# CS330 - Computer Organization and Assembly Language Programming

## Lecture 22

## Cache – Virtual Address

Professor : Mahmut Unan – UAB CS

# Agenda

- Locality of reference
- Caching in the memory hierarchy
- Address spaces
- VM as a tool for caching
- VM as a tool for memory management
- VM as a tool for memory protection
- Address translation

# 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.

- **Caching in the memory hierarchy**

# Example Memory Hierarchy

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

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

L0: Regs

CPU registers hold words retrieved from the L1 cache.

L1: L1 cache (SRAM)

L1 cache holds cache lines retrieved from the L2 cache.

L2: L2 cache (SRAM)

L2 cache holds cache lines retrieved from L3 cache

L3: L3 cache (SRAM)

L3 cache holds cache lines retrieved from main memory.

L4: Main memory (DRAM)

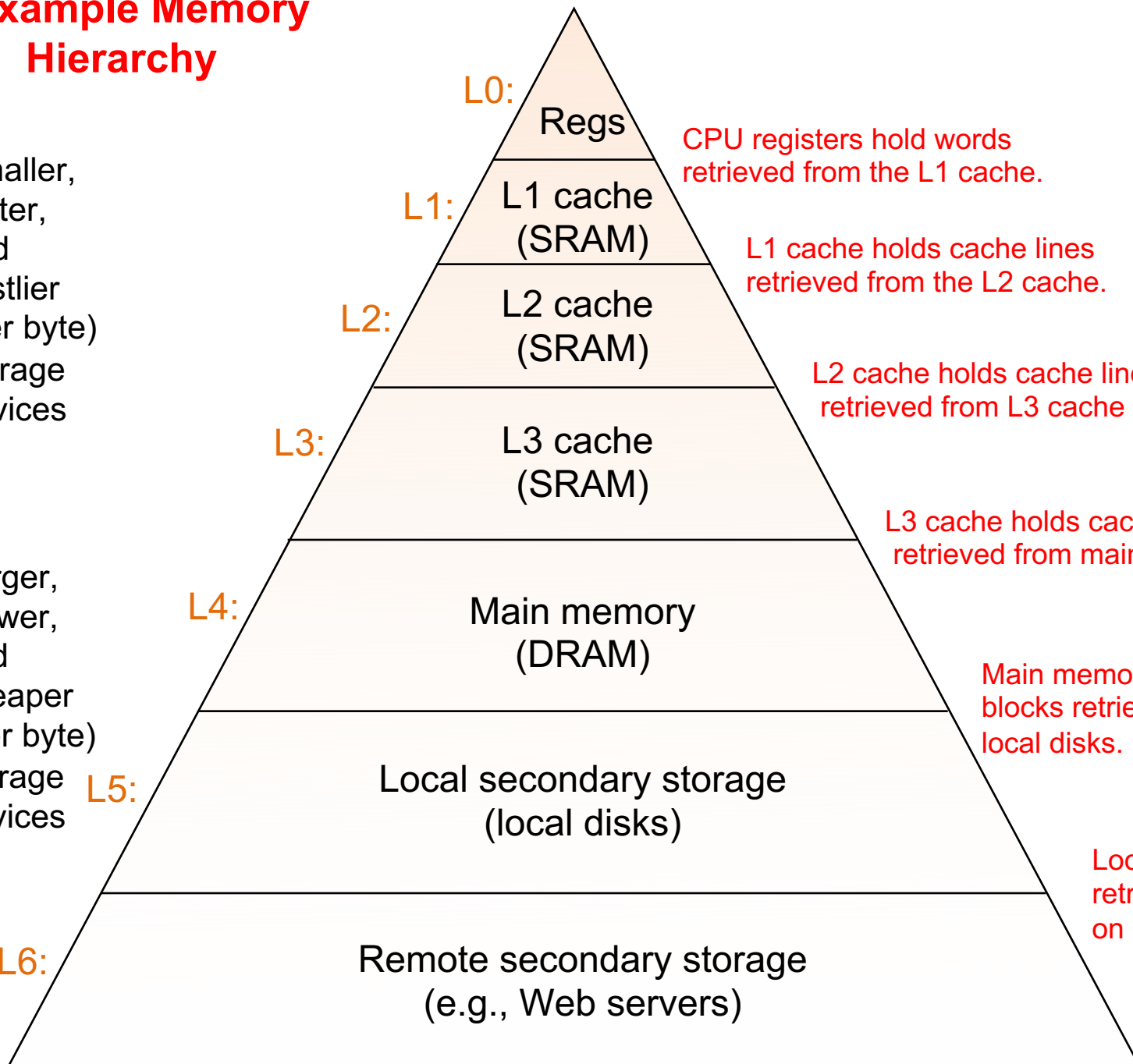Main memory holds disk blocks retrieved from local disks.

L5: Local secondary storage (local disks)

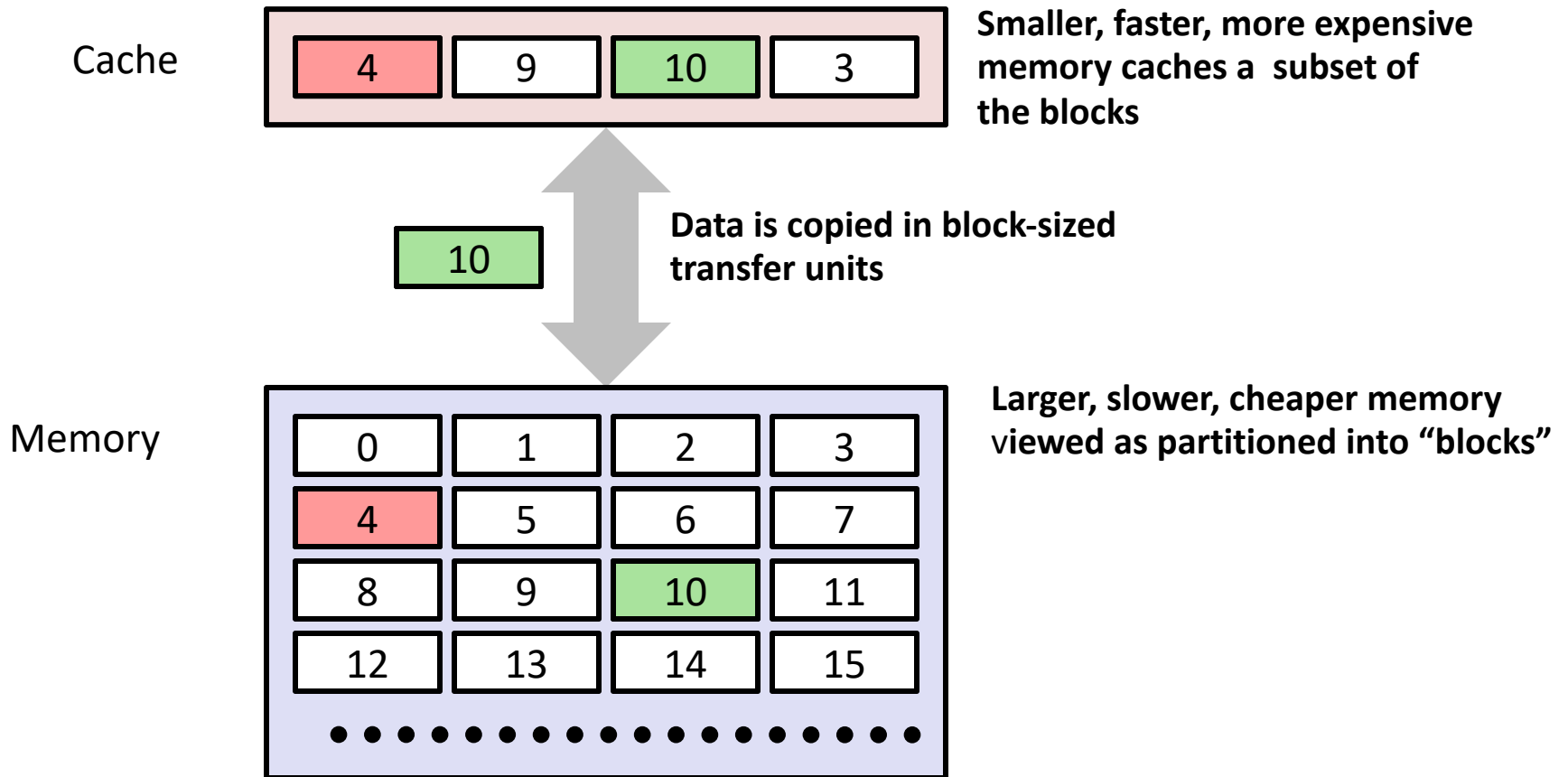Local disks hold files retrieved from disks on remote servers

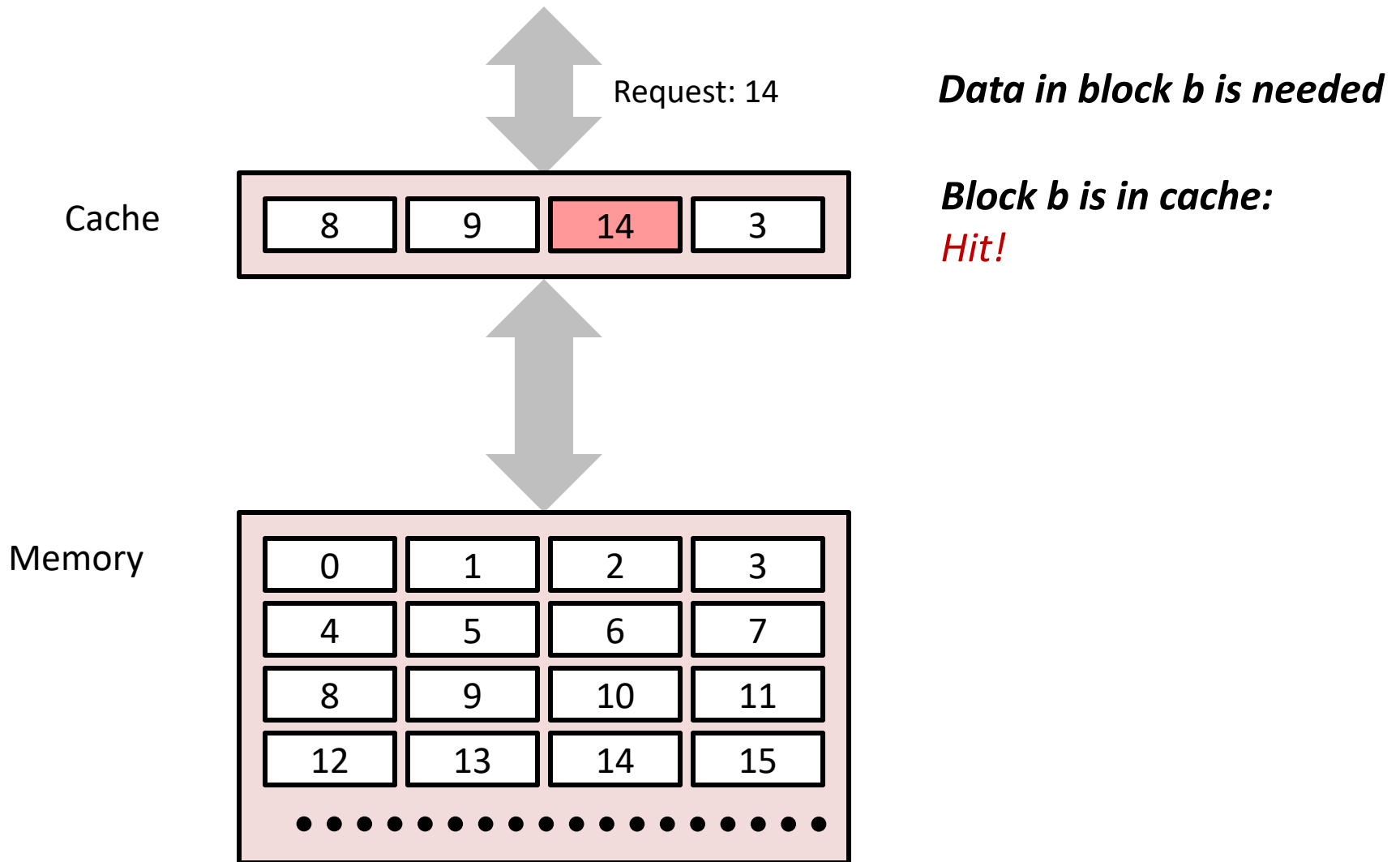L6: Remote secondary storage (e.g., Web servers)

# 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

Cache

| 8 | 12 | 14 | 3 |

12

Request: 12

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 not in cache:*
*Miss!*

*Block b is fetched from memory*

*Block b is stored in cache*
- Placement policy: determines where b goes
- Replacement policy: determines which block gets evicted (victim)
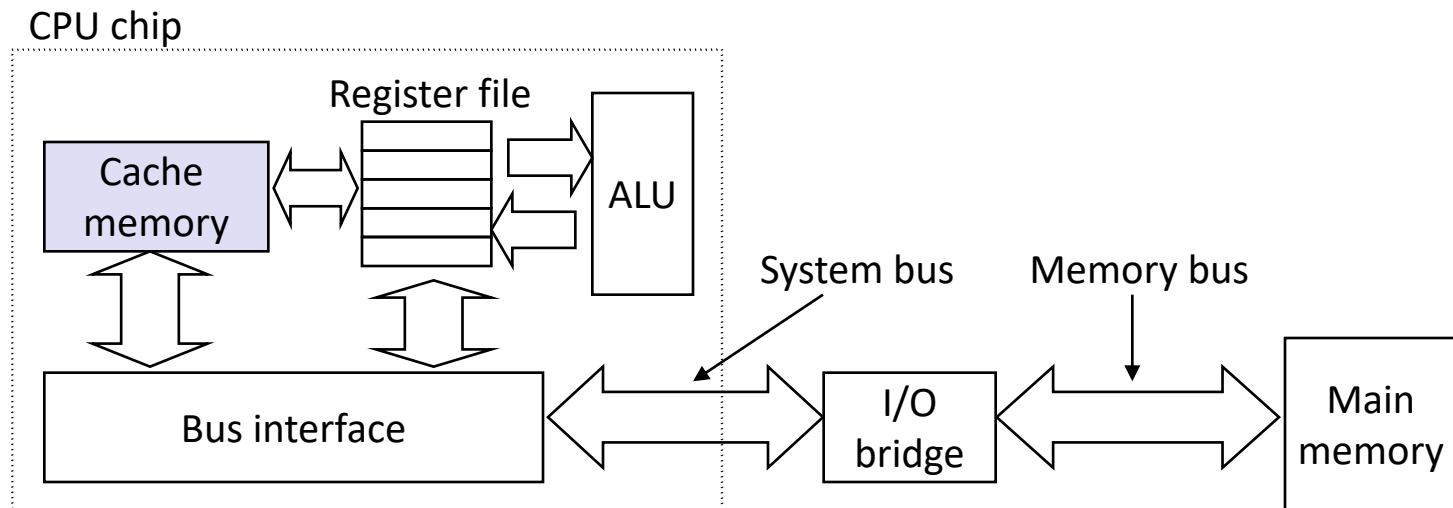
# General Caching Concepts: Types of Cache Misses

- ## 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.
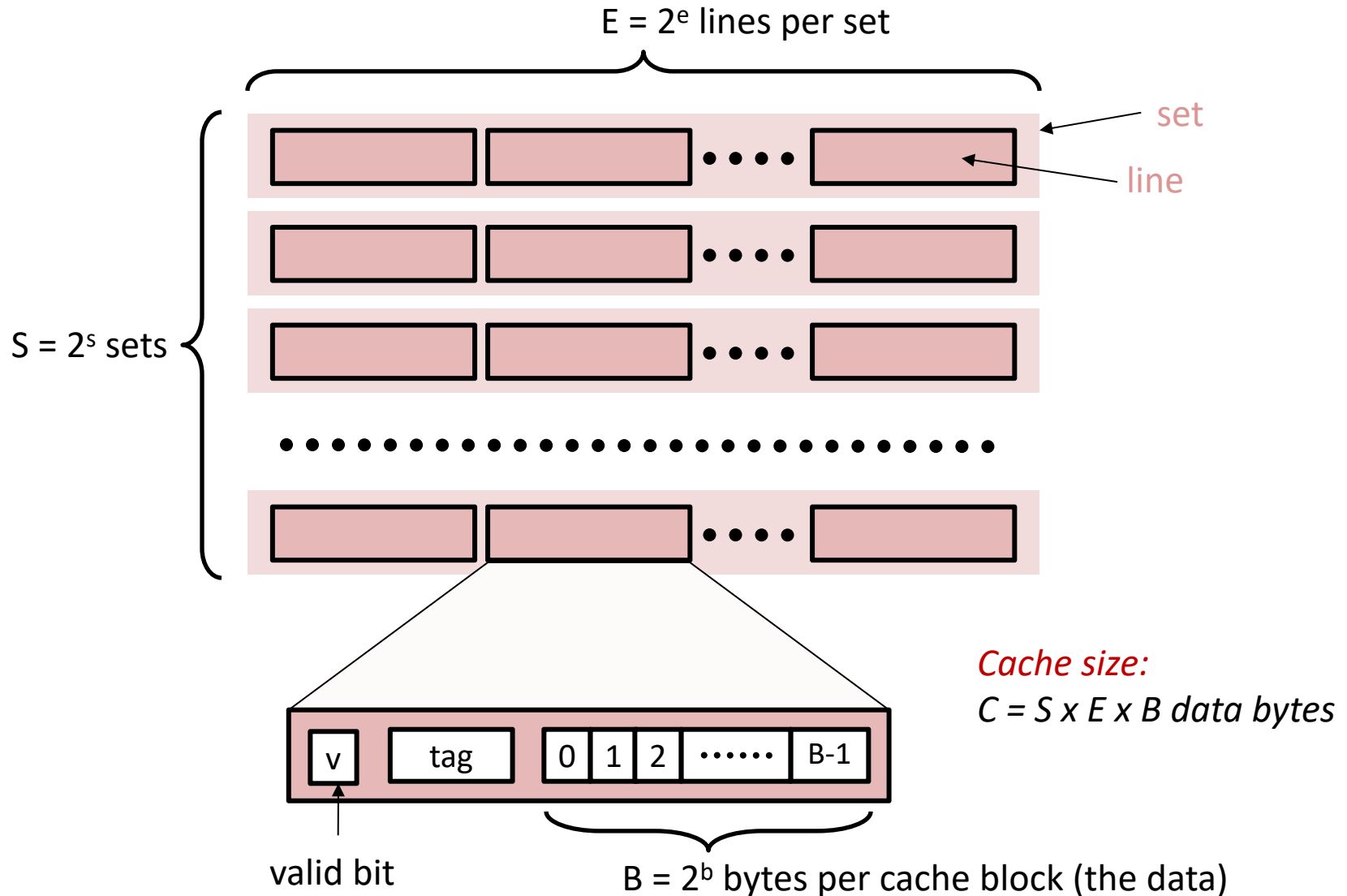
# Examples of Caching in the Mem. Hierarchy

