

```
*****
 28038 Sat Aug 1 08:08:32 2015
new/usr/src/common/avl/avl.c
6091 avl_add doesn't assert on non-debug builds
Reviewed by: Andy Stornmont <astormont@racktopysystems.com>
*****
1 /*
2 * CDDL HEADER START
3 *
4 * The contents of this file are subject to the terms of the
5 * Common Development and Distribution License (the "License").
6 * You may not use this file except in compliance with the License.
7 *
8 * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9 * or http://www.opensolaris.org/os/licensing.
10 * See the License for the specific language governing permissions
11 * and limitations under the License.
12 *
13 * When distributing Covered Code, include this CDDL HEADER in each
14 * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 * If applicable, add the following below this CDDL HEADER, with the
16 * fields enclosed by brackets "[]" replaced with your own identifying
17 * information: Portions Copyright [yyyy] [name of copyright owner]
18 *
19 * CDDL HEADER END
20 */
21 /*
22 * Copyright 2009 Sun Microsystems, Inc. All rights reserved.
23 * Use is subject to license terms.
24 */
25 /*
26 * Copyright (c) 2014 by Delphix. All rights reserved.
27 * Copyright 2015 Nexenta Systems, Inc. All rights reserved.
28 #endif /* ! codereview */
30 */
31 /*
32 * AVL - generic AVL tree implementation for kernel use
33 *
34 * A complete description of AVL trees can be found in many CS textbooks.
35 *
36 * Here is a very brief overview. An AVL tree is a binary search tree that is
37 * almost perfectly balanced. By "almost" perfectly balanced, we mean that at
38 * any given node, the left and right subtrees are allowed to differ in height
39 * by at most 1 level.
40 *
41 * This relaxation from a perfectly balanced binary tree allows doing
42 * insertion and deletion relatively efficiently. Searching the tree is
43 * still a fast operation, roughly O(log(N)).
44 *
45 * The key to insertion and deletion is a set of tree manipulations called
46 * rotations, which bring unbalanced subtrees back into the semi-balanced state.
47 *
48 * This implementation of AVL trees has the following peculiarities:
49 *
50 * - The AVL specific data structures are physically embedded as fields
51 *   in the "using" data structures. To maintain generality the code
52 *   must constantly translate between "avl_node_t *" and containing
53 *   data structure "void **"s by adding/subtracting the avl_offset.
54 *
55 * - Since the AVL data is always embedded in other structures, there is
56 *   no locking or memory allocation in the AVL routines. This must be
57 *   provided for by the enclosing data structure's semantics. Typically,
58 *   avl_insert()/_add()/_remove()/avl_insert_here() require some kind of
59 *   exclusive write lock. Other operations require a read lock.
60 *
```

```
61 /*
62 * - The implementation uses iteration instead of explicit recursion,
63 *   since it is intended to run on limited size kernel stacks. Since
64 *   there is no recursion stack present to move "up" in the tree,
65 *   there is an explicit "parent" link in the avl_node_t.
66 *
67 * - The left/right children pointers of a node are in an array.
68 * In the code, variables (instead of constants) are used to represent
69 * left and right indices. The implementation is written as if it only
70 * dealt with left handed manipulations. By changing the value assigned
71 * to "left", the code also works for right handed trees. The
72 * following variables/terms are frequently used:
73 *
74 *     int left;           // 0 when dealing with left children,
75 *                         // 1 for dealing with right children
76 *
77 *     int left_heavy;    // -1 when left subtree is taller at some node,
78 *                         // +1 when right subtree is taller
79 *
80 *     int right;          // will be the opposite of left (0 or 1)
81 *     int right_heavy; // will be the opposite of left_heavy (-1 or 1)
82 *
83 *     int direction;    // 0 for "<" (ie. left child); 1 for ">" (right)
84 *
85 * Though it is a little more confusing to read the code, the approach
86 * allows using half as much code (and hence cache footprint) for tree
87 * manipulations and eliminates many conditional branches.
