

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

1

```
*****
203650 Mon Aug 5 21:59:08 2013
new/usr/src/uts/common/fs/zfs/arc.c
3525 Persistent L2ARC
*****
1 /*
2  * CDDL HEADER START
3  *
4  * The contents of this file are subject to the terms of the
5  * Common Development and Distribution License (the "License").
6  * You may not use this file except in compliance with the License.
7  *
8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9  * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22 * Copyright (c) 2005, 2010, Oracle and/or its affiliates. All rights reserved.
23 * Copyright 2011 Nexenta Systems, Inc. All rights reserved.
24 * Copyright (c) 2013 by Delphix. All rights reserved.
25 * Copyright (c) 2013 by Saso Kiselkov. All rights reserved.
26 */
27
28 /*
29 * DVA-based Adjustable Replacement Cache
30 *
31 * While much of the theory of operation used here is
32 * based on the self-tuning, low overhead replacement cache
33 * presented by Megiddo and Modha at FAST 2003, there are some
34 * significant differences:
35 *
36 * 1. The Megiddo and Modha model assumes any page is evictable.
37 * Pages in its cache cannot be "locked" into memory. This makes
38 * the eviction algorithm simple: evict the last page in the list.
39 * This also make the performance characteristics easy to reason
40 * about. Our cache is not so simple. At any given moment, some
41 * subset of the blocks in the cache are un-evictable because we
42 * have handed out a reference to them. Blocks are only evictable
43 * when there are no external references active. This makes
44 * eviction far more problematic: we choose to evict the evictable
45 * blocks that are the "lowest" in the list.
46 *
47 * There are times when it is not possible to evict the requested
48 * space. In these circumstances we are unable to adjust the cache
49 * size. To prevent the cache growing unbounded at these times we
50 * implement a "cache throttle" that slows the flow of new data
51 * into the cache until we can make space available.
52 *
53 * 2. The Megiddo and Modha model assumes a fixed cache size.
54 * Pages are evicted when the cache is full and there is a cache
55 * miss. Our model has a variable sized cache. It grows with
56 * high use, but also tries to react to memory pressure from the
57 * operating system: decreasing its size when system memory is
58 * tight.
59 *
60 * 3. The Megiddo and Modha model assumes a fixed page size. All
61 * elements of the cache are therefore exactly the same size. So
```

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

2

```
62 * when adjusting the cache size following a cache miss, its simply
63 * a matter of choosing a single page to evict. In our model, we
64 * have variable sized cache blocks (ranging from 512 bytes to
65 * 128K bytes). We therefore choose a set of blocks to evict to make
66 * space for a cache miss that approximates as closely as possible
67 * the space used by the new block.
68 *
69 * See also: "ARC: A Self-Tuning, Low Overhead Replacement Cache"
70 * by N. Megiddo & D. Modha, FAST 2003
71 */
72
73 /*
74 * The locking model:
75 *
76 * A new reference to a cache buffer can be obtained in two
77 * ways: 1) via a hash table lookup using the DVA as a key,
78 * or 2) via one of the ARC lists. The arc_read() interface
79 * uses method 1, while the internal arc algorithms for
80 * adjusting the cache use method 2. We therefore provide two
81 * types of locks: 1) the hash table lock array, and 2) the
82 * arc list locks.
83 *
84 * Buffers do not have their own mutexes, rather they rely on the
85 * hash table mutexes for the bulk of their protection (i.e. most
86 * fields in the arc_buf_hdr_t are protected by these mutexes).
87 *
88 * buf_hash_find() returns the appropriate mutex (held) when it
89 * locates the requested buffer in the hash table. It returns
90 * NULL for the mutex if the buffer was not in the table.
91 *
92 * buf_hash_remove() expects the appropriate hash mutex to be
93 * already held before it is invoked.
94 *
95 * Each arc state also has a mutex which is used to protect the
96 * buffer list associated with the state. When attempting to
97 * obtain a hash table lock while holding an arc list lock you
98 * must use: mutex_tryenter() to avoid deadlock. Also note that
99 * the active state mutex must be held before the ghost state mutex.
100 *
101 * Arc buffers may have an associated eviction callback function.
102 * This function will be invoked prior to removing the buffer (e.g.
103 * in arc_do_user_evicts()). Note however that the data associated
104 * with the buffer may be evicted prior to the callback. The callback
105 * must be made with *no locks held* (to prevent deadlock). Additionally,
106 * the users of callbacks must ensure that their private data is
107 * protected from simultaneous callbacks from arc_buf_evict()
108 * and arc_do_user_evicts().
109 *
110 * Note that the majority of the performance stats are manipulated
111 * with atomic operations.
112 *
113 * The L2ARC uses the l2arc_buflist_mtx global mutex for the following:
114 *
115 * - L2ARC buflist creation
116 * - L2ARC buflist eviction
117 * - L2ARC write completion, which walks L2ARC buflists
118 * - ARC header destruction, as it removes from L2ARC buflists
119 * - ARC header release, as it removes from L2ARC buflists
120 */
121
122 #include <sys/spa.h>
123 #include <sys/zio.h>
124 #include <sys/zio_compress.h>
125 #include <sys/zfs_context.h>
126 #include <sys/arc.h>
127 #include <sys/refcount.h>
```

```

128 #include <sys/vdev.h>
129 #include <sys/vdev_impl.h>
130 #ifdef _KERNEL
131 #include <sys/vmsystem.h>
132 #include <vm/anon.h>
133 #include <sys/fs/swapnode.h>
134 #include <sys/dnlic.h>
135 #endif
136 #include <sys/callb.h>
137 #include <sys/kstat.h>
138 #include <zfs_fletcher.h>
139 #include <sys/byteorder.h>

141 #ifndef _KERNEL
142 /* set with ZFS_DEBUG=watch, to enable watchpoints on frozen buffers */
143 boolean_t arc_watch = B_FALSE;
144 int arc_procfid;
145 #endif

147 static kmutex_t      arc_reclaim_thr_lock;
148 static kcondvar_t    arc_reclaim_thr_cv; /* used to signal reclaim thr */
149 static uint8_t       arc_thread_exit;

151 extern int zfs_write_limit_shift;
152 extern uint64_t zfs_write_limit_max;
153 extern kmutex_t zfs_write_limit_lock;

155 #define ARC_REDUCE_DNLC_PERCENT 3
156 uint_t arc_reduce_dnlic_percent = ARC_REDUCE_DNLC_PERCENT;

158 typedef enum arc_reclaim_strategy {
159     ARC_RECLAIM_AGGR, /* Aggressive reclaim strategy */
160     ARC_RECLAIM_CONS /* Conservative reclaim strategy */
161 } arc_reclaim_strategy_t;
162 #define unchanged_portion_omitted

235 /* The 6 states: */
236 static arc_state_t ARC_anon;
237 static arc_state_t ARC_mru;
238 static arc_state_t ARC_mru_ghost;
239 static arc_state_t ARC_mfu;
240 static arc_state_t ARC_mfu_ghost;
241 static arc_state_t ARC_l2c_only;

243 typedef struct arc_stats {
244     kstat_named_t arcstat_hits;
245     kstat_named_t arcstat_misses;
246     kstat_named_t arcstat_demand_data_hits;
247     kstat_named_t arcstat_demand_data_misses;
248     kstat_named_t arcstat_demand_metadata_hits;
249     kstat_named_t arcstat_demand_metadata_misses;
250     kstat_named_t arcstat_prefetch_data_hits;
251     kstat_named_t arcstat_prefetch_data_misses;
252     kstat_named_t arcstat_prefetch_metadata_hits;
253     kstat_named_t arcstat_prefetch_metadata_misses;
254     kstat_named_t arcstat_mru_hits;
255     kstat_named_t arcstat_mru_ghost_hits;
256     kstat_named_t arcstat_mfu_hits;
257     kstat_named_t arcstat_mfu_ghost_hits;
258     kstat_named_t arcstat_deleted;
259     kstat_named_t arcstat_recycle_miss;
260     /*
261      * Number of buffers that could not be evicted because the hash lock
262      * was held by another thread. The lock may not necessarily be held
263      * by something using the same buffer, since hash locks are shared
264      * by multiple buffers.

```

```

265     */
266     kstat_named_t arcstat_mutex_miss;
267     /*
268      * Number of buffers skipped because they have I/O in progress, are
269      * indirect prefetch buffers that have not lived long enough, or are
270      * not from the spa we're trying to evict from.
271     */
272     kstat_named_t arcstat_evict_skip;
273     kstat_named_t arcstat_evict_l2_cached;
274     kstat_named_t arcstat_evict_l2_eligible;
275     kstat_named_t arcstat_evict_l2_ineligible;
276     kstat_named_t arcstat_hash_elements;
277     kstat_named_t arcstat_hash_elements_max;
278     kstat_named_t arcstat_hash_collisions;
279     kstat_named_t arcstat_hash_chains;
280     kstat_named_t arcstat_hash_chain_max;
281     kstat_named_t arcstat_p;
282     kstat_named_t arcstat_c;
283     kstat_named_t arcstat_c_min;
284     kstat_named_t arcstat_c_max;
285     kstat_named_t arcstat_size;
286     kstat_named_t arcstat_hdr_size;
287     kstat_named_t arcstat_data_size;
288     kstat_named_t arcstat_other_size;
289     kstat_named_t arcstat_l2_hits;
290     kstat_named_t arcstat_l2_misses;
291     kstat_named_t arcstat_l2_feeds;
292     kstat_named_t arcstat_l2_rw_clash;
293     kstat_named_t arcstat_l2_read_bytes;
294     kstat_named_t arcstat_l2_write_bytes;
295     kstat_named_t arcstat_l2_writes_sent;
296     kstat_named_t arcstat_l2_writes_done;
297     kstat_named_t arcstat_l2_writes_error;
298     kstat_named_t arcstat_l2_writes_hdr_miss;
299     kstat_named_t arcstat_l2_evict_lock_retry;
300     kstat_named_t arcstat_l2_evict_reading;
301     kstat_named_t arcstat_l2_free_on_write;
302     kstat_named_t arcstat_l2_abort_lowmem;
303     kstat_named_t arcstat_l2_cksum_bad;
304     kstat_named_t arcstat_l2_io_error;
305     kstat_named_t arcstat_l2_size;
306     kstat_named_t arcstat_l2_asize;
307     kstat_named_t arcstat_l2_hdr_size;
308     kstat_named_t arcstat_l2_compress_successes;
309     kstat_named_t arcstat_l2_compress_zeros;
310     kstat_named_t arcstat_l2_compress_failures;
311     kstat_named_t arcstat_l2_meta_writes;
312     kstat_named_t arcstat_l2_meta_avg_size;
313     kstat_named_t arcstat_l2_meta_avg_asize;
314     kstat_named_t arcstat_l2_asize_to_meta_ratio;
315     kstat_named_t arcstat_l2_rebuild_attempts;
316     kstat_named_t arcstat_l2_rebuild_successes;
317     kstat_named_t arcstat_l2_rebuild_unsupported;
318     kstat_named_t arcstat_l2_rebuild_timeout;
319     kstat_named_t arcstat_l2_rebuild_arc_bytes;
320     kstat_named_t arcstat_l2_rebuild_l2arc_bytes;
321     kstat_named_t arcstat_l2_rebuild_bufs;
322     kstat_named_t arcstat_l2_rebuild_bufs_precached;
323     kstat_named_t arcstat_l2_rebuild_metabufs;
324     kstat_named_t arcstat_l2_rebuild_uberblk_errors;
325     kstat_named_t arcstat_l2_rebuild_io_errors;
326     kstat_named_t arcstat_l2_rebuild_cksum_errors;
327     kstat_named_t arcstat_l2_rebuild_loop_errors;
328     kstat_named_t arcstat_l2_rebuild_abort_lowmem;
329     kstat_named_t arcstat_memory_throttle_count;
330     kstat_named_t arcstat_duplicate_buffers;

```

```

331     kstat_named_t arcstat_duplicate_buffers_size;
332     kstat_named_t arcstat_duplicate_reads;
333     kstat_named_t arcstat_meta_used;
334     kstat_named_t arcstat_meta_limit;
335     kstat_named_t arcstat_meta_max;
336 } arc_stats_t;

338 static arc_stats_t arc_stats = {
339     "hits", KSTAT_DATA_UINT64 },
340     "misses", KSTAT_DATA_UINT64 },
341     "demand_data_hits", KSTAT_DATA_UINT64 },
342     "demand_data_misses", KSTAT_DATA_UINT64 },
343     "demand_metadata_hits", KSTAT_DATA_UINT64 },
344     "demand_metadata_misses", KSTAT_DATA_UINT64 },
345     "prefetch_data_hits", KSTAT_DATA_UINT64 },
346     "prefetch_data_misses", KSTAT_DATA_UINT64 },
347     "prefetch_metadata_hits", KSTAT_DATA_UINT64 },
348     "prefetch_metadata_misses", KSTAT_DATA_UINT64 },
349     "mru_hits", KSTAT_DATA_UINT64 },
350     "mru_ghost_hits", KSTAT_DATA_UINT64 },
351     "mfu_hits", KSTAT_DATA_UINT64 },
352     "mfu_ghost_hits", KSTAT_DATA_UINT64 },
353     "deleted", KSTAT_DATA_UINT64 },
354     "recycle_miss", KSTAT_DATA_UINT64 },
355     "mutex_miss", KSTAT_DATA_UINT64 },
356     "evict_skip", KSTAT_DATA_UINT64 },
357     "evict_l2_cached", KSTAT_DATA_UINT64 },
358     "evict_l2_eligible", KSTAT_DATA_UINT64 },
359     "evict_l2_ineligible", KSTAT_DATA_UINT64 },
360     "hash_elements", KSTAT_DATA_UINT64 },
361     "hash_elements_max", KSTAT_DATA_UINT64 },
362     "hash_collisions", KSTAT_DATA_UINT64 },
363     "hash_chains", KSTAT_DATA_UINT64 },
364     "hash_chain_max", KSTAT_DATA_UINT64 },
365     "p", KSTAT_DATA_UINT64 },
366     "c", KSTAT_DATA_UINT64 },
367     "c_min", KSTAT_DATA_UINT64 },
368     "c_max", KSTAT_DATA_UINT64 },
369     "size", KSTAT_DATA_UINT64 },
370     "hdr_size", KSTAT_DATA_UINT64 },
371     "data_size", KSTAT_DATA_UINT64 },
372     "other_size", KSTAT_DATA_UINT64 },
373     "l2_hits", KSTAT_DATA_UINT64 },
374     "l2_misses", KSTAT_DATA_UINT64 },
375     "l2_feeds", KSTAT_DATA_UINT64 },
376     "l2_rw_clash", KSTAT_DATA_UINT64 },
377     "l2_read_bytes", KSTAT_DATA_UINT64 },
378     "l2_write_bytes", KSTAT_DATA_UINT64 },
379     "l2_writes_sent", KSTAT_DATA_UINT64 },
380     "l2_writes_done", KSTAT_DATA_UINT64 },
381     "l2_writes_error", KSTAT_DATA_UINT64 },
382     "l2_writes_hdr_miss", KSTAT_DATA_UINT64 },
383     "l2_evict_lock_retry", KSTAT_DATA_UINT64 },
384     "l2_evict_reading", KSTAT_DATA_UINT64 },
385     "l2_free_on_write", KSTAT_DATA_UINT64 },
386     "l2_abort_lowmem", KSTAT_DATA_UINT64 },
387     "l2_ckpt_bad", KSTAT_DATA_UINT64 },
388     "l2_io_error", KSTAT_DATA_UINT64 },
389     "l2_size", KSTAT_DATA_UINT64 },
390     "l2_asize", KSTAT_DATA_UINT64 },
391     "l2_hdr_size", KSTAT_DATA_UINT64 },
392     "l2_compress_successes", KSTAT_DATA_UINT64 },
393     "l2_compress_zeros", KSTAT_DATA_UINT64 },
394     "l2_compress_failures", KSTAT_DATA_UINT64 },
395     "l2_meta_writes", KSTAT_DATA_UINT64 },
396     "l2_meta_avg_size", KSTAT_DATA_UINT64 },

```

```

397     "l2_meta_avg_asize", KSTAT_DATA_UINT64 },
398     "l2_asize_to_meta_ratio", KSTAT_DATA_UINT64 },
399     "l2_rebuild_attempts", KSTAT_DATA_UINT64 },
400     "l2_rebuild_successes", KSTAT_DATA_UINT64 },
401     "l2_rebuild_unsupported", KSTAT_DATA_UINT64 },
402     "l2_rebuild_timeout", KSTAT_DATA_UINT64 },
403     "l2_rebuild_arc_bytes", KSTAT_DATA_UINT64 },
404     "l2_rebuild_l2arc_bytes", KSTAT_DATA_UINT64 },
405     "l2_rebuild_bufs", KSTAT_DATA_UINT64 },
406     "l2_rebuild_precached", KSTAT_DATA_UINT64 },
407     "l2_rebuild_metabufs", KSTAT_DATA_UINT64 },
408     "l2_rebuild_uberblk_errors", KSTAT_DATA_UINT64 },
409     "l2_rebuild_io_errors", KSTAT_DATA_UINT64 },
410     "l2_rebuild_cksum_errors", KSTAT_DATA_UINT64 },
411     "l2_rebuild_loop_errors", KSTAT_DATA_UINT64 },
412     "l2_rebuild_abort_lowmem", KSTAT_DATA_UINT64 },
413     "memory_throttle_count", KSTAT_DATA_UINT64 },
414     "duplicate_buffers", KSTAT_DATA_UINT64 },
415     "duplicate_buffers_size", KSTAT_DATA_UINT64 },
416     "duplicate_reads", KSTAT_DATA_UINT64 },
417     "arc_meta_used", KSTAT_DATA_UINT64 },
418     "arc_meta_limit", KSTAT_DATA_UINT64 },
419     "arc_meta_max", KSTAT_DATA_UINT64 },
420 };
    unchanged_portion_omitted

437 #define ARCSTAT_MAXSTAT(stat) \
438     ARCSTAT_MAX(stat##_max, arc_stats.stat.value.ui64)

440 /*
441  * We define a macro to allow ARC hits/misses to be easily broken down by
442  * two separate conditions, giving a total of four different subtypes for
443  * each of hits and misses (so eight statistics total).
444  */
445 #define ARCSTAT_CONDSTAT(cond1, stat1, notstat1, cond2, stat2, notstat2, stat) \
446     if (cond1) { \
447         if (cond2) { \
448             ARCSTAT_BUMP(arcstat_##stat1##_##stat2##_##stat); \
449         } else { \
450             ARCSTAT_BUMP(arcstat_##stat1##_##notstat2##_##stat); \
451         } \
452     } else { \
453         if (cond2) { \
454             ARCSTAT_BUMP(arcstat_##notstat1##_##stat2##_##stat); \
455         } else { \
456             ARCSTAT_BUMP(arcstat_##notstat1##_##notstat2##_##stat); \
457         } \
458     }

460 /*
461  * This macro allows us to use kstats as floating averages. Each time we
462  * update this kstat, we first factor it and the update value by
463  * ARCSTAT_AVG_FACTOR to shrink the new value's contribution to the overall
464  * average. This macro assumes that integer loads and stores are atomic, but
465  * is not safe for multiple writers updating the kstat in parallel (only the
466  * last writer's update will remain).
467  */
468 #define ARCSTAT_F_AVG_FACTOR 3
469 #define ARCSTAT_F_AVG(stat, value) \
470     do { \
471         uint64_t x = ARCSTAT(stat); \
472         x = x - x / ARCSTAT_F_AVG_FACTOR + \
473             (value) / ARCSTAT_F_AVG_FACTOR; \
474         ARCSTAT(stat) = x; \
475         _NOTE(NOTREACHED) \
476         _NOTE(CONSTCOND) \

```

```

477     } while (0)
478
479 kstat_t      *arc_ksp;
480 static arc_state_t *arc_anon;
481 static arc_state_t *arc_mru;
482 static arc_state_t *arc_mru_ghost;
483 static arc_state_t *arc_mfu;
484 static arc_state_t *arc_mfu_ghost;
485 static arc_state_t *arc_l2c_only;
486
487 /*
488  * There are several ARC variables that are critical to export as kstats --
489  * but we don't want to have to grovel around in the kstat whenever we wish to
490  * manipulate them.  For these variables, we therefore define them to be in
491  * terms of the statistic variable.  This assures that we are not introducing
492  * the possibility of inconsistency by having shadow copies of the variables,
493  * while still allowing the code to be readable.
494  */
495 #define arc_size      ARCSTAT(arcstat_size) /* actual total arc size */
496 #define arc_p        ARCSTAT(arcstat_p) /* target size of MRU */
497 #define arc_c        ARCSTAT(arcstat_c) /* target size of cache */
498 #define arc_c_min    ARCSTAT(arcstat_c_min) /* min target cache size */
499 #define arc_c_max    ARCSTAT(arcstat_c_max) /* max target cache size */
500 #define arc_meta_limit ARCSTAT(arcstat_meta_limit) /* max size for metadata */
501 #define arc_meta_used ARCSTAT(arcstat_meta_used) /* size of metadata */
502 #define arc_meta_max ARCSTAT(arcstat_meta_max) /* max size of metadata */
503
504 #define L2ARC_IS_VALID_COMPRESS(_c_) \
505     ((_c_) == ZIO_COMPRESS_LZ4 || (_c_) == ZIO_COMPRESS_EMPTY)
506
507 static int      arc_no_grow; /* Don't try to grow cache size */
508 static uint64_t arc_tempreserve;
509 static uint64_t arc_loaned_bytes;
510
511 typedef struct l2arc_buf_hdr l2arc_buf_hdr_t;
512
513 typedef struct arc_callback arc_callback_t;
514
515 struct arc_callback {
516     void *acb_private;
517     arc_done_func_t *acb_done;
518     arc_buf_t *acb_buf;
519     zio_t *acb_zio_dummy;
520     arc_callback_t *acb_next;
521 };
522
523 unchanged portion omitted
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544 static buf_hash_table_t buf_hash_table;
545
546 #define BUF_HASH_INDEX(spa, dva, birth) \
547     (buf_hash(spa, dva, birth) & buf_hash_table.ht_mask)
548 #define BUF_HASH_LOCK_NTRY(idx) (buf_hash_table.ht_locks[idx] & (BUF_LOCKS-1))
549 #define BUF_HASH_LOCK(idx) (&(BUF_HASH_LOCK_NTRY(idx).ht_lock))
550 #define HDR_LOCK(hdr) \
551     (BUF_HASH_LOCK(BUF_HASH_INDEX(hdr->b_spa, &hdr->b_dva, hdr->b_birth)))
552
553 uint64_t zfs_crc64_table[256];
554
555 /*
556  * Level 2 ARC
557  */
558
559 #define L2ARC_WRITE_SIZE      (8 * 1024 * 1024) /* initial write max */
560 #define L2ARC_HEADROOM      2 /* num of writes */
561 /*
562  * If we discover during ARC scan any buffers to be compressed, we boost

```

```

563  * our headroom for the next scanning cycle by this percentage multiple.
564  */
565 #define L2ARC_HEADROOM_BOOST      200
566 #define L2ARC_FEED_SECS      1 /* caching interval secs */
567 #define L2ARC_FEED_MIN_MS      200 /* min caching interval ms */
568
569 #define l2arc_writes_sent      ARCSTAT(arcstat_l2_writes_sent)
570 #define l2arc_writes_done      ARCSTAT(arcstat_l2_writes_done)
571
572 /* L2ARC Performance Tunables */
573 uint64_t l2arc_write_max = L2ARC_WRITE_SIZE; /* default max write size */
574 uint64_t l2arc_write_boost = L2ARC_WRITE_SIZE; /* extra write during warmup */
575 uint64_t l2arc_headroom = L2ARC_HEADROOM; /* number of dev writes */
576 uint64_t l2arc_headroom_boost = L2ARC_HEADROOM_BOOST;
577 uint64_t l2arc_feed_secs = L2ARC_FEED_SECS; /* interval seconds */
578 uint64_t l2arc_feed_min_ms = L2ARC_FEED_MIN_MS; /* min interval milliseconds */
579 boolean_t l2arc_noprefetch = B_TRUE; /* don't cache prefetch bufs */
580 boolean_t l2arc_feed_again = B_TRUE; /* turbo warmup */
581 boolean_t l2arc_norw = B_TRUE; /* no reads during writes */
582
583 /*
584  * L2ARC Internals
585  */
586 typedef struct l2arc_dev l2arc_dev_t;
587 typedef struct l2arc_dev {
588     vdev_t *l2ad_vdev; /* vdev */
589     spa_t *l2ad_spa; /* spa */
590     l2ad_hand_t *l2ad_hand; /* next write location */
591     uint64_t l2ad_start; /* first addr on device */
592     uint64_t l2ad_end; /* last addr on device */
593     uint64_t l2ad_evict; /* last addr eviction reached */
594     boolean_t l2ad_first; /* first sweep through */
595     boolean_t l2ad_writing; /* currently writing */
596     list_t *l2ad_buflist; /* buffer list */
597     list_node_t l2ad_node; /* device list node */
598 } l2arc_dev_t;
599
600
601
602
603
604
605
606
607 static list_t L2ARC_dev_list; /* device list */
608 static list_t *l2arc_dev_list; /* device list pointer */
609 static kmutex_t l2arc_dev_mtx; /* device list mutex */
610 static l2arc_dev_t *l2arc_dev_last; /* last device used */
611 static kmutex_t l2arc_buflist_mtx; /* mutex for all buflists */
612 static list_t L2ARC_free_on_write; /* free after write buf list */
613 static list_t *l2arc_free_on_write; /* free after write list ptr */
614 static kmutex_t l2arc_free_on_write_mtx; /* mutex for list */
615 static uint64_t l2arc_ndev; /* number of devices */
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999

