

18-447

Computer Architecture

Lecture 19: Memory Hierarchy and Caches

Prof. Onur Mutlu

Carnegie Mellon University

Spring 2013, 3/18/2014

Extra Credit Recognition for Lab 3

1. John Greth (13157 ns)
2. Kevin Bravo (91332 ns)
3. Elon Bauer (103071 ns)
4. Teng Fei Liao (111500 ns)
5. Albert Cho (127904 ns)
6. Bailey Forrest (130806 ns)

Reminders

- Lab 4: Due March 21
 - Please try to do the extra credit as well!
- Homework 5: Due March 26
- The course will move quickly... Keep your pace. Talk with the TAs and me if you are concerned about your performance.

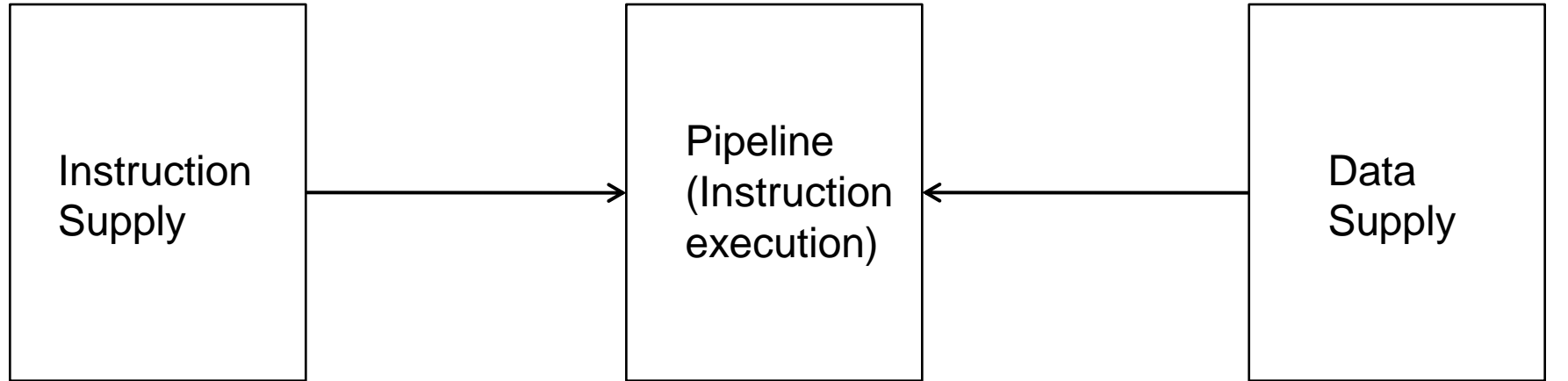
Readings for Today and Next Lecture

- Memory Hierarchy and Caches
- Cache chapters from P&H: 5.1-5.3
- Memory/cache chapters from Hamacher+: 8.1-8.7
- An early cache paper by Maurice Wilkes
 - Wilkes, “Slave Memories and Dynamic Storage Allocation,” IEEE Trans. On Electronic Computers, 1965.

Today

- The memory hierarchy
- Caches

Idealism



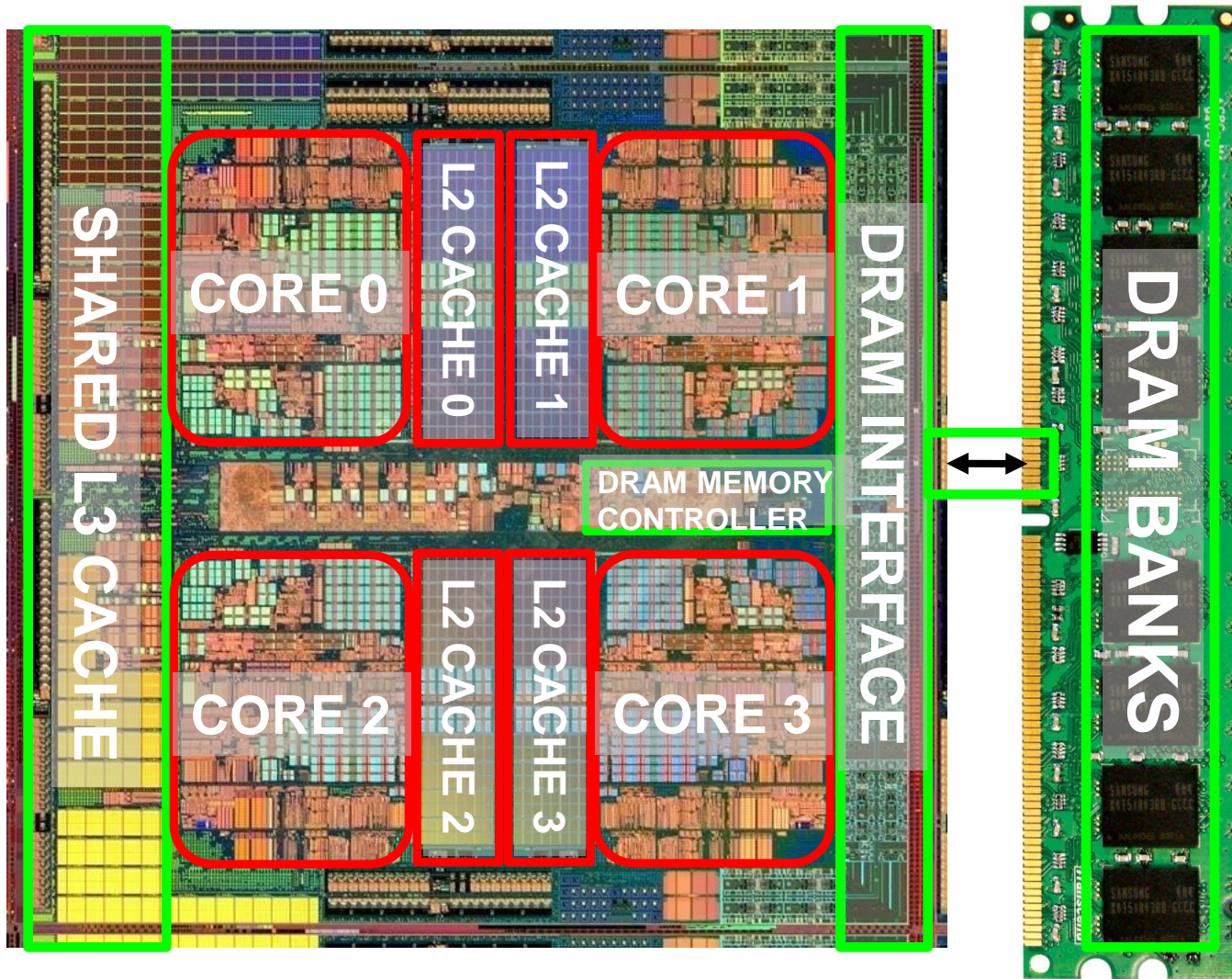
- Zero-cycle latency
- Infinite capacity
- Zero cost
- Perfect control flow

- No pipeline stalls
- Perfect data flow (reg/memory dependencies)
- Zero-cycle interconnect (operand communication)
- Enough functional units
- Zero latency compute

- Zero-cycle latency
- Infinite capacity
- Infinite bandwidth
- Zero cost

The Memory Hierarchy

Memory in a Modern System



Ideal Memory

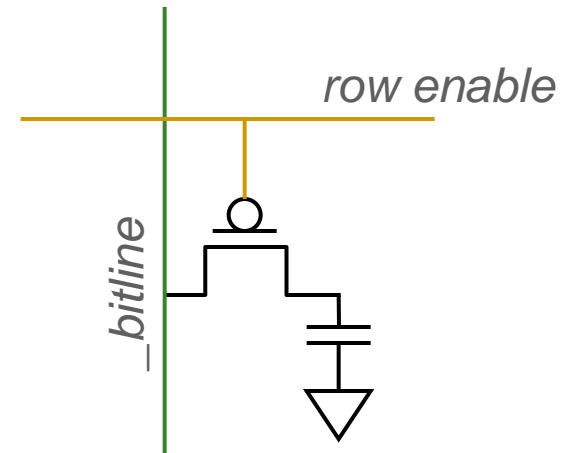
- Zero access time (latency)
- Infinite capacity
- Zero cost
- Infinite bandwidth (to support multiple accesses in parallel)

The Problem

- Ideal memory's requirements oppose each other
- Bigger is slower
 - Bigger → Takes longer to determine the location
- Faster is more expensive
 - Memory technology: SRAM vs. DRAM
- Higher bandwidth is more expensive
 - Need more banks, more ports, higher frequency, or faster technology

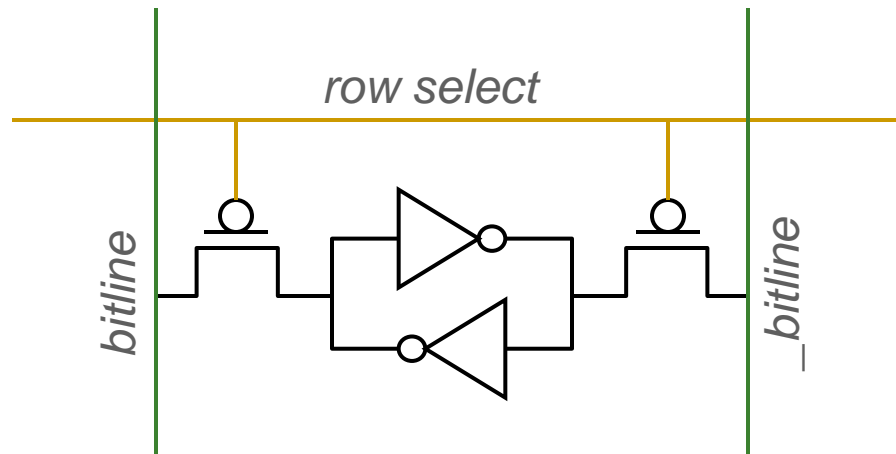
Memory Technology: DRAM

- Dynamic random access memory
- Capacitor charge state indicates stored value
 - Whether the capacitor is charged or discharged indicates storage of 1 or 0
 - 1 capacitor
 - 1 access transistor
- Capacitor leaks through the RC path
 - DRAM cell loses charge over time
 - DRAM cell needs to be refreshed

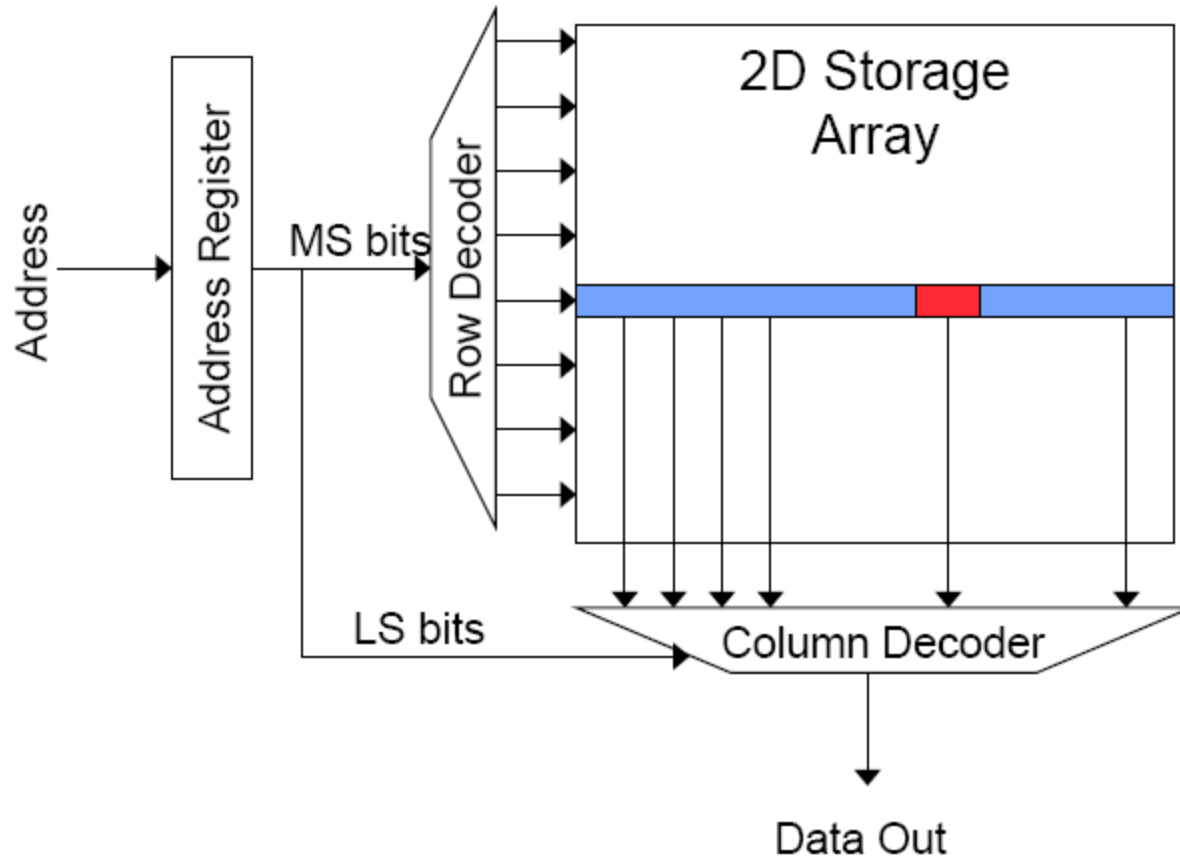


Memory Technology: SRAM

- Static random access memory
- Two cross coupled inverters store a single bit
 - ❑ Feedback path enables the stored value to persist in the “cell”
 - ❑ 4 transistors for storage
 - ❑ 2 transistors for access



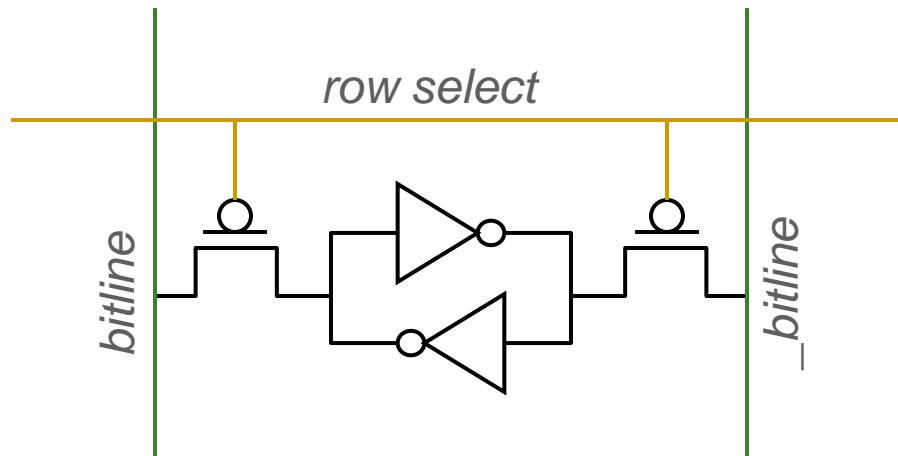
Memory Bank Organization and Operation



■ Read access sequence:

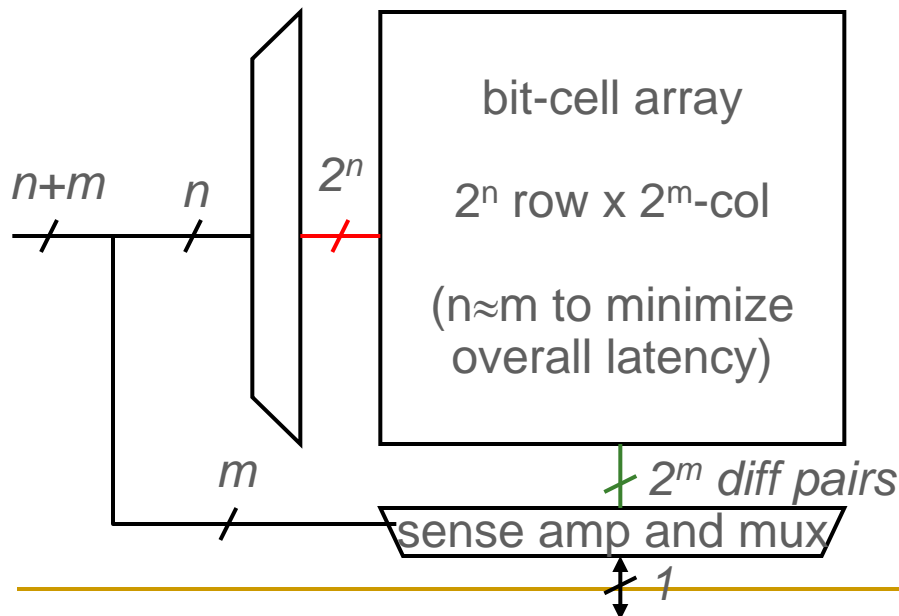
1. Decode row address & drive word-lines
2. Selected bits drive bit-lines
 - Entire row read
3. Amplify row data
4. Decode column address & select subset of row
 - Send to output
5. Precharge bit-lines
 - For next access

