# CSCE 5610
# Computer System Architecture
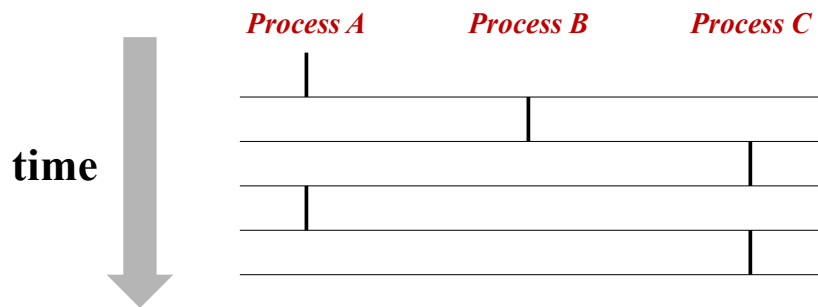
## Virtual Memory

## Processes

**Definition: A *process* is an instance of a running program**

- One of the most important ideas in computer science
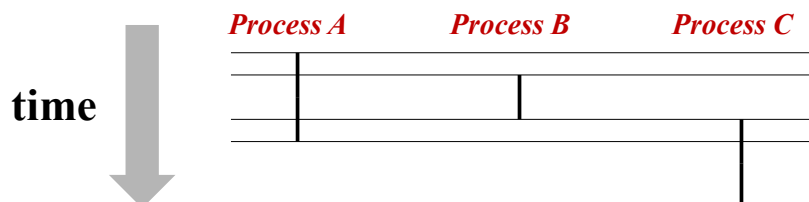- Not the same as "program" or "processor"

# Concurrent Processes

- **Two processes *run concurrently* (are concurrent) if their instruction executions (flows) overlap in time**
- **Otherwise, they are *sequential***
- **Examples:**
  - Concurrent: A & B, A & C
  - Sequential: B & C
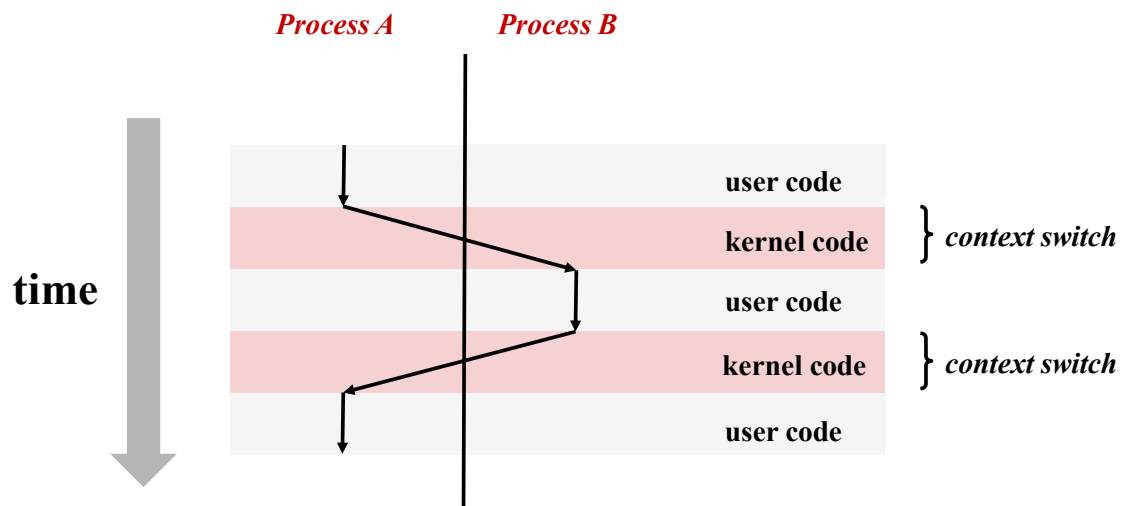
*Process A*   *Process B*   *Process C*

**time**

# User View of Concurrent Processes

- **Control flows for concurrent processes are physically disjoint in time**

- **However, we can think of concurrent processes as executing in parallel (only an illusion?)**

*Process A*   *Process B*   *Process C*

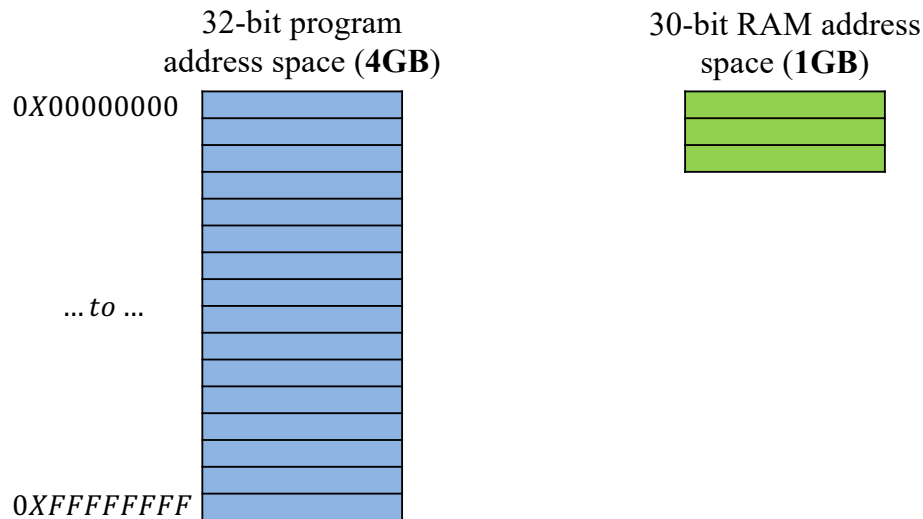**time**

# Context Switching

- **Processes are managed by a shared chunk of OS code called the *kernel***
  - Important: the kernel is not a separate process, but rather runs as part of a user process
- **Control flow passes from one process to another via a *context switch…* (how?)**

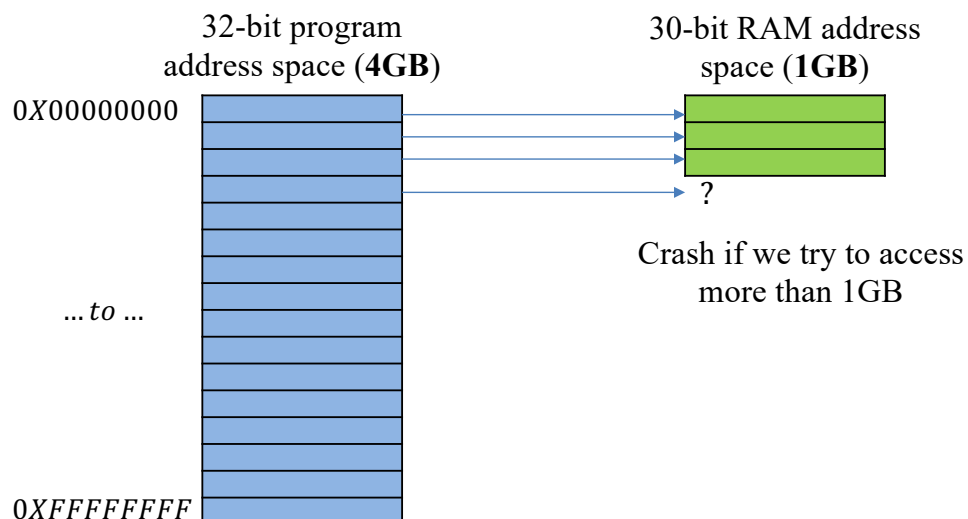| | | |
|---|---|---|
| *Process A* | *Process B* | |
| | | user code |
| | | kernel code } *context switch* |
| | | user code |
| | | kernel code } *context switch* |
| | | user code |

time

# Three problems of memory

# #1: What if we don't have enough memory

- MIPS gives each program its own 32-bit address space
- Programs can access any byte in their 32-bit address space
- How much memory can you access with a 32-bit address?
  - $2^{32} bytes = 4GB$
  - In practice, the OS reserves some of it. So, it is closer to 2GB of usable space.
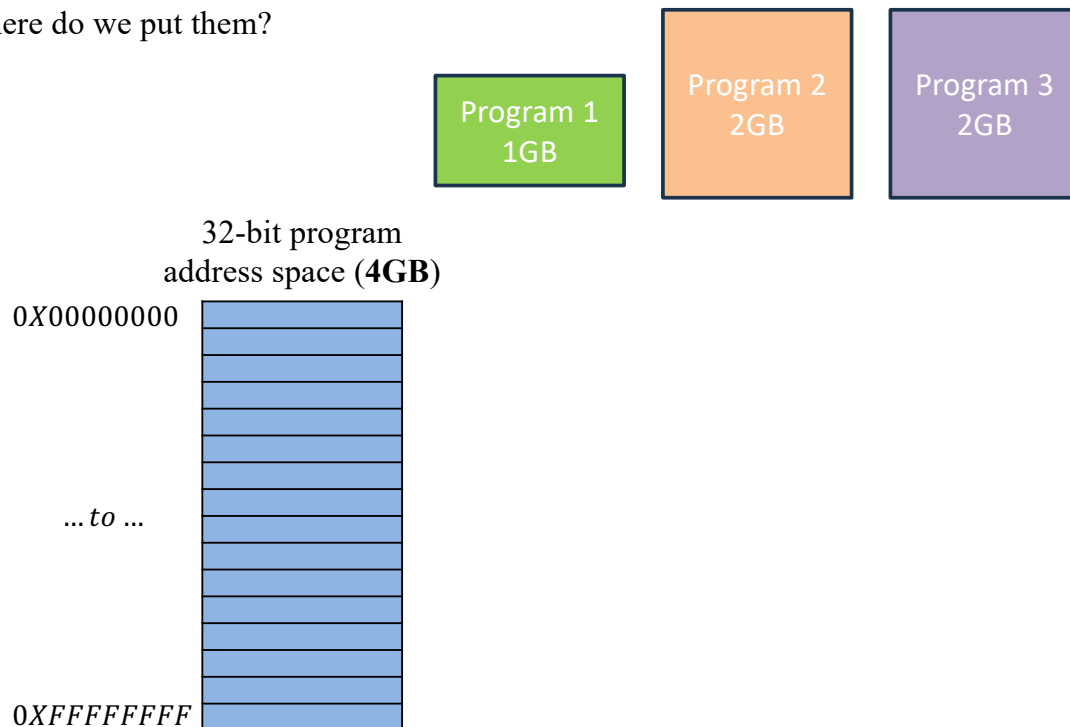- What if you don't have 4GB of memory?



