# ECE 5730
# Memory Systems

## Spring 2009

## Non-Blocking Caches
## Cache Content Management

Cornell University

# Announcements

- **Quiz on Tuesday**

- **Quiz 1**
  - **Average = 8.6**
  - **Scores are in the Gradebook on Blackboard**
  - Remind me to hand them back at the end of class!

*see lecture 3*

# Non-Blocking Caches

- A *blocking cache* stalls the pipeline on a cache miss

- A *non-blocking cache* permits additional cache accesses on a miss
  - Proposed by [Kroft81]; refined by [Farkas94], [Farkas97]
  - *Hit-under-miss*: the next miss causes a stall
  - *Miss-under-miss*: multiple misses may be present up to a limit (and then a stall)

- *Memory-level parallelism*: A miss-under-miss cache coupled with a parallel lower-level memory system
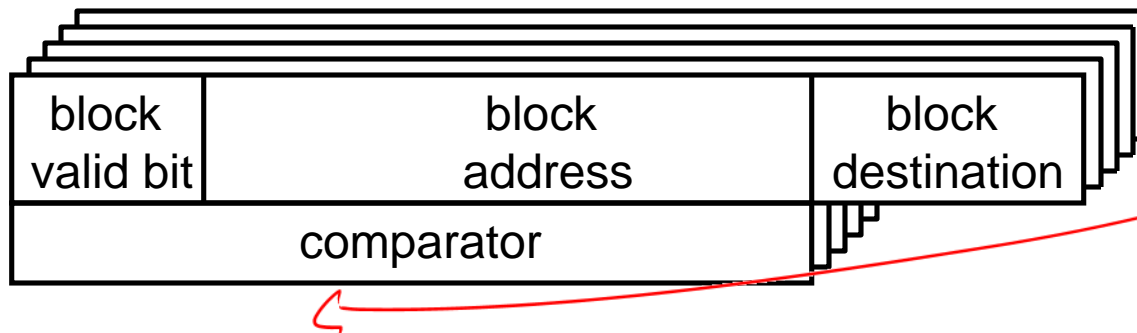
# Miss-Under-Miss Cache Design

- A *primary miss* occurs when there are no outstanding misses to the same cache line

- Otherwise, a *secondary miss* occurs and is merged with the primary miss (not issued to the next level of the memory hierarchy)

- Hardware required for a miss-under-miss design with an out-of-order processor
  - Miss Status Holding Registers (MSHRs)
  - Address Stack

*keeps track of parallel fetching*

# ↳ MSHRs

- **When a miss occurs, the MSHRs are looked up to determine if the cache block is already being fetched**

- **Each MSHR holds information for a primary miss**
  - **Block (line) address**
  - **Destination:  whether block corresponds to L1 Icache, L1 Dcache, or neither (uncached)**

| block valid bit | block address | block destination |
|---|---|---|
| comparator | | |

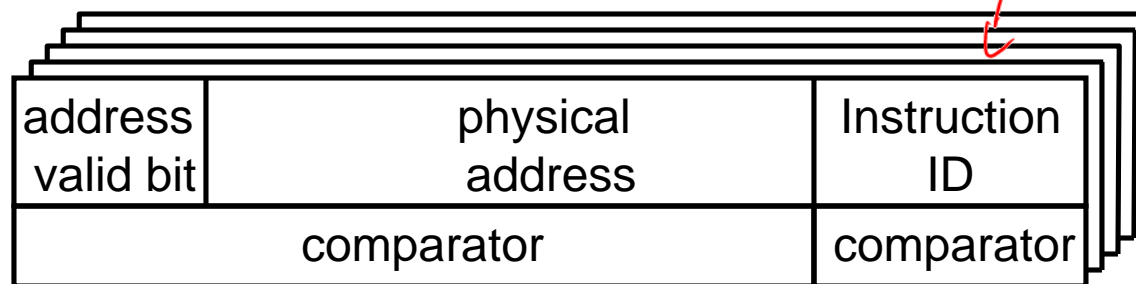- **If an MSHR hit, then a secondary miss has occurred (primary miss to the same line is outstanding)**

*multiple in-flight fetches, don't know order of return of data*

*{ you need to know what address came back from lower mem-hierarchy levels to make sure you associate the data w/ the right entry*

# Address Stack

- **Resolves primary and secondary data cache misses when a cache block is returned**

- **Each address stack entry holds information for a primary or secondary miss**
  - **Physical address**
  - **Instruction ID (e.g., the PC)**

points to the instruction to be replayed

| address valid bit | physical address | Instruction ID |
|---|---|---|
| | comparator | comparator |

# Address Stack

- **If the ~~address~~ _valid_ bit is set and the physical address matches that of the returned block, the instruction is replayed in the pipeline**

  _say we miss on a load. we get the data back, then we replay the load in the pipeline_

