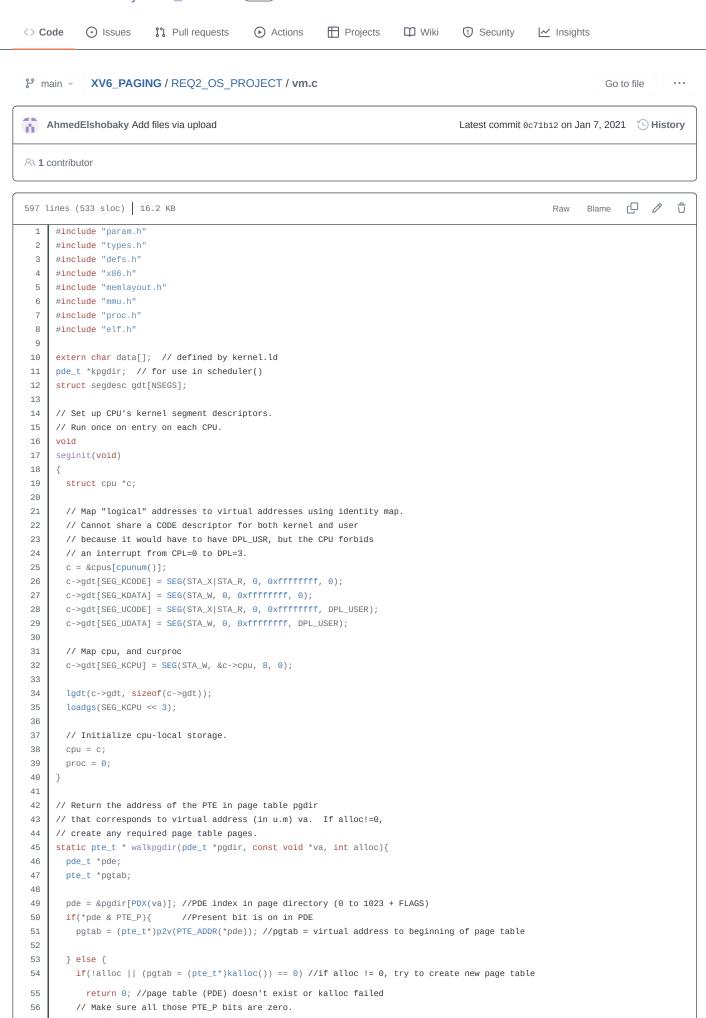
AhmedElshobaky / XV6_PAGING (Public)



```
memset(pqtab, 0, PGSIZE);
 58
          \ensuremath{//} The permissions here are overly generous, but they can
 59
          // be further restricted by the permissions in the page table
 60
          // entries, if necessary.
 61
          *pde = v2p(pgtab) | PTE_P | PTE_W | PTE_U; //link PDE to the new page table
 62
        return &pgtab[PTX(va)]; //return PTE in page table which corresponse to va address
 63
 64
 65
 66
 67
      // Create PTEs for virtual addresses starting at va (va in U.M) that refer to
 68
      // physical addresses starting at pa. va and size might not
 69
      // be page-aligned.
      static int mappages(pde_t *pgdir, void *va, uint size, uint pa, int perm){
 70
        char *a, *last;
 71
        pte_t *pte;
 72
 73
 74
        a = (char*)PGROUNDDOWN((uint)va):
 75
        last = (char*)PGROUNDDOWN(((uint)va) + size - 1);
 76
 77
         if((pte = walkpgdir(pgdir, a, 1)) == 0)
 78
           return -1;
 79
          if(*pte & PTE_P)
 80
                                  //PTE was already initialized for some reason
           panic("remap");
          *pte = pa | perm | PTE_P; //adds page physical address, flags, present bit
 81
          if(a == last)
 82
 83
           break;
 84
          a += PGSIZE;
 85
          pa += PGSIZE;
 86
 87
        return 0;
 88
 89
 90
      // There is one page table per process, plus one that's used when
 91
      // a CPU is not running any process (kpgdir). The kernel uses the
 92
      // current process's page table during system calls and interrupts;
      // page protection bits prevent user code from using the kernel's
 93
 94
      // mappings.
