

```

*****
49002 Fri Feb 27 16:32:49 2015
new/usr/src/uts/common/vm/page.h
vm: be careful about empty statements
*****
_____unchanged_portion_omitted_____

69 /*
70 * Callers of page_try_reclaim_lock and page_lock_es can use this flag
71 * to get SE_EXCL access before reader/writers are given access.
72 */
73 #define SE_EXCL_WANTED 0x02

75 /*
76 * All page_*lock() requests will be denied unless this flag is set in
77 * the 'es' parameter.
78 */
79 #define SE_RETIRED 0x04

81 #endif /* _KERNEL | _KMEMUSER */

83 typedef int selock_t;

85 /*
86 * Define VM_STATS to turn on all sorts of statistic gathering about
87 * the VM layer. By default, it is only turned on when DEBUG is
88 * also defined.
89 */
90 #ifdef DEBUG
91 #define VM_STATS
92 #endif /* DEBUG */

94 #ifdef VM_STATS
95 #define VM_STAT_ADD(stat) (stat)++
96 #define VM_STAT_COND_ADD(cond, stat) ((void) (!(cond) || (stat)++))
97 #else
98 #define VM_STAT_ADD(stat) do { } while (0)
99 #define VM_STAT_COND_ADD(cond, stat) do { } while (0)
100 #endif /* VM_STATS */

102 #ifdef _KERNEL

104 /*
105 * PAGE_LLOCK_SIZE is 2 * NCPU, but no smaller than 128.
106 * PAGE_LLOCK_SHIFT is log2(PAGE_LLOCK_SIZE).
107 *
108 * We use ? : instead of #if because <vm/page.h> is included everywhere;
109 * NCPU_P2 is only a constant in the "unix" module.
110 *
111 */
112 #define PAGE_LLOCK_SHIFT \
113 ((unsigned)((2*NCPU_P2) > 128) ? NCPU_LOG2 + 1 : 7))

115 #define PAGE_LLOCK_SIZE (1ul << PAGE_LLOCK_SHIFT)

117 /*
118 * The number of low order 0 (or less variable) bits in the page_t address.
119 */
120 #if defined(__sparc)
121 #define PP_SHIFT 7
122 #else
123 #define PP_SHIFT 6
124 #endif

```

```

126 /*
127 * pp may be the root of a large page, and many low order bits will be 0.
128 * Shift and XOR multiple times to capture the good bits across the range of
129 * possible page sizes.
130 */
131 #define PAGE_LLOCK_HASH(pp) \
132 (((((uintptr_t)(pp) >> PP_SHIFT) ^ \
133 ((uintptr_t)(pp) >> (PAGE_LLOCK_SHIFT + PP_SHIFT))) ^ \
134 ((uintptr_t)(pp) >> ((PAGE_LLOCK_SHIFT * 2) + PP_SHIFT))) ^ \
135 ((uintptr_t)(pp) >> ((PAGE_LLOCK_SHIFT * 3) + PP_SHIFT))) & \
136 (PAGE_LLOCK_SIZE - 1))

138 #define page_struct_lock(pp) \
139 mutex_enter(&page_llocks[PAGE_LLOCK_HASH(PP_PAGEROOT(pp))].pad_mutex)
140 #define page_struct_unlock(pp) \
141 mutex_exit(&page_llocks[PAGE_LLOCK_HASH(PP_PAGEROOT(pp))].pad_mutex)

143 #endif /* _KERNEL */

145 #include <sys/t_lock.h>

147 struct as;

149 /*
150 * Each physical page has a page structure, which is used to maintain
151 * these pages as a cache. A page can be found via a hashed lookup
152 * based on the [vp, offset]. If a page has an [vp, offset] identity,
153 * then it is entered on a doubly linked circular list off the
154 * vnode using the vnext/vprev pointers. If the p_free bit
155 * is on, then the page is also on a doubly linked circular free
156 * list using next/prev pointers. If the "p_selock" and "p_iolock"
157 * are held, then the page is currently being read in (exclusive p_selock)
158 * or written back (shared p_selock). In this case, the next/prev pointers
159 * are used to link the pages together for a consecutive i/o request. If
160 * the page is being brought in from its backing store, then other processes
161 * will wait for the i/o to complete before attaching to the page since it
162 * will have an "exclusive" lock.
163 *
164 * Each page structure has the locks described below along with
165 * the fields they protect:
166 *
167 * p_selock This is a per-page shared/exclusive lock that is
168 * used to implement the logical shared/exclusive
169 * lock for each page. The "shared" lock is normally
170 * used in most cases while the "exclusive" lock is
171 * required to destroy or retain exclusive access to
172 * a page (e.g., while reading in pages). The appropriate
173 * lock is always held whenever there is any reference
174 * to a page structure (e.g., during i/o).
175 * (Note that with the addition of the "writer-lock-wanted"
176 * semantics (via SE_EWANTED), threads must not acquire
177 * multiple reader locks or else a deadly embrace will
178 * occur in the following situation: thread 1 obtains a
179 * reader lock; next thread 2 fails to get a writer lock
180 * but specified SE_EWANTED so it will wait by either
181 * blocking (when using page_lock_es) or spinning while
182 * retrying (when using page_try_reclaim_lock) until the
183 * reader lock is released; then thread 1 attempts to
184 * get another reader lock but is denied due to
185 * SE_EWANTED being set, and now both threads are in a
186 * deadly embrace.)
187 *
188 * p_hash
189 * p_vnode
190 * p_offset
191 *

