# CSC 252: Computer Organization Spring 2018: Lecture 22

Instructor: Yuhao Zhu

Department of Computer Science
University of Rochester

**Action Items:**
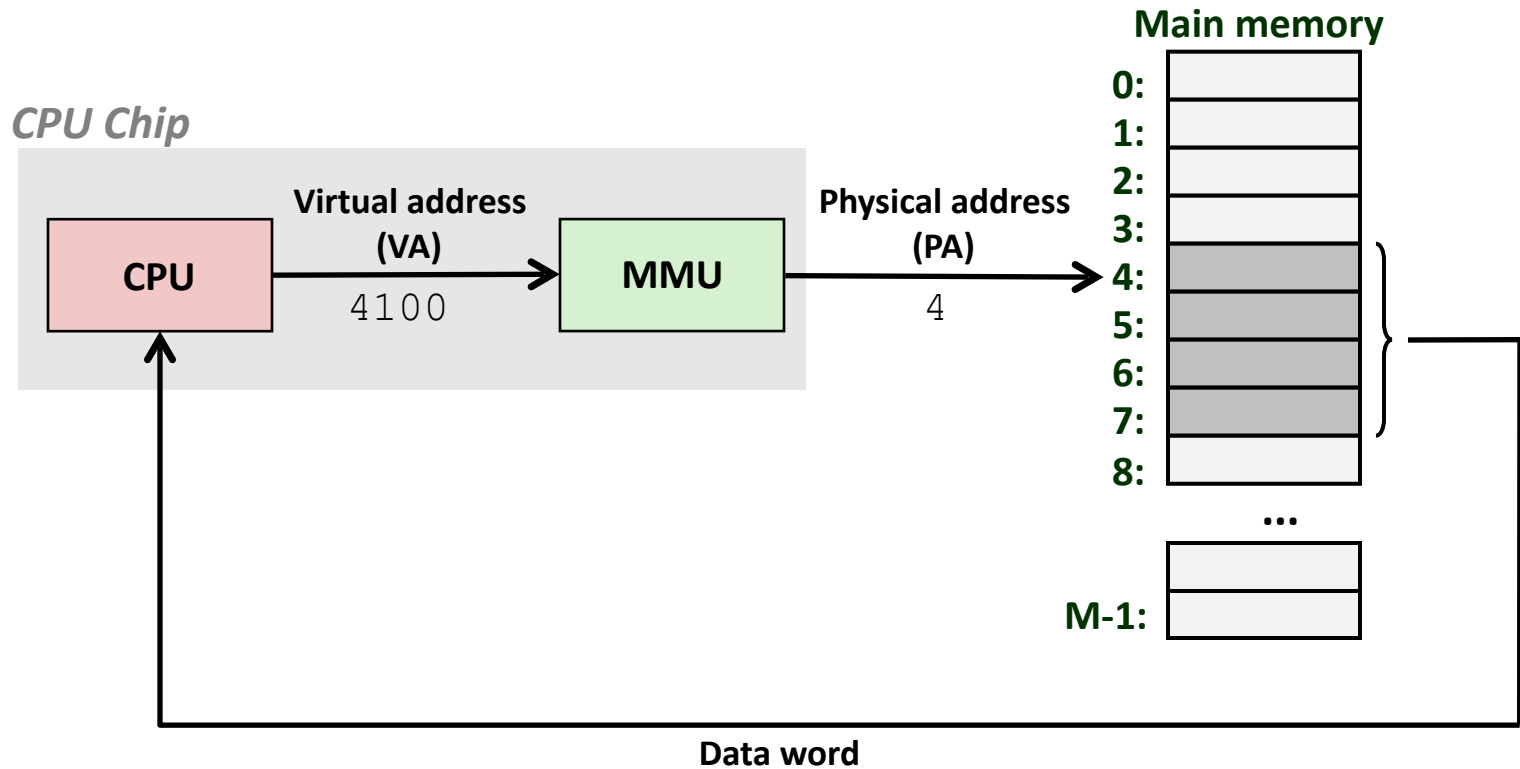- **Programming Assignment 5 is due Monday**

# Announcement

- Programming Assignment 5 is out
  - Main assignment: 11:59pm, **Monday, April 16**.
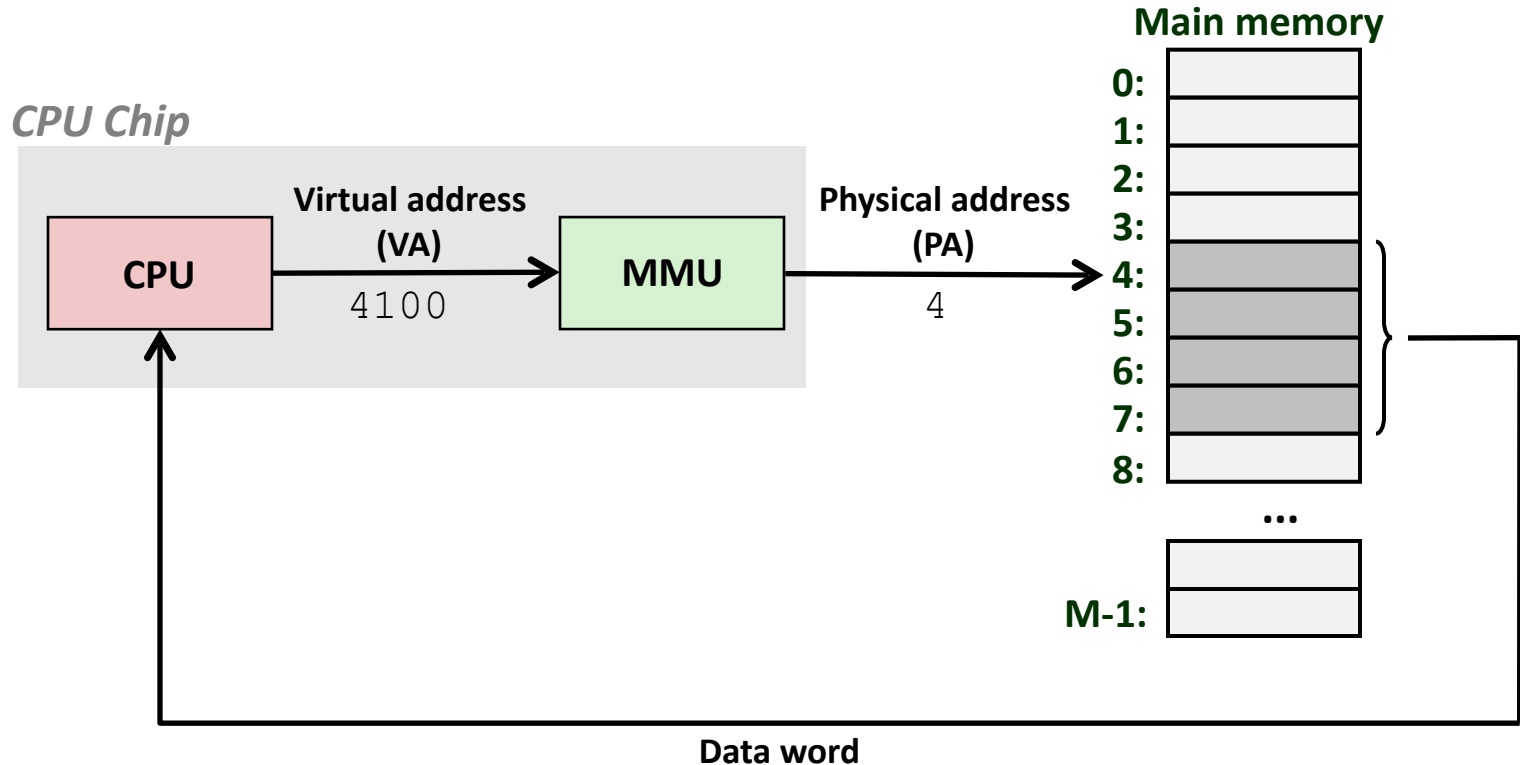- No office hours today. Had already moved to this Tuesday.

| 8 | 9 | 10 | 11 | 12 | 13 | 14 |
|---|---|----|----|----|----|----|
| 15 | 16 **Due** | 17 | 18 | 19 | 20 | 21 |

# A System Using Virtual Addressing

**Main memory**

*CPU Chip*

CPU → **Virtual address (VA)** `4100` → MMU → **Physical address (PA)** `4` → 

| | |
|---|---|
| **0:** | |
| **1:** | |
| **2:** | |
| **3:** | |
| **4:** | |
| **5:** | |
| **6:** | |
| **7:** | |
| **8:** | |
| **...** | |
| **M-1:** | |

**Data word**
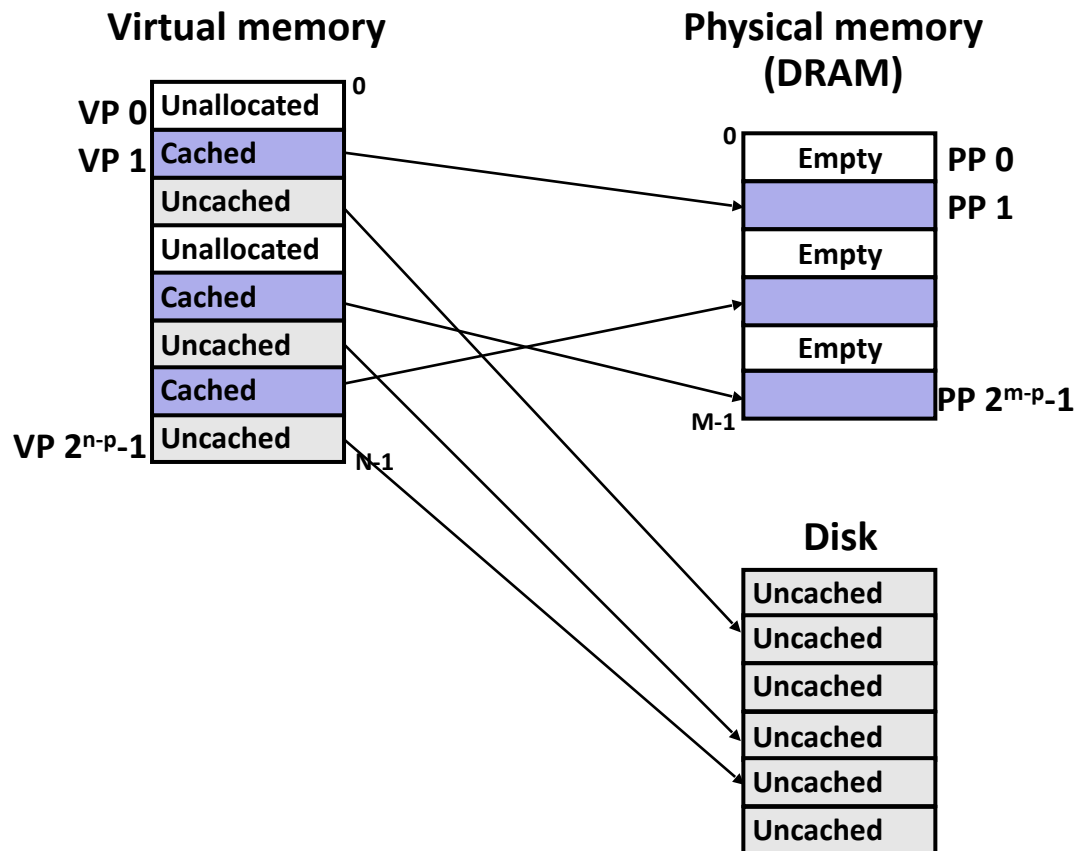
# A System Using Virtual Addressing



- On a 64-bit machine, virtual memory size = $2^{64}$
- Physical memory size is much much smaller:
  - iPhone 8: 2 GB ($2^{31}$)
  - 15-inch Macbook Pro 2017: 16 GB ($2^{34}$)

# VM Concepts

- Conceptually, *virtual memory* is an array of N *pages* stored on disk. The contents of the array on disk are cached in *physical memory*, which is an array of M pages (M $<<$ N).
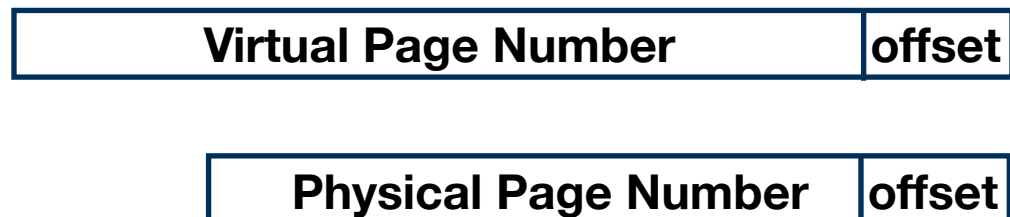
# VM Concepts

- Conceptually, *virtual memory* is an array of N *pages* stored on disk. The contents of the array on disk are cached in *physical memory*, which is an array of M pages (M << N).



**Virtual memory**

VP 0 | Unallocated | 0
VP 1 | Cached
| Uncached
| Unallocated
| Cached
| Uncached
| Cached
VP $2^{n-p}-1$ | Uncached | N-1

**Physical memory (DRAM)**

0
Empty | PP 0
| PP 1
Empty
Empty
| PP $2^{m-p}-1$
M-1

**Disk**

Uncached
Uncached
Uncached
Uncached
Uncached
Uncached

4

# VM Concepts

- Page size is the same in VM and PM

| Virtual Page Number | offset |
|:---:|:---:|

| Physical Page Number | offset |
|:---:|:---:|

# VM Concepts

- Page size is the same in VM and PM
- In a 64-bit machine, VA is 64-bit long. Assuming PM is 4 GB. Assuming 4KB page size.

| Virtual Page Number | offset |
|:---:|:---:|

| Physical Page Number | offset |
|:---:|:---:|

# VM Concepts

- Page size is the same in VM and PM
- In a 64-bit machine, VA is 64-bit long. Assuming PM is 4 GB. Assuming 4KB page size.
- How many bits for offset?
  - 12. Same for VM and PM

| Virtual Page Number | offset |
|---|---|

| Physical Page Number | offset |
|---|---|

# VM Concepts

- Page size is the same in VM and PM
- In a 64-bit machine, VA is 64-bit long. Assuming PM is 4 GB. Assuming 4KB page size.
- How many bits for offset?
  - 12. Same for VM and PM
- How many bits for Virtual Page Number?
  - 52

| Virtual Page Number | offset |
|---|---|

| Physical Page Number | offset |
|---|---|

# VM Concepts

- Page size is the same in VM and PM
- In a 64-bit machine, VA is 64-bit long. Assuming PM is 4 GB. Assuming 4KB page size.
- How many bits for offset?
  - 12. Same for VM and PM
- How many bits for Virtual Page Number?
  - 52
- How many bits for Physical Page Number?
  - 20