```

```

735 static kcondvar_t l2arc_feed_thr_cv;
736 static uint8_t l2arc_thread_exit;

738 static void l2arc_read_done(zio_t *zio);
739 static void l2arc_hdr_stat_add(boolean_t from_arc);
740 static void l2arc_hdr_stat_add(void);
741 static void l2arc_hdr_stat_remove(void);

742 static boolean_t l2arc_compress_buf(l2arc_buf_hdr_t *l2hdr);
743 static void l2arc_decompress_zio(zio_t *zio, arc_buf_hdr_t *hdr,
744     enum zio_compress c);
745 static void l2arc_release_cdata_buf(arc_buf_hdr_t *ab);

747 typedef enum {
748     L2UBLK_BIG_ENDIAN = (1 << 0), /* little endian assumed otherwise */
749     L2UBLK_EVICT_FIRST = (1 << 1) /* mirror of l2ad_first in l2dev */
750 } l2uberblock_flags_t;

752 typedef struct l2uberblock {
753     uint32_t          ub_magic;
754     uint8_t          ub_version;
755     l2uberblock_flags_t ub_flags;

757     uint64_t          ub_spa_guid;
758     uint64_t          ub_birth;
759     uint64_t          ub_evict_tail; /* current evict pointer */
760     uint64_t          ub_alloc_space; /* vdev space alloc status */
761     uint64_t          ub_pbuf_daddr; /* address of newest pbuf */
762     uint32_t          ub_pbuf_asize; /* size of newest pbuf */
763     zio_cksum_t       ub_pbuf_cksum; /* fletcher4 of newest pbuf */

765     zio_cksum_t       ub_cksum; /* cksum of uberblock */
766 } l2uberblock_t;

768 typedef enum {
769     L2PBUF_BIG_ENDIAN = (1 << 0), /* little endian assumed otherwise */
770     L2PBUF_COMPRESSED = (1 << 1) /* pbuf data items are compressed */
771 } l2pbuf_flags_t;

773 typedef struct l2pbuf {
774     uint32_t          pb_magic;
775     unsigned int      pb_version;
776     l2pbuf_flags_t    pb_flags;

778     uint64_t          pb_prev_daddr; /* address of previous pbuf */
779     uint32_t          pb_prev_asize; /* size of previous pbuf */
780     zio_cksum_t       pb_prev_cksum; /* fletcher4 of prev. pbuf */

782     /*
783     * This is a set of item lists that are contained in this pbuf. Each
784     * L2ARC write appends a new l2pbuf_buflist_t array of l2pbuf_buf_t's.
785     * This serves as a soft timeout feature - once the limit of the
786     * number of item lists that a pbuf can hold is reached, the pbuf is
787     * flushed to stable storage, regardless of its total size.
788     */
789     list_t            *pb_buflists_list;

791     /*
792     * Number of compressed bytes referenced by items in this pbuf and
793     * the number of lists present.
794     * This is not actually written to storage, it is only used by
795     * internal algorithms which check for when a pbuf reaches a
796     * certain size limit, after which it is flushed in a write.
797     */
798     uint64_t          pb_payload_asz;
799     /* Same thing for number of buflists */

```

```

800     int                pb_nbuflists;

802     /*
803     * Filled in by l2arc_pbuf_read to hold this pbuf's alloc'd size.
804     * This is then used by l2arc_pbuf_restore to update used space
805     * on the L2ARC vdev.
806     */
807     size_t             pb_asize;
808 } l2pbuf_t;

810 typedef struct l2pbuf_buf l2pbuf_buf_t;
811 typedef struct l2pbuf_buflist {
812     uint32_t          l2pbl_nbufls;
813     l2pbuf_buf_t      *l2pbl_bufls;
814     list_node_t       l2pbl_node;
815 } l2pbuf_buflist_t;

817 struct l2pbuf_buf {
818     dva_t             b_dva; /* dva of buffer */
819     uint64_t          b_birth; /* birth txg of buffer */
820     b_ckpt_t          b_ckpt0;
821     zio_cksum_t       b_freeze_cksum;
822     uint32_t          b_size; /* uncompressed buf size */
823     uint64_t          b_l2daddr; /* buf location on l2dev */
824     uint32_t          b_l2asize; /* actual buf data size */
825     enum zio_compress b_l2compress; /* compression applied */
826     uint16_t          b_contents_type;
827     uint32_t          b_flags;
828 };

830 struct l2arc_dev {
831     vdev_t            *l2ad_vdev; /* vdev */
832     spa_t              *l2ad_spa; /* spa */
833     uint64_t          l2ad_hand; /* next write location */
834     uint64_t          l2ad_start; /* first addr on device */
835     uint64_t          l2ad_end; /* last addr on device */
836     uint64_t          l2ad_evict; /* last addr eviction reached */
837     boolean_t         l2ad_first; /* first sweep through */
838     boolean_t         l2ad_writing; /* currently writing */
839     list_t             *l2ad_buflist; /* buffer list */
840     list_node_t       l2ad_node; /* device list node */
841     l2pbuf_t          l2ad_pbuf; /* currently open pbuf */
842     uint64_t          l2ad_pbuf_daddr; /* prev pbuf daddr */
843     uint64_t          l2ad_pbuf_asize; /* prev pbuf asize */
844     zio_cksum_t       l2ad_pbuf_cksum; /* prev pbuf cksum */
845     /* uberblock birth counter - incremented for each committed uberblk */
846     uint64_t          l2ad_uberblock_birth;
847     /* flag indicating whether a rebuild is currently going on */
848     boolean_t         l2ad_rebuilding;
849 };

851 /* Stores information about an L2ARC prefetch zio */
852 typedef struct l2arc_prefetch_info {
853     uint8_t           *pi_buf; /* where the zio writes to */
854     uint64_t          pi_buflen; /* length of 'buf' */
855     zio_t              *pi_hdr_io; /* see l2arc_pbuf_read below */
856 } l2arc_prefetch_info_t;

858 /* 256 x 4k of l2uberblocks */
859 #define L2UBERBLOCK_SIZE 4096
860 #define L2UBERBLOCK_MAGIC 0x12bab10c
861 #define L2UBERBLOCK_MAX_VERSION 1 /* our maximum uberblock version */
862 #define L2PBUF_MAGIC 0xdb0faba6
863 #define L2PBUF_MAX_VERSION 1 /* our maximum pbuf version */
864 #define L2PBUF_BUF_SIZE 88 /* size of one pbuf buf entry */
865 #define L2PBUF_HDR_SIZE 56 /* pbuf header excluding any payload */

```

```

866 #define L2PBUF_ENCODED_SIZE(_pb) \
867     (L2PBUF_HDR_SIZE + l2arc_pbuf_items_encoded_size(_pb))
868 /*
869  * Allocation limit for the payload of a pbuf. This also fundamentally
870  * limits the number of bufs we can reference in a pbuf.
871  */
872 #define L2PBUF_MAX_PAYLOAD_SIZE (24 * 1024 * 1024)
873 #define L2PBUF_MAX_BUFS (L2PBUF_MAX_PAYLOAD_SIZE / L2PBUF_BUF_SIZE)
874 #define L2PBUF_COMPRESS_MINSZ 8192 /* minimum size to compress a pbuf */
875 #define L2PBUF_MAXSZ 100 * 1024 * 1024 /* maximum pbuf size */
876 #define L2PBUF_MAX_BUFLISTS 128 /* max number of buflists per pbuf */
877 #define L2ARC_REBUILD_TIMEOUT 60 /* a rebuild may take at most 60s */
878 #define L2PBUF_IS_FULL(_pb) \
879     ((_pb)->pb_payload_asz > l2arc_pbuf_max_sz || \
880      (_pb)->pb_nbuflists + 1 >= l2arc_pbuf_max_buflists)
881 /*
882  * These are the flags we allow to persist in L2ARC pbufs. The other flags
883  * of an ARC buffer pertain to the buffer's runtime behavior.
884  */
885 #define L2ARC_PERSIST_FLAGS \
886     (ARC_IN_HASH_TABLE | ARC_L2CACHE | ARC_L2COMPRESS | ARC_PREFETCH)
887
888 /*
889  * Used during L2ARC rebuild after each read operation to check whether we
890  * haven't exceeded the rebuild timeout value.
891  */
892 #define L2ARC_CHK_REBUILD_TIMEOUT(_deadline, ...) \
893     do { \
894         if ((_deadline) != 0 && (_deadline) < ddi_get_lbolt64()) { \
895             __VA_ARGS__; \
896             ARCSTAT_BUMP(arcstat_l2_rebuild_timeout); \
897             cmn_err(CE_WARN, "L2ARC rebuild is taking too long, " \
898                  "dropping remaining L2ARC metadata."); \
899             return; \
900         } \
901         _NOTE(NOTREACHED) \
902         _NOTE(CONSTCOND) \
903     } while (0)
904
905 /*
906  * Performance tuning of L2ARC persistency:
907  */
908 * l2arc_pbuf_compress_minsz : Minimum size of a pbuf in order to attempt
909 * compressing it.
910 * l2arc_pbuf_max_sz : Upper bound on the physical size of L2ARC buffers
911 * referenced from a pbuf. Once a pbuf reaches this size, it is
912 * committed to stable storage. Ideally, there should be approx.
913 * l2arc_dev_size / l2arc_pbuf_max_sz pbufs on an L2ARC device.
914 * l2arc_pbuf_max_buflists : Maximum number of L2ARC feed cycles that will
915 * be buffered in a pbuf before it is committed to L2ARC. This
916 * puts a soft temporal upper bound on pbuf commit intervals.
917 * l2arc_rebuild_enabled : Controls whether L2ARC device adds (either at
918 * pool import or when adding one manually later) will attempt
919 * to rebuild L2ARC buffer contents. In special circumstances,
920 * the administrator may want to set this to B_FALSE, if they
921 * are having trouble importing a pool or attaching an L2ARC
922 * device (e.g. the L2ARC device is slow to read in stored pbuf
923 * metadata, or the metadata has become somehow
924 * fragmented/unusable).
925 * l2arc_rebuild_timeout : A hard timeout value on L2ARC rebuilding to help
926 * avoid a slow L2ARC device from preventing pool import. If we
927 * are not done rebuilding an L2ARC device by this time, we
928 * stop the rebuild and return immediately.
929 */
930 uint64_t l2arc_pbuf_compress_minsz = L2PBUF_COMPRESS_MINSZ;
931 uint64_t l2arc_pbuf_max_sz = L2PBUF_MAXSZ;

```

```

932 uint64_t l2arc_pbuf_max_buflists = L2PBUF_MAX_BUFLISTS;
933 boolean_t l2arc_rebuild_enabled = B_TRUE;
934 uint64_t l2arc_rebuild_timeout = L2ARC_REBUILD_TIMEOUT;
935
936 static void l2arc_rebuild_start(l2arc_dev_t *dev);
937 static void l2arc_rebuild(l2arc_dev_t *dev);
938 static void l2arc_pbuf_restore(l2arc_dev_t *dev, l2pbuf_t *pb);
939 static void l2arc_hdr_restore(const l2pbuf_buf_t *buf, l2arc_dev_t *dev,
940                             uint64_t guid);
941
942 static int l2arc_uberblock_find(l2arc_dev_t *dev, l2uberblock_t *ub);
943 static int l2arc_pbuf_read(l2arc_dev_t *dev, uint64_t daddr, uint32_t asize,
944                            zio_cksum_t cksum, l2pbuf_t *pb, zio_t *this_io, zio_t **next_io);
945 static int l2arc_pbuf_ptr_valid(l2arc_dev_t *dev, uint64_t daddr,
946                                uint32_t asize);
947 static zio_t *l2arc_pbuf_prefetch(vdev_t *vd, uint64_t daddr, uint32_t asize);
948 static void l2arc_pbuf_prefetch_abort(zio_t *zio);
949
950 static void l2arc_uberblock_encode(const l2uberblock_t *ub, uint8_t *buf);
951 static void l2arc_uberblock_decode(const uint8_t *buf, l2uberblock_t *ub);
952 static int l2arc_uberblock_verify(const uint8_t *buf, const l2uberblock_t *ub,
953                                 uint64_t guid);
954 static void l2arc_uberblock_update(l2arc_dev_t *dev, zio_t *pio,
955                                   l2arc_write_callback_t *cb);
956
957 static uint32_t l2arc_pbuf_encode(l2pbuf_t *pb, uint8_t *buf, uint32_t buflen);
958 static int l2arc_pbuf_decode(uint8_t *buf, uint32_t buflen,
959                              l2pbuf_t *pbuf);
960 static int l2arc_pbuf_decode_prev_ptr(const uint8_t *buf, size_t buflen,
961                                       uint64_t *daddr, uint32_t *asize, zio_cksum_t *cksum);
962 static void l2arc_pbuf_init(l2pbuf_t *pb);
963 static void l2arc_pbuf_destroy(l2pbuf_t *pb);
964 static void l2arc_pbuf_commit(l2arc_dev_t *dev, zio_t *pio,
965                               l2arc_write_callback_t *cb);
966 static l2pbuf_buflist_t *l2arc_pbuf_buflist_alloc(l2pbuf_t *pb, int nbufs);
967 static void l2arc_pbuflist_insert(l2pbuf_t *pb, l2pbuf_buflist_t *pbl,
968                                  const arc_buf_hdr_t *ab, int index);
969 static uint32_t l2arc_pbuf_items_encoded_size(l2pbuf_t *pb);
970
971 static uint64_t
972 buf_hash(uint64_t spa, const dva_t *dva, uint64_t birth)
973 {
974     uint8_t *vdva = (uint8_t *)dva;
975     uint64_t crc = -1ULL;
976     int i;
977
978     ASSERT(zfs_crc64_table[128] == ZFS_CRC64_POLY);
979
980     for (i = 0; i < sizeof (dva_t); i++)
981         crc = (crc >> 8) ^ zfs_crc64_table[(crc ^ vdva[i]) & 0xFF];
982
983     crc ^= (spa >> 8) ^ birth;
984
985     return (crc);
986 }
987
988 _____unchanged_portion_omitted_____
989
1425 /*
1426  * Move the supplied buffer to the indicated state. The mutex
1427  * for the buffer must be held by the caller.
1428  */
1429 static void
1430 arc_change_state(arc_state_t *new_state, arc_buf_hdr_t *ab, kmutex_t *hash_lock)
1431 {
1432     arc_state_t *old_state = ab->b_state;
1433     int64_t refcnt = refcount_count(&ab->b_refcnt);

```

```

1434     uint64_t from_delta, to_delta;

1436     ASSERT(MUTEX_HELD(hash_lock));
1437     ASSERT(new_state != old_state);
1438     ASSERT(refcnt == 0 || ab->b_datacnt > 0);
1439     ASSERT(ab->b_datacnt == 0 || !GHOST_STATE(new_state));
1440     ASSERT(ab->b_datacnt <= 1 || old_state != arc_anon);

1442     from_delta = to_delta = ab->b_datacnt * ab->b_size;

1444     /*
1445     * If this buffer is evictable, transfer it from the
1446     * old state list to the new state list.
1447     */
1448     if (refcnt == 0) {
1449         if (old_state != arc_anon) {
1450             int use_mutex = !MUTEX_HELD(&old_state->arcs_mtx);
1451             uint64_t *size = &old_state->arcs_lsize[ab->b_type];

1453             if (use_mutex)
1454                 mutex_enter(&old_state->arcs_mtx);

1456             ASSERT(list_link_active(&ab->b_arc_node));
1457             list_remove(&old_state->arcs_list[ab->b_type], ab);

1459             /*
1460             * If prefetching out of the ghost cache,
1461             * we will have a non-zero datacnt.
1462             */
1463             if (GHOST_STATE(old_state) && ab->b_datacnt == 0) {
1464                 /* ghost elements have a ghost size */
1465                 ASSERT(ab->b_buf == NULL);
1466                 from_delta = ab->b_size;
1467             }
1468             ASSERT3U(*size, >=, from_delta);
1469             atomic_add_64(size, -from_delta);

1471             if (use_mutex)
1472                 mutex_exit(&old_state->arcs_mtx);
1473         }
1474         if (new_state != arc_anon) {
1475             int use_mutex = !MUTEX_HELD(&new_state->arcs_mtx);
1476             uint64_t *size = &new_state->arcs_lsize[ab->b_type];

1478             if (use_mutex)
1479                 mutex_enter(&new_state->arcs_mtx);

1481             list_insert_head(&new_state->arcs_list[ab->b_type], ab);

1483             /* ghost elements have a ghost size */
1484             if (GHOST_STATE(new_state)) {
1485                 ASSERT(ab->b_datacnt == 0);
1486                 ASSERT(ab->b_buf == NULL);
1487                 to_delta = ab->b_size;
1488             }
1489             atomic_add_64(size, to_delta);

1491             if (use_mutex)
1492                 mutex_exit(&new_state->arcs_mtx);
1493         }
1494     }

1496     ASSERT(!BUF_EMPTY(ab));
1497     if (new_state == arc_anon && HDR_IN_HASH_TABLE(ab))
1498         buf_hash_remove(ab);

```

```

1500     /* adjust state sizes */
1501     if (to_delta)
1502         atomic_add_64(&new_state->arcs_size, to_delta);
1503     if (from_delta) {
1504         ASSERT3U(old_state->arcs_size, >=, from_delta);
1505         atomic_add_64(&old_state->arcs_size, -from_delta);
1506     }
1507     ab->b_state = new_state;

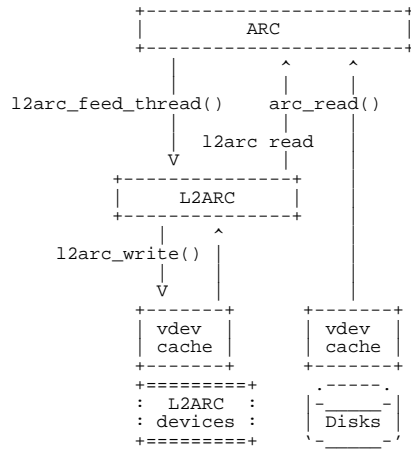
1509     /* adjust l2arc hdr stats */
1510     if (new_state == arc_l2c_only)
1511         l2arc_hdr_stat_add(old_state != arc_anon);
1512     else if (old_state == arc_l2c_only)
1513         l2arc_hdr_stat_remove();
1514 }
1515 }
1516 }
1517 }
1518 }
1519 }
1520 }
1521 }
1522 }
1523 }
1524 }
1525 }
1526 }
1527 }
1528 }
1529 }
1530 }
1531 }
1532 }
1533 }
1534 }
1535 }
1536 }
1537 }
1538 }
1539 }
1540 }
1541 }
1542 }
1543 }
1544 }
1545 }
1546 }
1547 }
1548 }
1549 }
1550 }
1551 }
1552 }
1553 }
1554 }
1555 }
1556 }
1557 }
1558 }
1559 }
1560 }
1561 }
1562 }
1563 }
1564 }
1565 }
1566 }
1567 }
1568 }
1569 }
1570 }
1571 }
1572 }
1573 }
1574 }
1575 }
1576 }
1577 }
1578 }
1579 }
1580 }
1581 }
1582 }
1583 }
1584 }
1585 }
1586 }
1587 }
1588 }
1589 }
1590 }
1591 }
1592 }
1593 }
1594 }
1595 }
1596 }
1597 }
1598 }
1599 }
1600 }
1601 }
1602 }
1603 }
1604 }
1605 }
1606 }
1607 }
1608 }
1609 }
1610 }
1611 }
1612 }
1613 }
1614 }
1615 }
1616 }
1617 }
1618 }
1619 }
1620 }
1621 }
1622 }
1623 }
1624 }
1625 }
1626 }
1627 }
1628 }
1629 }
1630 }
1631 }
1632 }
1633 }
1634 }
1635 }
1636 }
1637 }
1638 }
1639 }
1640 }
1641 }
1642 }
1643 }
1644 }
1645 }
1646 }
1647 }
1648 }
1649 }
1650 }
1651 }
1652 }
1653 }
1654 }
1655 }
1656 }
1657 }
1658 }
1659 }
1660 }
1661 }
1662 }
1663 }
1664 }
1665 }
1666 }
1667 }
1668 }
1669 }
1670 }
1671 }
1672 }
1673 }
1674 }
1675 }
1676 }
1677 }
1678 }
1679 }
1680 }
1681 }
1682 }
1683 }
1684 }
1685 }
1686 }
1687 }
1688 }
1689 }
1690 }
1691 }
1692 }
1693 }
1694 }
1695 }
1696 }
1697 }
1698 }
1699 }
1700 }
1701 }
1702 }
1703 }
1704 }
1705 }
1706 }
1707 }
1708 }
1709 }
1710 }
1711 }
1712 }
1713 }
1714 }
1715 }
1716 }
1717 }
1718 }
1719 }
1720 }
1721 }
1722 }
1723 }
1724 }
1725 }
1726 }
1727 }
1728 }
1729 }
1730 }
1731 }
1732 }
1733 }
1734 }
1735 }
1736 }
1737 }
1738 }
1739 }
1740 }
1741 }
1742 }
1743 }
1744 }
1745 }
1746 }
1747 }
1748 }
1749 }
1750 }
1751 }
1752 }
1753 }
1754 }
1755 }
1756 }
1757 }
1758 }
1759 }
1760 }
1761 }
1762 }
1763 }
1764 }
1765 }
1766 }
1767 }
1768 }
1769 }
1770 }
1771 }
1772 }
1773 }
1774 }
1775 }
1776 }
1777 }
1778 }
1779 }
1780 }
1781 }
1782 }
1783 }
1784 }
1785 }
1786 }
1787 }
1788 }
1789 }
1790 }
1791 }
1792 }
1793 }
1794 }
1795 }
1796 }
1797 }
1798 }
1799 }
1800 }
1801 }
1802 }
1803 }
1804 }
1805 }
1806 }
1807 }
1808 }
1809 }
1810 }
1811 }
1812 }
1813 }
1814 }
1815 }
1816 }
1817 }
1818 }
1819 }
1820 }
1821 }
1822 }
1823 }
1824 }
1825 }
1826 }
1827 }
1828 }
1829 }
1830 }
1831 }
1832 }
1833 }
1834 }
1835 }
1836 }
1837 }
1838 }
1839 }
1840 }
1841 }
1842 }
1843 }
1844 }
1845 }
1846 }
1847 }
1848 }
1849 }
1850 }
1851 }
1852 }
1853 }
1854 }
1855 }
1856 }
1857 }
1858 }
1859 }
1860 }
1861 }
1862 }
1863 }
1864 }
1865 }
1866 }
1867 }
1868 }
1869 }
1870 }
1871 }
1872 }
1873 }
1874 }
1875 }
1876 }
1877 }
1878 }
1879 }
1880 }
1881 }
1882 }
1883 }
1884 }
1885 }
1886 }
1887 }
1888 }
1889 }
1890 }
1891 }
1892 }
1893 }
1894 }
1895 }
1896 }
1897 }
1898 }
1899 }
1900 }
1901 }
1902 }
1903 }
1904 }
1905 }
1906 }
1907 }
1908 }
1909 }
1910 }
1911 }
1912 }
1913 }
1914 }
1915 }
1916 }
1917 }
1918 }
1919 }
1920 }
1921 }
1922 }
1923 }
1924 }
1925 }
1926 }
1927 }
1928 }
1929 }
1930 }
1931 }
1932 }
1933 }
1934 }
1935 }
1936 }
1937 }
1938 }
1939 }
1940 }
1941 }
1942 }
1943 }
1944 }
1945 }
1946 }
1947 }
1948 }
1949 }
1950 }
1951 }
1952 }
1953 }
1954 }
1955 }
1956 }
1957 }
1958 }
1959 }
1960 }
1961 }
1962 }
1963 }
1964 }
1965 }
1966 }
1967 }
1968 }
1969 }
1970 }
1971 }
1972 }
1973 }
1974 }
1975 }
1976 }
1977 }
1978 }
1979 }
1980 }
1981 }
1982 }
1983 }
1984 }
1985 }
1986 }
1987 }
1988 }
1989 }
1990 }
1991 }
1992 }
1993 }
1994 }
1995 }
1996 }
1997 }
1998 }
1999 }
2000 }

```

```

4132 * Level 2 ARC
4133 *
4134 * The level 2 ARC (L2ARC) is a cache layer in-between main memory and disk.
4135 * It uses dedicated storage devices to hold cached data, which are populated
4136 * using large infrequent writes. The main role of this cache is to boost
4137 * the performance of random read workloads. The intended L2ARC devices
4138 * include short-stroked disks, solid state disks, and other media with
4139 * substantially faster read latency than disk.

