

new/usr/src/uts/common/vm/vm_pvn.c

```
*****
31976 Thu Jan  8 09:14:29 2015
new/usr/src/uts/common/vm/vm_pvn.c
5384 pvn_getpages may assert in valid scenarios
*****
```

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) 1986, 2010, Oracle and/or its affiliates. All rights reserved.
23 * Copyright 2015 Nexenta Systems, Inc. All rights reserved.
24 #endif /* ! codereview */
25 */

27 /* Copyright (c) 1984, 1986, 1987, 1988, 1989 AT&T */
28 /* All Rights Reserved */

30 /*
31 * University Copyright- Copyright (c) 1982, 1986, 1988
32 * The Regents of the University of California
33 * All Rights Reserved
34 *
35 * University Acknowledgment- Portions of this document are derived from
36 * software developed by the University of California, Berkeley, and its
37 * contributors.
38 */

40 /*
41 * VM - paged vnode.
42 *
43 * This file supplies vm support for the vnode operations that deal with pages.
44 */
45 #include <sys/types.h>
46 #include <sys/t_lock.h>
47 #include <sys/param.h>
48 #include <sys/sysmacros.h>
49 #include <sys/sysm.h>
50 #include <sys/time.h>
51 #include <sys/buf.h>
52 #include <sys/vnode.h>
53 #include <sys/uio.h>
54 #include <sys/vmsystm.h>
55 #include <sys/rman.h>
56 #include <sys/vfs.h>
57 #include <sys/cred.h>
58 #include <sys/user.h>
59 #include <sys/kmem.h>
60 #include <sys/cmn_err.h>
61 #include <sys/debug.h>

1

new/usr/src/uts/common/vm/vm_pvn.c

```
62 #include <sys/cpuvar.h>  
63 #include <sys/vtrace.h>  
64 #include <sys/tnf_probe.h>  
  
66 #include <vm/hat.h>  
67 #include <vm/as.h>  
68 #include <vm/seg.h>  
69 #include <vm/rm.h>  
70 #include <vm/pvn.h>  
71 #include <vm/page.h>  
72 #include <vm/seg_map.h>  
73 #include <vm/seg_kmem.h>  
74 #include <sys/fs/swapnode.h>  
  
76 int pvn_nofodklust = 0;  
77 int pvn_write_noklust = 0;  
  
79 uint_t pvn_vmodsort_supported = 0; /* set if HAT supports VMODSORT */  
80 uint_t pvn_vmodsort_disable = 0; /* set in /etc/system to disable HAT */  
81 /* support for vmodsort for testing */  
  
83 static struct kmem_cache *marker_cache = NULL;  
  
85 /*  
86 * Find the largest contiguous block which contains 'addr' for file offset  
87 * 'offset' in it while living within the file system block sizes ('vp_off'  
88 * and 'vp_len') and the address space limits for which no pages currently  
89 * exist and which map to consecutive file offsets.  
90 */  
91 page_t *  
92 pvn_read_kluster(  
93     struct vnode *vp,  
94     u_offset_t off,  
95     struct seg *seg,  
96     caddr_t addr,  
97     u_offset_t *offp,  
98     size_t *lenp,  
99     u_offset_t vp_off,  
100    size_t vp_len,  
101    int isra)  
102 {  
103     ssize_t deltaf, deltab;  
104     page_t *pp;  
105     page_t *plist = NULL;  
106     spgcnt_t pagesavail;  
107     u_offset_t vp_end;  
  
109     ASSERT(off >= vp_off && off < vp_off + vp_len);  
  
111     /*  
112      * We only want to do klustering/read ahead if there  
113      * is more than minfree pages currently available.  
114      */  
115     pagesavail = freemem - minfree;  
  
117     if (pagesavail <= 0)  
118         if (isra)  
119             return ((page_t *)NULL); /* ra case - give up */  
120         else  
121             pagesavail = 1; /* must return a page */  
  
123     /* We calculate in pages instead of bytes due to 32-bit overflows */  
124     if (pagesavail < (spgcnt_t)btopr(vp_len)) {  
125         /*  
126          * Don't have enough free memory for the  
127          * max request, try sizing down vp request.
```