- **If an instruction in the address stack follows a mispredicted branch or exception, that entry's valid bit is reset**

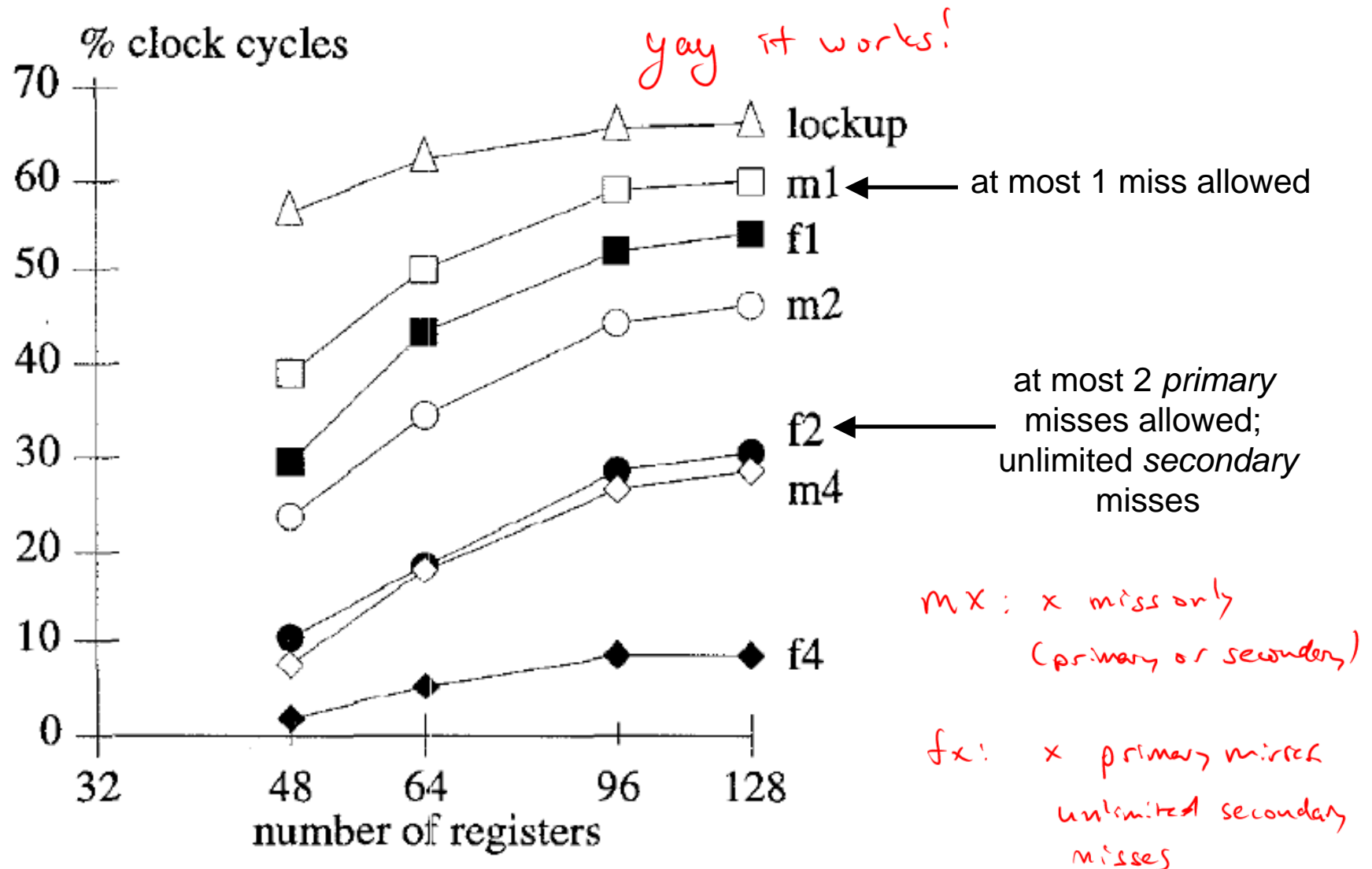  - **Instruction will not be found and thus not replayed**

  load &larr; primary miss, allocate MSHR

  branch &larr; mispredicted NOT taken, actually taken

  load &larr; secondary miss, should _NOT_ be replayed

  on this mispredicted branch, we reset the valid bits on all instruction ID's after the branch, so as to not replay the second load.

# Miss-Under-Miss Effectiveness



yay it works!

lockup
m1 ← at most 1 miss allowed
f1
m2
f2 ← at most 2 *primary* misses allowed; unlimited *secondary* misses
m4
f4

mx: x miss only (primary or secondary)

fx: x primary misses unlimited secondary misses

x MSHR's
∞ address stack entries

(a) Percentage of (simulated) clock cycles in which a load-miss induced stall was in effect

[Farkas97]

# Cache Content Management

- ## Partitioning heuristics

  *poorly worded, disregard*

  – ~~Which items get placed store where (cache level, cache sets/ways, buffers)~~

  what items do I put where?

