

Concurrent systems

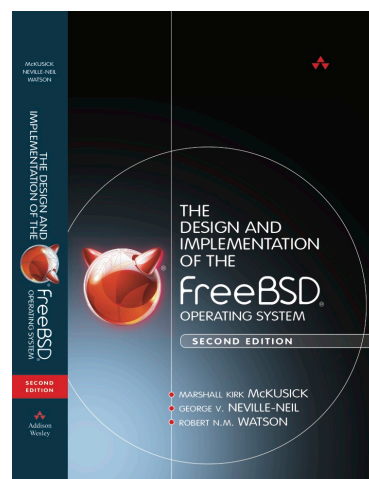
Case study: FreeBSD kernel concurrency

Dr Robert N. M. Watson

1

FreeBSD kernel

- Open-source OS kernel
 - **Large:** millions of LoC
 - **Complex:** thousands of subsystems, drivers, ...
 - **Very concurrent:** dozens or hundreds of CPU cores/threads
 - **Widely used:** NetApp, EMC, Dell, Apple, Juniper, Netflix, Sony, Cisco, Yahoo!, ...
- Why a case study?
 - Employs C&DS principles
 - Concurrency performance and composability at scale



In the library: Marshall Kirk McKusick, George V. Neville-Neil, and Robert N. M. Watson. *The Design and Implementation of the FreeBSD Operating System (2nd Edition)*, Pearson Education, 2014.

BSD + FreeBSD: a brief history

- 1980s Berkeley Standard Distribution (BSD)
 - ‘BSD’-style open-source license (MIT, ISC, CMU, ...)
 - UNIX Fast File System (UFS/FFS), sockets API, DNS, used TCP/IP stack, FTP, sendmail, BIND, cron, vi, ...
- Open-source FreeBSD operating system
 - 1993: FreeBSD 1.0 without support for multiprocessing
 - 1998: FreeBSD 3.0 with “giant-lock” multiprocessing
 - 2003: FreeBSD 5.0 with fine-grained locking
 - 2005: FreeBSD 6.0 with mature fine-grained locking
 - 2012: FreeBSD 9.0 with TCP scalability beyond 32 cores

3

FreeBSD: before multiprocessing (1)

- Concurrency model inherited from UNIX
- Userspace
 - **Preemptive multitasking between** processes
 - Later, **preemptive multithreading within** processes
- Kernel
 - ‘Just’ a C program running ‘bare metal’
 - Internally multithreaded
 - User threads ‘in kernel’ (e.g., in system calls)
 - Kernel services (e.g., async. work for VM, etc.)

4

FreeBSD: before multiprocessing (2)

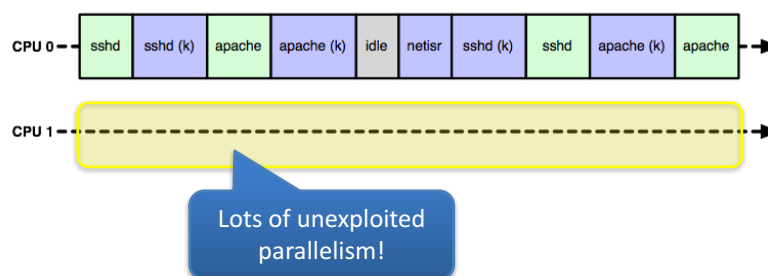
- **Cooperative multitasking** within kernel
 - Mutual exclusion as long as you don't `sleep()`
 - Implied global lock means local locks rarely required Except for interrupt handlers, non-preemptive kernel
 - **Critical sections** control interrupt-handler execution
- **Wait channels:** implied condition variable for every address


```
sleep(&x, ...);          // Wait for event on &x
wakeup(&x);             // Signal an event on &x
```

 - Must leave global state consistent when calling `sleep()`
 - Must reload any cached local state after `sleep()` returns
- Use to build higher-level synchronization primitives
 - E.g., `lockmgr()` reader-writer lock can be held over I/O (`sleep`)

