

xv6 is a re-implementation of Dennis Ritchie's and Ken Thompson's Unix Version 6 (v6). xv6 loosely follows the structure and style of v6, but is implemented for a modern x86-based multiprocessor using ANSI C.

ACKNOWLEDGMENTS

xv6 is inspired by John Lions's Commentary on UNIX 6th Edition (Peer to Peer Communications; ISBN: 1-57398-013-7; 1st edition (June 14, 2000)). See also <http://pdos.csail.mit.edu/6.828/2016/xv6.html>, which provides pointers to on-line resources for v6.

xv6 borrows code from the following sources:

- JOS (asm.h, elf.h, mmu.h, bootasm.S, ide.c, console.c, and others)
- Plan 9 (entryother.S, mp.h, mp.c, lapic.c)
- FreeBSD (ioapic.c)
- NetBSD (console.c)

The following people have made contributions: Russ Cox (context switching, locking), Cliff Frey (MP), Xiao Yu (MP), Nikolai Zeldovich, and Austin Clements.

We are also grateful for the bug reports and patches contributed by Silas Boyd-Wickizer, Anton Burtsev, Cody Cutler, Mike CAT, Tej Chajed, Nelson Elhage, Saar Ettinger, Alice Ferrazzi, Nathaniel Filardo, Peter Froehlich, Yakir Goaron, Shivam Handa, Bryan Henry, Jim Huang, Alexander Kapshuk, Anders Kaseorg, kehao95, Wolfgang Keller, Eddie Kohler, Austin Liew, Imbar Marinescu, Yandong Mao, Hitoshi Mitake, Carmi Merimovich, Joel Nider, Greg Price, Ayan Shafgat, Eldar Sehayek, Yongming Shen, Cam Tenny, Rafael Ubal, Warren Toomey, Stephen Tu, Pablo Ventura, Xi Wang, Keiichi Watanabe, Nicolas Wolovick, Grant Wu, Jindong Zhang, Icenowy Zheng, and Zou Chang Wei.

The code in the files that constitute xv6 is
Copyright 2006-2016 Frans Kaashoek, Robert Morris, and Russ Cox.

ERROR REPORTS

If you spot errors or have suggestions for improvement, please send email to Frans Kaashoek and Robert Morris (kaashoek,rtm@csail.mit.edu). If you have suggestions for improvements, please keep in mind that the main purpose of xv6 is as a teaching operating system for MIT's 6.828. For example, we are in particular interested in simplifications and clarifications, instead of suggestions for new systems calls, more portability, etc.

BUILDING AND RUNNING XV6

To build xv6 on an x86 ELF machine (like Linux or FreeBSD), run "make". On non-x86 or non-ELF machines (like OS X, even on x86), you will need to install a cross-compiler gcc suite capable of producing x86 ELF binaries. See <http://pdos.csail.mit.edu/6.828/2016/tools.html>. Then run "make TOOLPREFIX=i386-jos-elf-".

To run xv6, install the QEMU PC simulators. To run in QEMU, run "make qemu".

To create a typeset version of the code, run "make xv6.pdf". This requires the "mpage" utility. See <http://www.mesa.nl/pub/mpage/>.

The numbers to the left of the file names in the table are sheet numbers. The source code has been printed in a double column format with fifty lines per column, giving one hundred lines per sheet (or page). Thus there is a convenient relationship between line numbers and sheet numbers.

# basic headers	# system calls	68 pipe.c
01 types.h	32 traps.h	
01 param.h	33 vectors.pl	# string operations
02 memlayout.h	33 trapasm.S	70 string.c
02 defs.h	34 trap.c	
04 x86.h	35 syscall.h	# low-level hardware
06 asm.h	36 syscall.c	71 mp.h
07 mmu.h	38 sysproc.c	73 mp.c
10 elf.h		74 lapic.c
	# file system	77 ioapic.c
# entering xv6	39 buf.h	78 kbd.h
11 entry.S	39 sleeplock.h	80 kbd.c
12 entryother.S	40 fcntl.h	80 console.c
13 main.c	40 stat.h	84 uart.c
	41 fs.h	
# locks	42 file.h	# user-level
15 spinlock.h	43 ide.c	85 initcode.S
15 spinlock.c	45 bio.c	86 usys.S
	47 sleeplock.c	86 init.c
# processes	48 log.c	87 sh.c
17 vm.c	50 fs.c	
23 proc.h	59 file.c	# bootloader
24 proc.c	61 sysfile.c	92 bootasm.S
30 swtch.S	67 exec.c	93 bootmain.c
31 kalloc.c		
	# pipes	

The source listing is preceded by a cross-reference that lists every defined constant, struct, global variable, and function in xv6. Each entry gives, on the same line as the name, the line number (or, in a few cases, numbers) where the name is defined. Successive lines in an entry list the line numbers where the name is used. For example, this entry:

```
swtch 2658
    0374 2428 2466 2657 2658
```

indicates that swtch is defined on line 2658 and is mentioned on five lines on sheets 03, 24, and 26.

acquire 1574	3912 4394 4418 4423 4460
0380 1574 1578 2478 2548	4478 4586 4619 5039
2614 2649 2677 2769 2830	begin_op 4928
2891 2906 2966 2979 3176	0336 2644 4928 6033 6124
3193 3466 3872 3892 4409	6310 6411 6511 6556 6573
4465 4570 4633 4724 4736	6607 6721
4755 4930 4957 4976 5031	bfree 5152
5358 5391 5462 5475 5980	5152 5564 5574 5577
6004 6018 6913 6934 6955	bget 4566
8160 8331 8378 8414	4566 4596 4606
acquiresleep 4722	binit 4538
0389 4577 4592 4722 5411	0263 1330 4538
5460	bmap 5510
allocproc 2473	5254 5510 5536 5620 5670
2473 2525 2587	bootmain 9367
allocuvum 1927	9318 9367
0430 1927 1941 1947 2565	BPB 4157
6751 6765	4157 4160 5122 5124 5159
alltraps 3354	bread 4602
3309 3317 3330 3335 3353	0264 4602 4877 4878 4890
3354	4906 4990 4991 5085 5106
ALT 7860	5123 5158 5311 5335 5414
7860 7888 7890	5526 5570 5620 5670
argfd 6171	brelse 4626
6171 6223 6238 6257 6268	0265 4626 4629 4881 4882
6281	4897 4914 4994 4995 5087
argint 3652	5109 5129 5134 5165 5317
0404 3652 3666 3682 3833	5320 5344 5422 5532 5576
3856 3870 6176 6238 6257	5623 5674
6508 6575 6576 6632	BFSIZE 4105
argptr 3661	3909 4105 4124 4151 4157
0405 3661 6238 6257 6281	4380 4396 4419 4858 4879
6657	4992 5107 5620 5621 5622
argstr 3679	5666 5670 5671 5672
0406 3679 6307 6408 6508	buf 3900
6557 6574 6608 6632	0250 0264 0265 0266 0308
__attribute__ 1405	0335 2120 2123 2132 2134
0272 0367 1309 1405	3900 3906 3907 3908 4313
BACK 8711	4331 4334 4374 4406 4454
8711 8824 8970 9239	4456 4459 4526 4530 4534
backcmd 8746 8964	4540 4553 4565 4568 4601
8746 8759 8825 8964 8966	4604 4615 4626 4806 4877
9092 9205 9240	4878 4890 4891 4897 4906
BACKSPACE 8250	4907 4913 4914 4990 4991
8250 8267 8309 8342 8348	5022 5070 5083 5104 5119
ballocc 5116	5154 5307 5332 5405 5513
5116 5136 5517 5525 5529	5559 5606 5656 8080 8091
BBLOCK 4160	8095 8098 8318 8340 8354
4160 5123 5158	8388 8409 8416 8834 8837
B_DIRTY 3912	8838 8839 8853 8865 8866

```

8868 8869 8870 8874
B_VALID 3911
3911 4422 4460 4478 4607
bwrite 4615
0266 4615 4618 4880 4913
4993
bzero 5102
5102 5130
C 7881 8324
7881 7929 7954 7955 7956
7957 7958 7960 8324 8334
8338 8345 8356 8389
CAPSLOCK 7862
7862 7895 8036
cgaputc 8255
8255 8313
clearpteu 2022
0439 2022 2028 6767
cli 0557
0557 0559 1224 1660 8210
8304 9262
cmd 8715
8715 8727 8736 8737 8742
8743 8748 8752 8756 8765
8768 8773 8781 8787 8791
8801 8825 8827 8902 8905
8907 8908 8909 8910 8913
8914 8916 8918 8919 8920
8921 8922 8923 8924 8925
8926 8929 8930 8932 8934
8935 8936 8937 8938 8939
8950 8951 8953 8955 8956
8957 8958 8959 8960 8963
8964 8966 8968 8969 8970
8971 8972 9062 9063 9064
9065 9067 9071 9074 9080
9081 9084 9087 9089 9092
9096 9098 9100 9103 9105
9108 9110 9113 9114 9125
9128 9131 9135 9150 9153
9158 9162 9163 9166 9171
9172 9178 9187 9188 9194
9195 9201 9202 9211 9214
9216 9222 9223 9228 9234
9240 9241 9244
CMOS_PORT 7623
7623 7637 7638 7675
CMOS_RETURN 7624
7624 7678
CMOS_STATA 7662

```

```

7662 7713
CMOS_STATB 7663
7663 7706
CMOS_UIP 7664
7664 7713
COM1 8464
8464 8474 8477 8478 8479
8480 8481 8482 8485 8491
8492 8507 8509 8517 8519
commit 5001
4853 4975 5001
CONSOLE 4236
4236 8428 8429
consoleinit 8424
0269 1326 8424
consoleintr 8327
0271 8048 8327 8525
consoleread 8371
8371 8429
consolewrite 8409
8409 8428
consputc 8301
8067 8098 8168 8186 8189
8193 8194 8301 8342 8348
8355 8416
context 2326
0251 0377 2303 2326 2345
2509 2510 2511 2512 2781
2822 3028
CONV 7722
7722 7723 7724 7725 7726
7727 7728 7729
copyout 2118
0438 2118 6775 6786
copyuvm 2035
0435 2035 2046 2048 2592
cprintf 8152
0270 1354 1941 1947 3026
3030 3032 3490 3503 3508
3760 5253 6725 7809 8152
8213 8214 8215 8218
cpu 2301
0311 0363 1354 1368 1506
1566 1590 1608 1647 1717
2301 2312 2436 2458 2761
3490 3503 3508 7313
cpuid 2430
0358 1354 1723 2430 3465
3491 3504 3510
CRO_PE 0727

```

```

0727 1237 1270 9293
CRO_PG 0737
0737 1154 1270
CRO_WP 0733
0733 1154 1270
CR4_PSE 0739
0739 1147 1263
create 6457
6457 6477 6490 6494 6514
6557 6577
CRTPORT 8251
8251 8260 8261 8262 8263
8281 8282 8283 8284
CTL 7859
7859 7885 7889 8035
DAY 7669
7669 7686
deallocuvm 1961
0431 1942 1948 1961 2009
2568
DEVSPACE 0204
0204 1813 1826
devsw 4229
4229 4234 5609 5611 5659
5661 5962 8428 8429
dinode 4128
4128 4151 5308 5312 5333
5336 5406 5415
dirent 4165
4165 5714 5755 6355 6404
dirlink 5752
0288 5752 5767 5775 6330
6489 6493 6494
dirlookup 5711
0289 5711 5717 5721 5759
5875 6423 6467
DIRSIZ 4163
4163 4167 5705 5772 5828
5829 5892 6304 6405 6461
DPL_USER 0828
0828 1726 1727 2533 2534
3423 3518 3527
E0ESC 7866
7866 8020 8024 8025 8027
8030
elfhdr 1005
1005 6715 9369 9374
ELF_MAGIC 1002
1002 6734 9380
ELF_PROG_LOAD 1036

```

```

1036 6745
end_op 4953
0337 2646 4953 6035 6129
6312 6319 6337 6346 6413
6447 6452 6516 6521 6527
6536 6540 6558 6562 6578
6582 6609 6615 6620 6724
6759 6810
entry 1144
1011 1140 1143 1144 3302
3303 6799 7171 9371 9395
9396
EOI 7466
7466 7584 7613
ERROR 7487
7487 7577
ESR 7469
7469 7580 7581
exec 6710
0275 6648 6710 6725 8618
8679 8680 8776 8777
EXEC 8707
8707 8772 8909 9215
execcmd 8719 8903
8719 8760 8773 8903 8905
9171 9177 9178 9206 9216
exit 2627
0359 2627 2665 3455 3459
3519 3528 3818 8567 8570
8611 8676 8681 8766 8775
8785 8830 8877 8884
EXTMEM 0202
0202 0208 1810
fdalloc 6203
6203 6225 6532 6662
fetchint 3617
0407 3617 3654 6639
fetchstr 3631
0408 3631 3684 6645
file 4200
0252 0278 0279 0280 0282
0283 0284 0351 2348 4200
5071 5960 5965 5975 5978
5981 6001 6002 6014 6016
6052 6065 6102 6165 6171
6174 6203 6220 6234 6253
6266 6278 6505 6654 6858
6872 8061 8459 8728 8783
8784 8914 8922 9122
filealloc 5976

```

0278 5976 6532 6878
fileclose 6014
0279 2639 6014 6020 6271
6534 6665 6666 6904 6906
filedup 6002
0280 2607 6002 6006 6227
fileinit 5969
0281 1331 5969
fileread 6065
0282 6065 6080 6240
filestat 6052
0283 6052 6283
filewrite 6102
0284 6102 6134 6139 6259
FL_IF 0710
0710 1662 1669 2441 2537
2819
fork 2580
0360 2580 3812 8610 8673
8675 8892 8894
forkl 8888
8750 8792 8804 8811 8826
8873 8888
forkret 2853
2417 2512 2853
freerange 3151
3111 3135 3141 3151
freevm 2003
0432 1831 2003 2008 2060
2690 6802 6807
FSSIZE 0162
0162 4378
gatedesc 0951
0523 0526 0951 3411
getcallerpcs 1627
0381 1591 1627 3028 8216
getcnd 8834
8834 8865
gettoken 9006
9006 9091 9095 9107 9120
9121 9157 9161 9183
growproc 2558
0361 2558 3859
havedisk1 4333
4333 4363 4462
holding 1645
0382 1577 1604 1645 2813
holdingsleep 4751
0391 4458 4617 4628 4751
5433
HOURS 7668
7668 7685
ialloc 5304
0290 5304 5322 6476 6477
IBLOCK 4154
4154 5311 5335 5414
ICRHI 7480
7480 7587 7645 7656
ICRLO 7470
7470 7588 7589 7646 7648
7657
ID 7463
7463 7499 7605
IDE_BSY 4316
4316 4342
IDE_CMD_RDMUL 4323
4323 4382
IDE_CMD_READ 4321
4321 4382
IDE_CMD_WRITE 4322
4322 4383
IDE_CMD_WRMUL 4324
4324 4383
IDE_DF 4318
4318 4344
IDE_DRDY 4317
4317 4342
IDE_ERR 4319
4319 4344
ideinit 4351
0306 1332 4351
ideintr 4404
0307 3474 4404
idelock 4330
4330 4355 4409 4412 4430
4465 4479 4483
iderw 4454
0308 4454 4459 4461 4463
4608 4620
idestart 4374
4334 4374 4377 4385 4428
4475
idewait 4338
4338 4357 4387 4418
idtinit 3429
0415 1355 3429
idup 5389
0291 2608 5389 5862
iget 5354
5259 5318 5354 5374 5729

5860
iinit 5243
0292 2864 5243
ilock 5403
0293 5403 5409 5425 5865
6055 6074 6125 6316 6329
6342 6417 6425 6465 6469
6479 6524 6612 6728 8383
8403 8418
inb 0453
0453 4342 4362 7446 7678
8014 8017 8261 8263 8485
8491 8492 8507 8517 8519
9273 9281 9404
initlock 1562
0383 1562 2425 3133 3425
4355 4542 4715 4862 5247
5971 6886 8426
initlog 4856
0334 2865 4856 4859
initsleeplock 4713
0392 4556 4713 5249
inituvm 1885
0433 1885 1890 2530
inode 4212
0253 0288 0289 0290 0291
0293 0294 0295 0296 0297
0299 0300 0301 0302 0303
0434 1903 2349 4206 4212
4230 4231 5074 5239 5249
5259 5303 5330 5353 5356
5362 5388 5389 5403 5431
5458 5482 5510 5556 5588
5603 5653 5710 5711 5752
5756 5854 5857 5889 5900
6305 6352 6403 6456 6460
6506 6554 6569 6604 6716
8371 8409
INPUT_BUF 8316
8316 8318 8340 8352 8354
8356 8388
insl 0462
0462 0464 4419 9423
install_trans 4872
4872 4921 5006
INT_DISABLED 7769
7769 7814
ioapic 7777
7408 7425 7426 7774 7777
7786 7787 7793 7794 7805
IOAPIC 7758
7758 7805
ioapicenable 7820
0311 4356 7820 8432 8493
ioapicid 7315
0312 7315 7426 7808 7809
ioapicinit 7801
0313 1325 7801 7809
ioapicread 7784
7784 7806 7807
ioapicwrite 7791
7791 7814 7815 7825 7826
IPB 4151
4151 4154 5312 5336 5415
iput 5458
0294 2645 5458 5485 5760
5883 6034 6335 6619
IRQ_COM1 3283
3283 3484 8493
IRQ_ERROR 3285
3285 7577
IRQ_IDE 3284
3284 3473 3477 4356
IRQ_KBD 3282
3282 3480 8432
IRQ_SPURIOUS 3286
3286 3489 7557
IRQ_TIMER 3281
3281 3464 3523 7564
isdirempty 6352
6352 6359 6429
itrunc 5556
5074 5467 5556
iunlock 5431
0295 5431 5434 5484 5872
6057 6077 6128 6325 6539
6618 8376 8413
iunlockput 5482
0296 5482 5867 5876 5879
6318 6331 6334 6345 6430
6441 6445 6451 6468 6472
6496 6526 6535 6561 6581
6614 6758 6809
iupdate 5330
0297 5330 5469 5582 5679
6324 6344 6439 6444 6483
6487
kalloc 3188
0316 1384 1744 1823 1891
1939 2051 2494 3188 6880

KBDATAP 7854
 7854 8017
 kbdgetc 8006
 8006 8048
 kbdirtr 8046
 0322 3481 8046
 KBS_DIB 7853
 7853 8015
 KBSTATP 7852
 7852 8014
 KERNBASE 0207
 0207 0208 0210 0211 0213
 0214 1410 1634 1810 1932
 2009
 KERNLINK 0208
 0208 1811
 KEY_DEL 7878
 7878 7919 7941 7965
 KEY_DN 7872
 7872 7915 7937 7961
 KEY_END 7870
 7870 7918 7940 7964
 KEY_HOME 7869
 7869 7918 7940 7964
 KEY_INS 7877
 7877 7919 7941 7965
 KEY_LF 7873
 7873 7917 7939 7963
 KEY_PGDN 7876
 7876 7916 7938 7962
 KEY_PGUP 7875
 7875 7916 7938 7962
 KEY_RT 7874
 7874 7917 7939 7963
 KEY_UP 7871
 7871 7915 7937 7961
 kfree 3165
 0317 1949 1977 1979 2013
 2016 2593 2688 3156 3165
 3170 6902 6923
 kill 2975
 0362 2975 3509 3835 8617
 kinit1 3131
 0318 1319 3131
 kinit2 3139
 0319 1334 3139
 KSTACKSIZE 0151
 0151 1158 1167 1385 1873
 2498
 kvmalloc 1840

 0427 1320 1840
 lapiceoi 7610
 0328 3471 3475 3482 3486
 3492 7610
 lapicid 7601
 0326 2444 7601 8213
 lapicinit 7551
 0329 1322 1345 7551
 lapicstartap 7629
 0330 1389 7629
 lapicw 7496
 7496 7557 7563 7564 7565
 7568 7569 7574 7577 7580
 7581 7584 7587 7588 7593
 7613 7645 7646 7648 7656
 7657
 lcr3 0590
 0590 1855 1878
 lgdt 0512
 0512 0520 1235 1728 9291
 lidt 0526
 0526 0534 3431
 LINT0 7485
 7485 7568
 LINT1 7486
 7486 7569
 LIST 8710
 8710 8790 8957 9233
 listcmd 8740 8951
 8740 8761 8791 8951 8953
 9096 9207 9234
 loaduvm 1903
 0434 1903 1909 1912 6755
 log 4838 4850
 4838 4850 4862 4864 4865
 4866 4876 4877 4878 4890
 4893 4894 4895 4906 4909
 4910 4911 4922 4930 4932
 4933 4934 4936 4938 4939
 4957 4958 4959 4960 4961
 4963 4968 4970 4976 4977
 4978 4979 4989 4990 4991
 5003 5007 5026 5028 5031
 5032 5033 5036 5037 5038
 5040
 logheader 4833
 4833 4845 4858 4859 4891
 4907
 LOGSIZE 0160
 0160 4835 4934 5026 6117

log_write 5022
 0335 5022 5029 5108 5128
 5164 5316 5343 5530 5673
 ltr 0538
 0538 0540 1877
 mappages 1760
 1760 1829 1893 1946 2054
 MAXARG 0158
 0158 6628 6714 6772
 MAXARGS 8713
 8713 8721 8722 9190
 MAXFILE 4125
 4125 5666
 MAXOPBLOCKS 0159
 0159 0160 0161 4934
 memcmp 7015
 0395 7015 7337 7388 7716
 memmove 7031
 0396 1375 1894 2053 2132
 4879 4992 5086 5342 5421
 5622 5672 5829 5831 7031
 7054 8276
 memset 7004
 0397 1747 1825 1892 1945
 2511 2532 3173 5107 5314
 6434 6635 7004 8278 8837
 8908 8919 8935 8956 8969
 microdelay 7619
 0331 7619 7647 7649 7658
 7676 8508
 min 5073
 5073 5621 5671
 MINS 7667
 7667 7684
 MONTH 7670
 7670 7687
 mp 7152
 7152 7308 7329 7336 7337
 7338 7355 7360 7364 7365
 7368 7369 7380 7383 7385
 7387 7394 7405 7410 7442
 MPBUS 7202
 7202 7429
 mpconf 7163
 7163 7379 7382 7387 7406
 mpconfig 7380
 7380 7410
 mpenter 1341
 1341 1386
 mpinit 7401

 0341 1321 7401
 mpioapic 7189
 7189 7408 7425 7427
 MPIOAPIC 7203
 7203 7424
 MPIOINTR 7204
 7204 7430
 MPLINTR 7205
 7205 7431
 mpmain 1352
 1309 1336 1346 1352
 mpproc 7178
 7178 7407 7417 7422
 MPPROC 7201
 7201 7416
 mpsearch 7356
 7356 7385
 mpsearchl 7330
 7330 7364 7368 7371
 multiboot_header 1129
 1128 1129
 mycpu 2437
 0363 1356 1378 1590 1647
 1661 1662 1663 1671 1673
 1870 1871 1872 1873 1876
 2431 2437 2442 2461 2761
 2815 2821 2822 2823
 myproc 2457
 0364 2457 2561 2584 2629
 2675 2811 2831 2876 3454
 3456 3458 3501 3510 3512
 3518 3523 3527 3619 3634
 3654 3664 3754 3841 3858
 3875 4729 5862 6178 6206
 6270 6605 6664 6719 6937
 6957 8381
 namecmp 5703
 0298 5703 5724 6420
 namei 5890
 0299 2542 5890 6311 6520
 6608 6723
 nameiparent 5901
 0300 5855 5870 5882 5901
 6327 6412 6463
 namex 5855
 5855 5893 5903
 NBUF 0161
 0161 4530 4553
 ncpu 7314
 1377 2313 2447 4356 7314

7418 7419 7420
 NCPU 0152 7675 8260 8262 8281 8282
 0152 2312 7313 7418 8283 8284 8474 8477 8478
 NDEV 0156 8479 8480 8481 8482 8509
 0156 5609 5659 5962 9278 9286 9414 9415 9416
 NDIRECT 4123 9417 9418 9419
 4123 4125 4134 4224 5515 outsl 0483
 5520 5524 5525 5562 5569 0483 0485 4396
 5570 5577 5578 outw 0477
 NELEM 0442 0477 1280 1282 9324 9326
 0442 1828 3022 3757 6637 O_WRONLY 4001
 nextpid 2416 4001 6545 6546 9128 9131
 2416 2489 P2V 0211
 NFILE 0154 0211 1319 1334 1374 1742
 0154 5965 5981 1826 1918 1978 2012 2053
 NINDIRECT 4124 2111 7334 7362 7387 7639
 4124 4125 5522 5572 8252
 NINODE 0155 panic 8205 8881
 0155 5239 5248 5362 0272 1578 1605 1670 1672
 NO 7856 1771 1827 1863 1865 1867
 7856 7902 7905 7907 7908 1890 1909 1912 1977 2008
 7909 7910 7912 7924 7927 2028 2046 2048 2442 2451
 7929 7930 7931 7932 7934 2529 2634 2665 2814 2816
 7952 7953 7955 7956 7957 2818 2820 2879 2882 3170
 7958 3505 4377 4379 4385 4459
 NOFILE 0153 4461 4463 4596 4618 4629
 0153 2348 2605 2637 6178 4859 4960 5027 5029 5136
 6208 5162 5322 5374 5409 5425
 NPENTRIES 0870 5434 5536 5717 5721 5767
 0870 1406 2010 5775 6006 6020 6080 6134
 NPROC 0150 6139 6359 6428 6436 6477
 0150 2411 2480 2654 2681 6490 6494 7411 7440 8163
 2770 2957 2980 3019 8205 8213 8273 8751 8770
 NSEGS 0749 8803 8881 8894 9078 9122
 0749 2305 9156 9160 9186 9191
 nulterminate 9202 panicked 8069
 9065 9080 9202 9223 9229 8069 8219 8303
 9230 9235 9236 9241 parseblock 9151
 NUMLOCK 7863 9151 9156 9175
 7863 7896 parsecmd 9068
 O_CREATE 4003 8752 8874 9068
 4003 6513 9128 9131 parseexec 9167
 O_RDONLY 4000 9064 9105 9167
 4000 6525 9125 parseline 9085
 O_RDWR 4002 9062 9074 9085 9096 9158
 4002 6546 8664 8666 8857 parsepipe 9101
 outb 0471 9063 9089 9101 9108
 0471 4360 4369 4388 4389 parseredirs 9114
 4390 4391 4392 4393 4395 9114 9162 9181 9192
 4398 7445 7446 7637 7638 PCINT 7484
 7484 7574

pde_t 0103 PIPESIZE 6860
 0103 0428 0429 0430 0431 6860 6864 6936 6944 6966
 0432 0433 0434 0435 0438 pipewrite 6930
 0439 1310 1360 1406 1710 0354 6109 6930
 1735 1737 1760 1817 1820 popcli 1667
 1823 1885 1903 1927 1961 0386 1622 1667 1670 1672
 2003 2022 2034 2035 2037 1879 2463
 2102 2118 2339 6718 printint 8077
 PDX 0861 8077 8176 8180
 0861 1740 1973 proc 2337
 PDXSHIFT 0876 0255 0364 0369 0436 1305
 0861 0867 0876 1410 1558 1706 1860 2309 2337
 peek 9051 2343 2406 2411 2414 2456
 9051 9075 9090 9094 9106 2459 2462 2472 2475 2480
 9119 9155 9159 9174 9182 2522 2561 2583 2584 2629
 PGADDR 0867 2630 2654 2673 2675 2681
 0867 1973 2760 2762 2770 2777 2786
 PGROUNDDOWN 0879 2811 2876 2955 2957 2977
 0879 1765 1766 2125 2980 3015 3019 3405 3509
 PGROUNDUP 0878 3605 3619 3634 3664 3754
 0878 1937 1969 3154 6764 3807 4307 4708 5066 6161
 PGSIZE 0872 6206 6605 6704 6719 6854
 0872 0878 0879 1405 1747 7311 7407 7417 7419 8064
 1775 1776 1825 1889 1892 8461
 1893 1908 1910 1914 1917 procdump 3004
 1938 1945 1946 1970 1973 0366 3004 8366
 2044 2053 2054 2129 2135 proghdr 1024
 2531 2538 3155 3169 3173 1024 6717 9370 9384
 6753 6765 6767 PTE_ADDR 0893
 PHYSTOP 0203 0893 1742 1913 1975 2012
 0203 1334 1812 1826 1827 2049 2111
 3169 PTE_FLAGS 0894
 pinit 2423 0894 2050
 0365 1328 2423 PTE_P 0882
 pipe 6862 0882 1408 1410 1741 1751
 0254 0352 0353 0354 4205 1770 1772 1974 2011 2047
 6031 6072 6109 6862 6874 2107
 6880 6886 6890 6894 6911 PTE_PS 0889
 6930 6951 8613 8802 8803 0889 1408 1410
 PIPE 8709 pte_t 0897
 8709 8800 8936 9227 0897 1734 1738 1742 1744
 pipealloc 6872 1763 1906 1963 2024 2038
 0351 6659 6872 2104
 pipeclose 6911 PTE_U 0884
 0352 6031 6911 0884 1751 1893 1946 2029
 pipecmd 8734 8930 2109
 8734 8762 8801 8930 8932 PTE_W 0883
 9108 9208 9228 0883 1408 1410 1751 1810
 piperead 6951 1812 1813 1893 1946
 0353 6072 6951 PTX 0864

0864 1753
 PTXSHIFT 0875
 0864 0867 0875
 pushcli 1655
 0385 1576 1655 1869 2460
 rcr2 0582
 0582 3504 3511
 readeflags 0544
 0544 1659 1669 2441 2819
 read_head 4888
 4888 4920
 readi 5603
 0301 1918 5603 5720 5766
 6075 6358 6359 6732 6743
 readsb 5081
 0287 4863 5081 5157 5252
 readsect 9410
 9410 9445
 readseg 9429
 9364 9377 9388 9429
 recover_from_log 4918
 4852 4867 4918
 REDIR 8708
 8708 8780 8920 9221
 redircmd 8725 8914
 8725 8763 8781 8914 8916
 9125 9128 9131 9209 9222
 REG_ID 7760
 7760 7807
 REG_TABLE 7762
 7762 7814 7815 7825 7826
 REG_VER 7761
 7761 7806
 release 1602
 0384 1602 1605 2484 2491
 2552 2618 2696 2702 2788
 2833 2857 2892 2905 2968
 2986 2990 3181 3198 3469
 3876 3881 3894 4412 4430
 4483 4576 4591 4645 4730
 4740 4757 4939 4970 4979
 5040 5365 5381 5393 5464
 5477 5984 5988 6008 6022
 6028 6922 6925 6938 6947
 6958 6969 8201 8364 8382
 8402 8417
 releasesleep 4734
 0390 4631 4734 5436 5473
 ROOTDEV 0157
 0157 2864 2865 5860

ROOTINO 4104
 4104 5860
 run 3115
 3011 3115 3116 3122 3167
 3177 3190 7411
 runcmd 8756
 8756 8770 8787 8793 8795
 8809 8816 8827 8874
 RUNNING 2334
 2334 2779 2817 3011 3523
 safestrncpy 7082
 0398 2541 2610 6793 7082
 sb 5077
 0287 4154 4160 4861 4863
 4864 4865 5077 5081 5086
 5122 5123 5124 5157 5158
 5252 5253 5254 5255 5256
 5310 5311 5335 5414 7704
 7706 7708
 sched 2808
 0368 2664 2808 2814 2816
 2818 2820 2832 2898
 scheduler 2758
 0367 1357 2303 2758 2781
 2822
 SCROLLLOCK 7864
 7864 7897
 SECS 7666
 7666 7683
 SECTOR_SIZE 4315
 4315 4380
 SECTSIZE 9362
 9362 9423 9436 9439 9444
 SEG 0818
 0818 1724 1725 1726 1727
 SEG16 0822
 0822 1870
 SEG_ASM 0660
 0660 1289 1290 9334 9335
 segdesc 0801
 0509 0512 0801 0818 0822
 2305
 seginit 1715
 0426 1323 1344 1715
 SEG_KCODE 0742
 0742 1243 1724 3422 3423
 9303
 SEG_KDATA 0743
 0743 1253 1725 1872 3363
 9308

SEG_NULLASM 0654
 0654 1288 9333
 SEG_TSS 0746
 0746 1870 1871 1877
 SEG_UCODE 0744
 0744 1726 2533
 SEG_UDATA 0745
 0745 1727 2534
 SETGATE 0971
 0971 3422 3423
 setupkvm 1818
 0428 1818 1842 2042 2528
 6737
 SHIFT 7858
 7858 7886 7887 8035
 skipelem 5815
 5815 5864
 sleep 2874
 0370 2707 2874 2879 2882
 3009 3879 4479 4715 4726
 4933 4936 6942 6961 8386
 8629
 sleeplock 3951
 0258 0389 0390 0391 0392
 3904 3951 4216 4311 4524
 4710 4713 4722 4734 4751
 4804 5068 5959 6164 6857
 8059 8457
 spinlock 1501
 0257 0370 0380 0382 0383
 0384 0418 1501 1559 1562
 1574 1602 1645 2407 2410
 2874 3109 3120 3408 3413
 3953 4310 4330 4523 4529
 4709 4803 4839 5067 5238
 5958 5964 6163 6856 6863
 8058 8072 8456
 STA_R 0669 0835
 0669 0835 1289 1724 1726
 9334
 start 1223 8559 9261
 1222 1223 1266 1274 1276
 4840 4864 4877 4890 4906
 4990 5254 8558 8559 9260
 9261 9317
 startothers 1364
 1308 1333 1364
 stat 4054
 0259 0283 0302 4054 5064
 5588 6052 6159 6279 8653

stati 5588
 0302 5588 6056
 STA_W 0668 0834
 0668 0834 1290 1725 1727
 9335
 STA_X 0665 0831
 0665 0831 1289 1724 1726
 9334
 sti 0563
 0563 0565 1674 2766
 stosb 0492
 0492 0494 7010 9390
 stosl 0501
 0501 0503 7008
 strlen 7101
 0399 6774 6775 7101 8868
 9073
 strncmp 7058
 0400 5705 7058
 strncpy 7068
 0401 5772 7068
 STS_IG32 0849
 0849 0977
 STS_T32A 0846
 0846 1870
 STS_TG32 0850
 0850 0977
 sum 7318
 7318 7320 7322 7324 7325
 7337 7392
 superbblock 4113
 0260 0287 4113 4861 5077
 5081
 SVR 7467
 7467 7557
 switchkvm 1853
 0437 1343 1843 1853 2782
 switchvum 1860
 0436 1860 1863 1865 1867
 2572 2778 6801
 swtch 3058
 0377 2781 2822 3057 3058
 syscall 3751
 0409 3457 3607 3751
 SYSCALL 8603 8610 8611 8612 8613 86
 8610 8611 8612 8613 8614
 8615 8616 8617 8618 8619
 8620 8621 8622 8623 8624
 8625 8626 8627 8628 8629
 8630

```

sys_chdir 6601
  3700 3731 6601
SYS_chdir 3559
  3559 3731
sys_close 6263
  3701 3743 6263
SYS_close 3571
  3571 3743
sys_dup 6218
  3702 3732 6218
SYS_dup 3560
  3560 3732
sys_exec 6626
  3703 3729 6626
SYS_exec 3557
  3557 3729 8563
sys_exit 3816
  3704 3724 3816
SYS_exit 3552
  3552 3724 8568
sys_fork 3810
  3705 3723 3810
SYS_fork 3551
  3551 3723
sys_fstat 6276
  3706 3730 6276
SYS_fstat 3558
  3558 3730
sys_getpid 3839
  3707 3733 3839
SYS_getpid 3561
  3561 3733
sys_kill 3829
  3708 3728 3829
SYS_kill 3556
  3556 3728
sys_link 6302
  3709 3741 6302
SYS_link 3569
  3569 3741
sys_mkdir 6551
  3710 3742 6551
SYS_mkdir 3570
  3570 3742
sys_mknod 6567
  3711 3739 6567
SYS_mknod 3567
  3567 3739
sys_open 6501
  3712 3737 6501
SYS_open 3565
  3565 3737
sys_pipe 6651
  3713 3726 6651
SYS_pipe 3554
  3554 3726
sys_read 6232
  3714 3727 6232
SYS_read 3555
  3555 3727
sys_sbrk 3851
  3715 3734 3851
SYS_sbrk 3562
  3562 3734
sys_sleep 3865
  3716 3735 3865
SYS_sleep 3563
  3563 3735
sys_unlink 6401
  3717 3740 6401
SYS_unlink 3568
  3568 3740
sys_uptime 3888
  3720 3736 3888
SYS_uptime 3564
  3564 3736
sys_wait 3823
  3718 3725 3823
SYS_wait 3553
  3553 3725
sys_write 6251
  3719 3738 6251
SYS_write 3566
  3566 3738
taskstate 0901
  0901 2304
TDCR 7491
  7491 7563
T_DEV 4052
  4052 5608 5658 6577
T_DIR 4050
  4050 5716 5866 6317 6429
  6437 6485 6525 6557 6613
T_FILE 4051
  4051 6470 6514
ticks 3414
  0416 3414 3467 3468 3873
  3874 3879 3893
tickslock 3413
  0418 3413 3425 3466 3469

