

---

# **Chapter 5**

## **Large and Fast : Exploiting the Memory Hierarchy**

### **Virtual Memory**

# Why Virtual Memory (VM)?

Logical

Physical

Memory with  
infinite capacity



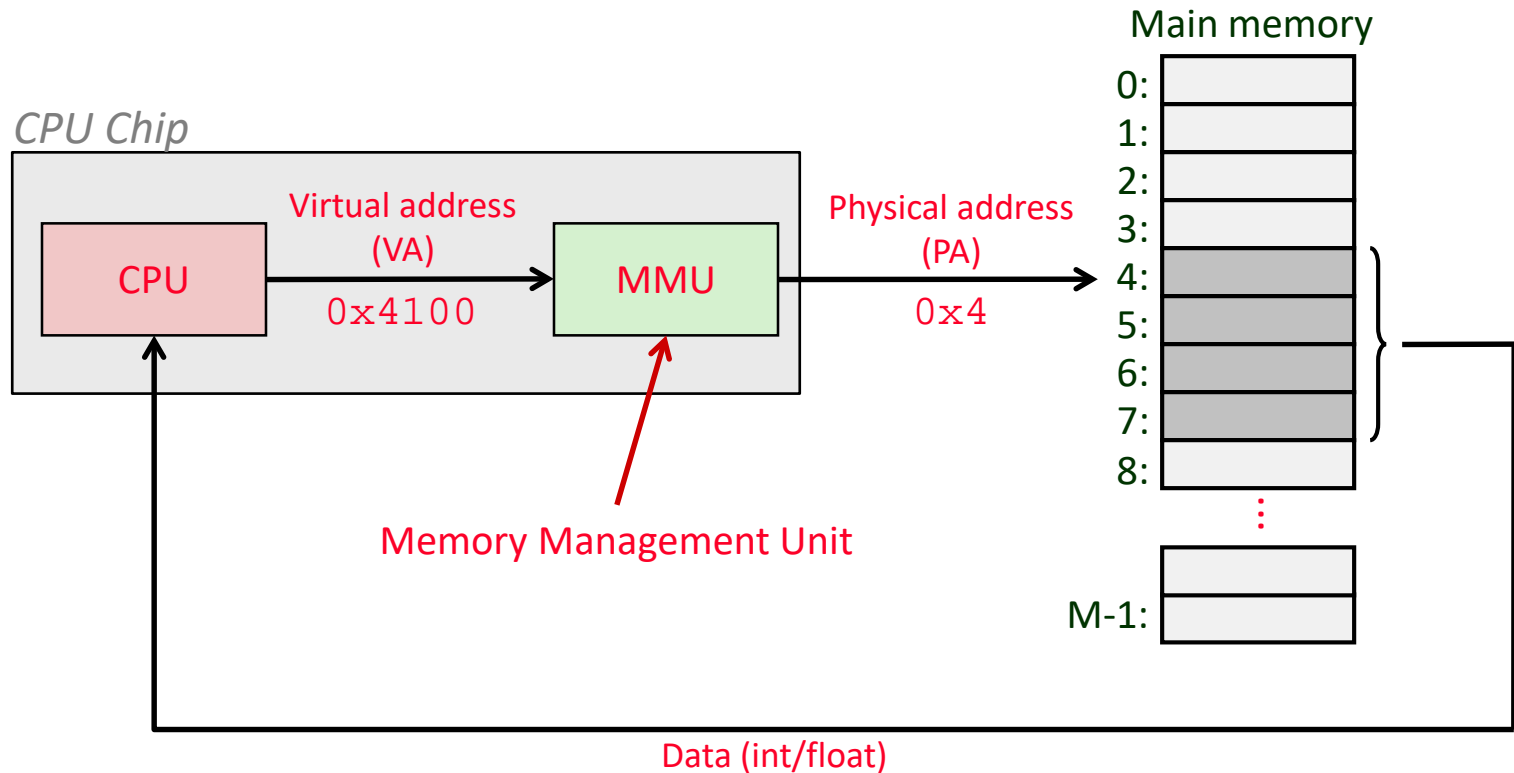
- ❑ Efficient use of limited main memory (RAM)
  - | Use RAM as a cache for the parts of a virtual address space
  - | Managed jointly by CPU hardware and the OS
  - | Keep only active areas of virtual address space in memory
    - Transfer data back and forth as needed
- ❑ Simplifies memory management for programmers
  - | Each process “gets” the same full, private linear address space
- ❑ Isolates address spaces (protection)
  - | One process can’t interfere with another’s memory
    - They operate in *different address spaces*
  - | User process cannot access privileged information
    - Different sections of address spaces have different permissions
- ❑ CPU and OS translate virtual addresses to physical addresses : a page, a page fault

