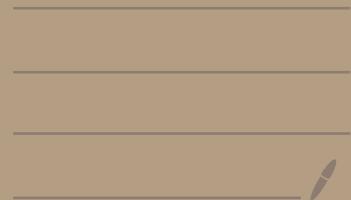


## —The Memory Hierarchy

---



## Random-Access Memory (RAM)

### (I) Key Features

- (i) RAM is traditionally packaged as a chip.
- (ii) Basic storage unit is normally a cell (one bit per cell)
- (iii) Multiple RAM chips form a memory

### (II) RAM comes in two varieties.

- (i) SRAM (static RAM)
- (ii) DRAM (Dynamic RAM)

|      | Trans.<br>per bit | Access<br>time | Needs<br>refresh? | Needs<br>EDC? | Cost | Applications                  |
|------|-------------------|----------------|-------------------|---------------|------|-------------------------------|
| SRAM | 4 or 6            | 1X             | NO                | Maybe         | 100X | Cache memory                  |
| DRAM | 1                 | 10X            | YES               | Yes           | 1X   | Main memories<br>frame buffer |

## Nonvolatile Memories

(I) DRAM and SRAM are volatile memory

(i) Lose information if powered off.

(II) Nonvolatile memories retain value even if powered off.

(i). Read-only memory (ROM): programmed during production.

(ii) Programmable ROM (PROM): can be programmed once.

(iii) Erasable PROM (EPROM): can be bulk erased (UV, X-ray)

(iv). Electrically erasable PROM (EEPROM): electronic erase capability.

(v). Flash memory: EEPROMs with partial block-level erase capability.

• Wears out after about 100,000 erasing.

(III) Uses for Nonvolatile Memories.

(i) Firmware programs stored in a ROM (BIOS, controllers for disks, network cards, graphics accelerators, security subsystems, ...)

(ii) Solid State disks (replace rotating disks in thumb drives, smart phones, mp3 players, tablets, laptops, ...)

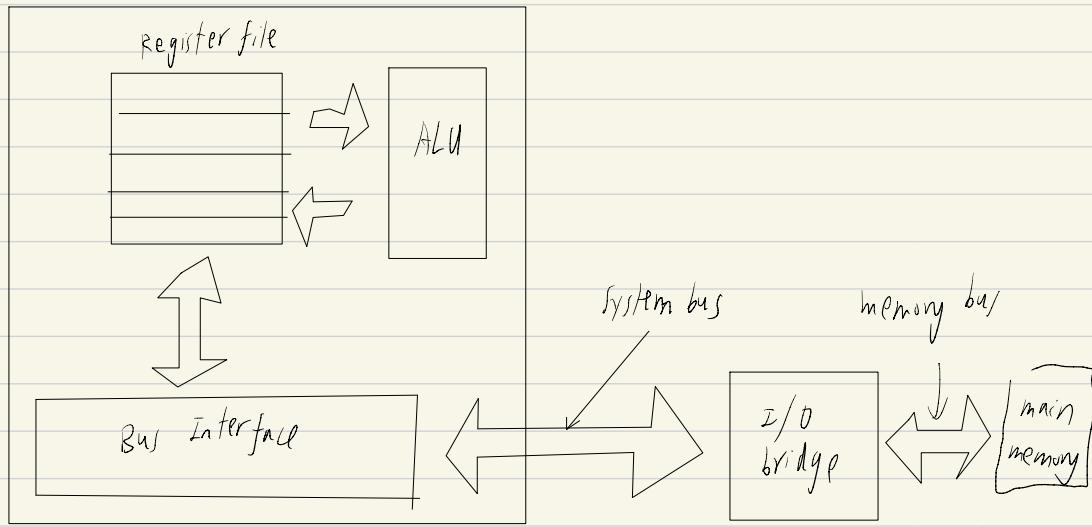
(iii) Disk caches

Traditional Bus Structure connecting  
CPU and memory

A **bus** is a collection of parallel wires that carry address,  
data and control signals.

Buses are typically shared by multiple devices.

CPU chip.



Main memory reads  $A$  from the memory bus, retrieves word  $X$  and places it on the bus.

CPU reads word  $X$  from the bus and copies it into register  $\%rax$ .

