

implemented in hardware. Another is the `started` variable in `main.c` (`kernel/main.c:7`), used to prevent other CPUs from running until CPU zero has finished initializing xv6; the `volatile` ensures that the compiler actually generates load and store instructions. A third are some uses of `p->parent` in `proc.c` (`kernel/proc.c:398`) (`kernel/proc.c:306`) where proper locking could deadlock, but it seems clear that no other process could be simultaneously modifying `p->parent`. A fourth example is `p->killed`, which is set while holding `p->lock` (`kernel/proc.c:611`), but checked without a holding lock (`kernel/trap.c:56`).

Xv6 contains cases in which one CPU or thread writes some data, and another CPU or thread reads the data, but there is no specific lock dedicated to protecting that data. For example, in `fork`, the parent writes the child's user memory pages, and the child (a different thread, perhaps on a different CPU) reads those pages; no lock explicitly protects those pages. This is not strictly a locking problem, since the child doesn't start executing until after the parent has finished writing. It is a potential memory ordering problem (see Chapter 6), since without a memory barrier there's no reason to expect one CPU to see another CPU's writes. However, since the parent releases locks, and the child acquires locks as it starts up, the memory barriers in `acquire` and `release` ensure that the child's CPU sees the parent's writes.

9.4 Parallelism

Locking is primarily about suppressing parallelism in the interests of correctness. Because performance is also important, kernel designers often have to think about how to use locks in a way that achieves both correctness and good parallelism. While xv6 is not systematically designed for high performance, it's still worth considering which xv6 operations can execute in parallel, and which might conflict on locks.

Pipes in xv6 are an example of fairly good parallelism. Each pipe has its own lock, so that different processes can read and write different pipes in parallel on different CPUs. For a given pipe, however, the writer and reader must wait for each other to release the lock; they can't read/write the same pipe at the same time. It is also the case that a read from an empty pipe (or a write to a full pipe) must block, but this is not due to the locking scheme.

Context switching is a more complex example. Two kernel threads, each executing on its own CPU, can call `yield`, `sched`, and `swtch` at the same time, and the calls will execute in parallel. The threads each hold a lock, but they are different locks, so they don't have to wait for each other. Once in `scheduler`, however, the two CPUs may conflict on locks while searching the table of processes for one that is `RUNNABLE`. That is, xv6 is likely to get a performance benefit from multiple CPUs during context switch, but perhaps not as much as it could.

Another example is concurrent calls to `fork` from different processes on different CPUs. The calls may have to wait for each other for `pid_lock` and `kmem.lock`, and for per-process locks needed to search the process table for an `UNUSED` process. On the other hand, the two forking processes can copy user memory pages and format page-table pages fully in parallel.

The locking scheme in each of the above examples sacrifices parallel performance in certain cases. In each case it's possible to obtain more parallelism using a more elaborate design. Whether it's worthwhile depends on details: how often the relevant operations are invoked, how long the

code spends with a contended lock held, how many CPUs might be running conflicting operations at the same time, whether other parts of the code are more restrictive bottlenecks. It can be difficult to guess whether a given locking scheme might cause performance problems, or whether a new design is significantly better, so measurement on realistic workloads is often required.

9.5 Exercises

1. Modify xv6's pipe implementation to allow a read and a write to the same pipe to proceed in parallel on different cores.
2. Modify xv6's `scheduler()` to reduce lock contention when different cores are looking for runnable processes at the same time.
3. Eliminate some of the serialization in xv6's `fork()`.

Chapter 10

Summary

This text introduced the main ideas in operating systems by studying one operating system, xv6, line by line. Some code lines embody the essence of the main ideas (e.g., context switching, user/kernel boundary, locks, etc.) and each line is important; other code lines provide an illustration of how to implement a particular operating system idea and could easily be done in different ways (e.g., a better algorithm for scheduling, better on-disk data structures to represent files, better logging to allow for concurrent transactions, etc.). All the ideas were illustrated in the context of one particular, very successful system call interface, the Unix interface, but those ideas carry over to the design of other operating systems.

