# CS 61C:
# Great Ideas in Computer Architecture
# *Virtual Memory*

Instructors:

John Wawrzynek & Vladimir Stojanovic

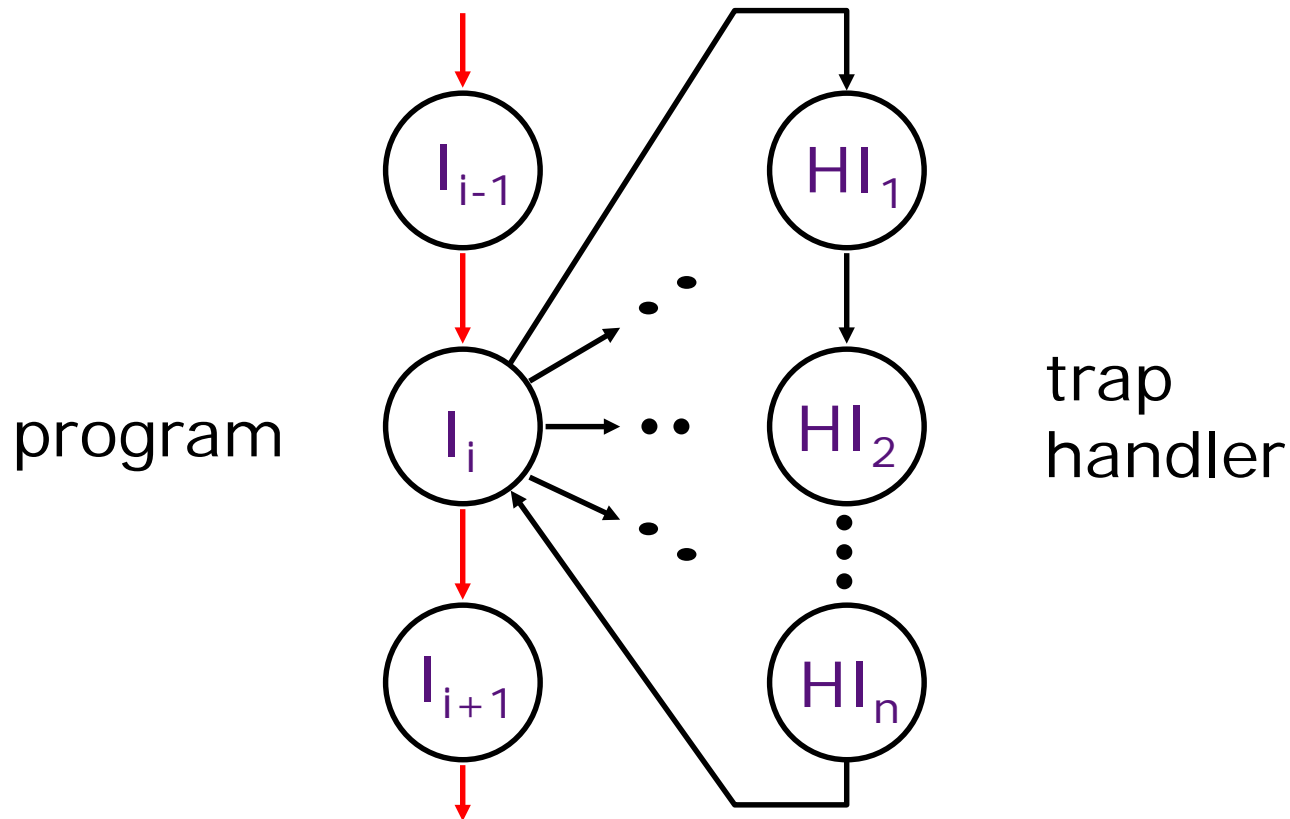http://inst.eecs.berkeley.edu/~cs61c/

# Review

- Programmed I/O
- Polling vs. Interrupts
- Booting a Computer
  - BIOS, Bootloader, OS Boot, Init
- Supervisor Mode, Syscalls
- Base and Bounds
  - Simple, but doesn't give us everything we want
- Intro to VM

# Traps/Interrupts/Execeptions:
## altering the normal flow of control

program

$I_{i-1}$

$I_i$

$I_{i+1}$

$HI_1$

$HI_2$

$HI_n$

trap handler

An *external or internal event* that needs to be processed - by another program – the OS. The event is often unexpected from original program's point of view.
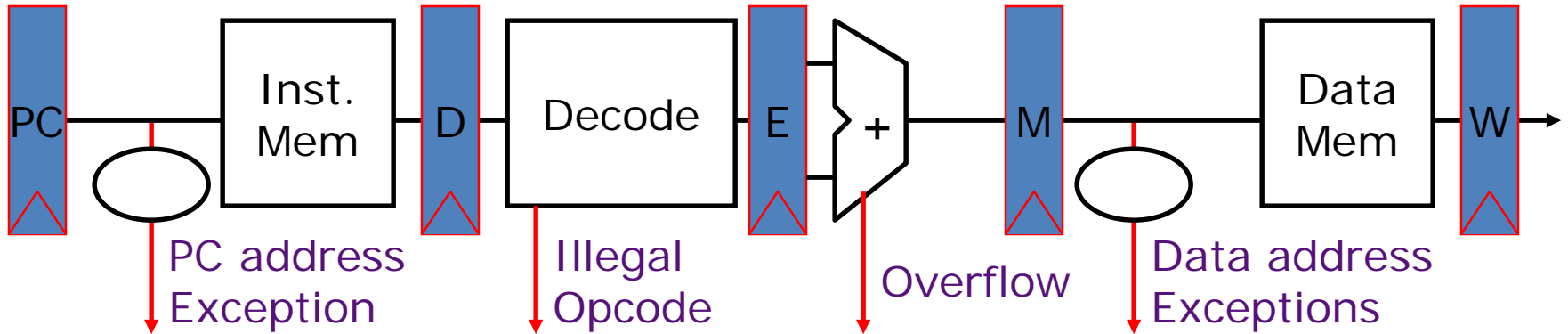
# Terminology

In CS61C (you'll see other definitions in use elsewhere):

- <u>Interrupt</u> – caused by an event *external* to current running program (e.g. key press, mouse activity)
    - Asynchronous to current program, can handle interrupt on any convenient instruction
- <u>Exception</u> – caused by some event during execution of one instruction of current running program (e.g., bus error, illegal instruction)
    - Synchronous, must handle exception on instruction that causes exception
- <u>Trap</u> – action of servicing interrupt or exception by hardware jump to "trap handler" code

# Precise Traps

- *Trap handler's view of machine state is that every instruction prior to the trapped one has completed, and no instruction after the trap has executed.*
- Implies that handler can return from an interrupt by restoring user registers and jumping back to interrupted instruction (EPC register will hold the instruction address)
  - Interrupt handler software doesn't need to understand the pipeline of the machine, or what program was doing!
  - More complex to handle trap caused by an exception than interrupt
- Providing precise traps is tricky in a pipelined superscalar out-of-order processor!
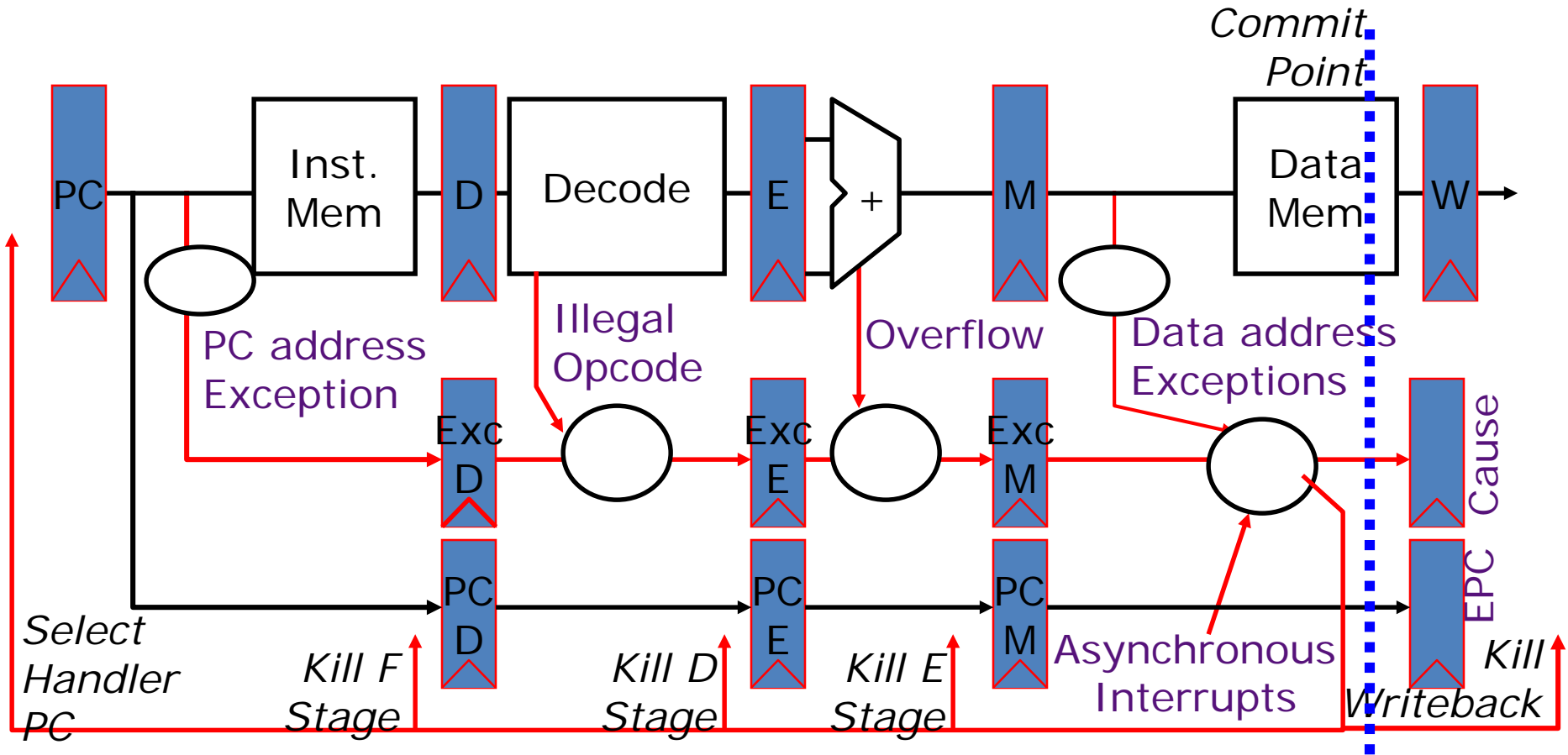  - But handling imprecise interrupts in software is even worse.

