

15-213

*"The course that gives CMU its Zip!"*

## Disk-based Storage Oct. 23, 2008

### Topics

- Storage technologies and trends
- Locality of reference
- Caching in the memory hierarchy

lecture-17.ppt

## Announcements

### Exam next Thursday

- style like exam #1: in class, open book/notes, no electronics

2

15-213, F'08

## Disk-based storage in computers

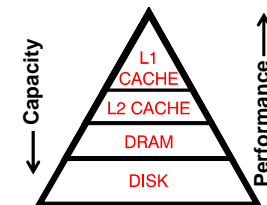
- **Memory/storage hierarchy**
  - Combining many technologies to balance costs/benefits
  - Recall the memory hierarchy and virtual memory lectures

3

15-213, F'08

## Memory/storage hierarchies

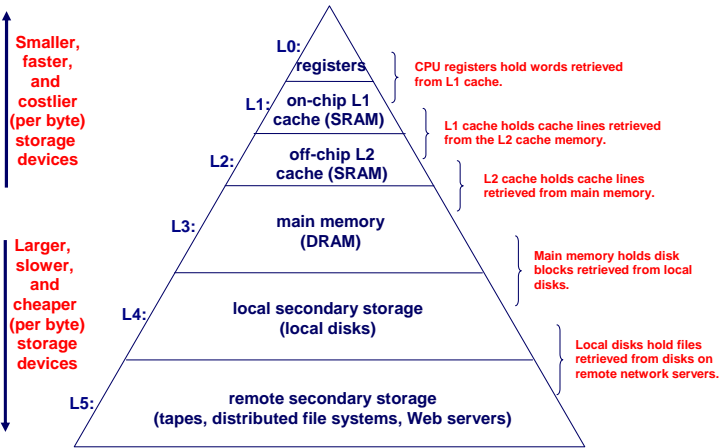
- **Balancing performance with cost**
  - Small memories are **fast but expensive**
  - Large memories are **slow but cheap**
- **Exploit locality to get the best of both worlds**
  - locality = re-use/nearness of accesses
  - allows most accesses to use small, fast memory



4

15-213, F'08

## An Example Memory Hierarchy



5

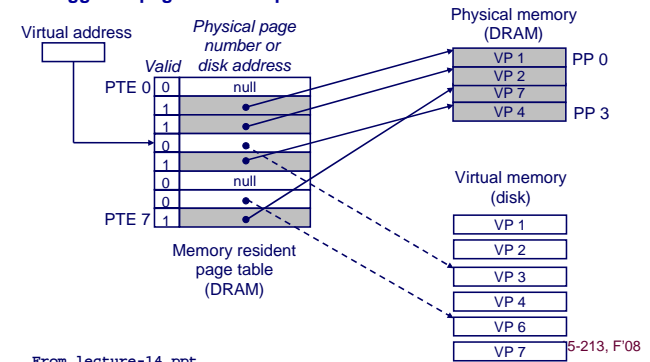
From lecture-9.ppt

15-213, F'08

## Page Faults

A page fault is caused by a reference to a VM word that is not in physical (main) memory

- Example: An instruction references a word contained in VP 3, a miss that triggers a page fault exception



6

From lecture-14.ppt

15-213, F'08

## Disk-based storage in computers

- Memory/storage hierarchy
  - Combining many technologies to balance costs/benefits
  - Recall the memory hierarchy and virtual memory lectures
- Persistence
  - Storing data for lengthy periods of time
    - DRAM/SRAM is "volatile": contents lost if power lost
    - Disks are "non-volatile": contents survive power outages
  - To be useful, it must also be possible to find it again later
    - this brings in many interesting data organization, consistency, and management issues
      - take 18-746/15-746 Storage Systems ©
    - we'll talk a bit about file systems next

7

15-213, F'08

## What's Inside A Disk Drive?

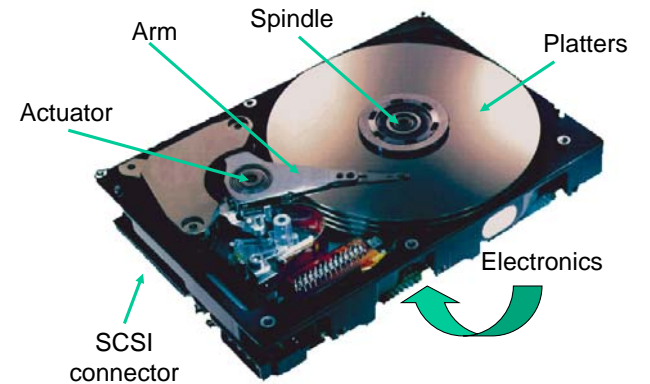


Image courtesy of Seagate Technology

8

15-213, F'08

## Disk Electronics

Quantum Viking (circa 1997)



