# Future of Computing

15-213/18-213/14-513/15-513/18-613: Introduction to Computer Systems
28th Lecture, December 5, 2019

# Today

- **Writeback-Aware Caching (guest lecture by Charles McGuffey)**
- Systems for Machine Learning (guest lecture by Angela Jiang)
- Prescriptive Memory

# Writeback Aware Caching
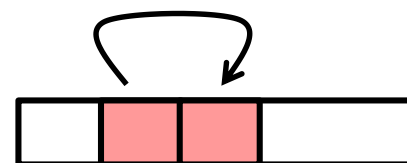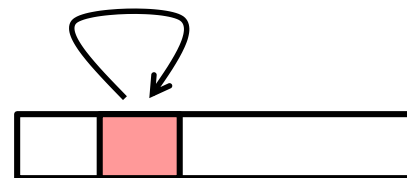
## Charles McGuffey

Nathan Beckmann, Phillip Gibbons
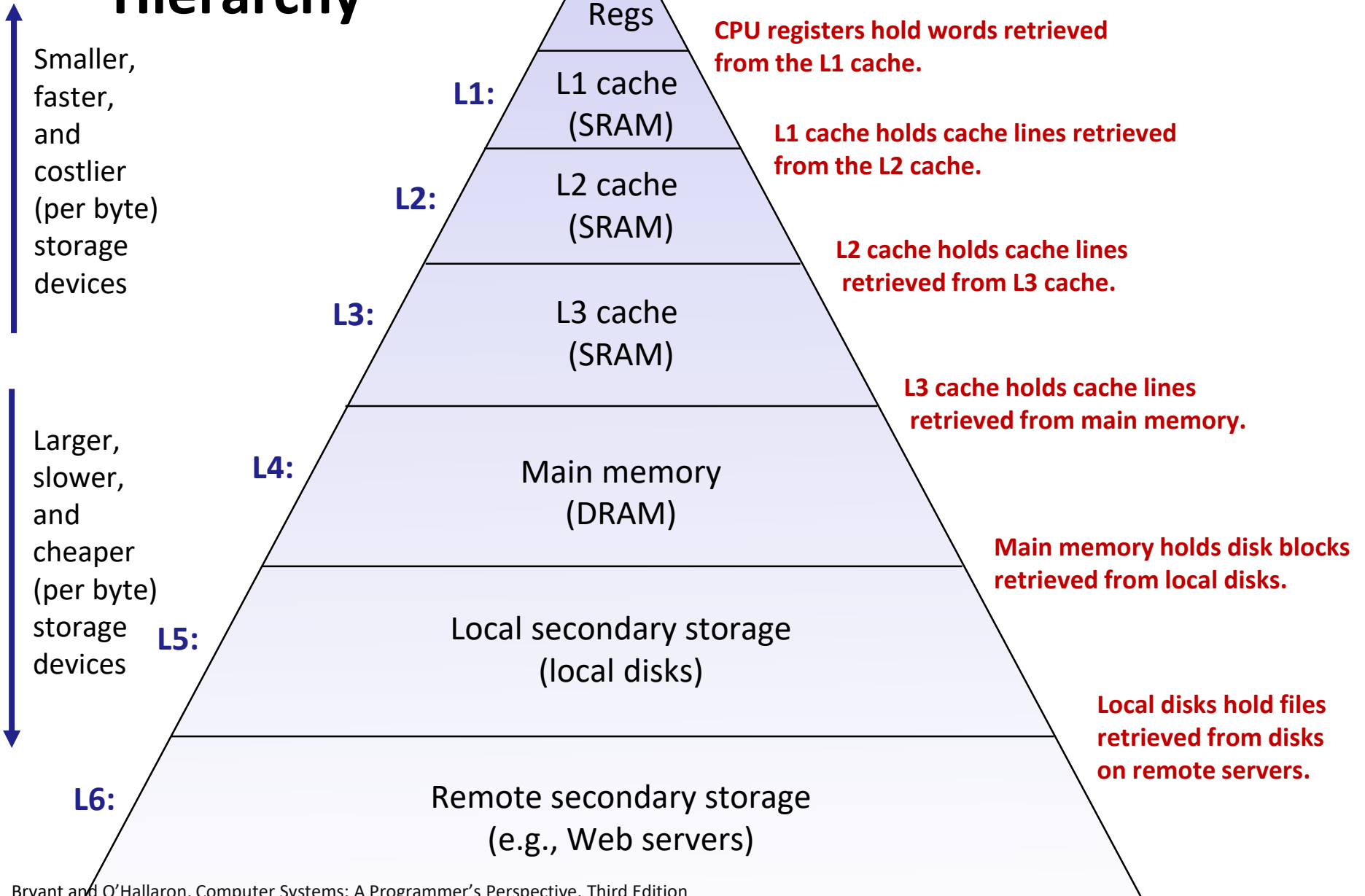
PARALLEL DATA LABORATORY
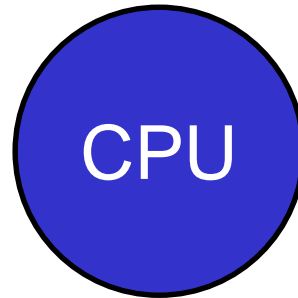
Carnegie Mellon University

# Recall: Locality

- **Principle of Locality: Programs tend to use data and instructions with addresses near or equal to those they have used recently**

- **Temporal locality:**
  - Recently referenced items are likely to be referenced again in the near future

