

```

*****
90876 Tue Oct 1 14:04:06 2013
new/usr/src/cmd/zdb/zdb.c
4171 clean up spa_feature_*() interfaces
4172 implement extensible_dataset feature for use by other zpool features
Reviewed by: Max Grossman <max.grossman@delphix.com>
Reviewed by: Christopher Siden <christopher.siden@delphix.com>
Reviewed by: George Wilson <george.wilson@delphix.com>
*****
_____unchanged_portion_omitted_____

563 static int
564 verify_spacemap_refcounts(spa_t *spa)
565 {
566     uint64_t expected_refcount = 0;
567     uint64_t actual_refcount;
568     int expected_refcount, actual_refcount;

569     (void) feature_get_refcount(spa,
570     &spa_feature_table[SPA_FEATURE_SPACEMAP_HISTOGRAM],
571     &expected_refcount);
572     expected_refcount = spa_feature_get_refcount(spa,
573     &spa_feature_table[SPA_FEATURE_SPACEMAP_HISTOGRAM]);
574     actual_refcount = get_dtl_refcount(spa->spa_root_vdev);
575     actual_refcount += get metaslab_refcount(spa->spa_root_vdev);

576     if (expected_refcount != actual_refcount) {
577         (void) printf("space map refcount mismatch: expected %lld != "
578         "actual %lld\n",
579         (longlong_t)expected_refcount,
580         (longlong_t)actual_refcount);
581         (void) printf("space map refcount mismatch: expected %d != "
582         "actual %d\n", expected_refcount, actual_refcount);
583     }
584     return (2);
585 }
586 return (0);
_____unchanged_portion_omitted_____

653 static void
654 dump metaslab(metaslab_t *msp)
655 {
656     vdev_t *vd = msp->ms_group->mg_vd;
657     spa_t *spa = vd->vdev_spa;
658     space_map_t *sm = msp->ms_sm;
659     char freebuf[32];

661     zdb_nicenum(msp->ms_size - space_map_allocated(sm), freebuf);

663     (void) printf(
664     "\tmetaslab %6llu offset %12llx spacemap %6llu free %5s\n",
665     (u_longlong_t)msp->ms_id, (u_longlong_t)msp->ms_start,
666     (u_longlong_t)space_map_object(sm), freebuf);

668     if (dump_opt['m'] > 2 && !dump_opt['L']) {
669         mutex_enter(&msp->ms_lock);
670         metaslab_load_wait(msp);
671         if (!msp->ms_loaded) {
672             VERIFY0(metaslab_load(msp));
673             range_tree_stat_verify(msp->ms_tree);
674         }
675         dump metaslab_stats(msp);
676         metaslab_unload(msp);
677         mutex_exit(&msp->ms_lock);
678     }

```

```

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

693     if (dump_opt['d'] > 5 || dump_opt['m'] > 3) {
694         ASSERT(msp->ms_size == (LULL << vd->vdev_ms_shift));
695         mutex_enter(&msp->ms_lock);
696         dump_spacemap(spa->spa_meta_objset, msp->ms_sm);
697         mutex_exit(&msp->ms_lock);
698     }
_____unchanged_portion_omitted_____

2437 static int
2438 dump_block_stats(spa_t *spa)
2439 {
2440     zdb_cb_t zcb = { 0 };
2441     zdb_blkstats_t *zb, *tzb;
2442     uint64_t norm_alloc, norm_space, total_alloc, total_found;
2443     int flags = TRAVERSE_PRE | TRAVERSE_PREFETCH_METADATA | TRAVERSE_HARD;
2444     int leaks = 0;

2446     (void) printf("\nTraversing all blocks %s%s%s%s..\n\n",
2447     (dump_opt['c'] || !dump_opt['L']) ? "to verify " : "",
2448     (dump_opt['c'] == 1) ? "metadata " : "",
2449     (dump_opt['c'] ? "checksums " : "",
2450     (dump_opt['c'] && !dump_opt['L']) ? "and verify " : "",
2451     !dump_opt['L'] ? "nothing leaked " : "");

2453     /*
2454     * Load all space maps as SM_ALLOC maps, then traverse the pool
2455     * claiming each block we discover. If the pool is perfectly
2456     * consistent, the space maps will be empty when we're done.
2457     * Anything left over is a leak; any block we can't claim (because
2458     * it's not part of any space map) is a double allocation,
2459     * reference to a freed block, or an unclaimed log block.
2460     */
2461     zdb_leak_init(spa, &zcb);

2463     /*
2464     * If there's a deferred-free bplist, process that first.
2465     */
2466     (void) bpobj_iterate_nofree(&spa->spa_deferred_bpobj,
2467     count_block_cb, &zcb, NULL);
2468     if (spa_version(spa) >= SPA_VERSION_DEADLISTS) {
2469         (void) bpobj_iterate_nofree(&spa->spa_dsl_pool->dp_free_bpobj,
2470         count_block_cb, &zcb, NULL);
2471     }
2472     if (spa_feature_is_active(spa, SPA_FEATURE_ASYNC_DESTROY)) {
2473         if (spa_feature_is_active(spa,
2474         &spa_feature_table[SPA_FEATURE_ASYNC_DESTROY])) {
2475             VERIFY3U(0, ==, bptree_iterate(spa->spa_meta_objset,
2476             spa->spa_dsl_pool->dp_bptree_obj, B_FALSE, count_block_cb,
2477             &zcb, NULL));
2478         }

```



```

2610             continue;
2611
2612             zdb_nicenum(zb->zb_count, csize);
2613             zdb_nicenum(zb->zb_lsize, lsize);
2614             zdb_nicenum(zb->zb_psize, psize);
2615             zdb_nicenum(zb->zb_asize, asize);
2616             zdb_nicenum(zb->zb_asize / zb->zb_count, avg);
2617             zdb_nicenum(zb->zb_gangs, gang);
2618
2619             (void) printf("%6s\t%5s\t%5s\t%5s\t%5s"
2620                "\t%5.2f\t%6.2f\t",
2621                csize, lsize, psize, asize, avg,
2622                (double)zb->zb_lsize / zb->zb_psize,
2623                100.0 * zb->zb_asize / tzb->zb_asize);
2624
2625             if (level == ZB_TOTAL)
2626                 (void) printf("%s\n", typename);
2627             else
2628                 (void) printf("    L%d %s\n",
2629                    level, typename);
2630
2631             if (dump_opt['b'] >= 3 && zb->zb_gangs > 0) {
2632                 (void) printf("\t number of ganged "
2633                    "blocks: %s\n", gang);
2634             }
2635
2636             if (dump_opt['b'] >= 4) {
2637                 (void) printf("psize "
2638                    "(in 512-byte sectors): "
2639                    "number of blocks\n");
2640                 dump_histogram(zb->zb_psize_histogram,
2641                    PSIZE_HISTO_SIZE, 0);
2642             }
2643         }
2644     }
2645 }
2646
2647 (void) printf("\n");
2648
2649 if (leaks)
2650     return (2);
2651
2652 if (zcb.zcb_haderrors)
2653     return (3);
2654
2655 return (0);
2656 }

```

unchanged\_portion\_omitted

```

2759 static void
2760 dump_zpool(spa_t *spa)
2761 {
2762     dsl_pool_t *dp = spa_get_dsl(spa);
2763     int rc = 0;
2764
2765     if (dump_opt['S']) {
2766         dump_simulated_ddt(spa);
2767         return;
2768     }
2769
2770     if (!dump_opt['e'] && dump_opt['C'] > 1) {
2771         (void) printf("\nCached configuration:\n");
2772         dump_nvlist(spa->spa_config, 8);
2773     }
2774
2775     if (dump_opt['C'])

```

```

2776         dump_config(spa);
2777
2778         if (dump_opt['u'])
2779             dump_uberblock(&spa->spa_uberblock, "\nUberblock:\n", "\n");
2780
2781         if (dump_opt['D'])
2782             dump_all_ddts(spa);
2783
2784         if (dump_opt['d'] > 2 || dump_opt['m'])
2785             dump metaslabs(spa);
2786
2787         if (dump_opt['d'] || dump_opt['i']) {
2788             dump_dir(dp->dp_meta_objset);
2789             if (dump_opt['d'] >= 3) {
2790                 dump_bpobj(&spa->spa_deferred_bpobj,
2791                    "Deferred frees", 0);
2792                 if (spa_version(spa) >= SPA_VERSION_DEADLISTS) {
2793                     dump_bpobj(&spa->spa_dsl_pool->dp_free_bpobj,
2794                        "Pool snapshot frees", 0);
2795                 }
2796             }
2797             if (spa_feature_is_active(spa,
2798                SPA_FEATURE_ASYNC_DESTROY)) {
2799                 &spa_feature_table[SPA_FEATURE_ASYNC_DESTROY]) {
2800                     dump_bptree(spa->spa_meta_objset,
2801                        spa->spa_dsl_pool->dp_bptree_obj,
2802                        "Pool dataset frees");
2803                 }
2804             }
2805             dump_dtl(spa->spa_root_vdev, 0);
2806             (void) dmu_objset_find(spa_name(spa), dump_one_dir,
2807                NULL, DS_FIND_SNAPSHOTS | DS_FIND_CHILDREN);
2808         }
2809         if (dump_opt['b'] || dump_opt['c'])
2810             rc = dump_block_stats(spa);
2811
2812         if (rc == 0)
2813             rc = verify_spacemap_refcounts(spa);
2814
2815         if (dump_opt['s'])
2816             show_pool_stats(spa);
2817
2818         if (dump_opt['h'])
2819             dump_history(spa);
2820
2821         if (rc != 0)
2822             exit(rc);
2823     }

```

unchanged\_portion\_omitted

```

*****
13331 Tue Oct 1 14:04:07 2013
new/usr/src/cmd/zhack/zhack.c
4171 clean up spa_feature_*( ) interfaces
4172 implement extensible_dataset feature for use by other zpool features
Reviewed by: Max Grossman <max.grossman@delphix.com>
Reviewed by: Christopher Siden <christopher.siden@delphix.com>
Reviewed by: George Wilson <george.wilson@delphix.com>
*****
_____unchanged_portion_omitted_____

285 static void
286 zhack_feature_enable_sync(void *arg, dmu_tx_t *tx)
287 feature_enable_sync(void *arg, dmu_tx_t *tx)
288 {
289     spa_t *spa = dmu_tx_pool(tx)->dp_spa;
290     zfeature_info_t *feature = arg;
291     feature_enable_sync(spa, feature, tx);
292
293     spa_feature_enable(spa, feature, tx);
294     spa_history_log_internal(spa, "zhack enable feature", tx,
295         "guid=%s can_readonly=%u",
296         feature->fi_guid, feature->fi_can_readonly);
297 }
298
299 static void
300 zhack_do_feature_enable(int argc, char **argv)
301 {
302     char c;
303     char *desc, *target;
304     spa_t *spa;
305     objset_t *mos;
306     zfeature_info_t feature;
307     spa_feature_t nodeps[] = { SPA_FEATURE_NONE };
308     zfeature_info_t *nodeps[] = { NULL };
309
310     /*
311      * Features are not added to the pool's label until their refcounts
312      * are incremented, so fi_mos can just be left as false for now.
313      */
314     desc = NULL;
315     feature.fi_uname = "zhack";
316     feature.fi_mos = B_FALSE;
317     feature.fi_can_readonly = B_FALSE;
318     feature.fi_depends = nodeps;
319
320     optind = 1;
321     while ((c = getopt(argc, argv, "rmd:")) != -1) {
322         switch (c) {
323             case 'r':
324                 feature.fi_can_readonly = B_TRUE;
325                 break;
326             case 'd':
327                 desc = strdup(optarg);
328                 break;
329             default:
330                 usage();
331                 break;
332         }
333     }
334
335     if (desc == NULL)
336         desc = strdup("zhack injected");
337     feature.fi_desc = desc;

```

```

337     argc -= optind;
338     argv += optind;
339
340     if (argc < 2) {
341         (void) fprintf(stderr, "error: missing feature or pool name\n");
342         usage();
343     }
344     target = argv[0];
345     feature.fi_guid = argv[1];
346
347     if (!zfeature_is_valid_guid(feature.fi_guid))
348         fatal(NULL, FTAG, "invalid feature guid: %s", feature.fi_guid);
349
350     zhack_spa_open(target, B_FALSE, FTAG, &spa);
351     mos = spa->spa_meta_objset;
352
353     if (zfeature_is_supported(feature.fi_guid))
354         if (0 == zfeature_lookup_guid(feature.fi_guid, NULL))
355             fatal(spa, FTAG, "%s' is a real feature, will not enable");
356     if (0 == zap_contains(mos, spa->spa_feat_desc_obj, feature.fi_guid))
357         fatal(spa, FTAG, "feature already enabled: %s",
358             feature.fi_guid);
359
360     VERIFY0(dsl_sync_task(spa_name(spa), NULL,
361         zhack_feature_enable_sync, &feature, 5));
362     feature_enable_sync, &feature, 5));
363
364     spa_close(spa, FTAG);
365     free(desc);
366 }
367
368 static void
369 feature_incr_sync(void *arg, dmu_tx_t *tx)
370 {
371     spa_t *spa = dmu_tx_pool(tx)->dp_spa;
372     zfeature_info_t *feature = arg;
373     uint64_t refcount;
374
375     VERIFY0(feature_get_refcount(spa, feature, &refcount));
376     feature_sync(spa, feature, refcount + 1, tx);
377     spa_feature_incr(spa, feature, tx);
378     spa_history_log_internal(spa, "zhack feature incr", tx,
379         "guid=%s", feature->fi_guid);
380 }
381
382 static void
383 feature_decr_sync(void *arg, dmu_tx_t *tx)
384 {
385     spa_t *spa = dmu_tx_pool(tx)->dp_spa;
386     zfeature_info_t *feature = arg;
387     uint64_t refcount;
388
389     VERIFY0(feature_get_refcount(spa, feature, &refcount));
390     feature_sync(spa, feature, refcount - 1, tx);
391     spa_feature_decr(spa, feature, tx);
392     spa_history_log_internal(spa, "zhack feature decr", tx,
393         "guid=%s", feature->fi_guid);
394 }
395
396 static void
397 zhack_do_feature_ref(int argc, char **argv)
398 {
399     char c;
400     char *target;
401     boolean_t decr = B_FALSE;

```

```

399     spa_t *spa;
400     objset_t *mos;
401     zfeature_info_t feature;
402     spa_feature_t nodeps[] = { SPA_FEATURE_NONE };
403     zfeature_info_t *nodeps[] = { NULL };
404
405     /*
406      * fi_desc does not matter here because it was written to disk
407      * when the feature was enabled, but we need to properly set the
408      * feature for read or write based on the information we read off
409      * disk later.
410      */
411     feature.fi_uname = "zhack";
412     feature.fi_mos = B_FALSE;
413     feature.fi_desc = NULL;
414     feature.fi_depends = nodeps;
415
416     optind = 1;
417     while ((c = getopt(argc, argv, "md")) != -1) {
418         switch (c) {
419             case 'm':
420                 feature.fi_mos = B_TRUE;
421                 break;
422             case 'd':
423                 decr = B_TRUE;
424                 break;
425             default:
426                 usage();
427                 break;
428         }
429     }
430     argc -= optind;
431     argv += optind;
432
433     if (argc < 2) {
434         (void) fprintf(stderr, "error: missing feature or pool name\n");
435         usage();
436     }
437     target = argv[0];
438     feature.fi_guid = argv[1];
439
440     if (!zfeature_is_valid_guid(feature.fi_guid))
441         fatal(NULL, FTAG, "invalid feature guid: %s", feature.fi_guid);
442
443     zhack_spa_open(target, B_FALSE, FTAG, &spa);
444     mos = spa->spa_meta_objset;
445
446     if (zfeature_is_supported(feature.fi_guid)) {
447         fatal(spa, FTAG,
448             "'%s' is a real feature, will not change refcount");
449     }
450     if (0 == zfeature_lookup_guid(feature.fi_guid, NULL))
451         fatal(spa, FTAG, "'%s' is a real feature, will not change "
452             "refcount");
453
454     if (0 == zap_contains(mos, spa->spa_feat_for_read_obj,
455         feature.fi_guid)) {
456         feature.fi_can_readonly = B_FALSE;
457     } else if (0 == zap_contains(mos, spa->spa_feat_for_write_obj,
458         feature.fi_guid)) {
459         feature.fi_can_readonly = B_TRUE;
460     } else {
461         fatal(spa, FTAG, "feature is not enabled: %s", feature.fi_guid);
462     }
463
464     if (decr) {

```

```

461         uint64_t count;
462         if (feature_get_refcount(spa, &feature, &count) == 0 &&
463             count != 0) {
464             if (decr && !spa_feature_is_active(spa, &feature))
465                 fatal(spa, FTAG, "feature refcount already 0: %s",
466                     feature.fi_guid);
467         }
468
469         VERIFY0(dsl_sync_task(spa_name(spa), NULL,
470             decr ? feature_decr_sync : feature_incr_sync, &feature, 5));
471
472         spa_close(spa, FTAG);
473     }
474 }
475
476 _____unchanged_portion_omitted_____

```

```

*****
129948 Tue Oct 1 14:04:08 2013
new/usr/src/cmd/zpool/zpool_main.c
4171 clean up spa_feature_*() interfaces
4172 implement extensible_dataset feature for use by other zpool features
Reviewed by: Max Grossman <max.grossman@delphix.com>
Reviewed by: Christopher Siden <christopher.siden@delphix.com>
Reviewed by: George Wilson <george.wilson@delphix.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */

22 /*
23  * Copyright (c) 2005, 2010, Oracle and/or its affiliates. All rights reserved.
24  * Copyright 2011 Nexenta Systems, Inc. All rights reserved.
25  * Copyright (c) 2013 by Delphix. All rights reserved.
26  * Copyright (c) 2012 by Delphix. All rights reserved.
27  * Copyright (c) 2012 by Frederik Wessels. All rights reserved.
28  * Copyright (c) 2013 by Prasad Joshi (sTec). All rights reserved.
29 */

30 #include <assert.h>
31 #include <ctype.h>
32 #include <dirent.h>
33 #include <errno.h>
34 #include <fcntl.h>
35 #include <libgen.h>
36 #include <libintl.h>
37 #include <libutil.h>
38 #include <locale.h>
39 #include <stdio.h>
40 #include <stdlib.h>
41 #include <string.h>
42 #include <strings.h>
43 #include <unistd.h>
44 #include <priv.h>
45 #include <pwd.h>
46 #include <zone.h>
47 #include <zfs_prop.h>
48 #include <sys/fs/zfs.h>
49 #include <sys/stat.h>

51 #include <libzfs.h>

53 #include "zpool_util.h"
54 #include "zfs_comutil.h"
55 #include "zfeature_common.h"

```

```

57 #include "statcommon.h"

59 static int zpool_do_create(int, char **);
60 static int zpool_do_destroy(int, char **);

62 static int zpool_do_add(int, char **);
63 static int zpool_do_remove(int, char **);

65 static int zpool_do_list(int, char **);
66 static int zpool_do_iostat(int, char **);
67 static int zpool_do_status(int, char **);

69 static int zpool_do_online(int, char **);
70 static int zpool_do_offline(int, char **);
71 static int zpool_do_clear(int, char **);
72 static int zpool_do_reopen(int, char **);

74 static int zpool_do_reguid(int, char **);

76 static int zpool_do_attach(int, char **);
77 static int zpool_do_detach(int, char **);
78 static int zpool_do_replace(int, char **);
79 static int zpool_do_split(int, char **);

81 static int zpool_do_scrub(int, char **);

83 static int zpool_do_import(int, char **);
84 static int zpool_do_export(int, char **);

86 static int zpool_do_upgrade(int, char **);

88 static int zpool_do_history(int, char **);

90 static int zpool_do_get(int, char **);
91 static int zpool_do_set(int, char **);

93 /*
94  * These libumem hooks provide a reasonable set of defaults for the allocator's
95  * debugging facilities.
96  */

98 #ifdef DEBUG
99 const char *
100 _umem_debug_init(void)
101 {
102     return ("default,verbose"); /* $UMEM_DEBUG setting */
103 }

unchanged portion omitted

618 /*
619  * zpool create [-fnd] [-o property=value] ...
620  *               [-O file-system-property=value] ...
621  *               [-R root] [-m mountpoint] <pool> <dev> ...
622  *
623  * -f      Force creation, even if devices appear in use
624  * -n      Do not create the pool, but display the resulting layout if it
625  *          were to be created.
626  * -R      Create a pool under an alternate root
627  * -m      Set default mountpoint for the root dataset. By default it's
628  *          '/<pool>'
629  * -o      Set property=value.
630  * -d      Don't automatically enable all supported pool features
631  *          (individual features can be enabled with -o).
632  * -O      Set fsproperty=value in the pool's root file system
633  *
634  * Creates the named pool according to the given vdev specification. The

```

```

635 * bulk of the vdev processing is done in get_vdev_spec() in zpool_vdev.c. Once
636 * we get the nvlist back from get_vdev_spec(), we either print out the contents
637 * (if '-n' was specified), or pass it to libzfs to do the creation.
638 */
639 int
640 zpool_do_create(int argc, char **argv)
641 {
642     boolean_t force = B_FALSE;
643     boolean_t dryrun = B_FALSE;
644     boolean_t enable_all_pool_feat = B_TRUE;
645     int c;
646     nvlist_t *nvroot = NULL;
647     char *poolname;
648     int ret = 1;
649     char *altroot = NULL;
650     char *mountpoint = NULL;
651     nvlist_t *fsprops = NULL;
652     nvlist_t *props = NULL;
653     char *propval;

655     /* check options */
656     while ((c = getopt(argc, argv, "fndR:m:o:O:")) != -1) {
657         switch (c) {
658             case 'f':
659                 force = B_TRUE;
660                 break;
661             case 'n':
662                 dryrun = B_TRUE;
663                 break;
664             case 'd':
665                 enable_all_pool_feat = B_FALSE;
666                 break;
667             case 'R':
668                 altroot = optarg;
669                 if (add_prop_list(zpool_prop_to_name(
670                     ZPOOL_PROP_ALTROOT), optarg, &props, B_TRUE))
671                     goto errout;
672             if (nvlist_lookup_string(props,
673                 zpool_prop_to_name(ZPOOL_PROP_CACHEFILE),
674                 &propval) == 0)
675                 break;
676             if (add_prop_list(zpool_prop_to_name(
677                 ZPOOL_PROP_CACHEFILE), "none", &props, B_TRUE))
678                 goto errout;
679             break;
680             case 'm':
681                 /* Equivalent to -O mountpoint=optarg */
682                 mountpoint = optarg;
683                 break;
684             case 'o':
685                 if ((propval = strchr(optarg, '=')) == NULL) {
686                     (void) fprintf(stderr, gettext("missing "
687                         "'=' for -o option\n"));
688                     goto errout;
689                 }
690                 *propval = '\0';
691                 propval++;

693                 if (add_prop_list(optarg, propval, &props, B_TRUE))
694                     goto errout;

696             /*
697              * If the user is creating a pool that doesn't support
698              * feature flags, don't enable any features.
699              */
700             if (zpool_name_to_prop(optarg) == ZPOOL_PROP_VERSION) {

```

```

701         char *end;
702         u_longlong_t ver;

704         ver = strtoull(propval, &end, 10);
705         if (*end == '\0' &&
706             ver < SPA_VERSION_FEATURES) {
707             enable_all_pool_feat = B_FALSE;
708         }
709     }
710     break;
711 case 'O':
712     if ((propval = strchr(optarg, '=')) == NULL) {
713         (void) fprintf(stderr, gettext("missing "
714             "'=' for -O option\n"));
715         goto errout;
716     }
717     *propval = '\0';
718     propval++;

720     /*
721      * Mountpoints are checked and then added later.
722      * Uniquely among properties, they can be specified
723      * more than once, to avoid conflict with -m.
724      */
725     if (0 == strcmp(optarg,
726         zfs_prop_to_name(ZFS_PROP_MOUNTPOINT))) {
727         mountpoint = propval;
728     } else if (add_prop_list(optarg, propval, &fsprops,
729         B_FALSE)) {
730         goto errout;
731     }
732     break;
733 case ':':
734     (void) fprintf(stderr, gettext("missing argument for "
735         "'%c' option\n"), optopt);
736     goto badusage;
737 case '?':
738     (void) fprintf(stderr, gettext("invalid option '%c'\n"),
739         optopt);
740     goto badusage;
741 }
742 }

744     argc -= optind;
745     argv += optind;

747     /* get pool name and check number of arguments */
748     if (argc < 1) {
749         (void) fprintf(stderr, gettext("missing pool name argument\n"));
750         goto badusage;
751     }
752     if (argc < 2) {
753         (void) fprintf(stderr, gettext("missing vdev specification\n"));
754         goto badusage;
755     }

757     poolname = argv[0];

759     /*
760      * As a special case, check for use of '/' in the name, and direct the
761      * user to use 'zfs create' instead.
762      */
763     if (strchr(poolname, '/') != NULL) {
764         (void) fprintf(stderr, gettext("cannot create '%s': invalid "
765             "character '/' in pool name\n"), poolname);
766         (void) fprintf(stderr, gettext("use 'zfs create' to "

```

```

767         "create a dataset\n"));
768         goto errout;
769     }

771     /* pass off to get_vdev_spec for bulk processing */
772     nvroot = make_root_vdev(NULL, force, !force, B_FALSE, dryrun,
773         argc - 1, argv + 1);
774     if (nvroot == NULL)
775         goto errout;

777     /* make_root_vdev() allows 0 toplevel children if there are spares */
778     if (!zfs_allocatable_devs(nvroot)) {
779         (void) fprintf(stderr, gettext("invalid vdev "
780             "specification: at least one toplevel vdev must be "
781             "specified\n"));
782         goto errout;
783     }

785     if (altroot != NULL && altroot[0] != '/') {
786         (void) fprintf(stderr, gettext("invalid alternate root '%s': "
787             "must be an absolute path\n"), altroot);
788         goto errout;
789     }

791     /*
792     * Check the validity of the mountpoint and direct the user to use the
793     * '-m' mountpoint option if it looks like its in use.
794     */
795     if (mountpoint == NULL ||
796         (strcmp(mountpoint, ZFS_MOUNTPOINT_LEGACY) != 0 &&
797         strcmp(mountpoint, ZFS_MOUNTPOINT_NONE) != 0)) {
798         char buf[MAXPATHLEN];
799         DIR *dirp;

801         if (mountpoint && mountpoint[0] != '/') {
802             (void) fprintf(stderr, gettext("invalid mountpoint "
803                 "'%s': must be an absolute path, 'legacy', or "
804                 "'none'\n"), mountpoint);
805             goto errout;
806         }

808         if (mountpoint == NULL) {
809             if (altroot != NULL)
810                 (void) snprintf(buf, sizeof (buf), "%s/%s",
811                     altroot, poolname);
812             else
813                 (void) snprintf(buf, sizeof (buf), "%s",
814                     poolname);
815         } else {
816             if (altroot != NULL)
817                 (void) snprintf(buf, sizeof (buf), "%s%s",
818                     altroot, mountpoint);
819             else
820                 (void) snprintf(buf, sizeof (buf), "%s",
821                     mountpoint);
822         }

824         if ((dirp = opendir(buf)) == NULL && errno != ENOENT) {
825             (void) fprintf(stderr, gettext("mountpoint '%s' : "
826                 "%s\n"), buf, strerror(errno));
827             (void) fprintf(stderr, gettext("use '-m' "
828                 "option to provide a different default\n"));
829             goto errout;
830         } else if (dirp) {
831             int count = 0;

```

```

833         while (count < 3 && readdir(dirp) != NULL)
834             count++;
835         (void) closedir(dirp);

837         if (count > 2) {
838             (void) fprintf(stderr, gettext("mountpoint "
839                 "'%s' exists and is not empty\n"), buf);
840             (void) fprintf(stderr, gettext("use '-m' "
841                 "option to provide a "
842                 "different default\n"));
843             goto errout;
844         }
845     }
846 }

848 /*
849 * Now that the mountpoint's validity has been checked, ensure that
850 * the property is set appropriately prior to creating the pool.
851 */
852 if (mountpoint != NULL) {
853     ret = add_prop_list(zfs_prop_to_name(ZFS_PROP_MOUNTPOINT),
854         mountpoint, &fsprops, B_FALSE);
855     if (ret != 0)
856         goto errout;
857 }

859 ret = 1;
860 if (dryrun) {
861     /*
862     * For a dry run invocation, print out a basic message and run
863     * through all the vdevs in the list and print out in an
864     * appropriate hierarchy.
865     */
866     (void) printf(gettext("would create '%s' with the "
867         "following layout:\n\n"), poolname);

869     print_vdev_tree(NULL, poolname, nvroot, 0, B_FALSE);
870     if (num_logs(nvroot) > 0)
871         print_vdev_tree(NULL, "logs", nvroot, 0, B_TRUE);

873     ret = 0;
874 } else {
875     /*
876     * Hand off to libzfs.
877     */
878     if (enable_all_pool_feat) {
879         spa_feature_t i;
880         int i;
881         for (i = 0; i < SPA_FEATURES; i++) {
882             char propname[MAXPATHLEN];
883             zfeature_info_t *feat = &spa_feature_table[i];

884             (void) snprintf(propname, sizeof (propname),
885                 "feature@%s", feat->fi_uname);

887             /*
888             * Skip feature if user specified it manually
889             * on the command line.
890             */
891             if (nvlist_exists(props, propname))
892                 continue;

894             ret = add_prop_list(propname,
895                 ZFS_FEATURE_ENABLED, &props, B_TRUE);
896             if (ret != 0)
897                 goto errout;

```



```
898         }
899     }
900
901     ret = 1;
902     if (zpool_create(g_zfs, poolname,
903         nvroot, props, fsprops) == 0) {
904         zfs_handle_t *pool = zfs_open(g_zfs, poolname,
905             ZFS_TYPE_FILESYSTEM);
906         if (pool != NULL) {
907             if (zfs_mount(pool, NULL, 0) == 0)
908                 ret = zfs_shareall(pool);
909             zfs_close(pool);
910         }
911     } else if (libzfs_errno(g_zfs) == EZFS_INVALIDNAME) {
912         (void) fprintf(stderr, gettext("pool name may have "
913             "been omitted\n"));
914     }
915 }
916
917 errout:
918     nvlist_free(nvroot);
919     nvlist_free(fsprops);
920     nvlist_free(props);
921     return (ret);
922 badusage:
923     nvlist_free(fsprops);
924     nvlist_free(props);
925     usage(B_FALSE);
926     return (2);
927 }
unchanged_portion_omitted
```

```

*****
4739 Tue Oct 1 14:04:10 2013
new/usr/src/common/zfs/zfeature_common.c
4171 clean up spa_feature_*() interfaces
4172 implement extensible_dataset feature for use by other zpool features
Reviewed by: Max Grossman <max.grossman@delphix.com>
Reviewed by: Christopher Siden <christopher.siden@delphix.com>
Reviewed by: George Wilson <george.wilson@delphix.com>
*****
unchanged_portion_omitted

89 boolean_t
90 zfeature_is_supported(const char *guid)
91 {
92     if (zfeature_checks_disable)
93         return (B_TRUE);

95     for (spa_feature_t i = 0; i < SPA_FEATURES; i++) {
96         return (0 == zfeature_lookup_guid(guid, NULL));
97     }

98     int
99     zfeature_lookup_guid(const char *guid, zfeature_info_t **res)
100 {
101     for (int i = 0; i < SPA_FEATURES; i++) {
102         zfeature_info_t *feature = &spa_feature_table[i];
103         if (strcmp(guid, feature->fi_guid) == 0)
104             return (B_TRUE);
105         if (strcmp(guid, feature->fi_guid) == 0) {
106             if (res != NULL)
107                 *res = feature;
108             return (0);
109         }
110     }
111     return (B_FALSE);
112 }

113

114 return (ENOENT);
115 }

116

117 int
118 zfeature_lookup_name(const char *name, spa_feature_t *res)
119 zfeature_lookup_name(const char *name, zfeature_info_t **res)
120 {
121     for (spa_feature_t i = 0; i < SPA_FEATURES; i++) {
122         for (int j = 0; j < SPA_FEATURES; j++) {
123             zfeature_info_t *feature = &spa_feature_table[j];
124             if (strcmp(name, feature->fi_uname) == 0) {
125                 if (res != NULL)
126                     *res = i;
127                 *res = feature;
128                 return (0);
129             }
130         }
131     }
132     return (ENOENT);
133 }

134

135 static void
136 zfeature_register(spa_feature_t fid, const char *guid, const char *name,
137                  const char *desc, boolean_t readonly, boolean_t mos,
138                  const spa_feature_t *deps)
139 zfeature_register(int fid, const char *guid, const char *name, const char *desc,
140                  boolean_t readonly, boolean_t mos, zfeature_info_t **deps)
141 {
142     zfeature_info_t *feature = &spa_feature_table[fid];
143     static spa_feature_t nodeps[] = { SPA_FEATURE_NONE };