# Trap Handling in 5-Stage Pipeline



PC address Exception

Illegal Opcode

Overflow

Data address Exceptions

→ Asynchronous Interrupts

- How to handle multiple simultaneous exceptions in different pipeline stages?
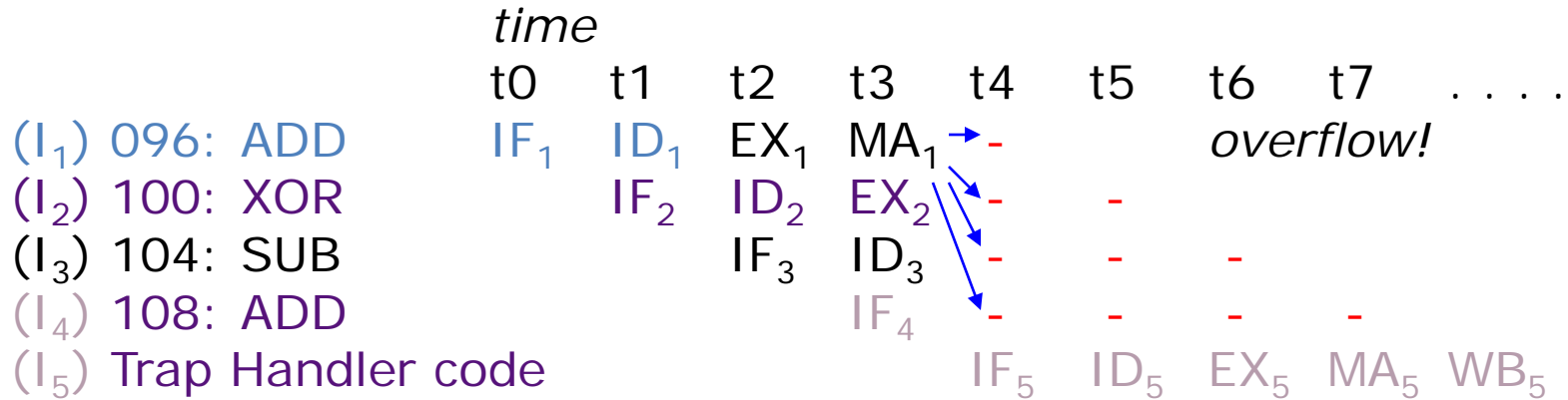- How and where to handle external asynchronous interrupts?

# Save Exceptions Until Commit

# Handling Traps in In-Order Pipeline

- Hold exception flags in pipeline until commit point (M stage)

- Exceptions in earlier instructions override exceptions in later instructions

- Exceptions in earlier pipe stages override later exceptions *for a given instruction*

- Inject external interrupts at commit point (override others)

- If exception/interrupt at commit: update Cause and EPC registers, kill all stages, inject handler PC into fetch stage
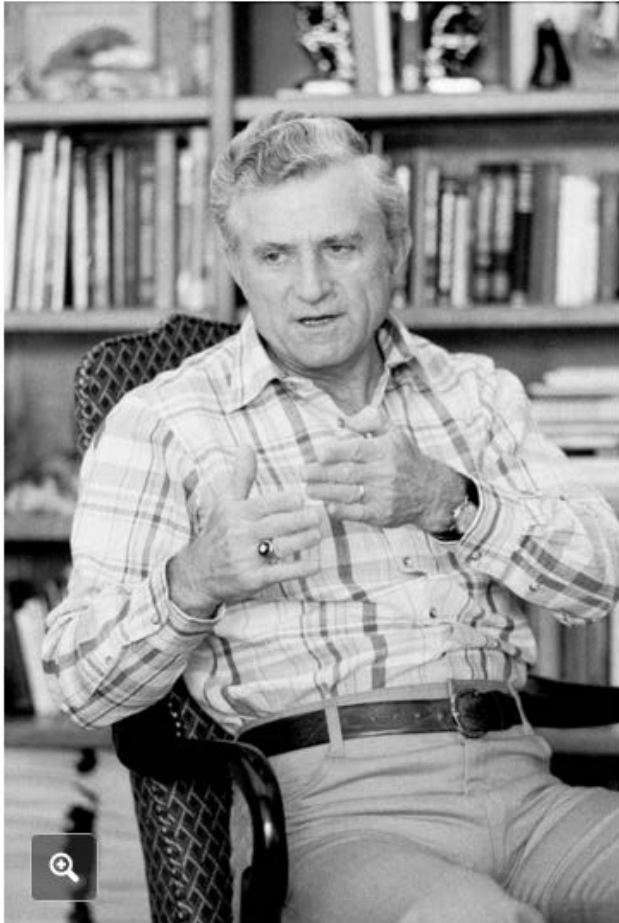
# Trap Pipeline Diagram

time

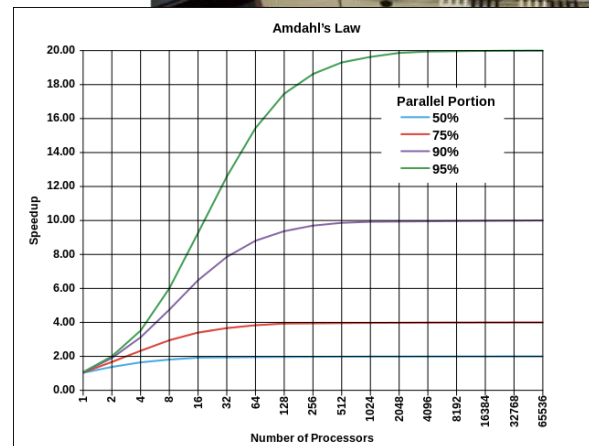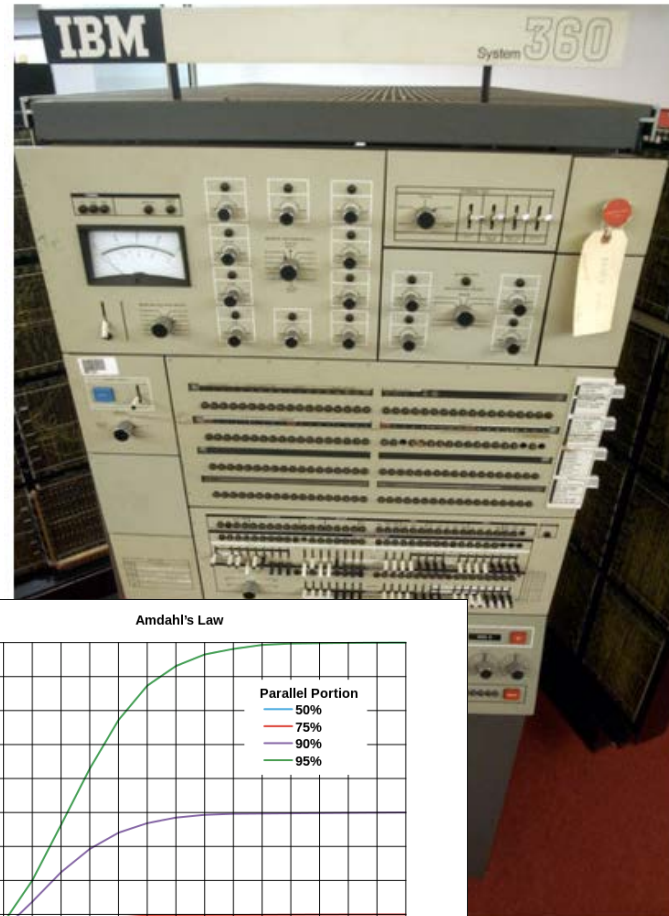|  | t0 | t1 | t2 | t3 | t4 | t5 | t6 | t7 | . . . . |
|---|---|---|---|---|---|---|---|---|---|
| $(I_1)$ 096: ADD | $IF_1$ | $ID_1$ | $EX_1$ | $MA_1$ | - | *overflow!* | | | |
| $(I_2)$ 100: XOR | | $IF_2$ | $ID_2$ | $EX_2$ | - | - | | | |
| $(I_3)$ 104: SUB | | | $IF_3$ | $ID_3$ | - | - | - | | |
| $(I_4)$ 108: ADD | | | | $IF_4$ | - | - | - | - | |
| $(I_5)$ Trap Handler code | | | | | $IF_5$ | $ID_5$ | $EX_5$ | $MA_5$ | $WB_5$ |

