🗄 **torvalds** / **linux**

<> Code          ⑂ Pull requests **321**          ▶ Actions          ▦ Projects          ⊘ Security          ⬋ Insights

🏷 **v4.4-rc1** ▾                                                                    ···

**linux** / **mm** / **madvise.c**

| | | |
|---|---|---|
| ⬤ | **Mike Kravetz** mm: madvise allow remove operation for hugetlbfs  ... | 🕘 **History** |

👥 **13 contributors**   ⬤ ⬤ ⬤ ⬤ ⬤ ⬤ ⬤ ⬤ ⬤ ⬤ ⬤ ⬤   **+1**

| Raw | Blame | Copy File | 30.46 KB ⬇ | | 🖥 ✏ 🗑 |
|---|---|---|---|---|---|

551 lines (491 sloc)   |   13.9 KB

```
 1   /*
 2    *      linux/mm/madvise.c
 3    *
 4    * Copyright (C) 1999  Linus Torvalds
 5    * Copyright (C) 2002  Christoph Hellwig
 6    */
 7
 8   #include <linux/mman.h>
 9   #include <linux/pagemap.h>
10   #include <linux/syscalls.h>
11   #include <linux/mempolicy.h>
12   #include <linux/page-isolation.h>
13   #include <linux/hugetlb.h>
14   #include <linux/falloc.h>
15   #include <linux/sched.h>
16   #include <linux/ksm.h>
17   #include <linux/fs.h>
18   #include <linux/file.h>
19   #include <linux/blkdev.h>
20   #include <linux/backing-dev.h>
21   #include <linux/swap.h>
22   #include <linux/swapops.h>
23
24   /*
25    * Any behaviour which results in changes to the vma->vm_flags needs to
26    * take mmap_sem for writing. Others, which simply traverse vmas, need
27    * to only take it for reading.
```

```
28     */
29    static int madvise_need_mmap_write(int behavior)
30    {
31            switch (behavior) {
32            case MADV_REMOVE:
33            case MADV_WILLNEED:
34            case MADV_DONTNEED:
35                    return 0;
36            default:
37                    /* be safe, default to 1. list exceptions explicitly */
38                    return 1;
39            }
40    }
41
42    /*
43     * We can potentially split a vm area into separate
44     * areas, each area with its own behavior.
45     */
46    static long madvise_behavior(struct vm_area_struct *vma,
47                         struct vm_area_struct **prev,
48                         unsigned long start, unsigned long end, int behavior)
49    {
50            struct mm_struct *mm = vma->vm_mm;
51            int error = 0;
52            pgoff_t pgoff;
53            unsigned long new_flags = vma->vm_flags;
54
55            switch (behavior) {
56            case MADV_NORMAL:
57                    new_flags = new_flags & ~VM_RAND_READ & ~VM_SEQ_READ;
58                    break;
59            case MADV_SEQUENTIAL:
60                    new_flags = (new_flags & ~VM_RAND_READ) | VM_SEQ_READ;
61                    break;
62            case MADV_RANDOM:
63                    new_flags = (new_flags & ~VM_SEQ_READ) | VM_RAND_READ;
64                    break;
65            case MADV_DONTFORK:
66                    new_flags |= VM_DONTCOPY;
67                    break;
68            case MADV_DOFORK:
69                    if (vma->vm_flags & VM_IO) {
70                            error = -EINVAL;
71                            goto out;
72                    }
73                    new_flags &= ~VM_DONTCOPY;
74                    break;
75            case MADV_DONTDUMP:
```

```
 76                        new_flags |= VM_DONTDUMP;
 77                        break;
 78                case MADV_DODUMP:
 79                        if (new_flags & VM_SPECIAL) {
 80                                error = -EINVAL;
 81                                goto out;
 82                        }
 83                        new_flags &= ~VM_DONTDUMP;
 84                        break;
 85                case MADV_MERGEABLE:
 86                case MADV_UNMERGEABLE:
 87                        error = ksm_madvise(vma, start, end, behavior, &new_flags);
 88                        if (error)
 89                                goto out;
 90                        break;
 91                case MADV_HUGEPAGE:
 92                case MADV_NOHUGEPAGE:
 93                        error = hugepage_madvise(vma, &new_flags, behavior);
 94                        if (error)
 95                                goto out;
 96                        break;
 97        }
 98
 99        if (new_flags == vma->vm_flags) {
100                *prev = vma;
101                goto out;
102        }
103
104        pgoff = vma->vm_pgoff + ((start - vma->vm_start) >> PAGE_SHIFT);
105        *prev = vma_merge(mm, *prev, start, end, new_flags, vma->anon_vma,
106                        vma->vm_file, pgoff, vma_policy(vma),
107                        vma->vm_userfaultfd_ctx);
108        if (*prev) {
109                vma = *prev;
110                goto success;
111        }
112
113        *prev = vma;
114
115        if (start != vma->vm_start) {
116                error = split_vma(mm, vma, start, 1);
117                if (error)
118                        goto out;
119        }
120
121        if (end != vma->vm_end) {
122                error = split_vma(mm, vma, end, 0);
123                if (error)
```

```
124                           goto out;
125               }
126
127   success:
128           /*
129            * vm_flags is protected by the mmap_sem held in write mode.
