```
**********************************************************
  131644 Mon Jun  4 22:08:28 2012
new/usr/src/uts/common/fs/zfs/arc.c
*** NO COMMENTS ***
**********************************************************
   1 /*
   2  * CDDL HEADER START
   3  *
   4  * The contents of this file are subject to the terms of the
   5  * Common Development and Distribution License (the "License").
   6  * You may not use this file except in compliance with the License.
   7  *
   8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
   9  * or http://www.opensolaris.org/os/licensing.
  10  * See the License for the specific language governing permissions
  11  * and limitations under the License.
  12  *
  13  * When distributing Covered Code, include this CDDL HEADER in each
  14  * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
  15  * If applicable, add the following below this CDDL HEADER, with the
  16  * fields enclosed by brackets "[]" replaced with your own identifying
  17  * information: Portions Copyright [yyyy] [name of copyright owner]
  18  *
  19  * CDDL HEADER END
  20  */
  21 /*
  22  * Copyright (c) 2005, 2010, Oracle and/or its affiliates. All rights reserved.
  23  * Copyright 2011 Nexenta Systems, Inc.  All rights reserved.
  24  * Copyright (c) 2012 by Delphix. All rights reserved.
  25  */
  
  27 /*
  28  * DVA-based Adjustable Replacement Cache
  29  *
  30  * While much of the theory of operation used here is
  31  * based on the self-tuning, low overhead replacement cache
  32  * presented by Megiddo and Modha at FAST 2003, there are some
  33  * significant differences:
  34  *
  35  * 1. The Megiddo and Modha model assumes any page is evictable.
  36  * Pages in its cache cannot be "locked" into memory.  This makes
  37  * the eviction algorithm simple: evict the last page in the list.
  38  * This also make the performance characteristics easy to reason
  39  * about.  Our cache is not so simple.  At any given moment, some
  40  * subset of the blocks in the cache are un-evictable because we
  41  * have handed out a reference to them.  Blocks are only evictable
  42  * when there are no external references active.  This makes
  43  * eviction far more problematic:  we choose to evict the evictable
  44  * blocks that are the "lowest" in the list.
  45  *
  46  * There are times when it is not possible to evict the requested
  47  * space.  In these circumstances we are unable to adjust the cache
  48  * size.  To prevent the cache growing unbounded at these times we
  49  * implement a "cache throttle" that slows the flow of new data
  50  * into the cache until we can make space available.
  51  *
  52  * 2. The Megiddo and Modha model assumes a fixed cache size.
  53  * Pages are evicted when the cache is full and there is a cache
  54  * miss.  Our model has a variable sized cache.  It grows with
  55  * high use, but also tries to react to memory pressure from the
  56  * operating system: decreasing its size when system memory is
  57  * tight.
  58  *
  59  * 3. The Megiddo and Modha model assumes a fixed page size. All
  60  * elements of the cache are therefor exactly the same size.  So
  61  * when adjusting the cache size following a cache miss, its simply
```

```
  62  * a matter of choosing a single page to evict.  In our model, we
  63  * have variable sized cache blocks (rangeing from 512 bytes to
  64  * 128K bytes).  We therefor choose a set of blocks to evict to make
  65  * space for a cache miss that approximates as closely as possible
  66  * the space used by the new block.
  67  *
  68  * See also:  "ARC: A Self-Tuning, Low Overhead Replacement Cache"
  69  * by N. Megiddo & D. Modha, FAST 2003
  70  */

  72 /*
  73  * The locking model:
  74  *
  75  * A new reference to a cache buffer can be obtained in two
  76  * ways: 1) via a hash table lookup using the DVA as a key,
  77  * or 2) via one of the ARC lists.  The arc_read() interface
  78  * uses method 1, while the internal arc algorithms for
  79  * adjusting the cache use method 2.  We therefor provide two
  80  * types of locks: 1) the hash table lock array, and 2) the
  81  * arc list locks.
  82  *
  83  * Buffers do not have their own mutexes, rather they rely on the
  84  * hash table mutexes for the bulk of their protection (i.e. most
  85  * fields in the arc_buf_hdr_t are protected by these mutexes).
  83  * Buffers do not have their own mutexs, rather they rely on the
  84  * hash table mutexs for the bulk of their protection (i.e. most
  85  * fields in the arc_buf_hdr_t are protected by these mutexs).
  86  *
  87  * buf_hash_find() returns the appropriate mutex (held) when it
  88  * locates the requested buffer in the hash table.  It returns
  89  * NULL for the mutex if the buffer was not in the table.
  90  *
  91  * buf_hash_remove() expects the appropriate hash mutex to be
  92  * already held before it is invoked.
  93  *
  94  * Each arc state also has a mutex which is used to protect the
  95  * buffer list associated with the state.  When attempting to
  96  * obtain a hash table lock while holding an arc list lock you
  97  * must use: mutex_tryenter() to avoid deadlock.  Also note that
  98  * the active state mutex must be held before the ghost state mutex.
