

new/usr/src/uts/common/vm/vm\_page.c

```
*****  
187466 Sun Nov  9 09:32:15 2014  
new/usr/src/uts/common/vm/vm_page.c  
5302 vm: remove 'nopageage' static global  
*****  
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 */  
24 /*      Copyright (c) 1983, 1984, 1985, 1986, 1987, 1988, 1989 AT&T      */  
25 /*          All Rights Reserved */  
26 /*  
27 */  
28 /*  
29 * University Copyright- Copyright (c) 1982, 1986, 1988  
30 * The Regents of the University of California  
31 * All Rights Reserved  
32 *  
33 * University Acknowledgment- Portions of this document are derived from  
34 * software developed by the University of California, Berkeley, and its  
35 * contributors.  
36 */  
37 /*  
38 * VM - physical page management.  
39 */  
40 /*  
41 #include <sys/types.h>  
42 #include <sys/t_lock.h>  
43 #include <sys/param.h>  
44 #include <sys/sysm.h>  
45 #include <sys/errno.h>  
46 #include <sys/time.h>  
47 #include <sys/vnode.h>  
48 #include <sys/vm.h>  
49 #include <sys/vtrace.h>  
50 #include <sys/swap.h>  
51 #include <sys/cmn_err.h>  
52 #include <sys/tunable.h>  
53 #include <sys/sysmacros.h>  
54 #include <sys/cpuvar.h>  
55 #include <sys/callb.h>  
56 #include <sys/debug.h>  
57 #include <sys/tnf_probe.h>  
58 #include <sys/condvar_impl.h>  
59 #include <sys/mem_config.h>  
60 #include <sys/mem_cage.h>
```

1

new/usr/src/uts/common/vm/vm\_page.c

```
62 #include <sys/kmem.h>  
63 #include <sys/atomic.h>  
64 #include <sys/strlog.h>  
65 #include <sys/mman.h>  
66 #include <sys/ontrap.h>  
67 #include <sys/lgrp.h>  
68 #include <sys/vfs.h>  
69  
70 #include <vm/hat.h>  
71 #include <vm/anon.h>  
72 #include <vm/page.h>  
73 #include <vm/seg.h>  
74 #include <vm/pvn.h>  
75 #include <vm/seg_kmem.h>  
76 #include <vm/vm_dep.h>  
77 #include <sys/vm_usage.h>  
78 #include <fs/fs_subr.h>  
79 #include <sys/ddi.h>  
80 #include <sys/modctl.h>  
81  
82 static int nopageage = 0;  
83  
84 static pgcnt_t max_page_get; /* max page_get request size in pages */  
85 static pgcnt_t total_pages = 0; /* total number of pages (used by /proc) */  
86 /*  
87  * freemem_lock protects all freemem variables:  
88  * availrmem. Also this lock protects the globals which track the  
89  * availrmem changes for accurate kernel footprint calculation.  
90  * See below for an explanation of these  
91  * globals.  
92  */  
93 kmutex_t freemem_lock;  
94 pgcnt_t availrmem;  
95 pgcnt_t availrmem_initial;  
96 /*  
97  * These globals track availrmem changes to get a more accurate  
98  * estimate of the kernel size. Historically pp_kernel is used for  
99  * kernel size and is based on availrmem. But availrmem is adjusted for  
100 * locked pages in the system not just for kernel locked pages.  
101 * These new counters will track the pages locked through segvn and  
102 * by explicit user locking.  
103 *  
104 * pages_locked : How many pages are locked because of user specified  
105 * locking through mlock or plock.  
106 *  
107 * pages_useclaim,pages_claimed : These two variables track the  
108 * claim adjustments because of the protection changes on a segvn segment.  
109 *  
110 * All these globals are protected by the same lock which protects availrmem.  
111 */  
112 pgcnt_t pages_locked = 0;  
113 pgcnt_t pages_useclaim = 0;  
114 pgcnt_t pages_claimed = 0;  
115  
116 /*  
117  * new_freemem_lock protects freemem, freemem_wait & freemem_cv.  
118  */  
119 static kmutex_t new_freemem_lock;  
120 static uint_t freemem_wait; /* someone waiting for freemem */  
121 static kcondvar_t freemem_cv;  
122  
123 /*  
124  * The logical page free list is maintained as two lists, the 'free'  
125 */
```

2

new/usr/src/uts/common/vm/vm\_page.c

3

```

126 * and the 'cache' lists.
127 * The free list contains those pages that should be reused first.
128 *
129 * The implementation of the lists is machine dependent.
130 * page_get_freelist(), page_get_cachelist(),
131 * page_list_sub(), and page_list_add()
132 * form the interface to the machine dependent implementation.
133 *
134 * Pages with p_free set are on the cache list.
135 * Pages with p_free and p_page set are on the free list,
136 *
137 * A page may be locked while on either list.
138 */
139
140 /*
141 * free list accounting stuff.
142 *
143 *
144 * Spread out the value for the number of pages on the
145 * page free and page cache lists. If there is just one
146 * value, then it must be under just one lock.
147 * The lock contention and cache traffic are a real bother.
148 *
149 * When we acquire and then drop a single pcf lock
150 * we can start in the middle of the array of pcf structures.
151 * If we acquire more than one pcf lock at a time, we need to
152 * start at the front to avoid deadlocking.
153 *
154 * pcf_count holds the number of pages in each pool.
155 *
156 * pcf_block is set when page_create_get_something() has asked the
157 * PSM page freelist and page cachelist routines without specifying
158 * a color and nothing came back. This is used to block anything
159 * else from moving pages from one list to the other while the
160 * lists are searched again. If a page is freed while pcf_block is
161 * set, then pcf_reserve is incremented. pcgs_unlock() takes care
162 * of clearing pcf_block, doing the wakeups, etc.
163 */
164
165 #define MAX_PCF_FANOUT NCPU
166 static uint_t pcf_fanout = 1; /* Will get changed at boot time */
167 static uint_t pcf_fanout_mask = 0;
168
169 struct pcf {
170     kmutex_t          pcf_lock;      /* protects the structure */
171     uint_t            pcf_count;    /* page count */
172     uint_t            pcf_wait;     /* number of waiters */
173     uint_t            pcf_block;    /* pcgs flag to page_free() */
174     uint_t            pcf_reserve;  /* pages freed after pcf_block set */
175     uint_t            pcf_fill[10]; /* to line up on the caches */
176 };


---


unchanged_portion_omitted
177
178
2625 /*
2626 * Put page on the "free" list.
2627 * The free list is really two lists maintained by
2628 * the PSM of whatever machine we happen to be on.
2629 */
2630 void
2631 page_free(page_t *pp, int dontneed)
2632 {
2633     struct pcf      *p;
2634     uint_t          pcf_index;
2635
2636     ASSERT((PAGE_EXCL(pp) &
2637             !page_iolock assert(pp)) || panicstr);

