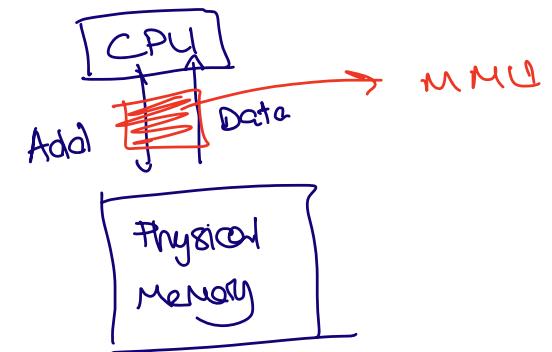


process is program in execution.

Virtual Memory: Concepts

CENG331 - Computer Organization

2^{48} bytes = 256 Tera byte → each process has this much space
 $2^8 \times 2^{40}$



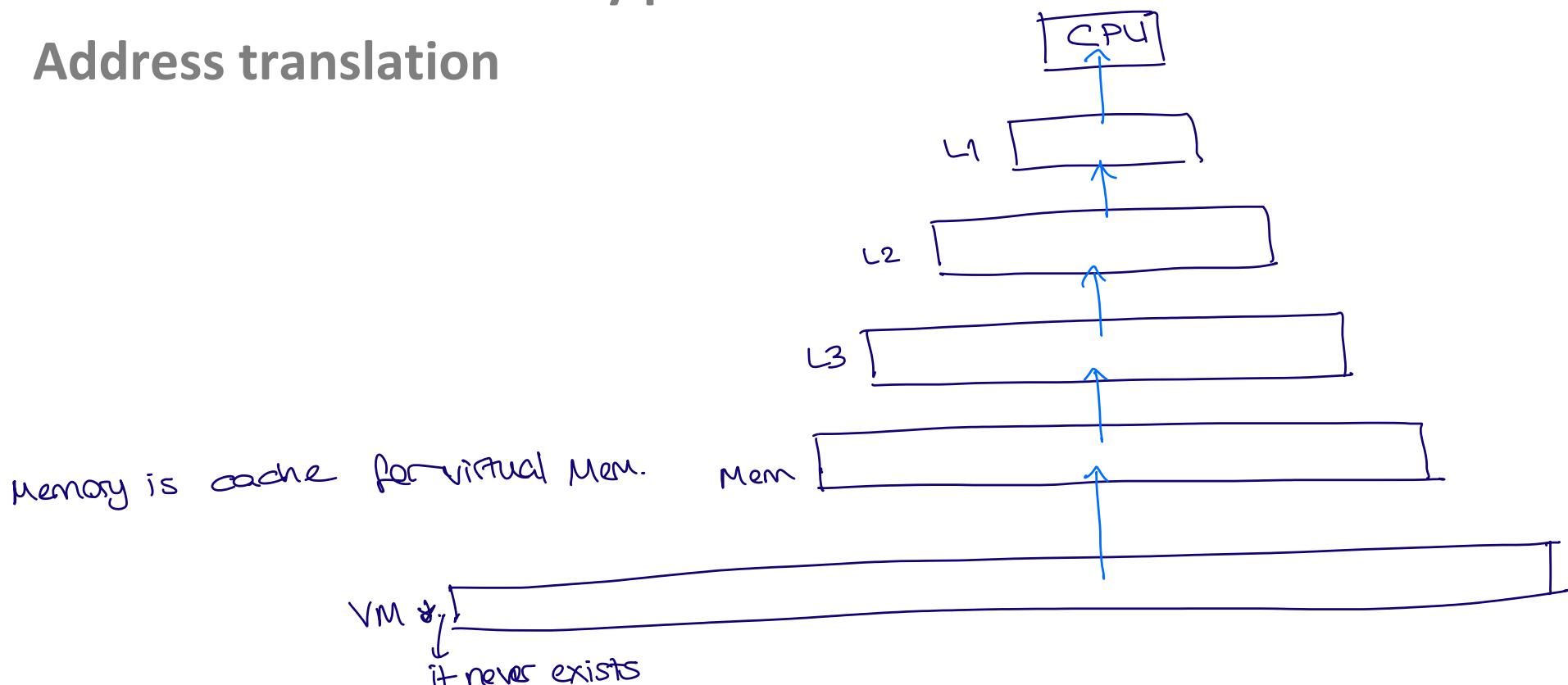
The world doesn't exist all at once, it will be partially exist when you step into that specific post

Adapted from slides of the textbook: <http://csapp.cs.cmu.edu/>

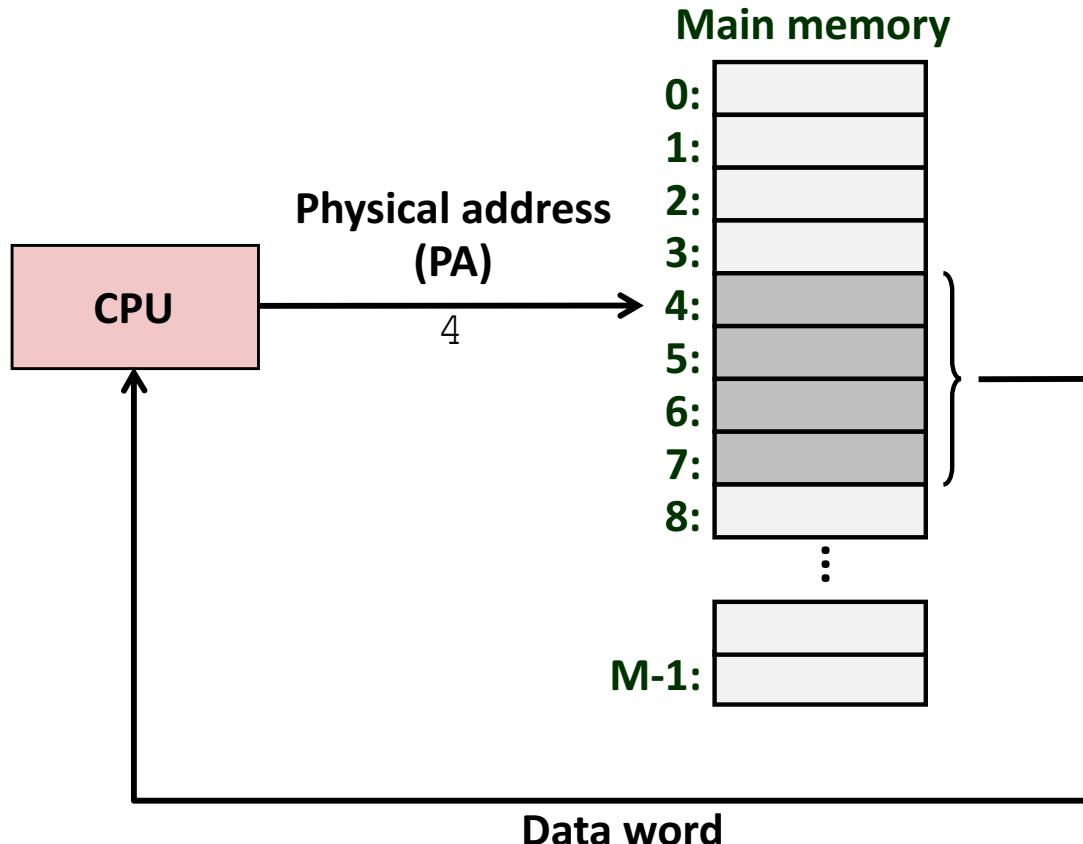
movq %.rax, 0x7fabcede

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

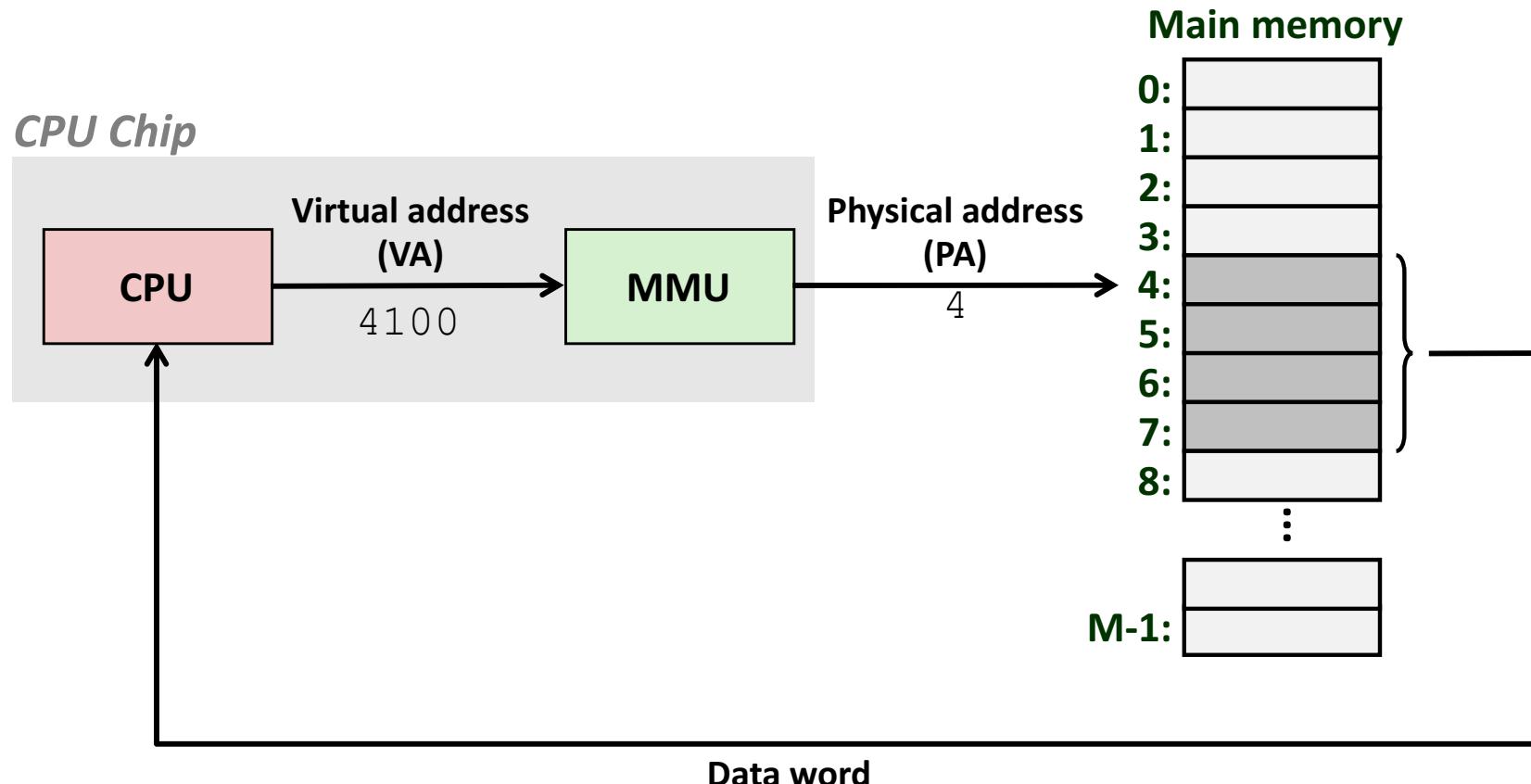


A System Using Physical Addressing



- Used in “simple” systems like embedded microcontrollers in devices like cars, 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

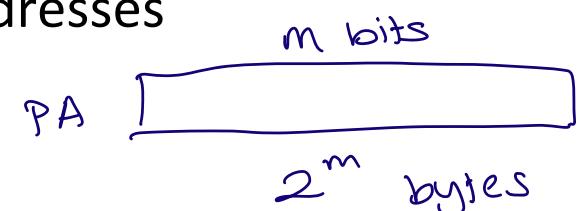
- **Linear address space:** Ordered set of contiguous non-negative integer addresses:

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

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

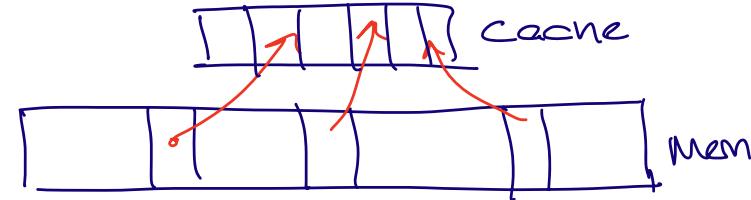
$$\{0, 1, 2, 3, \dots, N-1\}$$

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

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


- Clean distinction between data (bytes) and their attributes (addresses)
- Every byte in main memory has one physical address and zero or more virtual addresses

Why Virtual Memory (VM)?



■ Uses main memory efficiently

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

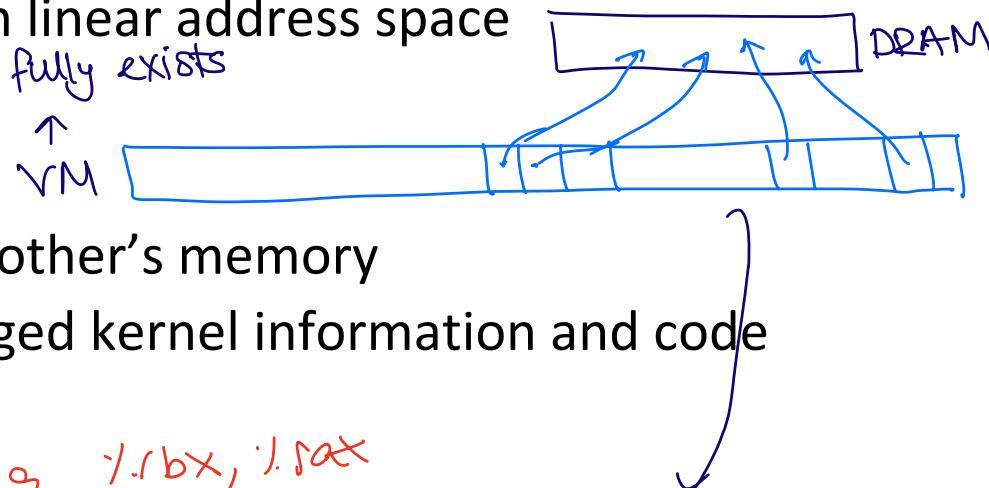
■ Simplifies memory management

- Each process gets the same uniform linear address space

*get fine value ← never fully exists
on demand*

■ Isolates address spaces

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



C - exc → *moveq \$.rbx, \$.sat*
C2 - exc

VM is used as a cache for DRAM

Today

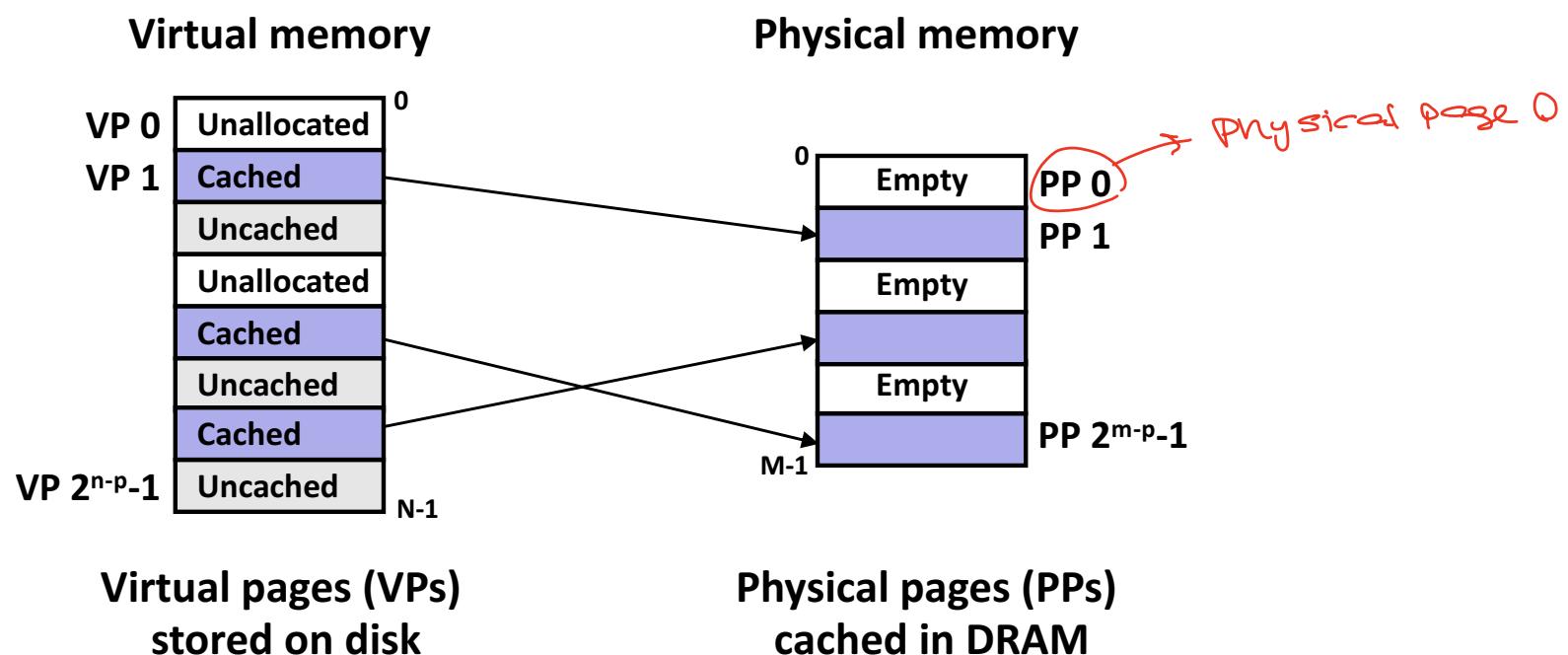
- Address spaces
- **VM as a tool for caching**
- VM as a tool for memory management
- VM as a tool for memory protection
- Address translation

VM as a Tool for Caching

Cache
↓
Data Block

VM
↓
page

- 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)



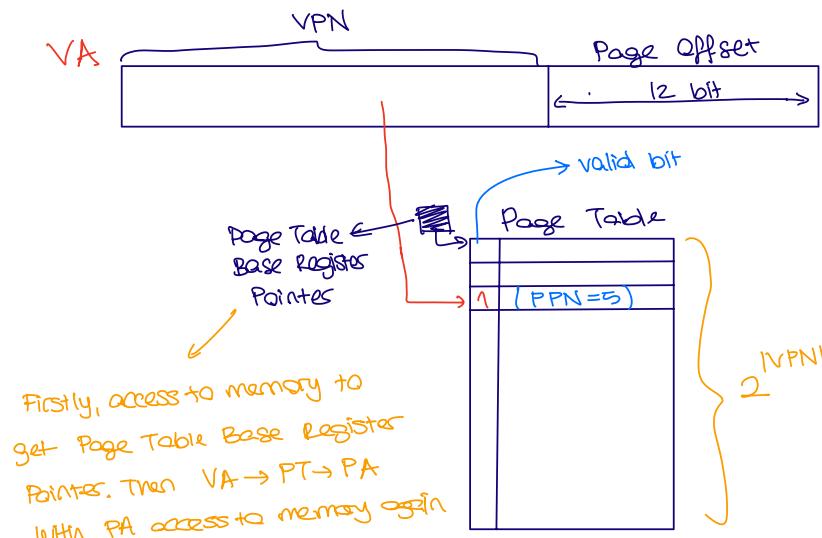
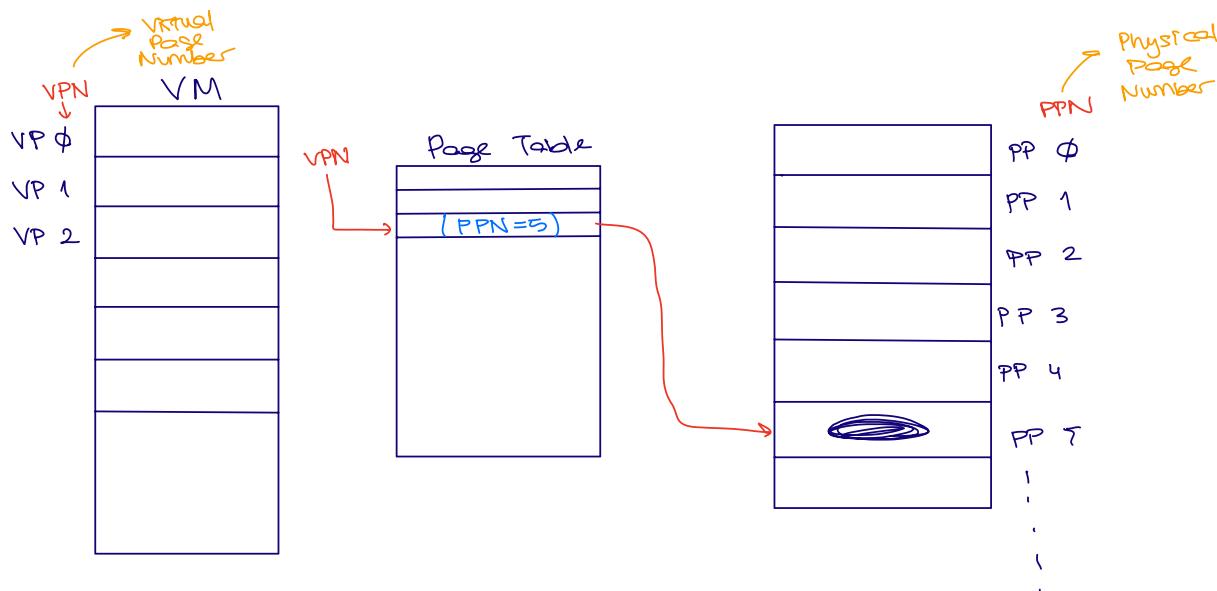
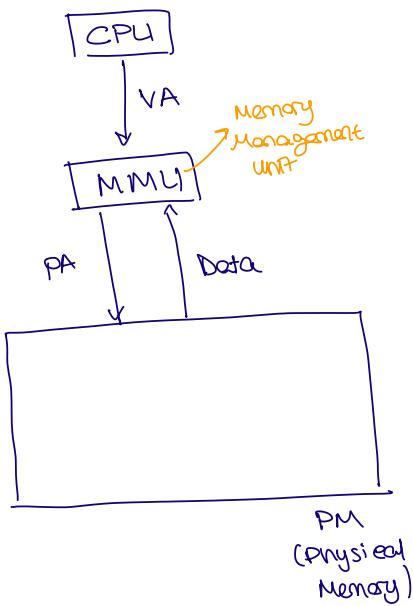
DRAM Cache Organization

■ DRAM cache organization driven by the enormous miss penalty

- DRAM is about **10x** slower than SRAM
- Disk is about **10,000x** slower than DRAM

■ Consequences

- Large page (block) size: typically 4-8 KB, sometimes 4 MB
- Fully associative
 - Any VP can be placed in any PP
 - Requires a “large” mapping function – different from CPU caches
- Highly sophisticated, expensive replacement algorithms
 - Too complicated and open-ended to be implemented in hardware
- Write-back rather than write-through

