inst.eecs.berkeley.edu/~cs61c

# CS61C : Machine Structures

## Lecture #23 Cache I
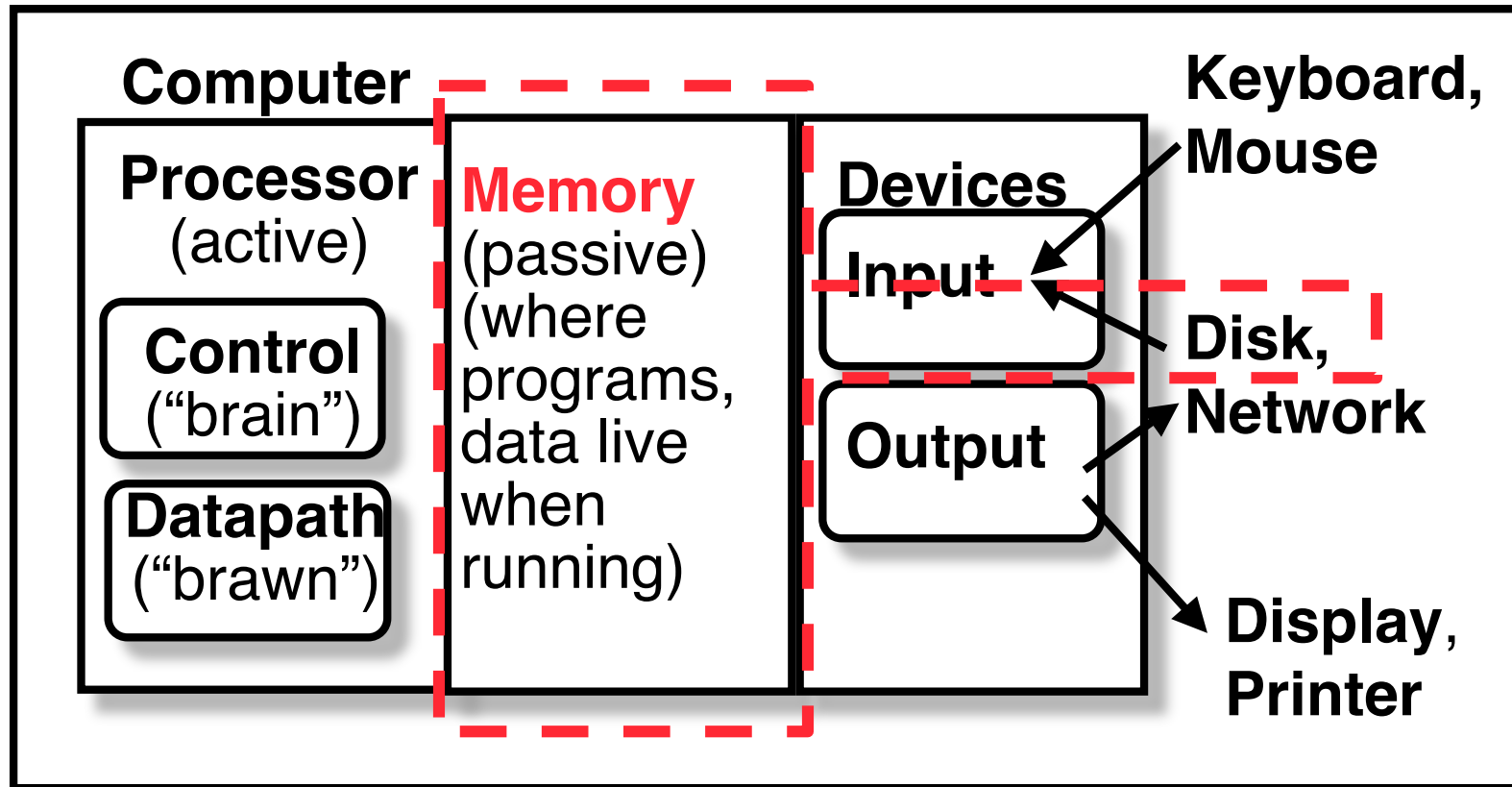
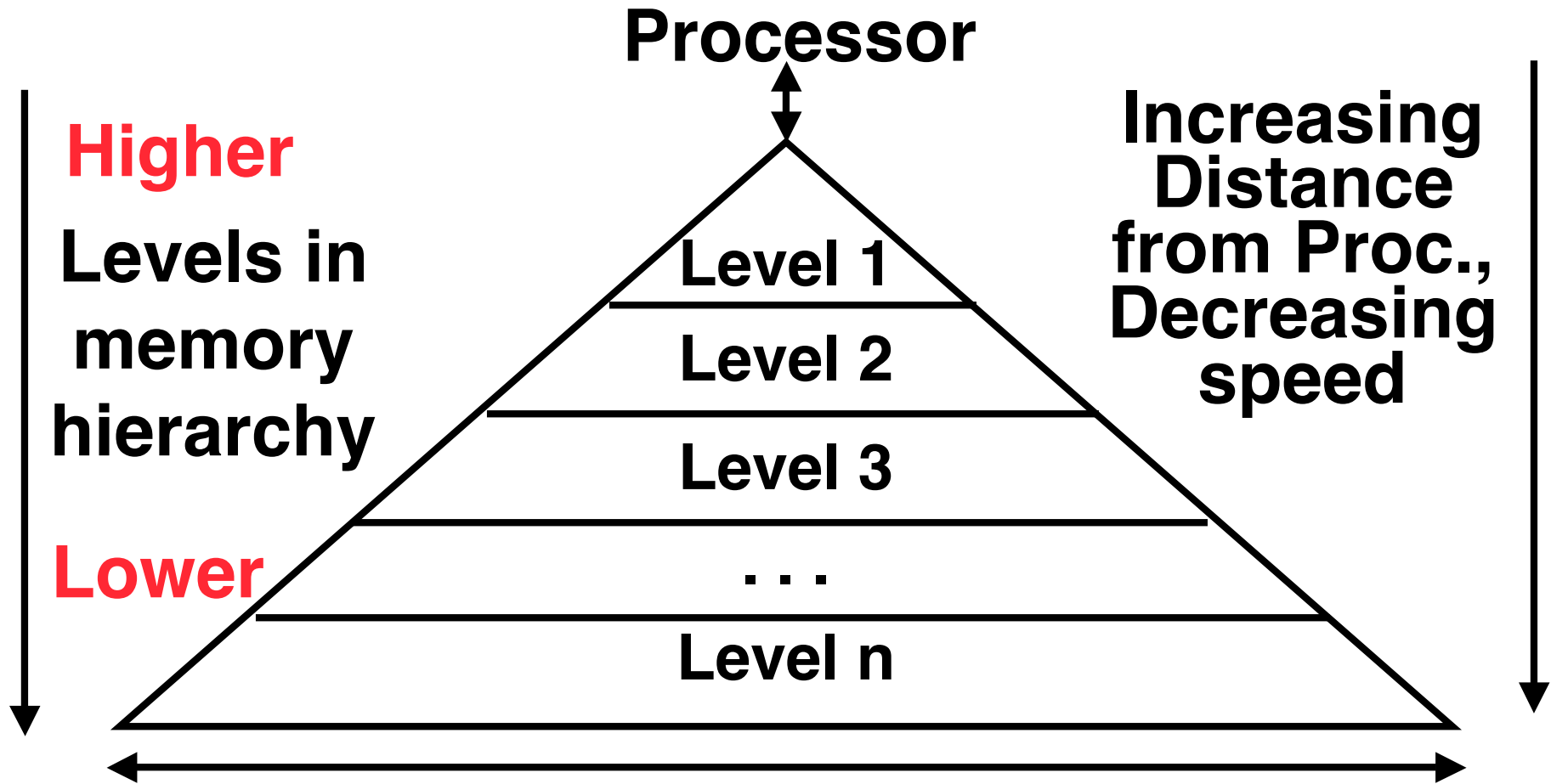## 2007-8-2

**Scott Beamer, Instructor**

# The Big Picture



**Computer**

**Processor** (active)

**Control** ("brain")

**Datapath** ("brawn")

**Memory** (passive) (where programs, data live when running)

**Devices**

**Input**

**Output**

**Keyboard, Mouse**

**Disk, Network**

**Display, Printer**

# Memory Hierarchy



**Processor**

**Higher**

**Levels in memory hierarchy**

**Lower**

Level 1

Level 2

Level 3

. . .

Level n

**Increasing Distance from Proc., Decreasing speed**

**Size of memory at each level**

*As we move to deeper levels the latency goes up and price per bit goes down.*

# Memory Hierarchy

- **If level closer to Processor, it is:**
  - **smaller**
  - **faster**
  - **subset of lower levels (contains most recently used data)**

- **Lowest Level (usually disk) contains all available data (or does it go beyond the disk?)**

- **Memory Hierarchy presents the processor with the illusion of a very large very fast memory.**

# Memory Hierarchy Analogy: Library (1/2)

- **You're writing a term paper (Processor) at a table in Doe**

- **Doe Library is equivalent to disk**
  - essentially limitless capacity
  - very slow to retrieve a book

- **Table is main memory**
  - smaller capacity: means you must return book when table fills up
  - easier and faster to find a book there once you've already retrieved it

# Memory Hierarchy Analogy: Library (2/2)

- **Open books on table are <u>cache</u>**

  - smaller capacity: can have very few open books fit on table; again, when table fills up, you must close a book

  - much, much faster to retrieve data

- **Illusion created: whole library open on the tabletop**

  - Keep as many recently used books open on table as possible since likely to use again

  - Also keep as many books on table as possible, since faster than going to library

# Memory Hierarchy Basis

- **Cache contains copies of data in memory that are being used.**

- **Memory contains copies of data on disk that are being used.**

- **Caches work on the principles of temporal and spatial locality.**

  - **Temporal Locality: if we use it now, chances are we'll want to use it again soon.**

  - **Spatial Locality: if we use a piece of memory, chances are we'll use the neighboring pieces soon.**