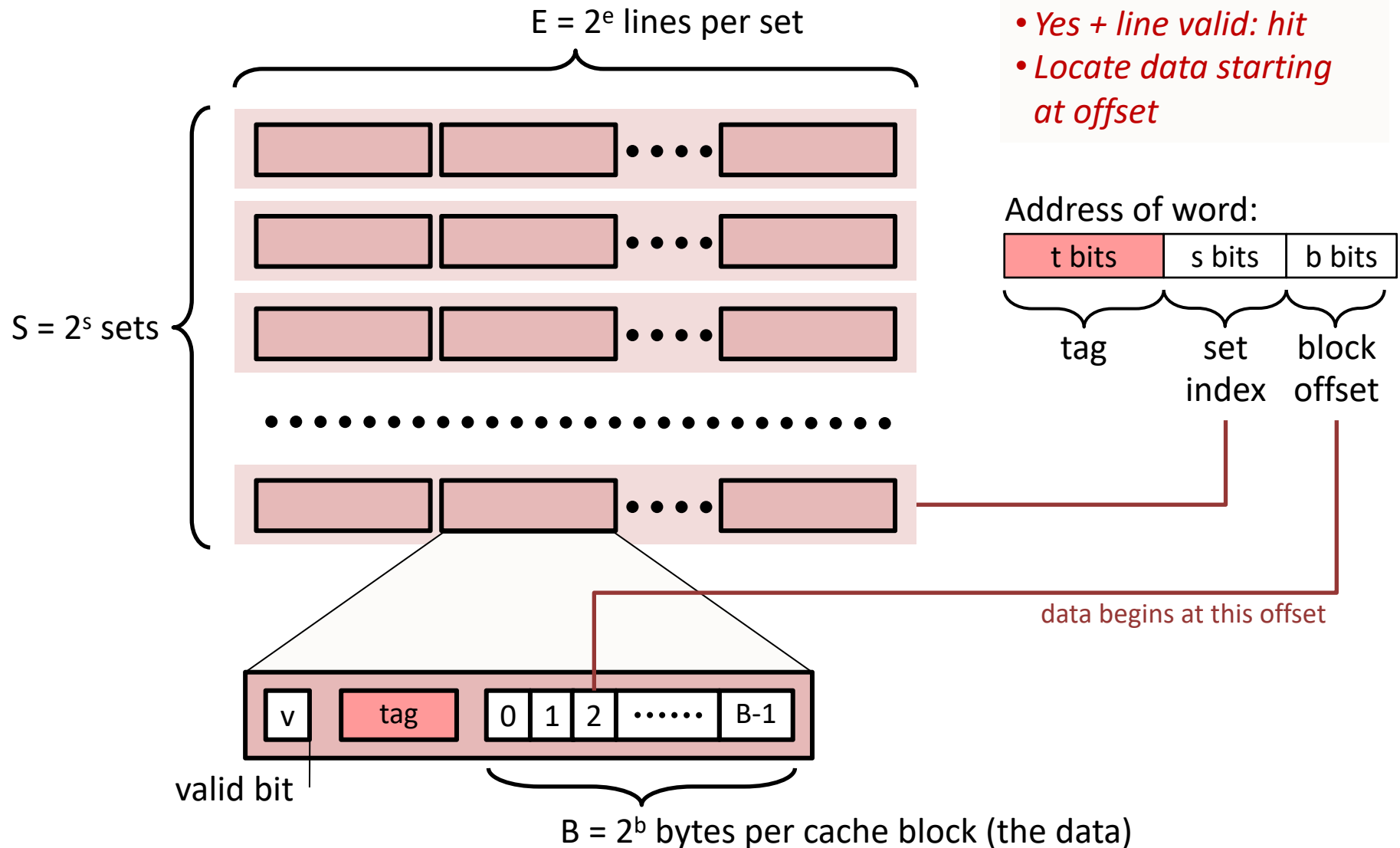| Cache Type | What is Cached? | Where is it Cached? | Latency (cycles) | Managed By |
|---|---|---|---|---|
| Registers | 4-8 bytes words | CPU core | 0 | Compiler |
| TLB | Address translations | On-Chip TLB | 0 | Hardware MMU |
| L1 cache | 64-byte blocks | On-Chip L1 | 4 | Hardware |
| L2 cache | 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 |

# 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
- Typical system structure:

CPU chip

Register file

Cache memory
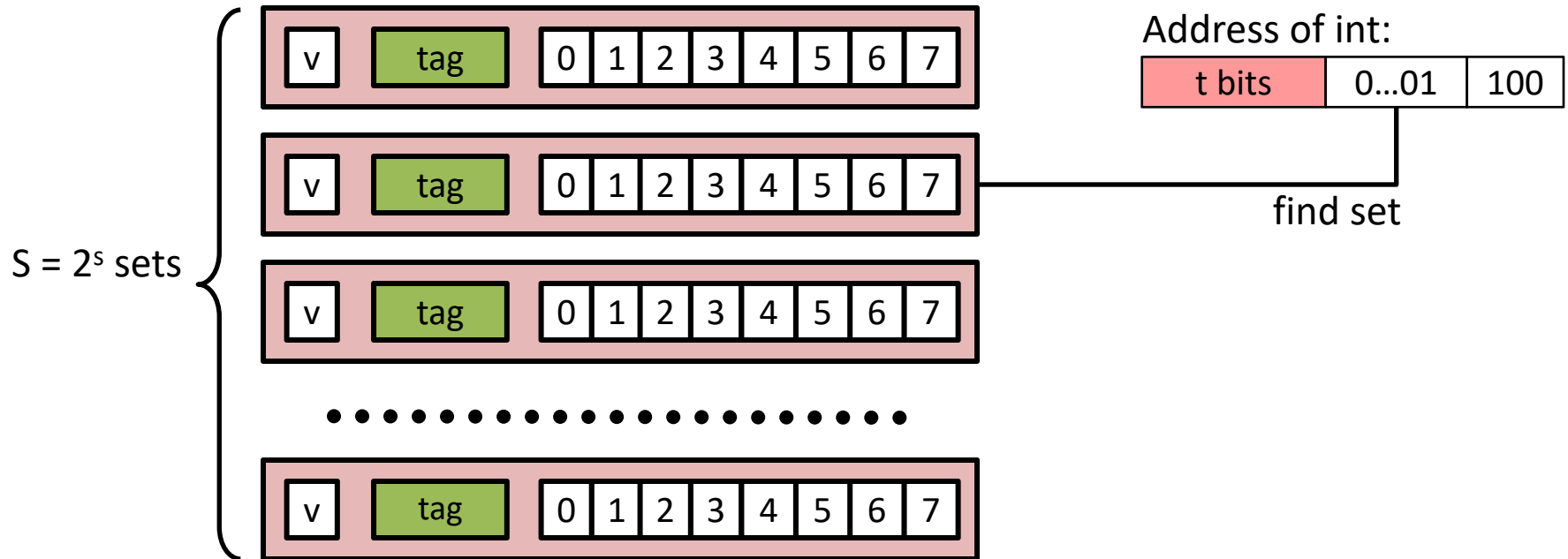
ALU

System bus

Memory bus

Bus interface

I/O bridge

Main memory

# General Cache Organization (S, E, B)

$E = 2^e$ lines per set

set

line

$S = 2^s$ sets

*Cache size:*
*C = S x E x B data bytes*

| v | tag | 0 | 1 | 2 | ...... | B-1 |

valid bit

$B = 2^b$ bytes per cache block (the data)

# Cache Read

- *Locate set*
- *Check if any line in set has matching tag*
- *Yes + line valid: hit*
- *Locate data starting at offset*

E = $2^e$ lines per set

S = $2^s$ sets

Address of word:

| t bits | s bits | b bits |
|--------|--------|--------|

tag      set index      block offset

data begins at this offset

| v | tag | 0 | 1 | 2 | ...... | B-1 |
|---|-----|---|---|---|--------|-----|

valid bit

B = $2^b$ bytes per cache block (the data)

# Example: Direct Mapped Cache (E = 1)

Direct mapped: One line per set
Assume: cache block size 8 bytes



S = $2^s$ sets

v | tag | 0 1 2 3 4 5 6 7
v | tag | 0 1 2 3 4 5 6 7
v | tag | 0 1 2 3 4 5 6 7
v | tag | 0 1 2 3 4 5 6 7

Address of int:

| t bits | 0...01 | 100 |

find set

# Example: Direct Mapped Cache (E = 1)

Direct mapped: One line per set
Assume: cache block size 8 bytes

valid?   +   match: assume yes = hit

Address of int:

| t bits | 0…01 | 100 |

| v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

block offset

# Example: Direct Mapped Cache (E = 1)

Direct mapped: One line per set
Assume: cache block size 8 bytes

valid?   +   match: assume yes = hit

Address of int:

| v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

| t bits | 0...01 | 100 |

block offset

int (4 Bytes) is here

If tag doesn't match: old line is evicted and replaced

# Direct-Mapped Cache Simulation

| t=1 | s=2 | b=1 |
|-----|-----|-----|
| x | xx | x |

M=16 bytes (4-bit addresses), B=2 bytes/block,
S=4 sets, E=1 Blocks/set

Address trace (reads, one byte per read):

|   |          |      |
|---|----------|------|
| 0 | [$0000_2$], | miss |
| 1 | [$0001_2$], | hit  |
| 7 | [$0111_2$], | miss |
| 8 | [$1000_2$], | miss |
| 0 | [$0000_2$]  | miss |

|       | v | Tag | Block |
|-------|---|-----|-------|
| Set 0 | 1 | 0 | M[0-1] |
| Set 1 |   |   |        |
| Set 2 |   |   |        |
| Set 3 | 1 | 0 | M[6-7] |

# E-way Set Associative Cache (Here: E = 2)

E = 2: Two lines per set
Assume: cache block size 8 bytes

Address of short int:

| t bits | 0…01 | 100 |
|--------|------|-----|

find set

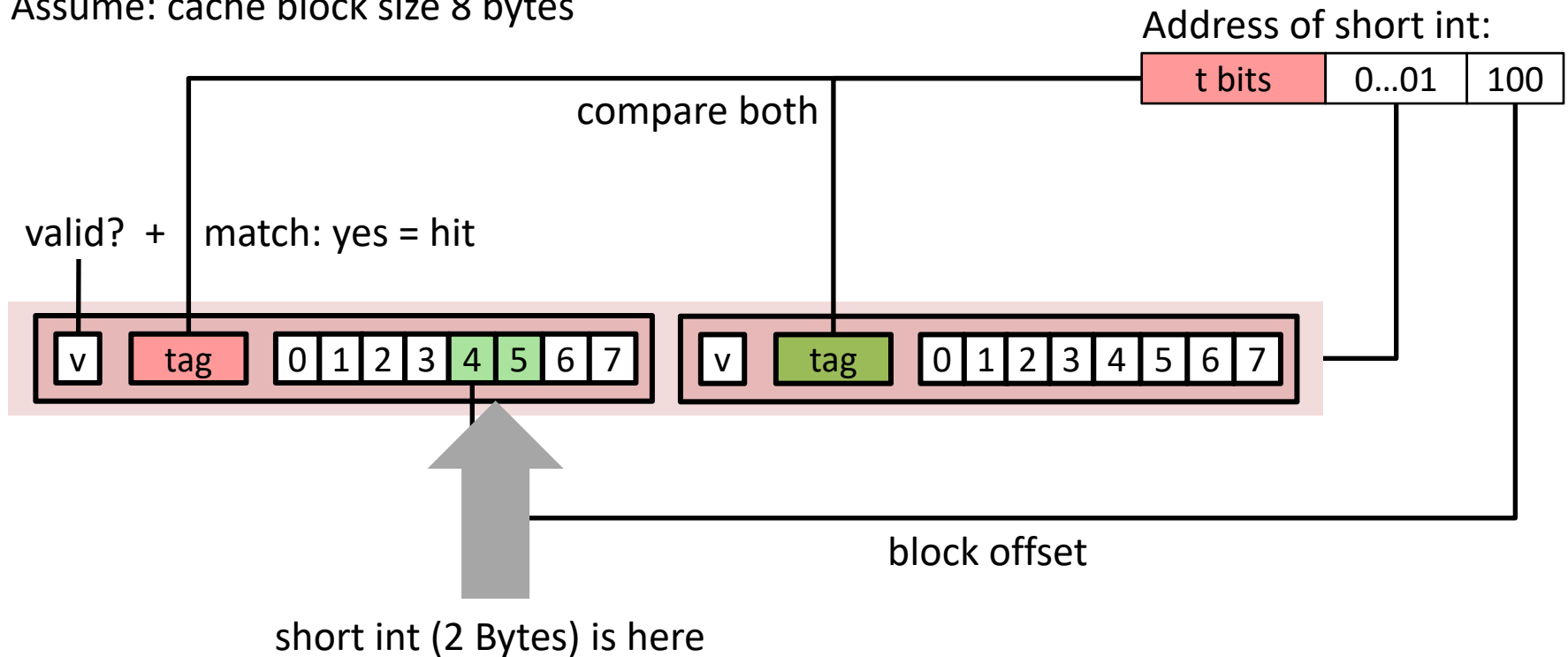# E-way Set Associative Cache (Here: E = 2)

E = 2: Two lines per set
Assume: cache block size 8 bytes

Address of short int:

| t bits | 0...01 | 100 |

compare both

valid? + match: yes = hit

| v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |     | v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

block offset

# E-way Set Associative Cache (Here: E = 2)

E = 2: Two lines per set
Assume: cache block size 8 bytes

Address of short int:

| t bits | 0...01 | 100 |

compare both

valid? + | match: yes = hit

| v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |   | v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

short int (2 Bytes) is here

block offset

No match:
- One line in set is selected for eviction and replacement
- Replacement policies: random, least recently used (LRU), ...

# 2-Way Set Associative Cache Simulation

t=2    s=1    b=1

| xx | x | x |
|----|---|---|

M=16 byte addresses, B=2 bytes/block,
S=2 sets, E=2 blocks/set

Address trace (reads, one byte per read):

| 0 | $[0000_2]$, | miss |
|---|-----------|------|
| 1 | $[0001_2]$, | hit |
| 7 | $[0111_2]$, | miss |
| 8 | $[1000_2]$, | miss |
| 0 | $[0000_2]$ | hit |

|       | v | Tag | Block |
|-------|---|-----|-------|
| Set 0 | 1 | 00 | M[0-1] |
|       | 1 | 10 | M[8-9] |
| Set 1 | 1 | 01 | M[6-7] |
|       | 0 |    |       |

# 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**

# Intel Core i7 Cache Hierarchy

Processor package

| | | |
|---|---|---|
| **Core 0** | ... | **Core 3** |

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 (<span style="color:red">temporal locality</span>)
  - Stride-1 reference patterns are good (<span style="color:red">spatial locality</span>)

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

# The Memory Mountain

- Read throughput (read bandwidth)
  - Number of bytes read from memory per second (MB/s)


- Memory mountain: Measured read throughput as a function of spatial and temporal locality.
  - Compact way to characterize memory system performance.

# Memory Mountain Test Function

```c
long data[MAXELEMS];  /* Global array to traverse */

/* test - Iterate over first "elems" elements of
 *        array "data" with stride of "stride", using
 *        using 4x4 loop unrolling.
 */
int test(int elems, int stride) {
    long i, sx2=stride*2, sx3=stride*3, sx4=stride*4;
    long acc0 = 0, acc1 = 0, acc2 = 0, acc3 = 0;
    long length = elems, limit = length - sx4;

    /* Combine 4 elements at a time */
    for (i = 0; i < limit; i += sx4) {
        acc0 = acc0 + data[i];
        acc1 = acc1 + data[i+stride];
        acc2 = acc2 + data[i+sx2];
        acc3 = acc3 + data[i+sx3];
    }

    /* Finish any remaining elements */
    for (; i < length; i++) {
        acc0 = acc0 + data[i];
    }
    return ((acc0 + acc1) + (acc2 + acc3));
}
```
*mountain/mountain.c*

Call `test()` with many combinations of `elems` and `stride`.

For each elems and stride:

1. Call test() once to warm up the caches.

2. Call test() again and measure the read throughput(MB/s)

# The Memory Mountain

Core i7 Haswell
2.1 GHz
32 KB L1 d-cache
256 KB L2 cache
8 MB L3 cache
64 B block size

# Rearranging loops to improve spatial locality

# Matrix Multiplication Example

- Description:
  - Multiply N x N matrices
  - Matrix elements are doubles (8 bytes)
  - O(N³) total operations
  - N reads per source element
  - N values summed per destination
    - but may be able to hold in register

Variable `sum` held in register

```
/* ijk */
for (i=0; i<n; i++)  {
  for (j=0; j<n; j++) {
    sum = 0.0;
    for (k=0; k<n; k++)
      sum += a[i][k] * b[k][j];
    c[i][j] = sum;
  }
}
                         matmult/mm.c
```
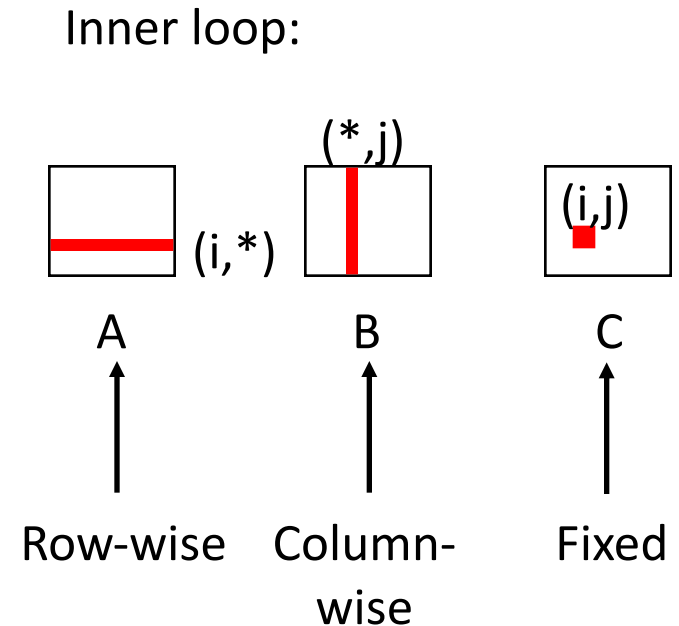
# Miss Rate Analysis for Matrix Multiply

- Assume:
  - Block size = 32B (big enough for four doubles)
  - Matrix dimension (N) is very large
    - Approximate 1/N as 0.0
  - Cache is not even big enough to hold multiple rows
- Analysis Method:
  - Look at access pattern of inner loop

# Layout of C Arrays in Memory (review)

- C arrays allocated in row-major order
  - each row in contiguous memory locations
- Stepping through columns in one row:
  - ```
    for (i = 0; i < N; i++)
       sum += a[0][i];
    ```
  - accesses successive elements
  - if block size (B) > sizeof($a_{ij}$) bytes, exploit spatial locality
    - miss rate = sizeof($a_{ij}$) / B
- Stepping through rows in one column:
  - ```
    for (i = 0; i < n; i++)
       sum += a[i][0];
    ```
  - accesses distant elements
  - no spatial locality!
    - miss rate = 1 (i.e. 100%)

# Matrix Multiplication (ijk)

```
/* ijk */
for (i=0; i<n; i++)  {
  for (j=0; j<n; j++) {
    sum = 0.0;
    for (k=0; k<n; k++)
      sum += a[i][k] * b[k][j];
    c[i][j] = sum;
  }
}                      matmult/mm.c
```
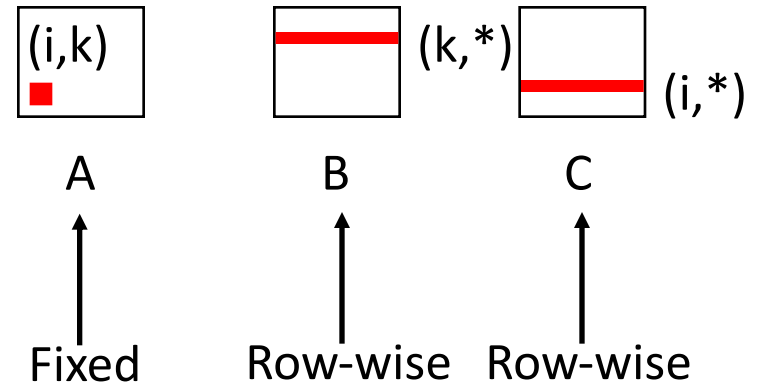
Inner loop:



| A | B | C |
|---|---|---|
| Row-wise | Column-wise | Fixed |

Misses per inner loop iteration:

| A | B | C |
|---|---|---|
| 0.25 | 1.0 | 0.0 |

# Matrix Multiplication (jik)

```
/* jik */
for (j=0; j<n; j++) {
  for (i=0; i<n; i++) {
    sum = 0.0;
    for (k=0; k<n; k++)
      sum += a[i][k] * b[k][j];
    c[i][j] = sum
  }
}                          matmult/mm.c
```

