



# Virtual Memory: Concepts

18-213/18-613: Introduction to Computer Systems  
12<sup>th</sup> Lecture, October 1, 2024

# Announcements

- **“Low-stakes take-home midterm”**  
**goes out on Monday evening (after small groups finish)**
  - **80 minutes self-timed**. Covers through virtual memort
  - Questions similar to homeworks, but **only one attempt**.
  - Tests what you’ve learned, as in a real midterm (and as in the Final).
  - Low-stakes: Only 4% of grade (could even be “half dropped”).

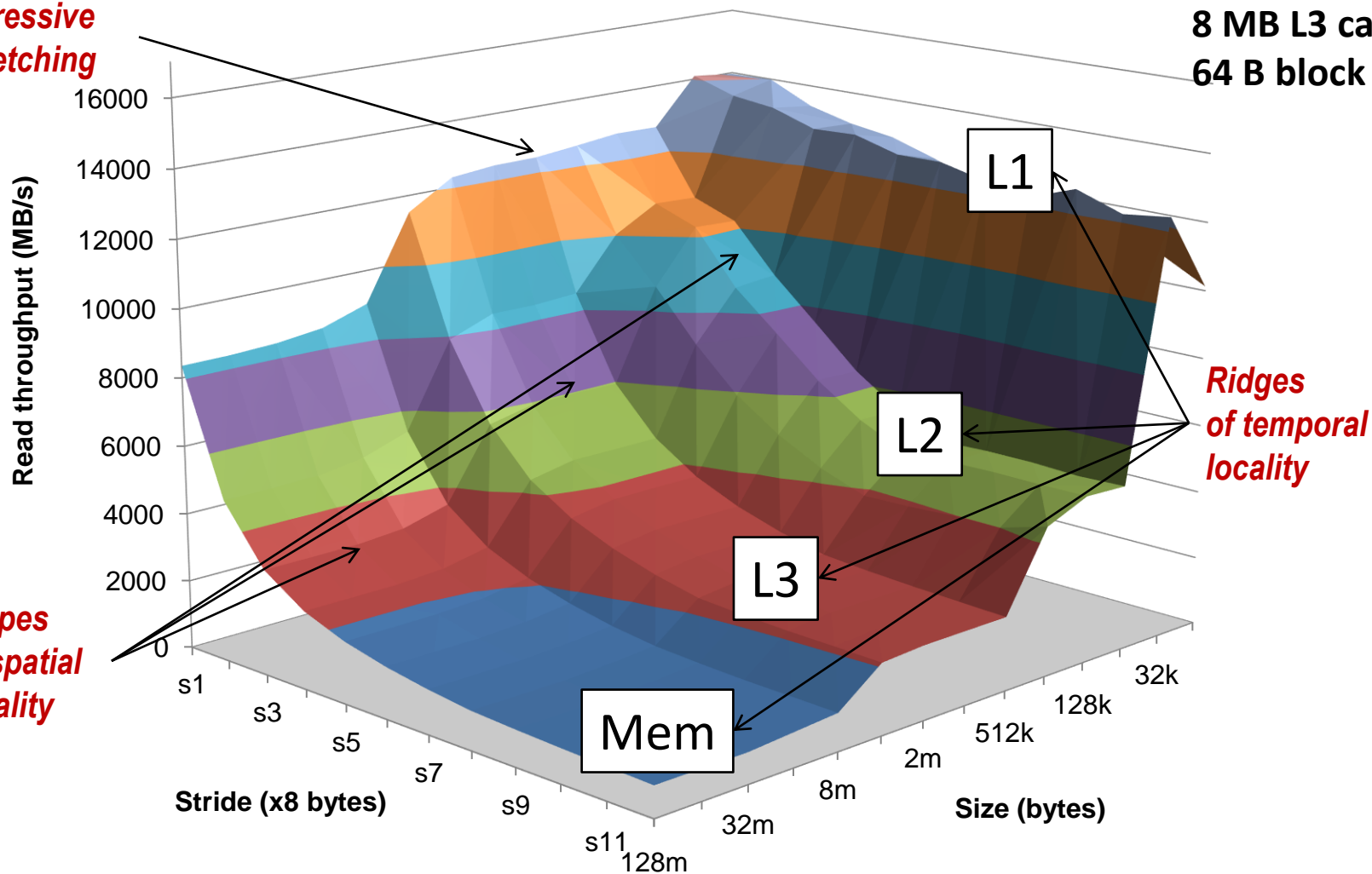
# Caching Wrap-Up

- Quick review
- Miss-Rate Analysis
- Blocked Operations

# 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

*Aggressive prefetching*



# Matrix Multiplication Example

## ■ Description:

- Multiply  $N \times N$  matrices
- Matrix elements are doubles (8 bytes)
- $O(N^3)$  total operations
- $N$  reads per source element
- $N$  values summed per destination
  - but may be able to hold 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;
    }
}
```

*Variable sum held in register*

*matmult/mm.c*

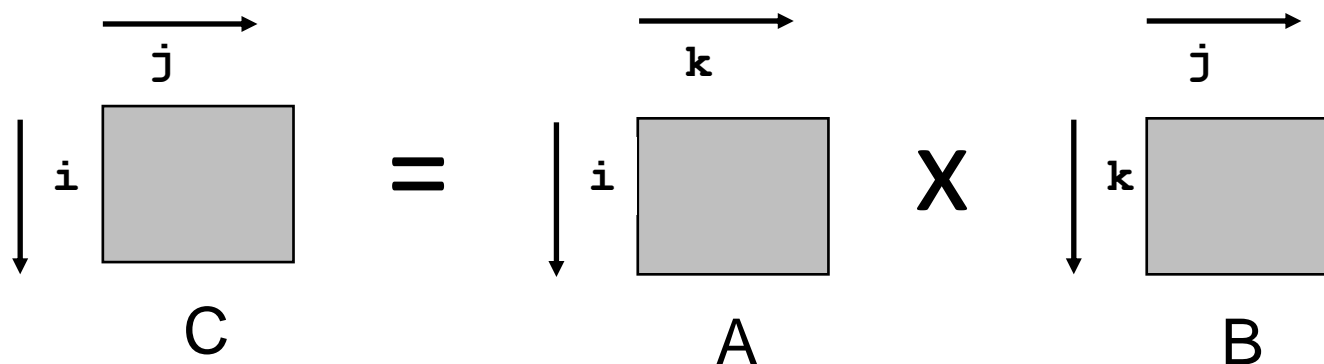
# Miss Rate Analysis for Matrix Multiply

## ■ Assume:

- Block size = 64B (big enough for eight 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



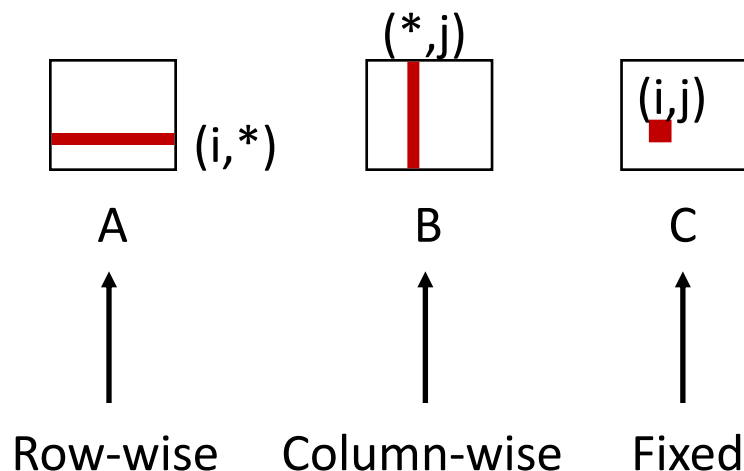
# Matrix Multiplication (ijk)

```

/* 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:



Miss rate for inner loop iterations:

<u>A</u>	<u>B</u>	<u>C</u>
0.125	1.0	0.0

**Avg misses/iter = 1.125**

**Block size = 64B (eight doubles)**



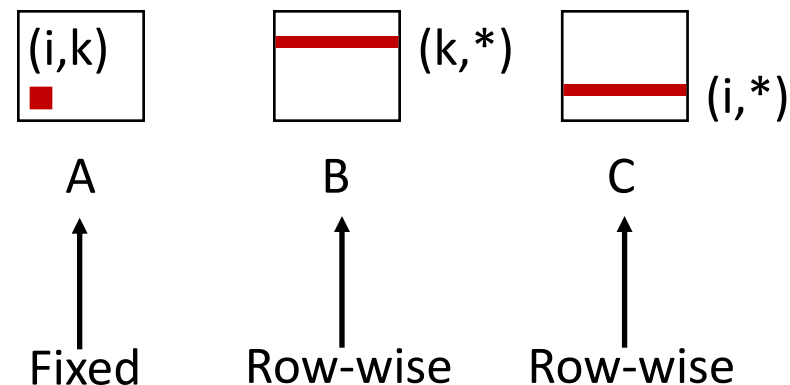
# Matrix Multiplication (kij)

```

/* 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:



Miss rate for inner loop iterations:

<u>A</u>	<u>B</u>	<u>C</u>
0.0	0.125	0.125

**Avg misses/iter = 0.25**

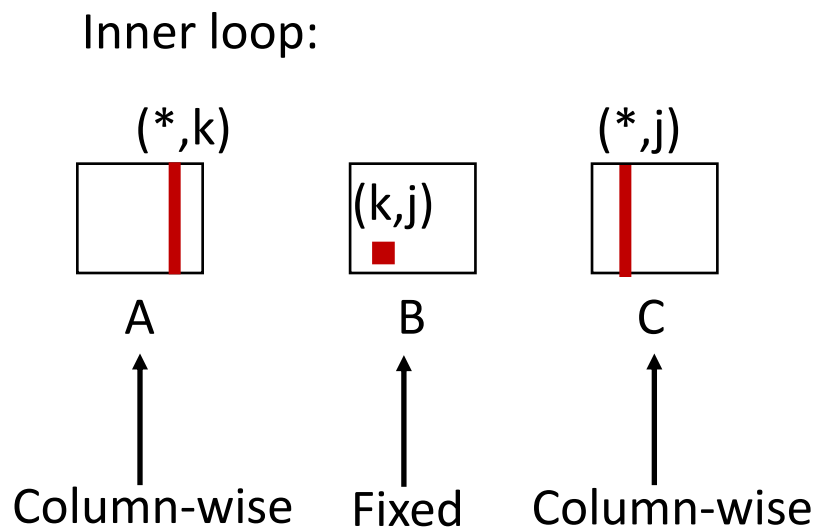
**Block size = 64B (eight doubles)**

# 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

```



Miss rate for inner loop iterations:

<u>A</u>	<u>B</u>	<u>C</u>
1.0	0.0	1.0

**Avg misses/iter = 2.0**

**Block size = 64B (eight doubles)**

# 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
- avg misses/iter = **1.125**

```

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
- avg misses/iter = **0.25**

```

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;
  }
}

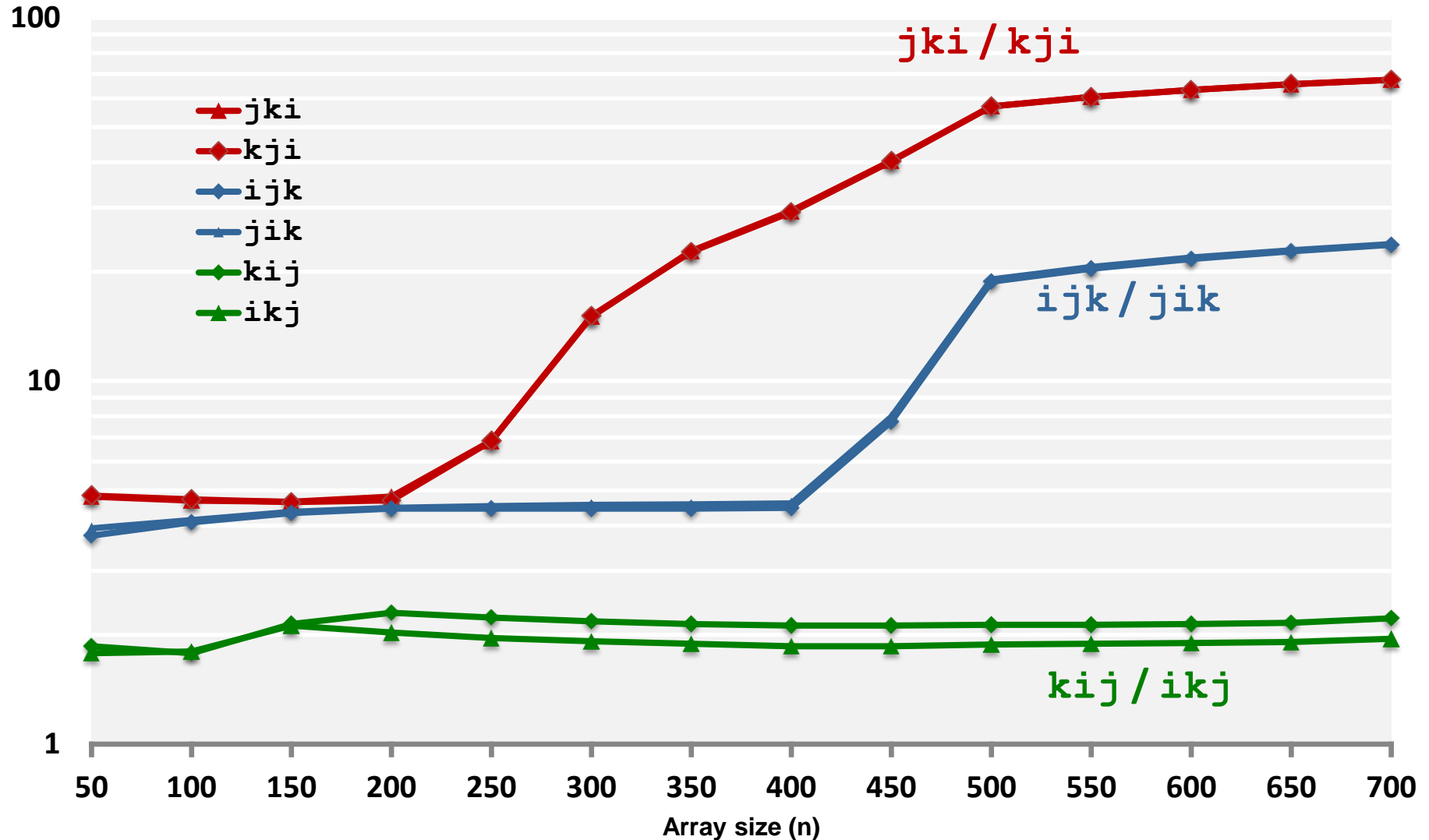
```

**jki (& kji):**

- 2 loads, 1 store
- avg misses/iter = **2.0**

# Core i7 Matrix Multiply Performance

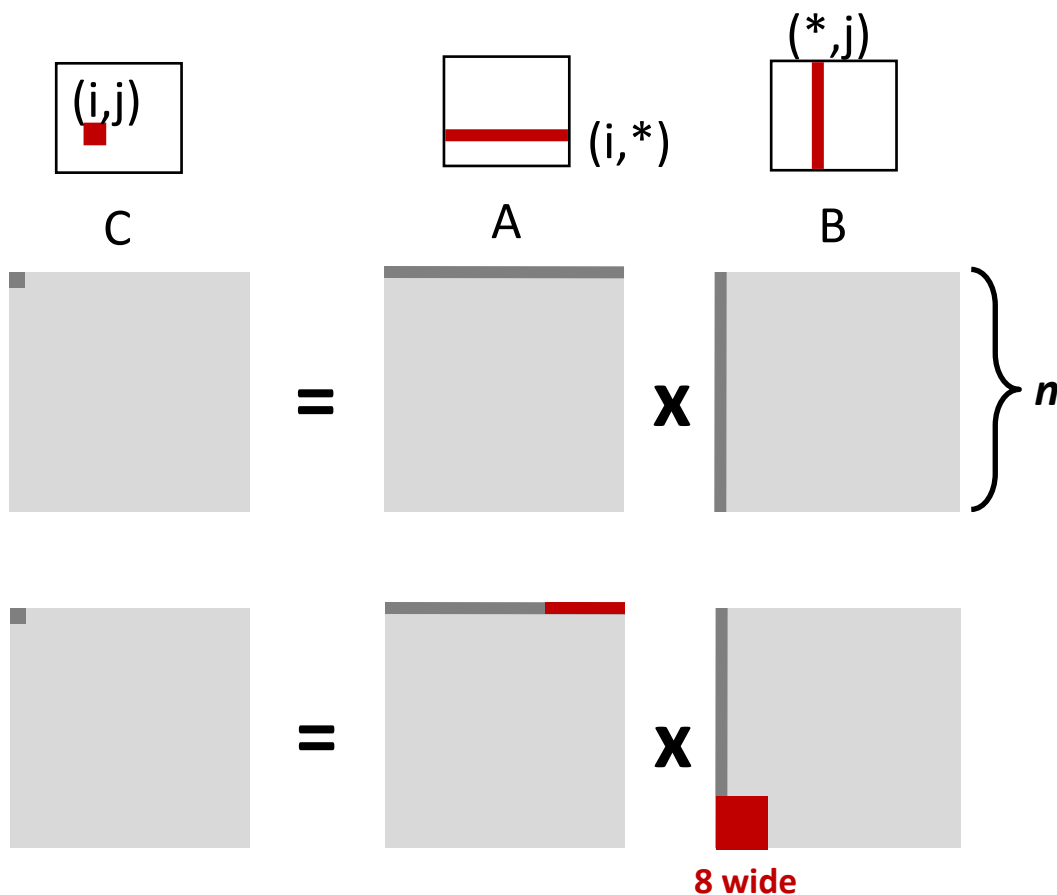
Cycles per inner loop iteration



# Matrix Multiplication Cache Miss Analysis

## ■ Assume:

- Matrix elements are doubles. Cache line = 8 doubles
- Cache size  $C \ll n$  (much smaller than  $n$ )



## ■ First iteration (ijk):

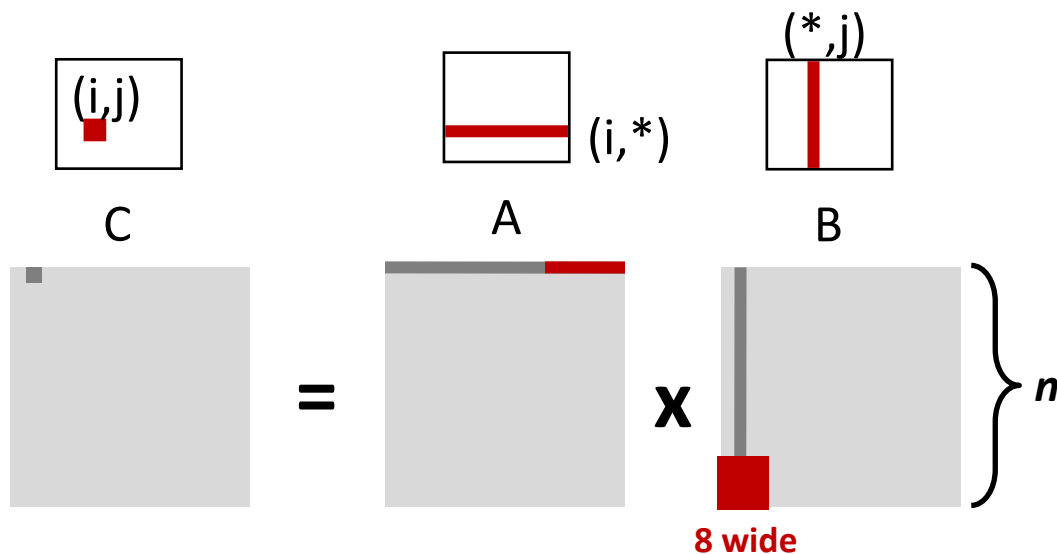
- $n/8 + n = 9n/8$  misses

- Afterwards **in cache:**  
(schematic)

# Cache Miss Analysis (cont)

## ■ Assume:

- Matrix elements are doubles. Cache line = 8 doubles
- Cache size  $C \ll n$  (much smaller than  $n$ )



## ■ Second iteration:

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

## ■ Total misses:

- $(9n/8) n^2 = (9/8) n^3$

# Blocked Matrix Multiplication

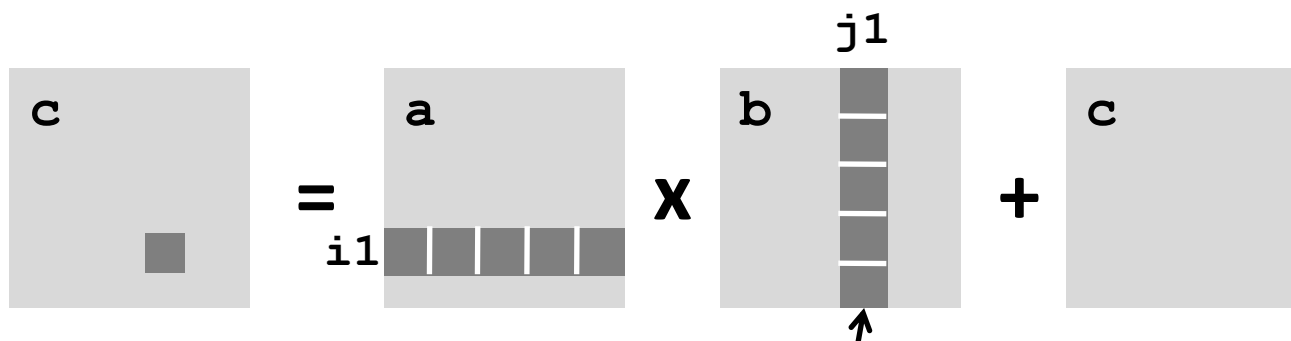
```

