



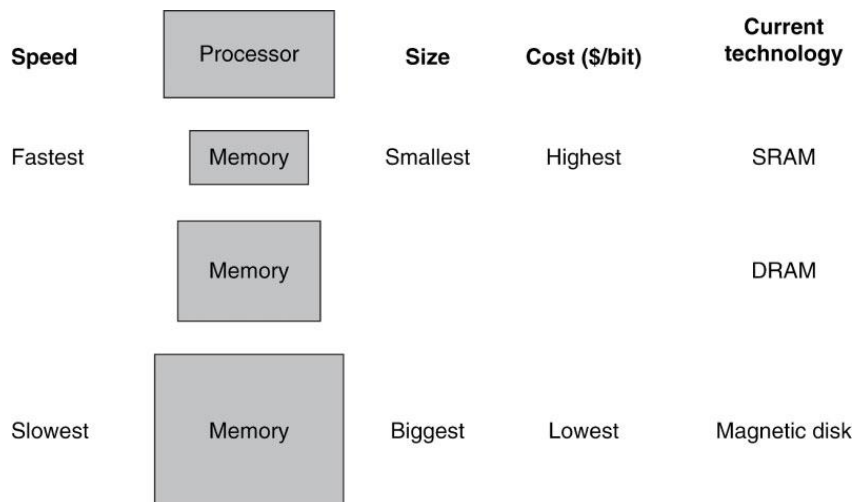
# Exploiting Memory Hierarchy

COMPUTER ORGANIZATION AND ARCHITECTURE

## Black Friday Deal



# Memory hierarchy



CSULB

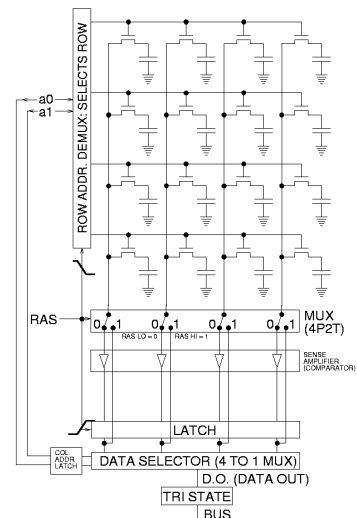
## Memory Technology

- Static RAM (SRAM)
  - 0.5ns – 2.5ns, \$2000 – \$5000 per GB
- Dynamic RAM (DRAM)
  - 50ns – 70ns, \$20 – \$75 per GB
- Magnetic disk
  - 5ms – 20ms, \$0.20 – \$2 per GB
- Ideal memory
  - Access time of SRAM
  - Capacity and cost/GB of disk

CSULB

# DRAM Technology

- Data stored as a charge in a capacitor
  - Single transistor used to access the charge
  - Must periodically be refreshed
    - Read contents and write back
    - Performed on a DRAM “row”
- Bits in a DRAM are organized as a rectangular array
  - DRAM accesses an entire row

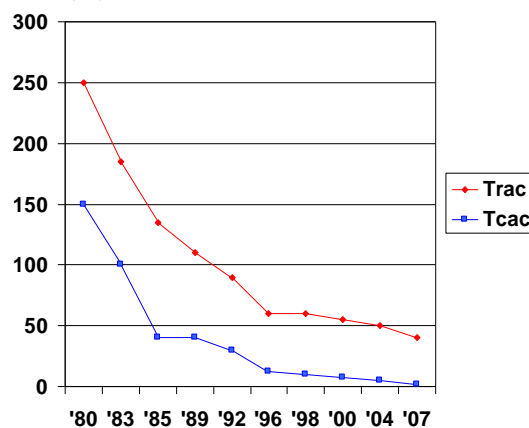


CSULB

## DRAM Generations

| Year | Capacity | \$/GB     |
|------|----------|-----------|
| 1980 | 64Kbit   | \$1500000 |
| 1983 | 256Kbit  | \$500000  |
| 1985 | 1Mbit    | \$200000  |
| 1989 | 4Mbit    | \$50000   |
| 1992 | 16Mbit   | \$15000   |
| 1996 | 64Mbit   | \$10000   |
| 1998 | 128Mbit  | \$4000    |
| 2000 | 256Mbit  | \$1000    |
| 2004 | 512Mbit  | \$250     |
| 2007 | 1Gbit    | \$50      |

access time (ns)

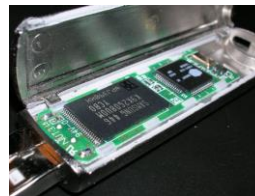


CSULB

# Flash Storage

---

- Nonvolatile semiconductor storage
  - 100× – 1000× faster than disk
  - Smaller, lower power, more robust
  - But more \$/GB (between disk and DRAM)



