

5/6/24

User Level Threads & Storage Device

→ managed & scheduled by the user libraries / runtime

→ why use them?

→ cheaper to create & schedule

→ custom scheduling policies

→ app specific decisions (important lock holders stay running, deadline driven policies)

→ cooperative scheduling

→ schedule only when threads voluntarily yield

→ how does a user level thread run?

→ on top of a kernel level thread (switches btwn user threads)

→ $N = 1$ Model

SSSS ^N user threads

⏟
}

1
kernel thread

(tasks) $N = M$ model (thread pool)
(workers)

SSSS SSSSS ^N user threads

⏟
} } }

$M = 3$
kernel thread

(proportional
to # of cores)

→ What happens if a user thread blocks?

→ how might it block?

→ synchronization (e.g. sleeplock) ^{user} ↗ transfers control to user scheduler & schedules a different user thread.

→ blocking syscalls, exceptions that generate I/O

→ the underlying kernel thread blocks, no user threads can run on the kernel thread in the meantime.

★ How can we mitigate this?

→ use nonblocking syscalls when possible

→ keep some back up kernel threads sleeping

wake one up to take over the rest of user threads when a user thread is about to make a blocking syscall.

→ class discussion: delay the blocking syscall user thread & run other user threads until close to the end of our time slice

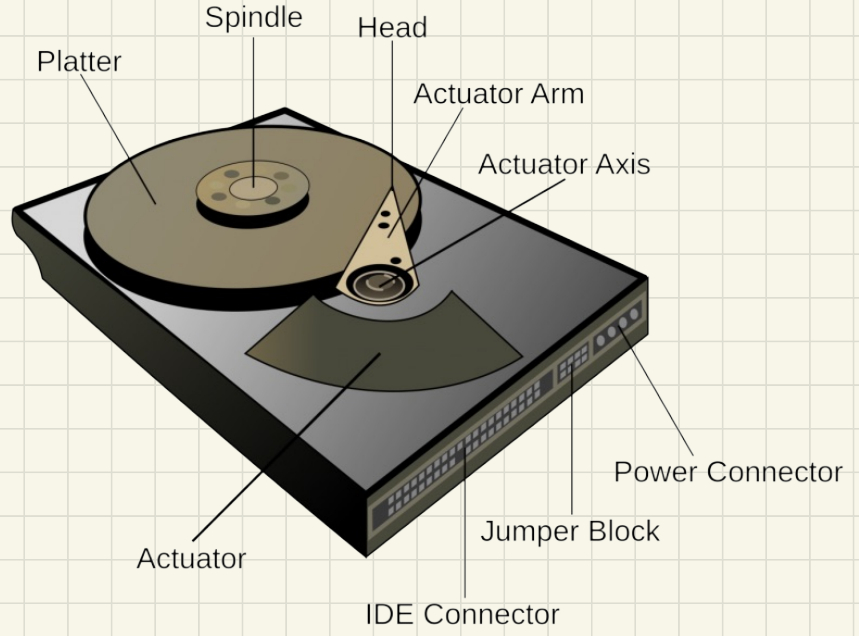
→ challenges: how do we know how much time slice is left? what if it changes (MUTQ)? what if syscall is on the crucial path of execution?

Golang
→ goroutines
go { ... }

Storage Devices

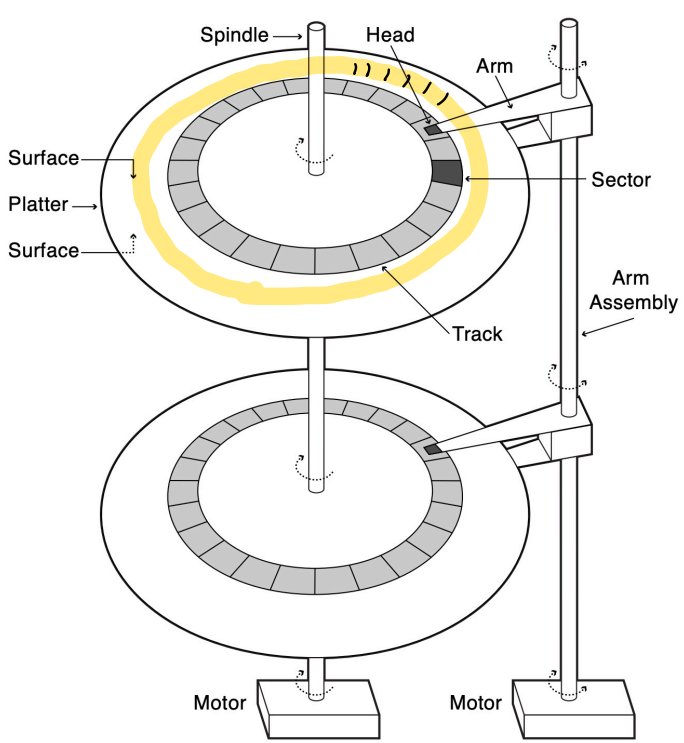
- persistent (non volatile), large capacity, block addressable
- hard drive / spinning disk (HDD)
 - cheap per GB (€ 10-20 per TB)
 - physical movement needed (slow access latency, 10-20 ms)
- solid state drive (SSD)
 - more expensive than HDD, cheaper than DRAM
 - no physical movement (faster access = 10 μ s - 100 μ s)

Hard Drive



Sectors

→ 512 bytes, unit of read & write



Disk Request

→ host (CPU) sends a request to the disk controller

→ disk moves arm to the specific track

* seek time: 1-20ms (depending on how far to move)
average 10ms

→ wait for the sector to spin under the disk head

* rotational time: depends on RPM
4-15ms
assume it takes half a rotation to reach the sector (avg)

→ transfer data back to host (for a read req).

* transfer time: depends on # of bytes transferred & disk bandwidth (80-160 MB/s)

Request Latency

→ total time = seek time + rotational time + transfer time
(ms)

→ cost of reading / writing 1 sector (512 bytes)

→ given 10ms seek time, 7200 RPM, 120 M.B/s

$$\rightarrow \frac{512 \text{ B}}{120 \text{ M.B/s}} \times 1024^2 \text{ B/M.B} \times 1000 \text{ ms/s} = 0.004 \text{ ms}$$

→ total latency = 10ms + 4ms

= 120 RPS

$$\hookrightarrow \frac{1}{0.12 \text{ RPS}} = 8.3 \text{ ms per rotation}$$

$$+ 0.004 \text{ ms} = 14.004 \text{ ms!}$$

→ 4 ms per half rotation (average)

→ cost of read/write 10 consecutive sectors

$$10 \text{ ms} + 4 \text{ ms} + 0.04 \text{ ms} = 14.04 \text{ ms}$$

1 seek 1 avg rotation 512 bytes

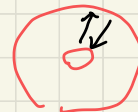
→ cost of read/write 10 random sectors

$$14.004 \text{ ms} \times 10 = 140.04 \text{ ms}$$

Single request 10 times (different tracks, sectors)

★ Can we improve random access performance?

→ disk head scheduling!



seek requests w/ less seek time.

→ CSKAN

Metrics for evaluating disk performance : IOPS

$$\text{I/O Operations Per Second} = \frac{\text{\# of I/O requests}}{\text{total latency}}$$

→ IOPS for 10 consecutive read requests

$$= \frac{10}{14.24 \text{ ms}} \times 1000 \text{ ms/s} = 712 \text{ IOPS}$$

→ IOPS for 10 random requests

$$= \frac{10}{140.04 \text{ ms}} \times 1000 \text{ ms/s} = 71 \text{ IOPS}$$