```

```

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

2639     if (PP_ISFREE(pp)) {
2640         panic("page_free: page %p is free", (void *)pp);
2641     }
2642
2643     if (pp->p_szc != 0) {
2644         if (pp->p_vnode == NULL || IS_SWAPFSVP(pp->p_vnode) ||
2645             PP_ISKAS(pp)) {
2646             panic("page_free: anon or kernel "
2647                   "or no vnode large page %p", (void *)pp);
2648         }
2649         page_demote_vp_pages(pp);
2650         ASSERT(pp->p_szc == 0);
2651     }
2652
2653     /*
2654      * The page_struct_lock need not be acquired to examine these
2655      * fields since the page has an "exclusive" lock.
2656      */
2657     if (hat_page_is_mapped(pp) || pp->p_lckcnt != 0 || pp->p_cowcnt != 0 ||
2658         pp->p_slckcnt != 0) {
2659         panic("page_free pp=%p, pfns=%lx, lckcnt=%d, cowcnt=%d "
2660               "slckcnt = %d", (void *)pp, page_pptonum(pp), pp->p_lckcnt,
2661               pp->p_cowcnt, pp->p_slckcnt);
2662         /*NOTREACHED*/
2663     }
2664
2665     ASSERT(!hat_page_getshare(pp));
2666
2667     PP_SETFREE(pp);
2668     ASSERT(pp->p_vnode == NULL || !IS_VMODSORT(pp->p_vnode) ||
2669            !hat_ismod(pp));
2670     page_clr_all_props(pp);
2671     ASSERT(!hat_page_getshare(pp));
2672
2673     /*
2674      * Now we add the page to the head of the free list.
2675      * But if this page is associated with a paged vnode
2676      * then we adjust the head forward so that the page is
2677      * effectively at the end of the list.
2678      */
2679     if (pp->p_vnode == NULL) {
2680         /*
2681          * Page has no identity, put it on the free list.
2682          */
2683         PP_SETAGED(pp);
2684         pp->p_offset = (u_offset_t)-1;
2685         page_list_add(pp, PG_FREE_LIST | PG_LIST_TAIL);
2686         VM_STAT_ADD(pagecnt.pc_free_free);
2687         TRACE_1(TR_FAC_VM, TR_PAGE_FREE_FREE,
2688                 "page_free_free:pp %p", pp);
2689     } else {
2690         PP_CLRAGED(pp);
2691
2692         if (!dontneed) {
2693             if (!dontneed || nophage) {
2694                 /* move it to the tail of the list */
2695                 page_list_add(pp, PG_CACHE_LIST | PG_LIST_TAIL);
2696
2697                 VM_STAT_ADD(pagecnt.pc_free_cache);
2698                 TRACE_1(TR_FAC_VM, TR_PAGE_FREE_CACHE_TAIL,
2699                         "page_free_cache_tail:pp %p", pp);
2700             } else {
2701                 page_list_add(pp, PG_CACHE_LIST | PG_LIST_HEAD);
2702
2703                 VM_STAT_ADD(pagecnt.pc_free_dontneed);
2704             }
2705         }
2706     }
2707 }

```

```
2703
2704         }
2705     }
2706     page_unlock(pp);
2707
2708     /*
2709      * Now do the 'freemem' accounting.
2710      */
2711     pcf_index = PCF_INDEX();
2712     p = &pcf[pcf_index];
2713
2714     mutex_enter(&p->pcf_lock);
2715     if (p->pcf_block) {
2716         p->pcf_reserve += 1;
2717     } else {
2718         p->pcf_count += 1;
2719         if (p->pcf_wait) {
2720             mutex_enter(&new_freemem_lock);
2721             /*
2722              * Check to see if some other thread
2723              * is actually waiting. Another bucket
2724              * may have woken it up by now. If there
2725              * are no waiters, then set our pcf_wait
2726              * count to zero to avoid coming in here
2727              * next time. Also, since only one page
2728              * was put on the free list, just wake
2729              * up one waiter.
2730             */
2731             if (freemem_wait) {
2732                 cv_signal(&freemem_cv);
2733                 p->pcf_wait--;
2734             } else {
2735                 p->pcf_wait = 0;
2736             }
2737             mutex_exit(&new_freemem_lock);
2738         }
2739     }
2740     mutex_exit(&p->pcf_lock);
2741
2742     /* freemem is approximate, so this test OK */
2743     if (!p->pcf_block)
2744         freemem += 1;
2745
2746 }
```

unchanged\_portion\_omitted