88 *
89 * - The avl_index_t is an opaque "cookie" used to find nodes at or
90 * adjacent to where a new value would be inserted in the tree. The value
91 * is a modified "avl_node_t *". The bottom bit (normally 0 for a
92 * pointer) is set to indicate if that the new node has a value greater
93 * than the value of the indicated "avl_node_t *".
94 *
95 * Note - in addition to userland (e.g. libavl and libutil) and the kernel
96 * (e.g. genunix), avl.c is compiled into ld.so and kmdb's genunix module,
97 * which each have their own compilation environments and subsequent
98 * requirements. Each of these environments must be considered when adding
99 * dependencies from avl.c.
100 */
101
102 #include <sys/types.h>
103 #include <sys/param.h>
104 #include <sys/debug.h>
105 #include <sys/avl.h>
106 #include <sys/cmn_err.h>
107
108 /*
109 * Small arrays to translate between balance (or diff) values and child indices.
110 *
111 * Code that deals with binary tree data structures will randomly use
112 * left and right children when examining a tree. C "if()" statements
113 * which evaluate randomly suffer from very poor hardware branch prediction.
114 * In this code we avoid some of the branch mispredictions by using the
115 * following translation arrays. They replace random branches with an
116 * additional memory reference. Since the translation arrays are both very
117 * small the data should remain efficiently in cache.
118 */
119 static const int avl_child2balance[2] = {-1, 1};
120 static const int avl_balance2child[] = {0, 0, 1};
121
122 /*
123 * Walk from one node to the previous valued node (ie. an infix walk
124 * towards the left). At any given node we do one of 2 things:
125 * 
```

```

127 * - If there is a left child, go to it, then to its rightmost descendant.
128 *
129 * - otherwise we return through parent nodes until we've come from a right
130 *   child.
131 *
132 * Return Value:
133 *   NULL - if at the end of the nodes
134 *   otherwise next node
135 */
136 void *
137 avl_walk(avl_tree_t *tree, void *oldnode, int left)
138 {
139     size_t off = tree->avl_offset;
140     avl_node_t *node = AVL_DATA2NODE(oldnode, off);
141     int right = 1 - left;
142     int was_child;

145     /*
146     * nowhere to walk to if tree is empty
147     */
148     if (node == NULL)
149         return (NULL);

151     /*
152     * Visit the previous valued node. There are two possibilities:
153     *
154     * If this node has a left child, go down one left, then all
155     * the way right.
156     */
157     if (node->avl_child[left] != NULL) {
158         for (node = node->avl_child[left];
159             node->avl_child[right] != NULL;
160             node = node->avl_child[right])
161             ;
162     /*
163     * Otherwise, return thru left children as far as we can.
164     */
165     } else {
166         for (;;) {
167             was_child = AVL_XCHILD(node);
168             node = AVL_XPARENT(node);
169             if (node == NULL)
170                 return (NULL);
171             if (was_child == right)
172                 break;
173         }
174     }

176     return (AVL_NODE2DATA(node, off));
177 }

179 /*
180 * Return the lowest valued node in a tree or NULL.
181 * (leftmost child from root of tree)
182 */
183 void *
184 avl_first(avl_tree_t *tree)
185 {
186     avl_node_t *node;
187     avl_node_t *prev = NULL;
188     size_t off = tree->avl_offset;

190     for (node = tree->avl_root; node != NULL; node = node->avl_child[0])
191         prev = node;

```

```

193     if (prev != NULL)
194         return (AVL_NODE2DATA(prev, off));
195     return (NULL);
196 }

198 /*
199 * Return the highest valued node in a tree or NULL.
200 * (rightmost child from root of tree)
201 */
202 void *
203 avl_last(avl_tree_t *tree)
204 {
205     avl_node_t *node;
206     avl_node_t *prev = NULL;
207     size_t off = tree->avl_offset;

209     for (node = tree->avl_root; node != NULL; node = node->avl_child[1])
210         prev = node;
211
212     if (prev != NULL)
213         return (AVL_NODE2DATA(prev, off));
214     return (NULL);
215 }

217 /*
218 * Access the node immediately before or after an insertion point.