SRAM (Static Random Access Memory)



Read Sequence

1. address decode
2. drive row select
3. selected bit-cells drive bitlines
(entire row is read together)
4. differential sensing and column select
(data is ready)
5. precharge all bitlines
(for next read or write)

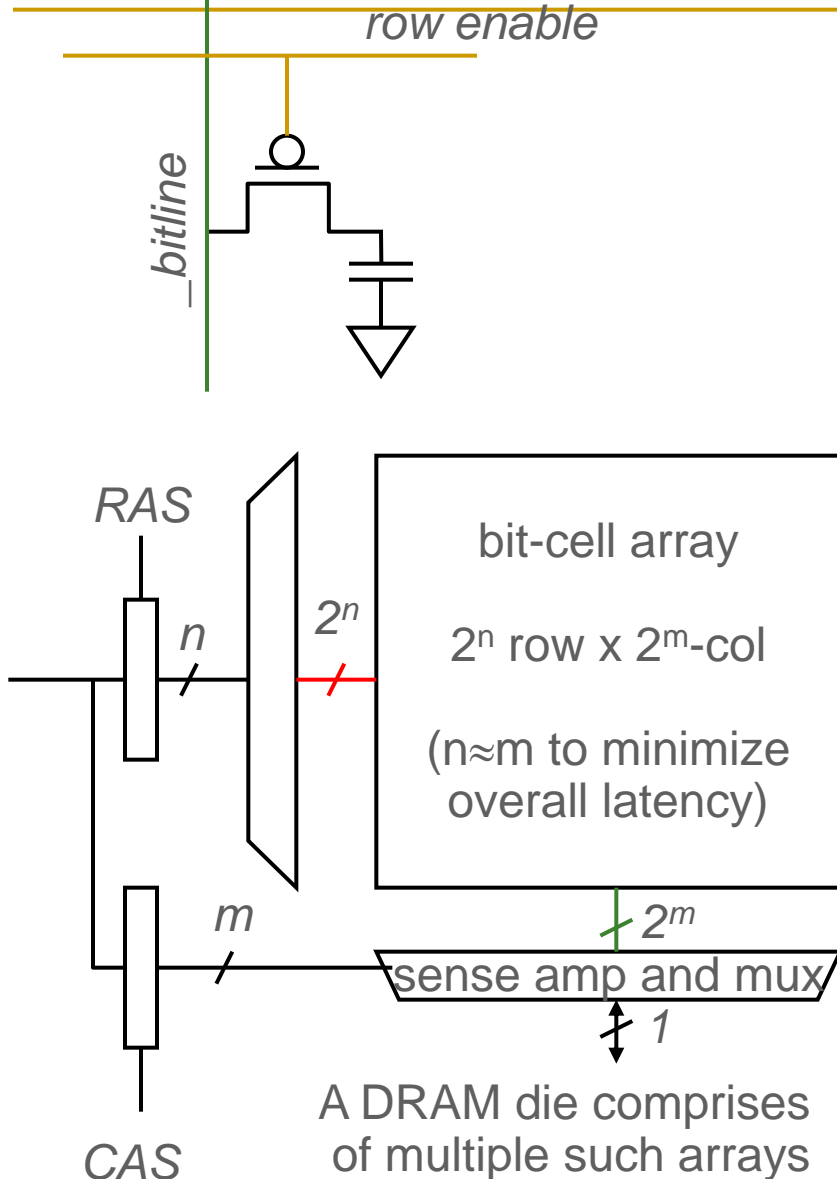


Access latency dominated by steps 2 and 3

Cycling time dominated by steps 2, 3 and 5

- step 2 proportional to 2^m
- step 3 and 5 proportional to 2^n

DRAM (Dynamic Random Access Memory)



Bits stored as charges on node capacitance (non-restorative)

- bit cell loses charge when read
- bit cell loses charge over time

Read Sequence

1~3 same as SRAM

4. a “flip-flopping” sense amp amplifies and regenerates the bitline, data bit is mux’ed out

5. precharge all bitlines

Destructive reads

Charge loss over time

Refresh: A DRAM controller must periodically read each row within the allowed refresh time (10s of ms) such that charge is restored

DRAM vs. SRAM

■ DRAM

- ❑ Slower access (capacitor)
- ❑ Higher density (1T 1C cell)
- ❑ Lower cost
- ❑ Requires refresh (power, performance, circuitry)
- ❑ Manufacturing requires putting capacitor and logic together

■ SRAM

- ❑ Faster access (no capacitor)
- ❑ Lower density (6T cell)
- ❑ Higher cost
- ❑ No need for refresh
- ❑ Manufacturing compatible with logic process (no capacitor)

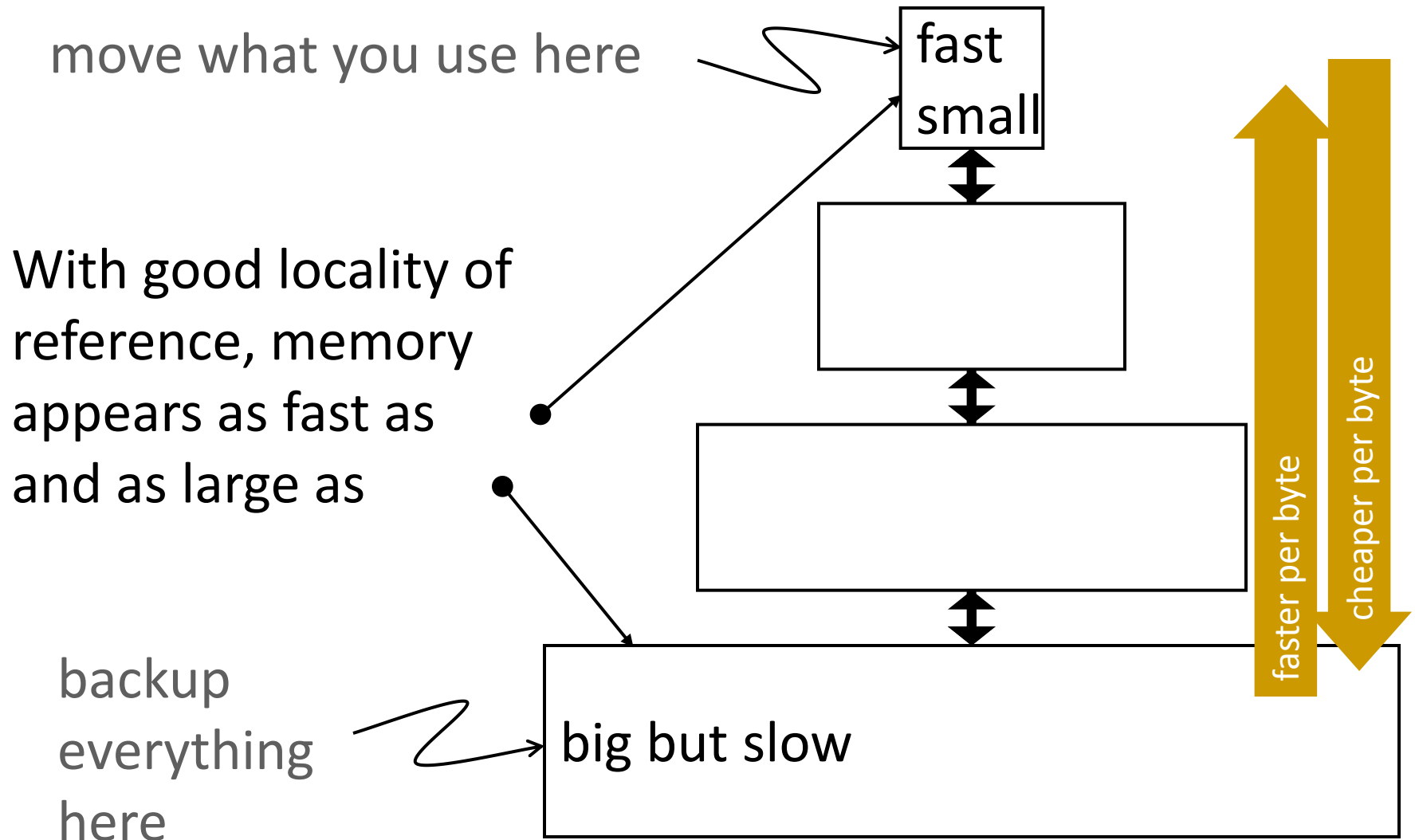
The Problem

- Bigger is slower
 - ❑ SRAM, 512 Bytes, sub-nanosec
 - ❑ SRAM, KByte~MByte, ~nanosec
 - ❑ DRAM, Gigabyte, ~50 nanosec
 - ❑ Hard Disk, Terabyte, ~10 millisec
- Faster is more expensive (dollars and chip area)
 - ❑ SRAM, < 10\$ per Megabyte
 - ❑ DRAM, < 1\$ per Megabyte
 - ❑ Hard Disk < 1\$ per Gigabyte
 - ❑ These sample values scale with time
- Other technologies have their place as well
 - ❑ Flash memory, Phase-change memory (not mature yet)

Why Memory Hierarchy?

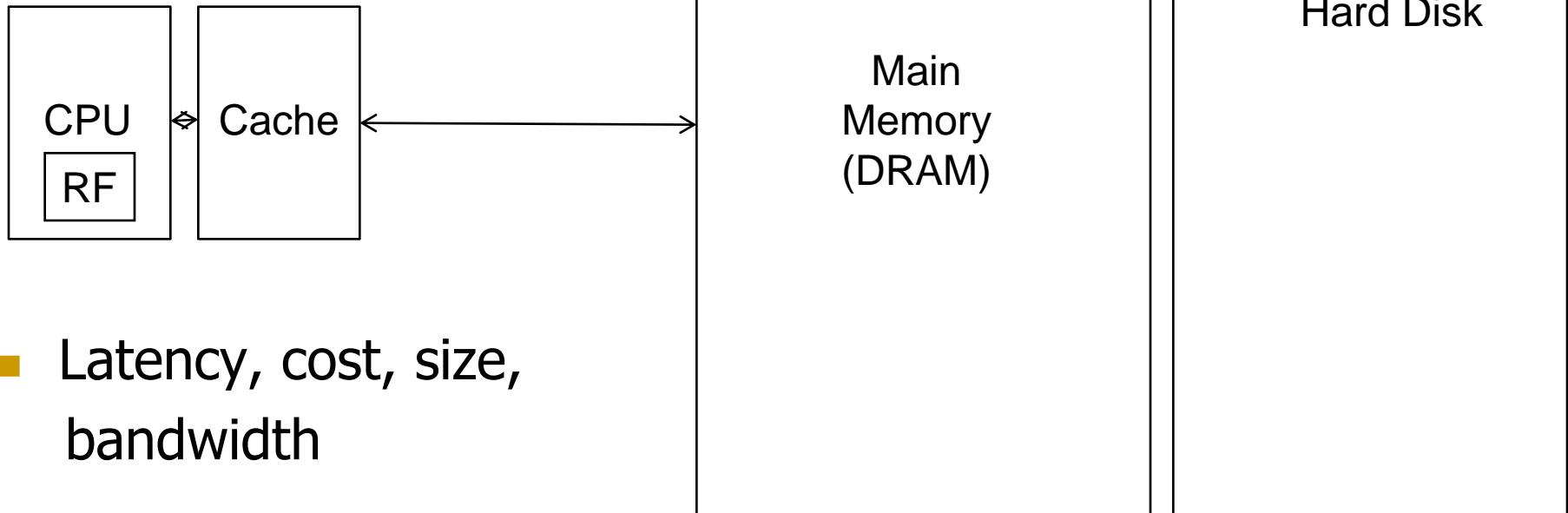
- We want both fast and large
- But we cannot achieve both with a single level of memory
- Idea: Have multiple levels of storage (progressively bigger and slower as the levels are farther from the processor) and ensure most of the data the processor needs is kept in the fast(er) level(s)

The Memory Hierarchy



Memory Hierarchy

- Fundamental tradeoff
 - Fast memory: small
 - Large memory: slow
- Idea: **Memory hierarchy**



- Latency, cost, size, bandwidth

Locality

- One's recent past is a very good predictor of his/her near future.
- **Temporal Locality**: If you just did something, it is very likely that you will do the same thing again soon
 - since you are here today, there is a good chance you will be here again and again regularly
- **Spatial Locality**: If you did something, it is very likely you will do something similar/related (in space)
 - every time I find you in this room, you are probably sitting close to the same people

Memory Locality

- A “typical” program has a lot of locality in memory references
 - typical programs are composed of “loops”
- **Temporal**: A program tends to reference the same memory location many times and all within a small window of time
- **Spatial**: A program tends to reference a cluster of memory locations at a time
 - most notable examples:
 - 1. instruction memory references
 - 2. array/data structure references

Caching Basics: Exploit Temporal Locality

- Idea: Store recently accessed data in automatically managed fast memory (called cache)
- Anticipation: the data will be accessed again soon
- Temporal locality principle
 - Recently accessed data will be again accessed in the near future
 - This is what Maurice Wilkes had in mind:
 - Wilkes, “Slave Memories and Dynamic Storage Allocation,” IEEE Trans. On Electronic Computers, 1965.
 - “The use is discussed of a fast core memory of, say 32000 words as a slave to a slower core memory of, say, one million words in such a way that in practical cases the effective access time is nearer that of the fast memory than that of the slow memory.”

Caching Basics: Exploit Spatial Locality

- Idea: Store addresses adjacent to the recently accessed one in automatically managed fast memory
 - Logically divide memory into equal size blocks
 - Fetch to cache the accessed block in its entirety
- Anticipation: nearby data will be accessed soon
- Spatial locality principle
 - Nearby data in memory will be accessed in the near future
 - E.g., sequential instruction access, array traversal
 - This is what IBM 360/85 implemented
 - 16 Kbyte cache with 64 byte blocks
 - Liptay, “Structural aspects of the System/360 Model 85 II: the cache,” IBM Systems Journal, 1968.

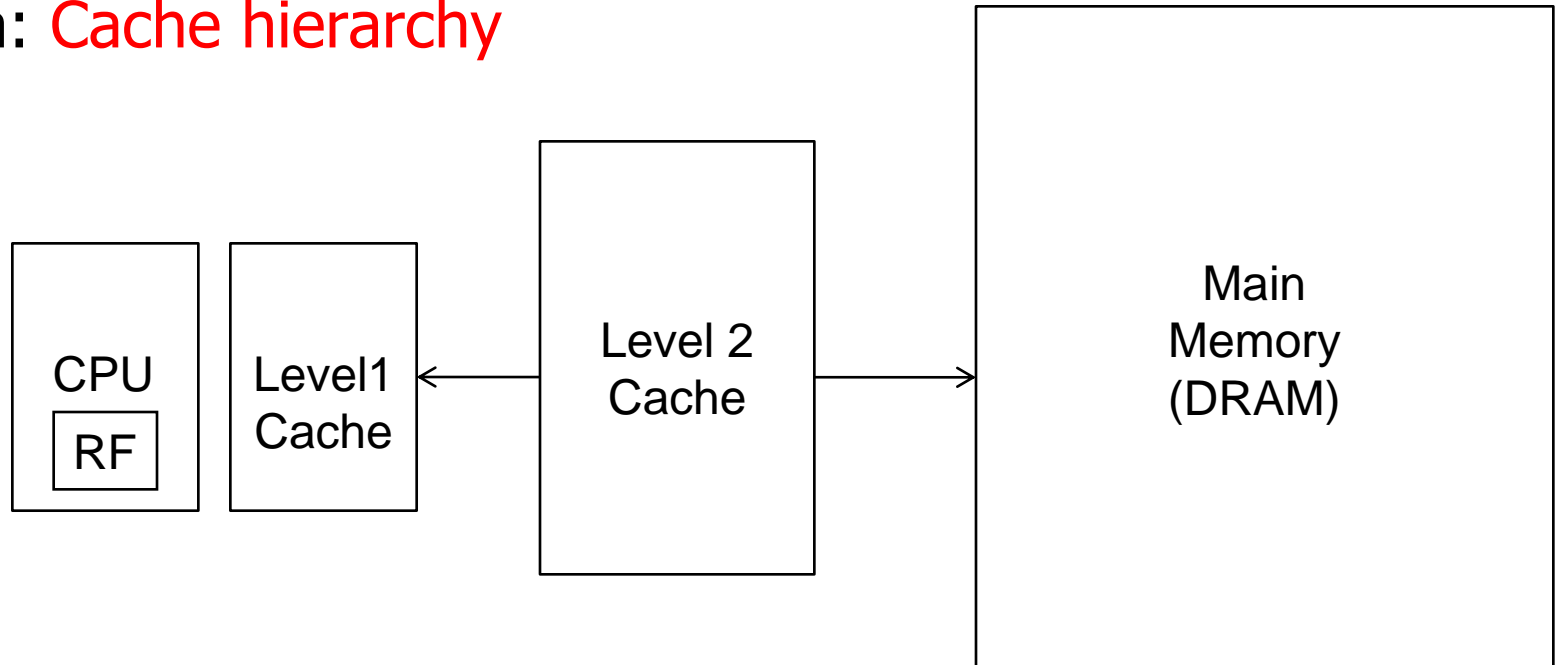
The Bookshelf Analogy

- Book in your hand
- Desk
- Bookshelf
- Boxes at home
- Boxes in storage

- Recently-used books tend to stay on desk
 - Comp Arch books, books for classes you are currently taking
 - Until the desk gets full
- Adjacent books in the shelf needed around the same time
 - If I have organized/categorized my books well in the shelf

Caching in a Pipelined Design

- The cache needs to be tightly integrated into the pipeline
 - Ideally, access in 1-cycle so that dependent operations do not stall
- High frequency pipeline → Cannot make the cache large
 - But, we want a large cache AND a pipelined design
- Idea: **Cache hierarchy**



A Note on Manual vs. Automatic Management

- **Manual:** Programmer manages data movement across levels
 - too painful for programmers on substantial programs
 - “core” vs “drum” memory in the 50’s
 - still done in some embedded processors (on-chip scratch pad SRAM in lieu of a cache)

- **Automatic:** Hardware manages data movement across levels, transparently to the programmer
 - ++ programmer’s life is easier
 - simple heuristic: keep most recently used items in cache
 - the average programmer doesn’t need to know about it
 - You don’t need to know how big the cache is and how it works to write a “correct” program! (What if you want a “fast” program?)

