# CS773-2022-Autumn: Computer Architecture for Performance and Security

## Lecture 4: Catch the cache-II
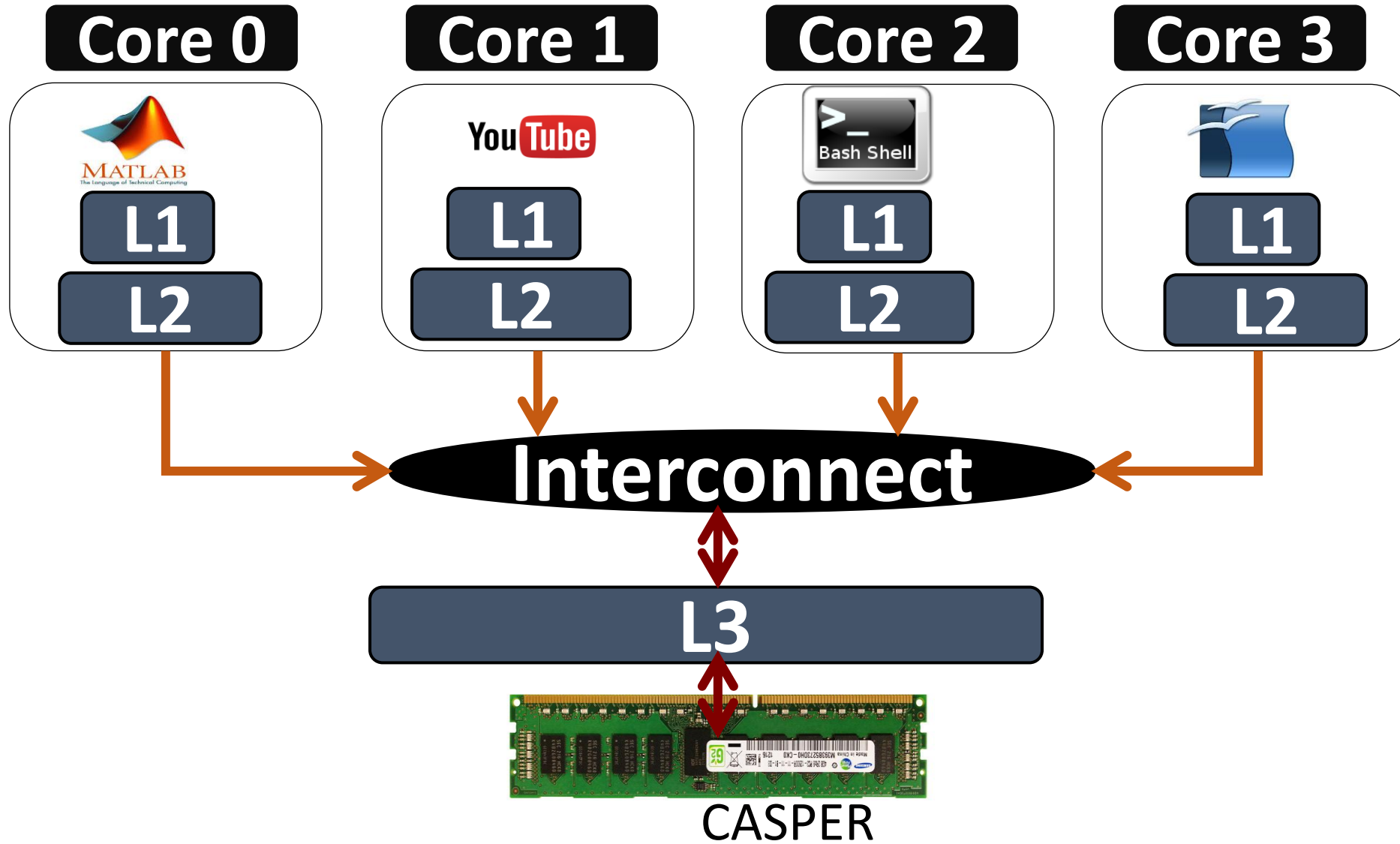
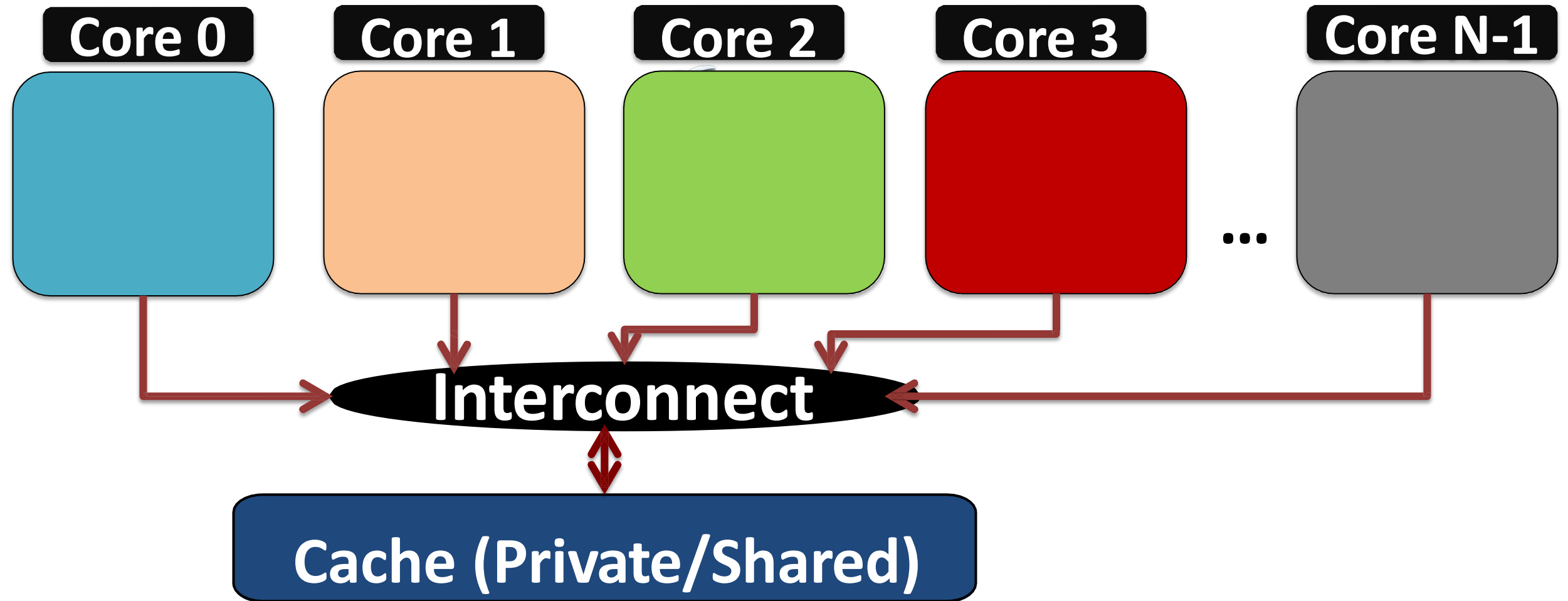*https://www.cse.iitb.ac.in/~biswa/*

