

new/usr/src/cmd/ztest/ztest.c

```
*****
160285 Thu Aug 1 22:39:34 2013
new/usr/src/cmd/ztest/ztest.c
3949 ztest fault injection should avoid resilvering devices
3950 ztest: deadman fires when we're doing a scan
3951 ztest hang when running dedup test
3952 ztest: ztest_reguid test and ztest_fault_inject don't place nice together
Reviewed by: Matthew Ahrens <mahrens@delphix.com>
Reviewed by: Adam Leventhal <ahl@delphix.com>
*****
_____unchanged_portion_omitted_____
185 extern uint64_t metaslab_gang_bang;
186 extern uint64_t metaslab_df_alloc_threshold;
187 extern uint64_t zfs_deadman_synctime;
188
189 static ztest_shared_opts_t *ztest_shared_opts;
190 static ztest_shared_opts_t ztest_opts;
191
192 typedef struct ztest_shared_ds {
193     uint64_t zd_seq;
194 } ztest_shared_ds_t;
195 _____unchanged_portion_omitted_____
196
197 static ztest_shared_callstate_t *ztest_shared_callstate;
198 #define ZTEST_GET_SHARED_CALLSTATE(c) (&ztest_shared_callstate[c])
199
200 /*
201  * Note: these aren't static because we want dladdr() to work.
202  */
203 ztest_func_t ztest_dmu_read_write;
204 ztest_func_t ztest_dmu_write_parallel;
205 ztest_func_t ztest_dmu_object_alloc_free;
206 ztest_func_t ztest_dmu_commit_callbacks;
207 ztest_func_t ztest_zap;
208 ztest_func_t ztest_zap_parallel;
209 ztest_func_t ztest_zil_commit;
210 ztest_func_t ztest_zil_remount;
211 ztest_func_t ztest_dmu_read_write_zcopy;
212 ztest_func_t ztest_dmu_objset_create_destroy;
213 ztest_func_t ztest_dmu_prealloc;
214 ztest_func_t ztest_fzap;
215 ztest_func_t ztest_dmu_snapshot_create_destroy;
216 ztest_func_t ztest_dsl_prop_get_set;
217 ztest_func_t ztest_spa_prop_get_set;
218 ztest_func_t ztest_spa_create_destroy;
219 ztest_func_t ztest_fault_inject;
220 ztest_func_t ztest_ddt_repair;
221 ztest_func_t ztest_dmu_snapshot_hold;
222 ztest_func_t ztest_spa_rename;
223 ztest_func_t ztest_scrub;
224 ztest_func_t ztest_dsl_dataset_promote_busy;
225 ztest_func_t ztest_vdev_attach_detach;
226 ztest_func_t ztest_vdev_LUN_growth;
227 ztest_func_t ztest_vdev_add_remove;
228 ztest_func_t ztest_vdev_aux_add_remove;
229 ztest_func_t ztest_dsl_dataset_promote_busy;
230 ztest_func_t ztest_vdev_attach_detach;
231 ztest_func_t ztest_vdev_LUN_growth;
232 ztest_func_t ztest_vdev_add_remove;
233 ztest_func_t ztest_vdev_aux_add_remove;
234 ztest_func_t ztest_split_pool;
235 ztest_func_t ztest_reguid;
236 ztest_func_t ztest_spa_upgrade;
237
238 uint64_t zopt_always = OULL * NANOSEC; /* all the time */
239 uint64_t zopt_incessant = 1ULL * NANOSEC / 10; /* every 1/10 second */
240 uint64_t zopt_often = 1ULL * NANOSEC; /* every second */
241 uint64_t zopt_sometimes = 10ULL * NANOSEC; /* every 10 seconds */
242 uint64_t zopt_rarely = 60ULL * NANOSEC; /* every 60 seconds */
```

1

new/usr/src/cmd/ztest/ztest.c

```
344 ztest_info_t ztest_info[] = {
345     { ztest_dmu_read_write,
346       ztest_dmu_write_parallel,
347       ztest_dmu_object_alloc_free,
348       ztest_dmu_commit_callbacks,
349       ztest_zap,
350       ztest_zap_parallel,
351       ztest_split_pool,
352       ztest_zil_commit,
353       ztest_zil_remount,
354       ztest_dmu_read_write_zcopy,
355       ztest_dmu_objset_create_destroy,
356       ztest_dsl_prop_get_set,
357       ztest_spa_prop_get_set,
358     #if 0
359     { ztest_dmu_prealloc,
360     #endif
361     { ztest_fzap,
362       ztest_dmu_snapshot_create_destroy,
363       ztest_spa_create_destroy,
364       ztest_fault_inject,
365       ztest_ddt_repair,
366       ztest_dmu_snapshot_hold,
367       ztest_reguid,
368       ztest_reguid,
369       ztest_spa_rename,
370       ztest_scrub,
371       ztest_spa_upgrade,
372       ztest_dsl_dataset_promote_busy,
373       ztest_vdev_attach_detach,
374       ztest_vdev_LUN_growth,
375       &ztest_opts.zo_vdevtime
376     { ztest_vdev_aux_add_remove,
377       &ztest_opts.zo_vdevtime
378   };
   _____unchanged_portion_omitted_____
4724 /*
4725  * Inject random faults into the on-disk data.
4726  */
4727 /* ARGSUSED */
4728 void
4729 ztest_fault_inject(ztest_ds_t *zd, uint64_t id)
4730 {
    ztest_shared_t *zs = ztest_shared;
4731     spa_t *spa = ztest_spa;
4732     int fd;
4733     uint64_t offset;
4734     uint64_t leaves;
4735     uint64_t bad = 0x1990c0ffeedecade;
4736     uint64_t top, leaf;
4737     char path0[MAXPATHLEN];
4738     char pathrand[MAXPATHLEN];
4739     size_t fsize;
4740     int bshift = SPA_MAXBLOCKSHIFT + 2; /* don't scrog all labels */
4741     int iters = 1000;
4742     int maxfaults;
4743     int mirror_save;
4744     vdev_t *vd0 = NULL;
4745     uint64_t guid0 = 0;
4746     boolean_t islog = B_FALSE;
4747
4748     VERIFY(mutex_lock(&ztest_vdev_lock) == 0);
4749     maxfaults = MAXFAULTS();
4750     leaves = MAX(zs->zs_mirrors, 1) * ztest_opts.zo_raidz;
```