  L1 vs L2, buffers, which set, etc

- ## Fetching heuristics

  – **When to bring an item into the cache**

- ## Locality optimizations

  – **Change layout or ordering to improve reuse**

  Decisions can be made at runtime (*on-line heuristics*),
  at design/compile time (*off-line heuristics*),
  or a combination of the two (*combined approaches*)

# On-line Partitioning Heuristics

- **Replacement policies** LRU, FIFO, etc
  MRU

- **Write policies** Write back vs Write through

  Simplistic, there are more complicated schemes

- **Cache/not-cache strategies**

  → may want to make more intelligent, more complex
  decisions about whether or not to store
  an item.

# Replacement Policy

- **Determines which item is replaced in the cache (e.g., for associative caches)**

*(ake MRU)*

- **Static policies** *not most recently used (don't kill the most recently used)*
  *less overhead than LRU, similar idea*

  *least recently used* →
  - LRU, NMRU, random, etc
  - Replacement policy is fixed at design time *(prob in hardware)*

- **Dynamic policies**
  - Policy changes with the workload characteristics
  - Recently proposed for microprocessor caches
  - Example: Set Dueling *lawl*

  *set A   set B*

# Set Dueling

- **Proposed for L2 caches by [Qureshi07]**

- **Replacement policy is LRU**
  - **The block marked as LRU is replaced**

- **But the LRU *position* in which a new block is inserted is *dynamic***

  *by definition, the block is the most recently used*

  - *LRU policy*: incoming block is placed into the MRU position (traditional LRU approach)

    *choose*

  - *Bimodal Insertion Policy (BIP)*: most blocks are placed into the *LRU* position → *set the new block as the least recently used*
    - These are replaced unless referenced again right away
  - Either LRU or BIP is chosen based on which one is doing better

    *it's not so much the LRU/MRU "position" as the state of the block.*
    *you set the appropriate bits to set LRU vs MRU*

