# Memory Hierarchy and Cache Memories

CENG331  - Computer Organization

**Instructor:**

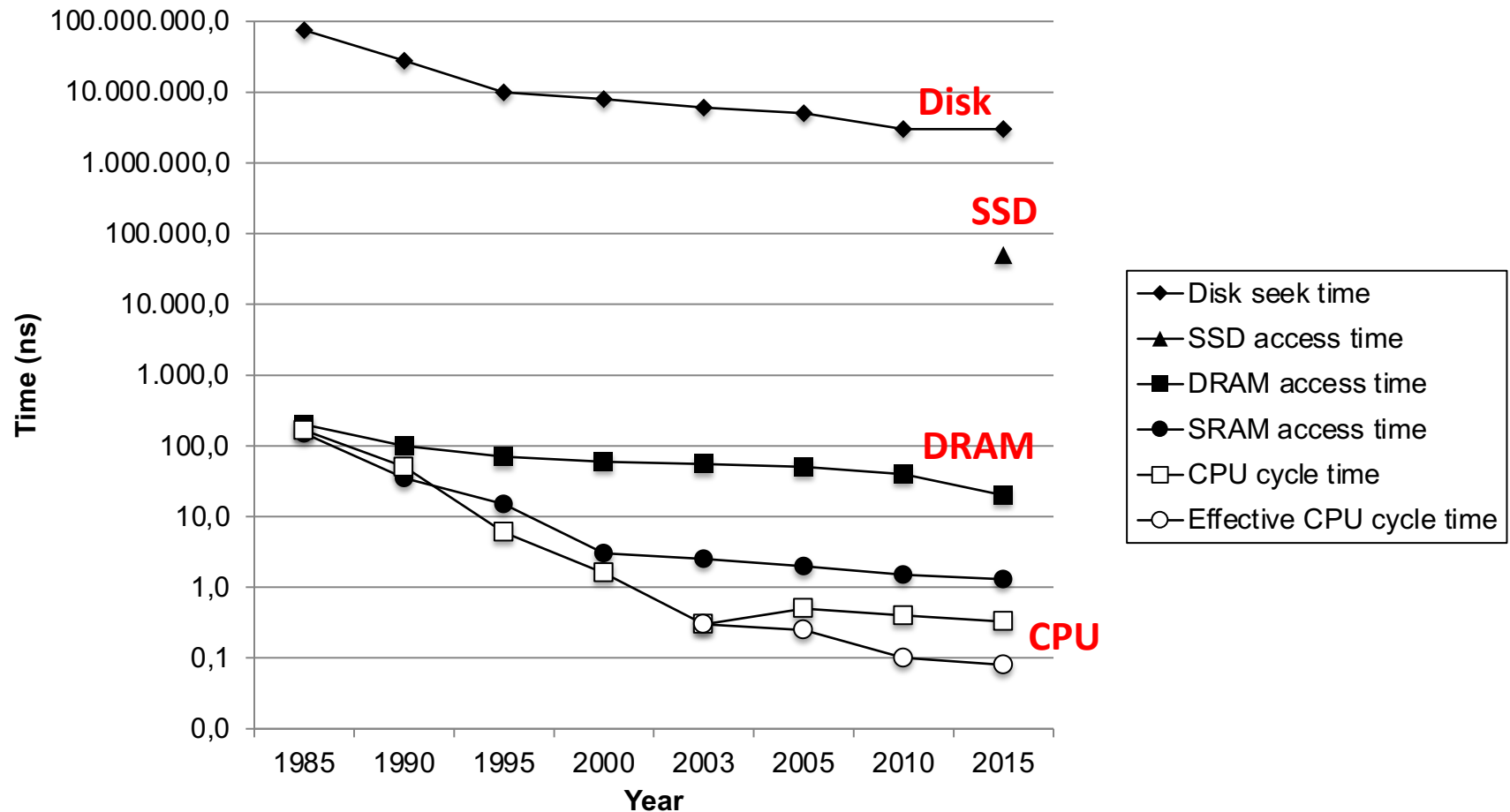Murat Manguoğlu          (Section 1)

# Today

- **Cache memory organization and operation**
- **Performance impact of caches**
  - The memory mountain
  - Rearranging loops to improve spatial locality
  - Using blocking to improve temporal locality

# The CPU-Memory Gap

**The gap widens between DRAM, disk, and CPU speeds.**

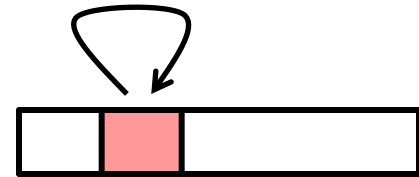# Locality to the Rescue!

**The key to bridging this CPU-Memory gap is a fundamental property of computer programs known as <span style="color:red">locality</span>**

# Locality

- **Principle of Locality: Programs tend to use data and instructions with addresses near or equal to those they have used recently**
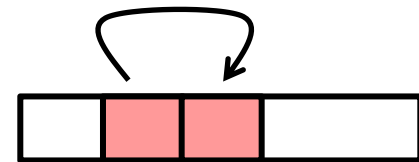
- **Temporal locality:**
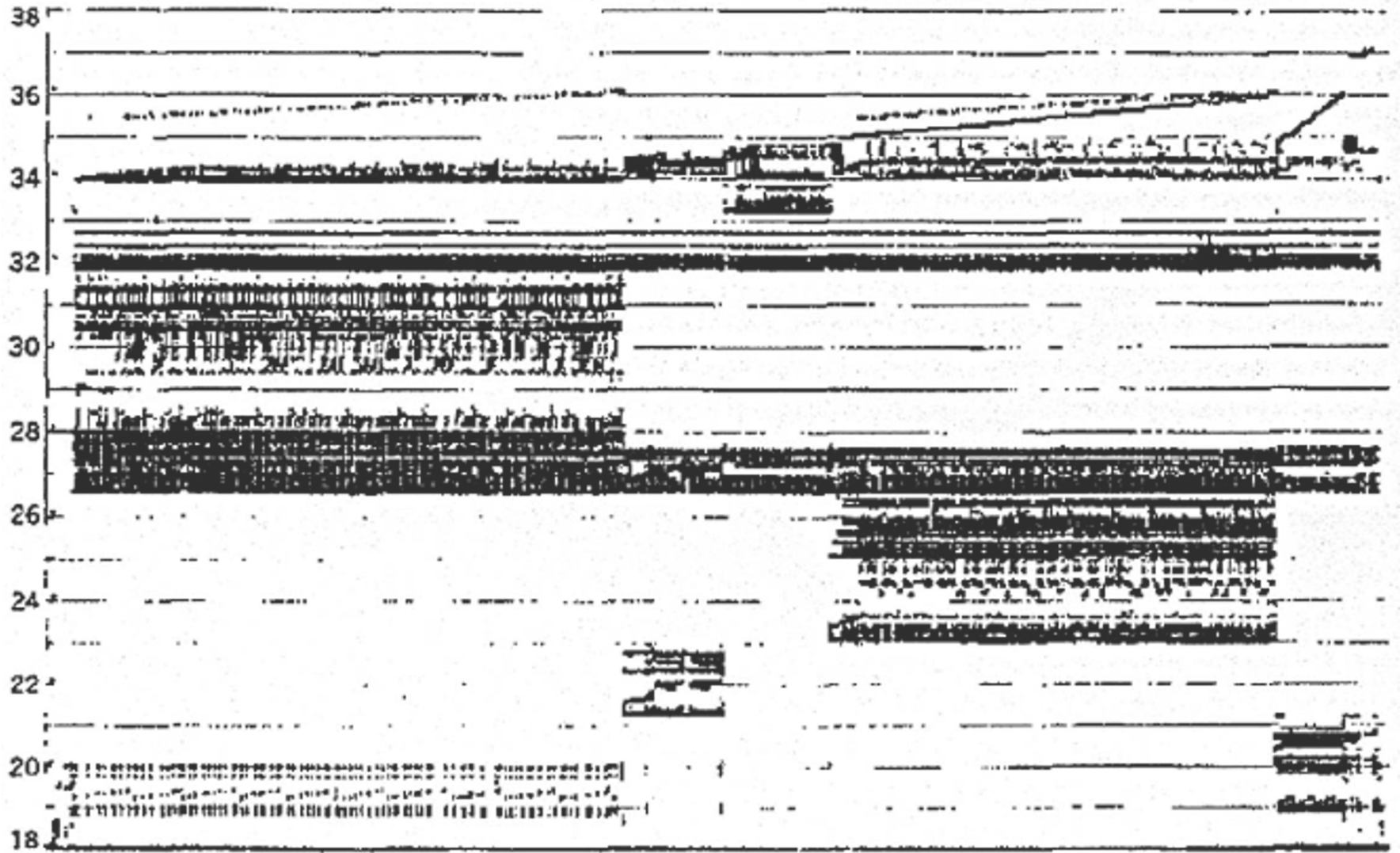    - Recently referenced items are likely to be referenced again in the near future