# Cache Design

- **How do we organize cache?**

- **Where does each memory address map to?**

   **(Remember that cache is subset of memory, so multiple memory addresses map to the same cache location.)**

- **How do we know which elements are in cache?**

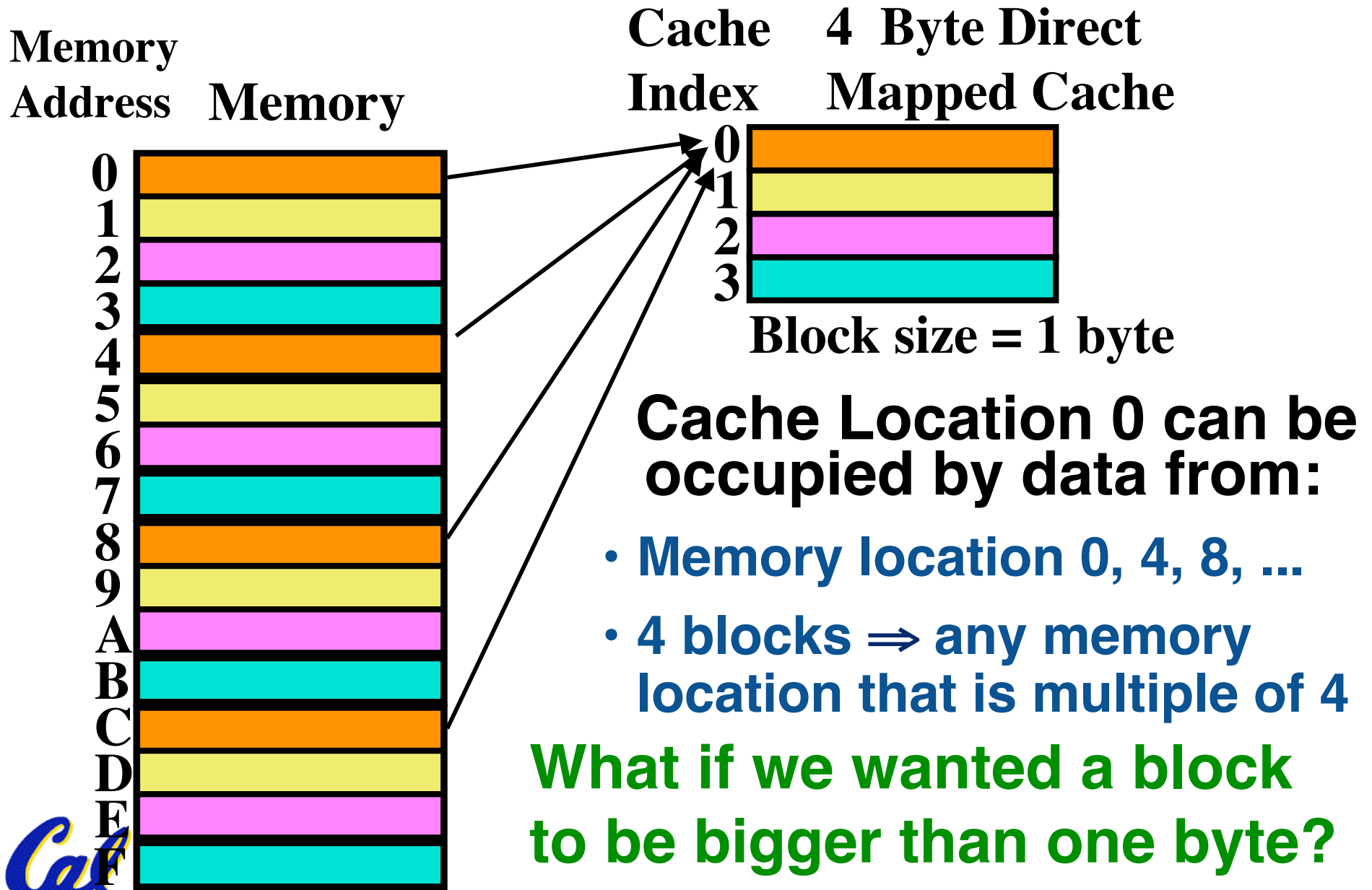- **How do we quickly locate them?**
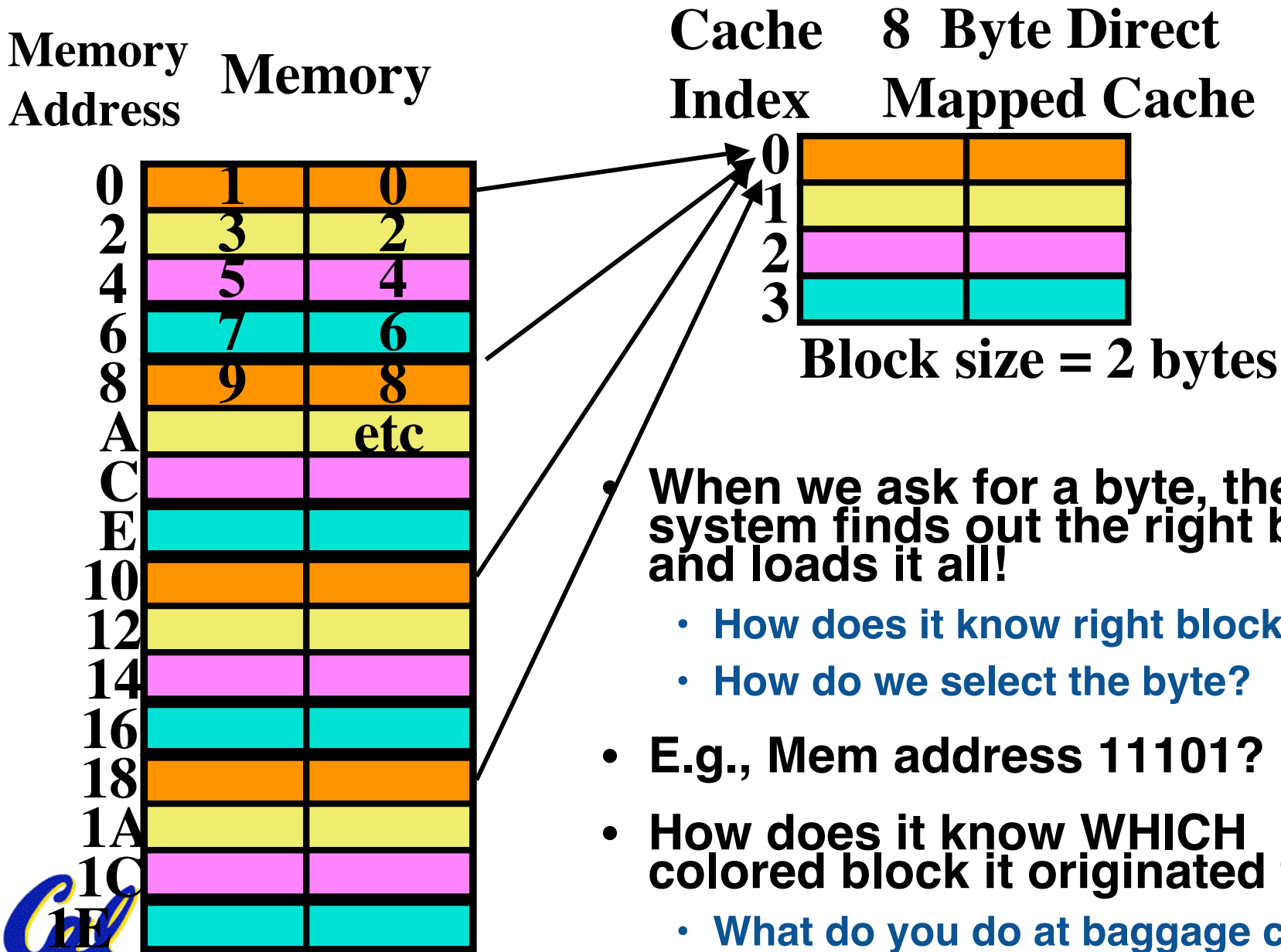
# Direct-Mapped Cache (1/4)

- **In a <u>direct-mapped cache</u>, each memory address is associated with one possible <u>block</u> within the cache**

  - **Therefore, we only need to look in a single location in the cache for the data if it exists in the cache**