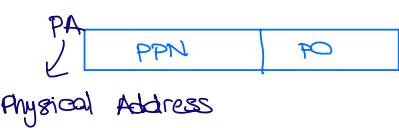


$$|Page| = 4 \text{ KB} = 2^2 \times 2^{10} = 2^{12}$$

$$|PO| = 12 \text{ bit}$$

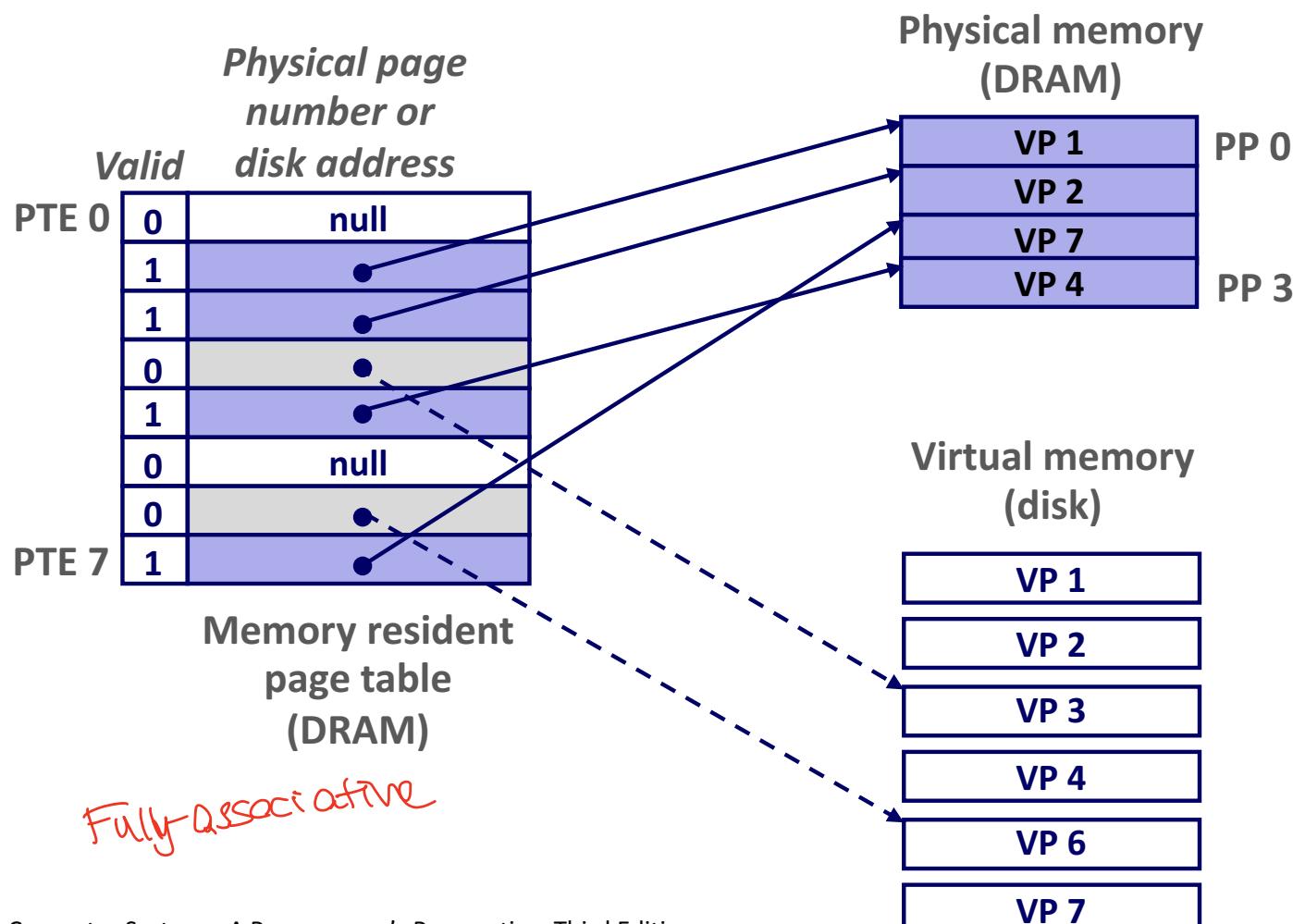
$$|VPN| = |VA| - |PO|$$

if valid bit == 1
copy PPN to the
beginning of the PO



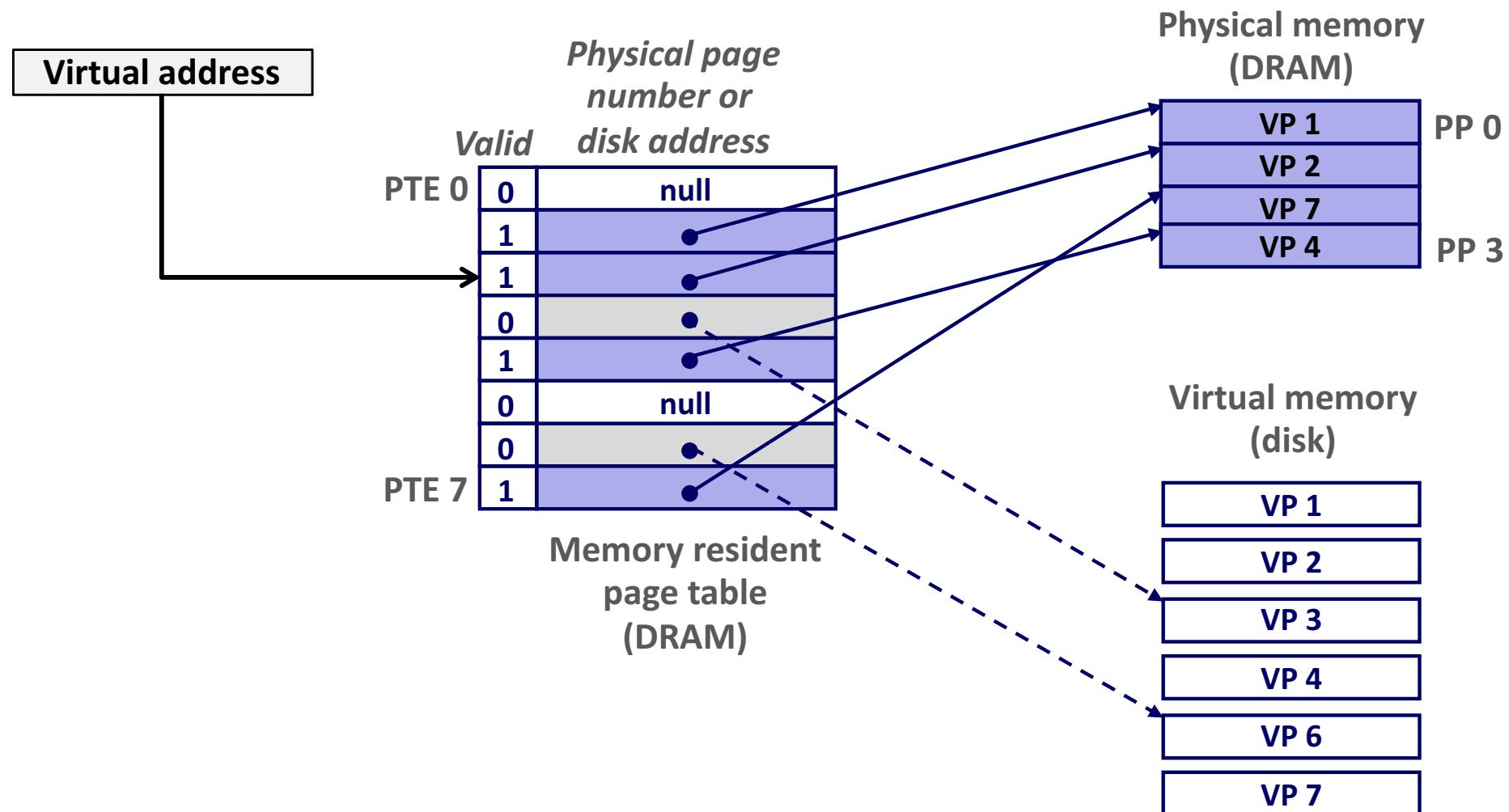
Enabling Data Structure: Page Table

- 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



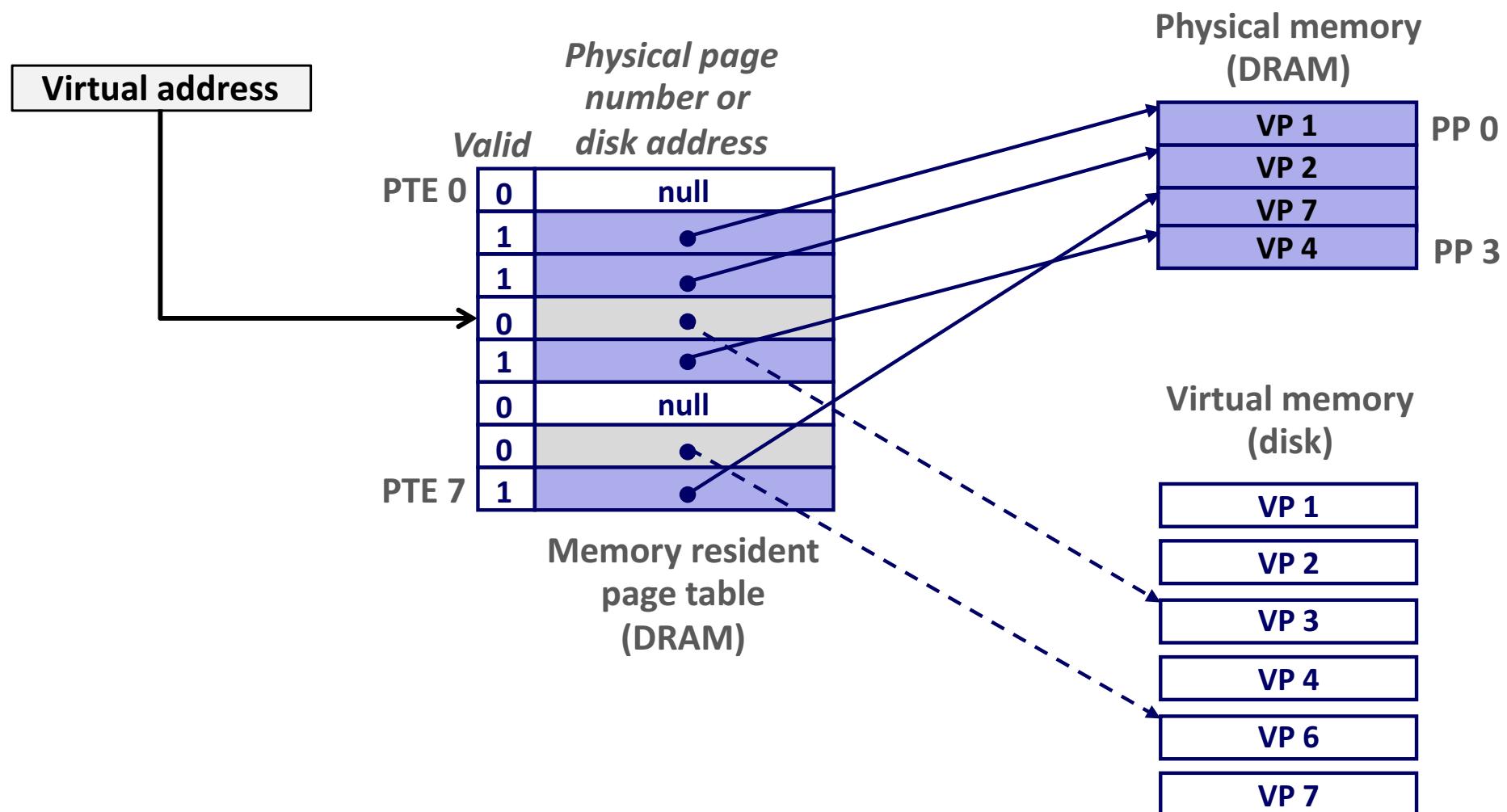
Page Hit

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



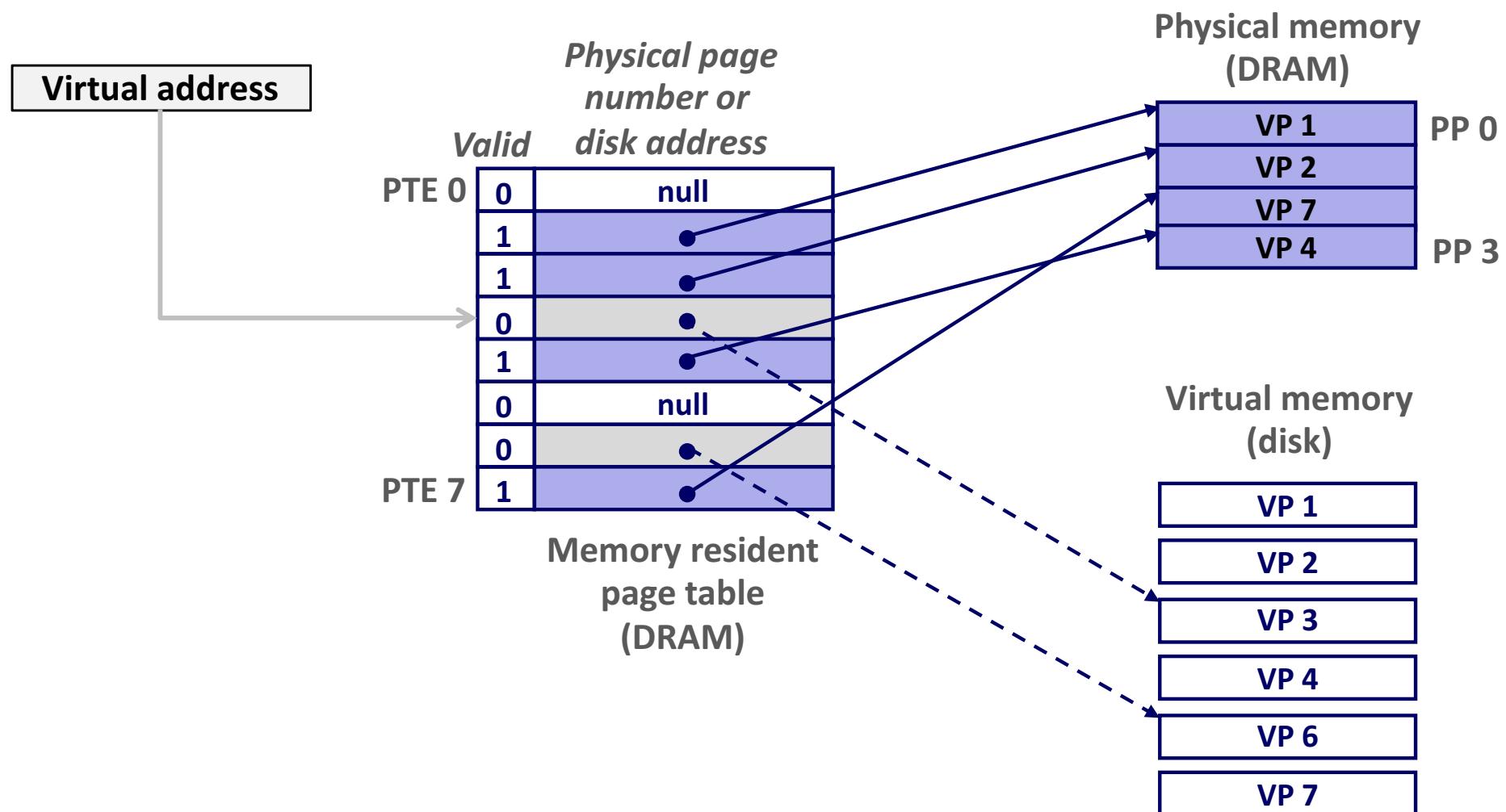
Page Fault

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



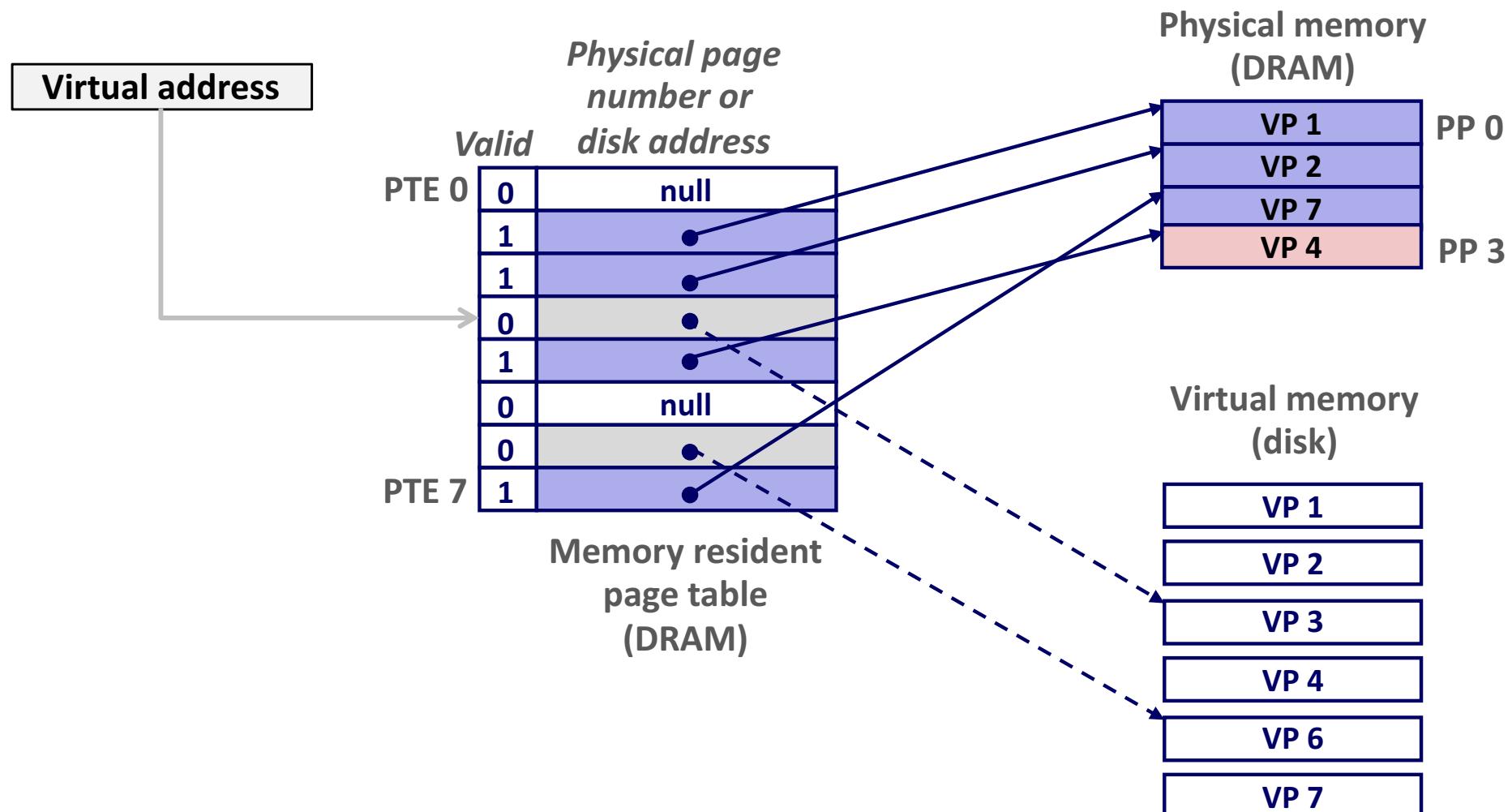
Handling Page Fault

- Page miss causes page fault (an exception)



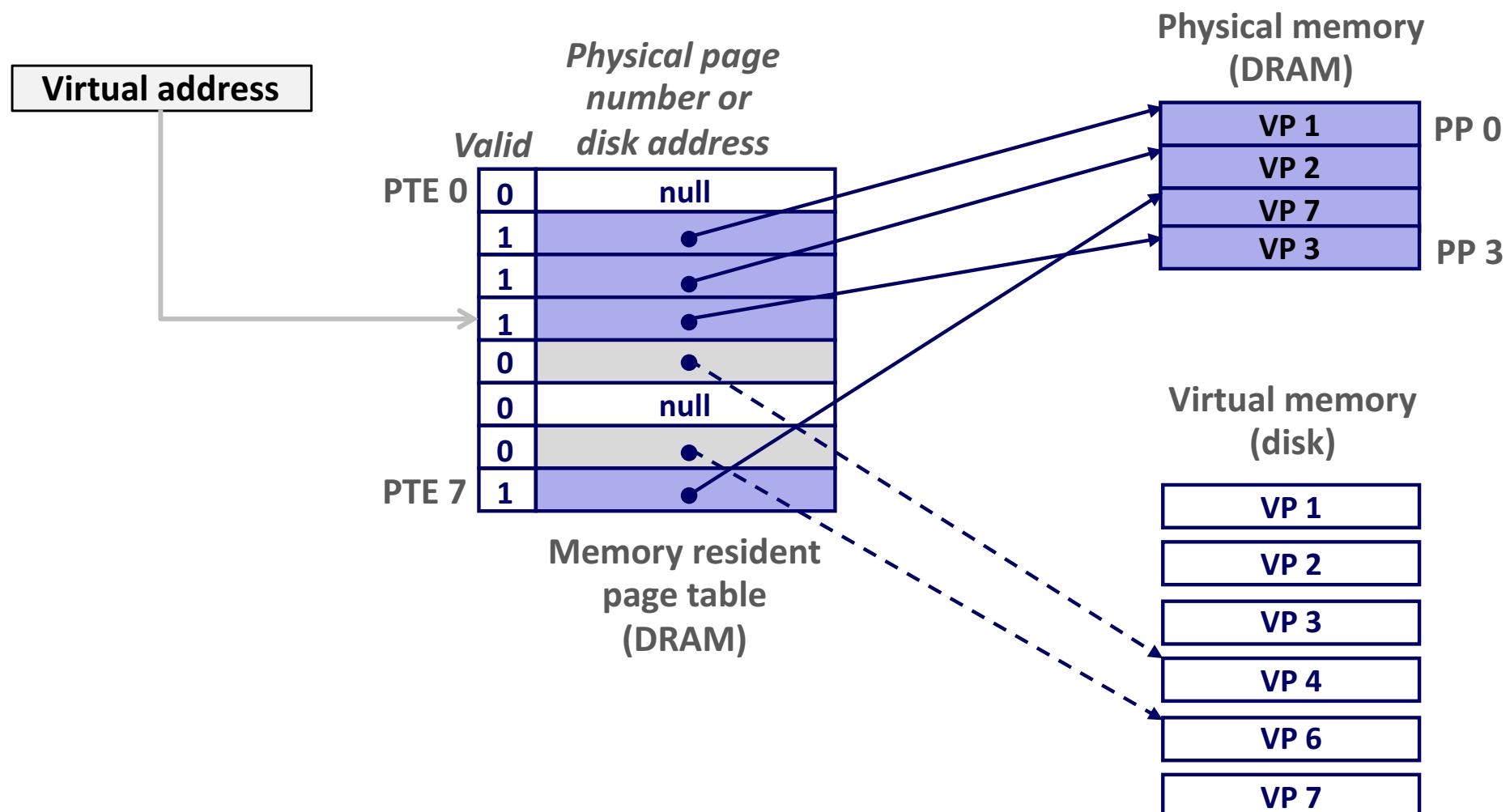
Handling Page Fault

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



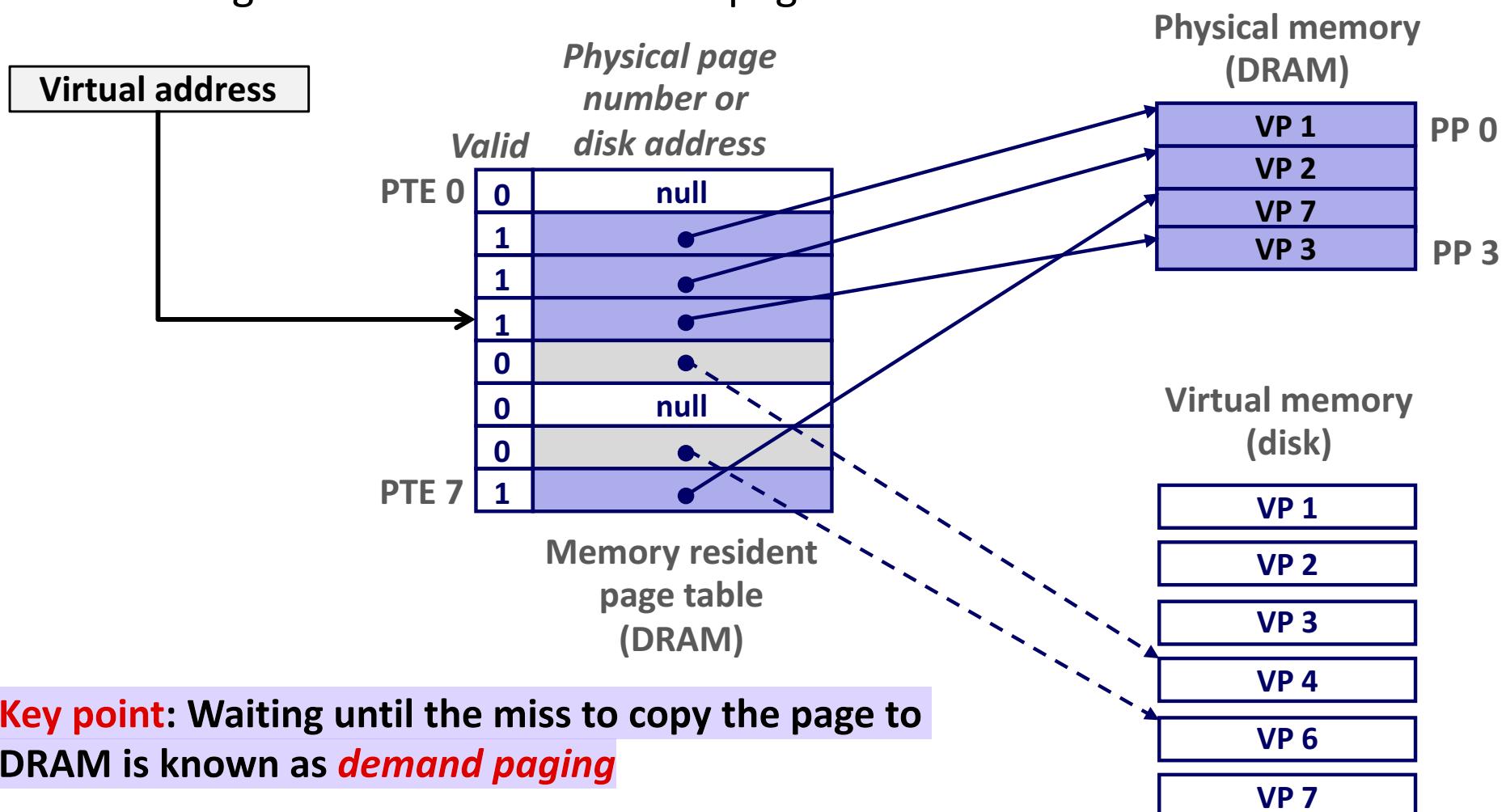
Handling Page Fault

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



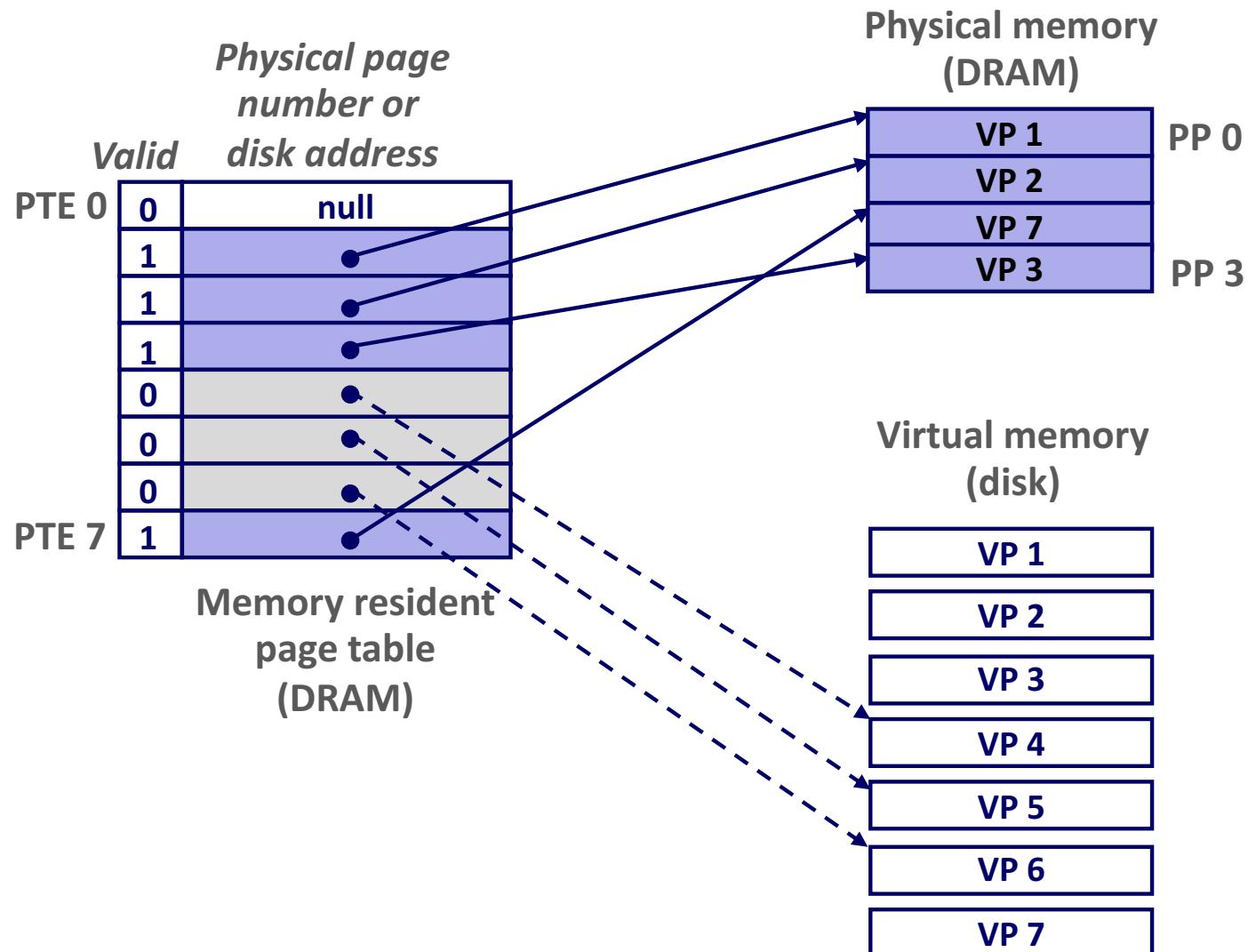
Handling Page Fault

- 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!