Inner loop:



| A | B | C |
|---|---|---|
| Row-wise | Column-wise | Fixed |

Misses per inner loop iteration:

| A | B | C |
|------|-----|-----|
| 0.25 | 1.0 | 0.0 |

# Matrix Multiplication (kij)

```
/* kij */
for (k=0; k<n; k++) {
  for (i=0; i<n; i++) {
    r = a[i][k];
    for (j=0; j<n; j++)
      c[i][j] += r * b[k][j];
  }
}
                        matmult/mm.c
```
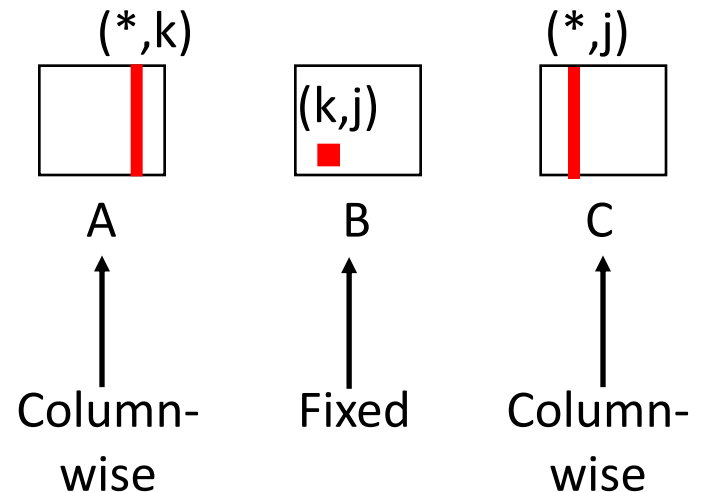
Inner loop:



|  A  |  B  |  C  |
| Fixed | Row-wise | Row-wise |

Misses per inner loop iteration:

| A | B | C |
|---|---|---|
| 0.0 | 0.25 | 0.25 |

# Matrix Multiplication (ikj)

```
/* ikj */
for (i=0; i<n; i++) {
  for (k=0; k<n; k++) {
    r = a[i][k];
    for (j=0; j<n; j++)
      c[i][j] += r * b[k][j];
  }
}
                        matmult/mm.c
```

Inner loop:



| (i,k) | (k,*) | (i,*) |
| A | B | C |
| Fixed | Row-wise | Row-wise |

Misses per inner loop iteration:

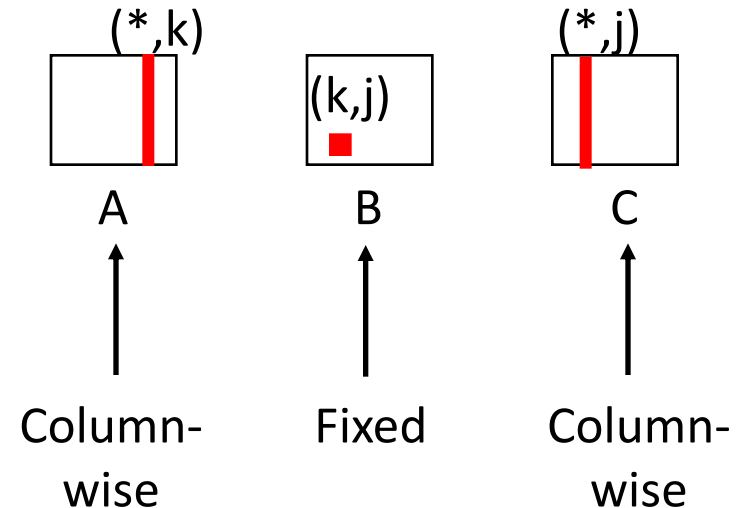| <u>A</u> | <u>B</u> | <u>C</u> |
| --- | --- | --- |
| 0.0 | 0.25 | 0.25 |

# Matrix Multiplication (jki)

```
/* jki */
for (j=0; j<n; j++) {
  for (k=0; k<n; k++) {
    r = b[k][j];
    for (i=0; i<n; i++)
      c[i][j] += a[i][k] * r;
  }
}
                        matmult/mm.c
```

Inner loop:



(*,k)     (k,j)     (*,j)

A         B         C

Column-    Fixed    Column-
wise                wise

Misses per inner loop iteration:

| A | B | C |
|---|---|---|
| 1.0 | 0.0 | 1.0 |

# Matrix Multiplication (kji)

```
/* kji */
for (k=0; k<n; k++) {
  for (j=0; j<n; j++) {
    r = b[k][j];
    for (i=0; i<n; i++)
      c[i][j] += a[i][k] * r;
  }
}
                          matmult/mm.c
```

Inner loop:



| | A | B | C |
|---|---|---|---|
| | Column-wise | Fixed | Column-wise |

Misses per inner loop iteration:

| A | B | C |
|---|---|---|
| 1.0 | 0.0 | 1.0 |

# Summary of Matrix Multiplication

```
for (i=0; i<n; i++) {
  for (j=0; j<n; j++) {
    sum = 0.0;
    for (k=0; k<n; k++)
      sum += a[i][k] * b[k][j];
    c[i][j] = sum;
  }
}
```

ijk (& jik):
* 2 loads, 0 stores
* misses/iter = 1.25

```
for (k=0; k<n; k++) {
 for (i=0; i<n; i++) {
  r = a[i][k];
  for (j=0; j<n; j++)
    c[i][j] += r * b[k][j];
 }
}
```
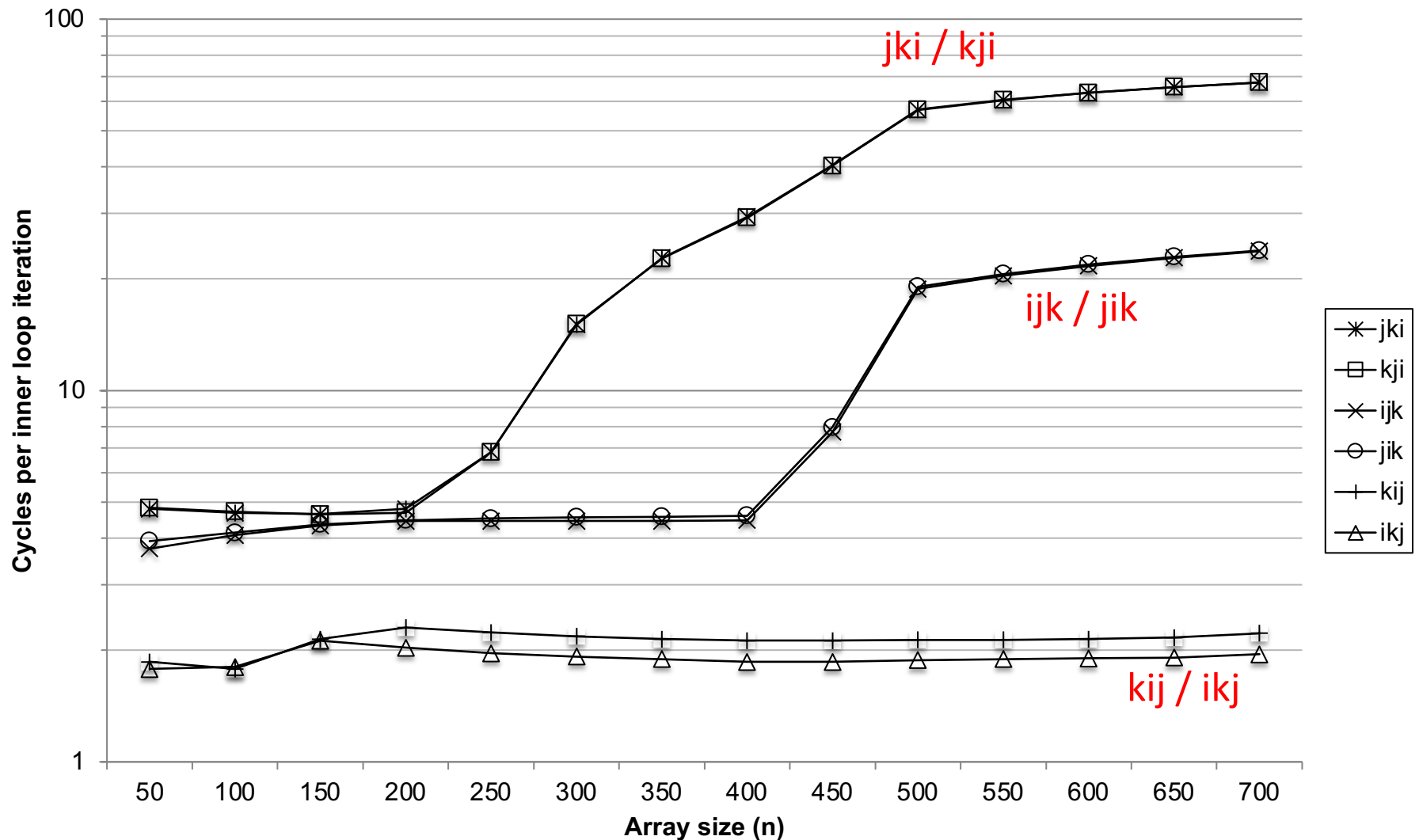
kij (& ikj):
* 2 loads, 1 store
* misses/iter = 0.5

```
for (j=0; j<n; j++) {
 for (k=0; k<n; k++) {
   r = b[k][j];
   for (i=0; i<n; i++)
     c[i][j] += a[i][k] * r;
 }
}
```

jki (& kji):
* 2 loads, 1 store
* misses/iter = 2.0
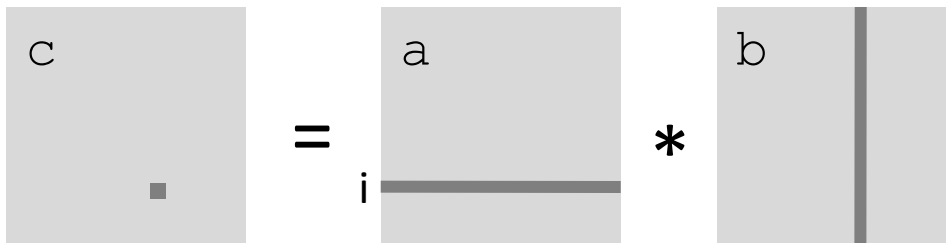
# Core i7 Matrix Multiply Performance

# Using blocking to improve temporal locality

# Example: Matrix Multiplication

```
c = (double *) calloc(sizeof(double), n*n);

/* Multiply n x n matrices a and b  */
void mmm(double *a, double *b, double *c, int n) {
    int i, j, k;
    for (i = 0; i < n; i++)
    for (j = 0; j < n; j++)
            for (k = 0; k < n; k++)
            c[i*n + j] += a[i*n + k] * b[k*n + j];
}
```
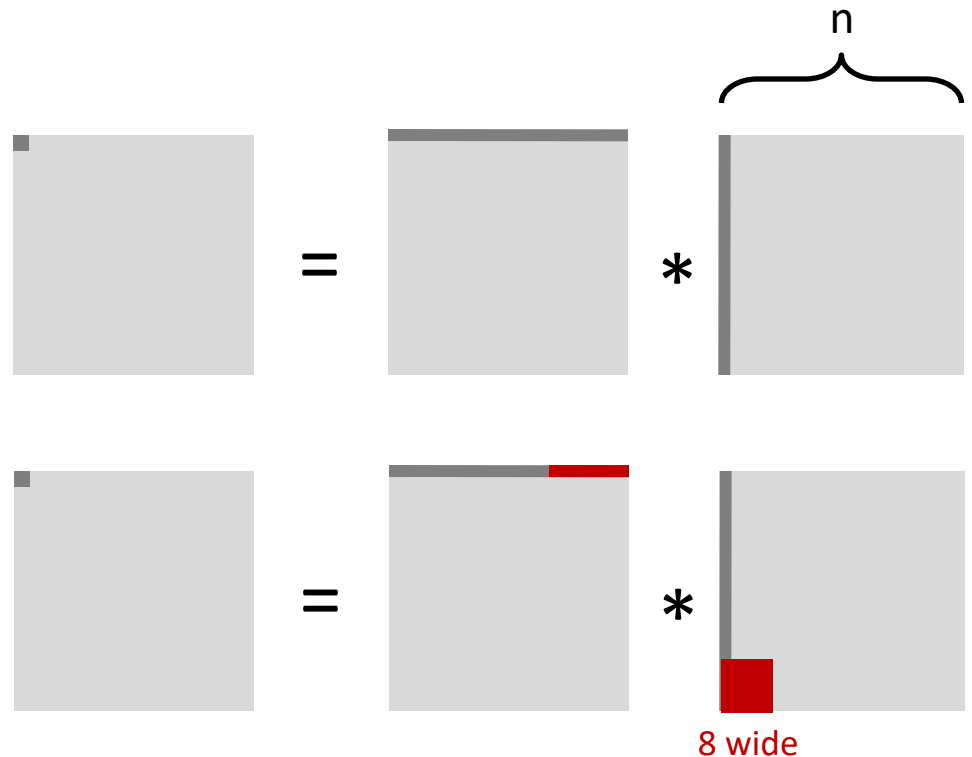
c  =  a  *  b

i

# Cache Miss Analysis

- Assume:
    - Matrix elements are doubles
    - Cache block = 8 doubles
    - Cache size C << n (much smaller than n)

- First iteration:
    - n/8 + n = 9n/8 misses
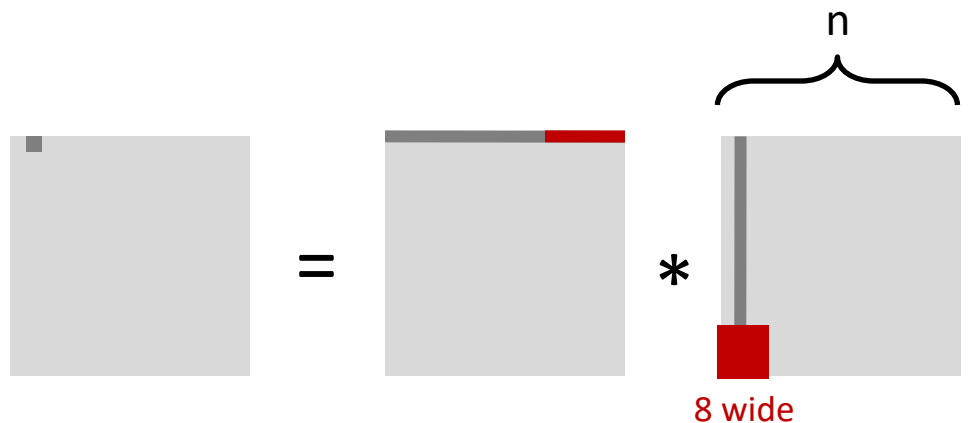
    - Afterwards in cache:
      (schematic)

# Cache Miss Analysis

- Assume:
  - Matrix elements are doubles
  - Cache block = 8 doubles
  - Cache size C << n (much smaller than n)

- Second iteration:
  - Again:
    $n/8 + n = 9n/8$ misses

- Total misses:
  - $9n/8 * n^2 = (9/8) * n^3$
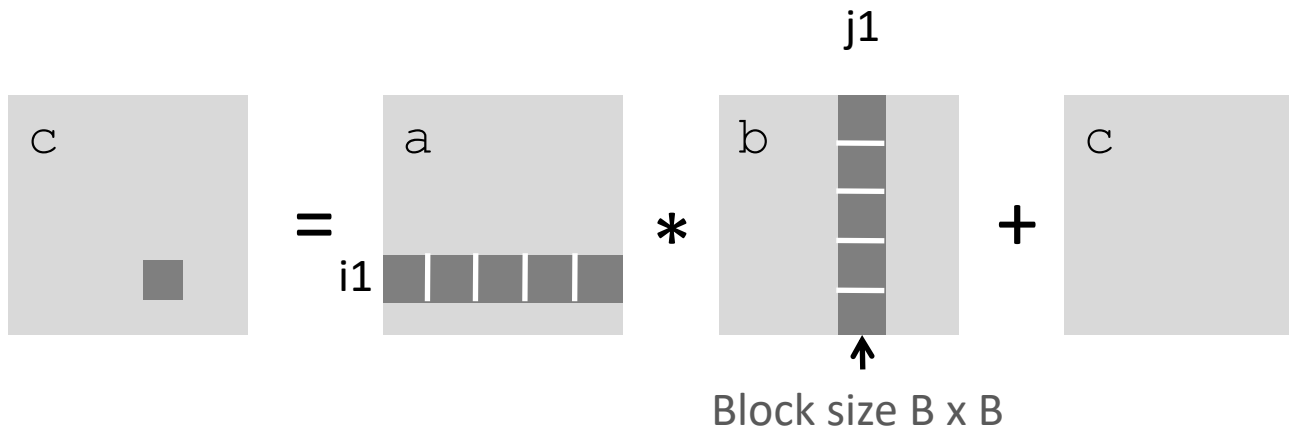


n

=

*

8 wide

# Blocked Matrix Multiplication

```c
c = (double *) calloc(sizeof(double), n*n);

/* Multiply n x n matrices a and b  */
void mmm(double *a, double *b, double *c, int n) {
    int i, j, k;
    for (i = 0; i < n; i+=B)
    for (j = 0; j < n; j+=B)
            for (k = 0; k < n; k+=B)
        /* B x B mini matrix multiplications */
                for (i1 = i; i1 < i+B; i++)
                    for (j1 = j; j1 < j+B; j++)
                      for (k1 = k; k1 < k+B; k++)
                        c[i1*n+j1] += a[i1*n + k1]*b[k1*n + j1];
}
```
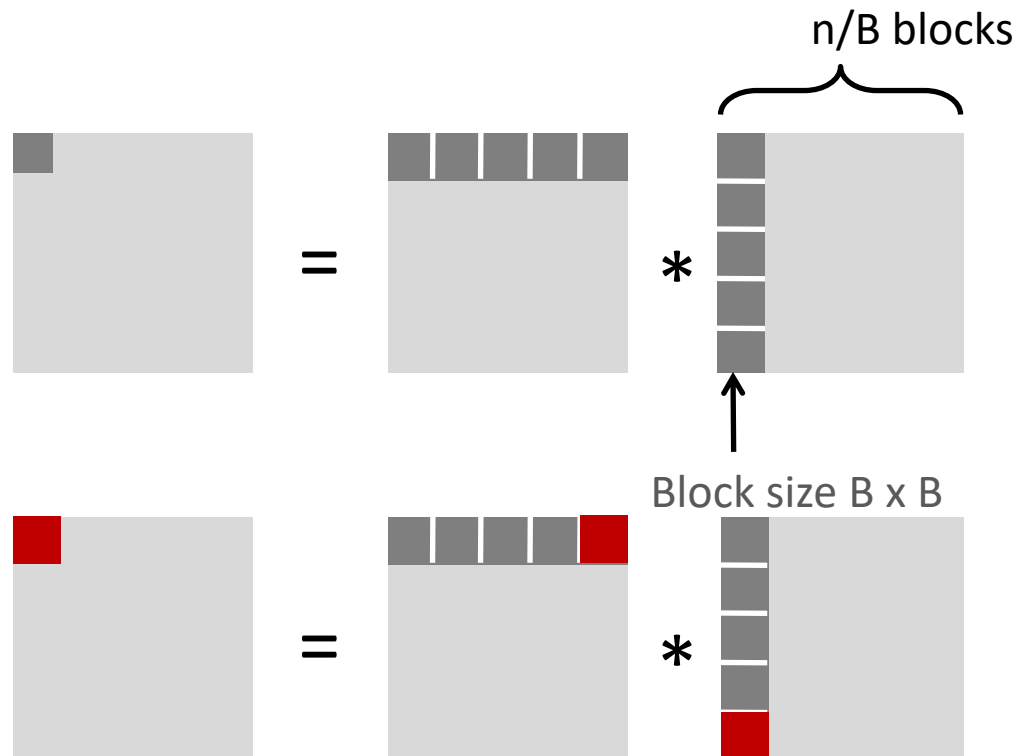*matmult/bmm.c*

j1



Block size B x B

# Cache Miss Analysis

- Assume:
  - Cache block = 8 doubles
  - Cache size C << n (much smaller than n)
  - Three blocks    fit into cache: $3B^2 < C$