Automatic Management in Memory Hierarchy

- Wilkes, “**Slave Memories and Dynamic Storage Allocation**,” IEEE Trans. On Electronic Computers, 1965.

Slave Memories and Dynamic Storage Allocation

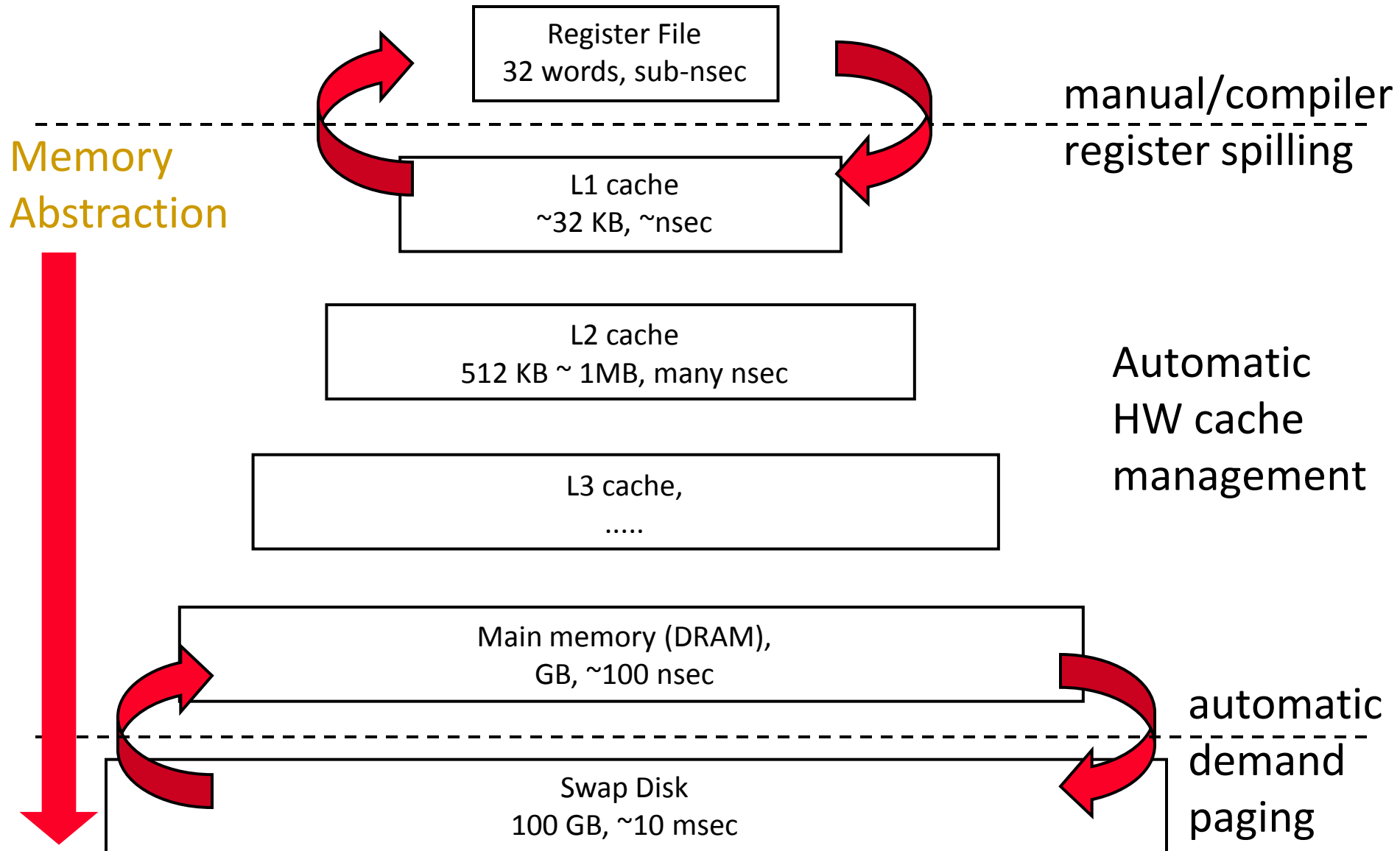
M. V. WILKES

SUMMARY

The use is discussed of a fast core memory of, say, 32 000 words as a slave to a slower core memory of, say, one million words in such a way that in practical cases the effective access time is nearer that of the fast memory than that of the slow memory.

- “By a slave memory I mean one which **automatically accumulates to itself words** that come from a slower main memory, and keeps them available for subsequent use without it being necessary for the penalty of main memory access to be incurred again.”

A Modern Memory Hierarchy



Hierarchical Latency Analysis

- For a given memory hierarchy level i it has a technology-intrinsic access time of t_i . The perceived access time T_i is longer than t_i
- Except for the outer-most hierarchy, when looking for a given address there is
 - a chance (hit-rate h_i) you “hit” and access time is t_i
 - a chance (miss-rate m_i) you “miss” and access time $t_i + T_{i+1}$
 - $h_i + m_i = 1$
- Thus

$$T_i = h_i \cdot t_i + m_i \cdot (t_i + T_{i+1})$$

$$T_i = t_i + m_i \cdot T_{i+1}$$

keep in mind, h_i and m_i are defined to be the hit-rate and miss-rate of just the references that missed at L_{i-1}

Hierarchy Design Considerations

- Recursive latency equation

$$T_i = t_i + m_i \cdot T_{i+1}$$

- The goal: achieve desired T_1 within allowed cost
- $T_i \approx t_i$ is desirable
- Keep m_i low
 - increasing capacity C_i lowers m_i , but beware of increasing t_i
 - lower m_i by smarter management (replacement::anticipate what you don't need, prefetching::anticipate what you will need)
- Keep T_{i+1} low
 - faster lower hierarchies, but beware of increasing cost
 - introduce intermediate hierarchies as a compromise

Intel Pentium 4 Example

- 90nm P4, 3.6 GHz

- L1 D-cache

- $C_1 = 16K$

- $t_1 = 4 \text{ cyc int} / 9 \text{ cycle fp}$

if $m_1=0.1, m_2=0.1$

$T_1=7.6, T_2=36$

- L2 D-cache

- $C_2 = 1024 \text{ KB}$

- $t_2 = 18 \text{ cyc int} / 18 \text{ cyc fp}$

if $m_1=0.01, m_2=0.01$

$T_1=4.2, T_2=19.8$

- Main memory

- $t_3 = \sim 50\text{ns or } 180 \text{ cyc}$

if $m_1=0.05, m_2=0.01$

$T_1=5.00, T_2=19.8$

- Notice

- best case latency is not 1

- worst case access latencies are into 500+ cycles

if $m_1=0.01, m_2=0.50$

$T_1=5.08, T_2=108$

Cache Basics and Operation

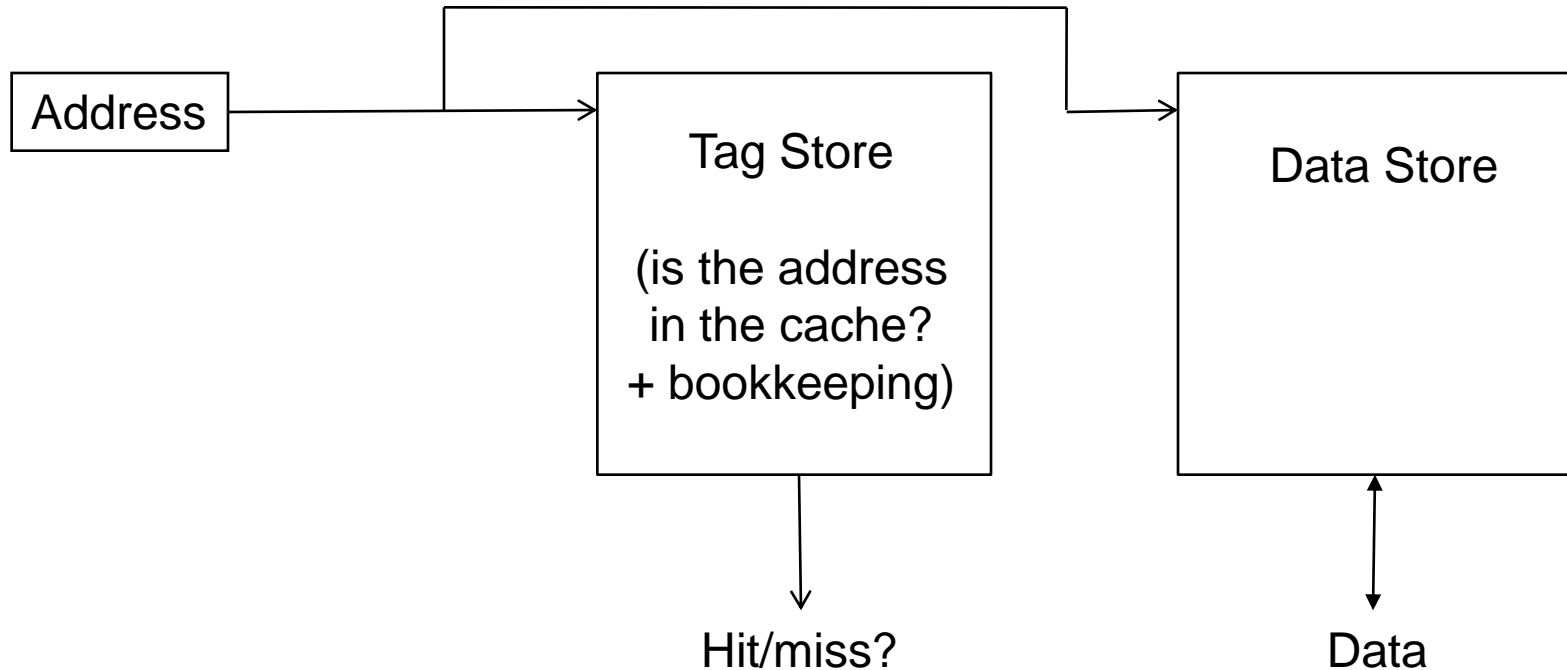
Cache

- Generically, any structure that “memoizes” frequently used results to avoid repeating the long-latency operations required to reproduce the results from scratch, e.g. a web cache
- Most commonly in the on-die context: an automatically-managed memory hierarchy based on SRAM
 - memoize in SRAM the most frequently accessed DRAM memory locations to avoid repeatedly paying for the DRAM access latency

Caching Basics

- Block (line): Unit of storage in the cache
 - Memory is logically divided into cache blocks that map to locations in the cache
- When data referenced
 - HIT: If in cache, use cached data instead of accessing memory
 - MISS: If not in cache, bring block into cache
 - Maybe have to kick something else out to do it
- Some important cache design decisions
 - Placement: where and how to place/find a block in cache?
 - Replacement: what data to remove to make room in cache?
 - Granularity of management: large, small, uniform blocks?
 - Write policy: what do we do about writes?
 - Instructions/data: Do we treat them separately?

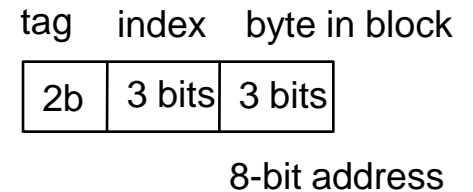
Cache Abstraction and Metrics



- Cache hit rate = $(\# \text{ hits}) / (\# \text{ hits} + \# \text{ misses}) = (\# \text{ hits}) / (\# \text{ accesses})$
- Average memory access time (AMAT)
= $(\text{hit-rate} * \text{hit-latency}) + (\text{miss-rate} * \text{miss-latency})$
- Aside: *Can reducing AMAT reduce performance?*

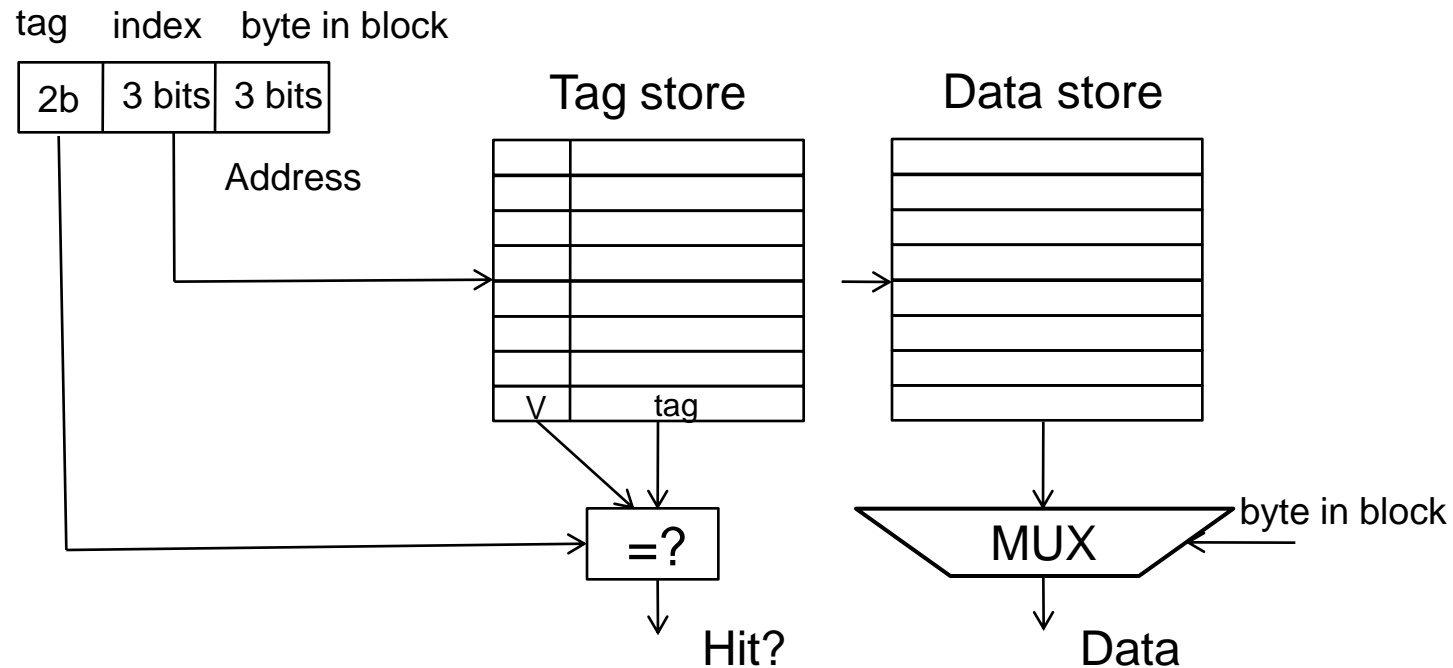
Blocks and Addressing the Cache

- Memory is logically divided into cache blocks
- Each block maps to a location in the cache, determined by the index bits in the address
 - used to index into the tag and data stores
- Cache access: index into the tag and data stores with index bits in address, check valid bit in tag store, compare tag bits in address with the stored tag in tag store
- If a block is in the cache (cache hit), the tag store should have the tag of the block stored in the index of the block



