

18-447: Computer Architecture

Lecture 23: Caches

Prof. Onur Mutlu

Carnegie Mellon University

Spring 2013, 3/29/2013

Reminder: Homework 5

- Due April 3
- Topics: Vector processing, VLIW, Virtual memory, Caching

Reminder: Lab Assignment 5

- Lab Assignment 5
 - Due Friday, April 5
 - Modeling caches and branch prediction at the microarchitectural level (cycle level) in C
 - Extra credit: Cache design optimization
 - Size, block size, associativity
 - Replacement and insertion policies
 - Cache indexing policies
 - Anything else you would like
- TAs will go over the baseline simulator in lab sessions

Last Lecture

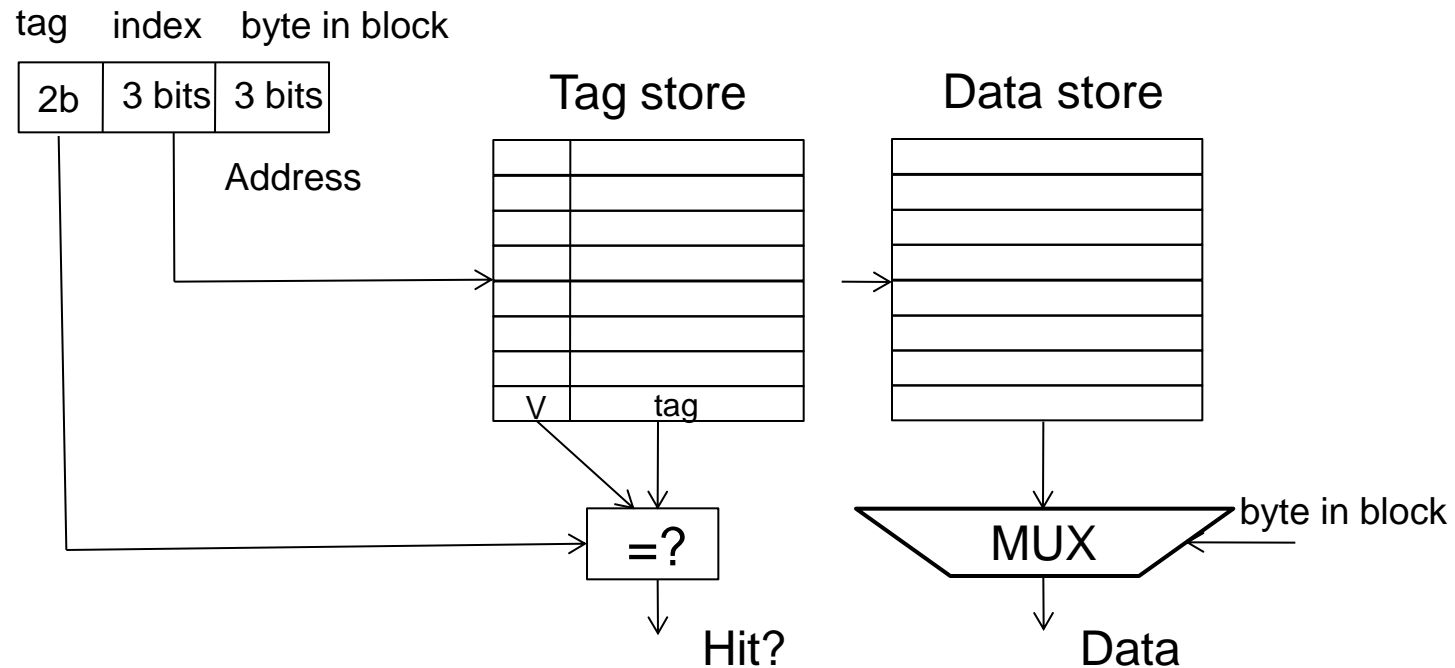
- The memory hierarchy
- Caches start
 - Structure
 - Associativity