32-bit program address space (**4GB**)

30-bit RAM address space (**1GB**)

0X00000000

0XFFFFFFFF

... to ...

---

# #1: What if we don't have enough memory

- What if you don't have 4GB of memory?



32-bit program address space (**4GB**)

30-bit RAM address space (**1GB**)

0X00000000

0XFFFFFFFF

... to ...

?

Crash if we try to access more than 1GB

# #2: Holes in our address space

How do programs share the memory?

Where do we put them?

Program 1
1GB

Program 2
2GB

Program 3
2GB

32-bit program
address space (**4GB**)

0X00000000

... to ...

0XFFFFFFFF

# #2: Holes in our address space

Program 3
2GB

32-bit program
address space (**4GB**)

0X00000000

Program 1
1GB

Program 2
2GB

... to ...

0XFFFFFFFF

1. Programs 1 and 2 fit
They use 3GB of memory, leaving **1GB free**

# #2: Holes in our address space

How do programs share the memory?

Where do we put them?

Program 3
2GB

32-bit program
address space (**4GB**)

0X00000000

... to ...

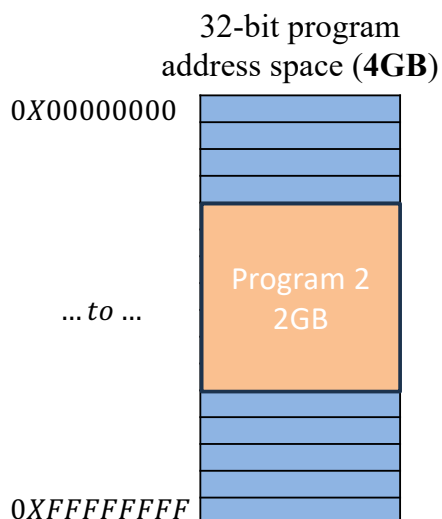Program 2
2GB

0XFFFFFFFF

1. Programs 1 and 2 fit
They use 3GB of memory, leaving **1GB free**

2. Quit Program 1
Program 2 uses 2GB of memory, leaving **2GB free**

---

# #2: Holes in our address space

How do programs share the memory?

Where do we put them?

Program 3
2GB

32-bit program
address space (**4GB**)

0X00000000

... to ...

Program 2
2GB

0XFFFFFFFF

1. Programs 1 and 2 fit
They use 3GB of memory, leaving **1GB free**

2. Quit Program 1
Program 2 uses 2GB of memory, leaving **2GB free**

3. Can't run Program 3
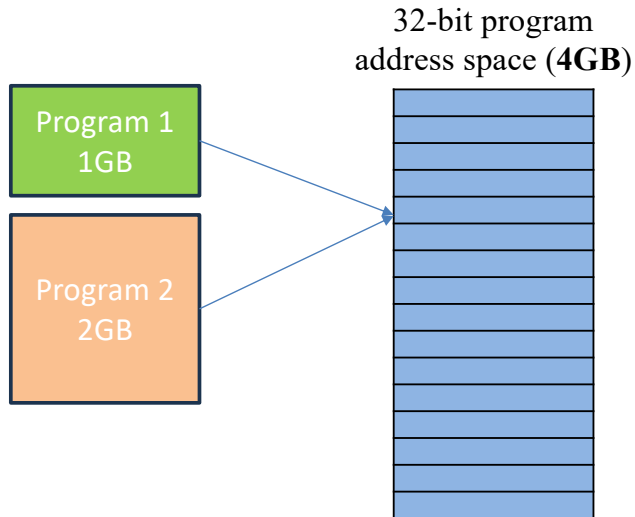Even though we have enough free space

Memory fragmentation

# #3: How do we keep programs secure?

Each program can access any 32-bit memory address

What if multiple programs access the same address?

   - sw $t0, 1024($0)

They can corrupt or crash each other: security and reliability

32-bit program
address space (**4GB**)



# Problems with memory

- **If all programs have access to the same 32-bit memory space:**
  - Can crash if less than 4GB of RAM memory in the system
  - Can run out of space if we run multiple programs
  - Can corrupt other programs' data

- **How do we solve this?**
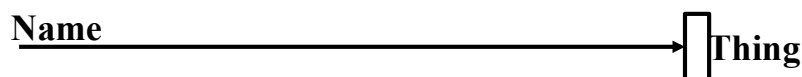
# Problems with memory

- **If all programs have access to the same 32-bit memory space:**
  - Can crash if less than 4GB of RAM memory in the system
  - Can run out of space if we run multiple programs
  - Can corrupt other programs' data

- **How do we solve this?**
  - Key to the problem: "Same memory space"

  - Can we give each program it's **own** virtual memory space?
  - If so, we can:
    - Separately **map** each program's memory space to the RAM memory space, and even move it to disk if we run out of memory
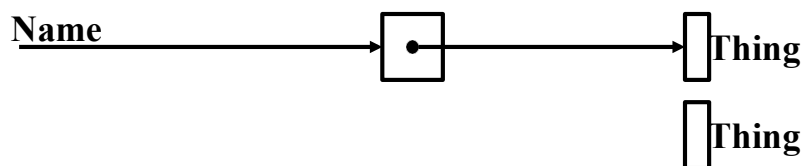
# What is Virtual Memory?

# Indirection

- **"Any problem in CS can be solved by adding a level of indirection"**

- **Without Indirection**

  **Name** ————————————————————→ ▯Thing
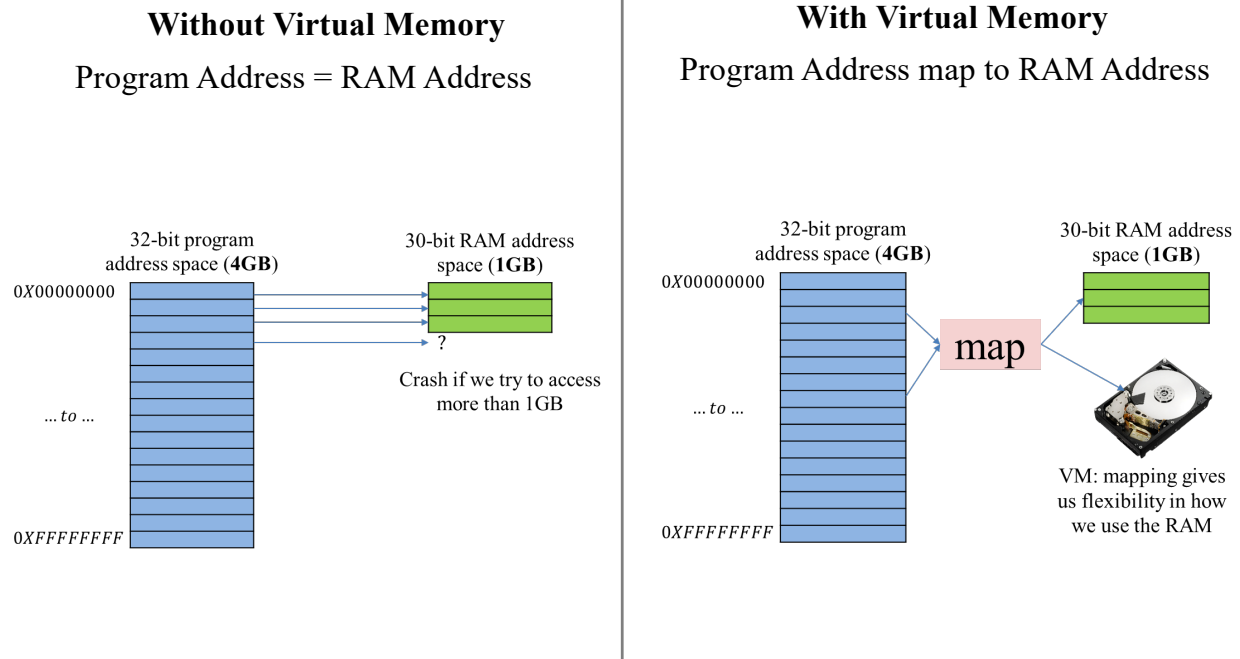
- **With Indirection**

  **Name** ————————————→ [•] ————→ ▯Thing

  ▯Thing

- Examples:

  Pointers, Domain Name Service (DNS) name->IP address, phone system (e.g., cell phone number portability), snail mail (e.g., mail forwarding), 911 (routed to local office), DHCP, call centers that route calls to available operators, etc.
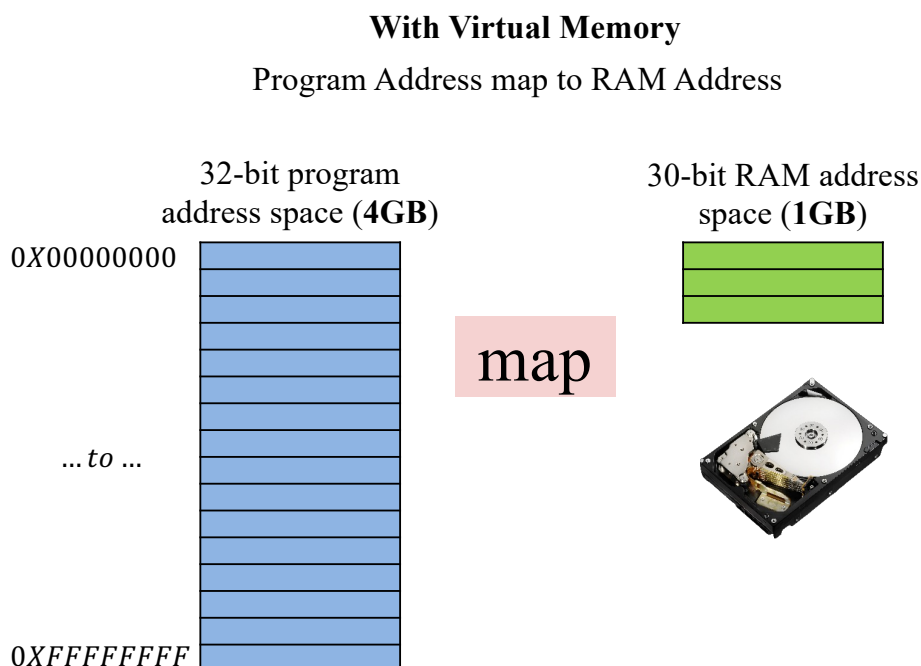
# Virtual memory is a layer of indirection

- "Any problem in Computer Science can be solved by adding indirection"
- Virtual memory takes program addresses and maps them to RAM addresses

**Without Virtual Memory**

Program Address = RAM Address

32-bit program address space (**4GB**)

30-bit RAM address space (**1GB**)

0X00000000

?

Crash if we try to access more than 1GB

... to ...

0XFFFFFFFF

**With Virtual Memory**

Program Address map to RAM Address

32-bit program address space (**4GB**)

30-bit RAM address space (**1GB**)

0X00000000

map

... to ...

0XFFFFFFFF

VM: mapping gives us flexibility in how we use the RAM

# Solving the problems: #1 not enough memory

- Map some of the program's address space to the disk
- When we need it, we bring it into memory

**With Virtual Memory**

Program Address map to RAM Address

32-bit program address space (**4GB**)

30-bit RAM address space (**1GB**)

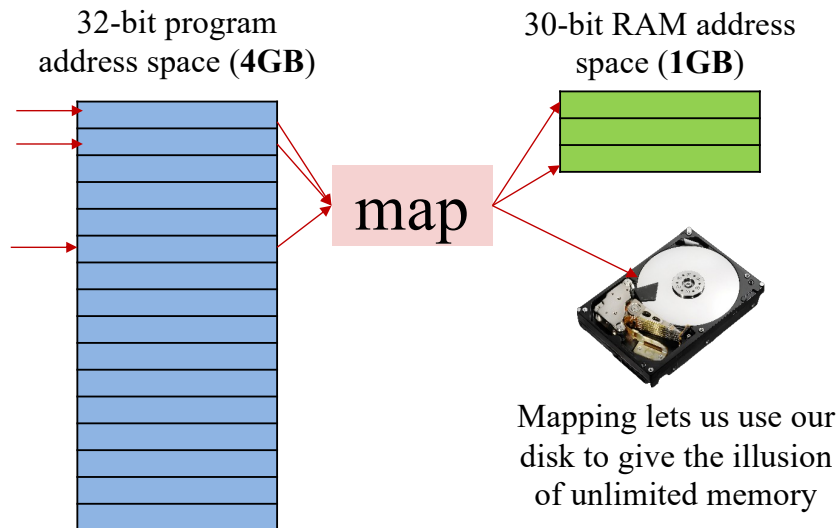0X00000000

... to ...

map

0XFFFFFFFF

# Solving the problems: #1 not enough memory

- Map some of the program's address space to the disk
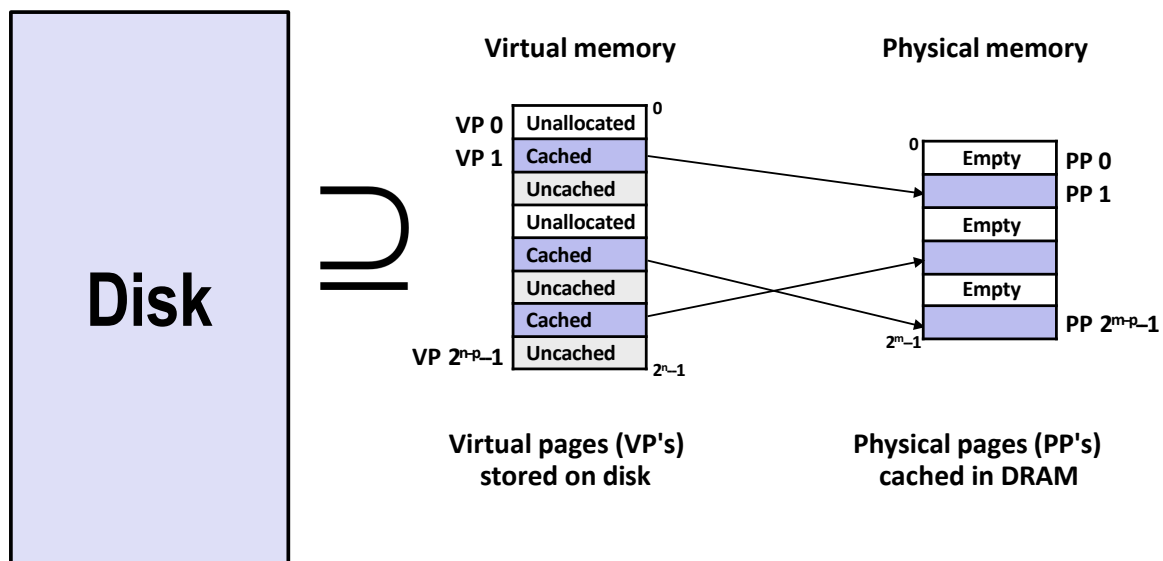
- When we need it, we bring it into memory
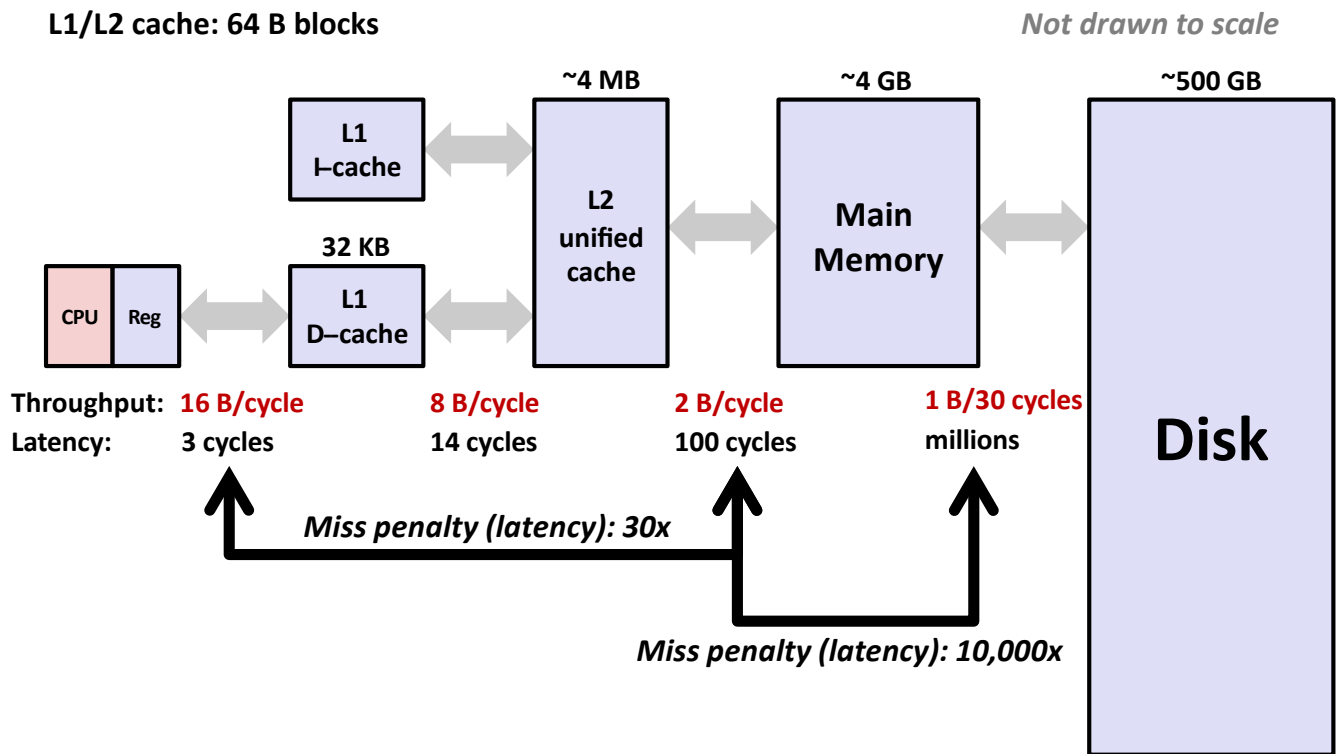
**With Virtual Memory**