Direct-Mapped Cache: Placement and Access

- Assume byte-addressable memory: 256 bytes, 8-byte blocks → 32 blocks
- Assume cache: 64 bytes, 8 blocks
 - ❑ **Direct-mapped: A block can go to only one location**



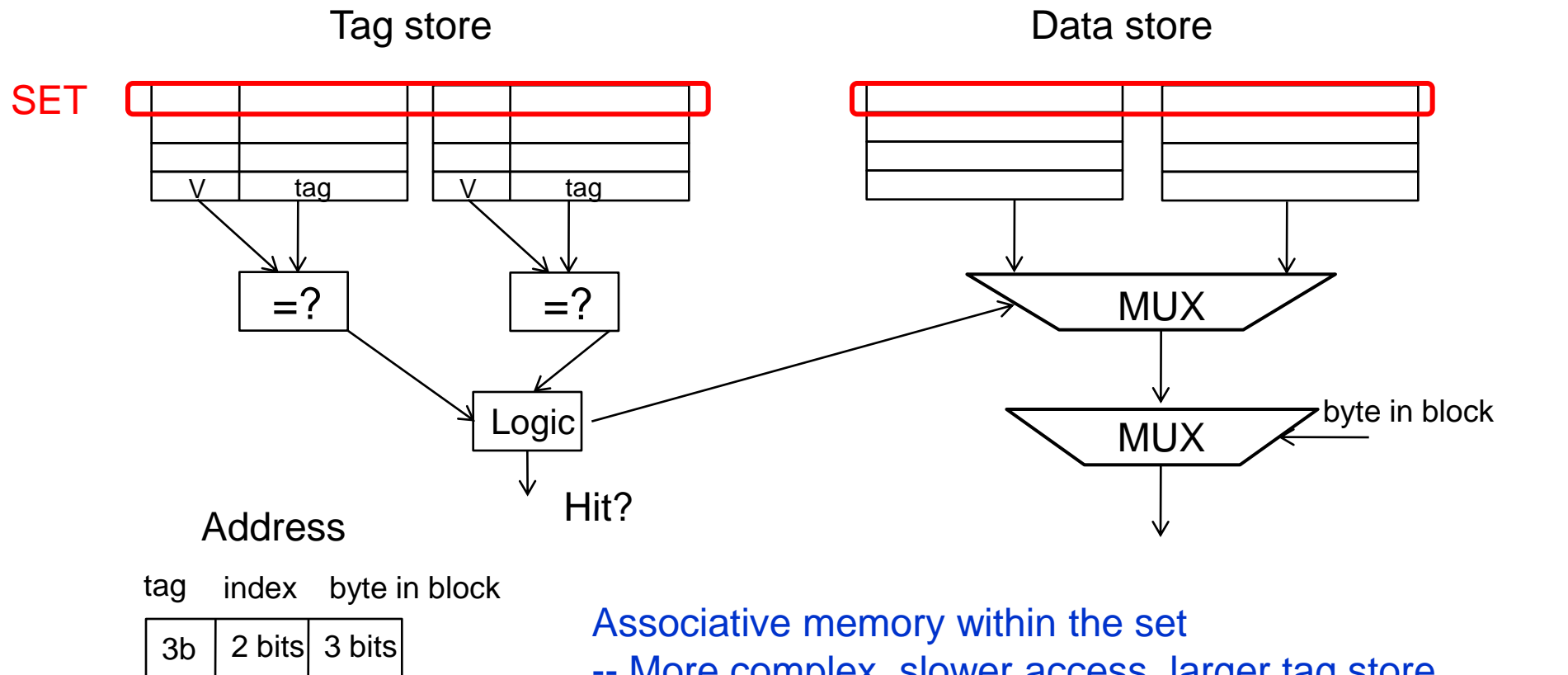
- ❑ **Addresses with same index contend for the same location**
 - **Cause conflict misses**

Direct-Mapped Caches

- **Direct-mapped cache:** Two blocks in memory that map to the same index in the cache cannot be present in the cache at the same time
 - One index \rightarrow one entry
- Can lead to 0% hit rate if more than one block accessed in an interleaved manner map to the same index
 - Assume addresses A and B have the same index bits but different tag bits
 - A, B, A, B, A, B, A, B, ... \rightarrow conflict in the cache index
 - All accesses are **conflict misses**

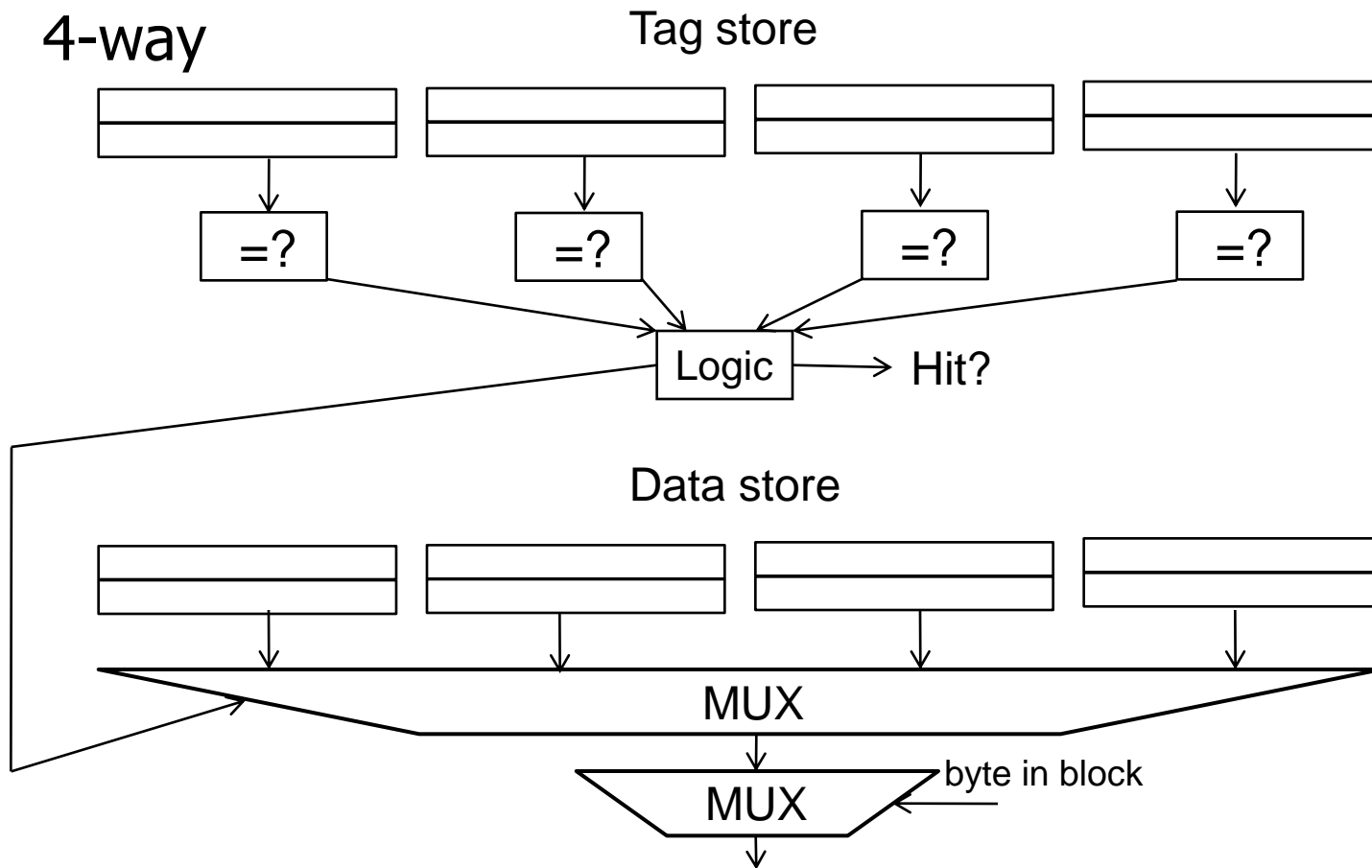
Set Associativity

- Addresses 0 and 8 always conflict in direct mapped cache
- Instead of having one column of 8, have 2 columns of 4 blocks



Higher Associativity

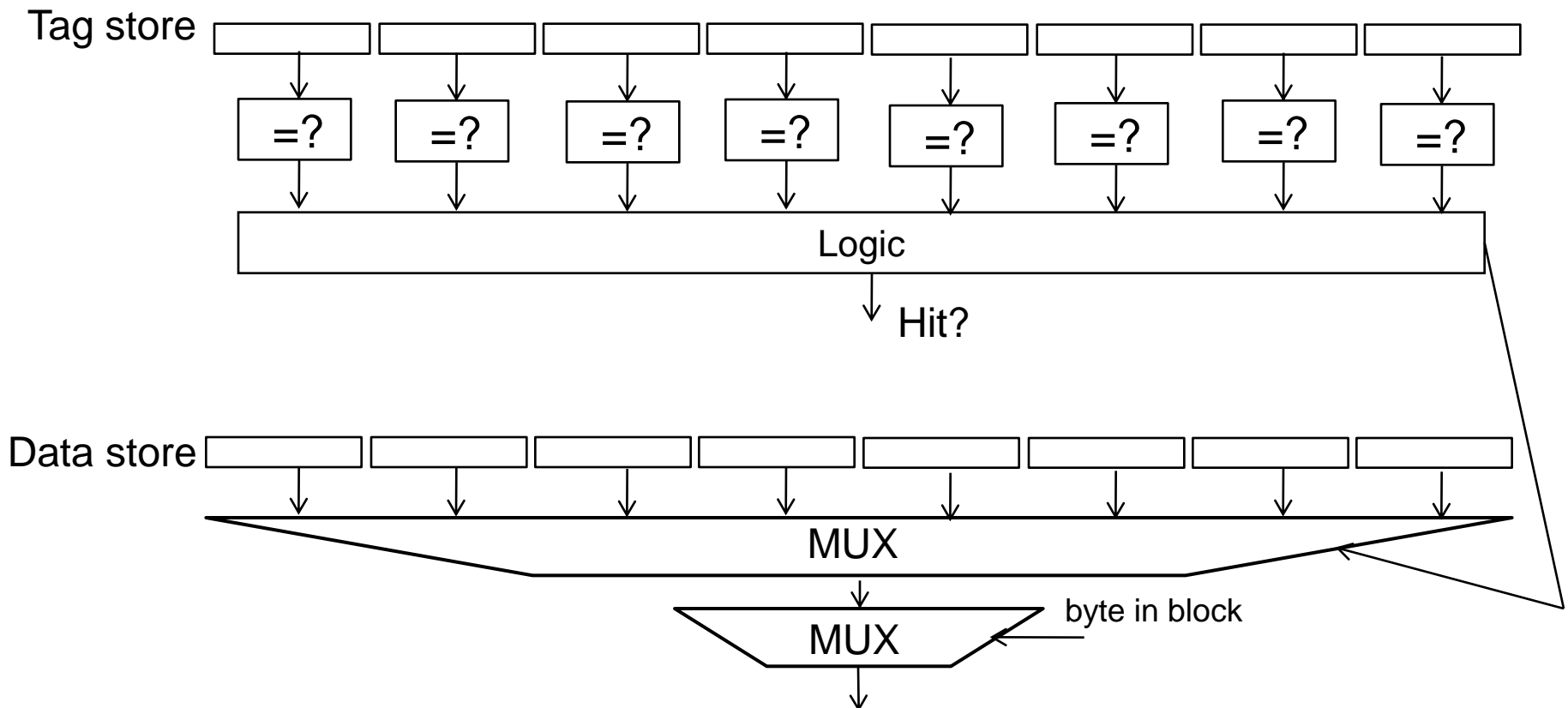
■ 4-way



- More tag comparators and wider data mux; larger tags
- + Likelihood of conflict misses even lower

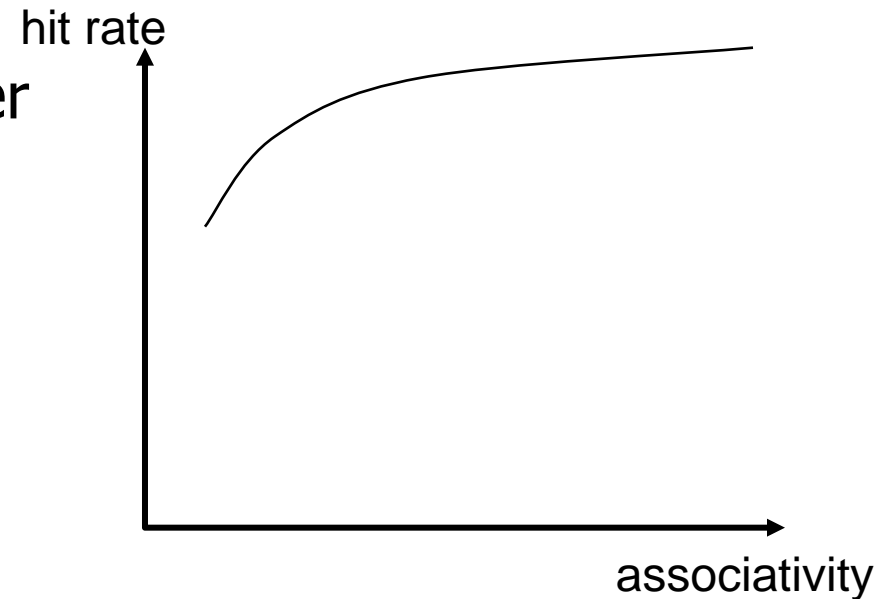
Full Associativity

- Fully associative cache
 - A block can be placed in any cache location



Associativity (and Tradeoffs)

- How many blocks can map to the same index (or set)?
- Higher associativity
 - ++ Higher hit rate
 - Slower cache access time (hit latency and data access latency)
 - More expensive hardware (more comparators)
- Diminishing returns from higher associativity



Set-Associative Caches (I)

- Diminishing returns in hit rate from higher associativity
- Longer access time with higher associativity
- Which block in the set to replace on a cache miss?
 - Any invalid block first
 - If all are valid, consult the replacement policy
 - Random
 - FIFO
 - Least recently used (how to implement?)
 - Not most recently used
 - Least frequently used?
 - Least costly to re-fetch?
 - Why would memory accesses have different cost?
 - Hybrid replacement policies
 - Optimal replacement policy?

Implementing LRU

- Idea: Evict the least recently accessed block
- Problem: Need to keep track of access ordering of blocks
- Question: 2-way set associative cache:
 - What do you need to implement LRU?
- Question: 4-way set associative cache:
 - How many different orderings possible for the 4 blocks in the set?
 - How many bits needed to encode the LRU order of a block?
 - What is the logic needed to determine the LRU victim?