6 Chips

- R/W Channel
- uProcessor  
32-bit, 25 MHz
- Power Array
- 2 MB DRAM
- Control ASIC  
SCSI, servo, ECC
- Motor/Spindle

Just like a small computer – processor, memory, network iface

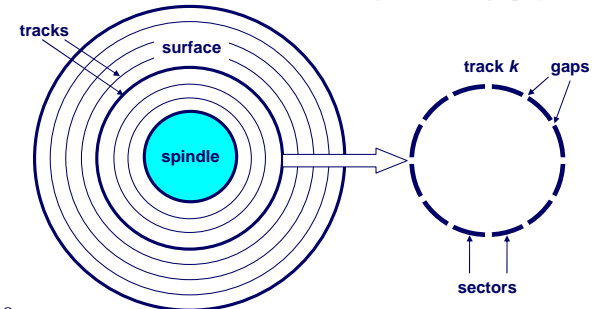
- Connect to disk
- Control processor
- Cache memory
- Control ASIC
- Connect to motor

9

15-213, F'08

## Disk “Geometry”

Disks contain **platters**, each with two **surfaces**  
 Each surface organized in concentric rings called **tracks**  
 Each track consists of **sectors** separated by **gaps**

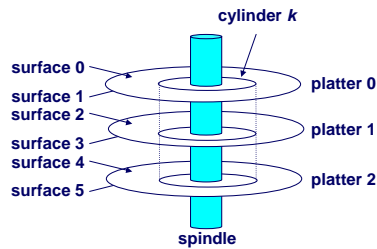


10

15-213, F'08

## Disk Geometry (Multiple-Platter View)

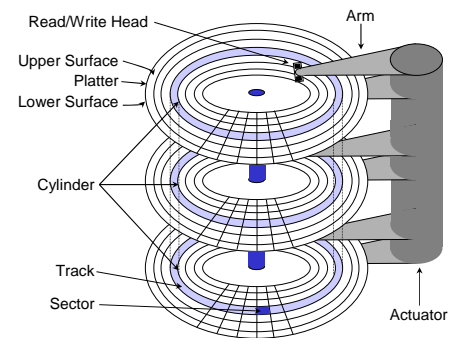
Aligned tracks form a cylinder



11

15-213, F'08

## Disk Structure

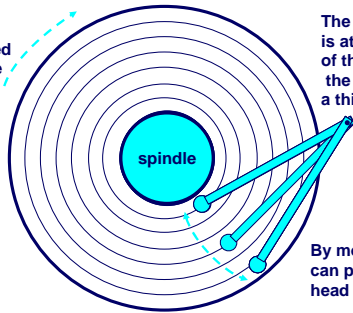


12

15-213, F'08

## Disk Operation (Single-Platter View)

The disk surface spins at a fixed rotational rate



The read/write head is attached to the end of the arm and flies over the disk surface on a thin cushion of air

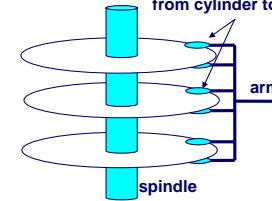
By moving radially, the arm can position the read/write head over any track

13

15-213, F'08

## Disk Operation (Multi-Platter View)

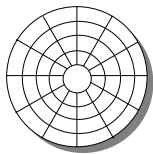
read/write heads move in unison from cylinder to cylinder



