entries_using_n_chunks, ZAP_HISTOGRAM_SIZE, 0);
332     dump_histogram(zs.zs_entries_using_n_chunks, ZAP_HISTOGRAM_SIZE);

334     (void) printf("\t\tBuckets with n entries:\n");
335     dump_histogram(zs.zs_buckets_with_n_entries, ZAP_HISTOGRAM_SIZE, 0);
335     dump_histogram(zs.zs_buckets_with_n_entries, ZAP_HISTOGRAM_SIZE);
336 }

unchanged_portion_omitted_

524 int
525 get_dtl_refcount(vdev_t *vd)
526 {
527     int refcount = 0;

529     if (vd->vdev_ops->vdev_op_leaf) {
530         space_map_t *sm = vd->vdev_dtl_sm;

532         if (sm != NULL &&
533             sm->sm_dbuf->db_size == sizeof (space_map_phys_t))
534             return (1);
535         return (0);
536     }

538     for (int c = 0; c < vd->vdev_children; c++)
539         refcount += get_dtl_refcount(vd->vdev_child[c]);

```

```

540     return (refcount);
541 }

543 int
544 get metaslab_refcount(vdev_t *vd)
545 {
546     int refcount = 0;

548     if (vd->vdev_top == vd) {
549         for (int m = 0; m < vd->vdev_ms_count; m++) {
550             space_map_t *sm = vd->vdev_ms[m]->ms_sm;

552             if (sm != NULL &&
553                 sm->sm_dbuf->db_size == sizeof (space_map_phys_t))
554                 refcount++;
555         }
556     }
557     for (int c = 0; c < vd->vdev_children; c++)
558         refcount += get metaslab_refcount(vd->vdev_child[c]);

560     return (refcount);
561 }

563 static int
564 verify spacemap_refcounts(spa_t *spa)
565 {
566     int expected_refcount, actual_refcount;

568     expected_refcount = spa_feature_get_refcount(spa,
569         &spa_feature_table[SPA_FEATURE_SPACEMAP_HISTOGRAM]);
570     actual_refcount = get_dtl_refcount(spa->spa_root_vdev);
571     actual_refcount += get metaslab_refcount(spa->spa_root_vdev);

573     if (expected_refcount != actual_refcount) {
574         (void) printf("space map refcount mismatch: expected %d != "
575             "actual %d\n", expected_refcount, actual_refcount);
576         return (2);
577     }
578     return (0);
579 }

581 static void
582 dump spacemap(objset_t *os, space_map_t *sm)
583 dump spacemap(objset_t *os, space_map_obj_t *smo, space_map_t *sm)
584 {
585     uint64_t alloc, offset, entry;
586     uint8_t mapshift = sm->sm_shift;
587     uint64_t mapstart = sm->sm_start;
588     char *ddata[] = { "ALLOC", "FREE", "CONDENSE", "INVALID",
589         "INVALID", "INVALID", "INVALID", "INVALID" };

591     if (sm == NULL)
592         return;
593     /*
594      * Print out the freelist entries in both encoded and decoded form.
595      */
596     for (offset = 0; offset < space_map_length(sm);
597         offset += sizeof (entry)) {
598         uint8_t mapshift = sm->sm_shift;

599         VERIFY0(dmu_read(os, space_map_object(sm), offset,
600             for (offset = 0; offset < sm->smo_objsize; offset += sizeof (entry)) {
601             VERIFY3U(0, ==, dmu_read(os, sm->smo_object, offset,

```

```

600         sizeof (entry), &entry, DMU_READ_PREFETCH));
601         if (SM_DEBUG_DECODE(entry)) {

603             (void) printf("\t [%6llu] %s: txg %llu, pass %llu\n",
604                 (u_longlong_t)(offset / sizeof (entry)),
605                 ddata[SM_DEBUG_ACTION_DECODE(entry)],
606                 (u_longlong_t)SM_DEBUG_TXG_DECODE(entry),
607                 (u_longlong_t)SM_DEBUG_SYNCPASS_DECODE(entry));
608         } else {
609             (void) printf("\t [%6llu] %c range:"
610                 " %010llx-%010llx size: %06llx\n",
611                 (u_longlong_t)(offset / sizeof (entry)),
612                 SM_TYPE_DECODE(entry) == SM_ALLOC ? 'A' : 'F',
613                 (u_longlong_t)((SM_OFFSET_DECODE(entry) <<
614                     mapshift) + sm->sm_start),
615                 mapshift + mapstart,
616                 (u_longlong_t)((SM_OFFSET_DECODE(entry) <<
617                     mapshift) + sm->sm_start +
618                     SM_RUN_DECODE(entry) << mapshift)),
619                 mapshift + mapstart + (SM_RUN_DECODE(entry) <<
620                     mapshift)),
621                 (u_longlong_t)(SM_RUN_DECODE(entry) << mapshift));
622             if (SM_TYPE_DECODE(entry) == SM_ALLOC)
623                 alloc += SM_RUN_DECODE(entry) << mapshift;
624             else
625                 alloc -= SM_RUN_DECODE(entry) << mapshift;
626         }
627     }
628     if (alloc != space_map_allocated(sm)) {
629         if (alloc != sm->smo_alloc) {
630             (void) printf("space_map_object alloc (%llu) INCONSISTENT "
631                 "with space map summary (%llu)\n",
632                 (u_longlong_t)space_map_allocated(sm), (u_longlong_t)alloc);
633             (u_longlong_t)smo->smo_alloc, (u_longlong_t)alloc);
634         }
635     }

636 static void
637 dump metaslab_stats(metaslab_t *msp)
638 {
639     char maxbuf[32];
640     range_tree_t *rt = msp->ms_tree;
641     avl_tree_t *t = &msp->ms_size_tree;
642     int free_pct = range_tree_space(rt) * 100 / msp->ms_size;
643     space_map_t *sm = msp->ms_map;
644     avl_tree_t *t = sm->sm_pp_root;
645     int free_pct = sm->sm_space * 100 / sm->sm_size;

646     zdb_nicenum(metaslab_block_maxsize(msp), maxbuf);
647     zdb_nicenum(space_map_maxsize(sm), maxbuf);

648     (void) printf("\t %25s %10lu %7s %6s %4s %4d%%\n",
649         "segments", avl_numnodes(t), "maxsize", maxbuf,
650         "freepct", free_pct);
651     (void) printf("\tIn-memory histogram:\n");
652     dump_histogram(rt->rt_histogram, RANGE_TREE_HISTOGRAM_SIZE, 0);
653 }

654 static void
655 dump metaslab(metaslab_t *msp)
656 {
657     vdev_t *vd = msp->ms_group->mg_vd;
658     spa_t *spa = vd->vdev_spa;
659     space_map_t *sm = msp->ms_sm;
660     space_map_t *sm = msp->ms_map;
661     space_map_obj_t *smo = &msp->ms_smo;

```

```

655     char freebuf[32];

657     zdb_nicenum(msp->ms_size - space_map_allocated(sm), freebuf);
659     zdb_nicenum(sm->sm_size - smo->smo_alloc, freebuf);

659     (void) printf(
660         "\tmetaslab %6llu offset %21llx spacemap %6llu free %5s\n",
661         (u_longlong_t)msp->ms_id, (u_longlong_t)msp->ms_start,
662         (u_longlong_t)space_map_object(sm), freebuf);
663         (u_longlong_t)(sm->sm_start / sm->sm_size),
664         (u_longlong_t)sm->sm_start, (u_longlong_t)smo->smo_object, freebuf);

664     if (dump_opt['m'] > 2 && !dump_opt['L']) {
665         if (dump_opt['m'] > 1 && !dump_opt['L']) {
666             mutex_enter(&msp->ms_lock);
667             metaslab_load_wait(msp);
668             if (!msp->ms_loaded) {
669                 VERIFY0(metaslab_load(msp));
670                 range_tree_stat_verify(msp->ms_tree);
671             }
672             space_map_load_wait(sm);
673             if (!sm->sm_loaded)
674                 VERIFY(space_map_load(sm, zfs_metaslab_ops,
675                     SM_FREE, smo, spa->spa_meta_objset) == 0);
676             dump_metaslab_stats(msp);
677             metaslab_unload(msp);
678             space_map_unload(sm);
679             mutex_exit(&msp->ms_lock);
680         }
681     }

682     if (dump_opt['m'] > 1 && sm != NULL &&
683         spa_feature_is_active(spa,
684             &spa_feature_table[SPA_FEATURE_SPACEMAP_HISTOGRAM])) {
685         /*
686          * The space map histogram represents free space in chunks
687          * of sm_shift (i.e. bucket 0 refers to 2^sm_shift).
688          */
689         (void) printf("\tOn-disk histogram:\n");
690         dump_histogram(sm->sm_phys->smp_histogram,
691             SPACE_MAP_HISTOGRAM_SIZE(sm), sm->sm_shift);
692     }

693     if (dump_opt['d'] > 5 || dump_opt['m'] > 2) {
694         ASSERT(sm->sm_size == (LULL << vd->vdev_ms_shift));
695     }

696     if (dump_opt['d'] > 5 || dump_opt['m'] > 3) {
697         ASSERT(msp->ms_size == (LULL << vd->vdev_ms_shift));
698     }

699     mutex_enter(&msp->ms_lock);
700     dump_spacemap(spa->spa_meta_objset, msp->ms_sm);
701     dump_spacemap(spa->spa_meta_objset, smo, sm);
702     mutex_exit(&msp->ms_lock);
703 }

704 unchanged_portion_omitted

705 static void
706 dump_dtl_seg(void *arg, uint64_t start, uint64_t size)
707 dump_dtl_seg(space_map_t *sm, uint64_t start, uint64_t size)
708 {
709     char *prefix = arg;
710     char *prefix = (void *)sm;

711     (void) printf("%s [%llu,%llu] length %llu\n",
712         prefix,
713         (u_longlong_t)start,
714         (u_longlong_t)(start + size),

```

```

887         (u_longlong_t)(size));
888     }

889 static void
890 dump_dtl(vdev_t *vd, int indent)
891 {
892     spa_t *spa = vd->vdev_spa;
893     boolean_t required;
894     char *name[DTL_TYPES] = { "missing", "partial", "scrub", "outage" };
895     char prefix[256];

896     spa_vdev_state_enter(spa, SCL_NONE);
897     required = vdev_dtl_required(vd);
898     (void) spa_vdev_state_exit(spa, NULL, 0);

902     if (indent == 0)
903         (void) printf("\nDirty time logs:\n\n");

904     (void) printf("\t%s [%s] [%s]\n", indent, "",
905         vd->vdev_path ? vd->vdev_path :
906         vd->vdev_parent ? vd->vdev_ops->vdev_op_type : spa_name(spa),
907         required ? "DTL-required" : "DTL-expendable");

908     for (int t = 0; t < DTL_TYPES; t++) {
909         range_tree_t *rt = vd->vdev_dtl[t];
910         if (range_tree_space(rt) == 0)
911             space_map_t *sm = &vd->vdev_dtl[t];
912             if (sm->sm_space == 0)
913                 continue;
914             (void) snprintf(prefix, sizeof(prefix), "\t%s [%s]",
915                 indent + 2, "", name[t]);
916             mutex_enter(rt->rt_lock);
917             range_tree_walk(rt, dump_dtl_seg, prefix);
918             mutex_exit(rt->rt_lock);
919             mutex_enter(sm->sm_lock);
920             space_map_walk(sm, dump_dtl_seg, (void *)prefix);
921             mutex_exit(sm->sm_lock);
922             if (dump_opt['d'] > 5 && vd->vdev_children == 0)
923                 dump_spacemap(spa->spa_meta_objset, vd->vdev_dtl_sm);
924                 dump_spacemap(spa->spa_meta_objset,
925                     &vd->vdev_dtl_smo, sm);
926     }

927     for (int c = 0; c < vd->vdev_children; c++)
928         dump_dtl(vd->vdev_child[c], indent + 4);
929 }

930 unchanged_portion_omitted

931 static void
932 zdb_leak(void *arg, uint64_t start, uint64_t size)
933 zdb_leak(space_map_t *sm, uint64_t start, uint64_t size)
934 {
935     vdev_t *vd = arg;
936     vdev_t *vd = sm->sm_ppd;

937     (void) printf("leaked space: vdev %llu, offset 0x%llx, size %llu\n",
938         (u_longlong_t)vd->vdev_id, (u_longlong_t)start, (u_longlong_t)size);
939 }

940 static metaslab_ops_t zdb_metaslab_ops = {
941     /* ARGSUSED */
942     static void
943     zdb_space_map_load(space_map_t *sm)
944     {
945     }
946 }

```

```

2238 static void
2239 zdb_space_map_unload(space_map_t *sm)
2240 {
2241     space_map_vacate(sm, zdb_leak, sm);
2242 }

2244 /* ARGSUSED */
2245 static void
2246 zdb_space_map_claim(space_map_t *sm, uint64_t start, uint64_t size)
2247 {
2248 }

2250 static space_map_ops_t zdb_space_map_ops = {
2251     zdb_space_map_load,
2252     zdb_space_map_unload,
2253     NULL, /* alloc */
2254     NULL, /* fragmented */
2255     zdb_space_map_claim,
2256     NULL, /* free */
2257     NULL, /* maxsize */
2258 };
2259 #ifndef unchanged_portion_omitted
2260
2261 static void
2262 zdb_leak_init(spa_t *spa, zdb_cb_t *zcb)
2263 {
2264     zcb->zcb_spa = spa;
2265
2266     if (!dump_opt['L']) {
2267         vdev_t *rvd = spa->spa_root_vdev;
2268         for (int c = 0; c < rvd->vdev_children; c++) {
2269             vdev_t *vd = rvd->vdev_child[c];
2270             for (int m = 0; m < vd->vdev_ms_count; m++) {
2271                 metaslab_t *msp = vd->vdev_ms[m];
2272                 mutex_enter(&msp->ms_lock);
2273                 metaslab_unload(msp);
2274             }
2275         }
2276
2277         /*
2278          * For leak detection, we overload the metaslab
2279          * ms_tree to contain allocated segments
2280          * instead of free segments. As a result,
2281          * we can't use the normal metaslab_load/unload
2282          * interfaces.
2283          */
2284         if (msp->ms_sm != NULL) {
2285             msp->ms_ops = &zdb_metaslab_ops;
2286             VERIFY0(space_map_load(msp->ms_sm,
2287                 msp->ms_tree, SM_ALLOC));
2288             msp->ms_loaded = B_TRUE;
2289         }
2290         space_map_unload(msp->ms_map);
2291         VERIFY(space_map_load(msp->ms_map,
2292             &zdb_space_map_ops, SM_ALLOC, &msp->ms_smo,
2293             spa->spa_meta_objset) == 0);
2294         msp->ms_map->sm_ppd = vd;
2295         mutex_exit(&msp->ms_lock);
2296     }
2297 }
2298
2299 spa_config_enter(spa, SCL_CONFIG, FTAG, RW_READER);
2300
2301 zdb_ddt_leak_init(spa, zcb);
2302
2303 spa_config_exit(spa, SCL_CONFIG, FTAG);
2304 }

```

```

2388 static void
2389 zdb_leak_fini(spa_t *spa)
2390 {
2391     if (!dump_opt['L']) {
2392         vdev_t *rvd = spa->spa_root_vdev;
2393         for (int c = 0; c < rvd->vdev_children; c++) {
2394             vdev_t *vd = rvd->vdev_child[c];
2395             for (int m = 0; m < vd->vdev_ms_count; m++) {
2396                 metaslab_t *msp = vd->vdev_ms[m];
2397                 mutex_enter(&msp->ms_lock);
2398
2399                 /*
2400                  * The ms_tree has been overloaded to
2401                  * contain allocated segments. Now that we
2402                  * finished traversing all blocks, any
2403                  * block that remains in the ms_tree
2404                  * represents an allocated block that we
2405                  * did not claim during the traversal.
2406                  * Claimed blocks would have been removed
2407                  * from the ms_tree.
2408                  */
2409                 range_tree_vacate(msp->ms_tree, zdb_leak, vd);
2410                 msp->ms_loaded = B_FALSE;
2411
2412                 space_map_unload(msp->ms_map);
2413                 mutex_exit(&msp->ms_lock);
2414             }
2415         }
2416     }
2417 #ifndef unchanged_portion_omitted
2418
2419 static int
2420 dump_block_stats(spa_t *spa)
2421 {
2422     zdb_cb_t zcb = { 0 };
2423     zdb_blkstats_t *zb, *tzb;
2424     uint64_t norm_alloc, norm_space, total_alloc, total_found;
2425     int flags = TRAVERSE_PRE | TRAVERSE_PREFETCH_METADATA | TRAVERSE_HARD;
2426     int leaks = 0;
2427
2428     (void) printf("\nTraversing all blocks %s%s%s%s...\n\n",
2429         (dump_opt['c'] || !dump_opt['L']) ? "to verify " : "",
2430         (dump_opt['c'] == 1) ? "metadata " : "",
2431         dump_opt['c'] ? "checksums " : "",
2432         (dump_opt['c'] && !dump_opt['L']) ? "and verify " : "",
2433         !dump_opt['L'] ? "nothing leaked " : "");
2434
2435     /*
2436      * Load all space maps as SM_ALLOC maps, then traverse the pool
2437      * claiming each block we discover. If the pool is perfectly
2438      * consistent, the space maps will be empty when we're done.
2439      * Anything left over is a leak; any block we can't claim (because
2440      * it's not part of any space map) is a double allocation,
2441      * reference to a freed block, or an unclaimed log block.
2442      */
2443     zdb_leak_init(spa, &zcb);
2444
2445     /*
2446      * If there's a deferred-free bplist, process that first.
2447      */
2448     (void) bpobj_iterate_nofree(&spa->spa_deferred_bpobj,
2449         count_block_cb, &zcb, NULL);
2450     if (spa_version(spa) >= SPA_VERSION_DEADLISTS) {
2451         (void) bpobj_iterate_nofree(&spa->spa_dsl_pool->dp_free_bpobj,

```

```

2467         count_block_cb, &zcb, NULL);
2468     }
2469     if (spa_feature_is_active(spa,
2470         &spa_feature_table[SPA_FEATURE_ASYNC_DESTROY])) {
2471         VERIFY3U(0, ==, bptree_iterate(spa->spa_meta_objset,
2472             spa->spa_dsl_pool->dp_bptree_obj, B_FALSE, count_block_cb,
2473             &zcb, NULL));
2474     }
2475
2476     if (dump_opt['c'] > 1)
2477         flags |= TRAVERSE_PREFETCH_DATA;
2478
2479     zcb.zcb_totalasize = metaslab_class_get_alloc(spa_normal_class(spa));
2480     zcb.zcb_start = zcb.zcb_lastprint = gethrtime();
2481     zcb.zcb_haderrors |= traverse_pool(spa, 0, flags, zdb_blkptr_cb, &zcb);
2482
2483     /*
2484     * If we've traversed the data blocks then we need to wait for those
2485     * I/Os to complete. We leverage "The Godfather" zio to wait on
2486     * all async I/Os to complete.
2487     */
2488     if (dump_opt['c']) {
2489         (void) zio_wait(spa->spa_async_zio_root);
2490         spa->spa_async_zio_root = zio_root(spa, NULL, NULL,
2491             ZIO_FLAG_CANFAIL | ZIO_FLAG_SPECULATIVE |
2492             ZIO_FLAG_GODFATHER);
2493     }
2494
2495     if (zcb.zcb_haderrors) {
2496         (void) printf("\nError counts:\n\n");
2497         (void) printf("\t%5s %s\n", "errno", "count");
2498         for (int e = 0; e < 256; e++) {
2499             if (zcb.zcb_errors[e] != 0) {
2500                 (void) printf("\t%5d %llu\n",
2501                     e, (u_longlong_t)zcb.zcb_errors[e]);
2502             }
2503         }
2504     }
2505
2506     /*
2507     * Report any leaked segments.
2508     */
2509     zdb_leak_fini(spa);
2510
2511     tzb = &zcb.zcb_type[ZB_TOTAL][ZDB_OT_TOTAL];
2512
2513     norm_alloc = metaslab_class_get_alloc(spa_normal_class(spa));
2514     norm_space = metaslab_class_get_space(spa_normal_class(spa));
2515
2516     total_alloc = norm_alloc + metaslab_class_get_alloc(spa_log_class(spa));
2517     total_found = tzb->z_b_asize - zcb.zcb_dedup_asize;
2518
2519     if (total_found == total_alloc) {
2520         if (!dump_opt['L'])
2521             (void) printf("\n\tNo leaks (block sum matches space"
2522                 " maps exactly)\n");
2523     } else {
2524         (void) printf("block traversal size %llu != alloc %llu "
2525             "(%s %lld)\n",
2526             (u_longlong_t)total_found,
2527             (u_longlong_t)total_alloc,
2528             (dump_opt['L'] ? "unreachable" : "leaked",
2529             (longlong_t)(total_alloc - total_found));
2530         leaks = 1;
2531     }

```

```

2533     if (tzb->z_b_count == 0)
2534         return (2);
2535
2536     (void) printf("\n");
2537     (void) printf("\tbp count: %10llu\n",
2538         (u_longlong_t)tzb->z_b_count);
2539     (void) printf("\tganged count: %10llu\n",
2540         (longlong_t)tzb->z_b_gangs);
2541     (void) printf("\tbp logical: %10llu avg: %6llu\n",
2542         (u_longlong_t)tzb->z_b_lsize,
2543         (u_longlong_t)(tzb->z_b_lsize / tzb->z_b_count));
2544     (void) printf("\tbp physical: %10llu avg:"
2545         " %6llu compression: %6.2f\n",
2546         (u_longlong_t)tzb->z_b_psize,
2547         (u_longlong_t)(tzb->z_b_psize / tzb->z_b_count),
2548         (double)tzb->z_b_lsize / tzb->z_b_psize);
2549     (void) printf("\tbp allocated: %10llu avg:"
2550         " %6llu compression: %6.2f\n",
2551         (u_longlong_t)tzb->z_b_asize,
2552         (u_longlong_t)(tzb->z_b_asize / tzb->z_b_count),
2553         (double)tzb->z_b_lsize / tzb->z_b_asize);
2554     (void) printf("\tbp deduped: %10llu ref>1:"
2555         " %6llu deduplication: %6.2f\n",
2556         (u_longlong_t)zcb.zcb_dedup_asize,
2557         (u_longlong_t)zcb.zcb_dedup_blocks,
2558         (double)zcb.zcb_dedup_asize / tzb->z_b_asize + 1.0);
2559     (void) printf("\tSPA allocated: %10llu used: %5.2f%%\n",
2560         (u_longlong_t)norm_alloc, 100.0 * norm_alloc / norm_space);
2561
2562     if (tzb->z_b_ditto_samevdev != 0) {
2563         (void) printf("\tDittoed blocks on same vdev: %llu\n",
2564             (longlong_t)tzb->z_b_ditto_samevdev);
2565     }
2566
2567     if (dump_opt['b'] >= 2) {
2568         int l, t, level;
2569         (void) printf("\nBlocks\tLSIZE\tPSIZE\tASIZE"
2570             "\t\t avg\t comp\t%%Total\tType\n");
2571
2572         for (t = 0; t <= ZDB_OT_TOTAL; t++) {
2573             char csize[32], lsize[32], psize[32], asize[32];
2574             char avg[32], gang[32];
2575             char *typename;
2576
2577             if (t < DMU_OT_NUMTYPES)
2578                 typename = dm_u_ot[t].ot_name;
2579             else
2580                 typename = zdb_ot_extname[t - DMU_OT_NUMTYPES];
2581
2582             if (zcb.zcb_type[ZB_TOTAL][t].z_b_asize == 0) {
2583                 (void) printf("%6s\t%5s\t%5s\t%5s"
2584                     "\t\t%5s\t%5s\t%6s\t%5s\n",
2585                     "-",
2586                     "-",
2587                     "-",
2588                     "-",
2589                     "-",
2590                     "-",
2591                     "-",
2592                     typename);
2593                 continue;
2594             }
2595
2596             for (l = ZB_TOTAL - 1; l >= -1; l--) {
2597                 level = (l == -1 ? ZB_TOTAL : l);
2598                 zb = &zcb.zcb_type[level][t];

```

```

2600         if (zb->zb_asize == 0)
2601             continue;

2603         if (dump_opt['b'] < 3 && level != ZB_TOTAL)
2604             continue;

2606         if (level == 0 && zb->zb_asize ==
2607             zcb.zcb_type[ZB_TOTAL][t].zb_asize)
2608             continue;

2610         zdb_nicenum(zb->zb_count, csize);
2611         zdb_nicenum(zb->zb_lsize, lsize);
2612         zdb_nicenum(zb->zb_psize, psize);
2613         zdb_nicenum(zb->zb_asize, asize);
2614         zdb_nicenum(zb->zb_asize / zb->zb_count, avg);
2615         zdb_nicenum(zb->zb_gangs, gang);

2617         (void) printf("%6s\t%5s\t%5s\t%5s\t%5s"
2618                      "\t%5.2f\t%6.2f\t",
2619                      csize, lsize, psize, asize, avg,
2620                      (double)zb->zb_lsize / zb->zb_psize,
2621                      100.0 * zb->zb_asize / tzb->zb_asize);

2623         if (level == ZB_TOTAL)
2624             (void) printf("%s\n", typename);
2625         else
2626             (void) printf("    L%d %s\n",
2627                          level, typename);

2629         if (dump_opt['b'] >= 3 && zb->zb_gangs > 0) {
2630             (void) printf("\t number of ganged "
2631                          "blocks: %s\n", gang);
2632         }

2634         if (dump_opt['b'] >= 4) {
2635             (void) printf("psize "
2636                          "(in 512-byte sectors): "
2637                          "number of blocks\n");
2638             dump_histogram(zb->zb_psize_histogram,
2639                          PSIZE_HISTO_SIZE, 0);
2640             dump_histogram(tzb->tzb_psize_histogram,
2641                          PSIZE_HISTO_SIZE);
2642         }
2643     }

2645     (void) printf("\n");

2647     if (leaks)
2648         return (2);

2650     if (zcb.zcb_haderrors)
2651         return (3);

2653     return (0);
2654 }

```

unchanged portion omitted

```

2757 static void
2758 dump_zpool(spa_t *spa)
2759 {
2760     dsl_pool_t *dp = spa_get_dsl(spa);
2761     int rc = 0;

2763     if (dump_opt['S']) {

```

```

2764         dump_simulated_ddt(spa);
2765         return;
2766     }

2768     if (!dump_opt['e'] && dump_opt['C'] > 1) {
2769         (void) printf("\nCached configuration:\n");
2770         dump_nvlist(spa->spa_config, 8);
2771     }

2773     if (dump_opt['C'])
2774         dump_config(spa);

2776     if (dump_opt['u'])
2777         dump_uberblock(&spa->spa_uberblock, "\nUberblock:\n", "\n");

2779     if (dump_opt['D'])
2780         dump_all_ddts(spa);

2782     if (dump_opt['d'] > 2 || dump_opt['m'])
2783         dump_metaslabs(spa);

2785     if (dump_opt['d'] || dump_opt['i']) {
2786         dump_dir(dp->dp_meta_objset);
2787         if (dump_opt['d'] >= 3) {
2788             dump_bpobj(&spa->spa_deferred_bpobj,
2789                      "Deferred frees", 0);
2790             if (spa_version(spa) >= SPA_VERSION_DEADLISTS) {
2791                 dump_bpobj(&spa->spa_dsl_pool->dp_free_bpobj,
2792                          "Pool snapshot frees", 0);
2793             }
2795             if (spa_feature_is_active(spa,
2796                                     &spa_feature_table[SPA_FEATURE_ASYNC_DESTROY])) {
2797                 dump_bptree(spa->spa_meta_objset,
2798                          spa->spa_dsl_pool->dp_bptree_obj,
2799                          "Pool dataset frees");
2800             }
2801             dump_dtl(spa->spa_root_vdev, 0);
2802         }
2803         (void) dmu_objset_find(spa_name(spa), dump_one_dir,
2804                              NULL, DS_FIND_SNAPSHOTS | DS_FIND_CHILDREN);
2805     }
2806     if (dump_opt['b'] || dump_opt['c'])
2807         rc = dump_block_stats(spa);

2809     if (rc == 0)
2810         rc = verify_spacemap_refcounts(spa);

2812     if (dump_opt['s'])
2813         show_pool_stats(spa);

2815     if (dump_opt['h'])
2816         dump_history(spa);

2818     if (rc != 0)
2819         exit(rc);
2820 }

```

unchanged portion omitted

new/usr/src/cmd/ztest/ztest.c

1

160708 Tue Sep 3 20:26:52 2013

new/usr/src/cmd/ztest/ztest.c

4101 metaslab_debug should allow for fine-grained control
4102 space_maps should store more information about themselves
4103 space_map object blocksize should be increased
4104 ::spa_space no longer works
4105 removing a mirrored log device results in a leaked object
4106 asynchronously load metaslab

Reviewed by: Matthew Ahrens <mahrens@delphix.com>

Reviewed by: Adam Leventhal <ahl@delphix.com>

Reviewed by: Sebastien Roy <seb@delphix.com>

unchanged_portion_omitted

```
5320 static void *
5321 ztest_deadman_thread(void *arg)
5322 {
5323     ztest_shared_t *zs = arg;
5324     spa_t *spa = ztest_spa;
5325     hrttime_t delta, total = 0;

5327     for (;;) {
5328         delta = zs->zs_thread_stop - zs->zs_thread_start +
5329             MSEC2NSEC(zfs_deadman_synctime_ms);

5331         (void) poll(NULL, 0, (int)NSEC2MSEC(delta));

5333         /*
5334          * If the pool is suspended then fail immediately. Otherwise,
5335          * check to see if the pool is making any progress. If
5336          * vdev_deadman() discovers that there hasn't been any recent
5337          * I/Os then it will end up aborting the tests.
5338          */
5339         if (spa_suspended(spa) || spa->spa_root_vdev == NULL) {
5339             if (spa_suspended(spa)) {
5340                 fatal(0, "aborting test after %llu seconds because "
5341                     "pool has transitioned to a suspended state.",
5342                     zfs_deadman_synctime_ms / 1000);
5343                 return (NULL);
5344             }
5345             vdev_deadman(spa->spa_root_vdev);

5347             total += zfs_deadman_synctime_ms/1000;
5348             (void) printf("ztest has been running for %lld seconds\n",
5349                 total);
5350         }
5351     }
```

unchanged_portion_omitted

new/usr/src/common/zfs/zfeature_common.c

1

```
*****
4629 Tue Sep 3 20:26:53 2013
new/usr/src/common/zfs/zfeature_common.c
4101 metaslab_debug should allow for fine-grained control
4102 space_maps should store more information about themselves
4103 space_map object blocksize should be increased
4104 ::spa_space no longer works
4105 removing a mirrored log device results in a leaked object
4106 asynchronously load metaslab
Reviewed by: Matthew Ahrens <mahrens@delphix.com>
Reviewed by: Adam Leventhal <ahl@delphix.com>
Reviewed by: Sebastien Roy <seb@delphix.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21
22 /*
23 * Copyright (c) 2013 by Delphix. All rights reserved.
24 * Copyright (c) 2012 by Delphix. All rights reserved.
25 * Copyright (c) 2013 by Saso Kiselkov. All rights reserved.
26 * Copyright (c) 2013, Joyent, Inc. All rights reserved.
27 */
28 #ifndef _KERNEL
29 #include <sys/system.h>
30 #else
31 #include <errno.h>
32 #include <string.h>
33 #endif
34 #include <sys/debug.h>
35 #include <sys/fs/zfs.h>
36 #include <sys/inttypes.h>
37 #include <sys/types.h>
38 #include "zfeature_common.h"
39
40 /*
41  * Set to disable all feature checks while opening pools, allowing pools with
42  * unsupported features to be opened. Set for testing only.
43  */
44 boolean_t zfeature_checks_disable = B_FALSE;
45
46 zfeature_info_t spa_feature_table[SPA_FEATURES];
47
48 /*
49  * Valid characters for feature guids. This list is mainly for aesthetic
50  * purposes and could be expanded in the future. There are different allowed
51  * characters in the guids reverse dns portion (before the colon) and its
52  * short name (after the colon).
```

new/usr/src/common/zfs/zfeature_common.c

2

```
53 */
54 static int
55 valid_char(char c, boolean_t after_colon)
56 {
57     return ((c >= 'a' && c <= 'z') ||
58            (c >= '0' && c <= '9') ||
59            c == (after_colon ? '_' : '.'));
60 }
61
62 unchanged_portion_omitted
63
64 void
65 zpool_feature_init(void)
66 {
67     zfeature_register(SPA_FEATURE_ASYNC_DESTROY,
68                     "com.delphix:async_destroy", "async_destroy",
69                     "Destroy filesystems asynchronously.", B_TRUE, B_FALSE, NULL);
70     zfeature_register(SPA_FEATURE_EMPTY_BPOBJ,
71                     "com.delphix:empty_bpobj", "empty_bpobj",
72                     "Snapshots use less space.", B_TRUE, B_FALSE, NULL);
73     zfeature_register(SPA_FEATURE_LZ4_COMPRESS,
74                     "org.illumos:lz4_compress", "lz4_compress",
75                     "LZ4 compression algorithm support.", B_FALSE, B_FALSE, NULL);
76     zfeature_register(SPA_FEATURE_MULTI_VDEV_CRASH_DUMP,
77                     "com.joyent:multi_vdev_crash_dump", "multi_vdev_crash_dump",
78                     "Crash dumps to multiple vdev pools.", B_FALSE, B_FALSE, NULL);
79     zfeature_register(SPA_FEATURE_SPACEMAP_HISTOGRAM,
80                     "com.delphix:spacemap_histogram", "spacemap_histogram",
81                     "Spacemaps maintain space histograms.", B_TRUE, B_FALSE, NULL);
82 }
83
84 unchanged_portion_omitted
```

new/usr/src/common/zfs/zfeature_common.h

1

```
*****
2304 Tue Sep 3 20:26:54 2013
new/usr/src/common/zfs/zfeature_common.h
4101 metaslab_debug should allow for fine-grained control
4102 space_maps should store more information about themselves
4103 space_map object blocksize should be increased
4104 ::spa_space no longer works
4105 removing a mirrored log device results in a leaked object
4106 asynchronously load metaslab
Reviewed by: Matthew Ahrens <mahrens@delphix.com>
Reviewed by: Adam Leventhal <ahl@delphix.com>
Reviewed by: Sebastien Roy <seb@delphix.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21
22 /*
23 * Copyright (c) 2013 by Delphix. All rights reserved.
24 * Copyright (c) 2012 by Delphix. All rights reserved.
25 * Copyright (c) 2013 by Saso Kiselkov. All rights reserved.
26 * Copyright (c) 2013, Joyent, Inc. All rights reserved.
27 */
28 #ifndef _ZFEATURE_COMMON_H
29 #define _ZFEATURE_COMMON_H
30
31 #include <sys/fs/zfs.h>
32 #include <sys/inttypes.h>
33 #include <sys/types.h>
34
35 #ifdef __cplusplus
36 extern "C" {
37 #endif
38
39 struct zfeature_info;
40
41 typedef struct zfeature_info {
42     const char *fi_uname; /* User-facing feature name */
43     const char *fi_guid; /* On-disk feature identifier */
44     const char *fi_desc; /* Feature description */
45     boolean_t fi_can_readonly; /* Can open pool readonly w/o support? */
46     boolean_t fi_mos; /* Is the feature necessary to read the MOS? */
47     struct zfeature_info **fi_depends; /* array; null terminated */
48 } zfeature_info_t;
49
50 typedef int (zfeature_func_t)(zfeature_info_t *fi, void *arg);
51
52 #define ZFS_FEATURE_DEBUG
```

new/usr/src/common/zfs/zfeature_common.h

2

```
54 enum spa_feature {
55     SPA_FEATURE_ASYNC_DESTROY,
56     SPA_FEATURE_EMPTY_BPOBJ,
57     SPA_FEATURE_LZ4_COMPRESS,
58     SPA_FEATURE_MULTI_VDEV_CRASH_DUMP,
59     SPA_FEATURE_SPACEMAP_HISTOGRAM,
60     SPA_FEATURES
61 } spa_feature_t;
_____ unchanged_portion_omitted
```

```

*****
  9200 Tue Sep  3 20:26:55 2013
new/usr/src/man/man5/zpool-features.5
4101 metaslab_debug should allow for fine-grained control
4102 space_maps should store more information about themselves
4103 space_map object blocksize should be increased
4104 ::spa_space no longer works
4105 removing a mirrored log device results in a leaked object
4106 asynchronously load metaslab
Reviewed by: Matthew Ahrens <mahrens@delphix.com>
Reviewed by: Adam Leventhal <ahl@delphix.com>
Reviewed by: Sebastien Roy <seb@delphix.com>
*****
1  \" te
2  .\" Copyright (c) 2012 by Delphix. All rights reserved.
3  .\" Copyright (c) 2013 by Saso Kiselkov. All rights reserved.
4  .\" Copyright (c) 2013, Joyent, Inc. All rights reserved.
5  .\" The contents of this file are subject to the terms of the Common Development
6  .\" and Distribution License (the \"License\").  You may not use this file except
7  .\" in compliance with the License.  You can obtain a copy of the license at
8  .\" usr/src/OPENSOLARIS.LICENSE or http://www.opensolaris.org/os/licensing.
9  .\"
10 .\" See the License for the specific language governing permissions and
11 .\" limitations under the License.  When distributing Covered Code, include this
12 .\" CDDL HEADER in each file and include the License file at
13 .\" usr/src/OPENSOLARIS.LICENSE.  If applicable, add the following below this
14 .\" CDDL HEADER, with the fields enclosed by brackets \"[]\" replaced with your
15 .\" own identifying information:
16 .\" Portions Copyright [yyyy] [name of copyright owner]
17 .TH ZPOOL-FEATURES 5 \"Aug 27, 2013\"
18 .TH ZPOOL-FEATURES 5 \"Mar 16, 2012\"
19 .SH NAME
20 zpool\|-features \- ZFS pool feature descriptions
21 .sp
22 .LP
23 ZFS pool on\|-disk format versions are specified via \"features\" which replace
24 the old on\|-disk format numbers (the last supported on\|-disk format number is
25 28).  To enable a feature on a pool use the \fBupgrade\fR subcommand of the
26 \fBzpool\fR(1M) command, or set the \fBfeature@\fR\fIfeature_name\fR property
27 to \fBenabled\fR.
28 .sp
29 .LP
30 The pool format does not affect file system version compatibility or the ability
31 to send file systems between pools.
32 .sp
33 .LP
34 Since most features can be enabled independently of each other the on\|-disk
35 format of the pool is specified by the set of all features marked as
36 \fBactive\fR on the pool.  If the pool was created by another software version
37 this set may include unsupported features.
38 .SS \"Identifying features\"
39 .sp
40 .LP
41 Every feature has a guid of the form \fIcom.example:feature_name\fR.  The reverse
42 DNS name ensures that the feature's guid is unique across all ZFS
43 implementations.  When unsupported features are encountered on a pool they will
44 be identified by their guids.  Refer to the documentation for the ZFS
45 implementation that created the pool for information about those features.
46 .sp
47 .LP
48 Each supported feature also has a short name.  By convention a feature's short
49 name is the portion of its guid which follows the ':' (e.g.
50 \fIcom.example:feature_name\fR would have the short name \fIfeature_name\fR),
51 however a feature's short name may differ across ZFS implementations if
52 following the convention would result in name conflicts.

```

```

53 .SS \"Feature states\"
54 .sp
55 .LP
56 Features can be in one of three states:
57 .sp
58 .ne 2
59 .na
60 \fB\fBactive\fR\fR
61 .ad
62 .RS 12n
63 This feature's on\|-disk format changes are in effect on the pool.  Support for
64 this feature is required to import the pool in read\|-write mode.  If this
65 feature is not read-only compatible, support is also required to import the pool
66 in read\|-only mode (see \"Read\|-only compatibility\").
67 .RE
68
69 .sp
70 .ne 2
71 .na
72 \fB\fBenabled\fR\fR
73 .ad
74 .RS 12n
75 An administrator has marked this feature as enabled on the pool, but the
76 feature's on\|-disk format changes have not been made yet.  The pool can still be
77 imported by software that does not support this feature, but changes may be made
78 to the on\|-disk format at any time which will move the feature to the
79 \fBactive\fR state.  Some features may support returning to the \fBenabled\fR
80 state after becoming \fBactive\fR.  See feature\|-specific documentation for
81 details.
82 .RE
83
84 .sp
85 .ne 2
86 .na
87 \fB\fBdisabled\fR\fR
88 .ad
89 .RS 12n
90 This feature's on\|-disk format changes have not been made and will not be made
91 unless an administrator moves the feature to the \fBenabled\fR state.  Features
92 cannot be disabled once they have been enabled.
93 .RE
94
95 .sp
96 .LP
97 The state of supported features is exposed through pool properties of the form
98 \fIfeature@short_name\fR.
99 .SS \"Read\|-only compatibility\"
100 .sp
101 .LP
102 Some features may make on\|-disk format changes that do not interfere with other
103 software's ability to read from the pool.  These features are referred to as
104 \"read\|-only compatible\".  If all unsupported features on a pool are read\|-only
105 compatible, the pool can be imported in read\|-only mode by setting the
106 \fBreadonly\fR property during import (see \fBzpool\fR(1M) for details on
107 importing pools).
108 .SS \"Unsupported features\"
109 .sp
110 .LP
111 For each unsupported feature enabled on an imported pool a pool property
112 named \fIunsupported@feature_guid\fR will indicate why the import was allowed
113 despite the unsupported feature.  Possible values for this property are:
114
115 .sp
116 .ne 2
117 .na
118 \fB\fBinactive\fR\fR

```

```

119 .ad
120 .RS 12n
121 The feature is in the \fBenabled\fR state and therefore the pool's on\~disk
122 format is still compatible with software that does not support this feature.
123 .RE

125 .sp
126 .ne 2
127 .na
128 \fB\fBreadonly\fR\fR
129 .ad
130 .RS 12n
131 The feature is read\~only compatible and the pool has been imported in
132 read\~only mode.
133 .RE

135 .SS "Feature dependencies"
136 .sp
137 .LP
138 Some features depend on other features being enabled in order to function
139 properly. Enabling a feature will automatically enable any features it
140 depends on.
141 .SH FEATURES
142 .sp
143 .LP
144 The following features are supported on this system:
145 .sp
146 .ne 2
147 .na
148 \fB\fBasync_destroy\fR\fR
149 .ad
150 .RS 4n
151 .TS
152 l l .
153 GUID      com.delphix:async_destroy
154 READ\~ONLY COMPATIBLE  yes
155 DEPENDENCIES  none
156 .TE

158 Destroying a file system requires traversing all of its data in order to
159 return its used space to the pool. Without \fBasync_destroy\fR the file system
160 is not fully removed until all space has been reclaimed. If the destroy
161 operation is interrupted by a reboot or power outage the next attempt to open
162 the pool will need to complete the destroy operation synchronously.

164 When \fBasync_destroy\fR is enabled the file system's data will be reclaimed
165 by a background process, allowing the destroy operation to complete without
166 traversing the entire file system. The background process is able to resume
167 interrupted destroys after the pool has been opened, eliminating the need
168 to finish interrupted destroys as part of the open operation. The amount
169 of space remaining to be reclaimed by the background process is available
170 through the \fBfreeing\fR property.

172 This feature is only \fBactive\fR while \fBfreeing\fR is non\~zero.
173 .RE

175 .sp
176 .ne 2
177 .na
178 \fB\fBempty_bpobj\fR\fR
179 .ad
180 .RS 4n
181 .TS
182 l l .
183 GUID      com.delphix:empty_bpobj
184 READ\~ONLY COMPATIBLE  yes

```

```

185 DEPENDENCIES  none
186 .TE

188 This feature increases the performance of creating and using a large
189 number of snapshots of a single filesystem or volume, and also reduces
190 the disk space required.

192 When there are many snapshots, each snapshot uses many Block Pointer
193 Objects (bpobj's) to track blocks associated with that snapshot.
194 However, in common use cases, most of these bpobj's are empty. This
195 feature allows us to create each bpobj on-demand, thus eliminating the
196 empty bpobj's.

198 This feature is \fBactive\fR while there are any filesystems, volumes,
199 or snapshots which were created after enabling this feature.
200 .RE

202 .sp
203 .ne 2
204 .na
205 \fB\fBlz4_compress\fR\fR
206 .ad
207 .RS 4n
208 .TS
209 l l .
210 GUID      org.illumos:lz4_compress
211 READ\~ONLY COMPATIBLE  no
212 DEPENDENCIES  none
213 .TE

215 \fBlz4\fR is a high-performance real-time compression algorithm that
216 features significantly faster compression and decompression as well as a
217 higher compression ratio than the older \fBlzjb\fR compression.
218 Typically, \fBlz4\fR compression is approximately 50% faster on
219 compressible data and 200% faster on incompressible data than
220 \fBlzjb\fR. It is also approximately 80% faster on decompression, while
221 giving approximately 10% better compression ratio.

