

```

*****
77503 Thu Apr 25 16:14:55 2013
new/usr/src/cmd/mdb/common/modules/zfs/zfs.c
3752 want more verifiable dbuf user eviction
Submitted by: Justin Gibbs <justing@spectralogic.com>
Submitted by: Will Andrews <willa@spectralogic.com>
*****
unchanged_portion_omitted

511 #define CHAIN_END 0xffff
512 /*
513 * ::zap_leaf [-v]
514 *
515 * Print a zap_leaf_phys_t, assumed to be 16k
516 */
517 /* ARGSUSED */
518 static int
519 zap_leaf(uintptr_t addr, uint_t flags, int argc, const mdb_arg_t *argv)
520 {
521     char buf[16*1024];
522     int verbose = B_FALSE;
523     int four = B_FALSE;
524     zap_leaf_t l = { 0 };
525     zap_leaf_phys_t *zlp = (void *)buf;
526     int i;

528     if (mdb_getopts(argc, argv,
529         'v', MDB_OPT_SETBITS, TRUE, &verbose,
530         '4', MDB_OPT_SETBITS, TRUE, &four,
531         NULL) != argc)
532         return (DCMD_USAGE);

534     l.l_phys = zlp;
535     l.l_bs = 14; /* assume 16k blocks */
536     if (four)
537         l.l_bs = 12;

539     if (!(flags & DCMD_ADDRSPEC)) {
540         return (DCMD_USAGE);
541     }

543     if (mdb_vread(buf, sizeof (buf), addr) == -1) {
544         mdb_warn("failed to read zap_leaf_phys_t at %p", addr);
545         return (DCMD_ERR);
546     }

548     if (zlp->l_hdr.lh_block_type != ZBT_LEAF ||
549         zlp->l_hdr.lh_magic != ZAP_LEAF_MAGIC) {
550         mdb_warn("This does not appear to be a zap_leaf_phys_t");
551         return (DCMD_ERR);
552     }

554     mdb_printf("zap_leaf_phys_t at %p:\n", addr);
555     mdb_printf("    lh_prefix_len = %u\n", zlp->l_hdr.lh_prefix_len);
556     mdb_printf("    lh_prefix = %llx\n", zlp->l_hdr.lh_prefix);
557     mdb_printf("    lh_nentries = %u\n", zlp->l_hdr.lh_nentries);
558     mdb_printf("    lh_nfree = %u\n", zlp->l_hdr.lh_nfree);
559     mdb_printf("    zlp->l_hdr.lh_nfree * 100 / (ZAP_LEAF_NUMCHUNKS(&l));");
560     mdb_printf("    lh_freelist = %u\n", zlp->l_hdr.lh_freelist);
561     mdb_printf("    lh_flags = %x (%s)\n", zlp->l_hdr.lh_flags,
562         zlp->l_hdr.lh_flags & ZLF_ENTRIES_CDSORTED ?
563         "ENTRIES_CDSORTED" : "");

565     if (verbose) {
566         mdb_printf(" hash table:\n");

```

```

567         for (i = 0; i < ZAP_LEAF_HASH_NUMENTRIES(&l); i++) {
568             if (zlp->l_hash[i] != CHAIN_END)
569                 mdb_printf("    %u: %u\n", i, zlp->l_hash[i]);
570         }
571     }

573     mdb_printf(" chunks:\n");
574     for (i = 0; i < ZAP_LEAF_NUMCHUNKS(&l); i++) {
575         /* LINTED: alignment */
576         zap_leaf_chunk_t *zlc = &ZAP_LEAF_CHUNK(&l, i);
577         switch (zlc->l_entry.le_type) {
578             case ZAP_CHUNK_FREE:
579                 if (verbose) {
580                     mdb_printf("    %u: free; lf_next = %u\n",
581                         i, zlc->l_free.lf_next);
582                 }
583                 break;
584             case ZAP_CHUNK_ENTRY:
585                 mdb_printf("    %u: entry\n", i);
586                 if (verbose) {
587                     mdb_printf("        le_next = %u\n",
588                         zlc->l_entry.le_next);
589                 }
590                 mdb_printf("        le_name_chunk = %u\n",
591                     zlc->l_entry.le_name_chunk);
592                 mdb_printf("        le_name_numints = %u\n",
593                     zlc->l_entry.le_name_numints);
594                 mdb_printf("        le_value_chunk = %u\n",
595                     zlc->l_entry.le_value_chunk);
596                 mdb_printf("        le_value_intlen = %u\n",
597                     zlc->l_entry.le_value_intlen);
598                 mdb_printf("        le_value_numints = %u\n",
599                     zlc->l_entry.le_value_numints);
600                 mdb_printf("        le_cd = %u\n",
601                     zlc->l_entry.le_cd);
602                 mdb_printf("        le_hash = %llx\n",
603                     zlc->l_entry.le_hash);
604                 break;
605             case ZAP_CHUNK_ARRAY:
606                 mdb_printf("    %u: array", i);
607                 if (strisprint((char *)zlc->l_array.la_array))
608                     mdb_printf(" \"%s\"", zlc->l_array.la_array);
609                 mdb_printf("\n");
610                 if (verbose) {
611                     int j;
612                     mdb_printf("        ");
613                     for (j = 0; j < ZAP_LEAF_ARRAY_BYTES; j++) {
614                         mdb_printf("%02x ",
615                             zlc->l_array.la_array[j]);
616                     }
617                     mdb_printf("\n");
618                 }
619                 if (zlc->l_array.la_next != CHAIN_END) {
620                     mdb_printf("        lf_next = %u\n",
621                         zlc->l_array.la_next);
622                 }
623                 break;
624             default:
625                 mdb_printf("    %u: undefined type %u\n",
626                     zlc->l_entry.le_type);
627         }
628     }

630     return (DCMD_OK);
631 }
unchanged_portion_omitted

```

```
1428 struct mdb_dsl_dir_phys;
1429 typedef struct mdb_dsl_dir_dbuf {
1430     uint8_t dddb_pad[offsetof(dmu_buf_t, db_data)];
1431     uintptr_t dddb_data;
1432 } mdb_dsl_dir_dbuf_t;

1434 #endif /* ! codereview */
1435 typedef struct mdb_dsl_dir {
1436     union {
1437         dmu_buf_t *dd_dmu_db;
1438         mdb_dsl_dir_dbuf_t *dd_db;
1439     } dd_db_u;
1428     uintptr_t dd_phys;
1440     int64_t dd_space_towrite[TXG_SIZE];
1441 } mdb_dsl_dir_t;
_____
    unchanged_portion_omitted_
```

new/usr/src/lib/libzfs/common/libzfs_impl.h

1

```
*****
6419 Thu Apr 25 16:14:56 2013
new/usr/src/lib/libzfs/common/libzfs_impl.h
3752 want more verifiable dbuf user eviction
Submitted by: Justin Gibbs <justing@spectrallogic.com>
Submitted by: Will Andrews <willa@spectrallogic.com>
*****
1 /*
2  * CDDL HEADER SART
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) 2005, 2010, Oracle and/or its affiliates. All rights reserved.
24  * Copyright (c) 2012 by Delphix. All rights reserved.
25  */

27 #ifndef _LIBZFS_IMPL_H
28 #define _LIBZFS_IMPL_H

30 #include <sys/dmu.h>
30 #include <sys/fs/zfs.h>
31 #include <sys/zfs_ioctl.h>
32 #include <sys/spa.h>
33 #include <sys/nvpair.h>

35 #include <libuutil.h>
36 #include <libzfs.h>
37 #include <libshare.h>
38 #include <libzfs_core.h>

40 #include <fm/libtopo.h>

42 #ifdef __cplusplus
43 extern "C" {
44 #endif

46 #ifdef VERIFY
47 #undef VERIFY
48 #endif
49 #define VERIFY verify

51 typedef struct libzfs_fru {
52     char *zf_device;
53     char *zf_fru;
54     struct libzfs_fru *zf_chain;
55     struct libzfs_fru *zf_next;
56 } libzfs_fru_t;
unchanged portion omitted
```

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

1

```
*****
75656 Thu Apr 25 16:14:56 2013
new/usr/src/uts/common/fs/zfs/dbuf.c
3752 want more verifiable dbuf user eviction
Submitted by: Justin Gibbs <justing@spectralogic.com>
Submitted by: Will Andrews <willa@spectralogic.com>
*****
unchanged_portion_omitted

204 static arc_evict_func_t dbuf_do_evict;

206 static void
207 dbuf_verify_user(dmu_buf_impl_t *db, boolean_t evicting)
208 {
209 #ifdef ZFS_DEBUG

211     if (db->db_level != 0)
212         ASSERT(db->db_user == NULL);

214     if (db->db_user == NULL)
215         return;

217     /* Clients must resolve a dbuf before attaching user data. */
218     ASSERT(db->db.db_data != NULL && db->db_state == DB_CACHED);
219     /*
220      * We can't check the hold count here, because they are modified
221      * independently of the dbuf mutex. But it would be nice to ensure
222      * that the user has the appropriate number.
223      */
224 #endif
225 }

227 /*
228  * Evict the dbuf's user, either immediately, or use a provided queue.
229  *
230  * Call dmu_buf_process_user_evicts or dmu_buf_destroy_user_evict_list
231  * on the list when finished generating it.
232  *
233  * NOTE: If db->db_immediate_evict is FALSE, evict_list_p must be provided.
234  * NOTE: See dmu_buf_user_t about how this process works.
235  */
236 static void
237 dbuf_evict_user(dmu_buf_impl_t *db, list_t *evict_list_p)
238 {
239     ASSERT(MUTEX_HELD(&db->db_mtx));
240     ASSERT(evict_list_p != NULL);
241     dbuf_verify_user(db, /*evicting*/B_TRUE);
242 #endif /* !codereview */

244     if (db->db_user == NULL)
210     if (db->db_level != 0 || db->db_evict_func == NULL)
245         return;

247     ASSERT(!list_link_active(&db->db_user->evict_queue_link));
248     list_insert_head(evict_list_p, db->db_user);
249     db->db_user = NULL;
250 }

252 /*
253  * Replace the current user of the dbuf. Requires that the caller knows who
254  * the old user is. Returns the old user, which may not necessarily be
255  * the same old_user provided by the caller.
256  */
257 dmu_buf_user_t *
258 dmu_buf_replace_user(dmu_buf_t *db_fake, dmu_buf_user_t *old_user,
```

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

2

```
259     dmu_buf_user_t *new_user)
260 {
261     dmu_buf_impl_t *db = (dmu_buf_impl_t *)db_fake;

263     mutex_enter(&db->db_mtx);
264     dbuf_verify_user(db, /*evicting*/B_FALSE);
265     if (db->db_user == old_user)
266         db->db_user = new_user;
267     else
268         old_user = db->db_user;
269     dbuf_verify_user(db, /*evicting*/B_FALSE);
270     mutex_exit(&db->db_mtx);

272     return (old_user);
273 }

275 /*
276  * Set the user eviction data for the DMU returns NULL on success,
277  * or the existing user if another user currently owns the buffer.
278  */
279 dmu_buf_user_t *
280 dmu_buf_set_user(dmu_buf_t *db_fake, dmu_buf_user_t *user)
281 {
282     return (dmu_buf_replace_user(db_fake, NULL, user));
283 }

285 dmu_buf_user_t *
286 dmu_buf_set_user_ie(dmu_buf_t *db_fake, dmu_buf_user_t *user)
287 {
288     dmu_buf_impl_t *db = (dmu_buf_impl_t *)db_fake;

290     db->db_immediate_evict = TRUE;
291     return (dmu_buf_set_user(db_fake, user));
292 }

294 /*
295  * Remove the user eviction data for the DMU buffer.
296  */
297 dmu_buf_user_t *
298 dmu_buf_remove_user(dmu_buf_t *db_fake, dmu_buf_user_t *user)
299 {
300     return (dmu_buf_replace_user(db_fake, user, NULL));
301 }

303 /*
304  * Returns the db_user set with dmu_buf_update_user(), or NULL if not set.
305  */
306 dmu_buf_user_t *
307 dmu_buf_get_user(dmu_buf_t *db_fake)
308 {
309     dmu_buf_impl_t *db = (dmu_buf_impl_t *)db_fake;

311     dbuf_verify_user(db, /*evicting*/B_FALSE);
312     return (db->db_user);
313 }

315 static void
316 dbuf_clear_data(dmu_buf_impl_t *db, list_t *evict_list_p)
317 {
318     ASSERT(MUTEX_HELD(&db->db_mtx));
319     ASSERT(db->db_buf == NULL || !arc_has_callback(db->db_buf));
320     dbuf_evict_user(db, evict_list_p);
321     db->db_buf = NULL;
322     db->db.db_data = NULL;
323     if (db->db_state != DB_NOFILL)
324         db->db_state = DB_UNCACHED;
```

```

325 }

327 static void
328 dbuf_set_data(dmu_buf_impl_t *db, arc_buf_t *buf)
329 {
330     ASSERT(MUTEX_HELD(&db->db_mtx));
331     ASSERT(db->db_buf == NULL || !arc_has_callback(db->db_buf));
332     ASSERT(buf != NULL);

334     db->db_buf = buf;
335     ASSERT(buf->b_data != NULL);
336     db->db.db_data = buf->b_data;
337     if (!arc_released(buf))
338         arc_set_callback(buf, dbuf_do_evict, db);
213     if (db->db_user_data_ptr_ptr)
214         *db->db_user_data_ptr_ptr = db->db.db_data;
215     db->db_evict_func(&db->db, db->db_user_ptr);
216     db->db_user_ptr = NULL;
217     db->db_user_data_ptr_ptr = NULL;
218     db->db_evict_func = NULL;
339 }

    unchanged_portion_omitted

357 void
358 dbuf_evict(dmu_buf_impl_t *db, list_t *evict_list_p)
359 {
360     ASSERT(MUTEX_HELD(&db->db_mtx));
361     ASSERT(db->db_buf == NULL);
362     ASSERT(db->db_data_pending == NULL);

364     dbuf_clear(db, evict_list_p);
244     dbuf_clear(db);
365     dbuf_destroy(db);
366 }

    unchanged_portion_omitted

525 #endif

407 static void
408 dbuf_update_data(dmu_buf_impl_t *db)
409 {
410     ASSERT(MUTEX_HELD(&db->db_mtx));
411     if (db->db_level == 0 && db->db_user_data_ptr_ptr) {
412         ASSERT(!refcount_is_zero(&db->db_holds));
413         *db->db_user_data_ptr_ptr = db->db.db_data;
414     }
415 }

417 static void
418 dbuf_set_data(dmu_buf_impl_t *db, arc_buf_t *buf)
419 {
420     ASSERT(MUTEX_HELD(&db->db_mtx));
421     ASSERT(db->db_buf == NULL || !arc_has_callback(db->db_buf));
422     db->db_buf = buf;
423     if (buf != NULL) {
424         ASSERT(buf->b_data != NULL);
425         db->db.db_data = buf->b_data;
426         if (!arc_released(buf))
427             arc_set_callback(buf, dbuf_do_evict, db);
428         dbuf_update_data(db);
429     } else {
430         dbuf_evict_user(db);
431         db->db.db_data = NULL;
432         if (db->db_state != DB_NOFILL)
433             db->db_state = DB_UNCACHED;
434     }

```

```

435 }

527 /*
528  * Loan out an arc_buf for read. Return the loaned arc_buf.
529  */
530 arc_buf_t *
531 dbuf_loan_arcbuf(dmu_buf_impl_t *db)
532 {
533     arc_buf_t *abuf;
534     list_t evict_list;

536     dmu_buf_create_user_evict_list(&evict_list);
537 #endif /* ! codereview */

539     mutex_enter(&db->db_mtx);
540     if (arc_released(db->db_buf) || refcount_count(&db->db_holds) > 1) {
541         int blksz = db->db.db_size;
542         spa_t *spa;

544         mutex_exit(&db->db_mtx);
545         DB_GET_SPA(&spa, db);
546         abuf = arc_loan_buf(spa, blksz);
547         bcopy(db->db.db_data, abuf->b_data, blksz);
548     } else {
549         abuf = db->db_buf;
550         arc_loan_inuse_buf(abuf, db);
551         dbuf_clear_data(db, &evict_list);
444         dbuf_set_data(db, NULL);
552         mutex_exit(&db->db_mtx);
553     }
554     dmu_buf_destroy_user_evict_list(&evict_list);
555 #endif /* ! codereview */
556     return (abuf);
557 }

559 uint64_t
560 dbuf_whichblock(dnode_t *dn, uint64_t offset)
561 {
562     if (dn->dn_datablkshift) {
563         return (offset >> dn->dn_datablkshift);
564     } else {
565         ASSERT3U(offset, <, dn->dn_datablksize);
566         return (0);
567     }
568 }

570 static void
571 dbuf_read_done(zio_t *zio, arc_buf_t *buf, void *vdb)
572 {
573     dmu_buf_impl_t *db = vdb;

575     mutex_enter(&db->db_mtx);
576     ASSERT3U(db->db_state, ==, DB_READ);
577     /*
578      * All reads are synchronous, so we must have a hold on the dbuf
579      */
580     ASSERT(refcount_count(&db->db_holds) > 0);
581     ASSERT(db->db_buf == NULL);
582     ASSERT(db->db.db_data == NULL);
583     if (db->db_level == 0 && db->db_freed_in_flight) {
584         /* we were freed in flight; disregard any error */
585         arc_release(buf, db);
586         bzero(buf->b_data, db->db.db_size);
587         arc_buf_freeze(buf);
588         db->db_freed_in_flight = FALSE;
589         dbuf_set_data(db, buf);

```

```

590         db->db_state = DB_CACHED;
591     } else if (zio == NULL || zio->io_error == 0) {
592         dbuf_set_data(db, buf);
593         db->db_state = DB_CACHED;
594     } else {
595         ASSERT(db->db_blkid != DMU_BONUS_BLKID);
596         ASSERT3P(db->db_buf, ==, NULL);
597         VERIFY(arc_buf_remove_ref(buf, db));
598         db->db_state = DB_UNCACHED;
599     }
600     cv_broadcast(&db->db_changed);
601     dbuf_rele_and_unlock(db, NULL);
602 }

604 static void
605 dbuf_read_impl(dmu_buf_impl_t *db, zio_t *zio, uint32_t *flags)
606 {
607     dnode_t *dn;
608     spa_t *spa;
609     zbookmark_t zb;
610     uint32_t aflags = ARC_NOWAIT;

612     DB_DNODE_ENTER(db);
613     dn = DB_DNODE(db);
614     ASSERT(!refcount_is_zero(&db->db_holds));
615     /* We need the struct_rwlock to prevent db_blkptr from changing. */
616     ASSERT(RW_LOCK_HELD(&dn->dn_struct_rwlock));
617     ASSERT(MUTEX_HELD(&db->db_mtx));
618     ASSERT(db->db_state == DB_UNCACHED);
619     ASSERT(db->db_buf == NULL);

621     if (db->db_blkid == DMU_BONUS_BLKID) {
622         int bonuslen = MIN(dn->dn_bonuslen, dn->dn_phys->dn_bonuslen);

624         ASSERT3U(bonuslen, <=, db->db.db_size);
625         db->db.db_data = zio_buf_alloc(DN_MAX_BONUSLEN);
626         arc_space_consume(DN_MAX_BONUSLEN, ARC_SPACE_OTHER);
627         if (bonuslen < DN_MAX_BONUSLEN)
628             bzero(db->db.db_data, DN_MAX_BONUSLEN);
629         if (bonuslen)
630             bcopy(DN_BONUS(dn->dn_phys), db->db.db_data, bonuslen);
631         DB_DNODE_EXIT(db);
632         dbuf_update_data(db);
633         db->db_state = DB_CACHED;
634         mutex_exit(&db->db_mtx);
635         return;
636     }

637     /*
638     * Recheck BP_IS_HOLE() after dnode_block_freed() in case dnode_sync()
639     * processes the delete record and clears the bp while we are waiting
640     * for the dn_mtx (resulting in a "no" from block_freed).
641     */
642     if (db->db_blkptr == NULL || BP_IS_HOLE(db->db_blkptr) ||
643         (db->db_level == 0 && (dnode_block_freed(dn, db->db_blkid) ||
644         BP_IS_HOLE(db->db_blkptr)))) {
645         arc_buf_contents_t type = DBUF_GET_BUFC_TYPE(db);

647         dbuf_set_data(db, arc_buf_alloc(dn->dn_objset->os_spa,
648         db->db.db_size, db, type));
649         DB_DNODE_EXIT(db);
650         bzero(db->db.db_data, db->db.db_size);
651         db->db_state = DB_CACHED;
652         *flags |= DB_RF_CACHED;
653         mutex_exit(&db->db_mtx);
654         return;

```

```

655     }

657     spa = dn->dn_objset->os_spa;
658     DB_DNODE_EXIT(db);

660     db->db_state = DB_READ;
661     mutex_exit(&db->db_mtx);

663     if (DBUF_IS_L2CACHEABLE(db))
664         aflags |= ARC_L2CACHE;

666     SET_BOOKMARK(&zb, db->db_objset->os_dsl_dataset ?
667         db->db_objset->os_dsl_dataset->ds_object : DMU_META_OBJSET,
668         db->db.db_object, db->db_level, db->db_blkid);

670     dbuf_add_ref(db, NULL);

672     (void) arc_read(zio, spa, db->db_blkptr,
673         dbuf_read_done, db, ZIO_PRIORITY_SYNC_READ,
674         (*flags & DB_RF_CANFAIL) ? ZIO_FLAG_CANFAIL : ZIO_FLAG_MUSTSUCCEED,
675         &aflags, &zb);
676     if (aflags & ARC_CACHED)
677         *flags |= DB_RF_CACHED;
678 }

unchanged_portion_omitted

761 static void
762 dbuf_noread(dmu_buf_impl_t *db)
763 {
764     list_t evict_list;

766 #endif /* ! codereview */
767     ASSERT(!refcount_is_zero(&db->db_holds));
768     ASSERT(db->db_blkid != DMU_BONUS_BLKID);
769     dmu_buf_create_user_evict_list(&evict_list);

771 #endif /* ! codereview */
772     mutex_enter(&db->db_mtx);
773     while (db->db_state == DB_READ || db->db_state == DB_FILL)
774         cv_wait(&db->db_changed, &db->db_mtx);
775     if (db->db_state == DB_UNCACHED) {
776         arc_buf_contents_t type = DBUF_GET_BUFC_TYPE(db);
777         spa_t *spa;

779         ASSERT(db->db_buf == NULL);
780         ASSERT(db->db.db_data == NULL);
781         DB_GET_SPA(&spa, db);
782         dbuf_set_data(db, arc_buf_alloc(spa, db->db.db_size, db, type));
783         db->db_state = DB_FILL;
784     } else if (db->db_state == DB_NOFILL) {
785         dbuf_clear_data(db, &evict_list);
786         dbuf_set_data(db, NULL);
787     } else {
788         ASSERT3U(db->db_state, ==, DB_CACHED);
789     }
790     mutex_exit(&db->db_mtx);
791     dmu_buf_destroy_user_evict_list(&evict_list);
792 #endif /* ! codereview */
793 }

794 /*
795  * This is our just-in-time copy function. It makes a copy of
796  * buffers, that have been modified in a previous transaction
797  * group, before we modify them in the current active group.
798  *
799  * This function is used in two places: when we are dirtying a

```

```

800 * buffer for the first time in a txg, and when we are freeing
801 * a range in a dnode that includes this buffer.
802 *
803 * Note that when we are called from dbuf_free_range() we do
804 * not put a hold on the buffer, we just traverse the active
805 * dbuf list for the dnode.
806 */
807 static void
808 dbuf_fix_old_data(dmu_buf_impl_t *db, uint64_t txg, list_t *evict_list_p)
809 dbuf_fix_old_data(dmu_buf_impl_t *db, uint64_t txg)
810 {
811     dbuf_dirty_record_t *dr = db->db_last_dirty;
812
813     ASSERT(MUTEX_HELD(&db->db_mtx));
814     ASSERT(db->db.db_data != NULL);
815     ASSERT(db->db.level == 0);
816     ASSERT(db->db.db_object != DMU_META_DNODE_OBJECT);
817
818     if (dr == NULL ||
819         (dr->dt.dl.dr_data !=
820          ((db->db_blkid == DMU_BONUS_BLKID) ? db->db.db_data : db->db_buf)))
821         return;
822
823     /*
824     * If the last dirty record for this dbuf has not yet synced
825     * and its referencing the dbuf data, either:
826     *   reset the reference to point to a new copy,
827     *   or (if there a no active holders)
828     *   just null out the current db_data pointer.
829     */
830     ASSERT(dr->dr_txg >= txg - 2);
831     if (db->db_blkid == DMU_BONUS_BLKID) {
832         /* Note that the data bufs here are zio_bufts */
833         dr->dt.dl.dr_data = zio_buf_alloc(DN_MAX_BONUSLEN);
834         arc_space_consume(DN_MAX_BONUSLEN, ARC_SPACE_OTHER);
835         bcopy(db->db.db_data, dr->dt.dl.dr_data, DN_MAX_BONUSLEN);
836     } else if (refcount_count(&db->db_holds) > db->db_dirtycnt) {
837         int size = db->db.db_size;
838         arc_buf_contents_t type = DBUF_GET_BUFC_TYPE(db);
839         spa_t *spa;
840
841         DB_GET_SPA(&spa, db);
842         dr->dt.dl.dr_data = arc_buf_alloc(spa, size, db, type);
843         bcopy(db->db.db_data, dr->dt.dl.dr_data->b_data, size);
844     } else {
845         dbuf_clear_data(db, evict_list_p);
846         dbuf_set_data(db, NULL);
847     }
848 }
849
850 unchanged_portion_omitted
851
852 /*
853 * Evict (if its unreferenced) or clear (if its referenced) any level-0
854 * data blocks in the free range, so that any future readers will find
855 * empty blocks. Also, if we happen accross any level-1 dbufs in the
856 * range that have not already been marked dirty, mark them dirty so
857 * they stay in memory.
858 */
859 void
860 dbuf_free_range(dnode_t *dn, uint64_t start, uint64_t end, dmu_tx_t *tx)
861 {
862     dmu_buf_impl_t *db, *db_next;
863     uint64_t txg = tx->tx_txg;
864     int epbs = dn->dn_indblkshift - SPA_BLKPTRSHIFT;
865     uint64_t first_ll = start >> epbs;
866     uint64_t last_ll = end >> epbs;

```

```

901     list_t evict_list;
902
903     dmu_buf_create_user_evict_list(&evict_list);
904 #endif /* !codereview */
905
906     if (end > dn->dn_maxblkid && (end != DMU_SPILL_BLKID)) {
907         end = dn->dn_maxblkid;
908         last_ll = end >> epbs;
909     }
910     dprintf_dnode(dn, "start=%llu end=%llu\n", start, end);
911     mutex_enter(&dn->dn_dbufs_mtx);
912     for (db = list_head(&dn->dn_dbufs); db; db = db_next) {
913         db_next = list_next(&dn->dn_dbufs, db);
914         ASSERT(db->db_blkid != DMU_BONUS_BLKID);
915
916         if (db->db.level == 1 &&
917             db->db_blkid >= first_ll && db->db_blkid <= last_ll) {
918             mutex_enter(&db->db_mtx);
919             if (db->db_last_dirty &&
920                 db->db_last_dirty->dr_txg < txg) {
921                 dbuf_add_ref(db, FTAG);
922                 mutex_exit(&db->db_mtx);
923                 dbuf_will_dirty(db, tx);
924                 dbuf_rele(db, FTAG);
925             } else {
926                 mutex_exit(&db->db_mtx);
927             }
928         }
929
930         if (db->db.level != 0)
931             continue;
932         dprintf_dbuf(db, "found buf %s\n", "");
933         if (db->db_blkid < start || db->db_blkid > end)
934             continue;
935
936         /* found a level 0 buffer in the range */
937         mutex_enter(&db->db_mtx);
938         if (dbuf_undirty(db, tx)) {
939             /* mutex has been dropped and dbuf destroyed */
940             continue;
941         }
942
943         if (db->db.state == DB_UNCACHED ||
944             db->db.state == DB_NOFILL ||
945             db->db.state == DB_EVICTING) {
946             ASSERT(db->db.db_data == NULL);
947             mutex_exit(&db->db_mtx);
948             continue;
949         }
950         if (db->db.state == DB_READ || db->db.state == DB_FILL) {
951             /* will be handled in dbuf_read_done or dbuf_rele */
952             db->db.freed_in_flight = TRUE;
953             mutex_exit(&db->db_mtx);
954             continue;
955         }
956         if (refcount_count(&db->db_holds) == 0) {
957             ASSERT(db->db_buf);
958             dbuf_clear(db, &evict_list);
959             dbuf_clear(db);
960             continue;
961         }
962         /* The dbuf is referenced */
963
964         if (db->db_last_dirty != NULL) {
965             dbuf_dirty_record_t *dr = db->db_last_dirty;

```

```

966         if (dr->dr_txg == txg) {
967             /*
968              * This buffer is "in-use", re-adjust the file
969              * size to reflect that this buffer may
970              * contain new data when we sync.
971              */
972             if (db->db_blkid != DMU_SPILL_BLKID &&
973                 db->db_blkid > dn->dn_maxblkid)
974                 dn->dn_maxblkid = db->db_blkid;
975             dbuf_unoverride(dr);
976         } else {
977             /*
978              * This dbuf is not dirty in the open context.
979              * Either uncache it (if its not referenced in
980              * the open context) or reset its contents to
981              * empty.
982              */
983             dbuf_fix_old_data(db, txg, &evict_list);
984             dbuf_fix_old_data(db, txg);
985         }
986         /* clear the contents if its cached */
987         if (db->db_state == DB_CACHED) {
988             ASSERT(db->db_data != NULL);
989             arc_release(db->db_buf, db);
990             bzero(db->db.db_data, db->db.db_size);
991             arc_buf_freeze(db->db_buf);
992         }
993     }
994     mutex_exit(&db->db_mtx);
995     dmu_buf_process_user_evicts(&evict_list);
996 #endif /* ! codereview */
997 }
998 mutex_exit(&dn->dn_dbufs_mtx);
999 dmu_buf_destroy_user_evict_list(&evict_list);
1000 #endif /* ! codereview */
1001 }

1003 static int
1004 dbuf_block_freeable(dmu_buf_impl_t *db)
1005 {
1006     dsl_dataset_t *ds = db->db_objset->os_dsl_dataset;
1007     uint64_t birth_txg = 0;

1009     /*
1010      * We don't need any locking to protect db_blkptr:
1011      * If it's syncing, then db_last_dirty will be set
1012      * so we'll ignore db_blkptr.
1013      */
1014     ASSERT(MUTEX_HELD(&db->db_mtx));
1015     if (db->db_last_dirty)
1016         birth_txg = db->db_last_dirty->dr_txg;
1017     else if (db->db_blkptr)
1018         birth_txg = db->db_blkptr->blk_birth;

1020     /*
1021      * If we don't exist or are in a snapshot, we can't be freed.
1022      * Don't pass the bp to dsl_dataset_block_freeable() since we
1023      * are holding the db_mtx lock and might deadlock if we are
1024      * prefetching a dedup-ed block.
1025      */
1026     if (birth_txg)
1027         return (ds == NULL ||
1028             dsl_dataset_block_freeable(ds, NULL, birth_txg));
1029     else
1030         return (FALSE);

```

```

1031 }

1033 void
1034 dbuf_new_size(dmu_buf_impl_t *db, int size, dmu_tx_t *tx)
1035 {
1036     arc_buf_t *buf, *obuf;
1037     int osize = db->db.db_size;
1038     arc_buf_contents_t type = DBUF_GET_BUFC_TYPE(db);
1039     dnode_t *dn;

1041     ASSERT(db->db_blkid != DMU_BONUS_BLKID);

1043     DB_DNODE_ENTER(db);
1044     dn = DB_DNODE(db);

1046     /* XXX does *this* func really need the lock? */
1047     ASSERT(RW_WRITE_HELD(&dn->dn_struct_rwlock));

1049     /*
1050      * This call to dbuf_will_dirty() with the dn_struct_rwlock held
1051      * is OK, because there can be no other references to the db
1052      * when we are changing its size, so no concurrent DB_FILL can
1053      * be happening.
1054      */
1055     /*
1056      * XXX we should be doing a dbuf_read, checking the return
1057      * value and returning that up to our callers
1058      */
1059     dbuf_will_dirty(db, tx);

1061     /* create the data buffer for the new block */
1062     buf = arc_buf_alloc(dn->dn_objset->os_spa, size, db, type);

1064     /* copy old block data to the new block */
1065     obuf = db->db_buf;
1066     bcopy(obuf->b_data, buf->b_data, MIN(osize, size));
1067     /* zero the remainder */
1068     if (size > osize)
1069         bzero((uint8_t *)buf->b_data + osize, size - osize);

1071     mutex_enter(&db->db_mtx);
1072     dbuf_set_data(db, buf);
1073     VERIFY(arc_buf_remove_ref(obuf, db));
1074     db->db.db_size = size;

1076     if (db->db_level == 0) {
1077         ASSERT3U(db->db_last_dirty->dr_txg, ==, tx->tx_txg);
1078         db->db_last_dirty->dt.dl.dr_data = buf;
1079     }
1080     mutex_exit(&db->db_mtx);

1082     dnode_willuse_space(dn, size-osize, tx);
1083     DB_DNODE_EXIT(db);
1084 }

1086 void
1087 dbuf_release_bp(dmu_buf_impl_t *db)
1088 {
1089     objset_t *os;

1091     DB_GET_OBJSET(&os, db);
1092     ASSERT(dsl_pool_sync_context(dmu_objset_pool(os));
1093         ASSERT(arc_released(os->os_phys_buf) ||
1094             list_link_active(&os->os_dsl_dataset->ds_synced_link));
1095     ASSERT(db->db_parent == NULL || arc_released(db->db_parent->db_buf));

```

```

1097     (void) arc_release(db->db_buf, db);
1098 }

1100 dbuf_dirty_record_t *
1101 dbuf_dirty(dmu_buf_impl_t *db, dmu_tx_t *tx)
1102 {
1103     dnode_t *dn;
1104     objset_t *os;
1105     dbuf_dirty_record_t **drp, *dr;
1106     int drop_struct_lock = FALSE;
1107     boolean_t do_free_accounting = B_FALSE;
1108     int txgoff = tx->tx_txg & TXG_MASK;
1109     list_t evict_list;

1111     dmu_buf_create_user_evict_list(&evict_list);
1112 #endif /* ! codereview */

1114     ASSERT(tx->tx_txg != 0);
1115     ASSERT(!refcount_is_zero(&db->db_holds));
1116     DMU_TX_DIRTY_BUF(tx, db);

1118     DB_DNODE_ENTER(db);
1119     dn = DB_DNODE(db);
1120     /*
1121      * Shouldn't dirty a regular buffer in syncing context. Private
1122      * objects may be dirtied in syncing context, but only if they
1123      * were already pre-dirtied in open context.
1124      */
1125     ASSERT(!dmu_tx_is_syncing(tx) ||
1126           BP_IS_HOLE(dn->dn_objset->os_rootbp) ||
1127           DMU_OBJECT_IS_SPECIAL(dn->dn_object) ||
1128           dn->dn_objset->os_dsl_dataset == NULL);
1129     /*
1130      * We make this assert for private objects as well, but after we
1131      * check if we're already dirty. They are allowed to re-dirty
1132      * in syncing context.
1133      */
1134     ASSERT(dn->dn_object == DMU_META_DNODE_OBJECT ||
1135           dn->dn_dirtyctx == DN_UNDIRTIED || dn->dn_dirtyctx ==
1136           (dmu_tx_is_syncing(tx) ? DN_DIRTY_SYNC : DN_DIRTY_OPEN));

1138     mutex_enter(&db->db_mtx);
1139     /*
1140      * XXX make this true for indirects too? The problem is that
1141      * transactions created with dmu_tx_create_assigned() from
1142      * syncing context don't bother holding ahead.
1143      */
1144     ASSERT(db->db_level != 0 ||
1145           db->db_state == DB_CACHED || db->db_state == DB_FILL ||
1146           db->db_state == DB_NOFILL);

1148     mutex_enter(&dn->dn_mtx);
1149     /*
1150      * Don't set dirtyctx to SYNC if we're just modifying this as we
1151      * initialize the objset.
1152      */
1153     if (dn->dn_dirtyctx == DN_UNDIRTIED &&
1154         !BP_IS_HOLE(dn->dn_objset->os_rootbp)) {
1155         dn->dn_dirtyctx =
1156             (dmu_tx_is_syncing(tx) ? DN_DIRTY_SYNC : DN_DIRTY_OPEN);
1157         ASSERT(dn->dn_dirtyctx_firstset == NULL);
1158         dn->dn_dirtyctx_firstset = kmem_alloc(1, KM_SLEEP);
1159     }
1160     mutex_exit(&dn->dn_mtx);

1162     if (db->db_blkid == DMU_SPILL_BLKID)

```

```

1163         dn->dn_have_spill = B_TRUE;

1165     /*
1166      * If this buffer is already dirty, we're done.
1167      */
1168     drp = &db->db_last_dirty;
1169     ASSERT(*drp == NULL || (*drp)->dr_txg <= tx->tx_txg ||
1170           db->db.db_object == DMU_META_DNODE_OBJECT);
1171     while ((dr = *drp) != NULL && dr->dr_txg > tx->tx_txg)
1172         drp = &dr->dr_next;
1173     if (dr && dr->dr_txg == tx->tx_txg) {
1174         DB_DNODE_EXIT(db);

1176         if (db->db_level == 0 && db->db_blkid != DMU_BONUS_BLKID) {
1177             /*
1178              * If this buffer has already been written out,
1179              * we now need to reset its state.
1180              */
1181             dbuf_unoverride(dr);
1182             if (db->db.db_object != DMU_META_DNODE_OBJECT &&
1183                 db->db_state != DB_NOFILL)
1184                 arc_buf_thaw(db->db_buf);
1185         }
1186         mutex_exit(&db->db_mtx);
1187         dmu_buf_destroy_user_evict_list(&evict_list);
1188 #endif /* ! codereview */
1189         return (dr);
1190     }

1192     /*
1193      * Only valid if not already dirty.
1194      */
1195     ASSERT(dn->dn_object == 0 ||
1196           dn->dn_dirtyctx == DN_UNDIRTIED || dn->dn_dirtyctx ==
1197           (dmu_tx_is_syncing(tx) ? DN_DIRTY_SYNC : DN_DIRTY_OPEN));

1199     ASSERT3U(dn->dn_nlevels, >, db->db_level);
1200     ASSERT((dn->dn_phys->dn_nlevels == 0 && db->db_level == 0) ||
1201           dn->dn_phys->dn_nlevels > db->db_level ||
1202           dn->dn_next_nlevels[txgoff] > db->db_level ||
1203           dn->dn_next_nlevels[(tx->tx_txg-1) & TXG_MASK] > db->db_level ||
1204           dn->dn_next_nlevels[(tx->tx_txg-2) & TXG_MASK] > db->db_level);

1206     /*
1207      * We should only be dirtying in syncing context if it's the
1208      * mos or we're initializing the os or it's a special object.
1209      * However, we are allowed to dirty in syncing context provided
1210      * we already dirtied it in open context. Hence we must make
1211      * this assertion only if we're not already dirty.
1212      */
1213     os = dn->dn_objset;
1214     ASSERT(!dmu_tx_is_syncing(tx) || DMU_OBJECT_IS_SPECIAL(dn->dn_object) ||
1215           os->os_dsl_dataset == NULL || BP_IS_HOLE(os->os_rootbp));
1216     ASSERT(db->db.db_size != 0);

1218     dprintf_dbuf(db, "size=%llx\n", (u_longlong_t)db->db_size);

1220     if (db->db_blkid != DMU_BONUS_BLKID) {
1221         /*
1222          * Update the accounting.
1223          * Note: we delay "free accounting" until after we drop
1224          * the db_mtx. This keeps us from grabbing other locks
1225          * (and possibly deadlocking) in bp_get_dsize() while
1226          * also holding the db_mtx.
1227          */
1228         dnode_willuse_space(dn, db->db.db_size, tx);