14

15-213, F'08

## Disk Structure - top view of single platter



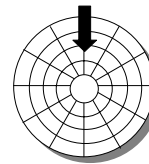
Surface organized into tracks

Tracks divided into sectors

15

15-213, F'08

## Disk Access

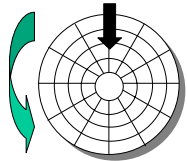


Head in position above a track

16

15-213, F'08

## Disk Access

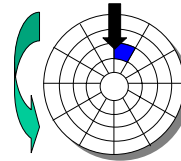


Rotation is counter-clockwise

17

15-213, F'08

## Disk Access – Read

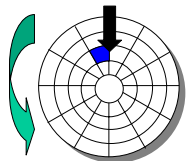


About to read blue sector

18

15-213, F'08

## Disk Access – Read



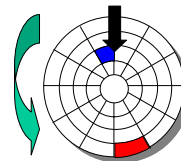
After BLUE read

After reading blue sector

19

15-213, F'08

## Disk Access – Read



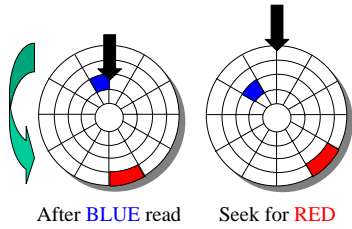
After BLUE read

Red request scheduled next

20

15-213, F'08

## Disk Access – Seek



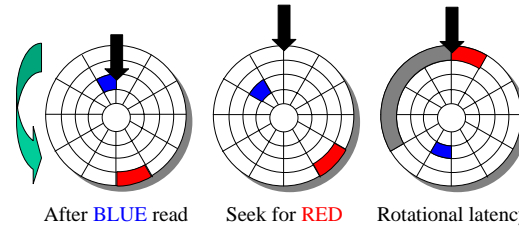
After BLUE read    Seek for RED

Seek to red's track

21

15-213, F'08

## Disk Access – Rotational Latency



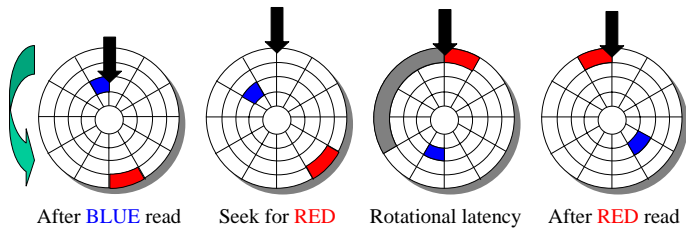
After BLUE read    Seek for RED    Rotational latency

Wait for red sector to rotate around

22

15-213, F'08

## Disk Access – Read



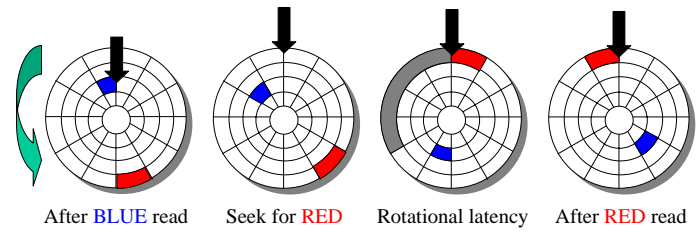
After BLUE read    Seek for RED    Rotational latency    After RED read

Complete read of red

23

15-213, F'08

## Disk Access – Service Time Components



After BLUE read    Seek for RED    Rotational latency    After RED read

Seek  
Rotational Latency  
Data Transfer

24

15-213, F'08

## Disk Access Time

Average time to access a specific sector approximated by:

- Taccess = Tavg seek + Tavg rotation + Tavg transfer

### Seek time (Tavg seek)

- Time to position heads over cylinder containing target sector
- Typical Tavg seek = 3-5 ms

### Rotational latency (Tavg rotation)

- Time waiting for first bit of target sector to pass under r/w head
- Tavg rotation =  $1/2 \times 1/\text{RPMs} \times 60 \text{ sec}/1 \text{ min}$ 
  - e.g., 3ms for 10,000 RPM disk

### Transfer time (Tavg transfer)

- Time to read the bits in the target sector
- Tavg transfer =  $1/\text{RPM} \times 1/(\text{avg \# sectors/track}) \times 60 \text{ secs}/1 \text{ min}$ 
  - e.g., 0.006ms for 10,000 RPM disk with 1,000 sectors/track

25

given 512-byte sectors, ~85 MB/s data transfer rate 15-213, F'08

## Disk Access Time Example

Given:

- Rotational rate = 7,200 RPM
- Average seek time = 5 ms
- Avg # sectors/track = 1000

Derived average time to access random sector:

- Tavg rotation =  $1/2 \times (60 \text{ secs}/7200 \text{ RPM}) \times 1000 \text{ ms}/\text{sec} = 4 \text{ ms}$
- Tavg transfer =  $60/7200 \text{ RPM} \times 1/400 \text{ secs/track} \times 1000 \text{ ms}/\text{sec} = 0.008 \text{ ms}$
- Taccess = 5 ms + 4 ms + 0.008 ms = 9.008 ms
  - Time to second sector: 0.008 ms

