

# Disk scheduling

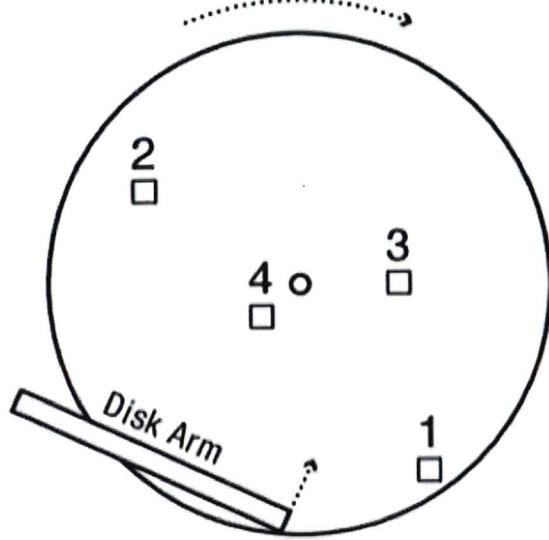
SCAN: - like an elevator

- first sweep in one direction, then other

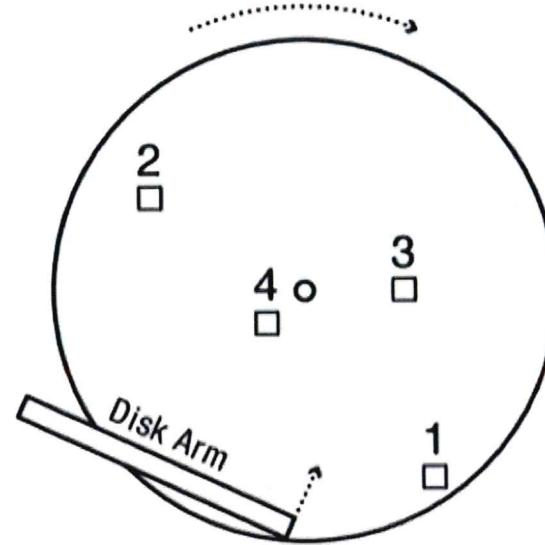
CSCAN: - always sweeps in same direction