Platter rotating

slower than DRAM  
SRAM

## Disk Geometry

- (i) Disks consist of platters, each with two surfaces.
- | (II) Each surface consists of concentric rings called tracks.
- | (III) Each track consists of sectors separated by gaps.

## Disk capacity:

(i) capacity: maximum number of bits that can be stored.

Vendors express capacity in units of gigabytes (GB), where

$$1 \text{ GB} = 10^9 \text{ Bytes}$$

- | (IV) capacity is determined by these technology factors:
- | (i) **Recording density** (bits/in): number of bits that can be squeezed into a 1 inch segment of track.

| (ii) **Track density** (tracks/in): number of tracks that can be squeezed into a 1 inch radial segment.

| (iii) **Areal density** (bits/in<sup>2</sup>): product of recording and track density

(II) Modern disks partition tracks . . . . into disjoint subsets called  
recording zones

## Disk Access Time

(I) Average time to access some target sector approximated by:

$$(i) T_{\text{access}} = T_{\text{avg seek}} + T_{\text{avg transfer}}$$

(II) Seek Time ( $T_{\text{avg seek}}$ )

Time to position heads over cylinder containing target sector.

Typical avg seek is 5-9 ms.

(III) Rotational Latency ( $T_{\text{avg rotation}}$ )

(read/write)

(i) Time waiting for first bit of target sector to pass under r/w head.

$$(ii) T_{\text{avg rotation}} = 1/2 \times 1/\text{RPMs} \times 60\text{sec/min.}$$

$$(iii) \text{Typical } T_{\text{avg rotation}} = 7200 \text{ RPMs.}$$

(IV) Transfer time ( $T_{\text{avg transfer}}$ ) ignorable compared to previous two.

(i) Time to read the bits into the target sector.

$$(ii) T_{\text{avg transfer}} = 1/\text{RPM} \times 1/(\text{avg #sectors/track}) \times 60 \text{ sec/min.}$$

## Logical Disk Blocks:

(i) Modern disks present a simpler abstract view of the complex sector geometry:

(ii) The set of available sectors is modelled by a sequence of b-sized logical blocks (0, 1, 2...)

(iii) Mapping between logical blocks and actual physical sectors.

(iv) Maintained by hardware/firmware device called disk controller.

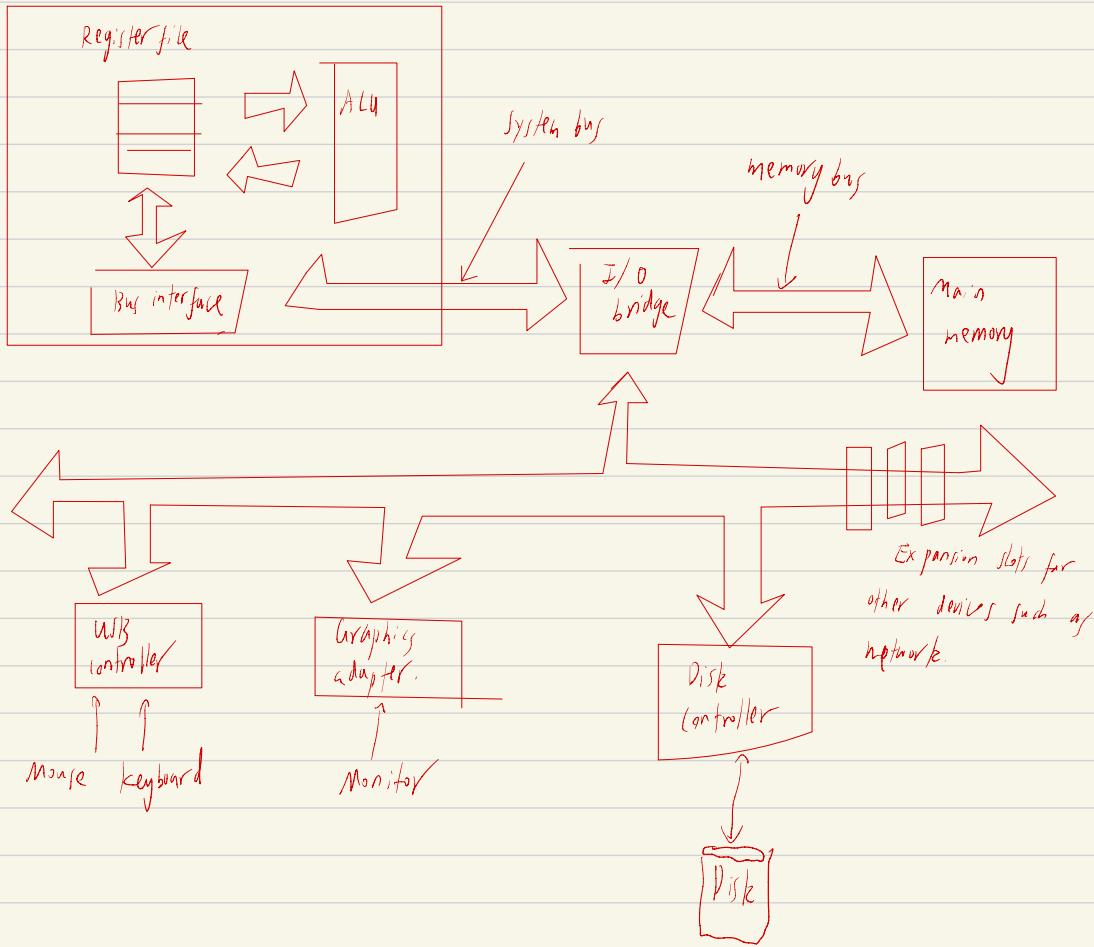
(v) Converts requests for logical blocks into (surface, track, sector) triplets