CSULB

## Flash Types

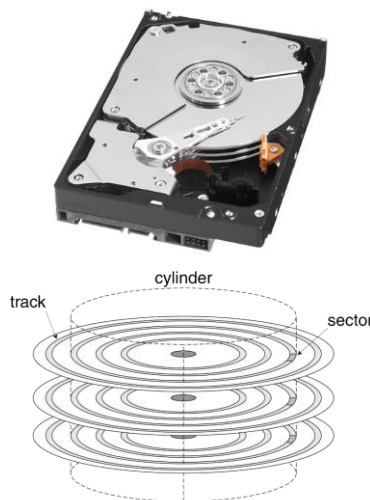
---

- NOR flash: bit cell like a NOR gate
  - Random read/write access
  - Used for instruction memory in embedded systems
- NAND flash: bit cell like a NAND gate
  - Denser (bits/area), but block-at-a-time access
  - Cheaper per GB
  - Used for USB keys, media storage, ...
- Flash bits wears out after 1000's of accesses
  - Not suitable for direct RAM or disk replacement
  - Wear leveling: remap data to less used blocks

CSULB

# Disk Sectors and Access

- Nonvolatile, rotating magnetic storage
- Each sector records
  - Sector ID
  - Data (512 bytes, 4096 bytes proposed)
  - Error correcting code (ECC)
    - Used to hide defects and recording errors
- Access to a sector involves
  - Queuing delay if other accesses are pending
  - Seek: move the heads
  - Rotational latency
  - Data transfer
  - Controller overhead



CSULB

## Disk Access Example

- Given
  - 512B sector, 15,000rpm, 4ms average seek time, 100MB/s transfer rate, 0.2ms controller overhead, idle disk
- Average read time
  - 4ms seek time
  - +  $\frac{1}{2} / (15,000/60) = 2\text{ms}$  rotational latency
  - +  $512 / 100\text{MB/s} = 0.005\text{ms}$  transfer time
  - + 0.2ms controller delay
  - = 6.2ms

CSULB

# Principle of Locality

- Programs access a small proportion of their address space at any time
- Temporal** locality
  - Items accessed recently are likely to be accessed again soon
  - e.g., instructions in a loop, induction variables
- Spatial** locality
  - Items near those accessed recently are likely to be accessed soon
  - E.g., sequential instruction access, array data

CSULB

## Taking Advantage of Locality

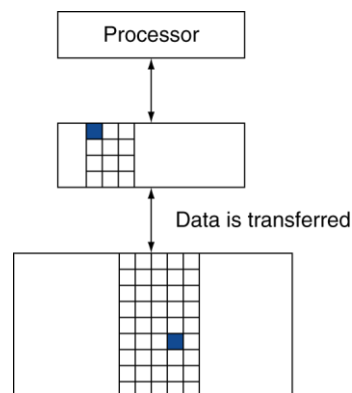
- Exploit memory hierarchy
- Store everything on disk
- Copy recently accessed (and nearby) items from disk to smaller DRAM memory
  - Main memory
- Copy more recently accessed (and nearby) items from DRAM to smaller SRAM memory
  - Cache memory attached to CPU

| Speed   |           | Size     | Cost (\$/bit) | Current technology |
|---------|-----------|----------|---------------|--------------------|
| Fastest | Processor |          |               |                    |
|         | Memory    | Smallest | Highest       | SRAM               |
|         | Memory    |          |               | DRAM               |
| Slowest | Memory    | Biggest  | Lowest        | Magnetic disk      |

