

# **CMPE110 Lecture 17**

## **Virtual Memory**

---

Heiner Litz

<https://canvas.ucsc.edu/courses/12652>



# Announcements

---

- Review Session on Friday 11<sup>th</sup>
- Call for topics:
  - <https://piazza.com/class/jf0roanugdd3pf?cid=342>

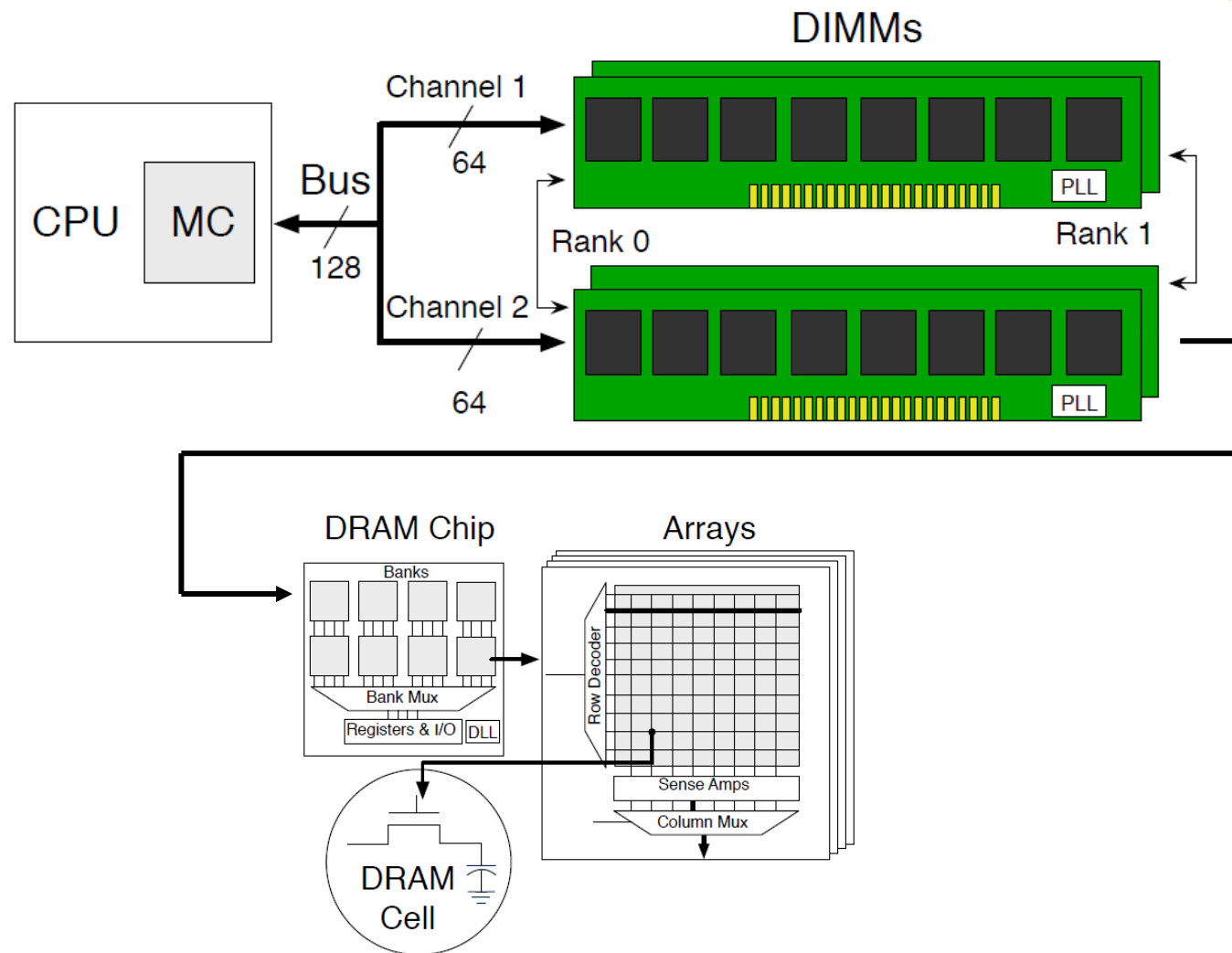
# Review

---



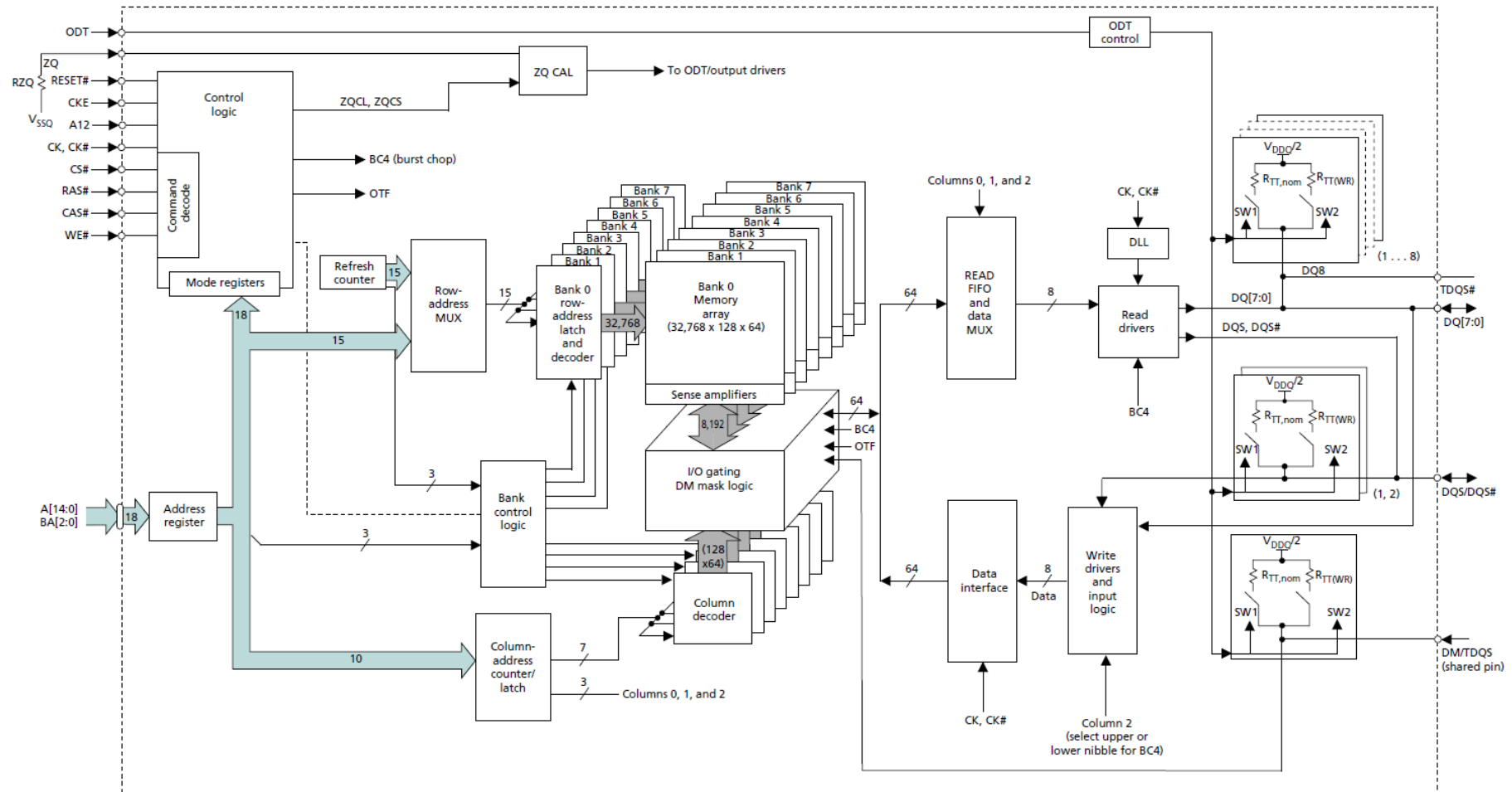


# Main Memory Overview





# 2Gb x8 DDR3 Chip [Micron]





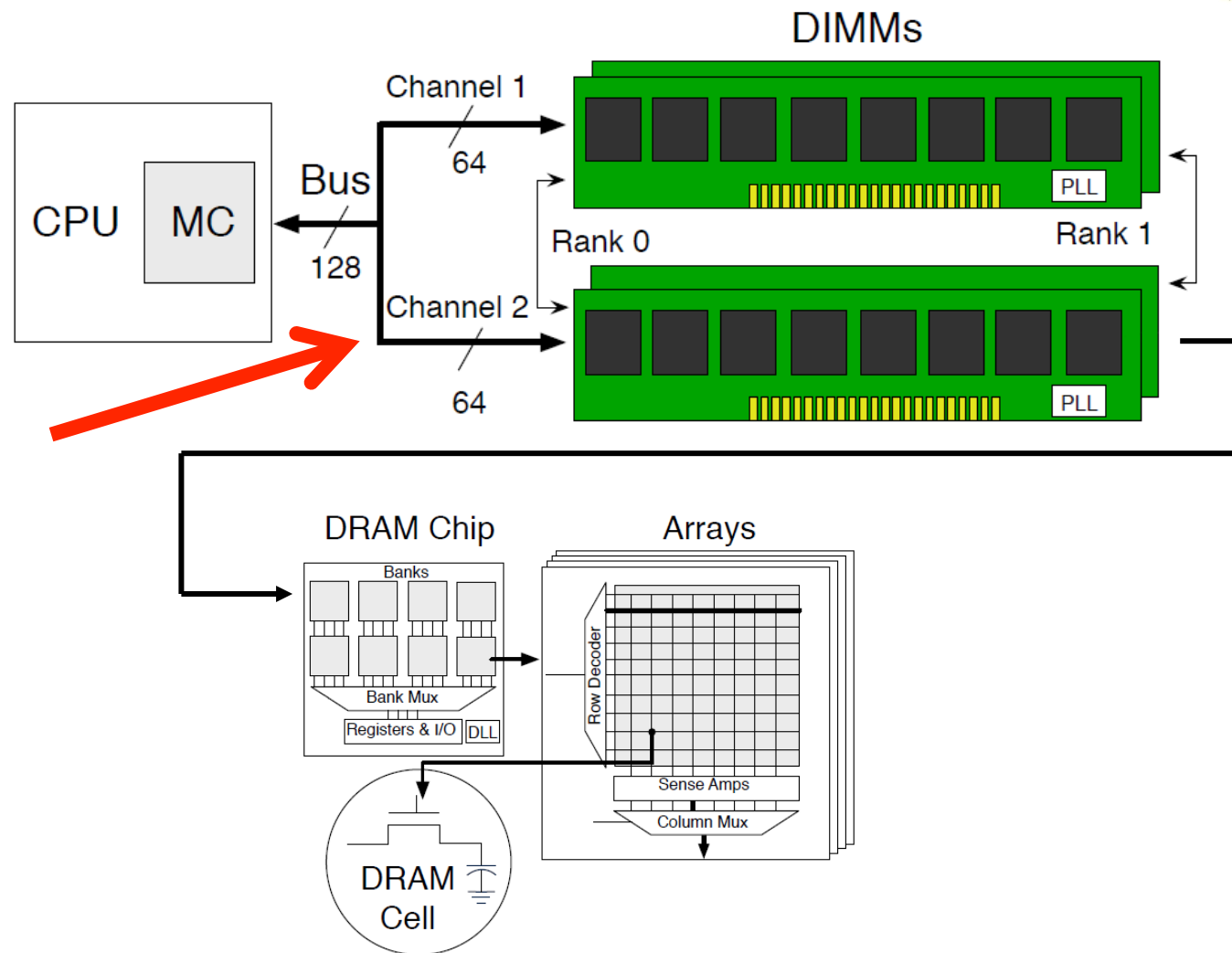


# DRAM Capacity

- What if we want even more capacity?
  - Can only put limited number of ranks on channel (signal and driver limitations)
- x2 chips? x1 chips?
  - What happens to power consumption?
