# Paging: inside the OS
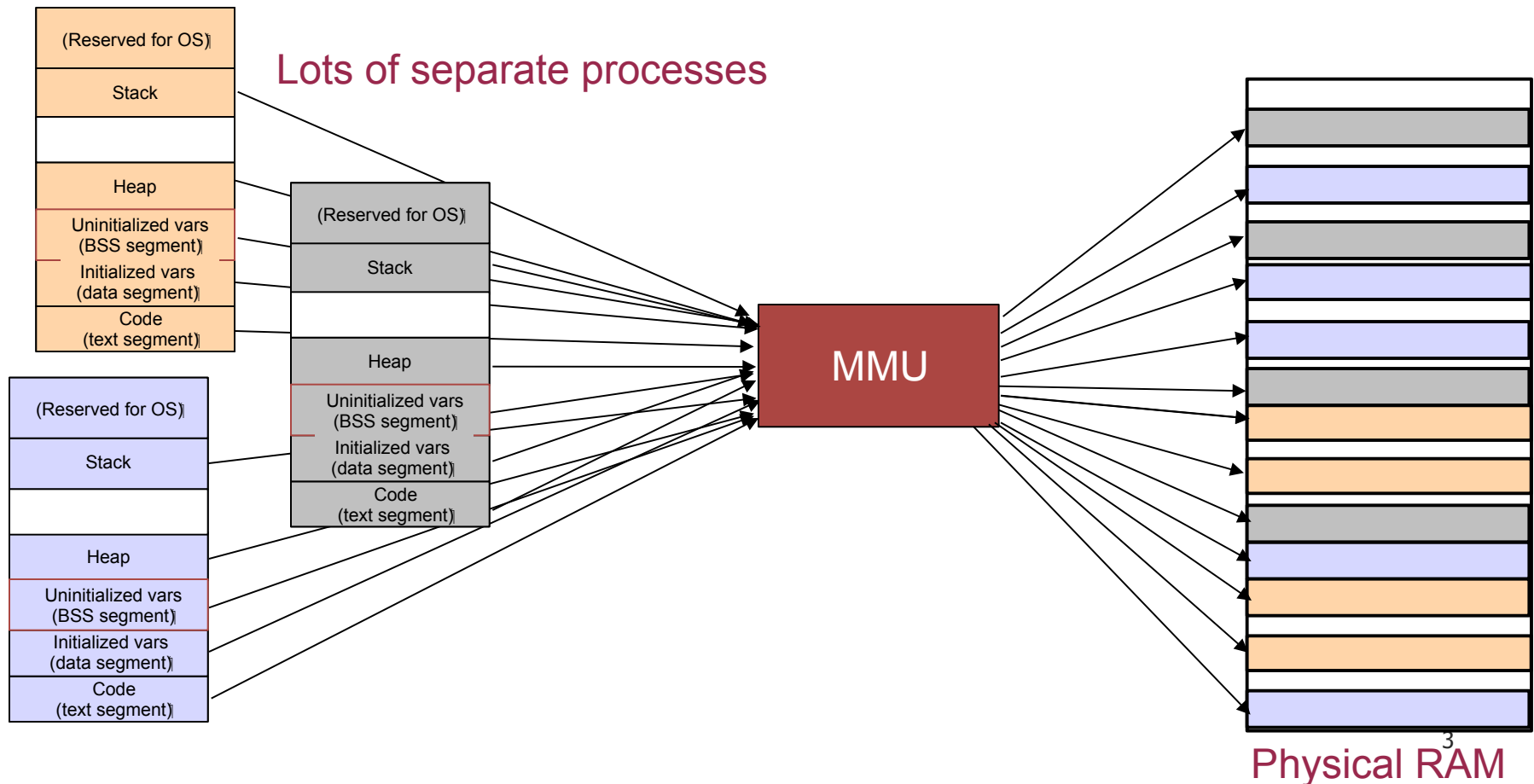
CS 241

February 8, 2012

# Paging

- Solve the external fragmentation problem by using **fixed-size chunks** of virtual and physical memory
  - Virtual memory unit called a **page**
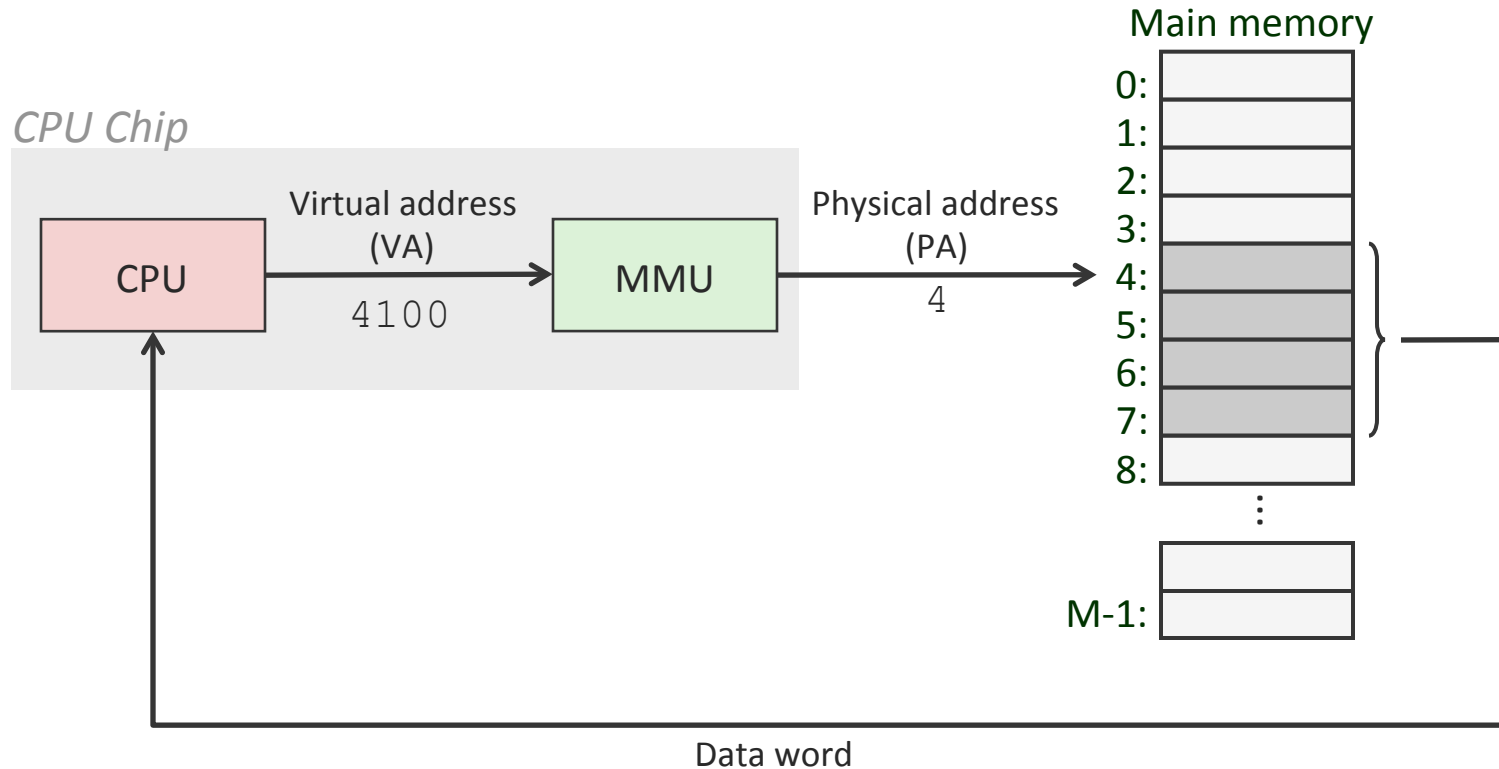  - Physical memory unit called a **frame** (or sometimes **page frame**)

**virtual memory
(for one process)**

| |
|---|
| page 0 |
| page 1 |
| page 2 |
| page 3 |
| ⋮ |
| page X |

**physical memory**

| |
|---|
| frame 0 |
| frame 1 |
| frame 2 |
| ⋮ |
| frame Y |

# Application Perspective

- Application believes it has a single, contiguous address space ranging from 0 to 2P – 1 bytes
  - Where P is the number of bits in a pointer (e.g., 32 bits)
- In reality, virtual pages are scattered across physical memory
  - This mapping is invisible to the program, and not even under it's control!

Lots of separate processes

| (Reserved for OS) |
| Stack |
| |
| Heap |
| Uninitialized vars (BSS segment) |
| Initialized vars (data segment) |
| Code (text segment) |

| (Reserved for OS) |
| Stack |
| |
| Heap |
| Uninitialized vars (BSS segment) |
| Initialized vars (data segment) |
| Code (text segment) |

| (Reserved for OS) |
| Stack |
| |
| Heap |
| Uninitialized vars (BSS segment) |
| Initialized vars (data segment) |
| Code (text segment) |

