

```

*****
12122 Mon Jun 17 22:40:14 2013
new/usr/src/lib/libzfs/common/libzfs_status.c
3818 zpool status -x should report pools with removed l2arc devices
Reviewed by: Saso Kiselkov <skiselkov.ml@gmail.com>
Reviewed by: George Wilson <gwilson@zfsmail.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) 2012 by Delphix. All rights reserved.
25  * Copyright (c) 2013 Steven Hartland. All rights reserved.
26 #endif /* ! codereview */
27 */

29 /*
30  * This file contains the functions which analyze the status of a pool. This
31  * include both the status of an active pool, as well as the status exported
32  * pools. Returns one of the ZPOOL_STATUS_* defines describing the status of
33  * the pool. This status is independent (to a certain degree) from the state of
34  * the pool. A pool's state describes only whether or not it is capable of
35  * providing the necessary fault tolerance for data. The status describes the
36  * overall status of devices. A pool that is online can still have a device
37  * that is experiencing errors.
38  *
39  * Only a subset of the possible faults can be detected using 'zpool status',
40  * and not all possible errors correspond to a FMA message ID. The explanation
41  * is left up to the caller, depending on whether it is a live pool or an
42  * import.
43  */

45 #include <libzfs.h>
46 #include <string.h>
47 #include <unistd.h>
48 #include "libzfs_impl.h"
49 #include "zfeature_common.h"

51 /*
52  * Message ID table. This must be kept in sync with the ZPOOL_STATUS_* defines
53  * in libzfs.h. Note that there are some status results which go past the end
54  * of this table, and hence have no associated message ID.
55  */
56 static char *zfs_msgid_table[] = {
57     "ZFS-8000-14",
58     "ZFS-8000-20",
59     "ZFS-8000-3C",

```

```

60     "ZFS-8000-4J",
61     "ZFS-8000-5E",
62     "ZFS-8000-6X",
63     "ZFS-8000-72",
64     "ZFS-8000-8A",
65     "ZFS-8000-9P",
66     "ZFS-8000-A5",
67     "ZFS-8000-EV",
68     "ZFS-8000-HC",
69     "ZFS-8000-JQ",
70     "ZFS-8000-K4",
71 };

73 #define NMSGID (sizeof (zfs_msgid_table) / sizeof (zfs_msgid_table[0]))

75 /* ARGSUSED */
76 static int
77 vdev_missing(uint64_t state, uint64_t aux, uint64_t errs)
78 {
79     return (state == VDEV_STATE_CANT_OPEN &&
80         aux == VDEV_AUX_OPEN_FAILED);
81 }

83 /* ARGSUSED */
84 static int
85 vdev_faulted(uint64_t state, uint64_t aux, uint64_t errs)
86 {
87     return (state == VDEV_STATE_FAULTED);
88 }

90 /* ARGSUSED */
91 static int
92 vdev_errors(uint64_t state, uint64_t aux, uint64_t errs)
93 {
94     return (state == VDEV_STATE_DEGRADED || errs != 0);
95 }

97 /* ARGSUSED */
98 static int
99 vdev_broken(uint64_t state, uint64_t aux, uint64_t errs)
100 {
101     return (state == VDEV_STATE_CANT_OPEN);
102 }

104 /* ARGSUSED */
105 static int
106 vdev_offlined(uint64_t state, uint64_t aux, uint64_t errs)
107 {
108     return (state == VDEV_STATE_OFFLINE);
109 }

111 /* ARGSUSED */
112 static int
113 vdev_removed(uint64_t state, uint64_t aux, uint64_t errs)
114 {
115     return (state == VDEV_STATE_REMOVED);
116 }

118 /*
119  * Detect if any leaf devices that have seen errors or could not be opened.
120  */
121 static boolean_t
122 find_vdev_problem(nvlist_t *vdev, int (*func)(uint64_t, uint64_t, uint64_t))
123 {
124     nvlist_t **child;
125     vdev_stat_t *vs;