- More channels?



# Main Memory Overview







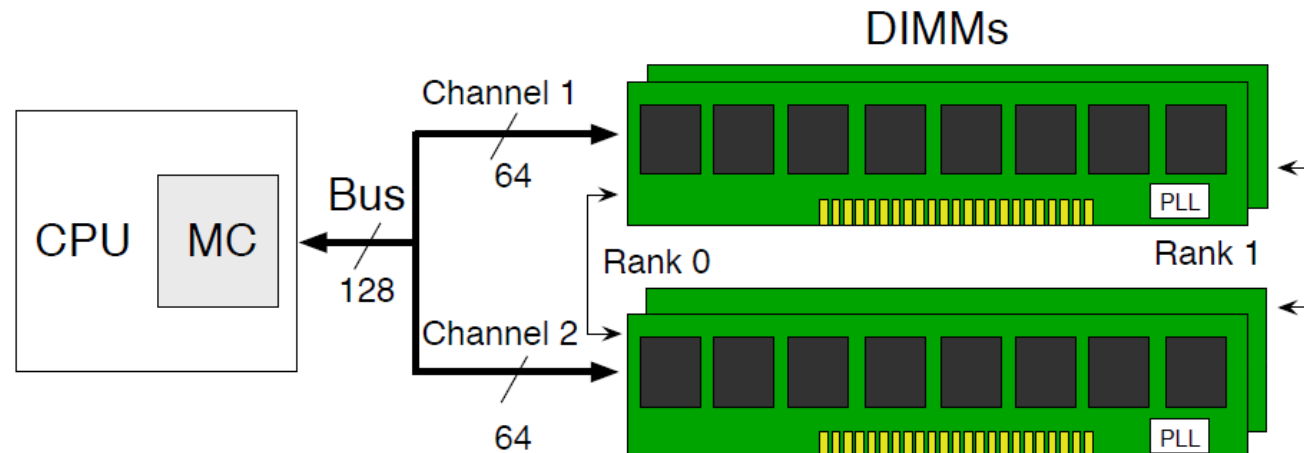
# DRAM Channels

- Channel: a set of DIMMs in series
  - All DIMMs get the same command, one of the ranks replies
- Multiple-channel options
  - Multiple lock-step channels
    - Single “channel” with wider interface (faster cache line refill!)
    - Sometimes called “Gang Mode”
    - Only works if DIMMs are identical (organization, timing)
  - Multiple independent channels
    - Requires multiple controllers
- Tradeoffs in having multiple channels
  - Cost: pins, wires, controllers
  - Benefit: higher bandwidth, capacity