MMU

Physical RAM
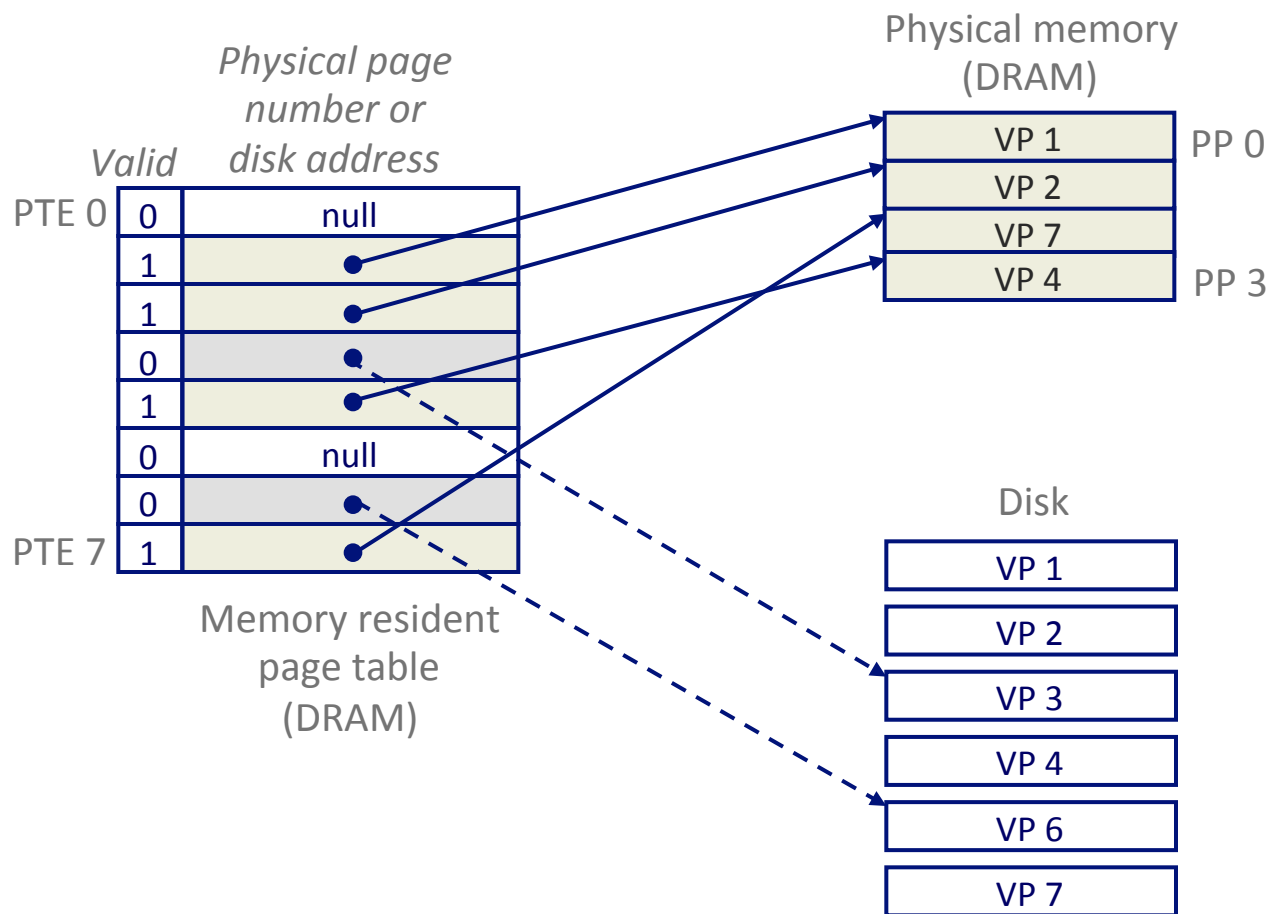
3

# Virtual addressing



- Used in all modern servers, desktops, and laptops
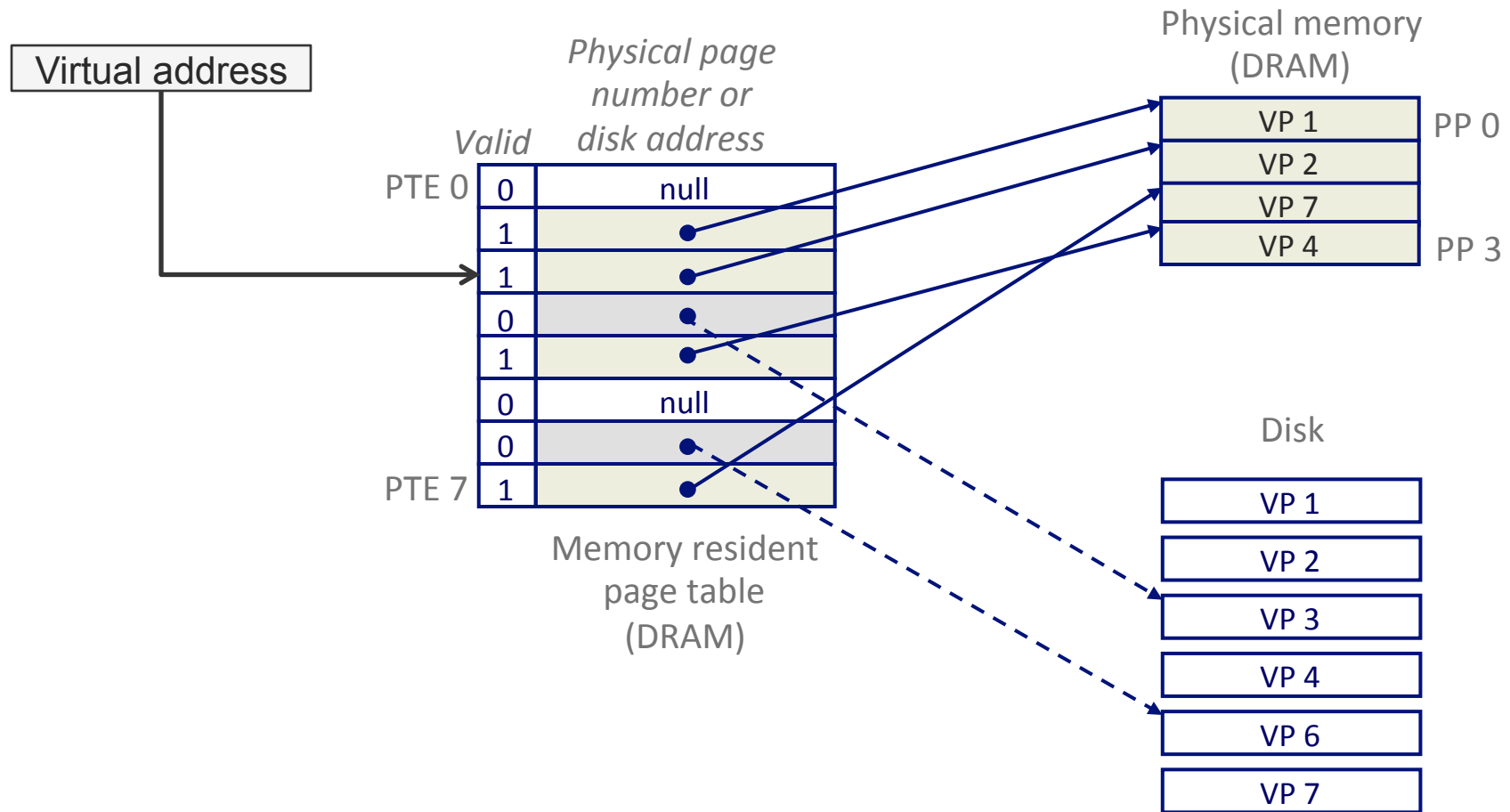- One of the great ideas in computer science

# Enabling data structure

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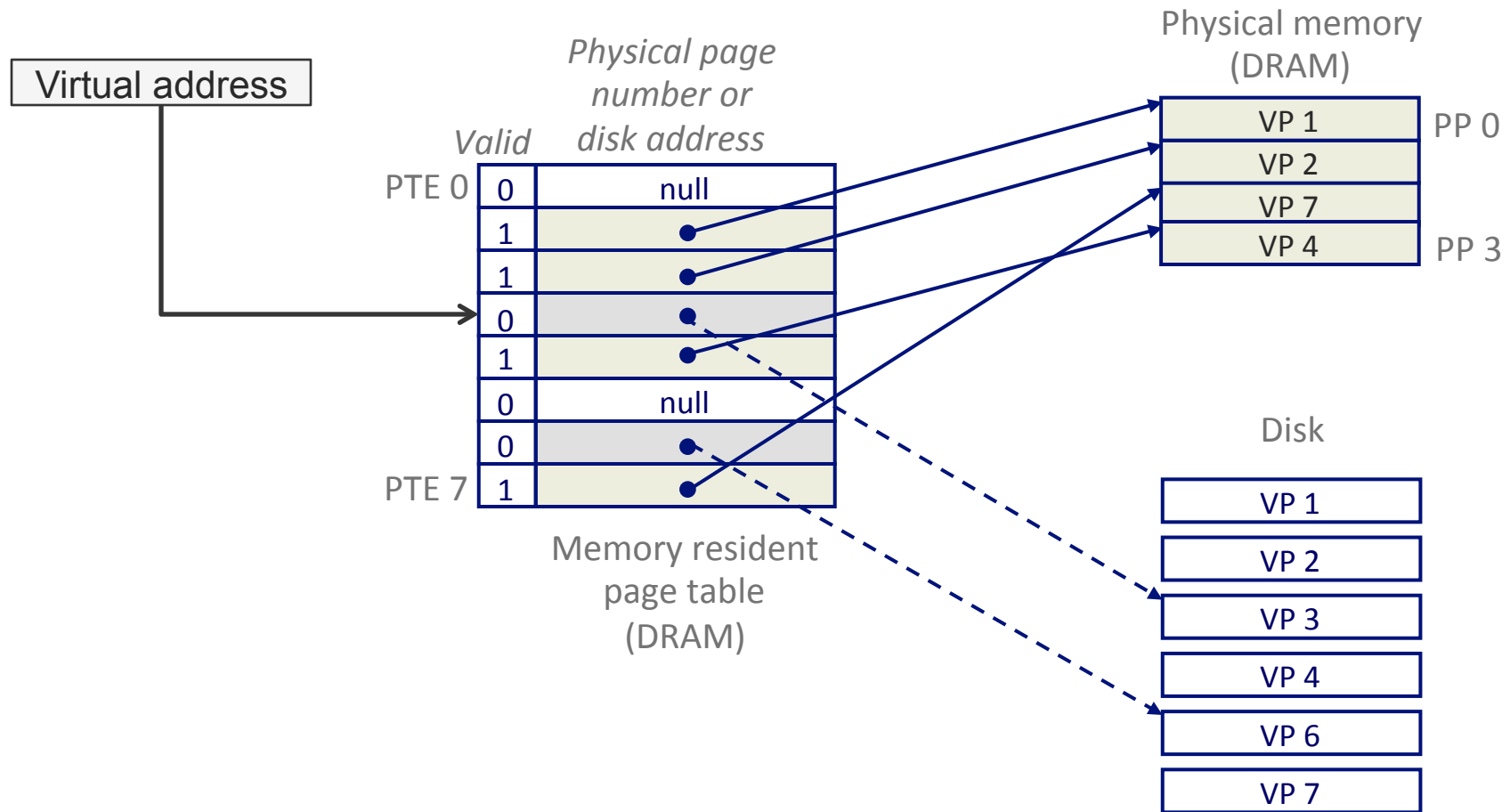  - Per-process kernel data structure in DRAM

