

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

1

35483 Thu Aug 16 22:11:07 2012

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

dmu_tx_count_free is doing a horrible over-estimation of used memory. It assumes that the file is fully non-sparse and calculates a worst-case estimate of how much memory is needed to hold all metadata for the file. If a large hole needs to be freed, the estimation goes into the TB-range, which obviously fails later on.

This patch tries to calculate a more realistic estimate by counting the 11 blocks (the loop for this is already present) and assumes a worst-case distribution of those blocks over the full length given.

unchanged_portion_omitted

```
423 static void
424 dmu_tx_count_free(dmu_tx_hold_t *txh, uint64_t off, uint64_t len)
425 {
426     uint64_t blkid, nblk, lastblk;
427     uint64_t space = 0, unref = 0, skipped = 0;
428     dnode_t *dn = txh->txh_dnode;
429     dsl_dataset_t *ds = dn->dn_objset->os_dsl_dataset;
430     spa_t *spa = txh->txh_tx->tx_pool->dp_spa;
431     int epbs;
432     uint64_t l0span = 0, n1blk = 0;
433 #endif /* ! codereview */
434
435     if (dn->dn_nlevels == 0)
436         return;
437
438     /*
439      * The struct_rwlock protects us against dn_nlevels
440      * changing, in case (against all odds) we manage to dirty &
441      * sync out the changes after we check for being dirty.
442      * Also, dbuf_hold_impl() wants us to have the struct_rwlock.
443     */
444     rw_enter(&dn->dn_struct_rwlock, RW_READER);
445     epbs = dn->dn_inblkshift - SPA_BLKPTRSHIFT;
446     if (dn->dn_maxblkid == 0) {
447         if (off == 0 && len >= dn->dn_datblksz) {
448             blkid = 0;
449             nblk = 1;
450         } else {
451             rw_exit(&dn->dn_struct_rwlock);
452             return;
453         }
454     } else {
455         blkid = off >> dn->dn_datblkshift;
456         nblk = (len + dn->dn_datblksz - 1) >> dn->dn_datblkshift;
457
458         if (blkid >= dn->dn_maxblkid) {
459             rw_exit(&dn->dn_struct_rwlock);
460             return;
461         }
462         if (blkid + nblk > dn->dn_maxblkid)
463             nblk = dn->dn_maxblkid - blkid;
464
465     }
466     l0span = nblk; /* save for later use to calc level > 1 overhead */
467 #endif /* ! codereview */
468     if (dn->dn_nlevels == 1) {
469         int i;
470         for (i = 0; i < nblk; i++) {
471             blkptr_t *bp = dn->dn_phys->dn_blkptr;
472             ASSERT3U(blkid + i, <, dn->dn_nblkptr);
473             bp += blkid + i;
474             if (dsl_dataset_block_freeable(ds, bp, bp->blk_birth)) {
```