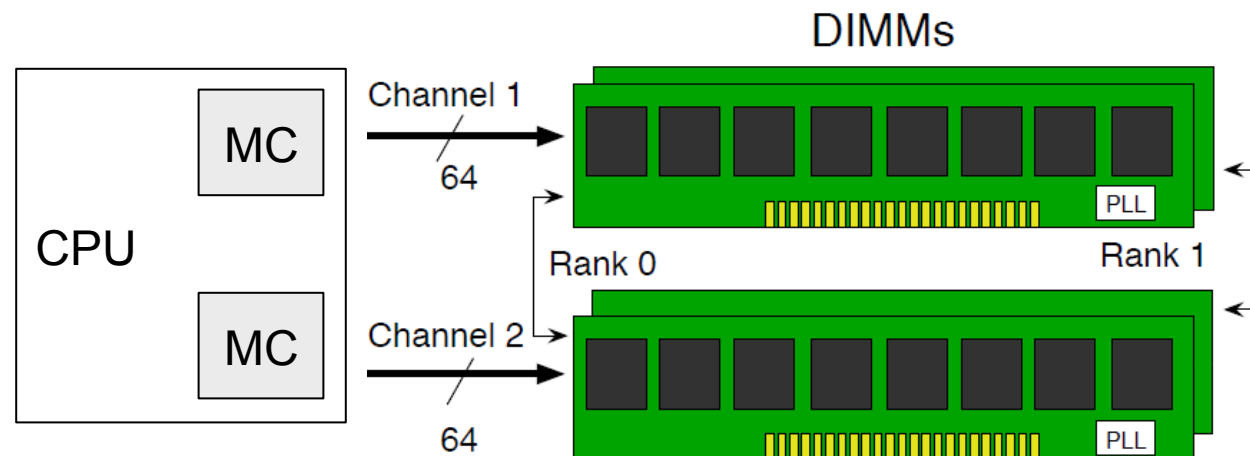


# DRAM Channel Options

Lock-step



Independent





# Virtual Memory

---

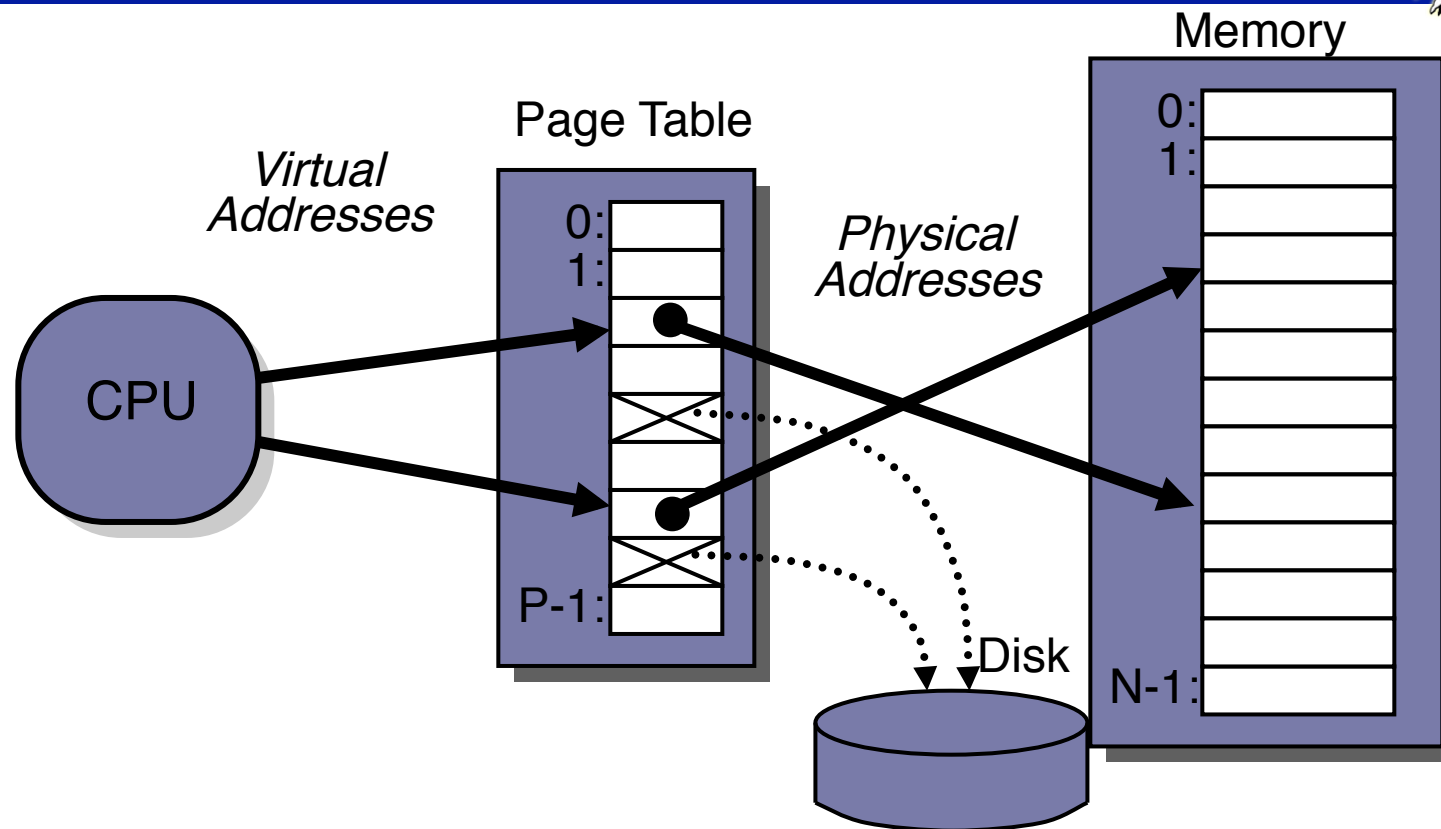


# Virtual Memory

- How to share/use physical memory by multiple applications?
- How to guarantee that applications can only access “their” memory (security)
- How to expose virtually unlimited amount of main memory to applications?