# Page hit

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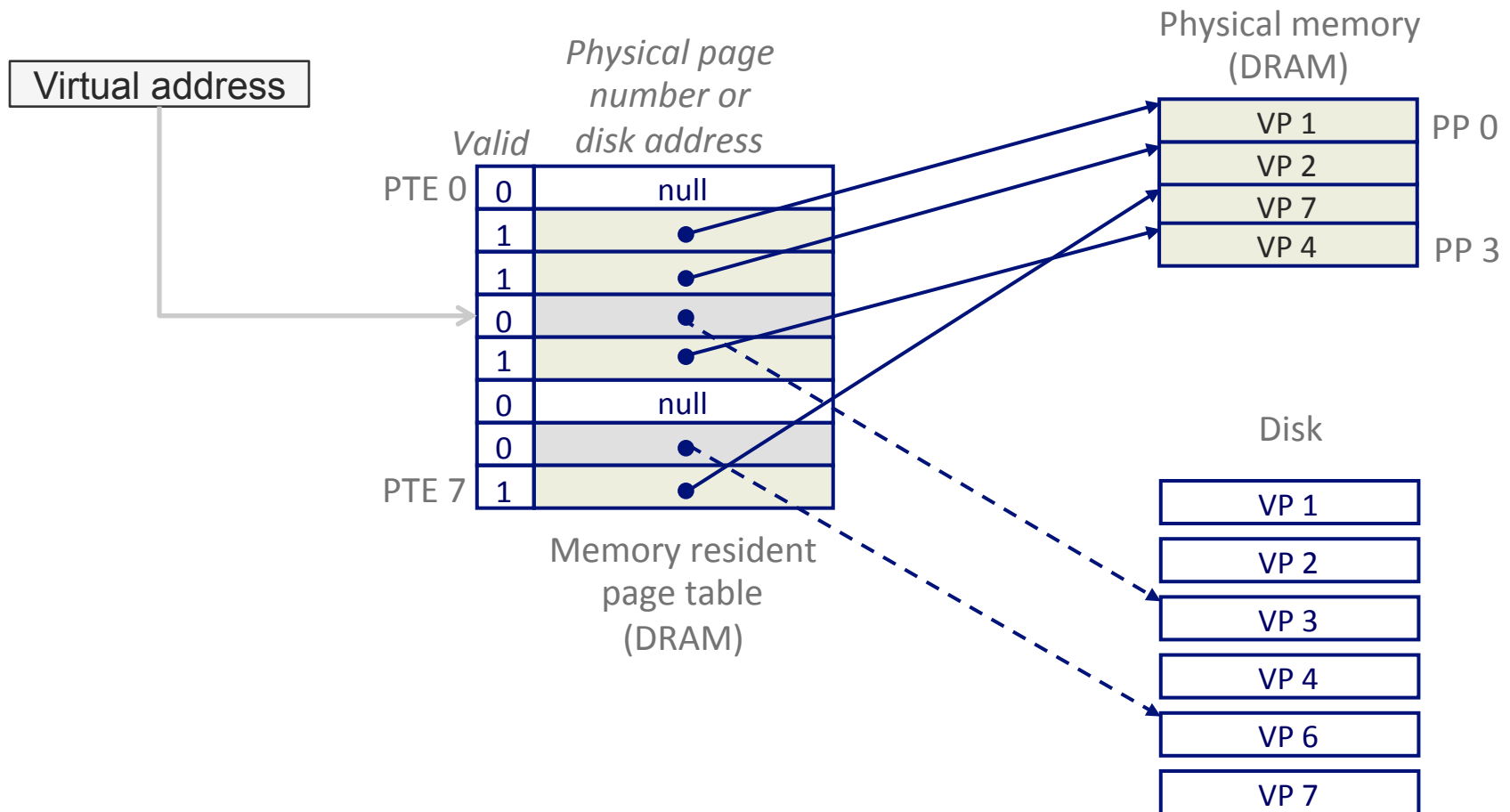
# Page fault

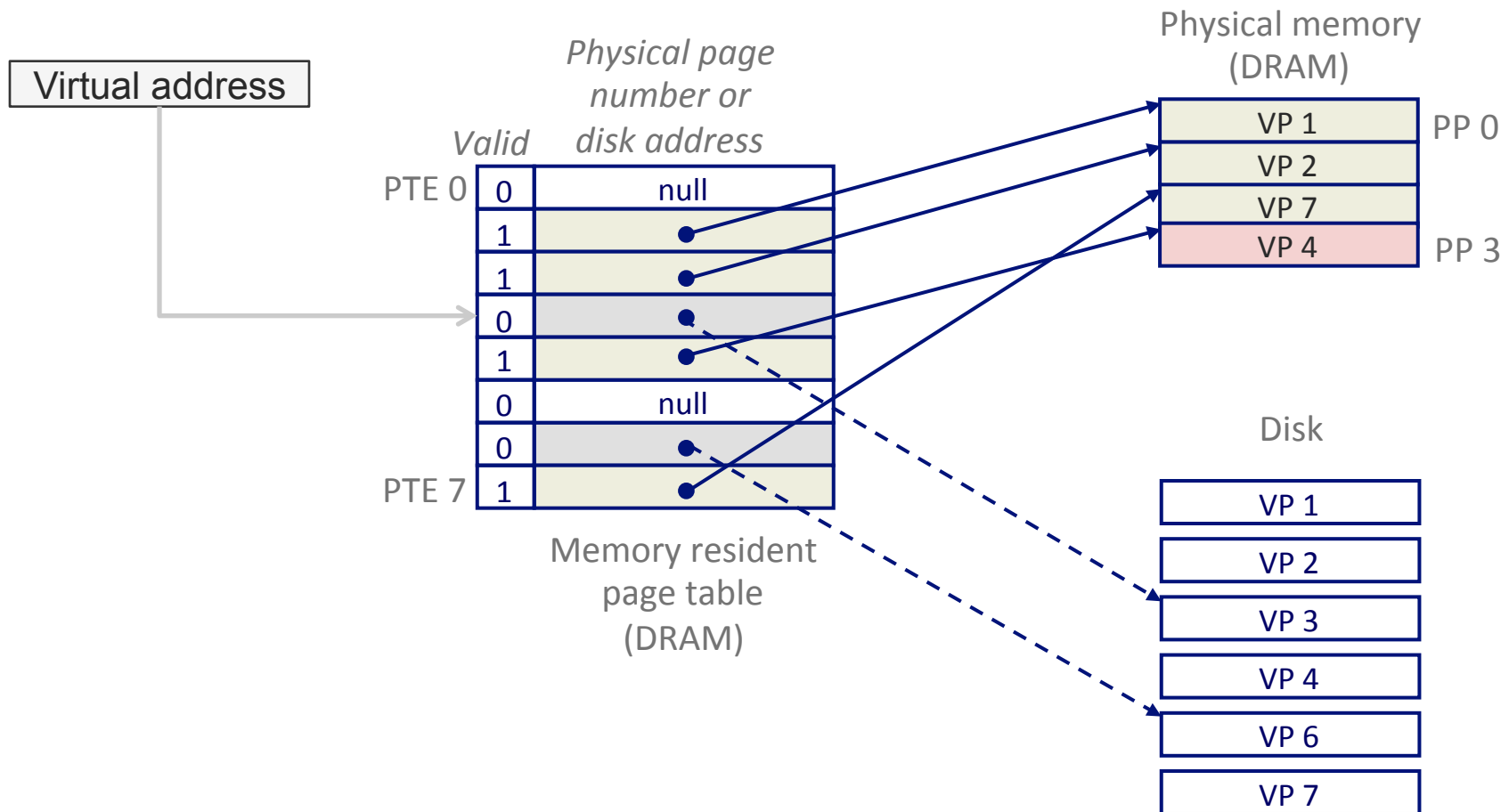- *Page fault:* reference to VM word that is not in physical memory (DRAM cache miss)

# Handling page fault

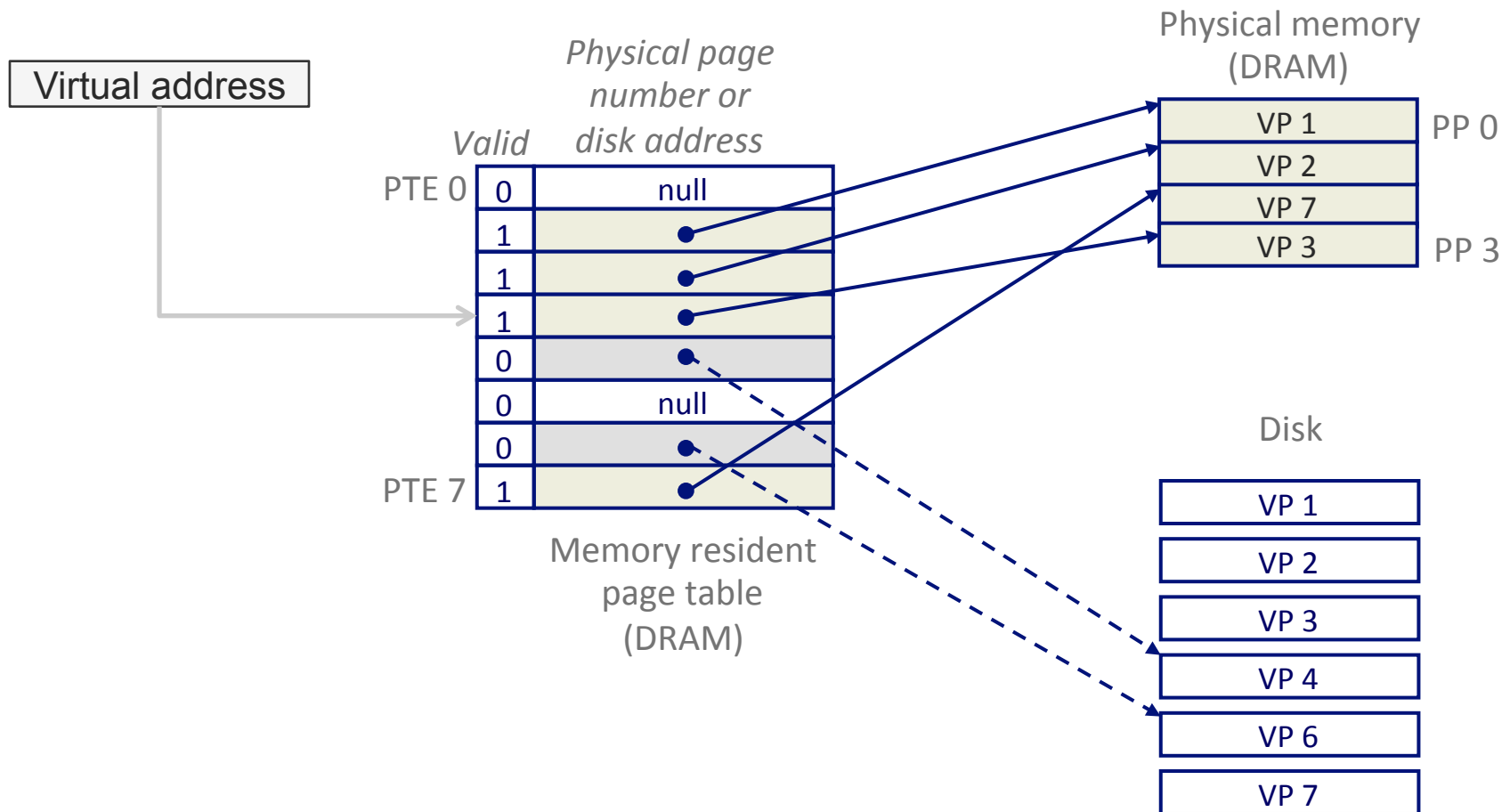- Page miss causes page fault (an exception)