2

```

128         */
129         deltab = (ssize_t)(off - vp_off);
130         vp_len -= deltab;
131         vp_off += deltab;
132         if (pagesavail < btopr(vp_len)) {
133             /*
134              * Still not enough memory, just settle for
135              * pagesavail which is at least 1.
136              */
137             vp_len = ptob(pagesavail);
138         }
139     }
140
141     vp_end = vp_off + vp_len;
142     ASSERT(off >= vp_off && off < vp_end);
143
144     if (isra && SEGOP_KLUSTER(seg, addr, 0))
145         return ((page_t *)NULL); /* segment driver says no */
146
147     if ((plist = page_create_va(vp, off,
148         PAGESIZE, PG_EXCL | PG_WAIT, seg, addr)) == NULL)
149         return ((page_t *)NULL);
150
151     if (vp_len <= PAGESIZE || pvn_nofodklust) {
152         *offp = off;
153         *lenp = MIN(vp_len, PAGESIZE);
154     } else {
155         /*
156          * Scan back from front by incrementing "deltab" and
157          * comparing "off" with "vp_off + deltab" to avoid
158          * "signed" versus "unsigned" conversion problems.
159          */
160         for (deltab = PAGESIZE; off >= vp_off + deltab;
161             deltab += PAGESIZE) {
162             /*
163               * Call back to the segment driver to verify that
164               * the klustering/read ahead operation makes sense.
165               */
166             if (SEGOP_KLUSTER(seg, addr, -deltab))
167                 break; /* page not eligible */
168             if ((pp = page_create_va(vp, off - deltab,
169                 PAGESIZE, PG_EXCL, seg, addr - deltab))
170                 == NULL)
171                 break; /* already have the page */
172             /*
173               * Add page to front of page list.
174               */
175             page_add(&plist, pp);
176         }
177         deltab -= PAGESIZE;
178
179         /* scan forward from front */
180         for (deltaf = PAGESIZE; off + deltaf < vp_end;
181             deltaf += PAGESIZE) {
182             /*
183               * Call back to the segment driver to verify that
184               * the klustering/read ahead operation makes sense.
185               */
186             if (SEGOP_KLUSTER(seg, addr, deltaf))
187                 break; /* page not file extension */
188             if ((pp = page_create_va(vp, off + deltaf,
189                 PAGESIZE, PG_EXCL, seg, addr + deltaf))
190                 == NULL)
191                 break; /* already have page */
192             /*

```

```

194                                         * Add page to end of page list.
195                                         */
196                                         page_add(&plist, pp);
197                                         plist = plist->p_next;
198                                     }
199                                     *offp = off = off - deltab;
200                                     *lenp = deltab + deltaf;
201                                     ASSERT(off >= vp_off);
202
203                                     /*
204                                      * If we ended up getting more than was actually
205                                      * requested, retract the returned length to only
206                                      * reflect what was requested. This might happen
207                                      * if we were allowed to kluster pages across a
208                                      * span of (say) 5 frags, and frag size is less
209                                      * than PAGESIZE. We need a whole number of
210                                      * pages to contain those frags, but the returned
211                                      * size should only allow the returned range to
212                                      * extend as far as the end of the frags.
213                                      */
214                                     if ((vp_off + vp_len) < (off + *lenp)) {
215                                         ASSERT(vp_end > off);
216                                         *lenp = vp_end - off;
217                                     }
218                                     TRACE_3(TR_FAC_VM, TR_PVN_READ_KLUSTER,
219                                         "pvn_read_kluster:seg %p addr %x isra %x",
220                                         seg, addr, isra);
221                                     return (plist);
222
223     /*
224      * Handle pages for this vnode on either side of the page "pp"
225      * which has been locked by the caller. This routine will also
226      * do klustering in the range [vp_off, vp_off + vp_len] up
227      * until a page which is not found. The offset and length
228      * of pages included is returned in "*offp" and "*lenp".
229      *
230      * Returns a list of dirty locked pages all ready to be
231      * written back.
232      */
233     page_t *
234     pvn_write_kluster(
235         struct vnode *vp,
236         page_t *pp,
237         u_offset_t *offp,
238         size_t *lenp,
239         u_offset_t vp_off,
240         size_t vp_len,
241         int flags)
242     {
243         u_offset_t off;
244         page_t *dirty;
245         size_t deltab, deltaf;
246         se_t se;
247         u_offset_t vp_end;
248
249         off = pp->p_offset;
250
251         /*
252          * Kustering should not be done if we are invalidating
253          * pages since we could destroy pages that belong to
254          * some other process if this is a swap vnode.
255          */
256         if (pvn_write_noklust || ((flags & B_INVAL) && IS_SWAPVP(vp))) {
257             *offp = off;
258
259

```

```

260             *lenp = PAGESIZE;
261             return (pp);
262     }
263
264     if (flags & (B_FREE | B_INVAL))
265         se = SE_EXCL;
266     else
267         se = SE_SHARED;
268
269     dirty = pp;
270
271     /* Scan backwards looking for pages to kluster by incrementing
272      * "deltab" and comparing "off" with "vp_off + deltab" to
273      * avoid "signed" versus "unsigned" conversion problems.
274      */
275     for (deltab = PAGESIZE; off >= vp_off + deltab; deltab += PAGESIZE) {
276         pp = page_lookup_nowait(vp, off - deltab, se);
277         if (pp == NULL)
278             break;           /* page not found */
279         if (pvn_getdirty(pp, flags | B_DELWRI) == 0)
280             break;
281         page_add(&dirty, pp);
282     }
283     deltab -= PAGESIZE;
284
285     vp_end = vp_off + vp_len;
286     /* now scan forwards looking for pages to kluster */
287     for (deltab = PAGESIZE; off + deltab < vp_end; deltab += PAGESIZE) {
288         pp = page_lookup_nowait(vp, off + deltab, se);
289         if (pp == NULL)
290             break;           /* page not found */
291         if (pvn_getdirty(pp, flags | B_DELWRI) == 0)
292             break;
293         page_add(&dirty, pp);
294         dirty = dirty->p_next;
295     }
296
297     *offp = off - deltab;
298     *lenp = deltab + deltabf;
299     return (dirty);
300 }
301
302 /*
303  * Generic entry point used to release the "shared/exclusive" lock
304  * and the "p_iolock" on pages after i/o is complete.
305  */
306 void
307 pvn_io_done(page_t *plist)
308 {
309     page_t *pp;
310
311     while (plist != NULL) {
312         pp = plist;
313         page_sub(&plist, pp);
314         page_io_unlock(pp);
315         page_unlock(pp);
316     }
317 }
318
319 /*
320  * Entry point to be used by file system getpage subr's and
321  * other such routines which either want to unlock pages (B_ASYNC
322  * request) or destroy a list of pages if an error occurred.
323  */
324 void
325 pvn_read_done(page_t *plist, int flags)