```

```

133     static zfeature_info_t *nodeps[] = { NULL };

134

135     ASSERT(name != NULL);
136     ASSERT(desc != NULL);
137     ASSERT(!readonly || !mos);
138     ASSERT3U(fid, <, SPA_FEATURES);
139     ASSERT(zfeature_is_valid_guid(guid));

140     if (deps == NULL)
141         deps = nodeps;

142     feature->fi_feature = fid;
143     feature->fi_guid = guid;
144     feature->fi_uname = name;
145     feature->fi_desc = desc;
146     feature->fi_can_readonly = readonly;
147     feature->fi_mos = mos;
148     feature->fi_depends = deps;
149 }

150 void
151 zpool_feature_init(void)
152 {
153     zfeature_register(SPA_FEATURE_ASYNC_DESTROY,
154                      "com.delphix:async_destroy", "async_destroy",
155                      "Destroy filesystems asynchronously.", B_TRUE, B_FALSE, NULL);
156     zfeature_register(SPA_FEATURE_EMPTY_BPOBJ,
157                      "com.delphix:empty_bpobj", "empty_bpobj",
158                      "Snapshots use less space.", B_TRUE, B_FALSE, NULL);
159     zfeature_register(SPA_FEATURE_LZ4_COMPRESS,
160                      "org.illumos:lz4_compress", "lz4_compress",
161                      "LZ4 compression algorithm support.", B_FALSE, B_FALSE, NULL);
162     zfeature_register(SPA_FEATURE_MULTI_VDEV_CRASH_DUMP,
163                      "com.joyent:multi_vdev_crash_dump", "multi_vdev_crash_dump",
164                      "Crash dumps to multiple vdev pools.", B_FALSE, B_FALSE, NULL);
165     zfeature_register(SPA_FEATURE_SPACEMAP_HISTOGRAM,
166                      "com.delphix:spacemap_histogram", "spacemap_histogram",
167                      "Spacemaps maintain space histograms.", B_TRUE, B_FALSE, NULL);
168     zfeature_register(SPA_FEATURE_EXTENSIBLE_DATASET,
169                      "com.delphix:extensible_dataset", "extensible_dataset",
170                      "Enhanced dataset functionality, used by other features.",
171                      B_FALSE, B_FALSE, NULL);
172 }

173 unchanged_portion_omitted

```

new/usr/src/common/zfs/zfeature\_common.h

1

```
*****
2357 Tue Oct 1 14:04:14 2013
new/usr/src/common/zfs/zfeature_common.h
4171 clean up spa_feature_*() interfaces
4172 implement extensible_dataset feature for use by other zpool features
Reviewed by: Max Grossman <max.grossman@delphix.com>
Reviewed by: Christopher Siden <christopher.siden@delphix.com>
Reviewed by: George Wilson <george.wilson@delphix.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */

22 /*
23  * Copyright (c) 2013 by Delphix. All rights reserved.
24  * Copyright (c) 2013 by Saso Kiselkov. All rights reserved.
25  * Copyright (c) 2013, Joyent, Inc. All rights reserved.
26  */

28 #ifndef _ZFEATURE_COMMON_H
29 #define _ZFEATURE_COMMON_H

31 #include <sys/fs/zfs.h>
32 #include <sys/inttypes.h>
33 #include <sys/types.h>

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

39 struct zfeature_info;

41 typedef enum spa_feature {
42     SPA_FEATURE_NONE = -1,
43     SPA_FEATURE_ASYNC_DESTROY,
44     SPA_FEATURE_EMPTY_BPOBJ,
45     SPA_FEATURE_LZ4_COMPRESS,
46     SPA_FEATURE_MULTI_VDEV_CRASH_DUMP,
47     SPA_FEATURE_SPACEMAP_HISTOGRAM,
48     SPA_FEATURE_EXTENSIBLE_DATASET,
49     SPA_FEATURES
50 } spa_feature_t;

52 typedef struct zfeature_info {
53     spa_feature_t fi_feature;
54     const char *fi_uname; /* User-facing feature name */
55     const char *fi_guid; /* On-disk feature identifier */
56     const char *fi_desc; /* Feature description */
57     boolean_t fi_can_readonly; /* Can open pool readonly w/o support? */

```

new/usr/src/common/zfs/zfeature\_common.h

2

```
58     boolean_t fi_mos; /* Is the feature necessary to read the MOS? */
59     /* array of dependencies, terminated by SPA_FEATURE_NONE */
60     const spa_feature_t *fi_depends;
61     struct zfeature_info **fi_depends; /* array; null terminated */
62 } zfeature_info_t;

63 typedef int (zfeature_func_t)(zfeature_info_t *fi, void *arg);

65 #define ZFS_FEATURE_DEBUG

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

63 extern zfeature_info_t spa_feature_table[SPA_FEATURES];

65 extern boolean_t zfeature_is_valid_guid(const char *);

71 extern boolean_t zfeature_is_supported(const char *);
72 extern int zfeature_lookup_name(const char *name, spa_feature_t *res);
68 extern int zfeature_lookup_guid(const char *, zfeature_info_t **res);
69 extern int zfeature_lookup_name(const char *, zfeature_info_t **res);

74 extern void zpool_feature_init(void);

76 #ifdef __cplusplus
77 }
_____unchanged_portion_omitted_____

```

```

*****
43611 Tue Oct 1 14:04:17 2013
new/usr/src/grub/grub-0.97/stage2/fsys_zfs.c
4171 clean up spa_feature_*() interfaces
4172 implement extensible_dataset feature for use by other zpool features
Reviewed by: Max Grossman <max.grossman@delphix.com>
Reviewed by: Christopher Siden <christopher.siden@delphix.com>
Reviewed by: George Wilson <george.wilson@delphix.com>
*****
1 /*
2  * GRUB -- GRand Unified Bootloader
3  * Copyright (C) 1999,2000,2001,2002,2003,2004 Free Software Foundation, Inc.
4  *
5  * This program is free software; you can redistribute it and/or modify
6  * it under the terms of the GNU General Public License as published by
7  * the Free Software Foundation; either version 2 of the License, or
8  * (at your option) any later version.
9  *
10 * This program is distributed in the hope that it will be useful,
11 * but WITHOUT ANY WARRANTY; without even the implied warranty of
12 * MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the
13 * GNU General Public License for more details.
14 *
15 * You should have received a copy of the GNU General Public License
16 * along with this program; if not, write to the Free Software
17 * Foundation, Inc., 675 Mass Ave, Cambridge, MA 02139, USA.
18 */
20 /*
21 * Copyright 2010 Sun Microsystems, Inc. All rights reserved.
22 * Use is subject to license terms.
23 */
25 /*
26 * Copyright (c) 2013 by Delphix. All rights reserved.
27 * Copyright (c) 2012 by Delphix. All rights reserved.
28 * Copyright (c) 2013 by Saso Kiselkov. All rights reserved.
29 */
30 /*
31 * The zfs plug-in routines for GRUB are:
32 *
33 * zfs_mount() - locates a valid uberblock of the root pool and reads
34 *               in its MOS at the memory address MOS.
35 *
36 * zfs_open() - locates a plain file object by following the MOS
37 *               and places its dnode at the memory address DNODE.
38 *
39 * zfs_read() - read in the data blocks pointed by the DNODE.
40 *
41 * ZFS_SCRATCH is used as a working area.
42 *
43 * (memory addr)  MOS          DNODE          ZFS_SCRATCH
44 *
45 * memory  +-----V-----V-----V-----+
46 *          |         |         |         |
47 *          |         |         |         |
48 *          |         |         |         |
49 *          +-----+-----+-----+
51 #ifndef  FSY_ZFS
53 #include "shared.h"
54 #include "fileys.h"
55 #include "fsys_zfs.h"

```

```

57 /* cache for a file block of the currently zfs_open()-ed file */
58 static void *file_buf = NULL;
59 static uint64_t file_start = 0;
60 static uint64_t file_end = 0;
62 /* cache for a dnode block */
63 static dnode_phys_t *dnode_buf = NULL;
64 static dnode_phys_t *dnode_mdn = NULL;
65 static uint64_t dnode_start = 0;
66 static uint64_t dnode_end = 0;
68 static uint64_t pool_guid = 0;
69 static uberblock_t current_uberblock;
70 static char *stackbase;
72 decomp_entry_t decomp_table[ZIO_COMPRESS_FUNCTIONS] =
73 {
74     {"inherit", 0}, /* ZIO_COMPRESS_INHERIT */
75     {"on", lzjb_decompress}, /* ZIO_COMPRESS_ON */
76     {"off", 0}, /* ZIO_COMPRESS_OFF */
77     {"lzjb", lzjb_decompress}, /* ZIO_COMPRESS_LZJB */
78     {"empty", 0}, /* ZIO_COMPRESS_EMPTY */
79     {"gzip-1", 0}, /* ZIO_COMPRESS_GZIP_1 */
80     {"gzip-2", 0}, /* ZIO_COMPRESS_GZIP_2 */
81     {"gzip-3", 0}, /* ZIO_COMPRESS_GZIP_3 */
82     {"gzip-4", 0}, /* ZIO_COMPRESS_GZIP_4 */
83     {"gzip-5", 0}, /* ZIO_COMPRESS_GZIP_5 */
84     {"gzip-6", 0}, /* ZIO_COMPRESS_GZIP_6 */
85     {"gzip-7", 0}, /* ZIO_COMPRESS_GZIP_7 */
86     {"gzip-8", 0}, /* ZIO_COMPRESS_GZIP_8 */
87     {"gzip-9", 0}, /* ZIO_COMPRESS_GZIP_9 */
88     {"zle", 0}, /* ZIO_COMPRESS_ZLE */
89     {"lz4", lz4_decompress} /* ZIO_COMPRESS_LZ4 */
90 };
91 /*
92  * unchanged portion omitted
93  */
959 /*
960 * List of pool features that the grub implementation of ZFS supports for
961 * read. Note that features that are only required for write do not need
962 * to be listed here since grub opens pools in read-only mode.
963 *
964 * When this list is updated the version number in usr/src/grub/capability
965 * must be incremented to ensure the new grub gets installed.
966 */
967 static const char *spa_feature_names[] = {
968     "org.illumos:lz4_compress",
969     "com.delphix:extensible_dataset",
970     NULL
971 };
972 /*
973  * unchanged portion omitted
974  */
1011 /*
1012 * Given a MOS metadnode, get the metadnode of a given filesystem name (fsname),
1013 * e.g. pool/rootfs, or a given object number (obj), e.g. the object number
1014 * of pool/rootfs.
1015 *
1016 * If no fsname and no obj are given, return the DSL_DIR metadnode.
1017 * If fsname is given, return its metadnode and its matching object number.
1018 * If only obj is given, return the metadnode for this object number.
1019 *
1020 * Return:
1021 * 0 - success
1022 * errnum - failure
1023 */
1024 static int
1025 get_objset_mdn(dnode_phys_t *mosmdn, char *fsname, uint64_t *obj,

```

```

1026     dnode_phys_t *mdn, char *stack)
1027 {
1028     uint64_t objnum, headobj;
1029     char *cname, ch;
1030     blkptr_t *bp;
1031     objset_phys_t *osp;
1032     int issnapshot = 0;
1033     char *snapname;
1034
1035     if (fsname == NULL && obj) {
1036         headobj = *obj;
1037         goto skip;
1038     }
1039
1040     if (errno = dnode_get(mosmdn, DMU_POOL_DIRECTORY_OBJECT,
1041         DMU_OT_OBJECT_DIRECTORY, mdn, stack))
1042         return (errno);
1043
1044     if (errno = zap_lookup(mdn, DMU_POOL_ROOT_DATASET, &objnum,
1045         stack))
1046         return (errno);
1047
1048     if (errno = dnode_get(mosmdn, objnum, 0, mdn, stack))
1049     if (errno = dnode_get(mosmdn, objnum, DMU_OT_DSL_DIR, mdn, stack))
1050         return (errno);
1051
1052     if (fsname == NULL) {
1053         headobj =
1054             ((dsl_dir_phys_t *)DN_BONUS(mdn))->dd_head_dataset_obj;
1055         goto skip;
1056     }
1057
1058     /* take out the pool name */
1059     while (*fsname && !grub_isspace(*fsname) && *fsname != '/')
1060         fsname++;
1061
1062     while (*fsname && !grub_isspace(*fsname)) {
1063         uint64_t childobj;
1064
1065         while (*fsname == '/')
1066             fsname++;
1067
1068         cname = fsname;
1069         while (*fsname && !grub_isspace(*fsname) && *fsname != '/')
1070             fsname++;
1071         ch = *fsname;
1072         *fsname = 0;
1073
1074         snapname = cname;
1075         while (*snapname && !grub_isspace(*snapname) && *snapname !=
1076             '@')
1077             snapname++;
1078         if (*snapname == '@') {
1079             issnapshot = 1;
1080             *snapname = 0;
1081         }
1082         childobj =
1083             ((dsl_dir_phys_t *)DN_BONUS(mdn))->dd_child_dir_zapobj;
1084         if (errno = dnode_get(mosmdn, childobj,
1085             DMU_OT_DSL_DIR_CHILD_MAP, mdn, stack))
1086             return (errno);
1087
1088         if (zap_lookup(mdn, cname, &objnum, stack))
1089             return (ERR_FILESYSTEM_NOT_FOUND);
1090
1091         if (errno = dnode_get(mosmdn, objnum, 0,

```

```

1089         if (errno = dnode_get(mosmdn, objnum, DMU_OT_DSL_DIR,
1090             mdn, stack))
1091             return (errno);
1092
1093         *fsname = ch;
1094         if (issnapshot)
1095             *snapname = '@';
1096     }
1097     headobj = ((dsl_dir_phys_t *)DN_BONUS(mdn))->dd_head_dataset_obj;
1098     if (obj)
1099         *obj = headobj;
1100
1101 skip:
1102     if (errno = dnode_get(mosmdn, headobj, 0, mdn, stack))
1103     if (errno = dnode_get(mosmdn, headobj, DMU_OT_DSL_DATASET, mdn, stack))
1104         return (errno);
1105     if (issnapshot) {
1106         uint64_t snapobj;
1107
1108         snapobj = ((dsl_dataset_phys_t *)DN_BONUS(mdn))->
1109             ds_snapnames_zapobj;
1110
1111         if (errno = dnode_get(mosmdn, snapobj,
1112             DMU_OT_DSL_DS_SNAP_MAP, mdn, stack))
1113             return (errno);
1114         if (zap_lookup(mdn, snapname + 1, &headobj, stack))
1115             return (ERR_FILESYSTEM_NOT_FOUND);
1116         if (errno = dnode_get(mosmdn, headobj, 0, mdn, stack))
1117         if (errno = dnode_get(mosmdn, headobj,
1118             DMU_OT_DSL_DATASET, mdn, stack))
1119             return (errno);
1120         if (obj)
1121             *obj = headobj;
1122
1123         bp = &((dsl_dataset_phys_t *)DN_BONUS(mdn))->ds_bp;
1124         osp = (objset_phys_t *)stack;
1125         stack += sizeof (objset_phys_t);
1126         if (errno = zio_read(bp, osp, stack))
1127             return (errno);
1128
1129         grub_memmove((char *)mdn, (char *)&osp->os_meta_dnode, DNODE_SIZE);
1130
1131     }
1132     return (0);
1133 }
1134 }
1135 }
1136 }
1137 }
1138 }
1139 }
1140 }
1141 }
1142 }
1143 }
1144 }
1145 }
1146 }
1147 }
1148 }
1149 }
1150 }
1151 }
1152 }
1153 }
1154 }
1155 }
1156 }
1157 }
1158 }
1159 }
1160 }
1161 }
1162 }
1163 }
1164 }
1165 }
1166 }
1167 }
1168 }
1169 }
1170 }
1171 }
1172 }
1173 }
1174 }
1175 }
1176 }
1177 }
1178 }
1179 }
1180 }
1181 }
1182 }
1183 }
1184 }
1185 }
1186 }
1187 }
1188 }
1189 }
1190 }
1191 }
1192 }
1193 }
1194 }
1195 }
1196 }
1197 }
1198 }
1199 }
1200 }
1201 }
1202 }
1203 }
1204 }
1205 }
1206 }
1207 }
1208 }
1209 }
1210 }
1211 }
1212 }
1213 }
1214 }
1215 }
1216 }
1217 }
1218 }
1219 }
1220 }
1221 }
1222 }
1223 }
1224 }
1225 }
1226 }
1227 }
1228 }
1229 }
1230 }
1231 }
1232 }
1233 }
1234 }
1235 }
1236 }
1237 }
1238 }
1239 }
1240 }
1241 }
1242 }
1243 }
1244 }
1245 }
1246 }
1247 }
1248 }
1249 }
1250 }
1251 }
1252 }
1253 }
1254 }
1255 }
1256 }
1257 }
1258 }
1259 }
1260 }
1261 }
1262 }
1263 }
1264 }
1265 }
1266 }
1267 }
1268 }
1269 }
1270 }
1271 }
1272 }
1273 }
1274 }
1275 }
1276 }
1277 }
1278 }
1279 }
1280 }
1281 }
1282 }
1283 }
1284 }
1285 }
1286 }
1287 }
1288 }
1289 }
1290 }
1291 }
1292 }
1293 }
1294 }
1295 }
1296 }
1297 }
1298 }
1299 }
1300 }
1301 }
1302 }
1303 }
1304 }
1305 }
1306 }
1307 }
1308 }
1309 }
1310 }
1311 }
1312 }
1313 }
1314 }
1315 }
1316 }
1317 }
1318 }
1319 }
1320 }
1321 }
1322 }
1323 }
1324 }
1325 }
1326 }
1327 }
1328 }
1329 }
1330 }
1331 }
1332 }
1333 }
1334 }
1335 }
1336 }
1337 }
1338 }
1339 }
1340 }
1341 }
1342 }
1343 }
1344 }
1345 }
1346 }
1347 }
1348 }
1349 }
1350 }
1351 }
1352 }
1353 }
1354 }
1355 }
1356 }
1357 }
1358 }
1359 }
1360 }
1361 }
1362 }
1363 }
1364 }
1365 }
1366 }
1367 }
1368 }
1369 }
1370 }
1371 }
1372 }
1373 }
1374 }
1375 }
1376 }
1377 }
1378 }
1379 }
1380 }
1381 }
1382 }
1383 }
1384 }
1385 }
1386 }
1387 }
1388 }
1389 }
1390 }
1391 }
1392 }
1393 }
1394 }
1395 }
1396 }
1397 }
1398 }
1399 }
1400 }
1401 }
1402 }
1403 }
1404 }
1405 }
1406 }
1407 }
1408 }
1409 }
1410 }
1411 }
1412 }
1413 }
1414 }
1415 }
1416 }
1417 }
1418 }
1419 }
1420 }
1421 }
1422 }
1423 }
1424 }
1425 }
1426 }
1427 }
1428 }
1429 }
1430 }
1431 }
1432 }
1433 }
1434 }
1435 }
1436 }
1437 }
1438 }
1439 }
1440 }
1441 }
1442 }
1443 }
1444 }
1445 }
1446 }
1447 }
1448 }
1449 }
1450 }
1451 }
1452 }
1453 }
1454 }
1455 }
1456 }
1457 }
1458 }
1459 }
1460 }
1461 }
1462 }
1463 }
1464 }
1465 }
1466 }
1467 }
1468 }
1469 }
1470 }
1471 }
1472 }
1473 }
1474 }
1475 }
1476 }
1477 }
1478 }
1479 }
1480 }
1481 }
1482 }
1483 }
1484 }
1485 }
1486 }
1487 }
1488 }
1489 }
1490 }
1491 }
1492 }
1493 }
1494 }
1495 }
1496 }
1497 }
1498 }
1499 }
1500 }

```

```

*****
102954 Tue Oct 1 14:04:19 2013
new/usr/src/lib/libzfs/common/libzfs_pool.c
4171 clean up spa_feature_*() interfaces
4172 implement extensible_dataset feature for use by other zpool features
Reviewed by: Max Grossman <max.grossman@delphix.com>
Reviewed by: Christopher Siden <christopher.siden@delphix.com>
Reviewed by: George Wilson <george.wilson@delphix.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */

22 /*
23  * Copyright (c) 2005, 2010, Oracle and/or its affiliates. All rights reserved.
24  * Copyright 2011 Nexenta Systems, Inc. All rights reserved.
25  * Copyright (c) 2013 by Delphix. All rights reserved.
26  * Copyright (c) 2012 by Delphix. All rights reserved.
27  * Copyright (c) 2013, Joyent, Inc. All rights reserved.
28 */

29 #include <ctype.h>
30 #include <errno.h>
31 #include <devid.h>
32 #include <fcntl.h>
33 #include <libintl.h>
34 #include <stdio.h>
35 #include <stdlib.h>
36 #include <strings.h>
37 #include <unistd.h>
38 #include <libgen.h>
39 #include <sys/efi_partition.h>
40 #include <sys/vtoc.h>
41 #include <sys/zfs_ioctl.h>
42 #include <dldfcn.h>

44 #include "zfs_namecheck.h"
45 #include "zfs_prop.h"
46 #include "libzfs_impl.h"
47 #include "zfs_comutil.h"
48 #include "zfeature_common.h"

50 static int read_efi_label(nvlist_t *config, diskaddr_t *sb);

52 #define DISK_ROOT      "/dev/dsk"
53 #define RDISK_ROOT    "/dev/rdsk"
54 #define BACKUP_SLICE  "s2"

56 typedef struct prop_flags {

```

```

57     int create:1; /* Validate property on creation */
58     int import:1; /* Validate property on import */
59 } prop_flags_t;
_____unchanged_portion_omitted

385 /*
386  * Given an nvlist of zpool properties to be set, validate that they are
387  * correct, and parse any numeric properties (index, boolean, etc) if they are
388  * specified as strings.
389  */
390 static nvlist_t *
391 zpool_valid_proplist(libzfs_handle_t *hdl, const char *poolname,
392     nvlist_t *props, uint64_t version, prop_flags_t flags, char *errbuf)
393 {
394     nvpair_t *elem;
395     nvlist_t *retprops;
396     zpool_prop_t prop;
397     char *strval;
398     uint64_t intval;
399     char *slash, *check;
400     struct stat64 statbuf;
401     zpool_handle_t *zhp;
402     nvlist_t *nvroot;

404     if (nvlist_alloc(&retprops, NV_UNIQUE_NAME, 0) != 0) {
405         (void) no_memory(hdl);
406         return (NULL);
407     }

409     elem = NULL;
410     while ((elem = nvlist_next_nvpair(props, elem)) != NULL) {
411         const char *propname = nvpair_name(elem);

413         prop = zpool_name_to_prop(propname);
414         if (prop == ZPROP_INVALID && zpool_prop_feature(propname)) {
415             int err;
416             zfeature_info_t *feature;
417             char *fname = strchr(propname, '@') + 1;

418             err = zfeature_lookup_name(fname, NULL);
419             err = zfeature_lookup_name(fname, &feature);
420             if (err != 0) {
421                 ASSERT3U(err, ==, ENOENT);
422                 zfs_error_aux(hdl, dgettext(TEXT_DOMAIN,
423                     "invalid feature '%s'", fname));
424                 (void) zfs_error(hdl, EZFS_BADPROP, errbuf);
425                 goto error;
426             }

427             if (nvpair_type(elem) != DATA_TYPE_STRING) {
428                 zfs_error_aux(hdl, dgettext(TEXT_DOMAIN,
429                     "'%s' must be a string"), propname);
430                 (void) zfs_error(hdl, EZFS_BADPROP, errbuf);
431                 goto error;
432             }

434             (void) nvpair_value_string(elem, &strval);
435             if (strcmp(strval, ZFS_FEATURE_ENABLED) != 0) {
436                 zfs_error_aux(hdl, dgettext(TEXT_DOMAIN,
437                     "property '%s' can only be set to "
438                     "'enabled'"), propname);
439                 (void) zfs_error(hdl, EZFS_BADPROP, errbuf);
440                 goto error;
441             }

```

```

443         if (nvlist_add_uint64(retprops, propname, 0) != 0) {
444             (void) no_memory(hdl);
445             goto error;
446         }
447         continue;
448     }
449
450     /*
451     * Make sure this property is valid and applies to this type.
452     */
453     if (prop == ZPROP_INVALID) {
454         zfs_error_aux(hdl, dgettext(TEXT_DOMAIN,
455             "invalid property '%s'", propname);
456         (void) zfs_error(hdl, EZFS_BADPROP, errbuf);
457         goto error;
458     }
459
460     if (zpool_prop_readonly(prop)) {
461         zfs_error_aux(hdl, dgettext(TEXT_DOMAIN, "'%s' "
462             "is readonly"), propname);
463         (void) zfs_error(hdl, EZFS_PROPREADONLY, errbuf);
464         goto error;
465     }
466
467     if (zprop_parse_value(hdl, elem, prop, ZFS_TYPE_POOL, retprops,
468         &strval, &intval, errbuf) != 0)
469         goto error;
470
471     /*
472     * Perform additional checking for specific properties.
473     */
474     switch (prop) {
475     case ZPOOL_PROP_VERSION:
476         if (intval < version ||
477             !SPA_VERSION_IS_SUPPORTED(intval)) {
478             zfs_error_aux(hdl, dgettext(TEXT_DOMAIN,
479                 "property '%s' number %d is invalid.",
480                 propname, intval);
481             (void) zfs_error(hdl, EZFS_BADVERSION, errbuf);
482             goto error;
483         }
484         break;
485
486     case ZPOOL_PROP_BOOTFS:
487         if (flags.create || flags.import) {
488             zfs_error_aux(hdl, dgettext(TEXT_DOMAIN,
489                 "property '%s' cannot be set at creation "
490                 "or import time"), propname);
491             (void) zfs_error(hdl, EZFS_BADPROP, errbuf);
492             goto error;
493         }
494
495         if (version < SPA_VERSION_BOOTFS) {
496             zfs_error_aux(hdl, dgettext(TEXT_DOMAIN,
497                 "pool must be upgraded to support "
498                 "'%s' property"), propname);
499             (void) zfs_error(hdl, EZFS_BADVERSION, errbuf);
500             goto error;
501         }
502
503         /*
504         * bootfs property value has to be a dataset name and
505         * the dataset has to be in the same pool as it sets to.
506         */
507         if (strval[0] != '\0' && !bootfs_name_valid(poolname,
508             strval)) {

```

```

509             zfs_error_aux(hdl, dgettext(TEXT_DOMAIN, "'%s' "
510                 "is an invalid name"), strval);
511             (void) zfs_error(hdl, EZFS_INVALIDNAME, errbuf);
512             goto error;
513         }
514
515         if ((zhp = zpool_open_canfail(hdl, poolname)) == NULL) {
516             zfs_error_aux(hdl, dgettext(TEXT_DOMAIN,
517                 "could not open pool '%s'", poolname);
518             (void) zfs_error(hdl, EZFS_OPENFAILED, errbuf);
519             goto error;
520         }
521         verify(nvlist_lookup_nvlist(zpool_get_config(zhp, NULL),
522             ZPOOL_CONFIG_VDEV_TREE, &nvroot) == 0);
523
524         /*
525         * bootfs property cannot be set on a disk which has
526         * been EFI labeled.
527         */
528         if (pool_uses_efi(nvroot)) {
529             zfs_error_aux(hdl, dgettext(TEXT_DOMAIN,
530                 "property '%s' not supported on "
531                 "EFI labeled devices"), propname);
532             (void) zfs_error(hdl, EZFS_POOL_NOTSUP, errbuf);
533             zpool_close(zhp);
534             goto error;
535         }
536         zpool_close(zhp);
537         break;
538
539     case ZPOOL_PROP_ALTROOT:
540         if (!flags.create && !flags.import) {
541             zfs_error_aux(hdl, dgettext(TEXT_DOMAIN,
542                 "property '%s' can only be set during pool "
543                 "creation or import"), propname);
544             (void) zfs_error(hdl, EZFS_BADPROP, errbuf);
545             goto error;
546         }
547
548         if (strval[0] != '/') {
549             zfs_error_aux(hdl, dgettext(TEXT_DOMAIN,
550                 "bad alternate root '%s'", strval);
551             (void) zfs_error(hdl, EZFS_BADPATH, errbuf);
552             goto error;
553         }
554         break;
555
556     case ZPOOL_PROP_CACHEFILE:
557         if (strval[0] == '\0')
558             break;
559
560         if (strcmp(strval, "none") == 0)
561             break;
562
563         if (strval[0] != '/') {
564             zfs_error_aux(hdl, dgettext(TEXT_DOMAIN,
565                 "property '%s' must be empty, an "
566                 "absolute path, or 'none'"), propname);
567             (void) zfs_error(hdl, EZFS_BADPATH, errbuf);
568             goto error;
569         }
570
571         slash = strrchr(strval, '/');
572
573         if (slash[1] == '\0' || strcmp(slash, "/.") == 0 ||
574             strcmp(slash, "/..") == 0) {

```

```

575         zfs_error_aux(hdl, dgettext(TEXT_DOMAIN,
576             "%s' is not a valid file"), strval);
577         (void) zfs_error(hdl, EZFS_BADPATH, errbuf);
578         goto error;
579     }
581     *slash = '\\0';
583     if (strval[0] != '\\0' &&
584         (stat64(strval, &statbuf) != 0 ||
585         !S_ISDIR(statbuf.st_mode))) {
586         zfs_error_aux(hdl, dgettext(TEXT_DOMAIN,
587             "%s' is not a valid directory"),
588             strval);
589         (void) zfs_error(hdl, EZFS_BADPATH, errbuf);
590         goto error;
591     }
593     *slash = '/';
594     break;
596 case ZPOOL_PROP_COMMENT:
597     for (check = strval; *check != '\\0'; check++) {
598         if (!isprint(*check)) {
599             zfs_error_aux(hdl,
600                 dgettext(TEXT_DOMAIN,
601                     "comment may only have printable "
602                     "characters"));
603             (void) zfs_error(hdl, EZFS_BADPROP,
604                 errbuf);
605             goto error;
606         }
607     }
608     if (strlen(strval) > ZPROP_MAX_COMMENT) {
609         zfs_error_aux(hdl, dgettext(TEXT_DOMAIN,
610             "comment must not exceed %d characters"),
611             ZPROP_MAX_COMMENT);
612         (void) zfs_error(hdl, EZFS_BADPROP, errbuf);
613         goto error;
614     }
615     break;
616 case ZPOOL_PROP_READONLY:
617     if (!flags.import) {
618         zfs_error_aux(hdl, dgettext(TEXT_DOMAIN,
619             "property '%s' can only be set at "
620             "import time"), propname);
621         (void) zfs_error(hdl, EZFS_BADPROP, errbuf);
622         goto error;
623     }
624     break;
625 }
626 }
628     return (retprops);
629 error:
630     nvlist_free(retprops);
631     return (NULL);
632 }

```

unchanged portion omitted

```

786 /*
787  * Get the state for the given feature on the given ZFS pool.
788  */
789 int
790 zpool_prop_get_feature(zpool_handle_t *zhp, const char *propname, char *buf,
791     size_t len)

```

```

792 {
793     uint64_t refcount;
794     boolean_t found = B_FALSE;
795     nvlist_t *features = zpool_get_features(zhp);
796     boolean_t supported;
797     const char *feature = strchr(propname, '@') + 1;
799     supported = zpool_prop_feature(propname);
800     ASSERT(supported || zfs_prop_unsupported(propname));
802     /*
803      * Convert from feature name to feature guid. This conversion is
804      * unnecessary for unsupported@... properties because they already
805      * use guids.
806      */
807     if (supported) {
808         int ret;
809         spa_feature_t fid;
810         zfeature_info_t *fi;
812         ret = zfeature_lookup_name(feature, &fid);
812         ret = zfeature_lookup_name(feature, &fi);
812         if (ret != 0) {
813             (void) strlcpy(buf, "-", len);
814             return (ENOTSUP);
815         }
816         feature = spa_feature_table[fid].fi_guid;
817         feature = fi->fi_guid;
819     if (nvlist_lookup_uint64(features, feature, &refcount) == 0)
820         found = B_TRUE;
822     if (supported) {
823         if (!found) {
824             (void) strlcpy(buf, ZFS_FEATURE_DISABLED, len);
825         } else {
826             if (refcount == 0)
827                 (void) strlcpy(buf, ZFS_FEATURE_ENABLED, len);
828             else
829                 (void) strlcpy(buf, ZFS_FEATURE_ACTIVE, len);
830         }
831     } else {
832         if (found) {
833             if (refcount == 0) {
834                 (void) strcpy(buf, ZFS_UNSUPPORTED_INACTIVE);
835             } else {
836                 (void) strcpy(buf, ZFS_UNSUPPORTED_READONLY);
837             }
838         } else {
839             (void) strlcpy(buf, "-", len);
840             return (ENOTSUP);
841         }
842     }
844     return (0);
845 }

```

unchanged portion omitted



```

*****
9655 Tue Oct 1 14:04:22 2013
new/usr/src/man/man5/zpool-features.5
4171 clean up spa_feature_*() interfaces
4172 implement extensible_dataset feature for use by other zpool features
Reviewed by: Max Grossman <max.grossman@delphix.com>
Reviewed by: Christopher Siden <christopher.siden@delphix.com>
Reviewed by: George Wilson <george.wilson@delphix.com>
*****
1 \" te
2.\" Copyright (c) 2013 by Delphix. All rights reserved.
3.\" Copyright (c) 2013 by Saso Kiselkov. All rights reserved.
4.\" Copyright (c) 2013, Joyent, Inc. All rights reserved.
5.\" The contents of this file are subject to the terms of the Common Development
6.\" and Distribution License (the \"License\"). You may not use this file except
7.\" in compliance with the License. You can obtain a copy of the license at
8.\" usr/src/OPENSOLARIS.LICENSE or http://www.opensolaris.org/os/licensing.
9.\"
10.\" See the License for the specific language governing permissions and
11.\" limitations under the License. When distributing Covered Code, include this
12.\" CDDL HEADER in each file and include the License file at
13.\" usr/src/OPENSOLARIS.LICENSE. If applicable, add the following below this
14.\" CDDL HEADER, with the fields enclosed by brackets \"[]\" replaced with your
15.\" own identifying information:
16.\" Portions Copyright [yyyy] [name of copyright owner]
17.TH ZPOOL-FEATURES 5 \"Aug 27, 2013\"
18.SH NAME
19 zpool\--features \- ZFS pool feature descriptions
20.SH DESCRIPTION
21 .sp
22 .LP
23 ZFS pool on\--disk format versions are specified via \"features\" which replace
24 the old on\--disk format numbers (the last supported on\--disk format number is
25 28). To enable a feature on a pool use the \fBbupgrade\fR subcommand of the
26 \fBzpool\fR(1M) command, or set the \fBfeature@\fR\fIfeature_name\fR property
27 to \fBenabled\fR.
28 .sp
29 .LP
30 The pool format does not affect file system version compatibility or the ability
31 to send file systems between pools.
32 .sp
33 .LP
34 Since most features can be enabled independently of each other the on\--disk
35 format of the pool is specified by the set of all features marked as
36 \fBactive\fR on the pool. If the pool was created by another software version
37 this set may include unsupported features.
38 .SS \"Identifying features\"
39 .sp
40 .LP
41 Every feature has a guid of the form \fIcom.example:feature_name\fR. The reverse
42 DNS name ensures that the feature's guid is unique across all ZFS
43 implementations. When unsupported features are encountered on a pool they will
44 be identified by their guids. Refer to the documentation for the ZFS
45 implementation that created the pool for information about those features.
46 .sp
47 .LP
48 Each supported feature also has a short name. By convention a feature's short
49 name is the portion of its guid which follows the ':' (e.g.
50 \fIcom.example:feature_name\fR would have the short name \fIfeature_name\fR),
51 however a feature's short name may differ across ZFS implementations if
52 following the convention would result in name conflicts.
53 .SS \"Feature states\"
54 .sp
55 .LP
56 Features can be in one of three states:
57 .sp

```

```

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

```

```

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

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

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

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

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

175 .sp
176 .ne 2
177 .na
178 \fB\fBempty_bpobj\fR\fR
179 .ad
180 .RS 4n
181 .TS
182 l l .
183 GUID      com.delphix:empty_bpobj
184 READ\ONLY COMPATIBLE  yes
185 DEPENDENCIES      none
186 .TE

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

```

```

190 the disk space required.

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

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

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

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

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

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

247 This features allows ZFS to maintain more information about how free space
248 is organized within the pool. If this feature is \fB\fBenabled\fR, ZFS will
249 set this feature to \fB\fBactive\fR when a new space map object is created or
250 an existing space map is upgraded to the new format. Once the feature is
251 \fB\fBactive\fR, it will remain in that state until the pool is destroyed.
252 .RE

254 .sp
255 .ne 2

```

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

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

271 When the \fBmulti\_vdev\_crash\_dump\fR feature is set to \fBEnabled\fR,  
272 the administrator can use the \fBdumpadm\fR(1M) command to configure a  
273 dump device on a pool comprised of multiple vdevs.  
274 **.RE**

```
276 .sp
277 .ne 2
278 .na
279 \fB\fBextensible_dataset\fR\fR
280 .ad
281 .RS 4n
282 .TS
283 l l .
284 GUID      com.delphix:extensible_dataset
285 READ\ -ONLY COMPATIBLE  no
286 DEPENDENCIES      none
287 .TE
```

289 This feature allows more flexible use of internal ZFS data structures,  
290 and exists for other features to depend on.