223 When the \fBlz4_compress\fR feature is set to \fBenabled\fR, the
224 administrator can turn on \fBlz4\fR compression on any dataset on the
225 pool using the \fBzfs\fR(1M) command. Please note that doing so will
226 immediately activate the \fBlz4_compress\fR feature on the underlying
227 pool (even before any data is written). Since this feature is not
228 read-only compatible, this operation will render the pool unimportable
229 on systems without support for the \fBlz4_compress\fR feature. At the
230 moment, this operation cannot be reversed. Booting off of
231 \fBlz4\fR-compressed root pools is supported.
232 .RE

234 .sp
235 .ne 2
236 .na
237 \fB\fBspacemap_histogram\fR\fR
238 .ad
239 .RS 4n
240 .TS
241 l l .
242 GUID      com.delphix:spacemap_histogram
243 READ\~ONLY COMPATIBLE  yes
244 DEPENDENCIES  none
245 .TE

247 This feature allows ZFS to maintain more information about how free space
248 is organized within the pool. If this feature is \fBenabled\fR, ZFS will
249 set this feature to \fBactive\fR when a new space map object is created or
250 an existing space map is upgraded to the new format. Once the feature is

```

```
251 \fBactive\fR, it will remain in that state until the pool is destroyed.  
252 .RE
```

```
254 .sp  
255 .ne 2  
256 .na  
257 \fB\fBmulti_vdev_crash_dump\fR\fR  
258 .ad  
259 .RS 4n  
260 .TS  
261 l l .  
262 GUID      com.joyent:multi_vdev_crash_dump  
263 READ\ -ONLY COMPATIBLE    no  
264 DEPENDENCIES      none  
265 .TE
```

```
267 This feature allows a dump device to be configured with a pool comprised  
268 of multiple vdevs. Those vdevs may be arranged in any mirrored or raidz  
269 configuration.
```

```
271 When the \fBmulti_vdev_crash_dump\fR feature is set to \fBenabled\fR,  
272 the administrator can use the \fBdumpadm\fR(1M) command to configure a  
273 dump device on a pool comprised of multiple vdevs.
```

```
275 .SH "SEE ALSO"  
276 \fBzpool\fR(1M)
```

```

*****
43700 Tue Sep 3 20:26:56 2013
new/usr/src/uts/common/Makefile.files
4101 metaslab_debug should allow for fine-grained control
4102 space_maps should store more information about themselves
4103 space_map object blocksize should be increased
4104 ::spa_space no longer works
4105 removing a mirrored log device results in a leaked object
4106 asynchronously load metaslab
Reviewed by: Matthew Ahrens <mahrens@delphix.com>
Reviewed by: Adam Leventhal <ahl@delphix.com>
Reviewed by: Sebastien Roy <seb@delphix.com>
*****
1 #
2 # CDDL HEADER START
3 #
4 # The contents of this file are subject to the terms of the
5 # Common Development and Distribution License (the "License").
6 # You may not use this file except in compliance with the License.
7 #
8 # You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9 # or http://www.opensolaris.org/os/licensing.
10 # See the License for the specific language governing permissions
11 # and limitations under the License.
12 #
13 # When distributing Covered Code, include this CDDL HEADER in each
14 # file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 # If applicable, add the following below this CDDL HEADER, with the
16 # fields enclosed by brackets "[]" replaced with your own identifying
17 # information: Portions Copyright [yyyy] [name of copyright owner]
18 #
19 # CDDL HEADER END
20 #
21 #
22 #
23 # Copyright (c) 1991, 2010, Oracle and/or its affiliates. All rights reserved.
24 # Copyright (c) 2012 Nexenta Systems, Inc. All rights reserved.
25 # Copyright (c) 2013 by Delphix. All rights reserved.
26 # Copyright (c) 2012 by Delphix. All rights reserved.
27 # Copyright (c) 2013 by Saso Kiselkov. All rights reserved.
28 #
29 #
30 # This Makefile defines all file modules for the directory uts/common
31 # and its children. These are the source files which may be considered
32 # common to all SunOS systems.
33 #
34 i386_CORE_OBJS += \
35     atomic.o      \
36     avintr.o      \
37     pic.o
38 #
39 sparc_CORE_OBJS +=
40 #
41 COMMON_CORE_OBJS += \
42     beep.o        \
43     bitset.o      \
44     bp_map.o      \
45     brand.o       \
46     cpucaps.o     \
47     cmt.o         \
48     cmt_policy.o  \
49     cpu.o         \
50     cpu_event.o   \
51     cpu_intr.o    \
52     cpu_pm.o      \

```

```

53     cpupart.o     \
54     cap_util.o    \
55     disp.o        \
56     group.o       \
57     kstat_fr.o    \
58     iscsiboot_prop.o \
59     lgrp.o         \
60     lgrp_topo.o   \
61     mmapobj.o     \
62     mutex.o       \
63     page_lock.o   \
64     page_retire.o \
65     panic.o       \
66     param.o       \
67     pg.o          \
68     pghw.o        \
69     putnext.o     \
70     rctl_proc.o   \
71     rwlock.o      \
72     seg_kmem.o    \
73     softint.o     \
74     string.o      \
75     strtol.o      \
76     strtoul.o     \
77     strtoll.o     \
78     strtoull.o   \
79     thread_intr.o \
80     vm_page.o     \
81     vm_pagelist.o \
82     zlib_obj.o    \
83     clock_tick.o
84 #
85 CORE_OBJS += $(COMMON_CORE_OBJS) $(MACH)_CORE_OBJS
86 #
87 ZLIB_OBJS = zutil.o zmod.o zmod_subr.o \
88     adler32.o crc32.o deflate.o inffast.o \
89     inflate.o inftrees.o trees.o
90 #
91 GENUNIX_OBJS += \
92     access.o      \
93     acl.o         \
94     acl_common.o \
95     adjtime.o     \
96     alarm.o       \
97     aio_subr.o    \
98     auditsys.o    \
99     audit_core.o  \
100    audit_zone.o  \
101    audit_memory.o \
102    autoconf.o    \
103    avl.o         \
104    bdev_dsort.o  \
105    bio.o         \
106    bitmap.o      \
107    blabel.o      \
108    brandsys.o    \
109    bz2blocksort.o \
110    bz2compress.o \
111    bz2decompress.o \
112    bz2randtable.o \
113    bz2zlib.o     \
114    bz2crctable.o \
115    bz2huffman.o  \
116    callb.o       \
117    callout.o     \
118    chdir.o       \

```

new/usr/src/uts/common/Makefile.files

```

119          chmod.o          \
120          chown.o          \
121          cladm.o          \
122          class.o          \
123          clock.o          \
124          clock_highres.o  \
125          clock_realtime.o \
126          close.o          \
127          compress.o       \
128          condvar.o        \
129          conf.o           \
130          console.o        \
131          contract.o       \
132          copyops.o        \
133          core.o           \
134          corectl.o        \
135          cred.o           \
136          cs_stubs.o       \
137          dacf.o           \
138          dacf_clnt.o      \
139          damap.o \
140          cyclic.o         \
141          ddi.o            \
142          ddifm.o          \
143          ddi_hp_impl.o    \
144          ddi_hp_ndi.o     \
145          ddi_intr.o       \
146          ddi_intr_impl.o  \
147          ddi_intr_irm.o   \
148          ddi_nodeid.o     \
149          ddi_periodic.o   \
150          devcfg.o         \
151          devcache.o       \
152          device.o         \
153          devid.o          \
154          devid_cache.o    \
155          devid_scsi.o     \
156          devid_smp.o      \
157          devpolicy.o      \
158          disp_lock.o      \
159          dnlc.o           \
160          driver.o         \
161          dumpsubr.o       \
162          driver_lyr.o     \
163          dtrace_subr.o    \
164          errorq.o         \
165          etheraddr.o     \
166          evchannels.o     \
167          exact.o          \
168          exact_core.o     \
169          exec.o           \
170          exit.o           \
171          fbio.o           \
172          fcntl.o          \
173          fdbuffer.o       \
174          fdsync.o         \
175          fem.o            \
176          ffs.o            \
177          fio.o            \
178          flock.o          \
179          fm.o             \
180          fork.o           \
181          vpm.o            \
182          fs_reparse.o     \
183          fs_subr.o        \
184          fsflush.o        \

```

3

new/usr/src/uts/common/Makefile.files

```

185          ftrace.o        \
186          getcwd.o         \
187          getdents.o       \
188          getloadavg.o     \
189          getpagesizes.o   \
190          getpid.o         \
191          gfs.o            \
192          rusagesys.o      \
193          gid.o            \
194          groups.o         \
195          grow.o           \
196          hat_refmod.o    \
197          id32.o           \
198          id_space.o       \
199          inet_ntop.o      \
200          instance.o       \
201          ioctl.o          \
202          ip_cksum.o       \
203          issetugid.o      \
204          ippconf.o        \
205          kcp.o            \
206          kdi.o            \
207          kiconv.o         \
208          klpd.o           \
209          kmem.o           \
210          ksyms_snapshot.o \
211          l_strplumb.o     \
212          labelsys.o       \
213          link.o           \
214          list.o           \
215          lockstat_subr.o  \
216          log_sysevent.o   \
217          logsubr.o        \
218          lookup.o         \
219          lseek.o          \
220          ltos.o           \
221          lwp.o            \
222          lwp_create.o     \
223          lwp_info.o       \
224          lwp_self.o       \
225          lwp_sobj.o       \
226          lwp_timer.o     \
227          lwpsys.o         \
228          main.o           \
229          mmapobjsys.o     \
230          memcntl.o        \
231          memstr.o         \
232          lgrpsys.o        \
233          mkdir.o          \
234          mknod.o          \
235          mount.o          \
236          move.o           \
237          msacct.o         \
238          multidata.o      \
239          nbmlck.o         \
240          ndifm.o          \
241          nice.o           \
242          netstack.o       \
243          ntptime.o        \
244          nvpair.o         \
245          nvpair_alloc_system.o \
246          nvpair_alloc_fixed.o \
247          fnvpair.o        \
248          octet.o          \
249          open.o           \
250          p_online.o       \

```

4

new/usr/src/uts/common/Makefile.files

```

251 pathconf.o \
252 pathname.o \
253 pause.o \
254 serializer.o \
255 pci_intr_lib.o \
256 pci_cap.o \
257 pcifm.o \
258 pgrp.o \
259 pgrpsys.o \
260 pid.o \
261 pkp_hash.o \
262 policy.o \
263 poll.o \
264 pool.o \
265 pool_pset.o \
266 port_subr.o \
267 ppriv.o \
268 printf.o \
269 priocntl.o \
270 priv.o \
271 priv_const.o \
272 proc.o \
273 procset.o \
274 processor_bind.o \
275 processor_info.o \
276 profil.o \
277 project.o \
278 qsort.o \
279 rctl.o \
280 rctlsys.o \
281 readlink.o \
282 refstr.o \
283 rename.o \
284 resolvepath.o \
285 retire_store.o \
286 process.o \
287 rlimit.o \
288 rmap.o \
289 rw.o \
290 rwstlock.o \
291 sad_conf.o \
292 sid.o \
293 sidsys.o \
294 sched.o \
295 schedctl.o \
296 sctp_crc32.o \
297 seg_dev.o \
298 seg_kp.o \
299 seg_kpm.o \
300 seg_map.o \
301 seg_vn.o \
302 seg_spt.o \
303 semaphore.o \
304 sendfile.o \
305 session.o \
306 share.o \
307 shuttle.o \
308 sig.o \
309 sigaction.o \
310 sigaltstack.o \
311 signotify.o \
312 sigpending.o \
313 sigprocmask.o \
314 sigqueue.o \
315 sigsendset.o \
316 sigsuspend.o \

```

5

new/usr/src/uts/common/Makefile.files

```

317 sigtimedwait.o \
318 sleepq.o \
319 sock_conf.o \
320 space.o \
321 sscanf.o \
322 stat.o \
323 statfs.o \
324 statvfs.o \
325 stol.o \
326 str_conf.o \
327 strcalls.o \
328 stream.o \
329 streamio.o \
330 strext.o \
331 strsubr.o \
332 strsun.o \
333 subr.o \
334 sunddi.o \
335 sunmdi.o \
336 sunndi.o \
337 sunpci.o \
338 sunpm.o \
339 sundlpi.o \
340 suntpi.o \
341 swap_subr.o \
342 swap_vnops.o \
343 symlink.o \
344 sync.o \
345 sysclass.o \
346 sysconfig.o \
347 sysent.o \
348 sysfs.o \
349 systeminfo.o \
350 task.o \
351 taskq.o \
352 tasksys.o \
353 time.o \
354 timer.o \
355 times.o \
356 timers.o \
357 thread.o \
358 tlabel.o \
359 tnf_res.o \
360 turnstile.o \
361 tty_common.o \
362 u8_textprep.o \
363 uadmin.o \
364 uconv.o \
365 ucredsys.o \
366 uid.o \
367 umask.o \
368 umount.o \
369 uname.o \
370 unix_bb.o \
371 unlink.o \
372 urw.o \
373 utime.o \
374 utssys.o \
375 uucopy.o \
376 vfs.o \
377 vfs_conf.o \
378 vmem.o \
379 vm_anon.o \
380 vm_as.o \
381 vm_meter.o \
382 vm_pageout.o \

```

6

new/usr/src/uts/common/Makefile.files

7

```

383          vm_pvn.o      \
384          vm_rm.o       \
385          vm_seg.o      \
386          vm_subr.o     \
387          vm_swap.o    \
388          vm_usage.o   \
389          vnode.o       \
390          vuid_queue.o  \
391          vuid_store.o  \
392          waitq.o       \
393          watchpoint.o \
394          yield.o       \
395          scsi_confdata.o \
396          xattr.o       \
397          xattr_common.o \
398          xdr_mblk.o    \
399          xdr_mem.o     \
400          xdr.o         \
401          xdr_array.o   \
402          xdr_refer.o   \
403          xhat.o        \
404          zone.o

406 #
407 #     Stubs for the stand-alone linker/loader
408 #
409 sparc_GENSTUBS_OBJS = \
410     kobj_stubs.o

412 i386_GENSTUBS_OBJS =

414 COMMON_GENSTUBS_OBJS =

416 GENSTUBS_OBJS += $(COMMON_GENSTUBS_OBJS) $($ (MACH)_GENSTUBS_OBJS)

418 #
419 #     DTrace and DTrace Providers
420 #
421 DTRACE_OBJS += dtrace.o dtrace_isa.o dtrace_asm.o

423 SDT_OBJS += sdt_subr.o

425 PROFILE_OBJS += profile.o

427 SYSTRACE_OBJS += systrace.o

429 LOCKSTAT_OBJS += lockstat.o

431 FASTTRAP_OBJS += fasttrap.o fasttrap_isa.o

433 DCPC_OBJS += dcpc.o

435 #
436 #     Driver (pseudo-driver) Modules
437 #
438 IPP_OBJS += ippctl.o

440 AUDIO_OBJS += audio_client.o audio_ddi.o audio_engine.o \
441     audio_fltdata.o audio_format.o audio_ctrl.o \
442     audio_grc3.o audio_output.o audio_input.o \
443     audio_oss.o audio_sun.o

445 AUDIOEMU10K_OBJS += audioemu10k.o

447 AUDIOENS_OBJS += audioens.o

```

new/usr/src/uts/common/Makefile.files

8

```

449 AUDIOVIA823X_OBJS += audiovia823x.o

451 AUDIOVIA97_OBJS += audiovia97.o

453 AUDIO1575_OBJS += audio1575.o

455 AUDIO810_OBJS += audio810.o

457 AUDIOCMI_OBJS += audiocmi.o

459 AUDIOCMIHD_OBJS += audiocmihd.o

461 AUDIOHD_OBJS += audiohd.o

463 AUDIOIXP_OBJS += audioixp.o

465 AUDIOLS_OBJS += audiols.o

467 AUDIOP16X_OBJS += audiop16x.o

469 AUDIOPCI_OBJS += audiopci.o

471 AUDIOSOLO_OBJS += audiosolo.o

473 AUDIOTS_OBJS += audiots.o

475 AC97_OBJS += ac97.o ac97_ad.o ac97_alc.o ac97_cmi.o

477 BLKDEV_OBJS += blkdev.o

479 CARDBUS_OBJS += cardbus.o cardbus_hp.o cardbus_cfg.o

481 CONSKBD_OBJS += conskbd.o

483 CONSMS_OBJS += consms.o

485 OLDPTY_OBJS += tty_ptyconf.o

487 PTC_OBJS += tty_pty.o

489 PTSL_OBJS += tty_pts.o

491 PTM_OBJS += ptm.o

493 MII_OBJS += mii.o mii_cicada.o mii_natsemi.o mii_intel.o mii_qualsemi.o \
494     mii_marvell.o mii_realtek.o mii_other.o

496 PTS_OBJS += pts.o

498 PTY_OBJS += ptms_conf.o

500 SAD_OBJS += sad.o

502 MD4_OBJS += md4.o md4_mod.o

504 MD5_OBJS += md5.o md5_mod.o

506 SHA1_OBJS += sha1.o sha1_mod.o

508 SHA2_OBJS += sha2.o sha2_mod.o

510 IPGPC_OBJS += classifierddi.o classifier.o filters.o trie.o table.o \
511     ba_table.o

513 DSCPMK_OBJS += dscpmk.o dscpmkddi.o

```

```

515 DLCOSMK_OBJS += dlcosmk.o dlcosmkddi.o
517 FLOWACCT_OBJS += flowacctddi.o flowacct.o
519 TOKENMT_OBJS += tokenmt.o tokenmtddi.o
521 TSWTCL_OBJS += tswtcl.o tswtclddi.o
523 ARP_OBJS += arpddi.o
525 ICMP_OBJS += icmpddi.o
527 ICMP6_OBJS += icmp6ddi.o
529 RTS_OBJS += rtsddi.o

531 IP_ICMP_OBJS = icmp.o icmp_opt_data.o
532 IP_RTS_OBJS = rts.o rts_opt_data.o
533 IP_TCP_OBJS = tcp.o tcp_fusion.o tcp_opt_data.o tcp_sack.o tcp_stats.o \
534 tcp_misc.o tcp_timers.o tcp_time_wait.o tcp_tpi.o tcp_output.o \
535 tcp_input.o tcp_socket.o tcp_bind.o tcp_cluster.o tcp_tunables.o
536 IP_UDP_OBJS = udp.o udp_opt_data.o udp_tunables.o udp_stats.o
537 IP_SCTP_OBJS = sctp.o sctp_opt_data.o sctp_output.o \
538 sctp_init.o sctp_input.o sctp_cookie.o \
539 sctp_conn.o sctp_error.o sctp_snmp.o \
540 sctp_tunables.o sctp_shutdown.o sctp_common.o \
541 sctp_timer.o sctp_heartbeat.o sctp_hash.o \
542 sctp_bind.o sctp_notify.o sctp_asconf.o \
543 sctp_addr.o tn_ipopt.o tnet.o ip_netinfo.o \
544 sctp_misc.o
545 IP_ILB_OBJS = ilb.o ilb_nat.o ilb_conn.o ilb_alg_hash.o ilb_alg_rr.o

547 IP_OBJS += igmp.o ipmp.o ip.o ip6.o ip6_asp.o ip6_if.o ip6_ire.o \
548 ip6_rts.o ip_if.o ip_ire.o ip_listutils.o ip_mroute.o \
549 ip_multi.o ip2mac.o ip_ndp.o ip_rts.o ip_srcid.o \
550 ipddi.o ipdrop.o mi.o nd.o tunables.o optcom.o snmpcom.o \
551 ipsec_loader.o spd.o ipclassifier.o inet_common.o ip_queue.o \
552 queue.o ip_sadb.o ip_ftable.o proto_set.o radix.o ip_dummy.o \
553 ip_helper_stream.o ip_tunables.o \
554 ip_output.o ip_input.o ip6_input.o ip6_output.o ip_arp.o \
555 conn_opt.o ip_attr.o ip_dce.o \
556 $(IP_ICMP_OBJS) \
557 $(IP_RTS_OBJS) \
558 $(IP_TCP_OBJS) \
559 $(IP_UDP_OBJS) \
560 $(IP_SCTP_OBJS) \
561 $(IP_ILB_OBJS)

563 IP6_OBJS += ip6ddi.o
565 HOOK_OBJS += hook.o
567 NETI_OBJS += neti_impl.o neti_mod.o neti_stack.o
569 KEYSOCK_OBJS += keysockddi.o keysock.o keysock_opt_data.o
571 IPNET_OBJS += ipnet.o ipnet_bpf.o
573 SPDSOCK_OBJS += spdsockddi.o spdsock.o spdsock_opt_data.o
575 IPSECESP_OBJS += ipsecespddi.o ipsecesp.o
577 IPSECAH_OBJS += ipsecahddi.o ipsecah.o sadb.o
579 SPPP_OBJS += sPPP.o sPPP_dlpi.o sPPP_mod.o sPPP_common.o

```

```

581 SPPPTUN_OBJS += sppptun.o sppptun_mod.o
583 SPPASYN_OBJS += spppasyn.o spppasyn_mod.o
585 SPPPCOMP_OBJS += sPPPcomp.o sPPPcomp_mod.o deflate.o bsd-comp.o vjcompress.o \
586 zlib.o
588 TCP_OBJS += tcpddi.o
590 TCP6_OBJS += tcp6ddi.o
592 NCA_OBJS += ncaddi.o
594 SDP SOCK_MOD_OBJS += sockmod_sdp.o socksdp.o socksdp_subr.o
596 SCTP SOCK_MOD_OBJS += sockmod_sctp.o socksctp.o socksctp_subr.o
598 PFP SOCK_MOD_OBJS += sockmod_pfp.o
600 RDS SOCK_MOD_OBJS += sockmod_rds.o
602 RDS_OBJS += rdsddi.o rdssubr.o rds_opt.o rds_ioctl.o
604 RDSIB_OBJS += rdsib.o rdsib_ib.o rdsib_cm.o rdsib_ep.o rdsib_buf.o \
605 rdsib_debug.o rdsib_sc.o
607 RDSV3_OBJS += af_rds.o rds_v3_ddi.o bind.o loop.o threads.o connection.o \
608 transport.o cong.o sysctl.o message.o rds_recv.o send.o \
609 stats.o info.o page.o rdma_transport.o ib_ring.o ib_rdma.o \
610 ib_recv.o ib.o ib_send.o ib_sysctl.o ib_stats.o ib_cm.o \
611 rds_v3_sc.o rds_v3_debug.o rds_v3_impl.o rdma.o rds_v3_af_thr.o
613 ISER_OBJS += iser.o iser_cm.o iser_cg.o iser_ib.o iser_idm.o \
614 iser_resource.o iser_xfer.o
616 UDP_OBJS += udpddi.o
618 UDP6_OBJS += udp6ddi.o
620 SY_OBJS += genty.o
622 TCO_OBJS += ticots.o
624 TCOO_OBJS += ticotsord.o
626 TCL_OBJS += ticlts.o
628 TL_OBJS += tl.o
630 DUMP_OBJS += dump.o
632 BPF_OBJS += bpf.o bpf_filter.o bpf_mod.o bpf_dlt.o bpf_mac.o
634 CLONE_OBJS += clone.o
636 CN_OBJS += cons.o
638 DLD_OBJS += dld_drv.o dld_proto.o dld_str.o dld_flow.o
640 DLS_OBJS += dls.o dls_link.o dls_mod.o dls_stat.o dls_mgmt.o
642 GLD_OBJS += gld.o gldutil.o
644 MAC_OBJS += mac.o mac_bcast.o mac_client.o mac_datapath_setup.o mac_flow.o
645 mac_hio.o mac_mod.o mac_ndd.o mac_provider.o mac_sched.o \
646 mac_protect.o mac_soft_ring.o mac_stat.o mac_util.o

```

```

648 MAC_6TO4_OBJS +=      mac_6to4.o
650 MAC_ETHER_OBJS +=     mac_ether.o
652 MAC_IPV4_OBJS +=     mac_ipv4.o
654 MAC_IPV6_OBJS +=     mac_ipv6.o
656 MAC_WIFI_OBJS +=     mac_wifi.o
658 MAC_IB_OBJS +=       mac_ib.o
660 IPTUN_OBJS +=        iptun_dev.o iptun_ctl.o iptun.o
662 AGGR_OBJS +=         aggr_dev.o aggr_ctl.o aggr_grp.o aggr_port.o \
663                       aggr_send.o aggr_recv.o aggr_lacp.o
665 SOFTMAC_OBJS +=      softmac_main.o softmac_ctl.o softmac_capab.o \
666                       softmac_dev.o softmac_stat.o softmac_pkt.o softmac_fp.o
668 NET80211_OBJS +=     net80211.o net80211_proto.o net80211_input.o \
669                       net80211_output.o net80211_node.o net80211_crypto.o \
670                       net80211_crypto_none.o net80211_crypto_wep.o net80211_ioctl.o \
671                       net80211_crypto_tkip.o net80211_crypto_ccmp.o \
672                       net80211_ht.o
674 VNIC_OBJS +=        vnic_ctl.o vnic_dev.o
676 SIMNET_OBJS +=      simnet.o
678 IB_OBJS +=          ibnex.o ibnex_ioctl.o ibnex_hca.o
680 IBCM_OBJS +=        ibcm_impl.o ibcm_sm.o ibcm_ti.o ibcm_utils.o ibcm_path.o \
681                       ibcm_arp.o ibcm_arp_link.o
683 IBDM_OBJS +=        ibdm.o
685 IBDMA_OBJS +=       ibdma.o
687 IBMF_OBJS +=        ibmf.o ibmf_impl.o ibmf_dr.o ibmf_wqe.o ibmf_ud_dest.o ibmf_mod.o
688                       ibmf_send.o ibmf_recv.o ibmf_handlers.o ibmf_trans.o \
689                       ibmf_timers.o ibmf_msg.o ibmf_utils.o ibmf_rmpp.o \
690                       ibmf_saa.o ibmf_saa_impl.o ibmf_saa_utils.o ibmf_saa_events.o
692 IBTL_OBJS +=        ibtl_impl.o ibtl_util.o ibtl_mem.o ibtl_handlers.o ibtl_qp.o \
693                       ibtl_cq.o ibtl_wr.o ibtl_hca.o ibtl_chan.o ibtl_cm.o \
694                       ibtl_mcg.o ibtl_ibnex.o ibtl_srql.o ibtl_part.o
696 TAVOR_OBJS +=       tavor.o tavor_agents.o tavor_cfg.o tavor_ci.o tavor_cmd.o \
697                       tavor_cq.o tavor_event.o tavor_ioctl.o tavor_misc.o \
698                       tavor_mr.o tavor_qp.o tavor_gpmod.o tavor_rsrc.o \
699                       tavor_srql.o tavor_stats.o tavor_umap.o tavor_wr.o
701 HERMON_OBJS +=      hermon.o hermon_agents.o hermon_cfg.o hermon_ci.o hermon_cmd.o \
702                       hermon_cq.o hermon_event.o hermon_ioctl.o hermon_misc.o \
703                       hermon_mr.o hermon_qp.o hermon_gpmod.o hermon_rsrc.o \
704                       hermon_srql.o hermon_stats.o hermon_umap.o hermon_wr.o \
705                       hermon_fcoib.o hermon_fm.o
707 DAPLT_OBJS +=       daplt.o
709 SOL_OFS_OBJS +=     sol_cma.o sol_ib_cma.o sol_uobj.o \
710                       sol_ofs_debug_util.o sol_ofs_gen_util.o \
711                       sol_kverbs.o

```

```

713 SOL_UCMA_OBJS +=    sol_ucma.o
715 SOL_UVERBS_OBJS +=  sol_uverbs.o sol_uverbs_comp.o sol_uverbs_event.o \
716                       sol_uverbs_hca.o sol_uverbs_qp.o
718 SOL_UMAD_OBJS +=    sol_umad.o
720 KSTAT_OBJS +=      kstat.o
722 KSYMS_OBJS +=      ksyms.o
724 INSTANCE_OBJS +=   inst_sync.o
726 IWSCN_OBJS +=      iwscons.o
728 LOFI_OBJS +=       lofi.o LzmaDec.o
730 FSSNAP_OBJS +=     fssnap.o
732 FSSNAPIF_OBJS +=   fssnap_if.o
734 MM_OBJS +=          mem.o
736 PHYSMEM_OBJS +=    physmem.o
738 OPTIONS_OBJS +=    options.o
740 WINLOCK_OBJS +=    winlockio.o
742 PM_OBJS +=          pm.o
743 SRN_OBJS +=          srn.o
745 PSEUDO_OBJS +=     pseudonex.o
747 RAMDISK_OBJS +=    ramdisk.o
749 LLC1_OBJS +=       llc1.o
751 USBKBM_OBJS +=     usbkbm.o
753 USBWCM_OBJS +=     usbwcm.o
755 BOFI_OBJS +=       bofi.o
757 HID_OBJS +=        hid.o
759 HWA_RC_OBJS +=     hwarc.o
761 USBSKEL_OBJS +=     usbskel.o
763 USBVC_OBJS +=       usbvc.o usbvc_v412.o
765 HIDPARSER_OBJS +=  hidparser.o
767 USB_AC_OBJS +=      usb_ac.o
769 USB_AS_OBJS +=      usb_as.o
771 USB_AH_OBJS +=      usb_ah.o
773 USBMS_OBJS +=       usbms.o
775 USBPRN_OBJS +=      usbprn.o
777 UGEN_OBJS +=        ugen.o

```

```

779 USBSER_OBJS += usbser.o usbser_rseq.o
781 USBSACM_OBJS += usbsacm.o
783 USBSER_KEYSPAN_OBJS += usbser_keyspan.o keyspan_dsd.o keyspan_pipe.o
785 USBS49_FW_OBJS += keyspan_49fw.o
787 USBSPRL_OBJS += usbser_pl2303.o pl2303_dsd.o
789 WUSB_CA_OBJS += wusb_ca.o
791 USBFTDI_OBJS += usbser_uftdi.o uftdi_dsd.o
793 USBECM_OBJS += usbecm.o
795 WC_OBJS += wscons.o vcons.o
797 VCONS_CONF_OBJS += vcons_conf.o
799 SCSI_OBJS +=      scsi_capabilities.o scsi_confsubr.o scsi_control.o \
800                scsi_data.o scsi_fm.o scsi_hba.o scsi_reset_notify.o \
801                scsi_resource.o scsi_subr.o scsi_transport.o scsi_watch.o \
802                smp_transport.o
804 SCSI_VHCI_OBJS +=      scsi_vhci.o mpapi_impl.o scsi_vhci_tpgs.o
806 SCSI_VHCI_F_SYM_OBJS +=      sym.o
808 SCSI_VHCI_F_TPGS_OBJS +=      tpgs.o
810 SCSI_VHCI_F_ASYM_SUN_OBJS +=  asym_sun.o
812 SCSI_VHCI_F_SYM_HDS_OBJS +=  sym_hds.o
814 SCSI_VHCI_F_TAPE_OBJS +=      tape.o
816 SCSI_VHCI_F_TPGS_TAPE_OBJS += tpgs_tape.o
818 SGEN_OBJS +=      sgen.o
820 SMP_OBJS +=      smp.o
822 SATA_OBJS +=      sata.o
824 USBA_OBJS +=      hcdi.o usba.o usbai.o hubdi.o parser.o genconsole.o \
825                usbai_pipe_mgmt.o usbai_req.o usbai_util.o usbai_register.o \
826                usba_devdb.o usbal0_calls.o usba_ugen.o whcdi.o wa.o
827 USBA_WITHOUT_WUSB_OBJS +=      hcdi.o usba.o usbai.o hubdi.o parser.o gencons
828                usbai_pipe_mgmt.o usbai_req.o usbai_util.o usbai_register.o \
829                usba_devdb.o usbal0_calls.o usba_ugen.o
831 USBA10_OBJS +=      usbal0.o
833 RSM_OBJS +=      rsm.o rsmka_pathmanager.o rsmka_util.o
835 RSMOPS_OBJS +=      rsmops.o
837 S1394_OBJS +=      t1394.o t1394_errmsg.o s1394.o s1394_addr.o s1394_async.o \
838                s1394_bus_reset.o s1394_cmp.o s1394_csr.o s1394_dev_disc.o \
839                s1394_fa.o s1394_fcp.o \
840                s1394_hotplug.o s1394_isoch.o s1394_misc.o h1394.o nx1394.o
842 HCI1394_OBJS +=      hcil1394.o hcil1394_async.o hcil1394_attach.o hcil1394_buf.o \
843                hcil1394_csr.o hcil1394_detach.o hcil1394_extern.o \
844                hcil1394_ioctl.o hcil1394_isoch.o hcil1394_isr.o \

```

```

845                hcil1394_ixl_comp.o hcil1394_ixl_isr.o hcil1394_ixl_misc.o \
846                hcil1394_ixl_update.o hcil1394_misc.o hcil1394_ohci.o \
847                hcil1394_g.o hcil1394_s1394if.o hcil1394_tlabel.o \
848                hcil1394_tlist.o hcil1394_vendor.o
850 AV1394_OBJS +=      av1394.o av1394_as.o av1394_async.o av1394_cfgrom.o \
851                av1394_cmp.o av1394_fcp.o av1394_isoch.o av1394_isoch_chan.o \
852                av1394_isoch_recv.o av1394_isoch_xmit.o av1394_list.o \
853                av1394_queue.o
855 DCAM1394_OBJS +=      dcam.o dcam_frame.o dcam_param.o dcam_reg.o \
856                dcam_ring_buff.o
858 SCSA1394_OBJS +=      hba.o sbp2_driver.o sbp2_bus.o
860 SBP2_OBJS +=      cfgrom.o sbp2.o
862 PMODEM_OBJS +=      pmodem.o pmodem_cis.o cis.o cis_callout.o cis_handlers.o cis_para
864 DSW_OBJS +=      dsw.o dsw_dev.o ii_tree.o
866 NCALL_OBJS +=      ncall.o \
867                ncall_stub.o
869 RDC_OBJS +=      rdc.o \
870                rdc_dev.o \
871                rdc_io.o \
872                rdc_clnt.o \
873                rdc_prot_xdr.o \
874                rdc_svc.o \
875                rdc_bitmap.o \
876                rdc_health.o \
877                rdc_subr.o \
878                rdc_diskq.o
880 RDCSRV_OBJS +=      rdcsrv.o
882 RDCSTUB_OBJS +=      rdc_stub.o
884 SDBC_OBJS +=      sd_bcache.o \
885                sd_bio.o \
886                sd_conf.o \
887                sd_ft.o \
888                sd_hash.o \
889                sd_io.o \
890                sd_misc.o \
891                sd_pcu.o \
892                sd_tdaemon.o \
893                sd_trace.o \
894                sd_iob_impl0.o \
895                sd_iob_impl1.o \
896                sd_iob_impl2.o \
897                sd_iob_impl3.o \
898                sd_iob_impl4.o \
899                sd_iob_impl5.o \
900                sd_iob_impl6.o \
901                sd_iob_impl7.o \
902                safestore.o \
903                safestore_ram.o
905 NSCTL_OBJS +=      nsctl.o \
906                nsc_cache.o \
907                nsc_disk.o \
908                nsc_dev.o \
909                nsc_freeze.o \
910                nsc_gen.o \

```

```

911         nsc_mem.o \
912         nsc_ncallio.o \
913         nsc_power.o \
914         nsc_resv.o \
915         nsc_rmspin.o \
916         nsc_solaris.o \
917         nsc_trap.o \
918         nsc_list.o
919 UNISTAT_OBJS += spuni.o \
920                spcs_s_k.o

922 NSKERN_OBJS += nsc_ddi.o \
923                nsc_proc.o \
924                nsc_raw.o \
925                nsc_thread.o \
926                nskernd.o

928 SV_OBJS += sv.o

930 PMCS_OBJS += pmcs_attach.o pmcs_ds.o pmcs_intr.o pmcs_nvram.o pmcs_sata.o \
931             pmcs_scsa.o pmcs_smhba.o pmcs_subr.o pmcs_fwlog.o

933 PMCS8001FW_C_OBJS += pmcs_fw_hdr.o
934 PMCS8001FW_OBJS += $(PMCS8001FW_C_OBJS) SPCBoot.o ila.o firmware.o

936 #
937 #   Build up defines and paths.

939 ST_OBJS += st.o st_conf.o

941 EMLXS_OBJS += emlxs_clock.o emlxs_dfc.o emlxs_dhchap.o emlxs_diag.o \
942             emlxs_download.o emlxs_dump.o emlxs_els.o emlxs_event.o \
943             emlxs_fcf.o emlxs_fcp.o emlxs_fct.o emlxs_hba.o emlxs_ip.o \
944             emlxs_mbox.o emlxs_mem.o emlxs_msg.o emlxs_node.o \
945             emlxs_pkt.o emlxs_sli3.o emlxs_sli4.o emlxs_solaris.o \
946             emlxs_thread.o

948 EMLXS_FW_OBJS += emlxs_fw.o

950 OCE_OBJS += oce_buf.o oce_fm.o oce_gld.o oce_hw.o oce_intr.o oce_main.o \
951            oce_mbx.o oce_mq.o oce_queue.o oce_rx.o oce_stat.o oce_tx.o \
952            oce_utils.o

954 FCT_OBJS += discovery.o fct.o

956 QLT_OBJS += 2400.o 2500.o 8100.o qlt.o qlt_dma.o

958 SRPT_OBJS += srpt_mod.o srpt_ch.o srpt_cm.o srpt_ioc.o srpt_stp.o

960 FCOE_OBJS += fcoe.o fcoe_eth.o fcoe_fc.o

962 FCOET_OBJS += fcoet.o fcoet_eth.o fcoet_fc.o

964 FCOEI_OBJS += fcoei.o fcoei_eth.o fcoei_lv.o

966 ISCSIT_SHARED_OBJS += \
967                    iscsit_common.o

969 ISCSIT_OBJS += $(ISCSIT_SHARED_OBJS) \
970              iscsit.o iscsit_tgt.o iscsit_sess.o iscsit_login.o \
971              iscsit_text.o iscsit_isns.o iscsit_radiusauth.o \
972              iscsit_radiuspacket.o iscsit_auth.o iscsit_authclient.o

974 PPPT_OBJS += alua_ic_if.o pppt.o pppt_msg.o pppt_tgt.o

976 STMF_OBJS += lun_map.o stmf.o

```

```

978 STMF_SBD_OBJS += sbd.o sbd_scsi.o sbd_pgr.o sbd_zvol.o

980 SYMSMSG_OBJS += sysmsg.o

982 SES_OBJS += ses.o ses_sen.o ses_safte.o ses_ses.o

984 TNF_OBJS += tnf_buf.o tnf_trace.o tnf_writer.o trace_init.o \
985            trace_funcs.o tnf_probe.o tnf.o

987 LOGINDMUX_OBJS += logindmux.o

989 DEVINFO_OBJS += devinfo.o

991 DEVPOLL_OBJS += devpoll.o

993 DEVPOOL_OBJS += devpool.o

995 I8042_OBJS += i8042.o

997 KB8042_OBJS += \
998             at_keyprocess.o \
999             kb8042.o \
1000            kb8042_keytables.o

1002 MOUSE8042_OBJS += mouse8042.o

1004 FDC_OBJS += fdc.o

1006 ASY_OBJS += asy.o

1008 ECPP_OBJS += ecpp.o

1010 VUIDM3P_OBJS += vuidmice.o vuidm3p.o

1012 VUIDM4P_OBJS += vuidmice.o vuidm4p.o

1014 VUIDM5P_OBJS += vuidmice.o vuidm5p.o

1016 VUIDPS2_OBJS += vuidmice.o vuidps2.o

1018 HPCSVCS_OBJS += hpcsvc.o

1020 PCIE_MISC_OBJS += pcie.o pcie_fault.o pcie_hp.o pciehpc.o pcishpc.o pcie_pwr.o p

1022 PCIHPNEXUS_OBJS += pcihp.o

1024 OPENEEPROM_OBJS += openprom.o

1026 RANDOM_OBJS += random.o

1028 PSHOT_OBJS += pshot.o

1030 GEN_DRV_OBJS += gen_drv.o

1032 TCLIENT_OBJS += tclient.o

1034 TPHCI_OBJS += tphci.o

1036 TVHCI_OBJS += tvhci.o

1038 EMUL64_OBJS += emul64.o emul64_bsd.o

1040 FCP_OBJS += fcp.o

1042 FCIP_OBJS += fcip.o

```

```

1044 FCSM_OBJS += fcsm.o
1046 FCTL_OBJS += fctl.o
1048 FP_OBJS += fp.o
1050 QLC_OBJS += ql_api.o ql_debug.o ql_hba_fru.o ql_init.o ql_iocb.o ql_ioctl.o \
1051     ql_isr.o ql_mbx.o ql_nx.o ql_xioctl.o ql_fw_table.o
1053 QLC_FW_2200_OBJS += ql_fw_2200.o
1055 QLC_FW_2300_OBJS += ql_fw_2300.o
1057 QLC_FW_2400_OBJS += ql_fw_2400.o
1059 QLC_FW_2500_OBJS += ql_fw_2500.o
1061 QLC_FW_6322_OBJS += ql_fw_6322.o
1063 QLC_FW_8100_OBJS += ql_fw_8100.o
1065 QLGE_OBJS += qlge.o qlge_dbg.o qlge_flash.o qlge_fm.o qlge_gld.o qlge_mpi.o
1067 ZCONS_OBJS += zcons.o
1069 NV_SATA_OBJS += nv_sata.o
1071 SI3124_OBJS += si3124.o
1073 AHCI_OBJS += ahci.o
1075 PCIIDE_OBJS += pci-ide.o
1077 PCEPP_OBJS += pcepp.o
1079 CPC_OBJS += cpc.o
1081 CPUID_OBJS += cpuid_drv.o
1083 SYSEVENT_OBJS += sysevent.o
1085 BL_OBJS += bl.o
1087 DRM_OBJS += drm_sunmod.o drm_kstat.o drm_agpsupport.o \
1088     drm_auth.o drm_bufs.o drm_context.o drm_dma.o \
1089     drm_drawable.o drm_drv.o drm_fops.o drm_ioctl.o drm_irq.o \
1090     drm_lock.o drm_memory.o drm_msg.o drm_pci.o drm_scatter.o \
1091     drm_cache.o drm_gem.o drm_mm.o ati_pigart.o
1093 FM_OBJS += devfm.o devfm_machdep.o
1095 RTLS_OBJS += rtls.o
1097 #
1098 #           exec modules
1099 #
1100 AOUTEXEC_OBJS += aout.o
1102 ELFEXEC_OBJS += elf.o elf_notes.o old_notes.o
1104 INTPEXEC_OBJS += intp.o
1106 SHBINEXEC_OBJS += shbin.o
1108 JAVAEXEC_OBJS += java.o

```

```

1110 #
1111 #           file system modules
1112 #
1113 AUTOFS_OBJS += auto_vfsops.o auto_vnops.o auto_subr.o auto_xdr.o auto_sys.o
1115 CACHEFS_OBJS += cachefs_cnode.o      cachefs_cod.o \
1116     cachefs_dir.o      cachefs_dlog.o  cachefs_filegrp.o \
1117     cachefs_fscache.o  cachefs_ioctl.o cachefs_log.o \
1118     cachefs_module.o \
1119     cachefs_noopc.o      cachefs_resource.o \
1120     cachefs_strict.o \
1121     cachefs_subr.o      cachefs_vfsops.o \
1122     cachefs_vnops.o
1124 DCFS_OBJS += dc_vnops.o
1126 DEVFS_OBJS += devfs_subr.o  devfs_vfsops.o  devfs_vnops.o
1128 DEV_OBJS += sdev_subr.o      sdev_vfsops.o  sdev_vnops.o \
1129     sdev_ptsops.o  sdev_zvolops.o  sdev_comm.o \
1130     sdev_profile.o sdev_ncache.o  sdev_netops.o \
1131     sdev_ipnetops.o \
1132     sdev_vtoprs.o
1134 CTFS_OBJS += ctfs_all.o ctfs_cdir.o ctfs_ctl.o ctfs_event.o \
1135     ctfs_latest.o ctfs_root.o ctfs_sym.o ctfs_tdir.o ctfs_tmpl.o
1137 OBJFS_OBJS += objfs_vfs.o      objfs_root.o  objfs_common.o \
1138     objfs_odir.o  objfs_data.o
1140 FDFS_OBJS += fdops.o
1142 FIFO_OBJS += fifosubr.o      fifovnops.o
1144 PIPE_OBJS += pipe.o
1146 HSFS_OBJS += hsfs_node.o      hsfs_subr.o  hsfs_vfsops.o  hsfs_vnops.o \
1147     hsfs_susp.o  hsfs_rrip.o  hsfs_susp_subr.o
1149 LOFS_OBJS += lofs_subr.o      lofs_vfsops.o  lofs_vnops.o
1151 NAMEFS_OBJS += namevfs.o      namevno.o
1153 NFS_OBJS += nfs_client.o      nfs_common.o  nfs_dump.o \
1154     nfs_subr.o      nfs_vfsops.o  nfs_vnops.o \
1155     nfs_xdr.o      nfs_sys.o      nfs_strerror.o \
1156     nfs3_vfsops.o  nfs3_vnops.o  nfs3_xdr.o \
1157     nfs_acl_vnops.o nfs_acl_xdr.o  nfs4_vfsops.o \
1158     nfs4_vnops.o  nfs4_xdr.o  nfs4_idmap.o \
1159     nfs4_shadow.o nfs4_subr.o \
1160     nfs4_attr.o  nfs4_rnode.o  nfs4_client.o \
1161     nfs4_acache.o nfs4_common.o  nfs4_client_state.o \
1162     nfs4_callback.o nfs4_recovery.o nfs4_client_secinfo.o \
1163     nfs4_client_debug.o  nfs_stats.o \
1164     nfs4_acl.o  nfs4_stub_vnops.o  nfs_cmd.o
1166 NFSSRV_OBJS += nfs_server.o  nfs_srv.o      nfs3_srv.o \
1167     nfs_acl_srv.o  nfs_auth.o  nfs_auth_xdr.o \
1168     nfs_export.o  nfs_log.o  nfs_log_xdr.o \
1169     nfs4_srv.o  nfs4_state.o  nfs4_srv_attr.o \
1170     nfs4_srv_ns.o  nfs4_db.o  nfs4_srv_deleg.o \
1171     nfs4_deleg_ops.o  nfs4_srv_readdir.o  nfs4_dispatch.o
1173 SMBSRV_SHARED_OBJS += \
1174     smb_inet.o \

```

new/usr/src/uts/common/Makefile.files

19

```

1175         smb_match.o \
1176         smb_msgbuf.o \
1177         smb_oem.o \
1178         smb_string.o \
1179         smb_utf8.o \
1180         smb_door_legacy.o \
1181         smb_xdr.o \
1182         smb_token.o \
1183         smb_token_xdr.o \
1184         smb_sid.o \
1185         smb_native.o \
1186         smb_netbios_util.o

1188 SMBSRV_OBJS += $(SMBSRV_SHARED_OBJS) \
1189         smb_acl.o \
1190         smb_alloc.o \
1191         smb_close.o \
1192         smb_common_open.o \
1193         smb_common_transact.o \
1194         smb_create.o \
1195         smb_delete.o \
1196         smb_directory.o \
1197         smb_dispatch.o \
1198         smb_echo.o \
1199         smb_fem.o \
1200         smb_find.o \
1201         smb_flush.o \
1202         smb_fsinfo.o \
1203         smb_fsops.o \
1204         smb_init.o \
1205         smb_kdoor.o \
1206         smb_kshare.o \
1207         smb_kutil.o \
1208         smb_lock.o \
1209         smb_lock_byte_range.o \
1210         smb_locking_andx.o \
1211         smb_logoff_andx.o \
1212         smb_mangle_name.o \
1213         smb_mbuf_marshall.o \
1214         smb_mbuf_util.o \
1215         smb_negotiate.o \
1216         smb_net.o \
1217         smb_node.o \
1218         smb_nt_cancel.o \
1219         smb_nt_create_andx.o \
1220         smb_nt_transact_create.o \
1221         smb_nt_transact_ioctl.o \
1222         smb_nt_transact_notify_change.o \
1223         smb_nt_transact_quota.o \
1224         smb_nt_transact_security.o \
1225         smb_odir.o \
1226         smb_ofile.o \
1227         smb_open_andx.o \
1228         smb_opipe.o \
1229         smb_oplock.o \
1230         smb_pathname.o \
1231         smb_print.o \
1232         smb_process_exit.o \
1233         smb_query_fileinfo.o \
1234         smb_read.o \
1235         smb_rename.o \
1236         smb_sd.o \
1237         smb_seek.o \
1238         smb_server.o \
1239         smb_session.o \
1240         smb_session_setup_andx.o

```

new/usr/src/uts/common/Makefile.files

20

```

1241         smb_set_fileinfo.o \
1242         smb_signing.o \
1243         smb_tree.o \
1244         smb_trans2_create_directory.o \
1245         smb_trans2_dfs.o \
1246         smb_trans2_find.o \
1247         smb_tree_connect.o \
1248         smb_unlock_byte_range.o \
1249         smb_user.o \
1250         smb_vfs.o \
1251         smb_vops.o \
1252         smb_vss.o \
1253         smb_write.o \
1254         smb_write_raw.o

1256 PCFS_OBJS += pc_alloc.o pc_dir.o pc_node.o pc_subr.o \
1257         pc_vfsops.o pc_vnops.o

1259 PROC_OBJS += prcontrol.o prioctl.o prsubr.o prusr.o \
1260         prvfops.o prvnops.o

1262 MNTFS_OBJS += mntvfops.o mntvnops.o

1264 SHAREFS_OBJS += sharetab.o sharefs_vfsops.o sharefs_vnops.o

1266 SPEC_OBJS += specsubr.o specvfops.o specvnops.o

1268 SOCK_OBJS += socksubr.o sockvfops.o sockparams.o \
1269         socksyscalls.o socktpi.o sockstr.o \
1270         sockcommon_vnops.o sockcommon_subr.o \
1271         sockcommon_sops.o sockcommon.o \
1272         sock_notsupp.o socknotify.o \
1273         nl7c.o nl7curi.o nl7chttp.o nl7clogd.o \
1274         nl7cnca.o sodirect.o sockfilter.o

1276 TMPFS_OBJS += tmp_dir.o tmp_subr.o tmp_tnode.o tmp_vfsops.o \
1277         tmp_vnops.o

1279 UDFS_OBJS += udf_alloc.o udf_bmap.o udf_dir.o \
1280         udf_inode.o udf_subr.o udf_vfsops.o \
1281         udf_vnops.o

1283 UFS_OBJS += ufs_alloc.o ufs_bmap.o ufs_dir.o ufs_xattr.o \
1284         ufs_inode.o ufs_subr.o ufs_tables.o ufs_vfsops.o \
1285         ufs_vnops.o quota.o quotacalls.o quota_ufs.o \
1286         ufs_filio.o ufs_lockfs.o ufs_thread.o ufs_trans.o \
1287         ufs_acl.o ufs_panic.o ufs_directio.o ufs_log.o \
1288         ufs_extvnops.o ufs_snap.o lufs.o lufs_thread.o \
1289         lufs_log.o lufs_map.o lufs_top.o lufs_debug.o

1290 VSCAN_OBJS += vscan_drv.o vscan_svc.o vscan_door.o

1292 NSMB_OBJS += smb_conn.o smb_dev.o smb_iod.o smb_pass.o \
1293         smb_rq.o smb_sign.o smb_smb.o smb_subrs.o \
1294         smb_time.o smb_tran.o smb_trantcp.o smb_usr.o \
1295         subr_mchain.o

1297 SMBFS_COMMON_OBJS += smbfs_ntacl.o
1298 SMBFS_OBJS += smbfs_vfsops.o smbfs_vnops.o smbfs_node.o \
1299         smbfs_acl.o smbfs_client.o smbfs_smb.o \
1300         smbfs_subr.o smbfs_subr2.o \
1301         smbfs_rwlock.o smbfs_xattr.o \
1302         $(SMBFS_COMMON_OBJS)

1305 #
1306 # LVM modules

```



```

1307 #
1308 MD_OBJS += md.o md_error.o md_ioctl.o md_mddb.o md_names.o \
1309     md_med.o md_rename.o md_subr.o

1311 MD_COMMON_OBJS = md_convert.o md_crc.o md_revchk.o

1313 MD_DERIVED_OBJS = metamed_xdr.o meta_basic_xdr.o

1315 SOFTPART_OBJS += sp.o sp_ioctl.o

1317 STRIPE_OBJS += stripe.o stripe_ioctl.o

1319 HOTSPARES_OBJS += hotspares.o

1321 RAID_OBJS += raid.o raid_ioctl.o raid_replay.o raid_resync.o raid_hotspare.o

1323 MIRROR_OBJS += mirror.o mirror_ioctl.o mirror_resync.o

1325 NOTIFY_OBJS += md_notify.o

1327 TRANS_OBJS += mdtrans.o trans_ioctl.o trans_log.o

1329 ZFS_COMMON_OBJS += \
1330     arc.o \
1331     bplist.o \
1332     bpobj.o \
1333     bptree.o \
1334     dbuf.o \
1335     ddt.o \
1336     ddt_zap.o \
1337     dmuf.o \
1338     dmuf_diff.o \
1339     dmuf_send.o \
1340     dmuf_object.o \
1341     dmuf_objset.o \
1342     dmuf_traverse.o \
1343     dmuf_tx.o \
1344     dnode.o \
1345     dnode_sync.o \
1346     dsl_dir.o \
1347     dsl_dataset.o \
1348     dsl_deadlist.o \
1349     dsl_destroy.o \
1350     dsl_pool.o \
1351     dsl_synctask.o \
1352     dsl_userhold.o \
1353     dmuf_zfetch.o \
1354     dsl_deleg.o \
1355     dsl_prop.o \
1356     dsl_scan.o \
1357     zfeature.o \
1358     gzip.o \
1359     lz4.o \
1360     lzjb.o \
1361     metaslab.o \
1362     range_tree.o \
1363     refcount.o \
1364     rrwlock.o \
1365     sa.o \
1366     sha256.o \
1367     spa.o \
1368     spa_config.o \
1369     spa_errlog.o \
1370     spa_history.o \
1371     spa_misc.o \
1372     space_map.o \

```

```

1373     space_reftree.o \
1374     txg.o \
1375     uberblock.o \
1376     unique.o \
1377     vdev.o \
1378     vdev_cache.o \
1379     vdev_file.o \
1380     vdev_label.o \
1381     vdev_mirror.o \
1382     vdev_missing.o \
1383     vdev_queue.o \
1384     vdev_raidz.o \
1385     vdev_root.o \
1386     zap.o \
1387     zap_leaf.o \
1388     zap_micro.o \
1389     zfs_byteswap.o \
1390     zfs_debug.o \
1391     zfs_fm.o \
1392     zfs_fuid.o \
1393     zfs_sa.o \
1394     zfs_znode.o \
1395     zil.o \
1396     zio.o \
1397     zio_checksum.o \
1398     zio_compress.o \
1399     zio_inject.o \
1400     zle.o \
1401     zrlock.o

1403 ZFS_SHARED_OBJS += \
1404     zfeature_common.o \
1405     zfs_comutil.o \
1406     zfs_deleg.o \
1407     zfs_fletcher.o \
1408     zfs_namecheck.o \
1409     zfs_prop.o \
1410     zpool_prop.o \
1411     zprop_common.o

1413 ZFS_OBJS += \
1414     $(ZFS_COMMON_OBJS) \
1415     $(ZFS_SHARED_OBJS) \
1416     vdev_disk.o \
1417     zfs_acl.o \
1418     zfs_ctldir.o \
1419     zfs_dir.o \
1420     zfs_ioctl.o \
1421     zfs_log.o \
1422     zfs_onexit.o \
1423     zfs_replay.o \
1424     zfs_rlock.o \
1425     zfs_vfsops.o \
1426     zfs_vnops.o \
1427     zvol.o

1429 ZUT_OBJS += \
1430     zut.o

1432 #
1433 #           streams modules
1434 #
1435 BUFMOD_OBJS += bufmod.o

1437 CONNLD_OBJS += connld.o

```

new/usr/src/uts/common/Makefile.files

23

```

1439 DEDUMP_OBJS += dedump.o
1441 DRCOMPAT_OBJS += drcompat.o
1443 LDLINUX_OBJS += ldlinux.o
1445 LDTERM_OBJS += ldterm.o uwidth.o
1447 PCKT_OBJS += pckt.o
1449 PFMOD_OBJS += pfmod.o
1451 PTEM_OBJS += ptem.o
1453 REDIRMOD_OBJS += strredirm.o
1455 TIMOD_OBJS += timod.o
1457 TIRDWR_OBJS += tirdwr.o
1459 TTCOMPAT_OBJS +=ttcompat.o
1461 LOG_OBJS += log.o
1463 PIPEMOD_OBJS += pipemod.o
1465 RPCMOD_OBJS += rpcmod.o clnt_cots.o clnt_clts.o \
1466 clnt_gen.o clnt_perr.o mt_rpcinit.o
1467 rpc_prot.o rpc_sztypes.o rpc_subr.o rpc_calmsg.o \
1468 svc.o svc_clts.o svc_gen.o svc_cots.o \
1469 rpcsys.o xdr_sizeof.o clnt_rdma.o svc_rdma.o \
1470 xdr_rdma.o rdma_subr.o xdrdma_sizeof.o
1472 KLMMOD_OBJS += klmmod.o \
1473 nlm_impl.o \
1474 nlm_rpc_handle.o \
1475 nlm_dispatch.o \
1476 nlm_rpc_svc.o \
1477 nlm_client.o \
1478 nlm_service.o \
1479 nlm_prot_clnt.o \
1480 nlm_prot_xdr.o \
1481 nlm_rpc_clnt.o \
1482 nsm_addr_clnt.o \
1483 nsm_addr_xdr.o \
1484 sm_inter_clnt.o \
1485 sm_inter_xdr.o
1487 KLMOPS_OBJS += klmops.o
1489 TLIMOD_OBJS += tlimod.o t_kalloc.o t_kbind.o t_kclose.o \
1490 t_kconnect.o t_kfree.o t_kgtstate.o
1491 t_krcvdat.o t_ksndudat.o t_kspoll.o t_kunbind.o \
1492 t_kutil.o
1494 RLMOD_OBJS += rlmmod.o
1496 TELMOD_OBJS += telmod.o
1498 CRYPTMOD_OBJS += cryptmod.o
1500 KB_OBJS += kbd.o keytables.o
1502 #
1503 # ID mapping module
1504 #

```

new/usr/src/uts/common/Makefile.files

24