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+=L)
        for (j = 0; j < n; j+=L)
            for (k = 0; k < n; k+=L)
                /* L x L mini matrix multiplications */
                for (i1 = i; i1 < i+L; i1++)
                    for (j1 = j; j1 < j+L; j1++)
                        for (k1 = k; k1 < k+L; k1++)
                            c[i1*n+j1] += a[i1*n + k1]*b[k1*n + j1];
}

```

*matmult/bmm.c*



Block size  $L \times L$

# Cache Miss Analysis

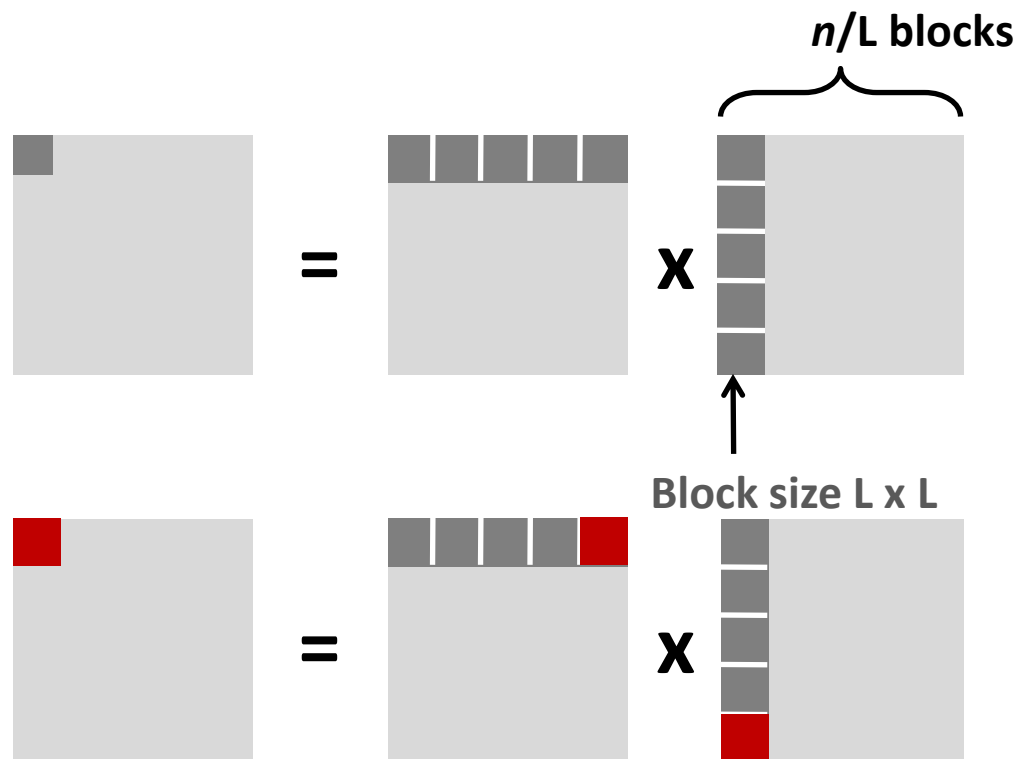
## ■ Assume:

- Cache line = 8 doubles. Blocking size  $L \geq 8$
- Cache size  $C \ll n$  (much smaller than  $n$ )
- Three blocks  $\blacksquare$  fit into cache:  $3L^2 < C$

## ■ First (block) iteration:

- Misses per block:  $L^2/8$
- Blocks per Iteration:  $2n/L$   
(omitting matrix  $c$ )
- Misses per Iteration:  
 $2n/L \times L^2/8 = nL/4$

- Afterwards in cache  
(schematic)





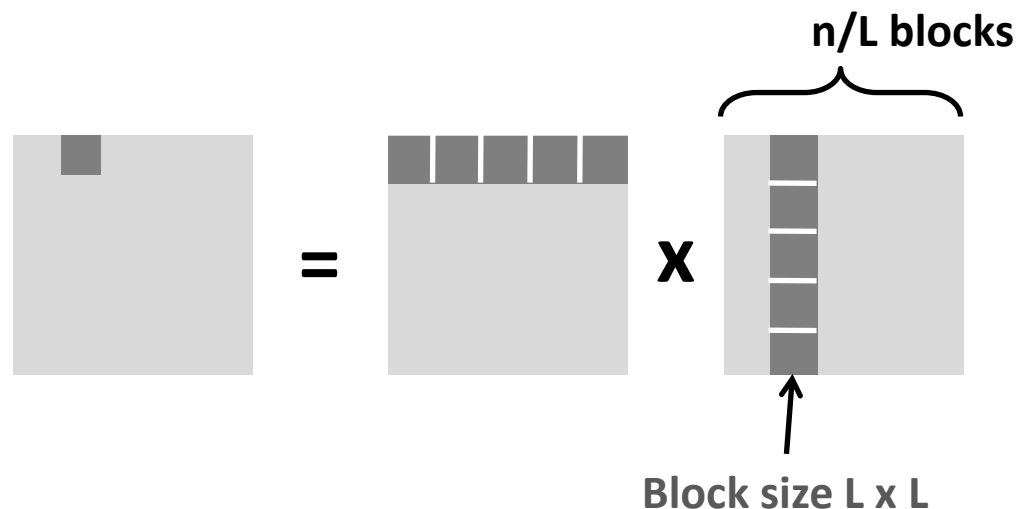
# Cache Miss Analysis

## ■ Assume:

- Cache line = 8 doubles. Blocking size  $L \geq 8$
- Cache size  $C \ll n$  (much smaller than  $n$ )
- Three blocks  $\blacksquare$  fit into cache:  $3L^2 < C$

## ■ Second (block) iteration:

- Same misses as first iteration
- $2n/L \times L^2/8 = nL/4$



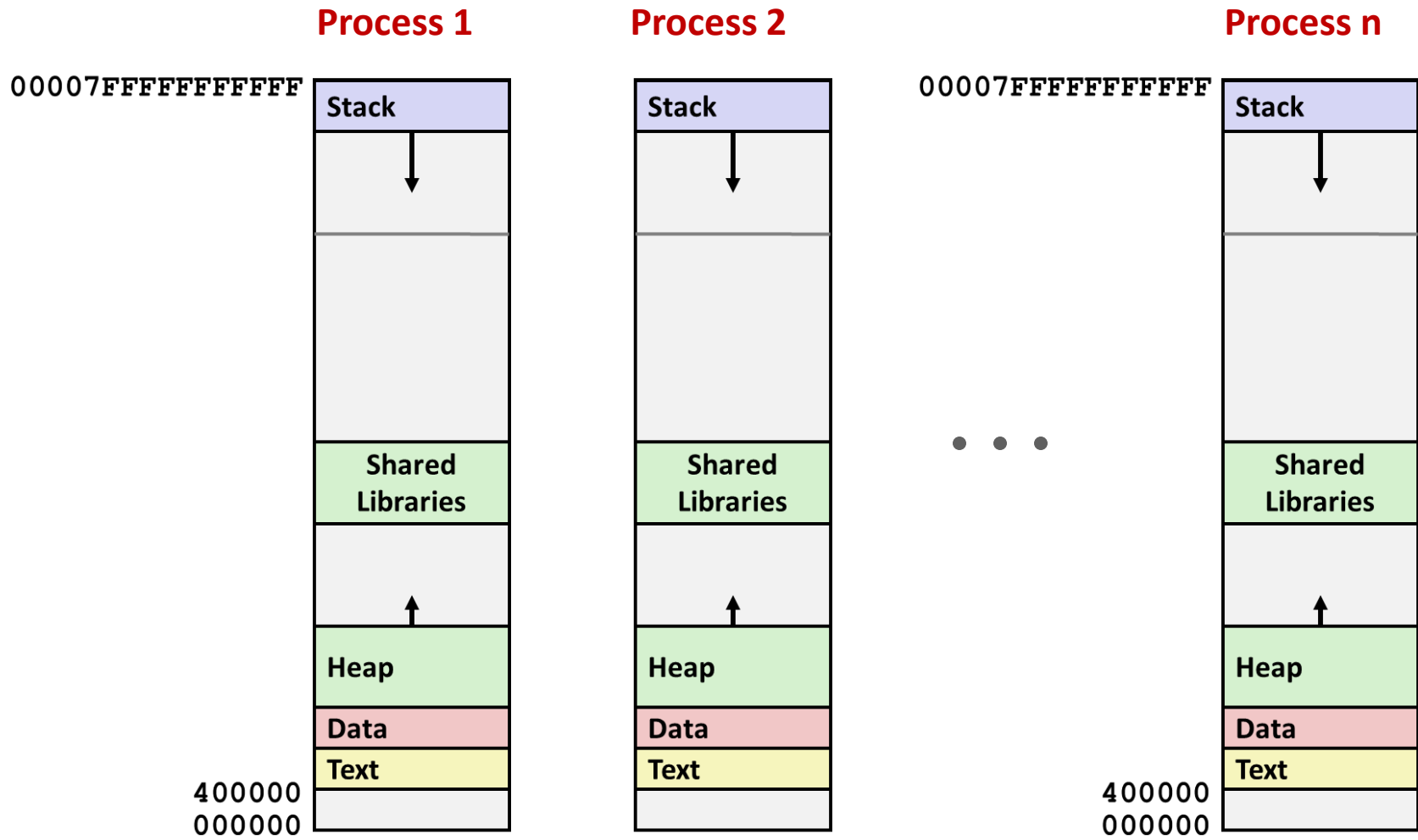
## ■ Total misses:

- $nL/4$  misses per iteration  $\times (n/L)^2$  iterations =  $n^3/(4L)$  misses

# Blocking Summary

- No blocking (ijk) :  $(9/8) n^3$  misses
- Blocking:  $(1/(4L)) n^3$  misses
  
- Use largest block size  $L$ , such that  $L$  satisfies  $3L^2 < C$ 
  - Fit three blocks in cache! Two input, one output.
  
- 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

# Hmmm, How Does This Work?!



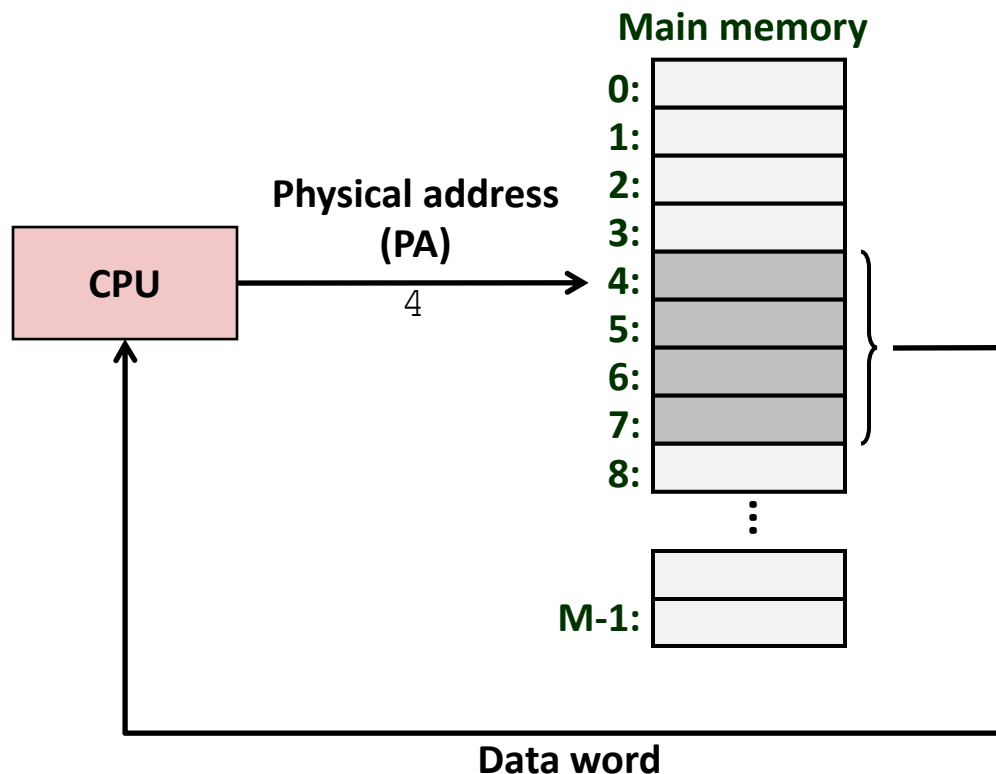
***Solution: Virtual Memory (today and next lecture)***

# Virtual Memory

- **Address spaces** CSAPP 9.1-9.2
- VM as a tool for caching CSAPP 9.3
- VM as a tool for memory management CSAPP 9.4
- VM as a tool for memory protection CSAPP 9.5
- Address translation CSAPP 9.6

