```
1   2. Locking brings a performance vs. complexity trade-off
2
3   /*
4    *       linux/mm/filemap.c
5    *
6    * Copyright (C) 1994-1999  Linus Torvalds
7    */
8
9   /*
10   * This file handles the generic file mmap semantics used by
11   * most "normal" filesystems (but you don't /have/ to use this:
12   * the NFS filesystem used to do this differently, for example)
13   */
14  #include <linux/export.h>
15  #include <linux/compiler.h>
16  #include <linux/dax.h>
17  #include <linux/fs.h>
18  #include <linux/sched/signal.h>
19  #include <linux/uaccess.h>
20  #include <linux/capability.h>
21  #include <linux/kernel_stat.h>
22  #include <linux/gfp.h>
23  #include <linux/mm.h>
24  #include <linux/swap.h>
25  #include <linux/mman.h>
26  #include <linux/pagemap.h>
27  #include <linux/file.h>
28  #include <linux/uio.h>
29  #include <linux/hash.h>
30  #include <linux/writeback.h>
31  #include <linux/backing-dev.h>
32  #include <linux/pagevec.h>
33  #include <linux/blkdev.h>
34  #include <linux/security.h>
35  #include <linux/cpuset.h>
36  #include <linux/hugetlb.h>
37  #include <linux/memcontrol.h>
38  #include <linux/cleancache.h>
39  #include <linux/shmem_fs.h>
40  #include <linux/rmap.h>
41  #include "internal.h"
42
43  #define CREATE_TRACE_POINTS
44  #include <trace/events/filemap.h>
45
46  /*
47   * FIXME: remove all knowledge of the buffer layer from the core VM
48   */
49  #include <linux/buffer_head.h> /* for try_to_free_buffers */
50
51  #include <asm/mman.h>
52
53  /*
54   * Shared mappings implemented 30.11.1994. It's not fully working yet,
55   * though.
56   *
57   * Shared mappings now work. 15.8.1995  Bruno.
58   *
59   * finished 'unifying' the page and buffer cache and SMP-threaded the
60   * page-cache, 21.05.1999, Ingo Molnar <mingo@redhat.com>
61   *
62   * SMP-threaded pagemap-LRU 1999, Andrea Arcangeli <andrea@suse.de>
63   */
64
65  /*
66   * Lock ordering:
67   *
68   *  ->i_mmap_rwsem              (truncate_pagecache)
69   *    ->private_lock            (__free_pte->__set_page_dirty_buffers)
70   *      ->swap_lock             (exclusive_swap_page, others)
71   *        ->i_pages lock
72   *
73   *  ->i_mutex
```

```
74   *    ->i_mmap_rwsem           (truncate->unmap_mapping_range)
75   *
76   *  ->mmap_sem
77   *    ->i_mmap_rwsem
78   *      ->page_table_lock or pte_lock   (various, mainly in memory.c)
79   *        ->i_pages lock        (arch-dependent flush_dcache_mmap_lock)
80   *
81   *  ->mmap_sem
82   *    ->lock_page               (access_process_vm)
83   *
84   *  ->i_mutex                   (generic_perform_write)
85   *    ->mmap_sem                (fault_in_pages_readable->do_page_fault)
86   *
87   *  bdi->wb.list_lock
88   *    sb_lock                   (fs/fs-writeback.c)
89   *    ->i_pages lock            (__sync_single_inode)
90   *
91   *  ->i_mmap_rwsem
92   *    ->anon_vma.lock           (vma_adjust)
93   *
94   *  ->anon_vma.lock
95   *    ->page_table_lock or pte_lock     (anon_vma_prepare and various)
96   *
97   *  ->page_table_lock or pte_lock
98   *    ->swap_lock               (try_to_unmap_one)
99   *    ->private_lock            (try_to_unmap_one)
100  *    ->i_pages lock            (try_to_unmap_one)
101  *    ->zone_lru_lock(zone)     (follow_page->mark_page_accessed)
102  *    ->zone_lru_lock(zone)     (check_pte_range->isolate_lru_page)
103  *    ->private_lock            (page_remove_rmap->set_page_dirty)
104  *    ->i_pages lock            (page_remove_rmap->set_page_dirty)
105  *    bdi.wb->list_lock         (page_remove_rmap->set_page_dirty)
106  *    ->inode->i_lock           (page_remove_rmap->set_page_dirty)
107  *    ->memcg->move_lock        (page_remove_rmap->lock_page_memcg)
108  *    bdi.wb->list_lock         (zap_pte_range->set_page_dirty)
109  *    ->inode->i_lock           (zap_pte_range->set_page_dirty)
110  *    ->private_lock            (zap_pte_range->__set_page_dirty_buffers)
111  *
112  * ->i_mmap_rwsem
113  *   ->tasklist_lock            (memory_failure, collect_procs_ao)
114  */
115
116 static int page_cache_tree_insert(struct address_space *mapping,
117                                   struct page *page, void **shadowp)
118 {
119         struct radix_tree_node *node;
120         .....