```

1505 IDMAP_OBJS += idmap_mod.o idmap_kapi.o idmap_xdr.o idmap_cache.o
1507 #
1508 # scheduling class modules
1509 #
1510 SDC_OBJS += sysdc.o
1512 RT_OBJS += rt.o
1513 RT_DPTBL_OBJS += rt_dptbl.o
1515 TS_OBJS += ts.o
1516 TS_DPTBL_OBJS += ts_dptbl.o
1518 IA_OBJS += ia.o
1520 FSS_OBJS += fss.o
1522 FX_OBJS += fx.o
1523 FX_DPTBL_OBJS += fx_dptbl.o
1525 #
1526 # Inter-Process Communication (IPC) modules
1527 #
1528 IPC_OBJS += ipc.o
1530 IPCMSG_OBJS += msg.o
1532 IPCSEM_OBJS += sem.o
1534 IPCSHM_OBJS += shm.o
1536 #
1537 # bignum module
1538 #
1539 COMMON_BIGNUM_OBJS += bignum_mod.o bignumimpl.o
1541 BIGNUM_OBJS += $(COMMON_BIGNUM_OBJS) $(BIGNUM_PSR_OBJS)
1543 #
1544 # kernel cryptographic framework
1545 #
1546 KCF_OBJS += kcf.o kcf_callprov.o kcf_cbufcall.o kcf_cipher.o kcf_crypto.o \
1547 kcf_cryptoadm.o kcf_ctxops.o kcf_digest.o kcf_dual.o \
1548 kcf_keys.o kcf_mac.o kcf_mech_tabs.o kcf_miscapi.o \
1549 kcf_object.o kcf_policy.o kcf_prov_lib.o kcf_prov_tabs.o \
1550 kcf_sched.o kcf_session.o kcf_sign.o kcf_spi.o kcf_verify.o \
1551 kcf_random.o modes.o ecb.o cbc.o ctr.o ccm.o gcm.o \
1552 fips_random.o
1554 CRYPTOADM_OBJS += cryptoadm.o
1556 CRYPTO_OBJS += crypto.o
1558 DPROV_OBJS += dprov.o
1560 DCA_OBJS += dca.o dca_3des.o dca_debug.o dca_dsa.o dca_kstat.o dca_rng.o \
1561 dca_rsa.o
1563 AESPROV_OBJS += aes.o aes_impl.o aes_modes.o
1565 ARCFOURPROV_OBJS += arcfour.o arcfour_crypt.o
1567 BLOWFISHPROV_OBJS += blowfish.o blowfish_impl.o
1569 ECCPROV_OBJS += ecc.o ec.o ec2_163.o ec2_mont.o ecdecode.o ecl_mult.o \
1570 ecp_384.o ecp_jac.o ec2_193.o ecl.o ecp_192.o ecp_521.o \

```

new/usr/src/uts/common/Makefile.files

25

```

1571          ecp_jm.o ec2_233.o ecl_curve.o ecp_224.o ecp_aff.o \
1572          ecp_mont.o ec2_aff.o ec_naf.o ecl_gf.o ecp_256.o mp_gf2m.o \
1573          mpi.o mplogic.o mpmontg.o mpprime.o oid.o \
1574          secitem.o ec2_test.o ecp_test.o

1576 RSAPROV_OBJS += rsa.o rsa_impl.o pkcs1.o

1578 SWRANDPROV_OBJS += swrand.o

1580 #
1581 #             kernel SSL
1582 #
1583 KSSL_OBJS += kssl.o ksslioc1.o

1585 KSSL_SOCKETMOD_OBJS += ksslfilter.o ksslapi.o ksslrec.o

1587 #
1588 #             misc. modules
1589 #

1591 C2AUDIT_OBJS += adr.o audit.o audit_event.o audit_io.o \
1592          audit_path.o audit_start.o audit_syscalls.o audit_token.o \
1593          audit_mem.o

1595 PCIC_OBJS += pcic.o

1597 RPCSEC_OBJS += secmod.o          sec_clnt.o          sec_svc.o          sec_gen.o \
1598          auth_des.o             auth_kern.o             auth_none.o             auth_loopb.o \
1599          authdesprt.o           authdesubr.o           authu_prot.o \
1600          key_call.o             key_prot.o             svc_authu.o             svcauthdes.o

1602 RPCSEC_GSS_OBJS += rpcsec_gssmod.o rpcsec_gss.o rpcsec_gss_misc.o \
1603          rpcsec_gss_utils.o svc_rpcsec_gss.o

1605 CONSCONFIG_OBJS += consconfig.o

1607 CONSCONFIG_DACF_OBJS += consconfig_dacf.o consplat.o

1609 TEM_OBJS += tem.o tem_safe.o 6x10.o 7x14.o 12x22.o

1611 KBTRANS_OBJS += \
1612          kbtrans.o \
1613          kbtrans_keytables.o \
1614          kbtrans_polled.o \
1615          kbtrans_streams.o \
1616          usb_keytables.o

1618 KGSSD_OBJS += gssd_clnt_stubs.o gssd_handle.o gssd_prot.o \
1619          gss_display_name.o gss_release_name.o gss_import_name.o \
1620          gss_release_buffer.o gss_release_oid_set.o gen_oids.o gssdmod.o

1622 KGSSD_DERIVED_OBJS = gssd_xdr.o

1624 KGSS_DUMMY_OBJS += dmech.o

1626 KSOCKET_OBJS += ksocket.o ksocket_mod.o

1628 CRYPTO= cksmtypes.o decrypt.o encrypt.o encrypt_length.o etypes.o \
1629          nfold.o verify_checksum.o prng.o block_size.o make_checksum.o \
1630          checksum_length.o hmac.o default_state.o mandatory_sumtype.o

1632 # crypto/des
1633 CRYPTO_DES= fcbc.o fcksum.o fparity.o weak_key.o d3cbc.o efcrypto.o

1635 CRYPTO_DK= checksum.o derive.o dk_decrypt.o dk_encrypt.o

```

new/usr/src/uts/common/Makefile.files

26

```

1637 CRYPTO_ARCFOUR= k5_arcfour.o

1639 # crypto/enc_provider
1640 CRYPTO_ENC= des.o des3.o arcfour_provider.o aes_provider.o

1642 # crypto/hash_provider
1643 CRYPTO_HASH= hash_kef_generic.o hash_kmd5.o hash_crc32.o hash_kshal.o

1645 # crypto/keyhash_provider
1646 CRYPTO_KEYHASH= descbc.o k5_kmd5des.o k_hmac_md5.o

1648 # crypto/crc32
1649 CRYPTO_CRC32= crc32.o

1651 # crypto/old
1652 CRYPTO_OLD= old_decrypt.o old_encrypt.o

1654 # crypto/raw
1655 CRYPTO_RAW= raw_decrypt.o raw_encrypt.o

1657 K5_KRB= kfree.o copy_key.o \
1658          parse.o init_ctx.o \
1659          ser_adata.o ser_addr.o \
1660          ser_auth.o ser_cksum.o \
1661          ser_key.o ser_princ.o \
1662          serialize.o unparse.o \
1663          ser_actx.o

1665 K5_OS= timeofday.o toffset.o \
1666          init_os_ctx.o c_ustime.o

1668 SEAL= seal.o unseal.o

1670 MECH= delete_sec_context.o \
1671          import_sec_context.o \
1672          gssapi_krb5.o \
1673          k5seal.o k5unseal.o k5sealv3.o \
1674          ser_sctx.o \
1675          sign.o \
1676          util_crypt.o \
1677          util_validate.o util_ordering.o \
1678          util_seqnum.o util_set.o util_seed.o \
1679          wrap_size_limit.o verify.o

1683 MECH_GEN= util_token.o

1686 KGSS_KRB5_OBJS += krb5mech.o \
1687          $(MECH) $(SEAL) $(MECH_GEN) \
1688          $(CRYPTO) $(CRYPTO_DES) $(CRYPTO_DK) $(CRYPTO_ARCFOUR) \
1689          $(CRYPTO_ENC) $(CRYPTO_HASH) \
1690          $(CRYPTO_KEYHASH) $(CRYPTO_CRC32) \
1691          $(CRYPTO_OLD) \
1692          $(CRYPTO_RAW) $(K5_KRB) $(K5_OS)

1694 DES_OBJS += des_crypt.o des_impl.o des_ks.o des_soft.o

1696 DLBOOT_OBJS += bootparam_xdr.o nfs_dlinet.o scan.o

1698 KRTLD_OBJS += kobj_bootflags.o getoptstr.o \
1699          kobj.o kobj_kdi.o kobj_lm.o kobj_subr.o

1701 MOD_OBJS += modctl.o modsubr.o modsysfile.o modconf.o modhash.o

```

new/usr/src/uts/common/Makefile.files

27

```

1703 STRPLUMB_OBJS += strplumb.o
1705 CPR_OBJS += cpr_driver.o cpr_dump.o \
1706 cpr_main.o cpr_misc.o cpr_mod.o cpr_stat.o \
1707 cpr_uthread.o
1709 PROF_OBJS += prf.o
1711 SE_OBJS += se_driver.o
1713 SYSACCT_OBJS += acct.o
1715 ACCTCTL_OBJS += acctctl.o
1717 EXACCTSYS_OBJS += exacctsys.o
1719 KAIO_OBJS += aio.o
1721 PCMCIA_OBJS += pcmcia.o cs.o cis.o cis_callout.o cis_handlers.o cis_params.o
1723 BUSRA_OBJS += busra.o
1725 PCS_OBJS += pcs.o
1727 PCAN_OBJS += pcan.o
1729 PCATA_OBJS += pcide.o pcdisk.o pclabel.o pcata.o
1731 PCSER_OBJS += pcser.o pcser_cis.o
1733 PCWL_OBJS += pcwl.o
1735 PSET_OBJS += pset.o
1737 OHCI_OBJS += ohci.o ohci_hub.o ohci_polled.o
1739 UHCI_OBJS += uhci.o uhciutil.o uhcitgt.o uhcihub.o uhcipolled.o
1741 EHCI_OBJS += ehci.o ehci_hub.o ehci_xfer.o ehci_intr.o ehci_util.o ehci_polled.o
1743 HUBD_OBJS += hubd.o
1745 USB_MID_OBJS += usb_mid.o
1747 USB_IA_OBJS += usb_ia.o
1749 UWBA_OBJS += uwba.o uwbai.o
1751 SCSA2USB_OBJS += scsa2usb.o usb_ms_bulkonly.o usb_ms_cbi.o
1753 HWAHC_OBJS += hwahc.o hwahc_util.o
1755 WUSB_DF_OBJS += wusb_df.o
1756 WUSB_FWMOD_OBJS += wusb_fwmod.o
1758 IPF_OBJS += ip_fil_solaris.o fil.o solaris.o ip_state.o ip_frag.o ip_nat.o \
1759 ip_proxy.o ip_auth.o ip_pool.o ip_hhtable.o ip_lookup.o \
1760 ip_log.o misc.o ip_compat.o ip_nat6.o drand48.o
1762 IBD_OBJS += ibd.o ibd_cm.o
1764 EIBNX_OBJS += enx_main.o enx_hdlrs.o enx_ibt.o enx_log.o enx_fip.o \
1765 enx_misc.o enx_q.o enx_ctl.o
1767 EOIB_OBJS += eib_adm.o eib_chan.o eib_cmn.o eib_ctl.o eib_data.o \
1768 eib_fip.o eib_ibt.o eib_log.o eib_mac.o eib_main.o \

```

new/usr/src/uts/common/Makefile.files

28

```

1769 eib_rsrc.o eib_svc.o eib_vnic.o
1771 DLPSTUB_OBJS += dlpistub.o
1773 SDP_OBJS += sdpddi.o
1775 TRILL_OBJS += trill.o
1777 CTF_OBJS += ctf_create.o ctf_decl.o ctf_error.o ctf_hash.o ctf_labels.o \
1778 ctf_lookup.o ctf_open.o ctf_types.o ctf_util.o ctf_subr.o ctf_mod.o
1780 SMBIOS_OBJS += smb_error.o smb_info.o smb_open.o smb_subr.o smb_dev.o
1782 RPCIB_OBJS += rpcib.o
1784 KMDB_OBJS += kdrv.o
1786 AFE_OBJS += afe.o
1788 BGE_OBJS += bge_main2.o bge_chip2.o bge_kstats.o bge_log.o bge_ndd.o \
1789 bge_atomic.o bge_mii.o bge_send.o bge_recv2.o bge_mii_5906.o
1791 DMFE_OBJS += dmfe_log.o dmfe_main.o dmfe_mii.o
1793 EFE_OBJS += efe.o
1795 ELXL_OBJS += elxl.o
1797 HME_OBJS += hme.o
1799 IXGB_OBJS += ixgb.o ixgb_atomic.o ixgb_chip.o ixgb_gld.o ixgb_kstats.o \
1800 ixgb_log.o ixgb_ndd.o ixgb_rx.o ixgb_tx.o ixgb_xmii.o
1802 NGE_OBJS += nge_main.o nge_atomic.o nge_chip.o nge_ndd.o nge_kstats.o \
1803 nge_log.o nge_rx.o nge_tx.o nge_xmii.o
1805 PCN_OBJS += pcn.o
1807 RGE_OBJS += rge_main.o rge_chip.o rge_ndd.o rge_kstats.o rge_log.o rge_rxtx.o
1809 URTW_OBJS += urtw.o
1811 ARN_OBJS += arn_hw.o arn_eeprom.o arn_mac.o arn_calib.o arn_ani.o arn_phy.o arn_
1812 arn_main.o arn_recv.o arn_xmit.o arn_rc.o
1814 ATH_OBJS += ath_aux.o ath_main.o ath_osdep.o ath_rate.o
1816 ATU_OBJS += atu.o
1818 IPW_OBJS += ipw2100_hw.o ipw2100.o
1820 IWI_OBJS += ipw2200_hw.o ipw2200.o
1822 IWH_OBJS += iwh.o
1824 IWK_OBJS += iw2.o
1826 IWP_OBJS += iwp.o
1828 MWL_OBJS += mwl.o
1830 MWLFW_OBJS += mwlfw_mode.o
1832 WPI_OBJS += wpi.o
1834 RAL_OBJS += rt2560.o ral_rate.o

```

```

1836 RUM_OBJS += rum.o
1838 RWD_OBJS += rt2661.o
1840 RWN_OBJS += rt2860.o
1842 UATH_OBJS += uath.o
1844 UATHFW_OBJS += uathfw_mod.o
1846 URAL_OBJS += ural.o
1848 RTW_OBJS += rtw.o smc93cx6.o rtwphy.o rtwphyio.o
1850 ZYD_OBJS += zyd.o zyd_usb.o zyd_hw.o zyd_fw.o
1852 MXFE_OBJS += mxfe.o
1854 MPTSAS_OBJS += mptsas.o mptsas_impl.o mptsas_init.o mptsas_raid.o mptsas_smhba.o
1856 SFE_OBJS += sfe.o sfe_util.o
1858 BFE_OBJS += bfe.o
1860 BRIDGE_OBJS += bridge.o
1862 IDM_SHARED_OBJS += base64.o
1864 IDM_OBJS += $(IDM_SHARED_OBJS) \
1865     idm.o idm_impl.o idm_text.o idm_conn_sm.o idm_so.o
1867 VR_OBJS += vr.o
1869 ATGE_OBJS += atge_main.o atge_lle.o atge_mii.o atge_ll.o atge_llc.o
1871 YGE_OBJS = yge.o
1873 #
1874 #     Build up defines and paths.
1875 #
1876 LINT_DEFS     += -Dunix
1878 #
1879 #     This duality can be removed when the native and target compilers
1880 #     are the same (or at least recognize the same command line syntax!)
1881 #     It is a bug in the current compilation system that the assembler
1882 #     can't process the -Y I, flag.
1883 #
1884 NATIVE_INC_PATH += $(INC_PATH) $(CCYFLAG)$(UTSBASE)/common
1885 AS_INC_PATH     += $(INC_PATH) -I$(UTSBASE)/common
1886 INCLUDE_PATH    += $(INC_PATH) $(CCYFLAG)$(UTSBASE)/common
1888 PCIEB_OBJS += pcieb.o
1890 #     Chelsio N110 10G NIC driver module
1891 #
1892 CH_OBJS = ch.o glue.o pe.o sge.o
1894 CH_COM_OBJS = ch_mac.o ch_subr.o cspi.o espi.o ixfl1010.o mc3.o mc4.o mc5.o \
1895     mv88elxxx.o mv88x201x.o my3126.o pm3393.o tp.o ulp.o \
1896     vsc7321.o vsc7326.o xpak.o
1898 #
1899 #     Chelsio Terminator 4 10G NIC nexus driver module
1900 #

```

```

1901 CXGBE_FW_OBJS =         t4_fw.o t4_cfg.o
1902 CXGBE_COM_OBJS =       t4_hw.o common.o
1903 CXGBE_NEX_OBJS =       t4_nexus.o t4_sge.o t4_mac.o t4_ioctl.o shared.o \
1904     t4_l2t.o adapter.o osdep.o
1906 #
1907 #     Chelsio Terminator 4 10G NIC driver module
1908 #
1909 CXGBE_OBJS =           cxgbe.o
1911 #
1912 #     PCI strings file
1913 #
1914 PCI_STRING_OBJS = pci_strings.o
1916 NET_DACF_OBJS += net_dacf.o
1918 #
1919 #     Xframe 10G NIC driver module
1920 #
1921 XGE_OBJS =             xge.o xgell.o
1923 XGE_HAL_OBJS =        xgehal-channel.o xgehal-fifo.o xgehal-ring.o xgehal-config.o \
1924     xgehal-driver.o xgehal-mm.o xgehal-stats.o xgehal-device.o \
1925     xge-queue.o xgehal-mgmt.o xgehal-mgmtaux.o
1927 #
1928 #     e1000g module
1929 #
1930 E1000G_OBJS +=        e1000_80003es2lan.o e1000_82540.o e1000_82541.o e1000_82542.o \
1931     e1000_82543.o e1000_82571.o e1000_api.o e1000_ich8lan.o \
1932     e1000_mac.o e1000_manage.o e1000_nvmm.o e1000_osdep.o \
1933     e1000_phy.o e1000g_debug.o e1000g_main.o e1000g_alloc.o \
1934     e1000g_tx.o e1000g_rx.o e1000g_stat.o
1936 #
1937 #     Intel 82575 1G NIC driver module
1938 #
1939 IGB_OBJS =            igb_82575.o igb_api.o igb_mac.o igb_manage.o \
1940     igb_nvmm.o igb_osdep.o igb_phy.o igb_buf.o \
1941     igb_debug.o igb_gld.o igb_log.o igb_main.o \
1942     igb_rx.o igb_stat.o igb_tx.o
1944 #
1945 #     Intel Pro/100 NIC driver module
1946 #
1947 IPRB_OBJS =          iprb.o
1949 #
1950 #     Intel 10GbE PCIE NIC driver module
1951 #
1952 IXGBE_OBJS =         ixgbe_82598.o ixgbe_82599.o ixgbe_api.o \
1953     ixgbe_common.o ixgbe_phy.o \
1954     ixgbe_buf.o ixgbe_debug.o ixgbe_gld.o \
1955     ixgbe_log.o ixgbe_main.o \
1956     ixgbe_osdep.o ixgbe_rx.o ixgbe_stat.o \
1957     ixgbe_tx.o ixgbe_x540.o ixgbe_mbx.o
1959 #
1960 #     NIU 10G/1G driver module
1961 #
1962 NXGE_OBJS =          nxge_mac.o nxge_ipp.o nxge_rxdma.o \
1963     nxge_txdma.o nxge_txc.o nxge_main.o \
1964     nxge_hw.o nxge_fzc.o nxge_virtual.o \
1965     nxge_send.o nxge_classify.o nxge_fflp.o \
1966     nxge_fflp_hash.o nxge_ndd.o nxge_kstats.o

```

new/usr/src/uts/common/Makefile.files

31

```

1967          nxge_zcp.o nxge_fm.o nxge_espc.o nxge_hv.o      \
1968          nxge_hio.o nxge_hio_guest.o nxge_intr.o

1970 NXGE_NPI_OBJS = \
1971          npi.o npi_mac.o npi_ipp.o                        \
1972          npi_txdma.o npi_rxdma.o npi_txc.o              \
1973          npi_zcp.o npi_espc.o npi_fflp.o                \
1974          npi_vir.o

1976 NXGE_HCALL_OBJS = \
1977          nxge_hcall.o

1979 #
1980 # Virtio modules
1981 #

1983 # Virtio core
1984 VIRTIO_OBJS = virtio.o

1986 # Virtio block driver
1987 VIOBLK_OBJS = vioblk.o

1989 #
1990 #          kiconv modules
1991 #
1992 KICONV_EMEA_OBJS += kiconv_emea.o

1994 KICONV_JA_OBJS += kiconv_ja.o

1996 KICONV_KO_OBJS += kiconv_cck_common.o kiconv_ko.o

1998 KICONV_SC_OBJS += kiconv_cck_common.o kiconv_sc.o

2000 KICONV_TC_OBJS += kiconv_cck_common.o kiconv_tc.o

2002 #
2003 #          AAC module
2004 #
2005 AAC_OBJS = aac.o aac_ioctl.o

2007 #
2008 #          sdc card modules
2009 #
2010 SDA_OBJS =          sda_cmd.o sda_host.o sda_init.o sda_mem.o sda_mod.o sda_slot.o
2011 SDHOST_OBJS =      sdhost.o

2013 #
2014 #          hxge 10G driver module
2015 #
2016 HXGE_OBJS =          hxge_main.o hxge_vmac.o hxge_send.o      \
2017          hxge_txdma.o hxge_rxdma.o hxge_virtual.o           \
2018          hxge_fm.o hxge_fzc.o hxge_hw.o hxge_kstats.o       \
2019          hxge_ndd.o hxge_pfc.o                               \
2020          hpi.o hpi_vmac.o hpi_rxdma.o hpi_txdma.o           \
2021          hpi_vir.o hpi_pfc.o

2023 #
2024 #          MEGARAID_SAS module
2025 #
2026 MEGA_SAS_OBJS = megaraid_sas.o

2028 #
2029 #          MR_SAS module
2030 #
2031 MR_SAS_OBJS = ld_pd_map.o mr_sas.o mr_sas_tbolt.o mr_sas_list.o

```

new/usr/src/uts/common/Makefile.files

32

```

2033 #
2034 #          ISCSI_INITIATOR module
2035 #
2036 ISCSI_INITIATOR_OBJS = chap.o iscsi_io.o iscsi_thread.o      \
2037          iscsi_ioctl.o iscsid.o iscsi.o                      \
2038          iscsi_login.o isns_client.o iscsiAuthClient.o      \
2039          iscsi_lun.o iscsiAuthClientGlue.o                   \
2040          iscsi_net.o nvfile.o iscsi_cmd.o                     \
2041          iscsi_queue.o persistent.o iscsi_conn.o             \
2042          iscsi_sess.o radius_auth.o iscsi_crc.o              \
2043          iscsi_stats.o radius_packet.o iscsi_doorclt.o       \
2044          iscsi_targetparam.o utils.o kifconf.o

2046 #
2047 #          ntxn 10Gb/1Gb NIC driver module
2048 #
2049 NTXN_OBJS =          unm_nic_init.o unm_gem.o unm_nic_hw.o unm_ndd.o \
2050          unm_nic_main.o unm_nic_isr.o unm_nic_ctx.o niu.o

2052 #
2053 #          Myricom 10Gb NIC driver module
2054 #
2055 MYRI10GE_OBJS = myril0ge.o myril0ge_lro.o

2057 #          nulldriver module
2058 #
2059 NULLDRIVER_OBJS =          nulldriver.o

2061 TPM_OBJS =          tpm.o tpm_hcall.o

```

new/usr/src/uts/common/fs/zfs/dnode.c

1

```
*****
56685 Tue Sep  3 20:26:58 2013
new/usr/src/uts/common/fs/zfs/dnode.c
4101 metaslab_debug should allow for fine-grained control
4102 space_maps should store more information about themselves
4103 space_map object blocksize should be increased
4104 ::spa_space no longer works
4105 removing a mirrored log device results in a leaked object
4106 asynchronously load metaslab
Reviewed by: Matthew Ahrens <mahrens@delphix.com>
Reviewed by: Adam Leventhal <ahl@delphix.com>
Reviewed by: Sebastien Roy <seb@delphix.com>
*****
_____unchanged_portion_omitted_____

1311 /*
1312  * Try to change the block size for the indicated dnode. This can only
1313  * succeed if there are no blocks allocated or dirty beyond first block
1314  */
1315 int
1316 dnode_set_blkisz(dnode_t *dn, uint64_t size, int ibs, dmu_tx_t *tx)
1317 {
1318     dmu_buf_impl_t *db, *db_next;
1319     int err;

1321     if (size == 0)
1322         size = SPA_MINBLOCKSIZE;
1323     if (size > SPA_MAXBLOCKSIZE)
1324         size = SPA_MAXBLOCKSIZE;
1325     else
1326         size = P2ROUNDUP(size, SPA_MINBLOCKSIZE);

1328     if (ibs == dn->dn_indblkshift)
1329         ibs = 0;

1331     if (size >> SPA_MINBLOCKSHIFT == dn->dn_datablkiszsec && ibs == 0)
1332         return (0);

1334     rw_enter(&dn->dn_struct_rwlock, RW_WRITER);

1336     /* Check for any allocated blocks beyond the first */
1337     if (dn->dn_maxblkid != 0)
1338     if (dn->dn_phys->dn_maxblkid != 0)
1339         goto fail;

1340     mutex_enter(&dn->dn_dbufs_mtx);
1341     for (db = list_head(&dn->dn_dbufs); db; db = db_next) {
1342         db_next = list_next(&dn->dn_dbufs, db);

1344         if (db->db_blkid != 0 && db->db_blkid != DMU_BONUS_BLKID &&
1345             db->db_blkid != DMU_SPILL_BLKID) {
1346             mutex_exit(&dn->dn_dbufs_mtx);
1347             goto fail;
1348         }
1349     }
1350     mutex_exit(&dn->dn_dbufs_mtx);

1352     if (ibs && dn->dn_nlevels != 1)
1353         goto fail;

1355     /* resize the old block */
1356     err = dbuf_hold_impl(dn, 0, 0, TRUE, FTAG, &db);
1357     if (err == 0)
1358         dbuf_new_size(db, size, tx);
1359     else if (err != ENOENT)
1360         goto fail;
```

new/usr/src/uts/common/fs/zfs/dnode.c

2

```
1362     dnode_setdblksz(dn, size);
1363     dnode_setdirty(dn, tx);
1364     dn->dn_next_blkisz[tx->tx_txg&TXG_MASK] = size;
1365     if (ibs) {
1366         dn->dn_indblkshift = ibs;
1367         dn->dn_next_indblkshift[tx->tx_txg&TXG_MASK] = ibs;
1368     }
1369     /* rele after we have fixed the blocksize in the dnode */
1370     if (db)
1371         dbuf_rele(db, FTAG);

1373     rw_exit(&dn->dn_struct_rwlock);
1374     return (0);

1376 fail:
1377     rw_exit(&dn->dn_struct_rwlock);
1378     return (SET_ERROR(ENOTSUP));
1379 }
_____unchanged_portion_omitted_____
```

new/usr/src/uts/common/fs/zfs/metablab.c

1

```
*****
62230 Tue Sep 3 20:26:59 2013
new/usr/src/uts/common/fs/zfs/metablab.c
4101 metablab_debug should allow for fine-grained control
4102 space_maps should store more information about themselves
4103 space_map object blocksize should be increased
4104 ::spa_space no longer works
4105 removing a mirrored log device results in a leaked object
4106 asynchronously load metablab
Reviewed by: Matthew Ahrens <mahrens@delphix.com>
Reviewed by: Adam Leventhal <ahl@delphix.com>
Reviewed by: Sebastien Roy <seb@delphix.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[ ]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22 * Copyright (c) 2005, 2010, Oracle and/or its affiliates. All rights reserved.
23 * Copyright (c) 2013 by Delphix. All rights reserved.
24 * Copyright (c) 2013 by Saso Kiselkov. All rights reserved.
25 */

27 #include <sys/zfs_context.h>
28 #include <sys/dmu.h>
29 #include <sys/dmu_tx.h>
30 #include <sys/space_map.h>
31 #include <sys/metablab_impl.h>
32 #include <sys/vdev_impl.h>
33 #include <sys/zio.h>
34 #include <sys/spa_impl.h>

36 /*
37  * Allow allocations to switch to gang blocks quickly. We do this to
38  * avoid having to load lots of space_maps in a given txg. There are,
39  * however, some cases where we want to avoid "fast" ganging and instead
40  * we want to do an exhaustive search of all metablabs on this device.
41  * Currently we don't allow any gang, zil, or dump device related allocations
42  * to "fast" gang.
43  */
44 #define CAN_FASTGANG(flags) \
45     (!(flags) & (METASLAB_GANG_CHILD | METASLAB_GANG_HEADER | \
46     METASLAB_GANG_AVOID))

48 #define METASLAB_WEIGHT_PRIMARY      (1ULL << 63)
49 #define METASLAB_WEIGHT_SECONDARY   (1ULL << 62)
50 #define METASLAB_ACTIVE_MASK        \
51     (METASLAB_WEIGHT_PRIMARY | METASLAB_WEIGHT_SECONDARY)

53 uint64_t metablab_aliquot = 512ULL << 10;
```

new/usr/src/uts/common/fs/zfs/metablab.c

2

```
54 uint64_t metablab_gang_bang = SPA_MAXBLOCKSIZE + 1; /* force gang blocks */

56 /*
57  * The in-core space map representation is more compact than its on-disk form.
58  * The zfs_condense_pct determines how much more compact the in-core
59  * space_map representation must be before we compact it on-disk.
60  * Values should be greater than or equal to 100.
61  */
62 int zfs_condense_pct = 200;

64 /*
65  * This value defines the number of allowed allocation failures per vdev.
66  * If a device reaches this threshold in a given txg then we consider skipping
67  * allocations on that device. The value of zfs_mg_alloc_failures is computed
68  * in zio_init() unless it has been overridden in /etc/system.
69  */
70 int zfs_mg_alloc_failures = 0;

72 /*
73  * The zfs_mg_noalloc_threshold defines which metablab groups should
74  * be eligible for allocation. The value is defined as a percentage of
75  * a free space. Metablab groups that have more free space than
76  * zfs_mg_noalloc_threshold are always eligible for allocations. Once
77  * a metablab group's free space is less than or equal to the
78  * zfs_mg_noalloc_threshold the allocator will avoid allocating to that
79  * group unless all groups in the pool have reached zfs_mg_noalloc_threshold.
80  * Once all groups in the pool reach zfs_mg_noalloc_threshold then all
81  * groups are allowed to accept allocations. Gang blocks are always
82  * eligible to allocate on any metablab group. The default value of 0 means
83  * no metablab group will be excluded based on this criterion.
84  */
85 int zfs_mg_noalloc_threshold = 0;

87 /*
88  * When set will load all metablabs when pool is first opened.
89  */
90 int metablab_debug_load = 0;
91 static int metablab_debug = 0;

92 /*
93  * When set will prevent metablabs from being unloaded.
94  */
95 int metablab_debug_unload = 0;

97 /*
98  * Minimum size which forces the dynamic allocator to change
99  * it's allocation strategy. Once the space map cannot satisfy
100 * an allocation of this size then it switches to using more
101 * aggressive strategy (i.e search by size rather than offset).
102 */
103 uint64_t metablab_df_alloc_threshold = SPA_MAXBLOCKSIZE;

105 /*
106  * The minimum free space, in percent, which must be available
107  * in a space map to continue allocations in a first-fit fashion.
108  * Once the space_map's free space drops below this level we dynamically
109  * switch to using best-fit allocations.
110  */
111 int metablab_df_free_pct = 4;

113 /*
114  * A metablab is considered "free" if it contains a contiguous
115  * segment which is greater than metablab_min_alloc_size.
116  */
117 uint64_t metablab_min_alloc_size = DMU_MAX_ACCESS;
```



```

119 /*
120 * Percentage of all cpus that can be used by the metablab taskq.
109 * Max number of space_maps to prefetch.
121 */
122 int metablab_load_pct = 50;
111 int metablab_prefetch_limit = SPA_DVAS_PER_BP;

124 /*
125 * Determines how many txgs a metablab may remain loaded without having any
126 * allocations from it. As long as a metablab continues to be used we will
127 * keep it loaded.
114 * Percentage bonus multiplier for metablabs that are in the bonus area.
128 */
129 int metablab_unload_delay = TXG_SIZE * 2;
116 int metablab_smo_bonus_pct = 150;

131 /*
132 * Should we be willing to write data to degraded vdevs?
133 */
134 boolean_t zfs_write_to_degraded = B_FALSE;

136 /*
137 * Max number of metablabs per group to preload.
138 */
139 int metablab_preload_limit = SPA_DVAS_PER_BP;

141 /*
142 * Enable/disable preloading of metablab.
143 */
144 boolean_t metablab_preload_enabled = B_TRUE;

146 /*
147 * Enable/disable additional weight factor for each metablab.
148 */
149 boolean_t metablab_weight_factor_enable = B_FALSE;

152 /*
153 * =====
154 * Metablab classes
155 * =====
156 */
157 metablab_class_t *
158 metablab_class_create(spa_t *spa, metablab_ops_t *ops)
129 metablab_class_create(spa_t *spa, space_map_ops_t *ops)
159 {
160     metablab_class_t *mc;

162     mc = kmem_zalloc(sizeof (metablab_class_t), KM_SLEEP);

164     mc->mc_spa = spa;
165     mc->mc_rotor = NULL;
166     mc->mc_ops = ops;

168     return (mc);
169 }
    unchanged_portion_omitted_

243 /*
244 * =====
245 * Metablab groups
246 * =====
247 */
248 static int
249 metablab_compare(const void *x1, const void *x2)

```

```

250 {
251     const metablab_t *m1 = x1;
252     const metablab_t *m2 = x2;

254     if (m1->ms_weight < m2->ms_weight)
255         return (1);
256     if (m1->ms_weight > m2->ms_weight)
257         return (-1);

259     /*
260      * If the weights are identical, use the offset to force uniqueness.
261      */
262     if (m1->ms_start < m2->ms_start)
233     if (m1->ms_map->sm_start < m2->ms_map->sm_start)
263         return (-1);
264     if (m1->ms_start > m2->ms_start)
235     if (m1->ms_map->sm_start > m2->ms_map->sm_start)
265         return (1);

267     ASSERT3P(m1, ==, m2);

269     return (0);
270 }
    unchanged_portion_omitted_

319 metablab_group_t *
320 metablab_group_create(metablab_class_t *mc, vdev_t *vd)
321 {
322     metablab_group_t *mg;

324     mg = kmem_zalloc(sizeof (metablab_group_t), KM_SLEEP);
325     mutex_init(&mg->mg_lock, NULL, MUTEX_DEFAULT, NULL);
326     avl_create(&mg->mg_metablab_tree, metablab_compare,
327         sizeof (metablab_t), offsetof(struct metablab, ms_group_node));
328     mg->mg_vd = vd;
329     mg->mg_class = mc;
330     mg->mg_activation_count = 0;

332     mg->mg_taskq = taskq_create("metablab_group_tasksq", metablab_load_pct,
333         minclsyspri, 10, INT_MAX, TASKQ_THREADS_CPU_PCT);

335     return (mg);
336 }
    unchanged_portion_omitted_

387 void
388 metablab_group_passivate(metablab_group_t *mg)
389 {
390     metablab_class_t *mc = mg->mg_class;
391     metablab_group_t *mgprev, *mgnext;

393     ASSERT(spa_config_held(mc->mc_spa, SCL_ALLOC, RW_WRITER));

395     if (--mg->mg_activation_count != 0) {
396         ASSERT(mc->mc_rotor != mg);
397         ASSERT(mg->mg_prev == NULL);
398         ASSERT(mg->mg_next == NULL);
399         ASSERT(mg->mg_activation_count < 0);
400         return;
401     }

403     taskq_wait(mg->mg_taskq);

405     mgprev = mg->mg_prev;
406     mgnext = mg->mg_next;

```

```

408     if (mg == mgnext) {
409         mc->mc_rotor = NULL;
410     } else {
411         mc->mc_rotor = mgnext;
412         mgprev->mg_next = mgnext;
413         mgnext->mg_prev = mgprev;
414     }

416     mg->mg_prev = NULL;
417     mg->mg_next = NULL;
418 }
    unchanged_portion_omitted

482 /*
483 * =====
484 * Range tree callbacks
485 * Common allocator routines
486 * =====
487 */

488 /*
489 * Comparison function for the private size-ordered tree. Tree is sorted
490 * by size, larger sizes at the end of the tree.
491 */
492 static int
493 metablab_rangesize_compare(const void *x1, const void *x2)
494 metablab_segsize_compare(const void *x1, const void *x2)
495 {
496     const range_seg_t *r1 = x1;
497     const range_seg_t *r2 = x2;
498     uint64_t rs_size1 = r1->rs_end - r1->rs_start;
499     uint64_t rs_size2 = r2->rs_end - r2->rs_start;
500     const space_seg_t *s1 = x1;
501     const space_seg_t *s2 = x2;
502     uint64_t ss_size1 = s1->ss_end - s1->ss_start;
503     uint64_t ss_size2 = s2->ss_end - s2->ss_start;

504     if (rs_size1 < rs_size2)
505     if (ss_size1 < ss_size2)
506         return (-1);
507     if (rs_size1 > rs_size2)
508     if (ss_size1 > ss_size2)
509         return (1);

510     return (0);
511 }
512 }

514 /*
515 * Create any block allocator specific components. The current allocators
516 * rely on using both a size-ordered range_tree_t and an array of uint64_t's.
517 * This is a helper function that can be used by the allocator to find
518 * a suitable block to allocate. This will search the specified AVL
519 * tree looking for a block that matches the specified criteria.
520 */
521 static void
522 metablab_rt_create(range_tree_t *rt, void *arg)
523 static uint64_t
524 metablab_block_picker(avl_tree_t *t, uint64_t *cursor, uint64_t size,

```

```

481     uint64_t align)
520 {
521     metablab_t *msp = arg;
522     space_seg_t *ss, ssearch;
523     avl_index_t where;

524     ASSERT3P(rt->rt_arg, ==, msp);
525     ASSERT(msp->ms_tree == NULL);
526     ssearch.ss_start = *cursor;
527     ssearch.ss_end = *cursor + size;

528     avl_create(&msp->ms_size_tree, metablab_rangesize_compare,
529             sizeof (range_seg_t), offsetof(range_seg_t, rs_pp_node));
530 }
531     ss = avl_find(t, &ssearch, &where);
532     if (ss == NULL)
533         ss = avl_nearest(t, where, AVL_AFTER);

534 /*
535 * Destroy the block allocator specific components.
536 while (ss != NULL) {
537     uint64_t offset = P2ROUNDUP(ss->ss_start, align);

538     if (offset + size <= ss->ss_end) {
539         *cursor = offset + size;
540         return (offset);
541     }
542     ss = AVL_NEXT(t, ss);
543 }

544 /*
545 * If we know we've searched the whole map (*cursor == 0), give up.
546 * Otherwise, reset the cursor to the beginning and try again.
547 */
548 if (*cursor == 0)
549     return (-1ULL);

550 *cursor = 0;
551 return (metabolab_block_picker(t, cursor, size, align));
552 }

553 static void
554 metablab_rt_destroy(range_tree_t *rt, void *arg)
555 metablab_pp_load(space_map_t *sm)
556 {
557     metablab_t *msp = arg;
558     space_seg_t *ss;

559     ASSERT3P(rt->rt_arg, ==, msp);
560     ASSERT3P(msp->ms_tree, ==, rt);
561     ASSERT0(avl_numnodes(&msp->ms_size_tree));
562     ASSERT(sm->sm_ppd == NULL);
563     sm->sm_ppd = kmem_zalloc(64 * sizeof (uint64_t), KM_SLEEP);

564     avl_destroy(&msp->ms_size_tree);
565     sm->sm_pp_root = kmem_alloc(sizeof (avl_tree_t), KM_SLEEP);
566     avl_create(sm->sm_pp_root, metablab_segsize_compare,
567             sizeof (space_seg_t), offsetof(struct space_seg, ss_pp_node));

568     for (ss = avl_first(&sm->sm_root); ss; ss = AVL_NEXT(&sm->sm_root, ss))
569         avl_add(sm->sm_pp_root, ss);
570 }

571 static void
572 metablab_rt_add(range_tree_t *rt, range_seg_t *rs, void *arg)
573 metablab_pp_unload(space_map_t *sm)

```

```

547 {
548     metaslab_t *msp = arg;
549     void *cookie = NULL;

550     ASSERT3P(rt->rt_arg, ==, msp);
551     ASSERT3P(msp->ms_tree, ==, rt);
552     VERIFY(!msp->ms_condensing);
553     avl_add(&msp->ms_size_tree, rs);
554     kmem_free(sm->sm_ppd, 64 * sizeof (uint64_t));
555     sm->sm_ppd = NULL;

556     while (avl_destroy_nodes(sm->sm_pp_root, &cookie) != NULL) {
557         /* tear down the tree */
558     }

559     avl_destroy(sm->sm_pp_root);
560     kmem_free(sm->sm_pp_root, sizeof (avl_tree_t));
561     sm->sm_pp_root = NULL;
562 }

563 /* ARGSUSED */
564 static void
565 metaslab_rt_remove(range_tree_t *rt, range_seg_t *rs, void *arg)
566 metaslab_pp_claim(space_map_t *sm, uint64_t start, uint64_t size)
567 {
568     metaslab_t *msp = arg;

569     ASSERT3P(rt->rt_arg, ==, msp);
570     ASSERT3P(msp->ms_tree, ==, rt);
571     VERIFY(!msp->ms_condensing);
572     avl_remove(&msp->ms_size_tree, rs);
573     /* No need to update cursor */
574 }

575 /* ARGSUSED */
576 static void
577 metaslab_rt_vacate(range_tree_t *rt, void *arg)
578 metaslab_pp_free(space_map_t *sm, uint64_t start, uint64_t size)
579 {
580     metaslab_t *msp = arg;

581     ASSERT3P(rt->rt_arg, ==, msp);
582     ASSERT3P(msp->ms_tree, ==, rt);

583     /*
584      * Normally one would walk the tree freeing nodes along the way.
585      * Since the nodes are shared with the range trees we can avoid
586      * walking all nodes and just reinitialize the avl tree. The nodes
587      * will be freed by the range tree, so we don't want to free them here.
588      */
589     avl_create(&msp->ms_size_tree, metaslab_rangesize_compare,
590             sizeof (range_seg_t), offsetof(range_seg_t, rs_pp_node));
591     /* No need to update cursor */
592 }

593 static range_tree_ops_t metaslab_rt_ops = {
594     metaslab_rt_create,
595     metaslab_rt_destroy,
596     metaslab_rt_add,
597     metaslab_rt_remove,
598     metaslab_rt_vacate
599 };

600 /*
601  * =====
602  * Metaslab block operations

```

```

603 * =====
604 */

605 /*
606  * Return the maximum contiguous segment within the metaslab.
607  */
608 uint64_t
609 metaslab_block_maxsize(metaslab_t *msp)
610 metaslab_pp_maxsize(space_map_t *sm)
611 {
612     avl_tree_t *t = &msp->ms_size_tree;
613     range_seg_t *rs;
614     avl_tree_t *t = sm->sm_pp_root;
615     space_seg_t *ss;

616     if (t == NULL || (rs = avl_last(t)) == NULL)
617     if (t == NULL || (ss = avl_last(t)) == NULL)
618         return (0ULL);

619     return (rs->rs_end - rs->rs_start);
620     return (ss->ss_end - ss->ss_start);
621 }

622 uint64_t
623 metaslab_block_alloc(metaslab_t *msp, uint64_t size)
624 {
625     uint64_t start;
626     range_tree_t *rt = msp->ms_tree;

627     VERIFY(!msp->ms_condensing);

628     start = msp->ms_ops->msop_alloc(msp, size);
629     if (start != -1ULL) {
630         vdev_t *vd = msp->ms_group->mg_vd;

631         VERIFY0(P2PHASE(start, 1ULL << vd->vdev_ashift));
632         VERIFY0(P2PHASE(size, 1ULL << vd->vdev_ashift));
633         VERIFY3U(range_tree_space(rt) - size, <=, msp->ms_size);
634         range_tree_remove(rt, start, size);
635     }
636     return (start);
637 }

638 /*
639  * =====
640  * Common allocator routines
641  * =====
642  */

643 /*
644  * This is a helper function that can be used by the allocator to find
645  * a suitable block to allocate. This will search the specified AVL
646  * tree looking for a block that matches the specified criteria.
647  */
648 static uint64_t
649 metaslab_block_picker(avl_tree_t *t, uint64_t *cursor, uint64_t size,
650                     uint64_t align)
651 {
652     range_seg_t *rs, rsearch;
653     avl_index_t where;

654     rsearch.rs_start = *cursor;
655     rsearch.rs_end = *cursor + size;

656     rs = avl_find(t, &rsearch, &where);
657     if (rs == NULL)

```

```

657     rs = avl_nearest(t, where, AVL_AFTER);
659     while (rs != NULL) {
660         uint64_t offset = P2ROUNDUP(rs->rs_start, align);
662         if (offset + size <= rs->rs_end) {
663             *cursor = offset + size;
664             return (offset);
665         }
666         rs = AVL_NEXT(t, rs);
667     }
669     /*
670     * If we know we've searched the whole map (*cursor == 0), give up.
671     * Otherwise, reset the cursor to the beginning and try again.
672     */
673     if (*cursor == 0)
674         return (-1ULL);
676     *cursor = 0;
677     return (metabolab_block_picker(t, cursor, size, align));
678 }
680 /*
681 * =====
682 * The first-fit block allocator
683 * =====
684 */
685 static uint64_t
686 metabolab_ff_alloc(metabolab_t *msp, uint64_t size)
687 {
688     /*
689     * Find the largest power of 2 block size that evenly divides the
690     * requested size. This is used to try to allocate blocks with similar
691     * alignment from the same area of the metabolab (i.e. same cursor
692     * bucket) but it does not guarantee that other allocations sizes
693     * may exist in the same region.
694     */
695     avl_tree_t *t = &msp->sm_root;
696     uint64_t align = size & ~size;
697     uint64_t *cursor = &msp->ms_lbas[highbit(align) - 1];
698     avl_tree_t *rt = &msp->ms_tree->rt_root;
699     uint64_t *cursor = (uint64_t *)msp->sm_ppd + highbit(align) - 1;
700 }
702 /* ARGSUSED */
703 static boolean_t
704 metabolab_ff_fragmented(metabolab_t *msp)
705 {
706     return (B_TRUE);
707 }
709 static metabolab_ops_t metabolab_ff_ops = {
710     static space_map_ops_t metabolab_ff_ops = {
711         metabolab_pp_load,
712         metabolab_pp_unload,
713         metabolab_ff_alloc,
714         metabolab_pp_claim,
715         metabolab_pp_free,
716         metabolab_pp_maxsize,
717         metabolab_ff_fragmented

```

```

712 };
714 /*
715 * =====
716 * Dynamic block allocator -
717 * Uses the first fit allocation scheme until space get low and then
718 * adjusts to a best fit allocation method. Uses metabolab_df_alloc_threshold
719 * and metabolab_df_free_pct to determine when to switch the allocation scheme.
720 * =====
721 */
722 static uint64_t
723 metabolab_df_alloc(metabolab_t *msp, uint64_t size)
724 {
725     /*
726     * Find the largest power of 2 block size that evenly divides the
727     * requested size. This is used to try to allocate blocks with similar
728     * alignment from the same area of the metabolab (i.e. same cursor
729     * bucket) but it does not guarantee that other allocations sizes
730     * may exist in the same region.
731     */
732     avl_tree_t *t = &msp->sm_root;
733     uint64_t align = size & ~size;
734     uint64_t *cursor = &msp->ms_lbas[highbit(align) - 1];
735     range_tree_t *rt = msp->ms_tree;
736     avl_tree_t *rt = &rt->rt_root;
737     uint64_t max_size = metabolab_block_maxsize(msp);
738     int free_pct = range_tree_space(rt) * 100 / msp->ms_size;
739     uint64_t *cursor = (uint64_t *)msp->sm_ppd + highbit(align) - 1;
740     uint64_t max_size = metabolab_pp_maxsize(msp);
741     int free_pct = msp->sm_space * 100 / msp->ms_size;
742     ASSERT(MUTEX_HELD(&msp->ms_lock));
743     ASSERT3U(avl_numnodes(t), ==, avl_numnodes(&msp->ms_size_tree));
744     ASSERT(MUTEX_HELD(msp->sm_lock));
745     ASSERT3U(avl_numnodes(&msp->sm_root), ==, avl_numnodes(msp->sm_pp_root));
747     if (max_size < size)
748         return (-1ULL);
750     /*
751     * If we're running low on space switch to using the size
752     * sorted AVL tree (best-fit).
753     */
754     if (max_size < metabolab_df_alloc_threshold ||
755         free_pct < metabolab_df_free_pct) {
756         t = &msp->ms_size_tree;
757         t = msp->sm_pp_root;
758         *cursor = 0;
759     }
760     return (metabolab_block_picker(t, cursor, size, 1ULL));
761 }
763 static boolean_t
764 metabolab_df_fragmented(metabolab_t *msp)
765 {
766     range_tree_t *rt = msp->ms_tree;
767     uint64_t max_size = metabolab_block_maxsize(msp);
768     int free_pct = range_tree_space(rt) * 100 / msp->ms_size;
769     uint64_t max_size = metabolab_pp_maxsize(msp);
770     int free_pct = msp->sm_space * 100 / msp->ms_size;
772     if (max_size >= metabolab_df_alloc_threshold &&
773         free_pct >= metabolab_df_free_pct)

```

```

767         return (B_FALSE);
769     return (B_TRUE);
770 }

772 static metabolab_ops_t metabolab_df_ops = {
657 static space_map_ops_t metabolab_df_ops = {
658     metabolab_pp_load,
659     metabolab_pp_unload,
773     metabolab_df_alloc,
661     metabolab_pp_claim,
662     metabolab_pp_free,
663     metabolab_pp_maxsize,
774     metabolab_df_fragmented
775 };