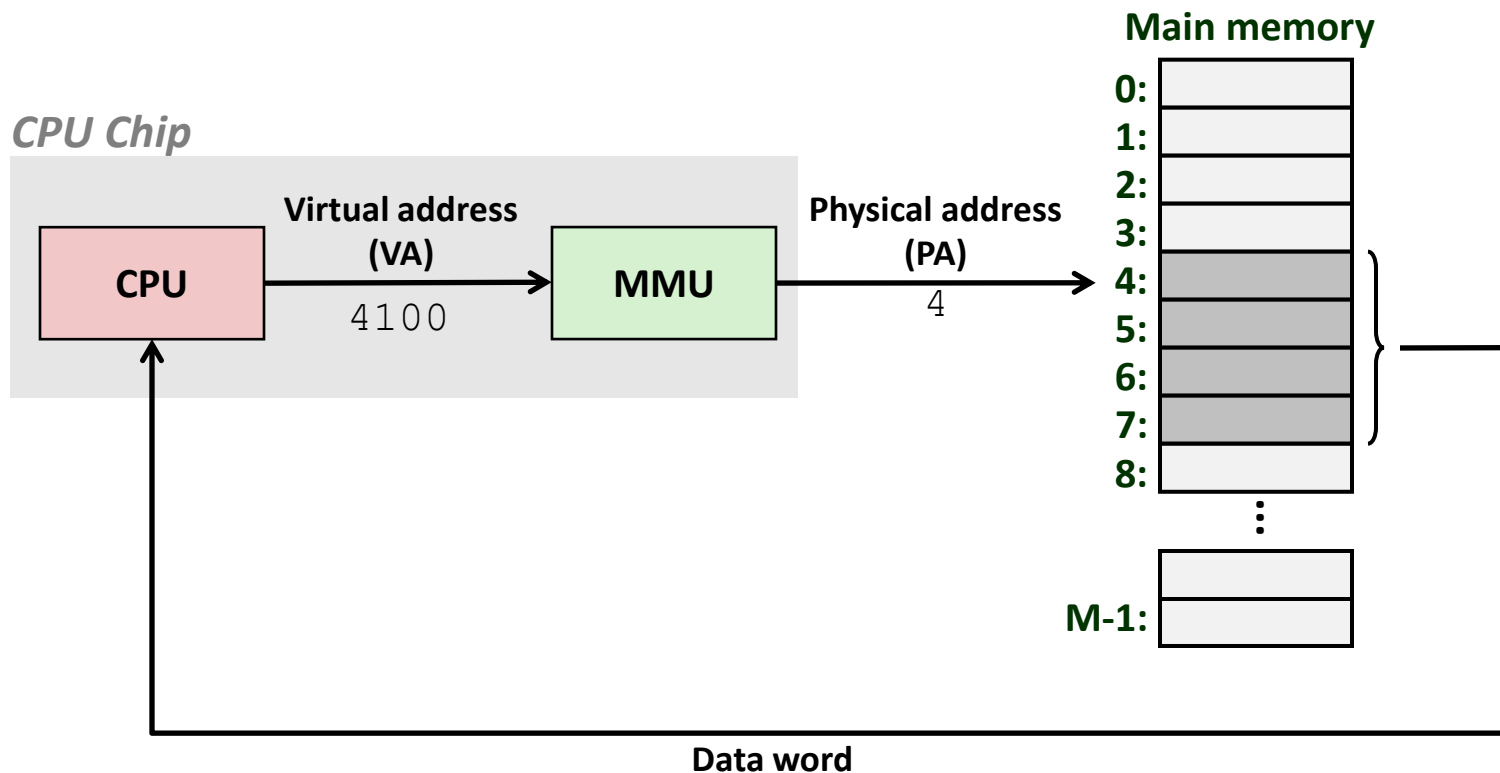
# Blank Slide for Intro Sketching

# A System Using Physical Addressing



- Used in “simple” systems like embedded microcontrollers in devices like elevators and digital picture frames

# A System Using Virtual Addressing



- 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)?

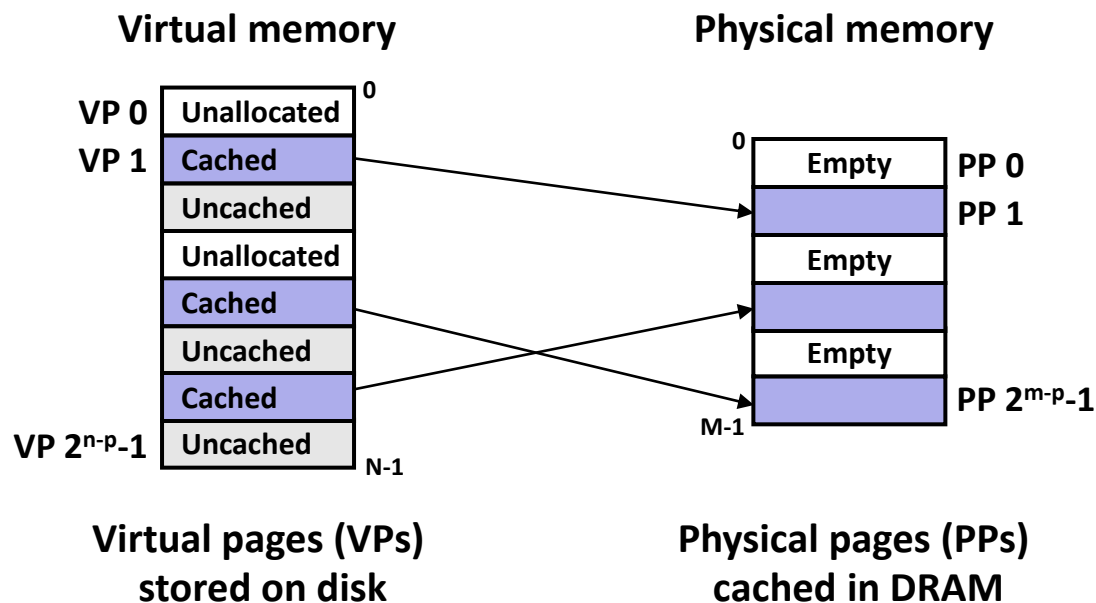
- **Uses main memory efficiently**
  - Use DRAM as a cache for parts of a virtual address space
- **Simplifies memory management**
  - Each process gets the same uniform linear address space
- **Isolates address spaces**
  - One process can't interfere with another's memory
  - User program cannot access privileged kernel information and code

# Today

- Address spaces
- **VM as a tool for caching**
- VM as a tool for memory management
- VM as a tool for memory protection
- Address translation

# 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)



# DRAM Cache Organization

## ■ DRAM cache organization driven by the enormous miss penalty

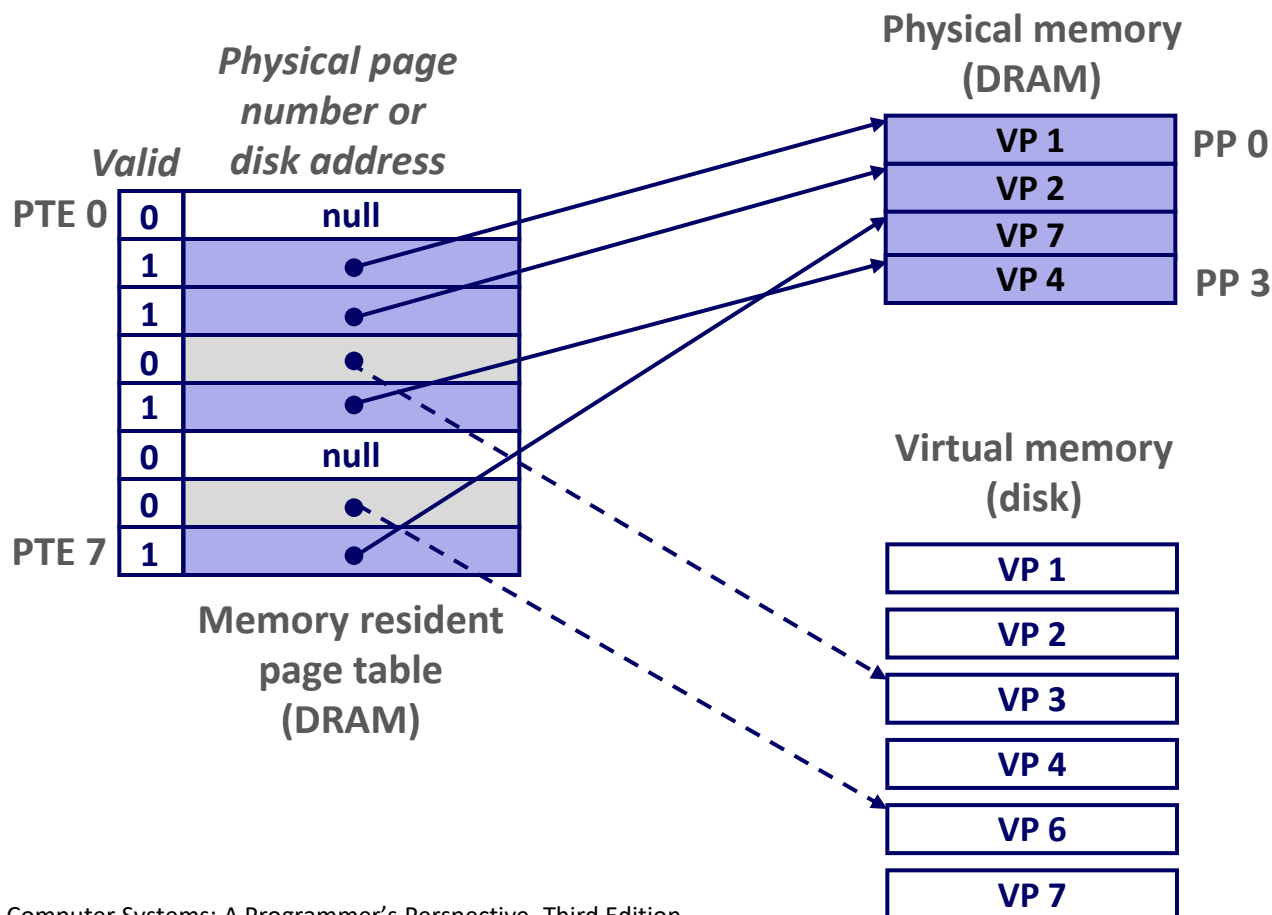
- DRAM is about **10x** slower than SRAM
- Disk is about **10,000x** slower than DRAM
- Time to load block from disk > 1ms (> 1 million clock cycles)
  - CPU can do a lot of computation during that time

## ■ Consequences

- Large page (block) size: typically 4 KB
  - Linux “huge pages” are 2 MB (default) to 1 GB
- Fully associative. *Why?*
  - Any VP can be placed in any PP
  - Requires a “large” mapping function – different from cache memories
- Highly sophisticated, expensive replacement algorithms. *Why?*
  - Too complicated and open-ended to be implemented in hardware
- Write-back rather than write-through. *Why?*

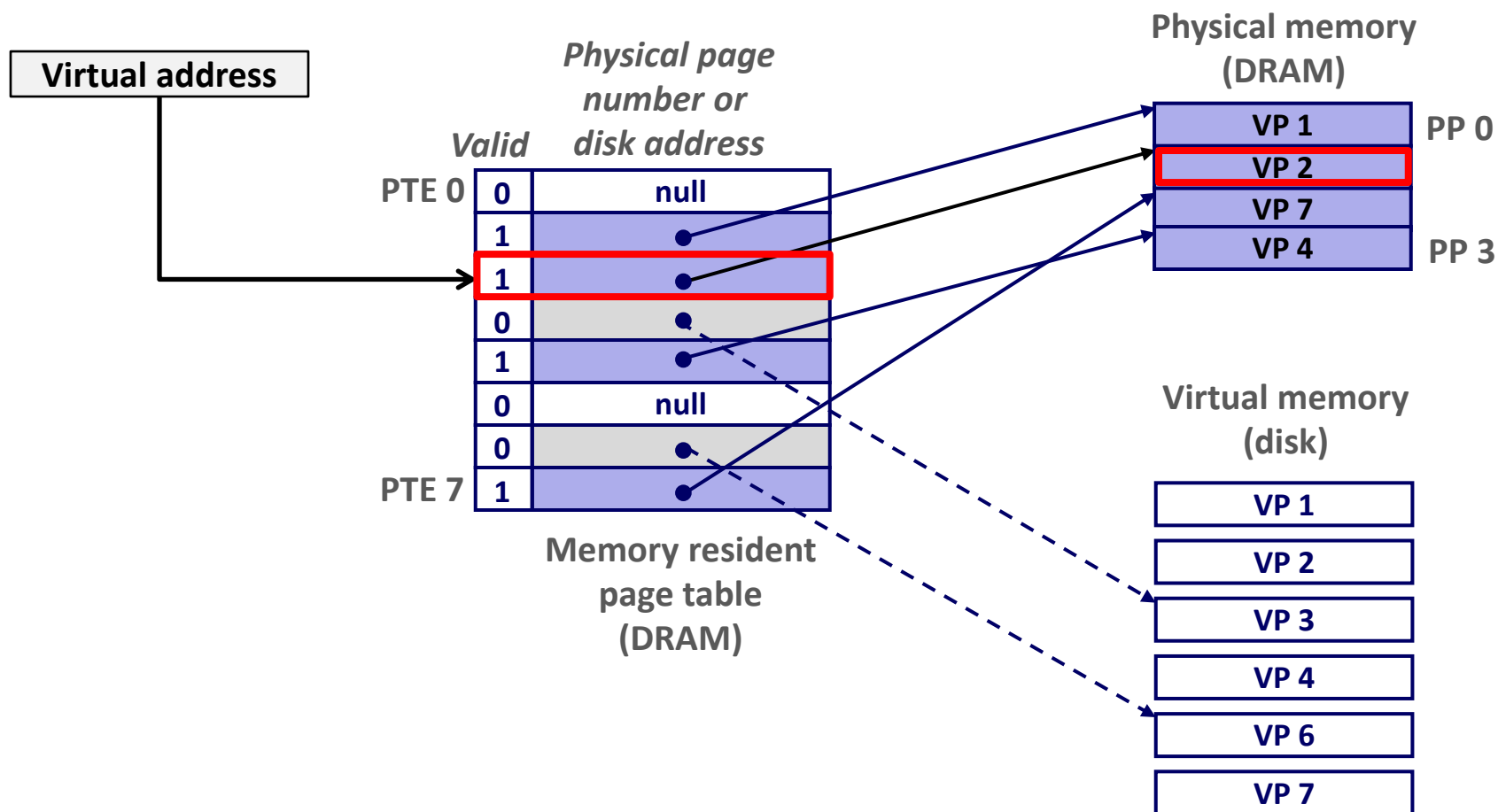
# Enabling Data Structure: Page Table

- 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



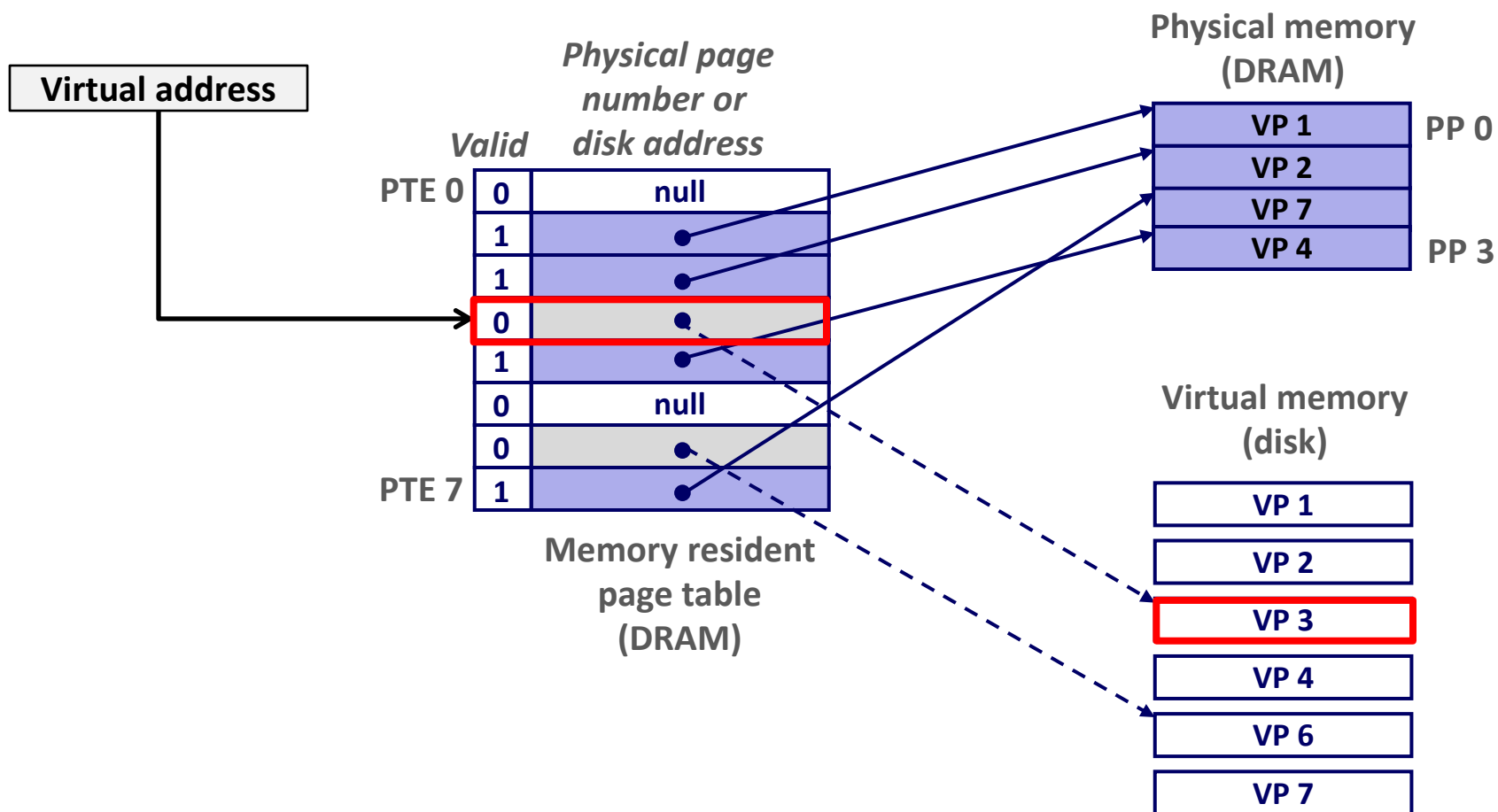
# Page Hit

- **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)



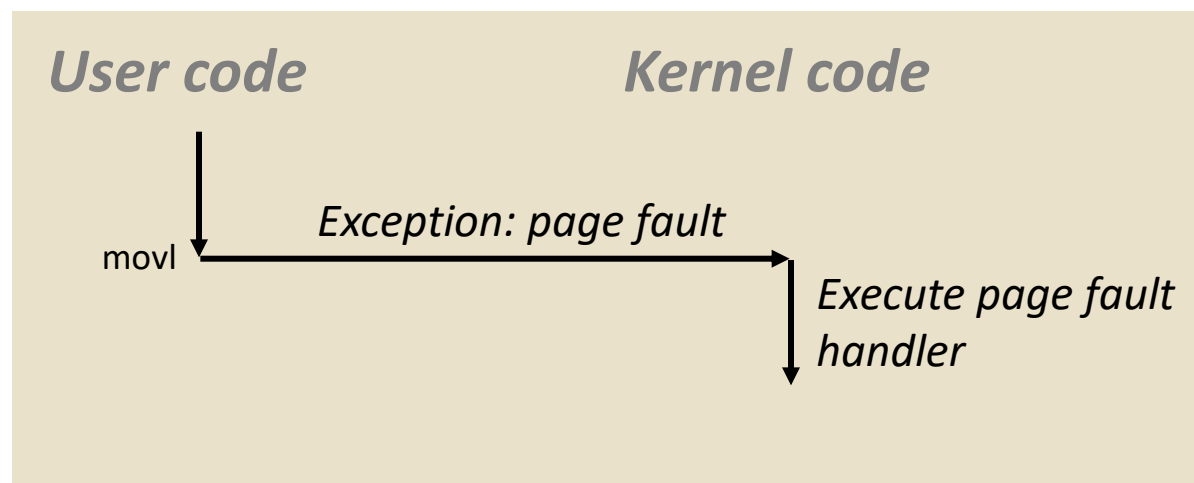
# Triggering a Page Fault

- User writes to memory location