# VM Design Consequences

---

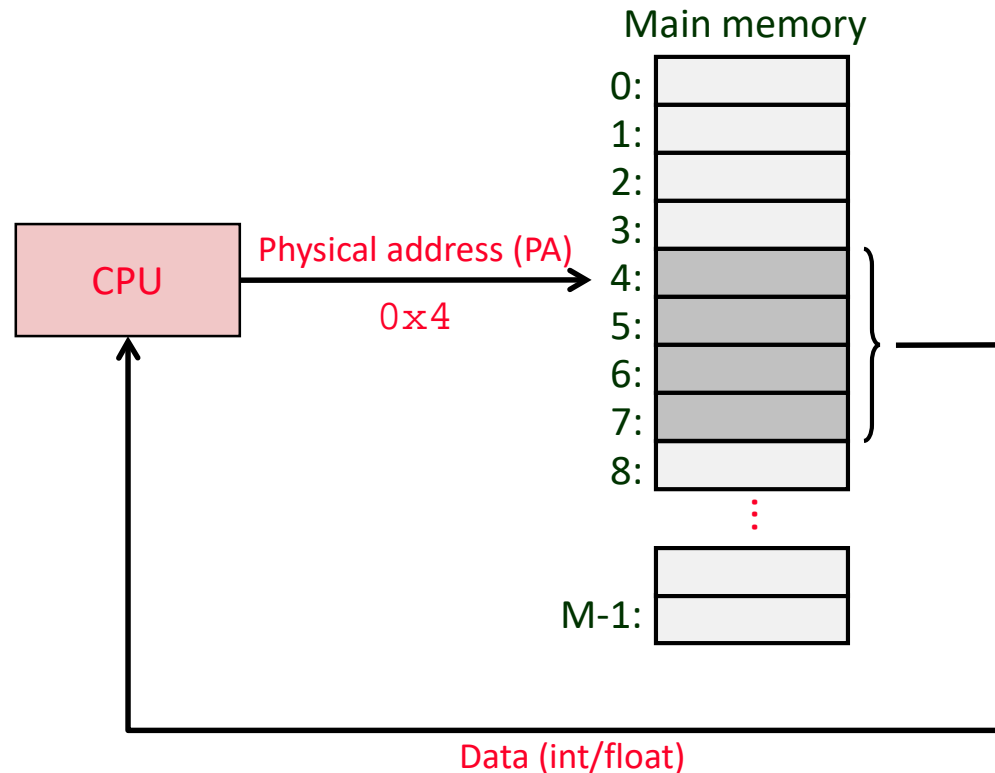
- ❑ **Large page size**: typically 4-8 KiB or 2-4 MiB
  - | Can be up to 1 GiB (for “Big Data” apps on big computers)
  - | Compared with 64-byte cache blocks
- ❑ **Fully associative**
  - | Any virtual page can be placed in any physical page
  - | Requires a “large” mapping function – different from CPU caches
- ❑ **Highly sophisticated, expensive replacement algorithms** in OS
  - | Too complicated and open-ended to be implemented in hardware
- ❑ ***Write-back*** rather than *write-through*
  - | *Really* don't want to write to disk every time we modify something in memory
  - | Some things may never end up on disk (e.g. stack for short-lived process)

# A System Using Virtual Addressing



- ❑ Physical addresses are *completely invisible to programs*
  - | Used in all modern desktops, laptops, servers, smartphones...
  - | One of the great ideas in computer science