777 /*
778 * =====
779 * Cursor fit block allocator -
780 * Select the largest region in the metabolab, set the cursor to the beginning
781 * of the range and the cursor_end to the end of the range. As allocations
782 * are made advance the cursor. Continue allocating from the cursor until
783 * the range is exhausted and then find a new range.
669 * Other experimental allocators
784 * =====
785 */
786 static uint64_t
787 metabolab_cf_alloc(metabolab_t *msp, uint64_t size)
673 metabolab_cdf_alloc(space_map_t *sm, uint64_t size)
788 {
789     range_tree_t *rt = msp->ms_tree;
790     avl_tree_t *t = &msp->ms_size_tree;
791     uint64_t *cursor = &msp->ms_lbas[0];
792     uint64_t *cursor_end = &msp->ms_lbas[1];
675     avl_tree_t *t = &sm->sm_root;
676     uint64_t *cursor = (uint64_t *)sm->sm_ppd;
677     uint64_t *extent_end = (uint64_t *)sm->sm_ppd + 1;
678     uint64_t max_size = metabolab_pp_maxsize(sm);
679     uint64_t rsize = size;
793     uint64_t offset = 0;

795     ASSERT(MUTEX_HELD(&msp->ms_lock));
796     ASSERT3U(avl_numnodes(t), ==, avl_numnodes(&rt->rt_root));
682     ASSERT(MUTEX_HELD(sm->sm_lock));
683     ASSERT3U(avl_numnodes(&sm->sm_root), ==, avl_numnodes(sm->sm_pp_root));

798     ASSERT3U(*cursor_end, >=, *cursor);

800     if ((*cursor + size) > *cursor_end) {
801         range_seg_t *rs;

803         rs = avl_last(&msp->ms_size_tree);
804         if (rs == NULL || (rs->rs_end - rs->rs_start) < size)
685             if (max_size < size)
805                 return (-1ULL);

807         *cursor = rs->rs_start;
808         *cursor_end = rs->rs_end;
809     }
688     ASSERT3U(*extent_end, >=, *cursor);

811     offset = *cursor;
812     *cursor += size;
690 /*
691  * If we're running low on space switch to using the size
692  * sorted AVL tree (best-fit).

```

```

693     */
694     if ((*cursor + size) > *extent_end) {

696         t = sm->sm_pp_root;
697         *cursor = *extent_end = 0;

699         if (max_size > 2 * SPA_MAXBLOCKSIZE)
700             rsize = MIN(metabolab_min_alloc_size, max_size);
701         offset = metabolab_block_picker(t, extent_end, rsize, 1ULL);
702         if (offset != -1)
703             *cursor = offset + size;
704     } else {
705         offset = metabolab_block_picker(t, cursor, rsize, 1ULL);
706     }
707     ASSERT3U(*cursor, <=, *extent_end);
814     return (offset);
815 }

817 static boolean_t
818 metabolab_cf_fragmented(metabolab_t *msp)
712 metabolab_cdf_fragmented(space_map_t *sm)
819 {
820     return (metabolab_block_maxsize(msp) < metabolab_min_alloc_size);
714     uint64_t max_size = metabolab_pp_maxsize(sm);

716     if (max_size > (metabolab_min_alloc_size * 10))
717         return (B_FALSE);
718     return (B_TRUE);
821 }

823 static metabolab_ops_t metabolab_cf_ops = {
824     metabolab_cf_alloc,
825     metabolab_cf_fragmented
721 static space_map_ops_t metabolab_cdf_ops = {
722     metabolab_pp_load,
723     metabolab_pp_unload,
724     metabolab_cdf_alloc,
725     metabolab_pp_claim,
726     metabolab_pp_free,
727     metabolab_pp_maxsize,
728     metabolab_cdf_fragmented
826 };

828 /*
829 * =====
830 * New dynamic fit allocator -
831 * Select a region that is large enough to allocate 2^metabolab_ndf_clump_shift
832 * contiguous blocks. If no region is found then just use the largest segment
833 * that remains.
834 * =====
835 */

837 /*
838 * Determines desired number of contiguous blocks (2^metabolab_ndf_clump_shift)
839 * to request from the allocator.
840 */
841 uint64_t metabolab_ndf_clump_shift = 4;

843 static uint64_t
844 metabolab_ndf_alloc(metabolab_t *msp, uint64_t size)
734 metabolab_ndf_alloc(space_map_t *sm, uint64_t size)
845 {
846     avl_tree_t *t = &msp->ms_tree->rt_root;
736     avl_tree_t *t = &sm->sm_root;
847     avl_index_t where;
848     range_seg_t *rs, rsearch;

```

```

738 space_seg_t *ss, ssearch;
849 uint64_t hbit = highbit(size);
850 uint64_t *cursor = &msp->ms_lbas[hbit - 1];
851 uint64_t max_size = metabolab_block_maxsize(msp);
740 uint64_t *cursor = (uint64_t *)sm->sm_ppd + hbit - 1;
741 uint64_t max_size = metabolab_pp_maxsize(sm);

853 ASSERT(MUTEX_HELD(&msp->ms_lock));
854 ASSERT3U(avl_numnodes(t), ==, avl_numnodes(&msp->ms_size_tree));
743 ASSERT(MUTEX_HELD(sm->sm_lock));
744 ASSERT3U(avl_numnodes(&sm->sm_root), ==, avl_numnodes(sm->sm_pp_root));

856 if (max_size < size)
857     return (-1ULL);

859 rsearch.rs_start = *cursor;
860 rsearch.rs_end = *cursor + size;
749 ssearch.ss_start = *cursor;
750 ssearch.ss_end = *cursor + size;

862 rs = avl_find(t, &rsearch, &where);
863 if (rs == NULL || (rs->rs_end - rs->rs_start) < size) {
864     t = &msp->ms_size_tree;
752 ss = avl_find(t, &ssearch, &where);
753 if (ss == NULL || (ss->ss_start + size > ss->ss_end)) {
754     t = sm->sm_pp_root;

866 rsearch.rs_start = 0;
867 rsearch.rs_end = MIN(max_size,
756 ssearch.ss_start = 0;
757 ssearch.ss_end = MIN(max_size,
868     1ULL << (hbit + metabolab_ndf_clump_shift));
869 rs = avl_find(t, &rsearch, &where);
870 if (rs == NULL)
871     rs = avl_nearest(t, where, AVL_AFTER);
872 ASSERT(rs != NULL);
759 ss = avl_find(t, &ssearch, &where);
760 if (ss == NULL)
761     ss = avl_nearest(t, where, AVL_AFTER);
762 ASSERT(ss != NULL);
873 }

875 if ((rs->rs_end - rs->rs_start) >= size) {
876     *cursor = rs->rs_start + size;
877     return (rs->rs_start);
765 if (ss != NULL) {
766     if (ss->ss_start + size <= ss->ss_end) {
767         *cursor = ss->ss_start + size;
768         return (ss->ss_start);
878     }
770 }
879 return (-1ULL);
880 }

882 static boolean_t
883 metabolab_ndf_fragmented(metabolab_t *msp)
775 metabolab_ndf_fragmented(space_map_t *sm)
884 {
885     return (metabolab_block_maxsize(msp) <=
886         (metabolab_min_alloc_size << metabolab_ndf_clump_shift));
777     uint64_t max_size = metabolab_pp_maxsize(sm);

779 if (max_size > (metabolab_min_alloc_size << metabolab_ndf_clump_shift))
780     return (B_FALSE);
781 return (B_TRUE);
887 }

```

```

889 static metabolab_ops_t metabolab_ndf_ops = {
785 static space_map_ops_t metabolab_ndf_ops = {
786     metabolab_pp_load,
787     metabolab_pp_unload,
890     metabolab_ndf_alloc,
789     metabolab_pp_claim,
790     metabolab_pp_free,
791     metabolab_pp_maxsize,
891     metabolab_ndf_fragmented
892 };

894 metabolab_ops_t *zfs_metabolab_ops = &metabolab_df_ops;
795 space_map_ops_t *zfs_space_map_ops = &metabolab_df_ops;

896 /*
897 * =====
898 * Metabolabs
899 * =====
900 */

902 /*
903 * Wait for any in-progress metabolab loads to complete.
904 */
905 void
906 metabolab_load_wait(metabolab_t *msp)
907 {
908     ASSERT(MUTEX_HELD(&msp->ms_lock));

910 while (msp->ms_loading) {
911     ASSERT(!msp->ms_loaded);
912     cv_wait(&msp->ms_load_cv, &msp->ms_lock);
913 }
914 }

916 int
917 metabolab_load(metabolab_t *msp)
918 {
919     int error = 0;

921     ASSERT(MUTEX_HELD(&msp->ms_lock));
922     ASSERT(!msp->ms_loaded);
923     ASSERT(!msp->ms_loading);

925     msp->ms_loading = B_TRUE;

927 /*
928 * If the space map has not been allocated yet, then treat
929 * all the space in the metabolab as free and add it to the
930 * ms_tree.
931 */
932 if (msp->ms_sm != NULL)
933     error = space_map_load(msp->ms_sm, msp->ms_tree, SM_FREE);
934 else
935     range_tree_add(msp->ms_tree, msp->ms_start, msp->ms_size);

937 msp->ms_loaded = (error == 0);
938 msp->ms_loading = B_FALSE;

940 if (msp->ms_loaded) {
941     for (int t = 0; t < TXG_DEFER_SIZE; t++) {
942         range_tree_walk(msp->ms_defertree[t],
943             range_tree_remove, msp->ms_tree);
944     }
945 }

```

```

946     cv_broadcast(&msp->ms_load_cv);
947     return (error);
948 }

950 void
951 metablab_unload(metablab_t *msp)
952 {
953     ASSERT(MUTEX_HELD(&msp->ms_lock));
954     range_tree_vacate(msp->ms_tree, NULL, NULL);
955     msp->ms_loaded = B_FALSE;
956     msp->ms_weight &= ~METASLAB_ACTIVE_MASK;
957 }

959 metablab_t *
960 metablab_init(metablab_group_t *mg, uint64_t id, uint64_t object, uint64_t txg)
961 metablab_init(metablab_group_t *mg, space_map_obj_t *smo,
962               uint64_t start, uint64_t size, uint64_t txg)
963 {
964     vdev_t *vd = mg->mg_vd;
965     objset_t *mos = vd->vdev_spa->spa_meta_objset;
966     metablab_t *msp;

967     msp = kmem_zalloc(sizeof (metablab_t), KM_SLEEP);
968     mutex_init(&msp->ms_lock, NULL, MUTEX_DEFAULT, NULL);
969     cv_init(&msp->ms_load_cv, NULL, CV_DEFAULT, NULL);
970     msp->ms_id = id;
971     msp->ms_start = id << vd->vdev_ms_shift;
972     msp->ms_size = LULL << vd->vdev_ms_shift;

973     /*
974      * We only open space map objects that already exist. All others
975      * will be opened when we finally allocate an object for it.
976      */
977     if (object != 0) {
978         VERIFY0(space_map_open(&msp->ms_sm, mos, object, msp->ms_start,
979                               msp->ms_size, vd->vdev_ashift, &msp->ms_lock));
980         ASSERT(msp->ms_sm != NULL);
981     }
982     msp->ms_smo_syncing = *smo;

983     /*
984      * We create the main range tree here, but we don't create the
985      * alloc tree and freetree until metablab_sync_done(). This serves
986      * We create the main space map here, but we don't create the
987      * alloc maps and freemaps until metablab_sync_done(). This serves
988      * two purposes: it allows metablab_sync_done() to detect the
989      * addition of new space; and for debugging, it ensures that we'd
990      * data fault on any attempt to use this metablab before it's ready.
991      */
992     msp->ms_tree = range_tree_create(&metablab_rt_ops, msp, &msp->ms_lock);
993     msp->ms_map = kmem_zalloc(sizeof (space_map_t), KM_SLEEP);
994     space_map_create(msp->ms_map, start, size,
995                     vd->vdev_ashift, &msp->ms_lock);

996     metablab_group_add(mg, msp);

997     msp->ms_ops = mg->mg_class->mc_ops;
998     if (metablab_debug && smo->smo_object != 0) {
999         mutex_enter(&msp->ms_lock);
1000         VERIFY(space_map_load(msp->ms_map, mg->mg_class->mc_ops,
1001                               SM_FREE, smo, spa_meta_objset(vd->vdev_spa)) == 0);
1002         mutex_exit(&msp->ms_lock);
1003     }

1004     /*
1005      * If we're opening an existing pool (txg == 0) or creating

```

```

1006     * a new one (txg == TXG_INITIAL), all space is available now.
1007     * If we're adding space to an existing pool, the new space
1008     * does not become available until after this txg has synced.
1009     */
1010     if (txg <= TXG_INITIAL)
1011         metablab_sync_done(msp, 0);

1012     /*
1013      * If metablab_debug_load is set and we're initializing a metablab
1014      * that has an allocated space_map object then load the its space
1015      * map so that can verify frees.
1016      */
1017     if (metablab_debug_load && msp->ms_sm != NULL) {
1018         mutex_enter(&msp->ms_lock);
1019         VERIFY0(metablab_load(msp));
1020         mutex_exit(&msp->ms_lock);
1021     }

1022     if (txg != 0) {
1023         vdev_dirty(vd, 0, NULL, txg);
1024         vdev_dirty(vd, VDD_METASLAB, msp, txg);
1025     }

1026     return (msp);
1027 }

1028 void
1029 metablab_fini(metablab_t *msp)
1030 {
1031     metablab_group_t *mg = msp->ms_group;

1032     vdev_space_update(mg->mg_vd,
1033                      -msp->ms_smo.smo_alloc, 0, -msp->ms_map->sm_size);

1034     metablab_group_remove(mg, msp);

1035     mutex_enter(&msp->ms_lock);

1036     VERIFY(msp->ms_group == NULL);
1037     vdev_space_update(mg->mg_vd, -space_map_allocated(msp->ms_sm),
1038                      0, -msp->ms_size);
1039     space_map_close(msp->ms_sm);
1040     space_map_unload(msp->ms_map);
1041     space_map_destroy(msp->ms_map);
1042     kmem_free(msp->ms_map, sizeof (*msp->ms_map));

1043     metablab_unload(msp);
1044     range_tree_destroy(msp->ms_tree);

1045     for (int t = 0; t < TXG_SIZE; t++) {
1046         range_tree_destroy(msp->ms_allocmap[t]);
1047         range_tree_destroy(msp->ms_freemap[t]);
1048         space_map_destroy(msp->ms_allocmap[t]);
1049         space_map_destroy(msp->ms_freemap[t]);
1050         kmem_free(msp->ms_allocmap[t], sizeof (*msp->ms_allocmap[t]));
1051         kmem_free(msp->ms_freemap[t], sizeof (*msp->ms_freemap[t]));
1052     }

1053     for (int t = 0; t < TXG_DEFER_SIZE; t++) {
1054         range_tree_destroy(msp->ms_defertree[t]);
1055         space_map_destroy(msp->ms_defermap[t]);
1056         kmem_free(msp->ms_defermap[t], sizeof (*msp->ms_defermap[t]));
1057     }

1058     ASSERT0(msp->ms_deferspace);

```

```

1051     mutex_exit(&msp->ms_lock);
1052     cv_destroy(&msp->ms_load_cv);
1053     mutex_destroy(&msp->ms_lock);

1055     kmem_free(msp, sizeof (metabolab_t));
1056 }

1058 /*
1059  * Apply a weighting factor based on the histogram information for this
1060  * metabolab. The current weighting factor is somewhat arbitrary and requires
1061  * additional investigation. The implementation provides a measure of
1062  * "weighted" free space and gives a higher weighting for larger contiguous
1063  * regions. The weighting factor is determined by counting the number of
1064  * sm_shift sectors that exist in each region represented by the histogram.
1065  * That value is then multiplied by the power of 2 exponent and the sm_shift
1066  * value.
1067  *
1068  * For example, assume the 2^21 histogram bucket has 4 2MB regions and the
1069  * metabolab has an sm_shift value of 9 (512B):
1070  *
1071  * 1) calculate the number of sm_shift sectors in the region:
1072  *    2^21 / 2^9 = 2^12 = 4096 * 4 (number of regions) = 16384
1073  * 2) multiply by the power of 2 exponent and the sm_shift value:
1074  *    16384 * 21 * 9 = 3096576
1075  * This value will be added to the weighting of the metabolab.
1076  */
1077 static uint64_t
1078 metabolab_weight_factor(metabolab_t *msp)
1079 {
1080     uint64_t factor = 0;
1081     uint64_t sectors;
1082     int i;
1083     #define METASLAB_WEIGHT_PRIMARY      (1ULL << 63)
1084     #define METASLAB_WEIGHT_SECONDARY   (1ULL << 62)
1085     #define METASLAB_ACTIVE_MASK       (
1086     #define METASLAB_WEIGHT_PRIMARY | METASLAB_WEIGHT_SECONDARY )
1087     #define METASLAB_ACTIVE_MASK

1088     /*
1089     * A null space map means that the entire metabolab is free,
1090     * calculate a weight factor that spans the entire size of the
1091     * metabolab.
1092     */
1093     if (msp->ms_sm == NULL) {
1094         vdev_t *vd = msp->ms_group->mg_vd;
1095
1096         i = highbit(msp->ms_size) - 1;
1097         sectors = msp->ms_size >> vd->vdev_ashift;
1098         return (sectors * i * vd->vdev_ashift);
1099     }

1100     if (msp->ms_sm->sm_dbuf->db_size != sizeof (space_map_phys_t))
1101         return (0);

1102     for (i = 0; i < SPACE_MAP_HISTOGRAM_SIZE(msp->ms_sm); i++) {
1103         if (msp->ms_sm->sm_phys->smp_histogram[i] == 0)
1104             continue;

1105         /*
1106         * Determine the number of sm_shift sectors in the region
1107         * indicated by the histogram. For example, given an
1108         * sm_shift value of 9 (512 bytes) and i = 4 then we know
1109         * that we're looking at an 8K region in the histogram
1110         * (i.e. 9 + 4 = 13, 2^13 = 8192). To figure out the
1111         * number of sm_shift sectors (512 bytes in this example),
1112         * we would take 8192 / 512 = 16. Since the histogram
1113         * is offset by sm_shift we can simply use the value of

```

```

1113         * of i to calculate this (i.e. 2^i = 16 where i = 4).
1114         */
1115         sectors = msp->ms_sm->sm_phys->smp_histogram[i] << i;
1116         factor += (i + msp->ms_sm->sm_shift) * sectors;
1117     }
1118     return (factor * msp->ms_sm->sm_shift);
1119 }

1121 static uint64_t
1122 metabolab_weight(metabolab_t *msp)
1123 {
1124     metabolab_group_t *mg = msp->ms_group;
1125     space_map_t *sm = msp->ms_map;
1126     space_map_obj_t *smo = &msp->ms_smo;
1127     vdev_t *vd = mg->mg_vd;
1128     uint64_t weight, space;

1129     ASSERT(MUTEX_HELD(&msp->ms_lock));

1130     /*
1131     * This vdev is in the process of being removed so there is nothing
1132     * for us to do here.
1133     */
1134     if (vd->vdev_removing) {
1135         ASSERT0(space_map_allocated(msp->ms_sm));
1136         ASSERT0(smo->smo_alloc);
1137         ASSERT0(vd->vdev_ms_shift);
1138         return (0);
1139     }

1140     /*
1141     * The baseline weight is the metabolab's free space.
1142     */
1143     space = msp->ms_size - space_map_allocated(msp->ms_sm);
1144     space = sm->sm_size - smo->smo_alloc;
1145     weight = space;

1146     /*
1147     * Modern disks have uniform bit density and constant angular velocity.
1148     * Therefore, the outer recording zones are faster (higher bandwidth)
1149     * than the inner zones by the ratio of outer to inner track diameter,
1150     * which is typically around 2:1. We account for this by assigning
1151     * higher weight to lower metabolabs (multiplier ranging from 2x to 1x).
1152     * In effect, this means that we'll select the metabolab with the most
1153     * free bandwidth rather than simply the one with the most free space.
1154     */
1155     weight = 2 * weight - (msp->ms_id * weight) / vd->vdev_ms_count;
1156     weight = 2 * weight -
1157         ((sm->sm_start >> vd->vdev_ms_shift) * weight) / vd->vdev_ms_count;
1158     ASSERT(weight >= space && weight <= 2 * space);

1159     msp->ms_factor = metabolab_weight_factor(msp);
1160     if (metabolab_weight_factor_enable)
1161         weight += msp->ms_factor;

1162     /*
1163     * For locality, assign higher weight to metabolabs which have
1164     * a lower offset than what we've already activated.
1165     */
1166     if (sm->sm_start <= mg->mg_bonus_area)
1167         weight *= (metabolab_smo_bonus_pct / 100);
1168     ASSERT(weight >= space &&
1169         weight <= 2 * (metabolab_smo_bonus_pct / 100) * space);

1170     if (msp->ms_loaded && !msp->ms_ops->msop_fragmented(msp)) {
1171         if (sm->sm_loaded && !sm->sm_ops->smop_fragmented(sm)) {
1172             /*

```



```

1164     * If this metaslab is one we're actively using, adjust its
1165     * weight to make it preferable to any inactive metaslab so
1166     * we'll polish it off.
1167     */
1168     weight |= (msp->ms_weight & METASLAB_ACTIVE_MASK);
1169 }
1171     return (weight);
1172 }

952 static void
953 metaslab_prefetch(metaslab_group_t *mg)
954 {
955     spa_t *spa = mg->mg_vd->vdev_spa;
956     metaslab_t *msp;
957     avl_tree_t *t = &mg->mg_metaslab_tree;
958     int m;

960     mutex_enter(&mg->mg_lock);

962     /*
963     * Prefetch the next potential metaslabs
964     */
965     for (msp = avl_first(t), m = 0; msp; msp = AVL_NEXT(t, msp), m++) {
966         space_map_t *sm = msp->ms_map;
967         space_map_obj_t *smo = &msp->ms_smo;

969         /* If we have reached our prefetch limit then we're done */
970         if (m >= metaslab_prefetch_limit)
971             break;

973         if (!sm->sm_loaded && smo->smo_object != 0) {
974             mutex_exit(&mg->mg_lock);
975             dmuf_prefetch(spa_meta_objset(spa), smo->smo_object,
976                 0ULL, smo->smo_objsize);
977             mutex_enter(&mg->mg_lock);
978         }
979     }
980     mutex_exit(&mg->mg_lock);
981 }

1174 static int
1175 metaslab_activate(metaslab_t *msp, uint64_t activation_weight)
1176 {
986     metaslab_group_t *mg = msp->ms_group;
987     space_map_t *sm = msp->ms_map;
988     space_map_ops_t *sm_ops = msp->ms_group->mg_class->mc_ops;

1177     ASSERT(MUTEX_HELD(&msp->ms_lock));

1179     if ((msp->ms_weight & METASLAB_ACTIVE_MASK) == 0) {
1180         metaslab_load_wait(msp);
1181         if (!msp->ms_loaded) {
1182             int error = metaslab_load(msp);
993             space_map_load_wait(sm);
994             if (!sm->sm_loaded) {
995                 space_map_obj_t *smo = &msp->ms_smo;

997                 int error = space_map_load(sm, sm_ops, SM_FREE, smo,
998                     spa_meta_objset(msp->ms_group->mg_vd->vdev_spa));
1183                 if (error) {
1184                     metaslab_group_sort(msp->ms_group, msp, 0);
1185                     return (error);
1186                 }
1003                 for (int t = 0; t < TXG_DEFER_SIZE; t++)
1004                     space_map_walk(msp->ms_defermap[t],

```

```

1005         space_map_claim, sm);

1187     }

1009     /*
1010     * Track the bonus area as we activate new metaslabs.
1011     */
1012     if (sm->sm_start > mg->mg_bonus_area) {
1013         mutex_enter(&mg->mg_lock);
1014         mg->mg_bonus_area = sm->sm_start;
1015         mutex_exit(&mg->mg_lock);
1016     }

1189     metaslab_group_sort(msp->ms_group, msp,
1190         msp->ms_weight | activation_weight);
1191 }
1192 ASSERT(msp->ms_loaded);
1021 ASSERT(sm->sm_loaded);
1193 ASSERT(msp->ms_weight & METASLAB_ACTIVE_MASK);

1195     return (0);
1196 }

1198 static void
1199 metaslab_passivate(metaslab_t *msp, uint64_t size)
1200 {
1201     /*
1202     * If size < SPA_MINBLOCKSIZE, then we will not allocate from
1203     * this metaslab again. In that case, it had better be empty,
1204     * or we would be leaving space on the table.
1205     */
1206     ASSERT(size >= SPA_MINBLOCKSIZE || range_tree_space(msp->ms_tree) == 0);
1035     ASSERT(size >= SPA_MINBLOCKSIZE || msp->ms_map->sm_space == 0);
1207     metaslab_group_sort(msp->ms_group, msp, MIN(msp->ms_weight, size));
1208     ASSERT((msp->ms_weight & METASLAB_ACTIVE_MASK) == 0);
1209 }

1211 static void
1212 metaslab_preload(void *arg)
1213 {
1214     metaslab_t *msp = arg;
1215     spa_t *spa = msp->ms_group->mg_vd->vdev_spa;

1217     mutex_enter(&msp->ms_lock);
1218     metaslab_load_wait(msp);
1219     if (!msp->ms_loaded)
1220         (void) metaslab_load(msp);

1222     /*
1223     * Set the ms_access_txx value so that we don't unload it right away.
1224     */
1225     msp->ms_access_txx = spa_syncing_txx(spa) + metaslab_unload_delay + 1;
1226     mutex_exit(&msp->ms_lock);
1227 }

1229 static void
1230 metaslab_group_preload(metaslab_group_t *mg)
1231 {
1232     spa_t *spa = mg->mg_vd->vdev_spa;
1233     metaslab_t *msp;
1234     avl_tree_t *t = &mg->mg_metaslab_tree;
1235     int m = 0;

1237     if (spa_shutting_down(spa) || !metaslab_preload_enabled) {
1238         taskq_wait(mg->mg_taskq);
1239         return;

```

```

1240     }
1241     mutex_enter(&mg->mg_lock);

1243     /*
1244     * Prefetch the next potential metabolabs
1245     */
1246     for (msp = avl_first(t); msp != NULL; msp = AVL_NEXT(t, msp)) {

1248         /* If we have reached our preload limit then we're done */
1249         if (++m > metabolab_preload_limit)
1250             break;

1252         VERIFY(taskq_dispatch(mg->mg_taskq, metabolab_preload,
1253             msp, TQ_SLEEP) != NULL);
1254     }
1255     mutex_exit(&mg->mg_lock);
1256 }

1258 /*
1259 * Determine if the space map's on-disk footprint is past our tolerance
1260 * for inefficiency. We would like to use the following criteria to make
1261 * our decision:
1262 * 1. Determine if the in-core space map representation can be condensed on-disk.
1263 * We would like to use the following criteria to make our decision:
1264 * 1. The size of the space map object should not dramatically increase as a
1265 * result of writing out the free space range tree.
1266 * 2. result of writing out our in-core free map.
1267 * 1. The minimal on-disk space map representation is zfs_condense_pct/100
1268 * times the size than the free space range tree representation
1269 * (i.e. zfs_condense_pct = 110 and in-core = 1MB, minimal = 1.1.MB).
1270 * times the size than the in-core representation (i.e. zfs_condense_pct = 110
1271 * and in-core = 1MB, minimal = 1.1.MB).
1272 * Checking the first condition is tricky since we don't want to walk
1273 * the entire AVL tree calculating the estimated on-disk size. Instead we
1274 * use the size-ordered range tree in the metabolab and calculate the
1275 * size required to write out the largest segment in our free tree. If the
1276 * size required for the largest segment in our in-core free map. If the
1277 * size required to represent that segment on disk is larger than the space
1278 * map object then we avoid condensing this map.
1279 * To determine the second criterion we use a best-case estimate and assume
1280 * each segment can be represented on-disk as a single 64-bit entry. We refer
1281 * to this best-case estimate as the space map's minimal form.
1282 */
1283 static boolean_t
1284 metabolab_should_condense(metabolab_t *msp)
1285 {
1286     space_map_t *sm = msp->ms_sm;
1287     range_seg_t *rs;
1288     space_map_t *sm = msp->ms_map;
1289     space_map_obj_t *smo = &msp->ms_smo_syncing;
1290     space_seg_t *ss;
1291     uint64_t size, entries, segsz;

1293     ASSERT(MUTEX_HELD(&msp->ms_lock));
1294     ASSERT(msp->ms_loaded);
1295     ASSERT(sm->sm_loaded);

1297     /*
1298     * Use the ms_size_tree range tree, which is ordered by size, to
1299     * obtain the largest segment in the free tree. If the tree is empty
1300     * then we should condense the map.

```

```

1074     * Use the sm_pp_root AVL tree, which is ordered by size, to obtain
1075     * the largest segment in the in-core free map. If the tree is
1076     * empty then we should condense the map.
1077     */
1078     rs = avl_last(&msp->ms_size_tree);
1079     if (rs == NULL)
1080         ss = avl_last(sm->sm_pp_root);
1081     if (ss == NULL)
1082         return (B_TRUE);

1300     /*
1301     * Calculate the number of 64-bit entries this segment would
1302     * require when written to disk. If this single segment would be
1303     * larger on-disk than the entire current on-disk structure, then
1304     * clearly condensing will increase the on-disk structure size.
1305     */
1306     size = (rs->rs_end - rs->rs_start) >> sm->sm_shift;
1307     size = (ss->ss_end - ss->ss_start) >> sm->sm_shift;
1308     entries = size / (MIN(size, SM_RUN_MAX));
1309     segsz = entries * sizeof (uint64_t);

1310     return (segsz <= space_map_length(msp->ms_sm) &&
1311         space_map_length(msp->ms_sm) >= (zfs_condense_pct *
1312             sizeof (uint64_t) * avl_numnodes(&msp->ms_tree->rt_root)) / 100);
1313     return (segsz <= smo->smo_objsize &&
1314         smo->smo_objsize >= (zfs_condense_pct *
1315             sizeof (uint64_t) * avl_numnodes(&sm->sm_root)) / 100);
1316 }

1317 /*
1318 * Condense the on-disk space map representation to its minimized form.
1319 * The minimized form consists of a small number of allocations followed by
1320 * the entries of the free range tree.
1321 * the in-core free map.
1322 */
1323 static void
1324 metabolab_condense(metabolab_t *msp, uint64_t txg, dmu_tx_t *tx)
1325 {
1326     spa_t *spa = msp->ms_group->mg_vd->vdev_spa;
1327     range_tree_t *freetree = msp->ms_freetree[txg & TXG_MASK];
1328     range_tree_t *condense_tree;
1329     space_map_t *sm = msp->ms_sm;
1330     space_map_t *freemap = msp->ms_freemap[txg & TXG_MASK];
1331     space_map_t condense_map;
1332     space_map_t *sm = msp->ms_map;
1333     objset_t *mos = spa_meta_objset(spa);
1334     space_map_obj_t *smo = &msp->ms_smo_syncing;

1336     ASSERT(MUTEX_HELD(&msp->ms_lock));
1337     ASSERT3U(spa_sync_pass(spa), ==, 1);
1338     ASSERT(msp->ms_loaded);
1339     ASSERT(sm->sm_loaded);

1340     spa_dbgmsg(spa, "condensing: txg %llu, msp[%llu] %p, "
1341         "smp size %llu, segments %lu", txg, msp->ms_id, msp,
1342         space_map_length(msp->ms_sm), avl_numnodes(&msp->ms_tree->rt_root));
1343     "smp size %llu, segments %lu", txg,
1344     (msp->ms_map->sm_start / msp->ms_map->sm_size), msp,
1345     smo->smo_objsize, avl_numnodes(&sm->sm_root));

1347     /*
1348     * Create an range tree that is 100% allocated. We remove segments
1349     * Create an map that is a 100% allocated map. We remove segments
1350     * that have been freed in this txg, any deferred frees that exist,
1351     * and any allocation in the future. Removing segments should be
1352     * a relatively inexpensive operation since we expect these trees to

```

```

1341      * have a small number of nodes.
1342      * a relatively inexpensive operation since we expect these maps to
1343      * a small number of nodes.
1344      */
1345      condense_tree = range_tree_create(NULL, NULL, &msp->ms_lock);
1346      range_tree_add(condense_tree, msp->ms_start, msp->ms_size);
1347      space_map_create(&condense_map, sm->sm_start, sm->sm_size,
1348                     sm->sm_shift, sm->sm_lock);
1349      space_map_add(&condense_map, condense_map.sm_start,
1350                  condense_map.sm_size);
1351
1352      /*
1353       * Remove what's been freed in this txg from the condense_tree.
1354       * Remove what's been freed in this txg from the condense_map.
1355       * Since we're in sync_pass 1, we know that all the frees from
1356       * this txg are in the freetree.
1357       * This txg are in the freemap.
1358       */
1359      range_tree_walk(freetree, range_tree_remove, condense_tree);
1360      space_map_walk(freemap, space_map_remove, &condense_map);
1361
1362      for (int t = 0; t < TXG_DEFER_SIZE; t++) {
1363          range_tree_walk(msp->ms_defertree[t],
1364                        range_tree_remove, condense_tree);
1365      }
1366      for (int t = 0; t < TXG_DEFER_SIZE; t++)
1367          space_map_walk(msp->ms_defermap[t],
1368                       space_map_remove, &condense_map);
1369
1370      for (int t = 1; t < TXG_CONCURRENT_STATES; t++) {
1371          range_tree_walk(msp->ms_alloctree[(txg + t) & TXG_MASK],
1372                        range_tree_remove, condense_tree);
1373      }
1374      for (int t = 1; t < TXG_CONCURRENT_STATES; t++)
1375          space_map_walk(msp->ms_allocmap[(txg + t) & TXG_MASK],
1376                       space_map_remove, &condense_map);
1377
1378      /*
1379       * We're about to drop the metabolab's lock thus allowing
1380       * other consumers to change it's content. Set the
1381       * metabolab's ms_condensing flag to ensure that
1382       * space_map's sm_condensing flag to ensure that
1383       * allocations on this metabolab do not occur while we're
1384       * in the middle of committing it to disk. This is only critical
1385       * for the ms_tree as all other range trees use per txg
1386       * for the ms_map as all other space_maps use per txg
1387       * views of their content.
1388       */
1389      msp->ms_condensing = B_TRUE;
1390      sm->sm_condensing = B_TRUE;
1391
1392      mutex_exit(&msp->ms_lock);
1393      space_map_truncate(sm, tx);
1394      space_map_truncate(smo, mos, tx);
1395      mutex_enter(&msp->ms_lock);
1396
1397      /*
1398       * While we would ideally like to create a space_map representation
1399       * that consists only of allocation records, doing so can be
1400       * prohibitively expensive because the in-core free tree can be
1401       * prohibitively expensive because the in-core free map can be
1402       * large, and therefore computationally expensive to subtract
1403       * from the condense_tree. Instead we sync out two trees, a cheap
1404       * allocation only tree followed by the in-core free tree. While not
1405       * from the condense_map. Instead we sync out two maps, a cheap
1406       * allocation only map followed by the in-core free map. While not

```

```

1385      * optimal, this is typically close to optimal, and much cheaper to
1386      * compute.
1387      */
1388      space_map_write(sm, condense_tree, SM_ALLOC, tx);
1389      range_tree_vacate(condense_tree, NULL, NULL);
1390      range_tree_destroy(condense_tree);
1391      space_map_sync(&condense_map, SM_ALLOC, smo, mos, tx);
1392      space_map_vacate(&condense_map, NULL, NULL);
1393      space_map_destroy(&condense_map);
1394
1395      space_map_write(sm, msp->ms_tree, SM_FREE, tx);
1396      msp->ms_condensing = B_FALSE;
1397      space_map_sync(sm, SM_FREE, smo, mos, tx);
1398      sm->sm_condensing = B_FALSE;
1399
1400      spa_dbgmsg(spa, "condensed: txg %llu, msp[%llu] %p, "
1401                "smo size %llu", txg,
1402                (msp->ms_map->sm_start / msp->ms_map->sm_size), msp,
1403                smo->sm_objsize);
1404      }
1405
1406      /*
1407       * Write a metabolab to disk in the context of the specified transaction group.
1408       */
1409      void
1410      metabolab_sync(metabolab_t *msp, uint64_t txg)
1411      {
1412          metabolab_group_t *mg = msp->ms_group;
1413          vdev_t *vd = mg->mg_vd;
1414          vdev_t *vdev = msp->ms_group->mg_vdev;
1415          spa_t *spa = vd->vdev_spa;
1416          objset_t *mos = spa_meta_objset(spa);
1417          range_tree_t *alloctree = msp->ms_alloctree[txg & TXG_MASK];
1418          range_tree_t *freetree = &msp->ms_freetree[txg & TXG_MASK];
1419          range_tree_t *freed_tree =
1420              &msp->ms_freetree[TXG_CLEAN(txg) & TXG_MASK];
1421          space_map_t *allocmap = msp->ms_allocmap[txg & TXG_MASK];
1422          space_map_t *freemap = &msp->ms_freemap[txg & TXG_MASK];
1423          space_map_t *freed_map = &msp->ms_freemap[TXG_CLEAN(txg) & TXG_MASK];
1424          space_map_t *sm = msp->ms_map;
1425          space_map_obj_t *smo = &msp->ms_smo_syncing;
1426          dmu_buf_t *db;
1427          dmu_tx_t *tx;
1428          uint64_t object = space_map_object(msp->ms_sm);
1429
1430          ASSERT(!vd->vdev_ishole);
1431
1432          /*
1433           * This metabolab has just been added so there's no work to do now.
1434           */
1435          if (*freetree == NULL) {
1436              ASSERT3P(alloctree, ==, NULL);
1437          }
1438          if (*freemap == NULL) {
1439              ASSERT3P(allocmap, ==, NULL);
1440              return;
1441          }
1442
1443          ASSERT3P(alloctree, !=, NULL);
1444          ASSERT3P(*freetree, !=, NULL);
1445          ASSERT3P(*freed_tree, !=, NULL);
1446          ASSERT3P(allocmap, !=, NULL);
1447          ASSERT3P(*freemap, !=, NULL);
1448          ASSERT3P(*freed_map, !=, NULL);
1449
1450          if (range_tree_space(alloctree) == 0 &&
1451              range_tree_space(*freetree) == 0)

```

```

1217     if (allocmap->sm_space == 0 && (*freemap)->sm_space == 0)
1429         return;

1431     /*
1432     * The only state that can actually be changing concurrently with
1433     * metablab_sync() is the metablab's ms_tree. No other thread can
1434     * be modifying this txg's allocmap, freemap, freed_map, or
1435     * space_map_phys_t. Therefore, we only hold ms_lock to satisfy
1436     * space_map ASSERTs. We drop it whenever we call into the DMU,
1437     * because the DMU can call down to us (e.g. via zio_free()) at
1438     * any time.
1439     *
1440     * metablab_sync() is the metablab's ms_map. No other thread can
1441     * be modifying this txg's allocmap, freemap, freed_map, or smo.
1442     * Therefore, we only hold ms_lock to satisfy space_map ASSERTs.
1443     * We drop it whenever we call into the DMU, because the DMU
1444     * can call down to us (e.g. via zio_free()) at any time.
1445     */

1441     tx = dmu_tx_create_assigned(spa_get_dsl(spa), txg);

1443     if (msp->ms_sm == NULL) {
1444         uint64_t new_object;

1446         new_object = space_map_alloc(mos, tx);
1447         VERIFY3U(new_object, !=, 0);

1449         VERIFY0(space_map_open(&msp->ms_sm, mos, new_object,
1450             msp->ms_start, msp->ms_size, vd->vdev_ashift,
1451             &msp->ms_lock));
1452         ASSERT(msp->ms_sm != NULL);
1453         if (smo->smo_object == 0) {
1454             ASSERT(smo->smo_objsize == 0);
1455             ASSERT(smo->smo_alloc == 0);
1456             smo->smo_object = dmu_object_alloc(mos,
1457                 DMU_OT_SPACE_MAP, 1 << SPACE_MAP_BLOCKSHIFT,
1458                 DMU_OT_SPACE_MAP_HEADER, sizeof (*smo), tx);
1459             ASSERT(smo->smo_object != 0);
1460             dmu_write(mos, vd->vdev_ms_array, sizeof (uint64_t) *
1461                 (sm->sm_start >> vd->vdev_ms_shift),
1462                 sizeof (uint64_t), &smo->smo_object, tx);
1463         }

1465         mutex_enter(&msp->ms_lock);

1467         if (msp->ms_loaded && spa_sync_pass(spa) == 1 &&
1468             if (sm->sm_loaded && spa_sync_pass(spa) == 1 &&
1469                 metablab_should_condense(msp)) {
1470             metablab_condense(msp, txg, tx);
1471         } else {
1472             space_map_write(msp->ms_sm, allocmap, SM_ALLOC, tx);
1473             space_map_write(msp->ms_sm, *freemap, SM_FREE, tx);
1474             space_map_sync(allocmap, SM_ALLOC, smo, mos, tx);
1475             space_map_sync(*freemap, SM_FREE, smo, mos, tx);
1476         }

1477         range_tree_vacate(allocmap, NULL, NULL);
1478         space_map_vacate(allocmap, NULL, NULL);

1479         if (msp->ms_loaded) {
1480             /*
1481             * When the space map is loaded, we have an accurate
1482             * histogram in the range tree. This gives us an opportunity
1483             * to bring the space map's histogram up-to-date so we clear
1484             * it first before updating it.
1485             */
1486             space_map_histogram_clear(msp->ms_sm);

```

```

1475         space_map_histogram_add(msp->ms_sm, msp->ms_tree, tx);
1476     } else {
1477         /*
1478         * Since the space map is not loaded we simply update the
1479         * existing histogram with what was freed in this txg. This
1480         * means that the on-disk histogram may not have an accurate
1481         * view of the free space but it's close enough to allow
1482         * us to make allocation decisions.
1483         */
1484         space_map_histogram_add(msp->ms_sm, *freemap, tx);
1485     }

1487     /*
1488     * For sync pass 1, we avoid traversing this txg's free range tree
1489     * and instead will just swap the pointers for freemap and
1490     * freed_map. We can safely do this since the freed_map is
1491     * guaranteed to be empty on the initial pass.
1492     */
1493     if (spa_sync_pass(spa) == 1) {
1494         range_tree_swap(freemap, freed_map);
1495         ASSERT0((*freemap)->sm_space);
1496         ASSERT0(avl_numnodes(&(*freemap)->sm_root));
1497         space_map_swap(freemap, freed_map);
1498     } else {
1499         range_tree_vacate(*freemap, range_tree_add, *freemap);
1500         space_map_vacate(*freemap, space_map_add, *freemap);
1501     }

1502     ASSERT0(range_tree_space(msp->ms_allocmap[txg & TXG_MASK]));
1503     ASSERT0(range_tree_space(msp->ms_freemap[txg & TXG_MASK]));
1504     ASSERT0(msp->ms_allocmap[txg & TXG_MASK]->sm_space);
1505     ASSERT0(msp->ms_freemap[txg & TXG_MASK]->sm_space);

1506     mutex_exit(&msp->ms_lock);

1507     if (object != space_map_object(msp->ms_sm)) {
1508         object = space_map_object(msp->ms_sm);
1509         dmu_write(mos, vd->vdev_ms_array, sizeof (uint64_t) *
1510             msp->ms_id, sizeof (uint64_t), &object, tx);
1511     }

1512     VERIFY0(dmu_bonus_hold(mos, smo->smo_object, FTAG, &db));
1513     dmu_buf_will_dirty(db, tx);
1514     ASSERT3U(db->db_size, >=, sizeof (*smo));
1515     bcopy(smo, db->db_data, sizeof (*smo));
1516     dmu_buf_rele(db, FTAG);

1517     dmu_tx_commit(tx);

1518     /*
1519     * Called after a transaction group has completely synced to mark
1520     * all of the metablab's free space as usable.
1521     */
1522     void
1523     metablab_sync_done(metablab_t *mslab, uint64_t txg)
1524     {
1525         space_map_obj_t *smo = &mslab->ms_smo;
1526         space_map_obj_t *smosync = &mslab->ms_smo_syncing;
1527         space_map_t *sm = mslab->ms_map;
1528         space_map_t **freed_map = &mslab->ms_freemap[TXG_CLEAN(txg) & TXG_MASK];
1529         space_map_t **defer_map = &mslab->ms_defermap[txg % TXG_DEFER_SIZE];
1530         metablab_group_t *mg = mslab->ms_group;
1531         vdev_t *vd = mg->mg_vd;

```

```

1521 range_tree_t **freed_tree;
1522 range_tree_t **defer_tree;
1523 int64_t alloc_delta, defer_delta;

1525 ASSERT(!vd->vdev_ishole);

1527 mutex_enter(&mmp->ms_lock);

1529 /*
1530  * If this metabolab is just becoming available, initialize its
1531  * alloctrees, freetrees, and defertree and add its capacity to
1532  * the vdev.
1533  * allocmaps, freemaps, and defermap and add its capacity to the vdev.
1534  */
1535 if (mmp->ms_freemap[TXG_CLEAN(txg) & TXG_MASK] == NULL) {
1536     if (*freed_map == NULL) {
1537         ASSERT(*defer_map == NULL);
1538         for (int t = 0; t < TXG_SIZE; t++) {
1539             ASSERT(mmp->ms_alloctree[t] == NULL);
1540             ASSERT(mmp->ms_freemap[t] == NULL);
1541             mmp->ms_alloctree[t] = range_tree_create(NULL, mmp,
1542             &mmp->ms_lock);
1543             mmp->ms_freemap[t] = range_tree_create(NULL, mmp,
1544             &mmp->ms_lock);
1545             mmp->ms_allocmap[t] = kmem_zalloc(sizeof (space_map_t),
1546             KM_SLEEP);
1547             space_map_create(mmp->ms_allocmap[t], sm->sm_start,
1548             sm->sm_size, sm->sm_shift, sm->sm_lock);
1549             mmp->ms_freemap[t] = kmem_zalloc(sizeof (space_map_t),
1550             KM_SLEEP);
1551             space_map_create(mmp->ms_freemap[t], sm->sm_start,
1552             sm->sm_size, sm->sm_shift, sm->sm_lock);
1553         }

1554         for (int t = 0; t < TXG_DEFER_SIZE; t++) {
1555             ASSERT(mmp->ms_defertree[t] == NULL);
1556             mmp->ms_defertree[t] = range_tree_create(NULL, mmp,
1557             &mmp->ms_lock);
1558             mmp->ms_defermap[t] = kmem_zalloc(sizeof (space_map_t),
1559             KM_SLEEP);
1560             space_map_create(mmp->ms_defermap[t], sm->sm_start,
1561             sm->sm_size, sm->sm_shift, sm->sm_lock);
1562         }

1563         vdev_space_update(vd, 0, 0, mmp->ms_size);
1564         freed_map = &mmp->ms_freemap[TXG_CLEAN(txg) & TXG_MASK];
1565         defer_map = &mmp->ms_defermap[txg % TXG_DEFER_SIZE];

1566         vdev_space_update(vd, 0, 0, sm->sm_size);
1567     }

1568     freed_tree = &mmp->ms_freemap[TXG_CLEAN(txg) & TXG_MASK];
1569     defer_tree = &mmp->ms_defermap[txg % TXG_DEFER_SIZE];
1570     alloc_delta = smosync->smo_alloc - smo->smo_alloc;
1571     defer_delta = (*freed_map->sm_space - (*defer_map->sm_space);

1572     alloc_delta = space_map_alloc_delta(mmp->ms_sm);
1573     defer_delta = range_tree_space(*freed_tree) -
1574     range_tree_space(*defer_tree);

1575     vdev_space_update(vd, alloc_delta + defer_delta, defer_delta, 0);

1576     ASSERT0(range_tree_space(mmp->ms_alloctree[txg & TXG_MASK]));
1577     ASSERT0(range_tree_space(mmp->ms_freemap[txg & TXG_MASK]));

```

```

1338     ASSERT(mmp->ms_allocmap[txg & TXG_MASK]->sm_space == 0);
1339     ASSERT(mmp->ms_freemap[txg & TXG_MASK]->sm_space == 0);

1567     /*
1568     * If there's a metabolab_load() in progress, wait for it to complete
1569     * If there's a space_map_load() in progress, wait for it to complete
1570     * so that we have a consistent view of the in-core space map.
1571     */
1572     metabolab_load_wait(mmp);
1573     space_map_load_wait(sm);

1574     /*
1575     * Move the frees from the defer_tree back to the free
1576     * range tree (if it's loaded). Swap the freed_tree and the
1577     * defer_tree -- this is safe to do because we've just emptied out
1578     * the defer_tree.
1579     * Move the frees from the defer_map to this map (if it's loaded).
1580     * Swap the freed_map and the defer_map -- this is safe to do
1581     * because we've just emptied out the defer_map.
1582     */
1583     range_tree_vacate(*defer_tree,
1584     mmp->ms_loaded ? range_tree_add : NULL, mmp->ms_tree);
1585     range_tree_swap(freed_tree, defer_tree);
1586     space_map_vacate(*defer_map, sm->sm_loaded ? space_map_free : NULL, sm);
1587     ASSERT0((*defer_map)->sm_space);
1588     ASSERT0(avl_numnodes(&(*defer_map)->sm_root));
1589     space_map_swap(freed_map, defer_map);

1590     space_map_update(mmp->ms_sm);
1591     *smo = *smosync;

1592     mmp->ms_deferspace += defer_delta;
1593     ASSERT3S(mmp->ms_deferspace, >=, 0);
1594     ASSERT3S(mmp->ms_deferspace, <=, mmp->ms_size);
1595     ASSERT3S(mmp->ms_deferspace, <=, sm->sm_size);
1596     if (mmp->ms_deferspace != 0) {
1597         /*
1598         * Keep syncing this metabolab until all deferred frees
1599         * are back in circulation.
1600         */
1601         vdev_dirty(vd, VDD_METASLAB, mmp, txg + 1);
1602     }

1603     if (mmp->ms_loaded && mmp->ms_access_txg < txg) {
1604         for (int t = 1; t < TXG_CONCURRENT_STATES; t++) {
1605             VERIFY0(range_tree_space(
1606             mmp->ms_alloctree[(txg + t) & TXG_MASK]));
1607         }
1608     }
1609     /*
1610     * If the map is loaded but no longer active, evict it as soon as all
1611     * future allocations have synced. (If we unloaded it now and then
1612     * loaded a moment later, the map wouldn't reflect those allocations.)
1613     */
1614     if (sm->sm_loaded && (mmp->ms_weight & METASLAB_ACTIVE_MASK) == 0) {
1615         int evictable = 1;

1616         if (!metabolab_debug_unload)
1617             metabolab_unload(mmp);
1618         for (int t = 1; t < TXG_CONCURRENT_STATES; t++)
1619             if (mmp->ms_allocmap[(txg + t) & TXG_MASK]->sm_space)
1620                 evictable = 0;

1621         if (evictable && !metabolab_debug)
1622             space_map_unload(sm);
1623     }

```

```

1606     metablab_group_sort(mg, msp, metablab_weight(msp));
1607     mutex_exit(&msp->ms_lock);

1388     mutex_exit(&msp->ms_lock);
1609 }

1611 void
1612 metablab_sync_reassess(metablab_group_t *mg)
1613 {
1394     vdev_t *vd = mg->mg_vd;
1614     int64_t failures = mg->mg_alloc_failures;

1616     metablab_group_alloc_update(mg);