R-CSCAN: - allow small backward seeks < rotational delay

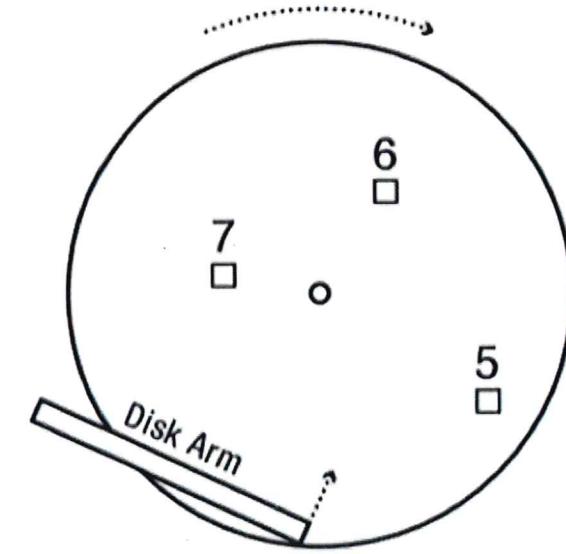
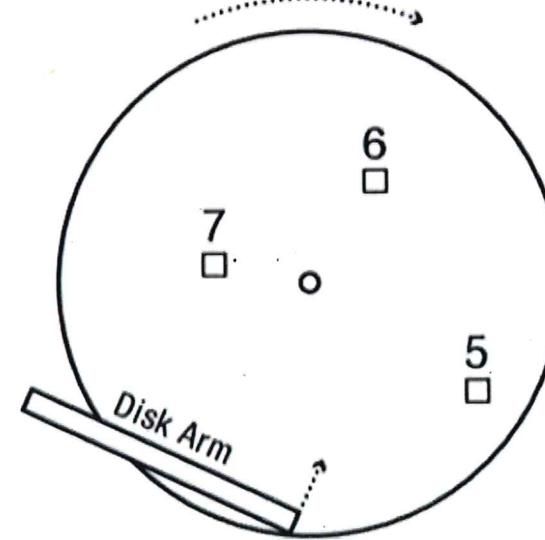
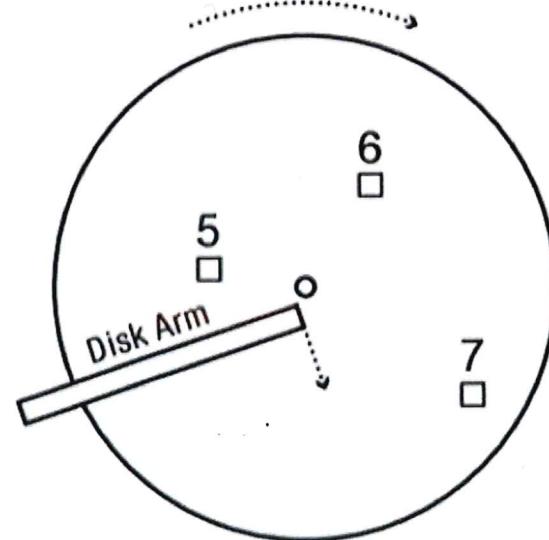
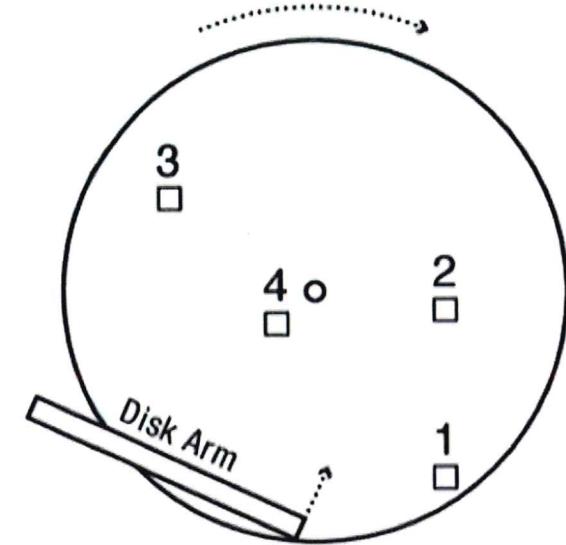
SCAN



CSCAN



R-CSCAN



18.9

# Example

avg seek = 8ms  
avg rotation = 4ms  
xfer = 250MB/s

How long to complete 100 random 4kb reads?

Can do this in one sweep!

All seeks will be minimal  $\approx$  1ms

rotation = 4ms avg

xfer:  $4\text{KB} / 250 \text{ MB/s} = 16 \text{ ms}$

$$\text{total} = 100 \cdot (1\text{ms} + 4\text{ms} + 0.01\text{ms}) = 0.5\text{s}$$

more than twice as fast as FIFO

---

How long to read whole disk?

$$32\text{TB} / 250 \text{ MB/s} = 32 \cdot 4000 = 128,000 \text{ s} = 35.5 \text{ hrs}$$

Chance 1 bit is corrupt? bit error rate =  $10^{-15}$

$$32\text{TB} = 256\text{Tbit} = 2.56 \times 10^{14} \text{ bit} \rightarrow \approx 23\% \text{ chance of } \geq 1 \text{ error}$$

18.10

Lec 19

# File Systems

# File System Abstraction

- persistent, named data
- hierarchical (directories/folders that nest)
- access control
- File: named collection of data
  - linear sequence of bytes
- crash and storage error tolerant
- good performance on typical block devices