130            */
131           vma->vm_flags = new_flags;
132
133   out:
134           if (error == -ENOMEM)
135                   error = -EAGAIN;
136           return error;
137   }
138
139   #ifdef CONFIG_SWAP
140   static int swapin_walk_pmd_entry(pmd_t *pmd, unsigned long start,
141           unsigned long end, struct mm_walk *walk)
142   {
143           pte_t *orig_pte;
144           struct vm_area_struct *vma = walk->private;
145           unsigned long index;
146
147           if (pmd_none_or_trans_huge_or_clear_bad(pmd))
148                   return 0;
149
150           for (index = start; index != end; index += PAGE_SIZE) {
151                   pte_t pte;
152                   swp_entry_t entry;
153                   struct page *page;
154                   spinlock_t *ptl;
155
156                   orig_pte = pte_offset_map_lock(vma->vm_mm, pmd, start, &ptl);
157                   pte = *(orig_pte + ((index - start) / PAGE_SIZE));
158                   pte_unmap_unlock(orig_pte, ptl);
159
160                   if (pte_present(pte) || pte_none(pte))
161                           continue;
162                   entry = pte_to_swp_entry(pte);
163                   if (unlikely(non_swap_entry(entry)))
164                           continue;
165
166                   page = read_swap_cache_async(entry, GFP_HIGHUSER_MOVABLE,
167                                                           vma, index);
168                   if (page)
169                           page_cache_release(page);
170           }
171
```

```
172             return 0;
173     }
174
175     static void force_swapin_readahead(struct vm_area_struct *vma,
176                     unsigned long start, unsigned long end)
177     {
178             struct mm_walk walk = {
179                     .mm = vma->vm_mm,
180                     .pmd_entry = swapin_walk_pmd_entry,
181                     .private = vma,
182             };
183
184             walk_page_range(start, end, &walk);
185
186             lru_add_drain();        /* Push any new pages onto the LRU now */
187     }
188
189     static void force_shm_swapin_readahead(struct vm_area_struct *vma,
190                     unsigned long start, unsigned long end,
191                     struct address_space *mapping)
192     {
193             pgoff_t index;
194             struct page *page;
195             swp_entry_t swap;
196
197             for (; start < end; start += PAGE_SIZE) {
198                     index = ((start - vma->vm_start) >> PAGE_SHIFT) + vma->vm_pgoff;
199
200                     page = find_get_entry(mapping, index);
201                     if (!radix_tree_exceptional_entry(page)) {
202                             if (page)
203                                     page_cache_release(page);
204                             continue;
205                     }
206                     swap = radix_to_swp_entry(page);
207                     page = read_swap_cache_async(swap, GFP_HIGHUSER_MOVABLE,
208                                                         NULL, 0);
209                     if (page)
210                             page_cache_release(page);
211             }
212
213             lru_add_drain();        /* Push any new pages onto the LRU now */
214     }
215     #endif           /* CONFIG_SWAP */
216
217     /*
218      * Schedule all required I/O operations.  Do not wait for completion.