# In the News …

# Gene Amdahl, Pioneer of Mainframe Computing, Dies at 92
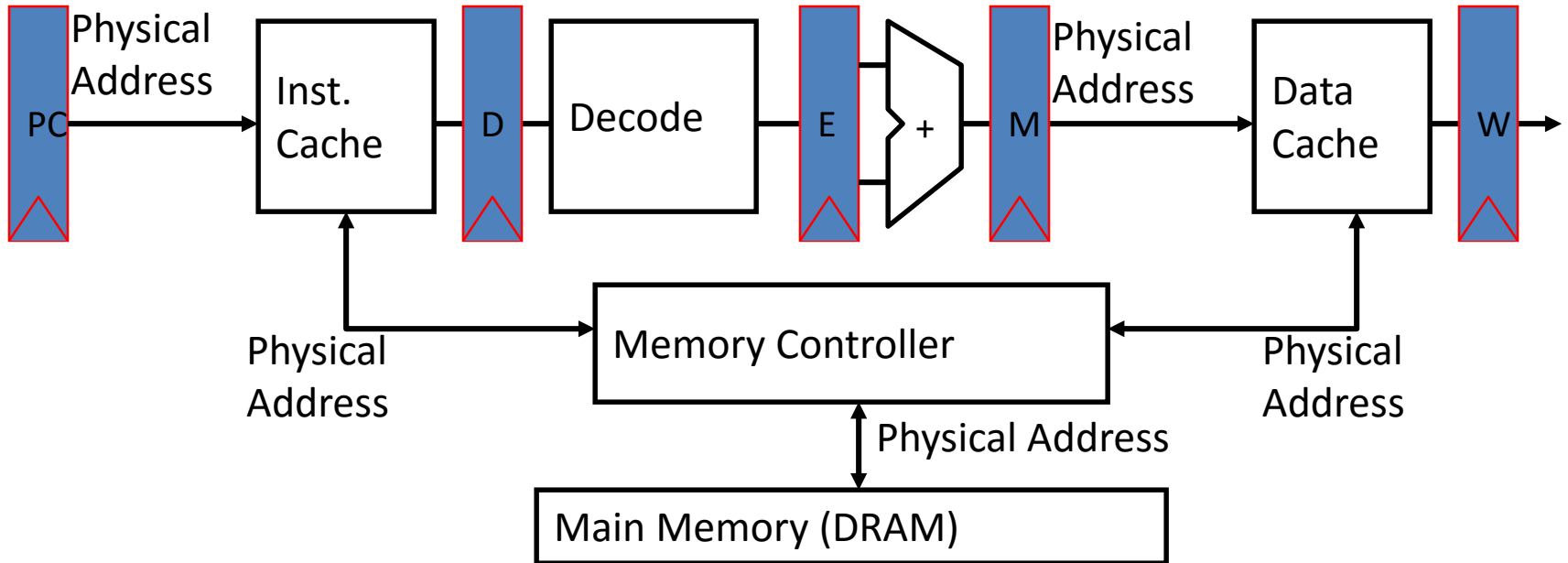
By **KATIE HAFNER**  NOV. 12, 2015



Gene Amdahl at home in Saratoga, Calif., in 1979. As a young computer scientist at IBM, Dr. Amdahl, who has a formula named after him, played a crucial role in the development of the System/360 series. Mickey Pfleger



The System/360 series, the most successful line of mainframe computers in IBM's history.
Jim McKnight/Associated Press.



10

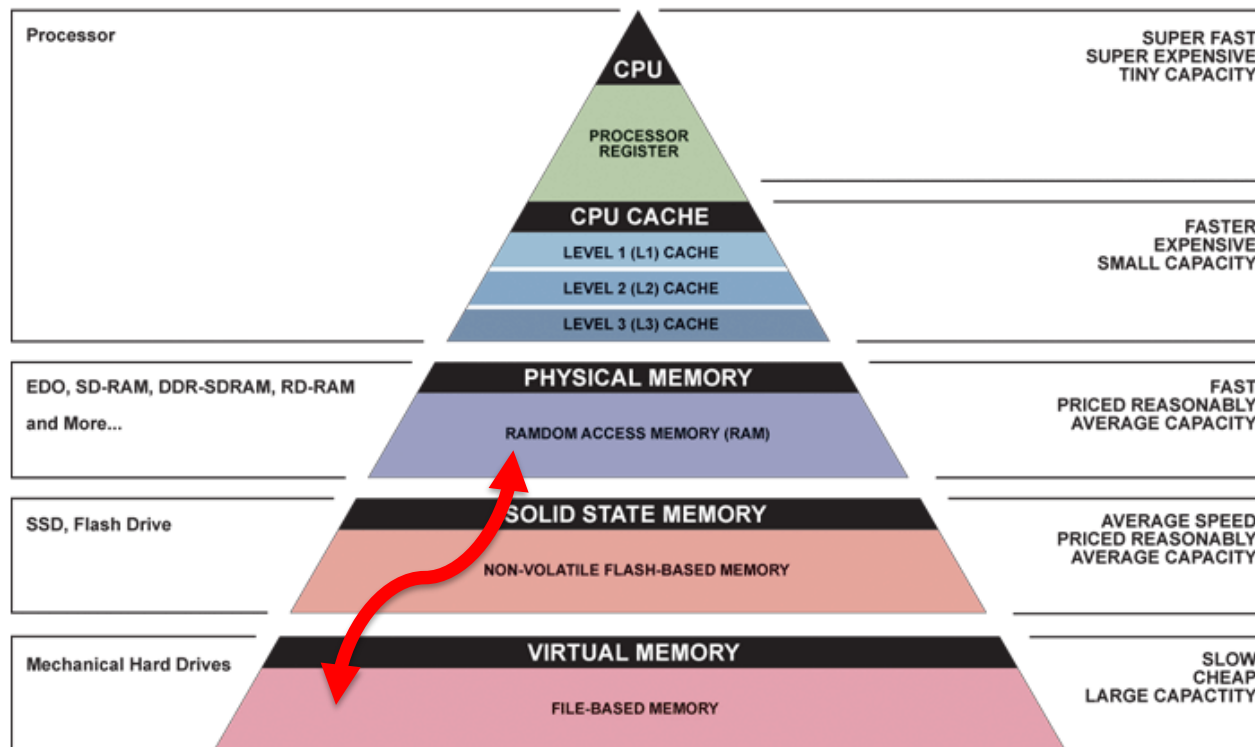# Virtual Memory

# "Bare" 5-Stage Pipeline



- In a bare machine, the only kind of address is a physical address

# What do we need Virtual Memory for? Reason 1: Adding Disks to Hierarchy

- Need to devise a mechanism to "connect" memory and disk in the memory hierarchy



▲ Simplified Computer Memory Hierarchy
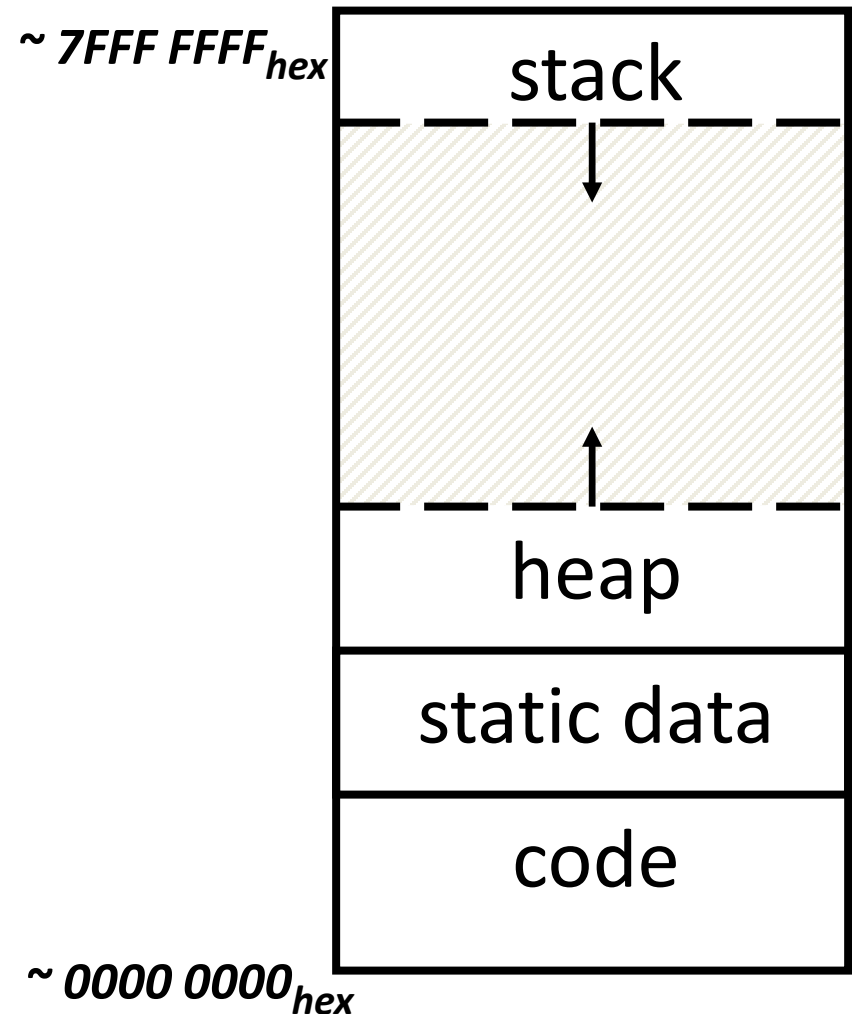Illustration: Ryan J. Leng

# What do we need Virtual Memory for? Reason 2: Simplifying Memory for Apps

- Applications should see the straightforward memory layout we saw earlier ->

- User-space applications should think they own all of memory

- So we give them a **virtual** view of memory

*~ 7FFF FFFF$_{hex}$*

| |
|---|
| stack |
| heap |
| static data |
| code |

*~ 0000 0000$_{hex}$*

# What do we need Virtual Memory for? Reason 3: Protection Between Processes

- With a bare system, addresses issued with loads/stores are real **physical** addresses

- This means any program can issue any address, therefore can access any part of memory, even areas which it doesn't own
  - Ex: The OS data structures