```

```

326 {
327     page_t *pp;
328
329     while (plist != NULL) {
330         pp = plist;
331         page_sub(&plist, pp);
332         page_io_unlock(pp);
333         if (flags & B_ERROR) {
334             /*LINTED: constant in conditional context*/
335             VN_DISPOSE(pp, B_INVAL, 0, kcred);
336         } else {
337             (void) page_release(pp, 0);
338         }
339     }
340 }
341
342 /*
343  * Automagic pageout.
344  * When memory gets tight, start freeing pages popping out of the
345  * write queue.
346  */
347 int     write_free = 1;
348 pgcnt_t pages_before_pager = 200;          /* LMXXX */
349
350 /*
351  * Routine to be called when page-out's complete.
352  * The caller, typically VOP_PUTPAGE, has to explicitly call this routine
353  * after waiting for i/o to complete (biowait) to free the list of
354  * pages associated with the buffer. These pages must be locked
355  * before i/o is initiated.
356  *
357  * If a write error occurs, the pages are marked as modified
358  * so the write will be re-tried later.
359 */
360
361 void
362 pvn_write_done(page_t *plist, int flags)
363 {
364     int dfree = 0;
365     int pgrec = 0;
366     int pcout = 0;
367     int ppgout = 0;
368     int anonpgout = 0;
369     int anonfree = 0;
370     int fspgout = 0;
371     int fsfree = 0;
372     int execpgout = 0;
373     int execfree = 0;
374     page_t *pp;
375     struct cpu *cpup;
376     struct vnode *vp = NULL;          /* for probe */
377     uint_t pattr;
378     kmutex_t *vphm = NULL;
379
380     ASSERT((flags & B_READ) == 0);
381
382     /*
383      * If we are about to start paging anyway, start freeing pages.
384      */
385     if (write_free && freemem < lotsfree + pages_before_pager &&
386         (flags & B_ERROR) == 0) {
387         flags |= B_FREE;
388     }
389
390     /*
391      * Handle each page involved in the i/o operation.
392