```

```

1229         do_free_accounting = dbuf_block_freeable(db);
1230     }
1231
1232     /*
1233     * If this buffer is dirty in an old transaction group we need
1234     * to make a copy of it so that the changes we make in this
1235     * transaction group won't leak out when we sync the older txg.
1236     */
1237     dr = kmem_malloc(sizeof (dbuf_dirty_record_t), KM_SLEEP);
1238     if (db->db_level == 0) {
1239         void *data_old = db->db_buf;
1240
1241         if (db->db_state != DB_NOFILL) {
1242             if (db->db_blkid == DMU_BONUS_BLKID) {
1243                 dbuf_fix_old_data(db, tx->tx_txg, &evict_list);
1244                 dbuf_fix_old_data(db, tx->tx_txg);
1245                 data_old = db->db.db_data;
1246             } else if (db->db.db_object != DMU_META_DNODE_OBJECT) {
1247                 /*
1248                 * Release the data buffer from the cache so
1249                 * that we can modify it without impacting
1250                 * possible other users of this cached data
1251                 * block. Note that indirect blocks and
1252                 * private objects are not released until the
1253                 * syncing state (since they are only modified
1254                 * then).
1255                 */
1256                 arc_release(db->db_buf, db);
1257                 dbuf_fix_old_data(db, tx->tx_txg, &evict_list);
1258                 dbuf_fix_old_data(db, tx->tx_txg);
1259                 data_old = db->db_buf;
1260             }
1261             ASSERT(data_old != NULL);
1262             dr->dt.dl.dr_data = data_old;
1263         } else {
1264             mutex_init(&dr->dt.di.dr_mtx, NULL, MUTEX_DEFAULT, NULL);
1265             list_create(&dr->dt.di.dr_children,
1266                 sizeof (dbuf_dirty_record_t),
1267                 offsetof(dbuf_dirty_record_t, dr_dirty_node));
1268         }
1269         dr->dr_dbuf = db;
1270         dr->dr_txg = tx->tx_txg;
1271         dr->dr_next = *drp;
1272         *drp = dr;
1273
1274     /*
1275     * We could have been freed_in_flight between the dbuf_noread
1276     * and dbuf_dirty. We win, as though the dbuf_noread() had
1277     * happened after the free.
1278     */
1279     if (db->db_level == 0 && db->db_blkid != DMU_BONUS_BLKID &&
1280         db->db_blkid != DMU_SPILL_BLKID) {
1281         mutex_enter(&dn->dn_mtx);
1282         dnode_clear_range(dn, db->db_blkid, 1, tx);
1283         mutex_exit(&dn->dn_mtx);
1284         db->db_freed_in_flight = FALSE;
1285     }
1286
1287     /*
1288     * This buffer is now part of this txg
1289     */
1290     dbuf_add_ref(db, (void *) (uintptr_t) tx->tx_txg);
1291     db->db_dirtycnt += 1;
1292     ASSERT3U(db->db_dirtycnt, <=, 3);

```

```

1293         mutex_exit(&db->db_mtx);
1294         dmu_buf_destroy_user_evict_list(&evict_list);
1295     #endif /* ! codereview */
1296
1297     if (db->db_blkid == DMU_BONUS_BLKID ||
1298         db->db_blkid == DMU_SPILL_BLKID) {
1299         mutex_enter(&dn->dn_mtx);
1300         ASSERT(!list_link_active(&dr->dr_dirty_node));
1301         list_insert_tail(&dn->dn_dirty_records[txgoff], dr);
1302         mutex_exit(&dn->dn_mtx);
1303         dnode_setdirty(dn, tx);
1304         DB_DNODE_EXIT(db);
1305         return (dr);
1306     } else if (do_free_accounting) {
1307         blkptr_t *bp = db->db_blkptr;
1308         int64_t willfree = (bp && !BP_IS_HOLE(bp)) ?
1309             bp_get_dsize(os->os_spa, bp) : db->db.db_size;
1310         /*
1311         * This is only a guess -- if the dbuf is dirty
1312         * in a previous txg, we don't know how much
1313         * space it will use on disk yet. We should
1314         * really have the struct_rwlock to access
1315         * db_blkptr, but since this is just a guess,
1316         * it's OK if we get an odd answer.
1317         */
1318         ddt_prefetch(os->os_spa, bp);
1319         dnode_willuse_space(dn, -willfree, tx);
1320     }
1321
1322     if (!RW_WRITE_HELD(&dn->dn_struct_rwlock)) {
1323         rw_enter(&dn->dn_struct_rwlock, RW_READER);
1324         drop_struct_lock = TRUE;
1325     }
1326
1327     if (db->db_level == 0) {
1328         dnode_new_blkid(dn, db->db_blkid, tx, drop_struct_lock);
1329         ASSERT(dn->dn_maxblkid >= db->db_blkid);
1330     }
1331
1332     if (db->db_level+1 < dn->dn_nlevels) {
1333         dmu_buf_impl_t *parent = db->db_parent;
1334         dbuf_dirty_record_t *di;
1335         int parent_held = FALSE;
1336
1337         if (db->db_parent == NULL || db->db_parent == dn->dn_dbuf) {
1338             int epbs = dn->dn_indblkshift - SPA_BLKPTRSHIFT;
1339
1340             parent = dbuf_hold_level(dn, db->db_level+1,
1341                 db->db_blkid >> epbs, FTAG);
1342             ASSERT(parent != NULL);
1343             parent_held = TRUE;
1344         }
1345         if (drop_struct_lock)
1346             rw_exit(&dn->dn_struct_rwlock);
1347         ASSERT3U(db->db_level+1, ==, parent->db_level);
1348         di = dbuf_dirty(parent, tx);
1349         if (parent_held)
1350             dbuf_rele(parent, FTAG);
1351
1352         mutex_enter(&db->db_mtx);
1353         /* possible race with dbuf_undirty() */
1354         if (db->db_last_dirty == dr ||
1355             dn->dn_object == DMU_META_DNODE_OBJECT) {
1356             mutex_enter(&di->dt.di.dr_mtx);
1357             ASSERT3U(di->dr_txg, ==, tx->tx_txg);
1358             ASSERT(!list_link_active(&dr->dr_dirty_node));

```

```

1359     list_insert_tail(&di->dt.di.dr_children, dr);
1360     mutex_exit(&di->dt.di.dr_mtx);
1361     dr->dr_parent = di;
1362 }
1363     mutex_exit(&db->db_mtx);
1364 } else {
1365     ASSERT(db->db_level+1 == dn->dn_nlevels);
1366     ASSERT(db->db_blkid < dn->dn_nblkptr);
1367     ASSERT(db->db_parent == NULL || db->db_parent == dn->dn_dbuf);
1368     mutex_enter(&dn->dn_mtx);
1369     ASSERT(!list_link_active(&dr->dr_dirty_node));
1370     list_insert_tail(&dn->dn_dirty_records[txgoff], dr);
1371     mutex_exit(&dn->dn_mtx);
1372     if (drop_struct_lock)
1373         rw_exit(&dn->dn_struct_rwlock);
1374 }
1376     dnode_setdirty(dn, tx);
1377     DB_DNODE_EXIT(db);
1378     return (dr);
1379 }
1381 /*
1382  * Return TRUE if this evicted the dbuf.
1383  */
1384 static boolean_t
1385 dbuf_undirty(dmu_buf_impl_t *db, dmu_tx_t *tx)
1386 {
1387     dnode_t *dn;
1388     uint64_t txg = tx->tx_txg;
1389     dbuf_dirty_record_t *dr, **drp;
1390     list_t evict_list;
1391 #endif /* !codereview */
1393     ASSERT(txg != 0);
1394     ASSERT(db->db_blkid != DMU_BONUS_BLKID);
1395     ASSERT0(db->db_level);
1396     ASSERT(MUTEX_HELD(&db->db_mtx));
1398     /*
1399      * If this buffer is not dirty, we're done.
1400      */
1401     for (drp = &db->db_last_dirty; (dr = *drp) != NULL; drp = &dr->dr_next)
1402         if (dr->dr_txg <= txg)
1403             break;
1404     if (dr == NULL || dr->dr_txg < txg)
1405         return (B_FALSE);
1406     ASSERT(dr->dr_txg == txg);
1407     ASSERT(dr->dr_dbuf == db);
1409     dmu_buf_create_user_evict_list(&evict_list);
1411 #endif /* !codereview */
1412     DB_DNODE_ENTER(db);
1413     dn = DB_DNODE(db);
1415     /*
1416      * Note: This code will probably work even if there are concurrent
1417      * holders, but it is untested in that scenario, as the ZPL and
1418      * ztest have additional locking (the range locks) that prevents
1419      * that type of concurrent access.
1420      */
1421     ASSERT3U(refcount_count(&db->db_holds), ==, db->db_dirtycnt);
1423     dprintf_dbuf(db, "size=%llx\n", (u_longlong_t)db->db_size);

```

```

1425     ASSERT(db->db_size != 0);
1427     /* XXX would be nice to fix up dn_towrite_space[] */
1429     *drp = dr->dr_next;
1431     /*
1432      * Note that there are three places in dbuf_dirty()
1433      * where this dirty record may be put on a list.
1434      * Make sure to do a list_remove corresponding to
1435      * every one of those list_insert calls.
1436      */
1437     if (dr->dr_parent) {
1438         mutex_enter(&dr->dr_parent->dt.di.dr_mtx);
1439         list_remove(&dr->dr_parent->dt.di.dr_children, dr);
1440         mutex_exit(&dr->dr_parent->dt.di.dr_mtx);
1441     } else if (db->db_blkid == DMU_SPILL_BLKID ||
1442         db->db_level+1 == dn->dn_nlevels) {
1443         ASSERT(db->db_blkptr == NULL || db->db_parent == dn->dn_dbuf);
1444         mutex_enter(&dn->dn_mtx);
1445         list_remove(&dn->dn_dirty_records[txg & TXG_MASK], dr);
1446         mutex_exit(&dn->dn_mtx);
1447     }
1448     DB_DNODE_EXIT(db);
1450     if (db->db_state != DB_NOFILL) {
1451         dbuf_unoverride(dr);
1453         ASSERT(db->db_buf != NULL);
1454         ASSERT(dr->dt.dl.dr_data != NULL);
1455         if (dr->dt.dl.dr_data != db->db_buf)
1456             VERIFY(arc_buf_remove_ref(dr->dt.dl.dr_data, db));
1457     }
1458     kmem_free(dr, sizeof (dbuf_dirty_record_t));
1460     ASSERT(db->db_dirtycnt > 0);
1461     db->db_dirtycnt --;
1463     if (refcount_remove(&db->db_holds, (void *) (uintptr_t)txg) == 0) {
1464         arc_buf_t *buf = db->db_buf;
1466         ASSERT(db->db_state == DB_NOFILL || arc_released(buf));
1467         dbuf_clear_data(db, &evict_list);
1468         dbuf_set_data(db, NULL);
1469         VERIFY(arc_buf_remove_ref(buf, db));
1470         dbuf_evict(db, &evict_list);
1471         dmu_buf_destroy_user_evict_list(&evict_list);
1472         dbuf_evict(db);
1473         return (B_TRUE);
1474     }
1475     dmu_buf_destroy_user_evict_list(&evict_list);
1476 #endif /* !codereview */
1477     return (B_FALSE);
1479 #pragma weak dmu_buf_will_dirty = dbuf_will_dirty
1480 void
1481 dbuf_will_dirty(dmu_buf_impl_t *db, dmu_tx_t *tx)
1482 {
1483     int rf = DB_RF_MUST_SUCCEED | DB_RF_NOPREFETCH;
1485     ASSERT(tx->tx_txg != 0);
1486     ASSERT(!refcount_is_zero(&db->db_holds));
1488     DB_DNODE_ENTER(db);

```

```

1489     if (RW_WRITE_HELD(&DB_DNODE(db)->dn_struct_rwlock))
1490         rf |= DB_RF_HAVESTRUCT;
1491     DB_DNODE_EXIT(db);
1492     (void) dbuf_read(db, NULL, rf);
1493     (void) dbuf_dirty(db, tx);
1494 }

1496 void
1497 dmu_buf_will_not_fill(dmu_buf_t *db_fake, dmu_tx_t *tx)
1498 {
1499     dmu_buf_impl_t *db = (dmu_buf_impl_t *)db_fake;

1501     db->db_state = DB_NOFILL;

1503     dmu_buf_will_fill(db_fake, tx);
1504 }

1506 void
1507 dmu_buf_will_fill(dmu_buf_t *db_fake, dmu_tx_t *tx)
1508 {
1509     dmu_buf_impl_t *db = (dmu_buf_impl_t *)db_fake;

1511     ASSERT(db->db_blkid != DMU_BONUS_BLKID);
1512     ASSERT(tx->tx_txg != 0);
1513     ASSERT(db->db_level == 0);
1514     ASSERT(!refcount_is_zero(&db->db_holds));

1516     ASSERT(db->db.db_object != DMU_META_DNODE_OBJECT ||
1517         dmu_tx_private_ok(tx));

1519     dbuf_noread(db);
1520     (void) dbuf_dirty(db, tx);
1521 }

1523 #pragma weak dmu_buf_fill_done = dbuf_fill_done
1524 /* ARGSUSED */
1525 void
1526 dbuf_fill_done(dmu_buf_impl_t *db, dmu_tx_t *tx)
1527 {
1528     mutex_enter(&db->db_mtx);
1529     DBUF_VERIFY(db);

1531     if (db->db_state == DB_FILL) {
1532         if (db->db_level == 0 && db->db_freed_in_flight) {
1533             ASSERT(db->db_blkid != DMU_BONUS_BLKID);
1534             /* we were freed while filling */
1535             /* XXX dbuf_undirty? */
1536             bzero(db->db.db_data, db->db.db_size);
1537             db->db_freed_in_flight = FALSE;
1538         }
1539         db->db_state = DB_CACHED;
1540         cv_broadcast(&db->db_changed);
1541     }
1542     mutex_exit(&db->db_mtx);
1543 }

1545 /*
1546 * Directly assign a provided arc buf to a given dbuf if it's not referenced
1547 * by anybody except our caller. Otherwise copy arcbuf's contents to dbuf.
1548 */
1549 void
1550 dbuf_assign_arcbuf(dmu_buf_impl_t *db, arc_buf_t *buf, dmu_tx_t *tx)
1551 {
1552     ASSERT(!refcount_is_zero(&db->db_holds));
1553     ASSERT(db->db_blkid != DMU_BONUS_BLKID);
1554     ASSERT(db->db_level == 0);

```

```

1555     ASSERT(DBUF_GET_BUFC_TYPE(db) == ARC_BUFC_DATA);
1556     ASSERT(buf != NULL);
1557     ASSERT(arc_buf_size(buf) == db->db.db_size);
1558     ASSERT(tx->tx_txg != 0);

1560     arc_return_buf(buf, db);
1561     ASSERT(arc_released(buf));

1563     mutex_enter(&db->db_mtx);

1565     while (db->db_state == DB_READ || db->db_state == DB_FILL)
1566         cv_wait(&db->db_changed, &db->db_mtx);

1568     ASSERT(db->db_state == DB_CACHED || db->db_state == DB_UNCACHED);

1570     if (db->db_state == DB_CACHED &&
1571         refcount_count(&db->db_holds) - 1 > db->db_dirtycnt) {
1572         mutex_exit(&db->db_mtx);
1573         (void) dbuf_dirty(db, tx);
1574         bcopy(buf->b_data, db->db.db_data, db->db.db_size);
1575         VERIFY(arc_buf_remove_ref(buf, db));
1576         xuio_stat_wbuf_copied();
1577         return;
1578     }

1580     xuio_stat_wbuf_nocopy();
1581     if (db->db_state == DB_CACHED) {
1582         dbuf_dirty_record_t *dr = db->db_last_dirty;

1584         ASSERT(db->db_buf != NULL);
1585         if (dr != NULL && dr->dr_txg == tx->tx_txg) {
1586             ASSERT(dr->dt.dl.dr_data == db->db_buf);
1587             if (!arc_released(db->db_buf)) {
1588                 ASSERT(dr->dt.dl.dr_override_state ==
1589                     DR_OVERRIDDEN);
1590                 arc_release(db->db_buf, db);
1591             }
1592             dr->dt.dl.dr_data = buf;
1593             VERIFY(arc_buf_remove_ref(db->db_buf, db));
1594         } else if (dr == NULL || dr->dt.dl.dr_data != db->db_buf) {
1595             arc_release(db->db_buf, db);
1596             VERIFY(arc_buf_remove_ref(db->db_buf, db));
1597         }
1598         db->db_buf = NULL;
1599     }
1600     ASSERT(db->db_buf == NULL);
1601     dbuf_set_data(db, buf);
1602     db->db_state = DB_FILL;
1603     mutex_exit(&db->db_mtx);
1604     (void) dbuf_dirty(db, tx);
1605     dbuf_fill_done(db, tx);
1606 }

1608 /*
1609 * "Clear" the contents of this dbuf. This will mark the dbuf
1610 * EVICTING and clear *most* of its references. Unfortunately,
1611 * when we are not holding the dn_dbufs_mtx, we can't clear the
1612 * entry in the dn_dbufs list. We have to wait until dbuf_destroy()
1613 * in this case. For callers from the DMU we will usually see:
1614 *   dbuf_clear()->arc_buf_evict()->dbuf_do_evict()->dbuf_destroy()
1615 * For the arc callback, we will usually see:
1616 *   dbuf_do_evict()->dbuf_clear();dbuf_destroy()
1617 * Sometimes, though, we will get a mix of these two:
1618 *   DMU: dbuf_clear()->arc_buf_evict()
1619 *   ARC: dbuf_do_evict()->dbuf_destroy()
1620 */

```

```

1621 void
1622 dbuf_clear(dmu_buf_impl_t *db, list_t *evict_list_p)
1623 {
1624     dnode_t *dn;
1625     dmu_buf_impl_t *parent = db->db_parent;
1626     dmu_buf_impl_t *dnadb;
1627     int dbuf_gone = FALSE;
1628
1629     ASSERT(MUTEX_HELD(&db->db_mtx));
1630     ASSERT(refcount_is_zero(&db->db_holds));
1631
1632     dbuf_evict_user(db, evict_list_p);
1633     dbuf_evict_user(db);
1634
1635     if (db->db_state == DB_CACHED) {
1636         ASSERT(db->db_data != NULL);
1637         if (db->db_blkid == DMU_BONUS_BLKID) {
1638             zio_buf_free(db->db_data, DN_MAX_BONUSLEN);
1639             arc_space_return(DN_MAX_BONUSLEN, ARC_SPACE_OTHER);
1640         }
1641         db->db_data = NULL;
1642         db->db_state = DB_UNCACHED;
1643     }
1644
1645     ASSERT(db->db_state == DB_UNCACHED || db->db_state == DB_NOFILL);
1646     ASSERT(db->db_data_pending == NULL);
1647
1648     db->db_state = DB_EVICTING;
1649     db->db_blkptr = NULL;
1650
1651     DB_DNODE_ENTER(db);
1652     dn = DB_DNODE(db);
1653     dnadb = dn->dn_dbuf;
1654     if (db->db_blkid != DMU_BONUS_BLKID && MUTEX_HELD(&dn->dn_dbufs_mtx)) {
1655         list_remove(&dn->dn_dbufs, db);
1656         (void) atomic_dec_32_nv(&dn->dn_dbufs_count);
1657         membar_producer();
1658         DB_DNODE_EXIT(db);
1659         /*
1660          * Decrementing the dbuf count means that the hold corresponding
1661          * to the removed dbuf is no longer discounted in dnode_move(),
1662          * so the dnode cannot be moved until after we release the hold.
1663          * The membar_producer() ensures visibility of the decremented
1664          * value in dnode_move(), since DB_DNODE_EXIT doesn't actually
1665          * release any lock.
1666          */
1667         dnode_rele(dn, db);
1668         db->db_dnode_handle = NULL;
1669     } else {
1670         DB_DNODE_EXIT(db);
1671     }
1672
1673     if (db->db_buf)
1674         dbuf_gone = arc_buf_evict(db->db_buf);
1675
1676     if (!dbuf_gone)
1677         mutex_exit(&db->db_mtx);
1678
1679     /*
1680      * If this dbuf is referenced from an indirect dbuf,
1681      * decrement the ref count on the indirect dbuf.
1682      */
1683     if (parent && parent != dnadb)
1684         dbuf_rele(parent, db);

```

unchanged_portion_omitted

```

1753 static dmu_buf_impl_t *
1754 dbuf_create(dnode_t *dn, uint8_t level, uint64_t blkid,
1755             dmu_buf_impl_t *parent, blkptr_t *blkptr)
1756 {
1757     objset_t *os = dn->dn_objset;
1758     dmu_buf_impl_t *db, *odb;
1759
1760     ASSERT(RW_LOCK_HELD(&dn->dn_struct_rwlock));
1761     ASSERT(dn->dn_type != DMU_OT_NONE);
1762
1763     db = kmem_cache_alloc(dbuf_cache, KM_SLEEP);
1764
1765     db->db_objset = os;
1766     db->db_object = dn->dn_object;
1767     db->db_level = level;
1768     db->db_blkid = blkid;
1769     db->db_last_dirty = NULL;
1770     db->db_dirtycnt = 0;
1771     db->db_dnode_handle = dn->dn_handle;
1772     db->db_parent = parent;
1773     db->db_blkptr = blkptr;
1774
1775     db->db_user = NULL;
1776     db->db_user_ptr = NULL;
1777     db->db_user_data_ptr_ptr = NULL;
1778     db->db_evict_func = NULL;
1779     db->db_immediate_evict = 0;
1780     db->db_freed_in_flight = 0;
1781
1782     if (blkid == DMU_BONUS_BLKID) {
1783         ASSERT3P(parent, ==, dn->dn_dbuf);
1784         db->db_size = DN_MAX_BONUSLEN -
1785             (dn->dn_nblkptr-1) * sizeof(blkptr_t);
1786         ASSERT3U(db->db_size, >=, dn->dn_bonuslen);
1787         db->db_offset = DMU_BONUS_BLKID;
1788         db->db_state = DB_UNCACHED;
1789         /* the bonus dbuf is not placed in the hash table */
1790         arc_space_consume(sizeof(dmu_buf_impl_t), ARC_SPACE_OTHER);
1791         return(db);
1792     } else if (blkid == DMU_SPILL_BLKID) {
1793         db->db_size = (blkptr != NULL) ?
1794             BP_GET_LSIZE(blkptr) : SPA_MINBLOCKSIZE;
1795         db->db_offset = 0;
1796     } else {
1797         int blocksize =
1798             db->db_level ? 1<<dn->dn_indblkshift : dn->dn_datablksz;
1799         db->db_size = blocksize;
1800         db->db_offset = db->db_blkid * blocksize;
1801     }
1802
1803     /*
1804      * Hold the dn_dbufs_mtx while we get the new dbuf
1805      * in the hash table *and* added to the dbufs list.
1806      * This prevents a possible deadlock with someone
1807      * trying to look up this dbuf before its added to the
1808      * dn_dbufs list.
1809      */
1810     mutex_enter(&dn->dn_dbufs_mtx);
1811     db->db_state = DB_EVICTING;
1812     if ((odb = dbuf_hash_insert(db)) != NULL) {
1813         /* someone else inserted it first */
1814         kmem_cache_free(dbuf_cache, db);
1815         mutex_exit(&dn->dn_dbufs_mtx);
1816         return(odb);
1817     }

```

```

1815     list_insert_head(&dn->dn_dbufs, db);
1816     db->db_state = DB_UNCACHED;
1817     mutex_exit(&dn->dn_dbufs_mtx);
1818     arc_space_consume(sizeof (dmu_buf_impl_t), ARC_SPACE_OTHER);

1820     if (parent && parent != dn->dn_dbuf)
1821         dbuf_add_ref(parent, db);

1823     ASSERT(dn->dn_object == DMU_META_DNODE_OBJECT ||
1824            refcount_count(&dn->dn_holds) > 0);
1825     (void) refcount_add(&dn->dn_holds, db);
1826     (void) atomic_inc_32_nv(&dn->dn_dbufs_count);

1828     dprintf_dbuf(db, "db=%p\n", db);

1830     return (db);
1831 }

1833 static int
1834 dbuf_do_evict(void *private)
1835 {
1836     arc_buf_t *buf = private;
1837     dmu_buf_impl_t *db = buf->b_private;
1838     list_t evict_list;

1840     dmu_buf_create_user_evict_list(&evict_list);
1841 #endif /* ! codereview */

1843     if (!MUTEX_HELD(&db->db_mtx))
1844         mutex_enter(&db->db_mtx);

1846     ASSERT(refcount_is_zero(&db->db_holds));

1848     if (db->db_state != DB_EVICTING) {
1849         ASSERT(db->db_state == DB_CACHED);
1850         DBUF_VERIFY(db);
1851         db->db_buf = NULL;
1852         dbuf_evict(db, &evict_list);
1853     } else {
1854         mutex_exit(&db->db_mtx);
1855         dbuf_destroy(db);
1856     }
1857     dmu_buf_destroy_user_evict_list(&evict_list);
1858 #endif /* ! codereview */
1859     return (0);
1860 }

1862 static void
1863 dbuf_destroy(dmu_buf_impl_t *db)
1864 {
1865     ASSERT(refcount_is_zero(&db->db_holds));

1867     if (db->db_blkid != DMU_BONUS_BLKID) {
1868         /*
1869          * If this dbuf is still on the dn_dbufs list,
1870          * remove it from that list.
1871          */
1872         if (db->db_dnode_handle != NULL) {
1873             dnode_t *dn;

1875             DB_DNODE_ENTER(db);
1876             dn = DB_DNODE(db);
1877             mutex_enter(&dn->dn_dbufs_mtx);
1878             list_remove(&dn->dn_dbufs, db);
1879             (void) atomic_dec_32_nv(&dn->dn_dbufs_count);

```

```

1880         mutex_exit(&dn->dn_dbufs_mtx);
1881         DB_DNODE_EXIT(db);
1882         /*
1883          * Decrementing the dbuf count means that the hold
1884          * corresponding to the removed dbuf is no longer
1885          * discounted in dnode_move(), so the dnode cannot be
1886          * moved until after we release the hold.
1887          */
1888         dnode_rele(dn, db);
1889         db->db_dnode_handle = NULL;
1890     }
1891     dbuf_hash_remove(db);
1892 }
1893 db->db_parent = NULL;
1894 db->db_buf = NULL;

1896     ASSERT(!list_link_active(&db->db_link));
1897     ASSERT(db->db_db_data == NULL);
1898     ASSERT(db->db_hash_next == NULL);
1899     ASSERT(db->db_blkptr == NULL);
1900     ASSERT(db->db_data_pending == NULL);

1902     kmem_cache_free(dbuf_cache, db);
1903     arc_space_return(sizeof (dmu_buf_impl_t), ARC_SPACE_OTHER);
1904 }

1906 void
1907 dbuf_prefetch(dnode_t *dn, uint64_t blkid)
1908 {
1909     dmu_buf_impl_t *db = NULL;
1910     blkptr_t *bp = NULL;

1912     ASSERT(blkid != DMU_BONUS_BLKID);
1913     ASSERT(RW_LOCK_HELD(&dn->dn_struct_rwlock));

1915     if (dnode_block_freed(dn, blkid))
1916         return;

1918     /* dbuf_find() returns with db_mtx held */
1919     if (db = dbuf_find(dn, 0, blkid)) {
1920         /*
1921          * This dbuf is already in the cache. We assume that
1922          * it is already CACHED, or else about to be either
1923          * read or filled.
1924          */
1925         mutex_exit(&db->db_mtx);
1926         return;
1927     }

1929     if (dbuf_findbp(dn, 0, blkid, TRUE, &db, &bp) == 0) {
1930         if (bp && !BP_IS_HOLE(bp)) {
1931             int priority = dn->dn_type == DMU_OT_DDT_ZAP ?
1932                 ZIO_PRIORITY_DDT_PREFETCH : ZIO_PRIORITY_ASYNC_READ;
1933             dsl_dataset_t *ds = dn->dn_objset->os_dsl_dataset;
1934             uint32_t aflags = ARC_NOWAIT | ARC_PREFETCH;
1935             zbookmark_t zb;

1937             SET_BOOKMARK(&zb, ds ? ds->ds_object : DMU_META_OBJSET,
1938                          dn->dn_object, 0, blkid);

1940             (void) arc_read(NULL, dn->dn_objset->os_spa,
1941                             bp, NULL, NULL, priority,
1942                             ZIO_FLAG_CANFAIL | ZIO_FLAG_SPECULATIVE,
1943                             &aflags, &zb);
1944         }
1945         if (db)

```

```

1946         dbuf_rele(db, NULL);
1947     }
1948 }

1950 /*
1951  * Returns with db_holds incremented, and db_mtx not held.
1952  * Note: dn_struct_rwlock must be held.
1953  */
1954 int
1955 dbuf_hold_impl(dnode_t *dn, uint8_t level, uint64_t blkid, int fail_sparse,
1956               void *tag, dmu_buf_impl_t **dbp)
1957 {
1958     dmu_buf_impl_t *db, *parent = NULL;
1959     list_t evict_list;
1960 #endif /* ! codereview */

1962     ASSERT(blkid != DMU_BONUS_BLKID);
1963     ASSERT(RW_LOCK_HELD(&dn->dn_struct_rwlock));
1964     ASSERT3U(dn->dn_nlevels, >, level);

1966     dmu_buf_create_user_evict_list(&evict_list);

1968 #endif /* ! codereview */
1969     *dbp = NULL;
1970 top:
1971     /* dbuf_find() returns with db_mtx held */
1972     db = dbuf_find(dn, level, blkid);

1974     if (db == NULL) {
1975         blkptr_t *bp = NULL;
1976         int err;

1978         ASSERT3P(parent, ==, NULL);
1979         err = dbuf_findbp(dn, level, blkid, fail_sparse, &parent, &bp);
1980         if (fail_sparse) {
1981             if (err == 0 && bp && BP_IS_HOLE(bp))
1982                 err = SET_ERROR(ENOENT);
1983             if (err) {
1984                 if (parent)
1985                     dbuf_rele(parent, NULL);
1986                 return (err);
1987             }
1988         }
1989         if (err && err != ENOENT)
1990             return (err);
1991         db = dbuf_create(dn, level, blkid, parent, bp);
1992     }

1994     if (db->db_buf && refcount_is_zero(&db->db_holds)) {
1995         arc_buf_add_ref(db->db_buf, db);
1996         if (db->db_buf->b_data == NULL) {
1997             dbuf_clear(db, &evict_list);
1998             dbuf_clear(db);
1999             if (parent) {
2000                 dbuf_rele(parent, NULL);
2001                 parent = NULL;
2002             }
2003             goto top;
2004         }
2005         ASSERT3P(db->db_data, ==, db->db_buf->b_data);
2006     }

2007     ASSERT(db->db_buf == NULL || arc_referenced(db->db_buf));

2009     /*
2010     * If this buffer is currently syncing out, and we are are

```

```

2011     * still referencing it from db_data, we need to make a copy
2012     * of it in case we decide we want to dirty it again in this txg.
2013     */
2014     if (db->db_level == 0 && db->db_blkid != DMU_BONUS_BLKID &&
2015         dn->dn_object != DMU_META_DNODE_OBJECT &&
2016         db->db_state == DB_CACHED && db->db_data_pending) {
2017         dbuf_dirty_record_t *dr = db->db_data_pending;

2019         if (dr->dt.dl.dr_data == db->db_buf) {
2020             arc_buf_contents_t type = DBUF_GET_BUFC_TYPE(db);

2022             dbuf_set_data(db,
2023                          arc_buf_alloc(dn->dn_objset->os_spa,
2024                                         db->db.db_size, db, type));
2025             bcopy(dr->dt.dl.dr_data->b_data, db->db.db_data,
2026                  db->db.db_size);
2027         }
2028     }

2030     (void) refcount_add(&db->db_holds, tag);
2031     dbuf_update_data(db);
2032     DBUF_VERIFY(db);
2033     mutex_exit(&db->db_mtx);

2034     dmu_buf_destroy_user_evict_list(&evict_list);

2036 #endif /* ! codereview */
2037     /* NOTE: we can't rele the parent until after we drop the db_mtx */
2038     if (parent)
2039         dbuf_rele(parent, NULL);

2041     ASSERT3P(DB_DNODE(db), ==, dn);
2042     ASSERT3U(db->db_blkid, ==, blkid);
2043     ASSERT3U(db->db_level, ==, level);
2044     *dbp = db;

2046     return (0);
2047 }

2049 dmu_buf_impl_t *
2050 dbuf_hold(dnode_t *dn, uint64_t blkid, void *tag)
2051 {
2052     dmu_buf_impl_t *db;
2053     int err = dbuf_hold_impl(dn, 0, blkid, FALSE, tag, &db);
2054     return (err ? NULL : db);
2055 }

2057 dmu_buf_impl_t *
2058 dbuf_hold_level(dnode_t *dn, int level, uint64_t blkid, void *tag)
2059 {
2060     dmu_buf_impl_t *db;
2061     int err = dbuf_hold_impl(dn, level, blkid, FALSE, tag, &db);
2062     return (err ? NULL : db);
2063 }

2065 void
2066 dbuf_create_bonus(dnode_t *dn)
2067 {
2068     ASSERT(RW_WRITE_HELD(&dn->dn_struct_rwlock));

2070     ASSERT(dn->dn_bonus == NULL);
2071     dn->dn_bonus = dbuf_create(dn, 0, DMU_BONUS_BLKID, dn->dn_dbuf, NULL);
2072 }

2074 int
2075 dbuf_spill_set_blkisz(dmu_buf_t *db_fake, uint64_t blkisz, dmu_tx_t *tx)

```

```

2076 {
2077     dmu_buf_impl_t *db = (dmu_buf_impl_t *)db_fake;
2078     dnode_t *dn;

2080     if (db->db_blkid != DMU_SPILL_BLKID)
2081         return (SET_ERROR(ENOTSUP));
2082     if (blkksz == 0)
2083         blkksz = SPA_MINBLOCKSIZE;
2084     if (blkksz > SPA_MAXBLOCKSIZE)
2085         blkksz = SPA_MAXBLOCKSIZE;
2086     else
2087         blkksz = P2ROUNDUP(blkksz, SPA_MINBLOCKSIZE);

2089     DB_DNODE_ENTER(db);
2090     dn = DB_DNODE(db);
2091     rw_enter(&dn->dn_struct_rwlock, RW_WRITER);
2092     dbuf_new_size(db, blkksz, tx);
2093     rw_exit(&dn->dn_struct_rwlock);
2094     DB_DNODE_EXIT(db);

2096     return (0);
2097 }

2099 void
2100 dbuf_rm_spill(dnode_t *dn, dmu_tx_t *tx)
2101 {
2102     dbuf_free_range(dn, DMU_SPILL_BLKID, DMU_SPILL_BLKID, tx);
2103 }

2105 #pragma weak dmu_buf_add_ref = dbuf_add_ref
2106 void
2107 dbuf_add_ref(dmu_buf_impl_t *db, void *tag)
2108 {
2109     int64_t holds = refcount_add(&db->db_holds, tag);
2110     ASSERT(holds > 1);
2111 }

2113 /*
2114  * If you call dbuf_rele() you had better not be referencing the dnode handle
2115  * unless you have some other direct or indirect hold on the dnode. (An indirect
2116  * hold is a hold on one of the dnode's dbufs, including the bonus buffer.)
2117  * Without that, the dbuf_rele() could lead to a dnode_rele() followed by the
2118  * dnode's parent dbuf evicting its dnode handles.
2119  */
2120 #pragma weak dmu_buf_rele = dbuf_rele
2121 void
2122 dbuf_rele(dmu_buf_impl_t *db, void *tag)
2123 {
2124     mutex_enter(&db->db_mtx);
2125     dbuf_rele_and_unlock(db, tag);
2126 }

2128 /*
2129  * dbuf_rele() for an already-locked dbuf. This is necessary to allow
2130  * db_dirtycnt and db_holds to be updated atomically.
2131  */
2132 void
2133 dbuf_rele_and_unlock(dmu_buf_impl_t *db, void *tag)
2134 {
2135     int64_t holds;
2136     list_t evict_list;
2137 #endif /* ! codereview */

2139     ASSERT(MUTEX_HELD(&db->db_mtx));
2140     DBUF_VERIFY(db);

```

```

2142     dmu_buf_create_user_evict_list(&evict_list);

2144 #endif /* ! codereview */
2145 /*
2146  * Remove the reference to the dbuf before removing its hold on the
2147  * dnode so we can guarantee in dnode_move() that a referenced bonus
2148  * buffer has a corresponding dnode hold.
2149  */
2150 holds = refcount_remove(&db->db_holds, tag);
2151 ASSERT(holds >= 0);

2153 /*
2154  * We can't freeze indirects if there is a possibility that they
2155  * may be modified in the current syncing context.
2156  */
2157 if (db->db_buf && holds == (db->db_level == 0 ? db->db_dirtycnt : 0))
2158     arc_buf_freeze(db->db_buf);

2160 if (holds == db->db_dirtycnt &&
2161     db->db_level == 0 && db->db_immediate_evict)
2162     dbuf_evict_user(db, &evict_list);
2163     dbuf_evict_user(db);

2164 if (holds == 0) {
2165     if (db->db_blkid == DMU_BONUS_BLKID) {
2166         mutex_exit(&db->db_mtx);

2168         /*
2169          * If the dnode moves here, we cannot cross this barrier
2170          * until the move completes.
2171          */
2172         DB_DNODE_ENTER(db);
2173         (void) atomic_dec_32_nv(&DB_DNODE(db)->dn_dbufs_count);
2174         DB_DNODE_EXIT(db);
2175         /*
2176          * The bonus buffer's dnode hold is no longer discounted
2177          * in dnode_move(). The dnode cannot move until after
2178          * the dnode_rele().
2179          */
2180         dnode_rele(DB_DNODE(db), db);
2181     } else if (db->db_buf == NULL) {
2182         /*
2183          * This is a special case: we never associated this
2184          * dbuf with any data allocated from the ARC.
2185          */
2186         ASSERT(db->db_state == DB_UNCACHED ||
2187             db->db_state == DB_NOFILL);
2188         dbuf_evict(db, &evict_list);
2189         dbuf_evict(db);
2190     } else if (arc_released(db->db_buf)) {
2191         arc_buf_t *buf = db->db_buf;
2192         /*
2193          * This dbuf has anonymous data associated with it.
2194          */
2195         dbuf_clear_data(db, &evict_list);
2196         dbuf_set_data(db, NULL);
2197         VERIFY(arc_buf_remove_ref(buf, db));
2198         dbuf_evict(db, &evict_list);
2199         dbuf_evict(db);
2200     } else {
2201         VERIFY(!arc_buf_remove_ref(db->db_buf, db));

2202         /*
2203          * A dbuf will be eligible for eviction if either the
2204          * 'primarycache' property is set or a duplicate
2205          * copy of this buffer is already cached in the arc.