Approximations of LRU

- Most modern processors do not implement “true LRU” in highly-associative caches
- Why?
 - True LRU is complex
 - LRU is an approximation to predict locality anyway (i.e., not the best possible replacement policy)
- Examples:
 - Not MRU (not most recently used)
 - Hierarchical LRU: divide the 4-way set into 2-way “groups”, track the MRU group and the MRU way in each group
 - Victim-NextVictim Replacement: Only keep track of the victim and the next victim

Hierarchical LRU (not MRU)

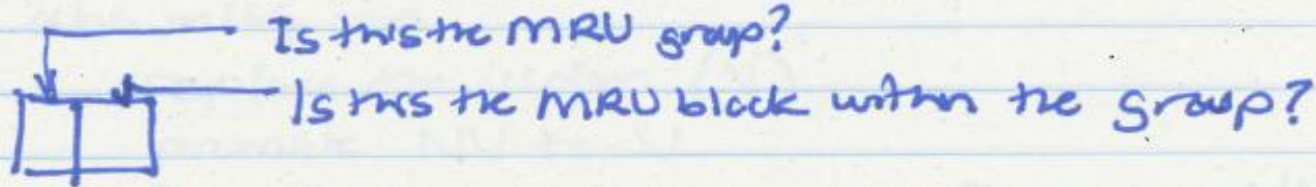
- Divide a set into multiple groups
- Keep track of the MRU group
- Keep track of the MRU block in each group
- On replacement, select victim as:
 - A not-MRU block in one of the not-MRU groups

Hierarchical LRU (not MRU) Example

Hierarchical LRU

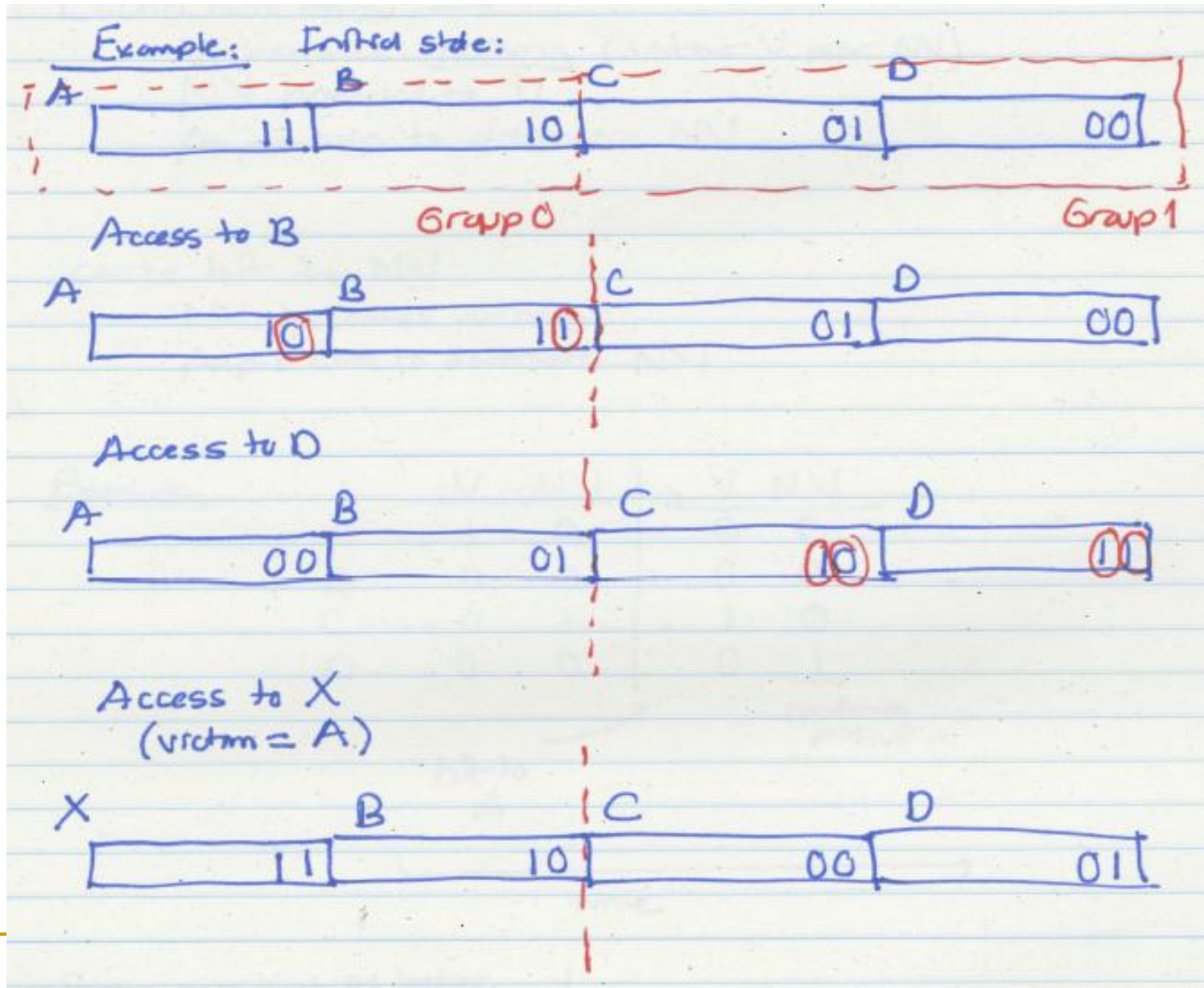
4-way cache

2 bits for replacement for each way in the tag store



Victim: The block that is not the MRU block and that is not in the MRU group

Hierarchical LRU (not MRU) Example



Hierarchical LRU (not MRU): Questions

- 8-way cache
- 2 4-way groups
- What is an access pattern that performs worse than true LRU?
- What is an access pattern that performs better than true LRU?

Victim/Next-Victim Policy

- Only 2 blocks' status tracked in each set:
 - victim (V), next victim (NV)
 - all other blocks denoted as (O) – Ordinary block

- On a cache miss
 - Replace V
 - Promote NV to V
 - Randomly pick an O block as NV

- On a cache hit to V
 - Promote NV to V
 - Randomly pick an O block as NV
 - Turn V to O

Victim/Next-Victim Policy (II)

- On a cache hit to NV
 - Randomly pick an O block as NV
 - Turn NV to O
- On a cache hit to O
 - Do nothing

Victim/Next-Victim Example

Example

	V	NV		V	NV
A	1	0		0	0
B	0	0		0	0
C	0	1		1	0
D	0	0		0	1

hit to A

randomly picked

time

Some questions as before

Replacement Policy

■ LRU vs. Random

- ❑ **Set thrashing:** When the “program working set” in a set is larger than set associativity
- ❑ 4-way: Cyclic references to A, B, C, D, E
 - 0% hit rate with LRU policy
- ❑ Random replacement policy is better when thrashing occurs

■ In practice:

- ❑ Depends on workload
- ❑ Average hit rate of LRU and Random are similar

■ Hybrid of LRU and Random

- ❑ How to choose between the two? **Set sampling**
 - See Qureshi et al., “**A Case for MLP-Aware Cache Replacement**,” ISCA 2006.

Optimal Replacement Policy?

- Belady's OPT
 - Replace the block that is going to be referenced furthest in the future by the program
 - Belady, “A study of replacement algorithms for a virtual-storage computer,” IBM Systems Journal, 1966.
 - How do we implement this? Simulate?
- Is this optimal for minimizing miss rate?
- Is this optimal for minimizing execution time?
 - No. Cache miss latency/cost varies from block to block!
 - Two reasons: Remote vs. local caches and miss overlapping
 - Qureshi et al. “A Case for MLP-Aware Cache Replacement,” ISCA 2006.

Aside: Cache versus Page Replacement

- Physical memory (DRAM) is a cache for disk
 - Usually managed by system software via the virtual memory subsystem
- Page replacement is similar to cache replacement
- Page table is the “tag store” for physical memory data store
- What is the difference?
 - Hardware versus software
 - Number of blocks in a cache versus physical memory
 - “Tolerable” amount of time to find a replacement candidate

What's In A Tag Store Entry?

- Valid bit
- Tag
- Replacement policy bits

- Dirty bit?
 - Write back vs. write through caches

Handling Writes (Stores)

- When do we write the modified data in a cache to the next level?
 - Write through: At the time the write happens
 - Write back: When the block is evicted
- Write-back
 - + Can consolidate multiple writes to the same block before eviction
 - Potentially saves bandwidth between cache levels + saves energy
 - Need a bit in the tag store indicating the block is “modified”
- Write-through
 - + Simpler
 - + All levels are up to date. Consistency: Simpler cache coherence because no need to check lower-level caches
 - More bandwidth intensive; no coalescing of writes

Handling Writes (Stores)

- Do we allocate a cache block on a write miss?
 - Allocate on write miss: Yes
 - No-allocate on write miss: No
- Allocate on write miss
 - + Can consolidate writes instead of writing each of them individually to next level
 - + Simpler because write misses can be treated the same way as read misses
 - Requires (?) transfer of the whole cache block
- No-allocate
 - + Conserves cache space if locality of writes is low (potentially better cache hit rate)

Sectored Caches

- Divide a block into subblocks (or sectors)
 - Have separate valid and dirty bits for each sector
 - When is this useful? (Think writes...)
 - How many subblocks do you transfer on a read?
- ++ No need to transfer the entire cache block into the cache
(A write simply validates and updates a subblock)
- ++ More freedom in transferring subblocks into the cache (a cache block does not need to be in the cache fully)
- More complex design
- May not exploit spatial locality fully when used for reads



Instruction vs. Data Caches

- Unified:
 - + Dynamic sharing of cache space: no overprovisioning that might happen with static partitioning (i.e., split I and D caches)
 - Instructions and data can thrash each other (i.e., no guaranteed space for either)
 - I and D are accessed in different places in the pipeline. Where do we place the unified cache for fast access?
- First level caches are almost always split
 - Mainly for the last reason above
- Second and higher levels are almost always unified

Multi-level Caching in a Pipelined Design

- First-level caches (instruction and data)
 - Decisions very much affected by cycle time
 - Small, lower associativity
 - Tag store and data store accessed in parallel
- Second-level caches
 - Decisions need to balance hit rate and access latency
 - Usually large and highly associative; latency not as important
 - Tag store and data store accessed serially
- Serial vs. Parallel access of levels
 - Serial: Second level cache accessed only if first-level misses
 - Second level does not see the same accesses as the first
 - First level acts as a filter

Virtual Memory and Cache Interaction

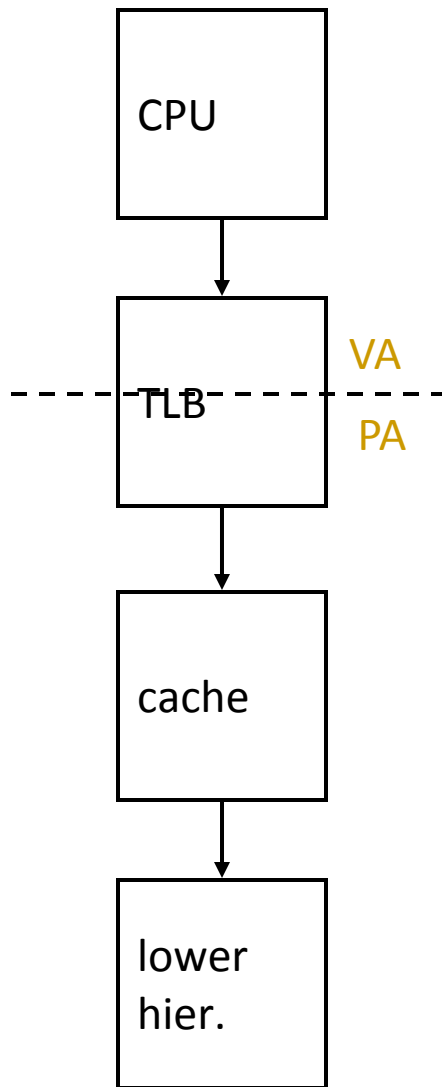
Address Translation and Caching

- When do we do the address translation?
 - Before or after accessing the L1 cache?
- In other words, is the cache virtually addressed or physically addressed?
 - Virtual versus physical cache
- What are the issues with a virtually addressed cache?
- **Synonym problem:**
 - Two different virtual addresses can map to the same physical address → same physical address can be present in multiple locations in the cache → can lead to inconsistency in data

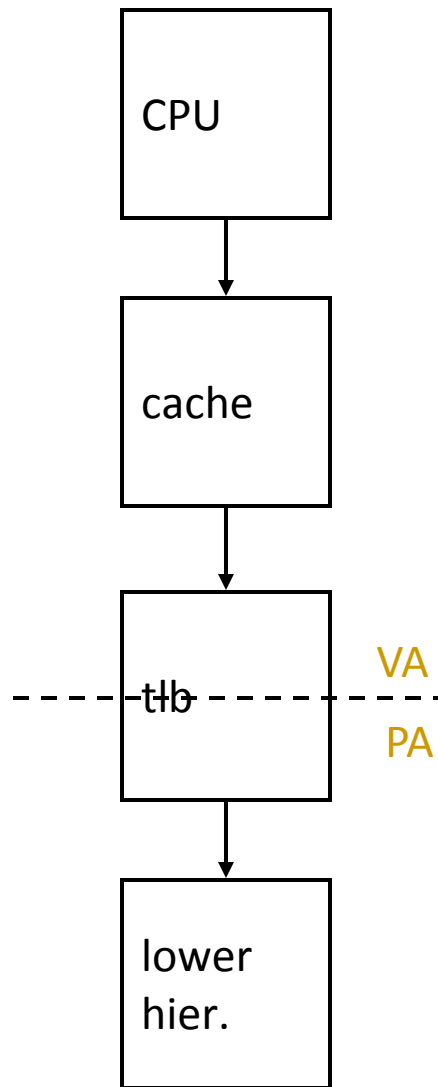
Homonyms and Synonyms

- **Homonym: Same VA can map to two different PAs**
 - Why?
 - VA is in different processes
- **Synonym: Different VAs can map to the same PA**
 - Why?
 - Different pages can share the same physical frame within or across processes
 - Reasons: shared libraries, shared data, copy-on-write pages within the same process, ...
- Do homonyms and synonyms create problems when we have a cache?
 - Is the cache virtually or physically addressed?