# Why Placement into the LRU Position?
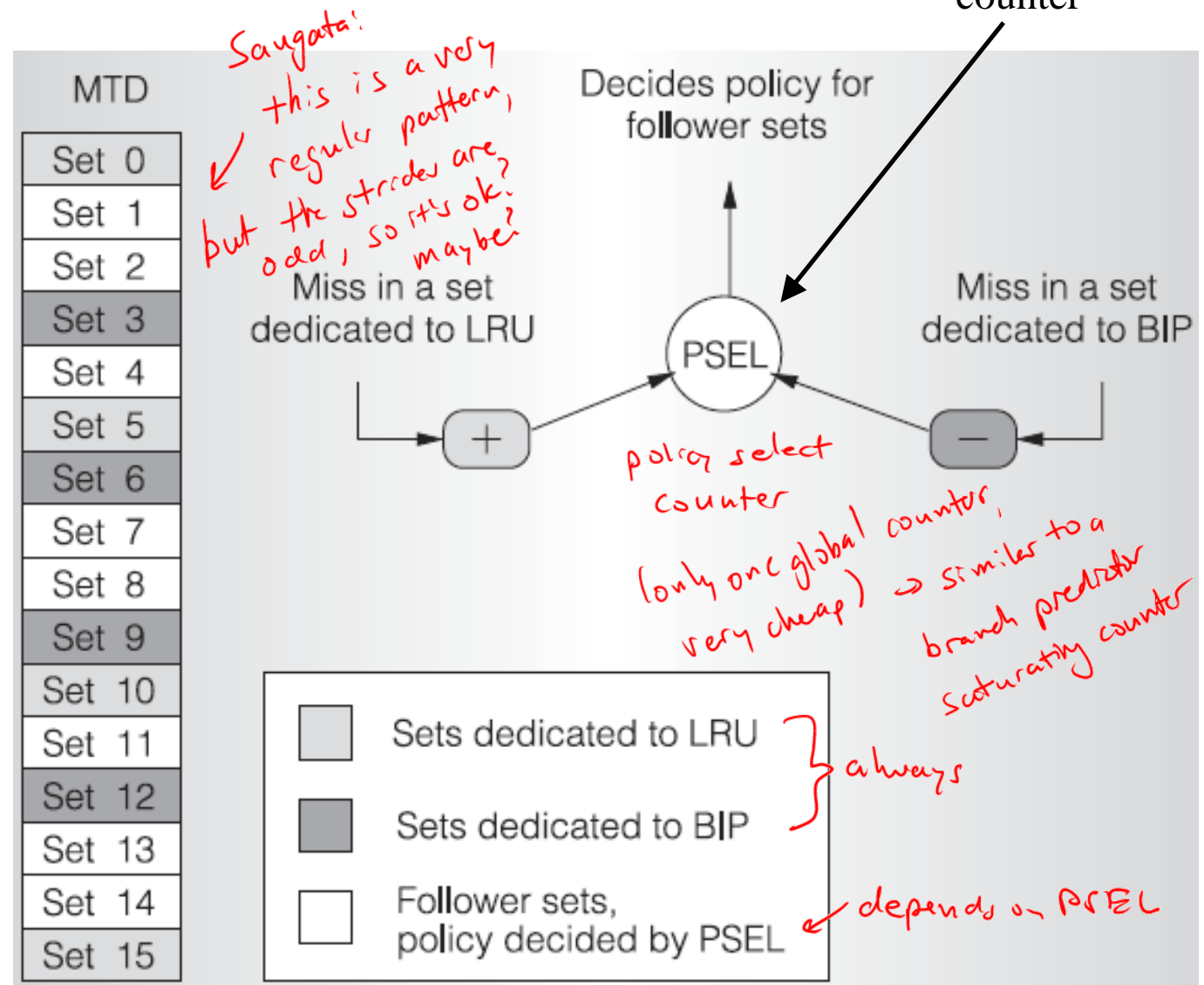
*(why BIP?)*

- **Traditional LRU algorithm works well for programs with good temporal locality**

- **For program with a working set larger than the cache, LRU causes thrashing**
  - **Blocks continually replace each other**     *oh noes!*

- **Can retain some portion of the working set by inserting most blocks into the LRU position**
  - **These will be quickly replaced**
  - **Other blocks in the same set will be retained**
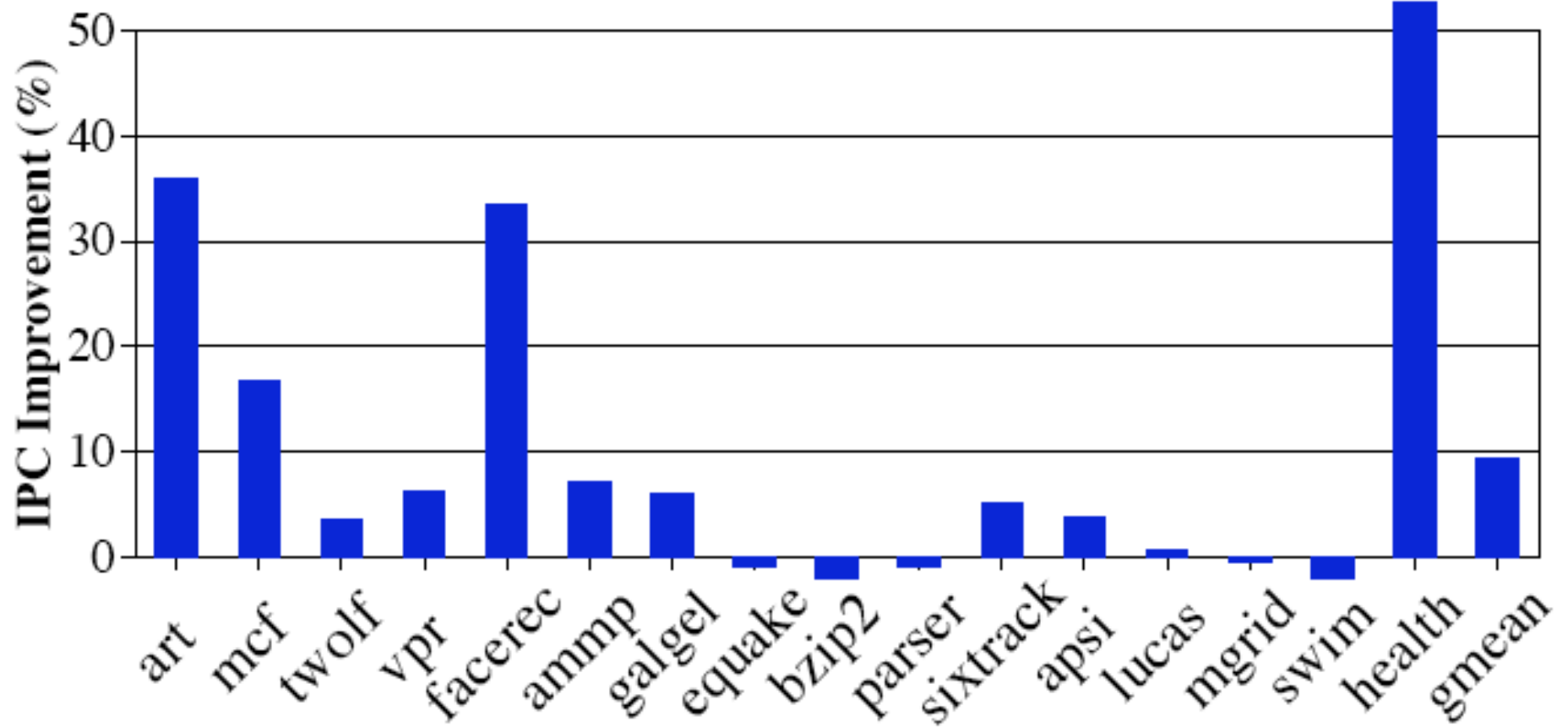
*can get up to 20% improvement with some workloads.*