# Phones on silence, please

# Thank You

Multicore

Core 0 | Core 1 | Core 2 | Core 3

L1 — L2 (for each core)

Interconnect

L3

CASPER

3

# Caches: Private/Shared

# Application behavior



**Core 0**

**Core 1**

L1/L2

L1/L2

Core-Cache Fitting

LLC Fitting/thrashing

**Interconnect**

LLC

LLC

CASPER

5

# Sliced/Banked LLC



All caches are
pipelined
In commercial
machines

CASPER

6

# Inclusive Cache Hierarchy



Core request

L1/L2

evict

fill

BackInval

LLC

fill

victim

memory

CASPER

7

# Non-inclusive (many commercial machines)



Core request

L1/L2

fill

LLC

fill

victim

memory

CASPER

# Exclusive hierarchy

Core request

fill

L1/L2

victim

LLC

fill

victim

memory

CASPER

# Cache misses

Cold Miss: cache starts empty and this is the first reference

Conflict Miss: Many mapped to the same index bits

Capacity Miss: Cache size is not sufficient

Coherence Miss: in Multi-core systems, only [not I/O coherence]

# Cache Performance

How good is the cache for a given application?

Hit rate

Miss rate

Misses per kilo instructions (MPKI)

But Why?

# On a Miss, Replace a block, which block?

Think of each block in a set having a "priority"
    Indicating how important it is to keep the block in the cache

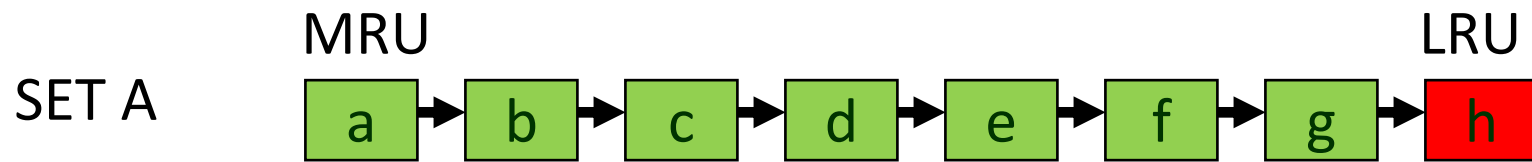Key issue: How do you determine/adjust block priorities?

Ideally: Belady's OPT policy, replace the block that will be used furthest in the future. No one knows the future though ☺
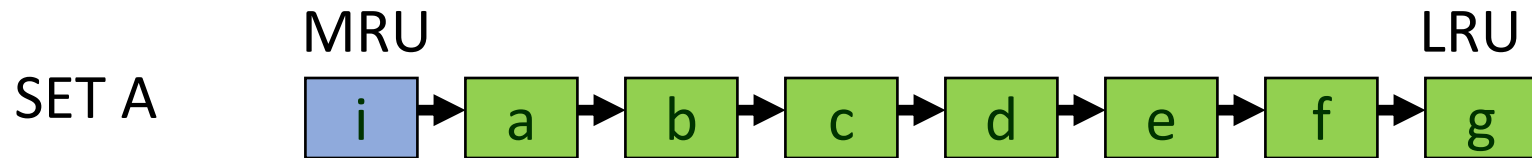
There are three key decisions in a set:
    Insertion, promotion, eviction (replacement)

# A simple LRU (Least-Recently-Used) Policy

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

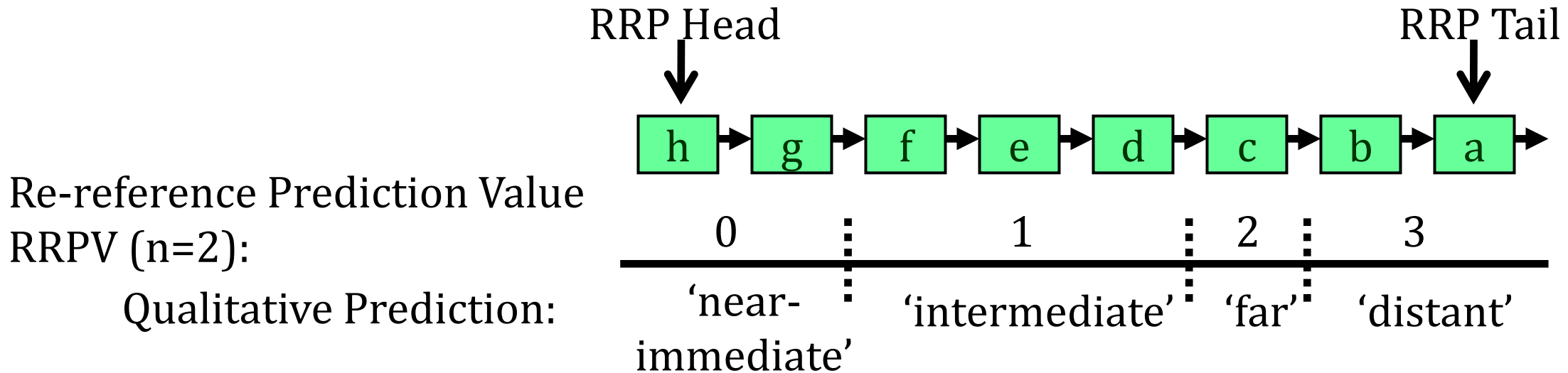i → a → b → c → d → e → f → g

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

*We need priority bits per block. For example, a 16-way cache will need four bit/block LRU causes thrashing when working set > cache size*

# Types of Applications



(a) Cache "Friendly" Workloads    (b) Cache "Fitting" Workloads    (c) Cache "Thrashing" Workloads    (d) Streaming Workloads

# LRU is not effective for Shared Caches

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. Only two bits per block.