- We should send all addresses through a mechanism that the OS controls, before they make it out to DRAM - **a translation mechanism**
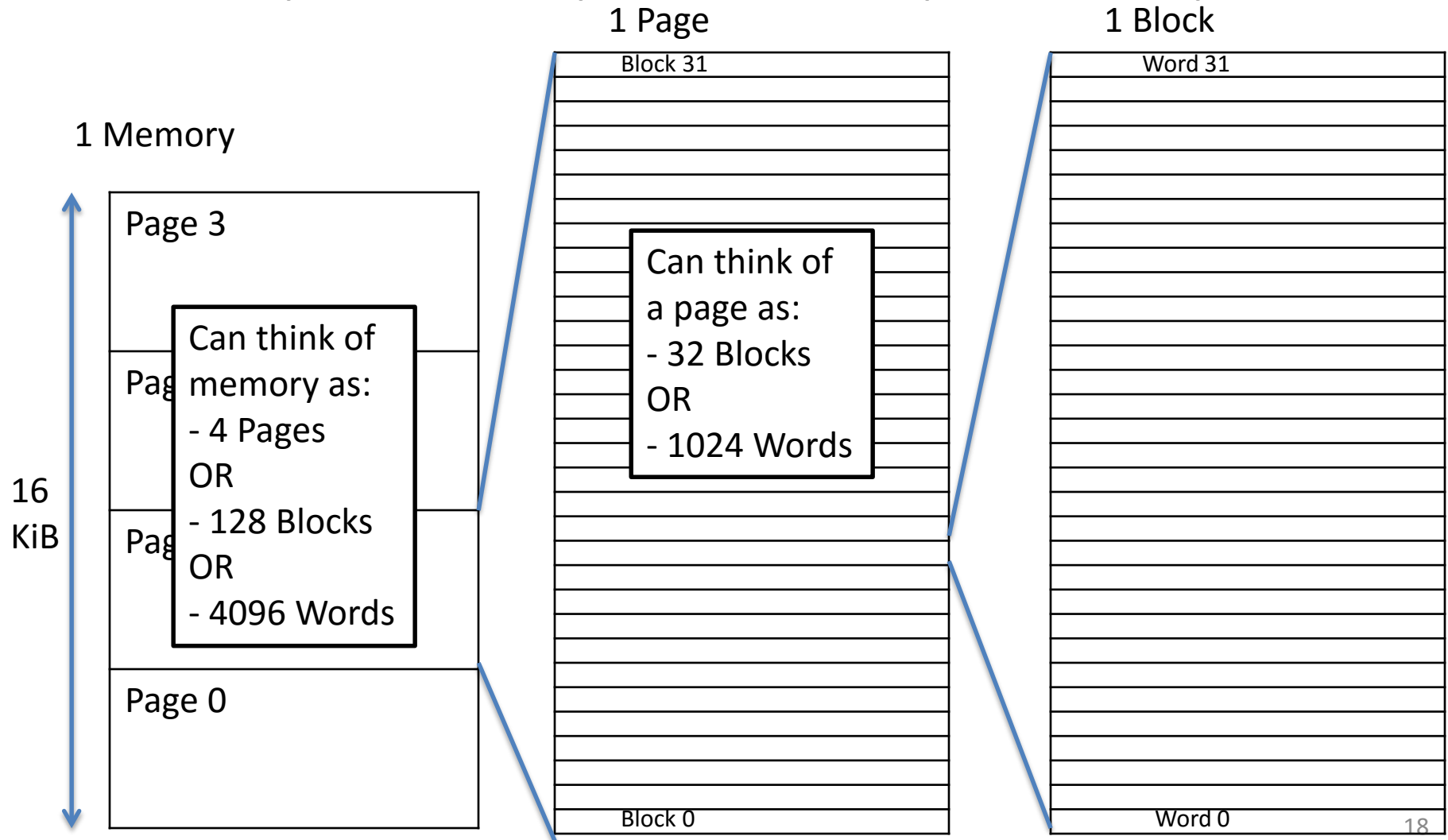
# Address Spaces

- The set of addresses labeling all of memory that we can access
- Now, 2 kinds:
  - **Virtual Address Space** - the set of addresses that the user program knows about
  - **Physical Address Space** - the set of addresses that map to actual physical cells in memory
    - Hidden from user applications
- So, we need a way to map between these two address spaces

# Blocks vs. Pages

- In caches, we dealt with individual *blocks*
  - Usually ~64B on modern systems
  - We could "divide" memory into a set of blocks
- In VM, we deal with individual *pages*
  - Usually ~4 KB on modern systems
  - Now, we'll "divide" memory into a set of pages
- Common point of confusion: Bytes, Words, Blocks, Pages are all just different ways of looking at memory!

# Bytes, Words, Blocks, Pages

Ex: 16 KiB DRAM, 4 KiB Pages (for VM), 128 B blocks (for caches), 4 B words (for lw/sw)

1 Memory

16 KiB

Page 3

Page 2

Page 1

Page 0

Can think of memory as:
- 4 Pages
OR
- 128 Blocks
OR
- 4096 Words

1 Page

Block 31

Block 0

Can think of a page as:
- 32 Blocks
OR
- 1024 Words

1 Block

Word 31

Word 0

# Address Translation

- So, what do we want to achieve at the hardware level?
  - Take a Virtual Address, that points to a spot in the Virtual Address Space of a particular program, and map it to a Physical Address, which points to a physical spot in DRAM of the whole machine

| Virtual Address | Virtual Page Number | Offset |
|---|---|---|

| Physical Address | Physical Page Number | Offset |
|---|---|---|

# Address Translation

Virtual Address

| Virtual Page Number | | Offset |
|---|---|---|

Address Translation

Copy Bits

Physical Address

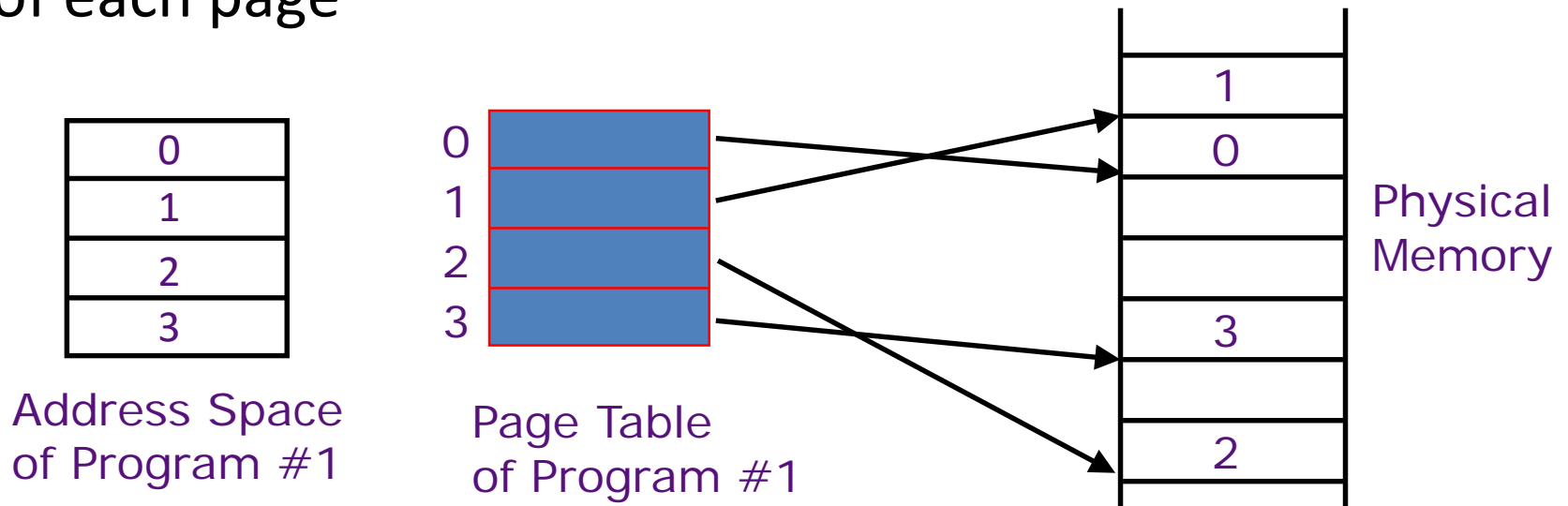| Physical Page Number | | Offset |
|---|---|---|

The rest of the lecture is all about implementing

# Paged Memory Systems

- Processor-generated address can be split into:

| Virtual Page Number | Offset |
|---|---|

- A *page table* contains the physical address of the base of each page



Address Space of Program #1

Page Table of Program #1

Physical Memory

*Page tables make it possible to store the pages of a program non-contiguously.*

# Private (Virtual) Address Space per Program

Prog 1

VA1

Page Table

Prog 2

VA1

Page Table

Prog 3

VA1

Page Table

OS pages

• • •

free

Physical Memory

- Each prog has a page table
- Page table contains an entry for each prog page
- Physical Memory acts like a "cache" of pages for currently running programs.  **Not recently used pages are stored in secondary memory, e.g. disk (in "swap partition")**