219      */
```

```
220    static long madvise_willneed(struct vm_area_struct *vma,
221                                 struct vm_area_struct **prev,
222                                 unsigned long start, unsigned long end)
223    {
224            struct file *file = vma->vm_file;
225
226    #ifdef CONFIG_SWAP
227            if (!file) {
228                    *prev = vma;
229                    force_swapin_readahead(vma, start, end);
230                    return 0;
231            }
232
233            if (shmem_mapping(file->f_mapping)) {
234                    *prev = vma;
235                    force_shm_swapin_readahead(vma, start, end,
236                                            file->f_mapping);
237                    return 0;
238            }
239    #else
240            if (!file)
241                    return -EBADF;
242    #endif
243
244            if (IS_DAX(file_inode(file))) {
245                    /* no bad return value, but ignore advice */
246                    return 0;
247            }
248
249            *prev = vma;
250            start = ((start - vma->vm_start) >> PAGE_SHIFT) + vma->vm_pgoff;
251            if (end > vma->vm_end)
252                    end = vma->vm_end;
253            end = ((end - vma->vm_start) >> PAGE_SHIFT) + vma->vm_pgoff;
254
255            force_page_cache_readahead(file->f_mapping, file, start, end - start);
256            return 0;
257    }
258
259    /*
260     * Application no longer needs these pages.  If the pages are dirty,
261     * it's OK to just throw them away.  The app will be more careful about
262     * data it wants to keep.  Be sure to free swap resources too.  The
263     * zap_page_range call sets things up for shrink_active_list to actually free
264     * these pages later if no one else has touched them in the meantime,
265     * although we could add these pages to a global reuse list for
266     * shrink_active_list to pick up before reclaiming other pages.
267     *
```

```
268    * NB: This interface discards data rather than pushes it out to swap,
269    * as some implementations do.  This has performance implications for
270    * applications like large transactional databases which want to discard
271    * pages in anonymous maps after committing to backing store the data
272    * that was kept in them.  There is no reason to write this data out to
273    * the swap area if the application is discarding it.
274    *
275    * An interface that causes the system to free clean pages and flush
276    * dirty pages is already available as msync(MS_INVALIDATE).
277    */
278   static long madvise_dontneed(struct vm_area_struct *vma,
279                           struct vm_area_struct **prev,
280                           unsigned long start, unsigned long end)
281   {
282         *prev = vma;
283         if (vma->vm_flags & (VM_LOCKED|VM_HUGETLB|VM_PFNMAP))
284               return -EINVAL;
285
286         zap_page_range(vma, start, end - start, NULL);
287         return 0;
288   }
289
290   /*
291    * Application wants to free up the pages and associated backing store.
292    * This is effectively punching a hole into the middle of a file.
293    */
294   static long madvise_remove(struct vm_area_struct *vma,
295                           struct vm_area_struct **prev,
296                           unsigned long start, unsigned long end)
297   {
298         loff_t offset;
299         int error;
300         struct file *f;
301
302         *prev = NULL;   /* tell sys_madvise we drop mmap_sem */
303
304         if (vma->vm_flags & VM_LOCKED)
305               return -EINVAL;
306
307         f = vma->vm_file;
308
309         if (!f || !f->f_mapping || !f->f_mapping->host) {
310                 return -EINVAL;
311         }
312
313         if ((vma->vm_flags & (VM_SHARED|VM_WRITE)) != (VM_SHARED|VM_WRITE))
314               return -EACCES;
315
```

```
316                 offset = (loff_t)(start - vma->vm_start)
317                          + ((loff_t)vma->vm_pgoff << PAGE_SHIFT);
318
319           /*
320            * Filesystem's fallocate may need to take i_mutex.  We need to
321            * explicitly grab a reference because the vma (and hence the
322            * vma's reference to the file) can go away as soon as we drop
323            * mmap_sem.