```

```

2204         *
2205         * In the case of the 'primarycache' a buffer
2206         * is considered for eviction if it matches the
2207         * criteria set in the property.
2208         *
2209         * To decide if our buffer is considered a
2210         * duplicate, we must call into the arc to determine
2211         * if multiple buffers are referencing the same
2212         * block on-disk. If so, then we simply evict
2213         * ourselves.
2214         */
2215         if (!DBUF_IS_CACHEABLE(db) ||
2216             arc_buf_eviction_needed(db->db_buf))
2217             dbuf_clear(db, &evict_list);
1088             dbuf_clear(db);
2218         else
2219             mutex_exit(&db->db_mtx);
2220     } else {
2221     }
2222         mutex_exit(&db->db_mtx);
2223     }
2224     dmu_buf_destroy_user_evict_list(&evict_list);
2225 #endif /* ! codereview */
2226 }

2228 #pragma weak dmu_buf_refcount = dbuf_refcount
2229 uint64_t
2230 dbuf_refcount(dmu_buf_impl_t *db)
2231 {
2232     return (refcount_count(&db->db_holds));
2233 }

1095 void *
1096 dmu_buf_set_user(dmu_buf_t *db_fake, void *user_ptr, void *user_data_ptr_ptr,
1097                 dmu_buf_evict_func_t *evict_func)
1098 {
1099     return (dmu_buf_update_user(db_fake, NULL, user_ptr,
1100                                user_data_ptr_ptr, evict_func));
1101 }

1103 void *
1104 dmu_buf_set_user_ie(dmu_buf_t *db_fake, void *user_ptr, void *user_data_ptr_ptr,
1105                    dmu_buf_evict_func_t *evict_func)
1106 {
1107     dmu_buf_impl_t *db = (dmu_buf_impl_t *)db_fake;

1109     db->db_immediate_evict = TRUE;
1110     return (dmu_buf_update_user(db_fake, NULL, user_ptr,
1111                                user_data_ptr_ptr, evict_func));
1112 }

1114 void *
1115 dmu_buf_update_user(dmu_buf_t *db_fake, void *old_user_ptr, void *user_ptr,
1116                    void *user_data_ptr_ptr, dmu_buf_evict_func_t *evict_func)
1117 {
1118     dmu_buf_impl_t *db = (dmu_buf_impl_t *)db_fake;
1119     ASSERT(db->db_level == 0);

1121     ASSERT((user_ptr == NULL) == (evict_func == NULL));

1123     mutex_enter(&db->db_mtx);

1125     if (db->db_user_ptr == old_user_ptr) {
1126         db->db_user_ptr = user_ptr;
1127         db->db_user_data_ptr_ptr = user_data_ptr_ptr;
1128         db->db_evict_func = evict_func;

```

```

1130         dbuf_update_data(db);
1131     } else {
1132         old_user_ptr = db->db_user_ptr;
1133     }

1135     mutex_exit(&db->db_mtx);
1136     return (old_user_ptr);
1137 }

1139 void *
1140 dmu_buf_get_user(dmu_buf_t *db_fake)
1141 {
1142     dmu_buf_impl_t *db = (dmu_buf_impl_t *)db_fake;
1143     ASSERT(!refcount_is_zero(&db->db_holds));

1145     return (db->db_user_ptr);
1146 }

2235 boolean_t
2236 dmu_buf_freeable(dmu_buf_t *dbuf)
2237 {
2238     boolean_t res = B_FALSE;
2239     dmu_buf_impl_t *db = (dmu_buf_impl_t *)dbuf;

2241     if (db->db_blkptr)
2242         res = dsl_dataset_block_freeable(db->db_objset->os_dsl_dataset,
2243                                         db->db_blkptr, db->db_blkptr->blk_birth);

2245     return (res);
2246 }
_____unchanged_portion_omitted_____

```

```

*****
56719 Thu Apr 25 16:14:57 2013
new/usr/src/uts/common/fs/zfs/dnode.c
3752 want more verifiable dbuf user eviction
Submitted by: Justin Gibbs <justing@spectrallogic.com>
Submitted by: Will Andrews <willa@spectrallogic.com>
*****
unchanged_portion_omitted_

433 /*
434  * Caller must be holding the dnode handle, which is released upon return.
435  */
436 static void
437 dnode_destroy(dnode_t *dn)
438 {
439     objset_t *os = dn->dn_objset;

441     ASSERT((dn->dn_id_flags & DN_ID_NEW_EXIST) == 0);

443     mutex_enter(&os->os_lock);
444     POINTER_INVALIDATE(&dn->dn_objset);
445     list_remove(&os->os_dnodes, dn);
446     mutex_exit(&os->os_lock);

448     /* the dnode can no longer move, so we can release the handle */
449     zrl_remove(&dn->dn_handle->dnh_zrlock);

451     dn->dn_allocated_txg = 0;
452     dn->dn_free_txg = 0;
453     dn->dn_assigned_txg = 0;

455     dn->dn_dirtyctx = 0;
456     if (dn->dn_dirtyctx_firstset != NULL) {
457         kmem_free(dn->dn_dirtyctx_firstset, 1);
458         dn->dn_dirtyctx_firstset = NULL;
459     }
460     if (dn->dn_bonus != NULL) {
461         list_t evict_list;

463         dmu_buf_create_user_evict_list(&evict_list);
464 #endif /* ! codereview */
465         mutex_enter(&dn->dn_bonus->db_mtx);
466         dbuf_evict(dn->dn_bonus, &evict_list);
467         dmu_buf_destroy_user_evict_list(&evict_list);
468         dbuf_evict(dn->dn_bonus);
469         dn->dn_bonus = NULL;
470     }
471     dn->dn_zio = NULL;

472     dn->dn_have_spill = B_FALSE;
473     dn->dn_oldused = 0;
474     dn->dn_oldflags = 0;
475     dn->dn_olduid = 0;
476     dn->dn_oldgid = 0;
477     dn->dn_newuid = 0;
478     dn->dn_newgid = 0;
479     dn->dn_id_flags = 0;

481     dmu_zfetch_rele(&dn->dn_zfetch);
482     kmem_cache_free(dnode_cache, dn);
483     arc_space_return(sizeof (dnode_t), ARC_SPACE_OTHER);
484 }
unchanged_portion_omitted_

964 static void
965 dnode_buf_pageout(dmu_buf_user_t *dbu)

```

```

959 dnode_buf_pageout(dmu_buf_t *db, void *arg)
960 {
961     dnode_children_t *children_dnodes = (dnode_children_t *)dbu;
962     dnode_children_t *children_dnodes = arg;
963     int i;
964     int epb = db->db_size >> DNODE_SHIFT;

965     ASSERT(epb == children_dnodes->dnc_count);

967     for (i = 0; i < children_dnodes->dnc_count; i++) {
968         for (i = 0; i < epb; i++) {
969             dnode_handle_t *dnh = &children_dnodes->dnc_children[i];
970             dnode_t *dn;

972             /*
973              * The dnode handle lock guards against the dnode moving to
974              * another valid address, so there is no need here to guard
975              * against changes to or from NULL.
976              */
977             if (dnh->dnh_dnode == NULL) {
978                 zrl_destroy(&dnh->dnh_zrlock);
979                 continue;
980             }

982             zrl_add(&dnh->dnh_zrlock);
983             dn = dnh->dnh_dnode;
984             /*
985              * If there are holds on this dnode, then there should
986              * be holds on the dnode's containing dbuf as well; thus
987              * it wouldn't be eligible for eviction and this function
988              * would not have been called.
989              */
990             ASSERT(refcount_is_zero(&dn->dn_holds));
991             ASSERT(refcount_is_zero(&dn->dn_tx_holds));

993             dnode_destroy(dn); /* implicit zrl_remove() */
994             zrl_destroy(&dnh->dnh_zrlock);
995             dnh->dnh_dnode = NULL;
996         }
997         kmem_free(children_dnodes, sizeof (dnode_children_t) +
1000             (children_dnodes->dnc_count - 1) * sizeof (dnode_handle_t));
1001     }

1003 /*
1004  * errors:
1005  * EINVAL - invalid object number.
1006  * EIO - i/o error.
1007  * succeeds even for free dnodes.
1008  */
1009 int
1010 dnode_hold_impl(objset_t *os, uint64_t object, int flag,
1011     void *tag, dnode_t **dnp)
1012 {
1013     int epb, idx, err;
1014     int drop_struct_lock = FALSE;
1015     int type;
1016     uint64_t blk;
1017     dnode_t *mdn, *dn;
1018     dmu_buf_impl_t *db;
1019     dnode_children_t *children_dnodes;
1020     dnode_handle_t *dnh;

1022     /*
1023      * If you are holding the spa config lock as writer, you shouldn't
1024      * be asking the DMU to do *anything* unless it's the root pool

```

```

1025     * which may require us to read from the root filesystem while
1026     * holding some (not all) of the locks as writer.
1027     */
1028 ASSERT(spa_config_held(os->os_spa, SCL_ALL, RW_WRITER) == 0 ||
1029     (spa_is_root(os->os_spa) &&
1030     spa_config_held(os->os_spa, SCL_STATE, RW_WRITER)));

1032 if (object == DMU_USERUSED_OBJECT || object == DMU_GROUPUSED_OBJECT) {
1033     dn = (object == DMU_USERUSED_OBJECT) ?
1034         DMU_USERUSED_DNODE(os) : DMU_GROUPUSED_DNODE(os);
1035     if (dn == NULL)
1036         return (SET_ERROR(ENOENT));
1037     type = dn->dn_type;
1038     if ((flag & DNODE_MUST_BE_ALLOCATED) && type == DMU_OT_NONE)
1039         return (SET_ERROR(ENOENT));
1040     if ((flag & DNODE_MUST_BE_FREE) && type != DMU_OT_NONE)
1041         return (SET_ERROR(EEXIST));
1042     DNODE_VERIFY(dn);
1043     (void) refcount_add(&dn->dn_holds, tag);
1044     *dnp = dn;
1045     return (0);
1046 }

1048 if (object == 0 || object >= DN_MAX_OBJECT)
1049     return (SET_ERROR(EINVAL));

1051 mdn = DMU_META_DNODE(os);
1052 ASSERT(mdn->dn_object == DMU_META_DNODE_OBJECT);

1054 DNODE_VERIFY(mdn);

1056 if (!RW_WRITE_HELD(&mdn->dn_struct_rwlock)) {
1057     rw_enter(&mdn->dn_struct_rwlock, RW_READER);
1058     drop_struct_lock = TRUE;
1059 }

1061 blk = dbuf_whichblock(mdn, object * sizeof (dnode_phys_t));

1063 db = dbuf_hold(mdn, blk, FTAG);
1064 if (drop_struct_lock)
1065     rw_exit(&mdn->dn_struct_rwlock);
1066 if (db == NULL)
1067     return (SET_ERROR(EIO));
1068 err = dbuf_read(db, NULL, DB_RF_CANFAIL);
1069 if (err) {
1070     dbuf_rele(db, FTAG);
1071     return (err);
1072 }

1074 ASSERT3U(db->db.db_size, >=, 1<<DNODE_SHIFT);
1075 epb = db->db.db_size >> DNODE_SHIFT;

1077 idx = object & (epb-1);

1079 ASSERT(DB_DNODE(db)->dn_type == DMU_OT_DNODE);
1080 children_dnodes = (dnode_children_t *)dmu_buf_get_user(&db->db);
1077 children_dnodes = dmu_buf_get_user(&db->db);
1081 if (children_dnodes == NULL) {
1082     int i;
1083     dnode_children_t *winner;
1084     children_dnodes = kmem_alloc(sizeof (dnode_children_t) +
1085     (epb - 1) * sizeof (dnode_handle_t), KM_SLEEP);
1086     children_dnodes->dnc_count = epb;
1087     dnh = &children_dnodes->dnc_children[0];
1088     for (i = 0; i < epb; i++) {
1089         zrl_init(&dnh[i].dnh_zrlock);

```

```

1090         dnh[i].dnh_dnode = NULL;
1091     }
1092     dmu_buf_init_user(&children_dnodes->db_evict,
1093     dnode_buf_pageout);
1094     winner = (dnode_children_t *)
1095     dmu_buf_set_user(&db->db, &children_dnodes->db_evict);
1096     if (winner) {
1097         if (winner = dmu_buf_set_user(&db->db, children_dnodes, NULL,
1098         dnode_buf_pageout)) {
1099             kmem_free(children_dnodes, sizeof (dnode_children_t) +
1100             (epb - 1) * sizeof (dnode_handle_t));
1101             children_dnodes = winner;
1102         }
1103     }
1104     ASSERT(children_dnodes->dnc_count == epb);

1104     dnh = &children_dnodes->dnc_children[idx];
1105     zrl_add(&dnh->dnh_zrlock);
1106     if ((dn = dnh->dnh_dnode) == NULL) {
1107         dnode_phys_t *phys = (dnode_phys_t *)db->db.data+idx;
1108         dnode_t *winner;

1110         dn = dnode_create(os, phys, db, object, dnh);
1111         winner = atomic_cas_ptr(&dnh->dnh_dnode, NULL, dn);
1112         if (winner != NULL) {
1113             zrl_add(&dnh->dnh_zrlock);
1114             dnode_destroy(dn); /* implicit zrl_remove() */
1115             dn = winner;
1116         }
1117     }

1119     mutex_enter(&dn->dn_mtx);
1120     type = dn->dn_type;
1121     if (dn->dn_free_txg ||
1122     ((flag & DNODE_MUST_BE_ALLOCATED) && type == DMU_OT_NONE) ||
1123     ((flag & DNODE_MUST_BE_FREE) &&
1124     (type != DMU_OT_NONE || !refcount_is_zero(&dn->dn_holds)))) {
1125         mutex_exit(&dn->dn_mtx);
1126         zrl_remove(&dnh->dnh_zrlock);
1127         dbuf_rele(db, FTAG);
1128         return (type == DMU_OT_NONE ? ENOENT : EEXIST);
1129     }
1130     mutex_exit(&dn->dn_mtx);

1132     if (refcount_add(&dn->dn_holds, tag) == 1)
1133         dbuf_add_ref(db, dnh);
1134     /* Now we can rely on the hold to prevent the dnode from moving. */
1135     zrl_remove(&dnh->dnh_zrlock);

1137     DNODE_VERIFY(dn);
1138     ASSERT3P(dn->dn_dbuf, ==, db);
1139     ASSERT3U(dn->dn_object, ==, object);
1140     dbuf_rele(db, FTAG);

1142     *dnp = dn;
1143     return (0);
1144 }
    unchanged_portion_omitted

```

```

*****
19546 Thu Apr 25 16:14:57 2013
new/usr/src/uts/common/fs/zfs/dnode_sync.c
3752 want more verifiable dbuf user eviction
Submitted by: Justin Gibbs <justing@spectralogic.com>
Submitted by: Will Andrews <willa@spectralogic.com>
*****
unchanged_portion_omitted

372 /*
373  * Try to kick all the dnodes dbufs out of the cache...
374  */
375 void
376 dnode_evict_dbufs(dnode_t *dn)
377 {
378     int progress;
379     int pass = 0;
380     list_t evict_list;

382     dmu_buf_create_user_evict_list(&evict_list);
383 #endif /* ! codereview */

385     do {
386         dmu_buf_impl_t *db, marker;
387         int evicting = FALSE;

389         progress = FALSE;
390         mutex_enter(&dn->dn_dbufs_mtx);
391         list_insert_tail(&dn->dn_dbufs, &marker);
392         db = list_head(&dn->dn_dbufs);
393         for (; db != &marker; db = list_head(&dn->dn_dbufs)) {
394             list_remove(&dn->dn_dbufs, db);
395             list_insert_tail(&dn->dn_dbufs, db);
396 #ifdef DEBUG
397             DB_DNODE_ENTER(db);
398             ASSERT3P(DB_DNODE(db), ==, dn);
399             DB_DNODE_EXIT(db);
400 #endif /* DEBUG */

402             mutex_enter(&db->db_mtx);
403             if (db->db_state == DB_EVICTING) {
404                 progress = TRUE;
405                 evicting = TRUE;
406                 mutex_exit(&db->db_mtx);
407             } else if (refcount_is_zero(&db->db_holds)) {
408                 progress = TRUE;
409                 dbuf_clear(db, &evict_list); /* exits db_mtx */
410                 dbuf_clear(db); /* exits db_mtx for us */
411             } else {
412                 mutex_exit(&db->db_mtx);
413             }
414             ASSERT(MUTEX_NOT_HELD(&db->db_mtx));
415             dmu_buf_process_user_evicts(&evict_list);
416         }
417         list_remove(&dn->dn_dbufs, &marker);
418         /*
419          * NB: we need to drop dn_dbufs_mtx between passes so
420          * that any DB_EVICTING dbufs can make progress.
421          * Ideally, we would have some cv we could wait on, but
422          * since we don't, just wait a bit to give the other
423          * thread a chance to run.
424          */
425         mutex_exit(&dn->dn_dbufs_mtx);
426         if (evicting)
427             delay(1);

```

```

427         pass++;
428         ASSERT(pass < 100); /* sanity check */
429     } while (progress);

431     rw_enter(&dn->dn_struct_rwlock, RW_WRITER);
432     if (dn->dn_bonus && refcount_is_zero(&dn->dn_bonus->db_holds)) {
433         mutex_enter(&dn->dn_bonus->db_mtx);
434         dbuf_evict(dn->dn_bonus, &evict_list);
435         dbuf_evict(dn->dn_bonus);
436         dn->dn_bonus = NULL;
437     }
438     rw_exit(&dn->dn_struct_rwlock);
439     dmu_buf_destroy_user_evict_list(&evict_list);
440 #endif /* ! codereview */
441 }

442 static void
443 dnode_undirty_dbufs(list_t *list)
444 {
445     dbuf_dirty_record_t *dr;

447     while (dr = list_head(list)) {
448         dmu_buf_impl_t *db = dr->dr_dbuf;
449         uint64_t txg = dr->dr_txg;

451         if (db->db_level != 0)
452             dnode_undirty_dbufs(&dr->dt.di.dr_children);

454         mutex_enter(&db->db_mtx);
455         /* XXX - use dbuf_undirty()? */
456         list_remove(list, dr);
457         ASSERT(db->db_last_dirty == dr);
458         db->db_last_dirty = NULL;
459         db->db_dirtycnt -= 1;
460         if (db->db_level == 0) {
461             ASSERT(db->db_blkid == DMU_BONUS_BLKID ||
462                 dr->dt.di.dr_data == db->db_buf);
463             dbuf_unoverride(dr);
464         }
465         kmem_free(dr, sizeof (dbuf_dirty_record_t));
466         dbuf_rele_and_unlock(db, (void *) (uintptr_t) txg);
467     }
468 }

470 static void
471 dnode_sync_free(dnode_t *dn, dmu_tx_t *tx)
472 {
473     int txgoff = tx->tx_txg & TXG_MASK;

475     ASSERT(dmu_tx_is_syncing(tx));

477     /*
478      * Our contents should have been freed in dnode_sync() by the
479      * free range record inserted by the caller of dnode_free().
480      */
481     ASSERT0(DN_USED_BYTES(dn->dn_phys));
482     ASSERT(BP_IS_HOLE(dn->dn_phys->dn_blkptr));

484     dnode_undirty_dbufs(&dn->dn_dirty_records[txgoff]);
485     dnode_evict_dbufs(dn);
486     ASSERT3P(list_head(&dn->dn_dbufs), ==, NULL);
487     ASSERT3P(dn->dn_bonus, ==, NULL);

489     /*
490      * XXX - It would be nice to assert this, but we may still
491      * have residual holds from async evictions from the arc...

```

```

492  *
493  * zfs_obj_to_path() also depends on this being
494  * commented out.
495  *
496  * ASSERT3U(refcount_count(&dn->dn_holds), ==, 1);
497  */
499  /* Undirty next bits */
500  dn->dn_next_nlevels[txgoff] = 0;
501  dn->dn_next_indblkshift[txgoff] = 0;
502  dn->dn_next_blkksz[txgoff] = 0;
504  /* ASSERT(blkptrs are zero); */
505  ASSERT(dn->dn_phys->dn_type != DMU_OT_NONE);
506  ASSERT(dn->dn_type != DMU_OT_NONE);
508  ASSERT(dn->dn_free_txg > 0);
509  if (dn->dn_allocated_txg != dn->dn_free_txg)
510      dbuf_will_dirty(dn->dn_dbuf, tx);
511  bzero(dn->dn_phys, sizeof (dnode_phys_t));
513  mutex_enter(&dn->dn_mtx);
514  dn->dn_type = DMU_OT_NONE;
515  dn->dn_maxblkid = 0;
516  dn->dn_allocated_txg = 0;
517  dn->dn_free_txg = 0;
518  dn->dn_have_spill = B_FALSE;
519  mutex_exit(&dn->dn_mtx);
521  ASSERT(dn->dn_object != DMU_META_DNODE_OBJECT);
523  dnode_rele(dn, (void *) (uintptr_t) tx->tx_txg);
524  /*
525   * Now that we've released our hold, the dnode may
526   * be evicted, so we musn't access it.
527   */
528 }
530 /*
531  * Write out the dnode's dirty buffers.
532  */
533 void
534 dnode_sync(dnode_t *dn, dmu_tx_t *tx)
535 {
536     free_range_t *rp;
537     dnode_phys_t *dnp = dn->dn_phys;
538     int txgoff = tx->tx_txg & TXG_MASK;
539     list_t *list = &dn->dn_dirty_records[txgoff];
540     static const dnode_phys_t zerodn = { 0 };
541     boolean_t kill_spill = B_FALSE;
543     ASSERT(dmu_tx_is_syncing(tx));
544     ASSERT(dnp->dn_type != DMU_OT_NONE || dn->dn_allocated_txg);
545     ASSERT(dnp->dn_type != DMU_OT_NONE ||
546            bcmp(dnp, &zerodn, DNODE_SIZE) == 0);
547     DNODE_VERIFY(dn);
549     ASSERT(dn->dn_dbuf == NULL || arc_released(dn->dn_dbuf->db_buf));
551     if (dmu_objset_userused_enabled(dn->dn_objset) &&
552         !DMU_OBJECT_IS_SPECIAL(dn->dn_object)) {
553         mutex_enter(&dn->dn_mtx);
554         dn->dn_oldused = DN_USED_BYTES(dn->dn_phys);
555         dn->dn_oldflags = dn->dn_phys->dn_flags;
556         dn->dn_phys->dn_flags |= DNODE_FLAG_USERUSED_ACCOUNTED;
557         mutex_exit(&dn->dn_mtx);

```

```

558     dmu_objset_userquota_get_ids(dn, B_FALSE, tx);
559 } else {
560     /* Once we account for it, we should always account for it. */
561     ASSERT(! (dn->dn_phys->dn_flags &
562              DNODE_FLAG_USERUSED_ACCOUNTED));
563 }
565 mutex_enter(&dn->dn_mtx);
566 if (dn->dn_allocated_txg == tx->tx_txg) {
567     /* The dnode is newly allocated or reallocated */
568     if (dnp->dn_type == DMU_OT_NONE) {
569         /* this is a first alloc, not a realloc */
570         dnp->dn_nlevels = 1;
571         dnp->dn_nblkptr = dn->dn_nblkptr;
572     }
574     dnp->dn_type = dn->dn_type;
575     dnp->dn_bonustype = dn->dn_bonustype;
576     dnp->dn_bonuslen = dn->dn_bonuslen;
577 }
579 ASSERT(dnp->dn_nlevels > 1 ||
580        BP_IS_HOLE(&dnp->dn_blkptr[0]) ||
581        BP_GET_LSIZE(&dnp->dn_blkptr[0]) ==
582        dnp->dn_datablkszsec << SPA_MINBLOCKSHIFT);
584 if (dn->dn_next_blkksz[txgoff]) {
585     ASSERT(P2PHASE(dn->dn_next_blkksz[txgoff],
586                  SPA_MINBLOCKSIZE) == 0);
587     ASSERT(BP_IS_HOLE(&dnp->dn_blkptr[0]) ||
588            dn->dn_maxblkid == 0 || list_head(list) != NULL ||
589            avl_last(&dn->dn_ranges[txgoff]) ||
590            dn->dn_next_blkksz[txgoff] >> SPA_MINBLOCKSHIFT ==
591            dnp->dn_datablkszsec);
592     dnp->dn_datablkszsec =
593         dn->dn_next_blkksz[txgoff] >> SPA_MINBLOCKSHIFT;
594     dn->dn_next_blkksz[txgoff] = 0;
595 }
597 if (dn->dn_next_bonuslen[txgoff]) {
598     if (dn->dn_next_bonuslen[txgoff] == DN_ZERO_BONUSLEN)
599         dnp->dn_bonuslen = 0;
600     else
601         dnp->dn_bonuslen = dn->dn_next_bonuslen[txgoff];
602     ASSERT(dnp->dn_bonuslen <= DN_MAX_BONUSLEN);
603     dn->dn_next_bonuslen[txgoff] = 0;
604 }
606 if (dn->dn_next_bonustype[txgoff]) {
607     ASSERT(DMU_OT_IS_VALID(dn->dn_next_bonustype[txgoff]));
608     dnp->dn_bonustype = dn->dn_next_bonustype[txgoff];
609     dn->dn_next_bonustype[txgoff] = 0;
610 }
612 /*
613  * We will either remove a spill block when a file is being removed
614  * or we have been asked to remove it.
615  */
616 if (dn->dn_rm_spillblk[txgoff] ||
617     ((dnp->dn_flags & DNODE_FLAG_SPILL_BLKPTR) &&
618     dn->dn_free_txg > 0 && dn->dn_free_txg <= tx->tx_txg)) {
619     if ((dnp->dn_flags & DNODE_FLAG_SPILL_BLKPTR))
620         kill_spill = B_TRUE;
621     dn->dn_rm_spillblk[txgoff] = 0;
622 }

```

```

624     if (dn->dn_next_indblkshift[txgoff]) {
625         ASSERT(dnp->dn_nlevels == 1);
626         dnp->dn_indblkshift = dn->dn_next_indblkshift[txgoff];
627         dn->dn_next_indblkshift[txgoff] = 0;
628     }
629
630     /*
631      * Just take the live (open-context) values for checksum and compress.
632      * Strictly speaking it's a future leak, but nothing bad happens if we
633      * start using the new checksum or compress algorithm a little early.
634      */
635     dnp->dn_checksum = dn->dn_checksum;
636     dnp->dn_compress = dn->dn_compress;
637
638     mutex_exit(&dn->dn_mtx);
639
640     if (kill_spill) {
641         (void) free_blocks(dn, &dn->dn_phys->dn_spill, 1, tx);
642         mutex_enter(&dn->dn_mtx);
643         dnp->dn_flags &= ~DNODE_FLAG_SPILL_BLKPTR;
644         mutex_exit(&dn->dn_mtx);
645     }
646
647     /* process all the "freed" ranges in the file */
648     while (rp = avl_last(&dn->dn_ranges[txgoff])) {
649         dnode_sync_free_range(dn, rp->fr_blkid, rp->fr_nblks, tx);
650         /* grab the mutex so we don't race with dnode_block_freed() */
651         mutex_enter(&dn->dn_mtx);
652         avl_remove(&dn->dn_ranges[txgoff], rp);
653         mutex_exit(&dn->dn_mtx);
654         kmem_free(rp, sizeof (free_range_t));
655     }
656
657     if (dn->dn_free_txg > 0 && dn->dn_free_txg <= tx->tx_txg) {
658         dnode_sync_free(dn, tx);
659         return;
660     }
661
662     if (dn->dn_next_nblkptr[txgoff]) {
663         /* this should only happen on a realloc */
664         ASSERT(dn->dn_allocated_txg == tx->tx_txg);
665         if (dn->dn_next_nblkptr[txgoff] > dnp->dn_nblkptr) {
666             /* zero the new blkptrs we are gaining */
667             bzero(dnp->dn_blkptr + dnp->dn_nblkptr,
668                 sizeof (blkptr_t) *
669                 (dn->dn_next_nblkptr[txgoff] - dnp->dn_nblkptr));
670 #ifdef ZFS_DEBUG
671         } else {
672             int i;
673             ASSERT(dn->dn_next_nblkptr[txgoff] < dnp->dn_nblkptr);
674             /* the blkptrs we are losing better be unallocated */
675             for (i = dn->dn_next_nblkptr[txgoff];
676                 i < dnp->dn_nblkptr; i++)
677                 ASSERT(BP_IS_HOLE(&dnp->dn_blkptr[i]));
678 #endif
679         }
680         mutex_enter(&dn->dn_mtx);
681         dnp->dn_nblkptr = dn->dn_next_nblkptr[txgoff];
682         dn->dn_next_nblkptr[txgoff] = 0;
683         mutex_exit(&dn->dn_mtx);
684     }
685
686     if (dn->dn_next_nlevels[txgoff]) {
687         dnode_increase_indirection(dn, tx);
688         dn->dn_next_nlevels[txgoff] = 0;
689     }

```

```

691     dbuf_sync_list(list, tx);
692
693     if (!DMU_OBJECT_IS_SPECIAL(dn->dn_object)) {
694         ASSERT3P(list_head(list), ==, NULL);
695         dnode_rele(dn, (void *) (uintptr_t) tx->tx_txg);
696     }
697
698     /*
699      * Although we have dropped our reference to the dnode, it
700      * can't be evicted until its written, and we haven't yet
701      * initiated the IO for the dnode's dbuf.
702      */
703 }

```

```

*****
80781 Thu Apr 25 16:14:57 2013
new/usr/src/uts/common/fs/zfs/dsl_dataset.c
3752 want more verifiable dbuf user eviction
Submitted by: Justin Gibbs <justing@spectrallogic.com>
Submitted by: Will Andrews <willa@spectrallogic.com>
*****
unchanged_portion_omitted

242 /* ARGSUSED */
243 static void
244 dsl_dataset_evict_impl(dsl_dataset_t *ds, boolean_t evict_deadlist)
245 {
246     dsl_dataset_t *ds = ds;
247
248     ASSERT(ds->ds_owner == NULL);
249
250     ds->ds_dbuf = NULL;
251
252 #endif /* ! codereview */
253     unique_remove(ds->ds_fsid_guid);
254
255     if (ds->ds_objset != NULL)
256         dmu_objset_evict(ds->ds_objset);
257
258     if (ds->ds_prev) {
259         dsl_dataset_rele(ds->ds_prev, ds);
260         ds->ds_prev = NULL;
261     }
262
263     bplist_destroy(&ds->ds_pending_deadlist);
264     if (evict_deadlist)
265         if (ds->ds_phys->ds_deadlist_obj != 0)
266             dsl_deadlist_close(&ds->ds_deadlist);
267     if (ds->ds_dir)
268         dsl_dir_rele(ds->ds_dir, ds);
269
270     ASSERT(!list_link_active(&ds->ds_synced_link));
271
272     mutex_destroy(&ds->ds_lock);
273     mutex_destroy(&ds->ds_opening_lock);
274     refcount_destroy(&ds->ds_longholds);
275
276     kmem_free(ds, sizeof (dsl_dataset_t));
277 }
278
279 /* ARGSUSED */
280 static void
281 dsl_dataset_evict(dmu_buf_user_t *dbu)
282 {
283     dsl_dataset_evict_impl((dsl_dataset_t *)dbu, B_TRUE);
284 }
285
286 #endif /* ! codereview */
287 int
288 dsl_dataset_get_snapname(dsl_dataset_t *ds)
289 {
290     dsl_dataset_phys_t *headphys;
291     int err;
292     dmu_buf_t *headdbuf;
293     dsl_pool_t *dp = ds->ds_dir->dd_pool;
294     objset_t *mos = dp->dp_meta_objset;
295
296     if (ds->ds_snapname[0])
297         return (0);

```

```

296     if (ds->ds_phys->ds_next_snap_obj == 0)
297         return (0);
298
299     err = dmu_bonus_hold(mos, ds->ds_dir->dd_phys->dd_head_dataset_obj,
300         FTAG, &headdbuf);
301     if (err != 0)
302         return (err);
303     headphys = headdbuf->db_data;
304     err = zap_value_search(dp->dp_meta_objset,
305         headphys->ds_snapnames_zapobj, ds->ds_object, 0, ds->ds_snapname);
306     dmu_buf_rele(headdbuf, FTAG);
307     return (err);
308 }
309
310 int
311 dsl_dataset_snap_lookup(dsl_dataset_t *ds, const char *name, uint64_t *value)
312 {
313     objset_t *mos = ds->ds_dir->dd_pool->dp_meta_objset;
314     uint64_t snapobj = ds->ds_phys->ds_snapnames_zapobj;
315     matchtype_t mt;
316     int err;
317
318     if (ds->ds_phys->ds_flags & DS_FLAG_CI_DATASET)
319         mt = MT_FIRST;
320     else
321         mt = MT_EXACT;
322
323     err = zap_lookup_norm(mos, snapobj, name, 8, 1,
324         value, mt, NULL, 0, NULL);
325     if (err == ENOTSUP && mt == MT_FIRST)
326         err = zap_lookup(mos, snapobj, name, 8, 1, value);
327     return (err);
328 }
329
330 int
331 dsl_dataset_snap_remove(dsl_dataset_t *ds, const char *name, dmu_tx_t *tx)
332 {
333     objset_t *mos = ds->ds_dir->dd_pool->dp_meta_objset;
334     uint64_t snapobj = ds->ds_phys->ds_snapnames_zapobj;
335     matchtype_t mt;
336     int err;
337
338     dsl_dir_snap_cmtime_update(ds->ds_dir);
339
340     if (ds->ds_phys->ds_flags & DS_FLAG_CI_DATASET)
341         mt = MT_FIRST;
342     else
343         mt = MT_EXACT;
344
345     err = zap_remove_norm(mos, snapobj, name, mt, tx);
346     if (err == ENOTSUP && mt == MT_FIRST)
347         err = zap_remove(mos, snapobj, name, tx);
348     return (err);
349 }
350
351 int
352 dsl_dataset_hold_obj(dsl_pool_t *dp, uint64_t dsobj, void *tag,
353     dsl_dataset_t **dsp)
354 {
355     objset_t *mos = dp->dp_meta_objset;
356     dmu_buf_t *dbuf;
357     dsl_dataset_t *ds;
358     int err;
359     dmu_object_info_t doi;
360
361     ASSERT(dsl_pool_config_held(dp));

```

```

363     err = dmu_bonus_hold(mos, dsobj, tag, &dbuf);
364     if (err != 0)
365         return (err);

367     /* Make sure dsobj has the correct object type. */
368     dmu_object_info_from_db(dbuf, &doi);
369     if (doi.doi_type != DMU_OT_DSL_DATASET)
370         return (SET_ERROR(EINVAL));

372     ds = (dsl_dataset_t *)dmu_buf_get_user(dbuf);
264     ds = dmu_buf_get_user(dbuf);
373     if (ds == NULL) {
374         dsl_dataset_t *winner = NULL;

376         ds = kmem_zalloc(sizeof (dsl_dataset_t), KM_SLEEP);
377         ds->ds_dbuf = dbuf;
378         ds->ds_object = dsobj;
271         ds->ds_phys = dbuf->db_data;

380         mutex_init(&ds->ds_lock, NULL, MUTEX_DEFAULT, NULL);
381         mutex_init(&ds->ds_opening_lock, NULL, MUTEX_DEFAULT, NULL);
382         mutex_init(&ds->ds_sendstream_lock, NULL, MUTEX_DEFAULT, NULL);
383         refcount_create(&ds->ds_longholds);

385         bplist_create(&ds->ds_pending_deadlist);
386         dsl_deadlist_open(&ds->ds_deadlist,
387             mos, ds->ds_phys->ds_deadlist_obj);

389         list_create(&ds->ds_sendstreams, sizeof (dmu_sendarg_t),
390             offsetof(dmu_sendarg_t, dsa_link));

392         if (err == 0) {
393             err = dsl_dir_hold_obj(dp,
394                 ds->ds_phys->ds_dir_obj, NULL, ds, &ds->ds_dir);
395         }
396         if (err != 0) {
397             mutex_destroy(&ds->ds_lock);
398             mutex_destroy(&ds->ds_opening_lock);
399             refcount_destroy(&ds->ds_longholds);
400             bplist_destroy(&ds->ds_pending_deadlist);
401             dsl_deadlist_close(&ds->ds_deadlist);
402             kmem_free(ds, sizeof (dsl_dataset_t));
403             dmu_buf_rele(dbuf, tag);
404             return (err);
405         }

407         if (!dsl_dataset_is_snapshot(ds)) {
408             ds->ds_snapname[0] = '\0';
409             if (ds->ds_phys->ds_prev_snap_obj != 0) {
410                 err = dsl_dataset_hold_obj(dp,
411                     ds->ds_phys->ds_prev_snap_obj,
412                     ds, &ds->ds_prev);
413             }
414         } else {
415             if (zfs_flags & ZFS_DEBUG_SNAPNAMES)
416                 err = dsl_dataset_get_snapname(ds);
417             if (err == 0 && ds->ds_phys->ds_userrefs_obj != 0) {
418                 err = zap_count(
419                     ds->ds_dir->dd_pool->dp_meta_objset,
420                     ds->ds_phys->ds_userrefs_obj,
421                     &ds->ds_userrefs);
422             }
423         }

425         if (err == 0 && !dsl_dataset_is_snapshot(ds)) {

```

```

426         err = dsl_prop_get_int(ds,
427             zfs_prop_to_name(ZFS_PROP_REFRESERVATION),
428             &ds->ds_reserved);
429         if (err == 0) {
430             err = dsl_prop_get_int(ds,
431                 zfs_prop_to_name(ZFS_PROP_REFQUOTA),
432                 &ds->ds_quota);
433         }
434     } else {
435         ds->ds_reserved = ds->ds_quota = 0;
436     }

438     dmu_buf_init_user(&ds->db_evict, dsl_dataset_evict);
439     if (err == 0)
440         winner = (dsl_dataset_t *)
441             dmu_buf_set_user_ie(dbuf, &ds->db_evict);

443     if (err || winner) {
331     if (err != 0 || (winner = dmu_buf_set_user_ie(dbuf, ds,
332         &ds->ds_phys, dsl_dataset_evict)) != NULL) {
444         bplist_destroy(&ds->ds_pending_deadlist);
445         dsl_deadlist_close(&ds->ds_deadlist);
446         if (ds->ds_prev)
447             dsl_dataset_rele(ds->ds_prev, ds);
448         dsl_dir_rele(ds->ds_dir, ds);
449         mutex_destroy(&ds->ds_lock);
450         mutex_destroy(&ds->ds_opening_lock);
451         refcount_destroy(&ds->ds_longholds);
452         kmem_free(ds, sizeof (dsl_dataset_t));
453         if (err != 0) {
454             dmu_buf_rele(dbuf, tag);
455             return (err);
456         }
457         ds = winner;
458     } else {
459         ds->ds_fsid_guid =
460             unique_insert(ds->ds_phys->ds_fsid_guid);
461     }
462     ASSERT3P(ds->ds_dbuf, ==, dbuf);
463     ASSERT3P(ds->ds_phys, ==, dbuf->db_data);
464     ASSERT(ds->ds_phys->ds_prev_snap_obj != 0 ||
465         spa_version(dp->dp_spa) < SPA_VERSION_ORIGIN ||
466         dp->dp_origin_snap == NULL || ds == dp->dp_origin_snap);
467     *dsp = ds;
468     return (0);
469 }
470 }
    unchanged portion omitted

636 void
637 dsl_dataset_disown(dsl_dataset_t *ds, void *tag)
638 {
639     ASSERT(ds->ds_owner == tag && ds->ds_dbuf != NULL);

641     mutex_enter(&ds->ds_lock);
642     ds->ds_owner = NULL;
643     mutex_exit(&ds->ds_lock);
644     dsl_dataset_long_rele(ds, tag);
645     if (ds->ds_dbuf != NULL)
646         dsl_dataset_rele(ds, tag);
647     else
648         dsl_dataset_evict_impl(ds, B_FALSE);
537     dsl_dataset_evict(NULL, ds);
649 }
    unchanged portion omitted

```

```

*****
35255 Thu Apr 25 16:14:58 2013
new/usr/src/uts/common/fs/zfs/dsl_dir.c
3752 want more verifiable dbuf user eviction
Submitted by: Justin Gibbs <justing@spectralogic.com>
Submitted by: Will Andrews <willa@spectralogic.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 */

26 #include <sys/dmu.h>
27 #include <sys/dmu_objset.h>
28 #include <sys/dmu_tx.h>
29 #include <sys/dsl_dataset.h>
30 #include <sys/dsl_dir.h>
31 #include <sys/dsl_prop.h>
32 #include <sys/dsl_synctask.h>
33 #include <sys/dsl_deleg.h>
34 #include <sys/spa.h>
35 #include <sys/metaslabs.h>
36 #include <sys/zap.h>
37 #include <sys/zio.h>
38 #include <sys/arc.h>
39 #include <sys/sunddi.h>
40 #include "zfs_namecheck.h"

42 static uint64_t dsl_dir_space_towrite(dsl_dir_t *dd);

44 /* ARGSUSED */
45 static void
46 dsl_dir_evict(dmu_buf_user_t *dbu)
47 dsl_dir_evict(dmu_buf_t *db, void *arg)
48 {
49     dsl_dir_t *dd = (dsl_dir_t *)dbu;
50     dsl_dir_t *dd = arg;
51     dsl_pool_t *dp = dd->dd_pool;
52     int t;

54     dd->dd_dbuf = NULL;

54 #endif /* ! codereview */
55     for (t = 0; t < TXG_SIZE; t++) {
56         ASSERT(!txg_list_member(&dp->dp_dirty_dirs, dd, t));
57         ASSERT(dd->dd_tempreserved[t] == 0);

```

```

58         ASSERT(dd->dd_space_towrite[t] == 0);
59     }

61     if (dd->dd_parent)
62         dsl_dir_rele(dd->dd_parent, dd);

64     spa_close(dd->dd_pool->dp_spa, dd);

66     /*
67      * The props callback list should have been cleaned up by
68      * objset_evict().
69      */
70     list_destroy(&dd->dd_prop_cbs);
71     mutex_destroy(&dd->dd_lock);
72     kmem_free(dd, sizeof (dsl_dir_t));
73 }

75 int
76 dsl_dir_hold_obj(dsl_pool_t *dp, uint64_t ddojb,
77                 const char *tail, void *tag, dsl_dir_t **ddp)
78 {
79     dmu_buf_t *dbuf;
80     dsl_dir_t *dd;
81     int err;

83     ASSERT(dsl_pool_config_held(dp));

85     err = dmu_bonus_hold(dp->dp_meta_objset, ddojb, tag, &dbuf);
86     if (err != 0)
87         return (err);
88     dd = (dsl_dir_t *)dmu_buf_get_user(dbuf);
89     dd = dmu_buf_get_user(dbuf);
89 #ifdef ZFS_DEBUG
90     {
91         dmu_object_info_t doi;
92         dmu_object_info_from_db(dbuf, &doi);
93         ASSERT3U(doi.doi_type, ==, DMU_OT_DSL_DIR);
94         ASSERT3U(doi.doi_bonus_size, >=, sizeof (dsl_dir_phys_t));
95     }
96 #endif
97     if (dd == NULL) {
98         dsl_dir_t *winner;

100         dd = kmem_zalloc(sizeof (dsl_dir_t), KM_SLEEP);
101         dd->dd_object = ddojb;
102         dd->dd_dbuf = dbuf;
103         dd->dd_pool = dp;
68         dd->dd_phys = dbuf->db_data;
104         mutex_init(&dd->dd_lock, NULL, MUTEX_DEFAULT, NULL);

106         list_create(&dd->dd_prop_cbs, sizeof (dsl_prop_cb_record_t),
107                   offsetof(dsl_prop_cb_record_t, cbr_node));

109         dsl_dir_snap_cmtime_update(dd);

111         if (dd->dd_phys->dd_parent_obj) {
112             err = dsl_dir_hold_obj(dp, dd->dd_phys->dd_parent_obj,
113                                   NULL, dd, &dd->dd_parent);
114             if (err != 0)
115                 goto errout;
116             if (tail) {
117 #ifdef ZFS_DEBUG
118                 uint64_t foundobj;

120                 err = zap_lookup(dp->dp_meta_objset,
121                                 dd->dd_parent->dd_phys->dd_child_dir_zapobj,