22

# Where Should Page Tables Reside?
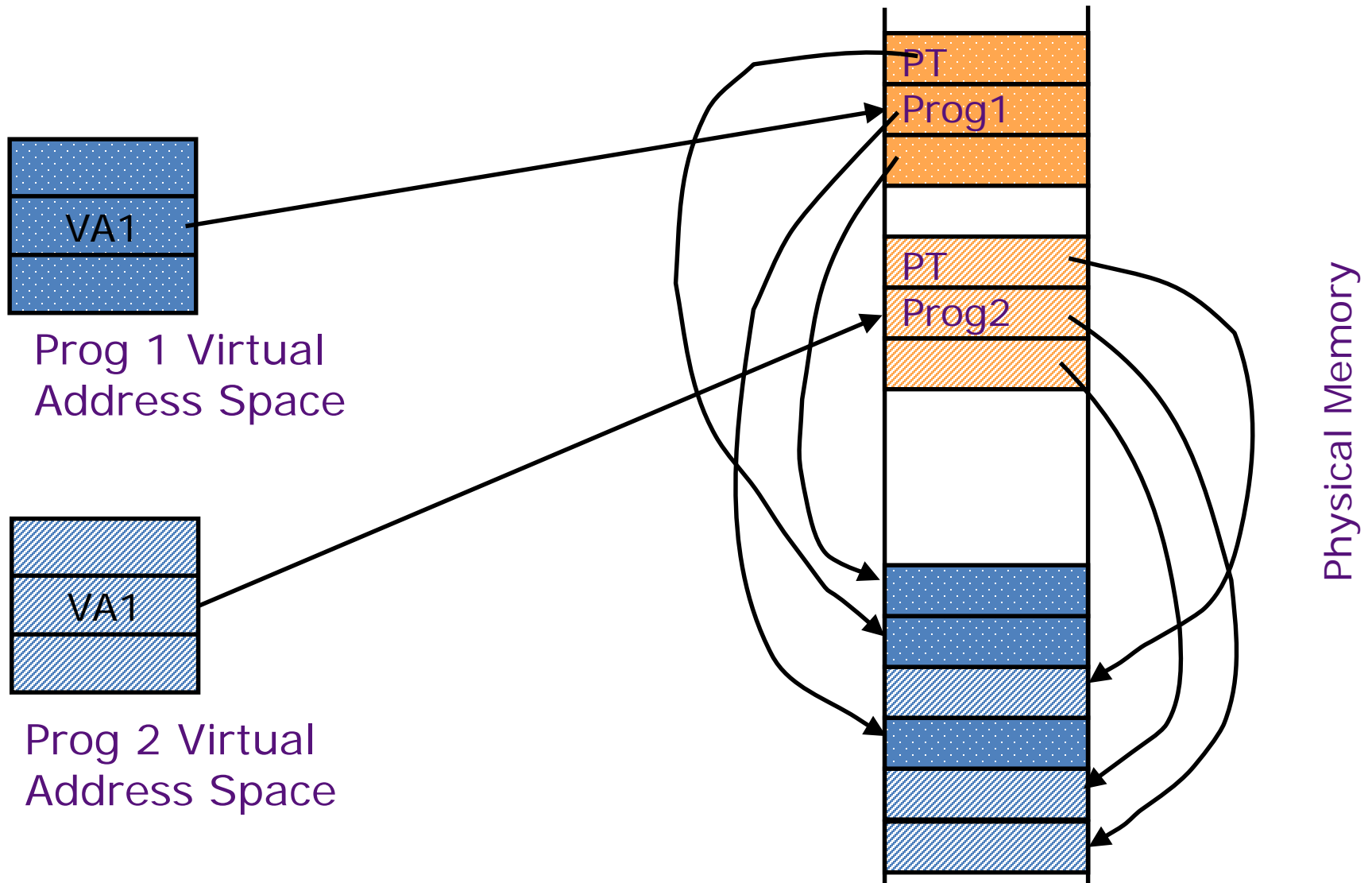
- Space required by the page tables (PT) is proportional to the address space, number of users, ...

    $\Rightarrow$ *Too large to keep in registers inside CPU*

- Idea: Keep page tables in the main memory

    – Needs one reference to retrieve the page base address and another to access the data word

    $\Rightarrow$ *doubles the number of memory references! (but we can fix this using something we already know about...)*

# Page Tables in Physical Memory

PT
Prog1

PT
Prog2

VA1

Prog 1 Virtual
Address Space

VA1

Prog 2 Virtual
Address Space

Physical Memory

# MT2 Grades Posted



| MINIMUM | MEDIAN | MAXIMUM | MEAN | STD DEV |
|---------|--------|---------|------|---------|
| 8.5 | 49.5 | 79.0 | 49.61 | 13.0 |

**62%**

25

# Administrivia
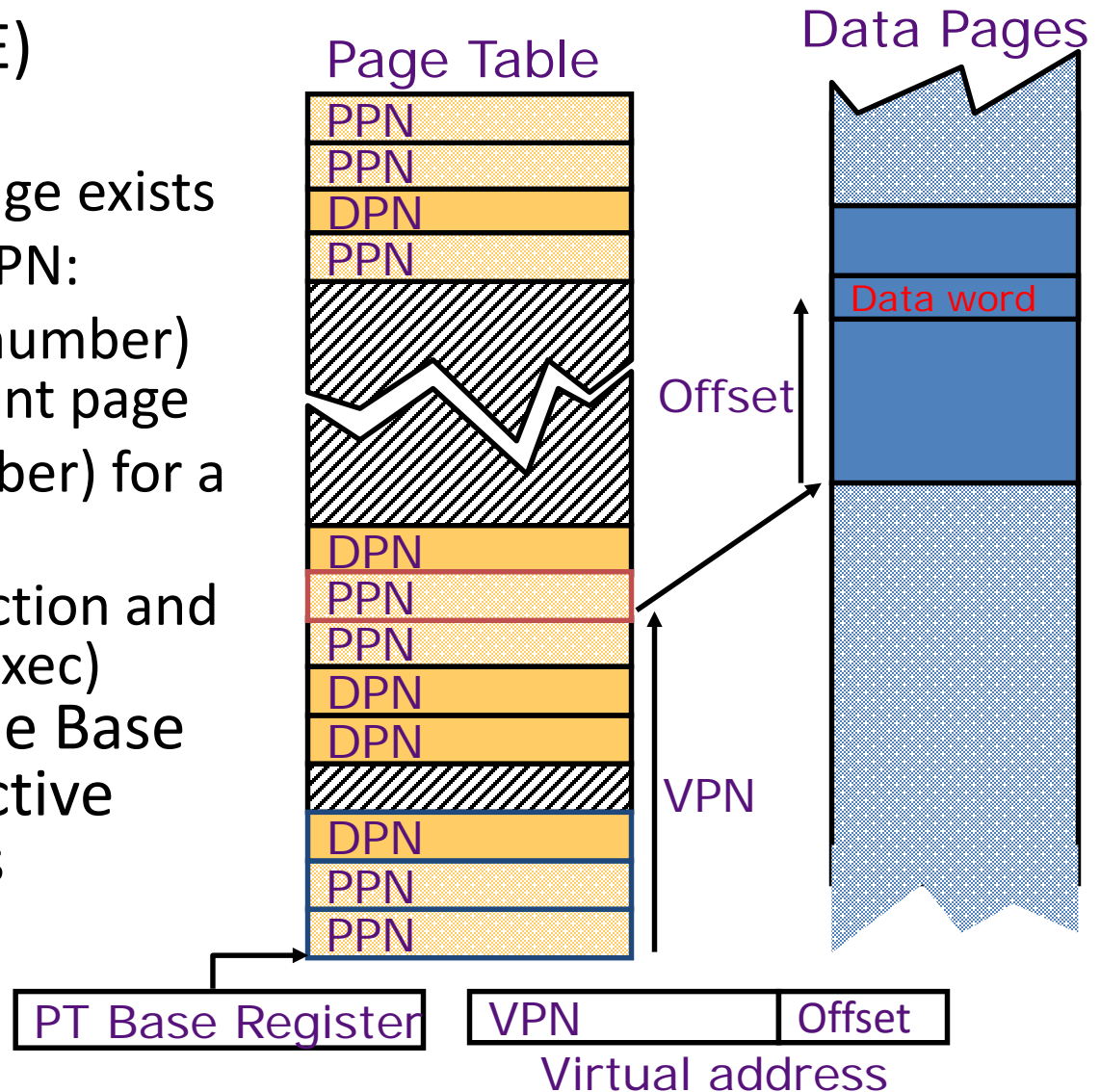
- Upcoming Lecture Schedule
  - 11/19: VM (today)
  - 11/24: I/O: DMA, Disks, Networking
  - 11/26: Thanksgiving Holiday (no class)
  - 12/01: Dependability: Parity, ECC, RAID
    - Last day of new material
  - 12/03: Summary, What's Next? (+ HKN reviews)

# Administrivia

- Okay to turn-in HW4 today (without slip charge)
- Project 4 programming competition rules posted.
  - Competition will end 11/30.
- Last HW (5) Virtual Memory
  - Due 12/06 (Sunday after end of classes)
- Last project (5): Spark
  - 5-1: posted, due 11/24 (next Wednesday)
  - 5-2: due 12/06 (Sunday after end of classes)
- <u>Final Exam</u> is Friday (12/18)
  - ***<u>7-10PM</u>***, RSF Field House
  - More info soon on review sessions, etc.

# Linear (simple) Page Table

- Page Table Entry (PTE) contains:
  - 1 bit to indicate if page exists
  - And either PPN or DPN:

    PPN (physical page number) for a memory-resident page

    DPN (disk page number) for a page on the disk
  - Status bits for protection and usage (read, write, exec)
- OS sets the Page Table Base Register whenever active user process changes

**Data Pages**

**Page Table**

PPN
PPN
DPN
PPN

DPN
PPN
PPN
DPN
DPN

DPN
PPN
PPN

Data word

Offset

VPN

PT Base Register | VPN | Offset

**Virtual address**

# Suppose an instruction references a memory page that isn't in DRAM?

- We get a exception of type "page fault"
- Page fault handler does the following:
  - If virtual page doesn't yet exist, assign an unused page in DRAM, or if page exists …
  - Initiate transfer of the page we're requesting from disk to DRAM, assigning to an unused page
  - If no unused page is left, a *page currently in DRAM is selected to be replaced* (based on usage)
  - The replaced page is written (back) to disk, page table entry that maps that VPN->PPN is marked as invalid/DPN
  - Page table entry of the page we're requesting is updated with a (now) valid PPN

# Size of Linear Page Table

With 32-bit memory addresses, 4-KB pages:

$\Rightarrow 2^{32} / 2^{12} = 2^{20}$ virtual pages per user, assuming 4-Byte PTEs,