```



```

4166 * Read requests are satisfied from the following sources, in order:

```

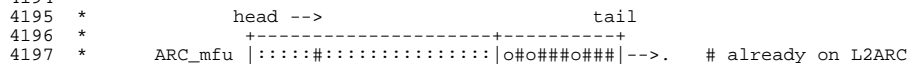
- 4167 * 1) ARC
- 4168 * 2) vdev cache of L2ARC devices
- 4169 * 3) L2ARC devices
- 4170 * 4) vdev cache of disks
- 4171 * 5) disks

```

4172 *
4173 *
4174 * Some L2ARC device types exhibit extremely slow write performance.
4175 * To accommodate for this there are some significant differences between
4176 * the L2ARC and traditional cache design:

```

- 4177 * 1. There is no eviction path from the ARC to the L2ARC. Evictions from
4178 * the ARC behave as usual, freeing buffers and placing headers on ghost
4179 * lists. The ARC does not send buffers to the L2ARC during eviction as
4180 * this would add inflated write latencies for all ARC memory pressure.
4181 *
4182 * 2. The L2ARC attempts to cache data from the ARC before it is evicted.
4183 * It does this by periodically scanning buffers from the eviction-end of
4184 * the MFU and MRU ARC lists, copying them to the L2ARC devices if they are
4185 * not already there. It scans until a headroom of buffers is satisfied,
4186 * which itself is a buffer for ARC eviction. If a compressible buffer is
4187 * found during scanning and selected for writing to an L2ARC device, we
4188 * temporarily boost scanning headroom during the next scan cycle to make
4189 * sure we adapt to compression effects (which might significantly reduce
4190 * the data volume we write to L2ARC). The thread that does this is
4191 * l2arc_feed_thread(), illustrated below; example sizes are included to
4192 * provide a better sense of ratio than this diagram:



```

4198 *
4199 *
4200 *
4201 *
4202 *
4203 *
4204 *
4205 *
4206 *
4207 *
4208 *
4209 *
4210 *
4211 *
4212 *
4213 *
4214 *
4215 *
4216 *
4217 *
4218 *
4219 *
4220 *
4221 *
4222 *
4223 *
4224 *
4225 *
4226 *
4227 *
4228 *
4229 *
4230 *
4231 *
4232 *
4233 *
4234 *
4235 *
4236 *
4237 *
4238 *
4239 *
4240 *
4241 *
4242 *
4243 *
4244 *
4245 *
4246 *
4247 *
4248 *
4249 *
4250 *
4251 *
4252 *
4253 *
4254 *
4255 *
4256 *
4257 *
4258 *
4259 *
4260 *
4261 *
4262 *
4263 *

```

o L2ARC eligible ARC buffer

15.9 Gbytes headroom ^ 32 Mbytes

l2arc_write hand <-- [oooo] --' 8 Mbyte write max

l2arc_feed_thread()

L2ARC dev |#####|#####|#####|#####|...|

32 Gbytes

3. If an ARC buffer is copied to the L2ARC but then hit instead of evicted, then the L2ARC has cached a buffer much sooner than it probably needed to, potentially wasting L2ARC device bandwidth and storage. It is safe to say that this is an uncommon case, since buffers at the end of the ARC lists have moved there due to inactivity.

4. If the ARC evicts faster than the L2ARC can maintain a headroom, then the L2ARC simply misses copying some buffers. This serves as a pressure valve to prevent heavy read workloads from both stalling the ARC with waits and clogging the L2ARC with writes. This also helps prevent the potential for the L2ARC to churn if it attempts to cache content too quickly, such as during backups of the entire pool.

5. After system boot and before the ARC has filled main memory, there are no evictions from the ARC and so the tails of the ARC_mfu and ARC_mru lists can remain mostly static. Instead of searching from tail of these lists as pictured, the l2arc_feed_thread() will search from the list heads for eligible buffers, greatly increasing its chance of finding them.

The L2ARC device write speed is also boosted during this time so that the L2ARC warms up faster. Since there have been no ARC evictions yet, there are no L2ARC reads, and no fear of degrading read performance through increased writes.

6. Writes to the L2ARC devices are grouped and sent in-sequence, so that the vdev queue can aggregate them into larger and fewer writes. Each device is written to in a rotor fashion, sweeping writes through available space then repeating.

7. The L2ARC does not store dirty content. It never needs to flush write buffers back to disk based storage.

8. If an ARC buffer is written (and dirtied) which also exists in the L2ARC, the now stale L2ARC buffer is immediately dropped.

The performance of the L2ARC can be tweaked by a number of tunables, which may be necessary for different workloads:

l2arc_write_max	max write bytes per interval
l2arc_write_boost	extra write bytes during device warmup
l2arc_noprefetch	skip caching prefetched buffers
l2arc_headroom	number of max device writes to precache
l2arc_headroom_boost	when we find compressed buffers during ARC scanning, we multiply headroom by this percentage factor for the next scan cycle, since more compressed buffers are likely to be present
l2arc_feed_secs	seconds between L2ARC writing

Tunables may be removed or added as future performance improvements are


```

4264 * integrated, and also may become zpool properties.
4265 *
4266 * There are three key functions that control how the L2ARC warms up:
4267 *
4268 *     l2arc_write_eligible() check if a buffer is eligible to cache
4269 *     l2arc_write_size()     calculate how much to write
4270 *     l2arc_write_interval() calculate sleep delay between writes
4271 *
4272 * These three functions determine what to write, how much, and how quickly
4273 * to send writes.
4274 *
4275 * L2ARC persistency:
4276 *
4277 * When writing buffers to L2ARC, we periodically add some metadata to
4278 * make sure we can pick them up after reboot, thus dramatically reducing
4279 * the impact that any downtime has on the performance of storage systems
4280 * with large caches.
4281 *
4282 * The implementation works fairly simply by integrating the following two
4283 * modifications:
4284 *
4285 * *) Every now and then, at end of an L2ARC feed cycle, we append a piece
4286 * of metadata (called a "pbuf", or "persistency buffer") to the L2ARC
4287 * write. This allows us to understand what what's been written, so that
4288 * we can rebuild the arc_buf_hdr_t structures of the main ARC buffers.
4289 * The pbuf also includes a "back-reference" pointer to the previous
4290 * pbuf, forming a linked list of pbufs on the L2ARC device.
4291 *
4292 * *) We reserve 4k of space at the start of each L2ARC device for our
4293 * header bookkeeping purposes. This contains a single 4k uberblock, which
4294 * contains our top-level reference structures. We update it on each pbuf
4295 * write. If this write results in an inconsistent uberblock (e.g. due to
4296 * power failure), we detect this by verifying the uberblock's checksum
4297 * and simply drop the entries from L2ARC. Once an L2ARC pbuf update
4298 * completes, we update the uberblock to point to it.
4299 *
4300 * Implementation diagram:
4301 *
4302 * +=== L2ARC device (not to scale) =====+
4303 * |                                         |
4304 * |-----newest pbuf pointer----->|
4305 * |                                         |
4306 * | |l2uberblock| --- |bufs|pbuf|bufs|pbuf|bufs|pbuf|bufs|pbuf| --- (empty) --- |
4307 * |                                         |
4308 * |     ^   ^   ^   ^   ^   ^   ^   ^   ^   ^   ^   ^   ^   ^   ^   ^   ^   ^   ^   ^
4309 * |     |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
4310 * |     |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
4311 * |-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
4312 *
4313 * On-device data structures:
4314 *
4315 * (L2ARC persistent uberblock)
4316 * struct l2uberblock {
4317 *     (these fields are in network byte order)
4318 *     uint32_t magic = 0x12bab10c;    l2-ber-block
4319 *     uint8_t  version = 0x1;
4320 *     uint8_t  reserved = 0x0;
4321 *     uint16_t ublk_flags;             see l2uberblock_flags_t
4322 *
4323 *     (byte order of fields below determined by 'ublk_flags')
4324 *     uint64_t spa_guid;               what pool this l2arc dev belongs to
4325 *     uint64_t birth_tngx;            ublk with highest birth_tngx is newest
4326 *     uint64_t evict_tail;            current evict pointer on l2arc dev
4327 *     uint64_t alloc_space;           how much space is alloc'd on the dev
4328 *     uint64_t pbuf_daddr;           dev addr of the newest l2pbuf_t
4329 *     uint32_t pbuf_asize;            size of newest pbuf
4330 *     uint64_t pbuf_cksum[4];        fletcher4 of newest pbuf

```

```

4330 *
4331 *     uint8_t  reserved[3996] = {0x0, 0x0, ... 0x0};
4332 *
4333 *     uint64_t ublk_cksum[4] = fletcher4(of the 4064 bytes above);
4334 * } l2dev_uberblock;
4335 *
4336 * (L2ARC persistent buffer list)
4337 * typedef struct l2pbuf_t {
4338 *     (these fields are in network byte order)
4339 *     uint32_t magic = 0xdb0faba6;    the-buffer-bag
4340 *     uint8_t  version = 0x1;
4341 *     uint8_t  reserved = 0x0;
4342 *     uint16_t pbuf_flags;            see l2pbuf_flags_t
4343 *
4344 *     (byte order of fields below determined by 'pbuf_flags')
4345 *     uint64_t prev_pbuf_daddr;       previous pbuf dev addr
4346 *     uint32_t prev_pbuf_asize;       previous pbuf size
4347 *     uint64_t prev_pbuf_cksum[4];    fletcher4(of previous pbuf)
4348 *
4349 *     uint32_t items_size;             uncompressed size of 'items' below
4350 *     (if (pbuf_flags & COMPRESS) decompress 'items' prior to decoding)
4351 *     struct l2pbuf_buf_item {
4352 *         (these fields mirror [l2]arc_buf_hdr fields)
4353 *         uint64_t dva[2];             buffer's DVA
4354 *         uint64_t birth;              buffer's birth TXG in ARC
4355 *         uint64_t cksum0;             lower 64-bits of buffer's cksum
4356 *         uint64_t freeze_cksum[4];    buffer's freeze cksum
4357 *         uint32_t size;               uncompressed buffer data size
4358 *         uint64_t l2daddr;            device address (offset) of buf
4359 *         uint32_t l2asize;            actual space occupied by buf
4360 *         uint8_t  compress;           compress algo used on data
4361 *         uint8_t  contents_type;      buffer's contents type
4362 *         uint16_t reserved = 0x0;     for alignment and future use
4363 *         uint32_t flags;              buffer's persistent flags
4364 *     } items[];                       continues for remainder of pbuf
4365 * } l2pbuf_t;
4366 *
4367 * L2ARC reconstruction:
4368 *
4369 * When writing data, we simply write in the standard rotary fashion,
4370 * evicting buffers as we go and simply writing new data over them (appending
4371 * an updated l2pbuf_t every now and then). This obviously means that once we
4372 * loop around the end of the device, we will start cutting into an already
4373 * committed l2pbuf (and its referenced data buffers), like so:
4374 *
4375 *     current write head --> |-----|  <-----| old tail
4376 *
4377 * |-----|-----|-----|-----|-----|-----|-----|-----|-----|
4378 * |<---|bufs|pbuf|bufs|pbuf|bufs|pbuf|bufs|pbuf|-----|bufs|pbuf|bufs|pbuf|---|
4379 * |                ^                ^                ^                ^                ^
4380 * |                |                |                |                |                |
4381 * |                |                |                |                |                |
4382 * |                |                |                |                |                |
4383 * |                |                |                |                |                |
4384 * |                |                |                |                |                |
4385 * |                |                |                |                |                |
4386 * |                |                |                |                |                |
4387 * |                |                |                |                |                |
4388 * |                |                |                |                |                |
4389 * |                |                |                |                |                |
4390 * |                |                |                |                |                |
4391 * |                |                |                |                |                |
4392 * |                |                |                |                |                |
4393 * |                |                |                |                |                |
4394 * |                |                |                |                |                |
4395 * |                |                |                |                |                |
4396 * |                |                |                |                |                |
4397 * |                |                |                |                |                |
4398 * |                |                |                |                |                |
4399 * |                |                |                |                |                |
4400 * |                |                |                |                |                |
4401 * |                |                |                |                |                |
4402 * |                |                |                |                |                |
4403 * |                |                |                |                |                |
4404 * |                |                |                |                |                |
4405 * |                |                |                |                |                |
4406 * |                |                |                |                |                |
4407 * |                |                |                |                |                |
4408 * |                |                |                |                |                |
4409 * |                |                |                |                |                |
4410 * |                |                |                |                |                |
4411 * |                |                |                |                |                |
4412 * |                |                |                |                |                |
4413 * |                |                |                |                |                |
4414 * |                |                |                |                |                |
4415 * |                |                |                |                |                |
4416 * |                |                |                |                |                |
4417 * |                |                |                |                |                |
4418 * |                |                |                |                |                |
4419 * |                |                |                |                |                |
4420 * |                |                |                |                |                |
4421 * |                |                |                |                |                |
4422 * |                |                |                |                |                |
4423 * |                |                |                |                |                |
4424 * |                |                |                |                |                |
4425 * |                |                |                |                |                |
4426 * |                |                |                |                |                |
4427 * |                |                |                |                |                |
4428 * |                |                |                |                |                |
4429 * |                |                |                |                |                |
4430 * |                |                |                |                |                |
4431 * |                |                |                |                |                |
4432 * |                |                |                |                |                |
4433 * |                |                |                |                |                |
4434 * |                |                |                |                |                |
4435 * |                |                |                |                |                |
4436 * |                |                |                |                |                |
4437 * |                |                |                |                |                |
4438 * |                |                |                |                |                |
4439 * |                |                |                |                |                |
4440 * |                |                |                |                |                |
4441 * |                |                |                |                |                |
4442 * |                |                |                |                |                |
4443 * |                |                |                |                |                |
4444 * |                |                |                |                |                |
4445 * |                |                |                |                |                |
4446 * |                |                |                |                |                |
4447 * |                |                |                |                |                |
4448 * |                |                |                |                |                |
4449 * |                |                |                |                |                |
4450 * |                |                |                |                |                |
4451 * |                |                |                |                |                |
4452 * |                |                |                |                |                |
4453 * |                |                |                |                |                |
4454 * |                |                |                |                |                |
4455 * |                |                |                |                |                |
4456 * |                |                |                |                |                |
4457 * |                |                |                |                |                |
4458 * |                |                |                |                |                |
4459 * |                |                |                |                |                |
4460 * |                |                |                |                |                |
4461 * |                |                |                |                |                |
4462 * |                |                |                |                |                |
4463 * |                |                |                |                |                |
4464 * |                |                |                |                |                |
4465 * |                |                |                |                |                |
4466 * |                |                |                |                |                |
4467 * |                |                |                |                |                |
4468 * |                |                |                |                |                |
4469 * |                |                |                |                |                |
4470 * |                |                |                |                |                |
4471 * |                |                |                |                |                |
4472 * |                |                |                |                |                |
4473 * |                |                |                |                |                |
4474 * |                |                |                |                |                |
4475 * |                |                |                |                |                |
4476 * |                |                |                |                |                |
4477 * |                |                |                |                |                |
4478 * |                |                |                |                |                |
4479 * |                |                |                |                |                |
4480 * |                |                |                |                |                |
4481 * |                |                |                |                |                |
4482 * |                |                |                |                |                |
4483 * |                |                |                |                |                |
4484 * |                |                |                |                |                |
4485 * |                |                |                |                |                |
4486 * |                |                |                |                |                |
4487 * |                |                |                |                |                |
4488 * |                |                |                |                |                |
4489 * |                |                |                |                |                |
4490 * |                |                |                |                |                |
4491 * |                |                |                |                |                |
4492 * |                |                |                |                |                |
4493 * |                |                |                |                |                |
4494 * |                |                |                |                |                |
4495 * |                |                |                |                |                |
4496 * |                |                |                |                |                |
4497 * |                |                |                |                |                |
4498 * |                |                |                |                |                |
4499 * |                |                |                |                |                |
4500 * |                |                |                |                |                |
4501 * |                |                |                |                |                |
4502 * |                |                |                |                |                |
4503 * |                |                |                |                |                |
4504 * |                |                |                |                |                |
4505 * |                |                |                |                |                |
4506 * |                |                |                |                |                |
4507 * |                |                |                |                |                |
4508 * |                |                |                |                |                |
4509 * |                |                |                |                |                |
4510 * |                |                |                |                |                |
4511 * |                |                |                |                |                |
4512 * |                |                |                |                |                |
4513 * |                |                |                |                |                |
4514 * |                |                |                |                |                |
4515 * |                |                |                |                |                |
4516 * |                |                |                |                |                |
4517 * |                |                |                |                |                |
4518 * |                |                |                |                |                |
4519 * |                |                |                |                |                |
4520 * |                |                |                |                |                |
4521 * |                |                |                |                |                |
4522 * |                |                |                |                |                |
4523 * |                |                |                |                |                |
4524 * |                |                |                |                |                |
4525 * |                |                |                |                |                |
4526 * |                |                |                |                |                |
4527 * |                |                |                |                |                |
4528 * |                |                |                |                |                |
4529 * |                |                |                |                |                |
4530 * |                |                |                |                |                |
4531 * |                |                |                |                |                |
4532 * |                |                |                |                |                |
4533 * |                |                |                |                |                |
4534 * |                |                |                |                |                |
4535 * |                |                |                |                |                |
4536 * |                |                |                |                |                |
4537 * |                |                |                |                |                |
4538 * |                |                |                |                |                |
4539 * |                |                |                |                |                |
4540 * |                |                |                |                |                |
4541 * |                |                |                |                |                |
4542 * |                |                |                |                |                |
4543 * |                |                |                |                |                |
4544 * |                |                |                |                |                |
4545 * |                |                |                |                |                |
4546 * |                |                |                |                |                |
4547 * |                |                |                |                |                |
4548 * |                |                |                |                |                |
4549 * |                |                |                |                |                |
4550 * |                |                |                |                |                |
4551 * |                |                |                |                |                |
4552 * |                |                |                |                |                |
4553 * |                |                |                |                |                |
4554 * |                |                |                |                |                |
4555 * |                |                |                |                |                |
4556 * |                |                |                |                |                |
4557 * |                |                |                |                |                |
4558 * |                |                |                |                |                |
4559 * |                |                |                |                |                |
4560 * |                |                |                |                |                |
4561 * |                |                |                |                |                |
4562 * |                |                |                |                |                |
4563 * |                |                |                |                |                |
4564 * |                |                |                |                |                |
4565 * |                |                |                |                |                |
4566 * |                |                |                |                |                |
4567 * |                |                |                |                |                |
4568 * |                |                |                |                |                |
4569 * |                |                |                |                |                |
4570 * |                |                |                |                |                |
4571 * |                |                |                |                |                |
4572 * |                |                |                |                |                |
4573 * |                |                |                |                |                |
4574 * |                |                |                |                |                |
4575 * |                |                |                |                |                |
4576 * |                |                |                |                |                |
4577 * |                |                |                |                |                |
4578 * |                |                |                |                |                |
4579 * |                |                |                |                |                |
4580 * |                |                |                |                |                |
4581 * |                |                |                |                |                |
4582 * |                |                |                |                |                |
4583 * |                |                |                |                |                |
4584 * |                |                |                |                |                |
4585 * |                |                |                |                |                |
4586 * |                |                |                |                |                |
4587 * |                |                |                |                |                |
4588 * |                |                |                |                |                |
4589 * |                |                |                |                |                |
4590 * |                |                |                |                |                |
4591 * |                |                |                |                |                |
4592 * |                |                |                |                |                |
4593 * |                |                |                |                |                |
4594 * |                |                |                |                |                |
4595 * |                |                |                |                |                |

```

```

4396 * from an L2ARC device: what about invalidated buffers? Given the above
4397 * construction, we cannot update pbufs which we've already written to amend
4398 * them to remove buffers which were invalidated. Thus, during reconstruction,
4399 * we might be populating the cache with buffers for data that's not on the
4400 * main pool anymore, or may have been overwritten!
4401 *
4402 * As it turns out, this isn't a problem. Every arc_read request includes
4403 * both the DVA and, crucially, the birth TXG of the BP the caller is
4404 * looking for. So even if the cache were populated by completely rotten
4405 * blocks for data that had been long deleted and/or overwritten, we'll
4406 * never actually return bad data from the cache, since the DVA with the
4407 * birth TXG uniquely identify a block in space and time - once created,
4408 * a block is immutable on disk. The worst thing we have done is wasted
4409 * some time and memory at l2arc rebuild to reconstruct outdated ARC
4410 * entries that will get dropped from the l2arc as it is being updated
4411 * with new blocks.
4412 */

4414 static boolean_t
4415 l2arc_write_eligible(uint64_t spa_guid, arc_buf_hdr_t *ab)
4416 {
4417     /*
4418      * A buffer is *not* eligible for the L2ARC if it:
4419      * 1. belongs to a different spa.
4420      * 2. is already cached on the L2ARC.
4421      * 3. has an I/O in progress (it may be an incomplete read).
4422      * 4. is flagged not eligible (zfs property).
4423      */
4424     if (ab->b_spa != spa_guid || ab->b_l2hdr != NULL ||
4425         HDR_IO_IN_PROGRESS(ab) || !HDR_L2CACHE(ab))
4426         return (B_FALSE);

4428     return (B_TRUE);
4429 }
    unchanged portion omitted

4477 static void
4478 l2arc_hdr_stat_add(boolean_t from_arc)
4479 {
4480     ARCSTAT_INCR(arcstat_l2_hdr_size, HDR_SIZE + L2HDR_SIZE);
4481     if (from_arc)
4482         ARCSTAT_INCR(arcstat_hdr_size, -HDR_SIZE);
4483 }
    unchanged portion omitted

4492 /*
4493  * Cycle through L2ARC devices. This is how L2ARC load balances.
4494  * If a device is returned, this also returns holding the spa config lock.
4495  */
4496 static l2arc_dev_t *
4497 l2arc_dev_get_next(void)
4498 {
4499     l2arc_dev_t *first, *next = NULL;

4501     /*
4502      * Lock out the removal of spas (spa_namespace_lock), then removal
4503      * of cache devices (l2arc_dev_mtx). Once a device has been selected,
4504      * both locks will be dropped and a spa config lock held instead.
4505      */
4506     mutex_enter(&spa_namespace_lock);
4507     mutex_enter(&l2arc_dev_mtx);

4509     /* if there are no vdevs, there is nothing to do */
4510     if (l2arc_ndev == 0)
4511         goto out;

```

```

4513     first = NULL;
4514     next = l2arc_dev_last;
4515     do {
4516         /*
4517          * Loop around the list looking for a non-faulted vdev
4518          * and one that isn't currently doing an L2ARC rebuild.
4519          */
4520         /* loop around the list looking for a non-faulted vdev */
4521         if (next == NULL) {
4522             next = list_head(l2arc_dev_list);
4523         } else {
4524             next = list_next(l2arc_dev_list, next);
4525             if (next == NULL)
4526                 next = list_head(l2arc_dev_list);
4527         }

4528         /* if we have come back to the start, bail out */
4529         if (first == NULL)
4530             first = next;
4531         else if (next == first)
4532             break;

4534     } while (vdev_is_dead(next->l2ad_vdev) || next->l2ad_rebuilding);
4535     } while (vdev_is_dead(next->l2ad_vdev));

4536     /* if we were unable to find any usable vdevs, return NULL */
4537     if (vdev_is_dead(next->l2ad_vdev) || next->l2ad_rebuilding)
4538         next = NULL;

4540     l2arc_dev_last = next;

4542 out:
4543     mutex_exit(&l2arc_dev_mtx);

4545     /*
4546      * Grab the config lock to prevent the 'next' device from being
4547      * removed while we are writing to it.
4548      */
4549     if (next != NULL)
4550         spa_config_enter(next->l2ad_spa, SCL_L2ARC, next, RW_READER);
4551     mutex_exit(&spa_namespace_lock);

4553     return (next);
4554 }
    unchanged portion omitted

4580 /*
4581  * A write to a cache device has completed. Update all headers to allow
4582  * reads from these buffers to begin.
4583  */
4584 static void
4585 l2arc_write_done(zio_t *zio)
4586 {
4587     l2arc_write_callback_t *cb;
4588     l2arc_dev_t *dev;
4589     list_t *buflist;
4590     arc_buf_hdr_t *head, *ab, *ab_prev;
4591     l2arc_buf_hdr_t *abl2;
4592     kmutex_t *hash_lock;

4594     cb = zio->io_private;
4595     ASSERT(cb != NULL);
4596     dev = cb->l2wcb_dev;
4597     ASSERT(dev != NULL);