- **Spatial locality:**
  - Items with nearby addresses tend to be referenced close together in time
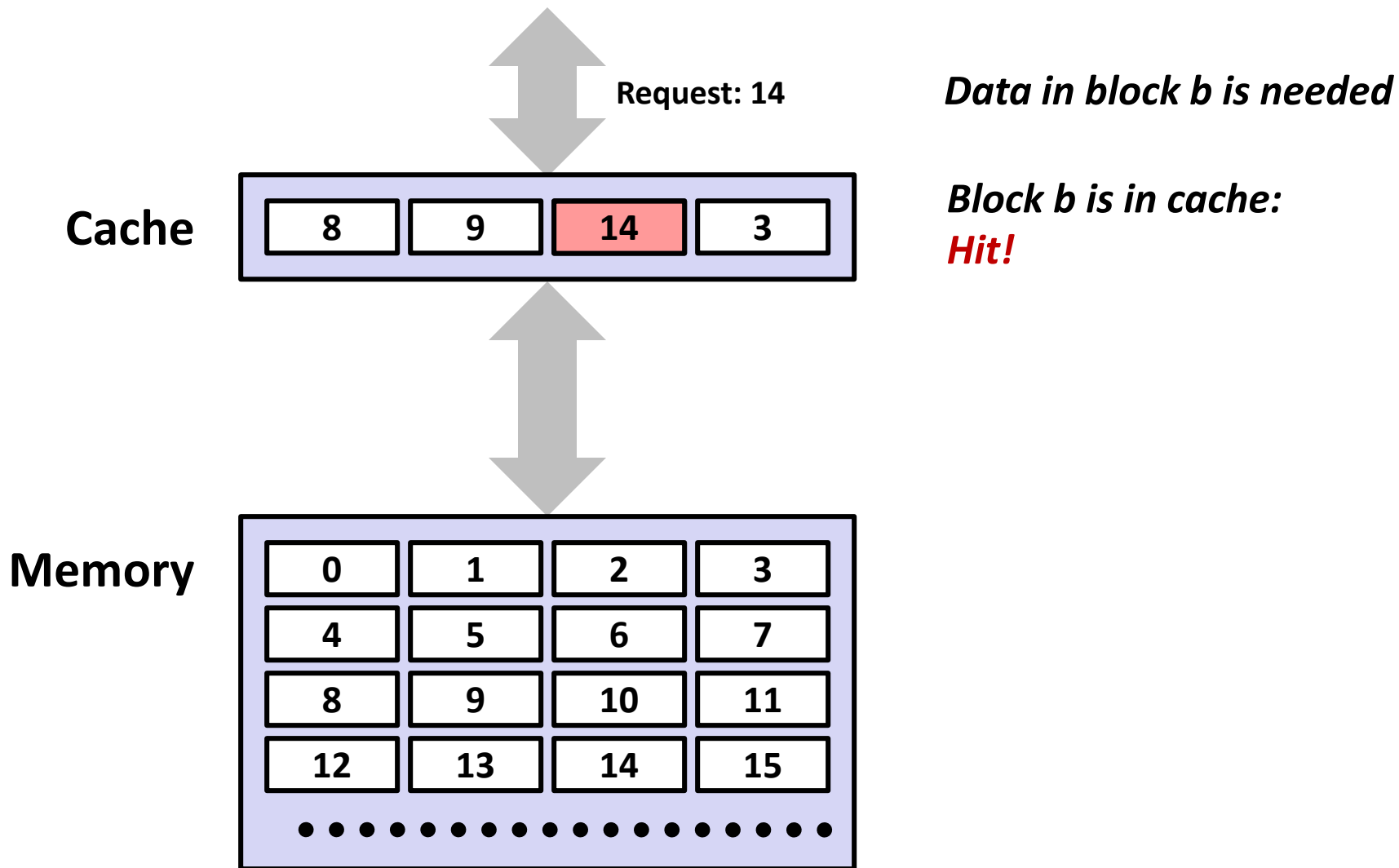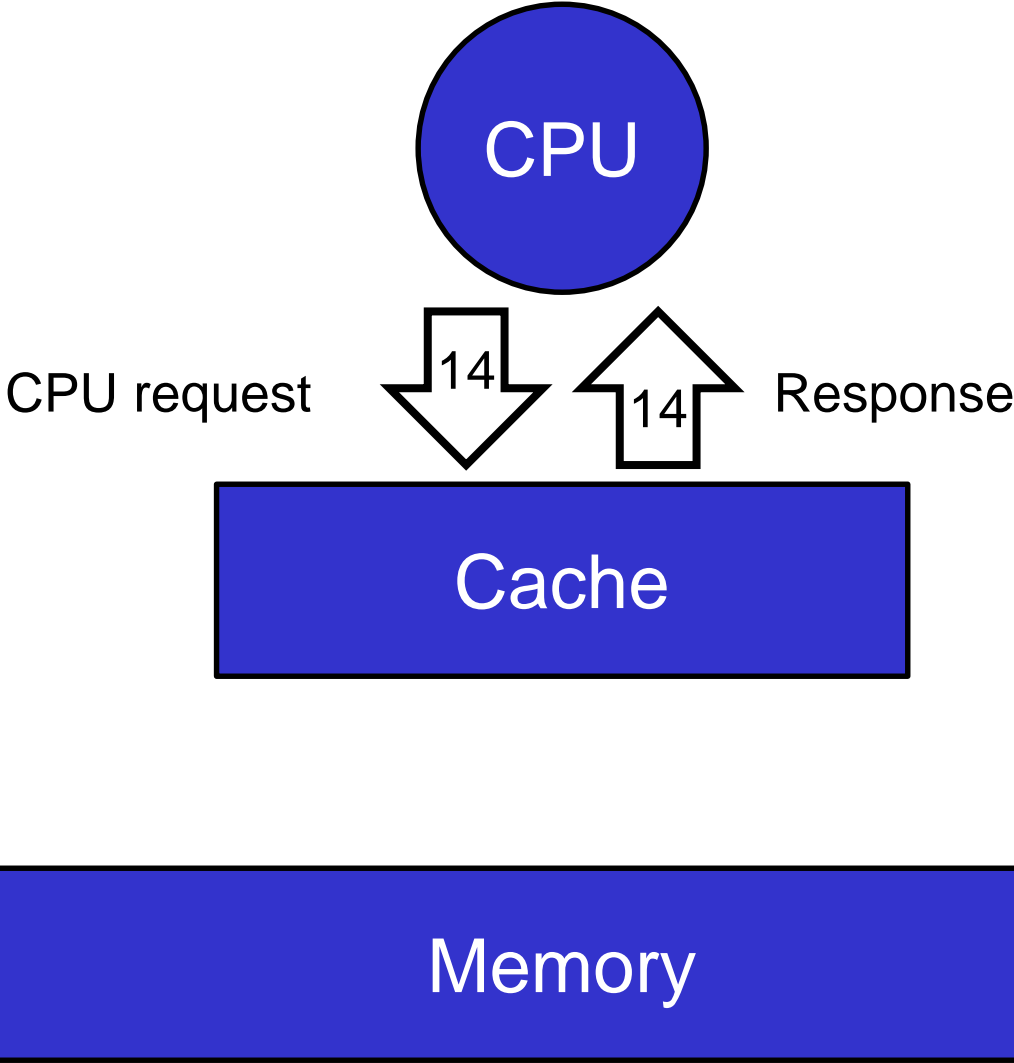
# Recall: Memory Hierarchy