(vi) Allows controller to set aside spare cylinders for each zone.

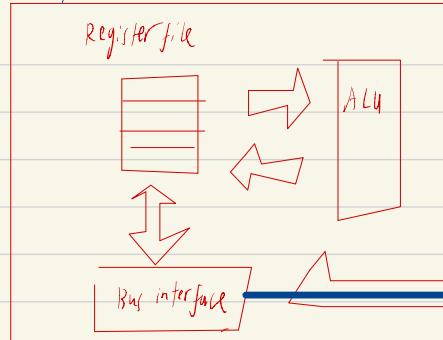
Allows for the difference in "formatted capacity" and "maximum capacity".

# I/O Bus

CPU chip

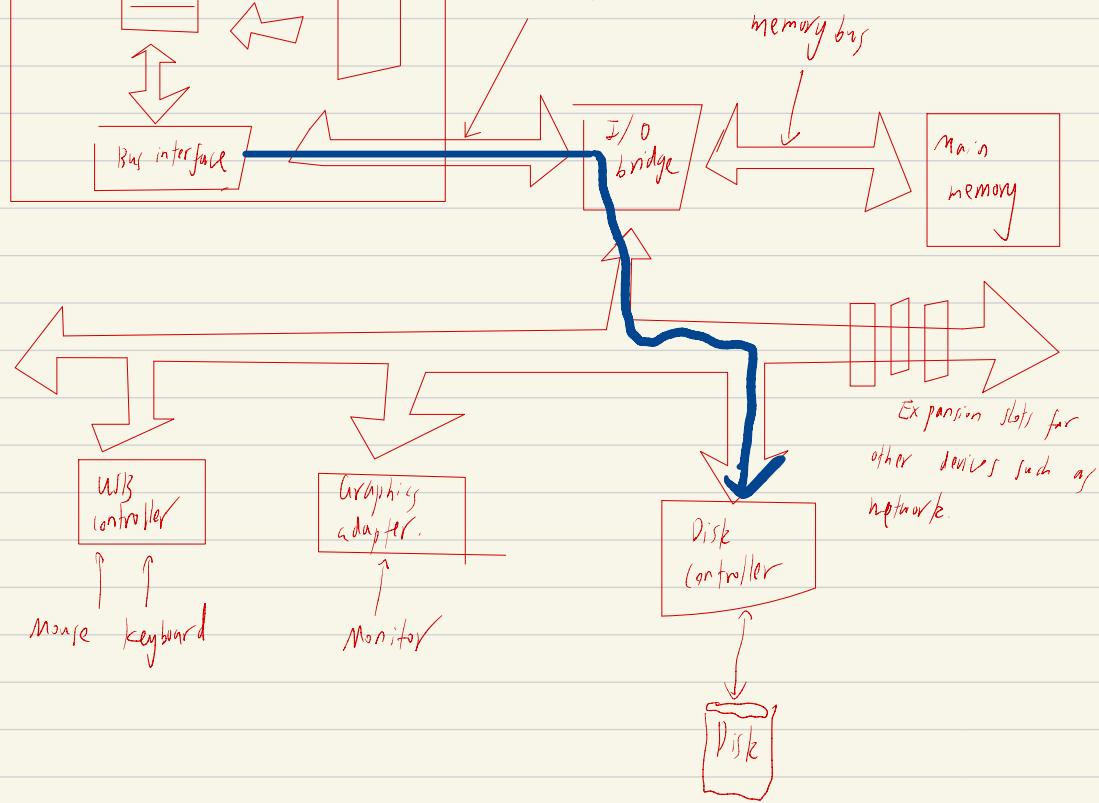