292 This feature will be \fBActive\fR when the first dependent feature uses it,  
293 and will be returned to the \fBEnabled\fR state when all datasets that use  
294 this feature are destroyed.

296 **.RE**

```
298 .SH "SEE ALSO"
299 \fBzpool\fR(1M)
```

```

*****
15097 Tue Oct 1 14:04:24 2013
new/usr/src/uts/common/fs/zfs/bpobj.c
4171 clean up spa_feature_*() interfaces
4172 implement extensible_dataset feature for use by other zpool features
Reviewed by: Max Grossman <max.grossman@delphix.com>
Reviewed by: Christopher Siden <christopher.siden@delphix.com>
Reviewed by: George Wilson <george.wilson@delphix.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22  * Copyright (c) 2005, 2010, Oracle and/or its affiliates. All rights reserved.
23  * Copyright (c) 2013 by Delphix. All rights reserved.
24 */

26 #include <sys/bpobj.h>
27 #include <sys/zfs_context.h>
28 #include <sys/refcount.h>
29 #include <sys/dsl_pool.h>
30 #include <sys/zfeature.h>
31 #include <sys/zap.h>

33 /*
34  * Return an empty bpobj, preferably the empty dummy one (dp_empty_bpobj).
35  */
36 uint64_t
37 bpobj_alloc_empty(objset_t *os, int blocksize, dmu_tx_t *tx)
38 {
39     zfeature_info_t *empty_bpobj_feat =
40         &spa_feature_table[SPA_FEATURE_EMPTY_BPOBJ];
41     spa_t *spa = dmu_objset_spa(os);
42     dsl_pool_t *dp = dmu_objset_pool(os);

44     if (spa_feature_is_enabled(spa, SPA_FEATURE_EMPTY_BPOBJ)) {
45         if (!spa_feature_is_active(spa, SPA_FEATURE_EMPTY_BPOBJ)) {
46             if (spa_feature_is_enabled(spa, empty_bpobj_feat)) {
47                 if (!spa_feature_is_active(spa, empty_bpobj_feat)) {
48                     ASSERT0(dp->dp_empty_bpobj);
49                     dp->dp_empty_bpobj =
50                         bpobj_alloc(os, SPA_MAXBLOCKSIZE, tx);
51                     VERIFY(zap_add(os,
52                         DMU_POOL_DIRECTORY_OBJECT,
53                         DMU_POOL_EMPTY_BPOBJ, sizeof (uint64_t), 1,
54                         &dp->dp_empty_bpobj, tx) == 0);
55                 }
56             }
57         }
58     }
59     spa_feature_incr(spa, SPA_FEATURE_EMPTY_BPOBJ, tx);
60     spa_feature_incr(spa, empty_bpobj_feat, tx);

```

```

53     ASSERT(dp->dp_empty_bpobj != 0);
54     return (dp->dp_empty_bpobj);
55 } else {
56     return (bpobj_alloc(os, blocksize, tx));
57 }
58 }

60 void
61 bpobj_decr_empty(objset_t *os, dmu_tx_t *tx)
62 {
63     zfeature_info_t *empty_bpobj_feat =
64         &spa_feature_table[SPA_FEATURE_EMPTY_BPOBJ];
65     dsl_pool_t *dp = dmu_objset_pool(os);

67     spa_feature_decr(dmu_objset_spa(os), SPA_FEATURE_EMPTY_BPOBJ, tx);
68     if (!spa_feature_is_active(dmu_objset_spa(os),
69         SPA_FEATURE_EMPTY_BPOBJ)) {
70         spa_feature_decr(dmu_objset_spa(os), empty_bpobj_feat, tx);
71         if (!spa_feature_is_active(dmu_objset_spa(os), empty_bpobj_feat)) {
72             VERIFY3U(0, ==, zap_remove(dp->dp_meta_objset,
73                 DMU_POOL_DIRECTORY_OBJECT,
74                 DMU_POOL_EMPTY_BPOBJ, tx));
75             VERIFY3U(0, ==, dmu_object_free(os, dp->dp_empty_bpobj, tx));
76             dp->dp_empty_bpobj = 0;
77         }
78     }
79 }
_____unchanged_portion_omitted_____

195 static int
196 bpobj_iterate_impl(bpobj_t *bpo, bpobj_itor_t func, void *arg, dmu_tx_t *tx,
197     boolean_t free)
198 {
199     dmu_object_info_t doi;
200     int epb;
201     int64_t i;
202     int err = 0;
203     dmu_buf_t *dbuf = NULL;

205     mutex_enter(&bpo->bpo_lock);

207     if (free)
208         dmu_buf_will_dirty(bpo->bpo_dbuf, tx);

210     for (i = bpo->bpo_phys->bpo_num_blkptrs - 1; i >= 0; i--) {
211         blkptr_t *bpararray;
212         blkptr_t *bp;
213         uint64_t offset, blkoff;

215         offset = i * sizeof (blkptr_t);
216         blkoff = P2PHASE(i, bpo->bpo_epb);

218         if (dbuf == NULL || dbuf->db_offset > offset) {
219             if (dbuf)
220                 dmu_buf_rele(dbuf, FTAG);
221             err = dmu_buf_hold(bpo->bpo_os, bpo->bpo_object, offset,
222                 FTAG, &dbuf, 0);
223             if (err)
224                 break;
225         }

227         ASSERT3U(offset, >=, dbuf->db_offset);
228         ASSERT3U(offset, <, dbuf->db_offset + dbuf->db_size);

230         bpararray = dbuf->db_data;
231         bp = &bpararray[blkoff];
232         err = func(arg, bp, tx);

```

```

233         if (err)
234             break;
235     if (free) {
236         bpo->bpo_phys->bpo_bytes -=
237             bp_get_dsize_sync(dmu_objset_spa(bpo->bpo_os), bp);
238         ASSERT3S(bpo->bpo_phys->bpo_bytes, >=, 0);
239         if (bpo->bpo_havecomp) {
240             bpo->bpo_phys->bpo_comp -= BP_GET_PSIZE(bp);
241             bpo->bpo_phys->bpo_uncomp -= BP_GET_UCSIZE(bp);
242         }
243         bpo->bpo_phys->bpo_num_blkptrs--;
244         ASSERT3S(bpo->bpo_phys->bpo_num_blkptrs, >=, 0);
245     }
246 }
247 if (dbuf) {
248     dmu_buf_rele(dbuf, FTAG);
249     dbuf = NULL;
250 }
251 if (free) {
252     i++;
253     VERIFY3U(0, ==, dmu_free_range(bpo->bpo_os, bpo->bpo_object,
254         i * sizeof(blkptr_t), -1ULL, tx));
255 }
256 if (err || !bpo->bpo_havesubobj || bpo->bpo_phys->bpo_subobjs == 0)
257     goto out;
258
259 ASSERT(bpo->bpo_havecomp);
260 err = dmu_object_info(bpo->bpo_os, bpo->bpo_phys->bpo_subobjs, &doi);
261 if (err) {
262     mutex_exit(&bpo->bpo_lock);
263     return (err);
264 }
265 ASSERT3U(doi.doi_type, ==, DMU_OT_BPOBJ_SUBOBJ);
266 epb = doi.doi_data_block_size / sizeof(uint64_t);
267
268 for (i = bpo->bpo_phys->bpo_num_subobjs - 1; i >= 0; i--) {
269     uint64_t *objarray;
270     uint64_t offset, blkoff;
271     bpobj_t sublist;
272     uint64_t used_before, comp_before, uncomp_before;
273     uint64_t used_after, comp_after, uncomp_after;
274
275     offset = i * sizeof(uint64_t);
276     blkoff = P2PHASE(i, epb);
277
278     if (dbuf == NULL || dbuf->db_offset > offset) {
279         if (dbuf)
280             dmu_buf_rele(dbuf, FTAG);
281         err = dmu_buf_hold(bpo->bpo_os,
282             bpo->bpo_phys->bpo_subobjs, offset, FTAG, &dbuf, 0);
283         if (err)
284             break;
285     }
286
287     ASSERT3U(offset, >=, dbuf->db_offset);
288     ASSERT3U(offset, <, dbuf->db_offset + dbuf->db_size);
289
290     objarray = dbuf->db_data;
291     err = bpobj_open(&sublist, bpo->bpo_os, objarray[blkoff]);
292     if (err)
293         break;
294     if (free) {
295         err = bpobj_space(&sublist,
296             &used_before, &comp_before, &uncomp_before);
297         if (err)
298             break;

```

```

299     }
300     err = bpobj_iterate_impl(&sublist, func, arg, tx, free);
301     if (free) {
302         VERIFY3U(0, ==, bpobj_space(&sublist,
303             &used_after, &comp_after, &uncomp_after));
304         bpo->bpo_phys->bpo_bytes -= used_before - used_after;
305         ASSERT3S(bpo->bpo_phys->bpo_bytes, >=, 0);
306         bpo->bpo_phys->bpo_comp -= comp_before - comp_after;
307         bpo->bpo_phys->bpo_uncomp -=
308             uncomp_before - uncomp_after;
309     }
310
311     bpobj_close(&sublist);
312     if (err)
313         break;
314     if (free) {
315         err = dmu_object_free(bpo->bpo_os,
316             objarray[blkoff], tx);
317         if (err)
318             break;
319         bpo->bpo_phys->bpo_num_subobjs--;
320         ASSERT3S(bpo->bpo_phys->bpo_num_subobjs, >=, 0);
321     }
322 }
323 if (dbuf) {
324     dmu_buf_rele(dbuf, FTAG);
325     dbuf = NULL;
326 }
327 if (free) {
328     VERIFY3U(0, ==, dmu_free_range(bpo->bpo_os,
329         bpo->bpo_phys->bpo_subobjs,
330         (i + 1) * sizeof(uint64_t), -1ULL, tx));
331 }
332
333 out:
334     /* If there are no entries, there should be no bytes. */
335     ASSERT(bpo->bpo_phys->bpo_num_blkptrs > 0 ||
336         (bpo->bpo_havesubobj && bpo->bpo_phys->bpo_num_subobjs > 0) ||
337         bpo->bpo_phys->bpo_bytes == 0);
338
339     mutex_exit(&bpo->bpo_lock);
340     return (err);
341 }

```

unchanged portion omitted

```

*****
6467 Tue Oct 1 14:04:26 2013
new/usr/src/uts/common/fs/zfs/dmu_object.c
4171 clean up spa_feature*() interfaces
4172 implement extensible_dataset feature for use by other zpool features
Reviewed by: Max Grossman <max.grossman@delphix.com>
Reviewed by: Christopher Siden <christopher.siden@delphix.com>
Reviewed by: George Wilson <george.wilson@delphix.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22  * Copyright (c) 2005, 2010, Oracle and/or its affiliates. All rights reserved.
23  * Copyright (c) 2013 by Delphix. All rights reserved.
24  */

26 #include <sys/dmu.h>
27 #include <sys/dmu_objset.h>
28 #include <sys/dmu_tx.h>
29 #include <sys/dnode.h>
30 #include <sys/zap.h>
31 #include <sys/zfeature.h>

33 uint64_t
34 dmu_object_alloc(objset_t *os, dmu_object_type_t ot, int blocksize,
35                 dmu_object_type_t bonustype, int bonuslen, dmu_tx_t *tx)
36 {
37     uint64_t object;
38     uint64_t L2_dnode_count = DNODES_PER_BLOCK <<
39         (DMU_META_DNODE(os)->dn_indblkshift - SPA_BLKPTRSHIFT);
40     dnode_t *dn = NULL;
41     int restarted = B_FALSE;

43     mutex_enter(&os->os_obj_lock);
44     for (;;) {
45         object = os->os_obj_next;
46         /*
47          * Each time we polish off an L2 bp worth of dnodes
48          * (2^13 objects), move to another L2 bp that's still
49          * reasonably sparse (at most 1/4 full). Look from the
50          * beginning once, but after that keep looking from here.
51          * If we can't find one, just keep going from here.
52          */
53         if (P2PHASE(object, L2_dnode_count) == 0) {
54             uint64_t offset = restarted ? object << DNODE_SHIFT : 0;
55             int error = dnode_next_offset(DMU_META_DNODE(os),
56                                         DNODE_FIND_HOLE,
57                                         &offset, 2, DNODES_PER_BLOCK >> 2, 0);

```

```

58         restarted = B_TRUE;
59         if (error == 0)
60             object = offset >> DNODE_SHIFT;
61     }
62     os->os_obj_next = ++object;

64     /*
65      * XXX We should check for an i/o error here and return
66      * up to our caller. Actually we should pre-read it in
67      * dmu_tx_assign(), but there is currently no mechanism
68      * to do so.
69      */
70     (void) dnode_hold_impl(os, object, DNODE_MUST_BE_FREE,
71                          FTAG, &dn);
72     if (dn)
73         break;

75     if (dmu_object_next(os, &object, B_TRUE, 0) == 0)
76         os->os_obj_next = object - 1;
77 }

79     dnode_allocate(dn, ot, blocksize, 0, bonustype, bonuslen, tx);
80     dnode_rele(dn, FTAG);

82     mutex_exit(&os->os_obj_lock);

84     dmu_tx_add_new_object(tx, os, object);
85     return (object);
86 }

unchanged_portion_omitted

187 int
188 dmu_object_next(objset_t *os, uint64_t *objectp, boolean_t hole, uint64_t txg)
189 {
190     uint64_t offset = (*objectp + 1) << DNODE_SHIFT;
191     int error;

193     error = dnode_next_offset(DMU_META_DNODE(os),
194                             (hole ? DNODE_FIND_HOLE : 0), &offset, 0, DNODES_PER_BLOCK, txg);

196     *objectp = offset >> DNODE_SHIFT;

198     return (error);
199 }

201 /*
202  * Turn this object from old_type into DMU_OTN_ZAP_METADATA, and bump the
203  * refcount on SPA_FEATURE_EXTENSIBLE_DATASET.
204  *
205  * Only for use from syncing context, on MOS objects.
206  */
207 void
208 dmu_object_zapify(objset_t *mos, uint64_t object, dmu_object_type_t old_type,
209                 dmu_tx_t *tx)
210 {
211     dnode_t *dn;

213     ASSERT(dmu_tx_is_syncing(tx));

215     VERIFY0(dnode_hold(mos, object, FTAG, &dn));
216     if (dn->dn_type == DMU_OTN_ZAP_METADATA) {
217         dnode_rele(dn, FTAG);
218         return;
219     }
220     ASSERT3U(dn->dn_type, ==, old_type);
221     ASSERT0(dn->dn_maxblkid);

```

```
222     dn->dn_next_type[tx->tx_txg & TXG_MASK] = dn->dn_type =
223         DMU_OTN_ZAP_METADATA;
224     dnode_setdirty(dn, tx);
225     dnode_rele(dn, FTAG);
226
227     mzap_create_impl(mos, object, 0, 0, tx);
228
229     spa_feature_incr(dmu_objset_spa(mos),
230         SPA_FEATURE_EXTENSIBLE_DATASET, tx);
231 }
232
233 void
234 dmu_object_free_zapified(objset_t *mos, uint64_t object, dmu_tx_t *tx)
235 {
236     dnode_t *dn;
237     dmu_object_type_t t;
238
239     ASSERT(dmu_tx_is_syncing(tx));
240
241     VERIFY0(dnode_hold(mos, object, FTAG, &dn));
242     t = dn->dn_type;
243     dnode_rele(dn, FTAG);
244
245     if (t == DMU_OTN_ZAP_METADATA) {
246         spa_feature_decr(dmu_objset_spa(mos),
247             SPA_FEATURE_EXTENSIBLE_DATASET, tx);
248     }
249     VERIFY0(dmu_object_free(mos, object, tx));
250 }
251
252 unchanged_portion_omitted
```

```

*****
16666 Tue Oct 1 14:04:28 2013
new/usr/src/uts/common/fs/zfs/dmu_traverse.c
4171 clean up spa_feature_*() interfaces
4172 implement extensible_dataset feature for use by other zpool features
Reviewed by: Max Grossman <max.grossman@delphix.com>
Reviewed by: Christopher Siden <christopher.siden@delphix.com>
Reviewed by: George Wilson <george.wilson@delphix.com>
*****
_____unchanged_portion_omitted_____

576 /*
577  * NB: pool must not be changing on-disk (eg, from zdb or sync context).
578  */
579 int
580 traverse_pool(spa_t *spa, uint64_t txg_start, int flags,
581             blkptr_cb_t func, void *arg)
582 {
583     int err, lasterr = 0;
584     uint64_t obj;
585     dsl_pool_t *dp = spa_get_dsl(spa);
586     objset_t *mos = dp->dp_meta_objset;
587     boolean_t hard = (flags & TRAVERSE_HARD);

589     /* visit the MOS */
590     err = traverse_impl(spa, NULL, 0, spa_get_rootblkptr(spa),
591                      txg_start, NULL, flags, func, arg);
592     if (err != 0)
593         return (err);

595     /* visit each dataset */
596     for (obj = 1; err == 0 || (err != ESRCH && hard);
597          err = dmu_object_next(mos, &obj, FALSE, txg_start)) {
598         dmu_object_info_t doi;

600         err = dmu_object_info(mos, obj, &doi);
601         if (err != 0) {
602             if (!hard)
603                 return (err);
604             lasterr = err;
605             continue;
606         }

608         if (doi.doi_bonus_type == DMU_OT_DSL_DATASET) {
609             if (doi.doi_type == DMU_OT_DSL_DATASET) {
610                 dsl_dataset_t *ds;
611                 uint64_t txg = txg_start;

612                 dsl_pool_config_enter(dp, FTAG);
613                 err = dsl_dataset_hold_obj(dp, obj, FTAG, &ds);
614                 dsl_pool_config_exit(dp, FTAG);
615                 if (err != 0) {
616                     if (!hard)
617                         return (err);
618                     lasterr = err;
619                     continue;
620                 }
621                 if (ds->ds_phys->ds_prev_snap_txg > txg)
622                     txg = ds->ds_phys->ds_prev_snap_txg;
623                 err = traverse_dataset(ds, txg, flags, func, arg);
624                 dsl_dataset_rele(ds, FTAG);
625                 if (err != 0) {
626                     if (!hard)
627                         return (err);
628                     lasterr = err;
629                 }

```

```

630     }
631     }
632     if (err == ESRCH)
633         err = 0;
634     return (err != 0 ? err : lasterr);
635 }
_____unchanged_portion_omitted_____

```



```

*****
19442 Tue Oct 1 14:04:30 2013
new/usr/src/uts/common/fs/zfs/dnode_sync.c
4171 clean up spa_feature_*() interfaces
4172 implement extensible_dataset feature for use by other zpool features
Reviewed by: Max Grossman <max.grossman@delphix.com>
Reviewed by: Christopher Siden <christopher.siden@delphix.com>
Reviewed by: George Wilson <george.wilson@delphix.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */

22 /*
23  * Copyright (c) 2005, 2010, Oracle and/or its affiliates. All rights reserved.
24  * Copyright (c) 2013 by Delphix. All rights reserved.
24  * Copyright (c) 2012 by Delphix. All rights reserved.
25  */

27 #include <sys/zfs_context.h>
28 #include <sys/dbuf.h>
29 #include <sys/dnode.h>
30 #include <sys/dmu.h>
31 #include <sys/dmu_tx.h>
32 #include <sys/dmu_objset.h>
33 #include <sys/dsl_dataset.h>
34 #include <sys/spa.h>

36 static void
37 dnode_increase_indirection(dnode_t *dn, dmu_tx_t *tx)
38 {
39     dmu_buf_impl_t *db;
40     int txgoff = tx->tx_txg & TXG_MASK;
41     int nblkptr = dn->dn_phys->dn_nblkptr;
42     int old_toplvl = dn->dn_phys->dn_nlevels - 1;
43     int new_level = dn->dn_next_nlevels[txgoff];
44     int i;

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

48     /* this dnode can't be paged out because it's dirty */
49     ASSERT(dn->dn_phys->dn_type != DMU_OT_NONE);
50     ASSERT(RW_WRITE_HELD(&dn->dn_struct_rwlock));
51     ASSERT(new_level > 1 && dn->dn_phys->dn_nlevels > 0);

53     db = dbuf_hold_level(dn, dn->dn_phys->dn_nlevels, 0, FTAG);
54     ASSERT(db != NULL);

56     dn->dn_phys->dn_nlevels = new_level;

```

```

57     dprintf("os=%p obj=%llu, increase to %d\n", dn->dn_objset,
58            dn->dn_object, dn->dn_phys->dn_nlevels);

60     /* check for existing blkptrs in the dnode */
61     for (i = 0; i < nblkptr; i++)
62         if (!BP_IS_HOLE(&dn->dn_phys->dn_blkptr[i]))
63             break;
64     if (i != nblkptr) {
65         /* transfer dnode's block pointers to new indirect block */
66         (void) dbuf_read(db, NULL, DB_RF_MUST_SUCCEED|DB_RF_HAVESTRUCT);
67         ASSERT(db->db_data);
68         ASSERT(arc_released(db->db_buf));
69         ASSERT3U(sizeof (blkptr_t) * nblkptr, <=, db->db.db_size);
70         bcopy(dn->dn_phys->dn_blkptr, db->db.db_data,
71              sizeof (blkptr_t) * nblkptr);
72         arc_buf_freeze(db->db_buf);
73     }

75     /* set dbuf's parent pointers to new indirect buf */
76     for (i = 0; i < nblkptr; i++) {
77         dmu_buf_impl_t *child = dbuf_find(dn, old_toplvl, i);

79         if (child == NULL)
80             continue;
81 #ifdef DEBUG
82         DB_DNODE_ENTER(child);
83         ASSERT3P(DB_DNODE(child), ==, dn);
84         DB_DNODE_EXIT(child);
85 #endif /* DEBUG */
86         if (child->db_parent && child->db_parent != dn->dn_dbuf) {
87             ASSERT(child->db_parent->db_level == db->db_level);
88             ASSERT(child->db_blkptr !=
89                    &dn->dn_phys->dn_blkptr[child->db_blkid]);
90             mutex_exit(&child->db_mtx);
91             continue;
92         }
93         ASSERT(child->db_parent == NULL ||
94                child->db_parent == dn->dn_dbuf);

96         child->db_parent = db;
97         dbuf_add_ref(db, child);
98         if (db->db_data)
99             child->db_blkptr = (blkptr_t *)db->db_data + i;
100        else
101            child->db_blkptr = NULL;
102        dprintf_dbuf_bp(child, child->db_blkptr,
103            "changed db_blkptr to new indirect %s", "");

105        mutex_exit(&child->db_mtx);
106    }

108    bzero(dn->dn_phys->dn_blkptr, sizeof (blkptr_t) * nblkptr);
110    dbuf_rele(db, FTAG);

112    rw_exit(&dn->dn_struct_rwlock);
113 }

unchanged portion omitted

523 /*
524  * Write out the dnode's dirty buffers.
525  */
526 void
527 dnode_sync(dnode_t *dn, dmu_tx_t *tx)
528 {
529     free_range_t *rp;

```

```

530     dnode_phys_t *dnp = dn->dn_phys;
531     int txgoff = tx->tx_txg & TXG_MASK;
532     list_t *list = &dn->dn_dirty_records[txgoff];
533     static const dnode_phys_t zerodn = { 0 };
534     boolean_t kill_spill = B_FALSE;

536     ASSERT(dmu_tx_is_syncing(tx));
537     ASSERT(dnp->dn_type != DMU_OT_NONE || dn->dn_allocated_txg);
538     ASSERT(dnp->dn_type != DMU_OT_NONE || dn->dn_allocated_txg);
539     bcmp(dnp, &zerodn, DNODE_SIZE) == 0);
540     DNODE_VERIFY(dn);

542     ASSERT(dn->dn_dbuf == NULL || arc_released(dn->dn_dbuf->db_buf));

544     if (dmu_objset_userused_enabled(dn->dn_objset) &&
545         !DMU_OBJECT_IS_SPECIAL(dn->dn_objset)) {
546         mutex_enter(&dn->dn_mtx);
547         dn->dn_oldused = DN_USED_BYTES(dn->dn_phys);
548         dn->dn_oldflags = dn->dn_phys->dn_flags;
549         dn->dn_phys->dn_flags |= DNODE_FLAG_USERUSED_ACCOUNTED;
550         mutex_exit(&dn->dn_mtx);
551         dmu_objset_userquota_get_ids(dn, B_FALSE, tx);
552     } else {
553         /* Once we account for it, we should always account for it. */
554         ASSERT(!(dn->dn_phys->dn_flags &
555             DNODE_FLAG_USERUSED_ACCOUNTED));
556     }

558     mutex_enter(&dn->dn_mtx);
559     if (dn->dn_allocated_txg == tx->tx_txg) {
560         /* The dnode is newly allocated or reallocated */
561         if (dnp->dn_type == DMU_OT_NONE) {
562             /* this is a first alloc, not a realloc */
563             dnp->dn_nlevels = 1;
564             dnp->dn_nblkptr = dn->dn_nblkptr;
565         }

567         dnp->dn_type = dn->dn_type;
568         dnp->dn_bonustype = dn->dn_bonustype;
569         dnp->dn_bonuslen = dn->dn_bonuslen;
570     }

572     ASSERT(dnp->dn_nlevels > 1 ||
573         BP_IS_HOLE(&dnp->dn_blkptr[0]) ||
574         BP_GET_LSIZE(&dnp->dn_blkptr[0]) ==
575         dnp->dn_datablkszsec << SPA_MINBLOCKSHIFT);

577     if (dn->dn_next_type[txgoff] != 0) {
578         dnp->dn_type = dn->dn_type;
579         dn->dn_next_type[txgoff] = 0;
580     }

582     if (dn->dn_next_blkisz[txgoff] != 0) {
577     if (dn->dn_next_blkisz[txgoff]) {
583         ASSERT(P2PHASE(dn->dn_next_blkisz[txgoff],
584             SPA_MINBLOCKSIZE) == 0);
585         ASSERT(BP_IS_HOLE(&dnp->dn_blkptr[0]) ||
586             dn->dn_maxblkid == 0 || list_head(list) != NULL ||
587             avl_last(&dn->dn_ranges[txgoff]) ||
588             dn->dn_next_blkisz[txgoff] >> SPA_MINBLOCKSHIFT ==
589             dnp->dn_datablkszsec);
590         dnp->dn_datablkszsec =
591             dn->dn_next_blkisz[txgoff] >> SPA_MINBLOCKSHIFT;
592         dn->dn_next_blkisz[txgoff] = 0;
593     }