219 *
220 * "avl_index_t" is a (avl_node_t *) with the bottom bit indicating a child
221 *
222 * Return value:
223 *   NULL: no node in the given direction
224 *   "void *" of the found tree node
225 */
226 void *
227 avl_nearest(avl_tree_t *tree, avl_index_t where, int direction)
228 {
229     int child = AVL_INDEX2CHILD(where);
230     avl_node_t *node = AVL_INDEX2NODE(where);
231     void *data;
232     size_t off = tree->avl_offset;

234     if (node == NULL) {
235         ASSERT(tree->avl_root == NULL);
236         return (NULL);
237     }
238     data = AVL_NODE2DATA(node, off);
239     if (child != direction)
240         return (data);

242     return (avl_walk(tree, data, direction));
243 }

246 /*
247 * Search for the node which contains "value". The algorithm is a
248 * simple binary tree search.
249 *
250 * Return value:
251 *   NULL: the value is not in the AVL tree
252 *           *where (if not NULL) is set to indicate the insertion point
253 *   "void *" of the found tree node
254 */
255 void *
256 avl_find(avl_tree_t *tree, const void *value, avl_index_t *where)
257 {
258     avl_node_t *node;

```

```

259     avl_node_t *prev = NULL;
260     int child = 0;
261     int diff;
262     size_t off = tree->avl_offset;
263
264     for (node = tree->avl_root; node != NULL;
265          node = node->avl_child[child]) {
266
267         prev = node;
268
269         diff = tree->avl_compar(value, AVL_NODE2DATA(node, off));
270         ASSERT(-1 <= diff && diff <= 1);
271         if (diff == 0) {
272 #ifdef DEBUG
273             if (where != NULL)
274                 *where = 0;
275 #endif
276             return (AVL_NODE2DATA(node, off));
277         }
278         child = avl_balance2child[1 + diff];
279
280     }
281
282     if (where != NULL)
283         *where = AVL_MKINDEX(prev, child);
284
285     return (NULL);
286 }
287
288 /*
289  * Perform a rotation to restore balance at the subtree given by depth.
290  *
291  * This routine is used by both insertion and deletion. The return value
292  * indicates:
293  *      0 : subtree did not change height
294  *      !0 : subtree was reduced in height
295  *
296  * The code is written as if handling left rotations, right rotations are
297  * symmetric and handled by swapping values of variables right/left[_heavy]
298  *
299  * On input balance is the "new" balance at "node". This value is either
300  * -2 or +2.
301  */
302
303 static int
304 avl_rotation(avl_tree_t *tree, avl_node_t *node, int balance)
305 {
306     int left = !(balance < 0);           /* when balance = -2, left will be 0 */
307     int right = 1 - left;
308     int left_heavy = balance >> 1;
309     int right_heavy = -left_heavy;
310     avl_node_t *parent = AVL_XPARENT(node);
311     avl_node_t *child = node->avl_child[left];
312     avl_node_t *cright;
313     avl_node_t *gchild;
314     avl_node_t *gright;
315     avl_node_t *gleft;
316     int which_child = AVL_XCHILD(node);
317     int child_bal = AVL_XBALANCE(child);
318
319     /* BEGIN CSTYLED */
320     /*
321      * case 1 : node is overly left heavy, the left child is balanced or
322      * also left heavy. This requires the following rotation.
323      *
324      *          (node bal:-2)

```

```

325
326
327
328
329
330
331
332 * becomes:
333
334
335
336
337
338
339
340
341
342 * we detect this situation by noting that child's balance is not
343 * right_heavy.