CPU chip. Reading a Disk Sector(s)



(CPU initiates a disk read by writing a command, logical block number, and destination memory address to a port (address) associated with disk controller.

System bus

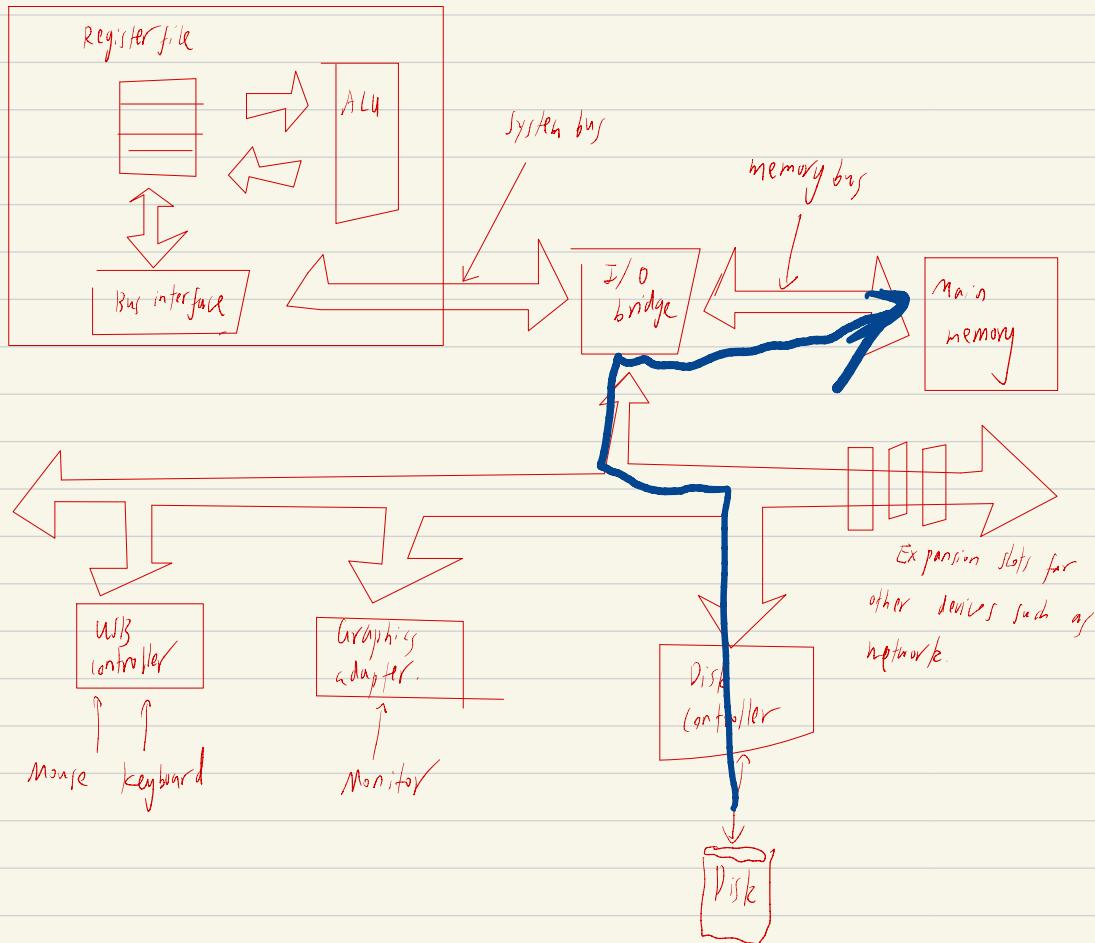


Reading a Disk Sector (4)