```

```

595     if (dn->dn_next_bonuslen[txgoff] != 0) {
596         if (dn->dn_next_bonuslen[txgoff]) {
597             if (dn->dn_next_bonuslen[txgoff] == DN_ZERO_BONUSLEN)
598                 dnp->dn_bonuslen = 0;
599             else
600                 dnp->dn_bonuslen = dn->dn_next_bonuslen[txgoff];
601             ASSERT(dnp->dn_bonuslen <= DN_MAX_BONUSLEN);
602             dn->dn_next_bonuslen[txgoff] = 0;
603         }

604     if (dn->dn_next_bonustype[txgoff] != 0) {
599     if (dn->dn_next_bonustype[txgoff]) {
605         ASSERT(DMU_OT_IS_VALID(dn->dn_next_bonustype[txgoff]));
606         dnp->dn_bonustype = dn->dn_next_bonustype[txgoff];
607         dn->dn_next_bonustype[txgoff] = 0;
608     }

610     /*
611     * We will either remove a spill block when a file is being removed
612     * or we have been asked to remove it.
613     */
614     if (dn->dn_rm_spillblk[txgoff] ||
615         ((dnp->dn_flags & DNODE_FLAG_SPILL_BLKPTR) &&
616             dn->dn_free_txg > 0 && dn->dn_free_txg <= tx->tx_txg)) {
617         if ((dnp->dn_flags & DNODE_FLAG_SPILL_BLKPTR))
618             kill_spill = B_TRUE;
619         dn->dn_rm_spillblk[txgoff] = 0;
620     }

622     if (dn->dn_next_indblkshift[txgoff] != 0) {
617     if (dn->dn_next_indblkshift[txgoff]) {
623         ASSERT(dnp->dn_nlevels == 1);
624         dnp->dn_indblkshift = dn->dn_next_indblkshift[txgoff];
625         dn->dn_next_indblkshift[txgoff] = 0;
626     }

628     /*
629     * Just take the live (open-context) values for checksum and compress.
630     * Strictly speaking it's a future leak, but nothing bad happens if we
631     * start using the new checksum or compress algorithm a little early.
632     */
633     dnp->dn_checksum = dn->dn_checksum;
634     dnp->dn_compress = dn->dn_compress;

636     mutex_exit(&dn->dn_mtx);

638     if (kill_spill) {
639         (void) free_blocks(dn, &dn->dn_phys->dn_spill, 1, tx);
640         mutex_enter(&dn->dn_mtx);
641         dnp->dn_flags &= ~DNODE_FLAG_SPILL_BLKPTR;
642         mutex_exit(&dn->dn_mtx);
643     }

645     /* process all the "freed" ranges in the file */
646     while (rp = avl_last(&dn->dn_ranges[txgoff])) {
647         dnode_sync_free_range(dn, rp->fr_blkid, rp->fr_nblks, tx);
648         /* grab the mutex so we don't race with dnode_block_freed() */
649         mutex_enter(&dn->dn_mtx);
650         avl_remove(&dn->dn_ranges[txgoff], rp);
651         mutex_exit(&dn->dn_mtx);
652         kmem_free(rp, sizeof (free_range_t));
653     }

655     if (dn->dn_free_txg > 0 && dn->dn_free_txg <= tx->tx_txg) {
656         dnode_sync_free(dn, tx);
657         return;

```

```
658     }
660     if (dn->dn_next_nblkptr[txgoff]) {
661         /* this should only happen on a realloc */
662         ASSERT(dn->dn_allocated_txg == tx->tx_txg);
663         if (dn->dn_next_nblkptr[txgoff] > dnp->dn_nblkptr) {
664             /* zero the new blkptrs we are gaining */
665             bzero(dnp->dn_blkptr + dnp->dn_nblkptr,
666                 sizeof (blkptr_t) *
667                 (dn->dn_next_nblkptr[txgoff] - dnp->dn_nblkptr));
668 #ifdef ZFS_DEBUG
669         } else {
670             int i;
671             ASSERT(dn->dn_next_nblkptr[txgoff] < dnp->dn_nblkptr);
672             /* the blkptrs we are losing better be unallocated */
673             for (i = dn->dn_next_nblkptr[txgoff];
674                 i < dnp->dn_nblkptr; i++)
675                 ASSERT(BP_IS_HOLE(&dnp->dn_blkptr[i]));
676 #endif
677         }
678         mutex_enter(&dn->dn_mtx);
679         dnp->dn_nblkptr = dn->dn_next_nblkptr[txgoff];
680         dn->dn_next_nblkptr[txgoff] = 0;
681         mutex_exit(&dn->dn_mtx);
682     }
684     if (dn->dn_next_nlevels[txgoff]) {
685         dnode_increase_indirection(dn, tx);
686         dn->dn_next_nlevels[txgoff] = 0;
687     }
689     dbuf_sync_list(list, tx);
691     if (!DMU_OBJECT_IS_SPECIAL(dn->dn_object)) {
692         ASSERT3P(list_head(list), ==, NULL);
693         dnode_rele(dn, (void *) (uintptr_t) tx->tx_txg);
694     }
696     /*
697     * Although we have dropped our reference to the dnode, it
698     * can't be evicted until its written, and we haven't yet
699     * initiated the IO for the dnode's dbuf.
700     */
701 }
```

unchanged\_portion\_omitted

```

*****
82608 Tue Oct 1 14:04:31 2013
new/usr/src/uts/common/fs/zfs/dsl_dataset.c
4171 clean up spa_feature_*() interfaces
4172 implement extensible_dataset feature for use by other zpool features
Reviewed by: Max Grossman <max.grossman@delphix.com>
Reviewed by: Christopher Siden <christopher.siden@delphix.com>
Reviewed by: George Wilson <george.wilson@delphix.com>
*****
_____unchanged_portion_omitted_____

337 int
338 dsl_dataset_hold_obj(dsl_pool_t *dp, uint64_t dsobj, void *tag,
339     dsl_dataset_t **dsp)
340 {
341     objset_t *mos = dp->dp_meta_objset;
342     dmu_buf_t *dbuf;
343     dsl_dataset_t *ds;
344     int err;
345     dmu_object_info_t doi;

347     ASSERT(dsl_pool_config_held(dp));

349     err = dmu_bonus_hold(mos, dsobj, tag, &dbuf);
350     if (err != 0)
351         return (err);

353     /* Make sure dsobj has the correct object type. */
354     dmu_object_info_from_db(dbuf, &doi);
355     if (doi.doi_bonus_type != DMU_OT_DSL_DATASET) {
356         if (doi.doi_type != DMU_OT_DSL_DATASET) {
357             dmu_buf_rele(dbuf, tag);
358             return (SET_ERROR(EINVAL));
359         }

360     ds = dmu_buf_get_user(dbuf);
361     if (ds == NULL) {
362         dsl_dataset_t *winner = NULL;

364         ds = kmem_zalloc(sizeof (dsl_dataset_t), KM_SLEEP);
365         ds->ds_dbuf = dbuf;
366         ds->ds_object = dsobj;
367         ds->ds_phys = dbuf->db_data;

369         mutex_init(&ds->ds_lock, NULL, MUTEX_DEFAULT, NULL);
370         mutex_init(&ds->ds_opening_lock, NULL, MUTEX_DEFAULT, NULL);
371         mutex_init(&ds->ds_sendstream_lock, NULL, MUTEX_DEFAULT, NULL);
372         refcount_create(&ds->ds_longholds);

374         bplist_create(&ds->ds_pending_deadlist);
375         dsl_deadlist_open(&ds->ds_deadlist,
376             mos, ds->ds_phys->ds_deadlist_obj);

378         list_create(&ds->ds_sendstreams, sizeof (dmu_sendarg_t),
379             offsetof(dmu_sendarg_t, dsa_link));

381         if (err == 0) {
382             err = dsl_dir_hold_obj(dp,
383                 ds->ds_phys->ds_dir_obj, NULL, ds, &ds->ds_dir);
384         }
385         if (err != 0) {
386             mutex_destroy(&ds->ds_lock);
387             mutex_destroy(&ds->ds_opening_lock);
388             refcount_destroy(&ds->ds_longholds);
389             bplist_destroy(&ds->ds_pending_deadlist);
390             dsl_deadlist_close(&ds->ds_deadlist);

```

```

391         kmem_free(ds, sizeof (dsl_dataset_t));
392         dmu_buf_rele(dbuf, tag);
393         return (err);
394     }

396     if (!dsl_dataset_is_snapshot(ds)) {
397         ds->ds_snapname[0] = '\0';
398         if (ds->ds_phys->ds_prev_snap_obj != 0) {
399             err = dsl_dataset_hold_obj(dp,
400                 ds->ds_phys->ds_prev_snap_obj,
401                 ds, &ds->ds_prev);
402         }
403     } else {
404         if (zfs_flags & ZFS_DEBUG_SNAPNAMES)
405             err = dsl_dataset_get_snapname(ds);
406         if (err == 0 && ds->ds_phys->ds_userrefs_obj != 0) {
407             err = zap_count(
408                 ds->ds_dir->dd_pool->dp_meta_objset,
409                 ds->ds_phys->ds_userrefs_obj,
410                 &ds->ds_userrefs);
411         }
412     }

414     if (err == 0 && !dsl_dataset_is_snapshot(ds)) {
415         err = dsl_prop_get_int_ds(ds,
416             zfs_prop_to_name(ZFS_PROP_REFRESERVATION),
417             &ds->ds_reserved);
418         if (err == 0) {
419             err = dsl_prop_get_int_ds(ds,
420                 zfs_prop_to_name(ZFS_PROP_REFQUOTA),
421                 &ds->ds_quota);
422         }
423     } else {
424         ds->ds_reserved = ds->ds_quota = 0;
425     }

427     if (err != 0 || (winner = dmu_buf_set_user_ie(dbuf, ds,
428         &ds->ds_phys, dsl_dataset_evict)) != NULL) {
429         bplist_destroy(&ds->ds_pending_deadlist);
430         dsl_deadlist_close(&ds->ds_deadlist);
431         if (ds->ds_prev)
432             dsl_dataset_rele(ds->ds_prev, ds);
433         dsl_dir_rele(ds->ds_dir, ds);
434         mutex_destroy(&ds->ds_lock);
435         mutex_destroy(&ds->ds_opening_lock);
436         refcount_destroy(&ds->ds_longholds);
437         kmem_free(ds, sizeof (dsl_dataset_t));
438         if (err != 0) {
439             dmu_buf_rele(dbuf, tag);
440             return (err);
441         }
442         ds = winner;
443     } else {
444         ds->ds_fsid_guid =
445             unique_insert(ds->ds_phys->ds_fsid_guid);
446     }

447     }
448     ASSERT3P(ds->ds_dbuf, ==, dbuf);
449     ASSERT3P(ds->ds_phys, ==, dbuf->db_data);
450     ASSERT(ds->ds_phys->ds_prev_snap_obj != 0 ||
451         spa_version(dp->dp_spa) < SPA_VERSION_ORIGIN ||
452         dp->dp_origin_snap == NULL || ds == dp->dp_origin_snap);
453     *dsp = ds;
454     return (0);
455 }
_____unchanged_portion_omitted_____

```

```
2939 /*
2940  * Return TRUE if 'earlier' is an earlier snapshot in 'later's timeline.
2941  * For example, they could both be snapshots of the same filesystem, and
2942  * 'earlier' is before 'later'. Or 'earlier' could be the origin of
2943  * 'later's filesystem. Or 'earlier' could be an older snapshot in the origin's
2944  * filesystem. Or 'earlier' could be the origin's origin.
2945  */
2946 boolean_t
2947 dsl_dataset_is_before(dsl_dataset_t *later, dsl_dataset_t *earlier)
2948 {
2949     dsl_pool_t *dp = later->ds_dir->dd_pool;
2950     int error;
2951     boolean_t ret;
2952
2953     ASSERT(dsl_pool_config_held(dp));
2954
2955     if (earlier->ds_phys->ds_creation_txg >=
2956         later->ds_phys->ds_creation_txg)
2957         return (B_FALSE);
2958
2959     if (later->ds_dir == earlier->ds_dir)
2960         return (B_TRUE);
2961     if (!dsl_dir_is_clone(later->ds_dir))
2962         return (B_FALSE);
2963
2964     if (later->ds_dir->dd_phys->dd_origin_obj == earlier->ds_object)
2965         return (B_TRUE);
2966     dsl_dataset_t *origin;
2967     error = dsl_dataset_hold_obj(dp,
2968         later->ds_dir->dd_phys->dd_origin_obj, FTAG, &origin);
2969     if (error != 0)
2970         return (B_FALSE);
2971     ret = dsl_dataset_is_before(origin, earlier);
2972     dsl_dataset_rele(origin, FTAG);
2973     return (ret);
2974 }
2975
2976
2977 void
2978 dsl_dataset_zapify(dsl_dataset_t *ds, dmu_tx_t *tx)
2979 {
2980     objset_t *mos = ds->ds_dir->dd_pool->dp_meta_objset;
2981     dmu_object_zapify(mos, ds->ds_object, DMU_OT_DSL_DATASET, tx);
2982 }
2983
2984 unchanged portion omitted
```

```

*****
25661 Tue Oct 1 14:04:33 2013
new/usr/src/uts/common/fs/zfs/dsl_destroy.c
4171 clean up spa_feature_*( ) interfaces
4172 implement extensible_dataset feature for use by other zpool features
Reviewed by: Max Grossman <max.grossman@delphix.com>
Reviewed by: Christopher Siden <christopher.siden@delphix.com>
Reviewed by: George Wilson <george.wilson@delphix.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22  * Copyright (c) 2005, 2010, Oracle and/or its affiliates. All rights reserved.
23  * Copyright (c) 2013 by Delphix. All rights reserved.
24  * Copyright (c) 2013 Steven Hartland. All rights reserved.
25  */

27 #include <sys/zfs_context.h>
28 #include <sys/dsl_userhold.h>
29 #include <sys/dsl_dataset.h>
30 #include <sys/dsl_synctask.h>
31 #include <sys/dmu_tx.h>
32 #include <sys/dsl_pool.h>
33 #include <sys/dsl_dir.h>
34 #include <sys/dmu_traverse.h>
35 #include <sys/dsl_scan.h>
36 #include <sys/dmu_objset.h>
37 #include <sys/zap.h>
38 #include <sys/zfeature.h>
39 #include <sys/zfs_ioctl.h>
40 #include <sys/dsl_deleg.h>
41 #include <sys/dmu_impl.h>

43 typedef struct dmu_snapshots_destroy_arg {
44     nvlist_t *dsda_snaps;
45     nvlist_t *dsda_successful_snaps;
46     boolean_t dsda_defer;
47     nvlist_t *dsda_errlist;
48 } dmu_snapshots_destroy_arg_t;
    unchanged portion omitted

232 void
233 dsl_destroy_snapshot_sync_impl(dsl_dataset_t *ds, boolean_t defer, dmu_tx_t *tx)
234 {
235     int err;
236     int after_branch_point = FALSE;
237     dsl_pool_t *dp = ds->ds_dir->dd_pool;
238     objset_t *mos = dp->dp_meta_objset;

```

```

239     dsl_dataset_t *ds_prev = NULL;
240     uint64_t obj;

242     ASSERT(RRW_WRITE_HELD(&dp->dp_config_rwlock));
243     ASSERT3U(ds->ds_phys->ds_bp.blk_birth, <=, tx->tx_txcg);
244     ASSERT(refcount_is_zero(&ds->ds_longholds));

246     if (defer &&
247         (ds->ds_userrefs > 0 || ds->ds_phys->ds_num_children > 1)) {
248         ASSERT(spa_version(dp->dp_spa) >= SPA_VERSION_USERREFS);
249         dmu_buf_will_dirty(ds->ds_dbuf, tx);
250         ds->ds_phys->ds_flags |= DS_FLAG_DEFER_DESTROY;
251         spa_history_log_internal_ds(ds, "defer_destroy", tx, "");
252         return;
253     }

255     ASSERT3U(ds->ds_phys->ds_num_children, <=, 1);

257     /* We need to log before removing it from the namespace. */
258     spa_history_log_internal_ds(ds, "destroy", tx, "");

260     dsl_scan_ds_destroyed(ds, tx);

262     obj = ds->ds_object;

264     if (ds->ds_phys->ds_prev_snap_obj != 0) {
265         ASSERT3P(ds->ds_prev, ==, NULL);
266         VERIFY0(dsl_dataset_hold_obj(dp,
267             ds->ds_phys->ds_prev_snap_obj, FTAG, &ds_prev));
268         after_branch_point =
269             (ds_prev->ds_phys->ds_next_snap_obj != obj);

271         dmu_buf_will_dirty(ds_prev->ds_dbuf, tx);
272         if (after_branch_point &&
273             ds_prev->ds_phys->ds_next_clones_obj != 0) {
274             dsl_dataset_remove_from_next_clones(ds_prev, obj, tx);
275             if (ds->ds_phys->ds_next_snap_obj != 0) {
276                 VERIFY0(zap_add_int(mos,
277                     ds_prev->ds_phys->ds_next_clones_obj,
278                     ds->ds_phys->ds_next_snap_obj, tx));
279             }
280         }
281         if (!after_branch_point) {
282             ds_prev->ds_phys->ds_next_snap_obj =
283                 ds->ds_phys->ds_next_snap_obj;
284         }
285     }

287     dsl_dataset_t *ds_next;
288     uint64_t old_unique;
289     uint64_t used = 0, comp = 0, uncomp = 0;

291     VERIFY0(dsl_dataset_hold_obj(dp,
292         ds->ds_phys->ds_next_snap_obj, FTAG, &ds_next));
293     ASSERT3U(ds_next->ds_phys->ds_prev_snap_obj, ==, obj);

295     old_unique = ds_next->ds_phys->ds_unique_bytes;

297     dmu_buf_will_dirty(ds_next->ds_dbuf, tx);
298     ds_next->ds_phys->ds_prev_snap_obj =
299         ds->ds_phys->ds_prev_snap_obj;
300     ds_next->ds_phys->ds_prev_snap_txg =
301         ds->ds_phys->ds_prev_snap_txg;
302     ASSERT3U(ds->ds_phys->ds_prev_snap_txg, ==,
303         ds_prev ? ds_prev->ds_phys->ds_creation_txg : 0);

```

```

305     if (ds_next->ds_deadlist.dl_oldfmt) {
306         process_old_deadlist(ds, ds_prev, ds_next,
307             after_branch_point, tx);
308     } else {
309         /* Adjust prev's unique space. */
310         if (ds_prev && !after_branch_point) {
311             dsl_deadlist_space_range(&ds_next->ds_deadlist,
312                 ds_prev->ds_phys->ds_prev_snap_tsg,
313                 ds->ds_phys->ds_prev_snap_tsg,
314                 &used, &comp, &uncomp);
315             ds_prev->ds_phys->ds_unique_bytes += used;
316         }
317
318         /* Adjust snapused. */
319         dsl_deadlist_space_range(&ds_next->ds_deadlist,
320             ds->ds_phys->ds_prev_snap_tsg, UINTE64_MAX,
321             &used, &comp, &uncomp);
322         dsl_dir_diduse_space(ds->ds_dir, DD_USED_SNAP,
323             -used, -comp, -uncomp, tx);
324
325         /* Move blocks to be freed to pool's free list. */
326         dsl_deadlist_move_bproj(&ds_next->ds_deadlist,
327             &dp->dp_free_bproj, ds->ds_phys->ds_prev_snap_tsg,
328             tx);
329         dsl_dir_diduse_space(tx->tx_pool->dp_free_dir,
330             DD_USED_HEAD, used, comp, uncomp, tx);
331
332         /* Merge our deadlist into next's and free it. */
333         dsl_deadlist_merge(&ds_next->ds_deadlist,
334             ds->ds_phys->ds_deadlist_obj, tx);
335     }
336     dsl_deadlist_close(&ds->ds_deadlist);
337     dsl_deadlist_free(mos, ds->ds_phys->ds_deadlist_obj, tx);
338     dmu_buf_will_dirty(ds->ds_dbuf, tx);
339     ds->ds_phys->ds_deadlist_obj = 0;
340
341     /* Collapse range in clone heads */
342     dsl_dataset_remove_clones_key(ds,
343         ds->ds_phys->ds_creation_tsg, tx);
344
345     if (dsl_dataset_is_snapshot(ds_next)) {
346         dsl_dataset_t *ds_nextnext;
347
348         /*
349          * Update next's unique to include blocks which
350          * were previously shared by only this snapshot
351          * and it. Those blocks will be born after the
352          * prev snap and before this snap, and will have
353          * died after the next snap and before the one
354          * after that (ie. be on the snap after next's
355          * deadlist).
356          */
357         VERIFY0(dsl_dataset_hold_obj(dp,
358             ds_next->ds_phys->ds_next_snap_obj, FTAG, &ds_nextnext));
359         dsl_deadlist_space_range(&ds_nextnext->ds_deadlist,
360             ds->ds_phys->ds_prev_snap_tsg,
361             ds->ds_phys->ds_creation_tsg,
362             &used, &comp, &uncomp);
363         ds_next->ds_phys->ds_unique_bytes += used;
364         dsl_dataset_rele(ds_nextnext, FTAG);
365         ASSERT3P(ds_next->ds_prev, ==, NULL);
366
367         /* Collapse range in this head. */
368         dsl_dataset_t *hds;
369         VERIFY0(dsl_dataset_hold_obj(dp,
370             ds->ds_dir->dd_phys->dd_head_dataset_obj, FTAG, &hds));

```

```

371         dsl_deadlist_remove_key(&hds->ds_deadlist,
372             ds->ds_phys->ds_creation_tsg, tx);
373         dsl_dataset_rele(hds, FTAG);
374
375     } else {
376         ASSERT3P(ds_next->ds_prev, ==, ds);
377         dsl_dataset_rele(ds_next->ds_prev, ds_next);
378         ds_next->ds_prev = NULL;
379         if (ds_prev) {
380             VERIFY0(dsl_dataset_hold_obj(dp,
381                 ds->ds_phys->ds_prev_snap_obj,
382                 ds_next, &ds_next->ds_prev));
383         }
384
385         dsl_dataset_recalc_head_uniq(ds_next);
386
387         /*
388          * Reduce the amount of our unconsumed reservation
389          * being charged to our parent by the amount of
390          * new unique data we have gained.
391          */
392         if (old_unique < ds_next->ds_reserved) {
393             int64_t mrsdelta;
394             uint64_t new_unique =
395                 ds_next->ds_phys->ds_unique_bytes;
396
397             ASSERT(old_unique <= new_unique);
398             mrsdelta = MIN(new_unique - old_unique,
399                 ds_next->ds_reserved - old_unique);
400             dsl_dir_diduse_space(ds->ds_dir,
401                 DD_USED_REFRESRV, -mrsdelta, 0, 0, tx);
402         }
403     }
404     dsl_dataset_rele(ds_next, FTAG);
405
406     /*
407      * This must be done after the dsl_traverse(), because it will
408      * re-open the objset.
409      */
410     if (ds->ds_objset) {
411         dmu_objset_evict(ds->ds_objset);
412         ds->ds_objset = NULL;
413     }
414
415     /* remove from snapshot namespace */
416     dsl_dataset_t *ds_head;
417     ASSERT(ds->ds_phys->ds_snapnames_zapobj == 0);
418     VERIFY0(dsl_dataset_hold_obj(dp,
419         ds->ds_dir->dd_phys->dd_head_dataset_obj, FTAG, &ds_head));
420     VERIFY0(dsl_dataset_get_snapname(ds));
421 #ifdef ZFS_DEBUG
422     {
423         uint64_t val;
424
425         err = dsl_dataset_snap_lookup(ds_head,
426             ds->ds_snapname, &val);
427         ASSERT0(err);
428         ASSERT3U(val, ==, obj);
429     }
430 #endif
431     VERIFY0(dsl_dataset_snap_remove(ds_head, ds->ds_snapname, tx));
432     dsl_dataset_rele(ds_head, FTAG);
433
434     if (ds_prev != NULL)
435         dsl_dataset_rele(ds_prev, FTAG);

```

```

437     spa_prop_clear_bootfs(dp->dp_spa, ds->ds_object, tx);
439     if (ds->ds_phys->ds_next_clones_obj != 0) {
440         uint64_t count;
441         ASSERT0(zap_count(mos,
442             ds->ds_phys->ds_next_clones_obj, &count) && count == 0);
443         VERIFY0(dmu_object_free(mos,
444             ds->ds_phys->ds_next_clones_obj, tx));
445     }
446     if (ds->ds_phys->ds_props_obj != 0)
447         VERIFY0(zap_destroy(mos, ds->ds_phys->ds_props_obj, tx));
448     if (ds->ds_phys->ds_userrefs_obj != 0)
449         VERIFY0(zap_destroy(mos, ds->ds_phys->ds_userrefs_obj, tx));
450     dsl_dir_rele(ds->ds_dir, ds);
451     ds->ds_dir = NULL;
452     dmu_object_free_zapified(mos, obj, tx);
453     VERIFY0(dmu_object_free(mos, obj, tx));
454 }
455 unchanged_portion_omitted
456
463 static void
464 dsl_dir_destroy_sync(uint64_t ddoobj, dmu_tx_t *tx)
465 {
466     dsl_dir_t *dd;
467     dsl_pool_t *dp = dmu_tx_pool(tx);
468     objset_t *mos = dp->dp_meta_objset;
469     dd_used_t t;
470
471     ASSERT(RRW_WRITE_HELD(&dmu_tx_pool(tx)->dp_config_rwlock));
472
473     VERIFY0(dsl_dir_hold_obj(dp, ddoobj, NULL, FTAG, &dd));
474
475     ASSERT0(dd->dd_phys->dd_head_dataset_obj);
476
477     /*
478      * Remove our reservation. The impl() routine avoids setting the
479      * actual property, which would require the (already destroyed) ds.
480      */
481     dsl_dir_set_reservation_sync_impl(dd, 0, tx);
482
483     ASSERT0(dd->dd_phys->dd_used_bytes);
484     ASSERT0(dd->dd_phys->dd_reserved);
485     for (t = 0; t < DD_USED_NUM; t++)
486         ASSERT0(dd->dd_phys->dd_used_breakdown[t]);
487
488     VERIFY0(zap_destroy(mos, dd->dd_phys->dd_child_dir_zapobj, tx));
489     VERIFY0(zap_destroy(mos, dd->dd_phys->dd_props_zapobj, tx));
490     VERIFY0(dsl_deleg_destroy(mos, dd->dd_phys->dd_deleg_zapobj, tx));
491     VERIFY0(zap_remove(mos,
492         dd->dd_parent->dd_phys->dd_child_dir_zapobj, dd->dd_myname, tx));
493
494     dsl_dir_rele(dd, FTAG);
495     dmu_object_free_zapified(mos, ddoobj, tx);
496     VERIFY0(dmu_object_free(mos, ddoobj, tx));
497 }
498
499 void
500 dsl_destroy_head_sync_impl(dsl_dataset_t *ds, dmu_tx_t *tx)
501 {
502     dsl_pool_t *dp = dmu_tx_pool(tx);
503     objset_t *mos = dp->dp_meta_objset;
504     uint64_t obj, ddoobj, prevobj = 0;
505     boolean_t rmorigin;
506
507     ASSERT3U(ds->ds_phys->ds_num_children, <=, 1);
508     ASSERT(ds->ds_prev == NULL ||

```

```

688         ds->ds_prev->ds_phys->ds_next_snap_obj != ds->ds_object);
689     ASSERT3U(ds->ds_phys->ds_bp.blk_birth, <=, tx->tx_txg);
690     ASSERT(RRW_WRITE_HELD(&dp->dp_config_rwlock));
691
692     /* We need to log before removing it from the namespace. */
693     spa_history_log_internal(ds, "destroy", tx, "");
694
695     rmorigin = (dsl_dir_is_clone(ds->ds_dir) &&
696         DS_IS_DEFER_DESTROY(ds->ds_prev) &&
697         ds->ds_prev->ds_phys->ds_num_children == 2 &&
698         ds->ds_prev->ds_userrefs == 0);
699
700     /* Remove our reservation */
701     if (ds->ds_reserved != 0) {
702         dsl_dataset_set_refreservation_sync_impl(ds,
703             (ZPROP_SRC_NONE | ZPROP_SRC_LOCAL | ZPROP_SRC_RECEIVED),
704             0, tx);
705         ASSERT0(ds->ds_reserved);
706     }
707
708     dsl_scan_ds_destroyed(ds, tx);
709
710     obj = ds->ds_object;
711
712     if (ds->ds_phys->ds_prev_snap_obj != 0) {
713         /* This is a clone */
714         ASSERT(ds->ds_prev != NULL);
715         ASSERT3U(ds->ds_prev->ds_phys->ds_next_snap_obj, !=, obj);
716         ASSERT0(ds->ds_phys->ds_next_snap_obj);
717
718         dmu_buf_will_dirty(ds->ds_prev->ds_dbuf, tx);
719         if (ds->ds_prev->ds_phys->ds_next_clones_obj != 0) {
720             dsl_dataset_remove_from_next_clones(ds->ds_prev,
721                 obj, tx);
722         }
723
724         ASSERT3U(ds->ds_prev->ds_phys->ds_num_children, >, 1);
725         ds->ds_prev->ds_phys->ds_num_children--;
726     }
727
728     zfeature_info_t *async_destroy =
729     &spa_feature_table[SPA_FEATURE_ASYNC_DESTROY];
730     objset_t *os;
731
732     /*
733      * Destroy the deadlist. Unless it's a clone, the
734      * deadlist should be empty. (If it's a clone, it's
735      * safe to ignore the deadlist contents.)
736      */
737     dsl_deadlist_close(&ds->ds_deadlist);
738     dsl_deadlist_free(mos, ds->ds_phys->ds_deadlist_obj, tx);
739     dmu_buf_will_dirty(ds->ds_dbuf, tx);
740     ds->ds_phys->ds_deadlist_obj = 0;
741
742     objset_t *os;
743     VERIFY0(dmu_objset_from_ds(ds, &os));
744
745     if (!spa_feature_is_enabled(dp->dp_spa, SPA_FEATURE_ASYNC_DESTROY)) {
746     if (!spa_feature_is_enabled(dp->dp_spa, async_destroy)) {
747         old_synchronous_dataset_destroy(ds, tx);
748     } else {
749         /*
750          * Move the bptree into the pool's list of trees to
751          * clean up and update space accounting information.
752          */
753         uint64_t used, comp, uncomp;

```



```

750 zil_destroy_sync(dmu_objset_zil(os), tx);

752 if (!spa_feature_is_active(dp->dp_spa,
753 SPA_FEATURE_ASYNC_DESTROY)) {
754 if (!spa_feature_is_active(dp->dp_spa, async_destroy)) {
755     dsl_scan_t *scn = dp->dp_scan;
756     spa_feature_incr(dp->dp_spa, SPA_FEATURE_ASYNC_DESTROY,
757 tx);

757     spa_feature_incr(dp->dp_spa, async_destroy, tx);
757     dp->dp_bptree_obj = bptree_alloc(mos, tx);
758     VERIFY0(zap_add(mos,
759 DMU_POOL_DIRECTORY_OBJECT,
760 DMU_POOL_BPTREE_OBJ, sizeof (uint64_t), 1,
761 &dp->dp_bptree_obj, tx));
762     ASSERT(!scn->scn_async_destroying);
763     scn->scn_async_destroying = B_TRUE;
764 }

766 used = ds->ds_dir->dd_phys->dd_used_bytes;
767 comp = ds->ds_dir->dd_phys->dd_compressed_bytes;
768 uncomp = ds->ds_dir->dd_phys->dd_uncompressed_bytes;

770 ASSERT(!DS_UNIQUE_IS_ACCURATE(ds) ||
771 ds->ds_phys->ds_unique_bytes == used);

773 bptree_add(mos, dp->dp_bptree_obj,
774 &ds->ds_phys->ds_bp, ds->ds_phys->ds_prev_snap_txg,
775 used, comp, uncomp, tx);
776 dsl_dir_diduse_space(ds->ds_dir, DD_USED_HEAD,
777 -used, -comp, -uncomp, tx);
778 dsl_dir_diduse_space(dp->dp_free_dir, DD_USED_HEAD,
779 used, comp, uncomp, tx);
780 }

782 if (ds->ds_prev != NULL) {
783     if (spa_version(dp->dp_spa) >= SPA_VERSION_DIR_CLONES) {
784         VERIFY0(zap_remove_int(mos,
785 ds->ds_prev->ds_dir->dd_phys->dd_clones,
786 ds->ds_object, tx));
787     }
788     prevobj = ds->ds_prev->ds_object;
789     dsl_dataset_rele(ds->ds_prev, ds);
790     ds->ds_prev = NULL;
791 }

793 /*
794  * This must be done after the dsl_traverse(), because it will
795  * re-open the objset.
796  */
797 if (ds->ds_objset) {
798     dmu_objset_evict(ds->ds_objset);
799     ds->ds_objset = NULL;
800 }

802 /* Erase the link in the dir */
803 dmu_buf_will_dirty(ds->ds_dir->dd_dbuf, tx);
804 ds->ds_dir->dd_phys->dd_head_dataset_obj = 0;
805 ddoobj = ds->ds_dir->dd_object;
806 ASSERT(ds->ds_phys->ds_snapnames_zapobj != 0);
807 VERIFY0(zap_destroy(mos, ds->ds_phys->ds_snapnames_zapobj, tx));

809 spa_prop_clear_bootfs(dp->dp_spa, ds->ds_object, tx);

811 ASSERT0(ds->ds_phys->ds_next_clones_obj);