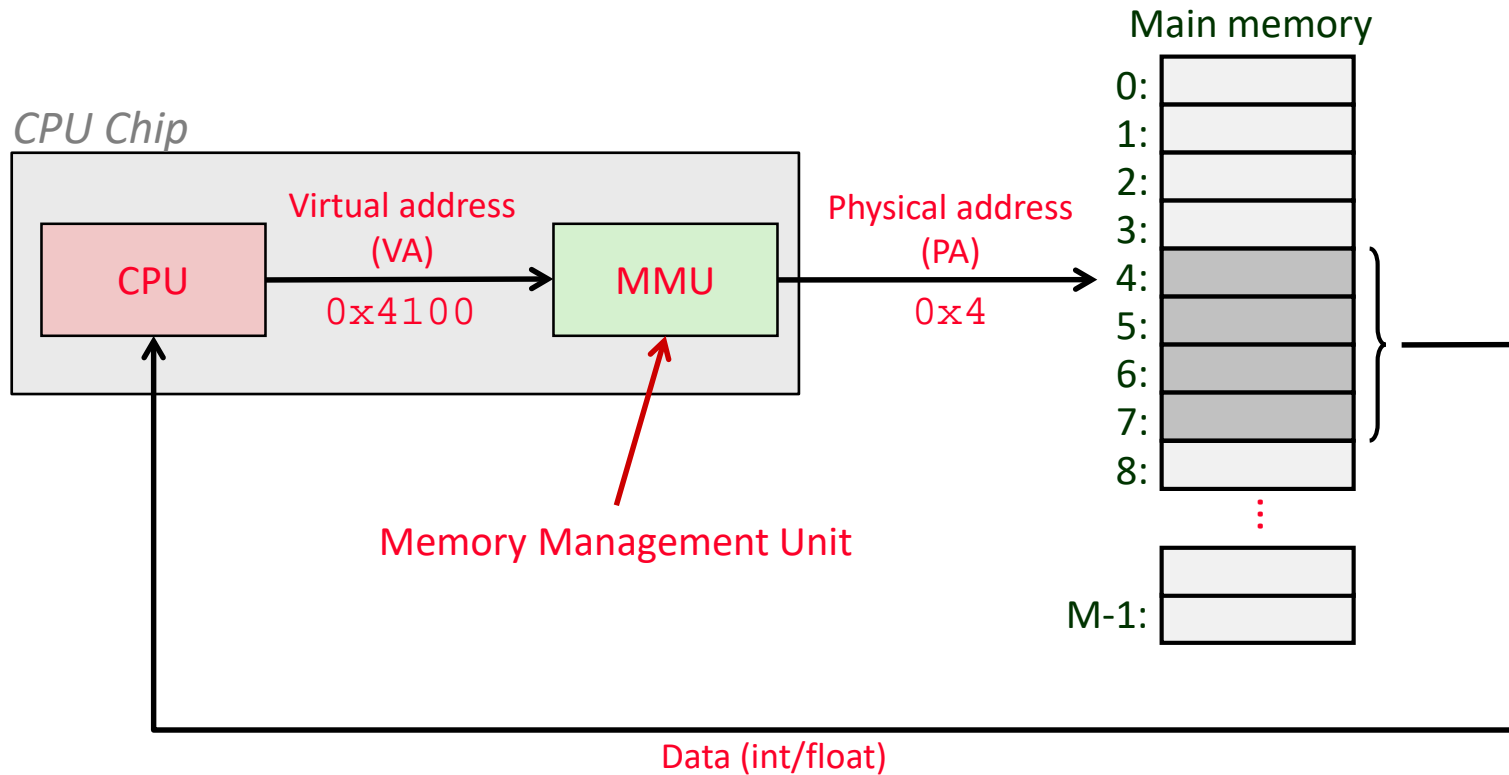
# A System Using Physical Addressing



- ❑ Used in “simple” systems with (usually) just one process:
  - | Embedded microcontrollers in devices like cars, elevators, and digital picture frames

# Address Translation

*How do we perform the virtual → physical address translation?  
Fig. 5. 25 and Fig. 5. 26*



# Address Translation: Page Tables

---

- ❑ CPU-generated address can be split into:

*n*-bit address:

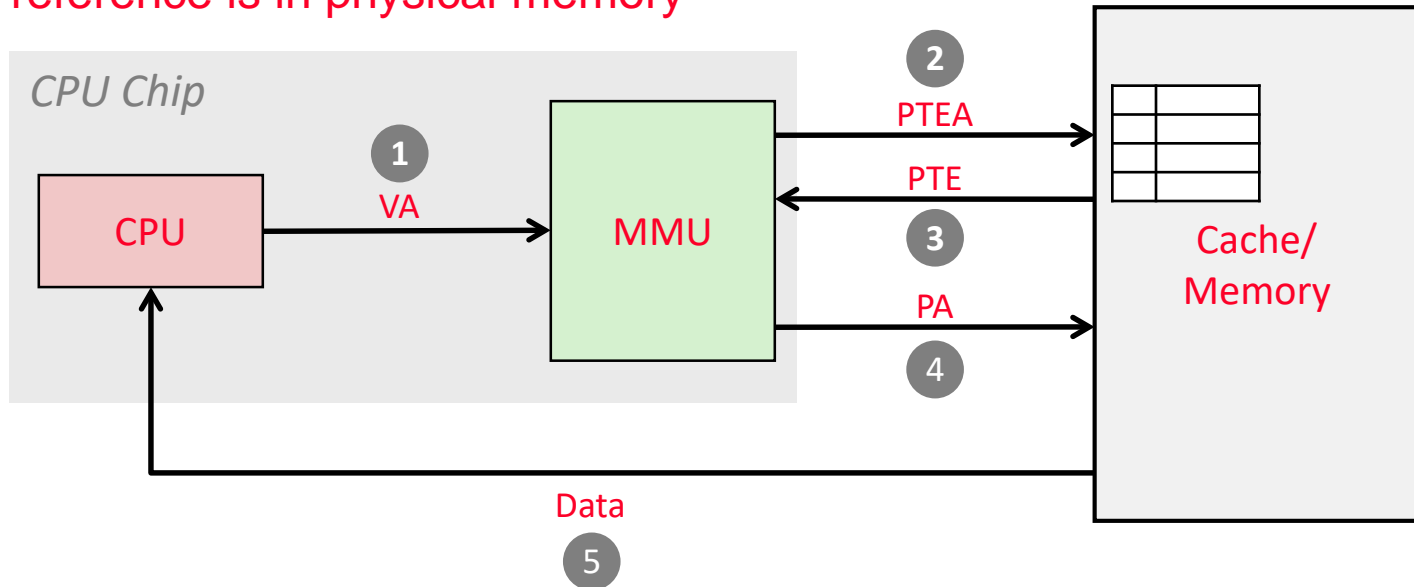
Virtual Page Number
---------------------

Page Offset
-------------

- | Request is Virtual Address (**VA**), want Physical Address (**PA**)
  - | Note that Physical Offset = Virtual Offset (page-aligned)
- 
- ❑ Use lookup table that we call the *page table* (**PT**)
- | Reading : pp. 432
  - | Fig. 5. 27

# Address Translation: Page Hit

VM reference is in physical memory



- 1) Processor sends *virtual* address to MMU (*memory management unit*)
- 2-3) MMU fetches PTE from page table in cache/memory  
(Uses PTBR to find beginning of page table for current process)
- 4) MMU sends *physical* address to cache/memory requesting data
- 5) Cache/memory sends data to processor