 95
 96
      // setupkvm() and exec() set up every page table like this:
 97
      //
 98
      //
           0..KERNBASE: user memory (text+data+stack+heap), mapped to
                        phys memory allocated by the kernel
 99
      //
100
      //
          KERNBASE..KERNBASE+EXTMEM: mapped to 0..EXTMEM (for I/O space)
           KERNBASE+EXTMEM..data: mapped to EXTMEM..V2P(data)
101
      //
102
      //
                       for the kernel's instructions and r/o data
103
           data..KERNBASE+PHYSTOP: mapped to V2P(data)..PHYSTOP,
104
      //
                                          rw data + free physical memory
           0xfe000000..0: mapped direct (devices such as ioapic)
105
      //
106
107
      // The kernel allocates physical memory for its heap and for user memory
108
      // between V2P(end) and the end of physical memory (PHYSTOP)
      // (directly addressable from end..P2V(PHYSTOP)).
109
110
111
      // This table defines the kernel's mappings, which are present in
112
      // every process's page table.
      static struct kmap {
113
114
        void *virt;
115
        uint phys_start;
        uint phys_end;
116
117
       int perm;
118
     } kmap[] = {
119
       { (void*)KERNBASE, 0,
                                        EXTMEM,
                                                  PTE_W}, // I/O space
       { (void*)KERNLINK, V2P(KERNLINK), V2P(data), 0}, // kern text+rodata
120
                       V2P(data),
121
       { (void*)data,
                                        PHYSTOP, PTE_W}, // kern data+memory
122
       { (void*)DEVSPACE, DEVSPACE,
                                         Θ,
                                                    PTE_W}, // more devices
123
124
125
     // Set up kernel part of a page table.
126
     pde_t* setupkvm(void){
127
       pde_t *pgdir;
128
        struct kmap *k:
129
130
        if((pgdir = (pde_t*)kalloc()) == 0)
```

```
131
          return 0;
132
        memset(pgdir, 0, PGSIZE);
        if (p2v(PHYSTOP) > (void*)DEVSPACE)
133
134
         panic("PHYSTOP too high");
135
        for(k = kmap; k < kmap[NELEM(kmap)]; k++){
         if(mappages(pgdir, k->virt, k->phys_end - k->phys_start, (uint)k->phys_start, k->perm) < 0)</pre>
136
137
            return 0;
138
139
        return pgdir;
140
141
142
      // Allocate one page table for the machine for the kernel address
143
      // space for scheduler processes.
144
      kvmalloc(void)
145
146
147
        kpgdir = setupkvm();
148
        switchkvm();
149
150
151
      // Switch h/w page table register to the kernel-only page table,
      // for when no process is running.
152
153
154
      switchkvm(void)
155
        lcr3(v2p(kpqdir)); // switch to the kernel page table
156
157
158
159
      // Switch TSS and h/w page table to correspond to process p.
160
161
      switchuvm(struct proc *p)
162
163
        pushcli();
        cpu->gdt[SEG_TSS] = SEG16(STS_T32A, &cpu->ts, sizeof(cpu->ts)-1, 0);
164
165
        cpu->gdt[SEG_TSS].s = 0;
166
        cpu->ts.ss0 = SEG_KDATA << 3;
        cpu->ts.esp0 = (uint)proc->kstack + KSTACKSIZE;
167
168
        ltr(SEG_TSS << 3);</pre>
        if(p->pgdir == 0)
169
170
         panic("switchuvm: no pgdir");
        lcr3(v2p(p->pgdir)); // switch to new address space
171
172
        popcli();
173
174
      // Load the initcode into address 0 of podir.
175
176
      // sz must be less than a page.
177
178
      inituvm(pde_t *pgdir, char *init, uint sz)
179
180
        char *mem;
181
182
        if(sz >= PGSIZE)
         panic("inituvm: more than a page");
183
184
        mem = kalloc();
185
        memset(mem, 0, PGSIZE);
        mappages(pgdir, 0, PGSIZE, v2p(mem), PTE_W|PTE_U);
186
187
        memmove(mem, init, sz);
188
189
      // Load a program segment into pgdir. addr must be page-aligned
190
191
      // and the pages from addr to addr+sz must already be mapped.
