



Main points

- File systems
 - Useful abstractions on top of physical devices
- Storage hardware characteristics
 - Disks and flash memory
- File system usage patterns

File systems

- Abstraction on top of persistent storage
 - Magnetic disk
 - Flash memory (e.g., USB thumb drive)
- Devices provide
 - Storage that (usually) survives across machine crashes
 - Block level (random) access
 - Large capacity at low cost
 - Relatively slow performance
 - * Magnetic disk read takes 10-20M processor instructions

File system as illusionist: hide limitations of physical storage

- Persistence of data stored in file system:
 - Even if crash happens during an update
 - Even if disk block becomes corrupted
 - Even if flash memory wears out
- Naming:
 - Named data instead of disk block numbers
 - Directories instead of flat storage
 - Byte addressable data even though devices are block-oriented
- Performance:
 - Cached data
 - Data placement and data structure organization
- Controlled access to shared data

File system abstraction

- File system
 - Persistent, named data
 - Hierarchical organization (directories, subdirectories)
 - Access control on data
- File: named collection of data
 - Linear sequence of bytes (or a set of sequences)
 - Read/write or memory mapped
- Crash and storage error tolerance
 - Operating system crashes (and disk errors) leave file system in a valid state
- Performance
 - Achieve close to the hardware limit in the average case

Storage devices

- Magnetic disks
 - Storage that rarely becomes corrupted
 - Large capacity at low cost
 - Block level random access
 - Slow performance for random access
 - Better performance for streaming access
- Flash memory
 - Storage that rarely becomes corrupted
 - Capacity at intermediate cost (50x disk)
 - Block level random access
 - Good performance for reads; worse for random writes



Caching inside a disk controller

Method

- Disk controller has DRAM to cache recently accessed blocks * Hitachi disk has 16MB
 - * Some of the RAM space stores "firmware" (an embedded OS)
- Blocks are replaced usually in an LRU order
- Pros
 - Good for reads if accesses have locality
- Cons
 - Expensive
 - Need to deal with reliable writes







Sectors

- Sectors contain sophisticated error correcting codes
 - Disk head magnet has a field wider than track
 - Hide corruptions due to neighboring track writes
- Sector sparing
 - Remap bad sectors transparently to spare sectors on the same surface
- Slip sparing
 - Remap all sectors (when there is a bad sector) to preserve sequential behavior
- Track skewing
 - Sector numbers offset from one track to the next, to allow for disk head movement for sequential ops



Disk cylinder and arm

- CD's and floppies come individually, but magnetic disks come organized in a disk pack
- Cylinder
 Certain track of the platter
- Disk arm
 A disk arm carries disk heads
- Read/write operation
 - Disk controller receives a command with <track#, sector#>
 - Seek the right cylinder (tracks)
 - Wait until the right sector comes
 - Perform read/write



<section-header><code-block><code-block><text><text><text><text></code></code>

hiba disk (2008)	
Size	
Platters/Heads	2/4
Capacity	320 GB
Performance	
Spindle speed	7200 RPM
Average seek time read/write	10.5 ms/ 12.0 ms
Maximum seek time	19 ms
Track-to-track seek time	1 ms
Transfer rate (surface to buffer)	54–128 MB/s
Transfer rate (buffer to host)	375 MB/s
Buffer memory	16 MB
Power	
Typical	16.35 W
Idle	11.68 W

How long to complete 500 random disk reads, in FIFO order?

Question How long to complete 500 random disk reads, in FIFO order? Seek: average 10.5 msec Rotation: average 4.15 msec Transfer: 5-10 usec 500 * (10.5 + 4.15 + 0.01)/1000 = 7.3 seconds

Question

How long to complete 500 sequential disk reads?





 How large a transfer is needed to achieve 80% of the max disk transfer rate?

Assume x rotations are needed, then solve for x: 0.8 (10.5 ms + (1ms + 8.5ms) x) = 8.5ms x

Total: x = 9.1 rotations, 9.8MB

Disk scheduling

- ♦ FIFO
 - Schedule disk operations in order they arrive
 - Downsides?

