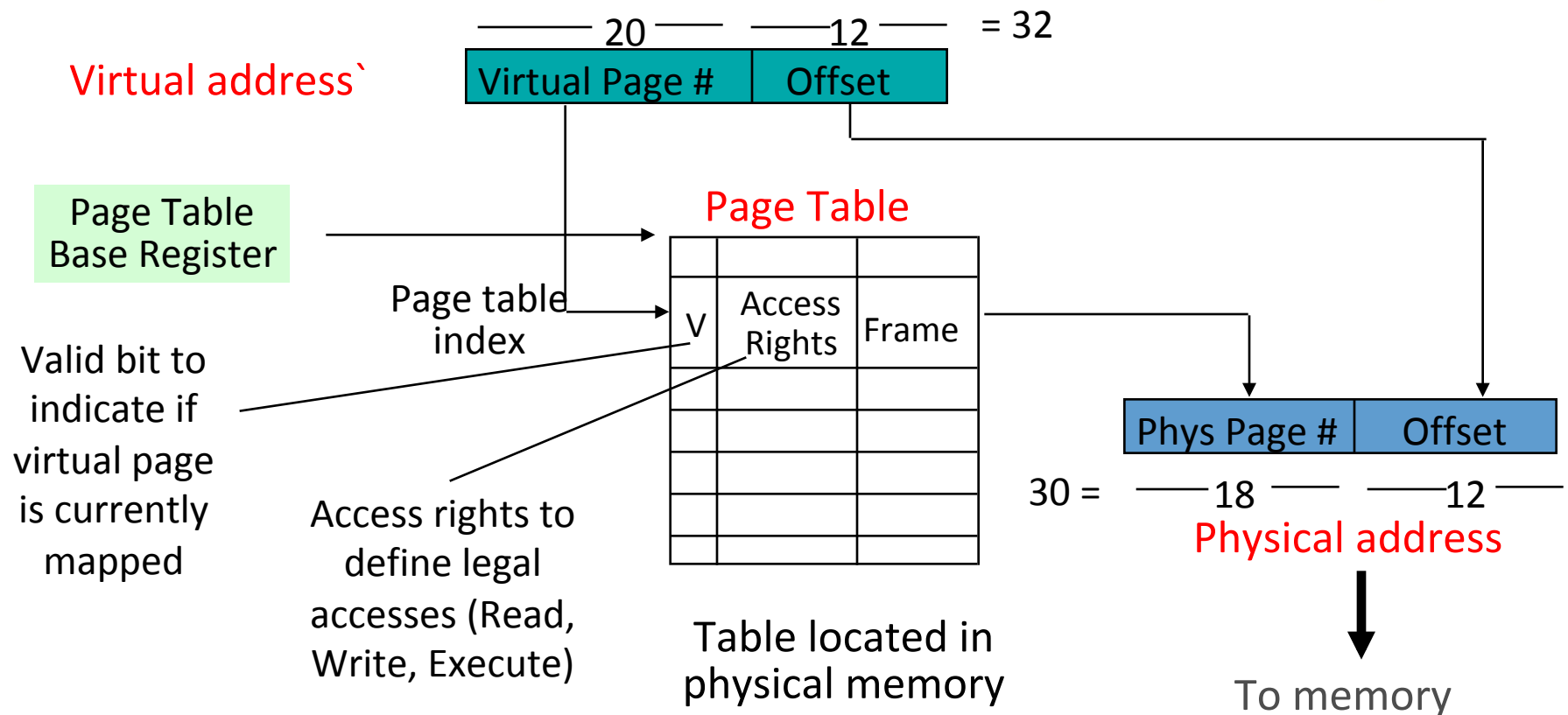
# Virtual Memory



- HW translates addresses using an OS-managed lookup table
- Indirection allows redirection and checks!



# Page Table

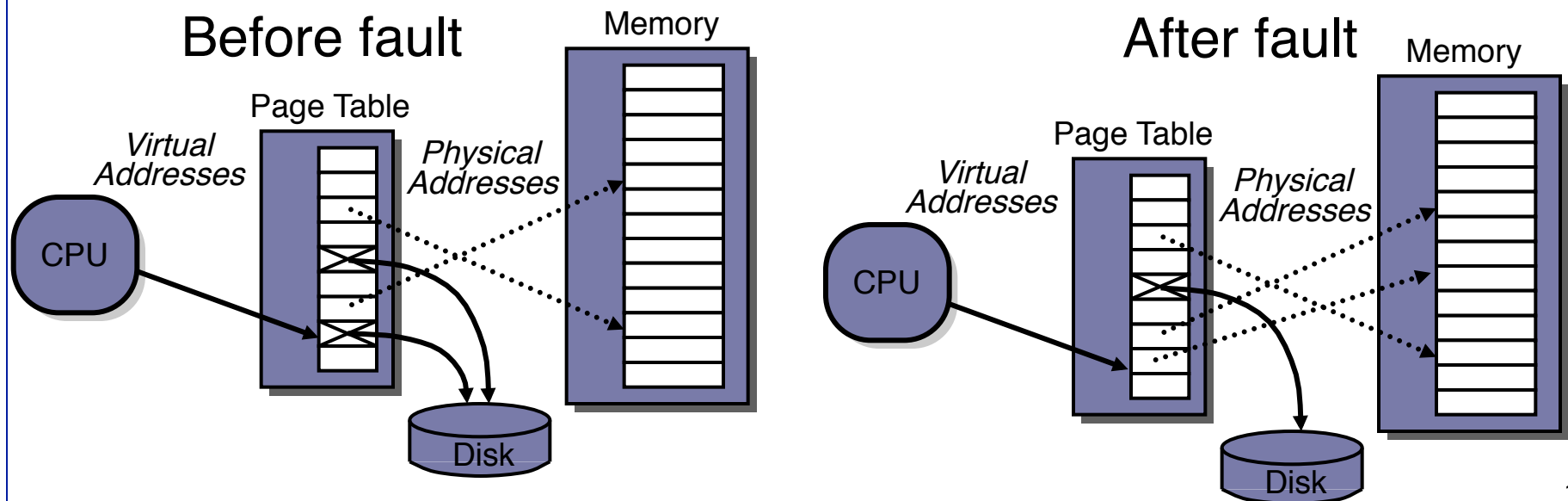


- One page table per process stored in memory
  - Process == instance of program
- One entry per virtual page number



# Page Faults

- Page fault → an exception case
  - HW indicates exception cause and problematic address
  - OS handler invoked to move data from disk into memory
    - Current process suspends, others can resume
    - OS has full control over placement
  - When process resumes, repeat load or store





# Wait, How About Performance?

- Virtual memory is great but
- We just doubled the memory accesses
  - A load requires an access to the page table first
  - Then an access to the actual data
- How can we do translation fast?
  - Without an additional memory access?