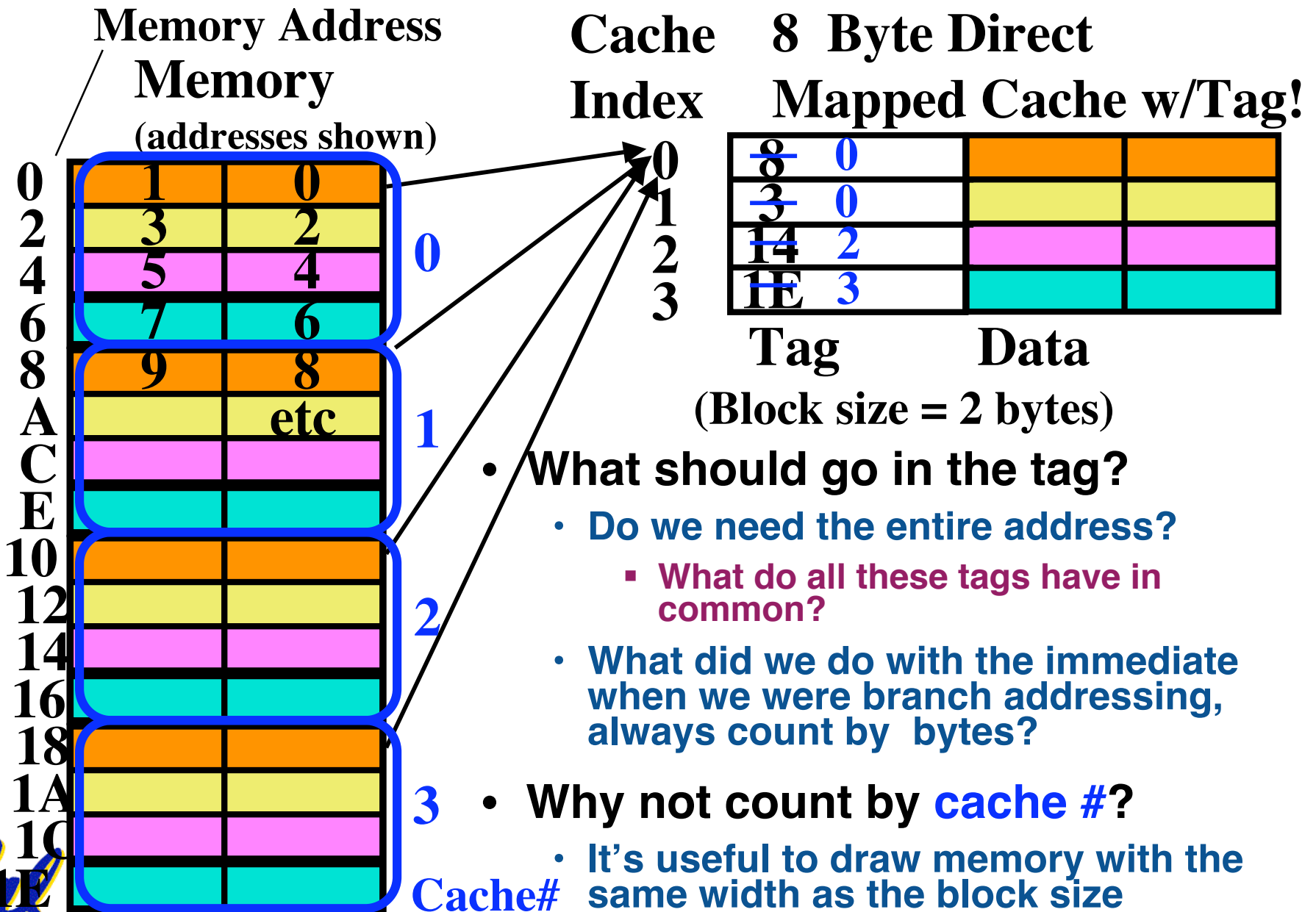  - **Block is the unit of transfer between cache and memory**

# Direct-Mapped Cache (2/4)

Memory
Address     Memory

Cache     4 Byte Direct
Index     Mapped Cache

0
1
2
3

Block size = 1 byte

## Cache Location 0 can be occupied by data from:

- Memory location 0, 4, 8, ...

- 4 blocks ⇒ any memory location that is multiple of 4

**What if we wanted a block to be bigger than one byte?**

# Direct-Mapped Cache (3/4)

**Memory Address**

**Memory**

| | |
|---|---|
| 1 | 0 |
| 3 | 2 |
| 5 | 4 |
| 7 | 6 |
| 9 | 8 |
| | etc |
| | |
| | |
| | |
| | |
| | |
| | |
| | |
| | |
| | |
| | |

0
2
4
6
8
A
C
E
10
12
14
16
18
1A
1C
1E

**Cache Index**

**8 Byte Direct Mapped Cache**

0
1
2
3

**Block size = 2 bytes**

- **When we ask for a byte, the system finds out the right block, and loads it all!**
  - **How does it know right block?**
  - **How do we select the byte?**

- **E.g., Mem address 11101?**

- **How does it know WHICH colored block it originated from?**
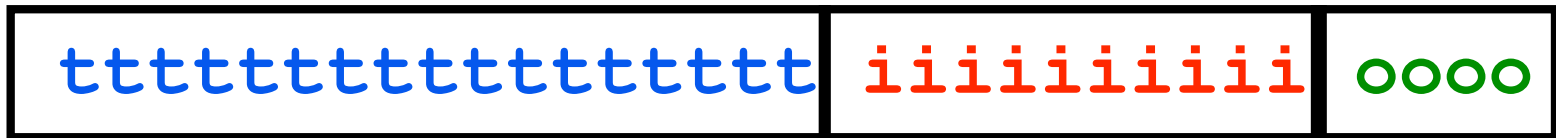  - **What do you do at baggage claim?**

# Direct-Mapped Cache (4/4)

**Memory Address**
**Memory**
(addresses shown)

| 1 | 0 |
|---|---|
| 3 | 2 |
| 5 | 4 |
| 7 | 6 |

**0**

| 9 | 8 |
|---|---|
|   | etc |
|   |   |
|   |   |

**1**

**2**

**3**

Memory addresses (left column): 0 2 4 6 8 A C E 10 12 14 16 18 1A 1C 1E

**Cache#**

**Cache Index**

| 0 |
|---|
| 1 |
| 2 |
| 3 |

**8 Byte Direct Mapped Cache w/Tag!**

| Tag | | Data | |
|---|---|---|---|
| 8 | 0 | | |
| 3 | 0 | | |
| 14 | 2 | | |
| 1E | 3 | | |

**Tag**          **Data**

(Block size = 2 bytes)

- **What should go in the tag?**
  - Do we need the entire address?
    - **What do all these tags have in common?**
  - What did we do with the immediate when we were branch addressing, always count by  bytes?

- **Why not count by cache #?**
  - It's useful to draw memory with the same width as the block size

# Issues with Direct-Mapped

- **Since multiple memory addresses map to same cache index, how do we tell which one is in there?**

- **What if we have a block size > 1 byte?**

- **Answer: divide memory address into three fields**

| tttttttttttttttttt | iiiiiiiiii | oooo |
|---|---|---|
| **tag** <br> **to check** <br> **if have** <br> **correct block** | **index** <br> **to** <br> **select** <br> **block** | **byte** <br> **offset** <br> **within** <br> **block** |

# Direct-Mapped Cache Terminology

- All fields are read as unsigned integers.

- **Index**: specifies the cache index (which "row"/block of the cache we should look in)

- **Offset**: once we've found correct block, specifies which byte within the block we want