```

```

4598     head = cb->l2wcb_head;
4599     ASSERT(head != NULL);
4600     buflist = dev->l2ad_buflist;
4601     ASSERT(buflist != NULL);
4602     DTRACE_PROBE2(l2arc_iodone, zio_t *, zio,
4603                 l2arc_write_callback_t *, cb);

4605     if (zio->io_error != 0)
4606         ARCSTAT_BUMP(arcstat_l2_writes_error);

4608     mutex_enter(&l2arc_buflist_mtx);

4610     /*
4611      * All writes completed, or an error was hit.
4612      */
4613     for (ab = list_prev(buflist, head); ab; ab = ab_prev) {
4614         ab_prev = list_prev(buflist, ab);
4615         abl2 = ab->b_l2hdr;

4617         /*
4618          * Release the temporary compressed buffer as soon as possible.
4619          */
4620         if (abl2->b_compress != ZIO_COMPRESS_OFF)
4621             l2arc_release_cdata_buf(ab);

4623         hash_lock = HDR_LOCK(ab);
4624         if (!mutex_tryenter(hash_lock)) {
4625             /*
4626              * This buffer misses out. It may be in a stage
4627              * of eviction. Its ARC_L2_WRITING flag will be
4628              * left set, denying reads to this buffer.
4629              */
4630             ARCSTAT_BUMP(arcstat_l2_writes_hdr_miss);
4631             continue;
4632         }

4634         abl2 = ab->b_l2hdr;

4636         /*
4637          * Release the temporary compressed buffer as soon as possible.
4638          */
4639         if (abl2->b_compress != ZIO_COMPRESS_OFF)
4640             l2arc_release_cdata_buf(ab);

4642         if (zio->io_error != 0) {
4643             /*
4644              * Error - drop L2ARC entry.
4645              */
4646             list_remove(buflist, ab);
4647             ARCSTAT_INCR(arcstat_l2_asize, -abl2->b_asize);
4648             ab->b_l2hdr = NULL;
4649             kmem_free(abl2, sizeof(l2arc_buf_hdr_t));
4650             ARCSTAT_INCR(arcstat_l2_size, -ab->b_size);
4651         }

4653         /*
4654          * Allow ARC to begin reads to this L2ARC entry.
4655          */
4656         ab->b_flags &= ~ARC_L2_WRITING;

4658         mutex_exit(hash_lock);

4660     }

4662     atomic_inc_64(&l2arc_writes_done);
4663     list_remove(buflist, head);
4664     kmem_cache_free(hdr_cache, head);

```

```

4656     mutex_exit(&l2arc_buflist_mtx);

4658     l2arc_do_free_on_write();

4660     if (cb->l2wcb_pbuf)
4661         kmem_free(cb->l2wcb_pbuf, cb->l2wcb_pbuf_size);
4662     if (cb->l2wcb_ub_buf)
4663         kmem_free(cb->l2wcb_ub_buf, L2UBERBLOCK_SIZE);
4664     kmem_free(cb, sizeof(l2arc_write_callback_t));
4665 }
    unchanged portion omitted

4918 /*
4919  * Find and write ARC buffers to the L2ARC device.
4920  */
4921  * An ARC_L2_WRITING flag is set so that the L2ARC buffers are not valid
4922  * for reading until they have completed writing.
4923  * The headroom_boost is an in-out parameter used to maintain headroom boost
4924  * state between calls to this function.
4925  *
4926  * Returns the number of bytes actually written (which may be smaller than
4927  * the delta by which the device hand has changed due to alignment).
4928  */
4929 static uint64_t
4930 l2arc_write_buffers(spa_t *spa, l2arc_dev_t *dev, uint64_t target_sz,
4931                   boolean_t *headroom_boost)
4932 {
4933     arc_buf_hdr_t *ab, *ab_prev, *head;
4934     list_t *list;
4935     uint64_t write_asize, write_psize, write_sz, headroom,
4936             buf_compress_minsz;
4937     void *buf_data;
4938     kmutex_t *list_lock;
4939     boolean_t full;
4940     l2arc_write_callback_t *cb;
4941     zio_t *pio, *wzio;
4942     uint64_t guid = spa_load_guid(spa);
4943     const boolean_t do_headroom_boost = *headroom_boost;

4945     /* persistency-related */
4946     l2pbuf_t *pb;
4947     l2pbuf_buflist_t *pb_buflist;
4948     int num_bufs, buf_index;

4950     ASSERT(dev->l2ad_vdev != NULL);

4952     /* Lower the flag now, we might want to raise it again later. */
4953     *headroom_boost = B_FALSE;

4955     pio = NULL;
4956     cb = NULL;
4957     write_sz = write_asize = write_psize = 0;
4958     full = B_FALSE;
4959     head = kmem_cache_alloc(hdr_cache, KM_PUSHPAGE);
4960     head->b_flags |= ARC_L2_WRITE_HEAD;

4962     /*
4963      * We will want to try to compress buffers that are at least 2x the
4964      * device sector size.
4965      */
4966     buf_compress_minsz = 2 << dev->l2ad_vdev->vdev_ashift;

4968     pb = &dev->l2ad_pbuf;
4969     num_bufs = 0;

4971     /*

```

```

4972  * We will want to try to compress buffers that are at least 2x the
4973  * device sector size.
4974  */
4975  buf_compress_minsz = 2 << dev->l2ad_vdev->vdev_ashift;

4977  /*
4978  * Copy buffers for L2ARC writing.
4979  */
4980  mutex_enter(&l2arc_buflist_mtx);
4981  for (int try = 0; try <= 3; try++) {
4982      uint64_t passed_sz = 0;

4984      list = l2arc_list_locked(try, &list_lock);

4986      /*
4987      * L2ARC fast warmup.
4988      *
4989      * Until the ARC is warm and starts to evict, read from the
4990      * head of the ARC lists rather than the tail.
4991      */
4992      if (arc_warm == B_FALSE)
4993          ab = list_head(list);
4994      else
4995          ab = list_tail(list);

4997      headroom = target_sz * l2arc_headroom;
4998      if (do_headroom_boost)
4999          headroom = (headroom * l2arc_headroom_boost) / 100;

5001      for (; ab; ab = ab_prev) {
5002          l2arc_buf_hdr_t *l2hdr;
5003          kmutex_t *hash_lock;
5004          uint64_t buf_sz;

5006          if (arc_warm == B_FALSE)
5007              ab_prev = list_next(list, ab);
5008          else
5009              ab_prev = list_prev(list, ab);

5011          hash_lock = HDR_LOCK(ab);
5012          if (!mutex_tryenter(hash_lock)) {
5013              /*
5014              * Skip this buffer rather than waiting.
5015              */
5016              continue;
5017          }

5019          passed_sz += ab->b_size;
5020          if (passed_sz > headroom) {
5021              /*
5022              * Searched too far.
5023              */
5024              mutex_exit(hash_lock);
5025              break;
5026          }

5028          if (!l2arc_write_eligible(guid, ab)) {
5029              mutex_exit(hash_lock);
5030              continue;
5031          }

5033          if ((write_sz + ab->b_size) > target_sz) {
5034              full = B_TRUE;
5035              mutex_exit(hash_lock);
5036              break;
5037          }

```

```

5039          if (pio == NULL) {
5040              /*
5041              * Insert a dummy header on the buflist so
5042              * l2arc_write_done() can find where the
5043              * write buffers begin without searching.
5044              */
5045              list_insert_head(dev->l2ad_buflist, head);

5047          cb = kmem_zalloc(
5048              cb = kmem_alloc(
5049                  sizeof(l2arc_write_callback_t), KM_SLEEP);
5050              cb->l2wcb_dev = dev;
5051              cb->l2wcb_head = head;
5052              pio = zio_root(spa, l2arc_write_done, cb,
5053                  ZIO_FLAG_CANFAIL);
5054          }

5055          /*
5056          * Create and add a new L2ARC header.
5057          */
5058          l2hdr = kmem_zalloc(sizeof(l2arc_buf_hdr_t), KM_SLEEP);
5059          l2hdr->b_dev = dev;
5060          ab->b_flags |= ARC_L2_WRITING;

5062          /*
5063          * Temporarily stash the data buffer in b_tmp_cdata.
5064          * The subsequent write step will pick it up from
5065          * there. This is because can't access ab->b_buf
5066          * without holding the hash_lock, which we in turn
5067          * can't access without holding the ARC list locks
5068          * (which we want to avoid during compression/writing).
5069          */
5070          l2hdr->b_compress = ZIO_COMPRESS_OFF;
5071          l2hdr->b_asize = ab->b_size;
5072          l2hdr->b_tmp_cdata = ab->b_buf->b_data;

5074          buf_sz = ab->b_size;
5075          ab->b_l2hdr = l2hdr;

5077          list_insert_head(dev->l2ad_buflist, ab);

5079          /*
5080          * Compute and store the buffer cksum before
5081          * writing. On debug the cksum is verified first.
5082          */
5083          arc_cksum_verify(ab->b_buf);
5084          arc_cksum_compute(ab->b_buf, B_TRUE);

5086          mutex_exit(hash_lock);

5088          write_sz += buf_sz;
5089          num_bufs++;
5090      }

5092      mutex_exit(list_lock);

5094      if (full == B_TRUE)
5095          break;
5096  }

5098  /* No buffers selected for writing? */
5099  if (pio == NULL) {
5100      ASSERT0(write_sz);
5101      mutex_exit(&l2arc_buflist_mtx);
5102      kmem_cache_free(hdr_cache, head);

```

```

5103         return (0);
5104     }

5106     /* expand the pbuff to include a new list */
5107     pb_buflist = l2arc_pbuff_buflist_alloc(pb, num_bufs);

5109     /*
5110      * Now start writing the buffers. We're starting at the write head
5111      * and work backwards, retracing the course of the buffer selector
5112      * loop above.
5113      */
5114     for (ab = list_prev(dev->l2ad_buflist, head), buf_index = 0; ab;
5115          ab = list_prev(dev->l2ad_buflist, ab), buf_index++) {
5116         for (ab = list_prev(dev->l2ad_buflist, head); ab;
5117              ab = list_prev(dev->l2ad_buflist, ab)) {
5118             l2arc_buf_hdr_t *l2hdr;
5119             uint64_t buf_sz;

5121             /*
5122              * We shouldn't need to lock the buffer here, since we flagged
5123              * it as ARC_L2_WRITING in the previous step, but we must take
5124              * care to only access its L2 cache parameters. In particular,
5125              * ab->b_buf may be invalid by now due to ARC eviction.
5126              */
5127             l2hdr = ab->b_l2hdr;
5128             l2hdr->b_daddr = dev->l2ad_hand;

5129             if ((ab->b_flags & ARC_L2COMPRESS) &&
5130                 l2hdr->b_asize >= buf_compress_minsz) {
5131                 if (l2arc_compress_buf(l2hdr)) {
5132                     /*
5133                      * If compression succeeded, enable headroom
5134                      * boost on the next scan cycle.
5135                      */
5136                     *headroom_boost = B_TRUE;
5137                 }
5138             }

5139             /*
5140              * Pick up the buffer data we had previously stashed away
5141              * (and now potentially also compressed).
5142              */
5143             buf_data = l2hdr->b_tmp_cdata;
5144             buf_sz = l2hdr->b_asize;

5145             /* Compression may have squashed the buffer to zero length. */
5146             if (buf_sz != 0) {
5147                 uint64_t buf_p_sz;

5148                 wzio = zio_write_phys(pio, dev->l2ad_vdev,
5149                                     dev->l2ad_hand, buf_sz, buf_data, ZIO_CHECKSUM_OFF,
5150                                     NULL, NULL, ZIO_PRIORITY_ASYNC_WRITE,
5151                                     ZIO_FLAG_CANFAIL, B_FALSE);

5152                 DTRACE_PROBE2(l2arc__write, vdev_t *, dev->l2ad_vdev,
5153                               zio_t *, wzio);
5154                 (void) zio_nowait(wzio);

5155                 write_asize += buf_sz;
5156                 /*
5157                  * Keep the clock hand suitably device-aligned.
5158                  */
5159                 buf_p_sz = vdev_psize_to_asize(dev->l2ad_vdev, buf_sz);
5160                 write_psize += buf_p_sz;
5161                 dev->l2ad_hand += buf_p_sz;
5162             }
5163         }
5164     }

```

```

5168         l2arc_pbuff_insert(pb, pb_buflist, ab, buf_index);
5169     }
5170     ASSERT(buf_index == num_bufs);

5171     mutex_exit(&l2arc_buflist_mtx);

5173     ASSERT3U(write_asize, <=, target_sz);
5174     ARCSTAT_BUMP(arcstat_l2_writes_sent);
5175     ARCSTAT_INCR(arcstat_l2_write_bytes, write_asize);
5176     ARCSTAT_INCR(arcstat_l2_size, write_sz);
5177     ARCSTAT_INCR(arcstat_l2_asize, write_asize);
5178     vdev_space_update(dev->l2ad_vdev, write_psize, 0, 0);

5180     /* Is it time to commit this pbuff? */
5181     if (L2PBUFF_IS_FULL(pb) &&
5182         dev->l2ad_hand + L2PBUFF_ENCODED_SIZE(pb) < dev->l2ad_end) {
5183         l2arc_pbuff_commit(dev, pio, cb);
5184         l2arc_pbuff_destroy(pb);
5185         l2arc_pbuff_init(pb);
5186     }

5188     /*
5189      * Bump device hand to the device start if it is approaching the end.
5190      * l2arc_evict() will already have evicted ahead for this case.
5191      */
5192     if (dev->l2ad_hand >= (dev->l2ad_end - target_sz)) {
5193         vdev_space_update(dev->l2ad_vdev,
5194                         dev->l2ad_end - dev->l2ad_hand, 0, 0);
5195         dev->l2ad_hand = dev->l2ad_start;
5196         dev->l2ad_evict = dev->l2ad_start;
5197         dev->l2ad_first = B_FALSE;
5198     }

5200     dev->l2ad_writing = B_TRUE;
5201     (void) zio_wait(pio);
5202     dev->l2ad_writing = B_FALSE;

5204     return (write_asize);
5205 }
_____unchanged_portion_omitted_____

5467 /*
5468  * Add a vdev for use by the L2ARC. By this point the spa has already
5469  * validated the vdev and opened it. The 'rebuild' flag indicates whether
5470  * we should attempt an L2ARC persistency rebuild.
5471  * validated the vdev and opened it.
5472  */
5473 void
5474 l2arc_add_vdev(spa_t *spa, vdev_t *vd, boolean_t rebuild)
5475 {
5476     l2arc_dev_t *adddev;

5477     ASSERT(!l2arc_vdev_present(vd));

5479     /*
5480      * Create a new l2arc device entry.
5481      */
5482     adddev = kmem_zalloc(sizeof (l2arc_dev_t), KM_SLEEP);
5483     adddev->l2ad_spa = spa;
5484     adddev->l2ad_vdev = vd;
5485     adddev->l2ad_start = VDEV_LABEL_START_SIZE + L2UBERBLOCK_SIZE;
5486     adddev->l2ad_start = VDEV_LABEL_START_SIZE;
5487     adddev->l2ad_end = VDEV_LABEL_START_SIZE + vdev_get_min_asize(vd);
5488     adddev->l2ad_hand = adddev->l2ad_start;

```

```

5488     adddev->l2ad_evict = adddev->l2ad_start;
5489     adddev->l2ad_first = B_TRUE;
5490     adddev->l2ad_writing = B_FALSE;
5491     l2arc_pbuf_init(&adddev->l2ad_pbuf);

5493     /*
5494      * This is a list of all ARC buffers that are still valid on the
5495      * device.
5496      */
5497     adddev->l2ad_buflist = kmem_zalloc(sizeof (list_t), KM_SLEEP);
5498     list_create(adddev->l2ad_buflist, sizeof (arc_buf_hdr_t),
5499               offsetof(arc_buf_hdr_t, b_l2node));

5501     vdev_space_update(vd, 0, 0, adddev->l2ad_end - adddev->l2ad_hand);

5503     /*
5504      * Add device to global list
5505      */
5506     mutex_enter(&l2arc_dev_mtx);
5507     list_insert_head(l2arc_dev_list, adddev);
5508     atomic_inc_64(&l2arc_ndev);
5509     if (rebuild && l2arc_rebuild_enabled) {
5510         adddev->l2ad_rebuilding = B_TRUE;
5511         (void) thread_create(NULL, 0, l2arc_rebuild_start, adddev,
5512                             0, &p0, TS_RUN, minclsyspri);
5513     }
5514     mutex_exit(&l2arc_dev_mtx);
5515 }

5517 /*
5518  * Remove a vdev from the L2ARC.
5519  */
5520 void
5521 l2arc_remove_vdev(vdev_t *vd)
5522 {
5523     l2arc_dev_t *dev, *nextdev, *remdev = NULL;

5525     /*
5526      * Find the device by vdev
5527      */
5528     mutex_enter(&l2arc_dev_mtx);
5529     for (dev = list_head(l2arc_dev_list); dev; dev = nextdev) {
5530         nextdev = list_next(l2arc_dev_list, dev);
5531         if (vd == dev->l2ad_vdev) {
5532             remdev = dev;
5533             break;
5534         }
5535     }
5536     ASSERT(remdev != NULL);

5538     /*
5539      * Remove device from global list
5540      */
5541     list_remove(l2arc_dev_list, remdev);
5542     l2arc_dev_last = NULL; /* may have been invalidated */
5543     atomic_dec_64(&l2arc_ndev);
5544     mutex_exit(&l2arc_dev_mtx);

5546     /*
5547      * Clear all buflists and ARC references. L2ARC device flush.
5548      */
5549     l2arc_pbuf_destroy(&remdev->l2ad_pbuf);
5550     l2arc_evict(remdev, 0, B_TRUE);
5551     list_destroy(remdev->l2ad_buflist);
5552     kmem_free(remdev->l2ad_buflist, sizeof (list_t));
5553     kmem_free(remdev, sizeof (l2arc_dev_t));

```

```

5554 }
      _____
      unchanged_portion_omitted

5609 void
5610 l2arc_stop(void)
5611 {
5612     if (!(spa_mode_global & FWRITE))
5613         return;

5615     mutex_enter(&l2arc_feed_thr_lock);
5616     cv_signal(&l2arc_feed_thr_cv); /* kick thread out of startup */
5617     l2arc_thread_exit = 1;
5618     while (l2arc_thread_exit != 0)
5619         cv_wait(&l2arc_feed_thr_cv, &l2arc_feed_thr_lock);
5620     mutex_exit(&l2arc_feed_thr_lock);
5621 }

5623 /*
5624  * Main entry point for L2ARC metadata rebuilding. This function must be
5625  * called via thread_create so that the L2ARC metadata rebuild doesn't block
5626  * pool import and may proceed in parallel on all available L2ARC devices.
5627  */
5628 static void
5629 l2arc_rebuild_start(l2arc_dev_t *dev)
5630 {
5631     vdev_t *vd = dev->l2ad_vdev;
5632     spa_t *spa = dev->l2ad_spa;

5634     /* Lock out device removal. */
5635     spa_config_enter(spa, SCL_L2ARC, vd, RW_READER);
5636     ASSERT(dev->l2ad_rebuilding == B_TRUE);
5637     l2arc_rebuild(dev);
5638     dev->l2ad_rebuilding = B_FALSE;
5639     spa_config_exit(spa, SCL_L2ARC, vd);
5640     thread_exit();
5641 }

5643 /*
5644  * This function implements the actual L2ARC metadata rebuild. It:
5645  *
5646  * 1) scans the device for valid l2uberblocks
5647  * 2) if it finds a good uberblock, starts reading the pbuf chain
5648  * 3) restores each pbuf's contents to memory
5649  *
5650  * Operation stops under any of the following conditions:
5651  *
5652  * 1) We reach the end of the pbuf chain (the previous-buffer reference
5653  *    in the pbuf is zero).
5654  * 2) We encounter *any* error condition (cksum errors, io errors, looped
5655  *    pbufs, etc.).
5656  * 3) The l2arc_rebuild_timeout is hit - this is a final resort to protect
5657  *    from making severely fragmented L2ARC pbufs or slow L2ARC devices
5658  *    prevent a machine from importing the pool (and letting the
5659  *    administrator take corrective action, e.g. by kicking the misbehaving
5660  *    L2ARC device out of the pool, or by reimporting the pool with L2ARC
5661  *    rebuilding disabled).
5662  */
5663 static void
5664 l2arc_rebuild(l2arc_dev_t *dev)
5665 {
5666     int err;
5667     l2uberblock_t ub;
5668     l2pbuf_t pb;
5669     zio_t *this_io = NULL, *next_io = NULL;
5670     int64_t deadline = ddi_get_lbolt64() + hz * l2arc_rebuild_timeout;

```

```

5672     if ((err = l2arc_uberblock_find(dev, &ub)) != 0)
5673         return;
5674     L2ARC_CHK_REBUILD_TIMEOUT(deadline, /* nop */);

5676     /* set up uberblock update info */
5677     dev->l2ad_uberblock_birth = ub.ub_birth + 1;

5679     /* initial sanity checks */
5680     l2arc_pbuf_init(&pb);
5681     if ((err = l2arc_pbuf_read(dev, ub.ub_pbuf_daddr, ub.ub_pbuf_asize,
5682         ub.ub_pbuf_cksum, &pb, NULL, &this_io)) != 0) {
5683         /* root pbuf is bad, we can't do anything about that */
5684         if (err == EINVAL) {
5685             ARCSTAT_BUMP(arcstat_l2_rebuild_cksum_errors);
5686         } else {
5687             ARCSTAT_BUMP(arcstat_l2_rebuild_io_errors);
5688         }
5689         l2arc_pbuf_destroy(&pb);
5690         return;
5691     }
5692     L2ARC_CHK_REBUILD_TIMEOUT(deadline, l2arc_pbuf_destroy(&pb));

5694     dev->l2ad_evict = ub.ub_evict_tail;

5696     /* keep on chaining in new blocks */
5697     dev->l2ad_pbuf_daddr = ub.ub_pbuf_daddr;
5698     dev->l2ad_pbuf_asize = ub.ub_pbuf_asize;
5699     dev->l2ad_pbuf_cksum = ub.ub_pbuf_cksum;
5700     dev->l2ad_hand = vdev_psize_to_asize(dev->l2ad_vdev,
5701         ub.ub_pbuf_daddr + ub.ub_pbuf_asize);
5702     dev->l2ad_first = ((ub.ub_flags & L2UBLK_EVICT_FIRST) != 0);

5704     /* start the rebuild process */
5705     for (;;) {
5706         l2pbuf_t pb_prev;

5708         l2arc_pbuf_init(&pb_prev);
5709         if ((err = l2arc_pbuf_read(dev, pb.pb_prev_daddr,
5710             pb.pb_prev_asize, pb.pb_prev_cksum, &pb_prev, this_io,
5711             &next_io)) != 0) {
5712             /*
5713              * We are done reading, discard the last good buffer.
5714              */
5715             if (pb.pb_prev_daddr > dev->l2ad_hand &&
5716                 pb.pb_prev_asize > L2PBUF_HDR_SIZE) {
5717                 /* this is an error, we stopped too early */
5718                 if (err == EINVAL) {
5719                     ARCSTAT_BUMP(
5720                         arcstat_l2_rebuild_cksum_errors);
5721                 } else {
5722                     ARCSTAT_BUMP(
5723                         arcstat_l2_rebuild_io_errors);
5724                 }
5725             }
5726             l2arc_pbuf_destroy(&pb_prev);
5727             l2arc_pbuf_destroy(&pb);
5728             break;
5729         }

5731         /*
5732          * Protection against infinite loops of pbufs. This is also
5733          * our primary termination mechanism - once the buffer list
5734          * loops around our starting pbuf, we can stop.
5735          */
5736         if (pb.pb_prev_daddr >= ub.ub_pbuf_daddr &&
5737             pb_prev.pb_prev_daddr <= ub.ub_pbuf_daddr) {