```
80483b7:      c7 05 10 9d 04 08 0d  movl   $0xd,0x8049d10
```

- That portion (page) of user's memory is currently on disk
- MMU triggers page fault exception
  - (More details in later lecture)
  - Raise privilege level to supervisor mode
  - Causes procedure call to software page fault handler

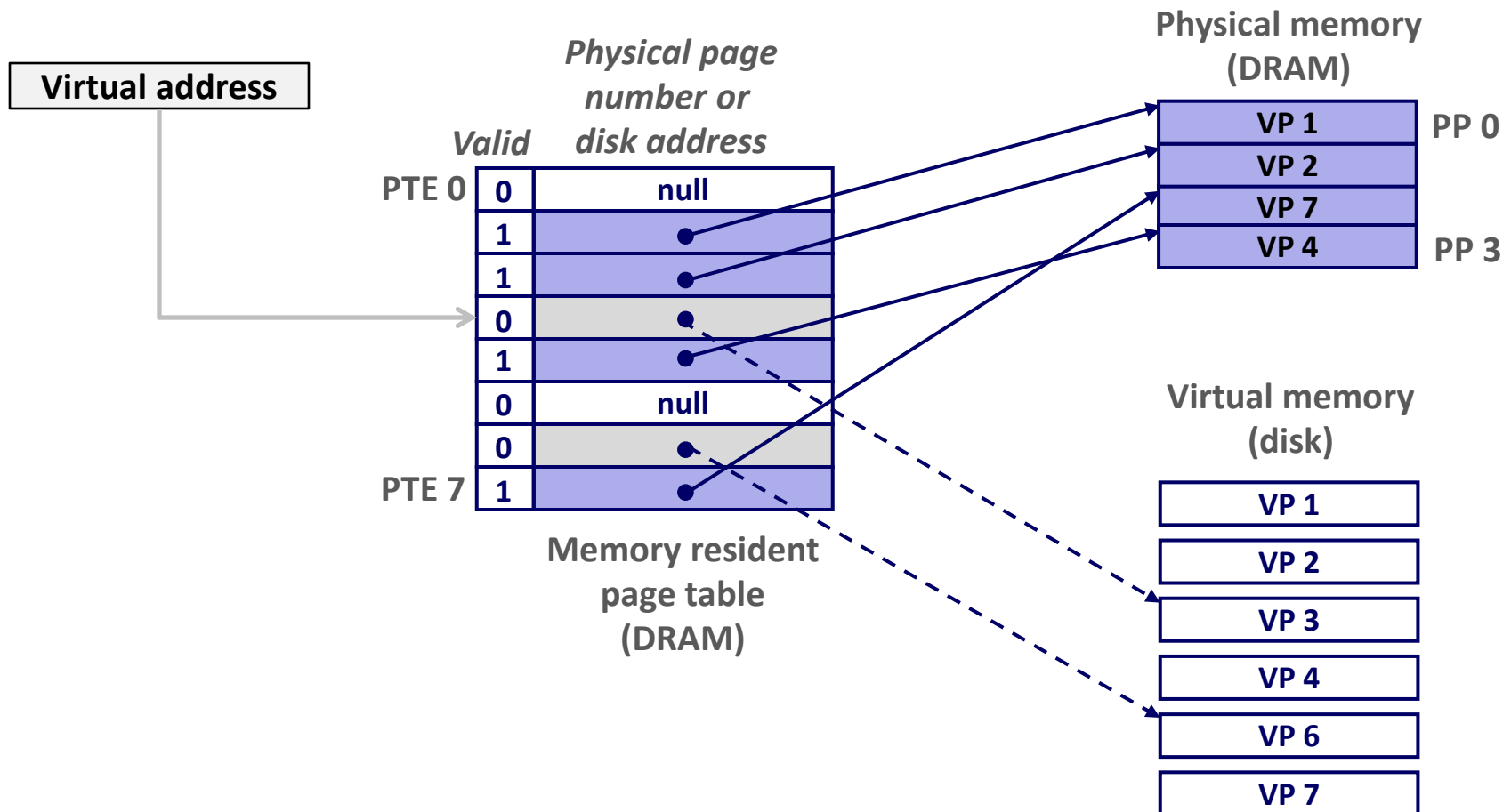
```
int a[1000];
main ()
{
    a[500] = 13;
}
```





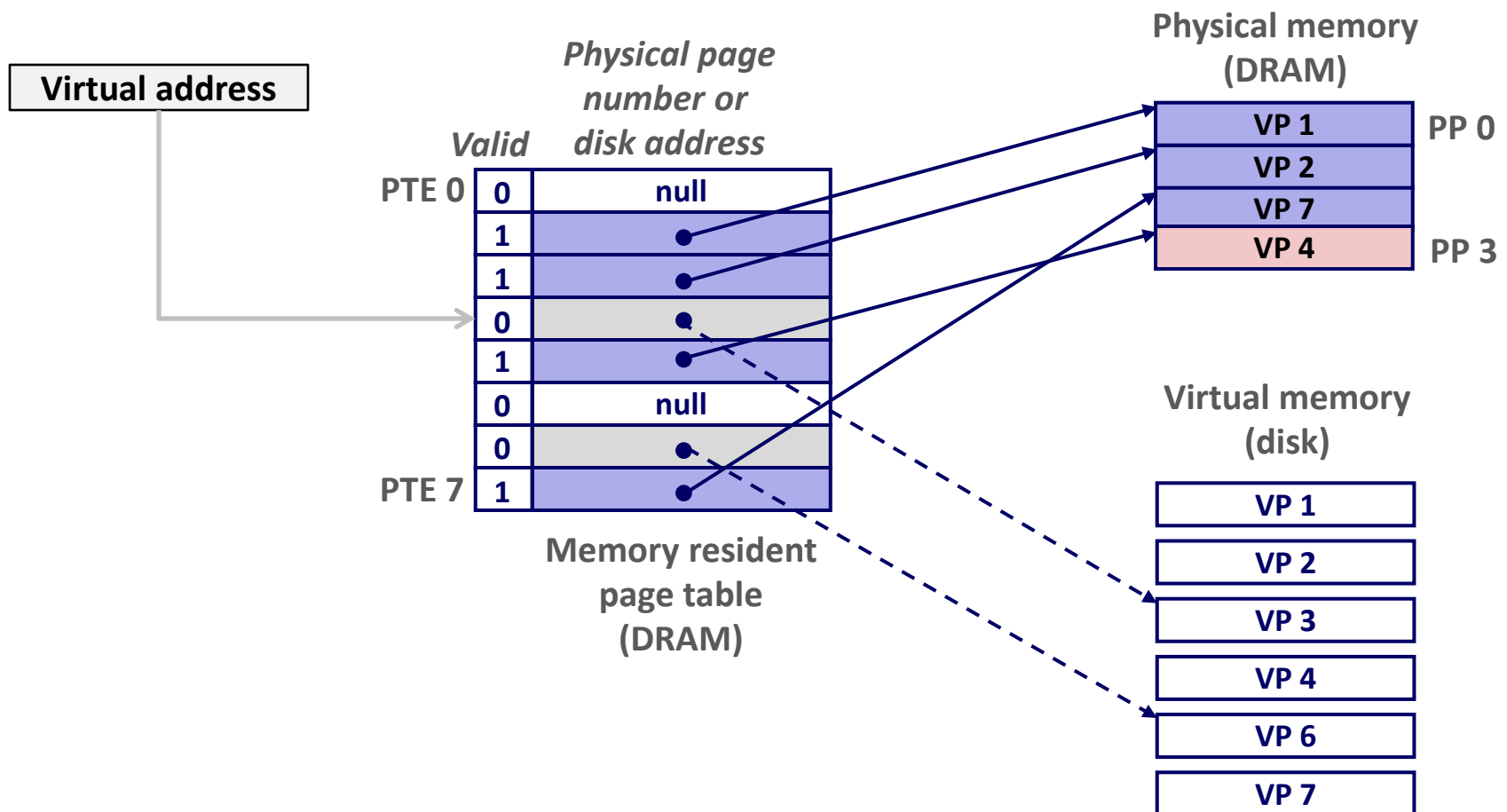
# Handling Page Fault

- Page miss causes page fault (an exception)



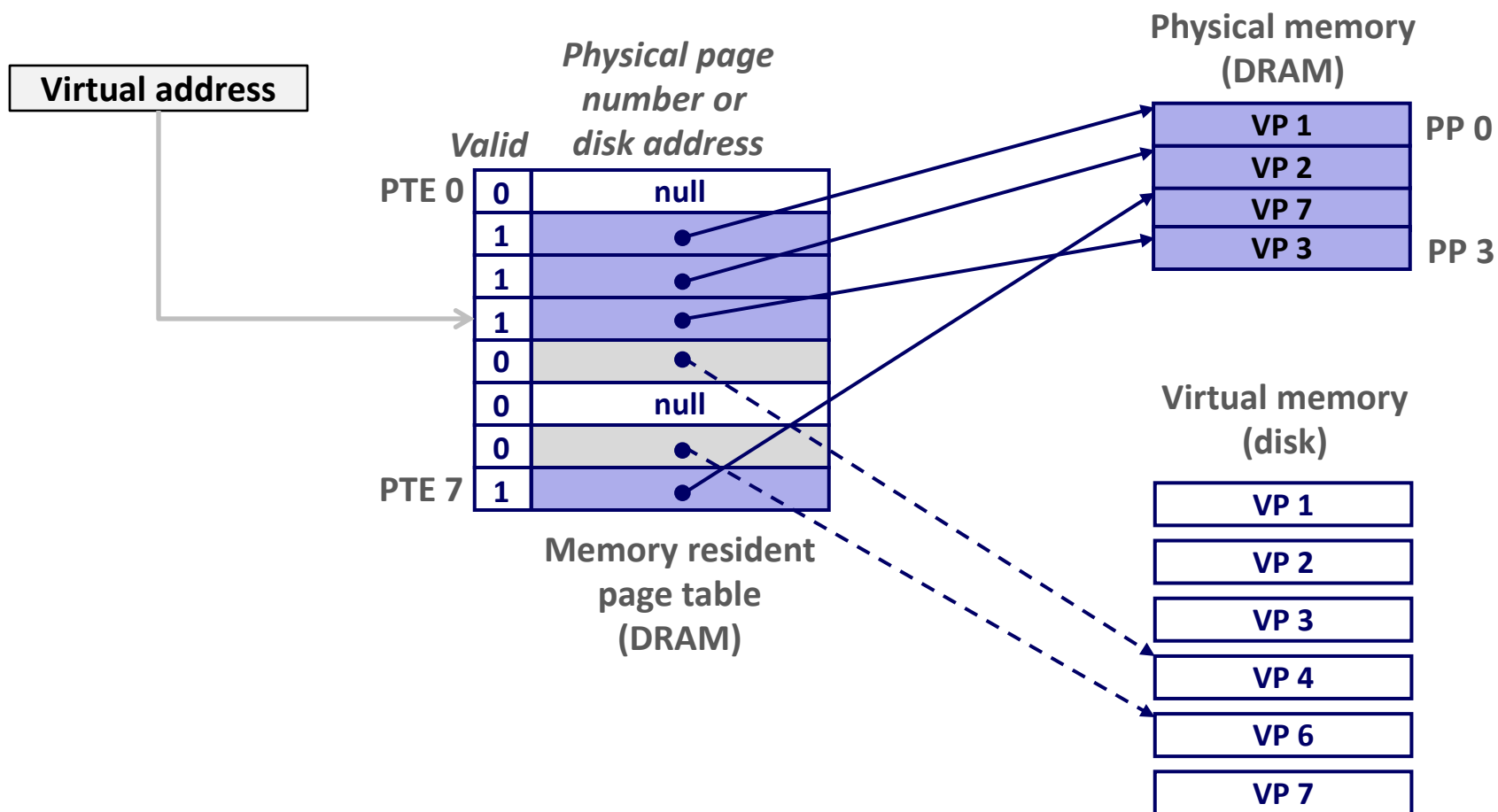
# Handling Page Fault

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



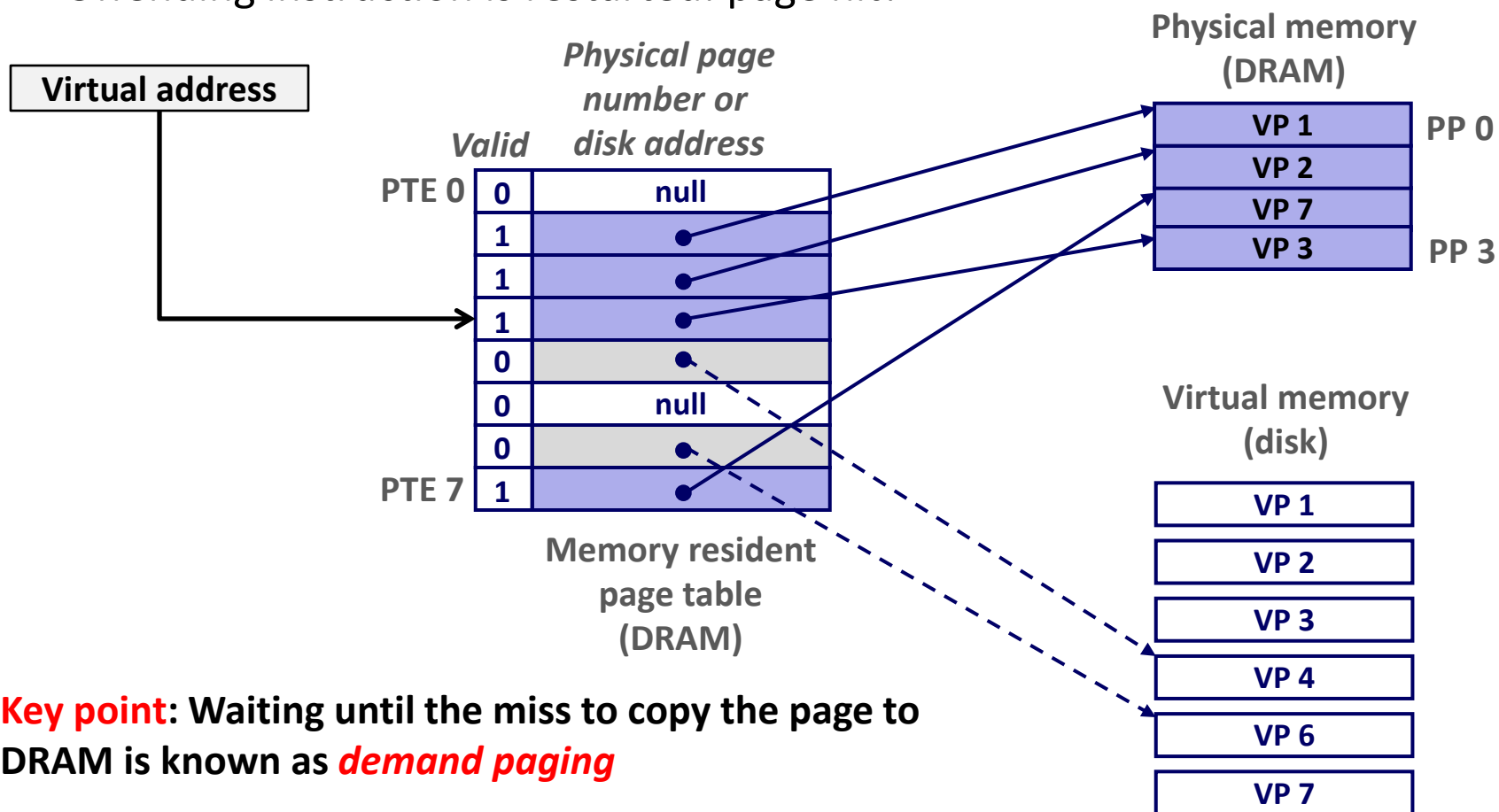
# Handling Page Fault

- 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!



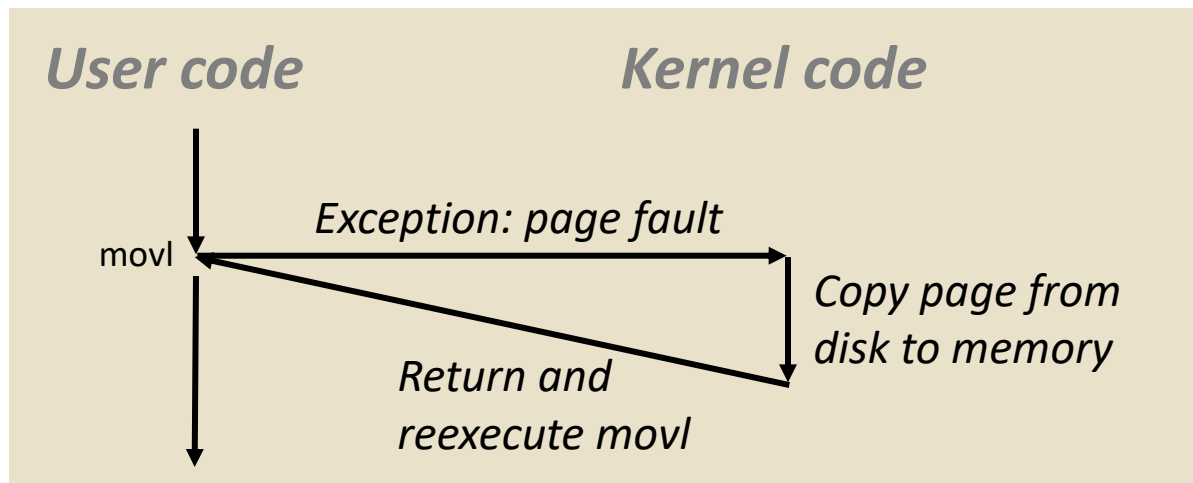
**Key point:** Waiting until the miss to copy the page to DRAM is known as **demand paging**

# Completing page fault

- Page fault handler executes return from interrupt (**iret**) instruction
  - Like **ret** instruction, but also restores privilege level
  - Return to instruction that caused fault
  - But, this time there is no page fault

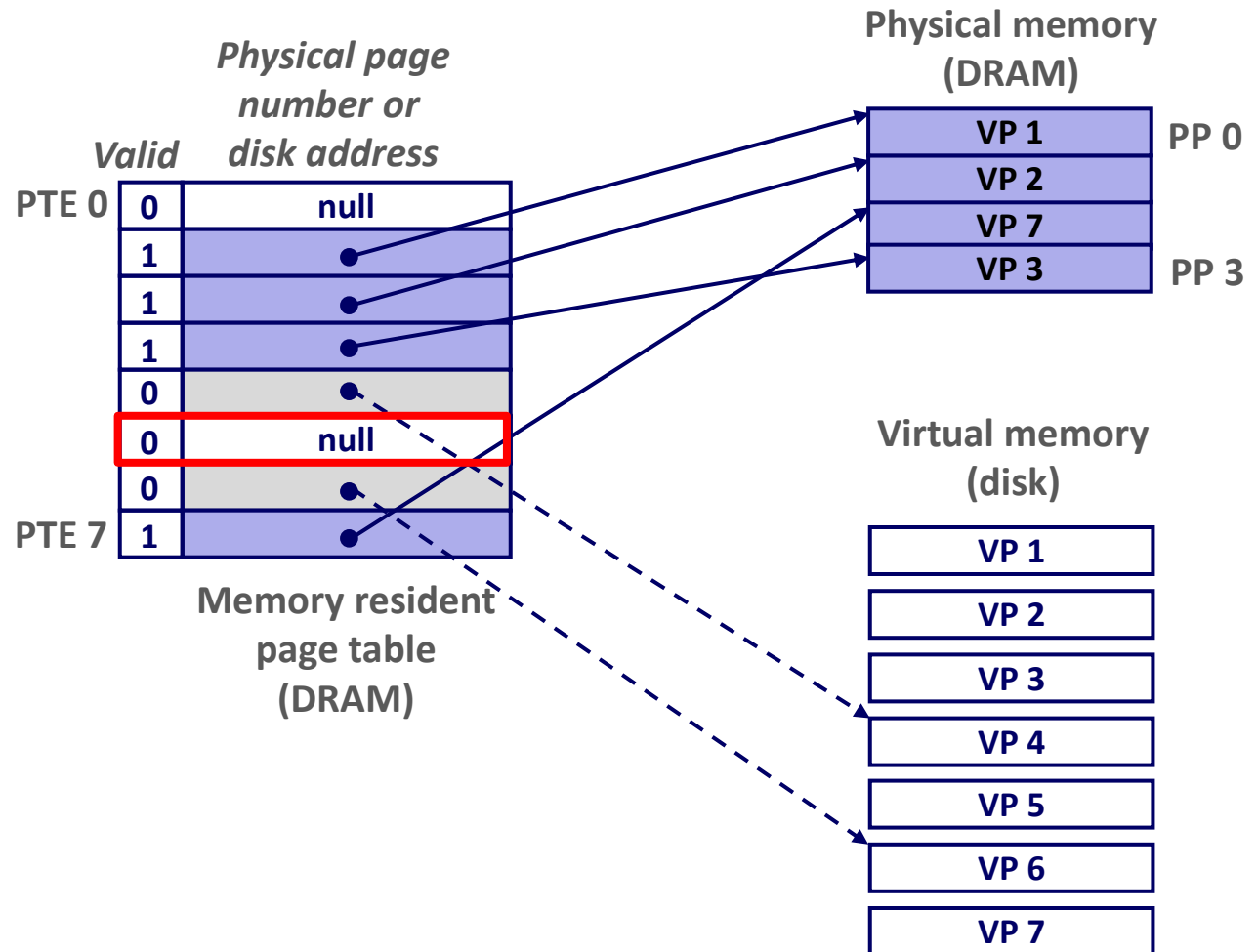
```
int a[1000];
main ()
{
    a[500] = 13;
}
```