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

# Real Memory Reference Patterns



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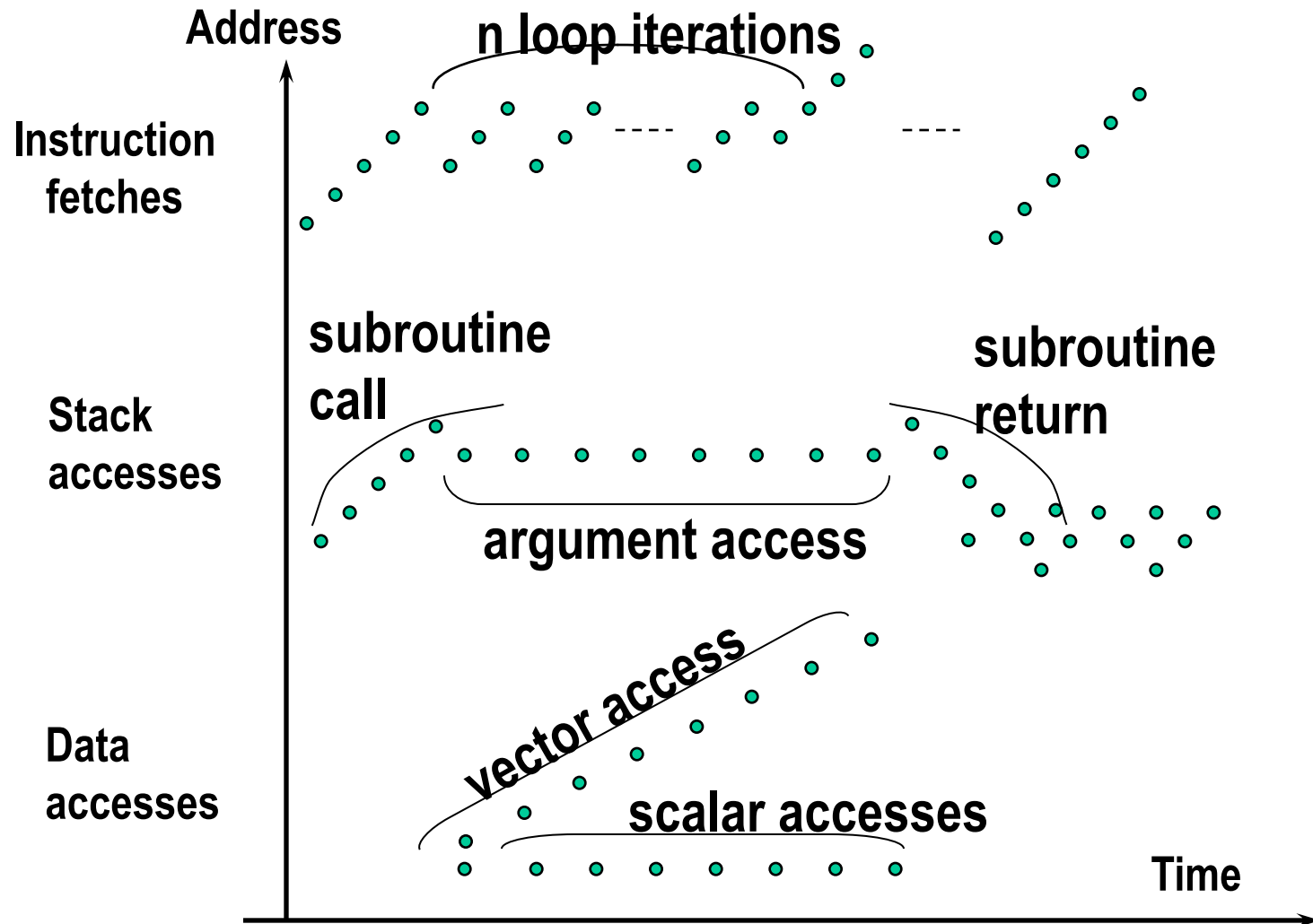
# Typical Memory Reference Patterns

# Memory Reference Patterns



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

# Locality Example

```
sum = 0;
for (i = 0; i < n; i++)
    sum += a[i];
return sum;
```

- **Data references**
  - Reference array elements in succession (stride-1 reference pattern).      **Spatial locality**
  - Reference variable `sum` each iteration.      **Temporal locality**
- **Instruction references**
  - Reference instructions in sequence.      **Spatial locality**
  - Cycle through loop repeatedly.      **Temporal locality**

# Locality Example

- **Question: Does this function have good locality with respect to array** $a$**?**

```
int sum_array_cols(int a[M][N])
{
    int i, j, sum = 0;

    for (j = 0; j < N; j++)
        for (i = 0; i < M; i++)
            sum += a[i][j];
    return sum;
}
```

# Locality Example

- **Question**: Can you permute the loops so that the function scans the 3-d array `a` with a stride-1 reference pattern (and thus has good spatial locality)?

```
int sum_array_3d(int a[M][N][N])
{
    int i, j, k, sum = 0;

    for (i = 0; i < M; i++)
        for (j = 0; j < N; j++)
            for (k = 0; k < N; k++)
                sum += a[k][i][j];
    return sum;
}
```

# Qualitative Estimates of Locality

■ **Claim: Being able to look at code and get a qualitative sense of its locality is a key skill for a professional programmer.**

■ **Question: Does this function have good locality with respect to array** $a$**?**

```
int sum_array_rows(int a[M][N])
{
    int i, j, sum = 0;

    for (i = 0; i < M; i++)
        for (j = 0; j < N; j++)
            sum += a[i][j];
    return sum;
}
```

# Memory Hierarchies

- **Some fundamental and enduring properties of hardware and software:**
  - Fast storage technologies cost more per byte, have less capacity, and require more power (heat!).
  - The gap between CPU and main memory speed is widening.
  - Well-written programs tend to exhibit good locality.

- **These fundamental properties complement each other beautifully.**

- **They suggest an approach for organizing memory and storage systems known as a memory hierarchy.**