VA = Virtual Address

PTEA = Page Table Entry Address

PTE = Page Table Entry

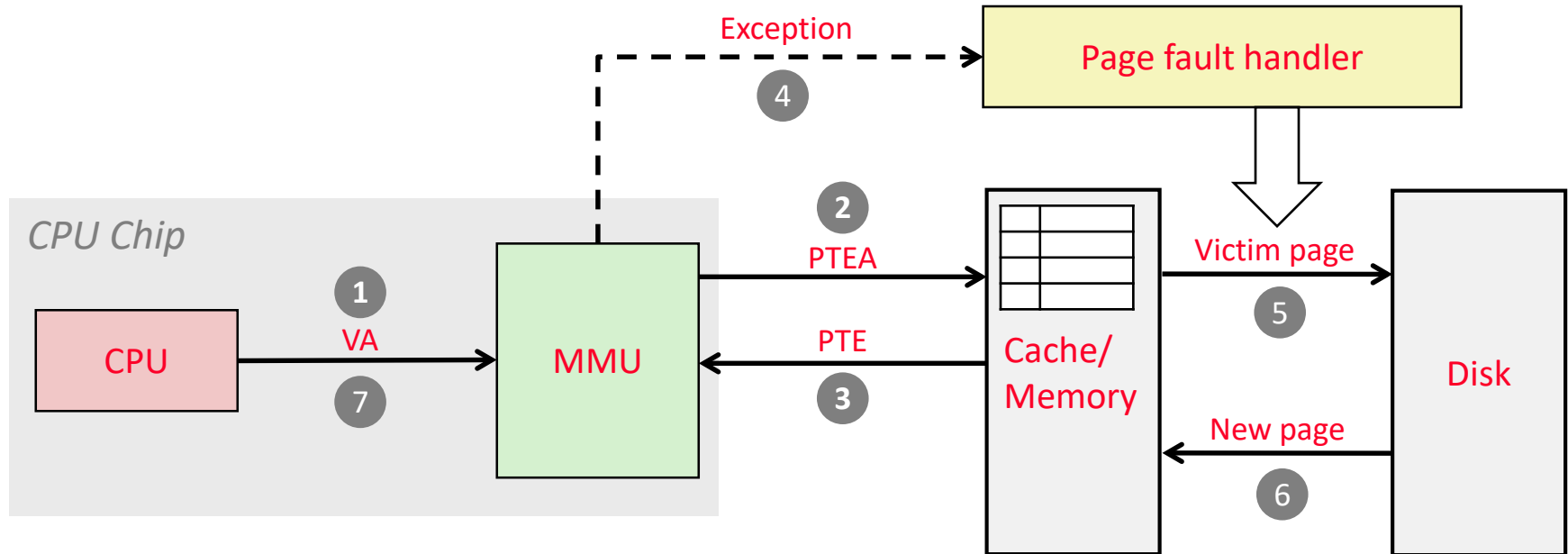
PA = Physical Address

Data = Contents of memory stored at VA originally requested by CPU



# Address Translation: Page Fault

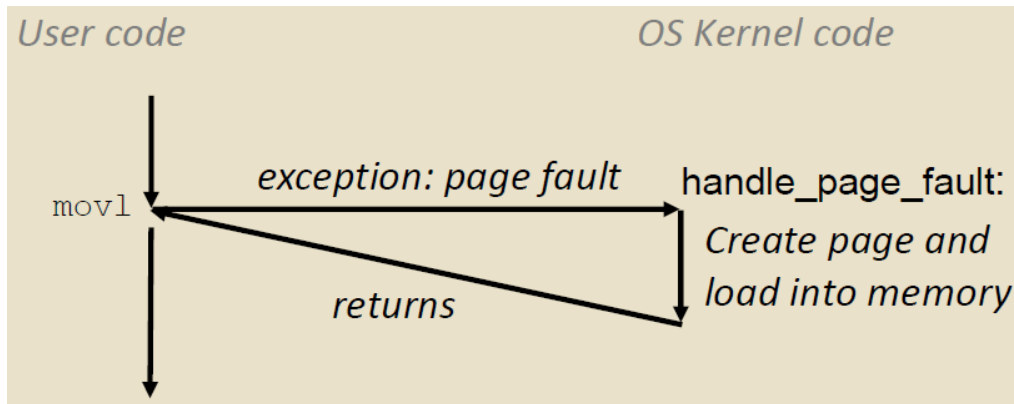
VM reference is NOT in physical memory



- 1) Processor sends virtual address to MMU
- 2-3) MMU fetches PTE from page table in cache/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

# Address Translation: Page Fault Penalty

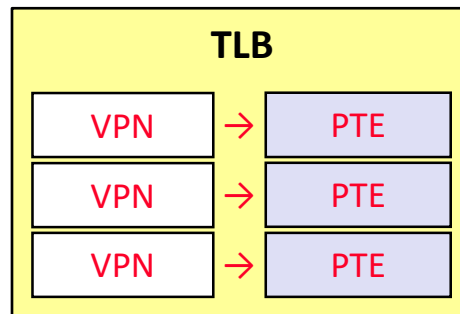
- ❑ On page fault, the page must be fetched from disk
  - | Takes millions of clock cycles
  - | Handled by OS code



- ❑ Try to minimize page fault rate
  - | Fully associative placement
  - | Smart replacement algorithms

