

CS 110
Computer Architecture
Lecture 14:
Caches Part 1

Instructor:
Sören Schwertfeger

<http://shitech.org/courses/ca/>

School of Information Science and Technology SIST

ShanghaiTech University

Slides based on UC Berkley's CS61C

New-School Machine Structures (It's a bit more complicated!)

Software

Hardware

- Parallel Requests
Assigned to computer
e.g., Search "Katz"
- Parallel Threads
Assigned to core
e.g., Lookup, Ads
- Parallel Instructions
>1 instruction @ one time
e.g., 5 pipelined instructions
- Parallel Data
>1 data item @ one time
e.g., Add of 4 pairs of words
- Hardware descriptions
All gates @ one time
- Programming Languages

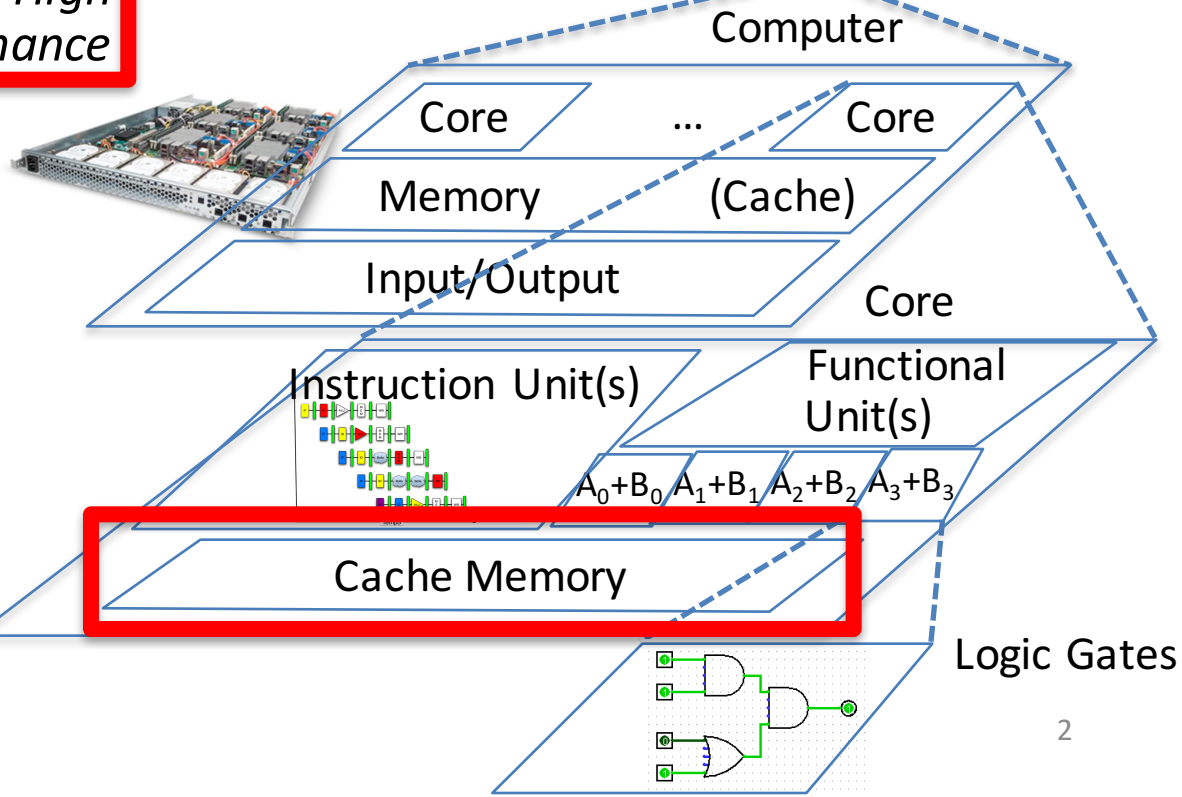
Harness
Parallelism &
**Achieve High
Performance**

Warehouse
Scale
Computer

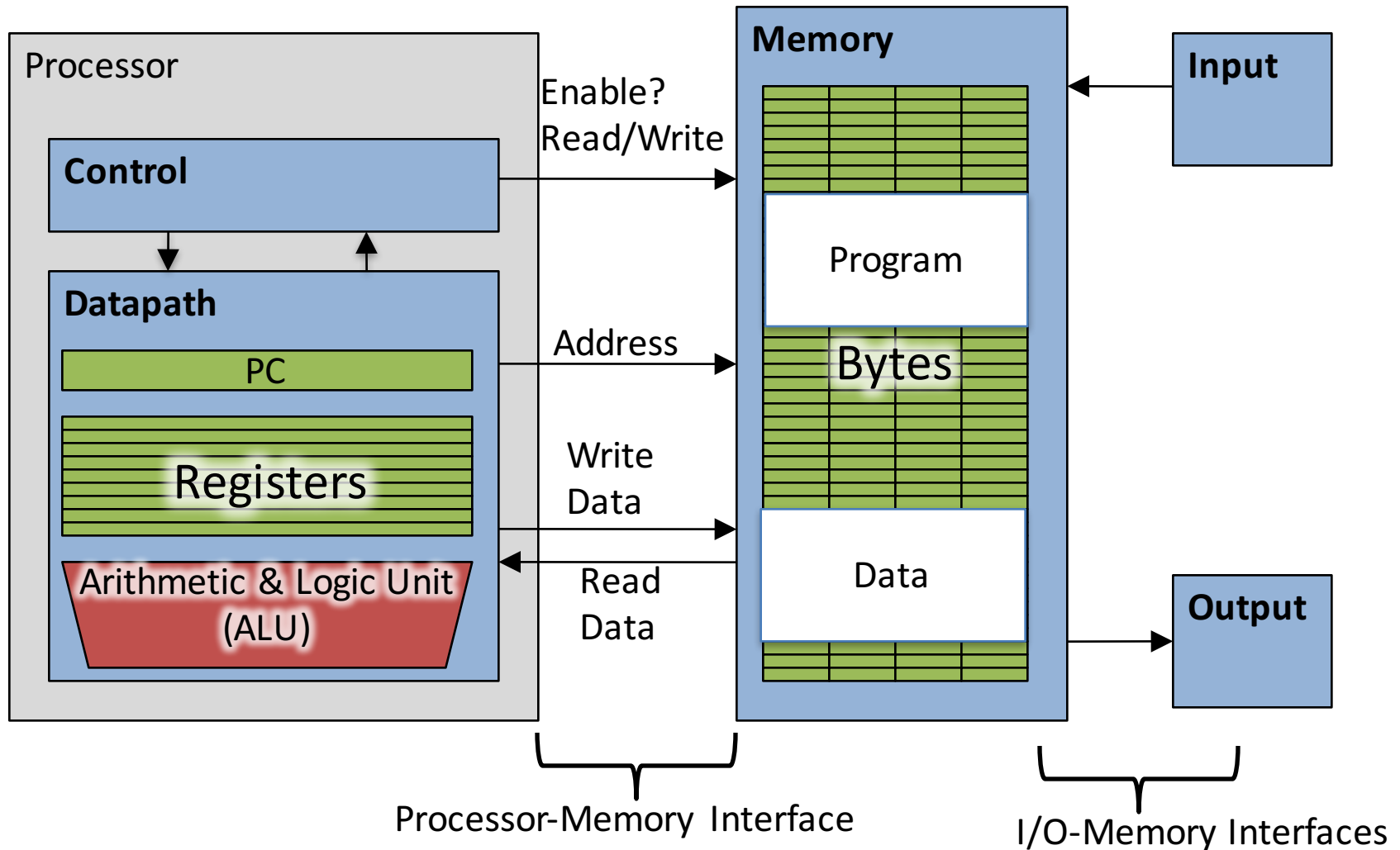
How do
we know?



Smart
Phone



Components of a Computer



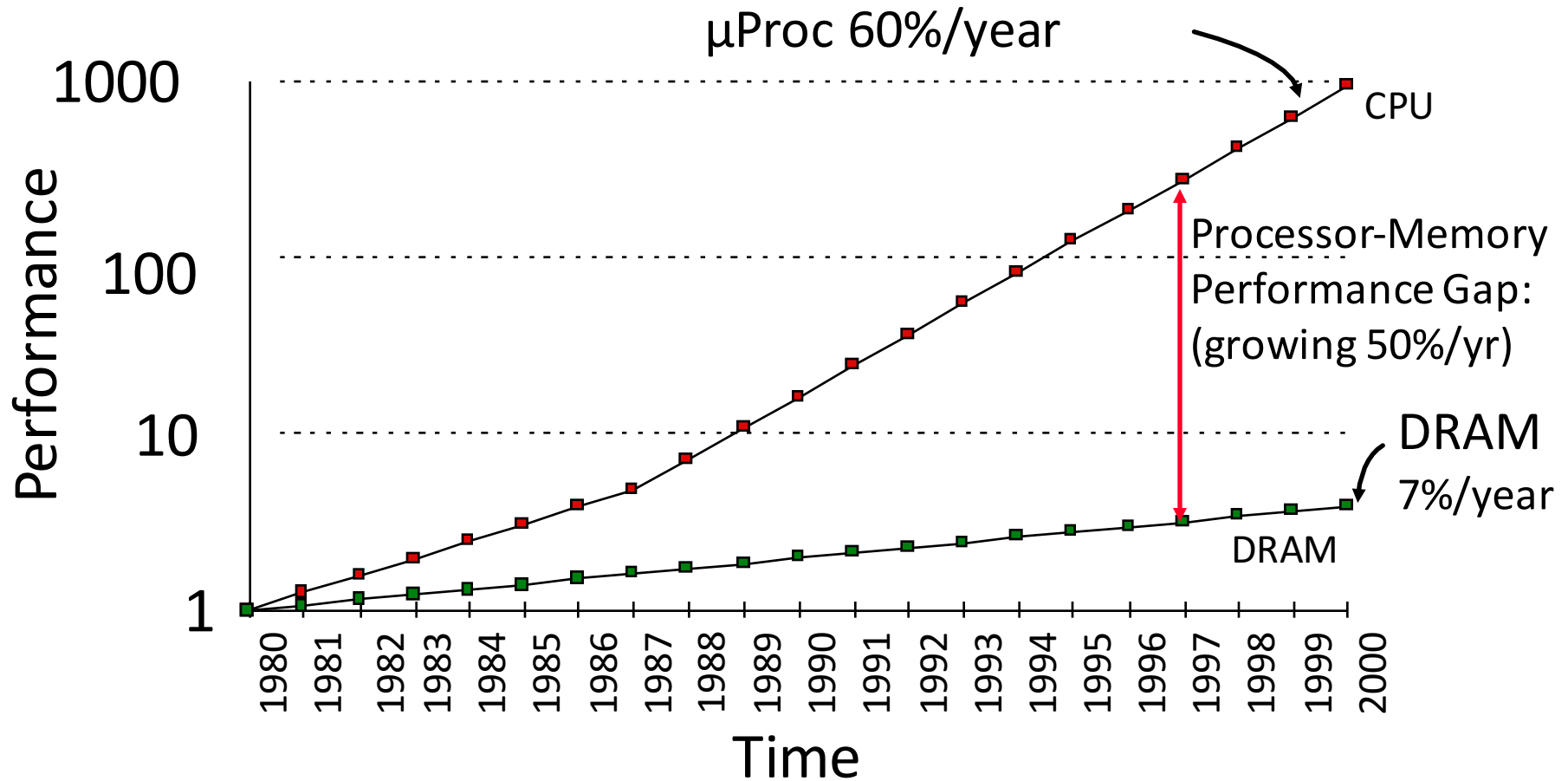
Problem: Large memories slow?