- First (block) iteration:
  - $B^2/8$ misses for each ▮ck
  - $2n/B * B^2/8 = nB/4$
    (omitting matrix c)

  - Afterwards in cache
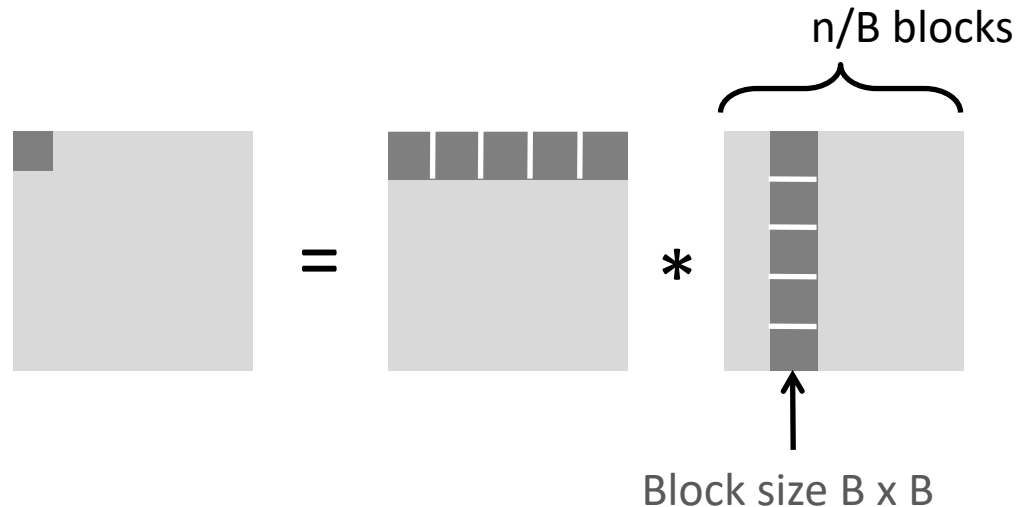    (schematic)

n/B blocks

=  ▮  *  ▮

Block size B x B

=  ▮  *  ▮

# Cache Miss Analysis

- Assume:
  - Cache block = 8 doubles
  - Cache size C << n (much smaller than n)
  - Three blocks ■ fit into cache: $3B^2 < C$

- Second (block) iteration:
  - Same as first iteration
  - $2n/B * B^2/8 = nB/4$

$$= \quad * \quad$$

n/B blocks

Block size B x B

- Total misses:
  - $nB/4 * (n/B)^2 = n^3/(4B)$
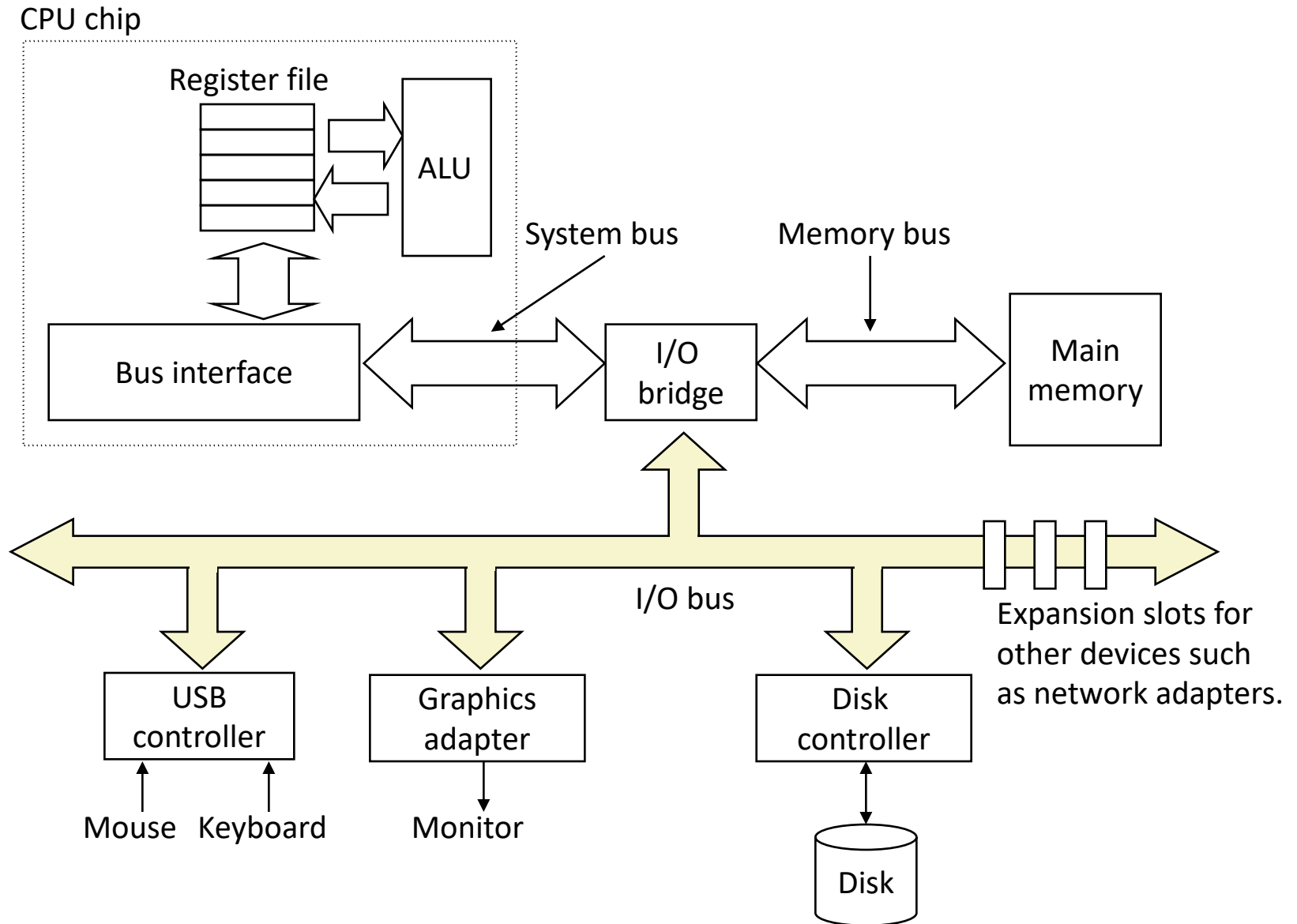
# Blocking Summary

- No blocking: $(9/8) * n^3$
- Blocking: $1/(4B) * n^3$

- Suggest largest possible block size B, but limit $3B^2 < C$!

- Reason for dramatic difference:
  - Matrix multiplication has inherent temporal locality:
    - Input data: $3n^2$, computation $2n^3$
    - Every array elements used $O(n)$ times!
  - But program has to be written properly

# Cache Summary
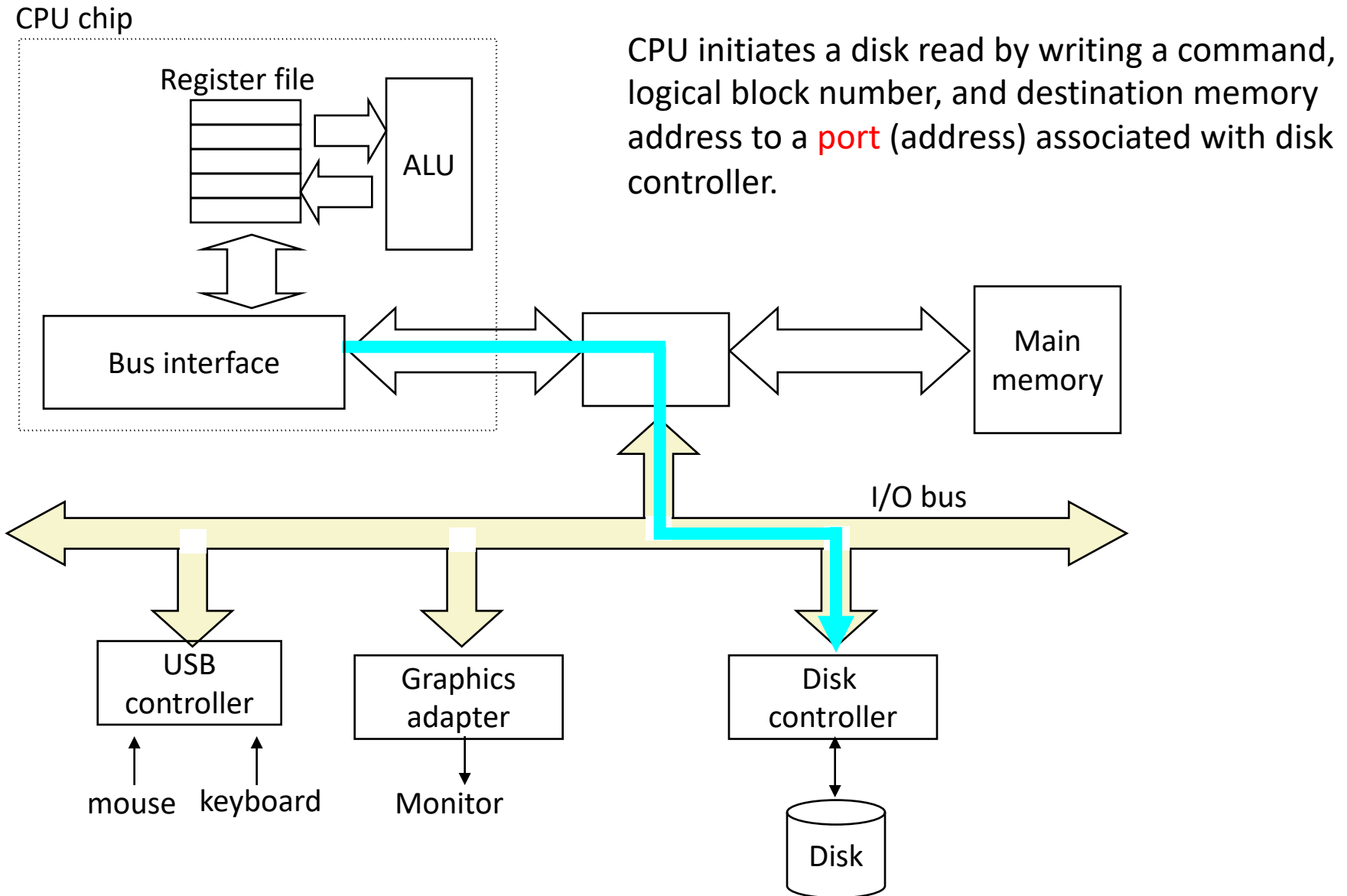
- Cache memories can have significant performance impact

- You can write your programs to exploit this!
  - Focus on the inner loops, where bulk of computations and memory accesses occur.
  - Try to maximize spatial locality by reading data objects with sequentially with stride 1.
  - Try to maximize temporal locality by using a data object as often as possible once it's read from memory.

# I/O Bus

CPU chip

Register file

ALU

System bus

Memory bus

Bus interface

I/O bridge

Main memory

I/O bus

USB controller

Graphics adapter

Disk controller

Expansion slots for other devices such as network adapters.

Mouse    Keyboard

Monitor

Disk

# Reading a Disk Sector (1)

CPU chip

Register file

ALU

Bus interface

Main memory

I/O bus

USB controller

Graphics adapter

Disk controller

mouse   keyboard

Monitor

Disk

CPU initiates a disk read by writing a command, logical block number, and destination memory address to a port (address) associated with disk controller.

# Reading a Disk Sector (2)

CPU chip

Register file

ALU

Bus interface

Disk controller reads the sector and performs a direct memory access (DMA) transfer into main memory.

Main memory

I/O bus

USB controller

Graphics adapter

Disk controller

Mouse    Keyboard

Monitor

Disk

# Reading a Disk Sector (3)

CPU chip

Register file

ALU

Bus interface

Main memory

I/O bus

USB controller

Graphics adapter

Disk controller

Mouse   Keyboard

Monitor

Disk

When the DMA transfer completes, the disk controller notifies the CPU with an *interrupt* (i.e., asserts a special "interrupt" pin on the CPU)

# Solid State Disks (SSDs)



- Pages: 512KB to 4KB, Blocks: 32 to 128 pages
- Data read/written in units of pages.
- Page can be written only after its block has been erased
- A block wears out after about 100,000 repeated writes.

# SSD Performance Characteristics

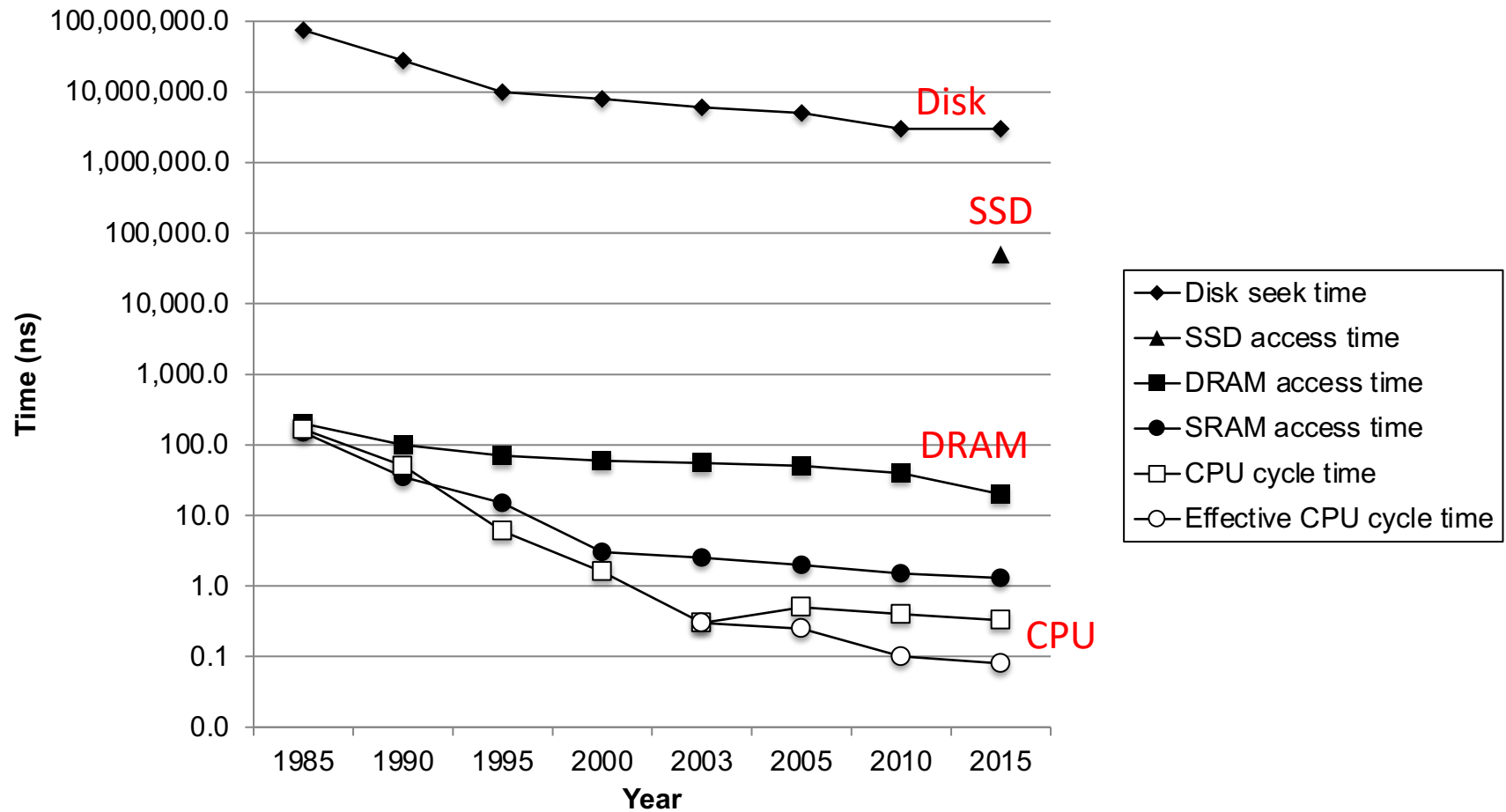| | | | |
|---|---|---|---|
| Sequential read tput | 550 MB/s | Sequential write tput | 470 MB/s |
| Random read tput | 365 MB/s | Random write tput | 303 MB/s |
| Avg seq read time | 50 us | Avg seq write time | 60 us |

- Sequential access faster than random access
  - Common theme in the memory hierarchy
- Random writes are somewhat slower
  - Erasing a block takes a long time (~1 ms)
  - Modifying a block page requires all other pages to be copied to new block
  - In earlier SSDs, the read/write gap was much larger.

Source: Intel SSD 730 product specification.

# SSD Tradeoffs vs Rotating Disks

- Advantages
  - No moving parts → faster, less power, more rugged

- Disadvantages
  - Have the potential to wear out
    - Mitigated by "wear leveling logic" in flash translation layer
    - E.g. Intel SSD 730 guarantees 128 petabyte (128 x $10^{15}$ bytes) of writes before they wear out
  - In 2015, about 30 times more expensive per byte

- Applications
  - MP3 players, smart phones, laptops
  - Beginning to appear in desktops and servers

# The CPU-Memory Gap

The gap widens between DRAM, disk, and CPU speeds.
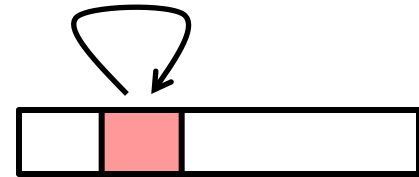
# Locality to the Rescue!