```

```

5738         ARCSTAT_BUMP(arcstat_l2_rebuild_loop_errors);
5739         l2arc_pbuf_destroy(&pb);
5740         l2arc_pbuf_destroy(&pb_prev);
5741         if (next_io)
5742             l2arc_pbuf_prefetch_abort(next_io);
5743         return;
5744     }

5746     /*
5747     * Our memory pressure valve. If the system is running low
5748     * on memory, rather than swamping memory with new ARC buf
5749     * hdrs, we opt not to reconstruct the L2ARC. At this point,
5750     * however, we have already set up our L2ARC dev to chain in
5751     * new metadata pbufs, so the user may choose to re-add the
5752     * L2ARC dev at a later time to reconstruct it (when there's
5753     * less memory pressure).
5754     */
5755     if (arc_reclaim_needed()) {
5756         ARCSTAT_BUMP(arcstat_l2_rebuild_abort_lowmem);
5757         cmn_err(CE_NOTE, "System running low on memory, "
5758             "aborting L2ARC rebuild.");
5759         l2arc_pbuf_destroy(&pb);
5760         l2arc_pbuf_destroy(&pb_prev);
5761         if (next_io)
5762             l2arc_pbuf_prefetch_abort(next_io);
5763         break;
5764     }

5766     /*
5767     * Now that we know that the prev_pbuf checks out alright, we
5768     * can start reconstruction from this pbuf - we can be sure
5769     * that the L2ARC write hand has not yet reached any of our
5770     * buffers.
5771     */
5772     l2arc_pbuf_restore(dev, &pb);

5774     /* pbuf restored, continue with next one in the list */
5775     l2arc_pbuf_destroy(&pb);
5776     pb = pb_prev;
5777     this_io = next_io;
5778     next_io = NULL;

5780     L2ARC_CHK_REBUILD_TIMEOUT(deadline, l2arc_pbuf_destroy(&pb));
5781 }

5783     ARCSTAT_BUMP(arcstat_l2_rebuild_successes);
5784 }

5786 /*
5787 * Restores the payload of a pbuf to ARC. This creates empty ARC hdr entries
5788 * which only contain an l2arc_hdr, essentially restoring the buffers to
5789 * their L2ARC evicted state. This function also updates space usage on the
5790 * L2ARC vdev to make sure it tracks restored buffers.
5791 */
5792 static void
5793 l2arc_pbuf_restore(l2arc_dev_t *dev, l2pbuf_t *pb)
5794 {
5795     spa_t *spa;
5796     uint64_t guid;
5797     list_t *buflists_list;
5798     l2pbuf_buflist_t *buflist;

5800     mutex_enter(&l2arc_buflist_mtx);
5801     spa = dev->l2ad_vdev->vdev_spa;
5802     guid = spa_load_guid(spa);
5803     buflists_list = pb->pb_buflists_list;

```

```

5804     for (buflist = list_head(buflists_list); buflist;
5805         buflist = list_next(buflists_list, buflist)) {
5806         int i;
5807         uint64_t size, asize, psize;

5809         size = asize = psize = 0;
5810         for (i = 0; i < buflist->l2pbl_nbufs; i++) {
5811             l2arc_hdr_restore(&buflist->l2pbl_bufs[i], dev,
5812                             guid);
5813             size += buflist->l2pbl_bufs[i].b_size;
5814             asize += buflist->l2pbl_bufs[i].b_l2asize;
5815             psize += vdev_psize_to_asize(dev->l2ad_vdev,
5816                                         buflist->l2pbl_bufs[i].b_l2asize);
5817         }
5818         ARCSTAT_INCR(arcstat_l2_rebuild_arc_bytes, size);
5819         ARCSTAT_INCR(arcstat_l2_rebuild_l2arc_bytes, asize);
5820         ARCSTAT_INCR(arcstat_l2_rebuild_bufs, buflist->l2pbl_nbufs);
5821         vdev_space_update(dev->l2ad_vdev, psize, 0, 0);
5822     }
5823     mutex_exit(&l2arc_buflist_mtx);
5824     ARCSTAT_BUMP(arcstat_l2_rebuild_metabufs);
5825     vdev_space_update(dev->l2ad_vdev, vdev_psize_to_asize(dev->l2ad_vdev,
5826                                                         pb->pb_asize), 0, 0);
5827 }

5829 /*
5830  * Restores a single ARC buf_hdr from a pbuf. The ARC buffer is put into
5831  * a state indicating that it has been evicted to L2ARC.
5832  * The 'guid' here is the ARC-load-guid from spa_load_guid.
5833  */
5834 static void
5835 l2arc_hdr_restore(const l2pbuf_buf_t *buf, l2arc_dev_t *dev, uint64_t guid)
5836 {
5837     arc_buf_hdr_t *hdr;
5838     kmutex_t *hash_lock;
5839     dva_t dva = {buf->b_dva.dva_word[0], buf->b_dva.dva_word[1]};

5841     hdr = buf_hash_find(guid, &dva, buf->b_birth, &hash_lock);
5842     if (hdr == NULL) {
5843         /* not in cache, try to insert */
5844         arc_buf_hdr_t *exists;
5845         arc_buf_contents_t type = buf->b_contents_type;
5846         l2arc_buf_hdr_t *l2hdr;

5848         hdr = arc_buf_hdr_alloc(guid, buf->b_size, type);
5849         hdr->b_dva = buf->b_dva;
5850         hdr->b_birth = buf->b_birth;
5851         hdr->b_cksum0 = buf->b_cksum0;
5852         hdr->b_size = buf->b_size;
5853         exists = buf_hash_insert(hdr, &hash_lock);
5854         if (exists) {
5855             /* somebody beat us to the hash insert */
5856             mutex_exit(hash_lock);
5857             arc_hdr_destroy(hdr);
5858             ARCSTAT_BUMP(arcstat_l2_rebuild_bufs_precached);
5859             return;
5860         }
5861         hdr->b_flags = buf->b_flags;
5862         mutex_enter(&hdr->b_freeze_lock);
5863         ASSERT(hdr->b_freeze_cksum == NULL);
5864         hdr->b_freeze_cksum = kmem_alloc(sizeof(zio_cksum_t),
5865                                       KM_SLEEP);
5866         *hdr->b_freeze_cksum = buf->b_freeze_cksum;
5867         mutex_exit(&hdr->b_freeze_lock);

5869         /* now rebuild the l2arc entry */

```

```

5870         ASSERT(hdr->b_l2hdr == NULL);
5871         l2hdr = kmem_zalloc(sizeof(l2arc_buf_hdr_t), KM_SLEEP);
5872         l2hdr->b_dev = dev;
5873         l2hdr->b_daddr = buf->b_l2daddr;
5874         l2hdr->b_asize = buf->b_l2asize;
5875         l2hdr->b_compress = buf->b_l2compress;
5876         hdr->b_l2hdr = l2hdr;
5877         list_insert_head(dev->l2ad_buflist, hdr);
5878         ARCSTAT_INCR(arcstat_l2_size, hdr->b_size);
5879         ARCSTAT_INCR(arcstat_l2_asize, l2hdr->b_asize);

5881         arc_change_state(arc_l2c_only, hdr, hash_lock);
5882     }
5883     mutex_exit(hash_lock);
5884 }

5886 /*
5887  * Attempts to locate and read the newest valid uberblock on the provided
5888  * L2ARC device and writes it to 'ub'. On success, this function returns 0,
5889  * otherwise the appropriate error code is returned.
5890  */
5891 static int
5892 l2arc_uberblock_find(l2arc_dev_t *dev, l2uberblock_t *ub)
5893 {
5894     int err = 0;
5895     uint8_t *ub_buf;
5896     uint64_t guid;

5898     ARCSTAT_BUMP(arcstat_l2_rebuild_attempts);
5899     ub_buf = kmem_alloc(L2UBERBLOCK_SIZE, KM_SLEEP);
5900     guid = spa_guid(dev->l2ad_vdev->vdev_spa);

5902     if ((err = zio_wait(zio_read_phys(NULL, dev->l2ad_vdev,
5903                                VDEV_LABEL_START_SIZE, L2UBERBLOCK_SIZE, ub_buf,
5904                                ZIO_CHECKSUM_OFF, NULL, NULL, ZIO_PRIORITY_ASYNC_READ,
5905                                ZIO_FLAG_DONT_CACHE | ZIO_FLAG_CANFAIL |
5906                                ZIO_FLAG_DONT_PROPAGATE | ZIO_FLAG_DONT_RETRY, B_FALSE))) != 0) {
5907         ARCSTAT_BUMP(arcstat_l2_rebuild_io_errors);
5908         goto cleanup;
5909     }

5911     /*
5912      * Initial peek - does the device even have any usable uberblocks?
5913      * If not, don't bother continuing.
5914      */
5915     l2arc_uberblock_decode(ub_buf, ub);
5916     if (ub->ub_magic != L2UBERBLOCK_MAGIC || ub->ub_version == 0 ||
5917         ub->ub_version > L2UBERBLOCK_MAX_VERSION ||
5918         ub->ub_spa_guid != guid) {
5919         err = ENOTSUP;
5920         ARCSTAT_BUMP(arcstat_l2_rebuild_unsupported);
5921         goto cleanup;
5922     }

5924     /* now check to make sure that what we selected is okay */
5925     if ((err = l2arc_uberblock_verify(ub_buf, ub, guid)) != 0) {
5926         if (err == EINVAL) {
5927             ARCSTAT_BUMP(arcstat_l2_rebuild_cksum_errors);
5928         } else {
5929             ARCSTAT_BUMP(arcstat_l2_rebuild_uberblk_errors);
5930         }
5931         goto cleanup;
5932     }

5934     /* this uberblock is valid */

```



```

5936 cleanup:
5937     kmem_free(ub_buf, L2UBERBLOCK_SIZE);
5938     return (err);
5939 }

5941 /*
5942  * Reads a pbuf from storage, decodes it and validates its contents against
5943  * the provided checksum. The result is placed in 'pb'.
5944  *
5945  * The 'this_io' and 'prefetch_io' arguments are used for pbuf prefetching.
5946  * When issuing the first pbuf IO during rebuild, you should pass NULL for
5947  * 'this_io'. This function will then issue a sync IO to read the pbuf and
5948  * also issue an async IO to fetch the next pbuf in the pbuf chain. The
5949  * prefetch IO is returned in 'prefetch_io'. On subsequent calls to this
5950  * function, pass the value returned in 'prefetch_io' from the previous
5951  * call as 'this_io' and a fresh 'prefetch_io' pointer to hold the next
5952  * prefetch IO. Prior to the call, you should initialize your 'prefetch_io'
5953  * pointer to be NULL. If no prefetch IO was issued, the pointer is left
5954  * set at NULL.
5955  *
5956  * Actual prefetching takes place in two steps: a header IO (pi_hdr_io)
5957  * and the main pbuf payload IO (placed in prefetch_io). The pi_hdr_io
5958  * IO is used internally in this function to be able to 'peek' at the next
5959  * buffer's header before the main IO to read it in completely has finished.
5960  * We can then begin to issue the IO for the next buffer in the chain before
5961  * we are done reading, keeping the L2ARC device's pipeline saturated with
5962  * reads (rather than issuing an IO, waiting for it to complete, validating
5963  * the returned buffer and issuing the next one). This will make sure that
5964  * the rebuild proceeds at maximum read throughput.
5965  *
5966  * On success, this function returns 0, otherwise it returns an appropriate
5967  * error code. On error the prefetching IO is aborted and cleared before
5968  * returning from this function. Therefore, if we return 'success', the
5969  * caller can assume that we have taken care of cleanup of prefetch IOs.
5970  */
5971 static int
5972 l2arc_pbuf_read(l2arc_dev_t *dev, uint64_t daddr, uint32_t asize,
5973     zio_cksum_t cksum, l2pbuf_t *pb, zio_t *this_io, zio_t **prefetch_io)
5974 {
5975     int err = 0;
5976     uint64_t prev_pb_start;
5977     uint32_t prev_pb_asize;
5978     zio_cksum_t calc_cksum, prev_pb_cksum;
5979     l2arc_prefetch_info_t *pi = NULL;

5981     ASSERT(dev != NULL);
5982     ASSERT(pb != NULL);
5983     ASSERT(*prefetch_io == NULL);

5985     if (!l2arc_pbuf_ptr_valid(dev, daddr, asize)) {
5986         /* We could not have issued a prefetch IO for this */
5987         ASSERT(this_io == NULL);
5988         return (EINVAL);
5989     }

5991     /*
5992      * Check to see if we have issued the IO for this pbuf in a previous
5993      * run. If not, issue it now.
5994      */
5995     if (this_io == NULL)
5996         this_io = l2arc_pbuf_prefetch(dev->l2ad_vdev, daddr, asize);

5998     /* Pick up the prefetch info buffer and read its contents */
5999     pi = this_io->io_private;
6000     ASSERT(pi != NULL);
6001     ASSERT(asize <= pi->pi_bufalen);

```

```

6003     /* Wait for the IO to read this pbuf's header to complete */
6004     if ((err = zio_wait(pi->pi_hdr_io)) != 0) {
6005         (void) zio_wait(this_io);
6006         goto cleanup;
6007     }

6009     /*
6010      * Peek to see if we can start issuing the next pbuf IO immediately.
6011      * At this point, only the current pbuf's header has been read.
6012      */
6013     if (l2arc_pbuf_decode_prev_ptr(pi->pi_buf, asize, &prev_pb_start,
6014         &prev_pb_asize, &prev_pb_cksum) == 0) {
6015         uint64_t this_pb_start, this_pb_end, prev_pb_end;
6016         /* Detect malformed pbuf references and loops */
6017         this_pb_start = daddr;
6018         this_pb_end = daddr + asize;
6019         prev_pb_end = prev_pb_start + prev_pb_asize;
6020         if ((prev_pb_start >= this_pb_start && prev_pb_start <
6021             this_pb_end) ||
6022             (prev_pb_end >= this_pb_start && prev_pb_end <
6023                 this_pb_end)) {
6024             ARCSTAT_BUMP(arcstat_l2_rebuild_loop_errors);
6025             cmn_err(CE_WARN, "Looping L2ARC metadata reference "
6026                 "detected, aborting rebuild.");
6027             err = EINVAL;
6028             goto cleanup;
6029         }
6030         /*
6031          * Start issuing IO for the next pbuf early - this should
6032          * help keep the L2ARC device busy while we read, decode
6033          * and restore this pbuf.
6034          */
6035         if (l2arc_pbuf_ptr_valid(dev, prev_pb_start, prev_pb_asize))
6036             *prefetch_io = l2arc_pbuf_prefetch(dev->l2ad_vdev,
6037                 prev_pb_start, prev_pb_asize);
6038     }

6040     /* Wait for the main pbuf IO to complete */
6041     if ((err = zio_wait(this_io)) != 0)
6042         goto cleanup;

6044     /* Make sure the buffer checks out ok */
6045     fletcher_4_native(pi->pi_buf, asize, &calc_cksum);
6046     if (!ZIO_CHECKSUM_EQUAL(calc_cksum, cksum)) {
6047         err = EINVAL;
6048         goto cleanup;
6049     }

6051     /* Now we can take our time decoding this buffer */
6052     if ((err = l2arc_pbuf_decode(pi->pi_buf, asize, pb)) != 0)
6053         goto cleanup;

6055     /* This will be used in l2arc_pbuf_restore for space accounting */
6056     pb->pb_asize = asize;

6058     ARCSTAT_F_AVG(arcstat_l2_meta_avg_size, L2PBUF_ENCODED_SIZE(pb));
6059     ARCSTAT_F_AVG(arcstat_l2_meta_avg_asize, asize);
6060     ARCSTAT_F_AVG(arcstat_l2_asize_to_meta_ratio,
6061         pb->pb_payload_asz / asize);

6063 cleanup:
6064     kmem_free(pi->pi_buf, pi->pi_bufalen);
6065     pi->pi_buf = NULL;
6066     kmem_free(pi, sizeof (l2arc_prefetch_info_t));
6067     /* Abort an in-flight prefetch in case of error */

```

```

6068     if (err != 0 && *prefetch_io != NULL) {
6069         l2arc_pbuf_prefetch_abort(*prefetch_io);
6070         *prefetch_io = NULL;
6071     }
6072     return (err);
6073 }

6075 /*
6076  * Validates a pbuf device address to make sure that it can be read
6077  * from the provided L2ARC device. Returns 1 if the address is within
6078  * the device's bounds, or 0 if not.
6079  */
6080 static int
6081 l2arc_pbuf_ptr_valid(l2arc_dev_t *dev, uint64_t daddr, uint32_t asize)
6082 {
6083     uint32_t psize;
6084     uint64_t end;

6086     psize = vdev_psize_to_asize(dev->l2ad_vdev, asize);
6087     end = daddr + psize;

6089     if (end > dev->l2ad_end || asize < L2PBUF_HDR_SIZE ||
6090         asize > L2PBUF_MAX_PAYLOAD_SIZE || daddr < dev->l2ad_start ||
6091         /* check that the buffer address is correctly aligned */
6092         (daddr & (vdev_psize_to_asize(dev->l2ad_vdev,
6093             SPA_MINBLOCKSIZE) - 1)) != 0)
6094         return (0);
6095     else
6096         return (1);
6097 }

6099 /*
6100  * Starts an asynchronous read IO to read a pbuf. This is used in pbuf
6101  * reconstruction to start reading the next pbuf before we are done
6102  * decoding and reconstructing the current pbuf, to keep the l2arc device
6103  * nice and hot with read IO to process.
6104  * The returned zio will contain a newly allocated memory buffers for the IO
6105  * data which should then be freed by the caller once the zio is no longer
6106  * needed (i.e. due to it having completed). If you wish to abort this
6107  * zio, you should do so using l2arc_pbuf_prefetch_abort, which takes care
6108  * of disposing of the allocated buffers correctly.
6109  */
6110 static zio_t *
6111 l2arc_pbuf_prefetch(vdev_t *vd, uint64_t daddr, uint32_t asize)
6112 {
6113     uint32_t i, psize;
6114     zio_t *pio, *hdr_io;
6115     uint64_t hdr_rsize;
6116     uint8_t *buf;
6117     l2arc_prefetch_info_t *pinfo;

6119     psize = vdev_psize_to_asize(vd, asize);
6120     buf = kmem_alloc(psize, KM_SLEEP);
6121     pinfo = kmem_alloc(sizeof (l2arc_prefetch_info_t), KM_SLEEP);
6122     pinfo->pi_buf = buf;
6123     pinfo->pi_buflen = psize;

6125     /*
6126     * We start issuing the IO for the pbuf header early. This
6127     * allows l2arc_pbuf_read to start issuing IO for the next
6128     * buffer before the current pbuf is read in completely.
6129     */

6131     hdr_rsize = vdev_psize_to_asize(vd, SPA_MINBLOCKSIZE);
6132     ASSERT(hdr_rsize <= psize);
6133     pinfo->pi_hdr_io = zio_root(vd->vdev_spa, NULL, NULL,

```

```

6134         ZIO_FLAG_DONT_CACHE | ZIO_FLAG_CANFAIL |
6135         ZIO_FLAG_DONT_PROPAGATE | ZIO_FLAG_DONT_RETRY);
6136     hdr_io = zio_read_phys(pinfo->pi_hdr_io, vd, daddr, hdr_rsize, buf,
6137         ZIO_CHECKSUM_OFF, NULL, NULL, ZIO_PRIORITY_SYNC_READ,
6138         ZIO_FLAG_DONT_CACHE | ZIO_FLAG_CANFAIL | ZIO_FLAG_DONT_PROPAGATE |
6139         ZIO_FLAG_DONT_RETRY, B_FALSE);
6140     (void) zio_nowait(hdr_io);

6142     /*
6143     * Read in the rest of the pbuf - this can take longer than just
6144     * having a peek at the header.
6145     */
6146     pio = zio_root(vd->vdev_spa, NULL, pinfo, ZIO_FLAG_DONT_CACHE |
6147         ZIO_FLAG_CANFAIL | ZIO_FLAG_DONT_PROPAGATE |
6148         ZIO_FLAG_DONT_RETRY);
6149     for (i = hdr_rsize; i < psize; ) {
6150         uint64_t rsize = psize - i;
6151         zio_t *rzio;

6153         if (psize - i > SPA_MAXBLOCKSIZE)
6154             rsize = SPA_MAXBLOCKSIZE;
6155         ASSERT(rsize >= SPA_MINBLOCKSIZE);
6156         rzio = zio_read_phys(pio, vd, daddr + i,
6157             rsize, buf + i, ZIO_CHECKSUM_OFF, NULL, NULL,
6158             ZIO_PRIORITY_ASYNC_READ, ZIO_FLAG_DONT_CACHE |
6159             ZIO_FLAG_CANFAIL | ZIO_FLAG_DONT_PROPAGATE |
6160             ZIO_FLAG_DONT_RETRY, B_FALSE);
6161         (void) zio_nowait(rzio);
6162         i += rsize;
6163     }

6165     return (pio);
6166 }

6168 /*
6169  * Aborts a zio returned from l2arc_pbuf_prefetch and frees the data
6170  * buffers allocated for it.
6171  */
6172 static void
6173 l2arc_pbuf_prefetch_abort(zio_t *zio)
6174 {
6175     l2arc_prefetch_info_t *pi;

6177     pi = zio->io_private;
6178     ASSERT(pi != NULL);
6179     if (pi->pi_hdr_io != NULL)
6180         (void) zio_wait(pi->pi_hdr_io);
6181     (void) zio_wait(zio);
6182     kmem_free(pi->pi_buf, pi->pi_buflen);
6183     pi->pi_buf = NULL;
6184     kmem_free(pi, sizeof (l2arc_prefetch_info_t));
6185 }

6187 /*
6188  * Encodes an l2uberblock_t structure into a destination buffer. This
6189  * buffer must be at least L2UBERBLOCK_SIZE bytes long. The resulting
6190  * uberblock is always of this constant size.
6191  */
6192 static void
6193 l2arc_uberblock_encode(const l2uberblock_t *ub, uint8_t *buf)
6194 {
6195     zio_cksum_t cksum;

6197     bzero(buf, L2UBERBLOCK_SIZE);

6199 #if defined(_BIG_ENDIAN)

```

```

6200 *(uint32_t *)buf = L2UBERBLOCK_MAGIC;
6201 *(uint16_t *) (buf + 6) = L2UB_BIG_ENDIAN;
6202 #else /* !defined( BIG_ENDIAN) */
6203 *(uint32_t *)buf = BSWAP_32(L2UBERBLOCK_MAGIC);
6204 /* zero flags is ok */
6205 #endif /* !defined( BIG_ENDIAN) */
6206 buf[4] = L2UBERBLOCK_MAX_VERSION;

6208 /* rest in native byte order */
6209 *(uint64_t *) (buf + 8) = ub->ub_spa_guid;
6210 *(uint64_t *) (buf + 16) = ub->ub_birth;
6211 *(uint64_t *) (buf + 24) = ub->ub_evict_tail;
6212 *(uint64_t *) (buf + 32) = ub->ub_alloc_space;
6213 *(uint64_t *) (buf + 40) = ub->ub_pbuf_daddr;
6214 *(uint32_t *) (buf + 48) = ub->ub_pbuf_asize;
6215 bcopy(&ub->ub_pbuf_cksum, buf + 52, 32);

6217 fletcher_4_native(buf, L2UBERBLOCK_SIZE - 32, &cksum);
6218 bcopy(&cksum, buf + L2UBERBLOCK_SIZE - 32, 32);
6219 }

6221 /*
6222 * Decodes an l2uberblock_t from an on-disk representation. Please note
6223 * that this function does not perform any uberblock validation and
6224 * checksumming - call l2arc_uberblock_verify() for that.
6225 */
6226 static void
6227 l2arc_uberblock_decode(const uint8_t *buf, l2uberblock_t *ub)
6228 {
6229     boolean_t bswap_needed;

6231     /* these always come in big endian */
6232     #if defined( BIG_ENDIAN)
6233     ub->ub_magic = *(uint32_t *)buf;
6234     ub->ub_flags = *(uint16_t *) (buf + 6);
6235     bswap_needed = ((ub->ub_flags & L2UBLK_BIG_ENDIAN) != 1);
6236     #else /* !defined( BIG_ENDIAN) */
6237     ub->ub_magic = BSWAP_32(*(uint32_t *)buf);
6238     ub->ub_flags = BSWAP_16(*(uint16_t *) (buf + 6));
6239     bswap_needed = ((ub->ub_flags & L2UBLK_BIG_ENDIAN) != 0);
6240     #endif /* !defined( BIG_ENDIAN) */
6241     ub->ub_version = buf[4];

6243     ub->ub_spa_guid = *(uint64_t *) (buf + 8);
6244     ub->ub_birth = *(uint64_t *) (buf + 16);
6245     ub->ub_evict_tail = *(uint64_t *) (buf + 24);
6246     ub->ub_alloc_space = *(uint64_t *) (buf + 32);
6247     ub->ub_pbuf_daddr = *(uint64_t *) (buf + 40);
6248     ub->ub_pbuf_asize = *(uint32_t *) (buf + 48);
6249     bcopy(buf + 52, &ub->ub_pbuf_cksum, 36);
6250     bcopy(buf + L2UBERBLOCK_SIZE - 32, &ub->ub_cksum, 32);

6252     /* swap the rest if endianness doesn't match us */
6253     if (bswap_needed) {
6254         ub->ub_spa_guid = BSWAP_64(ub->ub_spa_guid);
6255         ub->ub_birth = BSWAP_64(ub->ub_birth);
6256         ub->ub_evict_tail = BSWAP_64(ub->ub_evict_tail);
6257         ub->ub_alloc_space = BSWAP_64(ub->ub_alloc_space);
6258         ub->ub_pbuf_daddr = BSWAP_64(ub->ub_pbuf_daddr);
6259         ub->ub_pbuf_asize = BSWAP_32(ub->ub_pbuf_asize);
6260         ZIO_CHECKSUM_BSWAP(&ub->ub_pbuf_cksum);
6261         ZIO_CHECKSUM_BSWAP(&ub->ub_cksum);
6262     }
6263 }

6265 /*