# 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)
- Loads new frame into freed slot
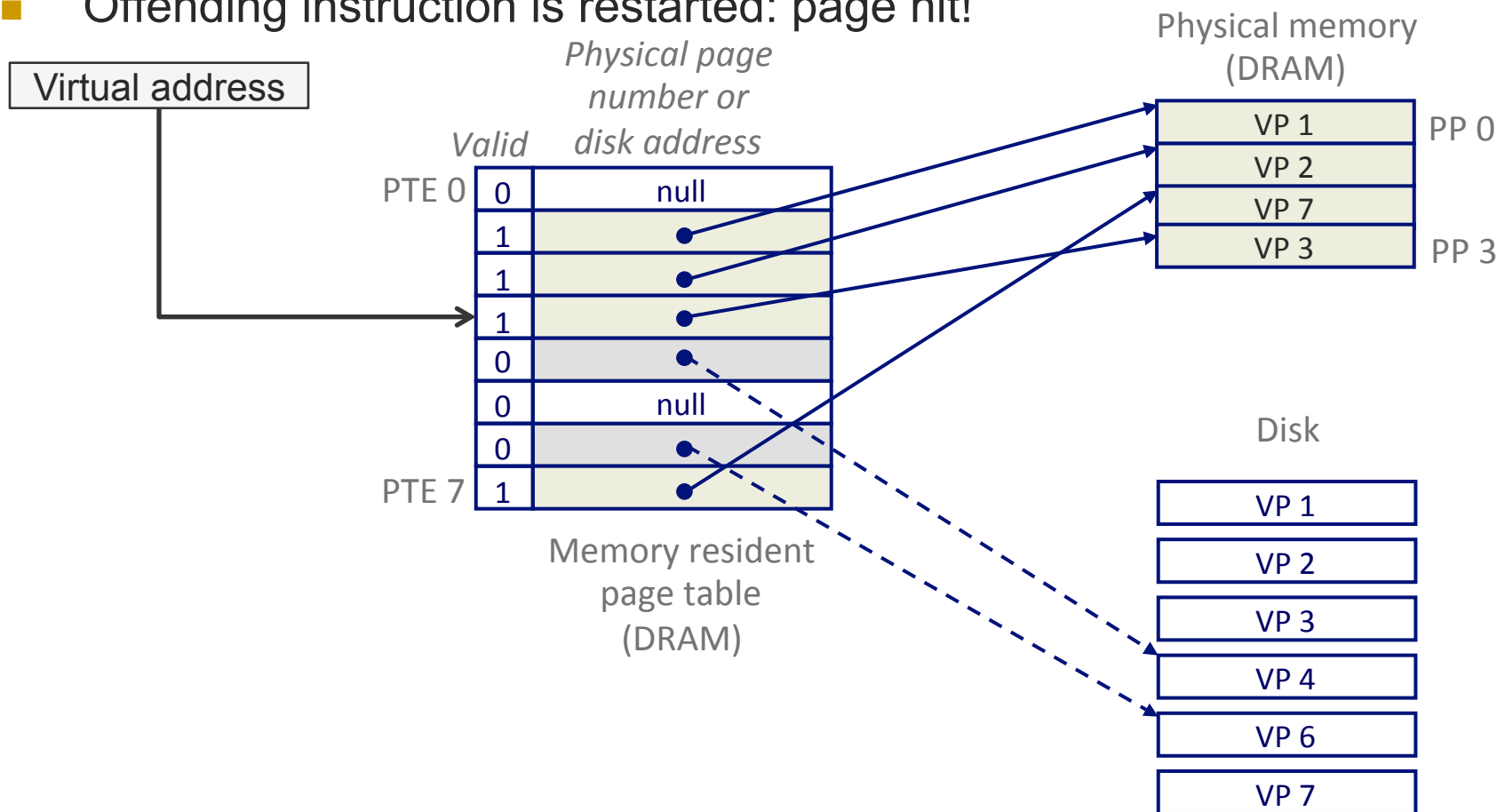
# Handling page fault

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

Virtual address

Physical page number or disk address

Physical memory (DRAM)

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

Memory resident page table (DRAM)

VP 1 — PP 0
VP 2
VP 7
VP 3 — PP 3

Disk

VP 1
VP 2
VP 3
VP 4
VP 6
VP 7

# Page Table Entry

- Typical PTE format (depends on CPU architecture!)

| 1 | 1 | 1 | 2 | 20 |
|---|---|---|---|---|
| M | R | V | prot | physical page (frame) number |

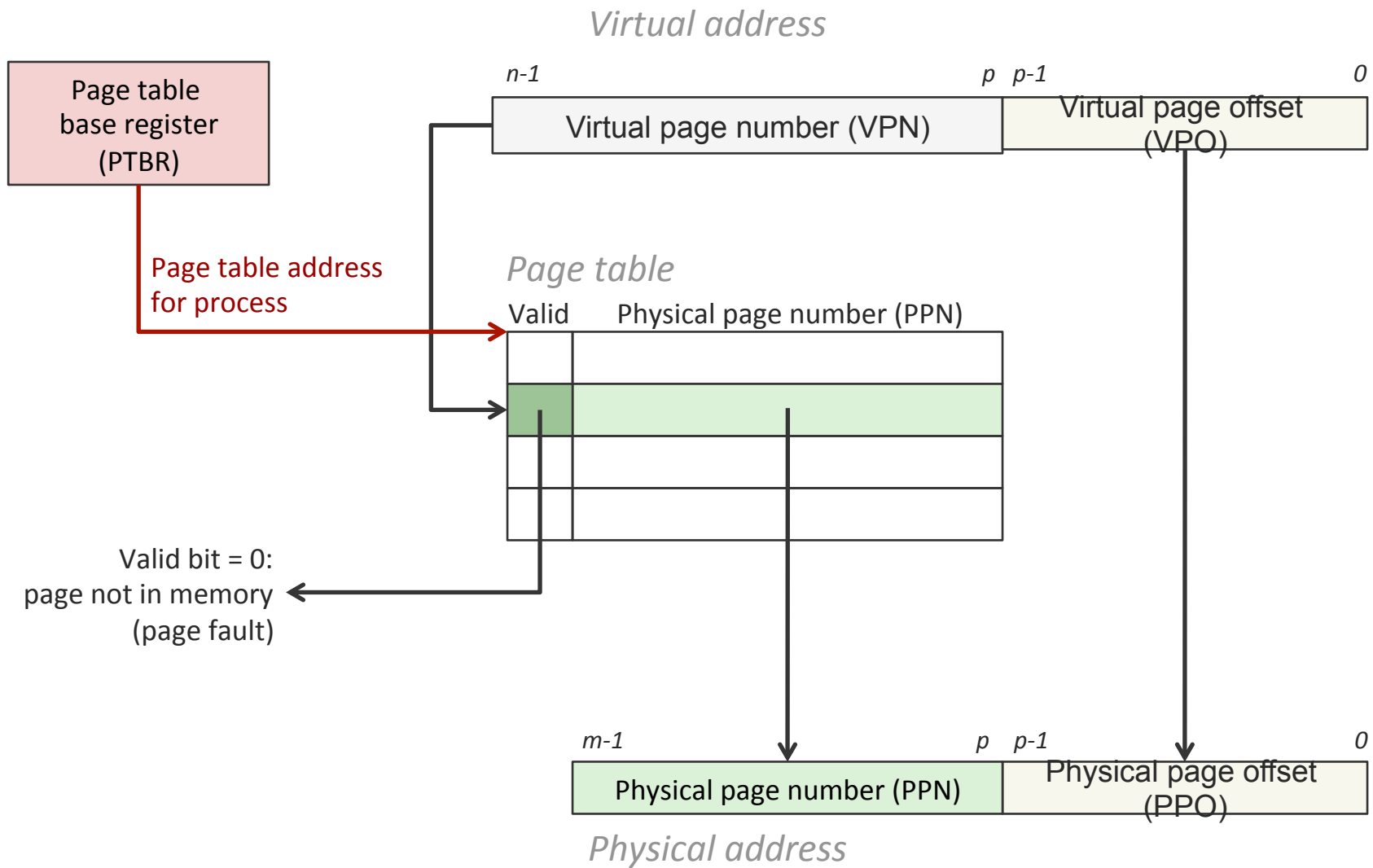- Various bits accessed by MMU on each page access:
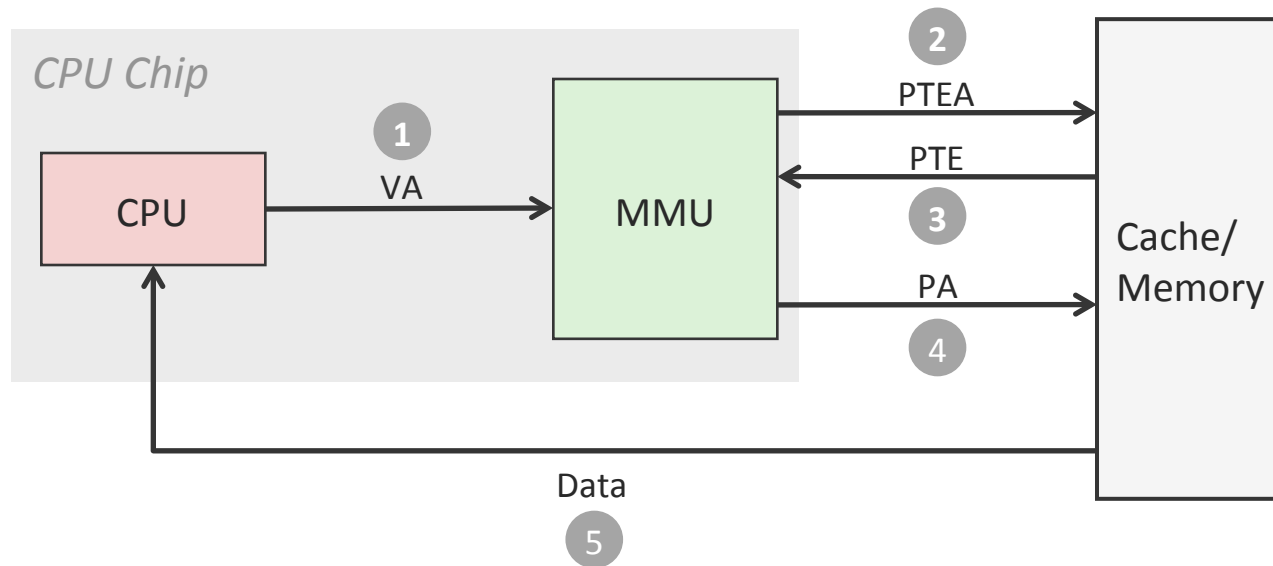
  - Modify bit: Indicates whether a page is "dirty" (modified)
  - Reference bit: Indicates whether a page has been accessed (read or written)
  - Valid bit: Whether the PTE represents a real memory mapping
  - Protection bits: Specify if page is readable, writable, or executable
  - Physical page number: Physical location of page in RAM
    - Why is this 20 bits wide in the above example?
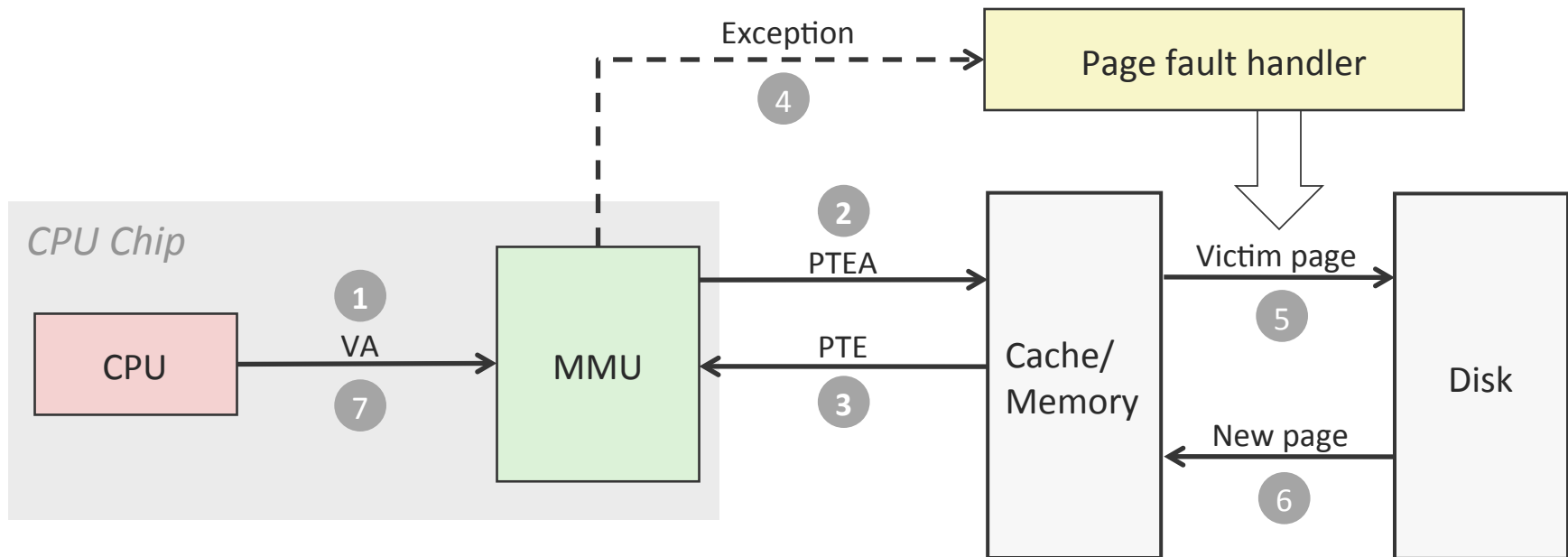
# Address translation with a P.T.