# Example 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

# Management of Memory Hierarchy

- **Small/fast storage, e.g., registers**
  - Address usually specified in instruction
  - Generally implemented directly as a register file
    - but hardware might do things behind software's back, e.g., stack management, register renaming
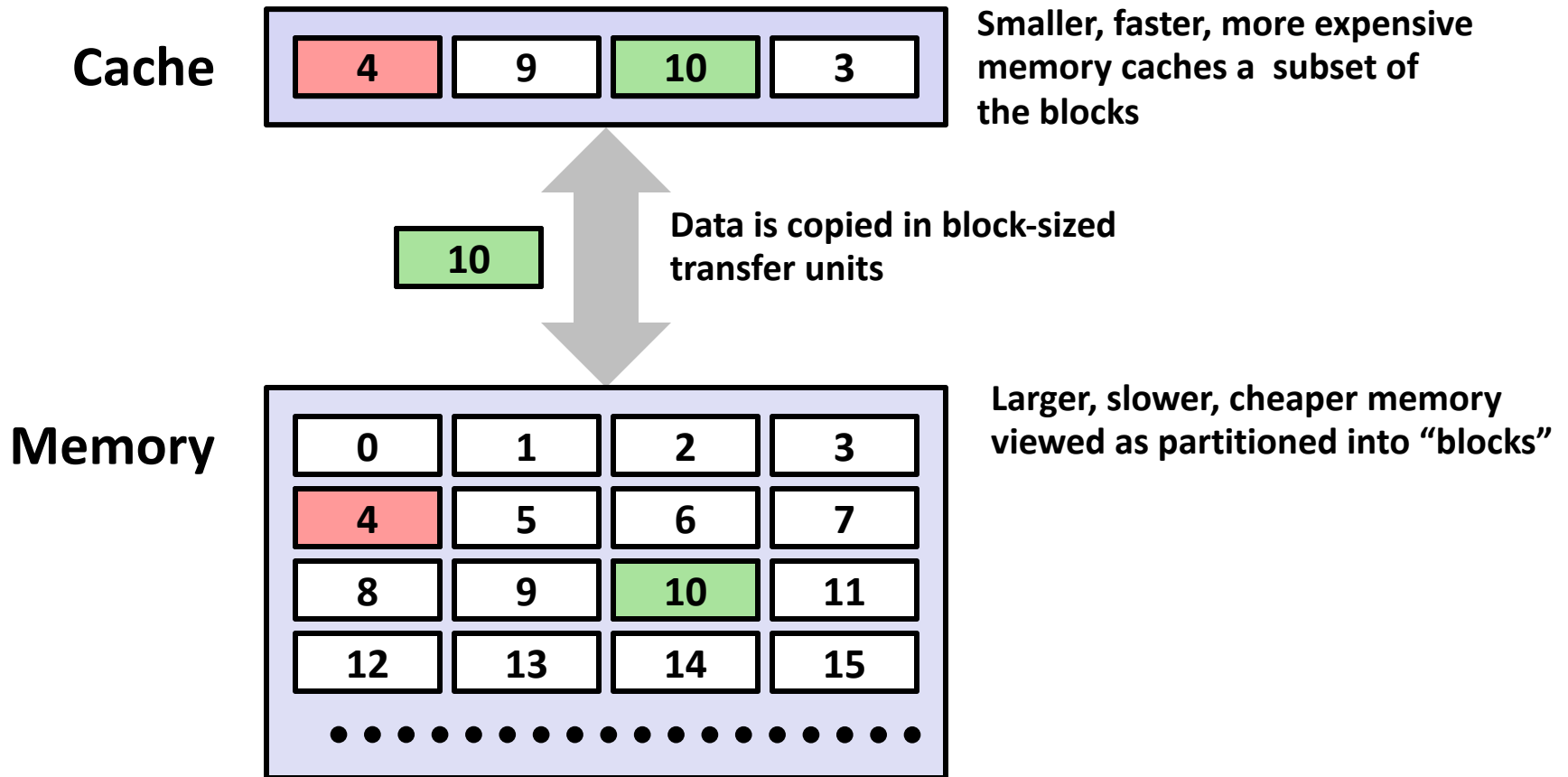
- **Larger/slower storage, e.g., main memory**
  - Address usually computed from values in register
  - Generally implemented as a hardware-managed cache hierarchy
    - hardware decides what is kept in fast memory
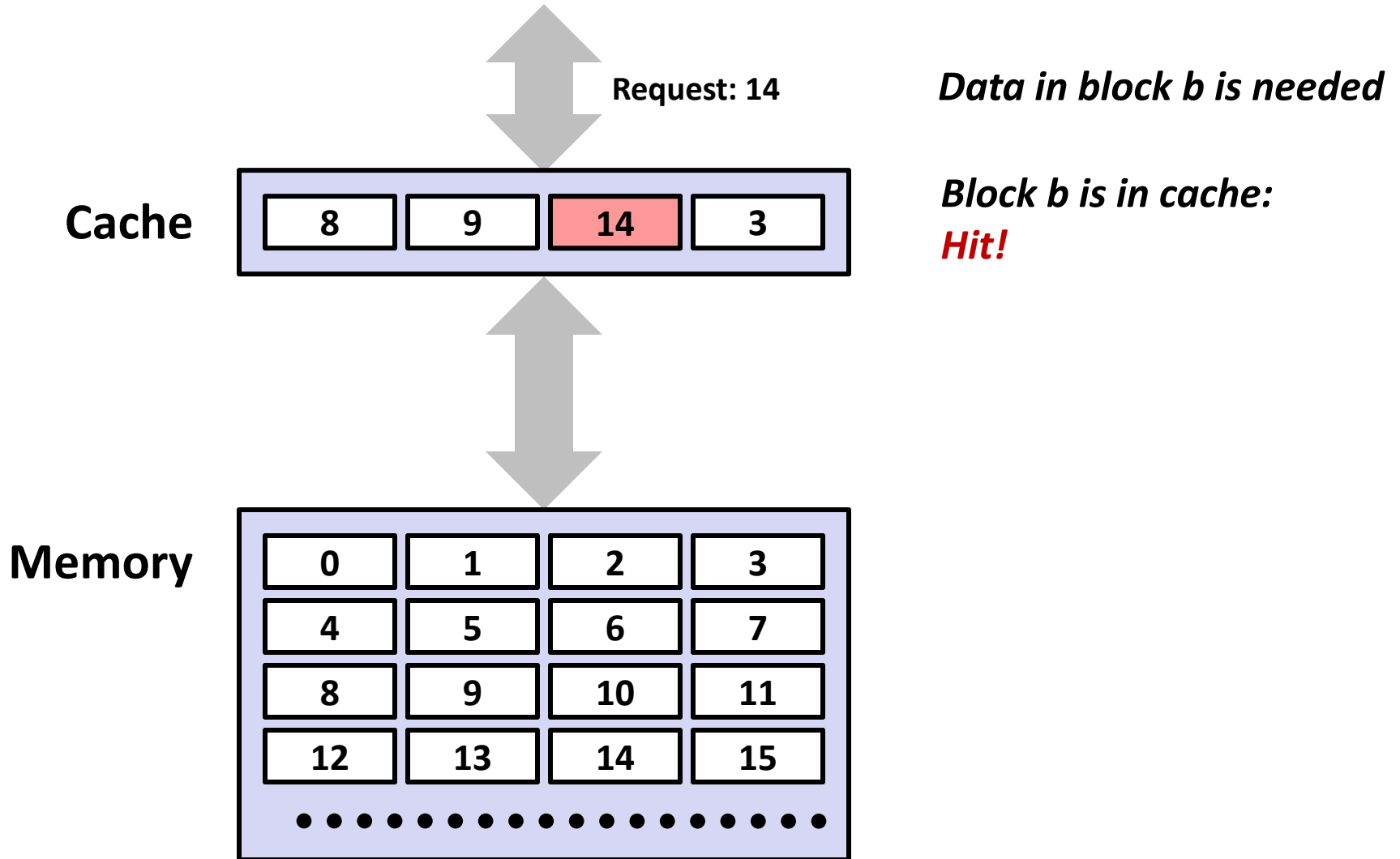    - but software may provide "hints", e.g., don't cache or prefetch

# Caches

- *Cache:* **A smaller, faster storage device that acts as a staging area for a subset of the data in a larger, slower device.**

- **Fundamental idea of a memory hierarchy:**
  - For each k, the faster, smaller device at level k serves as a cache for the larger, slower device at level k+1.

- **Why do memory hierarchies work?**
  - Because of locality, programs tend to access the data at level k more often than they access the data at level k+1.
  - Thus, the storage at level k+1 can be slower, and thus larger and cheaper per bit.

- *Big Idea:* **The memory hierarchy creates a large pool of storage that costs as much as the cheap storage near the bottom, but that serves data to programs at the rate of the fast storage near the top.**

# General Cache Concepts

**Cache**

| 4 | 9 | 10 | 3 |
|---|---|----|---|

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

| 10 |
|----|

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"

# General Cache Concepts: Hit

**Request: 14**

**Cache**

| 8 | 9 | 14 | 3 |
|---|---|----|---|

**Memory**

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

*Data in block b is needed*

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

# 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)

# General Caching Concepts: Types of Cache Misses

- **Cold (compulsory) miss**
  - Cold misses occur because the cache is empty.
- **Conflict miss**
  - Most caches limit blocks at level k+1 to a small subset (sometimes a singleton) of the block positions at level k.
    - E.g. Block i at level k+1 must be placed in block (i mod 4) at level k.
  - Conflict misses occur when the level k cache is large enough, but multiple data objects all map to the same level k block.
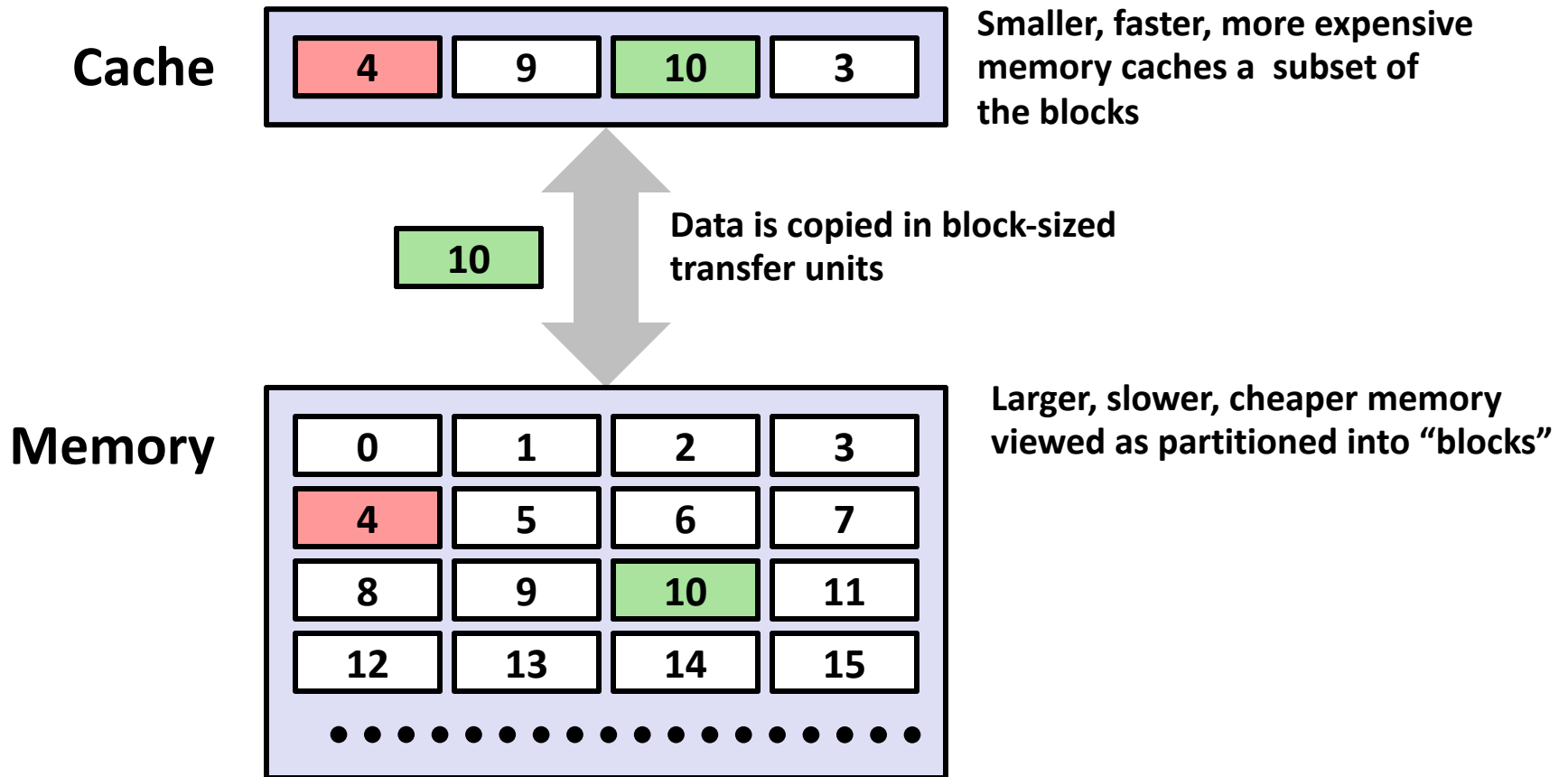    - E.g. Referencing blocks 0, 8, 0, 8, 0, 8, … would miss every time.
- **Capacity miss**
  - Occurs when the set of active cache blocks (working set) is larger than the cache.
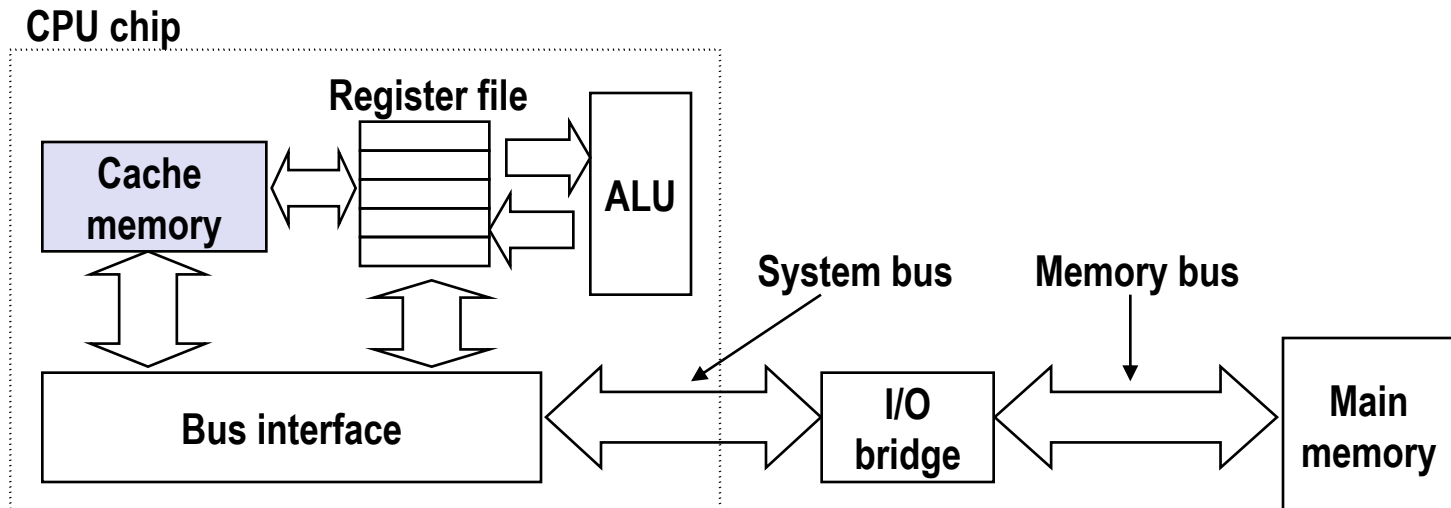
# Examples of Caching in the Mem. Hierarchy