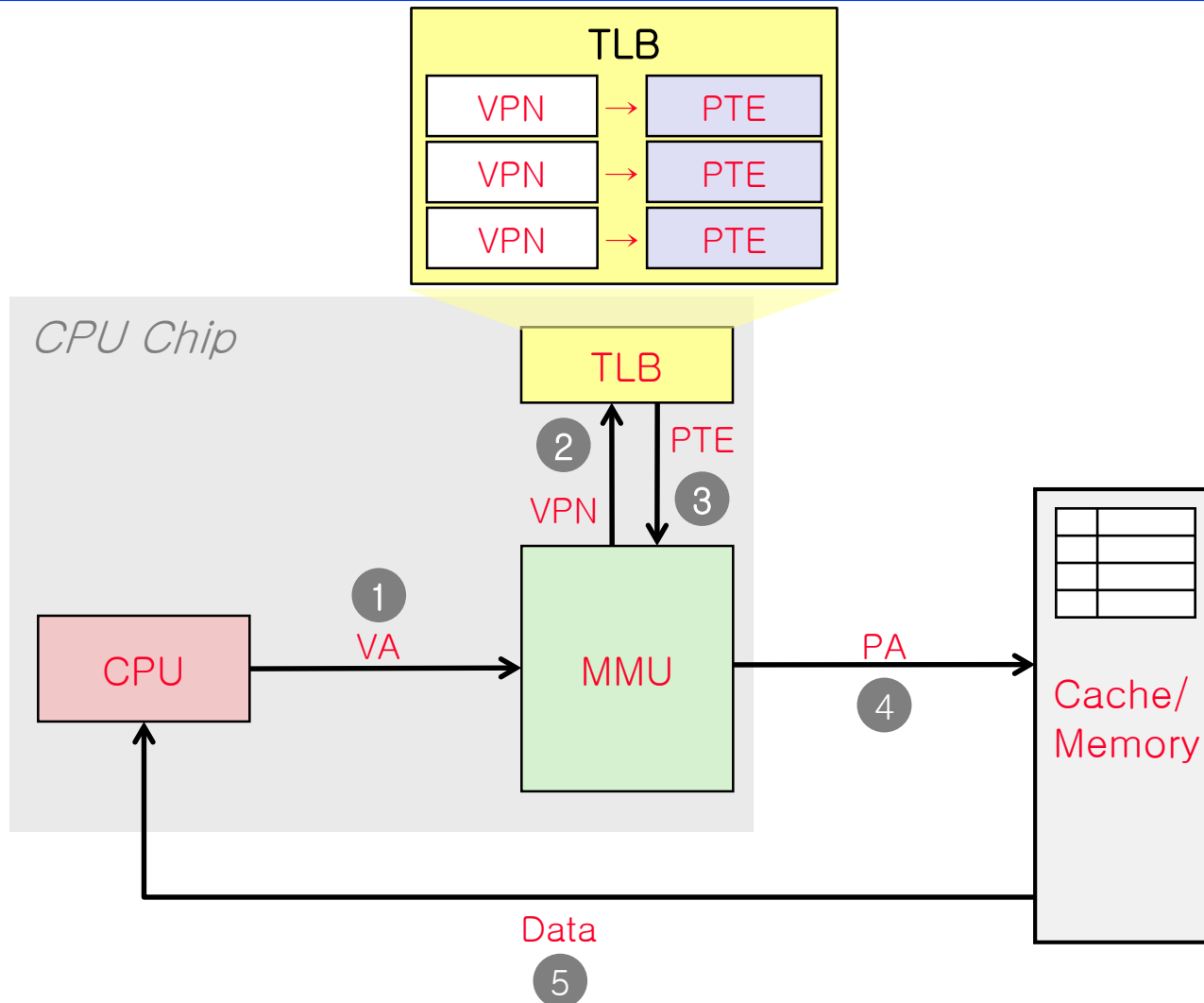
# Translation Sounds Slow

- ❑ The MMU accesses memory *twice*: once to get the PTE for translation, and then again for the actual memory request
  - | The PTEs *may* be cached in L1 like any other memory word
  - | Reading pp. 438
- ❑ Speeding up Translation with a TLB *Translation Lookaside Buffer*
  - | Small hardware **cache in MMU**
  - | Maps virtual page numbers to physical page numbers
  - | Contains complete *page table entries* for small number of pages
    - Modern Intel processors have 128 or 256 entries in TLB
  - | Much faster than a page table lookup in cache/memory