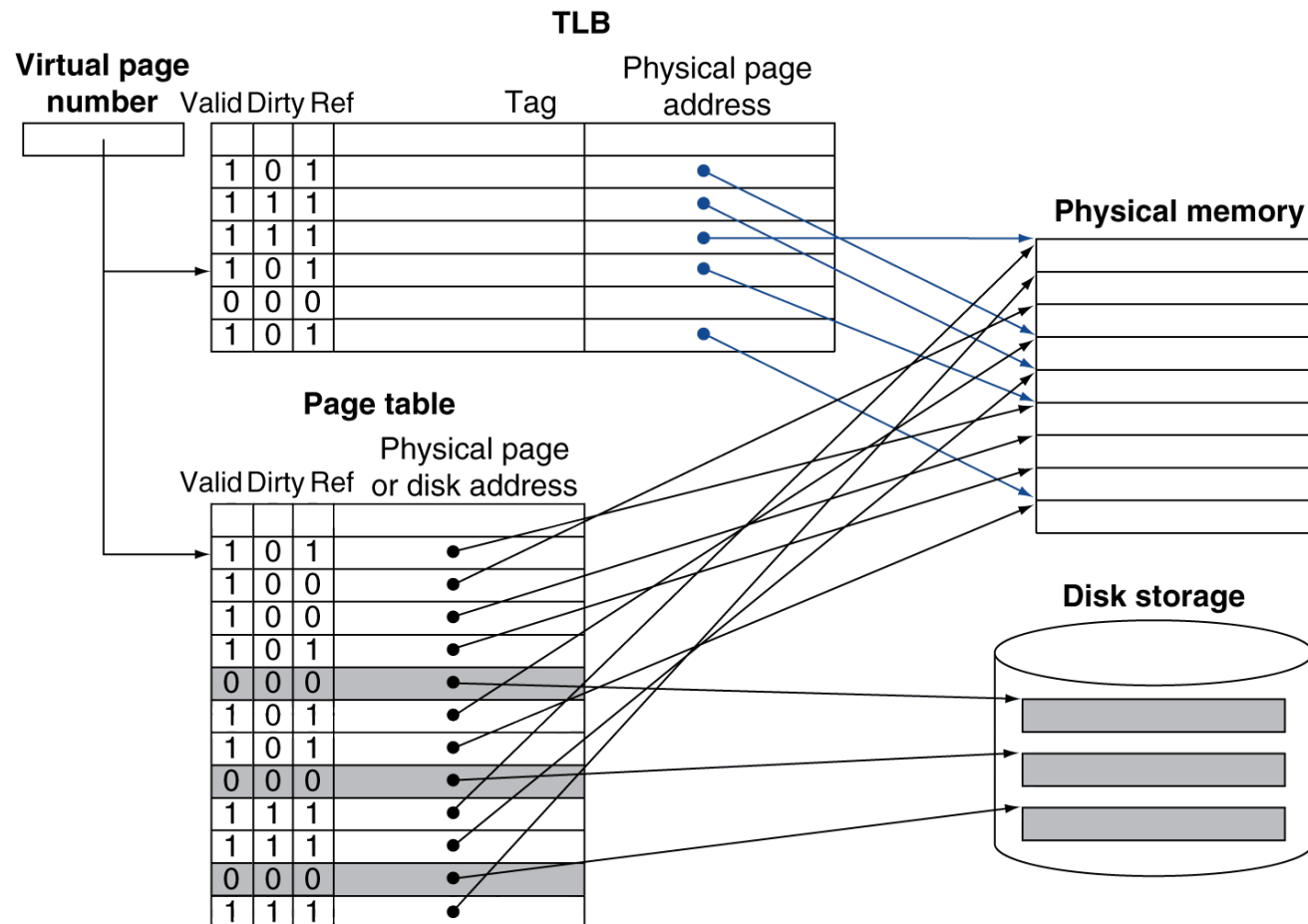


# Translation Look-aside Buffer

- TLB = a hardware cache just for translation entries
  - A hardware cache specializing in page table entries
- Key idea: locality in accesses → locality in translations
- TLB design: similar issues to all caches
  - Basic parameters: capacity, associativity replacement policy
  - Basic optimizations: instruction/data TLBs, multi-level TLBs, ...
  - Misses may be handled by HW or SW
    - x86: hardware services misses
    - MIPS: software services misses through exception



# TLB Organization





# TLB Entries

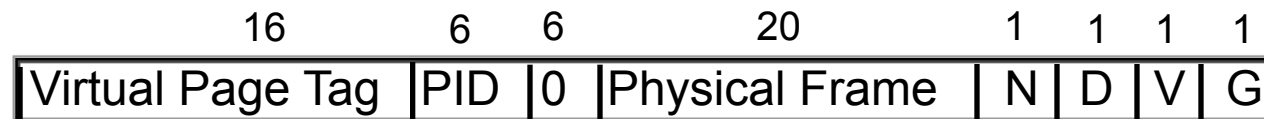
- Each TLB entry stores a page table entry (PTE)
- TLB entry data → PTE entry fields
  - Physical page numbers
  - Permission bits (RXW)
  - Other PTE information (dirty bit, etc)
- The TLB entry metadata
  - Tag: portion of virtual page # not used to index the TLB
    - Depends on the TLB associativity
  - Valid bit
  - LRU bits
    - If TLB is associative and LRU replacement is used



# Example TLB

## ■ TLB entry design

- Addresses are 32 bits with 4 KB pages (12 bit offset)
- TLB has 64 entries, 4-way set associative
- Each entry is 42 bits wide:



PID	Process ID
N	Do not cache memory address (optional)
D	Dirty bit
V	Valid bit
G	Global (valid regardless of PID)



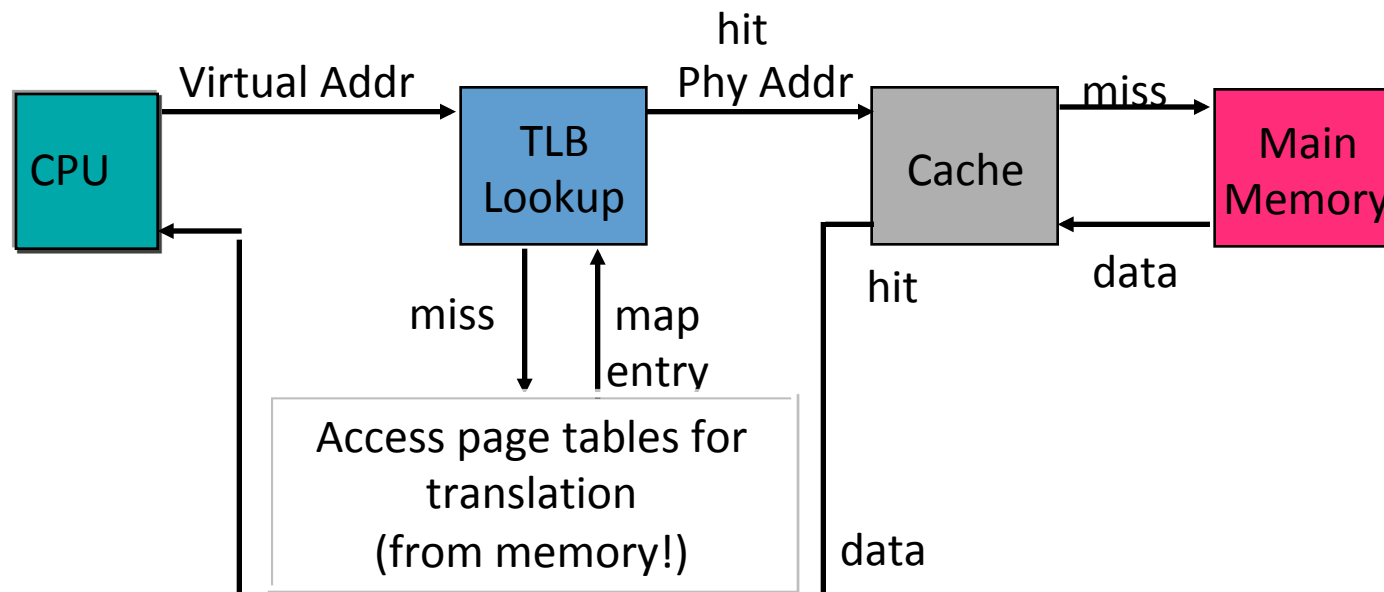
# TLB Caveats: Context Switching

- What happens to TLB when switching between processes?
- The OS must flush the entries in the TLB
  - Large number of TLB misses after every switch
- Alternatively, use a process ID each TLB entry
  - PID or ASID
  - Allows entries from multiple programs to co-exist
  - Gradual replacement



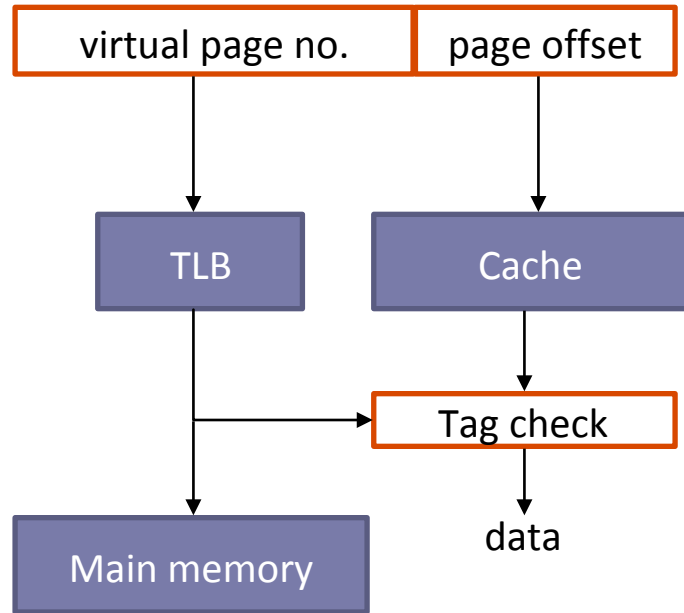
# TLB & Memory Hierarchies

- Basic process
  - Use TLB to get VA  $\rightarrow$  PA
  - Use PA to access caches and DRAM
- Question: can you ever access the TLB and the cache in parallel?