Disk controller reads the sector and performs a direct memory access

(DMA) transfer into memory

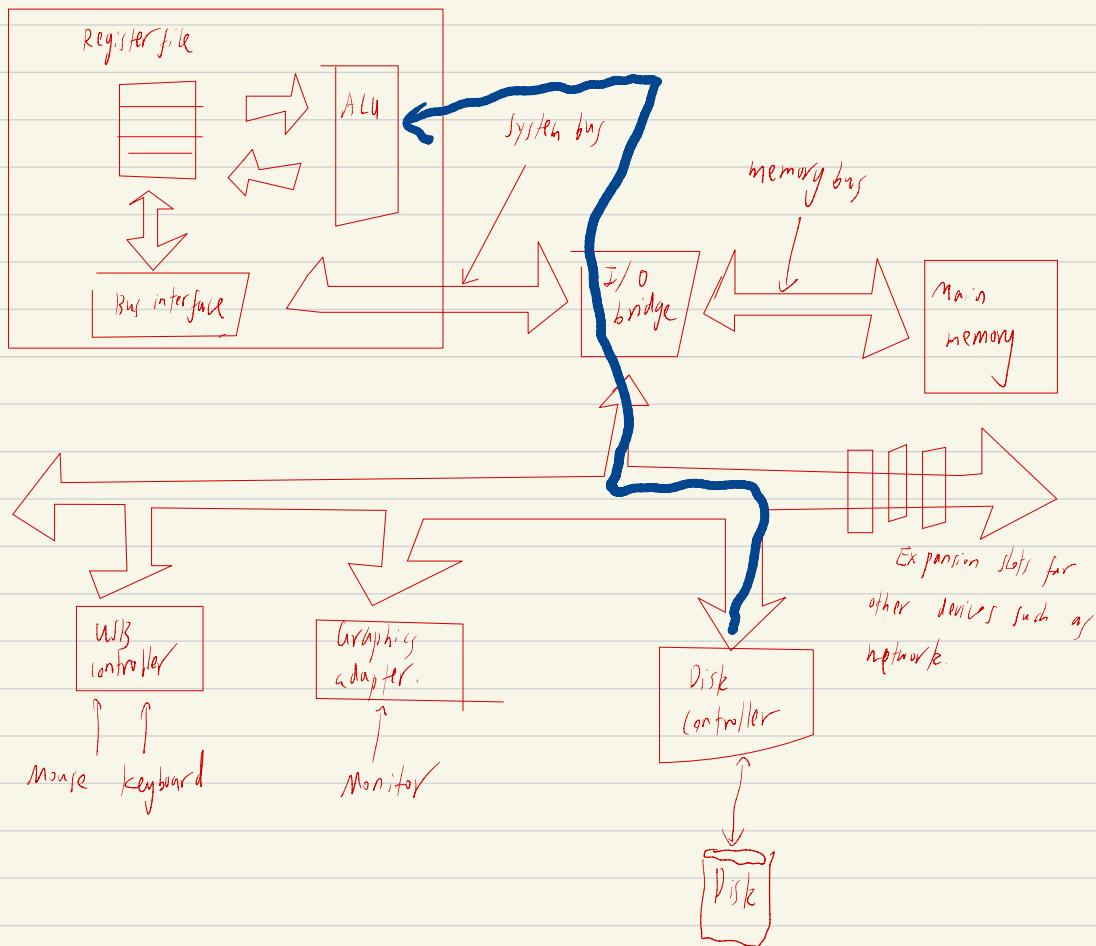
CPU chip



# Reading a Disk sector (3)

When the DMA transfer completes, the disk controller notifies the CPU with an interrupt (i.e. asserts a special "interrupt" pin on the CPU)

CPU chip



## SSD performance characteristics

|                      |            |                       |            |
|----------------------|------------|-----------------------|------------|
| Sequential read tput | 550 MB/s   | sequential write tput | 470 MB/s   |
| Random read tput     | 365 MB/s   | random write tput     | 303 MB/s   |
| Avg seq read time    | 50 $\mu$ s | Avg seq write time    | 60 $\mu$ s |

- (i) sequential always faster than random access
- (ii) common theme in the memory hierarchy.

- (iii) Random writes are somewhat slower.

- (iv) Erasing a block takes a long time (~1ms)
- (v) Modifying a block page requires all other pages to be copied to new block.
- (vi) In earlier SSDs, the read/write gap was much larger.