```

```

  3872 3876 3879 3881 3892
  3894
TICR 7489
  7489 7565
TIMER 7481
  7481 7564
T_IRQ0 3279
  3279 3464 3473 3477 3480
  3484 3488 3489 3523 7557
  7564 7577 7814 7825
TPR 7465
  7465 7593
trap 3451
  3302 3304 3369 3451 3503
  3505 3508
trapframe 0602
  0602 2344 2502 3451
trapret 3374
  2418 2507 3373 3374
T_SYSCALL 3276
  3276 3423 3453 8564 8569
  8607
tvinit 3417
  0417 1329 3417
uart 8466
  8466 8487 8505 8515
uartgetc 8513
  8513 8525
uartinit 8469
  0421 1327 8469
uartintr 8523
  0422 3485 8523
uartputc 8501
  0423 8310 8312 8497 8501
userinit 2520
  0371 1335 2520 2529
uva2ka 2102
  0429 2102 2126
V2P 0210
  0210 1387 1389 1751 1811
  1812 1855 1878 1893 1946
  2054 3169
V2P_WO 0213
  0213 1140 1150
VER 7464
  7464 7573
wait 2671
  0372 2671 3825 8612 8683
  8794 8820 8821 8875
waitdisk 9401
  9401 9413 9422
wakeup 2964
  0373 2964 3468 4424 4739
  4968 4978 6916 6919 6941
  6946 6968 8358
wakeupl 2953
  2420 2651 2658 2953 2967
walkpgdir 1735
  1735 1768 1911 1971 2026
  2045 2106
write_head 4904
  4904 4923 5005 5008
writei 5653
  0303 5653 5774 6126 6435
  6436
write_log 4985
  4985 5004
xchg 0569
  0569 1356 1581
YEAR 7671
  7671 7688
yield 2828
  0374 2828 3524

```



```
0100 typedef unsigned int    uint;
0101 typedef unsigned short  ushort;
0102 typedef unsigned char   uchar;
0103 typedef uint pde_t;
0104
0105
0106
0107
0108
0109
0110
0111
0112
0113
0114
0115
0116
0117
0118
0119
0120
0121
0122
0123
0124
0125
0126
0127
0128
0129
0130
0131
0132
0133
0134
0135
0136
0137
0138
0139
0140
0141
0142
0143
0144
0145
0146
0147
0148
0149
```

```

0150 #define NPROC      64 // maximum number of processes
0151 #define KSTACKSIZE 4096 // size of per-process kernel stack
0152 #define NCPU       8 // maximum number of CPUs
0153 #define NOFILE     16 // open files per process
0154 #define NFILE      100 // open files per system
0155 #define NINODE      50 // maximum number of active i-nodes
0156 #define NDEV       10 // maximum major device number
0157 #define ROOTDEV     1 // device number of file system root disk
0158 #define MAXARG      32 // max exec arguments
0159 #define MAXOPBLOCKS 10 // max # of blocks any FS op writes
0160 #define LOGSIZE      (MAXOPBLOCKS*3) // max data blocks in on-disk log
0161 #define NBUF        (MAXOPBLOCKS*3) // size of disk block cache
0162 #define FSSIZE      1000 // size of file system in blocks
0163
0164
0165
0166
0167
0168
0169
0170
0171
0172
0173
0174
0175
0176
0177
0178
0179
0180
0181
0182
0183
0184
0185
0186
0187
0188
0189
0190
0191
0192
0193
0194
0195
0196
0197
0198
0199

```

```

0200 // Memory layout
0201
0202 #define EXTMEM  0x100000 // Start of extended memory
0203 #define PHYSTOP 0xE000000 // Top physical memory
0204 #define DEVSPACE 0xFE000000 // Other devices are at high addresses
0205
0206 // Key addresses for address space layout (see kmap in vm.c for layout)
0207 #define KERNBASE 0x80000000 // First kernel virtual address
0208 #define KERNLINK (KERNBASE+EXTMEM) // Address where kernel is linked
0209
0210 #define V2P(a) (((uint) (a)) - KERNBASE)
0211 #define P2V(a) (((void *) (a)) + KERNBASE)
0212
0213 #define V2P_WO(x) ((x) - KERNBASE) // same as V2P, but without casts
0214 #define P2V_WO(x) ((x) + KERNBASE) // same as P2V, but without casts
0215
0216
0217
0218
0219
0220
0221
0222
0223
0224
0225
0226
0227
0228
0229
0230
0231
0232
0233
0234
0235
0236
0237
0238
0239
0240
0241
0242
0243
0244
0245
0246
0247
0248
0249

```

```

0250 struct buf;
0251 struct context;
0252 struct file;
0253 struct inode;
0254 struct pipe;
0255 struct proc;
0256 struct rtcdate;
0257 struct spinlock;
0258 struct sleeplock;
0259 struct stat;
0260 struct superblock;
0261
0262 // bio.c
0263 void      binit(void);
0264 struct buf* bread(uint, uint);
0265 void      brelse(struct buf*);
0266 void      bwrite(struct buf*);
0267
0268 // console.c
0269 void      consoleinit(void);
0270 void      cprintf(char*, ...);
0271 void      consoleintr(int (*)(void));
0272 void      panic(char*) __attribute__((noreturn));
0273
0274 // exec.c
0275 int      exec(char*, char**);
0276
0277 // file.c
0278 struct file* filealloc(void);
0279 void      fileclose(struct file*);
0280 struct file* filedup(struct file*);
0281 void      fileinit(void);
0282 int      fileread(struct file*, char*, int n);
0283 int      filestat(struct file*, struct stat*);
0284 int      filewrite(struct file*, char*, int n);
0285
0286 // fs.c
0287 void      readsb(int dev, struct superblock *sb);
0288 int      dirlink(struct inode*, char*, uint);
0289 struct inode* dirlookup(struct inode*, char*, uint*);
0290 struct inode* ialloc(uint, short);
0291 struct inode* idup(struct inode*);
0292 void      iinit(int dev);
0293 void      ilock(struct inode*);
0294 void      iput(struct inode*);
0295 void      iunlock(struct inode*);
0296 void      iunlockput(struct inode*);
0297 void      iupdate(struct inode*);
0298 int      namecmp(const char*, const char*);
0299 struct inode* namei(char*);

```

```

0300 struct inode* nameiparent(char*, char*);
0301 int      readi(struct inode*, char*, uint, uint);
0302 void      stati(struct inode*, struct stat*);
0303 int      writei(struct inode*, char*, uint, uint);
0304
0305 // ide.c
0306 void      ideinit(void);
0307 void      ideintr(void);
0308 void      iderw(struct buf*);
0309
0310 // ioapic.c
0311 void      ioapicenable(int irq, int cpu);
0312 extern uchar ioapicid;
0313 void      ioapicinit(void);
0314
0315 // kalloc.c
0316 char*      kalloc(void);
0317 void      kfree(char*);
0318 void      kinit1(void*, void*);
0319 void      kinit2(void*, void*);
0320
0321 // kbd.c
0322 void      kbdintr(void);
0323
0324 // lapic.c
0325 void      cmostime(struct rtcdate *r);
0326 int      lapicid(void);
0327 extern volatile uint* lapic;
0328 void      lapiceoi(void);
0329 void      lapicinit(void);
0330 void      lapicstartap(uchar, uint);
0331 void      microdelay(int);
0332
0333 // log.c
0334 void      initlog(int dev);
0335 void      log_write(struct buf*);
0336 void      begin_op();
0337 void      end_op();
0338
0339 // mp.c
0340 extern int ismp;
0341 void      mpinit(void);
0342
0343 // picirq.c
0344 void      picenable(int);
0345 void      picinit(void);
0346
0347
0348
0349

```

```

0350 // pipe.c
0351 int      pipealloc(struct file**, struct file**);
0352 void     pipeclose(struct pipe*, int);
0353 int      piperead(struct pipe*, char*, int);
0354 int      pipewrite(struct pipe*, char*, int);
0355
0356
0357 // proc.c
0358 int      cpuid(void);
0359 void     exit(void);
0360 int      fork(void);
0361 int      growproc(int);
0362 int      kill(int);
0363 struct cpu* mycpu(void);
0364 struct proc* myproc();
0365 void     pinit(void);
0366 void     procdump(void);
0367 void     scheduler(void) __attribute__((noreturn));
0368 void     sched(void);
0369 void     setproc(struct proc*);
0370 void     sleep(void*, struct spinlock*);
0371 void     userinit(void);
0372 int      wait(void);
0373 void     wakeup(void*);
0374 void     yield(void);
0375
0376 // swtch.S
0377 void     swtch(struct context**, struct context*);
0378
0379 // spinlock.c
0380 void     acquire(struct spinlock*);
0381 void     getcallerpcs(void*, uint*);
0382 int      holding(struct spinlock*);
0383 void     initlock(struct spinlock*, char*);
0384 void     release(struct spinlock*);
0385 void     pushcli(void);
0386 void     popcli(void);
0387
0388 // sleeplock.c
0389 void     acquiresleep(struct sleeplock*);
0390 void     releasesleep(struct sleeplock*);
0391 int      holdingsleep(struct sleeplock*);
0392 void     initsleeplock(struct sleeplock*, char*);
0393
0394 // string.c
0395 int      memcmp(const void*, const void*, uint);
0396 void*    memmove(void*, const void*, uint);
0397 void*    memset(void*, int, uint);
0398 char*    safestrcpy(char*, const char*, int);
0399 int      strlen(const char*);

```

```

0400 int      strncmp(const char*, const char*, uint);
0401 char*    strncpy(char*, const char*, int);
0402
0403 // syscall.c
0404 int      argint(int, int*);
0405 int      argptr(int, char**, int);
0406 int      argstr(int, char**);
0407 int      fetchint(uint, int*);
0408 int      fetchstr(uint, char**);
0409 void     syscall(void);
0410
0411 // timer.c
0412 void     timerinit(void);
0413
0414 // trap.c
0415 void     idtinit(void);
0416 extern uint ticks;
0417 void     tvinit(void);
0418 extern struct spinlock tickslock;
0419
0420 // uart.c
0421 void     uartinit(void);
0422 void     uartintr(void);
0423 void     uartputc(int);
0424
0425 // vm.c
0426 void     seginit(void);
0427 void     kvmalloc(void);
0428 pde_t*   setupkvm(void);
0429 char*    uva2ka(pde_t*, char*);
0430 int      allocvm(pde_t*, uint, uint);
0431 int      deallocvm(pde_t*, uint, uint);
0432 void     freevm(pde_t*);
0433 void     initvm(pde_t*, char*, uint);
0434 int      loadvm(pde_t*, char*, struct inode*, uint, uint);
0435 pde_t*   copyvm(pde_t*, uint);
0436 void     switchvm(struct proc*);
0437 void     switchkvm(void);
0438 int      copyout(pde_t*, uint, void*, uint);
0439 void     clearpteu(pde_t *pgdir, char *uva);
0440
0441 // number of elements in fixed-size array
0442 #define NELEM(x) (sizeof(x)/sizeof((x)[0]))
0443
0444
0445
0446
0447
0448
0449

```

```

0450 // Routines to let C code use special x86 instructions.
0451
0452 static inline uchar
0453 inb(ushort port)
0454 {
0455     uchar data;
0456
0457     asm volatile("in %1,%0 : "=a" (data) : "d" (port));
0458     return data;
0459 }
0460
0461 static inline void
0462 insl(int port, void *addr, int cnt)
0463 {
0464     asm volatile("cld; rep insl" :
0465                 "=D" (addr), "=c" (cnt) :
0466                 "d" (port), "0" (addr), "1" (cnt) :
0467                 "memory", "cc");
0468 }
0469
0470 static inline void
0471 outb(ushort port, uchar data)
0472 {
0473     asm volatile("out %0,%1" : : "a" (data), "d" (port));
0474 }
0475
0476 static inline void
0477 outw(ushort port, ushort data)
0478 {
0479     asm volatile("out %0,%1" : : "a" (data), "d" (port));
0480 }
0481
0482 static inline void
0483 outsl(int port, const void *addr, int cnt)
0484 {
0485     asm volatile("cld; rep outsl" :
0486                 "=S" (addr), "=c" (cnt) :
0487                 "d" (port), "0" (addr), "1" (cnt) :
0488                 "cc");
0489 }
0490
0491 static inline void
0492 stosb(void *addr, int data, int cnt)
0493 {
0494     asm volatile("cld; rep stosb" :
0495                 "=D" (addr), "=c" (cnt) :
0496                 "0" (addr), "1" (cnt), "a" (data) :
0497                 "memory", "cc");
0498 }
0499

```

```

0500 static inline void
0501 stosl(void *addr, int data, int cnt)
0502 {
0503     asm volatile("cld; rep stosl" :
0504                 "=D" (addr), "=c" (cnt) :
0505                 "0" (addr), "1" (cnt), "a" (data) :
0506                 "memory", "cc");
0507 }
0508
0509 struct segdesc;
0510
0511 static inline void
0512 lgdt(struct segdesc *p, int size)
0513 {
0514     volatile ushort pd[3];
0515
0516     pd[0] = size-1;
0517     pd[1] = (uint)p;
0518     pd[2] = (uint)p >> 16;
0519
0520     asm volatile("lgdt (%0)" : : "r" (pd));
0521 }
0522
0523 struct gatedesc;
0524
0525 static inline void
0526 lidt(struct gatedesc *p, int size)
0527 {
0528     volatile ushort pd[3];
0529
0530     pd[0] = size-1;
0531     pd[1] = (uint)p;
0532     pd[2] = (uint)p >> 16;
0533
0534     asm volatile("lidt (%0)" : : "r" (pd));
0535 }
0536
0537 static inline void
0538 ltr(ushort sel)
0539 {
0540     asm volatile("ltr %0" : : "r" (sel));
0541 }
0542
0543 static inline uint
0544 readeflags(void)
0545 {
0546     uint eflags;
0547     asm volatile("pushfl; popl %0" : "=r" (eflags));
0548     return eflags;
0549 }

```

```

0550 static inline void
0551 loadgs(ushort v)
0552 {
0553     asm volatile("movw %0, %%gs" : : "r" (v));
0554 }
0555
0556 static inline void
0557 cli(void)
0558 {
0559     asm volatile("cli");
0560 }
0561
0562 static inline void
0563 sti(void)
0564 {
0565     asm volatile("sti");
0566 }
0567
0568 static inline uint
0569 xchg(volatile uint *addr, uint newval)
0570 {
0571     uint result;
0572
0573     // The + in "+m" denotes a read-modify-write operand.
0574     asm volatile("lock; xchgl %0, %1" :
0575                 "+m" (*addr), "=a" (result) :
0576                 "l" (newval) :
0577                 "cc");
0578     return result;
0579 }
0580
0581 static inline uint
0582 rcr2(void)
0583 {
0584     uint val;
0585     asm volatile("movl %%cr2,%0" : "=r" (val));
0586     return val;
0587 }
0588
0589 static inline void
0590 lcr3(uint val)
0591 {
0592     asm volatile("movl %0,%%cr3" : : "r" (val));
0593 }
0594
0595
0596
0597
0598
0599

```

```

0600 // Layout of the trap frame built on the stack by the
0601 // hardware and by trapasm.S, and passed to trap().
0602 struct trapframe {
0603     // registers as pushed by pusha
0604     uint edi;
0605     uint esi;
0606     uint ebp;
0607     uint oesp;    // useless & ignored
0608     uint ebx;
0609     uint edx;
0610     uint ecx;
0611     uint eax;
0612
0613     // rest of trap frame
0614     ushort gs;
0615     ushort padding1;
0616     ushort fs;
0617     ushort padding2;
0618     ushort es;
0619     ushort padding3;
0620     ushort ds;
0621     ushort padding4;
0622     uint trapno;
0623
0624     // below here defined by x86 hardware
0625     uint err;
0626     uint eip;
0627     ushort cs;
0628     ushort padding5;
0629     uint eflags;
0630
0631     // below here only when crossing rings, such as from user to kernel
0632     uint esp;
0633     ushort ss;
0634     ushort padding6;
0635 };
0636
0637
0638
0639
0640
0641
0642
0643
0644
0645
0646
0647
0648
0649

```

```

0650 //
0651 // assembler macros to create x86 segments
0652 //
0653
0654 #define SEG_NULLASM          \
0655     .word 0, 0;              \
0656     .byte 0, 0, 0, 0
0657
0658 // The 0xC0 means the limit is in 4096-byte units
0659 // and (for executable segments) 32-bit mode.
0660 #define SEG_ASM(type,base,lim) \
0661     .word (((lim) >> 12) & 0xffff), ((base) & 0xffff); \
0662     .byte (((base) >> 16) & 0xff), (0x90 | (type)), \
0663     (0xC0 | (((lim) >> 28) & 0xf)), (((base) >> 24) & 0xff)
0664
0665 #define STA_X    0x8    // Executable segment
0666 #define STA_E    0x4    // Expand down (non-executable segments)
0667 #define STA_C    0x4    // Conforming code segment (executable only)
0668 #define STA_W    0x2    // Writeable (non-executable segments)
0669 #define STA_R    0x2    // Readable (executable segments)
0670 #define STA_A    0x1    // Accessed
0671
0672
0673
0674
0675
0676
0677
0678
0679
0680
0681
0682
0683
0684
0685
0686
0687
0688
0689
0690
0691
0692
0693
0694
0695
0696
0697
0698
0699

```

```

0700 // This file contains definitions for the
0701 // x86 memory management unit (MMU).
0702
0703 // Eflags register
0704 #define FL_CF        0x00000001    // Carry Flag
0705 #define FL_PF        0x00000004    // Parity Flag
0706 #define FL_AF        0x00000010    // Auxiliary carry Flag
0707 #define FL_ZF        0x00000040    // Zero Flag
0708 #define FL_SF        0x00000080    // Sign Flag
0709 #define FL_TF        0x00000100    // Trap Flag
0710 #define FL_IF        0x00000200    // Interrupt Enable
0711 #define FL_DF        0x00000400    // Direction Flag
0712 #define FL_OF        0x00000800    // Overflow Flag
0713 #define FL_IOPL_MASK 0x00003000    // I/O Privilege Level bitmask
0714 #define FL_IOPL_0    0x00000000    // IOPL == 0
0715 #define FL_IOPL_1    0x00001000    // IOPL == 1
0716 #define FL_IOPL_2    0x00002000    // IOPL == 2
0717 #define FL_IOPL_3    0x00003000    // IOPL == 3
0718 #define FL_NT        0x00004000    // Nested Task
0719 #define FL_RF        0x00010000    // Resume Flag
0720 #define FL_VM        0x00020000    // Virtual 8086 mode
0721 #define FL_AC        0x00040000    // Alignment Check
0722 #define FL_VIF       0x00080000    // Virtual Interrupt Flag
0723 #define FL_VIP       0x00100000    // Virtual Interrupt Pending
0724 #define FL_ID        0x00200000    // ID flag
0725
0726 // Control Register flags
0727 #define CR0_PE        0x00000001    // Protection Enable
0728 #define CR0_MP        0x00000002    // Monitor coProcessor
0729 #define CR0_EM        0x00000004    // Emulation
0730 #define CR0_TS        0x00000008    // Task Switched
0731 #define CR0_ET        0x00000010    // Extension Type
0732 #define CR0_NE        0x00000020    // Numeric Error
0733 #define CR0_WP        0x00010000    // Write Protect
0734 #define CR0_AM        0x00040000    // Alignment Mask
0735 #define CR0_NW        0x00200000    // Not Writethrough
0736 #define CR0_CD        0x00400000    // Cache Disable
0737 #define CR0_PG        0x00800000    // Paging
0738
0739 #define CR4_PSE       0x00000010    // Page size extension
0740
0741 // various segment selectors.
0742 #define SEG_KCODE 1 // kernel code
0743 #define SEG_KDATA 2 // kernel data+stack
0744 #define SEG_UCODE 3 // user code
0745 #define SEG_UDATA 4 // user data+stack
0746 #define SEG_TSS   5 // this process's task state
0747
0748 // cpu->gdt[NSEGS] holds the above segments.
0749 #define NSEGS     6

```

```

0750
0751 #ifndef __ASSEMBLER__
0752
0753
0754
0755
0756
0757
0758
0759
0760
0761
0762
0763
0764
0765
0766
0767
0768
0769
0770
0771
0772
0773
0774
0775
0776
0777
0778
0779
0780
0781
0782
0783
0784
0785
0786
0787
0788
0789
0790
0791
0792
0793
0794
0795
0796
0797
0798
0799

```

```

0800 // Segment Descriptor
0801 struct segdesc {
0802     uint lim_15_0 : 16; // Low bits of segment limit
0803     uint base_15_0 : 16; // Low bits of segment base address
0804     uint base_23_16 : 8; // Middle bits of segment base address
0805     uint type : 4; // Segment type (see STS_ constants)
0806     uint s : 1; // 0 = system, 1 = application
0807     uint dpl : 2; // Descriptor Privilege Level
0808     uint p : 1; // Present
0809     uint lim_19_16 : 4; // High bits of segment limit
0810     uint avl : 1; // Unused (available for software use)
0811     uint rsv1 : 1; // Reserved
0812     uint db : 1; // 0 = 16-bit segment, 1 = 32-bit segment
0813     uint g : 1; // Granularity: limit scaled by 4K when set
0814     uint base_31_24 : 8; // High bits of segment base address
0815 };
0816
0817 // Normal segment
0818 #define SEG(type, base, lim, dpl) (struct segdesc) \
0819 { ((lim) >> 12) & 0xffff, (uint)(base) & 0xffff, \
0820   ((uint)(base) >> 16) & 0xff, type, 1, dpl, 1, \
0821   (uint)(lim) >> 28, 0, 0, 1, 1, (uint)(base) >> 24 }
0822 #define SEG16(type, base, lim, dpl) (struct segdesc) \
0823 { (lim) & 0xffff, (uint)(base) & 0xffff, \
0824   ((uint)(base) >> 16) & 0xff, type, 1, dpl, 1, \
0825   (uint)(lim) >> 16, 0, 0, 1, 0, (uint)(base) >> 24 }
0826 #endif
0827
0828 #define DPL_USER 0x3 // User DPL
0829
0830 // Application segment type bits
0831 #define STA_X 0x8 // Executable segment
0832 #define STA_E 0x4 // Expand down (non-executable segments)
0833 #define STA_C 0x4 // Conforming code segment (executable only)
0834 #define STA_W 0x2 // Writeable (non-executable segments)
0835 #define STA_R 0x2 // Readable (executable segments)
0836 #define STA_A 0x1 // Accessed
0837
0838 // System segment type bits
0839 #define STS_T16A 0x1 // Available 16-bit TSS
0840 #define STS_LDT 0x2 // Local Descriptor Table
0841 #define STS_T16B 0x3 // Busy 16-bit TSS
0842 #define STS_CG16 0x4 // 16-bit Call Gate
0843 #define STS_TG 0x5 // Task Gate / Coum Transmissions
0844 #define STS_IG16 0x6 // 16-bit Interrupt Gate
0845 #define STS_TG16 0x7 // 16-bit Trap Gate
0846 #define STS_T32A 0x9 // Available 32-bit TSS
0847 #define STS_T32B 0xB // Busy 32-bit TSS
0848 #define STS_CG32 0xC // 32-bit Call Gate
0849 #define STS_IG32 0xE // 32-bit Interrupt Gate

```



```

0850 #define STS_TG32    0xF    // 32-bit Trap Gate
0851
0852 // A virtual address 'la' has a three-part structure as follows:
0853 //
0854 // +-----10-----+-----10-----+-----12-----+
0855 // | Page Directory | Page Table | Offset within Page |
0856 // |      Index      |      Index |                   |
0857 // +-----+-----+-----+
0858 // \--- PDX(va) --/ \--- PTX(va) --/
0859
0860 // page directory index
0861 #define PDX(va)      (((uint)(va) >> PDXSHIFT) & 0x3FF)
0862
0863 // page table index
0864 #define PTX(va)      (((uint)(va) >> PTXSHIFT) & 0x3FF)
0865
0866 // construct virtual address from indexes and offset
0867 #define PGADDR(d, t, o) (((uint)(d) << PDXSHIFT | (t) << PTXSHIFT | (o)))
0868
0869 // Page directory and page table constants.
0870 #define NPENTRIES    1024    // # directory entries per page directory
0871 #define NPENTRIES    1024    // # PTEs per page table
0872 #define PGSIZE       4096    // bytes mapped by a page
0873
0874 #define PGSHIFT      12      // log2(PGSIZE)
0875 #define PTXSHIFT     12      // offset of PTX in a linear address
0876 #define PDXSHIFT     22      // offset of PDX in a linear address
0877
0878 #define PGROUNDUP(sz) (((sz)+PGSIZE-1) & ~(PGSIZE-1))
0879 #define PGROUNDOWN(a) (((a) & ~(PGSIZE-1))
0880
0881 // Page table/directory entry flags.
0882 #define PTE_P        0x001    // Present
0883 #define PTE_W        0x002    // Writeable
0884 #define PTE_U        0x004    // User
0885 #define PTE_PWT      0x008    // Write-Through
0886 #define PTE_PCD      0x010    // Cache-Disable
0887 #define PTE_A        0x020    // Accessed
0888 #define PTE_D        0x040    // Dirty
0889 #define PTE_PS       0x080    // Page Size
0890 #define PTE_MBZ      0x180    // Bits must be zero
0891
0892 // Address in page table or page directory entry
0893 #define PTE_ADDR(pte) ((uint)(pte) & ~0xFFF)
0894 #define PTE_FLAGS(pte) ((uint)(pte) & 0xFFF)
0895
0896 #ifndef __ASSEMBLER__
0897 typedef uint pte_t;
0898
0899

```

```

0900 // Task state segment format
0901 struct taskstate {
0902     uint link;           // Old ts selector
0903     uint esp0;          // Stack pointers and segment selectors
0904     ushort ss0;         // after an increase in privilege level
0905     ushort padding1;
0906     uint *esp1;
0907     ushort ssl;
0908     ushort padding2;
0909     uint *esp2;
0910     ushort ss2;
0911     ushort padding3;
0912     void *cr3;          // Page directory base
0913     uint *eip;          // Saved state from last task switch
0914     uint eflags;
0915     uint eax;           // More saved state (registers)
0916     uint ecx;
0917     uint edx;
0918     uint ebx;
0919     uint *esp;
0920     uint *ebp;
0921     uint esi;
0922     uint edi;
0923     ushort es;          // Even more saved state (segment selectors)
0924     ushort padding4;
0925     ushort cs;
0926     ushort padding5;
0927     ushort ss;
0928     ushort padding6;
0929     ushort ds;
0930     ushort padding7;
0931     ushort fs;
0932     ushort padding8;
0933     ushort gs;
0934     ushort padding9;
0935     ushort ldt;
0936     ushort padding10;
0937     ushort t;           // Trap on task switch
0938     ushort iomb;       // I/O map base address
0939 };
0940
0941
0942
0943
0944
0945
0946
0947
0948
0949

```

```

0950 // Gate descriptors for interrupts and traps
0951 struct gatedesc {
0952     uint off_15_0 : 16;    // low 16 bits of offset in segment
0953     uint cs : 16;          // code segment selector
0954     uint args : 5;        // # args, 0 for interrupt/trap gates
0955     uint rsvl : 3;        // reserved(should be zero I guess)
0956     uint type : 4;        // type(STS_{TG,IG32,TG32})
0957     uint s : 1;          // must be 0 (system)
0958     uint dpl : 2;        // descriptor(meaning new) privilege level
0959     uint p : 1;          // Present
0960     uint off_31_16 : 16; // high bits of offset in segment
0961 };
0962
0963 // Set up a normal interrupt/trap gate descriptor.
0964 // - istrap: 1 for a trap (= exception) gate, 0 for an interrupt gate.
0965 // - interrupt gate clears FL_IF, trap gate leaves FL_IF alone
0966 // - sel: Code segment selector for interrupt/trap handler
0967 // - off: Offset in code segment for interrupt/trap handler
0968 // - dpl: Descriptor Privilege Level -
0969 //       the privilege level required for software to invoke
0970 //       this interrupt/trap gate explicitly using an int instruction.
0971 #define SETGATE(gate, istrap, sel, off, d) \
0972 { \
0973     (gate).off_15_0 = (uint)(off) & 0xffff; \
0974     (gate).cs = (sel); \
0975     (gate).args = 0; \
0976     (gate).rsvl = 0; \
0977     (gate).type = (istrap) ? STS_TG32 : STS_IG32; \
0978     (gate).s = 0; \
0979     (gate).dpl = (d); \
0980     (gate).p = 1; \
0981     (gate).off_31_16 = (uint)(off) >> 16; \
0982 }
0983
0984 #endif
0985
0986
0987
0988
0989
0990
0991
0992
0993
0994
0995
0996
0997
0998
0999

```

```

1000 // Format of an ELF executable file
1001
1002 #define ELF_MAGIC 0x464C457FU // "\x7FELF" in little endian
1003
1004 // File header
1005 struct elfhdr {
1006     uint magic; // must equal ELF_MAGIC
1007     uchar elf[12];
1008     ushort type;
1009     ushort machine;
1010     uint version;
1011     uint entry;
1012     uint phoff;
1013     uint shoff;
1014     uint flags;
1015     ushort ehsize;
1016     ushort phentsize;
1017     ushort phnum;
1018     ushort shentsize;
1019     ushort shnum;
1020     ushort shstrndx;
1021 };
1022
1023 // Program section header
1024 struct proghdr {
1025     uint type;
1026     uint off;
1027     uint vaddr;
1028     uint paddr;
1029     uint filesz;
1030     uint memsz;
1031     uint flags;
1032     uint align;
1033 };
1034
1035 // Values for Proghdr type
1036 #define ELF_PROG_LOAD 1
1037
1038 // Flag bits for Proghdr flags
1039 #define ELF_PROG_FLAG_EXEC 1
1040 #define ELF_PROG_FLAG_WRITE 2
1041 #define ELF_PROG_FLAG_READ 4
1042
1043
1044
1045
1046
1047
1048
1049

```

```

1050 // Blank page.
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079
1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099

```

```

1100 # The xv6 kernel starts executing in this file. This file is linked with
1101 # the kernel C code, so it can refer to kernel symbols such as main().
1102 # The boot block (bootasm.S and bootmain.c) jumps to entry below.
1103
1104 # Multiboot header, for multiboot boot loaders like GNU Grub.
1105 # http://www.gnu.org/software/grub/manual/multiboot/multiboot.html
1106 #
1107 # Using GRUB 2, you can boot xv6 from a file stored in a
1108 # Linux file system by copying kernel or kernelmemfs to /boot
1109 # and then adding this menu entry:
1110 #
1111 # menuentry "xv6" {
1112 #   insmod ext2
1113 #   set root='(hd0,msdos1)'
1114 #   set kernel='/boot/kernel'
1115 #   echo "Loading ${kernel}..."
1116 #   multiboot ${kernel} ${kernel}
1117 #   boot
1118 # }
1119
1120 #include "asm.h"
1121 #include "memlayout.h"
1122 #include "mmu.h"
1123 #include "param.h"
1124
1125 # Multiboot header. Data to direct multiboot loader.
1126 .p2align 2
1127 .text
1128 .globl multiboot_header
1129 multiboot_header:
1130   #define magic 0x1badb002
1131   #define flags 0
1132   .long magic
1133   .long flags
1134   .long (-magic-flags)
1135
1136 # By convention, the _start symbol specifies the ELF entry point.
1137 # Since we haven't set up virtual memory yet, our entry point is
1138 # the physical address of 'entry'.
1139 .globl _start
1140 _start = V2P_WO(entry)
1141
1142 # Entering xv6 on boot processor, with paging off.
1143 .globl entry
1144 entry:
1145   # Turn on page size extension for 4Mbyte pages
1146   movl    %cr4, %eax
1147   orl    $(CR4_PSE), %eax
1148   movl    %eax, %cr4
1149   # Set page directory

```

```

1150 movl    $(V2P_W0(entrypgdir)), %eax
1151 movl    %eax, %cr3
1152 # Turn on paging.
1153 movl    %cr0, %eax
1154 orl     $(CR0_PG|CR0_WP), %eax
1155 movl    %eax, %cr0
1156
1157 # Set up the stack pointer.
1158 movl    $(stack + KSTACKSIZE), %esp
1159
1160 # Jump to main(), and switch to executing at
1161 # high addresses. The indirect call is needed because
1162 # the assembler produces a PC-relative instruction
1163 # for a direct jump.
1164 mov     $main, %eax
1165 jmp    *%eax
1166
1167 .comm stack, KSTACKSIZE
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187
1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199