Program Address map to RAM Address

32-bit program
address space (**4GB**)

30-bit RAM address
space (**1GB**)

map

Mapping lets us use our
disk to give the illusion
of unlimited memory

# VM as a Tool for Caching

- *Virtual memory:* **array of N = $2^n$ contiguous bytes**
    - think of the array (allocated part) as being stored on disk
- **Physical main memory (DRAM) = cache for allocated virtual memory**
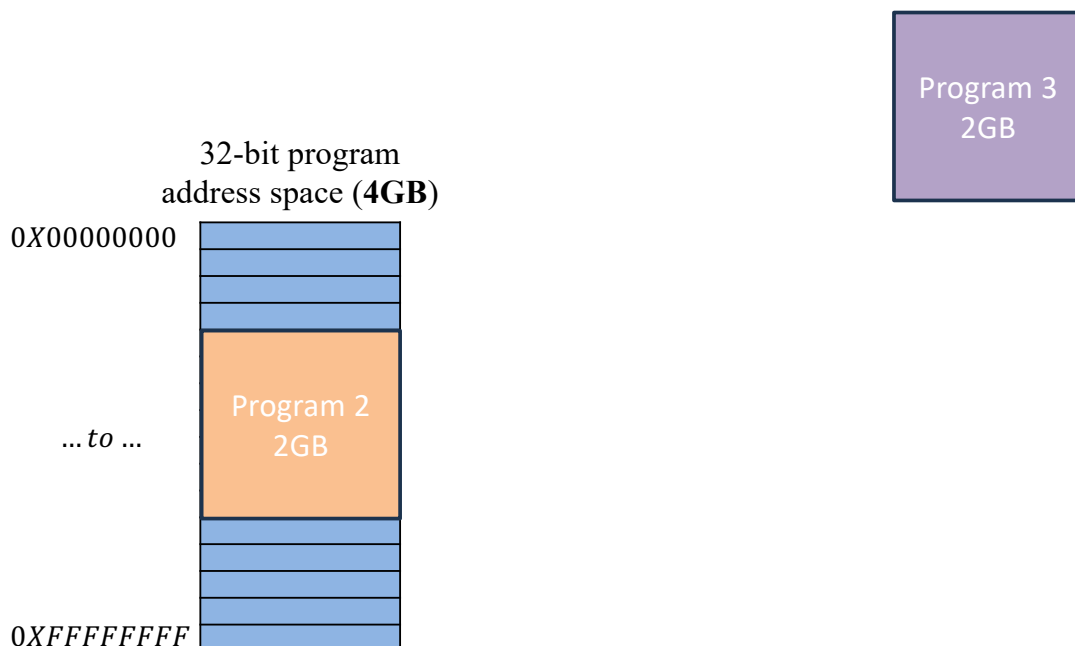- **Blocks are called pages; size $= 2^p$**

**Disk** $\supseteq$

**Virtual memory**

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

**Physical memory**

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

**Virtual pages (VP's)
stored on disk**

**Physical pages (PP's)
cached in DRAM**

# Memory Hierarchy: Core 2 Duo

**L1/L2 cache: 64 B blocks**                                                  *Not drawn to scale*

|  | ~4 MB | ~4 GB | ~500 GB |
|---|---|---|---|
| L1 I–cache | L2 unified cache | Main Memory | Disk |

**32 KB**

CPU | Reg → L1 D–cache

| | | | |
|---|---|---|---|
| **Throughput:** 16 B/cycle | 8 B/cycle | 2 B/cycle | 1 B/30 cycles |
| **Latency:** 3 cycles | 14 cycles | 100 cycles | millions |

*Miss penalty (latency): 30x*

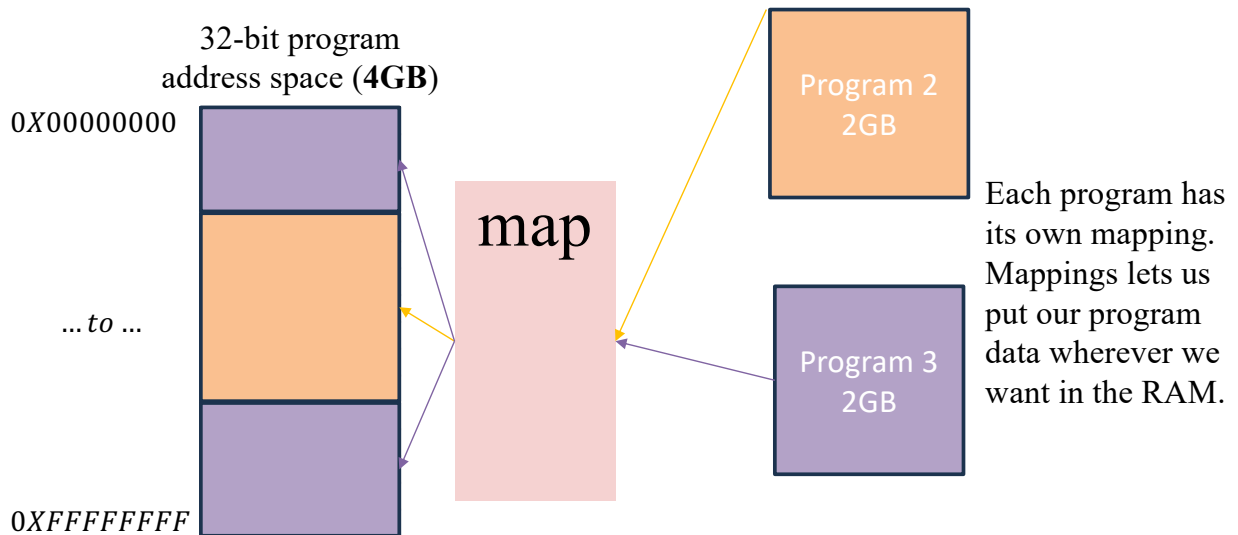*Miss penalty (latency): 10,000x*

# Solving the problems: #2 holes in address space

- How do we use the holes left when programs quit?
- We can map a program's addresses to RAM addresses however we like

Program 3
2GB

32-bit program
address space (**4GB**)

0*X*00000000

*... to ...*

Program 2
2GB

0*XFFFFFFFF*

# Solving the problems: #2 holes in address space

- How do we use the holes left when programs quit?
- We can map a program's addresses to RAM addresses however we like
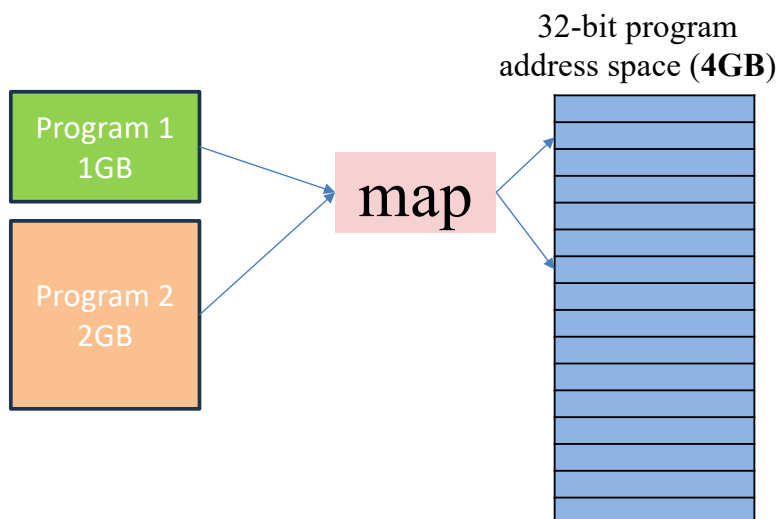
**With Virtual Memory**

Program Address map to RAM Address

32-bit program
address space (**4GB**)

$0X00000000$

... to ...

$0XFFFFFFFF$

map

Program 2
2GB

Program 3
2GB

Each program has its own mapping. Mappings lets us put our program data wherever we want in the RAM.

