# World of Predictors: Branch/Value Predictors
# CASS '18
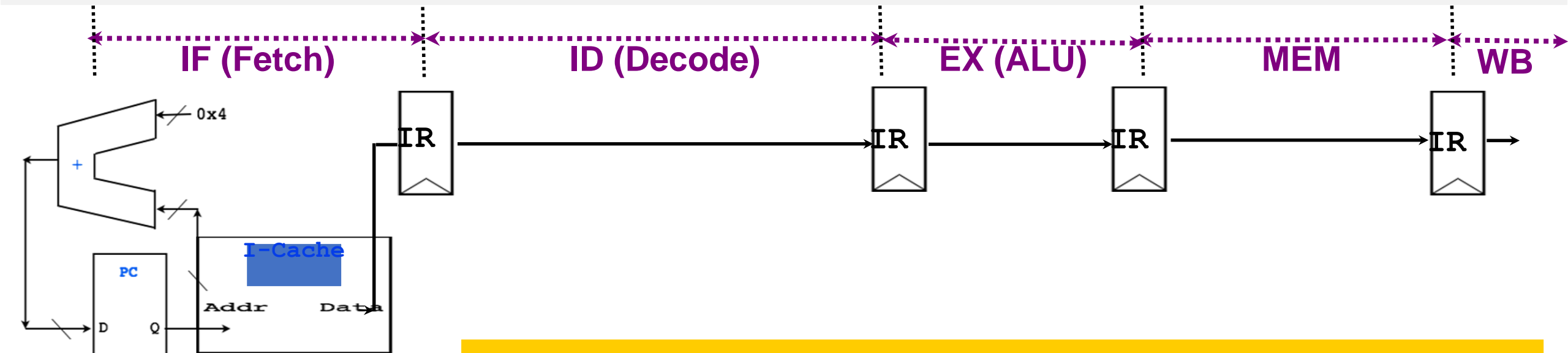
Biswa@CSE-IITK

# Basics First

0x4

PC

I-Cache

Addr    Data

IR    IR    IR    IR

**Sample Program (ISA w/o branch delay slot)**

I1:  BEQ R4,R3,25
I2:  AND R6,R5,R4
I3:  SUB R1,R9,R8

| Time: | t1 | t2 | t3 | t4 | t5 | t6 | t7 | t8 |
|-------|----|----|----|----|----|----|----|----|
| Inst |  |  |  |  |  |  |  |  |
| I1: | IF | ID | EX | MEM | WB |  |  |  |
| I2: |  | IF | ID |  |  |  |  |  |
| I3: |  |  | IF |  |  |  |  |  |
| I4: |  |  |  |  |  |  |  |  |
| I5: |  |  |  |  |  |  |  |  |
| I6: |  |  |  |  |  |  |  |  |

EX stage computes if branch is taken

If branch is taken, these instructions MUST NOT complete!

# Welcome to Branch Prediction



IF (Fetch) — ID (Decode) — EX (ALU) — MEM — WB

0x4

IR → IR → IR → IR

PC

I-Cache

Addr    Data

Branch Predictor

Predictions

A control instr?

Taken or Not Taken?

If taken, where to? What PC?

We update the PC based on the outputs of the branch predictor. If it is perfect, pipe stays full!

Dynamic Predictors: a cache of branch history

```
Time:  t1    t2    t3    t4    t5    t6    t7    t8
Inst
I1:    IF    ID    EX    MEM   WB
I2:          IF    ID
I3:                IF
I4:
I5:
I6:
```

EX stage computes if branch is taken

If we predicted incorrectly, these instructions MUST NOT complete!

# Branch Prediction

- **Idea:** Predict the next fetch address (to be used in the next cycle)

- Requires three things to be predicted at fetch stage:
  - Whether the fetched instruction is a branch
  - (Conditional) branch direction
  - Branch target address (if taken)

- Observation: Target address remains the same for a conditional direct branch across dynamic instances
  - Idea: Store the target address from previous instance and access it with the PC
  - Called Branch Target Buffer (BTB) or Branch Target Address Cache

# Static Branch Prediction

- **Always not-taken**
  - Simple to implement: no need for BTB, no direction prediction
  - Low accuracy: ~30-40%

- **Always taken**
  - No direction prediction
  - Better accuracy: ~60-70%
    - Backward branches (i.e. loop branches) are usually taken

- **Backward taken, forward not taken (BTFN)**
  - Predict backward (loop) branches as taken, others not-taken

# Static Branch Prediction

- **Profile-based**

  - Idea: Compiler determines likely direction for each branch using profile run. Encodes that direction as a hint bit in the branch instruction format.

+ Per branch prediction → accurate if profile is representative!

-- Requires hint bits in the branch instruction format

-- Accuracy depends on dynamic branch behavior:

    TTTTTTTTTTNNNNNNNNNN → 50% accuracy

    TNTNTNTNTNTNTNTNTNTN → 50% accuracy

-- Accuracy depends on the representativeness of profile input set

# Dynamic Branch Prediction

- Idea: Predict branches based on dynamic information (collected at run-time)

- Advantages

  + Prediction based on history of the execution of branches

    + It can adapt to dynamic changes in branch behavior

  + No need for static profiling: input set representativeness problem goes away

- Disadvantages

  -- More complex (requires additional hardware)

# Simplest One: Last-Time Predictor

- **Last time predictor**
  - ❑ Indicates which direction branch went last time it executed

    TTTTTTTTTTNNNNNNNNNN → 90% accuracy

- Always mis-predicts the last iteration and the first iteration of a loop branch
  - ❑ Accuracy for a loop with N iterations = (N-2)/N

+ Loop branches for loops with large number of iterations

-- Loop branches for loops will small number of iterations

   TNTNTNTNTNTNTNTNTN →   0% accuracy

*Last-time predictor CPI = [ 1 + (0.20\*0.15) \* 2 ]  = 1.06   (Assuming 85% accuracy)*

# Last-Time

# Last-time Predictor: The hardware

K bits of branch instruction address

Branch history table of $2^K$ entries, 1 bit per entry

①

Index

②

Use this entry to predict

0: predict not taken
1: predict taken

③

When branch direction resolved, go back into the table and update entry: 0 if not taken, 1 if taken

# Example: Predict!!

0xDC08:          for(i=0; i < 100000; i++)
          {
0xDC44:               if( ( i % 100) == 0 )
               tick( );

0xDC50:               if( (i & 1) == 1)
                    odd( );

          }

99.998%
Prediction
Rate

98.0%

0.0%

T

N

# Change Predictor after 2 Mistakes

# Is This Enough

- Control flow instructions (branches) are frequent
  - 15-25% of all instructions

- Problem: Next fetch address after a control-flow instruction is not determined after N cycles in a pipelined processor
  - N cycles: (minimum) branch resolution latency
  - Stalling on a branch wastes instruction processing bandwidth (i.e. reduces IPC)

- How do we keep the pipeline full after a branch?

- Problem: Need to determine the **next fetch address** when the branch is fetched (to avoid a pipeline bubble)

# Is This Enough?

- Assume a pipeline with 20-cycle branch resolution latency

- How long does it take to fetch 100 5-instruction blocks (500 instructions)?
    - Assume 1 out of 5 instructions is a branch, fetch width of 5, each 5 instruction block ends in a branch
    - 100% accuracy : 100 cycles (all instructions fetched on the correct path)
        - No wasted work
    - 99% accuracy: 100 (correct path) + 20 (wrong path) = 120 cycles
        - 20% extra instructions fetched
    - 98% accuracy: 100 (correct path) + 20 * 2 (wrong path) = 140 cycles
        - 40% extra instructions fetched
    - 95% accuracy: 100 (correct path) + 20 * 5 (wrong path) = 200 cycles
        - 100% extra instructions fetched

# Fetch Stage with BTB and Direction Prediction



Direction predictor (2-bit counters)

taken?

PC + inst size

Program Counter

Address of the current branch

hit?

Next Fetch Address

target address

Cache of Target Addresses (BTB: Branch Target Buffer)

# BTB (Why 30-bit Tag?)

Address of branch instruction
`0b0110[...]01001000`

30 bits

Branch instruction
`BNEZ R1 Loop`

Branch History Table (BHT)

4096 entries ...

## Branch Target Buffer (BTB)

| 30-bit address tag | target address |
| --- | --- |
| | |
| | |
| `0b0110[...]0010` | `PC + 4 + Loop` |
| | |

| |
| --- |
| |
| |
| **2 state bits** |
| |

"=" "Hit"

"Taken" Address ↓

"Taken" or "Not Taken" ↓

Drawn as fully associative to focus on the essentials.

In real designs, always direct-mapped.

At EX stage, update BTB/BHT, kill instructions, if necessary,

# No History based Branch Predictor

k bit

$(PC >> 2) \& (2^p - 1)$

$2^p$

Bimodal predictor: Good for biased branches

# Local History & Global History

- **Local Behavior**

  What is the predicted direction of Branch A given the outcomes of previous instances of Branch A?

- **Global Behavior**

  What is the predicted direction of Branch Z given the outcomes of *all** previous branches A, B, …, X and Y?

  *  Number of previous branches tracked limited by the history length

# Two Level Global Branch Prediction [MICRO '91]

- First level: Global branch history register (N bits)
    - The direction of last N branches
- Second level: Table of saturating counters for each history entry
    - The direction the branch took the last time the same history was seen

Pattern History Table (PHT)



GHR
(global history register)

previous one

index

00 .... 00
00 .... 01
00 .... 10

11 ....  11

- Table of saturating counters

$$k \text{ bit}$$

$$00 \dots 00$$

$$m \text{ bit}$$

$$00 \dots 01$$

GHR  $1\ 1\ \dots\ 1\ \boxed{0}$

$$00 \dots 10$$

$$2^m$$

$$11 \dots 11$$

# GHR per Branch (Gain/Loss?)

m bit

k bit

$2^p$

00 …. 00

00 …. 01

00 …. 10

1 1 ….. 1 | 0

$2^m$

11 …. 11

BHT

PHT

$2^m$

(PC >> 2) & ($2^p$ -1)

How large: k ?     Mostly K=2, m =12, how large m?

m bit

$2^p$

1 1 ..... 1 | 0

BHT

(PC % $2^p$ )

m bit

k bit

$2^p$

1 1 ..... 1    0

BHT

00 .... 00

00 .... 01

00 .... 10

11 ....  11

$2^m$

PHT

(PC >> 2) & ($2^p$ -1)

# Interference in Tables

- Sharing the PHTs between histories/branches leads to interference
    - Different branches map to the same PHT entry and modify it
    - Interference can be positive, negative, or neutral
- Interference can be eliminated by dedicating a PHT per branch
    -- Too much hardware cost
- How else can you eliminate

or reduce interference?



Figure 2: Interference in a two-level predictor.

# GShare

m bit

$$11 \ldots 1 \quad 0$$

k bit

00 .... 00

00 .... 01

00 .... 10

$2^m$

PC >>2 & $2^m$ -1

11 .... 11

For a given history and for a given branch (PC) counters are trained

# Y & P Classification [MICRO 91]



GBHR

GPHT

**GAg**

PABHR

GPHT

**PAg (SAg?)**

PABHR

PAPHT

**PAp**

- GAg: Global History Register, Global History Table
- PAg: Per-Address History Register, Global History Table
- PAp: Per-Address History Register, Per-Address History Table

# Tournament Predictor



table of 2-/3-bit counters

Final Prediction

If meta-counter MSB = 0, use $pred_0$ else use $pred_1$

| $Pred_0$ | $Pred_1$ | Meta Update |
|----------|----------|-------------|
| ✘ | ✘ | --- |
| ✘ | ✔ | Inc |
| ✔ | ✘ | Dec |
| ✔ | ✔ | --- |

# State of the art: Neural vs. TAGE

1970: Flynn
1972: Riseman/Foster
1979: Smith Predictor

1991: Two-level prediction
1993: gshare, tournament
1996: Confidence estimation
1996: Vary history length
1998: Cache exceptions

2001: Neural predictor
2004: PPM

2006: TAGE

2016: Still TAGE vs Neural

- Neural: AMD, Samsung

- TAGE: Intel?, ARM?