CSULB

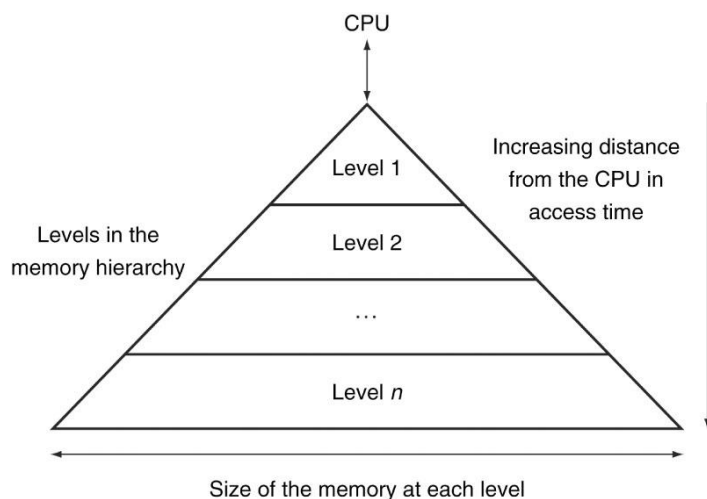
# Memory Hierarchy Levels

- Block (aka line): unit of copying
  - May be multiple words
- If accessed data is present in upper level
  - Hit: access satisfied by upper level
  - Hit ratio: hits/accesses
- If accessed data is absent
  - Miss: block copied from lower level
    - Time taken: miss penalty
    - Miss ratio: misses/accesses  
=  $1 - \text{hit ratio}$
  - Then accessed data supplied from upper level



CSULB

## Memory Hierarchy

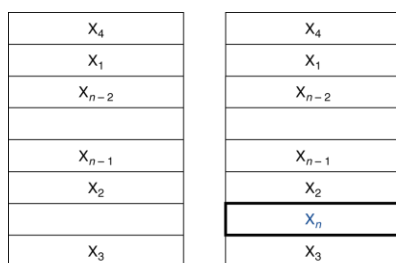


CSULB

# Cache Memory

- Cache memory
  - The level of the memory hierarchy closest to the CPU
- Given accesses  $X_1, \dots, X_{n-1}, X_n$

Mem [0xFFC]



a. Before the reference to  $X_n$

b. After the reference to  $X_n$

- How do we know if the data is present?
- Where do we look?

CSULB

## How is the hierarchy managed?

- Registers  $\Leftrightarrow$  memory
  - by compiler (or programmer)
- cache  $\Leftrightarrow$  main memory
  - by the cache controller hardware
- main memory  $\Leftrightarrow$  disks
  - by the operating system
  - by the programmer

Mem [0xFFC]

fread

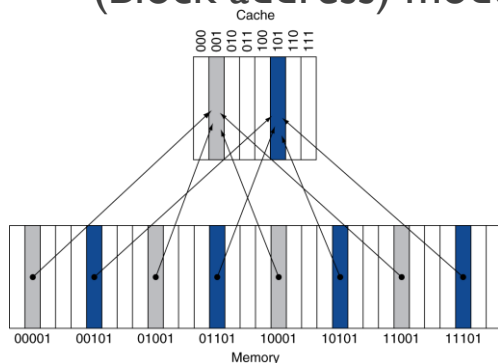
run.exe ./a.out

CSULB



# Direct Mapped Cache

- Location determined by address
- Direct mapped: only one choice
  - (Block address) modulo (#Blocks in cache)



- #Blocks is a power of 2.  
Why?
- Use low-order address bits

$$\begin{aligned}
 0011 &\% 2 = 1 \\
 0101 &\% 2 = 1 \\
 1010 &\% 2^2 = 10 \\
 1111 &\% 2^3 = 111
 \end{aligned}$$

CSULB

## Tags and Valid Bits

- How do we know which particular block is stored in a cache location?
  - Store block address as well as the data
  - Actually, only need the high-order bits
  - Called the tag
- What if there is no data in a location?
  - Valid bit: 1 = present, 0 = not present
  - Initially 0

| Index | V | Tag | Data |
|-------|---|-----|------|
| 000   | N |     |      |
| 001   | N |     |      |
| 010   | N |     |      |
| 011   | N |     |      |
| 100   | N |     |      |
| 101   | N |     |      |
| 110   | N |     |      |
| 111   | N |     |      |

CSULB

# Cache Example

## ■ 8-blocks, direct mapped

Cache Memory

| Index | V             | Tag           | Data       |
|-------|---------------|---------------|------------|
| 000   | <del>NY</del> | 10            | Mem[10000] |
| 001   | N             |               | 10010      |
| 010   | <del>NY</del> | <del>10</del> | Mem[11010] |
| 011   | <del>NY</del> | 00            | Mem[00011] |
| 100   | N             |               |            |
| 101   | N             |               |            |
| 110   | <del>NY</del> | 10            | Mem[10110] |
| 111   | N             |               |            |

Memory Access

| Order | Word addr | Binary addr | Hit/miss |
|-------|-----------|-------------|----------|
| 1     | 22        | 10110       | M        |
| 2     | 26        | 11010       | M        |
| 3     | 22        | 10110       | H        |
| 4     | 26        | 11010       | H        |
| 5     | 16        | 10000       | M        |
| 6     | 3         | 00011       | M        |
| 7     | 16        | 10000       | H        |
| 8     | 18        | 10010       | M        |

11/28