Insert with RRPV=2, Evict with RRPV=3, increment RRPVs till we get a block with RRPV=3, promote blocks with RRPV=0.

CASPER                                                                          15

# Average Access Time

On average, how much time it takes for a LOAD to complete
*Average memory access time (AMAT) =*
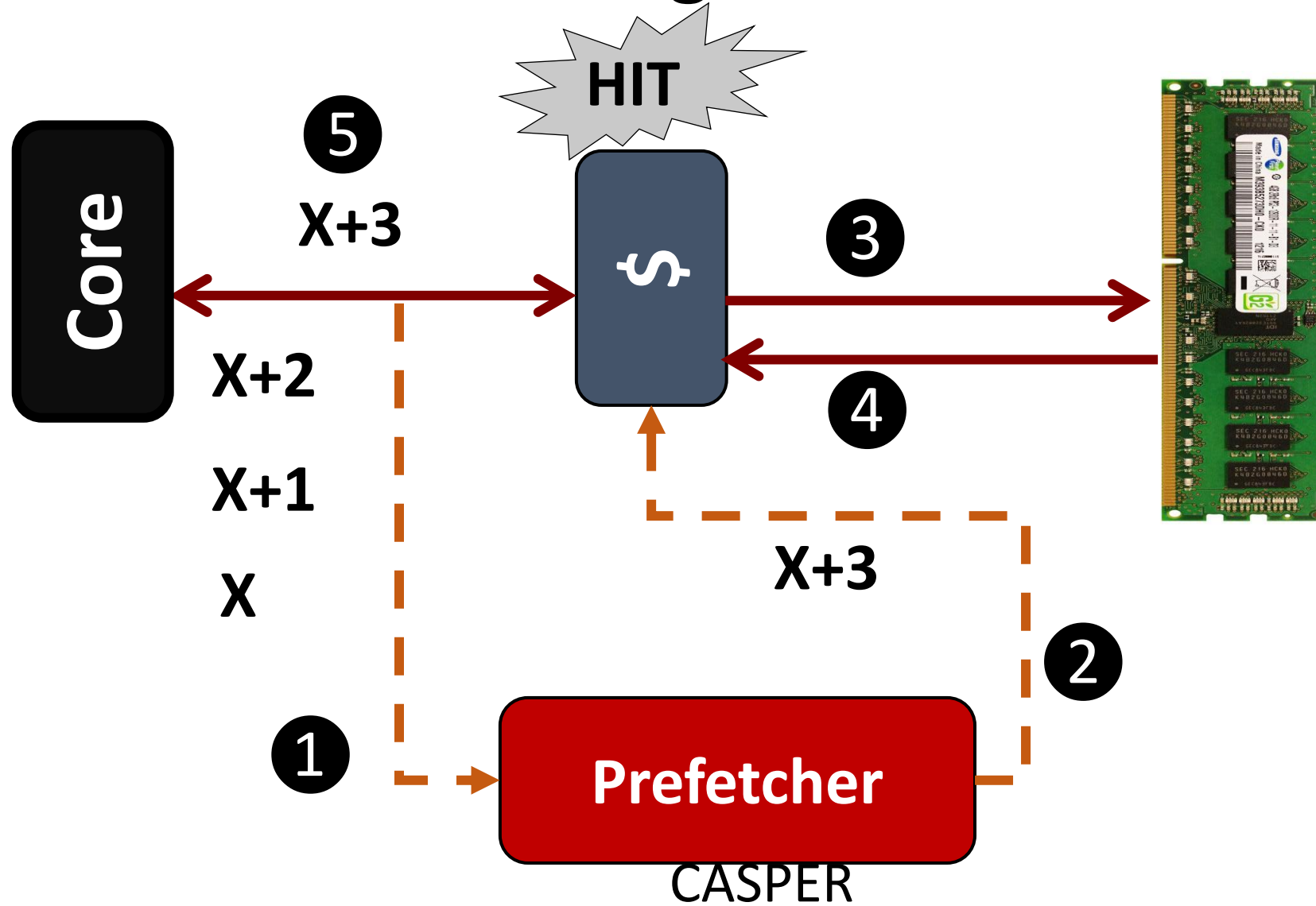*Hit time + Miss rate x Miss penalty*

*Hit time-L1 + Miss rate-L1 x Miss penalty-L1*
Miss penalty-L1 = Hit time-L2 + Miss rate-L2 x Miss penalty-L2
*Hit time: Low, Miss rate: Low , Miss penalty: Low*