Smaller, faster, and costlier (per byte) storage devices

Larger, slower, and cheaper (per byte) storage devices

**L0:** Regs

**L1:** L1 cache (SRAM)

**L2:** L2 cache (SRAM)

**L3:** L3 cache (SRAM)

**L4:** Main memory (DRAM)

**L5:** Local secondary storage (local disks)

**L6:** Remote secondary storage (e.g., Web servers)

CPU registers hold words retrieved from the L1 cache.

L1 cache holds cache lines retrieved from the L2 cache.

L2 cache holds cache lines retrieved from L3 cache.

L3 cache holds cache lines retrieved from main memory.

Main memory holds disk blocks retrieved from local disks.

Local disks hold files retrieved from disks on remote servers.

# Recall: General Cache Concepts

**Cache**

| 4 | 9 | 14 | 3 |
|---|---|----|---|

**Smaller, faster, more expensive memory caches a subset of the blocks**

| 4 |
|---|

**Data is copied in block-sized transfer units**

**Memory**

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

**Larger, slower, cheaper memory viewed as partitioned into "blocks"**

# Caching Model

**CPU**

**Cache**

**Memory**

**Carnegie Mellon**
**Parallel Data Laboratory**

# General Cache Concepts: Hit



**Request: 14**

*Data in block b is needed*

Cache

| 8 | 9 | 14 | 3 |

*Block b is in cache:*
*Hit!*

Memory

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

# Modeling Hits



**CPU**

CPU request    14    14    Response

**Cache**

**Memory**

# General Cache Concepts: Miss

**Request: 12**

*Data in block b is needed*

**Cache**

| 8 | 12 | 14 | 3 |
|---|----|----|---|

*Block b is not in cache:*
*Miss!*

| 12 |
|----|

**Request: 12**

*Block b is fetched from memory*

**Memory**

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

*Block b is stored in cache*
• Placement policy: determines where b goes
• Replacement policy: determines which block gets evicted (victim)

# Modeling Misses

# Cache Performance Metrics

- **Miss Rate**
  - Fraction of memory references not found in cache (misses / accesses) = 1 – hit rate
  - Typical numbers (in percentages):
    - 3-10% for L1
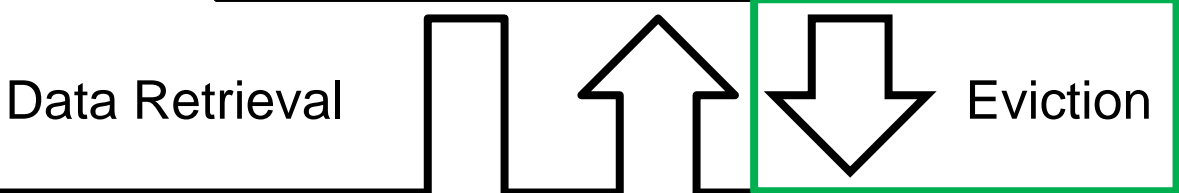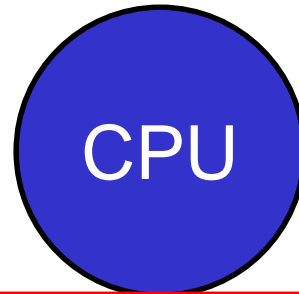    - can be quite small (e.g., < 1%) for L2, depending on size, etc.