```

```

122         tail, sizeof(foundobj), 1, &foundobj);
123         ASSERT(err || foundobj == ddoobj);
124 #endif
125         (void) strcpy(dd->dd_myname, tail);
126     } else {
127         err = zap_value_search(dp->dp_meta_objset,
128             dd->dd_parent->dd_phys->dd_child_dir_zapobj,
129             ddoobj, 0, dd->dd_myname);
130     }
131     if (err != 0)
132         goto errout;
133 } else {
134     (void) strcpy(dd->dd_myname, spa_name(dp->dp_spa));
135 }
136
137 if (dsl_dir_is_clone(dd)) {
138     dmu_buf_t *origin_bonus;
139     dsl_dataset_phys_t *origin_phys;
140
141     /*
142      * We can't open the origin dataset, because
143      * that would require opening this dsl_dir.
144      * Just look at its phys directly instead.
145      */
146     err = dmu_bonus_hold(dp->dp_meta_objset,
147         dd->dd_phys->dd_origin_obj, FTAG, &origin_bonus);
148     if (err != 0)
149         goto errout;
150     origin_phys = origin_bonus->db_data;
151     dd->dd_origin_txg =
152         origin_phys->ds_creation_txg;
153     dmu_buf_rele(origin_bonus, FTAG);
154 }
155
156 dmu_buf_init_user(&dd->db_evict, dsl_dir_evict);
157 winner = (dsl_dir_t *)dmu_buf_set_user_ie(dbuf, &dd->db_evict);
158 winner = dmu_buf_set_user_ie(dbuf, dd, &dd->dd_phys,
159     dsl_dir_evict);
160 if (winner) {
161     if (dd->dd_parent)
162         dsl_dir_rele(dd->dd_parent, dd);
163     mutex_destroy(&dd->dd_lock);
164     kmem_free(dd, sizeof(dsl_dir_t));
165     dd = winner;
166 } else {
167     spa_open_ref(dp->dp_spa, dd);
168 }
169
170 /*
171  * The dsl_dir_t has both open-to-close and instantiate-to-evict
172  * holds on the spa. We need the open-to-close holds because
173  * otherwise the spa_refcnt wouldn't change when we open a
174  * dir which the spa also has open, so we could incorrectly
175  * think it was OK to unload/export/destroy the pool. We need
176  * the instantiate-to-evict hold because the dsl_dir_t has a
177  * pointer to the dd_pool, which has a pointer to the spa_t.
178  */
179 spa_open_ref(dp->dp_spa, tag);
180 ASSERT3P(dd->dd_pool, ==, dp);
181 ASSERT3U(dd->dd_object, ==, ddoobj);
182 ASSERT3P(dd->dd_dbuf, ==, dbuf);
183 *ddp = dd;
184 return (0);
185 errout:

```

```

186     if (dd->dd_parent)
187         dsl_dir_rele(dd->dd_parent, dd);
188     mutex_destroy(&dd->dd_lock);
189     kmem_free(dd, sizeof(dsl_dir_t));
190     dmu_buf_rele(dbuf, tag);
191     return (err);
192 }

```

unchanged_portion_omitted

```

*****
29181 Thu Apr 25 16:14:58 2013
new/usr/src/uts/common/fs/zfs/dsl_prop.c
3752 want more verifiable dbuf user eviction
Submitted by: Justin Gibbs <justing@spectrallogic.com>
Submitted by: Will Andrews <willa@spectrallogic.com>
*****
unchanged_portion_omitted_

159 int
160 dsl_prop_get_ds(dsl_dataset_t *ds, const char *propname,
161                int intsz, int numints, void *buf, char *setpoint)
162 {
163     zfs_prop_t prop = zfs_name_to_prop(propname);
164     boolean_t inheritable;
165     boolean_t snapshot;
166     uint64_t zapobj;

168     ASSERT(dsl_pool_config_held(ds->ds_dir->dd_pool));
169     inheritable = (prop == ZPROP_INVALID || zfs_prop_inheritable(prop));
170     snapshot = (DS_HAS_PHYS(ds) && dsl_dataset_is_snapshot(ds));
171     zapobj = (DS_HAS_PHYS(ds) ? ds->ds_phys->ds_props_obj : 0);
170     snapshot = (ds->ds_phys != NULL && dsl_dataset_is_snapshot(ds));
171     zapobj = (ds->ds_phys == NULL ? 0 : ds->ds_phys->ds_props_obj);

173     if (zapobj != 0) {
174         objset_t *mos = ds->ds_dir->dd_pool->dp_meta_objset;
175         int err;

177         ASSERT(snapshot);

179         /* Check for a local value. */
180         err = zap_lookup(mos, zapobj, propname, intsz, numints, buf);
181         if (err != ENOENT) {
182             if (setpoint != NULL && err == 0)
183                 dsl_dataset_name(ds, setpoint);
184             return (err);
185         }

187         /*
188          * Skip the check for a received value if there is an explicit
189          * inheritance entry.
190          */
191         if (inheritable) {
192             char *inheritstr = kmem_asprintf("%s%s", propname,
193                                             ZPROP_INHERIT_SUFFIX);
194             err = zap_contains(mos, zapobj, inheritstr);
195             strfree(inheritstr);
196             if (err != 0 && err != ENOENT)
197                 return (err);
198         }

200         if (err == ENOENT) {
201             /* Check for a received value. */
202             char *recvdstr = kmem_asprintf("%s%s", propname,
203                                             ZPROP_RECVD_SUFFIX);
204             err = zap_lookup(mos, zapobj, recvdstr,
205                             intsz, numints, buf);
206             strfree(recvdstr);
207             if (err != ENOENT) {
208                 if (setpoint != NULL && err == 0)
209                     (void) strcpy(setpoint,
210                                   ZPROP_SOURCE_VAL_RECVD);
211                 return (err);
212             }
213         }
}

```

```

214     }

216     return (dsl_prop_get_dd(ds->ds_dir, propname,
217                             intsz, numints, buf, setpoint, snapshot));
218 }
unchanged_portion_omitted_

528 void
529 dsl_prop_set_sync_impl(dsl_dataset_t *ds, const char *propname,
530                        zprop_source_t source, int intsz, int numints, const void *value,
531                        dmu_tx_t *tx)
532 {
533     objset_t *mos = ds->ds_dir->dd_pool->dp_meta_objset;
534     uint64_t zapobj, intval, dummy;
535     int isint;
536     char valbuf[32];
537     const char *valstr = NULL;
538     char *inheritstr;
539     char *recvdstr;
540     char *tbuf = NULL;
541     int err;
542     uint64_t version = spa_version(ds->ds_dir->dd_pool->dp_spa);

544     isint = (dodefult(propname, 8, 1, &intval) == 0);

546     if (DS_HAS_PHYS(ds) && dsl_dataset_is_snapshot(ds)) {
546         if (ds->ds_phys != NULL && dsl_dataset_is_snapshot(ds)) {
547             ASSERT(version >= SPA_VERSION_SNAP_PROPS);
548             if (ds->ds_phys->ds_props_obj == 0) {
549                 dmu_buf_will_dirty(ds->ds_dbuf, tx);
550                 ds->ds_phys->ds_props_obj =
551                     zap_create(mos,
552                               DMU_OT_DSL_PROPS, DMU_OT_NONE, 0, tx);
553             }
554             zapobj = ds->ds_phys->ds_props_obj;
555         } else {
556             zapobj = ds->ds_dir->dd_phys->dd_props_zapobj;
557         }

559         if (version < SPA_VERSION_RECVD_PROPS) {
560             zfs_prop_t prop = zfs_name_to_prop(propname);
561             if (prop == ZFS_PROP_QUOTA || prop == ZFS_PROP_RESERVATION)
562                 return;

564             if (source & ZPROP_SRC_NONE)
565                 source = ZPROP_SRC_NONE;
566             else if (source & ZPROP_SRC_RECEIVED)
567                 source = ZPROP_SRC_LOCAL;
568         }

570         inheritstr = kmem_asprintf("%s%s", propname, ZPROP_INHERIT_SUFFIX);
571         recvdstr = kmem_asprintf("%s%s", propname, ZPROP_RECVD_SUFFIX);

573         switch (source) {
574             case ZPROP_SRC_NONE:
575                 /*
576                  * revert to received value, if any (inherit -S)
577                  * - remove propname
578                  * - remove propname$inherit
579                  */
580                 err = zap_remove(mos, zapobj, propname, tx);
581                 ASSERT(err == 0 || err == ENOENT);
582                 err = zap_remove(mos, zapobj, inheritstr, tx);
583                 ASSERT(err == 0 || err == ENOENT);
584                 break;
585             case ZPROP_SRC_LOCAL:

```

```

586     /*
587     * remove propname$inherit
588     * set propname -> value
589     */
590     err = zap_remove(mos, zapobj, inheritstr, tx);
591     ASSERT(err == 0 || err == ENOENT);
592     VERIFY0(zap_update(mos, zapobj, propname,
593         intsz, numints, value, tx));
594     break;
595 case ZPROP_SRC_INHERITED:
596     /*
597     * explicitly inherit
598     * - remove propname
599     * - set propname$inherit
600     */
601     err = zap_remove(mos, zapobj, propname, tx);
602     ASSERT(err == 0 || err == ENOENT);
603     if (version >= SPA_VERSION_RECVD_PROPS &&
604         dsl_prop_get_int_ds(ds, ZPROP_HAS_RECVD, &dummy) == 0) {
605         dummy = 0;
606         VERIFY0(zap_update(mos, zapobj, inheritstr,
607             8, 1, &dummy, tx));
608     }
609     break;
610 case ZPROP_SRC_RECEIVED:
611     /*
612     * set propname$recvd -> value
613     */
614     err = zap_update(mos, zapobj, recvdstr,
615         intsz, numints, value, tx);
616     ASSERT(err == 0);
617     break;
618 case (ZPROP_SRC_NONE | ZPROP_SRC_LOCAL | ZPROP_SRC_RECEIVED):
619     /*
620     * clear local and received settings
621     * - remove propname
622     * - remove propname$inherit
623     * - remove propname$recvd
624     */
625     err = zap_remove(mos, zapobj, propname, tx);
626     ASSERT(err == 0 || err == ENOENT);
627     err = zap_remove(mos, zapobj, inheritstr, tx);
628     ASSERT(err == 0 || err == ENOENT);
629     /* FALLTHRU */
630 case (ZPROP_SRC_NONE | ZPROP_SRC_RECEIVED):
631     /*
632     * remove propname$recvd
633     */
634     err = zap_remove(mos, zapobj, recvdstr, tx);
635     ASSERT(err == 0 || err == ENOENT);
636     break;
637 default:
638     cmn_err(CE_PANIC, "unexpected property source: %d", source);
639 }

641     strfree(inheritstr);
642     strfree(recvdstr);

644     if (isint) {
645         VERIFY0(dsl_prop_get_int_ds(ds, propname, &intval));

647         if (DS_HAS_PHYS(ds) && dsl_dataset_is_snapshot(ds)) {
648             if (ds->ds_phys != NULL && dsl_dataset_is_snapshot(ds)) {
649                 dsl_prop_cb_record_t *cbr;
650                 /*
651                  * It's a snapshot; nothing can inherit this

```

```

651         * property, so just look for callbacks on this
652         * ds here.
653         */
654         mutex_enter(&ds->ds_dir->dd_lock);
655         for (cbr = list_head(&ds->ds_dir->dd_prop_cbs); cbr;
656             cbr = list_next(&ds->ds_dir->dd_prop_cbs, cbr)) {
657             if (cbr->cbr_ds == ds &&
658                 strcmp(cbr->cbr_propname, propname) == 0)
659                 cbr->cbr_func(cbr->cbr_arg, intval);
660         }
661         mutex_exit(&ds->ds_dir->dd_lock);
662     } else {
663         dsl_prop_changed_notify(ds->ds_dir->dd_pool,
664             ds->ds_dir->dd_object, propname, intval, TRUE);
665     }

667     (void) snprintf(valbuf, sizeof(valbuf),
668         "%lld", (longlong_t)intval);
669     valstr = valbuf;
670 } else {
671     if (source == ZPROP_SRC_LOCAL) {
672         valstr = value;
673     } else {
674         tbuf = kmem_alloc(ZAP_MAXVALUELEN, KM_SLEEP);
675         if (dsl_prop_get_ds(ds, propname, 1,
676             ZAP_MAXVALUELEN, tbuf, NULL) == 0)
677             valstr = tbuf;
678     }
679 }

681     spa_history_log_internal_ds(ds, (source == ZPROP_SRC_NONE ||
682         source == ZPROP_SRC_INHERITED) ? "inherit" : "set", tx,
683         "%s=%s", propname, (valstr == NULL ? "" : valstr));

685     if (tbuf != NULL)
686         kmem_free(tbuf, ZAP_MAXVALUELEN);
687 }

```

unchanged_portion_omitted

```

*****
51823 Thu Apr 25 16:14:58 2013
new/usr/src/uts/common/fs/zfs/sa.c
3752 want more verifiable dbuf user eviction
Submitted by: Justin Gibbs <justing@spectralogic.com>
Submitted by: Will Andrews <willa@spectralogic.com>
*****
_____unchanged_portion_omitted_____

1300 /*ARGSUSED*/
1301 void
1302 sa_evict(dmu_buf_user_t *dbu)
1303 {
1304     panic("evicting sa dbuf\n");
1305     panic("evicting sa dbuf %p\n", (void *)db);
1306 }
_____unchanged_portion_omitted_____

1340 void
1341 sa_handle_destroy(sa_handle_t *hdl)
1342 {
1343     dmu_buf_t *db = hdl->sa_bonus;

1344     #endif /* ! codereview */
1345     mutex_enter(&hdl->sa_lock);
1346     (void) dmu_buf_remove_user(db, &hdl->db_evict);
1347     (void) dmu_buf_update_user((dmu_buf_t *)hdl->sa_bonus, hdl,
1348         NULL, NULL, NULL);

1349     if (hdl->sa_bonus_tab) {
1350         sa_idx_tab_rele(hdl->sa_os, hdl->sa_bonus_tab);
1351         hdl->sa_bonus_tab = NULL;
1352     }
1353     if (hdl->sa_spill_tab) {
1354         sa_idx_tab_rele(hdl->sa_os, hdl->sa_spill_tab);
1355         hdl->sa_spill_tab = NULL;
1356     }

1358     dmu_buf_rele(hdl->sa_bonus, NULL);

1360     if (hdl->sa_spill)
1361         dmu_buf_rele((dmu_buf_t *)hdl->sa_spill, NULL);
1362     mutex_exit(&hdl->sa_lock);

1364     kmem_cache_free(sa_cache, hdl);
1365 }

1367 int
1368 sa_handle_get_from_db(objset_t *os, dmu_buf_t *db, void *userp,
1369     sa_handle_type_t hdl_type, sa_handle_t **handlepp)
1370 {
1371     int error = 0;
1372     dmu_object_info_t doi;
1373     sa_handle_t *handle = NULL;
1374     sa_handle_t *handle;

1375 #ifdef ZFS_DEBUG
1376     dmu_object_info_from_db(db, &doi);
1377     ASSERT(doi.doi_bonus_type == DMU_OT_SA ||
1378         doi.doi_bonus_type == DMU_OT_ZNODE);
1379 #endif
1380     /* find handle, if it exists */
1381     /* if one doesn't exist then create a new one, and initialize it */

1383     if (hdl_type == SA_HDL_SHARED)

```

```

1384     handle = (sa_handle_t *)dmu_buf_get_user(db);

1380     handle = (hdl_type == SA_HDL_SHARED) ? dmu_buf_get_user(db) : NULL;
1386     if (handle == NULL) {
1387         sa_handle_t *winner = NULL;

1389         bzero(&handle->db_evict, sizeof(dmu_buf_user_t));
1382         sa_handle_t *newhandle;
1390         handle = kmem_cache_alloc(sa_cache, KM_SLEEP);
1391         handle->sa_userp = userp;
1392         handle->sa_bonus = db;
1393         handle->sa_os = os;
1394         handle->sa_spill = NULL;

1396         error = sa_build_index(handle, SA_BONUS);
1397         if (hdl_type == SA_HDL_SHARED) {
1398             dmu_buf_init_user(&handle->db_evict, sa_evict);
1399             winner = (sa_handle_t *)
1400                 dmu_buf_set_user_ie(db, &handle->db_evict);
1401         }
1402         if (winner != NULL) {
1390             newhandle = (hdl_type == SA_HDL_SHARED) ?
1391                 dmu_buf_set_user_ie(db, handle,
1392                 NULL, sa_evict) : NULL;

1394             if (newhandle != NULL) {
1403                 kmem_cache_free(sa_cache, handle);
1404                 handle = winner;
1405                 handle = newhandle;
1406             }
1407             *handlepp = handle;

1409             return (error);
1410 }
_____unchanged_portion_omitted_____

1914 void
1915 sa_update_user(sa_handle_t *newhdl, sa_handle_t *oldhdl)
1916 {
1917     dmu_buf_user_t *new_user = &newhdl->db_evict;
1918     dmu_buf_user_t *old_user = &oldhdl->db_evict;

1920     dmu_buf_init_user(new_user, sa_evict);
1921     VERIFY(dmu_buf_replace_user(newhdl->sa_bonus, old_user,
1922         new_user) == old_user);
1909     (void) dmu_buf_update_user((dmu_buf_t *)newhdl->sa_bonus,
1910         oldhdl, newhdl, NULL, sa_evict);
1923     oldhdl->sa_bonus = NULL;
1924 }
_____unchanged_portion_omitted_____

```

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

1

```
*****
9996 Thu Apr 25 16:14:58 2013
new/usr/src/uts/common/fs/zfs/sys/dbuf.h
3752 want more verifiable dbuf user eviction
Submitted by: Justin Gibbs <justing@spectralogic.com>
Submitted by: Will Andrews <willa@spectralogic.com>
*****
unchanged_portion_omitted_

139 typedef struct dmu_buf_impl {
140     /*
141      * The following members are immutable, with the exception of
142      * db.db_data, which is protected by db_mtx.
143      */
144
145     /* the publicly visible structure */
146     dmu_buf_t db;
147
148     /* the objset we belong to */
149     struct objset *db_objset;
150
151     /*
152      * handle to safely access the dnode we belong to (NULL when evicted)
153      */
154     struct dnode_handle *db_dnode_handle;
155
156     /*
157      * our parent buffer; if the dnode points to us directly,
158      * db_parent == db_dnode_handle->dnh_dnode->dn_dbuf
159      * only accessed by sync thread ???
160      * (NULL when evicted)
161      * May change from NULL to non-NULL under the protection of db_mtx
162      * (see dbuf_check_blkptr())
163      */
164     struct dmu_buf_impl *db_parent;
165
166     /*
167      * link for hash table of all dmu_buf_impl_t's
168      */
169     struct dmu_buf_impl *db_hash_next;
170
171     /* our block number */
172     uint64_t db_blkid;
173
174     /*
175      * Pointer to the blkptr_t which points to us. May be NULL if we
176      * don't have one yet. (NULL when evicted)
177      */
178     blkptr_t *db_blkptr;
179
180     /*
181      * Our indirection level. Data buffers have db_level==0.
182      * Indirect buffers which point to data buffers have
183      * db_level==1. etc. Buffers which contain dnodes have
184      * db_level==0, since the dnodes are stored in a file.
185      */
186     uint8_t db_level;
187
188     /* db_mtx protects the members below */
189     kmutex_t db_mtx;
190
191     /*
192      * Current state of the buffer
193      */
194     dbuf_states_t db_state;

```

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

2

```
196     /*
197      * Refcount accessed by dmu_buf_{hold,rele}.
198      * If nonzero, the buffer can't be destroyed.
199      * Protected by db_mtx.
200      */
201     refcount_t db_holds;
202
203     /* buffer holding our data */
204     arc_buf_t *db_buf;
205
206     kcondvar_t db_changed;
207     dbuf_dirty_record_t *db_data_pending;
208
209     /* pointer to most recent dirty record for this buffer */
210     dbuf_dirty_record_t *db_last_dirty;
211
212     /*
213      * Our link on the owner dnodes's dn_dbufs list.
214      * Protected by its dn_dbufs_mtx.
215      */
216     list_node_t db_link;
217
218     /* Data which is unique to data (leaf) blocks: */
219
220     /* User callback information. See dmu_buf_set_user(). */
221     dmu_buf_user_t *db_user;
222     /* stuff we store for the user (see dmu_buf_set_user) */
223     void *db_user_ptr;
224     void **db_user_data_ptr_ptr;
225     dmu_buf_evict_func_t *db_evict_func;
226
227     uint8_t db_immediate_evict;
228     uint8_t db_freed_in_flight;
229
230     uint8_t db_dirtycnt;
231 } dmu_buf_impl_t;
unchanged_portion_omitted_

232
233
234
235
236
237
238
239 uint64_t dbuf_whichblock(struct dnode *dn, uint64_t offset);
240
241 dmu_buf_impl_t *dbuf_create_tlib(struct dnode *dn, char *data);
242 void dbuf_create_bonus(struct dnode *dn);
243 int dbuf_spill_set_blksize(dmu_buf_t *db, uint64_t blksize, dmu_tx_t *tx);
244 void dbuf_spill_hold(struct dnode *dn, dmu_buf_impl_t **dbp, void *tag);
245
246 void dbuf_rm_spill(struct dnode *dn, dmu_tx_t *tx);
247
248 dmu_buf_impl_t *dbuf_hold(struct dnode *dn, uint64_t blkid, void *tag);
249 dmu_buf_impl_t *dbuf_hold_level(struct dnode *dn, int level, uint64_t blkid,
250     void *tag);
251 int dbuf_hold_impl(struct dnode *dn, uint8_t level, uint64_t blkid, int create,
252     void *tag, dmu_buf_impl_t **dbp);
253
254 void dbuf_prefetch(struct dnode *dn, uint64_t blkid);
255
256 void dbuf_add_ref(dmu_buf_impl_t *db, void *tag);
257 uint64_t dbuf_refcount(dmu_buf_impl_t *db);
258
259 void dbuf_rele(dmu_buf_impl_t *db, void *tag);
260 void dbuf_rele_and_unlock(dmu_buf_impl_t *db, void *tag);
261
262 dmu_buf_impl_t *dbuf_find(struct dnode *dn, uint8_t level, uint64_t blkid);
263
264 int dbuf_read(dmu_buf_impl_t *db, zio_t *zio, uint32_t flags);
265 void dbuf_will_dirty(dmu_buf_impl_t *db, dmu_tx_t *tx);

```

```
266 void dbuf_fill_done(dmu_buf_impl_t *db, dmu_tx_t *tx);
267 void dmu_buf_will_not_fill(dmu_buf_t *db, dmu_tx_t *tx);
268 void dmu_buf_will_fill(dmu_buf_t *db, dmu_tx_t *tx);
269 void dmu_buf_fill_done(dmu_buf_t *db, dmu_tx_t *tx);
270 void dbuf_assign_arcbuf(dmu_buf_impl_t *db, arc_buf_t *buf, dmu_tx_t *tx);
271 dbuf_dirty_record_t *dbuf_dirty(dmu_buf_impl_t *db, dmu_tx_t *tx);
272 arc_buf_t *dbuf_loan_arcbuf(dmu_buf_impl_t *db);
```

```
274 void dbuf_clear(dmu_buf_impl_t *db, list_t *evict_list_p);
275 void dbuf_evict(dmu_buf_impl_t *db, list_t *evict_list_p);
276 void dbuf_clear(dmu_buf_impl_t *db);
277 void dbuf_evict(dmu_buf_impl_t *db);
```

```
277 void dbuf_setdirty(dmu_buf_impl_t *db, dmu_tx_t *tx);
278 void dbuf_unoverride(dbuf_dirty_record_t *dr);
279 void dbuf_sync_list(list_t *list, dmu_tx_t *tx);
280 void dbuf_release_bp(dmu_buf_impl_t *db);
```

```
282 void dbuf_free_range(struct dnnode *dn, uint64_t start, uint64_t end,
283     struct dmu_tx *);
```

```
285 void dbuf_new_size(dmu_buf_impl_t *db, int size, dmu_tx_t *tx);
```

```
287 #define DB_DNODE(_db)          ((_db)->db_dnode_handle->dnh_dnode)
288 #define DB_DNODE_LOCK(_db)    ((_db)->db_dnode_handle->dnh_zrlock)
289 #define DB_DNODE_ENTER(_db)   (zrl_add(&DB_DNODE_LOCK(_db)))
290 #define DB_DNODE_EXIT(_db)    (zrl_remove(&DB_DNODE_LOCK(_db)))
291 #define DB_DNODE_HELD(_db)    (!zrl_is_zero(&DB_DNODE_LOCK(_db)))
292 #define DB_GET_SPA(_spa_p, _db) {
293     dnnode_t * _dn;
294     DB_DNODE_ENTER(_db);
295     _dn = DB_DNODE(_db);
296     *(_spa_p) = _dn->dn_objset->os_spa;
297     DB_DNODE_EXIT(_db);
298 }
```

unchanged portion omitted

```

*****
30813 Thu Apr 25 16:14:59 2013
new/usr/src/uts/common/fs/zfs/sys/dmu.h
3752 want more verifiable dbuf user eviction
Submitted by: Justin Gibbs <justing@spectrallogic.com>
Submitted by: Will Andrews <willa@spectrallogic.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) 2005, 2010, Oracle and/or its affiliates. All rights reserved.
24 * Copyright (c) 2012 by Delphix. All rights reserved.
25 * Copyright 2011 Nexenta Systems, Inc. All rights reserved.
26 * Copyright (c) 2012, Joyent, Inc. All rights reserved.
27 */

29 /* Portions Copyright 2010 Robert Milkowski */

31 #ifndef _SYS_DMU_H
32 #define _SYS_DMU_H

34 /*
35 * This file describes the interface that the DMU provides for its
36 * consumers.
37 *
38 * The DMU also interacts with the SPA. That interface is described in
39 * dmu_spa.h.
40 */

42 #include <sys/zfs_context.h>
43 #endif /* ! codereview */
44 #include <sys/inttypes.h>
45 #include <sys/types.h>
46 #include <sys/param.h>
47 #include <sys/cred.h>
48 #include <sys/time.h>
49 #include <sys/fs/zfs.h>

51 #ifdef __cplusplus
52 extern "C" {
53 #endif

55 struct uio;
56 struct xuio;
57 struct page;
58 struct vnode;
59 struct spa;

```

```

60 struct zilg;
61 struct zio;
62 struct blkptr;
63 struct zap_cursor;
64 struct dsl_dataset;
65 struct dsl_pool;
66 struct dnode;
67 struct drr_begin;
68 struct drr_end;
69 struct zbookmark;
70 struct spa;
71 struct nvlist;
72 struct arc_buf;
73 struct zio_prop;
74 struct sa_handle;

76 typedef struct objset objset_t;
77 typedef struct dmu_tx dmu_tx_t;
78 typedef struct dsl_dir dsl_dir_t;

80 typedef enum dmu_object_byteswap {
81     DMU_BSWAP_UINT8,
82     DMU_BSWAP_UINT16,
83     DMU_BSWAP_UINT32,
84     DMU_BSWAP_UINT64,
85     DMU_BSWAP_ZAP,
86     DMU_BSWAP_DNODE,
87     DMU_BSWAP_OBJSET,
88     DMU_BSWAP_ZNODE,
89     DMU_BSWAP_OLDACL,
90     DMU_BSWAP_ACL,
91     /*
92      * Allocating a new byteswap type number makes the on-disk format
93      * incompatible with any other format that uses the same number.
94      *
95      * Data can usually be structured to work with one of the
96      * DMU_BSWAP_UINT* or DMU_BSWAP_ZAP types.
97      */
98     DMU_BSWAP_NUMFUNCS
99 } dmu_object_byteswap_t;

101 #define DMU_OT_NEWTYPE 0x80
102 #define DMU_OT_METADATA 0x40
103 #define DMU_OT_BYTESWAP_MASK 0x3f

105 /*
106 * Defines a uint8_t object type. Object types specify if the data
107 * in the object is metadata (boolean) and how to byteswap the data
108 * (dmu_object_byteswap_t).
109 */
110 #define DMU_OT(byteswap, metadata) \
111     (DMU_OT_NEWTYPE | \
112      ((metadata) ? DMU_OT_METADATA : 0) | \
113      ((byteswap) & DMU_OT_BYTESWAP_MASK))

115 #define DMU_OT_IS_VALID(ot) (((ot) & DMU_OT_NEWTYPE) ? \
116     ((ot) & DMU_OT_BYTESWAP_MASK) < DMU_BSWAP_NUMFUNCS : \
117     (ot) < DMU_OT_NUMTYPES)

119 #define DMU_OT_IS_METADATA(ot) (((ot) & DMU_OT_NEWTYPE) ? \
120     ((ot) & DMU_OT_METADATA) : \
121     dmu_ot[(ot)].ot_metadata)

123 #define DMU_OT_BYTESWAP(ot) (((ot) & DMU_OT_NEWTYPE) ? \
124     ((ot) & DMU_OT_BYTESWAP_MASK) : \
125     dmu_ot[(ot)].ot_byteswap)

```

```

127 typedef enum dmu_object_type {
128     DMU_OT_NONE,
129     /* general: */
130     DMU_OT_OBJECT_DIRECTORY,      /* ZAP */
131     DMU_OT_OBJECT_ARRAY,          /* UINT64 */
132     DMU_OT_PACKED_NVLIST,         /* UINT8 (XDR by nvlist_pack/unpack) */
133     DMU_OT_PACKED_NVLIST_SIZE,    /* UINT64 */
134     DMU_OT_BPOBJ,                 /* UINT64 */
135     DMU_OT_BPOBJ_HDR,             /* UINT64 */
136     /* spa: */
137     DMU_OT_SPACE_MAP_HEADER,      /* UINT64 */
138     DMU_OT_SPACE_MAP,             /* UINT64 */
139     /* zil: */
140     DMU_OT_INTENT_LOG,            /* UINT64 */
141     /* dmu: */
142     DMU_OT_DNODE,                 /* DNODE */
143     DMU_OT_OBJSET,                /* OBJSET */
144     /* dsl: */
145     DMU_OT_DSL_DIR,               /* UINT64 */
146     DMU_OT_DSL_DIR_CHILD_MAP,     /* ZAP */
147     DMU_OT_DSL_DS_SNAP_MAP,       /* ZAP */
148     DMU_OT_DSL_PROPS,             /* ZAP */
149     DMU_OT_DSL_DATASET,           /* UINT64 */
150     /* zpl: */
151     DMU_OT_ZNODE,                 /* ZNODE */
152     DMU_OT_OLDACL,                /* Old ACL */
153     DMU_OT_PLAIN_FILE_CONTENTS,    /* UINT8 */
154     DMU_OT_DIRECTORY_CONTENTS,    /* ZAP */
155     DMU_OT_MASTER_NODE,           /* ZAP */
156     DMU_OT_UNLINKED_SET,          /* ZAP */
157     /* zvol: */
158     DMU_OT_ZVOL,                  /* UINT8 */
159     DMU_OT_ZVOL_PROP,             /* ZAP */
160     /* other; for testing only! */
161     DMU_OT_PLAIN_OTHER,            /* UINT8 */
162     DMU_OT_UINT64_OTHER,           /* UINT64 */
163     DMU_OT_ZAP_OTHER,             /* ZAP */
164     /* new object types: */
165     DMU_OT_ERROR_LOG,             /* ZAP */
166     DMU_OT_SPA_HISTORY,           /* UINT8 */
167     DMU_OT_SPA_HISTORY_OFFSETS,    /* spa_his_phys_t */
168     DMU_OT_POOL_PROPS,            /* ZAP */
169     DMU_OT_DSL_PERMS,             /* ZAP */
170     DMU_OT_ACL,                   /* ACL */
171     DMU_OT_SYSACL,                /* SYSACL */
172     DMU_OT_FUID,                  /* FUID table (Packed NVLIST UINT8) */
173     DMU_OT_FUID_SIZE,             /* FUID table size UINT64 */
174     DMU_OT_NEXT_CLONES,           /* ZAP */
175     DMU_OT_SCAN_QUEUE,            /* ZAP */
176     DMU_OT_USERGROUP_USED,        /* ZAP */
177     DMU_OT_USERGROUP_QUOTA,       /* ZAP */
178     DMU_OT_USERREFS,              /* ZAP */
179     DMU_OT_DDT_ZAP,               /* ZAP */
180     DMU_OT_DDT_STATS,             /* ZAP */
181     DMU_OT_SA,                    /* System attr */
182     DMU_OT_SA_MASTER_NODE,        /* ZAP */
183     DMU_OT_SA_ATTR_REGISTRATION,  /* ZAP */
184     DMU_OT_SA_ATTR_LAYOUTS,       /* ZAP */
185     DMU_OT_SCAN_XLATE,            /* ZAP */
186     DMU_OT_DEDUP,                 /* fake dedup BP from ddt_bp_create() */
187     DMU_OT_DEADLIST,              /* ZAP */
188     DMU_OT_DEADLIST_HDR,          /* UINT64 */
189     DMU_OT_DSL_CLONES,            /* ZAP */
190     DMU_OT_BPOBJ_SUBOBJ,          /* UINT64 */
191     /*

```

```

192     * Do not allocate new object types here. Doing so makes the on-disk
193     * format incompatible with any other format that uses the same object
194     * type number.
195     *
196     * When creating an object which does not have one of the above types
197     * use the DMU_OTN_* type with the correct byteswap and metadata
198     * values.
199     *
200     * The DMU_OTN_* types do not have entries in the dmu_ot table,
201     * use the DMU_OT_IS_METADATA() and DMU_OT_BYTESWAP() macros instead
202     * of indexing into dmu_ot directly (this works for both DMU_OT_* types
203     * and DMU_OTN_* types).
204     */
205     DMU_OT_NUMTYPES,
206
207     /*
208     * Names for valid types declared with DMU_OT().
209     */
210     DMU_OTN_UINT8_DATA = DMU_OT(DMU_BSWAP_UINT8, B_FALSE),
211     DMU_OTN_UINT8_METADATA = DMU_OT(DMU_BSWAP_UINT8, B_TRUE),
212     DMU_OTN_UINT16_DATA = DMU_OT(DMU_BSWAP_UINT16, B_FALSE),
213     DMU_OTN_UINT16_METADATA = DMU_OT(DMU_BSWAP_UINT16, B_TRUE),
214     DMU_OTN_UINT32_DATA = DMU_OT(DMU_BSWAP_UINT32, B_FALSE),
215     DMU_OTN_UINT32_METADATA = DMU_OT(DMU_BSWAP_UINT32, B_TRUE),
216     DMU_OTN_UINT64_DATA = DMU_OT(DMU_BSWAP_UINT64, B_FALSE),
217     DMU_OTN_UINT64_METADATA = DMU_OT(DMU_BSWAP_UINT64, B_TRUE),
218     DMU_OTN_ZAP_DATA = DMU_OT(DMU_BSWAP_ZAP, B_FALSE),
219     DMU_OTN_ZAP_METADATA = DMU_OT(DMU_BSWAP_ZAP, B_TRUE),
220 } dmu_object_type_t;
221
222 typedef enum txg_how {
223     TXG_WAIT = 1,
224     TXG_NOWAIT,
225 } txg_how_t;
226
227 void byteswap_uint64_array(void *buf, size_t size);
228 void byteswap_uint32_array(void *buf, size_t size);
229 void byteswap_uint16_array(void *buf, size_t size);
230 void byteswap_uint8_array(void *buf, size_t size);
231 void zap_byteswap(void *buf, size_t size);
232 void zfs_oldacl_byteswap(void *buf, size_t size);
233 void zfs_acl_byteswap(void *buf, size_t size);
234 void zfs_znode_byteswap(void *buf, size_t size);
235
236 #define DS_FIND_SNAPSHOTS      (1<<0)
237 #define DS_FIND_CHILDREN      (1<<1)
238
239 /*
240  * The maximum number of bytes that can be accessed as part of one
241  * operation, including metadata.
242  */
243 #define DMU_MAX_ACCESS (10<<20) /* 10MB */
244 #define DMU_MAX_DELETEBLKCNT (20480) /* ~5MB of indirect blocks */
245
246 #define DMU_USERUSED_OBJECT    (-1ULL)
247 #define DMU_GROUPUSED_OBJECT  (-2ULL)
248 #define DMU_DEADLIST_OBJECT   (-3ULL)
249
250 /*
251  * artificial blkids for bonus buffer and spill blocks
252  */
253 #define DMU_BONUS_BLKID        (-1ULL)
254 #define DMU_SPILL_BLKID        (-2ULL)
255 /*
256  * Public routines to create, destroy, open, and close objsets.
257  */

```

```

258 int dmu_objset_hold(const char *name, void *tag, objset_t **osp);
259 int dmu_objset_own(const char *name, dmu_objset_type_t type,
260     boolean_t readonly, void *tag, objset_t **osp);
261 void dmu_objset_rele(objset_t *os, void *tag);
262 void dmu_objset_disown(objset_t *os, void *tag);
263 int dmu_objset_open_ds(struct dsl_dataset *ds, objset_t **osp);

265 void dmu_objset_evict_dbufs(objset_t *os);
266 int dmu_objset_create(const char *name, dmu_objset_type_t type, uint64_t flags,
267     void (*func)(objset_t *os, void *arg, cred_t *cr, dmu_tx_t *tx), void *arg);
268 int dmu_objset_clone(const char *name, const char *origin);
269 int dsl_destroy_snapshots_nvlist(struct nvlist *snaps, boolean_t defer,
270     struct nvlist *errlist);
271 int dmu_objset_snapshot_one(const char *fsname, const char *snapname);
272 int dmu_objset_snapshot_tmp(const char *, const char *, int);
273 int dmu_objset_find(char *name, int func(const char *, void *), void *arg,
274     int flags);
275 void dmu_objset_byteswap(void *buf, size_t size);
276 int dsl_dataset_rename_snapshot(const char *fsname,
277     const char *oldsnapname, const char *newsnapname, boolean_t recursive);

279 typedef struct dmu_buf {
280     uint64_t db_object;           /* object that this buffer is part of */
281     uint64_t db_offset;         /* byte offset in this object */
282     uint64_t db_size;           /* size of buffer in bytes */
283     void *db_data;              /* data in buffer */
284 } dmu_buf_t;

42 typedef void dmu_buf_evict_func_t(struct dmu_buf *db, void *user_ptr);

286 /*
287  * The names of zap entries in the DIRECTORY_OBJECT of the MOS.
288  */
289 #define DMU_POOL_DIRECTORY_OBJECT 1
290 #define DMU_POOL_CONFIG "config"
291 #define DMU_POOL_FEATURES_FOR_WRITE "features_for_write"
292 #define DMU_POOL_FEATURES_FOR_READ "features_for_read"
293 #define DMU_POOL_FEATURE_DESCRIPTIONS "feature_descriptions"
294 #define DMU_POOL_ROOT_DATASET "root_dataset"
295 #define DMU_POOL_SYNC_BPOBJ "sync_bplist"
296 #define DMU_POOL_ERRLOG_SCRUB "errlog_scrub"
297 #define DMU_POOL_ERRLOG_LAST "errlog_last"
298 #define DMU_POOL_SPARES "spares"
299 #define DMU_POOL_DEFLATE "deflate"
300 #define DMU_POOL_HISTORY "history"
301 #define DMU_POOL_PROPS "pool_props"
302 #define DMU_POOL_L2CACHE "l2cache"
303 #define DMU_POOL_TMP_USERREFS "tmp_userrefs"
304 #define DMU_POOL_DDT "DDT-%s-%s-%s"
305 #define DMU_POOL_DDT_STATS "DDT-statistics"
306 #define DMU_POOL_CREATION_VERSION "creation_version"
307 #define DMU_POOL_SCAN "scan"
308 #define DMU_POOL_FREE_BPOBJ "free_bpobj"
309 #define DMU_POOL_BPTEE_OBJ "bptree_obj"
310 #define DMU_POOL_EMPTY_BPOBJ "empty_bpobj"

312 /*
313  * Allocate an object from this objset. The range of object numbers
314  * available is (0, DN_MAX_OBJECT). Object 0 is the meta-dnode.
315  *
316  * The transaction must be assigned to a txg. The newly allocated
317  * object will be "held" in the transaction (ie. you can modify the
318  * newly allocated object in this transaction).
319  *
320  * dmu_object_alloc() chooses an object and returns it in *objectp.
321  */