344 */
345 /* END CSTYLED */
346 if (child_bal != right_heavy) {

348
349     /*
350      * compute new balance of nodes
351      *
352      * If child used to be left heavy (now balanced) we reduced
353      * the height of this sub-tree -- used in "return...;" below
354     */
355     child_bal += right_heavy; /* adjust towards right */

356
357     /*
358      * move "cright" to be node's left child
359      */
360     cright = child->avl_child[right];
361     node->avl_child[left] = cright;
362     if (cright != NULL) {
363         AVL_SETPARENT(cright, node);
364         AVL_SETCHILD(cright, left);
365     }

366
367     /*
368      * move node to be child's right child
369      */
370     child->avl_child[right] = node;
371     AVL_SETBALANCE(node, -child_bal);
372     AVL_SETCHILD(node, right);
373     AVL_SETPARENT(node, child);

374
375     /*
376      * update the pointer into this subtree
377      */
378     AVL_SETBALANCE(child, child_bal);
379     AVL_SETCHILD(child, which_child);
380     AVL_SETPARENT(child, parent);
381     if (parent != NULL)
382         parent->avl_child[which_child] = child;
383     else
384         tree->avl_root = child;

385
386     return (child_bal == 0);
387 }

388 /* BEGIN CSTYLED */
389 /*
390  * case 2 : When node is left heavy, but child is right heavy we use

```

```

391     * a different rotation.
392     *
393     *           (node b:-2)
394     *           /   \
395     *           /   \
396     *           \   \
397     *           (child b:+1)
398     *           /   \
399     *           /   \
400     *           (gchild b: != 0)
401     *           /   \
402     *           /   \
403     *           gleft   gright
404
405     * becomes:
406
407     *           (gchild b:0)
408     *           /   \
409     *           /   \
410     *           (child b:?) (node b:?)
411     *           /   \
412     *           /   \
413     *           gleft   gright
414
415
416     * computing the new balances is more complicated. As an example:
417     *     if gchild was right_heavy, then child is now left heavy
418     *     else it is balanced
419 */
420 /* END CSTYLED */
421 gchild = child->avl_child[right];
422 gleft = gchild->avl_child[left];
423 gright = gchild->avl_child[right];
424
425     * move gright to left child of node and
426     *
427     * move gleft to right child of node
428     */
429 node->avl_child[left] = gright;
430 if (gright != NULL) {
431     AVL_SETPARENT(gright, node);
432     AVL_SETCHILD(gright, left);
433 }
434
435 child->avl_child[right] = gleft;
436 if (gleft != NULL) {
437     AVL_SETPARENT(gleft, child);
438     AVL_SETCHILD(gleft, right);
439 }
440
441     * move child to left child of gchild and
442     *
443     * move node to right child of gchild and
444     *
445     * fixup parent of all this to point to gchild
446     */
447 balance = AVL_XBALANCE(gchild);
448 gchild->avl_child[left] = child;
449 AVL_SETBALANCE(child, (balance == right_heavy ? left_heavy : 0));
450 AVL_SETPARENT(child, gchild);
451 AVL_SETCHILD(child, left);
452
453 gchild->avl_child[right] = node;
454 AVL_SETBALANCE(node, (balance == left_heavy ? right_heavy : 0));

```

```

455
456     AVL_SETPARENT(node, gchild);
457     AVL_SETCHILD(node, right);
458
459     AVL_SETBALANCE(gchild, 0);
460     AVL_SETPARENT(gchild, parent);
461     AVL_SETCHILD(gchild, which_child);
462     if (parent != NULL)
463         parent->avl_child[which_child] = gchild;
464     else
465         tree->avl_root = gchild;
466
467     return (1); /* the new tree is always shorter */
468 }
469
470 /*
471 * Insert a new node into an AVL tree at the specified (from avl_find()) place.
472 *
473 * Newly inserted nodes are always leaf nodes in the tree, since avl_find()
474 * searches out to the leaf positions. The avl_index_t indicates the node
475 * which will be the parent of the new node.
476 *
477 * After the node is inserted, a single rotation further up the tree may
478 * be necessary to maintain an acceptable AVL balance.