1399     /*
1400      * Re-evaluate all metablabs which have lower offsets than the
1401      * bonus area.
1402      */
1403     for (int m = 0; m < vd->vdev_ms_count; m++) {
1404         metablab_t *msp = vd->vdev_ms[m];

1406         if (msp->ms_map->sm_start > mg->mg_bonus_area)
1407             break;

1409         mutex_enter(&msp->ms_lock);
1410         metablab_group_sort(mg, msp, metablab_weight(msp));
1411         mutex_exit(&msp->ms_lock);
1412     }

1617     atomic_add_64(&mg->mg_alloc_failures, -failures);

1619     /*
1620      * Preload the next potential metablabs
1621      * Prefetch the next potential metablabs
1622      */
1622     metablab_group_preload(mg);
1419     metablab_prefetch(mg);
1623 }

1625 static uint64_t
1626 metablab_distance(metablab_t *msp, dva_t *dva)
1627 {
1628     uint64_t ms_shift = msp->ms_group->mg_vd->vdev_ms_shift;
1629     uint64_t offset = DVA_GET_OFFSET(dva) >> ms_shift;
1630     uint64_t start = msp->ms_id;
1427     uint64_t start = msp->ms_map->sm_start >> ms_shift;

1632     if (msp->ms_group->mg_vd->vdev_id != DVA_GET_VDEV(dva))
1633         return (LULL << 63);

1635     if (offset < start)
1636         return ((start - offset) << ms_shift);
1637     if (offset > start)
1638         return ((offset - start) << ms_shift);
1639     return (0);
1640 }

1642 static uint64_t
1643 metablab_group_alloc(metablab_group_t *mg, uint64_t psize, uint64_t asize,
1644     uint64_t txg, uint64_t min_distance, dva_t *dva, int d, int flags)
1645 {
1646     spa_t *spa = mg->mg_vd->vdev_spa;
1647     metablab_t *msp = NULL;
1648     uint64_t offset = -LULL;
1649     avl_tree_t *t = &mg->mg_metablab_tree;
1650     uint64_t activation_weight;

```

```

1651     uint64_t target_distance;
1652     int i;

1654     activation_weight = METASLAB_WEIGHT_PRIMARY;
1655     for (i = 0; i < d; i++) {
1656         if (DVA_GET_VDEV(&dva[i]) == mg->mg_vd->vdev_id) {
1657             activation_weight = METASLAB_WEIGHT_SECONDARY;
1658             break;
1659         }
1660     }

1662     for (;;) {
1663         boolean_t was_active;

1665         mutex_enter(&mg->mg_lock);
1666         for (msp = avl_first(t); msp; msp = AVL_NEXT(t, msp)) {
1667             if (msp->ms_weight < asize) {
1668                 spa_dbgmsg(spa, "%s: failed to meet weight ",
1669                     "requirement: vdev %llu, txg %llu, mg %p, "
1670                     "msp %p, psize %llu, asize %llu, "
1671                     "failures %llu, weight %llu",
1672                     spa_name(spa), mg->mg_vd->vdev_id, txg,
1673                     mg, msp, psize, asize,
1674                     mg->mg_alloc_failures, msp->ms_weight);
1675                 mutex_exit(&mg->mg_lock);
1676                 return (-LULL);
1677             }

1679             /*
1680              * If the selected metablab is condensing, skip it.
1681              */
1682             if (msp->ms_condensing)
1479                 if (msp->ms_map->sm_condensing)
1683                 continue;

1685             was_active = msp->ms_weight & METASLAB_ACTIVE_MASK;
1686             if (activation_weight == METASLAB_WEIGHT_PRIMARY)
1687                 break;

1689             target_distance = min_distance +
1690                 (space_map_allocated(msp->ms_sm) != 0 ? 0 :
1691                 min_distance >> 1);
1487             (msp->ms_smo_smo_alloc ? 0 : min_distance >> 1);

1693             for (i = 0; i < d; i++)
1694                 if (metablab_distance(msp, &dva[i]) <
1695                     target_distance)
1696                     break;

1697             if (i == d)
1698                 break;
1699         }
1700         mutex_exit(&mg->mg_lock);
1701         if (msp == NULL)
1702             return (-LULL);

1704         mutex_enter(&msp->ms_lock);

1706         /*
1707          * If we've already reached the allowable number of failed
1708          * allocation attempts on this metablab group then we
1709          * consider skipping it. We skip it only if we're allowed
1710          * to "fast" gang, the physical size is larger than
1711          * a gang block, and we're attempting to allocate from
1712          * the primary metablab.
1713          */
1714         if (mg->mg_alloc_failures > zfs_mg_alloc_failures &&

```

```

1715     CAN_FASTGANG(flags) && psize > SPA_GANGBLOCKSIZE &&
1716     activation_weight == METASLAB_WEIGHT_PRIMARY) {
1717         spa_dbgmsg(spa, "%s: skipping metablab group: "
1718             "vdev %llu, txg %llu, mg %p, msp[%llu] %p, "
1719             "psize %llu, asize %llu, failures %llu",
1720             spa_name(spa), mg->mg_vd->vdev_id, txg, mg,
1721             msp->ms_id, msp, psize, asize,
1722             "vdev %llu, txg %llu, mg %p, psize %llu, "
1723             "asize %llu, failures %llu", spa_name(spa),
1724             mg->mg_vd->vdev_id, txg, mg, psize, asize,
1725             mg->mg_alloc_failures);
1726         mutex_exit(&msp->ms_lock);
1727         return (-1ULL);
1728     }
1729
1730 /*
1731  * Ensure that the metablab we have selected is still
1732  * capable of handling our request. It's possible that
1733  * another thread may have changed the weight while we
1734  * were blocked on the metablab lock.
1735  */
1736 if (msp->ms_weight < asize || (was_active &&
1737     !(msp->ms_weight & METASLAB_ACTIVE_MASK) &&
1738     activation_weight == METASLAB_WEIGHT_PRIMARY)) {
1739     mutex_exit(&msp->ms_lock);
1740     continue;
1741 }
1742
1743 if ((msp->ms_weight & METASLAB_WEIGHT_SECONDARY) &&
1744     activation_weight == METASLAB_WEIGHT_PRIMARY) {
1745     metablab_passivate(msp,
1746         msp->ms_weight & ~METASLAB_ACTIVE_MASK);
1747     mutex_exit(&msp->ms_lock);
1748     continue;
1749 }
1750
1751 if (metablab_activate(msp, activation_weight) != 0) {
1752     mutex_exit(&msp->ms_lock);
1753     continue;
1754 }
1755
1756 /*
1757  * If this metablab is currently condensing then pick again as
1758  * we can't manipulate this metablab until it's committed
1759  * to disk.
1760  */
1761 if (msp->ms_condensing) {
1762     if (msp->ms_map->sm_condensing) {
1763         mutex_exit(&msp->ms_lock);
1764         continue;
1765     }
1766
1767     if ((offset = metablab_block_alloc(msp, asize)) != -1ULL)
1768     if ((offset = space_map_alloc(msp->ms_map, asize)) != -1ULL)
1769         break;
1770
1771     atomic_inc_64(&mg->mg_alloc_failures);
1772
1773     metablab_passivate(msp, metablab_block_maxsize(msp));
1774     metablab_passivate(msp, space_map_maxsize(msp->ms_map));
1775
1776     mutex_exit(&msp->ms_lock);
1777 }
1778
1779 if (range_tree_space(msp->ms_allocmap[txg & TXG_MASK]) == 0)
1780 if (msp->ms_allocmap[txg & TXG_MASK]->sm_space == 0)

```

```

1773         vdev_dirty(mg->mg_vd, VDD_METASLAB, msp, txg);
1774
1775     range_tree_add(msp->ms_allocmap[txg & TXG_MASK], offset, asize);
1776     msp->ms_access_txg = txg + metablab_unload_delay;
1777     space_map_add(msp->ms_allocmap[txg & TXG_MASK], offset, asize);
1778
1779     mutex_exit(&msp->ms_lock);
1780
1781     return (offset);
1782 }
1783
1784 unchanged_portion_omitted
1785
1786 /*
1787  * Free the block represented by DVA in the context of the specified
1788  * transaction group.
1789  */
1790 static void
1791 metablab_free_dva(spa_t *spa, const dva_t *dva, uint64_t txg, boolean_t now)
1792 {
1793     uint64_t vdev = DVA_GET_VDEV(dva);
1794     uint64_t offset = DVA_GET_OFFSET(dva);
1795     uint64_t size = DVA_GET_ASIZ(dva);
1796     vdev_t *vd;
1797     metablab_t *msp;
1798
1799     ASSERT(DVA_IS_VALID(dva));
1800
1801     if (txg > spa_freeze_txg(spa))
1802         return;
1803
1804     if ((vd = vdev_lookup_top(spa, vdev)) == NULL ||
1805         (offset >> vd->vdev_ms_shift) >= vd->vdev_ms_count) {
1806         cmn_err(CE_WARN, "metablab_free_dva(): bad DVA %llu:%llu",
1807             (u_longlong_t)vdev, (u_longlong_t)offset);
1808         ASSERT(0);
1809         return;
1810     }
1811
1812     msp = vd->vdev_ms[offset >> vd->vdev_ms_shift];
1813
1814     if (DVA_GET_GANG(dva))
1815         size = vdev_psize_to_asize(vd, SPA_GANGBLOCKSIZE);
1816
1817     mutex_enter(&msp->ms_lock);
1818
1819     if (now) {
1820         range_tree_remove(msp->ms_allocmap[txg & TXG_MASK],
1821             offset, size);
1822         space_map_remove(msp->ms_allocmap[txg & TXG_MASK],
1823             offset, size);
1824
1825         VERIFY(!msp->ms_condensing);
1826         VERIFY3U(offset, >=, msp->ms_start);
1827         VERIFY3U(offset + size, <=, msp->ms_start + msp->ms_size);
1828         VERIFY3U(range_tree_space(msp->ms_tree) + size, <=,
1829             msp->ms_size);
1830         VERIFY0(P2PHASE(offset, 1ULL << vd->vdev_ashift));
1831         VERIFY0(P2PHASE(size, 1ULL << vd->vdev_ashift));
1832         range_tree_add(msp->ms_tree, offset, size);
1833         space_map_free(msp->ms_map, offset, size);
1834     } else {
1835         if (range_tree_space(msp->ms_freemap[txg & TXG_MASK]) == 0)
1836         if (msp->ms_freemap[txg & TXG_MASK]->sm_space == 0)
1837             vdev_dirty(vd, VDD_METASLAB, msp, txg);
1838         range_tree_add(msp->ms_freemap[txg & TXG_MASK],
1839             offset, size);
1840         space_map_add(msp->ms_freemap[txg & TXG_MASK], offset, size);

```

```

2039     }
2041     mutex_exit(&mmp->ms_lock);
2042 }

2044 /*
2045  * Intent log support: upon opening the pool after a crash, notify the SPA
2046  * of blocks that the intent log has allocated for immediate write, but
2047  * which are still considered free by the SPA because the last transaction
2048  * group didn't commit yet.
2049  */
2050 static int
2051 metablab_claim_dva(spa_t *spa, const dva_t *dva, uint64_t txg)
2052 {
2053     uint64_t vdev = DVA_GET_VDEV(dva);
2054     uint64_t offset = DVA_GET_OFFSET(dva);
2055     uint64_t size = DVA_GET_ASIZE(dva);
2056     vdev_t *vd;
2057     metablab_t *mmp;
2058     int error = 0;

2060     ASSERT(DVA_IS_VALID(dva));

2062     if ((vd = vdev_lookup_top(spa, vdev)) == NULL ||
2063         (offset >> vd->vdev_ms_shift) >= vd->vdev_ms_count)
2064         return (SET_ERROR(ENXIO));

2066     mmp = vd->vdev_ms[offset >> vd->vdev_ms_shift];

2068     if (DVA_GET_GANG(dva))
2069         size = vdev_psize_to_asize(vd, SPA_GANGBLOCKSIZE);

2071     mutex_enter(&mmp->ms_lock);

2073     if ((txg != 0 && spa_writeable(spa)) || !mmp->ms_loaded)
1859     if ((txg != 0 && spa_writeable(spa)) || !mmp->ms_map->sm_loaded)
2074         error = metablab_activate(mmp, METASLAB_WEIGHT_SECONDARY);

2076     if (error == 0 && !range_tree_contains(mmp->ms_tree, offset, size))
1862     if (error == 0 && !space_map_contains(mmp->ms_map, offset, size))
2077         error = SET_ERROR(ENOENT);

2079     if (error || txg == 0) { /* txg == 0 indicates dry run */
2080         mutex_exit(&mmp->ms_lock);
2081         return (error);
2082     }

2084     VERIFY(!mmp->ms_condensing);
2085     VERIFY0(P2PHASE(offset, 1ULL << vd->vdev_ashift));
2086     VERIFY0(P2PHASE(size, 1ULL << vd->vdev_ashift));
2087     VERIFY3U(range_tree_space(mmp->ms_tree) - size, <=, mmp->ms_size);
2088     range_tree_remove(mmp->ms_tree, offset, size);
1870     space_map_claim(mmp->ms_map, offset, size);

2090     if (spa_writeable(spa)) { /* don't dirty if we're zdb(1M) */
2091         if (range_tree_space(mmp->ms_alloctree[txg & TXG_MASK]) == 0)
1873         if (mmp->ms_alloctree[txg & TXG_MASK]->sm_space == 0)
2092             vdev_dirty(vd, VDD_METASLAB, mmp, txg);
2093         range_tree_add(mmp->ms_alloctree[txg & TXG_MASK], offset, size);
1875         space_map_add(mmp->ms_alloctree[txg & TXG_MASK], offset, size);
2094     }

2096     mutex_exit(&mmp->ms_lock);

2098     return (0);
2099 }

```

```

2101 int
2102 metablab_alloc(spa_t *spa, metablab_class_t *mc, uint64_t psize, blkptr_t *bp,
2103               int ndvas, uint64_t txg, blkptr_t *hintbp, int flags)
2104 {
2105     dva_t *dva = bp->blk_dva;
2106     dva_t *hintdva = hintbp->blk_dva;
2107     int error = 0;

2109     ASSERT(bp->blk_birth == 0);
2110     ASSERT(BP_PHYSICAL_BIRTH(bp) == 0);

2112     spa_config_enter(spa, SCL_ALLOC, FTAG, RW_READER);

2114     if (mc->mc_rotor == NULL) { /* no vdevs in this class */
2115         spa_config_exit(spa, SCL_ALLOC, FTAG);
2116         return (SET_ERROR(ENOSPC));
2117     }

2119     ASSERT(ndvas > 0 && ndvas <= spa_max_replication(spa));
2120     ASSERT(BP_GET_NDVAS(bp) == 0);
2121     ASSERT(hintbp == NULL || ndvas <= BP_GET_NDVAS(hintbp));

2123     for (int d = 0; d < ndvas; d++) {
2124         error = metablab_alloc_dva(spa, mc, psize, dva, d, hintdva,
2125                                   txg, flags);
2126         if (error != 0) {
1908             if (error) {
2127                 for (d--; d >= 0; d--) {
2128                     metablab_free_dva(spa, &dva[d], txg, B_TRUE);
2129                     bzero(&dva[d], sizeof (dva_t));
2130                 }
2131                 spa_config_exit(spa, SCL_ALLOC, FTAG);
2132                 return (error);
2133             }
2134         }
2135         ASSERT(error == 0);
2136         ASSERT(BP_GET_NDVAS(bp) == ndvas);

2138     spa_config_exit(spa, SCL_ALLOC, FTAG);

2140     BP_SET_BIRTH(bp, txg, txg);

2142     return (0);
2143 }
_____ unchanged portion omitted _____

1975 static void
1976 checkmap(space_map_t *sm, uint64_t off, uint64_t size)
1977 {
1978     space_seg_t *ss;
1979     avl_index_t where;

1981     mutex_enter(sm->sm_lock);
1982     ss = space_map_find(sm, off, size, &where);
1983     if (ss != NULL)
1984         panic("freeing free block; ss=%p", (void *)ss);
1985     mutex_exit(sm->sm_lock);
1986 }

2193 void
2194 metablab_check_free(spa_t *spa, const blkptr_t *bp)
2195 {
2196     if ((zfs_flags & ZFS_DEBUG_ZIO_FREE) == 0)
2197         return;

```



```
2199     spa_config_enter(spa, SCL_VDEV, FTAG, RW_READER);
2200     for (int i = 0; i < BP_GET_NDVAS(bp); i++) {
2201         uint64_t vdev = DVA_GET_VDEV(&bp->blk_dva[i]);
2202         vdev_t *vd = vdev_lookup_top(spa, vdev);
2203         uint64_t offset = DVA_GET_OFFSET(&bp->blk_dva[i]);
2204         uint64_t vdid = DVA_GET_VDEV(&bp->blk_dva[i]);
2205         vdev_t *vd = vdev_lookup_top(spa, vdid);
2206         uint64_t off = DVA_GET_OFFSET(&bp->blk_dva[i]);
2207         uint64_t size = DVA_GET_ASIZE(&bp->blk_dva[i]);
2208         metaslab_t *msp = vd->vdev_ms[offset >> vd->vdev_ms_shift];
2209         metaslab_t *ms = vd->vdev_ms[off >> vd->vdev_ms_shift];
2210
2211         if (msp->ms_loaded)
2212             range_tree_verify(msp->ms_tree, offset, size);
2213         if (ms->ms_map->sm_loaded)
2214             checkmap(ms->ms_map, off, size);
2215
2216         for (int j = 0; j < TXG_SIZE; j++)
2217             range_tree_verify(msp->ms_freemap[j], offset, size);
2218         checkmap(ms->ms_freemap[j], off, size);
2219         for (int j = 0; j < TXG_DEFER_SIZE; j++)
2220             range_tree_verify(msp->ms_defertree[j], offset, size);
2221         checkmap(ms->ms_defermap[j], off, size);
2222     }
2223     spa_config_exit(spa, SCL_VDEV, FTAG);
2224 }
```

_____unchanged_portion_omitted_____

```

*****
9202 Tue Sep 3 20:27:00 2013
new/usr/src/uts/common/fs/zfs/range_tree.c
4101 metaslab_debug should allow for fine-grained control
4102 space_maps should store more information about themselves
4103 space_map object blocksize should be increased
4104 ::spa_space no longer works
4105 removing a mirrored log device results in a leaked object
4106 asynchronously load metaslab
Reviewed by: Matthew Ahrens <mahrens@delphix.com>
Reviewed by: Adam Leventhal <ahl@delphix.com>
Reviewed by: Sebastien Roy <seb@delphix.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22 * Copyright 2009 Sun Microsystems, Inc. All rights reserved.
23 * Use is subject to license terms.
24 */
25 /*
26 * Copyright (c) 2013 by Delphix. All rights reserved.
27 */
28
29 #include <sys/zfs_context.h>
30 #include <sys/spa.h>
31 #include <sys/dmu.h>
32 #include <sys/dnode.h>
33 #include <sys/zio.h>
34 #include <sys/range_tree.h>
35
36 static kmem_cache_t *range_seg_cache;
37
38 void
39 range_tree_init(void)
40 {
41     ASSERT(range_seg_cache == NULL);
42     range_seg_cache = kmem_cache_create("range_seg_cache",
43     sizeof (range_seg_t), 0, NULL, NULL, NULL, NULL, 0);
44 }
45
46 void
47 range_tree_fini(void)
48 {
49     kmem_cache_destroy(range_seg_cache);
50     range_seg_cache = NULL;
51 }
52
53 void

```

```

54 range_tree_stat_verify(range_tree_t *rt)
55 {
56     range_seg_t *rs;
57     uint64_t hist[RANGE_TREE_HISTOGRAM_SIZE] = { 0 };
58     int i;
59
60     for (rs = avl_first(&rt->rt_root); rs != NULL;
61          rs = AVL_NEXT(&rt->rt_root, rs)) {
62         uint64_t size = rs->rs_end - rs->rs_start;
63         int idx = highbit(size) - 1;
64
65         hist[idx]++;
66         ASSERT3U(hist[idx], !=, 0);
67     }
68
69     for (i = 0; i < RANGE_TREE_HISTOGRAM_SIZE; i++) {
70         if (hist[i] != rt->rt_histogram[i]) {
71             zfs_dbgmsg("i=%d, hist=%p, hist=%llu, rt_hist=%llu",
72                 i, hist, hist[i], rt->rt_histogram[i]);
73         }
74         VERIFY3U(hist[i], ==, rt->rt_histogram[i]);
75     }
76 }
77
78 static void
79 range_tree_stat_incr(range_tree_t *rt, range_seg_t *rs)
80 {
81     uint64_t size = rs->rs_end - rs->rs_start;
82     int idx = highbit(size) - 1;
83
84     ASSERT3U(idx, <,
85         sizeof (rt->rt_histogram) / sizeof (*rt->rt_histogram));
86
87     ASSERT(MUTEX_HELD(rt->rt_lock));
88     rt->rt_histogram[idx]++;
89     ASSERT3U(rt->rt_histogram[idx], !=, 0);
90 }
91
92 static void
93 range_tree_stat_decr(range_tree_t *rt, range_seg_t *rs)
94 {
95     uint64_t size = rs->rs_end - rs->rs_start;
96     int idx = highbit(size) - 1;
97
98     ASSERT3U(idx, <,
99         sizeof (rt->rt_histogram) / sizeof (*rt->rt_histogram));
100
101     ASSERT(MUTEX_HELD(rt->rt_lock));
102     ASSERT3U(rt->rt_histogram[idx], !=, 0);
103     rt->rt_histogram[idx]--;
104 }
105
106 /*
107  * NOTE: caller is responsible for all locking.
108  */
109 static int
110 range_tree_seg_compare(const void *x1, const void *x2)
111 {
112     const range_seg_t *r1 = x1;
113     const range_seg_t *r2 = x2;
114
115     if (r1->rs_start < r2->rs_start) {
116         if (r1->rs_end > r2->rs_start)
117             return (0);
118         return (-1);
119     }

```

```

120     if (r1->rs_start > r2->rs_start) {
121         if (r1->rs_start < r2->rs_end)
122             return (0);
123         return (1);
124     }
125     return (0);
126 }

128 range_tree_t *
129 range_tree_create(range_tree_ops_t *ops, void *arg, kmutex_t *lp)
130 {
131     range_tree_t *rt;

133     rt = kmem_zalloc(sizeof (range_tree_t), KM_SLEEP);

135     avl_create(&rt->rt_root, range_tree_seg_compare,
136             sizeof (range_seg_t), offsetof(range_seg_t, rs_node));

138     rt->rt_lock = lp;
139     rt->rt_ops = ops;
140     rt->rt_arg = arg;

142     if (rt->rt_ops != NULL)
143         rt->rt_ops->rtop_create(rt, rt->rt_arg);

145     return (rt);
146 }

148 void
149 range_tree_destroy(range_tree_t *rt)
150 {
151     VERIFY0(rt->rt_space);

153     if (rt->rt_ops != NULL)
154         rt->rt_ops->rtop_destroy(rt, rt->rt_arg);

156     avl_destroy(&rt->rt_root);
157     kmem_free(rt, sizeof (*rt));
158 }

160 void
161 range_tree_add(void *arg, uint64_t start, uint64_t size)
162 {
163     range_tree_t *rt = arg;
164     avl_index_t where;
165     range_seg_t rsearch, *rs_before, *rs_after, *rs;
166     uint64_t end = start + size;
167     boolean_t merge_before, merge_after;

169     ASSERT(MUTEX_HELD(rt->rt_lock));
170     VERIFY(size != 0);

172     rsearch.rs_start = start;
173     rsearch.rs_end = end;
174     rs = avl_find(&rt->rt_root, &rsearch, &where);

176     if (rs != NULL && rs->rs_start <= start && rs->rs_end >= end) {
177         zfs_panic_recover("zfs: allocating allocated segment"
178             "(offset=%llu size=%llu)\n",
179             (longlong_t)start, (longlong_t)size);
180         return;
181     }

183     /* Make sure we don't overlap with either of our neighbors */
184     VERIFY(rs == NULL);

```

```

186     rs_before = avl_nearest(&rt->rt_root, where, AVL_BEFORE);
187     rs_after = avl_nearest(&rt->rt_root, where, AVL_AFTER);

189     merge_before = (rs_before != NULL && rs_before->rs_end == start);
190     merge_after = (rs_after != NULL && rs_after->rs_start == end);

192     if (merge_before && merge_after) {
193         avl_remove(&rt->rt_root, rs_before);
194         if (rt->rt_ops != NULL) {
195             rt->rt_ops->rtop_remove(rt, rs_before, rt->rt_arg);
196             rt->rt_ops->rtop_remove(rt, rs_after, rt->rt_arg);
197         }

199         range_tree_stat_decr(rt, rs_before);
200         range_tree_stat_decr(rt, rs_after);

202         rs_after->rs_start = rs_before->rs_start;
203         kmem_cache_free(range_seg_cache, rs_before);
204         rs = rs_after;
205     } else if (merge_before) {
206         if (rt->rt_ops != NULL)
207             rt->rt_ops->rtop_remove(rt, rs_before, rt->rt_arg);

209         range_tree_stat_decr(rt, rs_before);

211         rs_before->rs_end = end;
212         rs = rs_before;
213     } else if (merge_after) {
214         if (rt->rt_ops != NULL)
215             rt->rt_ops->rtop_remove(rt, rs_after, rt->rt_arg);

217         range_tree_stat_decr(rt, rs_after);

219         rs_after->rs_start = start;
220         rs = rs_after;
221     } else {
222         rs = kmem_cache_alloc(range_seg_cache, KM_SLEEP);
223         rs->rs_start = start;
224         rs->rs_end = end;
225         avl_insert(&rt->rt_root, rs, where);
226     }

228     if (rt->rt_ops != NULL)
229         rt->rt_ops->rtop_add(rt, rs, rt->rt_arg);

231     range_tree_stat_incr(rt, rs);
232     rt->rt_space += size;
233 }

235 void
236 range_tree_remove(void *arg, uint64_t start, uint64_t size)
237 {
238     range_tree_t *rt = arg;
239     avl_index_t where;
240     range_seg_t rsearch, *rs, *newseg;
241     uint64_t end = start + size;
242     boolean_t left_over, right_over;

244     ASSERT(MUTEX_HELD(rt->rt_lock));
245     VERIFY3U(size, !=, 0);
246     VERIFY3U(size, <=, rt->rt_space);

248     rsearch.rs_start = start;
249     rsearch.rs_end = end;
250     rs = avl_find(&rt->rt_root, &rsearch, &where);

```

```

252  /* Make sure we completely overlap with someone */
253  if (rs == NULL) {
254      zfs_panic_recover("zfs: freeing free segment "
255                      "(offset=%llu size=%llu)",
256                      (longlong_t)start, (longlong_t)size);
257      return;
258  }
259  VERIFY3U(rs->rs_start, <=, start);
260  VERIFY3U(rs->rs_end, >=, end);
261
262  left_over = (rs->rs_start != start);
263  right_over = (rs->rs_end != end);
264
265  range_tree_stat_decr(rt, rs);
266
267  if (rt->rt_ops != NULL)
268      rt->rt_ops->rtop_remove(rt, rs, rt->rt_arg);
269
270  if (left_over && right_over) {
271      newseg = kmem_cache_alloc(range_seg_cache, KM_SLEEP);
272      newseg->rs_start = end;
273      newseg->rs_end = rs->rs_end;
274      range_tree_stat_incr(rt, newseg);
275
276      rs->rs_end = start;
277
278      avl_insert_here(&rt->rt_root, newseg, rs, AVL_AFTER);
279      if (rt->rt_ops != NULL)
280          rt->rt_ops->rtop_add(rt, newseg, rt->rt_arg);
281  } else if (left_over) {
282      rs->rs_end = start;
283  } else if (right_over) {
284      rs->rs_start = end;
285  } else {
286      avl_remove(&rt->rt_root, rs);
287      kmem_cache_free(range_seg_cache, rs);
288      rs = NULL;
289  }
290
291  if (rs != NULL) {
292      range_tree_stat_incr(rt, rs);
293
294      if (rt->rt_ops != NULL)
295          rt->rt_ops->rtop_add(rt, rs, rt->rt_arg);
296  }
297
298  rt->rt_space -= size;
299  }
300
301 static range_seg_t *
302 range_tree_find(range_tree_t *rt, uint64_t start, uint64_t size,
303               avl_index_t *wherep)
304 {
305     range_seg_t rsearch, *rs;
306     uint64_t end = start + size;
307
308     ASSERT(MUTEX_HELD(rt->rt_lock));
309     VERIFY(size != 0);
310
311     rsearch.rs_start = start;
312     rsearch.rs_end = end;
313     rs = avl_find(&rt->rt_root, &rsearch, wherep);
314
315     if (rs != NULL && rs->rs_start <= start && rs->rs_end >= end)
316         return (rs);
317     return (NULL);

```

```

318 }
319
320 void
321 range_tree_verify(range_tree_t *rt, uint64_t off, uint64_t size)
322 {
323     range_seg_t *rs;
324     avl_index_t where;
325
326     mutex_enter(rt->rt_lock);
327     rs = range_tree_find(rt, off, size, &where);
328     if (rs != NULL)
329         panic("freeing free block; rs=%p", (void *)rs);
330     mutex_exit(rt->rt_lock);
331 }
332
333 boolean_t
334 range_tree_contains(range_tree_t *rt, uint64_t start, uint64_t size)
335 {
336     avl_index_t where;
337
338     return (range_tree_find(rt, start, size, &where) != NULL);
339 }
340
341 void
342 range_tree_swap(range_tree_t **rtsrc, range_tree_t **rtdst)
343 {
344     range_tree_t *rt;
345
346     ASSERT(MUTEX_HELD((*rtsrc)->rt_lock));
347     ASSERT0(range_tree_space(*rtdst));
348     ASSERT0(avl_numnodes(&(*rtdst)->rt_root));
349
350     rt = *rtsrc;
351     *rtsrc = *rtdst;
352     *rtdst = rt;
353 }
354
355 void
356 range_tree_vacate(range_tree_t *rt, range_tree_func_t *func, void *arg)
357 {
358     range_seg_t *rs;
359     void *cookie = NULL;
360
361     ASSERT(MUTEX_HELD(rt->rt_lock));
362
363     if (rt->rt_ops != NULL)
364         rt->rt_ops->rtop_vacate(rt, rt->rt_arg);
365
366     while ((rs = avl_destroy_nodes(&rt->rt_root, &cookie)) != NULL) {
367         if (func != NULL)
368             func(arg, rs->rs_start, rs->rs_end - rs->rs_start);
369         kmem_cache_free(range_seg_cache, rs);
370     }
371
372     bzero(rt->rt_histogram, sizeof (rt->rt_histogram));
373     rt->rt_space = 0;
374 }
375
376 void
377 range_tree_walk(range_tree_t *rt, range_tree_func_t *func, void *arg)
378 {
379     range_seg_t *rs;
380
381     ASSERT(MUTEX_HELD(rt->rt_lock));
382
383     for (rs = avl_first(&rt->rt_root); rs; rs = AVL_NEXT(&rt->rt_root, rs))

```

```
384         func(arg, rs->rs_start, rs->rs_end - rs->rs_start);
385     }
387     uint64_t
388     range_tree_space(range_tree_t *rt)
389     {
390         return (rt->rt_space);
391     }
```

new/usr/src/uts/common/fs/zfs/spa.c

1

```
*****
176003 Tue Sep 3 20:27:01 2013
new/usr/src/uts/common/fs/zfs/spa.c
4101 metaslab_debug should allow for fine-grained control
4102 space_maps should store more information about themselves
4103 space_map object blocksize should be increased
4104 ::spa_space no longer works
4105 removing a mirrored log device results in a leaked object
4106 asynchronously load metaslab
Reviewed by: Matthew Ahrens <mahrens@delphix.com>
Reviewed by: Adam Leventhal <ahl@delphix.com>
Reviewed by: Sebastien Roy <seb@delphix.com>
*****
unchanged_portion_omitted
```

```
1197 /*
1198  * Opposite of spa_load().
1199  */
1200 static void
1201 spa_unload(spa_t *spa)
1202 {
1203     int i;
1204
1205     ASSERT(MUTEX_HELD(&spa_namespace_lock));
1206
1207     /*
1208      * Stop async tasks.
1209      */
1210     spa_async_suspend(spa);
1211
1212     /*
1213      * Stop syncing.
1214      */
1215     if (spa->spa_sync_on) {
1216         txg_sync_stop(spa->spa_dsl_pool);
1217         spa->spa_sync_on = B_FALSE;
1218     }
1219
1220     /*
1221      * Wait for any outstanding async I/O to complete.
1222      */
1223     if (spa->spa_async_zio_root != NULL) {
1224         (void) zio_wait(spa->spa_async_zio_root);
1225         spa->spa_async_zio_root = NULL;
1226     }
1227
1228     bpobj_close(&spa->spa_deferred_bpobj);
1229
1230     spa_config_enter(spa, SCL_ALL, FTAG, RW_WRITER);
1231
1232     /*
1233      * Close all vdevs.
1234      */
1235     if (spa->spa_root_vdev)
1236         vdev_free(spa->spa_root_vdev);
1237     ASSERT(spa->spa_root_vdev == NULL);
1238
1239     /*
1240      * Close the dsl pool.
1241      */
1242     if (spa->spa_dsl_pool) {
1243         dsl_pool_close(spa->spa_dsl_pool);
1244         spa->spa_dsl_pool = NULL;
1245         spa->spa_meta_objset = NULL;
1246     }

```

new/usr/src/uts/common/fs/zfs/spa.c

2

```
1248     ddt_unload(spa);
1249
1250     spa_config_enter(spa, SCL_ALL, FTAG, RW_WRITER);
1251
1252     /*
1253      * Drop and purge level 2 cache
1254      */
1255     spa_l2cache_drop(spa);
1256
1257     /*
1258      * Close all vdevs.
1259      */
1260     if (spa->spa_root_vdev)
1261         vdev_free(spa->spa_root_vdev);
1262     ASSERT(spa->spa_root_vdev == NULL);
1263
1264     for (i = 0; i < spa->spa_spares.sav_count; i++)
1265         vdev_free(spa->spa_spares.sav_vdevs[i]);
1266     if (spa->spa_spares.sav_vdevs) {
1267         kmem_free(spa->spa_spares.sav_vdevs,
1268             spa->spa_spares.sav_count * sizeof(void *));
1269         spa->spa_spares.sav_vdevs = NULL;
1270     }
1271     if (spa->spa_spares.sav_config) {
1272         nvlist_free(spa->spa_spares.sav_config);
1273         spa->spa_spares.sav_config = NULL;
1274     }
1275     spa->spa_spares.sav_count = 0;
1276
1277     for (i = 0; i < spa->spa_l2cache.sav_count; i++) {
1278         vdev_clear_stats(spa->spa_l2cache.sav_vdevs[i]);
1279         vdev_free(spa->spa_l2cache.sav_vdevs[i]);
1280     }
1281     if (spa->spa_l2cache.sav_vdevs) {
1282         kmem_free(spa->spa_l2cache.sav_vdevs,
1283             spa->spa_l2cache.sav_count * sizeof(void *));
1284         spa->spa_l2cache.sav_vdevs = NULL;
1285     }
1286     if (spa->spa_l2cache.sav_config) {
1287         nvlist_free(spa->spa_l2cache.sav_config);
1288         spa->spa_l2cache.sav_config = NULL;
1289     }
1290     spa->spa_l2cache.sav_count = 0;
1291
1292     spa->spa_async_suspended = 0;
1293
1294     if (spa->spa_comment != NULL) {
1295         spa_strfree(spa->spa_comment);
1296         spa->spa_comment = NULL;
1297     }
1298
1299     spa_config_exit(spa, SCL_ALL, FTAG);
1300 }
unchanged_portion_omitted
4322 /*
4323  * Attach a device to a mirror. The arguments are the path to any device
4324  * in the mirror, and the nvroot for the new device. If the path specifies
4325  * a device that is not mirrored, we automatically insert the mirror vdev.
4326  *
4327  * If 'replacing' is specified, the new device is intended to replace the
4328  * existing device; in this case the two devices are made into their own
4329  * mirror using the 'replacing' vdev, which is functionally identical to
4330  * the mirror vdev (it actually reuses all the same ops) but has a few
4331  * extra rules: you can't attach to it after it's been created, and upon
4332  * completion of resilvering, the first disk (the one being replaced)

```

```

4333 * is automatically detached.
4334 */
4335 int
4336 spa_vdev_attach(spa_t *spa, uint64_t guid, nvlist_t *nvroot, int replacing)
4337 {
4338     uint64_t txg, dtl_max_txg;
4339     vdev_t *rvd = spa->spa_root_vdev;
4340     vdev_t *oldvd, *newvd, *newrootvd, *pvd, *tvd;
4341     vdev_ops_t *pvops;
4342     char *oldvdpath, *newvdpath;
4343     int newvd_isspare;
4344     int error;
4345
4346     ASSERT(spa_writeable(spa));
4347
4348     txg = spa_vdev_enter(spa);
4349
4350     oldvd = spa_lookup_by_guid(spa, guid, B_FALSE);
4351
4352     if (oldvd == NULL)
4353         return (spa_vdev_exit(spa, NULL, txg, ENODEV));
4354
4355     if (!oldvd->vdev_ops->vdev_op_leaf)
4356         return (spa_vdev_exit(spa, NULL, txg, ENOTSUP));
4357
4358     pvd = oldvd->vdev_parent;
4359
4360     if ((error = spa_config_parse(spa, &newrootvd, nvroot, NULL, 0,
4361         VDEV_ALLOC_ATTACH)) != 0)
4362         return (spa_vdev_exit(spa, NULL, txg, EINVAL));
4363
4364     if (newrootvd->vdev_children != 1)
4365         return (spa_vdev_exit(spa, newrootvd, txg, EINVAL));
4366
4367     newvd = newrootvd->vdev_child[0];
4368
4369     if (!newvd->vdev_ops->vdev_op_leaf)
4370         return (spa_vdev_exit(spa, newrootvd, txg, EINVAL));
4371
4372     if ((error = vdev_create(newrootvd, txg, replacing)) != 0)
4373         return (spa_vdev_exit(spa, newrootvd, txg, error));
4374
4375     /*
4376      * Spares can't replace logs
4377      */
4378     if (oldvd->vdev_top->vdev_islog && newvd->vdev_isspare)
4379         return (spa_vdev_exit(spa, newrootvd, txg, ENOTSUP));
4380
4381     if (!replacing) {
4382         /*
4383          * For attach, the only allowable parent is a mirror or the root
4384          * vdev.
4385          */
4386         if (pvd->vdev_ops != &vdev_mirror_ops &&
4387             pvd->vdev_ops != &vdev_root_ops)
4388             return (spa_vdev_exit(spa, newrootvd, txg, ENOTSUP));
4389
4390         pvops = &vdev_mirror_ops;
4391     } else {
4392         /*
4393          * Active hot spares can only be replaced by inactive hot
4394          * spares.
4395          */
4396         if (pvd->vdev_ops == &vdev_spare_ops &&
4397             oldvd->vdev_isspare &&
4398             !spa_has_spare(spa, newvd->vdev_guid))

```

```

4399         return (spa_vdev_exit(spa, newrootvd, txg, ENOTSUP));
4400
4401     /*
4402      * If the source is a hot spare, and the parent isn't already a
4403      * spare, then we want to create a new hot spare. Otherwise, we
4404      * want to create a replacing vdev. The user is not allowed to
4405      * attach to a spared vdev child unless the 'isspare' state is
4406      * the same (spare replaces spare, non-spare replaces
4407      * non-spare).
4408      */
4409     if (pvd->vdev_ops == &vdev_replacing_ops &&
4410         spa_version(spa) < SPA_VERSION_MULTI_REPLACE) {
4411         return (spa_vdev_exit(spa, newrootvd, txg, ENOTSUP));
4412     } else if (pvd->vdev_ops == &vdev_spare_ops &&
4413         newvd->vdev_isspare != oldvd->vdev_isspare) {
4414         return (spa_vdev_exit(spa, newrootvd, txg, ENOTSUP));
4415     }
4416
4417     if (newvd->vdev_isspare)
4418         pvops = &vdev_spare_ops;
4419     else
4420         pvops = &vdev_replacing_ops;
4421 }
4422
4423 /*
4424  * Make sure the new device is big enough.
4425  */
4426 if (newvd->vdev_asize < vdev_get_min_asize(oldvd))
4427     return (spa_vdev_exit(spa, newrootvd, txg, EOVERFLOW));
4428
4429 /*
4430  * The new device cannot have a higher alignment requirement
4431  * than the top-level vdev.
4432  */
4433 if (newvd->vdev_ashift > oldvd->vdev_top->vdev_ashift)
4434     return (spa_vdev_exit(spa, newrootvd, txg, EDOM));
4435
4436 /*
4437  * If this is an in-place replacement, update oldvd's path and devid
4438  * to make it distinguishable from newvd, and unopenable from now on.
4439  */
4440 if (strcmp(oldvd->vdev_path, newvd->vdev_path) == 0) {
4441     spa_strfree(oldvd->vdev_path);
4442     oldvd->vdev_path = kmem_alloc(strlen(newvd->vdev_path) + 5,
4443         KM_SLEEP);
4444     (void) sprintf(oldvd->vdev_path, "%s/%s",
4445         newvd->vdev_path, "old");
4446     if (oldvd->vdev_devid != NULL) {
4447         spa_strfree(oldvd->vdev_devid);
4448         oldvd->vdev_devid = NULL;
4449     }
4450 }
4451
4452 /* mark the device being resilvered */
4453 newvd->vdev_resilver_txg = txg;
4454
4455 /*
4456  * If the parent is not a mirror, or if we're replacing, insert the new
4457  * mirror/replacing/spare vdev above oldvd.
4458  */
4459 if (pvd->vdev_ops != pvops)
4460     pvd = vdev_add_parent(oldvd, pvops);
4461
4462 ASSERT(pvd->vdev_top->vdev_parent == rvd);
4463 ASSERT(pvd->vdev_ops == pvops);
4464 ASSERT(oldvd->vdev_parent == pvd);

```

```

4466     /*
4467     * Extract the new device from its root and add it to pvd.
4468     */
4469     vdev_remove_child(newrootvd, newvd);
4470     newvd->vdev_id = pvd->vdev_children;
4471     newvd->vdev_crtxg = oldvd->vdev_crtxg;
4472     vdev_add_child(pvd, newvd);

4474     tvd = newvd->vdev_top;
4475     ASSERT(pvd->vdev_top == tvd);
4476     ASSERT(tvd->vdev_parent == rvd);

4478     vdev_config_dirty(tvd);

4480     /*
4481     * Set newvd's DTL to [TXG_INITIAL, dtl_max_txg] so that we account
4482     * for any dmu_sync-ed blocks. It will propagate upward when
4483     * spa_vdev_exit() calls vdev_dtl_reassess().
4484     */
4485     dtl_max_txg = txg + TXG_CONCURRENT_STATES;

4487     vdev_dtl_dirty(newvd, DTL_MISSING, TXG_INITIAL,
4488                   dtl_max_txg - TXG_INITIAL);

4490     if (newvd->vdev_isspare) {
4491         spa_spare_activate(newvd);
4492         spa_event_notify(spa, newvd, ESC_ZFS_VDEV_SPARE);
4493     }

4495     oldvdpath = spa_strdup(oldvd->vdev_path);
4496     newvdpath = spa_strdup(newvd->vdev_path);
4497     newvd_isspare = newvd->vdev_isspare;

4499     /*
4500     * Mark newvd's DTL dirty in this txg.
4501     */
4502     vdev_dirty(tvd, VDD_DTL, newvd, txg);

4504     /*
4505     * Schedule the resilver to restart in the future. We do this to
4506     * ensure that dmu_sync-ed blocks have been stitched into the
4507     * respective datasets.
4508     * Restart the resilver
4509     */
4509     dsl_resilver_restart(spa->spa_dsl_pool, dtl_max_txg);

4511     /*
4512     * Commit the config
4513     */
4514     (void) spa_vdev_exit(spa, newrootvd, dtl_max_txg, 0);

4516     spa_history_log_internal(spa, "vdev attach", NULL,
4517                             "%s vdev=%s %s vdev=%s",
4518                             replacing && newvd_isspare ? "spare in" :
4519                             replacing ? "replace" : "attach", newvdpath,
4520                             replacing ? "for" : "to", oldvdpath);

4522     spa_strfree(oldvdpath);
4523     spa_strfree(newvdpath);

4525     if (spa->spa_bootfs)
4526         spa_event_notify(spa, newvd, ESC_ZFS_BOOTFS_VDEV_ATTACH);

4528     return (0);
4529 }

```

unchanged portion omitted

```

5093     /*
5094     * Evacuate the device.
5095     */
5096     static int
5097     spa_vdev_remove_evacuate(spa_t *spa, vdev_t *vd)
5098     {
5099         uint64_t txg;
5100         int error = 0;

5102         ASSERT(MUTEX_HELD(&spa_namespace_lock));
5103         ASSERT(spa_config_held(spa, SCL_ALL, RW_WRITER) == 0);
5104         ASSERT(vd == vd->vdev_top);

5106         /*
5107         * Evacuate the device. We don't hold the config lock as writer
5108         * since we need to do I/O but we do keep the
5109         * spa_namespace_lock held. Once this completes the device
5110         * should no longer have any blocks allocated on it.
5111         */
5112         if (vd->vdev_islog) {
5113             if (vd->vdev_stat.vs_alloc != 0)
5114                 error = spa_offline_log(spa);
5115         } else {
5116             error = SET_ERROR(ENOTSUP);
5117         }

5119         if (error)
5120             return (error);

5122         /*
5123         * The evacuation succeeded. Remove any remaining MOS metadata
5124         * associated with this vdev, and wait for these changes to sync.
5125         */
5126         ASSERT0(vd->vdev_stat.vs_alloc);
5127         txg = spa_vdev_config_enter(spa);
5128         vd->vdev_removing = B_TRUE;
5129         vdev_dirty_leaves(vd, VDD_DTL, txg);
5126         vdev_dirty(vd, 0, NULL, txg);
5130         vdev_config_dirty(vd);
5131         spa_vdev_config_exit(spa, NULL, txg, 0, FTAG);

5133         return (0);
5134     }

```

unchanged portion omitted

```

5900     /*
5901     * Set zpool properties.
5902     */
5903     static void
5904     spa_sync_props(void *arg, dmu_tx_t *tx)
5905     {
5906         nvlist_t *nvp = arg;
5907         spa_t *spa = dmu_tx_pool(tx)->dp_spa;
5908         objset_t *mos = spa->spa_meta_objset;
5909         nvpair_t *elem = NULL;

5911         mutex_enter(&spa->spa_props_lock);

5913         while ((elem = nvlist_next_nvpair(nvp, elem))) {
5914             uint64_t intval;
5915             char *strval, *fname;
5916             zpool_prop_t prop;
5917             const char *propname;
5918             zprop_type_t proptype;
5919             zfeature_info_t *feature;

```



```

5921     switch (prop = zpool_name_to_prop(nvpair_name(elem))) {
5922     case ZPROP_INVAL:
5923         /*
5924          * We checked this earlier in spa_prop_validate().
5925          */
5926         ASSERT(zpool_prop_feature(nvpair_name(elem)));

5928     fname = strchr(nvpair_name(elem), '@') + 1;
5929     VERIFY0(zfeature_lookup_name(fname, &feature));
5926     VERIFY3U(0, ==, zfeature_lookup_name(fname, &feature));

5931     spa_feature_enable(spa, feature, tx);
5932     spa_history_log_internal(spa, "set", tx,
5933         "%s=enabled", nvpair_name(elem));
5934     break;

5936     case ZPOOL_PROP_VERSION:
5937         intval = fnvpair_value_uint64(elem);
5934     VERIFY(nvpair_value_uint64(elem, &intval) == 0);
5938     /*
5939      * The version is synced seperatly before other
5940      * properties and should be correct by now.
5941      */
5942     ASSERT3U(spa_version(spa), >=, intval);
5943     break;

5945     case ZPOOL_PROP_ALTROOT:
5946     /*
5947      * 'altroot' is a non-persistent property. It should
5948      * have been set temporarily at creation or import time.
5949      */
5950     ASSERT(spa->spa_root != NULL);
5951     break;

5953     case ZPOOL_PROP_READONLY:
5954     case ZPOOL_PROP_CACHEFILE:
5955     /*
5956      * 'readonly' and 'cachefile' are also non-persisitent
5957      * properties.
5958      */
5959     break;

5960     case ZPOOL_PROP_COMMENT:
5961     strval = fnvpair_value_string(elem);
5958     VERIFY(nvpair_value_string(elem, &strval) == 0);
5962     if (spa->spa_comment != NULL)
5963         spa_strfree(spa->spa_comment);
5964     spa->spa_comment = spa_strdup(strval);
5965     /*
5966      * We need to dirty the configuration on all the vdevs
5967      * so that their labels get updated. It's unnecessary
5968      * to do this for pool creation since the vdev's
5969      * configuratoin has already been dirtied.
5970      */
5971     if (tx->tx_txg != TXG_INITIAL)
5972         vdev_config_dirty(spa->spa_root_vdev);
5973     spa_history_log_internal(spa, "set", tx,
5974         "%s=%s", nvpair_name(elem), strval);
5975     break;

5976     default:
5977     /*
5978      * Set pool property values in the poolprops mos object.
5979      */
5980     if (spa->spa_pool_props_object == 0) {
5981         spa->spa_pool_props_object =
5982         zap_create_link(mos, DMU_OT_POOL_PROPS,

```

```

5983         DMU_POOL_DIRECTORY_OBJECT, DMU_POOL_PROPS,
5984         tx);
5985     }

5987     /* normalize the property name */
5988     propname = zpool_prop_to_name(prop);
5989     proptype = zpool_prop_get_type(prop);

5991     if (nvpair_type(elem) == DATA_TYPE_STRING) {
5992         ASSERT(proptype == PROP_TYPE_STRING);
5993         strval = fnvpair_value_string(elem);
5994         VERIFY0(zap_update(mos,
5990             VERIFY(nvpair_value_string(elem, &strval) == 0);
5991             VERIFY(zap_update(mos,
5995                 spa->spa_pool_props_object, propname,
5996                 1, strlen(strval) + 1, strval, tx));
5993                 1, strlen(strval) + 1, strval, tx) == 0);
5997             spa_history_log_internal(spa, "set", tx,
5998                 "%s=%s", nvpair_name(elem), strval);
5999         } else if (nvpair_type(elem) == DATA_TYPE_UINT64) {
6000             intval = fnvpair_value_uint64(elem);
5997             VERIFY(nvpair_value_uint64(elem, &intval) == 0);

6002         if (proptype == PROP_TYPE_INDEX) {
6003             const char *unused;
6004             VERIFY0(zpool_prop_index_to_string(
6005                 prop, intval, &unused));
6001             VERIFY(zpool_prop_index_to_string(
6002                 prop, intval, &unused) == 0);
6006         }
6007         VERIFY0(zap_update(mos,
6004             VERIFY(zap_update(mos,
6008                 spa->spa_pool_props_object, propname,
6009                 8, 1, &intval, tx));
6006                 8, 1, &intval, tx) == 0);
6010             spa_history_log_internal(spa, "set", tx,
6011                 "%s=%lld", nvpair_name(elem), intval);
6012         } else {
6013             ASSERT(0); /* not allowed */
6014         }