Today

- More (and more advanced) caching

Review: Direct-Mapped Cache Structure

- 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**



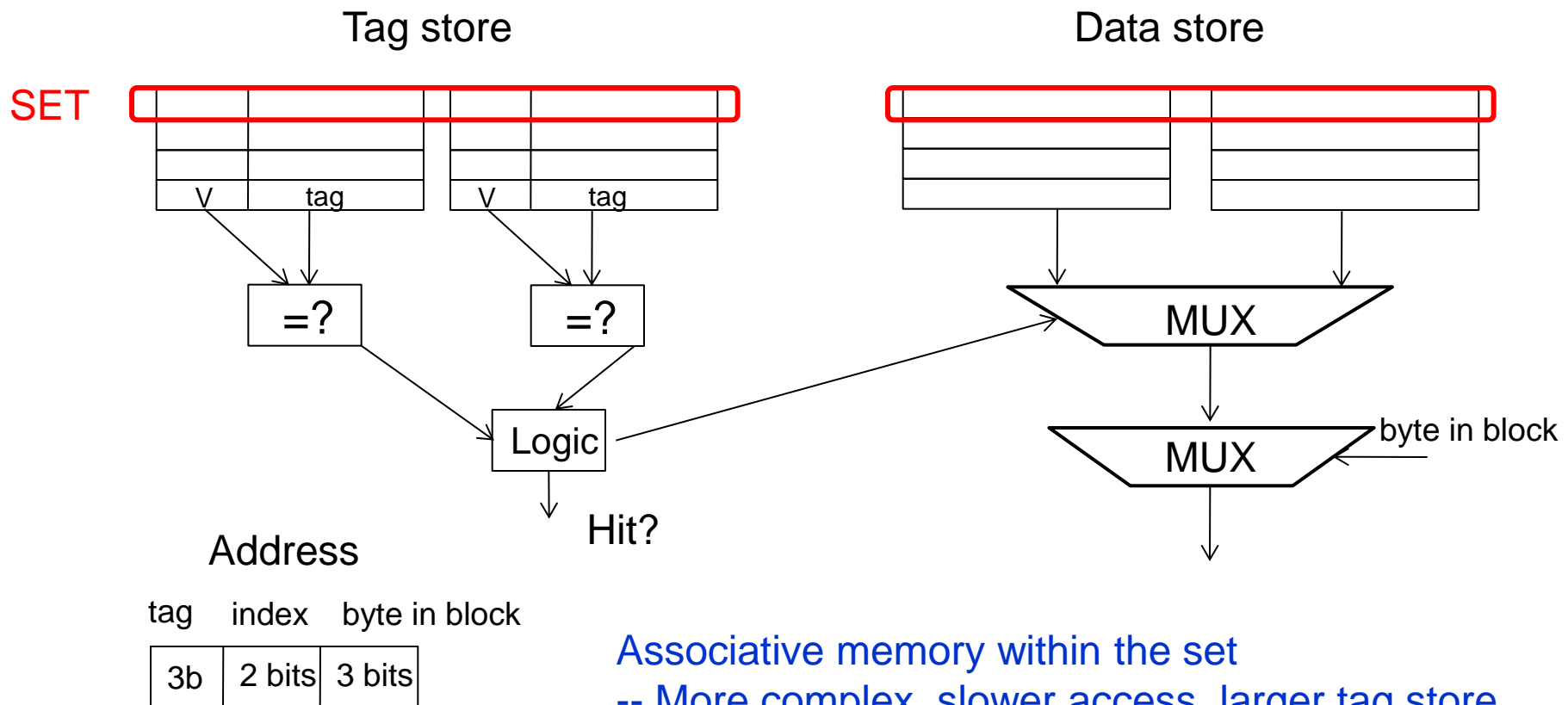
- **Addresses with same index contend for the same location**
 - **Cause conflict misses**

Review: Problem with Direct-Mapped