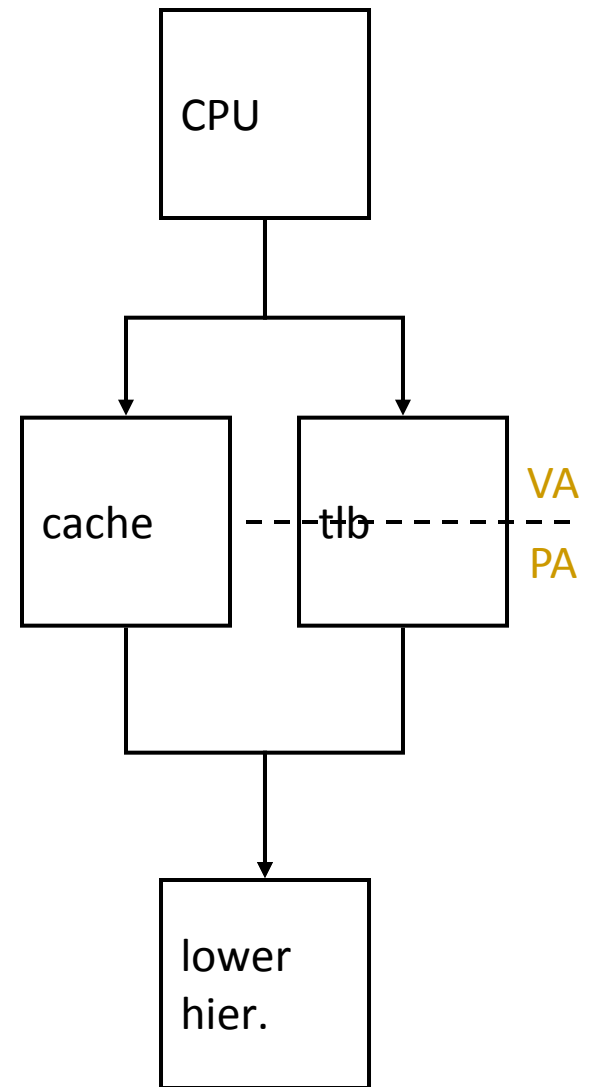
Cache-VM Interaction



physical cache

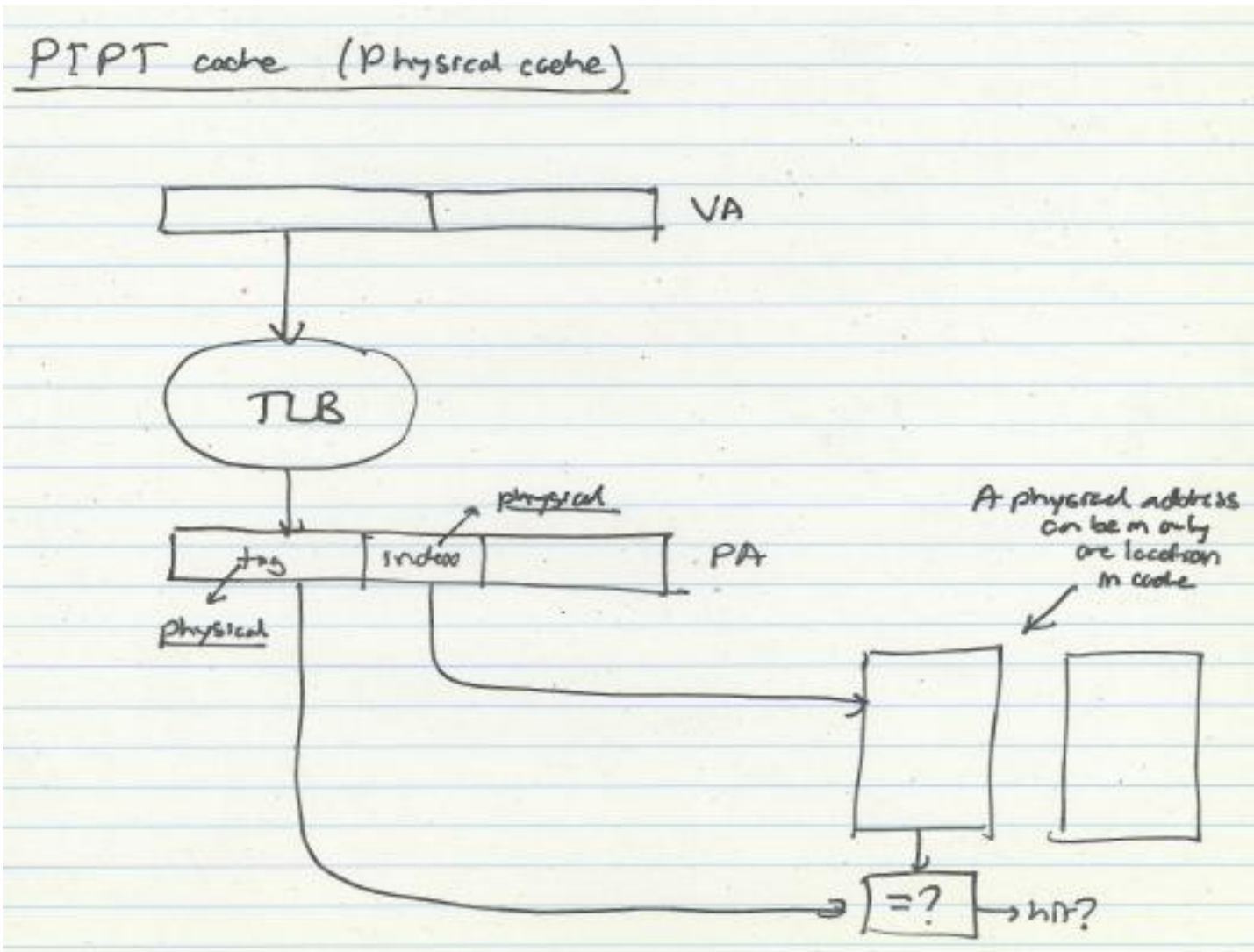


virtual (L1) cache



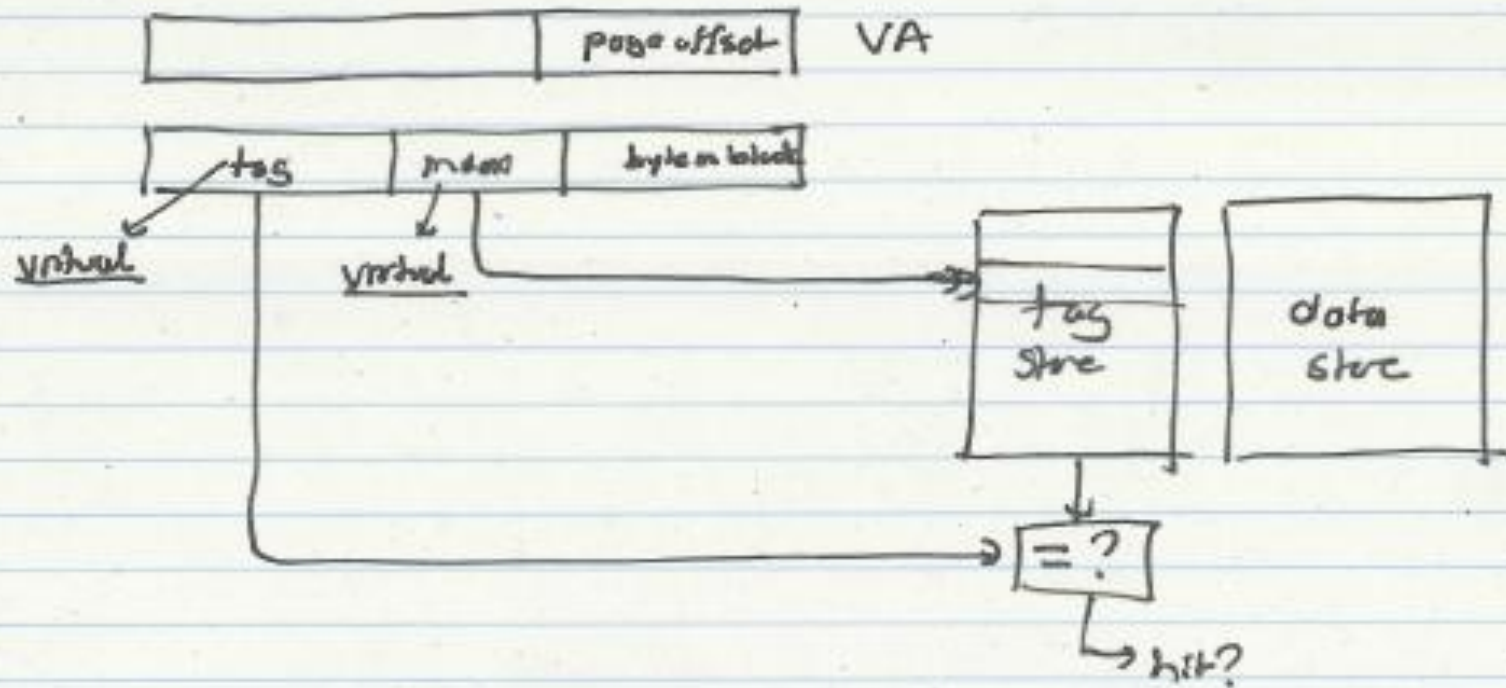
virtual-physical cache

Physical Cache



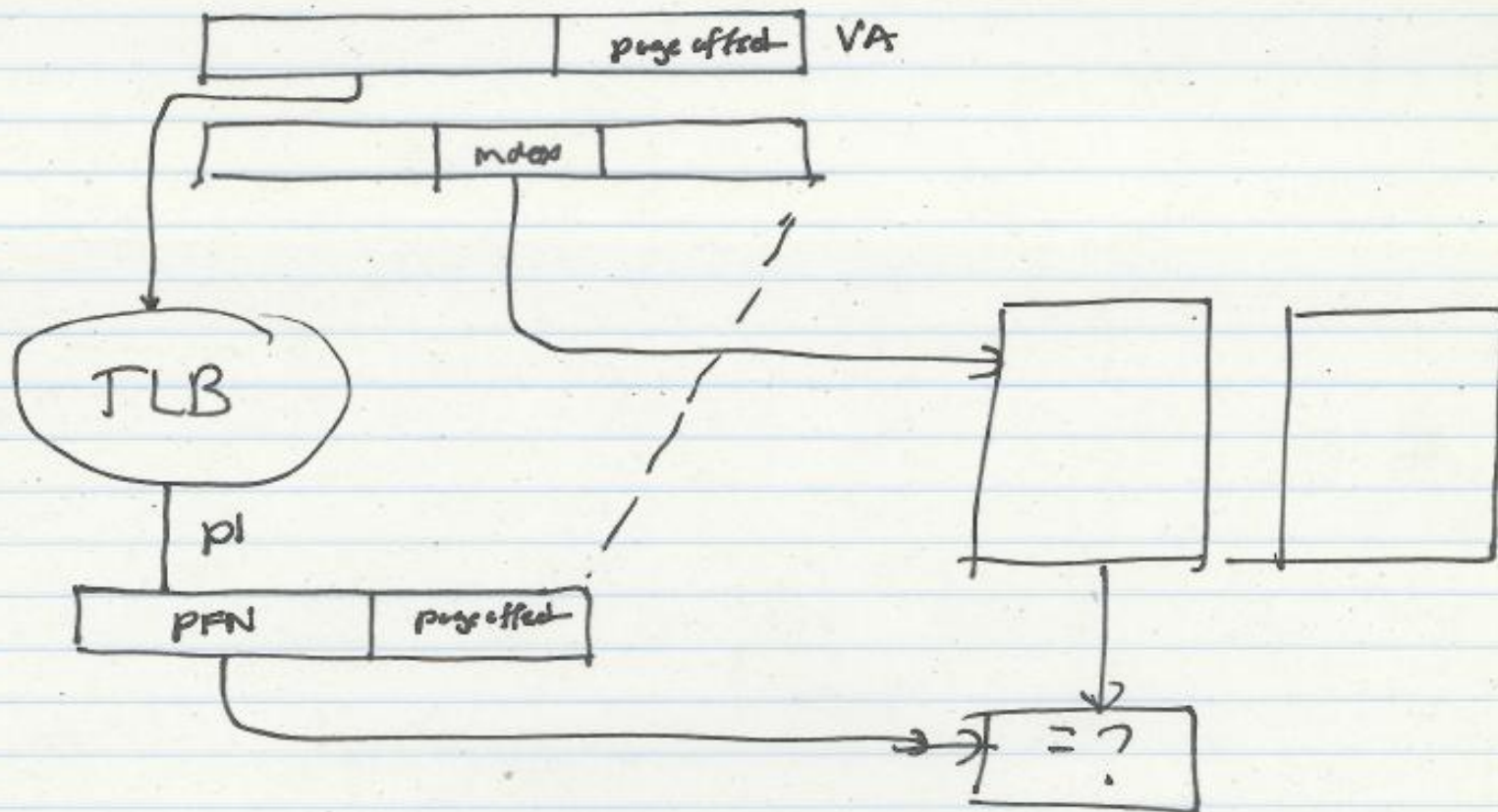
Virtual Cache

VINT cache (Virtual Cache)



Virtual-Physical Cache

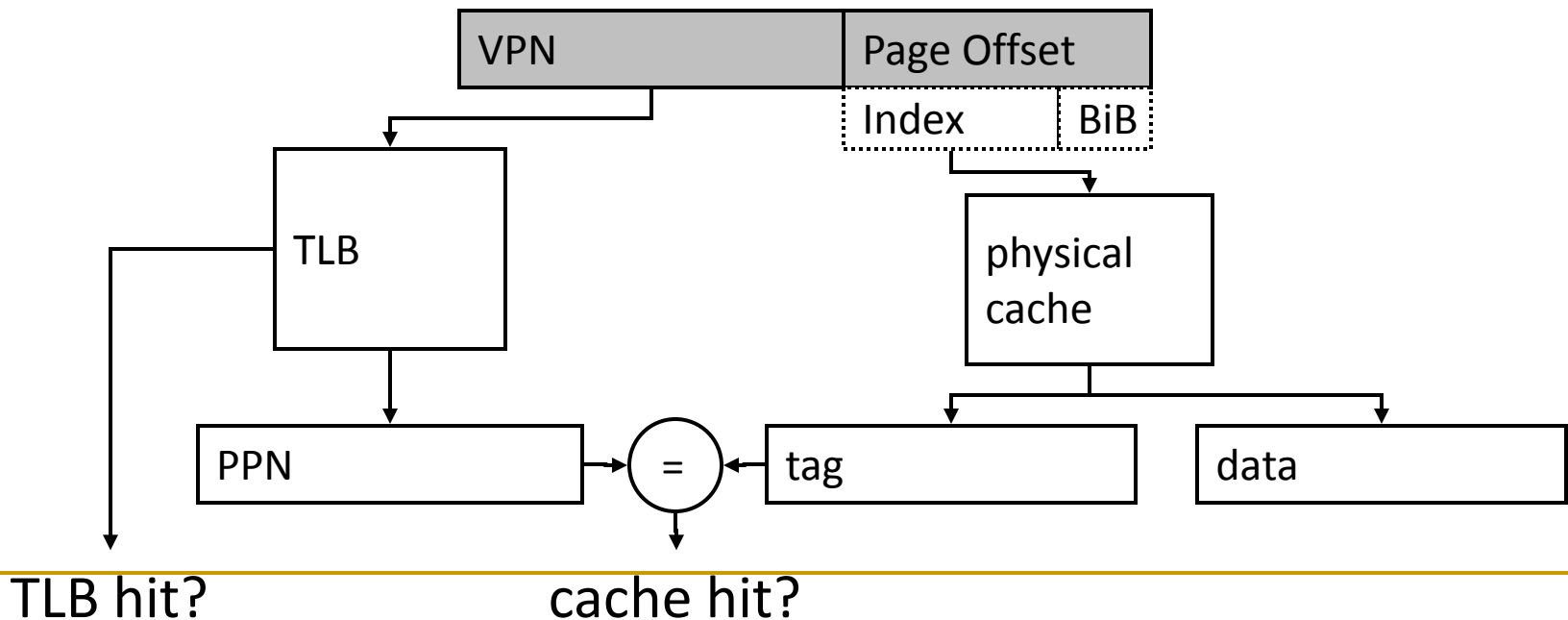
VIPT cache



Where can the same physical address be in the cache?

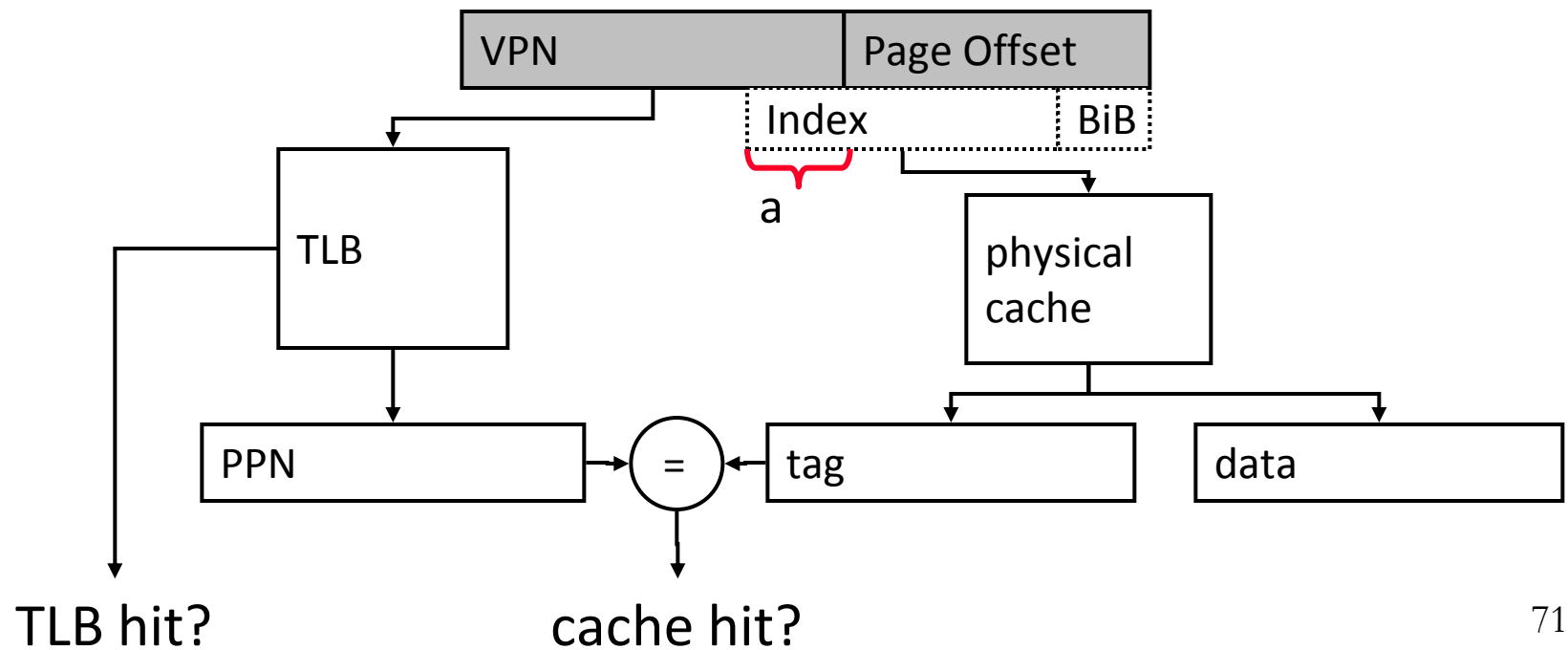
Virtually-Indexed Physically-Tagged

- If $C \leq (\text{page_size} \times \text{associativity})$, the cache index bits come only from page offset (same in VA and PA)
- If both cache and TLB are on chip
 - index both arrays concurrently using VA bits
 - check cache tag (physical) against TLB output at the end



Virtually-Indexed Physically-Tagged

- If $C > (\text{page_size} \times \text{associativity})$, the cache index bits include VPN
⇒ Synonyms can cause problems
 - The same physical address can exist in two locations
- Solutions?