479 */
480 void
481 avl_insert(avl_tree_t *tree, void *new_data, avl_index_t where)
482 {
483     avl_node_t *node;
484     avl_node_t *parent = AVL_INDEX2NODE(where);
485     int old_balance;
486     int new_balance;
487     int which_child = AVL_INDEX2CHILD(where);
488     size_t off = tree->avl_offset;
489
490     ASSERT(tree);
491 #ifdef LP64
492     ASSERT((uintptr_t)new_data & 0x7) == 0;
493 #endif
494
495     node = AVL_DATA2NODE(new_data, off);
496
497     /*
498      * First, add the node to the tree at the indicated position.
499      */
500     ++tree->avl_numnodes;
501
502     node->avl_child[0] = NULL;
503     node->avl_child[1] = NULL;
504
505     AVL_SETCHILD(node, which_child);
506     AVL_SETBALANCE(node, 0);
507     AVL_SETPARENT(node, parent);
508     if (parent != NULL) {
509         ASSERT(parent->avl_child[which_child] == NULL);
510         parent->avl_child[which_child] = node;
511     } else {
512         ASSERT(tree->avl_root == NULL);
513         tree->avl_root = node;
514     }
515
516     /*
517      * Now, back up the tree modifying the balance of all nodes above the
518      * insertion point. If we get to a highly unbalanced ancestor, we
519      * need to do a rotation. If we back out of the tree we are done.
520      * If we brought any subtree into perfect balance (0), we are also done.
521      */
522

```

```

523     for (;;) {
524         node = parent;
525         if (node == NULL)
526             return;
527
528         /*
529          * Compute the new balance
530          */
531         old_balance = AVL_XBALANCE(node);
532         new_balance = old_balance + avl_child2balance[which_child];
533
534         /*
535          * If we introduced equal balance, then we are done immediately
536          */
537         if (new_balance == 0) {
538             AVL_SETBALANCE(node, 0);
539             return;
540         }
541
542         /*
543          * If both old and new are not zero we went
544          * from -1 to -2 balance, do a rotation.
545          */
546         if (old_balance != 0)
547             break;
548
549         AVL_SETBALANCE(node, new_balance);
550         parent = AVL_XPARENT(node);
551         which_child = AVL_XCHILD(node);
552     }
553
554     /*
555      * perform a rotation to fix the tree and return
556      */
557     (void) avl_rotation(tree, node, new_balance);
558 }

560 /**
561  * Insert "new_data" in "tree" in the given "direction" either after or
562  * before (AVL_AFTER, AVL_BEFORE) the data "here".
563  *
564  * Insertions can only be done at empty leaf points in the tree, therefore
565  * if the given child of the node is already present we move to either
566  * the AVL_PREV or AVL_NEXT and reverse the insertion direction. Since
567  * every other node in the tree is a leaf, this always works.
568  *
569  * To help developers using this interface, we assert that the new node
570  * is correctly ordered at every step of the way in DEBUG kernels.
571 */
572 void
573 avl_insert_here(
574     avl_tree_t *tree,
575     void *new_data,
576     void *here,
577     int direction)
578 {
579     avl_node_t *node;
580     int child = direction; /* rely on AVL_BEFORE == 0, AVL_AFTER == 1 */
581 #ifdef DEBUG
582     int diff;
583 #endif
584
585     ASSERT(tree != NULL);
586     ASSERT(new_data != NULL);
587     ASSERT(here != NULL);
588     ASSERT(direction == AVL_BEFORE || direction == AVL_AFTER);

```

```

590     /*
591      * If corresponding child of node is not NULL, go to the neighboring
592      * node and reverse the insertion direction.