6016     switch (prop) {
6017     case ZPOOL_PROP_DELEGATION:
6018         spa->spa_delegation = intval;
6019         break;
6020     case ZPOOL_PROP_BOOTFS:
6021         spa->spa_bootfs = intval;
6022         break;
6023     case ZPOOL_PROP_FAILUREMODE:
6024         spa->spa_failmode = intval;
6025         break;
6026     case ZPOOL_PROP_AUTOEXPAND:
6027         spa->spa_autoexpand = intval;
6028         if (tx->tx_txg != TXG_INITIAL)
6029             spa_async_request(spa,
6030                 SPA_ASYNC_AUTOEXPAND);
6031         break;
6032     case ZPOOL_PROP_DEDUPDITTO:
6033         spa->spa_dedup_ditto = intval;
6034         break;

6035     default:
6036         break;
6037     }

6040     }

```

new/usr/src/uts/common/fs/zfs/spa.c

9

```
6042     mutex_exit(&spa->spa_props_lock);
6043 }
_____unchanged_portion_omitted_
```

new/usr/src/uts/common/fs/zfs/spa_misc.c

1

```
*****
46220 Tue Sep 3 20:27:03 2013
new/usr/src/uts/common/fs/zfs/spa_misc.c
4101 metaslab_debug should allow for fine-grained control
4102 space_maps should store more information about themselves
4103 space_map object blocksize should be increased
4104 ::spa_space no longer works
4105 removing a mirrored log device results in a leaked object
4106 asynchronously load metaslab
Reviewed by: Matthew Ahrens <mahrens@delphix.com>
Reviewed by: Adam Leventhal <ahl@delphix.com>
Reviewed by: Sebastien Roy <seb@delphix.com>
*****
_____unchanged_portion_omitted_____

450 /*
451  * Fires when spa_sync has not completed within zfs_deadman_synctime_ms.
452  * If the zfs_deadman_enabled flag is set then it inspects all vdev queues
453  * looking for potentially hung I/Os.
454  */
455 void
456 spa_deadman(void *arg)
457 {
458     spa_t *spa = arg;

460     /*
461      * Disable the deadman timer if the pool is suspended.
462      */
463     if (spa_suspended(spa)) {
464         VERIFY(cyclic_reprogram(spa->spa_deadman_cycid, CY_INFINITY));
465         return;
466     }

468     zfs_dbgmsg("slow spa_sync: started %llu seconds ago, calls %llu",
469               (gethrtime() - spa->spa_sync_starttime) / NANOSEC,
470               ++spa->spa_deadman_calls);
471     if (zfs_deadman_enabled)
472         vdev_deadman(spa->spa_root_vdev);
473 }
_____unchanged_portion_omitted_____

986 /*
987  * Used in combination with spa_vdev_config_enter() to allow the syncing
988  * of multiple transactions without releasing the spa_namespace_lock.
989  */
990 void
991 spa_vdev_config_exit(spa_t *spa, vdev_t *vd, uint64_t txg, int error, char *tag)
992 {
993     ASSERT(MUTEX_HELD(&spa_namespace_lock));

995     int config_changed = B_FALSE;

997     ASSERT(txg > spa_last_synced_txg(spa));

999     spa->spa_pending_vdev = NULL;

1001     /*
1002      * Reassess the DTLs.
1003      */
1004     vdev_dtl_reassess(spa->spa_root_vdev, 0, 0, B_FALSE);

1006     if (error == 0 && !list_is_empty(&spa->spa_config_dirty_list)) {
1007         config_changed = B_TRUE;
1008         spa->spa_config_generation++;
1009     }

```

new/usr/src/uts/common/fs/zfs/spa_misc.c

2

```
1011     /*
1012      * Verify the metaslab classes.
1013      */
1014     ASSERT(metaslab_class_validate(spa_normal_class(spa)) == 0);
1015     ASSERT(metaslab_class_validate(spa_log_class(spa)) == 0);

1017     spa_config_exit(spa, SCL_ALL, spa);

1019     /*
1020      * Panic the system if the specified tag requires it. This
1021      * is useful for ensuring that configurations are updated
1022      * transactionally.
1023      */
1024     if (zio_injection_enabled)
1025         zio_handle_panic_injection(spa, tag, 0);

1027     /*
1028      * Note: this txg_wait_synced() is important because it ensures
1029      * that there won't be more than one config change per txg.
1030      * This allows us to use the txg as the generation number.
1031      */
1032     if (error == 0)
1033         txg_wait_synced(spa->spa_dsl_pool, txg);

1035     if (vd != NULL) {
1036         ASSERT(!vd->vdev_detached || vd->vdev_dtl_sm == NULL);
1037         ASSERT(!vd->vdev_detached || vd->vdev_dtl_smo.smo_object == 0);
1038         spa_config_enter(spa, SCL_ALL, spa, RW_WRITER);
1039         vdev_free(vd);
1040         spa_config_exit(spa, SCL_ALL, spa);
1041     }

1042     /*
1043      * If the config changed, update the config cache.
1044      */
1045     if (config_changed)
1046         spa_config_sync(spa, B_FALSE, B_TRUE);
1047 }
_____unchanged_portion_omitted_____

1678 void
1679 spa_init(int mode)
1680 {
1681     mutex_init(&spa_namespace_lock, NULL, MUTEX_DEFAULT, NULL);
1682     mutex_init(&spa_spare_lock, NULL, MUTEX_DEFAULT, NULL);
1683     mutex_init(&spa_l2cache_lock, NULL, MUTEX_DEFAULT, NULL);
1684     cv_init(&spa_namespace_cv, NULL, CV_DEFAULT, NULL);

1686     avl_create(&spa_namespace_avl, spa_name_compare, sizeof (spa_t),
1687              offsetof(spa_t, spa_avl));

1689     avl_create(&spa_spare_avl, spa_spare_compare, sizeof (spa_aux_t),
1690              offsetof(spa_aux_t, aux_avl));

1692     avl_create(&spa_l2cache_avl, spa_l2cache_compare, sizeof (spa_aux_t),
1693              offsetof(spa_aux_t, aux_avl));

1695     spa_mode_global = mode;

1697 #ifdef _KERNEL
1698     spa_arch_init();
1699 #else
1700     if (spa_mode_global != FREAD && dprintf_find_string("watch")) {
1701         arc_procfid = open("/proc/self/ctl", O_WRONLY);
1702         if (arc_procfid == -1) {
1703             perror("could not enable watchpoints: ")

```

```
1704             "opening /proc/self/ctl failed: ");
1705         } else {
1706             arc_watch = B_TRUE;
1707         }
1708     }
1709 #endif

1711     refcount_init();
1712     unique_init();
1713     range_tree_init();
1714     space_map_init();
1715     zio_init();
1716     dmuf_init();
1717     zil_init();
1718     vdev_cache_stat_init();
1719     zfs_prop_init();
1720     zpool_prop_init();
1721     zpool_feature_init();
1722     spa_config_load();
1723     l2arc_start();
1724 }

1725 void
1726 spa_fini(void)
1727 {
1728     l2arc_stop();

1730     spa_evict_all();

1732     vdev_cache_stat_fini();
1733     zil_fini();
1734     dmuf_fini();
1735     zio_fini();
1736     range_tree_fini();
1737     space_map_fini();
1738     unique_fini();
1739     refcount_fini();

1740     avl_destroy(&spa_namespace_avl);
1741     avl_destroy(&spa_spare_avl);
1742     avl_destroy(&spa_l2cache_avl);

1744     cv_destroy(&spa_namespace_cv);
1745     mutex_destroy(&spa_namespace_lock);
1746     mutex_destroy(&spa_spare_lock);
1747     mutex_destroy(&spa_l2cache_lock);
1748 }
unchanged_portion_omitted
```

new/usr/src/uts/common/fs/zfs/space_map.c

1

```
*****
15953 Tue Sep 3 20:27:04 2013
new/usr/src/uts/common/fs/zfs/space_map.c
4101 metaslab_debug should allow for fine-grained control
4102 space_maps should store more information about themselves
4103 space_map object blocksize should be increased
4104 ::spa_space no longer works
4105 removing a mirrored log device results in a leaked object
4106 asynchronously load metaslab
Reviewed by: Matthew Ahrens <mahrens@delphix.com>
Reviewed by: Adam Leventhal <ahl@delphix.com>
Reviewed by: Sebastien Roy <seb@delphix.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22 * Copyright 2009 Sun Microsystems, Inc. All rights reserved.
23 * Use is subject to license terms.
24 */
25 /*
26 * Copyright (c) 2013 by Delphix. All rights reserved.
27 * Copyright (c) 2012 by Delphix. All rights reserved.
28 */
29 #include <sys/zfs_context.h>
30 #include <sys/spa.h>
31 #include <sys/dmu.h>
32 #include <sys/dmu_tx.h>
33 #include <sys/dnode.h>
34 #include <sys/dsl_pool.h>
35 #include <sys/zio.h>
36 #include <sys/space_map.h>
37 #include <sys/refcount.h>
38 #include <sys/zfeature.h>
39
40 static kmem_cache_t *space_seg_cache;
41
42 void
43 space_map_init(void)
44 {
45     ASSERT(space_seg_cache == NULL);
46     space_seg_cache = kmem_cache_create("space_seg_cache",
47         sizeof(space_seg_t), 0, NULL, NULL, NULL, NULL, 0);
48 }
49
50 void
51 space_map_fini(void)
52 {
53 }
```

new/usr/src/uts/common/fs/zfs/space_map.c

2

```
48     kmem_cache_destroy(space_seg_cache);
49     space_seg_cache = NULL;
50 }
51
52 /*
53  * This value controls how the space map's block size is allowed to grow.
54  * If the value is set to the same size as SPACE_MAP_INITIAL_BLOCKSIZE then
55  * the space map block size will remain fixed. Setting this value to something
56  * greater than SPACE_MAP_INITIAL_BLOCKSIZE will allow the space map to
57  * increase its block size as needed. To maintain backwards compatibility the
58  * space map's block size must be a power of 2 and SPACE_MAP_INITIAL_BLOCKSIZE
59  * or larger.
60  * Space map routines.
61  * NOTE: caller is responsible for all locking.
62 */
63 int space_map_max_blkisz = (1 << 12);
64 static int
65 space_map_seg_compare(const void *x1, const void *x2)
66 {
67     const space_seg_t *s1 = x1;
68     const space_seg_t *s2 = x2;
69
70     if (s1->ss_start < s2->ss_start) {
71         if (s1->ss_end > s2->ss_start)
72             return (0);
73         return (-1);
74     }
75     if (s1->ss_start > s2->ss_start) {
76         if (s1->ss_start < s2->ss_end)
77             return (0);
78         return (1);
79     }
80     return (0);
81 }
82
83 void
84 space_map_create(space_map_t *sm, uint64_t start, uint64_t size, uint8_t shift,
85     kmutex_t *lp)
86 {
87     bzero(sm, sizeof(*sm));
88     cv_init(&sm->sm_load_cv, NULL, CV_DEFAULT, NULL);
89
90     avl_create(&sm->sm_root, space_map_seg_compare,
91         sizeof(space_seg_t), offsetof(struct space_seg, ss_node));
92
93     sm->sm_start = start;
94     sm->sm_size = size;
95     sm->sm_shift = shift;
96     sm->sm_lock = lp;
97 }
98
99 void
100 space_map_destroy(space_map_t *sm)
101 {
102     ASSERT(!sm->sm_loaded && !sm->sm_loading);
103     VERIFY0(sm->sm_space);
104     avl_destroy(&sm->sm_root);
105     cv_destroy(&sm->sm_load_cv);
106 }
107
108 void
109 space_map_add(space_map_t *sm, uint64_t start, uint64_t size)
110 {
111     avl_index_t where;
112     space_seg_t *ss_before, *ss_after, *ss;
113 }
```

```

106 uint64_t end = start + size;
107 int merge_before, merge_after;

109 ASSERT(MUTEX_HELD(sm->sm_lock));
110 VERIFY(!sm->sm_condensing);
111 VERIFY(size != 0);
112 VERIFY3U(start, >=, sm->sm_start);
113 VERIFY3U(end, <=, sm->sm_start + sm->sm_size);
114 VERIFY(sm->sm_space + size <= sm->sm_size);
115 VERIFY(P2PHASE(start, 1ULL << sm->sm_shift) == 0);
116 VERIFY(P2PHASE(size, 1ULL << sm->sm_shift) == 0);

118 ss = space_map_find(sm, start, size, &where);
119 if (ss != NULL) {
120     zfs_panic_recover("zfs: allocating allocated segment"
121         "(offset=%llu size=%llu)\n",
122         (longlong_t)start, (longlong_t)size);
123     return;
124 }

126 /* Make sure we don't overlap with either of our neighbors */
127 VERIFY(ss == NULL);

129 ss_before = avl_nearest(&sm->sm_root, where, AVL_BEFORE);
130 ss_after = avl_nearest(&sm->sm_root, where, AVL_AFTER);

132 merge_before = (ss_before != NULL && ss_before->ss_end == start);
133 merge_after = (ss_after != NULL && ss_after->ss_start == end);

135 if (merge_before && merge_after) {
136     avl_remove(&sm->sm_root, ss_before);
137     if (sm->sm_pp_root) {
138         avl_remove(sm->sm_pp_root, ss_before);
139         avl_remove(sm->sm_pp_root, ss_after);
140     }
141     ss_after->ss_start = ss_before->ss_start;
142     kmem_cache_free(space_seg_cache, ss_before);
143     ss = ss_after;
144 } else if (merge_before) {
145     ss_before->ss_end = end;
146     if (sm->sm_pp_root)
147         avl_remove(sm->sm_pp_root, ss_before);
148     ss = ss_before;
149 } else if (merge_after) {
150     ss_after->ss_start = start;
151     if (sm->sm_pp_root)
152         avl_remove(sm->sm_pp_root, ss_after);
153     ss = ss_after;
154 } else {
155     ss = kmem_cache_alloc(space_seg_cache, KM_SLEEP);
156     ss->ss_start = start;
157     ss->ss_end = end;
158     avl_insert(&sm->sm_root, ss, where);
159 }

161 if (sm->sm_pp_root)
162     avl_add(sm->sm_pp_root, ss);

164 sm->sm_space += size;
165 }

167 void
168 space_map_remove(space_map_t *sm, uint64_t start, uint64_t size)
169 {
170     avl_index_t where;
171     space_seg_t *ss, *newseg;

```

```

172 uint64_t end = start + size;
173 int left_over, right_over;

175 VERIFY(!sm->sm_condensing);
176 ss = space_map_find(sm, start, size, &where);

178 /* Make sure we completely overlap with someone */
179 if (ss == NULL) {
180     zfs_panic_recover("zfs: freeing free segment "
181         "(offset=%llu size=%llu)",
182         (longlong_t)start, (longlong_t)size);
183     return;
184 }
185 VERIFY3U(ss->ss_start, <=, start);
186 VERIFY3U(ss->ss_end, >=, end);
187 VERIFY(sm->sm_space - size <= sm->sm_size);

189 left_over = (ss->ss_start != start);
190 right_over = (ss->ss_end != end);

192 if (sm->sm_pp_root)
193     avl_remove(sm->sm_pp_root, ss);

195 if (left_over && right_over) {
196     newseg = kmem_cache_alloc(space_seg_cache, KM_SLEEP);
197     newseg->ss_start = end;
198     newseg->ss_end = ss->ss_end;
199     ss->ss_end = start;
200     avl_insert_here(&sm->sm_root, newseg, ss, AVL_AFTER);
201     if (sm->sm_pp_root)
202         avl_add(sm->sm_pp_root, newseg);
203 } else if (left_over) {
204     ss->ss_end = start;
205 } else if (right_over) {
206     ss->ss_start = end;
207 } else {
208     avl_remove(&sm->sm_root, ss);
209     kmem_cache_free(space_seg_cache, ss);
210     ss = NULL;
211 }

213 if (sm->sm_pp_root && ss != NULL)
214     avl_add(sm->sm_pp_root, ss);

216 sm->sm_space -= size;
217 }

219 space_seg_t *
220 space_map_find(space_map_t *sm, uint64_t start, uint64_t size,
221     avl_index_t *wherep)
222 {
223     space_seg_t ssearch, *ss;

225     ASSERT(MUTEX_HELD(sm->sm_lock));
226     VERIFY(size != 0);
227     VERIFY(P2PHASE(start, 1ULL << sm->sm_shift) == 0);
228     VERIFY(P2PHASE(size, 1ULL << sm->sm_shift) == 0);

230     ssearch.ss_start = start;
231     ssearch.ss_end = start + size;
232     ss = avl_find(&sm->sm_root, &ssearch, wherep);

234     if (ss != NULL && ss->ss_start <= start && ss->ss_end >= start + size)
235         return (ss);
236     return (NULL);
237 }

```

```

239 boolean_t
240 space_map_contains(space_map_t *sm, uint64_t start, uint64_t size)
241 {
242     avl_index_t where;
243
244     return (space_map_find(sm, start, size, &where) != 0);
245 }
246
247 void
248 space_map_swap(space_map_t **msrc, space_map_t **mdst)
249 {
250     space_map_t *sm;
251
252     ASSERT(MUTEX_HELD((*msrc)->sm_lock));
253     ASSERT0((*mdst)->sm_space);
254     ASSERT0(avl_numnodes(&(*mdst)->sm_root));
255
256     sm = *msrc;
257     *msrc = *mdst;
258     *mdst = sm;
259 }
260
261 void
262 space_map_vacate(space_map_t *sm, space_map_func_t *func, space_map_t *mdest)
263 {
264     space_seg_t *ss;
265     void *cookie = NULL;
266
267     ASSERT(MUTEX_HELD(sm->sm_lock));
268
269     while ((ss = avl_destroy_nodes(&sm->sm_root, &cookie)) != NULL) {
270         if (func != NULL)
271             func(mdest, ss->ss_start, ss->ss_end - ss->ss_start);
272         kmem_cache_free(space_seg_cache, ss);
273     }
274     sm->sm_space = 0;
275 }
276
277 void
278 space_map_walk(space_map_t *sm, space_map_func_t *func, space_map_t *mdest)
279 {
280     space_seg_t *ss;
281
282     ASSERT(MUTEX_HELD(sm->sm_lock));
283
284     for (ss = avl_first(&sm->sm_root); ss; ss = AVL_NEXT(&sm->sm_root, ss))
285         func(mdest, ss->ss_start, ss->ss_end - ss->ss_start);
286 }
287
288 /*
289 * Load the space map disk into the specified range tree. Segments of matype
290 * are added to the range tree, other segment types are removed.
291 *
292 * Wait for any in-progress space_map_load() to complete.
293 */
294 void
295 space_map_load_wait(space_map_t *sm)
296 {
297     ASSERT(MUTEX_HELD(sm->sm_lock));
298
299     while (sm->sm_loading) {
300         ASSERT(!sm->sm_loaded);
301         cv_wait(&sm->sm_load_cv, sm->sm_lock);
302     }

```

```

302 /*
303 * Note: space_map_load() will drop sm_lock across dmu_read() calls.
304 * The caller must be OK with this.
305 */
306 int
307 space_map_load(space_map_t *sm, range_tree_t *rt, matype_t matype)
308 space_map_load(space_map_t *sm, space_map_ops_t *ops, uint8_t matype,
309                 space_map_obj_t *smo, objset_t *os)
310 {
311     uint64_t *entry, *entry_map, *entry_map_end;
312     uint64_t bufsize, size, offset, end, space;
313     uint64_t mapstart = sm->sm_start;
314     int error = 0;
315
316     ASSERT(MUTEX_HELD(sm->sm_lock));
317     ASSERT(!sm->sm_loaded);
318     ASSERT(!sm->sm_loading);
319
320     end = space_map_length(sm);
321     space = space_map_allocated(sm);
322     sm->sm_loading = B_TRUE;
323     end = smo->smo_objsize;
324     space = smo->smo_alloc;
325
326     VERIFY0(range_tree_space(rt));
327     ASSERT(sm->sm_ops == NULL);
328     VERIFY0(sm->sm_space);
329
330     if (matype == SM_FREE) {
331         range_tree_add(rt, sm->sm_start, sm->sm_size);
332         space_map_add(sm, sm->sm_start, sm->sm_size);
333         space = sm->sm_size - space;
334     }
335
336     bufsize = MAX(sm->sm_blkisz, SPA_MINBLOCKSIZE);
337     bufsize = IULL << SPACE_MAP_BLOCKSHIFT;
338     entry_map = zio_buf_alloc(bufsize);
339
340     mutex_exit(sm->sm_lock);
341     if (end > bufsize) {
342         dmu_prefetch(sm->sm_os, space_map_object(sm), bufsize,
343                     end - bufsize);
344     }
345     if (end > bufsize)
346         dmu_prefetch(os, smo->smo_object, bufsize, end - bufsize);
347     mutex_enter(sm->sm_lock);
348
349     for (offset = 0; offset < end; offset += bufsize) {
350         size = MIN(end - offset, bufsize);
351         VERIFY(P2PHASE(size, sizeof(uint64_t)) == 0);
352         VERIFY(size != 0);
353         ASSERT3U(sm->sm_blkisz, !=, 0);
354
355         dprintf("object=%llu offset=%llx size=%llx\n",
356                space_map_object(sm), offset, size);
357         smo->smo_object, offset, size);
358
359         mutex_exit(sm->sm_lock);
360         error = dmu_read(sm->sm_os, space_map_object(sm), offset, size,
361                        entry_map, DMU_READ_PREFETCH);
362         error = dmu_read(os, smo->smo_object, offset, size, entry_map,
363                        DMU_READ_PREFETCH);
364         mutex_enter(sm->sm_lock);
365         if (error != 0)
366             break;

```

```

103     entry_map_end = entry_map + (size / sizeof (uint64_t));
104     for (entry = entry_map; entry < entry_map_end; entry++) {
105         uint64_t e = *entry;
106         uint64_t offset, size;

108         if (SM_DEBUG_DECODE(e))          /* Skip debug entries */
109             continue;

111         offset = (SM_OFFSET_DECODE(e) << sm->sm_shift) +
112                 sm->sm_start;
113         size = SM_RUN_DECODE(e) << sm->sm_shift;

115         VERIFY0(P2PHASE(offset, LULL << sm->sm_shift));
116         VERIFY0(P2PHASE(size, LULL << sm->sm_shift));
117         VERIFY3U(offset, >=, sm->sm_start);
118         VERIFY3U(offset + size, <=, sm->sm_start + sm->sm_size);
119         if (SM_TYPE_DECODE(e) == maptype) {
120             VERIFY3U(range_tree_space(rt) + size, <=,
121                     sm->sm_size);
122             range_tree_add(rt, offset, size);
123         } else {
124             range_tree_remove(rt, offset, size);
125         }
126     }
127 }

129 if (error == 0)
130     VERIFY3U(range_tree_space(rt), ==, space);
131 else
132     range_tree_vacate(rt, NULL, NULL);
133
134 if (error == 0) {
135     VERIFY3U(sm->sm_space, ==, space);

137     sm->sm_loaded = B_TRUE;
138     sm->sm_ops = ops;
139     if (ops != NULL)
140         ops->smop_load(sm);
141 } else {
142     space_map_vacate(sm, NULL, NULL);
143 }

144 zio_buf_free(entry_map, bufsize);
145 return (error);
146 }

147 void
148 space_map_histogram_clear(space_map_t *sm)
149 {
150     if (sm->sm_dbuf->db_size != sizeof (space_map_phys_t))
151         return;
152     sm->sm_loading = B_FALSE;

154     bzero(sm->sm_phys->smp_histogram, sizeof (sm->sm_phys->smp_histogram));
155 }
156 cv_broadcast(&sm->sm_load_cv);

157 boolean_t
158 space_map_histogram_verify(space_map_t *sm, range_tree_t *rt)
159 {
160     /*
161      * Verify that the in-core range tree does not have any

```

```

152     * ranges smaller than our sm_shift size.
153     */
154     for (int i = 0; i < sm->sm_shift; i++) {
155         if (rt->rt_histogram[i] != 0)
156             return (B_FALSE);
157     }
158     return (B_TRUE);
159 }
160 return (error);
161 }

162 void
163 space_map_histogram_add(space_map_t *sm, range_tree_t *rt, dmu_tx_t *tx)
164 space_map_unload(space_map_t *sm)
165 {
166     int idx = 0;
167     ASSERT(MUTEX_HELD(sm->sm_lock));

168     ASSERT(MUTEX_HELD(rt->rt_lock));
169     ASSERT(dmu_tx_is_syncing(tx));
170     VERIFY3U(space_map_object(sm), !=, 0);
171     if (sm->sm_loaded && sm->sm_ops != NULL)
172         sm->sm_ops->smop_unload(sm);

173     if (sm->sm_dbuf->db_size != sizeof (space_map_phys_t))
174         return;
175     sm->sm_loaded = B_FALSE;
176     sm->sm_ops = NULL;

177     dmu_buf_will_dirty(sm->sm_dbuf, tx);
178     space_map_vacate(sm, NULL, NULL);
179 }

180 ASSERT(space_map_histogram_verify(sm, rt));

181 /*
182  * Transfer the content of the range tree histogram to the space
183  * map histogram. The space map histogram contains 32 buckets ranging
184  * between 2^sm_shift to 2^(32+sm_shift-1). The range tree,
185  * however, can represent ranges from 2^0 to 2^63. Since the space
186  * map only cares about allocatable blocks (minimum of sm_shift) we
187  * can safely ignore all ranges in the range tree smaller than sm_shift.
188  */
189 for (int i = sm->sm_shift; i < RANGE_TREE_HISTOGRAM_SIZE; i++) {

190     /*
191      * Since the largest histogram bucket in the space map is
192      * 2^(32+sm_shift-1), we need to normalize the values in
193      * the range tree for any bucket larger than that size. For
194      * example given an sm_shift of 9, ranges larger than 2^40
195      * would get normalized as if they were 1TB ranges. Assume
196      * the range tree had a count of 5 in the 2^44 (16TB) bucket,
197      * the calculation below would normalize this to 5 * 2^4 (16).
198      */
199     ASSERT3U(i, >=, idx + sm->sm_shift);
200     sm->sm_phys->smp_histogram[idx] +=
201         rt->rt_histogram[i] << (i - idx - sm->sm_shift);

202     /*
203      * Increment the space map's index as long as we haven't
204      * reached the maximum bucket size. Accumulate all ranges
205      * larger than the max bucket size into the last bucket.
206      */
207     if (idx < SPACE_MAP_HISTOGRAM_SIZE(sm) - 1) {
208         ASSERT3U(idx + sm->sm_shift, ==, i);
209         idx++;
210         ASSERT3U(idx, <, SPACE_MAP_HISTOGRAM_SIZE(sm));

```



```

209     }
210 }
402 uint64_t
403 space_map_maxsize(space_map_t *sm)
404 {
405     ASSERT(sm->sm_ops != NULL);
406     return (sm->sm_ops->smop_max(sm));
211 }

213 uint64_t
214 space_map_entries(space_map_t *sm, range_tree_t *rt)
410 space_map_alloc(space_map_t *sm, uint64_t size)
215 {
216     avl_tree_t *t = &rt->rt_root;
217     range_seg_t *rs;
218     uint64_t size, entries;
412     uint64_t start;

220     /*
221      * All space_maps always have a debug entry so account for it here.
222      */
223     entries = 1;

225     /*
226      * Traverse the range tree and calculate the number of space map
227      * entries that would be required to write out the range tree.
228      */
229     for (rs = avl_first(t); rs != NULL; rs = AVL_NEXT(t, rs)) {
230         size = (rs->rs_end - rs->rs_start) >> sm->sm_shift;
231         entries += howmany(size, SM_RUN_MAX);
232     }
233     return (entries);
414     start = sm->sm_ops->smop_alloc(sm, size);
415     if (start != -1ULL)
416         space_map_remove(sm, start, size);
417     return (start);
234 }

236 void
237 space_map_set_blocksize(space_map_t *sm, uint64_t size, dmu_tx_t *tx)
421 space_map_claim(space_map_t *sm, uint64_t start, uint64_t size)
238 {
239     uint32_t blkksz;
240     u_longlong_t blocks;
423     sm->sm_ops->smop_claim(sm, start, size);
424     space_map_remove(sm, start, size);
425 }

242     ASSERT3U(sm->sm_blkksz, !=, 0);
243     ASSERT3U(space_map_object(sm), !=, 0);
244     ASSERT(sm->sm_dbuf != NULL);
245     VERIFY(ISP2(space_map_max_blkksz));

247     if (sm->sm_blkksz >= space_map_max_blkksz)
248         return;

250     /*
251      * The object contains more than one block so we can't adjust
252      * its size.
253      */
254     if (sm->sm_phys->smp_objsize > sm->sm_blkksz)
255         return;

257     if (size > sm->sm_blkksz) {
258         uint64_t newsz;

```

```

260     /*
261      * Older software versions treat space map blocks as fixed
262      * entities. The DMU is capable of handling different block
263      * sizes making it possible for us to increase the
264      * block size and maintain backwards compatibility. The
265      * caveat is that the new block sizes must be a
266      * power of 2 so that old software can append to the file,
267      * adding more blocks. The block size can grow until it
268      * reaches space_map_max_blkksz.
269      */
270     newsz = ISP2(size) ? size : 1ULL << highbit(size);
271     if (newsz > space_map_max_blkksz)
272         newsz = space_map_max_blkksz;

274     VERIFY0(dmu_object_set_blocksize(sm->sm_os,
275         space_map_object(sm), newsz, 0, tx));
276     dmu_object_size_from_db(sm->sm_dbuf, &blkksz, &blocks);

278     zfs_dbgmsg("txg %llu, spa %s, increasing blkksz from %d to %d",
279         dmu_tx_get_txg(tx), spa_name(dmu_objset_spa(sm->sm_os)),
280         sm->sm_blkksz, blkksz);

282     VERIFY3U(newsz, ==, blkksz);
283     VERIFY3U(sm->sm_blkksz, <, blkksz);
284     sm->sm_blkksz = blkksz;
285 }
427 void
428 space_map_free(space_map_t *sm, uint64_t start, uint64_t size)
429 {
430     space_map_add(sm, start, size);
431     sm->sm_ops->smop_free(sm, start, size);
286 }

288 /*
289  * Note: space_map_write() will drop sm_lock across dmu_write() calls.
435  * Note: space_map_sync() will drop sm_lock across dmu_write() calls.
290  */
291 void
292 space_map_write(space_map_t *sm, range_tree_t *rt, mactype_t mactype,
293     dmu_tx_t *tx)
438 space_map_sync(space_map_t *sm, uint8_t mactype,
439     space_map_obj_t *smo, objset_t *os, dmu_tx_t *tx)
294 {
295     objset_t *os = sm->sm_os;
296     spa_t *spa = dmu_objset_spa(os);
297     avl_tree_t *t = &rt->rt_root;
298     range_seg_t *rs;
299     uint64_t size, total, rt_space, nodes;
442     avl_tree_t *t = &sm->sm_root;
443     space_seg_t *ss;
444     uint64_t bufsize, start, size, run_len, total, sm_space, nodes;
300     uint64_t *entry, *entry_map, *entry_map_end;
301     uint64_t newsz, expected_entries, actual_entries = 1;

303     ASSERT(MUTEX_HELD(rt->rt_lock));
304     ASSERT(dsl_pool_sync_context(dmu_objset_pool(os));
305     VERIFY3U(space_map_object(sm), !=, 0);
306     dmu_buf_will_dirty(sm->sm_dbuf, tx);
447     ASSERT(MUTEX_HELD(sm->sm_lock));

308     /*
309      * This field is no longer necessary since the in-core space map
310      * now contains the object number but is maintained for backwards
311      * compatibility.
312      */
313     sm->sm_phys->smp_object = sm->sm_object;

```

```

315     if (range_tree_space(rt) == 0) {
316         VERIFY3U(sm->sm_object, ==, sm->sm_phys->smp_object);
449     if (sm->sm_space == 0)
317         return;
318     }

452     dprintf("object %4llu, txg %llu, pass %d, %c, count %lu, space %llx\n",
453            smo->smo_object, dmu_tx_get_txg(tx), spa_sync_pass(spa),
454            maptype == SM_ALLOC ? 'A' : 'F', avl_numnodes(&sm->sm_root),
455            sm->sm_space);

320     if (maptype == SM_ALLOC)
321         sm->sm_phys->smp_alloc += range_tree_space(rt);
458         smo->smo_alloc += sm->sm_space;
322     else
323         sm->sm_phys->smp_alloc -= range_tree_space(rt);
460         smo->smo_alloc -= sm->sm_space;

325     expected_entries = space_map_entries(sm, rt);

327     /*
328     * Calculate the new size for the space map on-disk and see if
329     * we can grow the block size to accommodate the new size.
330     */
331     newsz = sm->sm_phys->smp_objsize + expected_entries * sizeof (uint64_t);
332     space_map_set_blocksize(sm, newsz, tx);

334     entry_map = zio_buf_alloc(sm->sm_blksize);
335     entry_map_end = entry_map + (sm->sm_blksize / sizeof (uint64_t));
462     bufsize = (8 + avl_numnodes(&sm->sm_root)) * sizeof (uint64_t);
463     bufsize = MIN(bufsize, 1ULL << SPACE_MAP_BLOCKSHIFT);
464     entry_map = zio_buf_alloc(bufsize);
465     entry_map_end = entry_map + (bufsize / sizeof (uint64_t));
336     entry = entry_map;

338     *entry++ = SM_DEBUG_ENCODE(1) |
339             SM_DEBUG_ACTION_ENCODE(maptype) |
340             SM_DEBUG_SYNCPASS_ENCODE(spa_sync_pass(spa)) |
341             SM_DEBUG_TXG_ENCODE(dmu_tx_get_txg(tx));

343     total = 0;
344     nodes = avl_numnodes(&rt->rt_root);
345     rt_space = range_tree_space(rt);
346     for (rs = avl_first(t); rs != NULL; rs = AVL_NEXT(t, rs)) {
347         uint64_t start;
474     nodes = avl_numnodes(&sm->sm_root);
475     sm_space = sm->sm_space;
476     for (ss = avl_first(t); ss != NULL; ss = AVL_NEXT(t, ss)) {
477         size = ss->ss_end - ss->ss_start;
478         start = (ss->ss_start - sm->sm_start) >> sm->sm_shift;

349         size = (rs->rs_end - rs->rs_start) >> sm->sm_shift;
350         start = (rs->rs_start - sm->sm_start) >> sm->sm_shift;
480         total += size;
481         size >>= sm->sm_shift;

352     total += size << sm->sm_shift;

354     while (size != 0) {
355         uint64_t run_len;

483     while (size) {
357         run_len = MIN(size, SM_RUN_MAX);

359         if (entry == entry_map_end) {

```

```

360         mutex_exit(rt->rt_lock);
361         dmu_write(os, space_map_object(sm),
362                sm->sm_phys->smp_objsize, sm->sm_blksize,
363                entry_map, tx);
364         mutex_enter(rt->rt_lock);
365         sm->sm_phys->smp_objsize += sm->sm_blksize;
487         mutex_exit(sm->sm_lock);
488         dmu_write(os, smo->smo_object, smo->smo_objsize,
489                bufsize, entry_map, tx);
490         mutex_enter(sm->sm_lock);
491         smo->smo_objsize += bufsize;
366         entry = entry_map;
367     }

369     *entry++ = SM_OFFSET_ENCODE(start) |
370             SM_TYPE_ENCODE(maptype) |
371             SM_RUN_ENCODE(run_len);

373     start += run_len;
374     size -= run_len;
375     actual_entries++;
376     }
377 }

379     if (entry != entry_map) {
380         size = (entry - entry_map) * sizeof (uint64_t);
381         mutex_exit(rt->rt_lock);
382         dmu_write(os, space_map_object(sm), sm->sm_phys->smp_objsize,
506         mutex_exit(sm->sm_lock);
507         dmu_write(os, smo->smo_object, smo->smo_objsize,
383                 size, entry_map, tx);
384         mutex_enter(rt->rt_lock);
385         sm->sm_phys->smp_objsize += size;
509         mutex_enter(sm->sm_lock);
510         smo->smo_objsize += size;
386     }
387     ASSERT3U(expected_entries, ==, actual_entries);

389     /*
390     * Ensure that the space_map's accounting wasn't changed
391     * while we were in the middle of writing it out.
392     */
393     VERIFY3U(nodes, ==, avl_numnodes(&rt->rt_root));
394     VERIFY3U(range_tree_space(rt), ==, rt_space);
395     VERIFY3U(range_tree_space(rt), ==, total);
517     VERIFY3U(nodes, ==, avl_numnodes(&sm->sm_root));
518     VERIFY3U(smo->smo_objsize, ==, sm_space);
519     VERIFY3U(smo->smo_objsize, ==, total);

397     zio_buf_free(entry_map, sm->sm_blksize);
521     zio_buf_free(entry_map, bufsize);
398 }

400 static int
401 space_map_open_impl(space_map_t *sm)
524 void
525 space_map_truncate(space_map_obj_t *smo, objset_t *os, dmu_tx_t *tx)
402 {
403     int error;
404     u_longlong_t blocks;
527     VERIFY(dmu_free_range(os, smo->smo_object, 0, -1ULL, tx) == 0);

406     error = dmu_bonus_hold(sm->sm_os, sm->sm_object, sm, &sm->sm_dbuf);
407     if (error)
408         return (error);

```

```

410     dmubuf_size_from_db(sm->sm_dbuf, &sm->sm_blksize, &blocks);
411     sm->sm_phys = sm->sm_dbuf->db_data;
412     return (0);
413 }
414
415 int
416 space_map_open(space_map_t **smp, objset_t *os, uint64_t object,
417               uint64_t start, uint64_t size, uint8_t shift, kmutex_t *lp)
418 /*
419  * Space map reference trees.
420  *
421  * A space map is a collection of integers. Every integer is either
422  * in the map, or it's not. A space map reference tree generalizes
423  * the idea: it allows its members to have arbitrary reference counts,
424  * as opposed to the implicit reference count of 0 or 1 in a space map.
425  * This representation comes in handy when computing the union or
426  * intersection of multiple space maps. For example, the union of
427  * N space maps is the subset of the reference tree with refcnt >= 1.
428  * The intersection of N space maps is the subset with refcnt >= N.
429  *
430  * [It's very much like a Fourier transform. Unions and intersections
431  * are hard to perform in the 'space map domain', so we convert the maps
432  * into the 'reference count domain', where it's trivial, then invert.]
433  *
434  * vdev_dtl_reassess() uses computations of this form to determine
435  * DTL_MISSING and DTL_OUTAGE for interior vdevs -- e.g. a RAID-Z vdev
436  * has an outage wherever refcnt >= vdev_nparity + 1, and a mirror vdev
437  * has an outage wherever refcnt >= vdev_children.
438  */
439 static int
440 space_map_ref_compare(const void *x1, const void *x2)
441 {
442     space_map_t *sm;
443     int error;
444     const space_ref_t *sr1 = x1;
445     const space_ref_t *sr2 = x2;
446
447     ASSERT(*smp == NULL);
448     ASSERT(os != NULL);
449     ASSERT(object != 0);
450     if (sr1->sr_offset < sr2->sr_offset)
451         return (-1);
452     if (sr1->sr_offset > sr2->sr_offset)
453         return (1);
454
455     sm = kmem_zalloc(sizeof (space_map_t), KM_SLEEP);
456     if (sr1 < sr2)
457         return (-1);
458     if (sr1 > sr2)
459         return (1);
460
461     sm->sm_start = start;
462     sm->sm_size = size;
463     sm->sm_shift = shift;
464     sm->sm_lock = lp;
465     sm->sm_os = os;
466     sm->sm_object = object;
467
468     error = space_map_open_impl(sm);
469     if (error != 0) {
470         space_map_close(sm);
471         return (error);
472     }
473 }

```

```

441     *smp = sm;
442
443     return (0);
444 }
445
446 void
447 space_map_close(space_map_t *sm)
448 {
449     space_map_ref_create(avl_tree_t *t)
450     {
451         if (sm == NULL)
452             return;
453         avl_create(t, space_map_ref_compare,
454                 sizeof (space_ref_t), offsetof(space_ref_t, sr_node));
455     }
456
457     if (sm->sm_dbuf != NULL)
458         dmubuf_rele(sm->sm_dbuf, sm);
459     sm->sm_dbuf = NULL;
460     sm->sm_phys = NULL;
461
462     void
463     space_map_ref_destroy(avl_tree_t *t)
464     {
465         space_ref_t *sr;
466         void *cookie = NULL;
467
468         kmem_free(sm, sizeof (*sm));
469         while ((sr = avl_destroy_nodes(t, &cookie)) != NULL)
470             kmem_free(sr, sizeof (*sr));
471
472         avl_destroy(t);
473     }
474
475     static void
476     space_map_reallocate(space_map_t *sm, dmubuf_t *tx)
477     {
478         space_map_ref_add_node(avl_tree_t *t, uint64_t offset, int64_t refcnt)
479         {
480             ASSERT(dmubuf_is_syncing(tx));
481             space_ref_t *sr;
482
483             space_map_free(sm, tx);
484             dmubuf_rele(sm->sm_dbuf, sm);
485             sr = kmem_alloc(sizeof (*sr), KM_SLEEP);
486             sr->sr_offset = offset;
487             sr->sr_refcnt = refcnt;
488
489             sm->sm_object = space_map_alloc(sm->sm_os, tx);
490             VERIFY0(space_map_open_impl(sm));
491             avl_add(t, sr);
492         }
493
494     void
495     space_map_truncate(space_map_t *sm, dmubuf_t *tx)
496     {
497         space_map_ref_add_seg(avl_tree_t *t, uint64_t start, uint64_t end,
498                             int64_t refcnt)
499         {
500             objset_t *os = sm->sm_os;
501             spa_t *spa = dmubuf_objset_spa(os);
502             zfeature_info_t *space_map_histogram =
503                 &spa_feature_table[SPA_FEATURE_SPACEMAP_HISTOGRAM];
504             dmubuf_objset_info_t doi;
505             int bonuslen;
506
507             ASSERT(dsl_pool_sync_context(dmubuf_objset_pool(os));
508                 ASSERT(dmubuf_is_syncing(tx));
509
510             VERIFY0(dmubuf_free_range(os, space_map_object(sm), 0, -1ULL, tx));

```

```

486     dmu_object_info_from_db(sm->sm_dbuf, &doi);
488     if (spa_feature_is_enabled(spa, space_map_histogram)) {
489         bonuslen = sizeof (space_map_phys_t);
490         ASSERT3U(bonuslen, <=, dmu_bonus_max());
491     } else {
492         bonuslen = SPACE_MAP_SIZE_V0;
493     }
495     if (bonuslen != doi.doi_bonus_size ||
496         doi.doi_data_block_size != SPACE_MAP_INITIAL_BLOCKSIZE) {
497         zfs_dbgmsg("txg %llu, spa %s, reallocating: "
498             "old bonus %u, old blocksz %u", dmu_tx_get_txg(tx),
499             spa_name(spa), doi.doi_bonus_size, doi.doi_data_block_size);
500         space_map_reallocate(sm, tx);
501         VERIFY3U(sm->sm_blkksz, ==, SPACE_MAP_INITIAL_BLOCKSIZE);
502     }
504     dmu_buf_will_dirty(sm->sm_dbuf, tx);
505     sm->sm_phys->smp_objsize = 0;
506     sm->sm_phys->smp_alloc = 0;
507     space_map_ref_add_node(t, start, refcnt);
508     space_map_ref_add_node(t, end, -refcnt);
509 }
510 /*
511 * Update the in-core space_map allocation and length values.
512 * Convert (or add) a space map into a reference tree.
513 */
514 void
515 space_map_update(space_map_t *sm)
516 {
517     space_map_ref_add_map(avl_tree_t *t, space_map_t *sm, int64_t refcnt)
518     {
519         if (sm == NULL)
520             return;
521         space_seg_t *ss;
522     }
523     ASSERT(MUTEX_HELD(sm->sm_lock));
524     sm->sm_alloc = sm->sm_phys->smp_alloc;
525     sm->sm_length = sm->sm_phys->smp_objsize;
526     for (ss = avl_first(&sm->sm_root); ss; ss = AVL_NEXT(&sm->sm_root, ss))
527         space_map_ref_add_seg(t, ss->ss_start, ss->ss_end, refcnt);
528 }
529
530 uint64_t
531 space_map_alloc(objset_t *os, dmu_tx_t *tx)
532 {
533     spa_t *spa = dmu_objset_spa(os);
534     zfeature_info_t *space_map_histogram =
535         &spa_feature_table[SPA_FEATURE_SPACEMAP_HISTOGRAM];
536     uint64_t object;
537     int bonuslen;
538
539     if (spa_feature_is_enabled(spa, space_map_histogram)) {
540         spa_feature_incr(spa, space_map_histogram, tx);
541         bonuslen = sizeof (space_map_phys_t);
542         ASSERT3U(bonuslen, <=, dmu_bonus_max());
543     } else {
544         bonuslen = SPACE_MAP_SIZE_V0;
545     }
546
547     object = dmu_object_alloc(os,
548         DMU_OT_SPACE_MAP, SPACE_MAP_INITIAL_BLOCKSIZE,
549         DMU_OT_SPACE_MAP_HEADER, bonuslen, tx);

```

```

545         return (object);
546     }
547
548 /*
549 * Convert a reference tree into a space map. The space map will contain
550 * all members of the reference tree for which refcnt >= minref.
551 */
552 void
553 space_map_free(space_map_t *sm, dmu_tx_t *tx)
554 {
555     space_map_ref_generate_map(avl_tree_t *t, space_map_t *sm, int64_t minref)
556     {
557         spa_t *spa;
558         zfeature_info_t *space_map_histogram =
559             &spa_feature_table[SPA_FEATURE_SPACEMAP_HISTOGRAM];
560         uint64_t start = -1ULL;
561         int64_t refcnt = 0;
562         space_ref_t *sr;
563
564         if (sm == NULL)
565             return;
566         ASSERT(MUTEX_HELD(sm->sm_lock));
567
568         spa = dmu_objset_spa(sm->sm_os);
569         if (spa_feature_is_enabled(spa, space_map_histogram)) {
570             dmu_object_info_t doi;
571             space_map_vacate(sm, NULL, NULL);
572
573             dmu_object_info_from_db(sm->sm_dbuf, &doi);
574             if (doi.doi_bonus_size != SPACE_MAP_SIZE_V0) {
575                 VERIFY(spa_feature_is_active(spa, space_map_histogram));
576                 spa_feature_decr(spa, space_map_histogram, tx);
577             }
578             for (sr = avl_first(t); sr != NULL; sr = AVL_NEXT(t, sr)) {
579                 refcnt += sr->sr_refcnt;
580                 if (refcnt >= minref) {
581                     if (start == -1ULL) {
582                         start = sr->sr_offset;
583                     } else {
584                         if (start != -1ULL) {
585                             uint64_t end = sr->sr_offset;
586                             ASSERT(start <= end);
587                             if (end > start)
588                                 space_map_add(sm, start, end - start);
589                             start = -1ULL;
590                         }
591                     }
592                 }
593             }
594             VERIFY3U(dmu_object_free(sm->sm_os, space_map_object(sm), tx), ==, 0);
595             sm->sm_object = 0;
596         }
597     }
598
599 uint64_t
600 space_map_object(space_map_t *sm)
601 {
602     return (sm != NULL ? sm->sm_object : 0);
603 }
604
605 /*
606 * Returns the already synced, on-disk allocated space.
607 */
608 uint64_t
609 space_map_allocated(space_map_t *sm)
610 {
611     return (sm != NULL ? sm->sm_alloc : 0);
612 }
613
614 /*

```

```
589 * Returns the already synced, on-disk length;
590 */
591 uint64_t
592 space_map_length(space_map_t *sm)
593 {
594     return (sm != NULL ? sm->sm_length : 0);
595 }

597 /*
598 * Returns the allocated space that is currently syncing.
599 */
600 int64_t
601 space_map_alloc_delta(space_map_t *sm)
602 {
603     if (sm == NULL)
604         return (0);
605     ASSERT(sm->sm_dbuf != NULL);
606     return (sm->sm_phys->smp_alloc - space_map_allocated(sm));
607 }
608     ASSERT(RefCount == 0);
609     ASSERT(start == -1ULL);
610 }
611     unchanged_portion_omitted_
```

```

*****
4308 Tue Sep 3 20:27:05 2013
new/usr/src/uts/common/fs/zfs/space_reftree.c
4101 metaslab_debug should allow for fine-grained control
4102 space_maps should store more information about themselves
4103 space_map object blocksize should be increased
4104 ::spa_space no longer works
4105 removing a mirrored log device results in a leaked object
4106 asynchronously load metaslab
Reviewed by: Matthew Ahrens <mahrens@delphix.com>
Reviewed by: Adam Leventhal <ahl@delphix.com>
Reviewed by: Sebastien Roy <seb@delphix.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22 * Copyright 2009 Sun Microsystems, Inc. All rights reserved.
23 * Use is subject to license terms.
24 */
25 /*
26 * Copyright (c) 2013 by Delphix. All rights reserved.
27 */
28
29 #include <sys/zfs_context.h>
30 #include <sys/range_tree.h>
31 #include <sys/space_reftree.h>
32
33 /*
34 * Space reference trees.
35 *
36 * A range tree is a collection of integers. Every integer is either
37 * in the tree, or it's not. A space reference tree generalizes
38 * the idea: it allows its members to have arbitrary reference counts,
39 * as opposed to the implicit reference count of 0 or 1 in a range tree.
40 * This representation comes in handy when computing the union or
41 * intersection of multiple space maps. For example, the union of
42 * N range trees is the subset of the reference tree with refcnt >= 1.
43 * The intersection of N range trees is the subset with refcnt >= N.
44 *
45 * [It's very much like a Fourier transform. Unions and intersections
46 * are hard to perform in the 'range tree domain', so we convert the trees
47 * into the 'reference count domain', where it's trivial, then invert.]
48 *
49 * vdev_dtl_reassess() uses computations of this form to determine
50 * DTL_MISSING and DTL_OUTAGE for interior vdevs -- e.g. a RAID-Z vdev
51 * has an outage wherever refcnt >= vdev_nparity + 1, and a mirror vdev
52 * has an outage wherever refcnt >= vdev_children.
53 */

```

```

54 static int
55 space_reftree_compare(const void *x1, const void *x2)
56 {
57     const space_ref_t *sr1 = x1;
58     const space_ref_t *sr2 = x2;
59
60     if (sr1->sr_offset < sr2->sr_offset)
61         return (-1);
62     if (sr1->sr_offset > sr2->sr_offset)
63         return (1);
64
65     if (sr1 < sr2)
66         return (-1);
67     if (sr1 > sr2)
68         return (1);
69
70     return (0);
71 }
72
73 void
74 space_reftree_create(avl_tree_t *t)
75 {
76     avl_create(t, space_reftree_compare,
77             sizeof (space_ref_t), offsetof(space_ref_t, sr_node));
78 }
79
80 void
81 space_reftree_destroy(avl_tree_t *t)
82 {
83     space_ref_t *sr;
84     void *cookie = NULL;
85
86     while ((sr = avl_destroy_nodes(t, &cookie)) != NULL)
87         kmem_free(sr, sizeof (*sr));
88
89     avl_destroy(t);
90 }
91
92 static void
93 space_reftree_add_node(avl_tree_t *t, uint64_t offset, int64_t refcnt)
94 {
95     space_ref_t *sr;
96
97     sr = kmem_alloc(sizeof (*sr), KM_SLEEP);
98     sr->sr_offset = offset;
99     sr->sr_refcnt = refcnt;
100
101     avl_add(t, sr);
102 }
103
104 void
105 space_reftree_add_seg(avl_tree_t *t, uint64_t start, uint64_t end,
106                     int64_t refcnt)
107 {
108     space_reftree_add_node(t, start, refcnt);
109     space_reftree_add_node(t, end, -refcnt);
110 }
111
112 /*
113  * Convert (or add) a range tree into a reference tree.
114  */
115 void
116 space_reftree_add_map(avl_tree_t *t, range_tree_t *rt, int64_t refcnt)
117 {
118     range_seg_t *rs;

```