- **Hit Time**
  - Time to deliver a line in the cache to the processor
    - includes time to determine whether the line is in the cache
  - Typical numbers:
    - 4 clock cycle for L1
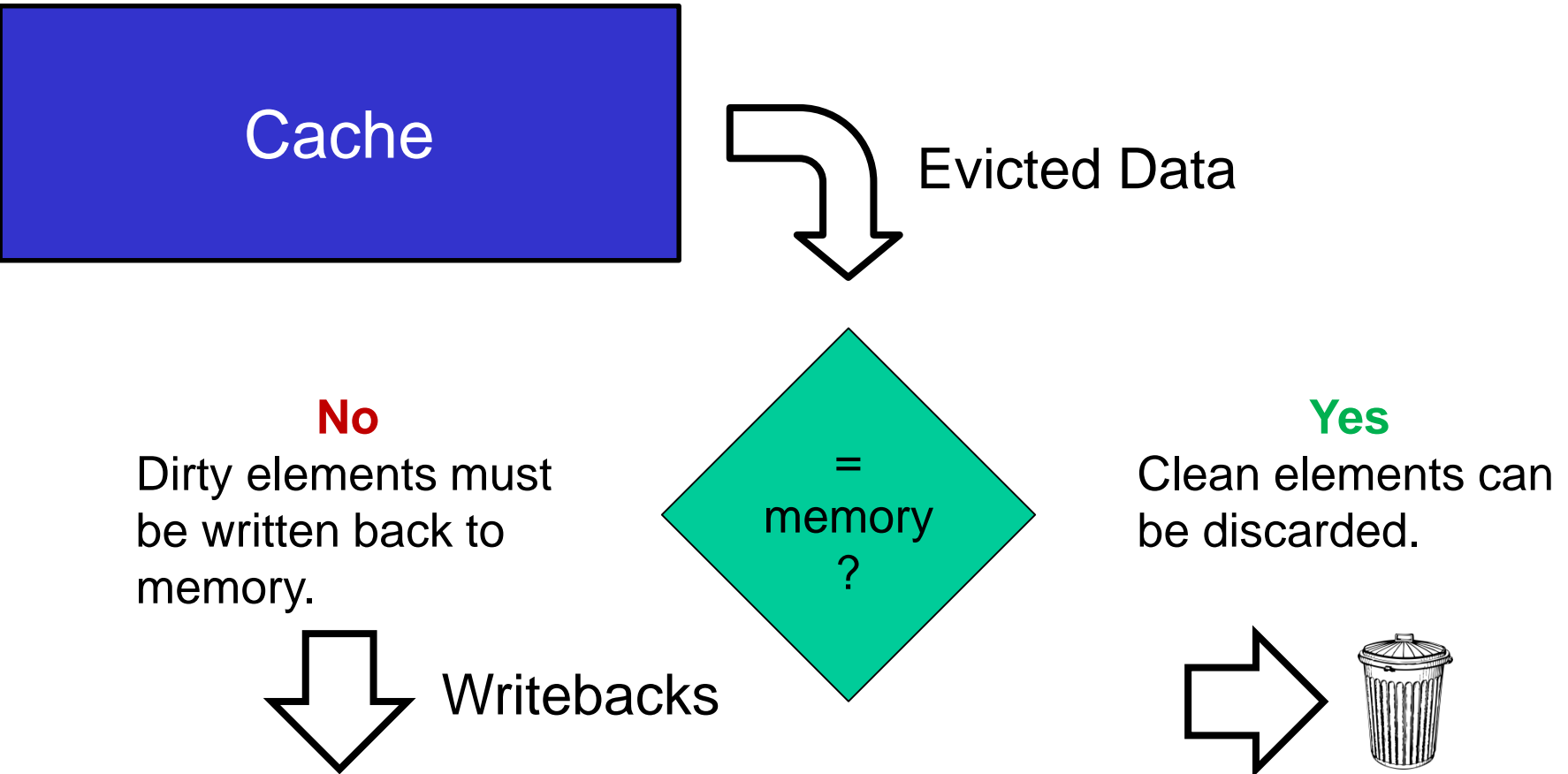    - 10 clock cycles for L2

- **Miss Penalty**
  - Additional time required because of a miss
    - typically 50-200 cycles for main memory (Trend: increasing!)

# Modeling Miss Rate



CPU

Unavoidable CPU request

Response

Cache

Data Retrieval

Eviction

Memory

**Carnegie Mellon**
**Parallel Data Laboratory**

Charles McGuffey © December 19

# Evictions

Cache

Evicted Data

**No**
Dirty elements must be written back to memory.

= memory ?

**Yes**
Clean elements can be discarded.

Writebacks

**Carnegie Mellon**
**Parallel Data Laboratory**

# What about writes?

- **Multiple copies of data exist:**
  - L1, L2, L3, Main Memory, Disk

| v | d | tag | 0 | 1 | 2 | ······ | B-1 |

**valid bit** **dirty bit**

$B = 2^b$ **bytes**

- **What to do on a write-hit?**
  - Write-through (write immediately to memory)
  - Write-back (defer write to memory until replacement of line)
    - Each cache line needs a dirty bit (set if data differs from memory)

- **What to do on a write-miss?**
  - Write-allocate (load into cache, update line in cache)
    - Good if more writes to the location will follow
  - No-write-allocate (writes straight to memory, does not load into cache)

- **Typical**
  - Write-through + No-write-allocate   **1 memory write / CPU write**
  - **Write-back + Write-allocate**

# Combining Writebacks

| R(1) | W(2) | R(3) | R(1) | W(2) | R(4) | Flush |
|------|------|------|------|------|------|-------|

Requests

| | 1 | 1 | 1 | 1 | 1 | 1 |
|---|---|---|---|---|---|---|
| | | 2 | 3 | 3 | 2 | 4 |

Loads
Writebacks

| | 1 | 1 | 3 | 1 | 1 | 4 |
|---|---|---|---|---|---|---|
| | | 2 | 2 | 2 | 2 | 2 |

Time

**Carnegie Mellon**
**Parallel Data Laboratory**

# Why Writebacks Matter

- Bandwidth
- Energy
- Wearout

These metrics are important to non-volatile memory and storage.

Cache

Memory

**Carnegie Mellon**
**Parallel Data Laboratory**

# A New Metric

A weighted sum of the cost due to loads and the cost due to writebacks.