- **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 → 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, ... → conflict in the cache index
 - All accesses are **conflict misses**

Review: Set Associativity

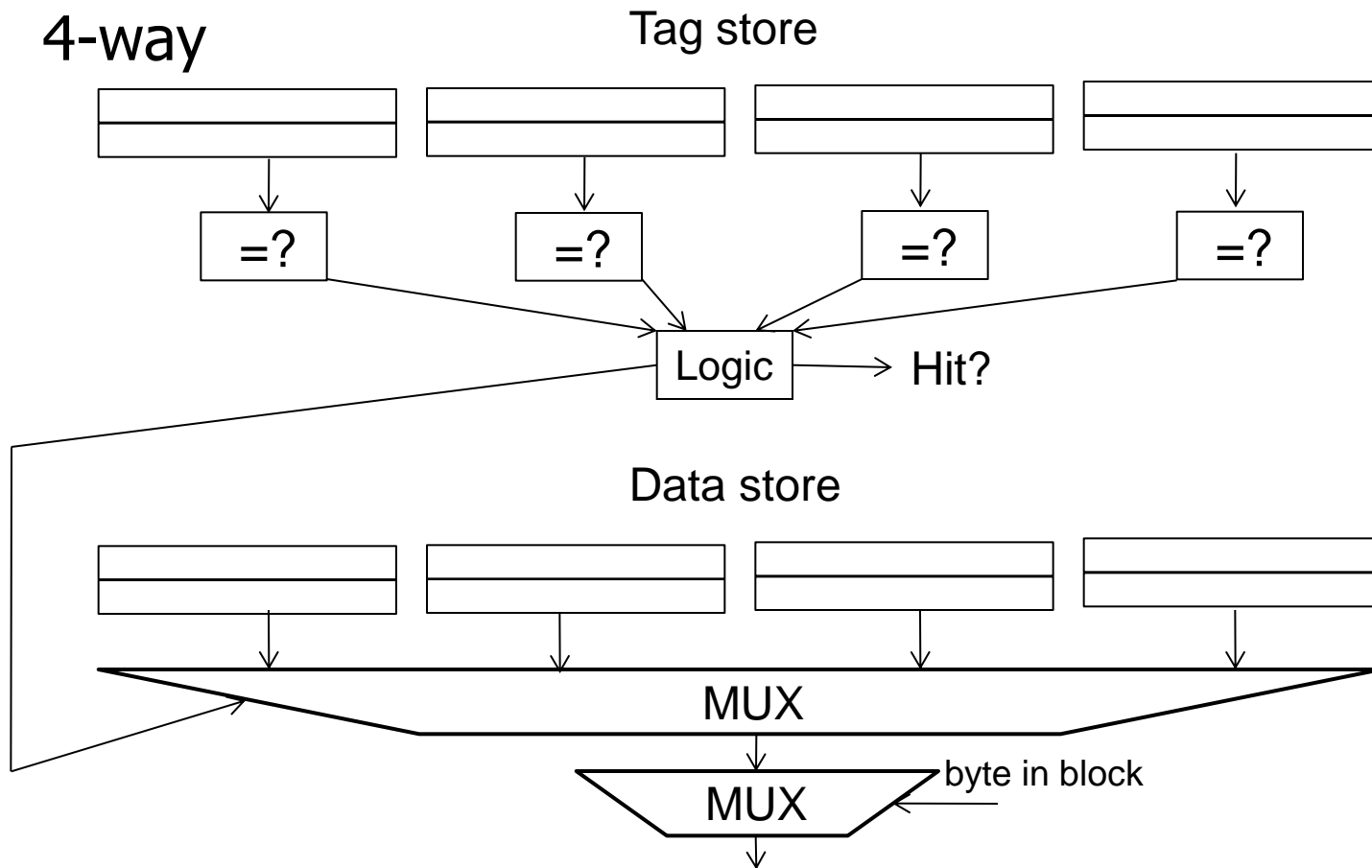
- Addresses 0 and 8 always conflict in direct mapped cache
- Instead of having one column of 8, have 2 columns of 4 blocks