192
193
      loaduvm(pde_t *pgdir, char *addr, struct inode *ip, uint offset, uint sz)
194
195
        uint i, pa, n;
196
        pte_t *pte;
197
        if((uint) addr % PGSIZE != 0)
198
199
         panic("loaduvm: addr must be page aligned");
200
        for(i = 0; i < sz; i += PGSIZE){</pre>
201
         if((pte = walkpgdir(pgdir, addr+i, 0)) == 0)
202
            panic("loaduvm: address should exist"):
203
          pa = PTE_ADDR(*pte);
204
          if(sz - i < PGSIZE)</pre>
```

```
205
            n = sz - 1;
206
          else
207
            n = PGSIZE;
208
          if(readi(ip, p2v(pa), offset+i, n) != n)
209
            return -1;
210
211
        return 0;
212
213
214
      int getPagePAddr(int userPageVAddr, pde_t * pgdir){
215
        pte_t *pte;
216
        pte = walkpgdir(pgdir, (int*)userPageVAddr, 0);
217
        if(!pte) //uninitialized page table
218
          return -1;
        return PTE_ADDR(*pte);
219
220
221
222
      void fixPagedOutPTE(int userPageVAddr, pde t * pgdir){
223
        pte t *pte;
224
        pte = walkpgdir(pgdir, (int*)userPageVAddr, 0);
225
          panic("PTE of swapped page is missing");
226
227
        *pte |= PTE_PG;
228
        *pte &= ~PTE_P;
229
        *pte &= PTE_FLAGS(*pte); //clear junk physical address
        lcr3(v2p(proc->pgdir)); //refresh CR3 register
230
231
232
233
      //This method cannot be replaced with mappages because mappages cannot turn off PTE_PG bit
234
      void fixPagedInPTE(int userPageVAddr, int pagePAddr, pde_t * pgdir){
235
        pte_t *pte;
236
        pte = walkpgdir(pgdir, (int*)userPageVAddr, 0);
237
        if (!pte)
238
         panic("PTE of swapped page is missing");
239
        if (*pte & PTE_P)
240
              panic("PAGE IN REMAP!");
        *pte |= PTE_P | PTE_W | PTE_U;
                                           //Turn on needed bits
241
242
        *pte &= ~PTE PG:
                                                                                        //Turn off inFile bit
243
        *pte |= pagePAddr;
                                                                                        //Map PTE to the new Page
244
        lcr3(v2p(proc->pgdir)); //refresh CR3 register
245
246
247
      int pageIsInFile(int userPageVAddr, pde_t * pgdir) {
248
        pte = walkpgdir(pgdir, (char *)userPageVAddr, 0);
249
250
        return (*pte & PTE_PG); //PAGE IS IN FILE
251
252
253
254
      int getFIFO(){
255
          pte_t * pte;
256
          int i = 0;
257
          int pageIndex;
258
         uint loadOrder;
259
        recheck:
260
         pageIndex = -1;
          loadOrder = 0xFFFFFFF;
261
262
          for (i = 0; i < MAX_PYSC_PAGES; i++) {</pre>
263
            if (proc->ramCtrlr[i].state == USED && proc->ramCtrlr[i].loadOrder <= loadOrder){</pre>
264
              pageIndex = i;
265
              loadOrder = proc->ramCtrlr[i].loadOrder;
266
            }
267
          }
268
          return pageIndex;
269
270
271
      int getLRU(){
272
273
        int i;
274
        int pageIndex = -1;
275
        uint minAccess = 0xffffffff;
276
277
        for (i = 0; i < MAX_PYSC_PAGES; i++) {</pre>
278
          if (proc->ramCtrlr[i].state == USED && proc->ramCtrlr[i].accessCount <= minAccess) {</pre>
```

```
279
                minAccess = proc->ramCtrlr[i].accessCount;
280
                pageIndex = i;
281
          }
282
        }
283
        return pageIndex;
284
285
286
      int getPageOutIndex(){
287
288
        #if FIF0
         return getFIFO();
289
290
291
        #if LRU
292
         return getLRU();
293
        #endif
294
        panic("Unrecognized paging machanism");
295
296
297
