```
*********************************************************
  149962 Thu Oct 17 00:34:49 2013
new/usr/src/uts/common/fs/zfs/arc.c
New ARC buf_hash architecture
*********************************************************
   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, Joyent, Inc. All rights reserved.
  25  * Copyright (c) 2013 by Delphix. All rights reserved.
  26  * Copyright (c) 2013 by Saso Kiselkov. All rights reserved.
  27  */

  29 /*
  30  * DVA-based Adjustable Replacement Cache
  31  *
  32  * While much of the theory of operation used here is
  33  * based on the self-tuning, low overhead replacement cache
  34  * presented by Megiddo and Modha at FAST 2003, there are some
  35  * significant differences:
  36  *
  37  * 1. The Megiddo and Modha model assumes any page is evictable.
  38  * Pages in its cache cannot be "locked" into memory.  This makes
  39  * the eviction algorithm simple: evict the last page in the list.
  40  * This also make the performance characteristics easy to reason
  41  * about.  Our cache is not so simple.  At any given moment, some
  42  * subset of the blocks in the cache are un-evictable because we
  43  * have handed out a reference to them.  Blocks are only evictable
  44  * when there are no external references active.  This makes
  45  * eviction far more problematic:  we choose to evict the evictable
  46  * blocks that are the "lowest" in the list.
  47  *
  48  * There are times when it is not possible to evict the requested
  49  * space.  In these circumstances we are unable to adjust the cache
  50  * size.  To prevent the cache growing unbounded at these times we
  51  * implement a "cache throttle" that slows the flow of new data
  52  * into the cache until we can make space available.
  53  *
  54  * 2. The Megiddo and Modha model assumes a fixed cache size.
  55  * Pages are evicted when the cache is full and there is a cache
  56  * miss.  Our model has a variable sized cache.  It grows with
  57  * high use, but also tries to react to memory pressure from the
  58  * operating system: decreasing its size when system memory is
  59  * tight.
  60  *
  61  * 3. The Megiddo and Modha model assumes a fixed page size. All
```

```
  62  * elements of the cache are therefore exactly the same size.  So
  63  * when adjusting the cache size following a cache miss, its simply
  64  * a matter of choosing a single page to evict.  In our model, we
  65  * have variable sized cache blocks (rangeing from 512 bytes to
  66  * 128K bytes).  We therefore choose a set of blocks to evict to make
  67  * space for a cache miss that approximates as closely as possible
  68  * the space used by the new block.
  69  *
  70  * See also:  "ARC: A Self-Tuning, Low Overhead Replacement Cache"
  71  * by N. Megiddo & D. Modha, FAST 2003
  72  */

  74 /*
  75  * External users typically access ARC buffers via a hash table
  76  * lookup, using the DVA, spa_t pointer value and the birth TXG
  77  * number as the key. The hash value is derived by buf_hash(),
  78  * which spits out a 64-bit hash index. This index is then masked
  79  * with ht_mask to obtain the final index into the hash table:
  80  *
  81  *                   ,---------------- & ht_mask ----------------,
  82  * 64-bit hash value |              (hash table index)           |
  83  * |XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX|
  84  *
  85  * Sizing of the hash table is done at boot from the amount of
  86  * physical memory. We start with a base value of 2^12 hash
  87  * buckets and then evaluate whether this number, multiplied by
  88  * 2^zfs_arc_ht_base_masklen (the minimum mask length), is
  89  * greater than or equal to the amount of physical memory. If not,
  90  * we double the number of hash buckets and repeat. Using the
  91  * default settings these values translate to ~1 MB of hash tables
  92  * for each 1 GB of physical memory.
  93  *
  94  * The locking model:
  95  *
  96  * A new reference to a cache buffer can be obtained in two
  97  * ways: 1) via a hash table lookup using the DVA as a key,
  98  * or 2) via one of the ARC lists.  The arc_read() interface
  99  * uses method 1, while the internal arc algorithms for
 100  * adjusting the cache use method 2.  We therefore provide two
 101  * types of locks: 1) the hash table lock array, and 2) the
 102  * arc list locks.