Page table
base register
(PTBR)

Page table address
for process

*Virtual address*

$n-1$                         $p$   $p-1$             $0$

| Virtual page number (VPN) | Virtual page offset (VPO) |
|---|---|

*Page table*

Valid     Physical page number (PPN)

Valid bit = 0:
page not in memory
(page fault)

$m-1$                    $p$   $p-1$             $0$

| Physical page number (PPN) | Physical page offset (PPO) |
|---|---|

*Physical address*

13

# 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 (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

# Question 1

- Isn't it slow to have to go to memory twice every time?
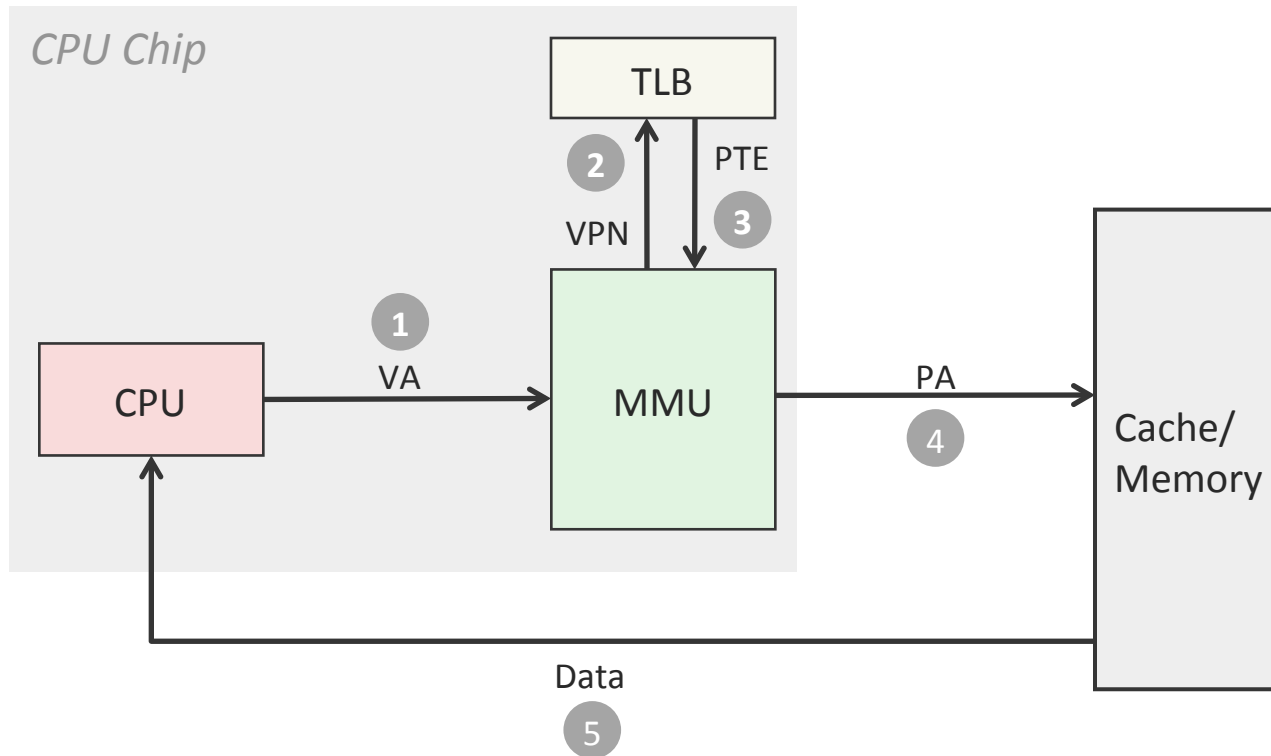
- Yes, it would be… so, real MMUs don't

# Speeding up translation with 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, dedicated, super-fast hardware cache of PTEs in MMU
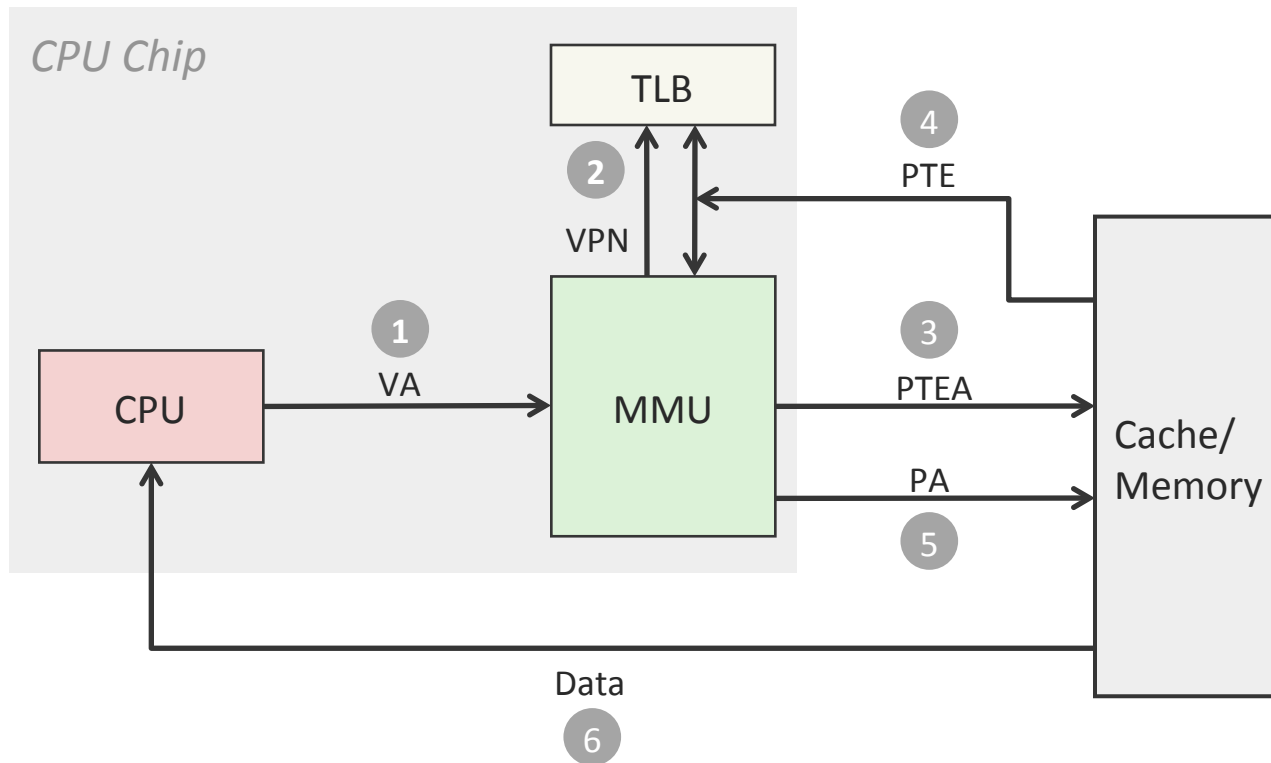  - Contains complete page table entries for small number of pages

# TLB hit



**A TLB hit eliminates a memory access**

# TLB miss



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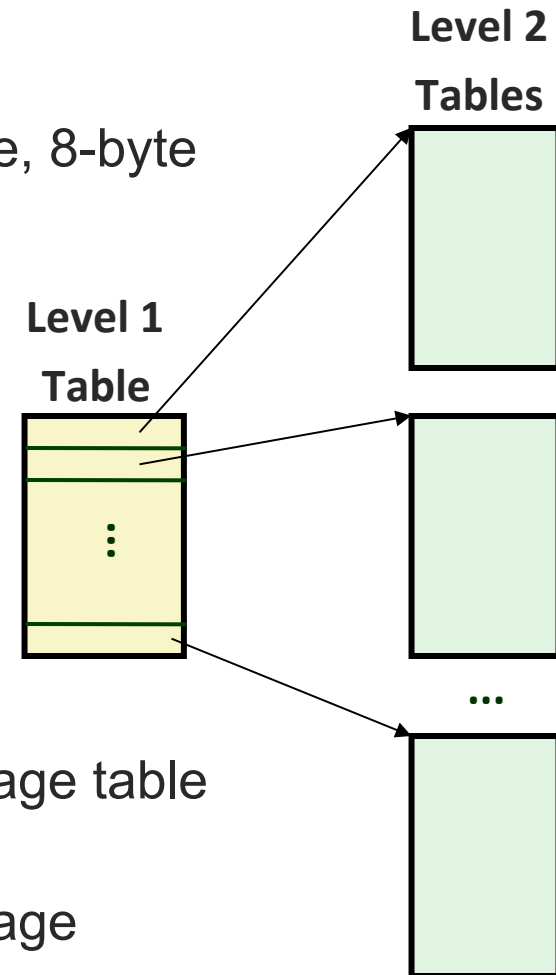
Fortunately, TLB misses are rare. Why?