Library Analogy

- Finding a book in a large library takes time
 - Takes time to search a large card catalog – (mapping title/author to index number)
 - Round-trip time to walk to the stacks and retrieve the desired book.
- Larger libraries makes both delays worse
- Electronic memories have the same issue, *plus* the technologies that we use to store an individual bit get slower as we increase density (SRAM versus DRAM versus Magnetic Disk)

However what we want is a large yet fast memory!

Processor-DRAM Gap (latency)

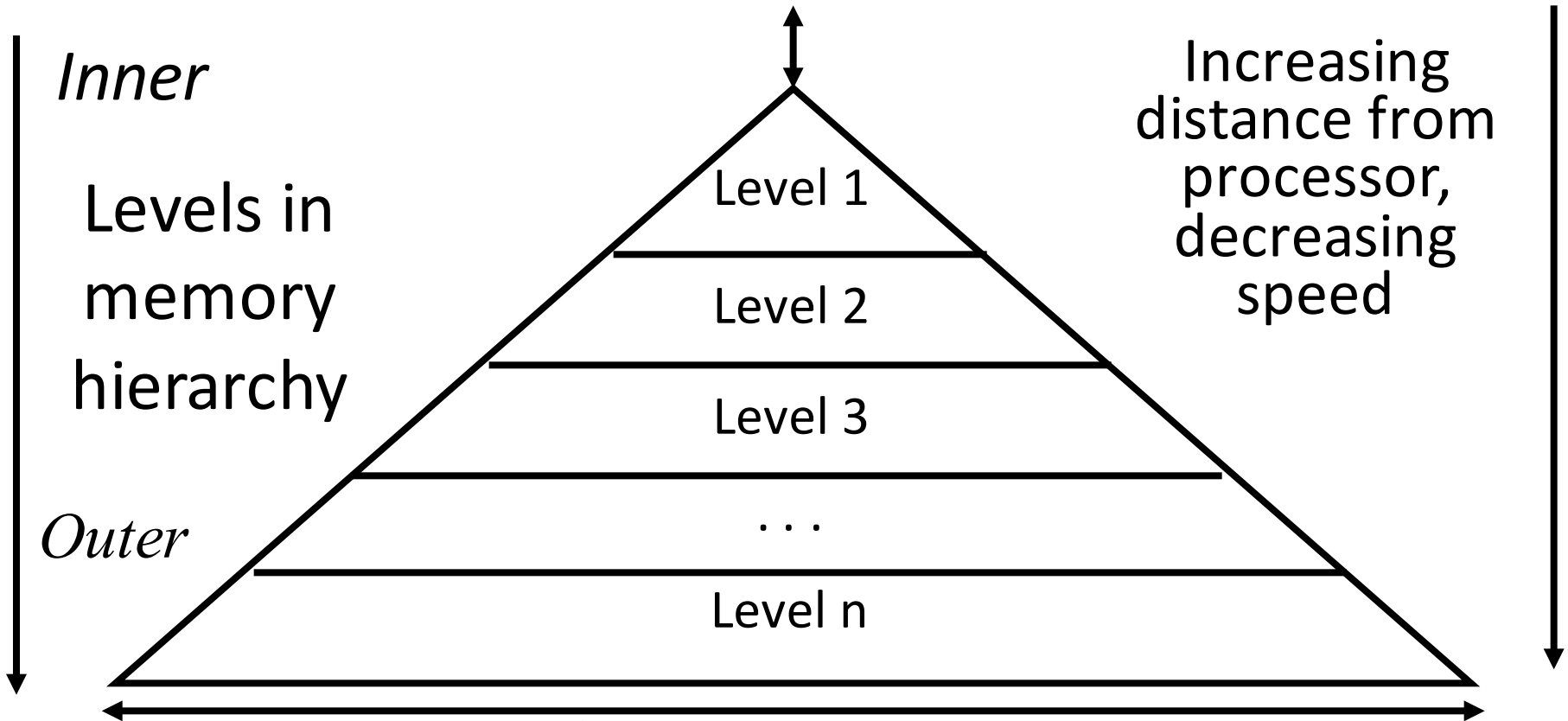


1980 microprocessor executes ~one instruction in same time as DRAM access
2015 microprocessor executes ~1000 instructions in same time as DRAM access

Slow DRAM access could have disastrous impact on CPU performance!

Big Idea: Memory Hierarchy

Processor

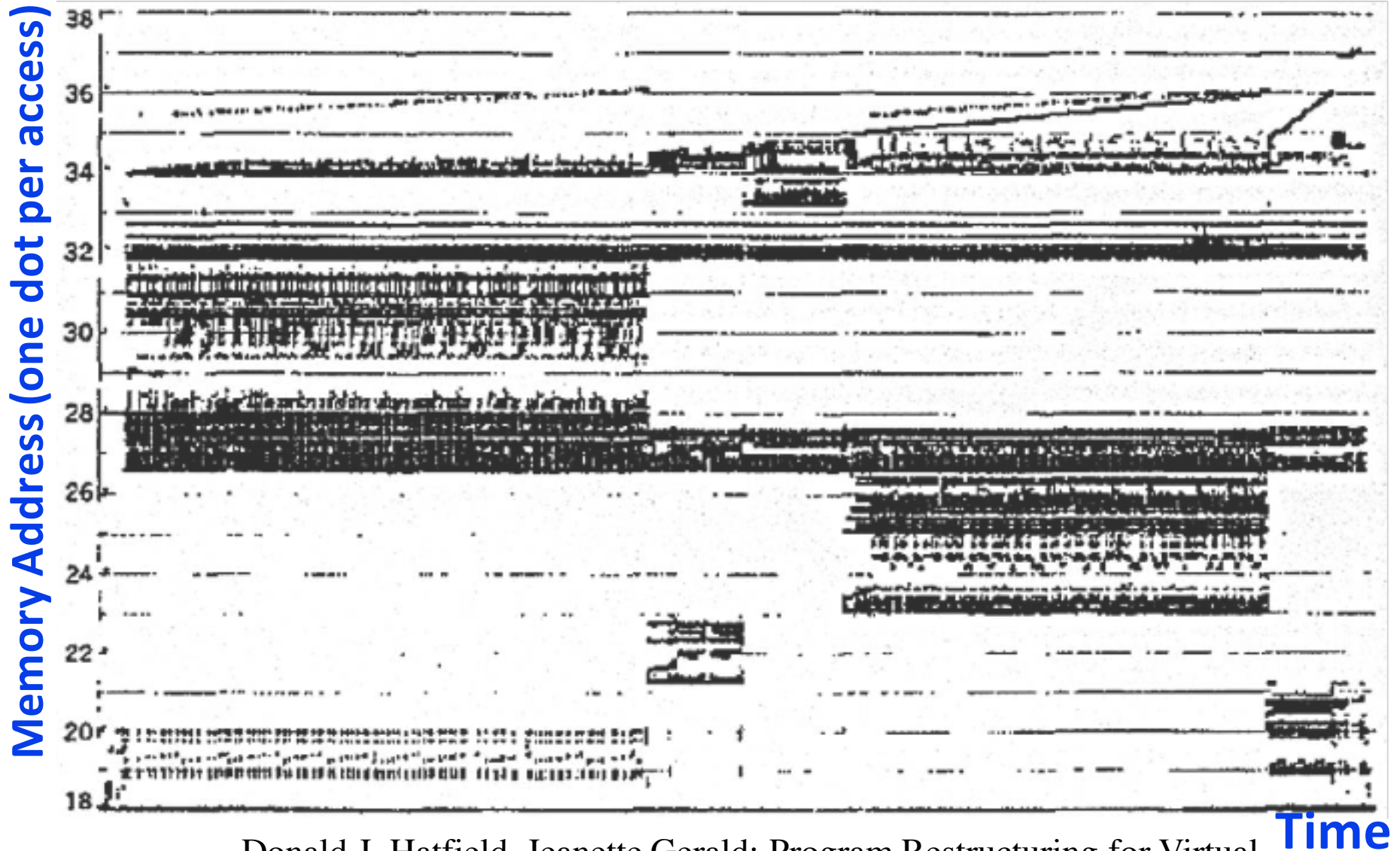


As we move to outer levels the latency goes up and price per bit goes down. Why?