# Virtually Indexed, Physically Tagged Caches

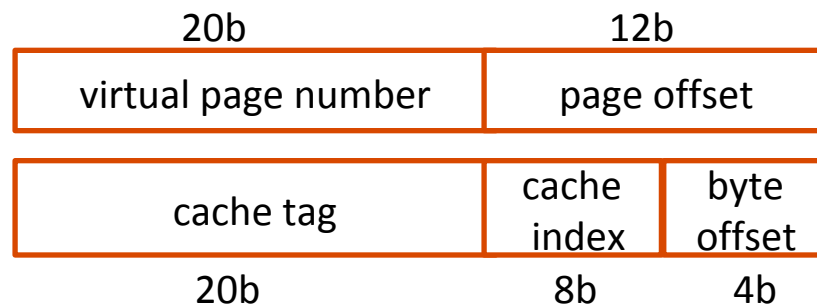


## ■ Translation & cache access in parallel

- Start access to cache with page offset
- Tag check used physical address

## ■ Only works when

- VPN bits not needed for cache lookup
- $\text{Cache Size} \leq \text{Page Size} * \text{Associativity}$ 
  - I.e.  $\text{Set Size} \leq \text{Page Size}$
- Ok, we want L1 to be small anyway





# Virtual Memory and Caches

TLB	Cache	DRAM	Possible?
hit	miss	hit	
miss	hit	hit	
miss	miss	hit	
miss	miss	miss	
hit	miss	miss	
hit	hit	miss	
miss	hit	miss	



# TLB Caveats: Limited Reach



- 64 entry TLB with 4KB pages maps 256 KB
  - Smaller than many L3 caches in most systems
  - TLB miss rate > L2 miss rate!
  
- Potential solutions
  - Larger pages
  - Multilevel TLBs (just like multi-level caches)



# Page Size Tradeoff

## Larger pages

Pros: smaller page tables, fewer page faults and more efficient transfer with larger applications, improved TLB coverage

Cons: higher internal fragmentation

## Smaller pages

Pros: improved time to start small processes with fewer pages, internal fragmentation is low (important for small programs)

Cons: high overhead in large page tables

## General trend toward larger pages

1978: 512 B, 1984: 4 KB, 1990: 16 KB, 2000: 64 KB, 2010: 4MB



# Multiple Page Sizes

Many machines support multiple page sizes

SPARC: 8KB, 64KB, 1 MB, 4MB

MIPS R4000: 4KB – 16 MB

x86: 4KB, 4MB, 1GB

Page size dependent upon application

OS kernel uses large pages

Most user applications use smaller pages

## Issues

Software complexity

TLB complexity

How do you do match VPN if not sure about the page size?

# Final Problem: Page Table Size



Page table size is proportional to size of address space

x86 example: virtual addresses are 32 bits, pages are 4 KB

Total number of pages:  $2^{32} / 2^{12} = 1$  Million

Page Table Entry (PTE) are 32 bits wide

Total page table size is therefore  $1\text{M} \times 4 \text{ bytes} = 4 \text{ MB}$

But, only a small fraction of those pages are actually used!

## Why is this a problem?

The page table must be resident in memory (why?)

What happens for the 64-bit version of x86?

What about running multiple programs?

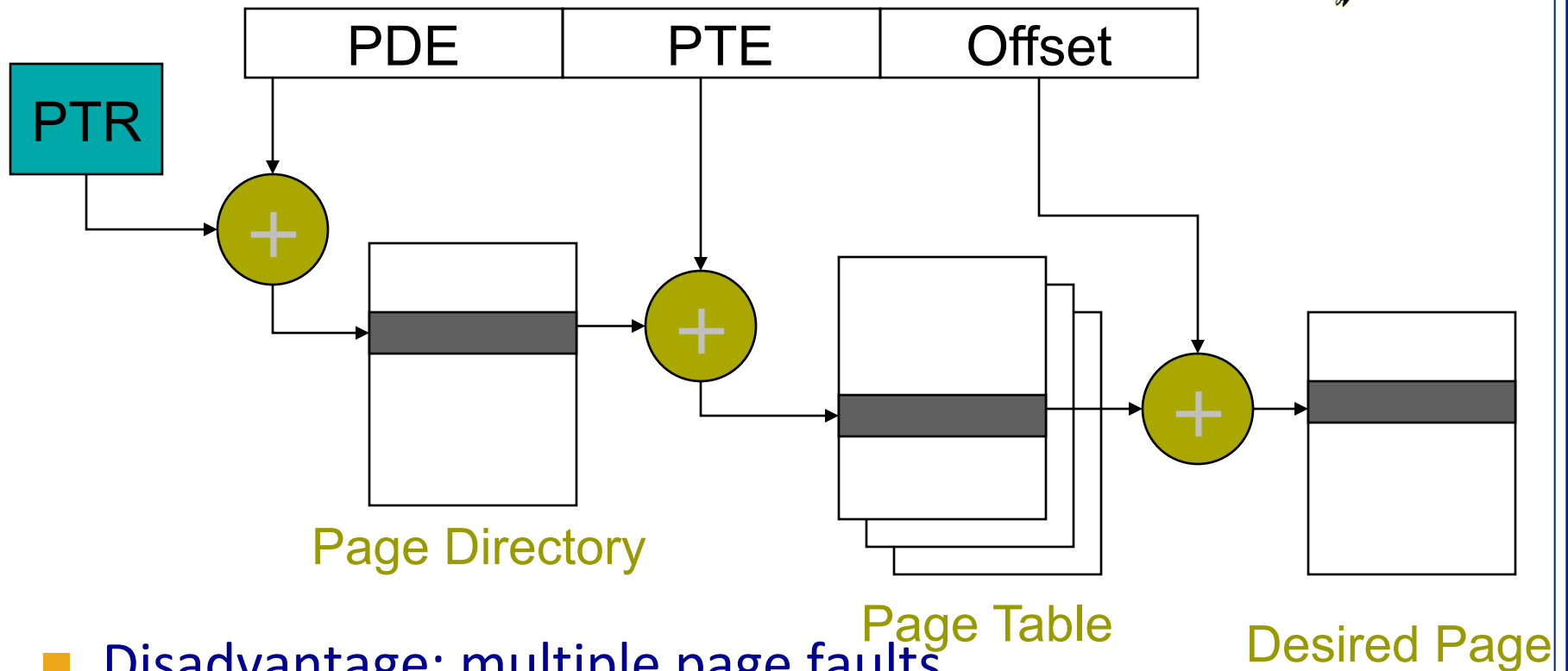


# Solution: Multi-Level Page Tables

- Use a hierarchical page table structure
- Example: Two levels
  - First level: directory entries
  - Second level: actual page table entries
- Only top level must be resident in memory
- Remaining levels can be in memory, on disk, or unallocated
  - Unallocated if corresponding ranges of the virtual address space are not used



# Two-level Page Tables



- Disadvantage: multiple page faults
  - Accessing a PTE page table can cause a page fault
  - Accessing the actual page can cause a second page fault
  - TLB plays an even more important role