```
80483b7:      c7 05 10 9d 04 08 0d  movl   $0xd,0x8049d10
```



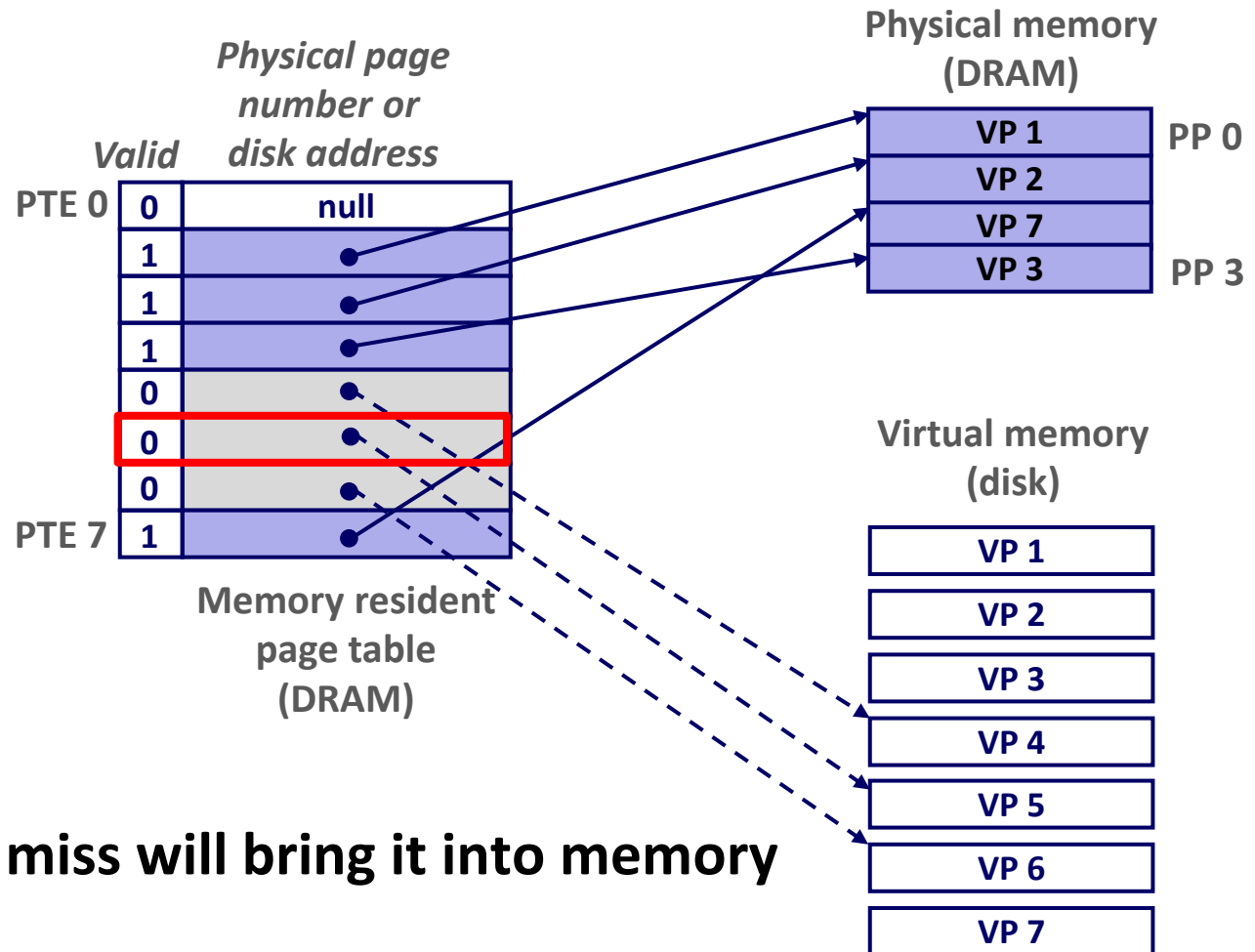
# Allocating Pages

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



# Allocating Pages

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



- Subsequent miss will bring it into memory

# 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 cold misses)
- If (working set size > main memory size )
  - *Thrashing*: Performance meltdown where pages are swapped (copied) in and out continuously
  - If multiple processes run at the same time, thrashing occurs if their total working set size > main memory size

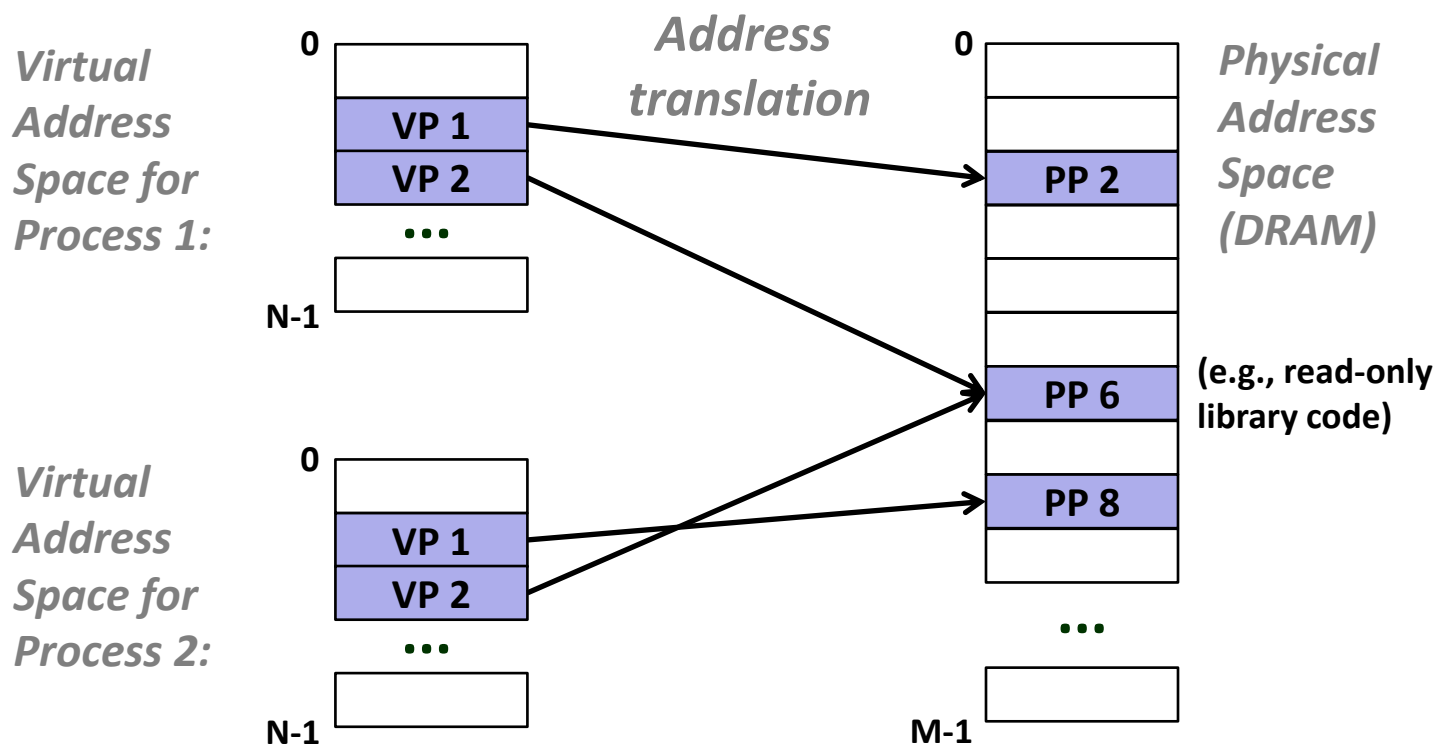


# Today

- Address spaces
- VM as a tool for caching
- **VM as a tool for memory management**
- VM as a tool for memory protection
- Address translation

# VM as a Tool for Memory Management

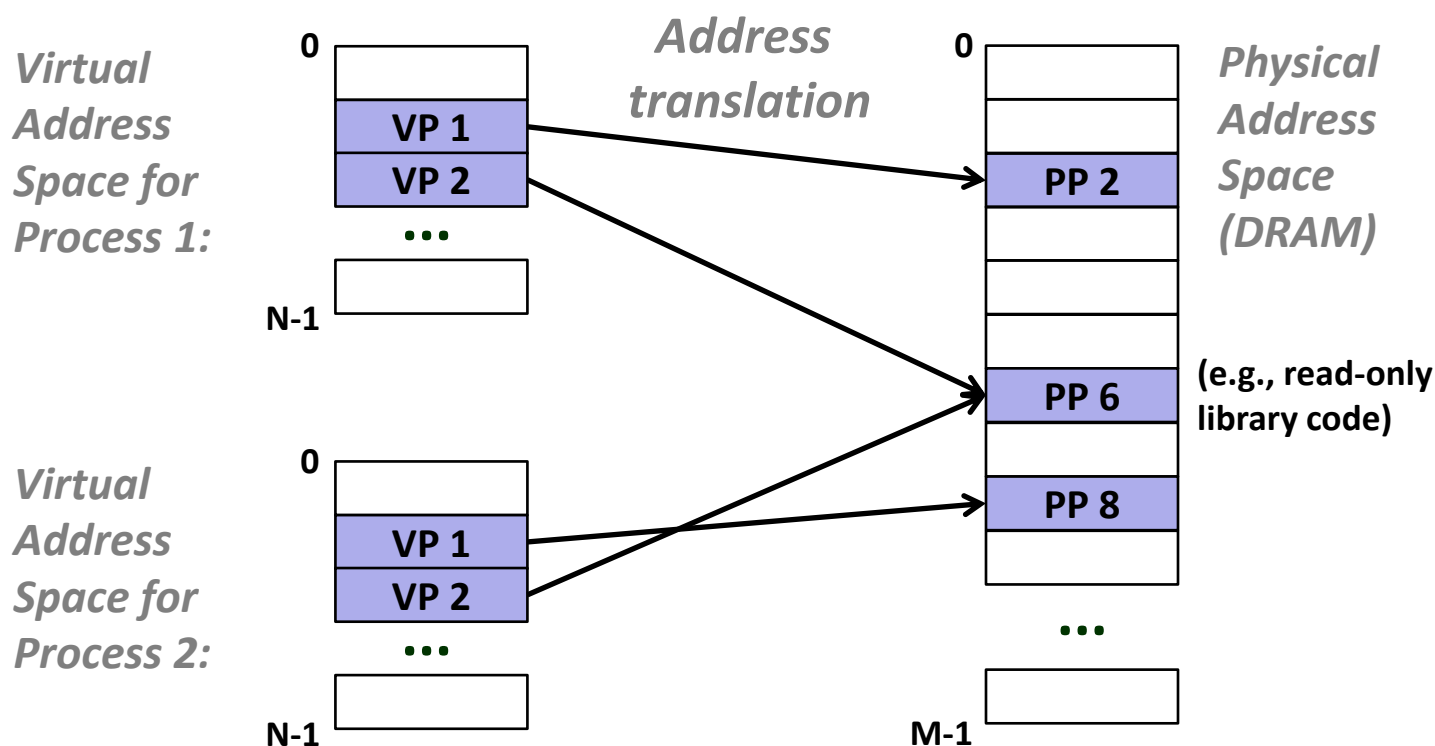
- **Key idea: each process has its own virtual address space**
  - It can view memory as a simple linear array
  - Mapping function scatters addresses through physical memory
    - Well-chosen mappings can improve locality



# VM as a Tool for Memory Management

## ■ Simplifying memory allocation

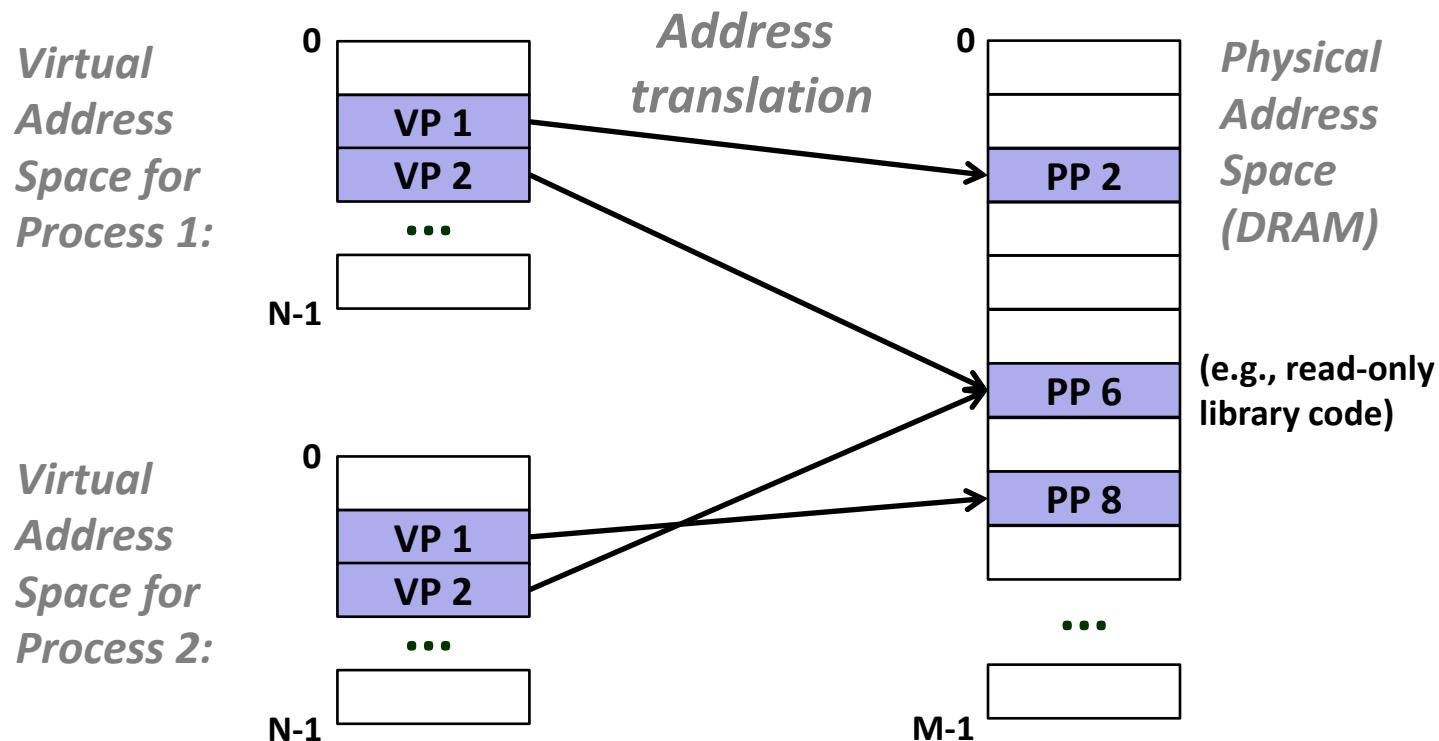