```

```

1200 #include "asm.h"
1201 #include "memlayout.h"
1202 #include "mmu.h"
1203
1204 # Each non-boot CPU ("AP") is started up in response to a STARTUP
1205 # IPI from the boot CPU. Section B.4.2 of the Multi-Processor
1206 # Specification says that the AP will start in real mode with CS:IP
1207 # set to XY00:0000, where XY is an 8-bit value sent with the
1208 # STARTUP. Thus this code must start at a 4096-byte boundary.
1209 #
1210 # Because this code sets DS to zero, it must sit
1211 # at an address in the low 2^16 bytes.
1212 #
1213 # Startothers (in main.c) sends the STARTUPS one at a time.
1214 # It copies this code (start) at 0x7000. It puts the address of
1215 # a newly allocated per-core stack in start-4, the address of the
1216 # place to jump to (mpenter) in start-8, and the physical address
1217 # of entrypgdir in start-12.
1218 #
1219 # This code combines elements of bootasm.S and entry.S.
1220
1221 .code16
1222 .globl start
1223 start:
1224     cli
1225
1226 # Zero data segment registers DS, ES, and SS.
1227     xorw    %ax,%ax
1228     movw    %ax,%ds
1229     movw    %ax,%es
1230     movw    %ax,%ss
1231
1232 # Switch from real to protected mode. Use a bootstrap GDT that makes
1233 # virtual addresses map directly to physical addresses so that the
1234 # effective memory map doesn't change during the transition.
1235     lgdt    gtdesc
1236     movl    %cr0, %eax
1237     orl     $CR0_PE, %eax
1238     movl    %eax, %cr0
1239
1240 # Complete the transition to 32-bit protected mode by using a long jmp
1241 # to reload %cs and %eip. The segment descriptors are set up with no
1242 # translation, so that the mapping is still the identity mapping.
1243     ljmpl   $(SEG_KCODE<<3), $(start32)
1244
1245
1246
1247
1248
1249

```

```

1250 .code32 # Tell assembler to generate 32-bit code now.
1251 start32:
1252 # Set up the protected-mode data segment registers
1253 movw $(SEG_KDATA<<3), %ax # Our data segment selector
1254 movw %ax, %ds # -> DS: Data Segment
1255 movw %ax, %es # -> ES: Extra Segment
1256 movw %ax, %ss # -> SS: Stack Segment
1257 movw $0, %ax # Zero segments not ready for use
1258 movw %ax, %fs # -> FS
1259 movw %ax, %gs # -> GS
1260
1261 # Turn on page size extension for 4Mbyte pages
1262 movl %cr4, %eax
1263 orl $(CR4_PSE), %eax
1264 movl %eax, %cr4
1265 # Use entrypgdir as our initial page table
1266 movl (start-12), %eax
1267 movl %eax, %cr3
1268 # Turn on paging.
1269 movl %cr0, %eax
1270 orl $(CR0_PE|CR0_PG|CR0_WP), %eax
1271 movl %eax, %cr0
1272
1273 # Switch to the stack allocated by startothers()
1274 movl (start-4), %esp
1275 # Call mptenter()
1276 call *(start-8)
1277
1278 movw $0x8a00, %ax
1279 movw %ax, %dx
1280 outw %ax, %dx
1281 movw $0x8ae0, %ax
1282 outw %ax, %dx
1283 spin:
1284 jmp spin
1285
1286 .p2align 2
1287 gdt:
1288 SEG_NULLASM
1289 SEG_ASM(STA_X|STA_R, 0, 0xffffffff)
1290 SEG_ASM(STA_W, 0, 0xffffffff)
1291
1292
1293 gdtdesc:
1294 .word (gdtdesc - gdt - 1)
1295 .long gdt
1296
1297
1298
1299

```

```

1300 #include "types.h"
1301 #include "defs.h"
1302 #include "param.h"
1303 #include "memlayout.h"
1304 #include "mmu.h"
1305 #include "proc.h"
1306 #include "x86.h"
1307
1308 static void startothers(void);
1309 static void mpmain(void) __attribute__((noreturn));
1310 extern pde_t *kpgdir;
1311 extern char end[]; // first address after kernel loaded from ELF file
1312
1313 // Bootstrap processor starts running C code here.
1314 // Allocate a real stack and switch to it, first
1315 // doing some setup required for memory allocator to work.
1316 int
1317 main(void)
1318 {
1319     kinit1(end, P2V(4*1024*1024)); // phys page allocator
1320     kvmalloc(); // kernel page table
1321     mpinit(); // detect other processors
1322     lapicinit(); // interrupt controller
1323     seginit(); // segment descriptors
1324     picinit(); // disable pic
1325     ioapicinit(); // another interrupt controller
1326     consoleinit(); // console hardware
1327     uartinit(); // serial port
1328     pinit(); // process table
1329     tvinit(); // trap vectors
1330     binit(); // buffer cache
1331     fileinit(); // file table
1332     ideinit(); // disk
1333     startothers(); // start other processors
1334     kinit2(P2V(4*1024*1024), P2V(PHYSTOP)); // must come after startothers()
1335     userinit(); // first user process
1336     mpmain(); // finish this processor's setup
1337 }
1338
1339 // Other CPUs jump here from entryother.S.
1340 static void
1341 mptenter(void)
1342 {
1343     switchkvm();
1344     seginit();
1345     lapicinit();
1346     mpmain();
1347 }
1348
1349

```

```

1350 // Common CPU setup code.
1351 static void
1352 mpmain(void)
1353 {
1354     cprintf("cpu%d: starting %d\n", cpuid(), cpuid());
1355     idtinit(); // load idt register
1356     xchg(&mycpu()->started, 1); // tell startothers() we're up
1357     scheduler(); // start running processes
1358 }
1359
1360 pde_t entrypgdir[]; // For entry.S
1361
1362 // Start the non-boot (AP) processors.
1363 static void
1364 startothers(void)
1365 {
1366     extern uchar _binary_entryother_start[], _binary_entryother_size[];
1367     uchar *code;
1368     struct cpu *c;
1369     char *stack;
1370
1371     // Write entry code to unused memory at 0x7000.
1372     // The linker has placed the image of entryother.S in
1373     // _binary_entryother_start.
1374     code = P2V(0x7000);
1375     memmove(code, _binary_entryother_start, (uint)_binary_entryother_size);
1376
1377     for(c = cpus; c < cpus+ncpu; c++){
1378         if(c == mycpu()) // We've started already.
1379             continue;
1380
1381         // Tell entryother.S what stack to use, where to enter, and what
1382         // pgdir to use. We cannot use kpgdir yet, because the AP processor
1383         // is running in low memory, so we use entrypgdir for the APs too.
1384         stack = kalloc();
1385         *(void**)(code-4) = stack + KSTACKSIZE;
1386         *(void**)(code-8) = mpenter;
1387         *(int**)(code-12) = (void *) V2P(entrypgdir);
1388
1389         lapicstartap(c->apicid, V2P(code));
1390
1391         // wait for cpu to finish mpmain()
1392         while(c->started == 0)
1393             ;
1394     }
1395 }
1396
1397
1398
1399

```

```

1400 // The boot page table used in entry.S and entryother.S.
1401 // Page directories (and page tables) must start on page boundaries,
1402 // hence the __aligned__ attribute.
1403 // PTE_PS in a page directory entry enables 4Mbyte pages.
1404
1405 __attribute__((__aligned__(PGSIZE)))
1406 pde_t entrypgdir[NPDENTRIES] = {
1407     // Map VA's [0, 4MB) to PA's [0, 4MB)
1408     [0] = (0) | PTE_P | PTE_W | PTE_PS,
1409     // Map VA's [KERNBASE, KERNBASE+4MB) to PA's [0, 4MB)
1410     [KERNBASE>>PDXSHIFT] = (0) | PTE_P | PTE_W | PTE_PS,
1411 };
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449

```

```
1450 // Blank page.
1451
1452
1453
1454
1455
1456
1457
1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
```

```
1500 // Mutual exclusion lock.
1501 struct spinlock {
1502     uint locked;      // Is the lock held?
1503
1504     // For debugging:
1505     char *name;       // Name of lock.
1506     struct cpu *cpu;  // The cpu holding the lock.
1507     uint pcs[10];    // The call stack (an array of program counters)
1508                     // that locked the lock.
1509 };
1510
1511
1512
1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528
1529
1530
1531
1532
1533
1534
1535
1536
1537
1538
1539
1540
1541
1542
1543
1544
1545
1546
1547
1548
1549
```

```

1550 // Mutual exclusion spin locks.
1551
1552 #include "types.h"
1553 #include "defs.h"
1554 #include "param.h"
1555 #include "x86.h"
1556 #include "memlayout.h"
1557 #include "mmu.h"
1558 #include "proc.h"
1559 #include "spinlock.h"
1560
1561 void
1562 initlock(struct spinlock *lk, char *name)
1563 {
1564     lk->name = name;
1565     lk->locked = 0;
1566     lk->cpu = 0;
1567 }
1568
1569 // Acquire the lock.
1570 // Loops (spins) until the lock is acquired.
1571 // Holding a lock for a long time may cause
1572 // other CPUs to waste time spinning to acquire it.
1573 void
1574 acquire(struct spinlock *lk)
1575 {
1576     pushcli(); // disable interrupts to avoid deadlock.
1577     if(holding(lk))
1578         panic("acquire");
1579
1580     // The xchg is atomic.
1581     while(xchg(&lk->locked, 1) != 0)
1582         ;
1583
1584     // Tell the C compiler and the processor to not move loads or stores
1585     // past this point, to ensure that the critical section's memory
1586     // references happen after the lock is acquired.
1587     __sync_synchronize();
1588
1589     // Record info about lock acquisition for debugging.
1590     lk->cpu = mycpu();
1591     getcallerpcs(&lk, lk->pcs);
1592 }
1593
1594
1595
1596
1597
1598
1599

```

```

1600 // Release the lock.
1601 void
1602 release(struct spinlock *lk)
1603 {
1604     if(!holding(lk))
1605         panic("release");
1606
1607     lk->pcs[0] = 0;
1608     lk->cpu = 0;
1609
1610     // Tell the C compiler and the processor to not move loads or stores
1611     // past this point, to ensure that all the stores in the critical
1612     // section are visible to other cores before the lock is released.
1613     // Both the C compiler and the hardware may re-order loads and
1614     // stores; __sync_synchronize() tells them both not to.
1615     __sync_synchronize();
1616
1617     // Release the lock, equivalent to lk->locked = 0.
1618     // This code can't use a C assignment, since it might
1619     // not be atomic. A real OS would use C atomics here.
1620     asm volatile("movl $0, %0" : "+m" (lk->locked) : );
1621
1622     popcli();
1623 }
1624
1625 // Record the current call stack in pcs[] by following the %ebp chain.
1626 void
1627 getcallerpcs(void *v, uint pcs[])
1628 {
1629     uint *ebp;
1630     int i;
1631
1632     ebp = (uint*)v - 2;
1633     for(i = 0; i < 10; i++){
1634         if(ebp == 0 || ebp < (uint*)KERNBASE || ebp == (uint*)0xffffffff)
1635             break;
1636         pcs[i] = ebp[1]; // saved %eip
1637         ebp = (uint*)ebp[0]; // saved %ebp
1638     }
1639     for(; i < 10; i++)
1640         pcs[i] = 0;
1641 }
1642
1643 // Check whether this cpu is holding the lock.
1644 int
1645 holding(struct spinlock *lock)
1646 {
1647     return lock->locked && lock->cpu == mycpu();
1648 }
1649

```



```

1650 // Pushcli/popcli are like cli/sti except that they are matched:
1651 // it takes two popcli to undo two pushcli. Also, if interrupts
1652 // are off, then pushcli, popcli leaves them off.
1653
1654 void
1655 pushcli(void)
1656 {
1657     int eflags;
1658
1659     eflags = readeflags();
1660     cli();
1661     if(mycpu()->ncli == 0)
1662         mycpu()->intena = eflags & FL_IF;
1663     mycpu()->ncli += 1;
1664 }
1665
1666 void
1667 popcli(void)
1668 {
1669     if(readeflags() & FL_IF)
1670         panic("popcli - interruptible");
1671     if(--mycpu()->ncli < 0)
1672         panic("popcli");
1673     if(mycpu()->ncli == 0 && mycpu()->intena)
1674         sti();
1675 }
1676
1677
1678
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1699

```

```

1700 #include "param.h"
1701 #include "types.h"
1702 #include "defs.h"
1703 #include "x86.h"
1704 #include "memlayout.h"
1705 #include "mmu.h"
1706 #include "proc.h"
1707 #include "elf.h"
1708
1709 extern char data[]; // defined by kernel.ld
1710 pde_t *kpgdir; // for use in scheduler()
1711
1712 // Set up CPU's kernel segment descriptors.
1713 // Run once on entry on each CPU.
1714 void
1715 seginit(void)
1716 {
1717     struct cpu *c;
1718
1719     // Map "logical" addresses to virtual addresses using identity map.
1720     // Cannot share a CODE descriptor for both kernel and user
1721     // because it would have to have DPL_USR, but the CPU forbids
1722     // an interrupt from CPL=0 to DPL=3.
1723     c = &cpu[cuid()];
1724     c->gdt[SEG_KCODE] = SEG(STA_X|STA_R, 0, 0xffffffff, 0);
1725     c->gdt[SEG_KDATA] = SEG(STA_W, 0, 0xffffffff, 0);
1726     c->gdt[SEG_UCODE] = SEG(STA_X|STA_R, 0, 0xffffffff, DPL_USER);
1727     c->gdt[SEG_UDATA] = SEG(STA_W, 0, 0xffffffff, DPL_USER);
1728     lgdt(c->gdt, sizeof(c->gdt));
1729 }
1730
1731 // Return the address of the PTE in page table pgdir
1732 // that corresponds to virtual address va. If alloc!=0,
1733 // create any required page table pages.
1734 static pte_t *
1735 walkpgdir(pde_t *pgdir, const void *va, int alloc)
1736 {
1737     pde_t *pde;
1738     pte_t *pgtab;
1739
1740     pde = &pgdir[PDX(va)];
1741     if(*pde & PTE_P){
1742         pgtab = (pte_t*)P2V(PTE_ADDR(*pde));
1743     } else {
1744         if(!alloc || (pgtab = (pte_t*)kalloc()) == 0)
1745             return 0;
1746         // Make sure all those PTE_P bits are zero.
1747         memset(pgtab, 0, PGSIZE);
1748         // The permissions here are overly generous, but they can
1749         // be further restricted by the permissions in the page table

```

```

1750 // entries, if necessary.
1751 *pde = V2P(pgtab) | PTE_P | PTE_W | PTE_U;
1752 }
1753 return &pgtab[PTX(va)];
1754 }
1755
1756 // Create PTEs for virtual addresses starting at va that refer to
1757 // physical addresses starting at pa. va and size might not
1758 // be page-aligned.
1759 static int
1760 mappages(pde_t *pgdir, void *va, uint size, uint pa, int perm)
1761 {
1762     char *a, *last;
1763     pte_t *pte;
1764
1765     a = (char*)PGROUNDDOWN((uint)va);
1766     last = (char*)PGROUNDDOWN((uint)va) + size - 1;
1767     for(;;){
1768         if((pte = walkpgdir(pgdir, a, 1)) == 0)
1769             return -1;
1770         if(*pte & PTE_P)
1771             panic("remap");
1772         *pte = pa | perm | PTE_P;
1773         if(a == last)
1774             break;
1775         a += PGSIZE;
1776         pa += PGSIZE;
1777     }
1778     return 0;
1779 }
1780
1781 // There is one page table per process, plus one that's used when
1782 // a CPU is not running any process (kpgdir). The kernel uses the
1783 // current process's page table during system calls and interrupts;
1784 // page protection bits prevent user code from using the kernel's
1785 // mappings.
1786 //
1787 // setupkvm() and exec() set up every page table like this:
1788 //
1789 // 0..KERNBASE: user memory (text+data+stack+heap), mapped to
1790 // phys memory allocated by the kernel
1791 // KERNBASE..KERNBASE+EXTMEM: mapped to 0..EXTMEM (for I/O space)
1792 // KERNBASE+EXTMEM..data: mapped to EXTMEM..V2P(data)
1793 // for the kernel's instructions and r/o data
1794 // data..KERNBASE+PHYSTOP: mapped to V2P(data)..PHYSTOP,
1795 // rw data + free physical memory
1796 // 0xfe000000..0: mapped direct (devices such as ioapic)
1797 //
1798 // The kernel allocates physical memory for its heap and for user memory
1799 // between V2P(end) and the end of physical memory (PHYSTOP)

```

```

1800 // (directly addressable from end..P2V(PHYSTOP)).
1801
1802 // This table defines the kernel's mappings, which are present in
1803 // every process's page table.
1804 static struct kmap {
1805     void *virt;
1806     uint phys_start;
1807     uint phys_end;
1808     int perm;
1809 } kmap[] = {
1810     { (void*)KERNBASE, 0,          EXTMEM,   PTE_W}, // I/O space
1811     { (void*)KERNLINK, V2P(KERNLINK), V2P(data), 0}, // kern text+rodata
1812     { (void*)data,     V2P(data),   PHYSTOP,  PTE_W}, // kern data+memory
1813     { (void*)DEVSPACE, DEVSPACE,    0,       PTE_W}, // more devices
1814 };
1815
1816 // Set up kernel part of a page table.
1817 pde_t*
1818 setupkvm(void)
1819 {
1820     pde_t *pgdir;
1821     struct kmap *k;
1822
1823     if((pgdir = (pde_t*)kalloc()) == 0)
1824         return 0;
1825     memset(pgdir, 0, PGSIZE);
1826     if (P2V(PHYSTOP) > (void*)DEVSPACE)
1827         panic("PHYSTOP too high");
1828     for(k = kmap; k < &kmap[NELEM(kmap)]; k++)
1829         if(mappages(pgdir, k->virt, k->phys_end - k->phys_start,
1830             (uint)k->phys_start, k->perm) < 0) {
1831             freevm(pgdir);
1832             return 0;
1833         }
1834     return pgdir;
1835 }
1836
1837 // Allocate one page table for the machine for the kernel address
1838 // space for scheduler processes.
1839 void
1840 kvmalloc(void)
1841 {
1842     kpgdir = setupkvm();
1843     switchkvm();
1844 }
1845
1846
1847
1848
1849

```

```

1850 // Switch h/w page table register to the kernel-only page table,
1851 // for when no process is running.
1852 void
1853 switchkvmm(void)
1854 {
1855     lcr3(V2P(kpgdir)); // switch to the kernel page table
1856 }
1857
1858 // Switch TSS and h/w page table to correspond to process p.
1859 void
1860 switchuvm(struct proc *p)
1861 {
1862     if(p == 0)
1863         panic("switchuvm: no process");
1864     if(p->kstack == 0)
1865         panic("switchuvm: no kstack");
1866     if(p->pgdir == 0)
1867         panic("switchuvm: no pgdir");
1868
1869     pushcli();
1870     mycpu()->gdt[SEG_TSS] = SEG16(STS_T32A, &mycpu()->ts, sizeof(mycpu()->ts)-1);
1871     mycpu()->gdt[SEG_TSS].s = 0;
1872     mycpu()->ts.ss0 = SEG_KDATA << 3;
1873     mycpu()->ts.esp0 = (uint)p->kstack + KSTACKSIZE;
1874     // setting IOPL=0 in eflags *and* iomb beyond the tss segment limit
1875     // forbids I/O instructions (e.g., inb and outb) from user space
1876     mycpu()->ts.iomb = (ushort) 0xFFFF;
1877     ltr(SEG_TSS << 3);
1878     lcr3(V2P(p->pgdir)); // switch to process's address space
1879     popcli();
1880 }
1881
1882 // Load the initcode into address 0 of pgdir.
1883 // sz must be less than a page.
1884 void
1885 inituvm(pde_t *pgdir, char *init, uint sz)
1886 {
1887     char *mem;
1888
1889     if(sz >= PGSIZE)
1890         panic("inituvm: more than a page");
1891     mem = kalloc();
1892     memset(mem, 0, PGSIZE);
1893     mappages(pgdir, 0, PGSIZE, V2P(mem), PTE_W|PTE_U);
1894     memmove(mem, init, sz);
1895 }
1896
1897
1898
1899

```

```

1900 // Load a program segment into pgdir.  addr must be page-aligned
1901 // and the pages from addr to addr+sz must already be mapped.
1902 int
1903 loaduvm(pde_t *pgdir, char *addr, struct inode *ip, uint offset, uint sz)
1904 {
1905     uint i, pa, n;
1906     pte_t *pte;
1907
1908     if((uint) addr % PGSIZE != 0)
1909         panic("loaduvm: addr must be page aligned");
1910     for(i = 0; i < sz; i += PGSIZE){
1911         if((pte = walkpgdir(pgdir, addr+i, 0)) == 0)
1912             panic("loaduvm: address should exist");
1913         pa = PTE_ADDR(*pte);
1914         if(sz - i < PGSIZE)
1915             n = sz - i;
1916         else
1917             n = PGSIZE;
1918         if(readi(ip, P2V(pa), offset+i, n) != n)
1919             return -1;
1920     }
1921     return 0;
1922 }
1923
1924 // Allocate page tables and physical memory to grow process from oldsz to
1925 // newsz, which need not be page aligned. Returns new size or 0 on error.
1926 int
1927 allocuvm(pde_t *pgdir, uint oldsz, uint newsz)
1928 {
1929     char *mem;
1930     uint a;
1931
1932     if(newsz >= KERNBASE)
1933         return 0;
1934     if(newsz < oldsz)
1935         return oldsz;
1936
1937     a = PGROUNDUP(oldsz);
1938     for(i = a < newsz; i += PGSIZE){
1939         mem = kalloc();
1940         if(mem == 0){
1941             cprintf("allocuvm out of memory\n");
1942             deallocuvm(pgdir, newsz, oldsz);
1943             return 0;
1944         }
1945         memset(mem, 0, PGSIZE);
1946         if(mappages(pgdir, (char*)a, PGSIZE, V2P(mem), PTE_W|PTE_U) < 0){
1947             cprintf("allocuvm out of memory (2)\n");
1948             deallocuvm(pgdir, newsz, oldsz);
1949             kfree(mem);

```

```

1950     return 0;
1951   }
1952 }
1953 return newsz;
1954 }
1955
1956 // Deallocate user pages to bring the process size from oldsz to
1957 // newsz.  oldsz and newsz need not be page-aligned, nor does newsz
1958 // need to be less than oldsz.  oldsz can be larger than the actual
1959 // process size.  Returns the new process size.
1960 int
1961 deallocvm(pde_t *pgdir, uint oldsz, uint newsz)
1962 {
1963   pte_t *pte;
1964   uint a, pa;
1965
1966   if(newsz >= oldsz)
1967     return oldsz;
1968
1969   a = PGROUNDUP(newsz);
1970   for(; a < oldsz; a += PGSIZE){
1971     pte = walkpgdir(pgdir, (char*)a, 0);
1972     if(!pte)
1973       a = PGADDR(PDX(a) + 1, 0, 0) - PGSIZE;
1974     else if((*pte & PTE_P) != 0){
1975       pa = PTE_ADDR(*pte);
1976       if(pa == 0)
1977         panic("kfree");
1978       char *v = P2V(pa);
1979       kfree(v);
1980       *pte = 0;
1981     }
1982   }
1983   return newsz;
1984 }
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997
1998
1999

```

```

2000 // Free a page table and all the physical memory pages
2001 // in the user part.
2002 void
2003 freevm(pde_t *pgdir)
2004 {
2005   uint i;
2006
2007   if(pgdir == 0)
2008     panic("freevm: no pgdir");
2009   deallocvm(pgdir, KERNBASE, 0);
2010   for(i = 0; i < NPENTRIES; i++){
2011     if(pgdir[i] & PTE_P){
2012       char *v = P2V(PTE_ADDR(pgdir[i]));
2013       kfree(v);
2014     }
2015   }
2016   kfree((char*)pgdir);
2017 }
2018
2019 // Clear PTE_U on a page. Used to create an inaccessible
2020 // page beneath the user stack.
2021 void
2022 clearpteu(pde_t *pgdir, char *uva)
2023 {
2024   pte_t *pte;
2025
2026   pte = walkpgdir(pgdir, uva, 0);
2027   if(pte == 0)
2028     panic("clearpteu");
2029   *pte &= ~PTE_U;
2030 }
2031
2032 // Given a parent process's page table, create a copy
2033 // of it for a child.
2034 pde_t*
2035 copyvm(pde_t *pgdir, uint sz)
2036 {
2037   pde_t *d;
2038   pte_t *pte;
2039   uint pa, i, flags;
2040   char *mem;
2041
2042   if((d = setupkvm()) == 0)
2043     return 0;
2044   for(i = 0; i < sz; i += PGSIZE){
2045     if((pte = walkpgdir(pgdir, (void *) i, 0)) == 0)
2046       panic("copyvm: pte should exist");
2047     if(!(*pte & PTE_P))
2048       panic("copyvm: page not present");
2049     pa = PTE_ADDR(*pte);

```

```

2050     flags = PTE_FLAGS(*pte);
2051     if((mem = kalloc()) == 0)
2052         goto bad;
2053     memmove(mem, (char*)P2V(pa), PGSIZE);
2054     if(mappages(d, (void*)i, PGSIZE, V2P(mem), flags) < 0)
2055         goto bad;
2056     }
2057     return d;
2058
2059 bad:
2060     freevm(d);
2061     return 0;
2062 }
2063
2064
2065
2066
2067
2068
2069
2070
2071
2072
2073
2074
2075
2076
2077
2078
2079
2080
2081
2082
2083
2084
2085
2086
2087
2088
2089
2090
2091
2092
2093
2094
2095
2096
2097
2098
2099

```

```

2100 // Map user virtual address to kernel address.
2101 char*
2102 uva2ka(pde_t *pgdir, char *uva)
2103 {
2104     pte_t *pte;
2105
2106     pte = walkpgdir(pgdir, uva, 0);
2107     if((*pte & PTE_P) == 0)
2108         return 0;
2109     if((*pte & PTE_U) == 0)
2110         return 0;
2111     return (char*)P2V(PTE_ADDR(*pte));
2112 }
2113
2114 // Copy len bytes from p to user address va in page table pgdir.
2115 // Most useful when pgdir is not the current page table.
2116 // uva2ka ensures this only works for PTE_U pages.
2117 int
2118 copyout(pde_t *pgdir, uint va, void *p, uint len)
2119 {
2120     char *buf, *pa0;
2121     uint n, va0;
2122
2123     buf = (char*)p;
2124     while(len > 0){
2125         va0 = (uint)PGROUNDDOWN(va);
2126         pa0 = uva2ka(pgdir, (char*)va0);
2127         if(pa0 == 0)
2128             return -1;
2129         n = PGSIZE - (va - va0);
2130         if(n > len)
2131             n = len;
2132         memmove(pa0 + (va - va0), buf, n);
2133         len -= n;
2134         buf += n;
2135         va = va0 + PGSIZE;
2136     }
2137     return 0;
2138 }
2139
2140
2141
2142
2143
2144
2145
2146
2147
2148
2149

```

2150 // Blank page.
2151
2152
2153
2154
2155
2156
2157
2158
2159
2160
2161
2162
2163
2164
2165
2166
2167
2168
2169
2170
2171
2172
2173
2174
2175
2176
2177
2178
2179
2180
2181
2182
2183
2184
2185
2186
2187
2188
2189
2190
2191
2192
2193
2194
2195
2196
2197
2198
2199

2200 // Blank page.
2201
2202
2203
2204
2205
2206
2207
2208
2209
2210
2211
2212
2213
2214
2215
2216
2217
2218
2219
2220
2221
2222
2223
2224
2225
2226
2227
2228
2229
2230
2231
2232
2233
2234
2235
2236
2237
2238
2239
2240
2241
2242
2243
2244
2245
2246
2247
2248
2249

```

2250 // Blank page.
2251
2252
2253
2254
2255
2256
2257
2258
2259
2260
2261
2262
2263
2264
2265
2266
2267
2268
2269
2270
2271
2272
2273
2274
2275
2276
2277
2278
2279
2280
2281
2282
2283
2284
2285
2286
2287
2288
2289
2290
2291
2292
2293
2294
2295
2296
2297
2298
2299

```

```

2300 // Per-CPU state
2301 struct cpu {
2302     uchar apicid;           // Local APIC ID
2303     struct context *scheduler; // swtch() here to enter scheduler
2304     struct taskstate ts;    // Used by x86 to find stack for interrupt
2305     struct segdesc gdt[NSEGS]; // x86 global descriptor table
2306     volatile uint started;  // Has the CPU started?
2307     int ncli;               // Depth of pushcli nesting.
2308     int intena;             // Were interrupts enabled before pushcli?
2309     struct proc *proc;     // The process running on this cpu or null
2310 };
2311
2312 extern struct cpu cpus[NCPU];
2313 extern int ncpu;
2314
2315
2316 // Saved registers for kernel context switches.
2317 // Don't need to save all the segment registers (%cs, etc),
2318 // because they are constant across kernel contexts.
2319 // Don't need to save %eax, %ecx, %edx, because the
2320 // x86 convention is that the caller has saved them.
2321 // Contexts are stored at the bottom of the stack they
2322 // describe; the stack pointer is the address of the context.
2323 // The layout of the context matches the layout of the stack in swtch.S
2324 // at the "Switch stacks" comment. Switch doesn't save eip explicitly,
2325 // but it is on the stack and allocproc() manipulates it.
2326 struct context {
2327     uint edi;
2328     uint esi;
2329     uint ebx;
2330     uint ebp;
2331     uint eip;
2332 };
2333
2334 enum procstate { UNUSED, EMBRYO, SLEEPING, RUNNABLE, RUNNING, ZOMBIE };
2335
2336 // Per-process state
2337 struct proc {
2338     uint sz;                 // Size of process memory (bytes)
2339     pde_t* pgdir;           // Page table
2340     char *kstack;           // Bottom of kernel stack for this process
2341     enum procstate state;   // Process state
2342     int pid;                 // Process ID
2343     struct proc *parent;    // Parent process
2344     struct trapframe *tf;   // Trap frame for current syscall
2345     struct context *context; // swtch() here to run process
2346     void *chan;             // If non-zero, sleeping on chan
2347     int killed;             // If non-zero, have been killed
2348     struct file *ofile[NOFILE]; // Open files
2349     struct inode *cwd;      // Current directory

```

```

2350 char name[16];           // Process name (debugging)
2351 };
2352
2353 // Process memory is laid out contiguously, low addresses first:
2354 // text
2355 // original data and bss
2356 // fixed-size stack
2357 // expandable heap
2358
2359
2360
2361
2362
2363
2364
2365
2366
2367
2368
2369
2370
2371
2372
2373
2374
2375
2376
2377
2378
2379
2380
2381
2382
2383
2384
2385
2386
2387
2388
2389
2390
2391
2392
2393
2394
2395
2396
2397
2398
2399

```

```

2400 #include "types.h"
2401 #include "defs.h"
2402 #include "param.h"
2403 #include "memlayout.h"
2404 #include "mmu.h"
2405 #include "x86.h"
2406 #include "proc.h"
2407 #include "spinlock.h"
2408
2409 struct {
2410     struct spinlock lock;
2411     struct proc proc[NPROC];
2412 } ptable;
2413
2414 static struct proc *initproc;
2415
2416 int nextpid = 1;
2417 extern void forkret(void);
2418 extern void trapret(void);
2419
2420 static void wakeup1(void *chan);
2421
2422 void
2423 pinit(void)
2424 {
2425     initlock(&ptable.lock, "ptable");
2426 }
2427
2428 // Must be called with interrupts disabled
2429 int
2430 cpuid() {
2431     return mycpu() - cpus;
2432 }
2433
2434 // Must be called with interrupts disabled to avoid the caller being resched
2435 // between reading lapicid and running through the loop.
2436 struct cpu*
2437 mycpu(void)
2438 {
2439     int apicid, i;
2440
2441     if(readeflags() & FL_IF)
2442         panic("mycpu called with interrupts enabled\n");
2443
2444     apicid = lapicid();
2445     // APIC IDs are not guaranteed to be contiguous. Maybe we should have
2446     // a reverse map, or reserve a register to store &cpus[i].
2447     for (i = 0; i < ncpu; ++i) {
2448         if (cpus[i].apicid == apicid)
2449             return &cpus[i];

```



```

2450 }
2451 panic("unknown apicid\n");
2452 }
2453
2454 // Disable interrupts so that we are not rescheduled
2455 // while reading proc from the cpu structure
2456 struct proc*
2457 myproc(void) {
2458     struct cpu *c;
2459     struct proc *p;
2460     pushcli();
2461     c = mycpu();
2462     p = c->proc;
2463     popcli();
2464     return p;
2465 }
2466
2467
2468 // Look in the process table for an UNUSED proc.
2469 // If found, change state to EMBRYO and initialize
2470 // state required to run in the kernel.
2471 // Otherwise return 0.
2472 static struct proc*
2473 allocproc(void)
2474 {
2475     struct proc *p;
2476     char *sp;
2477
2478     acquire(&ptable.lock);
2479
2480     for(p = ptable.proc; p < &ptable.proc[NPROC]; p++)
2481         if(p->state == UNUSED)
2482             goto found;
2483
2484     release(&ptable.lock);
2485     return 0;
2486
2487 found:
2488     p->state = EMBRYO;
2489     p->pid = nextpid++;
2490
2491     release(&ptable.lock);
2492
2493     // Allocate kernel stack.
2494     if((p->kstack = kalloc()) == 0){
2495         p->state = UNUSED;
2496         return 0;
2497     }
2498     sp = p->kstack + KSTACKSIZE;
2499

```

```

2500 // Leave room for trap frame.
2501 sp -= sizeof *p->tf;
2502 p->tf = (struct trapframe*)sp;
2503
2504 // Set up new context to start executing at forkret,
2505 // which returns to trapret.
2506 sp -= 4;
2507 *(uint*)sp = (uint)trapret;
2508
2509 sp -= sizeof *p->context;
2510 p->context = (struct context*)sp;
2511 memset(p->context, 0, sizeof *p->context);
2512 p->context->eip = (uint)forkret;
2513
2514 return p;
2515 }
2516
2517
2518 // Set up first user process.
2519 void
2520 userinit(void)
2521 {
2522     struct proc *p;
2523     extern char _binary_initcode_start[], _binary_initcode_size[];
2524
2525     p = allocproc();
2526
2527     initproc = p;
2528     if((p->pgdir = setupkvm()) == 0)
2529         panic("userinit: out of memory?");
2530     initvm(p->pgdir, _binary_initcode_start, (int)_binary_initcode_size);
2531     p->sz = PGSIZE;
2532     memset(p->tf, 0, sizeof(*p->tf));
2533     p->tf->cs = (SEG_UCODE << 3) | DPL_USER;
2534     p->tf->ds = (SEG_UDATA << 3) | DPL_USER;
2535     p->tf->es = p->tf->ds;
2536     p->tf->ss = p->tf->ds;
2537     p->tf->eflags = FL_IF;
2538     p->tf->esp = PGSIZE;
2539     p->tf->eip = 0; // beginning of initcode.S
2540
2541     safestrcpy(p->name, "initcode", sizeof(p->name));
2542     p->cwd = namei("/");
2543
2544     // this assignment to p->state lets other cores
2545     // run this process. the acquire forces the above
2546     // writes to be visible, and the lock is also needed
2547     // because the assignment might not be atomic.
2548     acquire(&ptable.lock);
2549