What to do: Library Analogy

- Want to write a report using library books
- Go to library, look up relevant books, fetch from stacks, and place on desk in library
- If need more, check them out and keep on desk
 - But don't return earlier books since might need them
- You hope this collection of ~10 books on desk enough to write report, despite 10 being only a tiny fraction of books available

Real Memory Reference Patterns



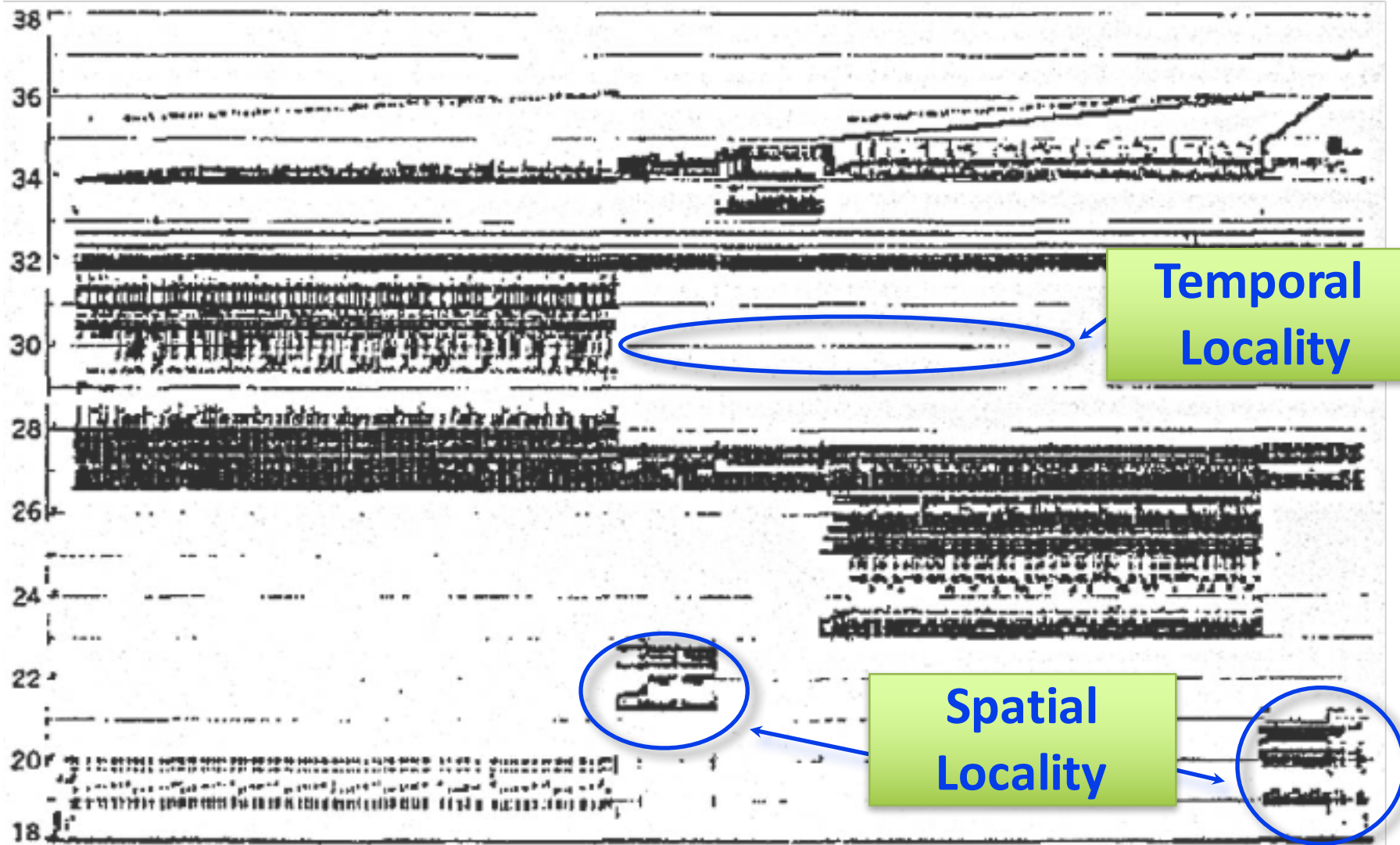
Donald J. Hatfield, Jeanette Gerald: Program Restructuring for Virtual Memory. IBM Systems Journal 10(3): 168-192 (1971)

Big Idea: Locality

- *Temporal Locality* (locality in time)
 - Go back to same book on desktop multiple times
 - If a memory location is referenced, then it will tend to be referenced again soon
- *Spatial Locality* (locality in space)
 - When go to book shelf, pick up multiple books on J.D. Salinger since library stores related books together
 - If a memory location is referenced, the locations with nearby addresses will tend to be referenced soon

Memory Reference Patterns

Memory Address (one dot per access)



Temporal
Locality

Spatial
Locality

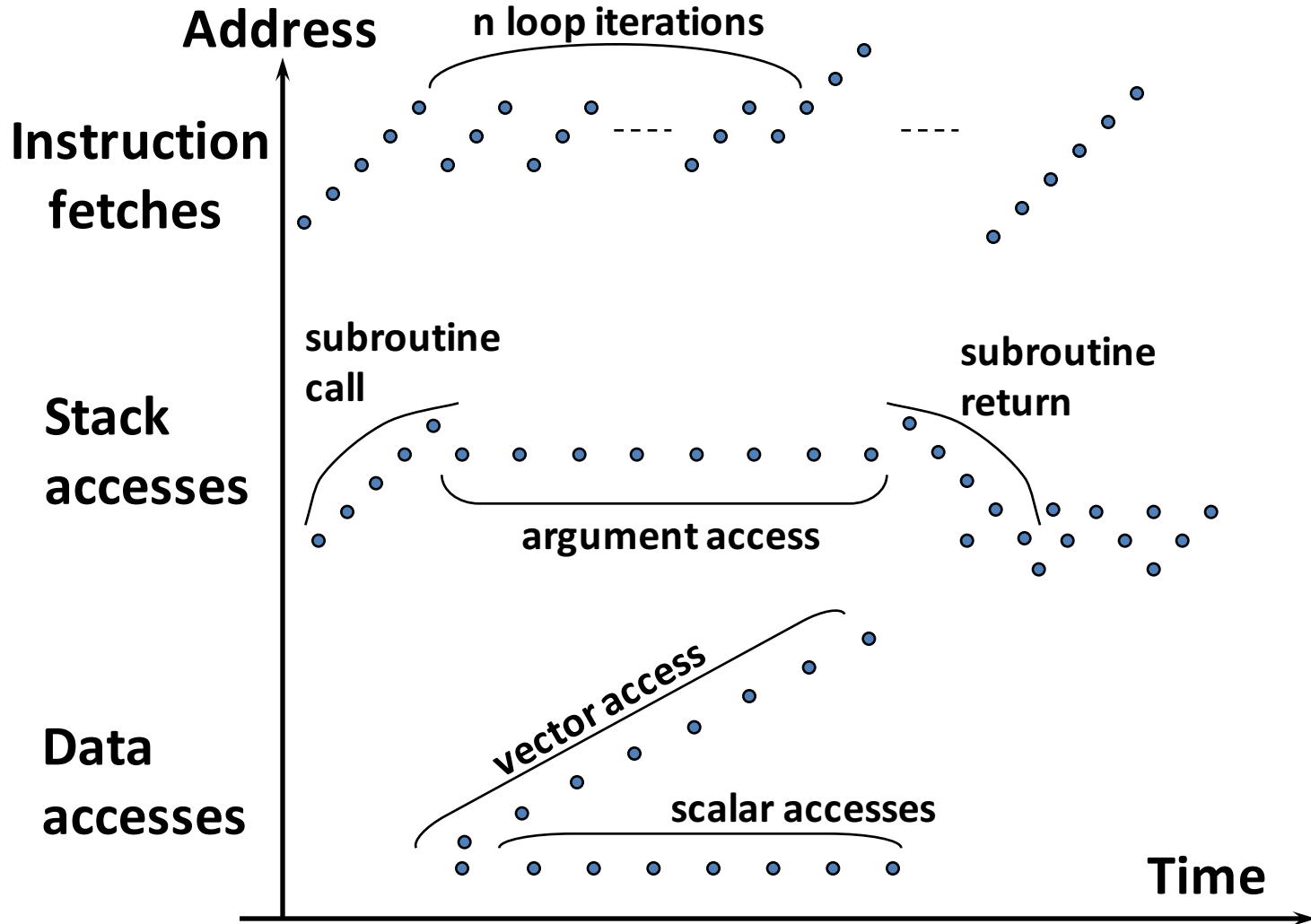
Time

Donald J. Hatfield, Jeanette Gerald: Program Restructuring for Virtual Memory. IBM Systems Journal 10(3): 168-192 (1971)

Principle of Locality

- *Principle of Locality*: Programs access small portion of address space at any instant of time (spatial locality) and repeatedly access that portion (temporal locality)
- What program structures lead to **temporal** and **spatial locality** in instruction accesses?
- In data accesses?