      void updateAccessCounters(struct proc * p){
298
        pte t * pte;
        int i;
299
        for (i = 0; i < MAX PYSC PAGES; i++) {</pre>
300
301
          if (p->ramCtrlr[i].state == USED){
302
            pte = walkpgdir(p->ramCtrlr[i].pgdir, (char*)p->ramCtrlr[i].userPageVAddr,0);
303
            if (*pte & PTE_A) {
              *pte &= ~PTE_A; // turn off PTE_A flag
304
305
               p->ramCtrlr[i].accessCount++;
306
307
308
        }
309
      }
310
311
      int getFreeRamCtrlrIndex() {
312
        if (proc == 0)
313
         return -1;
314
        for (i = 0; i < MAX_PYSC_PAGES; i++) {</pre>
315
316
         if (proc->ramCtrlr[i].state == NOTUSED)
317
318
        return -1; //NO ROOM IN RAMCTRLR
319
320
321
322
      static char buff[PGSIZE]; //buffer used to store swapped page in getPageFromFile method
323
324
      int getPageFromFile(int cr2){
325
        proc->faultCounter++;
326
        int userPageVAddr = PGROUNDDOWN(cr2);
327
        char * newPg = kalloc();
328
        memset(newPg, 0, PGSIZE);
329
        int outIndex = getFreeRamCtrlrIndex();
330
        lcr3(v2p(proc->pgdir)); //refresh CR3 register
331
        if (outIndex >= 0) { //Free location in RamCtrlr is available, no need for swapping
332
          fixPagedInPTE(userPageVAddr, v2p(newPg), proc->pgdir);
333
          readPageFromFile(proc, outIndex, userPageVAddr, (char*)userPageVAddr);
334
         return 1; //Operation was successful
335
        }
336
        proc->countOfPagedOut++;
337
        //If reached here - Swapping is needed.
338
        outIndex = getPageOutIndex(); //select a page to swap to file
339
        struct pagecontroller outPage = proc->ramCtrlr[outIndex];
340
        fixPagedInPTE(userPageVAddr, v2p(newPg), proc->pgdir);
341
        readPageFromFile(proc, outIndex, userPageVAddr, buff); //automatically adds to ramctrlr
342
        int outPagePAddr = getPagePAddr(outPage.userPageVAddr, outPage.pgdir);
343
        memmove(newPg, buff, PGSIZE);
344
        writePageToFile(proc, outPage.userPageVAddr, outPage.pgdir);
345
        fixPagedOutPTE(outPage.userPageVAddr, outPage.pgdir);
        char *v = p2v(outPagePAddr);
346
347
        kfree(v); //free swapped page
348
        return 1;
349
      }
350
351
      void addToRamCtrlr(pde_t *pgdir, uint userPageVAddr) {
352
        int freeLocation = getFreeRamCtrlrIndex();
```

```
proc->ramCtrlr[treeLocation].state = USED;
354
        proc->ramCtrlr[freeLocation].pgdir = pgdir;
355
        proc->ramCtrlr[freeLocation].userPageVAddr = userPageVAddr;
356
        proc->ramCtrlr[freeLocation].loadOrder = proc->loadOrderCounter++;
357
        proc->ramCtrlr[freeLocation].accessCount = 0;
358
359
360
361
      void swap(pde_t *pgdir, uint userPageVAddr){
362
        proc->countOfPagedOut++;
363
        int outIndex = getPageOutIndex();
364
        int outPagePAddr = getPagePAddr(proc->ramCtrlr[outIndex].userPageVAddr, proc->ramCtrlr[outIndex].pgdir);
365
        writePageToFile(proc, proc->ramCtrlr[outIndex].userPageVAddr, proc->ramCtrlr[outIndex].pgdir);
366
        char *v = p2v(outPagePAddr);
367
        kfree(v); //free swapped page
        proc->ramCtrlr[outIndex].state = NOTUSED;
368
369
        fixPagedOutPTE(proc->ramCtrlr[outIndex].userPageVAddr, proc->ramCtrlr[outIndex].pgdir);
370
        addToRamCtrlr(pgdir, userPageVAddr);
371
372
373
374
      int isNONEpolicy(){
375
              #if NONE
376
                       return 1;
377
              #endif
378
              return 0;
379
380
      // Allocate page tables and physical memory to grow process from oldsz to
381
      \ensuremath{//} newsz, which need not be page aligned. Returns new size or 0 on error.