2

```

4752     mirror_save = zs->zs_mirrors;
4753     VERIFY(mutex_unlock(&ztest_vdev_lock) == 0);
4755     ASSERT(leaves >= 1);
4757     /*
4758      * Grab the name lock as reader. There are some operations
4759      * which don't like to have their vdevs changed while
4760      * they are in progress (i.e. spa_change_guid). Those
4761      * operations will have grabbed the name lock as writer.
4762      */
4763     (void) rw_rdlock(&ztest_name_lock);
4765     /*
4766      * We need SCL_STATE here because we're going to look at vd0->vdev_tsd.
4767      */
4768     spa_config_enter(spa, SCL_STATE, FTAG, RW_READER);
4770     if (ztest_random(2) == 0) {
4771         /*
4772          * Inject errors on a normal data device or slog device.
4773          */
4774         top = ztest_random_vdev_top(spa, B_TRUE);
4775         leaf = ztest_random(leaves) + zs->zs_splits;
4777         /*
4778          * Generate paths to the first leaf in this top-level vdev,
4779          * and to the random leaf we selected. We'll induce transient
4780          * write failures and random online/offline activity on leaf 0,
4781          * and we'll write random garbage to the randomly chosen leaf.
4782          */
4783         (void) snprintf(path0, sizeof(path0), ztest_dev_template,
4784                         ztest_opts.zo_dir, ztest_opts.zo_pool,
4785                         top * leaves + zs->zs_splits);
4786         (void) snprintf(pathrand, sizeof(pathrand), ztest_dev_template,
4787                         ztest_opts.zo_dir, ztest_opts.zo_pool,
4788                         top * leaves + leaf);
4789         vd0 = vdev_lookup_by_path(spa->spa_root_vdev, path0);
4790         if (vd0 != NULL && vd0->vdev_top->vdev_islog)
4792             islog = B_TRUE;
4794         if (vd0 != NULL && maxfaults != 1) {
4795             /*
4796              * If the top-level vdev needs to be resilvered
4797              * then we only allow faults on the device that is
4798              * resilvering.
4799              */
4800             if (vd0 != NULL && maxfaults != 1 &&
4801                 (!vdev_resilver_needed(vd0->vdev_top, NULL, NULL) ||
4802                  vd0->vdev_resilvering)) {
4803                 /*
4804                   * Make vd0 explicitly claim to be unreadable,
4805                   * or unwriteable, or reach behind its back
4806                   * and close the underlying fd. We can do this if
4807                   * maxfaults == 0 because we'll fail and reexecute,
4808                   * and we can do it if maxfaults >= 2 because we'll
4809                   * have enough redundancy. If maxfaults == 1, the
4810                   * combination of this with injection of random data
4811                   * corruption below exceeds the pool's fault tolerance.
4812                   */
4813                 vdev_file_t *vf = vd0->vdev_tsdb;
4814                 if (vf != NULL && ztest_random(3) == 0) {
4815                     (void) close(vf->vf vnode->v_fd);
4816                     vf->vf vnode->v_fd = -1;

```

```

4817             } else if (ztest_random(2) == 0) {
4818                 vd0->vdev_cant_read = B_TRUE;
4819             } else {
4820                 vd0->vdev_cant_write = B_TRUE;
4821             }
4822             guid0 = vd0->vdev_guid;
4823         } else {
4824             /*
4825              * Inject errors on an l2cache device.
4826              */
4827             spa_aux_vdev_t *sav = &spa->spa_l2cache;
4828             if (sav->sav_count == 0) {
4829                 spa_config_exit(spa, SCL_STATE, FTAG);
4830                 (void) rw_unlock(&ztest_name_lock);
4831             }
4832             vd0 = sav->sav_vdevs[ztest_random(sav->sav_count)];
4833             guid0 = vd0->vdev_guid;
4834             (void) strcpy(path0, vd0->vdev_path);
4835             (void) strcpy(pathrand, vd0->vdev_path);
4836             leaf = 0;
4837             leaves = 1;
4838             maxfaults = INT_MAX; /* no limit on cache devices */
4839         }
4840         spa_config_exit(spa, SCL_STATE, FTAG);
4841         (void) rw_unlock(&ztest_name_lock);
4842         /*
4843          * If we can tolerate two or more faults, or we're dealing
4844          * with a slog, randomly online/offline vd0.
4845          */
4846         if ((maxfaults >= 2 || islog) && guid0 != 0) {
4847             if (ztest_random(10) < 6) {
4848                 int flags = (ztest_random(2) == 0 ?
4849                             ZFS_OFFLINE_TEMPORARY : 0);
4850                 /*
4851                   * We have to grab the zs_name_lock as writer to
4852                   * prevent a race between offlineing a slog and
4853                   * destroying a dataset. Offlineing the slog will
4854                   * grab a reference on the dataset which may cause
4855                   * dmu_objset_destroy() to fail with EBUSY thus
4856                   * leaving the dataset in an inconsistent state.
4857                   */
4858                 if (islog)
4859                     (void) rw_wrlock(&ztest_name_lock);
4860                 VERIFY(vdev_offline(spa, guid0, flags) != EBUSY);
4861                 if (islog)
4862                     (void) rw_unlock(&ztest_name_lock);
4863             } else {
4864                 /*
4865                   * Ideally we would like to be able to randomly
4866                   * call vdev_[on|off]line without holding locks
4867                   * to force unpredictable failures but the side
4868                   * effects of vdev_[on|off]line prevent us from
4869                   * doing so. We grab the ztest_vdev_lock here to
4870                   * prevent a race between injection testing and
4871                   * aux_vdev removal.
4872                   */
4873                 VERIFY(mutex_lock(&ztest_vdev_lock) == 0);