```

```

392     */
393     while (plist != NULL) {
394         pp = plist;
395         ASSERT(PAGE_LOCKED(pp) && page_iolock_assert(pp));
396         page_sub(&plist, pp);
397
398         /* Kernel probe support */
399         if (vp == NULL)
400             vp = pp->p_vnode;
401
402         if (((flags & B_ERROR) == 0) && IS_VMODSORT(vp)) {
403             /*
404             * Move page to the top of the v_page list.
405             * Skip pages modified during IO.
406             */
407             vphm = page_vnode_mutex(vp);
408             mutex_enter(vphm);
409             if ((pp->p_vpnext != pp) && !hat_ismod(pp)) {
410                 page_vpsub(&vp->v_pages, pp);
411                 page_vpadd(&vp->v_pages, pp);
412             }
413             mutex_exit(vphm);
414         }
415
416         if (flags & B_ERROR) {
417             /*
418             * Write operation failed. We don't want
419             * to destroy (or free) the page unless B_FORCE
420             * is set. We set the mod bit again and release
421             * all locks on the page so that it will get written
422             * back again later when things are hopefully
423             * better again.
424             * If B_INVAL and B_FORCE is set we really have
425             * to destroy the page.
426             */
427             if ((flags & (B_INVAL|B_FORCE)) == (B_INVAL|B_FORCE)) {
428                 page_io_unlock(pp);
429                 /*LINTED: constant in conditional context*/
430                 VN_DISPOSE(pp, B_INVAL, 0, kcred);
431             } else {
432                 hat_setmod_only(pp);
433                 page_io_unlock(pp);
434                 page_unlock(pp);
435             }
436         } else if (flags & B_INVAL) {
437             /*
438             * XXX - Failed writes with B_INVAL set are
439             * not handled appropriately.
440             */
441             page_io_unlock(pp);
442             /*LINTED: constant in conditional context*/
443             VN_DISPOSE(pp, B_INVAL, 0, kcred);
444         } else if (flags & B_FREE || !hat_page_is_mapped(pp)) {
445             /*
446             * Update statistics for pages being paged out
447             */
448             if (pp->p_vnode) {
449                 if (IS_SWAPFSVP(pp->p_vnode)) {
450                     anonpgout++;
451                 } else {
452                     if (pp->p_vnode->v_flag & VVMEXEC) {
453                         execpgout++;
454                     } else {
455                         fspgout++;
456                     }
457                 }
458             }
459         }
460     }
461     page_io_unlock(pp);
462     pgout = 1;
463     ppgout++;
464     TRACE_1(TR_FAC_VM, TR_PAGE_WS_OUT,
465             "page_ws_out:pp %p", pp);
466
467     /*
468      * The page_struct_lock need not be acquired to
469      * examine "p_lckcnt" and "p_cowcnt" since we'll
470      * have an "exclusive" lock if the upgrade succeeds.
471      */
472     if (page_tryupgrade(pp) &&
473         pp->p_lckcnt == 0 && pp->p_cowcnt == 0) {
474         /*
475          * Check if someone has reclaimed the
476          * page. If ref and mod are not set, no
477          * one is using it so we can free it.
478          * The rest of the system is careful
479          * to use the NOSYNC flag to unload
480          * translations set up for i/o w/o
481          * affecting ref and mod bits.
482          */
483          /*
484          * Obtain a copy of the real hardware
485          * mod bit using hat_pagesync(pp, HAT_DONTZERO)
486          * to avoid having to flush the cache.
487          */
488          ppattr = hat_pagesync(pp, HAT_SYNC_DONTZERO |
489                               HAT_SYNC_STOPON_MOD);
490
491         ck_refmod:
492             if (!(ppattr & (P_REF | P_MOD))) {
493                 if (hat_page_is_mapped(pp)) {
494                     /*
495                      * Doesn't look like the page
496                      * was modified so now we
497                      * really have to unload the
498                      * translations. Meanwhile
499                      * another CPU could've
500                      * modified it so we have to
501                      * check again. We don't loop
502                      * forever here because now
503                      * the translations are gone
504                      * and no one can get a new one
505                      * since we have the "exclusive"
506                      * lock on the page.
507                     */
508                     (void) hat_pageunload(pp,
509                                         HAT_FORCE_PGUNLOAD);
510                     ppattr = hat_page_getattr(pp,
511                                         P_REF | P_MOD);
512                     goto ck_refmod;
513
514             /*
515              * Update statistics for pages being
516              * freed
517              */
518             if (pp->p_vnode) {
519                 if (IS_SWAPFSVP(pp->p_vnode)) {
520                     anonfree++;
521                 } else {
522                     if (pp->p_vnode->v_flag
523                         & VVMEXEC) {
524                         execfree++;
525                     } else {
526                         fsfree++;
527                     }
528                 }
529             }
530         }
531     }
532 }
```

```

458     }
459     page_io_unlock(pp);
460     pgout = 1;
461     ppgout++;
462     TRACE_1(TR_FAC_VM, TR_PAGE_WS_OUT,
463             "page_ws_out:pp %p", pp);
464
465     /*
466      * The page_struct_lock need not be acquired to
467      * examine "p_lckcnt" and "p_cowcnt" since we'll
468      * have an "exclusive" lock if the upgrade succeeds.
469      */
470     if (page_tryupgrade(pp) &&
471         pp->p_lckcnt == 0 && pp->p_cowcnt == 0) {
472         /*
473          * Check if someone has reclaimed the
474          * page. If ref and mod are not set, no
475          * one is using it so we can free it.
476          * The rest of the system is careful
477          * to use the NOSYNC flag to unload
478          * translations set up for i/o w/o
479          * affecting ref and mod bits.
480          */
481          /*
482          * Obtain a copy of the real hardware
483          * mod bit using hat_pagesync(pp, HAT_DONTZERO)
484          * to avoid having to flush the cache.
485          */
486          ppattr = hat_pagesync(pp, HAT_SYNC_DONTZERO |
487                               HAT_SYNC_STOPON_MOD);
488
489         ck_refmod:
490             if (!(ppattr & (P_REF | P_MOD))) {
491                 if (hat_page_is_mapped(pp)) {
492                     /*
493                      * Doesn't look like the page
494                      * was modified so now we
495                      * really have to unload the
496                      * translations. Meanwhile
497                      * another CPU could've
498                      * modified it so we have to
499                      * check again. We don't loop
500                      * forever here because now
501                      * the translations are gone
502                      * and no one can get a new one
503                      * since we have the "exclusive"
504                      * lock on the page.
505                     */
506                     (void) hat_pageunload(pp,
507                                         HAT_FORCE_PGUNLOAD);
508                     ppattr = hat_page_getattr(pp,
509                                         P_REF | P_MOD);
510                     goto ck_refmod;
511
512             /*
513              * Update statistics for pages being
514              * freed
515              */
516             if (pp->p_vnode) {
517                 if (IS_SWAPFSVP(pp->p_vnode)) {
518                     anonfree++;
519                 } else {
520                     if (pp->p_vnode->v_flag
521                         & VVMEXEC) {
522                         execfree++;
523                     } else {
524                         fsfree++;
525                     }
526                 }
527             }
528         }
529     }
530 }
```

```

524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574 }

575 /*
576 * Flags are composed of {B_ASYNC, B_INVAL, B_FREE, B_DONTNEED, B_DELWRI,
577 * B_TRUNC, B_FORCE}. B_DELWRI indicates that this page is part of a kluster
578 * operation and is only to be considered if it doesn't involve any
579 * waiting here. B_TRUNC indicates that the file is being truncated
580 * and so no i/o needs to be done. B_FORCE indicates that the page
581 * must be destroyed so don't try writing it out.
582 *
583 * The caller must ensure that the page is locked. Returns 1, if
584 * the page should be written back (the "iolock" is held in this
585 * case), or 0 if the page has been dealt with or has been
586 * unlocked.
587 */
588 int
589