$\Rightarrow 2^{20}$ PTEs, i.e, 4 MB page table per user!
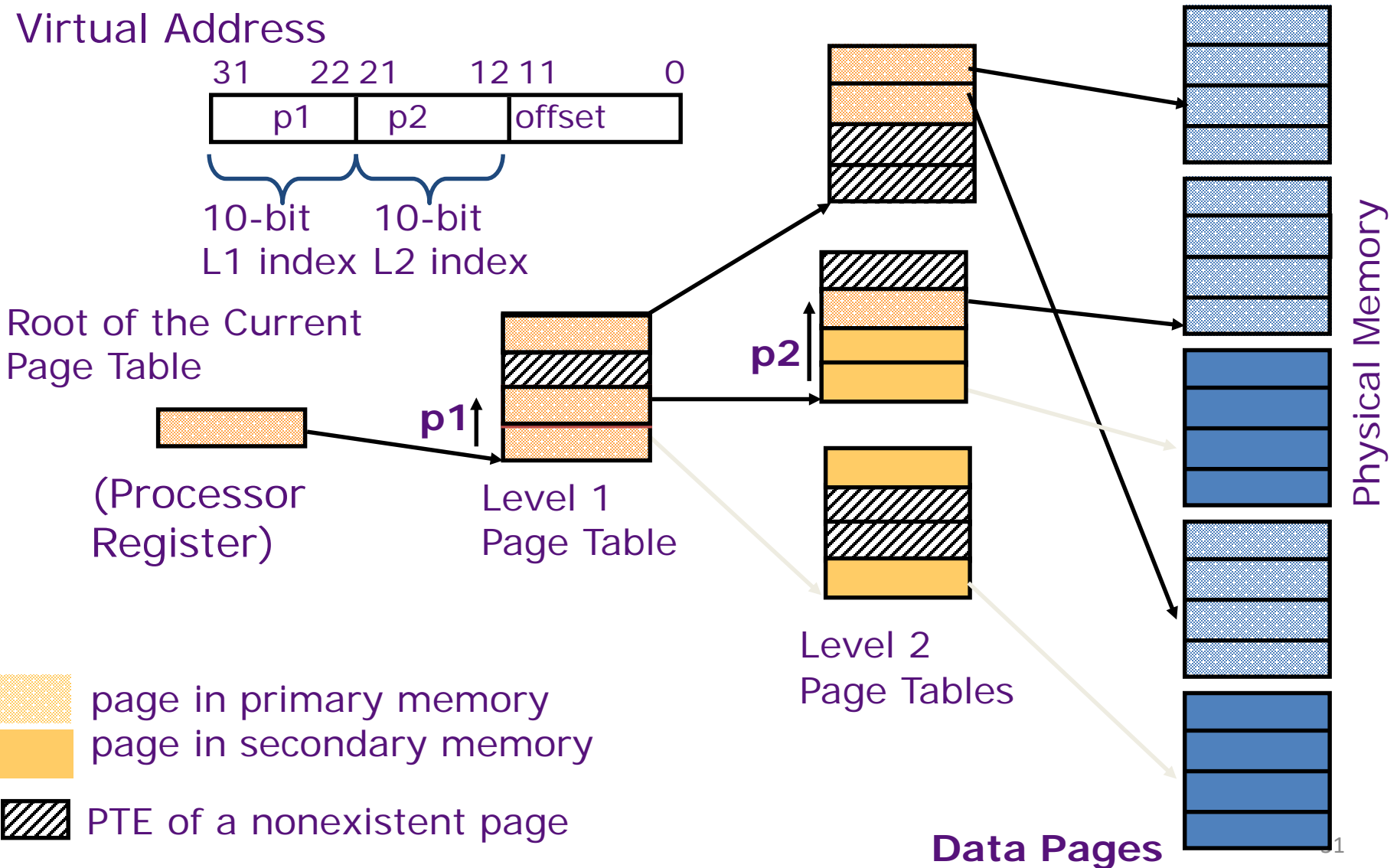
Larger pages?

- Internal fragmentation (Not all memory in page gets used)
- Larger page fault penalty (more time to read from disk)
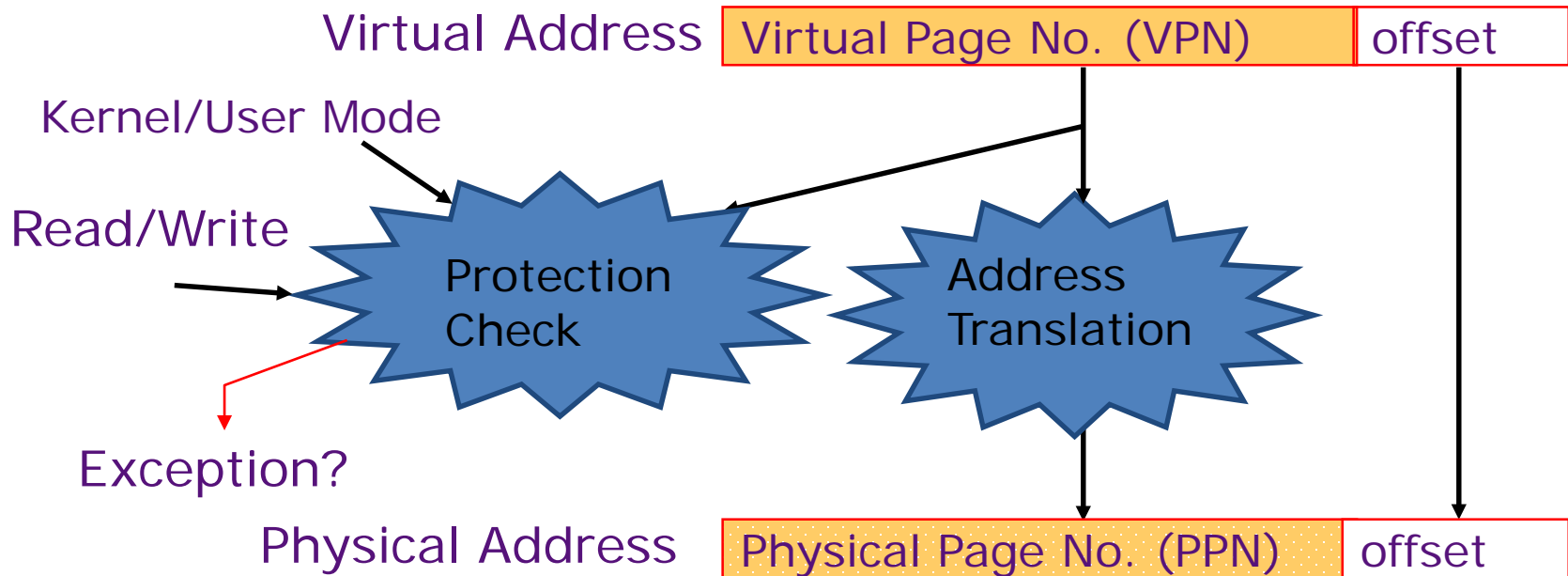
What about 64-bit virtual address space???

- Even 1MB pages would require $2^{44}$ 8-Byte PTEs (35 TB!)

*What is the "saving grace" ? Most processes only use a set of high address (stack), and a set of low address (instructions, heap)*

# *Hierarchical Page Table* – exploits sparcity of virtual address space use



Virtual Address

| 31 | 22 | 21 | 12 | 11 | 0 |
|---|---|---|---|---|---|
| p1 | | p2 | | offset | |

10-bit
L1 index

10-bit
L2 index

Root of the Current
Page Table

(Processor
Register)

**p1**

Level 1
Page Table

**p2**

Level 2
Page Tables

Physical Memory

page in primary memory
page in secondary memory

PTE of a nonexistent page

**Data Pages**

# Address Translation & Protection

Virtual Address    | Virtual Page No. (VPN) | offset |

Kernel/User Mode

Read/Write

Protection Check

Address Translation

Exception?

Physical Address    | Physical Page No. (PPN) | offset |

- Every instruction and data access needs address translation and protection checks

*A good VM design needs to be fast (~ one cycle) and space efficient*
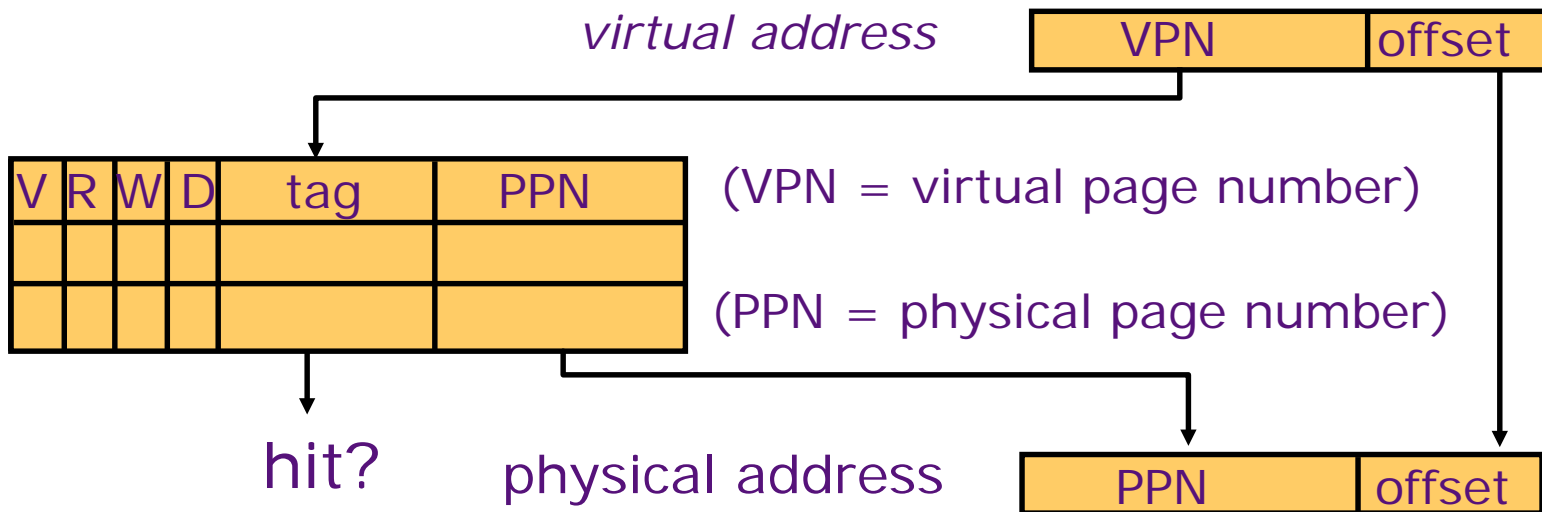
# Translation Lookaside Buffers (TLB)