Memory Reference Patterns



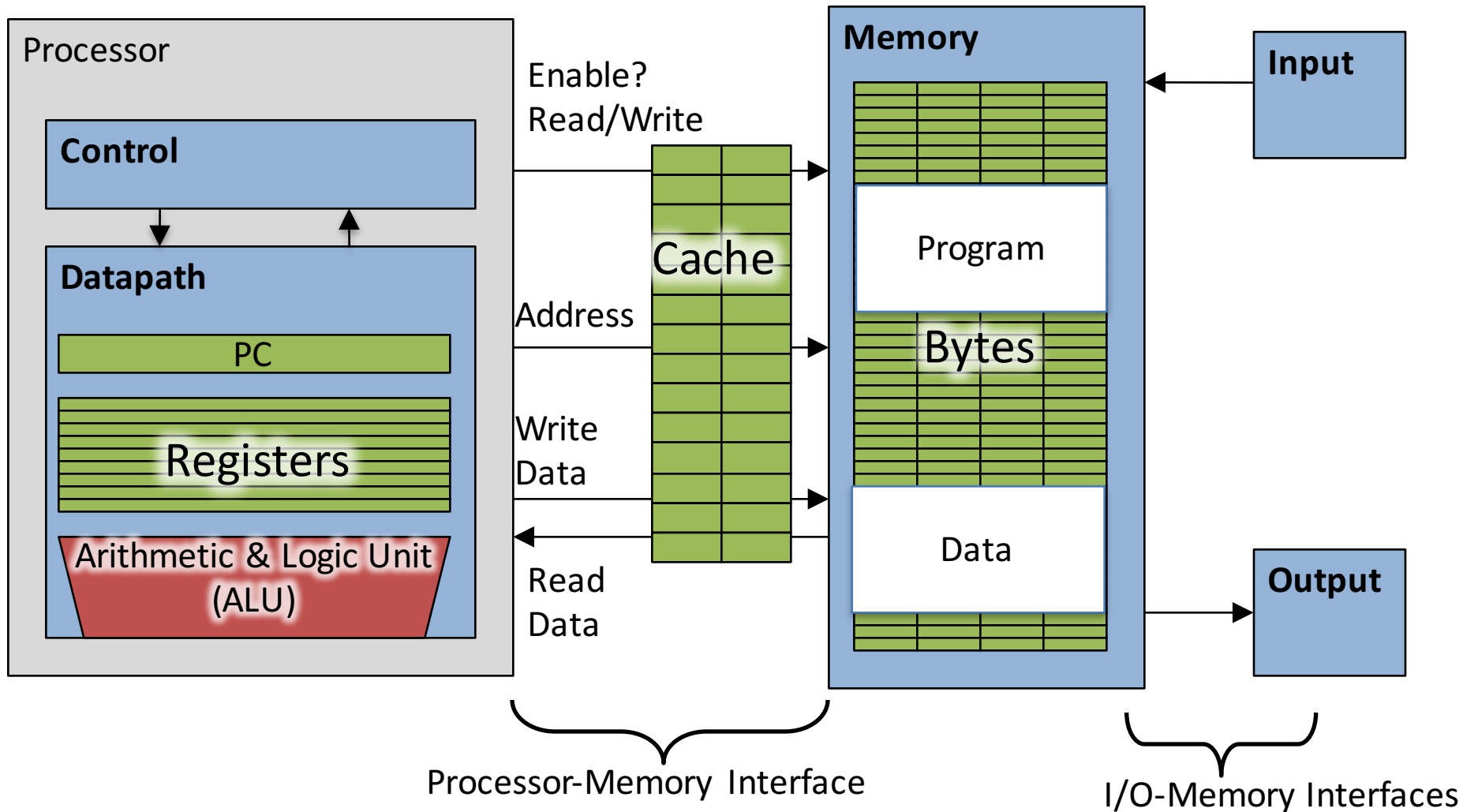
Cache Philosophy

- Programmer-invisible hardware mechanism to give illusion of speed of fastest memory with size of largest memory
 - Works fine even if programmer has no idea what a cache is
 - However, performance-oriented programmers today sometimes “reverse engineer” cache design to design data structures to match cache

Memory Access without Cache

- Load word instruction: `lw $t0, 0($t1)`
- `$t1` contains 1022_{ten} , `Memory[1022] = 99`
 1. Processor issues address 1022_{ten} to Memory
 2. Memory reads word at address 1022_{ten} (99)
 3. Memory sends 99 to Processor
 4. Processor loads 99 into register `$t0`

Adding Cache to Computer

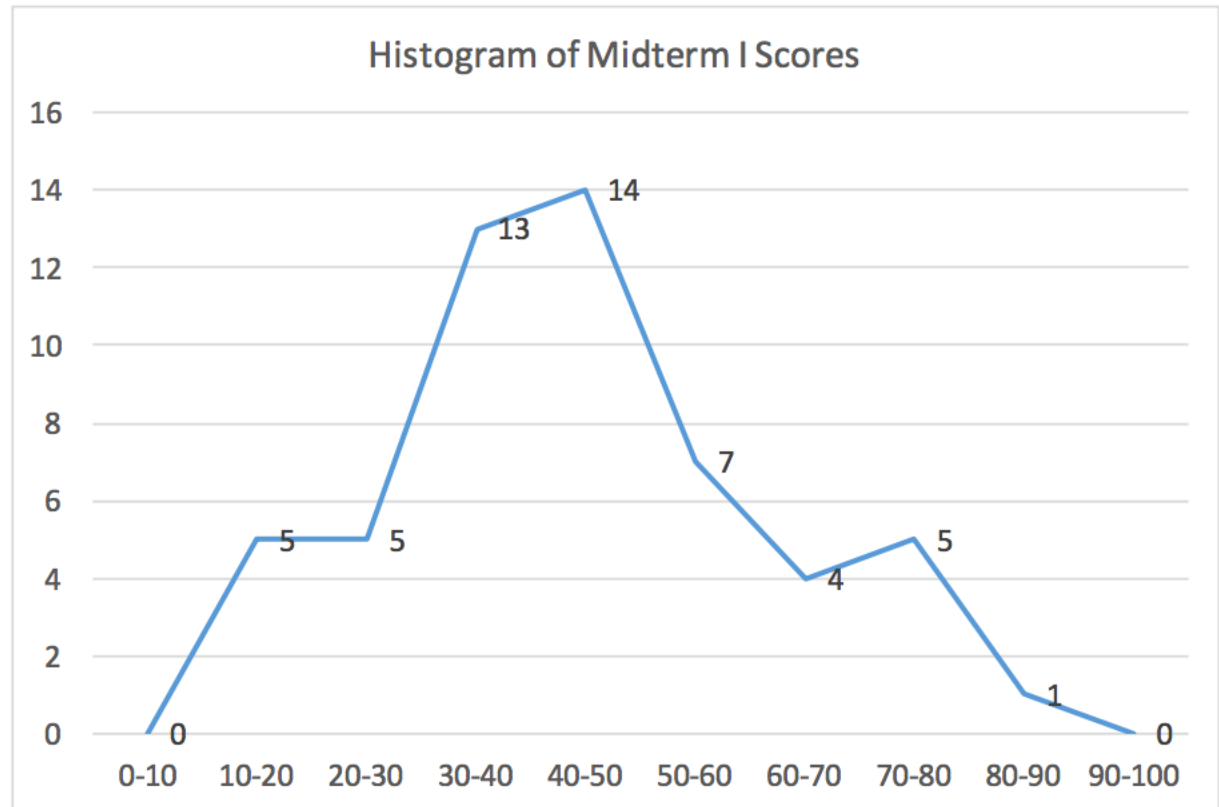


Memory Access with Cache

- Load word instruction: `lw $t0, 0($t1)`
- `$t1` contains 1022_{ten} , `Memory[1022] = 99`
- With cache: Processor issues address 1022_{ten} to Cache
 1. Cache checks to see if has copy of data at address 1022_{ten}
 - 2a. If finds a match (Hit): cache reads 99, sends to processor
 - 2b. No match (Miss): cache sends address 1022 to Memory
 - I. Memory reads 99 at address 1022_{ten}
 - II. Memory sends 99 to Cache
 - III. Cache replaces word with new 99
 - IV. Cache sends 99 to processor
 2. Processor loads 99 into register `$t0`

Administrivia

- Midterm 1
- Go through all questions at today's discussion
- Grading for the course will be relative, not absolute



- Postpone HW5 or Project 1.2?

Cache “Tags”

- Need way to tell if have copy of location in memory so that can decide on hit or miss
- On cache miss, put memory address of block in “tag address” of cache block

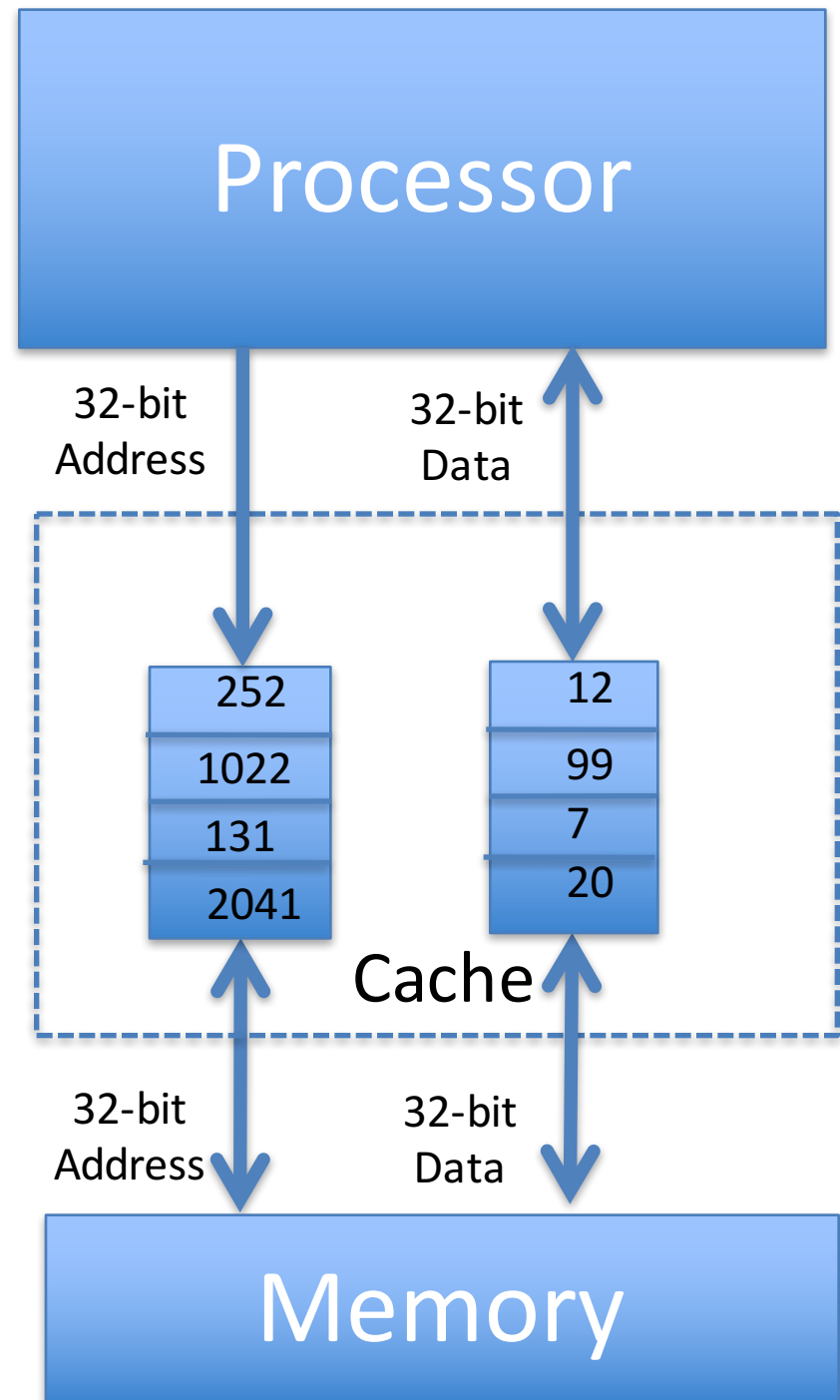
1022 placed in tag next to data from memory (99)