# Putting it All Together

---



# Example: Intel P6 Processor

- 32 bit address space with 4KB pages
- 2-level cache hierarchy
  - L1 I/D: 16KB, 4-way, 128 sets, 32B blocks
  - L2: 128KB – 2MB
- TLBs
  - I-TLB: 32 entries, 4-way, 8 sets
  - D-TLB: 64 entries, 4-way, 16 sets
  - 32 entries, 8 sets

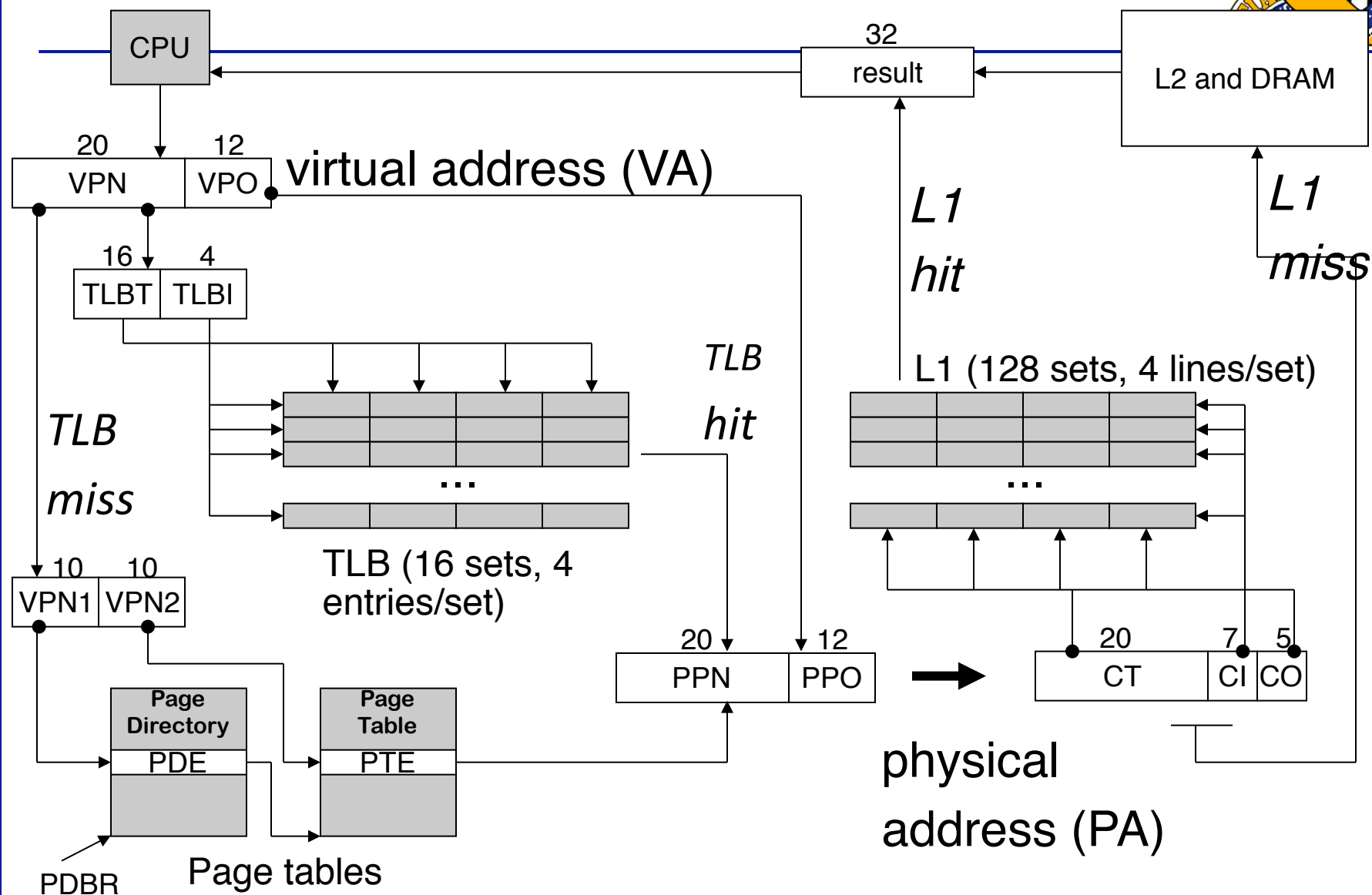




# Abbreviations

- MMU
  - Memory management unit: controls TLB, handles TLB misses
- 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

# Overview of P6 Address Translation





# P6 2-level Page Table Structure

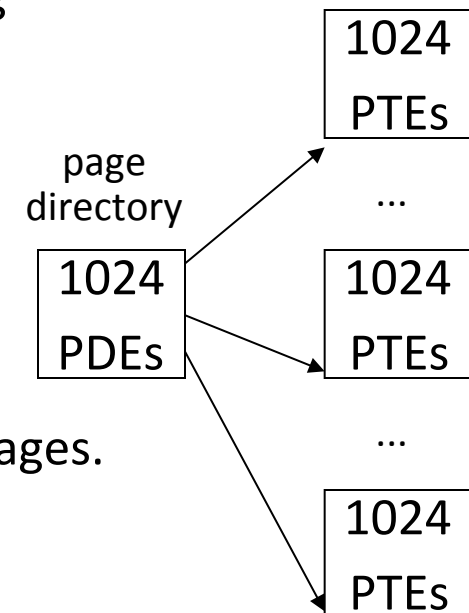
## ■ Page directory

- One page directory per process.
- 1024 4-byte page directory entries (PDEs) that point to page tables
- Page directory must be in memory if process running
- Always pointed to by PDBR