Address translation is very expensive!
In a two-level page table, each reference becomes several memory accesses

Solution: *Cache some translations in TLB*

TLB hit   $\Rightarrow$ *Single-Cycle Translation*
TLB miss  $\Rightarrow$ *Page-Table Walk to refill*

*virtual address*

| | | | | | |
|---|---|---|---|---|---|
| | | | | VPN | offset |

| V | R | W | D | tag | PPN |
|---|---|---|---|---|---|
| | | | | | |
| | | | | | |

(VPN = virtual page number)

(PPN = physical page number)

hit?    physical address
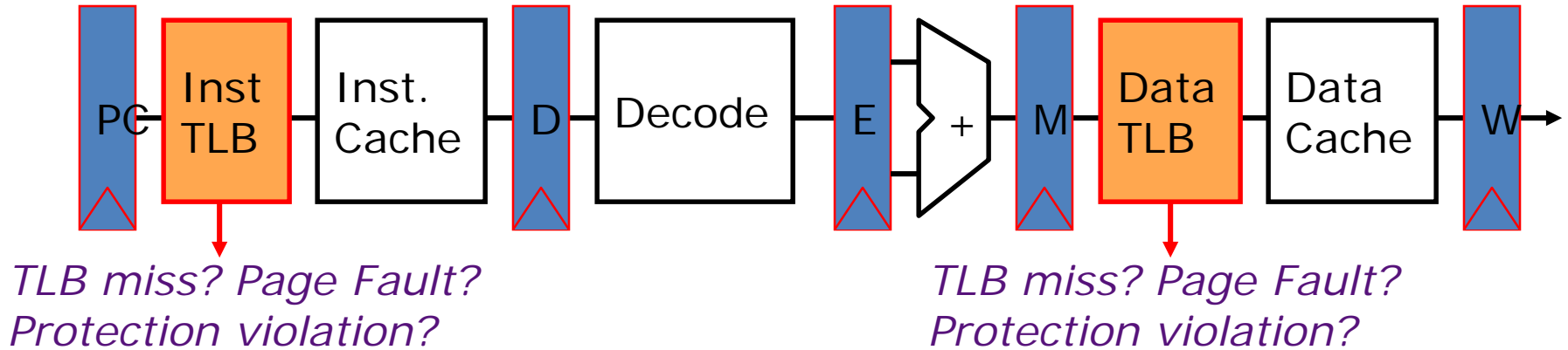
| PPN | offset |
|---|---|

# TLB Designs

- Typically 32-128 entries, usually fully associative
  - Each entry maps a large page, hence less spatial locality across pages => more likely that two entries conflict
  - Sometimes larger TLBs (256-512 entries) are 4-8 way set-associative
  - Larger systems sometimes have multi-level (L1 and L2) TLBs
- Random or FIFO replacement policy
- "TLB Reach": Size of largest virtual address space that can be simultaneously mapped by TLB

Example: 64 TLB entries, 4KB pages, one page per entry
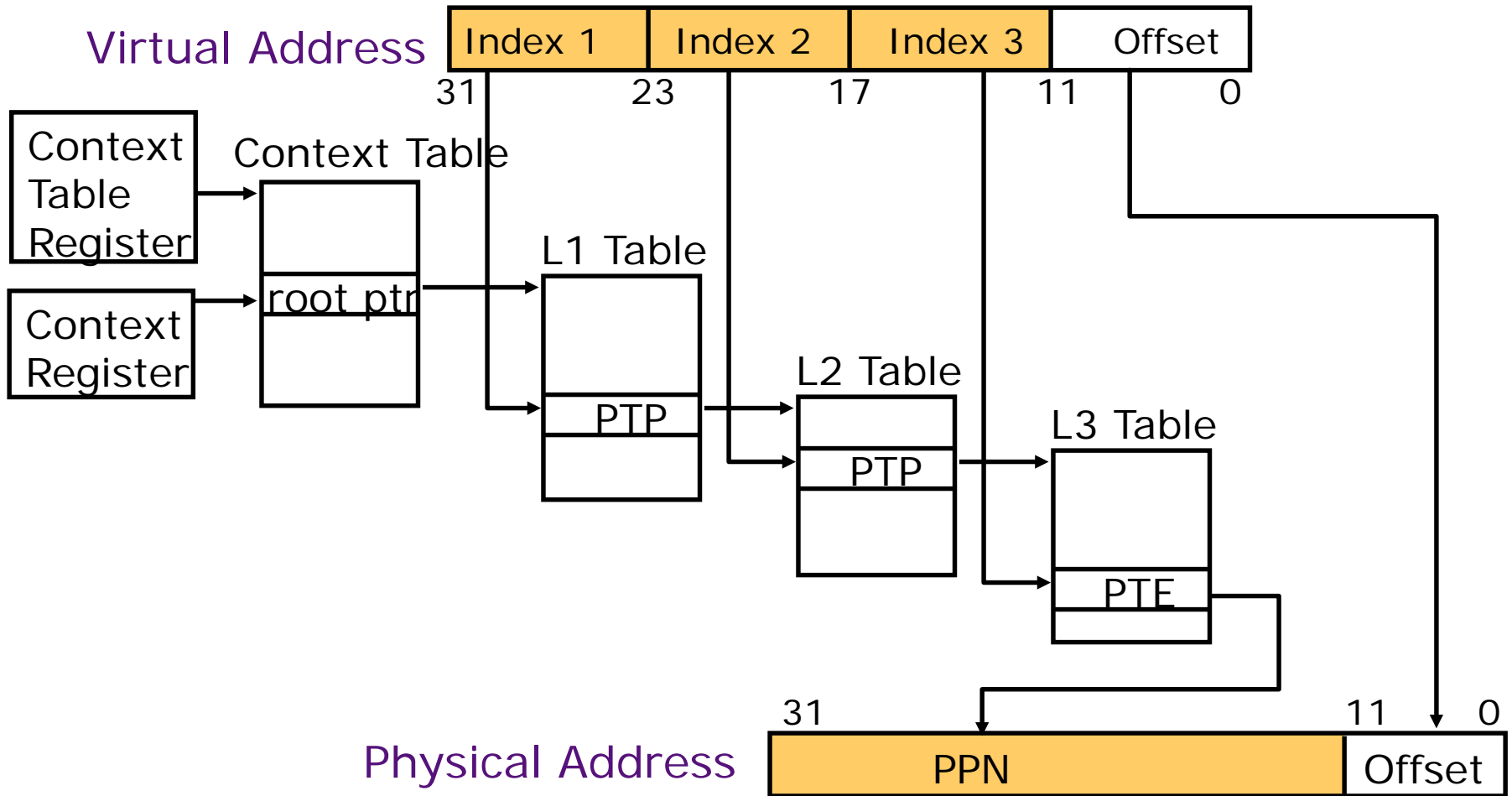
TLB Reach = _____?

# VM-related events in pipeline



PC — Inst TLB — Inst. Cache — D — Decode — E — + — M — Data TLB — Data Cache — W

*TLB miss? Page Fault? Protection violation?*

*TLB miss? Page Fault? Protection violation?*

- Handling a TLB miss needs a hardware or software mechanism to refill TLB
  - usually done in hardware now
- Handling a page fault (e.g., page is on disk) needs a *precise* trap so software handler can easily resume after retrieving page
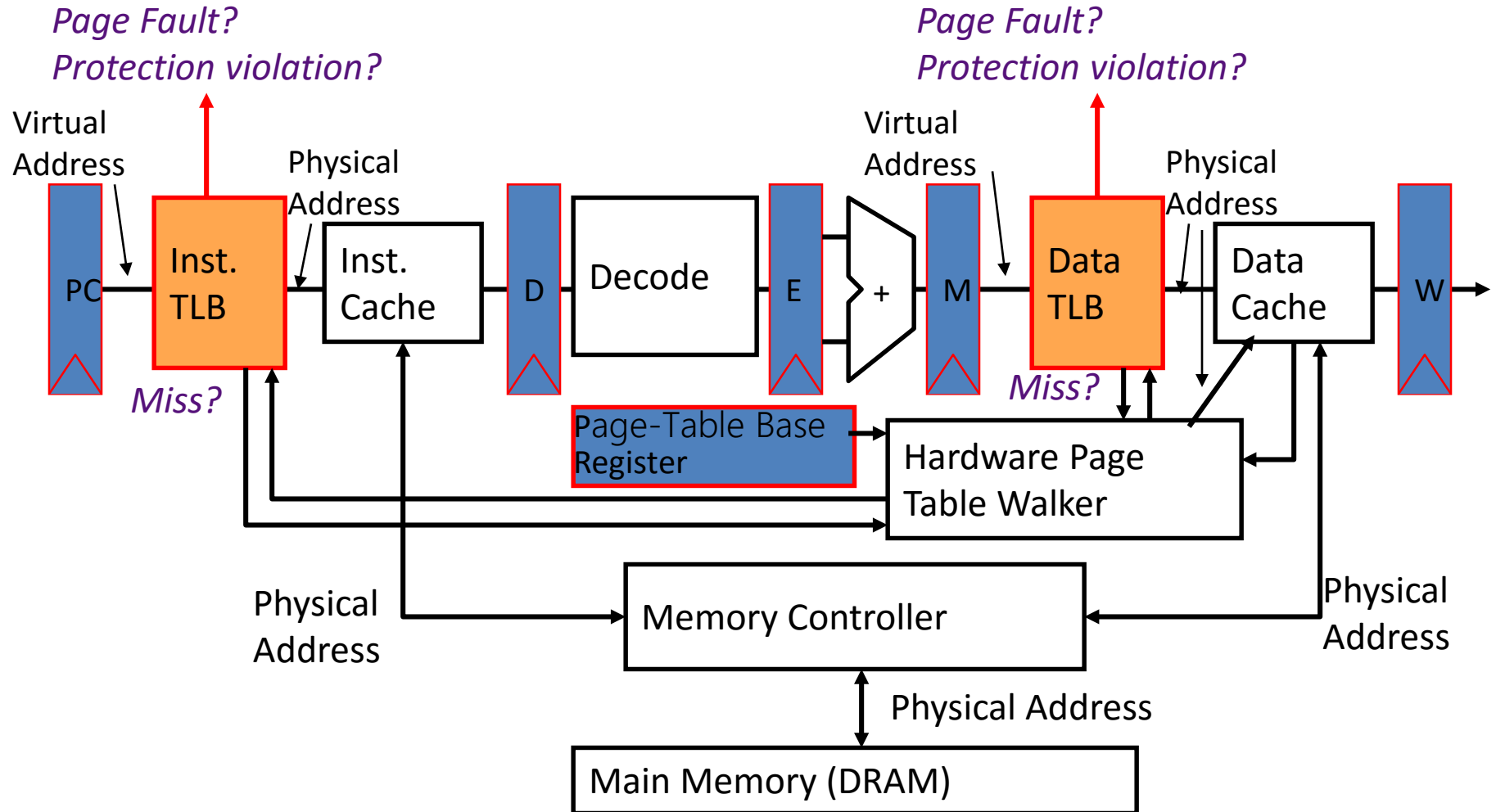- Handling protection violation may abort process

