

```

*****
74754 Thu Jan 10 23:20:01 2013
new/usr/src/uts/common/fs/zfs/dbuf.c
3469 dbuf_read_impl shows too much enthusiasm
Reviewed by: Bryan Cantrill <bryan@joyent.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 2011 Nexenta Systems, Inc. All rights reserved.
24 * Copyright (c) 2012 by Delphix. All rights reserved.
25 * Copyright (c) 2013, Joyent, Inc. All rights reserved.
26 */

28 #include <sys/zfs_context.h>
29 #include <sys/dmu.h>
30 #include <sys/dmu_impl.h>
31 #include <sys/dbuf.h>
32 #include <sys/dmu_objset.h>
33 #include <sys/dsl_dataset.h>
34 #include <sys/dsl_dir.h>
35 #include <sys/dmu_tx.h>
36 #include <sys/spa.h>
37 #include <sys/zio.h>
38 #include <sys/dmu_zfetch.h>
39 #include <sys/sa.h>
40 #include <sys/sa_impl.h>

42 static void dbuf_destroy(dmu_buf_impl_t *db);
43 static int dbuf_undirty(dmu_buf_impl_t *db, dmu_tx_t *tx);
44 static void dbuf_write(dbuf_dirty_record_t *dr, arc_buf_t *data, dmu_tx_t *tx);

46 /*
47  * Global data structures and functions for the dbuf cache.
48  */
49 static kmem_cache_t *dbuf_cache;

51 /* ARGSUSED */
52 static int
53 dbuf_cons(void *vdb, void *unused, int kmflag)
54 {
55     dmu_buf_impl_t *db = vdb;
56     bzero(db, sizeof (dmu_buf_impl_t));

58     mutex_init(&db->db_mtx, NULL, MUTEX_DEFAULT, NULL);
59     cv_init(&db->db_changed, NULL, CV_DEFAULT, NULL);
60     refcount_create(&db->db_holds);

```

```

61         return (0);
62     }
    unchanged_portion_omitted

510 static void
511 dbuf_read_impl(dmu_buf_impl_t *db, zio_t *zio, uint32_t *flags)
512 {
513     dnode_t *dn;
514     spa_t *spa;
515     zbookmark_t zb;
516     uint32_t aflags = ARC_NOWAIT;
517     arc_buf_t *pbuf;

519     DB_DNODE_ENTER(db);
520     dn = DB_DNODE(db);
521     ASSERT(!refcount_is_zero(&db->db_holds));
522     /* We need the struct_rwlock to prevent db_blkptr from changing. */
523     ASSERT(RW_LOCK_HELD(&dn->dn_struct_rwlock));
524     ASSERT(MUTEX_HELD(&db->db_mtx));
525     ASSERT(db->db_state == DB_UNCACHED);
526     ASSERT(db->db_buf == NULL);

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

531         ASSERT3U(bonuslen, <=, db->db.db_size);
532         db->db.db_data = zio_buf_alloc(DN_MAX_BONUSLEN);
533         arc_space_consume(DN_MAX_BONUSLEN, ARC_SPACE_OTHER);
534         if (bonuslen < DN_MAX_BONUSLEN)
535             bzero(db->db.db_data, DN_MAX_BONUSLEN);

537         if (bonuslen) {
538             /*
539              * Absent byzantine on-disk corruption, we fully expect
540              * our bonuslen to be no more than DN_MAX_BONUSLEN --
541              * but we nonetheless explicitly clamp it on the bcopy()
542              * to prevent any on-disk corruption from becoming
543              * rampant in-kernel corruption.
544              */
545             bcopy(DN_BONUS(dn->dn_phys), db->db.db_data,
546                 MIN(bonuslen, DN_MAX_BONUSLEN));
547         }

535         if (bonuslen)
536             bcopy(DN_BONUS(dn->dn_phys), db->db.db_data, bonuslen);
549         DB_DNODE_EXIT(db);
550         dbuf_update_data(db);
551         db->db_state = DB_CACHED;
552         mutex_exit(&db->db_mtx);
553         return;
554     }

556     /*
557      * Recheck BP_IS_HOLE() after dnode_block_freed() in case dnode_sync()
558      * processes the delete record and clears the bp while we are waiting
559      * for the dn_mtx (resulting in a "no" from block_freed).
560      */
561     if (db->db_blkptr == NULL || BP_IS_HOLE(db->db_blkptr) ||
562         (db->db_level == 0 && (dnode_block_freed(dn, db->db_blkid) ||
563             BP_IS_HOLE(db->db_blkptr)))) {
564         arc_buf_contents_t type = DBUF_GET_BUFC_TYPE(db);

566         dbuf_set_data(db, arc_buf_alloc(dn->dn_objset->os_spa,
567             db->db.db_size, db, type));
568         DB_DNODE_EXIT(db);
569         bzero(db->db.db_data, db->db.db_size);

```

```
570         db->db_state = DB_CACHED;
571         *flags |= DB_RF_CACHED;
572         mutex_exit(&db->db_mtx);
573         return;
574     }

576     spa = dn->dn_objset->os_spa;
577     DB_DNODE_EXIT(db);

579     db->db_state = DB_READ;
580     mutex_exit(&db->db_mtx);

582     if (DBUF_IS_L2CACHEABLE(db))
583         aflags |= ARC_L2CACHE;

585     SET_BOOKMARK(&zsb, db->db_objset->os_dsl_dataset ?
586         db->db_objset->os_dsl_dataset->ds_object : DMU_META_OBJSET,
587         db->db.db_object, db->db_level, db->db_blkid);

589     dbuf_add_ref(db, NULL);
590     /* ZIO_FLAG_CANFAIL callers have to check the parent zio's error */

592     if (db->db_parent)
593         pbuf = db->db_parent->db_buf;
594     else
595         pbuf = db->db_objset->os_phys_buf;

597     (void) dsl_read(zio, spa, db->db_blkptr, pbuf,
598         dbuf_read_done, db, ZIO_PRIORITY_SYNC_READ,
599         (*flags & DB_RF_CANFAIL) ? ZIO_FLAG_CANFAIL : ZIO_FLAG_MUSTSUCCEED,
600         &aflags, &zsb);
601     if (aflags & ARC_CACHED)
602         *flags |= DB_RF_CACHED;
603 }
```

unchanged portion omitted