  99  *
 100  * Arc buffers may have an associated eviction callback function.
 101  * This function will be invoked prior to removing the buffer (e.g.
 102  * in arc_do_user_evicts()).  Note however that the data associated
 103  * with the buffer may be evicted prior to the callback.  The callback
 104  * must be made with *no locks held* (to prevent deadlock).  Additionally,
 105  * the users of callbacks must ensure that their private data is
 106  * protected from simultaneous callbacks from arc_buf_evict()
 107  * and arc_do_user_evicts().
 108  *
 109  * Note that the majority of the performance stats are manipulated
 110  * with atomic operations.
 111  *
 112  * The L2ARC uses the l2arc_buflist_mtx global mutex for the following:
 113  *
 114  *      - L2ARC buflist creation
 115  *      - L2ARC buflist eviction
 116  *      - L2ARC write completion, which walks L2ARC buflists
 117  *      - ARC header destruction, as it removes from L2ARC buflists
 118  *      - ARC header release, as it removes from L2ARC buflists
 119  */

 121 #include <sys/spa.h>
 122 #include <sys/zio.h>
 123 #include <sys/zfs_context.h>
 124 #include <sys/arc.h>
```

```
 125 #include <sys/refcount.h>
 126 #include <sys/vdev.h>
 127 #include <sys/vdev_impl.h>
 128 #ifdef _KERNEL
 129 #include <sys/vmsystm.h>
 130 #include <vm/anon.h>
 131 #include <sys/fs/swapnode.h>
 132 #include <sys/dnlc.h>
 133 #endif
 134 #include <sys/callb.h>
 135 #include <sys/kstat.h>
 136 #include <zfs_fletcher.h>

 138 static kmutex_t          arc_reclaim_thr_lock;
 139 static kcondvar_t        arc_reclaim_thr_cv;     /* used to signal reclaim thr */
 140 static uint8_t           arc_thread_exit;

 142 extern int zfs_write_limit_shift;
 143 extern uint64_t zfs_write_limit_max;
 144 extern kmutex_t zfs_write_limit_lock;

 146 #define ARC_REDUCE_DNLC_PERCENT 3
 147 uint_t arc_reduce_dnlc_percent = ARC_REDUCE_DNLC_PERCENT;

 149 typedef enum arc_reclaim_strategy {
 150         ARC_RECLAIM_AGGR,                /* Aggressive reclaim strategy */
 151         ARC_RECLAIM_CONS                 /* Conservative reclaim strategy */
 152 } arc_reclaim_strategy_t;
_____unchanged_portion_omitted_

2638 /*
2639  * **"Read" the block at the specified DVA (in bp) via the**
2639  * *"Read" the block block at the specified DVA (in bp) via the*
2640  * cache.  If the block is found in the cache, invoke the provided
2641  * callback immediately and return.  Note that the `zio` parameter
2642  * in the callback will be NULL in this case, since no IO was
2643  * required.  If the block is not in the cache pass the read request
2644  * on to the spa with a substitute callback function, so that the
2645  * requested block will be added to the cache.
2646  *
2647  * If a read request arrives for a block that has a read in-progress,
2648  * either wait for the in-progress read to complete (and return the
2649  * results); or, if this is a read with a "done" func, add a record
2650  * to the read to invoke the "done" func when the read completes,
2651  * and return; or just return.
2652  *
2653  * arc_read_done() will invoke all the requested "done" functions
2654  * for readers of this block.
2655  *
2656  * Normal callers should use arc_read and pass the arc buffer and offset
2657  * for the bp.  But if you know you don't need locking, you can use
2658  * arc_read_bp.
2659  */
2660 int
2661 arc_read(zio_t *pio, spa_t *spa, const blkptr_t *bp, arc_buf_t *pbuf,
2662     arc_done_func_t *done, void *private, int priority, int zio_flags,
2663     uint32_t *arc_flags, const zbookmark_t *zb)
2664 {
2665         int err;

2667         if (pbuf == NULL) {
2668                 /*
2669                  * XXX This happens from traverse callback funcs, for
2670                  * the objset_phys_t block.
2671                  */
2672                 return (arc_read_nolock(pio, spa, bp, done, private, priority,
```

```
2673                     zio_flags, arc_flags, zb));
2674         }

2676         ASSERT(!refcount_is_zero(&pbuf->b_hdr->b_refcnt));
2677         ASSERT3U((char *)bp - (char *)pbuf->b_data, <, pbuf->b_hdr->b_size);
2678         rw_enter(&pbuf->b_data_lock, RW_READER);

2680         err = arc_read_nolock(pio, spa, bp, done, private, priority,
2681             zio_flags, arc_flags, zb);
2682         rw_exit(&pbuf->b_data_lock);

2684         return (err);
2685 }
_____unchanged_portion_omitted_
```