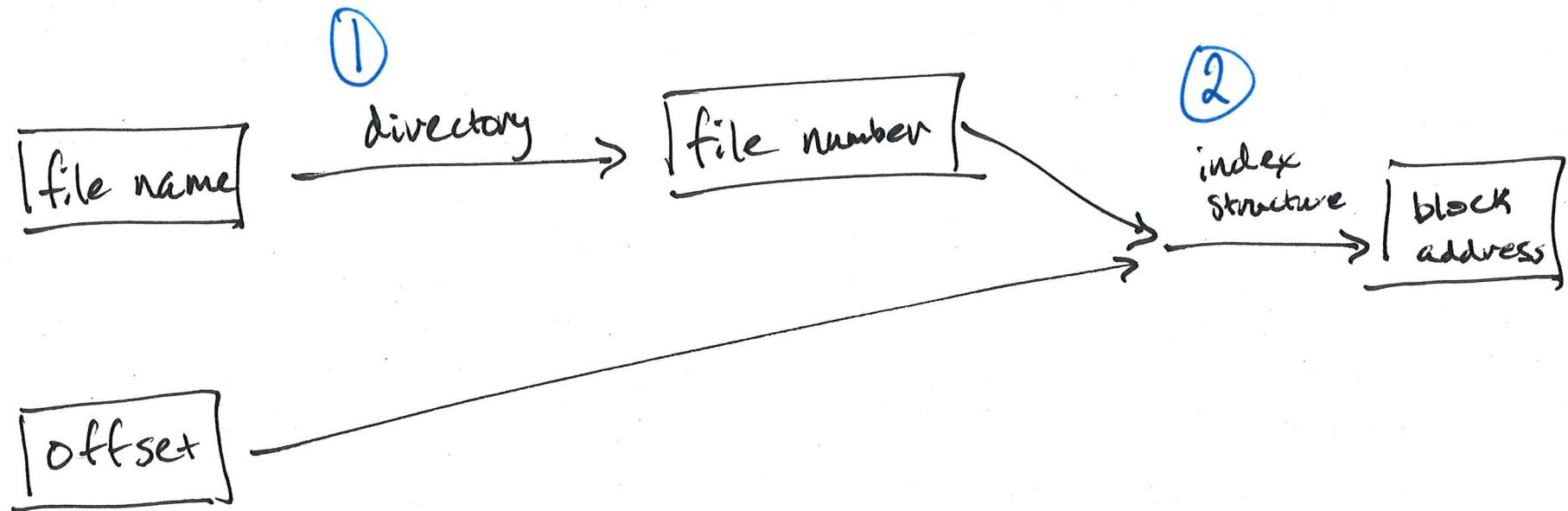
# More Abstractions

- Directory
  - group of named files + subdirectories
  - maps name to metadata location
- Path : string that names a file / directory
  - eg /cse/web/courses/cse451/25sp/index.html
- Links:
  - hard link : two filenames point to same file
  - soft link : one filename points to another file name
- Mount: Map a path in one FS to root of another FS

# Unix File System API

- create, link, unlink, createdir, rmdir
- open, close, read, write, seek
- fsync
  - Write does not actually send to disk - cached
  - fsync writes back the cache to disk
- **Warning:** most programmers use fsync incorrectly

# Implementing Names



`read("/foo", offset = 100)`

index structure maps  
"virtual" offset to "physical" offset

# File System Workload

file size: are most files large or small?

which accounts for more total storage?

file access: are most accesses to large or small files?

which accounts for more total I/O?

are most accesses sequential or random?

are file sizes known at creation time?

# File System design constraints

- for small files:
  - small blocks (low overhead)
  - locality across files
- for large files:
  - large blocks (low overhead)
  - locality within file
  - efficient random access
- unknown at creation time:
  - file size, file lifetime, access pattern

# File System Design Options

	FAT	FFS	NTFS
index structure	linked list	fixed shape asymmetric tree	dynamic tree
granularity	block	block	extent
free space management	FAT array	bitmap (fixed location)	bitmap (file)
locality	defragmentation	block groups + reserve space	- extents - best bit - defragmentation