```

```

812 ASSERT0(ds->ds_phys->ds_props_obj);
813 ASSERT0(ds->ds_phys->ds_userrefs_obj);
814 dsl_dir_rele(ds->ds_dir, ds);
815 ds->ds_dir = NULL;
816 dmu_object_free_zapified(mos, obj, tx);
817 VERIFY0(dmu_object_free(mos, obj, tx));

818 dsl_dir_destroy_sync(ddobj, tx);

820 if (rmorigin) {
821     dsl_dataset_t *prev;
822     VERIFY0(dsl_dataset_hold_obj(dp, prevobj, FTAG, &prev));
823     dsl_destroy_snapshot_sync_impl(prev, B_FALSE, tx);
824     dsl_dataset_rele(prev, FTAG);
825 }
826 }
unchanged_portion_omitted

857 int
858 dsl_destroy_head(const char *name)
859 {
860     dsl_destroy_head_arg_t ddha;
861     int error;
862     spa_t *spa;
863     boolean_t isenabled;

865 #ifdef _KERNEL
866     zfs_destroy_unmount_origin(name);
867 #endif

869 error = spa_open(name, &spa, FTAG);
870 if (error != 0)
871     return (error);
872 isenabled = spa_feature_is_enabled(spa, SPA_FEATURE_ASYNC_DESTROY);
873 isenabled = spa_feature_is_enabled(spa,
874 &spa_feature_table[SPA_FEATURE_ASYNC_DESTROY]);
875 spa_close(spa, FTAG);

877 ddha.ddha_name = name;

877 if (!isenabled) {
878     objset_t *os;

880     error = dsl_sync_task(name, dsl_destroy_head_check,
881 dsl_destroy_head_begin_sync, &ddha, 0);
882     if (error != 0)
883         return (error);

885     /*
886      * Head deletion is processed in one txg on old pools;
887      * remove the objects from open context so that the txg sync
888      * is not too long.
889      */
890     error = dmu_objset_own(name, DMU_OST_ANY, B_FALSE, FTAG, &os);
891     if (error == 0) {
892         uint64_t prev_snap_txg =
893             dmu_objset_ds(os)->ds_phys->ds_prev_snap_txg;
894         for (uint64_t obj = 0; error == 0;
895             error = dmu_object_next(os, &obj, FALSE,
896 prev_snap_txg))
897             (void) dmu_free_long_object(os, obj);
898         /* sync out all frees */
899         txg_wait_synced(dmu_objset_pool(os), 0);
900         dmu_objset_disown(os, FTAG);
901     }
902 }

```

new/usr/src/uts/common/fs/zfs/dsl\_destroy.c

9

```
904     return (dsl_sync_task(name, dsl_destroy_head_check,  
905         dsl_destroy_head_sync, &ddha, 0));  
906 }  
_____unchanged_portion_omitted_____
```

```

*****
35887 Tue Oct 1 14:04:35 2013
new/usr/src/uts/common/fs/zfs/dsl_dir.c
4171 clean up spa_feature_*() interfaces
4172 implement extensible_dataset feature for use by other zpool features
Reviewed by: Max Grossman <max.grossman@delphix.com>
Reviewed by: Christopher Siden <christopher.siden@delphix.com>
Reviewed by: George Wilson <george.wilson@delphix.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22  * Copyright (c) 2005, 2010, Oracle and/or its affiliates. All rights reserved.
23  * Copyright (c) 2013 by Delphix. All rights reserved.
24  * Copyright (c) 2013 Martin Matuska. All rights reserved.
25  */

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

44 static uint64_t dsl_dir_space_towrite(dsl_dir_t *dd);

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

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

```

```

58     }
59
60     if (dd->dd_parent)
61         dsl_dir_rele(dd->dd_parent, dd);
62
63     spa_close(dd->dd_pool->dp_spa, dd);
64
65     /*
66      * The props callback list should have been cleaned up by
67      * objset_evict().
68      */
69     list_destroy(&dd->dd_prop_cbs);
70     mutex_destroy(&dd->dd_lock);
71     kmem_free(dd, sizeof (dsl_dir_t));
72 }

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

82     ASSERT(dsl_pool_config_held(dp));

84     err = dmu_bonus_hold(dp->dp_meta_objset, ddojb, tag, &dbuf);
85     if (err != 0)
86         return (err);
87     dd = dmu_buf_get_user(dbuf);
88 #ifdef ZFS_DEBUG
89     {
90         dmu_object_info_t doi;
91         dmu_object_info_from_db(dbuf, &doi);
92         ASSERT3U(doi.doi_bonus_type, ==, DMU_OT_DSL_DIR);
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         dd = kmem_zalloc(sizeof (dsl_dir_t), KM_SLEEP);
99         dd->dd_object = ddojb;
100         dd->dd_dbuf = dbuf;
101         dd->dd_pool = dp;
102         dd->dd_phys = dbuf->db_data;
103         mutex_init(&dd->dd_lock, NULL, MUTEX_DEFAULT, NULL);

106         list_create(&dd->dd_prop_cbs, sizeof (dsl_dir_phys_t),
107                   offsetof(dsl_dir_phys_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                 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     winner = dmu_buf_set_user_ie(dbuf, dd, &dd->dd_phys,
157         dsl_dir_evict);
158     if (winner) {
159         if (dd->dd_parent)
160             dsl_dir_rele(dd->dd_parent, dd);
161         mutex_destroy(&dd->dd_lock);
162         kmem_free(dd, sizeof (dsl_dir_t));
163         dd = winner;
164     } else {
165         spa_open_ref(dp->dp_spa, dd);
166     }
167
168     /*
169     * The dsl_dir_t has both open-to-close and instantiate-to-evict
170     * holds on the spa. We need the open-to-close holds because
171     * otherwise the spa_refcnt wouldn't change when we open a
172     * dir which the spa also has open, so we could incorrectly
173     * think it was OK to unload/export/destroy the pool. We need
174     * the instantiate-to-evict hold because the dsl_dir_t has a
175     * pointer to the dd_pool, which has a pointer to the spa_t.
176     */
177     spa_open_ref(dp->dp_spa, tag);
178     ASSERT3P(dd->dd_pool, ==, dp);
179     ASSERT3U(dd->dd_object, ==, ddoobj);
180     ASSERT3P(dd->dd_dbuf, ==, dbuf);
181     *ddp = dd;
182     return (0);
183
184 errout:
185     if (dd->dd_parent)
186         dsl_dir_rele(dd->dd_parent, dd);
187     mutex_destroy(&dd->dd_lock);
188

```

```

189         kmem_free(dd, sizeof (dsl_dir_t));
190         dmu_buf_rele(dbuf, tag);
191         return (err);
192     }
193     _____ unchanged_portion_omitted _____
194
195 void
196 dsl_dir_snap_cmtime_update(dsl_dir_t *dd)
197 {
198     timestruc_t t;
199
200     gethrestime(&t);
201     mutex_enter(&dd->dd_lock);
202     dd->dd_snap_cmtime = t;
203     mutex_exit(&dd->dd_lock);
204 }
205
206 void
207 dsl_dir_zapify(dsl_dir_t *dd, dmu_tx_t *tx)
208 {
209     objset_t *mos = dd->dd_pool->dp_meta_objset;
210     dmu_object_zapify(mos, dd->dd_object, DMU_OT_DSL_DIR, tx);
211 }
212     _____ unchanged_portion_omitted _____

```

```

*****
30598 Tue Oct 1 14:04:35 2013
new/usr/src/uts/common/fs/zfs/dsl_pool.c
4171 clean up spa_feature_*() interfaces
4172 implement extensible_dataset feature for use by other zpool features
Reviewed by: Max Grossman <max.grossman@delphix.com>
Reviewed by: Christopher Siden <christopher.siden@delphix.com>
Reviewed by: George Wilson <george.wilson@delphix.com>
*****
_____unchanged_portion_omitted_____

198 int
199 dsl_pool_open(dsl_pool_t *dp)
200 {
201     int err;
202     dsl_dir_t *dd;
203     dsl_dataset_t *ds;
204     uint64_t obj;

206     rrw_enter(&dp->dp_config_rwlock, RW_WRITER, FTAG);
207     err = zap_lookup(dp->dp_meta_objset, DMU_POOL_DIRECTORY_OBJECT,
208         DMU_POOL_ROOT_DATASET, sizeof (uint64_t), 1,
209         &dp->dp_root_dir_obj);
210     if (err)
211         goto out;

213     err = dsl_dir_hold_obj(dp, dp->dp_root_dir_obj,
214         NULL, dp, &dp->dp_root_dir);
215     if (err)
216         goto out;

218     err = dsl_pool_open_special_dir(dp, MOS_DIR_NAME, &dp->dp_mos_dir);
219     if (err)
220         goto out;

222     if (spa_version(dp->dp_spa) >= SPA_VERSION_ORIGIN) {
223         err = dsl_pool_open_special_dir(dp, ORIGIN_DIR_NAME, &dd);
224         if (err)
225             goto out;
226         err = dsl_dataset_hold_obj(dp, dd->dd_phys->dd_head_dataset_obj,
227             FTAG, &ds);
228         if (err == 0) {
229             err = dsl_dataset_hold_obj(dp,
230                 ds->ds_phys->ds_prev_snap_obj, dp,
231                 &dp->dp_origin_snap);
232             dsl_dataset_rele(ds, FTAG);
233         }
234         dsl_dir_rele(dd, dp);
235         if (err)
236             goto out;
237     }

239     if (spa_version(dp->dp_spa) >= SPA_VERSION_DEADLISTS) {
240         err = dsl_pool_open_special_dir(dp, FREE_DIR_NAME,
241             &dp->dp_free_dir);
242         if (err)
243             goto out;

245         err = zap_lookup(dp->dp_meta_objset, DMU_POOL_DIRECTORY_OBJECT,
246             DMU_POOL_FREE_BPOBJ, sizeof (uint64_t), 1, &obj);
247         if (err)
248             goto out;
249         VERIFY0(bpobj_open(&dp->dp_free_bpobj,
250             dp->dp_meta_objset, obj));
251     }

```

```

253     if (spa_feature_is_active(dp->dp_spa, SPA_FEATURE_ASYNC_DESTROY)) {
254         if (spa_feature_is_active(dp->dp_spa,
255             &spa_feature_table[SPA_FEATURE_ASYNC_DESTROY])) {
256             err = zap_lookup(dp->dp_meta_objset, DMU_POOL_DIRECTORY_OBJECT,
257                 DMU_POOL_BPTREE_OBJ, sizeof (uint64_t), 1,
258                 &dp->dp_bptree_obj);
259             if (err != 0)
260                 goto out;
261         }
262     }

261     if (spa_feature_is_active(dp->dp_spa, SPA_FEATURE_EMPTY_BPOBJ)) {
262         if (spa_feature_is_active(dp->dp_spa,
263             &spa_feature_table[SPA_FEATURE_EMPTY_BPOBJ])) {
264             err = zap_lookup(dp->dp_meta_objset, DMU_POOL_DIRECTORY_OBJECT,
265                 DMU_POOL_EMPTY_BPOBJ, sizeof (uint64_t), 1,
266                 &dp->dp_empty_bpobj);
267             if (err != 0)
268                 goto out;
269         }
270     }

269     err = zap_lookup(dp->dp_meta_objset, DMU_POOL_DIRECTORY_OBJECT,
270         DMU_POOL_TMP_USERREFS, sizeof (uint64_t), 1,
271         &dp->dp_tmp_userrefs_obj);
272     if (err == ENOENT)
273         err = 0;
274     if (err)
275         goto out;

277     err = dsl_scan_init(dp, dp->dp_tx.tx_open_txg);

279 out:
280     rrw_exit(&dp->dp_config_rwlock, FTAG);
281     return (err);
282 }
_____unchanged_portion_omitted_____

```

```

*****
50297 Tue Oct 1 14:04:37 2013
new/usr/src/uts/common/fs/zfs/dsl_scan.c
4171 clean up spa_feature_*( ) interfaces
4172 implement extensible_dataset feature for use by other zpool features
Reviewed by: Max Grossman <max.grossman@delphix.com>
Reviewed by: Christopher Siden <christopher.siden@delphix.com>
Reviewed by: George Wilson <george.wilson@delphix.com>
*****
_unchanged_portion_omitted_

87 int
88 dsl_scan_init(dsl_pool_t *dp, uint64_t txg)
89 {
90     int err;
91     dsl_scan_t *scn;
92     spa_t *spa = dp->dp_spa;
93     uint64_t f;

95     scn = dp->dp_scan = kmem_zalloc(sizeof (dsl_scan_t), KM_SLEEP);
96     scn->scn_dp = dp;

98     /*
99     * It's possible that we're resuming a scan after a reboot so
100    * make sure that the scan_async_destroying flag is initialized
101    * appropriately.
102    */
103    ASSERT(!scn->scn_async_destroying);
104    scn->scn_async_destroying = spa_feature_is_active(dp->dp_spa,
105    SPA_FEATURE_ASYNC_DESTROY);
105    &spa_feature_table[SPA_FEATURE_ASYNC_DESTROY]);

107    err = zap_lookup(dp->dp_meta_objset, DMU_POOL_DIRECTORY_OBJECT,
108    "scrub_func", sizeof (uint64_t), 1, &f);
109    if (err == 0) {
110        /*
111        * There was an old-style scrub in progress. Restart a
112        * new-style scrub from the beginning.
113        */
114        scn->scn_restart_txg = txg;
115        zfs_dbgmsg("old-style scrub was in progress; "
116        "restarting new-style scrub in txg %llu",
117        scn->scn_restart_txg);

119        /*
120        * Load the queue obj from the old location so that it
121        * can be freed by dsl_scan_done().
122        */
123        (void) zap_lookup(dp->dp_meta_objset, DMU_POOL_DIRECTORY_OBJECT,
124        "scrub_queue", sizeof (uint64_t), 1,
125        &scn->scn_phys.scn_queue_obj);
126    } else {
127        err = zap_lookup(dp->dp_meta_objset, DMU_POOL_DIRECTORY_OBJECT,
128        DMU_POOL_SCAN, sizeof (uint64_t), SCAN_PHYS_NUMINTS,
129        &scn->scn_phys);
130        if (err == ENOENT)
131            return (0);
132        else if (err)
133            return (err);

135        if (scn->scn_phys.scn_state == DSS_SCANNING &&
136        spa_prev_software_version(dp->dp_spa) < SPA_VERSION_SCAN) {
137            /*
138            * A new-type scrub was in progress on an old
139            * pool, and the pool was accessed by old
140            * software. Restart from the beginning, since

```

```

141         * the old software may have changed the pool in
142         * the meantime.
143         */
144         scn->scn_restart_txg = txg;
145         zfs_dbgmsg("new-style scrub was modified "
146         "by old software; restarting in txg %llu",
147         scn->scn_restart_txg);
148     }
149 }

151     spa_scan_stat_init(spa);
152     return (0);
153 }
_unchanged_portion_omitted_

1346 boolean_t
1347 dsl_scan_active(dsl_scan_t *scn)
1348 {
1349     spa_t *spa = scn->scn_dp->dp_spa;
1350     uint64_t used = 0, comp, uncomp;

1352     if (spa->spa_load_state != SPA_LOAD_NONE)
1353         return (B_FALSE);
1354     if (spa_shutting_down(spa))
1355         return (B_FALSE);

1356     if (scn->scn_phys.scn_state == DSS_SCANNING ||
1357         scn->scn_async_destroying)
1358         return (B_TRUE);

1360     if (spa_version(scn->scn_dp->dp_spa) >= SPA_VERSION_DEADLISTS) {
1361         (void) bpobj_space(&scn->scn_dp->dp_free_bpobj,
1362         &used, &comp, &uncomp);
1363     }
1364     return (used != 0);
1365 }

1367 void
1368 dsl_scan_sync(dsl_pool_t *dp, dmu_tx_t *tx)
1369 {
1370     dsl_scan_t *scn = dp->dp_scan;
1371     spa_t *spa = dp->dp_spa;
1372     int err;

1374     /*
1375     * Check for scn_restart_txg before checking spa_load_state, so
1376     * that we can restart an old-style scan while the pool is being
1377     * imported (see dsl_scan_init).
1378     */
1379     if (scn->scn_restart_txg != 0 &&
1380         scn->scn_restart_txg <= tx->tx_txg) {
1381         pool_scan_func_t func = POOL_SCAN_SCRUB;
1382         dsl_scan_done(scn, B_FALSE, tx);
1383         if (vdev_resilver_needed(spa->spa_root_vdev, NULL, NULL))
1384             func = POOL_SCAN_RESILVER;
1385         zfs_dbgmsg("restarting scan func=%u txg=%llu",
1386         func, tx->tx_txg);
1387         dsl_scan_setup_sync(&func, tx);
1388     }

1390     if (!dsl_scan_active(scn) ||
1391         spa_sync_pass(dp->dp_spa) > 1)
1392         return;

1394     scn->scn_visited_this_txg = 0;
1395     scn->scn_pausing = B_FALSE;

```

```

1396     scn->scn_sync_start_time = gethrtime();
1397     spa->spa_scrub_active = B_TRUE;

1399     /*
1400     * First process the free list.  If we pause the free, don't do
1401     * any scanning.  This ensures that there is no free list when
1402     * we are scanning, so the scan code doesn't have to worry about
1403     * traversing it.
1404     */
1405     if (spa_version(dp->dp_spa) >= SPA_VERSION_DEADLISTS) {
1406         scn->scn_is_bptree = B_FALSE;
1407         scn->scn_zio_root = zio_root(dp->dp_spa, NULL,
1408             NULL, ZIO_FLAG_MUSTSUCCEED);
1409         err = bpobj_iterate(&dp->dp_free_bpobj,
1410             dsl_scan_free_block_cb, scn, tx);
1411         VERIFY3U(0, ==, zio_wait(scn->scn_zio_root));

1413         if (err == 0 && spa_feature_is_active(spa,
1414             SPA_FEATURE_ASYNC_DESTROY)) {
1415             &spa_feature_table[SPA_FEATURE_ASYNC_DESTROY])) {
1416                 ASSERT(scn->scn_async_destroying);
1417                 scn->scn_is_bptree = B_TRUE;
1418                 scn->scn_zio_root = zio_root(dp->dp_spa, NULL,
1419                     NULL, ZIO_FLAG_MUSTSUCCEED);
1420                 err = bptree_iterate(dp->dp_meta_objset,
1421                     dp->dp_bptree_obj, B_TRUE, dsl_scan_free_block_cb,
1422                     scn, tx);
1423                 VERIFY0(zio_wait(scn->scn_zio_root));

1424                 if (err == 0) {
1425                     zfeature_info_t *feat = &spa_feature_table
1426                         [SPA_FEATURE_ASYNC_DESTROY];
1427                     /* finished; deactivate async destroy feature */
1428                     spa_feature_decr(spa, SPA_FEATURE_ASYNC_DESTROY,
1429                         tx);
1430                     ASSERT(!spa_feature_is_active(spa, feat));
1431                     VERIFY0(zap_remove(dp->dp_meta_objset,
1432                         DMU_POOL_DIRECTORY_OBJECT,
1433                         DMU_POOL_BPTREE_OBJ, tx));
1434                     VERIFY0(bptree_free(dp->dp_meta_objset,
1435                         dp->dp_bptree_obj, tx));
1436                     dp->dp_bptree_obj = 0;
1437                     scn->scn_async_destroying = B_FALSE;
1438                 }
1439             }
1440             if (scn->scn_visited_this_txg) {
1441                 zfs_dbgmsg("freed %llu blocks in %llums from "
1442                     "free_bpobj/bptree txg %llu",
1443                     (longlong_t)scn->scn_visited_this_txg,
1444                     (longlong_t)
1445                         NSEC2MSEC(gethrtime() - scn->scn_sync_start_time),
1446                     (longlong_t)tx->tx_txg);
1447                 scn->scn_visited_this_txg = 0;
1448                 /*
1449                 * Re-sync the ddt so that we can further modify
1450                 * it when doing bprewrite.
1451                 */
1452                 ddt_sync(spa, tx->tx_txg);
1453             }
1454             if (err == ERESTART)
1455                 return;
1456         }

```

```

1457         if (scn->scn_phys.scn_state != DSS_SCANNING)
1458             return;

1460         if (scn->scn_done_txg == tx->tx_txg) {
1461             ASSERT(!scn->scn_pausing);
1462             /* finished with scan. */
1463             zfs_dbgmsg("txg %llu scan complete", tx->tx_txg);
1464             dsl_scan_done(scn, B_TRUE, tx);
1465             ASSERT3U(spa->spa_scrub_inflight, ==, 0);
1466             dsl_scan_sync_state(scn, tx);
1467             return;
1468         }

1470         if (scn->scn_phys.scn_ddt_bookmark.ddb_class <=
1471             scn->scn_phys.scn_ddt_class_max) {
1472             zfs_dbgmsg("doing scan sync txg %llu; "
1473                 "ddt bm=%llu/%llu/%llu/%llx",
1474                 (longlong_t)tx->tx_txg,
1475                 (longlong_t)scn->scn_phys.scn_ddt_bookmark.ddb_class,
1476                 (longlong_t)scn->scn_phys.scn_ddt_bookmark.ddb_type,
1477                 (longlong_t)scn->scn_phys.scn_ddt_bookmark.ddb_checks,
1478                 (longlong_t)scn->scn_phys.scn_ddt_bookmark.ddb_cursor);
1479             ASSERT(scn->scn_phys.scn_bookmark.zb_objset == 0);
1480             ASSERT(scn->scn_phys.scn_bookmark.zb_object == 0);
1481             ASSERT(scn->scn_phys.scn_bookmark.zb_level == 0);
1482             ASSERT(scn->scn_phys.scn_bookmark.zb_blkid == 0);
1483         } else {
1484             zfs_dbgmsg("doing scan sync txg %llu; bm=%llu/%llu/%llu/%llu",
1485                 (longlong_t)tx->tx_txg,
1486                 (longlong_t)scn->scn_phys.scn_bookmark.zb_objset,
1487                 (longlong_t)scn->scn_phys.scn_bookmark.zb_object,
1488                 (longlong_t)scn->scn_phys.scn_bookmark.zb_level,
1489                 (longlong_t)scn->scn_phys.scn_bookmark.zb_blkid);
1490         }

1492         scn->scn_zio_root = zio_root(dp->dp_spa, NULL,
1493             NULL, ZIO_FLAG_CANFAIL);
1494         dsl_pool_config_enter(dp, FTAG);
1495         dsl_scan_visit(scn, tx);
1496         dsl_pool_config_exit(dp, FTAG);
1497         (void) zio_wait(scn->scn_zio_root);
1498         scn->scn_zio_root = NULL;

1500         zfs_dbgmsg("visited %llu blocks in %llums",
1501             (longlong_t)scn->scn_visited_this_txg,
1502             (longlong_t)NSEC2MSEC(gethrtime() - scn->scn_sync_start_time));

1504         if (!scn->scn_pausing) {
1505             scn->scn_done_txg = tx->tx_txg + 1;
1506             zfs_dbgmsg("txg %llu traversal complete, waiting till txg %llu",
1507                 tx->tx_txg, scn->scn_done_txg);
1508         }

1510         if (DSL_SCAN_IS_SCRUB_RESILVER(scn)) {
1511             mutex_enter(&spa->spa_scrub_lock);
1512             while (spa->spa_scrub_inflight > 0) {
1513                 cv_wait(&spa->spa_scrub_io_cv,
1514                     &spa->spa_scrub_lock);
1515             }
1516             mutex_exit(&spa->spa_scrub_lock);
1517         }

1519         dsl_scan_sync_state(scn, tx);
1520     }

```

unchanged portion omitted

```

*****
175806 Tue Oct 1 14:04:39 2013
new/usr/src/uts/common/fs/zfs/spa.c
4171 clean up spa_feature_*() interfaces
4172 implement extensible_dataset feature for use by other zpool features
Reviewed by: Max Grossman <max.grossman@delphix.com>
Reviewed by: Christopher Siden <christopher.siden@delphix.com>
Reviewed by: George Wilson <george.wilson@delphix.com>
*****
_____unchanged_portion_omitted_____

2070 /*
2071  * Load an existing storage pool, using the pool's builtin spa_config as a
2072  * source of configuration information.
2073  */
2074 static int
2075 spa_load_impl(spa_t *spa, uint64_t pool_guid, nvlist_t *config,
2076              spa_load_state_t state, spa_import_type_t type, boolean_t mosconfig,
2077              char **ereport)
2078 {
2079     int error = 0;
2080     nvlist_t *nvroot = NULL;
2081     nvlist_t *label;
2082     vdev_t *rvd;
2083     uberblock_t *ub = &spa->spa_uberblock;
2084     uint64_t children, config_cache_txg = spa->spa_config_txg;
2085     int orig_mode = spa->spa_mode;
2086     int parse;
2087     uint64_t obj;
2088     boolean_t missing_feat_write = B_FALSE;

2090     /*
2091      * If this is an untrusted config, access the pool in read-only mode.
2092      * This prevents things like resilvering recently removed devices.
2093      */
2094     if (!mosconfig)
2095         spa->spa_mode = FREAD;

2097     ASSERT(MUTEX_HELD(&spa_namespace_lock));

2099     spa->spa_load_state = state;

2101     if (nvlist_lookup_nvlist(config, ZPOOL_CONFIG_VDEV_TREE, &nvroot))
2102         return (SET_ERROR(EINVAL));

2104     parse = (type == SPA_IMPORT_EXISTING ?
2105             VDEV_ALLOC_LOAD : VDEV_ALLOC_SPLIT);

2107     /*
2108      * Create "The Godfather" zio to hold all async IOs
2109      */
2110     spa->spa_async_zio_root = zio_root(spa, NULL, NULL,
2111                                     ZIO_FLAG_CANFAIL | ZIO_FLAG_SPECULATIVE | ZIO_FLAG_GODFATHER);

2113     /*
2114      * Parse the configuration into a vdev tree. We explicitly set the
2115      * value that will be returned by spa_version() since parsing the
2116      * configuration requires knowing the version number.
2117      */
2118     spa_config_enter(spa, SCL_ALL, FTAG, RW_WRITER);
2119     error = spa_config_parse(spa, &rvd, nvroot, NULL, 0, parse);
2120     spa_config_exit(spa, SCL_ALL, FTAG);

2122     if (error != 0)
2123         return (error);

```

```

2125     ASSERT(spa->spa_root_vdev == rvd);

2127     if (type != SPA_IMPORT_ASSEMBLE) {
2128         ASSERT(spa_guid(spa) == pool_guid);
2129     }

2131     /*
2132      * Try to open all vdevs, loading each label in the process.
2133      */
2134     spa_config_enter(spa, SCL_ALL, FTAG, RW_WRITER);
2135     error = vdev_open(rvd);
2136     spa_config_exit(spa, SCL_ALL, FTAG);
2137     if (error != 0)
2138         return (error);

2140     /*
2141      * We need to validate the vdev labels against the configuration that
2142      * we have in hand, which is dependent on the setting of mosconfig. If
2143      * mosconfig is true then we're validating the vdev labels based on
2144      * that config. Otherwise, we're validating against the cached config
2145      * (zpool.cache) that was read when we loaded the zfs module, and then
2146      * later we will recursively call spa_load() and validate against
2147      * the vdev config.
2148      *
2149      * If we're assembling a new pool that's been split off from an
2150      * existing pool, the labels haven't yet been updated so we skip
2151      * validation for now.
2152      */
2153     if (type != SPA_IMPORT_ASSEMBLE) {
2154         spa_config_enter(spa, SCL_ALL, FTAG, RW_WRITER);
2155         error = vdev_validate(rvd, mosconfig);
2156         spa_config_exit(spa, SCL_ALL, FTAG);

2158         if (error != 0)
2159             return (error);

2161         if (rvd->vdev_state <= VDEV_STATE_CANT_OPEN)
2162             return (SET_ERROR(ENXIO));
2163     }

2165     /*
2166      * Find the best uberblock.
2167      */
2168     vdev_uberblock_load(rvd, ub, &label);

2170     /*
2171      * If we weren't able to find a single valid uberblock, return failure.
2172      */
2173     if (ub->ub_txg == 0) {
2174         nvlist_free(label);
2175         return (spa_vdev_err(rvd, VDEV_AUX_CORRUPT_DATA, ENXIO));
2176     }

2178     /*
2179      * If the pool has an unsupported version we can't open it.
2180      */
2181     if (!SPA_VERSION_IS_SUPPORTED(ub->ub_version)) {
2182         nvlist_free(label);
2183         return (spa_vdev_err(rvd, VDEV_AUX_VERSION_NEWER, ENOTSUP));
2184     }

2186     if (ub->ub_version >= SPA_VERSION_FEATURES) {
2187         nvlist_t *features;

2189         /*
2190          * If we weren't able to find what's necessary for reading the

```



```

2191         * MOS in the label, return failure.
2192         */
2193         if (label == NULL || nvlist_lookup_nvlist(label,
2194             ZPOOL_CONFIG_FEATURES_FOR_READ, &features) != 0) {
2195             nvlist_free(label);
2196             return (spa_vdev_err(rvd, VDEV_AUX_CORRUPT_DATA,
2197                 ENXIO));
2198         }
2199
2200         /*
2201          * Update our in-core representation with the definitive values
2202          * from the label.
2203          */
2204         nvlist_free(spa->spa_label_features);
2205         VERIFY(nvlist_dup(features, &spa->spa_label_features, 0) == 0);
2206     }
2207
2208     nvlist_free(label);
2209
2210     /*
2211      * Look through entries in the label nvlist's features_for_read. If
2212      * there is a feature listed there which we don't understand then we
2213      * cannot open a pool.
2214      */
2215     if (ub->ub_version >= SPA_VERSION_FEATURES) {
2216         nvlist_t *unsup_feat;
2217
2218         VERIFY(nvlist_alloc(&unsup_feat, NV_UNIQUE_NAME, KM_SLEEP) ==
2219             0);
2220
2221         for (nvpair_t *nvp = nvlist_next_nvpair(spa->spa_label_features,
2222             NULL); nvp != NULL;
2223             nvp = nvlist_next_nvpair(spa->spa_label_features, nvp)) {
2224             if (!zfeature_is_supported(nvpair_name(nvp))) {
2225                 VERIFY(nvlist_add_string(unsup_feat,
2226                     nvpair_name(nvp), "") == 0);
2227             }
2228         }
2229
2230         if (!nvlist_empty(unsup_feat)) {
2231             VERIFY(nvlist_add_nvlist(spa->spa_load_info,
2232                 ZPOOL_CONFIG_UNSUP_FEAT, unsup_feat) == 0);
2233             nvlist_free(unsup_feat);
2234             return (spa_vdev_err(rvd, VDEV_AUX_UNSUP_FEAT,
2235                 ENOTSUP));
2236         }
2237
2238         nvlist_free(unsup_feat);
2239     }
2240
2241     /*
2242      * If the vdev guid sum doesn't match the uberblock, we have an
2243      * incomplete configuration. We first check to see if the pool
2244      * is aware of the complete config (i.e ZPOOL_CONFIG_VDEV_CHILDREN).
2245      * If it is, defer the vdev_guid_sum check till later so we
2246      * can handle missing vdevs.
2247      */
2248     if (nvlist_lookup_uint64(config, ZPOOL_CONFIG_VDEV_CHILDREN,
2249         &children) != 0 && mosconfig && type != SPA_IMPORT_ASSEMBLE &&
2250         rvd->vdev_guid_sum != ub->ub_guid_sum)
2251         return (spa_vdev_err(rvd, VDEV_AUX_BAD_GUID_SUM, ENXIO));
2252
2253     if (type != SPA_IMPORT_ASSEMBLE && spa->spa_config_splitting) {
2254         spa_config_enter(spa, SCL_ALL, FTAG, RW_WRITER);
2255         spa_try_repair(spa, config);
2256         spa_config_exit(spa, SCL_ALL, FTAG);