 103  *
 104  * Buffers do not have their own mutexes, rather they rely on the
 105  * hash table mutexes for the bulk of their protection (i.e. most
 106  * fields in the arc_buf_hdr_t are protected by these mutexes). The
 107  * specific mutex is selected by taking its hash table value and masking
 108  * it by ht_lock_mask, which then produces an index into the mutex
 109  * table. The size of the lock table is derived from the amount of
 110  * physical memory, which is simply divided by
 111  * 2^zfs_arc_ht_lock_shift, giving the number of locks, with a
 112  * minimum of MIN_BUF_LOCKS.
 113  * fields in the arc_buf_hdr_t are protected by these mutexes).
 114  *
 115  * buf_hash_find() returns the appropriate mutex (held) when it
 116  * locates the requested buffer in the hash table.  It returns
 117  * NULL for the mutex if the buffer was not in the table.
 118  *
 119  * buf_hash_remove() expects the appropriate hash mutex to be
 120  * already held before it is invoked.
 121  *
 122  * Each arc state also has a mutex which is used to protect the
 123  * buffer list associated with the state.  When attempting to
 124  * obtain a hash table lock while holding an arc list lock you
 125  * must use: mutex_tryenter() to avoid deadlock.  Also note that
 126  * the active state mutex must be held before the ghost state mutex.
 127  *
```

```
127  * Arc buffers may have an associated eviction callback function.
128  * This function will be invoked prior to removing the buffer (e.g.
129  * in arc_do_user_evicts()).  Note however that the data associated
130  * with the buffer may be evicted prior to the callback.  The callback
131  * must be made with *no locks held* (to prevent deadlock).  Additionally,
132  * the users of callbacks must ensure that their private data is
133  * protected from simultaneous callbacks from arc_buf_evict()
134  * and arc_do_user_evicts().
135  *
136  * Note that the majority of the performance stats are manipulated
137  * with atomic operations.
138  *
139  * The L2ARC uses the l2arc_buflist_mtx global mutex for the following:
140  *
141  *      - L2ARC buflist creation
142  *      - L2ARC buflist eviction
143  *      - L2ARC write completion, which walks L2ARC buflists
144  *      - ARC header destruction, as it removes from L2ARC buflists
145  *      - ARC header release, as it removes from L2ARC buflists
146  */

148 #include <sys/spa.h>
149 #include <sys/zio.h>
150 #include <sys/zio_compress.h>
151 #include <sys/zfs_context.h>
152 #include <sys/arc.h>
153 #include <sys/refcount.h>
154 #include <sys/vdev.h>
155 #include <sys/vdev_impl.h>
156 #include <sys/dsl_pool.h>
157 #ifdef _KERNEL
158 #include <sys/vmsystm.h>
159 #include <vm/anon.h>
160 #include <sys/fs/swapnode.h>
161 #include <sys/dnlc.h>
162 #endif
163 #include <sys/callb.h>
164 #include <sys/kstat.h>
165 #include <zfs_fletcher.h>

167 #ifndef _KERNEL
168 /* set with ZFS_DEBUG=watch, to enable watchpoints on frozen buffers */
169 boolean_t arc_watch = B_FALSE;
170 int arc_procfd;
171 #endif

173 static kmutex_t         arc_reclaim_thr_lock;
174 static kcondvar_t       arc_reclaim_thr_cv;     /* used to signal reclaim thr */
175 static uint8_t          arc_thread_exit;

177 #define ARC_REDUCE_DNLC_PERCENT 3
178 uint_t arc_reduce_dnlc_percent = ARC_REDUCE_DNLC_PERCENT;

180 typedef enum arc_reclaim_strategy {
181         ARC_RECLAIM_AGGR,               /* Aggressive reclaim strategy */
182         ARC_RECLAIM_CONS                /* Conservative reclaim strategy */
183 } arc_reclaim_strategy_t;

185 /*
186  * The number of iterations through arc_evict_*() before we
187  * drop & reacquire the lock.