# Set Dueling

- **Some sets use LRU, others BIP**
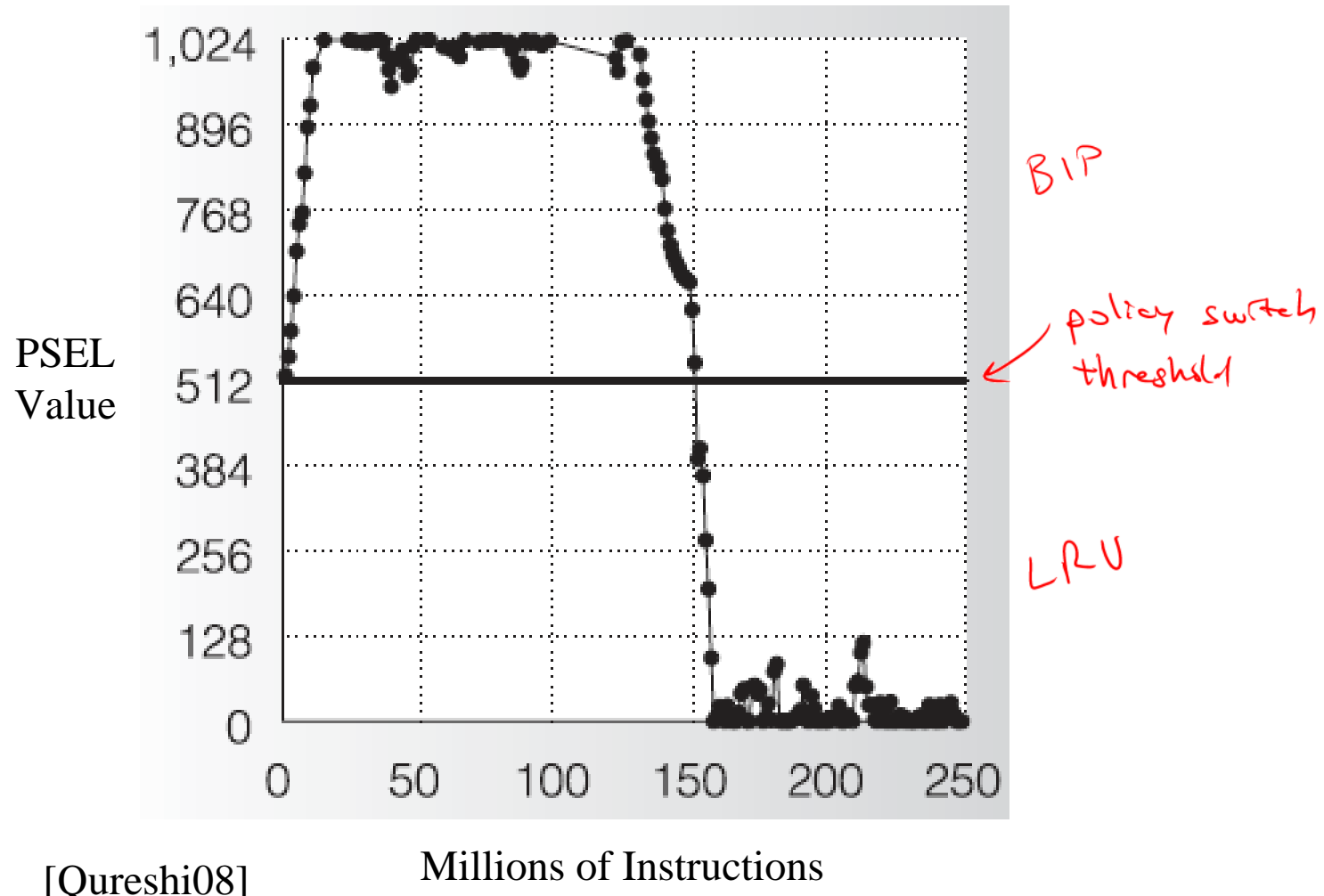
- ***Follower sets* follow the policy that does best**



[Qureshi08]

10 bit saturating counter

MTD

Set 0
Set 1
Set 2
Set 3
Set 4
Set 5
Set 6
Set 7
Set 8
Set 9
Set 10
Set 11
Set 12
Set 13
Set 14
Set 15

Decides policy for follower sets

Miss in a set dedicated to LRU

Miss in a set dedicated to BIP

PSEL

Sets dedicated to LRU
Sets dedicated to BIP
Follower sets, policy decided by PSEL