- How long to complete 500 random disk reads, in any order?
 - Disk seek: 1ms (most will be short)
 - Rotation: 4.15ms
 - Transfer: 5-10usec
- ◆ Total: 500 * (1 + 4.15 + 0.01) = 2.2 seconds
 - Would be a bit shorter with R-CSCAN
 - vs. 7.3 seconds if FIFO order

Question

How long to read all of the bytes off of a disk?





Flash memory

- Writes must be to "clean" cells; no update in place
 - Large block erasure required before write
 - Erasure block: 128 512 KB
 - Erasure time: Several milliseconds
- Write/read page (2-4KB)
 - 50-100 usec

-lash arive (2011)		
Size		
Capacity	300 GB	
Page Size	4KB	
Performance		
Bandwidth (Sequential Reads)	270 MB/s	
Bandwidth (Sequential Writes)	210 MB/s	
Read/Write Latency	75 μs	
Random Reads Per Second	38,500	
Random Writes Per Second	2,000 (2,400 with 20% space reserve)	
Interface	SATA 3 Gb/s	
Endurance		
Endurance	1.1 PB (1.5 PB with 20% space reserve)	
Power		
Power Consumption Active/Idle	3.7 W / 0.7 W	

- Why are random writes so slow?
 - Random write: 2000/sec
 - Random read: 38500/sec

Flash translation layer

- Flash device firmware maps logical page # to a physical location
 - Garbage collect erasure block by copying live pages to new location, then erase
 - * More efficient if blocks stored at same time are deleted at same time (e.g., keep blocks of a file together)
 - Wear-levelling: only write each physical page a limited number of times
 - Remap pages that no longer work (sector sparing)
- Transparent to the device user

File system - flash

- How does Flash device know which blocks are live?
 - Live blocks must be remapped to a new location during erasure
- TRIM command
 - File system tells device when blocks are no longer in use

File system workload

- File sizes
 - Are most files small or large?
 - Which accounts for more total storage: small or large files?

File system workload

- File sizes
 - Are most files small or large?
 * SMALL
 - Which accounts for more total storage: small or large files?
 * LARGE

File system workload

- File access
 - Are most accesses to small or large files?
 - Which accounts for more total I/O bytes: small or large files?

File system workload

- File access
 - Are most accesses to small or large files?
 * SMALL
 - Which accounts for more total I/O bytes: small or large files? * LARGE

File system workload

- + How are files used?
 - Most files are read/written sequentially
 - Some files are read/written randomly * Ex: database files, swap files
 - Some files have a pre-defined size at creation
 - Some files start small and grow over time
 - * Ex: program stdout, system logs

File system design

- For small files:
 - Small blocks for storage efficiency
 - Concurrent ops more efficient than sequential
 - Files used together should be stored together
- For large files:
 - Storage efficient (large blocks)
 - Contiguous allocation for sequential access
 - Efficient lookup for random access
- May not know at file creation
 - Whether file will become small or large
 - Whether file is persistent or temporary
 - Whether file will be used sequentially or randomly

File system abstraction

- Directory
 - Group of named files or subdirectories
 - Mapping from file name to file metadata location
- Path
 - String that uniquely identifies file or directory
 - Ex: /cse/www/education/courses/cse451/12au
- Links
 - Hard link: link from name to metadata location
 - Soft link: link from name to alternate name
- Mount
 - Mapping from name in one file system to root of another

UNIX file system API

- create, link, unlink, createdir, rmdir
 - Create file, link to file, remove link
 - Create directory, remove directory
- open, close, read, write, seek
 - Open/close a file for reading/writing
 - Seek resets current position
- fsync
 - File modifications can be cached
 - fsync forces modifications to disk (like a memory barrier)



- UNIX file open is a Swiss Army knife:
 - Open the file, return file descriptor
 - Options:
 - * if file doesn't exist, return an error
 - * If file doesn't exist, create file and open it
 - * If file does exist, return an error
 - * If file does exist, open file
 - * If file exists but isn't empty, nix it then open
 - * If file exists but isn't empty, return an error
 - * ...



Why not separate syscalls for open/create/exists?
 Would be more modular!

```
if (!exists(name))
```

- create(name); // can create fail?
- fd = open(name); // does the file exist?