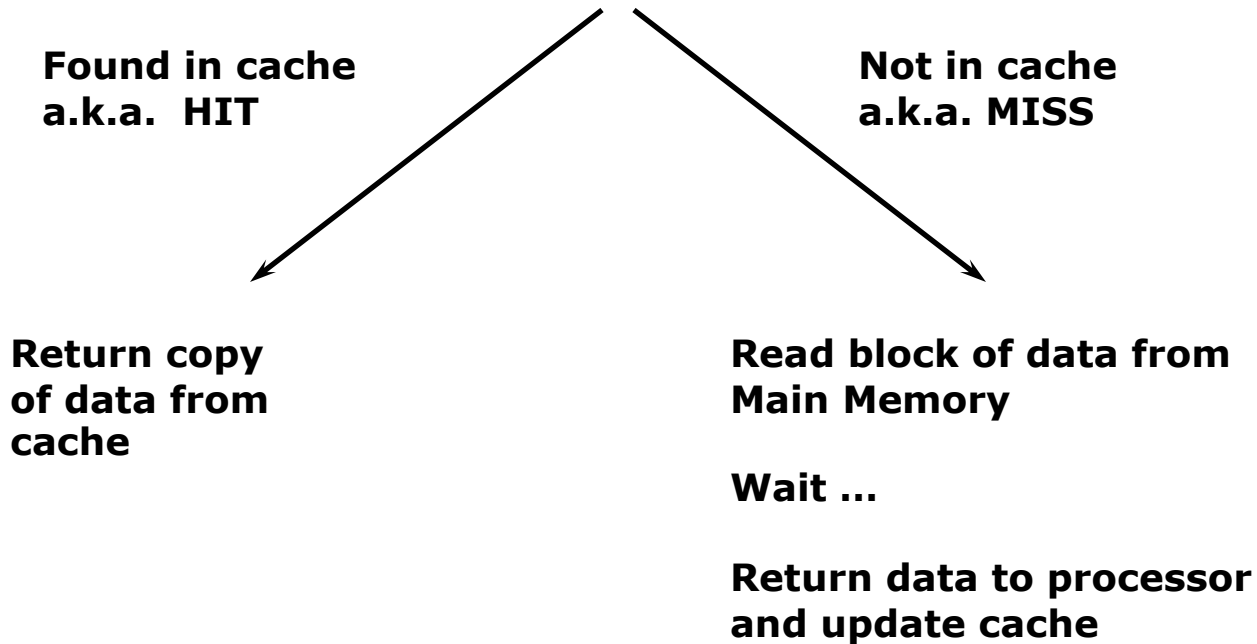| Cache Type | What is Cached? | Where is it Cached? | Latency (cycles) | Managed By |
|---|---|---|---|---|
| Registers | 4-8 bytes words | CPU core | 0 | Compiler |
| TLB | Address translations | On-Chip TLB | 0 | Hardware MMU |
| L1 cache | 64-byte blocks | On-Chip L1 | 4 | Hardware |
| L2 cache | 64-byte blocks | On-Chip L2 | 10 | Hardware |
| Virtual Memory | 4-KB pages | Main memory | 100 | Hardware + OS |
| Buffer cache | Parts of files | Main memory | 100 | OS |
| Disk cache | Disk sectors | Disk controller | 100,000 | Disk firmware |
| Network buffer cache | Parts of files | Local disk | 10,000,000 | NFS client |
| Browser cache | Web pages | Local disk | 10,000,000 | Web browser |
| Web cache | Web pages | Remote server disks | 1,000,000,000 | Web proxy server |

# General Cache Concept

**Cache**

| 4 | 9 | 10 | 3 |
|---|---|----|---|

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

| 10 |
|----|

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"

# Cache Memories

- **Cache memories** are small, fast SRAM-based memories managed automatically in hardware
  - Hold frequently accessed blocks of main memory
- **CPU looks first for data in cache, then in main memory**
- **Typical system structure:**

# Cache Algorithm (Read)

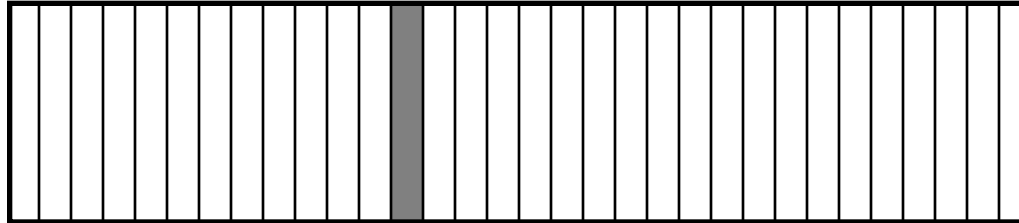Look at Processor Address, search cache tags to find match.  Then either

**Found in cache
a.k.a.  HIT**

**Not in cache
a.k.a. MISS**

**Return copy
of data from
cache**

**Read block of data from
Main Memory**

**Wait ...**

**Return data to processor
and update cache**

**Q: Which line do we replace?**

# Placement Policy

**Block Number**

0 1 2 ........................... 12 ..................................................... 31