```

```

322 * dmu_object_claim() allocates a specific object number. If that
323 * number is already allocated, it fails and returns EEXIST.
324 *
325 * Return 0 on success, or ENOSPC or EEXIST as specified above.
326 */
327 uint64_t dmu_object_alloc(objset_t *os, dmu_object_type_t ot,
328     int blocksize, dmu_object_type_t bonus_type, int bonus_len, dmu_tx_t *tx);
329 int dmu_object_claim(objset_t *os, uint64_t object, dmu_object_type_t ot,
330     int blocksize, dmu_object_type_t bonus_type, int bonus_len, dmu_tx_t *tx);
331 int dmu_object_reclaim(objset_t *os, uint64_t object, dmu_object_type_t ot,
332     int blocksize, dmu_object_type_t bonustype, int bonuslen);

334 /*
335  * Free an object from this objset.
336  *
337  * The object's data will be freed as well (ie. you don't need to call
338  * dmu_free(object, 0, -1, tx)).
339  *
340  * The object need not be held in the transaction.
341  *
342  * If there are any holds on this object's buffers (via dmu_buf_hold()),
343  * or tx holds on the object (via dmu_tx_hold_object()), you can not
344  * free it; it fails and returns EBUSY.
345  *
346  * If the object is not allocated, it fails and returns ENOENT.
347  *
348  * Return 0 on success, or EBUSY or ENOENT as specified above.
349  */
350 int dmu_object_free(objset_t *os, uint64_t object, dmu_tx_t *tx);

352 /*
353  * Find the next allocated or free object.
354  *
355  * The objectp parameter is in-out. It will be updated to be the next
356  * object which is allocated. Ignore objects which have not been
357  * modified since txg.
358  *
359  * XXX Can only be called on a objset with no dirty data.
360  *
361  * Returns 0 on success, or ENOENT if there are no more objects.
362  */
363 int dmu_object_next(objset_t *os, uint64_t *objectp,
364     boolean_t hole, uint64_t txg);

366 /*
367  * Set the data blocksize for an object.
368  *
369  * The object cannot have any blocks allocated beyond the first. If
370  * the first block is allocated already, the new size must be greater
371  * than the current block size. If these conditions are not met,
372  * ENOTSUP will be returned.
373  *
374  * Returns 0 on success, or EBUSY if there are any holds on the object
375  * contents, or ENOTSUP as described above.
376  */
377 int dmu_object_set_blocksize(objset_t *os, uint64_t object, uint64_t size,
378     int ibs, dmu_tx_t *tx);

380 /*
381  * Set the checksum property on a dnode. The new checksum algorithm will
382  * apply to all newly written blocks; existing blocks will not be affected.
383  */
384 void dmu_object_set_checksum(objset_t *os, uint64_t object, uint8_t checksum,
385     dmu_tx_t *tx);

387 /*

```

```

388 * Set the compress property on a dnode. The new compression algorithm will
389 * apply to all newly written blocks; existing blocks will not be affected.
390 */
391 void dmu_object_set_compress(objset_t *os, uint64_t object, uint8_t compress,
392     dmu_tx_t *tx);

394 /*
395 * Decide how to write a block: checksum, compression, number of copies, etc.
396 */
397 #define WP_NOFILL      0x1
398 #define WP_DMU_SYNC    0x2
399 #define WP_SPILL       0x4

401 void dmu_write_policy(objset_t *os, struct dnode *dn, int level, int wp,
402     struct zio_prop *zp);
403 /*
404 * The bonus data is accessed more or less like a regular buffer.
405 * You must dmu_bonus_hold() to get the buffer, which will give you a
406 * dmu_buf_t with db_offset==1ULL, and db_size = the size of the bonus
407 * data. As with any normal buffer, you must call dmu_buf_read() to
408 * read db_data, dmu_buf_will_dirty() before modifying it, and the
409 * object must be held in an assigned transaction before calling
410 * dmu_buf_will_dirty. You may use dmu_buf_set_user() on the bonus
411 * buffer as well. You must release your hold with dmu_buf_rele().
412 */
413 int dmu_bonus_hold(objset_t *os, uint64_t object, void *tag, dmu_buf_t **);
414 int dmu_bonus_max(void);
415 int dmu_set_bonus(dmu_buf_t *, int, dmu_tx_t *);
416 int dmu_set_bonustype(dmu_buf_t *, dmu_object_type_t, dmu_tx_t *);
417 dmu_object_type_t dmu_get_bonustype(dmu_buf_t *);
418 int dmu_rm_spill(objset_t *, uint64_t, dmu_tx_t *);

420 /*
421 * Special spill buffer support used by "SA" framework
422 */

424 int dmu_spill_hold_by_bonus(dmu_buf_t *bonus, void *tag, dmu_buf_t **dbp);
425 int dmu_spill_hold_by_dnode(struct dnode *dn, uint32_t flags,
426     void *tag, dmu_buf_t **dbp);
427 int dmu_spill_hold_existing(dmu_buf_t *bonus, void *tag, dmu_buf_t **dbp);

429 /*
430 * Obtain the DMU buffer from the specified object which contains the
431 * specified offset. dmu_buf_hold() puts a "hold" on the buffer, so
432 * that it will remain in memory. You must release the hold with
433 * dmu_buf_rele(). You musn't access the dmu_buf_t after releasing your
434 * hold. You must have a hold on any dmu_buf_t* you pass to the DMU.
435 *
436 * You must call dmu_buf_read, dmu_buf_will_dirty, or dmu_buf_will_fill
437 * on the returned buffer before reading or writing the buffer's
438 * db_data. The comments for those routines describe what particular
439 * operations are valid after calling them.
440 *
441 * The object number must be a valid, allocated object number.
442 */
443 int dmu_buf_hold(objset_t *os, uint64_t object, uint64_t offset,
444     void *tag, dmu_buf_t **, int flags);
445 void dmu_buf_add_ref(dmu_buf_t *db, void *tag);
446 void dmu_buf_rele(dmu_buf_t *db, void *tag);
447 uint64_t dmu_buf_refcount(dmu_buf_t *db);

449 /*
450 * dmu_buf_hold_array holds the DMU buffers which contain all bytes in a
451 * range of an object. A pointer to an array of dmu_buf_t*'s is
452 * returned (in *dbpp).
453 */

```

```

454 * dmu_buf_rele_array releases the hold on an array of dmu_buf_t*'s, and
455 * frees the array. The hold on the array of buffers MUST be released
456 * with dmu_buf_rele_array. You can NOT release the hold on each buffer
457 * individually with dmu_buf_rele.
458 */
459 int dmu_buf_hold_array_by_bonus(dmu_buf_t *db, uint64_t offset,
460     uint64_t length, int read, void *tag, int *numbufsp, dmu_buf_t ***dbpp);
461 void dmu_buf_rele_array(dmu_buf_t **, int numbufs, void *tag);

463 struct dmu_buf_user;

465 typedef void dmu_buf_evict_func_t(struct dmu_buf_user *);

467 #endif /* ! codereview */
468 /*
469 * The DMU buffer user object is used to allow private data to be
470 * associated with a dbuf for the duration of its lifetime. This private
471 * data must include a dmu_buf_user_t as its first object, which is passed
472 * into the DMU user data API and can be attached to a dbuf. Clients can
473 * regain access to their private data structure with a cast.
474 *
475 * DMU buffer users can be notified via a callback when their associated
476 * dbuf has been evicted. This is typically used to free the user's
477 * private data. The eviction callback is executed without the dbuf
478 * mutex held or any other type of mechanism to guarantee that the
479 * dbuf is still available. For this reason, users must assume the dbuf
480 * has already been freed and not reference the dbuf from the callback
481 * context.
482 *
483 * Users requestion "immediate eviction" are notified as soon as the dbuf
484 * is only referenced by dirty records (dirties == holds). Otherwise the
485 * eviction callback occurs after the last reference to the dbuf is dropped.
486 *
487 * Eviction Callback Processing
488 * =====
489 * In any context where a dbuf reference drop may trigger an eviction, an
490 * eviction queue object must be provided. This queue must then be
491 * processed while not holding any dbuf locks. In this way, the user can
492 * perform any work needed in their eviction function without fear of
493 * lock order reversals.
494 *
495 * Implementation Note
496 * =====
497 * Some users will occasionally want to map a structure directly onto the
498 * backing dbuf. Using an union with an name alias macro to access these
499 * overlays reduces the ugliness of code that accesses them. Initial work on
500 * user objects involved using a macro that took the user object as an
501 * argument to access the fields, which resulted in hundreds of lines of
502 * needless diffs and wasn't any easier to read.
503 * Returns NULL on success, or the existing user ptr if it's already
504 * been set.
505 *
506 * user_ptr is for use by the user and can be obtained via dmu_buf_get_user().
507 *
508 * user_data_ptr_ptr should be NULL, or a pointer to a pointer which
509 * will be set to db->db_data when you are allowed to access it. Note
510 * that db->db_data (the pointer) can change when you do dmu_buf_read(),
511 * dmu_buf_tryupgrade(), dmu_buf_will_dirty(), or dmu_buf_will_fill().
512 * *user_data_ptr_ptr will be set to the new value when it changes.
513 *
514 * If non-NULL, pageout func will be called when this buffer is being
515 * excised from the cache, so that you can clean up the data structure
516 * pointed to by user_ptr.
517 *
518 * dmu_evict_user() will call the pageout func for all buffers in a
519 * objset with a given pageout func.

```

```

503 */
504 typedef struct dmu_buf_user {
505     /*
506      * This instance's link in the eviction queue. Set when the buffer
507      * has evicted and the callback needs to be called.
508      */
509     void *dmu_buf_set_user(dmu_buf_t *db, void *user_ptr, void *user_data_ptr_ptr,
510         dmu_buf_evict_func_t *pageout_func);
511     /*
512      * set_user_ie is the same as set_user, but request immediate eviction
513      * when hold count goes to zero.
514      */
515     list_node_t evict_queue_link;
516     /** This instance's eviction function pointer. */
517     dmu_buf_evict_func_t *evict_func;
518 } dmu_buf_user_t;
519
520 /*
521 * Initialize the given dmu_buf_user_t instance with the eviction function
522 * evict_func, to be called when the user is evicted.
523 *
524 * NOTE: This function should only be called once on a given object. To
525 * help enforce this, dbu should already be zeroed on entry.
526 */
527 static inline void
528 dmu_buf_init_user(dmu_buf_user_t *dbu, dmu_buf_evict_func_t *evict_func)
529 {
530     ASSERT(dbu->evict_func == NULL);
531     ASSERT(!list_link_active(&dbu->evict_queue_link));
532     dbu->evict_func = evict_func;
533 }
534
535 static inline void
536 dmu_buf_create_user_evict_list(list_t *evict_list_p)
537 {
538     list_create(evict_list_p, sizeof(dmu_buf_user_t),
539         offsetof(dmu_buf_user_t, evict_queue_link));
540 }
541
542 static inline void
543 dmu_buf_process_user_evicts(list_t *evict_list_p)
544 {
545     dmu_buf_user_t *dbu, *next;
546
547     for (dbu = (dmu_buf_user_t *)list_head(evict_list_p); dbu != NULL;
548         dbu = next) {
549         next = (dmu_buf_user_t *)list_next(evict_list_p, dbu);
550         list_remove(evict_list_p, dbu);
551         dbu->evict_func(dbu);
552     }
553 }
554
555 void *dmu_buf_set_user_ie(dmu_buf_t *db, void *user_ptr,
556     void *user_data_ptr_ptr, dmu_buf_evict_func_t *pageout_func);
557 void *dmu_buf_update_user(dmu_buf_t *db_fake, void *old_user_ptr,
558     void *user_ptr, void *user_data_ptr_ptr,
559     dmu_buf_evict_func_t *pageout_func);
560 void dmu_evict_user(objset_t *os, dmu_buf_evict_func_t *func);
561
562 static inline void
563 dmu_buf_destroy_user_evict_list(list_t *evict_list_p)
564 {
565     list_destroy(evict_list_p);
566 }
567
568 dmu_buf_user_t *dmu_buf_set_user(dmu_buf_t *db, dmu_buf_user_t *user);
569 dmu_buf_user_t *dmu_buf_set_user_ie(dmu_buf_t *db, dmu_buf_user_t *user);

```

```

558 dmu_buf_user_t *dmu_buf_replace_user(dmu_buf_t *db,
559     dmu_buf_user_t *old_user, dmu_buf_user_t *new_user);
560 dmu_buf_user_t *dmu_buf_remove_user(dmu_buf_t *db, dmu_buf_user_t *user);
561 dmu_buf_user_t *dmu_buf_get_user(dmu_buf_t *db);
562 /*
563 * Returns the user_ptr set with dmu_buf_set_user(), or NULL if not set.
564 */
565 void *dmu_buf_get_user(dmu_buf_t *db);
566
567 /*
568 * Returns the blkptr associated with this dbuf, or NULL if not set.
569 */
570 struct blkptr *dmu_buf_get_blkptr(dmu_buf_t *db);
571
572 /*
573 * Indicate that you are going to modify the buffer's data (db_data).
574 *
575 * The transaction (tx) must be assigned to a txg (ie. you've called
576 * dmu_tx_assign()). The buffer's object must be held in the tx
577 * (ie. you've called dmu_tx_hold_object(tx, db->db_object)).
578 */
579 void dmu_buf_will_dirty(dmu_buf_t *db, dmu_tx_t *tx);
580
581 /*
582 * Tells if the given dbuf is freeable.
583 */
584 boolean_t dmu_buf_freeable(dmu_buf_t *);
585
586 /*
587 * You must create a transaction, then hold the objects which you will
588 * (or might) modify as part of this transaction. Then you must assign
589 * the transaction to a transaction group. Once the transaction has
590 * been assigned, you can modify buffers which belong to held objects as
591 * part of this transaction. You can't modify buffers before the
592 * transaction has been assigned; you can't modify buffers which don't
593 * belong to objects which this transaction holds; you can't hold
594 * objects once the transaction has been assigned. You may hold an
595 * object which you are going to free (with dmu_object_free()), but you
596 * don't have to.
597 *
598 * You can abort the transaction before it has been assigned.
599 *
600 * Note that you may hold buffers (with dmu_buf_hold) at any time,
601 * regardless of transaction state.
602 */
603 #define DMU_NEW_OBJECT (-1ULL)
604 #define DMU_OBJECT_END (-1ULL)
605
606 dmu_tx_t *dmu_tx_create(objset_t *os);
607 void dmu_tx_hold_write(dmu_tx_t *tx, uint64_t object, uint64_t off, int len);
608 void dmu_tx_hold_free(dmu_tx_t *tx, uint64_t object, uint64_t off,
609     uint64_t len);
610 void dmu_tx_hold_zap(dmu_tx_t *tx, uint64_t object, int add, const char *name);
611 void dmu_tx_hold_bonus(dmu_tx_t *tx, uint64_t object);
612 void dmu_tx_hold_spill(dmu_tx_t *tx, uint64_t object);
613 void dmu_tx_hold_sa(dmu_tx_t *tx, struct sa_handle *hdl, boolean_t may_grow);
614 void dmu_tx_hold_sa_create(dmu_tx_t *tx, int total_size);
615 void dmu_tx_abort(dmu_tx_t *tx);
616 int dmu_tx_assign(dmu_tx_t *tx, enum txg_how txg_how);
617 void dmu_tx_wait(dmu_tx_t *tx);
618 void dmu_tx_commit(dmu_tx_t *tx);
619
620 /*
621 * To register a commit callback, dmu_tx_callback_register() must be called.
622 */

```

```

620 * dcb_data is a pointer to caller private data that is passed on as a
621 * callback parameter. The caller is responsible for properly allocating and
622 * freeing it.
623 *
624 * When registering a callback, the transaction must be already created, but
625 * it cannot be committed or aborted. It can be assigned to a txg or not.
626 *
627 * The callback will be called after the transaction has been safely written
628 * to stable storage and will also be called if the dmu_tx is aborted.
629 * If there is any error which prevents the transaction from being committed to
630 * disk, the callback will be called with a value of error != 0.
631 */
632 typedef void dmu_tx_callback_func_t(void *dcb_data, int error);

634 void dmu_tx_callback_register(dmu_tx_t *tx, dmu_tx_callback_func_t *dcb_func,
635 void *dcb_data);

637 /*
638 * Free up the data blocks for a defined range of a file. If size is
639 * -1, the range from offset to end-of-file is freed.
640 */
641 int dmu_free_range(objset_t *os, uint64_t object, uint64_t offset,
642 uint64_t size, dmu_tx_t *tx);
643 int dmu_free_long_range(objset_t *os, uint64_t object, uint64_t offset,
644 uint64_t size);
645 int dmu_free_object(objset_t *os, uint64_t object);

647 /*
648 * Convenience functions.
649 *
650 * Canfail routines will return 0 on success, or an errno if there is a
651 * nonrecoverable I/O error.
652 */
653 #define DMU_READ_PREFETCH 0 /* prefetch */
654 #define DMU_READ_NO_PREFETCH 1 /* don't prefetch */
655 int dmu_read(objset_t *os, uint64_t object, uint64_t offset, uint64_t size,
656 void *buf, uint32_t flags);
657 void dmu_write(objset_t *os, uint64_t object, uint64_t offset, uint64_t size,
658 const void *buf, dmu_tx_t *tx);
659 void dmu_prealloc(objset_t *os, uint64_t object, uint64_t offset, uint64_t size,
660 dmu_tx_t *tx);
661 int dmu_read_uio(objset_t *os, uint64_t object, struct uio *uio, uint64_t size);
662 int dmu_write_uio(objset_t *os, uint64_t object, struct uio *uio, uint64_t size,
663 dmu_tx_t *tx);
664 int dmu_write_uio_dbuf(dmu_buf_t *zdb, struct uio *uio, uint64_t size,
665 dmu_tx_t *tx);
666 int dmu_write_pages(objset_t *os, uint64_t object, uint64_t offset,
667 uint64_t size, struct page *pp, dmu_tx_t *tx);
668 struct arc_buf *dmu_request_arcbuf(dmu_buf_t *handle, int size);
669 void dmu_return_arcbuf(struct arc_buf *buf);
670 void dmu_assign_arcbuf(dmu_buf_t *handle, uint64_t offset, struct arc_buf *buf,
671 dmu_tx_t *tx);
672 int dmu_xuio_init(struct xuio *uio, int niow);
673 void dmu_xuio_fini(struct xuio *uio);
674 int dmu_xuio_add(struct xuio *uio, struct arc_buf *abuf, offset_t off,
675 size_t n);
676 int dmu_xuio_cnt(struct xuio *uio);
677 struct arc_buf *dmu_xuio_arcbuf(struct xuio *uio, int i);
678 void dmu_xuio_clear(struct xuio *uio, int i);
679 void xuio_stat_wbuf_copied();
680 void xuio_stat_wbuf_nocopy();

682 extern int zfs_prefetch_disable;

684 /*
685 * Asynchronously try to read in the data.

```

```

686 */
687 void dmu_prefetch(objset_t *os, uint64_t object, uint64_t offset,
688 uint64_t len);

690 typedef struct dmu_object_info {
691 /* All sizes are in bytes unless otherwise indicated. */
692 uint32_t doi_data_block_size;
693 uint32_t doi_metadata_block_size;
694 dmu_object_type_t doi_type;
695 dmu_object_type_t doi_bonus_type;
696 uint64_t doi_bonus_size;
697 uint8_t doi_indirection; /* 2 = dnode->indirect->data */
698 uint8_t doi_checksum;
699 uint8_t doi_compress;
700 uint8_t doi_pad[5];
701 uint64_t doi_physical_blocks_512; /* data + metadata, 512b blks */
702 uint64_t doi_max_offset;
703 uint64_t doi_fill_count; /* number of non-empty blocks */
704 } dmu_object_info_t;

```

unchanged portion omitted

```

*****
8239 Thu Apr 25 16:14:59 2013
new/usr/src/uts/common/fs/zfs/sys/dmu_impl.h
3752 want more verifiable dbuf user eviction
Submitted by: Justin Gibbs <justing@spectrallogic.com>
Submitted by: Will Andrews <willa@spectrallogic.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 2010 Sun Microsystems, Inc. All rights reserved.
23 * Use is subject to license terms.
24 * Copyright (c) 2012, Joyent, Inc. All rights reserved.
25 */

27 #ifndef _SYS_DMU_IMPL_H
28 #define _SYS_DMU_IMPL_H

30 #include <sys/txg_impl.h>
31 #include <sys/zio.h>
32 #include <sys/dnode.h>
33 #include <sys/zfs_context.h>
34 #include <sys/zfs_ioctl.h>
35 #include <sys/dmu.h>
36 #endif /* !codereview */

38 #ifdef __cplusplus
39 extern "C" {
40 #endif

42 /*
43  * This is the locking strategy for the DMU. Numbers in parenthesis are
44  * cases that use that lock order, referenced below:
45  *
46  * ARC is self-contained
47  * bplist is self-contained
48  * refcount is self-contained
49  * txg is self-contained (hopefully!)
50  * zst_lock
51  * zf_rwlock
52  *
53  * XXX try to improve evicting path?
54  *
55  * dp_config_rwlock > os_obj_lock > dn_struct_rwlock >
56  *   dn_dbufs_mtx > hash_mutexes > db_mtx > dd_lock > leafs
57  *
58  * dp_config_rwlock
59  *   must be held before: everything

```

```

60 *   protects dd namespace changes
61 *   protects property changes globally
62 * held from:
63 *   dsl_dir_open/r:
64 *   dsl_dir_create_sync/w:
65 *   dsl_dir_sync_destroy/w:
66 *   dsl_dir_rename_sync/w:
67 *   dsl_prop_changed_notify/r:
68 *
69 * os_obj_lock
70 *   must be held before:
71 *   everything except dp_config_rwlock
72 *   protects os_obj_next
73 * held from:
74 *   dmu_object_alloc: dn_dbufs_mtx, db_mtx, hash_mutexes, dn_struct_rwlock
75 *
76 * dn_struct_rwlock
77 *   must be held before:
78 *   everything except dp_config_rwlock and os_obj_lock
79 *   protects structure of dnode (eg. nlevels)
80 *   db_blkptr can change when syncing out change to nlevels
81 *   dn_maxblkid
82 *   dn_nlevels
83 *   dn_*blkz*
84 *   phys_nlevels, maxblkid, physical blkptr_t's (?)
85 * held from:
86 *   callers of dbuf_read_impl, dbuf_hold[_impl], dbuf_prefetch
87 *   dmu_object_info_from_dnode: dn_dirty_mtx (dn_datablkz)
88 *   dmu_tx_count_free:
89 *   dbuf_read_impl: db_mtx, dmu_zfetch()
90 *   dmu_zfetch: zf_rwlock/r, zst_lock, dbuf_prefetch()
91 *   dbuf_new_size: db_mtx
92 *   dbuf_dirty: db_mtx
93 *   dbuf_findbp: (callers, phys? - the real need)
94 *   dbuf_create: dn_dbufs_mtx, hash_mutexes, db_mtx (phys?)
95 *   dbuf_prefetch: dn_dirty_mtx, hash_mutexes, db_mtx, dn_dbufs_mtx
96 *   dbuf_hold_impl: hash_mutexes, db_mtx, dn_dbufs_mtx, dbuf_findbp()
97 *   dnode_sync/w (increase_indirection): db_mtx (phys)
98 *   dnode_set_blkz/w: dn_dbufs_mtx (dn_*blkz*)
99 *   dnode_new_blkid/w: (dn_maxblkid)
100 *   dnode_free_range/w: dn_dirty_mtx (dn_maxblkid)
101 *   dnode_next_offset: (phys)
102 *
103 * dn_dbufs_mtx
104 *   must be held before:
105 *   db_mtx, hash_mutexes
106 *   protects:
107 *   dn_dbufs
108 *   dn_evicted
109 * held from:
110 *   dmu_evict_user: db_mtx (dn_dbufs)
111 *   dbuf_free_range: db_mtx (dn_dbufs)
112 *   dbuf_remove_ref: db_mtx, callees:
113 *     dbuf_hash_remove: hash_mutexes, db_mtx
114 *   dbuf_create: hash_mutexes, db_mtx (dn_dbufs)
115 *   dnode_set_blkz: (dn_dbufs)
116 *
117 * hash_mutexes (global)
118 *   must be held before:
119 *   db_mtx
120 *   protects dbuf_hash_table (global) and db_hash_next
121 * held from:
122 *   dbuf_find: db_mtx
123 *   dbuf_hash_insert: db_mtx
124 *   dbuf_hash_remove: db_mtx
125 *

```

```

126 * db_mtx (meta-leaf)
127 * must be held before:
128 *   dn_mtx, dn_dirty_mtx, dd_lock (leaf mutexes)
129 * protects:
130 *   db_state
131 *   db_holds
132 *   db_buf
133 *   db_changed
134 *   db_data_pending
135 *   db_dirtied
136 *   db_link
137 *   db_dirty_node (??)
138 *   db_dirtycnt
139 *   db_d.*
140 *   db.*
141 * held from:
142 *   dbuf_dirty: dn_mtx, dn_dirty_mtx
143 *   dbuf_dirty->dsl_dir_willuse_space: dd_lock
144 *   dbuf_dirty->dbuf_new_block->dsl_dataset_block_freeable: dd_lock
145 *   dbuf_undirty: dn_dirty_mtx (db_d)
146 *   dbuf_write_done: dn_dirty_mtx (db_state)
147 *   dbuf.*
148 *   dmu_buf_update_user: none (db_d)
149 *   dmu_evict_user: none (db_d) (maybe can eliminate)
150 *   dbuf_find: none (db_holds)
151 *   dbuf_hash_insert: none (db_holds)
152 *   dmu_buf_read_array_impl: none (db_state, db_changed)
153 *   dmu_sync: none (db_dirty_node, db_d)
154 *   dnode_reallocate: none (db)
155 *
156 * dn_mtx (leaf)
157 * protects:
158 *   dn_dirty_dbufs
159 *   dn_ranges
160 *   phys accounting
161 *   dn_allocated_txg
162 *   dn_free_txg
163 *   dn_assigned_txg
164 *   dd_assigned_tx
165 *   dn_notxholds
166 *   dn_dirtyctx
167 *   dn_dirtyctx_firstset
168 *   (dn_phys copy fields?)
169 *   (dn_phys contents?)
170 * held from:
171 *   dnode.*
172 *   dbuf_dirty: none
173 *   dbuf_sync: none (phys accounting)
174 *   dbuf_undirty: none (dn_ranges, dn_dirty_dbufs)
175 *   dbuf_write_done: none (phys accounting)
176 *   dmu_object_info_from_dnode: none (accounting)
177 *   dmu_tx_commit: none
178 *   dmu_tx_hold_object_impl: none
179 *   dmu_tx_try_assign: dn_notxholds(cv)
180 *   dmu_tx_unassign: none
181 *
182 * dd_lock
183 * must be held before:
184 *   ds_lock
185 *   ancestors' dd_lock
186 * protects:
187 *   dd_prop_cbs
188 *   dd_sync_*
189 *   dd_used_bytes
190 *   dd_tempreserved
191 *   dd_space_towrite

```

```

192 *   dd_myname
193 *   dd_phys accounting?
194 * held from:
195 *   dsl_dir_*
196 *   dsl_prop_changed_notify: none (dd_prop_cbs)
197 *   dsl_prop_register: none (dd_prop_cbs)
198 *   dsl_prop_unregister: none (dd_prop_cbs)
199 *   dsl_dataset_block_freeable: none (dd_sync_*)
200 *
201 * os_lock (leaf)
202 * protects:
203 *   os_dirty_dnodes
204 *   os_free_dnodes
205 *   os_dnodes
206 *   os_downgraded_dbufs
207 *   dn_dirtyblksz
208 *   dn_dirty_link
209 * held from:
210 *   dnode_create: none (os_dnodes)
211 *   dnode_destroy: none (os_dnodes)
212 *   dnode_setdirty: none (dn_dirtyblksz, os_*_dnodes)
213 *   dnode_free: none (dn_dirtyblksz, os_*_dnodes)
214 *
215 * ds_lock
216 * protects:
217 *   ds_objset
218 *   ds_open_refcount
219 *   ds_snapname
220 *   ds_phys accounting
221 *   ds_phys userrefs zapobj
222 *   ds_reserved
223 * held from:
224 *   dsl_dataset_*
225 *
226 * dr_mtx (leaf)
227 * protects:
228 *   dr_children
229 * held from:
230 *   dbuf_dirty
231 *   dbuf_undirty
232 *   dbuf_sync_indirect
233 *   dnode_new_blkid
234 */

236 struct objset;
237 struct dmu_pool;

239 typedef struct dmu_xuio {
240     int next;
241     int cnt;
242     struct arc_buf **bufs;
243     iovec_t *iovp;
244 } dmu_xuio_t;

246 typedef struct xuio_stats {
247     /* loaned yet not returned arc_buf */
248     kstat_named_t xuiostat_onloan_rbuf;
249     kstat_named_t xuiostat_onloan_wbuf;
250     /* whether a copy is made when loaning out a read buffer */
251     kstat_named_t xuiostat_rbuf_copied;
252     kstat_named_t xuiostat_rbuf_nocopy;
253     /* whether a copy is made when assigning a write buffer */
254     kstat_named_t xuiostat_wbuf_copied;
255     kstat_named_t xuiostat_wbuf_nocopy;
256 } xuio_stats_t;

```

```
258 static xuiostat_t xuiostat = {
259     { "onloan_read_buf",    KSTAT_DATA_UINT64 },
260     { "onloan_write_buf",   KSTAT_DATA_UINT64 },
261     { "read_buf_copied",    KSTAT_DATA_UINT64 },
262     { "read_buf_nocopy",    KSTAT_DATA_UINT64 },
263     { "write_buf_copied",   KSTAT_DATA_UINT64 },
264     { "write_buf_nocopy",   KSTAT_DATA_UINT64 }
265 };

267 #define XUIOSTAT_INCR(stat, val) \
268     atomic_add_64(&xuiostat.stat.value.ui64, (val))
269 #define XUIOSTAT_BUMP(stat)      XUIOSTAT_INCR(stat, 1)

271 /*
272  * The list of data whose inclusion in a send stream can be pending from
273  * one call to backup_cb to another. Multiple calls to dump_free() and
274  * dump_freeobjects() can be aggregated into a single DRR_FREE or
275  * DRR_FREEOBJECTS replay record.
276  */
277 typedef enum {
278     PENDING_NONE,
279     PENDING_FREE,
280     PENDING_FREEOBJECTS
281 } dmu_pending_t;

283 typedef struct dmu_sendarg {
284     list_node_t dsa_link;
285     dmu_replay_record_t *dsa_drr;
286     vnode_t *dsa_vp;
287     int dsa_outfd;
288     struct proc *dsa_proc;
289     offset_t *dsa_off;
290     objset_t *dsa_os;
291     zio_cksum_t dsa_zc;
292     uint64_t dsa_toguid;
293     int dsa_err;
294     dmu_pending_t dsa_pending_op;
295 } dmu_sendarg_t;

298 #ifdef __cplusplus
299 }
300 #endif

302 #endif /* _SYS_DMU_IMPL_H */
```

```

*****
10615 Thu Apr 25 16:14:59 2013
new/usr/src/uts/common/fs/zfs/sys/dnode.h
3752 want more verifiable dbuf user eviction
Submitted by: Justin Gibbs <justing@spectrallogic.com>
Submitted by: Will Andrews <willa@spectrallogic.com>
*****
unchanged_portion_omitted_

244 typedef struct dnode_children {
245     dmu_buf_user_t db_evict;          /* User evict data */
246 #endif /* ! codereview */
247     size_t dnc_count;                /* number of children */
248     dnode_handle_t dnc_children[1]; /* sized dynamically */
249 } dnode_children_t;

251 typedef struct free_range {
252     avl_node_t fr_node;
253     uint64_t fr_blkid;
254     uint64_t fr_nblks;
255 } free_range_t;

257 dnode_t *dnode_special_open(struct objset *dd, dnode_phys_t *dnp,
258     uint64_t object, dnode_handle_t *dnh);
259 void dnode_special_close(dnode_handle_t *dnh);

261 void dnode_setbonuslen(dnode_t *dn, int newsize, dmu_tx_t *tx);
262 void dnode_setbonus_type(dnode_t *dn, dmu_object_type_t, dmu_tx_t *tx);
263 void dnode_rm_spill(dnode_t *dn, dmu_tx_t *tx);

265 int dnode_hold(struct objset *dd, uint64_t object,
266     void *ref, dnode_t **dnp);
267 int dnode_hold_impl(struct objset *dd, uint64_t object, int flag,
268     void *ref, dnode_t **dnp);
269 boolean_t dnode_add_ref(dnode_t *dn, void *ref);
270 void dnode_rele(dnode_t *dn, void *ref);
271 void dnode_setdirty(dnode_t *dn, dmu_tx_t *tx);
272 void dnode_sync(dnode_t *dn, dmu_tx_t *tx);
273 void dnode_allocate(dnode_t *dn, dmu_object_type_t ot, int blocksize, int ibs,
274     dmu_object_type_t bonustype, int bonuslen, dmu_tx_t *tx);
275 void dnode_reallocate(dnode_t *dn, dmu_object_type_t ot, int blocksize,
276     dmu_object_type_t bonustype, int bonuslen, dmu_tx_t *tx);
277 void dnode_free(dnode_t *dn, dmu_tx_t *tx);
278 void dnode_byteswap(dnode_phys_t *dnp);
279 void dnode_buf_byteswap(void *buf, size_t size);
280 void dnode_verify(dnode_t *dn);
281 int dnode_set_blkisz(dnode_t *dn, uint64_t size, int ibs, dmu_tx_t *tx);
282 void dnode_free_range(dnode_t *dn, uint64_t off, uint64_t len, dmu_tx_t *tx);
283 void dnode_clear_range(dnode_t *dn, uint64_t blkid,
284     uint64_t nblks, dmu_tx_t *tx);
285 void dnode_diduse_space(dnode_t *dn, int64_t space);
286 void dnode_willuse_space(dnode_t *dn, int64_t space, dmu_tx_t *tx);
287 void dnode_new_blkid(dnode_t *dn, uint64_t blkid, dmu_tx_t *tx, boolean_t);
288 uint64_t dnode_block_freed(dnode_t *dn, uint64_t blkid);
289 void dnode_init(void);
290 void dnode_fini(void);
291 int dnode_next_offset(dnode_t *dn, int flags, uint64_t *off,
292     int minlvl, uint64_t blkfill, uint64_t txg);
293 void dnode_evict_dbufs(dnode_t *dn);

295 #ifdef ZFS_DEBUG

297 /*
298  * There should be a ## between the string literal and fmt, to make it
299  * clear that we're joining two strings together, but that piece of shit
300  * gcc doesn't support that preprocessor token.

```

```

301 */
302 #define dprintf_dnode(dn, fmt, ...) do { \
303     if (zfs_flags & ZFS_DEBUG_DPRINTF) { \
304         char __db_buf[32]; \
305         uint64_t __db_obj = (dn)->dn_object; \
306         if (__db_obj == DMU_META_DNODE_OBJECT) \
307             (void) strcpy(__db_buf, "mdn"); \
308         else \
309             (void) snprintf(__db_buf, sizeof (__db_buf), "%lld", \
310                 (u_longlong_t) __db_obj); \
311         dprintf_ds((dn)->dn_objset->os_dsl_dataset, "obj=%s " fmt, \
312             __db_buf, __VA_ARGS__); \
313     } \
314     _NOTE(CONSTCOND) } while (0)

316 #define DNODE_VERIFY(dn)                dnode_verify(dn)
317 #define FREE_VERIFY(db, start, end, tx) free_verify(db, start, end, tx)

319 #else

321 #define dprintf_dnode(db, fmt, ...)
322 #define DNODE_VERIFY(dn)
323 #define FREE_VERIFY(db, start, end, tx)

325 #endif

327 #ifdef __cplusplus
328 }
329 #endif

331 #endif /* _SYS_DNODE_H */

```

```

*****
11005 Thu Apr 25 16:14:59 2013
new/usr/src/uts/common/fs/zfs/sys/dsl_dataset.h
3752 want more verifiable dbuf user eviction
Submitted by: Justin Gibbs <justing@spectralogic.com>
Submitted by: Will Andrews <willa@spectralogic.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) 2012 by Delphix. All rights reserved.
24 * Copyright (c) 2012, Joyent, Inc. All rights reserved.
25 */

27 #ifndef _SYS_DSL_DATASET_H
28 #define _SYS_DSL_DATASET_H

30 #include <sys/dmu.h>
31 #include <sys/spa.h>
32 #include <sys/txg.h>
33 #include <sys/zio.h>
34 #include <sys/bplist.h>
35 #include <sys/dsl_synctask.h>
36 #include <sys/zfs_context.h>
37 #include <sys/dsl_deadlist.h>
38 #include <sys/refcount.h>

40 #ifdef __cplusplus
41 extern "C" {
42 #endif

44 struct dsl_dataset;
45 struct dsl_dir;
46 struct dsl_pool;

48 #define DS_HAS_PHYS(ds) \
49     ((ds)->ds_dbuf != NULL && (ds)->ds_dbuf->db_data != NULL)

51 #endif /* ! codereview */
52 #define DS_FLAG_INCONSISTENT (1ULL<<0)
53 #define DS_IS_INCONSISTENT(ds) \
54     ((ds)->ds_phys->ds_flags & DS_FLAG_INCONSISTENT)
55 /*
56 * Note: nopromote can not yet be set, but we want support for it in this
57 * on-disk version, so that we don't need to upgrade for it later.
58 */
59 #define DS_FLAG_NOPROMOTE (1ULL<<1)

```

```

61 /*
62  * DS_FLAG_UNIQUE_ACCURATE is set if ds_unique_bytes has been correctly
63  * calculated for head datasets (starting with SPA_VERSION_UNIQUE_ACCURATE,
64  * refquota/refreservations).
65  */
66 #define DS_FLAG_UNIQUE_ACCURATE (1ULL<<2)

68 /*
69  * DS_FLAG_DEFER_DESTROY is set after 'zfs destroy -d' has been called
70  * on a dataset. This allows the dataset to be destroyed using 'zfs release'.
71  */
72 #define DS_FLAG_DEFER_DESTROY (1ULL<<3)
73 #define DS_IS_DEFER_DESTROY(ds) \
74     ((ds)->ds_phys->ds_flags & DS_FLAG_DEFER_DESTROY)

76 /*
77  * DS_FLAG_CI_DATASET is set if the dataset contains a file system whose
78  * name lookups should be performed case-insensitively.
79  */
80 #define DS_FLAG_CI_DATASET (1ULL<<16)

82 #define DS_CREATE_FLAG_NODIRTY (1ULL<<24)

84 typedef struct dsl_dataset_phys {
85     uint64_t ds_dir_obj; /* DMU_OT_DSL_DIR */
86     uint64_t ds_prev_snap_obj; /* DMU_OT_DSL_DATASET */
87     uint64_t ds_prev_snap_txg;
88     uint64_t ds_next_snap_obj; /* DMU_OT_DSL_DATASET */
89     uint64_t ds_snapnames_zapobj; /* DMU_OT_DSL_DS_SNAP_MAP 0 for snaps */
90     uint64_t ds_num_children; /* clone/snap children; ==0 for head */
91     uint64_t ds_creation_time; /* seconds since 1970 */
92     uint64_t ds_creation_txg;
93     uint64_t ds_deadlist_obj; /* DMU_OT_DEADLIST */
94     /*
95      * ds_referenced_bytes, ds_compressed_bytes, and ds_uncompressed_bytes
96      * include all blocks referenced by this dataset, including those
97      * shared with any other datasets.
98      */
99     uint64_t ds_referenced_bytes;
100     uint64_t ds_compressed_bytes;
101     uint64_t ds_uncompressed_bytes;
102     uint64_t ds_unique_bytes; /* only relevant to snapshots */
103     /*
104      * The ds_fsid_guid is a 56-bit ID that can change to avoid
105      * collisions. The ds_guid is a 64-bit ID that will never
106      * change, so there is a small probability that it will collide.
107      */
108     uint64_t ds_fsid_guid;
109     uint64_t ds_guid;
110     uint64_t ds_flags; /* DS_FLAG */
111     blkptr_t ds_bp;
112     uint64_t ds_next_clones_obj; /* DMU_OT_DSL_CLONES */
113     uint64_t ds_props_obj; /* DMU_OT_DSL_PROPS for snaps */
114     uint64_t ds_userrefs_obj; /* DMU_OT_USERREFS */
115     uint64_t ds_pad[5]; /* pad out to 320 bytes for good measure */
116 } dsl_dataset_phys_t;

118 typedef struct dsl_dataset_dbuf {
119     uint8_t dsdb_pad[offsetof(dmu_buf_t, db_data)];
120     dsl_dataset_phys_t *dsdb_data;
121 } dsl_dataset_dbuf_t;

123 #endif /* ! codereview */
124 typedef struct dsl_dataset {
125     dmu_buf_user_t db_evict;