| Virtual Page Number | offset |
|---|---|

| Physical Page Number | offset |
|---|---|

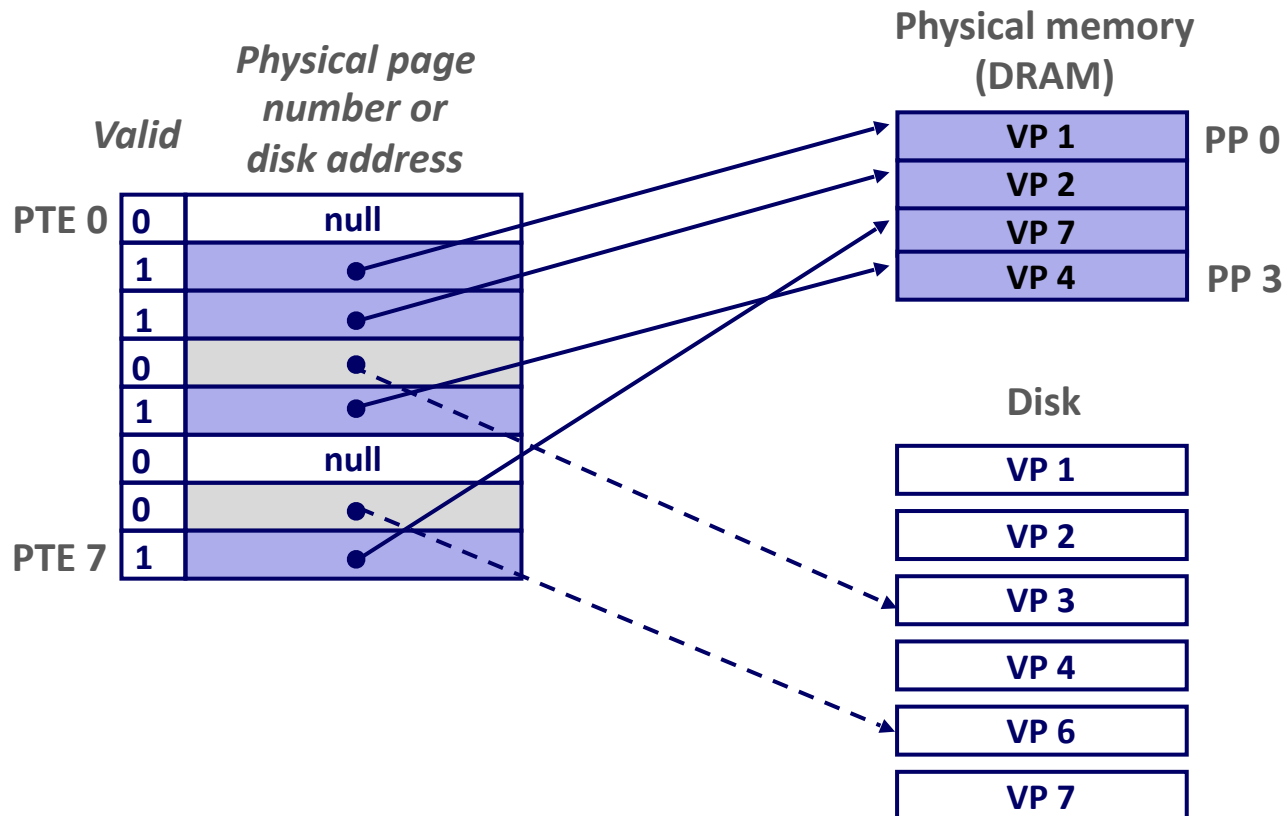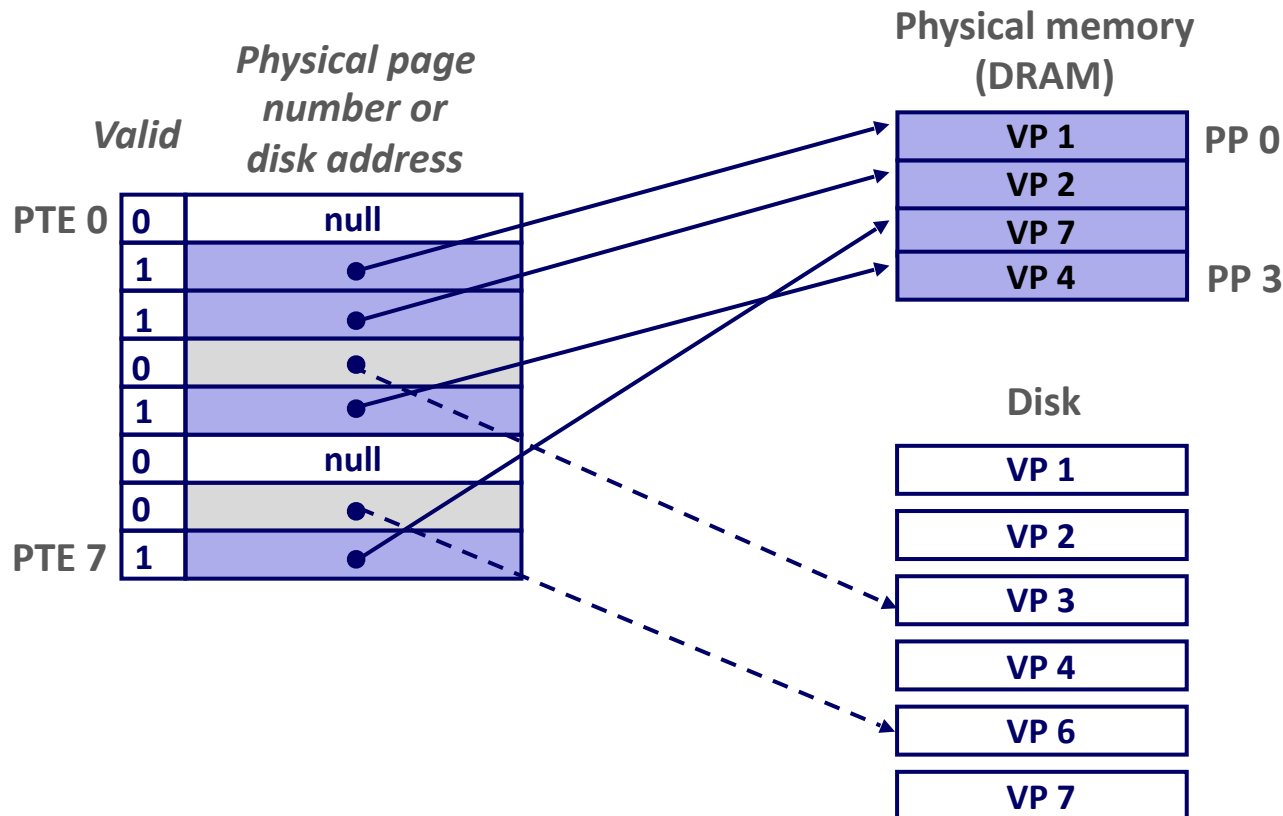# Page Table: Enabling VA to PA Translation

- A page table is an array of page table entries (PTEs) that maps every virtual page to its physical page.

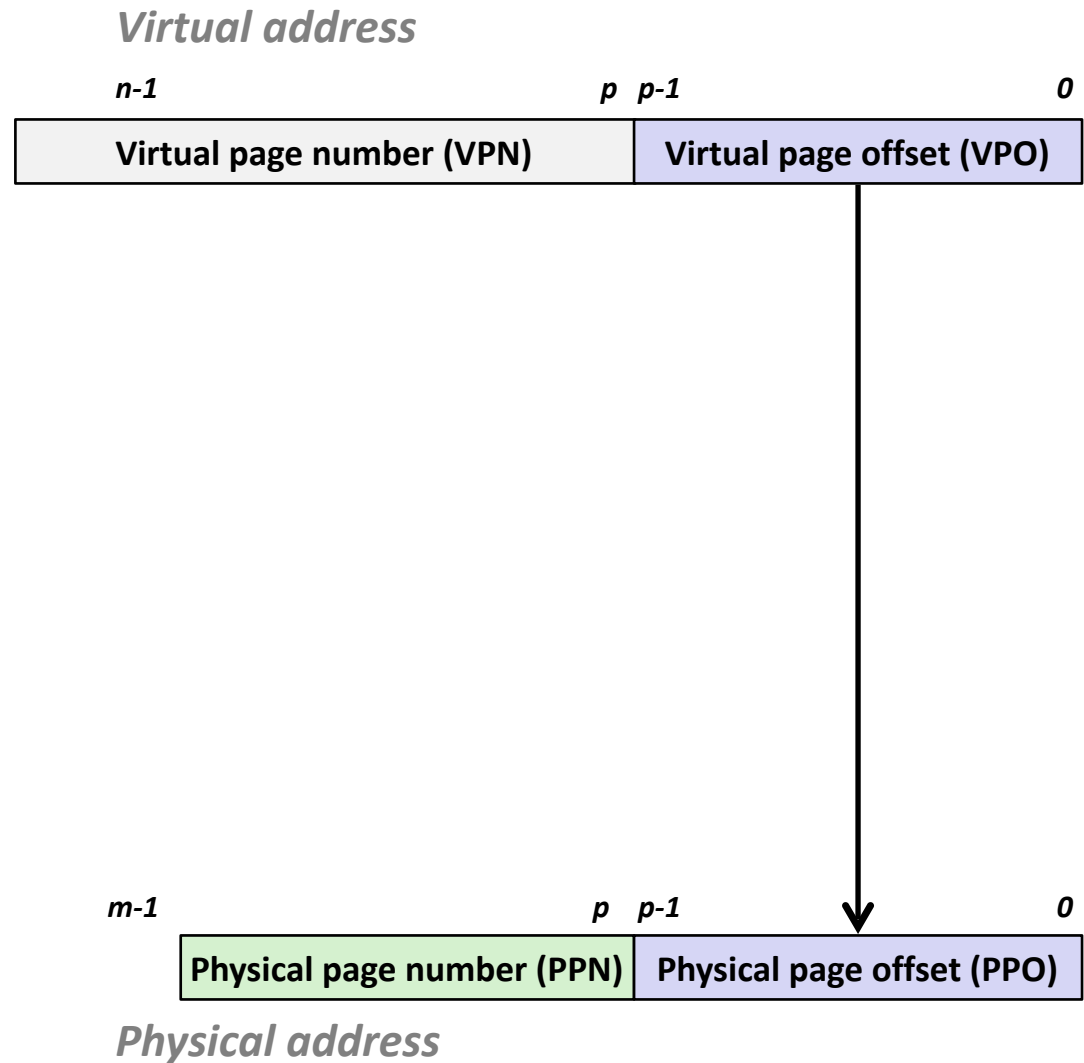| | Valid | Physical page number or disk address |
|---|---|---|
| PTE 0 | 0 | null |
| | 1 | ● |
| | 1 | ● |
| | 0 | ● |
| | 1 | ● |
| | 0 | null |
| | 0 | ● |
| PTE 7 | 1 | ● |

# Page Table: Enabling VA to PA Translation

- A page table is an array of page table entries (PTEs) that maps every virtual page to its physical page.

# Page Table: Enabling VA to PA Translation

- A page table is an array of page table entries (PTEs) that maps every virtual page to its physical page.
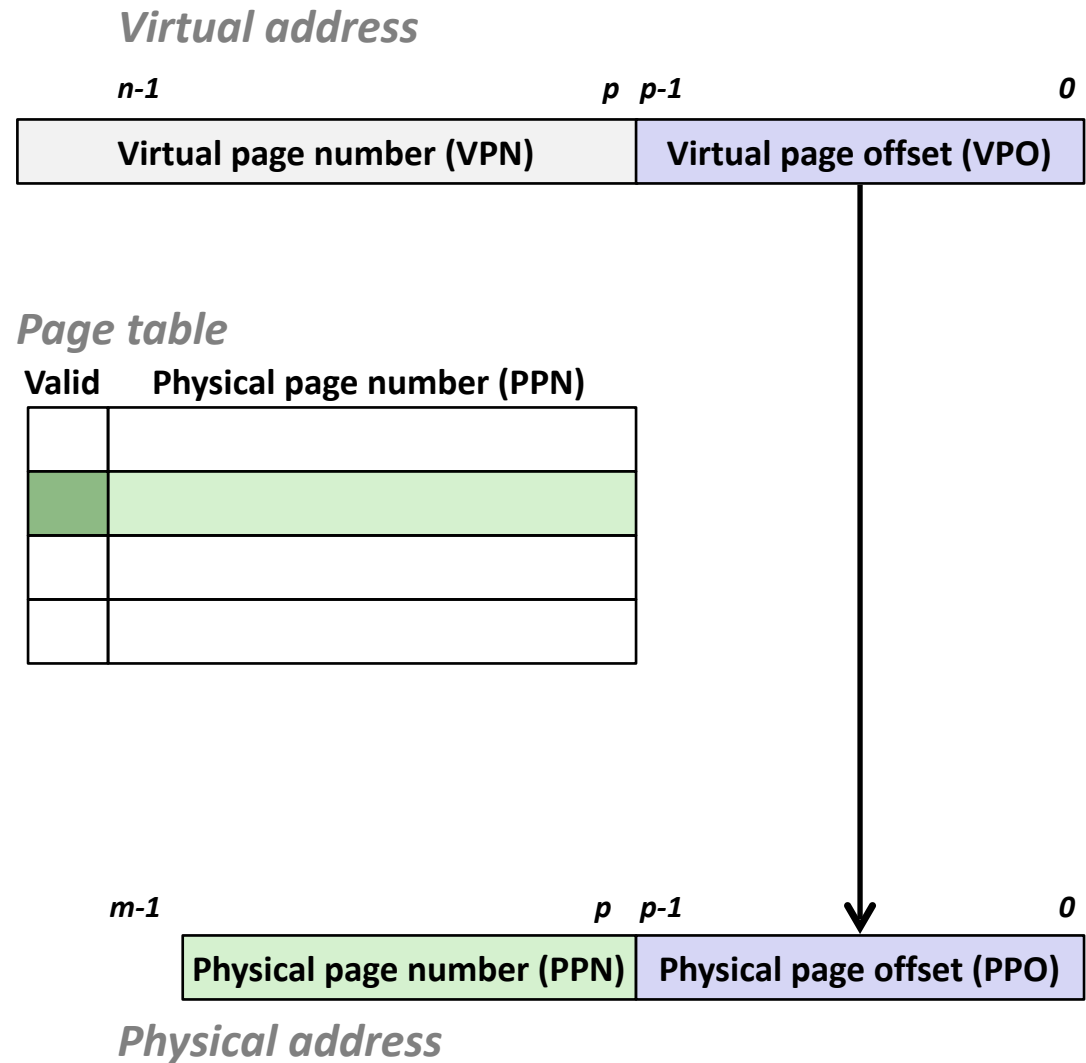
# Page Table: Enabling VA to PA Translation

- A page table is an array of page table entries (PTEs) that maps every virtual page to its physical page.

# Page Table: Enabling VA to PA Translation

- A page table is an array of page table entries (PTEs) that maps every virtual page to its physical page.