- Each virtual page can be mapped to any physical page
- A virtual page can be stored in different physical pages at different times
- Can allocate the same virtual addresses on the heap for multiple processes



# VM as a Tool for Memory Management

## ■ Sharing code and data among processes

- Map virtual pages to the same physical page (here: PP 6)



# Simplifying Linking and Loading

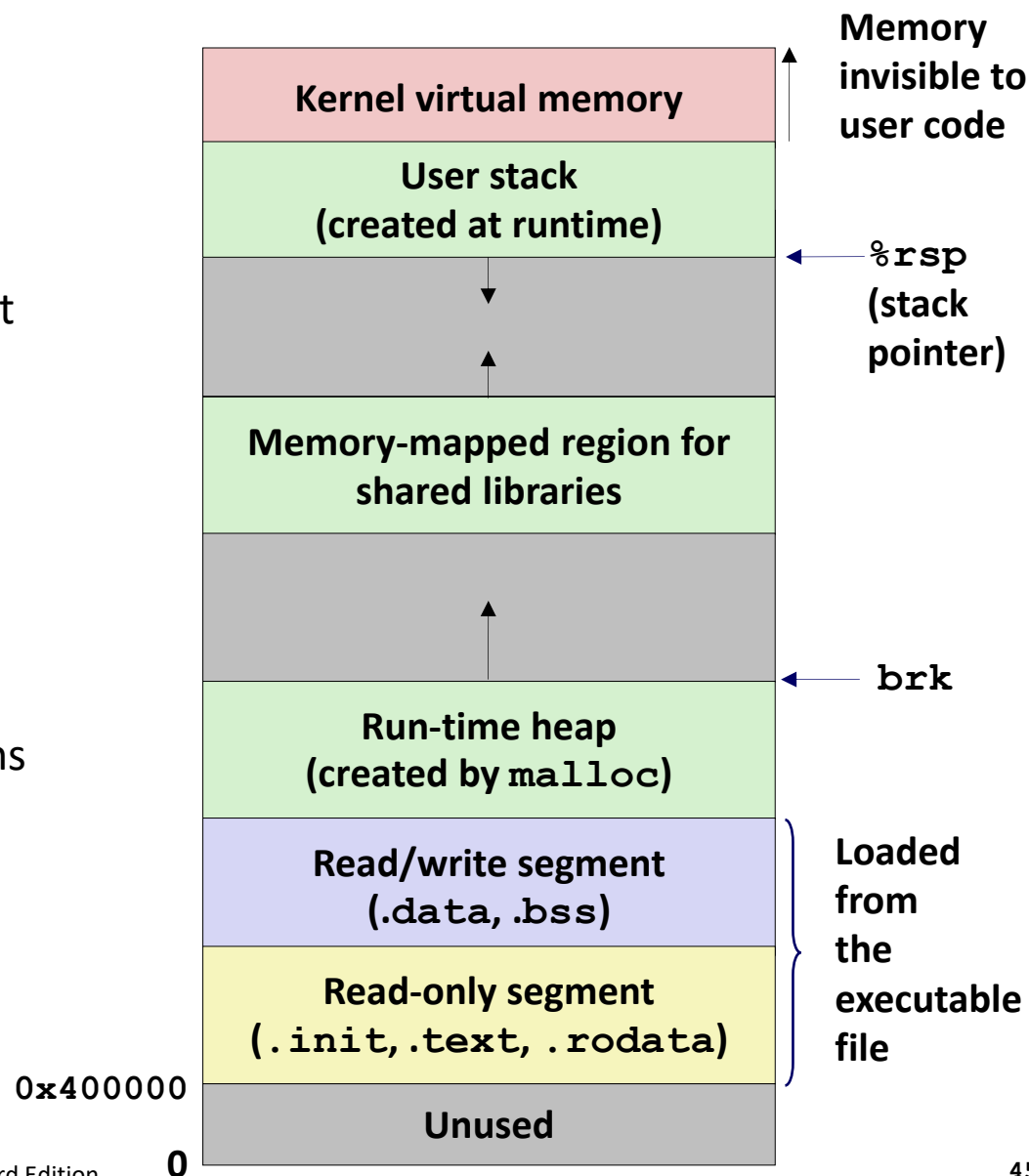
## ■ Linking

- 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

## ■ Discussed later in lecture on Linking and Loading

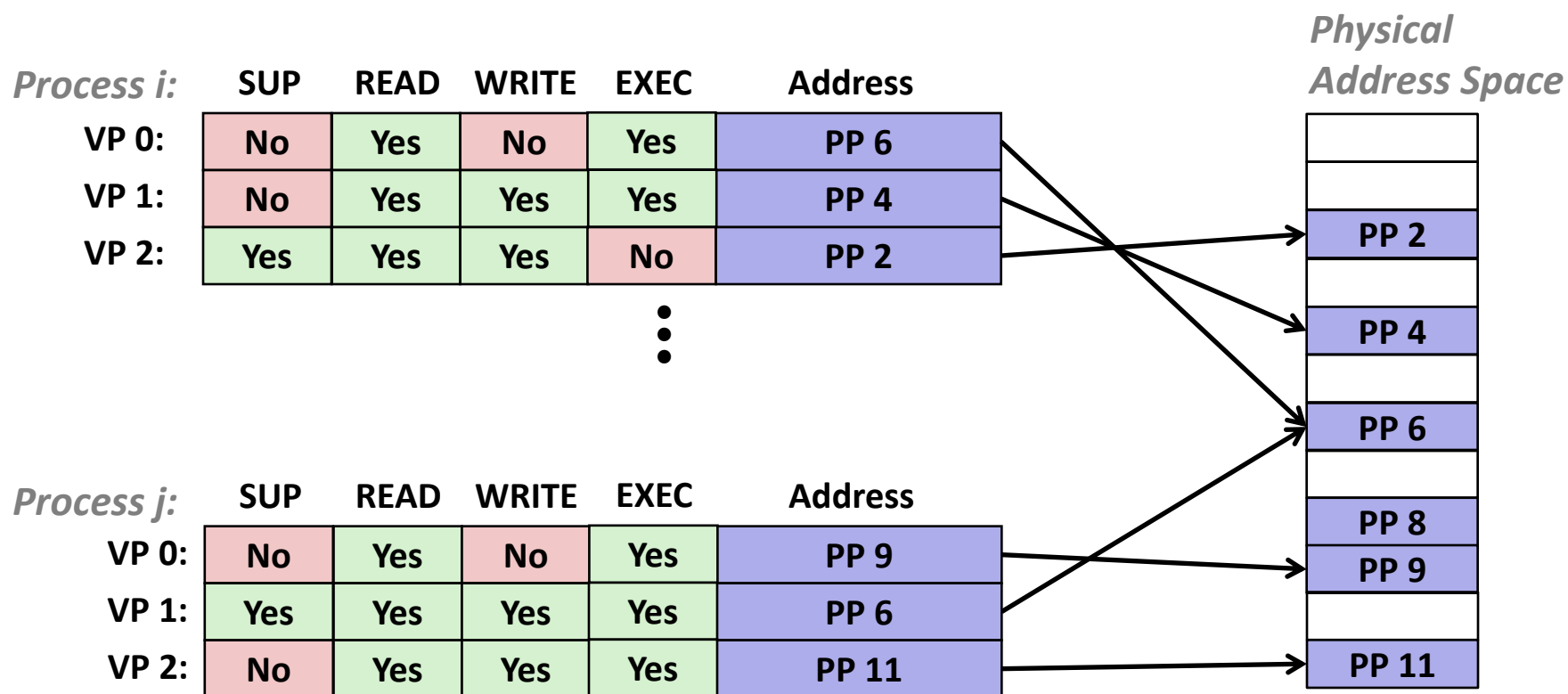


# Today

- Address spaces
- VM as a tool for caching
- VM as a tool for memory management
- **VM as a tool for memory protection**
- Address translation

# VM as a Tool for Memory Protection

- Extend page table entries (PTEs) with permission bits
- MMU checks these bits on each access



**SUP:** requires kernel mode

# Today

- Address spaces
- VM as a tool for caching
- VM as a tool for memory management
- VM as a tool for memory protection
- **Address translation**



# VM Address Translation

## ■ Virtual Address Space

- $V = \{0, 1, \dots, N-1\}$

## ■ Physical Address Space

- $P = \{0, 1, \dots, M-1\}$

## ■ Address Translation

- $MAP: V \rightarrow P \cup \{\emptyset\}$