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

2

```
475             dprintf_bp(bp, "can free old%s", "");
476             space += bp_get_dsize(spa, bp);
477         }
478         unref += BP_GET_ASIZE(bp);
479     }
480     n1blk = 1;
481 #endif /* ! codereview */
482     nblk = 0;
483 }
484
485 /*
486  * Add in memory requirements of higher-level indirects.
487  * This assumes a worst-possible scenario for dn_nlevels.
488 */
489 {
490     uint64_t blkcnt = 1 + ((nblk >> epbs) >> epbs);
491     int level = (dn->dn_nlevels > 1) ? 2 : 1;
492
493     while (level++ < DN_MAX_LEVELS) {
494         txh->txh_memory_tohold += blkcnt << dn->dn_inblkshift;
495         blkcnt = 1 + (blkcnt >> epbs);
496     }
497     ASSERT(blkcnt <= dn->dn_nblkptr);
498 }
499
500 lastblk = blkid + nblk - 1;
501 while (nblk) {
502     dbuf_t *dbuf;
503     uint64_t ibyte, new_blkid;
504     int epb = 1 << epbs;
505     int err, blkoff, tochk;
506     blkptr_t *bp;
507
508     ibyte = blkid << dn->dn_datblkshift;
509     err = dnode_next_offset(dn,
510                             DNODE_FIND_HAVELOCK, &ibyte, 2, 1, 0);
511     new_blkid = ibyte >> dn->dn_datblkshift;
512     if (err == ESRCH) {
513         skipped += (lastblk >> epbs) - (blkid >> epbs) + 1;
514         break;
515     }
516     if (err) {
517         txh->txh_tx->tx_err = err;
518         break;
519     }
520     if (new_blkid > lastblk) {
521         skipped += (lastblk >> epbs) - (blkid >> epbs) + 1;
522         break;
523     }
524
525     if (new_blkid > blkid) {
526         ASSERT((new_blkid >> epbs) > (blkid >> epbs));
527         skipped += (new_blkid >> epbs) - (blkid >> epbs) - 1;
528         nblk -= new_blkid - blkid;
529         blkid = new_blkid;
530     }
531     blkoff = P2PHASE(blkid, epb);
532     tochk = MIN(epb - blkoff, nblk);
533
534     err = dbuf_hold_impl(dn, 1, blkid >> epbs, FALSE, FTAG, &dbuf);
535     if (err) {
536         txh->txh_tx->tx_err = err;
537         break;
538     }
539
540     txh->txh_memory_tohold += dbuf->db.db_size;
```

```

527
528     /*
529      * We don't check memory_tohold against DMU_MAX_ACCESS because
530      * memory_tohold is an over-estimation (especially the >L1
531      * indirect blocks), so it could fail. Callers should have
532      * already verified that they will not be holding too much
533      * memory.
534     */
535
536     err =dbuf_read(dbuf, NULL, DB_RF_HAVESTRUCT | DB_RF_CANFAIL);
537     if (err != 0) {
538         txh->txh_tx->tx_err = err;
539         dbuf_rele(dbuf, FTAG);
540         break;
541     }
542
543     bp =dbuf->db.db_data;
544     bp += blkoff;
545
546     for (i = 0; i < tochk; i++) {
547         if (ds1_dataset_block_freeable(ds, &bp[i],
548             bp[i].blk_birth)) {
549             dprintf_bp(&bp[i], "can free old%s", "");
550             space += bp_get_dsize(spa, &bp[i]);
551         }
552         unref += BP_GET_ASIZE(bp);
553     }
554     dbuf_rele(dbuf, FTAG);
555
556 #endif /* ! codereview */
557     blkid += tochk;
558     nblkks -= tochk;
559 }
560 rw_exit(&dn->dn_struct_rwlock);
561
562 /*
563  * Add in memory requirements of higher-level indirects.
564  * This assumes a worst-possible scenario for dn_nlevels and a
565  * worst-possible distribution of l1-blocks over the region to free.
566  */
567 {
568     uint64_t blkcnt = 1 + ((10span >> epbs) >> epbs);
569     int level = 2;
570
571     /*
572      * Here we don't use DN_MAX_LEVEL, but calculate it with the
573      * given datablkshift and indblkshift. This makes the
574      * difference between 19 and 8 on large files.
575     */
576     int maxlevel = 2 + (DN_MAX_OFFSET_SHIFT - dn->dn_datablkshift) /
577                     (dn->dn_indblksshift - SPA_BLKPTRSHIFT);
578
579     while (level++ < maxlevel) {
580         txh->txh_memory_tohold += MIN(blkcnt, (nblkks >> epbs))
581             << dn->dn_indblksshift;
582         blkcnt = 1 + (blkcnt >> epbs);
583     }
584 }
585 #endif /* ! codereview */
586 /* account for new level 1 indirect blocks that might show up */
587 if (skipped > 0) {
588     txh->txh_fudge += skipped << dn->dn_indblksshift;
589     skipped = MIN(skipped, DMU_MAX_DELETEBLKCNT >> epbs);
590     txh->txh_memory_tohold += skipped << dn->dn_indblksshift;
591 }

```

```

592     txh->txh_space_tofree += space;
593     txh->txh_space_tounref += unref;
594 }
595
596 void
597 dmu_tx_hold_free(dmu_tx_t *tx, uint64_t object, uint64_t off, uint64_t len)
598 {
599     dmu_tx_hold_t *txh;
600     dnode_t *dn;
601     uint64_t start, end, i;
602     int err, shift;
603     zio_t *zio;
604
605     ASSERT(tx->tx_txg == 0);
606
607     txh = dmu_tx_hold_objectImpl(tx, tx->tx_objset,
608                                  object, THT_FREE, off, len);
609     if (txh == NULL)
610         return;
611     dn = txh->txh_dnode;
612
613     /* first block */
614     if (off != 0)
615         dmu_tx_count_write(txh, off, 1);
616     /* last block */
617     if (len != DMU_OBJECT_END)
618         dmu_tx_count_write(txh, off+len, 1);
619
620     dmu_tx_count_dnode(txh);
621
622     if (off >= (dn->dn_maxblkid+1) * dn->dn_datblkksz)
623         return;
624     if (len == DMU_OBJECT_END)
625         len = (dn->dn_maxblkid+1) * dn->dn_datblkksz - off;
626
627     /*
628      * For i/o error checking, read the first and last level-0
629      * blocks, and all the level-1 blocks. The above count_write's
630      * have already taken care of the level-0 blocks.
631     */
632     if (dn->dn_nlevels > 1) {
633         shift = dn->dn_datblkshift + dn->dn_indblksshift -
634                 SPA_BLKPTRSHIFT;
635         start = off >> shift;
636         end = dn->dn_datblkshift ? ((off+len) >> shift) : 0;
637
638         zio = zio_root(tx->tx_pool->dp_spa,
639                        NULL, NULL, ZIO_FLAG_CANFAIL);
640         for (i = start; i <= end; i++) {
641             uint64_t ibyte = i << shift;
642             err = dnode_next_offset(dn, 0, &ibyte, 2, 1, 0);
643             i = ibyte >> shift;
644             if (err == ESRCH)
645                 break;
646             if (err) {
647                 tx->tx_err = err;
648                 return;
649             }
650
651             err = dmu_tx_check_iocerr(zio, dn, 1, i);
652             if (err) {
653                 tx->tx_err = err;
654                 return;
655             }
656         }
657     }
658
659     err = zio_wait(zio);