- 64-bit machine, 4KB page size, how many PTEs?
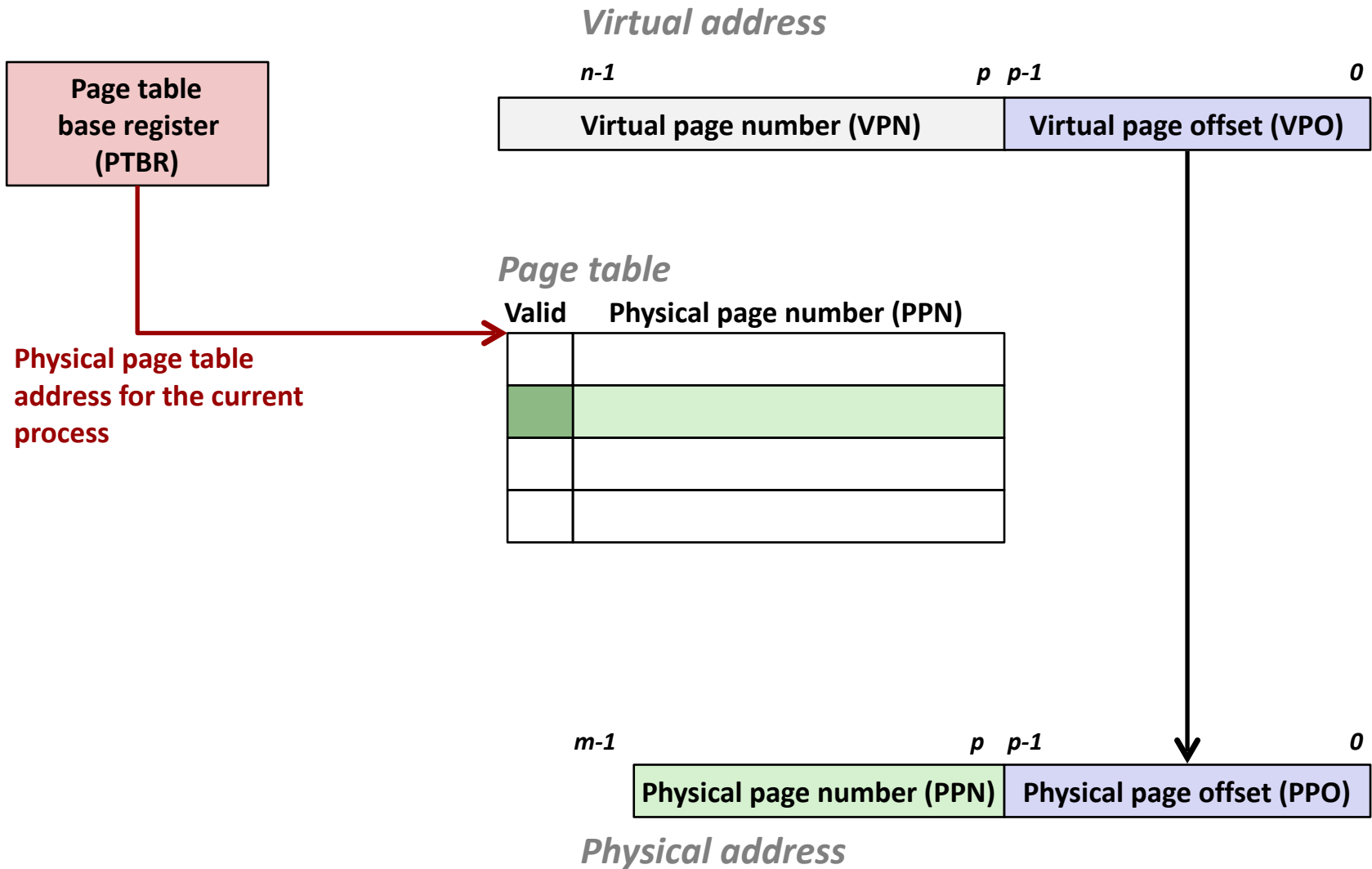  - Every page has a PTE, so $2^{52}$ PTEs.
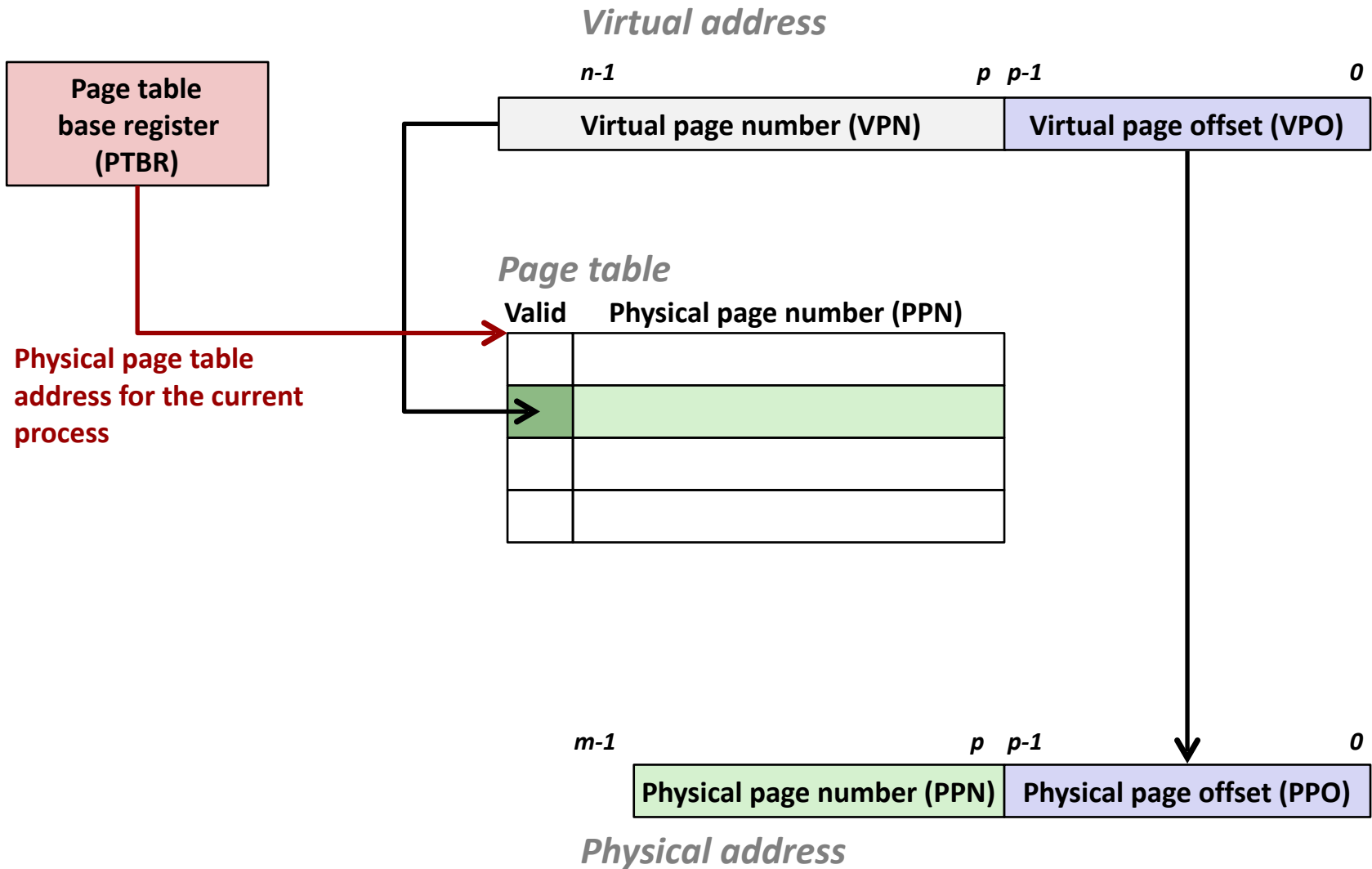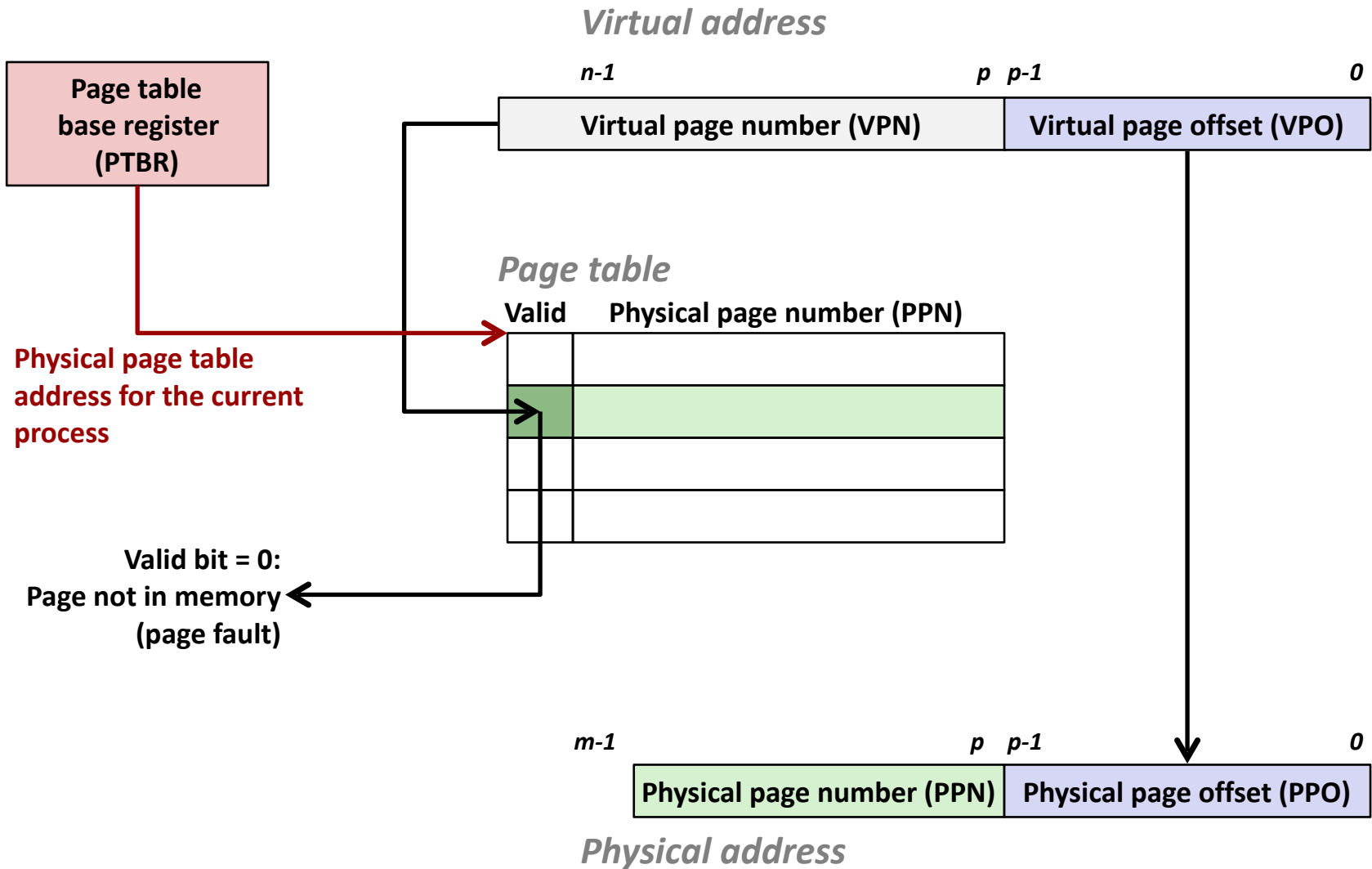
# Address Translation With a Page Table

*Virtual address*

| | | |
|---|---|---|
| *n-1* | *p  p-1* | *0* |
| Virtual page number (VPN) | Virtual page offset (VPO) | |

| | | |
|---|---|---|
| *m-1* | *p   p-1* | *0* |
| Physical page number (PPN) | Physical page offset (PPO) | |

*Physical address*

# Address Translation With a Page Table

*Virtual address*

| *n-1* | *p  p-1* | *0* |
|---|---|---|
| Virtual page number (VPN) | Virtual page offset (VPO) | |

*Page table*

| Valid | Physical page number (PPN) |
|---|---|
| | |
| | |
| | |
| | |

*Physical address*

| *m-1* | *p  p-1* | *0* |
|---|---|---|
| Physical page number (PPN) | Physical page offset (PPO) | |

# Address Translation With a Page Table

*Virtual address*

| n-1 | p | p-1 | 0 |
|---|---|---|---|
| **Virtual page number (VPN)** | | **Virtual page offset (VPO)** | |

**Page table
base register
(PTBR)**

Physical page table
address for the current
process

*Page table*

**Valid    Physical page number (PPN)**

*Physical address*

| m-1 | p | p-1 | 0 |
|---|---|---|---|
| **Physical page number (PPN)** | | **Physical page offset (PPO)** | |

# Address Translation With a Page Table

*Virtual address*

| Page table base register (PTBR) | | n-1 | Virtual page number (VPN) | p | p-1 | Virtual page offset (VPO) | 0 |

**Physical page table address for the current process**

*Page table*

| Valid | Physical page number (PPN) |

*Physical address*

| m-1 | Physical page number (PPN) | p | p-1 | Physical page offset (PPO) | 0 |

# Address Translation With a Page Table

*Virtual address*

| *n-1* | *p* *p-1* | *0* |
|---|---|---|
| **Virtual page number (VPN)** | **Virtual page offset (VPO)** | |

**Page table base register (PTBR)**

**Physical page table address for the current process**

*Page table*

**Valid**   **Physical page number (PPN)**

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

*Physical address*

| *m-1* | *p* *p-1* | *0* |
|---|---|---|
| **Physical page number (PPN)** | **Physical page offset (PPO)** | |

# Address Translation With a Page Table

*Virtual address*

| Page table base register (PTBR) | | *n-1*                         *p  p-1*                        *0* |
|---|---|---|
| | | | Virtual page number (VPN) | Virtual page offset (VPO) |

**Physical page table address for the current process**

*Page table*

**Valid**  **Physical page number (PPN)**

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

**Valid bit = 1**

*m-1*                         *p  p-1*                        *0*

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

*Physical address*

# Address Translation: Page Hit

# Address Translation: Page Hit



1) Processor sends virtual address to MMU

# Address Translation: Page Hit



1) Processor sends virtual address to MMU

# Address Translation: Page Hit



1) Processor sends virtual address to MMU

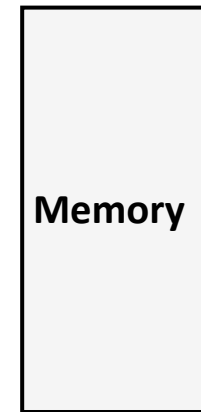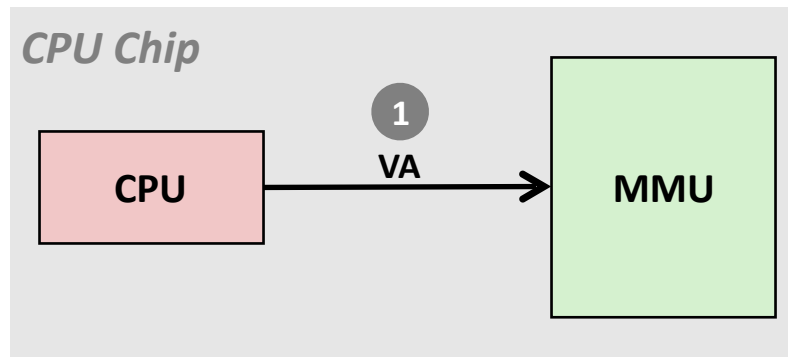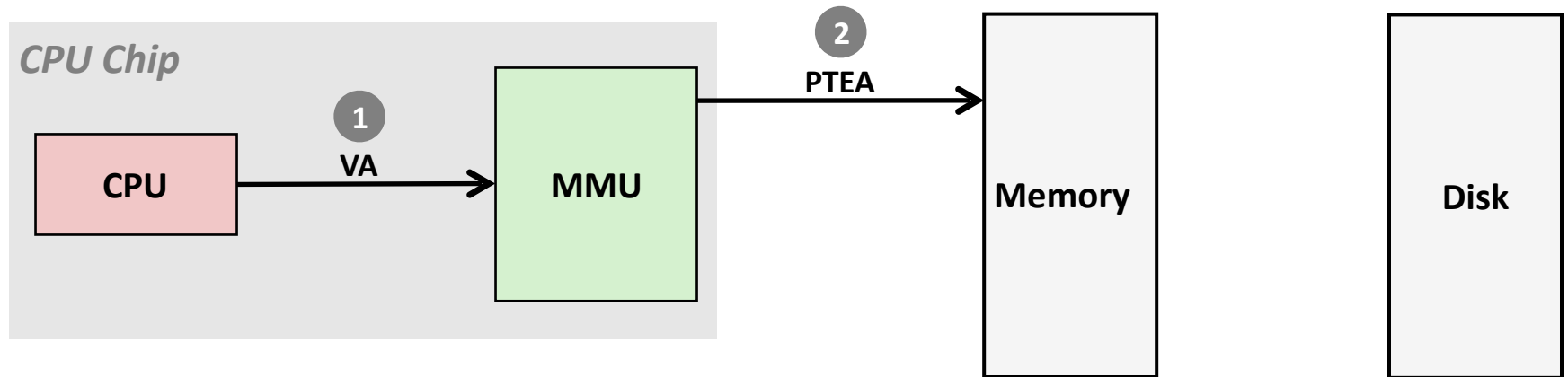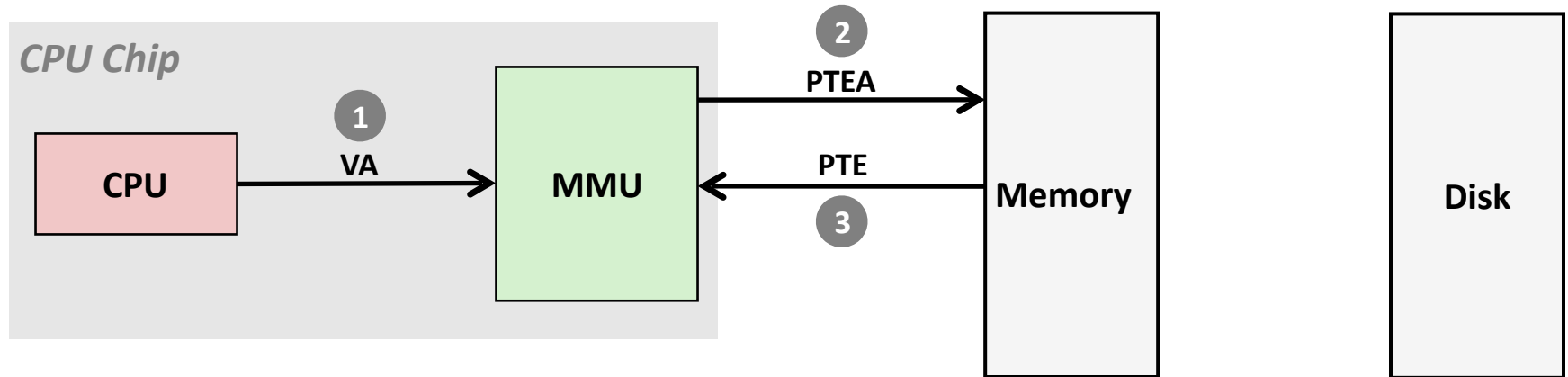2-3) MMU fetches PTE from page table in memory

# 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

# 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

# Address Translation: Page Fault



1) Processor sends virtual address to MMU

# Address Translation: Page Fault
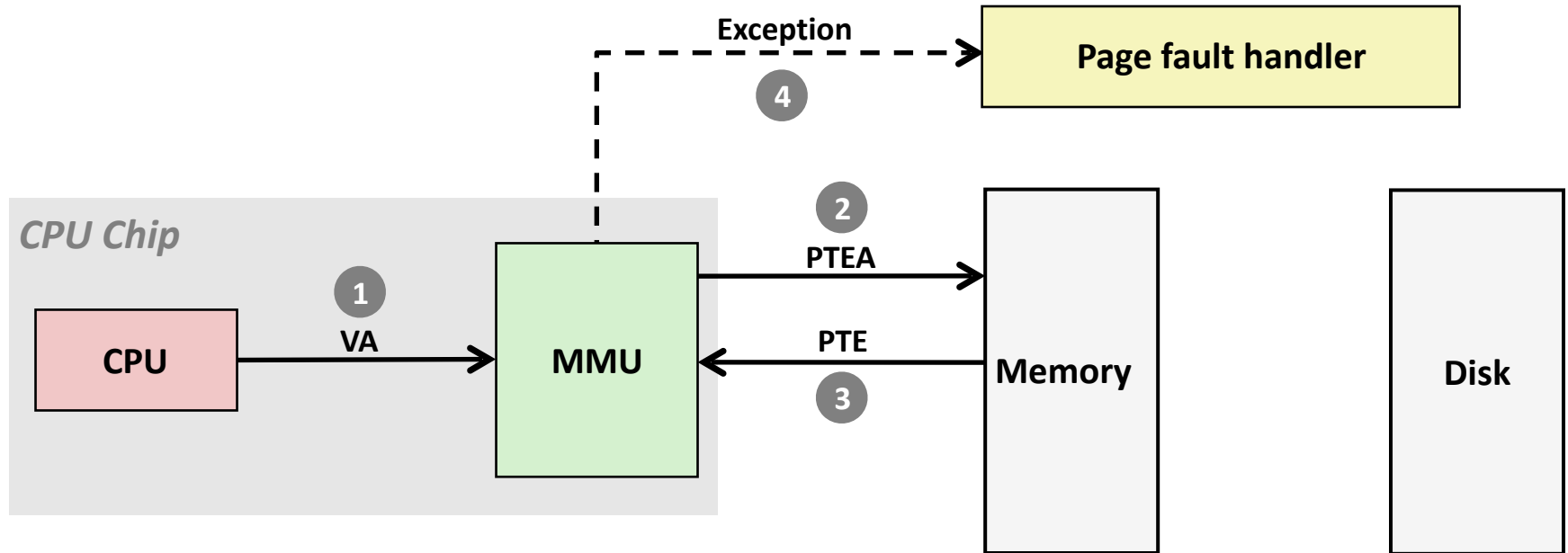


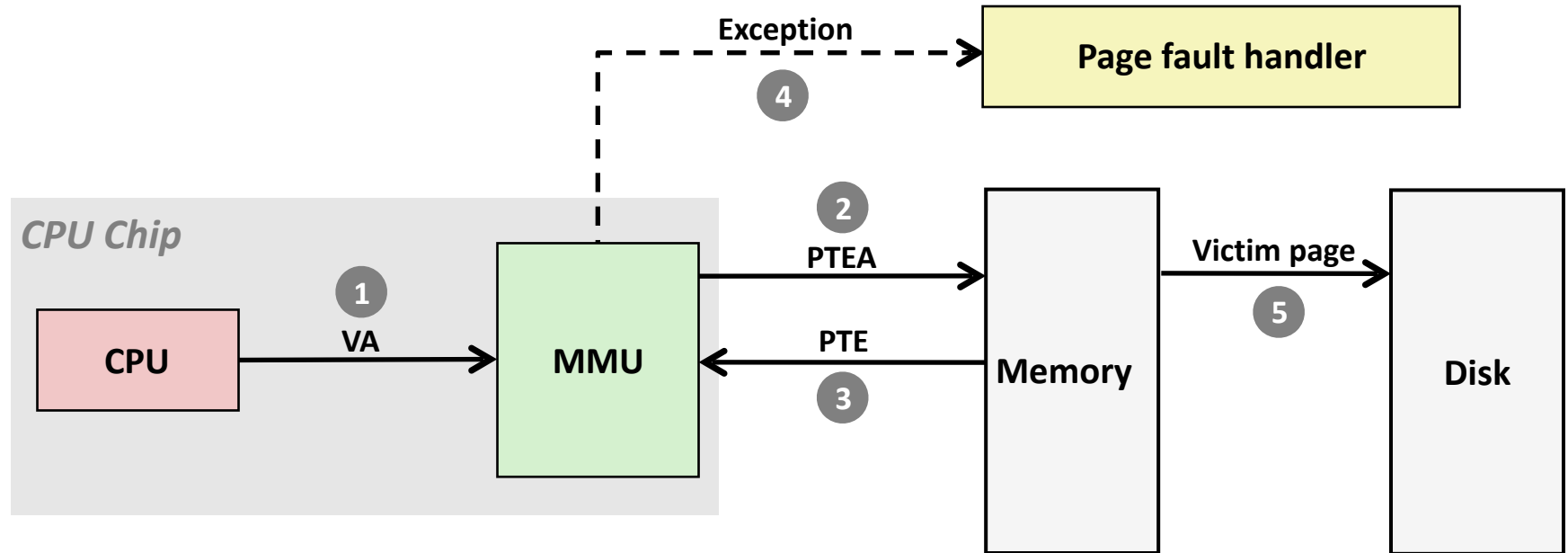1) Processor sends virtual address to MMU