- Similarity
  - Many sources or "features"

- Key difference: how to combine them
  - TAGE: Override via partial match
  - Neural: integrate + threshold

- Every CBP is a cage match
  - Andre Seznec vs. Daniel Jimenez

# Spectre and Meltdown

https://meltdownattack.com/

# What about Values? And Why Values?

ILP = 1.3

ILP = 4

Predict

A B C D

Verify

A B

C

D

What is value prediction?
1. Generate a speculative value (predict)
2. Consume speculative value (execute)
3. Verify speculative value (compare/recover)

Goal: performance, i.e. expose more ILP

# Branch vs Value

- Predict result values of executing instruction
- Comparison to Branch Prediction
  - Branch Prediction
    - Single bit prediction (taken or not-taken)
    - Speed-up by solving control dependencies
    - Small penalty when prediction fails

  - Value Prediction
    - Multi-bit (32 or 64) prediction
    - Speed-up by solving data-dependencies
    - Large penalty when prediction fails

# Branch vs Value

- Most branch predictors are history-based
  - Simple predictors (e.g. last two values) have high hit-ratio.
  - SOTA predictor shows very high hit-ratio.

- Value prediction is much difficult
  - Locality of result values exists but far less than that of branch
  - About 50% executed instructions output last values and 80% instructions output one of 16 latest results [Lipasti+,96]
  - Much lower hit-ratio than branch predictors

# Branch vs Value

- Miss penalty

| Branch prediction | << | Value prediction |

- Latency of normal instruction < branch instruction

- Number of canceled instruction is much larger than that in branch prediction

# ACKS

- Some of the slides are adapted and modified from Joel Emer, Arvind, Yale Patt, John Kubiatowicz, Onur Mutlu, Krste Asanovic, Mattan Erez, Rajeev Balasubramonian, and Mainak Chaudhuri, Mikko Lipasti, and Kei Hiraki

# CASH @CASS '18:

Biswa@CSE-IITK

# CACHES
@CASS '18:



Biswa@CSE-IITK

# Memory Wall Problem

**Core** ↔ **L1**
4 Cycles

**L2**
12 Cycles

**L3**
30 Cycles

**DRAM Contr.**

100s of Cycles

**Latency wall**

**Bandwidth wall**

40 Years of Microprocessor Trend Data

# BASICS FIRST: Size Affects Latency

CPU

CPU

Small Memory

## Big Memory

- Signals have further to travel
- Fan out to more locations

# Memory Hierarchy



- *capacity*:  Register << SRAM << DRAM
- *latency*:   Register << SRAM << DRAM
- *bandwidth:* on-chip >> off-chip

On a data access:

*if* data $\in$ fast memory $\Rightarrow$ low latency access *(SRAM)*
*if* data $\notin$ fast memory $\Rightarrow$ high latency access *(DRAM)*

# Access Patterns



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

# Examples

# Locality of Reference

- **Temporal Locality**: If a location is referenced it is likely to be referenced again in the near future.

- **Spatial Locality**: If a location is referenced it is likely that locations near it will be referenced in the near future.
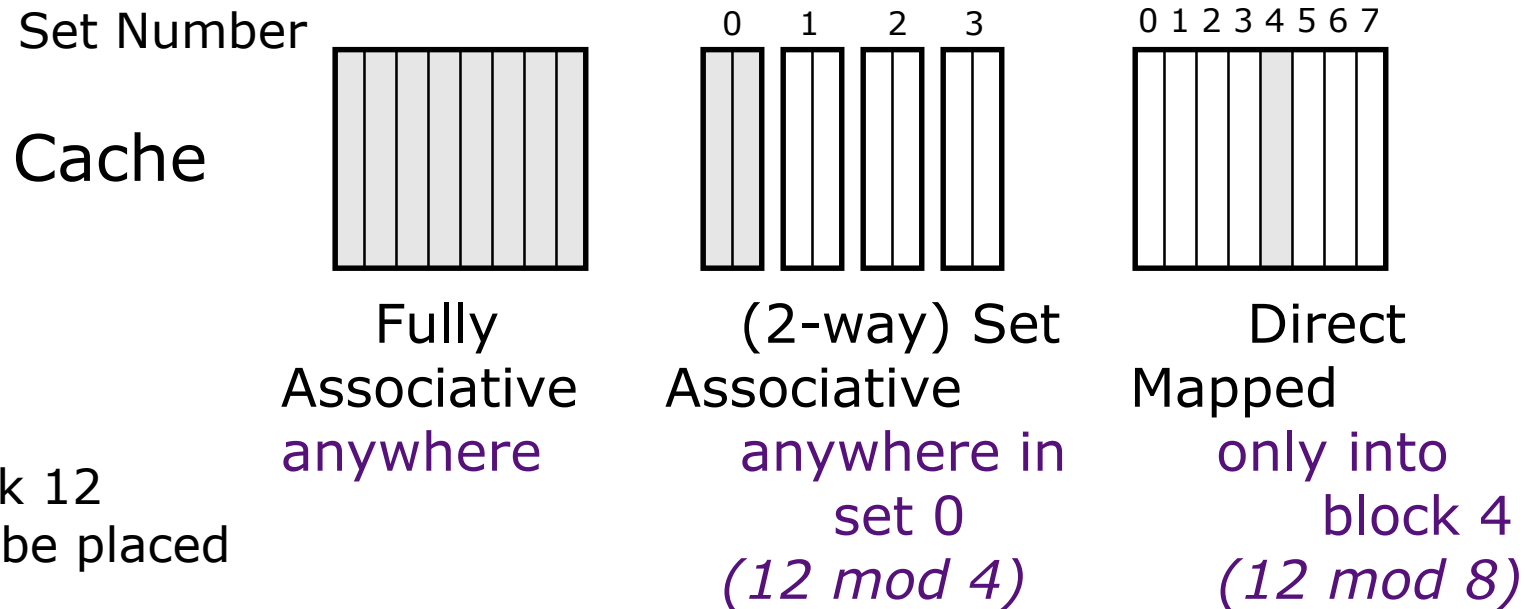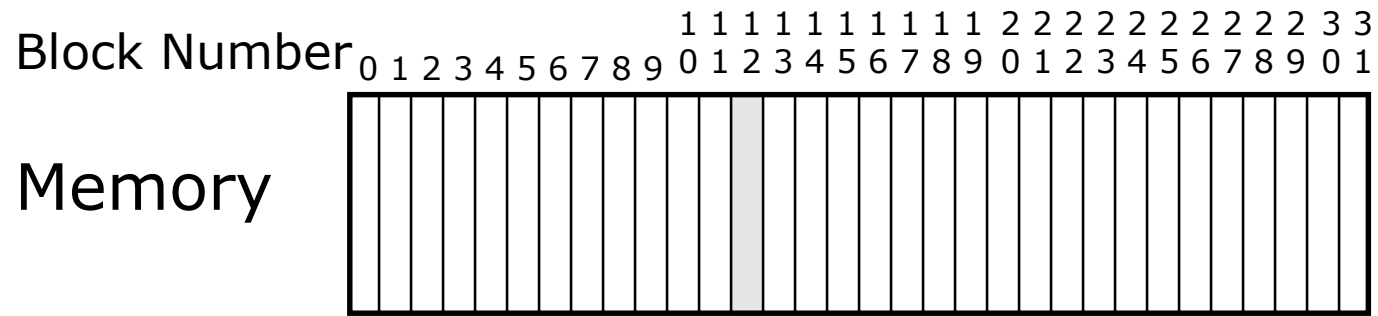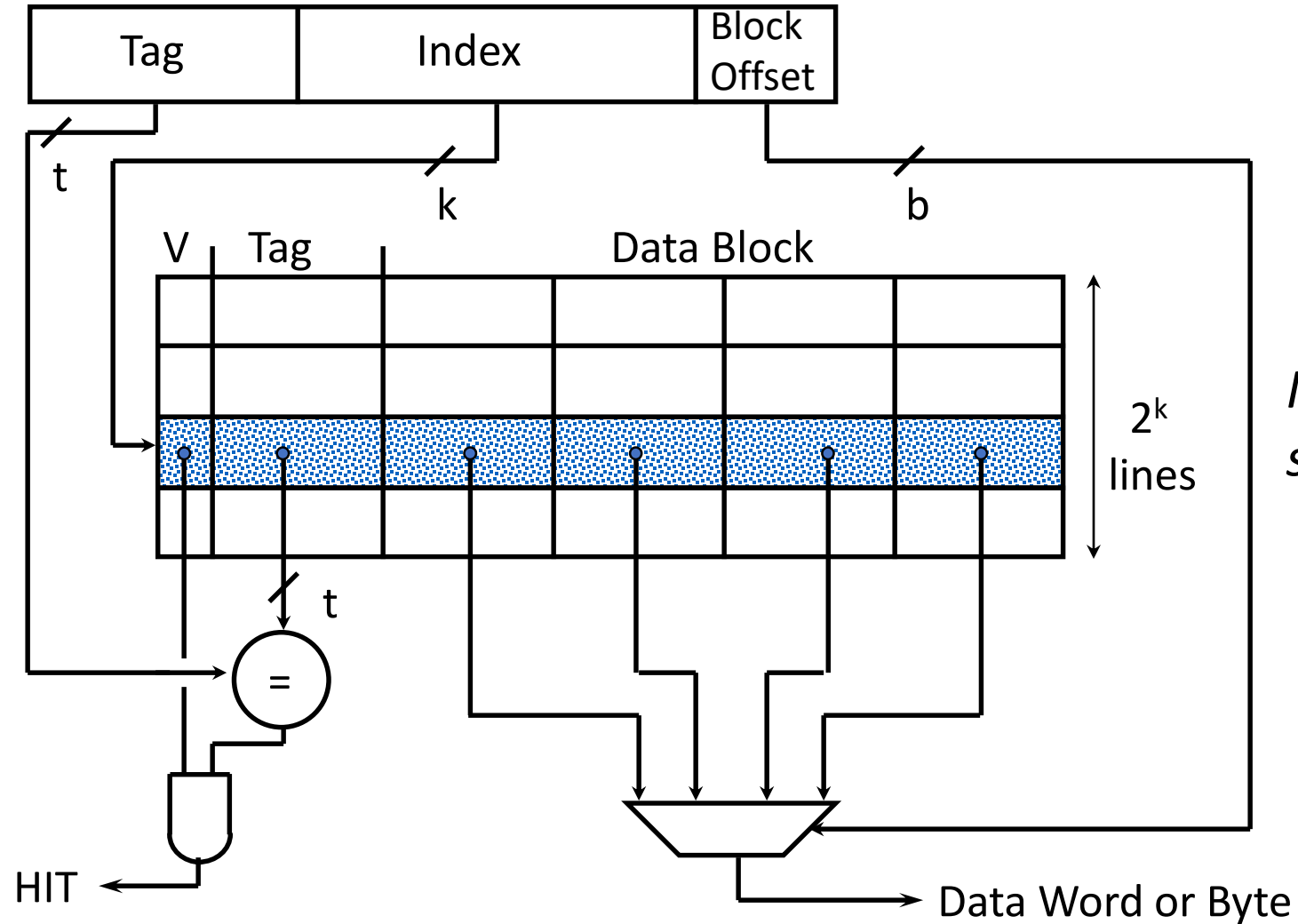
# Inside a Cache

# Placement Policy



Block Number

1 1 1 1 1 1 1 1 1 1 2 2 2 2 2 2 2 2 2 2 3 3
0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1