```

```

4883             (void) vdev_online(spa, guid0, 0, NULL);
4884         }
4885     }
4886 }
4887 if (maxfaults == 0)
4888     return;
4889 /*
4890 * We have at least single-fault tolerance, so inject data corruption.
4891 */
4892 fd = open(pathrand, O_RDWR);
4893
4894 if (fd == -1) /* we hit a gap in the device namespace */
4895     return;
4896
4897 fsize = lseek(fd, 0, SEEK_END);
4898
4899 while (--iters != 0) {
4900     offset = ztest_random(fsize / (leaves << bshift)) *
4901             (leaves << bshift) + (leaf << bshift) +
4902             (ztest_random(1ULL << (bshift - 1)) & -8ULL);
4903
4904     if (offset >= fsize)
4905         continue;
4906
4907     VERIFY(mutex_lock(&ztest_vdev_lock) == 0);
4908     if (mirror_save != zs->zs_mirrors) {
4909         VERIFY(mutex_unlock(&ztest_vdev_lock) == 0);
4910         (void) close(fd);
4911         return;
4912     }
4913
4914     if (pwrite(fd, &bad, sizeof(bad), offset) != sizeof(bad))
4915         fatal(1, "can't inject bad word at 0x%llx in %s",
4916               offset, pathrand);
4917
4918     VERIFY(mutex_unlock(&ztest_vdev_lock) == 0);
4919
4920     if (ztest_opts.zo_verbose >= 7)
4921         (void) printf("injected bad word into %s,"
4922                       " offset 0x%llx\n", pathrand, (u_longlong_t)offset);
4923
4924 }
4925
4926 (void) close(fd);
4927
4928 }
_____unchanged_portion_omitted
5307 static void *
5308 ztest_deadman_thread(void *arg)
5309 {
5310     ztest_shared_t *zs = arg;
5311     spa_t *spa = ztest_spa;
5312     hrtimetime_t delta, total = 0;
5313     int grace = 300;
5314     hrtimetime_t delta;
5315
5316     for (;;) {
5317         delta = (zs->zs_thread_stop - zs->zs_thread_start) /
5318             NANOSEC + zfs_deadman_synctime;
5319         delta = (zs->zs_thread_stop - zs->zs_thread_start) / NANOSEC + grace;
5320
5321         (void) poll(NULL, 0, (int)(1000 * delta));
5322
5323         /*
5324          * If the pool is suspended then fail immediately. Otherwise,

```

```

5322             * check to see if the pool is making any progress. If
5323             * vdev_deadman() discovers that there hasn't been any recent
5324             * I/Os then it will end up aborting the tests.
5325             */
5326     if (spa_suspended(spa)) {
5327         fatal(0, "aborting test after %llu seconds because "
5328               "pool has transitioned to a suspended state.", zfs_deadman_synctime);
5329     }
5330     vdev_deadman(spa->spa_root_vdev);
5331
5332     fatal(0, "failed to complete within %d seconds of deadline", grace);
5333
5334     total += zfs_deadman_synctime;
5335     (void) printf("ztest has been running for %lld seconds\n",
5336                   total);
5337 }
5338 }_____unchanged_portion_omitted
6042 int
6043 main(int argc, char **argv)
6044 {
6045     int kills = 0;
6046     int iters = 0;
6047     int older = 0;
6048     int newer = 0;
6049     ztest_shared_t *zs;
6050     ztest_info_t *zi;
6051     ztest_shared_callstate_t *zc;
6052     char timebuf[100];
6053     char numbuf[6];
6054     spa_t *spa;
6055     char *cmd;
6056     boolean_t hasalt;
6057     char *fd_data_str = getenv("ZTEST_FD_DATA");
6058
6059     (void) setvbuf(stdout, NULL, _IOLBF, 0);
6060
6061     dprintf_setup(&argc, argv);
6062     zfs_deadman_synctime = 300;
6063
6064     ztest_fd_rand = open("/dev/urandom", O_RDONLY);
6065     ASSERT3S(ztest_fd_rand, >=, 0);
6066
6067     if (!fd_data_str) {
6068         process_options(argc, argv);
6069
6070         setup_data_fd();
6071         setup_hdr();
6072         setup_data();
6073         bcopy(&ztest_opts, ztest_shared_opts,
6074               sizeof(*ztest_shared_opts));
6075     } else {
6076         ztest_fd_data = atoi(fd_data_str);
6077         setup_data();
6078         bcopy(ztest_shared_opts, &ztest_opts, sizeof(ztest_opts));
6079     }
6080     ASSERT3U(ztest_opts.zo_datasets, ==, ztest_shared_hdr->zh_ds_count);
6081
6082     /* Override location of zpool.cache */
6083     VERIFY3U(asprintf((char **)&spa_config_path, "%s/zpool.cache",
6084                       ztest_opts.zo_dir), !=, -1);
6085
6086     ztest_ds = umem_alloc(ztest_opts.zo_datasets * sizeof(ztest_ds_t),