Tag	Data
252	12
1022	99
131	7
2041	20

From earlier instructions

Anatomy of a 16 Byte Cache, 4 Byte Block

- Operations:
 1. Cache Hit
 2. Cache Miss
 3. Refill cache from memory
- Cache needs Address Tags to decide if Processor Address is a Cache Hit or Cache Miss
 - Compares all 4 tags



Cache Replacement

- Suppose processor now requests location 511, which contains 11?
- Doesn't match any cache block, so must "evict" one resident block to make room
 - Which block to evict?
- Replace "victim" with new memory block at address 511

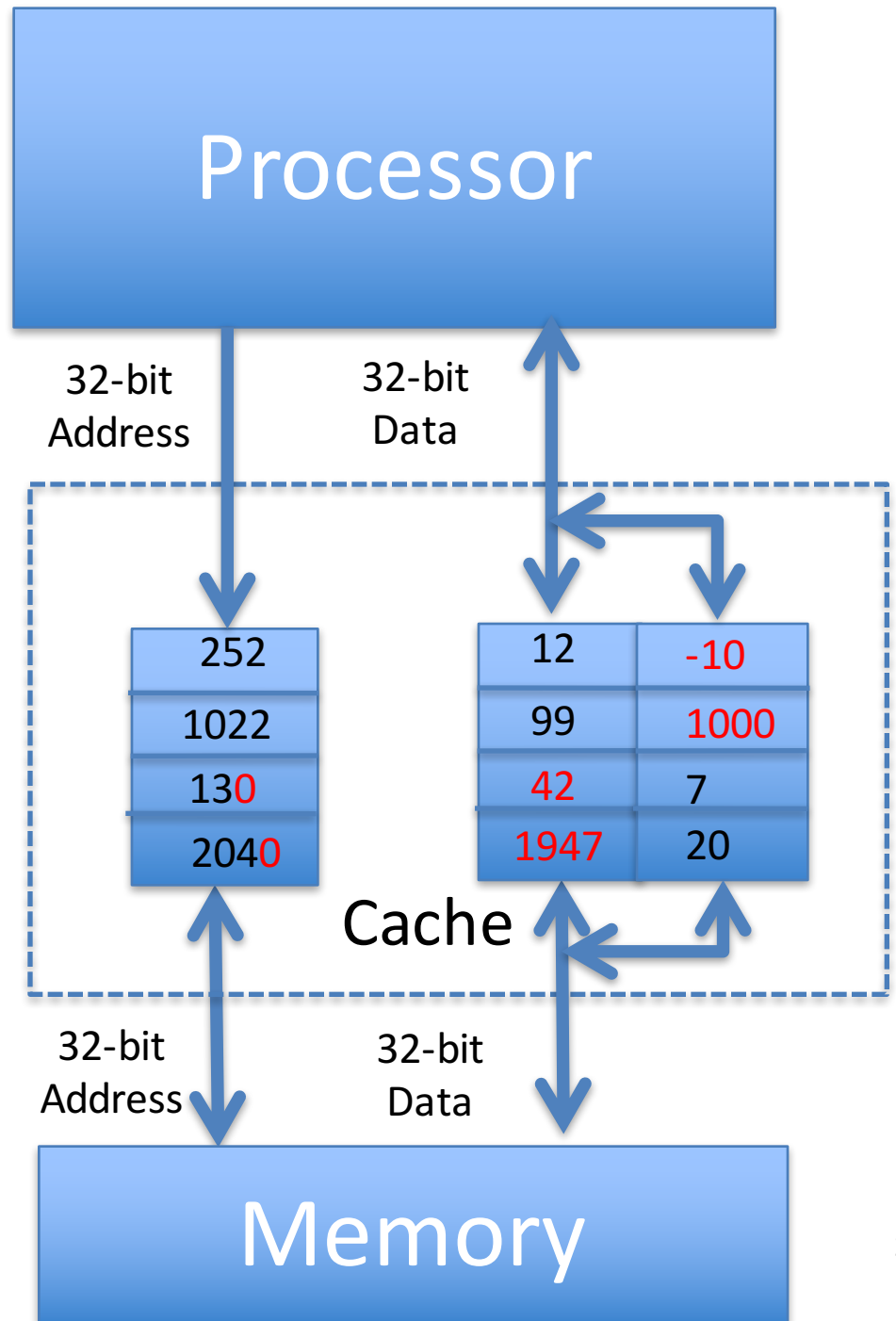
Tag	Data
252	12
1022	99
511	11
2041	20

Block Must be Aligned in Memory

- Word blocks are aligned, so binary address of all words in cache always ends in 00_{two}
 - How to take advantage of this to save hardware and energy?
 - Don't need to compare last 2 bits of 32-bit byte address (comparator can be narrower)
- => Don't need to store last 2 bits of 32-bit byte address in Cache Tag (Tag can be narrower)

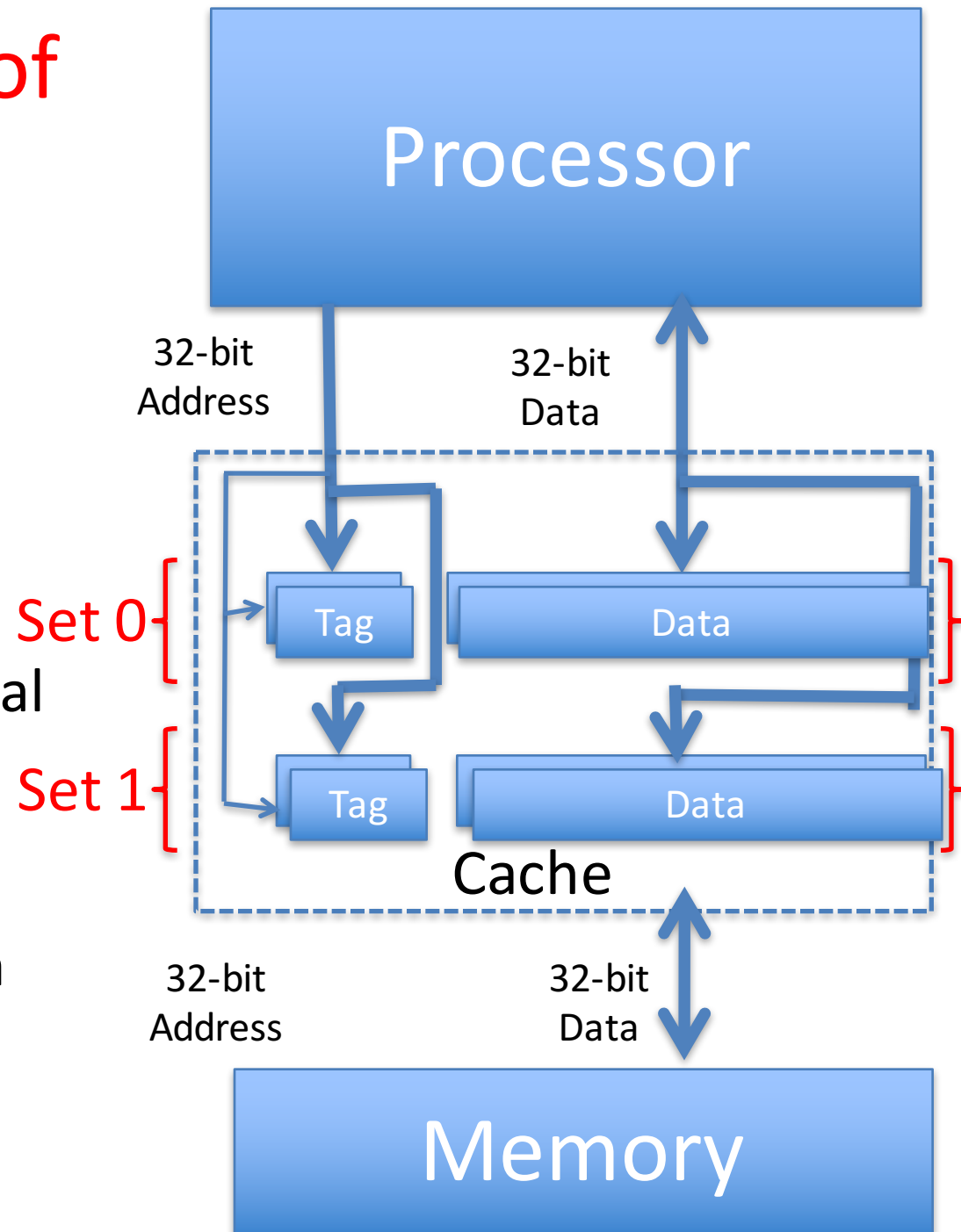
Anatomy of a 32B Cache, 8B Block

- Blocks must be aligned in pairs, otherwise could get same word twice in cache
 - Tags only have even-numbered words
 - Last 3 bits of address always 000_{two}
 - Tags, comparators can be narrower
- Can get hit for either word in block



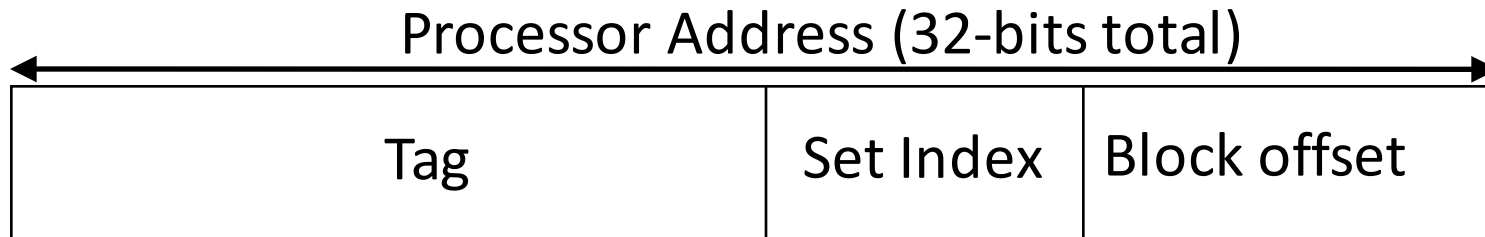
Hardware Cost of Cache

- Need to compare every tag to the Processor address
- Comparators are expensive
- Optimization: use 2 “sets” of data with a total of only 2 comparators
- 1 Address bit selects which set
- Compare only tags from selected set
- Generalize to more sets