```

```

2550 p->state = RUNNABLE;
2551
2552 release(&ptable.lock);
2553 }
2554
2555 // Grow current process's memory by n bytes.
2556 // Return 0 on success, -1 on failure.
2557 int
2558 growproc(int n)
2559 {
2560     uint sz;
2561     struct proc *curproc = myproc();
2562
2563     sz = curproc->sz;
2564     if(n > 0){
2565         if((sz = allocuvn(curproc->pgdir, sz, sz + n)) == 0)
2566             return -1;
2567     } else if(n < 0){
2568         if((sz = deallocuvn(curproc->pgdir, sz, sz + n)) == 0)
2569             return -1;
2570     }
2571     curproc->sz = sz;
2572     switchvm(curproc);
2573     return 0;
2574 }
2575
2576 // Create a new process copying p as the parent.
2577 // Sets up stack to return as if from system call.
2578 // Caller must set state of returned proc to RUNNABLE.
2579 int
2580 fork(void)
2581 {
2582     int i, pid;
2583     struct proc *np;
2584     struct proc *curproc = myproc();
2585
2586     // Allocate process.
2587     if((np = allocproc()) == 0){
2588         return -1;
2589     }
2590
2591     // Copy process state from proc.
2592     if((np->pgdir = copyvm(curproc->pgdir, curproc->sz)) == 0){
2593         kfree(np->kstack);
2594         np->kstack = 0;
2595         np->state = UNUSED;
2596         return -1;
2597     }
2598     np->sz = curproc->sz;
2599     np->parent = curproc;

```

```

2600 *np->tf = *curproc->tf;
2601
2602 // Clear %eax so that fork returns 0 in the child.
2603 np->tf->eax = 0;
2604
2605 for(i = 0; i < NOFILE; i++){
2606     if(curproc->ofile[i])
2607         np->ofile[i] = filedup(curproc->ofile[i]);
2608     np->cwd = idup(curproc->cwd);
2609
2610     safestrcpy(np->name, curproc->name, sizeof(curproc->name));
2611
2612     pid = np->pid;
2613
2614     acquire(&ptable.lock);
2615     np->state = RUNNABLE;
2616     release(&ptable.lock);
2617
2618     return pid;
2619 }
2620
2621 // Exit the current process. Does not return.
2622 // An exited process remains in the zombie state
2623 // until its parent calls wait() to find out it exited.
2624 void
2625 exit(void)
2626 {
2627     struct proc *curproc = myproc();
2628     struct proc *p;
2629     int fd;
2630
2631     if(curproc == initproc)
2632         panic("init exiting");
2633
2634     // Close all open files.
2635     for(fd = 0; fd < NOFILE; fd++){
2636         if(curproc->ofile[fd]){
2637             fileclose(curproc->ofile[fd]);
2638             curproc->ofile[fd] = 0;
2639         }
2640     }
2641
2642     begin_op();
2643     iput(curproc->cwd);
2644     end_op();
2645     curproc->cwd = 0;
2646
2647     acquire(&ptable.lock);

```

```

2650 // Parent might be sleeping in wait().
2651 wakeup1(curproc->parent);
2652
2653 // Pass abandoned children to init.
2654 for(p = ptable.proc; p < &ptable.proc[NPROC]; p++){
2655     if(p->parent == curproc){
2656         p->parent = initproc;
2657         if(p->state == ZOMBIE)
2658             wakeup1(initproc);
2659     }
2660 }
2661
2662 // Jump into the scheduler, never to return.
2663 curproc->state = ZOMBIE;
2664 sched();
2665 panic("zombie exit");
2666 }
2667
2668 // Wait for a child process to exit and return its pid.
2669 // Return -1 if this process has no children.
2670 int
2671 wait(void)
2672 {
2673     struct proc *p;
2674     int havekids, pid;
2675     struct proc *curproc = myproc();
2676
2677     acquire(&ptable.lock);
2678     for(;;){
2679         // Scan through table looking for exited children.
2680         havekids = 0;
2681         for(p = ptable.proc; p < &ptable.proc[NPROC]; p++){
2682             if(p->parent != curproc)
2683                 continue;
2684             havekids = 1;
2685             if(p->state == ZOMBIE){
2686                 // Found one.
2687                 pid = p->pid;
2688                 kfree(p->kstack);
2689                 p->kstack = 0;
2690                 freevm(p->pgdir);
2691                 p->pid = 0;
2692                 p->parent = 0;
2693                 p->name[0] = 0;
2694                 p->killed = 0;
2695                 p->state = UNUSED;
2696                 release(&ptable.lock);
2697                 return pid;
2698             }
2699         }

```

```

2700 // No point waiting if we don't have any children.
2701 if(!havekids || curproc->killed){
2702     release(&ptable.lock);
2703     return -1;
2704 }
2705
2706 // Wait for children to exit. (See wakeup1 call in proc_exit.)
2707 sleep(curproc, &ptable.lock);
2708 }
2709 }
2710
2711
2712
2713
2714
2715
2716
2717
2718
2719
2720
2721
2722
2723
2724
2725
2726
2727
2728
2729
2730
2731
2732
2733
2734
2735
2736
2737
2738
2739
2740
2741
2742
2743
2744
2745
2746
2747
2748
2749

```

```

2750 // Per-CPU process scheduler.
2751 // Each CPU calls scheduler() after setting itself up.
2752 // Scheduler never returns. It loops, doing:
2753 // - choose a process to run
2754 // - swtch to start running that process
2755 // - eventually that process transfers control
2756 //   via swtch back to the scheduler.
2757 void
2758 scheduler(void)
2759 {
2760     struct proc *p;
2761     struct cpu *c = mycpu();
2762     c->proc = 0;
2763
2764     for(;;){
2765         // Enable interrupts on this processor.
2766         sti();
2767
2768         // Loop over process table looking for process to run.
2769         acquire(&ptable.lock);
2770         for(p = ptable.proc; p < &ptable.proc[NPROC]; p++){
2771             if(p->state != RUNNABLE)
2772                 continue;
2773
2774             // Switch to chosen process. It is the process's job
2775             // to release ptable.lock and then reacquire it
2776             // before jumping back to us.
2777             c->proc = p;
2778             switchvm(p);
2779             p->state = RUNNING;
2780
2781             swtch(&(c->scheduler), p->context);
2782             switchkvm();
2783
2784             // Process is done running for now.
2785             // It should have changed its p->state before coming back.
2786             c->proc = 0;
2787         }
2788         release(&ptable.lock);
2789     }
2790 }
2791 }
2792
2793
2794
2795
2796
2797
2798
2799

```

```

2800 // Enter scheduler. Must hold only ptable.lock
2801 // and have changed proc->state. Saves and restores
2802 // intena because intena is a property of this
2803 // kernel thread, not this CPU. It should
2804 // be proc->intena and proc->ncli, but that would
2805 // break in the few places where a lock is held but
2806 // there's no process.
2807 void
2808 sched(void)
2809 {
2810     int intena;
2811     struct proc *p = myproc();
2812
2813     if(!holding(&ptable.lock))
2814         panic("sched ptable.lock");
2815     if(mycpu()->ncli != 1)
2816         panic("sched locks");
2817     if(p->state == RUNNING)
2818         panic("sched running");
2819     if(readeflags() & FL_IF)
2820         panic("sched interruptible");
2821     intena = mycpu()->intena;
2822     swtch(&p->context, mycpu()->scheduler);
2823     mycpu()->intena = intena;
2824 }
2825
2826 // Give up the CPU for one scheduling round.
2827 void
2828 yield(void)
2829 {
2830     acquire(&ptable.lock);
2831     myproc()->state = RUNNABLE;
2832     sched();
2833     release(&ptable.lock);
2834 }
2835
2836
2837
2838
2839
2840
2841
2842
2843
2844
2845
2846
2847
2848
2849

```

```

2850 // A fork child's very first scheduling by scheduler()
2851 // will swtch here. "Return" to user space.
2852 void
2853 forkret(void)
2854 {
2855     static int first = 1;
2856     // Still holding ptable.lock from scheduler.
2857     release(&ptable.lock);
2858
2859     if (first) {
2860         // Some initialization functions must be run in the context
2861         // of a regular process (e.g., they call sleep), and thus cannot
2862         // be run from main().
2863         first = 0;
2864         iinit(ROOTDEV);
2865         initlog(ROOTDEV);
2866     }
2867
2868     // Return to "caller", actually trapret (see allocproc).
2869 }
2870
2871 // Atomically release lock and sleep on chan.
2872 // Reacquires lock when awakened.
2873 void
2874 sleep(void *chan, struct spinlock *lk)
2875 {
2876     struct proc *p = myproc();
2877
2878     if(p == 0)
2879         panic("sleep");
2880
2881     if(lk == 0)
2882         panic("sleep without lk");
2883
2884     // Must acquire ptable.lock in order to
2885     // change p->state and then call sched.
2886     // Once we hold ptable.lock, we can be
2887     // guaranteed that we won't miss any wakeup
2888     // (wakeup runs with ptable.lock locked),
2889     // so it's okay to release lk.
2890     if(lk != &ptable.lock){
2891         acquire(&ptable.lock);
2892         release(lk);
2893     }
2894     // Go to sleep.
2895     p->chan = chan;
2896     p->state = SLEEPING;
2897
2898     sched();
2899

```

```

2900 // Tidy up.
2901 p->chan = 0;
2902
2903 // Reacquire original lock.
2904 if(lk != &ptable.lock){
2905     release(&ptable.lock);
2906     acquire(lk);
2907 }
2908 }
2909
2910
2911
2912
2913
2914
2915
2916
2917
2918
2919
2920
2921
2922
2923
2924
2925
2926
2927
2928
2929
2930
2931
2932
2933
2934
2935
2936
2937
2938
2939
2940
2941
2942
2943
2944
2945
2946
2947
2948
2949

```

```

2950 // Wake up all processes sleeping on chan.
2951 // The ptable lock must be held.
2952 static void
2953 wakeup1(void *chan)
2954 {
2955     struct proc *p;
2956
2957     for(p = ptable.proc; p < &ptable.proc[NPROC]; p++)
2958         if(p->state == SLEEPING && p->chan == chan)
2959             p->state = RUNNABLE;
2960 }
2961
2962 // Wake up all processes sleeping on chan.
2963 void
2964 wakeup(void *chan)
2965 {
2966     acquire(&ptable.lock);
2967     wakeup1(chan);
2968     release(&ptable.lock);
2969 }
2970
2971 // Kill the process with the given pid.
2972 // Process won't exit until it returns
2973 // to user space (see trap in trap.c).
2974 int
2975 kill(int pid)
2976 {
2977     struct proc *p;
2978
2979     acquire(&ptable.lock);
2980     for(p = ptable.proc; p < &ptable.proc[NPROC]; p++){
2981         if(p->pid == pid){
2982             p->killed = 1;
2983             // Wake process from sleep if necessary.
2984             if(p->state == SLEEPING)
2985                 p->state = RUNNABLE;
2986             release(&ptable.lock);
2987             return 0;
2988         }
2989     }
2990     release(&ptable.lock);
2991     return -1;
2992 }
2993
2994
2995
2996
2997
2998
2999

```

```

3000 // Print a process listing to console. For debugging.
3001 // Runs when user types ^P on console.
3002 // No lock to avoid wedging a stuck machine further.
3003 void
3004 procdump(void)
3005 {
3006     static char *states[] = {
3007         [UNUSED]    "unused",
3008         [EMBRYO]    "embryo",
3009         [SLEEPING]  "sleep ",
3010         [RUNNABLE]  "runble",
3011         [RUNNING]   "run  ",
3012         [ZOMBIE]    "zombie"
3013     };
3014     int i;
3015     struct proc *p;
3016     char *state;
3017     uint pc[10];
3018
3019     for(p = ptable.proc; p < &ptable.proc[NPROC]; p++){
3020         if(p->state == UNUSED)
3021             continue;
3022         if(p->state >= 0 && p->state < NELEM(states) && states[p->state])
3023             state = states[p->state];
3024         else
3025             state = "???";
3026         cprintf("%d %s %s", p->pid, state, p->name);
3027         if(p->state == SLEEPING){
3028             getcallerpcs((uint*)p->context->ebp+2, pc);
3029             for(i=0; i<10 && pc[i] != 0; i++)
3030                 cprintf(" %p", pc[i]);
3031         }
3032         cprintf("\n");
3033     }
3034 }
3035
3036
3037
3038
3039
3040
3041
3042
3043
3044
3045
3046
3047
3048
3049

```

```

3050 # Context switch
3051 #
3052 # void swtch(struct context **old, struct context *new);
3053 #
3054 # Save current register context in old
3055 # and then load register context from new.
3056
3057 .globl swtch
3058 swtch:
3059     movl 4(%esp), %eax
3060     movl 8(%esp), %edx
3061
3062 # Save old callee-save registers
3063     pushl %ebp
3064     pushl %ebx
3065     pushl %esi
3066     pushl %edi
3067
3068 # Switch stacks
3069     movl %esp, (%eax)
3070     movl %edx, %esp
3071
3072 # Load new callee-save registers
3073     popl %edi
3074     popl %esi
3075     popl %ebx
3076     popl %ebp
3077     ret
3078
3079
3080
3081
3082
3083
3084
3085
3086
3087
3088
3089
3090
3091
3092
3093
3094
3095
3096
3097
3098
3099

```

```

3100 // Physical memory allocator, intended to allocate
3101 // memory for user processes, kernel stacks, page table pages,
3102 // and pipe buffers. Allocates 4096-byte pages.
3103
3104 #include "types.h"
3105 #include "defs.h"
3106 #include "param.h"
3107 #include "memlayout.h"
3108 #include "mmu.h"
3109 #include "spinlock.h"
3110
3111 void freerange(void *vstart, void *vend);
3112 extern char end[]; // first address after kernel loaded from ELF file
3113 // defined by the kernel linker script in kernel.ld
3114
3115 struct run {
3116     struct run *next;
3117 };
3118
3119 struct {
3120     struct spinlock lock;
3121     int use_lock;
3122     struct run *freelist;
3123 } kmem;
3124
3125 // Initialization happens in two phases.
3126 // 1. main() calls kinit1() while still using entrypgdir to place just
3127 // the pages mapped by entrypgdir on free list.
3128 // 2. main() calls kinit2() with the rest of the physical pages
3129 // after installing a full page table that maps them on all cores.
3130 void
3131 kinit1(void *vstart, void *vend)
3132 {
3133     initlock(&kmem.lock, "kmem");
3134     kmem.use_lock = 0;
3135     freerange(vstart, vend);
3136 }
3137
3138 void
3139 kinit2(void *vstart, void *vend)
3140 {
3141     freerange(vstart, vend);
3142     kmem.use_lock = 1;
3143 }
3144
3145
3146
3147
3148
3149

```

```

3150 void
3151 freerange(void *vstart, void *vend)
3152 {
3153     char *p;
3154     p = (char*)PGROUNDUP((uint)vstart);
3155     for(; p + PGSIZE <= (char*)vend; p += PGSIZE)
3156         kfree(p);
3157 }
3158
3159
3160 // Free the page of physical memory pointed at by v,
3161 // which normally should have been returned by a
3162 // call to kalloc(). (The exception is when
3163 // initializing the allocator; see kinit above.)
3164 void
3165 kfree(char *v)
3166 {
3167     struct run *r;
3168
3169     if((uint)v % PGSIZE || v < end || V2P(v) >= PHYSTOP)
3170         panic("kfree");
3171
3172     // Fill with junk to catch dangling refs.
3173     memset(v, 1, PGSIZE);
3174
3175     if(kmem.use_lock)
3176         acquire(&kmem.lock);
3177     r = (struct run*)v;
3178     r->next = kmem.freelist;
3179     kmem.freelist = r;
3180     if(kmem.use_lock)
3181         release(&kmem.lock);
3182 }
3183
3184 // Allocate one 4096-byte page of physical memory.
3185 // Returns a pointer that the kernel can use.
3186 // Returns 0 if the memory cannot be allocated.
3187 char*
3188 kalloc(void)
3189 {
3190     struct run *r;
3191
3192     if(kmem.use_lock)
3193         acquire(&kmem.lock);
3194     r = kmem.freelist;
3195     if(r)
3196         kmem.freelist = r->next;
3197     if(kmem.use_lock)
3198         release(&kmem.lock);
3199     return (char*)r;

```

```

3200 }
3201
3202
3203
3204
3205
3206
3207
3208
3209
3210
3211
3212
3213
3214
3215
3216
3217
3218
3219
3220
3221
3222
3223
3224
3225
3226
3227
3228
3229
3230
3231
3232
3233
3234
3235
3236
3237
3238
3239
3240
3241
3242
3243
3244
3245
3246
3247
3248
3249

```



```

3250 // x86 trap and interrupt constants.
3251
3252 // Processor-defined:
3253 #define T_DIVIDE      0    // divide error
3254 #define T_DEBUG      1    // debug exception
3255 #define T_NMI        2    // non-maskable interrupt
3256 #define T_BRKPT      3    // breakpoint
3257 #define T_OFLOW      4    // overflow
3258 #define T_BOUND      5    // bounds check
3259 #define T_ILLOP      6    // illegal opcode
3260 #define T_DEVICE      7    // device not available
3261 #define T_DBLFLT     8    // double fault
3262 // #define T_COPROC   9    // reserved (not used since 486)
3263 #define T_TSS        10   // invalid task switch segment
3264 #define T_SEGNP     11   // segment not present
3265 #define T_STACK     12   // stack exception
3266 #define T_GPFLT     13   // general protection fault
3267 #define T_PGFLT     14   // page fault
3268 // #define T_RES      15   // reserved
3269 #define T_FPEERR    16   // floating point error
3270 #define T_ALIGN     17   // alignment check
3271 #define T_MCHK      18   // machine check
3272 #define T_SIMDERR   19   // SIMD floating point error
3273
3274 // These are arbitrarily chosen, but with care not to overlap
3275 // processor defined exceptions or interrupt vectors.
3276 #define T_SYSCALL    64   // system call
3277 #define T_DEFAULT    500  // catchall
3278
3279 #define T_IRQ0       32   // IRQ 0 corresponds to int T_IRQ
3280
3281 #define IRQ_TIMER    0
3282 #define IRQ_KBD      1
3283 #define IRQ_COM1     4
3284 #define IRQ_IDE     14
3285 #define IRQ_ERROR    19
3286 #define IRQ_SPURIOUS 31
3287
3288
3289
3290
3291
3292
3293
3294
3295
3296
3297
3298
3299

```

```

3300 #!/usr/bin/perl -w
3301
3302 # Generate vectors.S, the trap/interrupt entry points.
3303 # There has to be one entry point per interrupt number
3304 # since otherwise there's no way for trap() to discover
3305 # the interrupt number.
3306
3307 print "# generated by vectors.pl - do not edit\n";
3308 print "# handlers\n";
3309 print ".globl alltraps\n";
3310 for(my $i = 0; $i < 256; $i++){
3311     print ".globl vector$i\n";
3312     print "vector$i:\n";
3313     if(!($i == 8 || ($i >= 10 && $i <= 14) || $i == 17)){
3314         print "    pushl \\\$0\n";
3315     }
3316     print "    pushl \\\$i\n";
3317     print "    jmp alltraps\n";
3318 }
3319
3320 print "\n# vector table\n";
3321 print ".data\n";
3322 print ".globl vectors\n";
3323 print "vectors:\n";
3324 for(my $i = 0; $i < 256; $i++){
3325     print "    .long vector$i\n";
3326 }
3327
3328 # sample output:
3329 # # handlers
3330 # .globl alltraps
3331 # .globl vector0
3332 # vector0:
3333 #     pushl $0
3334 #     pushl $0
3335 #     jmp alltraps
3336 # ...
3337 #
3338 # # vector table
3339 # .data
3340 # .globl vectors
3341 # vectors:
3342 #     .long vector0
3343 #     .long vector1
3344 #     .long vector2
3345 # ...
3346
3347
3348
3349

```

```

3350 #include "mmu.h"
3351
3352 # vectors.S sends all traps here.
3353 .globl alltraps
3354 alltraps:
3355 # Build trap frame.
3356 pushl %ds
3357 pushl %es
3358 pushl %fs
3359 pushl %gs
3360 pushal
3361
3362 # Set up data segments.
3363 movw $(SEG_KDATA<<3), %ax
3364 movw %ax, %ds
3365 movw %ax, %es
3366
3367 # Call trap(tf), where tf=%esp
3368 pushl %esp
3369 call trap
3370 addl $4, %esp
3371
3372 # Return falls through to trapret...
3373 .globl trapret
3374 trapret:
3375 popal
3376 popl %gs
3377 popl %fs
3378 popl %es
3379 popl %ds
3380 addl $0x8, %esp # trapno and errcode
3381 iret
3382
3383
3384
3385
3386
3387
3388
3389
3390
3391
3392
3393
3394
3395
3396
3397
3398
3399

```

```

3400 #include "types.h"
3401 #include "defs.h"
3402 #include "param.h"
3403 #include "memlayout.h"
3404 #include "mmu.h"
3405 #include "proc.h"
3406 #include "x86.h"
3407 #include "traps.h"
3408 #include "spinlock.h"
3409
3410 // Interrupt descriptor table (shared by all CPUs).
3411 struct gatedesc idt[256];
3412 extern uint vectors[]; // in vectors.S: array of 256 entry pointers
3413 struct spinlock tickslock;
3414 uint ticks;
3415
3416 void
3417 tvinit(void)
3418 {
3419     int i;
3420
3421     for(i = 0; i < 256; i++)
3422         SETGATE(idt[i], 0, SEG_KCODE<<3, vectors[i], 0);
3423     SETGATE(idt[T_SYSCALL], 1, SEG_KCODE<<3, vectors[T_SYSCALL], DPL_USER);
3424
3425     initlock(&tickslock, "time");
3426 }
3427
3428 void
3429 idtinit(void)
3430 {
3431     lidt(idt, sizeof(idt));
3432 }
3433
3434
3435
3436
3437
3438
3439
3440
3441
3442
3443
3444
3445
3446
3447
3448
3449

```

```

3450 void
3451 trap(struct trapframe *tf)
3452 {
3453     if(tf->trapno == T_SYSCALL){
3454         if(myproc()->killed)
3455             exit();
3456         myproc()->tf = tf;
3457         syscall();
3458         if(myproc()->killed)
3459             exit();
3460         return;
3461     }
3462     switch(tf->trapno){
3463     case T_IRQ0 + IRQ_TIMER:
3464         if(cpuid() == 0){
3465             acquire(&tickslock);
3466             ticks++;
3467             wakeup(&ticks);
3468             release(&tickslock);
3469         }
3470         lapiceoi();
3471         break;
3472     case T_IRQ0 + IRQ_IDE:
3473         ideintr();
3474         lapiceoi();
3475         break;
3476     case T_IRQ0 + IRQ_IDE+1:
3477         // Bochs generates spurious IDE1 interrupts.
3478         break;
3479     case T_IRQ0 + IRQ_KBD:
3480         kbdintr();
3481         lapiceoi();
3482         break;
3483     case T_IRQ0 + IRQ_COM1:
3484         uartintr();
3485         lapiceoi();
3486         break;
3487     case T_IRQ0 + 7:
3488     case T_IRQ0 + IRQ_SPURIOUS:
3489         cprintf("cpu%d: spurious interrupt at %x:%x\n",
3490             cpuid(), tf->cs, tf->eip);
3491         lapiceoi();
3492         break;
3493     }
3494 }
3495
3496
3497
3498
3499

```

```

3500 default:
3501     if(myproc() == 0 || (tf->cs&3) == 0){
3502         // In kernel, it must be our mistake.
3503         cprintf("unexpected trap %d from cpu %d eip %x (cr2=0x%x)\n",
3504             tf->trapno, cpuid(), tf->eip, rcr2());
3505         panic("trap");
3506     }
3507     // In user space, assume process misbehaved.
3508     cprintf("pid %d %s: trap %d err %d on cpu %d "
3509         "eip 0x%x addr 0x%x--kill proc\n",
3510         myproc()->pid, myproc()->name, tf->trapno, tf->err, cpuid(), tf->
3511         rcr2());
3512     myproc()->killed = 1;
3513 }
3514
3515 // Force process exit if it has been killed and is in user space.
3516 // (If it is still executing in the kernel, let it keep running
3517 // until it gets to the regular system call return.)
3518 if(myproc() && myproc()->killed && (tf->cs&3) == DPL_USER)
3519     exit();
3520
3521 // Force process to give up CPU on clock tick.
3522 // If interrupts were on while locks held, would need to check nlock.
3523 if(myproc() && myproc()->state == RUNNING && tf->trapno == T_IRQ0+IRQ_TIMER)
3524     yield();
3525
3526 // Check if the process has been killed since we yielded
3527 if(myproc() && myproc()->killed && (tf->cs&3) == DPL_USER)
3528     exit();
3529 }
3530
3531
3532
3533
3534
3535
3536
3537
3538
3539
3540
3541
3542
3543
3544
3545
3546
3547
3548
3549

```

```

3550 // System call numbers
3551 #define SYS_fork 1
3552 #define SYS_exit 2
3553 #define SYS_wait 3
3554 #define SYS_pipe 4
3555 #define SYS_read 5
3556 #define SYS_kill 6
3557 #define SYS_exec 7
3558 #define SYS_fstat 8
3559 #define SYS_chdir 9
3560 #define SYS_dup 10
3561 #define SYS_getpid 11
3562 #define SYS_sbrk 12
3563 #define SYS_sleep 13
3564 #define SYS_uptime 14
3565 #define SYS_open 15
3566 #define SYS_write 16
3567 #define SYS_mknod 17
3568 #define SYS_unlink 18
3569 #define SYS_link 19
3570 #define SYS_mkdir 20
3571 #define SYS_close 21
3572
3573
3574
3575
3576
3577
3578
3579
3580
3581
3582
3583
3584
3585
3586
3587
3588
3589
3590
3591
3592
3593
3594
3595
3596
3597
3598
3599

```

```

3600 #include "types.h"
3601 #include "defs.h"
3602 #include "param.h"
3603 #include "memlayout.h"
3604 #include "mmu.h"
3605 #include "proc.h"
3606 #include "x86.h"
3607 #include "syscall.h"
3608
3609 // User code makes a system call with INT T_SYSCALL.
3610 // System call number in %eax.
3611 // Arguments on the stack, from the user call to the C
3612 // library system call function. The saved user %esp points
3613 // to a saved program counter, and then the first argument.
3614
3615 // Fetch the int at addr from the current process.
3616 int
3617 fetchint(uint addr, int *ip)
3618 {
3619     struct proc *curproc = myproc();
3620
3621     if(addr >= curproc->sz || addr+4 > curproc->sz)
3622         return -1;
3623     *ip = *(int*)(addr);
3624     return 0;
3625 }
3626
3627 // Fetch the nul-terminated string at addr from the current process.
3628 // Doesn't actually copy the string - just sets *pp to point at it.
3629 // Returns length of string, not including nul.
3630 int
3631 fetchstr(uint addr, char **pp)
3632 {
3633     char *s, *ep;
3634     struct proc *curproc = myproc();
3635
3636     if(addr >= curproc->sz)
3637         return -1;
3638     *pp = (char*)addr;
3639     ep = (char*)curproc->sz;
3640     for(s = *pp; s < ep; s++){
3641         if(*s == 0)
3642             return s - *pp;
3643     }
3644     return -1;
3645 }
3646
3647
3648
3649

```

```

3650 // Fetch the nth 32-bit system call argument.
3651 int
3652 argint(int n, int *ip)
3653 {
3654     return fetchint((myproc()->tf->esp) + 4 + 4*n, ip);
3655 }
3656
3657 // Fetch the nth word-sized system call argument as a pointer
3658 // to a block of memory of size bytes. Check that the pointer
3659 // lies within the process address space.
3660 int
3661 argptr(int n, char **pp, int size)
3662 {
3663     int i;
3664     struct proc *curproc = myproc();
3665
3666     if(argint(n, &i) < 0)
3667         return -1;
3668     if(size < 0 || (uint)i >= curproc->sz || (uint)i+size > curproc->sz)
3669         return -1;
3670     *pp = (char*)i;
3671     return 0;
3672 }
3673
3674 // Fetch the nth word-sized system call argument as a string pointer.
3675 // Check that the pointer is valid and the string is nul-terminated.
3676 // (There is no shared writable memory, so the string can't change
3677 // between this check and being used by the kernel.)
3678 int
3679 argstr(int n, char **pp)
3680 {
3681     int addr;
3682     if(argint(n, &addr) < 0)
3683         return -1;
3684     return fetchstr(addr, pp);
3685 }
3686
3687
3688
3689
3690
3691
3692
3693
3694
3695
3696
3697
3698
3699

```

```

3700 extern int sys_chdir(void);
3701 extern int sys_close(void);
3702 extern int sys_dup(void);
3703 extern int sys_exec(void);
3704 extern int sys_exit(void);
3705 extern int sys_fork(void);
3706 extern int sys_fstat(void);
3707 extern int sys_getpid(void);
3708 extern int sys_kill(void);
3709 extern int sys_link(void);
3710 extern int sys_mkdir(void);
3711 extern int sys_mknod(void);
3712 extern int sys_open(void);
3713 extern int sys_pipe(void);
3714 extern int sys_read(void);
3715 extern int sys_sbrk(void);
3716 extern int sys_sleep(void);
3717 extern int sys_unlink(void);
3718 extern int sys_wait(void);
3719 extern int sys_write(void);
3720 extern int sys_uptime(void);
3721
3722 static int (*syscalls[])(void) = {
3723     [SYS_fork]    sys_fork,
3724     [SYS_exit]   sys_exit,
3725     [SYS_wait]   sys_wait,
3726     [SYS_pipe]   sys_pipe,
3727     [SYS_read]   sys_read,
3728     [SYS_kill]   sys_kill,
3729     [SYS_exec]   sys_exec,
3730     [SYS_fstat]  sys_fstat,
3731     [SYS_chdir]  sys_chdir,
3732     [SYS_dup]    sys_dup,
3733     [SYS_getpid] sys_getpid,
3734     [SYS_sbrk]   sys_sbrk,
3735     [SYS_sleep]  sys_sleep,
3736     [SYS_uptime] sys_uptime,
3737     [SYS_open]   sys_open,
3738     [SYS_write]  sys_write,
3739     [SYS_mknod]  sys_mknod,
3740     [SYS_unlink] sys_unlink,
3741     [SYS_link]   sys_link,
3742     [SYS_mkdir]  sys_mkdir,
3743     [SYS_close]  sys_close,
3744 };
3745
3746
3747
3748
3749

```

```

3750 void
3751 syscall(void)
3752 {
3753     int num;
3754     struct proc *curproc = myproc();
3755
3756     num = curproc->tf->eax;
3757     if(num > 0 && num < NELEM(syscalls) && syscalls[num]) {
3758         curproc->tf->eax = syscalls[num];
3759     } else {
3760         cprintf("%d %s: unknown sys call %d\n",
3761             curproc->pid, curproc->name, num);
3762         curproc->tf->eax = -1;
3763     }
3764 }
3765
3766
3767
3768
3769
3770
3771
3772
3773
3774
3775
3776
3777
3778
3779
3780
3781
3782
3783
3784
3785
3786
3787
3788
3789
3790
3791
3792
3793
3794
3795
3796
3797
3798
3799

```

```

3800 #include "types.h"
3801 #include "x86.h"
3802 #include "defs.h"
3803 #include "date.h"
3804 #include "param.h"
3805 #include "memlayout.h"
3806 #include "mmu.h"
3807 #include "proc.h"
3808
3809 int
3810 sys_fork(void)
3811 {
3812     return fork();
3813 }
3814
3815 int
3816 sys_exit(void)
3817 {
3818     exit();
3819     return 0; // not reached
3820 }
3821
3822 int
3823 sys_wait(void)
3824 {
3825     return wait();
3826 }
3827
3828 int
3829 sys_kill(void)
3830 {
3831     int pid;
3832
3833     if(argint(0, &pid) < 0)
3834         return -1;
3835     return kill(pid);
3836 }
3837
3838 int
3839 sys_getpid(void)
3840 {
3841     return myproc()->pid;
3842 }
3843
3844
3845
3846
3847
3848
3849

```

```

3850 int
3851 sys_sbrk(void)
3852 {
3853     int addr;
3854     int n;
3855
3856     if(argint(0, &n) < 0)
3857         return -1;
3858     addr = myproc()->sz;
3859     if(growproc(n) < 0)
3860         return -1;
3861     return addr;
3862 }
3863
3864 int
3865 sys_sleep(void)
3866 {
3867     int n;
3868     uint ticks0;
3869
3870     if(argint(0, &n) < 0)
3871         return -1;
3872     acquire(&tickslock);
3873     ticks0 = ticks;
3874     while(ticks - ticks0 < n){
3875         if(myproc()->killed){
3876             release(&tickslock);
3877             return -1;
3878         }
3879         sleep(&ticks, &tickslock);
3880     }
3881     release(&tickslock);
3882     return 0;
3883 }
3884
3885 // return how many clock tick interrupts have occurred
3886 // since start.
3887 int
3888 sys_uptime(void)
3889 {
3890     uint xticks;
3891
3892     acquire(&tickslock);
3893     xticks = ticks;
3894     release(&tickslock);
3895     return xticks;
3896 }
3897
3898
3899

```

```

3900 struct buf {
3901     int flags;
3902     uint dev;
3903     uint blockno;
3904     struct sleeplock lock;
3905     uint refcnt;
3906     struct buf *prev; // LRU cache list
3907     struct buf *next;
3908     struct buf *qnext; // disk queue
3909     uchar data[BSIZE];
3910 };
3911 #define B_VALID 0x2 // buffer has been read from disk
3912 #define B_DIRTY 0x4 // buffer needs to be written to disk
3913
3914
3915
3916
3917
3918
3919
3920
3921
3922
3923
3924
3925
3926
3927
3928
3929
3930
3931
3932
3933
3934
3935
3936
3937
3938
3939
3940
3941
3942
3943
3944
3945
3946
3947
3948
3949