      int allocuvm(pde_t *pgdir, uint oldsz, uint newsz){
382
383
        char *mem;
384
        uint a;
385
        if(newsz >= KERNBASE)
386
          return 0:
387
        if(newsz < oldsz)</pre>
388
          return oldsz;
389
390
        if (!isNONEpolicv()){
           if (PGROUNDUP(newsz)/PGSIZE > MAX_TOTAL_PAGES && proc->pid > 2) {
391
392
                          cprintf("proc is too big\n", PGROUNDUP(newsz)/PGSIZE);
393
                           return 0;
394
                        }
395
              }
396
397
        a = PGROUNDUP(oldsz);
398
        int i = 0; //debugging
        for(; a < newsz; a += PGSIZE){</pre>
399
400
          mem = kalloc();
401
          i++;
402
          if(mem == 0){
403
            cprintf("allocuvm out of memory\n");
404
            deallocuvm(pgdir, newsz, oldsz);
405
            return 0:
406
          }
407
          memset(mem, 0, PGSIZE);
          mappages(pgdir, (char*)a, PGSIZE, v2p(mem), PTE_W|PTE_U);
408
409
          if (!isNONEpolicy() && proc->pid > 2){
410
            if (PGROUNDUP(oldsz)/PGSIZE + i > MAX_PYSC_PAGES)
411
              swap(pgdir, a);
412
            else //there's room
              addToRamCtrlr(pgdir, a);
413
414
415
        }
416
        return newsz;
417
418
419
      //This must use userVaddress+pgdir addresses!
420
421
      //(The proc has identical vAddresses on different page directories until exec finish executing)
422
      void removeFromRamCtrlr(uint userPageVAddr, pde_t *pgdir){
423
        if (proc == 0)
424
          return:
425
        int i;
        for (i = 0; i < MAX_PYSC_PAGES; i++) {</pre>
426
```

```
it (proc->ramCtrlr[1].state == USED
427
428
              && proc->ramCtrlr[i].userPageVAddr == userPageVAddr
              && proc->ramCtrlr[i].pgdir == pgdir){
429
430
            proc->ramCtrlr[i].state = NOTUSED;
431
            return;
432
        }
433
434
      }
435
436
      void removeFromFileCtrlr(uint userPageVAddr, pde_t *pgdir){
437
        if (proc == 0)
438
          return;
439
        int i;
440
        for (i = 0; i < MAX TOTAL PAGES-MAX PYSC PAGES; i++) {</pre>
441
          if (proc->fileCtrlr[i].state == USED
442
              && proc->fileCtrlr[i].userPageVAddr == userPageVAddr
443
              && proc->fileCtrlr[i].pgdir == pgdir){
            proc->fileCtrlr[i].state = NOTUSED;
444
445
            return;
446
447
        }
448
449
      // Deallocate user pages to bring the process size from oldsz to
450
      // newsz. oldsz and newsz need not be page-aligned, nor does newsz
451
      // need to be less than oldsz. oldsz can be larger than the actual
452
      // process size. Returns the new process size.
453
      int deallocuvm(pde_t *pgdir, uint oldsz, uint newsz){
454
        pte_t *pte;
455
        uint a, pa;
456
457
        if(newsz >= oldsz)
458
          return oldsz;
459
        a = PGROUNDUP(newsz);
460
461
        int i = 0; //debugging
462
        for(; a < oldsz; a += PGSIZE){</pre>
463
          pte = walkpgdir(pgdir, (char*)a, 0);
464
          if(!pte) //uninitialized page table
465
            a += (NPTENTRIES - 1) * PGSIZE; //jump to next page table
466
          else if((*pte & PTE_P) != 0){ //page table exists and page is present
                                            //pa = beginning of page physical address
467
            pa = PTE_ADDR(*pte);
468
            if(pa == 0)
469
              panic("kfree");
470
            char v = p2v(pa);
471
            kfree(v); //free page
472
            if (!isNONEpolicy())
473
              removeFromRamCtrlr(a, pgdir);
474
475
            i++;
476
            *pte = 0;
477
478
479
        return newsz:
480
481
      // Free a page table and all the physical memory pages
482
483
      // in the user part.