- **Tag**: the remaining bits after offset and index are determined; these are used to distinguish between all the memory addresses that map to the same location
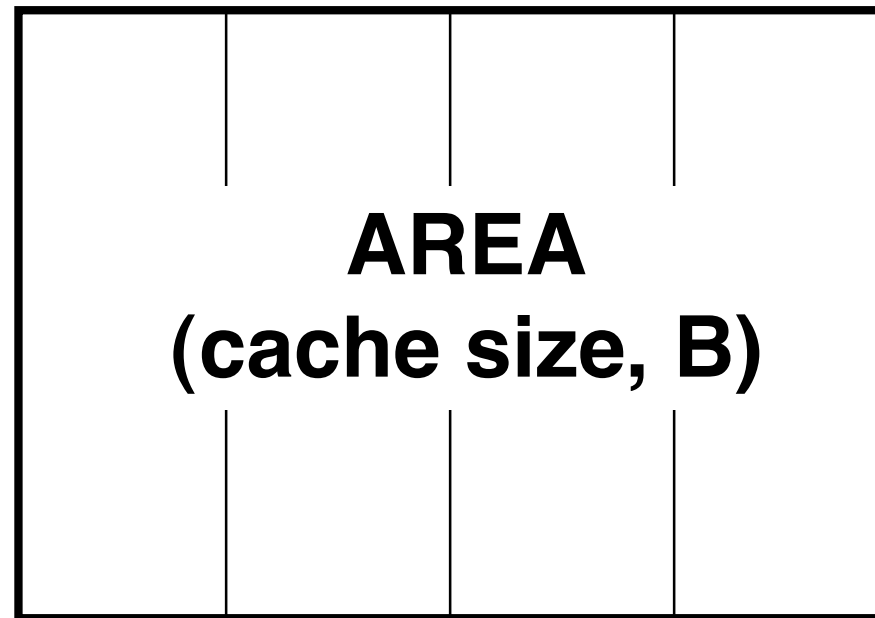
# TIO Dan's great cache mnemonic

AREA (cache size, B)
 = HEIGHT (# of blocks)
   * WIDTH (size of one block, B/block)

$$2^{(H+W)} = 2^H * 2^W$$

| Tag | Index | Offset |
|-----|-------|--------|

WIDTH
(size of one block, B/block)

HEIGHT
(# of blocks)

AREA
(cache size, B)

# Direct-Mapped Cache Example (1/3)

- **Suppose we have a 16KB of data in a direct-mapped cache with 4 word blocks**

- **Determine the size of the tag, index and offset fields if we're using a 32-bit architecture**

- **Offset**

  - **need to specify correct byte within a block**

  - **block contains 4 words**

    = 16 bytes

    = $2^4$ bytes

  - **need <u>4 bits</u> to specify correct byte**

# Direct-Mapped Cache Example (2/3)

- **Index: (~index into an "array of blocks")**

  - **need to specify correct block in cache**

  - **cache contains 16 KB = $2^{14}$ bytes**

  - **block contains $2^4$ bytes (4 words)**

  - **# blocks/cache**

    $$= \frac{\text{bytes/cache}}{\text{bytes/block}}$$

    $$= \frac{2^{14} \text{ bytes/cache}}{2^4 \text{ bytes/block}}$$

    $$= 2^{10} \text{ blocks/cache}$$

  - **need 10 bits to specify this many blocks**

# Direct-Mapped Cache Example (3/3)

- **Tag: use remaining bits as tag**

  - tag length = addr length – offset - index
    = 32 - 4 - 10 bits
    = 18 bits

  - so tag is leftmost **18 bits** of memory address

- **Why not full 32 bit address as tag?**

  - All bytes within block need same address (4b)

  - Index must be same for every address within a block, so it's redundant in tag check, thus can leave off to save memory (here 10 bits)

# Caching Terminology

- When we try to read memory,
  3 things can happen:

1. **cache hit**:
   cache block is valid and contains
   proper address, so read desired word

2. **cache miss**:
   nothing in cache in appropriate block,
   so fetch from memory

3. **cache miss, block replacement**:
   wrong data is in cache at appropriate
   block, so discard it and fetch desired
   data from memory (cache always copy)

# Administrivia

- **Assignments**
  - **HW7 due Tonight 8/2**
  - **Proj3 due Sunday 8/5**

- **Proj3 will have face-to-face grading**
  - **You will be able to sign up online tonight or tomorrow for timeslots next week**

- **Today's lab is "non-trivial," so work in groups and make sure you understand it**

- **Course Survey during last lecture**

# Block Size Tradeoff (1/3)

- **Benefits of Larger Block Size**
  - **Spatial Locality**: if we access a given word, we're likely to access other nearby words soon
  - Very applicable with Stored-Program Concept: if we execute a given instruction, it's likely that we'll execute the next few as well
  - Works nicely in sequential array accesses too

# Block Size Tradeoff (2/3)

- **Drawbacks of Larger Block Size**
  - **Larger block size means larger miss penalty**
    - on a miss, takes longer time to load a new block from next level
  - **If block size is too big relative to cache size, then there are too few blocks**
    - Result: miss rate goes up

- **In general, minimize Average Memory Access Time (AMAT)**
  - **= Hit Time**
    - **+ Miss Penalty x Miss Rate**

# Block Size Tradeoff (3/3)

- **Hit Time** = time to find and retrieve data from current level cache

- **Miss Penalty** = average time to retrieve data on a current level miss (includes the possibility of misses on successive levels of memory hierarchy)