5

Pre-multiprocessor scheduling



6

Hardware parallelism, synchronization

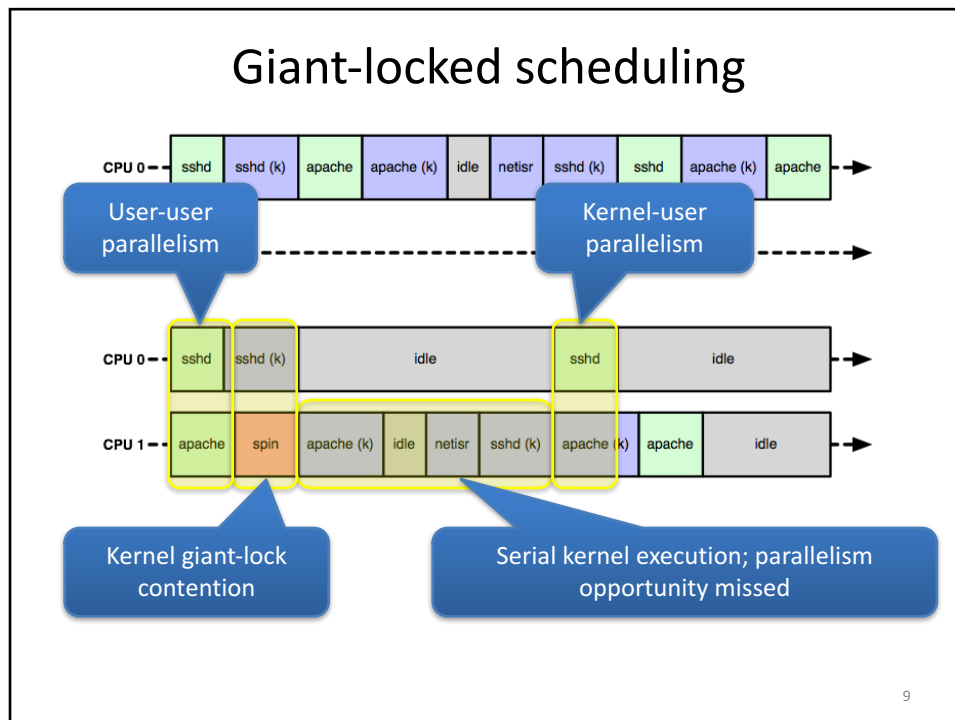
- Late 1990s: multi-CPU begins to move down market
 - In 2000s: 2-processor a big deal
 - In 2010s: 64-core is increasingly common
- **Coherent, symmetric, shared memory** systems
 - Instructions for **atomic memory access**
 - Compare-and-swap, test-and-set, load linked/store conditional
- Signaling via **Inter-Processor Interrupts** (IPIs)
 - CPUs can trigger an interrupt handler on each another
- Vendor extensions for performance, programmability
 - MIPS inter-thread message passing
 - Intel TM support: TSX (Whoops: HSW136!)

7

Giant locking the kernel

- FreeBSD follows footsteps of Cray, Sun, ...
- First, allow user programs to run in parallel
 - One instance of kernel code/data shared by all CPUs
 - Different user processes/threads on different CPUs
- **Giant spinlock** around kernel
 - Acquire on syscall/trap to kernel; drop on return
 - In effect: kernel runs on at most once CPU at a time; 'migrates' between CPUs on demand
- **Interrupts**
 - If interrupt delivered on CPU X while kernel is on CPU Y, forward interrupt to Y using an IPI