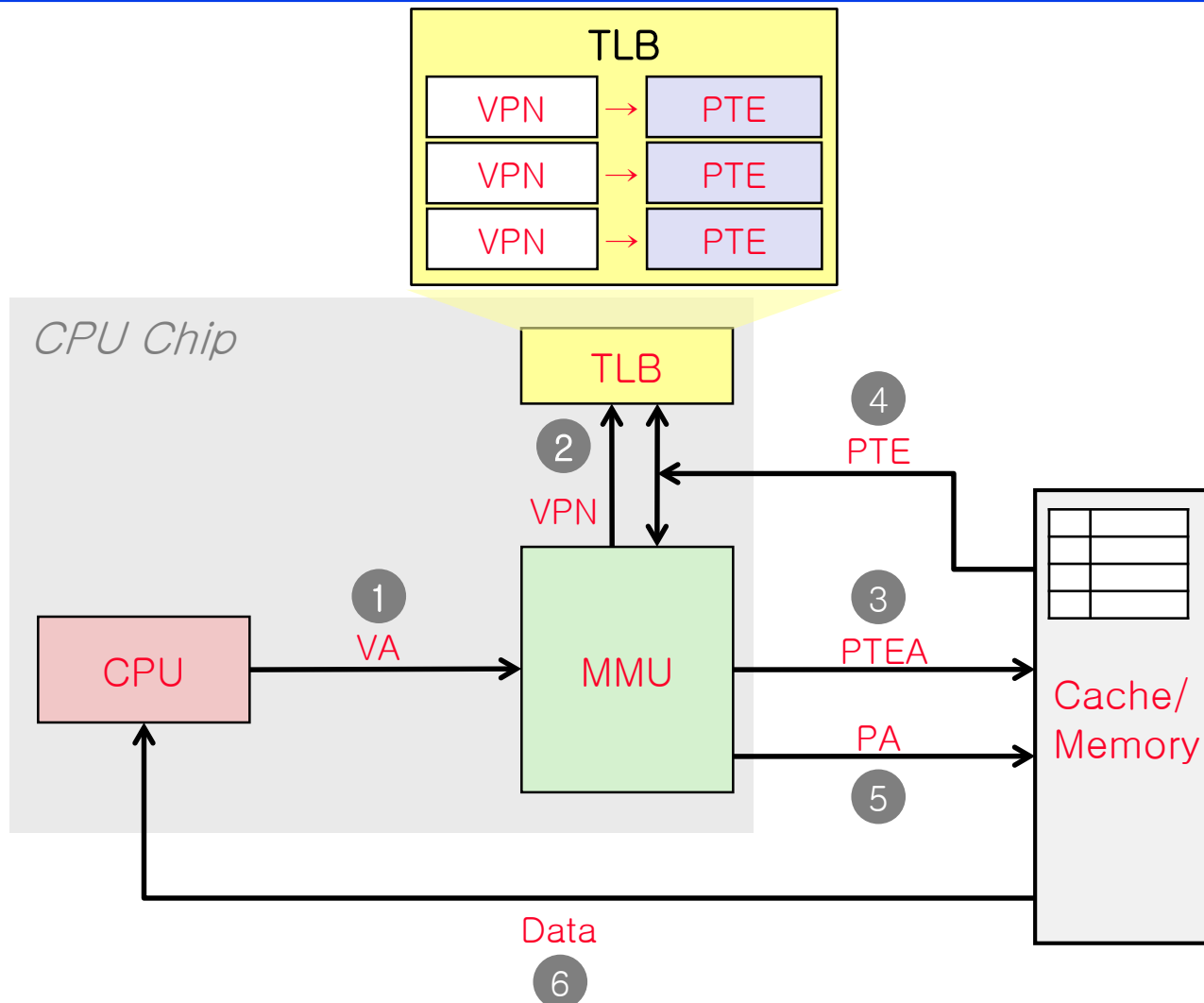
| Fig. 5.29

# TLB Hit



❑ A TLB hit eliminates a memory access!

# TLB Miss

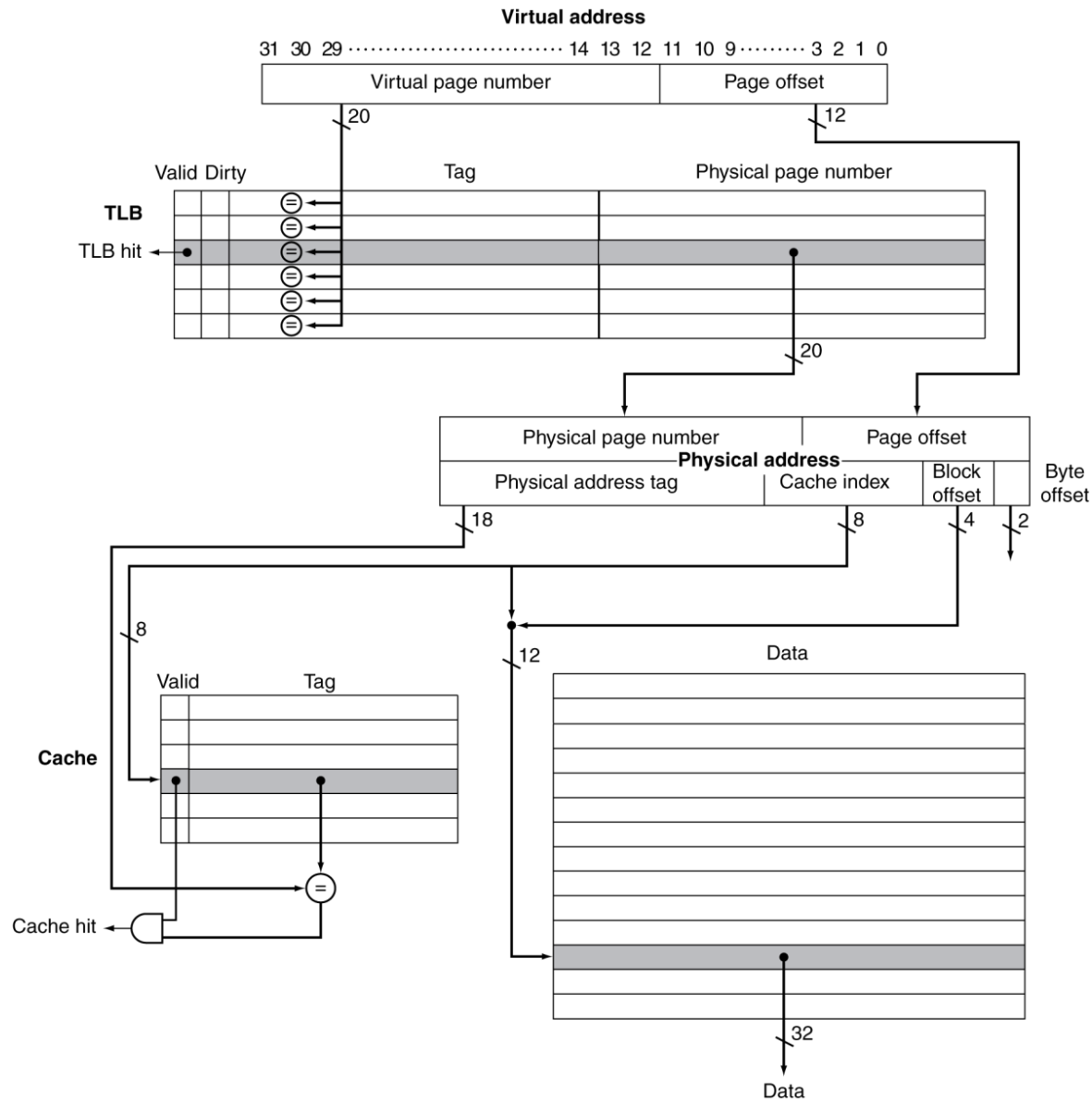


- ❑ A TLB miss incurs an additional memory access (the PTE)
  - | Fortunately, TLB misses are rare