Up to  
1024  
page  
tables

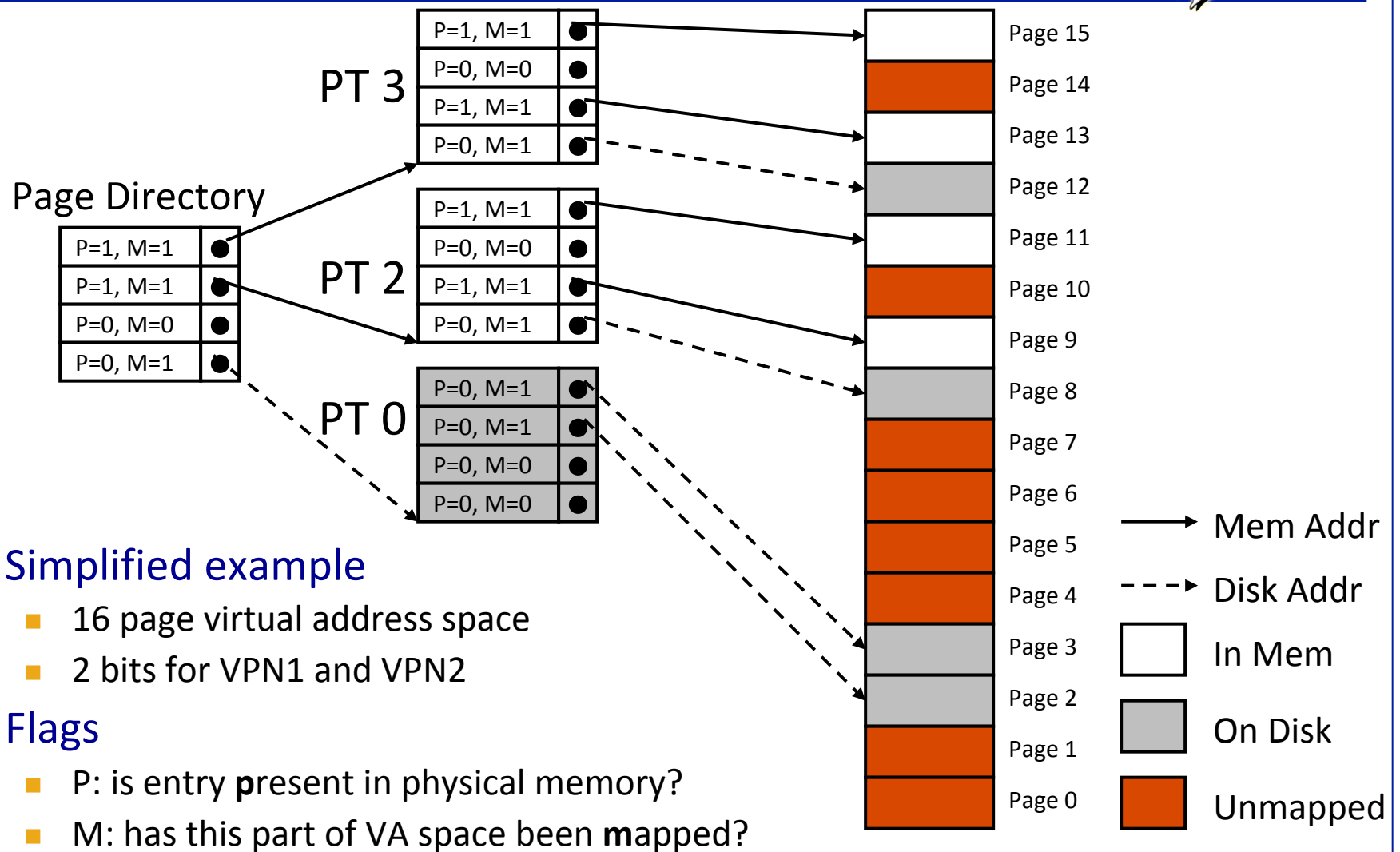
## ■ Page tables

- 1024 4-byte page table entries (PTEs) that point to pages.
- Page tables can be paged in and out





# Representation of Virtual Address Space



## ■ Simplified example

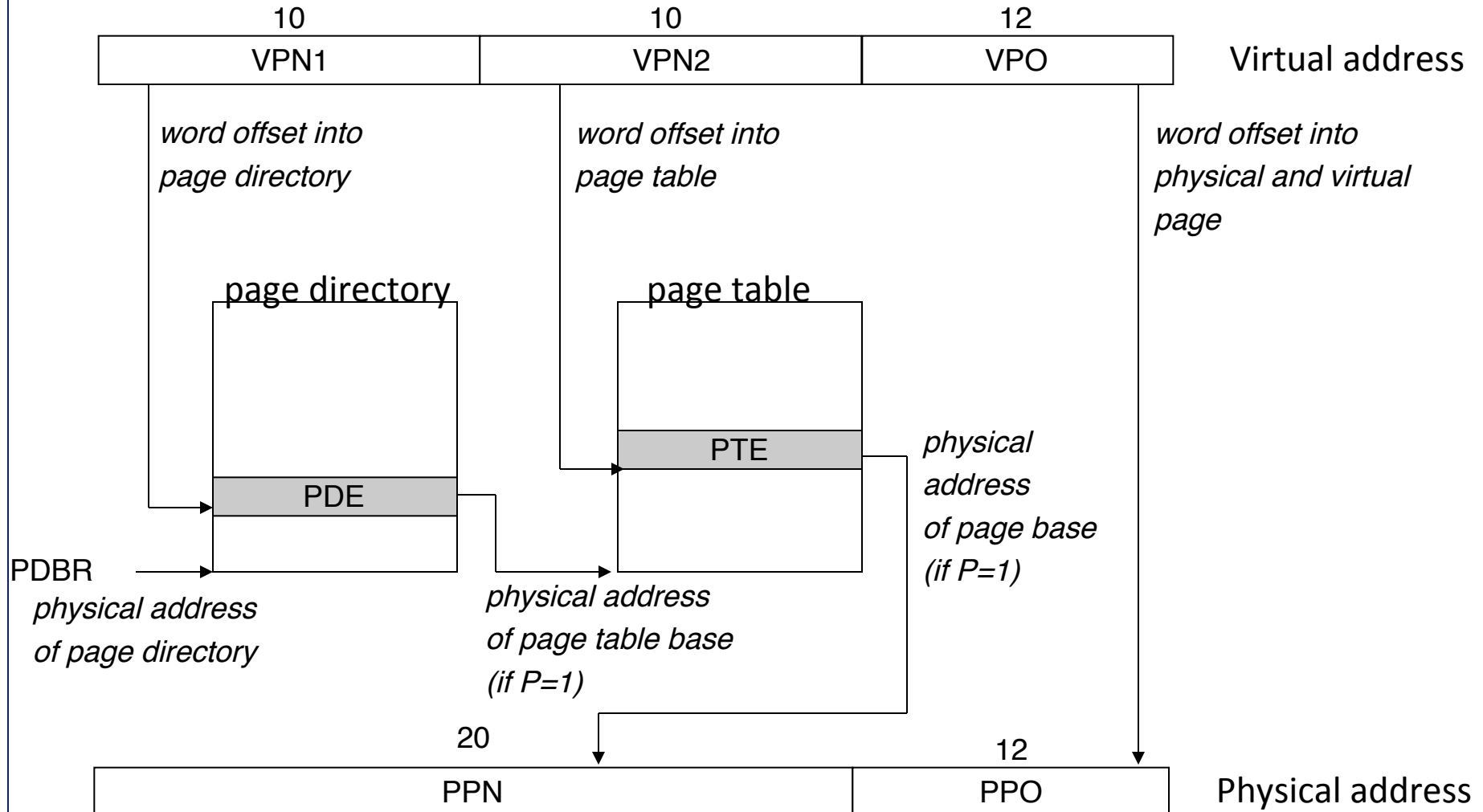
- 16 page virtual address space
- 2 bits for VPN1 and VPN2

## ■ Flags

- P: is entry **p**resent in physical memory?
- M: has this part of VA space been **m**apped?



# P6 Translation





# P6 Page Directory Entry (PDE)

31	12	11	9	8	7	6	5	4	3	2	1	0
Page table physical base addr			Avail		G	PS		A	CD	WT	U/S	R/W P=1

Page table physical base address: 20 most significant bits of physical page table address (forces page tables to be 4KB aligned)

Avail: available for system programmers

G: global page (don't evict from TLB on task switch)

PS: page size 4K (0) or 4M (1)

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

CD: cache disabled (1) or enabled (0)

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

U/S: user or supervisor mode access

R/W: read-only or read-write access

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

# P6 page table entry (PTE)



31	12	11	9	8	7	6	5	4	3	2	1	0	
Page physical base address			Avail		G	0	D	A	CD	WT	U/S	R/W	P=1

Page base address: 20 most significant bits of physical page address (forces pages to be 4 KB aligned)

Avail: available for system programmers

G: global page (don't evict from TLB on task switch)

**D: dirty (set by MMU on writes)**

A: accessed (set by MMU on reads and writes)

CD: cache disabled or enabled

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

U/S: user/supervisor

R/W: read/write

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

31	1	0
Available for OS (page location in secondary storage)		P=0