- **Hit Rate** = % of requests that are found in current level cache
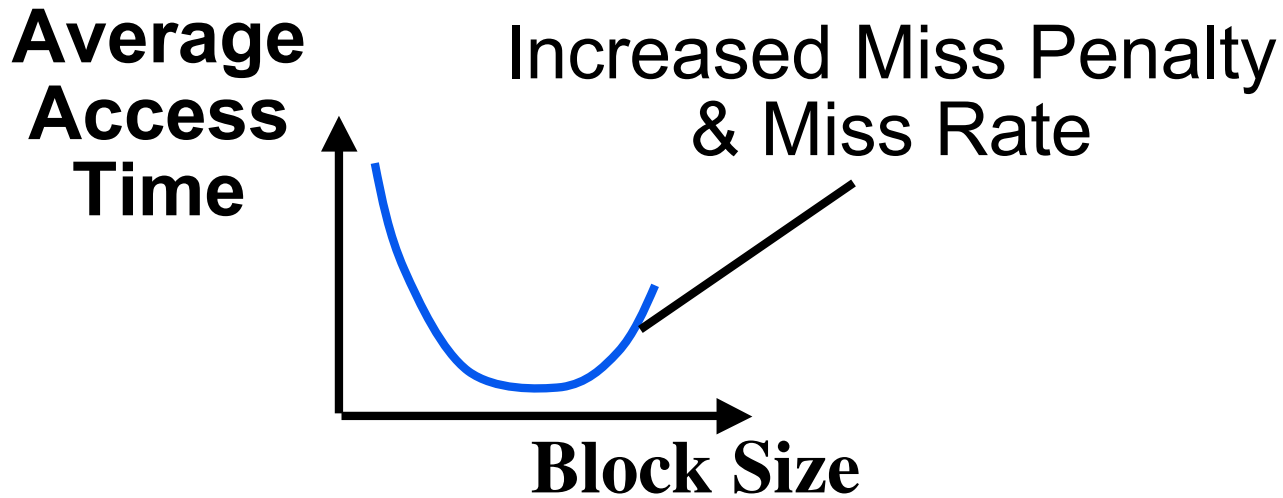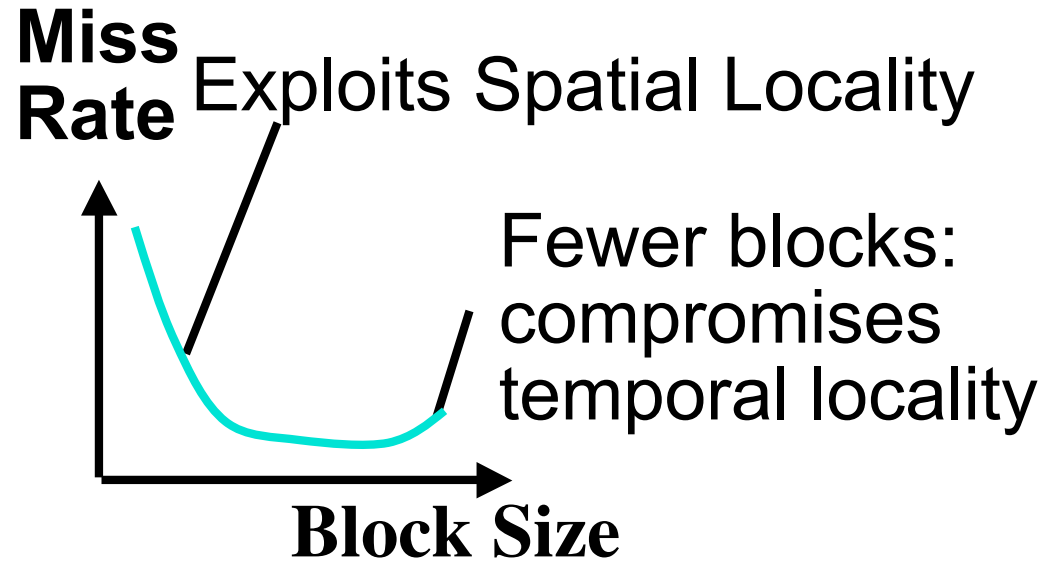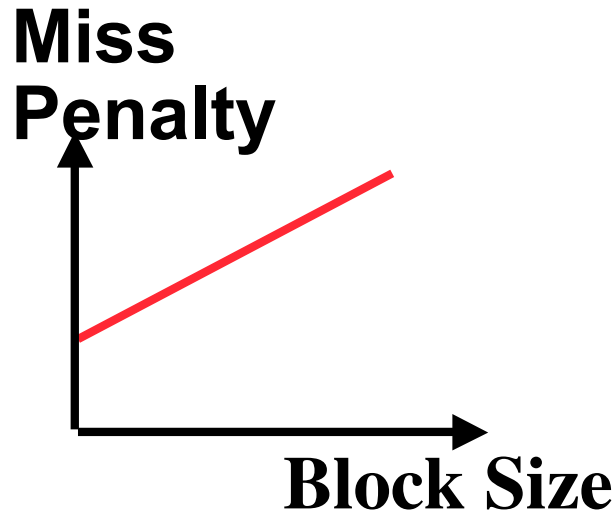
- **Miss Rate** = 1 - Hit Rate

# Extreme Example: One Big Block

| Valid Bit | Tag | Cache Data |
|---|---|---|
| ☐ | | B 3 \| B 2 \| B 1 \| B 0 |

- **Cache Size = 4 bytes    Block Size = 4 bytes**
  - Only **ONE** entry (row) in the cache!

- **If item accessed, likely accessed again soon**
  - But unlikely will be accessed again immediately!

- **The next access will likely to be a miss again**
  - Continually loading data into the cache but discard data (force out) before use it again
  - Nightmare for cache designer: **Ping Pong Effect**

# Block Size Tradeoff Conclusions

**Miss Penalty**



**Block Size**

**Miss Rate** Exploits Spatial Locality

Fewer blocks: compromises temporal locality