Memory

Set Number

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

Cache

Fully
Associative
anywhere

(2-way) Set
Associative
anywhere in
set 0
*(12 mod 4)*

Direct
Mapped
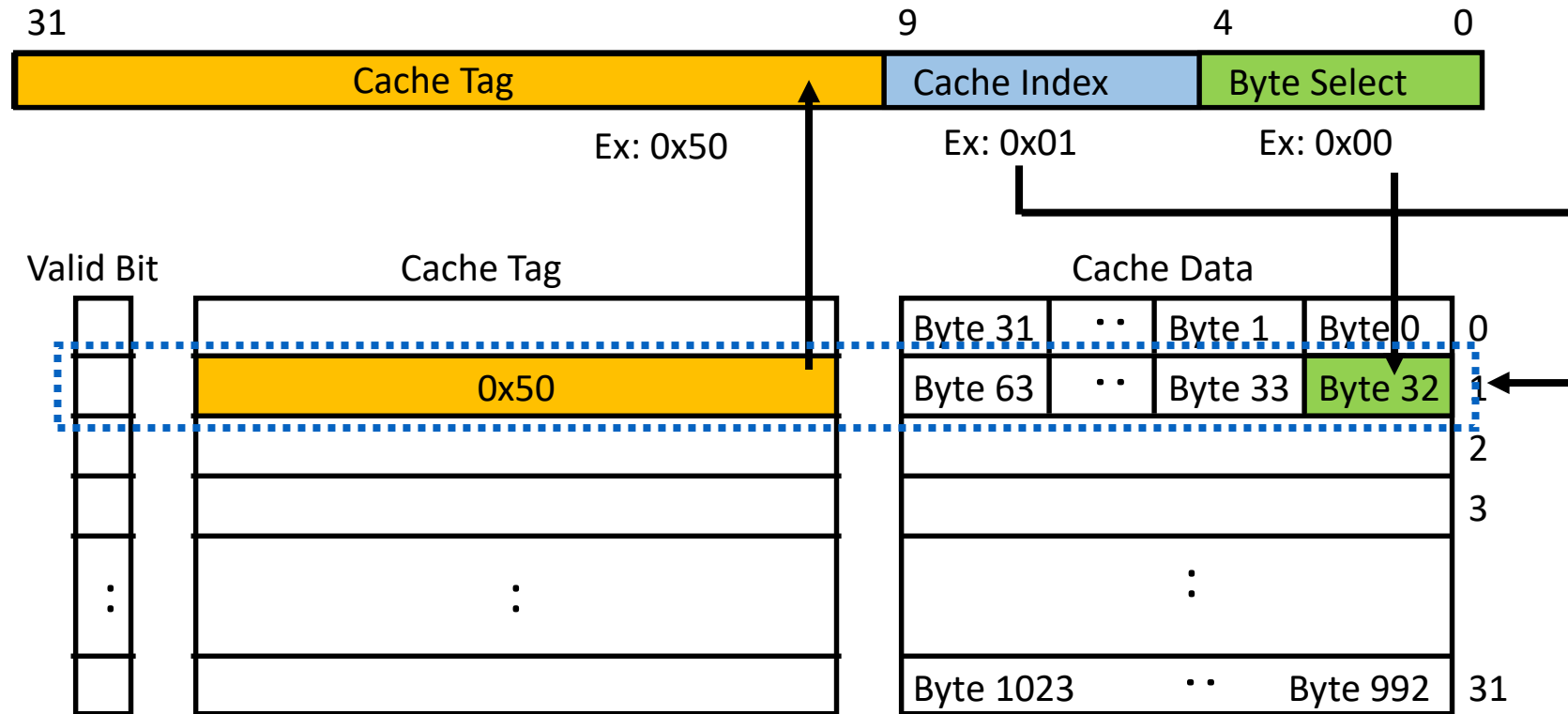only into
block 4
*(12 mod 8)*

block 12
can be placed

# Direct Mapped
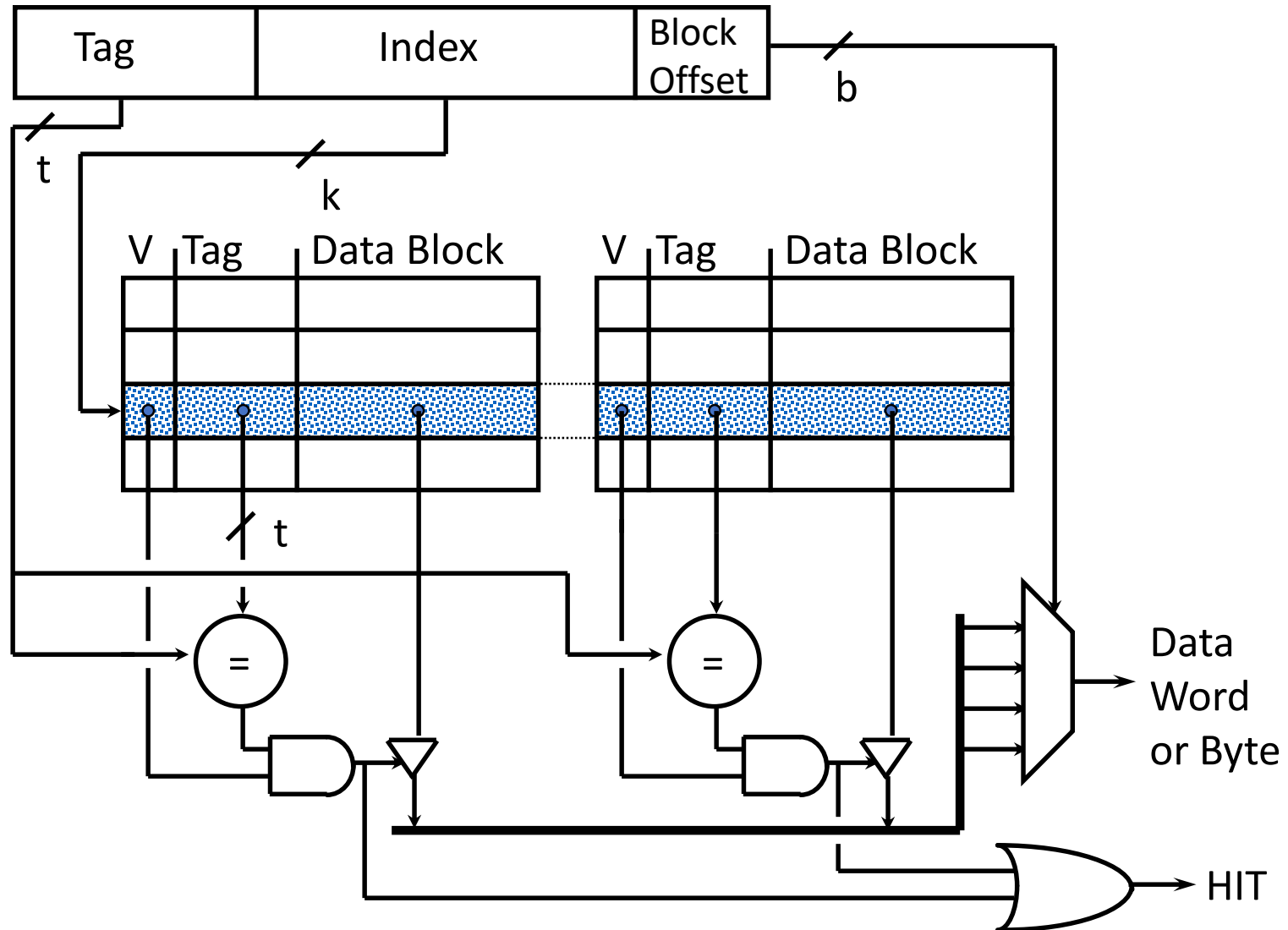


In reality, tag-store is placed separately

# An Example

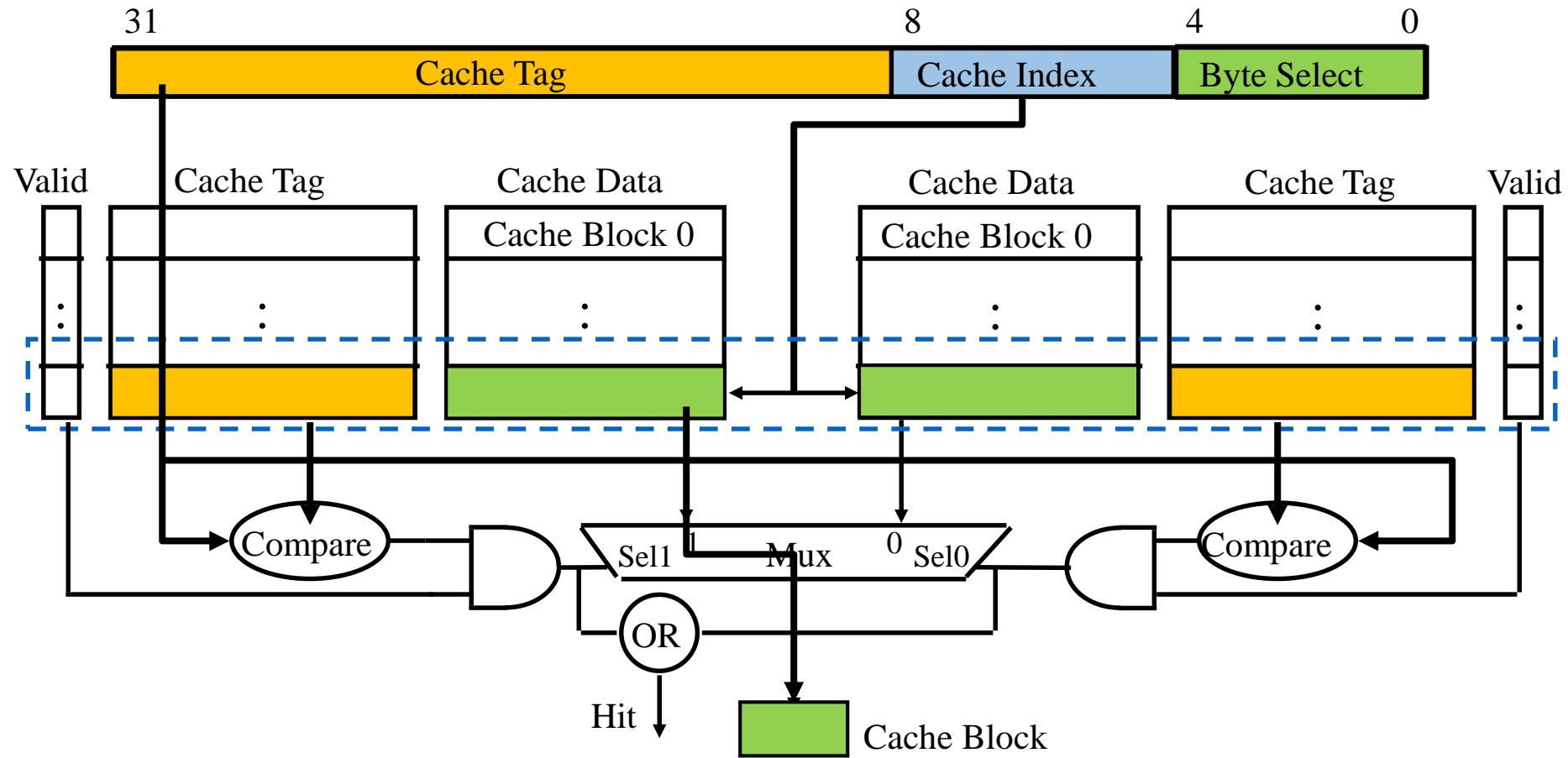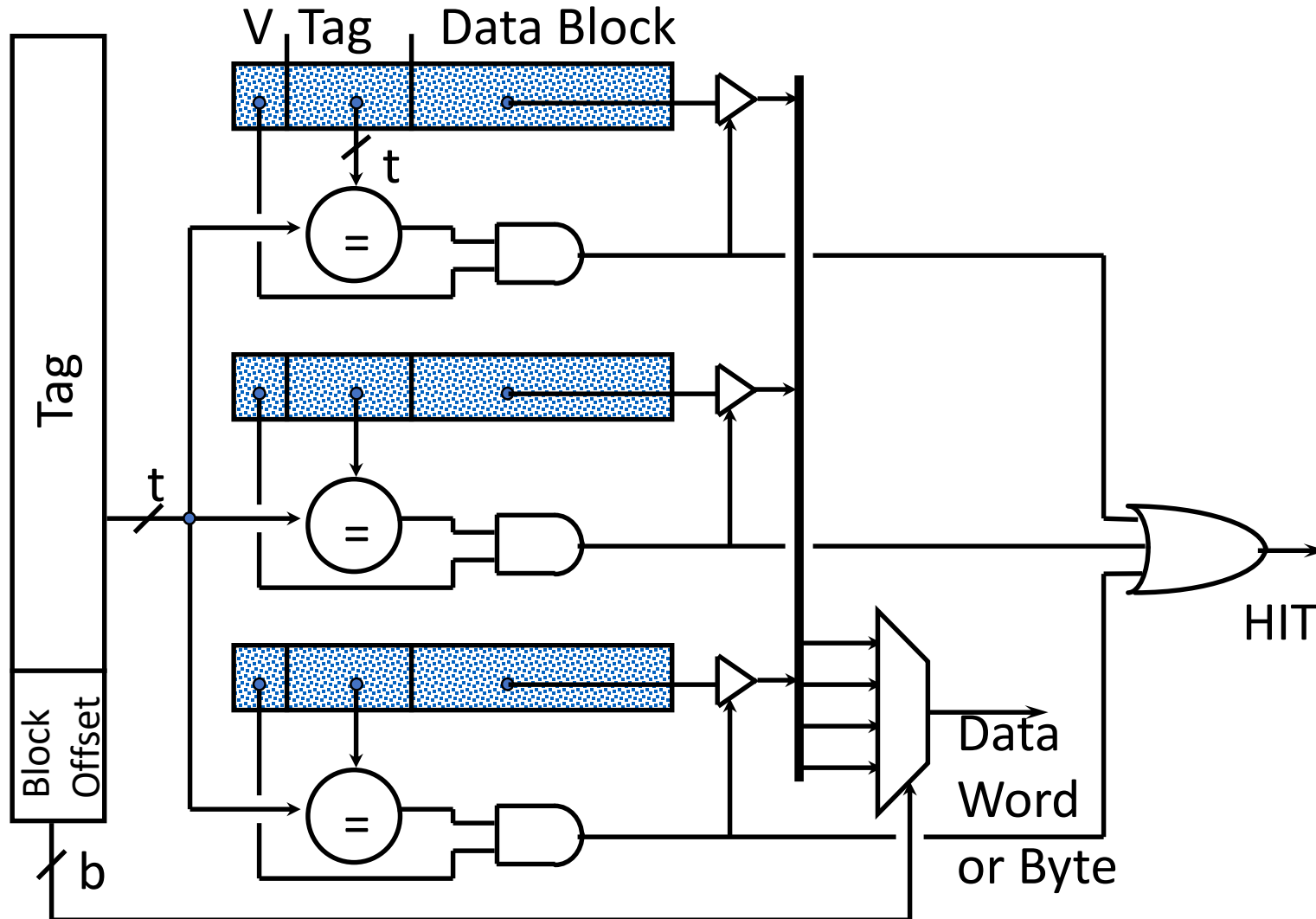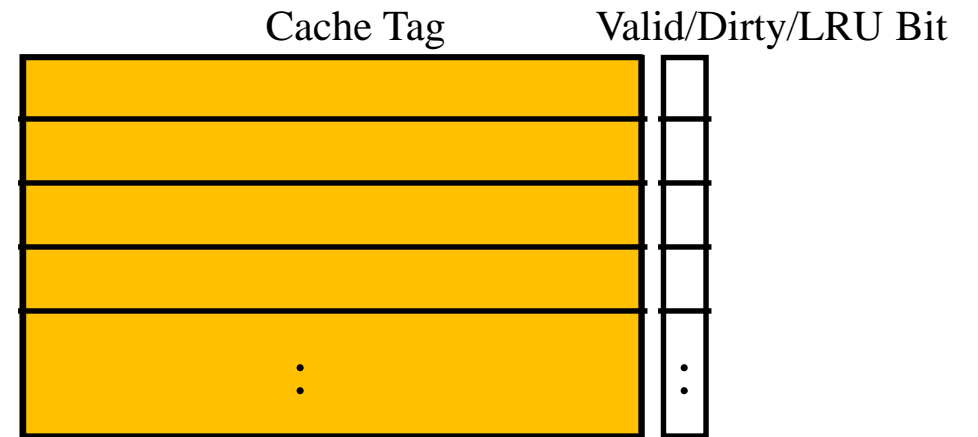# High bits or Low bits