```

```

126     uint_t c, children;
127     char *type;

129     /*
130      * Ignore problems within a 'replacing' vdev, since we're presumably in
131      * the process of repairing any such errors, and don't want to call them
132      * out again. We'll pick up the fact that a resilver is happening
133      * later.
134      */
135     verify(nvlist_lookup_string(vdev, ZPOOL_CONFIG_TYPE, &type) == 0);
136     if (strcmp(type, VDEV_TYPE_REPLACING) == 0)
137         return (B_FALSE);

139     if (nvlist_lookup_nvlist_array(vdev, ZPOOL_CONFIG_CHILDREN, &child,
140     &children) == 0) {
141         for (c = 0; c < children; c++)
142             if (find_vdev_problem(child[c], func))
143                 return (B_TRUE);
144     } else {
145         verify(nvlist_lookup_uint64_array(vdev, ZPOOL_CONFIG_VDEV_STATS,
146         (uint64_t **)&vs, &c) == 0);

148         if (func(vs->vs_state, vs->vs_aux,
149         vs->vs_read_errors +
150         vs->vs_write_errors +
151         vs->vs_checksum_errors))
152             return (B_TRUE);
153     }

155     /*
156      * Check any L2 cache devs
157      */
158     if (nvlist_lookup_nvlist_array(vdev, ZPOOL_CONFIG_L2CACHE, &child,
159     &children) == 0) {
160         for (c = 0; c < children; c++)
161             if (find_vdev_problem(child[c], func))
162                 return (B_TRUE);
163     }

165 #endif /* ! codereview */
166     return (B_FALSE);
167 }

169 /*
170  * Active pool health status.
171  *
172  * To determine the status for a pool, we make several passes over the config,
173  * picking the most egregious error we find. In order of importance, we do the
174  * following:
175  *
176  * - Check for a complete and valid configuration
177  * - Look for any faulted or missing devices in a non-replicated config
178  * - Check for any data errors
179  * - Check for any faulted or missing devices in a replicated config
180  * - Look for any devices showing errors
181  * - Check for any resilvering devices
182  *
183  * There can obviously be multiple errors within a single pool, so this routine
184  * only picks the most damaging of all the current errors to report.
185  */
186 static zpool_status_t
187 check_status(nvlist_t *config, boolean_t isimport)
188 {
189     nvlist_t *nvroot;
190     vdev_stat_t *vs;
191     pool_scan_stat_t *ps = NULL;

```

```

192     uint_t vsc, psc;
193     uint64_t nerr;
194     uint64_t version;
195     uint64_t stateval;
196     uint64_t suspended;
197     uint64_t hostid = 0;

199     verify(nvlist_lookup_uint64(config, ZPOOL_CONFIG_VERSION,
200     &version) == 0);
201     verify(nvlist_lookup_nvlist(config, ZPOOL_CONFIG_VDEV_TREE,
202     &nvroot) == 0);
203     verify(nvlist_lookup_uint64_array(nvroot, ZPOOL_CONFIG_VDEV_STATS,
204     (uint64_t **)&vs, &vsc) == 0);
205     verify(nvlist_lookup_uint64(config, ZPOOL_CONFIG_POOL_STATE,
206     &stateval) == 0);

208     /*
209      * Currently resilvering a vdev
210      */
211     (void) nvlist_lookup_uint64_array(nvroot, ZPOOL_CONFIG_SCAN_STATS,
212     (uint64_t **)&ps, &psc);
213     if (ps && ps->pss_func == POOL_SCAN_RESILVER &&
214     ps->pss_state == DSS_SCANNING)
215         return (ZPOOL_STATUS_RESILVERING);

217     /*
218      * Pool last accessed by another system.
219      */
220     (void) nvlist_lookup_uint64(config, ZPOOL_CONFIG_HOSTID, &hostid);
221     if (hostid != 0 && (unsigned long)hostid != gethostid() &&
222     stateval == POOL_STATE_ACTIVE)
223         return (ZPOOL_STATUS_HOSTID_MISMATCH);

225     /*
226      * Newer on-disk version.
227      */
228     if (vs->vs_state == VDEV_STATE_CANT_OPEN &&
229     vs->vs_aux == VDEV_AUX_VERSION_NEWER)
230         return (ZPOOL_STATUS_VERSION_NEWER);

232     /*
233      * Unsupported feature(s).
234      */
235     if (vs->vs_state == VDEV_STATE_CANT_OPEN &&
236     vs->vs_aux == VDEV_AUX_UNSUP_FEAT) {
237         nvlist_t *nvinfo;

239         verify(nvlist_lookup_nvlist(config, ZPOOL_CONFIG_LOAD_INFO,
240         &nvinfo) == 0);
241         if (nvlist_exists(nvinfo, ZPOOL_CONFIG_CAN_RDONLY))
242             return (ZPOOL_STATUS_UNSUP_FEAT_WRITE);
243         return (ZPOOL_STATUS_UNSUP_FEAT_READ);
244     }

246     /*
247      * Check that the config is complete.
248      */
249     if (vs->vs_state == VDEV_STATE_CANT_OPEN &&
250     vs->vs_aux == VDEV_AUX_BAD_GUID_SUM)
251         return (ZPOOL_STATUS_BAD_GUID_SUM);

253     /*
254      * Check whether the pool has suspended due to failed I/O.
255      */
256     if (nvlist_lookup_uint64(config, ZPOOL_CONFIG_SUSPENDED,
257     &suspended) == 0) {