```

```

6266 * Verifies whether a decoded uberblock (via l2arc_uberblock_decode()) is
6267 * valid and matches its cksum.
6268 */
6269 static int
6270 l2arc_uberblock_verify(const uint8_t *buf, const l2uberblock_t *ub,
6271     uint64_t guid)
6272 {
6273     zio_cksum_t cksum;

6275     if (ub->ub_magic != L2UBERBLOCK_MAGIC ||
6276         ub->ub_version == 0 || ub->ub_version > L2UBERBLOCK_MAX_VERSION)
6277         /*
6278          * bad magic or invalid version => persistent l2arc not
6279          * supported
6280          */
6281         return (ENOTSUP);

6283     if (ub->ub_spa_guid != guid)
6284         /* this l2arc dev isn't ours */
6285         return (EINVAL);

6287     fletcher_4_native(buf, L2UBERBLOCK_SIZE - 32, &cksum);
6288     if (!ZIO_CHECKSUM_EQUAL(cksum, ub->ub_cksum))
6289         /* bad checksum, corrupt uberblock */
6290         return (EINVAL);

6292     return (0);
6293 }

6295 /*
6296 * Schedules a zio to update the uberblock on an l2arc device. The zio is
6297 * initiated as a child of 'pio' and 'cb' is filled with the information
6298 * needed to free the uberblock data buffer after writing.
6299 */
6300 static void
6301 l2arc_uberblock_update(l2arc_dev_t *dev, zio_t *pio, l2arc_write_callback_t *cb)
6302 {
6303     uint8_t *ub_buf;
6304     l2uberblock_t ub;
6305     zio_t *wzio;
6306     vdev_stat_t st;

6308     ASSERT(cb->l2wcb_ub_buf == NULL);
6309     vdev_get_stats(dev->l2ad_vdev, &st);

6311     bzero(&ub, sizeof(ub));
6312     ub.ub_spa_guid = spa_guid(dev->l2ad_vdev->vdev_spa);
6313     ub.ub_birth = dev->l2ad_uberblock_birth++;
6314     ub.ub_evict_tail = dev->l2ad_evict;
6315     ub.ub_alloc_space = st.vs_alloc;
6316     ub.ub_pbuf_daddr = dev->l2ad_pbuf_daddr;
6317     ub.ub_pbuf_asize = dev->l2ad_pbuf_asize;
6318     ub.ub_pbuf_cksum = dev->l2ad_pbuf_cksum;
6319     if (dev->l2ad_first)
6320         ub.ub_flags |= L2UBLK_EVICT_FIRST;

6322     ub_buf = kmem_alloc(L2UBERBLOCK_SIZE, KM_SLEEP);
6323     cb->l2wcb_ub_buf = ub_buf;
6324     l2arc_uberblock_encode(&ub, ub_buf);
6325     wzio = zio_write_phys(pio, dev->l2ad_vdev, VDEV_LABEL_START_SIZE,
6326         L2UBERBLOCK_SIZE, ub_buf, ZIO_CHECKSUM_OFF, NULL, NULL,
6327         ZIO_PRIORITY_ASYNC_WRITE, ZIO_FLAG_CANFAIL, B_FALSE);
6328     DTRACE_PROBE2(l2arc__write, vdev_t *, dev->l2ad_vdev,
6329         zio_t *, wzio);
6330     (void) zio_nowait(wzio);
6331 }

```

```

6333 /*
6334  * Encodes a l2pbuf_t structure into the portable on-disk format. The
6335  * 'buf' buffer must be suitably sized to hold the entire uncompressed
6336  * structure (use L2PBUF_ENCODED_SIZE()). If requested, this function
6337  * also compresses the buffer.
6338  *
6339  * The return value is the length of the resulting encoded pbuf structure.
6340  * This can be either equal to L2PBUF_ENCODED_SIZE(pb) if no compression
6341  * was applied, or smaller if compression was applied. In either case,
6342  * prior to writing to disk, the caller must suitably pad the output
6343  * buffer so that it is aligned on a multiple of the underlying storage
6344  * system's block size.
6345  */
6346 static uint32_t
6347 l2arc_pbuf_encode(l2pbuf_t *pb, uint8_t *buf, uint32_t buflen)
6348 {
6349     uint16_t flags = 0;
6350     uint8_t *dst_buf;
6351     uint32_t enclen;
6352     l2pbuf_buflist_t *buflist;
6353
6354     enclen = L2PBUF_ENCODED_SIZE(pb);
6355     ASSERT(buflen >= enclen);
6356     bzero(buf, enclen);
6357
6358     /* non-header portions of pbufs are in native byte order */
6359     *(uint64_t *) (buf + 8) = pb->pb_prev_daddr;
6360     *(uint32_t *) (buf + 16) = pb->pb_prev_asize;
6361     bcopy(&pb->pb_prev_cksum, buf + 20, 32);
6362     *(uint32_t *) (buf + 52) = enclen - L2PBUF_HDR_SIZE;
6363
6364     /* first we encode the buflists uncompressed */
6365     dst_buf = buf + L2PBUF_HDR_SIZE;
6366     for (buflist = list_head(pb->pb_buflists_list); buflist;
6367          buflist = list_next(pb->pb_buflists_list, buflist)) {
6368         int i;
6369
6370         ASSERT(buflist->l2pbl_nbufs != 0);
6371         for (i = 0; i < buflist->l2pbl_nbufs; i++) {
6372             l2pbuf_buf_t *pbl_buf = &buflist->l2pbl_bufs[i];
6373
6374             ASSERT(pbl_buf->b_size != 0);
6375             *(uint64_t *) dst_buf = pbl_buf->b_dva.dva_word[0];
6376             *(uint64_t *) (dst_buf + 8) = pbl_buf->b_dva.dva_word[1];
6377             *(uint64_t *) (dst_buf + 16) = pbl_buf->b_birth;
6378             *(uint64_t *) (dst_buf + 24) = pbl_buf->b_cksum0;
6379             bcopy(&pbl_buf->b_freeze_cksum, dst_buf + 32, 32);
6380             *(uint32_t *) (dst_buf + 64) = pbl_buf->b_size;
6381             *(uint64_t *) (dst_buf + 68) = pbl_buf->b_l2daddr;
6382             *(uint32_t *) (dst_buf + 76) = pbl_buf->b_l2asize;
6383             dst_buf[80] = pbl_buf->b_l2compress;
6384             dst_buf[81] = pbl_buf->b_contents_type;
6385             *(uint32_t *) (dst_buf + 84) = pbl_buf->b_flags;
6386             dst_buf += L2PBUF_BUF_SIZE;
6387         }
6388     }
6389     ASSERT((uint32_t)(dst_buf - buf) == enclen);
6390
6391     /* and then compress them if necessary */
6392     if (enclen >= l2arc_pbuf_compress_minsz) {
6393         uint8_t *cbuf;
6394         size_t slen, clen;
6395
6396         slen = l2arc_pbuf_items_encoded_size(pb);
6397         cbuf = kmem_alloc(slen, KM_SLEEP);

```

```

6398         clen = lz4_compress(buf + L2PBUF_HDR_SIZE, cbuf, slen, slen, 0);
6399         ASSERT(clen != 0);
6400         if (clen < slen) {
6401             bcopy(cbuf, buf + L2PBUF_HDR_SIZE, clen);
6402             flags |= L2PBUF_COMPRESSED;
6403             /* zero out the rest of the input buffer */
6404             bzero(buf + L2PBUF_HDR_SIZE + clen,
6405                  buflen - (L2PBUF_HDR_SIZE + clen));
6406             /* adjust our buffer length now that it's shortened */
6407             enclen = L2PBUF_HDR_SIZE + clen;
6408         }
6409         kmem_free(cbuf, slen);
6410     }
6411
6412     /* the header goes last since 'flags' may change due to compression */
6413     #if defined(_BIG_ENDIAN)
6414     *(uint32_t *) buf = L2PBUF_MAGIC;
6415     flags |= L2PBUF_BIG_ENDIAN;
6416     *(uint16_t *) (buf + 6) = flags;
6417     #else /* !defined(_BIG_ENDIAN) */
6418     *(uint32_t *) buf = BSWAP_32(L2PBUF_MAGIC);
6419     *(uint16_t *) (buf + 6) = BSWAP_16(flags);
6420     #endif /* !defined(_BIG_ENDIAN) */
6421     buf[4] = L2PBUF_MAX_VERSION;
6422
6423     return (enclen);
6424 }
6425
6426 /*
6427  * Decodes a stored l2pbuf_t structure previously encoded using
6428  * l2arc_pbuf_encode. The source buffer is not modified. The passed pbuf
6429  * must be initialized by l2arc_pbuf_init by the caller beforehand, but
6430  * must not have been used to store any buffers yet.
6431  *
6432  * Please note that we don't do checksum verification here, as we don't
6433  * know our own checksum (that's know by the previous block in the linked
6434  * list, or by the uberblock). This should be performed by the caller
6435  * prior to calling l2arc_pbuf_decode.
6436  */
6437 static int
6438 l2arc_pbuf_decode(uint8_t *input_buf, uint32_t buflen, l2pbuf_t *pb)
6439 {
6440     boolean_t bswap_needed;
6441     uint32_t payload_sz, payload_asz;
6442     uint8_t *src_bufs;
6443     l2pbuf_buflist_t *buflist;
6444     int i, nbufs;
6445
6446     ASSERT(input_buf != NULL);
6447     ASSERT(pb != NULL);
6448     ASSERT(pb->pb_version != 0);
6449     ASSERT(pb->pb_nbuflists == 0);
6450
6451     /* no valid buffer can be this small */
6452     if (buflen < L2PBUF_HDR_SIZE)
6453         return (EINVAL);
6454
6455     /* these always come in big endian */
6456     #if defined(_BIG_ENDIAN)
6457     pb->pb_magic = *(uint32_t *) input_buf;
6458     pb->pb_flags = *(uint16_t *) (input_buf + 6);
6459     bswap_needed = ((pb->pb_flags & L2PBUF_BIG_ENDIAN) != 1);
6460     #else /* !defined(_BIG_ENDIAN) */
6461     pb->pb_magic = BSWAP_32(*(uint32_t *) input_buf);
6462     pb->pb_flags = BSWAP_16(*(uint16_t *) (input_buf + 6));
6463     bswap_needed = ((pb->pb_flags & L2PBUF_BIG_ENDIAN) != 0);

```

```

6464 #endif /* !defined(_BIG_ENDIAN) */
6465 pb->pb_version = input_buf[4];

6467 if (pb->pb_magic != L2PBUF_MAGIC || pb->pb_version == 0)
6468     return (EINVAL);
6469 if (pb->pb_version > L2PBUF_MAX_VERSION)
6470     return (ENOTSUP);

6472 /* remainder of pbuf may need bswap'ping */
6473 pb->pb_prev_daddr = *(uint64_t *)(input_buf + 8);
6474 pb->pb_prev_asize = *(uint64_t *)(input_buf + 16);
6475 bcopy(input_buf + 20, &pb->pb_prev_cksum, 32);
6476 payload_sz = *(uint32_t *)(input_buf + 52);
6477 payload_asz = buflen - L2PBUF_HDR_SIZE;

6479 if (bswap_needed) {
6480     pb->pb_prev_daddr = BSWAP_64(pb->pb_prev_daddr);
6481     pb->pb_prev_asize = BSWAP_64(pb->pb_prev_asize);
6482     ZIO_CHECKSUM_BSWAP(&pb->pb_prev_cksum);
6483     payload_sz = BSWAP_32(payload_sz);
6484 }

6486 /* check for sensible buffer allocation limits */
6487 if (((pb->pb_flags & L2PBUF_COMPRESSED) && payload_sz <= payload_asz) ||
6488     (payload_sz > L2PBUF_MAX_PAYLOAD_SIZE) ||
6489     (payload_sz % L2PBUF_BUF_SIZE) != 0 || payload_sz == 0)
6490     return (EINVAL);
6491 nbufs = payload_sz / L2PBUF_BUF_SIZE;

6493 /* decompression might be needed */
6494 if (pb->pb_flags & L2PBUF_COMPRESSED) {
6495     src_bufs = kmem_alloc(payload_sz, KM_SLEEP);
6496     if (lz4_decompress(input_buf + L2PBUF_HDR_SIZE, src_bufs,
6497         payload_asz, payload_sz, 0) != 0) {
6498         kmem_free(src_bufs, payload_sz);
6499         return (EINVAL);
6500     }
6501 } else {
6502     src_bufs = input_buf + L2PBUF_HDR_SIZE;
6503 }

6505 /* Decode individual pbuf items from our source buffer. */
6506 buflist = l2arc_pbuf_buflist_alloc(pb, nbufs);
6507 for (i = 0; i < nbufs; i++) {
6508     l2pbuf_buf_t *pbl_buf = &buflist->l2pbl_bufs[i];
6509     const uint8_t *src = src_bufs + i * L2PBUF_BUF_SIZE;

6511     pbl_buf->b_dva.dva_word[0] = *(uint64_t *)src;
6512     pbl_buf->b_dva.dva_word[1] = *(uint64_t *)src + 8;
6513     pbl_buf->b_birth = *(uint64_t *)src + 16;
6514     pbl_buf->b_cksum0 = *(uint64_t *)src + 24;
6515     bcopy(src + 32, &pbl_buf->b_freeze_cksum, 32);
6516     pbl_buf->b_size = *(uint32_t *)src + 64;
6517     pbl_buf->b_l2daddr = *(uint64_t *)src + 68;
6518     pbl_buf->b_l2asize = *(uint32_t *)src + 76;
6519     pbl_buf->b_l2compress = src[80];
6520     pbl_buf->b_contents_type = src[81];
6521     pbl_buf->b_flags = *(uint32_t *)src + 84;

6523     if (bswap_needed) {
6524         pbl_buf->b_dva.dva_word[0] =
6525             BSWAP_64(pbl_buf->b_dva.dva_word[0]);
6526         pbl_buf->b_dva.dva_word[1] =
6527             BSWAP_64(pbl_buf->b_dva.dva_word[1]);
6528         pbl_buf->b_birth = BSWAP_64(pbl_buf->b_birth);
6529         pbl_buf->b_cksum0 = BSWAP_64(pbl_buf->b_cksum0);

```

```

6530     ZIO_CHECKSUM_BSWAP(&pbl_buf->b_freeze_cksum);
6531     pbl_buf->b_size = BSWAP_32(pbl_buf->b_size);
6532     pbl_buf->b_l2daddr = BSWAP_64(pbl_buf->b_l2daddr);
6533     pbl_buf->b_l2asize = BSWAP_32(pbl_buf->b_l2asize);
6534     pbl_buf->b_flags = BSWAP_32(pbl_buf->b_flags);
6535 }

6537     pb->pb_payload_asz += pbl_buf->b_l2asize;
6538 }

6540 if (pb->pb_flags & L2PBUF_COMPRESSED)
6541     kmem_free(src_bufs, payload_sz);

6543     return (0);
6544 }

6546 /*
6547  * Decodes the previous buffer pointer encoded in a pbuf. This is used
6548  * during L2ARC reconstruction to "peek" at the next buffer and start
6549  * issuing IO to fetch it early, before decoding of the current buffer
6550  * is done (which can take time due to decompression).
6551  * Returns 0 on success (and fills in the return parameters 'daddr',
6552  * 'asize' and 'cksum' with the info of the previous pbuf), and an errno
6553  * on error.
6554  */
6555 static int
6556 l2arc_pbuf_decode_prev_ptr(const uint8_t *buf, size_t buflen, uint64_t *daddr,
6557     uint32_t *asize, zio_cksum_t *cksum)
6558 {
6559     boolean_t bswap_needed;
6560     uint16_t version, flags;
6561     uint32_t magic;

6563     ASSERT(buf != NULL);

6565     /* no valid buffer can be this small */
6566     if (buflen <= L2PBUF_HDR_SIZE)
6567         return (EINVAL);

6569     /* these always come in big endian */
6570 #if defined(_BIG_ENDIAN)
6571     magic = *(uint32_t *)buf;
6572     flags = *(uint16_t *)buf + 6;
6573     bswap_needed = ((flags & L2PBUF_BIG_ENDIAN) != 1);
6574 #else
6575     /* !defined(_BIG_ENDIAN) */
6576     magic = BSWAP_32(*(uint32_t *)buf);
6577     flags = BSWAP_16(*(uint16_t *)buf + 6);
6578     bswap_needed = ((flags & L2PBUF_BIG_ENDIAN) != 0);
6579 #endif /* !defined(_BIG_ENDIAN) */
6580     version = buf[4];

6581     if (magic != L2PBUF_MAGIC || version == 0)
6582         return (EINVAL);
6583     if (version > L2PBUF_MAX_VERSION)
6584         return (ENOTSUP);

6586     *daddr = *(uint64_t *)buf + 4;
6587     *asize = *(uint64_t *)buf + 12;
6588     bcopy(buf + 16, cksum, 32);

6590     if (bswap_needed) {
6591         *daddr = BSWAP_64(*daddr);
6592         *asize = BSWAP_64(*asize);
6593         ZIO_CHECKSUM_BSWAP(cksum);
6594     }

```

```

6596     return (0);
6597 }

6599 /*
6600 * Initializes a pbuf structure into a clean state. All version and flags
6601 * fields are filled in as appropriate for this architecture.
6602 * If the structure was used before, first call l2arc_pbuf_destroy on it,
6603 * as this function assumes the structure is uninitialized.
6604 */
6605 static void
6606 l2arc_pbuf_init(l2pbuf_t *pb)
6607 {
6608     bzero(pb, sizeof (l2pbuf_t));
6609     pb->pb_version = L2PBUF_MAX_VERSION;
6610 #if defined(_BIG_ENDIAN)
6611     pb->pb_flags |= L2PB_BIG_ENDIAN;
6612 #endif
6613     pb->pb_buflists_list = kmem_zalloc(sizeof (list_t), KM_SLEEP);
6614     list_create(pb->pb_buflists_list, sizeof (l2pbuf_buflist_t),
6615         offsetof(l2pbuf_buflist_t, l2pbl_node));
6616 }

6618 /*
6619 * Destroys a pbuf structure and puts it into a clean state ready to be
6620 * initialized by l2arc_pbuf_init. All buflists created by
6621 * l2arc_pbuf_buflist_alloc are released as well.
6622 */
6623 static void
6624 l2arc_pbuf_destroy(l2pbuf_t *pb)
6625 {
6626     list_t *buflist_list = pb->pb_buflists_list;
6627     l2pbuf_buflist_t *buflist;

6629     while ((buflist = list_head(buflist_list)) != NULL) {
6630         ASSERT(buflist->l2pbl_nbufs > 0);
6631         kmem_free(buflist->l2pbl_bufs, sizeof (l2pbuf_buf_t) *
6632             buflist->l2pbl_nbufs);
6633         list_remove(buflist_list, buflist);
6634         kmem_free(buflist, sizeof (l2pbuf_buflist_t));
6635     }
6636     pb->pb_nbuflists = 0;
6637     list_destroy(pb->pb_buflists_list);
6638     kmem_free(pb->pb_buflists_list, sizeof (list_t));
6639     bzero(pb, sizeof (l2pbuf_t));
6640 }

6642 /*
6643 * Allocates a new buflist inside of a pbuf, which can hold up to 'nbufs'
6644 * buffers. This is used during the buffer write cycle - each cycle allocates
6645 * a new buflist and fills it with buffers it writes. Then, when the pbuf
6646 * reaches its buflist limit, it is committed to stable storage.
6647 */
6648 static l2pbuf_buflist_t *
6649 l2arc_pbuf_buflist_alloc(l2pbuf_t *pb, int nbufs)
6650 {
6651     l2pbuf_buflist_t *buflist;

6653     ASSERT(pb->pb_buflists_list != NULL);
6654     buflist = kmem_zalloc(sizeof (l2pbuf_buflist_t), KM_SLEEP);
6655     buflist->l2pbl_nbufs = nbufs;
6656     buflist->l2pbl_bufs = kmem_zalloc(sizeof (l2pbuf_buf_t) * nbufs,
6657         KM_SLEEP);
6658     list_insert_tail(pb->pb_buflists_list, buflist);
6659     pb->pb_nbuflists++;

6661     return (buflist);

```

```

6662 }

6664 /*
6665 * Inserts ARC buffer 'ab' into the pbuf 'pb' buflist 'pbl' at index 'idx'.
6666 * The buffer being inserted must be present in L2ARC.
6667 */
6668 static void
6669 l2arc_pbuflist_insert(l2pbuf_t *pb, l2pbuf_buflist_t *pbl,
6670     const arc_buf_hdr_t *ab, int index)
6671 {
6672     l2pbuf_buf_t *pb_buf;
6673     const l2arc_buf_hdr_t *l2hdr;

6675     l2hdr = ab->b_l2hdr;
6676     ASSERT(l2hdr != NULL);
6677     ASSERT(pbl->l2pbl_nbufs > index);

6679     pb_buf = &pbl->l2pbl_bufs[index];
6680     pb_buf->b_dva = ab->b_dva;
6681     pb_buf->b_birth = ab->b_birth;
6682     pb_buf->b_cksum0 = ab->b_cksum0;
6683     pb_buf->b_freeze_cksum = *ab->b_freeze_cksum;
6684     pb_buf->b_size = ab->b_size;
6685     pb_buf->b_l2daddr = l2hdr->b_daddr;
6686     pb_buf->b_l2asize = l2hdr->b_asize;
6687     pb_buf->b_l2compress = l2hdr->b_compress;
6688     pb_buf->b_contents_type = ab->b_type;
6689     pb_buf->b_flags = ab->b_flags & L2ARC_PERSIST_FLAGS;
6690     pb->pb_payload_asz += l2hdr->b_asize;
6691 }

6693 /*
6694 * Commits a pbuf to stable storage. This routine is invoked when writing
6695 * ARC buffers to an L2ARC device. When the pbuf associated with the device
6696 * has reached its limits (either in size or in number of writes), it is
6697 * scheduled here for writing.
6698 * This function allocates some memory to temporarily hold the serialized
6699 * buffer to be written. This is then released in l2arc_write_done.
6700 */
6701 static void
6702 l2arc_pbuf_commit(l2arc_dev_t *dev, zio_t *pio, l2arc_write_callback_t *cb)
6703 {
6704     l2pbuf_t *pb = &dev->l2ad_pbuf;
6705     uint64_t i, est_encsize, bufsize, encsize, io_size;
6706     uint8_t *pb_buf;

6708     pb->pb_prev_daddr = dev->l2ad_pbuf_daddr;
6709     pb->pb_prev_asize = dev->l2ad_pbuf_asize;
6710     pb->pb_prev_cksum = dev->l2ad_pbuf_cksum;

6712     est_encsize = L2PBUF_ENCODED_SIZE(pb);
6713     bufsize = vdev_psize_to_asize(dev->l2ad_vdev, est_encsize);
6714     pb_buf = kmem_zalloc(bufsize, KM_SLEEP);
6715     encsize = l2arc_pbuf_encode(pb, pb_buf, bufsize);
6716     cb->l2wcb_pbuf = pb_buf;
6717     cb->l2wcb_pbuf_size = bufsize;

6719     dev->l2ad_pbuf_daddr = dev->l2ad_hand;
6720     dev->l2ad_pbuf_asize = encsize;
6721     fletcher_4_native(pb_buf, encsize, &dev->l2ad_pbuf_cksum);

6723     io_size = vdev_psize_to_asize(dev->l2ad_vdev, encsize);
6724     for (i = 0; i < io_size; ) {
6725         zio_t *wzio;
6726         uint64_t wsize = io_size - i;