```

```
3950 // Long-term locks for processes
3951 struct sleeplock {
3952     uint locked;           // Is the lock held?
3953     struct spinlock lk;   // spinlock protecting this sleep lock
3954
3955     // For debugging:
3956     char *name;           // Name of lock.
3957     int pid;              // Process holding lock
3958 };
3959
3960
3961
3962
3963
3964
3965
3966
3967
3968
3969
3970
3971
3972
3973
3974
3975
3976
3977
3978
3979
3980
3981
3982
3983
3984
3985
3986
3987
3988
3989
3990
3991
3992
3993
3994
3995
3996
3997
3998
3999
```

```
4000 #define O_RDONLY 0x000
4001 #define O_WRONLY 0x001
4002 #define O_RDWR 0x002
4003 #define O_CREATE 0x200
4004
4005
4006
4007
4008
4009
4010
4011
4012
4013
4014
4015
4016
4017
4018
4019
4020
4021
4022
4023
4024
4025
4026
4027
4028
4029
4030
4031
4032
4033
4034
4035
4036
4037
4038
4039
4040
4041
4042
4043
4044
4045
4046
4047
4048
4049
```



```

4050 #define T_DIR 1 // Directory
4051 #define T_FILE 2 // File
4052 #define T_DEV 3 // Device
4053
4054 struct stat {
4055     short type; // Type of file
4056     int dev; // File system's disk device
4057     uint ino; // Inode number
4058     short nlink; // Number of links to file
4059     uint size; // Size of file in bytes
4060 };
4061
4062
4063
4064
4065
4066
4067
4068
4069
4070
4071
4072
4073
4074
4075
4076
4077
4078
4079
4080
4081
4082
4083
4084
4085
4086
4087
4088
4089
4090
4091
4092
4093
4094
4095
4096
4097
4098
4099

```

```

4100 // On-disk file system format.
4101 // Both the kernel and user programs use this header file.
4102
4103
4104 #define ROOTINO 1 // root i-number
4105 #define BSIZE 512 // block size
4106
4107 // Disk layout:
4108 // [ boot block | super block | log | inode blocks |
4109 //                               free bit map | data blocks]
4110 //
4111 // mkfs computes the super block and builds an initial file system. The
4112 // super block describes the disk layout:
4113 struct superblock {
4114     uint size; // Size of file system image (blocks)
4115     uint nblocks; // Number of data blocks
4116     uint ninodes; // Number of inodes.
4117     uint nlog; // Number of log blocks
4118     uint logstart; // Block number of first log block
4119     uint inodestart; // Block number of first inode block
4120     uint bmapstart; // Block number of first free map block
4121 };
4122
4123 #define NDIRECT 12
4124 #define NINDIRECT (BSIZE / sizeof(uint))
4125 #define MAXFILE (NDIRECT + NINDIRECT)
4126
4127 // On-disk inode structure
4128 struct dinode {
4129     short type; // File type
4130     short major; // Major device number (T_DEV only)
4131     short minor; // Minor device number (T_DEV only)
4132     short nlink; // Number of links to inode in file system
4133     uint size; // Size of file (bytes)
4134     uint addrs[NDIRECT+1]; // Data block addresses
4135 };
4136
4137
4138
4139
4140
4141
4142
4143
4144
4145
4146
4147
4148
4149

```

```

4150 // Inodes per block.
4151 #define IPB          (BSIZE / sizeof(struct dinode))
4152
4153 // Block containing inode i
4154 #define IBLOCK(i, sb)  ((i) / IPB + sb.inodestart)
4155
4156 // Bitmap bits per block
4157 #define BPB          (BSIZE*8)
4158
4159 // Block of free map containing bit for block b
4160 #define BBLOCK(b, sb) (b/BPB + sb.bmapstart)
4161
4162 // Directory is a file containing a sequence of dirent structures.
4163 #define DIRSIZ 14
4164
4165 struct dirent {
4166     ushort inum;
4167     char name[DIRSIZ];
4168 };
4169
4170
4171
4172
4173
4174
4175
4176
4177
4178
4179
4180
4181
4182
4183
4184
4185
4186
4187
4188
4189
4190
4191
4192
4193
4194
4195
4196
4197
4198
4199

```

```

4200 struct file {
4201     enum { FD_NONE, FD_PIPE, FD_INODE } type;
4202     int ref; // reference count
4203     char readable;
4204     char writable;
4205     struct pipe *pipe;
4206     struct inode *ip;
4207     uint off;
4208 };
4209
4210
4211 // in-memory copy of an inode
4212 struct inode {
4213     uint dev;           // Device number
4214     uint inum;         // Inode number
4215     int ref;           // Reference count
4216     struct sleeplock lock; // protects everything below here
4217     int valid;         // inode has been read from disk?
4218
4219     short type;        // copy of disk inode
4220     short major;
4221     short minor;
4222     short nlink;
4223     uint size;
4224     uint addrs[NDIRECT+1];
4225 };
4226
4227 // table mapping major device number to
4228 // device functions
4229 struct devsw {
4230     int (*read)(struct inode*, char*, int);
4231     int (*write)(struct inode*, char*, int);
4232 };
4233
4234 extern struct devsw devsw[];
4235
4236 #define CONSOLE 1
4237
4238
4239
4240
4241
4242
4243
4244
4245
4246
4247
4248
4249

```

```

4250 // Blank page.
4251
4252
4253
4254
4255
4256
4257
4258
4259
4260
4261
4262
4263
4264
4265
4266
4267
4268
4269
4270
4271
4272
4273
4274
4275
4276
4277
4278
4279
4280
4281
4282
4283
4284
4285
4286
4287
4288
4289
4290
4291
4292
4293
4294
4295
4296
4297
4298
4299

```

```

4300 // Simple PIO-based (non-DMA) IDE driver code.
4301
4302 #include "types.h"
4303 #include "defs.h"
4304 #include "param.h"
4305 #include "memlayout.h"
4306 #include "mmu.h"
4307 #include "proc.h"
4308 #include "x86.h"
4309 #include "traps.h"
4310 #include "spinlock.h"
4311 #include "sleeplock.h"
4312 #include "fs.h"
4313 #include "buf.h"
4314
4315 #define SECTOR_SIZE 512
4316 #define IDE_BSY 0x80
4317 #define IDE_DRDY 0x40
4318 #define IDE_DF 0x20
4319 #define IDE_ERR 0x01
4320
4321 #define IDE_CMD_READ 0x20
4322 #define IDE_CMD_WRITE 0x30
4323 #define IDE_CMD_RDMUL 0xc4
4324 #define IDE_CMD_WRMUL 0xc5
4325
4326 // idequeue points to the buf now being read/written to the disk.
4327 // idequeue->qnext points to the next buf to be processed.
4328 // You must hold idelock while manipulating queue.
4329
4330 static struct spinlock idelock;
4331 static struct buf *idequeue;
4332
4333 static int havedisk1;
4334 static void idestart(struct buf*);
4335
4336 // Wait for IDE disk to become ready.
4337 static int
4338 idewait(int checkerr)
4339 {
4340     int r;
4341
4342     while(((r = inb(0x1f7)) & (IDE_BSY|IDE_DRDY)) != IDE_DRDY)
4343         ;
4344     if(checkerr && (r & (IDE_DF|IDE_ERR)) != 0)
4345         return -1;
4346     return 0;
4347 }
4348
4349

```

```

4350 void
4351 ideinit(void)
4352 {
4353     int i;
4354
4355     initlock(&idelock, "ide");
4356     ioapicenable(IRQ_IDE, ncpu - 1);
4357     idewait(0);
4358
4359     // Check if disk 1 is present
4360     outb(0x1f6, 0xe0 | (1<<4));
4361     for(i=0; i<1000; i++){
4362         if(inb(0x1f7) != 0){
4363             havedisk1 = 1;
4364             break;
4365         }
4366     }
4367
4368     // Switch back to disk 0.
4369     outb(0x1f6, 0xe0 | (0<<4));
4370 }
4371
4372 // Start the request for b. Caller must hold idelock.
4373 static void
4374 idestart(struct buf *b)
4375 {
4376     if(b == 0)
4377         panic("idestart");
4378     if(b->blockno >= FSSIZE)
4379         panic("incorrect blockno");
4380     int sector_per_block = BSIZE/SECTOR_SIZE;
4381     int sector = b->blockno * sector_per_block;
4382     int read_cmd = (sector_per_block == 1) ? IDE_CMD_READ : IDE_CMD_RDMDL;
4383     int write_cmd = (sector_per_block == 1) ? IDE_CMD_WRITE : IDE_CMD_WRMDL;
4384
4385     if (sector_per_block > 7) panic("idestart");
4386
4387     idewait(0);
4388     outb(0x3f6, 0); // generate interrupt
4389     outb(0x1f2, sector_per_block); // number of sectors
4390     outb(0x1f3, sector & 0xff);
4391     outb(0x1f4, (sector >> 8) & 0xff);
4392     outb(0x1f5, (sector >> 16) & 0xff);
4393     outb(0x1f6, 0xe0 | ((b->dev&1)<<4) | ((sector>>24)&0x0f));
4394     if(b->flags & B_DIRTY){
4395         outb(0x1f7, write_cmd);
4396         outsl(0x1f0, b->data, BSIZE/4);
4397     } else {
4398         outb(0x1f7, read_cmd);
4399     }

```

```

4400 }
4401
4402 // Interrupt handler.
4403 void
4404 ideintr(void)
4405 {
4406     struct buf *b;
4407
4408     // First queued buffer is the active request.
4409     acquire(&idelock);
4410
4411     if((b = idequeue) == 0){
4412         release(&idelock);
4413         return;
4414     }
4415     idequeue = b->qnext;
4416
4417     // Read data if needed.
4418     if(!(b->flags & B_DIRTY) && idewait(1) >= 0)
4419         insl(0x1f0, b->data, BSIZE/4);
4420
4421     // Wake process waiting for this buf.
4422     b->flags |= B_VALID;
4423     b->flags &= ~B_DIRTY;
4424     wakeup(b);
4425
4426     // Start disk on next buf in queue.
4427     if(idequeue != 0)
4428         idestart(idequeue);
4429
4430     release(&idelock);
4431 }
4432
4433
4434
4435
4436
4437
4438
4439
4440
4441
4442
4443
4444
4445
4446
4447
4448
4449

```

```

4450 // Sync buf with disk.
4451 // If B_DIRTY is set, write buf to disk, clear B_DIRTY, set B_VALID.
4452 // Else if B_VALID is not set, read buf from disk, set B_VALID.
4453 void
4454 iderw(struct buf *b)
4455 {
4456     struct buf **pp;
4457
4458     if(!holdingsleep(&b->lock))
4459         panic("iderw: buf not locked");
4460     if((b->flags & (B_VALID|B_DIRTY)) == B_VALID)
4461         panic("iderw: nothing to do");
4462     if(b->dev != 0 && !havedisk1)
4463         panic("iderw: ide disk 1 not present");
4464
4465     acquire(&idelock);
4466
4467     // Append b to idequeue.
4468     b->qnext = 0;
4469     for(pp=&idequeue; *pp; pp=&(*pp)->qnext)
4470         ;
4471     *pp = b;
4472
4473     // Start disk if necessary.
4474     if(idequeue == b)
4475         idestart(b);
4476
4477     // Wait for request to finish.
4478     while((b->flags & (B_VALID|B_DIRTY)) != B_VALID){
4479         sleep(b, &idelock);
4480     }
4481
4482     release(&idelock);
4483 }
4484
4485
4486
4487
4488
4489
4490
4491
4492
4493
4494
4495
4496
4497
4498
4499

```

```

4500 // Buffer cache.
4501 //
4502 // The buffer cache is a linked list of buf structures holding
4503 // cached copies of disk block contents. Caching disk blocks
4504 // in memory reduces the number of disk reads and also provides
4505 // a synchronization point for disk blocks used by multiple processes.
4506 //
4507 // Interface:
4508 // * To get a buffer for a particular disk block, call bread.
4509 // * After changing buffer data, call bwrite to write it to disk.
4510 // * When done with the buffer, call brelse.
4511 // * Do not use the buffer after calling brelse.
4512 // * Only one process at a time can use a buffer,
4513 //   so do not keep them longer than necessary.
4514 //
4515 // The implementation uses two state flags internally:
4516 // * B_VALID: the buffer data has been read from the disk.
4517 // * B_DIRTY: the buffer data has been modified
4518 //   and needs to be written to disk.
4519
4520 #include "types.h"
4521 #include "defs.h"
4522 #include "param.h"
4523 #include "spinlock.h"
4524 #include "sleeplock.h"
4525 #include "fs.h"
4526 #include "buf.h"
4527
4528 struct {
4529     struct spinlock lock;
4530     struct buf buf[NBUF];
4531
4532     // Linked list of all buffers, through prev/next.
4533     // head.next is most recently used.
4534     struct buf head;
4535 } bcache;
4536
4537 void
4538 binit(void)
4539 {
4540     struct buf *b;
4541
4542     initlock(&bcache.lock, "bcache");
4543
4544
4545
4546
4547
4548
4549

```

```

4550 // Create linked list of buffers
4551 bcache.head.prev = &bcache.head;
4552 bcache.head.next = &bcache.head;
4553 for(b = bcache.buf; b < bcache.buf+NBUF; b++){
4554     b->next = bcache.head.next;
4555     b->prev = &bcache.head;
4556     initsleeplock(&b->lock, "buffer");
4557     bcache.head.next->prev = b;
4558     bcache.head.next = b;
4559 }
4560 }
4561
4562 // Look through buffer cache for block on device dev.
4563 // If not found, allocate a buffer.
4564 // In either case, return locked buffer.
4565 static struct buf*
4566 bget(uint dev, uint blockno)
4567 {
4568     struct buf *b;
4569
4570     acquire(&bcache.lock);
4571
4572     // Is the block already cached?
4573     for(b = bcache.head.next; b != &bcache.head; b = b->next){
4574         if(b->dev == dev && b->blockno == blockno){
4575             b->refcnt++;
4576             release(&bcache.lock);
4577             acquiresleep(&b->lock);
4578             return b;
4579         }
4580     }
4581
4582     // Not cached; recycle an unused buffer.
4583     // Even if refcnt==0, B_DIRTY indicates a buffer is in use
4584     // because log.c has modified it but not yet committed it.
4585     for(b = bcache.head.prev; b != &bcache.head; b = b->prev){
4586         if(b->refcnt == 0 && (b->flags & B_DIRTY) == 0) {
4587             b->dev = dev;
4588             b->blockno = blockno;
4589             b->flags = 0;
4590             b->refcnt = 1;
4591             release(&bcache.lock);
4592             acquiresleep(&b->lock);
4593             return b;
4594         }
4595     }
4596     panic("bget: no buffers");
4597 }
4598
4599

```

```

4600 // Return a locked buf with the contents of the indicated block.
4601 struct buf*
4602 bread(uint dev, uint blockno)
4603 {
4604     struct buf *b;
4605
4606     b = bget(dev, blockno);
4607     if((b->flags & B_VALID) == 0) {
4608         iderw(b);
4609     }
4610     return b;
4611 }
4612
4613 // Write b's contents to disk. Must be locked.
4614 void
4615 bwrite(struct buf *b)
4616 {
4617     if(!holdingsleep(&b->lock))
4618         panic("bwrite");
4619     b->flags |= B_DIRTY;
4620     iderw(b);
4621 }
4622
4623 // Release a locked buffer.
4624 // Move to the head of the MRU list.
4625 void
4626 brelse(struct buf *b)
4627 {
4628     if(!holdingsleep(&b->lock))
4629         panic("brelse");
4630
4631     releasesleep(&b->lock);
4632
4633     acquire(&bcache.lock);
4634     b->refcnt--;
4635     if (b->refcnt == 0) {
4636         // no one is waiting for it.
4637         b->next->prev = b->prev;
4638         b->prev->next = b->next;
4639         b->next = bcache.head.next;
4640         b->prev = &bcache.head;
4641         bcache.head.next->prev = b;
4642         bcache.head.next = b;
4643     }
4644
4645     release(&bcache.lock);
4646 }
4647
4648
4649

```

```
4650 // Blank page.
4651
4652
4653
4654
4655
4656
4657
4658
4659
4660
4661
4662
4663
4664
4665
4666
4667
4668
4669
4670
4671
4672
4673
4674
4675
4676
4677
4678
4679
4680
4681
4682
4683
4684
4685
4686
4687
4688
4689
4690
4691
4692
4693
4694
4695
4696
4697
4698
4699
```

```
4700 // Sleeping locks
4701
4702 #include "types.h"
4703 #include "defs.h"
4704 #include "param.h"
4705 #include "x86.h"
4706 #include "memlayout.h"
4707 #include "mmu.h"
4708 #include "proc.h"
4709 #include "spinlock.h"
4710 #include "sleeplock.h"
4711
4712 void
4713 initsleeplock(struct sleeplock *lk, char *name)
4714 {
4715     initlock(&lk->lk, "sleep lock");
4716     lk->name = name;
4717     lk->locked = 0;
4718     lk->pid = 0;
4719 }
4720
4721 void
4722 acquiresleep(struct sleeplock *lk)
4723 {
4724     acquire(&lk->lk);
4725     while (lk->locked) {
4726         sleep(lk, &lk->lk);
4727     }
4728     lk->locked = 1;
4729     lk->pid = myproc()->pid;
4730     release(&lk->lk);
4731 }
4732
4733 void
4734 releasesleep(struct sleeplock *lk)
4735 {
4736     acquire(&lk->lk);
4737     lk->locked = 0;
4738     lk->pid = 0;
4739     wakeup(lk);
4740     release(&lk->lk);
4741 }
4742
4743
4744
4745
4746
4747
4748
4749
```

```

4750 int
4751 holdingsleep(struct sleeplock *lk)
4752 {
4753     int r;
4754
4755     acquire(&lk->lk);
4756     r = lk->locked;
4757     release(&lk->lk);
4758     return r;
4759 }
4760
4761
4762
4763
4764
4765
4766
4767
4768
4769
4770
4771
4772
4773
4774
4775
4776
4777
4778
4779
4780
4781
4782
4783
4784
4785
4786
4787
4788
4789
4790
4791
4792
4793
4794
4795
4796
4797
4798
4799

```

```

4800 #include "types.h"
4801 #include "defs.h"
4802 #include "param.h"
4803 #include "spinlock.h"
4804 #include "sleeplock.h"
4805 #include "fs.h"
4806 #include "buf.h"
4807
4808 // Simple logging that allows concurrent FS system calls.
4809 //
4810 // A log transaction contains the updates of multiple FS system
4811 // calls. The logging system only commits when there are
4812 // no FS system calls active. Thus there is never
4813 // any reasoning required about whether a commit might
4814 // write an uncommitted system call's updates to disk.
4815 //
4816 // A system call should call begin_op()/end_op() to mark
4817 // its start and end. Usually begin_op() just increments
4818 // the count of in-progress FS system calls and returns.
4819 // But if it thinks the log is close to running out, it
4820 // sleeps until the last outstanding end_op() commits.
4821 //
4822 // The log is a physical re-do log containing disk blocks.
4823 // The on-disk log format:
4824 //   header block, containing block #s for block A, B, C, ...
4825 //   block A
4826 //   block B
4827 //   block C
4828 //   ...
4829 // Log appends are synchronous.
4830
4831 // Contents of the header block, used for both the on-disk header block
4832 // and to keep track in memory of logged block# before commit.
4833 struct logheader {
4834     int n;
4835     int block[LOGSIZE];
4836 };
4837
4838 struct log {
4839     struct spinlock lock;
4840     int start;
4841     int size;
4842     int outstanding; // how many FS sys calls are executing.
4843     int committing; // in commit(), please wait.
4844     int dev;
4845     struct logheader lh;
4846 };
4847
4848
4849

```



```

4850 struct log log;
4851
4852 static void recover_from_log(void);
4853 static void commit();
4854
4855 void
4856 initlog(int dev)
4857 {
4858     if (sizeof(struct logheader) >= BSIZE)
4859         panic("initlog: too big logheader");
4860
4861     struct superblock sb;
4862     initlock(&log.lock, "log");
4863     readsb(dev, &sb);
4864     log.start = sb.logstart;
4865     log.size = sb.nlog;
4866     log.dev = dev;
4867     recover_from_log();
4868 }
4869
4870 // Copy committed blocks from log to their home location
4871 static void
4872 install_trans(void)
4873 {
4874     int tail;
4875
4876     for (tail = 0; tail < log.lh.n; tail++) {
4877         struct buf *lbuf = bread(log.dev, log.start+tail+1); // read log block
4878         struct buf *dbuf = bread(log.dev, log.lh.block[tail]); // read dst
4879         memmove(dbuf->data, lbuf->data, BSIZE); // copy block to dst
4880         bwrite(dbuf); // write dst to disk
4881         brelse(lbuf);
4882         brelse(dbuf);
4883     }
4884 }
4885
4886 // Read the log header from disk into the in-memory log header
4887 static void
4888 read_head(void)
4889 {
4890     struct buf *buf = bread(log.dev, log.start);
4891     struct logheader *lh = (struct logheader *) (buf->data);
4892     int i;
4893     log.lh.n = lh->n;
4894     for (i = 0; i < log.lh.n; i++) {
4895         log.lh.block[i] = lh->block[i];
4896     }
4897     brelse(buf);
4898 }
4899

```

```

4900 // Write in-memory log header to disk.
4901 // This is the true point at which the
4902 // current transaction commits.
4903 static void
4904 write_head(void)
4905 {
4906     struct buf *buf = bread(log.dev, log.start);
4907     struct logheader *hb = (struct logheader *) (buf->data);
4908     int i;
4909     hb->n = log.lh.n;
4910     for (i = 0; i < log.lh.n; i++) {
4911         hb->block[i] = log.lh.block[i];
4912     }
4913     bwrite(buf);
4914     brelse(buf);
4915 }
4916
4917 static void
4918 recover_from_log(void)
4919 {
4920     read_head();
4921     install_trans(); // if committed, copy from log to disk
4922     log.lh.n = 0;
4923     write_head(); // clear the log
4924 }
4925
4926 // called at the start of each FS system call.
4927 void
4928 begin_op(void)
4929 {
4930     acquire(&log.lock);
4931     while(1){
4932         if(log.committing){
4933             sleep(&log, &log.lock);
4934         } else if(log.lh.n + (log.outstanding+1)*MAXOPBLOCKS > LOGSIZE){
4935             // this op might exhaust log space; wait for commit.
4936             sleep(&log, &log.lock);
4937         } else {
4938             log.outstanding += 1;
4939             release(&log.lock);
4940             break;
4941         }
4942     }
4943 }
4944
4945
4946
4947
4948
4949

```

```

4950 // called at the end of each FS system call.
4951 // commits if this was the last outstanding operation.
4952 void
4953 end_op(void)
4954 {
4955     int do_commit = 0;
4956
4957     acquire(&log.lock);
4958     log.outstanding -= 1;
4959     if(log.committing)
4960         panic("log.committing");
4961     if(log.outstanding == 0){
4962         do_commit = 1;
4963         log.committing = 1;
4964     } else {
4965         // begin_op() may be waiting for log space,
4966         // and decrementing log.outstanding has decreased
4967         // the amount of reserved space.
4968         wakeup(&log);
4969     }
4970     release(&log.lock);
4971
4972     if(do_commit){
4973         // call commit w/o holding locks, since not allowed
4974         // to sleep with locks.
4975         commit();
4976         acquire(&log.lock);
4977         log.committing = 0;
4978         wakeup(&log);
4979         release(&log.lock);
4980     }
4981 }
4982
4983 // Copy modified blocks from cache to log.
4984 static void
4985 write_log(void)
4986 {
4987     int tail;
4988
4989     for (tail = 0; tail < log.lh.n; tail++) {
4990         struct buf *to = bread(log.dev, log.start+tail+1); // log block
4991         struct buf *from = bread(log.dev, log.lh.block[tail]); // cache block
4992         memmove(to->data, from->data, BSIZE);
4993         bwrite(to); // write the log
4994         brelse(from);
4995         brelse(to);
4996     }
4997 }
4998
4999

```

```

5000 static void
5001 commit()
5002 {
5003     if (log.lh.n > 0) {
5004         write_log(); // Write modified blocks from cache to log
5005         write_head(); // Write header to disk -- the real commit
5006         install_trans(); // Now install writes to home locations
5007         log.lh.n = 0;
5008         write_head(); // Erase the transaction from the log
5009     }
5010 }
5011
5012 // Caller has modified b->data and is done with the buffer.
5013 // Record the block number and pin in the cache with B_DIRTY.
5014 // commit()/write_log() will do the disk write.
5015 //
5016 // log_write() replaces bwrite(); a typical use is:
5017 // bp = bread(...)
5018 // modify bp->data[]
5019 // log_write(bp)
5020 // brelse(bp)
5021 void
5022 log_write(struct buf *b)
5023 {
5024     int i;
5025
5026     if (log.lh.n >= LOGSIZE || log.lh.n >= log.size - 1)
5027         panic("too big a transaction");
5028     if (log.outstanding < 1)
5029         panic("log_write outside of trans");
5030
5031     acquire(&log.lock);
5032     for (i = 0; i < log.lh.n; i++) {
5033         if (log.lh.block[i] == b->blockno) // log absorption
5034             break;
5035     }
5036     log.lh.block[i] = b->blockno;
5037     if (i == log.lh.n)
5038         log.lh.n++;
5039     b->flags |= B_DIRTY; // prevent eviction
5040     release(&log.lock);
5041 }
5042
5043
5044
5045
5046
5047
5048
5049

```

```

5050 // File system implementation. Five layers:
5051 //   + Blocks: allocator for raw disk blocks.
5052 //   + Log: crash recovery for multi-step updates.
5053 //   + Files: inode allocator, reading, writing, metadata.
5054 //   + Directories: inode with special contents (list of other inodes!)
5055 //   + Names: paths like /usr/rtn/xv6/fs.c for convenient naming.
5056 //
5057 // This file contains the low-level file system manipulation
5058 // routines. The (higher-level) system call implementations
5059 // are in sysfile.c.
5060
5061 #include "types.h"
5062 #include "defs.h"
5063 #include "param.h"
5064 #include "stat.h"
5065 #include "mmu.h"
5066 #include "proc.h"
5067 #include "spinlock.h"
5068 #include "sleeplock.h"
5069 #include "fs.h"
5070 #include "buf.h"
5071 #include "file.h"
5072
5073 #define min(a, b) ((a) < (b) ? (a) : (b))
5074 static void itrunc(struct inode*);
5075 // there should be one superblock per disk device, but we run with
5076 // only one device
5077 struct superblock sb;
5078
5079 // Read the super block.
5080 void
5081 readsb(int dev, struct superblock *sb)
5082 {
5083   struct buf *bp;
5084
5085   bp = bread(dev, 1);
5086   memmove(sb, bp->data, sizeof(*sb));
5087   brelse(bp);
5088 }
5089
5090
5091
5092
5093
5094
5095
5096
5097
5098
5099

```

```

5100 // Zero a block.
5101 static void
5102 bzero(int dev, int bno)
5103 {
5104   struct buf *bp;
5105
5106   bp = bread(dev, bno);
5107   memset(bp->data, 0, BSIZE);
5108   log_write(bp);
5109   brelse(bp);
5110 }
5111
5112 // Blocks.
5113
5114 // Allocate a zeroed disk block.
5115 static uint
5116 balloc(uint dev)
5117 {
5118   int b, bi, m;
5119   struct buf *bp;
5120
5121   bp = 0;
5122   for(b = 0; b < sb.size; b += BPB){
5123     bp = bread(dev, BBLOCK(b, sb));
5124     for(bi = 0; bi < BPB && b + bi < sb.size; bi++){
5125       m = 1 << (bi % 8);
5126       if((bp->data[bi/8] & m) == 0){ // Is block free?
5127         bp->data[bi/8] |= m; // Mark block in use.
5128         log_write(bp);
5129         brelse(bp);
5130         bzero(dev, b + bi);
5131         return b + bi;
5132       }
5133     }
5134     brelse(bp);
5135   }
5136   panic("balloc: out of blocks");
5137 }
5138
5139
5140
5141
5142
5143
5144
5145
5146
5147
5148
5149

```

```

5150 // Free a disk block.
5151 static void
5152 bfree(int dev, uint b)
5153 {
5154     struct buf *bp;
5155     int bi, m;
5156
5157     readsb(dev, &sb);
5158     bp = bread(dev, BBLOCK(b, sb));
5159     bi = b % BPB;
5160     m = 1 << (bi % 8);
5161     if((bp->data[bi/8] & m) == 0)
5162         panic("freeing free block");
5163     bp->data[bi/8] &= ~m;
5164     log_write(bp);
5165     brelse(bp);
5166 }
5167
5168 // Inodes.
5169 //
5170 // An inode describes a single unnamed file.
5171 // The inode disk structure holds metadata: the file's type,
5172 // its size, the number of links referring to it, and the
5173 // list of blocks holding the file's content.
5174 //
5175 // The inodes are laid out sequentially on disk at
5176 // sb.startinode. Each inode has a number, indicating its
5177 // position on the disk.
5178 //
5179 // The kernel keeps a cache of in-use inodes in memory
5180 // to provide a place for synchronizing access
5181 // to inodes used by multiple processes. The cached
5182 // inodes include book-keeping information that is
5183 // not stored on disk: ip->ref and ip->valid.
5184 //
5185 // An inode and its in-memory representation go through a
5186 // sequence of states before they can be used by the
5187 // rest of the file system code.
5188 //
5189 // * Allocation: an inode is allocated if its type (on disk)
5190 //   is non-zero. ialloc() allocates, and iput() frees if
5191 //   the reference and link counts have fallen to zero.
5192 //
5193 // * Referencing in cache: an entry in the inode cache
5194 //   is free if ip->ref is zero. Otherwise ip->ref tracks
5195 //   the number of in-memory pointers to the entry (open
5196 //   files and current directories). iget() finds or
5197 //   creates a cache entry and increments its ref; iput()
5198 //   decrements ref.
5199 //

```

```

5200 // * Valid: the information (type, size, &c) in an inode
5201 //   cache entry is only correct when ip->valid is 1.
5202 //   ilock() reads the inode from
5203 //   the disk and sets ip->valid, while iput() clears
5204 //   ip->valid if ip->ref has fallen to zero.
5205 //
5206 // * Locked: file system code may only examine and modify
5207 //   the information in an inode and its content if it
5208 //   has first locked the inode.
5209 //
5210 // Thus a typical sequence is:
5211 //   ip = iget(dev, inum)
5212 //   ilock(ip)
5213 //   ... examine and modify ip->xxx ...
5214 //   iunlock(ip)
5215 //   iput(ip)
5216 //
5217 // ilock() is separate from iget() so that system calls can
5218 // get a long-term reference to an inode (as for an open file)
5219 // and only lock it for short periods (e.g., in read()).
5220 // The separation also helps avoid deadlock and races during
5221 // pathname lookup. iget() increments ip->ref so that the inode
5222 // stays cached and pointers to it remain valid.
5223 //
5224 // Many internal file system functions expect the caller to
5225 // have locked the inodes involved; this lets callers create
5226 // multi-step atomic operations.
5227 //
5228 // The icache.lock spin-lock defends the allocation of icache
5229 // entries. Since ip->ref indicates whether an entry is free,
5230 // and ip->dev and ip->inum indicate which i-node an entry
5231 // holds, one must hold icache.lock while using any of those fields.
5232 //
5233 // An ip->lock sleep-lock defends all ip-> fields other than ref,
5234 // dev, and inum. One must hold ip->lock in order to
5235 // read or write that inode's ip->valid, ip->size, ip->type, &c.
5236
5237 struct {
5238     struct spinlock lock;
5239     struct inode inode[NINODE];
5240 } icache;
5241
5242 void
5243 iinit(int dev)
5244 {
5245     int i = 0;
5246
5247     initlock(&icache.lock, "icache");
5248     for(i = 0; i < NINODE; i++) {
5249         initsleeplock(&icache.inode[i].lock, "inode");

```

```

5250 }
5251
5252 readsb(dev, &sb);
5253 cprintf("sb: size %d nblocks %d ninodes %d nlog %d logstart %d\
5254 inodestart %d bmap start %d\n", sb.size, sb.nblocks,
5255         sb.ninodes, sb.nlog, sb.logstart, sb.inodestart,
5256         sb.bmapstart);
5257 }
5258
5259 static struct inode* iget(uint dev, uint inum);
5260
5261
5262
5263
5264
5265
5266
5267
5268
5269
5270
5271
5272
5273
5274
5275
5276
5277
5278
5279
5280
5281
5282
5283
5284
5285
5286
5287
5288
5289
5290
5291
5292
5293
5294
5295
5296
5297
5298
5299

```

```

5300 // Allocate an inode on device dev.
5301 // Mark it as allocated by giving it type type.
5302 // Returns an unlocked but allocated and referenced inode.
5303 struct inode*
5304 ialloc(uint dev, short type)
5305 {
5306     int inum;
5307     struct buf *bp;
5308     struct dinode *dip;
5309
5310     for(inum = 1; inum < sb.ninodes; inum++){
5311         bp = bread(dev, IBLOCK(inum, sb));
5312         dip = (struct dinode*)bp->data + inum%IPB;
5313         if(dip->type == 0){ // a free inode
5314             memset(dip, 0, sizeof(*dip));
5315             dip->type = type;
5316             log_write(bp); // mark it allocated on the disk
5317             brelse(bp);
5318             return iget(dev, inum);
5319         }
5320         brelse(bp);
5321     }
5322     panic("ialloc: no inodes");
5323 }
5324
5325 // Copy a modified in-memory inode to disk.
5326 // Must be called after every change to an ip->xxx field
5327 // that lives on disk, since i-node cache is write-through.
5328 // Caller must hold ip->lock.
5329 void
5330 iupdate(struct inode *ip)
5331 {
5332     struct buf *bp;
5333     struct dinode *dip;
5334
5335     bp = bread(ip->dev, IBLOCK(ip->inum, sb));
5336     dip = (struct dinode*)bp->data + ip->inum%IPB;
5337     dip->type = ip->type;
5338     dip->major = ip->major;
5339     dip->minor = ip->minor;
5340     dip->nlink = ip->nlink;
5341     dip->size = ip->size;
5342     memmove(dip->addrs, ip->addrs, sizeof(ip->addrs));
5343     log_write(bp);
5344     brelse(bp);
5345 }
5346
5347
5348
5349

```

```

5350 // Find the inode with number inum on device dev
5351 // and return the in-memory copy. Does not lock
5352 // the inode and does not read it from disk.
5353 static struct inode*
5354 iget(uint dev, uint inum)
5355 {
5356     struct inode *ip, *empty;
5357
5358     acquire(&icache.lock);
5359
5360     // Is the inode already cached?
5361     empty = 0;
5362     for(ip = &icache.inode[0]; ip < &icache.inode[NINODE]; ip++){
5363         if(ip->ref > 0 && ip->dev == dev && ip->inum == inum){
5364             ip->ref++;
5365             release(&icache.lock);
5366             return ip;
5367         }
5368         if(empty == 0 && ip->ref == 0)    // Remember empty slot.
5369             empty = ip;
5370     }
5371
5372     // Recycle an inode cache entry.
5373     if(empty == 0)
5374         panic("iget: no inodes");
5375
5376     ip = empty;
5377     ip->dev = dev;
5378     ip->inum = inum;
5379     ip->ref = 1;
5380     ip->valid = 0;
5381     release(&icache.lock);
5382
5383     return ip;
5384 }
5385
5386 // Increment reference count for ip.
5387 // Returns ip to enable ip = idup(ip1) idiom.
5388 struct inode*
5389 idup(struct inode *ip)
5390 {
5391     acquire(&icache.lock);
5392     ip->ref++;
5393     release(&icache.lock);
5394     return ip;
5395 }
5396
5397
5398
5399