Allocating Pages

- Allocating a new page (VP 5) of virtual 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 compulsory misses
- If ($\text{SUM}(\text{working set sizes}) > \text{main memory size}$)
 - *Thrashing*: Performance meltdown where pages are swapped (copied) in and out continuously

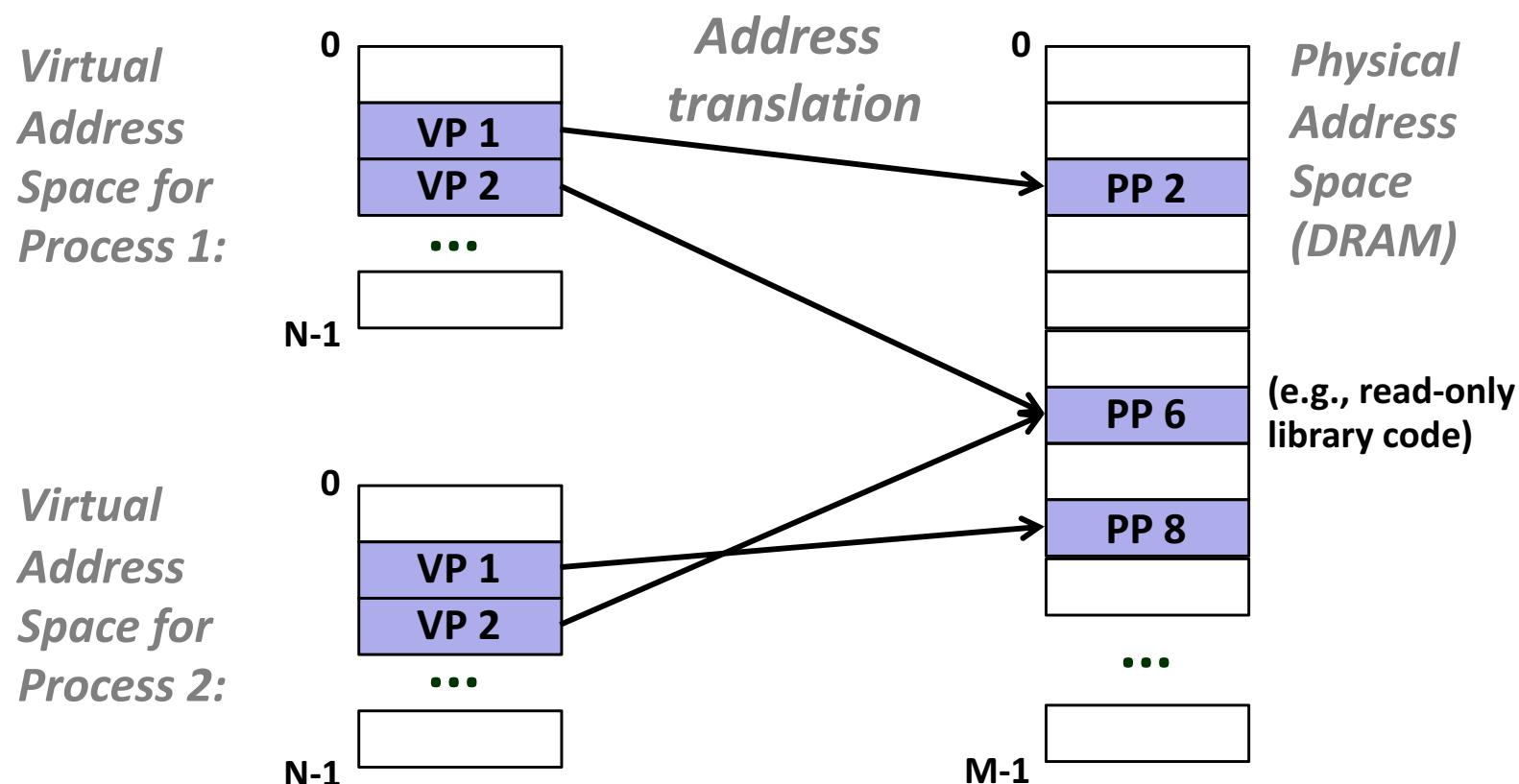
Today

- 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



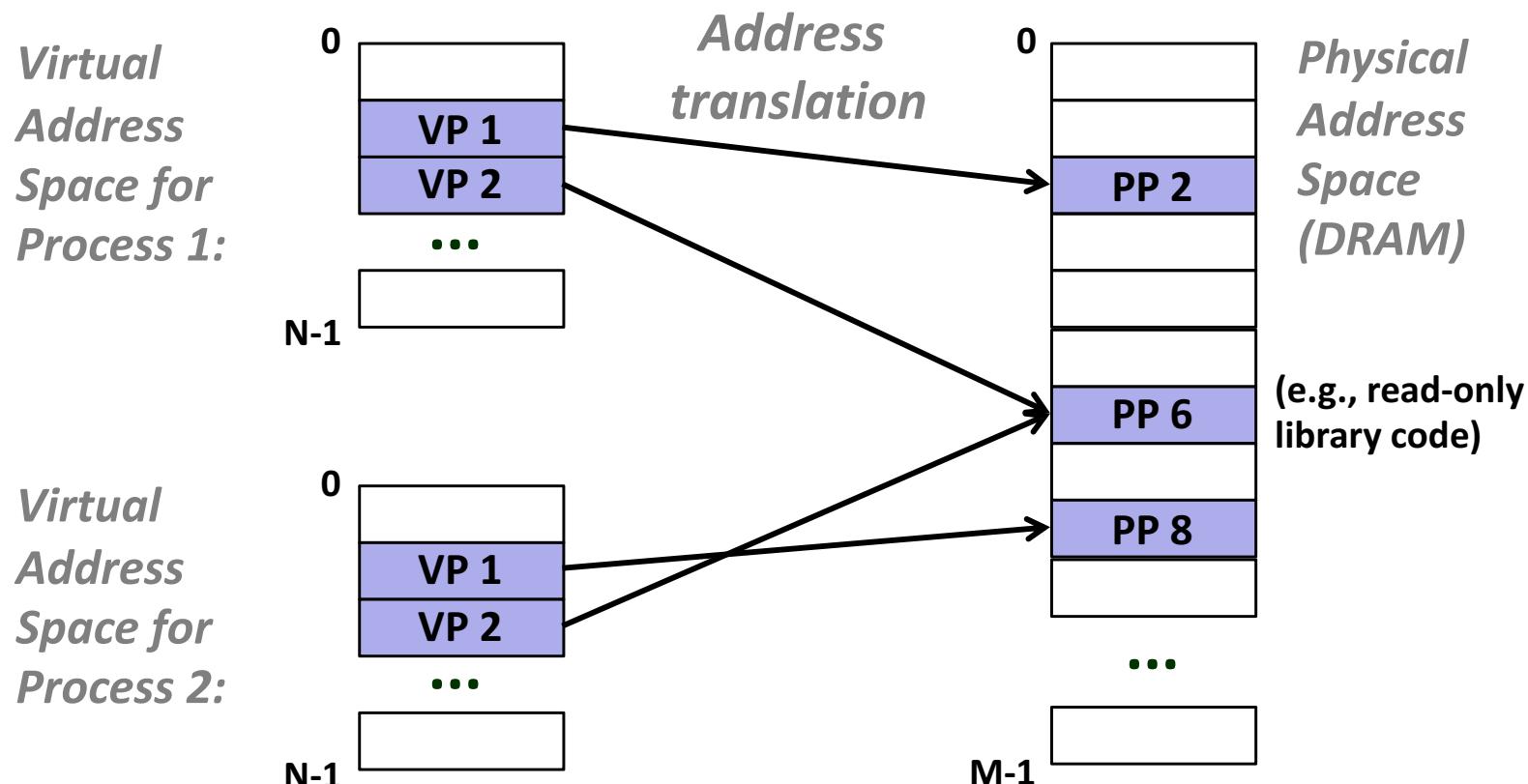
VM as a Tool for Memory Management

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)



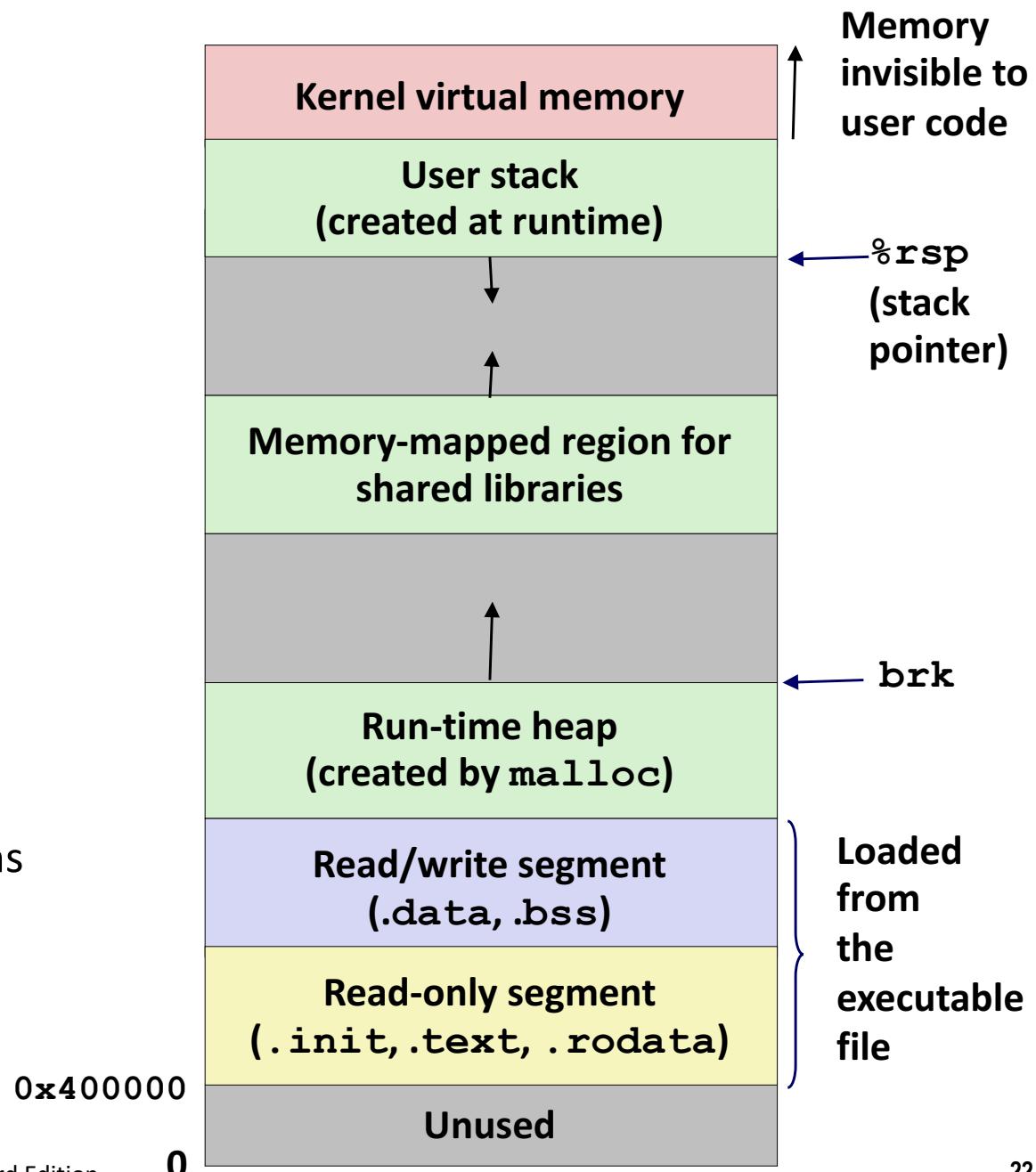
Simplifying Linking and Loading

■ Linking

- Each program has similar virtual address space
- Code, stack, and shared libraries always start at the same address

■ 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



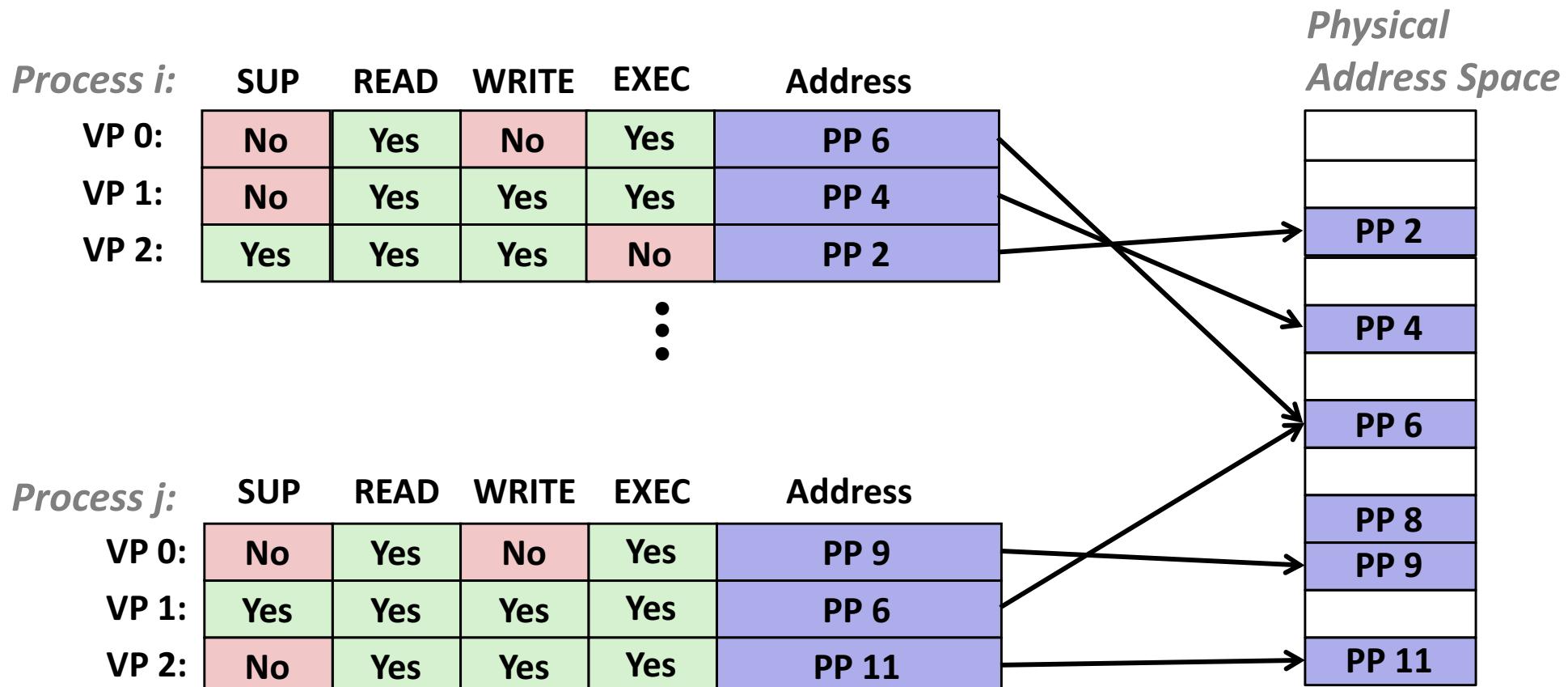
Today

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

SUP bit set to 1 indicates that the process must be running in SUPervisor mode to access. Processes running in user mode can only access pages with SUP set to 0.



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)

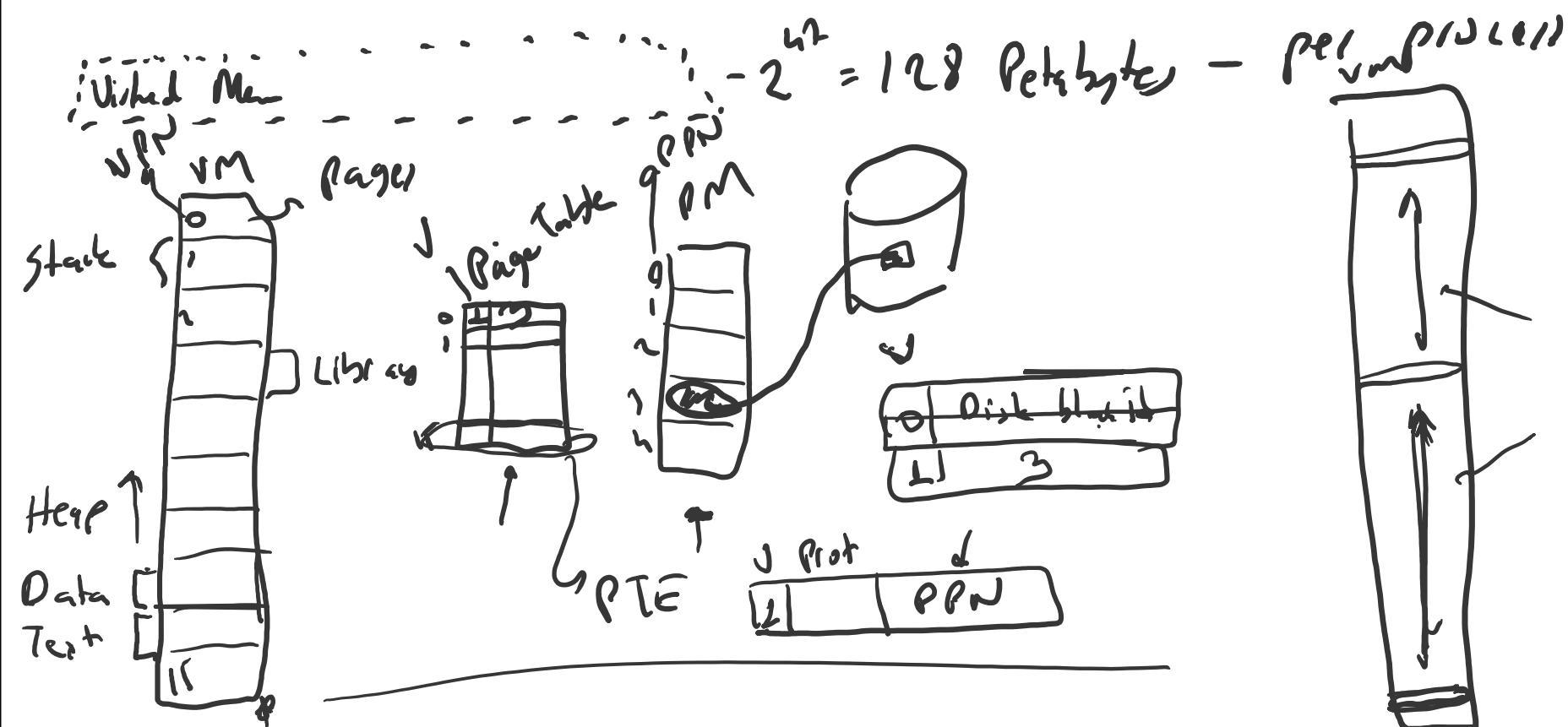
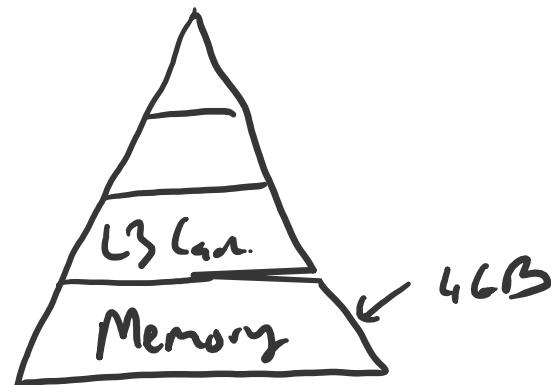
- **TLBI**: TLB index
- **TLBT**: TLB tag
- **VPO**: Virtual page offset
- **VPN**: Virtual page number

■ Components of the physical address (PA)