```
120     ASSERT(MUTEX_HELD(rt->rt_lock));
121
122     for (rs = avl_first(&rt->rt_root); rs; rs = AVL_NEXT(&rt->rt_root, rs))
123         space_reftree_add_seg(t, rs->rs_start, rs->rs_end, refcnt);
124 }
125
126 /*
127  * Convert a reference tree into a range tree. The range tree will contain
128  * all members of the reference tree for which refcnt >= minref.
129  */
130 void
131 space_reftree_generate_map(avl_tree_t *t, range_tree_t *rt, int64_t minref)
132 {
133     uint64_t start = -1ULL;
134     int64_t refcnt = 0;
135     space_ref_t *sr;
136
137     ASSERT(MUTEX_HELD(rt->rt_lock));
138
139     range_tree_vacate(rt, NULL, NULL);
140
141     for (sr = avl_first(t); sr != NULL; sr = AVL_NEXT(t, sr)) {
142         refcnt += sr->sr_refcnt;
143         if (refcnt >= minref) {
144             if (start == -1ULL) {
145                 start = sr->sr_offset;
146             }
147         } else {
148             if (start != -1ULL) {
149                 uint64_t end = sr->sr_offset;
150                 ASSERT(start <= end);
151                 if (end > start)
152                     range_tree_add(rt, start, end - start);
153                 start = -1ULL;
154             }
155         }
156     }
157     ASSERT(refcnt == 0);
158     ASSERT(start == -1ULL);
159 }
```

```

*****
3187 Tue Sep 3 20:27:06 2013
new/usr/src/uts/common/fs/zfs/sys/metaslabs.h
4101 metaslab_debug should allow for fine-grained control
4102 space_maps should store more information about themselves
4103 space_map object blocksize should be increased
4104 ::spa_space no longer works
4105 removing a mirrored log device results in a leaked object
4106 asynchronously load metaslab
Reviewed by: Matthew Ahrens <mahrens@delphix.com>
Reviewed by: Adam Leventhal <ahl@delphix.com>
Reviewed by: Sebastien Roy <seb@delphix.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22 * Copyright (c) 2005, 2010, Oracle and/or its affiliates. All rights reserved.
23 * Copyright (c) 2013 by Delphix. All rights reserved.
24 * Copyright (c) 2012 by Delphix. All rights reserved.
25 */

26 #ifndef _SYS_METASLAB_H
27 #define _SYS_METASLAB_H

28
29 #include <sys/spa.h>
30 #include <sys/space_map.h>
31 #include <sys/txg.h>
32 #include <sys/zio.h>
33 #include <sys/avl.h>

34
35 #ifdef __cplusplus
36 extern "C" {
37 #endif

38
39 typedef struct metaslab_ops {
40     uint64_t (*msop_alloc)(metaslab_t *msp, uint64_t size);
41     boolean_t (*msop_fragmented)(metaslab_t *msp);
42 } metaslab_ops_t;
43 extern space_map_ops_t *zfs_metaslabs_ops;

44 extern metaslab_ops_t *zfs_metaslabs_ops;
45 extern metaslab_t *metaslabs_init(metaslab_group_t *mg, space_map_obj_t *smo,
46     uint64_t start, uint64_t size, uint64_t txg);
47 extern void metaslab_fini(metaslab_t *msp);
48 extern void metaslab_sync(metaslab_t *msp, uint64_t txg);
49 extern void metaslab_sync_done(metaslab_t *msp, uint64_t txg);
50 extern void metaslab_sync_reassess(metaslab_group_t *mg);

```

```

46 metaslab_t *metaslabs_init(metaslab_group_t *mg, uint64_t id,
47     uint64_t object, uint64_t txg);
48 void metaslab_fini(metaslab_t *msp);

49
50 void metaslab_load_wait(metaslab_t *msp);
51 int metaslab_load(metaslab_t *msp);
52 void metaslab_unload(metaslab_t *msp);

53
54 void metaslab_sync(metaslab_t *msp, uint64_t txg);
55 void metaslab_sync_done(metaslab_t *msp, uint64_t txg);
56 void metaslab_sync_reassess(metaslab_group_t *mg);
57 uint64_t metaslab_block_maxsize(metaslab_t *msp);

58
59 #define METASLAB_HINTBP_FAVOR    0x0
60 #define METASLAB_HINTBP_AVOID    0x1
61 #define METASLAB_GANG_HEADER    0x2
62 #define METASLAB_GANG_CHILD    0x4
63 #define METASLAB_GANG_AVOID    0x8

64
65 int metaslab_alloc(spa_t *spa, metaslab_class_t *mc, uint64_t psize,
66     blkptr_t *bp, int ncopies, uint64_t txg, blkptr_t *hintbp, int flags);
67 void metaslab_free(spa_t *spa, const blkptr_t *bp, uint64_t txg, boolean_t now);
68 int metaslab_claim(spa_t *spa, const blkptr_t *bp, uint64_t txg);
69 void metaslab_check_free(spa_t *spa, const blkptr_t *bp);
70 extern void metaslab_free(spa_t *spa, const blkptr_t *bp, uint64_t txg,
71     boolean_t now);
72 extern int metaslab_claim(spa_t *spa, const blkptr_t *bp, uint64_t txg);
73 extern void metaslab_check_free(spa_t *spa, const blkptr_t *bp);

74
75 metaslab_class_t *metaslabs_class_create(spa_t *spa, metaslab_ops_t *ops);
76 void metaslab_class_destroy(metaslab_class_t *mc);
77 int metaslab_class_validate(metaslab_class_t *mc);
78 extern metaslab_class_t *metaslabs_class_create(spa_t *spa,
79     space_map_ops_t *ops);
80 extern void metaslab_class_destroy(metaslab_class_t *mc);
81 extern int metaslab_class_validate(metaslab_class_t *mc);

82
83 void metaslab_class_space_update(metaslab_class_t *mc,
84     extern void metaslab_class_space_update(metaslab_class_t *mc,
85     int64_t alloc_delta, int64_t defer_delta,
86     int64_t space_delta, int64_t dspace_delta);
87 uint64_t metaslab_class_get_alloc(metaslab_class_t *mc);
88 uint64_t metaslab_class_get_space(metaslab_class_t *mc);
89 uint64_t metaslab_class_get_dspace(metaslab_class_t *mc);
90 uint64_t metaslab_class_get_deferred(metaslab_class_t *mc);
91 extern uint64_t metaslab_class_get_alloc(metaslab_class_t *mc);
92 extern uint64_t metaslab_class_get_space(metaslab_class_t *mc);
93 extern uint64_t metaslab_class_get_dspace(metaslab_class_t *mc);
94 extern uint64_t metaslab_class_get_deferred(metaslab_class_t *mc);

95
96 metaslab_group_t *metaslabs_group_create(metaslab_class_t *mc, vdev_t *vd);
97 void metaslab_group_destroy(metaslab_group_t *mg);
98 void metaslab_group_activate(metaslab_group_t *mg);
99 void metaslab_group_passivate(metaslab_group_t *mg);
100 extern metaslab_group_t *metaslabs_group_create(metaslab_class_t *mc,
101     vdev_t *vd);
102 extern void metaslab_group_destroy(metaslab_group_t *mg);
103 extern void metaslab_group_activate(metaslab_group_t *mg);
104 extern void metaslab_group_passivate(metaslab_group_t *mg);

105
106 #ifdef __cplusplus
107 }
108 #endif

109 #endif /* _SYS_METASLAB_H */

```



```

*****
6085 Tue Sep 3 20:27:07 2013
new/usr/src/uts/common/fs/zfs/sys/metaslub_impl.h
4101 metaslab_debug should allow for fine-grained control
4102 space_maps should store more information about themselves
4103 space_map object blocksize should be increased
4104 ::spa_space no longer works
4105 removing a mirrored log device results in a leaked object
4106 asynchronously load metaslab
Reviewed by: Matthew Ahrens <mahrens@delphix.com>
Reviewed by: Adam Leventhal <ahl@delphix.com>
Reviewed by: Sebastien Roy <seb@delphix.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22 * Copyright 2009 Sun Microsystems, Inc. All rights reserved.
23 * Use is subject to license terms.
24 */
25
26 /*
27 * Copyright (c) 2013 by Delphix. All rights reserved.
28 */
29
30 #ifndef _SYS_METASLAB_IMPL_H
31 #define _SYS_METASLAB_IMPL_H
32
33 #include <sys/metaslub.h>
34 #include <sys/space_map.h>
35 #include <sys/range_tree.h>
36 #include <sys/vdev.h>
37 #include <sys/txg.h>
38 #include <sys/avl.h>
39
40 #ifdef __cplusplus
41 extern "C" {
42 #endif
43
44 struct metaslab_class {
45     spa_t          *mc_spa;
46     metaslab_group_t *mc_rotor;
47     metaslab_ops_t *mc_ops;
48     space_map_ops_t *mc_ops;
49     uint64_t       mc_aliquot;
50     uint64_t       mc_alloc_groups; /* # of allocatable groups */
51     uint64_t       mc_alloc; /* total allocated space */
52     uint64_t       mc_deferred; /* total deferred frees */
53     uint64_t       mc_space; /* total space (alloc + free) */

```

```

53     uint64_t       mc_dspace; /* total deflated space */
54 };
55
56 struct metaslab_group {
57     kmutex_t       mg_lock;
58     avl_tree_t     mg_metaslub_tree;
59     uint64_t       mg_aliquot;
60     uint64_t       mg_bonus_area;
61     uint64_t       mg_alloc_failures;
62     boolean_t      mg_allocatable; /* can we allocate? */
63     uint64_t       mg_free_capacity; /* percentage free */
64     int64_t        mg_bias;
65     int64_t        mg_activation_count;
66     metaslab_class_t *mg_class;
67     vdev_t         *mg_vd;
68     taskq_t        *mg_taskq;
69     metaslab_group_t *mg_prev;
70     metaslab_group_t *mg_next;
71 };
72
73 /*
74 * This value defines the number of elements in the ms_lbas array. The value
75 * of 64 was chosen as it covers to cover all power of 2 buckets up to
76 * UINT64_MAX. This is the equivalent of highbit(UINT64_MAX).
77 */
78 #define MAX_LBAS 64
79
80 /*
81 * Each metaslab maintains a set of in-core trees to track metaslab operations.
82 * The in-core free tree (ms_tree) contains the current list of free segments.
83 * As blocks are allocated, the allocated segment are removed from the ms_tree
84 * and added to a per txg allocation tree (ms_alloctree). As blocks are freed,
85 * they are added to the per txg free tree (ms_freetree). These per txg
86 * trees allow us to process all allocations and frees in syncing context
87 * where it is safe to update the on-disk space maps. One additional in-core
88 * tree is maintained to track deferred frees (ms_defertree). Once a block
89 * is freed it will move from the ms_freetree to the ms_defertree. A deferred
90 * free means that a block has been freed but cannot be used by the pool
91 * until TXG_DEFER_SIZE transactions groups later. For example, a block
92 * that is freed in txg 50 will not be available for reallocation until
93 * txg 52 (50 + TXG_DEFER_SIZE). This provides a safety net for uberblock
94 * rollback. A pool could be safely rolled back TXG_DEFERS_SIZE
95 * transactions groups and ensure that no block has been reallocated.
96 * Each metaslab maintains an in-core free map (ms_map) that contains the
97 * current list of free segments. As blocks are allocated, the allocated
98 * segment is removed from the ms_map and added to a per txg allocation map.
99 * As blocks are freed, they are added to the per txg free map. These per
100 * txg maps allow us to process all allocations and frees in syncing context
101 * where it is safe to update the on-disk space maps.
102 */
103
104 /*
105 * The simplified transition diagram looks like this:
106 */
107
108
109
110
111

```

```

112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
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

```

```

112 * Each metabolab's space is tracked in a single space map in the MOS,
79 * Each metabolab's free space is tracked in a space map object in the MOS,
113 * which is only updated in syncing context. Each time we sync a txg,
114 * we append the allocs and frees from that txg to the space map.
115 * The pool space is only updated once all metabolabs have finished syncing.
81 * we append the allocs and frees from that txg to the space map object.
82 * When the txg is done syncing, metabolab_sync_done() updates ms_smo
83 * to ms_smo_syncing. Everything in ms_smo is always safe to allocate.
116 *
117 * To load the in-core free tree we read the space map from disk.
85 * To load the in-core free map we read the space map object from disk.
118 * This object contains a series of alloc and free records that are
119 * combined to make up the list of all free segments in this metabolab. These
120 * segments are represented in-core by the ms_tree and are stored in an
88 * segments are represented in-core by the ms_map and are stored in an
121 * AVL tree.
122 *
123 * As the space map grows (as a result of the appends) it will
124 * eventually become space-inefficient. When the metabolab's in-core free tree
125 * is zfs_condense_pct/100 times the size of the minimal on-disk
126 * representation, we rewrite it in its minimized form. If a metabolab
127 * needs to condense then we must set the ms_condensing flag to ensure
128 * that allocations are not performed on the metabolab that is being written.
91 * As the space map objects grows (as a result of the appends) it will
92 * eventually become space-inefficient. When the space map object is
93 * zfs_condense_pct/100 times the size of the minimal on-disk representation,
94 * we rewrite it in its minimized form.
129 */
130 struct metabolab {
131     kmutex_t      ms_lock;
132     kcondvar_t    ms_load_cv;
133     space_map_t   *ms_sm;
134     metabolab_ops_t *ms_ops;
135     uint64_t      ms_id;
136     uint64_t      ms_start;
137     uint64_t      ms_size;

139     range_tree_t  *ms_alloctree[TXG_SIZE];
140     range_tree_t  *ms_freemap[TXG_SIZE];
141     range_tree_t  *ms_defertree[TXG_DEFER_SIZE];
142     range_tree_t  *ms_tree;

144     boolean_t     ms_condensing; /* condensing? */
145     boolean_t     ms_loaded;
146     boolean_t     ms_loading;

97     kmutex_t      ms_lock; /* metabolab lock */
98     space_map_obj_t ms_smo; /* synced space map object */
99     space_map_obj_t ms_smo_syncing; /* syncing space map object */
100     space_map_t    *ms_allocmap[TXG_SIZE]; /* allocated this txg */
101     space_map_t    *ms_freemap[TXG_SIZE]; /* freed this txg */
102     space_map_t    *ms_defermap[TXG_DEFER_SIZE]; /* deferred frees */
103     space_map_t    *ms_map; /* in-core free space map */
148     int64_t       ms_deferspace; /* sum of ms_defermap[] space */
149     uint64_t      ms_weight; /* weight vs. others in group */
150     uint64_t      ms_factor;
151     uint64_t      ms_access_txg;

153     /*
154     * The metabolab block allocators can optionally use a size-ordered
155     * range tree and/or an array of LBAs. Not all allocators use
156     * this functionality. The ms_size_tree should always contain the
157     * same number of segments as the ms_tree. The only difference
158     * is that the ms_size_tree is ordered by segment sizes.
159     */
160     avl_tree_t     ms_size_tree;

```

```

161     uint64_t      ms_lbas[MAX_LBAS];

163     metabolab_group_t *ms_group; /* metabolab group */
164     avl_node_t      ms_group_node; /* node in metabolab group tree */
165     txg_node_t      ms_txg_node; /* per-txg dirty metabolab links */
166 };
    unchanged_portion_omitted

```

```

*****
3197 Tue Sep 3 20:27:08 2013
new/usr/src/uts/common/fs/zfs/sys/range_tree.h
4101 metaslab_debug should allow for fine-grained control
4102 space_maps should store more information about themselves
4103 space_map object blocksize should be increased
4104 ::spa_space no longer works
4105 removing a mirrored log device results in a leaked object
4106 asynchronously load metaslab
Reviewed by: Matthew Ahrens <mahrens@delphix.com>
Reviewed by: Adam Leventhal <ahl@delphix.com>
Reviewed by: Sebastien Roy <seb@delphix.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22 * Copyright 2009 Sun Microsystems, Inc. All rights reserved.
23 * Use is subject to license terms.
24 */

26 /*
27 * Copyright (c) 2013 by Delphix. All rights reserved.
28 */

30 #ifndef _SYS_RANGE_TREE_H
31 #define _SYS_RANGE_TREE_H

33 #include <sys/avl.h>
34 #include <sys/dmu.h>

36 #ifdef __cplusplus
37 extern "C" {
38 #endif

40 #define RANGE_TREE_HISTOGRAM_SIZE    64

42 typedef struct range_tree_ops range_tree_ops_t;

44 typedef struct range_tree {
45     avl_tree_t    rt_root;        /* offset-ordered segment AVL tree */
46     uint64_t      rt_space;       /* sum of all segments in the map */
47     range_tree_ops_t *rt_ops;
48     void          *rt_arg;

50     /*
51      * The rt_histogram maintains a histogram of ranges. Each bucket,
52      * rt_histogram[i], contains the number of ranges whose size is:
53      * 2^i <= size of range in bytes < 2^(i+1)

```

```

54     */
55     uint64_t      rt_histogram[RANGE_TREE_HISTOGRAM_SIZE];
56     kmutex_t      *rt_lock;      /* pointer to lock that protects map */
57 } range_tree_t;

59 typedef struct range_seg {
60     avl_node_t    rs_node;       /* AVL node */
61     avl_node_t    rs_pp_node;    /* AVL picker-private node */
62     uint64_t      rs_start;      /* starting offset of this segment */
63     uint64_t      rs_end;        /* ending offset (non-inclusive) */
64 } range_seg_t;

66 struct range_tree_ops {
67     void          (*rtop_create)(range_tree_t *rt, void *arg);
68     void          (*rtop_destroy)(range_tree_t *rt, void *arg);
69     void          (*rtop_add)(range_tree_t *rt, range_seg_t *rs, void *arg);
70     void          (*rtop_remove)(range_tree_t *rt, range_seg_t *rs, void *arg);
71     void          (*rtop_vacate)(range_tree_t *rt, void *arg);
72 };

74 typedef void range_tree_func_t(void *arg, uint64_t start, uint64_t size);

76 void range_tree_init(void);
77 void range_tree_fini(void);
78 range_tree_t *range_tree_create(range_tree_ops_t *ops, void *arg, kmutex_t *lp);
79 void range_tree_destroy(range_tree_t *rt);
80 boolean_t range_tree_contains(range_tree_t *rt, uint64_t start, uint64_t size);
81 uint64_t range_tree_space(range_tree_t *rt);
82 void range_tree_verify(range_tree_t *rt, uint64_t start, uint64_t size);
83 void range_tree_swap(range_tree_t **rtsrc, range_tree_t **rtdst);
84 void range_tree_stat_verify(range_tree_t *rt);

86 void range_tree_add(void *arg, uint64_t start, uint64_t size);
87 void range_tree_remove(void *arg, uint64_t start, uint64_t size);

89 void range_tree_vacate(range_tree_t *rt, range_tree_func_t *func, void *arg);
90 void range_tree_walk(range_tree_t *rt, range_tree_func_t *func, void *arg);

92 #ifdef __cplusplus
93 }
94 #endif

96 #endif /* _SYS_RANGE_TREE_H */

```

```

*****
6178 Tue Sep 3 20:27:09 2013
new/usr/src/uts/common/fs/zfs/sys/space_map.h
4101 metaslab_debug should allow for fine-grained control
4102 space_maps should store more information about themselves
4103 space_map object blocksize should be increased
4104 ::spa_space no longer works
4105 removing a mirrored log device results in a leaked object
4106 asynchronously load metaslab
Reviewed by: Matthew Ahrens <mahrens@delphix.com>
Reviewed by: Adam Leventhal <ahl@delphix.com>
Reviewed by: Sebastien Roy <seb@delphix.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22 * Copyright 2009 Sun Microsystems, Inc. All rights reserved.
23 * Use is subject to license terms.
24 */
25
26 /*
27 * Copyright (c) 2013 by Delphix. All rights reserved.
28 * Copyright (c) 2012 by Delphix. All rights reserved.
29 */
30 #ifndef _SYS_SPACE_MAP_H
31 #define _SYS_SPACE_MAP_H
32
33 #include <sys/avl.h>
34 #include <sys/range_tree.h>
35 #include <sys/dmu.h>
36
37 #ifdef __cplusplus
38 extern "C" {
39 #endif
40
41 /*
42 * The size of the space map object has increased to include a histogram.
43 * The SPACE_MAP_SIZE_V0 designates the original size and is used to
44 * maintain backward compatibility.
45 */
46 #define SPACE_MAP_SIZE_V0 (3 * sizeof(uint64_t))
47 #define SPACE_MAP_HISTOGRAM_SIZE(sm) \
48     (sizeof((sm)->sm_phys->smp_histogram) / \
49     sizeof((sm)->sm_phys->smp_histogram[0]))
40 typedef struct space_map_ops space_map_ops_t;
51 */

```

```

52 * The space_map_phys is the on-disk representation of the space map.
53 * Consumers of space maps should never reference any of the members of this
54 * structure directly. These members may only be updated in syncing context.
55 *
56 * Note the smp_object is no longer used but remains in the structure
57 * for backward compatibility.
58 */
59 typedef struct space_map_phys {
60     uint64_t smp_object; /* on-disk space map object */
61     uint64_t smp_objsize; /* size of the object */
62     uint64_t smp_alloc; /* space allocated from the map */
63     uint64_t smp_pad[5]; /* reserved */
64
65     /*
66      * The smp_histogram maintains a histogram of free regions. Each
67      * bucket, smp_histogram[i], contains the number of free regions
68      * whose size is:
69      * 2^(i+sm_shift) <= size of free region in bytes < 2^(i+sm_shift+1)
70      */
71     uint64_t smp_histogram[32]; /* histogram of free space */
72 } space_map_phys_t;
73
74 /*
75 * The space map object defines a region of space, its size, how much is
76 * allocated, and the on-disk object that stores this information.
77 * Consumers of space maps may only access the members of this structure.
78 */
79 typedef struct space_map {
80     avl_tree_t sm_root; /* offset-ordered segment AVL tree */
81     uint64_t sm_space; /* sum of all segments in the map */
82     uint64_t sm_start; /* start of map */
83     uint64_t sm_size; /* size of map */
84     uint8_t sm_shift; /* unit shift */
85     uint64_t sm_length; /* synced length */
86     uint64_t sm_alloc; /* synced space allocated */
87     objset_t *sm_os; /* objset for this map */
88     uint64_t sm_object; /* object id for this map */
89     uint32_t sm_blkisz; /* block size for space map */
90     dmu_buf_t *sm_dbuf; /* space_map_phys_t dbuf */
91     space_map_phys_t *sm_phys; /* on-disk space map */
92     uint8_t sm_loaded; /* map loaded? */
93     uint8_t sm_loading; /* map loading? */
94     uint8_t sm_condensing; /* map condensing? */
95     kcondvar_t sm_load_cv; /* map load completion */
96     space_map_ops_t *sm_ops; /* space map block picker ops vector */
97     avl_tree_t *sm_pp_root; /* size-ordered, picker-private tree */
98     void *sm_ppd; /* picker-private data */
99     kmutex_t *sm_lock; /* pointer to lock that protects map */
100 } space_map_t;
101
102 typedef struct space_seg {
103     avl_node_t ss_node; /* AVL node */
104     avl_node_t ss_pp_node; /* AVL picker-private node */
105     uint64_t ss_start; /* starting offset of this segment */
106     uint64_t ss_end; /* ending offset (non-inclusive) */
107 } space_seg_t;
108
109 typedef struct space_ref {
110     avl_node_t sr_node; /* AVL node */
111     uint64_t sr_offset; /* offset (start or end) */
112     int64_t sr_refcnt; /* associated reference count */
113 } space_ref_t;
114
115 typedef struct space_map_obj {
116     uint64_t smo_object; /* on-disk space map object */
117     uint64_t smo_objsize; /* size of the object */

```

```

74     uint64_t      smo_alloc;    /* space allocated from the map */
75 } space_map_obj_t;

77 struct space_map_ops {
78     void (*smop_load)(space_map_t *sm);
79     void (*smop_unload)(space_map_t *sm);
80     uint64_t (*smop_alloc)(space_map_t *sm, uint64_t size);
81     void (*smop_claim)(space_map_t *sm, uint64_t start, uint64_t size);
82     void (*smop_free)(space_map_t *sm, uint64_t start, uint64_t size);
83     uint64_t (*smop_max)(space_map_t *sm);
84     boolean_t (*smop_fragmented)(space_map_t *sm);
85 };

93 /*
94 * debug entry
95 *
96 *      1      3      10      50
97 *  +-----+-----+-----+-----+
98 *  | 1 | action | syncpass | txg (lower bits) |
99 *  +-----+-----+-----+-----+
100 * 63 62 60 59 50 49 0
101 *
102 *
103 * non-debug entry
104 *
105 *      1      47      1      15
106 *  +-----+-----+-----+-----+
107 *  | 0 | offset (sm_shift units) | type | run |
108 *  +-----+-----+-----+-----+
109 * 63 62 17 16 15 0
110 */

112 /* All this stuff takes and returns bytes */
113 #define SM_RUN_DECODE(x) (BF64_DECODE(x, 0, 15) + 1)
114 #define SM_RUN_ENCODE(x) BF64_ENCODE(x) - 1, 0, 15)
115 #define SM_TYPE_DECODE(x) BF64_DECODE(x, 15, 1)
116 #define SM_TYPE_ENCODE(x) BF64_ENCODE(x, 15, 1)
117 #define SM_OFFSET_DECODE(x) BF64_DECODE(x, 16, 47)
118 #define SM_OFFSET_ENCODE(x) BF64_ENCODE(x, 16, 47)
119 #define SM_DEBUG_DECODE(x) BF64_DECODE(x, 63, 1)
120 #define SM_DEBUG_ENCODE(x) BF64_ENCODE(x, 63, 1)

122 #define SM_DEBUG_ACTION_DECODE(x) BF64_DECODE(x, 60, 3)
123 #define SM_DEBUG_ACTION_ENCODE(x) BF64_ENCODE(x, 60, 3)

125 #define SM_DEBUG_SYNCPASS_DECODE(x) BF64_DECODE(x, 50, 10)
126 #define SM_DEBUG_SYNCPASS_ENCODE(x) BF64_ENCODE(x, 50, 10)

128 #define SM_DEBUG_TXG_DECODE(x) BF64_DECODE(x, 0, 50)
129 #define SM_DEBUG_TXG_ENCODE(x) BF64_ENCODE(x, 0, 50)

131 #define SM_RUN_MAX SM_RUN_DECODE(~0ULL)

133 typedef enum {
134     SM_ALLOC,
135     SM_FREE
136 } maptype_t;
137 #define SM_ALLOC 0x0
138 #define SM_FREE 0x1

139 /*
140 * The data for a given space map can be kept on blocks of any size.
141 * Larger blocks entail fewer i/o operations, but they also cause the
142 * DMU to keep more data in-core, and also to waste more i/o bandwidth
143 * when only a few blocks have changed since the last transaction group.
144 * Rather than having a fixed block size for all space maps the block size

```

```

144 * can adjust as needed (see space_map_max_blkksz). Set the initial block
145 * size for the space map to 4k.
146 * This could use a lot more research, but for now, set the freelist
147 * block size to 4k (2^12).
148 */
149 #define SPACE_MAP_INITIAL_BLOCKSIZE (1ULL << 12)
150 #define SPACE_MAP_BLOCKSHIFT 12

151 int space_map_load(space_map_t *sm, range_tree_t *rt, maptype_t maptype);
152 typedef void space_map_func_t(space_map_t *sm, uint64_t start, uint64_t size);

153 void space_map_histogram_clear(space_map_t *sm);
154 void space_map_histogram_add(space_map_t *sm, range_tree_t *rt,
155     dmu_tx_t *tx);
156 extern void space_map_init(void);
157 extern void space_map_fini(void);
158 extern void space_map_create(space_map_t *sm, uint64_t start, uint64_t size,
159     uint8_t shift, kmutex_t *lp);
160 extern void space_map_destroy(space_map_t *sm);
161 extern void space_map_add(space_map_t *sm, uint64_t start, uint64_t size);
162 extern void space_map_remove(space_map_t *sm, uint64_t start, uint64_t size);
163 extern boolean_t space_map_contains(space_map_t *sm,
164     uint64_t start, uint64_t size);
165 extern space_seg_t *space_map_find(space_map_t *sm, uint64_t start,
166     uint64_t size, avl_index_t *wherep);
167 extern void space_map_swap(space_map_t **msrc, space_map_t **mdest);
168 extern void space_map_vacate(space_map_t *sm,
169     space_map_func_t *func, space_map_t *mdest);
170 extern void space_map_walk(space_map_t *sm,
171     space_map_func_t *func, space_map_t *mdest);

172 void space_map_update(space_map_t *sm);
173 extern void space_map_load_wait(space_map_t *sm);
174 extern int space_map_load(space_map_t *sm, space_map_ops_t *ops,
175     uint8_t maptype, space_map_obj_t *smo, objset_t *os);
176 extern void space_map_unload(space_map_t *sm);

177 uint64_t space_map_object(space_map_t *sm);
178 uint64_t space_map_allocated(space_map_t *sm);
179 uint64_t space_map_length(space_map_t *sm);
180 extern uint64_t space_map_alloc(space_map_t *sm, uint64_t size);
181 extern void space_map_claim(space_map_t *sm, uint64_t start, uint64_t size);
182 extern void space_map_free(space_map_t *sm, uint64_t start, uint64_t size);
183 extern uint64_t space_map_maxsize(space_map_t *sm);

184 void space_map_write(space_map_t *sm, range_tree_t *rt, maptype_t maptype,
185     dmu_tx_t *tx);
186 void space_map_truncate(space_map_t *sm, dmu_tx_t *tx);
187 uint64_t space_map_alloc(objset_t *os, dmu_tx_t *tx);
188 void space_map_free(space_map_t *sm, dmu_tx_t *tx);
189 extern void space_map_sync(space_map_t *sm, uint8_t maptype,
190     space_map_obj_t *smo, objset_t *os, dmu_tx_t *tx);
191 extern void space_map_truncate(space_map_obj_t *smo,
192     objset_t *os, dmu_tx_t *tx);

193 int space_map_open(space_map_t **smp, objset_t *os, uint64_t object,
194     uint64_t start, uint64_t size, uint8_t shift, kmutex_t *lp);
195 void space_map_close(space_map_t *sm);
196 extern void space_map_ref_create(avl_tree_t *t);
197 extern void space_map_ref_destroy(avl_tree_t *t);
198 extern void space_map_ref_add_seg(avl_tree_t *t,
199     uint64_t start, uint64_t end, int64_t refcnt);
200 extern void space_map_ref_add_map(avl_tree_t *t,
201     space_map_t *sm, int64_t refcnt);
202 extern void space_map_ref_generate_map(avl_tree_t *t,
203     space_map_t *sm, int64_t minref);

```

```
171 int64_t space_map_alloc_delta(space_map_t *sm);
```

```
173 #ifdef __cplusplus
```

```
174 }
```

```
_____unchanged_portion_omitted_
```

new/usr/src/uts/common/fs/zfs/sys/space_reftree.h

1

1711 Tue Sep 3 20:27:10 2013
new/usr/src/uts/common/fs/zfs/sys/space_reftree.h
4101 metaslab_debug should allow for fine-grained control
4102 space_maps should store more information about themselves
4103 space_map object blocksize should be increased
4104 ::spa_space no longer works
4105 removing a mirrored log device results in a leaked object
4106 asynchronously load metaslab
Reviewed by: Matthew Ahrens <mahrens@delphix.com>
Reviewed by: Adam Leventhal <ahl@delphix.com>
Reviewed by: Sebastien Roy <seb@delphix.com>

```
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22 * Copyright 2009 Sun Microsystems, Inc. All rights reserved.
23 * Use is subject to license terms.
24 */
25
26 /*
27 * Copyright (c) 2013 by Delphix. All rights reserved.
28 */
29
30 #ifndef _SYS_SPACE_REFTREE_H
31 #define _SYS_SPACE_REFTREE_H
32
33 #include <sys/range_tree.h>
34
35 #ifdef __cplusplus
36 extern "C" {
37 #endif
38
39 typedef struct space_ref {
40     avl_node_t      sr_node;          /* AVL node */
41     uint64_t        sr_offset;        /* range offset (start or end) */
42     int64_t         sr_refcnt;        /* associated reference count */
43 } space_ref_t;
44
45 void space_reftree_create(avl_tree_t *t);
46 void space_reftree_destroy(avl_tree_t *t);
47 void space_reftree_add_seg(avl_tree_t *t, uint64_t start, uint64_t end,
48     int64_t refcnt);
49 void space_reftree_add_map(avl_tree_t *t, range_tree_t *rt, int64_t refcnt);
50 void space_reftree_generate_map(avl_tree_t *t, range_tree_t *rt,
51     int64_t minref);
52
53 #endif
54 #endif
55 #endif
```

new/usr/src/uts/common/fs/zfs/sys/space_reftree.h

2

```
54 }
55 #endif
56
57 #endif /* _SYS_SPACE_REFTREE_H */
```

new/usr/src/uts/common/fs/zfs/sys/vdev_impl.h

1

```
*****
11807 Tue Sep 3 20:27:10 2013
new/usr/src/uts/common/fs/zfs/sys/vdev_impl.h
4101 metaslab_debug should allow for fine-grained control
4102 space_maps should store more information about themselves
4103 space_map object blocksize should be increased
4104 ::spa_space no longer works
4105 removing a mirrored log device results in a leaked object
4106 asynchronously load metaslab
Reviewed by: Matthew Ahrens <mahrens@delphix.com>
Reviewed by: Adam Leventhal <ahl@delphix.com>
Reviewed by: Sebastien Roy <seb@delphix.com>
*****
unchanged_portion_omitted_

121 /*
122 * Virtual device descriptor
123 */
124 struct vdev {
125     /*
126     * Common to all vdev types.
127     */
128     uint64_t    vdev_id;        /* child number in vdev parent */
129     uint64_t    vdev_guid;     /* unique ID for this vdev */
130     uint64_t    vdev_guid_sum; /* self guid + all child guids */
131     uint64_t    vdev_orig_guid; /* orig. guid prior to remove */
132     uint64_t    vdev_asize;    /* allocatable device capacity */
133     uint64_t    vdev_min_asize; /* min acceptable asize */
134     uint64_t    vdev_max_asize; /* max acceptable asize */
135     uint64_t    vdev_ashift;   /* block alignment shift */
136     uint64_t    vdev_state;    /* see VDEV_STATE_* #defines */
137     uint64_t    vdev_prevstate; /* used when reopening a vdev */
138     vdev_ops_t  *vdev_ops;     /* vdev operations */
139     spa_t       *vdev_spa;     /* spa for this vdev */
140     void        *vdev_tsd;     /* type-specific data */
141     vnode_t     *vdev_name_vp; /* vnode for pathname */
142     vnode_t     *vdev_devid_vp; /* vnode for devid */
143     vdev_t      *vdev_top;     /* top-level vdev */
144     vdev_t      *vdev_parent;  /* parent vdev */
145     vdev_t      **vdev_child; /* array of children */
146     uint64_t    vdev_children; /* number of children */
147     space_map_t vdev_dtl[DTL_TYPES]; /* in-core dirty time logs */
148     vdev_stat_t vdev_stat;     /* virtual device statistics */
149     boolean_t   vdev_expanding; /* expand the vdev? */
150     boolean_t   vdev_reopening; /* reopen in progress? */
151     int         vdev_open_error; /* error on last open */
152     kthread_t   *vdev_open_thread; /* thread opening children */
153     uint64_t    vdev_crtxg;    /* txg when top-level was added */

154     /*
155     * Top-level vdev state.
156     */
157     uint64_t    vdev_ms_array; /* metaslab array object */
158     uint64_t    vdev_ms_shift; /* metaslab size shift */
159     uint64_t    vdev_ms_count; /* number of metaslabs */
160     metaslab_group_t *vdev_mg; /* metaslab group */
161     metaslab_t  **vdev_ms;    /* metaslab array */
162     txg_list_t  vdev_ms_list; /* per-txg dirty metaslab lists */
163     txg_list_t  vdev_dtl_list; /* per-txg dirty DTL lists */
164     txg_node_t  vdev_txg_node; /* per-txg dirty vdev linkage */
165     boolean_t   vdev_remove_wanted; /* async remove wanted? */
166     boolean_t   vdev_probe_wanted; /* async probe wanted? */
167     uint64_t    vdev_removing; /* device is being removed? */
168     list_node_t vdev_config_dirty_node; /* config dirty list */
169     list_node_t vdev_state_dirty_node; /* state dirty list */
170     uint64_t    vdev_deflate_ratio; /* deflation ratio (x512) */

```

new/usr/src/uts/common/fs/zfs/sys/vdev_impl.h

2

```
170     uint64_t    vdev_islog;   /* is an intent log device */
171     uint64_t    vdev_removing; /* device is being removed? */
172     boolean_t   vdev_ishole;  /* is a hole in the namespace */
173     uint64_t    vdev_ishole;  /* is a hole in the namespace */

174     /*
175     * Leaf vdev state.
176     */
177     range_tree_t *vdev_dtl[DTL_TYPES]; /* dirty time logs */
178     space_map_t  *vdev_dtl_sm; /* dirty time log space map */
179     txg_node_t   vdev_dtl_node; /* per-txg dirty DTL linkage */
180     uint64_t     vdev_dtl_object; /* DTL object */
181     uint64_t     vdev_psize;     /* physical device capacity */
182     space_map_obj_t vdev_dtl_smo; /* dirty time log space map obj */
183     txg_node_t   vdev_dtl_node; /* per-txg dirty DTL linkage */
184     uint64_t     vdev_whole_disk; /* true if this is a whole disk */
185     uint64_t     vdev_offline;  /* persistent offline state */
186     uint64_t     vdev_faulted;  /* persistent faulted state */
187     uint64_t     vdev_degraded; /* persistent degraded state */
188     uint64_t     vdev_removed;  /* persistent removed state */
189     uint64_t     vdev_resilver_txg; /* persistent resilvering state */
190     uint64_t     vdev_nparity;  /* number of parity devices for raidz */
191     char         *vdev_path;    /* vdev path (if any) */
192     char         *vdev_devid;   /* vdev devid (if any) */
193     char         *vdev_physpath; /* vdev device path (if any) */
194     char         *vdev_fru;     /* physical FRU location */
195     uint64_t     vdev_not_present; /* not present during import */
196     uint64_t     vdev_unspare;  /* unspare when resilvering done */
197     hrtime_t     vdev_last_try; /* last reopen time */
198     boolean_t    vdev_nowritecache; /* true if flushwritecache failed */
199     boolean_t    vdev_checkremove; /* temporary online test */
200     boolean_t    vdev_forcefault; /* force online fault */
201     boolean_t    vdev_splitting; /* split or repair in progress */
202     boolean_t    vdev_delayed_close; /* delayed device close? */
203     boolean_t    vdev_tmppoffline; /* device taken offline temporarily? */
204     boolean_t    vdev_detached; /* device detached? */
205     boolean_t    vdev_cant_read; /* vdev is failing all reads */
206     boolean_t    vdev_cant_write; /* vdev is failing all writes */
207     boolean_t    vdev_isspare;  /* was a hot spare */
208     boolean_t    vdev_isl2cache; /* was a l2cache device */
209     uint8_t     vdev_tmppoffline; /* device taken offline temporarily? */
210     uint8_t     vdev_detached; /* device detached? */
211     uint8_t     vdev_cant_read; /* vdev is failing all reads */
212     uint8_t     vdev_cant_write; /* vdev is failing all writes */
213     uint64_t    vdev_isspare;  /* was a hot spare */
214     uint64_t    vdev_isl2cache; /* was a l2cache device */
215     vdev_queue_t vdev_queue;    /* I/O deadline schedule queue */
216     vdev_cache_t vdev_cache;    /* physical block cache */
217     spa_aux_vdev_t *vdev_aux;   /* for l2cache vdevs */
218     zio_t        *vdev_probe_zio; /* root of current probe */
219     vdev_aux_t   vdev_label_aux; /* on-disk aux state */

220     /*
221     * For DTrace to work in userland (libzpool) context, these fields must
222     * remain at the end of the structure. DTrace will use the kernel's
223     * CTF definition for 'struct vdev', and since the size of a kmutex_t is
224     * larger in userland, the offsets for the rest of the fields would be
225     * incorrect.
226     */
227     kmutex_t    vdev_dtl_lock; /* vdev_dtl_{map,resilver} */
228     kmutex_t    vdev_stat_lock; /* vdev_stat */
229     kmutex_t    vdev_probe_lock; /* protects vdev_probe_zio */
230 };
unchanged_portion_omitted_

252 /*

```



```

253 * vdev_dirty() flags
254 */
255 #define VDD_METASLAB    0x01
256 #define VDD_DTL        0x02

258 /* Offset of embedded boot loader region on each label */
259 #define VDEV_BOOT_OFFSET    (2 * sizeof (vdev_label_t))
260 /*
261  * Size of embedded boot loader region on each label.
262  * The total size of the first two labels plus the boot area is 4MB.
263  */
264 #define VDEV_BOOT_SIZE      (7ULL << 19)          /* 3.5M */

266 /*
267  * Size of label regions at the start and end of each leaf device.
268  */
269 #define VDEV_LABEL_START_SIZE    (2 * sizeof (vdev_label_t) + VDEV_BOOT_SIZE)
270 #define VDEV_LABEL_END_SIZE      (2 * sizeof (vdev_label_t))
271 #define VDEV_LABELS              4
272 #define VDEV_BEST_LABEL          VDEV_LABELS

274 #define VDEV_ALLOC_LOAD          0
275 #define VDEV_ALLOC_ADD          1
276 #define VDEV_ALLOC_SPARE        2
277 #define VDEV_ALLOC_L2CACHE      3
278 #define VDEV_ALLOC_ROOTPOOL     4
279 #define VDEV_ALLOC_SPLIT        5
280 #define VDEV_ALLOC_ATTACH       6

282 /*
283  * Allocate or free a vdev
284  */
285 extern vdev_t *vdev_alloc_common(spa_t *spa, uint_t id, uint64_t guid,
286     vdev_ops_t *ops);
287 extern int vdev_alloc(spa_t *spa, vdev_t **vdp, nvlist_t *config,
288     vdev_t *parent, uint_t id, int alloctype);
289 extern void vdev_free(vdev_t *vd);

291 /*
292  * Add or remove children and parents
293  */
294 extern void vdev_add_child(vdev_t *pvd, vdev_t *cvd);
295 extern void vdev_remove_child(vdev_t *pvd, vdev_t *cvd);
296 extern void vdev_compact_children(vdev_t *pvd);
297 extern vdev_t *vdev_add_parent(vdev_t *cvd, vdev_ops_t *ops);
298 extern void vdev_remove_parent(vdev_t *cvd);

300 /*
301  * vdev sync load and sync
302  */
303 extern void vdev_load_log_state(vdev_t *nvd, vdev_t *ovd);
304 extern boolean_t vdev_log_state_valid(vdev_t *vd);
305 extern void vdev_load(vdev_t *vd);
306 extern int vdev_dtl_load(vdev_t *vd);
307 extern void vdev_sync(vdev_t *vd, uint64_t txg);
308 extern void vdev_sync_done(vdev_t *vd, uint64_t txg);
309 extern void vdev_dirty(vdev_t *vd, int flags, void *arg, uint64_t txg);
310 extern void vdev_dirty_leaves(vdev_t *vd, int flags, uint64_t txg);

312 /*
313  * Available vdev types.
314  */
315 extern vdev_ops_t vdev_root_ops;
316 extern vdev_ops_t vdev_mirror_ops;
317 extern vdev_ops_t vdev_replacing_ops;
318 extern vdev_ops_t vdev RAIDZ_ops;

```

```

319 extern vdev_ops_t vdev_disk_ops;
320 extern vdev_ops_t vdev_file_ops;
321 extern vdev_ops_t vdev_missing_ops;
322 extern vdev_ops_t vdev_hole_ops;
323 extern vdev_ops_t vdev_spare_ops;

325 /*
326  * Common size functions
327  */
328 extern uint64_t vdev_default_asize(vdev_t *vd, uint64_t psize);
329 extern uint64_t vdev_get_min_asize(vdev_t *vd);
330 extern void vdev_set_min_asize(vdev_t *vd);

332 /*
333  * Global variables
334  */
335 /* zdb uses this tunable, so it must be declared here to make lint happy. */
336 extern int zfs_vdev_cache_size;

338 /*
339  * The vdev_buf_t is used to translate between zio_t and buf_t, and back again.
340  */
341 typedef struct vdev_buf {
342     buf_t    vb_buf;          /* buffer that describes the io */
343     zio_t    *vb_io;         /* pointer back to the original zio_t */
344 } vdev_buf_t;
_____ unchanged_portion_omitted _____

```

new/usr/src/uts/common/fs/zfs/sys/zfeature.h

1

```
*****
1809 Tue Sep 3 20:27:12 2013
new/usr/src/uts/common/fs/zfs/sys/zfeature.h
4101 metaslab_debug should allow for fine-grained control
4102 space_maps should store more information about themselves
4103 space_map object blocksize should be increased
4104 ::spa_space no longer works
4105 removing a mirrored log device results in a leaked object
4106 asynchronously load metaslab
Reviewed by: Matthew Ahrens <mahrens@delphix.com>
Reviewed by: Adam Leventhal <ahl@delphix.com>
Reviewed by: Sebastien Roy <seb@delphix.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */

22 /*
23  * Copyright (c) 2013 by Delphix. All rights reserved.
23  * Copyright (c) 2012 by Delphix. All rights reserved.
24 */

26 #ifndef _SYS_ZFEATURE_H
27 #define _SYS_ZFEATURE_H

29 #include <sys/nvpair.h>
30 #include "zfeature_common.h"

32 #ifdef __cplusplus
33 extern "C" {
34 #endif

36 struct spa;
37 struct dmu_tx;
38 struct objset;

40 extern boolean_t feature_is_supported(struct objset *os, uint64_t obj,
41     uint64_t desc_obj, nvlist_t *unsup_feat, nvlist_t *enabled_feat);

43 extern void spa_feature_create_zap_objects(struct spa *, struct dmu_tx *);
44 extern void spa_feature_enable(struct spa *, zfeature_info_t *,
45     struct dmu_tx *);
46 extern void spa_feature_incr(struct spa *, zfeature_info_t *, struct dmu_tx *);
47 extern void spa_feature_decr(struct spa *, zfeature_info_t *, struct dmu_tx *);
48 extern boolean_t spa_feature_is_enabled(struct spa *, zfeature_info_t *);
49 extern boolean_t spa_feature_is_active(struct spa *, zfeature_info_t *);
50 extern int spa_feature_get_refcount(struct spa *, zfeature_info_t *);

52 #ifdef __cplusplus
```

new/usr/src/uts/common/fs/zfs/sys/zfeature.h

2

```
53 }
    _____
    unchanged_portion_omitted
```

```

*****
89088 Tue Sep 3 20:27:13 2013
new/usr/src/uts/common/fs/zfs/vdev.c
4101 metaslab_debug should allow for fine-grained control
4102 space_maps should store more information about themselves
4103 space_map object blocksize should be increased
4104 ::spa_space no longer works
4105 removing a mirrored log device results in a leaked object
4106 asynchronously load metaslab
Reviewed by: Matthew Ahrens <mahrens@delphix.com>
Reviewed by: Adam Leventhal <ahl@delphix.com>
Reviewed by: Sebastien Roy <seb@delphix.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21
22 /*
23  * Copyright (c) 2005, 2010, Oracle and/or its affiliates. All rights reserved.
24  * Copyright 2011 Nexenta Systems, Inc. All rights reserved.
25  * Copyright (c) 2013 by Delphix. All rights reserved.
26 */
27
28 #include <sys/zfs_context.h>
29 #include <sys/fm/fs/zfs.h>
30 #include <sys/spa.h>
31 #include <sys/spa_impl.h>
32 #include <sys/dmu.h>
33 #include <sys/dmu_tx.h>
34 #include <sys/vdev_impl.h>
35 #include <sys/uberblock_impl.h>
36 #include <sys/metaslab.h>
37 #include <sys/metaslab_impl.h>
38 #include <sys/space_map.h>
39 #include <sys/space_reftree.h>
40 #include <sys/zio.h>
41 #include <sys/zap.h>
42 #include <sys/fs/zfs.h>
43 #include <sys/arc.h>
44 #include <sys/zil.h>
45 #include <sys/dsl_scan.h>
46
47 /*
48  * Virtual device management.
49  */
50
51 static vdev_ops_t *vdev_ops_table[] = {
52     &vdev_root_ops,
53     &vdev_raidz_ops,

```

```

54     &vdev_mirror_ops,
55     &vdev_replacing_ops,
56     &vdev_spare_ops,
57     &vdev_disk_ops,
58     &vdev_file_ops,
59     &vdev_missing_ops,
60     &vdev_hole_ops,
61     NULL
62 };
63
64 _____ unchanged_portion_omitted _____
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79 /*
80  * Allocate and minimally initialize a vdev_t.
81  */
82 vdev_t *
83 vdev_alloc_common(spa_t *spa, uint_t id, uint64_t guid, vdev_ops_t *ops)
84 {
85     vdev_t *vd;
86
87     vd = kmem_zalloc(sizeof (vdev_t), KM_SLEEP);
88
89     if (spa->spa_root_vdev == NULL) {
90         ASSERT(ops == &vdev_root_ops);
91         spa->spa_root_vdev = vd;
92         spa->spa_load_guid = spa_generate_guid(NULL);
93     }
94
95     if (guid == 0 && ops != &vdev_hole_ops) {
96         if (spa->spa_root_vdev == vd) {
97             /*
98              * The root vdev's guid will also be the pool guid,
99              * which must be unique among all pools.
100            */
101            guid = spa_generate_guid(NULL);
102        } else {
103            /*
104             * Any other vdev's guid must be unique within the pool.
105            */
106            guid = spa_generate_guid(spa);
107        }
108        ASSERT(!spa_guid_exists(spa_guid(spa), guid));
109    }
110
111    vd->vdev_spa = spa;
112    vd->vdev_id = id;
113    vd->vdev_guid = guid;
114    vd->vdev_guid_sum = guid;
115    vd->vdev_ops = ops;
116    vd->vdev_state = VDEV_STATE_CLOSED;
117    vd->vdev_ishole = (ops == &vdev_hole_ops);
118
119    mutex_init(&vd->vdev_dtl_lock, NULL, MUTEX_DEFAULT, NULL);
120    mutex_init(&vd->vdev_stat_lock, NULL, MUTEX_DEFAULT, NULL);
121    mutex_init(&vd->vdev_probe_lock, NULL, MUTEX_DEFAULT, NULL);
122    for (int t = 0; t < DTL_TYPES; t++) {
123        vd->vdev_dtl[t] = range_tree_create(NULL, NULL,
124            space_map_create(&vd->vdev_dtl[t], 0, -1ULL, 0,
125                &vd->vdev_dtl_lock));
126    }
127    txg_list_create(&vd->vdev_ms_list,
128        offsetof(struct metaslab, ms_txg_node));
129    txg_list_create(&vd->vdev_dtl_list,
130        offsetof(struct vdev, vdev_dtl_node));
131    vd->vdev_stat.vs_timestamp = gethrtime();
132    vdev_queue_init(vd);
133    vdev_cache_init(vd);

```