**Block Size**

**Average Access Time**

Increased Miss Penalty & Miss Rate



**Block Size**

# Types of Cache Misses (1/2)

- "Three Cs" Model of Misses

- **1st C: Compulsory Misses**

  - occur when a program is first started

  - cache does not contain any of that program's data yet, so misses are bound to occur

  - can't be avoided easily, so won't focus on these in this course

# Types of Cache Misses (2/2)

- **2nd C: Conflict Misses**
    - miss that occurs because two distinct memory addresses map to the same cache location
    - two blocks (which happen to map to the same location) can keep overwriting each other
    - big problem in direct-mapped caches
    - how do we lessen the effect of these?

- **Dealing with Conflict Misses**
    - Solution 1: Make the cache size bigger
        - Fails at some point
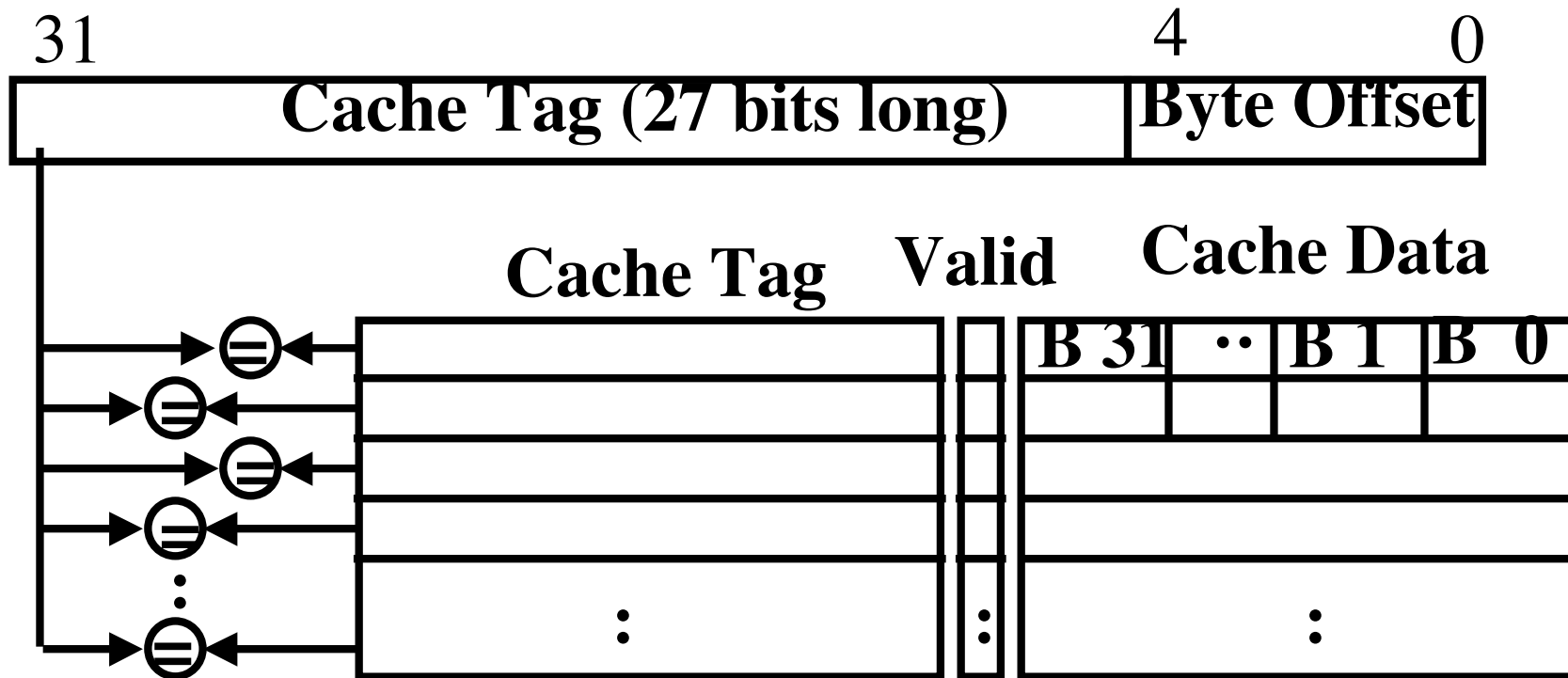    - Solution 2: Multiple distinct blocks can fit in the same cache Index?