Associative memory within the set
 -- More complex, slower access, larger tag store
 + Accommodates conflicts better (fewer conflict misses)

Review: Higher Associativity

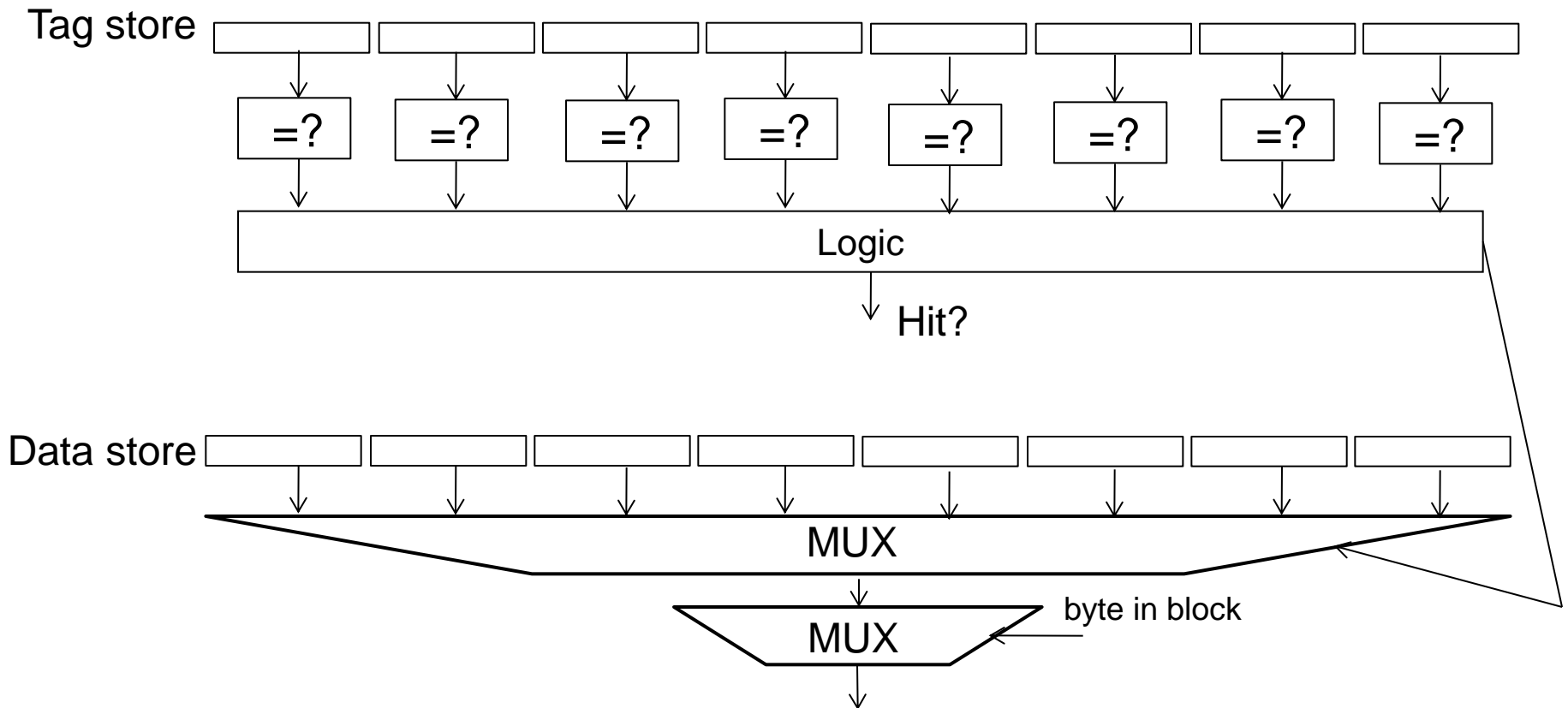
■ 4-way



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

Review: Full Associativity

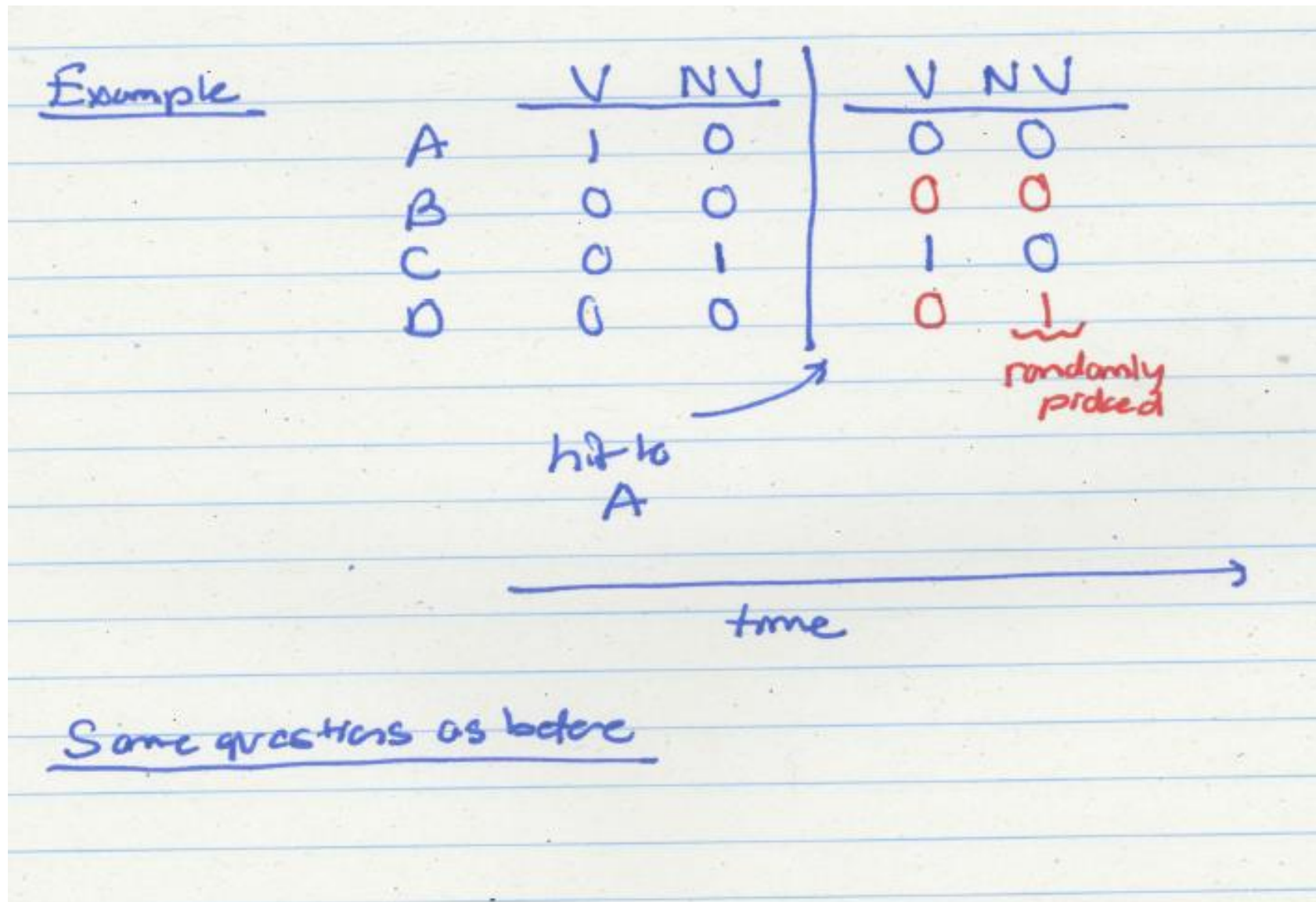
- Fully associative cache
 - A block can be placed in any cache location