```

```

590 pvn_getdirty(page_t *pp, int flags)
591 {
592     ASSERT((flags & (B_INVAL | B_FREE)) ?
593             PAGE_EXCL(pp) : PAGE_SHARED(pp));
594     ASSERT(PP_ISFREE(pp) == 0);
595
596     /*
597      * If trying to invalidate or free a logically 'locked' page,
598      * forget it. Don't need page_struct_lock to check p_lckcnt and
599      * p_cowcnt as the page is exclusively locked.
600      */
601     if ((flags & (B_INVAL | B_FREE)) && !(flags & (B_TRUNC|B_FORCE)) &&
602         (pp->p_lckcnt != 0 || pp->p_cowcnt != 0)) {
603         page_unlock(pp);
604         return (0);
605     }
606
607     /*
608      * Now acquire the i/o lock so we can add it to the dirty
609      * list (if necessary). We avoid blocking on the i/o lock
610      * in the following cases:
611      *
612      * If B_DELWRI is set, which implies that this request is
613      * due to a klustering operation.
614      *
615      * If this is an async (B_ASYNC) operation and we are not doing
616      * invalidation (B_INVAL) [The current i/o or fsflush will ensure
617      * that the page is written out].
618      */
619     if ((flags & B_DELWRI) || ((flags & (B_INVAL | B_ASYNC)) == B_ASYNC)) {
620         if (!page_io_trylock(pp)) {
621             page_unlock(pp);
622             return (0);
623         }
624         page_io_lock(pp);
625     }
626
627     /*
628      * If we want to free or invalidate the page then
629      * we need to unload it so that anyone who wants
630      * it will have to take a minor fault to get it.
631      * Otherwise, we're just writing the page back so we
632      * need to sync up the hardware and software mod bit to
633      * detect any future modifications. We clear the
634      * software mod bit when we put the page on the dirty
635      * list.
636      */
637     if (flags & (B_INVAL | B_FREE)) {
638         (void) hat_pageunload(pp, HAT_FORCE_PGUNLOAD);
639     } else {
640         (void) hat_pagesync(pp, HAT_SYNC_ZERORM);
641     }
642
643     if (!hat_ismod(pp) || (flags & B_TRUNC)) {
644         /*
645          * Don't need to add it to the
646          * list after all.
647          */
648         page_io_unlock(pp);
649         if (flags & B_INVAL) {
650             /*LINTED: constant in conditional context*/
651             VN_DISPOSE(pp, B_INVAL, 0, kcred);
652         } else if (flags & B_FREE) {
653             /*LINTED: constant in conditional context*/
654             VN_DISPOSE(pp, B_FREE, (flags & B_DONTNEED), kcred);
655         }
656     }
657 }

```

```

656     } else {
657         /*
658          * This is advisory path for the callers
659          * of VOP_PUTPAGE() who prefer freeing the
660          * page _only_ if no one else is accessing it.
661          * E.g. segmap_release()
662          *
663          * The above hat_ismod() check is useless because:
664          * (1) we may not be holding SE_EXCL lock;
665          * (2) we've not unloaded _all_ translations
666          *
667          * Let page_release() do the heavy-lifting.
668          */
669     (void) page_release(pp, 1);
670 }
671 return (0);
672 }

674 /*
675  * Page is dirty, get it ready for the write back
676  * and add page to the dirty list.
677  */
678 hat_clrrefmod(pp);

680 /*
681  * If we're going to free the page when we're done
682  * then we can let others try to use it starting now.
683  * We'll detect the fact that they used it when the
684  * i/o is done and avoid freeing the page.
685  */
686 if (flags & B_FREE)
687     page_downgrade(pp);

690 TRACE_1(TR_FAC_VM, TR_PVN_GETDIRTY, "pvn_getdirty:pp %p", pp);
692 return (1);
693 }

696 /*ARGSUSED*/
697 static int
698 marker_constructor(void *buf, void *cdrarg, int kmfflags)
699 {
700     page_t *mark = buf;
701     bzero(mark, sizeof (page_t));
702     mark->p_hash = PVN_VPLIST_HASH_TAG;
703     return (0);
704 }

706 void
707 pvn_init()
708 {
709     if (pvn_vmodsort_disable == 0)
710         pvn_vmodsort_supported = hat_supported(HAT_VMODSORT, NULL);
711     marker_cache = kmem_cache_create("marker_cache",
712         sizeof (page_t), 0, marker_constructor,
713         NULL, NULL, NULL, NULL, 0);
714 }

717 /*
718  * Process a vnode's page list for all pages whose offset is >= off.
719  * Pages are to either be free'd, invalidated, or written back to disk.
720  *
721  * An "exclusive" lock is acquired for each page if B_INVAL or B_FREE