8



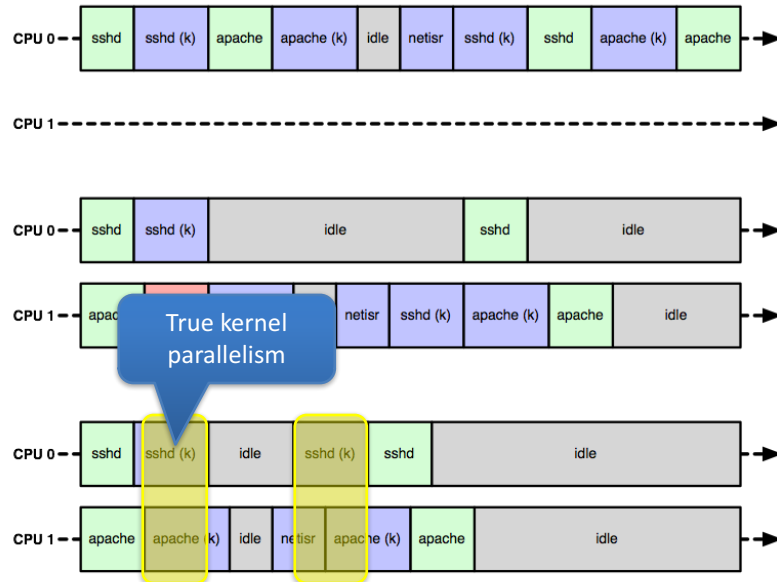
9

Fine-grained locking

- Giant locking is OK for user-program parallelism
- Kernel-centered workloads trigger **Giant contention**
 - Scheduler, IPC-intensive workloads
 - TCP/buffer cache on high-load web servers
 - Process-model contention with multithreading (VM, ...)
- Motivates migration to **fine-grained locking**
 - Greater granularity (may) afford greater parallelism
- Mutexes/condition variables rather than semaphores
 - Increasing consensus on pthreads-like synchronization
 - Explicit locks are easier to debug than semaphores
 - Support for **priority inheritance** + **priority propagation**
 - E.g., Linux is also now migrating away from semaphores

10

Fine-grained scheduling



11

Kernel synchronization primitives

- **Spin locks**
 - Used to implement the scheduler, interrupt handlers
- **Mutexes, reader-writer, read-mostly locks**
 - Most heavily used – different optimization tradeoffs
 - Can be held only over “bounded” computations
 - **Adaptive**: on contention, sleep is expensive; spin first
 - Sleeping depends on scheduler, and hence on spinlocks...
- **Shared-eXclusive (SX) locks, condition variables**
 - Can be held over I/O and other unbounded waits
- **Condition variables** usable with any lock type
- Most primitives support **priority propagation**

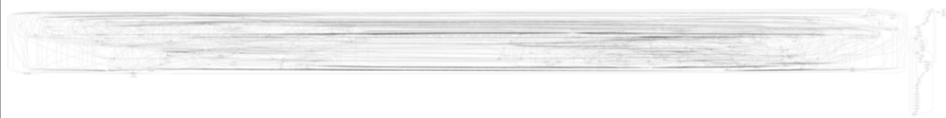
12

WITNESS lock-order checker

- Kernel relies on **partial lock order** to prevent deadlock (Recall dining philosophers)
 - In-field lock-related deadlocks are (very) rare
- WITNESS is a **lock-order debugging tool**
 - Warns when lock cycles (could) arise by tracking edges
 - Only in debugging kernels due to overhead (15%+)
- Tracks both statically declared, dynamic lock orders
 - **Static orders** most commonly intra-module
 - **Dynamic orders** most commonly inter-module
- Deadlocks for condition variables remain hard to debug
 - What thread should have woken up a CV being waited on?
 - Similar to semaphore problem

13

WITNESS: global lock-order graph*



* Turns out that the global lock-order graph is pretty complicated.

14



*

* Commentary on WITNESS full-system lock-order graph complexity; courtesy Scott Long, Netflix

15

Excerpt from global lock-order graph*

This bit mostly has to do with networking

Local clusters: e.g., related locks from the firewall: two leaf nodes; one is held over calls to other subsystems

Network interface locks: "transmit" occurs at the bottom of call stacks via many layers holding locks

Memory allocator locks follow most other locks, since most kernel components require memory allocation

* The local lock-order graph is **also** complicated.

WITNESS debug output

```
1st 0xffffffff80025207f0 run0_node_lock (run0_node_lock) @
/usr/src/sys/net80211/ieee80211_ioctl.c:1341
2nd 0xffffffff80025142a8 run0 (network driver) @
/usr/src/sys/modules/usb/run/../../dev/usb/wlan/if_run.c:3368
```

KDB: stack backtrace:

```
db_trace_self_wrapper() at db_trace_self_wrapper+0x2a
kdb_backtrace() at kdb_backtrace+0x37
_witness_debugger() at _witness_debugger+0x2c
witness_checkorder() at witness_checkorder+0x853
_mtx_lock_flags() at _mtx_lock_flags+0x85
run_raw_xmit() at run_raw_xmit+0x58
ieee80211_send_mgmt() at ieee80211_send_mgmt+0x4d5
domlme() at domlme+0x95
setmlme_common() at setmlme_common+0x2f0
ieee80211_ioctl_setmlme() at ieee80211_ioctl_setmlme+0x7e
ieee80211_ioctl_set80211() at ieee80211_ioctl_set80211+0x46f
in_control() at in_control+0xad
ifioctl() at ifioctl+0xece
kern_ioctl() at kern_ioctl+0xcd
sys_ioctl() at sys_ioctl+0xf0
amd64_syscall() at amd64_syscall+0x380
Xfast_syscall() at Xfast_syscall+0xf7
--- syscall (54, FreeBSD ELF64, sys_ioctl), rip = 0x
0x7ffffffffffd848, rbp = 0x2a ---
```

Lock names and source code locations of acquisitions adding the offending graph edge

Stack trace to acquisition that triggered cycle:
802.11 called USB;
previously, perhaps USB called 802.11?