Review: 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

Victim/Next-Victim Example



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

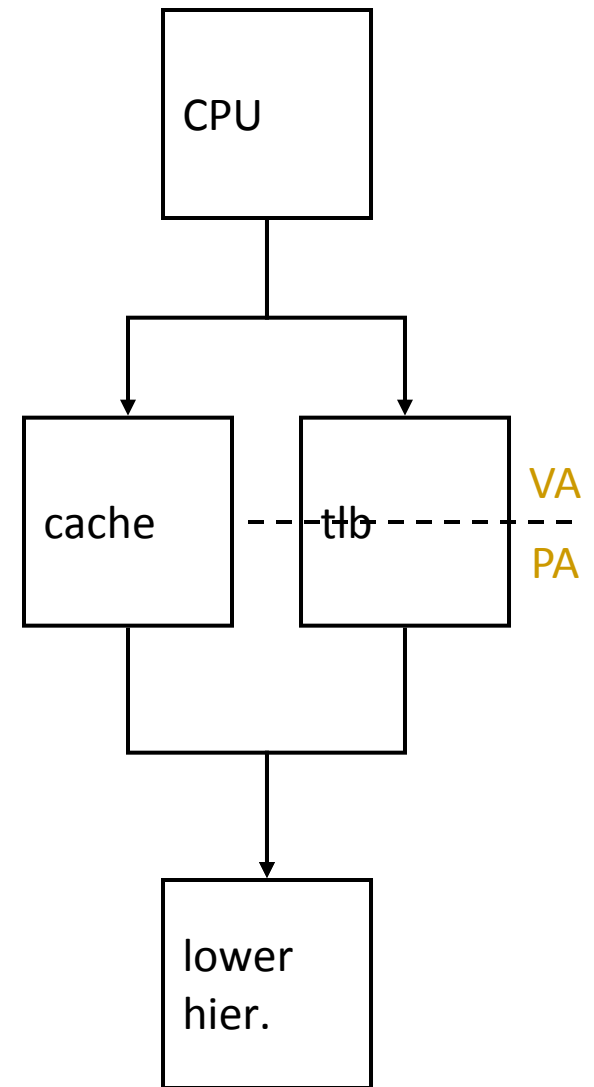
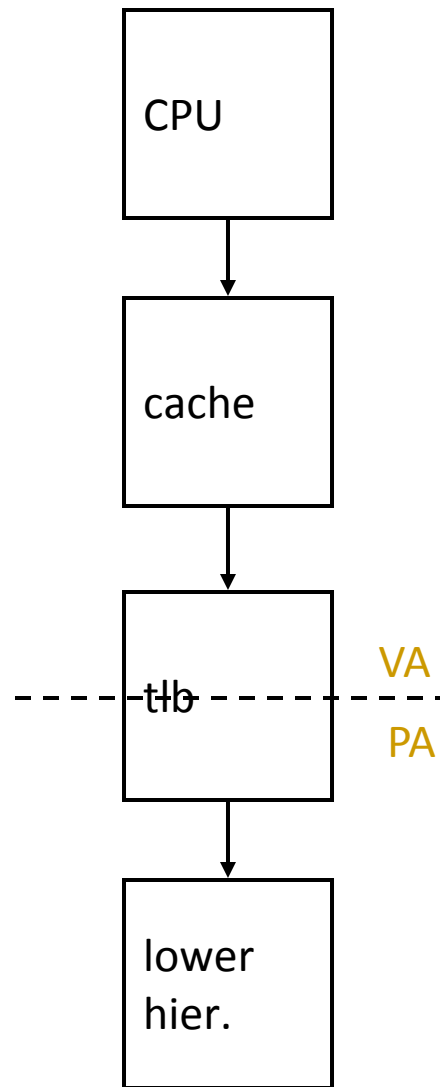
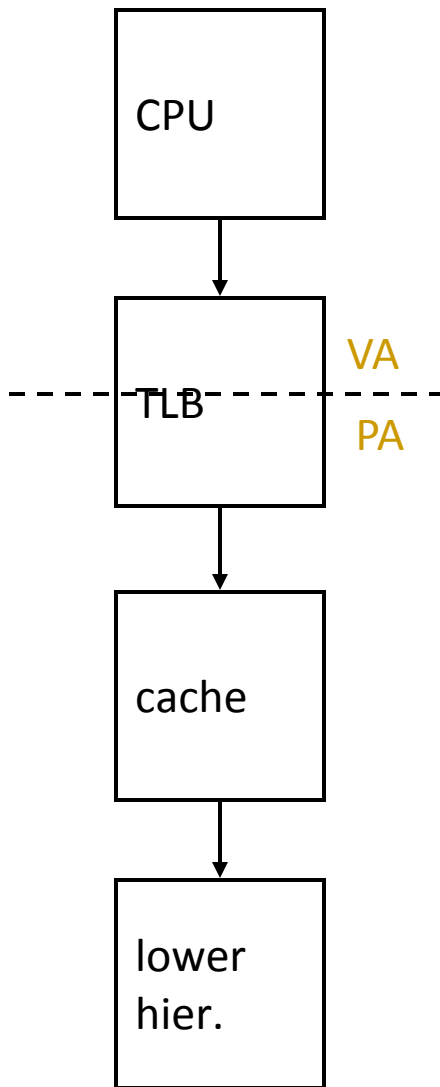
Address Translation and Caching