```

```

658         if (err) {
659             tx->tx_err = err;
660             return;
661         }
662     }
663
664     dmu_tx_count_free(txh, off, len);
665 }
666
667 void
668 dmu_tx_hold_zap(dmu_tx_t *tx, uint64_t object, int add, const char *name)
669 {
670     dmu_tx_hold_t *txh;
671     dnode_t *dn;
672     uint64_t nblocks;
673     int epbs, err;
674
675     ASSERT(tx->tx_txg == 0);
676
677     txh = dmu_tx_hold_object_impl(tx, tx->tx_objset,
678                                   object, THT_ZAP, add, (uintptr_t)name);
679     if (txh == NULL)
680         return;
681     dn = txh->txh_dnode;
682
683     dmu_tx_count_dnode(txh);
684
685     if (dn == NULL) {
686         /*
687          * We will be able to fit a new object's entries into one leaf
688          * block. So there will be at most 2 blocks total,
689          * including the header block.
690         */
691         dmu_tx_count_write(txh, 0, 2 << fzap_default_block_shift);
692         return;
693     }
694
695     ASSERT3P(DMU_OT_BYTESWAP(dn->dn_type), ==, DMU_BSWAP_ZAP);
696
697     if (dn->dn_maxblkid == 0 && !add) {
698         blkptr_t *bp;
699
700         /*
701          * If there is only one block (i.e. this is a micro-zap)
702          * and we are not adding anything, the accounting is simple.
703         */
704         err = dmu_tx_check_ioerr(NULL, dn, 0, 0);
705         if (err) {
706             tx->tx_err = err;
707             return;
708         }
709
710         /*
711          * Use max block size here, since we don't know how much
712          * the size will change between now and the dbuf dirty call.
713         */
714         bp = &dn->dn_phys->dn_blkptr[0];
715         if (dsl_dataset_block_freeable(dn->dn_objset->os_dsl_dataset,
716                                       bp, bp->blk_birth))
717             txh->txh_space_tooverwrite += SPA_MAXBLOCKSIZE;
718         else
719             txh->txh_space_towrite += SPA_MAXBLOCKSIZE;
720         if (!BP_IS_HOLE(bp))
721             txh->txh_space_tounref += SPA_MAXBLOCKSIZE;
722         return;
723     }

```

```

725     if (dn->dn_maxblkid > 0 && name) {
726         /*
727          * access the name in this fat-zap so that we'll check
728          * for i/o errors to the leaf blocks, etc.
729         */
730         err = zap_lookup(dn->dn_objset, dn->dn_object, name,
731                         8, 0, NULL);
732         if (err == EIO) {
733             tx->tx_err = err;
734             return;
735         }
736     }
737
738     err = zap_count_write(dn->dn_objset, dn->dn_object, name, add,
739                           &txh->txh_space_towrite, &txh->txh_space_tooverwrite);
740
741     /*
742      * If the modified blocks are scattered to the four winds,
743      * we'll have to modify an indirect twig for each.
744     */
745     epbs = dn->dn_indblkshift - SPA_BLKPTRSHIFT;
746     for (nblocks = dn->dn_maxblkid >> epbs; nblocks != 0; nblocks >>= epbs)
747         if (dn->dn_objset->os_dsl_dataset->ds_phys->ds_prev_snap_obj)
748             txh->txh_space_towrite += 3 << dn->dn_indblkshift;
749         else
750             txh->txh_space_tooverwrite += 3 << dn->dn_indblkshift;
751 }
752
753 void
754 dmu_tx_hold_bonus(dmu_tx_t *tx, uint64_t object)
755 {
756     dmu_tx_hold_t *txh;
757
758     ASSERT(tx->tx_txg == 0);
759
760     txh = dmu_tx_hold_object_impl(tx, tx->tx_objset,
761                                   object, THT_BONUS, 0, 0);
762     if (txh)
763         dmu_tx_count_dnode(txh);
764 }
765
766 void
767 dmu_tx_hold_space(dmu_tx_t *tx, uint64_t space)
768 {
769     dmu_tx_hold_t *txh;
770     ASSERT(tx->tx_txg == 0);
771
772     txh = dmu_tx_hold_object_impl(tx, tx->tx_objset,
773                                 DMU_NEW_OBJECT, THT_SPACE, space, 0);
774
775     txh->txh_space_towrite += space;
776 }
777
778 int
779 dmu_tx_holds(dmu_tx_t *tx, uint64_t object)
780 {
781     dmu_tx_hold_t *txh;
782     int holds = 0;
783
784     /*
785      * By asserting that the tx is assigned, we're counting the
786      * number of dn_tx_holds, which is the same as the number of
787      * dn_holds. Otherwise, we'd be counting dn_holds, but
788      * dn_tx_holds could be 0.
789     */