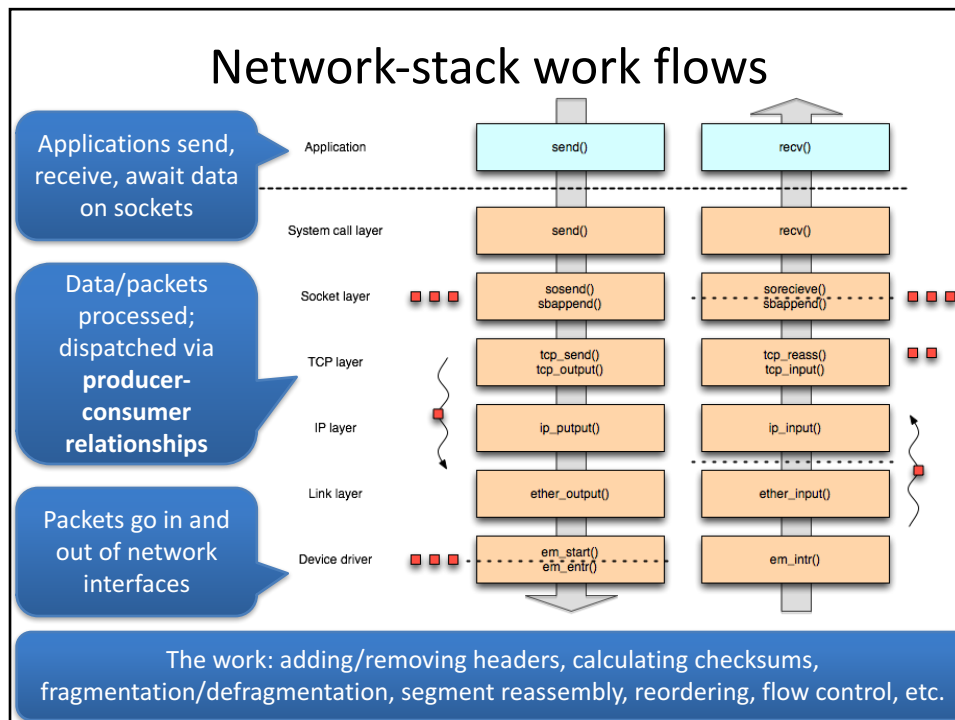
How does this work in practice?

- Kernel is heavily multi-threaded
- Each user thread has a corresponding kernel thread
 - Represents user thread when in syscall, page fault, etc.
- Kernels services often execute in asynchronous threads
 - Interrupts, timers, I/O, networking, etc.
- Therefore extensive synchronization
 - Locking model is almost always data-oriented
 - Think ‘monitors’ rather than ‘critical sections’
 - Reference counting or reader-writer locks used for stability
 - Higher-level patterns (producer-consumer, active objects, etc.) used frequently