```

```

722     * is specified, otherwise they are "shared" locked.
723     *
724     * Flags are {B_ASYNC, B_INVAL, B_FREE, B_DONTNEED, B_TRUNC}
725     *
726     * Special marker page_t's are inserted in the list in order
727     * to keep track of where we are in the list when locks are dropped.
728     *
729     * Note the list is circular and insertions can happen only at the
730     * head and tail of the list. The algorithm ensures visiting all pages
731     * on the list in the following way:
732     *
733     * Drop two marker pages at the end of the list.
734     *
735     * Move one marker page backwards towards the start of the list until
736     * it is at the list head, processing the pages passed along the way.
737     *
738     * Due to race conditions when the vphm mutex is dropped, additional pages
739     * can be added to either end of the list, so we'll continue to move
740     * the marker and process pages until it is up against the end marker.
741     *
742     * There is one special exit condition. If we are processing a VMODSORT
743     * vnode and only writing back modified pages, we can stop as soon as
744     * we run into an unmodified page. This makes fsync(3) operations fast.
745     */
746 int
747 pvn_vplist_dirty(
748     vnode_t           *vp,
749     u_offset_t        off,
750     int               (*putapage)(vnode_t *, page_t *, u_offset_t *,
751                               size_t *, int, cred_t *),
752     int               flags,
753     cred_t            *cred)
754 {
755     page_t             *pp;
756     page_t             *mark;           /* marker page that moves toward head */
757     page_t             *end;            /* marker page at end of list */
758     int                err = 0;
759     int                error;
760     kmutex_t           *vphm;
761     se_t               se;
762     page_t             **where_to_move;

764     ASSERT(vp->v_type != VCHR);

766     if (vp->v_pages == NULL)
767         return (0);

770     /*
771      * Serialize vplist_dirty operations on this vnode by setting VVMLOCK.
772      *
773      * Don't block on VVMLOCK if B_ASYNC is set. This prevents sync()
774      * from getting blocked while flushing pages to a dead NFS server.
775      */
776     mutex_enter(&vp->v_lock);
777     if ((vp->v_flag & VVMLOCK) && (flags & B_ASYNC)) {
778         mutex_exit(&vp->v_lock);
779         return (EAGAIN);
780     }

782     while (vp->v_flag & VVMLOCK)
783         cv_wait(&vp->v_cv, &vp->v_lock);

785     if (vp->v_pages == NULL) {
786         mutex_exit(&vp->v_lock);
787         return (0);

```

```

788         }
790
791         vp->v_flag |= VVMLOCK;
792         mutex_exit(&vp->v_lock);
793
794         /*
795          * Set up the marker pages used to walk the list
796          */
797         end = kmem_cache_alloc(marker_cache, KM_SLEEP);
798         end->p_vnode = vp;
799         end->p_offset = (u_offset_t)-2;
800         mark = kmem_cache_alloc(marker_cache, KM_SLEEP);
801         mark->p_vnode = vp;
802         mark->p_offset = (u_offset_t)-1;
803
804         /*
805          * Grab the lock protecting the vnode's page list
806          * note that this lock is dropped at times in the loop.
807          */
808         vphm = page_vnode_mutex(vp);
809         mutex_enter(vphm);
810         if (vp->v_pages == NULL)
811             goto leave;
812
813         /*
814          * insert the markers and loop through the list of pages
815          */
816         page_vpadd(&vp->v_pages->p_vpprev->p_vpnext, mark);
817         page_vpadd(smark->p_vpnext, end);
818         for (;;) {
819
820             /*
821              * If only doing an async write back, then we can
822              * stop as soon as we get to start of the list.
823              */
824             if (flags == B_ASYNC && vp->v_pages == mark)
825                 break;
826
827             /*
828              * otherwise stop when we've gone through all the pages
829              */
830             if (mark->p_vpprev == end)
831                 break;
832
833             pp = mark->p_vpprev;
834             if (vp->v_pages == pp)
835                 where_to_move = &vp->v_pages;
836             else
837                 where_to_move = &pp->p_vpprev->p_vpnext;
838
839             ASSERT(pp->p_vnode == vp);
840
841             /*
842              * If just flushing dirty pages to disk and this vnode
843              * is using a sorted list of pages, we can stop processing
844              * as soon as we find an unmodified page. Since all the
845              * modified pages are visited first.
846              */
847             if (IS_VMODSORT(vp) &&
848                 !(flags & (B_INVAL | B_FREE | B_TRUNC))) {
849                 if (!hat_ismod(pp) && !page_io_locked(pp)) {
850 #ifdef DEBUG
851
852                 /*
853                  * For debug kernels examine what should be
854                  * all the remaining clean pages, asserting