# Question 2

- Isn't the page table huge?  How can it be stored in RAM?

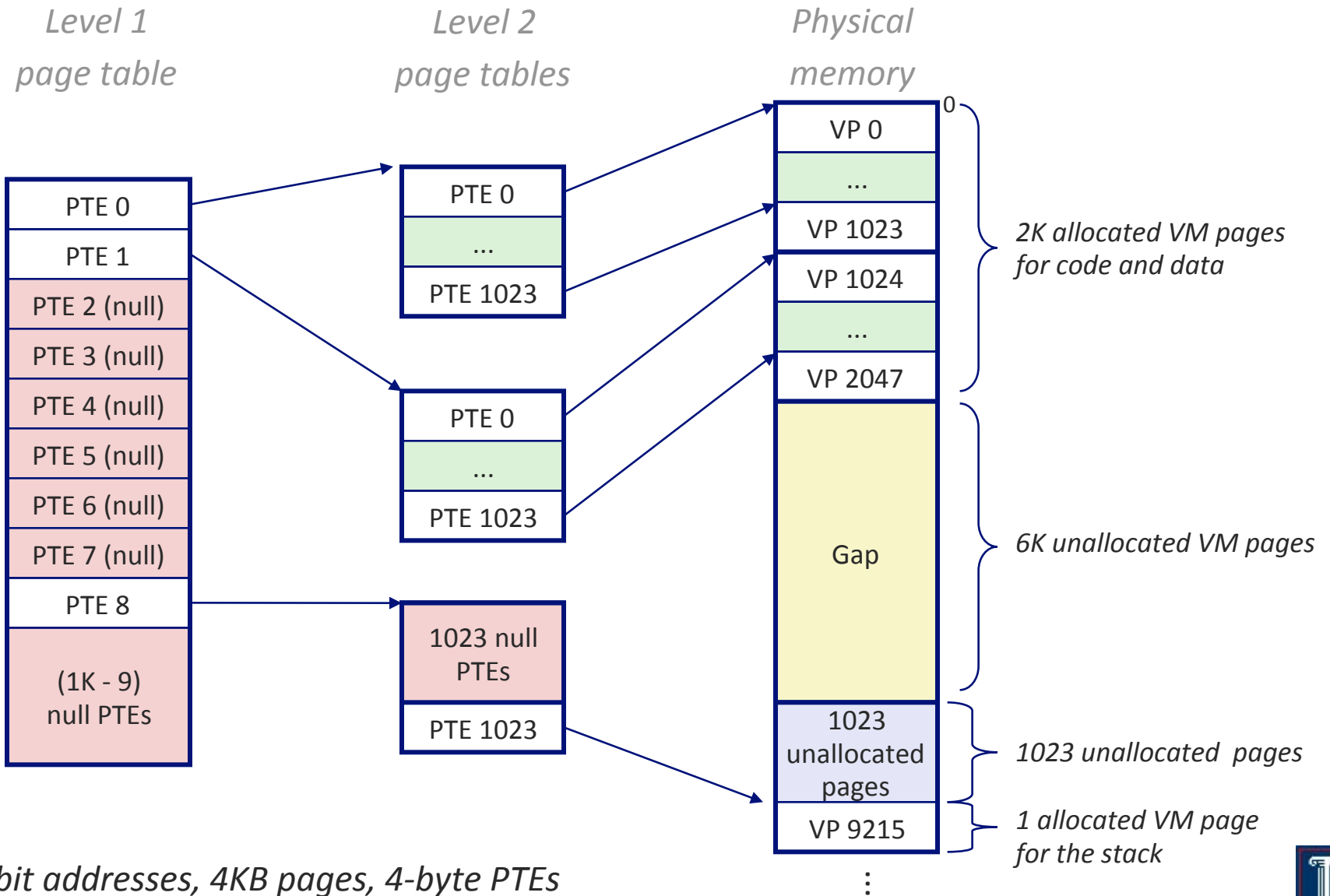- Yes, it would be… so, real page tables aren't simple arrays

# Multi-Level Page Tables

- ## Suppose:
  - ○ 4KB ($2^{12}$) page size, 64-bit address space, 8-byte PTE

- ## Problem:
  - ○ Would need a 32,000 TB page table!
  - ○ $2^{64} * 2^{-12} * 2^3 = 2^{55}$ bytes

- ## Common solution:
  - ○ Multi-level page tables
  - ○ Example: 2-level page table
    - ■ Level 1 table: each PTE points to a page table (always memory resident)
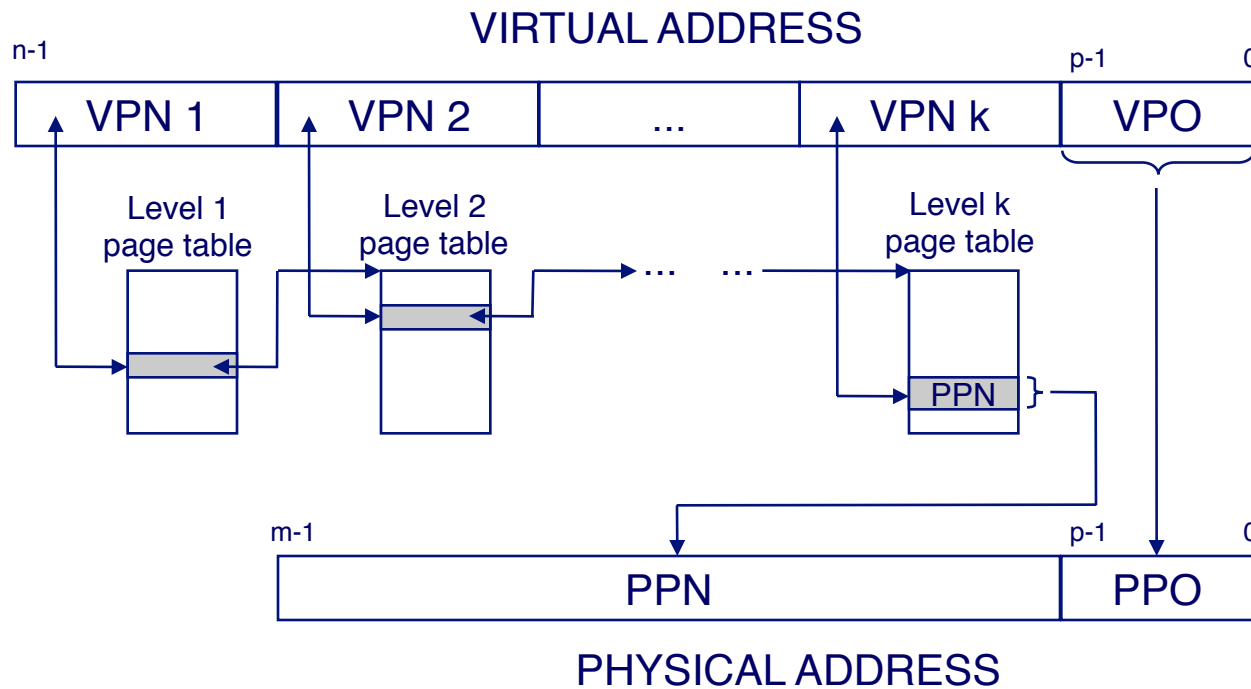    - ■ Level 2 table: each PTE points to a page (paged in and out like any other data)

**Level 2 Tables**

**Level 1 Table**

# 2-level page table hierarchy

**Level 1 page table**

**Level 2 page tables**

**Physical 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 |

| VP 0 | 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*

# Addr. translation with k-level PT

# Multilevel Page Tables

- With two levels of page tables, how big is each table?
    - Say we allocate 10 bits to the primary page, 10 bits to the secondary page, 12 bits to the page offset
    - Primary page table is then $2^{10}$ * 4 bytes per PTE == 4 KB
    - Secondary page table is also 4 KB
    - Hey ... that's exactly the size of a page on most systems ... cool

- What happens on a page fault?
    - MMU looks up index in primary page table to get secondary page table
    - MMU tries to access secondary page table
        - May result in another page fault to load the secondary table!
    - MMU looks up index in secondary page table to get physical frame #
    - CPU can then access physical memory address

- Issues
    - Page translation has very high overhead
        - Up to three memory accesses plus two disk I/Os!!
    - TLB usage is clearly very important

# Problem (from Tanenbaum)

- Suppose:
  - 32-bit address
  - Two-level page table
  - Virtual addresses split into a 9-bit top-level page table field, an 11-bit second-level page table field, and an offset

- Question: How large are the pages and how many are there in the address space?