# Hierarchical Page Table Walk: SPARC v8

Virtual Address

| Index 1 | Index 2 | Index 3 | Offset |
|---------|---------|---------|--------|

31      23      17      11      0

Context Table Register

Context Register

Context Table

root ptr

L1 Table

PTP

L2 Table

PTP

L3 Table

PTE

Physical Address

31      11      0

| PPN | Offset |
|-----|--------|

MMU does this table walk in hardware on a TLB miss
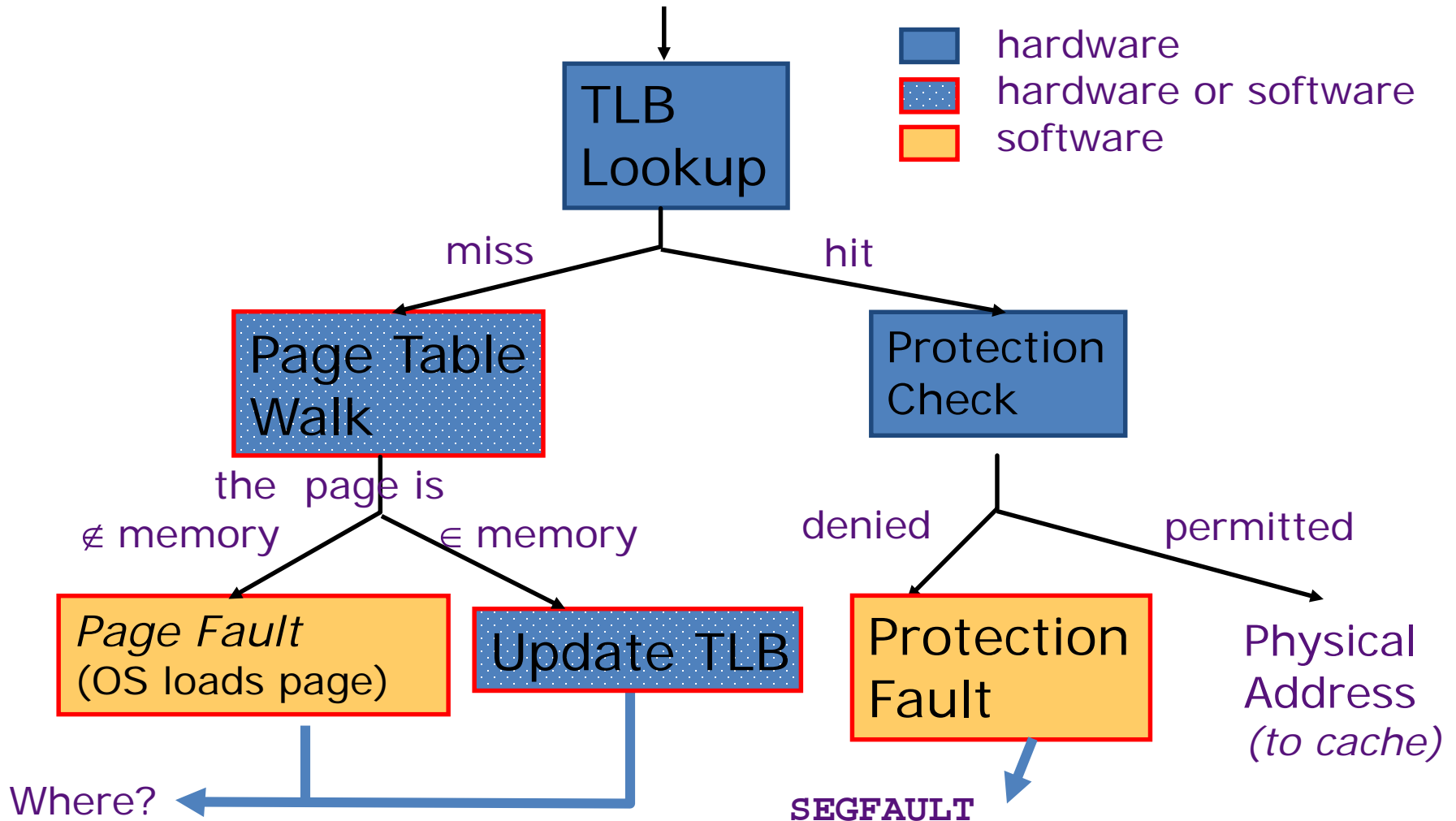
# Page-Based Virtual-Memory Machine

## (Hardware Page-Table Walk)



- Assumes page tables held in untranslated physical memory

# Address Translation:
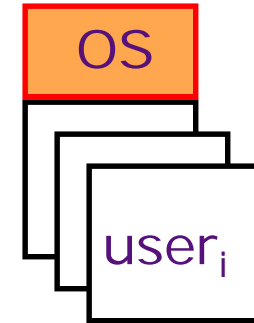
*putting it all together*

Virtual Address



hardware
hardware or software
software

TLB
Lookup

miss                    hit

Page Table
Walk

Protection
Check

the page is

∉ memory        ∈ memory        denied              permitted

*Page Fault*
(OS loads page)

Update TLB

Protection
Fault

Physical
Address
*(to cache)*

Where?

SEGFAULT

# Modern Virtual Memory Systems
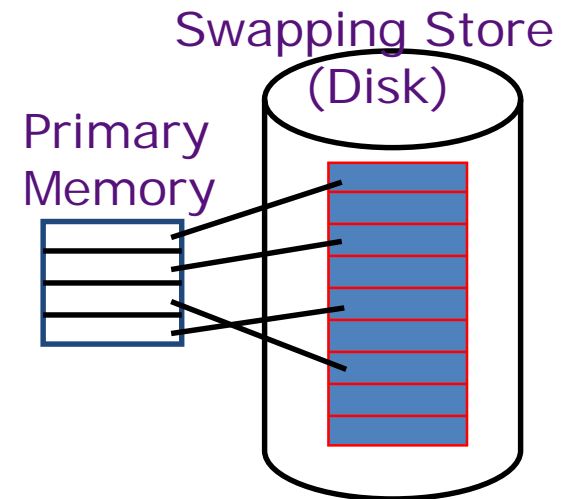
*Illusion of a large, private, uniform store*

## Protection & Privacy
several users, each with their private address space and one or more shared address spaces
page table $\equiv$ name space

OS

$user_i$

## Demand Paging
Provides the ability to run programs larger than the primary memory

Hides differences in machine configurations

Swapping Store (Disk)

Primary Memory

*The price is address translation on each memory reference*

VA → mapping  TLB → PA

# Clicker Question

Let's try to extrapolate from caches… Which one is false?

A. # offset bits in V.A. = log2(page size)

B. # offset bits in P.A. = log2(page size)

C. # VPN bits in V.A. = log2(# of physical pages)

D. # PPN bits in P.A. = log2(# of physical pages)

E. A single-level page table contains a PTE for every possible VPN in the system

# Conclusion: VM features track historical uses

- **Bare machine, only physical addresses**
  - One program owned entire machine
- **Batch-style multiprogramming**
  - Several programs sharing CPU while waiting for I/O
  - Base & bound: translation and protection between programs (not virtual memory)
  - Problem with external fragmentation (holes in memory), needed occasional memory defragmentation as new jobs arrived
- **Time sharing**
  - More interactive programs, waiting for user.  Also, more jobs/second.
  - Motivated move to fixed-size page translation and protection, no external fragmentation (but now internal fragmentation, wasted bytes in page)
  - Motivated adoption of virtual memory to allow more jobs to share limited physical memory resources while holding working set in memory
- **Virtual Machine Monitors**
  - Run multiple operating systems on one machine
  - Idea from 1970s IBM mainframes, now common on laptops
    - e.g., run Windows on top of Mac OS X
  - Hardware support for two levels of translation/protection
    - Guest OS virtual -> Guest OS physical -> Host machine physical
  - Also basis of Cloud Computing
    - Virtual machine instances on EC2 for Lab 13