**Memory**

**Set Number**

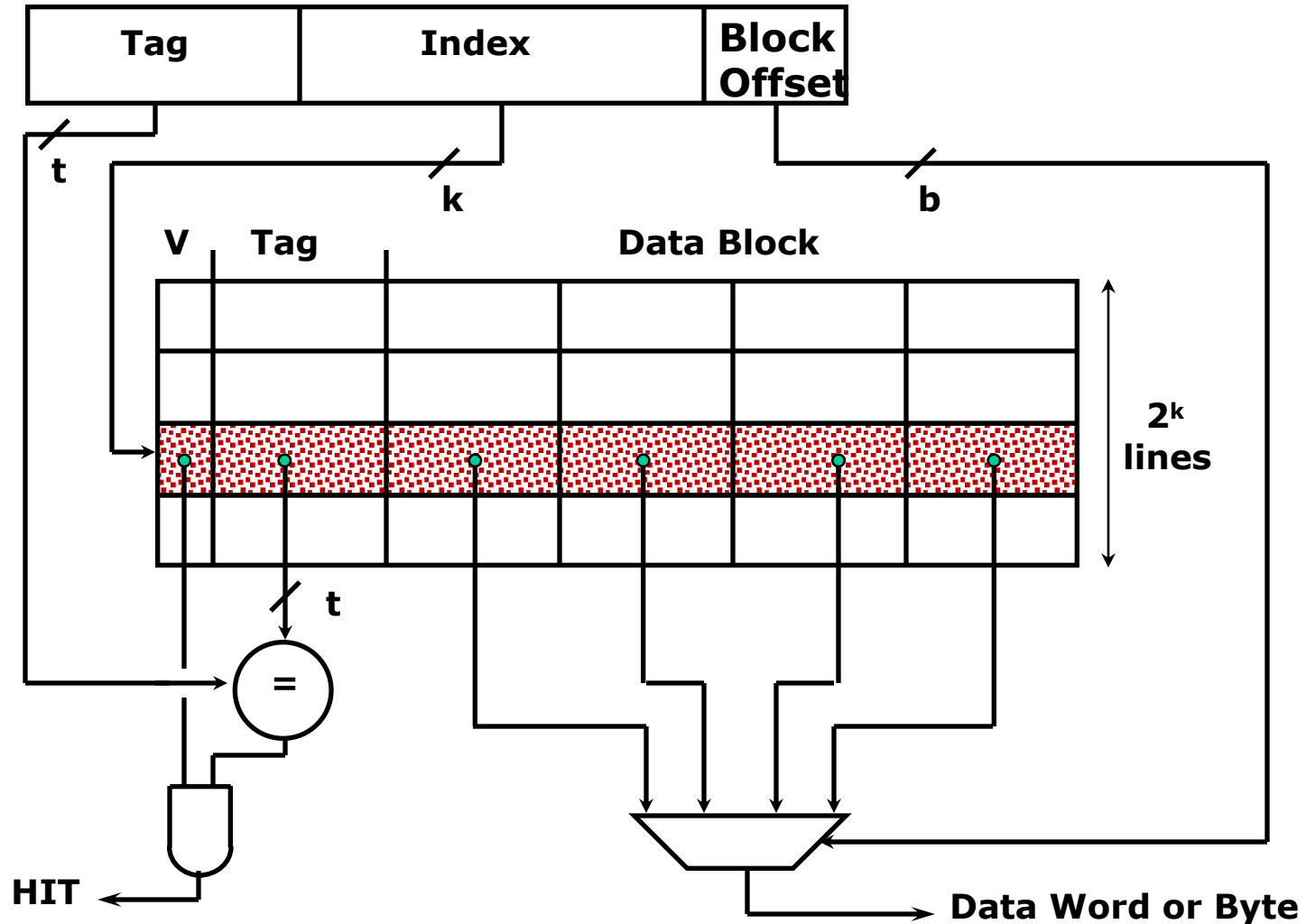0    1    2    3       0 1 2 3 4 5 6 7

**Cache**

**Fully Associative**
**anywhere**

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

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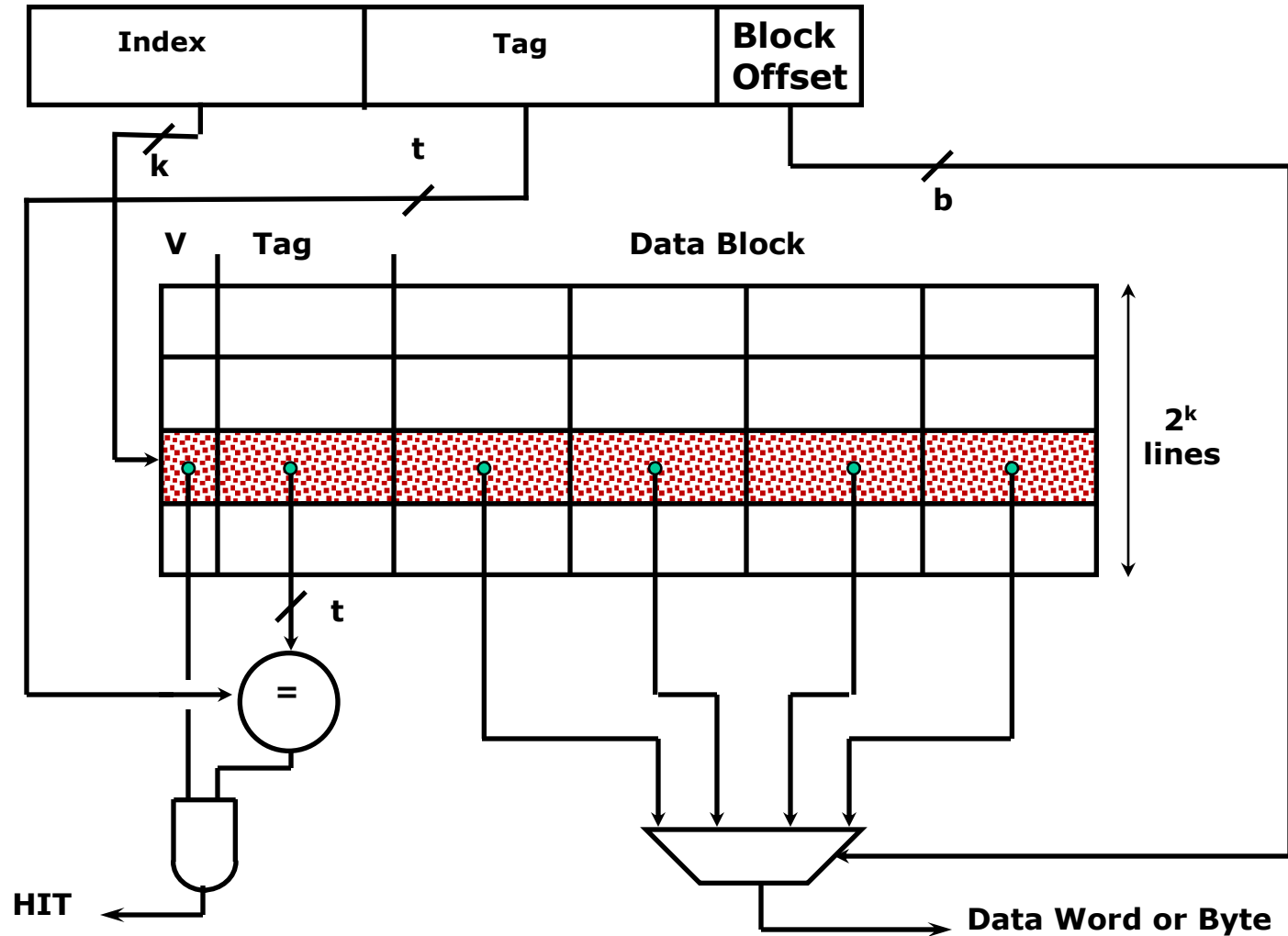**only into block 4**
*(12 mod 8)*