Some Solutions to the Synonym Problem

- Limit cache size to (page size times associativity)
 - get index from page offset
- On a write to a block, search all possible indices that can contain the same physical block, and update/invalidate
 - Used in Alpha 21264, MIPS R10K
- Restrict page placement in OS
 - make sure $\text{index}(\text{VA}) = \text{index}(\text{PA})$
 - Called page coloring
 - Used in many SPARC processors

An Exercise

- Problem 5 from
 - ECE 741 midterm exam Problem 5, Spring 2009
 - http://www.ece.cmu.edu/~ece740/f11/lib/exe/fetch.php?media=wiki:midterm:midterm_s09.pdf

An Exercise (I)

We have a byte-addressable toy computer that has a physical address space of 512 bytes. The computer uses a simple, one-level virtual memory system. The page table is always in physical memory. The page size is specified as 8 bytes and the virtual address space is 2 KB.

Part A.

i. (1 point)

How many bits of each virtual address is the virtual page number?

ii. (1 point)

How many bits of each physical address is the physical frame number?

We would like to add a 128-byte *write-through* cache to enhance the performance of this computer. However, we would like the cache access and address translation to be performed simultaneously. In other words, we would like to index our cache using a virtual address, but do the tag comparison using the physical addresses (virtually-indexed physically-tagged). The cache we would like to add is direct-mapped, and has a block size of 2 bytes. The replacement policy is LRU. Answer the following questions:

iii. (1 point)

How many bits of a virtual address are used to determine which byte in a block is accessed?

iv. (2 point)

How many bits of a virtual address are used to index into the cache? Which bits exactly?

v. (1 point)

How many bits of the virtual page number are used to index into the cache?

vi. (5 points)

What is the size of the tag store in bits? Show your work.

Part B.

Suppose we have two processes sharing our toy computer. These processes share some portion of the physical memory. Some of the virtual page-physical frame mappings of each process are given below:

PROCESS 0	
Virtual Page	Physical Frame
Page 0	Frame 0
Page 3	Frame 7
Page 7	Frame 1
Page 15	Frame 3

PROCESS 1	
Virtual Page	Physical Frame
Page 0	Frame 4
Page 1	Frame 5
Page 7	Frame 3
Page 11	Frame 2

vii. (2 points)

Give a complete physical address whose data can exist in two different locations in the cache.

viii. (3 points)

Give the indexes of those two different locations in the cache.

An Exercise (Concluded)

ix. (5 points)

We do not want the same physical address stored in two different locations in the 128-byte cache. We can prevent this by increasing the associativity of our virtually-indexed physically-tagged cache. What is the minimum associativity required?

x. (4 points)

Assume we would like to use a direct-mapped cache. Describe a solution that ensures that the same physical address is never stored in two different locations in the 128-byte cache.

Solutions to the Exercise

- http://www.ece.cmu.edu/~ece740/f11/lib/exe/fetch.php?media=wiki:midterm:midterm_s09_solution.pdf
- And, more exercises are in past exams and in your homeworks...

Review: Solutions to the Synonym Problem

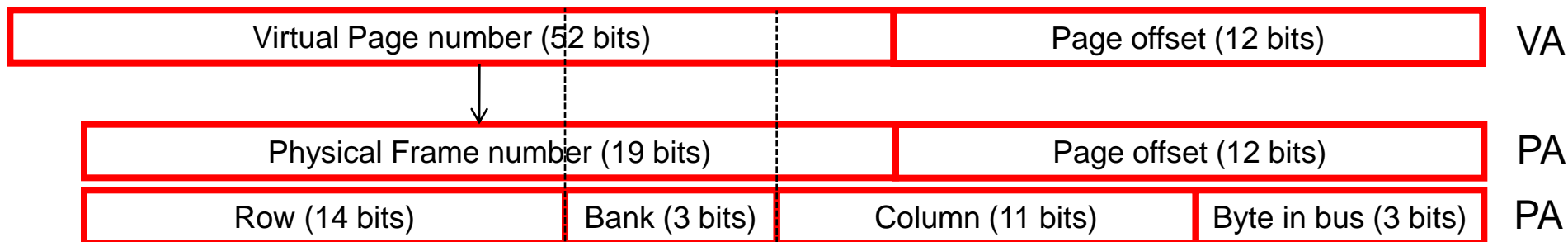
- Limit cache size to (page size times associativity)
 - get index from page offset
- On a write to a block, search all possible indices that can contain the same physical block, and update/invalidate
 - Used in Alpha 21264, MIPS R10K
- Restrict page placement in OS
 - make sure $\text{index}(\text{VA}) = \text{index}(\text{PA})$
 - Called page coloring
 - Used in many SPARC processors

Some Questions to Ponder

- At what cache level should we worry about the synonym and homonym problems?
- What levels of the memory hierarchy does the system software's page mapping algorithms influence?
- What are the potential benefits and downsides of page coloring?

Virtual Memory – DRAM Interaction

- Operating System influences where an address maps to in DRAM



- Operating system can control which bank/channel/rank a virtual page is mapped to.
- It can perform page coloring to minimize bank conflicts
- Or to minimize inter-application interference

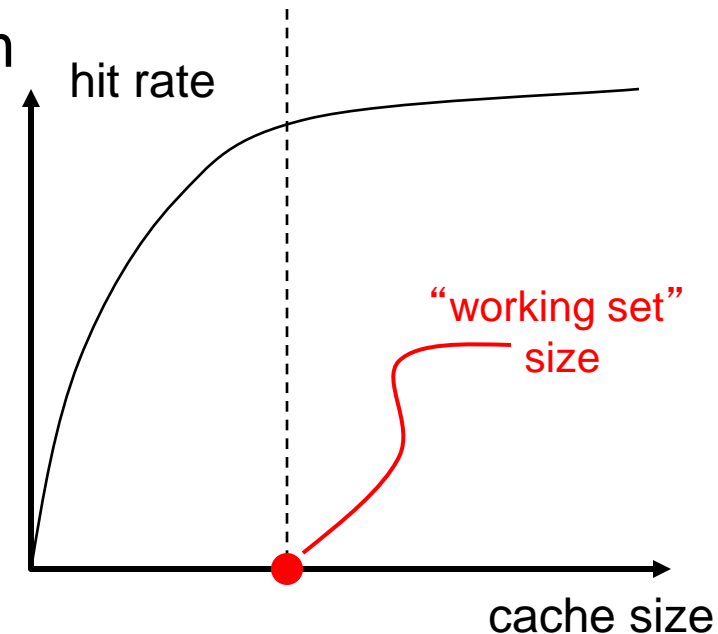
Cache Performance

Cache Parameters vs. Miss Rate

- Cache size
- Block size
- Associativity
- Replacement policy
- Insertion/Placement policy

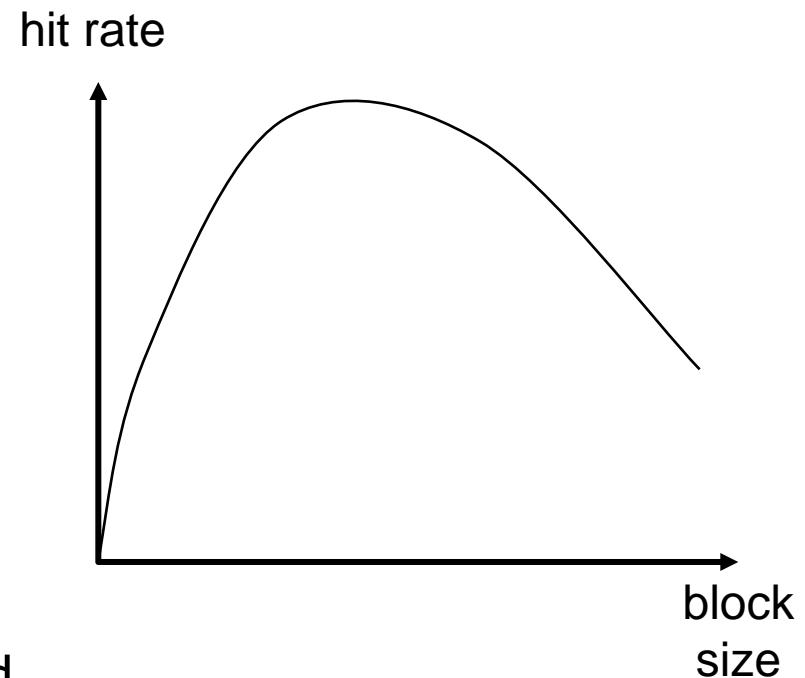
Cache Size

- Cache size: total data (not including tag) capacity
 - bigger can exploit temporal locality better
 - not ALWAYS better
- Too large a cache adversely affects hit and miss latency
 - smaller is faster => bigger is slower
 - access time may degrade critical path
- Too small a cache
 - doesn't exploit temporal locality well
 - useful data replaced often
- **Working set**: the whole set of data the executing application references
 - Within a time interval



Block Size

- Block size is the data that is associated with an address tag
 - not necessarily the unit of transfer between hierarchies
 - Sub-blocking: A block divided into multiple pieces (each with V bit)
 - Can improve “write” performance
- Too small blocks
 - don't exploit spatial locality well
 - have larger tag overhead
- Too large blocks
 - too few total # of blocks
 - likely-useless data transferred
 - Extra bandwidth/energy consumed



Large Blocks: Critical-Word and Subblocking

- Large cache blocks can take a long time to fill into the cache
 - ❑ fill cache line **critical word first**
 - ❑ restart cache access before complete fill
- Large cache blocks can waste bus bandwidth
 - ❑ divide a block into subblocks
 - ❑ associate separate valid bits for each subblock
 - ❑ **When is this useful?**



Associativity

- How many blocks can map to the same index (or set)?

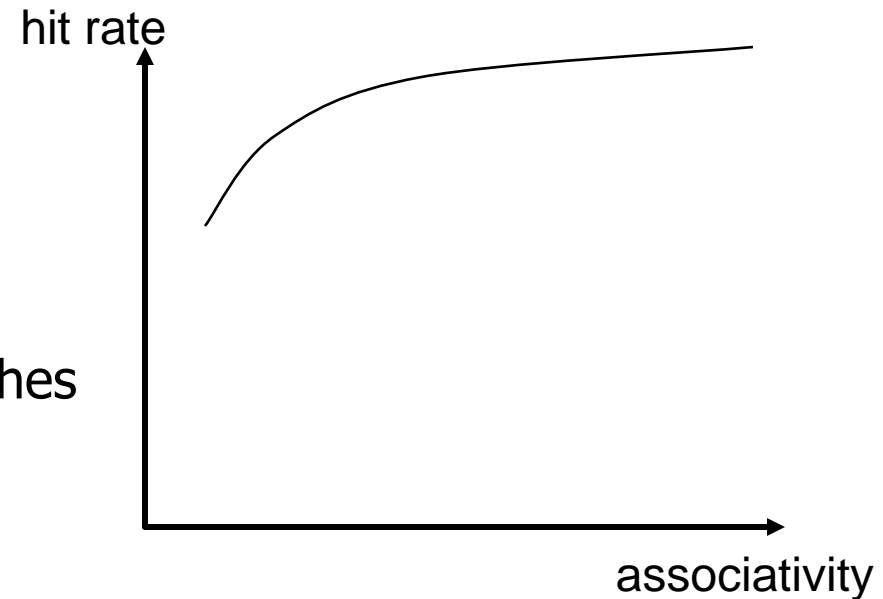
- Larger associativity

- lower miss rate, less variation among programs
- diminishing returns, higher hit latency

- Smaller associativity

- lower cost
- lower hit latency
 - Especially important for L1 caches

- Power of 2 associativity?



Classification of Cache Misses

- Compulsory miss
 - ❑ first reference to an address (block) always results in a miss
 - ❑ subsequent references should hit unless the cache block is displaced for the reasons below
 - ❑ dominates when locality is poor
- Capacity miss
 - ❑ cache is too small to hold everything needed
 - ❑ defined as the misses that would occur even in a fully-associative cache (with optimal replacement) of the same capacity
- Conflict miss
 - ❑ defined as any miss that is neither a compulsory nor a capacity miss

How to Reduce Each Miss Type

- Compulsory
 - Caching cannot help
 - Prefetching
- Conflict
 - More associativity
 - Other ways to get more associativity without making the cache associative
 - Victim cache
 - Hashing
 - Software hints?
- Capacity
 - Utilize cache space better: keep blocks that will be referenced
 - Software management: divide working set such that each “phase” fits in cache

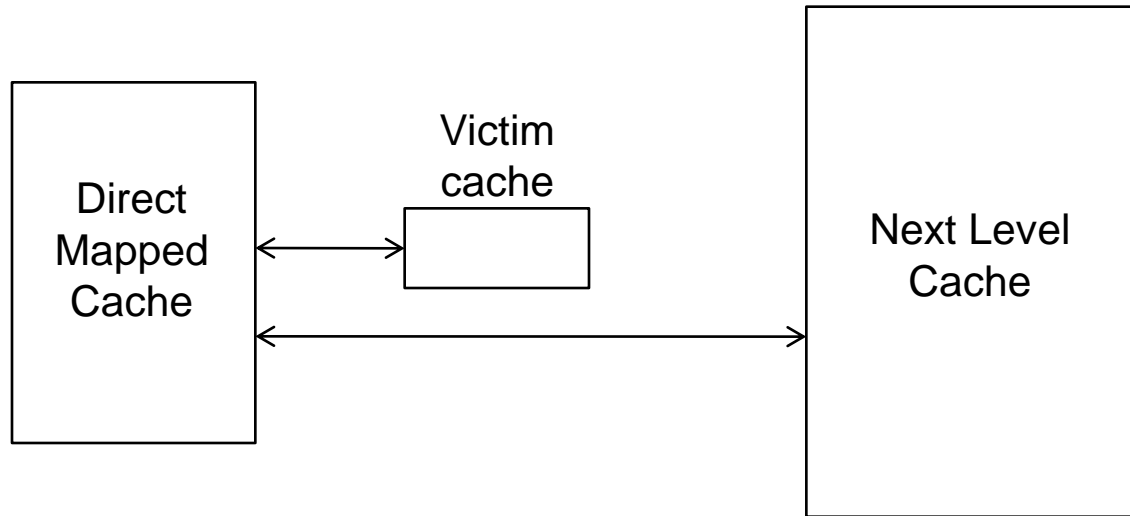
Improving Cache “Performance”

- Remember
 - Average memory access time (AMAT)
$$= (\text{hit-rate} * \text{hit-latency}) + (\text{miss-rate} * \text{miss-latency})$$
- Reducing miss rate
 - Caveat: reducing miss rate can reduce performance if more costly-to-refetch blocks are evicted
- Reducing miss latency/cost
- Reducing hit latency

Improving Basic Cache Performance

- Reducing miss rate
 - ❑ More associativity
 - ❑ Alternatives/enhancements to associativity
 - Victim caches, hashing, pseudo-associativity, skewed associativity
 - ❑ Better replacement/insertion policies
 - ❑ Software approaches
- Reducing miss latency/cost
 - ❑ Multi-level caches
 - ❑ Critical word first
 - ❑ Subblocking/sectoring
 - ❑ Better replacement/insertion policies
 - ❑ Non-blocking caches (multiple cache misses in parallel)
 - ❑ Multiple accesses per cycle
 - ❑ Software approaches