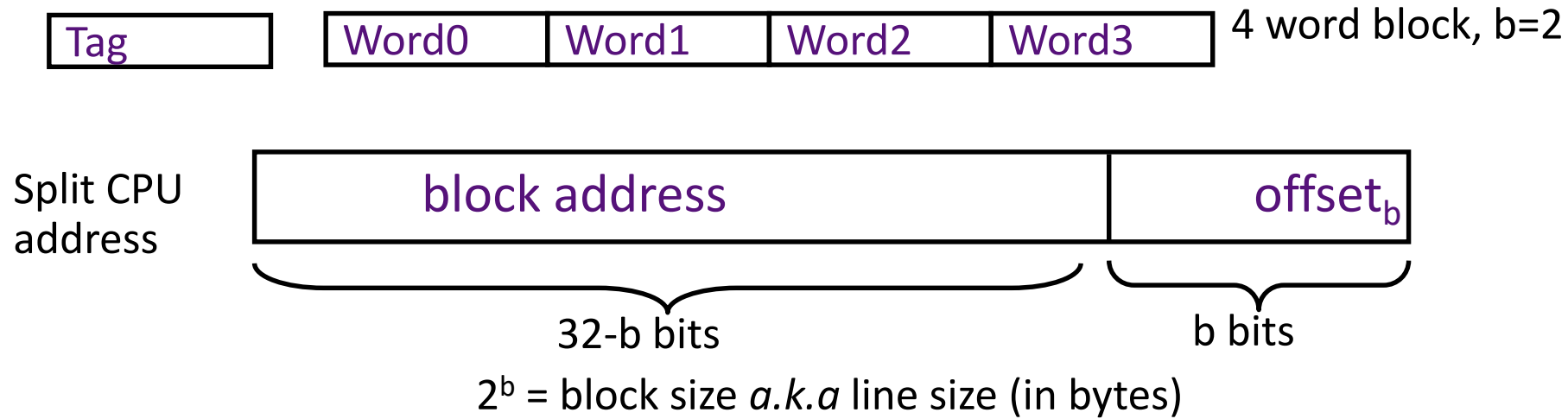# Set-Associative

# An Example

# Fully-associative

# An Example

# What's in Tag Store?

- Valid bit

- Tag

- Replacement policy bits


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

Cache Tag          Valid/Dirty/LRU Bit



Cache Data

| Byte 31 | ·· | Byte 1 | Byte 0 |
| Byte 63 | ·· | Byte 33 | Byte 32 |
| | | | |
| | | | |
| | | : | |

# Block (line) Size ?

| Tag |  | Word0 | Word1 | Word2 | Word3 | 4 word block, b=2 |
|-----|--|-------|-------|-------|-------|-------------------|

Split CPU address

| block address | $offset_b$ |
|---------------|------------|

32-b bits      b bits

$2^b$ = block size *a.k.a* line size (in bytes)

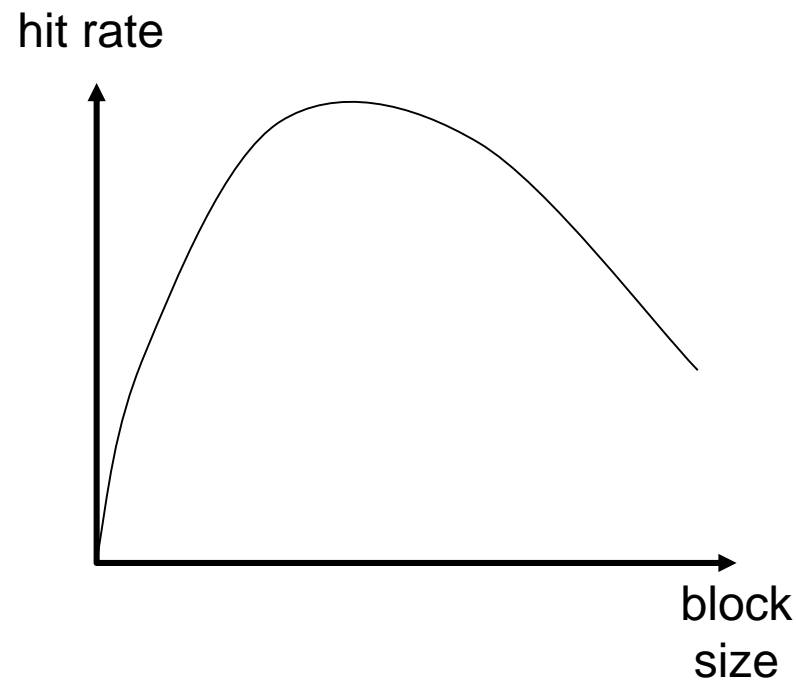Larger block size has distinct hardware advantages
- less tag overhead
- exploit fast burst transfers from DRAM
- exploit fast burst transfers over wide busses

*What are the disadvantages of increasing block size?*

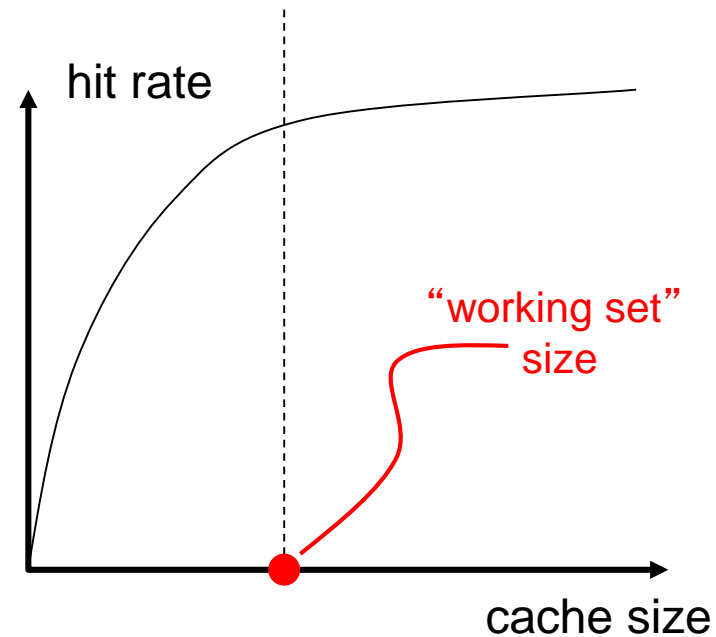*Fewer blocks => more conflicts.  Can waste bandwidth.*

# Block Size?

- Block size is the data that is associated with an address tag
  - not necessarily the unit of transfer between hierarchies

- 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

hit rate

block size

# 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

hit rate

"working set" size

cache size

# Associativity

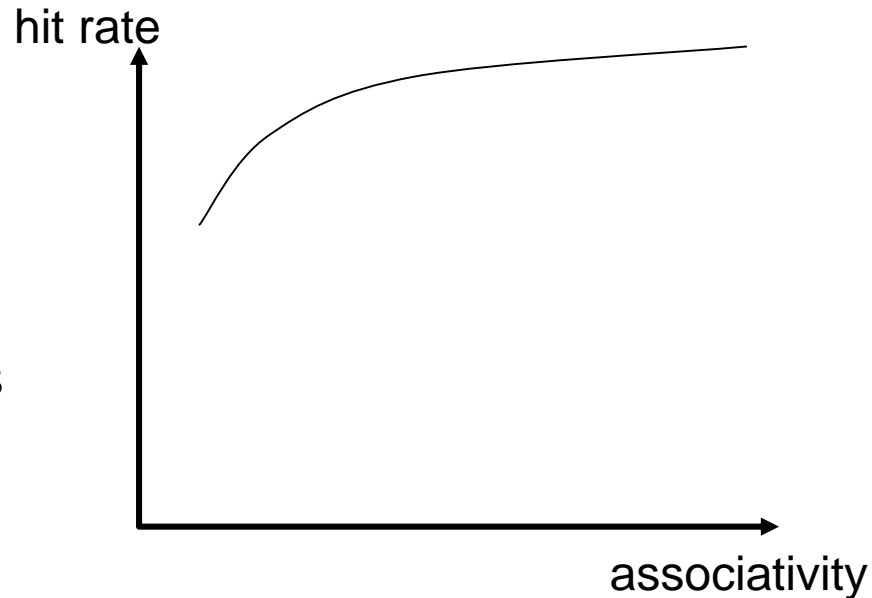- How many blocks can map to the same index (or set)?