```

```

258     if (suspended == ZIO_FAILURE_MODE_CONTINUE)
259         return (ZPOOL_STATUS_IO_FAILURE_CONTINUE);
260     return (ZPOOL_STATUS_IO_FAILURE_WAIT);
261 }
262
263 /*
264  * Could not read a log.
265  */
266 if (vs->vs_state == VDEV_STATE_CANT_OPEN &&
267     vs->vs_aux == VDEV_AUX_BAD_LOG) {
268     return (ZPOOL_STATUS_BAD_LOG);
269 }
270
271 /*
272  * Bad devices in non-replicated config.
273  */
274 if (vs->vs_state == VDEV_STATE_CANT_OPEN &&
275     find_vdev_problem(nvroot, vdev_faulted))
276     return (ZPOOL_STATUS_FAULTED_DEV_NR);
277
278 if (vs->vs_state == VDEV_STATE_CANT_OPEN &&
279     find_vdev_problem(nvroot, vdev_missing))
280     return (ZPOOL_STATUS_MISSING_DEV_NR);
281
282 if (vs->vs_state == VDEV_STATE_CANT_OPEN &&
283     find_vdev_problem(nvroot, vdev_broken))
284     return (ZPOOL_STATUS_CORRUPT_LABEL_NR);
285
286 /*
287  * Corrupted pool metadata
288  */
289 if (vs->vs_state == VDEV_STATE_CANT_OPEN &&
290     vs->vs_aux == VDEV_AUX_CORRUPT_DATA)
291     return (ZPOOL_STATUS_CORRUPT_POOL);
292
293 /*
294  * Persistent data errors.
295  */
296 if (!isimport) {
297     if (nvlist_lookup_uint64(config, ZPOOL_CONFIG_ERRCOUNT,
298         &nerr) == 0 && nerr != 0)
299         return (ZPOOL_STATUS_CORRUPT_DATA);
300 }
301
302 /*
303  * Missing devices in a replicated config.
304  */
305 if (find_vdev_problem(nvroot, vdev_faulted))
306     return (ZPOOL_STATUS_FAULTED_DEV_R);
307 if (find_vdev_problem(nvroot, vdev_missing))
308     return (ZPOOL_STATUS_MISSING_DEV_R);
309 if (find_vdev_problem(nvroot, vdev_broken))
310     return (ZPOOL_STATUS_CORRUPT_LABEL_R);
311
312 /*
313  * Devices with errors
314  */
315 if (!isimport && find_vdev_problem(nvroot, vdev_errors))
316     return (ZPOOL_STATUS_FAILING_DEV);
317
318 /*
319  * Offlined devices
320  */
321 if (find_vdev_problem(nvroot, vdev_offlined))
322     return (ZPOOL_STATUS_OFFLINE_DEV);