324            */
325           get_file(f);
326           up_read(&current->mm->mmap_sem);
327           error = vfs_fallocate(f,
328                               FALLOC_FL_PUNCH_HOLE | FALLOC_FL_KEEP_SIZE,
329                               offset, end - start);
330           fput(f);
331           down_read(&current->mm->mmap_sem);
332           return error;
333 }
334
335 #ifdef CONFIG_MEMORY_FAILURE
336 /*
337  * Error injection support for memory error handling.
338  */
339 static int madvise_hwpoison(int bhv, unsigned long start, unsigned long end)
340 {
341           struct page *p;
342           if (!capable(CAP_SYS_ADMIN))
343                   return -EPERM;
344           for (; start < end; start += PAGE_SIZE <<
345                               compound_order(compound_head(p))) {
346                 int ret;
347
348                 ret = get_user_pages_fast(start, 1, 0, &p);
349                 if (ret != 1)
350                         return ret;
351
352                 if (PageHWPoison(p)) {
353                         put_page(p);
354                         continue;
355                 }
356                 if (bhv == MADV_SOFT_OFFLINE) {
357                         pr_info("Soft offlining page %#lx at %#lx\n",
358                                 page_to_pfn(p), start);
359                         ret = soft_offline_page(p, MF_COUNT_INCREASED);
360                         if (ret)
361                                 return ret;
362                         continue;
363                 }
```

```
364                     pr_info("Injecting memory failure for page %#lx at %#lx\n",
365                             page_to_pfn(p), start);
366                     /* Ignore return value for now */
367                     memory_failure(page_to_pfn(p), 0, MF_COUNT_INCREASED);
368             }
369             return 0;
370     }
371     #endif
372
373     static long
374     madvise_vma(struct vm_area_struct *vma, struct vm_area_struct **prev,
375                     unsigned long start, unsigned long end, int behavior)
376     {
377             switch (behavior) {
378             case MADV_REMOVE:
379                     return madvise_remove(vma, prev, start, end);
380             case MADV_WILLNEED:
381                     return madvise_willneed(vma, prev, start, end);
382             case MADV_DONTNEED:
383                     return madvise_dontneed(vma, prev, start, end);
384             default:
385                     return madvise_behavior(vma, prev, start, end, behavior);
386             }
387     }
388
389     static bool
390     madvise_behavior_valid(int behavior)
391     {
392             switch (behavior) {
393             case MADV_DOFORK:
394             case MADV_DONTFORK:
395             case MADV_NORMAL:
396             case MADV_SEQUENTIAL:
397             case MADV_RANDOM:
398             case MADV_REMOVE:
399             case MADV_WILLNEED:
400             case MADV_DONTNEED:
401     #ifdef CONFIG_KSM
402             case MADV_MERGEABLE:
403             case MADV_UNMERGEABLE:
404     #endif
405     #ifdef CONFIG_TRANSPARENT_HUGEPAGE
406             case MADV_HUGEPAGE:
407             case MADV_NOHUGEPAGE:
408     #endif
409             case MADV_DONTDUMP:
410             case MADV_DODUMP:
411                     return true;
```

```
412
413         default:
414                 return false;
415         }
416  }
417
418  /*
419   * The madvise(2) system call.
420   *
421   * Applications can use madvise() to advise the kernel how it should
422   * handle paging I/O in this VM area.  The idea is to help the kernel
423   * use appropriate read-ahead and caching techniques.  The information
424   * provided is advisory only, and can be safely disregarded by the
425   * kernel without affecting the correct operation of the application.
426   *
427   * behavior values:
428   *  MADV_NORMAL - the default behavior is to read clusters.  This
429   *              results in some read-ahead and read-behind.
430   *  MADV_RANDOM - the system should read the minimum amount of data
431   *              on any access, since it is unlikely that the appli-
432   *              cation will need more than what it asks for.
433   *  MADV_SEQUENTIAL - pages in the given range will probably be accessed
434   *              once, so they can be aggressively read ahead, and
435   *              can be freed soon after they are accessed.
436   *  MADV_WILLNEED - the application is notifying the system to read
437   *              some pages ahead.
438   *  MADV_DONTNEED - the application is finished with the given range,
439   *              so the kernel can free resources associated with it.
440   *  MADV_REMOVE - the application wants to free up the given range of
441   *              pages and associated backing store.