Important points:

- Access time dominated by seek time and rotational latency
- First bit in a sector is the most expensive, the rest are free
- SRAM access time is about 4 ns/doubleword, DRAM about 60 ns
  - ~100,000 times longer to access a word on disk than in DRAM

26

15-213, F'08

## Disk storage as array of blocks



OS's view of storage device  
(as exposed by SCSI or IDE/ATA protocols)

- Common "logical block" size: 512 bytes
- Number of blocks: device capacity / block size
- Common OS-to-storage requests defined by few fields
  - R/W, block #, # of blocks, memory source/dest

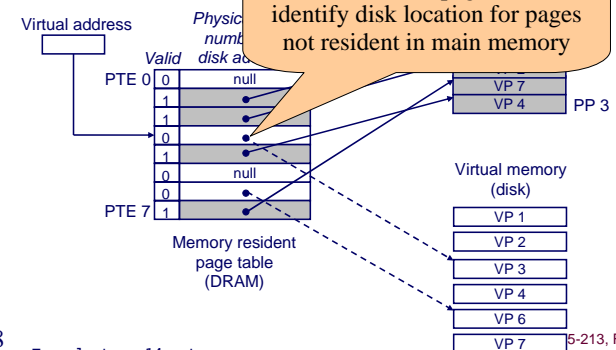
27

15-213, F'08

## Page Faults

A page fault is caused by a reference to a VM word that is not in physical (main) memory

- Example: An instruction reference that triggers a page fault

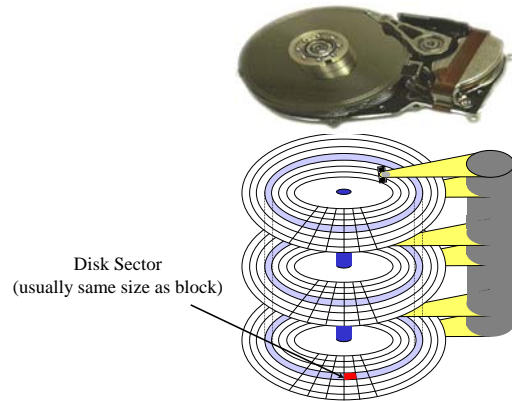


28

From lecture-14.ppt

15-213, F'08

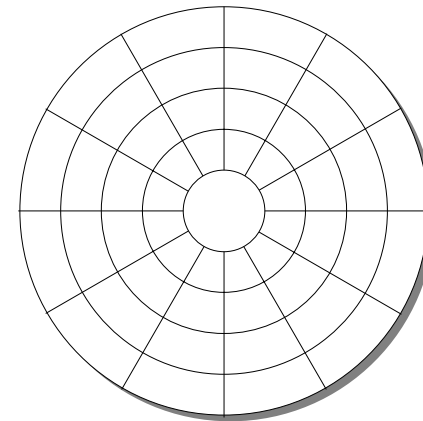
In device, "blocks" mapped to physical store



29

15-213, F'08

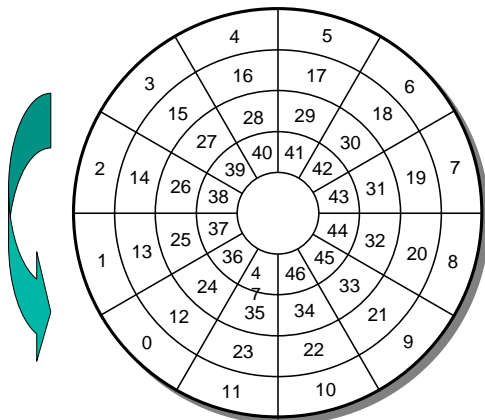
## Physical sectors of a single-surface disk



30

15-213, F'08

## LBN-to-physical for a single-surface disk



31

15-213, F'08

## Disk Capacity

**Capacity:** maximum number of bits that can be stored

- Vendors express capacity in units of gigabytes (GB), where 1 GB =  $10^9$  Bytes (Lawsuit pending! Claims deceptive advertising)

**Capacity is determined by these technology factors:**

- **Recording density** (bits/in): number of bits that can be squeezed into a 1 inch linear segment of a track
- **Track density** (tracks/in): number of tracks that can be squeezed into a 1 inch radial segment
- **Areal density** (bits/in<sup>2</sup>): product of recording and track density