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

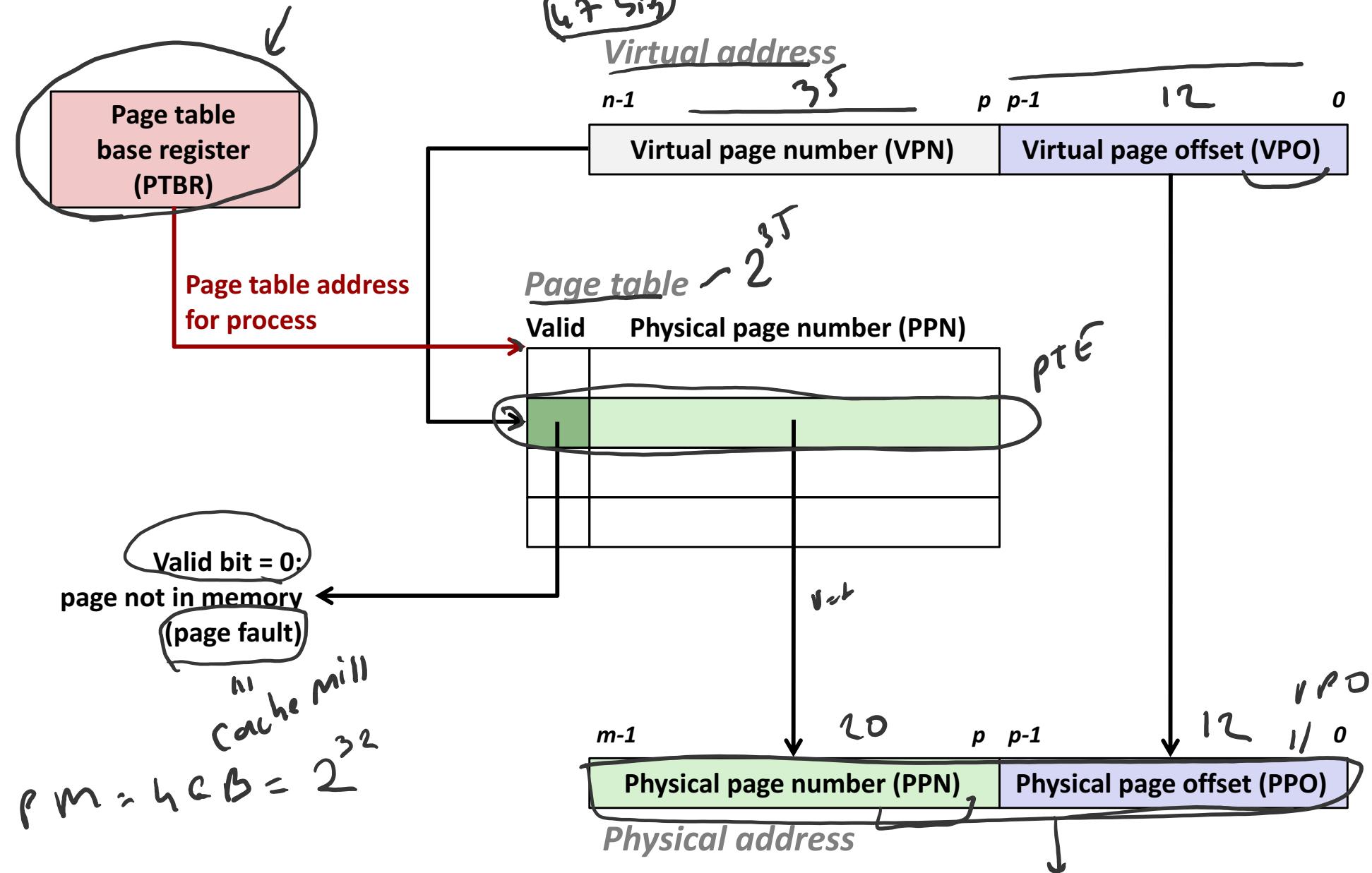
Virtual Memory

PT per process

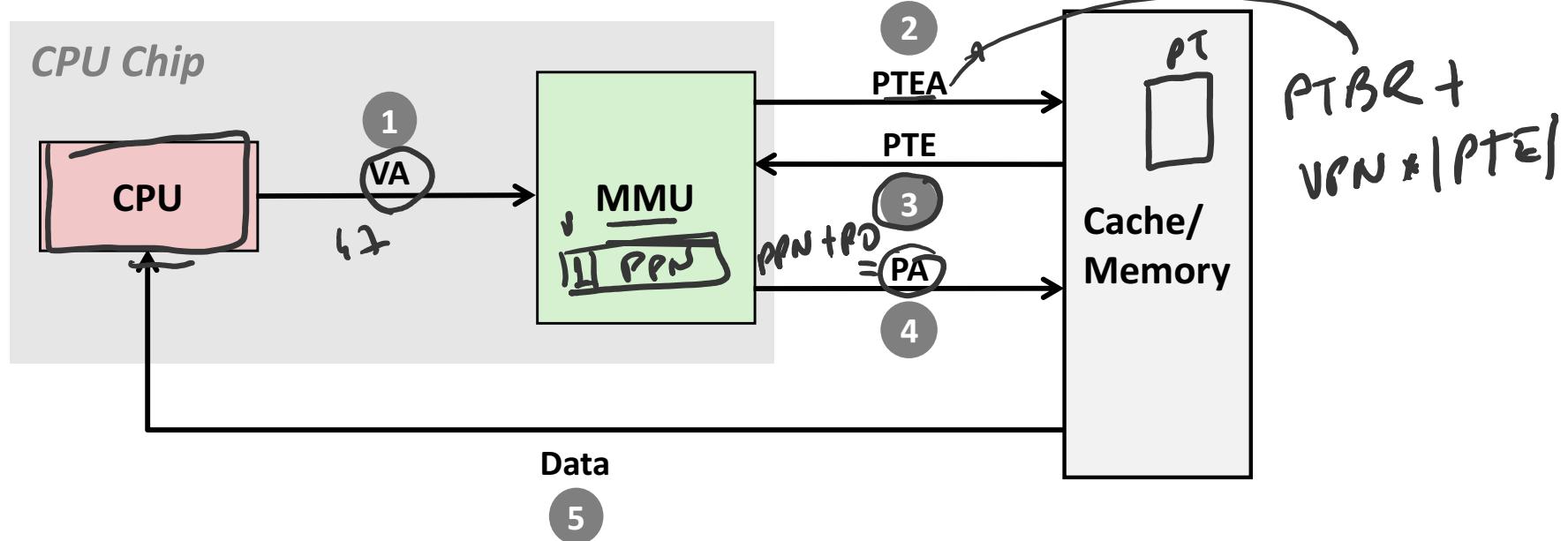


$$\text{Page size} = 4 \text{ KB} = 2^{12}$$

Address Translation With a Page Table

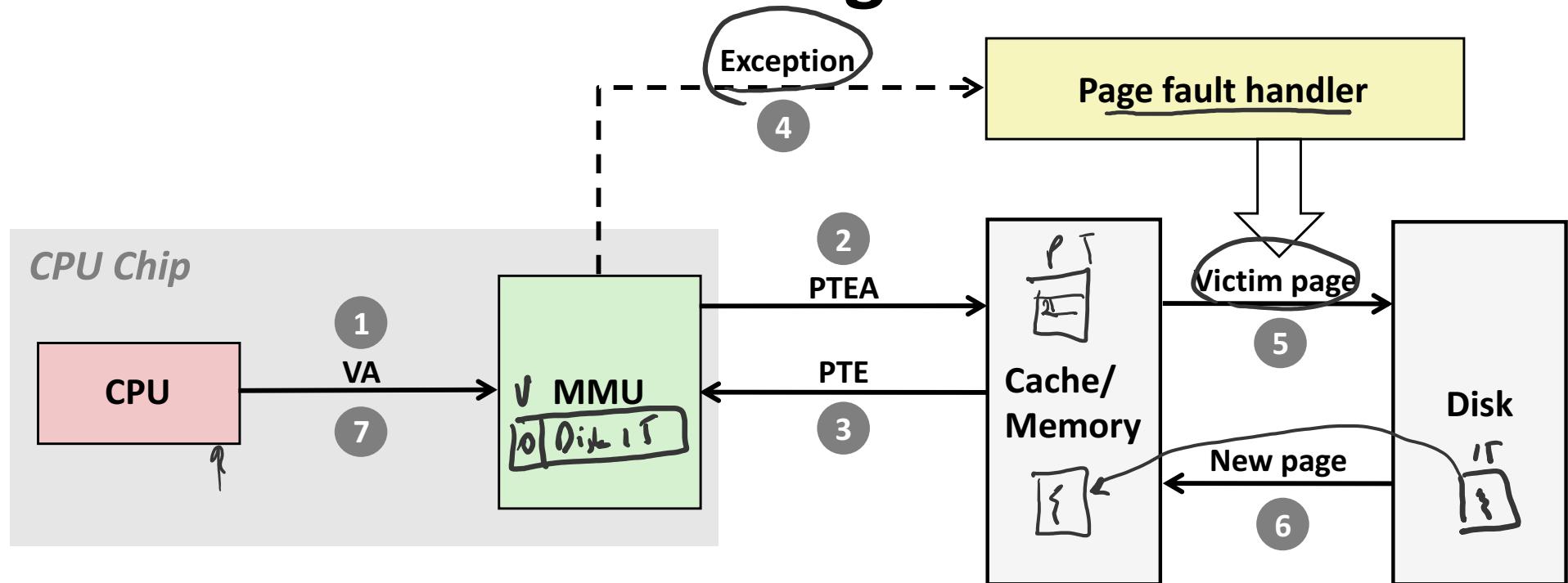


Address Translation: Page Hit



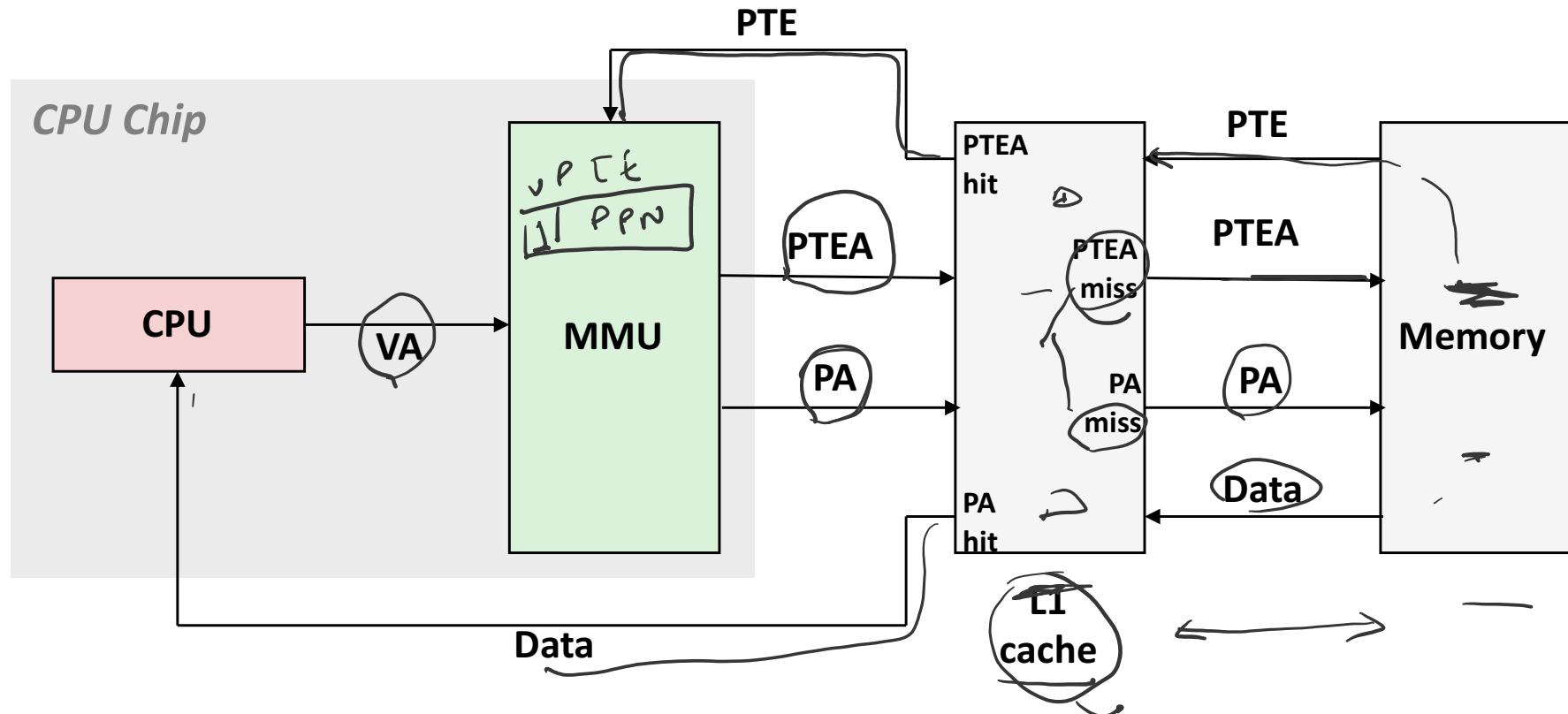
- 1) Processor sends virtual address to MMU
- 2-3) MMU fetches PTE from page table in memory
- 4) MMU constructs physical address
- 5) MMU sends physical address to cache/memory
- 5) Cache/memory sends data word to processor

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

Integrating VM and Cache



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

Speeding up Translation with a TLB

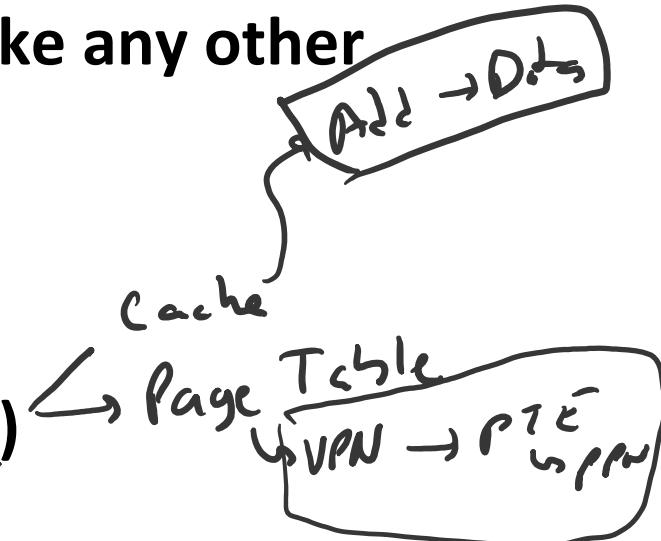
cache for
Page Table
every time we
shouldn't have
to get this

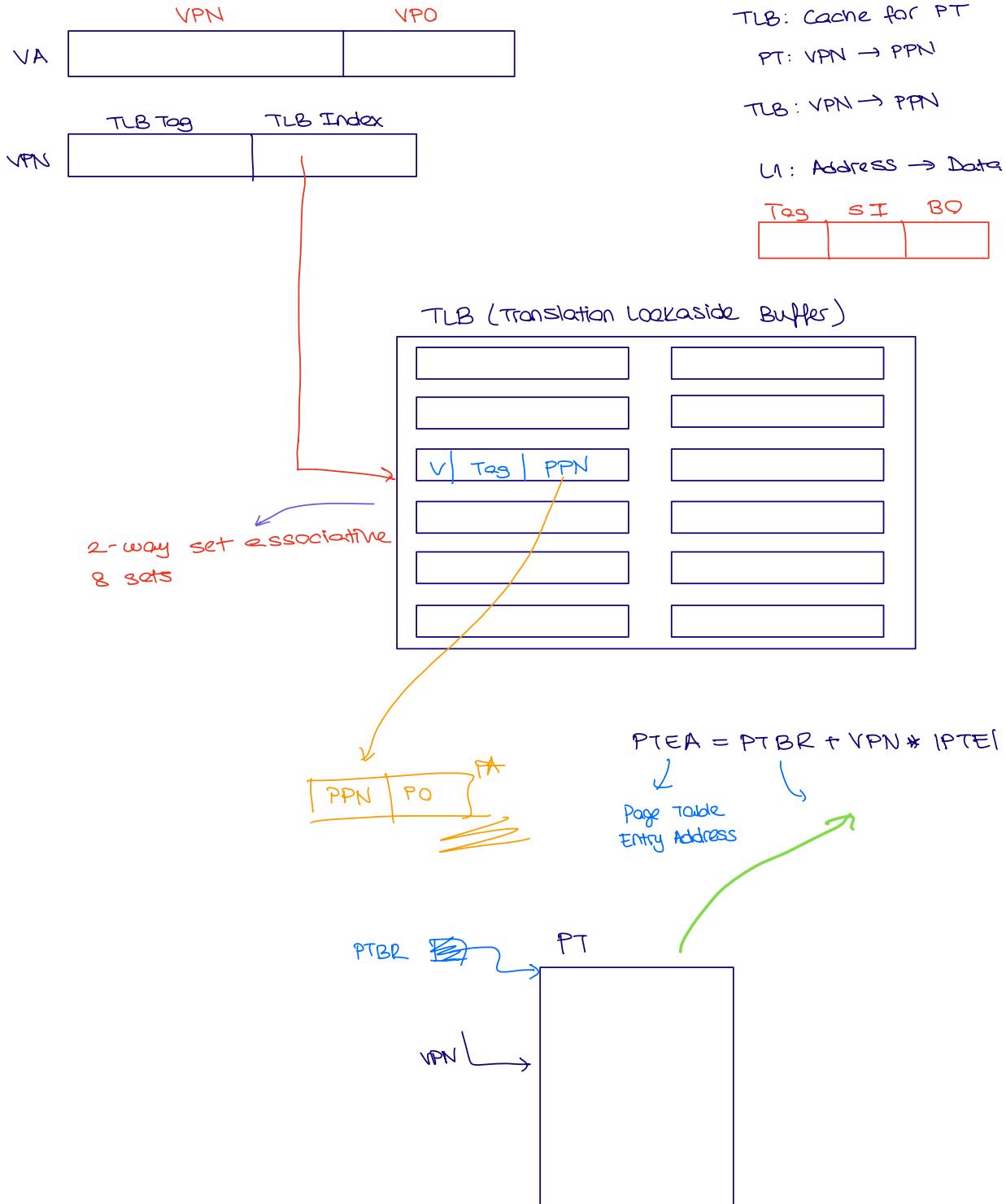
- Page table entries (PTEs) are cached in L1 like any other memory word

- 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





32-bit system IA32

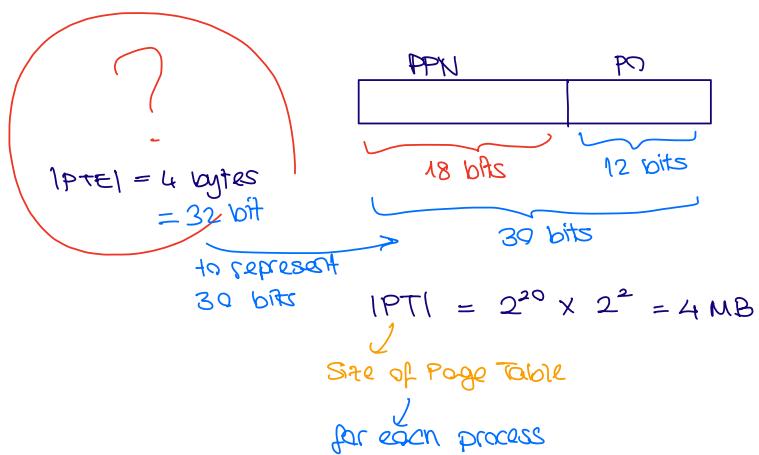
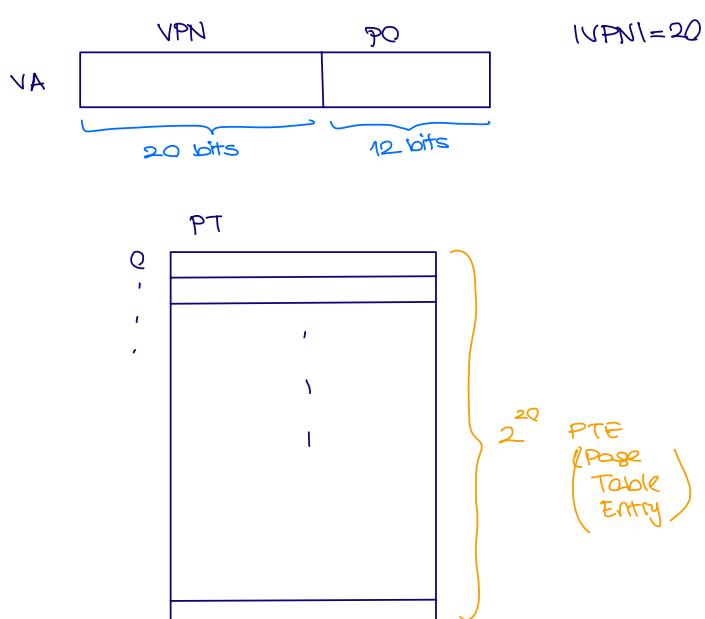
$$|VMI| = 2^{32} = 4 \text{ GB}$$

$$|VA| = 32 \text{ bits}$$

$$|Page| = 4 \text{ KB} = 2^{12} \text{ bytes}$$

$$|PM| = 1 \text{ GB}$$

$$|PA| = 30 \text{ bits}$$



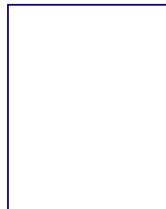
x86-64

$$|VA| = 48 \text{ bits}$$

$$|Page| = 4 \text{ KB} = 2^{12}$$



PT

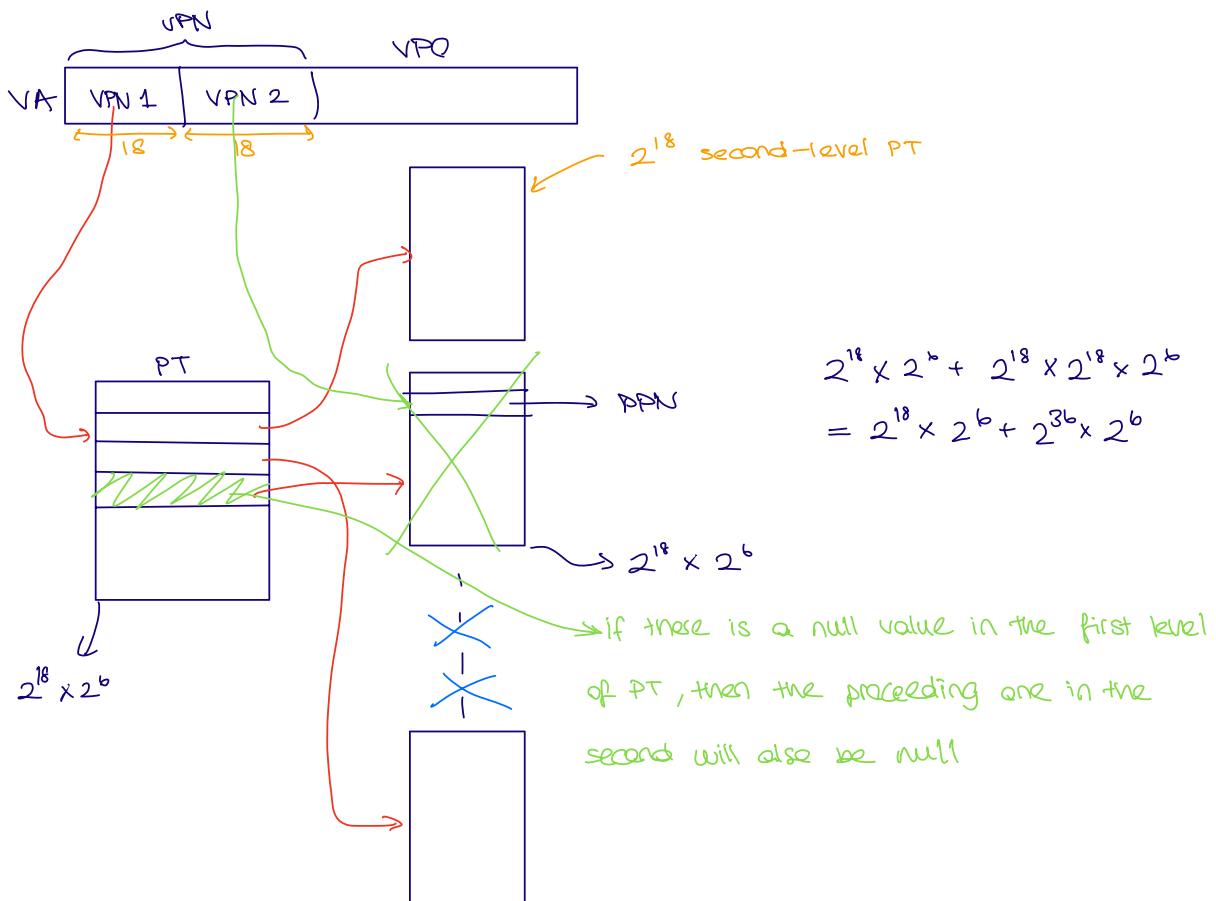


$$\begin{aligned} PT &= 2^{36} \times |PTE| \\ &= 2^{39} \text{ bytes} \end{aligned}$$

infeasible

$$|PM| = 1 \text{ TB} = 2^{40} \Rightarrow |PA| = 40 \text{ bits}$$

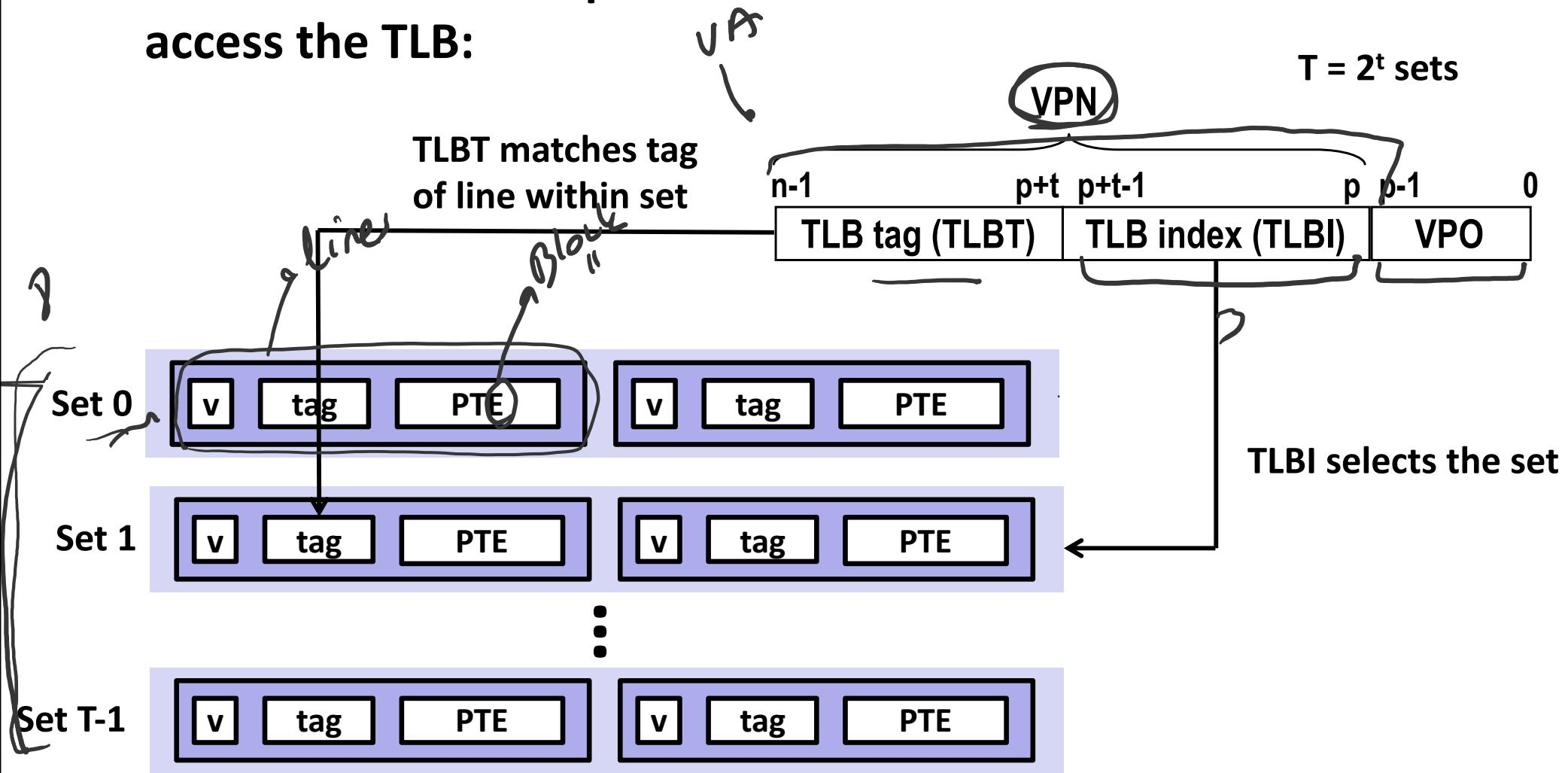
$$|PTE|$$



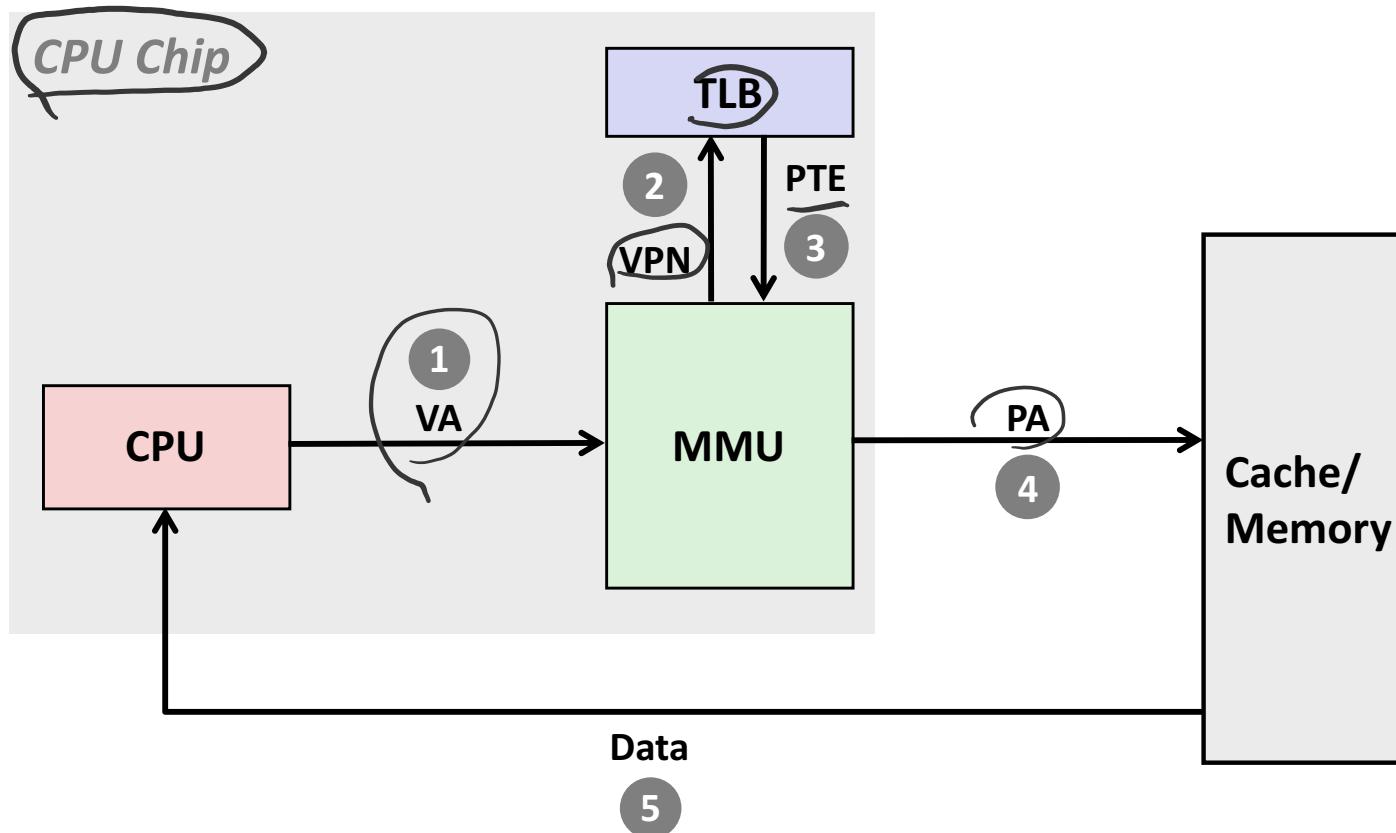
Accessing the TLB

*Cache - Block → Entry
by kis*

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

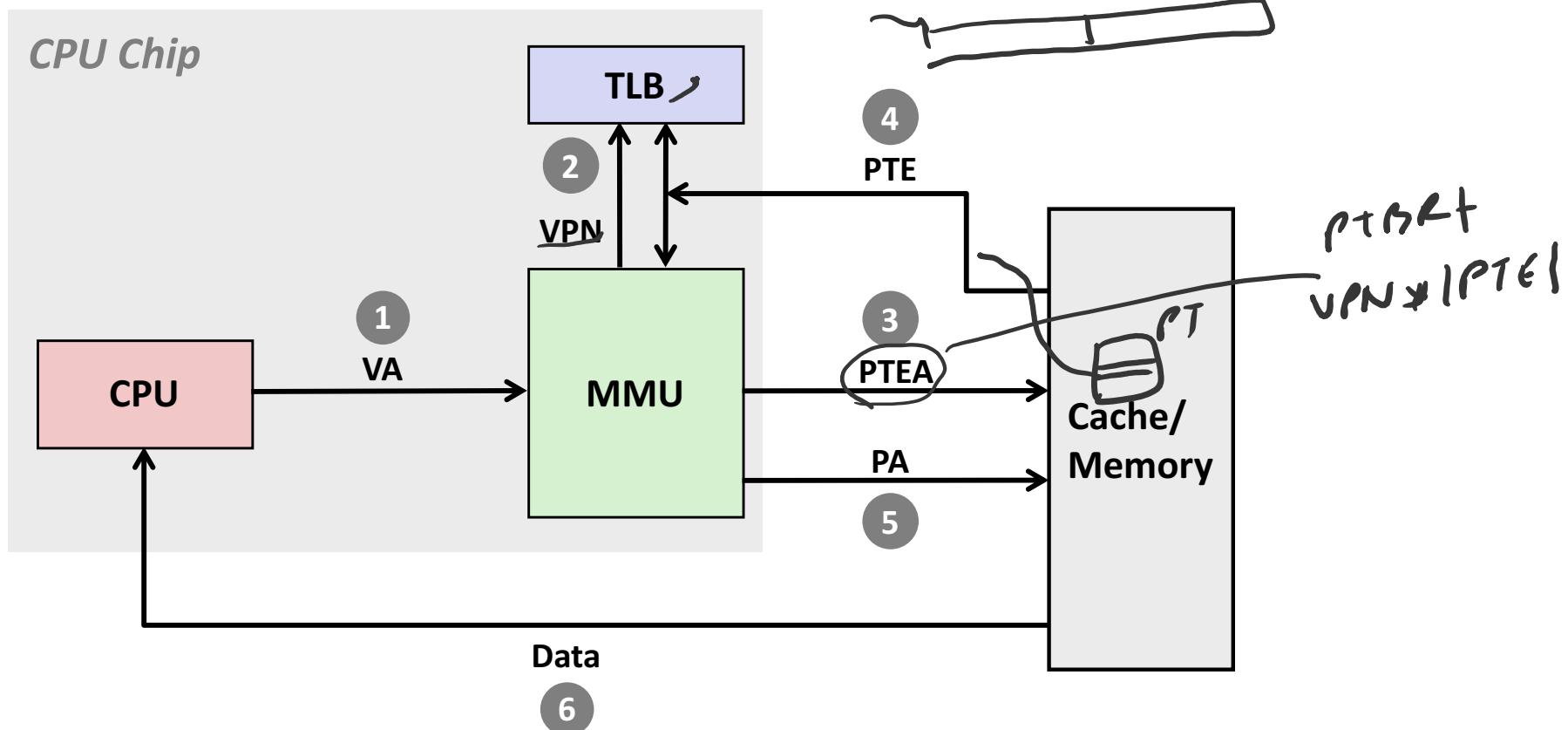


TLB Hit + Page Hit



A TLB hit eliminates a memory access

TLB Miss + Page Hit



A TLB miss incurs an additional 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