```

```

6087     UMEM_NOFAIL);
6088     zs = ztest_shared;
6089
6090     if (fd_data_str) {
6091         metaslab_gang_bang = ztest_opts.zo_metaslab_gang_bang;
6092         metaslab_df_alloc_threshold =
6093             zs->zs_metaslab_df_alloc_threshold;
6094
6095         if (zs->zs_do_init)
6096             ztest_run_init();
6097         else
6098             ztest_run(zs);
6099         exit(0);
6100     }
6102
6103     hasalt = (strlen(ztest_opts.zo_alt_ztest) != 0);
6104
6105     if (ztest_opts.zo_verbose >= 1) {
6106         (void) printf("%llu vdevs, %d datasets, %d threads,\n",
6107                     " %llu seconds...\n",
6108                     (u_longlong_t)ztest_opts.zo_vdevs,
6109                     ztest_opts.zo_datasets,
6110                     ztest_opts.zo_threads,
6111                     (u_longlong_t)ztest_opts.zo_time);
6112     }
6113
6114     cmd = umem_alloc(MAXNAMELEN, UMEM_NOFAIL);
6115     (void) strlcpy(cmd, getexecname(), MAXNAMELEN);
6116
6117     zs->zs_do_init = B_TRUE;
6118     if (strlen(ztest_opts.zo_alt_ztest) != 0) {
6119         if (ztest_opts.zo_verbose >= 1) {
6120             (void) printf("Executing older ztest for "
6121                         "initialization: %s\n", ztest_opts.zo_alt_ztest);
6122         }
6123         VERIFY(!exec_child(ztest_opts.zo_alt_ztest,
6124                             ztest_opts.zo_alt_libpath, B_FALSE, NULL));
6125     } else {
6126         VERIFY(!exec_child(NULL, NULL, B_FALSE, NULL));
6127     }
6128     zs->zs_do_init = B_FALSE;
6129
6130     zs->zs_proc_start = gethrtime();
6131     zs->zs_proc_stop = zs->zs_proc_start + ztest_opts.zo_time * NANOSEC;
6132
6133     for (int f = 0; f < ZTEST_FUNCS; f++) {
6134         zi = &ztest_info[f];
6135         zc = ZTEST_GET_SHARED_CALLSTATE(f);
6136         if (zs->zs_proc_start + zi->zi_interval[0] > zs->zs_proc_stop)
6137             zc->zc_next = UINT64_MAX;
6138         else
6139             zc->zc_next = zs->zs_proc_start +
6140                           ztest_random(2 * zi->zi_interval[0] + 1);
6141     }
6142
6143     /*
6144      * Run the tests in a loop. These tests include fault injection
6145      * to verify that self-healing data works, and forced crashes
6146      * to verify that we never lose on-disk consistency.
6147      */
6148     while (gethrtime() < zs->zs_proc_stop) {
6149         int status;
6150         boolean_t killed;
6151
6152         /*
6153          * Initialize the workload counters for each function.
6154        */

```

```

6153
6154     */
6155     for (int f = 0; f < ZTEST_FUNCS; f++) {
6156         zc = ZTEST_GET_SHARED_CALLSTATE(f);
6157         zc->zc_count = 0;
6158         zc->zc_time = 0;
6159     }
6160
6161     /* Set the allocation switch size */
6162     zs->zs_metaslab_df_alloc_threshold =
6163         ztest_random(zs->zs_metaslab_sz / 4) + 1;
6164
6165     if (!hasalt || ztest_random(2) == 0) {
6166         if (hasalt && ztest_opts.zo_verbose >= 1) {
6167             (void) printf("Executing newer ztest: %s\n",
6168                         cmd);
6169         }
6170         newer++;
6171         killed = exec_child(cmd, NULL, B_TRUE, &status);
6172     } else {
6173         if (hasalt && ztest_opts.zo_verbose >= 1) {
6174             (void) printf("Executing older ztest: %s\n",
6175                         ztest_opts.zo_alt_ztest);
6176         }
6177         older++;
6178         killed = exec_child(ztest_opts.zo_alt_ztest,
6179                             ztest_opts.zo_alt_libpath, B_TRUE, &status);
6180     }
6181
6182     if (killed)
6183         kills++;
6184     iters++;
6185
6186     if (ztest_opts.zo_verbose >= 1) {
6187         hrtimer_t now = gethrtime();
6188
6189         now = MIN(now, zs->zs_proc_stop);
6190         print_time(zs->zs_proc_stop - now, timebuf);
6191         niceenum(zs->zs_space, numbuf);
6192
6193         (void) printf("Pass %d, %s, %llu ENOSPC, "
6194                     "%4.1f%% of %s used, %3.0f%% done, %s to go\n",
6195                     iters,
6196                     WIFEXITED(status) ? "Complete" : "SIGKILL",
6197                     (u_longlong_t)zs->zs_enospc_count,
6198                     100.0 * zs->zs_alloc / zs->zs_space,
6199                     numbuf,
6200                     100.0 * (now - zs->zs_proc_start) /
6201                     (ztest_opts.zo_time * NANOSEC), timebuf);
6202
6203     if (ztest_opts.zo_verbose >= 2) {
6204         (void) printf("\nWorkload summary:\n\n");
6205         (void) printf("%7s %9s %s\n",
6206                     "Calls", "Time", "Function");
6207         (void) printf("%7s %9s %s\n",
6208                     "----", "----", "-----");
6209         for (int f = 0; f < ZTEST_FUNCS; f++) {
6210             Dl_info dli;
6211
6212             zi = &ztest_info[f];
6213             zc = ZTEST_GET_SHARED_CALLSTATE(f);
6214             print_time(zc->zc_time, timebuf);
6215             (void) dladdr((void *)zi->zi_func, &dli);
6216             (void) printf("%7llu %9s %s\n",
6217                         (u_longlong_t)zc->zc_count, timebuf,
6218                         dli.dli_sname);
6219         }
6220     }