```

```

2257         nvlist_free(spa->spa_config_splitting);
2258         spa->spa_config_splitting = NULL;
2259     }
2260
2261     /*
2262      * Initialize internal SPA structures.
2263      */
2264     spa->spa_state = POOL_STATE_ACTIVE;
2265     spa->spa_ubsync = spa->spa_uberblock;
2266     spa->spa_verify_min_txg = spa->spa_extreme_rewind ?
2267         TXG_INITIAL - 1 : spa_last_synced_txg(spa) - TXG_DEFER_SIZE - 1;
2268     spa->spa_first_txg = spa->spa_last_ubsync_txg ?
2269         spa->spa_last_ubsync_txg : spa_last_synced_txg(spa) + 1;
2270     spa->spa_claim_max_txg = spa->spa_first_txg;
2271     spa->spa_prev_software_version = ub->ub_software_version;
2272
2273     error = dsl_pool_init(spa, spa->spa_first_txg, &spa->spa_dsl_pool);
2274     if (error)
2275         return (spa_vdev_err(rvd, VDEV_AUX_CORRUPT_DATA, EIO));
2276     spa->spa_meta_objset = spa->spa_dsl_pool->dp_meta_objset;
2277
2278     if (spa_dir_prop(spa, DMU_POOL_CONFIG, &spa->spa_config_object) != 0)
2279         return (spa_vdev_err(rvd, VDEV_AUX_CORRUPT_DATA, EIO));
2280
2281     if (spa_version(spa) >= SPA_VERSION_FEATURES) {
2282         boolean_t missing_feat_read = B_FALSE;
2283         nvlist_t *unsup_feat, *enabled_feat;
2284
2285         if (spa_dir_prop(spa, DMU_POOL_FEATURES_FOR_READ,
2286             &spa->spa_feat_for_read_obj) != 0) {
2287             return (spa_vdev_err(rvd, VDEV_AUX_CORRUPT_DATA, EIO));
2288         }
2289
2290         if (spa_dir_prop(spa, DMU_POOL_FEATURES_FOR_WRITE,
2291             &spa->spa_feat_for_write_obj) != 0) {
2292             return (spa_vdev_err(rvd, VDEV_AUX_CORRUPT_DATA, EIO));
2293         }
2294
2295         if (spa_dir_prop(spa, DMU_POOL_FEATURE_DESCRIPTIONS,
2296             &spa->spa_feat_desc_obj) != 0) {
2297             return (spa_vdev_err(rvd, VDEV_AUX_CORRUPT_DATA, EIO));
2298         }
2299
2300         enabled_feat = fnvlist_alloc();
2301         unsup_feat = fnvlist_alloc();
2302
2303         if (!spa_features_check(spa, B_FALSE,
2304             if (!feature_is_supported(spa->spa_meta_objset,
2305                 spa->spa_feat_for_read_obj, spa->spa_feat_desc_obj,
2306                     unsup_feat, enabled_feat)) {
2307                 missing_feat_read = B_TRUE;
2308             }
2309
2310         if (spa_writeable(spa) || state == SPA_LOAD_TRYIMPORT) {
2311             if (!spa_features_check(spa, B_TRUE,
2312                 if (!feature_is_supported(spa->spa_meta_objset,
2313                     spa->spa_feat_for_write_obj, spa->spa_feat_desc_obj,
2314                         unsup_feat, enabled_feat)) {
2315                     missing_feat_write = B_TRUE;
2316                 }
2317             }
2318         }
2319
2320         fnvlist_add_nvlist(spa->spa_load_info,
2321             ZPOOL_CONFIG_ENABLED_FEAT, enabled_feat);
2322
2323         if (!nvlist_empty(unsup_feat)) {
2324             fnvlist_add_nvlist(spa->spa_load_info,

```

```

2319         ZPOOL_CONFIG_UNSUP_FEAT, unsup_feat);
2320     }
2322     fnvlist_free(enabled_feat);
2323     fnvlist_free(unsup_feat);
2325     if (!missing_feat_read) {
2326         fnvlist_add_boolean(spa->spa_load_info,
2327             ZPOOL_CONFIG_CAN_RDONLY);
2328     }
2330     /*
2331     * If the state is SPA_LOAD_TRYIMPORT, our objective is
2332     * twofold: to determine whether the pool is available for
2333     * import in read-write mode and (if it is not) whether the
2334     * pool is available for import in read-only mode. If the pool
2335     * is available for import in read-write mode, it is displayed
2336     * as available in userland; if it is not available for import
2337     * in read-only mode, it is displayed as unavailable in
2338     * userland. If the pool is available for import in read-only
2339     * mode but not read-write mode, it is displayed as unavailable
2340     * in userland with a special note that the pool is actually
2341     * available for open in read-only mode.
2342     *
2343     * As a result, if the state is SPA_LOAD_TRYIMPORT and we are
2344     * missing a feature for write, we must first determine whether
2345     * the pool can be opened read-only before returning to
2346     * userland in order to know whether to display the
2347     * abovementioned note.
2348     */
2349     if (missing_feat_read || (missing_feat_write &&
2350         spa_writeable(spa))) {
2351         return (spa_vdev_err(rvd, VDEV_AUX_UNSUP_FEAT,
2352             ENOTSUP));
2353     }
2354 }
2356 spa->spa_is_initializing = B_TRUE;
2357 error = dsl_pool_open(spa->spa_dsl_pool);
2358 spa->spa_is_initializing = B_FALSE;
2359 if (error != 0)
2360     return (spa_vdev_err(rvd, VDEV_AUX_CORRUPT_DATA, EIO));
2362 if (!mosconfig) {
2363     uint64_t hostid;
2364     nvlist_t *policy = NULL, *nvconfig;
2366     if (load_nvlist(spa, spa->spa_config_object, &nvconfig) != 0)
2367         return (spa_vdev_err(rvd, VDEV_AUX_CORRUPT_DATA, EIO));
2369     if (!spa_is_root(spa) && nvlist_lookup_uint64(nvconfig,
2370         ZPOOL_CONFIG_HOSTID, &hostid) == 0) {
2371         char *hostname;
2372         unsigned long myhostid = 0;
2374         VERIFY(nvlist_lookup_string(nvconfig,
2375             ZPOOL_CONFIG_HOSTNAME, &hostname) == 0);
2377 #ifdef _KERNEL
2378         myhostid = zone_get_hostid(NULL);
2379 #else /* _KERNEL */
2380         /*
2381         * We're emulating the system's hostid in userland, so
2382         * we can't use zone_get_hostid().
2383         */
2384         (void) ddi_strtoul(hw_serial, NULL, 10, &myhostid);

```

```

2385 #endif /* _KERNEL */
2386         if (hostid != 0 && myhostid != 0 &&
2387             hostid != myhostid) {
2388             nvlist_free(nvconfig);
2389             cmn_err(CE_WARN, "pool '%s' could not be "
2390                 "loaded as it was last accessed by "
2391                 "another system (host: %s hostid: 0x%lx). "
2392                 "See: http://illumos.org/msg/ZFS-8000-EY",
2393                 spa_name(spa), hostname,
2394                 (unsigned long)hostid);
2395             return (SET_ERROR(EBADF));
2396         }
2397     }
2398     if (nvlist_lookup_nvlist(spa->spa_config,
2399         ZPOOL_REWIND_POLICY, &policy) == 0)
2400         VERIFY(nvlist_add_nvlist(nvconfig,
2401             ZPOOL_REWIND_POLICY, policy) == 0);
2403     spa_config_set(spa, nvconfig);
2404     spa_unload(spa);
2405     spa_deactivate(spa);
2406     spa_activate(spa, orig_mode);
2408     return (spa_load(spa, state, SPA_IMPORT_EXISTING, B_TRUE));
2409 }
2411 if (spa_dir_prop(spa, DMU_POOL_SYNC_BPOBJ, &obj) != 0)
2412     return (spa_vdev_err(rvd, VDEV_AUX_CORRUPT_DATA, EIO));
2413 error = bpobj_open(&spa->spa_deferred_bpobj, spa->spa_meta_objset, obj);
2414 if (error != 0)
2415     return (spa_vdev_err(rvd, VDEV_AUX_CORRUPT_DATA, EIO));
2417 /*
2418 * Load the bit that tells us to use the new accounting function
2419 * (raid-z deflation). If we have an older pool, this will not
2420 * be present.
2421 */
2422 error = spa_dir_prop(spa, DMU_POOL_DEFLATE, &spa->spa_deflate);
2423 if (error != 0 && error != ENOENT)
2424     return (spa_vdev_err(rvd, VDEV_AUX_CORRUPT_DATA, EIO));
2426 error = spa_dir_prop(spa, DMU_POOL_CREATION_VERSION,
2427     &spa->spa_creation_version);
2428 if (error != 0 && error != ENOENT)
2429     return (spa_vdev_err(rvd, VDEV_AUX_CORRUPT_DATA, EIO));
2431 /*
2432 * Load the persistent error log. If we have an older pool, this will
2433 * not be present.
2434 */
2435 error = spa_dir_prop(spa, DMU_POOL_ERRLOG_LAST, &spa->spa_errlog_last);
2436 if (error != 0 && error != ENOENT)
2437     return (spa_vdev_err(rvd, VDEV_AUX_CORRUPT_DATA, EIO));
2439 error = spa_dir_prop(spa, DMU_POOL_ERRLOG_SCRUB,
2440     &spa->spa_errlog_scrub);
2441 if (error != 0 && error != ENOENT)
2442     return (spa_vdev_err(rvd, VDEV_AUX_CORRUPT_DATA, EIO));
2444 /*
2445 * Load the history object. If we have an older pool, this
2446 * will not be present.
2447 */
2448 error = spa_dir_prop(spa, DMU_POOL_HISTORY, &spa->spa_history);
2449 if (error != 0 && error != ENOENT)
2450     return (spa_vdev_err(rvd, VDEV_AUX_CORRUPT_DATA, EIO));

```

```

2452 /*
2453  * If we're assembling the pool from the split-off vdevs of
2454  * an existing pool, we don't want to attach the spares & cache
2455  * devices.
2456  */
2458 /*
2459  * Load any hot spares for this pool.
2460  */
2461 error = spa_dir_prop(spa, DMU_POOL_SPARES, &spa->spa_spares.sav_object);
2462 if (error != 0 && error != ENOENT)
2463     return (spa_vdev_err(rvd, VDEV_AUX_CORRUPT_DATA, EIO));
2464 if (error == 0 && type != SPA_IMPORT_ASSEMBLE) {
2465     ASSERT(spa_version(spa) >= SPA_VERSION_SPARES);
2466     if (load_nvlist(spa, spa->spa_spares.sav_object,
2467         &spa->spa_spares.sav_config) != 0)
2468         return (spa_vdev_err(rvd, VDEV_AUX_CORRUPT_DATA, EIO));
2470     spa_config_enter(spa, SCL_ALL, FTAG, RW_WRITER);
2471     spa_load_spares(spa);
2472     spa_config_exit(spa, SCL_ALL, FTAG);
2473 } else if (error == 0) {
2474     spa->spa_spares.sav_sync = B_TRUE;
2475 }
2477 /*
2478  * Load any level 2 ARC devices for this pool.
2479  */
2480 error = spa_dir_prop(spa, DMU_POOL_L2CACHE,
2481     &spa->spa_l2cache.sav_object);
2482 if (error != 0 && error != ENOENT)
2483     return (spa_vdev_err(rvd, VDEV_AUX_CORRUPT_DATA, EIO));
2484 if (error == 0 && type != SPA_IMPORT_ASSEMBLE) {
2485     ASSERT(spa_version(spa) >= SPA_VERSION_L2CACHE);
2486     if (load_nvlist(spa, spa->spa_l2cache.sav_object,
2487         &spa->spa_l2cache.sav_config) != 0)
2488         return (spa_vdev_err(rvd, VDEV_AUX_CORRUPT_DATA, EIO));
2490     spa_config_enter(spa, SCL_ALL, FTAG, RW_WRITER);
2491     spa_load_l2cache(spa);
2492     spa_config_exit(spa, SCL_ALL, FTAG);
2493 } else if (error == 0) {
2494     spa->spa_l2cache.sav_sync = B_TRUE;
2495 }
2497 spa->spa_delegation = zpool_prop_default_numeric(ZPOOL_PROP_DELEGATION);
2499 error = spa_dir_prop(spa, DMU_POOL_PROPS, &spa->spa_pool_props_object);
2500 if (error && error != ENOENT)
2501     return (spa_vdev_err(rvd, VDEV_AUX_CORRUPT_DATA, EIO));
2503 if (error == 0) {
2504     uint64_t autoreplace;
2506     spa_prop_find(spa, ZPOOL_PROP_BOOTFS, &spa->spa_bootfs);
2507     spa_prop_find(spa, ZPOOL_PROP_AUTOREPLACE, &autoreplace);
2508     spa_prop_find(spa, ZPOOL_PROP_DELEGATION, &spa->spa_delegation);
2509     spa_prop_find(spa, ZPOOL_PROP_FAILUREMODE, &spa->spa_failmode);
2510     spa_prop_find(spa, ZPOOL_PROP_AUTOEXPAND, &spa->spa_autoexpand);
2511     spa_prop_find(spa, ZPOOL_PROP_DEDUPDITTO,
2512         &spa->spa_dedup_ditto);
2514     spa->spa_autoreplace = (autoreplace != 0);
2515 }

```

```

2517 /*
2518  * If the 'autoreplace' property is set, then post a resource notifying
2519  * the ZFS DE that it should not issue any faults for unopenable
2520  * devices. We also iterate over the vdevs, and post a sysevent for any
2521  * unopenable vdevs so that the normal autoreplace handler can take
2522  * over.
2523  */
2524 if (spa->spa_autoreplace && state != SPA_LOAD_TRYIMPORT) {
2525     spa_check_removed(spa->spa_root_vdev);
2526     /*
2527      * For the import case, this is done in spa_import(), because
2528      * at this point we're using the spare definitions from
2529      * the MOS config, not necessarily from the userland config.
2530      */
2531     if (state != SPA_LOAD_IMPORT) {
2532         spa_aux_check_removed(&spa->spa_spares);
2533         spa_aux_check_removed(&spa->spa_l2cache);
2534     }
2535 }
2537 /*
2538  * Load the vdev state for all toplevel vdevs.
2539  */
2540 vdev_load(rvd);
2542 /*
2543  * Propagate the leaf DTLs we just loaded all the way up the tree.
2544  */
2545 spa_config_enter(spa, SCL_ALL, FTAG, RW_WRITER);
2546 vdev_dtl_reassess(rvd, 0, 0, B_FALSE);
2547 spa_config_exit(spa, SCL_ALL, FTAG);
2549 /*
2550  * Load the DDTs (dedup tables).
2551  */
2552 error = ddt_load(spa);
2553 if (error != 0)
2554     return (spa_vdev_err(rvd, VDEV_AUX_CORRUPT_DATA, EIO));
2556 spa_update_dspace(spa);
2558 /*
2559  * Validate the config, using the MOS config to fill in any
2560  * information which might be missing. If we fail to validate
2561  * the config then declare the pool unfit for use. If we're
2562  * assembling a pool from a split, the log is not transferred
2563  * over.
2564  */
2565 if (type != SPA_IMPORT_ASSEMBLE) {
2566     nvlist_t *nvconfig;
2568     if (load_nvlist(spa, spa->spa_config_object, &nvconfig) != 0)
2569         return (spa_vdev_err(rvd, VDEV_AUX_CORRUPT_DATA, EIO));
2571     if (!spa_config_valid(spa, nvconfig)) {
2572         nvlist_free(nvconfig);
2573         return (spa_vdev_err(rvd, VDEV_AUX_BAD_GUID_SUM,
2574             ENXIO));
2575     }
2576     nvlist_free(nvconfig);
2578     /*
2579      * Now that we've validated the config, check the state of the
2580      * root vdev. If it can't be opened, it indicates one or
2581      * more toplevel vdevs are faulted.
2582      */

```

```

2583     if (rvd->vdev_state <= VDEV_STATE_CANT_OPEN)
2584         return (SET_ERROR(ENXIO));

2586     if (spa_check_logs(spa)) {
2587         *ereport = FM_EREPORT_ZFS_LOG_REPLAY;
2588         return (spa_vdev_err(rvd, VDEV_AUX_BAD_LOG, ENXIO));
2589     }
2590 }

2592 if (missing_feat_write) {
2593     ASSERT(state == SPA_LOAD_TRYIMPORT);

2595     /*
2596      * At this point, we know that we can open the pool in
2597      * read-only mode but not read-write mode. We now have enough
2598      * information and can return to userland.
2599      */
2600     return (spa_vdev_err(rvd, VDEV_AUX_UNSUP_FEAT, ENOTSUP));
2601 }

2603 /*
2604  * We've successfully opened the pool, verify that we're ready
2605  * to start pushing transactions.
2606  */
2607 if (state != SPA_LOAD_TRYIMPORT) {
2608     if (error = spa_load_verify(spa))
2609         return (spa_vdev_err(rvd, VDEV_AUX_CORRUPT_DATA,
2610             error));
2611 }

2613 if (spa_writeable(spa) && (state == SPA_LOAD_RECOVER ||
2614     spa->spa_load_max_txc == UINT64_MAX)) {
2615     dmu_tx_t *tx;
2616     int need_update = B_FALSE;

2618     ASSERT(state != SPA_LOAD_TRYIMPORT);

2620     /*
2621      * Claim log blocks that haven't been committed yet.
2622      * This must all happen in a single txg.
2623      * Note: spa_claim_max_txc is updated by spa_claim_notify(),
2624      * invoked from zil_claim_log_block()'s i/o done callback.
2625      * Price of rollback is that we abandon the log.
2626      */
2627     spa->spa_claiming = B_TRUE;

2629     tx = dmu_tx_create_assigned(spa_get_dsl(spa),
2630         spa_first_txc(spa));
2631     (void) dmu_objset_find(spa_name(spa),
2632         zil_claim, tx, DS_FIND_CHILDREN);
2633     dmu_tx_commit(tx);

2635     spa->spa_claiming = B_FALSE;

2637     spa_set_log_state(spa, SPA_LOG_GOOD);
2638     spa->spa_sync_on = B_TRUE;
2639     txg_sync_start(spa->spa_dsl_pool);

2641     /*
2642      * Wait for all claims to sync. We sync up to the highest
2643      * claimed log block birth time so that claimed log blocks
2644      * don't appear to be from the future. spa_claim_max_txc
2645      * will have been set for us by either zil_check_log_chain()
2646      * (invoked from spa_check_logs()) or zil_claim() above.
2647      */
2648     txg_wait_synced(spa->spa_dsl_pool, spa->spa_claim_max_txc);

```

```

2650     /*
2651      * If the config cache is stale, or we have uninitialized
2652      * metaslabs (see spa_vdev_add()), then update the config.
2653      */
2654     * If this is a verbatim import, trust the current
2655     * in-core spa_config and update the disk labels.
2656     */
2657     if (config_cache_txc != spa->spa_config_txc ||
2658         state == SPA_LOAD_IMPORT ||
2659         state == SPA_LOAD_RECOVER ||
2660         (spa->spa_import_flags & ZFS_IMPORT_VERBATIM))
2661         need_update = B_TRUE;

2663     for (int c = 0; c < rvd->vdev_children; c++)
2664         if (rvd->vdev_child[c]->vdev_ms_array == 0)
2665             need_update = B_TRUE;

2667     /*
2668      * Update the config cache asynchronously in case we're the
2669      * root pool, in which case the config cache isn't writable yet.
2670      */
2671     if (need_update)
2672         spa_async_request(spa, SPA_ASYNC_CONFIG_UPDATE);

2674     /*
2675      * Check all DTLs to see if anything needs resilvering.
2676      */
2677     if (!dsl_scan_resilvering(spa->spa_dsl_pool) &&
2678         vdev_resilver_needed(rvd, NULL, NULL))
2679         spa_async_request(spa, SPA_ASYNC_RESILVER);

2681     /*
2682      * Log the fact that we booted up (so that we can detect if
2683      * we rebooted in the middle of an operation).
2684      */
2685     spa_history_log_version(spa, "open");

2687     /*
2688      * Delete any inconsistent datasets.
2689      */
2690     (void) dmu_objset_find(spa_name(spa),
2691         dsl_destroy_inconsistent, NULL, DS_FIND_CHILDREN);

2693     /*
2694      * Clean up any stale temporary dataset userrefs.
2695      */
2696     dsl_pool_clean_tmp_userrefs(spa->spa_dsl_pool);
2697 }

2699     return (0);
2700 }

unchanged_portion_omitted_

5896 /*
5897  * Set zpool properties.
5898  */
5899 static void
5900 spa_sync_props(void *arg, dmu_tx_t *tx)
5901 {
5902     nvlist_t *nvp = arg;
5903     spa_t *spa = dmu_tx_pool(tx)->dp_spa;
5904     objset_t *mos = spa->spa_meta_objset;
5905     nvpair_t *elem = NULL;

5907     mutex_enter(&spa->spa_props_lock);

```

```

5909     while ((elem = nvlist_next_nvpair(nvp, elem)) {
5910         uint64_t intval;
5911         char *strval, *fname;
5912         zpool_prop_t prop;
5913         const char *propname;
5914         zprop_type_t proptype;
5915         spa_feature_t fid;
5916         zfeature_info_t *feature;
5917
5918         switch (prop = zpool_name_to_prop(nvpair_name(elem)) {
5919             case ZPROP_INVALID:
5920                 /*
5921                  * We checked this earlier in spa_prop_validate().
5922                  */
5923                 ASSERT(zpool_prop_feature(nvpair_name(elem)));
5924
5925                 fname = strchr(nvpair_name(elem), '@') + 1;
5926                 VERIFY0(zfeature_lookup_name(fname, &fid));
5927                 VERIFY0(zfeature_lookup_name(fname, &feature));
5928
5929                 spa_feature_enable(spa, fid, tx);
5930                 spa_feature_enable(spa, feature, tx);
5931                 spa_history_log_internal(spa, "set", tx,
5932                     "%s=enabled", nvpair_name(elem));
5933                 break;
5934
5935             case ZPOOL_PROP_VERSION:
5936                 intval = fnvpair_value_uint64(elem);
5937                 /*
5938                  * The version is synced seperatly before other
5939                  * properties and should be correct by now.
5940                  */
5941                 ASSERT3U(spa_version(spa), >=, intval);
5942                 break;
5943
5944             case ZPOOL_PROP_ALTROOT:
5945                 /*
5946                  * 'altroot' is a non-persistent property. It should
5947                  * have been set temporarily at creation or import time.
5948                  */
5949                 ASSERT(spa->spa_root != NULL);
5950                 break;
5951
5952             case ZPOOL_PROP_READONLY:
5953             case ZPOOL_PROP_CACHEFILE:
5954                 /*
5955                  * 'readonly' and 'cachefile' are also non-persistent
5956                  * properties.
5957                  */
5958                 break;
5959             case ZPOOL_PROP_COMMENT:
5960                 strval = fnvpair_value_string(elem);
5961                 if (spa->spa_comment != NULL)
5962                     spa_strfree(spa->spa_comment);
5963                 spa->spa_comment = spa_strdup(strval);
5964                 /*
5965                  * We need to dirty the configuration on all the vdevs
5966                  * so that their labels get updated. It's unnecessary
5967                  * to do this for pool creation since the vdev's
5968                  * configuratoin has already been dirtied.
5969                  */
5970                 if (tx->tx_txg != TXG_INITIAL)
5971                     vdev_config_dirty(spa->spa_root_vdev);
5972                 spa_history_log_internal(spa, "set", tx,
5973                     "%s=%s", nvpair_name(elem), strval);

```

```

5971         break;
5972     default:
5973         /*
5974          * Set pool property values in the poolprops mos object.
5975          */
5976         if (spa->spa_pool_props_object == 0) {
5977             spa->spa_pool_props_object =
5978                 zap_create_link(mos, DMU_OT_POOL_PROPS,
5979                     DMU_POOL_DIRECTORY_OBJECT, DMU_POOL_PROPS,
5980                     tx);
5981         }
5982
5983         /* normalize the property name */
5984         propname = zpool_prop_to_name(prop);
5985         proptype = zpool_prop_get_type(prop);
5986
5987         if (nvpair_type(elem) == DATA_TYPE_STRING) {
5988             ASSERT(proptype == PROP_TYPE_STRING);
5989             strval = fnvpair_value_string(elem);
5990             VERIFY0(zap_update(mos,
5991                 spa->spa_pool_props_object, propname,
5992                 1, strlen(strval) + 1, strval, tx));
5993             spa_history_log_internal(spa, "set", tx,
5994                 "%s=%s", nvpair_name(elem), strval);
5995         } else if (nvpair_type(elem) == DATA_TYPE_UINT64) {
5996             intval = fnvpair_value_uint64(elem);
5997
5998             if (proptype == PROP_TYPE_INDEX) {
5999                 const char *unused;
6000                 VERIFY0(zpool_prop_index_to_string(
6001                     prop, intval, &unused));
6002             }
6003             VERIFY0(zap_update(mos,
6004                 spa->spa_pool_props_object, propname,
6005                 8, 1, &intval, tx));
6006             spa_history_log_internal(spa, "set", tx,
6007                 "%s=%lld", nvpair_name(elem), intval);
6008         } else {
6009             ASSERT(0); /* not allowed */
6010         }
6011
6012         switch (prop) {
6013             case ZPOOL_PROP_DELEGATION:
6014                 spa->spa_delegation = intval;
6015                 break;
6016             case ZPOOL_PROP_BOOTFS:
6017                 spa->spa_bootfs = intval;
6018                 break;
6019             case ZPOOL_PROP_FAILUREMODE:
6020                 spa->spa_failmode = intval;
6021                 break;
6022             case ZPOOL_PROP_AUTOEXPAND:
6023                 spa->spa_autoexpand = intval;
6024                 if (tx->tx_txg != TXG_INITIAL)
6025                     spa_async_request(spa,
6026                         SPA_ASYNC_AUTOEXPAND);
6027                 break;
6028             case ZPOOL_PROP_DEDUPDITTO:
6029                 spa->spa_dedup_ditto = intval;
6030                 break;
6031         }
6032     default:
6033         break;
6034 }
6035 }
6036 }

```

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

13

```
6038         mutex_exit(&spa->spa_props_lock);
6039     }
_____unchanged_portion_omitted_
```

```
*****
46279 Tue Oct 1 14:04:41 2013
new/usr/src/uts/common/fs/zfs/spa_misc.c
4171 clean up spa_feature_*() interfaces
4172 implement extensible_dataset feature for use by other zpool features
Reviewed by: Max Grossman <max.grossman@delphix.com>
Reviewed by: Christopher Siden <christopher.siden@delphix.com>
Reviewed by: George Wilson <george.wilson@delphix.com>
*****
_____unchanged_portion_omitted_

1137 /*
1138 * =====
1139 * Miscellaneous functions
1140 * =====
1141 */

1143 void
1144 spa_activate_mos_feature(spa_t *spa, const char *feature)
1145 {
1146     if (!nvlist_exists(spa->spa_label_features, feature)) {
1147         nvlist_add_boolean(spa->spa_label_features, feature);
1148         (void) nvlist_add_boolean(spa->spa_label_features, feature);
1149         vdev_config_dirty(spa->spa_root_vdev);
1150     }

1152 void
1153 spa_deactivate_mos_feature(spa_t *spa, const char *feature)
1154 {
1155     if (nvlist_remove_all(spa->spa_label_features, feature) == 0)
1156         (void) nvlist_remove_all(spa->spa_label_features, feature);
1157     vdev_config_dirty(spa->spa_root_vdev);
1158 }
_____unchanged_portion_omitted_
```

```

*****
15742 Tue Oct 1 14:04:43 2013
new/usr/src/uts/common/fs/zfs/space_map.c
4171 clean up spa_feature_*() interfaces
4172 implement extensible_dataset feature for use by other zpool features
Reviewed by: Max Grossman <max.grossman@delphix.com>
Reviewed by: Christopher Siden <christopher.siden@delphix.com>
Reviewed by: George Wilson <george.wilson@delphix.com>
*****
_____unchanged_portion_omitted_____

472 void
473 space_map_truncate(space_map_t *sm, dmu_tx_t *tx)
474 {
475     objset_t *os = sm->sm_os;
476     spa_t *spa = dmu_objset_spa(os);
477     zfeature_info_t *space_map_histogram =
478         &spa_feature_table[SPA_FEATURE_SPACEMAP_HISTOGRAM];
477     dmu_object_info_t doi;
478     int bonuslen;

480     ASSERT(dsl_pool_sync_context(dmu_objset_pool(os)));
481     ASSERT(dmu_tx_is_syncing(tx));

483     VERIFY0(dmu_free_range(os, space_map_object(sm), 0, -1ULL, tx));
484     dmu_object_info_from_db(sm->sm_dbuf, &doi);

486     if (spa_feature_is_enabled(spa, SPA_FEATURE_SPACEMAP_HISTOGRAM)) {
488         if (spa_feature_is_enabled(spa, space_map_histogram)) {
487             bonuslen = sizeof(space_map_phys_t);
488             ASSERT3U(bonuslen, <=, dmu_bonus_max());
489         } else {
490             bonuslen = SPACE_MAP_SIZE_V0;
491         }

493         if (bonuslen != doi.doi_bonus_size ||
494             doi.doi_data_block_size != SPACE_MAP_INITIAL_BLOCKSIZE) {
495             zfs_dbgmsg("txg %llu, spa %s, reallocating: "
496                 "old bonus %u, old blocksz %u", dmu_tx_get_txg(tx),
497                 spa_name(spa), doi.doi_bonus_size, doi.doi_data_block_size);
498             space_map_reallocate(sm, tx);
499             VERIFY3U(sm->sm_blkisz, ==, SPACE_MAP_INITIAL_BLOCKSIZE);
500         }

502         dmu_buf_will_dirty(sm->sm_dbuf, tx);
503         sm->sm_phys->smp_objsize = 0;
504         sm->sm_phys->smp_alloc = 0;
505     }
_____unchanged_portion_omitted_____

522 uint64_t
523 space_map_alloc(objset_t *os, dmu_tx_t *tx)
524 {
525     spa_t *spa = dmu_objset_spa(os);
528     zfeature_info_t *space_map_histogram =
529         &spa_feature_table[SPA_FEATURE_SPACEMAP_HISTOGRAM];
526     uint64_t object;
527     int bonuslen;

529     if (spa_feature_is_enabled(spa, SPA_FEATURE_SPACEMAP_HISTOGRAM)) {
530         spa_feature_incr(spa, SPA_FEATURE_SPACEMAP_HISTOGRAM, tx);
533     if (spa_feature_is_enabled(spa, space_map_histogram)) {
534         spa_feature_incr(spa, space_map_histogram, tx);
531         bonuslen = sizeof(space_map_phys_t);
532         ASSERT3U(bonuslen, <=, dmu_bonus_max());
533     } else {

```

```

534         bonuslen = SPACE_MAP_SIZE_V0;
535     }

537     object = dmu_object_alloc(os,
538         DMU_OT_SPACE_MAP, SPACE_MAP_INITIAL_BLOCKSIZE,
539         DMU_OT_SPACE_MAP_HEADER, bonuslen, tx);

541     return (object);
542 }

544 void
545 space_map_free(space_map_t *sm, dmu_tx_t *tx)
546 {
547     spa_t *spa;
552     zfeature_info_t *space_map_histogram =
553         &spa_feature_table[SPA_FEATURE_SPACEMAP_HISTOGRAM];

549     if (sm == NULL)
550         return;

552     spa = dmu_objset_spa(sm->sm_os);
553     if (spa_feature_is_enabled(spa, SPA_FEATURE_SPACEMAP_HISTOGRAM)) {
559         if (spa_feature_is_enabled(spa, space_map_histogram)) {
554             dmu_object_info_t doi;

556             dmu_object_info_from_db(sm->sm_dbuf, &doi);
557             if (doi.doi_bonus_size != SPACE_MAP_SIZE_V0) {
558                 VERIFY(spa_feature_is_active(spa,
559                     SPA_FEATURE_SPACEMAP_HISTOGRAM));
560                 spa_feature_decr(spa,
561                     SPA_FEATURE_SPACEMAP_HISTOGRAM, tx);
564                 VERIFY(spa_feature_is_active(spa, space_map_histogram));
565                 spa_feature_decr(spa, space_map_histogram, tx);
562             }
563         }

565     VERIFY3U(dmu_object_free(sm->sm_os, space_map_object(sm), tx), ==, 0);
566     sm->sm_object = 0;
567 }
_____unchanged_portion_omitted_____

```



new/usr/src/uts/common/fs/zfs/sys/dmu\_impl.h

1

\*\*\*\*\*

8514 Tue Oct 1 14:04:44 2013

new/usr/src/uts/common/fs/zfs/sys/dmu\_impl.h

4171 clean up spa\_feature\_\*() interfaces

4172 implement extensible\_dataset feature for use by other zpool features

Reviewed by: Max Grossman <max.grossman@delphix.com>

Reviewed by: Christopher Siden <christopher.siden@delphix.com>

Reviewed by: George Wilson <george.wilson@delphix.com>

\*\*\*\*\*

\_\_\_\_\_ unchanged\_portion\_omitted\_

301 void dmu\_object\_zapify(objset\_t \*, uint64\_t, dmu\_object\_type\_t, dmu\_tx\_t \*);

302 void dmu\_object\_free\_zapified(objset\_t \*, uint64\_t, dmu\_tx\_t \*);