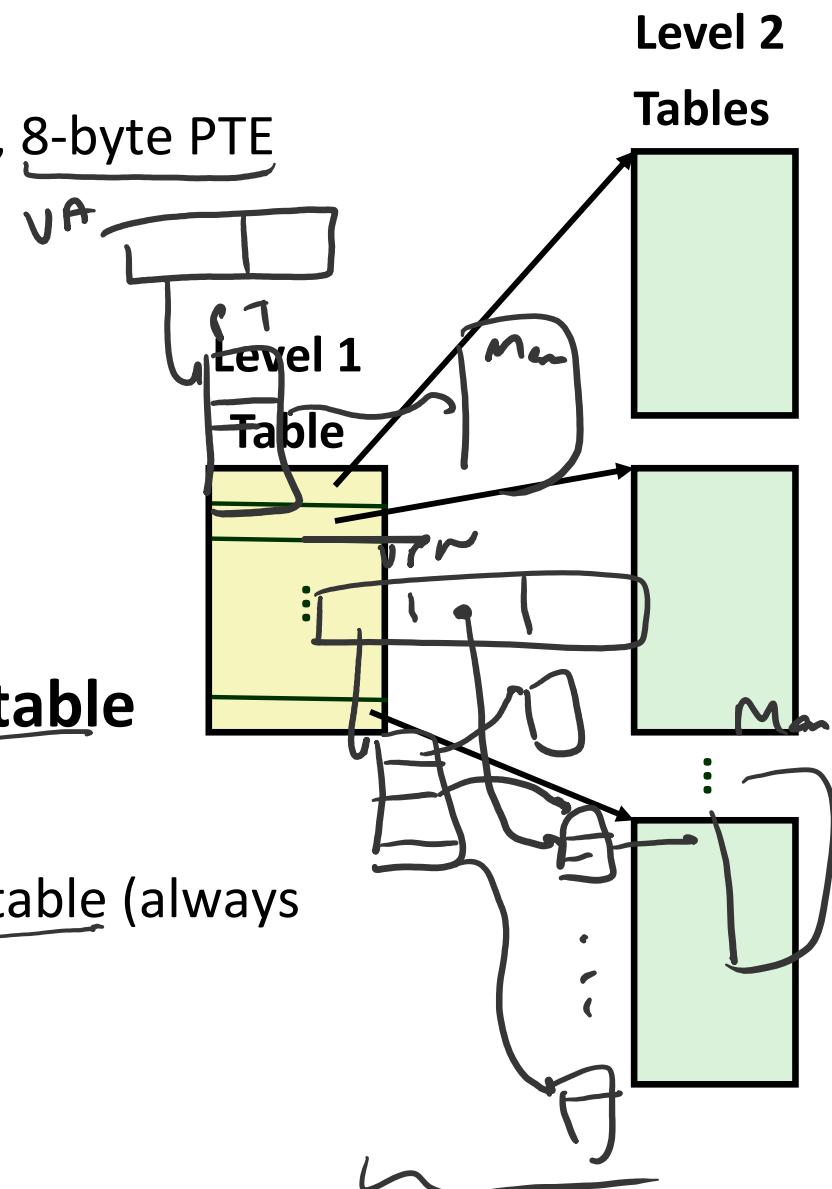
■ Problem:

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

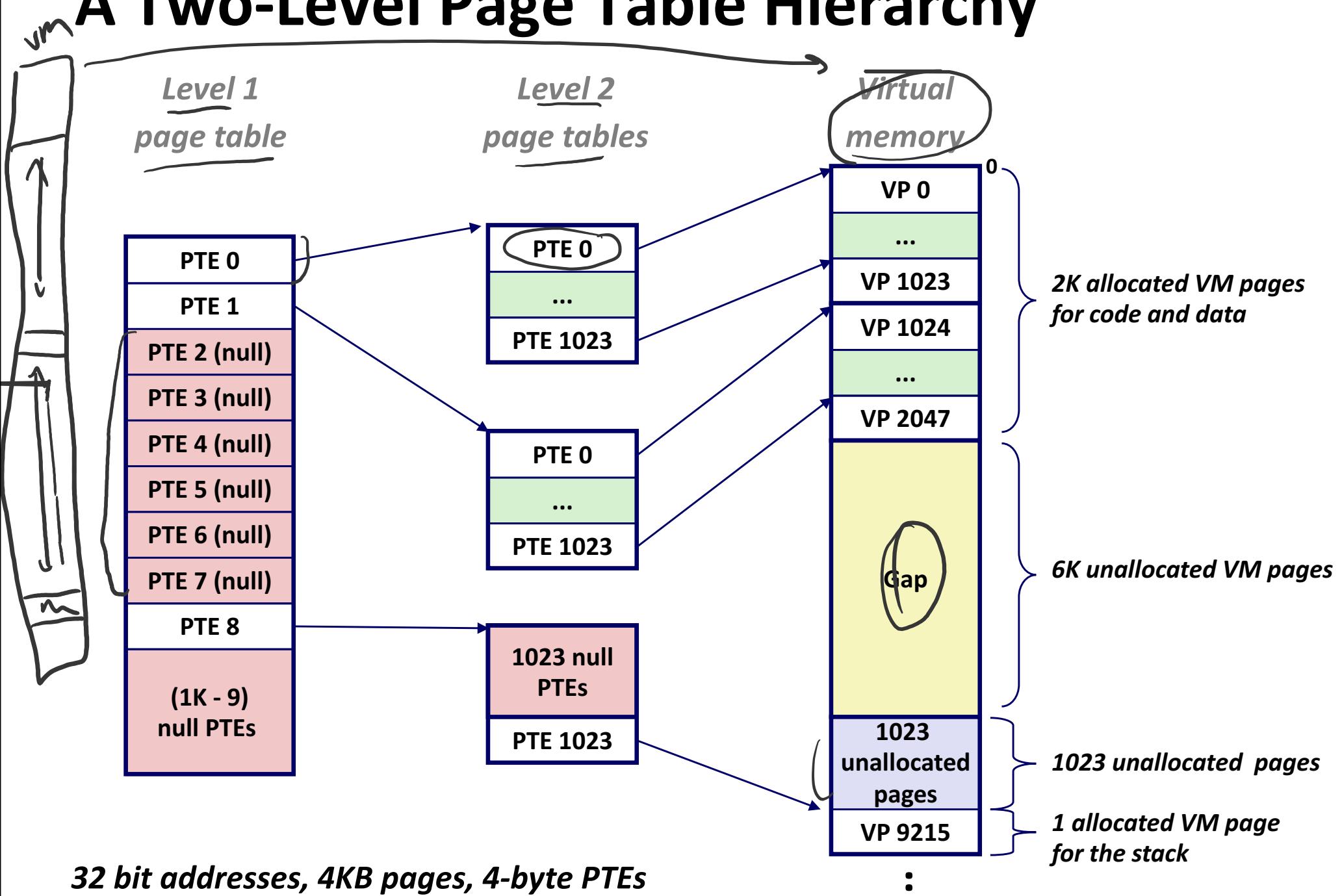
■ Common solution: Multi-level page table

■ Example: 2-level page table

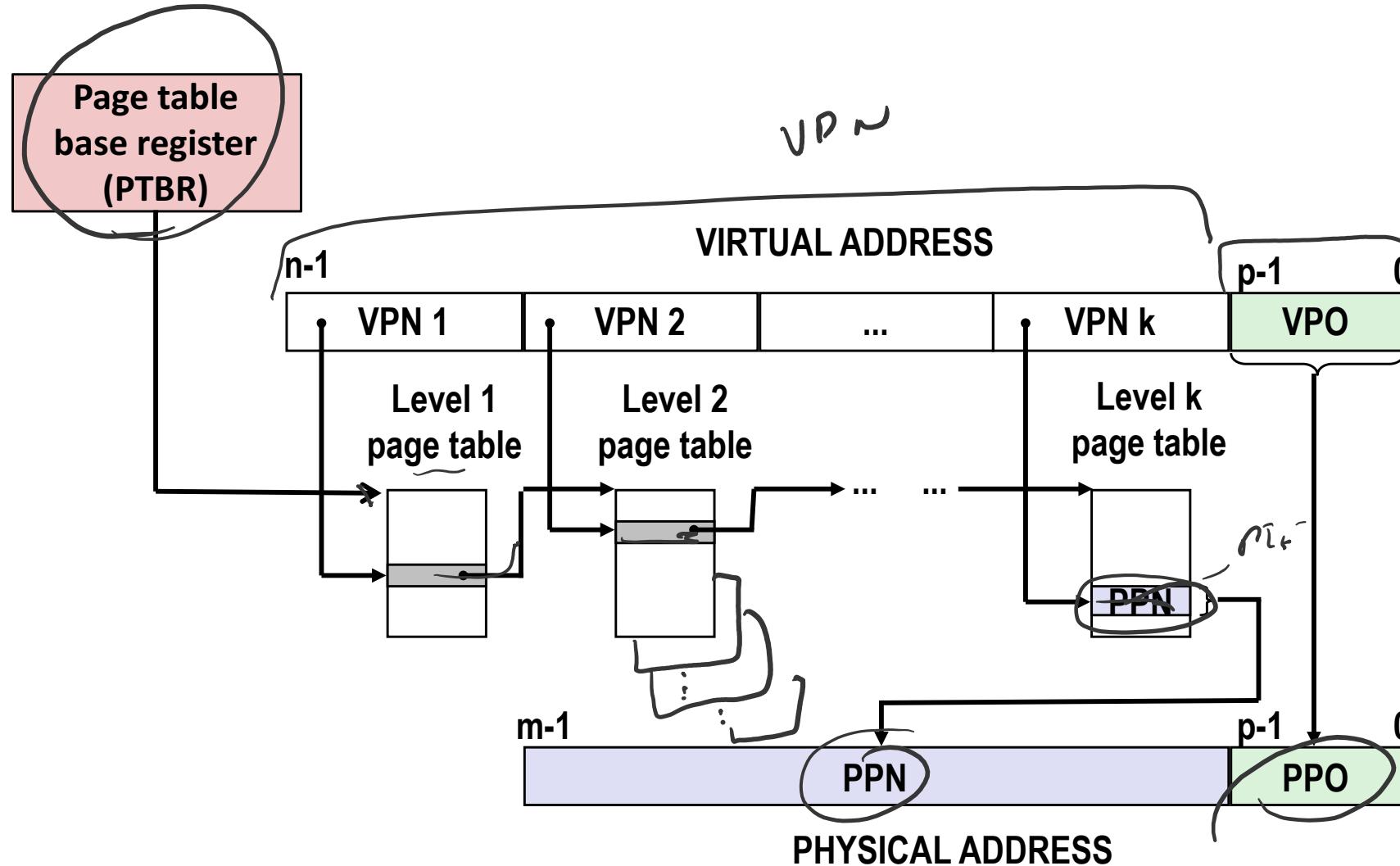
- Level 1 table: each PTE points to a page table (always memory resident)
- Level 2 table: each PTE points to a page (paged in and out like any other data)



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

Virtual Memory: Systems

CENG331 - Computer Organization

Instructors:

Murat Manguoglu (Section 1)

Erol Sahin (Section 2)

Adapted from slides of the textbook: <http://csapp.cs.cmu.edu/>

Today

- **Simple memory system example**
- Case study: Core i7/Linux memory system
- Memory mapping

Review of Symbols

■ Basic Parameters

- $N = 2^n$: Number of addresses in virtual address space
- $M = 2^m$: Number of addresses in physical address space
- $\bar{P} = 2^p$: Page size (bytes)

■ Components of the virtual address (VA)

- **TLBI**: TLB index
- **TLBT**: TLB tag
- **VPO**: Virtual page offset
- **VPN**: Virtual page number

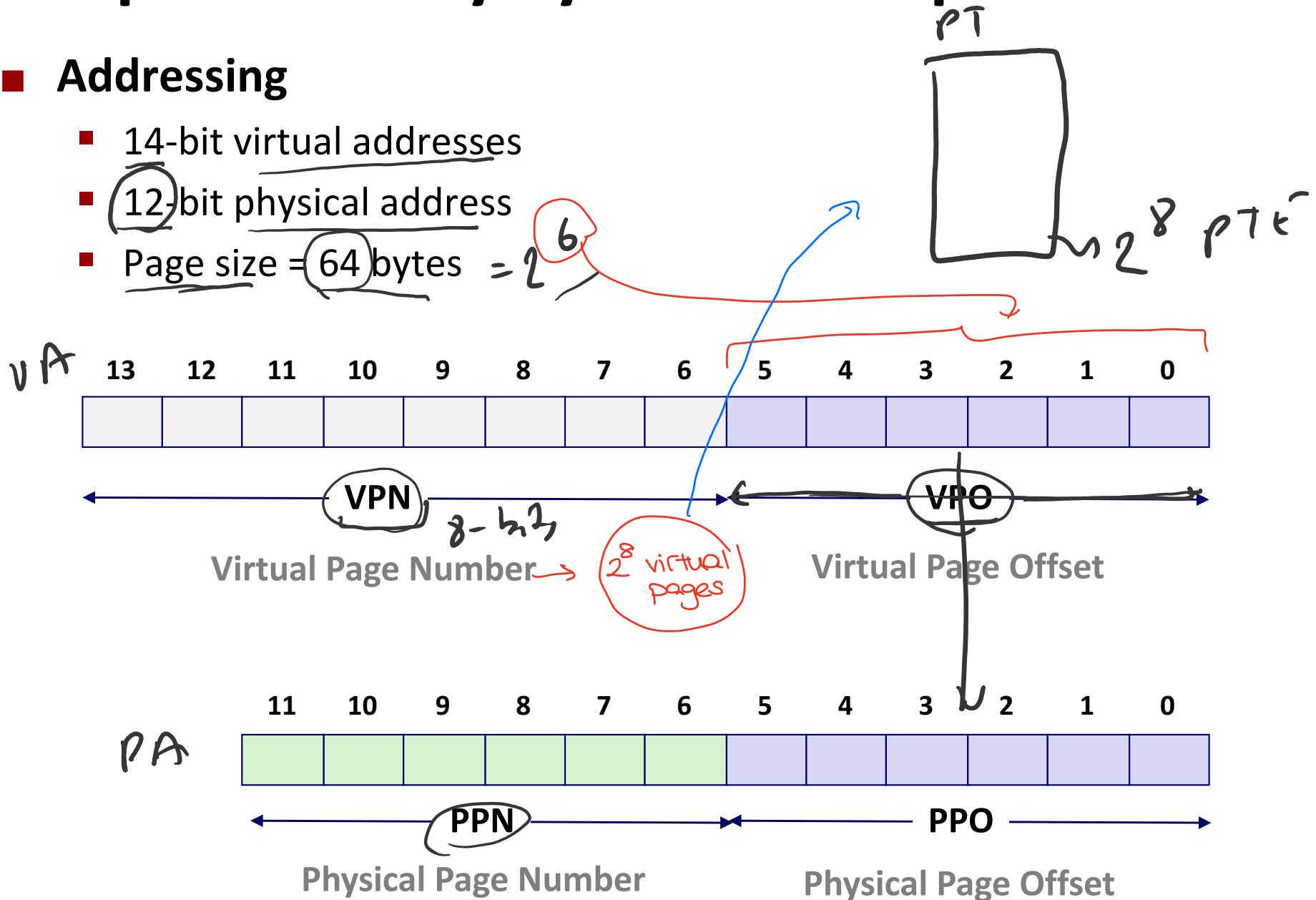
■ Components of the physical address (PA)

- **PPO**: Physical page offset (same as VPO)
- **PPN**: Physical page number
- **CO**: Byte offset within cache line
- **CI**: Cache index
- **CT**: Cache tag

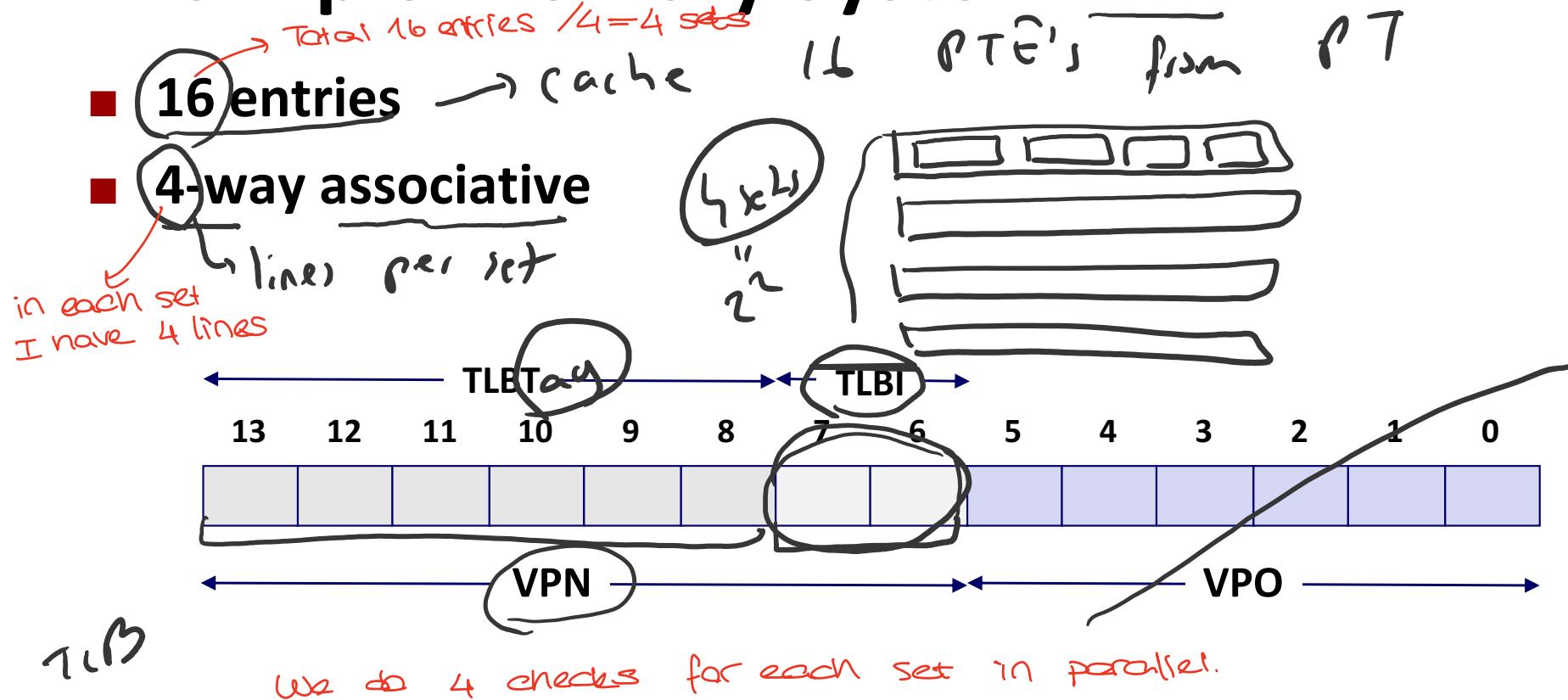
Simple Memory System Example

■ Addressing

- 14-bit virtual addresses
- 12-bit physical address
- Page size = 64 bytes = 2^6



1. Simple Memory System TLB

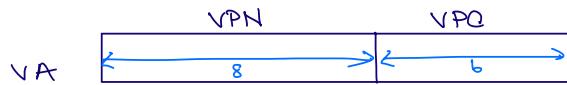


Set	Tag	PPN	Valid									
0	03	2D	1	09	0D	1	00	-	0	07	02	1
1	03	2D	1	02	-	0	04	-	0	0A	-	0
2	02	-	0	08	-	0	06	-	0	03	-	0
3	07	-	0	03	0D	1	0A	34	1	02	-	0

$$|VA| = 14 \text{ bits} \Rightarrow |VM| = 2^{14} = 16 \text{ kB}$$

$$|PA| = 12 \text{ bits} \Rightarrow |PM| = 2^{12} = 4 \text{ kB}$$

$$|Page| = 64 \text{ bytes} \Rightarrow |PO| = 6 \text{ bits} \Rightarrow 64 = 2^6$$



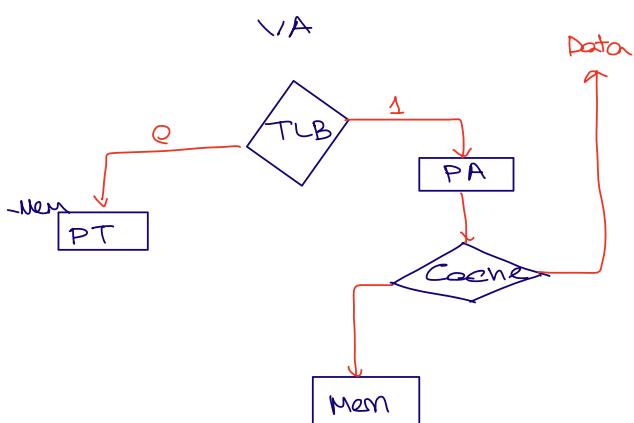
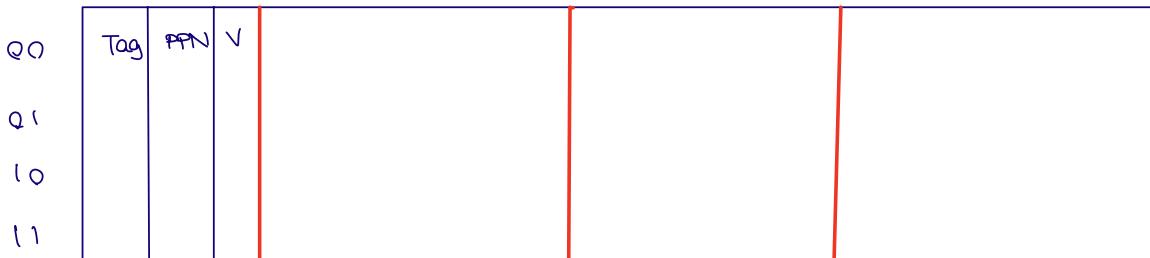
$$|VPN| = 8 \text{ bits} \Rightarrow \# \text{ of virtual Page} = 2^8 = 256 = \# \text{ of PTE in PT}$$



$$|PFN| = 6 \text{ bits} \Rightarrow 2^6 = 64 \text{ Physical Page}$$

TLB : 16 entries , 4-way associative
 (PTE)

$$16 = 4 \times S \Rightarrow S=4$$



2. Simple Memory System Page Table

Only show first 16 entries (out of 256)