```

```

127 #endif /* ! codereview */
128 /* Immutable: */
129 struct dsl_dir *ds_dir;
130 union {
131     dmu_buf_t *ds_dmu_db;
132     dsl_dataset_dbuf_t *ds_db;
133 } ds_db_u;
134 dsl_dataset_phys_t *ds_phys;
135 dmu_buf_t *ds_dbuf;
136 uint64_t ds_object;
137 uint64_t ds_fsid_guid;
138
139 /* only used in syncing context, only valid for non-snapshots: */
140 struct dsl_dataset *ds_prev;
141
142 /* has internal locking: */
143 dsl_deadlist_t ds_deadlist;
144 bplist_t ds_pending_deadlist;
145
146 /* protected by lock on pool's dp_dirty_datasets list */
147 txg_node_t ds_dirty_link;
148 list_node_t ds_synced_link;
149
150 /*
151  * ds_phys->ds_accounting is also protected by ds_lock.
152  * Protected by ds_lock:
153  */
154 kmutex_t ds_lock;
155 objset_t *ds_objset;
156 uint64_t ds_userrefs;
157 void *ds_owner;
158
159 /*
160  * Long holds prevent the ds from being destroyed; they allow the
161  * ds to remain held even after dropping the dp_config_rwlock.
162  * Owning counts as a long hold. See the comments above
163  * dsl_pool_hold() for details.
164  */
165 refcount_t ds_longholds;
166
167 /* no locking; only for making guesses */
168 uint64_t ds_trysnap_txg;
169
170 /* for objset_open() */
171 kmutex_t ds_opening_lock;
172
173 uint64_t ds_reserved; /* cached reservation */
174 uint64_t ds_quota; /* cached refquota */
175
176 kmutex_t ds_sendstream_lock;
177 list_t ds_sendstreams;
178
179 /* Protected by ds_lock; keep at end of struct for better locality */
180 char ds_snapname[MAXNAMELEN];
181 } dsl_dataset_t;
182
183 /* See sys/dmu.h:dmu_buf_user_t for why we have these. */
184 #define ds_dbuf ds_db_u.ds_dmu_db
185 #define ds_phys ds_db_u.ds_db->dsdb_data
186
187 struct dsl_ds_destroyarg {
188     dsl_dataset_t *ds; /* ds to destroy */
189     dsl_dataset_t *rm_origin; /* also remove our origin? */
190     boolean_t is_origin_rm; /* set if removing origin snap */
191     boolean_t defer; /* destroy -d requested? */

```

```

190     boolean_t releasing; /* destroying due to release? */
191     boolean_t need_prep; /* do we need to retry due to EBUSY? */
192 };
193
194 #endif /* ! codereview */
195 /*
196  * The max length of a temporary tag prefix is the number of hex digits
197  * required to express UINTE64_MAX plus one for the hyphen.
198  */
199 #define MAX_TAG_PREFIX_LEN 17
200
201 #define dsl_dataset_is_snapshot(ds) \
202     ((ds)->ds_phys->ds_num_children != 0)
203
204 #define DS_UNIQUE_IS_ACCURATE(ds) \
205     (((ds)->ds_phys->ds_flags & DS_FLAG_UNIQUE_ACCURATE) != 0)
206
207 int dsl_dataset_hold(struct dsl_pool *dp, const char *name, void *tag,
208     dsl_dataset_t **dsp);
209 int dsl_dataset_hold_obj(struct dsl_pool *dp, uint64_t dsobj, void *tag,
210     dsl_dataset_t **);
211 void dsl_dataset_rele(dsl_dataset_t *ds, void *tag);
212 int dsl_dataset_own(struct dsl_pool *dp, const char *name,
213     void *tag, dsl_dataset_t **dsp);
214 int dsl_dataset_own_obj(struct dsl_pool *dp, uint64_t dsobj,
215     void *tag, dsl_dataset_t **dsp);
216 void dsl_dataset_disown(dsl_dataset_t *ds, void *tag);
217 void dsl_dataset_name(dsl_dataset_t *ds, char *name);
218 boolean_t dsl_dataset_tryown(dsl_dataset_t *ds, void *tag);
219 void dsl_register_onexit_hold_cleanup(dsl_dataset_t *ds, const char *htag,
220     minor_t minor);
221 uint64_t dsl_dataset_create_sync(dsl_dir_t *pds, const char *lastname,
222     dsl_dataset_t *origin, uint64_t flags, cred_t *, dmu_tx_t *);
223 uint64_t dsl_dataset_create_sync_dd(dsl_dir_t *dd, dsl_dataset_t *origin,
224     uint64_t flags, dmu_tx_t *);
225 int dsl_dataset_snapshot(nvlist_t *snaps, nvlist_t *props, nvlist_t *errors);
226 int dsl_dataset_promote(const char *name, char *conflsnap);
227 int dsl_dataset_clone_swap(dsl_dataset_t *clone, dsl_dataset_t *origin_head,
228     boolean_t force);
229 int dsl_dataset_rename_snapshot(const char *fsname,
230     const char *oldsnapname, const char *newsnapname, boolean_t recursive);
231 int dsl_dataset_snapshot_tmp(const char *fsname, const char *snapname,
232     minor_t cleanup_minor, const char *htag);
233
234 blkptr_t *dsl_dataset_get_blkptr(dsl_dataset_t *ds);
235 void dsl_dataset_set_blkptr(dsl_dataset_t *ds, blkptr_t *bp, dmu_tx_t *tx);
236
237 spa_t *dsl_dataset_get_spa(dsl_dataset_t *ds);
238
239 boolean_t dsl_dataset_modified_since_lastsnap(dsl_dataset_t *ds);
240
241 void dsl_dataset_sync(dsl_dataset_t *os, zio_t *zio, dmu_tx_t *tx);
242
243 void dsl_dataset_block_born(dsl_dataset_t *ds, const blkptr_t *bp,
244     dmu_tx_t *tx);
245 int dsl_dataset_block_kill(dsl_dataset_t *ds, const blkptr_t *bp,
246     dmu_tx_t *tx, boolean_t async);
247 boolean_t dsl_dataset_block_freeable(dsl_dataset_t *ds, const blkptr_t *bp,
248     uint64_t blk_birth);
249 uint64_t dsl_dataset_prev_snap_txg(dsl_dataset_t *ds);
250
251 void dsl_dataset_dirty(dsl_dataset_t *ds, dmu_tx_t *tx);
252 void dsl_dataset_stats(dsl_dataset_t *os, nvlist_t *nv);
253 void dsl_dataset_fast_stat(dsl_dataset_t *ds, dmu_objset_stats_t *stat);
254 void dsl_dataset_space(dsl_dataset_t *ds,
255     uint64_t *refdbytesp, uint64_t *availbytesp,

```

```
256 uint64_t *usedobjsp, uint64_t *availobjsp);
257 uint64_t dsl_dataset_fsid_guid(dsl_dataset_t *ds);
258 int dsl_dataset_space_written(dsl_dataset_t *oldsnap, dsl_dataset_t *new,
259 uint64_t *usedp, uint64_t *compp, uint64_t *uncompp);
260 int dsl_dataset_space_wouldfree(dsl_dataset_t *firstsnap, dsl_dataset_t *last,
261 uint64_t *usedp, uint64_t *compp, uint64_t *uncompp);
262 boolean_t dsl_dataset_is_dirty(dsl_dataset_t *ds);

264 int dsl_dsobj_to_dsname(char *pname, uint64_t obj, char *buf);

266 int dsl_dataset_check_quota(dsl_dataset_t *ds, boolean_t check_quota,
267 uint64_t asize, uint64_t inflight, uint64_t *used,
268 uint64_t *ref_rsrv);
269 int dsl_dataset_set_refquota(const char *dsname, zprop_source_t source,
270 uint64_t quota);
271 int dsl_dataset_set_refreservation(const char *dsname, zprop_source_t source,
272 uint64_t reservation);

274 boolean_t dsl_dataset_is_before(dsl_dataset_t *later, dsl_dataset_t *earlier);
275 void dsl_dataset_long_hold(dsl_dataset_t *ds, void *tag);
276 void dsl_dataset_long_rele(dsl_dataset_t *ds, void *tag);
277 boolean_t dsl_dataset_long_held(dsl_dataset_t *ds);

279 int dsl_dataset_clone_swap_check_impl(dsl_dataset_t *clone,
280 dsl_dataset_t *origin_head, boolean_t force);
281 void dsl_dataset_clone_swap_sync_impl(dsl_dataset_t *clone,
282 dsl_dataset_t *origin_head, dmu_tx_t *tx);
283 int dsl_dataset_snapshot_check_impl(dsl_dataset_t *ds, const char *snapname,
284 dmu_tx_t *tx);
285 void dsl_dataset_snapshot_sync_impl(dsl_dataset_t *ds, const char *snapname,
286 dmu_tx_t *tx);

288 void dsl_dataset_remove_from_next_clones(dsl_dataset_t *ds, uint64_t obj,
289 dmu_tx_t *tx);
290 void dsl_dataset_recalc_head_uniq(dsl_dataset_t *ds);
291 int dsl_dataset_get_snapname(dsl_dataset_t *ds);
292 int dsl_dataset_snap_lookup(dsl_dataset_t *ds, const char *name,
293 uint64_t *value);
294 int dsl_dataset_snap_remove(dsl_dataset_t *ds, const char *name, dmu_tx_t *tx);
295 void dsl_dataset_set_refreservation_sync_impl(dsl_dataset_t *ds,
296 zprop_source_t source, uint64_t value, dmu_tx_t *tx);
297 int dsl_dataset_rollback(const char *fsname);

299 #ifdef ZFS_DEBUG
300 #define dprintf_ds(ds, fmt, ...) do { \
301 if (zfs_flags & ZFS_DEBUG_DPRINTF) { \
302 char *__ds_name = kmem_alloc(MAXNAMELEN, KM_SLEEP); \
303 dsl_dataset_name(ds, __ds_name); \
304 dprintf("ds=%s " fmt, __ds_name, __VA_ARGS__); \
305 kmem_free(__ds_name, MAXNAMELEN); \
306 } \
307 _NOTE(CONSTCOND) } while (0)
308 #else
309 #define dprintf_ds(dd, fmt, ...)
310 #endif

312 #ifdef __cplusplus
313 }
314 #endif

316 #endif /* _SYS_DSL_DATASET_H */
```

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

1

```
*****
2790 Thu Apr 25 16:15:00 2013
new/usr/src/uts/common/fs/zfs/sys/dsl_deleg.h
3752 want more verifiable dbuf user eviction
Submitted by: Justin Gibbs <justing@spectrallogic.com>
Submitted by: Will Andrews <willa@spectrallogic.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) 2007, 2010, Oracle and/or its affiliates. All rights reserved.
23 * Copyright (c) 2012 by Delphix. All rights reserved.
24 */

26 #ifndef _SYS_DSL_DELEG_H
27 #define _SYS_DSL_DELEG_H

29 #include <sys/zfs_context.h>
30 #endif /* ! codereview */
31 #include <sys/dmu.h>
32 #include <sys/dsl_pool.h>
33 #include <sys/zfs_context.h>

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

38 #define ZFS_DELEG_PERM_NONE          ""
39 #define ZFS_DELEG_PERM_CREATE        "create"
40 #define ZFS_DELEG_PERM_DESTROY       "destroy"
41 #define ZFS_DELEG_PERM_SNAPSHOT      "snapshot"
42 #define ZFS_DELEG_PERM_ROLLBACK     "rollback"
43 #define ZFS_DELEG_PERM_CLONE         "clone"
44 #define ZFS_DELEG_PERM_PROMOTE       "promote"
45 #define ZFS_DELEG_PERM_RENAME        "rename"
46 #define ZFS_DELEG_PERM_MOUNT         "mount"
47 #define ZFS_DELEG_PERM_SHARE         "share"
48 #define ZFS_DELEG_PERM_SEND          "send"
49 #define ZFS_DELEG_PERM_RECEIVE       "receive"
50 #define ZFS_DELEG_PERM_ALLOW         "allow"
51 #define ZFS_DELEG_PERM_USERPROP      "userprop"
52 #define ZFS_DELEG_PERM_VSCAN         "vscan"
53 #define ZFS_DELEG_PERM_USERQUOTA     "userquota"
54 #define ZFS_DELEG_PERM_GROUPQUOTA    "groupquota"
55 #define ZFS_DELEG_PERM_USERUSED      "userused"
56 #define ZFS_DELEG_PERM_GROUPUSED     "groupused"
57 #define ZFS_DELEG_PERM_HOLD          "hold"
58 #define ZFS_DELEG_PERM_RELEASE       "release"
```

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

2

```
59 #define ZFS_DELEG_PERM_DIFF        "diff"

61 /*
62  * Note: the names of properties that are marked delegatable are also
63  * valid delegated permissions
64  */

66 int dsl_deleg_get(const char *ddname, nvlist_t **nvp);
67 int dsl_deleg_set(const char *ddname, nvlist_t *nvp, boolean_t unset);
68 int dsl_deleg_access(const char *ddname, const char *perm, cred_t *cr);
69 int dsl_deleg_access_impl(struct dsl_dataset *ds, const char *perm, cred_t *cr);
70 void dsl_deleg_set_create_perms(dsl_dir_t *dd, dmu_tx_t *tx, cred_t *cr);
71 int dsl_deleg_can_allow(char *ddname, nvlist_t *nvp, cred_t *cr);
72 int dsl_deleg_can_unallow(char *ddname, nvlist_t *nvp, cred_t *cr);
73 int dsl_deleg_destroy(objset_t *os, uint64_t zapobj, dmu_tx_t *tx);
74 boolean_t dsl_delegation_on(objset_t *os);

76 #ifdef __cplusplus
77 }
_____
unchanged_portion_omitted
```

```

*****
5809 Thu Apr 25 16:15:00 2013
new/usr/src/uts/common/fs/zfs/sys/dsl_dir.h
3752 want more verifiable dbuf user eviction
Submitted by: Justin Gibbs <justing@spectrallogic.com>
Submitted by: Will Andrews <willa@spectrallogic.com>
*****
unchanged_portion_omitted

77 typedef struct dsl_dir_dbuf {
78     uint8_t dddb_pad[offsetof(dmu_buf_t, db_data)];
79     dsl_dir_phys_t *ddb_data;
80 } dsl_dir_dbuf_t;

82 #endif /* ! codereview */
83 struct dsl_dir {
84     dmu_buf_user_t db_evict;

86 #endif /* ! codereview */
87     /* These are immutable; no lock needed: */
88     uint64_t dd_object;
89     union {
90         dmu_buf_t *dd_dmu_db;
91         dsl_dir_dbuf_t *dd_db;
92     } dd_db_u;
93     dsl_dir_phys_t *dd_phys;
94     dmu_buf_t *dd_dbuf;
95     dsl_pool_t *dd_pool;

96     /* protected by lock on pool's dp_dirty_dirs list */
97     txg_node_t dd_dirty_link;

98     /* protected by dp_config_rwlock */
99     dsl_dir_t *dd_parent;

101     /* Protected by dd_lock */
102     kmutex_t dd_lock;
103     list_t dd_prop_cbs; /* list of dsl_prop_cb_record_t's */
104     timestruc_t dd_snap_cmtime; /* last time snapshot namespace changed */
105     uint64_t dd_origin_txg;

107     /* gross estimate of space used by in-flight tx's */
108     uint64_t dd_tempreserved[TXG_SIZE];
109     /* amount of space we expect to write; == amount of dirty data */
110     int64_t dd_space_towrite[TXG_SIZE];

112     /* protected by dd_lock; keep at end of struct for better locality */
113     char dd_myname[MAXNAMELEN];
114 };

116 /* See sys/dmu.h:dmu_buf_user_t for why we have these. */
117 #define dd_dbuf dd_db_u.dd_dmu_db
118 #define dd_phys dd_db_u.dd_db->ddb_data

120 #endif /* ! codereview */
121 void dsl_dir_rele(dsl_dir_t *dd, void *tag);
122 int dsl_dir_hold(dsl_pool_t *dp, const char *name, void *tag,
123     dsl_dir_t **, const char **tail);
124 int dsl_dir_hold_obj(dsl_pool_t *dp, uint64_t ddbobj,
125     const char *tail, void *tag, dsl_dir_t **);
126 void dsl_dir_name(dsl_dir_t *dd, char *buf);
127 int dsl_dir_namelen(dsl_dir_t *dd);
128 uint64_t dsl_dir_create_sync(dsl_pool_t *dp, dsl_dir_t *pds,
129     const char *name, dmu_tx_t *tx);
130 void dsl_dir_stats(dsl_dir_t *dd, nvlist_t *nv);
131 uint64_t dsl_dir_space_available(dsl_dir_t *dd,

```

```

132     dsl_dir_t *ancestor, int64_t delta, int ondiskonly);
133 void dsl_dir_dirty(dsl_dir_t *dd, dmu_tx_t *tx);
134 void dsl_dir_sync(dsl_dir_t *dd, dmu_tx_t *tx);
135 int dsl_dir_tempreserve_space(dsl_dir_t *dd, uint64_t mem,
136     uint64_t asize, uint64_t fsize, uint64_t usize, void **tr_cookiep,
137     dmu_tx_t *tx);
138 void dsl_dir_tempreserve_clear(void *tr_cookie, dmu_tx_t *tx);
139 void dsl_dir_willuse_space(dsl_dir_t *dd, int64_t space, dmu_tx_t *tx);
140 void dsl_dir_diduse_space(dsl_dir_t *dd, dd_used_t type,
141     int64_t used, int64_t compressed, int64_t uncompressed, dmu_tx_t *tx);
142 void dsl_dir_transfer_space(dsl_dir_t *dd, int64_t delta,
143     dd_used_t oldtype, dd_used_t newtype, dmu_tx_t *tx);
144 int dsl_dir_set_quota(const char *ddname, zprop_source_t source,
145     uint64_t quota);
146 int dsl_dir_set_reservation(const char *ddname, zprop_source_t source,
147     uint64_t reservation);
148 int dsl_dir_rename(const char *oldname, const char *newname);
149 int dsl_dir_transfer_possible(dsl_dir_t *sdd, dsl_dir_t *tdd, uint64_t space);
150 boolean_t dsl_dir_is_clone(dsl_dir_t *dd);
151 void dsl_dir_new_reservation(dsl_dir_t *dd, struct dsl_dataset *ds,
152     uint64_t reservation, cred_t *cr, dmu_tx_t *tx);
153 void dsl_dir_snap_cmtime_update(dsl_dir_t *dd);
154 timestruc_t dsl_dir_snap_cmtime(dsl_dir_t *dd);
155 void dsl_dir_set_reservation_sync_impl(dsl_dir_t *dd, uint64_t value,
156     dmu_tx_t *tx);

158 /* internal reserved dir name */
159 #define MOS_DIR_NAME "$MOS"
160 #define ORIGIN_DIR_NAME "$ORIGIN"
161 #define XLATION_DIR_NAME "$XLATION"
162 #define FREE_DIR_NAME "$FREE"

164 #ifdef ZFS_DEBUG
165 #define dprintf_dd(dd, fmt, ...) do { \
166     if (zfs_flags & ZFS_DEBUG_DPRINTF) { \
167         char *__ds_name = kmem_alloc(MAXNAMELEN + strlen(MOS_DIR_NAME) + 1, \
168             KM_SLEEP); \
169         dsl_dir_name(dd, __ds_name); \
170         dprintf("dd=%s " fmt, __ds_name, __VA_ARGS__); \
171         kmem_free(__ds_name, MAXNAMELEN + strlen(MOS_DIR_NAME) + 1); \
172     } \
173     _NOTE(CONSTCOND) } while (0)
174 #else
175 #define dprintf_dd(dd, fmt, ...)
176 #endif

178 #ifdef __cplusplus
179 }
180 #endif

182 #endif /* _SYS_DSL_DIR_H */

```

```

*****
      8415 Thu Apr 25 16:15:00 2013
new/usr/src/uts/common/fs/zfs/sys/sa_impl.h
3752 want more verifiable dbuf user eviction
Submitted by: Justin Gibbs <justing@spectrallogic.com>
Submitted by: Will Andrews <willa@spectrallogic.com>
*****
_____unchanged_portion_omitted_____

205 /*
206  * Opaque handle used for most sa functions
207  *
208  * This needs to be kept as small as possible.
209  */

211 struct sa_handle {
212     dmu_buf_user_t db_evict;
213 #endif /* ! codereview */
214     kmutex_t sa_lock;
215     dmu_buf_t *sa_bonus;
216     dmu_buf_t *sa_spill;
217     objset_t *sa_os;
218     void *sa_userp;
219     sa_idx_tab_t *sa_bonus_tab; /* idx of bonus */
220     sa_idx_tab_t *sa_spill_tab; /* only present if spill activated */
221 };

223 #define SA_GET_DB(hdl, type) \
224     (dmu_buf_impl_t *)((type == SA_BONUS) ? hdl->sa_bonus : hdl->sa_spill)

226 #define SA_GET_HDR(hdl, type) \
227     ((sa_hdr_phys_t *)((dmu_buf_impl_t *) (SA_GET_DB(hdl, \
228     type))->db.db_data))

230 #define SA_IDX_TAB_GET(hdl, type) \
231     (type == SA_BONUS ? hdl->sa_bonus_tab : hdl->sa_spill_tab)

233 #define IS_SA_BONUSTYPE(a) \
234     ((a == DMU_OT_SA) ? B_TRUE : B_FALSE)

236 #define SA_BONUSTYPE_FROM_DB(db) \
237     (dmu_get_bonustype((dmu_buf_t *)db))

239 #define SA_BLKPTR_SPACE (DN_MAX_BONUSLEN - sizeof(blkptr_t))

241 #define SA_LAYOUT_NUM(x, type) \
242     (((IS_SA_BONUSTYPE(type) ? 0 : ((IS_SA_BONUSTYPE(type)) && \
243     ((SA_HDR_LAYOUT_NUM(x) == 0)) ? 1 : SA_HDR_LAYOUT_NUM(x))))

246 #define SA_REGISTERED_LEN(sa, attr) sa->sa_attr_table[attr].sa_length

248 #define SA_ATTR_LEN(sa, idx, attr, hdr) ((SA_REGISTERED_LEN(sa, attr) == 0) ? \
249     hdr->sa_lengths[TOC_LEN_IDX(idx->sa_idx_tab[attr])] : \
250     SA_REGISTERED_LEN(sa, attr))

252 #define SA_SET_HDR(hdr, num, size) \
253     { \
254         hdr->sa_magic = SA_MAGIC; \
255         SA_HDR_LAYOUT_INFO_ENCODE(hdr->sa_layout_info, num, size); \
256     }

258 #define SA_ATTR_INFO(sa, idx, hdr, attr, bulk, type, hdl) \
259     { \
260         bulk.sa_size = SA_ATTR_LEN(sa, idx, attr, hdr); \
261         bulk.sa_buftype = type; \

```

```

262         bulk.sa_addr = \
263             (void *)((uintptr_t)TOC_OFF(idx->sa_idx_tab[attr]) + \
264             (uintptr_t)hdr); \
265     }

267 #define SA_HDR_SIZE_MATCH_LAYOUT(hdr, tb) \
268     (SA_HDR_SIZE(hdr) == (sizeof(sa_hdr_phys_t) + \
269     (tb->lot_var_sizes > 1 ? P2ROUNDUP((tb->lot_var_sizes - 1) * \
270     sizeof(uint16_t), 8) : 0)))

272 int sa_add_impl(sa_handle_t *, sa_attr_type_t,
273     uint32_t, sa_data_locator_t, void *, dmu_tx_t *);

275 void sa_register_update_callback_locked(objset_t *, sa_update_cb_t *);
276 int sa_size_locked(sa_handle_t *, sa_attr_type_t, int *);

278 void sa_default_locator(void **, uint32_t *, uint32_t, boolean_t, void *);
279 int sa_attr_size(sa_os_t *, sa_idx_tab_t *, sa_attr_type_t,
280     uint16_t *, sa_hdr_phys_t *);

282 #ifdef __cplusplus
283 extern "C" {
284 #endif

286 #ifdef __cplusplus
287 }
288 #endif

290 #endif /* _SYS_SA_IMPL_H */

```

```

*****
7196 Thu Apr 25 16:15:01 2013
new/usr/src/uts/common/fs/zfs/sys/zap_impl.h
3752 want more verifiable dbuf user eviction
Submitted by: Justin Gibbs <justing@spectrallogic.com>
Submitted by: Will Andrews <willa@spectrallogic.com>
*****
unchanged_portion_omitted

65 typedef struct mzap_dbuf {
66     uint8_t mzdb_pad[offsetof(dmu_buf_t, db_data)];
67     mzap_phys_t *mzdb_data;
68 } mzap_dbuf_t;

70 #endif /* ! codereview */
71 typedef struct mzap_ent {
72     avl_node_t mze_node;
73     int mze_chunkid;
74     uint64_t mze_hash;
75     uint32_t mze_cd; /* copy from mze_phys->mze_cd */
76 } mzap_ent_t;

78 #define MZE_PHYS(zap, mze) \
79     (&(zap)->zap_m_phys->mz_chunk[(mze)->mze_chunkid])
65     (&(zap)->zap_m.zap_phys->mz_chunk[(mze)->mze_chunkid])

81 /*
82  * The (fat) zap is stored in one object. It is an array of
83  * 1<FZAP_BLOCK_SHIFT byte blocks. The layout looks like one of:
84  *
85  * ptrtbl fits in first block:
86  *     [zap_phys_t zap_ptrtbl_shift < 6] [zap_leaf_t] ...
87  *
88  * ptrtbl too big for first block:
89  *     [zap_phys_t zap_ptrtbl_shift >= 6] [zap_leaf_t] [ptrtbl] ...
90  *
91  */

93 struct dmu_buf;
94 struct zap_leaf;

96 #define ZBT_LEAF                ((1ULL << 63) + 0)
97 #define ZBT_HEADER              ((1ULL << 63) + 1)
98 #define ZBT_MICRO               ((1ULL << 63) + 3)
99 /* any other values are ptrtbl blocks */

101 /*
102  * the embedded pointer table takes up half a block:
103  * block size / entry size (2^3) / 2
104  */
105 #define ZAP_EMBEDDED_PTRTBL_SHIFT(zap) (FZAP_BLOCK_SHIFT(zap) - 3 - 1)

107 /*
108  * The embedded pointer table starts half-way through the block. Since
109  * the pointer table itself is half the block, it starts at (64-bit)
110  * word number (1<<ZAP_EMBEDDED_PTRTBL_SHIFT(zap)).
111  */
112 #define ZAP_EMBEDDED_PTRTBL_ENT(zap, idx) \
113     ((uint64_t *) (zap)->zap_f_phys) \
99     ((uint64_t *) (zap)->zap_f.zap_phys) \
114     [(idx) + (1<<ZAP_EMBEDDED_PTRTBL_SHIFT(zap))]

116 /*
117  * TAKE NOTE:
118  * If zap_phys_t is modified, zap_byteswap() must be modified.
119  */

```

```

120 typedef struct zap_phys {
121     uint64_t zap_block_type; /* ZBT_HEADER */
122     uint64_t zap_magic; /* ZAP_MAGIC */

124     struct zap_table_phys {
125         uint64_t zt_blk; /* starting block number */
126         uint64_t zt_numblks; /* number of blocks */
127         uint64_t zt_shift; /* bits to index it */
128         uint64_t zt_nextblk; /* next (larger) copy start block */
129         uint64_t zt_blks_copied; /* number source blocks copied */
130     } zap_ptrtbl;

132     uint64_t zap_freeblk; /* the next free block */
133     uint64_t zap_num_leafs; /* number of leafs */
134     uint64_t zap_num_entries; /* number of entries */
135     uint64_t zap_salt; /* salt to stir into hash function */
136     uint64_t zap_normflags; /* flags for u8_textprep_str() */
137     uint64_t zap_flags; /* zap_flags_t */
138     /*
139      * This structure is followed by padding, and then the embedded
140      * pointer table. The embedded pointer table takes up second
141      * half of the block. It is accessed using the
142      * ZAP_EMBEDDED_PTRTBL_ENT() macro.
143     */
144 } zap_phys_t;

146 typedef struct zap_table_phys zap_table_phys_t;

148 typedef struct fzap_dbuf {
149     uint8_t fzdb_pad[offsetof(dmu_buf_t, db_data)];
150     zap_phys_t *fzap_data;
151 } fzap_dbuf_t;

153 #endif /* ! codereview */
154 typedef struct zap {
155     dmu_buf_user_t db_evict;
156 #endif /* ! codereview */
157     objset_t *zap_objset;
158     uint64_t zap_object;
159     union {
160         dmu_buf_t *zap_dmu_db;
161         mzap_dbuf_t *mzap_db;
162         fzap_dbuf_t *fzap_db;
163     } zap_db_u;
164     struct dmu_buf *zap_dbuf;
165     krwlock_t zap_rwlock;
166     boolean_t zap_ismicro;
167     int zap_normflags;
168     uint64_t zap_salt;
169     union {
170         struct {
171             /* protects zap_num_entries */
172             zap_phys_t *zap_phys;

173             /*
174              * zap_num_entries_mtx protects
175              * zap_num_entries
176              */
177             kmutex_t zap_num_entries_mtx;
178             int zap_block_shift;
179         } zap_fat;
180         struct {
181             mzap_phys_t *zap_phys;
182             int16_t zap_num_entries;
183             int16_t zap_num_chunks;
184             int16_t zap_alloc_next;

```

```

178             avl_tree_t zap_avl;
179         } zap_micro;
180     } zap_u;
181 } zap_t;

183 /* See sys/dmu.h:dmu_buf_user_t for why we have these. */
184 #define zap_dbuf      zap_db_u.zap_dmu_db
185 #define zap_f        zap_u.zap_fat
186 #define zap_m        zap_u.zap_micro
187 #define zap_f_phys   zap_db_u.fzap_db->fzdb_data
188 #define zap_m_phys   zap_db_u.mzap_db->mzdb_data

190 #endif /* ! codereview */
191 typedef struct zap_name {
192     zap_t *zn_zap;
193     int zn_key_intlen;
194     const void *zn_key_orig;
195     int zn_key_orig_numints;
196     const void *zn_key_norm;
197     int zn_key_norm_numints;
198     uint64_t zn_hash;
199     matchtype_t zn_matchtype;
200     char zn_normbuf[ZAP_MAXNAMELEN];
201 } zap_name_t;

160 #define zap_f      zap_u.zap_fat
161 #define zap_m      zap_u.zap_micro

203 boolean_t zap_match(zap_name_t *zn, const char *matchname);
204 int zap_lockdir(objset_t *os, uint64_t obj, dmu_tx_t *tx,
205     krw_t lti, boolean_t fatreader, boolean_t adding, zap_t **zap);
206 void zap_unlockdir(zap_t *zap);
207 void zap_evict(dmu_buf_user_t *dbu);
208 void zap_evict(dmu_buf_t *db, void *vmzap);
209 zap_name_t *zap_name_alloc(zap_t *zap, const char *key, matchtype_t mt);
210 void zap_name_free(zap_name_t *zn);
211 int zap_hashbits(zap_t *zap);
212 uint32_t zap_maxcd(zap_t *zap);
213 uint64_t zap_getflags(zap_t *zap);

214 #define ZAP_HASH_IDX(hash, n) (((n) == 0) ? 0 : ((hash) >> (64 - (n))))

216 void fzap_byteswap(void *buf, size_t size);
217 int fzap_count(zap_t *zap, uint64_t *count);
218 int fzap_lookup(zap_name_t *zn,
219     uint64_t integer_size, uint64_t num_integers, void *buf,
220     char *realname, int rn_len, boolean_t *normalization_conflict);
221 void fzap_prefetch(zap_name_t *zn);
222 int fzap_count_write(zap_name_t *zn, int add, uint64_t *towrite,
223     uint64_t *tooverwrite);
224 int fzap_add(zap_name_t *zn, uint64_t integer_size, uint64_t num_integers,
225     const void *val, dmu_tx_t *tx);
226 int fzap_update(zap_name_t *zn,
227     int integer_size, uint64_t num_integers, const void *val, dmu_tx_t *tx);
228 int fzap_length(zap_name_t *zn,
229     uint64_t *integer_size, uint64_t *num_integers);
230 int fzap_remove(zap_name_t *zn, dmu_tx_t *tx);
231 int fzap_cursor_retrieve(zap_t *zap, zap_cursor_t *zc, zap_attribute_t *za);
232 void fzap_get_stats(zap_t *zap, zap_stats_t *zs);
233 void zap_put_leaf(struct zap_leaf *l);

235 int fzap_add_cd(zap_name_t *zn,
236     uint64_t integer_size, uint64_t num_integers,
237     const void *val, uint32_t cd, dmu_tx_t *tx);
238 void fzap_upgrade(zap_t *zap, dmu_tx_t *tx, zap_flags_t flags);
239 int fzap_cursor_move_to_key(zap_cursor_t *zc, zap_name_t *zn);

```

```

241 #ifdef __cplusplus
242 }
_____unchanged_portion_omitted_

```

```

*****
7867 Thu Apr 25 16:15:01 2013
new/usr/src/uts/common/fs/zfs/sys/zap_leaf.h
3752 want more verifiable dbuf user eviction
Submitted by: Justin Gibbs <justing@spectralogic.com>
Submitted by: Will Andrews <willa@spectralogic.com>
*****
unchanged_portion_omitted

130 typedef struct zap_leaf_dbuf {
131     uint8_t zldb_pad[offsetof(dmu_buf_t, db_data)];
132     zap_leaf_phys_t *zldb_data;
133 } zap_leaf_dbuf_t;

135 #endif /* ! codereview */
136 typedef union zap_leaf_chunk {
137     struct zap_leaf_entry {
138         uint8_t le_type;           /* always ZAP_CHUNK_ENTRY */
139         uint8_t le_value_intlen;  /* size of value's ints */
140         uint16_t le_next;         /* next entry in hash chain */
141         uint16_t le_name_chunk;   /* first chunk of the name */
142         uint16_t le_name_numints; /* ints in name (incl null) */
143         uint16_t le_value_chunk;  /* first chunk of the value */
144         uint16_t le_value_numints; /* value length in ints */
145         uint32_t le_cd;           /* collision differentiator */
146         uint64_t le_hash;        /* hash value of the name */
147     } l_entry;
148     struct zap_leaf_array {
149         uint8_t la_type;           /* always ZAP_CHUNK_ARRAY */
150         uint8_t la_array[ZAP_LEAF_ARRAY_BYTES];
151         uint16_t la_next;         /* next blk or CHAIN_END */
152     } l_array;
153     struct zap_leaf_free {
154         uint8_t lf_type;           /* always ZAP_CHUNK_FREE */
155         uint8_t lf_pad[ZAP_LEAF_ARRAY_BYTES];
156         uint16_t lf_next;        /* next in free list, or CHAIN_END */
157     } l_free;
158 } zap_leaf_chunk_t;

160 typedef struct zap_leaf {
161     dmu_buf_user_t db_evict;
162 #endif /* ! codereview */
163     krwlock_t l_rwlock;
164     uint64_t l_blkid;           /* 1<<ZAP_BLOCK_SHIFT byte block off */
165     int l_bs;                  /* block size shift */
166     union {
167         dmu_buf_t *l_dmu_db;
168         zap_leaf_dbuf_t *l_db;
169     } zl_db_u;
170     dmu_buf_t *l_dbuf;
171     zap_leaf_phys_t *l_phys;
172 } zap_leaf_t;

172 #define l_dbuf zl_db_u.l_dmu_db
173 #define l_phys zl_db_u.l_db->zldb_data
174 #endif /* ! codereview */

176 typedef struct zap_entry_handle {
177     /* below is set by zap_leaf.c and is public to zap.c */
178     uint64_t zeh_num_integers;
179     uint64_t zeh_hash;
180     uint32_t zeh_cd;
181     uint8_t zeh_integer_size;

183     /* below is private to zap_leaf.c */
184     uint16_t zeh_fakechunk;

```

```

185     uint16_t *zeh_chunkp;
186     zap_leaf_t *zeh_leaf;
187 } zap_entry_handle_t;

189 /*
190 * Return a handle to the named entry, or ENOENT if not found. The hash
191 * value must equal zap_hash(name).
192 */
193 extern int zap_leaf_lookup(zap_leaf_t *l,
194     struct zap_name *zn, zap_entry_handle_t *zeh);

196 /*
197 * Return a handle to the entry with this hash+cd, or the entry with the
198 * next closest hash+cd.
199 */
200 extern int zap_leaf_lookup_closest(zap_leaf_t *l,
201     uint64_t hash, uint32_t cd, zap_entry_handle_t *zeh);

203 /*
204 * Read the first num_integers in the attribute. Integer size
205 * conversion will be done without sign extension. Return EINVAL if
206 * integer_size is too small. Return EOVERFLOW if there are more than
207 * num_integers in the attribute.
208 */
209 extern int zap_entry_read(const zap_entry_handle_t *zeh,
210     uint8_t integer_size, uint64_t num_integers, void *buf);

212 extern int zap_entry_read_name(struct zap *zap, const zap_entry_handle_t *zeh,
213     uint16_t buflen, char *buf);

215 /*
216 * Replace the value of an existing entry.
217 *
218 * zap_entry_update may fail if it runs out of space (ENOSPC).
219 */
220 extern int zap_entry_update(zap_entry_handle_t *zeh,
221     uint8_t integer_size, uint64_t num_integers, const void *buf);

223 /*
224 * Remove an entry.
225 */
226 extern void zap_entry_remove(zap_entry_handle_t *zeh);

228 /*
229 * Create an entry. An equal entry must not exist, and this entry must
230 * belong in this leaf (according to its hash value). Fills in the
231 * entry handle on success. Returns 0 on success or ENOSPC on failure.
232 */
233 extern int zap_entry_create(zap_leaf_t *l, struct zap_name *zn, uint32_t cd,
234     uint8_t integer_size, uint64_t num_integers, const void *buf,
235     zap_entry_handle_t *zeh);

237 /*
238 * Return true if there are additional entries with the same normalized
239 * form.
240 */
241 extern boolean_t zap_entry_normalization_conflict(zap_entry_handle_t *zeh,
242     struct zap_name *zn, const char *name, struct zap *zap);

244 /*
245 * Other stuff.
246 */

248 extern void zap_leaf_init(zap_leaf_t *l, boolean_t sort);
249 extern void zap_leaf_byteswap(zap_leaf_phys_t *buf, int len);
250 extern void zap_leaf_split(zap_leaf_t *l, zap_leaf_t *nl, boolean_t sort);

```

```
251 extern void zap_leaf_stats(struct zap *zap, zap_leaf_t *l,  
252     struct zap_stats *zs);  
  
254 #ifdef __cplusplus  
255 }  
256 #endif  
  
258 #endif /* _SYS_ZAP_LEAF_H */
```

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

1

```
*****
2133 Thu Apr 25 16:15:01 2013
new/usr/src/uts/common/fs/zfs/sys/zfs_context.h
3752 want more verifiable dbuf user eviction
Submitted by: Justin Gibbs <justing@spectrallogic.com>
Submitted by: Will Andrews <willa@spectrallogic.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 2011 Nexenta Systems, Inc. All rights reserved.
28 * Copyright (c) 2012 by Delphix. All rights reserved.
29 */

31 #ifndef _SYS_ZFS_CONTEXT_H
32 #define _SYS_ZFS_CONTEXT_H

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

38 #ifndef _KERNEL
39 #include <stddef.h>
40 #endif
41 #endif /* ! codereview */
42 #include <sys/note.h>
43 #include <sys/types.h>
44 #include <sys/t_lock.h>
45 #include <sys/atomic.h>
46 #include <sys/sysmacros.h>
47 #include <sys/bitmap.h>
48 #include <sys/cmn_err.h>
49 #include <sys/kmem.h>
50 #include <sys/taskq.h>
51 #include <sys/taskq_impl.h>
52 #include <sys/buf.h>
53 #include <sys/param.h>
54 #include <sys/system.h>
55 #include <sys/cpuvar.h>
56 #include <sys/kobj.h>
57 #include <sys/conf.h>
58 #include <sys/disp.h>
59 #include <sys/debug.h>
```

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

2

```
60 #include <sys/random.h>
61 #include <sys/byteorder.h>
62 #include <sys/system.h>
63 #include <sys/list.h>
64 #include <sys/uiso.h>
65 #include <sys/dirent.h>
66 #include <sys/time.h>
67 #include <vm/seg_kmem.h>
68 #include <sys/zone.h>
69 #include <sys/uiso.h>
70 #include <sys/zfs_debug.h>
71 #include <sys/sysevent.h>
72 #include <sys/sysevent/eventdefs.h>
73 #include <sys/sysevent/dev.h>
74 #include <sys/fm/util.h>
75 #include <sys/sunddi.h>
76 #include <sys/cyclic.h>

78 #define CPU_SEQID      (CPU->cpu_seqid)

80 #ifdef __cplusplus
81 }
82 #endif

84 #endif /* _SYS_ZFS_CONTEXT_H */
```

```

*****
10120 Thu Apr 25 16:15:01 2013
new/usr/src/uts/common/fs/zfs/sys/zfs_ioctl.h
3752 want more verifiable dbuf user eviction
Submitted by: Justin Gibbs <justing@spectrallogic.com>
Submitted by: Will Andrews <willa@spectrallogic.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) 2012 by Delphix. All rights reserved.
24 */

26 #ifndef _SYS_ZFS_IOCTL_H
27 #define _SYS_ZFS_IOCTL_H

29 #ifndef _KERNEL
30 #include <stddef.h>
31 #endif
32 #endif /* ! codereview */
33 #include <sys/cred.h>
34 #include <sys/dmu.h>
35 #include <sys/zio.h>
36 #include <sys/dsl_deleg.h>
37 #include <sys/spa.h>
38 #include <sys/zfs_stat.h>

40 #ifdef _KERNEL
41 #include <sys/nvpair.h>
42 #endif /* _KERNEL */

44 #ifdef __cplusplus
45 extern "C" {
46 #endif

48 /*
49  * The structures in this file are passed between userland and the
50  * kernel. Userland may be running a 32-bit process, while the kernel
51  * is 64-bit. Therefore, these structures need to compile the same in
52  * 32-bit and 64-bit. This means not using type "long", and adding
53  * explicit padding so that the 32-bit structure will not be packed more
54  * tightly than the 64-bit structure (which requires 64-bit alignment).
55 */

57 /*
58  * Property values for snapdir
59 */