```

```
6728     if (wsize > SPA_MAXBLOCKSIZE)
6729         wsize = SPA_MAXBLOCKSIZE;
6730     ASSERT(wsize >= SPA_MINBLOCKSIZE);
6731     wzio = zio_write_phys(pio, dev->l2ad_vdev, dev->l2ad_hand + i,
6732         wsize, pb_buf + i, ZIO_CHECKSUM_OFF, NULL, NULL,
6733         ZIO_PRIORITY_ASYNC_WRITE, ZIO_FLAG_CANFAIL, B_FALSE);
6734     DTRACE_PROBE2(l2arc__write, vdev_t *, dev->l2ad_vdev,
6735         zio_t *, wzio);
6736     (void) zio_nowait(wzio);
6737     i += wsize;
6738 }
6740 dev->l2ad_hand += io_size;
6741 vdev_space_update(dev->l2ad_vdev, io_size, 0, 0);
6742 l2arc_uberblock_update(dev, pio, cb);
6744     ARCSTAT_INCR(arcstat_l2_write_bytes, io_size);
6745     ARCSTAT_BUMP(arcstat_l2_meta_writes);
6746     ARCSTAT_F_AVG(arcstat_l2_meta_avg_size, est_encsize);
6747     ARCSTAT_F_AVG(arcstat_l2_meta_avg_asize, encsize);
6748     ARCSTAT_F_AVG(arcstat_l2_asize_to_meta_ratio,
6749         pb->pb_payload_asz / encsize);
6750 }
6752 /*
6753  * Returns the number of bytes occupied by the payload buffer items of
6754  * a pbuf in portable (on-disk) encoded form, i.e. the bytes following
6755  * L2PBUF_HDR_SIZE.
6756  */
6757 static uint32_t
6758 l2arc_pbuf_items_encoded_size(l2pbuf_t *pb)
6759 {
6760     uint32_t size = 0;
6761     l2pbuf_buflist_t *buflist;
6763     for (buflist = list_head(pb->pb_buflists_list); buflist != NULL;
6764         buflist = list_next(pb->pb_buflists_list, buflist))
6765         size += L2PBUF_BUF_SIZE * buflist->l2pbl_nbufs;
6767     return (size);
6768 }
unchanged_portion_omitted
```

```

*****
175566 Mon Aug 5 21:59:08 2013
new/usr/src/uts/common/fs/zfs/spa.c
3525 Persistent L2ARC
*****
_____unchanged_portion_omitted_____

1403 /*
1404 * Load (or re-load) the current list of vdevs describing the active l2cache for
1405 * this pool. When this is called, we have some form of basic information in
1406 * 'spa_l2cache.sav_config'. We parse this into vdevs, try to open them, and
1407 * then re-generate a more complete list including status information.
1408 * Devices which are already active have their details maintained, and are
1409 * not re-opened.
1410 */
1411 static void
1412 spa_load_l2cache(spa_t *spa)
1413 {
1414     nvlist_t **l2cache;
1415     uint_t nl2cache;
1416     int i, j, oldnvdevs;
1417     uint64_t guid;
1418     vdev_t *vd, **oldvdevs, **newvdevs;
1419     spa_aux_vdev_t *sav = &spa->spa_l2cache;

1421     ASSERT(spa_config_held(spa, SCL_ALL, RW_WRITER) == SCL_ALL);

1423     if (sav->sav_config != NULL) {
1424         VERIFY(nvlist_lookup_nvlist_array(sav->sav_config,
1425             ZPOOL_CONFIG_L2CACHE, &l2cache, &nl2cache) == 0);
1426         newvdevs = kmem_alloc(nl2cache * sizeof(void *), KM_SLEEP);
1427     } else {
1428         nl2cache = 0;
1429         newvdevs = NULL;
1430     }

1432     oldvdevs = sav->sav_vdevs;
1433     oldnvdevs = sav->sav_count;
1434     sav->sav_vdevs = NULL;
1435     sav->sav_count = 0;

1437     /*
1438     * Process new nvlist of vdevs.
1439     */
1440     for (i = 0; i < nl2cache; i++) {
1441         VERIFY(nvlist_lookup_uint64(l2cache[i], ZPOOL_CONFIG_GUID,
1442             &guid) == 0);

1444         newvdevs[i] = NULL;
1445         for (j = 0; j < oldnvdevs; j++) {
1446             vd = oldvdevs[j];
1447             if (vd != NULL && guid == vd->vdev_guid) {
1448                 /*
1449                  * Retain previous vdev for add/remove ops.
1450                  */
1451                 newvdevs[i] = vd;
1452                 oldvdevs[j] = NULL;
1453                 break;
1454             }
1455         }

1457         if (newvdevs[i] == NULL) {
1458             /*
1459             * Create new vdev
1460             */
1461             VERIFY(spa_config_parse(spa, &vd, l2cache[i], NULL, 0,

```

```

1462         VDEV_ALLOC_L2CACHE) == 0);
1463         ASSERT(vd != NULL);
1464         newvdevs[i] = vd;

1466         /*
1467         * Commit this vdev as an l2cache device,
1468         * even if it fails to open.
1469         */
1470         spa_l2cache_add(vd);

1472         vd->vdev_top = vd;
1473         vd->vdev_aux = sav;

1475         spa_l2cache_activate(vd);

1477         if (vdev_open(vd) != 0)
1478             continue;

1480         (void) vdev_validate_aux(vd);

1482         if (!vdev_is_dead(vd)) {
1483             boolean_t persist = B_FALSE;

1485             if (spa->spa_load_state != SPA_LOAD_TRYIMPORT) {
1486                 /*
1487                 * Only allow to do the L2ARC rebuild
1488                 * when not doing a spa try-load.
1489                 */
1490                 (void) nvlist_lookup_boolean_value(
1491                     l2cache[i],
1492                     ZPOOL_CONFIG_L2CACHE_PERSISTENT,
1493                     &persist);
1482                 if (!vdev_is_dead(vd))
1483                     l2arc_add_vdev(spa, vd);
1494             }
1495             l2arc_add_vdev(spa, vd, persist);
1496         }
1497     }
1498 }

1500     /*
1501     * Purge vdevs that were dropped
1502     */
1503     for (i = 0; i < oldnvdevs; i++) {
1504         uint64_t pool;

1506         vd = oldvdevs[i];
1507         if (vd != NULL) {
1508             ASSERT(vd->vdev_isl2cache);

1510             if (spa_l2cache_exists(vd->vdev_guid, &pool) &&
1511                 pool != 0ULL && l2arc_vdev_present(vd))
1512                 l2arc_remove_vdev(vd);
1513             vdev_clear_stats(vd);
1514             vdev_free(vd);
1515         }
1516     }

1518     if (oldvdevs)
1519         kmem_free(oldvdevs, oldnvdevs * sizeof(void *));

1521     if (sav->sav_config == NULL)
1522         goto out;

1524     sav->sav_vdevs = newvdevs;
1525     sav->sav_count = (int)nl2cache;

```



```

1527      /*
1528       * Recompute the stashed list of l2cache devices, with status
1529       * information this time.
1530       */
1531      VERIFY(nvlist_remove(sav->sav_config, ZPOOL_CONFIG_L2CACHE,
1532          DATA_TYPE_NVLIST_ARRAY) == 0);

1534      l2cache = kmem_alloc(sav->sav_count * sizeof(void *), KM_SLEEP);
1535      for (i = 0; i < sav->sav_count; i++)
1536          l2cache[i] = vdev_config_generate(spa,
1537              sav->sav_vdevs[i], B_TRUE, VDEV_CONFIG_L2CACHE);
1538      VERIFY(nvlist_add_nvlist_array(sav->sav_config,
1539          ZPOOL_CONFIG_L2CACHE, l2cache, sav->sav_count) == 0);
1540  out:
1541      for (i = 0; i < sav->sav_count; i++)
1542          nvlist_free(l2cache[i]);
1543      if (sav->sav_count)
1544          kmem_free(l2cache, sav->sav_count * sizeof(void *));
1545  }
  
```

unchanged portion omitted

```

3801 #endif

3803 /*
3804  * Import a non-root pool into the system.
3805  */
3806  int
3807  spa_import(const char *pool, nvlist_t *config, nvlist_t *props, uint64_t flags)
3808  {
3809      spa_t *spa;
3810      char *altroot = NULL;
3811      spa_load_state_t state = SPA_LOAD_IMPORT;
3812      zpool_rewind_policy_t policy;
3813      uint64_t mode = spa_mode_global;
3814      uint64_t readonly = B_FALSE;
3815      int error;
3816      nvlist_t *nvroot;
3817      nvlist_t **spares, **l2cache;
3818      uint_t nspares, nl2cache;

3820      /*
3821       * If a pool with this name exists, return failure.
3822       */
3823      mutex_enter(&spa_namespace_lock);
3824      if (spa_lookup(pool) != NULL) {
3825          mutex_exit(&spa_namespace_lock);
3826          return (SET_ERROR(EEXIST));
3827      }

3829      /*
3830       * Create and initialize the spa structure.
3831       */
3832      (void) nvlist_lookup_string(props,
3833          zpool_prop_to_name(ZPOOL_PROP_ALTROOT), &altroot);
3834      (void) nvlist_lookup_uint64(props,
3835          zpool_prop_to_name(ZPOOL_PROP_READONLY), &readonly);
3836      if (readonly)
3837          mode = FREAD;
3838      spa = spa_add(pool, config, altroot);
3839      spa->spa_import_flags = flags;

3841      /*
3842       * Verbatim import - Take a pool and insert it into the namespace
3843       * as if it had been loaded at boot.
3844       */
  
```

```

3845      if (spa->spa_import_flags & ZFS_IMPORT_VERBATIM) {
3846          if (props != NULL)
3847              spa_configfile_set(spa, props, B_FALSE);

3849          spa_config_sync(spa, B_FALSE, B_TRUE);

3851          mutex_exit(&spa_namespace_lock);
3852          spa_history_log_version(spa, "import");

3854          return (0);
3855      }

3857      spa_activate(spa, mode);

3859      /*
3860       * Don't start async tasks until we know everything is healthy.
3861       */
3862      spa_async_suspend(spa);

3864      zpool_get_rewind_policy(config, &policy);
3865      if (policy.zrp_request & ZPOOL_DO_REWIND)
3866          state = SPA_LOAD_RECOVER;

3868      /*
3869       * Pass off the heavy lifting to spa_load(). Pass TRUE for mosconfig
3870       * because the user-supplied config is actually the one to trust when
3871       * doing an import.
3872       */
3873      if (state != SPA_LOAD_RECOVER)
3874          spa->spa_last_ubsync_txg = spa->spa_load_txg = 0;

3876      error = spa_load_best(spa, state, B_TRUE, policy.zrp_txg,
3877          policy.zrp_request);

3879      /*
3880       * Propagate anything learned while loading the pool and pass it
3881       * back to caller (i.e. rewind info, missing devices, etc).
3882       */
3883      VERIFY(nvlist_add_nvlist(config, ZPOOL_CONFIG_LOAD_INFO,
3884          spa->spa_load_info) == 0);

3886      spa_config_enter(spa, SCL_ALL, FTAG, RW_WRITER);
3887      /*
3888       * Toss any existing sparelist, as it doesn't have any validity
3889       * anymore, and conflicts with spa_has_spare().
3890       */
3891      if (spa->spa_spares.sav_config) {
3892          nvlist_free(spa->spa_spares.sav_config);
3893          spa->spa_spares.sav_config = NULL;
3894          spa_load_spares(spa);
3895      }
3883      if (spa->spa_l2cache.sav_config) {
3884          nvlist_free(spa->spa_l2cache.sav_config);
3885          spa->spa_l2cache.sav_config = NULL;
3886          spa_load_l2cache(spa);
3887      }

3897      VERIFY(nvlist_lookup_nvlist(config, ZPOOL_CONFIG_VDEV_TREE,
3898          &nvroot) == 0);
3899      if (error == 0)
3900          error = spa_validate_aux(spa, nvroot, -1ULL,
3901              VDEV_ALLOC_SPARE);
3902      if (error == 0)
3903          error = spa_validate_aux(spa, nvroot, -1ULL,
3904              VDEV_ALLOC_L2CACHE);
3905      spa_config_exit(spa, SCL_ALL, FTAG);
  
```

```

3907     if (props != NULL)
3908         spa_configfile_set(spa, props, B_FALSE);

3910     if (error != 0 || (props && spa_writeable(spa) &&
3911         (error = spa_prop_set(spa, props)))) {
3912         spa_unload(spa);
3913         spa_deactivate(spa);
3914         spa_remove(spa);
3915         mutex_exit(&spa_namespace_lock);
3916         return (error);
3917     }

3919     spa_async_resume(spa);

3921     /*
3922     * Override any spares and level 2 cache devices as specified by
3923     * the user, as these may have correct device names/devids, etc.
3924     */
3925     if (nvlist_lookup_nvlist_array(nvroot, ZPOOL_CONFIG_SPARES,
3926         &spares, &nspares) == 0) {
3927         if (spa->spa_spares.sav_config)
3928             VERIFY(nvlist_remove(spa->spa_spares.sav_config,
3929                 ZPOOL_CONFIG_SPARES, DATA_TYPE_NVLIST_ARRAY) == 0);
3930         else
3931             VERIFY(nvlist_alloc(&spa->spa_spares.sav_config,
3932                 NV_UNIQUE_NAME, KM_SLEEP) == 0);
3933         VERIFY(nvlist_add_nvlist_array(spa->spa_spares.sav_config,
3934             ZPOOL_CONFIG_SPARES, spares, nspares) == 0);
3935         spa_config_enter(spa, SCL_ALL, FTAG, RW_WRITER);
3936         spa_load_spares(spa);
3937         spa_config_exit(spa, SCL_ALL, FTAG);
3938         spa->spa_spares.sav_sync = B_TRUE;
3939     }
3940     if (nvlist_lookup_nvlist_array(nvroot, ZPOOL_CONFIG_L2CACHE,
3941         &l2cache, &nl2cache) == 0) {
3942         if (spa->spa_l2cache.sav_config)
3943             VERIFY(nvlist_remove(spa->spa_l2cache.sav_config,
3944                 ZPOOL_CONFIG_L2CACHE, DATA_TYPE_NVLIST_ARRAY) == 0);
3945         else
3946             VERIFY(nvlist_alloc(&spa->spa_l2cache.sav_config,
3947                 NV_UNIQUE_NAME, KM_SLEEP) == 0);
3948         VERIFY(nvlist_add_nvlist_array(spa->spa_l2cache.sav_config,
3949             ZPOOL_CONFIG_L2CACHE, l2cache, nl2cache) == 0);
3950         spa_config_enter(spa, SCL_ALL, FTAG, RW_WRITER);
3951         spa_load_l2cache(spa);
3952         spa_config_exit(spa, SCL_ALL, FTAG);
3953         spa->spa_l2cache.sav_sync = B_TRUE;
3954     }

3956     /*
3957     * Check for any removed devices.
3958     */
3959     if (spa->spa_autoreplace) {
3960         spa_aux_check_removed(&spa->spa_spares);
3961         spa_aux_check_removed(&spa->spa_l2cache);
3962     }

3964     if (spa_writeable(spa)) {
3965         /*
3966         * Update the config cache to include the newly-imported pool.
3967         */
3968         spa_config_update(spa, SPA_CONFIG_UPDATE_POOL);
3969     }

3971     /*

```

```

3972         * It's possible that the pool was expanded while it was exported.
3973         * We kick off an async task to handle this for us.
3974         */
3975         spa_async_request(spa, SPA_ASYNC_AUTOEXPAND);

3977         mutex_exit(&spa_namespace_lock);
3978         spa_history_log_version(spa, "import");

3980         return (0);
3981     }
    unchanged_portion_omitted

```

new/usr/src/uts/common/fs/zfs/sys/arc.h

1

```
*****
4516 Mon Aug 5 21:59:09 2013
new/usr/src/uts/common/fs/zfs/sys/arc.h
3525 Persistent L2ARC
*****
_____unchanged_portion_omitted_

84 void arc_space_consume(uint64_t space, arc_space_type_t type);
85 void arc_space_return(uint64_t space, arc_space_type_t type);
86 void *arc_data_buf_alloc(uint64_t space);
87 void arc_data_buf_free(void *buf, uint64_t space);
88 arc_buf_t *arc_buf_alloc(spa_t *spa, int size, void *tag,
89     arc_buf_contents_t type);
90 arc_buf_t *arc_loan_buf(spa_t *spa, int size);
91 void arc_return_buf(arc_buf_t *buf, void *tag);
92 void arc_loan_inuse_buf(arc_buf_t *buf, void *tag);
93 void arc_buf_add_ref(arc_buf_t *buf, void *tag);
94 boolean_t arc_buf_remove_ref(arc_buf_t *buf, void *tag);
95 int arc_buf_size(arc_buf_t *buf);
96 void arc_release(arc_buf_t *buf, void *tag);
97 int arc_released(arc_buf_t *buf);
98 int arc_has_callback(arc_buf_t *buf);
99 void arc_buf_freeze(arc_buf_t *buf);
100 void arc_buf_thaw(arc_buf_t *buf);
101 boolean_t arc_buf_eviction_needed(arc_buf_t *buf);
102 #ifdef ZFS_DEBUG
103 int arc_referenced(arc_buf_t *buf);
104 #endif

106 int arc_read(zio_t *pio, spa_t *spa, const blkptr_t *bp,
107     arc_done_func_t *done, void *private, int priority, int flags,
108     uint32_t *arc_flags, const zbookmark_t *zb);
109 zio_t *arc_write(zio_t *pio, spa_t *spa, uint64_t txg,
110     blkptr_t *bp, arc_buf_t *buf, boolean_t l2arc, boolean_t l2arc_compress,
111     const zio_prop_t *zp, arc_done_func_t *ready, arc_done_func_t *done,
112     void *private, int priority, int zio_flags, const zbookmark_t *zb);
113 void arc_freed(spa_t *spa, const blkptr_t *bp);

115 void arc_set_callback(arc_buf_t *buf, arc_evict_func_t *func, void *private);
116 int arc_buf_evict(arc_buf_t *buf);

118 void arc_flush(spa_t *spa);
119 void arc_tempreserve_clear(uint64_t reserve);
120 int arc_tempreserve_space(uint64_t reserve, uint64_t txg);

122 void arc_init(void);
123 void arc_fini(void);

125 /*
126  * Level 2 ARC
127  */

129 void l2arc_add_vdev(spa_t *spa, vdev_t *vd, boolean_t rebuild);
129 void l2arc_add_vdev(spa_t *spa, vdev_t *vd);
130 void l2arc_remove_vdev(vdev_t *vd);
131 boolean_t l2arc_vdev_present(vdev_t *vd);
132 void l2arc_init(void);
133 void l2arc_fini(void);
134 void l2arc_start(void);
135 void l2arc_stop(void);

137 #ifndef _KERNEL
138 extern boolean_t arc_watch;
139 extern int arc_procfid;
140 #endif
```

new/usr/src/uts/common/fs/zfs/sys/arc.h

2

```
142 #ifdef __cplusplus
143 }
_____unchanged_portion_omitted_
```

new/usr/src/uts/common/fs/zfs/sys/spa.h

1

```
*****
25969 Mon Aug 5 21:59:09 2013
new/usr/src/uts/common/fs/zfs/sys/spa.h
3525 Persistent L2ARC
*****
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) 2012 by Delphix. All rights reserved.
24 * Copyright 2011 Nexenta Systems, Inc. All rights reserved.
25 * Copyright 2013 Saso Kiselkov. All rights reserved.
26 */

28 #ifndef _SYS_SPA_H
29 #define _SYS_SPA_H

31 #include <sys/avl.h>
32 #include <sys/zfs_context.h>
33 #include <sys/nvpair.h>
34 #include <sys/sysmacros.h>
35 #include <sys/types.h>
36 #include <sys/fs/zfs.h>

38 #ifdef __cplusplus
39 extern "C" {
40 #endif

42 /*
43  * Forward references that lots of things need.
44  */
45 typedef struct spa spa_t;
46 typedef struct vdev vdev_t;
47 typedef struct metaslab metaslab_t;
48 typedef struct metaslab_group metaslab_group_t;
49 typedef struct metaslab_class metaslab_class_t;
50 typedef struct zio zio_t;
51 typedef struct zillog zillog_t;
52 typedef struct spa_aux_vdev spa_aux_vdev_t;
53 typedef struct ddt ddt_t;
54 typedef struct ddt_entry ddt_entry_t;
55 struct dsl_pool;
56 struct dsl_dataset;

58 /*
59  * General-purpose 32-bit and 64-bit bitfield encodings.
60  */
61 #define BF32_DECODE(x, low, len)      P2PHASE((x) >> (low), 1U << (len))
```

new/usr/src/uts/common/fs/zfs/sys/spa.h

2

```
62 #define BF64_DECODE(x, low, len)      P2PHASE((x) >> (low), 1ULL << (len))
63 #define BF32_ENCODE(x, low, len)      (P2PHASE((x), 1U << (len)) << (low))
64 #define BF64_ENCODE(x, low, len)      (P2PHASE((x), 1ULL << (len)) << (low))

66 #define BF32_GET(x, low, len)          BF32_DECODE(x, low, len)
67 #define BF64_GET(x, low, len)          BF64_DECODE(x, low, len)

69 #define BF32_SET(x, low, len, val)     \
70 ((x) ^= BF32_ENCODE((x) >> low) ^ (val), low, len)
71 #define BF64_SET(x, low, len, val)     \
72 ((x) ^= BF64_ENCODE((x) >> low) ^ (val), low, len)

74 #define BF32_GET_SB(x, low, len, shift, bias) \
75 ((BF32_GET(x, low, len) + (bias)) << (shift))
76 #define BF64_GET_SB(x, low, len, shift, bias) \
77 ((BF64_GET(x, low, len) + (bias)) << (shift))

79 #define BF32_SET_SB(x, low, len, shift, bias, val) \
80 BF32_SET(x, low, len, ((val) >> (shift)) - (bias))
81 #define BF64_SET_SB(x, low, len, shift, bias, val) \
82 BF64_SET(x, low, len, ((val) >> (shift)) - (bias))

84 /*
85  * We currently support nine block sizes, from 512 bytes to 128K.
86  * We could go higher, but the benefits are near-zero and the cost
87  * of COWing a giant block to modify one byte would become excessive.
88  */
89 #define SPA_MINBLOCKSHIFT      9
90 #define SPA_MAXBLOCKSHIFT     17
91 #define SPA_MINBLOCKSIZE      (1ULL << SPA_MINBLOCKSHIFT)
92 #define SPA_MAXBLOCKSIZE      (1ULL << SPA_MAXBLOCKSHIFT)

94 #define SPA_BLOCKSIZE          (SPA_MAXBLOCKSHIFT - SPA_MINBLOCKSHIFT + 1)

96 /*
97  * Size of block to hold the configuration data (a packed nvlist)
98  */
99 #define SPA_CONFIG_BLOCKSIZE  (1ULL << 14)

101 /*
102  * The DVA size encodings for LSIZE and PSIZE support blocks up to 32MB.
103  * The ASIZE encoding should be at least 64 times larger (6 more bits)
104  * to support up to 4-way RAID-Z mirror mode with worst-case gang block
105  * overhead, three DVAs per bp, plus one more bit in case we do anything
106  * else that expands the ASIZE.
107  */
108 #define SPA_LSIZEBITS          16      /* LSIZE up to 32M (2^16 * 512) */
109 #define SPA_PSIZEBITS          16      /* PSIZE up to 32M (2^16 * 512) */
110 #define SPA_ASIZEBITS          24      /* ASIZE up to 64 times larger */

112 /*
113  * All SPA data is represented by 128-bit data virtual addresses (DVAs).
114  * The members of the dva_t should be considered opaque outside the SPA.
115  */
116 typedef struct dva {
117     uint64_t          dva_word[2];
118 } dva_t;
119 unchanged portion omitted

262 #define BP_GET_ASIZE(bp)          \
263 (DVA_GET_ASIZE(&(bp)->blk_dva[0]) + DVA_GET_ASIZE(&(bp)->blk_dva[1]) + \
264 DVA_GET_ASIZE(&(bp)->blk_dva[2]))

266 #define BP_GET_UCSIZE(bp) \
267 ((BP_GET_LEVEL(bp) > 0 || DMU_OT_IS_METADATA(BP_GET_TYPE(bp))) ? \
268 BP_GET_PSIZE(bp) : BP_GET_LSIZE(bp))
```

```

270 #define BP_GET_NDVAS(bp) \
271     (!!DVA_GET_ASIZE(&(bp)->blk_dva[0]) + \
272     !!DVA_GET_ASIZE(&(bp)->blk_dva[1]) + \
273     !!DVA_GET_ASIZE(&(bp)->blk_dva[2]))
274
275 #define BP_COUNT_GANG(bp) \
276     (DVA_GET_GANG(&(bp)->blk_dva[0]) + \
277     DVA_GET_GANG(&(bp)->blk_dva[1]) + \
278     DVA_GET_GANG(&(bp)->blk_dva[2]))
279
280 #define DVA_EQUAL(dva1, dva2) \
281     ((dva1)->dva_word[1] == (dva2)->dva_word[1] && \
282     (dva1)->dva_word[0] == (dva2)->dva_word[0])
283
284 #define BP_EQUAL(bp1, bp2) \
285     (BP_PHYSICAL_BIRTH(bp1) == BP_PHYSICAL_BIRTH(bp2) && \
286     DVA_EQUAL(&(bp1)->blk_dva[0], &(bp2)->blk_dva[0]) && \
287     DVA_EQUAL(&(bp1)->blk_dva[1], &(bp2)->blk_dva[1]) && \
288     DVA_EQUAL(&(bp1)->blk_dva[2], &(bp2)->blk_dva[2]))
289
290 #define ZIO_CHECKSUM_EQUAL(zc1, zc2) \
291     (0 == (((zc1).zc_word[0] - (zc2).zc_word[0]) | \
292     ((zc1).zc_word[1] - (zc2).zc_word[1]) | \
293     ((zc1).zc_word[2] - (zc2).zc_word[2]) | \
294     ((zc1).zc_word[3] - (zc2).zc_word[3])))
295
296 #define ZIO_CHECKSUM_BSWAP(_zc) \
297     do { \
298         zio_cksum_t *zc = (_zc); \
299         zc->zc_word[0] = BSWAP_64(zc->zc_word[0]); \
300         zc->zc_word[1] = BSWAP_64(zc->zc_word[1]); \
301         zc->zc_word[2] = BSWAP_64(zc->zc_word[2]); \
302         zc->zc_word[3] = BSWAP_64(zc->zc_word[3]); \
303         _NOTE(NOTREACHED) \
304         _NOTE(CONSTCOND) \
305     } while (0)
306
307 #define DVA_IS_VALID(dva) (DVA_GET_ASIZE(dva) != 0)
308
309 #define ZIO_SET_CHECKSUM(zcp, w0, w1, w2, w3) \
310 { \
311     (zcp)->zc_word[0] = w0; \
312     (zcp)->zc_word[1] = w1; \
313     (zcp)->zc_word[2] = w2; \
314     (zcp)->zc_word[3] = w3; \
315 }