```

```

5400 // Lock the given inode.
5401 // Reads the inode from disk if necessary.
5402 void
5403 ilock(struct inode *ip)
5404 {
5405     struct buf *bp;
5406     struct dinode *dip;
5407
5408     if(ip == 0 || ip->ref < 1)
5409         panic("ilock");
5410
5411     acquiresleep(&ip->lock);
5412
5413     if(ip->valid == 0){
5414         bp = bread(ip->dev, IBLOCK(ip->inum, sb));
5415         dip = (struct dinode*)bp->data + ip->inum%IPB;
5416         ip->type = dip->type;
5417         ip->major = dip->major;
5418         ip->minor = dip->minor;
5419         ip->nlink = dip->nlink;
5420         ip->size = dip->size;
5421         memmove(ip->addrs, dip->addrs, sizeof(ip->addrs));
5422         brelse(bp);
5423         ip->valid = 1;
5424         if(ip->type == 0)
5425             panic("ilock: no type");
5426     }
5427 }
5428
5429 // Unlock the given inode.
5430 void
5431 iunlock(struct inode *ip)
5432 {
5433     if(ip == 0 || !holdingsleep(&ip->lock) || ip->ref < 1)
5434         panic("iunlock");
5435
5436     releasesleep(&ip->lock);
5437 }
5438
5439
5440
5441
5442
5443
5444
5445
5446
5447
5448
5449

```

```

5450 // Drop a reference to an in-memory inode.
5451 // If that was the last reference, the inode cache entry can
5452 // be recycled.
5453 // If that was the last reference and the inode has no links
5454 // to it, free the inode (and its content) on disk.
5455 // All calls to iput() must be inside a transaction in
5456 // case it has to free the inode.
5457 void
5458 iput(struct inode *ip)
5459 {
5460     acquiresleep(&ip->lock);
5461     if(ip->valid && ip->nlink == 0){
5462         acquire(&icache.lock);
5463         int r = ip->ref;
5464         release(&icache.lock);
5465         if(r == 1){
5466             // inode has no links and no other references: truncate and free.
5467             itrunc(ip);
5468             ip->type = 0;
5469             iupdate(ip);
5470             ip->valid = 0;
5471         }
5472     }
5473     releasesleep(&ip->lock);
5474
5475     acquire(&icache.lock);
5476     ip->ref--;
5477     release(&icache.lock);
5478 }
5479
5480 // Common idiom: unlock, then put.
5481 void
5482 iunlockput(struct inode *ip)
5483 {
5484     iunlock(ip);
5485     iput(ip);
5486 }
5487
5488
5489
5490
5491
5492
5493
5494
5495
5496
5497
5498
5499

```

```

5500 // Inode content
5501 //
5502 // The content (data) associated with each inode is stored
5503 // in blocks on the disk. The first NDIRECT block numbers
5504 // are listed in ip->addrs[]. The next NINDIRECT blocks are
5505 // listed in block ip->addrs[NDIRECT].
5506
5507 // Return the disk block address of the nth block in inode ip.
5508 // If there is no such block, bmap allocates one.
5509 static uint
5510 bmap(struct inode *ip, uint bn)
5511 {
5512     uint addr, *a;
5513     struct buf *bp;
5514
5515     if(bn < NDIRECT){
5516         if((addr = ip->addrs[bn]) == 0)
5517             ip->addrs[bn] = addr = balloc(ip->dev);
5518         return addr;
5519     }
5520     bn -= NDIRECT;
5521
5522     if(bn < NINDIRECT){
5523         // Load indirect block, allocating if necessary.
5524         if((addr = ip->addrs[NDIRECT]) == 0)
5525             ip->addrs[NDIRECT] = addr = balloc(ip->dev);
5526         bp = bread(ip->dev, addr);
5527         a = (uint*)bp->data;
5528         if((addr = a[bn]) == 0){
5529             a[bn] = addr = balloc(ip->dev);
5530             log_write(bp);
5531         }
5532         brelse(bp);
5533         return addr;
5534     }
5535
5536     panic("bmap: out of range");
5537 }
5538
5539
5540
5541
5542
5543
5544
5545
5546
5547
5548
5549

```

```

5550 // Truncate inode (discard contents).
5551 // Only called when the inode has no links
5552 // to it (no directory entries referring to it)
5553 // and has no in-memory reference to it (is
5554 // not an open file or current directory).
5555 static void
5556 itrunc(struct inode *ip)
5557 {
5558     int i, j;
5559     struct buf *bp;
5560     uint *a;
5561
5562     for(i = 0; i < NDIRECT; i++){
5563         if(ip->addrs[i]){
5564             bfree(ip->dev, ip->addrs[i]);
5565             ip->addrs[i] = 0;
5566         }
5567     }
5568
5569     if(ip->addrs[NDIRECT]){
5570         bp = bread(ip->dev, ip->addrs[NDIRECT]);
5571         a = (uint*)bp->data;
5572         for(j = 0; j < NINDIRECT; j++){
5573             if(a[j])
5574                 bfree(ip->dev, a[j]);
5575         }
5576         brelse(bp);
5577         bfree(ip->dev, ip->addrs[NDIRECT]);
5578         ip->addrs[NDIRECT] = 0;
5579     }
5580
5581     ip->size = 0;
5582     iupdate(ip);
5583 }
5584
5585 // Copy stat information from inode.
5586 // Caller must hold ip->lock.
5587 void
5588 stati(struct inode *ip, struct stat *st)
5589 {
5590     st->dev = ip->dev;
5591     st->ino = ip->inum;
5592     st->type = ip->type;
5593     st->nlink = ip->nlink;
5594     st->size = ip->size;
5595 }
5596
5597
5598
5599

```

```

5600 // Read data from inode.
5601 // Caller must hold ip->lock.
5602 int
5603 readi(struct inode *ip, char *dst, uint off, uint n)
5604 {
5605     uint tot, m;
5606     struct buf *bp;
5607
5608     if(ip->type == T_DEV){
5609         if(ip->major < 0 || ip->major >= NDEV || !devsw[ip->major].read)
5610             return -1;
5611         return devsw[ip->major].read(ip, dst, n);
5612     }
5613
5614     if(off > ip->size || off + n < off)
5615         return -1;
5616     if(off + n > ip->size)
5617         n = ip->size - off;
5618
5619     for(tot=0; tot<n; tot+=m, off+=m, dst+=m){
5620         bp = bread(ip->dev, bmap(ip, off/BSIZE));
5621         m = min(n - tot, BSIZE - off%BSIZE);
5622         memmove(dst, bp->data + off%BSIZE, m);
5623         brelse(bp);
5624     }
5625     return n;
5626 }
5627
5628
5629
5630
5631
5632
5633
5634
5635
5636
5637
5638
5639
5640
5641
5642
5643
5644
5645
5646
5647
5648
5649

```



```

5650 // Write data to inode.
5651 // Caller must hold ip->lock.
5652 int
5653 writei(struct inode *ip, char *src, uint off, uint n)
5654 {
5655     uint tot, m;
5656     struct buf *bp;
5657
5658     if(ip->type == T_DEV){
5659         if(ip->major < 0 || ip->major >= NDEV || !devsw[ip->major].write)
5660             return -1;
5661         return devsw[ip->major].write(ip, src, n);
5662     }
5663
5664     if(off > ip->size || off + n < off)
5665         return -1;
5666     if(off + n > MAXFILE*BSIZE)
5667         return -1;
5668
5669     for(tot=0; tot<n; tot+=m, off+=m, src+=m){
5670         bp = bread(ip->dev, bmap(ip, off/BSIZE));
5671         m = min(n - tot, BSIZE - off%BSIZE);
5672         memmove(bp->data + off%BSIZE, src, m);
5673         log_write(bp);
5674         brelse(bp);
5675     }
5676
5677     if(n > 0 && off > ip->size){
5678         ip->size = off;
5679         iupdate(ip);
5680     }
5681     return n;
5682 }
5683
5684
5685
5686
5687
5688
5689
5690
5691
5692
5693
5694
5695
5696
5697
5698
5699

```

```

5700 // Directories
5701
5702 int
5703 namecmp(const char *s, const char *t)
5704 {
5705     return strncmp(s, t, DIRSIZ);
5706 }
5707
5708 // Look for a directory entry in a directory.
5709 // If found, set *poff to byte offset of entry.
5710 struct inode*
5711 dirlookup(struct inode *dp, char *name, uint *poff)
5712 {
5713     uint off, inum;
5714     struct dirent de;
5715
5716     if(dp->type != T_DIR)
5717         panic("dirlookup not DIR");
5718
5719     for(off = 0; off < dp->size; off += sizeof(de)){
5720         if(readi(dp, (char*)&de, off, sizeof(de)) != sizeof(de))
5721             panic("dirlookup read");
5722         if(de.inum == 0)
5723             continue;
5724         if(namecmp(name, de.name) == 0){
5725             // entry matches path element
5726             if(poff)
5727                 *poff = off;
5728             inum = de.inum;
5729             return iget(dp->dev, inum);
5730         }
5731     }
5732
5733     return 0;
5734 }
5735
5736
5737
5738
5739
5740
5741
5742
5743
5744
5745
5746
5747
5748
5749

```

```

5750 // Write a new directory entry (name, inum) into the directory dp.
5751 int
5752 dirlink(struct inode *dp, char *name, uint inum)
5753 {
5754     int off;
5755     struct dirent de;
5756     struct inode *ip;
5757
5758     // Check that name is not present.
5759     if((ip = dirlookup(dp, name, 0)) != 0){
5760         iput(ip);
5761         return -1;
5762     }
5763
5764     // Look for an empty dirent.
5765     for(off = 0; off < dp->size; off += sizeof(de)){
5766         if(readi(dp, (char*)&de, off, sizeof(de)) != sizeof(de))
5767             panic("dirlink read");
5768         if(de.inum == 0)
5769             break;
5770     }
5771
5772     strncpy(de.name, name, DIRSIZ);
5773     de.inum = inum;
5774     if(writei(dp, (char*)&de, off, sizeof(de)) != sizeof(de))
5775         panic("dirlink");
5776
5777     return 0;
5778 }
5779
5780
5781
5782
5783
5784
5785
5786
5787
5788
5789
5790
5791
5792
5793
5794
5795
5796
5797
5798
5799

```

```

5800 // Paths
5801
5802 // Copy the next path element from path into name.
5803 // Return a pointer to the element following the copied one.
5804 // The returned path has no leading slashes,
5805 // so the caller can check *path=='\0' to see if the name is the last one.
5806 // If no name to remove, return 0.
5807 //
5808 // Examples:
5809 //   skipelem("a/bb/c", name) = "bb/c", setting name = "a"
5810 //   skipelem("///a//bb", name) = "bb", setting name = "a"
5811 //   skipelem("a", name) = "", setting name = "a"
5812 //   skipelem("", name) = skipelem("///", name) = 0
5813 //
5814 static char*
5815 skipelem(char *path, char *name)
5816 {
5817     char *s;
5818     int len;
5819
5820     while(*path == '/')
5821         path++;
5822     if(*path == 0)
5823         return 0;
5824     s = path;
5825     while(*path != '/' && *path != 0)
5826         path++;
5827     len = path - s;
5828     if(len >= DIRSIZ)
5829         memmove(name, s, DIRSIZ);
5830     else {
5831         memmove(name, s, len);
5832         name[len] = 0;
5833     }
5834     while(*path == '/')
5835         path++;
5836     return path;
5837 }
5838
5839
5840
5841
5842
5843
5844
5845
5846
5847
5848
5849

```

```

5850 // Look up and return the inode for a path name.
5851 // If parent != 0, return the inode for the parent and copy the final
5852 // path element into name, which must have room for DIRSIZ bytes.
5853 // Must be called inside a transaction since it calls iput().
5854 static struct inode*
5855 nameix(char *path, int nameparent, char *name)
5856 {
5857     struct inode *ip, *next;
5858
5859     if(*path == '/')
5860         ip = iget(ROOTDEV, ROOTINO);
5861     else
5862         ip = idup(myproc()->cwd);
5863
5864     while((path = skipelem(path, name)) != 0){
5865         ilock(ip);
5866         if(ip->type != T_DIR){
5867             iunlockput(ip);
5868             return 0;
5869         }
5870         if(nameparent && *path == '\0'){
5871             // Stop one level early.
5872             iunlock(ip);
5873             return ip;
5874         }
5875         if((next = dirlookup(ip, name, 0)) == 0){
5876             iunlockput(ip);
5877             return 0;
5878         }
5879         iunlockput(ip);
5880         ip = next;
5881     }
5882     if(nameparent){
5883         iput(ip);
5884         return 0;
5885     }
5886     return ip;
5887 }
5888
5889 struct inode*
5890 namei(char *path)
5891 {
5892     char name[DIRSIZ];
5893     return nameix(path, 0, name);
5894 }
5895
5896
5897
5898
5899

```

```

5900 struct inode*
5901 nameiparent(char *path, char *name)
5902 {
5903     return nameix(path, 1, name);
5904 }
5905
5906
5907
5908
5909
5910
5911
5912
5913
5914
5915
5916
5917
5918
5919
5920
5921
5922
5923
5924
5925
5926
5927
5928
5929
5930
5931
5932
5933
5934
5935
5936
5937
5938
5939
5940
5941
5942
5943
5944
5945
5946
5947
5948
5949

```

```

5950 //
5951 // File descriptors
5952 //
5953
5954 #include "types.h"
5955 #include "defs.h"
5956 #include "param.h"
5957 #include "fs.h"
5958 #include "spinlock.h"
5959 #include "sleeplock.h"
5960 #include "file.h"
5961
5962 struct devsw devsw[NDEV];
5963 struct {
5964     struct spinlock lock;
5965     struct file file[NFILE];
5966 } ftable;
5967
5968 void
5969 fileinit(void)
5970 {
5971     initlock(&ftable.lock, "ftable");
5972 }
5973
5974 // Allocate a file structure.
5975 struct file*
5976 filealloc(void)
5977 {
5978     struct file *f;
5979
5980     acquire(&ftable.lock);
5981     for(f = ftable.file; f < ftable.file + NFILE; f++){
5982         if(f->ref == 0){
5983             f->ref = 1;
5984             release(&ftable.lock);
5985             return f;
5986         }
5987     }
5988     release(&ftable.lock);
5989     return 0;
5990 }
5991
5992
5993
5994
5995
5996
5997
5998
5999

```

```

6000 // Increment ref count for file f.
6001 struct file*
6002 filedup(struct file *f)
6003 {
6004     acquire(&ftable.lock);
6005     if(f->ref < 1)
6006         panic("filedup");
6007     f->ref++;
6008     release(&ftable.lock);
6009     return f;
6010 }
6011
6012 // Close file f. (Decrement ref count, close when reaches 0.)
6013 void
6014 fileclose(struct file *f)
6015 {
6016     struct file ff;
6017
6018     acquire(&ftable.lock);
6019     if(f->ref < 1)
6020         panic("fileclose");
6021     if(--f->ref > 0){
6022         release(&ftable.lock);
6023         return;
6024     }
6025     ff = *f;
6026     f->ref = 0;
6027     f->type = FD_NONE;
6028     release(&ftable.lock);
6029
6030     if(ff.type == FD_PIPE)
6031         pipeclose(ff.pipe, ff.writable);
6032     else if(ff.type == FD_INODE){
6033         begin_op();
6034         iput(ff.ip);
6035         end_op();
6036     }
6037 }
6038
6039
6040
6041
6042
6043
6044
6045
6046
6047
6048
6049

```

```

6050 // Get metadata about file f.
6051 int
6052 filestat(struct file *f, struct stat *st)
6053 {
6054     if(f->type == FD_INODE){
6055         ilock(f->ip);
6056         stati(f->ip, st);
6057         iunlock(f->ip);
6058         return 0;
6059     }
6060     return -1;
6061 }
6062
6063 // Read from file f.
6064 int
6065 fileread(struct file *f, char *addr, int n)
6066 {
6067     int r;
6068
6069     if(f->readable == 0)
6070         return -1;
6071     if(f->type == FD_PIPE)
6072         return piperead(f->pipe, addr, n);
6073     if(f->type == FD_INODE){
6074         ilock(f->ip);
6075         if((r = readi(f->ip, addr, f->off, n)) > 0)
6076             f->off += r;
6077         iunlock(f->ip);
6078         return r;
6079     }
6080     panic("fileread");
6081 }
6082
6083
6084
6085
6086
6087
6088
6089
6090
6091
6092
6093
6094
6095
6096
6097
6098
6099

```

```

6100 // Write to file f.
6101 int
6102 filewrite(struct file *f, char *addr, int n)
6103 {
6104     int r;
6105
6106     if(f->writable == 0)
6107         return -1;
6108     if(f->type == FD_PIPE)
6109         return pipewrite(f->pipe, addr, n);
6110     if(f->type == FD_INODE){
6111         // write a few blocks at a time to avoid exceeding
6112         // the maximum log transaction size, including
6113         // i-node, indirect block, allocation blocks,
6114         // and 2 blocks of slop for non-aligned writes.
6115         // this really belongs lower down, since writei()
6116         // might be writing a device like the console.
6117         int max = ((LOGSIZE-1-1-2) / 2) * 512;
6118         int i = 0;
6119         while(i < n){
6120             int nl = n - i;
6121             if(nl > max)
6122                 nl = max;
6123
6124             begin_op();
6125             ilock(f->ip);
6126             if ((r = writei(f->ip, addr + i, f->off, nl)) > 0)
6127                 f->off += r;
6128             iunlock(f->ip);
6129             end_op();
6130
6131             if(r < 0)
6132                 break;
6133             if(r != nl)
6134                 panic("short filewrite");
6135             i += r;
6136         }
6137         return i == n ? n : -1;
6138     }
6139     panic("filewrite");
6140 }
6141
6142
6143
6144
6145
6146
6147
6148
6149

```

```

6150 //
6151 // File-system system calls.
6152 // Mostly argument checking, since we don't trust
6153 // user code, and calls into file.c and fs.c.
6154 //
6155
6156 #include "types.h"
6157 #include "defs.h"
6158 #include "param.h"
6159 #include "stat.h"
6160 #include "mmu.h"
6161 #include "proc.h"
6162 #include "fs.h"
6163 #include "spinlock.h"
6164 #include "sleeplock.h"
6165 #include "file.h"
6166 #include "fcntl.h"
6167
6168 // Fetch the nth word-sized system call argument as a file descriptor
6169 // and return both the descriptor and the corresponding struct file.
6170 static int
6171 argfd(int n, int *pfd, struct file **pf)
6172 {
6173     int fd;
6174     struct file *f;
6175
6176     if(argint(n, &fd) < 0)
6177         return -1;
6178     if(fd < 0 || fd >= NOFILE || (f=myproc()->ofile[fd]) == 0)
6179         return -1;
6180     if(pfd)
6181         *pfd = fd;
6182     if(pf)
6183         *pf = f;
6184     return 0;
6185 }
6186
6187
6188
6189
6190
6191
6192
6193
6194
6195
6196
6197
6198
6199

```

```

6200 // Allocate a file descriptor for the given file.
6201 // Takes over file reference from caller on success.
6202 static int
6203 fdalloc(struct file *f)
6204 {
6205     int fd;
6206     struct proc *curproc = myproc();
6207
6208     for(fd = 0; fd < NOFILE; fd++){
6209         if(curproc->ofile[fd] == 0){
6210             curproc->ofile[fd] = f;
6211             return fd;
6212         }
6213     }
6214     return -1;
6215 }
6216
6217 int
6218 sys_dup(void)
6219 {
6220     struct file *f;
6221     int fd;
6222
6223     if(argfd(0, 0, &f) < 0)
6224         return -1;
6225     if((fd=fdalloc(f)) < 0)
6226         return -1;
6227     filedup(f);
6228     return fd;
6229 }
6230
6231 int
6232 sys_read(void)
6233 {
6234     struct file *f;
6235     int n;
6236     char *p;
6237
6238     if(argfd(0, 0, &f) < 0 || argint(2, &n) < 0 || argptr(1, &p, n) < 0)
6239         return -1;
6240     return fileread(f, p, n);
6241 }
6242
6243
6244
6245
6246
6247
6248
6249

```

```

6250 int
6251 sys_write(void)
6252 {
6253     struct file *f;
6254     int n;
6255     char *p;
6256
6257     if(argfd(0, 0, &f) < 0 || argint(2, &n) < 0 || argptr(1, &p, n) < 0)
6258         return -1;
6259     return filewrite(f, p, n);
6260 }
6261
6262 int
6263 sys_close(void)
6264 {
6265     int fd;
6266     struct file *f;
6267
6268     if(argfd(0, &fd, &f) < 0)
6269         return -1;
6270     myproc()->ofile[fd] = 0;
6271     fileclose(f);
6272     return 0;
6273 }
6274
6275 int
6276 sys_fstat(void)
6277 {
6278     struct file *f;
6279     struct stat *st;
6280
6281     if(argfd(0, 0, &f) < 0 || argptr(1, (void*)&st, sizeof(*st)) < 0)
6282         return -1;
6283     return filestat(f, st);
6284 }
6285
6286
6287
6288
6289
6290
6291
6292
6293
6294
6295
6296
6297
6298
6299

```

```

6300 // Create the path new as a link to the same inode as old.
6301 int
6302 sys_link(void)
6303 {
6304     char name[DIRSIZ], *new, *old;
6305     struct inode *dp, *ip;
6306
6307     if(argstr(0, &old) < 0 || argstr(1, &new) < 0)
6308         return -1;
6309
6310     begin_op();
6311     if((ip = namei(old)) == 0){
6312         end_op();
6313         return -1;
6314     }
6315
6316     ilock(ip);
6317     if(ip->type == T_DIR){
6318         iunlockput(ip);
6319         end_op();
6320         return -1;
6321     }
6322
6323     ip->nlink++;
6324     iupdate(ip);
6325     iunlock(ip);
6326
6327     if((dp = nameiparent(new, name)) == 0)
6328         goto bad;
6329     ilock(dp);
6330     if(dp->dev != ip->dev || dirlink(dp, name, ip->inum) < 0){
6331         iunlockput(dp);
6332         goto bad;
6333     }
6334     iunlockput(dp);
6335     iput(ip);
6336
6337     end_op();
6338
6339     return 0;
6340
6341 bad:
6342     ilock(ip);
6343     ip->nlink--;
6344     iupdate(ip);
6345     iunlockput(ip);
6346     end_op();
6347     return -1;
6348 }
6349

```

```

6350 // Is the directory dp empty except for "." and ".." ?
6351 static int
6352 isdirempty(struct inode *dp)
6353 {
6354     int off;
6355     struct dirent de;
6356
6357     for(off=2*sizeof(de); off<dp->size; off+=sizeof(de)){
6358         if(readi(dp, (char*)&de, off, sizeof(de)) != sizeof(de))
6359             panic("isdirempty: readi");
6360         if(de.inum != 0)
6361             return 0;
6362     }
6363     return 1;
6364 }
6365
6366
6367
6368
6369
6370
6371
6372
6373
6374
6375
6376
6377
6378
6379
6380
6381
6382
6383
6384
6385
6386
6387
6388
6389
6390
6391
6392
6393
6394
6395
6396
6397
6398
6399

```

```

6400 int
6401 sys_unlink(void)
6402 {
6403     struct inode *ip, *dp;
6404     struct dirent de;
6405     char name[DIRSIZ], *path;
6406     uint off;
6407
6408     if(argstr(0, &path) < 0)
6409         return -1;
6410
6411     begin_op();
6412     if((dp = nameiparent(path, name)) == 0){
6413         end_op();
6414         return -1;
6415     }
6416
6417     ilock(dp);
6418
6419     // Cannot unlink "." or "..".
6420     if(namecmp(name, ".") == 0 || namecmp(name, "..") == 0)
6421         goto bad;
6422
6423     if((ip = dirlookup(dp, name, &off)) == 0)
6424         goto bad;
6425     ilock(ip);
6426
6427     if(ip->nlink < 1)
6428         panic("unlink: nlink < 1");
6429     if(ip->type == T_DIR && !isdirempty(ip)){
6430         iunlockput(ip);
6431         goto bad;
6432     }
6433
6434     memset(&de, 0, sizeof(de));
6435     if(writei(dp, (char*)&de, off, sizeof(de)) != sizeof(de))
6436         panic("unlink: writei");
6437     if(ip->type == T_DIR){
6438         dp->nlink--;
6439         iupdate(dp);
6440     }
6441     iunlockput(dp);
6442
6443     ip->nlink--;
6444     iupdate(ip);
6445     iunlockput(ip);
6446
6447     end_op();
6448
6449     return 0;

```



```

6450 bad:
6451  iunlockput(dp);
6452  end_op();
6453  return -1;
6454 }
6455
6456 static struct inode*
6457 create(char *path, short type, short major, short minor)
6458 {
6459  uint off;
6460  struct inode *ip, *dp;
6461  char name[DIRSIZ];
6462
6463  if((dp = nameiparent(path, name)) == 0)
6464    return 0;
6465  ilock(dp);
6466
6467  if((ip = dirlookup(dp, name, &off)) != 0){
6468    iunlockput(dp);
6469    ilock(ip);
6470    if(type == T_FILE && ip->type == T_FILE)
6471      return ip;
6472    iunlockput(ip);
6473    return 0;
6474  }
6475
6476  if((ip = ialloc(dp->dev, type)) == 0)
6477    panic("create: ialloc");
6478
6479  ilock(ip);
6480  ip->major = major;
6481  ip->minor = minor;
6482  ip->nlink = 1;
6483  iupdate(ip);
6484
6485  if(type == T_DIR){ // Create . and .. entries.
6486    dp->nlink++; // for ".."
6487    iupdate(dp);
6488    // No ip->nlink++ for ".": avoid cyclic ref count.
6489    if(dirlink(ip, ".", ip->inum) < 0 || dirlink(ip, "..", dp->inum) < 0)
6490      panic("create dots");
6491  }
6492
6493  if(dirlink(dp, name, ip->inum) < 0)
6494    panic("create: dirlink");
6495
6496  iunlockput(dp);
6497
6498  return ip;
6499 }

```

```

6500 int
6501 sys_open(void)
6502 {
6503  char *path;
6504  int fd, omode;
6505  struct file *f;
6506  struct inode *ip;
6507
6508  if(argstr(0, &path) < 0 || argint(1, &omode) < 0)
6509    return -1;
6510
6511  begin_op();
6512
6513  if(omode & O_CREATE){
6514    ip = create(path, T_FILE, 0, 0);
6515    if(ip == 0){
6516      end_op();
6517      return -1;
6518    }
6519  } else {
6520    if((ip = namei(path)) == 0){
6521      end_op();
6522      return -1;
6523    }
6524    ilock(ip);
6525    if(ip->type == T_DIR && omode != O_RDONLY){
6526      iunlockput(ip);
6527      end_op();
6528      return -1;
6529    }
6530  }
6531
6532  if((f = filealloc()) == 0 || (fd = fdalloc(f)) < 0){
6533    if(f)
6534      fileclose(f);
6535    iunlockput(ip);
6536    end_op();
6537    return -1;
6538  }
6539  iunlock(ip);
6540  end_op();
6541
6542  f->type = FD_INODE;
6543  f->ip = ip;
6544  f->off = 0;
6545  f->readable = !(omode & O_WRONLY);
6546  f->writable = (omode & O_WRONLY) || (omode & O_RDWR);
6547  return fd;
6548 }
6549

```

```

6550 int
6551 sys_mkdir(void)
6552 {
6553     char *path;
6554     struct inode *ip;
6555
6556     begin_op();
6557     if(argstr(0, &path) < 0 || (ip = create(path, T_DIR, 0, 0)) == 0){
6558         end_op();
6559         return -1;
6560     }
6561     iunlockput(ip);
6562     end_op();
6563     return 0;
6564 }
6565
6566 int
6567 sys_mknod(void)
6568 {
6569     struct inode *ip;
6570     char *path;
6571     int major, minor;
6572
6573     begin_op();
6574     if((argstr(0, &path)) < 0 ||
6575        argint(1, &major) < 0 ||
6576        argint(2, &minor) < 0 ||
6577        (ip = create(path, T_DEV, major, minor)) == 0){
6578         end_op();
6579         return -1;
6580     }
6581     iunlockput(ip);
6582     end_op();
6583     return 0;
6584 }
6585
6586
6587
6588
6589
6590
6591
6592
6593
6594
6595
6596
6597
6598
6599

```

```

6600 int
6601 sys_chdir(void)
6602 {
6603     char *path;
6604     struct inode *ip;
6605     struct proc *curproc = myproc();
6606
6607     begin_op();
6608     if(argstr(0, &path) < 0 || (ip = namei(path)) == 0){
6609         end_op();
6610         return -1;
6611     }
6612     ilock(ip);
6613     if(ip->type != T_DIR){
6614         iunlockput(ip);
6615         end_op();
6616         return -1;
6617     }
6618     iunlock(ip);
6619     iput(curproc->cwd);
6620     end_op();
6621     curproc->cwd = ip;
6622     return 0;
6623 }
6624
6625 int
6626 sys_exec(void)
6627 {
6628     char *path, *argv[MAXARG];
6629     int i;
6630     uint uargv, uarg;
6631
6632     if(argstr(0, &path) < 0 || argint(1, (int*)&uargv) < 0){
6633         return -1;
6634     }
6635     memset(argv, 0, sizeof(argv));
6636     for(i=0;; i++){
6637         if(i >= NELEM(argv))
6638             return -1;
6639         if(fetchint(uargv+4*i, (int*)&uarg) < 0)
6640             return -1;
6641         if(uarg == 0){
6642             argv[i] = 0;
6643             break;
6644         }
6645         if(fetchstr(uarg, &argv[i]) < 0)
6646             return -1;
6647     }
6648     return exec(path, argv);
6649 }

```

```

6650 int
6651 sys_pipe(void)
6652 {
6653     int *fd;
6654     struct file *rf, *wf;
6655     int fd0, fd1;
6656
6657     if(argptr(0, (void*)&fd, 2*sizeof(fd[0])) < 0)
6658         return -1;
6659     if(pipealloc(&rf, &wf) < 0)
6660         return -1;
6661     fd0 = -1;
6662     if((fd0 = fdalloc(rf)) < 0 || (fd1 = fdalloc(wf)) < 0){
6663         if(fd0 >= 0)
6664             myproc()->ofile[fd0] = 0;
6665         fileclose(rf);
6666         fileclose(wf);
6667         return -1;
6668     }
6669     fd[0] = fd0;
6670     fd[1] = fd1;
6671     return 0;
6672 }
6673
6674
6675
6676
6677
6678
6679
6680
6681
6682
6683
6684
6685
6686
6687
6688
6689
6690
6691
6692
6693
6694
6695
6696
6697
6698
6699

```

```

6700 #include "types.h"
6701 #include "param.h"
6702 #include "memlayout.h"
6703 #include "mmu.h"
6704 #include "proc.h"
6705 #include "defs.h"
6706 #include "x86.h"
6707 #include "elf.h"
6708
6709 int
6710 exec(char *path, char **argv)
6711 {
6712     char *s, *last;
6713     int i, off;
6714     uint argc, sz, sp, ustack[3+MAXARG+1];
6715     struct elfhdr elf;
6716     struct inode *ip;
6717     struct proghdr ph;
6718     pde_t *pgdir, *oldpgdir;
6719     struct proc *curproc = myproc();
6720
6721     begin_op();
6722
6723     if((ip = namei(path)) == 0){
6724         end_op();
6725         cprintf("exec: fail\n");
6726         return -1;
6727     }
6728     ilock(ip);
6729     pgdir = 0;
6730
6731     // Check ELF header
6732     if(readi(ip, (char*)&elf, 0, sizeof(elf)) != sizeof(elf))
6733         goto bad;
6734     if(elf.magic != ELF_MAGIC)
6735         goto bad;
6736
6737     if((pgdir = setupkvm()) == 0)
6738         goto bad;
6739
6740     // Load program into memory.
6741     sz = 0;
6742     for(i=0, off=elf.phoff; i<elf.phnum; i++, off+=sizeof(ph)){
6743         if(readi(ip, (char*)&ph, off, sizeof(ph)) != sizeof(ph))
6744             goto bad;
6745         if(ph.type != ELF_PROG_LOAD)
6746             continue;
6747         if(ph.memsz < ph.filesz)
6748             goto bad;
6749         if(ph.vaddr + ph.memsz < ph.vaddr)

```

```

6750     goto bad;
6751     if((sz = allocuvm(pgdir, sz, ph.vaddr + ph.memsz)) == 0)
6752         goto bad;
6753     if(ph.vaddr % PGSIZE != 0)
6754         goto bad;
6755     if(loaduvm(pgdir, (char*)ph.vaddr, ip, ph.off, ph.filesz) < 0)
6756         goto bad;
6757 }
6758 iunlockput(ip);
6759 end_op();
6760 ip = 0;
6761
6762 // Allocate two pages at the next page boundary.
6763 // Make the first inaccessible. Use the second as the user stack.
6764 sz = PGROUNDUP(sz);
6765 if((sz = allocuvm(pgdir, sz, sz + 2*PGSIZE)) == 0)
6766     goto bad;
6767 clearpteu(pgdir, (char*)(sz - 2*PGSIZE));
6768 sp = sz;
6769
6770 // Push argument strings, prepare rest of stack in ustack.
6771 for(argc = 0; argv[argc]; argc++) {
6772     if(argc >= MAXARG)
6773         goto bad;
6774     sp = (sp - (strlen(argv[argc]) + 1)) & ~3;
6775     if(copyout(pgdir, sp, argv[argc], strlen(argv[argc]) + 1) < 0)
6776         goto bad;
6777     ustack[3+argc] = sp;
6778 }
6779 ustack[3+argc] = 0;
6780
6781 ustack[0] = 0xffffffff; // fake return PC
6782 ustack[1] = argc;
6783 ustack[2] = sp - (argc+1)*4; // argv pointer
6784
6785 sp -= (3+argc+1) * 4;
6786 if(copyout(pgdir, sp, ustack, (3+argc+1)*4) < 0)
6787     goto bad;
6788
6789 // Save program name for debugging.
6790 for(last=s=path; *s; s++)
6791     if(*s == '/')
6792         last = s+1;
6793 safestrcpy(curproc->name, last, sizeof(curproc->name));
6794
6795 // Commit to the user image.
6796 oldpgdir = curproc->pgdir;
6797 curproc->pgdir = pgdir;
6798 curproc->sz = sz;
6799 curproc->tf->eip = elf.entry; // main

```

```

6800     curproc->tf->esp = sp;
6801     switchuvm(curproc);
6802     freevm(oldpgdir);
6803     return 0;
6804
6805 bad:
6806     if(pgdir)
6807         freevm(pgdir);
6808     if(ip){
6809         iunlockput(ip);
6810         end_op();
6811     }
6812     return -1;
6813 }
6814
6815
6816
6817
6818
6819
6820
6821
6822
6823
6824
6825
6826
6827
6828
6829
6830
6831
6832
6833
6834
6835
6836
6837
6838
6839
6840
6841
6842
6843
6844
6845
6846
6847
6848
6849