304 #ifdef \_\_cplusplus

305 }

\_\_\_\_\_ unchanged\_portion\_omitted\_

```

*****
10682 Tue Oct 1 14:04:45 2013
new/usr/src/uts/common/fs/zfs/sys/dnode.h
4171 clean up spa_feature*() interfaces
4172 implement extensible_dataset feature for use by other zpool features
Reviewed by: Max Grossman <max.grossman@delphix.com>
Reviewed by: Christopher Siden <christopher.siden@delphix.com>
Reviewed by: George Wilson <george.wilson@delphix.com>
*****
_____ unchanged_portion_omitted_

146 typedef struct dnode {
147     /*
148     * Protects the structure of the dnode, including the number of levels
149     * of indirection (dn_nlevels), dn_maxblkid, and dn_next_*
150     */
151     krwlock_t dn_struct_rwlock;

153     /* Our link on dn_objset->os_dnodes list; protected by os_lock. */
154     list_node_t dn_link;

156     /* immutable: */
157     struct objset *dn_objset;
158     uint64_t dn_object;
159     struct dmu_buf_impl *dn_dbuf;
160     struct dnode_handle *dn_handle;
161     dnode_phys_t *dn_phys; /* pointer into dn->dn_dbuf->db.db_data */

163     /*
164     * Copies of stuff in dn_phys. They're valid in the open
165     * context (eg. even before the dnode is first synced).
166     * Where necessary, these are protected by dn_struct_rwlock.
167     */
168     dmu_object_type_t dn_type; /* object type */
169     uint16_t dn_bonuslen; /* bonus length */
170     uint8_t dn_bonustype; /* bonus type */
171     uint8_t dn_nblkptr; /* number of blkptrs (immutable) */
172     uint8_t dn_checksum; /* ZIO_CHECKSUM type */
173     uint8_t dn_compress; /* ZIO_COMPRESS type */
174     uint8_t dn_nlevels;
175     uint8_t dn_indblkshift;
176     uint8_t dn_datablkshift; /* zero if blksize not power of 2! */
177     uint8_t dn_moved; /* Has this dnode been moved? */
178     uint16_t dn_datablkszsec; /* in 512b sectors */
179     uint32_t dn_datablksz; /* in bytes */
180     uint64_t dn_maxblkid;
181     uint8_t dn_next_type[TXG_SIZE];
182     uint8_t dn_next_nblkptr[TXG_SIZE];
183     uint8_t dn_next_nlevels[TXG_SIZE];
184     uint8_t dn_next_indblkshift[TXG_SIZE];
185     uint8_t dn_next_bonustype[TXG_SIZE];
186     uint8_t dn_rm_spillblk[TXG_SIZE]; /* for removing spill blk */
187     uint16_t dn_next_bonuslen[TXG_SIZE];
188     uint32_t dn_next_blksize[TXG_SIZE]; /* next block size in bytes */

190     /* protected by dn_dbufs_mtx; declared here to fill 32-bit hole */
191     uint32_t dn_dbufs_count; /* count of dn_dbufs */
192     /* There are no level-0 blocks of this blkid or higher in dn_dbufs */
193     uint64_t dn_unlisted_l0_blkid;

195     /* protected by os_lock: */
196     list_node_t dn_dirty_link[TXG_SIZE]; /* next on dataset's dirty */

198     /* protected by dn_mtx: */
199     kmutex_t dn_mtx;
200     list_t dn_dirty_records[TXG_SIZE];

```

```

201     avl_tree_t dn_ranges[TXG_SIZE];
202     uint64_t dn_allocated_txg;
203     uint64_t dn_free_txg;
204     uint64_t dn_assigned_txg;
205     kcondvar_t dn_notxholds;
206     enum dnode_dirtycontext dn_dirtyctx;
207     uint8_t *dn_dirtyctx_firstset; /* dbg: contents meaningless */

209     /* protected by own devices */
210     refcount_t dn_tx_holds;
211     refcount_t dn_holds;

213     kmutex_t dn_dbufs_mtx;
214     list_t dn_dbufs; /* descendent dbufs */

216     /* protected by dn_struct_rwlock */
217     struct dmu_buf_impl *dn_bonus; /* bonus buffer dbuf */

219     boolean_t dn_have_spill; /* have spill or are spilling */

221     /* parent IO for current sync write */
222     zio_t *dn_zio;

224     /* used in syncing context */
225     uint64_t dn_oldused; /* old phys used bytes */
226     uint64_t dn_oldflags; /* old phys dn_flags */
227     uint64_t dn_olduid, dn_oldgid;
228     uint64_t dn_newuid, dn_newgid;
229     int dn_id_flags;

231     /* holds prefetch structure */
232     struct zfetch dn_zfetch;
233 } dnode_t;
_____ unchanged_portion_omitted_

```

```

*****
10351 Tue Oct 1 14:04:46 2013
new/usr/src/uts/common/fs/zfs/sys/dsl_dataset.h
4171 clean up spa_feature_*() interfaces
4172 implement extensible_dataset feature for use by other zpool features
Reviewed by: Max Grossman <max.grossman@delphix.com>
Reviewed by: Christopher Siden <christopher.siden@delphix.com>
Reviewed by: George Wilson <george.wilson@delphix.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22  * Copyright (c) 2005, 2010, Oracle and/or its affiliates. All rights reserved.
23  * Copyright (c) 2013 by Delphix. All rights reserved.
24  * Copyright (c) 2012, Joyent, Inc. All rights reserved.
25  * Copyright (c) 2013 Steven Hartland. All rights reserved.
26  */

28 #ifndef _SYS_DSL_DATASET_H
29 #define _SYS_DSL_DATASET_H

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

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

45 struct dsl_dataset;
46 struct dsl_dir;
47 struct dsl_pool;

49 #define DS_FLAG_INCONSISTENT (1ULL<<0)
50 #define DS_IS_INCONSISTENT(ds) \
51     ((ds)->ds_phys->ds_flags & DS_FLAG_INCONSISTENT)

53 /*
54  * Do not allow this dataset to be promoted.
55  * Note: nopromote can not yet be set, but we want support for it in this
56  * on-disk version, so that we don't need to upgrade for it later.
57  */

```

```

56 #define DS_FLAG_NOPROMOTE (1ULL<<1)

58 /*
59  * DS_FLAG_UNIQUE_ACCURATE is set if ds_unique_bytes has been correctly
60  * calculated for head datasets (starting with SPA_VERSION_UNIQUE_ACCURATE,
61  * refquota/refreservations).
62  */
63 #define DS_FLAG_UNIQUE_ACCURATE (1ULL<<2)

65 /*
66  * DS_FLAG_DEFER_DESTROY is set after 'zfs destroy -d' has been called
67  * on a dataset. This allows the dataset to be destroyed using 'zfs release'.
68  */
69 #define DS_FLAG_DEFER_DESTROY (1ULL<<3)
70 #define DS_IS_DEFER_DESTROY(ds) \
71     ((ds)->ds_phys->ds_flags & DS_FLAG_DEFER_DESTROY)

73 /*
74  * DS_FIELD_* are strings that are used in the "extensified" dataset zap object.
75  * They should be of the format <reverse-dns>:<field>.
76  */

78 /*
79  * DS_FLAG_CI_DATASET is set if the dataset contains a file system whose
80  * name lookups should be performed case-insensitively.
81  */
82 #define DS_FLAG_CI_DATASET (1ULL<<16)

84 #define DS_CREATE_FLAG_NODIRTY (1ULL<<24)

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

```

unchanged portion omitted

```

*****
17686 Tue Oct 1 14:04:47 2013
new/usr/src/uts/common/fs/zfs/sys/zap.h
4171 clean up spa_feature_*() interfaces
4172 implement extensible_dataset feature for use by other zpool features
Reviewed by: Max Grossman <max.grossman@delphix.com>
Reviewed by: Christopher Siden <christopher.siden@delphix.com>
Reviewed by: George Wilson <george.wilson@delphix.com>
*****
unchanged portion omitted

120 /*
121 * Create a new zapobj with no attributes and return its object number.
122 * MT_EXACT will cause the zap object to only support MT_EXACT lookups,
123 * otherwise any matchtype can be used for lookups.
124 *
125 * normflags specifies what normalization will be done. values are:
126 * 0: no normalization (legacy on-disk format, supports MT_EXACT matching
127 * only)
128 * U8_TEXTPREP_TOLOWER: case normalization will be performed.
129 * MT_FIRST/MT_BEST matching will find entries that match without
130 * regard to case (eg. looking for "foo" can find an entry "Foo").
131 * Eventually, other flags will permit unicode normalization as well.
132 */
133 uint64_t zap_create(objset_t *ds, dmu_object_type_t ot,
134 dmu_object_type_t bonustype, int bonuslen, dmu_tx_t *tx);
135 uint64_t zap_create_norm(objset_t *ds, int normflags, dmu_object_type_t ot,
136 dmu_object_type_t bonustype, int bonuslen, dmu_tx_t *tx);
137 uint64_t zap_create_flags(objset_t *os, int normflags, zap_flags_t flags,
138 dmu_object_type_t ot, int leaf_blockshift, int indirect_blockshift,
139 dmu_object_type_t bonustype, int bonuslen, dmu_tx_t *tx);
140 uint64_t zap_create_link(objset_t *os, dmu_object_type_t ot,
141 uint64_t parent_obj, const char *name, dmu_tx_t *tx);

143 /*
144 * Initialize an already-allocated object.
145 */
146 void mzap_create_impl(objset_t *os, uint64_t obj, int normflags,
147 zap_flags_t flags, dmu_tx_t *tx);

149 /*
150 * Create a new zapobj with no attributes from the given (unallocated)
151 * object number.
152 */
153 int zap_create_claim(objset_t *ds, uint64_t obj, dmu_object_type_t ot,
154 dmu_object_type_t bonustype, int bonuslen, dmu_tx_t *tx);
155 int zap_create_claim_norm(objset_t *ds, uint64_t obj,
156 int normflags, dmu_object_type_t ot,
157 dmu_object_type_t bonustype, int bonuslen, dmu_tx_t *tx);

159 /*
160 * The zapobj passed in must be a valid ZAP object for all of the
161 * following routines.
162 */

164 /*
165 * Destroy this zapobj and all its attributes.
166 *
167 * Frees the object number using dmu_object_free.
168 */
169 int zap_destroy(objset_t *ds, uint64_t zapobj, dmu_tx_t *tx);

171 /*
172 * Manipulate attributes.
173 *
174 * 'integer_size' is in bytes, and must be 1, 2, 4, or 8.

```

```

175 */

177 /*
178 * Retrieve the contents of the attribute with the given name.
179 *
180 * If the requested attribute does not exist, the call will fail and
181 * return ENOENT.
182 *
183 * If 'integer_size' is smaller than the attribute's integer size, the
184 * call will fail and return EINVAL.
185 *
186 * If 'integer_size' is equal to or larger than the attribute's integer
187 * size, the call will succeed and return 0.
188 *
189 * When converting to a larger integer size, the integers will be treated as
190 * unsigned (ie. no sign-extension will be performed).
191 *
192 * 'num_integers' is the length (in integers) of 'buf'.
193 *
194 * If the attribute is longer than the buffer, as many integers as will
195 * fit will be transferred to 'buf'. If the entire attribute was not
196 * transferred, the call will return EOVERFLOW.
197 */
198 int zap_lookup(objset_t *ds, uint64_t zapobj, const char *name,
199 uint64_t integer_size, uint64_t num_integers, void *buf);

201 /*
202 * If rn_len is nonzero, realname will be set to the name of the found
203 * entry (which may be different from the requested name if matchtype is
204 * not MT_EXACT).
205 *
206 * If normalization_conflict is not NULL, it will be set if there is
207 * another name with the same case/unicode normalized form.
208 */
209 int zap_lookup_norm(objset_t *ds, uint64_t zapobj, const char *name,
210 uint64_t integer_size, uint64_t num_integers, void *buf,
211 matchtype_t mt, char *realname, int rn_len,
212 boolean_t *normalization_conflict);
213 int zap_lookup_uint64(objset_t *os, uint64_t zapobj, const uint64_t *key,
214 int key_numints, uint64_t integer_size, uint64_t num_integers, void *buf);
215 int zap_contains(objset_t *ds, uint64_t zapobj, const char *name);
216 int zap_prefetch_uint64(objset_t *os, uint64_t zapobj, const uint64_t *key,
217 int key_numints);

219 int zap_count_write(objset_t *os, uint64_t zapobj, const char *name,
220 int add, uint64_t *towrite, uint64_t *tooverwrite);

222 /*
223 * Create an attribute with the given name and value.
224 *
225 * If an attribute with the given name already exists, the call will
226 * fail and return EEXIST.
227 */
228 int zap_add(objset_t *ds, uint64_t zapobj, const char *key,
229 int integer_size, uint64_t num_integers,
230 const void *val, dmu_tx_t *tx);
231 int zap_add_uint64(objset_t *ds, uint64_t zapobj, const uint64_t *key,
232 int key_numints, int integer_size, uint64_t num_integers,
233 const void *val, dmu_tx_t *tx);

235 /*
236 * Set the attribute with the given name to the given value. If an
237 * attribute with the given name does not exist, it will be created. If
238 * an attribute with the given name already exists, the previous value
239 * will be overwritten. The integer_size may be different from the
240 * existing attribute's integer size, in which case the attribute's

```

```

241 * integer size will be updated to the new value.
242 */
243 int zap_update(objset_t *ds, uint64_t zapobj, const char *name,
244 int integer_size, uint64_t num_integers, const void *val, dmu_tx_t *tx);
245 int zap_update_uint64(objset_t *os, uint64_t zapobj, const uint64_t *key,
246 int key_numints,
247 int integer_size, uint64_t num_integers, const void *val, dmu_tx_t *tx);
248
249 /*
250 * Get the length (in integers) and the integer size of the specified
251 * attribute.
252 *
253 * If the requested attribute does not exist, the call will fail and
254 * return ENOENT.
255 */
256 int zap_length(objset_t *ds, uint64_t zapobj, const char *name,
257 uint64_t *integer_size, uint64_t *num_integers);
258 int zap_length_uint64(objset_t *os, uint64_t zapobj, const uint64_t *key,
259 int key_numints, uint64_t *integer_size, uint64_t *num_integers);
260
261 /*
262 * Remove the specified attribute.
263 *
264 * If the specified attribute does not exist, the call will fail and
265 * return ENOENT.
266 */
267 int zap_remove(objset_t *ds, uint64_t zapobj, const char *name, dmu_tx_t *tx);
268 int zap_remove_norm(objset_t *ds, uint64_t zapobj, const char *name,
269 matchtype_t mt, dmu_tx_t *tx);
270 int zap_remove_uint64(objset_t *os, uint64_t zapobj, const uint64_t *key,
271 int key_numints, dmu_tx_t *tx);
272
273 /*
274 * Returns (in *count) the number of attributes in the specified zap
275 * object.
276 */
277 int zap_count(objset_t *ds, uint64_t zapobj, uint64_t *count);
278
279 /*
280 * Returns (in name) the name of the entry whose (value & mask)
281 * (za_first_integer) is value, or ENOENT if not found. The string
282 * pointed to by name must be at least 256 bytes long. If mask==0, the
283 * match must be exact (ie, same as mask=-1ULL).
284 */
285 int zap_value_search(objset_t *os, uint64_t zapobj,
286 uint64_t value, uint64_t mask, char *name);
287
288 /*
289 * Transfer all the entries from fromobj into intoobj. Only works on
290 * int_size=8 num_integers=1 values. Fails if there are any duplicated
291 * entries.
292 */
293 int zap_join(objset_t *os, uint64_t fromobj, uint64_t intoobj, dmu_tx_t *tx);
294
295 /* Same as zap_join, but set the values to 'value'. */
296 int zap_join_key(objset_t *os, uint64_t fromobj, uint64_t intoobj,
297 uint64_t value, dmu_tx_t *tx);
298
299 /* Same as zap_join, but add together any duplicated entries. */
300 int zap_join_increment(objset_t *os, uint64_t fromobj, uint64_t intoobj,
301 dmu_tx_t *tx);
302
303 /*
304 * Manipulate entries where the name + value are the "same" (the name is
305 * a stringified version of the value).
306 */

```

```

307 int zap_add_int(objset_t *os, uint64_t obj, uint64_t value, dmu_tx_t *tx);
308 int zap_remove_int(objset_t *os, uint64_t obj, uint64_t value, dmu_tx_t *tx);
309 int zap_lookup_int(objset_t *os, uint64_t obj, uint64_t value);
310 int zap_increment_int(objset_t *os, uint64_t obj, uint64_t key, int64_t delta,
311 dmu_tx_t *tx);
312
313 /* Here the key is an int and the value is a different int. */
314 int zap_add_int_key(objset_t *os, uint64_t obj,
315 uint64_t key, uint64_t value, dmu_tx_t *tx);
316 int zap_update_int_key(objset_t *os, uint64_t obj,
317 uint64_t key, uint64_t value, dmu_tx_t *tx);
318 int zap_lookup_int_key(objset_t *os, uint64_t obj,
319 uint64_t key, uint64_t *valuep);
320
321 int zap_increment(objset_t *os, uint64_t obj, const char *name, int64_t delta,
322 dmu_tx_t *tx);
323
324 struct zap;
325 struct zap_leaf;
326 typedef struct zap_cursor {
327 /* This structure is opaque! */
328 objset_t *zc_objset;
329 struct zap *zc_zap;
330 struct zap_leaf *zc_leaf;
331 uint64_t zc_zapobj;
332 uint64_t zc_serialized;
333 uint64_t zc_hash;
334 uint32_t zc_cd;
335 } zap_cursor_t;
336
337 unchanged portion omitted

```

```

*****
2121 Tue Oct 1 14:04:47 2013
new/usr/src/uts/common/fs/zfs/sys/zfeature.h
4171 clean up spa_feature_*( ) interfaces
4172 implement extensible_dataset feature for use by other zpool features
Reviewed by: Max Grossman <max.grossman@delphix.com>
Reviewed by: Christopher Siden <christopher.siden@delphix.com>
Reviewed by: George Wilson <george.wilson@delphix.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */

22 /*
23  * Copyright (c) 2013 by Delphix. All rights reserved.
24 */

26 #ifndef _SYS_ZFEATURE_H
27 #define _SYS_ZFEATURE_H

29 #include <sys/nvpair.h>
30 #include <sys/txg.h>
31 #include "zfeature_common.h"

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

37 struct spa;
38 struct dmu_tx;
39 struct objset;

40 extern boolean_t feature_is_supported(struct objset *os, uint64_t obj,
41     uint64_t desc_obj, nvlist_t *unsup_feat, nvlist_t *enabled_feat);

42 extern void spa_feature_create_zap_objects(struct spa *, struct dmu_tx *);
43 extern void spa_feature_enable(struct spa *, spa_feature_t,
44     struct dmu_tx *);
45 extern void spa_feature_enable(struct spa *, zfeature_info_t *,
46     struct dmu_tx *);
47 extern void spa_feature_incr(struct spa *, spa_feature_t, struct dmu_tx *);
48 extern void spa_feature_decr(struct spa *, spa_feature_t, struct dmu_tx *);
49 extern boolean_t spa_feature_is_enabled(struct spa *, spa_feature_t);
50 extern boolean_t spa_feature_is_active(struct spa *, spa_feature_t);
51 extern uint64_t spa_feature_refcount(spa_t *, spa_feature_t, uint64_t);
52 extern boolean_t spa_features_check(spa_t *, boolean_t, nvlist_t *, nvlist_t *);
53 extern void spa_feature_incr(struct spa *, zfeature_info_t *, struct dmu_tx *);
54 extern void spa_feature_decr(struct spa *, zfeature_info_t *, struct dmu_tx *);
55 extern boolean_t spa_feature_is_enabled(struct spa *, zfeature_info_t *);
56 extern boolean_t spa_feature_is_active(struct spa *, zfeature_info_t *);

```

```

50 extern int spa_feature_get_refcount(struct spa *, zfeature_info_t *);

51 /*
52  * These functions are only exported for zhack and zdb; normal callers should
53  * use the above interfaces.
54  */
55 extern int feature_get_refcount(struct spa *, zfeature_info_t *, uint64_t *);
56 extern void feature_enable_sync(struct spa *, zfeature_info_t *,
57     struct dmu_tx *);
58 extern void feature_sync(struct spa *, zfeature_info_t *, uint64_t,
59     struct dmu_tx *);

61 #ifdef __cplusplus
62 }
_____unchanged_portion_omitted_____

```

```

*****
34195 Tue Oct 1 14:04:48 2013
new/usr/src/uts/common/fs/zfs/zap_micro.c
4171 clean up spa_feature_*( ) interfaces
4172 implement extensible_dataset feature for use by other zpool features
Reviewed by: Max Grossman <max.grossman@delphix.com>
Reviewed by: Christopher Siden <christopher.siden@delphix.com>
Reviewed by: George Wilson <george.wilson@delphix.com>
*****
_____unchanged_portion_omitted_____

575 void
576 static void
577 mzap_create_impl(objset_t *os, uint64_t obj, int normflags, zap_flags_t flags,
578     dmu_tx_t *tx)
579 {
580     dmu_buf_t *db;
581     mzap_phys_t *zp;
582
583     VERIFY(0 == dmu_buf_hold(os, obj, 0, FTAG, &db, DMU_READ_NO_PREFETCH));
584
585 #ifdef ZFS_DEBUG
586     {
587         dmu_object_info_t doi;
588         dmu_object_info_from_db(db, &doi);
589         ASSERT3U(DMU_OT_BYTESWAP(doi.doi_type), ==, DMU_BSWAP_ZAP);
590     }
591 #endif
592
593     dmu_buf_will_dirty(db, tx);
594     zp = db->db_data;
595     zp->mz_block_type = ZBT_MICRO;
596     zp->mz_salt = ((uintptr_t)db ^ (obj << 1)) | 1ULL;
597     zp->mz_normflags = normflags;
598     dmu_buf_rele(db, FTAG);
599
600     if (flags != 0) {
601         zap_t *zap;
602         /* Only fat zap supports flags; upgrade immediately. */
603         VERIFY(0 == zap_lockdir(os, obj, tx, RW_WRITER,
604             B_FALSE, B_FALSE, &zap));
605         VERIFY3U(0, ==, mzap_upgrade(&zap, tx, flags));
606         zap_unlockdir(zap);
607     }
608 }
_____unchanged_portion_omitted_____

862 int
863 zap_contains(objset_t *os, uint64_t zapobj, const char *name)
864 {
865     int err = zap_lookup_norm(os, zapobj, name, 0,
866         0, NULL, MT_EXACT, NULL, 0, NULL);
867     int err = (zap_lookup_norm(os, zapobj, name, 0,
868         0, NULL, MT_EXACT, NULL, 0, NULL));
869     if (err == EOVERFLOW || err == EINVAL)
870         err = 0; /* found, but skipped reading the value */
871     return (err);
872 }
_____unchanged_portion_omitted_____

```

```

*****
14027 Tue Oct 1 14:04:49 2013
new/usr/src/uts/common/fs/zfs/zfeature.c
4171 clean up spa_feature_*() interfaces
4172 implement extensible_dataset feature for use by other zpool features
Reviewed by: Max Grossman <max.grossman@delphix.com>
Reviewed by: Christopher Siden <christopher.siden@delphix.com>
Reviewed by: George Wilson <george.wilson@delphix.com>
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */

22 /*
23  * Copyright (c) 2013 by Delphix. All rights reserved.
24 */

26 #include <sys/zfs_context.h>
27 #include <sys/zfeature.h>
28 #include <sys/dmu.h>
29 #include <sys/nvpair.h>
30 #include <sys/zap.h>
31 #include <sys/dmu_tx.h>
32 #include "zfeature_common.h"
33 #include <sys/spa_impl.h>

35 /*
36  * ZFS Feature Flags
37  * -----
38  *
39  * ZFS feature flags are used to provide fine-grained versioning to the ZFS
40  * on-disk format. Once enabled on a pool feature flags replace the old
41  * spa_version() number.
42  *
43  * Each new on-disk format change will be given a uniquely identifying string
44  * guid rather than a version number. This avoids the problem of different
45  * organizations creating new on-disk formats with the same version number. To
46  * keep feature guids unique they should consist of the reverse dns name of the
47  * organization which implemented the feature and a short name for the feature,
48  * separated by a colon (e.g. com.delphix:async_destroy).
49  *
50  * Reference Counts
51  * -----
52  *
53  * Within each pool features can be in one of three states: disabled, enabled,
54  * or active. These states are differentiated by a reference count stored on
55  * disk for each feature:
56  *
57  * 1) If there is no reference count stored on disk the feature is disabled.

```

```

58  * 2) If the reference count is 0 a system administrator has enabled the
59  * feature, but the feature has not been used yet, so no on-disk
60  * format changes have been made.
61  * 3) If the reference count is greater than 0 the feature is active.
62  * The format changes required by the feature are currently on disk.
63  * Note that if the feature's format changes are reversed the feature
64  * may choose to set its reference count back to 0.
65  *
66  * Feature flags makes no differentiation between non-zero reference counts
67  * for an active feature (e.g. a reference count of 1 means the same thing as a
68  * reference count of 27834721), but feature implementations may choose to use
69  * the reference count to store meaningful information. For example, a new RAID
70  * implementation might set the reference count to the number of vdevs using
71  * it. If all those disks are removed from the pool the feature goes back to
72  * having a reference count of 0.
73  *
74  * It is the responsibility of the individual features to maintain a non-zero
75  * reference count as long as the feature's format changes are present on disk.
76  *
77  * Dependencies
78  * -----
79  *
80  * Each feature may depend on other features. The only effect of this
81  * relationship is that when a feature is enabled all of its dependencies are
82  * automatically enabled as well. Any future work to support disabling of
83  * features would need to ensure that features cannot be disabled if other
84  * enabled features depend on them.
85  *
86  * On-disk Format
87  * -----
88  *
89  * When feature flags are enabled spa_version() is set to SPA_VERSION_FEATURES
90  * (5000). In order for this to work the pool is automatically upgraded to
91  * SPA_VERSION_BEFORE_FEATURES (28) first, so all pre-feature flags on disk
92  * format changes will be in use.
93  *
94  * Information about features is stored in 3 ZAP objects in the pool's MOS.
95  * These objects are linked to by the following names in the pool directory
96  * object:
97  *
98  * 1) features_for_read: feature guid -> reference count
99  * Features needed to open the pool for reading.
100 * 2) features_for_write: feature guid -> reference count
101 * Features needed to open the pool for writing.
102 * 3) feature_descriptions: feature guid -> descriptive string
103 * A human readable string.
104 *
105 * All enabled features appear in either features_for_read or
106 * features_for_write, but not both.
107 *
108 * To open a pool in read-only mode only the features listed in
109 * features_for_read need to be supported.
110 *
111 * To open the pool in read-write mode features in both features_for_read and
112 * features_for_write need to be supported.
113 *
114 * Some features may be required to read the ZAP objects containing feature
115 * information. To allow software to check for compatibility with these features
116 * before the pool is opened their names must be stored in the label in a
117 * new "features_for_read" entry (note that features that are only required
118 * to write to a pool never need to be stored in the label since the
119 * features_for_write ZAP object can be read before the pool is written to).
120 * To save space in the label features must be explicitly marked as needing to
121 * be written to the label. Also, reference counts are not stored in the label,
122 * instead any feature whose reference count drops to 0 is removed from the
123 * label.

```



```

124 *
125 * Adding New Features
126 * -----
127 *
128 * Features must be registered in zpool_feature_init() function in
129 * zfeature_common.c using the zfeature_register() function. This function
130 * has arguments to specify if the feature should be stored in the
131 * features_for_read or features_for_write ZAP object and if it needs to be
132 * written to the label when active.
133 *
134 * Once a feature is registered it will appear as a "feature@<feature name>"
135 * property which can be set by an administrator. Feature implementors should
136 * use the spa_feature_is_enabled() and spa_feature_is_active() functions to
137 * query the state of a feature and the spa_feature_incr() and
138 * spa_feature_decr() functions to change an enabled feature's reference count.
139 * Reference counts may only be updated in the syncing context.
140 *
141 * Features may not perform enable-time initialization. Instead, any such
142 * initialization should occur when the feature is first used. This design
143 * enforces that on-disk changes be made only when features are used. Code
144 * should only check if a feature is enabled using spa_feature_is_enabled(),
145 * not by relying on any feature specific metadata existing. If a feature is
146 * enabled, but the feature's metadata is not on disk yet then it should be
147 * created as needed.
148 *
149 * As an example, consider the com.delphix:async_destroy feature. This feature
150 * relies on the existence of a bptree in the MOS that store blocks for
151 * asynchronous freeing. This bptree is not created when async_destroy is
152 * enabled. Instead, when a dataset is destroyed spa_feature_is_enabled() is
153 * called to check if async_destroy is enabled. If it is and the bptree object
154 * does not exist yet, the bptree object is created as part of the dataset
155 * destroy and async_destroy's reference count is incremented to indicate it
156 * has made an on-disk format change. Later, after the destroyed dataset's
157 * blocks have all been asynchronously freed there is no longer any use for the
158 * bptree object, so it is destroyed and async_destroy's reference count is
159 * decremented back to 0 to indicate that it has undone its on-disk format
160 * changes.
161 */

163 typedef enum {
164     FEATURE_ACTION_ENABLE,
165     FEATURE_ACTION_INCR,
166     FEATURE_ACTION_DECR,
167 } feature_action_t;

168 /*
169 * Checks that the active features in the pool are supported by
170 * Checks that the features active in the specified object are supported by
171 * this software. Adds each unsupported feature (name -> description) to
172 * the supplied nvlist.
173 */
174 boolean_t
175 spa_features_check(spa_t *spa, boolean_t for_write,
176     feature_is_supported(objset_t *os, uint64_t obj, uint64_t desc_obj,
177     nvlist_t *unsup_feat, nvlist_t *enabled_feat)
178 {
179     objset_t *os = spa->spa_meta_objset;
180     boolean_t supported;
181     zap_cursor_t zc;
182     zap_attribute_t za;
183     uint64_t obj = for_write ?
184         spa->spa_feat_for_write_obj : spa->spa_feat_for_read_obj;

184     supported = B_TRUE;
185     for (zap_cursor_init(&zc, os, obj);
186         zap_cursor_retrieve(&zc, &za) == 0;

```

```

187     zap_cursor_advance(&zc)) {
188         ASSERT(za.za_integer_length == sizeof (uint64_t) &&
189             za.za_num_integers == 1);

191         if (NULL != enabled_feat) {
192             fnvlist_add_uint64(enabled_feat, za.za_name,
193                 za.za_first_integer);
194         }

196         if (za.za_first_integer != 0 &&
197             !zfeature_is_supported(za.za_name)) {
198             supported = B_FALSE;

200             if (NULL != unsup_feat) {
201                 char *desc = "";
202                 char buf[MAXPATHLEN];

204                 if (zap_lookup(os, spa->spa_feat_desc_obj,
205                     za.za_name, 1, sizeof (buf), buf) == 0)
206                     if (zap_lookup(os, desc_obj, za.za_name,
207                         1, sizeof (buf), buf) == 0)
208                         desc = buf;

208                 VERIFY(nvlist_add_string(unsup_feat, za.za_name,
209                     desc) == 0);
210             }
211         }
212     }
213     zap_cursor_fini(&zc);

215     return (supported);
216 }

218 /*
219 * Note: well-designed features will not need to use this; they should
220 * use spa_feature_is_enabled() and spa_feature_is_active() instead.
221 * However, this is non-static for zdb and zhack.
222 */
223 int
224 feature_get_refcount(spa_t *spa, zfeature_info_t *feature, uint64_t *res)
225 static int
226 feature_get_refcount(objset_t *os, uint64_t read_obj, uint64_t write_obj,
227     zfeature_info_t *feature, uint64_t *res)
228 {
229     int err;
230     uint64_t refcount;
231     uint64_t zapobj = feature->fi_can_readonly ?
232         spa->spa_feat_for_write_obj : spa->spa_feat_for_read_obj;
233     uint64_t zapobj = feature->fi_can_readonly ? write_obj : read_obj;

234     /*
235      * If the pool is currently being created, the feature objects may not
236      * have been allocated yet. Act as though all features are disabled.
237      */
238     if (zapobj == 0)
239         return (SET_ERROR(ENOTSUP));

240     err = zap_lookup(spa->spa_meta_objset, zapobj,
241         feature->fi_guid, sizeof (uint64_t), 1, &refcount);
242     err = zap_lookup(os, zapobj, feature->fi_guid, sizeof (uint64_t), 1,
243         &refcount);
244     if (err != 0) {
245         if (err == ENOENT)
246             return (SET_ERROR(ENOTSUP));
247         else
248             return (err);