```

```

854
855
856
857
858         * that they are not modified.
859         */
860         page_t *chk = pp;
861         int attr;
862
863         page_vpsub(&vp->v_pages, mark);
864         page_vpadd(where_to_move, mark);
865         do {
866             chk = chk->p_vpprev;
867             ASSERT(chk != end);
868             if (chk == mark)
869                 continue;
870             attr = hat_page_getattr(chk, P_MOD | P_REF);
871             if ((attr & P_MOD) == 0)
872                 continue;
873             panic("v_pages list not all clean: "
874                  "page_t=%p vnode=%p off=%lx "
875                  "attr=0x%x last clean page_t=%p\n",
876                  (void *)chk, (void *)chk->p_vnode,
877                  (long)chk->p_offset, attr,
878                  (void *)pp);
879             } while (chk != vp->v_pages);
880             break;
881         } else if (!(flags & B_ASYNC) && !hat_ismod(pp)) {
882
883             /*
884              * Couldn't get io lock, wait until IO is done.
885              * Block only for sync IO since we don't want
886              * to block async IO.
887              */
888             mutex_exit(vphm);
889             page_io_wait(pp);
890             mutex_enter(vphm);
891             continue;
892         }
893
894         /*
895          * Skip this page if the offset is out of the desired range.
896          * Just move the marker and continue.
897          */
898         if (pp->p_offset < off) {
899             page_vpsub(&vp->v_pages, mark);
900             page_vpadd(where_to_move, mark);
901             continue;
902         }
903
904         /*
905          * If we are supposed to invalidate or free this
906          * page, then we need an exclusive lock.
907          */
908         se = (flags & (B_INVAL | B_FREE)) ? SE_EXCL : SE_SHARED;
909
910         /*
911          * We must acquire the page lock for all synchronous
912          * operations (invalidate, free and write).
913          */
914         if ((flags & B_INVAL) != 0 || (flags & B_ASYNC) == 0) {
915             /*
916              * If the page_lock() drops the mutex
917              * we must retry the loop.
918              */
919             if (!page_lock(pp, se, vphm, P_NO_RECLAIM))
920                 continue;

```

```

920             /*
921              * It's ok to move the marker page now.
922              */
923             page_vpsub(&vp->v_pages, mark);
924             page_vpadd(where_to_move, mark);
925         } else {
926
927             /*
928              * update the marker page for all remaining cases
929              */
930             page_vpsub(&vp->v_pages, mark);
931             page_vpadd(where_to_move, mark);
932
933             /*
934              * For write backs, If we can't lock the page, it's
935              * invalid or in the process of being destroyed. Skip
936              * it, assuming someone else is writing it.
937              */
938             if (!page_trylock(pp, se))
939                 continue;
940
941             ASSERT(pp->p_vnode == vp);
942
943             /*
944              * Successfully locked the page, now figure out what to
945              * do with it. Free pages are easily dealt with, invalidate
946              * if desired or just go on to the next page.
947              */
948
949             if (PP_ISFREE(pp)) {
950                 if ((flags & B_INVAL) == 0) {
951                     page_unlock(pp);
952                     continue;
953
954                 /*
955                    * Invalidate (destroy) the page.
956                    */
957                     mutex_exit(vphm);
958                     page_destroy_free(pp);
959                     mutex_enter(vphm);
960                     continue;
961
962             /*
963              * pvn_getdirty() figures out what do do with a dirty page.
964              * If the page is dirty, the putapage() routine will write it
965              * and will kluster any other adjacent dirty pages it can.
966              *
967              * pvn_getdirty() and `(*putapage)' unlock the page.
968              */
969
970             mutex_exit(vphm);
971             if (pvn_getdirty(pp, flags)) {
972                 error = (*putapage)(vp, pp, NULL, NULL, flags, cred);
973                 if (!err)
974                     err = error;
975             }
976             mutex_enter(vphm);
977
978             page_vpsub(&vp->v_pages, mark);
979             page_vpsub(&vp->v_pages, end);
980
981 leave:
982             /*
983              * Release v_pages mutex, also VVLOCK and wakeup blocked thrds
984              */

```

```

986             mutex_exit(vphm);
987             kmem_cache_free(marker_cache, mark);
988             kmem_cache_free(marker_cache, end);
989             mutex_enter(&vp->v_lock);
990             vp->v_flag &= ~VVLOCK;
991             cv_broadcast(&vp->v_cv);
992             mutex_exit(&vp->v_lock);
993
994         }
995
996         /*
997          * Walk the vp->v_pages list, for every page call the callback function
998          * pointed by *page_check. If page_check returns non-zero, then mark the
999          * page as modified and if VMODSORT is set, move it to the end of v_pages
1000         * list. Moving makes sense only if we have at least two pages - this also
1001         * avoids having v_pages temporarily being NULL after calling page_vpsub()
1002         * if there was just one page.
1003         */
1004         void
1005         pvn_vplist_setdirty(vnode_t *vp, int (*page_check)(page_t *))
1006     {
1007         page_t *pp, *next, *end;
1008         kmutex_t *vphm;
1009         int shuffle;
1010
1011         vphm = page_vnode_mutex(vp);
1012         mutex_enter(vphm);
1013
1014         if (vp->v_pages == NULL) {
1015             mutex_exit(vphm);
1016             return;
1017         }
1018
1019         end = vp->v_pages->p_vpprev;
1020         shuffle = IS_VMODSORT(vp) && (vp->v_pages != end);
1021         pp = vp->v_pages;
1022
1023         for (;;) {
1024             next = pp->p_vpnext;
1025             if (pp->p_hash != PVN_VPLIST_HASH_TAG && page_check(pp)) {
1026                 /*
1027                  * hat_setmod_only() in contrast to hat_setmod() does
1028                  * not shuffle the pages and does not grab the mutex
1029                  * page_vnode_mutex. Exactly what we need.
1030                  */
1031                 hat_setmod_only(pp);
1032                 if (shuffle) {
1033                     page_vpsub(&vp->v_pages, pp);
1034                     ASSERT(vp->v_pages != NULL);
1035                     page_vpadd(&vp->v_pages->p_vpprev->p_vpnext,
1036                     pp);
1037                 }
1038             }
1039             /* Stop if we have just processed the last page. */
1040             if (pp == end)
1041                 break;
1042             pp = next;
1043         }
1044         mutex_exit(vphm);
1045
1046     }
1047
1048     /*
1049      * Zero out zbytes worth of data. Caller should be aware that this
1050      * routine may enter back into the fs layer (xxx_getpage). Locks
1051      * that the xxx_getpage routine may need should not be held while