$$\sum_{\forall\ misses\ p} Cost_{Load}(p) + weight \times \sum_{\forall\ WBs\ q} Cost_{WB}(q)$$
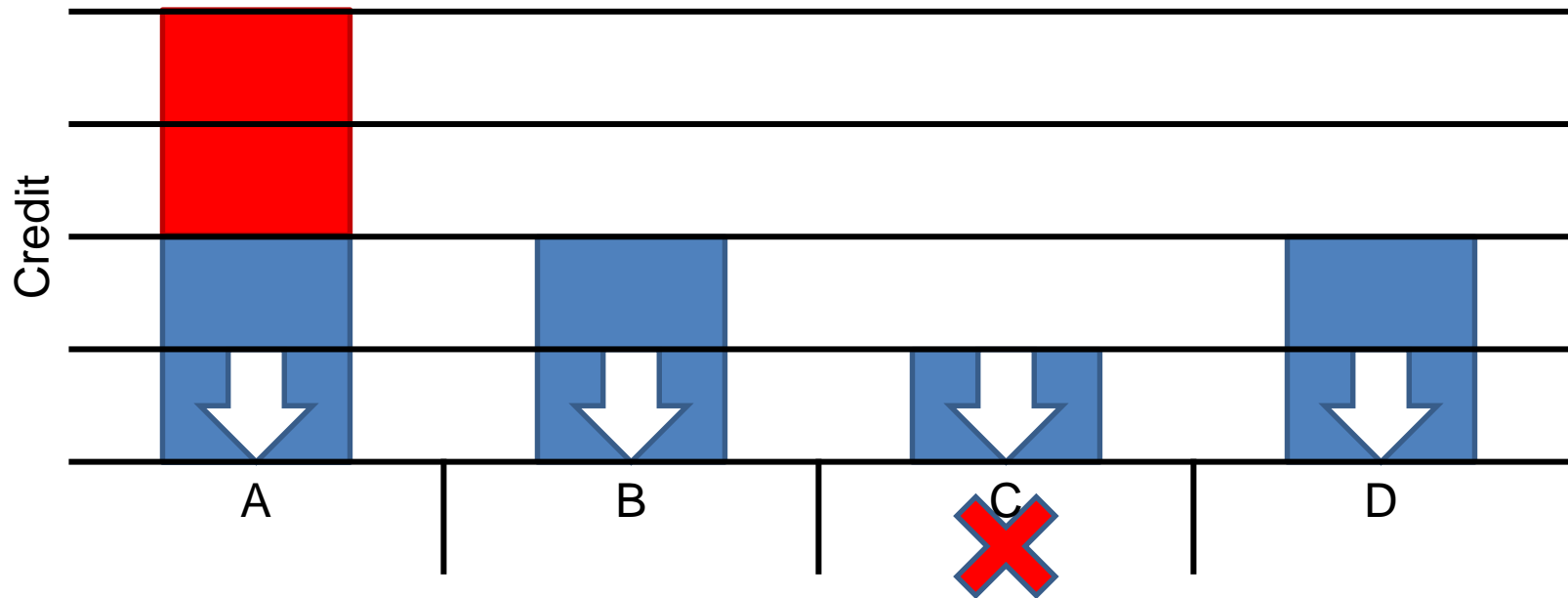
# Versions of the Problem

| R(a) | R(b) | W(c) |
|------|------|------|

Requests

Cache

a

**Online Caching:**
A request must be served before the next appears.

**Offline Caching:**
All requests are known in advance.

**Carnegie Mellon**
**Parallel Data Laboratory**

# Summary of Results

- Modeling Writebacks

  - Theoretical model that generalizes caching

- Offline Problem:

  - Analysis of writeback-oblivious policies

  - Complexity results

  - Approximation algorithms

- Online Policy: Writeback-Aware Landlord

  - Optimal worst-case analysis

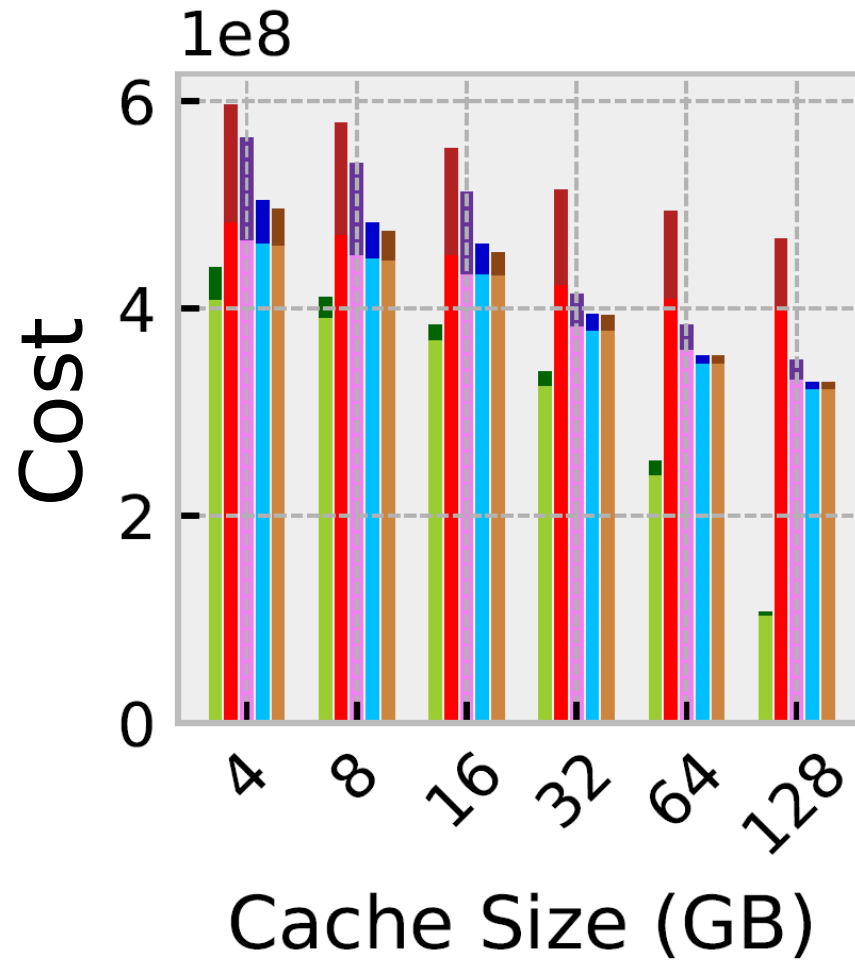  - Good empirical performance

**Carnegie Mellon**
**Parallel Data Laboratory**

# Writeback-Aware Landlord

# Cache Simulation

- ## MSR storage traces
  - ### Disk requests for real MSR servers

- ## Competing policies:
  - ### Offline approximation (WAPFOO-L)
  - ### Least Recently Used (LRU)
  - ### Greedy-Dual Size (GDS)

- ## Metric:
  - ### Loads have unit cost
  - ### Writes have 10x the cost

# Simulation Results on src1_1

# More Results

- ## Other MSR traces

  - ### Shows broader impact

- ## Applying the frequency metric

  - ### Traditionally useful heuristic helps here

- ## Sensitivity to write/read cost ratio

  - ### Generally performs well

**Carnegie Mellon**
**Parallel Data Laboratory**

# Summary of Results

- Modeling Writebacks

  - Theoretical model that generalizes caching

- Offline Problem:

  - Analysis of writeback-oblivious policies

  - Complexity results

  - Approximation algorithms

- Online Policy: Writeback-Aware Landlord

  - Optimal worst-case analysis

  - Good empirical performance

**Carnegie Mellon**
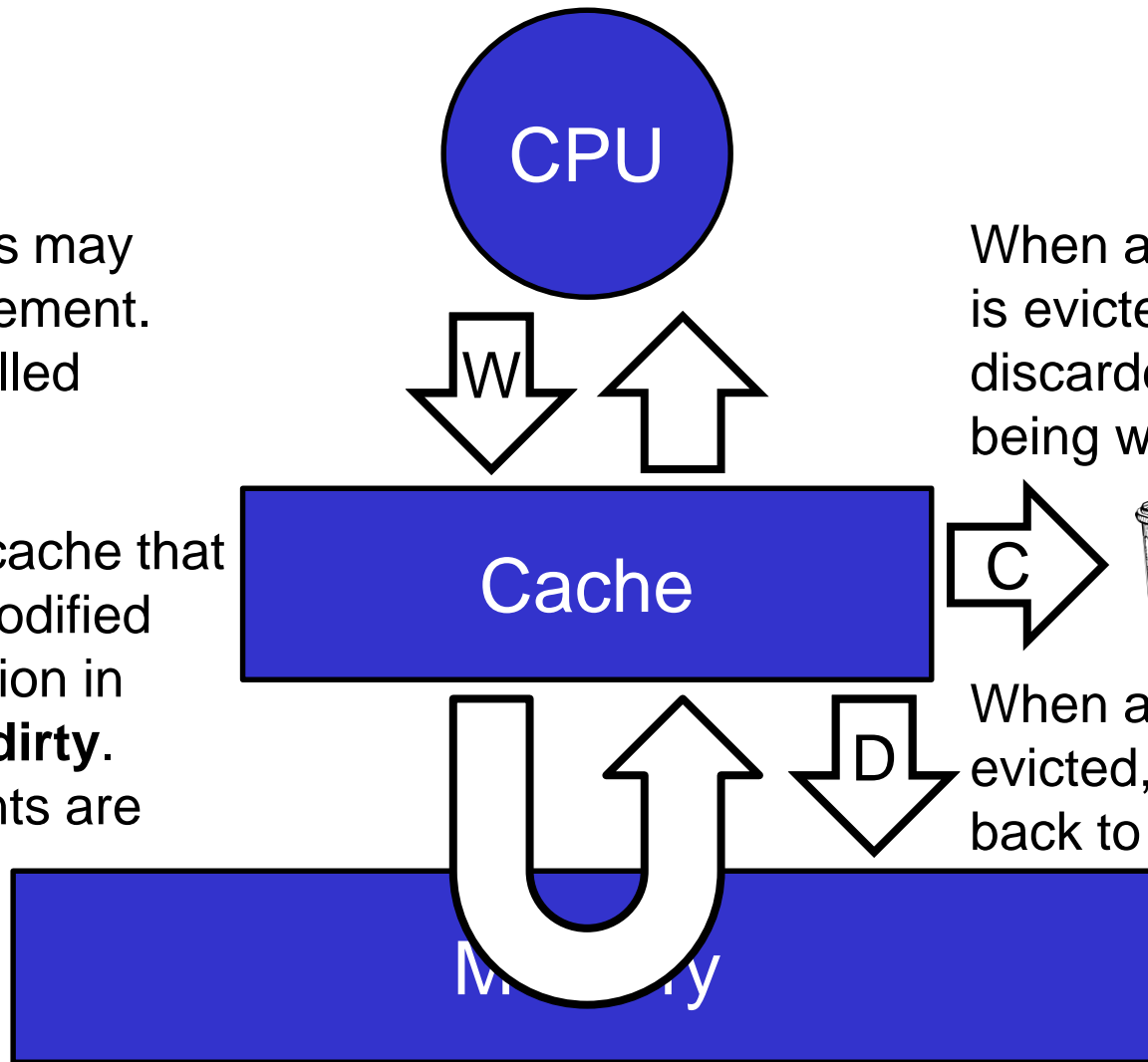**Parallel Data Laboratory**

# More Caching Research

- This work

- Cache coherence
  - Ensuring data consistency

- Parallel Caching
  - Multiple CPUs sharing a cache

- Distributed Caching
  - Data can be in different locations

- Web Caching
  - Objects have different size and cost

# Backup Slides

**Carnegie Mellon**
**Parallel Data Laboratory**

# Writebacks

CPU

CPU requests may modify the element. These are called **writes**.

Elements in cache that have been modified from the version in memory are **dirty**. Other elements are **clean**.
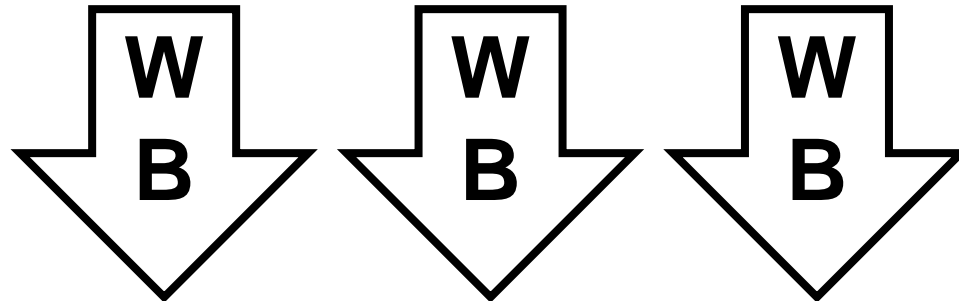
W

Cache

C

When a clean element is evicted, it can be discarded without being written back.

D

When a dirty element is evicted, it is written back to memory.