```

```

245     }
246     *res = refcount;
247     return (0);
248 }

250 /*
251  * This function is non-static for zhack; it should otherwise not be used
252  * outside this file.
253  */
254 void
255 feature_sync(spa_t *spa, zfeature_info_t *feature, uint64_t refcount,
256 static int
257 feature_do_action(objset_t *os, uint64_t read_obj, uint64_t write_obj,
258 uint64_t desc_obj, zfeature_info_t *feature, feature_action_t action,
259 dmdu_tx_t *tx)
260 {
261     uint64_t zapobj = feature->fi_can_readonly ?
262     spa->spa_feat_for_write_obj : spa->spa_feat_for_read_obj;
263     int error;
264     uint64_t refcount;
265     uint64_t zapobj = feature->fi_can_readonly ? write_obj : read_obj;
266
267     VERIFY0(zap_update(spa->spa_meta_objset, zapobj, feature->fi_guid,
268     sizeof (uint64_t), 1, &refcount, tx));
269
270     if (refcount == 0)
271         spa_deactivate_mos_feature(spa, feature->fi_guid);
272     else if (feature->fi_mos)
273         spa_activate_mos_feature(spa, feature->fi_guid);
274 }
275
276 /*
277  * This function is non-static for zhack; it should otherwise not be used
278  * outside this file.
279  */
280 void
281 feature_enable_sync(spa_t *spa, zfeature_info_t *feature, dmdu_tx_t *tx)
282 {
283     uint64_t zapobj = feature->fi_can_readonly ?
284     spa->spa_feat_for_write_obj : spa->spa_feat_for_read_obj;
285
286     ASSERT(0 != zapobj);
287     ASSERT(zfeature_is_valid_guid(feature->fi_guid));
288     ASSERT3U(spa_version(spa), >=, SPA_VERSION_FEATURES);
289
290     error = zap_lookup(os, zapobj, feature->fi_guid,
291     sizeof (uint64_t), 1, &refcount);
292
293     /*
294      * If the feature is already enabled, ignore the request.
295      * If we can't ascertain the status of the specified feature, an I/O
296      * error occurred.
297      */
298     if (zap_contains(spa->spa_meta_objset, zapobj, feature->fi_guid) == 0)
299         return;
300     if (error != 0 && error != ENOENT)
301         return (error);
302
303     for (int i = 0; feature->fi_depends[i] != SPA_FEATURE_NONE; i++)
304         spa_feature_enable(spa, feature->fi_depends[i], tx);
305
306     VERIFY0(zap_update(spa->spa_meta_objset, spa->spa_feat_desc_obj,
307     feature->fi_guid, 1, strlen(feature->fi_desc) + 1,
308     feature->fi_desc, tx));
309     feature_sync(spa, feature, 0, tx);
310 }

```

```

299 static void
300 feature_do_action(spa_t *spa, spa_feature_t fid, feature_action_t action,
301 dmdu_tx_t *tx)
302 {
303     uint64_t refcount;
304     zfeature_info_t *feature = &spa_feature_table[fid];
305     uint64_t zapobj = feature->fi_can_readonly ?
306     spa->spa_feat_for_write_obj : spa->spa_feat_for_read_obj;
307
308     ASSERT3U(fid, <, SPA_FEATURES);
309     ASSERT(0 != zapobj);
310     ASSERT(zfeature_is_valid_guid(feature->fi_guid));
311
312     ASSERT(dmdu_tx_is_syncing(tx));
313     ASSERT3U(spa_version(spa), >=, SPA_VERSION_FEATURES);
314
315     VERIFY0(zap_lookup(spa->spa_meta_objset, zapobj, feature->fi_guid,
316     sizeof (uint64_t), 1, &refcount));
317
318     switch (action) {
319     case FEATURE_ACTION_ENABLE:
320         /*
321          * If the feature is already enabled, ignore the request.
322          */
323         if (error == 0)
324             return (0);
325         refcount = 0;
326         break;
327     case FEATURE_ACTION_INCR:
328         VERIFY3U(refcount, !=, UINT64_MAX);
329         if (error == ENOENT)
330             return (SET_ERROR(ENOTSUP));
331         if (refcount == UINT64_MAX)
332             return (SET_ERROR(EOVERFLOW));
333         refcount++;
334         break;
335     case FEATURE_ACTION_DECR:
336         VERIFY3U(refcount, !=, 0);
337         if (error == ENOENT)
338             return (SET_ERROR(ENOTSUP));
339         if (refcount == 0)
340             return (SET_ERROR(EOVERFLOW));
341         refcount--;
342         break;
343     default:
344         ASSERT(0);
345         break;
346     }
347
348     feature_sync(spa, feature, refcount, tx);
349     if (action == FEATURE_ACTION_ENABLE) {
350         int i;
351
352         for (i = 0; feature->fi_depends[i] != NULL; i++) {
353             zfeature_info_t *dep = feature->fi_depends[i];
354
355             error = feature_do_action(os, read_obj, write_obj,
356             desc_obj, dep, FEATURE_ACTION_ENABLE, tx);
357             if (error != 0)
358                 return (error);
359         }
360     }
361
362     error = zap_update(os, zapobj, feature->fi_guid,
363     sizeof (uint64_t), 1, &refcount, tx);

```

```

308     if (error != 0)
309         return (error);

311     if (action == FEATURE_ACTION_ENABLE) {
312         error = zap_update(os, desc_obj,
313             feature->fi_guid, 1, strlen(feature->fi_desc) + 1,
314             feature->fi_desc, tx);
315         if (error != 0)
316             return (error);
317     }

319     if (action == FEATURE_ACTION_INCR && refcount == 1 && feature->fi_mos) {
320         spa_activate_mos_feature(dmu_objset_spa(os), feature->fi_guid);
321     }

323     if (action == FEATURE_ACTION_DECR && refcount == 0) {
324         spa_deactivate_mos_feature(dmu_objset_spa(os),
325             feature->fi_guid);
326     }

328     return (0);
333 }
    unchanged_portion_omitted

356 /*
357  * Enable any required dependencies, then enable the requested feature.
358  */
359 void
360 spa_feature_enable(spa_t *spa, spa_feature_t fid, dmu_tx_t *tx)
361 spa_feature_enable(spa_t *spa, zfeature_info_t *feature, dmu_tx_t *tx)
362 {
363     ASSERT3U(spa_version(spa), >=, SPA_VERSION_FEATURES);
364     ASSERT3U(fid, <, SPA_FEATURES);
365     feature_enable_sync(spa, &spa_feature_table[fid], tx);
366     VERIFY3U(0, ==, feature_do_action(spa->spa_meta_objset,
367         spa->spa_feat_for_read_obj, spa->spa_feat_for_write_obj,
368         spa->spa_feat_desc_obj, feature, FEATURE_ACTION_ENABLE, tx));
369 }

370 void
371 spa_feature_incr(spa_t *spa, spa_feature_t fid, dmu_tx_t *tx)
372 spa_feature_incr(spa_t *spa, zfeature_info_t *feature, dmu_tx_t *tx)
373 {
374     feature_do_action(spa, fid, FEATURE_ACTION_INCR, tx);
375     ASSERT(dmu_tx_is_syncing(tx));
376     ASSERT3U(spa_version(spa), >=, SPA_VERSION_FEATURES);
377     VERIFY3U(0, ==, feature_do_action(spa->spa_meta_objset,
378         spa->spa_feat_for_read_obj, spa->spa_feat_for_write_obj,
379         spa->spa_feat_desc_obj, feature, FEATURE_ACTION_INCR, tx));
380 }

381 void
382 spa_feature_decr(spa_t *spa, spa_feature_t fid, dmu_tx_t *tx)
383 spa_feature_decr(spa_t *spa, zfeature_info_t *feature, dmu_tx_t *tx)
384 {
385     feature_do_action(spa, fid, FEATURE_ACTION_DECR, tx);
386     ASSERT(dmu_tx_is_syncing(tx));
387     ASSERT3U(spa_version(spa), >=, SPA_VERSION_FEATURES);
388     VERIFY3U(0, ==, feature_do_action(spa->spa_meta_objset,
389         spa->spa_feat_for_read_obj, spa->spa_feat_for_write_obj,
390         spa->spa_feat_desc_obj, feature, FEATURE_ACTION_DECR, tx));
391 }

392 /*
393  * This interface is for debugging only. Normal consumers should use
394  * spa_feature_is_enabled/spa_feature_is_active.

```

```

387  */
388 int
389 spa_feature_get_refcount(spa_t *spa, zfeature_info_t *feature)
390 {
391     int err;
392     uint64_t refcount;

394     if (spa_version(spa) < SPA_VERSION_FEATURES)
395         return (B_FALSE);

397     err = feature_get_refcount(spa->spa_meta_objset,
398         spa->spa_feat_for_read_obj, spa->spa_feat_for_write_obj,
399         feature, &refcount);
400     ASSERT(err == 0 || err == ENOTSUP);
401     return (err == 0 ? refcount : 0);
402 }

399 boolean_t
400 spa_feature_is_enabled(spa_t *spa, spa_feature_t fid)
401 spa_feature_is_enabled(spa_t *spa, zfeature_info_t *feature)
402 {
403     int err;
404     uint64_t refcount;

406     ASSERT3U(fid, <, SPA_FEATURES);
407     if (spa_version(spa) < SPA_VERSION_FEATURES)
408         return (B_FALSE);

410     err = feature_get_refcount(spa, &spa_feature_table[fid], &refcount);
411     err = feature_get_refcount(spa->spa_meta_objset,
412         spa->spa_feat_for_read_obj, spa->spa_feat_for_write_obj,
413         feature, &refcount);
414     ASSERT(err == 0 || err == ENOTSUP);
415     return (err == 0);
416 }

417 boolean_t
418 spa_feature_is_active(spa_t *spa, spa_feature_t fid)
419 spa_feature_is_active(spa_t *spa, zfeature_info_t *feature)
420 {
421     int err;
422     uint64_t refcount;

424     ASSERT3U(fid, <, SPA_FEATURES);
425     if (spa_version(spa) < SPA_VERSION_FEATURES)
426         return (B_FALSE);

428     err = feature_get_refcount(spa, &spa_feature_table[fid], &refcount);
429     err = feature_get_refcount(spa->spa_meta_objset,
430         spa->spa_feat_for_read_obj, spa->spa_feat_for_write_obj,
431         feature, &refcount);
432     ASSERT(err == 0 || err == ENOTSUP);
433     return (err == 0 && refcount > 0);
434 }
    unchanged_portion_omitted

```

```

*****
144359 Tue Oct 1 14:04:49 2013
new/usr/src/uts/common/fs/zfs/zfs_ioctl.c
4171 clean up spa_feature_*( ) interfaces
4172 implement extensible_dataset feature for use by other zpool features
Reviewed by: Max Grossman <max.grossman@delphix.com>
Reviewed by: Christopher Siden <christopher.siden@delphix.com>
Reviewed by: George Wilson <george.wilson@delphix.com>
*****
_unchanged_portion_omitted_

235 static int zfs_ioc_userspace_upgrade(zfs_cmd_t *zc);
236 static int zfs_check_settable(const char *name, nvpair_t *property,
237     cred_t *cr);
238 static int zfs_check_clearable(char *dataset, nvlist_t *props,
239     nvlist_t **errors);
240 static int zfs_fill_zplprops_root(uint64_t, nvlist_t *, nvlist_t *,
241     boolean_t *);
242 int zfs_set_prop_nvlist(const char *, zprop_source_t, nvlist_t *, nvlist_t *);
243 static int get_nvlist(uint64_t nvl, uint64_t size, int iflag, nvlist_t **nvp);

245 static int zfs_prop_activate_feature(spa_t *spa, spa_feature_t feature);
245 static int zfs_prop_activate_feature(spa_t *spa, zfeature_info_t *feature);

247 /* _NOTE(PRINTFLIKE(4)) - this is printf-like, but lint is too whiney */
248 void
249 _dprintf(const char *file, const char *func, int line, const char *fmt, ...)
250 {
251     const char *newfile;
252     char buf[512];
253     va_list adx;

255     /*
256      * Get rid of annoying "../common/" prefix to filename.
257      */
258     newfile = strrchr(file, '/');
259     if (newfile != NULL) {
260         newfile = newfile + 1; /* Get rid of leading / */
261     } else {
262         newfile = file;
263     }

265     va_start(adx, fmt);
266     (void) vsnprintf(buf, sizeof (buf), fmt, adx);
267     va_end(adx);

269     /*
270      * To get this data, use the zfs-dprintf probe as so:
271      * dtrace -q -n 'zfs-dprintf \
272      *     /stringof(arg0) == "dbuf.c" / \
273      *     {printf("%s: %s", stringof(arg1), stringof(arg3))}'
274      * arg0 = file name
275      * arg1 = function name
276      * arg2 = line number
277      * arg3 = message
278      */
279     DTRACE_PROBE4(zfs__dprintf,
280         char *, newfile, char *, func, int, line, char *, buf);
281 }
_unchanged_portion_omitted_

2311 /*
2312  * If the named property is one that has a special function to set its value,
2313  * return 0 on success and a positive error code on failure; otherwise if it is
2314  * not one of the special properties handled by this function, return -1.
2315  */

```

```

2316  * XXX: It would be better for callers of the property interface if we handled
2317  * these special cases in dsl_prop.c (in the dsl layer).
2318  */
2319 static int
2320 zfs_prop_set_special(const char *dsname, zprop_source_t source,
2321     nvpair_t *pair)
2322 {
2323     const char *propname = nvpair_name(pair);
2324     zfs_prop_t prop = zfs_name_to_prop(propname);
2325     uint64_t intval;
2326     int err;

2328     if (prop == ZPROP_INVALID) {
2329         if (zfs_prop_userquota(propname))
2330             return (zfs_prop_set_userquota(dsname, pair));
2331         return (-1);
2332     }

2334     if (nvpair_type(pair) == DATA_TYPE_NVLIST) {
2335         nvlist_t *attrs;
2336         VERIFY(nvpair_value_nvlist(pair, &attrs) == 0);
2337         VERIFY(nvlist_lookup_nvpair(attrs, ZPROP_VALUE,
2338             &pair) == 0);
2339     }

2341     if (zfs_prop_get_type(prop) == PROP_TYPE_STRING)
2342         return (-1);

2344     VERIFY(0 == nvpair_value_uint64(pair, &intval));

2346     switch (prop) {
2347     case ZFS_PROP_QUOTA:
2348         err = dsl_dir_set_quota(dsname, source, intval);
2349         break;
2350     case ZFS_PROP_REFQUOTA:
2351         err = dsl_dataset_set_refquota(dsname, source, intval);
2352         break;
2353     case ZFS_PROP_RESERVATION:
2354         err = dsl_dir_set_reservation(dsname, source, intval);
2355         break;
2356     case ZFS_PROP_REFRESERVATION:
2357         err = dsl_dataset_set_refreservation(dsname, source, intval);
2358         break;
2359     case ZFS_PROP_VOLSIZE:
2360         err = zvol_set_volsize(dsname, intval);
2361         break;
2362     case ZFS_PROP_VERSION:
2363     {
2364         zfsvfs_t *zfsvfs;

2366         if ((err = zfsvfs_hold(dsname, FTAG, &zfsvfs, B_TRUE)) != 0)
2367             break;

2369         err = zfs_set_version(zfsvfs, intval);
2370         zfsvfs_rele(zfsvfs, FTAG);

2372         if (err == 0 && intval >= ZPL_VERSION_USERSPACE) {
2373             zfs_cmd_t *zc;

2375             zc = kmem_zalloc(sizeof (zfs_cmd_t), KM_SLEEP);
2376             (void) strcpy(zc->zc_name, dsname);
2377             (void) zfs_ioc_userspace_upgrade(zc);
2378             kmem_free(zc, sizeof (zfs_cmd_t));
2379         }
2380         break;
2381     }

```

```

2382     case ZFS_PROP_COMPRESSION:
2383     {
2384         if (intval == ZIO_COMPRESS_LZ4) {
2385             zfeature_info_t *feature =
2386                 &spa_feature_table[SPA_FEATURE_LZ4_COMPRESS];
2385             spa_t *spa;
2387
2388             if ((err = spa_open(dsname, &spa, FTAG)) != 0)
2389                 return (err);
2390
2391             /*
2392              * Setting the LZ4 compression algorithm activates
2393              * the feature.
2394              */
2395             if (!spa_feature_is_active(spa,
2396                 SPA_FEATURE_LZ4_COMPRESS)) {
2397                 if (!spa_feature_is_active(spa, feature)) {
2398                     if ((err = zfs_prop_activate_feature(spa,
2399                         SPA_FEATURE_LZ4_COMPRESS) != 0) {
2400                         if ((err = zfs_prop_activate_feature(spa,
2401                             feature)) != 0) {
2402                             spa_close(spa, FTAG);
2403                             return (err);
2404                         }
2405                     }
2406                 }
2407             }
2408             spa_close(spa, FTAG);
2409         }
2410     }
2411
2412     /*
2413      * We still want the default set action to be performed in the
2414      * caller, we only performed zfeature settings here.
2415      */
2416     err = -1;
2417     break;
2418 }
2419
2420 default:
2421     err = -1;
2422 }
2423
2424 return (err);
2425 }

```

unchanged portion omitted

```

3583 static int
3584 zfs_check_settable(const char *dsname, nvpair_t *pair, cred_t *cr)
3585 {
3586     const char *propname = nvpair_name(pair);
3587     boolean_t issnap = (strchr(dsname, '@') != NULL);
3588     zfs_prop_t prop = zfs_name_to_prop(propname);
3589     uint64_t intval;
3590     int err;
3591
3592     if (prop == ZPROP_INVALID) {
3593         if (zfs_prop_user(propname)) {
3594             if (err = zfs_secpolicy_write_perms(dsname,
3595                 ZFS_DELEG_PERM_USERPROP, cr))
3596                 return (err);
3597             return (0);
3598         }
3599     }
3600
3601     if (!issnap && zfs_prop_userquota(propname)) {
3602         const char *perm = NULL;
3603         const char *uq_prefix =
3604             zfs_userquota_prop_prefixes[ZFS_PROP_USERQUOTA];
3605         const char *gq_prefix =
3606             zfs_userquota_prop_prefixes[ZFS_PROP_GROUPQUOTA];

```

```

3607         if (strncmp(propname, uq_prefix,
3608             strlen(uq_prefix)) == 0) {
3609             perm = ZFS_DELEG_PERM_USERQUOTA;
3610         } else if (strncmp(propname, gq_prefix,
3611             strlen(gq_prefix)) == 0) {
3612             perm = ZFS_DELEG_PERM_GROUPQUOTA;
3613         } else {
3614             /* USERUSED and GROUPUSED are read-only */
3615             return (SET_ERROR(EINVAL));
3616         }
3617
3618         if (err = zfs_secpolicy_write_perms(dsname, perm, cr))
3619             return (err);
3620         return (0);
3621     }
3622
3623     return (SET_ERROR(EINVAL));
3624 }
3625
3626 if (issnap)
3627     return (SET_ERROR(EINVAL));
3628
3629 if (nvpair_type(pair) == DATA_TYPE_NVLIST) {
3630     /*
3631      * dsl_prop_get_all_impl() returns properties in this
3632      * format.
3633      */
3634     nvlist_t *attrs;
3635     VERIFY(nvpair_value_nvlist(pair, &attrs) == 0);
3636     VERIFY(nvlist_lookup_nvpair(attrs, ZPROP_VALUE,
3637         &pair) == 0);
3638 }
3639
3640 /*
3641  * Check that this value is valid for this pool version
3642  */
3643 switch (prop) {
3644 case ZFS_PROP_COMPRESSION:
3645     /*
3646      * If the user specified gzip compression, make sure
3647      * the SPA supports it. We ignore any errors here since
3648      * we'll catch them later.
3649      */
3650     if (nvpair_type(pair) == DATA_TYPE_UINT64 &&
3651         nvpair_value_uint64(pair, &intval) == 0) {
3652         if (intval >= ZIO_COMPRESS_GZIP_1 &&
3653             intval <= ZIO_COMPRESS_GZIP_9 &&
3654             zfs_earlier_version(dsname,
3655                 SPA_VERSION_GZIP_COMPRESSION)) {
3656             return (SET_ERROR(ENOTSUP));
3657         }
3658     }
3659
3660     if (intval == ZIO_COMPRESS_ZLE &&
3661         zfs_earlier_version(dsname,
3662             SPA_VERSION_ZLE_COMPRESSION))
3663         return (SET_ERROR(ENOTSUP));
3664
3665     if (intval == ZIO_COMPRESS_LZ4) {
3666         zfeature_info_t *feature =
3667             &spa_feature_table[
3668                 SPA_FEATURE_LZ4_COMPRESS];
3669         spa_t *spa;
3670
3671         if ((err = spa_open(dsname, &spa, FTAG)) != 0)
3672             return (err);

```

```

3670         if (!spa_feature_is_enabled(spa,
3671             SPA_FEATURE_LZ4_COMPRESS)) {
3672             if (!spa_feature_is_enabled(spa, feature)) {
3673                 spa_close(spa, FTAG);
3674                 return (SET_ERROR(ENOTSUP));
3675             }
3676             spa_close(spa, FTAG);
3677         }
3678         /*
3679          * If this is a bootable dataset then
3680          * verify that the compression algorithm
3681          * is supported for booting. We must return
3682          * something other than ENOTSUP since it
3683          * implies a downrev pool version.
3684          */
3685         if (zfs_is_bootfs(dsname) &&
3686             !BOOTFS_COMPRESS_VALID(intval)) {
3687             return (SET_ERROR(ERANGE));
3688         }
3689     }
3690     break;
3691
3692     case ZFS_PROP_COPIES:
3693         if (zfs_earlier_version(dsname, SPA_VERSION_DITTO_BLOCKS))
3694             return (SET_ERROR(ENOTSUP));
3695         break;
3696
3697     case ZFS_PROP_DEDUP:
3698         if (zfs_earlier_version(dsname, SPA_VERSION_DEDUP))
3699             return (SET_ERROR(ENOTSUP));
3700         break;
3701
3702     case ZFS_PROP_SHARESMB:
3703         if (zpl_earlier_version(dsname, ZPL_VERSION_FUID))
3704             return (SET_ERROR(ENOTSUP));
3705         break;
3706
3707     case ZFS_PROP_ACLINHERIT:
3708         if (nvpair_type(pair) == DATA_TYPE_UINT64 &&
3709             nvpair_value_uint64(pair, &intval) == 0) {
3710             if (intval == ZFS_ACL_PASSTHROUGH_X &&
3711                 zfs_earlier_version(dsname,
3712                     SPA_VERSION_PASSTHROUGH_X))
3713                 return (SET_ERROR(ENOTSUP));
3714         }
3715         break;
3716     }
3717
3718     return (zfs_secpolicy_setprop(dsname, prop, pair, CRED()));
3719 }
3720
3721 /*
3722  * Checks for a race condition to make sure we don't increment a feature flag
3723  * multiple times.
3724  */
3725 static int
3726 zfs_prop_activate_feature_check(void *arg, dmu_tx_t *tx)
3727 {
3728     spa_t *spa = dmu_tx_pool(tx)->dp_spa;
3729     spa_feature_t *featurep = arg;
3730     zfeature_info_t *feature = arg;
3731
3732     if (!spa_feature_is_active(spa, *featurep))
3733         if (!spa_feature_is_active(spa, feature))

```

```

3732         return (0);
3733     else
3734         return (SET_ERROR(EBUSY));
3735 }
3736
3737 /*
3738  * The callback invoked on feature activation in the sync task caused by
3739  * zfs_prop_activate_feature.
3740  */
3741 static void
3742 zfs_prop_activate_feature_sync(void *arg, dmu_tx_t *tx)
3743 {
3744     spa_t *spa = dmu_tx_pool(tx)->dp_spa;
3745     spa_feature_t *featurep = arg;
3746     zfeature_info_t *feature = arg;
3747
3748     spa_feature_incr(spa, *featurep, tx);
3749     spa_feature_incr(spa, feature, tx);
3750 }
3751
3752 /*
3753  * Activates a feature on a pool in response to a property setting. This
3754  * creates a new sync task which modifies the pool to reflect the feature
3755  * as being active.
3756  */
3757 static int
3758 zfs_prop_activate_feature(spa_t *spa, spa_feature_t feature)
3759 {
3760     zfs_prop_activate_feature(spa, zfeature_info_t *feature)
3761     {
3762         int err;
3763
3764         /* EBUSY here indicates that the feature is already active */
3765         err = dsl_sync_task(spa_name(spa),
3766             zfs_prop_activate_feature_check, zfs_prop_activate_feature_sync,
3767             &feature, 2);
3768         if (err != 0 && err != EBUSY)
3769             return (err);
3770         else
3771             return (0);
3772     }
3773 }
3774
3775 unchanged_portion_omitted

```

```

*****
52748 Tue Oct 1 14:04:50 2013
new/usr/src/uts/common/fs/zfs/zvol.c
4171 clean up spa_feature_*( ) interfaces
4172 implement extensible_dataset feature for use by other zpool features
Reviewed by: Max Grossman <max.grossman@delphix.com>
Reviewed by: Christopher Siden <christopher.siden@delphix.com>
Reviewed by: George Wilson <george.wilson@delphix.com>
*****
_____unchanged_portion_omitted_____

1831 /*ARGSUSED*/
1832 static int
1833 zfs_mvdev_dump_feature_check(void *arg, dmu_tx_t *tx)
1834 {
1835     spa_t *spa = dmu_tx_pool(tx)->dp_spa;
1836
1837     if (spa_feature_is_active(spa, SPA_FEATURE_MULTI_VDEV_CRASH_DUMP))
1838         if (spa_feature_is_active(spa,
1839             &spa_feature_table[SPA_FEATURE_MULTI_VDEV_CRASH_DUMP]))
1840             return (1);
1841     return (0);
1842 }
1843
1844 /*ARGSUSED*/
1845 static void
1846 zfs_mvdev_dump_activate_feature_sync(void *arg, dmu_tx_t *tx)
1847 {
1848     spa_t *spa = dmu_tx_pool(tx)->dp_spa;
1849
1850     spa_feature_incr(spa, SPA_FEATURE_MULTI_VDEV_CRASH_DUMP, tx);
1851     &spa_feature_table[SPA_FEATURE_MULTI_VDEV_CRASH_DUMP], tx);
1852 }
1853
1854 static int
1855 zvol_dump_init(zvol_state_t *zv, boolean_t resize)
1856 {
1857     dmu_tx_t *tx;
1858     int error;
1859     objset_t *os = zv->zv_objset;
1860     spa_t *spa = dmu_objset_spa(os);
1861     vdev_t *vd = spa->spa_root_vdev;
1862     nvlist_t *nv = NULL;
1863     uint64_t version = spa_version(spa);
1864     enum zio_checksum checksum;
1865
1866     ASSERT(MUTEX_HELD(&zfsdev_state_lock));
1867     ASSERT(vd->vdev_ops == &vdev_root_ops);
1868
1869     error = dmu_free_long_range(zv->zv_objset, ZVOL_OBJ, 0,
1870         DMU_OBJECT_END);
1871     /* wait for dmu_free_long_range to actually free the blocks */
1872     txg_wait_synced(dmu_objset_pool(zv->zv_objset), 0);
1873
1874     /*
1875     * If the pool on which the dump device is being initialized has more
1876     * than one child vdev, check that the MULTI_VDEV_CRASH_DUMP feature is
1877     * enabled. If so, bump that feature's counter to indicate that the
1878     * feature is active. We also check the vdev type to handle the
1879     * following case:
1880     * # zpool create test raidz disk1 disk2 disk3
1881     * Now have spa_root_vdev->vdev_children == 1 (the raidz vdev),
1882     * the raidz vdev itself has 3 children.
1883     */
1884     if (vd->vdev_children > 1 || vd->vdev_ops == &vdev_raidz_ops) {

```

```

1885     if (!spa_feature_is_enabled(spa,
1886         SPA_FEATURE_MULTI_VDEV_CRASH_DUMP))
1887         &spa_feature_table[SPA_FEATURE_MULTI_VDEV_CRASH_DUMP]))
1888         return (SET_ERROR(ENOTSUP));
1889     (void) dsl_sync_task(spa_name(spa),
1890         zfs_mvdev_dump_feature_check,
1891         zfs_mvdev_dump_activate_feature_sync, NULL, 2);
1892 }
1893
1894 tx = dmu_tx_create(os);
1895 dmu_tx_hold_zap(tx, ZVOL_ZAP_OBJ, TRUE, NULL);
1896 dmu_tx_hold_bonus(tx, ZVOL_OBJ);
1897 error = dmu_tx_assign(tx, TXG_WAIT);
1898 if (error) {
1899     dmu_tx_abort(tx);
1900     return (error);
1901 }
1902
1903 /*
1904 * If MULTI_VDEV_CRASH_DUMP is active, use the NOPARITY checksum
1905 * function. Otherwise, use the old default -- OFF.
1906 */
1907 checksum = spa_feature_is_active(spa,
1908     SPA_FEATURE_MULTI_VDEV_CRASH_DUMP) ? ZIO_CHECKSUM_NOPARITY :
1909     ZIO_CHECKSUM_OFF;
1910 &spa_feature_table[SPA_FEATURE_MULTI_VDEV_CRASH_DUMP]) ?
1911     ZIO_CHECKSUM_NOPARITY : ZIO_CHECKSUM_OFF;
1912 }
1913
1914 /*
1915 * If we are resizing the dump device then we only need to
1916 * update the reservation to match the newly updated
1917 * zvolsize. Otherwise, we save off the original state of the
1918 * zvol so that we can restore them if the zvol is ever undumped.
1919 */
1920 if (resize) {
1921     error = zap_update(os, ZVOL_ZAP_OBJ,
1922         zfs_prop_to_name(ZFS_PROP_REFRESERVATION), 8, 1,
1923         &zv->zv_volsize, tx);
1924 } else {
1925     uint64_t checksum, compress, refresrv, vbs, dedup;
1926
1927     error = dsl_prop_get_integer(zv->zv_name,
1928         zfs_prop_to_name(ZFS_PROP_COMPRESSION), &compress, NULL);
1929     error = error ? error : dsl_prop_get_integer(zv->zv_name,
1930         zfs_prop_to_name(ZFS_PROP_CHECKSUM), &checksum, NULL);
1931     error = error ? error : dsl_prop_get_integer(zv->zv_name,
1932         zfs_prop_to_name(ZFS_PROP_REFRESERVATION), &refresrv, NULL);
1933     error = error ? error : dsl_prop_get_integer(zv->zv_name,
1934         zfs_prop_to_name(ZFS_PROP_VOLBLOCKSIZE), &vbs, NULL);
1935     if (version >= SPA_VERSION_DEDUP) {
1936         error = error ? error :
1937             dsl_prop_get_integer(zv->zv_name,
1938                 zfs_prop_to_name(ZFS_PROP_DEDUP), &dedup, NULL);
1939     }
1940
1941     error = error ? error : zap_update(os, ZVOL_ZAP_OBJ,
1942         zfs_prop_to_name(ZFS_PROP_COMPRESSION), 8, 1,
1943         &compress, tx);
1944     error = error ? error : zap_update(os, ZVOL_ZAP_OBJ,
1945         zfs_prop_to_name(ZFS_PROP_CHECKSUM), 8, 1, &checksum, tx);
1946     error = error ? error : zap_update(os, ZVOL_ZAP_OBJ,
1947         zfs_prop_to_name(ZFS_PROP_REFRESERVATION), 8, 1,
1948         &refresrv, tx);
1949     error = error ? error : zap_update(os, ZVOL_ZAP_OBJ,
1950         zfs_prop_to_name(ZFS_PROP_VOLBLOCKSIZE), 8, 1,
1951         &vbs, tx);

```

```
1945         error = error ? error : dmu_object_set_blocksize(
1946             os, ZVOL_OBJ, SPA_MAXBLOCKSIZE, 0, tx);
1947         if (version >= SPA_VERSION_DEDUP) {
1948             error = error ? error : zap_update(os, ZVOL_ZAP_OBJ,
1949                 zfs_prop_to_name(ZFS_PROP_DEDUP), 8, 1,
1950                 &dedup, tx);
1951         }
1952         if (error == 0)
1953             zv->zv_volblocksize = SPA_MAXBLOCKSIZE;
1954     }
1955     dmu_tx_commit(tx);
1956
1957     /*
1958     * We only need update the zvol's property if we are initializing
1959     * the dump area for the first time.
1960     */
1961     if (!resize) {
1962         VERIFY(nvlist_alloc(&nv, NV_UNIQUE_NAME, KM_SLEEP) == 0);
1963         VERIFY(nvlist_add_uint64(nv,
1964             zfs_prop_to_name(ZFS_PROP_REFRESERVATION), 0) == 0);
1965         VERIFY(nvlist_add_uint64(nv,
1966             zfs_prop_to_name(ZFS_PROP_COMPRESSION),
1967             ZIO_COMPRESS_OFF) == 0);
1968         VERIFY(nvlist_add_uint64(nv,
1969             zfs_prop_to_name(ZFS_PROP_CHECKSUM),
1970             checksum) == 0);
1971         if (version >= SPA_VERSION_DEDUP) {
1972             VERIFY(nvlist_add_uint64(nv,
1973                 zfs_prop_to_name(ZFS_PROP_DEDUP),
1974                 ZIO_CHECKSUM_OFF) == 0);
1975         }
1976
1977         error = zfs_set_prop_nvlist(zv->zv_name, ZPROP_SRC_LOCAL,
1978             nv, NULL);
1979         nvlist_free(nv);
1980
1981         if (error)
1982             return (error);
1983     }
1984
1985     /* Allocate the space for the dump */
1986     error = zvol_prealloc(zv);
1987     return (error);
1988 }
```

\_\_\_\_\_unchanged\_portion\_omitted\_\_\_\_\_