188  */
189 int arc_evict_iterations = 100;

191 /* number of seconds before growing cache again */
192 static int              arc_grow_retry = 60;
```

```
194 /* shift of arc_c for calculating both min and max arc_p */
195 static int              arc_p_min_shift = 4;

197 /* log2(fraction of arc to reclaim) */
198 static int              arc_shrink_shift = 5;

200 /*
201  * minimum lifespan of a prefetch block in clock ticks
202  * (initialized in arc_init())
203  */
204 static int              arc_min_prefetch_lifespan;

206 /*
207  * If this percent of memory is free, don't throttle.
208  */
209 int arc_lotsfree_percent = 10;

211 static int arc_dead;

213 /*
214  * The arc has filled available memory and has now warmed up.
215  */
216 static boolean_t arc_warm;

218 /*
219  * These tunables are for performance analysis.
220  */
221 uint64_t zfs_arc_max;
222 uint64_t zfs_arc_min;
223 uint64_t zfs_arc_meta_limit = 0;
224 int zfs_arc_grow_retry = 0;
225 int zfs_arc_shrink_shift = 0;
226 int zfs_arc_p_min_shift = 0;
227 int zfs_disable_dup_eviction = 0;

229 /*
230  * Used to calculate the size of ARC hash tables and number of hash locks.
231  * See big theory block comment at the start of this file.
232  */
233 uint64_t zfs_arc_ht_base_masklen = 13;
234 /*
235  * We want to allocate one hash lock for every 4GB of memory with a minimum
236  * of MIN_BUF_LOCKS.
237  */
238 uint64_t zfs_arc_ht_lock_shift = 32;
239 #define MIN_BUF_LOCKS   256

241 /*
242  * Note that buffers can be in one of 6 states:
243  *      ARC_anon        - anonymous (discussed below)
244  *      ARC_mru         - recently used, currently cached
245  *      ARC_mru_ghost   - recently used, no longer in cache
246  *      ARC_mfu         - frequently used, currently cached
247  *      ARC_mfu_ghost   - frequently used, no longer in cache
248  *      ARC_l2c_only    - exists in L2ARC but not other states
249  * When there are no active references to the buffer, they are
250  * are linked onto a list in one of these arc states.  These are
251  * the only buffers that can be evicted or deleted.  Within each
252  * state there are multiple lists, one for meta-data and one for
253  * non-meta-data.  Meta-data (indirect blocks, blocks of dnodes,
254  * etc.) is tracked separately so that it can be managed more
255  * explicitly: favored over data, limited explicitly.
256  *
257  * Anonymous buffers are buffers that are not associated with
258  * a DVA.  These are buffers that hold dirty block copies
```

```
259  * before they are written to stable storage.  By definition,
260  * they are "ref'd" and are considered part of arc_mru
261  * that cannot be freed.  Generally, they will aquire a DVA
262  * as they are written and migrate onto the arc_mru list.
263  *
264  * The ARC_l2c_only state is for buffers that are in the second
265  * level ARC but no longer in any of the ARC_m* lists.  The second
266  * level ARC itself may also contain buffers that are in any of
267  * the ARC_m* states - meaning that a buffer can exist in two
268  * places.  The reason for the ARC_l2c_only state is to keep the
269  * buffer header in the hash table, so that reads that hit the
270  * second level ARC benefit from these fast lookups.