```

```

6850 #include "types.h"
6851 #include "defs.h"
6852 #include "param.h"
6853 #include "mmu.h"
6854 #include "proc.h"
6855 #include "fs.h"
6856 #include "spinlock.h"
6857 #include "sleeplock.h"
6858 #include "file.h"
6859
6860 #define PIPESIZE 512
6861
6862 struct pipe {
6863     struct spinlock lock;
6864     char data[PIPESIZE];
6865     uint nread; // number of bytes read
6866     uint nwrite; // number of bytes written
6867     int readopen; // read fd is still open
6868     int writeopen; // write fd is still open
6869 };
6870
6871 int
6872 pipealloc(struct file **f0, struct file **f1)
6873 {
6874     struct pipe *p;
6875
6876     p = 0;
6877     *f0 = *f1 = 0;
6878     if((*f0 = filealloc()) == 0 || (*f1 = filealloc()) == 0)
6879         goto bad;
6880     if((p = (struct pipe*)kalloc()) == 0)
6881         goto bad;
6882     p->readopen = 1;
6883     p->writeopen = 1;
6884     p->nwrite = 0;
6885     p->nread = 0;
6886     initlock(&p->lock, "pipe");
6887     (*f0)->type = FD_PIPE;
6888     (*f0)->readable = 1;
6889     (*f0)->writable = 0;
6890     (*f0)->pipe = p;
6891     (*f1)->type = FD_PIPE;
6892     (*f1)->readable = 0;
6893     (*f1)->writable = 1;
6894     (*f1)->pipe = p;
6895     return 0;
6896
6897
6898
6899

```

```

6900 bad:
6901     if(p)
6902         kfree((char*)p);
6903     if(*f0)
6904         fileclose(*f0);
6905     if(*f1)
6906         fileclose(*f1);
6907     return -1;
6908 }
6909
6910 void
6911 pipeclose(struct pipe *p, int writable)
6912 {
6913     acquire(&p->lock);
6914     if(writable){
6915         p->writeopen = 0;
6916         wakeup(&p->nread);
6917     } else {
6918         p->readopen = 0;
6919         wakeup(&p->nwrite);
6920     }
6921     if(p->readopen == 0 && p->writeopen == 0){
6922         release(&p->lock);
6923         kfree((char*)p);
6924     } else
6925         release(&p->lock);
6926 }
6927
6928
6929 int
6930 pipewrite(struct pipe *p, char *addr, int n)
6931 {
6932     int i;
6933
6934     acquire(&p->lock);
6935     for(i = 0; i < n; i++){
6936         while(p->nwrite == p->nread + PIPESIZE){
6937             if(p->readopen == 0 || myproc()->killed){
6938                 release(&p->lock);
6939                 return -1;
6940             }
6941             wakeup(&p->nread);
6942             sleep(&p->nwrite, &p->lock);
6943         }
6944         p->data[p->nwrite++ % PIPESIZE] = addr[i];
6945     }
6946     wakeup(&p->nread);
6947     release(&p->lock);
6948     return n;
6949 }

```

```

6950 int
6951 piperead(struct pipe *p, char *addr, int n)
6952 {
6953     int i;
6954
6955     acquire(&p->lock);
6956     while(p->nread == p->nwrite && p->writeopen){
6957         if(myproc()->killed){
6958             release(&p->lock);
6959             return -1;
6960         }
6961         sleep(&p->nread, &p->lock);
6962     }
6963     for(i = 0; i < n; i++){
6964         if(p->nread == p->nwrite)
6965             break;
6966         addr[i] = p->data[p->nread++ % PIPESIZE];
6967     }
6968     wakeup(&p->nwrite);
6969     release(&p->lock);
6970     return i;
6971 }
6972
6973
6974
6975
6976
6977
6978
6979
6980
6981
6982
6983
6984
6985
6986
6987
6988
6989
6990
6991
6992
6993
6994
6995
6996
6997
6998
6999

```

```

7000 #include "types.h"
7001 #include "x86.h"
7002
7003 void*
7004 memset(void *dst, int c, uint n)
7005 {
7006     if ((int)dst%4 == 0 && n%4 == 0){
7007         c &= 0xFF;
7008         stosl(dst, (c<<24)|(c<<16)|(c<<8)|c, n/4);
7009     } else
7010         stosb(dst, c, n);
7011     return dst;
7012 }
7013
7014 int
7015 memcmp(const void *v1, const void *v2, uint n)
7016 {
7017     const uchar *s1, *s2;
7018
7019     s1 = v1;
7020     s2 = v2;
7021     while(n-- > 0){
7022         if(*s1 != *s2)
7023             return *s1 - *s2;
7024         s1++, s2++;
7025     }
7026
7027     return 0;
7028 }
7029
7030 void*
7031 memmove(void *dst, const void *src, uint n)
7032 {
7033     const char *s;
7034     char *d;
7035
7036     s = src;
7037     d = dst;
7038     if(s < d && s + n > d){
7039         s += n;
7040         d += n;
7041         while(n-- > 0)
7042             *--d = *--s;
7043     } else
7044         while(n-- > 0)
7045             *d++ = *s++;
7046
7047     return dst;
7048 }
7049

```

```

7050 // memcpy exists to placate GCC. Use memmove.
7051 void*
7052 memcpy(void *dst, const void *src, uint n)
7053 {
7054     return memmove(dst, src, n);
7055 }
7056
7057 int
7058 strncmp(const char *p, const char *q, uint n)
7059 {
7060     while(n > 0 && *p && *p == *q)
7061         n--, p++, q++;
7062     if(n == 0)
7063         return 0;
7064     return (uchar)*p - (uchar)*q;
7065 }
7066
7067 char*
7068 strncpy(char *s, const char *t, int n)
7069 {
7070     char *os;
7071
7072     os = s;
7073     while(n-- > 0 && (*s++ = *t++) != 0)
7074         ;
7075     while(n-- > 0)
7076         *s++ = 0;
7077     return os;
7078 }
7079
7080 // Like strncpy but guaranteed to NUL-terminate.
7081 char*
7082 safestrcpy(char *s, const char *t, int n)
7083 {
7084     char *os;
7085
7086     os = s;
7087     if(n <= 0)
7088         return os;
7089     while(--n > 0 && (*s++ = *t++) != 0)
7090         ;
7091     *s = 0;
7092     return os;
7093 }
7094
7095
7096
7097
7098
7099

```

```

7100 int
7101 strlen(const char *s)
7102 {
7103     int n;
7104
7105     for(n = 0; s[n]; n++)
7106         ;
7107     return n;
7108 }
7109
7110
7111
7112
7113
7114
7115
7116
7117
7118
7119
7120
7121
7122
7123
7124
7125
7126
7127
7128
7129
7130
7131
7132
7133
7134
7135
7136
7137
7138
7139
7140
7141
7142
7143
7144
7145
7146
7147
7148
7149

```

```

7150 // See MultiProcessor Specification Version 1.[14]
7151
7152 struct mp {           // floating pointer
7153     uchar signature[4]; // "_MP_"
7154     void *physaddr;    // phys addr of MP config table
7155     uchar length;     // 1
7156     uchar specrev;    // [14]
7157     uchar checksum;   // all bytes must add up to 0
7158     uchar type;       // MP system config type
7159     uchar imcrp;
7160     uchar reserved[3];
7161 };
7162
7163 struct mpconf {       // configuration table header
7164     uchar signature[4]; // "PCMP"
7165     ushort length;     // total table length
7166     uchar version;     // [14]
7167     uchar checksum;   // all bytes must add up to 0
7168     uchar product[20]; // product id
7169     uint *oemtable;   // OEM table pointer
7170     ushort oemlength; // OEM table length
7171     ushort entry;     // entry count
7172     uint *lapicaddr;  // address of local APIC
7173     ushort xlength;   // extended table length
7174     uchar xchecksum;  // extended table checksum
7175     uchar reserved;
7176 };
7177
7178 struct mpproc {       // processor table entry
7179     uchar type;       // entry type (0)
7180     uchar apicid;    // local APIC id
7181     uchar version;   // local APIC verison
7182     uchar flags;     // CPU flags
7183     #define MPBOOT 0x02 // This proc is the bootstrap processor.
7184     uchar signature[4]; // CPU signature
7185     uint feature;     // feature flags from CPUID instruction
7186     uchar reserved[8];
7187 };
7188
7189 struct mpioapic {    // I/O APIC table entry
7190     uchar type;       // entry type (2)
7191     uchar apicno;    // I/O APIC id
7192     uchar version;   // I/O APIC version
7193     uchar flags;     // I/O APIC flags
7194     uint *addr;      // I/O APIC address
7195 };
7196
7197
7198
7199

```

```

7200 // Table entry types
7201 #define MPPROC 0x00 // One per processor
7202 #define MPBUS 0x01 // One per bus
7203 #define MPIOAPIC 0x02 // One per I/O APIC
7204 #define MPIOINTR 0x03 // One per bus interrupt source
7205 #define MPLINTR 0x04 // One per system interrupt source
7206
7207
7208
7209
7210
7211
7212
7213
7214
7215
7216
7217
7218
7219
7220
7221
7222
7223
7224
7225
7226
7227
7228
7229
7230
7231
7232
7233
7234
7235
7236
7237
7238
7239
7240
7241
7242
7243
7244
7245
7246
7247
7248
7249

```



```
7250 // Blank page.
7251
7252
7253
7254
7255
7256
7257
7258
7259
7260
7261
7262
7263
7264
7265
7266
7267
7268
7269
7270
7271
7272
7273
7274
7275
7276
7277
7278
7279
7280
7281
7282
7283
7284
7285
7286
7287
7288
7289
7290
7291
7292
7293
7294
7295
7296
7297
7298
7299
```

```
7300 // Multiprocessor support
7301 // Search memory for MP description structures.
7302 // http://developer.intel.com/design/pentium/datashts/24201606.pdf
7303
7304 #include "types.h"
7305 #include "defs.h"
7306 #include "param.h"
7307 #include "memlayout.h"
7308 #include "mp.h"
7309 #include "x86.h"
7310 #include "mmu.h"
7311 #include "proc.h"
7312
7313 struct cpu cpus[NCPU];
7314 int ncpu;
7315 uchar ioapicid;
7316
7317 static uchar
7318 sum(uchar *addr, int len)
7319 {
7320     int i, sum;
7321
7322     sum = 0;
7323     for(i=0; i<len; i++)
7324         sum += addr[i];
7325     return sum;
7326 }
7327
7328 // Look for an MP structure in the len bytes at addr.
7329 static struct mp*
7330 mpsearch1(uint a, int len)
7331 {
7332     uchar *e, *p, *addr;
7333
7334     addr = P2V(a);
7335     e = addr+len;
7336     for(p = addr; p < e; p += sizeof(struct mp))
7337         if(memcmp(p, "_MP_", 4) == 0 && sum(p, sizeof(struct mp)) == 0)
7338             return (struct mp*)p;
7339     return 0;
7340 }
7341
7342
7343
7344
7345
7346
7347
7348
7349
```

```

7350 // Search for the MP Floating Pointer Structure, which according to the
7351 // spec is in one of the following three locations:
7352 // 1) in the first KB of the EBDA;
7353 // 2) in the last KB of system base memory;
7354 // 3) in the BIOS ROM between 0xE0000 and 0xFFFFF.
7355 static struct mp*
7356 mpsearch(void)
7357 {
7358     uchar *bda;
7359     uint p;
7360     struct mp *mp;
7361
7362     bda = (uchar *) P2V(0x400);
7363     if((p = ((bda[0x0F]<<8) | bda[0x0E]) << 4)){
7364         if((mp = mpsearch1(p, 1024)))
7365             return mp;
7366     } else {
7367         p = ((bda[0x14]<<8) | bda[0x13])*1024;
7368         if((mp = mpsearch1(p-1024, 1024)))
7369             return mp;
7370     }
7371     return mpsearch1(0xF0000, 0x10000);
7372 }
7373
7374 // Search for an MP configuration table. For now,
7375 // don't accept the default configurations (physaddr == 0).
7376 // Check for correct signature, calculate the checksum and,
7377 // if correct, check the version.
7378 // To do: check extended table checksum.
7379 static struct mpconf*
7380 mpconfig(struct mp **pmp)
7381 {
7382     struct mpconf *conf;
7383     struct mp *mp;
7384
7385     if((mp = mpsearch()) == 0 || mp->physaddr == 0)
7386         return 0;
7387     conf = (struct mpconf*) P2V((uint) mp->physaddr);
7388     if(memcmp(conf, "PCMP", 4) != 0)
7389         return 0;
7390     if(conf->version != 1 && conf->version != 4)
7391         return 0;
7392     if(sum((uchar*)conf, conf->length) != 0)
7393         return 0;
7394     *pmp = mp;
7395     return conf;
7396 }
7397
7398
7399

```

```

7400 void
7401 mpinit(void)
7402 {
7403     uchar *p, *e;
7404     int ismp;
7405     struct mp *mp;
7406     struct mpconf *conf;
7407     struct mpproc *proc;
7408     struct mpioapic *ioapic;
7409
7410     if((conf = mpconfig(&mp)) == 0)
7411         panic("Expect to run on an SMP");
7412     ismp = 1;
7413     lapic = (uint*)conf->lapicaddr;
7414     for(p=(uchar*)(conf+1), e=(uchar*)conf+conf->length; p<e; ){
7415         switch(*p){
7416             case MPPROC:
7417                 proc = (struct mpproc*)p;
7418                 if(ncpu < NCPU) {
7419                     cpus[ncpu].apicid = proc->apicid; // apicid may differ from ncpu
7420                     ncpu++;
7421                 }
7422                 p += sizeof(struct mpproc);
7423                 continue;
7424             case MPIOAPIC:
7425                 ioapic = (struct mpioapic*)p;
7426                 ioapicid = ioapic->apicno;
7427                 p += sizeof(struct mpioapic);
7428                 continue;
7429             case MPBUS:
7430             case MPIOINTR:
7431             case MPLINTR:
7432                 p += 8;
7433                 continue;
7434             default:
7435                 ismp = 0;
7436                 break;
7437         }
7438     }
7439     if(!ismp)
7440         panic("Didn't find a suitable machine");
7441
7442     if(mp->imcrp){
7443         // Bochs doesn't support IMCR, so this doesn't run on Bochs.
7444         // But it would on real hardware.
7445         outb(0x22, 0x70); // Select IMCR
7446         outb(0x23, inb(0x23) | 1); // Mask external interrupts.
7447     }
7448 }
7449

```

```

7450 // The local APIC manages internal (non-I/O) interrupts.
7451 // See Chapter 8 & Appendix C of Intel processor manual volume 3.
7452
7453 #include "param.h"
7454 #include "types.h"
7455 #include "defs.h"
7456 #include "date.h"
7457 #include "memlayout.h"
7458 #include "traps.h"
7459 #include "mmu.h"
7460 #include "x86.h"
7461
7462 // Local APIC registers, divided by 4 for use as uint[] indices.
7463 #define ID      (0x0020/4) // ID
7464 #define VER     (0x0030/4) // Version
7465 #define TPR    (0x0080/4) // Task Priority
7466 #define EOI    (0x00B0/4) // EOI
7467 #define SVR    (0x00F0/4) // Spurious Interrupt Vector
7468 #define ENABLE 0x00000100 // Unit Enable
7469 #define ESR    (0x0280/4) // Error Status
7470 #define ICRLO (0x0300/4) // Interrupt Command
7471 #define INIT   0x00000500 // INIT/RESET
7472 #define STARTUP 0x00000600 // Startup IPI
7473 #define DELIVS 0x00001000 // Delivery status
7474 #define ASSERT 0x00004000 // Assert interrupt (vs deassert)
7475 #define DEASSERT 0x00000000
7476 #define LEVEL  0x00008000 // Level triggered
7477 #define BCAST  0x00008000 // Send to all APICs, including self.
7478 #define BUSY   0x00001000
7479 #define FIXED  0x00000000
7480 #define ICRHI  (0x0310/4) // Interrupt Command [63:32]
7481 #define TIMER  (0x0320/4) // Local Vector Table 0 (TIMER)
7482 #define X1     0x0000000B // divide counts by 1
7483 #define PERIODIC 0x00020000 // Periodic
7484 #define PCINT  (0x0340/4) // Performance Counter LVT
7485 #define LINT0  (0x0350/4) // Local Vector Table 1 (LINT0)
7486 #define LINT1  (0x0360/4) // Local Vector Table 2 (LINT1)
7487 #define ERROR  (0x0370/4) // Local Vector Table 3 (ERROR)
7488 #define MASKED 0x00010000 // Interrupt masked
7489 #define TICC  (0x0380/4) // Timer Initial Count
7490 #define TCCR  (0x0390/4) // Timer Current Count
7491 #define TDCR  (0x03E0/4) // Timer Divide Configuration
7492
7493 volatile uint *lapic; // Initialized in mp.c
7494
7495 static void
7496 lapicw(int index, int value)
7497 {
7498     lapic[index] = value;
7499     lapic[ID]; // wait for write to finish, by reading

```

```

7500 }
7501
7502
7503
7504
7505
7506
7507
7508
7509
7510
7511
7512
7513
7514
7515
7516
7517
7518
7519
7520
7521
7522
7523
7524
7525
7526
7527
7528
7529
7530
7531
7532
7533
7534
7535
7536
7537
7538
7539
7540
7541
7542
7543
7544
7545
7546
7547
7548
7549

```

```

7550 void
7551 lapicinit(void)
7552 {
7553     if(!lapic)
7554         return;
7555
7556     // Enable local APIC; set spurious interrupt vector.
7557     lapicw(SVR, ENABLE | (T_IRQ0 + IRQ_SPURIOUS));
7558
7559     // The timer repeatedly counts down at bus frequency
7560     // from lapic[TICR] and then issues an interrupt.
7561     // If xv6 cared more about precise timekeeping,
7562     // TICR would be calibrated using an external time source.
7563     lapicw(TDCR, X1);
7564     lapicw(TIMER, PERIODIC | (T_IRQ0 + IRQ_TIMER));
7565     lapicw(TICR, 10000000);
7566
7567     // Disable logical interrupt lines.
7568     lapicw(LINT0, MASKED);
7569     lapicw(LINT1, MASKED);
7570
7571     // Disable performance counter overflow interrupts
7572     // on machines that provide that interrupt entry.
7573     if(((lapic[VER]>>16) & 0xFF) >= 4)
7574         lapicw(PCINT, MASKED);
7575
7576     // Map error interrupt to IRQ_ERROR.
7577     lapicw(ERROR, T_IRQ0 + IRQ_ERROR);
7578
7579     // Clear error status register (requires back-to-back writes).
7580     lapicw(ESR, 0);
7581     lapicw(ESR, 0);
7582
7583     // Ack any outstanding interrupts.
7584     lapicw(EOI, 0);
7585
7586     // Send an Init Level De-Assert to synchronise arbitration ID's.
7587     lapicw(ICRHI, 0);
7588     lapicw(ICRLO, BCAST | INIT | LEVEL);
7589     while(lapic[ICRLO] & DELIVS)
7590         ;
7591
7592     // Enable interrupts on the APIC (but not on the processor).
7593     lapicw(TPR, 0);
7594 }
7595
7596
7597
7598
7599

```

```

7600 int
7601 lapicid(void)
7602 {
7603     if (!lapic)
7604         return 0;
7605     return lapic[ID] >> 24;
7606 }
7607
7608 // Acknowledge interrupt.
7609 void
7610 lapiceoi(void)
7611 {
7612     if(lapic)
7613         lapicw(EOI, 0);
7614 }
7615
7616 // Spin for a given number of microseconds.
7617 // On real hardware would want to tune this dynamically.
7618 void
7619 microdelay(int us)
7620 {
7621 }
7622
7623 #define CMOS_PORT    0x70
7624 #define CMOS_RETURN  0x71
7625
7626 // Start additional processor running entry code at addr.
7627 // See Appendix B of MultiProcessor Specification.
7628 void
7629 lapicstartap(uchar apicid, uint addr)
7630 {
7631     int i;
7632     ushort *wrv;
7633
7634     // "The BSP must initialize CMOS shutdown code to 0AH
7635     // and the warm reset vector (DWORD based at 40:67) to point at
7636     // the AP startup code prior to the [universal startup algorithm]."
7637     outb(CMOS_PORT, 0xF); // offset 0xF is shutdown code
7638     outb(CMOS_PORT+1, 0x0A);
7639     wrv = (ushort*)P2V((0x40<<4 | 0x67)); // Warm reset vector
7640     wrv[0] = 0;
7641     wrv[1] = addr >> 4;
7642
7643     // "Universal startup algorithm."
7644     // Send INIT (level-triggered) interrupt to reset other CPU.
7645     lapicw(ICRHI, apicid<<24);
7646     lapicw(ICRLO, INIT | LEVEL | ASSERT);
7647     microdelay(200);
7648     lapicw(ICRLO, INIT | LEVEL);
7649     microdelay(100); // should be 10ms, but too slow in Bochs!

```

```

7650 // Send startup IPI (twice!) to enter code.
7651 // Regular hardware is supposed to only accept a STARTUP
7652 // when it is in the halted state due to an INIT. So the second
7653 // should be ignored, but it is part of the official Intel algorithm.
7654 // Bochs complains about the second one. Too bad for Bochs.
7655 for(i = 0; i < 2; i++){
7656     lapicw(ICRHI, apicid<<24);
7657     lapicw(ICRLO, STARTUP | (addr>>12));
7658     microdelay(200);
7659 }
7660 }
7661
7662 #define CMOS_STATA  0x0a
7663 #define CMOS_STATB  0x0b
7664 #define CMOS_UIP    (1 << 7)      // RTC update in progress
7665
7666 #define SECS        0x00
7667 #define MINS        0x02
7668 #define HOURS       0x04
7669 #define DAY         0x07
7670 #define MONTH       0x08
7671 #define YEAR        0x09
7672
7673 static uint cmos_read(uint reg)
7674 {
7675     outb(CMOS_PORT, reg);
7676     microdelay(200);
7677
7678     return inb(CMOS_RETURN);
7679 }
7680
7681 static void fill_rtcdate(struct rtcdate *r)
7682 {
7683     r->second = cmos_read(SECS);
7684     r->minute = cmos_read(MINS);
7685     r->hour   = cmos_read(HOURS);
7686     r->day    = cmos_read(DAY);
7687     r->month  = cmos_read(MONTH);
7688     r->year   = cmos_read(YEAR);
7689 }
7690
7691
7692
7693
7694
7695
7696
7697
7698
7699

```

```

7700 // qemu seems to use 24-hour GWT and the values are BCD encoded
7701 void cmostime(struct rtcdate *r)
7702 {
7703     struct rtcdate t1, t2;
7704     int sb, bcd;
7705
7706     sb = cmos_read(CMOS_STATB);
7707
7708     bcd = (sb & (1 << 2)) == 0;
7709
7710     // make sure CMOS doesn't modify time while we read it
7711     for(;;) {
7712         fill_rtcdate(&t1);
7713         if(cmos_read(CMOS_STATA) & CMOS_UIP)
7714             continue;
7715         fill_rtcdate(&t2);
7716         if(memcmp(&t1, &t2, sizeof(t1)) == 0)
7717             break;
7718     }
7719
7720     // convert
7721     if(bcd) {
7722 #define CONV(x)      (t1.x = ((t1.x >> 4) * 10) + (t1.x & 0xf))
7723         CONV(second);
7724         CONV(minute);
7725         CONV(hour );
7726         CONV(day );
7727         CONV(month );
7728         CONV(year );
7729 #undef CONV
7730     }
7731
7732     *r = t1;
7733     r->year += 2000;
7734 }
7735
7736
7737
7738
7739
7740
7741
7742
7743
7744
7745
7746
7747
7748
7749

```

```

7750 // The I/O APIC manages hardware interrupts for an SMP system.
7751 // http://www.intel.com/design/chipsets/datashts/29056601.pdf
7752 // See also picirq.c.
7753
7754 #include "types.h"
7755 #include "defs.h"
7756 #include "traps.h"
7757
7758 #define IOAPIC 0xFEC00000 // Default physical address of IO APIC
7759
7760 #define REG_ID 0x00 // Register index: ID
7761 #define REG_VER 0x01 // Register index: version
7762 #define REG_TABLE 0x10 // Redirection table base
7763
7764 // The redirection table starts at REG_TABLE and uses
7765 // two registers to configure each interrupt.
7766 // The first (low) register in a pair contains configuration bits.
7767 // The second (high) register contains a bitmask telling which
7768 // CPUs can serve that interrupt.
7769 #define INT_DISABLED 0x00010000 // Interrupt disabled
7770 #define INT_LEVEL 0x00008000 // Level-triggered (vs edge-)
7771 #define INT_ACTIVELOW 0x00002000 // Active low (vs high)
7772 #define INT_LOGICAL 0x00000800 // Destination is CPU id (vs APIC ID)
7773
7774 volatile struct ioapic *ioapic;
7775
7776 // IO APIC MMIO structure: write reg, then read or write data.
7777 struct ioapic {
7778     uint reg;
7779     uint pad[3];
7780     uint data;
7781 };
7782
7783 static uint
7784 ioapicread(int reg)
7785 {
7786     ioapic->reg = reg;
7787     return ioapic->data;
7788 }
7789
7790 static void
7791 ioapicwrite(int reg, uint data)
7792 {
7793     ioapic->reg = reg;
7794     ioapic->data = data;
7795 }
7796
7797
7798
7799

```

```

7800 void
7801 ioapicinit(void)
7802 {
7803     int i, id, maxintr;
7804
7805     ioapic = (volatile struct ioapic*)IOAPIC;
7806     maxintr = (ioapicread(REG_VER) >> 16) & 0xFF;
7807     id = ioapicread(REG_ID) >> 24;
7808     if(id != ioapicid)
7809         printf("ioapicinit: id isn't equal to ioapicid; not a MP\n");
7810
7811     // Mark all interrupts edge-triggered, active high, disabled,
7812     // and not routed to any CPUs.
7813     for(i = 0; i <= maxintr; i++){
7814         ioapicwrite(REG_TABLE+2*i, INT_DISABLED | (T_IRQ0 + i));
7815         ioapicwrite(REG_TABLE+2*i+1, 0);
7816     }
7817 }
7818
7819 void
7820 ioapicenable(int irq, int cpunum)
7821 {
7822     // Mark interrupt edge-triggered, active high,
7823     // enabled, and routed to the given cpunum,
7824     // which happens to be that cpu's APIC ID.
7825     ioapicwrite(REG_TABLE+2*irq, T_IRQ0 + irq);
7826     ioapicwrite(REG_TABLE+2*irq+1, cpunum << 24);
7827 }
7828
7829
7830
7831
7832
7833
7834
7835
7836
7837
7838
7839
7840
7841
7842
7843
7844
7845
7846
7847
7848
7849

```

```

7850 // PC keyboard interface constants
7851
7852 #define KBSTAP      0x64    // kbd controller status port(I)
7853 #define KBS_DIB     0x01    // kbd data in buffer
7854 #define KBDATAP     0x60    // kbd data port(I)
7855
7856 #define NO          0
7857
7858 #define SHIFT       (1<<0)
7859 #define CTL         (1<<1)
7860 #define ALT         (1<<2)
7861
7862 #define CAPSLOCK    (1<<3)
7863 #define NUMLOCK     (1<<4)
7864 #define SCROLLLOCK (1<<5)
7865
7866 #define EOESC       (1<<6)
7867
7868 // Special keycodes
7869 #define KEY_HOME    0xE0
7870 #define KEY_END     0xE1
7871 #define KEY_UP      0xE2
7872 #define KEY_DN      0xE3
7873 #define KEY_LF      0xE4
7874 #define KEY_RT      0xE5
7875 #define KEY_PGUP    0xE6
7876 #define KEY_PGDN    0xE7
7877 #define KEY_INS     0xE8
7878 #define KEY_DEL     0xE9
7879
7880 // C('A') == Control-A
7881 #define C(x) (x - '@')
7882
7883 static uchar shiftcode[256] =
7884 {
7885     [0x1D] CTL,
7886     [0x2A] SHIFT,
7887     [0x36] SHIFT,
7888     [0x38] ALT,
7889     [0x9D] CTL,
7890     [0xB8] ALT
7891 };
7892
7893 static uchar togglecode[256] =
7894 {
7895     [0x3A] CAPSLOCK,
7896     [0x45] NUMLOCK,
7897     [0x46] SCROLLLOCK
7898 };
7899

```

```

7900 static uchar normalmap[256] =
7901 {
7902     NO,    0x1B, '1', '2', '3', '4', '5', '6', // 0x00
7903     '7', '8', '9', '0', '-', '=', '\b', '\t',
7904     'q', 'w', 'e', 'r', 't', 'y', 'u', 'i', // 0x10
7905     'o', 'p', '[', ']', '\n', NO, 'a', 's',
7906     'd', 'f', 'g', 'h', 'j', 'k', 'l', ';', // 0x20
7907     '\'', '`', NO, '\\', 'z', 'x', 'c', 'v',
7908     'b', 'n', 'm', ',', '.', '/', NO, '*', // 0x30
7909     NO, ' ', NO, NO, NO, NO, NO, NO,
7910     NO, NO, NO, NO, NO, NO, NO, '7', // 0x40
7911     '8', '9', '-', '4', '5', '6', '+', '1',
7912     '2', '3', '0', '.', NO, NO, NO, NO, // 0x50
7913     [0x9C] '\n', // KP_Enter
7914     [0xB5] '/', // KP_Div
7915     [0xC8] KEY_UP, [0xD0] KEY_DN,
7916     [0xC9] KEY_PGUP, [0xD1] KEY_PGDN,
7917     [0xCB] KEY_LF, [0xCD] KEY_RT,
7918     [0x97] KEY_HOME, [0xCF] KEY_END,
7919     [0xD2] KEY_INS, [0xD3] KEY_DEL
7920 };
7921
7922 static uchar shiftmap[256] =
7923 {
7924     NO,    033, '!', '@', '#', '$', '%', '^', // 0x00
7925     '&', '*', '(', ')', '_', '+', '\b', '\t',
7926     'Q', 'W', 'E', 'R', 'T', 'Y', 'U', 'I', // 0x10
7927     'O', 'P', '{', '}', '\n', NO, 'A', 'S',
7928     'D', 'F', 'G', 'H', 'J', 'K', 'L', ':', // 0x20
7929     '"', '~', NO, '|', 'Z', 'X', 'C', 'V',
7930     'B', 'N', 'M', '<', '>', '?', NO, '*', // 0x30
7931     NO, ' ', NO, NO, NO, NO, NO, NO,
7932     NO, NO, NO, NO, NO, NO, NO, '7', // 0x40
7933     '8', '9', '-', '4', '5', '6', '+', '1',
7934     '2', '3', '0', '.', NO, NO, NO, NO, // 0x50
7935     [0x9C] '\n', // KP_Enter
7936     [0xB5] '/', // KP_Div
7937     [0xC8] KEY_UP, [0xD0] KEY_DN,
7938     [0xC9] KEY_PGUP, [0xD1] KEY_PGDN,
7939     [0xCB] KEY_LF, [0xCD] KEY_RT,
7940     [0x97] KEY_HOME, [0xCF] KEY_END,
7941     [0xD2] KEY_INS, [0xD3] KEY_DEL
7942 };
7943
7944
7945
7946
7947
7948
7949

```

```

7950 static uchar ctlmap[256] =
7951 {
7952     NO,      NO,      NO,      NO,      NO,      NO,      NO,      NO,
7953     NO,      NO,      NO,      NO,      NO,      NO,      NO,      NO,
7954     C('Q'),  C('W'),  C('E'),  C('R'),  C('T'),  C('Y'),  C('U'),  C('I'),
7955     C('O'),  C('P'),  NO,      NO,      '\r',    NO,      C('A'),  C('S'),
7956     C('D'),  C('F'),  C('G'),  C('H'),  C('J'),  C('K'),  C('L'),  NO,
7957     NO,      NO,      NO,      C('\'),  C('Z'),  C('X'),  C('C'),  C('V'),
7958     C('B'),  C('N'),  C('M'),  NO,      NO,      C('/'),  NO,      NO,
7959     [0x9C] '\r',    // KP_Enter
7960     [0xB5] C('/'),  // KP_Div
7961     [0xC8] KEY_UP,  [0xD0] KEY_DN,
7962     [0xC9] KEY_PGUP, [0xD1] KEY_PGDN,
7963     [0xCB] KEY_LF,  [0xCD] KEY_RT,
7964     [0x97] KEY_HOME, [0xCF] KEY_END,
7965     [0xD2] KEY_INS, [0xD3] KEY_DEL
7966 };
7967
7968
7969
7970
7971
7972
7973
7974
7975
7976
7977
7978
7979
7980
7981
7982
7983
7984
7985
7986
7987
7988
7989
7990
7991
7992
7993
7994
7995
7996
7997
7998
7999

```

```

8000 #include "types.h"
8001 #include "x86.h"
8002 #include "defs.h"
8003 #include "kbd.h"
8004
8005 int
8006 kbdgetc(void)
8007 {
8008     static uint shift;
8009     static uchar *charcode[4] = {
8010         normalmap, shiftmap, ctlmap, ctlmap
8011     };
8012     uint st, data, c;
8013
8014     st = inb(KBSTATP);
8015     if((st & KBS_DIB) == 0)
8016         return -1;
8017     data = inb(KBDATAP);
8018
8019     if(data == 0xE0){
8020         shift |= EOESC;
8021         return 0;
8022     } else if(data & 0x80){
8023         // Key released
8024         data = (shift & EOESC ? data : data & 0x7F);
8025         shift &= ~(shiftcode[data] | EOESC);
8026         return 0;
8027     } else if(shift & EOESC){
8028         // Last character was an E0 escape; or with 0x80
8029         data |= 0x80;
8030         shift &= ~EOESC;
8031     }
8032
8033     shift |= shiftcode[data];
8034     shift ^= togglecode[data];
8035     c = charcode[shift & (CTL | SHIFT)][data];
8036     if(shift & CAPSLOCK){
8037         if('a' <= c && c <= 'z')
8038             c += 'A' - 'a';
8039         else if('A' <= c && c <= 'Z')
8040             c += 'a' - 'A';
8041     }
8042     return c;
8043 }
8044
8045 void
8046 kbdtintr(void)
8047 {
8048     consoleintr(kbdgetc);
8049 }

```



```

8050 // Console input and output.
8051 // Input is from the keyboard or serial port.
8052 // Output is written to the screen and serial port.
8053
8054 #include "types.h"
8055 #include "defs.h"
8056 #include "param.h"
8057 #include "traps.h"
8058 #include "spinlock.h"
8059 #include "sleeplock.h"
8060 #include "fs.h"
8061 #include "file.h"
8062 #include "memlayout.h"
8063 #include "mmu.h"
8064 #include "proc.h"
8065 #include "x86.h"
8066
8067 static void consputc(int);
8068
8069 static int panicked = 0;
8070
8071 static struct {
8072     struct spinlock lock;
8073     int locking;
8074 } cons;
8075
8076 static void
8077 printint(int xx, int base, int sign)
8078 {
8079     static char digits[] = "0123456789abcdef";
8080     char buf[16];
8081     int i;
8082     uint x;
8083
8084     if(sign && (sign = xx < 0))
8085         x = -xx;
8086     else
8087         x = xx;
8088
8089     i = 0;
8090     do{
8091         buf[i++] = digits[x % base];
8092     }while((x /= base) != 0);
8093
8094     if(sign)
8095         buf[i++] = '-';
8096
8097     while(--i >= 0)
8098         consputc(buf[i]);
8099 }

