

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

```
*****
82421 Thu Aug 15 17:31:03 2013
new/usr/src/uts/common/fs/zfs/dsl_dataset.c
4046 dsl_dataset_t ds_dir->dd_lock is highly contended
Reviewed by: Eric Schrock <eric.schrock@delphix.com>
Reviewed by: George Wilson <george.wilson@delphix.com>
*****
unchanged_portion_omitted_
```

```
82 void
83 dsl_dataset_block_born(dsl_dataset_t *ds, const blkptr_t *bp, dmu_tx_t *tx)
84 {
85     int used = bp_get_dsize_sync(tx->tx_pool->dp_spa, bp);
86     int compressed = BP_GET_PSIZE(bp);
87     int uncompressed = BP_GET_UCSIZE(bp);
88     int64_t delta;
89
90     dprintf_bp(bp, "ds=%p", ds);
91
92     ASSERT(dmu_tx_is_syncing(tx));
93     /* It could have been compressed away to nothing */
94     if (BP_IS_HOLE(bp))
95         return;
96     ASSERT(BP_GET_TYPE(bp) != DMU_OT_NONE);
97     ASSERT(DMU_OT_IS_VALID(BP_GET_TYPE(bp)));
98     if (ds == NULL) {
99         dsl_pool_mos_diduse_space(tx->tx_pool,
100             used, compressed, uncompressed);
101         return;
102     }
103
104     dmu_buf_will_dirty(ds->dsdbuf, tx);
105
106     mutex_enter(&ds->ds_dir->dd_lock);
107     mutex_enter(&ds->ds_lock);
108     delta = parent_delta(ds, used);
109     ds->ds_phys->ds_referenced_bytes += used;
110     ds->ds_phys->ds_compressed_bytes += compressed;
111     ds->ds_phys->ds_uncompressed_bytes += uncompressed;
112     ds->ds_phys->ds_unique_bytes += used;
113     mutex_exit(&ds->ds_lock);
114     dsl_dir_diduse_space(ds->ds_dir, DD_USED_HEAD, delta,
115         compressed, uncompressed, tx);
116     dsl_dir_transfer_space(ds->ds_dir, used - delta,
117         DD_USED_REFRSRV, DD_USED_HEAD, tx);
118 }
119
120 int
121 dsl_dataset_block_kill(dsl_dataset_t *ds, const blkptr_t *bp, dmu_tx_t *tx,
122     boolean_t async)
123 {
124     if (BP_IS_HOLE(bp))
125         return (0);
126
127     ASSERT(dmu_tx_is_syncing(tx));
128     ASSERT(bp->blk_birth <= tx->tx_txg);
129
130     int used = bp_get_dsize_sync(tx->tx_pool->dp_spa, bp);
131     int compressed = BP_GET_PSIZE(bp);
132     int uncompressed = BP_GET_UCSIZE(bp);
133
134     ASSERT(used > 0);
135     if (ds == NULL) {
136         dsl_free(tx->tx_pool, tx->tx_txg, bp);
137         dsl_pool_mos_diduse_space(tx->tx_pool,
```

1

```
136         -used, -compressed, -uncompressed);
137         return (used);
138     }
139     ASSERT3P(tx->tx_pool, ==, ds->ds_dir->dd_pool);
140
141     ASSERT(!dsl_dataset_is_snapshot(ds));
142     dmu_buf_will_dirty(ds->dsdbuf, tx);
143
144     if (bp->blk_birth > ds->ds_phys->ds_prev_snap_txg) {
145         int64_t delta;
146
147         dprintf_bp(bp, "freeing ds=%llu", ds->ds_object);
148         dsl_free(tx->tx_pool, tx->tx_txg, bp);
149
150         mutex_enter(&ds->ds_dir->dd_lock);
151         mutex_enter(&ds->ds_lock);
152         ASSERT(ds->ds_phys->ds_unique_bytes >= used ||
153             !DS_UNIQUE_IS_ACCURATE(ds));
154         delta = parent_delta(ds, -used);
155         ds->ds_phys->ds_unique_bytes -= used;
156         mutex_exit(&ds->ds_lock);
157         dsl_dir_diduse_space(ds->ds_dir, DD_USED_HEAD,
158             delta, -compressed, -uncompressed, tx);
159         dsl_dir_transfer_space(ds->ds_dir, -used - delta,
160             DD_USED_REFRSRV, DD_USED_HEAD, tx);
161         mutex_exit(&ds->ds_dir->dd_lock);
162     } else {
163         dprintf_bp(bp, "putting on dead list: %s", "");
164         if (async) {
165             /*
166             * We are here as part of zio's write done callback,
167             * which means we're a zio interrupt thread. We can't
168             * call dsl_deadlist_insert() now because it may block
169             * waiting for I/O. Instead, put bp on the deferred
170             * queue and let dsl_pool_sync() finish the job.