```

```

60 #define ZFS_SNAPDIR_HIDDEN 0
61 #define ZFS_SNAPDIR_VISIBLE 1

63 /*
64  * Field manipulation macros for the drr_versioninfo field of the
65  * send stream header.
66 */

68 /*
69  * Header types for zfs send streams.
70 */
71 typedef enum drr_headertype {
72     DMU_SUBSTREAM = 0x1,
73     DMU_COMPOUNDSTREAM = 0x2
74 } drr_headertype_t;

76 #define DMU_GET_STREAM_HDRTYPE(vi) BF64_GET((vi), 0, 2)
77 #define DMU_SET_STREAM_HDRTYPE(vi, x) BF64_SET((vi), 0, 2, x)

79 #define DMU_GET_FEATUREFLAGS(vi) BF64_GET((vi), 2, 30)
80 #define DMU_SET_FEATUREFLAGS(vi, x) BF64_SET((vi), 2, 30, x)

82 /*
83  * Feature flags for zfs send streams (flags in drr_versioninfo)
84 */

86 #define DMU_BACKUP_FEATURE_DEDUP (0x1)
87 #define DMU_BACKUP_FEATURE_DEDUPPROPS (0x2)
88 #define DMU_BACKUP_FEATURE_SA_SPILL (0x4)

90 /*
91  * Mask of all supported backup features
92 */
93 #define DMU_BACKUP_FEATURE_MASK (DMU_BACKUP_FEATURE_DEDUP | \
94     DMU_BACKUP_FEATURE_DEDUPPROPS | DMU_BACKUP_FEATURE_SA_SPILL)

96 /* Are all features in the given flag word currently supported? */
97 #define DMU_STREAM_SUPPORTED(x) (!(x) & ~DMU_BACKUP_FEATURE_MASK)

99 /*
100  * The drr_versioninfo field of the dmu_replay_record has the
101  * following layout:
102  *
103  *      64      56      48      40      32      24      16      8      0
104  * +-----+-----+-----+-----+-----+-----+-----+-----+
105  * |                                     | feature-flags |C|S|
106  * +-----+-----+-----+-----+-----+-----+-----+
107  *
108  * The low order two bits indicate the header type: SUBSTREAM (0x1)
109  * or COMPOUNDSTREAM (0x2). Using two bits for this is historical:
110  * this field used to be a version number, where the two version types
111  * were 1 and 2. Using two bits for this allows earlier versions of
112  * the code to be able to recognize send streams that don't use any
113  * of the features indicated by feature flags.
114 */

116 #define DMU_BACKUP_MAGIC 0x2F5bacbacULL

118 #define DRR_FLAG_CLONE (1<<0)
119 #define DRR_FLAG_CI_DATA (1<<1)

121 /*
122  * flags in the drr_checksumflags field in the DRR_WRITE and
123  * DRR_WRITE_BYREF blocks
124 */
125 #define DRR_CHECKSUM_DEDUP (1<<0)

```

```

127 #define DRR_IS_DEDUP_CAPABLE(flags)    ((flags) & DRR_CHECKSUM_DEDUP)
128
129 /*
130 * zfs ioctl command structure
131 */
132 typedef struct dmuf_replay_record {
133     enum {
134         DRR_BEGIN, DRR_OBJECT, DRR_FREEOBJECTS,
135         DRR_WRITE, DRR_FREE, DRR_END, DRR_WRITE_BYREF,
136         DRR_SPILL, DRR_NUMTYPES
137     } drr_type;
138     uint32_t drr_payloadlen;
139     union {
140         struct drr_begin {
141             uint64_t drr_magic;
142             uint64_t drr_versioninfo; /* was drr_version */
143             uint64_t drr_creation_time;
144             dmuf_objset_type_t drr_type;
145             uint32_t drr_flags;
146             uint64_t drr_toguid;
147             uint64_t drr_fromguid;
148             char drr_toname[MAXNAMELEN];
149         } drr_begin;
150         struct drr_end {
151             zio_cksum_t drr_checksum;
152             uint64_t drr_toguid;
153         } drr_end;
154         struct drr_object {
155             uint64_t drr_object;
156             dmuf_object_type_t drr_type;
157             dmuf_object_type_t drr_bonustype;
158             uint32_t drr_blkisz;
159             uint32_t drr_bonuslen;
160             uint8_t drr_checksumtype;
161             uint8_t drr_compress;
162             uint8_t drr_pad[6];
163             uint64_t drr_toguid;
164             /* bonus content follows */
165         } drr_object;
166         struct drr_freeobjects {
167             uint64_t drr_firstobj;
168             uint64_t drr_numobjs;
169             uint64_t drr_toguid;
170         } drr_freeobjects;
171         struct drr_write {
172             uint64_t drr_object;
173             dmuf_object_type_t drr_type;
174             uint32_t drr_pad;
175             uint64_t drr_offset;
176             uint64_t drr_length;
177             uint64_t drr_toguid;
178             uint8_t drr_checksumtype;
179             uint8_t drr_checksumflags;
180             uint8_t drr_pad2[6];
181             ddt_key_t drr_key; /* deduplication key */
182             /* content follows */
183         } drr_write;
184         struct drr_free {
185             uint64_t drr_object;
186             uint64_t drr_offset;
187             uint64_t drr_length;
188             uint64_t drr_toguid;
189         } drr_free;
190         struct drr_write_byref {
191             /* where to put the data */

```

```

192         uint64_t drr_object;
193         uint64_t drr_offset;
194         uint64_t drr_length;
195         uint64_t drr_toguid;
196         /* where to find the prior copy of the data */
197         uint64_t drr_refguid;
198         uint64_t drr_refobject;
199         uint64_t drr_refoffset;
200         /* properties of the data */
201         uint8_t drr_checksumtype;
202         uint8_t drr_checksumflags;
203         uint8_t drr_pad2[6];
204         ddt_key_t drr_key; /* deduplication key */
205     } drr_write_byref;
206     struct drr_spill {
207         uint64_t drr_object;
208         uint64_t drr_length;
209         uint64_t drr_toguid;
210         uint64_t drr_pad[4]; /* needed for crypto */
211         /* spill data follows */
212     } drr_spill;
213     } drr_u;
214 } dmuf_replay_record_t;
215
216 /* diff record range types */
217 typedef enum diff_type {
218     DDR_NONE = 0x1,
219     DDR_INUSE = 0x2,
220     DDR_FREE = 0x4
221 } diff_type_t;
222
223 /*
224 * The diff reports back ranges of free or in-use objects.
225 */
226 typedef struct dmuf_diff_record {
227     uint64_t ddr_type;
228     uint64_t ddr_first;
229     uint64_t ddr_last;
230 } dmuf_diff_record_t;
231
232 typedef struct zinject_record {
233     uint64_t zi_objset;
234     uint64_t zi_object;
235     uint64_t zi_start;
236     uint64_t zi_end;
237     uint64_t zi_guid;
238     uint32_t zi_level;
239     uint32_t zi_error;
240     uint64_t zi_type;
241     uint32_t zi_freq;
242     uint32_t zi_failfast;
243     char zi_func[MAXNAMELEN];
244     uint32_t zi_iotype;
245     int32_t zi_duration;
246     uint64_t zi_timer;
247     uint32_t zi_cmd;
248     uint32_t zi_pad;
249 } zinject_record_t;
250
251 #define ZINJECT_NULL          0x1
252 #define ZINJECT_FLUSH_ARC    0x2
253 #define ZINJECT_UNLOAD_SPA    0x4
254
255 typedef enum zinject_type {
256     ZINJECT_UNINITIALIZED,
257     ZINJECT_DATA_FAULT,

```

```

258     ZINJECT_DEVICE_FAULT,
259     ZINJECT_LABEL_FAULT,
260     ZINJECT_IGNORED_WRITES,
261     ZINJECT_PANIC,
262     ZINJECT_DELAY_IO,
263 } zinject_type_t;

265 typedef struct zfs_share {
266     uint64_t      z_exportdata;
267     uint64_t      z_sharedata;
268     uint64_t      z_sharetype; /* 0 = share, 1 = unshare */
269     uint64_t      z_sharemax; /* max length of share string */
270 } zfs_share_t;

272 /*
273  * ZFS file systems may behave the usual, POSIX-compliant way, where
274  * name lookups are case-sensitive. They may also be set up so that
275  * all the name lookups are case-insensitive, or so that only some
276  * lookups, the ones that set an FIGNORECASE flag, are case-insensitive.
277  */
278 typedef enum zfs_case {
279     ZFS_CASE_SENSITIVE,
280     ZFS_CASE_INSENSITIVE,
281     ZFS_CASE_MIXED
282 } zfs_case_t;

284 typedef struct zfs_cmd {
285     char          zc_name[MAXPATHLEN]; /* name of pool or dataset */
286     uint64_t      zc_nvlist_src; /* really (char *) */
287     uint64_t      zc_nvlist_src_size;
288     uint64_t      zc_nvlist_dst; /* really (char *) */
289     uint64_t      zc_nvlist_dst_size;
290     boolean_t     zc_nvlist_dst_filled; /* put an nvlist in dst? */
291     int           zc_pad2;

293     /*
294      * The following members are for legacy ioctls which haven't been
295      * converted to the new method.
296      */
297     uint64_t      zc_history; /* really (char *) */
298     char          zc_value[MAXPATHLEN * 2];
299     char          zc_string[MAXNAMELEN];
300     uint64_t      zc_guid;
301     uint64_t      zc_nvlist_conf; /* really (char *) */
302     uint64_t      zc_nvlist_conf_size;
303     uint64_t      zc_cookie;
304     uint64_t      zc_objset_type;
305     uint64_t      zc_perm_action;
306     uint64_t      zc_history_len;
307     uint64_t      zc_history_offset;
308     uint64_t      zc_obj;
309     uint64_t      zc_iflags; /* internal to zfs(7fs) */
310     zfs_share_t  zc_share;
311     dmu_objset_stats_t zc_objset_stats;
312     struct drr_begin zc_begin_record;
313     zinject_record_t zc_inject_record;
314     boolean_t     zc_defer_destroy;
315     boolean_t     zc_temphold;
316     uint64_t      zc_action_handle;
317     int           zc_cleanup_fd;
318     uint8_t       zc_pad[4]; /* alignment */
319     uint64_t      zc_sendobj;
320     uint64_t      zc_fromobj;
321     uint64_t      zc_createtxg;
322     zfs_stat_t    zc_stat;
323 } zfs_cmd_t;

```

```

325 typedef struct zfs_useracct {
326     char zu_domain[256];
327     uid_t zu_rid;
328     uint32_t zu_pad;
329     uint64_t zu_space;
330 } zfs_useracct_t;

332 #define ZFSDEV_MAX_MINOR      (1 << 16)
333 #define ZFS_MIN_MINOR        (ZFSDEV_MAX_MINOR + 1)

335 #define ZPOOL_EXPORT_AFTER_SPLIT 0x1

337 #ifndef _KERNEL

339 typedef struct zfs_creat {
340     nvlist_t      *zct_zplprops;
341     nvlist_t      *zct_props;
342 } zfs_creat_t;

344 extern dev_info_t *zfs_dip;

346 extern int zfs_secpolicy_snapshot_perms(const char *name, cred_t *cr);
347 extern int zfs_secpolicy_rename_perms(const char *from,
348     const char *to, cred_t *cr);
349 extern int zfs_secpolicy_destroy_perms(const char *name, cred_t *cr);
350 extern int zfs_busy(void);
351 extern void zfs_unmount_snap(const char *);
352 extern void zfs_destroy_unmount_origin(const char *);

354 /*
355  * ZFS minor numbers can refer to either a control device instance or
356  * a zvol. Depending on the value of zss_type, zss_data points to either
357  * a zvol_state_t or a zfs_onexit_t.
358  */
359 enum zfs_soft_state_type {
360     ZSST_ZVOL,
361     ZSST_CTLDEV
362 };

364 typedef struct zfs_soft_state {
365     enum zfs_soft_state_type zss_type;
366     void *zss_data;
367 } zfs_soft_state_t;

369 extern void *zfsdev_get_soft_state(minor_t minor,
370     enum zfs_soft_state_type which);
371 extern minor_t zfsdev_minor_alloc(void);

373 extern void *zfsdev_state;
374 extern kmutex_t zfsdev_state_lock;

376 #endif /* _KERNEL */

378 #ifdef __cplusplus
379 }
380 #endif

382 #endif /* _SYS_ZFS_IOCTL_H */

```

```

*****
33655 Thu Apr 25 16:15:02 2013
new/usr/src/uts/common/fs/zfs/zap.c
3752 want more verifiable dbuf user eviction
Submitted by: Justin Gibbs <justing@spectrallogic.com>
Submitted by: Will Andrews <willa@spectrallogic.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 */

26 /*
27 * This file contains the top half of the zfs directory structure
28 * implementation. The bottom half is in zap_leaf.c.
29 *
30 * The zdir is an extendable hash data structure. There is a table of
31 * pointers to buckets (zap_t->zdata->zleafs). The buckets are
32 * each a constant size and hold a variable number of directory entries.
33 * The buckets (aka "leaf nodes") are implemented in zap_leaf.c.
34 *
35 * The pointer table holds a power of 2 number of pointers.
36 * (1<<zap_t->zdata->zphys->zprefix_len). The bucket pointed to
37 * by the pointer at index i in the table holds entries whose hash value
38 * has a zprefix_len - bit prefix
39 */

41 #include <sys/spa.h>
42 #include <sys/dmu.h>
43 #include <sys/zfs_context.h>
44 #include <sys/zfs_znode.h>
45 #include <sys/fs/zfs.h>
46 #include <sys/zap.h>
47 #include <sys/refcount.h>
48 #include <sys/zap_impl.h>
49 #include <sys/zap_leaf.h>

51 int fzap_default_block_shift = 14; /* 16k blocksize */

53 static void zap_leaf_pageout(dmu_buf_t *db, void *v1);
53 static uint64_t zap_allocate_blocks(zap_t *zap, int nblocks);

55 void
56 fzap_byteswap(void *vbuf, size_t size)
57 {

```

```

58     uint64_t block_type;

60     block_type = *(uint64_t *)vbuf;

62     if (block_type == ZBT_LEAF || block_type == BSWAP_64(ZBT_LEAF))
63         zap_leaf_byteswap(vbuf, size);
64     else {
65         /* it's a ptrtbl block */
66         byteswap_uint64_array(vbuf, size);
67     }
68 }

70 void
71 fzap_upgrade(zap_t *zap, dmu_tx_t *tx, zap_flags_t flags)
72 {
73     dmu_buf_t *db;
74     zap_leaf_t *l;
75     int i;
76     zap_phys_t *zp;

78     ASSERT(RW_WRITE_HELD(&zap->zap_rwlock));
79     zap->zap_ismicro = FALSE;

81     zap->db_evict.evict_func = zap_evict;
83     (void) dmu_buf_update_user(zap->zap_dbuf, zap, zap,
84         &zap->zap_f.zap_phys, zap_evict);

83     mutex_init(&zap->zap_f.zap_num_entries_mtx, 0, 0, 0);
84     zap->zap_f.zap_block_shift = highbit(zap->zap_dbuf->db_size) - 1;

86     zp = zap->zap_f_phys;
89     zp = zap->zap_f.zap_phys;
87     /*
88      * explicitly zero it since it might be coming from an
89      * initialized microzap
90      */
91     bzero(zap->zap_dbuf->db_data, zap->zap_dbuf->db_size);
92     zap->zap_block_type = ZBT_HEADER;
93     zap->zap_magic = ZAP_MAGIC;

95     zap->zap_ptrtbl.zt_shift = ZAP_EMBEDDED_PTRTBL_SHIFT(zap);

97     zap->zap_freeblk = 2; /* block 1 will be the first leaf */
98     zap->zap_num_leafs = 1;
99     zap->zap_num_entries = 0;
100    zap->zap_salt = zap->zap_salt;
101    zap->zap_normflags = zap->zap_normflags;
102    zap->zap_flags = flags;

104    /* block 1 will be the first leaf */
105    for (i = 0; i < (1<<zap->zap_ptrtbl.zt_shift); i++)
106        ZAP_EMBEDDED_PTRTBL_ENT(zap, i) = 1;

108    /*
109     * set up block 1 - the first leaf
110     */
111    VERIFY(0 == dmu_buf_hold(zap->zap_objset, zap->zap_object,
112        1<<FZAP_BLOCK_SHIFT(zap), FTAG, &db, DMU_READ_NO_PREFETCH));
113    dmu_buf_will_dirty(db, tx);

115    l = kmem_zalloc(sizeof (zap_leaf_t), KM_SLEEP);
116    l->l_dbuf = db;
117    l->l_phys = db->db_data;

118    zap_leaf_init(l, zap->zap_normflags != 0);

```

```

120     kmem_free(l, sizeof (zap_leaf_t));
121     dmu_buf_rele(db, FTAG);
122 }
    unchanged_portion_omitted

314 static int
315 zap_grow_ptrtbl(zap_t *zap, dmu_tx_t *tx)
316 {
317     /*
318      * The pointer table should never use more hash bits than we
319      * have (otherwise we'd be using useless zero bits to index it).
320      * If we are within 2 bits of running out, stop growing, since
321      * this is already an aberrant condition.
322      */
323     if (zap->zap_f_phys->zap_ptrtbl.zt_shift >= zap_hashbits(zap) - 2)
324         if (zap->zap_f.zap_phys->zap_ptrtbl.zt_shift >= zap_hashbits(zap) - 2)
325             return (SET_ERROR(ENOSPC));

326     if (zap->zap_f_phys->zap_ptrtbl.zt_numblks == 0) {
327         if (zap->zap_f.zap_phys->zap_ptrtbl.zt_numblks == 0) {
328             /*
329              * We are outgrowing the "embedded" ptrtbl (the one
330              * stored in the header block). Give it its own entire
331              * block, which will double the size of the ptrtbl.
332              */
333             uint64_t newblk;
334             dmu_buf_t *db_new;
335             int err;

336             ASSERT3U(zap->zap_f_phys->zap_ptrtbl.zt_shift, ==,
337                 ASSERT3U(zap->zap_f.zap_phys->zap_ptrtbl.zt_shift, ==,
338                     ZAP_EMBEDDED_PTRTBL_SHIFT(zap)));
339             ASSERT0(zap->zap_f_phys->zap_ptrtbl.zt_blk);
340             ASSERT0(zap->zap_f.zap_phys->zap_ptrtbl.zt_blk);

341             newblk = zap_allocate_blocks(zap, 1);
342             err = dmu_buf_hold(zap->zap_objset, zap->zap_object,
343                 newblk << FZAP_BLOCK_SHIFT(zap), FTAG, &db_new,
344                 DMU_READ_NO_PREFETCH);
345             if (err)
346                 return (err);
347             dmu_buf_will_dirty(db_new, tx);
348             zap_ptrtbl_transfer(&ZAP_EMBEDDED_PTRTBL_ENT(zap, 0),
349                 db_new->db_data, 1 << ZAP_EMBEDDED_PTRTBL_SHIFT(zap));
350             dmu_buf_rele(db_new, FTAG);

351             zap->zap_f_phys->zap_ptrtbl.zt_blk = newblk;
352             zap->zap_f_phys->zap_ptrtbl.zt_numblks = 1;
353             zap->zap_f_phys->zap_ptrtbl.zt_shift++;
354             zap->zap_f.zap_phys->zap_ptrtbl.zt_blk = newblk;
355             zap->zap_f.zap_phys->zap_ptrtbl.zt_numblks = 1;
356             zap->zap_f.zap_phys->zap_ptrtbl.zt_shift++;

357             ASSERT3U(1ULL << zap->zap_f_phys->zap_ptrtbl.zt_shift, ==,
358                 zap->zap_f_phys->zap_ptrtbl.zt_numblks <<
359                 ASSERT3U(1ULL << zap->zap_f.zap_phys->zap_ptrtbl.zt_shift, ==,
360                     zap->zap_f.zap_phys->zap_ptrtbl.zt_numblks <<
361                     (FZAP_BLOCK_SHIFT(zap)-3)));

362             return (0);
363         } else {
364             return (zap_table_grow(zap, &zap->zap_f_phys->zap_ptrtbl,
365                 zap_table_grow(zap, &zap->zap_f.zap_phys->zap_ptrtbl,
366                     zap_ptrtbl_transfer, tx)));
367         }
368     }
369 }

```

```

366 static void
367 zap_increment_num_entries(zap_t *zap, int delta, dmu_tx_t *tx)
368 {
369     dmu_buf_will_dirty(zap->zap_dbuf, tx);
370     mutex_enter(&zap->zap_f.zap_num_entries_mtx);
371     ASSERT(delta > 0 || zap->zap_f_phys->zap_num_entries >= -delta);
372     zap->zap_f_phys->zap_num_entries += delta;
373     ASSERT(delta > 0 || zap->zap_f.zap_phys->zap_num_entries >= -delta);
374     zap->zap_f.zap_phys->zap_num_entries += delta;
375     mutex_exit(&zap->zap_f.zap_num_entries_mtx);
376 }

376 static uint64_t
377 zap_allocate_blocks(zap_t *zap, int nblocks)
378 {
379     uint64_t newblk;
380     ASSERT(RW_WRITE_HELD(&zap->zap_rwlock));
381     newblk = zap->zap_f_phys->zap_freeblk;
382     zap->zap_f_phys->zap_freeblk += nblocks;
383     newblk = zap->zap_f.zap_phys->zap_freeblk;
384     zap->zap_f.zap_phys->zap_freeblk += nblocks;
385     return (newblk);
386 }

386 static void
387 zap_leaf_pageout(dmu_buf_user_t *dbu)
388 {
389     zap_leaf_t *l = (zap_leaf_t *)dbu;

390     rw_destroy(&l->l_rwlock);
391     kmem_free(l, sizeof (zap_leaf_t));
392 }

395 #endif /* ! codereview */
396 static zap_leaf_t *
397 zap_create_leaf(zap_t *zap, dmu_tx_t *tx)
398 {
399     void *winner;
400     zap_leaf_t *l = kmem_alloc(sizeof (zap_leaf_t), KM_SLEEP);

401     ASSERT(RW_WRITE_HELD(&zap->zap_rwlock));

402     rw_init(&l->l_rwlock, 0, 0, 0);
403     rw_enter(&l->l_rwlock, RW_WRITER);
404     l->l_blkid = zap_allocate_blocks(zap, 1);
405     l->l_dbuf = NULL;
406     l->l_phys = NULL;

407     VERIFY(0 == dmu_buf_hold(zap->zap_objset, zap->zap_object,
408         l->l_blkid << FZAP_BLOCK_SHIFT(zap), NULL, &l->l_dbuf,
409         DMU_READ_NO_PREFETCH));
410     dmu_buf_init_user(&l->db_evict, zap_leaf_pageout);
411     winner = (zap_leaf_t *)dmu_buf_set_user(l->l_dbuf, &l->db_evict);
412     winner = dmu_buf_set_user(l->l_dbuf, 1, &l->l_phys, zap_leaf_pageout);
413     ASSERT(winner == NULL);
414     dmu_buf_will_dirty(l->l_dbuf, tx);

415     zap_leaf_init(l, zap->zap_normflags != 0);

416     zap->zap_f_phys->zap_num_leafs++;
417     zap->zap_f.zap_phys->zap_num_leafs++;

418     return (l);
419 }

```

```

424 int
425 fzap_count(zap_t *zap, uint64_t *count)
426 {
427     ASSERT(!zap->zap_ismicro);
428     mutex_enter(&zap->zap_f.zap_num_entries_mtx); /* unnecessary */
429     *count = zap->zap_f.phys->zap_num_entries;
430     *count = zap->zap_f.phys->zap_num_entries;
431     mutex_exit(&zap->zap_f.zap_num_entries_mtx);
432     return (0);
433 }
434
435 unchanged_portion_omitted
436
437 _NOTE(ARGUNUSED(0))
438 static void
439 zap_leaf_pageout(dmu_buf_t *db, void *vl)
440 {
441     zap_leaf_t *l = vl;
442
443     rw_destroy(&l->l_rwlock);
444     kmem_free(l, sizeof (zap_leaf_t));
445 }
446
447 static zap_leaf_t *
448 zap_open_leaf(uint64_t blkid, dmu_buf_t *db)
449 {
450     zap_leaf_t *l, *winner;
451
452     ASSERT(blkid != 0);
453
454     l = kmem_alloc(sizeof (zap_leaf_t), KM_SLEEP);
455     rw_init(&l->l_rwlock, 0, 0, 0);
456     rw_enter(&l->l_rwlock, RW_WRITER);
457     l->l_blkid = blkid;
458     l->l_bs = highbit(db->db_size)-1;
459     l->l_dbuf = db;
460     l->l_phys = NULL;
461
462     dmu_buf_init_user(&l->db_evict, zap_leaf_pageout);
463     winner = (zap_leaf_t *)dmu_buf_set_user(db, &l->db_evict);
464     winner = dmu_buf_set_user(db, l, &l->l_phys, zap_leaf_pageout);
465
466     rw_exit(&l->l_rwlock);
467     if (winner != NULL) {
468         /* someone else set it first */
469         zap_leaf_pageout(&l->db_evict);
470         zap_leaf_pageout(NULL, l);
471         l = winner;
472     }
473
474     /*
475     * lhr_pad was previously used for the next leaf in the leaf
476     * chain. There should be no chained leafs (as we have removed
477     * support for them).
478     */
479     ASSERT0(l->l_phys->l_hdr.lh_pad1);
480
481     /*
482     * There should be more hash entries than there can be
483     * chunks to put in the hash table
484     */
485     ASSERT3U(ZAP_LEAF_HASH_NUMENTRIES(1), >, ZAP_LEAF_NUMCHUNKS(1) / 3);
486
487     /* The chunks should begin at the end of the hash table */
488     ASSERT3P(&ZAP_LEAF_CHUNK(1, 0), ==,
489             &l->l_phys->l_hash[ZAP_LEAF_HASH_NUMENTRIES(1)]);

```

```

486     /* The chunks should end at the end of the block */
487     ASSERT3U((uintptr_t)&ZAP_LEAF_CHUNK(1, ZAP_LEAF_NUMCHUNKS(1)) -
488             (uintptr_t)l->l_phys, ==, l->l_dbuf->db_size);
489
490     return (l);
491 }
492
493 static int
494 zap_get_leaf_byblk(zap_t *zap, uint64_t blkid, dmu_tx_t *tx, krw_t lt,
495     zap_leaf_t **lp)
496 {
497     dmu_buf_t *db;
498     zap_leaf_t *l;
499     int bs = FZAP_BLOCK_SHIFT(zap);
500     int err;
501
502     ASSERT(RW_LOCK_HELD(&zap->zap_rwlock));
503
504     err = dmu_buf_hold(zap->zap_objset, zap->zap_object,
505         blkid << bs, NULL, &db, DMU_READ_NO_PREFETCH);
506     if (err)
507         return (err);
508
509     ASSERT3U(db->db_object, ==, zap->zap_object);
510     ASSERT3U(db->db_offset, ==, blkid << bs);
511     ASSERT3U(db->db_size, ==, 1 << bs);
512     ASSERT(blkid != 0);
513
514     l = (zap_leaf_t *)dmu_buf_get_user(db);
515     l = dmu_buf_get_user(db);
516
517     if (l == NULL)
518         l = zap_open_leaf(blkid, db);
519
520     rw_enter(&l->l_rwlock, lt);
521     /*
522     * Must lock before dirtying, otherwise l->l_phys could change,
523     * causing ASSERT below to fail.
524     */
525     if (lt == RW_WRITER)
526         dmu_buf_will_dirty(db, tx);
527     ASSERT3U(l->l_blkid, ==, blkid);
528     ASSERT3P(l->l_dbuf, ==, db);
529     ASSERT3P(l->l_phys, ==, l->l_dbuf->db_data);
530     ASSERT3U(l->l_phys->l_hdr.lh_block_type, ==, ZBT_LEAF);
531     ASSERT3U(l->l_phys->l_hdr.lh_magic, ==, ZAP_LEAF_MAGIC);
532
533     *lp = l;
534     return (0);
535 }
536
537 static int
538 zap_idx_to_blk(zap_t *zap, uint64_t idx, uint64_t *valp)
539 {
540     ASSERT(RW_LOCK_HELD(&zap->zap_rwlock));
541
542     if (zap->zap_f.phys->zap_ptrtbl.zt_numblks == 0) {
543         if (zap->zap_f.zap_phys->zap_ptrtbl.zt_numblks == 0) {
544             ASSERT3U(idx, <,
545                 (1ULL << zap->zap_f.phys->zap_ptrtbl.zt_shift));
546             (1ULL << zap->zap_f.zap_phys->zap_ptrtbl.zt_shift);
547             *valp = ZAP_EMBEDDED_PTRTBL_ENT(zap, idx);
548             return (0);
549         } else {
550             return (zap_table_load(zap, &zap->zap_f.phys->zap_ptrtbl,
551                 zap_table_load(zap, &zap->zap_f.zap_phys->zap_ptrtbl,

```

```

548         idx, valp));
549     }
550 }

552 static int
553 zap_set_idx_to_blk(zap_t *zap, uint64_t idx, uint64_t blk, dmu_tx_t *tx)
554 {
555     ASSERT(tx != NULL);
556     ASSERT(RW_WRITE_HELD(&zap->zap_rwlock));

558     if (zap->zap_f_phys->zap_ptrtbl.zt_blk == 0) {
559         if (zap->zap_f.zap_phys->zap_ptrtbl.zt_blk == 0) {
560             ZAP_EMBEDDED_PTRTBL_ENT(zap, idx) = blk;
561             return (0);
562         } else {
563             return (zap_table_store(zap, &zap->zap_f_phys->zap_ptrtbl,
564                 zap_table_store(zap, &zap->zap_f.zap_phys->zap_ptrtbl,
565                     idx, blk, tx));
566         }
567     }

568     static int
569     zap_deref_leaf(zap_t *zap, uint64_t h, dmu_tx_t *tx, krw_t lt, zap_leaf_t **lp)
570     {
571         uint64_t idx, blk;
572         int err;

573         ASSERT(zap->zap_dbuf == NULL ||
574             zap->zap_f_phys == zap->zap_dbuf->db_data);
575         ASSERT3U(zap->zap_f_phys->zap_magic, ==, ZAP_MAGIC);
576         idx = ZAP_HASH_IDX(h, zap->zap_f_phys->zap_ptrtbl.zt_shift);
577         zap->zap_f.zap_phys == zap->zap_dbuf->db_data);
578         ASSERT3U(zap->zap_f.zap_phys->zap_magic, ==, ZAP_MAGIC);
579         idx = ZAP_HASH_IDX(h, zap->zap_f.zap_phys->zap_ptrtbl.zt_shift);
580         err = zap_idx_to_blk(zap, idx, &blk);
581         if (err != 0)
582             return (err);
583         err = zap_get_leaf_byblk(zap, blk, tx, lt, lp);

584         ASSERT(err || ZAP_HASH_IDX(h, (*lp)->l_phys->l_hdr.lh_prefix_len) ==
585             (*lp)->l_phys->l_hdr.lh_prefix);
586         return (err);
587     }

588     static int
589     zap_expand_leaf(zap_name_t *zn, zap_leaf_t *l, dmu_tx_t *tx, zap_leaf_t **lp)
590     {
591         zap_t *zap = zn->zn_zap;
592         uint64_t hash = zn->zn_hash;
593         zap_leaf_t *nl;
594         int prefix_diff, i, err;
595         uint64_t sibling;
596         int old_prefix_len = l->l_phys->l_hdr.lh_prefix_len;

597         ASSERT3U(old_prefix_len, <=, zap->zap_f_phys->zap_ptrtbl.zt_shift);
598         ASSERT3U(old_prefix_len, <=, zap->zap_f.zap_phys->zap_ptrtbl.zt_shift);
599         ASSERT(RW_LOCK_HELD(&zap->zap_rwlock));

600         ASSERT3U(ZAP_HASH_IDX(hash, old_prefix_len), ==,
601             l->l_phys->l_hdr.lh_prefix);

602         if (zap_tryupgradedir(zap, tx) == 0 ||
603             old_prefix_len == zap->zap_f_phys->zap_ptrtbl.zt_shift) {
604             old_prefix_len == zap->zap_f.zap_phys->zap_ptrtbl.zt_shift) {
605                 /* We failed to upgrade, or need to grow the pointer table */
606                 objset_t *os = zap->zap_objset;

```

```

607         uint64_t object = zap->zap_object;

608         zap_put_leaf(l);
609         zap_unlockdir(zap);
610         err = zap_lockdir(os, object, tx, RW_WRITER,
611             FALSE, FALSE, &zn->zn_zap);
612         zap = zn->zn_zap;
613         if (err)
614             return (err);
615         ASSERT(!zap->zap_ismicro);

616         while (old_prefix_len ==
617             zap->zap_f_phys->zap_ptrtbl.zt_shift) {
618             zap->zap_f.zap_phys->zap_ptrtbl.zt_shift) {
619                 err = zap_grow_ptrtbl(zap, tx);
620                 if (err)
621                     return (err);
622             }

623         err = zap_deref_leaf(zap, hash, tx, RW_WRITER, &l);
624         if (err)
625             return (err);

626         if (l->l_phys->l_hdr.lh_prefix_len != old_prefix_len) {
627             /* it split while our locks were down */
628             *lp = l;
629             return (0);
630         }
631         ASSERT(RW_WRITE_HELD(&zap->zap_rwlock));
632         ASSERT3U(old_prefix_len, <, zap->zap_f_phys->zap_ptrtbl.zt_shift);
633         ASSERT3U(old_prefix_len, <, zap->zap_f.zap_phys->zap_ptrtbl.zt_shift);
634         ASSERT3U(ZAP_HASH_IDX(hash, old_prefix_len), ==,
635             l->l_phys->l_hdr.lh_prefix);

636         prefix_diff = zap->zap_f_phys->zap_ptrtbl.zt_shift -
637             prefix_diff = zap->zap_f.zap_phys->zap_ptrtbl.zt_shift -
638             (old_prefix_len + 1);
639         sibling = (ZAP_HASH_IDX(hash, old_prefix_len + 1) | 1) << prefix_diff;

640         /* check for i/o errors before doing zap_leaf_split */
641         for (i = 0; i < (1ULL << prefix_diff); i++) {
642             uint64_t blk;
643             err = zap_idx_to_blk(zap, sibling+i, &blk);
644             if (err)
645                 return (err);
646             ASSERT3U(blk, ==, l->l_blkid);

647         nl = zap_create_leaf(zap, tx);
648         zap_leaf_split(l, nl, zap->zap_normflags != 0);

649         /* set sibling pointers */
650         for (i = 0; i < (1ULL << prefix_diff); i++) {
651             err = zap_set_idx_to_blk(zap, sibling+i, nl->l_blkid, tx);
652             ASSERT0(err); /* we checked for i/o errors above */
653         }

654         if (hash & (1ULL << (64 - l->l_phys->l_hdr.lh_prefix_len))) {
655             /* we want the sibling */
656             zap_put_leaf(l);
657             *lp = nl;
658         } else {
659             zap_put_leaf(nl);
660             *lp = l;
661         }

```

```

671     return (0);
672 }

674 static void
675 zap_put_leaf_maybe_grow_ptrtbl(zap_name_t *zn, zap_leaf_t *l, dmu_tx_t *tx)
676 {
677     zap_t *zap = zn->zn_zap;
678     int shift = zap->zap_f_phys->zap_ptrtbl.zt_shift;
679     int shift = zap->zap_f.zap_phys->zap_ptrtbl.zt_shift;
680     int leaffull = (l->l_phys->l_hdr.lh_prefix_len == shift &&
681         l->l_phys->l_hdr.lh_nfree < ZAP_LEAF_LOW_WATER);

682     zap_put_leaf(l);

684     if (leaffull || zap->zap_f_phys->zap_ptrtbl.zt_nextblk) {
685         if (leaffull || zap->zap_f.zap_phys->zap_ptrtbl.zt_nextblk) {
686             int err;

687             /*
688              * We are in the middle of growing the pointer table, or
689              * this leaf will soon make us grow it.
690              */
691             if (zap_tryupgradedir(zap, tx) == 0) {
692                 objset_t *os = zap->zap_objset;
693                 uint64_t zapobj = zap->zap_object;

694                 zap_unlockdir(zap);
695                 err = zap_lockdir(os, zapobj, tx,
696                     RW_WRITER, FALSE, FALSE, &zn->zn_zap);
697                 zap = zn->zn_zap;
698                 if (err)
699                     return;
700             }
701         }

702         /* could have finished growing while our locks were down */
703         if (zap->zap_f_phys->zap_ptrtbl.zt_shift == shift)
704             if (zap->zap_f.zap_phys->zap_ptrtbl.zt_shift == shift)
705                 (void) zap_grow_ptrtbl(zap, tx);
706     }
707 }

unchanged_portion_omitted

927 void
928 fzap_prefetch(zap_name_t *zn)
929 {
930     uint64_t idx, blk;
931     zap_t *zap = zn->zn_zap;
932     int bs;

933     idx = ZAP_HASH_IDX(zn->zn_hash,
934         zap->zap_f_phys->zap_ptrtbl.zt_shift);
935     zap->zap_f.zap_phys->zap_ptrtbl.zt_shift);
936     if (zap_idx_to_blk(zap, idx, &blk) != 0)
937         return;
938     bs = FZAP_BLOCK_SHIFT(zap);
939     dmu_prefetch(zap->zap_objset, zap->zap_object, blk << bs, 1 << bs);
940 }

unchanged_portion_omitted

1266 void
1267 fzap_get_stats(zap_t *zap, zap_stats_t *zs)
1268 {
1269     int bs = FZAP_BLOCK_SHIFT(zap);
1270     zs->zs_blocksize = 1ULL << bs;

```

```

1272     /*
1273      * Set zap_phys_t fields
1274      */
1275     zs->zs_num_leafs = zap->zap_f_phys->zap_num_leafs;
1276     zs->zs_num_entries = zap->zap_f_phys->zap_num_entries;
1277     zs->zs_num_blocks = zap->zap_f_phys->zap_freeblk;
1278     zs->zs_block_type = zap->zap_f_phys->zap_block_type;
1279     zs->zs_magic = zap->zap_f_phys->zap_magic;
1280     zs->zs_salt = zap->zap_f_phys->zap_salt;
1281     zs->zs_num_leafs = zap->zap_f.zap_phys->zap_num_leafs;
1282     zs->zs_num_entries = zap->zap_f.zap_phys->zap_num_entries;
1283     zs->zs_num_blocks = zap->zap_f.zap_phys->zap_freeblk;
1284     zs->zs_block_type = zap->zap_f.zap_phys->zap_block_type;
1285     zs->zs_magic = zap->zap_f.zap_phys->zap_magic;
1286     zs->zs_salt = zap->zap_f.zap_phys->zap_salt;

1287     /*
1288      * Set zap_ptrtbl fields
1289      */
1290     zs->zs_ptrtbl_len = 1ULL << zap->zap_f_phys->zap_ptrtbl.zt_shift;
1291     zs->zs_ptrtbl_nextblk = zap->zap_f_phys->zap_ptrtbl.zt_nextblk;
1292     zs->zs_ptrtbl_len = 1ULL << zap->zap_f.zap_phys->zap_ptrtbl.zt_shift;
1293     zs->zs_ptrtbl_nextblk = zap->zap_f.zap_phys->zap_ptrtbl.zt_nextblk;
1294     zs->zs_ptrtbl_blks_copied =
1295         zap->zap_f_phys->zap_ptrtbl.zt_blks_copied;
1296     zs->zs_ptrtbl_zt_blk = zap->zap_f_phys->zap_ptrtbl.zt_blk;
1297     zs->zs_ptrtbl_zt_numblks = zap->zap_f_phys->zap_ptrtbl.zt_numblks;
1298     zs->zs_ptrtbl_zt_shift = zap->zap_f_phys->zap_ptrtbl.zt_shift;
1299     zap->zap_f.zap_phys->zap_ptrtbl.zt_blks_copied;
1300     zs->zs_ptrtbl_zt_blk = zap->zap_f.zap_phys->zap_ptrtbl.zt_blk;
1301     zs->zs_ptrtbl_zt_numblks = zap->zap_f.zap_phys->zap_ptrtbl.zt_numblks;
1302     zs->zs_ptrtbl_zt_shift = zap->zap_f.zap_phys->zap_ptrtbl.zt_shift;

1303     if (zap->zap_f_phys->zap_ptrtbl.zt_numblks == 0) {
1304         if (zap->zap_f.zap_phys->zap_ptrtbl.zt_numblks == 0) {
1305             /* the ptrtbl is entirely in the header block. */
1306             zap_stats_ptrtbl(zap, &ZAP_EMBEDDED_PTRTBL_ENT(zap, 0),
1307                 1 << ZAP_EMBEDDED_PTRTBL_SHIFT(zap), zs);
1308         } else {
1309             int b;

1310             dmu_prefetch(zap->zap_objset, zap->zap_object,
1311                 zap->zap_f_phys->zap_ptrtbl.zt_blk << bs,
1312                 zap->zap_f_phys->zap_ptrtbl.zt_numblks << bs);
1313             zap->zap_f.zap_phys->zap_ptrtbl.zt_numblks << bs);

1314             for (b = 0; b < zap->zap_f_phys->zap_ptrtbl.zt_numblks;
1315                 for (b = 0; b < zap->zap_f.zap_phys->zap_ptrtbl.zt_numblks;
1316                     b++) {
1317                 dmu_buf_t *db;
1318                 int err;

1319                 err = dmu_buf_hold(zap->zap_objset, zap->zap_object,
1320                     (zap->zap_f_phys->zap_ptrtbl.zt_blk + b) << bs,
1321                     (zap->zap_f.zap_phys->zap_ptrtbl.zt_blk + b) << bs,
1322                     FTAG, &db, DMU_READ_NO_PREFETCH);
1323                 if (err == 0) {
1324                     zap_stats_ptrtbl(zap, db->db_data,
1325                         1 << (bs-3), zs);
1326                     dmu_buf_rele(db, FTAG);
1327                 }
1328             }
1329         }
1330     }
1331 }
1332 }
1333 }
1334 }
1335 }
1336 }
1337 }
1338 }
1339 }