The key to bridging this CPU-Memory gap is a fundamental property of computer programs known as locality
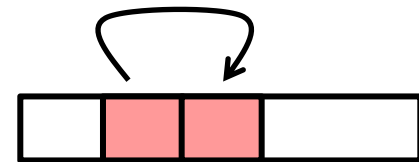
- **Locality of reference**

# 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

# 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

# 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;
}
```

# 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;
}
```
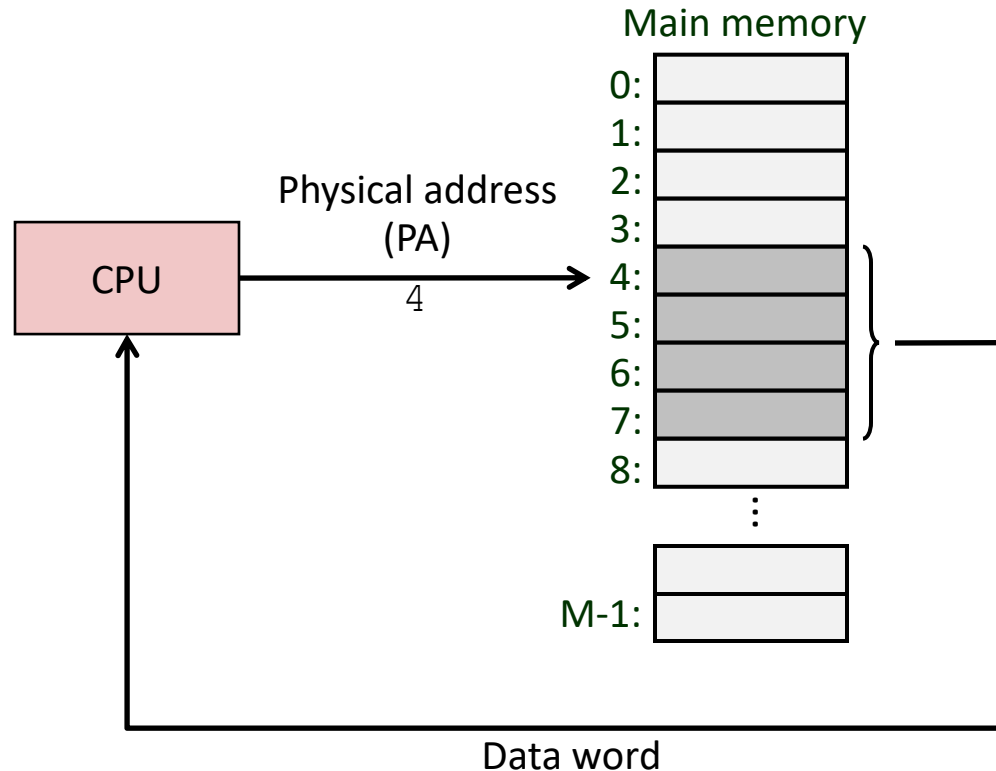
# Locality Example