# Address Translation: Page Fault



1) Processor sends virtual address to MMU

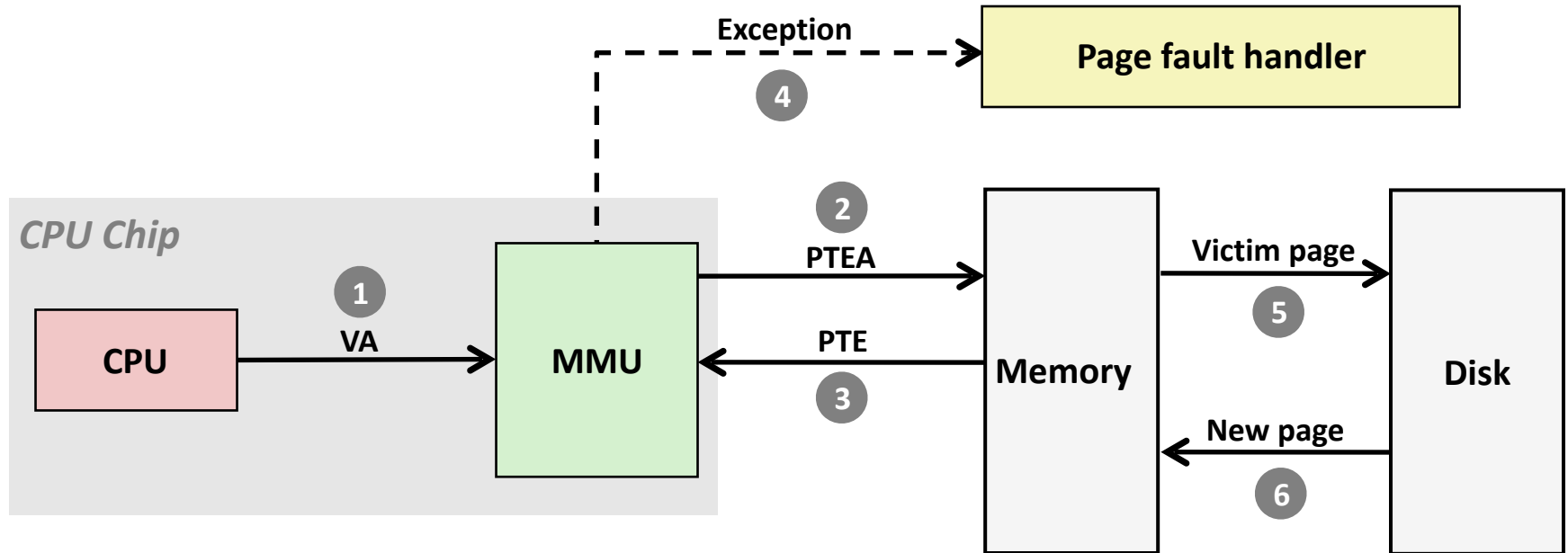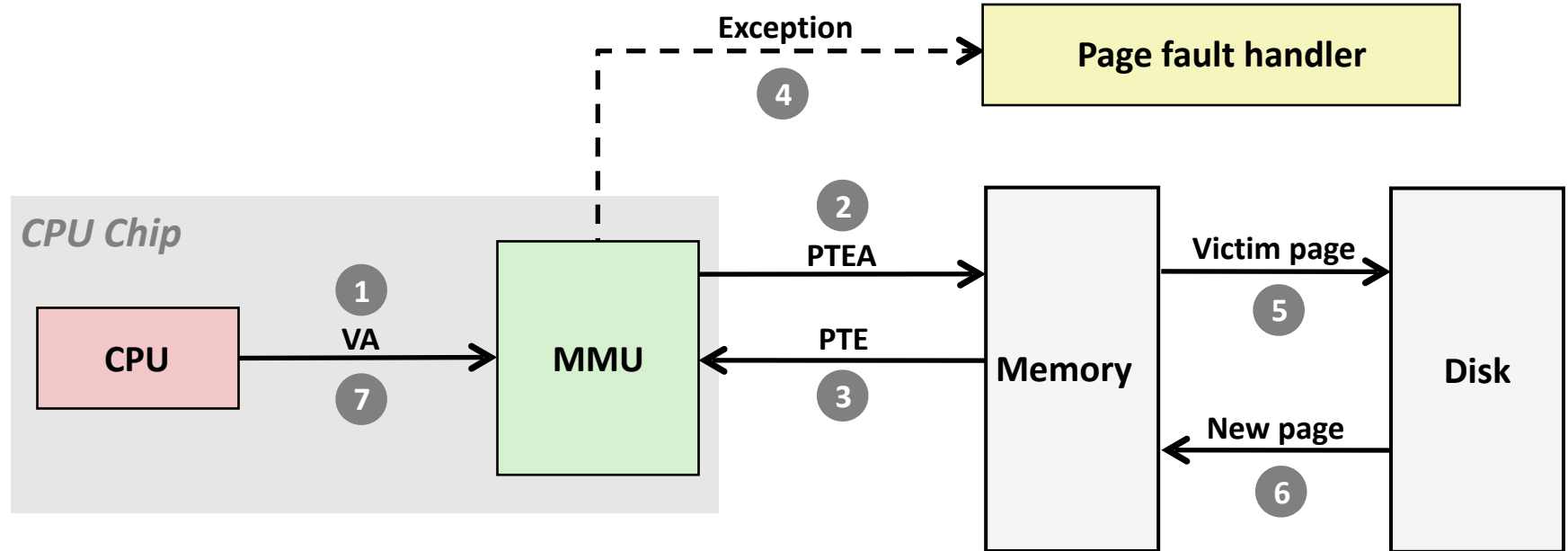2-3) MMU fetches PTE from page table in memory

# 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

# 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)
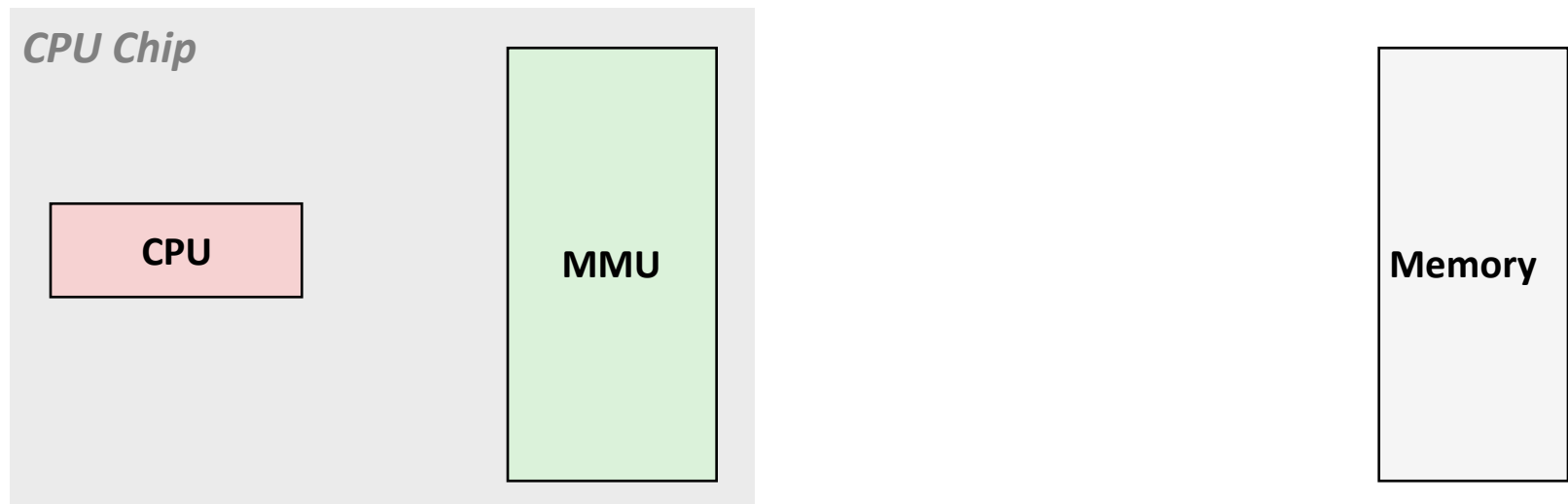
# Address Translation: Page Fault



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
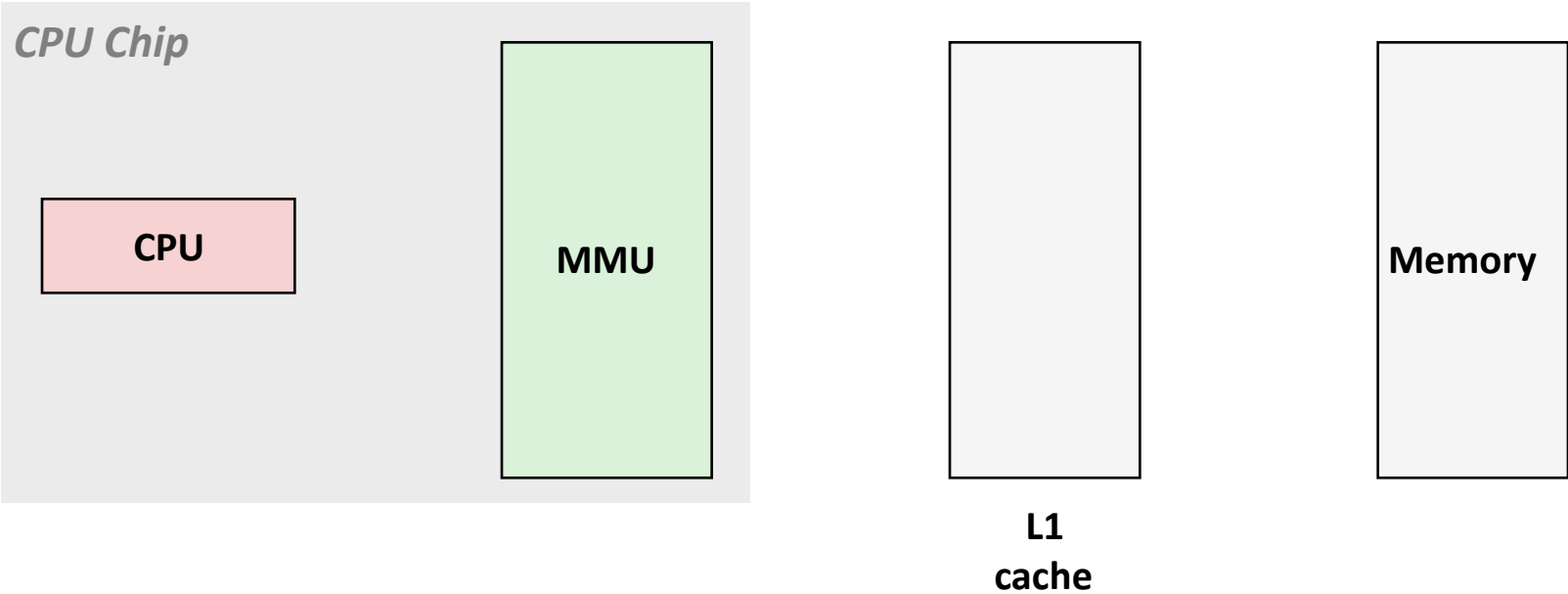
# Address Translation: Page Fault



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
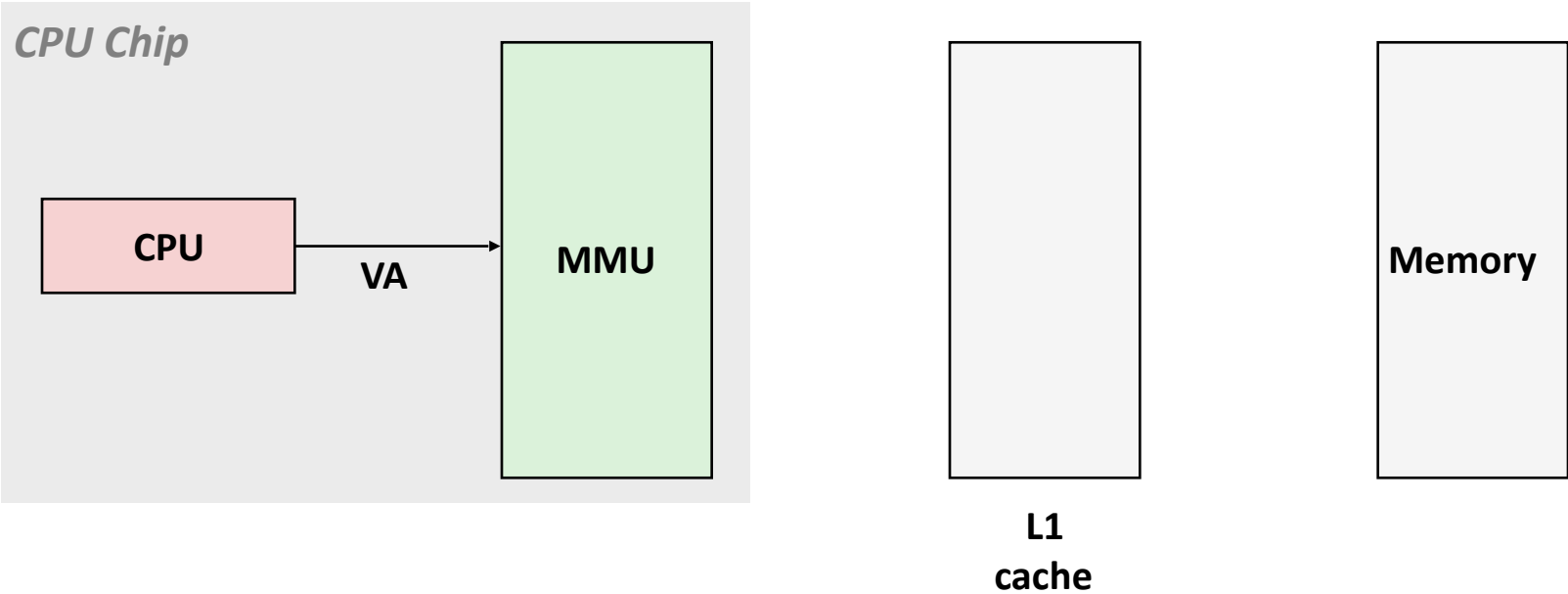
# Integrating VM and Cache