```

```

1321 int
1322 fzap_count_write(zap_name_t *zn, int add, uint64_t *towrite,
1323                uint64_t *tooverwrite)
1324 {
1325     zap_t *zap = zn->zn_zap;
1326     zap_leaf_t *l;
1327     int err;
1328
1329     /*
1330      * Account for the header block of the fatzap.
1331      */
1332     if (!add && dmu_buf_freeable(zap->zap_dbuf)) {
1333         *tooverwrite += zap->zap_dbuf->db_size;
1334     } else {
1335         *towrite += zap->zap_dbuf->db_size;
1336     }
1337
1338     /*
1339      * Account for the pointer table blocks.
1340      * If we are adding we need to account for the following cases :
1341      * - If the pointer table is embedded, this operation could force an
1342      *   external pointer table.
1343      * - If this already has an external pointer table this operation
1344      *   could extend the table.
1345      */
1346     if (add) {
1347         if (zap->zap_f_phys->zap_ptrtbl.zt_blk == 0)
1348             if (zap->zap_f.zap_phys->zap_ptrtbl.zt_blk == 0)
1349                 *towrite += zap->zap_dbuf->db_size;
1350         else
1351             *towrite += (zap->zap_dbuf->db_size * 3);
1352     }
1353
1354     /*
1355      * Now, check if the block containing leaf is freeable
1356      * and account accordingly.
1357      */
1358     err = zap_deref_leaf(zap, zn->zn_hash, NULL, RW_READER, &l);
1359     if (err != 0) {
1360         return (err);
1361     }
1362
1363     if (!add && dmu_buf_freeable(l->l_dbuf)) {
1364         *tooverwrite += l->l_dbuf->db_size;
1365     } else {
1366         /*
1367          * If this an add operation, the leaf block could split.
1368          * Hence, we need to account for an additional leaf block.
1369          */
1370         *towrite += (add ? 2 : 1) * l->l_dbuf->db_size;
1371     }
1372
1373     zap_put_leaf(l);
1374     return (0);
1375 }

```

_____unchanged_portion_omitted_____

```

*****
23138 Thu Apr 25 16:15:02 2013
new/usr/src/uts/common/fs/zfs/zap_leaf.c
3752 want more verifiable dbuf user eviction
Submitted by: Justin Gibbs <justing@spectralogic.com>
Submitted by: Will Andrews <willa@spectralogic.com>
*****
unchanged_portion_omitted_

```

```

103 void
104 zap_leaf_byteswap(zap_leaf_phys_t *buf, int size)
105 {
106     int i;
107     zap_leaf_t l = { 0 };
108
109     zap_leaf_t l;
110     l.l_bs = highbit(size)-1;
111     l.l_phys = buf;
112
113     buf->l_hdr.lh_block_type = BSWAP_64(buf->l_hdr.lh_block_type);
114     buf->l_hdr.lh_prefix = BSWAP_64(buf->l_hdr.lh_prefix);
115     buf->l_hdr.lh_magic = BSWAP_32(buf->l_hdr.lh_magic);
116     buf->l_hdr.lh_nfree = BSWAP_16(buf->l_hdr.lh_nfree);
117     buf->l_hdr.lh_nentries = BSWAP_16(buf->l_hdr.lh_nentries);
118     buf->l_hdr.lh_prefix_len = BSWAP_16(buf->l_hdr.lh_prefix_len);
119     buf->l_hdr.lh_freelist = BSWAP_16(buf->l_hdr.lh_freelist);
120
121     for (i = 0; i < ZAP_LEAF_HASH_NUMENTRIES(&l); i++)
122         buf->l_hash[i] = BSWAP_16(buf->l_hash[i]);
123
124     for (i = 0; i < ZAP_LEAF_NUMCHUNKS(&l); i++) {
125         zap_leaf_chunk_t *lc = &ZAP_LEAF_CHUNK(&l, i);
126         struct zap_leaf_entry *le;
127
128         switch (lc->l_free.lf_type) {
129             case ZAP_CHUNK_ENTRY:
130                 le = &lc->l_entry;
131                 le->le_type = BSWAP_8(le->le_type);
132                 le->le_value_intlen = BSWAP_8(le->le_value_intlen);
133                 le->le_next = BSWAP_16(le->le_next);
134                 le->le_name_chunk = BSWAP_16(le->le_name_chunk);
135                 le->le_name_numints = BSWAP_16(le->le_name_numints);
136                 le->le_value_chunk = BSWAP_16(le->le_value_chunk);
137                 le->le_value_numints = BSWAP_16(le->le_value_numints);
138                 le->le_cd = BSWAP_32(le->le_cd);
139                 le->le_hash = BSWAP_64(le->le_hash);
140                 break;
141             case ZAP_CHUNK_FREE:
142                 lc->l_free.lf_type = BSWAP_8(lc->l_free.lf_type);
143                 lc->l_free.lf_next = BSWAP_16(lc->l_free.lf_next);
144                 break;
145             case ZAP_CHUNK_ARRAY:
146                 lc->l_array.la_type = BSWAP_8(lc->l_array.la_type);
147                 lc->l_array.la_next = BSWAP_16(lc->l_array.la_next);
148                 /* la_array doesn't need swapping */
149                 break;
150             default:
151                 ASSERT(!"bad leaf type");
152         }
153     }
154
155     unchanged_portion_omitted_

```

```

830 void
831 zap_leaf_stats(zap_t *zap, zap_leaf_t *l, zap_stats_t *zs)

```

```

832 {
833     int i, n;
834
835     n = zap->zap_f_phys->zap_ptrtbl.zt_shift -
836     n = zap->zap_f.zap_phys->zap_ptrtbl.zt_shift -
837     l->l_phys->l_hdr.lh_prefix_len;
838     n = MIN(n, ZAP_HISTOGRAM_SIZE-1);
839     zs->zs_leafs_with_2n_pointers[n]++;
840
841     n = l->l_phys->l_hdr.lh_nentries/5;
842     n = MIN(n, ZAP_HISTOGRAM_SIZE-1);
843     zs->zs_blocks_with_n5_entries[n]++;
844
845     n = ((1<<FZAP_BLOCK_SHIFT(zap)) -
846     l->l_phys->l_hdr.lh_nfree * (ZAP_LEAF_ARRAY_BYTES+1))*10 /
847     (1<<FZAP_BLOCK_SHIFT(zap));
848     n = MIN(n, ZAP_HISTOGRAM_SIZE-1);
849     zs->zs_blocks_n_tenths_full[n]++;
850
851     for (i = 0; i < ZAP_LEAF_HASH_NUMENTRIES(l); i++) {
852         int nentries = 0;
853         int chunk = l->l_phys->l_hash[i];
854
855         while (chunk != CHAIN_END) {
856             struct zap_leaf_entry *le =
857                 ZAP_LEAF_ENTRY(l, chunk);
858
859             n = 1 + ZAP_LEAF_ARRAY_NCHUNKS(le->le_name_numints) +
860                 ZAP_LEAF_ARRAY_NCHUNKS(le->le_value_numints *
861                 le->le_value_intlen);
862             n = MIN(n, ZAP_HISTOGRAM_SIZE-1);
863             zs->zs_entries_using_n_chunks[n]++;
864
865             chunk = le->le_next;
866             nentries++;
867         }
868
869         n = nentries;
870         n = MIN(n, ZAP_HISTOGRAM_SIZE-1);
871         zs->zs_buckets_with_n_entries[n]++;
872     }
873 }

```

unchanged_portion_omitted_

```

*****
34972 Thu Apr 25 16:15:02 2013
new/usr/src/uts/common/fs/zfs/zap_micro.c
3752 want more verifiable dbuf user eviction
Submitted by: Justin Gibbs <justing@spectrallogic.com>
Submitted by: Will Andrews <willa@spectrallogic.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 */

26 #include <sys/zio.h>
27 #include <sys/spa.h>
28 #include <sys/dmu.h>
29 #include <sys/zfs_context.h>
30 #include <sys/zap.h>
31 #include <sys/refcount.h>
32 #include <sys/zap_impl.h>
33 #include <sys/zap_leaf.h>
34 #include <sys/avl.h>
35 #include <sys/arc.h>

37 #ifdef _KERNEL
38 #include <sys/sunddi.h>
39 #endif

41 static int mzap_upgrade(zap_t **zapp, dmu_tx_t *tx, zap_flags_t flags);

43 uint64_t
44 zap_getflags(zap_t *zap)
45 {
46     if (zap->zap_ismicro)
47         return (0);
48     return (zap->zap_f_phys->zap_flags);
49     return (zap->zap_u.zap_fat.zap_phys->zap_flags);
49 }
unchanged_portion_omitted

359 static zap_t *
360 mzap_open(objset_t *os, uint64_t obj, dmu_buf_t *db)
361 {
362     zap_t *winner;
363     zap_t *zap;
364     int i;

```

```

366     ASSERT3U(MZAP_ENT_LEN, ==, sizeof (mzap_ent_phys_t));

368     zap = kmem_zalloc(sizeof (zap_t), KM_SLEEP);
369     rw_init(&zap->zap_rwlock, 0, 0, 0);
370     rw_enter(&zap->zap_rwlock, RW_WRITER);
371     zap->zap_objset = os;
372     zap->zap_object = obj;
373     zap->zap_dbuf = db;

375     if (*(uint64_t *)db->db_data != ZBT_MICRO) {
376         mutex_init(&zap->zap_f.zap_num_entries_mtx, 0, 0, 0);
377         zap->zap_f.zap_block_shift = highbit(db->db_size) - 1;
378     } else {
379         zap->zap_ismicro = TRUE;
380     }

382     /*
383     * Make sure that zap_ismicro is set before we let others see
384     * it, because zap_lockdir() checks zap_ismicro without the lock
385     * held.
386     */
387     dmu_buf_init_user(&zap->db_evict, zap_evict);
388     winner = (zap_t *)dmu_buf_set_user(db, &zap->db_evict);
389     winner = dmu_buf_set_user(db, zap, &zap->zap_m.zap_phys, zap_evict);

390     if (winner != NULL) {
391         rw_exit(&zap->zap_rwlock);
392         rw_destroy(&zap->zap_rwlock);
393         if (!zap->zap_ismicro)
394             mutex_destroy(&zap->zap_f.zap_num_entries_mtx);
395         kmem_free(zap, sizeof (zap_t));
396         return (winner);
397     }

399     if (zap->zap_ismicro) {
400         zap->zap_salt = zap->zap_m_phys->mz_salt;
401         zap->zap_normflags = zap->zap_m_phys->mz_normflags;
399         zap->zap_salt = zap->zap_m.zap_phys->mz_salt;
400         zap->zap_normflags = zap->zap_m.zap_phys->mz_normflags;
402         zap->zap_m.zap_num_chunks = db->db_size / MZAP_ENT_LEN - 1;
403         avl_create(&zap->zap_m.zap_avl, mze_compare,
404                 sizeof (mzap_ent_t), offsetof(mzap_ent_t, mze_node));

406         for (i = 0; i < zap->zap_m.zap_num_chunks; i++) {
407             mzap_ent_phys_t *mze =
408                 &zap->zap_m_phys->mz_chunk[i];
409             &zap->zap_m.zap_phys->mz_chunk[i];
409             if (mze->mze_name[0]) {
410                 zap_name_t *zn;

412                 zap->zap_m.zap_num_entries++;
413                 zn = zap_name_alloc(zap, mze->mze_name,
414                                     MT_EXACT);
415                 mze_insert(zap, i, zn->zn_hash);
416                 zap_name_free(zn);
417             }
418         }
419     } else {
420         zap->zap_salt = zap->zap_f_phys->zap_salt;
421         zap->zap_normflags = zap->zap_f_phys->zap_normflags;
422         zap->zap_salt = zap->zap_f.zap_phys->zap_salt;
423         zap->zap_normflags = zap->zap_f.zap_phys->zap_normflags;

423     ASSERT3U(sizeof (struct zap_leaf_header), ==,
424             2*ZAP_LEAF_CHUNKSIZE);

```

```

426      /*
427       * The embedded pointer table should not overlap the
428       * other members.
429       */
430      ASSERT3P(&ZAP_EMBEDDED_PTR_TBL_ENT(zap, 0), >,
431              &zap->zap_f_phys->zap_salt);
432      &zap->zap_f.zap_phys->zap_salt);
433
434      /*
435       * The embedded pointer table should end at the end of
436       * the block
437       */
438      ASSERT3U((uintptr_t)&ZAP_EMBEDDED_PTR_TBL_ENT(zap,
439              1<<ZAP_EMBEDDED_PTR_TBL_SHIFT(zap)) -
440              (uintptr_t)zap->zap_f_phys, ==,
441              (uintptr_t)zap->zap_f.zap_phys, ==,
442              zap->zap_dbuf->db_size);
443      }
444      rw_exit(&zap->zap_rwlock);
445      return (zap);
446 }
447
448 int
449 zap_lockdir(objset_t *os, uint64_t obj, dmu_tx_t *tx,
450             krw_t lti, boolean_t fatreader, boolean_t adding, zap_t **zapp)
451 {
452     zap_t *zap;
453     dmu_buf_t *db;
454     krw_t lt;
455     int err;
456
457     *zapp = NULL;
458
459     err = dmu_buf_hold(os, obj, 0, NULL, &db, DMU_READ_NO_PREFETCH);
460     if (err)
461         return (err);
462
463 #ifdef ZFS_DEBUG
464     {
465         dmu_object_info_t doi;
466         dmu_object_info_from_db(db, &doi);
467         ASSERT3U(DMU_OT_BYTESWAP(doi.doi_type), ==, DMU_BSWAP_ZAP);
468     }
469 #endif
470
471     zap = (zap_t *)dmu_buf_get_user(db);
472     zap = dmu_buf_get_user(db);
473     if (zap == NULL)
474         zap = mzap_open(os, obj, db);
475
476     /*
477      * We're checking zap_ismicro without the lock held, in order to
478      * tell what type of lock we want. Once we have some sort of
479      * lock, see if it really is the right type. In practice this
480      * can only be different if it was upgraded from micro to fat,
481      * and micro wanted WRITER but fat only needs READER.
482      */
483     lt = (!zap->zap_ismicro && fatreader) ? RW_READER : lti;
484     rw_enter(&zap->zap_rwlock, lt);
485     if (lt != ((!zap->zap_ismicro && fatreader) ? RW_READER : lti)) {
486         /* it was upgraded, now we only need reader */
487         ASSERT(lt == RW_WRITER);
488         ASSERT(RW_READER ==
489              (!zap->zap_ismicro && fatreader) ? RW_READER : lti);
490         rw_downgrade(&zap->zap_rwlock);
491         lt = RW_READER;
492     }

```

```

489     }
490
491     zap->zap_objset = os;
492
493     if (lt == RW_WRITER)
494         dmu_buf_will_dirty(db, tx);
495
496     ASSERT3P(zap->zap_dbuf, ==, db);
497
498     ASSERT(!zap->zap_ismicro ||
499           zap->zap_m.zap_num_entries <= zap->zap_m.zap_num_chunks);
500     if (zap->zap_ismicro && tx && adding &&
501         zap->zap_m.zap_num_entries == zap->zap_m.zap_num_chunks) {
502         uint64_t newsz = db->db_size + SPA_MINBLOCKSIZE;
503         if (newsz > MZAP_MAX_BLKSIZE) {
504             dprintf("upgrading obj %llu: num_entries=%u\n",
505                   obj, zap->zap_m.zap_num_entries);
506             *zapp = zap;
507             return (mzap_upgrade(zapp, tx, 0));
508         }
509         err = dmu_object_set_blocksize(os, obj, newsz, 0, tx);
510         ASSERT0(err);
511         zap->zap_m.zap_num_chunks =
512             db->db_size / MZAP_ENT_LEN - 1;
513     }
514
515     *zapp = zap;
516     return (0);
517 }
518
519 unchanged portion omitted
520
521 void
522 zap_evict(dmu_buf_user_t *dbu)
523 {
524     zap_t *zap = (zap_t *)dbu;
525     zap_t *zap = vzap;
526
527     rw_destroy(&zap->zap_rwlock);
528
529     if (zap->zap_ismicro)
530         mze_destroy(zap);
531     else
532         mutex_destroy(&zap->zap_f.zap_num_entries_mtx);
533
534     kmem_free(zap, sizeof (zap_t));
535 }
536
537 unchanged portion omitted
538
539 static void
540 mzap_addent(zap_name_t *zn, uint64_t value)
541 {
542     int i;
543     zap_t *zap = zn->zn_zap;
544     int start = zap->zap_m.zap_alloc_next;
545     uint32_t cd;
546
547     ASSERT(RW_WRITE_HELD(&zap->zap_rwlock));
548
549 #ifdef ZFS_DEBUG
550     for (i = 0; i < zap->zap_m.zap_num_chunks; i++) {
551         mzap_ent_phys_t *mze = &zap->zap_m_phys->mz_chunk[i];
552         mzap_ent_phys_t *mze = &zap->zap_m.zap_phys->mz_chunk[i];
553         ASSERT(strcmp(zn->zn_key_orig, mze->mze_name) != 0);
554     }

```

```

944 #endif
945
946     cd = mze_find_unused_cd(zap, zn->zn_hash);
947     /* given the limited size of the microzap, this can't happen */
948     ASSERT(cd < zap_maxcd(zap));
949
950 again:
951     for (i = start; i < zap->zap_m.zap_num_chunks; i++) {
952         mzap_ent_phys_t *mze = &zap->zap_m_phys->mz_chunk[i];
953         mzap_ent_phys_t *mze = &zap->zap_m.zap_phys->mz_chunk[i];
954         if (mze->mze_name[0] == 0) {
955             mze->mze_value = value;
956             mze->mze_cd = cd;
957             (void) strcpy(mze->mze_name, zn->zn_key_orig);
958             zap->zap_m.zap_num_entries++;
959             zap->zap_m.zap_alloc_next = i+1;
960             if (zap->zap_m.zap_alloc_next ==
961                 zap->zap_m.zap_num_chunks)
962                 zap->zap_m.zap_alloc_next = 0;
963             mze_insert(zap, i, zn->zn_hash);
964             return;
965         }
966     }
967     if (start != 0) {
968         start = 0;
969         goto again;
970     }
971     ASSERT(!"out of entries!");

```

unchanged portion omitted

```

1128 int
1129 zap_remove_norm(objset_t *os, uint64_t zapobj, const char *name,
1130     matchtype_t mt, dmu_tx_t *tx)
1131 {
1132     zap_t *zap;
1133     int err;
1134     mzap_ent_t *mze;
1135     zap_name_t *zn;
1136
1137     err = zap_lockdir(os, zapobj, tx, RW_WRITER, TRUE, FALSE, &zap);
1138     if (err)
1139         return (err);
1140     zn = zap_name_alloc(zap, name, mt);
1141     if (zn == NULL) {
1142         zap_unlockdir(zap);
1143         return (SET_ERROR(ENOTSUP));
1144     }
1145     if (!zap->zap_ismicro) {
1146         err = fzap_remove(zn, tx);
1147     } else {
1148         mze = mze_find(zn);
1149         if (mze == NULL) {
1150             err = SET_ERROR(ENOENT);
1151         } else {
1152             zap->zap_m.zap_num_entries--;
1153             bzero(&zap->zap_m_phys->mz_chunk[mze->mze_chunkid],
1154                 bzero(&zap->zap_m.zap_phys->mz_chunk[mze->mze_chunkid],
1155                     sizeof (mzap_ent_phys_t));
1156             mze_remove(zap, mze);
1157         }
1158     }
1159     zap_name_free(zn);
1160     zap_unlockdir(zap);
1161     return (err);

```

unchanged portion omitted

```

*****
10563 Thu Apr 25 16:15:02 2013
new/usr/src/uts/common/fs/zfs/zfs_sa.c
3752 want more verifiable dbuf user eviction
Submitted by: Justin Gibbs <justing@spectrallogic.com>
Submitted by: Will Andrews <willa@spectrallogic.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) 2010, Oracle and/or its affiliates. All rights reserved.
23 */

25 #include <sys/zfs_context.h>
26 #endif /* ! codereview */
27 #include <sys/types.h>
28 #include <sys/param.h>
29 #include <sys/vnode.h>
30 #include <sys/sa.h>
31 #include <sys/zfs_acl.h>
32 #include <sys/zfs_sa.h>

34 /*
35  * ZPL attribute registration table.
36  * Order of attributes doesn't matter
37  * a unique value will be assigned for each
38  * attribute that is file system specific
39  *
40  * This is just the set of ZPL attributes that this
41  * version of ZFS deals with natively. The file system
42  * could have other attributes stored in files, but they will be
43  * ignored. The SA framework will preserve them, just that
44  * this version of ZFS won't change or delete them.
45  */

47 sa_attr_reg_t zfs_attr_table[ZPL_END+1] = {
48     {"ZPL_ATIME", sizeof (uint64_t) * 2, SA_UINT64_ARRAY, 0},
49     {"ZPL_MTIME", sizeof (uint64_t) * 2, SA_UINT64_ARRAY, 1},
50     {"ZPL_CTIME", sizeof (uint64_t) * 2, SA_UINT64_ARRAY, 2},
51     {"ZPL_CRTIME", sizeof (uint64_t) * 2, SA_UINT64_ARRAY, 3},
52     {"ZPL_GEN", sizeof (uint64_t), SA_UINT64_ARRAY, 4},
53     {"ZPL_MODE", sizeof (uint64_t), SA_UINT64_ARRAY, 5},
54     {"ZPL_SIZE", sizeof (uint64_t), SA_UINT64_ARRAY, 6},
55     {"ZPL_PARENT", sizeof (uint64_t), SA_UINT64_ARRAY, 7},
56     {"ZPL_LINKS", sizeof (uint64_t), SA_UINT64_ARRAY, 8},
57     {"ZPL_XATTR", sizeof (uint64_t), SA_UINT64_ARRAY, 9},
58     {"ZPL_RDEV", sizeof (uint64_t), SA_UINT64_ARRAY, 10},
59     {"ZPL_FLAGS", sizeof (uint64_t), SA_UINT64_ARRAY, 11},

```

```

60     {"ZPL_UID", sizeof (uint64_t), SA_UINT64_ARRAY, 12},
61     {"ZPL_GID", sizeof (uint64_t), SA_UINT64_ARRAY, 13},
62     {"ZPL_PAD", sizeof (uint64_t) * 4, SA_UINT64_ARRAY, 14},
63     {"ZPL_ZNODE_ACL", 88, SA_UINT8_ARRAY, 15},
64     {"ZPL_DACL_COUNT", sizeof (uint64_t), SA_UINT64_ARRAY, 0},
65     {"ZPL_SYMLINK", 0, SA_UINT8_ARRAY, 0},
66     {"ZPL_SCANSTAMP", 32, SA_UINT8_ARRAY, 0},
67     {"ZPL_DACL_ACES", 0, SA_ACL, 0},
68     {NULL, 0, 0, 0}
69 };

71 #ifdef _KERNEL

73 int
74 zfs_sa_readlink(znode_t *zp, uio_t *uio)
75 {
76     dmu_buf_t *db = sa_get_db(zp->z_sa_hdl);
77     size_t bufsz;
78     int error;

80     bufsz = zp->z_size;
81     if (bufsz + ZFS_OLD_ZNODE_PHYS_SIZE <= db->db_size) {
82         error = uiomove((caddr_t)db->db_data +
83             ZFS_OLD_ZNODE_PHYS_SIZE,
84             MIN((size_t)bufsz, uio->uio_resid), UIO_READ, uio);
85     } else {
86         dmu_buf_t *dbp;
87         if ((error = dmu_buf_hold(zp->z_zfsvfs->z_os, zp->z_id,
88             0, FTAG, &dbp, DMU_READ_NO_PREFETCH)) == 0) {
89             error = uiomove(dbp->db_data,
90                 MIN((size_t)bufsz, uio->uio_resid), UIO_READ, uio);
91             dmu_buf_rele(dbp, FTAG);
92         }
93     }
94     return (error);
95 }

97 void
98 zfs_sa_symlink(znode_t *zp, char *link, int len, dmu_tx_t *tx)
99 {
100     dmu_buf_t *db = sa_get_db(zp->z_sa_hdl);

102     if (ZFS_OLD_ZNODE_PHYS_SIZE + len <= dmu_bonus_max()) {
103         VERIFY(dmu_set_bonus(db,
104             len + ZFS_OLD_ZNODE_PHYS_SIZE, tx) == 0);
105         if (len) {
106             bcopy(link, (caddr_t)db->db_data +
107                 ZFS_OLD_ZNODE_PHYS_SIZE, len);
108         }
109     } else {
110         dmu_buf_t *dbp;

112         zfs_grow_blocksize(zp, len, tx);
113         VERIFY(0 == dmu_buf_hold(zp->z_zfsvfs->z_os,
114             zp->z_id, 0, FTAG, &dbp, DMU_READ_NO_PREFETCH));

116         dmu_buf_will_dirty(dbp, tx);

118         ASSERT3U(len, <=, dbp->db_size);
119         bcopy(link, dbp->db_data, len);
120         dmu_buf_rele(dbp, FTAG);
121     }
122 }

124 void
125 zfs_sa_get_scanstamp(znode_t *zp, xvattr_t *xvap)

```

```

126 {
127     zfsvfs_t *zfsvfs = zp->z_zfsvfs;
128     xoptattr_t *xoap;

130     ASSERT(MUTEX_HELD(&zp->z_lock));
131     VERIFY((xoap = xva_getxoptattr(xvap)) != NULL);
132     if (zp->z_is_sa) {
133         if (sa_lookup(zp->z_sa_hdl, SA_ZPL_SCANSTAMP(zfsvfs),
134             &xoap->xoa_av_scanstamp,
135             sizeof(xoap->xoa_av_scanstamp)) != 0)
136             return;
137     } else {
138         dmu_object_info_t doi;
139         dmu_buf_t *db = sa_get_db(zp->z_sa_hdl);
140         int len;

142         if (!(zp->z_pflags & ZFS_BONUS_SCANSTAMP))
143             return;

145         sa_object_info(zp->z_sa_hdl, &doi);
146         len = sizeof(xoap->xoa_av_scanstamp) +
147             ZFS_OLD_ZNODE_PHYS_SIZE;

149         if (len <= doi.doi_bonus_size) {
150             (void) memcpy(xoap->xoa_av_scanstamp,
151                 (caddr_t)db->db_data + ZFS_OLD_ZNODE_PHYS_SIZE,
152                 sizeof(xoap->xoa_av_scanstamp));
153         }
154     }
155     XVA_SET_RTN(xvap, XAT_AV_SCANSTAMP);
156 }

158 void
159 zfs_sa_set_scanstamp(znode_t *zp, xvattr_t *xvap, dmu_tx_t *tx)
160 {
161     zfsvfs_t *zfsvfs = zp->z_zfsvfs;
162     xoptattr_t *xoap;

164     ASSERT(MUTEX_HELD(&zp->z_lock));
165     VERIFY((xoap = xva_getxoptattr(xvap)) != NULL);
166     if (zp->z_is_sa)
167         VERIFY(0 == sa_update(zp->z_sa_hdl, SA_ZPL_SCANSTAMP(zfsvfs),
168             &xoap->xoa_av_scanstamp,
169             sizeof(xoap->xoa_av_scanstamp), tx));
170     else {
171         dmu_object_info_t doi;
172         dmu_buf_t *db = sa_get_db(zp->z_sa_hdl);
173         int len;

175         sa_object_info(zp->z_sa_hdl, &doi);
176         len = sizeof(xoap->xoa_av_scanstamp) +
177             ZFS_OLD_ZNODE_PHYS_SIZE;
178         if (len > doi.doi_bonus_size)
179             VERIFY(dmu_set_bonus(db, len, tx) == 0);
180         (void) memcpy((caddr_t)db->db_data + ZFS_OLD_ZNODE_PHYS_SIZE,
181             xoap->xoa_av_scanstamp, sizeof(xoap->xoa_av_scanstamp));

183         zp->z_pflags |= ZFS_BONUS_SCANSTAMP;
184         VERIFY(0 == sa_update(zp->z_sa_hdl, SA_ZPL_FLAGS(zfsvfs),
185             &zp->z_pflags, sizeof(uint64_t), tx));
186     }
187 }

189 /*
190  * I'm not convinced we should do any of this upgrade.
191  * since the SA code can read both old/new znode formats

```

```

192  * with probably little to know performance difference.
193  *
194  * All new files will be created with the new format.
195  */

197 void
198 zfs_sa_upgrade(sa_handle_t *hdl, dmu_tx_t *tx)
199 {
200     dmu_buf_t *db = sa_get_db(hdl);
201     znode_t *zp = sa_get_userdata(hdl);
202     zfsvfs_t *zfsvfs = zp->z_zfsvfs;
203     sa_bulk_attr_t bulk[20];
204     int count = 0;
205     sa_bulk_attr_t sa_attrs[20] = { 0 };
206     zfs_acl_locator_cb_t locate = { 0 };
207     uint64_t uid, gid, mode, rdev, xattr, parent;
208     uint64_t crtime[2], mtime[2], ctime[2];
209     zfs_acl_phys_t znode_acl;
210     char scanstamp[AV_SCANSTAMP_SZ];
211     boolean_t drop_lock = B_FALSE;

213     /*
214      * No upgrade if ACL isn't cached
215      * since we won't know which locks are held
216      * and ready the ACL would require special "locked"
217      * interfaces that would be messy
218      */
219     if (zp->z_acl_cached == NULL || ZTOV(zp)->v_type == VLNK)
220         return;

222     /*
223      * If the z_lock is held and we aren't the owner
224      * the just return since we don't want to deadlock
225      * trying to update the status of z_is_sa. This
226      * file can then be upgraded at a later time.
227      *
228      * Otherwise, we know we are doing the
229      * sa_update() that caused us to enter this function.
230      */
231     if (mutex_owner(&zp->z_lock) != curthread) {
232         if (mutex_tryenter(&zp->z_lock) == 0)
233             return;
234         else
235             drop_lock = B_TRUE;
236     }

238     /* First do a bulk query of the attributes that aren't cached */
239     SA_ADD_BULK_ATTR(bulk, count, SA_ZPL_MTIME(zfsvfs), NULL, &mtime, 16);
240     SA_ADD_BULK_ATTR(bulk, count, SA_ZPL_CTIME(zfsvfs), NULL, &ctime, 16);
241     SA_ADD_BULK_ATTR(bulk, count, SA_ZPL_CRTIME(zfsvfs), NULL, &crtime, 16);
242     SA_ADD_BULK_ATTR(bulk, count, SA_ZPL_MODE(zfsvfs), NULL, &mode, 8);
243     SA_ADD_BULK_ATTR(bulk, count, SA_ZPL_PARENT(zfsvfs), NULL, &parent, 8);
244     SA_ADD_BULK_ATTR(bulk, count, SA_ZPL_XATTR(zfsvfs), NULL, &xattr, 8);
245     SA_ADD_BULK_ATTR(bulk, count, SA_ZPL_RDEV(zfsvfs), NULL, &rdev, 8);
246     SA_ADD_BULK_ATTR(bulk, count, SA_ZPL_UID(zfsvfs), NULL, &uid, 8);
247     SA_ADD_BULK_ATTR(bulk, count, SA_ZPL_GID(zfsvfs), NULL, &gid, 8);
248     SA_ADD_BULK_ATTR(bulk, count, SA_ZPL_ZNODE_ACL(zfsvfs), NULL,
249         &znode_acl, 88);

251     if (sa_bulk_lookup_locked(hdl, bulk, count) != 0)
252         goto done;

255     /*
256      * While the order here doesn't matter its best to try and organize
257      * it is such a way to pick up an already existing layout number

```

```

258  */
259  count = 0;
260  SA_ADD_BULK_ATTR(sa_attrs, count, SA_ZPL_MODE(zfsvfs), NULL, &mode, 8);
261  SA_ADD_BULK_ATTR(sa_attrs, count, SA_ZPL_SIZE(zfsvfs), NULL,
262  &zp->z_size, 8);
263  SA_ADD_BULK_ATTR(sa_attrs, count, SA_ZPL_GEN(zfsvfs),
264  NULL, &zp->z_gen, 8);
265  SA_ADD_BULK_ATTR(sa_attrs, count, SA_ZPL_UID(zfsvfs), NULL, &uid, 8);
266  SA_ADD_BULK_ATTR(sa_attrs, count, SA_ZPL_GID(zfsvfs), NULL, &gid, 8);
267  SA_ADD_BULK_ATTR(sa_attrs, count, SA_ZPL_PARENT(zfsvfs),
268  NULL, &parent, 8);
269  SA_ADD_BULK_ATTR(sa_attrs, count, SA_ZPL_FLAGS(zfsvfs), NULL,
270  &zp->z_pflags, 8);
271  SA_ADD_BULK_ATTR(sa_attrs, count, SA_ZPL_ETIME(zfsvfs), NULL,
272  zp->z_etime, 16);
273  SA_ADD_BULK_ATTR(sa_attrs, count, SA_ZPL_MTIME(zfsvfs), NULL,
274  &mtime, 16);
275  SA_ADD_BULK_ATTR(sa_attrs, count, SA_ZPL_CTIME(zfsvfs), NULL,
276  &ctime, 16);
277  SA_ADD_BULK_ATTR(sa_attrs, count, SA_ZPL_CRTIME(zfsvfs), NULL,
278  &crttime, 16);
279  SA_ADD_BULK_ATTR(sa_attrs, count, SA_ZPL_LINKS(zfsvfs), NULL,
280  &zp->z_links, 8);
281  if (zp->z_vnode->v_type == VBLK || zp->z_vnode->v_type == VCHR)
282  SA_ADD_BULK_ATTR(sa_attrs, count, SA_ZPL_RDEV(zfsvfs), NULL,
283  &rdev, 8);
284  SA_ADD_BULK_ATTR(sa_attrs, count, SA_ZPL_DACL_COUNT(zfsvfs), NULL,
285  &zp->z_acl_cached->z_acl_count, 8);
287  if (zp->z_acl_cached->z_version < ZFS_ACL_VERSION_FUID)
288  zfs_acl_xform(zp, zp->z_acl_cached, CRED());
290  locate.cb_aclp = zp->z_acl_cached;
291  SA_ADD_BULK_ATTR(sa_attrs, count, SA_ZPL_DACL_ACES(zfsvfs),
292  zfs_acl_data_locator, &locate, zp->z_acl_cached->z_acl_bytes);
294  if (xattr)
295  SA_ADD_BULK_ATTR(sa_attrs, count, SA_ZPL_XATTR(zfsvfs),
296  NULL, &xattr, 8);
298  /* if scanstamp then add scanstamp */
300  if (zp->z_pflags & ZFS_BONUS_SCANSTAMP) {
301  bcopy((caddr_t)db->db_data + ZFS_OLD_ZNODE_PHYS_SIZE,
302  scanstamp, AV_SCANSTAMP_SZ);
303  SA_ADD_BULK_ATTR(sa_attrs, count, SA_ZPL_SCANSTAMP(zfsvfs),
304  NULL, scanstamp, AV_SCANSTAMP_SZ);
305  zp->z_pflags &= ~ZFS_BONUS_SCANSTAMP;
306  }
308  VERIFY(dmu_set_bonustype(db, DMU_OT_SA, tx) == 0);
309  VERIFY(sa_replace_all_by_template_locked(hdl, sa_attrs,
310  count, tx) == 0);
311  if (znode_acl.z_acl_extern_obj)
312  VERIFY(0 == dmu_object_free(zfsvfs->z_os,
313  znode_acl.z_acl_extern_obj, tx));
315  zp->z_is_sa = B_TRUE;
316  done:
317  if (drop_lock)
318  mutex_exit(&zp->z_lock);
319  }
321  void
322  zfs_sa_upgrade_txholds(dmu_tx_t *tx, znode_t *zp)
323  {

```

```

324  if (!zp->z_zfsvfs->z_use_sa || zp->z_is_sa)
325  return;
328  dmu_tx_hold_sa(tx, zp->z_sa_hdl, B_TRUE);
330  if (zfs_external_acl(zp)) {
331  dmu_tx_hold_free(tx, zfs_external_acl(zp), 0,
332  DMU_OBJECT_END);
333  }
334  }
336  #endif

```

```

*****
53536 Thu Apr 25 16:15:03 2013
new/usr/src/uts/common/fs/zfs/zfs_znode.c
3752 want more verifiable dbuf user eviction
Submitted by: Justin Gibbs <justing@spectralogic.com>
Submitted by: Will Andrews <willa@spectralogic.com>
*****
_____unchanged_portion_omitted_____

1104 int
1105 zfs_zget(zfsvfs_t *zfsvfs, uint64_t obj_num, znode_t **zpp)
1106 {
1107     dmu_object_info_t doi;
1108     dmu_buf_t *db;
1109     znode_t *zp;
1110     int err;
1111     sa_handle_t *hdl;
1112
1113     *zpp = NULL;
1114
1115     ZFS_OBJ_HOLD_ENTER(zfsvfs, obj_num);
1116
1117     err = sa_buf_hold(zfsvfs->z_os, obj_num, NULL, &db);
1118     if (err) {
1119         ZFS_OBJ_HOLD_EXIT(zfsvfs, obj_num);
1120         return (err);
1121     }
1122
1123     dmu_object_info_from_db(db, &doi);
1124     if (doi.doi_bonus_type != DMU_OT_SA &&
1125         (doi.doi_bonus_type != DMU_OT_ZNODE ||
1126          (doi.doi_bonus_type == DMU_OT_ZNODE &&
1127           doi.doi_bonus_size < sizeof (znode_phys_t)))) {
1128         sa_buf_rele(db, NULL);
1129         ZFS_OBJ_HOLD_EXIT(zfsvfs, obj_num);
1130         return (SET_ERROR(EINVAL));
1131     }
1132
1133     hdl = (sa_handle_t *)dmu_buf_get_user(db);
1134     hdl = dmu_buf_get_user(db);
1135     if (hdl != NULL) {
1136         zp = sa_get_userdata(hdl);
1137     }
1138
1139     /*
1140     * Since "SA" does immediate eviction we
1141     * should never find a sa handle that doesn't
1142     * know about the znode.
1143     */
1144
1145     ASSERT3P(zp, !=, NULL);
1146
1147     mutex_enter(&zp->z_lock);
1148     ASSERT3U(zp->z_id, ==, obj_num);
1149     if (zp->z_unlinked) {
1150         err = SET_ERROR(ENOENT);
1151     } else {
1152         VN_HOLD(ZTOV(zp));
1153         *zpp = zp;
1154         err = 0;
1155     }
1156     sa_buf_rele(db, NULL);
1157     mutex_exit(&zp->z_lock);
1158     ZFS_OBJ_HOLD_EXIT(zfsvfs, obj_num);
1159     return (err);
1160 }

```

```

1161     /*
1162     * Not found create new znode/vnode
1163     * but only if file exists.
1164     */
1165     * There is a small window where zfs_vget() could
1166     * find this object while a file create is still in
1167     * progress. This is checked for in zfs_znode_alloc()
1168     *
1169     * if zfs_znode_alloc() fails it will drop the hold on the
1170     * bonus buffer.
1171     */
1172     zp = zfs_znode_alloc(zfsvfs, db, doi.doi_data_block_size,
1173         doi.doi_bonus_type, NULL);
1174     if (zp == NULL) {
1175         err = SET_ERROR(ENOENT);
1176     } else {
1177         *zpp = zp;
1178     }
1179     ZFS_OBJ_HOLD_EXIT(zfsvfs, obj_num);
1180     return (err);
1181 }
_____unchanged_portion_omitted_____

```