Processor Address Fields used by Cache Controller

- **Block Offset**: Byte address within block
- **Set Index**: Selects which set
- **Tag**: Remaining portion of processor address



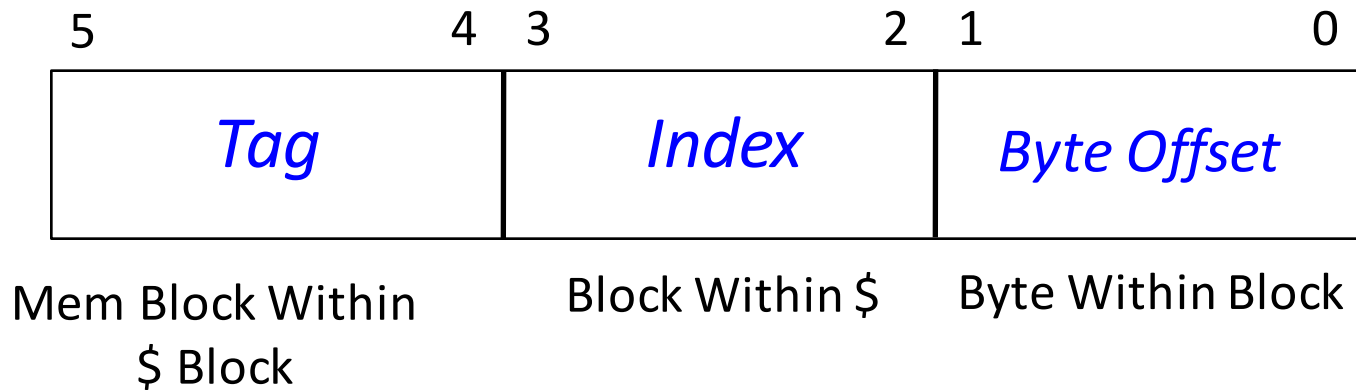
- Size of Index = \log_2 (number of sets)
- Size of Tag = Address size – Size of Index – \log_2 (number of bytes/block)

What is limit to number of sets?

- For a given total number of blocks, we can save more comparators if have more than 2 sets
- Limit: As Many Sets as Cache Blocks => only one block per set – only needs one comparator!
- Called “Direct-Mapped” Design



Direct Mapped Cache Ex: Mapping a 6-bit Memory Address

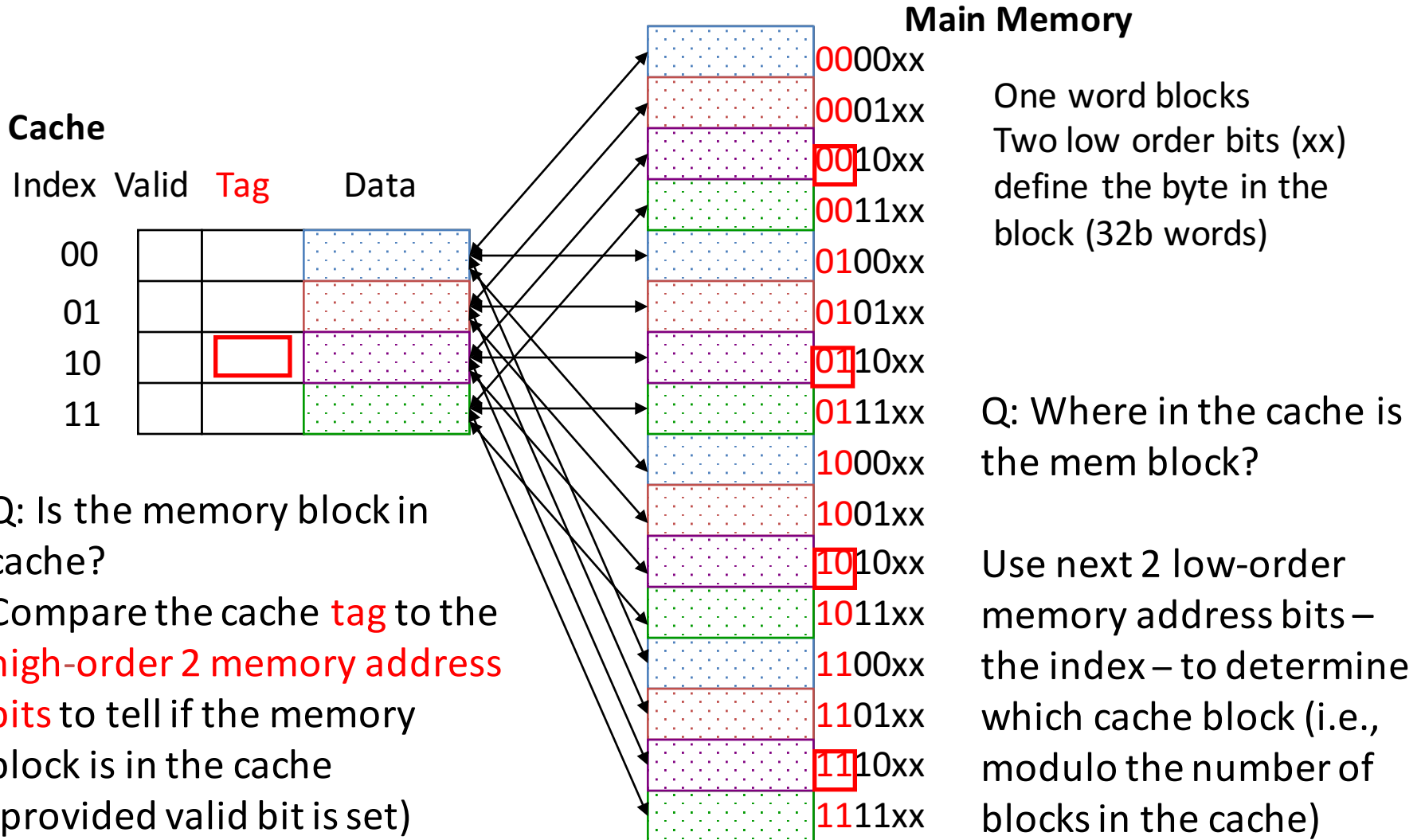


- In example, block size is 4 bytes/1 word
- Memory and cache blocks always the same size, unit of transfer between memory and cache
- # Memory blocks >> # Cache blocks
 - 16 Memory blocks = 16 words = 64 bytes => 6 bits to address all bytes
 - 4 Cache blocks, 4 bytes (1 word) per block
 - 4 Memory blocks map to each cache block
- Memory block to cache block, aka *index*: middle two bits
- Which memory block is in a given cache block, aka *tag*: top two bits

One More Detail: Valid Bit

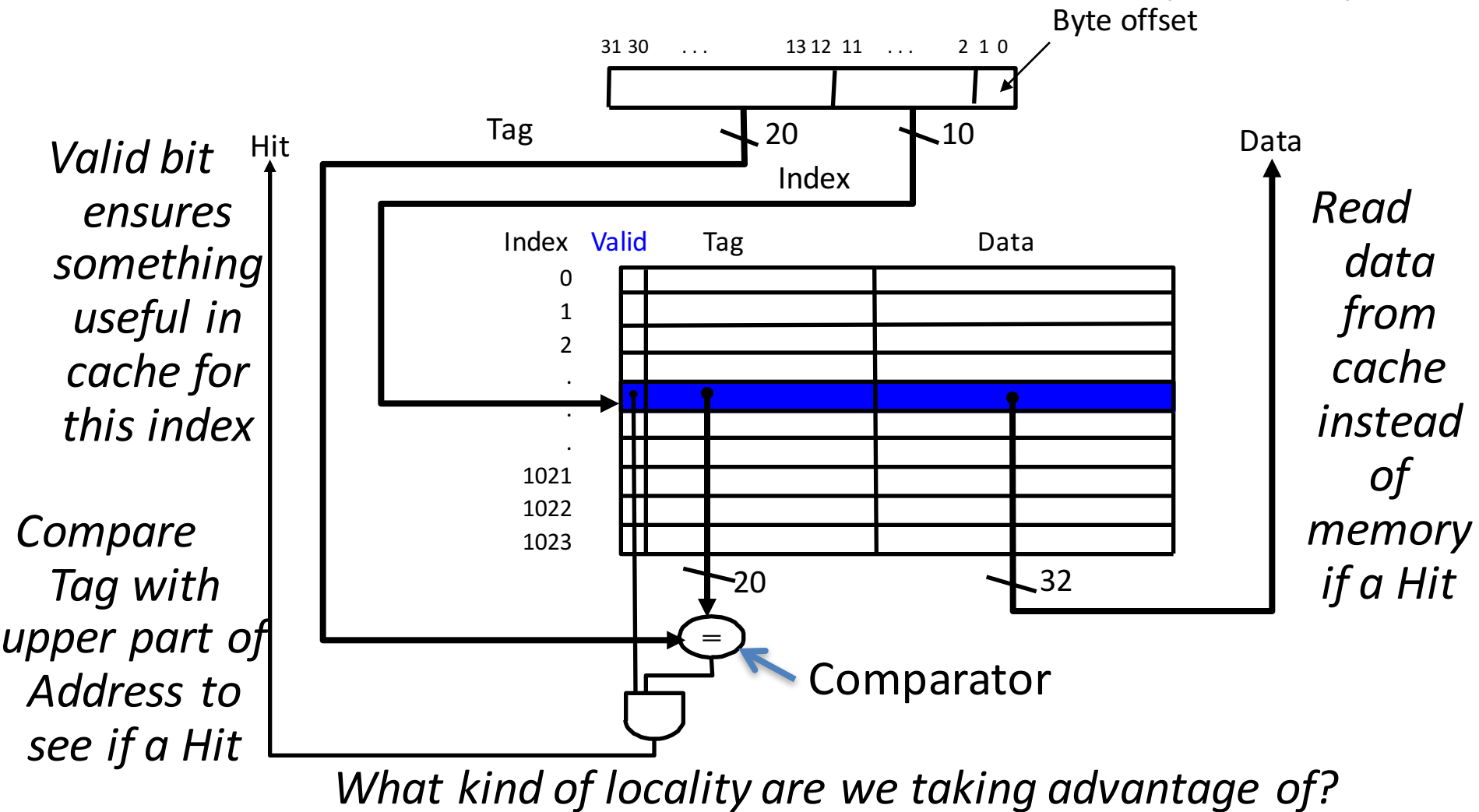
- When start a new program, cache does not have valid information for this program
- Need an indicator whether this tag entry is valid for this program
- Add a “valid bit” to the cache tag entry
 - 0 => cache miss, even if by chance, address = tag
 - 1 => cache hit, if processor address = tag