Memory

# Combining Writebacks



D F 5 6 2 9 E B

# Combining Writebacks

R(a), W(x), R(y), R(b), W(x), R(y), R(c), W(x), R(y)

Requests

| x | a | x | x | b | x | x | c | x | x |
|---|---|---|---|---|---|---|---|---|---|
| y | y | y | y | y | y | y | y | y | y |

| x | x | x | x | x | x | x | x | x | x |
|---|---|---|---|---|---|---|---|---|---|
| y | a | a | y | b | b | y | c | c | y |

Writebacks

Loads

Time

**Carnegie Mellon**
**Parallel Data Laboratory**

# Why Writebacks Matter

**Bandwidth**
Memory has to take time to service writebacks.

**Energy**
The system must spend energy to service writebacks.

CPU

Cache

Memory

**Wearout**
Reducing writebacks reduces wear on persistent storage.

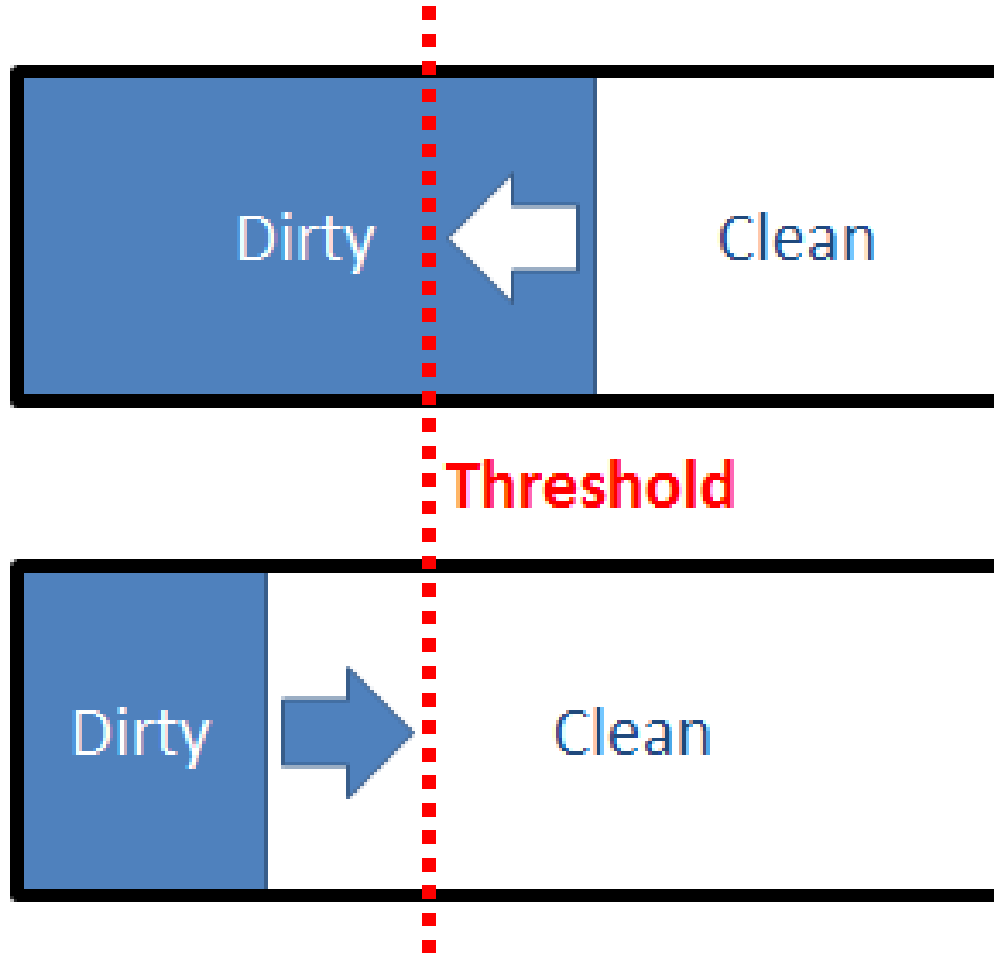Non-Volatile Memory makes these more important!

# Considering Loads AND Writebacks
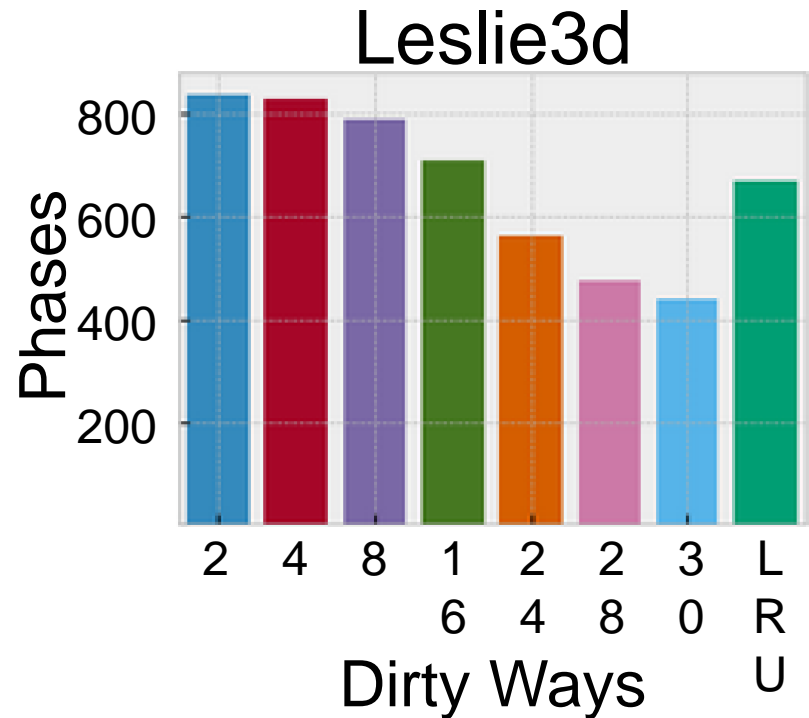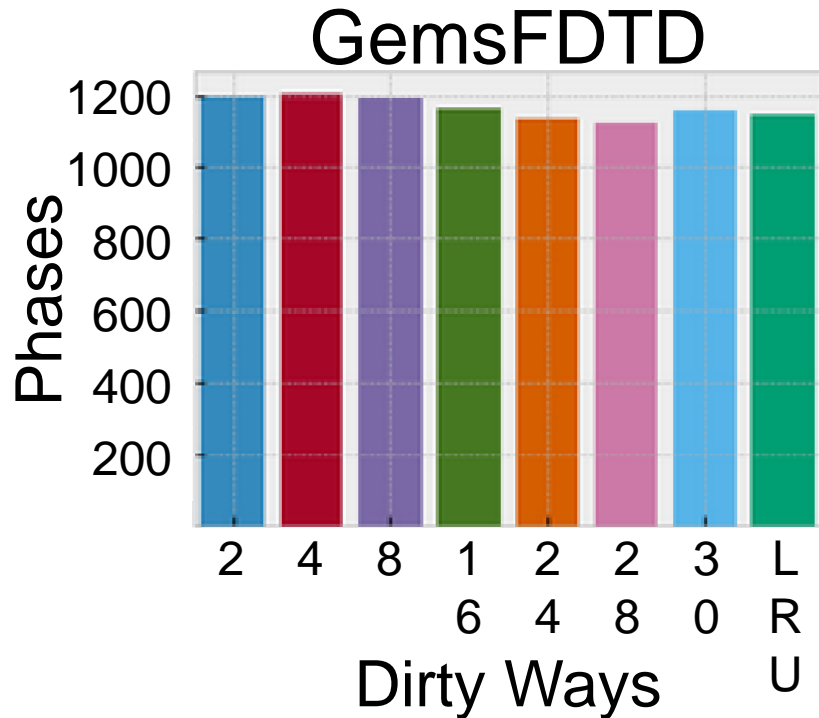
Traditional caches minimize:

$$\sum_{\forall\ misses\ p} Cost_{Load}(p)$$

Our work seeks to minimize:

$$\sum_{\forall\ misses\ p} Cost_{Load}(p) + weight \times \sum_{\forall\ WBs\ q} Cost_{WB}(q)$$

# Cache Partitioning

**Carnegie Mellon**
**Parallel Data Laboratory**

# Initial Results

## GemsFDTD



## Leslie3d

**Carnegie Mellon**
**Parallel Data Laboratory**
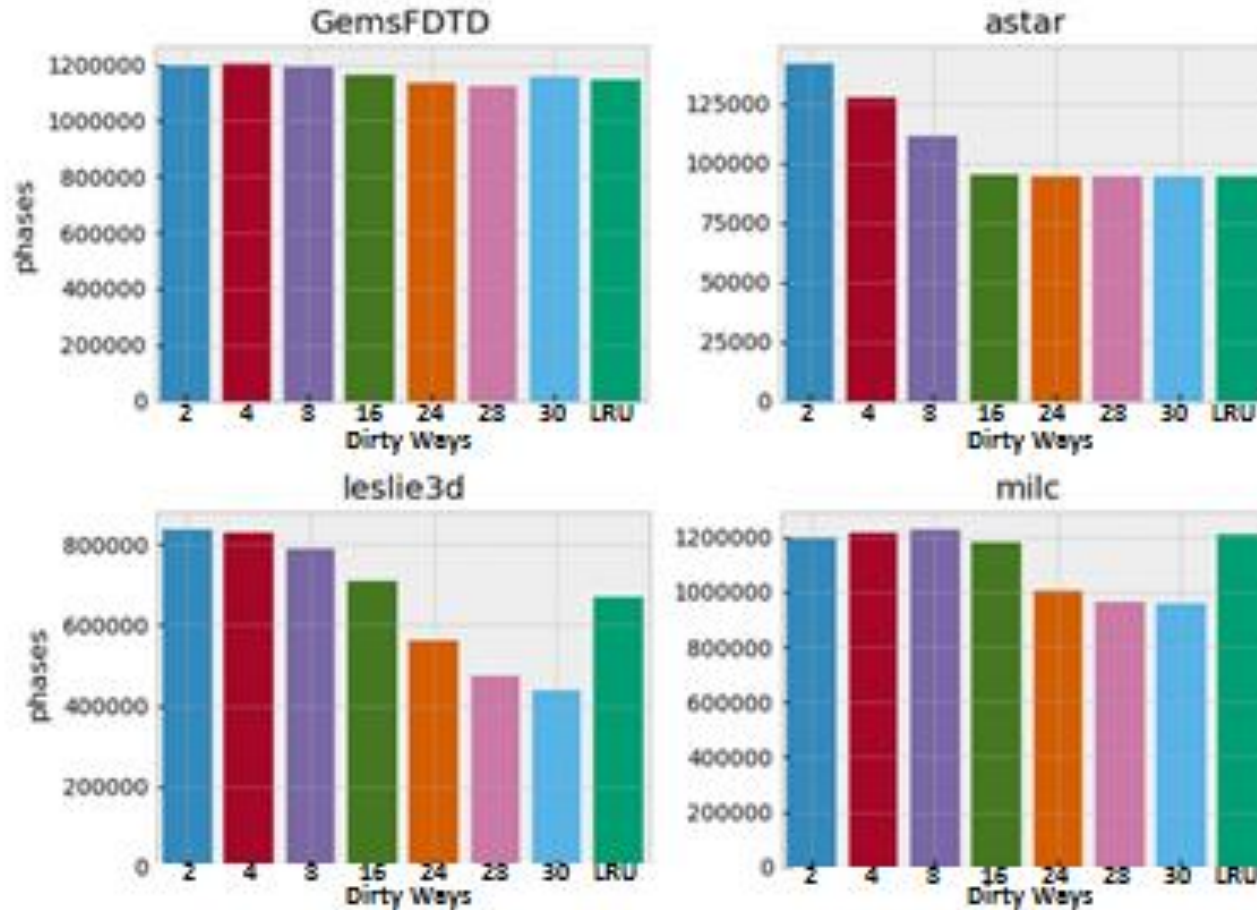
# Initial Results

# Summary of Results

- Offline Problem:
  - Non-optimality of traditional policies
  - Proof of NP-Completeness
  - Proof of APX-Completeness
  - Approximation Algorithms
- Online Problem:
  - Dirty/Clean Cache Partitioning
  - WA-RRIP
  - WA Landlord

**Come see our poster for more details!**