```

```
6219         }
6220         (void) printf("\n");
6221     }

6222     /*
6223      * It's possible that we killed a child during a rename test,
6224      * in which case we'll have a 'ztest_tmp' pool lying around
6225      * instead of 'ztest'. Do a blind rename in case this happened.
6226      */
6227     kernel_init(FREAD);
6228     if (spa_open(ztest_opts.zo_pool, &spa, FTAG) == 0) {
6229         spa_close(spa, FTAG);
6230     } else {
6231         char tmpname[MAXNAMELEN];
6232         kernel_fini();
6233         kernel_init(FREAD | FWRITE);
6234         (void) snprintf(tmpname, sizeof (tmpname), "%s_tmp",
6235                         ztest_opts.zo_pool);
6236         (void) spa_rename(tmpname, ztest_opts.zo_pool);
6237     }
6238     kernel_fini();

6241     ztest_run_zdb(ztest_opts.zo_pool);
6242 }

6244     if (ztest_opts.zo_verbose >= 1) {
6245         if (hasalt) {
6246             (void) printf("%d runs of older ztest: %s\n", older,
6247                         ztest_opts.zo_alt_ztest);
6248             (void) printf("%d runs of newer ztest: %s\n", newer,
6249                         cmd);
6250         }
6251         (void) printf("%d killed, %d completed, %.0f%% kill rate\n",
6252                     kills, iters - kills, (100.0 * kills) / MAX(1, iters));
6253     }

6255     umem_free(cmd, MAXNAMELEN);

6257     return (0);
6258 }
```

unchanged portion omitted

```
new/usr/src/lib/libzpool/common/llib-lzpool
```

```
*****
1937 Thu Aug 1 22:39:37 2013
new/usr/src/lib/libzpool/common/llib-lzpool
3949 ztest fault injection should avoid resilvering devices
3950 ztest: deadman fires when we're doing a scan
3951 ztest hang when running dedup test
3952 ztest: ztest_reguid test and ztest_fault_inject don't place nice together
Reviewed by: Matthew Ahrens <mahrens@delphix.com>
Reviewed by: Adam Leventhal <ahl@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 2010 Sun Microsystems, Inc. All rights reserved.
23 * Use is subject to license terms.
24 */

26 /*
27 * Copyright (c) 2012 by Delphix. All rights reserved.
28 */

30 /* LINTLIBRARY */
31 /* PROTOLIB1 */

33 #include <sys/zfs_context.h>
34 #include <sys/list.h>
35 #include <sys/list_impl.h>
36 #include <sys/sysmacros.h>
37 #include <sys/debug.h>
38 #include <sys/dmu_traverse.h>
39 #include <sys/dnode.h>
40 #include <sys/dsl_prop.h>
41 #include <sys/dsl_dataset.h>
42 #include <sys/spa.h>
43 #include <sys/spa_impl.h>
44 #include <sys/space_map.h>
45 #include <sys/vdev.h>
46 #include <sys/vdev_impl.h>
47 #include <sys/zap.h>
48 #include <sys/zio.h>
49 #include <sys/zio_compress.h>
50 #include <sys/zil.h>
51 #include <sys/bplist.h>
52 #include <sys/zfs_znode.h>
53 #include <sys/arc.h>
54 #include <sys/dbuf.h>
55 #include <sys/zio_checksum.h>
56 #include <sys/ddt.h>
```

```
1
```

```
new/usr/src/lib/libzpool/common/llib-lzpool
```

```
57 #include <sys/sa.h>
58 #include <sys/zfs_sa.h>
59 #include <sys/zfeature.h>
60 #include <sys/dmu_tx.h>
61 #include <sys/dsl_destroy.h>
62 #include <sys/dsl_userhold.h>
64 extern uint64_t metaslab_gang_bang;
65 extern uint64_t metaslab_df_alloc_threshold;
66 extern boolean_t zfeature_checks_disable;
67 extern uint64_t zfs_deadman_syntime;
```

```
2
```

```
*****
175636 Thu Aug 1 22:39:39 2013
new/usr/src/uts/common/fs/zfs/spa.c
3949 ztest fault injection should avoid resilvering devices
3950 ztest: deadman fires when we're doing a scan
3951 ztest hang when running dedup test
3952 ztest: ztest_reguid test and ztest_fault_inject don't place nice together
Reviewed by: Matthew Ahrens <mahrens@delphix.com>
Reviewed by: Adam Leventhal <ahl@delphix.com>
*****
unchanged_portion_omitted_

745 /*
746  * Change the GUID for the pool. This is done so that we can later
747  * re-import a pool built from a clone of our own vdevs. We will modify
748  * the root vdev's guid, our own pool guid, and then mark all of our
749  * vdevs dirty. Note that we must make sure that all our vdevs are
750  * online when we do this, or else any vdevs that weren't present
751  * would be orphaned from our pool. We are also going to issue a
752  * sysevent to update any watchers.
753 */
754 int
755 spa_change_guid(spa_t *spa)
756 {
757     int error;
758     uint64_t guid;

760     mutex_enter(&spa->spa_vdev_top_lock);
761     mutex_enter(&spa_namespace_lock);
762     guid = spa_generate_guid(NULL);

764     error = dsl_sync_task(spa->spa_name, spa_change_guid_check,
765                           spa_change_guid_sync, &guid, 5);

767     if (error == 0) {
768         spa_config_sync(spa, B_FALSE, B_TRUE);
769         spa_event_notify(spa, NULL, ESC_ZFS_POOL_REGUID);
770     }

772     mutex_exit(&spa_namespace_lock);
773     mutex_exit(&spa->spa_vdev_top_lock);
775 }
return (error);
unchanged_portion_omitted_

4525 /*
4526  * Detach a device from a mirror or replacing vdev.
4527  *
4528  * If 'replace_done' is specified, only detach if the parent
4529  * is a replacing vdev.
4530 */
4531 int
4532 spa_vdev_detach(spa_t *spa, uint64_t guid, uint64_t pguid, int replace_done)
4533 {
4534     uint64_t txg;
4535     int error;
4536     vdev_t *rvd = spa->spa_root_vdev;
4537     vdev_t *vd, *pvд, *cvd, *tvд;
4538     boolean_t unspare = B_FALSE;
4539     uint64_t unspare_guid = 0;
4540     char *vdpath;

4542     ASSERT(spa_writeable(spa));
4544     txg = spa_vdev_enter(spa);