block 12 can be placed

# Direct-Mapped Cache

| Tag | Index | Block Offset |
|-----|-------|--------------|

t

k

b

V | Tag | Data Block
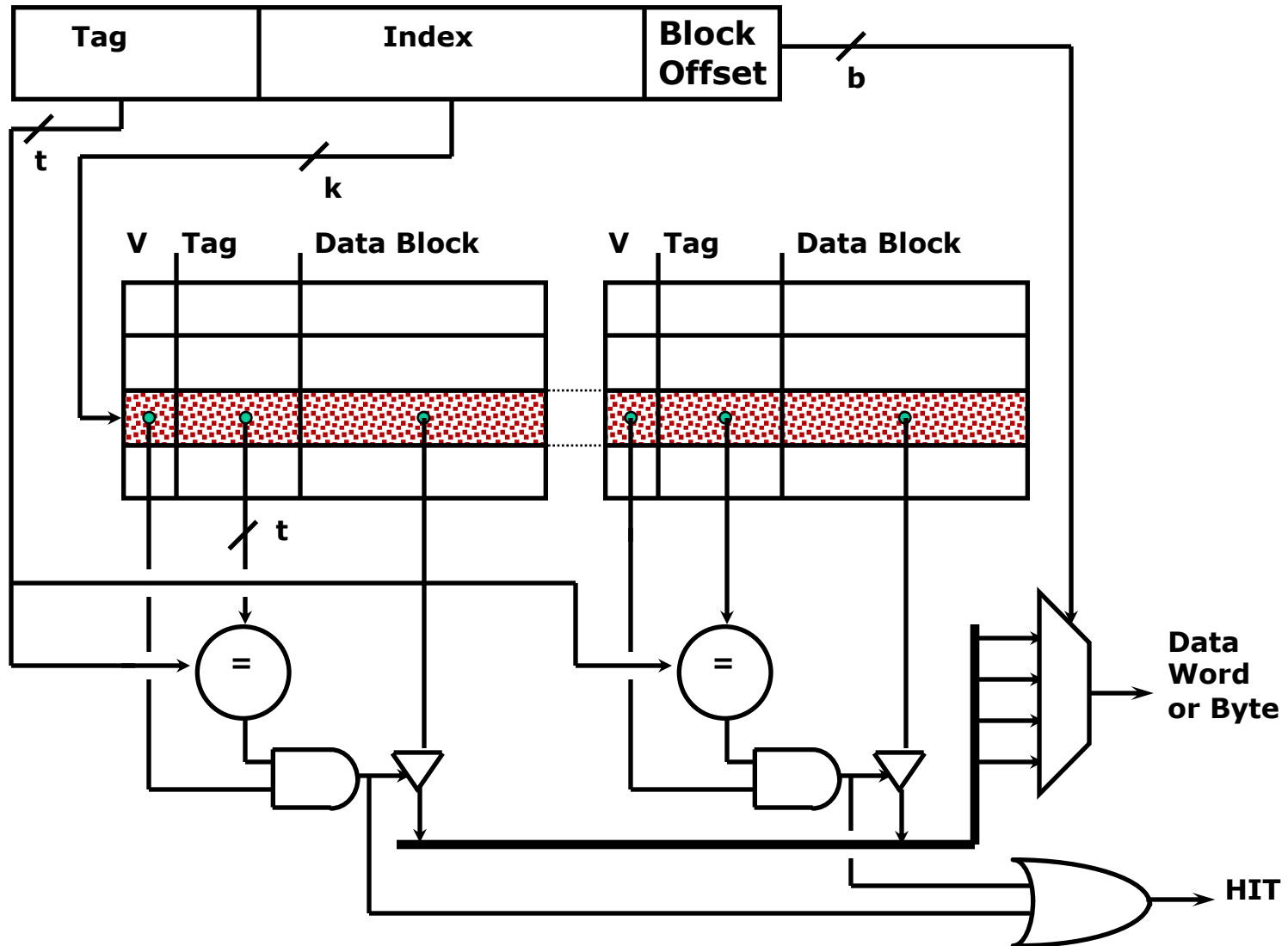
$2^k$ lines

t

=

HIT

Data Word or Byte

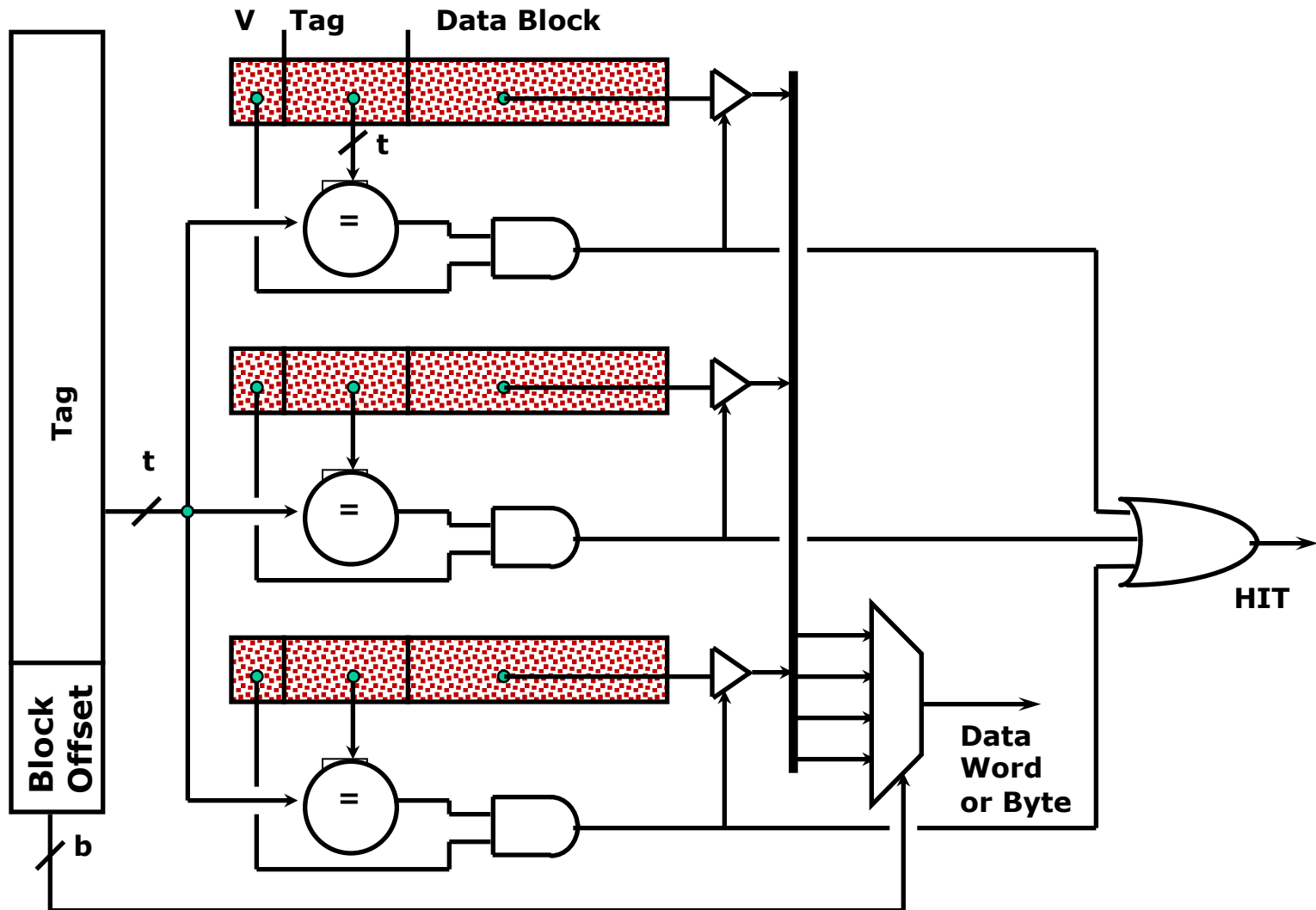# Direct Map Address Selection

*higher-order vs. lower-order address bits*

# 2-Way Set-Associative Cache

# Fully Associative Cache
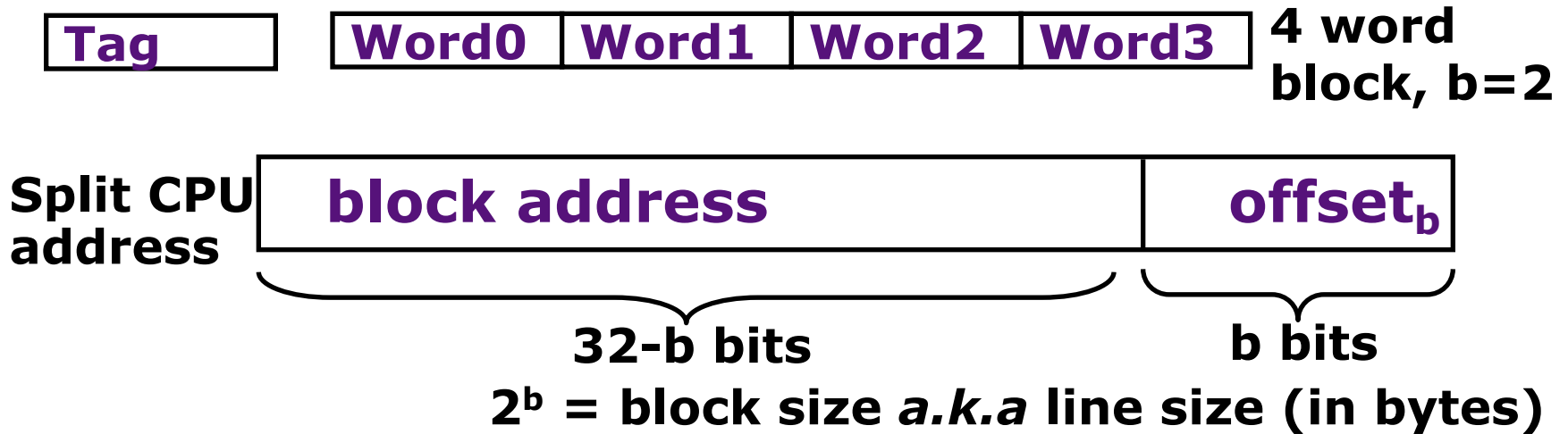
# Replacement Policy

**In an associative cache, which block from a set should be evicted when the set becomes full?**

- **Random**

- **Least-Recently Used (LRU)**
  - **LRU cache state must be updated on every access**
  - **true implementation only feasible for small sets (2-way)**
  - **pseudo-LRU binary tree often used for 4-8 way**

- **First-In, First-Out (FIFO) a.k.a. Round-Robin**
  - **used in highly associative caches**

- **Not-Most-Recently Used (NMRU)**
  - **FIFO with exception for most-recently used block or blocks**

- **Most Recently Used  (?)**     *Replacement only happens on misses*

*This is a second-order effect.  Why?*

# Block Size and Spatial Locality

**Block is unit of transfer between the cache and memory**

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

**Split CPU address**

| block address | offset$_b$ |
|---|---|

$$32\text{-}b \text{ bits} \qquad b \text{ bits}$$
$$2^b = \text{block size } a.k.a \text{ line size (in bytes)}$$