- 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

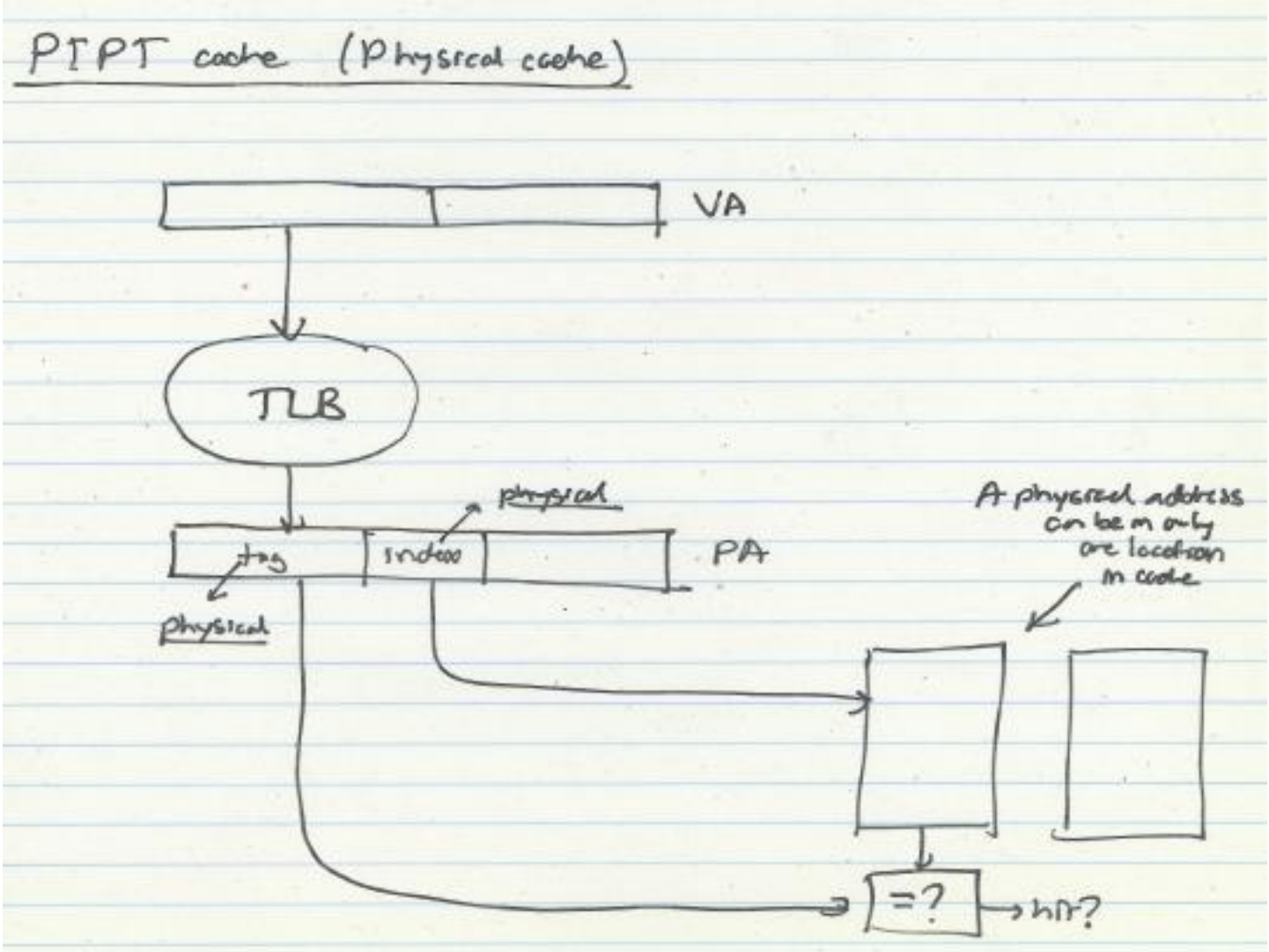


physical cache

virtual (L1) cache