```

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

7

```

990 ASSERT(tx->tx_txg != 0);

991 /* if (tx->tx_anyobj == TRUE) */
992 /*     return (0); */

993 for (txh = list_head(&tx->tx_holds); txh;
994      txh = list_next(&tx->tx_holds, txh)) {
995     if (txh->txh_dnnode && txh->txh_dnnode->dn_object == object)
996         holds++;
997 }

801     return (holds);
802 }

803 #ifdef ZFS_DEBUG
804 void
805 dmu_tx_dirty_buf(dmu_tx_t *tx, dmu_buf_impl_t *db)
806 {
807     dmu_tx_hold_t *txh;
808     int match_object = FALSE, match_offset = FALSE;
809     dnode_t *dn;
810
811     DB_DNODE_ENTER(db);
812     dn = DB_DNODE(db);
813     ASSERT(tx->tx_txg != 0);
814     ASSERT(tx->tx_objset == NULL || dn->dn_objset == tx->tx_objset);
815     ASSERT3U(dn->dn_object, ==, db->db.db_object);
816
817     if (tx->tx_anyobj) {
818         DB_DNODE_EXIT(db);
819         return;
820     }
821
822     /* XXX No checking on the meta dnode for now */
823     if (db->db_object == DMU_META_DNODE_OBJECT) {
824         DB_DNODE_EXIT(db);
825         return;
826     }
827
828     for (txh = list_head(&tx->tx_holds); txh;
829          txh = list_next(&tx->tx_holds, txh)) {
830         ASSERT(dn == NULL || dn->dn_assigned_txg == tx->tx_txg);
831         if (txh->txh_dnnode == dn && txh->txh_type != THT_NEWOBJECT)
832             match_object = TRUE;
833         if (txh->txh_dnnode == NULL || txh->txh_dnnode == dn) {
834             int datablkshift = dn->dn_datablkshift ?
835                 dn->dn_datablkshift : SPA_MAXBLOCKSHIFT;
836             int epbs = dn->dn_inblkshift - SPA_BLKPTRSHIFT;
837             int shift = datablkshift + epbs * db->db_level;
838             uint64_t beginblk = shift >= 64 ? 0 :
839                 (txh->txh_arg1 >> shift);
840             uint64_t endblk = shift >= 64 ? 0 :
841                 ((txh->txh_arg1 + txh->txh_arg2 - 1) >> shift);
842             uint64_t blkid = db->db_blkid;
843
844             /* XXX txh_arg2 better not be zero... */
845
846             dprintf("found txh type %x beginblk=%llx endblk=%llx\n",
847                    txh->txh_type, beginblk, endblk);
848
849             switch (txh->txh_type) {
850             case THT_WRITE:
851                 if (blkid >= beginblk && blkid <= endblk)
852                     match_offset = TRUE;
853                 /*
854                  * We will let this hold work for the bonus

```

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

```

856                                     * or spill buffer so that we don't need to
857                                     * hold it when creating a new object.
858                                     */
859     if (blkid == DMU_BONUS_BLKID ||
860         blkid == DMU_SPILL_BLKID)
861         match_offset = TRUE;
862     /*
863      * They might have to increase nlevels,
864      * thus dirtying the new TLIBS. Or the
865      * might have to change the block size,
866      * thus dirying the new lvl=0 blk=0.
867      */
868     if (blkid == 0)
869         match_offset = TRUE;
870     break;
871 case THT_FREE:
872     /*
873      * We will dirty all the level 1 blocks in
874      * the free range and perhaps the first and
875      * last level 0 block.
876      */
877     if (blkid >= beginblk && (blkid <= endblk ||
878         txh->txh_arg2 == DMU_OBJECT_END))
879         match_offset = TRUE;
880     break;
881 case THT_SPILL:
882     if (blkid == DMU_SPILL_BLKID)
883         match_offset = TRUE;
884     break;
885 case THT_BONUS:
886     if (blkid == DMU_BONUS_BLKID)
887         match_offset = TRUE;
888     break;
889 case THT_ZAP:
890     match_offset = TRUE;
891     break;
892 case THT_NEWORDOBJECT:
893     match_object = TRUE;
894     break;
895 default:
896     ASSERT(!"bad txh_type");
897 }
898 }
899 if (match_object && match_offset) {
900     DB_DNODE_EXIT(db);
901     return;
902 }
903 DB_DNODE_EXIT(db);
904 panic("dirtying dbuf obj=%llx lvl=%u blkid=%llx but not tx_hold\n",
905       (u_longlong_t)db->db_object, db->db_level,
906       (u_longlong_t)db->db_blkid);
907 }
908 }
909 #endiff

911 static int
912 dmu_tx_try_assign(dmu_tx_t *tx, uint64_t txg_how)
913 {
914     dmu_tx_hold_t *txh;
915     spa_t *spa = tx->tx_pool->dp_spa;
916     uint64_t memory, asize, fsize, usize;
917     uint64_t towrite, tofree, tooverwrite, tounref, tohold, fudge;

919     ASSERT3U(tx->tx_txg, ==, 0);

921     if (tx->tx_err)