# Solving the problems: #3 keeping program secure

- Program 1's and Program 2's addresses map to different RAM addresses

    sw $t0, 1024($0)

**With Virtual Memory**

Program Address map to RAM Address

32-bit program
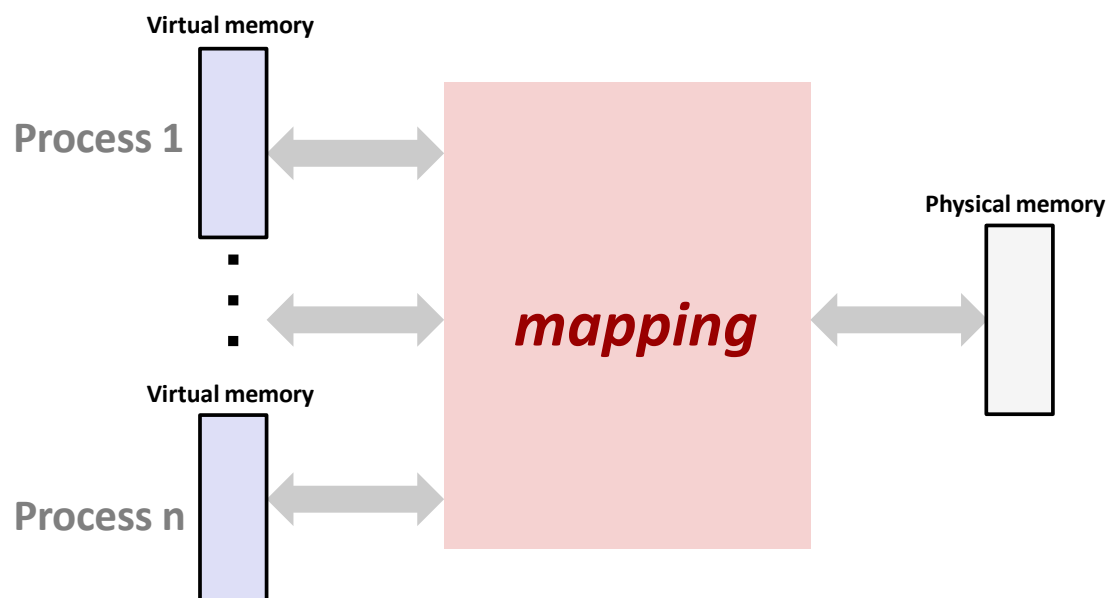address space (**4GB**)

Program 1
1GB

Program 2
2GB

map

# Solving the problems: #3 keeping program secure

- Because different processes will have different mappings from virtual to physical addresses, two programs can freely use the same virtual address.
- By allocating distinct regions of physical memory to A and B, they are prevented from reading/writing each others data.
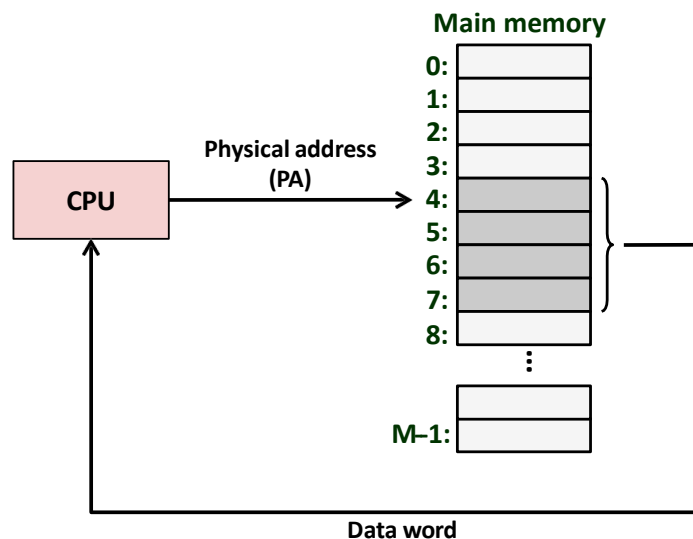


# Solution: Level Of Indirection



- **Each process gets its own private memory space**
- **Solves the previous problems**
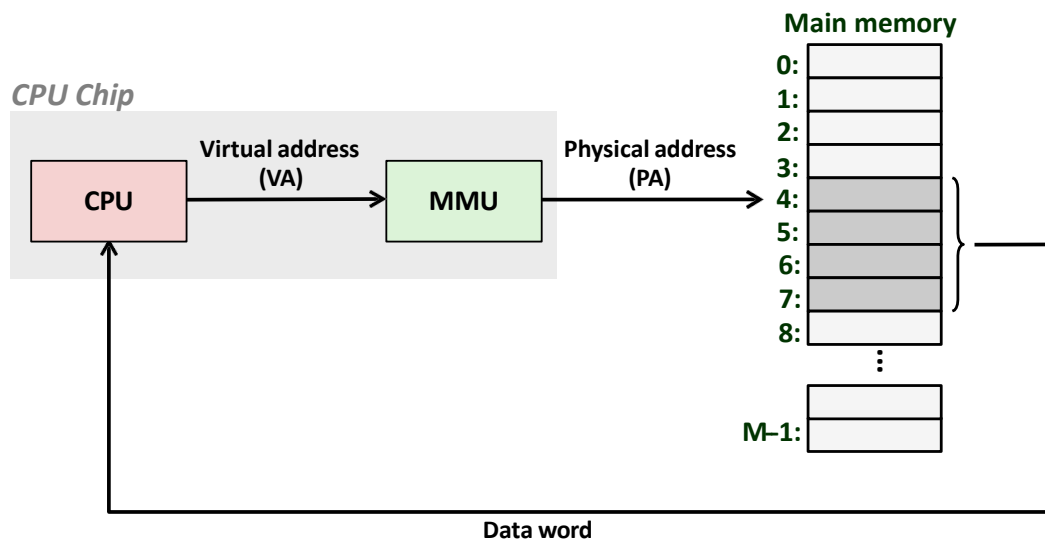
# How does VM work?

- Basic idea: separate memory spaces
  - Virtual memory: what the program sees
  - Physical memory: the physical RAM in the computer

- Virtual Addresses (VA):
  - What the program uses
  - In MIPS, this is the full 32-bit address space: 0 to $2^{32} - 1$
- Physical Address (PA):
  - What the hardware uses to talk to the RAM
  - Address space determined by how much RAM is installed

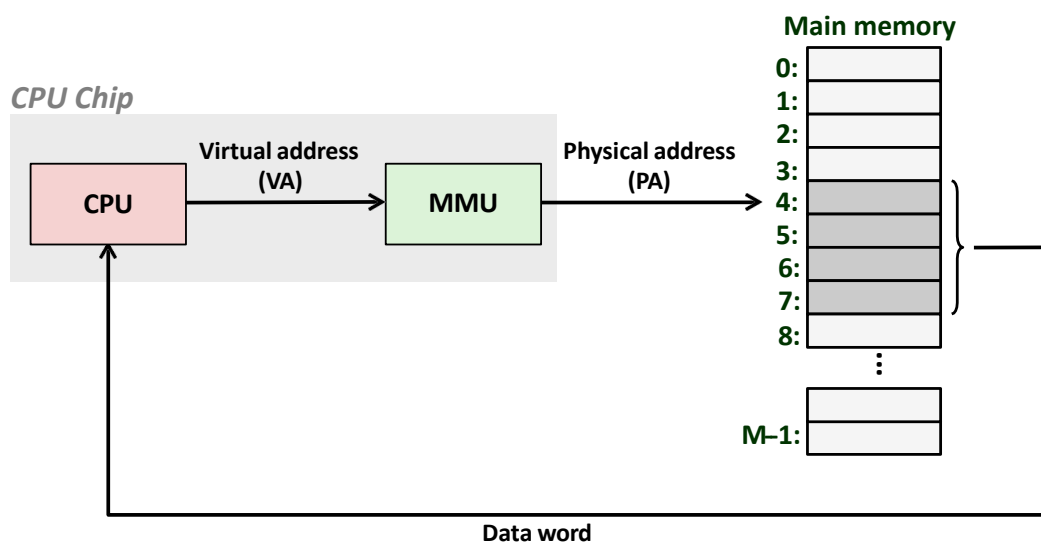# A System Using Physical Addressing

**Main memory**



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 desktops, laptops, workstations
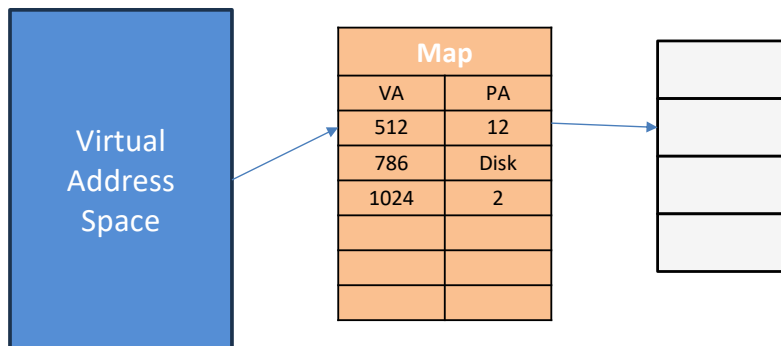- One of the great ideas in computer science

# A System Using Virtual Addressing



*How would you do the VA -> PA translation?*
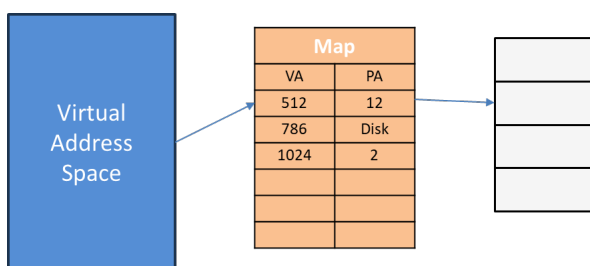
# Page tables

- The map from VA to PA is the page table
- So far we have had one Page Table Entry (PTE) for every Virtual Address
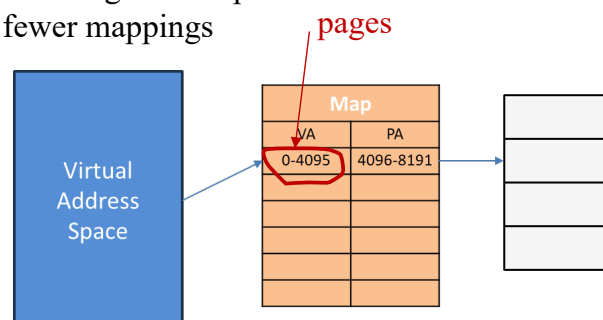- How many entries do we need in our page table?