Kernel threads in action

robot@lemongrass-freebsd64:~\$ procstat -at

PID	TID	COMM	TNAME	CPU	PRI	STATE	WCHAN
0	100000	kernel	swapper	1	84	sleep	sched
0	100009	kernel	firmware taskq	0	108	sleep	-
0	100014	kernel	kqueue taskq	0	108	sleep	-
0	100016	kernel	thread taskq	0	108	sleep	-
0	100020	kernel	acpi_task	0	32	sleep	-
0	100021	kernel	acpi_task	0	32	sleep	-
0	100022	kernel	acpi_task	0	32	sleep	-
0	100023	kernel	ffs_trim taskq	0	32	sleep	-
0	100033	kernel	em0 taskq	0	84	sleep	-
12	100037	int					
12	100038	int					
13	100010	ged					
13	100011	ged					
13	100012	ged					
13	100013	ged					
13	100014	ged					
13	100015	ged					
13	100016	ged					
13	100017	ged					
13	100018	ged					
13	100019	ged					
13	100020	ged					
13	100021	ged					
13	100022	ged					
13	100023	ged					
13	100024	ged					
13	100025	ged					
13	100026	ged					
13	100027	ged					
13	100028	ged					
13	100029	ged					
13	100030	ged					
13	100031	ged					
13	100032	ged					
13	100033	ged					
13	100034	ged					
13	100035	ged					
13	100036	ged					
13	100037	ged					
13	100038	ged					
13	100039	ged					
13	100040	ged					
13	100041	ged					
13	100042	ged					
13	100043	ged					
13	100044	ged					
13	100045	ged					
13	100046	ged					
13	100047	ged					
13	100048	ged					
13	100049	ged					
13	100050	ged					
13	100051	ged					
13	100052	ged					
13	100053	ged					
13	100054	ged					
13	100055	ged					
13	100056	ged					
13	100057	ged					
13	100058	ged					
13	100059	ged					
13	100060	ged					
13	100061	ged					
13	100062	ged					
13	100063	ged					
13	100064	ged					
13	100065	ged					
13	100066	ged					
13	100067	ged					
13	100068	ged					
13	100069	ged					
13	100070	ged					
13	100071	ged					
13	100072	ged					
13	100073	ged					
13	100074	ged					
13	100075	ged					
13	100076	ged					
13	100077	ged					
13	100078	ged					
13	100079	ged					
13	100080	ged					
13	100081	ged					
13	100082	ged					
13	100083	ged					
13	100084	ged					
13	100085	ged					
13	100086	ged					
13	100087	ged					
13	100088	ged					
13	100089	ged					
13	100090	ged					
13	100091	ged					
13	100092	ged					
13	100093	ged					
13	100094	ged					
13	100095	ged					
13	100096	ged					
13	100097	ged					
13	100098	ged					
13	100099	ged					
13	100100	ged					
13	100101	ged					
13	100102	ged					
13	100103	ged					
13	100104	ged					
13	100105	ged					
13	100106	ged					
13	100107	ged					
13	100108	ged					
13	100109	ged					
13	100110	ged					
13	100111	ged					
13	100112	ged					
13	100113	ged					
13	100114	ged					
13	100115	ged					
13	100116	ged					
13	100117	ged					
13	100118	ged					
13	100119	ged					
13	100120	ged					
13	100121	ged					
13	100122	ged					
13	100123	ged					
13	100124	ged					
13	100125	ged					
13	100126	ged					
13	100127	ged					
13	100128	ged					
13	100129	ged					
13	100130	ged					
13	100131	ged					
13	100132	ged					
13	100133	ged					
13	100134	ged					
13	100135	ged					
13	100136	ged					
13	100137	ged					
13	100138	ged					
13	100139	ged					
13	100140	ged					
13	100141	ged					
13	100142	ged					
13	100143	ged					
13	100144	ged					
13	100145	ged					
13	100146	ged					
13	100147	ged					
13	100148	ged					
13	100149	ged					
13	100150	ged					
13	100151	ged					
13	100152	ged					
13	100153	ged					
13	100154	ged					
13	100155	ged					
13	100156	ged					
13	100157	ged					
13	100158	ged					
13	100159	ged					
13	100160	ged					
13	100161	ged					
13	100162	ged					
13	100163	ged					
13	100164	ged					
13	100165	ged					
13	100166	ged					
13	100167	ged					
13	100168	ged					
13	100169	ged					
13	100170	ged					
13	100171	ged					
13	100172	ged					
13	100173	ged					
13	100174	ged					
13	100175	ged					
13	100176	ged					
13	100177	ged					
13	100178	ged					
13	100179	ged					
13	100180	ged					
13	100181	ged					
13	100182	ged					
13	100183	ged					
13	100184	ged					
13	100185	ged					
13	100186	ged					
13	100187	ged					
13	100188	ged					
13	100189	ged					
13	100190	ged					
13	100191	ged					
13	100192	ged					
13	100193	ged					
13	100194	ged					
13	100195	ged					
13	100196	ged					
13	100197	ged					
13	100198	ged					
13	100199	ged					
13	100200	ged					
13	100201	ged					
13	100202	ged					
13	100203	ged					



Case study: the network stack (2)