121
122 [the point is: fine-grained locking leads to complexity.]
```

```
66  Finally, here are some references on this topic:
67
68      --http://www.aristeia.com/Papers/DDJ_Jul_Aug_2004_revised.pdf
69       explores issues with this pattern in C++
70
71      --The "Double-Checked Locking is Broken" Declaration:
72      http://www.cs.umd.edu/~pugh/java/memoryModel/DoubleCheckedLocking.html
73
74      --C++11 provides a way to implement the pattern correctly and
75      portably (again, using memory barriers):
76      https://preshing.com/20130930/double-checked-locking-is-fixed-in-cpp11/
```

```
1   CS 202, Spring 2024
2   Handout 8 (Class 9)
3
4   Therac-25
5
6   1. Software problem #1 (our best guess)
7
8       A. Three threads:
9
10          --Hand: sets the collimator/turntable position
11
12          --Treat: sets a bunch of other parameters. Part of its job takes
13          eight seconds, during which time it's ignoring everything else.
14
15          --Vtkbp (keyboard handler): invoked when user types. It parses
16          the input, and writes to a two-byte shared variable, "MEOS" (mode/energy
17          offset)
18              --"Treat" reads top byte, sets current and energy
19              --"Hand" reads bottom byte, sets the collimator/turntable position
20
21      B. Pseudocode:
22
23          Vtkbp (gets and parses keyboard input):
24
25              data_completion_flag = 0
26
27              while (1) {
28                  wait_for_keyboard_activity();
29                  /* there was some keyboard activity; let's check it */
30                  if (cursor_in_bottom_right) {
31                      parse_the_input();
32                      set the MEOS variable
33                      set data_completion_flag = 1;
34                      signal hand thread
35                      signal treat thread
36                  } else {
37                      /* operator still typing */
38                      data_completion_flag = 0;
39                  }
40                  yield();
41              }
42
43
44          Hand (sets the turntable position):
45
46              while (1) {
47                  wait until signalled
48                  read bottom byte of MEOS variable
49                  /* next line executes quickly */
50                  set turntable position
51                  yield();
52              }
53
54          Treat (sets a bunch of parameters and delivers treatment):
55
56              dataent() { /* this is a subroutine that was called */
57
58                  while (1) {
59                      wait until signalled
60                      read top byte of MEOS variable
61                      set_energy_and_current();
62                      set_bending_magnets(); /* this takes eight seconds */
63                      if (data_completion_flag == 1)
64                          break;
65                  }
66                  /*
67                   * now we leave the subroutine and progress to a state in
68                   * which the machine will accept a "beam on" command
69                   */
70                  return;
71              }
72
```

```
73  2. Software problem #2 (simplified)
74
75      [Simplifying here and condensing to one thread of control; in
76      reality, the functions below are spread over two different threads,
77      but that is not actually the problem, despite what the paper
78      sometimes says. The problem appears to be given by the following
79      simplified description.]
80
81      class3 = 0;
82
83      while (1) {
84
85          if (in field light position) {
86              increment class3;
87          }
88
89          check whether operator pressed "set"
90
91          if (operator pressed set) {
92              if (class3 != 0) {
93                  move turntable out of field light position;
94              }
95              break;
96          }
97      }
98
99      What's the issue here? (Hint: class3 is only one byte.)
100
```