```

```

922         return (tx->tx_err);
923
924     if (spa_suspended(spa)) {
925         /*
926          * If the user has indicated a blocking failure mode
927          * then return ERESTART which will block in dmu_tx_wait().
928          * Otherwise, return EIO so that an error can get
929          * propagated back to the VOP calls.
930          *
931          * Note that we always honor the txg_how flag regardless
932          * of the failuremode setting.
933          */
934     if (spa_get_failmode(spa) == ZIO_FAILURE_MODE_CONTINUE &&
935         txg_how != TXG_WAIT)
936         return (EIO);
937
938     return (ERESTART);
939 }
940
941 tx->tx_txg = txg_hold_open(tx->tx_pool, &tx->tx_txgh);
942 tx->tx_needassign_txh = NULL;
943
944 /*
945  * NB: No error returns are allowed after txg_hold_open, but
946  * before processing the dnode holds, due to the
947  * dmu_tx_unassign() logic.
948 */
949
950 towrite = tofree = tooverwrite = tounref = tohold = fudge = 0;
951 for (txh = list_head(&tx->tx_holds); txh;
952      txh = list_next(&tx->tx_holds, txh)) {
953     dnode_t *dn = txh->txh_dnode;
954     if (dn != NULL) {
955         mutex_enter(&dn->dn_mtx);
956         if (dn->dn_assigned_txg == tx->tx_txg - 1) {
957             mutex_exit(&dn->dn_mtx);
958             tx->tx_needassign_txh = txh;
959             return (ERESTART);
960         }
961         if (dn->dn_assigned_txg == 0)
962             dn->dn_assigned_txg = tx->tx_txg;
963         ASSERT3U(dn->dn_assigned_txg, ==, tx->tx_txg);
964         (void) refcount_add(&dn->dn_tx_holds, tx);
965         mutex_exit(&dn->dn_mtx);
966     }
967     towrite += txh->txh_space_towrite;
968     tofree += txh->txh_space_tofree;
969     tooverwrite += txh->txh_space_tooverwrite;
970     tounref += txh->txh_space_tounref;
971     tohold += txh->txh_memory_tohold;
972     fudge += txh->txh_fudge;
973 }
974
975 /*
976  * NB: This check must be after we've held the dnodes, so that
977  * the dmu_tx_unassign() logic will work properly
978  */
979 if (txg_how >= TXG_INITIAL && txg_how != tx->tx_txg)
980     return (ERESTART);
981
982 /*
983  * If a snapshot has been taken since we made our estimates,
984  * assume that we won't be able to free or overwrite anything.
985  */
986 if (tx->tx_objset &&
987     dsl_dataset_prev_snap_txg(tx->tx_objset->os_dsl_dataset) >

```

```

988         tx->tx_lastsnap_txg) {
989             towrite += tooverwrite;
990             tooverwrite = tofree = 0;
991         }
992
993         /* needed allocation: worst-case estimate of write space */
994         asize = spa_get_asize(tx->tx_pool->dp_spa, towrite + tooverwrite);
995         /* freed space estimate: worst-case overwrite + free estimate */
996         fsize = spa_get_asize(tx->tx_pool->dp_spa, tooverwrite) + tofree;
997         /* convert unrefd space to worst-case estimate */
998         usize = spa_get_asize(tx->tx_pool->dp_spa, tounref);
999         /* calculate memory footprint estimate */
1000        memory = towrite + tooverwrite + tohold;
1001
1002 #ifdef ZFS_DEBUG
1003 /*
1004  * Add in 'tohold' to account for our dirty holds on this memory
1005  * XXX - the "fudge" factor is to account for skipped blocks that
1006  * we missed because dnode_next_offset() misses in-core-only blocks.
1007  */
1008        tx->tx_space_towrite = asize +
1009            spa_get_asize(tx->tx_pool->dp_spa, tohold + fudge);
1010        tx->tx_space_tofree = tofree;
1011        tx->tx_space_tooverwrite = tooverwrite;
1012        tx->tx_space_tounref = tounref;
1013 #endif
1014
1015        if (tx->tx_dir && asize != 0) {
1016            int err = dsl_dir_tempreserve_space(tx->tx_dir, memory,
1017                asize, fsize, usize, &tx->tx_tempreserve_cookie, tx);
1018            if (err)
1019                return (err);
1020        }
1021
1022        return (0);
1023    }
1024
1025 static void
1026 dmu_tx_unassign(dmu_tx_t *tx)
1027 {
1028     dmu_tx_hold_t *txh;
1029
1030     if (tx->tx_txg == 0)
1031         return;
1032
1033     txg_rele_to_quiesce(&tx->tx_txgh);
1034
1035     for (txh = list_head(&tx->tx_holds); txh != tx->tx_needassign_txh;
1036          txh = list_next(&tx->tx_holds, txh)) {
1037         dnode_t *dn = txh->txh_dnode;
1038
1039         if (dn == NULL)
1040             continue;
1041         mutex_enter(&dn->dn_mtx);
1042         ASSERT3U(dn->dn_assigned_txg, ==, tx->tx_txg);
1043
1044         if (refcount_remove(&dn->dn_tx_holds, tx) == 0) {
1045             dn->dn_assigned_txg = 0;
1046             cv_broadcast(&dn->dn_notxholds);
1047         }
1048         mutex_exit(&dn->dn_mtx);
1049     }
1050     txg_rele_to_sync(&tx->tx_txgh);
1051
1052     tx->tx_lasttried_txg = tx->tx_txg;