271  */

273 typedef struct arc_state {
274         list_t  arcs_list[ARC_BUFC_NUMTYPES];   /* list of evictable buffers */
275         uint64_t arcs_lsize[ARC_BUFC_NUMTYPES]; /* amount of evictable data */
276         uint64_t arcs_size;     /* total amount of data in this state */
277         kmutex_t arcs_mtx;
278 } arc_state_t;
_____unchanged_portion_omitted_

591 #define BUF_LOCKS 256
628 typedef struct buf_hash_table {
629         uint64_t        ht_mask;
630         arc_buf_hdr_t   **ht_table;
631         struct ht_lock  *ht_locks;
632         uint64_t        ht_num_locks, ht_lock_mask;
595         struct ht_lock ht_locks[BUF_LOCKS];
633 } buf_hash_table_t;

635 static buf_hash_table_t buf_hash_table;

637 #define BUF_HASH_INDEX(spa, dva, birth) \
638         (buf_hash(spa, dva, birth) & buf_hash_table.ht_mask)
639 #define BUF_HASH_LOCK_NTRY(idx) \
640         (buf_hash_table.ht_locks[idx & buf_hash_table.ht_lock_mask])
602 #define BUF_HASH_LOCK_NTRY(idx) (buf_hash_table.ht_locks[idx & (BUF_LOCKS-1)])
641 #define BUF_HASH_LOCK(idx)      (&(BUF_HASH_LOCK_NTRY(idx).ht_lock))
642 #define HDR_LOCK(hdr) \
643         (BUF_HASH_LOCK(BUF_HASH_INDEX(hdr->b_spa, &hdr->b_dva, hdr->b_birth)))

645 uint64_t zfs_crc64_table[256];

647 /*
648  * Level 2 ARC
649  */

651 #define L2ARC_WRITE_SIZE        (8 * 1024 * 1024)       /* initial write max */
652 #define L2ARC_HEADROOM          2                       /* num of writes */
653 /*
654  * If we discover during ARC scan any buffers to be compressed, we boost
655  * our headroom for the next scanning cycle by this percentage multiple.
656  */
657 #define L2ARC_HEADROOM_BOOST    200
658 #define L2ARC_FEED_SECS         1                       /* caching interval secs */
659 #define L2ARC_FEED_MIN_MS       200                     /* min caching interval ms */

661 #define l2arc_writes_sent       ARCSTAT(arcstat_l2_writes_sent)
662 #define l2arc_writes_done       ARCSTAT(arcstat_l2_writes_done)

664 /* L2ARC Performance Tunables */
665 uint64_t l2arc_write_max = L2ARC_WRITE_SIZE;    /* default max write size */
666 uint64_t l2arc_write_boost = L2ARC_WRITE_SIZE;  /* extra write during warmup */
667 uint64_t l2arc_headroom = L2ARC_HEADROOM;       /* number of dev writes */
668 uint64_t l2arc_headroom_boost = L2ARC_HEADROOM_BOOST;
```

```
669 uint64_t l2arc_feed_secs = L2ARC_FEED_SECS;     /* interval seconds */
670 uint64_t l2arc_feed_min_ms = L2ARC_FEED_MIN_MS; /* min interval milliseconds */
671 boolean_t l2arc_noprefetch = B_TRUE;            /* don't cache prefetch bufs */
672 boolean_t l2arc_feed_again = B_TRUE;            /* turbo warmup */
673 boolean_t l2arc_norw = B_TRUE;                  /* no reads during writes */

675 /*
676  * L2ARC Internals
677  */
678 typedef struct l2arc_dev {
679         vdev_t                  *l2ad_vdev;     /* vdev */
680         spa_t                   *l2ad_spa;      /* spa */
681         uint64_t                l2ad_hand;      /* next write location */
682         uint64_t                l2ad_start;     /* first addr on device */
683         uint64_t                l2ad_end;       /* last addr on device */
684         uint64_t                l2ad_evict;     /* last addr eviction reached */
685         boolean_t               l2ad_first;     /* first sweep through */
686         boolean_t               l2ad_writing;   /* currently writing */
687         list_t                  *l2ad_buflist;  /* buffer list */
688         list_node_t             l2ad_node;      /* device list node */
689 } l2arc_dev_t;
_____unchanged_portion_omitted_

735 static kmutex_t l2arc_feed_thr_lock;
736 static kcondvar_t l2arc_feed_thr_cv;
737 static uint8_t l2arc_thread_exit;

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

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

748 static inline uint64_t
710 static uint64_t
749 buf_hash(uint64_t spa, const dva_t *dva, uint64_t birth)
750 {
751         uint8_t *vdva = (uint8_t *)dva;
752         uint64_t crc = -1ULL;
753         int i;

755         ASSERT(zfs_crc64_table[128] == ZFS_CRC64_POLY);

757         for (i = 0; i < sizeof (dva_t); i++)
758                 crc = (crc >> 8) ^ zfs_crc64_table[(crc ^ vdva[i]) & 0xFF];

760         crc ^= (spa>>8) ^ birth;

762         return (crc);
763 }
_____unchanged_portion_omitted_

872 /*
873  * Global data structures and functions for the buf kmem cache.