*Handwritten annotations:*
Saugata: this is a very regular pattern, but the strides are odd, so it's ok? maybe?

policy select counter (only one global counter, very cheap) → similar to a branch predictor saturating counter

} always

← depends on PSEL

# Performance Improvement



[Qureshi07]

# Dynamic Behavior of *ammp*



[Qureshi08]    Millions of Instructions

# Cache Write Policies

- **Cache management decisions involving writes must balance several tradeoffs**
  - **The time delay to store the data on a hit or miss**
  - **The amount of traffic generated to the next level**
    - **The future usage of the rest of the line containing the write address and the line that is replaced in the cache**

- **For write hits, need to decide whether or not to write to the next level of the hierarchy**
  - **Writethrough versus writeback**

- **For write misses, there are a number of options…**

# Cache Write Miss Options

- **Fetch-on-write?** → *and put it into the cache, then write?*
  - **Do we fetch the block that is being written?**

- **Write-allocate?** → *do we write to the cache at all?*
  - **Do we allocate space in the cache for the write?**

- **Write-before-hit?** → *can we write in parallel with the hit check to save time? assume hit*
  - **Do we write the block before checking for a cache hit?**

# Write Miss Alternatives

## Fetch-on-write?

*fetch first, then write*

*don't fetch on writes but make room and write into the cache*

| Write-allocate? | | Yes | No | | Write-before-hit? |
|---|---|---|---|---|---|
| | **Yes** | Fetch-on-write | Write-validate | No | |
| | **Yes** | Fetch-on-write | Write-validate | Yes | |
| | **No** | | Write-around | No | |
| | **No** | | Write-invalidate | Yes | |

*allocate space in cache*

[Jouppi93]   *as of 2009, he's @ HP, apparently he's super awesome*

*write anyway, before I know if it's a hit or miss, fix it later*

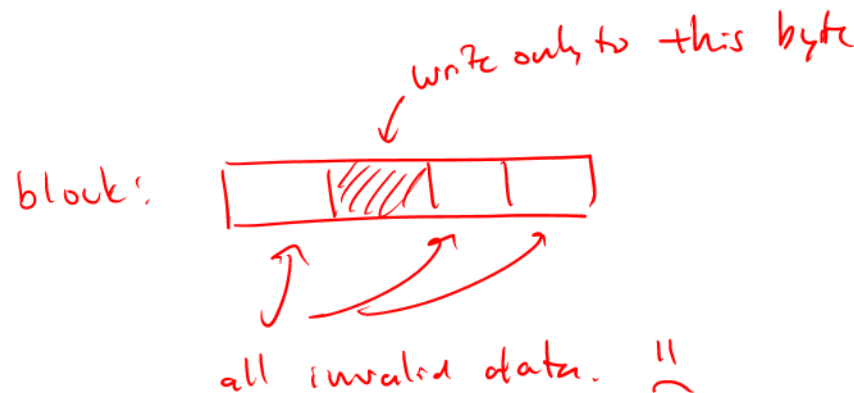*write to the next lower mem level*

Lecture 4: 19

# Write Miss Alternatives

- **Fetch-on-write**
  - Line is fetched and data is written into the cache

- **Write-validate**
  - Line is not fetched and data is written into the cache
  - Valid bits for all but the written data are turned off
  - If used with write-before-hit, entire line is invalidated on a miss

# Write Miss Alternatives

- ## Write-around
  - **On a miss, write data bypasses the cache and is written to the next level of the memory hierarchy**

- ## Write-invalidate
  - **Data is written into the cache before hit detection**
  - **If a miss occurred, line is invalidated**

# Which Approach is Best?

_additional "miss" due to fetch_

- **Fetch-on-write is simple and useful if other words in the fetched block are often accessed**

  → Spatial locality

- **Write-invalidate avoids misses when**

  _we don't fetch if we're just writing and not reading the block_

  - **neither the line containing the data being written,**
  - **or the old contents of the cache line**
  - **are read before being replaced**