442   *  MADV_DONTFORK - omit this area from child's address space when forking:
443   *              typically, to avoid COWing pages pinned by get_user_pages().
444   *  MADV_DOFORK - cancel MADV_DONTFORK: no longer omit this area when forking.
445   *  MADV_MERGEABLE - the application recommends that KSM try to merge pages in
446   *              this area with pages of identical content from other such areas.
447   *  MADV_UNMERGEABLE- cancel MADV_MERGEABLE: no longer merge pages with others.
448   *
449   * return values:
450   *  zero    - success
451   *  -EINVAL - start + len < 0, start is not page-aligned,
452   *              "behavior" is not a valid value, or application
453   *              is attempting to release locked or shared pages.
454   *  -ENOMEM - addresses in the specified range are not currently
455   *              mapped, or are outside the AS of the process.
456   *  -EIO    - an I/O error occurred while paging in data.
457   *  -EBADF  - map exists, but area maps something that isn't a file.
458   *  -EAGAIN - a kernel resource was temporarily unavailable.
459   */
```

```c
460    SYSCALL_DEFINE3(madvise, unsigned long, start, size_t, len_in, int, behavior)
461    {
462            unsigned long end, tmp;
463            struct vm_area_struct *vma, *prev;
464            int unmapped_error = 0;
465            int error = -EINVAL;
466            int write;
467            size_t len;
468            struct blk_plug plug;
469
470    #ifdef CONFIG_MEMORY_FAILURE
471            if (behavior == MADV_HWPOISON || behavior == MADV_SOFT_OFFLINE)
472                    return madvise_hwpoison(behavior, start, start+len_in);
473    #endif
474            if (!madvise_behavior_valid(behavior))
475                    return error;
476
477            if (start & ~PAGE_MASK)
478                    return error;
479            len = (len_in + ~PAGE_MASK) & PAGE_MASK;
480
481            /* Check to see whether len was rounded up from small -ve to zero */
482            if (len_in && !len)
483                    return error;
484
485            end = start + len;
486            if (end < start)
487                    return error;
488
489            error = 0;
490            if (end == start)
491                    return error;
492
493            write = madvise_need_mmap_write(behavior);
494            if (write)
495                    down_write(&current->mm->mmap_sem);
496            else
497                    down_read(&current->mm->mmap_sem);
498
499            /*
500             * If the interval [start,end) covers some unmapped address
501             * ranges, just ignore them, but return -ENOMEM at the end.
502             * - different from the way of handling in mlock etc.
503             */
504            vma = find_vma_prev(current->mm, start, &prev);
505            if (vma && start > vma->vm_start)
506                    prev = vma;
507
```

```
508                    blk_start_plug(&plug);
509                    for (;;) {
510                            /* Still start < end. */
511                            error = -ENOMEM;
512                            if (!vma)
513                                    goto out;
514
515                            /* Here start < (end|vma->vm_end). */
516                            if (start < vma->vm_start) {
517                                    unmapped_error = -ENOMEM;
518                                    start = vma->vm_start;
519                                    if (start >= end)
520                                            goto out;
521                            }
522
523                            /* Here vma->vm_start <= start < (end|vma->vm_end) */
524                            tmp = vma->vm_end;
525                            if (end < tmp)
526                                    tmp = end;
527
528                            /* Here vma->vm_start <= start < tmp <= (end|vma->vm_end). */
529                            error = madvise_vma(vma, &prev, start, tmp, behavior);
530                            if (error)
531                                    goto out;
532                            start = tmp;
533                            if (prev && start < prev->vm_end)
534                                    start = prev->vm_end;
535                            error = unmapped_error;
536                            if (start >= end)
537                                    goto out;
538                            if (prev)
539                                    vma = prev->vm_next;
540                            else    /* madvise_remove dropped mmap_sem */
541                                    vma = find_vma(current->mm, start);
542                    }
543    out:
544            blk_finish_plug(&plug);
545            if (write)
546                    up_write(&current->mm->mmap_sem);
547            else
548                    up_read(&current->mm->mmap_sem);
549
550            return error;
551    }
```