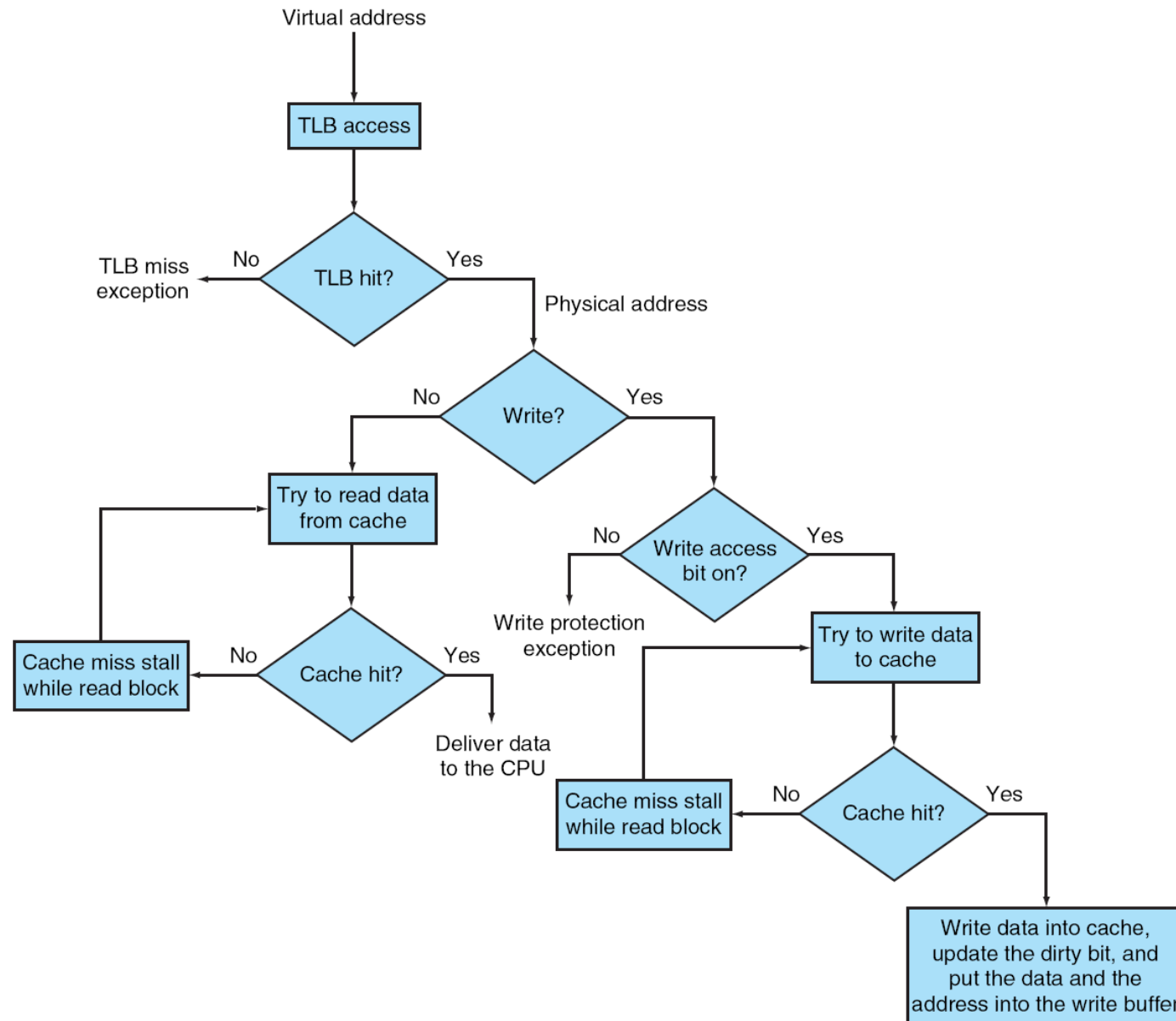
# Fast Translation Using a TLB

- ❑ Access to page tables has good locality
  - | So use a fast cache of PTEs within the CPU
  - | Called a **Translation Look-aside Buffer (TLB)**
  - | Typical:
    - 16–512 PTEs
    - 0.5–1 cycle for hit, 10–100 cycles for miss
    - 0.01%–1% miss rate
  - | Misses could be handled by hardware or software

# TLB and Cache Interaction (FastMATH, Fig. 5.30)



# Typical Memory Hierarchy - The Big Picture





# A Case of Memory System State

TLB:

Set	Tag	PPN	V	Tag	PPN	V	Tag	PPN	V	Tag	PPN	V
0	03	–	0	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

Page table (partial):