- 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;
}
```

# A System Using Physical Addressing



Main memory

Physical address (PA)

CPU

4

0:
1:
2:
3:
4:
5:
6:
7:
8:
⋮
M-1:

Data word

- Used in "simple" systems like embedded microcontrollers in devices like cars, elevators, and digital picture frames

# A System Using Virtual Addressing

Main memory

CPU Chip

CPU → Virtual address (VA) `4100` → MMU → Physical address (PA) `4` → Main memory

0:
1:
2:
3:
4:
5:
6:
7:
8:
⋮
M-1:

Data word

- Used in all modern servers, laptops, and smart phones
- One of the great ideas in computer science
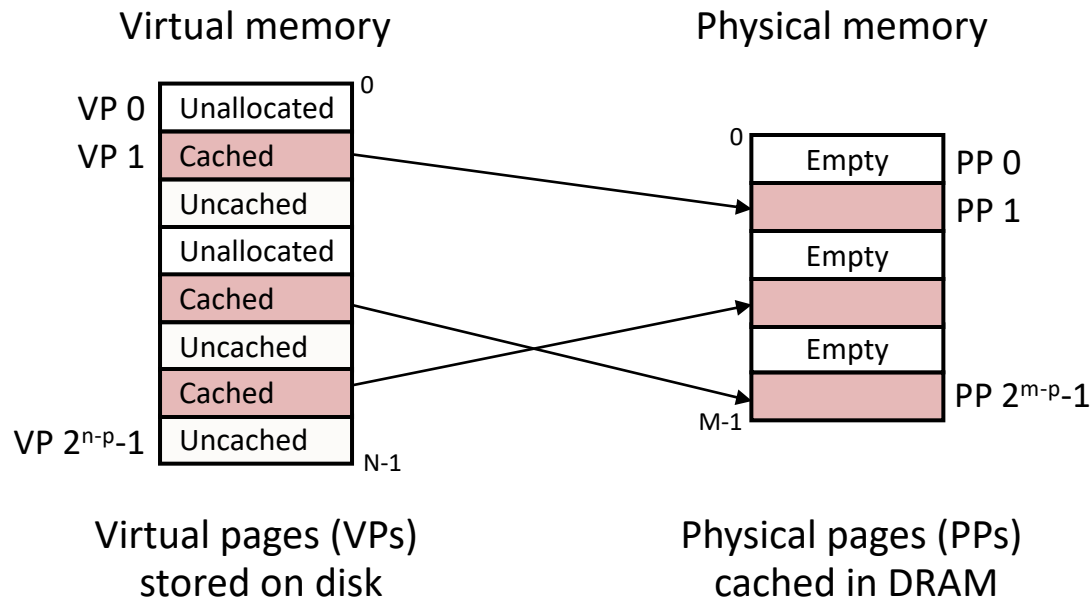
# Address Spaces

- **Linear address space:** Ordered set of contiguous non-negative integer addresses:

$$\{0, 1, 2, 3 \dots \}$$

- **Virtual address space:** Set of $N = 2^n$ virtual addresses

$$\{0, 1, 2, 3, \dots, N-1\}$$

- **Physical address space:** Set of $M = 2^m$ physical addresses

$$\{0, 1, 2, 3, \dots, M-1\}$$

# Why Virtual Memory (VM)?

- Uses main memory efficiently
  - Use DRAM as a cache for parts of a virtual address space

- Simplifies memory management
  - Each process gets the same uniform linear address space

- Isolates address spaces
  - One process can't interfere with another's memory
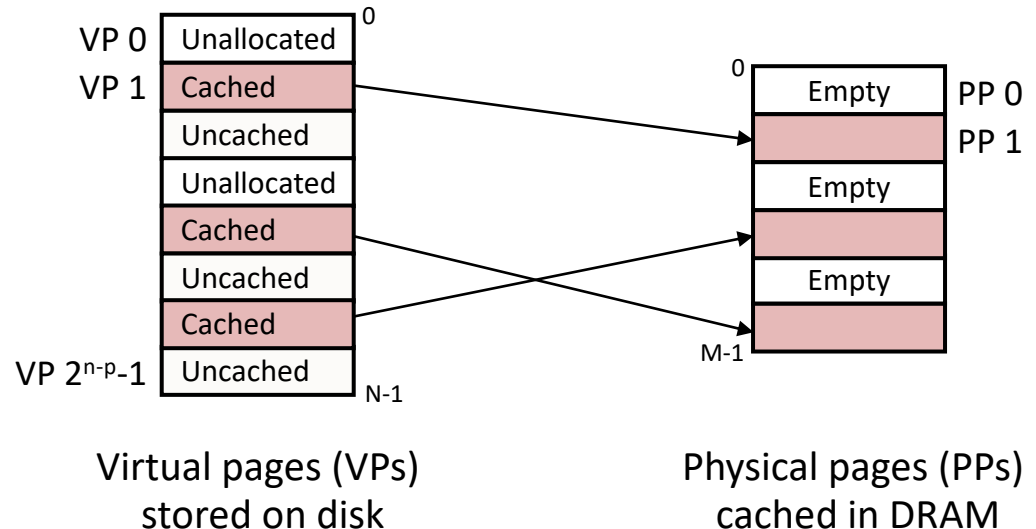  - User program cannot access privileged kernel information and code

# VM as a Tool for Caching

- Conceptually, *virtual memory* is an array of N contiguous bytes stored on disk.

- The contents of the array on disk are cached in *physical memory* (*DRAM cache*)
  - These cache blocks are called *pages* (size is $P = 2^p$ bytes)

Virtual memory

| | |
|---|---|
| VP 0 | Unallocated |
| VP 1 | Cached |
| | Uncached |
| | Unallocated |
| | Cached |
| | Uncached |
| | Cached |
| VP $2^{n-p}-1$ | Uncached |

0

N-1

Physical memory

| | |
|---|---|
| | Empty |
| | |
| | Empty |
| | |
| | Empty |
| | |

0

PP 0

PP 1

PP $2^{m-p}-1$

M-1

Virtual pages (VPs)
stored on disk

Physical pages (PPs)
cached in DRAM

# Virtual Pages

- At any point in time, the set of virtual pages is partitioned into three disjoint subsets;

- **Unallocated**: pages that have not been allocated by VM system. Unallocated blocks do not have any data associated with them, thus do not occupy any space on the disk

- **Cached**: allocated pages that are currently cached in physical memory

- **Uncached**: allocated pages that are not cached in physical memory

| | | 0 |
|---|---|---|
| VP 0 | Unallocated | |
| VP 1 | Cached | |
| | Uncached | |
| | Unallocated | |
| | Cached | |
| | Uncached | |
| | Cached | |
| VP $2^{n-p}-1$ | Uncached | N-1 |

| | | |
|---|---|---|
| 0 | Empty | PP 0 |
| | | PP 1 |
| | Empty | |
| | | |
| | Empty | |
| | | |
| M-1 | | |

Virtual pages (VPs) stored on disk

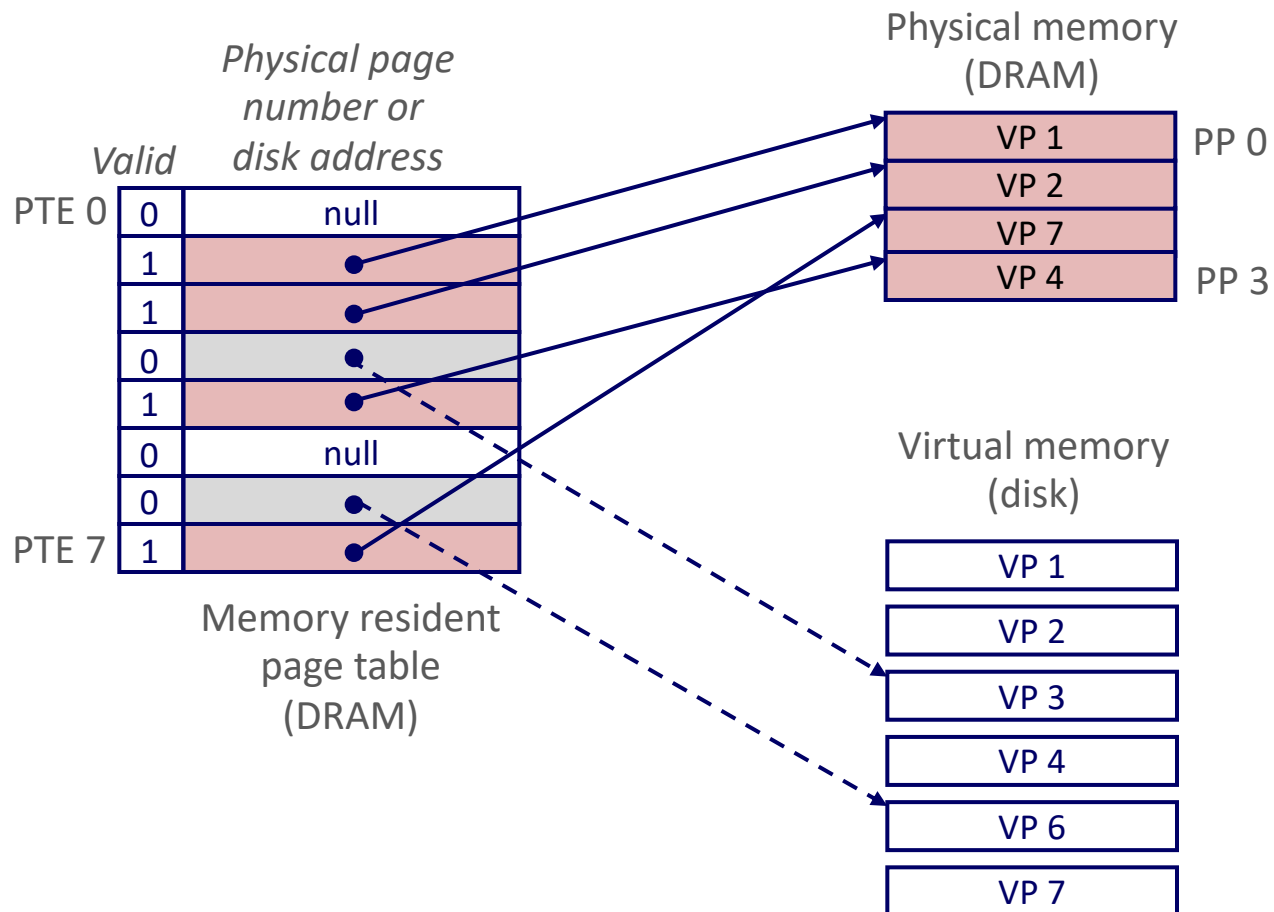Physical pages (PPs) cached in DRAM

# DRAM Cache Organization

- DRAM cache organization driven by the enormous miss penalty
  - DRAM is about **10x** slower than SRAM
  - Disk is about **10,000x** slower than DRAM

- Consequences
  - Large page (block) size: typically 4 KB, sometimes 4 MB
  - Fully associative
    - Any VP can be placed in any PP
    - Requires a "large" mapping function – different from cache memories
  - Highly sophisticated, expensive replacement algorithms
    - Too complicated and open-ended to be implemented in hardware
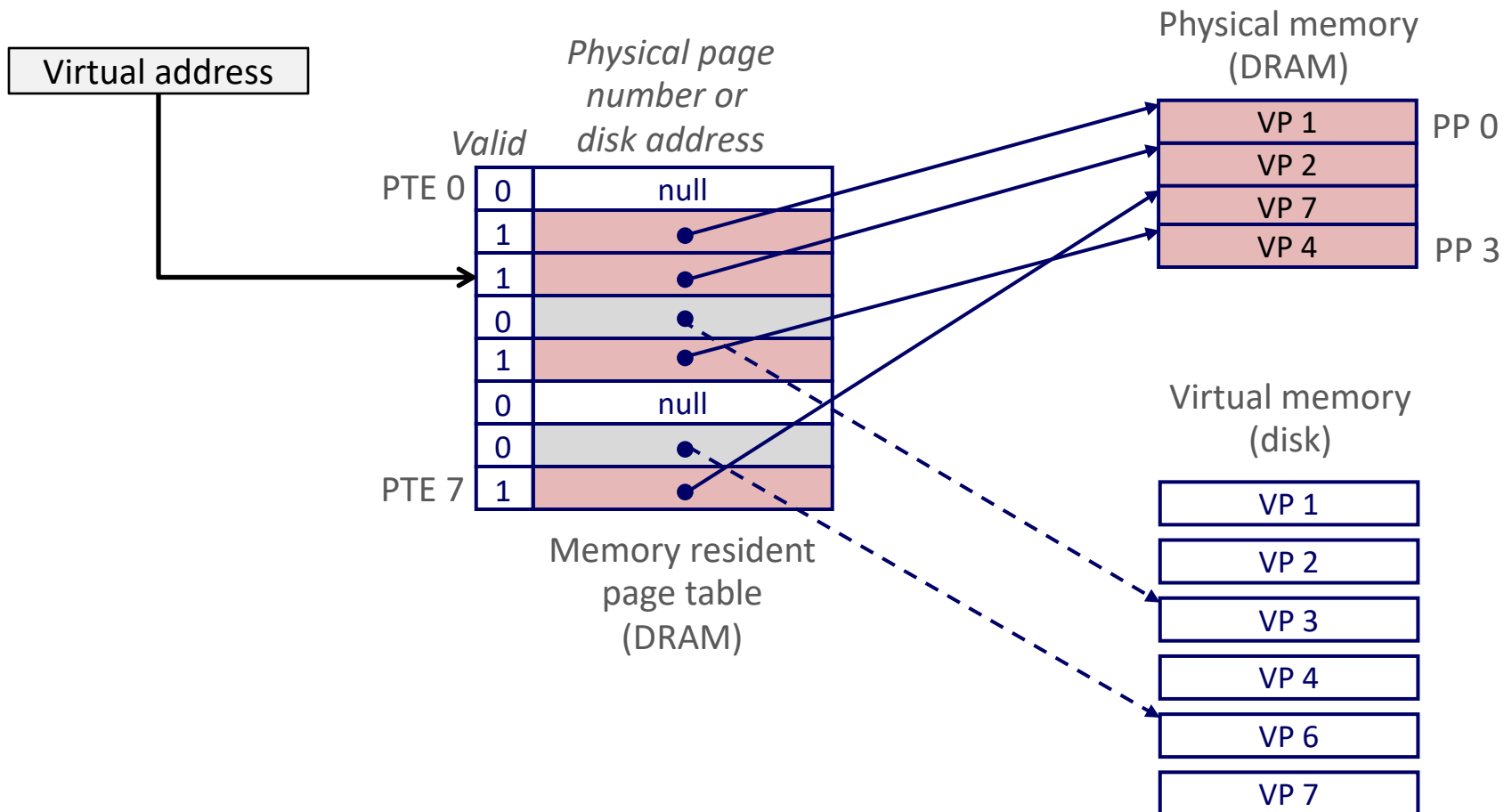  - Write-back rather than write-through

# Enabling Data Structure: Page Table

- A *page table* is an array of page table entries (PTEs) that maps virtual pages to physical pages.
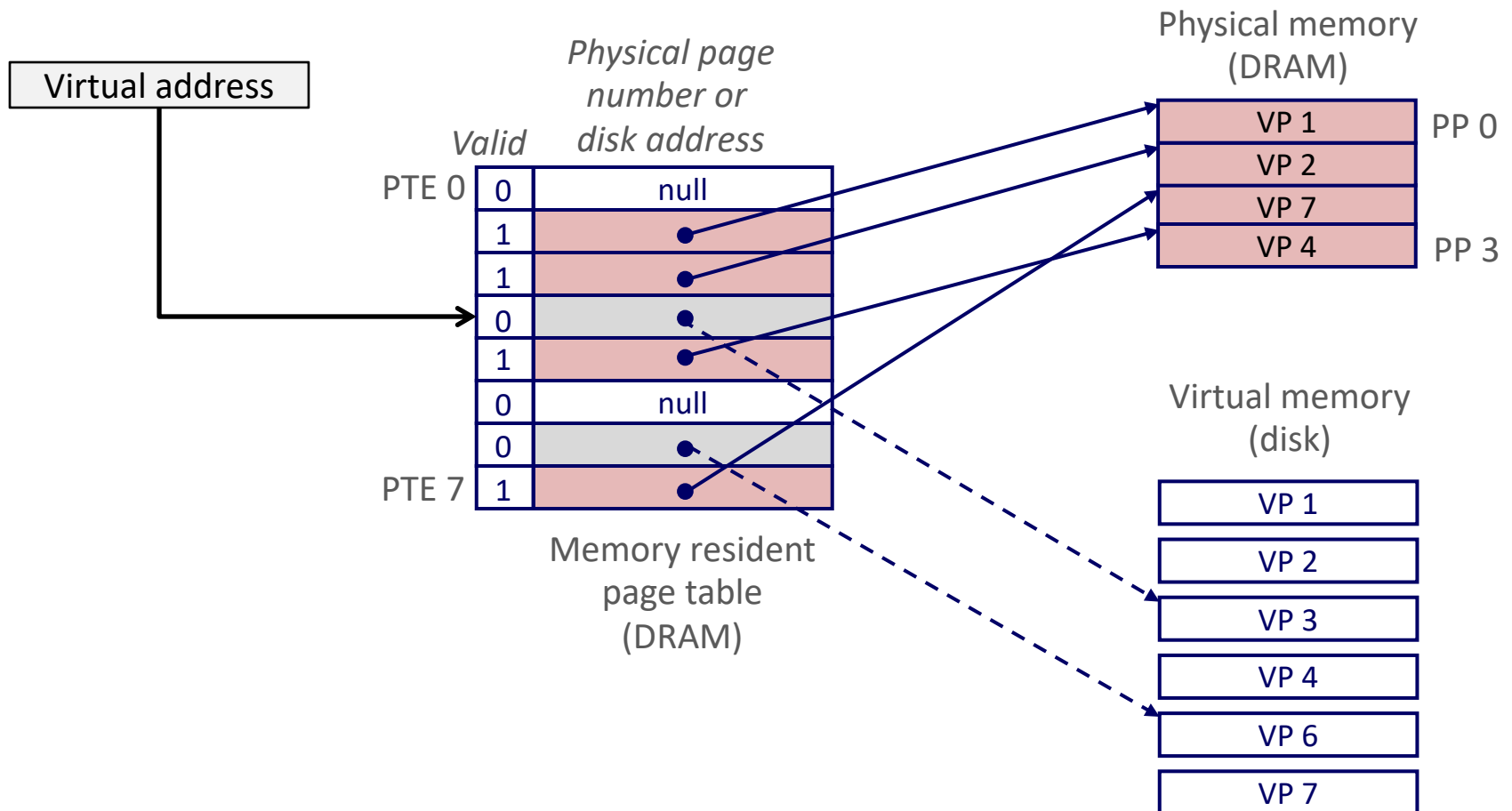  - Per-process kernel data structure in DRAM

# Page Hit

- *Page hit:* reference to VM word that is in physical memory (DRAM cache hit)
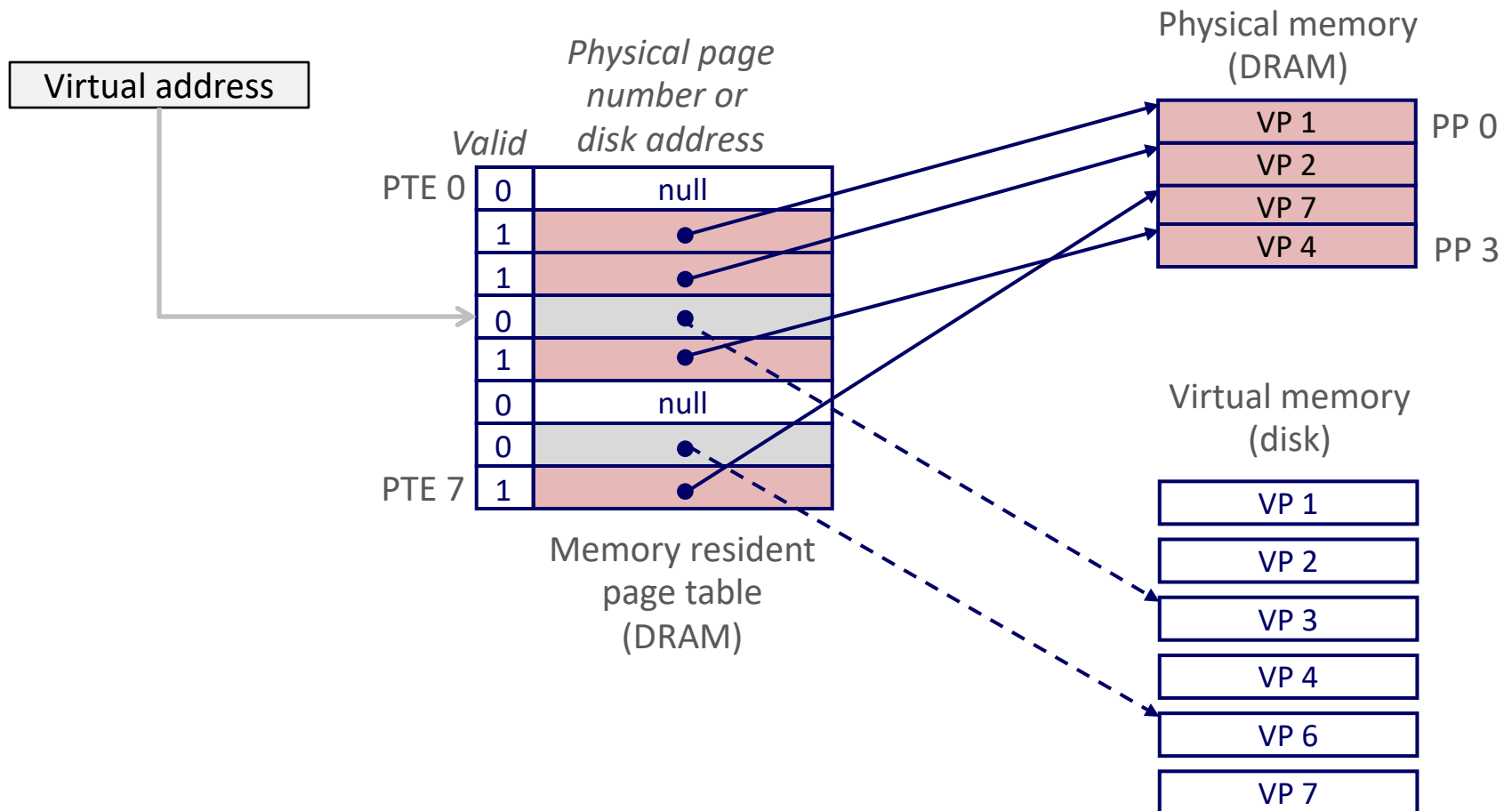
# Page Fault

- *Page fault:* reference to VM word that is not in physical memory (DRAM cache miss)
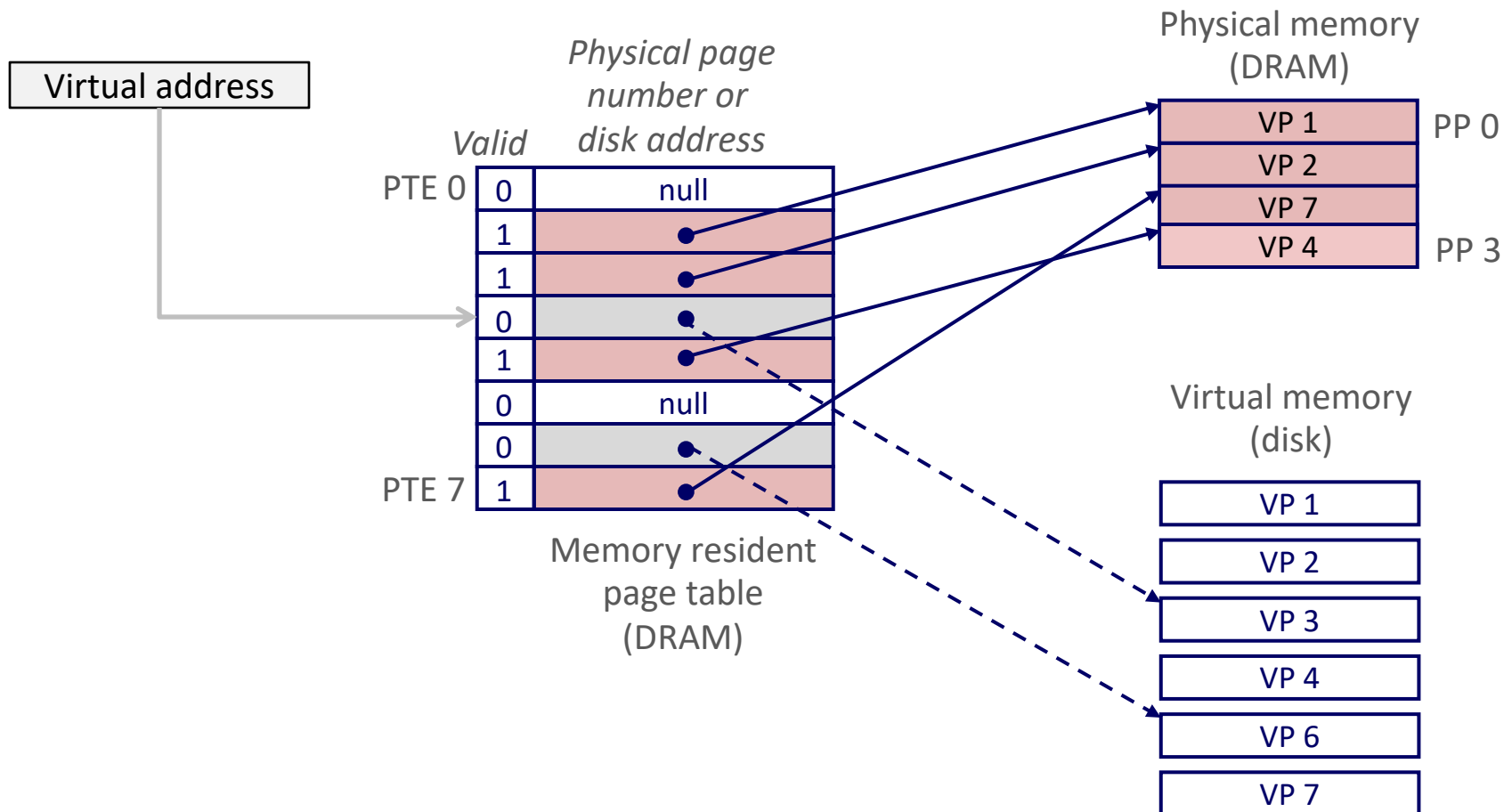
# Handling Page Fault

- Page miss causes page fault (an exception)

# Handling Page Fault

- Page miss causes page fault (an exception)
- Page fault handler selects a victim to be evicted (here VP 4)

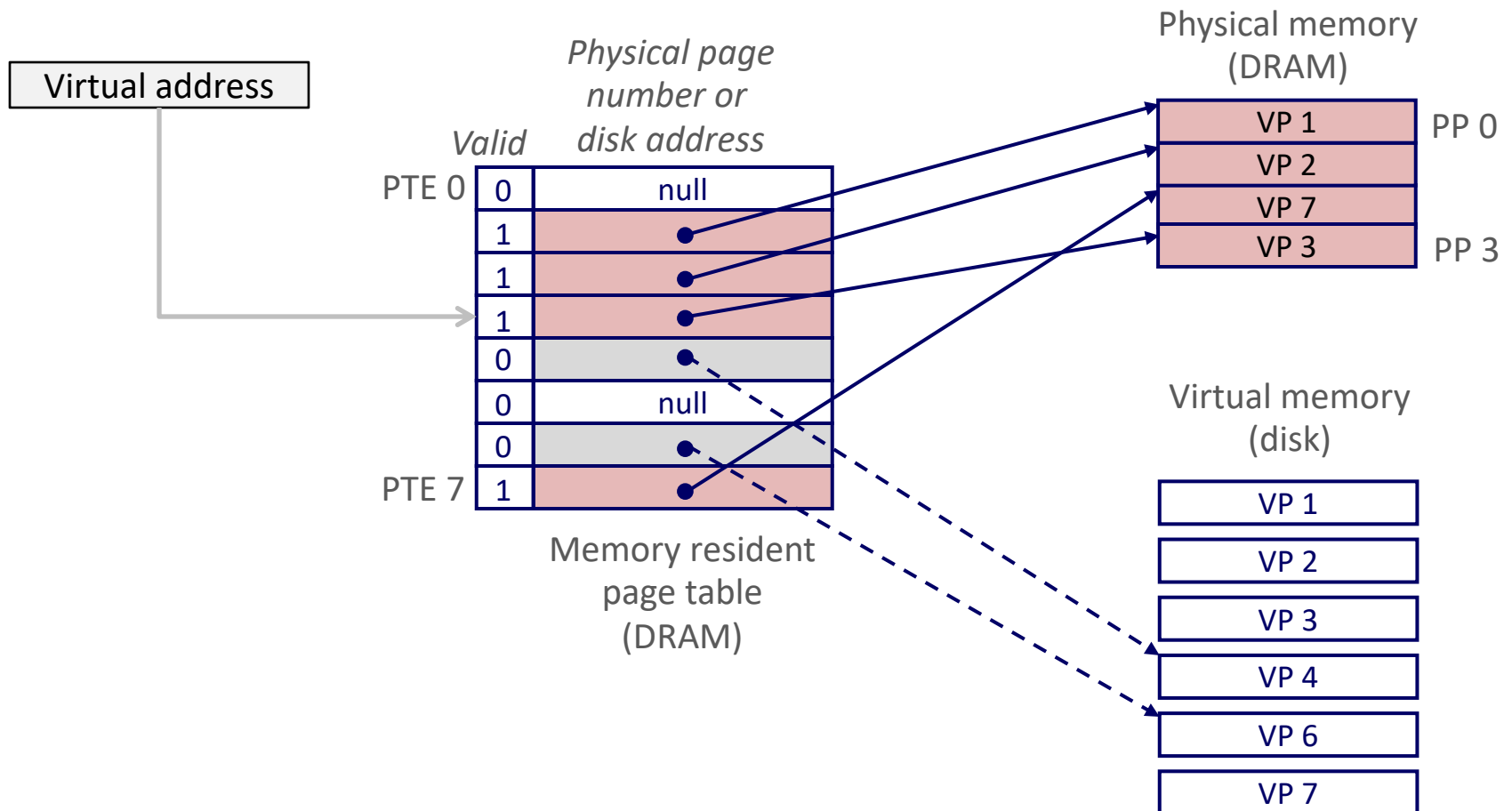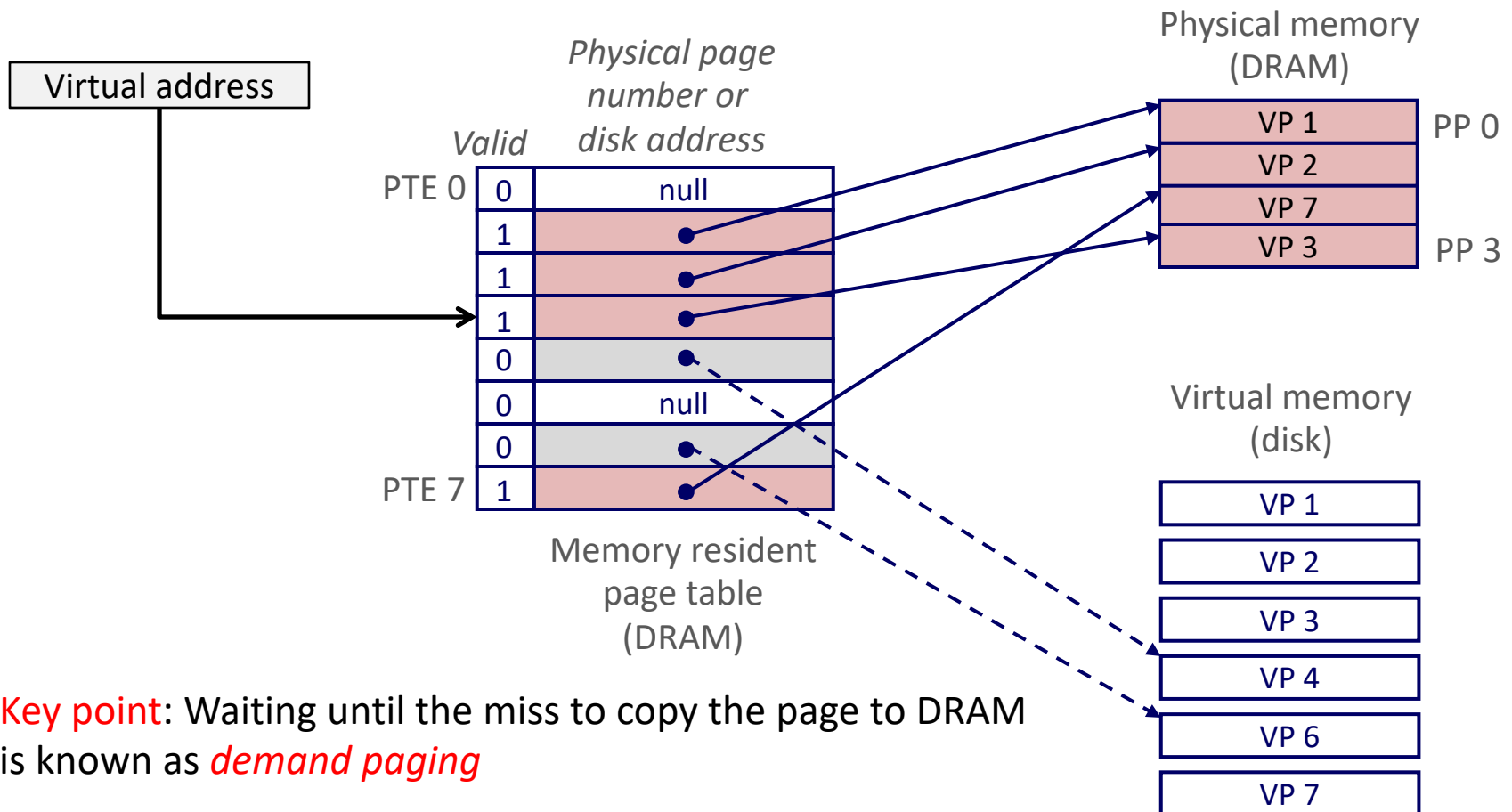# Handling Page Fault

- Page miss causes page fault (an exception)
- Page fault handler selects a victim to be evicted (here VP 4)

# Handling Page Fault

- Page miss causes page fault (an exception)