```

```

1054     tx->tx_txg = 0;
1055 }

1057 /*
1058 * Assign tx to a transaction group. txg_how can be one of:
1059 *
1060 * (1) TXG_WAIT. If the current open txg is full, waits until there's
1061 * a new one. This should be used when you're not holding locks.
1062 * If will only fail if we're truly out of space (or over quota).
1063 *
1064 * (2) TXG_NOWAIT. If we can't assign into the current open txg without
1065 * blocking, returns immediately with ERESTART. This should be used
1066 * whenever you're holding locks. On an ERESTART error, the caller
1067 * should drop locks, do a dmu_tx_wait(tx), and try again.
1068 *
1069 * (3) A specific txg. Use this if you need to ensure that multiple
1070 * transactions all sync in the same txg. Like TXG_NOWAIT, it
1071 * returns ERESTART if it can't assign you into the requested txg.
1072 */
1073 int
1074 dmu_tx_assign(dmu_tx_t *tx, uint64_t txg_how)
1075 {
1076     int err;

1078     ASSERT(tx->tx_txg == 0);
1079     ASSERT(txg_how != 0);
1080     ASSERT(!dsl_pool_sync_context(tx->tx_pool));

1082     while ((err = dmu_tx_try_assign(tx, txg_how)) != 0) {
1083         dmu_tx_unassign(tx);

1085         if (err != ERESTART || txg_how != TXG_WAIT)
1086             return (err);

1088         dmu_tx_wait(tx);
1089     }

1091     txg_rele_to_quiesce(&tx->tx_txgh);

1093     return (0);
1094 }

1096 void
1097 dmu_tx_wait(dmu_tx_t *tx)
1098 {
1099     spa_t *spa = tx->tx_pool->dp_spa;
1100
1101     ASSERT(tx->tx_txg == 0);

1103     /*
1104      * It's possible that the pool has become active after this thread
1105      * has tried to obtain a tx. If that's the case then his
1106      * tx_lasttry_txg would not have been assigned.
1107      */
1108     if (spa_suspended(spa) || tx->tx_lasttry_txg == 0) {
1109         txg_wait_synced(tx->tx_pool, spa_last_synced_txg(spa) + 1);
1110     } else if (tx->tx_needassign_txh) {
1111         dnode_t *dn = tx->tx_needassign_txh->txh_dnode;

1113         mutex_enter(&dn->dn_mtx);
1114         while (dn->dn_assigned_txg == tx->tx_lasttry_txg - 1)
1115             cv_wait(&dn->dn_notxholds, &dn->dn_mtx);
1116         mutex_exit(&dn->dn_mtx);
1117         tx->tx_needassign_txh = NULL;
1118     } else {
1119         txg_wait_open(tx->tx_pool, tx->tx_lasttry_txg + 1);

```

```

1120         }
1121     }

1123     void
1124     dmu_tx_willuse_space(dmu_tx_t *tx, int64_t delta)
1125     {
1126 #ifdef ZFS_DEBUG
1127         if (tx->tx_dir == NULL || delta == 0)
1128             return;

1130         if (delta > 0) {
1131             ASSERT3U(refcount_count(&tx->tx_space_written) + delta, <=,
1132                     tx->tx_space_towrite);
1133             (void) refcount_add_many(&tx->tx_space_written, delta, NULL);
1134         } else {
1135             (void) refcount_add_many(&tx->tx_space_freed, -delta, NULL);
1136         }
1137 #endif
1138     }

1140     void
1141     dmu_tx_commit(dmu_tx_t *tx)
1142     {
1143         dmu_tx_hold_t *txh;

1145         ASSERT(tx->tx_txg != 0);

1147         while (txh = list_head(&tx->tx_holds)) {
1148             dnode_t *dn = txh->txh_dnode;

1150             list_remove(&tx->tx_holds, txh);
1151             kmem_free(txh, sizeof (dmu_tx_hold_t));
1152             if (dn == NULL)
1153                 continue;
1154             mutex_enter(&dn->dn_mtx);
1155             ASSERT3U(dn->dn_assigned_txg, ==, tx->tx_txg);

1157             if (refcount_remove(&dn->dn_tx_holds, tx) == 0) {
1158                 dn->dn_assigned_txg = 0;
1159                 cv_broadcast(&dn->dn_notxholds);
1160             }
1161             mutex_exit(&dn->dn_mtx);
1162             dnode_rele(dn, tx);
1163         }

1165         if (tx->tx_tempreserve_cookie)
1166             dsl_dir_tempreserve_clear(tx->tx_tempreserve_cookie, tx);

1168         if (!list_is_empty(&tx->tx_callbacks))
1169             txg_register_callbacks(&tx->tx_txgh, &tx->tx_callbacks);

1171         if (tx->tx_anyobj == FALSE)
1172             txg_rele_to_sync(&tx->tx_txgh);

1174         list_destroy(&tx->tx_callbacks);
1175         list_destroy(&tx->tx_holds);
1176 #ifdef ZFS_DEBUG
1177         dprintf("towrite=%llu written=%llu tofree=%llu freed=%llu\n",
1178                 tx->tx_space_towrite, refcount_count(&tx->tx_space_written),
1179                 tx->tx_space_tofree, refcount_count(&tx->tx_space_freed));
1180         refcount_destroy_many(&tx->tx_space_written,
1181                               refcount_count(&tx->tx_space_written));
1182         refcount_destroy_many(&tx->tx_space_freed,
1183                               refcount_count(&tx->tx_space_freed));
1184 #endif
1185         kmem_free(tx, sizeof (dmu_tx_t));