- Larger associativity
  - lower miss rate, less variation among programs
  - diminishing returns, higher hit latency
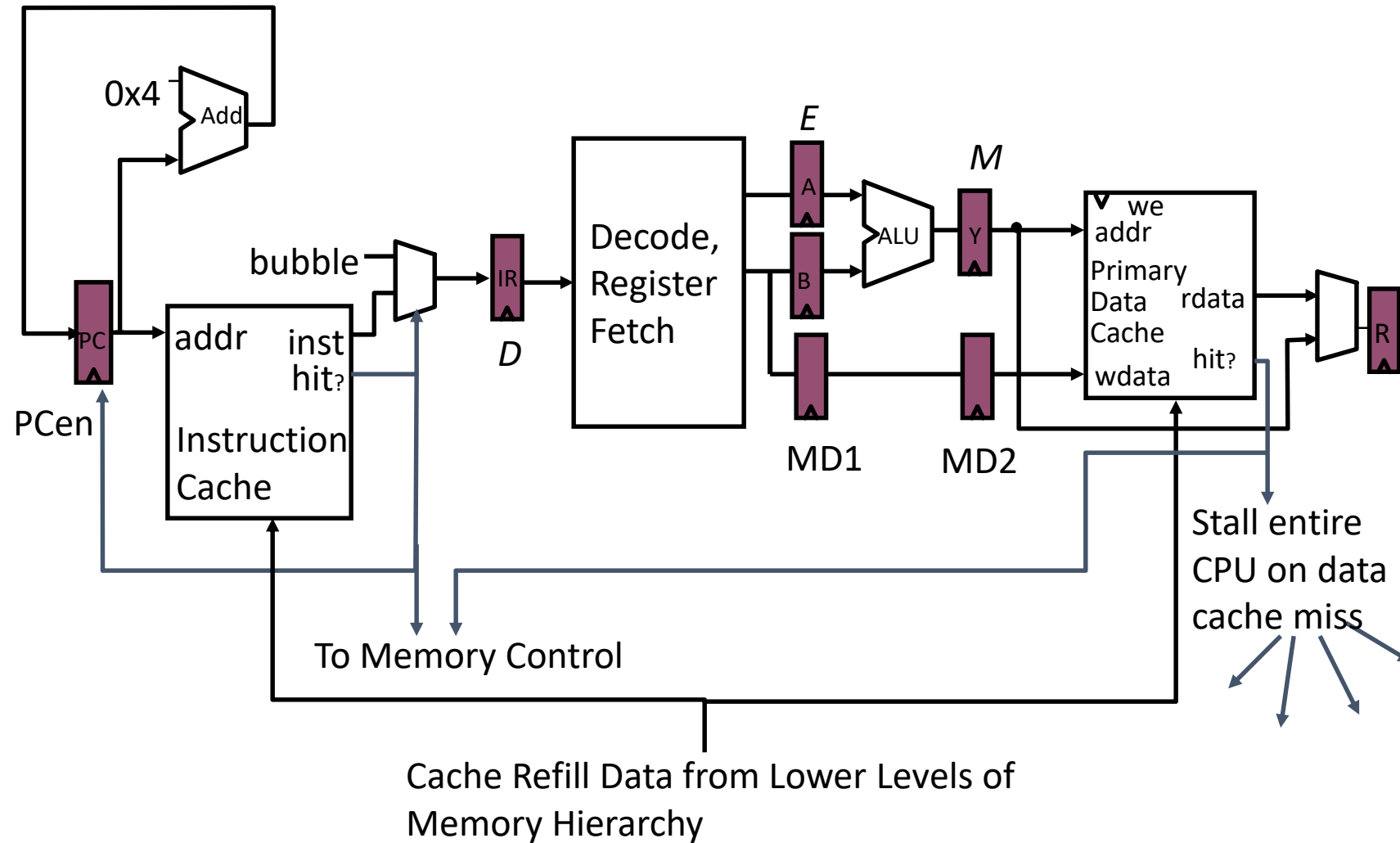
- Smaller associativity
  - lower cost
  - lower hit latency
    - Especially important for L1 caches

- Power of 2 associativity?

# Let's Pause!!

# CPU – Cache Interaction

# Design Issues: Unified vs Split

- Unified:
    - \+ Dynamic sharing of cache space: no overprovisioning

    - \-- Instructions and data can thrash each other
    - \-- I and D are accessed in different places in the pipeline.

- First level caches are almost always split
    - for the reasons above

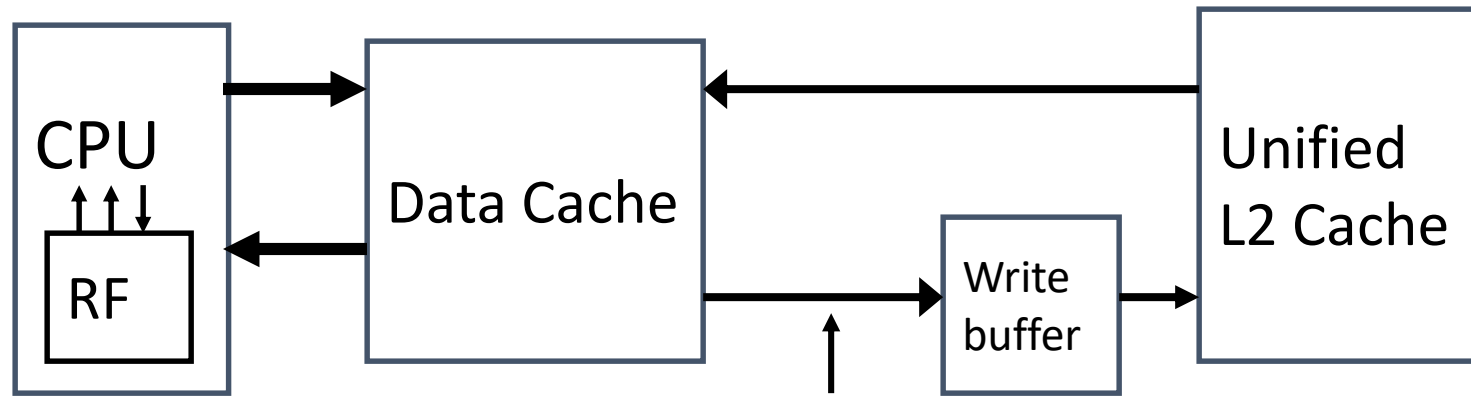- Second and higher levels are almost always unified

# Reads are not Writes

- If a write enters the cache, what happens if
    - There is a cache miss
        - Does the cache need to bring in the cache line?
    - There is a cache hit
        - Does the cache need to write back to memory?

# Write Policies

- Cache hit:
  - ***write through***: write both cache & memory
    - Generally higher traffic but simpler pipeline & cache design
  - ***write back***: write cache only, memory is written only when the entry is evicted
    - A dirty bit per line further reduces write-back traffic
    - Must handle 0, 1, or 2 accesses to memory for each load/store
- Cache miss:
  - ***no write allocate***: only write to main memory
  - ***write allocate*** (aka fetch on write):  fetch into cache

- Common combinations:
  - write through and no write allocate
  - write back with write allocate

# Write Buffers



Evicted dirty lines for writeback cache
OR
All writes in writethrough cache

Processor is not stalled on writes, and read misses can go ahead of write to main memory

**Problem:** Write buffer may hold updated value of location needed by a read miss

**Simple solution:** on a read miss, wait for the write buffer to go empty

**Faster solution:** Check write buffer addresses against read miss addresses, if no match, allow read miss to go ahead of writes, else, return value in write buffer

# Performance: AMAT

Average memory access time (AMAT) = Hit time + Miss rate x Miss penalty

Average memory access time (AMAT) = Hit time + Miss rate$_1$ x Miss penalty$_1$
+ Miss rate$_2$ x Miss penalty$_2$

# Improving Cache Performance

Average memory access time (AMAT) =

Hit time + Miss rate x Miss penalty

To improve performance:
- reduce the hit time
- reduce the miss rate
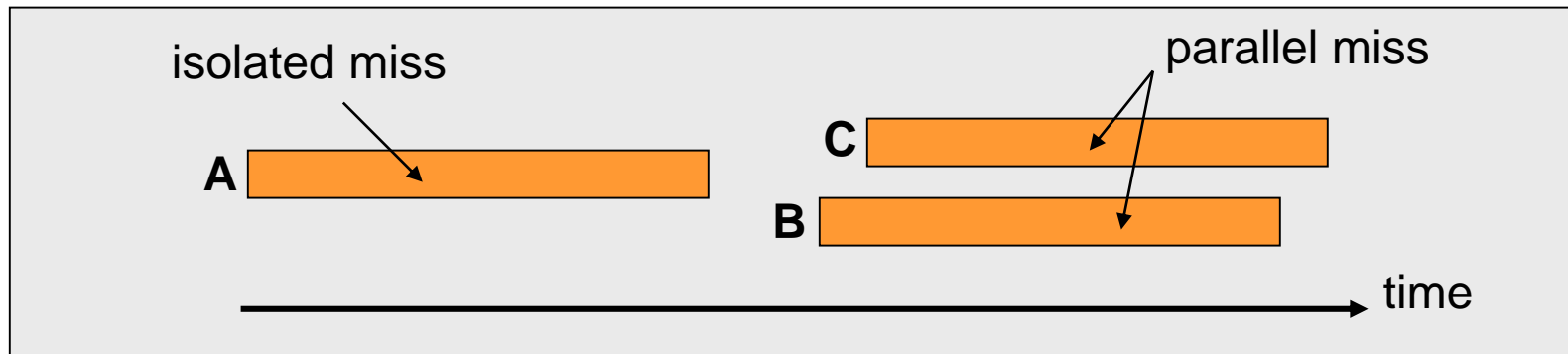- reduce the miss penalty

*Biggest cache that doesn't increase hit time past 1 cycle (approx 8-32KB in modern technology)*

*[ design issues more complex with deeper pipelines and/or out-of-order superscalar processors]*

# Cache Optimizations: Refer H&P

# Non-blocking Cache

- Enable cache access when there is a pending miss

- Enable multiple misses in parallel
  - Memory-level parallelism (MLP)
    - generating and servicing multiple memory accesses in parallel
  - Why generate multiple misses?



  - Enables latency tolerance: overlaps latency of different misses
- How to generate multiple misses?
  - Out-of-order execution, multithreading, prefetching

# Miss-Status Holding Registers

- Also called "miss buffer"
- Keeps track of
  - Outstanding cache misses
  - Pending load/store accesses that refer to the missing cache block
- Fields of a single MSHR
  - Valid bit
  - Cache block address (to match incoming accesses)
  - Control/status bits (prefetch, issued to memory, which subblocks have arrived, etc)
  - Data for each subblock
  - For each pending load/store
    - Valid, type, data size, byte in block, destination register or store buffer entry address

# MSHRs

| Valid | Block Address | Issued |
|-------|---------------|--------|
| 1 | 27 | 1 |

| Valid | Type | Block Offset | Destination | |
|-------|------|--------------|-------------|---|
| 1 | 3 | 5 | 5 | |
| Valid | Type | Block Offset | Destination | Load/store 0 |
| Valid | Type | Block Offset | Destination | Load/store 1 |
| Valid | Type | Block Offset | Destination | Load/store 2 |
| Valid | Type | Block Offset | Destination | Load/store 3 |

# MSHR in Action

- On a cache miss:
  - Search MSHR for a pending access to the same block
    - Found: Allocate a load/store entry in the same MSHR entry
    - Not found: Allocate a new MSHR
    - No free entry: stall

- When a subblock returns from the next level in memory
  - Check which loads/stores waiting for it
    - Forward data to the load/store unit
    - Deallocate load/store entry in the MSHR entry
  - Write subblock in cache or MSHR
  - If last subblock, dellaocate MSHR (after writing the block in cache)

# The 3Cs

**Compulsory:**

first reference to a line (a.k.a. cold start misses)

- *misses that would occur even with infinite cache*

**Capacity:**

cache is too small to hold all data needed by the program

- *misses that would occur even under perfect replacement policy*

**Conflict:**

misses that occur because of collisions due to line-placement strategy

- *misses that would not occur with ideal full associativity*

# Cache Knobs and Performance

- Larger cache size
  + reduces capacity and conflict misses
  - hit time will increase