This is in Physical memory

VPN	PPN	Valid
00	28	1
01	-	0
02	33	1
03	02	1
04	-	0
05	16	1
06	-	0
07	-	0

VPN	PPN	Valid
08	13	1
09	17	1
0A	09	1
0B	-	0
0C	-	0
0D	2D	1
0E	11	1
0F	0D	1

→ Valid = 0 means the entry

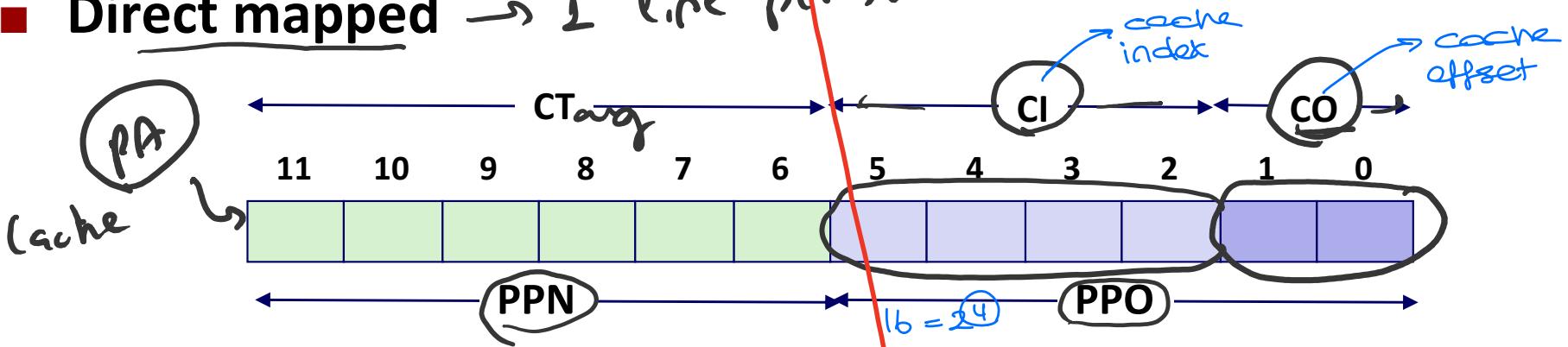
has not yet
moved on to the
memory



Page Fault

3. Simple Memory System Cache

- 16 lines, 4-byte block size
- Physically addressed \Rightarrow 12 bits for Physical Address (16 sets)
- Direct mapped \rightarrow 1 line per set

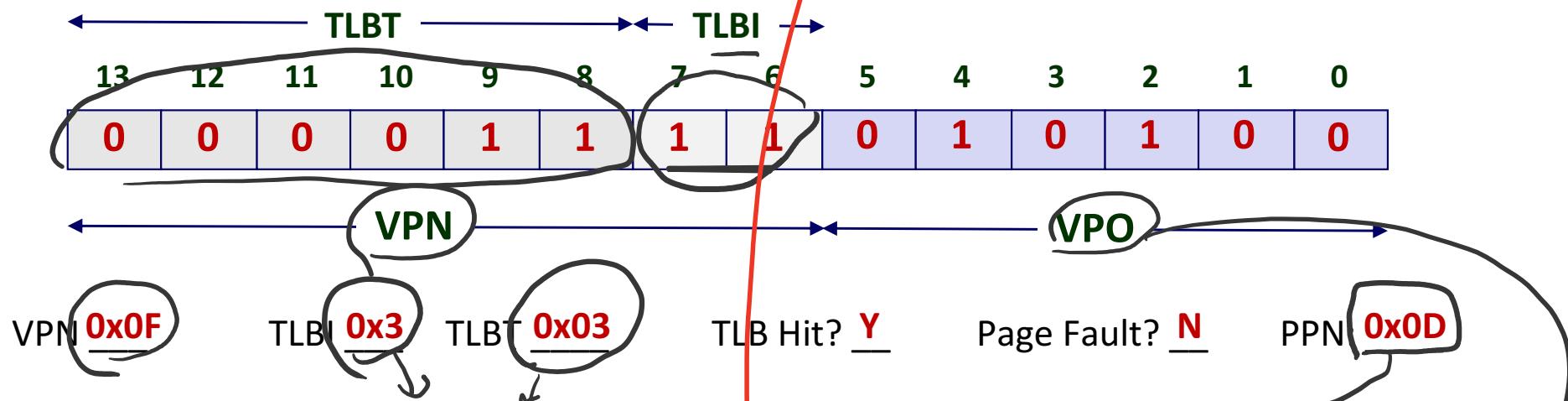


<i>Idx</i>	<i>Tag</i>	<i>Valid</i>	<i>B0</i>	<i>B1</i>	<i>B2</i>	<i>B3</i>
0	19	1	99	11	23	11
1	15	0	-	-	-	-
2	1B	1	00	02	04	08
3	36	0	-	-	-	-
4	32	1	43	6D	8F	09
5	0D	1	36	72	F0	1D
6	31	0	-	-	-	-
7	16	1	11	C2	DF	03

<i>Idx</i>	<i>Tag</i>	<i>Valid</i>	<i>B0</i>	<i>B1</i>	<i>B2</i>	<i>B3</i>
8	24	1	3A	00	51	89
9	2D	0	-	-	-	-
A	2D	1	93	15	DA	3B
B	0B	0	-	-	-	-
C	12	0	-	-	-	-
D	16	1	04	96	34	15
E	13	1	83	77	1B	D3
F	14	0	-	-	-	-

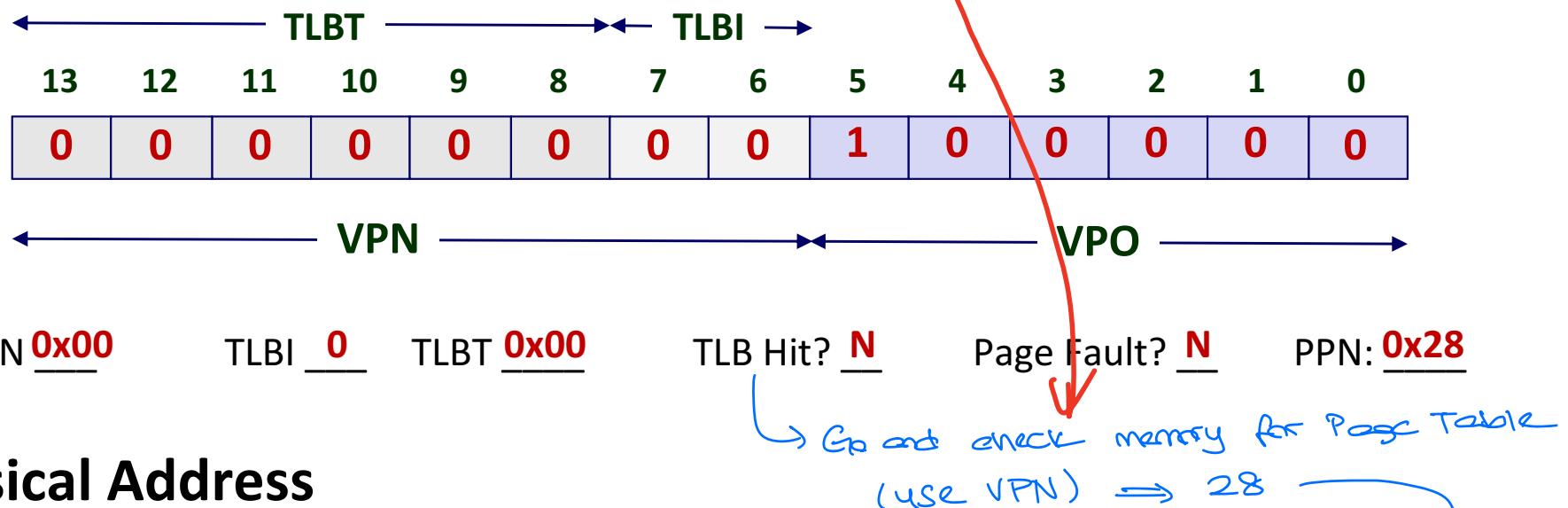
Address Translation Example #1

Virtual Address: 0x03D4

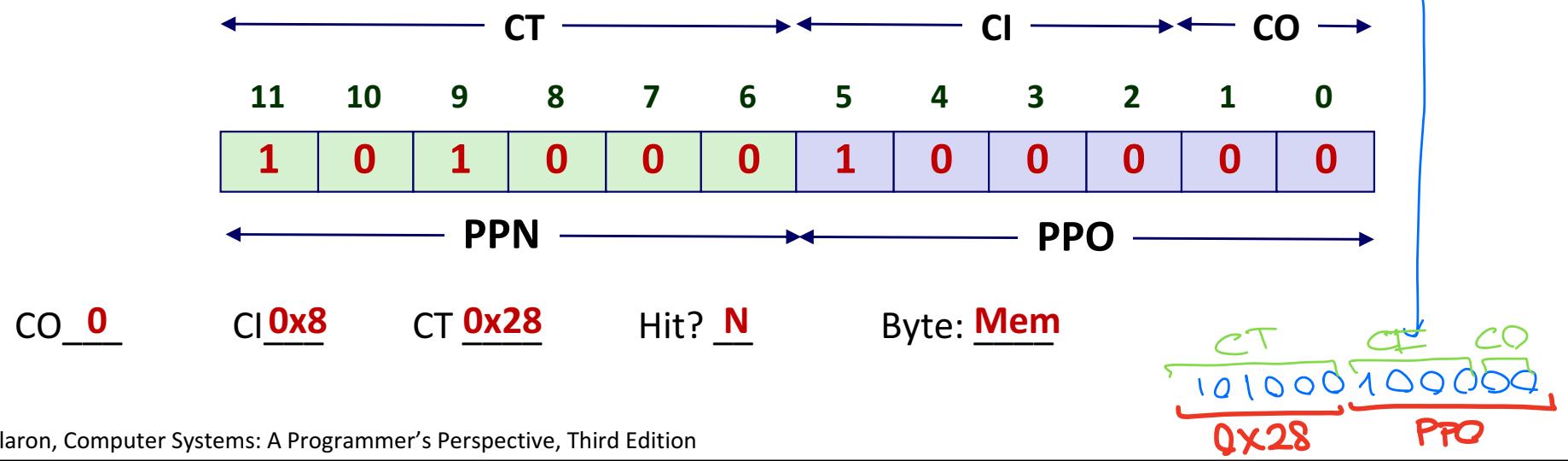


Address Translation Example #2

Virtual Address: 0x0020

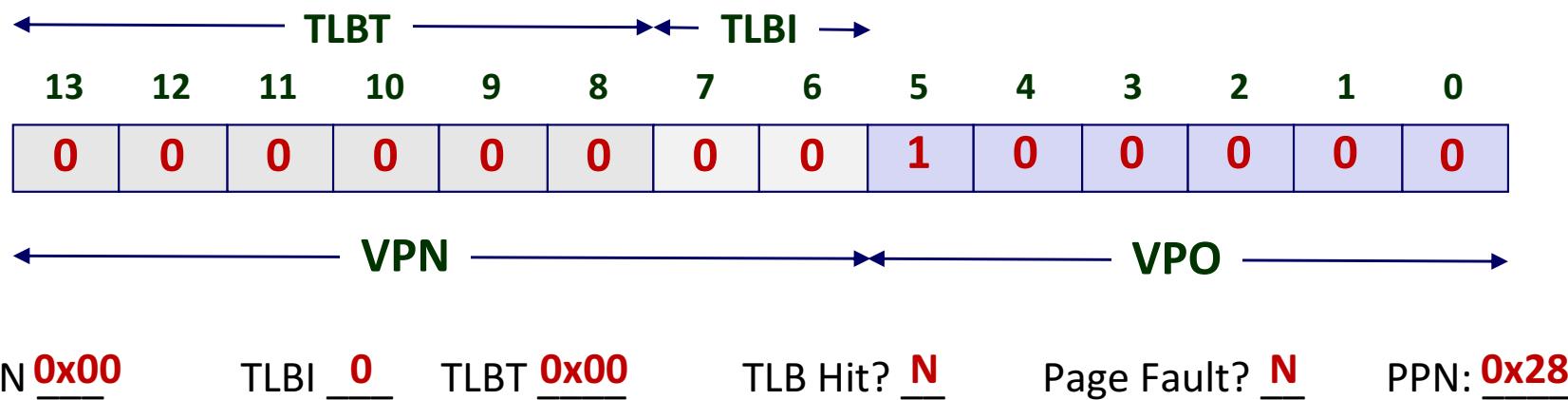


Physical Address

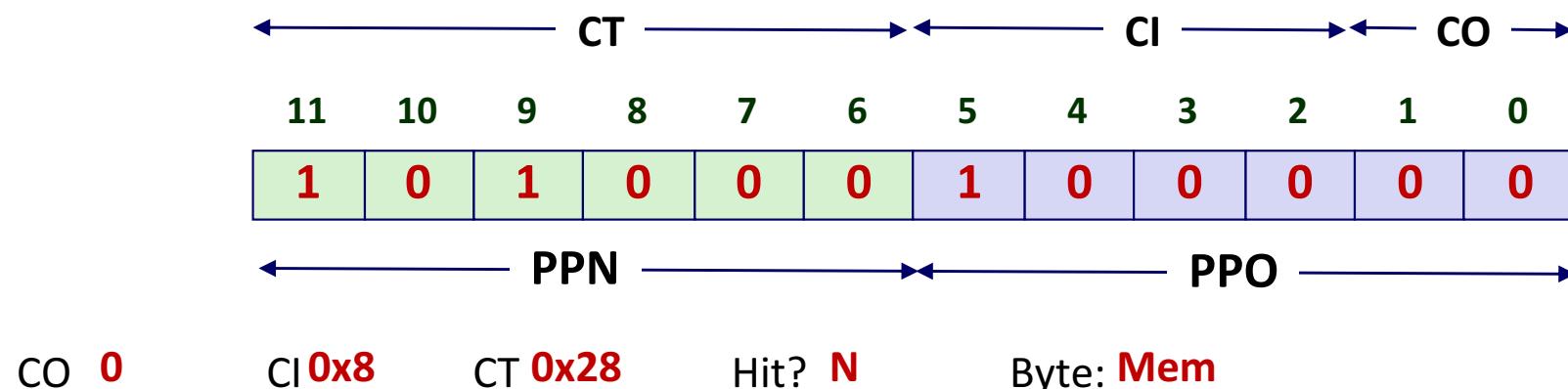


Address Translation Example #3

Virtual Address: 0x0020

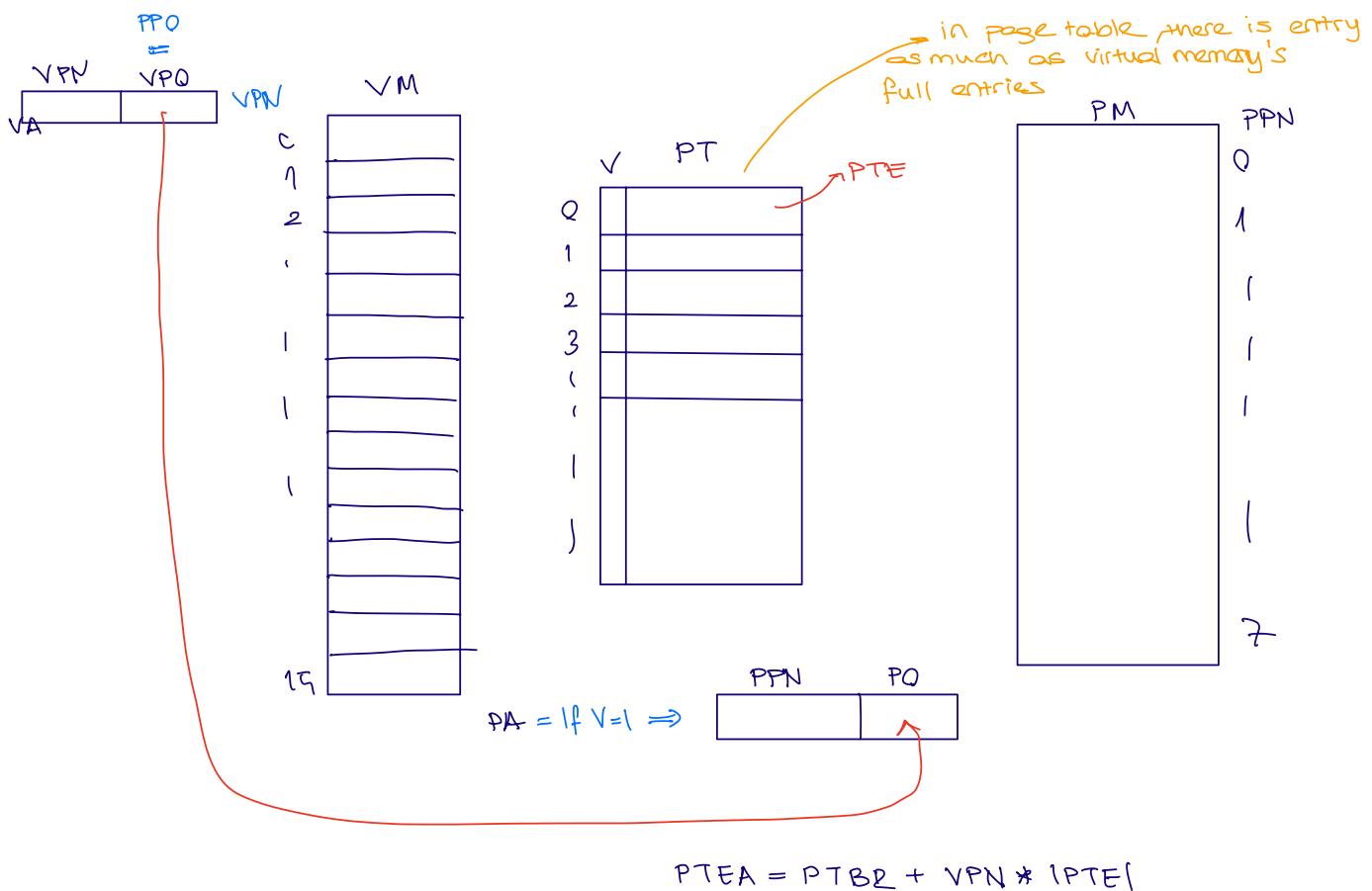


Physical Address



Today

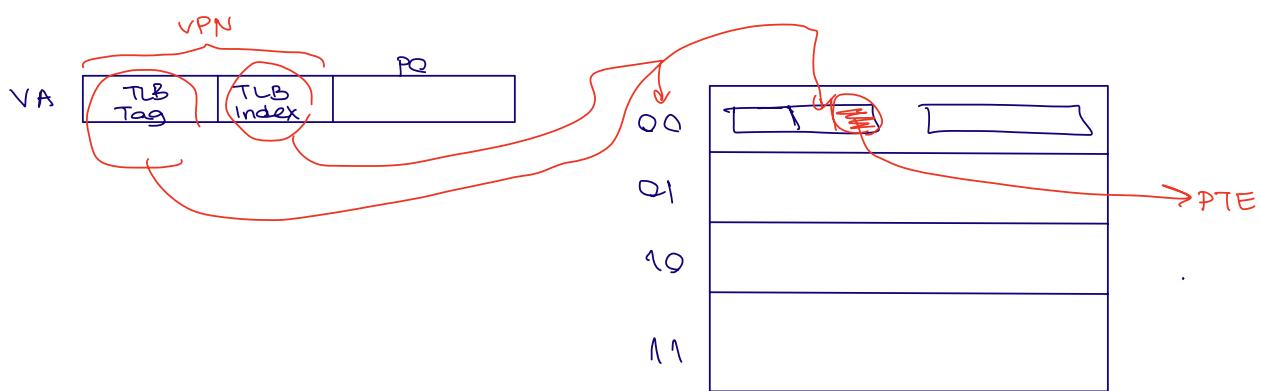
- Simple memory system example
- Case study: Core i7/Linux memory system
- Memory mapping



For each process, there is a distinct page Table (PT) and virtual memory (VM)

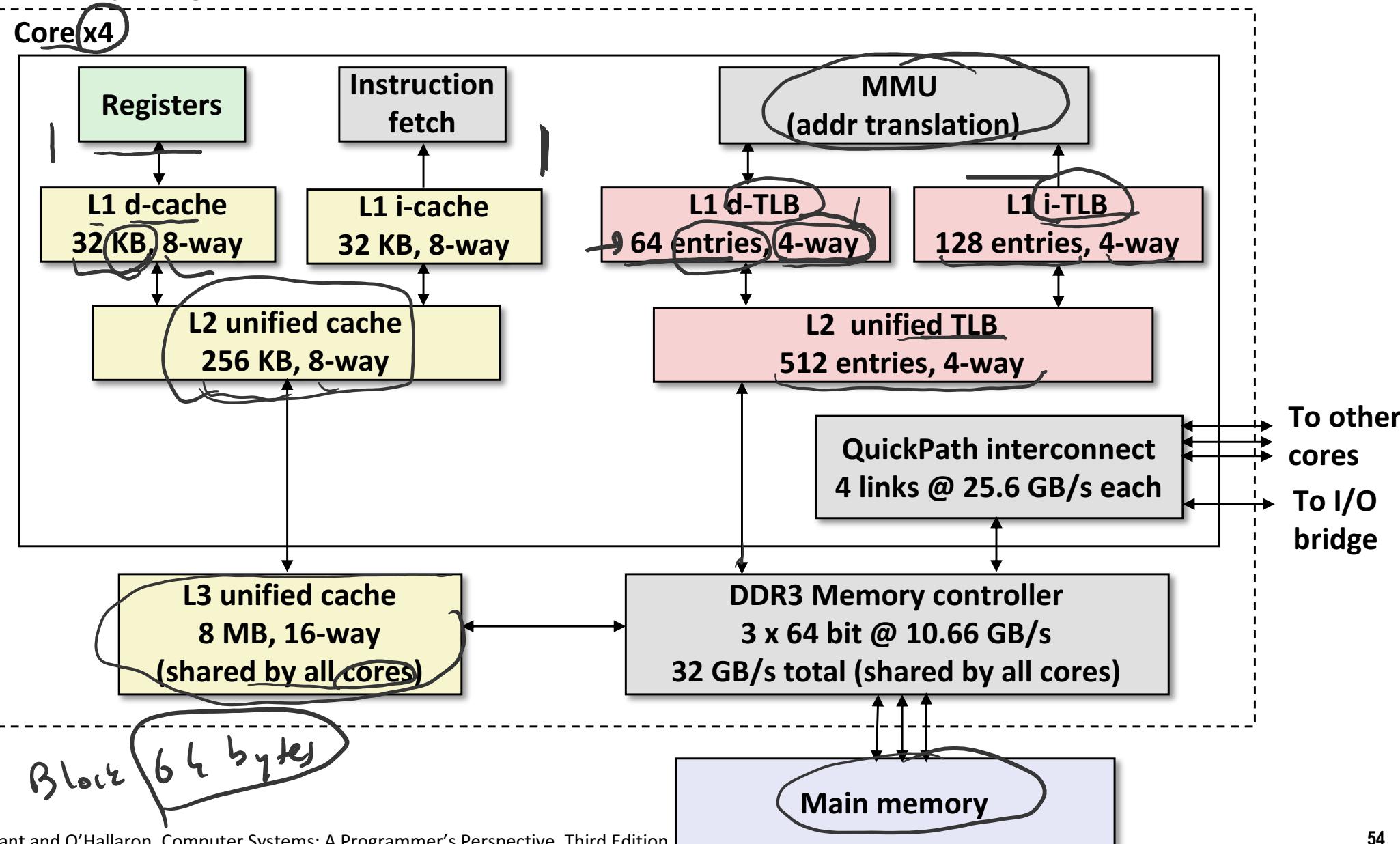
↳ PTBase changes accordingly.

4GB physical yarin also bile 16 GB'lik bir process'in dosyaları. PM
değilimiz səy sonucta virtual memory ianın bir cache (ilə 4 GB kənarı
isləm yaxşı sonra digər data üzərində...)