*Ideally, miss rate = 0.00% and hit time should be one cycle, so all LOADs will take just one cycle* ☺

# Hardware Prefetching



CASPER

17

# 10K Feet View

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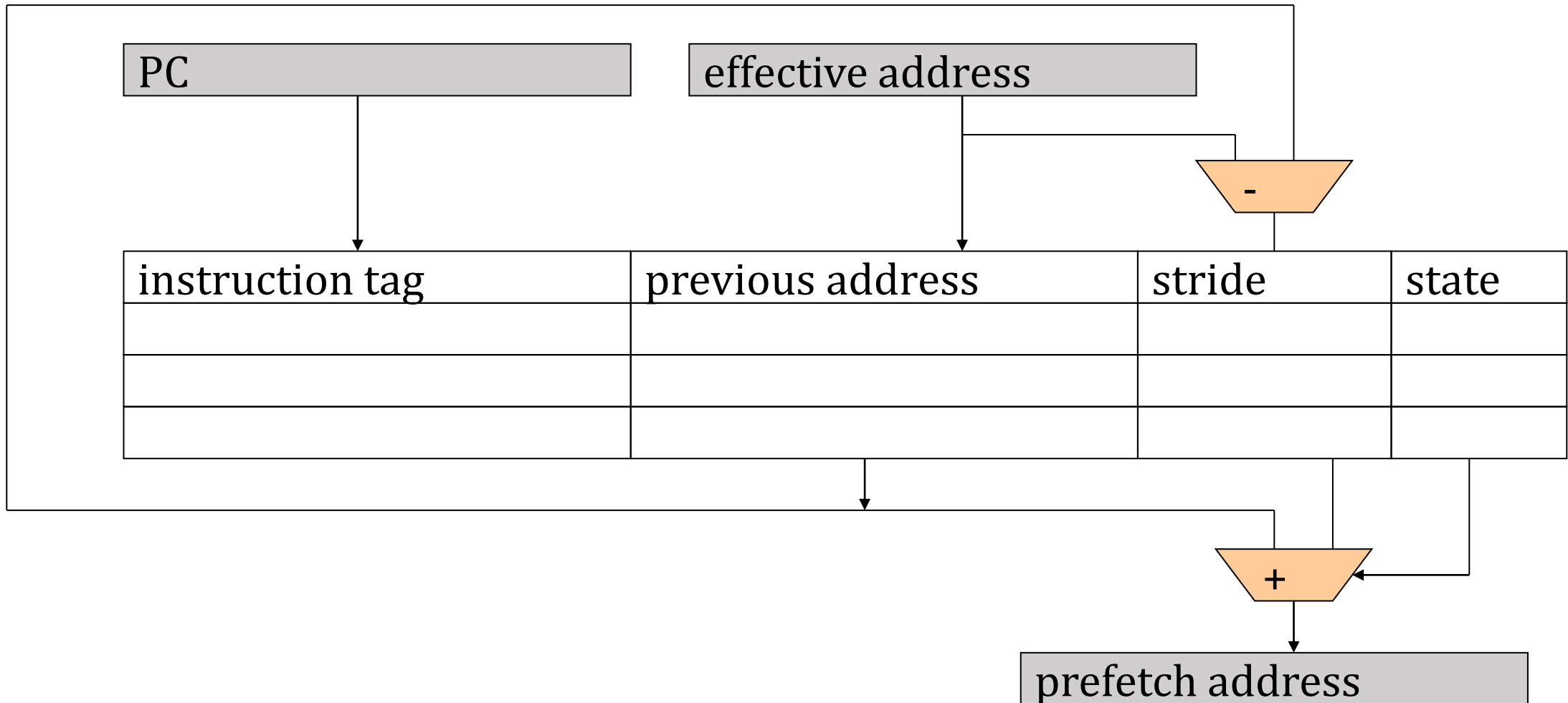Latency-hiding technique - Fetches data before the core demands.

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

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

# IP-stride prefetcher



| instruction tag | previous address | stride | state |
|---|---|---|---|
|  |  |  |  |
|  |  |  |  |
|  |  |  |  |

PC

effective address

-

+

prefetch address

CASPER                                                                 19

# Time for Assignment-1 (Behind every urgency there is lack of planning)