- Higher associativity
  + reduces conflict misses
  - may increase hit time

- Larger line size
  + reduces compulsory and capacity (reload) misses
  - increases conflict misses and miss penalty
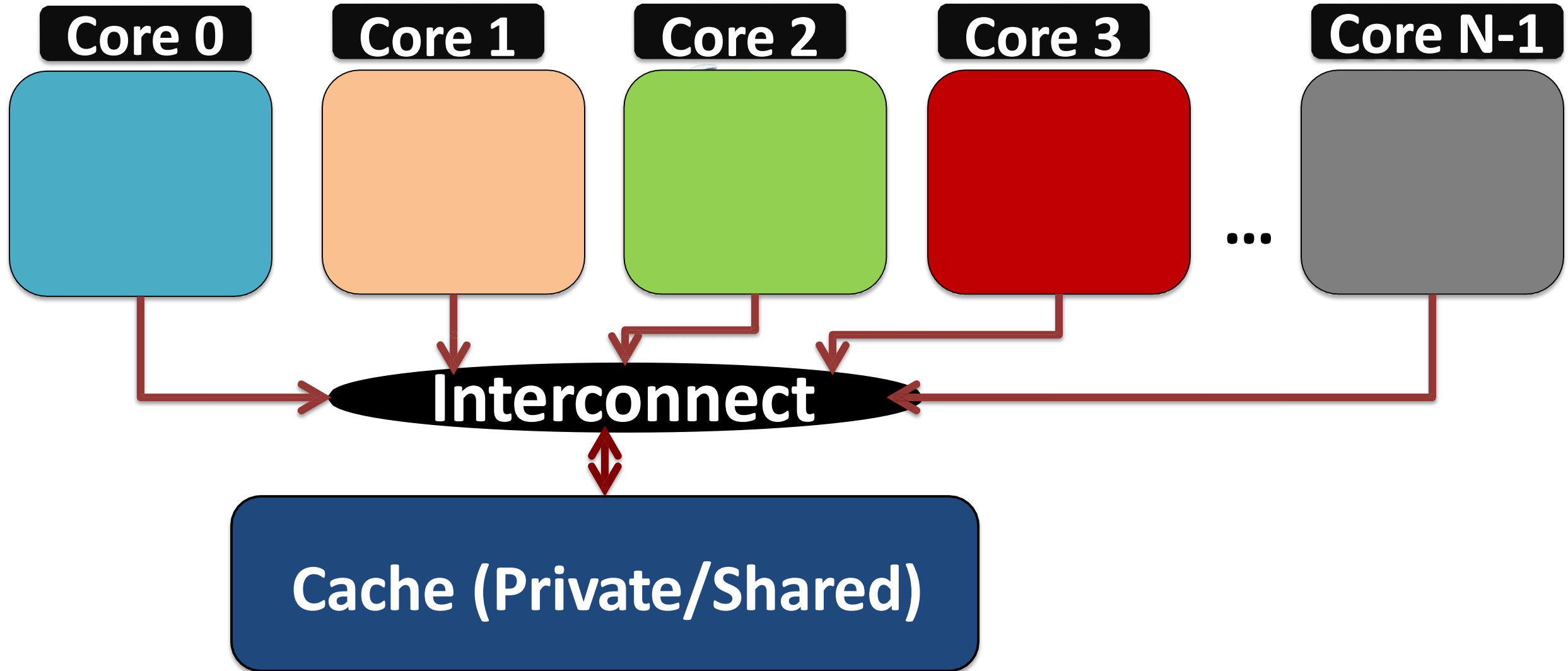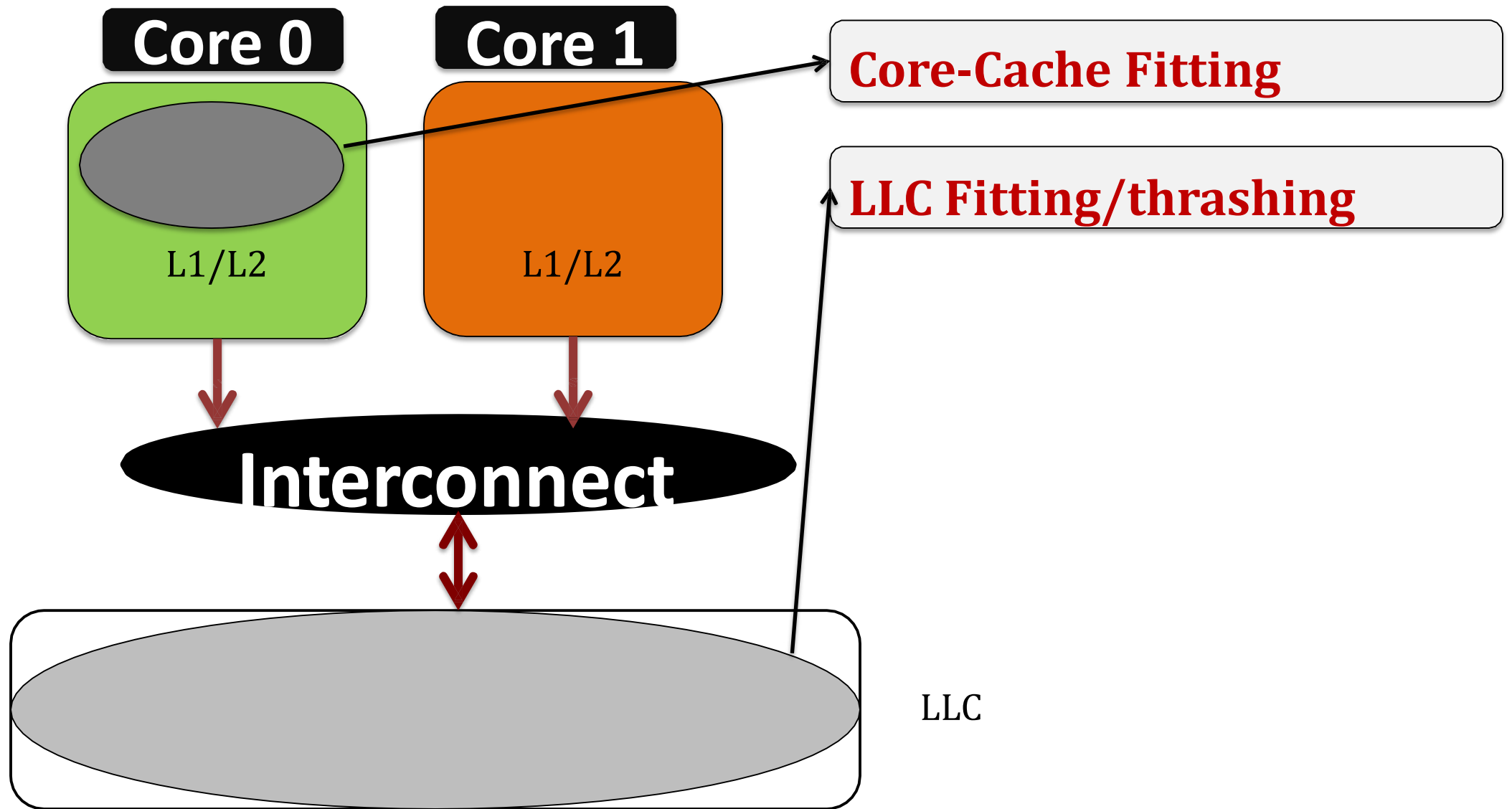
# Cache Hierarchy

# Inclusive

Core request

L1/L2

fill

evict

BackInval

fill

LLC

victim

memory

# Non-inclusive

# Exclusive

# LLC: Shared or Private?

# Application Behavior

# Shared LLC

Shared LLC provides a good tradeoff for all kinds of apps.

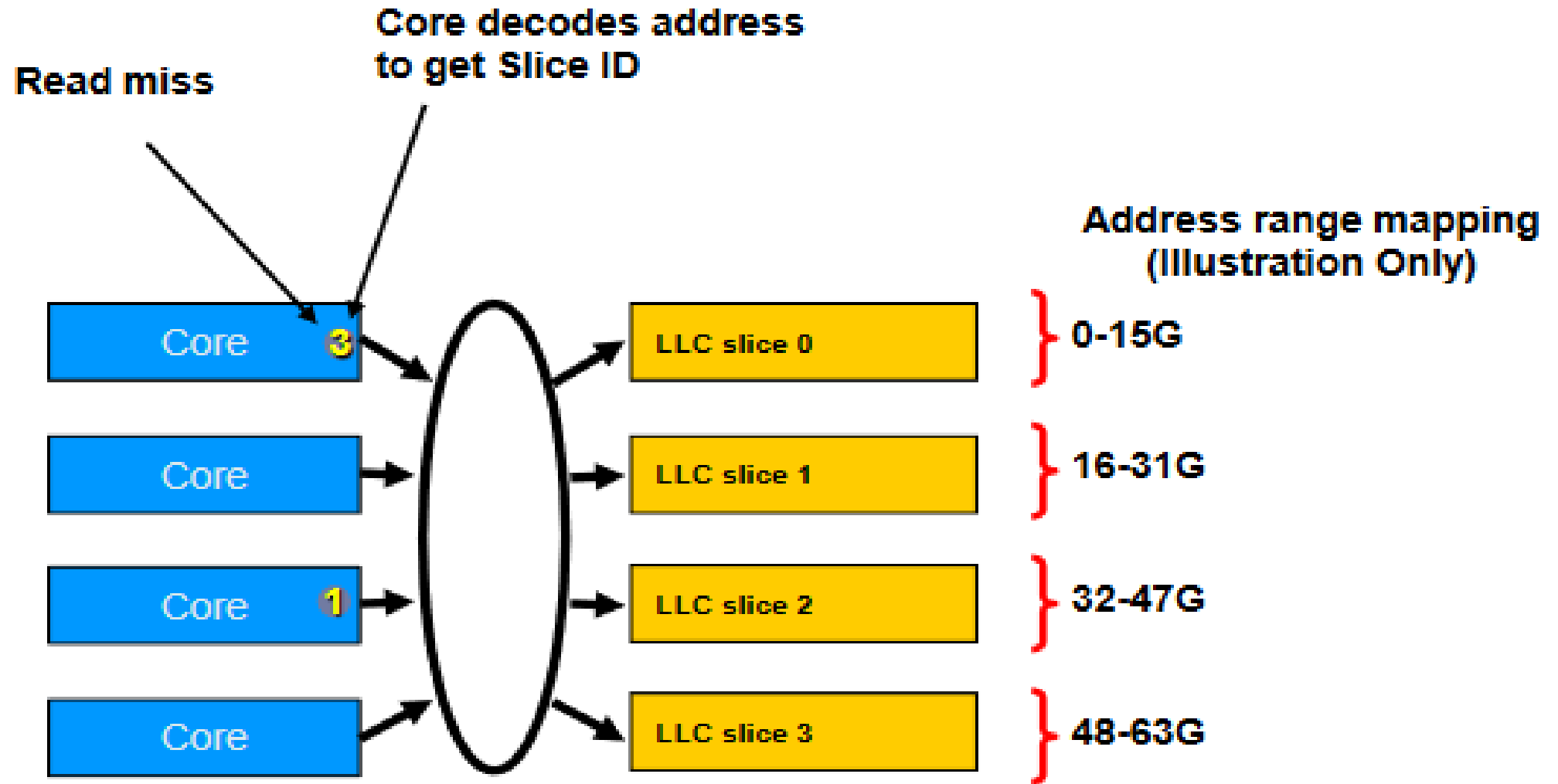Space unutilized by one app. can be utilized by other apps.