Intel Core i7 Memory System

Processor package



Review of 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)

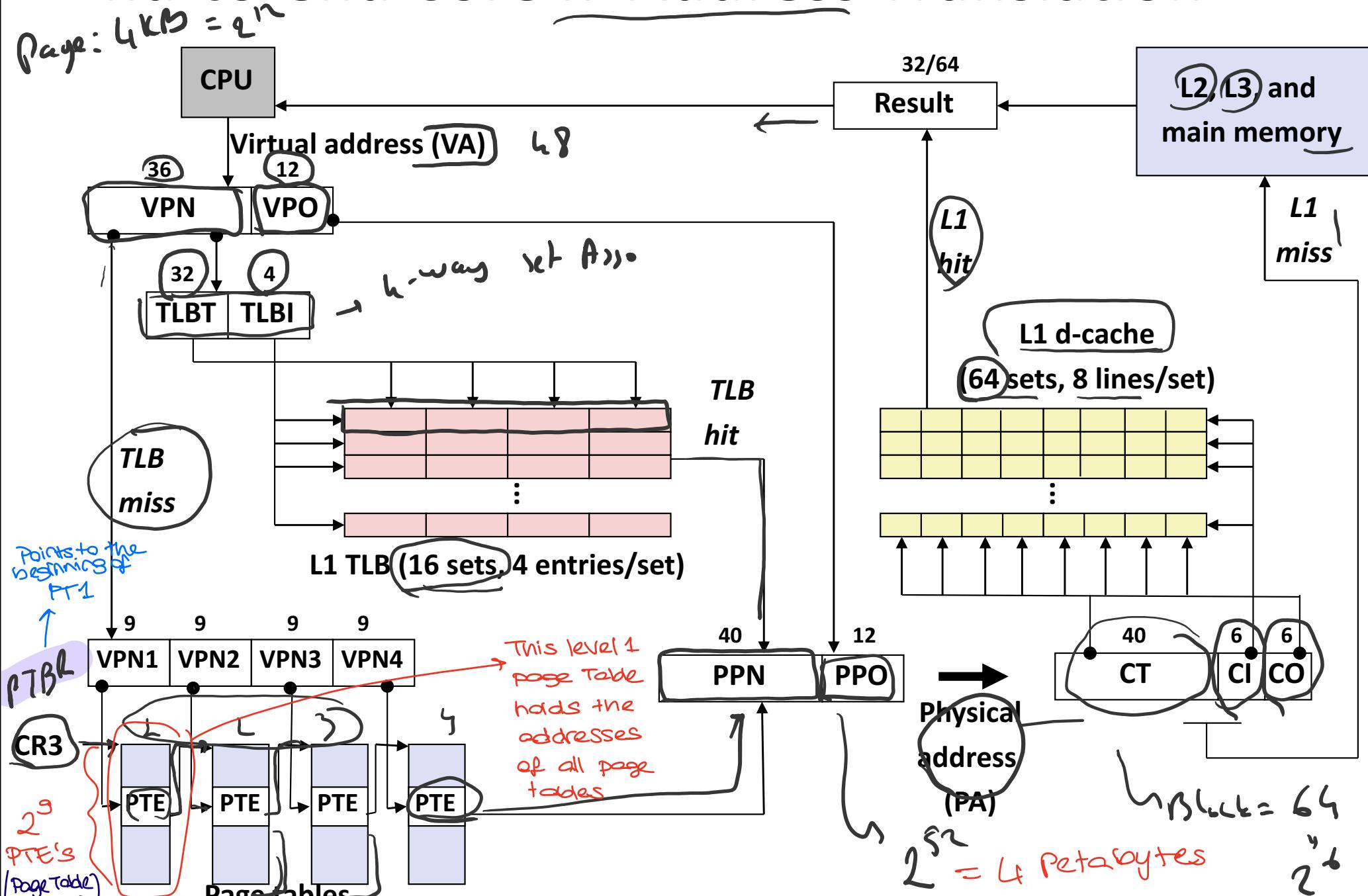
■ Components of the virtual address (VA)

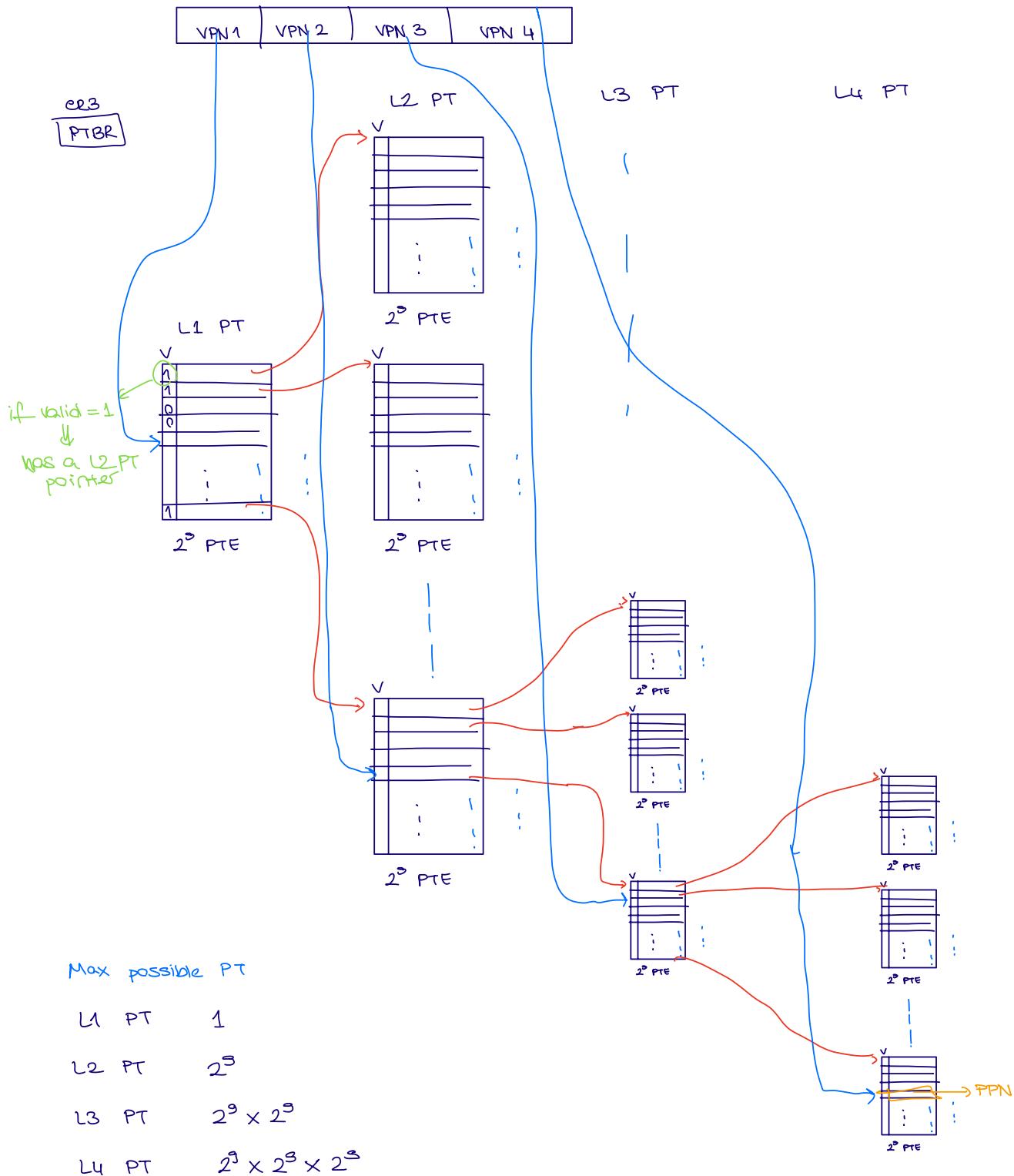
- **TLBI**: TLB index
- **TLBT**: TLB tag
- **VPO**: Virtual page offset
- **VPN**: Virtual page number

■ Components of the physical address (PA)

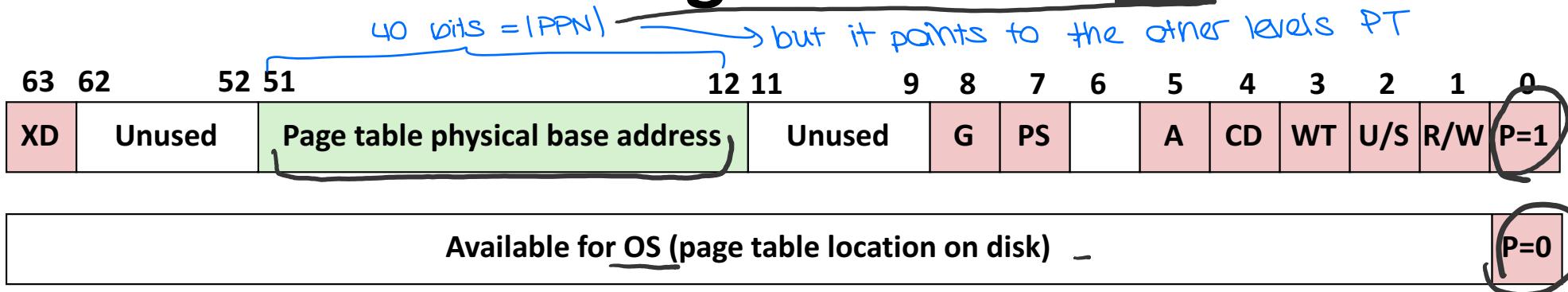
- **PPO**: Physical page offset (same as VPO)
- **PPN**: Physical page number
- **CO**: Byte offset within cache line
- **CI**: Cache index
- **CT**: Cache tag

End-to-end Core i7 Address Translation





Core i7 Level 1-3 Page Table Entries PTE



Each entry references a 4K child page table. Significant fields:

P: Child page table present in physical memory (1) or not (0). P=1

R/W: Read-only or read-write access access permission for all reachable pages.

U/S: user or supervisor (kernel) mode access permission for all reachable pages.

WT: Write-through or write-back cache policy for the child page table.

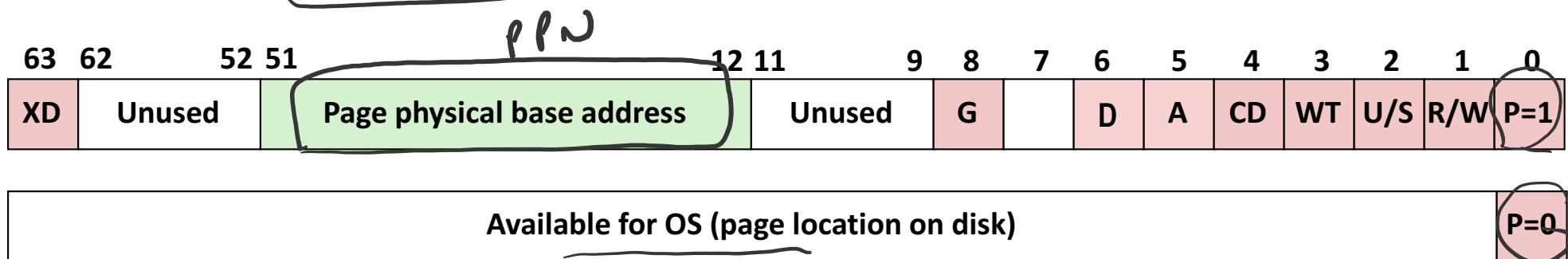
A: Reference bit (set by MMU on reads and writes, cleared by software).

PS: Page size either 4 KB or 4 MB (defined for Level 1 PTEs only).

Page table physical base address: 40 most significant bits of physical page table address (forces page tables to be 4KB aligned) \open\

XD: Disable or enable instruction fetches from all pages reachable from this PTE.

Core i7 Level 4 Page Table Entries



Each entry references a 4K child page. Significant fields:

P: Child page is present in memory (1) or not (0)

R/W: Read-only or read-write access permission for child page

U/S: User or supervisor mode access

WT: Write-through or write-back cache policy for this page

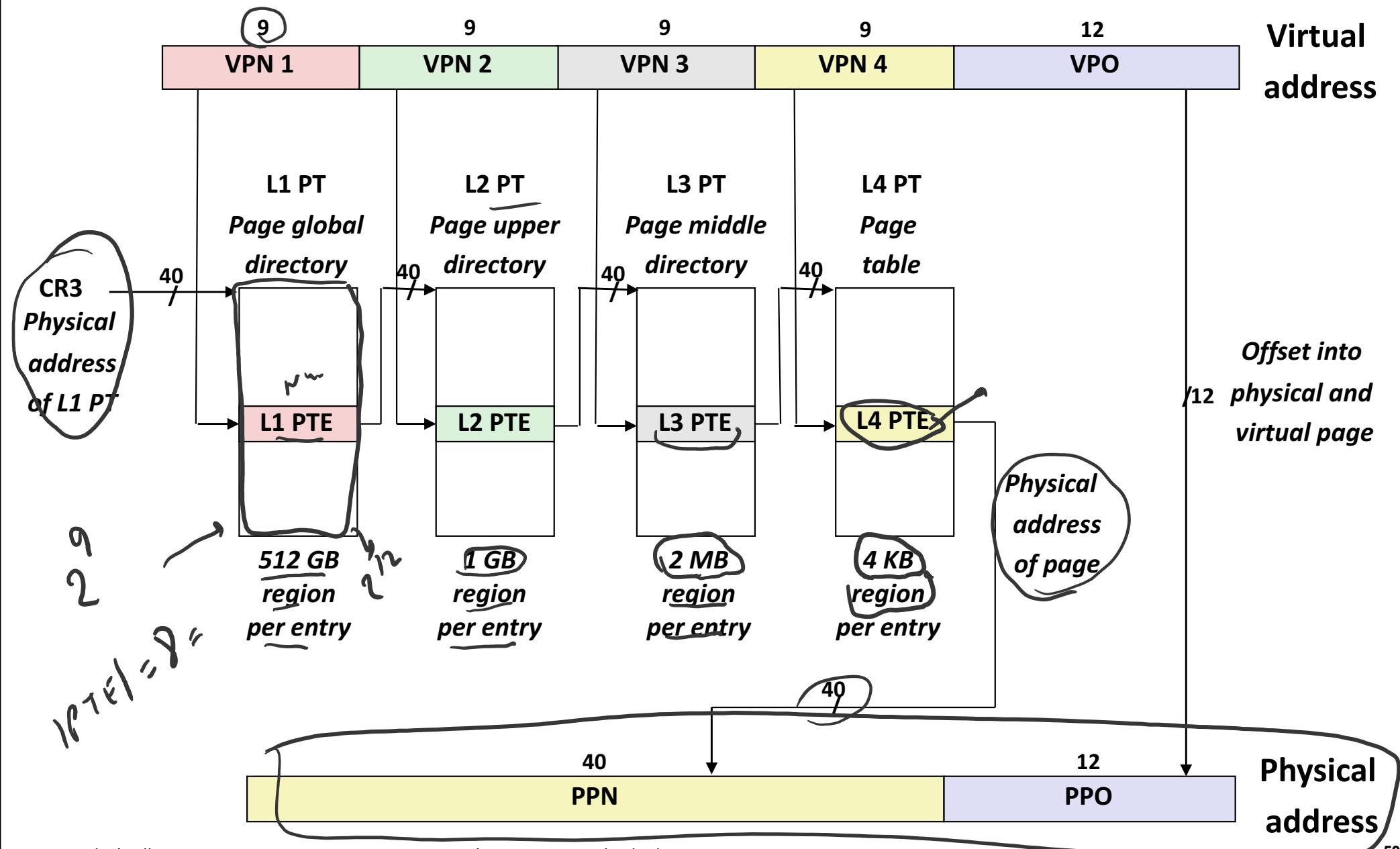
A: Reference bit (set by MMU on reads and writes, cleared by software)

D: Dirty bit (set by MMU on writes, cleared by software)

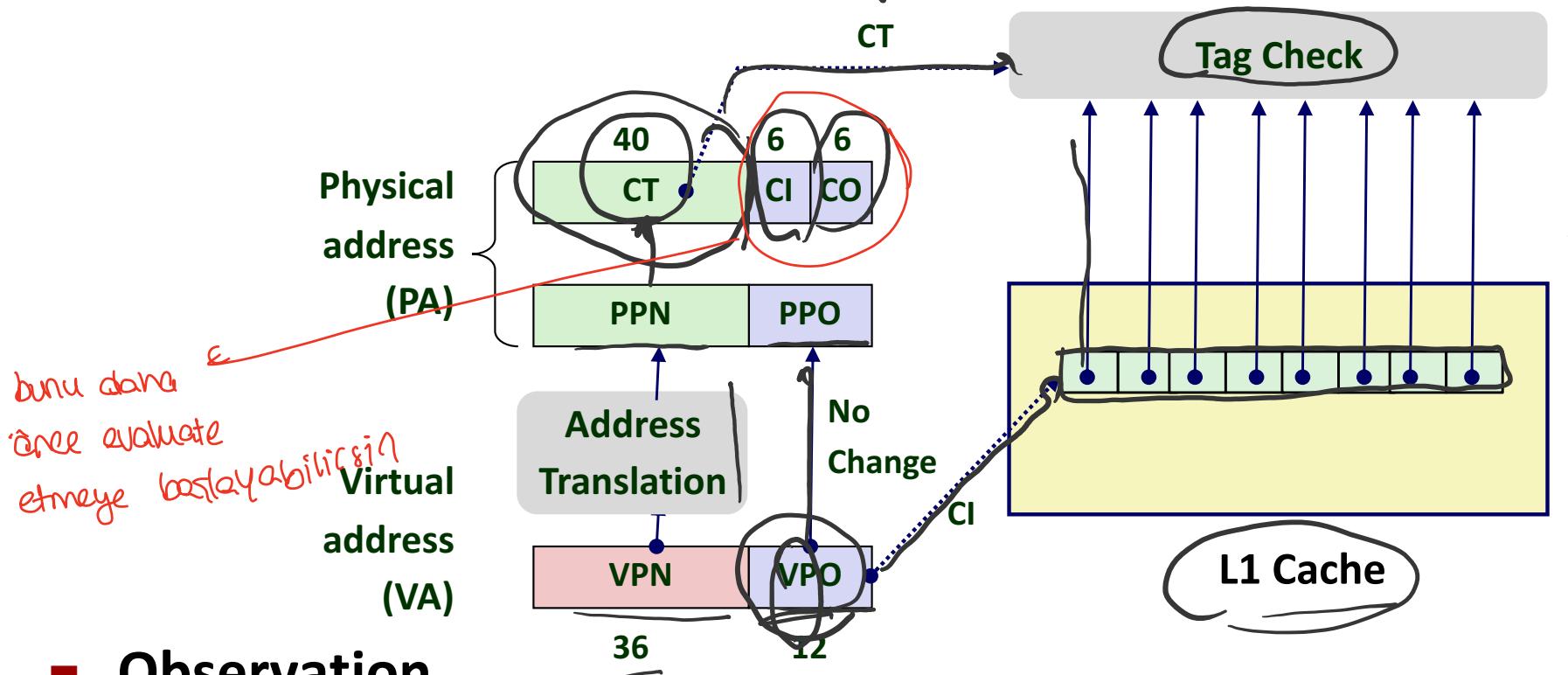
Page physical base address: 40 most significant bits of physical page address
(forces pages to be 4KB aligned)

XD: Disable or enable instruction fetches from this page.

Core i7 Page Table Translation



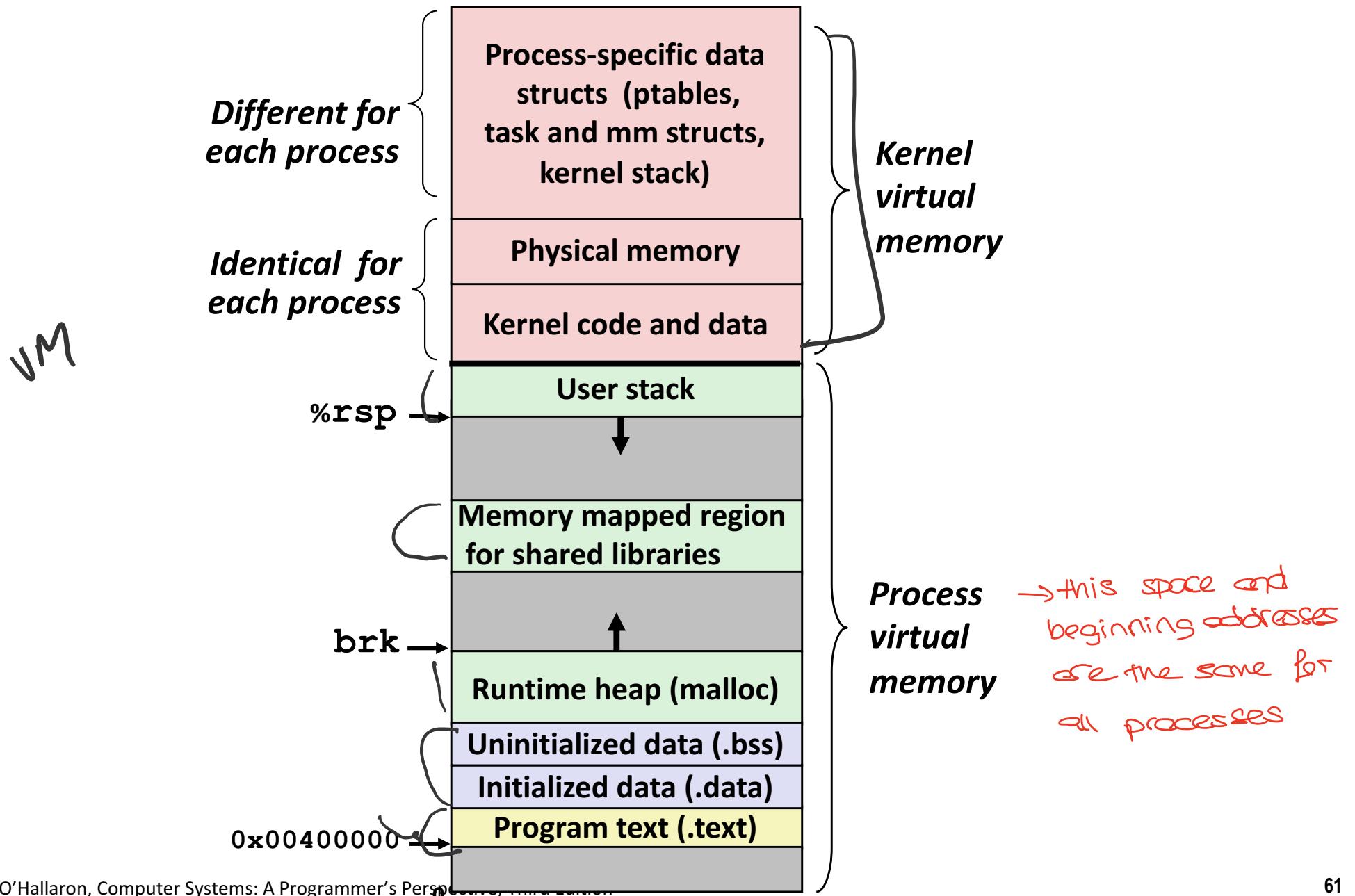
Cute Trick for Speeding Up L1 Access



Observation

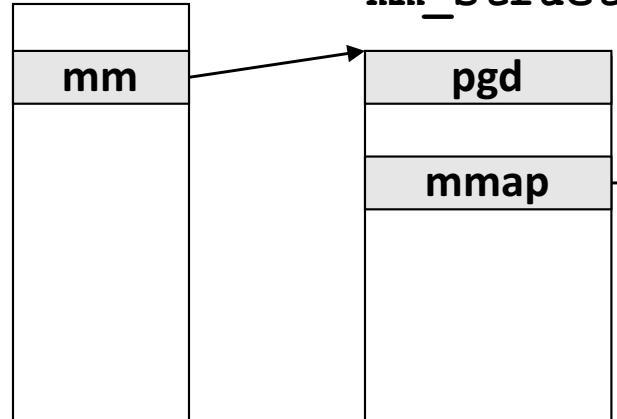
- Bits that determine CI identical in virtual and physical address
- Can index into cache while address translation taking place
- Generally we hit in TLB, so PPN bits (CT bits) available next
- "Virtually indexed, physically tagged" L1 cache
- Cache carefully sized to make this possible

Virtual Address Space of a Linux Process

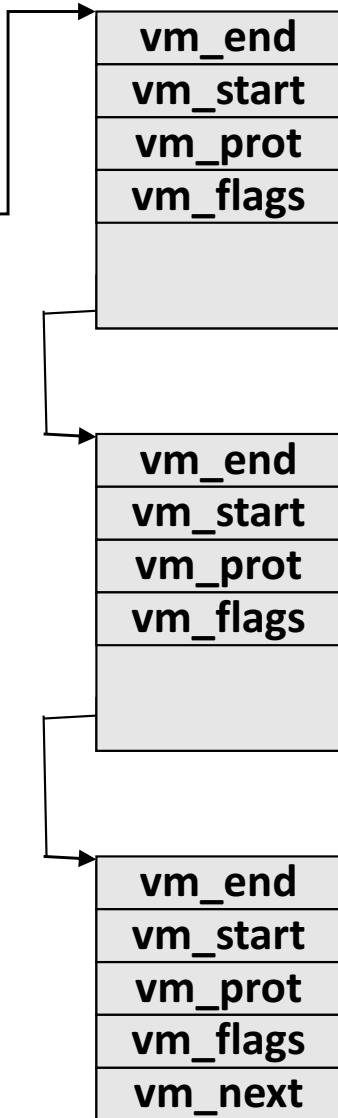


Linux Organizes VM as Collection of “Areas”

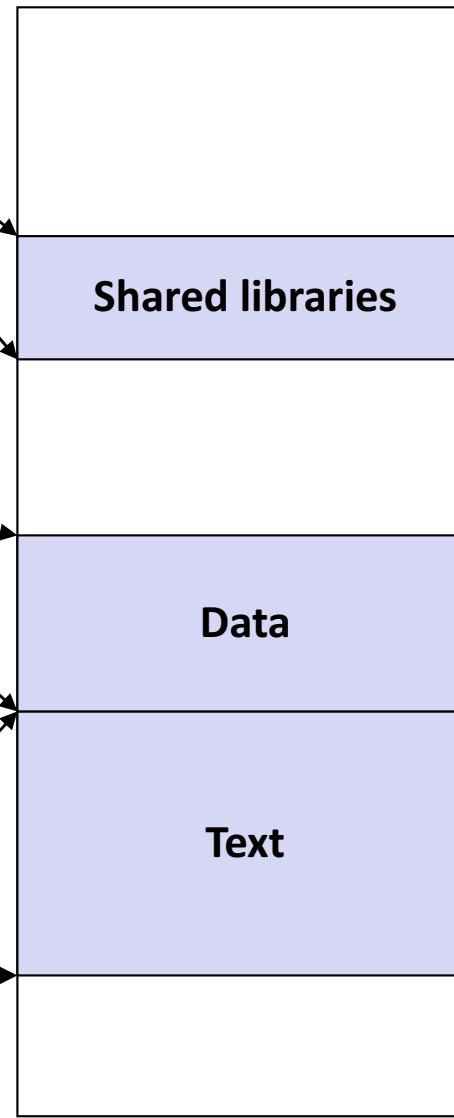
`task_struct`



`vm_area_struct`



Process virtual memory



■ `pgd`:

- Page global directory address
- Points to L1 page table

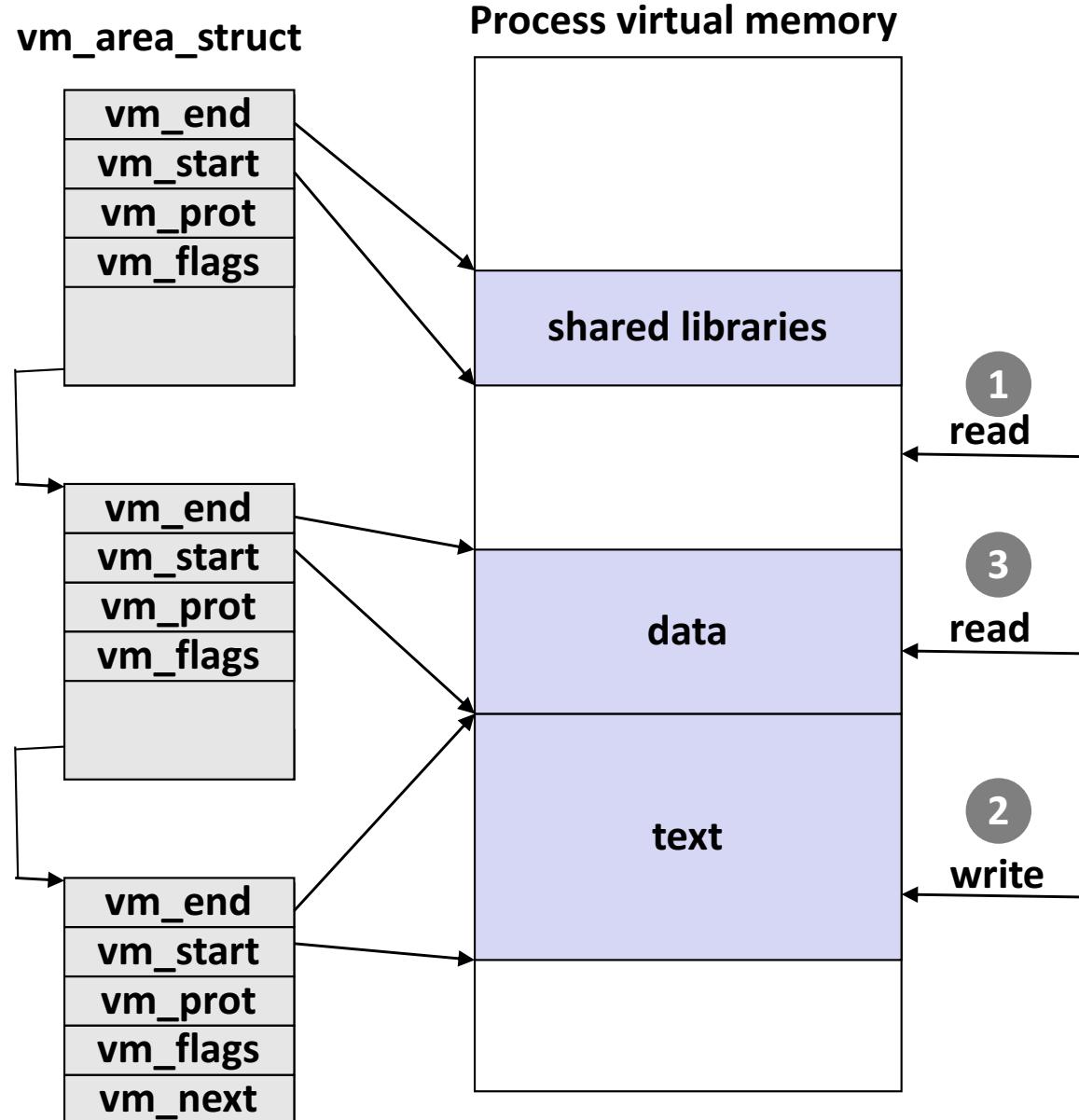
■ `vm_prot`:

- Read/write permissions for this area

■ `vm_flags`

- Pages **shared** with other processes or **private** to this process

Linux Page Fault Handling



Segmentation fault:
accessing a non-existing page

Normal page fault

Protection exception:
e.g., violating permission by
writing to a read-only page (Linux
reports as Segmentation fault)

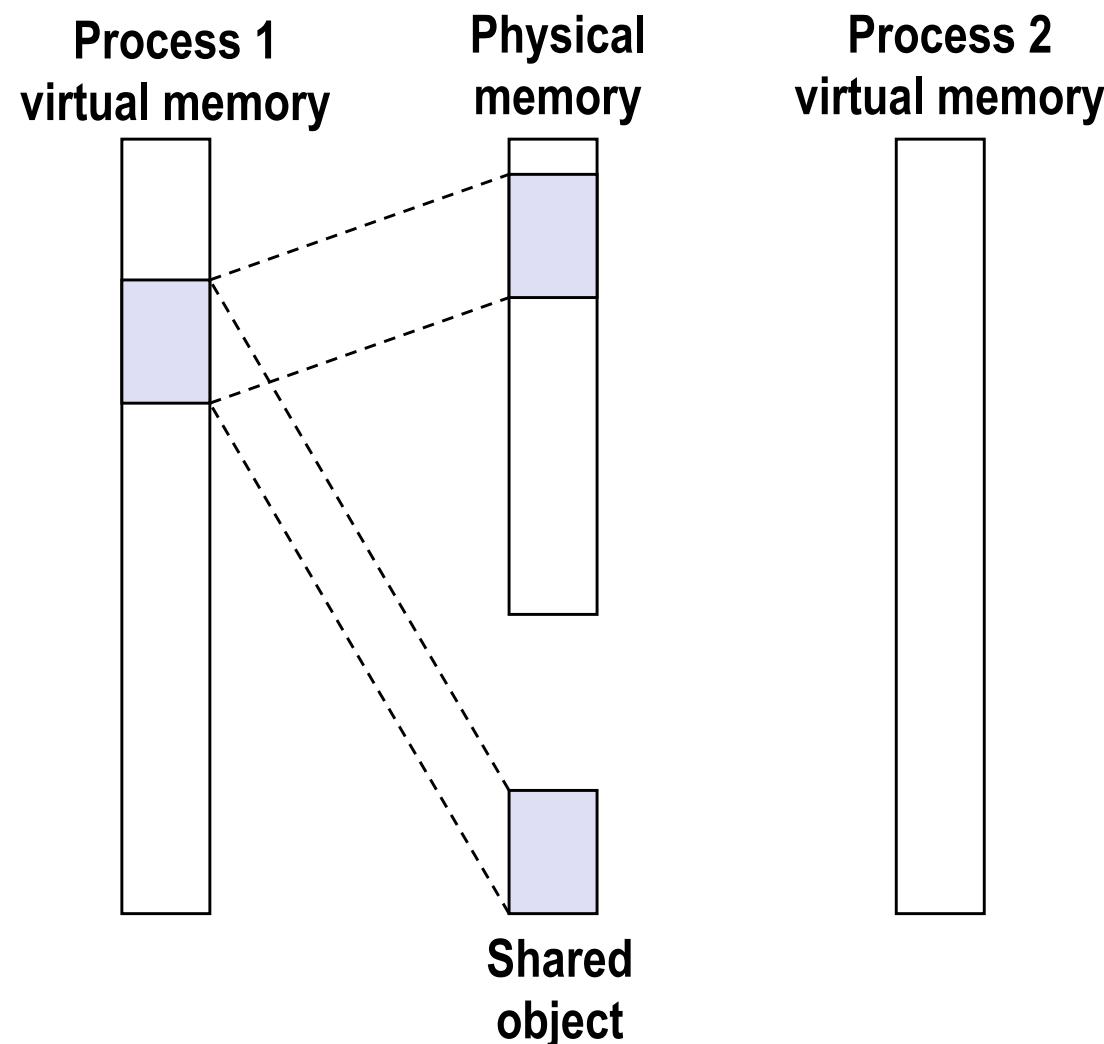
Today

- Simple memory system example
- Case study: Core i7/Linux memory system
- Memory mapping

Memory Mapping

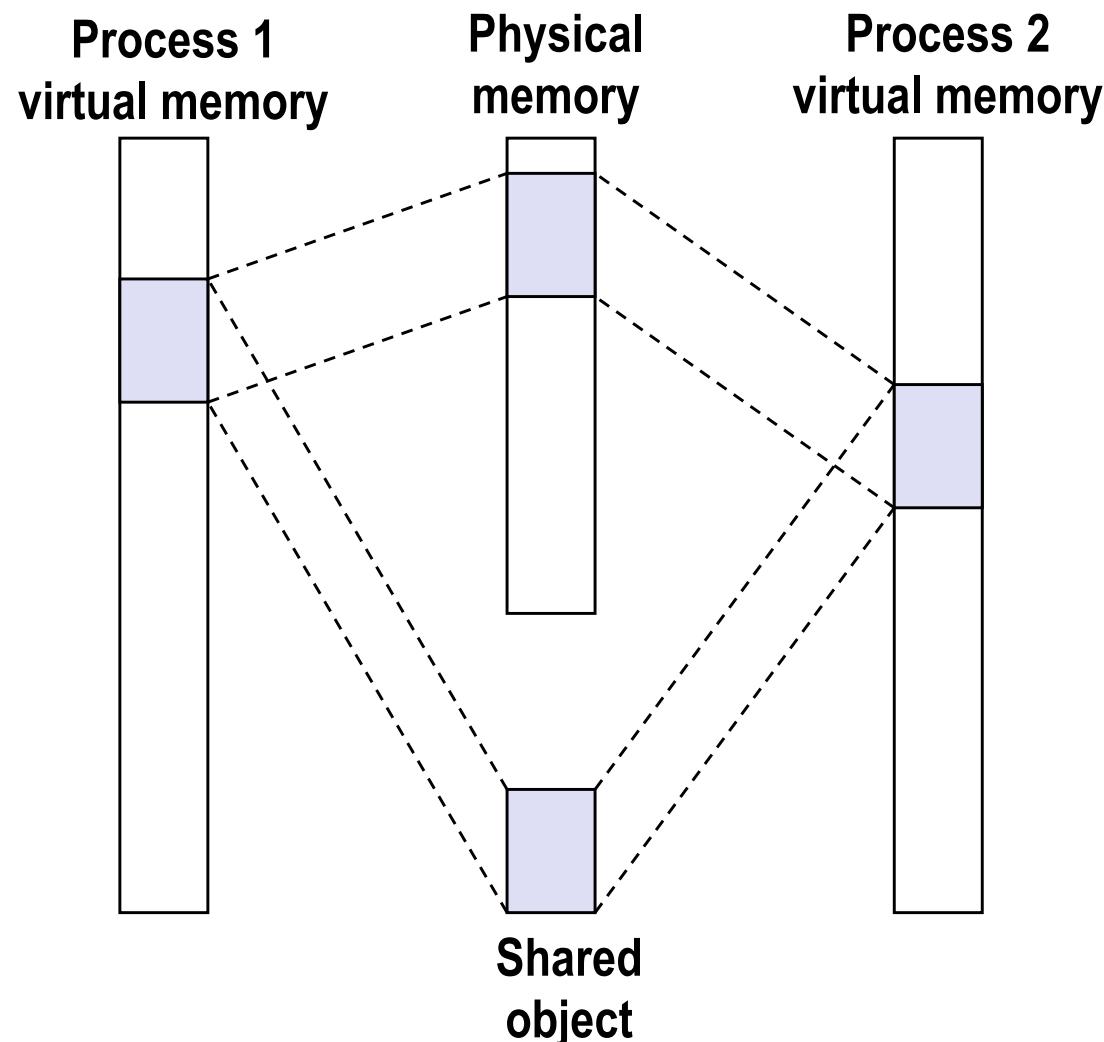
- VM areas initialized by associating them with disk objects.
 - Process is known as *memory mapping*.
- Area can be *backed by* (i.e., get its initial values from) :
 - *Regular file* on disk (e.g., an executable object file)
 - Initial page bytes come from a section of a file
 - *Anonymous file* (e.g., nothing)
 - First fault will allocate a physical page full of 0's (*demand-zero page*)
 - Once the page is written to (*dirtied*), it is like any other page
- Dirty pages are copied back and forth between memory and a special *swap file*.

Sharing Revisited: Shared Objects



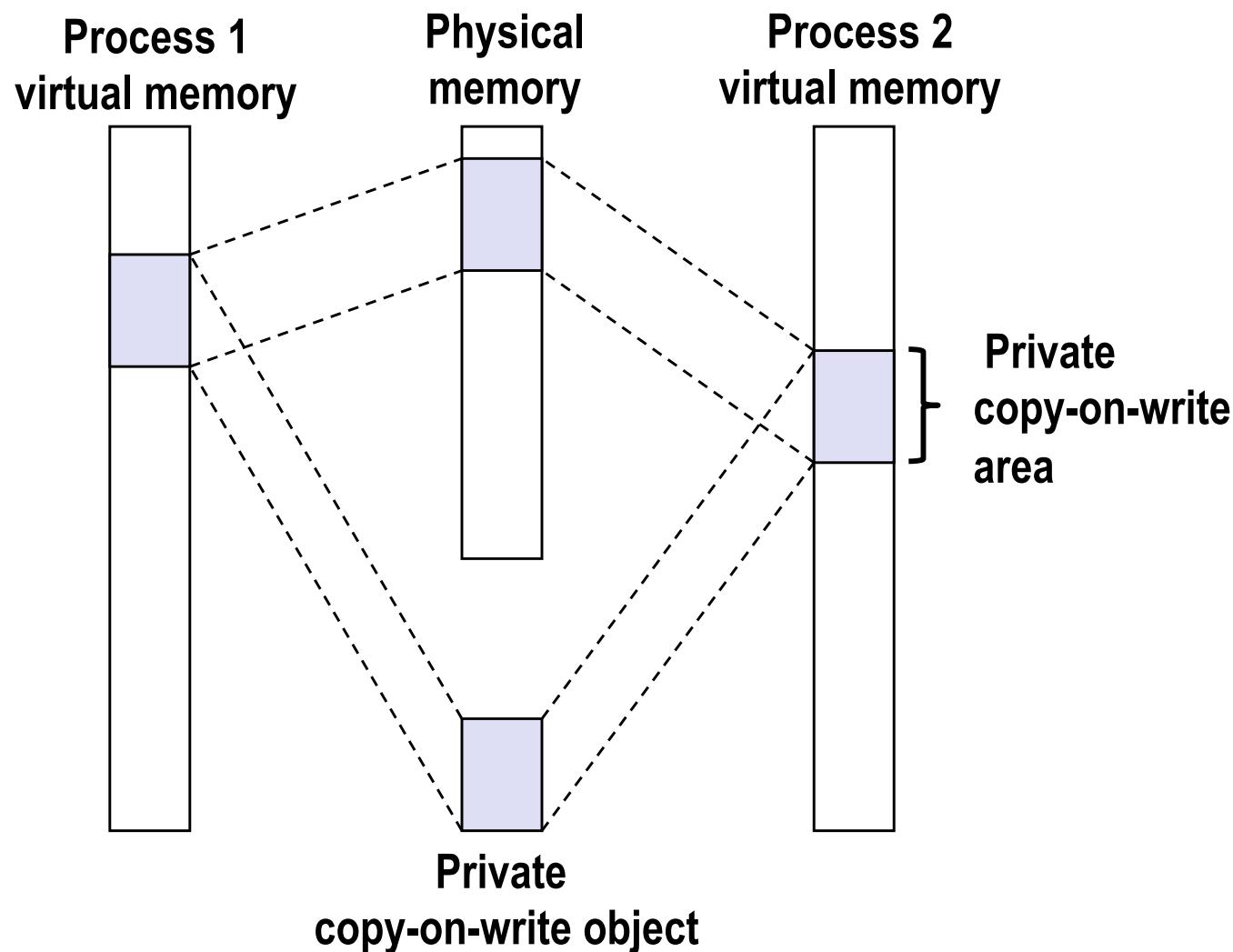
- **Process 1 maps the shared object.**

Sharing Revisited: Shared Objects



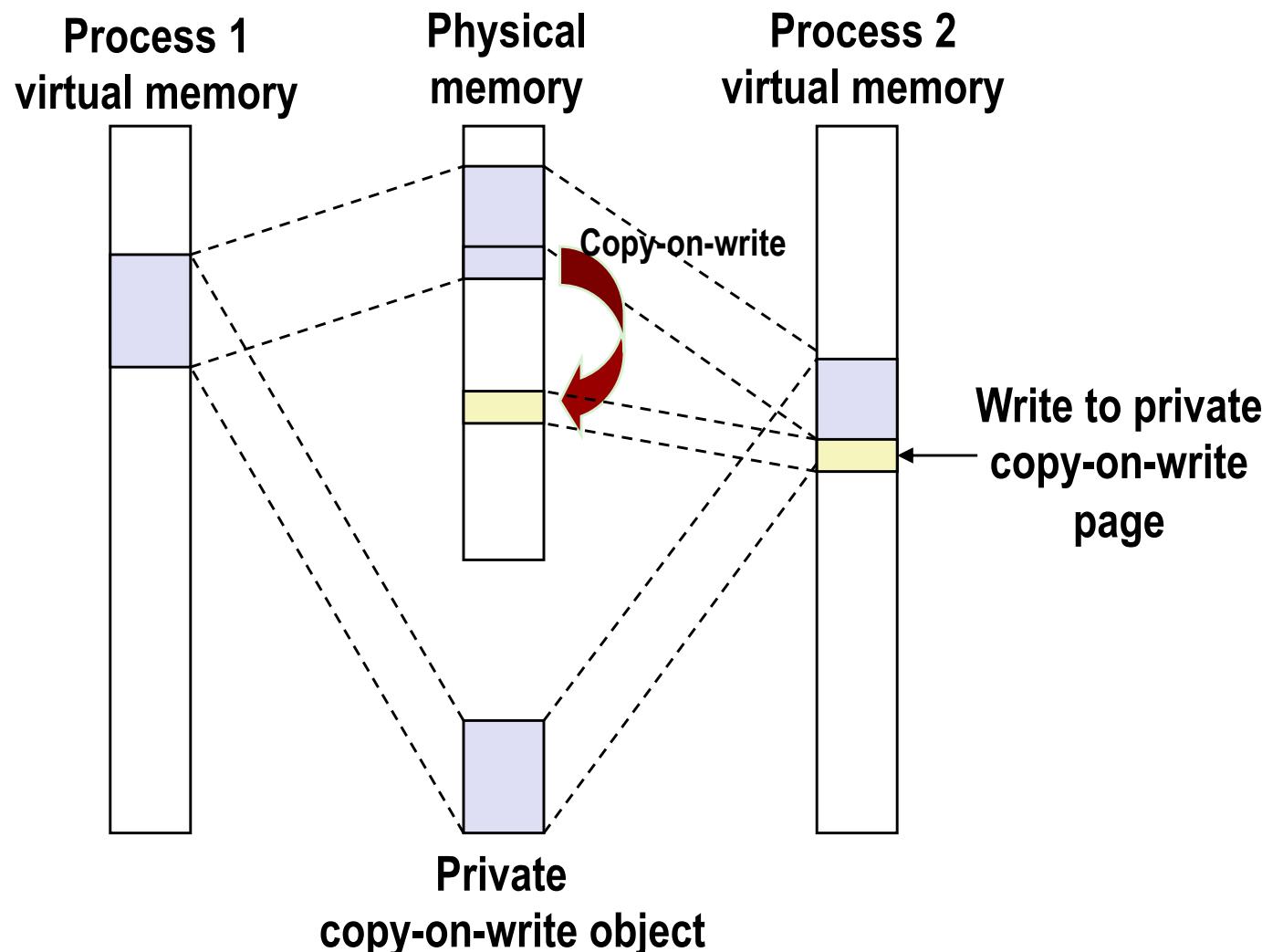
- **Process 2 maps the shared object.**
- **Notice how the virtual addresses can be different.**

Sharing Revisited: Private Copy-on-write (COW) Objects



- Two processes mapping a *private copy-on-write (COW)* object.
- Area flagged as private copy-on-write
- PTEs in private areas are flagged as read-only

Sharing Revisited: Private Copy-on-write (COW) Objects

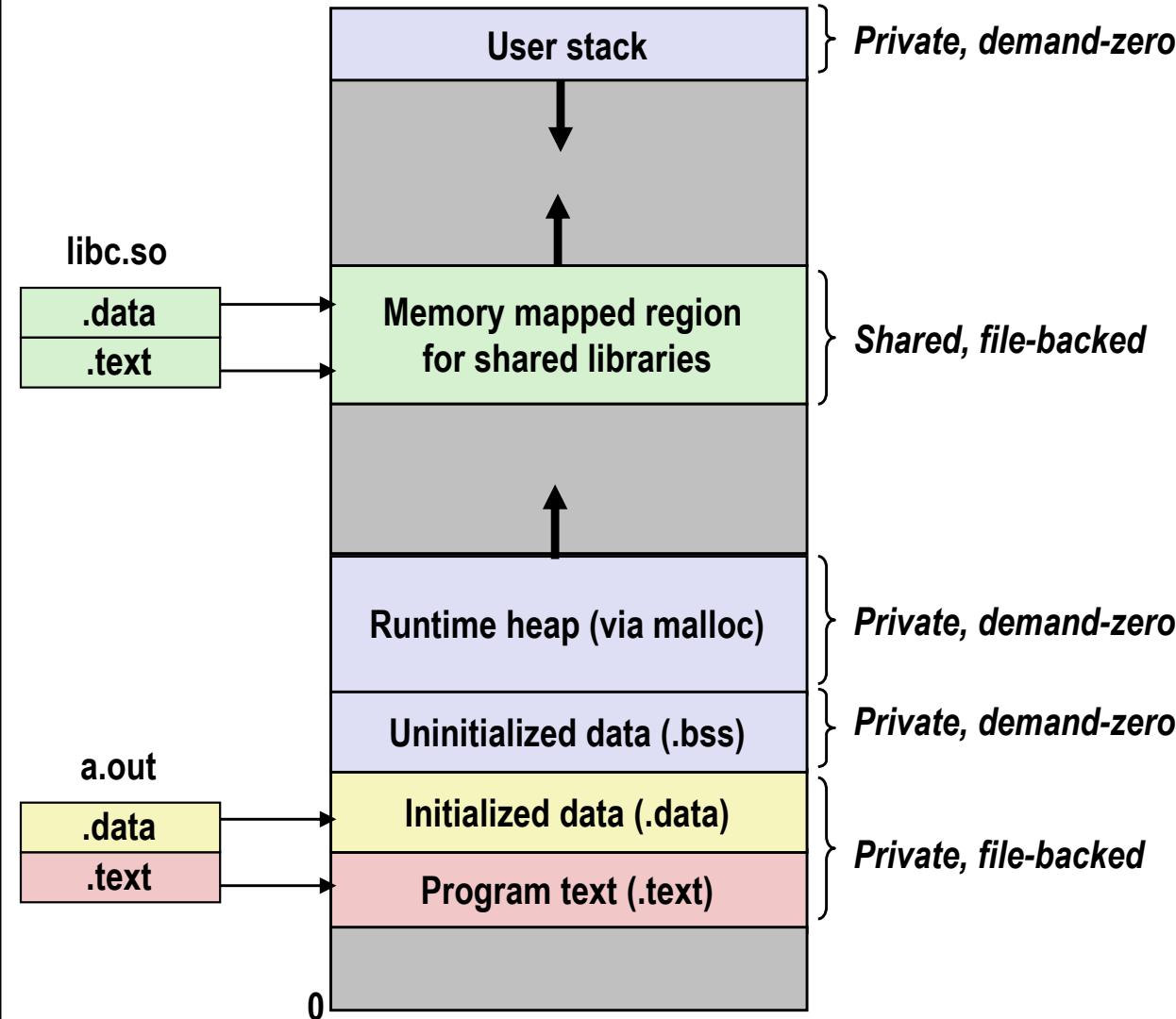


- Instruction writing to private page triggers protection fault.
- Handler creates new R/W page.
- Instruction restarts upon handler return.
- Copying deferred as long as possible!

The `fork` Function Revisited

- VM and memory mapping explain how `fork` provides private address space for each process.
- To create virtual address for new new process
 - Create exact copies of current `mm_struct`, `vm_area_struct`, and page tables.
 - Flag each page in both processes as read-only
 - Flag each `vm_area_struct` in both processes as private COW
- On return, each process has exact copy of virtual memory
- Subsequent writes create new pages using COW mechanism.

The execve Function Revisited



- To load and run a new program `a.out` in the current process using `execve`:
- Free `vm_area_struct`'s and `page tables` for old areas
- Create `vm_area_struct`'s and `page tables` for new areas
 - Programs and initialized data backed by object files.
 - `.bss` and stack backed by anonymous files .
- Set PC to entry point in `.text`
 - Linux will fault in code and data pages as needed.

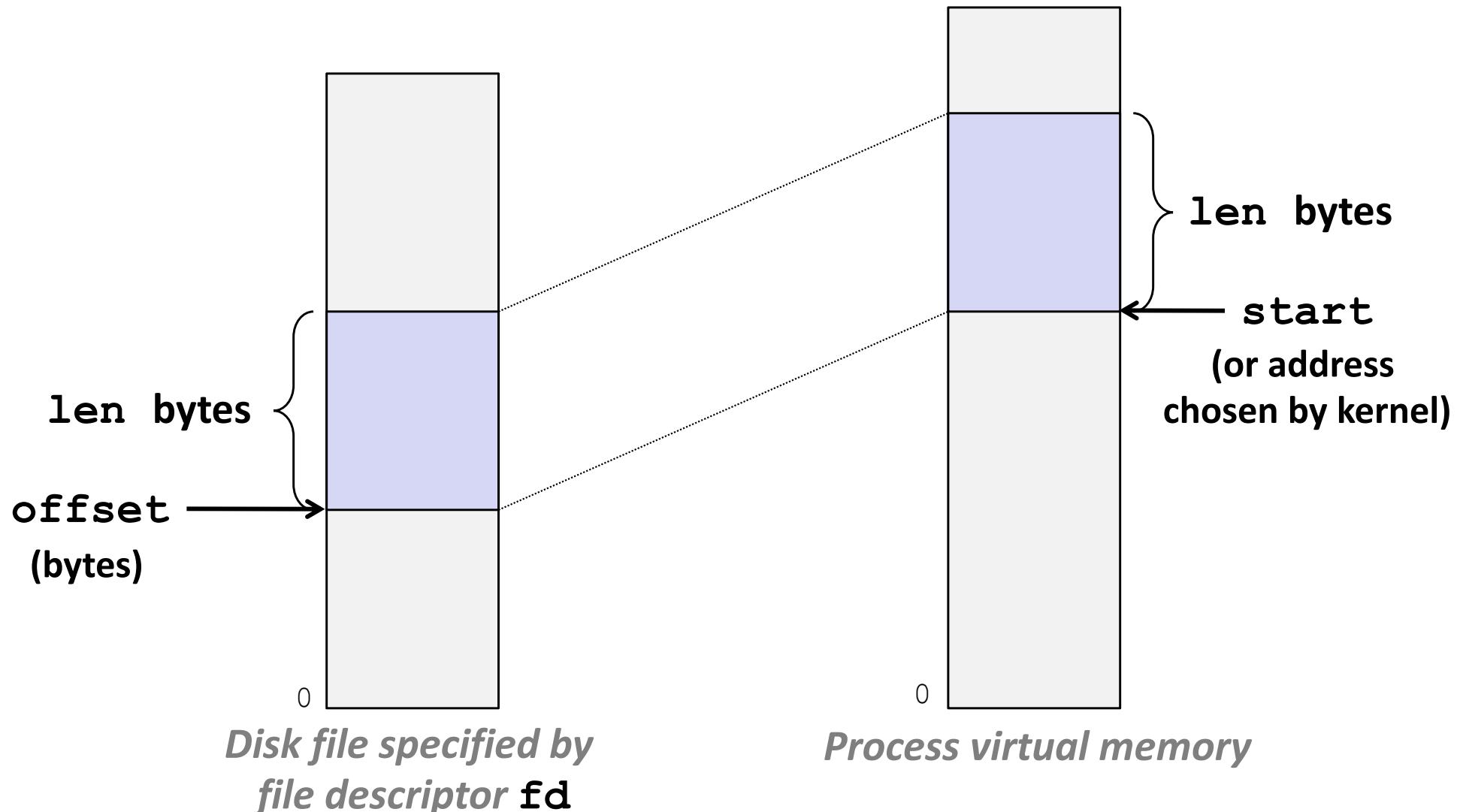
User-Level Memory Mapping

```
void *mmap(void *start, int len,  
           int prot, int flags, int fd, int offset)
```

- Map **len** bytes starting at offset **offset** of the file specified by file description **fd**, preferably at address **start**
 - **start**: may be 0 for “pick an address”
 - **prot**: PROT_READ, PROT_WRITE, ...
 - **flags**: MAP_ANON, MAP_PRIVATE, MAP_SHARED, ...
- Return a pointer to start of mapped area (may not be **start**)

User-Level Memory Mapping

```
void *mmap(void *start, int len,  
           int prot, int flags, int fd, int offset)
```



Example: Using mmap to Copy Files

- Copying a file to stdout without transferring data to user space .

```
#include "csapp.h"

void mmapcopy(int fd, int size)
{
    /* Ptr to memory mapped area */
    char *bufp;

    bufp = Mmap(NULL, size,
                PROT_READ,
                MAP_PRIVATE,
                fd, 0);
    Write(1, bufp, size);
    return;
}
```

mmapcopy.c

```
/* mmapcopy driver */
int main(int argc, char **argv)
{
    struct stat stat;
    int fd;

    /* Check for required cmd line arg */
    if (argc != 2) {
        printf("usage: %s <filename>\n",
               argv[0]);
        exit(0);
    }

    /* Copy input file to stdout */
    fd = Open(argv[1], O_RDONLY, 0);
    Fstat(fd, &stat);
    mmapcopy(fd, stat.st_size);
    exit(0);
}
```

mmapcopy.c