VPN	PPN	V	VPN	PPN	V
0	28	1	8	13	1
1	–	0	9	17	1
2	33	1	A	09	1
3	02	1	B	–	0
4	–	0	C	–	0
5	16	1	D	2D	1
6	–	0	E	–	0
7	–	0	F	0D	1

Cache:

Index	Tag	V	B0	B1	B2	B3
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

Index	Tag	V	B0	B1	B2	B3
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	–	–	–	–

# Summary

---

## ❑ Virtual Memory

- | Supports many OS-related functions
  - Process creation, task switching, protection
- | Operating System (software)
  - Allocates/shares physical memory among processes
  - Maintains high-level tables tracking memory type, source, sharing
  - Handles exceptions, fills in hardware-defined mapping tables
- | Hardware
  - Translates virtual addresses via mapping tables, enforcing permissions
  - Accelerates mapping via translation cache (TLB)

❑ Reading pp. 452 ~ 454

❑ Check Yourself pp. 454



# A Common Framework for The Memory Hierarchy

- ❑ Common principles apply at all levels of the memory hierarchy
  - | Based on notions of caching
  
- ❑ At each level in the hierarchy
  - | Block placement
  
  - | Finding a block
  
  - | Replacement on a miss
  
  - | Write policy

# Block Placement

---

- ❑ Determined by associativity
  - | Direct mapped (1-way associative)
    - One choice for placement
  - | n-way set associative
    - n choices within a set
  - | Fully associative
    - Any location
  
- ❑ Higher associativity reduces miss rate
  - | Increases complexity, cost, and access time

# Finding a Block

Associativity	Location method	Comparisons required
Direct mapped	Index	1
n-way set associative	Set index, then search entries within the set	n
Fully associative	Search all entries	#entries
	Full lookup table	0

## ❑ Hardware caches

- | Reduce comparisons to reduce cost

## ❑ Virtual memory

- | Full table lookup makes full associativity feasible
- | Benefit in reduced miss rate

# Replacement

---

- ❑ Choice of entry to replace on a miss
  - | Least recently used (LRU)
    - Complex and costly hardware for high associativity
  - | Random
    - Close to LRU, easier to implement
  
- ❑ Virtual memory
  - | LRU approximation with hardware support

# Write Policy

---

## ❑ Write-through

- | Update both upper and lower levels
- | Simplifies replacement, but may require write buffer

## ❑ Write-back

- | Update upper level only
- | Update lower level when block is replaced
- | Need to keep more state

## ❑ Virtual memory

- | Only write-back is feasible, given disk write latency



# Sources of Misses

---

## ❑ Compulsory misses (aka cold start misses)

- | First access to a block

## ❑ Capacity misses

- | Due to finite cache size
- | A replaced block is later accessed again

## ❑ Conflict misses (aka collision misses)

- | In a non-fully associative cache
- | Due to competition for entries in a set
- | Would not occur in a fully associative cache of the same total size

# Cache Design Trade-offs

---

Design change	Effect on miss rate	Negative performance effect
Increase cache size	Decrease capacity misses	May increase access time
Increase associativity	Decrease conflict misses	May increase access time
Increase block size	Decrease compulsory misses	Increases miss penalty. For very large block size, may increase miss rate due to pollution.

# Cache Coherence Protocols

---

- ❑ Operations performed by **caches in multiprocessors** to ensure coherence
  - | Migration of data to local caches
    - Reduces bandwidth for shared memory
  - | Replication of read-shared data
    - Reduces contention for access
  
- ❑ Snooping protocols
  - | Each cache monitors bus reads/writes
  
- ❑ Directory-based protocols
  - | Caches and memory record sharing status of blocks in a directory

# Invalidating Snooping Protocols

- ❑ Cache gets exclusive access to a block when it is to be written
  - | Broadcasts an invalidate message on the bus
  - | Subsequent read in another cache misses
    - Owning cache supplies updated value

❑ Fig. 5. 42

CPU activity	Bus activity	CPU A's cache	CPU B's cache	Memory
				0
CPU A reads X	Cache miss for X	0		0
CPU B reads X	Cache miss for X	0	0	0
CPU A writes 1 to X	Invalidate for X	1		0
CPU B read X	Cache miss for X	1	1	1

## 2-Level TLB Organization

Characteristic	ARM Cortex-A8	Intel Core i7
Virtual address	32 bits	48 bits
Physical address	32 bits	44 bits
Page size	Variable: 4, 16, 64 KiB, 1, 16 MiB	Variable: 4 KiB, 2/4 MiB
TLB organization	<p>1 TLB for instructions and 1 TLB for data</p> <p>Both TLBs are fully associative, with 32 entries, round robin replacement</p> <p>TLB misses handled in hardware</p>	<p>1 TLB for instructions and 1 TLB for data per core</p> <p>Both L1 TLBs are four-way set associative, LRU replacement</p> <p>L1 I-TLB has 128 entries for small pages, 7 per thread for large pages</p> <p>L1 D-TLB has 64 entries for small pages, 32 for large pages</p> <p>The L2 TLB is four-way set associative, LRU replacement</p> <p>The L2 TLB has 512 entries</p> <p>TLB misses handled in hardware</p>

# Concluding Remarks

---

- ❑ Fast memories are small, large memories are slow
  - | We really want fast, large memories
  - | Caching gives this illusion
  
- ❑ Principle of locality
  - | Programs use a small part of their memory space frequently
  
- ❑ Memory hierarchy
  - | L1 cache ↔ L2 cache ↔ ... ↔ DRAM memory ↔ disk
  
- ❑ Memory system design is critical for multiprocessors