171             */
172             bplist_append(&ds->ds_pending_deadlist, bp);
173         } else {
174             dsl_deadlist_insert(&ds->ds_deadlist, bp, tx);
175         }
176         ASSERT3U(ds->ds_prev->ds_object, ==,
177             ds->ds_phys->ds_prev_snap_obj);
178         ASSERT(ds->ds_prev->ds_phys->ds_num_children > 0);
179         /* if (bp->blk_birth > prev prev snap txg) prev unique += bs */
180         if (ds->ds_prev->ds_phys->ds_next_snap_obj ==
181             ds->ds_object && bp->blk_birth >
182             ds->ds_phys->ds_prev_snap_txg) {
183             dmu_buf_will_dirty(ds->ds_phys->dsdbuf, tx);
184             mutex_enter(&ds->ds_phys->ds_phys->ds_lock);
185             ds->ds_phys->ds_phys->ds_unique_bytes += used;
186             mutex_exit(&ds->ds_phys->ds_phys->ds_lock);
187         }
188         if (bp->blk_birth > ds->ds_dir->dd_origin_txg) {
189             dsl_dir_transfer_space(ds->ds_dir, used,
190             DD_USED_HEAD, DD_USED_SNAP, tx);
191         }
192         mutex_enter(&ds->ds_lock);
193         ASSERT3U(ds->ds_phys->ds_referenced_bytes, >=, used);
194         ds->ds_phys->ds_referenced_bytes -= used;
195         ASSERT3U(ds->ds_phys->ds_compressed_bytes, >=, compressed);
196         ds->ds_phys->ds_compressed_bytes -= compressed;
197         ASSERT3U(ds->ds_phys->ds_uncompressed_bytes, >=, uncompressed);
198         ds->ds_phys->ds_uncompressed_bytes -= uncompressed;
199         mutex_exit(&ds->ds_lock);
200     }
```

2

```
200         return (used);
201 }


---

unchanged portion omitted
594 static int
595 dsl_dataset_nameLEN(dsl_dataset_t *ds)
596 {
597     int result;
598
599     if (ds == NULL) {
600         result = 3;      /* "mos" */
601     } else {
602         result = dsl_dir_nameLEN(ds->ds_dir);
603         VERIFY0(dsl_dataset_get_snapname(ds));
604         if (ds->ds_snapname[0]) {
605             ++result; /* adding one for the @-sign */
606             if (!MUTEX_HELD(&ds->ds_lock)) {
607                 mutex_enter(&ds->ds_lock);
608                 result += strlen(ds->ds_snapname);
609                 mutex_exit(&ds->ds_lock);
610             } else {
611                 result += strlen(ds->ds_snapname);
612             }
613         }
614     }
615
616     return (result);
617 }


---


590 void
591 dsl_dataset_rele(dsl_dataset_t *ds, void *tag)
592 {
593     dmuf_buf_rele(ds->ds_dbuf, tag);
594 }


---

unchanged portion omitted
```

```
*****
35916 Thu Aug 15 17:31:07 2013
new/usr/src/uts/common/fs/zfs/dsl_dir.c
4046 dsl_dataset_t ds_dir->dd_lock is highly contended
Reviewed by: Eric Schrock <eric.schrock@delphix.com>
Reviewed by: George Wilson <george.wilson@delphix.com>
*****
unchanged_portion_omitted

838 /* call from syncing context when we actually write/free space for this dd */
839 void
840 dsl_dir_diduse_space(dsl_dir_t *dd, dd_used_t type,
841     int64_t used, int64_t compressed, int64_t uncompressed, dmu_tx_t *tx)
842 {
843     int64_t accounted_delta;

845     /*
846      * dsl_dataset_set_reservation_sync_impl() calls this with
847      * dd_lock held, so that it can atomically update
848      * ds->ds_reserved and the dsl_dir accounting, so that
849      * dsl_dataset_check_quota() can see dataset and dir accounting
850      * consistently.