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

# Integrating VM and Cache



*VA: virtual address, PA: physical address, PTE: page table entry, PTEA = PTE address*
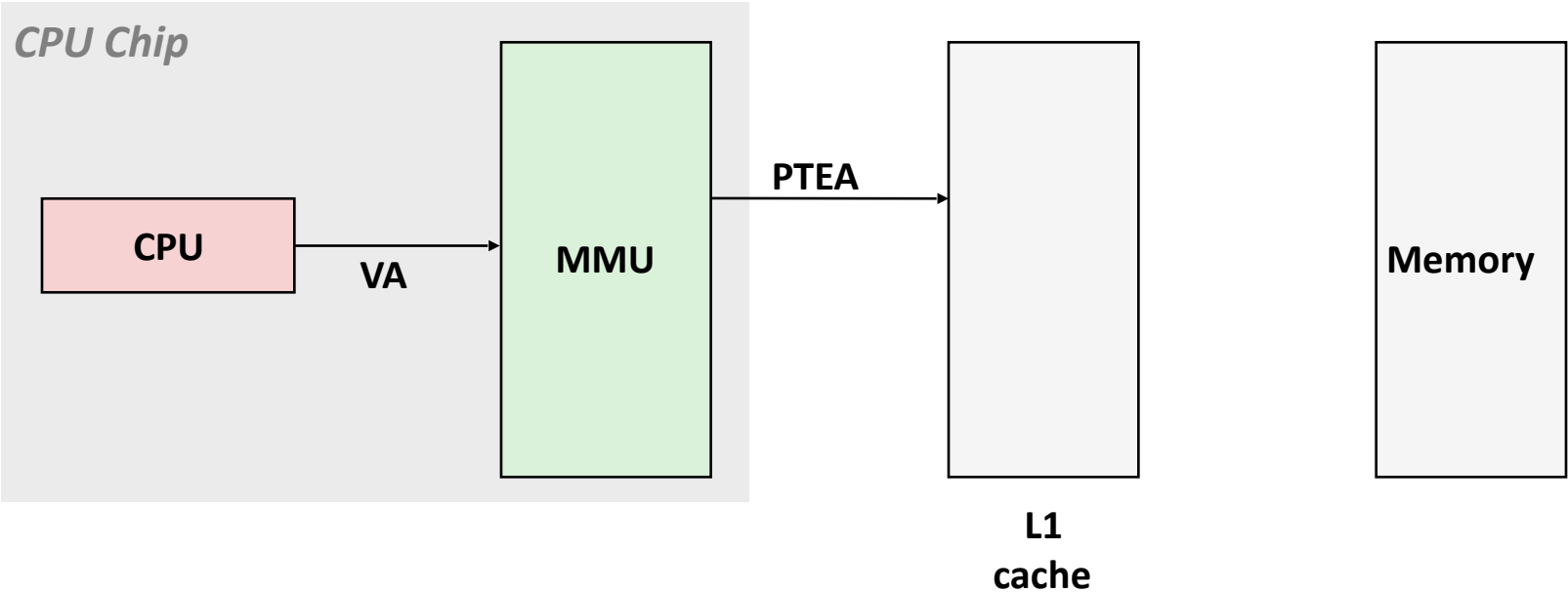
# Integrating VM and Cache



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

# Integrating VM and Cache



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

# Integrating VM and Cache



*VA: virtual address, PA: physical address, PTE: page table entry, PTEA = PTE address*
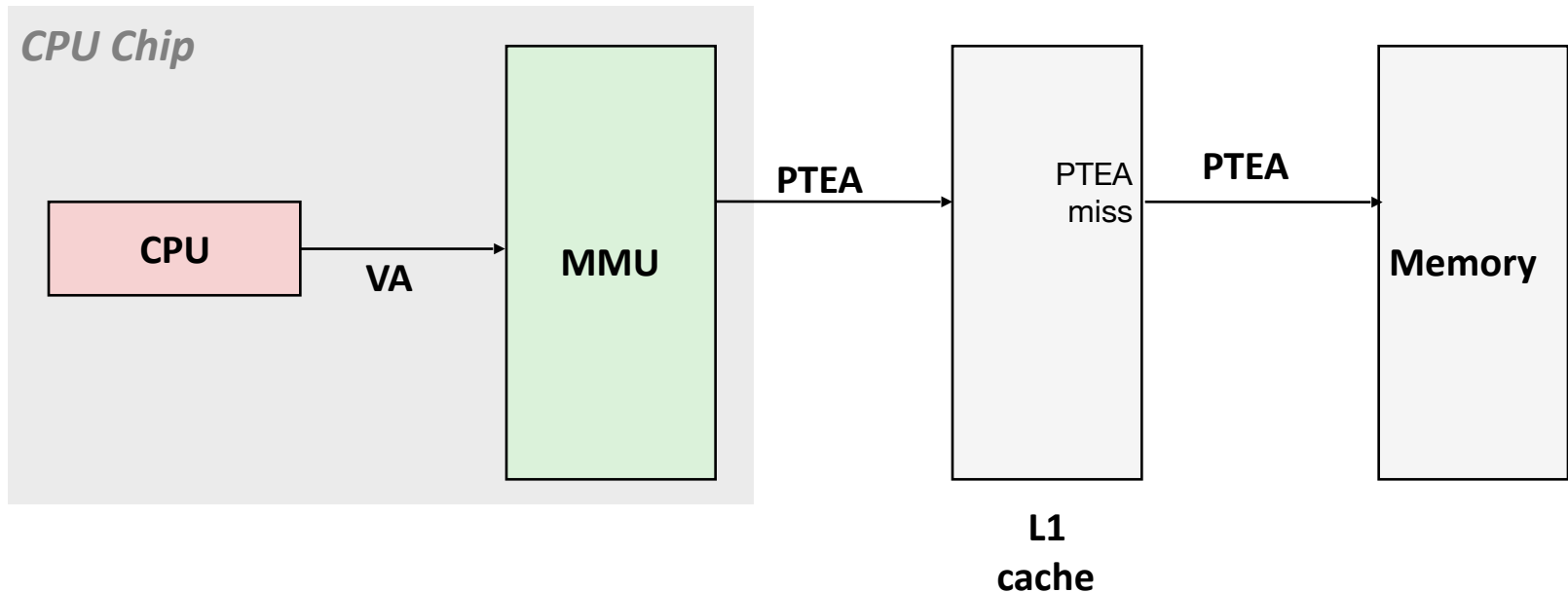
# Integrating VM and Cache



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

# Integrating VM and Cache



*VA: virtual address, PA: physical address, PTE: page table entry, PTEA = PTE address*
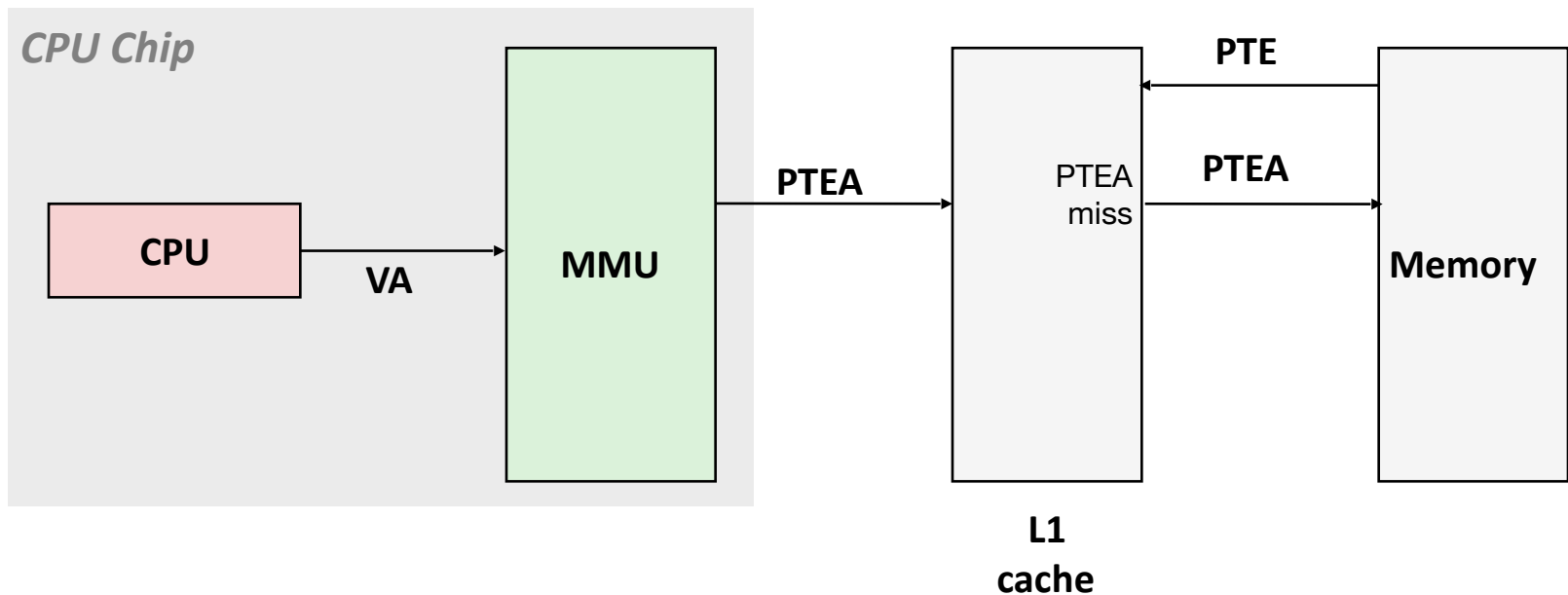
# Integrating VM and Cache



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

# Integrating VM and Cache



*VA: virtual address, PA: physical address, PTE: page table entry, PTEA = PTE address*
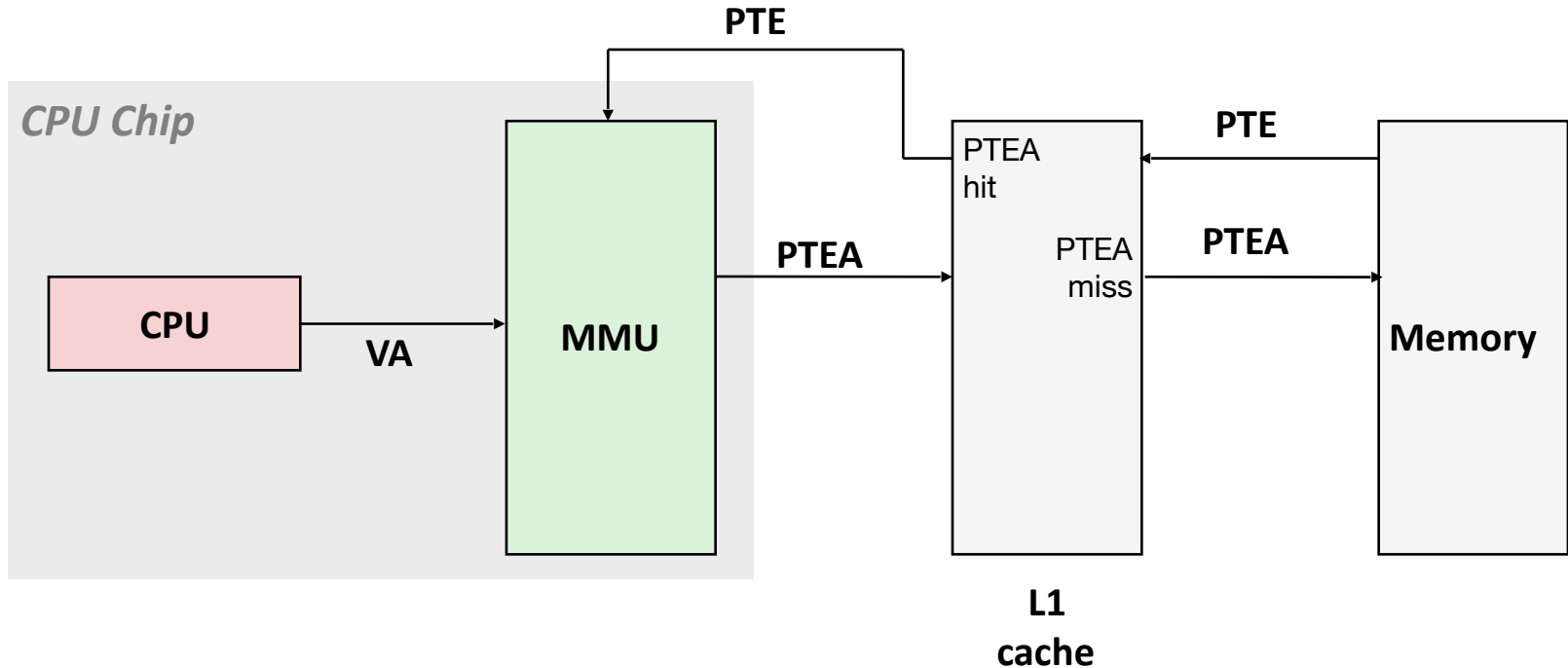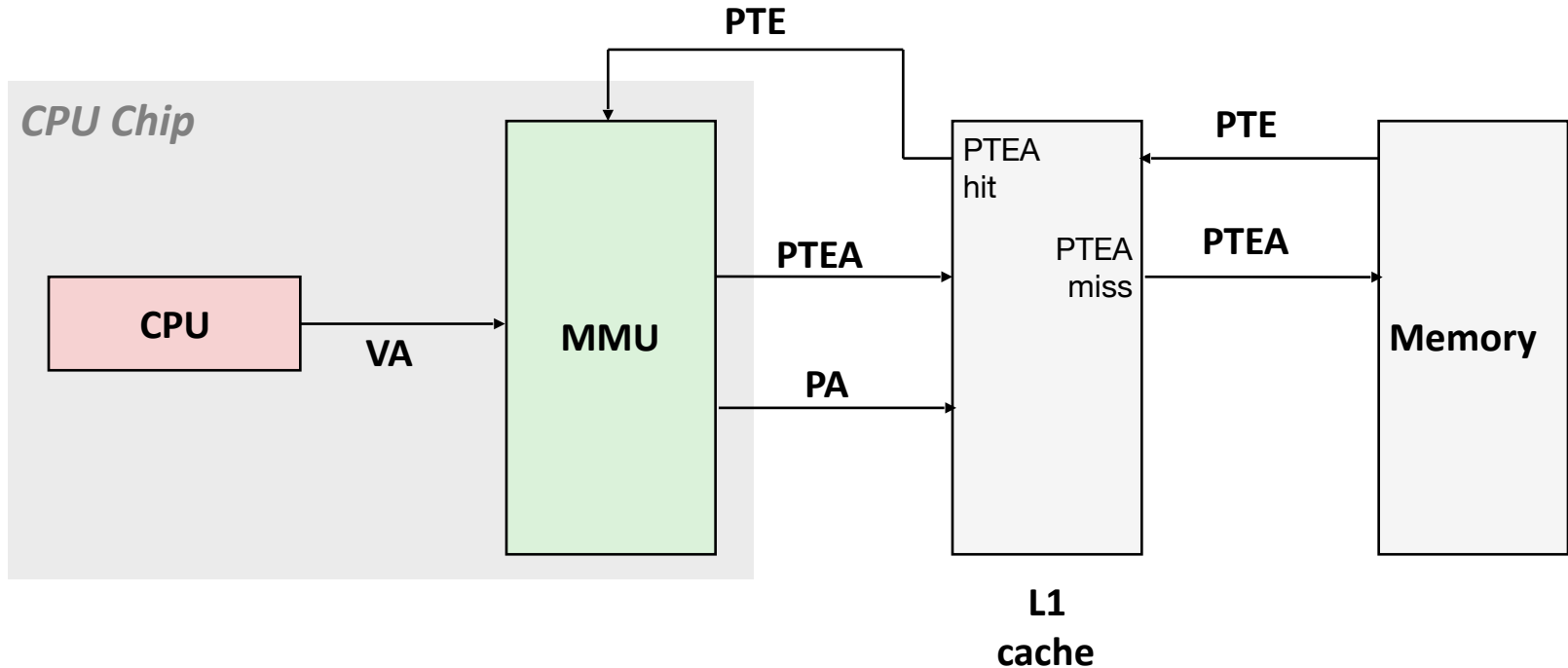
# Integrating VM and Cache



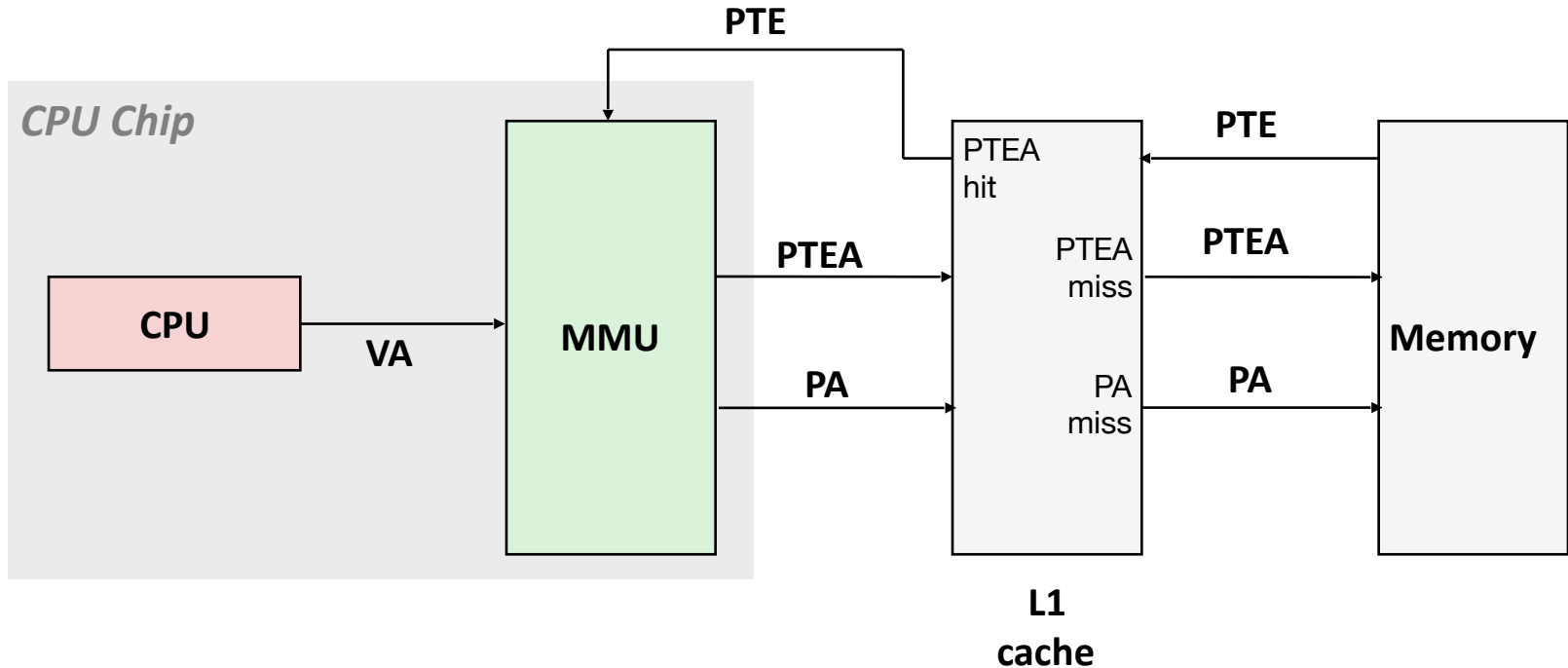*VA: virtual address, PA: physical address, PTE: page table entry, PTEA = PTE address*

# Integrating VM and Cache



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

# Today

- Three Virtual Memory Optimizations
  - TLB
  - Page the page table (a.k.a., multi-level page table)
  - Virtually-indexed, physically-tagged cache
- Case-study: Intel Core i7/Linux example

# Speeding up Address Translation

# Speeding up Address Translation

- Problem: Every memory load/store requires two memory accesses: one for PTE, another for real

  - The PTE access is kind of an overhead

  - Can we speed it up?