484
      void freevm(pde_t *pgdir){
485
        uint i;
        if(pgdir == 0)
486
487
         panic("freevm: no pgdir");
488
        deallocuvm(pgdir, KERNBASE, 0);
489
        int j = 0;
        for(i = 0; i < NPDENTRIES; i++){</pre>
490
491
          if(pgdir[i] & PTE_P){ //PDE exists
492
            char * v = p2v(PTE_ADDR(pgdir[i]));
493
            kfree(v); //free page table
            j++;
494
495
         }
496
497
        kfree((char*)pgdir); //free page directory
498
499
500
      // Clear PTE_U on a page. Used to create an inaccessible
```

```
// page beneath the user stack.
502
      void
503
      clearpteu(pde_t *pgdir, char *uva)
504
505
        pte_t *pte;
506
507
        pte = walkpgdir(pgdir, uva, 0);
508
        if(pte == 0)
509
         panic("clearpteu");
510
        *pte &= ~PTE_U;
511
512
513
      // Given a parent process's page table, create a copy
514
      // of it for a child.
515
      pde_t* copyuvm(pde_t *pgdir, uint sz){
516
        pde_t *d;
517
        pte_t *pte;
518
        uint pa, i, flags;
519
        char *mem;
520
521
        if((d = setupkvm()) == 0)
522
         return 0;
523
        int j = 0;
524
        for(i = 0; i < sz; i += PGSIZE){</pre>
          if((pte = walkpgdir(pgdir, (void *) i, 0)) == 0)
525
            panic("copyuvm: pte should exist");
526
527
          if (*pte & PTE_PG){
528
              fixPagedOutPTE(i, d);
529
              continue;
530
          }
531
532
          if(!(*pte & PTE_P))
533
           panic("copyuvm: page not present");
          pa = PTE_ADDR(*pte);
534
535
          flags = PTE_FLAGS(*pte);
536
          if((mem = kalloc()) == 0)
537
           goto bad;
538
          memmove(mem, (char*)p2v(pa), PGSIZE);
539
540
          if(mappages(d, (void*)i, PGSIZE, v2p(mem), flags) < 0)</pre>
541
           goto bad;
542
        }
543
        return d;
544
545
      bad:
546
        freevm(d);
547
        return 0;
548
549
550
      //PAGEBREAK!
551
      // Map user virtual address to kernel address.
552
      char*
553
      uva2ka(pde_t *pgdir, char *uva)
554
555
        pte_t *pte;
556
557
        pte = walkpgdir(pgdir, uva, 0);
558
        if((*pte & PTE_P) == 0)
559
          return 0;
560
        if((*pte & PTE_U) == 0)
561
         return 0;
562
        return (char*)p2v(PTE_ADDR(*pte));
563
564
565
      // Copy len bytes from p to user address va in page table pgdir.
566
      // Most useful when pgdir is not the current page table.
567
      // uva2ka ensures this only works for PTE_U pages.
568
      int
      copyout(pde_t *pgdir, uint va, void *p, uint len)
569
570
        char *buf, *pa0;
571
572
        uint n, va0;
573
574
        buf = (char*)p;
```

2/12/22, 8:57 PM

```
575
       while(len > 0){
         va0 = (uint)PGROUNDDOWN(va);
576
577
         pa0 = uva2ka(pgdir, (char*)va0);
578
         if(pa0 == 0)
579
          return -1;
         n = PGSIZE - (va - va0);
580
         if(n > len)
581
582
          n = len;
583
         memmove(pa0 + (va - va0), buf, n);
584
         len -= n;
585
         buf += n;
586
         va = va0 + PGSIZE;
587
588
        return 0;
589
     }
590
591
     //PAGEBREAK!
592
     // Blank page.
593
     //PAGEBREAK!
594
     // Blank page.
595
     //PAGEBREAK!
     // Blank page.
596
```