874  */
875 static kmem_cache_t *hdr_cache;
876 static kmem_cache_t *buf_cache;

878 static void
879 buf_fini(void)
880 {
881         int i;
```

```
 883            kmem_free(buf_hash_table.ht_table,
 884                    (buf_hash_table.ht_mask + 1) * sizeof (void *));

 886            for (i = 0; i < buf_hash_table.ht_num_locks; i++)
 847            for (i = 0; i < BUF_LOCKS; i++)
 887                    mutex_destroy(&buf_hash_table.ht_locks[i].ht_lock);
 888            kmem_free(buf_hash_table.ht_locks, sizeof (struct ht_lock) *
 889                buf_hash_table.ht_num_locks);
 890            kmem_cache_destroy(hdr_cache);
 891            kmem_cache_destroy(buf_cache);
 892  }
_____unchanged_portion_omitted_

 969  static void
 970  buf_init(void)
 971  {
 972            uint64_t          *ct;
 973            uint64_t          ht_masklen = 12;
 932            uint64_t hsize = 1ULL << 12;
 974            int               i, j;

 976            while ((1ULL << (ht_masklen + zfs_arc_ht_base_masklen)) <
 977                physmem * PAGESIZE)
 978                    ht_masklen++;
 979            buf_hash_table.ht_mask = (1ULL << ht_masklen) - 1;
 935            /*
 936             * The hash table is big enough to fill all of physical memory
 937             * with an average 64K block size.  The table will take up
 938             * totalmem*sizeof(void*)/64K (eg. 128KB/GB with 8-byte pointers).
 939             */
 940            while (hsize * 65536 < physmem * PAGESIZE)
 941                    hsize <<= 1;
 942  retry:
 943            buf_hash_table.ht_mask = hsize - 1;
 980            buf_hash_table.ht_table =
 981                kmem_zalloc((1ULL << ht_masklen) * sizeof (void *), KM_SLEEP);

 983            buf_hash_table.ht_num_locks = MAX((physmem * PAGESIZE) >>
 984                zfs_arc_ht_lock_shift, MIN_BUF_LOCKS);
 985            buf_hash_table.ht_lock_mask = buf_hash_table.ht_num_locks - 1;
 986            buf_hash_table.ht_locks = kmem_zalloc(sizeof (struct ht_lock) *
 987                buf_hash_table.ht_num_locks, KM_SLEEP);
 988            for (i = 0; i < buf_hash_table.ht_num_locks; i++) {
 989                    mutex_init(&buf_hash_table.ht_locks[i].ht_lock,
 990                        NULL, MUTEX_DEFAULT, NULL);
 945                kmem_zalloc(hsize * sizeof (void*), KM_NOSLEEP);
 946            if (buf_hash_table.ht_table == NULL) {
 947                    ASSERT(hsize > (1ULL << 8));
 948                    hsize >>= 1;
 949                    goto retry;
 991            }

 993            hdr_cache = kmem_cache_create("arc_buf_hdr_t", sizeof (arc_buf_hdr_t),
 994                0, hdr_cons, hdr_dest, hdr_recl, NULL, NULL, 0);
 995            buf_cache = kmem_cache_create("arc_buf_t", sizeof (arc_buf_t),
 996                0, buf_cons, buf_dest, NULL, NULL, NULL, 0);

 998            for (i = 0; i < 256; i++)
 999                    for (ct = zfs_crc64_table + i, *ct = i, j = 8; j > 0; j--)
1000                            *ct = (*ct >> 1) ^ (-(*ct & 1) & ZFS_CRC64_POLY);

 961            for (i = 0; i < BUF_LOCKS; i++) {
 962                    mutex_init(&buf_hash_table.ht_locks[i].ht_lock,
 963                        NULL, MUTEX_DEFAULT, NULL);
 964            }
1001  }
_____unchanged_portion_omitted_
```