Bibliography

- [1] The RISC-V instruction set manual: privileged architecture. <https://riscv.org/specifications/privileged-isa/>, 2019.
- [2] The RISC-V instruction set manual: user-level ISA. <https://riscv.org/specifications/isa-spec-pdf/>, 2019.
- [3] Hans-J Boehm. Threads cannot be implemented as a library. *ACM PLDI Conference*, 2005.
- [4] Edsger Dijkstra. Cooperating sequential processes. <https://www.cs.utexas.edu/users/EWD/transcriptions/EWD01xx/EWD123.html>, 1965.
- [5] Maurice Herlihy and Nir Shavit. *The Art of Multiprocessor Programming, Revised Reprint*. 2012.
- [6] Brian W. Kernighan. *The C Programming Language*. Prentice Hall Professional Technical Reference, 2nd edition, 1988.
- [7] Donald Knuth. *Fundamental Algorithms. The Art of Computer Programming. (Second ed.)*, volume 1. 1997.
- [8] L Lamport. A new solution of dijkstra’s concurrent programming problem. *Communications of the ACM*, 1974.
- [9] John Lions. *Commentary on UNIX 6th Edition*. Peer to Peer Communications, 2000.
- [10] Paul E. Mckenney, Silas Boyd-wickizer, and Jonathan Walpole. RCU usage in the linux kernel: One decade later, 2013.
- [11] Martin Michael and Daniel Durich. The NS16550A: UART design and application considerations. http://bitsavers.trailing-edge.com/components/national/_appNotes/AN-0491.pdf, 1987.
- [12] David Patterson and Andrew Waterman. *The RISC-V Reader: an open architecture Atlas*. Strawberry Canyon, 2017.
- [13] Dave Presotto, Rob Pike, Ken Thompson, and Howard Trickey. Plan 9, a distributed system. In *In Proceedings of the Spring 1991 EurOpen Conference*, pages 43–50, 1991.

- [14] Dennis M. Ritchie and Ken Thompson. The UNIX time-sharing system. *Commun. ACM*, 17(7):365–375, July 1974.

Index

., 92, 94
.., 92, 94
/init, 28, 37
_entry, 27

absorption, 86
acquire, 59, 62
address space, 25
argc, 37
argv, 37
atomic, 59

ballocc, 87, 89
batching, 85
bcache.head, 83
begin_op, 86
bfree, 87
bget, 83
binit, 83
bmap, 91
bottom half, 49
bread, 82, 84
brelse, 82, 84
BSIZE, 91
buf, 82
busy waiting, 72
bwrite, 82, 84, 86

chan, 72, 74
child process, 10
commit, 84
concurrency, 55
concurrency control, 55
condition lock, 73

conditional synchronization, 71
conflict, 58
contention, 58
contexts, 68
convoys, 78
copy-on-write (COW) fork, 46
copyinstr, 45
copyout, 37
coroutines, 69
CPU, 9
cpu->scheduler, 68, 69
crash recovery, 81
create, 94
critical section, 58
current directory, 17

deadlock, 60
direct blocks, 91
direct memory access (DMA), 52
dirlink, 92
dirlookup, 91, 92, 94
DIRSIZ, 91
disk, 83
driver, 49
dup, 93

ecall, 23, 26
ELF format, 37
ELF_MAGIC, 37
end_op, 86
exception, 41
exec, 12, 14, 27, 37, 44
exit, 11, 70, 76

file descriptor, 13

- filealloc, 93
- fileclose, 93
- filedup, 93
- fileread, 93, 96
- filestat, 93
- filewrite, 87, 93, 96
- fork, 10, 12, 14, 93
- forkret, 69
- freerange, 35
- fsck, 95
- fsinit, 86
- ftable, 93

- getcmd, 12
- group commit, 85
- guard page, 33

- hartid, 70

- I/O, 13
- I/O concurrency, 51
- I/O redirection, 14
- ialloc, 89, 94
- iget, 88, 89, 92
- ilock, 88, 89, 92
- indirect block, 91
- initcode.S, 27, 44
- initlog, 86
- inode, 18, 82, 87
- install_trans, 86
- interface design, 9
- interrupt, 41
- iput, 88, 89
- isolation, 21
- itrunc, 89, 91
- iunlock, 89