# Fully Associative Cache (1/3)

- **Memory address fields:**
  - Tag: same as before
  - Offset: same as before
  - Index: non-existant

- **What does this mean?**
  - no "rows": any block can go anywhere in the cache
  - must compare with all tags in entire cache to see if data is there

# Fully Associative Cache (2/3)

- **Fully Associative Cache (e.g., 32 B block)**
  - **compare tags in parallel**

# Fully Associative Cache (3/3)

- **Benefit of Fully Assoc Cache**

  - **No Conflict Misses (since data can go anywhere)**

- **Drawbacks of Fully Assoc Cache**

  - **Need hardware comparator for every single entry: if we have a 64KB of data in cache with 4B entries, we need 16K comparators: infeasible**

# Third Type of Cache Miss

- **Capacity Misses**
  - miss that occurs because the cache has a limited size
  - miss that would not occur if we increase the size of the cache
  - sketchy definition, so just get the general idea

- This is the primary type of miss for Fully Associative caches.
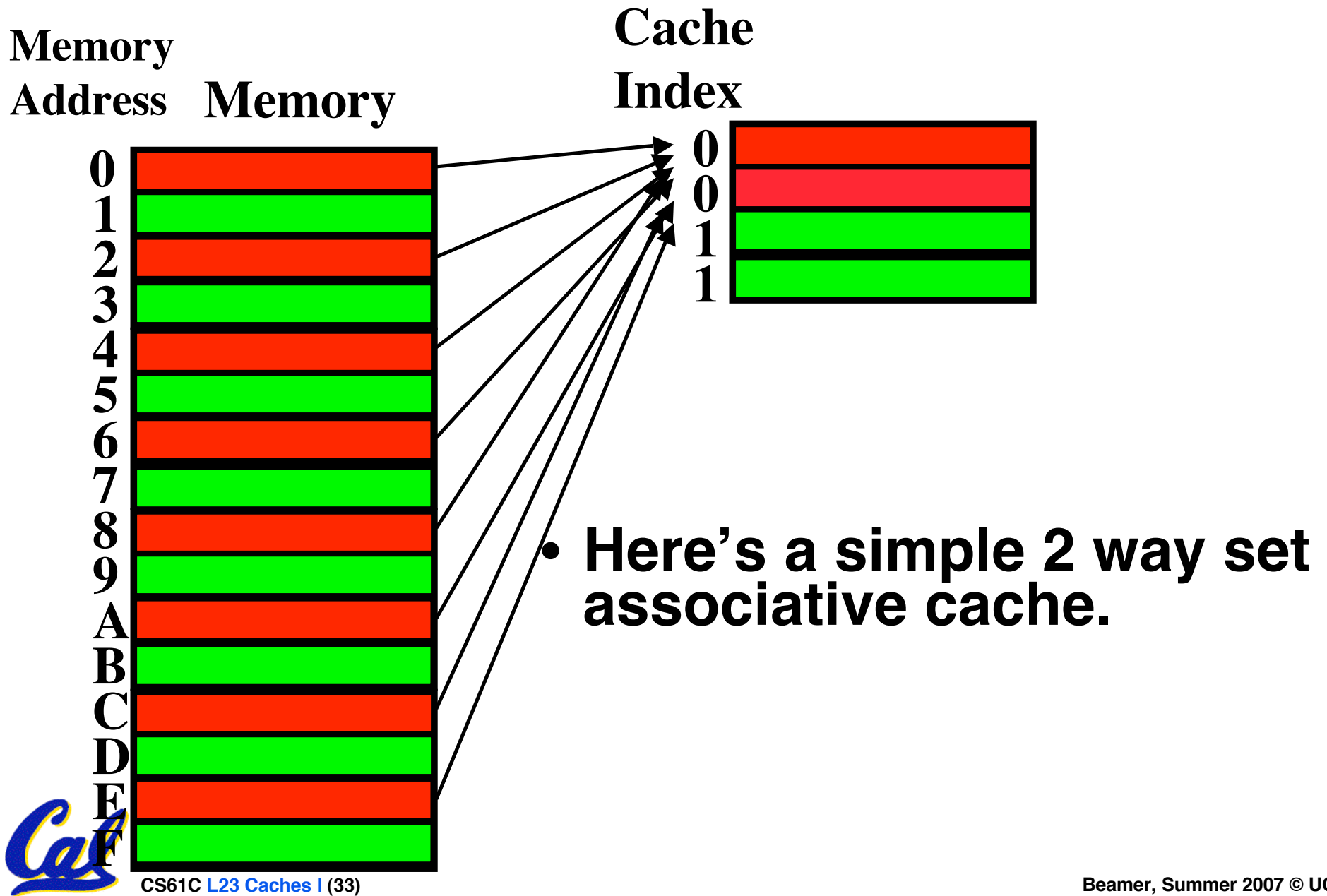
# N-Way Set Associative Cache (1/3)

- **Memory address fields:**
  - **Tag**: same as before
  - **Offset**: same as before
  - **Index**: points us to the correct "row" (called a <u>set</u> in this case)

- **So what's the difference?**
  - each set contains multiple blocks
  - once we've found correct set, must compare with all tags in that set to find our data