# Speeding up Address Translation

- Problem: Every memory load/store requires two memory accesses: one for PTE, another for real

  - The PTE access is kind of an overhead

  - Can we speed it up?

- Page table entries (PTEs) are already cached in L1 data cache like any other memory data. But:

  - PTEs may be evicted by other data references

  - PTE hit still requires a small L1 delay

# Speeding up Translation with a TLB

- Solution: *Translation Lookaside Buffer* (TLB)
    - Think of it as a dedicated cache for page table
    - Small set-associative hardware cache in MMU
    - Contains complete page table entries for a small number of pages

# Speeding up Translation with a TLB

- Solution: *Translation Lookaside Buffer* (TLB)
    - Think of it as a dedicated cache for page table
    - Small set-associative hardware cache in MMU
    - Contains complete page table entries for a small number of pages

| Tag | Set Index |
|-----|-----------|

# Speeding up Translation with a TLB

- Solution: *Translation Lookaside Buffer* (TLB)
  - Think of it as a dedicated cache for page table
  - Small set-associative hardware cache in MMU
  - Contains complete page table entries for a small number of pages

| Tag | Set Index |
|-----|-----------|

**Set 0**  | v | tag | Data |    | v | tag | Data |

**Set 1**  | v | tag | Data |    | v | tag | Data |

⋮

**Set T-1**  | v | tag | Data |    | v | tag | Data |
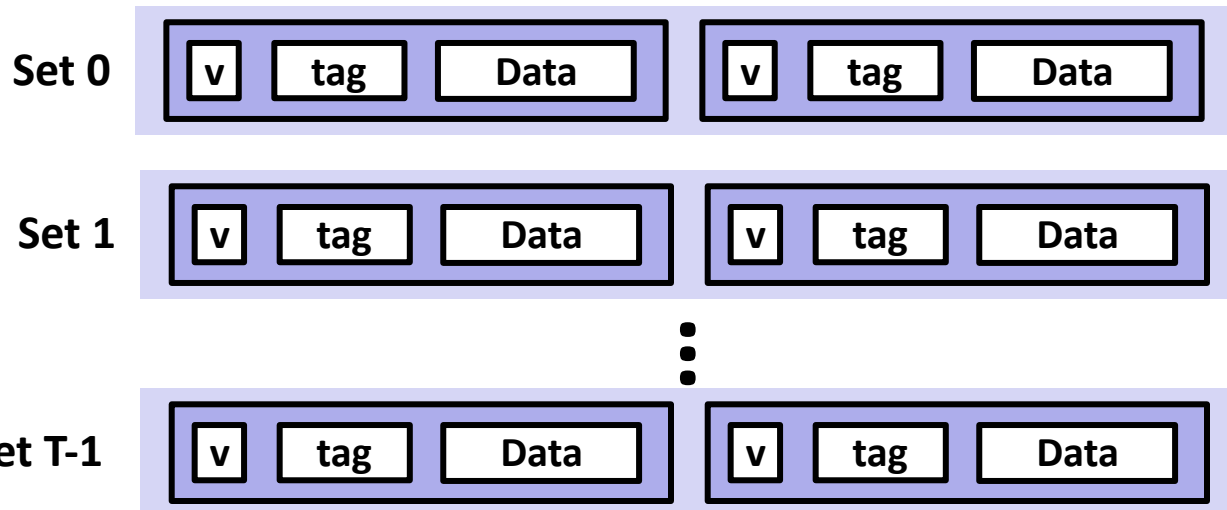
**A Conventional Data Cache**

# Speeding up Translation with a TLB

- Solution: *Translation Lookaside Buffer* (TLB)
  - Think of it as a dedicated cache for page table
  - Small set-associative hardware cache in MMU
  - Contains complete page table entries for a small number of pages

| Tag | Set Index |
|-----|-----------|

**Set 0** | v | tag | Data | v | tag | Data |

**Set 1** | v | tag | Data | v | tag | Data |

Set Index selects a set

⋮

**Set T-1** | v | tag | Data | v | tag | Data |

**A Conventional Data Cache**

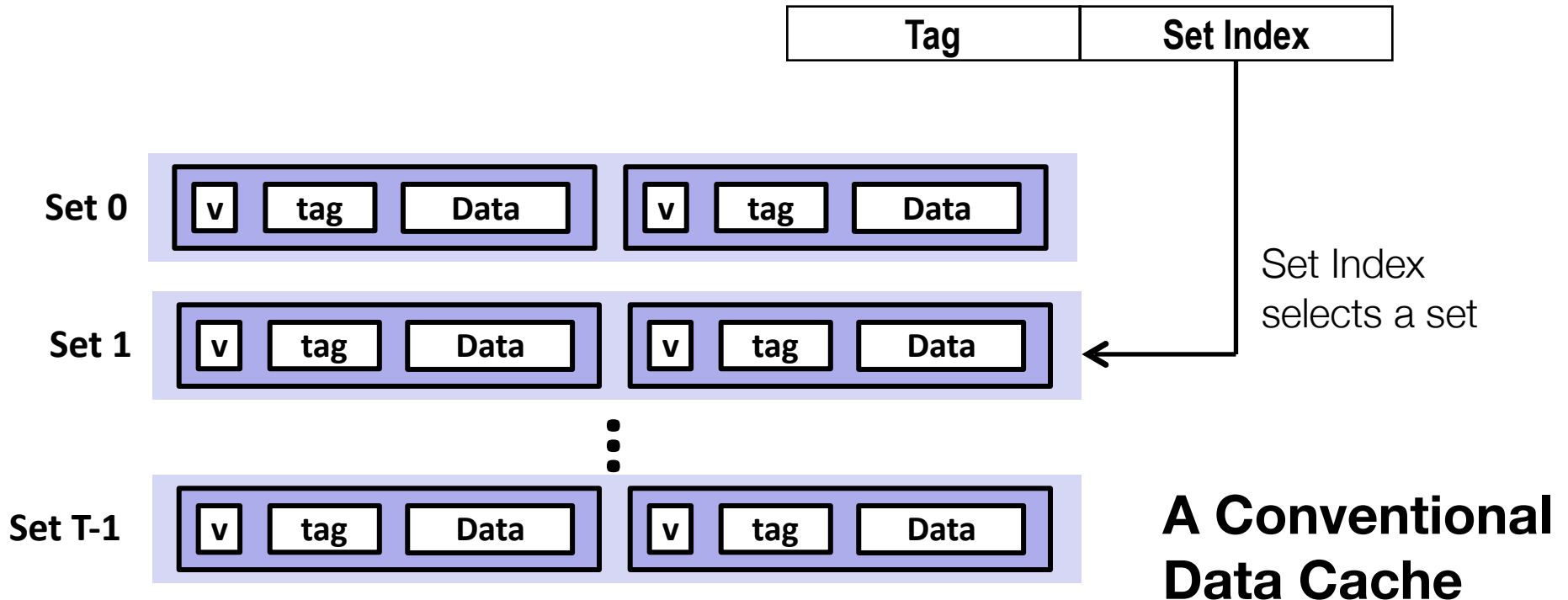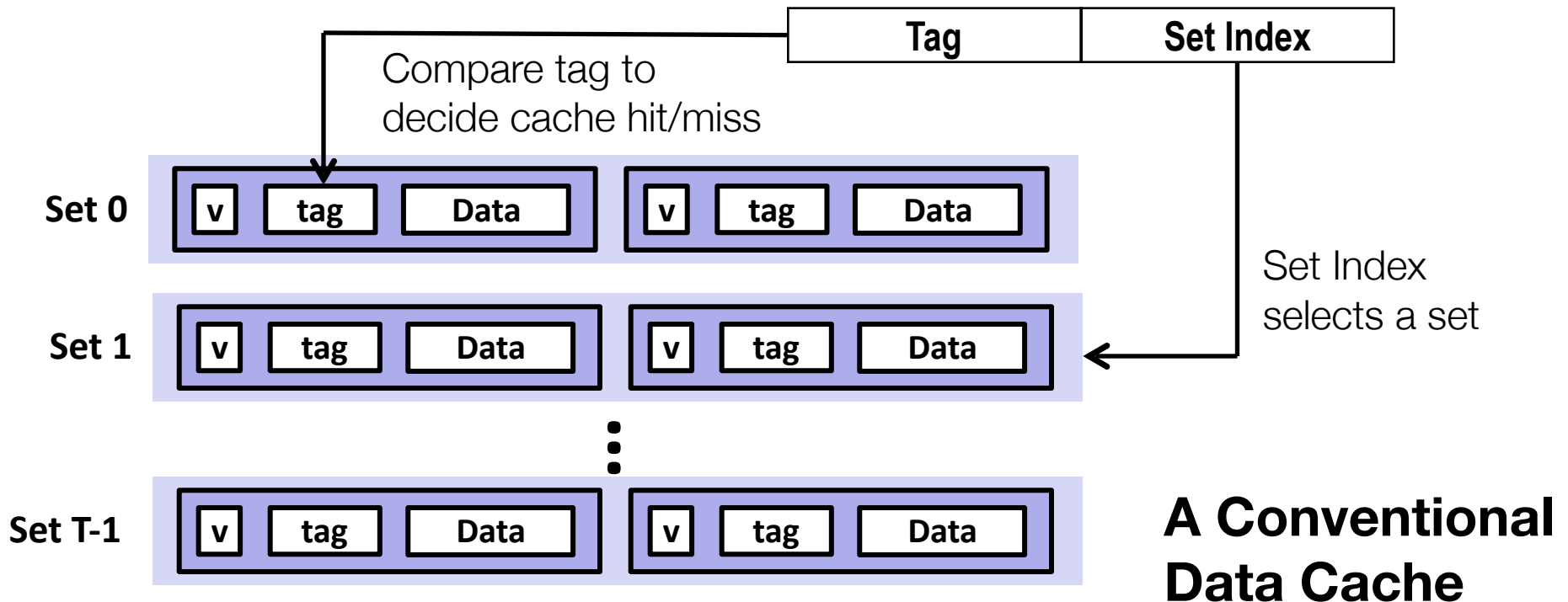# Speeding up Translation with a TLB

- Solution: *Translation Lookaside Buffer* (TLB)
  - Think of it as a dedicated cache for page table
  - Small set-associative hardware cache in MMU
  - Contains complete page table entries for a small number of pages



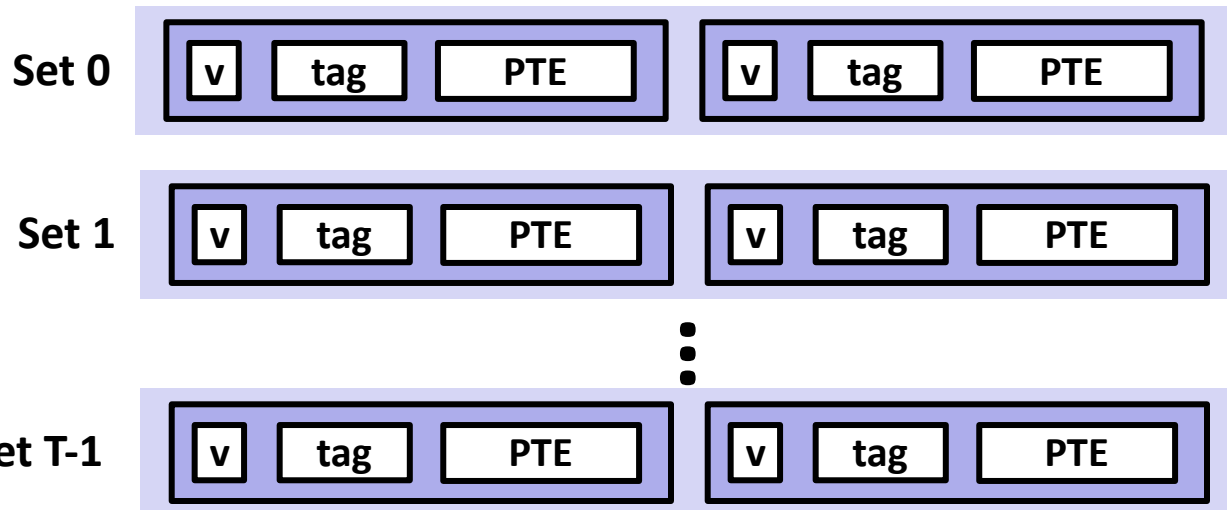| Tag | Set Index |
|-----|-----------|

Compare tag to decide cache hit/miss

**Set 0** | v | tag | Data | v | tag | Data |

Set Index selects a set

**Set 1** | v | tag | Data | v | tag | Data |

**Set T-1** | v | tag | Data | v | tag | Data |

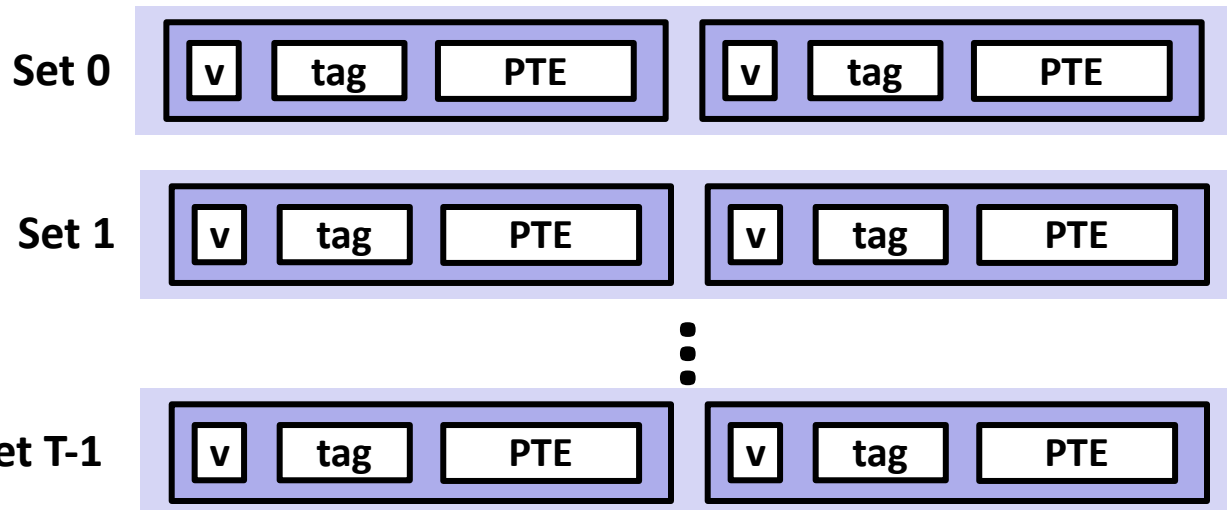**A Conventional Data Cache**

# Accessing the TLB

- MMU uses the Virtual Page Number portion of the virtual address to access the TLB:



**A Page Table Cache**

# Accessing the TLB

- MMU uses the Virtual Page Number portion of the virtual address to access the TLB:

**Virtual Page Number**

| n-1 | p+t | p+t-1 | p | p-1 | 0 |
|---|---|---|---|---|---|
| TLB tag (TLBT) | | TLB index (TLBI) | | Offset | |



**Set 0** | v | tag | PTE | v | tag | PTE |

**Set 1** | v | tag | PTE | v | tag | PTE |

...