```

```
4546     vd = spa_lookup_by_guid(spa, guid, B_FALSE);
4548     if (vd == NULL)
4549         return (spa_vdev_exit(spa, NULL, txg, ENODEV));
4551     if (!vd->vdev_ops->vdev_op_leaf)
4552         return (spa_vdev_exit(spa, NULL, txg, ENOTSUP));
4554     pvd = vd->vdev_parent;

4556     /*
4557      * If the parent/child relationship is not as expected, don't do it.
4558      * Consider M(A,R(B,C)) -- that is, a mirror of A with a replacing
4559      * vdev that's replacing B with C. The user's intent in replacing
4560      * is to go from M(A,B) to M(A,C). If the user decides to cancel
4561      * the replace by detaching C, the expected behavior is to end up
4562      * M(A,B). But suppose that right after deciding to detach C,
4563      * the replacement of B completes. We would have M(A,C), and then
4564      * ask to detach C, which would leave us with just A -- not what
4565      * the user wanted. To prevent this, we make sure that the
4566      * parent/child relationship hasn't changed -- in this example,
4567      * that C's parent is still the replacing vdev R.
4568     */
4569     if (pvд->vdev_guid != pguid && pguid != 0)
4570         return (spa_vdev_exit(spa, NULL, txg, EBUSY));

4572     /*
4573      * Only 'replacing' or 'spare' vdevs can be replaced.
4574      */
4575     if (replace_done && pvд->vdev_ops != &vdev_replacing_ops &&
4576         pvд->vdev_ops != &vdev_spare_ops)
4577         return (spa_vdev_exit(spa, NULL, txg, ENOTSUP));

4579     ASSERT(pvd->vdev_ops != &vdev_spare_ops ||
4580           spa_version(spa) >= SPA_VERSION_SPARES);

4582     /*
4583      * Only mirror, replacing, and spare vdevs support detach.
4584      */
4585     if (pvд->vdev_ops != &vdev_replacing_ops &&
4586         pvд->vdev_ops != &vdev_mirror_ops &&
4587         pvд->vdev_ops != &vdev_spare_ops)
4588         return (spa_vdev_exit(spa, NULL, txg, ENOTSUP));

4590     /*
4591      * If this device has the only valid copy of some data,
4592      * we cannot safely detach it.
4593      */
4594     if (vdev_dtl_required(vd))
4595         return (spa_vdev_exit(spa, NULL, txg, EBUSY));

4597     ASSERT(pvd->vdev_children >= 2);

4599     /*
4600      * If we are detaching the second disk from a replacing vdev, then
4601      * check to see if we changed the original vdev's path to have "/old"
4602      * at the end in spa_vdev_attach(). If so, undo that change now.
4603      */
4604     if (pvд->vdev_ops == &vdev_replacing_ops && vd->vdev_id > 0 &&
4605         vd->vdev_path != NULL) {
4606         size_t len = strlen(vd->vdev_path);

4608         for (int c = 0; c < pvd->vdev_children; c++) {
4609             cvd = pvd->vdev_child[c];