- For virtual address  $a$ :

- $MAP(a) = a'$  if data at virtual address  $a$  is at physical address  $a'$  in  $P$

- $MAP(a) = \emptyset$  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)

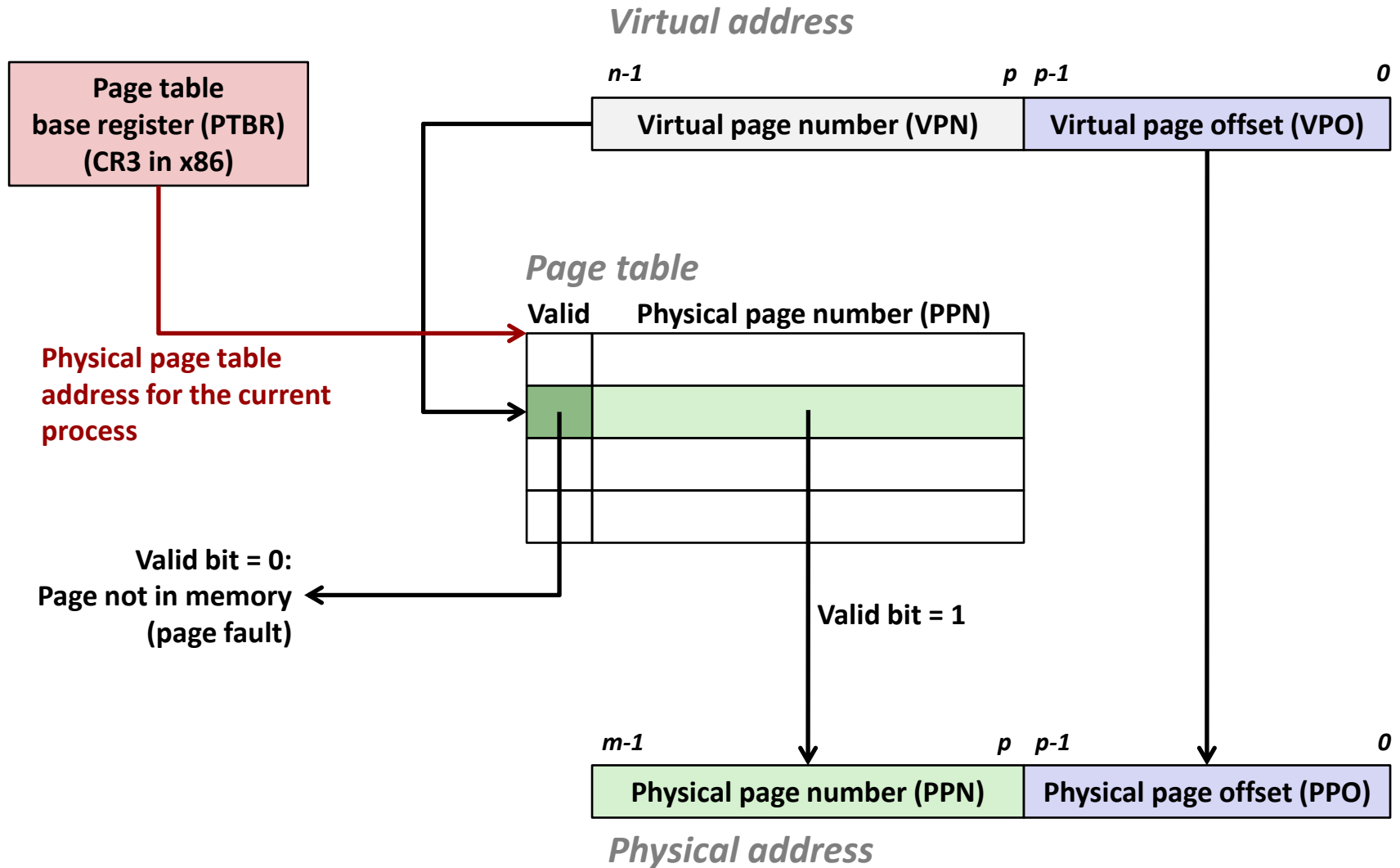
## ■ Components of the virtual address (VA)

- VPO: Virtual page offset
- VPN: Virtual page number

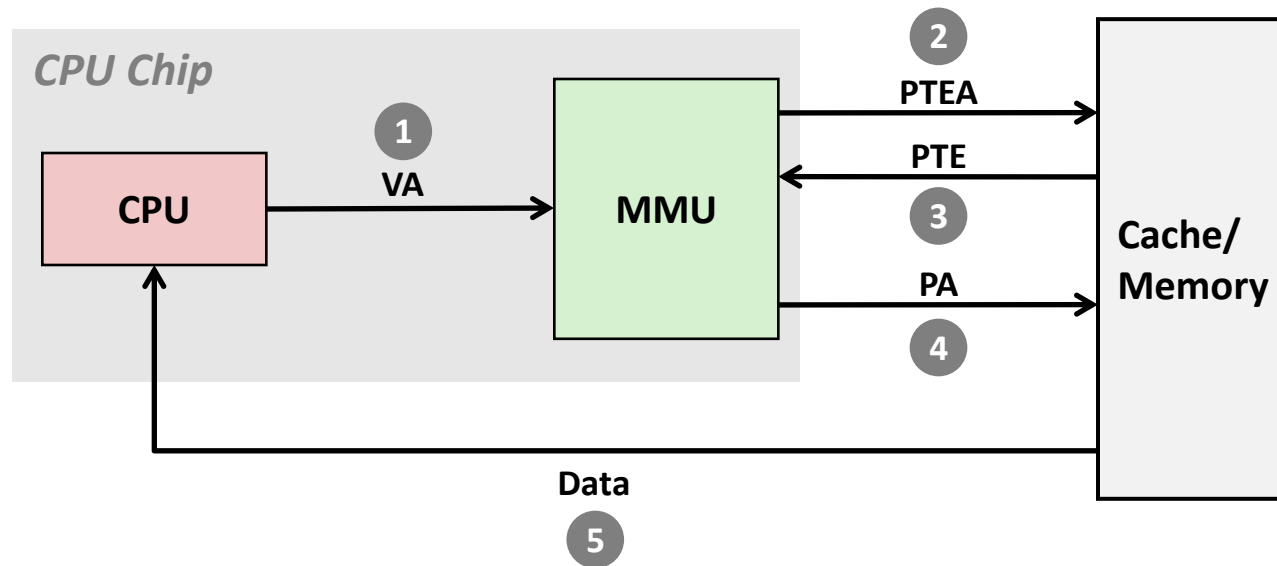
## ■ Components of the physical address (PA)

- PPO: Physical page offset (same as VPO)
- PPN: Physical page number

# Address Translation With a Page Table

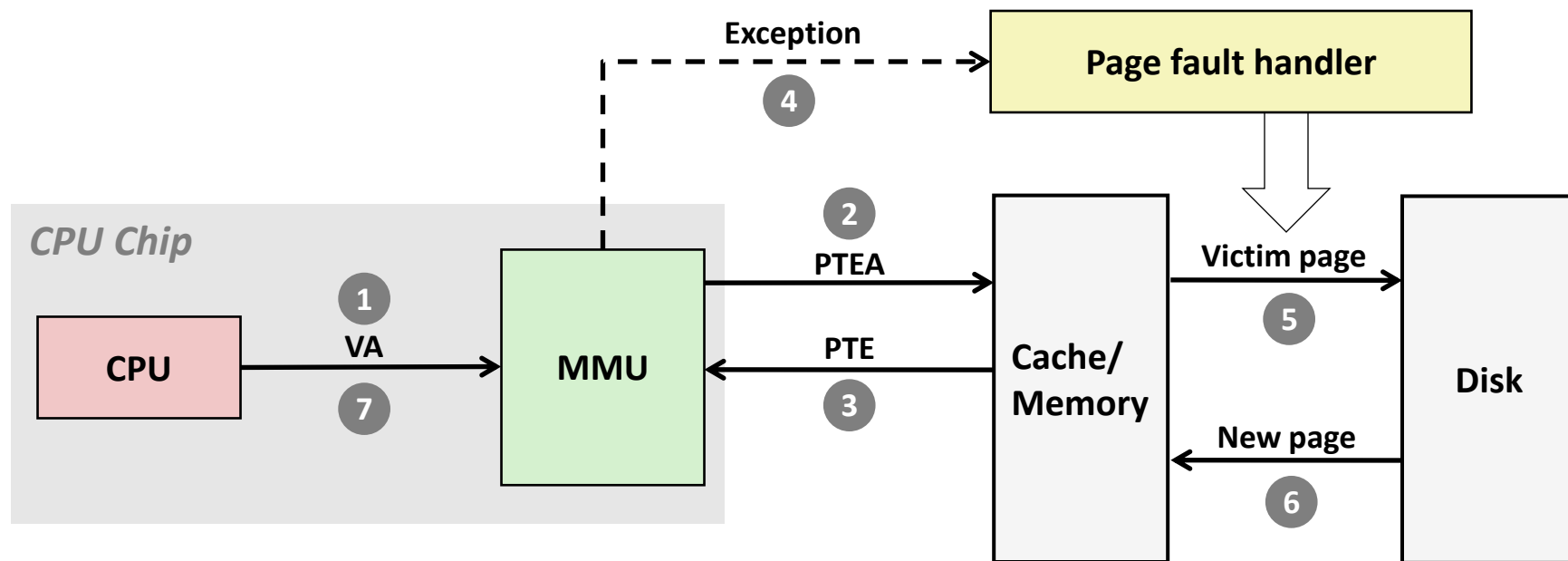


# Address Translation: Page Hit



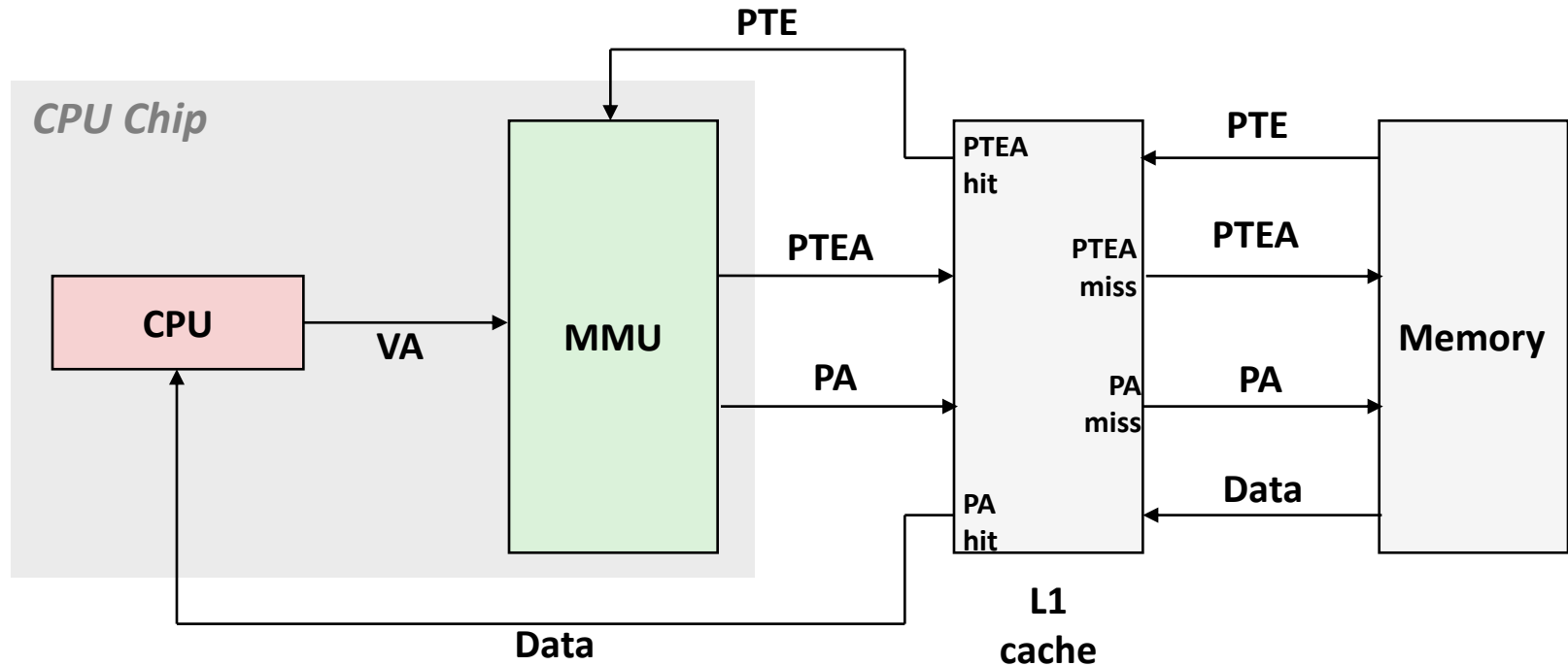
- 1) Processor sends virtual address to MMU
- 2-3) MMU fetches PTE from page table in memory
- 4) MMU sends physical address to cache/memory
- 5) Cache/memory sends data word to processor

# 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 to page out (if dirty, writes pages to disk)
- 6) Handler pages in new page and updates PTE in memory
- 7) Handler returns to original process, restarting faulting instruction

# Integrating VM and 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

# 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

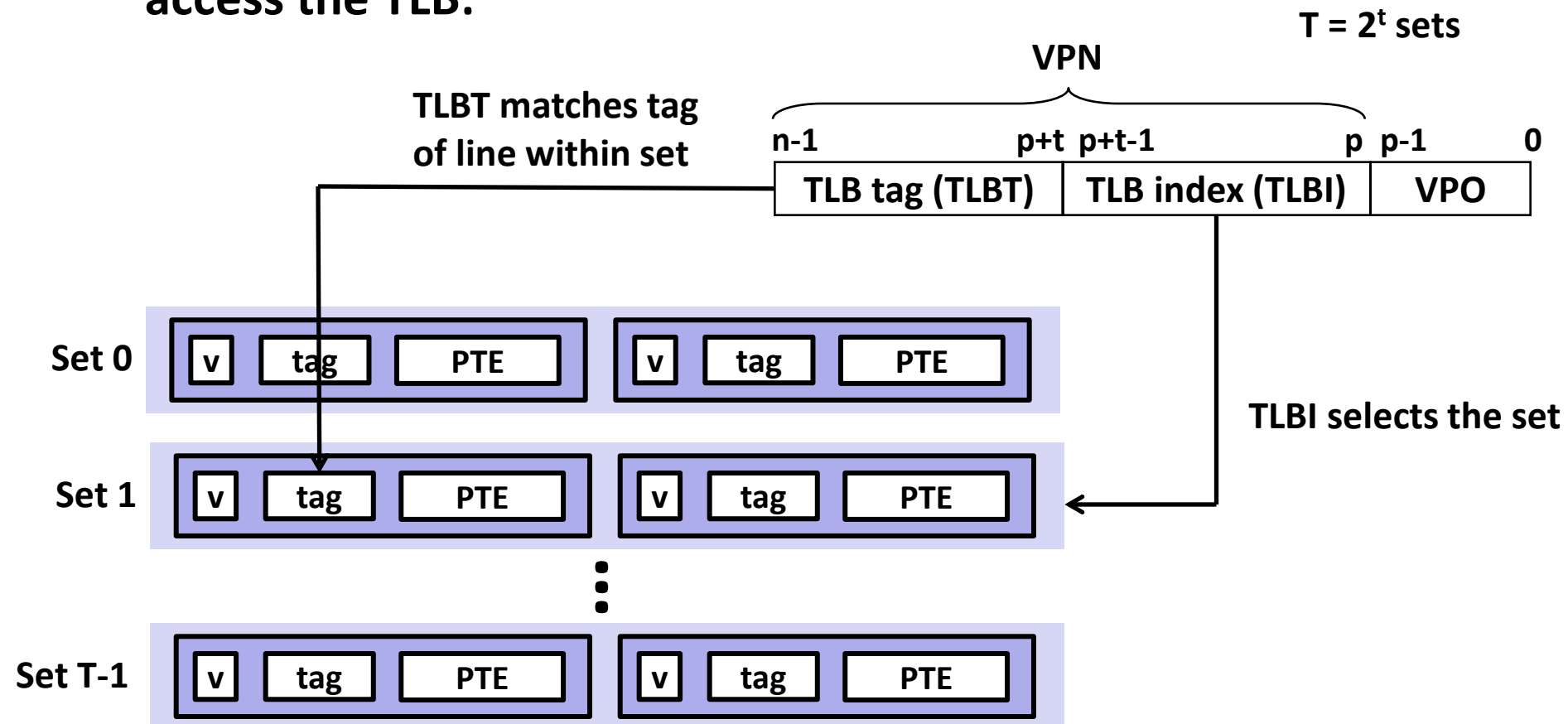