```

```

192 *           p_free
193 *           p_age
194 *
195 *   p_iolock   This is a binary semaphore lock that provides
196 *             exclusive access to the i/o list links in each
197 *             page structure. It is always held while the page
198 *             is on an i/o list (i.e., involved in i/o). That is,
199 *             even though a page may be only 'shared' locked
200 *             while it is doing a write, the following fields may
201 *             change anyway. Normally, the page must be
202 *             'exclusively' locked to change anything in it.
203 *
204 *           p_next
205 *           p_prev
206 *
207 * The following fields are protected by the global page_llocks[]:
208 *
209 *           p_lckcnt
210 *           p_cowcnt
211 *
212 * The following lists are protected by the global page_freelock:
213 *
214 *           page_cachelist
215 *           page_freelist
216 *
217 * The following, for our purposes, are protected by
218 * the global freemem_lock:
219 *
220 *           freemem
221 *           freemem_wait
222 *           freemem_cv
223 *
224 * The following fields are protected by hat layer lock(s). When a page
225 * structure is not mapped and is not associated with a vnode (after a call
226 * to page_hashout() for example) the p_nrm field may be modified with out
227 * holding the hat layer lock:
228 *
229 *           p_nrm
230 *           p_mapping
231 *           p_share
232 *
233 * The following field is file system dependent. How it is used and
234 * the locking strategies applied are up to the individual file system
235 * implementation.
236 *
237 *           p_fsdata
238 *
239 * The page structure is used to represent and control the system's
240 * physical pages. There is one instance of the structure for each
241 * page that is not permanently allocated. For example, the pages that
242 * hold the page structures are permanently held by the kernel
243 * and hence do not need page structures to track them. The array
244 * of page structures is allocated early on in the kernel's life and
245 * is based on the amount of available physical memory.
246 *
247 * Each page structure may simultaneously appear on several linked lists.
248 * The lists are: hash list, free or in i/o list, and a vnode's page list.
249 * Each type of list is protected by a different group of mutexes as described
250 * below:
251 *
252 * The hash list is used to quickly find a page when the page's vnode and
253 * offset within the vnode are known. Each page that is hashed is
254 * connected via the 'p_hash' field. The anchor for each hash is in the
255 * array 'page_hash'. An array of mutexes, 'ph_mutex', protects the
256 * lists anchored by page_hash[]. To either search or modify a given hash
257 * list, the appropriate mutex in the ph_mutex array must be held.