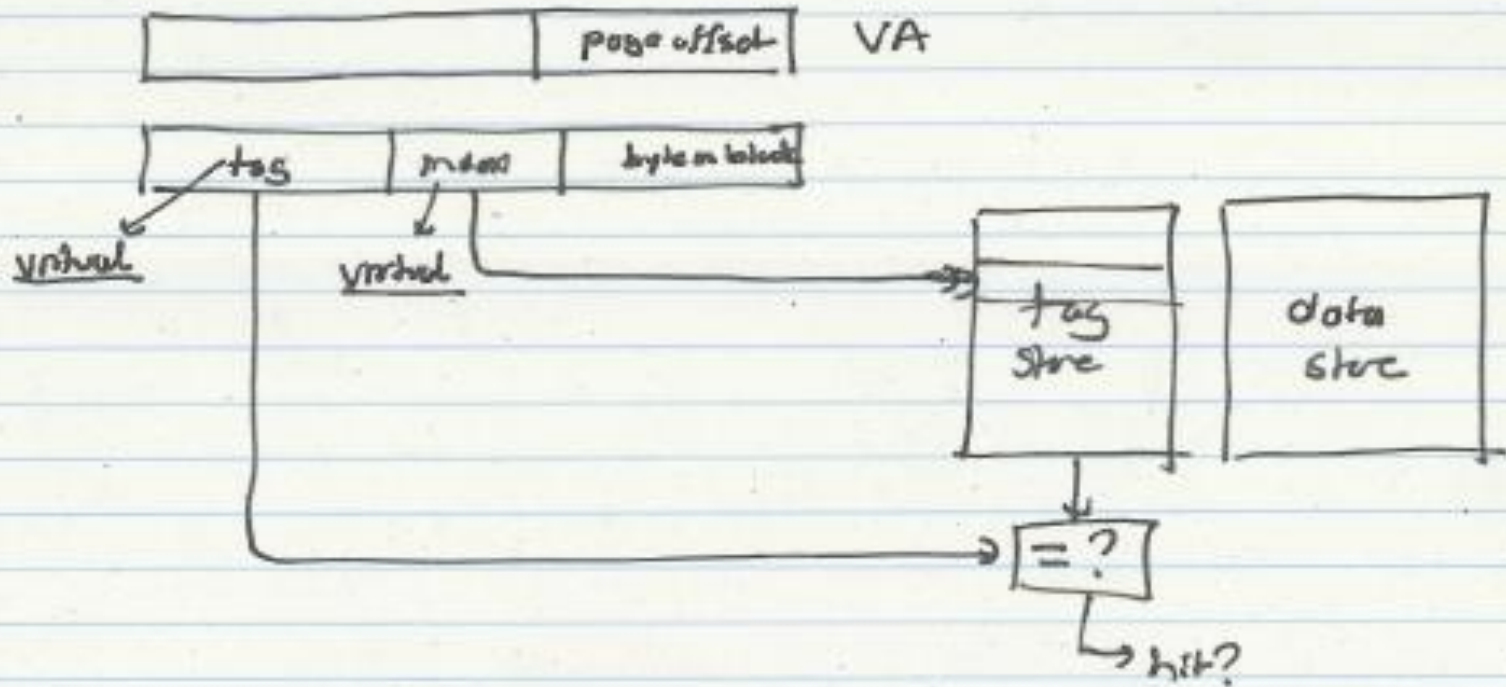
virtual-physical cache 25

Physical Cache



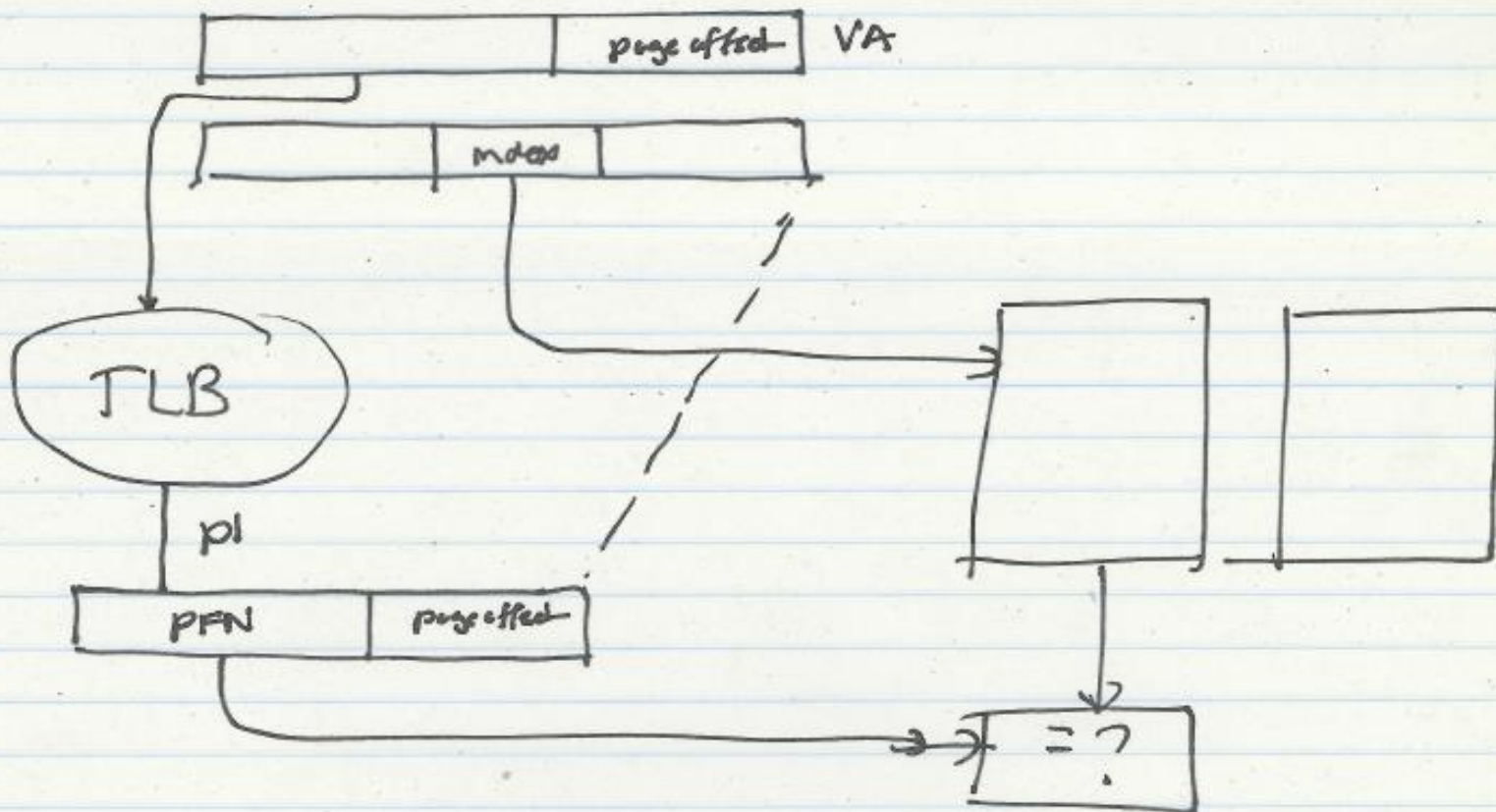
Virtual Cache

VINT cache (Virtual Cache)