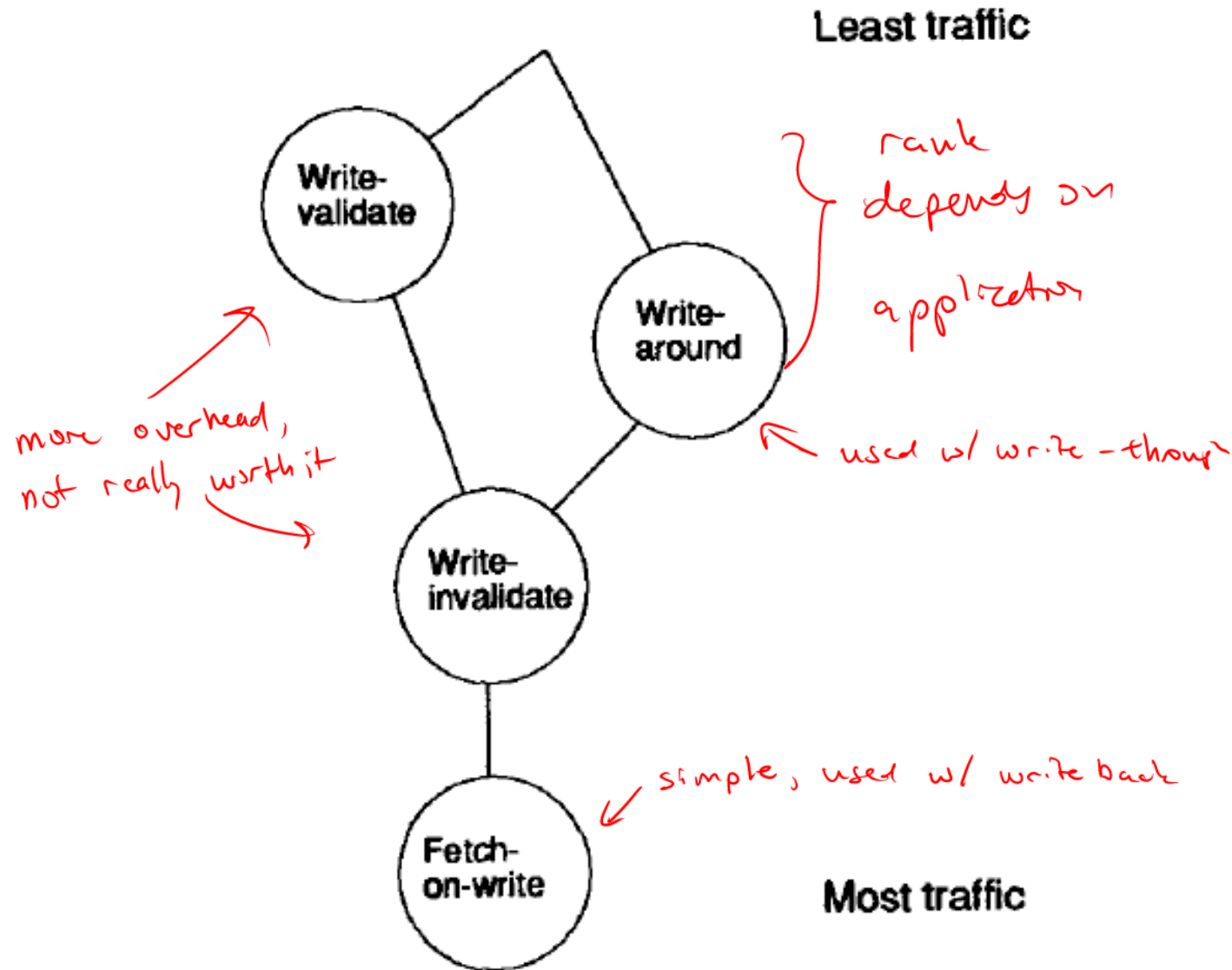- **Write-around also avoids misses when the old contents of the cache line are read next**

- **Write-validate also avoids misses when the data that was just written is accessed next**

  → writes over what we just wrote

# Relative Fetch Traffic



[Jouppi93]

# Writethrough Versus Writeback

- ## Advantages of writethrough

  - ### Simpler to maintain L1-L2 coherence  *bigger problem when L2 was off-die*

  - ### Parity instead of ECC for error tolerance (bigger issue)

    *caches have soft errors due to bit flips, fix w/ parity*

    - #### Can generate byte parity on a byte write  *← overhead not terrible*

    - #### If an error occurs, can force a miss and read from L2

    - #### For writeback, need ECC, requiring a read-modify-write operation on a byte write  *, only copy I have could be corrupted*

      *high overhead (up to 50%)*

- ## Disadvantage of writethrough  *↳ L1 + L2, 2 copies! ☺*

  - ### Higher write traffic to L2

# Coalescing Write Buffer

*"Merging" Write Buffer*

- **Hardware buffer between L1 and L2 that can accumulate write bursts to L2**   *L1 ⟷ CWB ⟷ L2*

- **Narrow writes are merged into a wider entry**
  - **Reduces the number of L2 writes**

byte writes from L1 ← *lots of byte writes get aggregated, then written*

*not super helpful, because all the byte writes have to be from the same cache line* — *addr*

| V | address | | | | | | | | | |
|---|---------|---|---|---|---|---|---|---|---|---|
|   |         |   |   |   |   |   |   |   |   |   |
|   |         |   |   |   |   |   |   |   |   |   |
|   |         |   |   |   |   |   |   |   |   |   |
|   |         |   |   |   |   |   |   |   |   |   |