# Page tables size

- We need to translate every possible address:
- One program has 32-bit Virtual Memory spaces
  - That's $2^{30}$ words that need Page Table Entries (1 billion entries!)
  - If they don't have a Page Table Entry then we can't access them because we can't find the physical address.
- How can we make this more manageable?
  - What if we divided memory up into chunks (pages) instead of words?
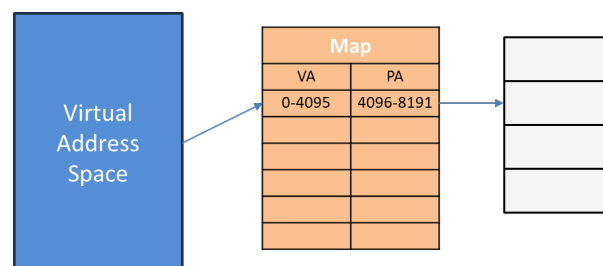
Fine-grain: Maps each word address: $2^{30}$ words

Coarse-grain: Maps chunks of address: fewer mappings

pages

# Coarse-grained: pages instead of words

- The page table manages larger chunk (page) of data
  - Fewer Page Table Entries needed to cover the whole address space
  - But, less flexibility in how to use the RAM (have to move a page at a time)
- Today:
  - Typically, 4kB pages (1024 words per page)
    - Question: How many entries do we need in our Page Table with 4kB pages on a 32-bit machine?
  - Sometimes 2MB pages
    - Question: How many entries do we need in our Page Table with 4kB pages on a 32-bit machine?
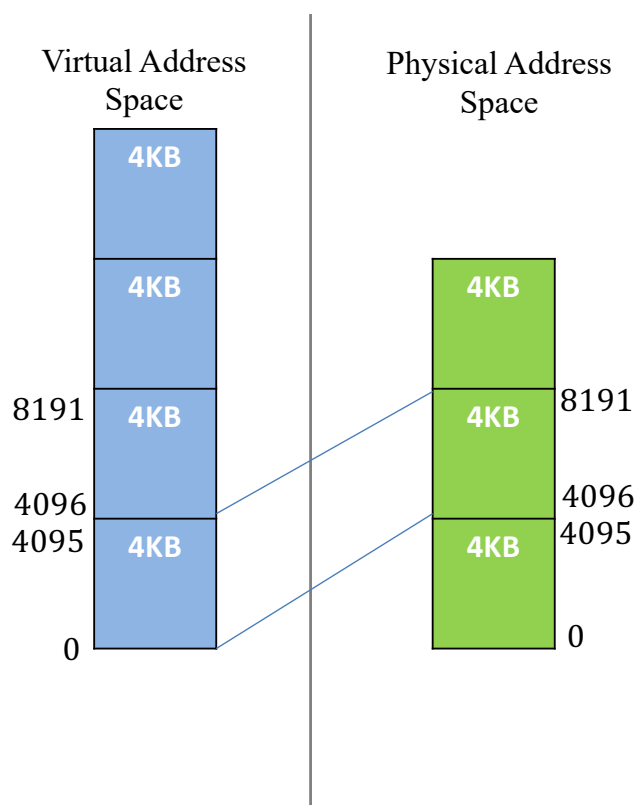


| Map | |
|---|---|
| VA | PA |
| 0-4095 | 4096-8191 |
| | |
| | |
| | |
| | |

Virtual Address Space

---

# How do we map address with pages?

Coarse-grain:
Maps chunks of
address: fewer
mappings

| Map | |
|---|---|
| VA | PA |
| 0-4095 | 4096-8191 |
| | |
| | |
| | |
| | |

Question: What is the
Physical Address for
Virtual Address 4?



Virtual Address Space

4KB
4KB
8191    4KB
4096
4095    4KB
0

Physical Address Space
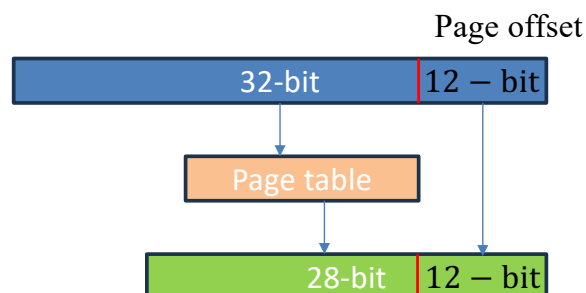
4KB
8191    4KB
4096
4095    4KB
0

# Address translation

- What happens on a 32-bit machine with 256MB RAM and 4kB pages?
  - 32-bit Virtual Address
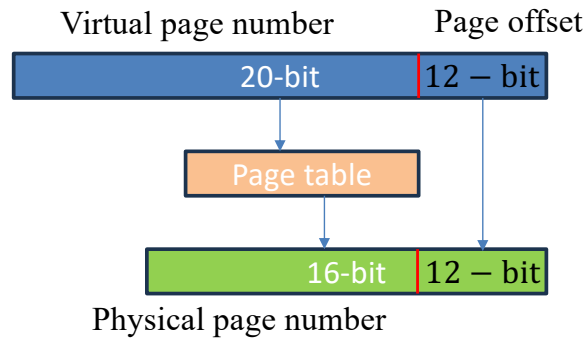  - 28-bit Physical Address



# Address translation

- What happens on a 32-bit machine with 256MB RAM and 4kB pages?
  - 32-bit Virtual Address
  - 28-bit Physical Address

Page offset



Each PTE contains 4kB of address space…
For every page we have 4096 addresses (12 bits)
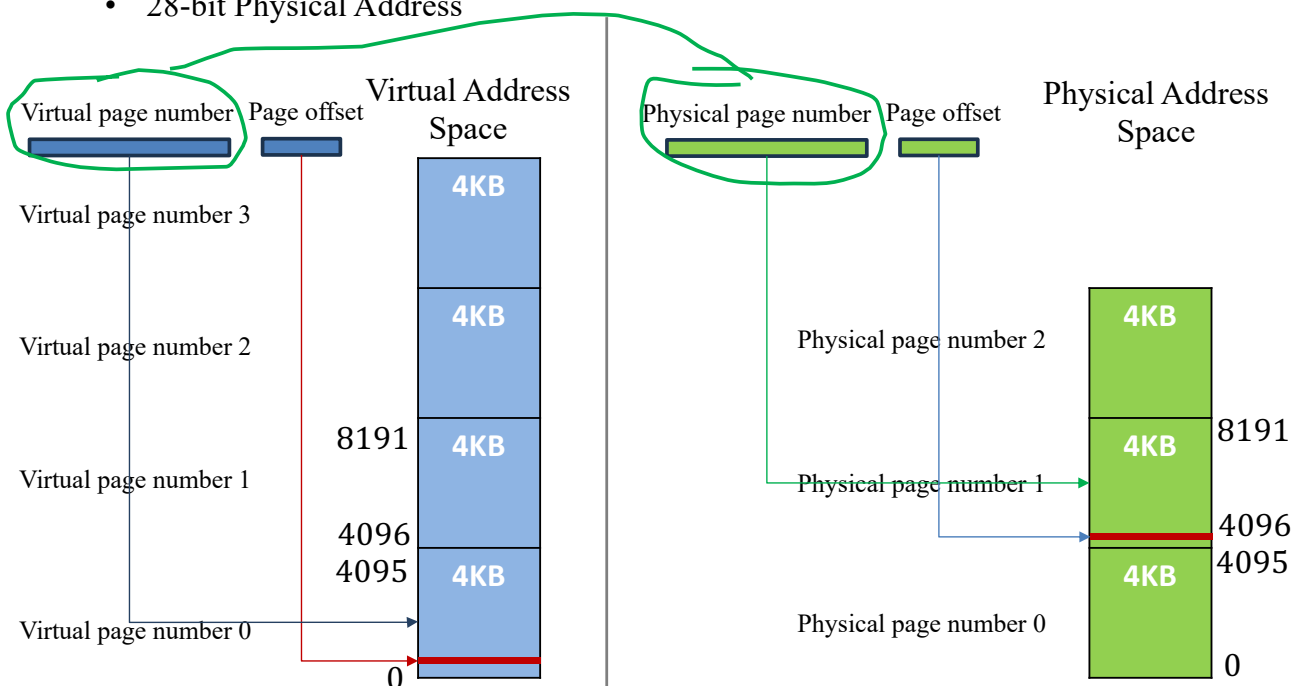that don't get translated

# Address translation

- What happens on a 32-bit machine with 256MB RAM and 4kB pages?
  - 32-bit Virtual Address
  - 28-bit Physical Address

Virtual page number      Page offset

| 20-bit | 12 − bit |
|--------|----------|

Page table

| 16-bit | 12 − bit |
|--------|----------|

Physical page number

But each PTE contains 4kB of address space…
For every page we have 4096 addresses (12 bits)
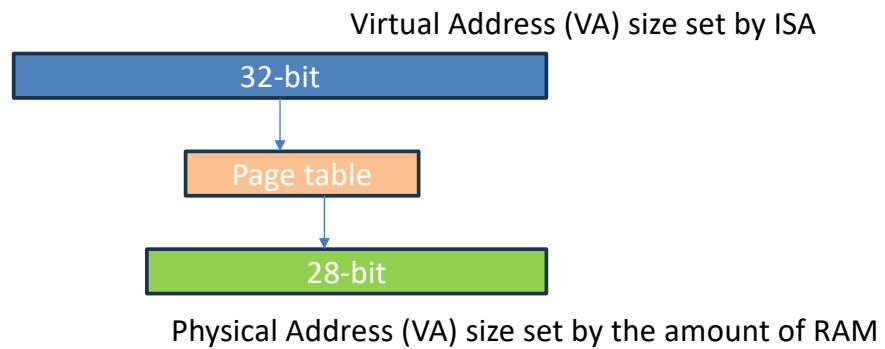that don't get translated