```

```

4611
4612     if (cvd == vd || cvd->vdev_path == NULL)
4613         continue;
4614
4615     if (strncmp(cvd->vdev_path, vd->vdev_path, len) == 0 &&
4616         strcmp(cvd->vdev_path + len, "/old") == 0) {
4617         spa_strfree(cvd->vdev_path);
4618         cvd->vdev_path = spa_strdup(vd->vdev_path);
4619         break;
4620     }
4621 }
4622
4623 /*
4624 * If we are detaching the original disk from a spare, then it implies
4625 * that the spare should become a real disk, and be removed from the
4626 * active spare list for the pool.
4627 */
4628 if (pdev->vdev_ops == &vdev_spare_ops &&
4629     vd->vdev_id == 0 &&
4630     pvd->vdev_child[pvd->vdev_children - 1]->vdev_isspare)
4631     unspare = B_TRUE;
4632
4633 /*
4634 * Erase the disk labels so the disk can be used for other things.
4635 * This must be done after all other error cases are handled,
4636 * but before we disembowel vd (so we can still do I/O to it).
4637 * But if we can't do it, don't treat the error as fatal --
4638 * it may be that the unwritability of the disk is the reason
4639 * it's being detached!
4640 */
4641 error = vdev_label_init(vd, 0, VDEV_LABEL_REMOVE);
4642
4643 /*
4644 * Remove vd from its parent and compact the parent's children.
4645 */
4646 vdev_remove_child(pdev, vd);
4647 vdev_compact_children(pdev);
4648
4649 /*
4650 * Remember one of the remaining children so we can get tvd below.
4651 */
4652 cvd = pvd->vdev_child[pvd->vdev_children - 1];
4653
4654 /*
4655 * If we need to remove the remaining child from the list of hot spares,
4656 * do it now, marking the vdev as no longer a spare in the process.
4657 * We must do this before vdev_remove_parent(), because that can
4658 * change the GUID if it creates a new toplevel GUID. For a similar
4659 * reason, we must remove the spare now, in the same txg as the detach;
4660 * otherwise someone could attach a new sibling, change the GUID, and
4661 * the subsequent attempt to spa_vdev_remove(unspare_guid) would fail.
4662 */
4663 if (unspare) {
4664     ASSERT(cvd->vdev_isspare);
4665     spa_spare_remove(cvd);
4666     unspare_guid = cvd->vdev_guid;
4667     (void) spa_vdev_remove(spa, unspare_guid, B_TRUE);
4668     cvd->vdev_unspare = B_TRUE;
4669 }
4670
4671 /*
4672 * If the parent mirror/replacing vdev only has one child,
4673 * the parent is no longer needed. Remove it from the tree.
4674 */
4675 if (pdev->vdev_children == 1) {
4676     if (pdev->vdev_ops == &vdev_spare_ops)

```

```

4677             cvd->vdev_unspare = B_FALSE;
4678             vdev_remove_parent(cvd);
4679             cvd->vdev_resilvering = B_FALSE;
4680         }
4681
4682         /*
4683          * We don't set tvd until now because the parent we just removed
4684          * may have been the previous top-level vdev.
4685          */
4686         tvd = cvd->vdev_top;
4687         ASSERT(tvd->vdev_parent == rvd);
4688
4689         /*
4690          * Reevaluate the parent vdev state.
4691          */
4692         vdev_propagate_state(cvd);
4693
4694         /*
4695          * If the 'autoexpand' property is set on the pool then automatically
4696          * try to expand the size of the pool. For example if the device we
4697          * just detached was smaller than the others, it may be possible to
4698          * add metaslabs (i.e. grow the pool). We need to reopen the vdev
4699          * first so that we can obtain the updated sizes of the leaf vdevs.
4700          */
4701         if (spa->spa_autoexpand) {
4702             vdev_reopen(tvd);
4703             vdev_expand(tvd, txg);
4704         }
4705
4706         vdev_config_dirty(tvd);
4707
4708         /*
4709          * Mark vd's DTL as dirty in this txg. vdev_dtl_sync() will see that
4710          * vd->vdev_detached is set and free vd's DTL object in syncing context.
4711          * But first make sure we're not on any *other* txg's DTL list, to
4712          * prevent vd from being accessed after it's freed.
4713          */
4714         vdpath = spa_strdup(vd->vdev_path);
4715         for (int t = 0; t < TXG_SIZE; t++)
4716             (void) txg_list_remove(this(&tvd->vdev_dtl_list, vd, t));
4717         vd->vdev_detached = B_TRUE;
4718         vdev_dirty(tvd, VDD_DTL, vd, txg);
4719
4720         spa_event_notify(spa, vd, ESC_ZFS_VDEV_REMOVE);
4721
4722         /* hang on to the spa before we release the lock */
4723         spa_open_ref(spa, FTAG);
4724
4725         error = spa_vdev_exit(spa, vd, txg, 0);
4726
4727         spa_history_log_internal(spa, "detach", NULL,
4728                                 "vdev=%s", vdpath);
4729         spa_strfree(vdpath);
4730
4731         /*
4732          * If this was the removal of the original device in a hot spare vdev,
4733          * then we want to go through and remove the device from the hot spare
4734          * list of every other pool.
4735          */
4736         if (unspare) {
4737             spa_t *altspa = NULL;
4738
4739             mutex_enter(&spa_namespace_lock);
4740             while ((altspa = spa_next(altspa)) != NULL) {
4741                 if (altspa->spa_state != POOL_STATE_ACTIVE ||
```

```

4742         altspa == spa)
4743             continue;
4744
4745             spa_open_ref(altspa, FTAG);
4746             mutex_exit(&spa_namespace_lock);
4747             (void) spa_vdev_remove(altspa, unspare_guid, B_TRUE);
4748             mutex_enter(&spa_namespace_lock);
4749             spa_close(altspa, FTAG);
4750
4751             mutex_exit(&spa_namespace_lock);
4752
4753             /* search the rest of the vdevs for spares to remove */
4754             spa_vdev_resilver_done(spa);
4755
4756         } /* all done with the spa; OK to release */
4757         mutex_enter(&spa_namespace_lock);
4758         spa_close(spa, FTAG);
4759         mutex_exit(&spa_namespace_lock);
4760
4761     return (error);
4762 }
unchanged_portion_omitted

5291 /*
5292  * Find any device that's done replacing, or a vdev marked 'unspare' that's
5293  * currently spared, so we can detach it.
5294  */
5295 static vdev_t *
5296 spa_vdev_resilver_done_hunt(vdev_t *vd)
5297 {
5298     vdev_t *newvd, *oldvd;
5299
5300     for (int c = 0; c < vd->vdev_children; c++) {
5301         oldvd = spa_vdev_resilver_done_hunt(vd->vdev_child[c]);
5302         if (oldvd != NULL)
5303             return (oldvd);
5304     }
5305
5306     if (vd->vdev_resilvering && vdev_dtl_empty(vd, DTL_MISSING) &&
5307         vdev_dtl_empty(vd, DTL_OUTAGE)) {
5308         ASSERT(vd->vdev_ops->vdev_op_leaf);
5309         vd->vdev_resilvering = B_FALSE;
5310         vdev_config_dirty(vd->vdev_top);
5311     }
5312
5313     /*
5314      * Check for a completed replacement. We always consider the first
5315      * vdev in the list to be the oldest vdev, and the last one to be
5316      * the newest (see spa_vdev_attach() for how that works). In
5317      * the case where the newest vdev is faulted, we will not automatically
5318      * remove it after a resilver completes. This is OK as it will require
5319      * user intervention to determine which disk the admin wishes to keep.
5320      */
5321     if (vd->vdev_ops == &vdev_replacing_ops) {
5322         ASSERT(vd->vdev_children > 1);
5323
5324         newvd = vd->vdev_child[vd->vdev_children - 1];
5325         oldvd = vd->vdev_child[0];
5326
5327         if (vdev_dtl_empty(newvd, DTL_MISSING) &&
5328             vdev_dtl_empty(newvd, DTL_OUTAGE) &&
5329             !vdev_dtl_required(oldvd))
5330             return (oldvd);
5331     }