**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.*

# What about writes?

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

- **What to do on a write-hit?**
  - Write-through (write immediately to memory)
  - Write-back (defer write to memory until replacement of line)
    - Need a dirty bit (line different from memory or not)
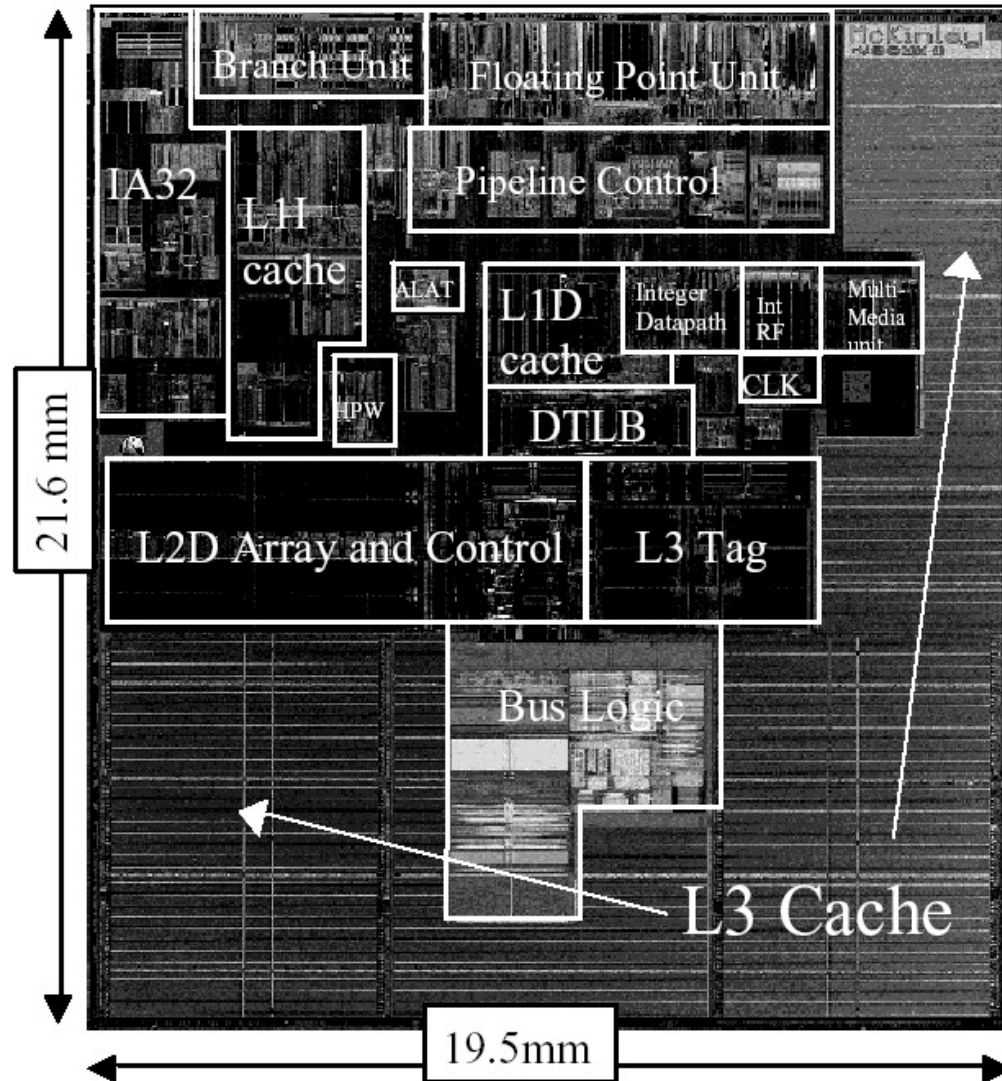
- **What to do on a write-miss?**
  - Write-allocate (load into cache, update line in cache)
    - Good if more writes to the location follow
  - No-write-allocate (writes straight to memory, does not load into cache)

- **Typical**
  - Write-through + No-write-allocate
  - **Write-back + Write-allocate**

# Itanium-2 On-Chip Caches

## (Intel/HP, 2002)



**Level 1: 16KB, 4-way s.a., 64B line, quad-port (2 load+2 store), single cycle latency**

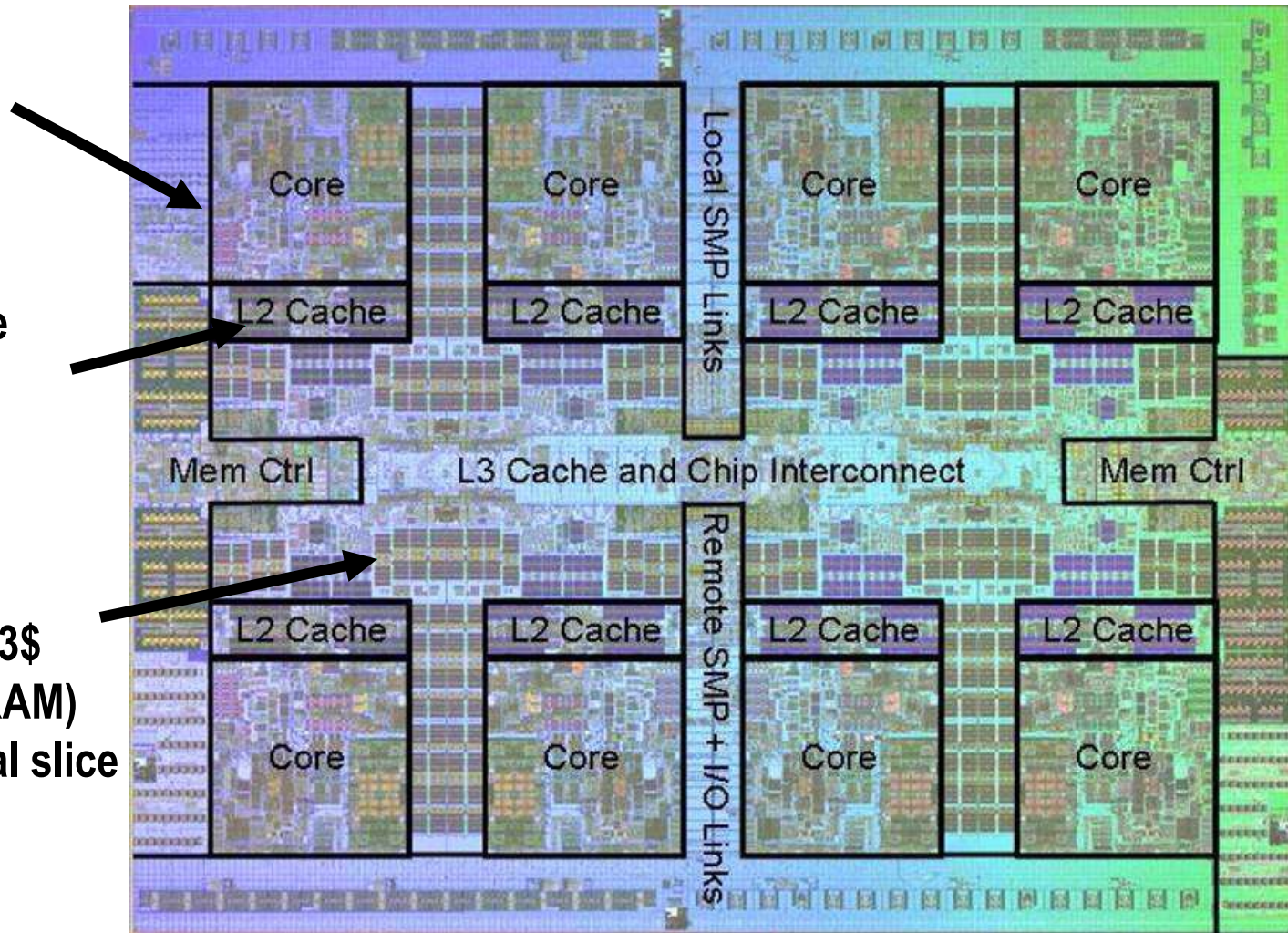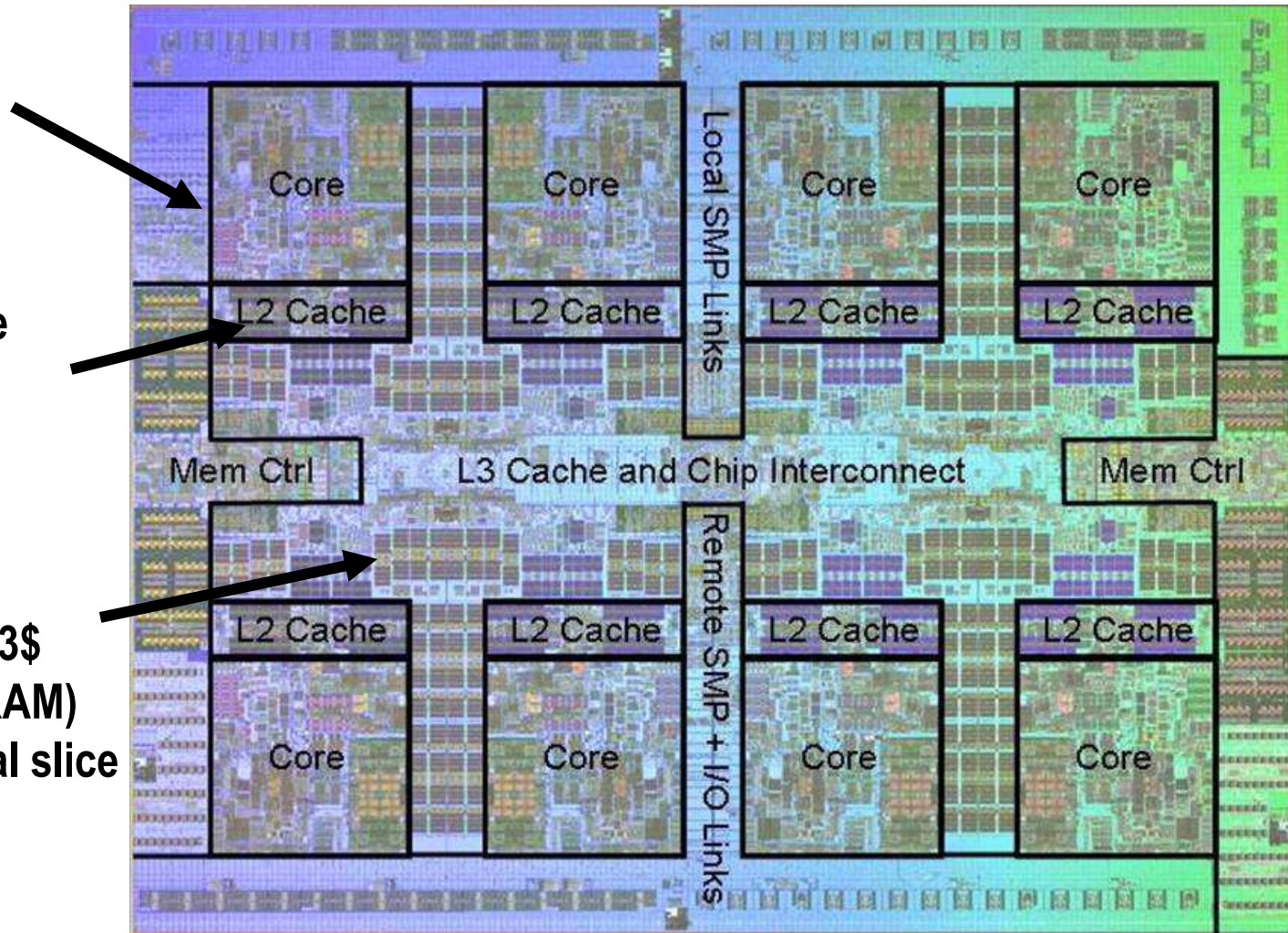**Level 2: 256KB, 4-way s.a, 128B line, quad-port (4 load or 4 store), five cycle latency**