However bandwidth is an issue ☹

1000 monkeys: one banana

# Banked/Sliced NUCA

# Intel's Sandybridge

# Cache Replacement-101

- Think of each block in a set having a "priority"
  - Indicating how important it is to keep the block in the cache
- Key issue: How do you determine/adjust block priorities?
- There are three key decisions in a set:
  - Insertion, promotion, eviction (replacement)

- Insertion: What happens to priorities on a cache fill?
  - Where to insert the incoming block, whether or not to insert the block
- Promotion: What happens to priorities on a cache hit?
  - Whether and how to change block priority
- Eviction/replacement: What happens to priorities on a cache miss?
  - Which block to evict and how to adjust priorities

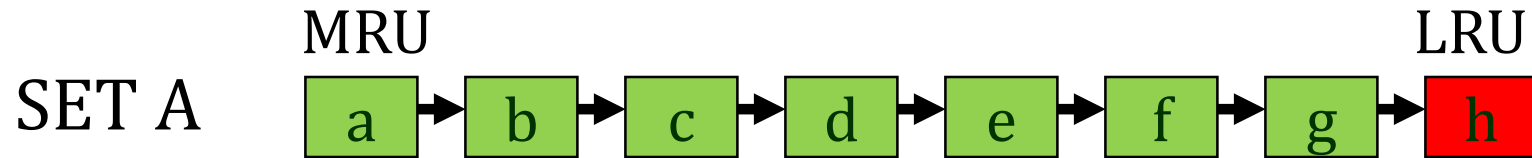# Eviction (Replacement) Policy?

- 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?
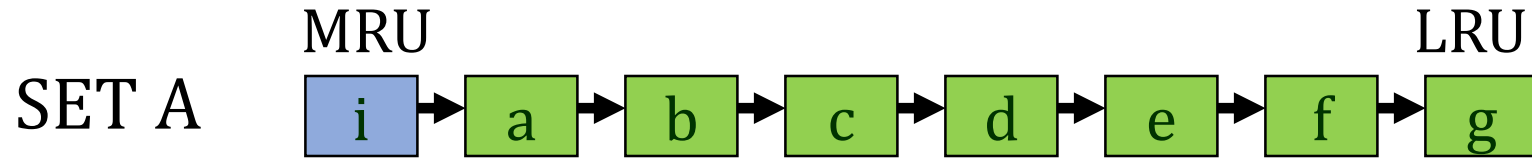    - Hybrid replacement policies
    - Optimal replacement policy?

# Belady

- 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.

**Cache Eviction Policy: On a miss (block *i*), which block to evict (replace) ?**

SET A    MRU                                                    LRU
         a → b → c → d → e → f → g → h

**Cache Insertion Policy: New block *i* inserted into MRU.**

SET A    MRU                                                    LRU
         i → a → b → c → d → e → f → g

**Cache Promotion Policy: On a future hit (block *i*), promote to MRU**

**LRU causes thrashing when working set > cache size**

# Access Patterns

**Recency friendly** $(a_1, a_2,....a_k, a_{k-1}, ....a_2, a_1)^N$

**Thrashing** $(a_1, a_2,....a_k)^N$  **[k > cache size]**

**Streaming** $(a_1, a_2,....a_\infty)^N$

**Combination of above three**

# Types of Workloads – 4MB cache



(a) Cache "Friendly" Workloads

(b) Cache "Fitting" Workloads

(c) Cache "Thrashing" Workloads

(d) Streaming Workloads

# Limitations of LRU

LRU exploits **temporal locality**

**Streaming data ($a_1$, $a_2$, $a_3$,....$a_\infty$):**
No temporal locality,
No temporal reuse

**Thrashing data ($a_1$, $a_2$, $a_3$,....$a_n$) [n>c]**
Temporal locality exists. However, LRU fails to capture.

# Still Miles to Go

$LLC_{size}$

$W_{size}$

Working set larger than the cache causes thrashing

· · · miss miss miss miss miss · · ·

References to non-temporal data (*scans*) discards frequently referenced working set

· · · hit hit hit ***scan*** miss hit ***scan*** miss hit ***scan*** miss · · ·

# Still Miles to Go

Working set larger than the cache ➔
Preserve some of working set in the cache

W$_{size}$

LLC$_{size}$

Recurring *scans (bursts of non-temporal data)* ➔ Preserve frequently referenced working set in the cache
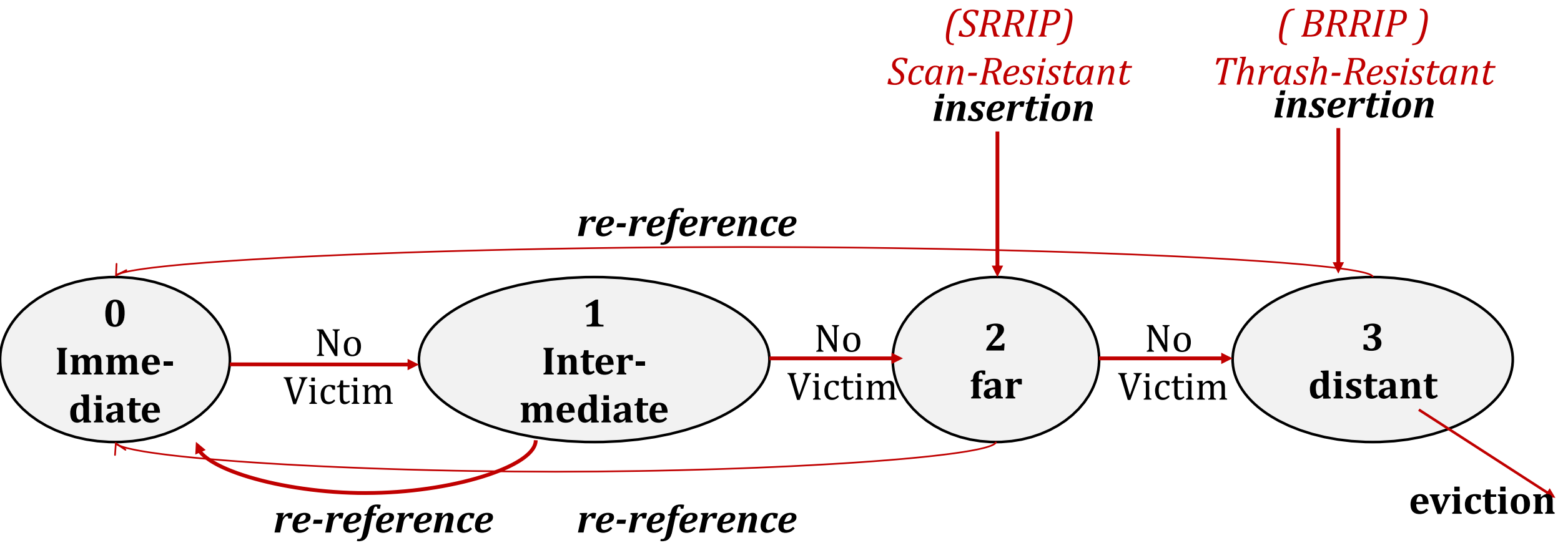
# RRIP – ISCA '10

RRP head
~~MRU position~~

RRP tail
~~LRU position~~

| h | g | f | e | d | c | b | a |

LRU chain position stored with each cache block →

0  1  2  3  4  5  6  7

RRP: Re-reference prediction

# RRIP

RRP Head

RRP Tail

| h | g | f | e | d | c | b | a |

Re-reference Prediction Value
RRPV (n=2):

0     1     2     3

Qualitative Prediction:

'near-immediate'    'intermediate'    'far'    'distant'

Intuition: New cache block will not be re-referenced soon. Replaces block with distant RRPV.

Insert with RRPV=2, Evict with RRPV=3 promote blocks with RRPV=0.

Static RRIP (Single core) and Thread-Aware Dynamic RRIP (SRRIP+BRRIP, multi-core, based on SDMs).

# Hardware Prefetching

*What?*
Latency-hiding technique - Fetches data before the core demands.

*Why?*
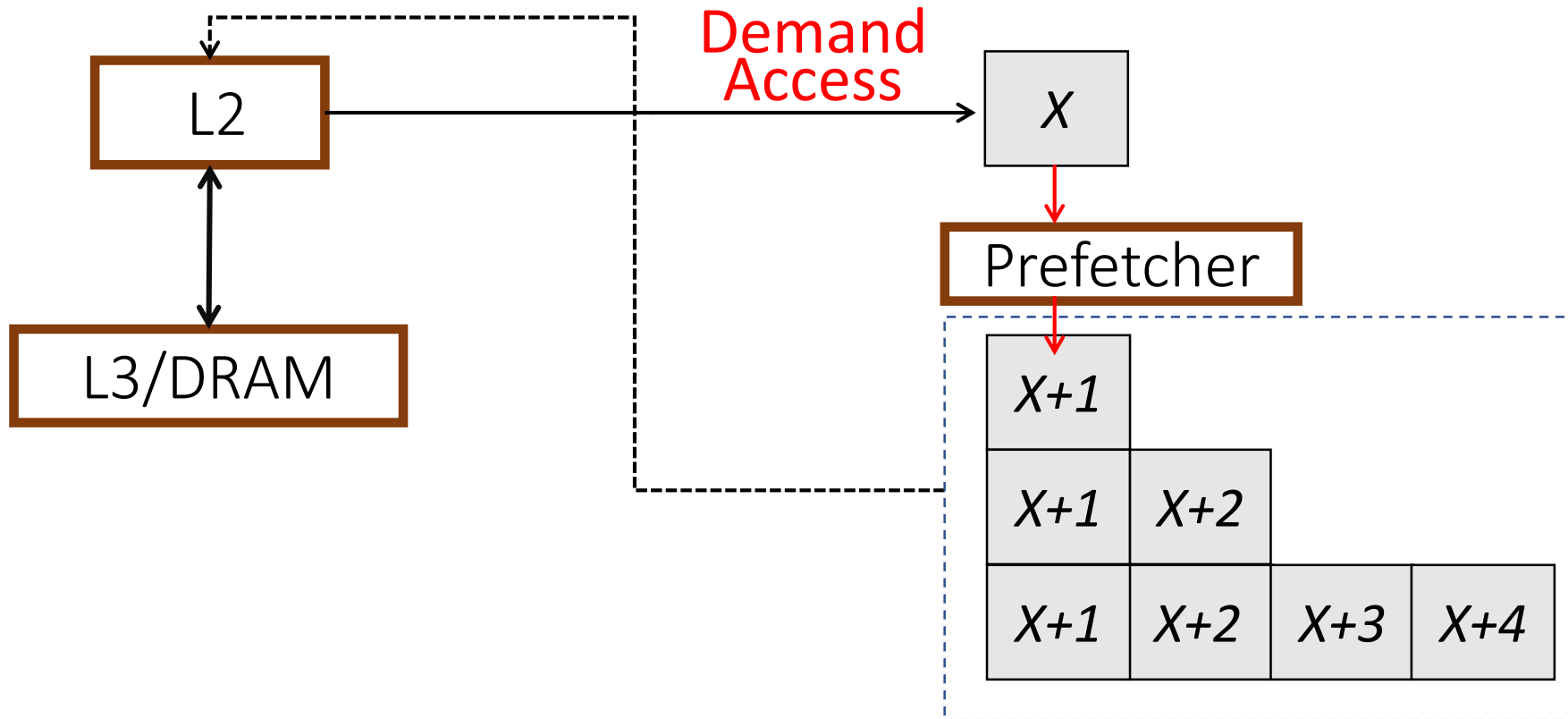Off-chip DRAM latency has grown up to 400 to 800 cycles.