**Set T-1** | v | tag | PTE | v | tag | PTE |

**A Page Table Cache**

# Accessing the TLB

- MMU uses the Virtual Page Number portion of the virtual address to access the TLB:



**Virtual Page Number**

| n-1 | p+t | p+t-1 | p | p-1 | 0 |
|---|---|---|---|---|---|
| TLB tag (TLBT) | | TLB index (TLBI) | | Offset | |

**TLBI selects the set**

Set 0    | v | tag | PTE |    | v | tag | PTE |

Set 1    | v | tag | PTE |    | v | tag | PTE |

Set T-1    | v | tag | PTE |    | v | tag | PTE |

**A Page Table Cache**

# Accessing the TLB

- MMU uses the Virtual Page Number portion of the virtual address to access the TLB:
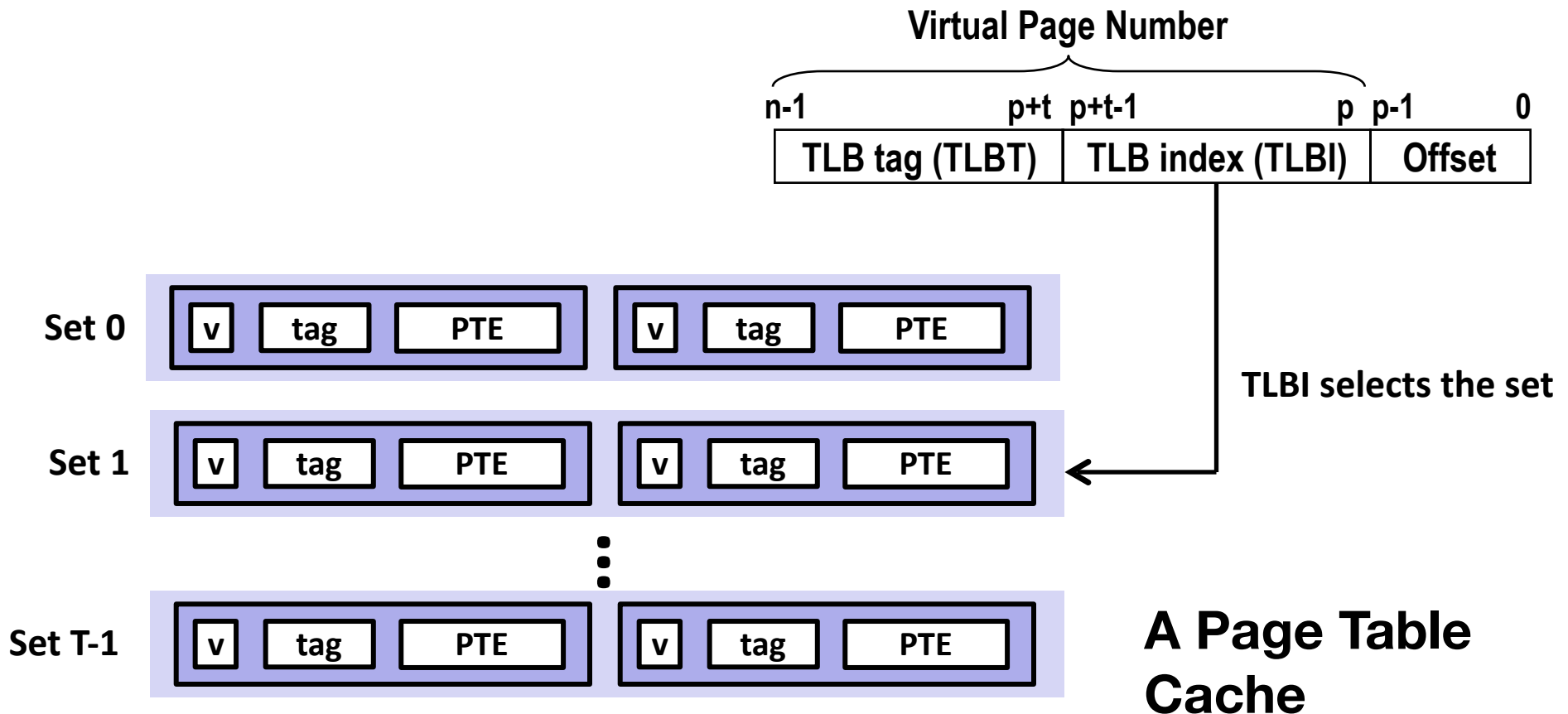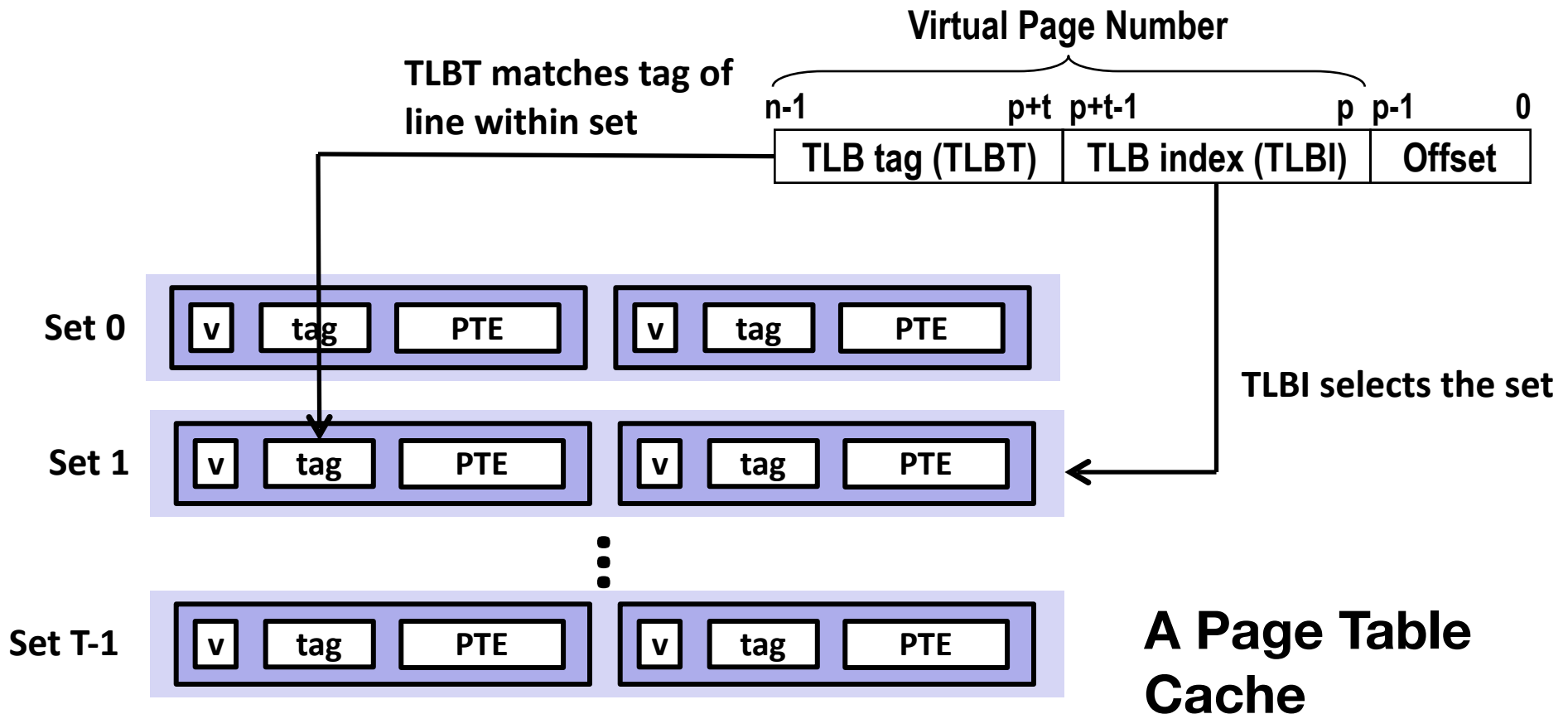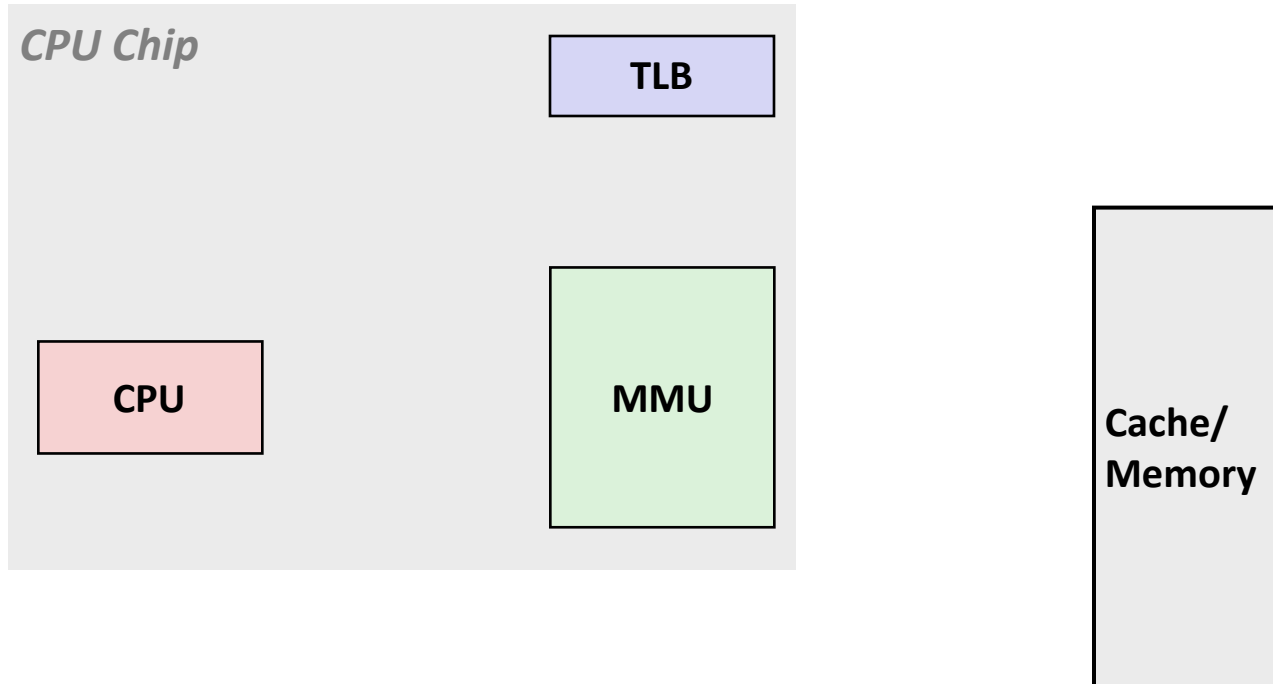


**Virtual Page Number**

**TLBT matches tag of line within set**

| n-1 | p+t | p+t-1 | p | p-1 | 0 |
|---|---|---|---|---|---|
| TLB tag (TLBT) | | TLB index (TLBI) | | Offset | |

**Set 0** | v | tag | PTE | v | tag | PTE |

**TLBI selects the set**

**Set 1** | v | tag | PTE | v | tag | PTE |

**Set T-1** | v | tag | PTE | v | tag | PTE |

**A Page Table Cache**

# TLB Hit

# TLB Hit

# TLB Hit

# TLB Hit

# TLB Hit

# TLB Hit

# TLB Hit



A TLB hit eliminates a memory access

# TLB Miss

# TLB Miss

# TLB Miss

# TLB Miss

# TLB Miss

# TLB Miss



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

# Today

- Three Virtual Memory Optimizations
  - TLB
  - Page the page table (a.k.a., multi-level page table)
  - Virtually-indexed, physically-tagged cache
- Case-study: Intel Core i7/Linux example

# Where Does Page Table Live?

# Where Does Page Table Live?

- It needs to be at a specific location where we can find it
  - In main memory, with its start address stored in a special register (PTBR)

# Where Does Page Table Live?

- It needs to be at a specific location where we can find it
  - In main memory, with its start address stored in a special register (PTBR)
- Assume 4KB page, 48-bit virtual memory, each PTE is 8 Bytes
  - $2^{36}$ PTEs in a page table
  - 512 GB total size per page table??!!

# Where Does Page Table Live?

- It needs to be at a specific location where we can find it
  - In main memory, with its start address stored in a special register (PTBR)
- Assume 4KB page, 48-bit virtual memory, each PTE is 8 Bytes
  - $2^{36}$ PTEs in a page table
  - 512 GB total size per page table??!!
- Problem: Page tables are huge
  - One table per process!
  - Storing them all in main memory wastes space

# Solution: Page the Page Table

- Observation: Only a small number of pages (working set) are accessed during a certain period of time, due to locality
- Put only the relevant page table entires in main memory
- Idea: Put Page Table in Virtual Memory and swap it just like data

**PM**

**VM**

# Solution: Page the Page Table

- Observation: Only a small number of pages (working set) are accessed during a certain period of time, due to locality
- Put only the relevant page table entires in main memory
- Idea: Put Page Table in Virtual Memory and swap it just like data



PM

VM

# Solution: Page the Page Table

- Observation: Only a small number of pages (working set) are accessed during a certain period of time, due to locality
- Put only the relevant page table entires in main memory
- Idea: Put Page Table in Virtual Memory and swap it just like data

**VM**

**PM**

# Solution: Page the Page Table

- Observation: Only a small number of pages (working set) are accessed during a certain period of time, due to locality
- Put only the relevant page table entires in main memory
- Idea: Put Page Table in Virtual Memory and swap it just like data

# Solution: Page the Page Table

- Observation: Only a small number of pages (working set) are accessed during a certain period of time, due to locality
- Put only the relevant page table entires in main memory
- Idea: Put Page Table in Virtual Memory and swap it just like data

**PM**

**VM**

# Solution: Page the Page Table

- Observation: Only a small number of pages (working set) are accessed during a certain period of time, due to locality
- Put only the relevant page table entires in main memory
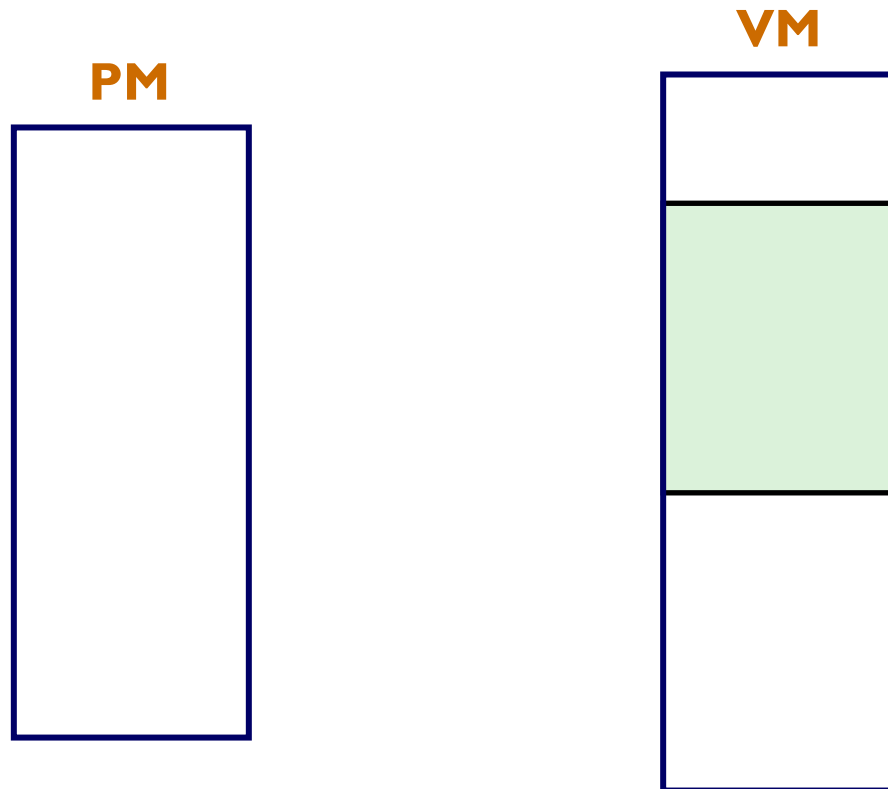- Idea: Put Page Table in Virtual Memory and swap it just like data

# Solution: Page the Page Table

- Observation: Only a small number of pages (working set) are accessed during a certain period of time, due to locality
- Put only the relevant page table entires in main memory
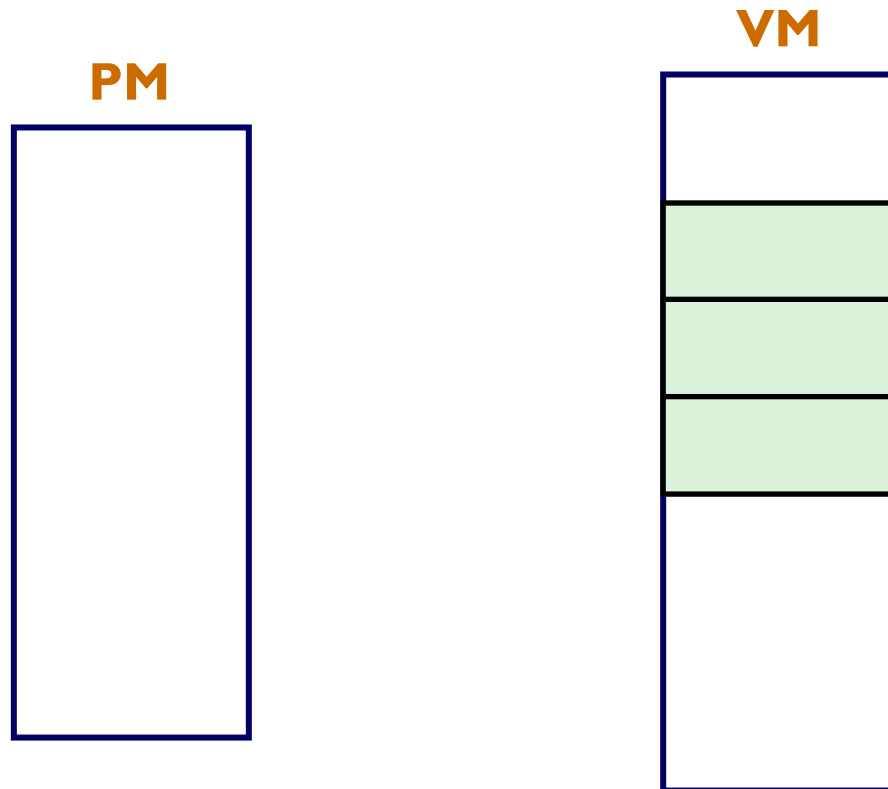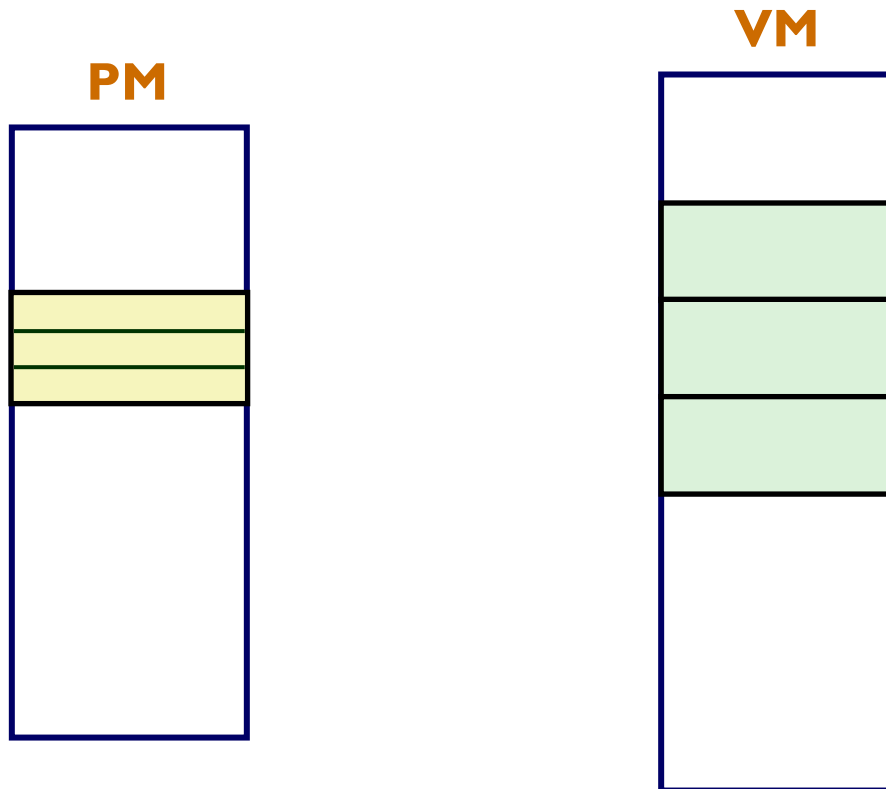- Idea: Put Page Table in Virtual Memory and swap it just like data

# Solution: Page the Page Table

- Observation: Only a small number of pages (working set) are accessed during a certain period of time, due to locality
- Put only the relevant page table entires in main memory
- Idea: Put Page Table in Virtual Memory and swap it just like data

**VM**

**PM**

Virtual address

# Effectively: A 2-Level Page Table

- Level 1 table:
  - Always in memory at a known location.
  - Each L1 PTE points to the start address of a L2 page table.
  - Bring that table to memory on-demand.
- Level 2 table:
  - Each PTE points to an actual data page

**Level 2**

**Tables**

**Level 1**

**Table**

...