# Associative Cache Example

**Memory Address**  **Memory**

**Cache Index**

0
1
2
3
4
5
6
7
8
9
A
B
C
D
E
F

0
0
1
1

- **Here's a simple 2 way set associative cache.**

# N-Way Set Associative Cache (2/3)

- **Basic Idea**
  - cache is direct-mapped w/respect to sets
  - each set is fully associative
  - basically N direct-mapped caches working in parallel: each has its own valid bit and data

- **Given memory address:**
  - Find correct set using Index value.
  - Compare Tag with all Tag values in the determined set.
  - If a match occurs, hit!, otherwise a miss.
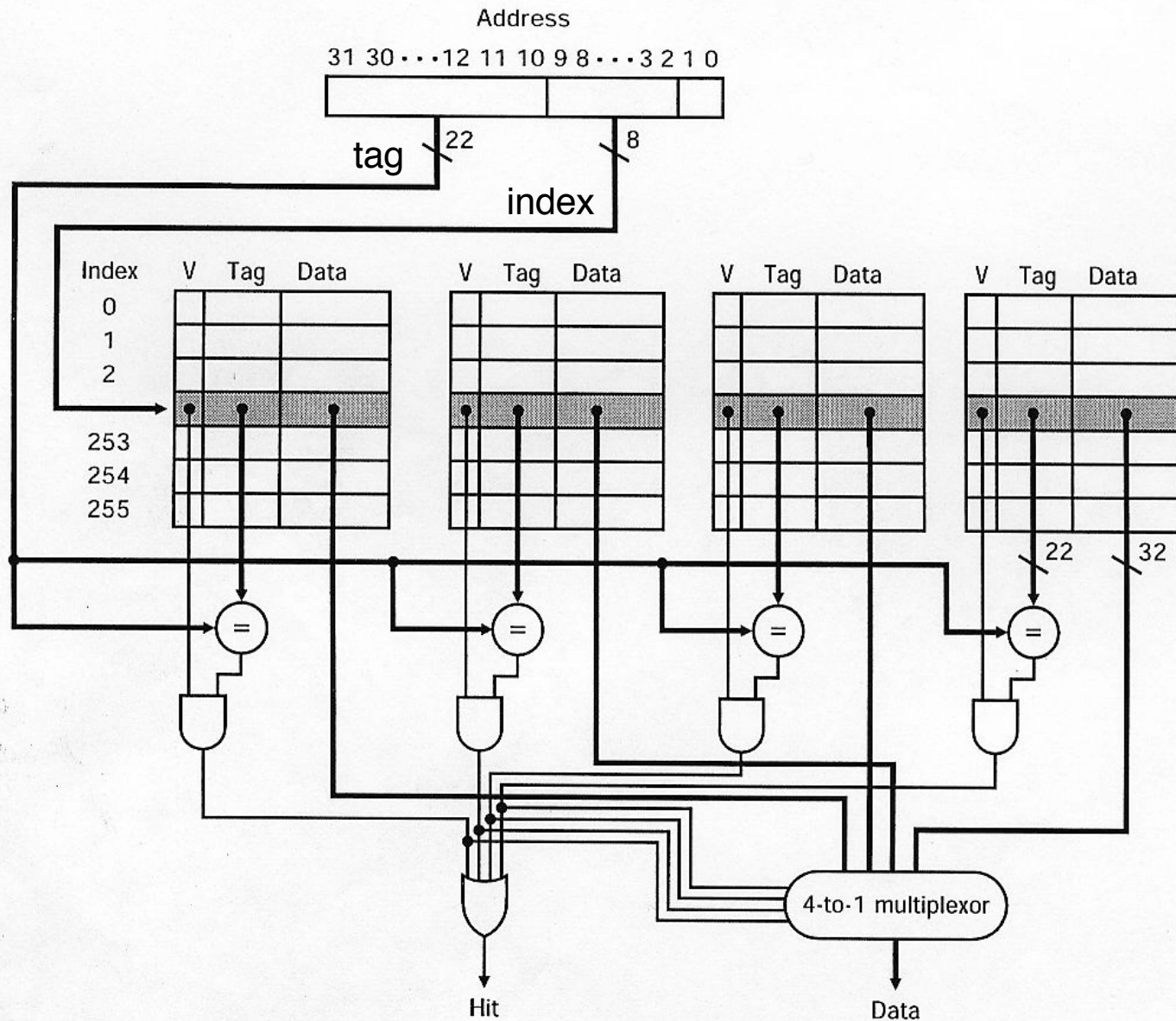  - Finally, use the offset field as usual to find the desired data within the block.

# N-Way Set Associative Cache (3/3)

- **What's so great about this?**
  - **even a 2-way set assoc cache avoids a lot of conflict misses**
  - **hardware cost isn't that bad: only need N comparators**

- **In fact, for a cache with M blocks,**
  - **it's Direct-Mapped if it's 1-way set assoc**
  - **it's Fully Assoc if it's M-way set assoc**
  - **so these two are just special cases of the more general set associative design**

# 4-Way Set Associative Cache Circuit

# Peer Instruction

A. **Mem hierarchies were invented before 1950. (UNIVAC I wasn't delivered 'til 1951)**

B. **If you know your computer's cache size, you can often make your code run faster.**

C. **Memory hierarchies take advantage of spatial locality by keeping the most recent data items closer to the processor.**

| | ABC |
|---|---|
| 0: | FFF |
| 1: | FFT |
| 2: | FTF |
| 3: | FTT |
| 4: | TFF |
| 5: | TFT |
| 6: | TTF |
| 7: | TTT |

# Peer Instruction Answer

**A.** **"We are…forced to recognize the possibility of constructing a hierarchy of memories, each of which has greater capacity than the preceding but which is less accessible." – von Neumann, 1946**

**B.** **Certainly! That's call "tuning"**

**C.** **"Most Recent" items ⇒ Temporal locality**

**A.** **Mem hierarchies were invented before 1950. (UNIVAC I wasn't delivered 'til 1951)**

**B.** **If you know your computer's cache size, you can often make your code run faster.**

**C.** **Memory hierarchies take advantage of spatial locality by keeping the most recent data items closer to the processor.**

|     | ABC |
| --- | --- |
| 0: | FFF |
| 1: | FFT |
| 2: | FTF |
| 3: | FTT |
| 4: | TFF |
| 5: | TFT |
| **6:** | **TTF** |
| 7: | TTT |

# And in Conclusion…

- **We would like to have the capacity of disk at the speed of the processor: unfortunately this is not feasible.**

- **So we create a memory hierarchy:**
  - **each successively lower level contains "most used" data from next higher level**
  - **exploits <span style="color:red">temporal & spatial locality</span>**
  - **do the common case fast, worry less about the exceptions (design principle of MIPS)**

- **Locality of reference is a Big Idea**