```

```
new/usr/src/uts/common/fs/zfs/dmu_tx.c
1186 }

1188 void
1189 dmu_tx_abort(dmu_tx_t *tx)
1190 {
1191     dmu_tx_hold_t *txh;
1193     ASSERT(tx->tx_txg == 0);
1195     while (txh = list_head(&tx->tx_holds)) {
1196         dnode_t *dn = txh->txh_dnode;
1198         list_remove(&tx->tx_holds, txh);
1199         kmem_free(txh, sizeof (dmu_tx_hold_t));
1200         if (dn != NULL)
1201             dnode_rele(dn, tx);
1202     }
1204     /*
1205      * Call any registered callbacks with an error code.
1206      */
1207     if (!list_is_empty(&tx->tx_callbacks))
1208         dmu_tx_do_callbacks(&tx->tx_callbacks, ECANCELED);
1210     list_destroy(&tx->tx_callbacks);
1211     list_destroy(&tx->tx_holds);
1212 #ifdef ZFS_DEBUG
1213     refcount_destroy_many(&tx->tx_space_written,
1214                           refcount_count(&tx->tx_space_written));
1215     refcount_destroy_many(&tx->tx_space_freed,
1216                           refcount_count(&tx->tx_space_freed));
1217 #endif
1218     kmem_free(tx, sizeof (dmu_tx_t));
1219 }

1221 uint64_t
1222 dmu_tx_get_txg(dmu_tx_t *tx)
1223 {
1224     ASSERT(tx->tx_txg != 0);
1225     return (tx->tx_txg);
1226 }

1228 void
1229 dmu_tx_callback_register(dmu_tx_t *tx, dmu_tx_callback_func_t *func, void *data)
1230 {
1231     dmu_tx_callback_t *dcb;
1233     dcb = kmem_alloc(sizeof (dmu_tx_callback_t), KM_SLEEP);
1235     dcb->dcb_func = func;
1236     dcb->dcb_data = data;
1238     list_insert_tail(&tx->tx_callbacks, dcb);
1239 }

1241 /*
1242  * Call all the commit callbacks on a list, with a given error code.
1243  */
1244 void
1245 dmu_tx_do_callbacks(list_t *cb_list, int error)
1246 {
1247     dmu_tx_callback_t *dcb;
1249     while (dcb = list_head(cb_list)) {
1250         list_remove(cb_list, dcb);
1251         dcb->dcb_func(dcb->dcb_data, error);
1252     }
1253 }
```

```
new/usr/src/uts/common/fs/zfs/dmu_tx.c

1252         kmem_free(dcb, sizeof (dmu_tx_callback_t));
1253     }
1254 }

1255 /*
1256  * Interface to hold a bunch of attributes.
1257  * used for creating new files.
1258  * attrsizes is the total size of all attributes
1259  * to be added during object creation
1260  *
1261  */
1262  * For updating/adding a single attribute dmu_tx_hold_sa() should be used.
1263 */

1264 /*
1265  * hold necessary attribute name for attribute registration.
1266  * should be a very rare case where this is needed. If it does
1267  * happen it would only happen on the first write to the file system.
1268  */
1269 static void
1270 dmu_tx_sa_registration_hold(sa_os_t *sa, dmu_tx_t *tx)
1271 {
1272     int i;
1273
1274     if (!sa->sa_need_attr_registration)
1275         return;
1276
1277     for (i = 0; i != sa->sa_num_attrs; i++) {
1278         if (!sa->sa_attr_table[i].sa_registered) {
1279             if (sa->sa_reg_attr_obj)
1280                 dmu_tx_hold_zap(tx, sa->sa_reg_attr_obj,
1281                                 B_TRUE, sa->sa_attr_table[i].sa_name);
1282             else
1283                 dmu_tx_hold_zap(tx, DMU_NEW_OBJECT,
1284                                 B_TRUE, sa->sa_attr_table[i].sa_name);
1285         }
1286     }
1287 }
1288 }

1289 void
1290 dmu_tx_hold_spill(dmu_tx_t *tx, uint64_t object)
1291 {
1292     dnode_t *dn;
1293     dmu_tx_hold_t *txh;
1294
1295     txh = dmu_tx_hold_object_impl(tx, tx->tx_objset, object,
1296                                  THT_SPILL, 0, 0);
1297
1298     dn = txh->txh_dnode;
1299
1300     if (dn == NULL)