```

```

5333         /*
5334          * Check for a completed resilver with the 'unspare' flag set.
5335          */
5336         if (vd->vdev_ops == &vdev_spare_ops) {
5337             vdev_t *first = vd->vdev_child[0];
5338             vdev_t *last = vd->vdev_child[vd->vdev_children - 1];
5339
5340             if (last->vdev_unspare) {
5341                 oldvd = first;
5342                 newvd = last;
5343             } else if (first->vdev_unspare) {
5344                 oldvd = last;
5345                 newvd = first;
5346             } else {
5347                 oldvd = NULL;
5348             }
5349
5350             if (oldvd != NULL &&
5351                 vdev_dtl_empty(newvd, DTL_MISSING) &&
5352                 vdev_dtl_empty(newvd, DTL_OUTAGE) &&
5353                 !vdev_dtl_required(oldvd))
5354                 return (oldvd);
5355
5356             /*
5357              * If there are more than two spares attached to a disk,
5358              * and those spares are not required, then we want to
5359              * attempt to free them up now so that they can be used
5360              * by other pools. Once we're back down to a single
5361              * disk+spare, we stop removing them.
5362              */
5363             if (vd->vdev_children > 2) {
5364                 newvd = vd->vdev_child[1];
5365
5366                 if (newvd->vdev_issspare && last->vdev_issspare &&
5367                     vdev_dtl_empty(last, DTL_MISSING) &&
5368                     vdev_dtl_empty(last, DTL_OUTAGE) &&
5369                     !vdev_dtl_required(newvd))
5370                     return (newvd);
5371             }
5372         }
5373
5374     return (NULL);
5375 }
unchanged_portion_omitted

```

```
*****
6126 Thu Aug 1 22:39:42 2013
new/usr/src/uts/common/fs/zfs/vdev_file.c
3949 ztest fault injection should avoid resilvering devices
3950 ztest: deadman fires when we're doing a scan
3951 ztest hang when running dedup test
3952 ztest: ztest_reguid test and ztest_fault_inject don't place nice together
Reviewed by: Matthew Ahrens <mahrens@delphix.com>
Reviewed by: Adam Leventhal <ahl@delphix.com>
*****
_____unchanged_portion_omitted_____
185 static int
186 vdev_file_io_start(zio_t *zio)
187 {
188     spa_t *spa = zio->io_spa;
189     vdev_t *vd = zio->io_vd;
190     vdev_file_t *vf = vd->vdev_tsds;
191     vdev_buf_t *vb;
192     buf_t *bp;
193
194     if (zio->io_type == ZIO_TYPE_IOCTL) {
195         /* XXPOLICY */
196         if (!vdev_readable(vd)) {
197             zio->io_error = SET_ERROR(ENXIO);
198             return (ZIO_PIPELINE_CONTINUE);
199
200         switch (zio->io_cmd) {
201             case DKIOCFLUSHWRITECACHE:
202                 zio->io_error = VOP_FSYNC(vf->vf_vnode, FSYNC | FDSYNC,
203                                         kcred, NULL);
204                 break;
205             default:
206                 zio->io_error = SET_ERROR(ENOTSUP);
207             }
208
209         return (ZIO_PIPELINE_CONTINUE);
210     }
211
212     vb = kmem_alloc(sizeof(vdev_buf_t), KM_SLEEP);
213
214     vb->vb_io = zio;
215     bp = &vb->vb_buf;
216
217     bioinit(bp);
218     bp->b_flags = (zio->io_type == ZIO_TYPE_READ ? B_READ : B_WRITE);
219     bp->b_bcount = zio->io_size;
220     bp->b_un.b_addr = zio->io_data;
221     bp->b_lblkno = lbtodb(zio->io_offset);
222     bp->b_bufsize = zio->io_size;
223     bp->b_private = vf->vf_vnode;
224     bp->b_iodone = (int (*)())vdev_file_io_intr;
225
226     VERIFY3U(taskq_dispatch(system_taskq, vdev_file_io_strategy, bp,
227                             TQ_SLEEP), !=, 0);
228     spa_taskq_dispatch_ent(spa, ZIO_TYPE_FREE, ZIO_TASKQ_ISSUE,
229                            vdev_file_io_strategy, bp, 0, &zio->io_tqent);
230
231     return (ZIO_PIPELINE_STOP);
232 }
_____unchanged_portion_omitted_____

```