```

466  /*
467  * Get the alignment requirement.
468  */
469  (void) nvlist_lookup_uint64(nv, ZPOOL_CONFIG_ASHIFT, &vdev->vdev_ashift);

471  /*
472  * Retrieve the vdev creation time.
473  */
474  (void) nvlist_lookup_uint64(nv, ZPOOL_CONFIG_CREATE_TXG,
475  &vdev->vdev_crtxg);

477  /*
478  * If we're a top-level vdev, try to load the allocation parameters.
479  */
480  if (parent && !parent->vdev_parent &&
481  (alloctype == VDEV_ALLOC_LOAD || alloctype == VDEV_ALLOC_SPLIT)) {
482  (void) nvlist_lookup_uint64(nv, ZPOOL_CONFIG_METASLAB_ARRAY,
483  &vdev->vdev_ms_array);
484  (void) nvlist_lookup_uint64(nv, ZPOOL_CONFIG_METASLAB_SHIFT,
485  &vdev->vdev_ms_shift);
486  (void) nvlist_lookup_uint64(nv, ZPOOL_CONFIG_ASIZE,
487  &vdev->vdev_asize);
488  (void) nvlist_lookup_uint64(nv, ZPOOL_CONFIG_REMOVING,
489  &vdev->vdev_removing);
490  }

492  if (parent && !parent->vdev_parent && alloctype != VDEV_ALLOC_ATTACH) {
493  ASSERT(alloctype == VDEV_ALLOC_LOAD ||
494  alloctype == VDEV_ALLOC_ADD ||
495  alloctype == VDEV_ALLOC_SPLIT ||
496  alloctype == VDEV_ALLOC_ROOTPOOL);
497  vd->vdev_mg = metaslab_group_create(islog ?
498  spa_log_class(spa) : spa_normal_class(spa), vd);
499  }

501  /*
502  * If we're a leaf vdev, try to load the DTL object and other state.
503  */
504  if (vd->vdev_ops->vdev_op_leaf &&
505  (alloctype == VDEV_ALLOC_LOAD || alloctype == VDEV_ALLOC_L2CACHE ||
506  alloctype == VDEV_ALLOC_ROOTPOOL)) {
507  if (alloctype == VDEV_ALLOC_LOAD) {
508  (void) nvlist_lookup_uint64(nv, ZPOOL_CONFIG_DTL,
509  &vdev->vdev_dtl_object);
510  &vdev->vdev_dtl_smo.smo_object);
511  (void) nvlist_lookup_uint64(nv, ZPOOL_CONFIG_UNSPARE,
512  &vdev->vdev_unspare);
513  }

514  if (alloctype == VDEV_ALLOC_ROOTPOOL) {
515  uint64_t spare = 0;

517  if (nvlist_lookup_uint64(nv, ZPOOL_CONFIG_IS_SPARE,
518  &spare) == 0 && spare)
519  spa_spare_add(vd);
520  }

522  (void) nvlist_lookup_uint64(nv, ZPOOL_CONFIG_OFFLINE,
523  &vdev->vdev_offline);

525  (void) nvlist_lookup_uint64(nv, ZPOOL_CONFIG_RESILVER_TXG,
526  &vdev->vdev_resilver_txg);

528  /*
529  * When importing a pool, we want to ignore the persistent fault

```

```

530  * state, as the diagnosis made on another system may not be
531  * valid in the current context. Local vdevs will
532  * remain in the faulted state.
533  */
534  if (spa_load_state(spa) == SPA_LOAD_OPEN) {
535  (void) nvlist_lookup_uint64(nv, ZPOOL_CONFIG_FAULTED,
536  &vdev->vdev_faulted);
537  (void) nvlist_lookup_uint64(nv, ZPOOL_CONFIG_DEGRADED,
538  &vdev->vdev_degraded);
539  (void) nvlist_lookup_uint64(nv, ZPOOL_CONFIG_REMOVED,
540  &vdev->vdev_removed);

542  if (vd->vdev_faulted || vd->vdev_degraded) {
543  char *aux;

545  vd->vdev_label_aux =
546  VDEV_AUX_ERR_EXCEEDED;
547  if (nvlist_lookup_string(nv,
548  ZPOOL_CONFIG_AUX_STATE, &aux) == 0 &&
549  strcmp(aux, "external") == 0)
550  vd->vdev_label_aux = VDEV_AUX_EXTERNAL;
551  }
552  }
553  }

555  /*
556  * Add ourselves to the parent's list of children.
557  */
558  vdev_add_child(parent, vd);

560  *vdp = vd;

562  return (0);
563  }

565  void
566  vdev_free(vdev_t *vd)
567  {
568  spa_t *spa = vd->vdev_spa;

570  /*
571  * vdev_free() implies closing the vdev first. This is simpler than
572  * trying to ensure complicated semantics for all callers.
573  */
574  vdev_close(vd);

576  ASSERT(!list_link_active(&vdev->vdev_config_dirty_node));
577  ASSERT(!list_link_active(&vdev->vdev_state_dirty_node));

579  /*
580  * Free all children.
581  */
582  for (int c = 0; c < vd->vdev_children; c++)
583  vdev_free(vd->vdev_child[c]);

585  ASSERT(vd->vdev_child == NULL);
586  ASSERT(vd->vdev_guid_sum == vd->vdev_guid);

588  /*
589  * Discard allocation state.
590  */
591  if (vd->vdev_mg != NULL) {
592  vdev metaslab_fini(vd);
593  metaslab_group_destroy(vd->vdev_mg);
594  }

```

```

596     ASSERT0(vd->vdev_stat.vs_space);
597     ASSERT0(vd->vdev_stat.vs_dspace);
598     ASSERT0(vd->vdev_stat.vs_alloc);

600     /*
601      * Remove this vdev from its parent's child list.
602      */
603     vdev_remove_child(vd->vdev_parent, vd);

605     ASSERT(vd->vdev_parent == NULL);

607     /*
608      * Clean up vdev structure.
609      */
610     vdev_queue_fini(vd);
611     vdev_cache_fini(vd);

613     if (vd->vdev_path)
614         spa_strfree(vd->vdev_path);
615     if (vd->vdev_devid)
616         spa_strfree(vd->vdev_devid);
617     if (vd->vdev_physpath)
618         spa_strfree(vd->vdev_physpath);
619     if (vd->vdev_fru)
620         spa_strfree(vd->vdev_fru);

622     if (vd->vdev_isspare)
623         spa_spare_remove(vd);
624     if (vd->vdev_isl2cache)
625         spa_l2cache_remove(vd);

627     txg_list_destroy(&vd->vdev_ms_list);
628     txg_list_destroy(&vd->vdev_dtl_list);

630     mutex_enter(&vd->vdev_dtl_lock);
631     space_map_close(vd->vdev_dtl_sm);
632     for (int t = 0; t < DTL_TYPES; t++) {
633         range_tree_vacate(vd->vdev_dtl[t], NULL, NULL);
634         range_tree_destroy(vd->vdev_dtl[t]);
635         space_map_unload(&vd->vdev_dtl[t]);
636         space_map_destroy(&vd->vdev_dtl[t]);
637     }
638     mutex_exit(&vd->vdev_dtl_lock);

639     mutex_destroy(&vd->vdev_dtl_lock);
640     mutex_destroy(&vd->vdev_stat_lock);
641     mutex_destroy(&vd->vdev_probe_lock);

642     if (vd == spa->spa_root_vdev)
643         spa->spa_root_vdev = NULL;

644     kmem_free(vd, sizeof (vdev_t));
645 }
646 }

```

unchanged portion omitted

```

802 int
803 vdev metaslab_init(vdev_t *vd, uint64_t txg)
804 {
805     spa_t *spa = vd->vdev_spa;
806     objset_t *mos = spa->spa_meta_objset;
807     uint64_t m;
808     uint64_t oldc = vd->vdev_ms_count;
809     uint64_t newc = vd->vdev_asize >> vd->vdev_ms_shift;
810     metaslab_t **mspp;
811     int error;

```

```

813     ASSERT(txg == 0 || spa_config_held(spa, SCL_ALLOC, RW_WRITER));

815     /*
816      * This vdev is not being allocated from yet or is a hole.
817      */
818     if (vd->vdev_ms_shift == 0)
819         return (0);

821     ASSERT(!vd->vdev_ishole);

823     /*
824      * Compute the raidz-deflation ratio. Note, we hard-code
825      * in 128k (1 << 17) because it is the current "typical" blocksize.
826      * Even if SPA_MAXBLOCKSIZE changes, this algorithm must never change,
827      * or we will inconsistently account for existing bp's.
828      */
829     vd->vdev_deflate_ratio = (1 << 17) /
830         (vdev_psize_to_asize(vd, 1 << 17) >> SPA_MINBLOCKSHIFT);

832     ASSERT(oldc <= newc);

834     mspp = kmem_zalloc(newc * sizeof (*mspp), KM_SLEEP);

836     if (oldc != 0) {
837         bcopy(vd->vdev_ms, mspp, oldc * sizeof (*mspp));
838         kmem_free(vd->vdev_ms, oldc * sizeof (*mspp));
839     }

841     vd->vdev_ms = mspp;
842     vd->vdev_ms_count = newc;

844     for (m = oldc; m < newc; m++) {
845         uint64_t object = 0;

846         space_map_obj_t smo = { 0, 0, 0 };
847         if (txg == 0) {
848             uint64_t object = 0;
849             error = dmu_read(mos, vd->vdev_ms_array,
850                 m * sizeof (uint64_t), sizeof (uint64_t), &object,
851                 DMU_READ_PREFETCH);
852             if (error)
853                 return (error);
854             if (object != 0) {
855                 dmu_buf_t *db;
856                 error = dmu_bonus_hold(mos, object, FTAG, &db);
857                 if (error)
858                     return (error);
859                 ASSERT3U(db->db_size, >=, sizeof (smo));
860                 bcopy(db->db_data, &smo, sizeof (smo));
861                 ASSERT3U(smo.smo_object, ==, object);
862                 dmu_buf_rele(db, FTAG);
863             }
864             vd->vdev_ms[m] = metaslab_init(vd->vdev_mg, m, object, txg);
865         }
866         vd->vdev_ms[m] = metaslab_init(vd->vdev_mg, &smo,
867             m << vd->vdev_ms_shift, 1ULL << vd->vdev_ms_shift, txg);
868     }

869     if (txg == 0)
870         spa_config_enter(spa, SCL_ALLOC, FTAG, RW_WRITER);

871     /*
872      * If the vdev is being removed we don't activate
873      * the metaslabs since we want to ensure that no new
874      * allocations are performed on this device.
875      */

```

```

865     if (oldc == 0 && !vd->vdev_removing)
866         metaslab_group_activate(vd->vdev_mg);
868     if (txg == 0)
869         spa_config_exit(spa, SCL_ALLOC, FTAG);
871     return (0);
872 }

874 void
875 vdev_metaslab_fini(vdev_t *vd)
876 {
877     uint64_t m;
878     uint64_t count = vd->vdev_ms_count;

880     if (vd->vdev_ms != NULL) {
881         metaslab_group_passivate(vd->vdev_mg);
882         for (m = 0; m < count; m++) {
883             metaslab_t *msp = vd->vdev_ms[m];

885                 if (msp != NULL)
886                     metaslab_fini(msp);
887             }
889             for (m = 0; m < count; m++)
890                 if (vd->vdev_ms[m] != NULL)
891                     metaslab_fini(vd->vdev_ms[m]);
888             kmem_free(vd->vdev_ms, count * sizeof (metaslab_t *));
889             vd->vdev_ms = NULL;
890         }
891     }

```

unchanged portion omitted

```

1520 int
1521 vdev_create(vdev_t *vd, uint64_t txg, boolean_t isreplacing)
1522 {
1523     int error;

1525     /*
1526      * Normally, partial opens (e.g. of a mirror) are allowed.
1527      * For a create, however, we want to fail the request if
1528      * there are any components we can't open.
1529      */
1530     error = vdev_open(vd);

1532     if (error || vd->vdev_state != VDEV_STATE_HEALTHY) {
1533         vdev_close(vd);
1534         return (error ? error : ENXIO);
1535     }

1537     /*
1538      * Recursively load DTLs and initialize all labels.
1539      * Recursively initialize all labels.
1540      */
1541     if ((error = vdev_dtl_load(vd)) != 0 ||
1542         (error = vdev_label_init(vd, txg, isreplacing ?
1543             if ((error = vdev_label_init(vd, txg, isreplacing ?
1544                 VDEV_LABEL_REPLACE : VDEV_LABEL_CREATE)) != 0) {
1545                 vdev_close(vd);
1546                 return (error);
1547             }
1548         }

1547     return (0);
1548 }

```

unchanged portion omitted

```
1577 void
```

```

1578 vdev_dirty_leaves(vdev_t *vd, int flags, uint64_t txg)
1579 {
1580     for (int c = 0; c < vd->vdev_children; c++)
1581         vdev_dirty_leaves(vd->vdev_child[c], flags, txg);

1583     if (vd->vdev_ops->vdev_op_leaf)
1584         vdev_dirty(vd->vdev_top, flags, vd, txg);
1585 }

1587 /*
1588  * DTLs.
1589  *
1590  * A vdev's DTL (dirty time log) is the set of transaction groups for which
1591  * the vdev has less than perfect replication. There are four kinds of DTL:
1592  *
1593  * DTL_MISSING: txgs for which the vdev has no valid copies of the data
1594  *
1595  * DTL_PARTIAL: txgs for which data is available, but not fully replicated
1596  *
1597  * DTL_SCRUB: the txgs that could not be repaired by the last scrub; upon
1598  * scrub completion, DTL_SCRUB replaces DTL_MISSING in the range of
1599  * txgs that was scrubbed.
1600  *
1601  * DTL_OUTAGE: txgs which cannot currently be read, whether due to
1602  * persistent errors or just some device being offline.
1603  * Unlike the other three, the DTL_OUTAGE map is not generally
1604  * maintained; it's only computed when needed, typically to
1605  * determine whether a device can be detached.
1606  *
1607  * For leaf vdevs, DTL_MISSING and DTL_PARTIAL are identical: the device
1608  * either has the data or it doesn't.
1609  *
1610  * For interior vdevs such as mirror and RAID-Z the picture is more complex.
1611  * A vdev's DTL_PARTIAL is the union of its children's DTL_PARTIALS, because
1612  * if any child is less than fully replicated, then so is its parent.
1613  * A vdev's DTL_MISSING is a modified union of its children's DTL_MISSINGs,
1614  * comprising only those txgs which appear in 'maxfaults' or more children;
1615  * those are the txgs we don't have enough replication to read. For example,
1616  * double-parity RAID-Z can tolerate up to two missing devices (maxfaults == 2);
1617  * thus, its DTL_MISSING consists of the set of txgs that appear in more than
1618  * two child DTL_MISSING maps.
1619  *
1620  * It should be clear from the above that to compute the DTLs and outage maps
1621  * for all vdevs, it suffices to know just the leaf vdevs' DTL_MISSING maps.
1622  * Therefore, that is all we keep on disk. When loading the pool, or after
1623  * a configuration change, we generate all other DTLs from first principles.
1624  */
1625 void
1626 vdev_dtl_dirty(vdev_t *vd, vdev_dtl_type_t t, uint64_t txg, uint64_t size)
1627 {
1628     range_tree_t *rt = vd->vdev_dtl[t];
1629     space_map_t *sm = &vd->vdev_dtl[t];

1630     ASSERT(t < DTL_TYPES);
1631     ASSERT(vd != vd->vdev_spa->spa_root_vdev);
1632     ASSERT(spa_writable(vd->vdev_spa));

1634     mutex_enter(rt->rt_lock);
1635     if (!range_tree_contains(rt, txg, size))
1636         range_tree_add(rt, txg, size);
1637     mutex_exit(rt->rt_lock);
1638     mutex_enter(sm->sm_lock);
1639     if (!space_map_contains(sm, txg, size))
1640         space_map_add(sm, txg, size);
1641     mutex_exit(sm->sm_lock);
1642 }

```

```

1640 boolean_t
1641 vdev_dtl_contains(vdev_t *vd, vdev_dtl_type_t t, uint64_t txg, uint64_t size)
1642 {
1643     range_tree_t *rt = vd->vdev_dtl[t];
1644     space_map_t *sm = &vd->vdev_dtl[t];
1645     boolean_t dirty = B_FALSE;
1646
1647     ASSERT(t < DTL_TYPES);
1648     ASSERT(vd != vd->vdev_spa->spa_root_vdev);
1649
1650     mutex_enter(rt->rt_lock);
1651     if (range_tree_space(rt) != 0)
1652         dirty = range_tree_contains(rt, txg, size);
1653     mutex_exit(rt->rt_lock);
1654     mutex_enter(sm->sm_lock);
1655     if (sm->sm_space != 0)
1656         dirty = space_map_contains(sm, txg, size);
1657     mutex_exit(sm->sm_lock);
1658
1659     return (dirty);
1660 }
1661
1662 boolean_t
1663 vdev_dtl_empty(vdev_t *vd, vdev_dtl_type_t t)
1664 {
1665     range_tree_t *rt = vd->vdev_dtl[t];
1666     space_map_t *sm = &vd->vdev_dtl[t];
1667     boolean_t empty;
1668
1669     mutex_enter(rt->rt_lock);
1670     empty = (range_tree_space(rt) == 0);
1671     mutex_exit(rt->rt_lock);
1672     mutex_enter(sm->sm_lock);
1673     empty = (sm->sm_space == 0);
1674     mutex_exit(sm->sm_lock);
1675
1676     return (empty);
1677 }
1678
1679 /*
1680  * Returns the lowest txg in the DTL range.
1681  */
1682 static uint64_t
1683 vdev_dtl_min(vdev_t *vd)
1684 {
1685     range_seg_t *rs;
1686     space_seg_t *ss;
1687
1688     ASSERT(MUTEX_HELD(&vd->vdev_dtl_lock));
1689     ASSERT3U(range_tree_space(vd->vdev_dtl[DTL_MISSING]), !=, 0);
1690     ASSERT3U(vd->vdev_dtl[DTL_MISSING].sm_space, !=, 0);
1691     ASSERT0(vd->vdev_children);
1692
1693     rs = avl_first(&vd->vdev_dtl[DTL_MISSING]->rt_root);
1694     return (rs->rs_start - 1);
1695 }
1696
1697 /*
1698  * Returns the highest txg in the DTL.
1699  */
1700 static uint64_t
1701 vdev_dtl_max(vdev_t *vd)
1702 {

```

```

1692     range_seg_t *rs;
1693     space_seg_t *ss;
1694
1695     ASSERT(MUTEX_HELD(&vd->vdev_dtl_lock));
1696     ASSERT3U(range_tree_space(vd->vdev_dtl[DTL_MISSING]), !=, 0);
1697     ASSERT3U(vd->vdev_dtl[DTL_MISSING].sm_space, !=, 0);
1698     ASSERT0(vd->vdev_children);
1699
1700     rs = avl_last(&vd->vdev_dtl[DTL_MISSING]->rt_root);
1701     return (rs->rs_end);
1702 }
1703
1704 /*
1705  * Determine if a resilvering vdev should remove any DTL entries from
1706  * its range. If the vdev was resilvering for the entire duration of the
1707  * scan then it should excise that range from its DTLs. Otherwise, this
1708  * vdev is considered partially resilvered and should leave its DTL
1709  * entries intact. The comment in vdev_dtl_reassess() describes how we
1710  * excise the DTLs.
1711  */
1712 static boolean_t
1713 vdev_dtl_should_excise(vdev_t *vd)
1714 {
1715     spa_t *spa = vd->vdev_spa;
1716     dsl_scan_t *scn = spa->spa_dsl_pool->dp_scan;
1717
1718     ASSERT0(scn->scn_phys.scn_errors);
1719     ASSERT0(vd->vdev_children);
1720
1721     if (vd->vdev_resilver_txg == 0 ||
1722         range_tree_space(vd->vdev_dtl[DTL_MISSING]) == 0 ||
1723         vd->vdev_dtl[DTL_MISSING].sm_space == 0)
1724         return (B_TRUE);
1725
1726     /*
1727      * When a resilver is initiated the scan will assign the scn_max_txg
1728      * value to the highest txg value that exists in all DTLs. If this
1729      * device's max DTL is not part of this scan (i.e. it is not in
1730      * the range (scn_min_txg, scn_max_txg) then it is not eligible
1731      * for excision.
1732      */
1733     if (vdev_dtl_max(vd) <= scn->scn_phys.scn_max_txg) {
1734         ASSERT3U(scn->scn_phys.scn_min_txg, <=, vdev_dtl_min(vd));
1735         ASSERT3U(scn->scn_phys.scn_min_txg, <, vd->vdev_resilver_txg);
1736         ASSERT3U(vd->vdev_resilver_txg, <=, scn->scn_phys.scn_max_txg);
1737         return (B_TRUE);
1738     }
1739     return (B_FALSE);
1740 }
1741
1742 /*
1743  * Reassess DTLs after a config change or scrub completion.
1744  */
1745 void
1746 vdev_dtl_reassess(vdev_t *vd, uint64_t txg, uint64_t scrub_txg, int scrub_done)
1747 {
1748     spa_t *spa = vd->vdev_spa;
1749     avl_tree_t reftree;
1750     int minref;
1751
1752     ASSERT(spa_config_held(spa, SCL_ALL, RW_READER) != 0);
1753
1754     for (int c = 0; c < vd->vdev_children; c++)
1755         vdev_dtl_reassess(vd->vdev_child[c], txg,

```



```

1753         scrub_txg, scrub_done);
1755     if (vd == spa->spa_root_vdev || vd->vdev_ishole || vd->vdev_aux)
1756         return;
1758     if (vd->vdev_ops->vdev_op_leaf) {
1759         dsl_scan_t *scn = spa->spa_dsl_pool->dp_scan;
1761         mutex_enter(&vd->vdev_dtl_lock);
1763         /*
1764          * If we've completed a scan cleanly then determine
1765          * if this vdev should remove any DTLs. We only want to
1766          * excise regions on vdevs that were available during
1767          * the entire duration of this scan.
1768          */
1769         if (scrub_txg != 0 &&
1770             (spa->spa_scrub_started ||
1771              (scn != NULL && scn->scn_phys.scn_errors == 0)) &&
1772             vdev_dtl_should_excise(vd)) {
1773             /*
1774              * We completed a scrub up to scrub_txg. If we
1775              * did it without rebooting, then the scrub dtl
1776              * will be valid, so excise the old region and
1777              * fold in the scrub dtl. Otherwise, leave the
1778              * dtl as-is if there was an error.
1779              *
1780              * There's little trick here: to excise the beginning
1781              * of the DTL_MISSING map, we put it into a reference
1782              * tree and then add a segment with refcnt -1 that
1783              * covers the range [0, scrub_txg). This means
1784              * that each txg in that range has refcnt -1 or 0.
1785              * We then add DTL_SCRUB with a refcnt of 2, so that
1786              * entries in the range [0, scrub_txg) will have a
1787              * positive refcnt -- either 1 or 2. We then convert
1788              * the reference tree into the new DTL_MISSING map.
1789              */
1790             space_reftree_create(&reftree);
1791             space_reftree_add_map(&reftree,
1792                                 vd->vdev_dtl[DTL_MISSING], 1);
1793             space_reftree_add_seg(&reftree, 0, scrub_txg, -1);
1794             space_reftree_add_map(&reftree,
1795                                 vd->vdev_dtl[DTL_SCRUB], 2);
1796             space_reftree_generate_map(&reftree,
1797                                       vd->vdev_dtl[DTL_MISSING], 1);
1798             space_reftree_destroy(&reftree);
1799             space_map_ref_create(&reftree);
1800             space_map_ref_add_map(&reftree,
1801                                  &vd->vdev_dtl[DTL_MISSING], 1);
1802             space_map_ref_add_seg(&reftree, 0, scrub_txg, -1);
1803             space_map_ref_add_map(&reftree,
1804                                  &vd->vdev_dtl[DTL_SCRUB], 2);
1805             space_map_ref_generate_map(&reftree,
1806                                       &vd->vdev_dtl[DTL_MISSING], 1);
1807             space_map_ref_destroy(&reftree);
1808         }
1809         range_tree_vacate(vd->vdev_dtl[DTL_PARTIAL], NULL, NULL);
1810         range_tree_walk(vd->vdev_dtl[DTL_MISSING],
1811                        range_tree_add, vd->vdev_dtl[DTL_PARTIAL]);
1812         space_map_vacate(&vd->vdev_dtl[DTL_PARTIAL], NULL, NULL);
1813         space_map_walk(&vd->vdev_dtl[DTL_MISSING],
1814                       space_map_add, &vd->vdev_dtl[DTL_PARTIAL]);
1815         if (scrub_done)
1816             range_tree_vacate(vd->vdev_dtl[DTL_SCRUB], NULL, NULL);
1817         range_tree_vacate(vd->vdev_dtl[DTL_OUTAGE], NULL, NULL);
1818         space_map_vacate(&vd->vdev_dtl[DTL_SCRUB], NULL, NULL);

```

```

1800         space_map_vacate(&vd->vdev_dtl[DTL_OUTAGE], NULL, NULL);
1801         if (!vdev_readable(vd))
1802             range_tree_add(vd->vdev_dtl[DTL_OUTAGE], 0, -1ULL);
1803         space_map_add(&vd->vdev_dtl[DTL_OUTAGE], 0, -1ULL);
1804     } else
1805         range_tree_walk(vd->vdev_dtl[DTL_MISSING],
1806                        range_tree_add, vd->vdev_dtl[DTL_OUTAGE]);
1807         space_map_walk(&vd->vdev_dtl[DTL_MISSING],
1808                       space_map_add, &vd->vdev_dtl[DTL_OUTAGE]);
1809     }
1810     /*
1811      * If the vdev was resilvering and no longer has any
1812      * DTLs then reset its resilvering flag.
1813      */
1814     if (vd->vdev_resilver_txg != 0 &&
1815         range_tree_space(vd->vdev_dtl[DTL_MISSING]) == 0 &&
1816         range_tree_space(vd->vdev_dtl[DTL_OUTAGE]) == 0)
1817         vd->vdev_dtl[DTL_MISSING].sm_space == 0 &&
1818         vd->vdev_dtl[DTL_OUTAGE].sm_space == 0)
1819         vd->vdev_resilver_txg = 0;
1820     mutex_exit(&vd->vdev_dtl_lock);
1821     if (txg != 0)
1822         vdev_dirty(vd->vdev_top, VDD_DTL, vd, txg);
1823     return;
1824 }
1825
1826 mutex_enter(&vd->vdev_dtl_lock);
1827 for (int t = 0; t < DTL_TYPES; t++) {
1828     /* account for child's outage in parent's missing map */
1829     int s = (t == DTL_MISSING) ? DTL_OUTAGE : t;
1830     if (t == DTL_SCRUB)
1831         continue; /* leaf vdevs only */
1832     if (t == DTL_PARTIAL)
1833         minref = 1; /* i.e. non-zero */
1834     else if (vd->vdev_nparity != 0)
1835         minref = vd->vdev_nparity + 1; /* RAID-Z */
1836     else
1837         minref = vd->vdev_children; /* any kind of mirror */
1838     space_reftree_create(&reftree);
1839     space_map_ref_create(&reftree);
1840     for (int c = 0; c < vd->vdev_children; c++) {
1841         vdev_t *cvd = vd->vdev_child[c];
1842         mutex_enter(&cvd->vdev_dtl_lock);
1843         space_reftree_add_map(&reftree, cvd->vdev_dtl[s], 1);
1844         space_map_ref_add_map(&reftree, &cvd->vdev_dtl[s], 1);
1845         mutex_exit(&cvd->vdev_dtl_lock);
1846     }
1847     space_reftree_generate_map(&reftree, vd->vdev_dtl[t], minref);
1848     space_reftree_destroy(&reftree);
1849     space_map_ref_generate_map(&reftree, &vd->vdev_dtl[t], minref);
1850     space_map_ref_destroy(&reftree);
1851 }
1852 mutex_exit(&vd->vdev_dtl_lock);
1853
1854 int
1855 static int
1856 vdev_dtl_load(vdev_t *vd)
1857 {
1858     spa_t *spa = vd->vdev_spa;
1859     space_map_obj_t *smo = &vd->vdev_dtl_smo;
1860     objset_t *mos = spa->spa_meta_objset;
1861     int error = 0;
1862     dmu_buf_t *db;

```

```

1855     int error;

1860     if (vd->vdev_ops->vdev_op_leaf && vd->vdev_dtl_object != 0) {
1857         ASSERT(vd->vdev_children == 0);

1859         if (smo->smo_object == 0)
1860             return (0);

1861         ASSERT(!vd->vdev_ishole);

1863         error = space_map_open(&vd->vdev_dtl_sm, mos,
1864                               vd->vdev_dtl_object, 0, -1ULL, 0, &vd->vdev_dtl_lock);
1865         if (error)
1864             if ((error = dmu_bonus_hold(mos, smo->smo_object, FTAG, &db)) != 0)
1866                 return (error);
1867         ASSERT(vd->vdev_dtl_sm != NULL);

1869         mutex_enter(&vd->vdev_dtl_lock);
1867         ASSERT3U(db->db_size, >=, sizeof (*smo));
1868         bcopy(db->db_data, smo, sizeof (*smo));
1869         dmu_buf_rele(db, FTAG);

1871         /*
1872          * Now that we've opened the space_map we need to update
1873          * the in-core DTL.
1874          */
1875         space_map_update(vd->vdev_dtl_sm);

1877         error = space_map_load(vd->vdev_dtl_sm,
1878                               vd->vdev_dtl[DTL_MISSING], SM_ALLOC);
1871         mutex_enter(&vd->vdev_dtl_lock);
1872         error = space_map_load(&vd->vdev_dtl[DTL_MISSING],
1873                               NULL, SM_ALLOC, smo, mos);
1879         mutex_exit(&vd->vdev_dtl_lock);

1881         return (error);
1882     }

1884     for (int c = 0; c < vd->vdev_children; c++) {
1885         error = vdev_dtl_load(vd->vdev_child[c]);
1886         if (error != 0)
1887             break;
1888     }

1890     return (error);
1891 }

1893 void
1894 vdev_dtl_sync(vdev_t *vd, uint64_t txg)
1895 {
1896     spa_t *spa = vd->vdev_spa;
1897     range_tree_t *rt = vd->vdev_dtl[DTL_MISSING];
1883     space_map_obj_t *smo = &vd->vdev_dtl_smo;
1884     space_map_t *sm = &vd->vdev_dtl[DTL_MISSING];
1898     objset_t *mos = spa->spa_meta_objset;
1899     range_tree_t *rtsync;
1900     kmutex_t rtlock;
1886     space_map_t smsync;
1887     kmutex_t smlock;
1888     dmu_buf_t *db;
1901     dmu_tx_t *tx;
1902     uint64_t object = space_map_object(vd->vdev_dtl_sm);

1904     ASSERT(!vd->vdev_ishole);
1905     ASSERT(vd->vdev_ops->vdev_op_leaf);

```

```

1907     tx = dmu_tx_create_assigned(spa->spa_dsl_pool, txg);

1909     if (vd->vdev_detached || vd->vdev_top->vdev_removing) {
1910         mutex_enter(&vd->vdev_dtl_lock);
1911         space_map_free(vd->vdev_dtl_sm, tx);
1912         space_map_close(vd->vdev_dtl_sm);
1913         vd->vdev_dtl_sm = NULL;
1914         mutex_exit(&vd->vdev_dtl_lock);
1895     if (vd->vdev_detached) {
1896         if (smo->smo_object != 0) {
1897             int err = dmu_object_free(mos, smo->smo_object, tx);
1898             ASSERT0(err);
1899             smo->smo_object = 0;
1900         }
1915         dmu_tx_commit(tx);
1916         return;
1917     }

1919     if (vd->vdev_dtl_sm == NULL) {
1920         uint64_t new_object;

1922         new_object = space_map_alloc(mos, tx);
1923         VERIFY3U(new_object, !=, 0);

1925         VERIFY0(space_map_open(&vd->vdev_dtl_sm, mos, new_object,
1926                               0, -1ULL, 0, &vd->vdev_dtl_lock));
1927         ASSERT(vd->vdev_dtl_sm != NULL);
1905     if (smo->smo_object == 0) {
1906         ASSERT(smo->smo_objsize == 0);
1907         ASSERT(smo->smo_alloc == 0);
1908         smo->smo_object = dmu_object_alloc(mos,
1909                                           DMU_OT_SPACE_MAP, 1 << SPACE_MAP_BLOCKSHIFT,
1910                                           DMU_OT_SPACE_MAP_HEADER, sizeof (*smo), tx);
1911         ASSERT(smo->smo_object != 0);
1912         vdev_config_dirty(vd->vdev_top);
1928     }

1930     mutex_init(&rtlock, NULL, MUTEX_DEFAULT, NULL);
1915     mutex_init(&smlock, NULL, MUTEX_DEFAULT, NULL);

1932     rtsync = range_tree_create(NULL, NULL, &rtlock);
1917     space_map_create(&smsync, sm->sm_start, sm->sm_size, sm->sm_shift,
1918                    &smlock);

1934     mutex_enter(&rtlock);
1920     mutex_enter(&smlock);

1936     mutex_enter(&vd->vdev_dtl_lock);
1937     range_tree_walk(rt, range_tree_add, rtsync);
1923     space_map_walk(sm, space_map_add, &smsync);
1938     mutex_exit(&vd->vdev_dtl_lock);

1940     space_map_truncate(vd->vdev_dtl_sm, tx);
1941     space_map_write(vd->vdev_dtl_sm, rtsync, SM_ALLOC, tx);
1942     range_tree_vacate(rtsync, NULL, NULL);
1926     space_map_truncate(smo, mos, tx);
1927     space_map_sync(&smsync, SM_ALLOC, smo, mos, tx);
1928     space_map_vacate(&smsync, NULL, NULL);

1944     range_tree_destroy(rtsync);
1930     space_map_destroy(&smsync);

1946     mutex_exit(&rtlock);
1947     mutex_destroy(&rtlock);
1932     mutex_exit(&smlock);
1933     mutex_destroy(&smlock);

```

```

1949  /*
1950  * If the object for the space map has changed then dirty
1951  * the top level so that we update the config.
1952  */
1953  if (object != space_map_object(vd->vdev_dtl_sm)) {
1954      zfs_dbgmsg("txg %llu, spa %s, DTL old object %llu, "
1955               "new object %llu", txg, spa_name(spa), object,
1956               space_map_object(vd->vdev_dtl_sm));
1957      vdev_config_dirty(vd->vdev_top);
1958  }
1959  VERIFY(0 == dmu_bonus_hold(mos, smo->smo_object, FTAG, &db));
1960  dmu_buf_will_dirty(db, tx);
1961  ASSERT3U(db->db_size, >=, sizeof (*smo));
1962  bcopy(smo, db->db_data, sizeof (*smo));
1963  dmu_buf_rele(db, FTAG);
1964
1965  dmu_tx_commit(tx);
1966
1967  mutex_enter(&vd->vdev_dtl_lock);
1968  space_map_update(vd->vdev_dtl_sm);
1969  mutex_exit(&vd->vdev_dtl_lock);
1970 }
1971
1972 unchanged portion omitted
1973
2001 /*
2002 * Determine if resilver is needed, and if so the txg range.
2003 */
2004 boolean_t
2005 vdev_resilver_needed(vdev_t *vd, uint64_t *minp, uint64_t *maxp)
2006 {
2007     boolean_t needed = B_FALSE;
2008     uint64_t thismin = UINT64_MAX;
2009     uint64_t thismax = 0;
2010
2011     if (vd->vdev_children == 0) {
2012         mutex_enter(&vd->vdev_dtl_lock);
2013         if (range_tree_space(vd->vdev_dtl[DTL_MISSING]) != 0 &&
2014             if (vd->vdev_dtl[DTL_MISSING].sm_space != 0 &&
2015                 vdev_writable(vd)) {
2016             thismin = vdev_dtl_min(vd);
2017             thismax = vdev_dtl_max(vd);
2018             needed = B_TRUE;
2019         }
2020         mutex_exit(&vd->vdev_dtl_lock);
2021     } else {
2022         for (int c = 0; c < vd->vdev_children; c++) {
2023             vdev_t *cvd = vd->vdev_child[c];
2024             uint64_t cmin, cmax;
2025
2026             if (vdev_resilver_needed(cvd, &cmin, &cmax)) {
2027                 thismin = MIN(thismin, cmin);
2028                 thismax = MAX(thismax, cmax);
2029                 needed = B_TRUE;
2030             }
2031         }
2032     }
2033
2034     if (needed && minp) {
2035         *minp = thismin;
2036         *maxp = thismax;
2037     }
2038     return (needed);
2039 }
2040
2041 unchanged portion omitted

```

```

2109 void
2110 vdev_remove(vdev_t *vd, uint64_t txg)
2111 {
2112     spa_t *spa = vd->vdev_spa;
2113     objset_t *mos = spa->spa_meta_objset;
2114     dmu_tx_t *tx;
2115
2116     tx = dmu_tx_create_assigned(spa_get_dsl(spa), txg);
2117
2118     if (vd->vdev_dtl_smo.smo_object) {
2119         ASSERT0(vd->vdev_dtl_smo.smo_alloc);
2120         (void) dmu_object_free(mos, vd->vdev_dtl_smo.smo_object, tx);
2121         vd->vdev_dtl_smo.smo_object = 0;
2122     }
2123
2124     if (vd->vdev_ms != NULL) {
2125         for (int m = 0; m < vd->vdev_ms_count; m++) {
2126             metaslab_t *msp = vd->vdev_ms[m];
2127
2128             if (msp == NULL || msp->ms_sm == NULL)
2129                 continue;
2130             if (msp == NULL || msp->ms_smo.smo_object == 0)
2131                 continue;
2132
2133             mutex_enter(&msp->ms_lock);
2134             VERIFY0(space_map_allocated(msp->ms_sm));
2135             space_map_free(msp->ms_sm, tx);
2136             space_map_close(msp->ms_sm);
2137             msp->ms_sm = NULL;
2138             mutex_exit(&msp->ms_lock);
2139             ASSERT0(msp->ms_smo.smo_alloc);
2140             (void) dmu_object_free(mos, msp->ms_smo.smo_object, tx);
2141             msp->ms_smo.smo_object = 0;
2142         }
2143     }
2144
2145     if (vd->vdev_ms_array) {
2146         (void) dmu_object_free(mos, vd->vdev_ms_array, tx);
2147         vd->vdev_ms_array = 0;
2148         vd->vdev_ms_shift = 0;
2149     }
2150     dmu_tx_commit(tx);
2151 }
2152
2153 unchanged portion omitted

```

```

*****
37492 Tue Sep 3 20:27:14 2013
new/usr/src/uts/common/fs/zfs/vdev_label.c
4101 metaslab_debug should allow for fine-grained control
4102 space_maps should store more information about themselves
4103 space_map object blocksize should be increased
4104 ::spa_space no longer works
4105 removing a mirrored log device results in a leaked object
4106 asynchronously load metaslab
Reviewed by: Matthew Ahrens <mahrens@delphix.com>
Reviewed by: Adam Leventhal <ahl@delphix.com>
Reviewed by: Sebastien Roy <seb@delphix.com>
*****
unchanged_portion_omitted

210 /*
211  * Generate the nvlist representing this vdev's config.
212  */
213 nvlist_t *
214 vdev_config_generate(spa_t *spa, vdev_t *vd, boolean_t getstats,
215     vdev_config_flag_t flags)
216 {
217     nvlist_t *nv = NULL;
218
219     nv = fnvlist_alloc();
220
221     fnvlist_add_string(nv, ZPOOL_CONFIG_TYPE, vd->vdev_ops->vdev_op_type);
222     if (!(flags & (VDEV_CONFIG_SPARE | VDEV_CONFIG_L2CACHE)))
223         fnvlist_add_uint64(nv, ZPOOL_CONFIG_ID, vd->vdev_id);
224     fnvlist_add_uint64(nv, ZPOOL_CONFIG_GUID, vd->vdev_guid);
225
226     if (vd->vdev_path != NULL)
227         fnvlist_add_string(nv, ZPOOL_CONFIG_PATH, vd->vdev_path);
228
229     if (vd->vdev_devid != NULL)
230         fnvlist_add_string(nv, ZPOOL_CONFIG_DEVID, vd->vdev_devid);
231
232     if (vd->vdev_physpath != NULL)
233         fnvlist_add_string(nv, ZPOOL_CONFIG_PHYS_PATH,
234             vd->vdev_physpath);
235
236     if (vd->vdev_fru != NULL)
237         fnvlist_add_string(nv, ZPOOL_CONFIG_FRU, vd->vdev_fru);
238
239     if (vd->vdev_nparity != 0) {
240         ASSERT(strcmp(vd->vdev_ops->vdev_op_type,
241             VDEV_TYPE_RAIDZ) == 0);
242
243         /*
244          * Make sure someone hasn't managed to sneak a fancy new vdev
245          * into a crufty old storage pool.
246          */
247         ASSERT(vd->vdev_nparity == 1 ||
248             (vd->vdev_nparity <= 2 &&
249             spa_version(spa) >= SPA_VERSION_RAIDZ2) ||
250             (vd->vdev_nparity <= 3 &&
251             spa_version(spa) >= SPA_VERSION_RAIDZ3));
252
253         /*
254          * Note that we'll add the nparity tag even on storage pools
255          * that only support a single parity device -- older software
256          * will just ignore it.
257          */
258         fnvlist_add_uint64(nv, ZPOOL_CONFIG_NPARITY, vd->vdev_nparity);
259     }

```

```

261     if (vd->vdev_wholedisk != -1ULL)
262         fnvlist_add_uint64(nv, ZPOOL_CONFIG_WHOLE_DISK,
263             vd->vdev_wholedisk);
264
265     if (vd->vdev_not_present)
266         fnvlist_add_uint64(nv, ZPOOL_CONFIG_NOT_PRESENT, 1);
267
268     if (vd->vdev_isspare)
269         fnvlist_add_uint64(nv, ZPOOL_CONFIG_IS_SPARE, 1);
270
271     if (!(flags & (VDEV_CONFIG_SPARE | VDEV_CONFIG_L2CACHE)) &&
272         vd == vd->vdev_top) {
273         fnvlist_add_uint64(nv, ZPOOL_CONFIG_METASLAB_ARRAY,
274             vd->vdev_ms_array);
275         fnvlist_add_uint64(nv, ZPOOL_CONFIG_METASLAB_SHIFT,
276             vd->vdev_ms_shift);
277         fnvlist_add_uint64(nv, ZPOOL_CONFIG_ASHIFT, vd->vdev_ashift);
278         fnvlist_add_uint64(nv, ZPOOL_CONFIG_ASIZE,
279             vd->vdev_asize);
280         fnvlist_add_uint64(nv, ZPOOL_CONFIG_IS_LOG, vd->vdev_islog);
281         if (vd->vdev_removing)
282             fnvlist_add_uint64(nv, ZPOOL_CONFIG_REMOVING,
283                 vd->vdev_removing);
284     }
285
286     if (vd->vdev_dtl_sm != NULL) {
287         if (vd->vdev_dtl_smo.smo_object != 0)
288             fnvlist_add_uint64(nv, ZPOOL_CONFIG_DTL,
289                 space_map_object(vd->vdev_dtl_sm));
290         fnvlist_add_uint64(nv, ZPOOL_CONFIG_DTL_SMO,
291             vd->vdev_dtl_smo.smo_object);
292     }
293
294     if (vd->vdev_crtxg)
295         fnvlist_add_uint64(nv, ZPOOL_CONFIG_CREATE_TXG, vd->vdev_crtxg);
296
297     if (getstats) {
298         vdev_stat_t vs;
299         pool_scan_stat_t ps;
300
301         vdev_get_stats(vd, &vs);
302         fnvlist_add_uint64_array(nv, ZPOOL_CONFIG_VDEV_STATS,
303             (uint64_t *)&vs, sizeof (vs) / sizeof (uint64_t));
304
305         /* provide either current or previous scan information */
306         if (spa_scan_get_stats(spa, &ps) == 0) {
307             fnvlist_add_uint64_array(nv,
308                 ZPOOL_CONFIG_SCAN_STATS, (uint64_t *)&ps,
309                 sizeof (pool_scan_stat_t) / sizeof (uint64_t));
310         }
311     }
312
313     if (!vd->vdev_ops->vdev_op_leaf) {
314         nvlist_t **child;
315         int c, idx;
316
317         ASSERT(!vd->vdev_ishole);
318
319         child = kmem_alloc(vd->vdev_children * sizeof (nvlist_t *),
320             KM_SLEEP);
321
322         for (c = 0, idx = 0; c < vd->vdev_children; c++) {
323             vdev_t *cvd = vd->vdev_child[c];
324
325             /*
326              * If we're generating an nvlist of removing
327              * vdevs then skip over any device which is

```

```

325         * not being removed.
326         */
327         if ((flags & VDEV_CONFIG_REMOVING) &&
328             !cvd->vdev_removing)
329             continue;
331
332         child[idx++] = vdev_config_generate(spa, cvd,
333             getstats, flags);
334     }
335
336     if (idx) {
337         fnvlist_add_nvlist_array(nv, ZPOOL_CONFIG_CHILDREN,
338             child, idx);
339     }
340
341     for (c = 0; c < idx; c++)
342         nvlist_free(child[c]);
343
344     kmem_free(child, vd->vdev_children * sizeof (nvlist_t *));
345 } else {
346     const char *aux = NULL;
347
348     if (vd->vdev_offline && !vd->vdev_tmpoffline)
349         fnvlist_add_uint64(nv, ZPOOL_CONFIG_OFFLINE, B_TRUE);
350     if (vd->vdev_resilver_txc != 0)
351         fnvlist_add_uint64(nv, ZPOOL_CONFIG_RESILVER_TXG,
352             vd->vdev_resilver_txc);
353     if (vd->vdev_faulted)
354         fnvlist_add_uint64(nv, ZPOOL_CONFIG_FAULTED, B_TRUE);
355     if (vd->vdev_degraded)
356         fnvlist_add_uint64(nv, ZPOOL_CONFIG_DEGRADED, B_TRUE);
357     if (vd->vdev_removed)
358         fnvlist_add_uint64(nv, ZPOOL_CONFIG_REMOVED, B_TRUE);
359     if (vd->vdev_unspare)
360         fnvlist_add_uint64(nv, ZPOOL_CONFIG_UNSPARE, B_TRUE);
361     if (vd->vdev_ishole)
362         fnvlist_add_uint64(nv, ZPOOL_CONFIG_IS_HOLE, B_TRUE);
363
364     switch (vd->vdev_stat.vs_aux) {
365     case VDEV_AUX_ERR_EXCEEDED:
366         aux = "err_exceeded";
367         break;
368
369     case VDEV_AUX_EXTERNAL:
370         aux = "external";
371         break;
372     }
373
374     if (aux != NULL)
375         fnvlist_add_string(nv, ZPOOL_CONFIG_AUX_STATE, aux);
376
377     if (vd->vdev_splitting && vd->vdev_orig_guid != 0LL) {
378         fnvlist_add_uint64(nv, ZPOOL_CONFIG_ORIG_GUID,
379             vd->vdev_orig_guid);
380     }
381 }
382
383     return (nv);
384 }

```

unchanged portion omitted

```
*****
```

```
14617 Tue Sep 3 20:27:15 2013
```

```
new/usr/src/uts/common/fs/zfs/zfeature.c
```

```
4101 metaslab_debug should allow for fine-grained control
```

```
4102 space_maps should store more information about themselves
```

```
4103 space_map object blocksize should be increased
```

```
4104 ::spa_space no longer works
```

```
4105 removing a mirrored log device results in a leaked object
```

```
4106 asynchronously load metaslab
```

```
Reviewed by: Matthew Ahrens <mahrens@delphix.com>
```

```
Reviewed by: Adam Leventhal <ahl@delphix.com>
```

```
Reviewed by: Sebastien Roy <seb@delphix.com>
```

```
*****
```

```
unchanged_portion_omitted
```

```
364 /*
365  * If the specified feature has not yet been enabled, this function returns
366  * ENOTSUP; otherwise, this function increments the feature's refcount (or
367  * returns EOVERFLOW if the refcount cannot be incremented). This function must
368  * be called from syncing context.
369  */
370 void
371 spa_feature_incr(spa_t *spa, zfeature_info_t *feature, dmu_tx_t *tx)
372 {
373     ASSERT(dmu_tx_is_syncing(tx));
374     ASSERT3U(spa_version(spa), >=, SPA_VERSION_FEATURES);
375     VERIFY3U(0, ==, feature_do_action(spa->spa_meta_objset,
376         spa->spa_feat_for_read_obj, spa->spa_feat_for_write_obj,
377         spa->spa_feat_desc_obj, feature, FEATURE_ACTION_INCR, tx));
378 }
```

```
379 /*
380  * If the specified feature has not yet been enabled, this function returns
381  * ENOTSUP; otherwise, this function decrements the feature's refcount (or
382  * returns EOVERFLOW if the refcount is already 0). This function must
383  * be called from syncing context.
384  */
385 void
386 spa_feature_decr(spa_t *spa, zfeature_info_t *feature, dmu_tx_t *tx)
387 {
388     ASSERT(dmu_tx_is_syncing(tx));
389     ASSERT3U(spa_version(spa), >=, SPA_VERSION_FEATURES);
390     VERIFY3U(0, ==, feature_do_action(spa->spa_meta_objset,
391         spa->spa_feat_for_read_obj, spa->spa_feat_for_write_obj,
392         spa->spa_feat_desc_obj, feature, FEATURE_ACTION_DECR, tx));
393 }
```

```
384 /*
385  * This interface is for debugging only. Normal consumers should use
386  * spa_feature_is_enabled/spa_feature_is_active.
387  */
388 int
389 spa_feature_get_refcount(spa_t *spa, zfeature_info_t *feature)
390 {
391     int err;
392     uint64_t refcount;
393
394     if (spa_version(spa) < SPA_VERSION_FEATURES)
395         return (B_FALSE);
396
397     err = feature_get_refcount(spa->spa_meta_objset,
398         spa->spa_feat_for_read_obj, spa->spa_feat_for_write_obj,
399         feature, &refcount);
400     ASSERT(err == 0 || err == ENOTSUP);
401     return (err == 0 ? refcount : 0);
402 }
```

```
404 boolean_t
405 spa_feature_is_enabled(spa_t *spa, zfeature_info_t *feature)
406 {
407     int err;
408     uint64_t refcount;
409
410     if (spa_version(spa) < SPA_VERSION_FEATURES)
411         return (B_FALSE);
412
413     err = feature_get_refcount(spa->spa_meta_objset,
414         spa->spa_feat_for_read_obj, spa->spa_feat_for_write_obj,
415         feature, &refcount);
416     ASSERT(err == 0 || err == ENOTSUP);
417     return (err == 0);
418 }
```

```
unchanged_portion_omitted
```