1301         return;
1302
1303     /* If blkptr doesn't exist then add space to towrite */
1304     if (!(dn->dn_phys->dn_flags & DNODE_FLAG_SPILL_BLKPTR))
1305         txh->txh_space_towrite += SPA_MAXBLOCKSIZE;
1306     } else {
1307         blkptr_t *bp;
1308
1309         bp = &dn->dn_phys->dn_spill;
1310         if (dsl_dataset_block_freeable(dn->dn_objset->os_dsl_dataset,
1311                                       bp, bp->blk_birth))
1312             txh->txh_space_tooverwrite += SPA_MAXBLOCKSIZE;
1313         else
1314             txh->txh_space_towrite += SPA_MAXBLOCKSIZE;
1315         if (!BP_IS_WOLF(bp))
1316             if (BP_IS_WOLF(bp))
```

```

1318             txh->txh_space_tounref += SPA_MAXBLOCKSIZE;
1319     }
1320 }
1322 void
1323 dmu_tx_hold_sa_create(dmu_tx_t *tx, int attrszie)
1324 {
1325     sa_os_t *sa = tx->tx_objset->os_sa;
1327     dmu_tx_hold_bonus(tx, DMU_NEW_OBJECT);
1329     if (tx->tx_objset->os_sa->sa_master_obj == 0)
1330         return;
1332     if (tx->tx_objset->os_sa->sa_layout_attr_obj)
1333         dmu_tx_hold_zap(tx, sa->sa_layout_attr_obj, B_TRUE, NULL);
1334     else {
1335         dmu_tx_hold_zap(tx, sa->sa_master_obj, B_TRUE, SA_LAYOUTS);
1336         dmu_tx_hold_zap(tx, sa->sa_master_obj, B_TRUE, SA_REGISTRY);
1337         dmu_tx_hold_zap(tx, DMU_NEW_OBJECT, B_TRUE, NULL);
1338         dmu_tx_hold_zap(tx, DMU_NEW_OBJECT, B_TRUE, NULL);
1339     }
1341     dmu_tx_sa_registration_hold(sa, tx);
1343     if (attrszie <= DN_MAX_BONUSLEN && !sa->sa_force_spill)
1344         return;
1346     (void) dmu_tx_hold_object_impl(tx, tx->tx_objset, DMU_NEW_OBJECT,
1347         THT_SPILL, 0, 0);
1348 }

1350 /*
1351 * Hold SA attribute
1352 *
1353 * dmu_tx_hold_sa(dmu_tx_t *tx, sa_handle_t *, attribute, add, size)
1354 *
1355 * variable_size is the total size of all variable sized attributes
1356 * passed to this function. It is not the total size of all
1357 * variable size attributes that *may* exist on this object.
1358 */
1359 void
1360 dmu_tx_hold_sa(dmu_tx_t *tx, sa_handle_t *hdl, boolean_t may_grow)
1361 {
1362     uint64_t object;
1363     sa_os_t *sa = tx->tx_objset->os_sa;
1365     ASSERT(hdl != NULL);
1367     object = sa_handle_object(hdl);
1369     dmu_tx_hold_bonus(tx, object);
1371     if (tx->tx_objset->os_sa->sa_master_obj == 0)
1372         return;
1374     if (tx->tx_objset->os_sa->sa_reg_attr_obj == 0 ||
1375         tx->tx_objset->os_sa->sa_layout_attr_obj == 0) {
1376         dmu_tx_hold_zap(tx, sa->sa_master_obj, B_TRUE, SA_LAYOUTS);
1377         dmu_tx_hold_zap(tx, sa->sa_master_obj, B_TRUE, SA_REGISTRY);
1378         dmu_tx_hold_zap(tx, DMU_NEW_OBJECT, B_TRUE, NULL);
1379         dmu_tx_hold_zap(tx, DMU_NEW_OBJECT, B_TRUE, NULL);
1380     }
1382     dmu_tx_sa_registration_hold(sa, tx);

```

```

1384     if (may_grow && tx->tx_objset->os_sa->sa_layout_attr_obj)
1385         dmu_tx_hold_zap(tx, sa->sa_layout_attr_obj, B_TRUE, NULL);
1387     if (sa->sa_force_spill || may_grow || hdl->sa_spill) {
1388         ASSERT(tx->tx_txg == 0);
1389         dmu_tx_hold_spill(tx, object);
1390     } else {
1391         dmu_buf_impl_t *db = (dmu_buf_impl_t *)hdl->sa_bonus;
1392         dnode_t *dn;
1394         DB_DNODE_ENTER(db);
1395         dn = DB_DNODE(db);
1396         if (dn->dn_have_spill) {
1397             ASSERT(tx->tx_txg == 0);
1398             dmu_tx_hold_spill(tx, object);
1399         }
1400     }
1401     DB_DNODE_EXIT(db);
1402 }

```