## ■ Components of the physical address (PA)

- **PPO**: Physical page offset (same as VPO)
- **PPN**: Physical page number

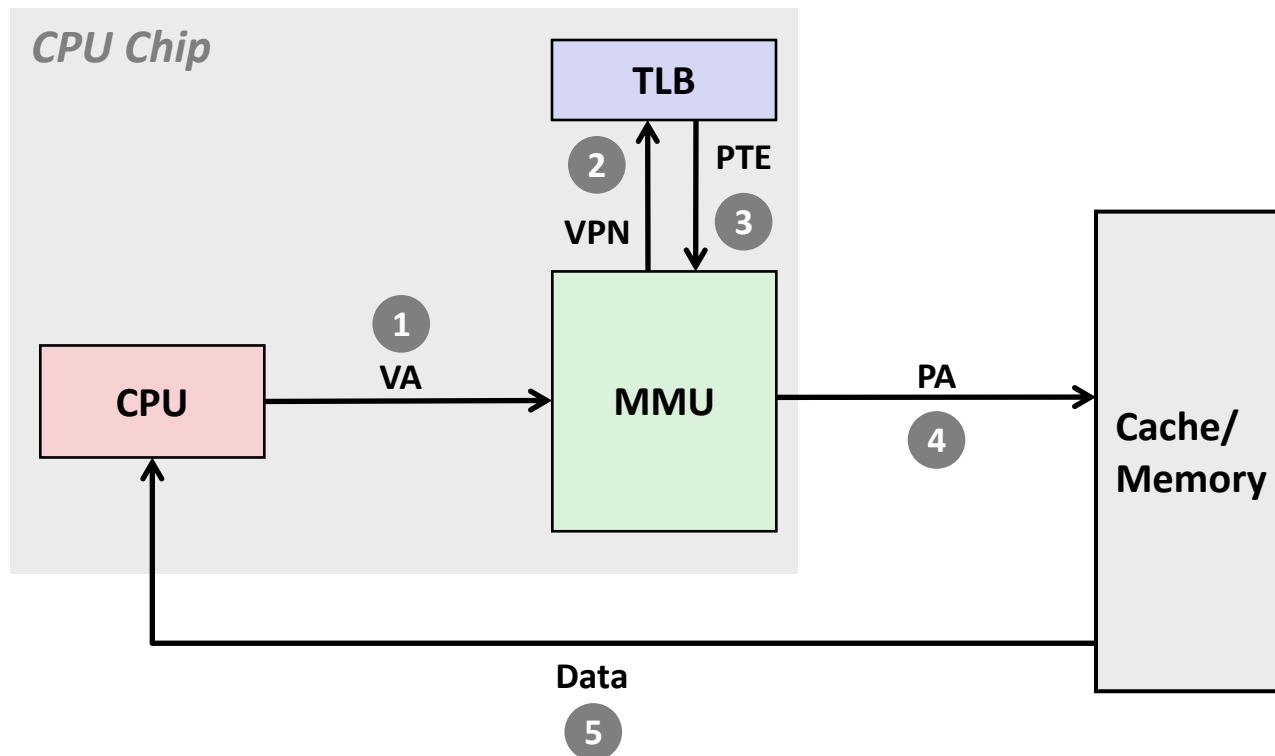


# Accessing the TLB

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

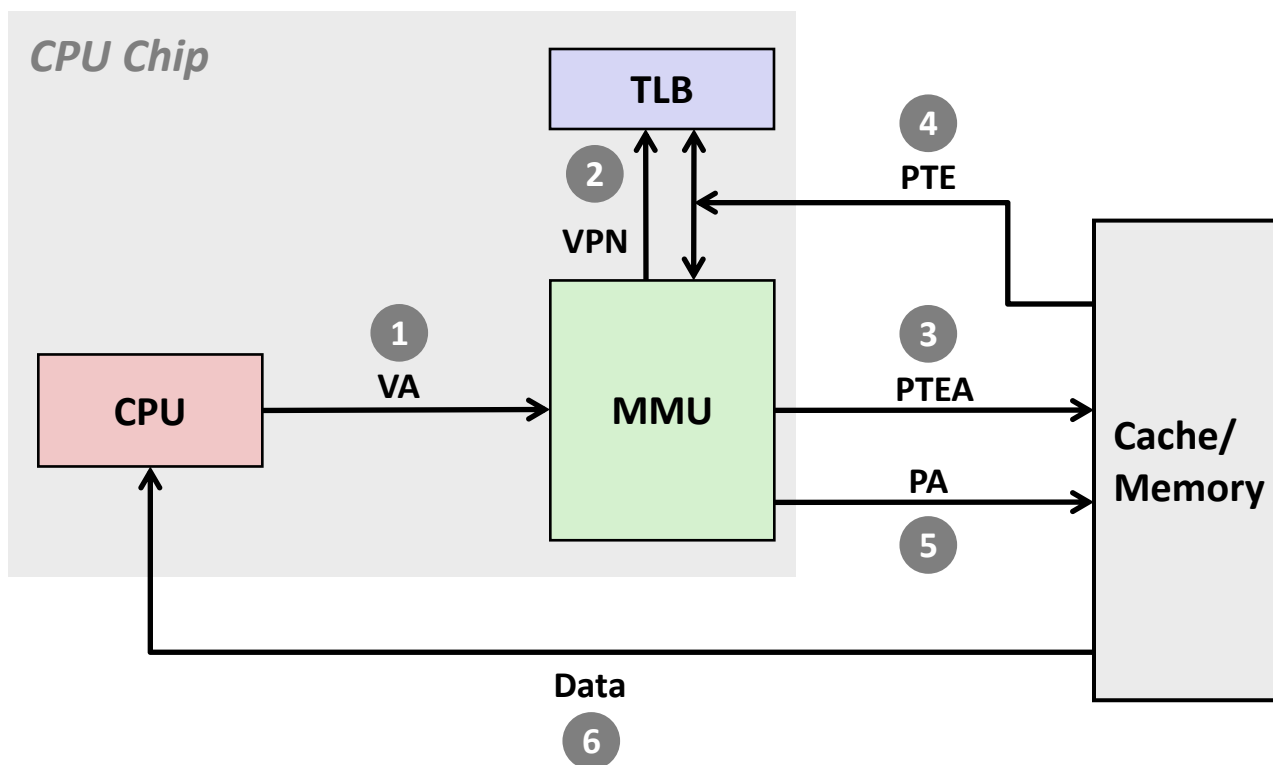


# TLB Hit



**A TLB hit eliminates a cache/memory access**

# TLB Miss



**A TLB miss incurs an additional cache/memory access (the PTE)**

Fortunately, TLB misses are rare. *Why?*

# Multi-Level Page Tables

## ■ Suppose:

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

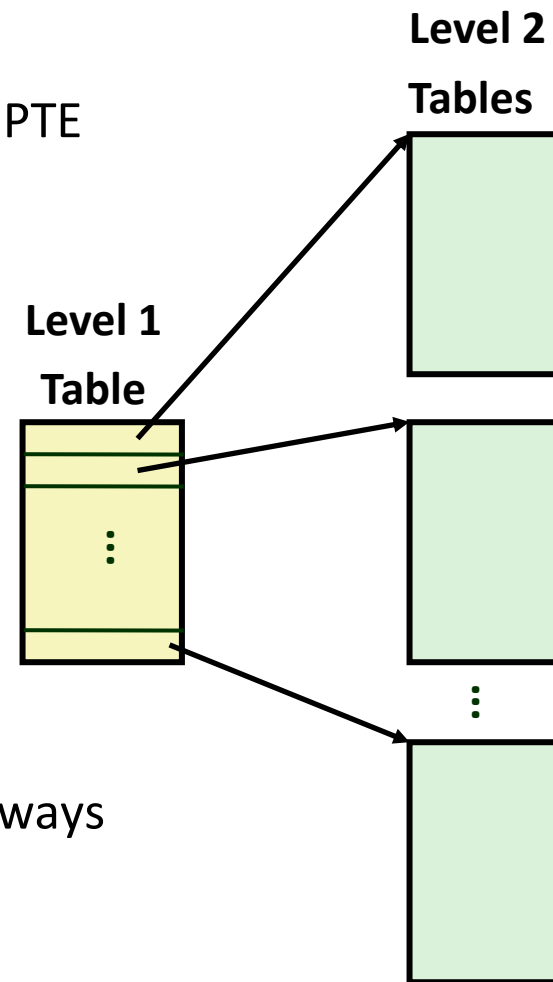
## ■ Problem:

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

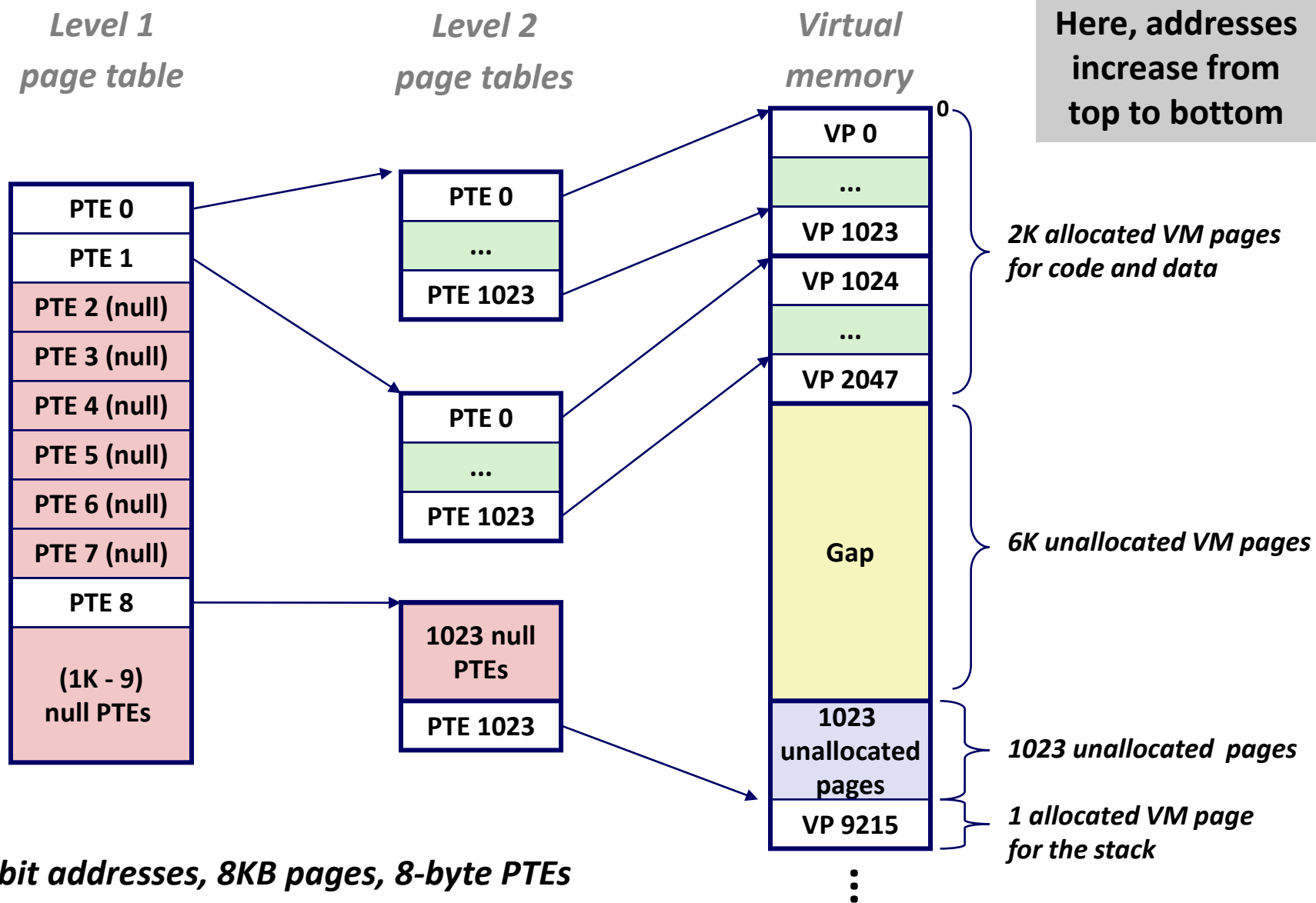
## ■ Common solution: Multi-level page table

## ■ Example: 2-level page table

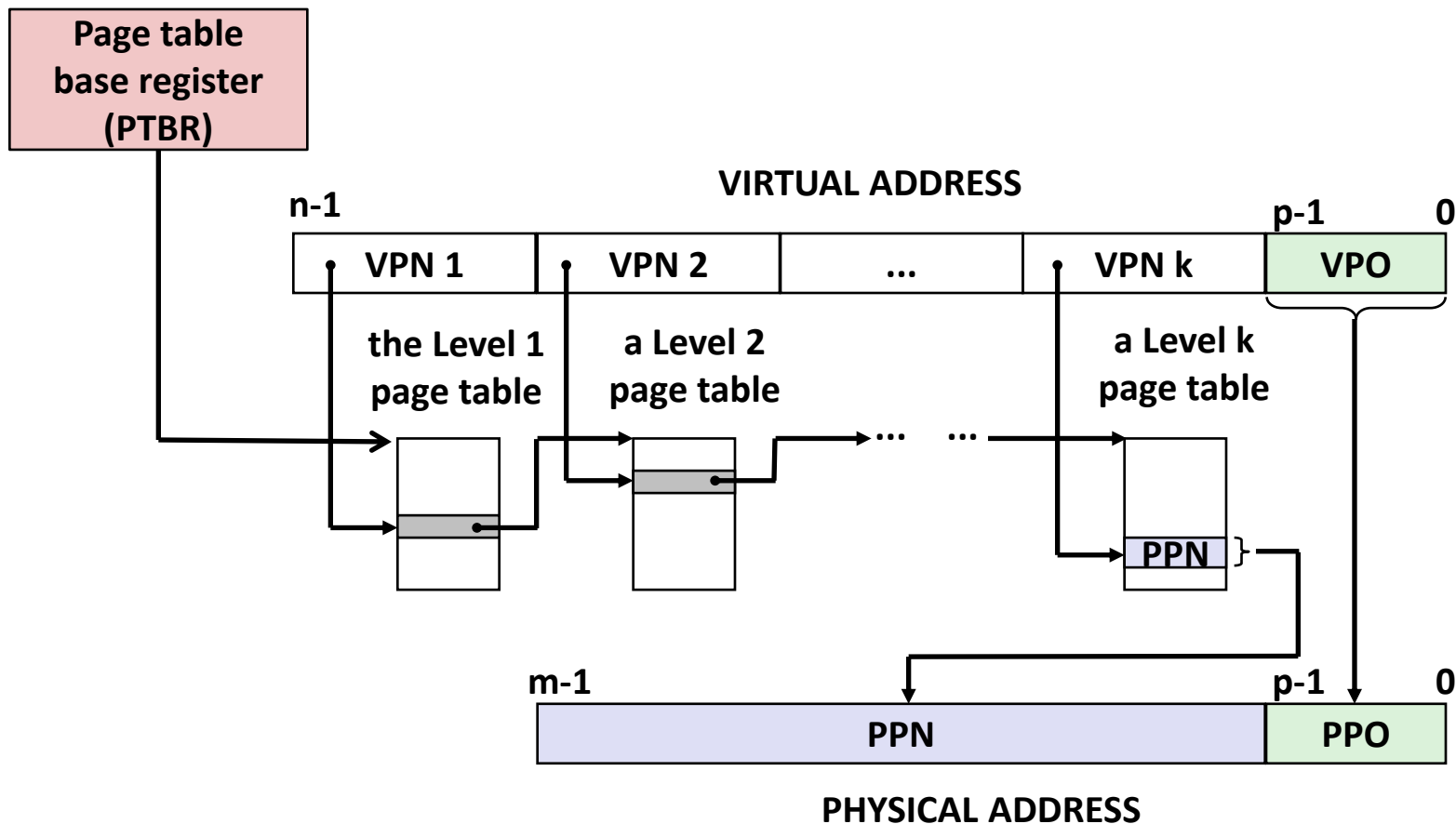
- Level 1 table: each PTE points to a page table (always memory resident)
- Level 2 table: each PTE points to a page (paged in and out like any other data)



# A Two-Level Page Table Hierarchy



# 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

## ■ Implemented via combination of hardware & software

- MMU, TLB, exception handling mechanisms part of hardware
- Page fault handlers, TLB management performed in software