```

```

258 *
259 * The free list contains pages that are 'free to be given away'. For
260 * efficiency reasons, pages on this list are placed in two categories:
261 * pages that are still associated with a vnode, and pages that are not
262 * associated with a vnode. Free pages always have their 'p_free' bit set,
263 * free pages that are still associated with a vnode also have their
264 * 'p_age' bit set. Pages on the free list are connected via their
265 * 'p_next' and 'p_prev' fields. When a page is involved in some sort
266 * of i/o, it is not free and these fields may be used to link associated
267 * pages together. At the moment, the free list is protected by a
268 * single mutex 'page_freelock'. The list of free pages still associated
269 * with a vnode is anchored by 'page_cachelist' while other free pages
270 * are anchored in architecture dependent ways (to handle page coloring etc.).
271 *
272 * Pages associated with a given vnode appear on a list anchored in the
273 * vnode by the 'v_pages' field. They are linked together with
274 * 'p_vpnext' and 'p_vpprev'. The field 'p_offset' contains a page's
275 * offset within the vnode. The pages on this list are not kept in
276 * offset order. These lists, in a manner similar to the hash lists,
277 * are protected by an array of mutexes called 'vph_hash'. Before
278 * searching or modifying this chain the appropriate mutex in the
279 * vph_hash[] array must be held.
280 *
281 * Again, each of the lists that a page can appear on is protected by a
282 * mutex. Before reading or writing any of the fields comprising the
283 * list, the appropriate lock must be held. These list locks should only
284 * be held for very short intervals.
285 *
286 * In addition to the list locks, each page structure contains a
287 * shared/exclusive lock that protects various fields within it.
288 * To modify one of these fields, the 'p_selock' must be exclusively held.
289 * To read a field with a degree of certainty, the lock must be at least
290 * held shared.
291 *
292 * Removing a page structure from one of the lists requires holding
293 * the appropriate list lock and the page's p_selock. A page may be
294 * prevented from changing identity, being freed, or otherwise modified
295 * by acquiring p_selock shared.
296 *
297 * To avoid deadlocks, a strict locking protocol must be followed. Basically
298 * there are two cases: In the first case, the page structure in question
299 * is known ahead of time (e.g., when the page is to be added or removed
300 * from a list). In the second case, the page structure is not known and
301 * must be found by searching one of the lists.
302 *
303 * When adding or removing a known page to one of the lists, first the
304 * page must be exclusively locked (since at least one of its fields
305 * will be modified), second the lock protecting the list must be acquired,
306 * third the page inserted or deleted, and finally the list lock dropped.
307 *
308 * The more interesting case occurs when the particular page structure
309 * is not known ahead of time. For example, when a call is made to
310 * page_lookup(), it is not known if a page with the desired (vnode and
311 * offset pair) identity exists. So the appropriate mutex in ph_mutex is
312 * acquired, the hash list searched, and if the desired page is found
313 * an attempt is made to lock it. The attempt to acquire p_selock must
314 * not block while the hash list lock is held. A deadlock could occur
315 * if some other process was trying to remove the page from the list.
316 * The removing process (following the above protocol) would have exclusively
317 * locked the page, and be spinning waiting to acquire the lock protecting
318 * the hash list. Since the searching process holds the hash list lock
319 * and is waiting to acquire the page lock, a deadlock occurs.
320 *
321 * The proper scheme to follow is: first, lock the appropriate list,
322 * search the list, and if the desired page is found either use
323 * page_trylock() (which will not block) or pass the address of the

```

```

324 * list lock to page_lock(). If page_lock() can not acquire the page's
325 * lock, it will drop the list lock before going to sleep. page_lock()
326 * returns a value to indicate if the list lock was dropped allowing the
327 * calling program to react appropriately (i.e., retry the operation).
328 *
329 * If the list lock was dropped before the attempt at locking the page
330 * was made, checks would have to be made to ensure that the page had
331 * not changed identity before its lock was obtained. This is because
332 * the interval between dropping the list lock and acquiring the page
333 * lock is indeterminate.
334 *
335 * In addition, when both a hash list lock (ph_mutex[]) and a vnode list
336 * lock (vph_mutex[]) are needed, the hash list lock must be acquired first.
337 * The routine page_hashin() is a good example of this sequence.
338 * This sequence is ASSERTed by checking that the vph_mutex[] is not held
339 * just before each acquisition of one of the mutexes in ph_mutex[].
340 *
341 * So, as a quick summary:
342 *
343 *     pse_mutex[]'s protect the p_selock and p_cv fields.
344 *
345 *     p_selock protects the p_free, p_age, p_vnode, p_offset and p_hash,
346 *
347 *     ph_mutex[]'s protect the page_hash[] array and its chains.
348 *
349 *     vph_mutex[]'s protect the v_pages field and the vp page chains.
350 *
351 * First lock the page, then the hash chain, then the vnode chain. When
352 * this is not possible 'trylocks' must be used. Sleeping while holding
353 * any of these mutexes (p_selock is not a mutex) is not allowed.
354 *
355 *
356 *     field          reading          writing          ordering
357 *     =====
358 *     p_vnode        p_selock(E,S)    p_selock(E)
359 *     p_offset
360 *     p_free
361 *     p_age
362 *     =====
363 *     p_hash         p_selock(E,S)    p_selock(E) &&    p_selock, ph_mutex
364 *                   p_selock(E,S)    ph_mutex[]
365 *     =====
366 *     p_vpnext       p_selock(E,S)    p_selock(E) &&    p_selock, vph_mutex
367 *     p_vpprev
368 *                   vph_mutex[]
369 *     =====
370 *     When the p_free bit is set:
371 *
372 *     p_next         p_selock(E,S)    p_selock(E) &&    p_selock,
373 *     p_prev         page_freelock    page_freelock
374 *
375 *     When the p_free bit is not set:
376 *
377 *     p_next         p_selock(E,S)    p_selock(E) &&    p_selock, p_iolock
378 *     p_prev         p_iolock
379 *     =====
380 *     p_selock       pse_mutex[]      pse_mutex[]      can't acquire any
381 *     p_cv            other mutexes or  sleep while holding
382 *                   this lock.
383 *     =====
384 *     p_lckcnt       p_selock(E,S)    p_selock(E)
385 *                   OR
386 *                   p_selock(S) &&
387 *                   page_llocks[]
388 *     p_cowcnt
389 *     =====