- kalloc, 35
- kernel, 9, 23
- kernel space, 9, 23
- kfree, 35
- kinit, 35
- kvminit, 33
- kvminithart, 34

- kvmmmap, 33

- lazy allocation, 47
- links, 18
- loadseg, 37
- lock, 55
- log, 84
- log_write, 86
- lost wake-up, 72

- machine mode, 23
- main, 33–35, 83
- malloc, 13
- mappages, 33
- memory barrier, 63
- memory model, 63
- memory-mapped, 31, 49
- metadata, 18
- microkernel, 24
- mkdev, 94
- mkdir, 94
- mkfs, 82
- monolithic kernel, 21, 23
- multi-core, 21
- multiplexing, 67
- multiprocessor, 21
- mutual exclusion, 57
- mycpu, 70
- myproc, 71

- namei, 37, 94
- nameiparent, 92, 94
- namex, 92
- NBUF, 83
- NDIRECT, 90, 91
- NINDIRECT, 91

- O_CREATE, 94
- open, 93, 94

- p->context, 70
- p->killed, 77, 101
- p->kstack, 26
- p->lock, 69, 70, 74

- p->pagetable, 26, 27
- p->state, 27
- p->xxx, 26
- page, 29
- page table entries (PTEs), 29
- page-fault exception, 30, 47
- paging from disk, 47
- parent process, 10
- path, 17
- persistence, 81
- PGROUNDUP, 35
- physical address, 25
- PHYSTOP, 33, 34
- PID, 10
- pipe, 15
- piperead, 75
- pipewrite, 75
- polling, 52, 72
- pop_off, 62
- printf, 12
- priority inversion, 78
- privileged instructions, 23
- proc_pagetable, 37
- process, 9, 24
- procinit, 34
- programmed I/O, 52
- PTE_R, 30
- PTE_U, 30
- PTE_V, 30
- PTE_W, 30
- PTE_X, 30
- push_off, 62

- race condition, 57
- read, 93
- readi, 37, 91
- recover_from_log, 86
- release, 59, 62
- root, 17
- round robin, 77
- RUNNABLE, 70, 74, 76

- satp, 30

- sbrk, 13
- scause, 42
- sched, 68, 69, 74
- scheduler, 69, 70
- semaphore, 71
- sepc, 42
- sequence coordination, 71
- serializing, 58
- sfence.vma, 34
- shell, 10
- signal, 78
- skipelem, 92
- sleep, 72–74
- sleep-locks, 63
- SLEEPING, 74
- sret, 27
- sscratch, 42
- sstatus, 42
- stat, 91, 93
- stati, 91, 93
- struct context, 68
- struct cpu, 70
- struct dinode, 87, 90
- struct dirent, 91
- struct elfhdr, 37
- struct file, 93
- struct inode, 88
- struct pipe, 75
- struct proc, 26
- struct run, 34
- struct spinlock, 58
- stval, 47
- stvec, 42
- superblock, 82
- supervisor mode, 23
- swtch, 68–70
- SYS_exec, 45
- sys_link, 94
- sys_mkdir, 94
- sys_mknod, 94
- sys_open, 94
- sys_pipe, 95

- sys_sleep, 62
- sys_unlink, 94
- syscall, 45
- system call, 9
- T_DEV, 91
- T_DIR, 91
- T_FILE, 94
- thread, 26
- thundering herd, 78
- ticks, 62
- tickslock, 62
- time-share, 10, 21
- top half, 49
- TRAMPOLINE, 43
- trampoline, 26, 43
- transaction, 81
- Translation Look-aside Buffer (TLB), 34
- transmit complete, 50
- trap, 41
- trapframe, 26
- type cast, 35
- UART, 49
- unlink, 85
- user memory, 25
- user mode, 23
- user space, 9, 23
- usertrap, 68
- ustack, 37
- uvmalloc, 37
- valid, 83
- virtio_disk_rw, 83, 84
- virtual address, 25
- wait, 11, 12, 70, 76
- wait channel, 72
- wakeup, 61, 72, 74
- walk, 33
- walkaddr, 37
- write, 85, 93
- write-through, 88
- writel, 87, 91
- yield, 68–70
- ZOMBIE, 76