- Page fault handler selects a victim to be evicted (here VP 4)
- Offending instruction is restarted: page hit!



**Key point**: Waiting until the miss to copy the page to DRAM is known as *demand paging*

# Allocating Pages

- Allocating a new page (VP 5) of virtual memory.

Physical page number or disk address

Physical memory (DRAM)

Valid

| | | | |
|---|---|---|---|
| PTE 0 | 0 | null | |
| | 1 | | |
| | 1 | | |
| | 1 | | |
| | 0 | | |
| | 0 | | |
| | 0 | | |
| PTE 7 | 1 | | |

Memory resident page table (DRAM)

Physical memory (DRAM)

| | |
|---|---|
| VP 1 | PP 0 |
| VP 2 | |
| VP 7 | |
| VP 3 | PP 3 |

Virtual memory (disk)

| |
|---|
| VP 1 |
| VP 2 |
| VP 3 |
| VP 4 |
| VP 5 |
| VP 6 |
| VP 7 |

# Locality to the Rescue Again!

- Virtual memory seems terribly inefficient, but it works because of locality.

- At any point in time, programs tend to access a set of active virtual pages called the *working set*
  - Programs with better temporal locality will have smaller working sets

- If (working set size < main memory size)
  - Good performance for one process after compulsory misses

- If ( SUM(working set sizes) > main memory size )
  - *Thrashing:* Performance meltdown where pages are swapped (copied) in and out continuously

# VM as a Tool for Memory Management

- Key idea: each process has its own virtual address space
  - It can view memory as a simple linear array
  - Mapping function scatters addresses through physical memory
    - Well-chosen mappings can improve locality



*Virtual Address Space for Process 1:*

0

VP 1
VP 2
…

N-1

*Address translation*

0

PP 2

PP 6

PP 8

…

M-1

*Physical Address Space (DRAM)*

**(e.g., read-only library code)**

*Virtual Address Space for Process 2:*

0

VP 1
VP 2
…

N-1

# VM as a Tool for Memory Management

- **Simplifying memory allocation**
  - Each virtual page can be mapped to any physical page
  - A virtual page can be stored in different physical pages at different times
- **Sharing code and data among processes**
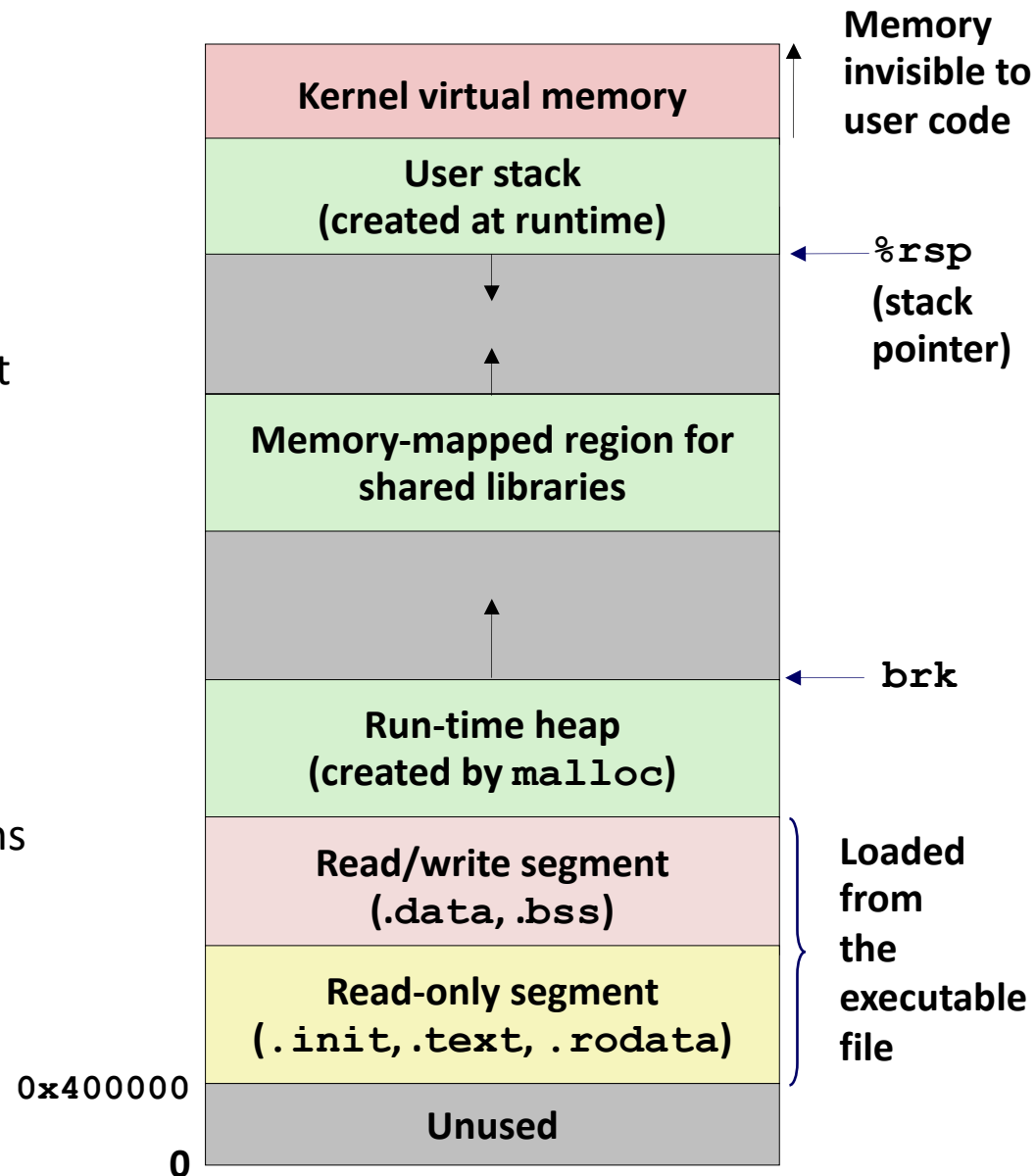  - Map virtual pages to the same physical page (here: PP 6)



Virtual Address Space for Process 1:

Virtual Address Space for Process 2:

Address translation

Physical Address Space (DRAM)

(e.g., read-only library code)

VP 1
VP 2
...
PP 2
PP 6
PP 8
0
N-1
0
N-1
0
M-1

# Simplifying Linking and Loading

- ## Linking
  - Each program has similar virtual address space
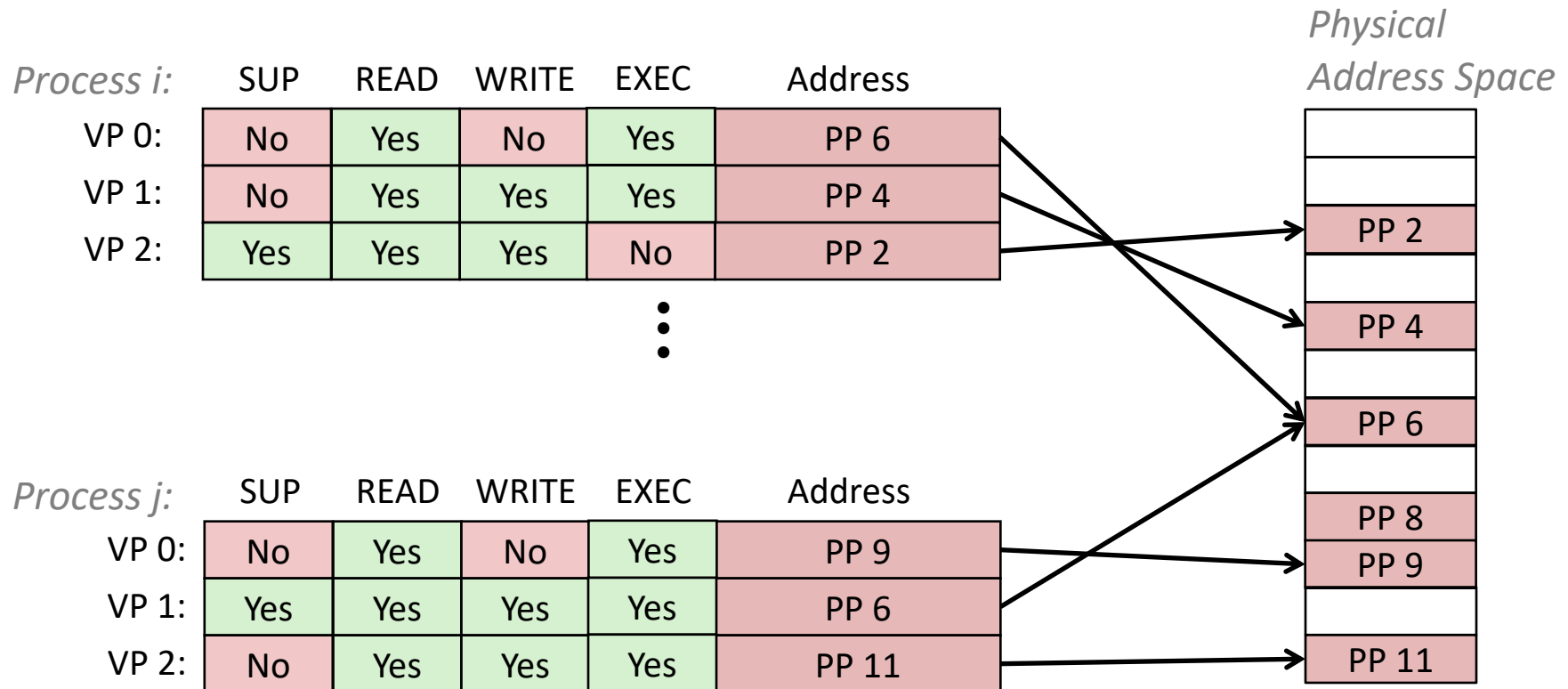  - Code, data, and heap always start at the same addresses.

- ## Loading
  - `execve` allocates virtual pages for .text and .data sections & creates PTEs marked as invalid
  - The `.text` and `.data` sections are copied, page by page, on demand by the virtual memory system

| | |
|---|---|
| **Kernel virtual memory** | **Memory invisible to user code** ↑ |
| **User stack (created at runtime)** | ← `%rsp` (stack pointer) |
| ↓ ↑ | |
| **Memory-mapped region for shared libraries** | |
| ↑ | |
| | ← `brk` |
| **Run-time heap (created by `malloc`)** | |
| **Read/write segment (`.data`, `.bss`)** | Loaded from the executable file |
| **Read-only segment (`.init`, `.text`, `.rodata`)** | |

`0x400000`

| |
|---|
| **Unused** |

`0`

# VM as a Tool for Memory Protection

- Extend PTEs with permission bits
- MMU checks these bits on each access

| Process i: | SUP | READ | WRITE | EXEC | Address |
|---|---|---|---|---|---|
| VP 0: | No | Yes | No | Yes | PP 6 |
| VP 1: | No | Yes | Yes | Yes | PP 4 |
| VP 2: | Yes | Yes | Yes | No | PP 2 |

| Process j: | SUP | READ | WRITE | EXEC | Address |
|---|---|---|---|---|---|
| VP 0: | No | Yes | No | Yes | PP 9 |
| VP 1: | Yes | Yes | Yes | Yes | PP 6 |
| VP 2: | No | Yes | Yes | Yes | PP 11 |

*Physical Address Space*

PP 2
PP 4
PP 6
PP 8
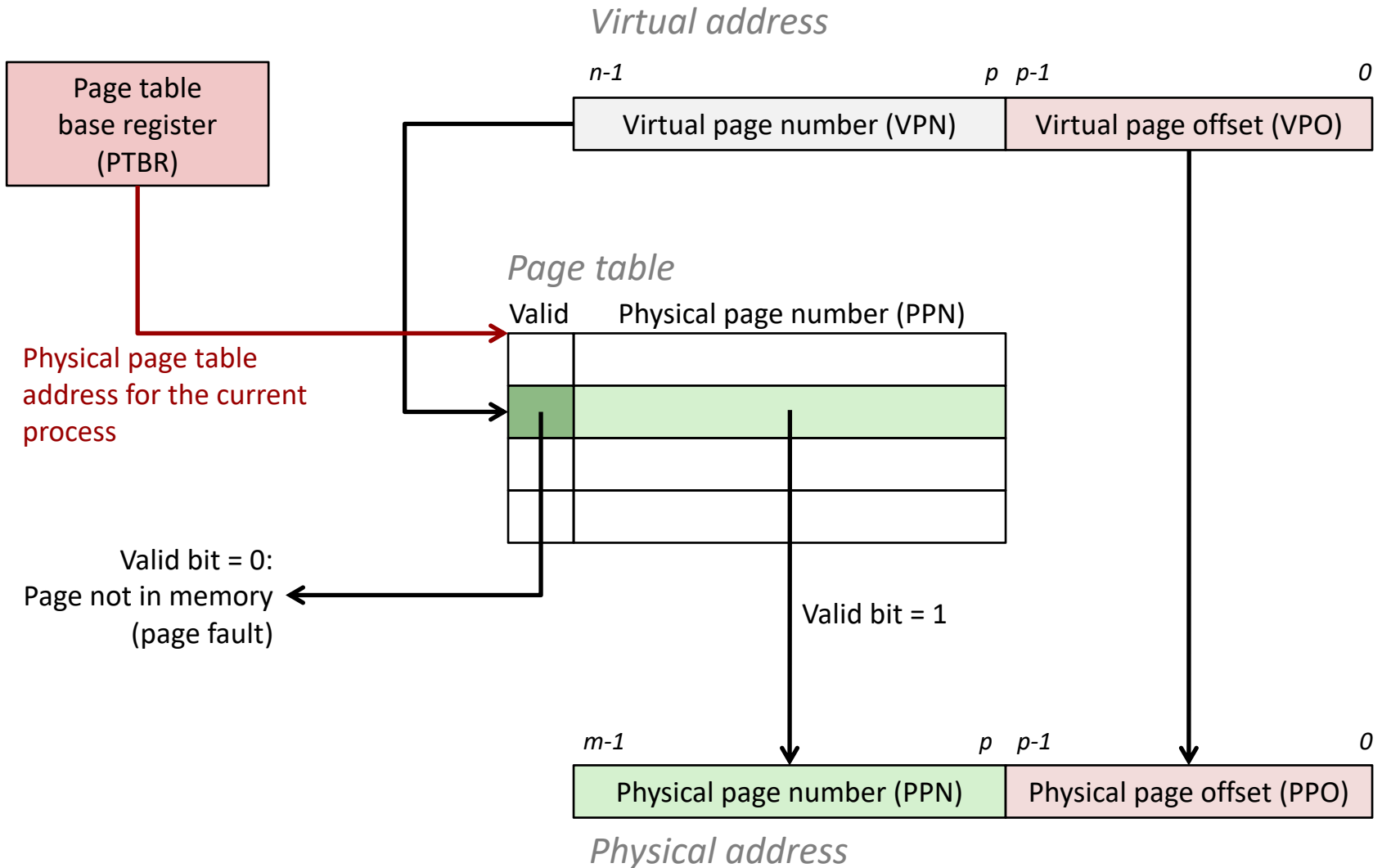PP 9
PP 11

# VM Address Translation

- Virtual Address Space
  - $V = \{0, 1, \ldots, N-1\}$
- Physical Address Space
  - $P = \{0, 1, \ldots, M-1\}$
- Address Translation
  - **MAP: $V \rightarrow P$ $U$ $\{\varnothing\}$**
  - For virtual address **a**:
    - **MAP(a) = a'** if data at virtual address **a** is at physical address **a'** in **P**
    - **MAP(a) = $\varnothing$** if data at virtual address **a** is not in physical memory
      - Either invalid or stored on disk
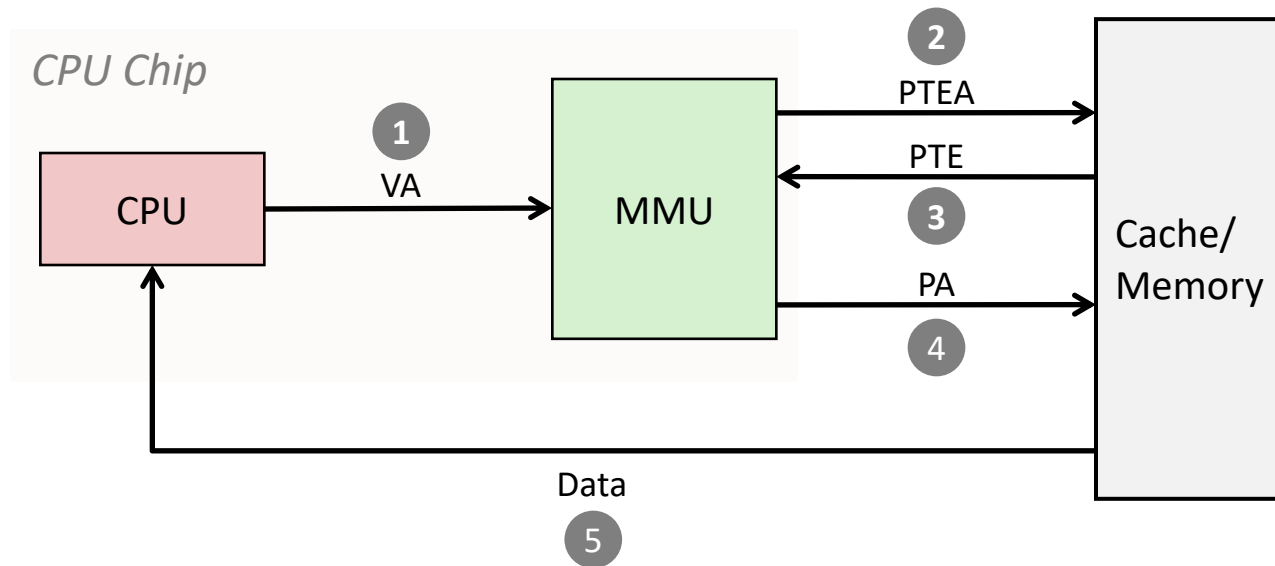
# Summary of Address Translation Symbols

- Basic Parameters
  - **N = $2^n$** : Number of addresses in virtual address space