Caching: A Simple First Example



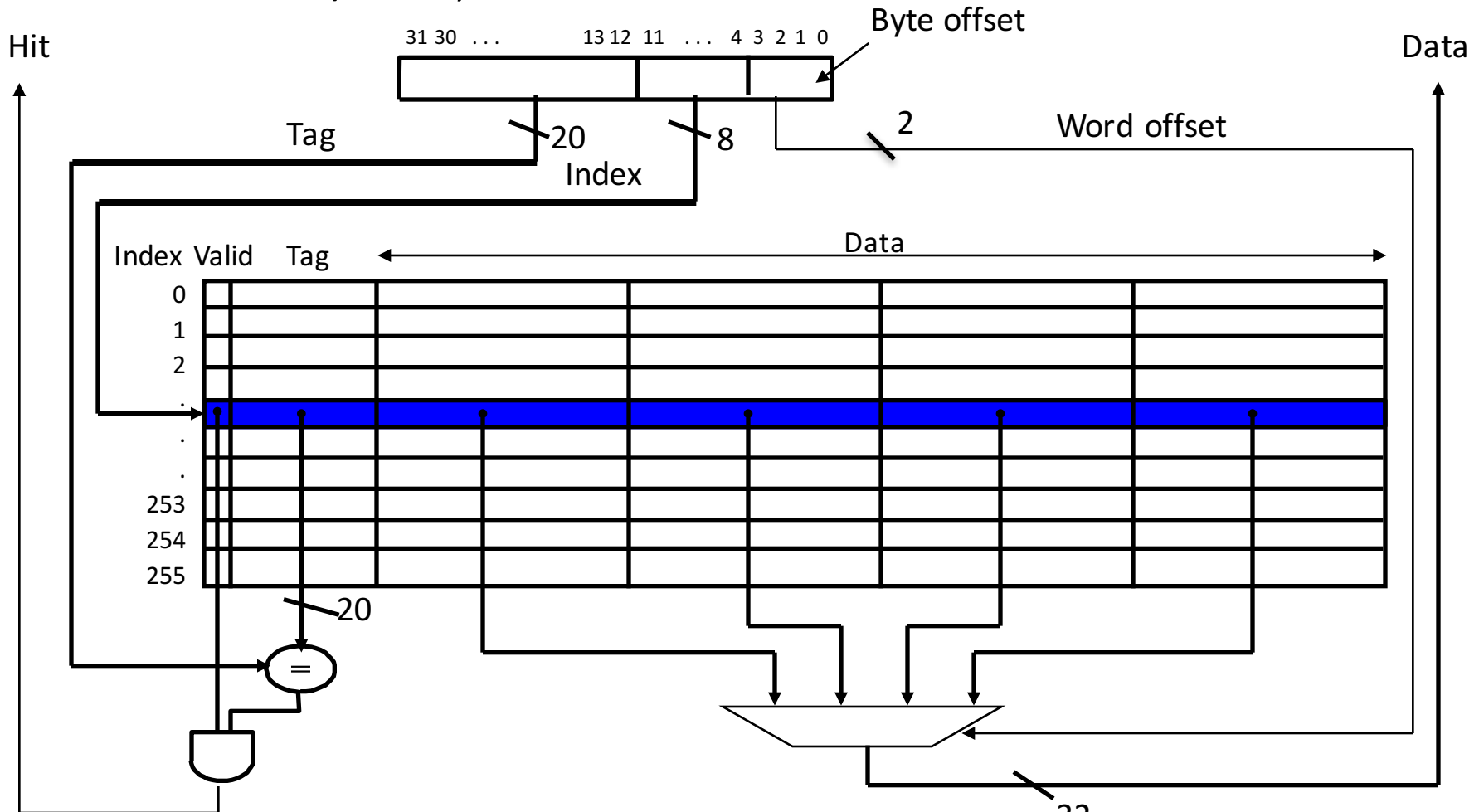
Direct-Mapped Cache Example

- One word blocks, cache size = 1K words (or 4KB)



Multiword-Block Direct-Mapped Cache

- Four words/block, cache size = 1K words



What kind of locality are we taking advantage of?

Cache Names for Each Organization

- “Fully Associative”: Block can go anywhere
 - First design in lecture
 - Note: No Index field, but 1 comparator/block
- “Direct Mapped”: Block goes one place
 - Note: Only 1 comparator
 - Number of sets = number blocks
- “N-way Set Associative”: N places for a block
 - Number of sets = number of blocks / N
 - N comparators
 - **Fully Associative: $N = \text{number of blocks}$**
 - **Direct Mapped: $N = 1$**

Range of Set-Associative Caches

- For a fixed-size cache, and a given block size, each increase by a factor of 2 in associativity doubles the number of blocks per set (i.e., the number of “ways”) and halves the number of sets –
 - decreases the size of the index by 1 bit and increases the size of the tag by 1 bit

More Associativity (more ways)



What if we can also change the block size?

Question

- For a cache with constant total capacity, if we increase the number of ways by a factor of 2, which statement is false:
- A: The number of sets could be doubled
- B: The tag width could decrease
- C: The block size could stay the same
- D: The block size could be halved
- E: Tag width must increase

Total Cash Capacity =

Associativity * # of sets * block_size

*Bytes = blocks/set * sets * Bytes/block*

$$C = N * S * B$$



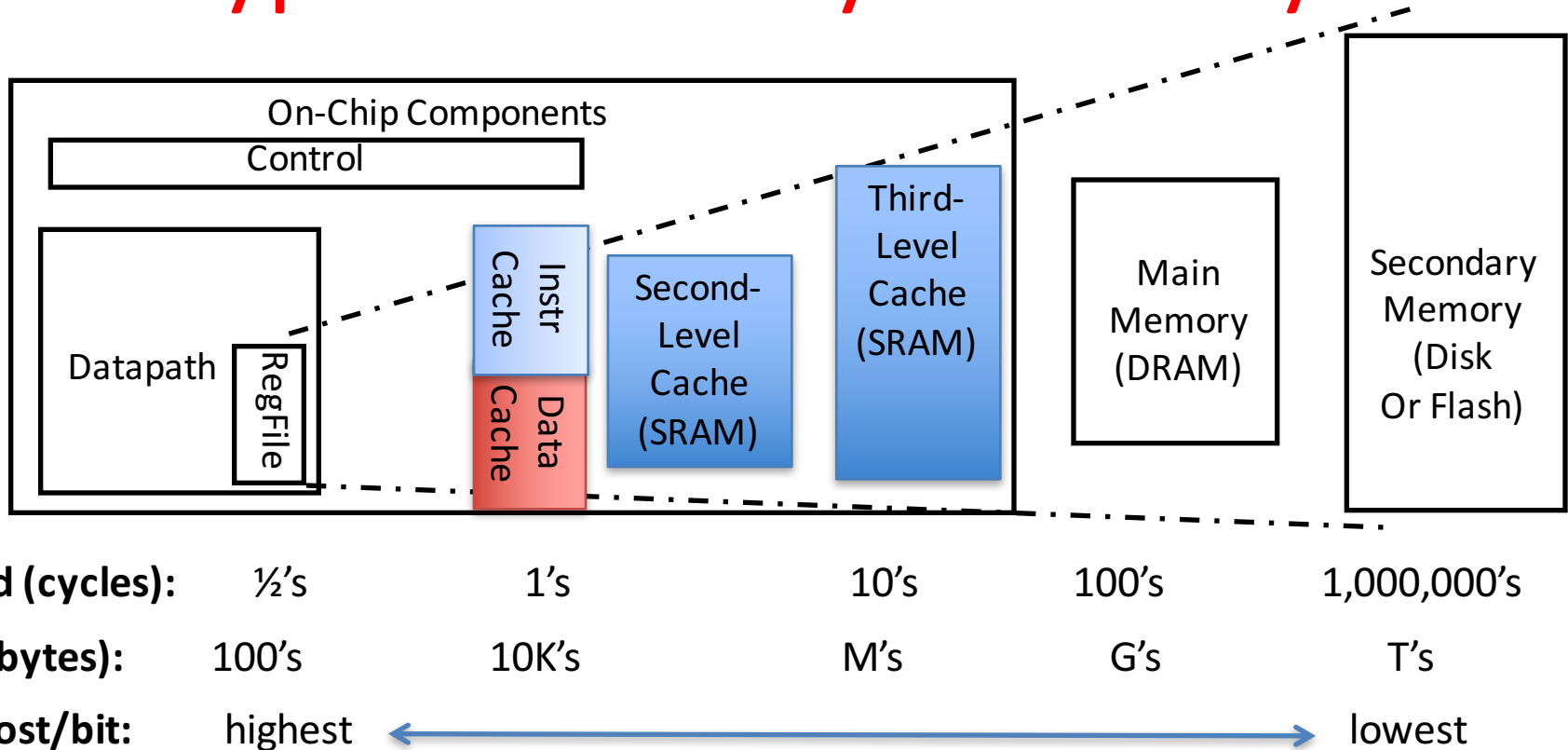
$$\begin{aligned} \text{address_size} &= \text{tag_size} + \text{index_size} + \text{offset_size} \\ &= \text{tag_size} + \log_2(S) + \log_2(B) \end{aligned}$$