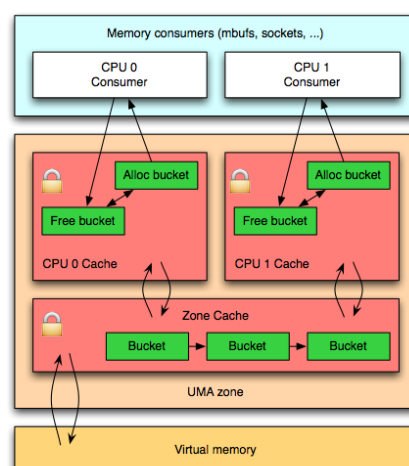
- First, make it **safe** without the Giant lock
 - Lots of data structures require locks
 - Condition signaling already exists but will be added to
 - Establish key work flows, lock orders
- Then, make it **fast**
 - Especially locking primitives themselves
 - Increase locking granularity where there is contention
- As hardware becomes more parallel, identify and exploit further concurrency opportunities
 - Add more threads, distribute more work

What to lock and how?

- Fine-grained locking **overhead** vs. **contention**
 - Some contention is inherent: reflects necessary communication
 - Some contention is **false sharing**: side effect of structure choices
- Principle: lock data, not code (i.e., not critical sections)
 - Key structures: network interfaces, sockets, work queues
 - Independent structure instances often have their own locks
- Horizontal vs. vertical parallelism
 - H: Different locks for different connections (e.g., TCP1 vs. TCP2)
 - H: Different locks within a layer (e.g., receive vs. send buffers)
 - V: Different locks at different layers (e.g., socket vs. TCP state)
- Things not to lock: packets in flight - mbufs ('work')

23

Example: Universal Memory Allocator (UMA)



- Key kernel service
- Slab allocator
 - (Bonwick 1994)
- Per-CPU caches
 - Individually locked
 - Amortise (or avoid) global lock contention
- Some allocation patterns only per-CPU caches
- Others require dipping into the global pool

24

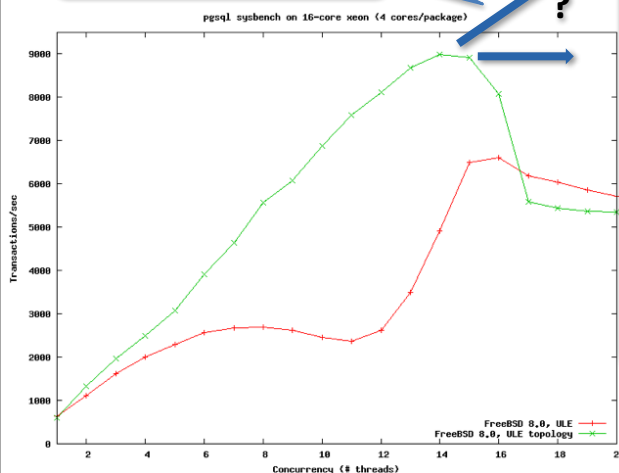
Work distribution

- Packets (mbufs) are units of work
- Parallel work requires distribution to threads
 - Must keep packets ordered – or TCP gets cranky!
- Implication: strong per-flow serialization
 - I.e., no generalized producer-consumer/round robin
 - Various strategies to keep work ordered; e.g.:
 - Process in a single thread
 - Multiple threads in a 'pipeline' linked by a queue
 - Misordering allowed between flows, just not within them
- Establish flow-CPU affinity can both order processing and utilize caches well

25

Scalability

What might we expect if we didn't hit contention?



Key idea:
speedup

As we add more parallelism, we would like the system to get faster.

Key idea:
performance collapse

Sometimes parallelism hurts performance more than it helps due to work-distribution overheads, contention.

26

Longer-term strategies

- Hardware change motivates continuing work
 - Optimize inevitable contention
 - Lockless primitives
 - rmlocks, read-copy-update (RCU)
 - Per-CPU data structures
 - Distribute work to more threads .. to utilise growing core count
- Optimise for locality, not just contention: cache, NUMA, and I/O affinity

27

Conclusions

- FreeBSD employs many of C&DS techniques
 - Mutual exclusion, condition synchronization
 - Producer-consumer, lockless primitives
 - Also Write-Ahead Logging (WAL) in filesystems
- Real-world systems are really complicated
 - Hopefully, you will mostly consume, rather than produce, concurrency primitives like these
 - Composition is not straightforward
 - Parallelism performance wins are a lot of work
 - Hardware continues to evolve
- See you in Distributed Systems!

28