Virtual-Physical Cache

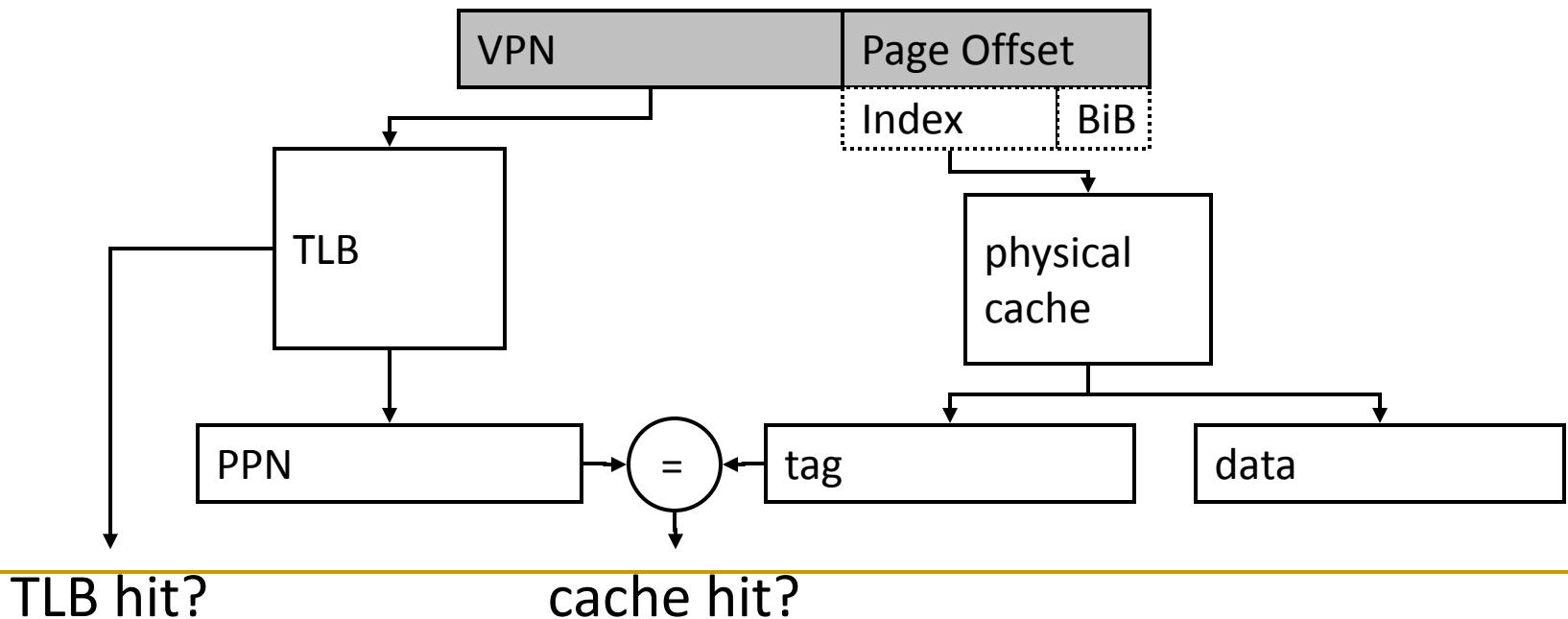
VIPT cache



Where can the same physical address be in the cache?

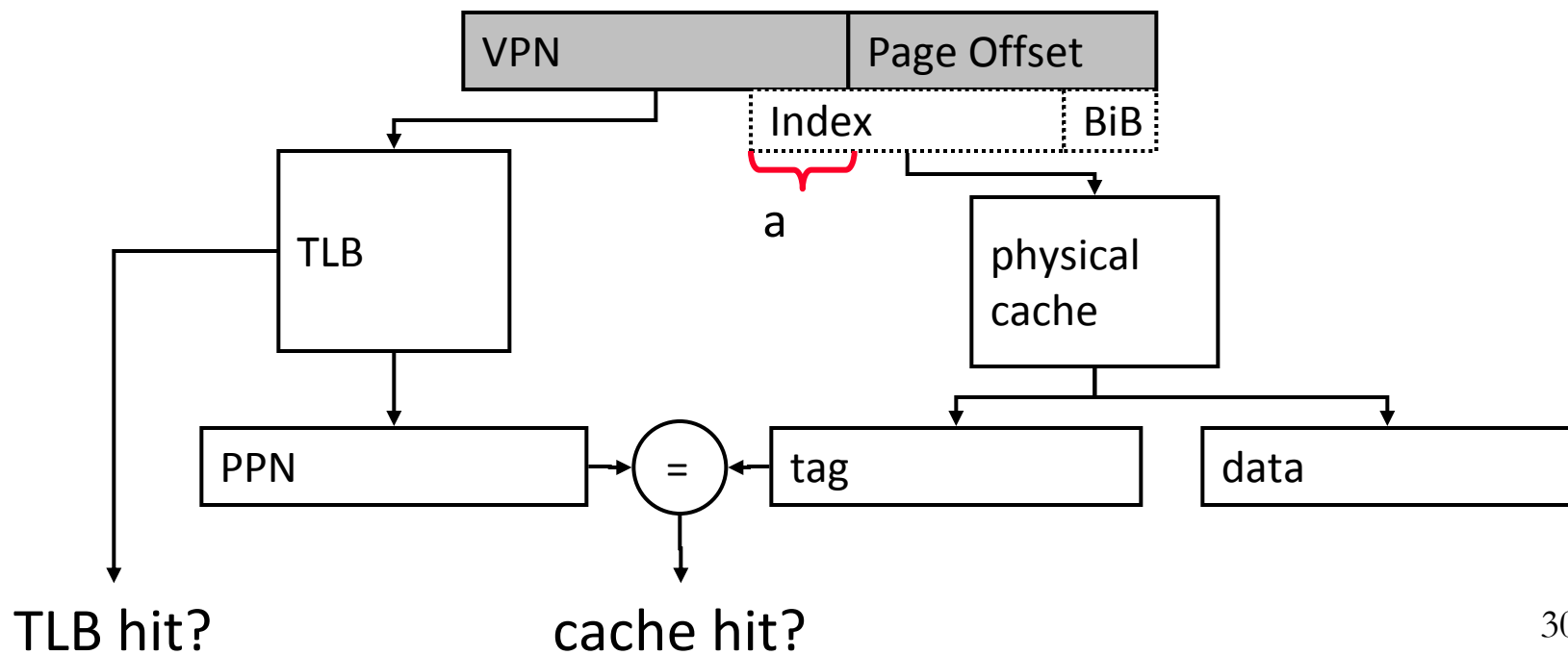
Virtually-Indexed Physically-Tagged

- 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.