up to 64-bit write to L2

*dunno what that is?*

*there are memory consistency problems, so there might be a fewer operations*

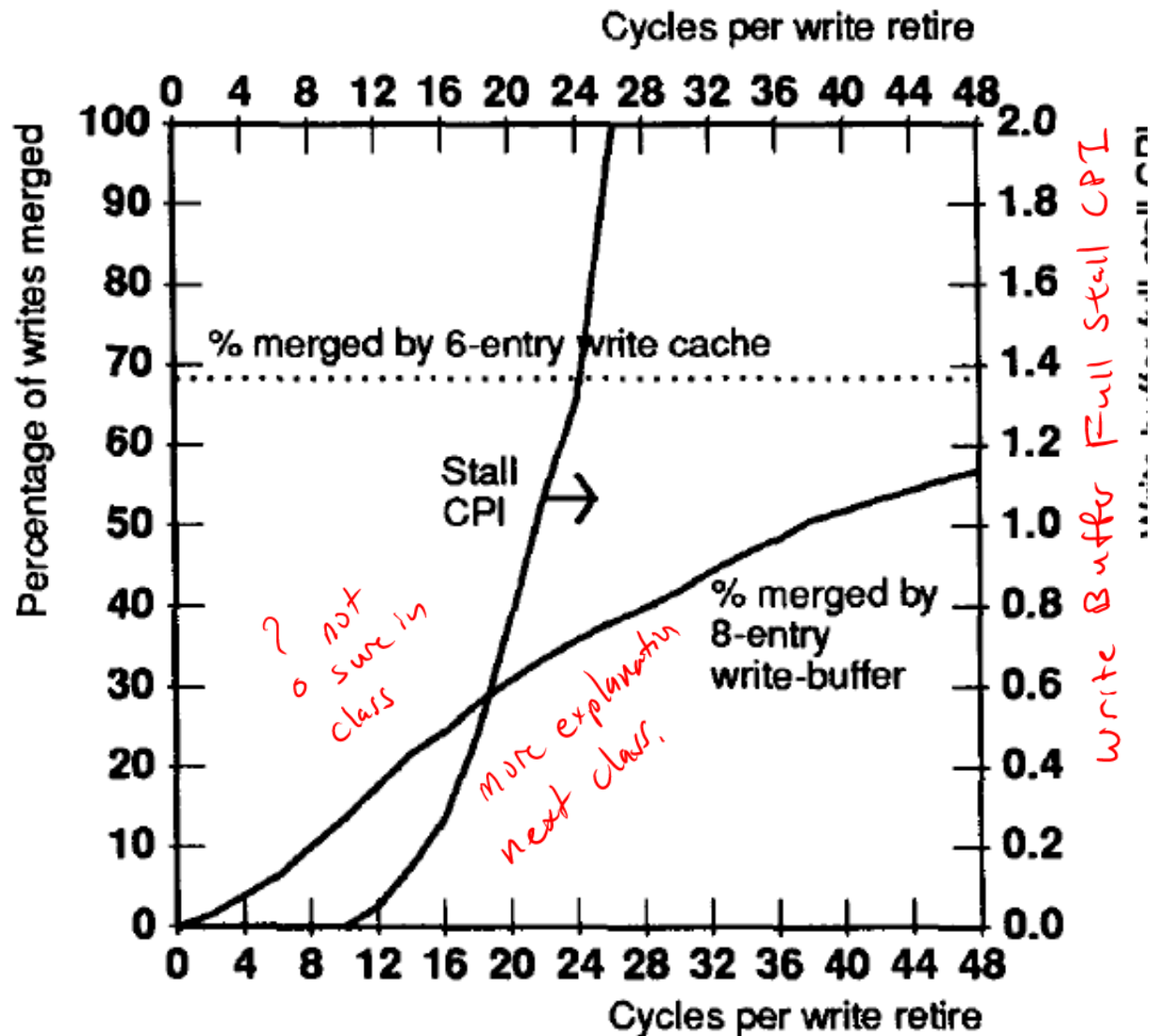# When to Empty the Write Buffer?

- **Want to retain data to increase the merging of data into a wider L2 write**

- **Want to empty the buffer to prevent it from getting full and stalling the pipeline**



[Jouppi93]

# Write Cache

Address from processor     Data from processor

Data to processor

Tags     Data     Data cache

Data to cache if
miss in data cache
but hit in write
cache or buffer

*basically a single
buffer that aggressively
empties itself if
we get over 4 entries*

empty only
if need room

| MRU entry | Tag and comparator | 8B of data | Fully-associative write cache |
|---|---|---|---|
| | Tag and comparator | 8B of data | |
| | Tag and comparator | 8B of data | |
| LRU entry | Tag and comparator | 8B of data | |

empty
aggressively

| Tag and comparator | 8B of data | Write buffer |
|---|---|---|
| Tag and comparator | 8B of data | |
| Tag and comparator | 8B of data | |

To next lower cache     To next lower cache

[Jouppi93]

# Next Time

**Cache Content Management**