  - **M = $2^m$** : Number of addresses in physical address space
  - **P = $2^p$** : Page size (bytes)
- Components of the virtual address (VA)
  - **TLBI**: TLB index
  - **TLBT**: TLB tag
  - **VPO**: Virtual page offset
  - **VPN**: Virtual page number
- Components of the physical address (PA)
  - **PPO**: Physical page offset (same as VPO)
  - **PPN:** Physical page number
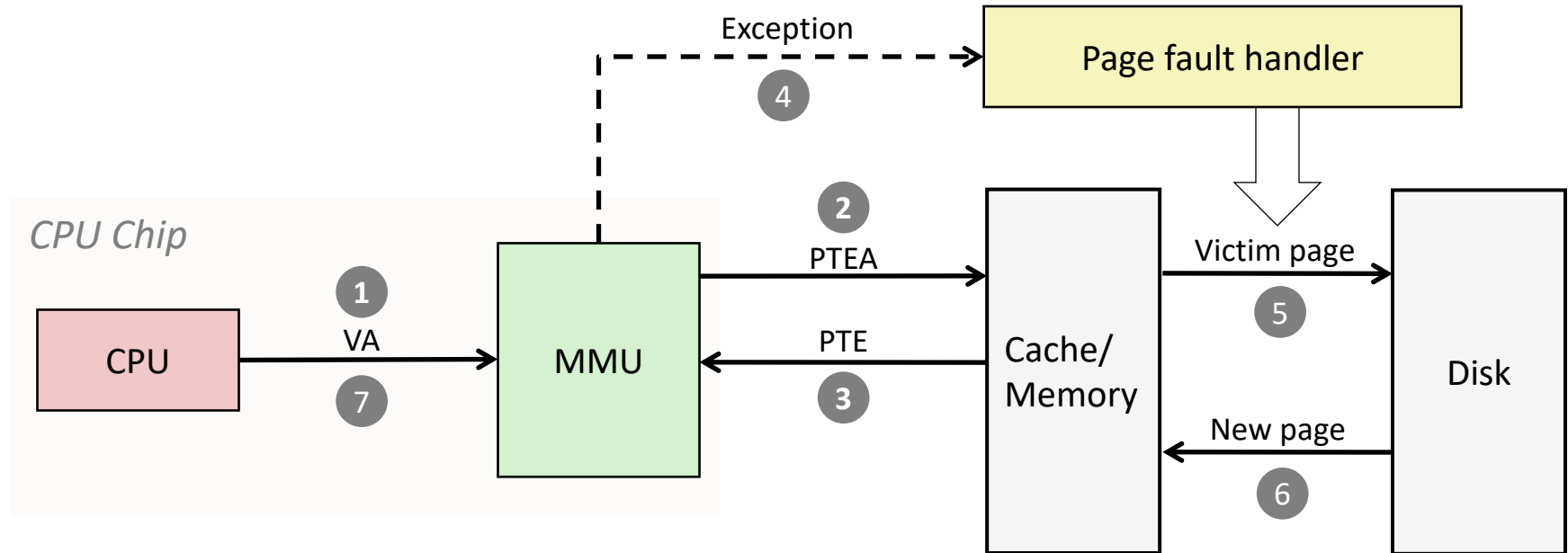
# Address Translation With a Page Table

*Virtual address*

# Address Translation: Page Hit



1) Processor sends virtual address to MMU

2-3) MMU fetches PTE from page table in memory

4) MMU sends physical address to cache/memory

5) Cache/memory sends data word to processor

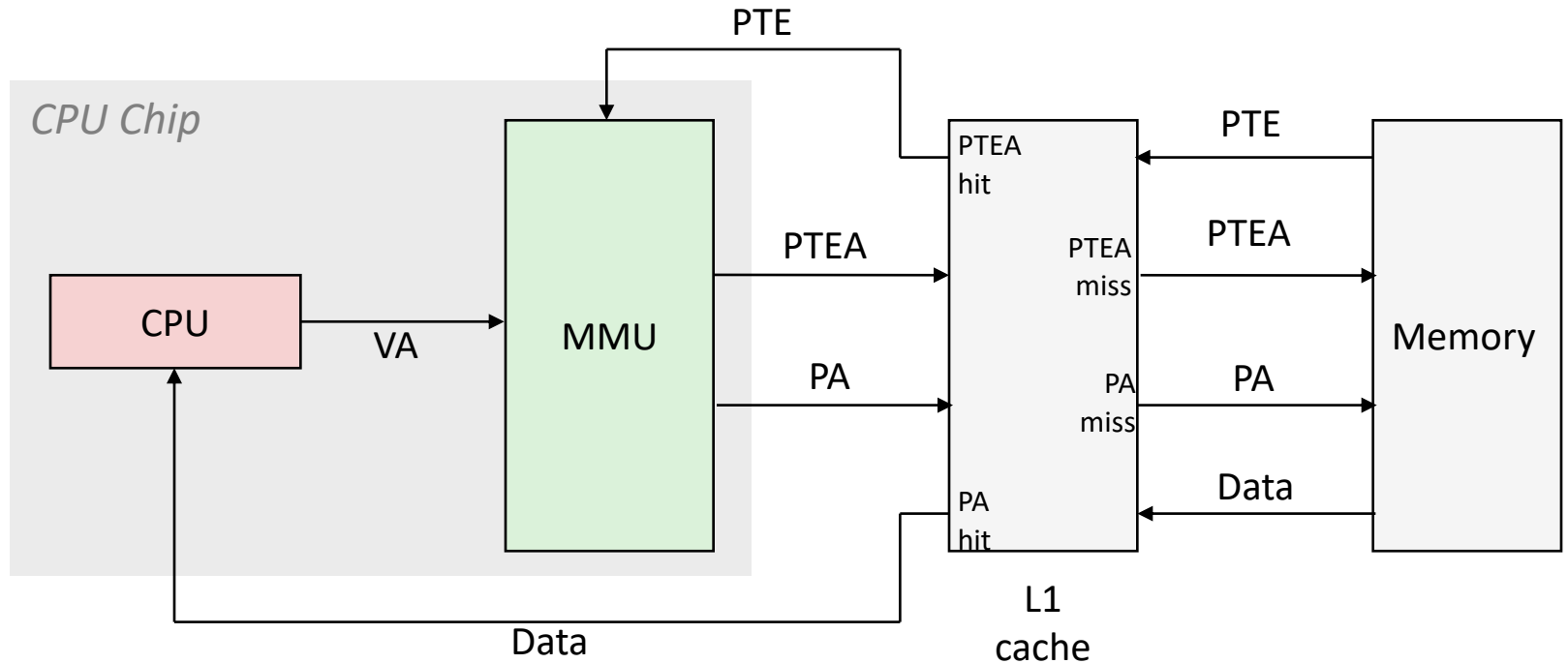*VA: virtual address, PA: physical address, PTE: page table entry, PTEA = PTE address*

# Address Translation: Page Fault



1) Processor sends virtual address to MMU

2-3) MMU fetches PTE from page table in memory

4) Valid bit is zero, so MMU triggers page fault exception

5) Handler identifies victim (and, if dirty, pages it out to disk)

6) Handler pages in new page and updates PTE in memory

7) Handler returns to original process, restarting faulting instruction

*VA: virtual address, PA: physical address, PTE: page table entry, PTEA = PTE address*

# Integrating VM and Cache

PTE

CPU Chip

PTE

CPU

VA

MMU

PTEA

PTEA
hit

PTEA
miss

PTEA

Memory

PA

PA
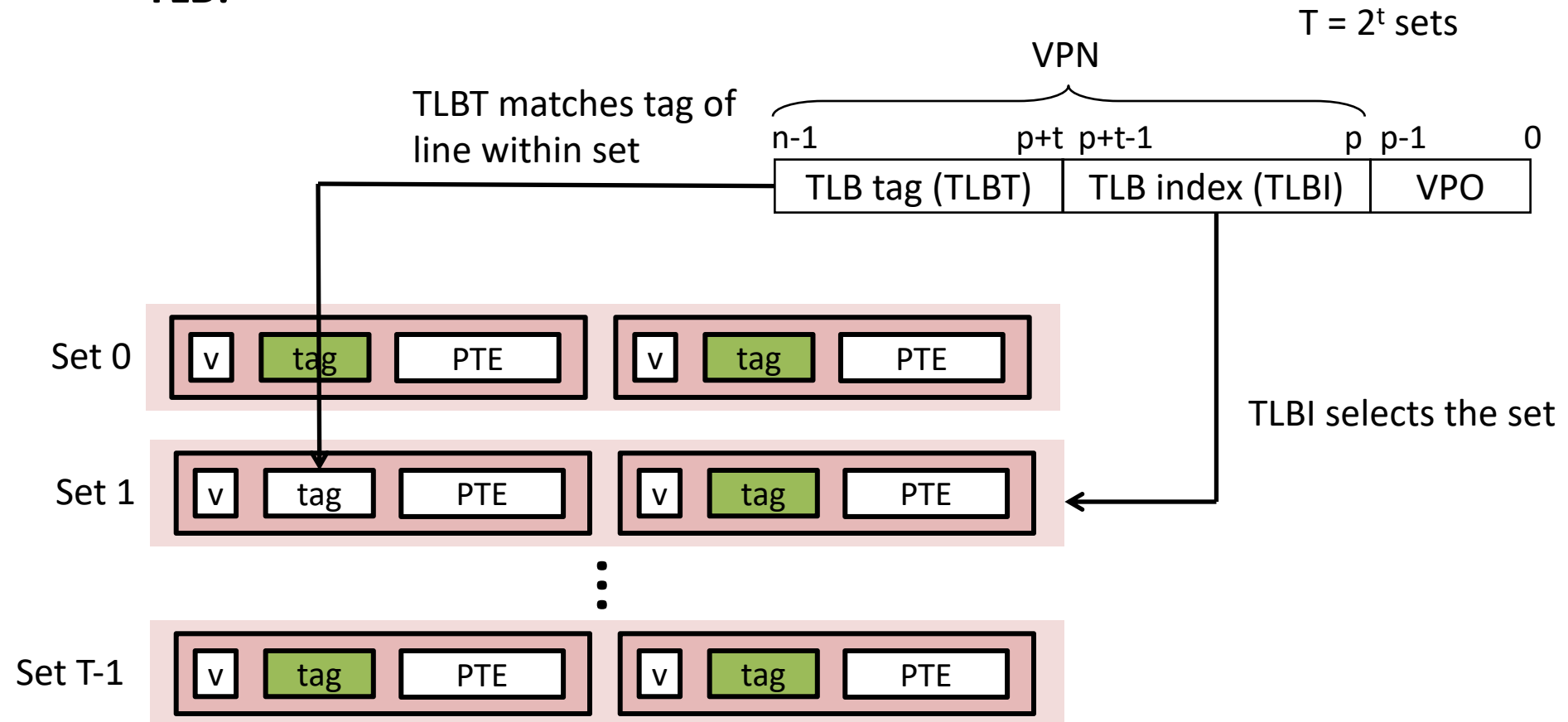miss

PA

PA
hit

Data

Data

L1
cache

*VA: virtual address, PA: physical address, PTE: page table entry, PTEA = PTE address*
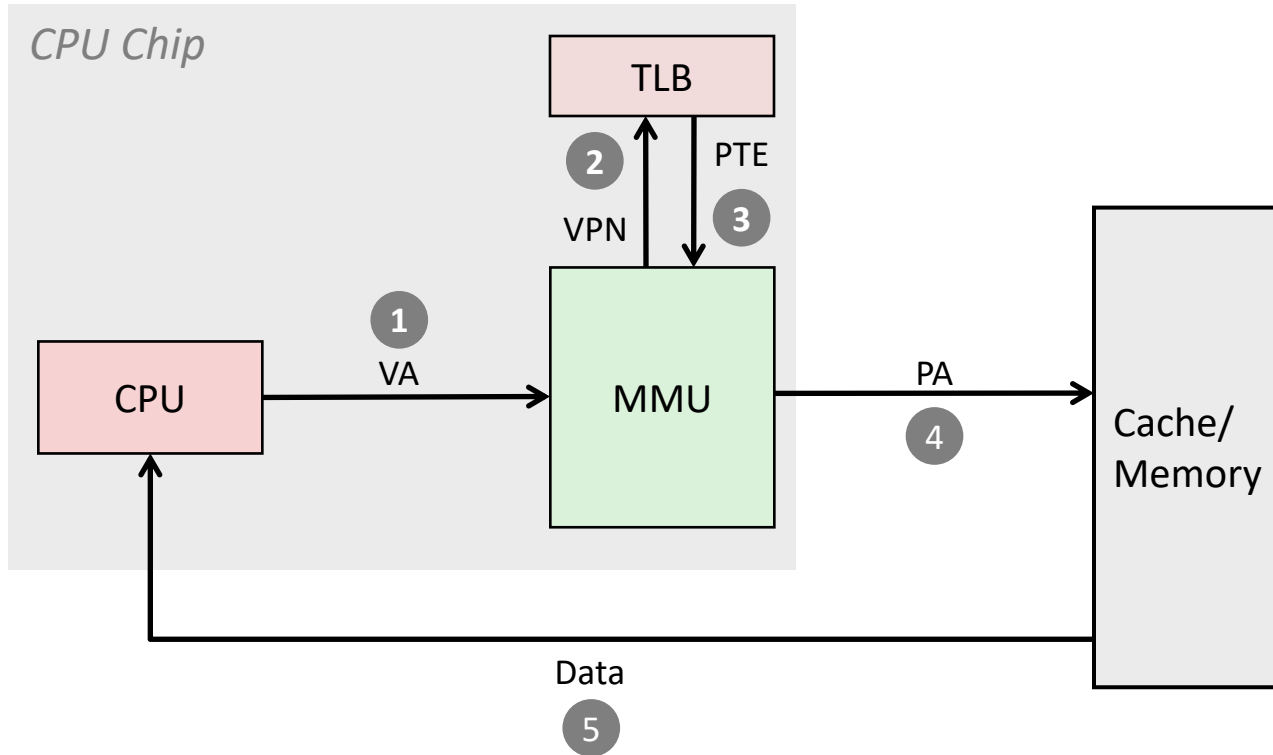
# Speeding up Translation with a TLB

- **Page table entries (PTEs) are cached in L1 like any other memory word**
  - PTEs may be evicted by other data references
  - PTE hit still requires a small L1 delay

- **Solution: *Translation Lookaside Buffer* (TLB)**
  - Small set-associative hardware cache in MMU
  - Maps virtual page numbers to  physical page numbers
  - Contains complete page table entries for small number of pages

# Accessing the TLB

- **MMU uses the VPN portion of the virtual address to access the TLB:**

$T = 2^t$ sets

VPN

TLBT matches tag of line within set

| $n-1$ | | $p+t$ $p+t-1$ | | $p$ $p-1$ | $0$ |
|---|---|---|---|---|---|
| TLB tag (TLBT) | | TLB index (TLBI) | | VPO | |

Set 0

| v | tag | PTE | | v | tag | PTE |

Set 1

| v | tag | PTE | | v | tag | PTE |

TLBI selects the set

⋮

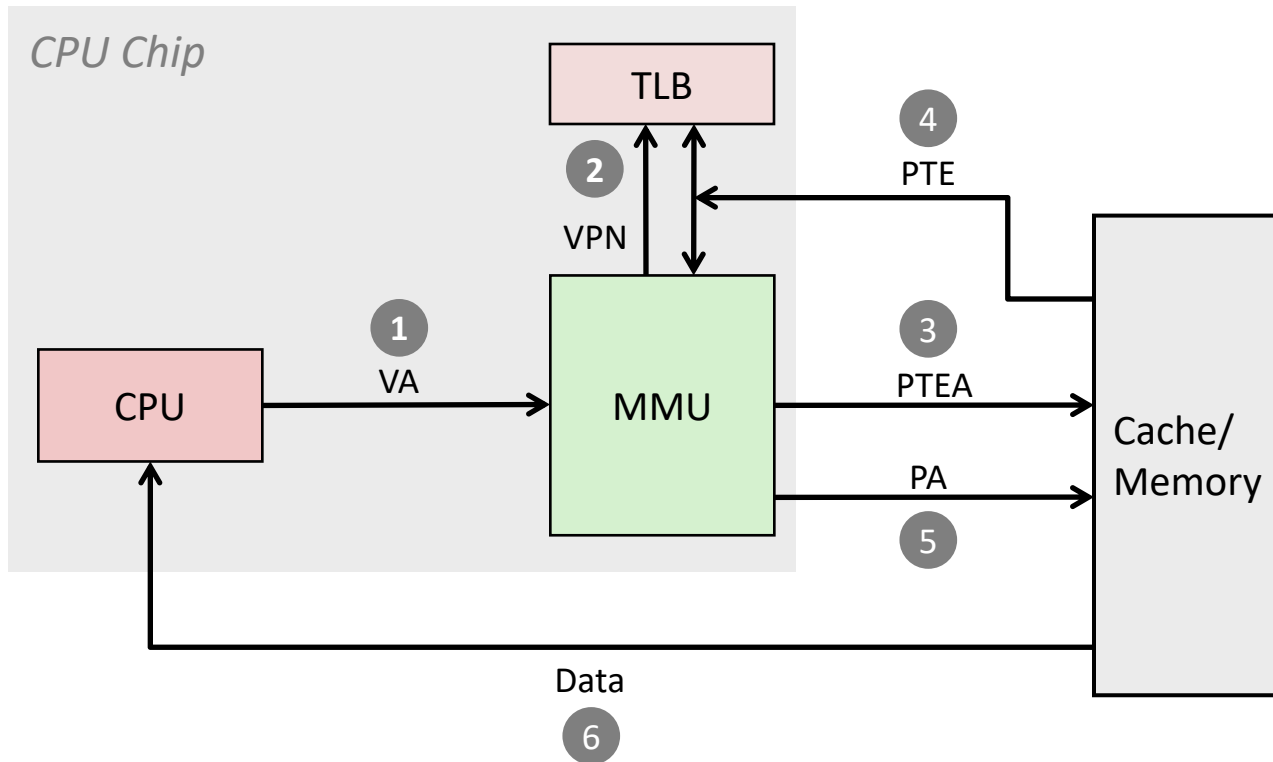Set T-1

| v | tag | PTE | | v | tag | PTE |

# TLB Hit

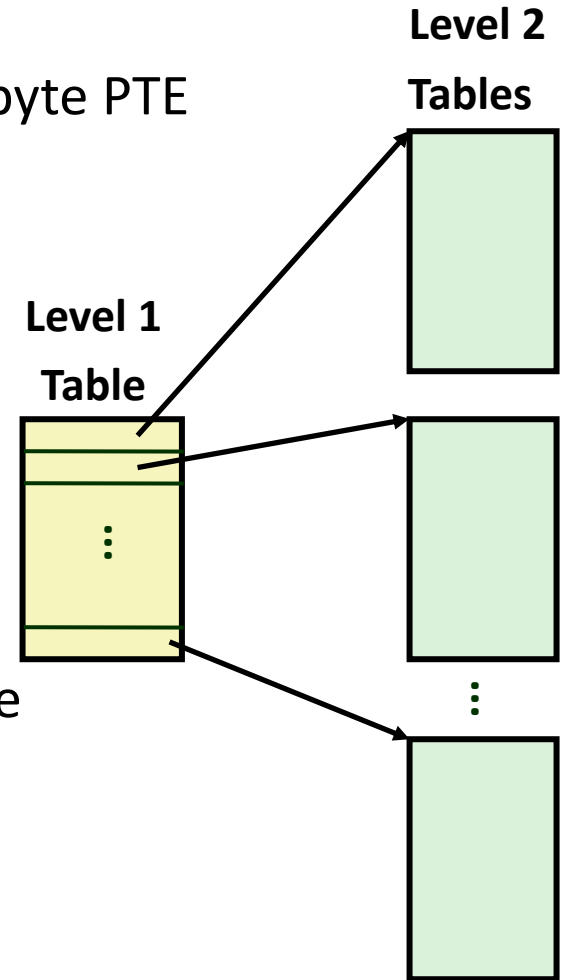

**A TLB hit eliminates a memory access**

# TLB Miss



**A TLB miss incurs an additional memory access (the PTE)**
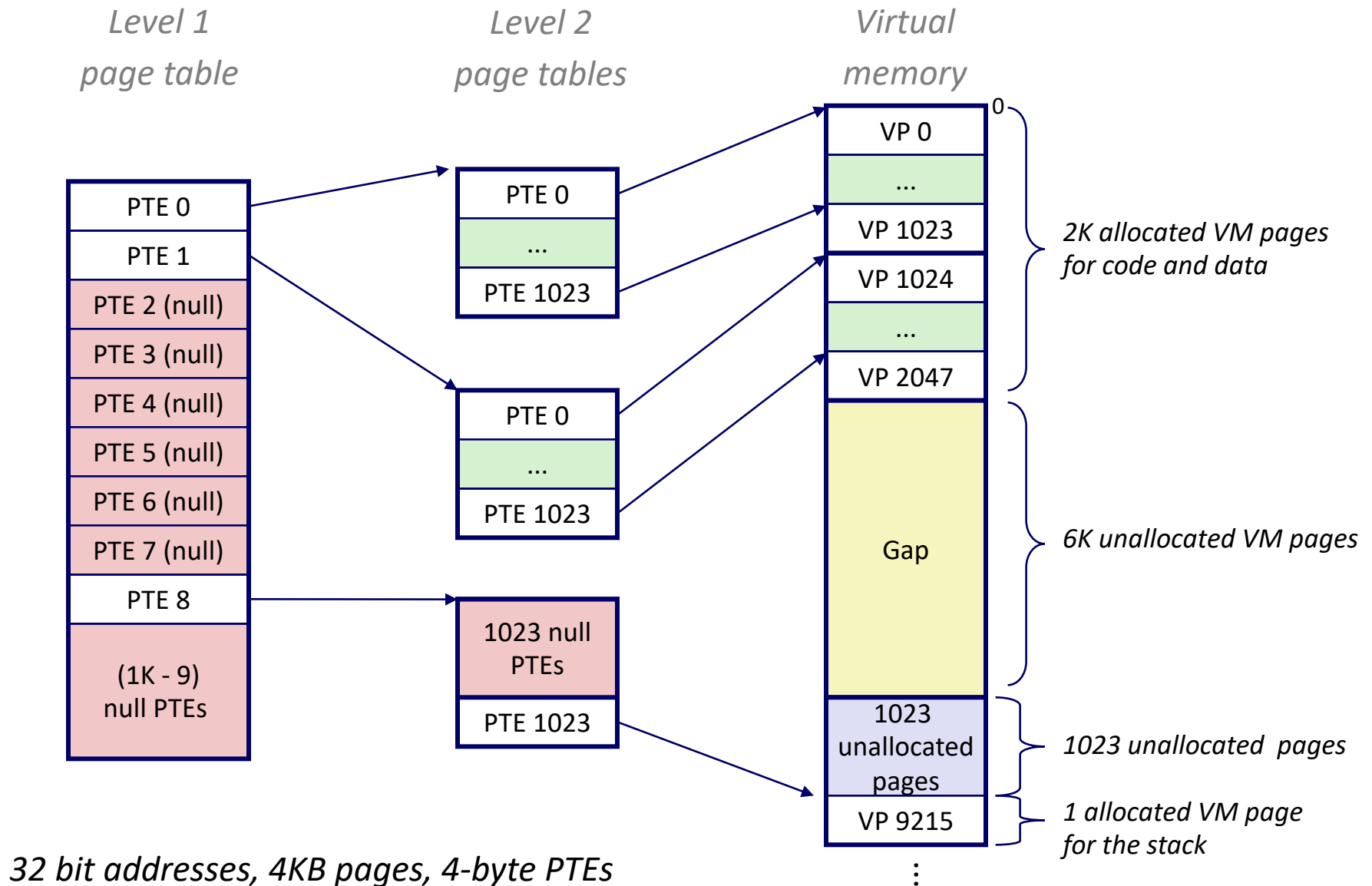Fortunately, TLB misses are rare.

# Multi-Level Page Tables

- Suppose:
  - 4KB ($2^{12}$) page size, 48-bit address space, 8-byte PTE

- Problem:
  - Would need a 512 GB page table!
    - $2^{48} * 2^{-12} * 2^{3} = 2^{39}$ bytes

- Common solution: Multi-level page table

- Example: 2-level page table
  - Level 1 table: each PTE points to a page table (always memory resident)
  - Level 2 table: each PTE points to a page (paged in and out like any other data)

**Level 2**

**Tables**

**Level 1**

**Table**

# A Two-Level Page Table Hierarchy

Level 1
page table

Level 2
page tables

Virtual
memory

| | |
|---|---|
| PTE 0 | |
| PTE 1 | |
| PTE 2 (null) | |
| PTE 3 (null) | |
| PTE 4 (null) | |
| PTE 5 (null) | |
| PTE 6 (null) | |
| PTE 7 (null) | |
| PTE 8 | |
| (1K - 9) null PTEs | |

| |
|---|
| PTE 0 |
| ... |
| PTE 1023 |

| |
|---|
| PTE 0 |
| ... |
| PTE 1023 |

| |
|---|
| 1023 null PTEs |
| PTE 1023 |

0

| |
|---|
| VP 0 |
| ... |
| VP 1023 |
| VP 1024 |
| ... |
| VP 2047 |

2K allocated VM pages
for code and data

| |
|---|
| Gap |

6K unallocated VM pages

| |
|---|
| 1023 unallocated pages |

1023 unallocated  pages

| |
|---|
| VP 9215 |

1 allocated VM page
for the stack

*32 bit addresses, 4KB pages, 4-byte PTEs*

⋮

# Translating with a k-level Page Table

# Summary

- Programmer's view of virtual memory
  - Each process has its own private linear address space
  - Cannot be corrupted by other processes

- System view of virtual memory
  - Uses memory efficiently by caching virtual memory pages
    - Efficient only because of locality
  - Simplifies memory management and programming
  - Simplifies protection by providing a convenient interpositioning point to check permissions