---

# Pages, offsets, and translation

- What happens on a 32-bit machine with 256MB RAM and 4kB pages?
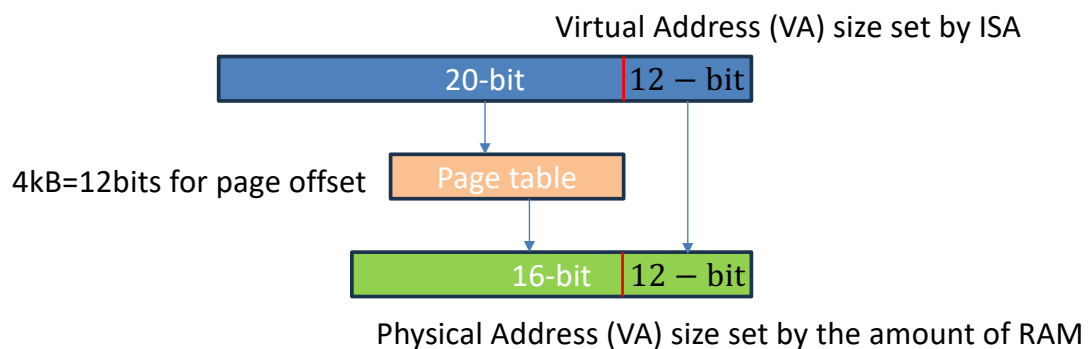  - 32-bit Virtual Address
  - 28-bit Physical Address

Virtual page number   Page offset   Virtual Address Space

Virtual page number 3

Virtual page number 2

8191

Virtual page number 1

4096
4095

Virtual page number 0

0

4KB
4KB
4KB
4KB

Physical page number   Page offset   Physical Address Space

Physical page number 2

Physical page number 1

Physical page number 0

8191

4096
4095

0

4KB
4KB
4KB

# How to do a page table lookup?

- What happens on a 32-bit machine with 256MB RAM and 4kB pages?
  - 32-bit Virtual Address
  - 28-bit Physical Address

Virtual Address (VA) size set by ISA

| 32-bit |
| --- |

Page table

| 28-bit |
| --- |

Physical Address (VA) size set by the amount of RAM

---

# How to do a page table lookup?

- What happens on a 32-bit machine with 256MB RAM and 4kB pages?
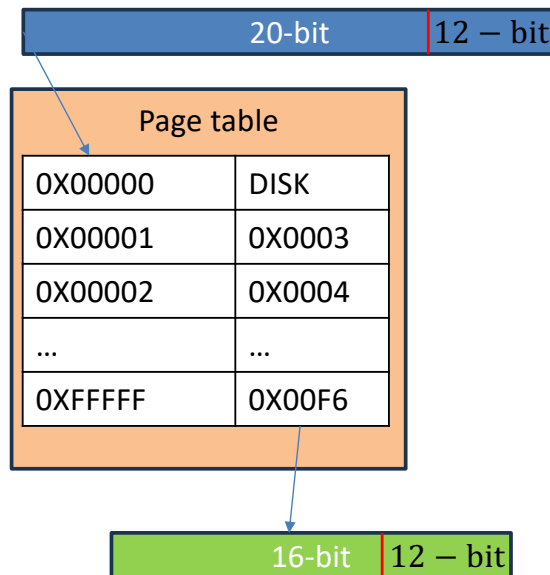  - 32-bit Virtual Address
  - 28-bit Physical Address

Virtual Address (VA) size set by ISA

| 20-bit | 12 − bit |
| --- | --- |

4kB=12bits for page offset          Page table

| 16-bit | 12 − bit |
| --- | --- |

Physical Address (VA) size set by the amount of RAM

# How to do a page table lookup?

- What happens on a 32-bit machine with 256MB RAM and 4kB pages?
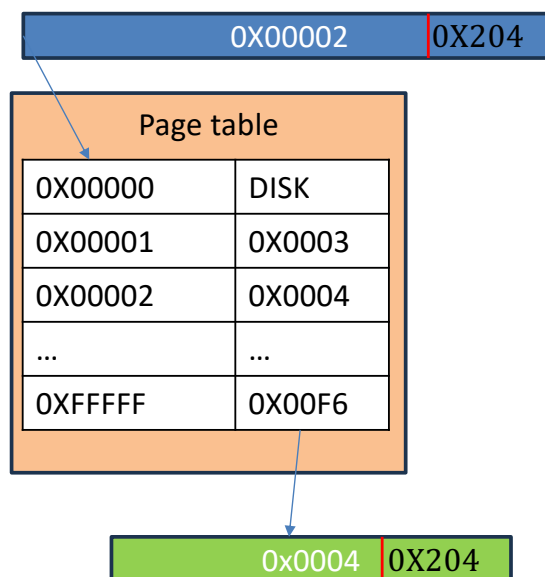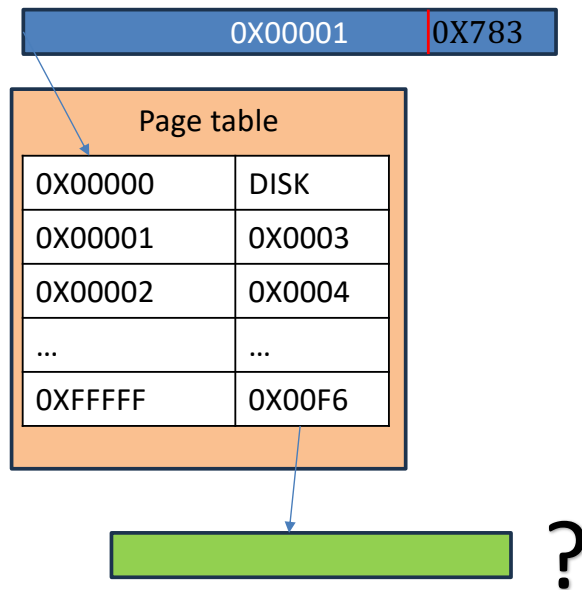  - 32-bit Virtual Address
  - 28-bit Physical Address

| 20-bit | 12 − bit |
|--------|----------|

### Page table

| 0X00000 | DISK   |
|---------|--------|
| 0X00001 | 0X0003 |
| 0X00002 | 0X0004 |
| ...     | ...    |
| 0XFFFFF | 0X00F6 |

| 16-bit | 12 − bit |
|--------|----------|

# Example Translation

- What happens on a 32-bit machine with 256MB RAM and 4kB pages?
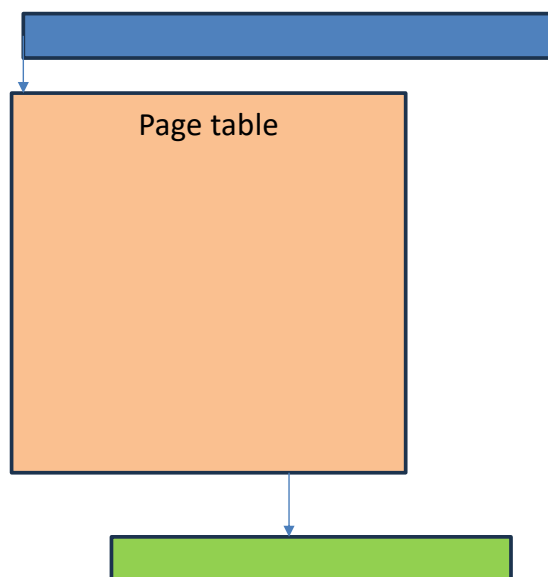  - 32-bit Virtual Address
  - 28-bit Physical Address

| 0X00002 | 0X204 |
|---------|-------|

### Page table

| 0X00000 | DISK   |
|---------|--------|
| 0X00001 | 0X0003 |
| 0X00002 | 0X0004 |
| ...     | ...    |
| 0XFFFFF | 0X00F6 |