## Solid-state drive

SSD trade off v.s. rotating Disks

### (i) Advantages

- (i) No moving part → faster, less power, more rugged

### (ii) Disadvantages

- (i) Have the potential to wear out
- (ii) In 2011, about 30 times more expensive per byte

### (iii) Applications:

- (i) MP3 players, smart phones, laptops
- (ii) Beginning to appear in desktops and servers

## Locality to the Rescue!

The key to bridging this CPU-Memory gap is a fundamental property of computer programs known as **locality**.

**Principle of Locality:** Programs tend to use data and instructions with addresses near or equal to those they have used recently

### Temporal Locality:

- o Recently referenced items are likely to be referenced again in the near future.



### Spatial Locality:

- items with nearby addresses tend to be referenced together in long

## Locality Example

```
sum = 0;  
for(i=0; i<n; i++)  
    sum += a[i];  
return sum;
```

spatial locality  
temporal locality.

### (I) Data references

- (i) Reference array elements in succession  
(stride-1 reference pattern)

- (ii) Reference variable sum each iteration.

### (II) Instruction references

- (i) Reference instructions in sequence

- (ii) Cycle through loop repeatedly

Qualitative estimates of locality:

(claim: Being able to look at code and get a qualitative sense of its locality is a key skill for a professional programmer.

Question: Does the function have good locality with respect to array  $a$ ?

```
int sum_array_rows(int a[M][N])
```

}

```
int i, j, sum = 0;
```

```
for (i = 0; i < M; i++)
```

```
    for (j = 0; j < N; j++)
```

sum += a[i][j]; temporal and spatial locality

```
return sum;
```



both good.

}

if i & j exchanged awful (?)

skipping in memory



## Memory Hierarchies

(I) Some fundamental and enduring properties of hardware and software:

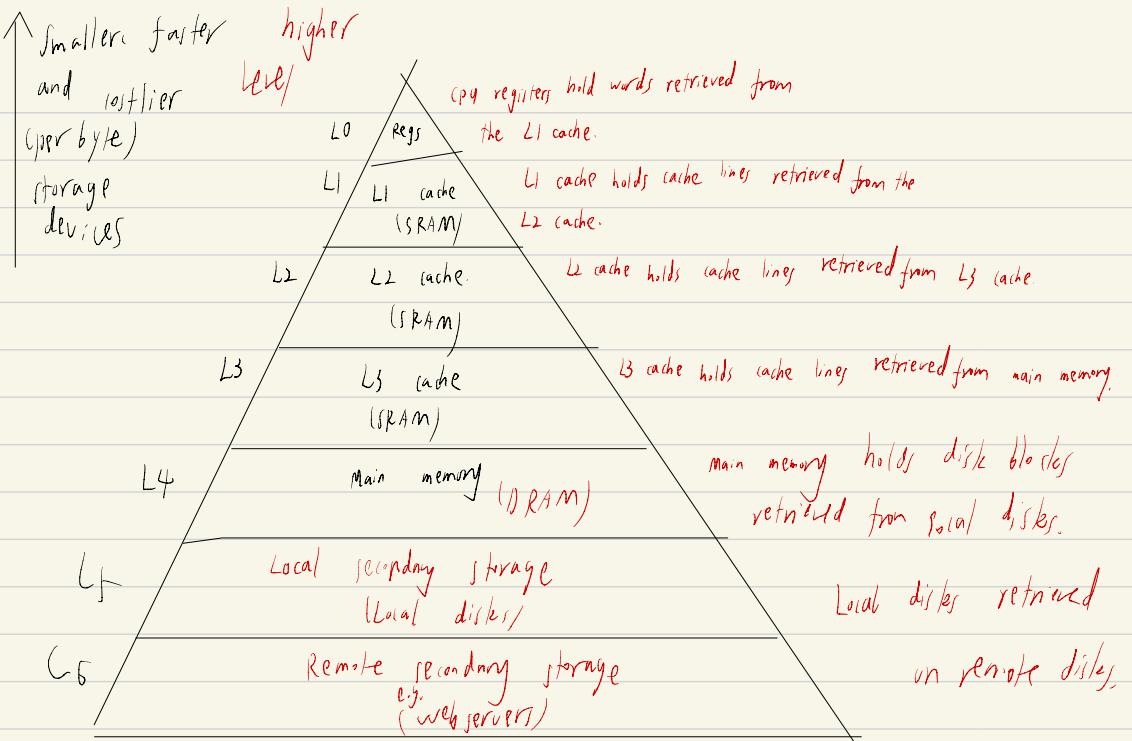
(i) fast storage technologies cost more per byte, have less capacity and require more power (heat).