Victim Cache: Reducing Conflict Misses



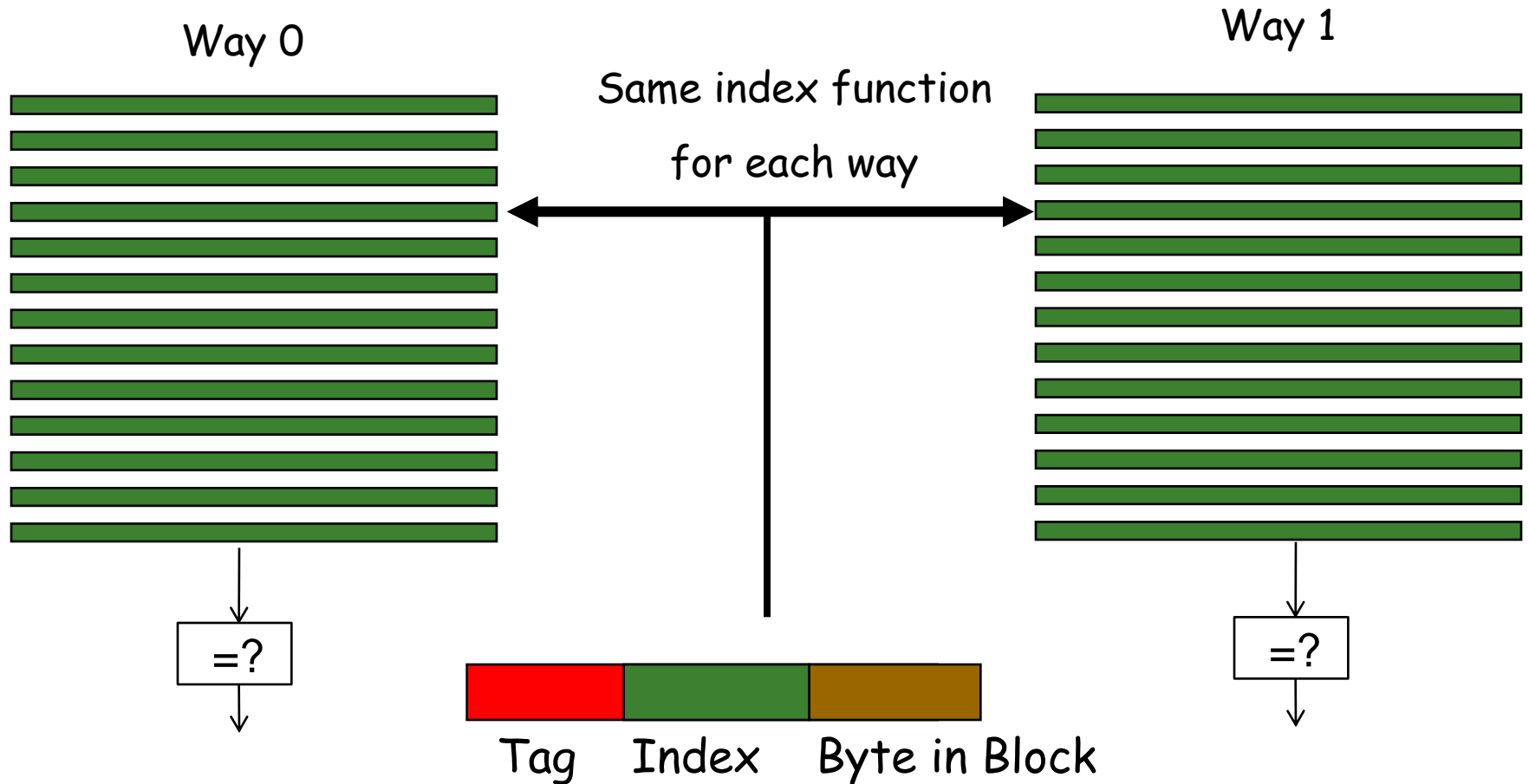
- Jouppi, “Improving Direct-Mapped Cache Performance by the Addition of a Small Fully-Associative Cache and Prefetch Buffers,” ISCA 1990.
- Idea: Use a small fully associative buffer (victim cache) to store evicted blocks
 - + Can avoid ping ponging of cache blocks mapped to the same set (if two cache blocks continuously accessed in nearby time conflict with each other)
 - Increases miss latency if accessed serially with L2

Hashing and Pseudo-Associativity

- Hashing: Better “randomizing” index functions
 - + can reduce conflict misses
 - by distributing the accessed memory blocks more evenly to sets
 - Example: stride where stride value equals cache size
 - More complex to implement: can lengthen critical path
- Pseudo-associativity (Poor Man’s associative cache)
 - Serial lookup: On a miss, use a different index function and access cache again
 - Given a direct-mapped array with K cache blocks
 - Implement K/N sets
 - Given address Addr, sequentially look up: $\{0, \text{Addr}[\lg(K/N)-1: 0]\}$, $\{1, \text{Addr}[\lg(K/N)-1: 0]\}$, ... , $\{N-1, \text{Addr}[\lg(K/N)-1: 0]\}$

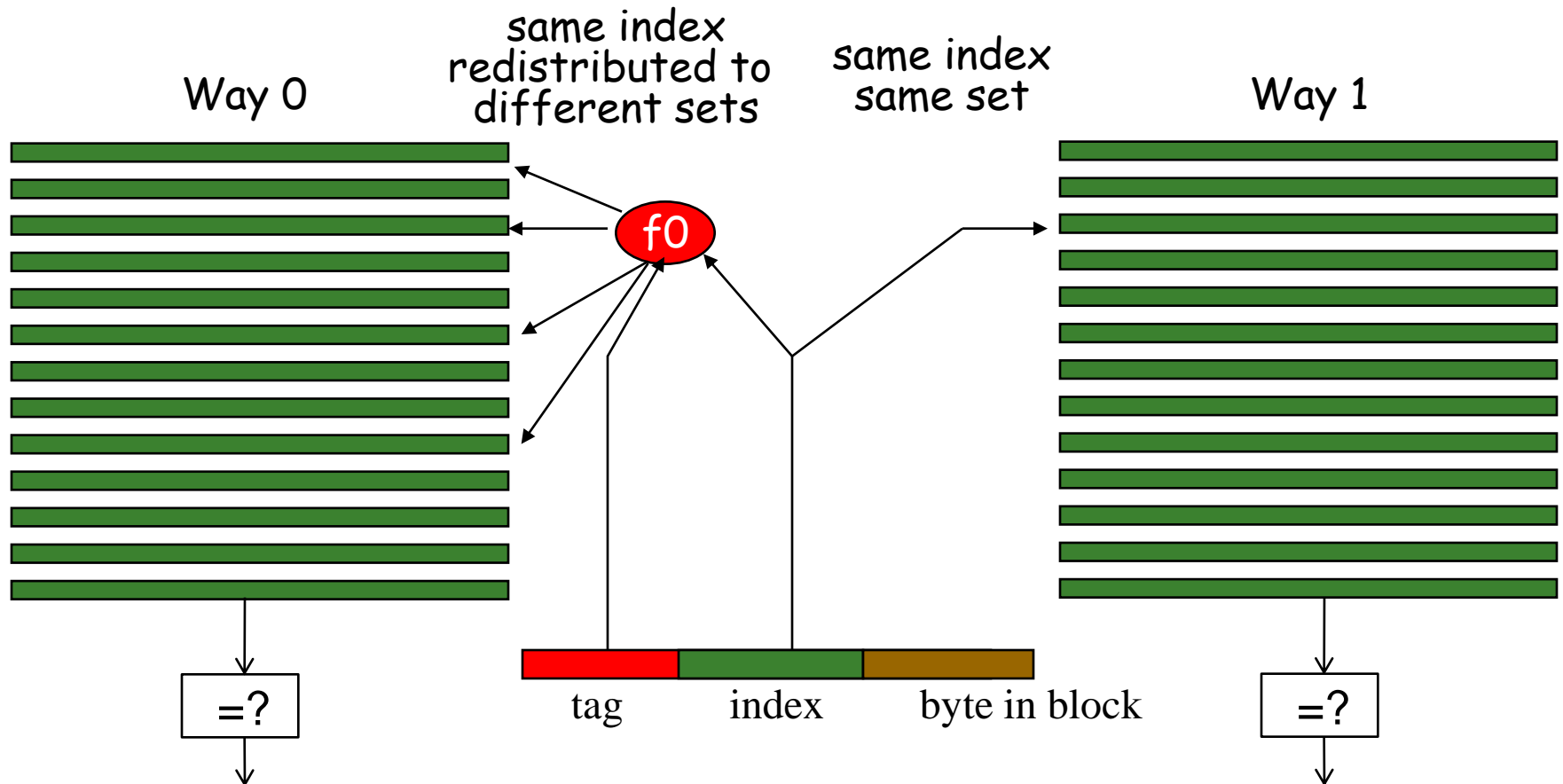
Skewed Associative Caches (I)

- Basic 2-way associative cache structure



Skewed Associative Caches (II)

- Skewed associative caches
 - Each bank has a different index function



Skewed Associative Caches (III)

- Idea: Reduce conflict misses by using different index functions for each cache way
- Benefit: indices are randomized
 - Less likely two blocks have same index
 - Reduced conflict misses
 - May be able to reduce associativity
- Cost: additional latency of hash function
- Seznec, “A Case for Two-Way Skewed-Associative Caches,” ISCA 1993.

Improving Hit Rate via Software (I)

- Restructuring data layout
- Example: If column-major
 - $x[i+1,j]$ follows $x[i,j]$ in memory
 - $x[i,j+1]$ is far away from $x[i,j]$

Poor code

```
for i = 1, rows
  for j = 1, columns
    sum = sum + x[i,j]
```

Better code

```
for j = 1, columns
  for i = 1, rows
    sum = sum + x[i,j]
```

- This is called **loop interchange**
- Other optimizations can also increase hit rate
 - Loop fusion, array merging, ...
- What if multiple arrays? Unknown array size at compile time?

More on Data Structure Layout

```
struct Node {  
    struct Node* node;  
    int key;  
    char [256] name;  
    char [256] school;  
}  
  
while (node) {  
    if (node→key == input-key) {  
        // access other fields of node  
    }  
    node = node→next;  
}
```

- Pointer based traversal (e.g., of a linked list)
- Assume a huge linked list (1M nodes) and unique keys
- Why does the code on the left have poor cache hit rate?
 - “Other fields” occupy most of the cache line even though rarely accessed!

How Do We Make This Cache-Friendly?

```
struct Node {  
    struct Node* node;  
    int key;  
    struct Node-data* node-data;  
}  
  
struct Node-data {  
    char [256] name;  
    char [256] school;  
}  
  
while (node) {  
    if (node→key == input-key) {  
        // access node→node-data  
    }  
    node = node→next;  
}
```

- Idea: separate frequently-used fields of a data structure and pack them into a separate data structure
- Who should do this?
 - ❑ Programmer
 - ❑ Compiler
 - Profiling vs. dynamic
 - ❑ Hardware?
 - ❑ Who can determine what is frequently used?

Improving Hit Rate via Software (II)

■ Blocking

- ❑ Divide loops operating on arrays into computation chunks so that each chunk can hold its data in the cache
- ❑ Avoids cache conflicts between different chunks of computation
- ❑ Essentially: Divide the working set so that each piece fits in the cache

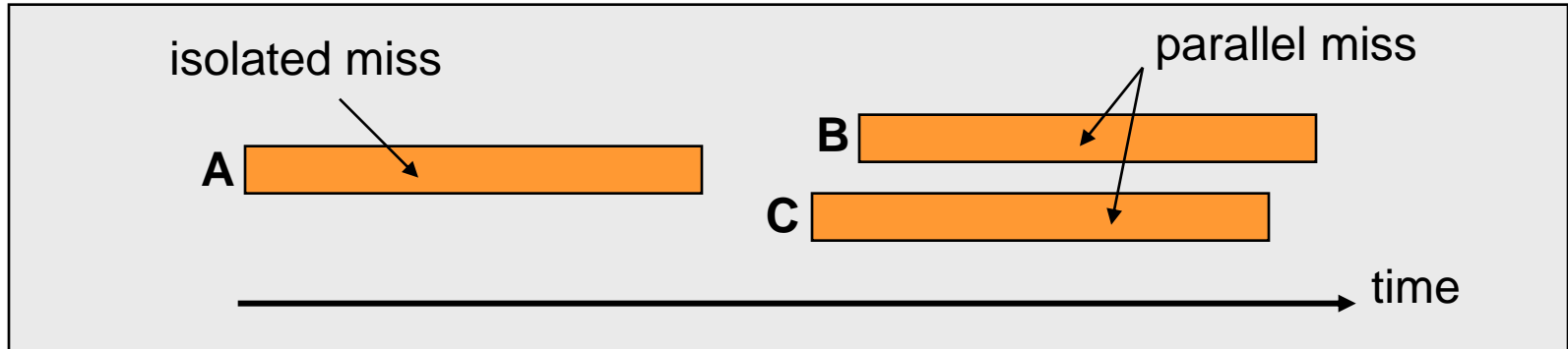
■ But, there are still self-conflicts in a block

1. there can be conflicts among different arrays
2. array sizes may be unknown at compile/programming time

Improving Basic Cache Performance

- Reducing miss rate
 - ❑ More associativity
 - ❑ Alternatives/enhancements to associativity
 - Victim caches, hashing, pseudo-associativity, skewed associativity
 - ❑ Better replacement/insertion policies
 - ❑ Software approaches
- Reducing miss latency/cost
 - ❑ Multi-level caches
 - ❑ Critical word first
 - ❑ Subblocking/sectoring
 - ❑ Better replacement/insertion policies
 - ❑ Non-blocking caches (multiple cache misses in parallel)
 - ❑ Multiple accesses per cycle
 - ❑ Software approaches

Memory Level Parallelism (MLP)



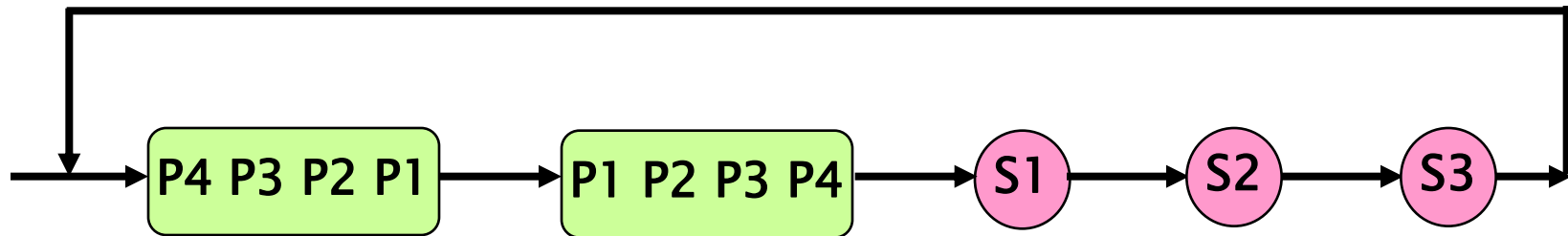
- ❑ Memory Level Parallelism (MLP) means generating and servicing multiple memory accesses in parallel [Glew' 98]
- ❑ Several techniques to improve MLP (e.g., out-of-order execution)
- ❑ MLP varies. Some misses are isolated and some parallel

How does this affect cache replacement?

Traditional Cache Replacement Policies

- ❑ Traditional cache replacement policies try to reduce miss count
- ❑ **Implicit assumption**: Reducing miss count reduces memory-related stall time
- ❑ Misses with varying cost/MLP **breaks** this assumption!
- ❑ Eliminating an isolated miss helps performance more than eliminating a parallel miss
- ❑ Eliminating a higher-latency miss could help performance more than eliminating a lower-latency miss

An Example



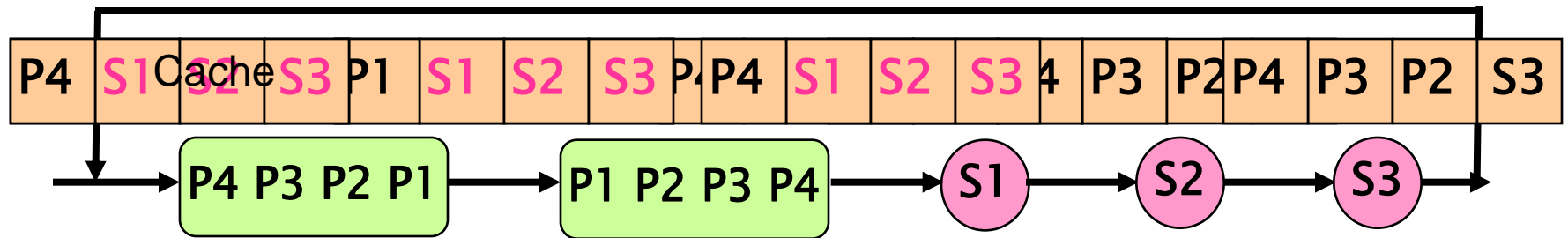
Misses to blocks P1, P2, P3, P4 can be parallel
Misses to blocks S1, S2, and S3 are isolated

Two replacement algorithms:

1. Minimizes miss count (Belady's OPT)
2. Reduces isolated miss (MLP-Aware)

For a fully associative cache containing 4 blocks

Fewest Misses \neq Best Performance



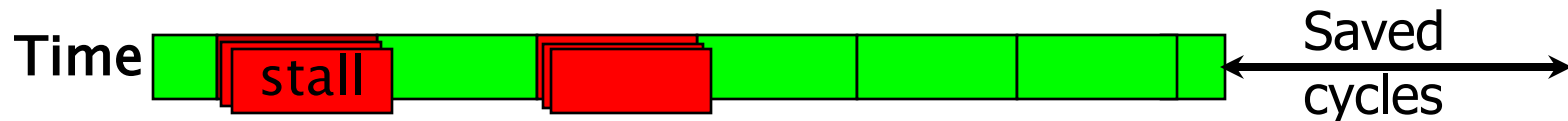
Hit/Miss H H H M H H H H M M M



Misses=4
Stalls=4

Belady's OPT replacement

Hit/Miss H M M M H M M M H H H



Misses=6
Stalls=2

MLP-Aware replacement

MLP-Aware Cache Replacement

- How do we incorporate MLP into replacement decisions?
- Qureshi et al., “A Case for MLP-Aware Cache Replacement,” ISCA 2006.
 - Required reading for this week