593      */
594     node = AVL_DATA2NODE(here, tree->avl_offset);

595 #ifdef DEBUG
596     diff = tree->avl_compar(new_data, here);
597     ASSERT(-1 <= diff && diff <= 1);
598     ASSERT(diff != 0);
599     ASSERT(diff > 0 ? child == 1 : child == 0);
600 #endif

601     if (node->avl_child[child] != NULL) {
602         node = node->avl_child[child];
603         child = 1 - child;
604         while (node->avl_child[child] != NULL) {
605             #ifdef DEBUG
606                 diff = tree->avl_compar(new_data,
607                     AVL_NODE2DATA(node, tree->avl_offset));
608                 ASSERT(-1 <= diff && diff <= 1);
609                 ASSERT(diff != 0);
610                 ASSERT(diff > 0 ? child == 1 : child == 0);
611             #endif
612             node = node->avl_child[child];
613         }
614     }
615     #ifdef DEBUG
616         diff = tree->avl_compar(new_data,
617             AVL_NODE2DATA(node, tree->avl_offset));
618         ASSERT(-1 <= diff && diff <= 1);
619         ASSERT(diff != 0);
620         ASSERT(diff > 0 ? child == 1 : child == 0);
621     #endif
622     ASSERT(node->avl_child[child] == NULL);

623     avl_insert(tree, new_data, AVL_MKINDEX(node, child));

624     /*
625      * Add a new node to an AVL tree.
626      */
627     void
628     avl_add(avl_tree_t *tree, void *new_node)
629     {
630         avl_index_t where;
631
632         /*
633          * This is unfortunate. We want to call panic() here, even for
634          * non-DEBUG kernels. In userland, however, we can't depend on anything
635          * in libc or else the rtld build process gets confused.
636          * Thankfully, rtld provides us with its own assfail() so we can use
637          * that here. We use assfail() directly to get a nice error message
638          * in the core - much like what panic() does for crashdumps.
639          * in libc or else the rtld build process gets confused. So, all we can
640          * do in userland is resort to a normal ASSERT().
641         */
642         if (avl_find(tree, new_node, &where) != NULL)
643             panic("avl_find() succeeded inside avl_add()");
644 #ifdef _KERNEL
645             panic("avl_find() succeeded inside avl_add()");
646 #else
647             (void) assfail("avl_find() succeeded inside avl_add()", __FILE__, __LINE__);
648 #endif
649             ASSERT(0);
650     35
651 #endif

```

```
new/usr/src/common/avl/avl.c  
652     avl_insert(tree, new_node, where);  
653 }  
unchanged portion omitted
```

11

```
*****
1341 Sat Aug 1 08:08:32 2015
new/usr/src/lib/libavl/Makefile.com
6091 avl_add doesn't assert on non-debug builds
Reviewed by: Andy Storment <astormont@racktopysystems.com>
*****
1 #
2 # CDDL HEADER START
3 #
4 # The contents of this file are subject to the terms of the
5 # Common Development and Distribution License (the "License").
6 # You may not use this file except in compliance with the License.
7 #
8 # You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
9 # or http://www.opensolaris.org/os/licensing.
10 # See the License for the specific language governing permissions
11 # and limitations under the License.
12 #
13 # When distributing Covered Code, include this CDDL HEADER in each
14 # file and include the License file at usr/src/OPENSOLARIS.LICENSE.
15 # If applicable, add the following below this CDDL HEADER, with the
16 # fields enclosed by brackets "[]" replaced with your own identifying
17 # information: Portions Copyright [yyyy] [name of copyright owner]
18 #
19 # CDDL HEADER END
20 #
21 #
22 # Copyright 2008 Sun Microsystems, Inc. All rights reserved.
23 # Use is subject to license terms.
24 #
25 # ident "%Z%%M% %I%     %E% SMI"
26 #

28 LIBRARY=      libavl.a
29 VERS=         .1
30 OBJECTS=      avl.o

32 include ../../Makefile.lib
33 include ../../Makefile.rootfs

35 LIBS =        $(DYNLIB) $(LINTLIB)
36 SRCS =        $(COMDIR)/avl.c

38 $(LINTLIB) := SRCS = $(SRCDIR)/$(LINTSRC)

40 COMDIR =      $(SRC)/common/avl

42 LDLIBS +=     -lc
43 #endif /* ! codereview */
44 CFLAGS +=     $(CCVERBOSE)

46 .KEEP_STATE:

48 all : $(LIBS)

50 lint : lintcheck

52 pics/%.o:      $(COMDIR)/%.c
53           $(COMPILE.c) -o $@ $<
54           $(POST_PROCESS_O)

56 include ../../Makefile.targ
```