```

unchanged_portion_omitted

```
*****
86544 Mon Aug  5 21:59:09 2013
new/usr/src/uts/common/fs/zfs/vdev.c
3525 Persistent L2ARC
*****
_unchanged_portion_omitted_

1487 /*
1488  * Reopen all interior vdevs and any unopened leaves. We don't actually
1489  * reopen leaf vdevs which had previously been opened as they might deadlock
1490  * on the spa_config_lock. Instead we only obtain the leaf's physical size.
1491  * If the leaf has never been opened then open it, as usual.
1492  */
1493 void
1494 vdev_reopen(vdev_t *vd)
1495 {
1496     spa_t *spa = vd->vdev_spa;

1498     ASSERT(spa_config_held(spa, SCL_STATE_ALL, RW_WRITER) == SCL_STATE_ALL);

1500     /* set the reopening flag unless we're taking the vdev offline */
1501     vd->vdev_reopening = !vd->vdev_offline;
1502     vdev_close(vd);
1503     (void) vdev_open(vd);

1505     /*
1506      * Call vdev_validate() here to make sure we have the same device.
1507      * Otherwise, a device with an invalid label could be successfully
1508      * opened in response to vdev_reopen().
1509      */
1510     if (vd->vdev_aux) {
1511         (void) vdev_validate_aux(vd);
1512         if (vdev_readable(vd) && vdev_writeable(vd) &&
1513             vd->vdev_aux == &spa->spa_l2cache &&
1514             !l2arc_vdev_present(vd)) {
1515             /*
1516              * When reopening we can assume persistent L2ARC is
1517              * supported, since we've already opened the device
1518              * in the past and prepended an L2 uberblock.
1519              */
1520             l2arc_add_vdev(spa, vd, B_TRUE);
1521         }
1522         !l2arc_vdev_present(vd)
1523         l2arc_add_vdev(spa, vd);
1524     } else {
1525         (void) vdev_validate(vd, B_TRUE);
1526     }

1527     /*
1528      * Reassess parent vdev's health.
1529      */
1530     vdev_propagate_state(vd);
1531 }
_unchanged_portion_omitted_
```

```

*****
38163 Mon Aug 5 21:59:10 2013
new/usr/src/uts/common/fs/zfs/vdev_label.c
3525 Persistent L2ARC
*****
_____unchanged_portion_omitted_____

210 /*
211  * Generate the nvlist representing this vdev's config.
212  */
213 nvlist_t *
214 vdev_config_generate(spa_t *spa, vdev_t *vd, boolean_t getstats,
215     vdev_config_flag_t flags)
216 {
217     nvlist_t *nv = NULL;

219     VERIFY(nvlist_alloc(&nv, NV_UNIQUE_NAME, KM_SLEEP) == 0);

221     VERIFY(nvlist_add_string(nv, ZPOOL_CONFIG_TYPE,
222         vd->vdev_ops->vdev_op_type) == 0);
223     if (!(flags & (VDEV_CONFIG_SPARE | VDEV_CONFIG_L2CACHE)))
224         VERIFY(nvlist_add_uint64(nv, ZPOOL_CONFIG_ID, vd->vdev_id)
225             == 0);
226     VERIFY(nvlist_add_uint64(nv, ZPOOL_CONFIG_GUID, vd->vdev_guid) == 0);

228     if (vd->vdev_path != NULL)
229         VERIFY(nvlist_add_string(nv, ZPOOL_CONFIG_PATH,
230             vd->vdev_path) == 0);

232     if (vd->vdev_devid != NULL)
233         VERIFY(nvlist_add_string(nv, ZPOOL_CONFIG_DEVID,
234             vd->vdev_devid) == 0);

236     if (vd->vdev_physpath != NULL)
237         VERIFY(nvlist_add_string(nv, ZPOOL_CONFIG_PHYS_PATH,
238             vd->vdev_physpath) == 0);

240     if (vd->vdev_fru != NULL)
241         VERIFY(nvlist_add_string(nv, ZPOOL_CONFIG_FRU,
242             vd->vdev_fru) == 0);

244     if (vd->vdev_nparity != 0) {
245         ASSERT(strcmp(vd->vdev_ops->vdev_op_type,
246             VDEV_TYPE_RAIDZ) == 0);

248         /*
249          * Make sure someone hasn't managed to sneak a fancy new vdev
250          * into a crufy old storage pool.
251          */
252         ASSERT(vd->vdev_nparity == 1 ||
253             (vd->vdev_nparity <= 2 &&
254             spa_version(spa) >= SPA_VERSION_RAIDZ2) ||
255             (vd->vdev_nparity <= 3 &&
256             spa_version(spa) >= SPA_VERSION_RAIDZ3));

258         /*
259          * Note that we'll add the nparity tag even on storage pools
260          * that only support a single parity device -- older software
261          * will just ignore it.
262          */
263         VERIFY(nvlist_add_uint64(nv, ZPOOL_CONFIG_NPARITY,
264             vd->vdev_nparity) == 0);
265     }

267     if (vd->vdev_wholedisk != -1ULL)
268         VERIFY(nvlist_add_uint64(nv, ZPOOL_CONFIG_WHOLE_DISK,

```

```

269         vd->vdev_wholedisk) == 0);

271     if (vd->vdev_not_present)
272         VERIFY(nvlist_add_uint64(nv, ZPOOL_CONFIG_NOT_PRESENT, 1) == 0);

274     if (vd->vdev_isspare)
275         VERIFY(nvlist_add_uint64(nv, ZPOOL_CONFIG_IS_SPARE, 1) == 0);

277     if (!(flags & (VDEV_CONFIG_SPARE | VDEV_CONFIG_L2CACHE))) &&
278         vd == vd->vdev_top) {
279         VERIFY(nvlist_add_uint64(nv, ZPOOL_CONFIG_METASLAB_ARRAY,
280             vd->vdev_ms_array) == 0);
281         VERIFY(nvlist_add_uint64(nv, ZPOOL_CONFIG_METASLAB_SHIFT,
282             vd->vdev_ms_shift) == 0);
283         VERIFY(nvlist_add_uint64(nv, ZPOOL_CONFIG_ASHIFT,
284             vd->vdev_ashift) == 0);
285         VERIFY(nvlist_add_uint64(nv, ZPOOL_CONFIG_ASIZE,
286             vd->vdev_asize) == 0);
287         VERIFY(nvlist_add_uint64(nv, ZPOOL_CONFIG_IS_LOG,
288             vd->vdev_islog) == 0);
289         if (vd->vdev_removing)
290             VERIFY(nvlist_add_uint64(nv, ZPOOL_CONFIG_REMOVING,
291                 vd->vdev_removing) == 0);
292     }

294     if (flags & VDEV_CONFIG_L2CACHE)
295         /* indicate that we support L2ARC persistency */
296         VERIFY(nvlist_add_boolean_value(nv,
297             ZPOOL_CONFIG_L2CACHE_PERSISTENT, B_TRUE) == 0);

299     if (vd->vdev_dtl_smo.smo_object != 0)
300         VERIFY(nvlist_add_uint64(nv, ZPOOL_CONFIG_DTL,
301             vd->vdev_dtl_smo.smo_object) == 0);

303     if (vd->vdev_crtxg)
304         VERIFY(nvlist_add_uint64(nv, ZPOOL_CONFIG_CREATE_TXG,
305             vd->vdev_crtxg) == 0);

307     if (getstats) {
308         vdev_stat_t vs;
309         pool_scan_stat_t ps;

311         vdev_get_stats(vd, &vs);
312         VERIFY(nvlist_add_uint64_array(nv, ZPOOL_CONFIG_VDEV_STATS,
313             (uint64_t *)&vs, sizeof (vs) / sizeof (uint64_t)) == 0);

315         /* provide either current or previous scan information */
316         if (spa_scan_get_stats(spa, &ps) == 0) {
317             VERIFY(nvlist_add_uint64_array(nv,
318                 ZPOOL_CONFIG_SCAN_STATS, (uint64_t *)&ps,
319                 sizeof (pool_scan_stat_t) / sizeof (uint64_t))
320                 == 0);
321         }
322     }

324     if (!vd->vdev_ops->vdev_op_leaf) {
325         nvlist_t **child;
326         int c, idx;

328         ASSERT(!vd->vdev_ishole);

330         child = kmem_alloc(vd->vdev_children * sizeof (nvlist_t *),
331             KM_SLEEP);

333         for (c = 0, idx = 0; c < vd->vdev_children; c++) {
334             vdev_t *cvd = vd->vdev_child[c];

```

```

336         /*
337          * If we're generating an nvlist of removing
338          * vdevs then skip over any device which is
339          * not being removed.
340          */
341         if ((flags & VDEV_CONFIG_REMOVING) &&
342             !cvd->vdev_removing)
343             continue;
344
345         child[idx++] = vdev_config_generate(spa, cvd,
346             getstats, flags);
347     }
348
349     if (idx) {
350         VERIFY(nvlist_add_nvlist_array(nv,
351             ZPOOL_CONFIG_CHILDREN, child, idx) == 0);
352     }
353
354     for (c = 0; c < idx; c++)
355         nvlist_free(child[c]);
356
357     kmem_free(child, vd->vdev_children * sizeof (nvlist_t *));
358
359 } else {
360     const char *aux = NULL;
361
362     if (vd->vdev_offline && !vd->vdev_tmpoffline)
363         VERIFY(nvlist_add_uint64(nv, ZPOOL_CONFIG_OFFLINE,
364             B_TRUE) == 0);
365     if (vd->vdev_resilvering)
366         VERIFY(nvlist_add_uint64(nv, ZPOOL_CONFIG_RESILVERING,
367             B_TRUE) == 0);
368     if (vd->vdev_faulted)
369         VERIFY(nvlist_add_uint64(nv, ZPOOL_CONFIG_FAULTED,
370             B_TRUE) == 0);
371     if (vd->vdev_degraded)
372         VERIFY(nvlist_add_uint64(nv, ZPOOL_CONFIG_DEGRADED,
373             B_TRUE) == 0);
374     if (vd->vdev_removed)
375         VERIFY(nvlist_add_uint64(nv, ZPOOL_CONFIG_REMOVED,
376             B_TRUE) == 0);
377     if (vd->vdev_unspare)
378         VERIFY(nvlist_add_uint64(nv, ZPOOL_CONFIG_UNSPARE,
379             B_TRUE) == 0);
380     if (vd->vdev_ishole)
381         VERIFY(nvlist_add_uint64(nv, ZPOOL_CONFIG_IS_HOLE,
382             B_TRUE) == 0);
383
384     switch (vd->vdev_stat.vs_aux) {
385     case VDEV_AUX_ERR_EXCEEDED:
386         aux = "err_exceeded";
387         break;
388
389     case VDEV_AUX_EXTERNAL:
390         aux = "external";
391         break;
392     }
393
394     if (aux != NULL)
395         VERIFY(nvlist_add_string(nv, ZPOOL_CONFIG_AUX_STATE,
396             aux) == 0);
397
398     if (vd->vdev_splitting && vd->vdev_orig_guid != 0LL) {
399         VERIFY(nvlist_add_uint64(nv, ZPOOL_CONFIG_ORIG_GUID,
400             vd->vdev_orig_guid) == 0);

```

```

401     }
402 }
403
404     return (nv);
405 }

```

unchanged_portion_omitted


```

*****
29031 Mon Aug 5 21:59:10 2013
new/usr/src/uts/common/sys/fs/zfs.h
3525 Persistent L2ARC
*****
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 2011 Nexenta Systems, Inc. All rights reserved.
26  * Copyright (c) 2012, Joyent, Inc. All rights reserved.
27  * Copyright (c) 2013, Saso Kiselkov. All rights reserved.
28 */

30 /* Portions Copyright 2010 Robert Milkowski */

32 #ifndef _SYS_FS_ZFS_H
33 #define _SYS_FS_ZFS_H

35 #include <sys/time.h>

37 #ifdef __cplusplus
38 extern "C" {
39 #endif

41 /*
42  * Types and constants shared between userland and the kernel.
43  */

45 /*
46  * Each dataset can be one of the following types. These constants can be
47  * combined into masks that can be passed to various functions.
48  */
49 typedef enum {
50     ZFS_TYPE_FILESYSTEM      = 0x1,
51     ZFS_TYPE_SNAPSHOT       = 0x2,
52     ZFS_TYPE_VOLUME         = 0x4,
53     ZFS_TYPE_POOL           = 0x8
54 } zfs_type_t;
    unchanged_portion_omitted

477 /*
478  * The following are configuration names used in the nvlist describing a pool's
479  * configuration.
480  */
481 #define ZPOOL_CONFIG_VERSION          "version"

```

```

482 #define ZPOOL_CONFIG_POOL_NAME      "name"
483 #define ZPOOL_CONFIG_POOL_STATE     "state"
484 #define ZPOOL_CONFIG_POOL_TXG      "txg"
485 #define ZPOOL_CONFIG_POOL_GUID     "pool_guid"
486 #define ZPOOL_CONFIG_CREATE_TXG    "create_txg"
487 #define ZPOOL_CONFIG_TOP_GUID      "top_guid"
488 #define ZPOOL_CONFIG_VDEV_TREE     "vdev_tree"
489 #define ZPOOL_CONFIG_TYPE          "type"
490 #define ZPOOL_CONFIG_CHILDREN      "children"
491 #define ZPOOL_CONFIG_ID            "id"
492 #define ZPOOL_CONFIG_GUID          "guid"
493 #define ZPOOL_CONFIG_PATH          "path"
494 #define ZPOOL_CONFIG_DEVID         "devid"
495 #define ZPOOL_CONFIG_METASLAB_ARRAY "metaslab_array"
496 #define ZPOOL_CONFIG_METASLAB_SHIFT "metaslab_shift"
497 #define ZPOOL_CONFIG_ASHIFT        "ashift"
498 #define ZPOOL_CONFIG_ASIZE         "asize"
499 #define ZPOOL_CONFIG_DTL           "DTL"
500 #define ZPOOL_CONFIG_SCAN_STATS    "scan_stats" /* not stored on disk */
501 #define ZPOOL_CONFIG_VDEV_STATS    "vdev_stats" /* not stored on disk */
502 #define ZPOOL_CONFIG_WHOLE_DISK    "whole_disk"
503 #define ZPOOL_CONFIG_ERRCOUNT     "error_count"
504 #define ZPOOL_CONFIG_NOT_PRESENT    "not_present"
505 #define ZPOOL_CONFIG_SPARES        "spares"
506 #define ZPOOL_CONFIG_IS_SPARE      "is_spare"
507 #define ZPOOL_CONFIG_NPARITY       "nparity"
508 #define ZPOOL_CONFIG_HOSTID        "hostid"
509 #define ZPOOL_CONFIG_HOSTNAME      "hostname"
510 #define ZPOOL_CONFIG_LOADED_TIME   "initial_load_time"
511 #define ZPOOL_CONFIG_UNSPARE       "unspare"
512 #define ZPOOL_CONFIG_PHYS_PATH     "phys_path"
513 #define ZPOOL_CONFIG_IS_LOG        "is_log"
514 #define ZPOOL_CONFIG_L2CACHE       "l2cache"
515 #define ZPOOL_CONFIG_L2CACHE_PERSISTENT "l2cache_persistent"
516 #define ZPOOL_CONFIG_HOLE_ARRAY    "hole_array"
517 #define ZPOOL_CONFIG_VDEV_CHILDREN "vdev_children"
518 #define ZPOOL_CONFIG_IS_HOLE       "is_hole"
519 #define ZPOOL_CONFIG_DDT_HISTOGRAM "ddt_histogram"
520 #define ZPOOL_CONFIG_DDT_OBJ_STATS "ddt_object_stats"
521 #define ZPOOL_CONFIG_DDT_STATS     "ddt_stats"
522 #define ZPOOL_CONFIG_SPLIT         "splitcfg"
523 #define ZPOOL_CONFIG_ORIG_GUID     "orig_guid"
524 #define ZPOOL_CONFIG_SPLIT_GUID    "split_guid"
525 #define ZPOOL_CONFIG_SPLIT_LIST    "guid_list"
526 #define ZPOOL_CONFIG_REMOVING      "removing"
527 #define ZPOOL_CONFIG_RESILVERING   "resilvering"
528 #define ZPOOL_CONFIG_COMMENT       "comment"
529 #define ZPOOL_CONFIG_SUSPENDED     "suspended" /* not stored on disk */
530 #define ZPOOL_CONFIG_TIMESTAMP     "timestamp" /* not stored on disk */
531 #define ZPOOL_CONFIG_BOOTFS        "bootfs" /* not stored on disk */
532 #define ZPOOL_CONFIG_MISSING_DEVICES "missing_vdevs" /* not stored on disk */
533 #define ZPOOL_CONFIG_LOAD_INFO     "load_info" /* not stored on disk */
534 #define ZPOOL_CONFIG_REWIND_INFO   "rewind_info" /* not stored on disk */
535 #define ZPOOL_CONFIG_UNSUP_FEAT    "unsup_feat" /* not stored on disk */
536 #define ZPOOL_CONFIG_ENABLED_FEAT  "enabled_feat" /* not stored on disk */
537 #define ZPOOL_CONFIG_CAN_RDONLY    "can_rdonly" /* not stored on disk */
538 #define ZPOOL_CONFIG_FEATURES_FOR_READ "features_for_read"
539 #define ZPOOL_CONFIG_FEATURE_STATS "feature_stats" /* not stored on disk */
540 /*
541  * The persistent vdev state is stored as separate values rather than a single
542  * 'vdev_state' entry. This is because a device can be in multiple states, such
543  * as offline and degraded.
544  */
545 #define ZPOOL_CONFIG_OFFLINE        "offline"
546 #define ZPOOL_CONFIG_FAULTED       "faulted"
547 #define ZPOOL_CONFIG_DEGRADED       "degraded"

```

```
548 #define ZPOOL_CONFIG_REMOVED          "removed"
549 #define ZPOOL_CONFIG_FRU               "fru"
550 #define ZPOOL_CONFIG_AUX_STATE        "aux_state"

552 /* Rewind policy parameters */
553 #define ZPOOL_REWIND_POLICY            "rewind-policy"
554 #define ZPOOL_REWIND_REQUEST          "rewind-request"
555 #define ZPOOL_REWIND_REQUEST_TXG      "rewind-request-txg"
556 #define ZPOOL_REWIND_META_THRESH      "rewind-meta-thresh"
557 #define ZPOOL_REWIND_DATA_THRESH      "rewind-data-thresh"

559 /* Rewind data discovered */
560 #define ZPOOL_CONFIG_LOAD_TIME         "rewind_txg_ts"
561 #define ZPOOL_CONFIG_LOAD_DATA_ERRORS "verify_data_errors"
562 #define ZPOOL_CONFIG_REWIND_TIME      "seconds_of_rewind"

564 #define VDEV_TYPE_ROOT                 "root"
565 #define VDEV_TYPE_MIRROR               "mirror"
566 #define VDEV_TYPE_REPLACING            "replacing"
567 #define VDEV_TYPE_RAIDZ                "raidz"
568 #define VDEV_TYPE_DISK                 "disk"
569 #define VDEV_TYPE_FILE                 "file"
570 #define VDEV_TYPE_MISSING              "missing"
571 #define VDEV_TYPE_HOLE                  "hole"
572 #define VDEV_TYPE_SPARE                 "spare"
573 #define VDEV_TYPE_LOG                  "log"
574 #define VDEV_TYPE_L2CACHE              "l2cache"

576 /*
577  * This is needed in userland to report the minimum necessary device size.
578  */
579 #define SPA_MINDEVSIZE                 (64ULL << 20)

581 /*
582  * The location of the pool configuration repository, shared between kernel and
583  * userland.
584  */
585 #define ZPOOL_CACHE                     "/etc/zfs/zpool.cache"

587 /*
588  * vdev states are ordered from least to most healthy.
589  * A vdev that's CANT_OPEN or below is considered unusable.
590  */
591 typedef enum vdev_state {
592     VDEV_STATE_UNKNOWN = 0, /* Uninitialized vdev          */
593     VDEV_STATE_CLOSED,     /* Not currently open    */
594     VDEV_STATE_OFFLINE,    /* Not allowed to open   */
595     VDEV_STATE_REMOVED,    /* Explicitly removed from system */
596     VDEV_STATE_CANT_OPEN,  /* Tried to open, but failed */
597     VDEV_STATE_FAULTED,    /* External request to fault device */
598     VDEV_STATE_DEGRADED,   /* Replicated vdev with unhealthy kids */
599     VDEV_STATE_HEALTHY     /* Presumed good         */
600 } vdev_state_t;
unchanged portion omitted
```

new/usr/src/uts/common/sys/list.h

1

```
*****
2456 Mon Aug 5 21:59:10 2013
new/usr/src/uts/common/sys/list.h
3525 Persistent L2ARC
*****
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 2008 Sun Microsystems, Inc. All rights reserved.
23 * Use is subject to license terms.
24 */
25 /*
26 * Copyright 2013 Saso Kiselkov, All rights reserved.
27 */

29 #ifndef _SYS_LIST_H
30 #define _SYS_LIST_H

29 #pragma ident      "%Z%M% %I%      %E% SMI"

32 #include <sys/list_impl.h>

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

38 /*
39 * Please note that a list_node_t contains pointers back to its parent list_t
40 * so you cannot copy the list_t around once it has been initialized. In
41 * particular, this kind of construct won't work:
42 *
43 * struct { list_t l; } a, b;
44 * list_create(&a.l, ...);
45 * b = a;      <= This will break the list in 'b', as the 'l' element in 'a'
46 *             got copied to a different memory address.
47 *
48 * When copying structures with lists use list_move_tail() to move the list
49 * from the src to dst (the source reference will then become invalid).
50 */
51 typedef struct list_node list_node_t;
52 typedef struct list list_t;

54 void list_create(list_t *, size_t, size_t);
55 void list_destroy(list_t *);

57 void list_insert_after(list_t *, void *, void *);
58 void list_insert_before(list_t *, void *, void *);
59 void list_insert_head(list_t *, void *);
```

new/usr/src/uts/common/sys/list.h

2

```
60 void list_insert_tail(list_t *, void *);
61 void list_remove(list_t *, void *);
62 void *list_remove_head(list_t *);
63 void *list_remove_tail(list_t *);
64 void list_move_tail(list_t *, list_t *);

66 void *list_head(list_t *);
67 void *list_tail(list_t *);
68 void *list_next(list_t *, void *);
69 void *list_prev(list_t *, void *);
70 int list_is_empty(list_t *);

72 void list_link_init(list_node_t *);
73 void list_link_replace(list_node_t *, list_node_t *);

75 int list_link_active(list_node_t *);

77 #ifdef __cplusplus
78 }
_____unchanged_portion_omitted_____
```