851     */
852     boolean_t needlock = !MUTEX_HELD(&dd->dd_lock);

854     ASSERT(dmu_tx_is_syncing(tx));
855     ASSERT(type < DD_USED_NUM);

857     dmu_buf_will_dirty(dd->dd_dbuf, tx);

859     if (needlock)
860         mutex_enter(&dd->dd_lock);
861     accounted_delta = parent_delta(dd, dd->dd_phys->dd_used_bytes, used);
862     ASSERT(used >= 0 || dd->dd_phys->dd_used_bytes >= -used);
863     ASSERT(compressed >= 0 ||
864             dd->dd_phys->dd_compressed_bytes >= -compressed);
865     ASSERT(uncompressed >= 0 ||
866             dd->dd_phys->dd_uncompressed_bytes >= -uncompressed);
857     dmu_buf_will_dirty(dd->dd_dbuf, tx);
867     dd->dd_phys->dd_used_bytes += used;
868     dd->dd_phys->dd_uncompressed_bytes += uncompressed;
869     dd->dd_phys->dd_compressed_bytes += compressed;

871     if (dd->dd_phys->dd_flags & DD_FLAG_USED_BREAKDOWN) {
872         ASSERT(used > 0 ||
873                 dd->dd_phys->dd_used_breakdown[type] >= -used);
874         dd->dd_phys->dd_used_breakdown[type] += used;
875 #ifdef DEBUG
876         dd_used_t t;
877         uint64_t u = 0;
878         for (t = 0; t < DD_USED_NUM; t++)
879             u += dd->dd_phys->dd_used_breakdown[t];
880         ASSERT3U(u, ==, dd->dd_phys->dd_used_bytes);
881 #endif
882     }
883     if (needlock)
884         mutex_exit(&dd->dd_lock);

886     if (dd->dd_parent != NULL) {
887         dsl_dir_diduse_space(dd->dd_parent, DD_USED_CHILD,
888             accounted_delta, compressed, uncompressed, tx);
889         dsl_dir_transfer_space(dd->dd_parent,
890             used - accounted_delta,
891             DD_USED_CHILD_RSRV, DD_USED_CHILD, tx);
892     }
893 }
```

```
895 void
896 dsl_dir_transfer_space(dsl_dir_t *dd, int64_t delta,
897     dd_used_t oldtype, dd_used_t newtype, dmu_tx_t *tx)
898 {
899     boolean_t needlock = !MUTEX_HELD(&dd->dd_lock);
900
901     ASSERT(dmu_tx_is_syncing(tx));
902     ASSERT(oldtype < DD_USED_NUM);
903     ASSERT(newtype < DD_USED_NUM);
904
905     if (delta == 0 || !(dd->dd_phys->dd_flags & DD_FLAG_USED_BREAKDOWN))
906         return;
907
908     dmu_buf_will_dirty(dd->dd_dbuf, tx);
909     if (needlock)
910         mutex_enter(&dd->dd_lock);
911     ASSERT(delta > 0 ?
912             dd->dd_phys->dd_used_breakdown[oldtype] >= delta :
913             dd->dd_phys->dd_used_breakdown[newtype] >= -delta);
914     ASSERT((dd->dd_phys->dd_used_bytes >= ABS(delta));
915     dmu_buf_will_dirty(dd->dd_dbuf, tx);
916     dd->dd_phys->dd_used_breakdown[oldtype] -= delta;
917     dd->dd_phys->dd_used_breakdown[newtype] += delta;
918     if (needlock)
919         mutex_exit(&dd->dd_lock);
920 }
```

unchanged_portion_omitted