*How?*
By observing/predicting the demand access (LOAD/STORE) patterns.
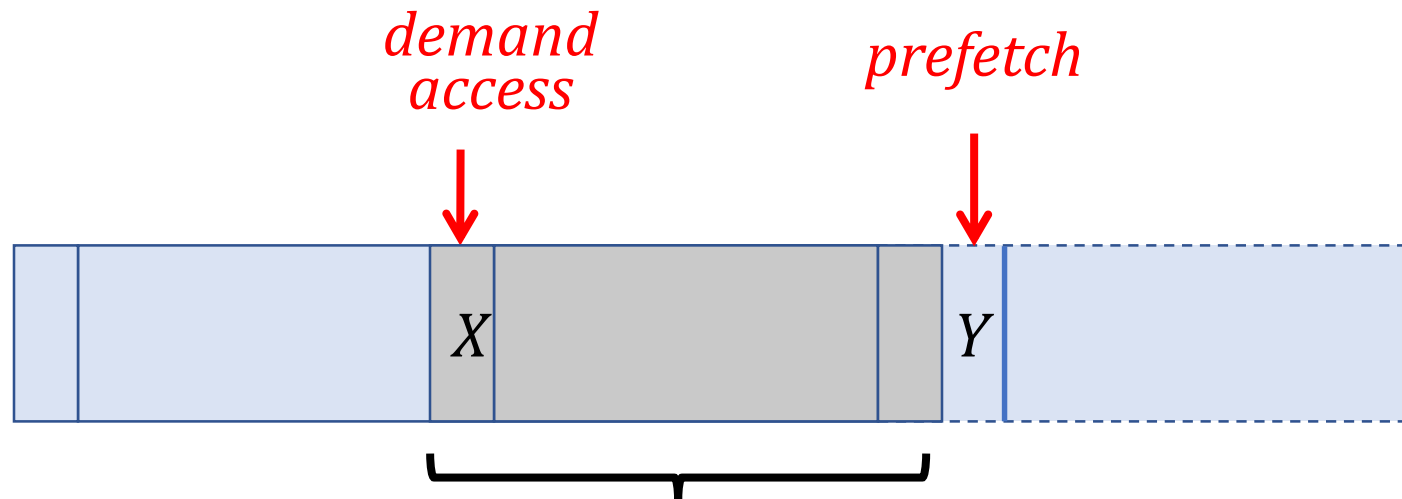
# Prefetch Engine

# Prefetch Degree

Prefetch Degree: Number of prefetch requests to issue at a given time.

# Prefetch Distance

Prefetch Distance: How far ahead of the demand access stream are the prefetch requests issued?



*demand access*

*prefetch*

$X$

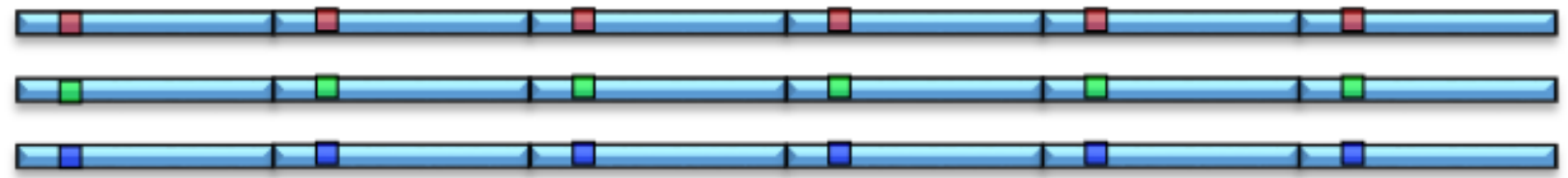$Y$

*Prefetch-distance*

$Y = X + 4$

$Y = X + 8$

$Y = X + 16$

# Next-line Prefetcher

Next Line: Miss to cache block X , prefetch X+1. Degree=1, Distance=1

Works well for L1 Icache and L1 Dcache.
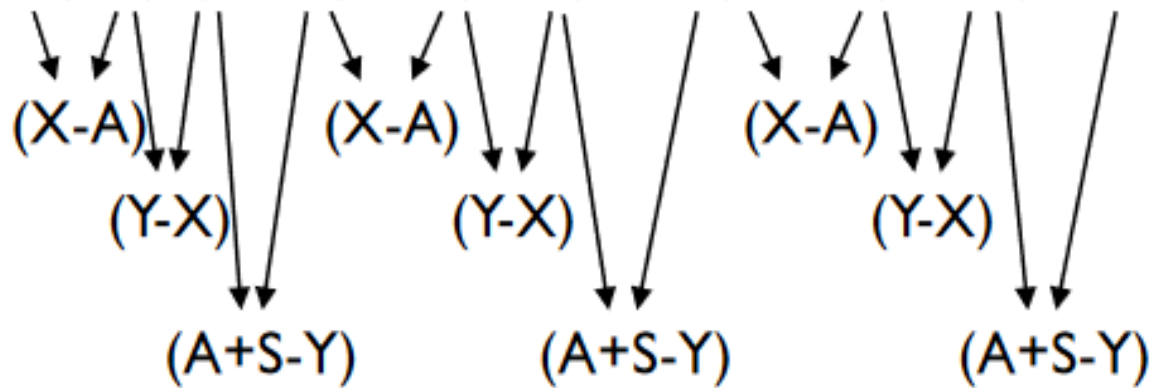
# What About This?



$$Y = A + X?$$
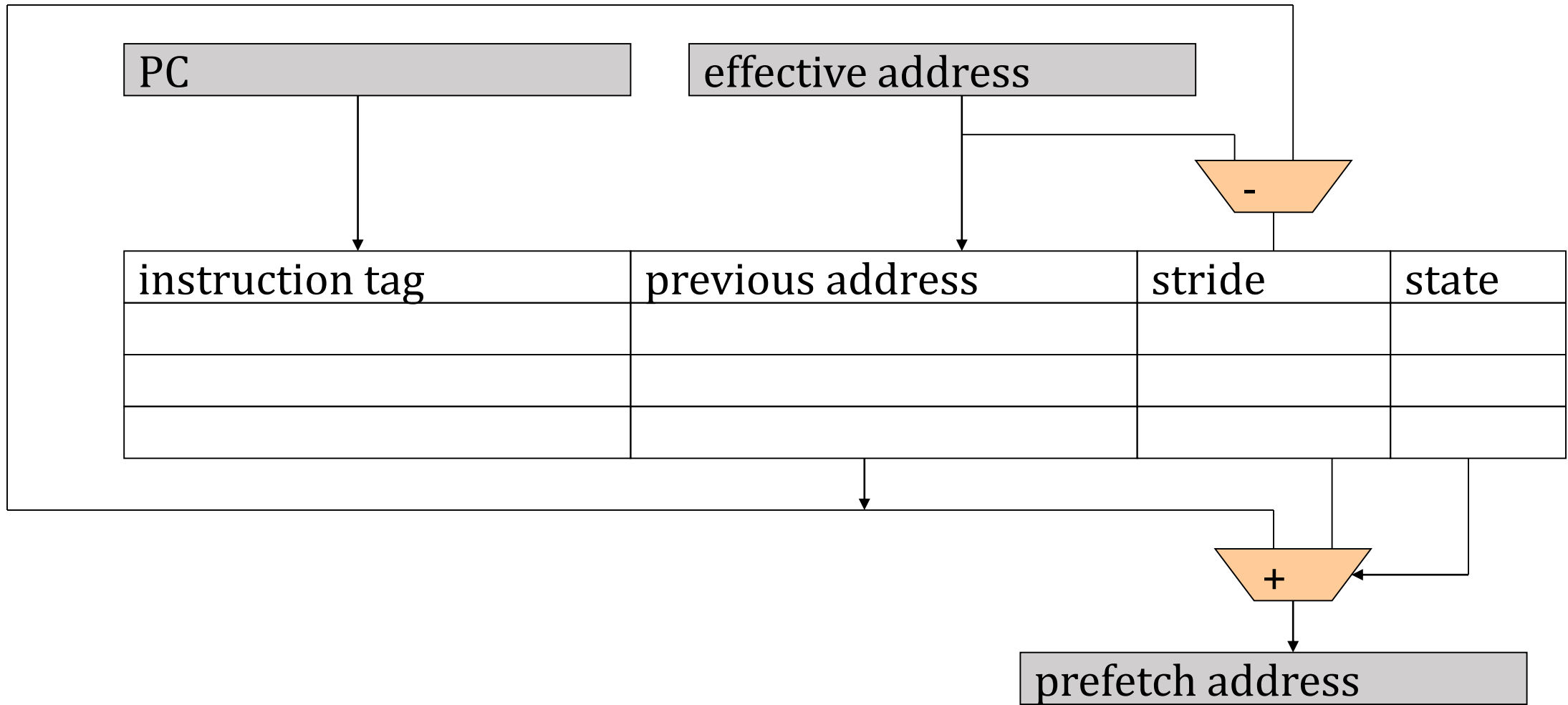
Load R1 = [R2]
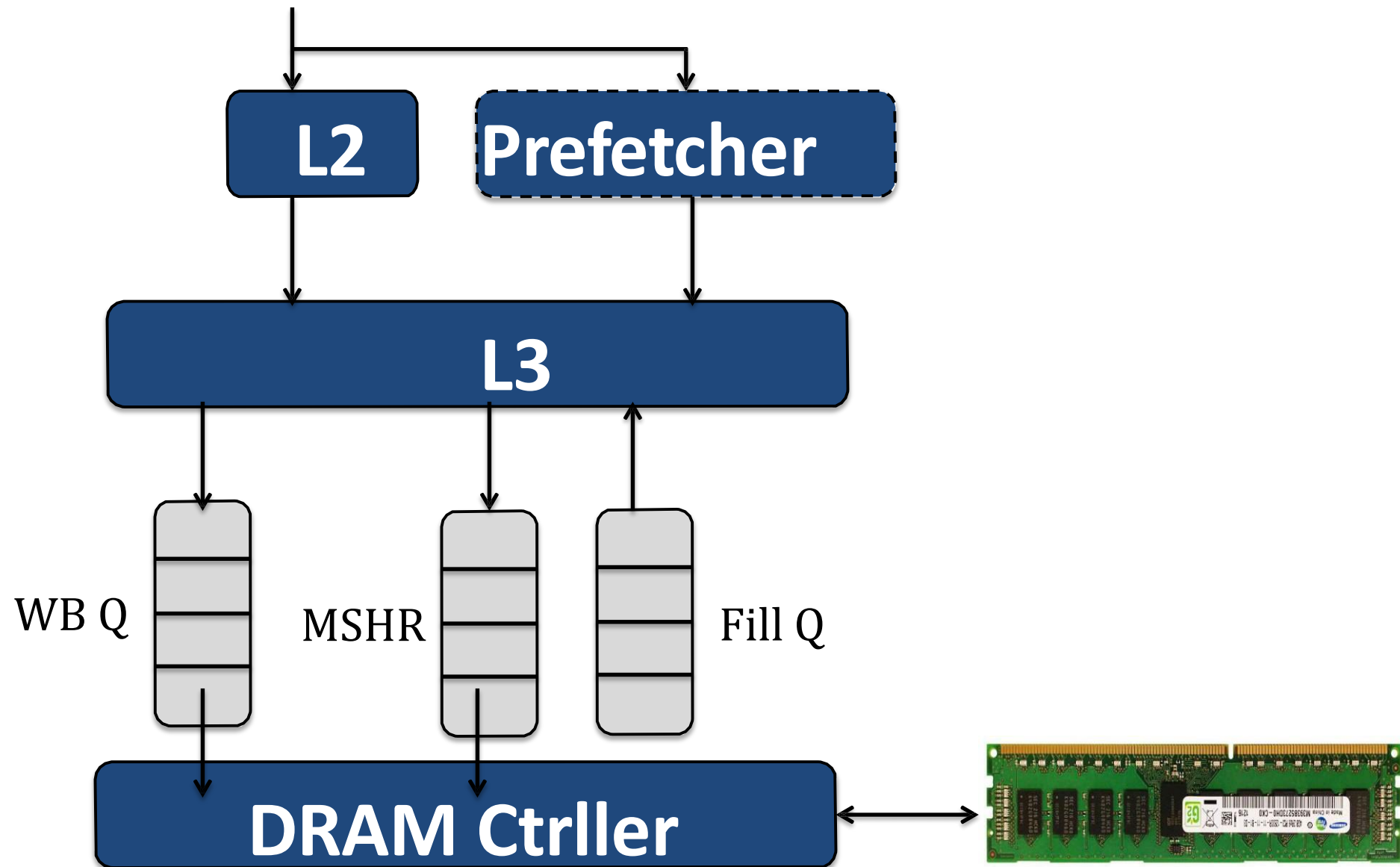Load R3 = [R4]
Add R5, R1, R3
Store [R6] = R5

A, X, Y, A+S, X+S, Y+S, A+2S, X+2S, Y+2S, ...

(X-A)   (X-A)   (X-A)
(Y-X)   (Y-X)   (Y-X)
(A+S-Y)   (A+S-Y)   (A+S-Y)

# Stride Prefetching

# ACKS

- Some of the slides are adapted and modified from Joel Emer, Arvind, Yale Patt, John Kubiatowicz, Onur Mutlu, Krste Asanovic, Mattan Erez, Rajeev Balasubramonian, and Mainak Chaudhuri