| 0x0004 | 0X204 |
|--------|-------|

# Question

- What happens on a 32-bit machine with 256MB RAM and 4kB pages?
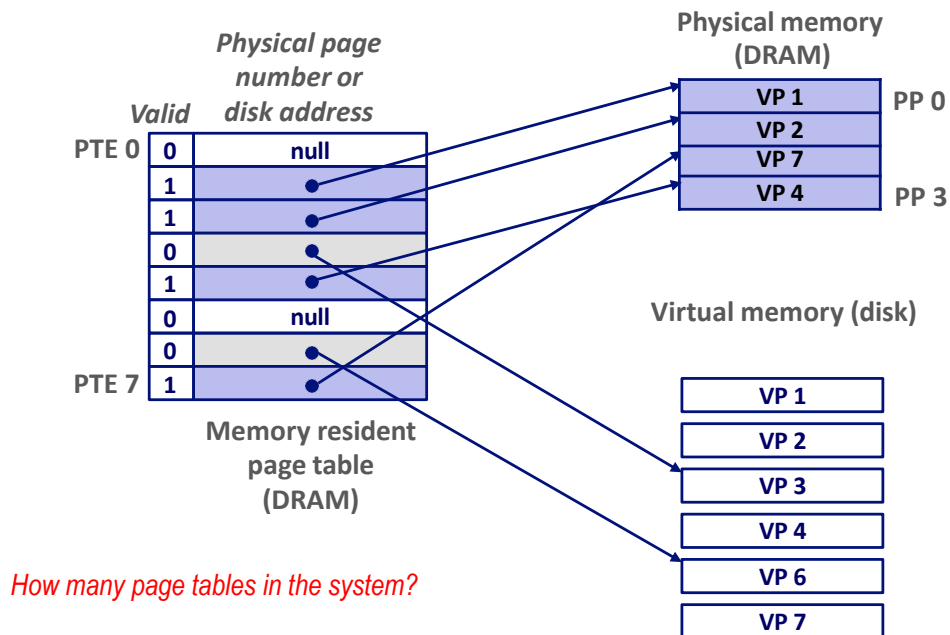    - 32-bit Virtual Address
    - 28-bit Physical Address

| 0X00001 | 0X783 |
|---------|-------|

**Page table**

| 0X00000 | DISK |
|---------|------|
| 0X00001 | 0X0003 |
| 0X00002 | 0X0004 |
| ... | ... |
| 0XFFFFF | 0X00F6 |

**?**

# Question

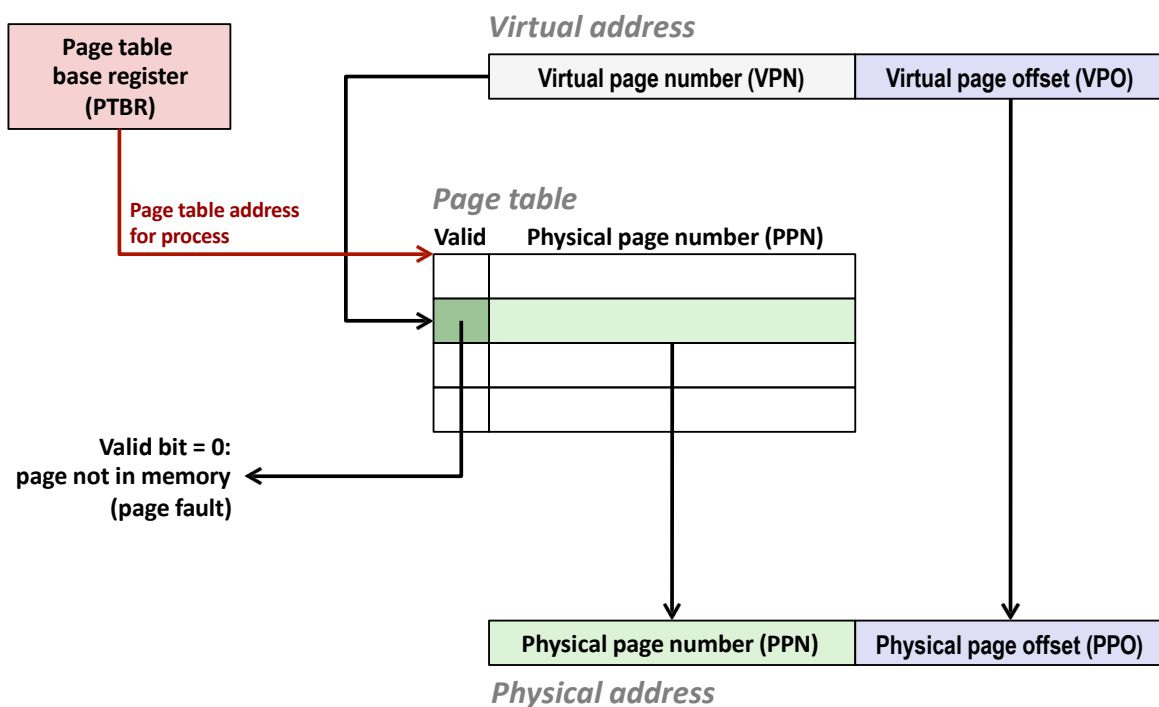- If we have 64kB pages, how many bits do we use for the page offset?

**Page table**

# Address Translation: Page Tables

A *page table* is an array of page table entries (PTEs) that maps virtual pages to physical pages. Here: 8 VPs
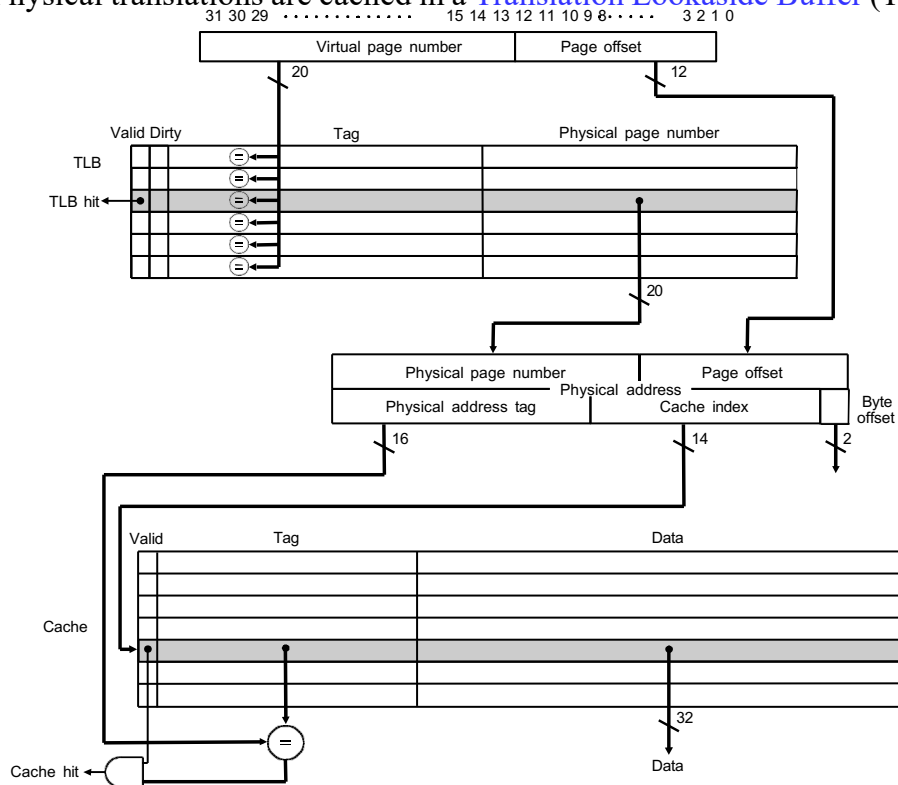


*How many page tables in the system?*

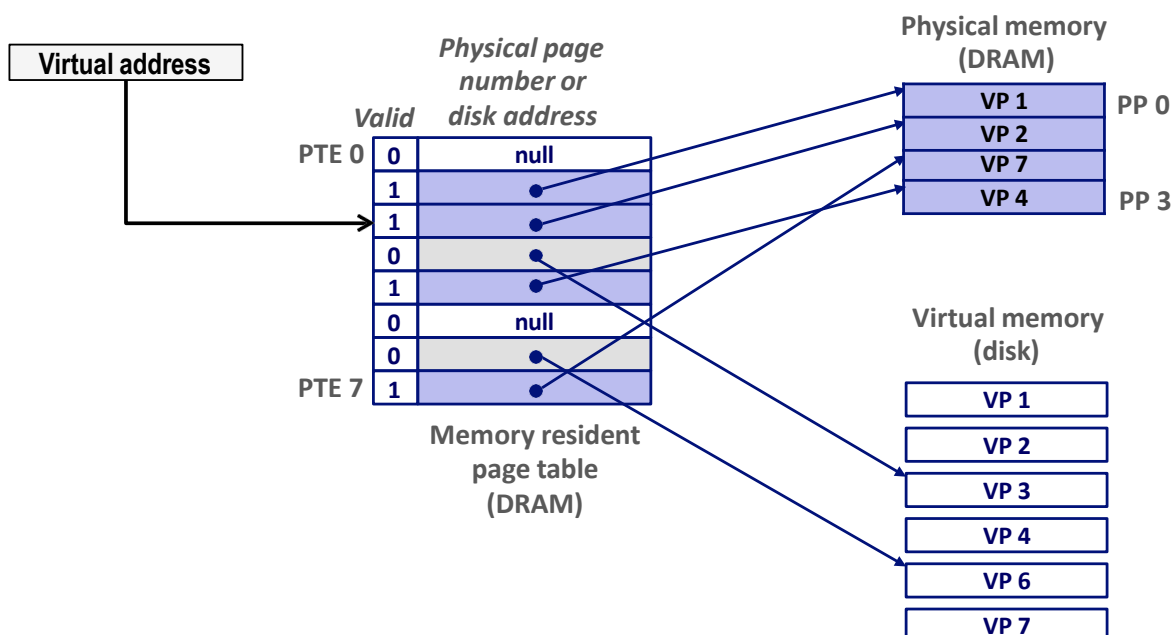# Address Translation With a Page Table

# Caching Translations

- Virtual to Physical translations are cached in a Translation Lookaside Buffer (TLB).



# Page Hit

*Page hit:* **reference to VM word that is in physical memory**

# Page Miss

*Page miss:* reference to VM word that is not in physical memory (shOOt!)