```

```

324 /*
325  * Removed device
326  */
327 if (find_vdev_problem(nvroot, vdev_removed))
328     return (ZPOOL_STATUS_REMOVED_DEV);
329
330 /*
331  * Outdated, but usable, version
332  */
333 if (SPA_VERSION_IS_SUPPORTED(version) && version != SPA_VERSION)
334     return (ZPOOL_STATUS_VERSION_OLDER);
335
336 /*
337  * Usable pool with disabled features
338  */
339 if (version >= SPA_VERSION_FEATURES) {
340     int i;
341     nvlist_t *feat;
342
343     if (isimport) {
344         feat = fnvlist_lookup_nvlist(config,
345             ZPOOL_CONFIG_LOAD_INFO);
346         feat = fnvlist_lookup_nvlist(feat,
347             ZPOOL_CONFIG_ENABLED_FEAT);
348     } else {
349         feat = fnvlist_lookup_nvlist(config,
350             ZPOOL_CONFIG_FEATURE_STATS);
351     }
352
353     for (i = 0; i < SPA_FEATURES; i++) {
354         zfeature_info_t *fi = &spa_feature_table[i];
355         if (!nvlist_exists(feat, fi->fi_guid))
356             return (ZPOOL_STATUS_FEAT_DISABLED);
357     }
358 }
359
360 return (ZPOOL_STATUS_OK);
361 }
362
363 zpool_status_t
364 zpool_get_status(zpool_handle_t *zhp, char **msgid)
365 {
366     zpool_status_t ret = check_status(zhp->zpool_config, B_FALSE);
367
368     if (ret >= NMSGID)
369         *msgid = NULL;
370     else
371         *msgid = zfs_msgid_table[ret];
372
373     return (ret);
374 }
375
376 zpool_status_t
377 zpool_import_status(nvlist_t *config, char **msgid)
378 {
379     zpool_status_t ret = check_status(config, B_TRUE);
380
381     if (ret >= NMSGID)
382         *msgid = NULL;
383     else
384         *msgid = zfs_msgid_table[ret];
385
386     return (ret);
387 }
388
389 static void

```

```

390 dump_ddt_stat(const ddt_stat_t *dds, int h)
391 {
392     char refcnt[6];
393     char blocks[6], lsize[6], psize[6], dsize[6];
394     char ref_blocks[6], ref_lsize[6], ref_psize[6], ref_dsize[6];
395
396     if (dds == NULL || dds->dds_blocks == 0)
397         return;
398
399     if (h == -1)
400         (void) strcpy(refcnt, "Total");
401     else
402         zfs_nicenum(1ULL << h, refcnt, sizeof (refcnt));
403
404     zfs_nicenum(dds->dds_blocks, blocks, sizeof (blocks));
405     zfs_nicenum(dds->dds_lsize, lsize, sizeof (lsize));
406     zfs_nicenum(dds->dds_psize, psize, sizeof (psize));
407     zfs_nicenum(dds->dds_dsize, dsize, sizeof (dsize));
408     zfs_nicenum(dds->dds_ref_blocks, ref_blocks, sizeof (ref_blocks));
409     zfs_nicenum(dds->dds_ref_lsize, ref_lsize, sizeof (ref_lsize));
410     zfs_nicenum(dds->dds_ref_psize, ref_psize, sizeof (ref_psize));
411     zfs_nicenum(dds->dds_ref_dsize, ref_dsize, sizeof (ref_dsize));
412
413     (void) printf("%6s %6s %5s %5s %5s %6s %5s %5s %5s\n",
414                 refcnt,
415                 blocks, lsize, psize, dsize,
416                 ref_blocks, ref_lsize, ref_psize, ref_dsize);
417 }
418
419 /*
420  * Print the DDT histogram and the column totals.
421  */
422 void
423 zpool_dump_ddt(const ddt_stat_t *dds_total, const ddt_histogram_t *ddh)
424 {
425     int h;
426
427     (void) printf("\n");
428
429     (void) printf("bucket  "
430                 "      allocated      "
431                 "      referenced      \n");
432     (void) printf("-----"
433                 "-----"
434                 "-----\n");
435
436     (void) printf("%6s %6s %5s %5s %5s %6s %5s %5s %5s\n",
437                 "refcnt",
438                 "blocks", "LSIZE", "PSIZE", "DSIZE",
439                 "blocks", "LSIZE", "PSIZE", "DSIZE");
440
441     (void) printf("%6s %6s %5s %5s %5s %6s %5s %5s %5s\n",
442                 "-----",
443                 "-----", "-----", "-----",
444                 "-----", "-----", "-----", "-----");
445
446     for (h = 0; h < 64; h++)
447         dump_ddt_stat(&ddh->ddh_stat[h], h);
448
449     dump_ddt_stat(dds_total, -1);
450
451     (void) printf("\n");
452 }

```