```

```

390 *     p_nrm          hat layer lock  hat layer lock
391 *     p_mapping
392 *     p_pagenum
393 *     =====
394 *
395 *     where:
396 *     E----> exclusive version of p_selock.
397 *     S----> shared version of p_selock.
398 *
399 *
400 *     Global data structures and variable:
401 *
402 *     field          reading          writing          ordering
403 *     =====
404 *     page_hash[]    ph_mutex[]      ph_mutex[]      can hold this lock
405 *                   before acquiring
406 *                   a vph_mutex or
407 *                   pse_mutex.
408 *     =====
409 *     vp->v_pages    vph_mutex[]      vph_mutex[]      can only acquire
410 *                   a pse_mutex while
411 *                   holding this lock.
412 *     =====
413 *     page_cachelist page_freelock    page_freelock    can't acquire any
414 *     page_freelist  page_freelock    page_freelock
415 *     =====
416 *     freemem        freemem_lock    freemem_lock     can't acquire any
417 *     freemem_wait   other mutexes while
418 *     freemem_cv     holding this mutex.
419 *     =====
420 *
421 *     Page relocation, PG_NORELOC and P_NORELOC.
422 *
423 *     Pages may be relocated using the page_relocate() interface. Relocation
424 *     involves moving the contents and identity of a page to another, free page.
425 *     To relocate a page, the SE_EXCL lock must be obtained. The way to prevent
426 *     a page from being relocated is to hold the SE_SHARED lock (the SE_EXCL
427 *     lock must not be held indefinitely). If the page is going to be held
428 *     SE_SHARED indefinitely, then the PG_NORELOC hint should be passed
429 *     to page_create_va so that pages that are prevented from being relocated
430 *     can be managed differently by the platform specific layer.
431 *
432 *     Pages locked in memory using page_pp_lock (p_lckcnt/p_cowcnt != 0)
433 *     are guaranteed to be held in memory, but can still be relocated
434 *     providing the SE_EXCL lock can be obtained.
435 *
436 *     The P_NORELOC bit in the page_t.p_state field is provided for use by
437 *     the platform specific code in managing pages when the PG_NORELOC
438 *     hint is used.
439 *
440 *     Memory delete and page locking.
441 *
442 *     The set of all usable pages is managed using the global page list as
443 *     implemented by the memseg structure defined below. When memory is added
444 *     or deleted this list changes. Additions to this list guarantee that the
445 *     list is never corrupt. In order to avoid the necessity of an additional
446 *     lock to protect against failed accesses to the memseg being deleted and,
447 *     more importantly, the page_ts, the memseg structure is never freed and the
448 *     page_t virtual address space is remapped to a page (or pages) of
449 *     zeros. If a page_t is manipulated while it is p_selock'd, or if it is
450 *     locked indirectly via a hash or freelist lock, it is not possible for
451 *     memory delete to collect the page and so that part of the page list is
452 *     prevented from being deleted. If the page is referenced outside of one
453 *     of these locks, it is possible for the page_t being referenced to be
454 *     deleted. Examples of this are page_t pointers returned by
455 *     page_numtopp_nolock, page_first and page_next. Providing the page_t