**Level 3: 3MB, 12-way s.a., 128B line, single 32B port, twelve cycle latency**

# Power 7 On-Chip Caches [IBM 2009]

32KB L1 I$/core
32KB L1 D$/core
3-cycle latency

256KB Unified L2$/core
8-cycle latency

32MB Unified Shared L3$
Embedded DRAM (eDRAM)
25-cycle latency to local slice



Core

Core

Core

Core

L2 Cache

L2 Cache

L2 Cache

L2 Cache

Local SMP Links

Mem Ctrl

L3 Cache and Chip Interconnect

Mem Ctrl

Remote SMP + I/O Links

L2 Cache
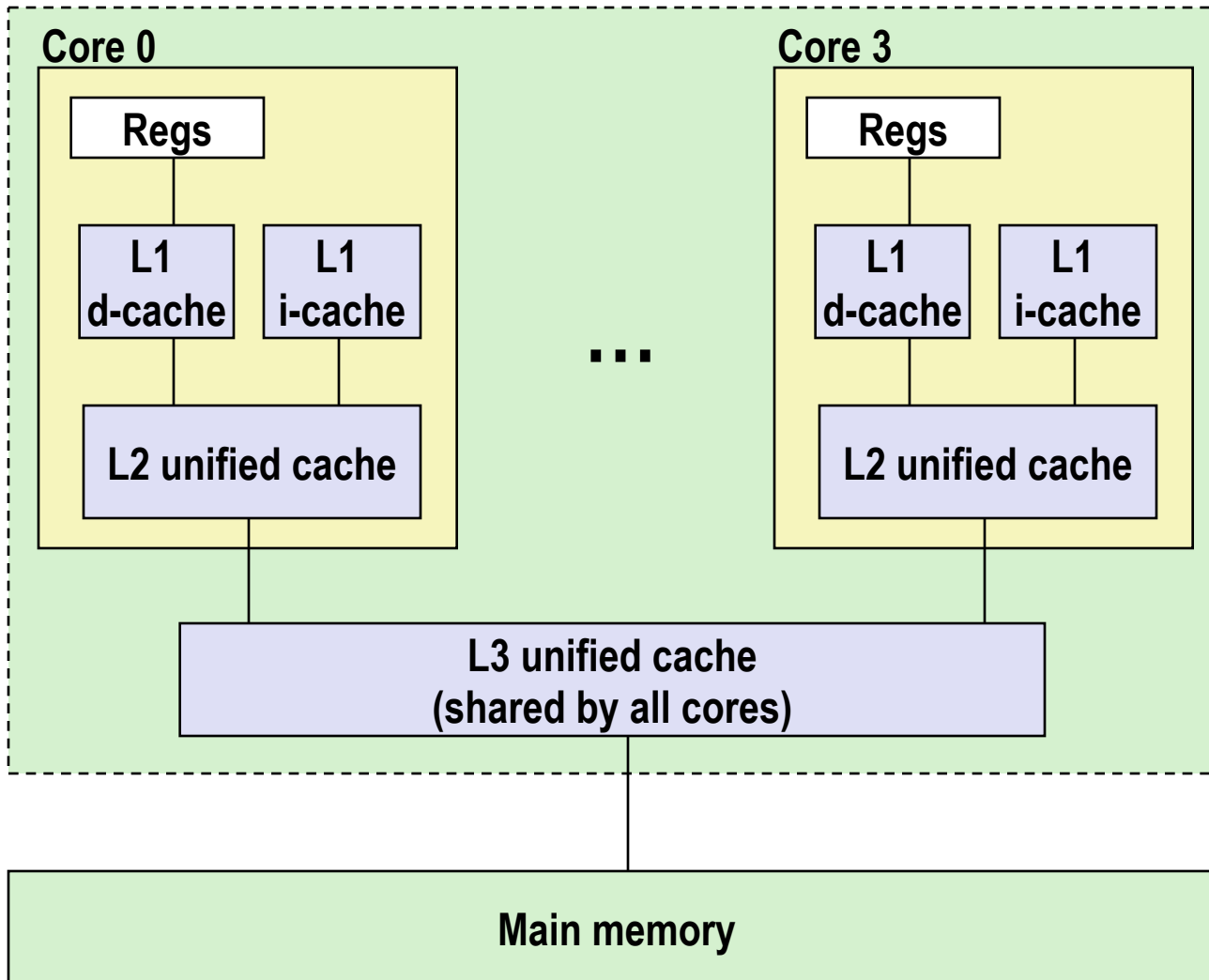
L2 Cache

L2 Cache

L2 Cache

Core

Core

Core

Core

# IBM z196 Mainframe Caches 2010

- **96 cores (4 cores/chip, 24 chips/system)**
  - Out-of-order, 3-way superscalar @ 5.2GHz
- **L1: 64KB I-$/core + 128KB D-$/core**
- **L2: 1.5MB private/core (144MB total)**
- **L3: 24MB shared/chip (eDRAM) (576MB total)**
- **L4: 768MB shared/system (eDRAM)**

# Intel Core i7 Cache Hierarchy



**Processor package**

Core 0
- Regs
- L1 d-cache
- L1 i-cache
- L2 unified cache

. . .

Core 3
- Regs
- L1 d-cache
- L1 i-cache
- L2 unified cache

L3 unified cache
(shared by all cores)

Main memory

**L1 i-cache and d-cache:**
32 KB, 8-way,
Access: 4 cycles

**L2 unified cache:**
256 KB, 8-way,
Access: 10 cycles

**L3 unified cache:**
8 MB, 16-way,
Access: 40-75 cycles

**Block size**: 64 bytes for all caches.

# 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!)

# Let's think about those numbers

- **Huge difference between a hit and a miss**
    - Could be 100x, if just L1 and main memory

- **Would you believe 99% hits is twice as good as 97%?**
    - Consider:
      cache hit time of 1 cycle
      miss penalty of 100 cycles

    - Average access time:
      97% hits:  1 cycle + 0.03 * 100 cycles = **4 cycles**
      99% hits:  1 cycle + 0.01 * 100 cycles = **2 cycles**

- **This is why "miss rate" is used instead of "hit rate"**

# Writing Cache Friendly Code

- **Make the common case go fast**
  - Focus on the inner loops of the core functions

- **Minimize the misses in the inner loops**
  - Repeated references to variables are good (temporal locality)
  - Stride-1 reference patterns are good (spatial locality)

**Key idea: Our qualitative notion of locality is quantified through our understanding of cache memories**

# Acknowledgements

- **These slides contain material developed and copyright by:**
  - Arvind (MIT)
  - Krste Asanovic (MIT/UCB)
  - Joel Emer (Intel/MIT)
  - James Hoe (CMU)
  - John Kubiatowicz (UCB)
  - David Patterson (UCB)

- **MIT material derived from course 6.823**
- **UCB material derived from course CS252**