```

```

1052 * calling this.
1053 */
1054 void
1055 pvn_vpzero(struct vnode *vp, u_offset_t vplen, size_t zbytes)
1056 {
1057     caddr_t addr;
1058
1059     ASSERT(vp->v_type != VCHR);
1060
1061     if (vp->v_pages == NULL)
1062         return;
1063
1064     /*
1065      * zbytes may be zero but there still may be some portion of
1066      * a page which needs clearing (since zbytes is a function
1067      * of filesystem block size, not pagesize.)
1068     */
1069     if (zbytes == 0 && (PAGESIZE - (vplen & PAGEOFFSET)) == 0)
1070         return;
1071
1072     /*
1073      * We get the last page and handle the partial
1074      * zeroing via kernel mappings. This will make the page
1075      * dirty so that we know that when this page is written
1076      * back, the zeroed information will go out with it. If
1077      * the page is not currently in memory, then the kzero
1078      * operation will cause it to be brought it. We use kzero
1079      * instead of bzero so that if the page cannot be read in
1080      * for any reason, the system will not panic. We need
1081      * to zero out a minimum of the fs given zbytes, but we
1082      * might also have to do more to get the entire last page.
1083     */
1084
1085     if ((zbytes + (vplen & MAXBOFFSET)) > MAXBSIZE)
1086         panic("pvn_vptrunc zbytes");
1087     addr = segmap_getmapflit(segkmap, vp, vplen,
1088         MAX(zbytes, PAGESIZE - (vplen & PAGEOFFSET)), 1, S_WRITE);
1089     (void) kzero(addr + (vplen & MAXBOFFSET),
1090         MAX(zbytes, PAGESIZE - (vplen & PAGEOFFSET)));
1091     (void) segmap_release(segkmap, addr, SM_WRITE | SM_ASYNC);
1092 }
1093 */
1094 /* Handles common work of the VOP_GETPAGE routines by iterating page by page
1095 * calling the getpage helper for each.
1096 * Handles common work of the VOP_GETPAGE routines when more than
1097 * one page must be returned by calling a file system specific operation
1098 * to do most of the work. Must be called with the vp already locked
1099 * by the VOP_GETPAGE routine.
1100 */
1101 int pvn_getpages(
1102     int (*getpage)(vnode_t *, u_offset_t, size_t, uint_t *, page_t *[],
1103                     size_t, struct seg *, caddr_t, enum seg_rw, cred_t *),
1104     struct vnode *vp,
1105     u_offset_t off,
1106     size_t len,
1107     uint_t *prot,
1108     page_t *pl[],
1109     size_t plsz,
1110     struct seg *seg,
1111     caddr_t addr,
1112     enum seg_rw rw,
1113     struct cred *cred)
1114 {
1115     page_t **ppp;

```

```

1114     u_offset_t o, eoff;
1115     size_t sz, xlen;
1116     int err;
1117
1118     /* ensure that we have enough space */
1119     ASSERT(pl == NULL || plsz >= len);
1120     ASSERT(plsz >= len); /* insure that we have enough space */
1121
1122     /*
1123      * Loop one page at a time and let getpage function fill
1124      * in the next page in array. We only allow one page to be
1125      * returned at a time (except for the last page) so that we
1126      * don't have any problems with duplicates and other such
1127      * painful problems. This is a very simple minded algorithm,
1128      * but it does the job correctly. We hope that the cost of a
1129      * getpage call for a resident page that we might have been
1130      * able to get from an earlier call doesn't cost too much.
1131     */
1132     ppp = pl;
1133     sz = (pl != NULL) ? PAGESIZE : 0;
1134     sz = PAGESIZE;
1135     eoff = off + len;
1136     xlen = len;
1137     for (o = off; o < eoff; o += PAGESIZE, addr += PAGESIZE,
1138           xlen -= PAGESIZE) {
1139         if (o + PAGESIZE >= eoff && pl != NULL) {
1140             if (o + PAGESIZE >= eoff) {
1141                 /*
1142                  * Last time through - allow the all of
1143                  * what's left of the pl[] array to be used.
1144                 */
1145                 sz = plsz - (o - off);
1146             }
1147             err = (*getpage)(vp, o, xlen, protp, ppp, sz, seg, addr,
1148                             rw, cred);
1149             if (err) {
1150                 /*
1151                  * Release any pages we already got.
1152                 */
1153                 if (o > off && pl != NULL) {
1154                     for (ppp = pl; *ppp != NULL; *ppp++ = NULL)
1155                         (void) page_release(*ppp, 1);
1156                 }
1157             }
1158         }
1159     }
1160 } unchanged_portion_omitted

```