# Problem (from Tanenbaum)

- Suppose:
  - 32-bit address
  - Two-level page table
  - Virtual addresses split into a 9-bit top-level page table field, an 11-bit second-level page table field, and an offset

- Question: How large are the pages and how many are there in the address space?
  - Offset is 12 bits
  - Page size $2^{12}$ bytes = 4KB
  - # Virtual pages = $(2^{32} / 2^{12}) = 2^{20}$
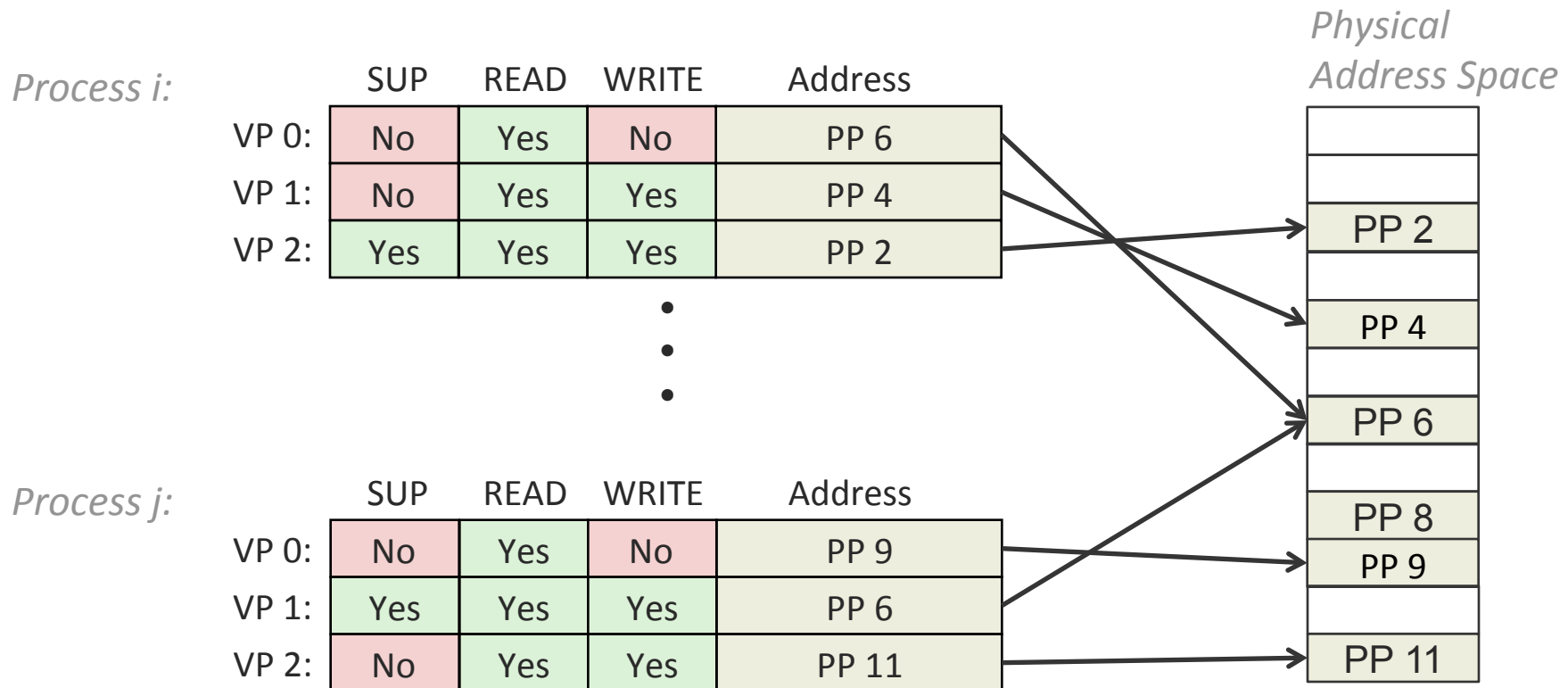  - Note: driven by number of bits in offset
    - Independent of size of top and 2nd level

# Question 3

- Is there any other super slick stuff can I do with page tables?


- Yes!

# Paging as a tool for protection

- Extend PTEs with permission bits
- Page fault handler checks these before remapping
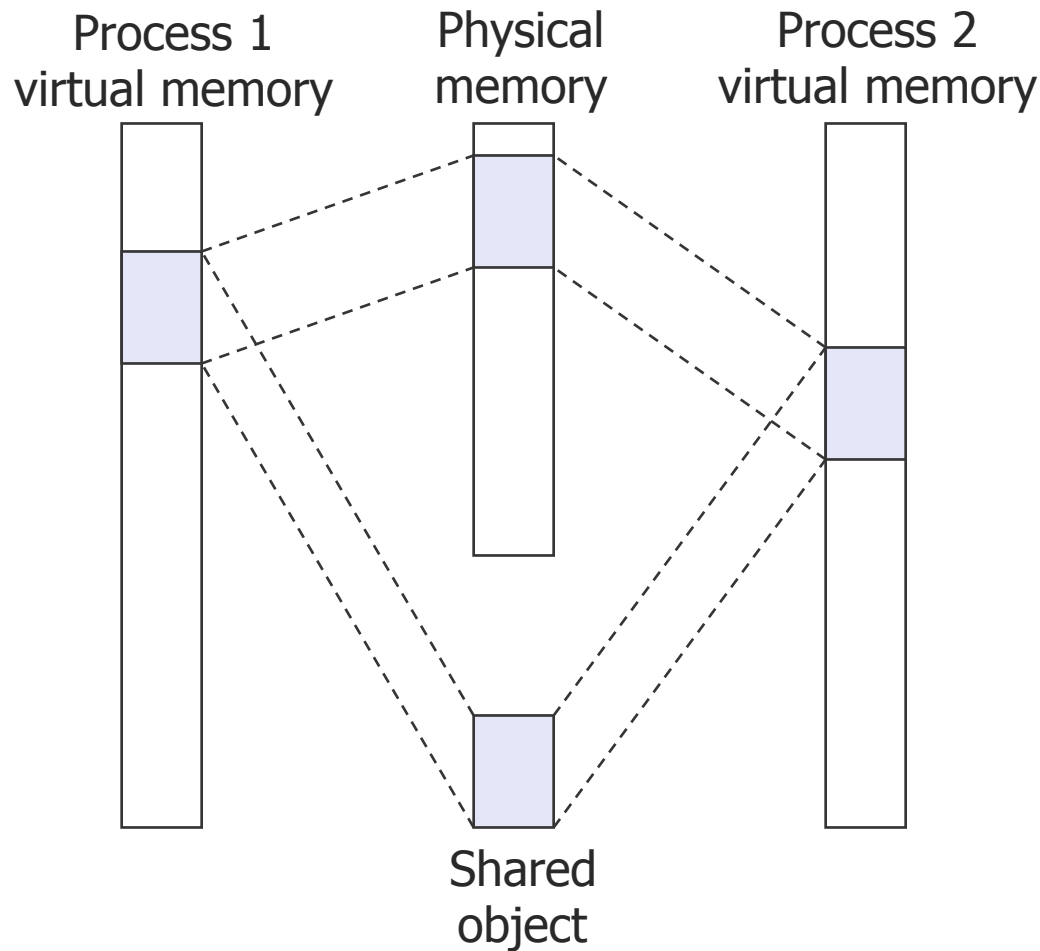    - If violated, send process SIGSEGV (segmentation fault)

*Process i:*

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

*Process j:*

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

*Physical Address Space*

| |
|---|
| |
| |
| PP 2 |
| |
| PP 4 |
| |
| PP 6 |
| |
| PP 8 |
| PP 9 |
| |
| PP 11 |

# VM as a tool for sharing

Process 1
virtual memory

Physical
memory

Process 2
virtual memory

Shared
object

- Process 1 maps the shared object.

# VM as a tool for sharing



Process 1 virtual memory

Physical memory

Process 2 virtual memory

Shared object

- Process 2 maps the shared object.
- Notice how the virtual addresses can be different.

# Protection + sharing example

- `fork()` creates exact copy of a process
  - Lots more on this next week…

- When we fork a new process, does it make sense to make a copy of all of its memory?
  - Why or why not?

- What if the child process doesn't end up touching most of the memory the parent was using?
  - `exec()` replaces a process with a new one
  - Extreme example *and common case*: What happens if a process does an `exec()` immediately after `fork()`?

# Copy-on-write

- Idea: Give the child process access to the same memory, but don't let it write to any of the pages directly!
  - 1) Parent forks a child process
  - 2) Child gets a copy of the parent's page tables
    - They point to the same physical frames!!!

# Copy-on-write

- All pages (both parent and child) marked read-only
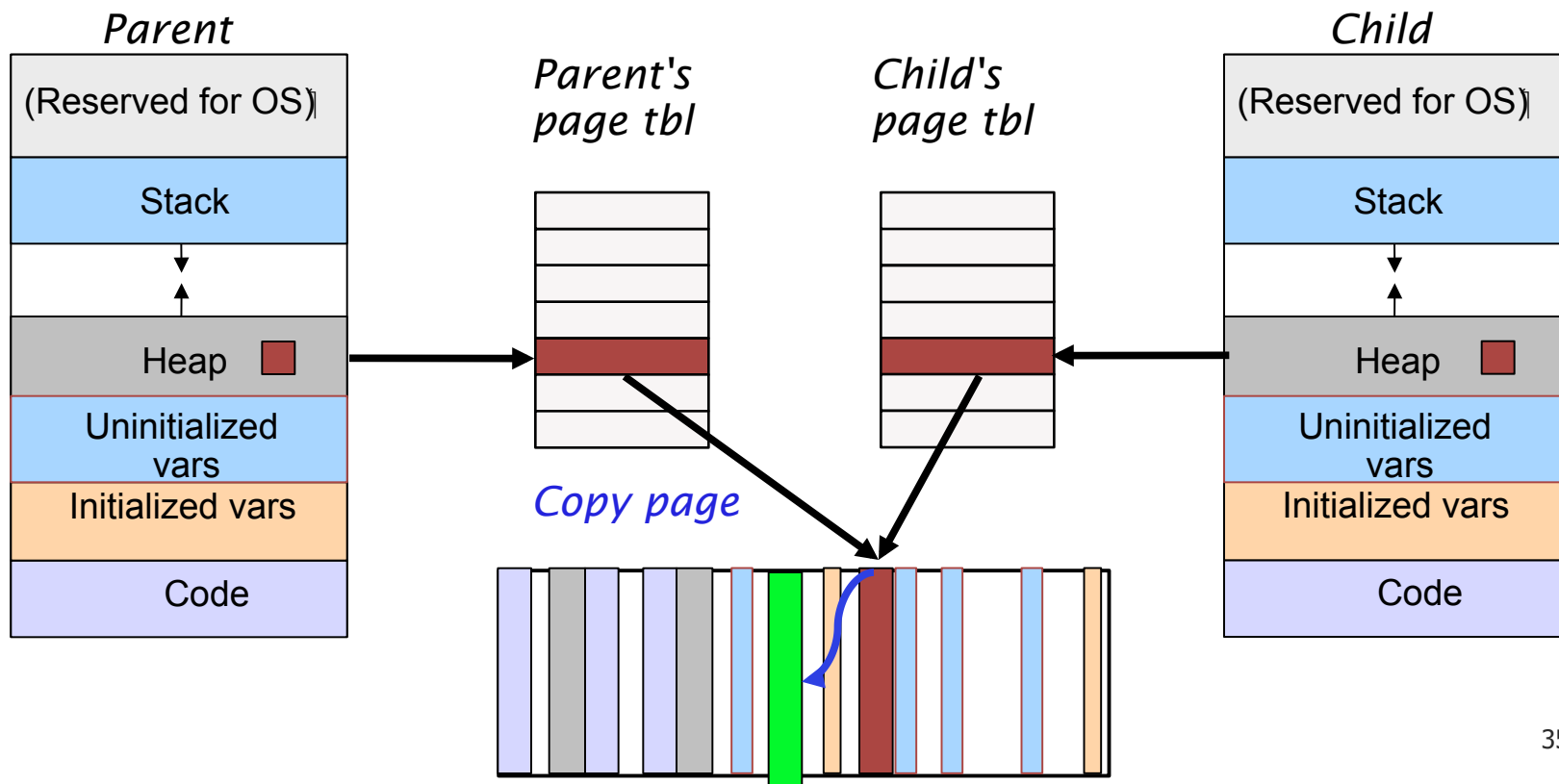  - **Why?**