```

```

456 * is re-checked after taking the p_slock (for p_vnode != NULL), the
457 * remapping to the zero pages will be detected.
458 *
459 *
460 * Page size (p_szc field) and page locking.
461 *
462 * p_szc field of free pages is changed by free list manager under freelist
463 * locks and is of no concern to the rest of VM subsystem.
464 *
465 * p_szc changes of allocated anonymous (swapfs) can only be done only after
466 * exclusively locking all constituent pages and calling hat_pageunload() on
467 * each of them. To prevent p_szc changes of non free anonymous (swapfs) large
468 * pages it's enough to either lock SHARED any of constituent pages or prevent
469 * hat_pageunload() by holding hat level lock that protects mapping lists (this
470 * method is for hat code only)
471 *
472 * To increase (promote) p_szc of allocated non anonymous file system pages
473 * one has to first lock exclusively all involved constituent pages and call
474 * hat_pageunload() on each of them. To prevent p_szc promote it's enough to
475 * either lock SHARED any of constituent pages that will be needed to make a
476 * large page or prevent hat_pageunload() by holding hat level lock that
477 * protects mapping lists (this method is for hat code only).
478 *
479 * To decrease (demote) p_szc of an allocated non anonymous file system large
480 * page one can either use the same method as used for changing p_szc of
481 * anonymous large pages or if it's not possible to lock all constituent pages
482 * exclusively a different method can be used. In the second method one only
483 * has to exclusively lock one of constituent pages but then one has to
484 * acquire further locks by calling page_szc_lock() and
485 * hat_page_demote(). hat_page_demote() acquires hat level locks and then
486 * demotes the page. This mechanism relies on the fact that any code that
487 * needs to prevent p_szc of a file system large page from changing either
488 * locks all constituent large pages at least SHARED or locks some pages at
489 * least SHARED and calls page_szc_lock() or uses hat level page locks.
490 * Demotion using this method is implemented by page_demote_vp_pages().
491 * Please see comments in front of page_demote_vp_pages(), hat_page_demote()
492 * and page_szc_lock() for more details.
493 *
494 * Lock order: p_slock, page_szc_lock, ph_mutex/vph_mutex/freelist,
495 * hat level locks.
496 */

498 typedef struct page {
499     u_offset_t    p_offset;    /* offset into vnode for this page */
500     struct vnode  *p_vnode;    /* vnode that this page is named by */
501     selock_t      p_slock;    /* shared/exclusive lock on the page */
502 #if defined(_LP64)
503     uint_t        p_vpmlref;   /* vpm ref - index of the vmap_t */
504 #endif
505     struct page   *p_hash;     /* hash by [vnode, offset] */
506     struct page   *p_vpnext;   /* next page in vnode list */
507     struct page   *p_vpprev;   /* prev page in vnode list */
508     struct page   *p_next;     /* next page in free/intrans lists */
509     struct page   *p_prev;     /* prev page in free/intrans lists */
510     ushort_t      p_lckcnt;    /* number of locks on page data */
511     ushort_t      p_cowcnt;    /* number of copy on write lock */
512     kcondvar_t    p_cv;       /* page struct's condition var */
513     kcondvar_t    p_io_cv;    /* for iolock */
514     uchar_t       p_iolock_state; /* replaces p_iolock */
515     volatile uchar_t p_szc;    /* page size code */
516     uchar_t       p_fsdata;   /* file system dependent byte */
517     uchar_t       p_state;    /* p_free, p_noreloc */
518     uchar_t       p_nrm;     /* non-cache, ref, mod readonly bits */
519 #if defined(__sparc)
520     uchar_t       p_vcolor;   /* virtual color */
521 #else

```

```

522     uchar_t       p_embed;    /* x86 - changes p_mapping & p_index */
523 #endif
524     uchar_t       p_index;    /* MPSS mapping info. Not used on x86 */
525     uchar_t       p_toxic;    /* page has an unrecoverable error */
526     void          *p_mapping; /* hat specific translation info */
527     pfn_t         p_pagenum;  /* physical page number */

529     uint_t        p_share;    /* number of translations */
530 #if defined(_LP64)
531     uint_t        p_sharepad; /* pad for growing p_share */
532 #endif
533     uint_t        p_slckcnt;  /* number of softlocks */
534 #if defined(__sparc)
535     uint_t        p_kpmlref;  /* number of kpm mapping sharers */
536     struct kpme   *p_kpmlist; /* kpm specific mapping info */
537 #else
538     /* index of entry in p_map when p_embed is set */
539     uint_t        p_mlentry;
540 #endif
541 #if defined(_LP64)
542     kmutex_t      p_ilock;    /* protects p_vpmlref */
543 #else
544     uint64_t      p_msresv_2; /* page allocation debugging */
545 #endif
546 } page_t;

```

unchanged portion omitted