(ii) The gap between CPU and main memory speed is widening.

(iii) Well-written programs tend to exhibit good locality.

(II) These fundamental properties complement each other beautifully.

(IV) They suggest an approach for organizing memory and storage systems known as a **memory hierarchy**.



Larger,  
slower,  
and cheaper  
(per byte)  
storage dev's

↓

↓ lower cost

## Caches.

(comparison: *your backpack in school*)

(i) Cache: A smaller, faster storage device that acts as a staging area for a subset of the data in a larger, slower device.

(ii) Fundamental idea of a memory hierarchy:

(i) For each  $k$ , the faster, smaller device at level  $k$  serves as a cache for the larger, slower device at level  $k+1$ .

(iii) Why do memory hierarchies work?

- (i) Because of locality, programs tend to access the data at level  $k$  more often than they access the data at level  $k+1$ .
- (ii) Thus, the storage at level  $k+1$  can be slower, and thus larger and cheaper per unit.

**Big idea:** the memory hierarchy creates a large pool of storage that costs as much as the cheap storage near the bottom, but that serves data to programs at the rate of level  $k$ . *fast storage near the top*

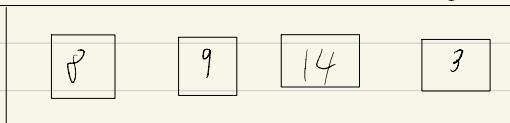
| : | : | : |

General Cache concepts.



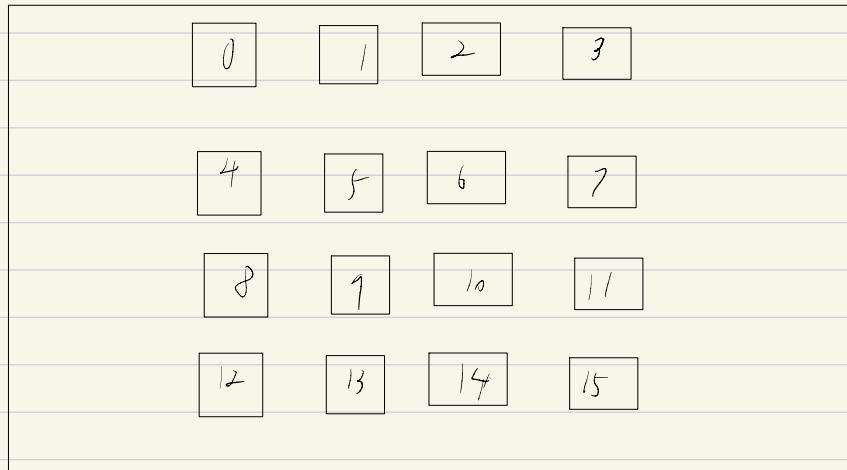
request 14

Cache



smaller, faster, more expensive memory caches a subset  
of blocks

Data is copied in block-sized transfer unit.



Larger, slower  
cheaper memory  
viewed as  
"partitioned"  
into blocks

## General Caching concepts:

### Types of Cache Misses

#### (I) cold (compulsory) miss:

(i) cold misses occur because the cache is empty.

#### (II) conflict miss:

(ii) Most caches limit blocks at level  $k+1$  to a small subset (sometimes a singleton) of the block positions at level  $k$ .

e.g. block  $i$  at level  $k+1$  must be placed in block  $(i \bmod 4)$  at level  $k$ .

(iii) conflict misses occur when the level  $k$  cache is large enough, but multiple data objects all map to the same level  $k$  block.

e.g. Referencing blocks 0, 8, 0, 8, 0, 8, ... would miss every time.

#### (III) capacity miss

(iv) occurs when the set of active cache blocks (working set) is larger than the cache

Summary:

- (i) The speed gap between CPU, memory and mass storage continues to widen.
- (ii) Well-written programs exhibit a property called locality.
- (iii) Memory hierarchy based on caching close the gap by exploiting locality.