32

15-213, F'08



## Computing Disk Capacity

Capacity = (# bytes/sector) x (avg. # sectors/track) x  
 (# tracks/surface) x (# surfaces/platter) x  
 (# platters/disk)

Example:

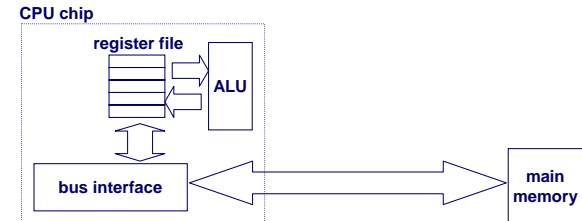
- 512 bytes/sector
- 1000 sectors/track (on average)
- 20,000 tracks/surface
- 2 surfaces/platter
- 5 platters/disk

Capacity = 512 x 1000 x 80000 x 2 x 5  
 = 409,600,000,000  
 = 409.6 GB

33

15-213, F'08

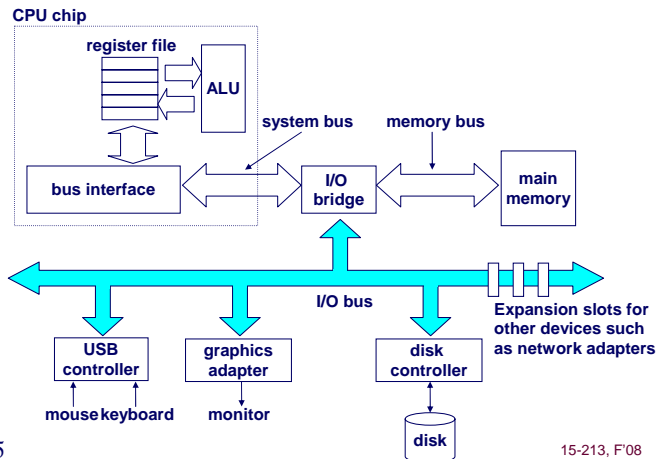
## Looking back at the hardware



34

15-213, F'08

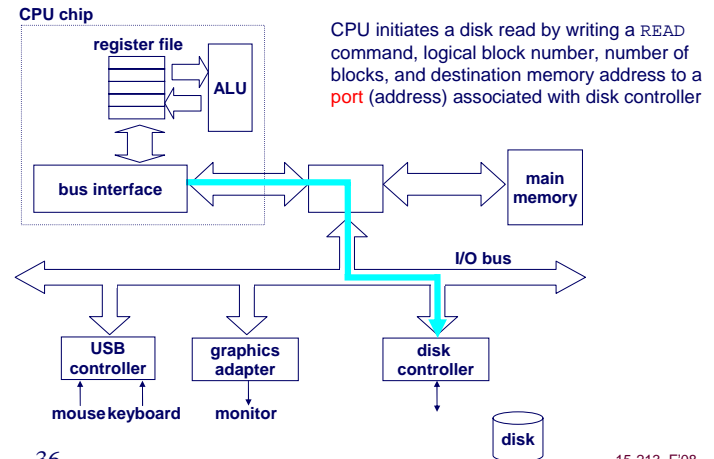
## Connecting I/O devices: the I/O Bus



35

15-213, F'08

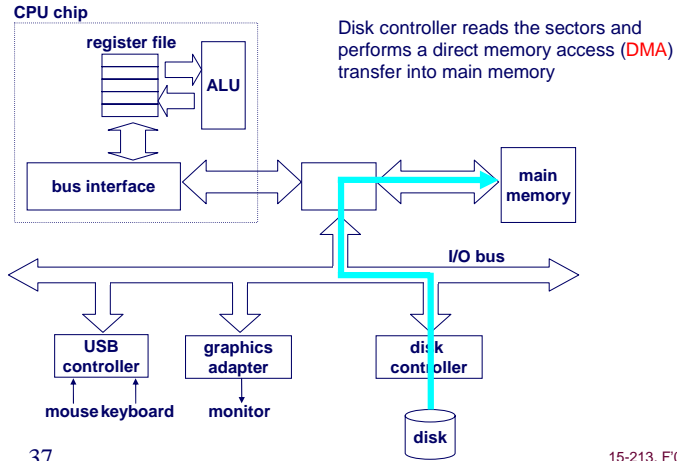
## Reading from disk (1)



36

15-213, F'08

## Reading from disk (2)



## Reading from disk (3)