```

```

8100
8101
8102
8103
8104
8105
8106
8107
8108
8109
8110
8111
8112
8113
8114
8115
8116
8117
8118
8119
8120
8121
8122
8123
8124
8125
8126
8127
8128
8129
8130
8131
8132
8133
8134
8135
8136
8137
8138
8139
8140
8141
8142
8143
8144
8145
8146
8147
8148
8149

```

```

8150 // Print to the console. only understands %d, %x, %p, %s.
8151 void
8152 cprintf(char *fmt, ...)
8153 {
8154     int i, c, locking;
8155     uint *argp;
8156     char *s;
8157
8158     locking = cons.locking;
8159     if(locking)
8160         acquire(&cons.lock);
8161
8162     if (fmt == 0)
8163         panic("null fmt");
8164
8165     argp = (uint*)(void*)&fmt + 1;
8166     for(i = 0; (c = fmt[i] & 0xff) != 0; i++){
8167         if(c != '%'){
8168             consputc(c);
8169             continue;
8170         }
8171         c = fmt[++i] & 0xff;
8172         if(c == 0)
8173             break;
8174         switch(c){
8175             case 'd':
8176                 printint(*argp++, 10, 1);
8177                 break;
8178             case 'x':
8179             case 'p':
8180                 printint(*argp++, 16, 0);
8181                 break;
8182             case 's':
8183                 if((s = (char*)*argp++) == 0)
8184                     s = "(null)";
8185                 for(; *s; s++)
8186                     consputc(*s);
8187                 break;
8188             case '%':
8189                 consputc('%');
8190                 break;
8191             default:
8192                 // Print unknown % sequence to draw attention.
8193                 consputc('%');
8194                 consputc(c);
8195                 break;
8196         }
8197     }
8198 }
8199

```

```

8200     if(locking)
8201         release(&cons.lock);
8202 }
8203
8204 void
8205 panic(char *s)
8206 {
8207     int i;
8208     uint pcs[10];
8209
8210     cli();
8211     cons.locking = 0;
8212     // use lapiccpunum so that we can call panic from mycpu()
8213     cprintf("lapicid %d: panic: ", lapicid());
8214     cprintf(s);
8215     cprintf("\n");
8216     getcallerpcs(&s, pcs);
8217     for(i=0; i<10; i++)
8218         cprintf(" %p", pcs[i]);
8219     panicked = 1; // freeze other CPU
8220     for(;;)
8221         ;
8222 }
8223
8224
8225
8226
8227
8228
8229
8230
8231
8232
8233
8234
8235
8236
8237
8238
8239
8240
8241
8242
8243
8244
8245
8246
8247
8248
8249

```

```

8250 #define BACKSPACE 0x100
8251 #define CRTPORT 0x3d4
8252 static ushort *crt = (ushort*)P2V(0xb8000); // CGA memory
8253
8254 static void
8255 cgaputc(int c)
8256 {
8257     int pos;
8258
8259     // Cursor position: col + 80*row.
8260     outb(CRTPORT, 14);
8261     pos = inb(CRTPORT+1) << 8;
8262     outb(CRTPORT, 15);
8263     pos |= inb(CRTPORT+1);
8264
8265     if(c == '\n')
8266         pos += 80 - pos%80;
8267     else if(c == BACKSPACE){
8268         if(pos > 0) --pos;
8269     } else
8270         crt[pos++] = (c&0xff) | 0x0700; // black on white
8271
8272     if(pos < 0 || pos > 25*80)
8273         panic("pos under/overflow");
8274
8275     if((pos/80) >= 24){ // Scroll up.
8276         memmove(crt, crt+80, sizeof(crt[0])*23*80);
8277         pos -= 80;
8278         memset(crt+pos, 0, sizeof(crt[0])*(24*80 - pos));
8279     }
8280
8281     outb(CRTPORT, 14);
8282     outb(CRTPORT+1, pos>>8);
8283     outb(CRTPORT, 15);
8284     outb(CRTPORT+1, pos);
8285     crt[pos] = ' ' | 0x0700;
8286 }
8287
8288
8289
8290
8291
8292
8293
8294
8295
8296
8297
8298
8299

```

```

8300 void
8301 consputc(int c)
8302 {
8303     if(panicked){
8304         cli();
8305         for(;;)
8306             ;
8307     }
8308
8309     if(c == BACKSPACE){
8310         uartputc('\b'); uartputc(' '); uartputc('\b');
8311     } else
8312         uartputc(c);
8313     cgaputc(c);
8314 }
8315
8316 #define INPUT_BUF 128
8317 struct {
8318     char buf[INPUT_BUF];
8319     uint r; // Read index
8320     uint w; // Write index
8321     uint e; // Edit index
8322 } input;
8323
8324 #define C(x) ((x)-'@') // Control-x
8325
8326 void
8327 consoleintr(int (*getc)(void))
8328 {
8329     int c, doprocdump = 0;
8330
8331     acquire(&cons.lock);
8332     while((c = getc()) >= 0){
8333         switch(c){
8334             case C('P'): // Process listing.
8335                 // procdump() locks cons.lock indirectly; invoke later
8336                 doprocdump = 1;
8337                 break;
8338             case C('U'): // Kill line.
8339                 while(input.e != input.w &&
8340                        input.buf[(input.e-1) % INPUT_BUF] != '\n'){
8341                     input.e--;
8342                     consputc(BACKSPACE);
8343                 }
8344                 break;
8345             case C('H'): case '\x7f': // Backspace
8346                 if(input.e != input.w){
8347                     input.e--;
8348                     consputc(BACKSPACE);
8349                 }

```

```

8350     break;
8351     default:
8352     if(c != 0 && input.e-input.r < INPUT_BUF){
8353         c = (c == '\r') ? '\n' : c;
8354         input.buf[input.e++ % INPUT_BUF] = c;
8355         consputc(c);
8356         if(c == '\n' || c == C('D') || input.e == input.r+INPUT_BUF){
8357             input.w = input.e;
8358             wakeup(&input.r);
8359         }
8360     }
8361     break;
8362 }
8363 }
8364 release(&cons.lock);
8365 if(doprocDump) {
8366     procdump(); // now call procdump() wo. cons.lock held
8367 }
8368 }
8369
8370 int
8371 consoleread(struct inode *ip, char *dst, int n)
8372 {
8373     uint target;
8374     int c;
8375
8376     iunlock(ip);
8377     target = n;
8378     acquire(&cons.lock);
8379     while(n > 0){
8380         while(input.r == input.w){
8381             if(myproc()->killed){
8382                 release(&cons.lock);
8383                 ilock(ip);
8384                 return -1;
8385             }
8386             sleep(&input.r, &cons.lock);
8387         }
8388         c = input.buf[input.r++ % INPUT_BUF];
8389         if(c == C('D')){ // EOF
8390             if(n < target){
8391                 // Save ^D for next time, to make sure
8392                 // caller gets a 0-byte result.
8393                 input.r--;
8394             }
8395             break;
8396         }
8397         *dst++ = c;
8398         --n;
8399         if(c == '\n')

```

```

8400     break;
8401 }
8402 release(&cons.lock);
8403 ilock(ip);
8404
8405 return target - n;
8406 }
8407
8408 int
8409 consolewrite(struct inode *ip, char *buf, int n)
8410 {
8411     int i;
8412
8413     iunlock(ip);
8414     acquire(&cons.lock);
8415     for(i = 0; i < n; i++){
8416         consputc(buf[i] & 0xff);
8417     }
8418     release(&cons.lock);
8419     ilock(ip);
8420
8421     return n;
8422 }
8423 void
8424 consoleinit(void)
8425 {
8426     initlock(&cons.lock, "console");
8427
8428     devsw[CONSOLE].write = consolewrite;
8429     devsw[CONSOLE].read = consoleread;
8430     cons.locking = 1;
8431
8432     ioapicenable(IRQ_KBD, 0);
8433 }
8434
8435
8436
8437
8438
8439
8440
8441
8442
8443
8444
8445
8446
8447
8448
8449

```

```

8450 // Intel 8250 serial port (UART).
8451
8452 #include "types.h"
8453 #include "defs.h"
8454 #include "param.h"
8455 #include "traps.h"
8456 #include "spinlock.h"
8457 #include "sleeplock.h"
8458 #include "fs.h"
8459 #include "file.h"
8460 #include "mmu.h"
8461 #include "proc.h"
8462 #include "x86.h"
8463
8464 #define COM1    0x3f8
8465
8466 static int uart;    // is there a uart?
8467
8468 void
8469 uartinit(void)
8470 {
8471     char *p;
8472
8473     // Turn off the FIFO
8474     outb(COM1+2, 0);
8475
8476     // 9600 baud, 8 data bits, 1 stop bit, parity off.
8477     outb(COM1+3, 0x80);    // Unlock divisor
8478     outb(COM1+0, 115200/9600);
8479     outb(COM1+1, 0);
8480     outb(COM1+3, 0x03);    // Lock divisor, 8 data bits.
8481     outb(COM1+4, 0);
8482     outb(COM1+1, 0x01);    // Enable receive interrupts.
8483
8484     // If status is 0xFF, no serial port.
8485     if(inb(COM1+5) == 0xFF)
8486         return;
8487     uart = 1;
8488
8489     // Acknowledge pre-existing interrupt conditions;
8490     // enable interrupts.
8491     inb(COM1+2);
8492     inb(COM1+0);
8493     ioapicenable(IRQ_COM1, 0);
8494
8495     // Announce that we're here.
8496     for(p="xv6...\n"; *p; p++)
8497         uartputc(*p);
8498 }
8499

```

```

8500 void
8501 uartputc(int c)
8502 {
8503     int i;
8504
8505     if(!uart)
8506         return;
8507     for(i = 0; i < 128 && !(inb(COM1+5) & 0x20); i++)
8508         microdelay(10);
8509     outb(COM1+0, c);
8510 }
8511
8512 static int
8513 uartgetc(void)
8514 {
8515     if(!uart)
8516         return -1;
8517     if(!(inb(COM1+5) & 0x01))
8518         return -1;
8519     return inb(COM1+0);
8520 }
8521
8522 void
8523 uartintr(void)
8524 {
8525     consoleintr(uartgetc);
8526 }
8527
8528
8529
8530
8531
8532
8533
8534
8535
8536
8537
8538
8539
8540
8541
8542
8543
8544
8545
8546
8547
8548
8549

```

```

8550 # Initial process execs /init.
8551 # This code runs in user space.
8552
8553 #include "syscall.h"
8554 #include "traps.h"
8555
8556
8557 # exec(init, argv)
8558 .globl start
8559 start:
8560     pushl $argv
8561     pushl $init
8562     pushl $0 // where caller pc would be
8563     movl $SYS_exec, %eax
8564     int $T_SYSCALL
8565
8566 # for(;;) exit();
8567 exit:
8568     movl $SYS_exit, %eax
8569     int $T_SYSCALL
8570     jmp exit
8571
8572 # char init[] = "/init\0";
8573 init:
8574     .string "/init\0"
8575
8576 # char *argv[] = { init, 0 };
8577 .p2align 2
8578 argv:
8579     .long init
8580     .long 0
8581
8582
8583
8584
8585
8586
8587
8588
8589
8590
8591
8592
8593
8594
8595
8596
8597
8598
8599

```

```

8600 #include "syscall.h"
8601 #include "traps.h"
8602
8603 #define SYSCALL(name) \
8604     .globl name; \
8605     name: \
8606     movl $SYS_ ## name, %eax; \
8607     int $T_SYSCALL; \
8608     ret
8609
8610 SYSCALL(fork)
8611 SYSCALL(exit)
8612 SYSCALL(wait)
8613 SYSCALL(pipe)
8614 SYSCALL(read)
8615 SYSCALL(write)
8616 SYSCALL(close)
8617 SYSCALL(kill)
8618 SYSCALL(exec)
8619 SYSCALL(open)
8620 SYSCALL(mknod)
8621 SYSCALL(unlink)
8622 SYSCALL(fstat)
8623 SYSCALL(link)
8624 SYSCALL(mkdir)
8625 SYSCALL(chdir)
8626 SYSCALL(dup)
8627 SYSCALL(getpid)
8628 SYSCALL(sbrk)
8629 SYSCALL(sleep)
8630 SYSCALL(uptime)
8631
8632
8633
8634
8635
8636
8637
8638
8639
8640
8641
8642
8643
8644
8645
8646
8647
8648
8649

```

```

8650 // init: The initial user-level program
8651
8652 #include "types.h"
8653 #include "stat.h"
8654 #include "user.h"
8655 #include "fcntl.h"
8656
8657 char *argv[] = { "sh", 0 };
8658
8659 int
8660 main(void)
8661 {
8662     int pid, wpid;
8663
8664     if(open("console", O_RDWR) < 0){
8665         mknod("console", 1, 1);
8666         open("console", O_RDWR);
8667     }
8668     dup(0); // stdout
8669     dup(0); // stderr
8670
8671     for(;;){
8672         printf(1, "init: starting sh\n");
8673         pid = fork();
8674         if(pid < 0){
8675             printf(1, "init: fork failed\n");
8676             exit();
8677         }
8678         if(pid == 0){
8679             exec("sh", argv);
8680             printf(1, "init: exec sh failed\n");
8681             exit();
8682         }
8683         while((wpid=wait()) >= 0 && wpid != pid)
8684             printf(1, "zombie!\n");
8685     }
8686 }
8687
8688
8689
8690
8691
8692
8693
8694
8695
8696
8697
8698
8699

```

```

8700 // Shell.
8701
8702 #include "types.h"
8703 #include "user.h"
8704 #include "fcntl.h"
8705
8706 // Parsed command representation
8707 #define EXEC 1
8708 #define REDIR 2
8709 #define PIPE 3
8710 #define LIST 4
8711 #define BACK 5
8712
8713 #define MAXARGS 10
8714
8715 struct cmd {
8716     int type;
8717 };
8718
8719 struct execcmd {
8720     int type;
8721     char *argv[MAXARGS];
8722     char *eargv[MAXARGS];
8723 };
8724
8725 struct redircmd {
8726     int type;
8727     struct cmd *cmd;
8728     char *file;
8729     char *efile;
8730     int mode;
8731     int fd;
8732 };
8733
8734 struct pipecmd {
8735     int type;
8736     struct cmd *left;
8737     struct cmd *right;
8738 };
8739
8740 struct listcmd {
8741     int type;
8742     struct cmd *left;
8743     struct cmd *right;
8744 };
8745
8746 struct backcmd {
8747     int type;
8748     struct cmd *cmd;
8749 };

```

```

8750 int fork1(void); // Fork but panics on failure.
8751 void panic(char*);
8752 struct cmd *parsecmd(char*);
8753
8754 // Execute cmd. Never returns.
8755 void
8756 runcmd(struct cmd *cmd)
8757 {
8758     int p[2];
8759     struct backcmd *bcmd;
8760     struct execcmd *ecmd;
8761     struct listcmd *lcmd;
8762     struct pipecmd *pcmd;
8763     struct redircmd *rcmd;
8764
8765     if(cmd == 0)
8766         exit();
8767
8768     switch(cmd->type){
8769     default:
8770         panic("runcmd");
8771
8772     case EXEC:
8773         ecmd = (struct execcmd*)cmd;
8774         if(ecmd->argv[0] == 0)
8775             exit();
8776         exec(ecmd->argv[0], ecmd->argv);
8777         printf(2, "exec %s failed\n", ecmd->argv[0]);
8778         break;
8779
8780     case REDIR:
8781         rcmd = (struct redircmd*)cmd;
8782         close(rcmd->fd);
8783         if(open(rcmd->file, rcmd->mode) < 0){
8784             printf(2, "open %s failed\n", rcmd->file);
8785             exit();
8786         }
8787         runcmd(rcmd->cmd);
8788         break;
8789
8790     case LIST:
8791         lcmd = (struct listcmd*)cmd;
8792         if(fork1() == 0)
8793             runcmd(lcmd->left);
8794         wait();
8795         runcmd(lcmd->right);
8796         break;
8797
8798
8799

```

```

8800     case PIPE:
8801         pcmd = (struct pipecmd*)cmd;
8802         if(pipe(p) < 0)
8803             panic("pipe");
8804         if(fork1() == 0){
8805             close(1);
8806             dup(p[1]);
8807             close(p[0]);
8808             close(p[1]);
8809             runcmd(pcmd->left);
8810         }
8811         if(fork1() == 0){
8812             close(0);
8813             dup(p[0]);
8814             close(p[0]);
8815             close(p[1]);
8816             runcmd(pcmd->right);
8817         }
8818         close(p[0]);
8819         close(p[1]);
8820         wait();
8821         wait();
8822         break;
8823
8824     case BACK:
8825         bcmd = (struct backcmd*)cmd;
8826         if(fork1() == 0)
8827             runcmd(bcmd->cmd);
8828         break;
8829     }
8830     exit();
8831 }
8832
8833 int
8834 getcmd(char *buf, int nbuf)
8835 {
8836     printf(2, "$ ");
8837     memset(buf, 0, nbuf);
8838     gets(buf, nbuf);
8839     if(buf[0] == 0) // EOF
8840         return -1;
8841     return 0;
8842 }
8843
8844
8845
8846
8847
8848
8849

```



```

8850 int
8851 main(void)
8852 {
8853     static char buf[100];
8854     int fd;
8855
8856     // Ensure that three file descriptors are open.
8857     while((fd = open("console", O_RDWR)) >= 0){
8858         if(fd >= 3){
8859             close(fd);
8860             break;
8861         }
8862     }
8863
8864     // Read and run input commands.
8865     while(getcmd(buf, sizeof(buf)) >= 0){
8866         if(buf[0] == 'c' && buf[1] == 'd' && buf[2] == ' '){
8867             // Chdir must be called by the parent, not the child.
8868             buf[strlen(buf)-1] = 0; // chop \n
8869             if(chdir(buf+3) < 0)
8870                 printf(2, "cannot cd %s\n", buf+3);
8871             continue;
8872         }
8873         if(fork1() == 0)
8874             runcmd(parsecmd(buf));
8875         wait();
8876     }
8877     exit();
8878 }
8879
8880 void
8881 panic(char *s)
8882 {
8883     printf(2, "%s\n", s);
8884     exit();
8885 }
8886
8887 int
8888 fork1(void)
8889 {
8890     int pid;
8891
8892     pid = fork();
8893     if(pid == -1)
8894         panic("fork");
8895     return pid;
8896 }
8897
8898
8899

```

```

8900 // Constructors
8901
8902 struct cmd*
8903 execcmd(void)
8904 {
8905     struct execcmd *cmd;
8906
8907     cmd = malloc(sizeof(*cmd));
8908     memset(cmd, 0, sizeof(*cmd));
8909     cmd->type = EXEC;
8910     return (struct cmd*)cmd;
8911 }
8912
8913 struct cmd*
8914 redircmd(struct cmd *subcmd, char *file, char *efile, int mode, int fd)
8915 {
8916     struct redircmd *cmd;
8917
8918     cmd = malloc(sizeof(*cmd));
8919     memset(cmd, 0, sizeof(*cmd));
8920     cmd->type = REDIR;
8921     cmd->cmd = subcmd;
8922     cmd->file = file;
8923     cmd->efile = efile;
8924     cmd->mode = mode;
8925     cmd->fd = fd;
8926     return (struct cmd*)cmd;
8927 }
8928
8929 struct cmd*
8930 pipecmd(struct cmd *left, struct cmd *right)
8931 {
8932     struct pipecmd *cmd;
8933
8934     cmd = malloc(sizeof(*cmd));
8935     memset(cmd, 0, sizeof(*cmd));
8936     cmd->type = PIPE;
8937     cmd->left = left;
8938     cmd->right = right;
8939     return (struct cmd*)cmd;
8940 }
8941
8942
8943
8944
8945
8946
8947
8948
8949

```

```

8950 struct cmd*
8951 listcmd(struct cmd *left, struct cmd *right)
8952 {
8953     struct listcmd *cmd;
8954
8955     cmd = malloc(sizeof(*cmd));
8956     memset(cmd, 0, sizeof(*cmd));
8957     cmd->type = LIST;
8958     cmd->left = left;
8959     cmd->right = right;
8960     return (struct cmd*)cmd;
8961 }
8962
8963 struct cmd*
8964 backcmd(struct cmd *subcmd)
8965 {
8966     struct backcmd *cmd;
8967
8968     cmd = malloc(sizeof(*cmd));
8969     memset(cmd, 0, sizeof(*cmd));
8970     cmd->type = BACK;
8971     cmd->cmd = subcmd;
8972     return (struct cmd*)cmd;
8973 }
8974
8975
8976
8977
8978
8979
8980
8981
8982
8983
8984
8985
8986
8987
8988
8989
8990
8991
8992
8993
8994
8995
8996
8997
8998
8999

```

```

9000 // Parsing
9001
9002 char whitespace[] = " \t\r\n\v";
9003 char symbols[] = "<|>&()";
9004
9005 int
9006 gettoken(char **ps, char *es, char **q, char **eq)
9007 {
9008     char *s;
9009     int ret;
9010
9011     s = *ps;
9012     while(s < es && strchr(whitespace, *s))
9013         s++;
9014     if(q)
9015         *q = s;
9016     ret = *s;
9017     switch(*s){
9018     case 0:
9019         break;
9020     case '|':
9021     case '(':
9022     case ')':
9023     case ';':
9024     case '&':
9025     case '<':
9026         s++;
9027         break;
9028     case '>':
9029         s++;
9030         if(*s == '>'){
9031             ret = '+';
9032             s++;
9033         }
9034         break;
9035     default:
9036         ret = 'a';
9037         while(s < es && !strchr(whitespace, *s) && !strchr(symbols, *s))
9038             s++;
9039         break;
9040     }
9041     if(eq)
9042         *eq = s;
9043
9044     while(s < es && strchr(whitespace, *s))
9045         s++;
9046     *ps = s;
9047     return ret;
9048 }
9049

```

```

9050 int
9051 peek(char **ps, char *es, char *toks)
9052 {
9053     char *s;
9054
9055     s = *ps;
9056     while(s < es && strchr(whitespace, *s))
9057         s++;
9058     *ps = s;
9059     return *s && strchr(toks, *s);
9060 }
9061
9062 struct cmd *parseline(char**, char*);
9063 struct cmd *parsepipe(char**, char*);
9064 struct cmd *parseexec(char**, char*);
9065 struct cmd *nulterminate(struct cmd*);
9066
9067 struct cmd*
9068 parsecmd(char *s)
9069 {
9070     char *es;
9071     struct cmd *cmd;
9072
9073     es = s + strlen(s);
9074     cmd = parseline(&s, es);
9075     peek(&s, es, "");
9076     if(s != es){
9077         printf(2, "leftovers: %s\n", s);
9078         panic("syntax");
9079     }
9080     nulterminate(cmd);
9081     return cmd;
9082 }
9083
9084 struct cmd*
9085 parseline(char **ps, char *es)
9086 {
9087     struct cmd *cmd;
9088
9089     cmd = parsepipe(ps, es);
9090     while(peek(ps, es, "&")){
9091         gettoken(ps, es, 0, 0);
9092         cmd = backcmd(cmd);
9093     }
9094     if(peek(ps, es, ";")){
9095         gettoken(ps, es, 0, 0);
9096         cmd = listcmd(cmd, parseline(ps, es));
9097     }
9098     return cmd;
9099 }

```

```

9100 struct cmd*
9101 parsepipe(char **ps, char *es)
9102 {
9103     struct cmd *cmd;
9104
9105     cmd = parseexec(ps, es);
9106     if(peek(ps, es, "|")){
9107         gettoken(ps, es, 0, 0);
9108         cmd = pipecmd(cmd, parsepipe(ps, es));
9109     }
9110     return cmd;
9111 }
9112
9113 struct cmd*
9114 parseredirs(struct cmd *cmd, char **ps, char *es)
9115 {
9116     int tok;
9117     char *q, *eq;
9118
9119     while(peek(ps, es, "<>")){
9120         tok = gettoken(ps, es, 0, 0);
9121         if(gettoken(ps, es, &q, &eq) != 'a')
9122             panic("missing file for redirection");
9123         switch(tok){
9124             case '<':
9125                 cmd = redircmd(cmd, q, eq, O_RDONLY, 0);
9126                 break;
9127             case '>':
9128                 cmd = redircmd(cmd, q, eq, O_WRONLY|O_CREATE, 1);
9129                 break;
9130             case '+': // >>
9131                 cmd = redircmd(cmd, q, eq, O_WRONLY|O_CREATE, 1);
9132                 break;
9133         }
9134     }
9135     return cmd;
9136 }
9137
9138
9139
9140
9141
9142
9143
9144
9145
9146
9147
9148
9149

```

```

9150 struct cmd*
9151 parseblock(char **ps, char *es)
9152 {
9153     struct cmd *cmd;
9154
9155     if(!peek(ps, es, "("))
9156         panic("parseblock");
9157     gettoken(ps, es, 0, 0);
9158     cmd = parseline(ps, es);
9159     if(!peek(ps, es, "|"))
9160         panic("syntax - missing ");
9161     gettoken(ps, es, 0, 0);
9162     cmd = parseredirs(cmd, ps, es);
9163     return cmd;
9164 }
9165
9166 struct cmd*
9167 parseexec(char **ps, char *es)
9168 {
9169     char *q, *eq;
9170     int tok, argc;
9171     struct execcmd *cmd;
9172     struct cmd *ret;
9173
9174     if(peek(ps, es, "("))
9175         return parseblock(ps, es);
9176
9177     ret = execcmd();
9178     cmd = (struct execcmd*)ret;
9179
9180     argc = 0;
9181     ret = parseredirs(ret, ps, es);
9182     while(!peek(ps, es, "|)&");){
9183         if((tok=gettoken(ps, es, &q, &eq)) == 0)
9184             break;
9185         if(tok != 'a')
9186             panic("syntax");
9187         cmd->argv[argc] = q;
9188         cmd->eargv[argc] = eq;
9189         argc++;
9190         if(argc >= MAXARGS)
9191             panic("too many args");
9192         ret = parseredirs(ret, ps, es);
9193     }
9194     cmd->argv[argc] = 0;
9195     cmd->eargv[argc] = 0;
9196     return ret;
9197 }
9198
9199

```

```

9200 // NUL-terminate all the counted strings.
9201 struct cmd*
9202 nulterminate(struct cmd *cmd)
9203 {
9204     int i;
9205     struct backcmd *bcmd;
9206     struct execcmd *ecmd;
9207     struct listcmd *lcmd;
9208     struct pipecmd *pcmd;
9209     struct redircmd *rcmd;
9210
9211     if(cmd == 0)
9212         return 0;
9213
9214     switch(cmd->type){
9215     case EXEC:
9216         ecmd = (struct execcmd*)cmd;
9217         for(i=0; ecmd->argv[i]; i++)
9218             *ecmd->eargv[i] = 0;
9219         break;
9220
9221     case REDIR:
9222         rcmd = (struct redircmd*)cmd;
9223         nulterminate(rcmd->cmd);
9224         *rcmd->efile = 0;
9225         break;
9226
9227     case PIPE:
9228         pcmd = (struct pipecmd*)cmd;
9229         nulterminate(pcmd->left);
9230         nulterminate(pcmd->right);
9231         break;
9232
9233     case LIST:
9234         lcmd = (struct listcmd*)cmd;
9235         nulterminate(lcmd->left);
9236         nulterminate(lcmd->right);
9237         break;
9238
9239     case BACK:
9240         bcmd = (struct backcmd*)cmd;
9241         nulterminate(bcmd->cmd);
9242         break;
9243     }
9244     return cmd;
9245 }
9246
9247
9248
9249

```

```

9250 #include "asm.h"
9251 #include "memlayout.h"
9252 #include "mmu.h"
9253
9254 # Start the first CPU: switch to 32-bit protected mode, jump into C.
9255 # The BIOS loads this code from the first sector of the hard disk into
9256 # memory at physical address 0x7c00 and starts executing in real mode
9257 # with %cs=0 %ip=7c00.
9258
9259 .code16                # Assemble for 16-bit mode
9260 .globl start
9261 start:
9262     cli                # BIOS enabled interrupts; disable
9263
9264     # Zero data segment registers DS, ES, and SS.
9265     xorw    %ax,%ax    # Set %ax to zero
9266     movw   %ax,%ds    # -> Data Segment
9267     movw   %ax,%es    # -> Extra Segment
9268     movw   %ax,%ss    # -> Stack Segment
9269
9270     # Physical address line A20 is tied to zero so that the first PCs
9271     # with 2 MB would run software that assumed 1 MB. Undo that.
9272 seta20.1:
9273     inb    $0x64,%al    # Wait for not busy
9274     testb  $0x2,%al
9275     jnz    seta20.1
9276
9277     movb   $0xd1,%al    # 0xd1 -> port 0x64
9278     outb   %al,$0x64
9279
9280 seta20.2:
9281     inb    $0x64,%al    # Wait for not busy
9282     testb  $0x2,%al
9283     jnz    seta20.2
9284
9285     movb   $0xdf,%al    # 0xdf -> port 0x60
9286     outb   %al,$0x60
9287
9288     # Switch from real to protected mode. Use a bootstrap GDT that makes
9289     # virtual addresses map directly to physical addresses so that the
9290     # effective memory map doesn't change during the transition.
9291     lgdt   gdtdesc
9292     movl   %cr0, %eax
9293     orl   $CR0_PE, %eax
9294     movl   %eax, %cr0
9295
9296
9297
9298
9299

```

```

9300     # Complete the transition to 32-bit protected mode by using a long jmp
9301     # to reload %cs and %eip. The segment descriptors are set up with no
9302     # translation, so that the mapping is still the identity mapping.
9303     ljmp   $(SEG_KCODE<<3), $start32
9304
9305 .code32 # Tell assembler to generate 32-bit code now.
9306 start32:
9307     # Set up the protected-mode data segment registers
9308     movw   $(SEG_KDATA<<3), %ax    # Our data segment selector
9309     movw   %ax, %ds                # -> DS: Data Segment
9310     movw   %ax, %es                # -> ES: Extra Segment
9311     movw   %ax, %ss                # -> SS: Stack Segment
9312     movw   $0, %ax                 # Zero segments not ready for use
9313     movw   %ax, %fs                # -> FS
9314     movw   %ax, %gs                # -> GS
9315
9316     # Set up the stack pointer and call into C.
9317     movl   $start, %esp
9318     call   bootmain
9319
9320     # If bootmain returns (it shouldn't), trigger a Bochs
9321     # breakpoint if running under Bochs, then loop.
9322     movw   $0x8a00, %ax            # 0x8a00 -> port 0x8a00
9323     movw   %ax, %dx
9324     outw   %ax, %dx
9325     movw   $0x8ae0, %ax           # 0x8ae0 -> port 0x8a00
9326     outw   %ax, %dx
9327 spin:
9328     jmp    spin
9329
9330 # Bootstrap GDT
9331 .p2align 2                # force 4 byte alignment
9332 gdt:
9333     SEG_NULLASM            # null seg
9334     SEG_ASM(STA_X|STA_R, 0x0, 0xffffffff) # code seg
9335     SEG_ASM(STA_W, 0x0, 0xffffffff)      # data seg
9336
9337 gdtdesc:
9338     .word   (gdtdesc - gdt - 1)         # sizeof(gdt) - 1
9339     .long   gdt                          # address gdt
9340
9341
9342
9343
9344
9345
9346
9347
9348
9349

```

```

9350 // Boot loader.
9351 //
9352 // Part of the boot block, along with bootasm.S, which calls bootmain().
9353 // bootasm.S has put the processor into protected 32-bit mode.
9354 // bootmain() loads an ELF kernel image from the disk starting at
9355 // sector 1 and then jumps to the kernel entry routine.
9356
9357 #include "types.h"
9358 #include "elf.h"
9359 #include "x86.h"
9360 #include "memlayout.h"
9361
9362 #define SECTSIZE 512
9363
9364 void readseg(uchar*, uint, uint);
9365
9366 void
9367 bootmain(void)
9368 {
9369     struct elfhdr *elf;
9370     struct proghdr *ph, *eph;
9371     void (*entry)(void);
9372     uchar* pa;
9373
9374     elf = (struct elfhdr*)0x10000; // scratch space
9375
9376     // Read 1st page off disk
9377     readseg((uchar*)elf, 4096, 0);
9378
9379     // Is this an ELF executable?
9380     if(elf->magic != ELF_MAGIC)
9381         return; // let bootasm.S handle error
9382
9383     // Load each program segment (ignores ph flags).
9384     ph = (struct proghdr*)((uchar*)elf + elf->phoff);
9385     eph = ph + elf->phnum;
9386     for(; ph < eph; ph++){
9387         pa = (uchar*)ph->paddr;
9388         readseg(pa, ph->filesz, ph->off);
9389         if(ph->memsz > ph->filesz)
9390             stosb(pa + ph->filesz, 0, ph->memsz - ph->filesz);
9391     }
9392
9393     // Call the entry point from the ELF header.
9394     // Does not return!
9395     entry = (void(*) (void))(elf->entry);
9396     entry();
9397 }
9398
9399

```

```

9400 void
9401 waitdisk(void)
9402 {
9403     // Wait for disk ready.
9404     while((inb(0x1F7) & 0xC0) != 0x40)
9405         ;
9406 }
9407
9408 // Read a single sector at offset into dst.
9409 void
9410 readsect(void *dst, uint offset)
9411 {
9412     // Issue command.
9413     waitdisk();
9414     outb(0x1F2, 1); // count = 1
9415     outb(0x1F3, offset);
9416     outb(0x1F4, offset >> 8);
9417     outb(0x1F5, offset >> 16);
9418     outb(0x1F6, (offset >> 24) | 0xE0);
9419     outb(0x1F7, 0x20); // cmd 0x20 - read sectors
9420
9421     // Read data.
9422     waitdisk();
9423     insl(0x1F0, dst, SECTSIZE/4);
9424 }
9425
9426 // Read 'count' bytes at 'offset' from kernel into physical address 'pa'.
9427 // Might copy more than asked.
9428 void
9429 readseg(uchar* pa, uint count, uint offset)
9430 {
9431     uchar* epa;
9432
9433     epa = pa + count;
9434
9435     // Round down to sector boundary.
9436     pa -= offset % SECTSIZE;
9437
9438     // Translate from bytes to sectors; kernel starts at sector 1.
9439     offset = (offset / SECTSIZE) + 1;
9440
9441     // If this is too slow, we could read lots of sectors at a time.
9442     // We'd write more to memory than asked, but it doesn't matter --
9443     // we load in increasing order.
9444     for(; pa < epa; pa += SECTSIZE, offset++){
9445         readsect(pa, offset);
9446     }
9447
9448
9449

```