Parent

| (Reserved for OS) |
| Stack |
| |
| Heap |
| Uninitialized vars |
| Initialized vars |
| Code |

*Parent's page tbl*

*Child's page tbl*

Child

| (Reserved for OS) |
| Stack |
| |
| Heap |
| Uninitialized vars |
| Initialized vars |
| Code |

# Copy-on-write

- What happens when the child *reads* the page?
  - Just accesses same memory as parent .... niiiiiice
- What happens when the child *writes* the page?
  - Protection fault occurs (page is read-only!)
  - OS copies the page and maps it R/W into the child's addr space
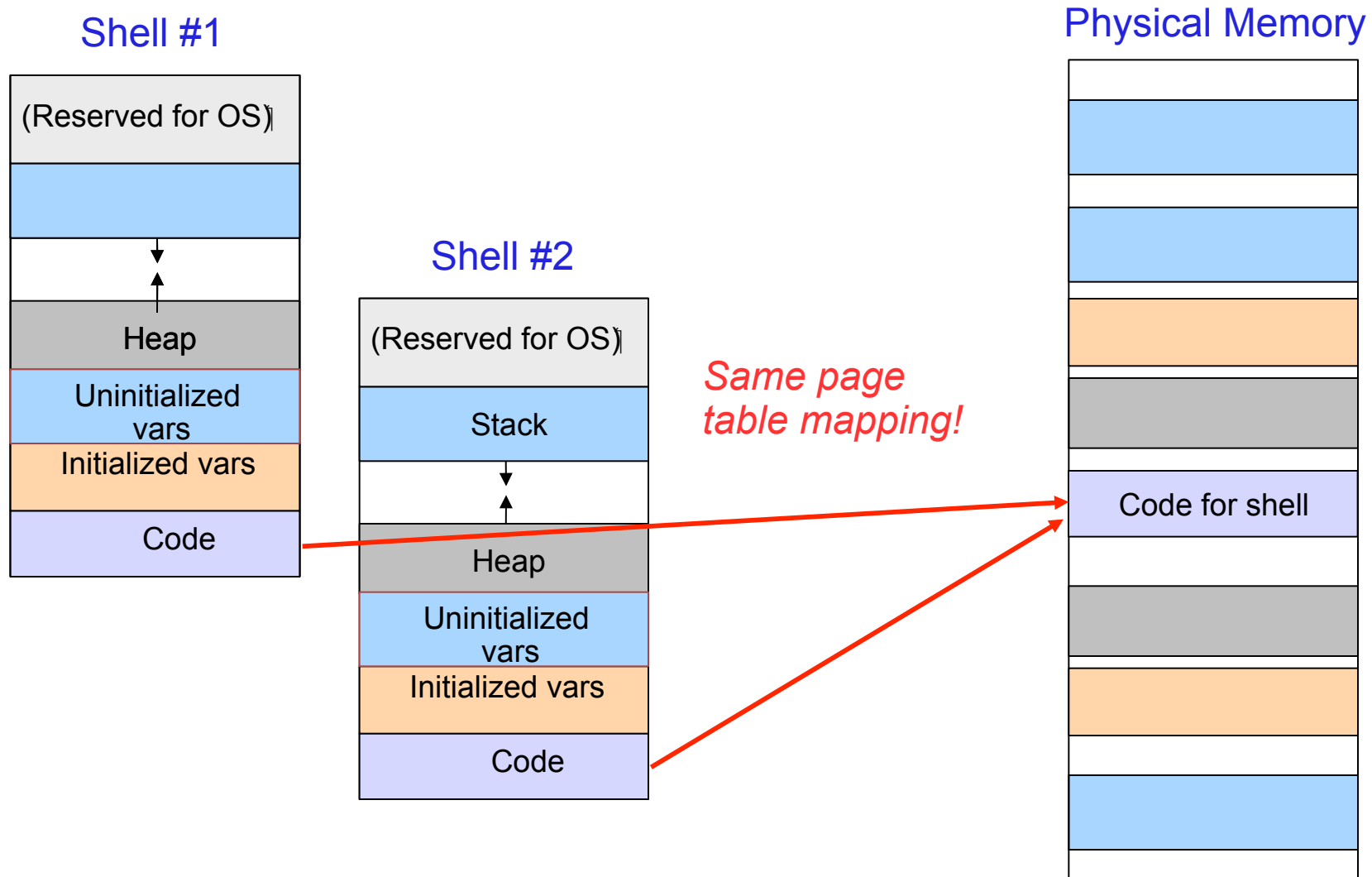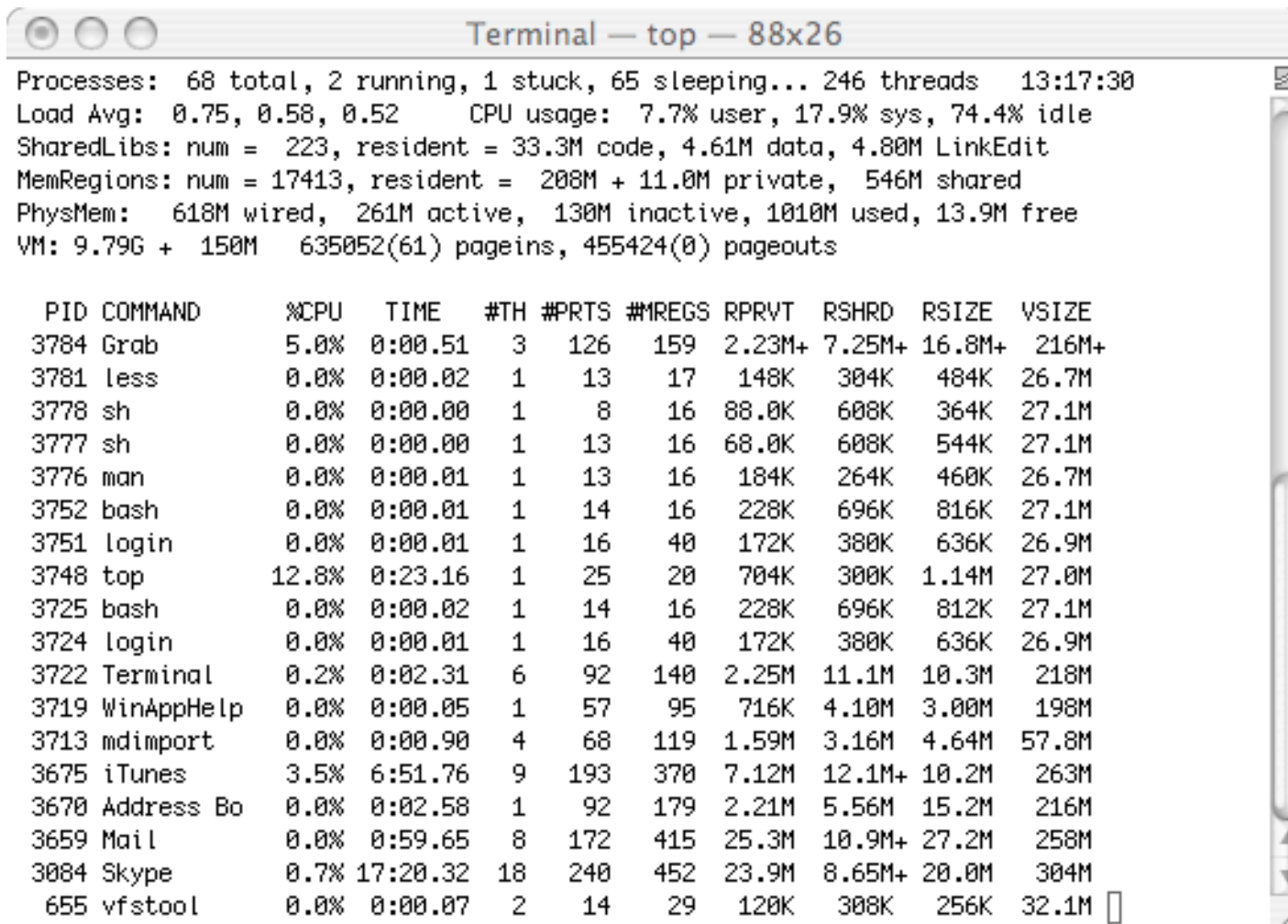
| *Parent* | *Parent's page tbl* | *Child's page tbl* | *Child* |
|---|---|---|---|
| (Reserved for OS) | | | (Reserved for OS) |
| Stack | | | Stack |
| | | | |
| Heap ■ | | | Heap ■ |
| Uninitialized vars | | | Uninitialized vars |
| Initialized vars | | | Initialized vars |
| Code | | | Code |

# Copy-on-write

- What happens when the child *reads* the page?
  - Just accesses same memory as parent .... niiiiiice
- What happens when the child *writes* the page?
  - Protection fault occurs (page is read-only!)
  - OS copies the page and maps it R/W into the child's addr space

# Copy-on-write

- What happens when the child *reads* the page?
  - Just accesses same memory as parent .... niiiiiice
- What happens when the child *writes* the page?
  - Protection fault occurs (page is read-only!)
  - OS copies the page and maps it R/W into the child's addr space

# Another sharing example

- Can also share code segment

**Shell #1**

| (Reserved for OS) |
| --- |
| (blue) |
| (white, with up/down arrows) |
| Heap |
| Uninitialized vars |
| Initialized vars |
| Code |

**Shell #2**

| (Reserved for OS) |
| --- |
| Stack |
| (white, with up/down arrows) |
| Heap |
| Uninitialized vars |
| Initialized vars |
| Code |

*Same page table mapping!*

**Physical Memory**

| |
| --- |
| (blue) |
| (blue) |
| (orange) |
| (gray) |
| Code for shell |
| (gray) |
| (orange) |
| (blue) |

# Benefits of sharing pages

- How much memory savings do we get from sharing pages across identical processes?
  - A lot! Use the "top" command...

# Summary

- Paging implementation
  - Basics: get page off disk if necessary (*page fault*) and then map virtual to physical address
  - Problem: Mapping requires extra memory access (solution?)
  - Problem: Page table can get huge (solution?)

- Paging enables flexible use of memory
  - Protection
  - Sharing (e.g., copy-on-write defers writes as long as possible)
  - Caching
    - Q: How do I choose which page to evict when swapping?