Clicker Question: C remains constant, S and/or B can change such that

$$C = 2N * (SB)' \Rightarrow (SB)' = SB/2$$

$$\begin{aligned} \text{Tag_size} &= \text{address_size} - (\log_2(S) + \log_2(B)) = \text{address_size} - \log_2(SB) \\ &= \text{address_size} - (\log_2(SB) - 1) \end{aligned}$$

Typical Memory Hierarchy



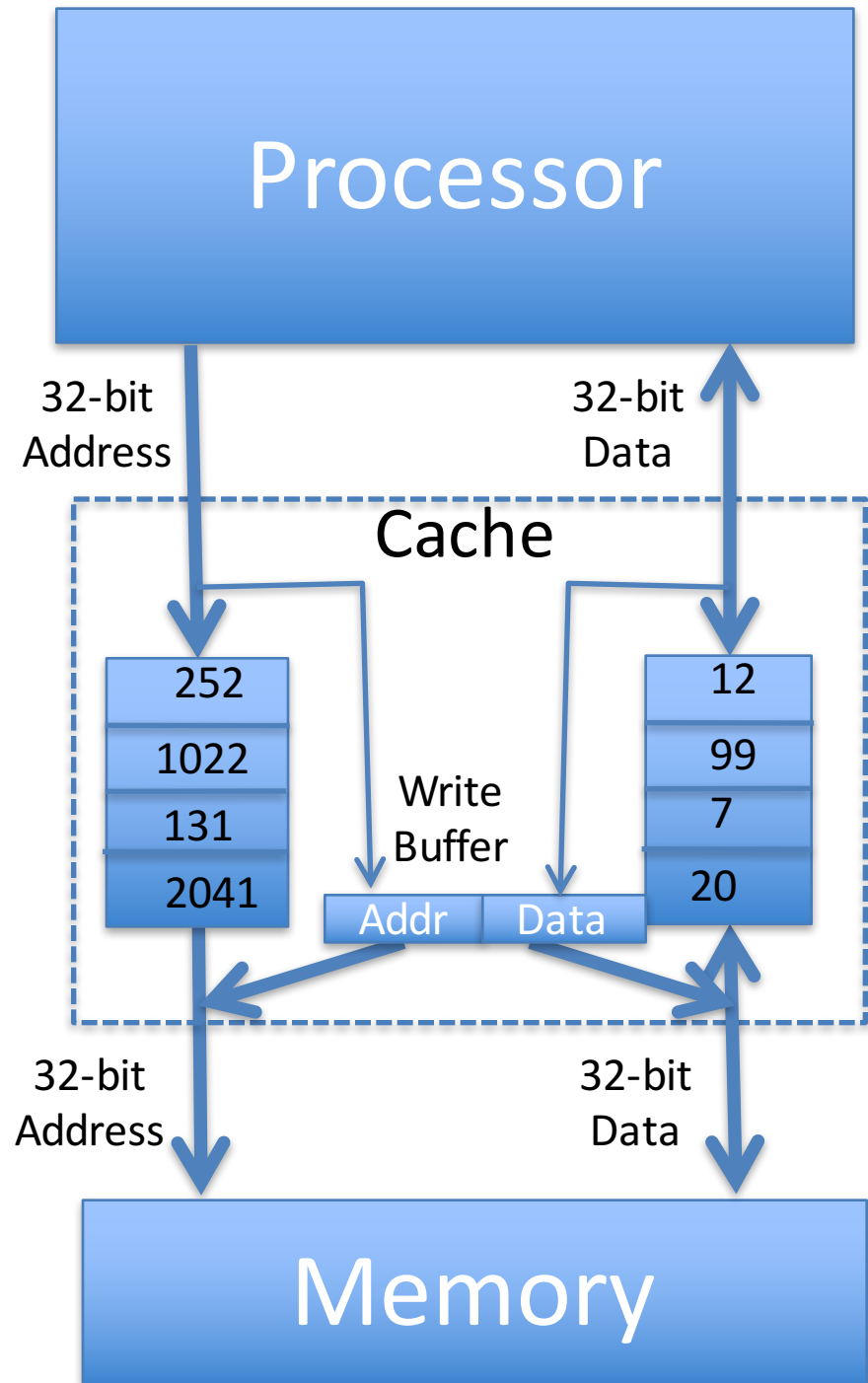
- **Principle of locality + memory hierarchy** presents programmer with \approx as much memory as is available in the *cheapest* technology at the \approx speed offered by the *fastest* technology

Handling Stores with Write-Through

- Store instructions write to memory, changing values
 - Need to make sure cache and memory have same values on writes: 2 policies
- 1) **Write-Through Policy**: write cache and write *through* the cache to memory
- Every write eventually gets to memory
 - Too slow, so include Write Buffer to allow processor to continue once data in Buffer
 - Buffer updates memory in parallel to processor

Write-Through Cache

- Write both values in cache and in memory
- Write buffer stops CPU from stalling if memory cannot keep up
- Write buffer may have multiple entries to absorb bursts of writes
- What if store misses in cache?

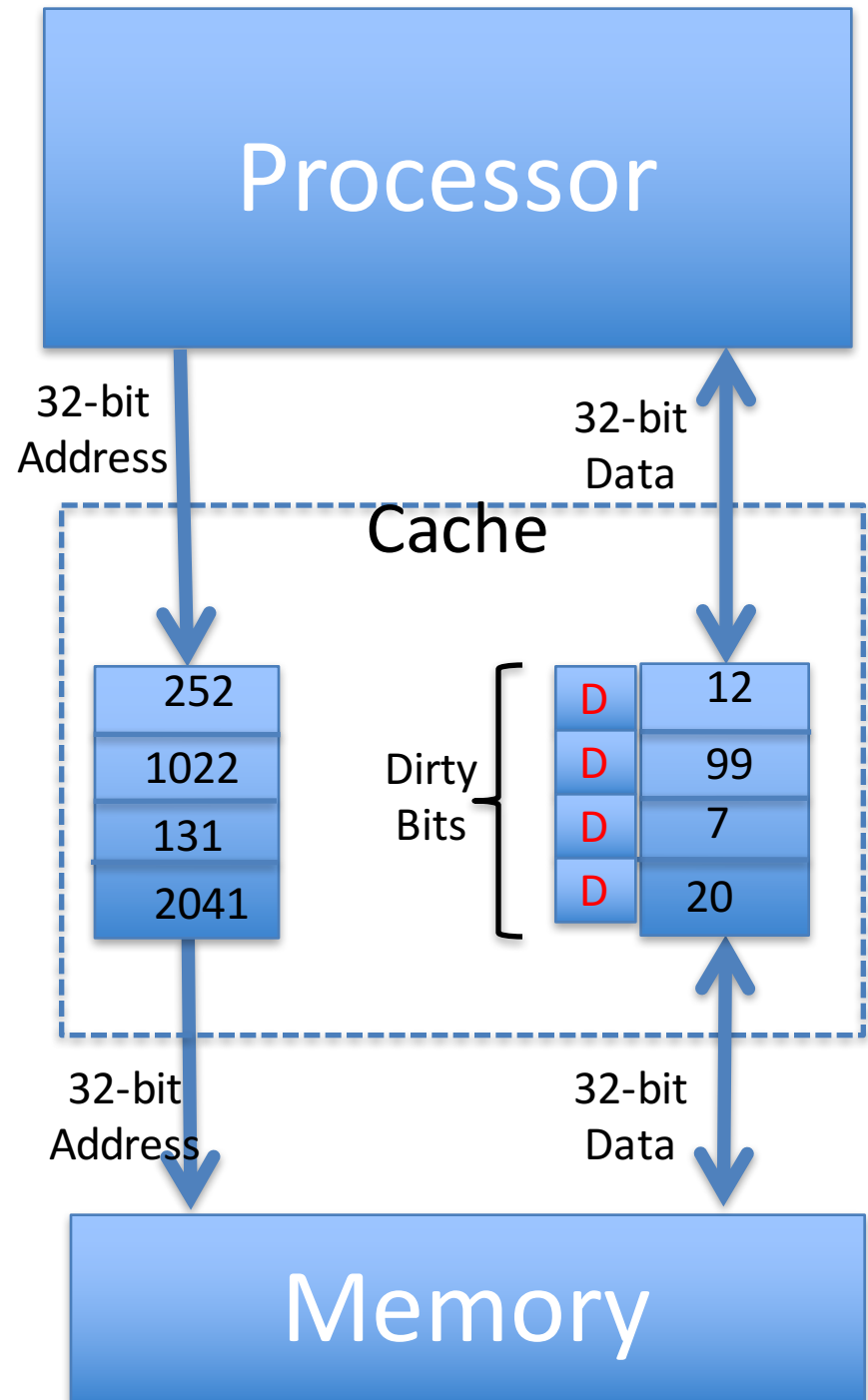


Handling Stores with Write-Back

- 2) **Write-Back Policy**: write only to cache and then write cache block *back* to memory when evict block from cache
- Writes collected in cache, only single write to memory per block
 - Include bit to see if wrote to block or not, and then only write back if bit is set
 - Called “**Dirty**” bit (writing makes it “dirty”)

Write-Back Cache

- Store/cache hit, write data in cache *only* & set dirty bit
 - Memory has stale value
- Store/cache miss, read data from memory, then update and set dirty bit
 - “Write-allocate” policy
- Load/cache hit, use value from cache
- On any miss, write back evicted block, only if dirty. Update cache with new block and clear dirty bit.



Write-Through vs. Write-Back

- Write-Through:
 - Simpler control logic
 - More predictable timing
simplifies processor control logic
 - Easier to make reliable, since memory always has copy of data (big idea: Redundancy!)
- Write-Back
 - More complex control logic
 - More variable timing (0,1,2 memory accesses per cache access)
 - Usually reduces write traffic
 - Harder to make reliable, sometimes cache has only copy of data

And In Conclusion, ...

- Principle of Locality for Libraries /Computer Memory
- Hierarchy of Memories (speed/size/cost per bit) to Exploit Locality
- Cache – copy of data lower level in memory hierarchy
- Direct Mapped to find block in cache using Tag field and Valid bit for Hit
- Cache design choice:
 - Write-Through vs. Write-Back