# A Two-Level Page Table Hierarchy

Virtual
memory

| |
|:---:|
| **VP 0** |
| **...** |
| **VP 1023** |
| **VP 1024** |
| **...** |
| **VP 2047** |
| **unallocated pages** |
| **unallocated pages** |
| **VP 9215** |

*32 bit addresses, 4KB pages, 4-byte PTEs*

**...**

# A Two-Level Page Table Hierarchy



*32 bit addresses, 4KB pages, 4-byte PTEs*

**...**

# A Two-Level Page Table Hierarchy

Level 1
page table

Level 2
page tables

Virtual
memory



*32 bit addresses, 4KB pages, 4-byte PTEs*

...

# A Two-Level Page Table Hierarchy

Level 1
page table

Level 2
page tables

Virtual
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 |
| ... |
| VP 1023 |
| VP 1024 |
| ... |
| VP 2047 |
| unallocated pages |
| unallocated pages |
| VP 9215 |

- Level 2 page table size:
  - $2^{32} / 2^{12} * 4 = 4$ MB
- Level 1 page table size:
  - $(2^{32} / 2^{12} * 4) / 2^{12} * 4 = 4$ KB

*32 bit addresses, 4KB pages, 4-byte PTEs*

...

# How to Access a 2-Level Page Table?

| VPN | VPO |
|-----|-----|

**Page Table**

**0001**

0000
0001
0010
0011

0100
0101
0110
0111

⋮

1100
1101
1110
1111

# How to Access a 2-Level Page Table?

| VPN | VPO |
|---|---|

**Level 2 Tables**

**0001**

**Level 1 Table**

| | |
|---|---|
| 00 | |
| 01 | |
| | ... |
| 11 | |

| | |
|---|---|
| | 00 |
| | 01 |
| | 10 |
| | 11 |

| | |
|---|---|
| | 00 |
| | 01 |
| | 10 |
| | 11 |

| | |
|---|---|
| | 00 |
| | 01 |
| | 10 |
| | 11 |

# How to Access a 2-Level Page Table?

| VPN 1 | VPN 2 | VPO |
|-------|-------|-----|

**Level 2 Tables**

**0001**

**Level 1 Table**

Level 1 Table entries:
00
01
...
11

Level 2 Table (top):
00
01
10
11

Level 2 Table (middle):
00
01
10
11

Level 2 Table (bottom):
00
01
10
11

# How to Access a 2-Level Page Table?

| VPN 1 | VPN 2 | VPO |
|-------|-------|-----|



Level 2 Tables

**01**

Level 1
Table

Level 1 Table entries:
00
01
...
11

Level 2 Tables (top):
00
01
10
11

Level 2 Tables (middle):
00
01
10
11

Level 2 Tables (bottom):
00
01
10
11

# How to Access a 2-Level Page Table?

# How to Access a 2-Level Page Table?

**Page table base register (PTBR)**

**VIRTUAL ADDRESS**

n-1 | VPN | p-1 | VPO | 0

**page table**

PPN

m-1 | PPN | p-1 | PPO | 0

**PHYSICAL ADDRESS**

# How to Access a 2-Level Page Table?

# Translating with a k-level Page Table

# Today

- Three Virtual Memory Optimizations
  - TLB
  - Page the page table (a.k.a., multi-level page table)
  - Virtually-indexed, physically-tagged cache
- Case-study: Intel Core i7/Linux example

# Performance Issue in VM

- Address translation and cache accesses are serialized
  - First translate from VA to PA
  - Then use PA to access cache
  - Slow! Can we speed it up?

# Performance Issue in VM

**Virtual Address**

| Virtual page number (VPN) | Page Offset |
|---|---|

**Physical Address**

| Physical page number (PPN) | Page Offset |
|---|---|

| Tag | Set Index | Cache Line Offset |
|---|---|---|

**L1 cache**

# Performance Issue in VM



**Virtual Address**

| Virtual page number (VPN) | Page Offset |
|---|---|

**Unchanged!!**

**Physical Address**

| Physical page number (PPN) | Page Offset |
|---|---|

| Tag | Set Index | Cache Line Offset |
|---|---|---|

**L1 cache**

# Performance Issue in VM

# Performance Issue in VM

**Virtual Address**

| Virtual page number (VPN) | Page Offset |
|---|---|

**Unchanged!!**

**Physical Address**

| Physical page number (PPN) | Page Offset |
|---|---|

**=**

| Tag | Set Index | Cache Line Offset |
|---|---|---|

**L1 cache**

- Set Index + Cache Line Offset = Page Offset
- Indexing into cache in parallel with translation (TLB access)
- If TLB hits, can get the data back in one cycle

# Performance Issue in VM

| Virtual Address | Virtual page number (VPN) | Page Offset |
|---|---|---|

**Unchanged!!**

| Physical Address | Physical page number (PPN) | Page Offset |
|---|---|---|

Virtually-Indexed, Physically-Tagged Cache

$=$

| Tag | Set Index | Cache Line Offset |
|---|---|---|

**L1 cache**

- Set Index + Cache Line Offset = Page Offset
- Indexing into cache in parallel with translation (TLB access)
- If TLB hits, can get the data back in one cycle

# Any Implications?

Virtual
Address

| Virtual page number (VPN) | Page Offset |
|---|---|

Physical
Address

| Tag | Set Index | Cache Line Offset |
|---|---|---|

# Any Implications?

12 bits

**Virtual Address**

| Virtual page number (VPN) | Page Offset |
|---|---|

4 bits

**Physical Address**

| Tag | Set Index | Cache Line Offset |
|---|---|---|

- Assuming 4K page size, cache line size is 16 bytes.

# Any Implications?



- Assuming 4K page size, cache line size is 16 bytes.
- Set Index = 8 bits. Can only have 256 Sets => Limit cache size

# Any Implications?



- Assuming 4K page size, cache line size is 16 bytes.
- Set Index = 8 bits. Can only have 256 Sets => Limit cache size
- Increasing cache size then requires increasing associativity

# Any Implications?



- Assuming 4K page size, cache line size is 16 bytes.
- Set Index = 8 bits. Can only have 256 Sets => Limit cache size
- Increasing cache size then requires increasing associativity
  - Not ideal because that requires comparing more tags

# Any Implications?



- Assuming 4K page size, cache line size is 16 bytes.
- Set Index = 8 bits. Can only have 256 Sets => Limit cache size
- Increasing cache size then requires increasing associativity
  - Not ideal because that requires comparing more tags
- Solutions?

# Any Implications?



- What if we use 9 bits for Set Index? More Sets now.

# Any Implications?



- What if we use 9 bits for Set Index? More Sets now.
- How can this still work???

# Any Implications?



- What if we use 9 bits for Set Index? More Sets now.
- How can this still work???
- The least significant bit in VPN and PPN must be the same

# Any Implications?



12 bits

**Virtual Address**

| Virtual page number (VPN) | Page Offset |

**9 bits**     4 bits

**Physical Address**

| Tag | Set Index | Cache Line Offset |

- What if we use 9 bits for Set Index? More Sets now.
- How can this still work???
- The least significant bit in VPN and PPN must be the same
- That is: an even VA must be mapped to an even PA, and an odd VA must be mapped to an odd PA

# Today

- Three Virtual Memory Optimizations
  - TLB
  - Page the page table (a.k.a., multi-level page table)
  - Virtually-indexed, physically-tagged cache
- **Case-study: Intel Core i7/Linux example**

# Intel Core i7 Memory System

**Processor package**

**Core x4**

| Registers | Instruction fetch | | MMU (addr translation) |
|---|---|---|---|

| L1 d-cache 32 KB, 8-way | L1 i-cache 32 KB, 8-way | L1 d-TLB 64 entries, 4-way | L1 i-TLB 128 entries, 4-way |
|---|---|---|---|

**L2 unified cache 256 KB, 8-way**

**L2 unified TLB 512 entries, 4-way**

**QuickPath interconnect 4 links @ 25.6 GB/s each**

To other cores

To I/O bridge

**L3 unified cache 8 MB, 16-way (shared by all cores)**

**DDR3 Memory controller 3 x 64 bit @ 10.66 GB/s 32 GB/s total (shared by all cores)**

**Main memory**

# End-to-End Core i7 Address Translation

**CPU**

**Virtual address (VA)**

| 36 | 12 |
|-----|------|
| VPN | VPO |

# End-to-End Core i7 Address Translation



CPU

Virtual address (VA)

36 | 12
VPN | VPO

32 | 4
TLBT | TLBI

# End-to-End Core i7 Address Translation



L1 TLB (16 sets, 4 entries/set)

# End-to-End Core i7 Address Translation



CPU

Virtual address (VA)

36        12

VPN      VPO

32    4

TLBT  TLBI

TLB hit

L1 TLB (16 sets, 4 entries/set)

40

PPN

# End-to-End Core i7 Address Translation



**CPU**

Virtual address (VA)

36 — VPN
12 — VPO

32 — TLBT
4 — TLBI

*TLB hit*

*TLB miss*

L1 TLB (16 sets, 4 entries/set)

9 — VPN1
9 — VPN2
9 — VPN3
9 — VPN4

40 — PPN

CR3

PTE  PTE  PTE  PTE

Page tables

34

# End-to-End Core i7 Address Translation



CPU

Virtual address (VA)

36     12

VPN   VPO

32   4

TLBT   TLBI

TLB hit

TLB miss

L1 TLB (16 sets, 4 entries/set)

9   9   9   9

VPN1   VPN2   VPN3   VPN4

40   12

PPN   PPO

CR3

PTE   PTE   PTE   PTE

Page tables

34

# End-to-End Core i7 Address Translation



**CPU**

Virtual address (VA)

36 / 12

**VPN** | **VPO**

32 / 4

**TLBT** | **TLBI**

*TLB hit*

*TLB miss*

**L1 TLB (16 sets, 4 entries/set)**

9 / 9 / 9 / 9

**VPN1** | **VPN2** | **VPN3** | **VPN4**

40 / 12

**PPN** | **PPO**

**Physical address (PA)**

**CR3**

**PTE** **PTE** **PTE** **PTE**

**Page tables**

# End-to-End Core i7 Address Translation



CPU

Virtual address (VA)

36 VPN  12 VPO

32 TLBT  4 TLBI

TLB hit

TLB miss

L1 TLB (16 sets, 4 entries/set)

9 VPN1  9 VPN2  9 VPN3  9 VPN4

40 PPN  12 PPO

40 CT  6 CI  6 CO

CR3

PTE  PTE  PTE  PTE

Page tables

Physical address (PA)

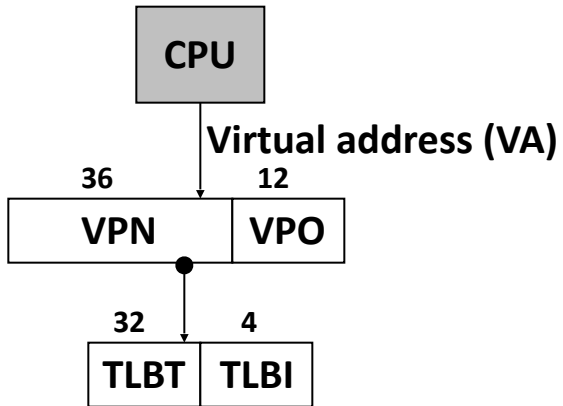# End-to-End Core i7 Address Translation

# End-to-End Core i7 Address Translation



32/64
CPU
Result

Virtual address (VA)

36 VPN  12 VPO

L1 hit

32 TLBT  4 TLBI

L1 d-cache
(64 sets, 8 lines/set)

TLB hit

TLB miss

L1 TLB (16 sets, 4 entries/set)

9 VPN1  9 VPN2  9 VPN3  9 VPN4

40 PPN  12 PPO

40 CT  6 CI  6 CO

CR3

PTE  PTE  PTE  PTE

Physical address (PA)

Page tables

# End-to-End Core i7 Address Translation



CPU

Virtual address (VA)

**36** VPN  **12** VPO

**32** TLBT  **4** TLBI

*TLB miss*

*TLB hit*

L1 TLB (16 sets, 4 entries/set)

**9** VPN1  **9** VPN2  **9** VPN3  **9** VPN4

CR3

PTE  PTE  PTE  PTE

**Page tables**

**32/64**
Result

L2, L3, and main memory

*L1 hit*

*L1 miss*

L1 d-cache
(64 sets, 8 lines/set)

**40** PPN  **12** PPO

**Physical address (PA)**

**40** CT  **6** CI  **6** CO

# Core i7 Level 4 Page Table Entries

| 63 | 62 | 52 | 51 | | | 12 | 11 | | 9 | 8 | 7 | 6 | 5 | 4 | 3 | 2 | 1 | 0 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| XD | Unused | | Page physical base address | | | | Unused | | | G | | D | A | CD | WT | U/S | R/W | P=1 |

| | |
|---|---|
| Available for OS (page location on disk) | P=0 |

## 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 Level 1-3 Page Table Entries

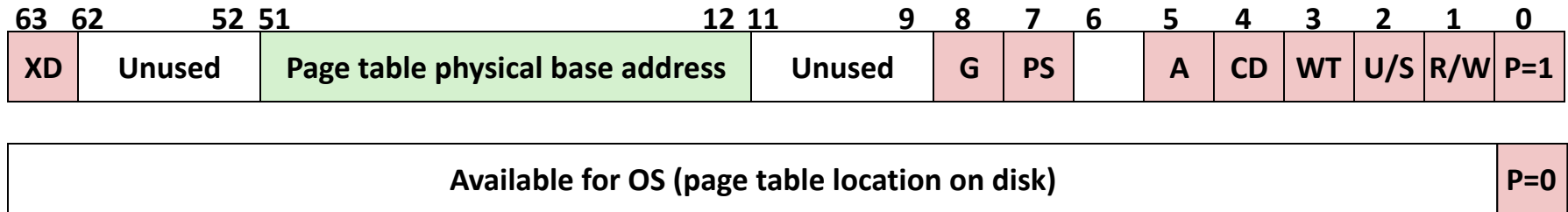| 63 | 62 | 52 | 51 | | | 12 | 11 | | 9 | 8 | 7 | 6 | 5 | 4 | 3 | 2 | 1 | 0 |
|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| XD | Unused | | Page table physical base address | | | | Unused | | | G | PS | | A | CD | WT | U/S | R/W | P=1 |

| | P=0 |
|---|---|
| Available for OS (page table location on disk) | P=0 |

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

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

**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)

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

# Core i7 Page Table Translation



**Virtual address**

| 9 | 9 | 9 | 9 | 12 |
|---|---|---|---|---|
| VPN 1 | VPN 2 | VPN 3 | VPN 4 | VPO |

**L1 PT** *Page global directory* — 512 GB region per entry

**L2 PT** *Page upper directory* — 1 GB region per entry

**L3 PT** *Page middle directory* — 2 MB region per entry

**L4 PT** *Page table* — 4 KB region per entry

CR3 *Physical address of L1 PT*

40 L1 PTE → 40 L2 PTE → 40 L3 PTE → 40 L4 PTE

*Physical address of page*

*Offset into physical and virtual page*

**Physical address**

| 40 | 12 |
|---|---|
| PPN | PPO |