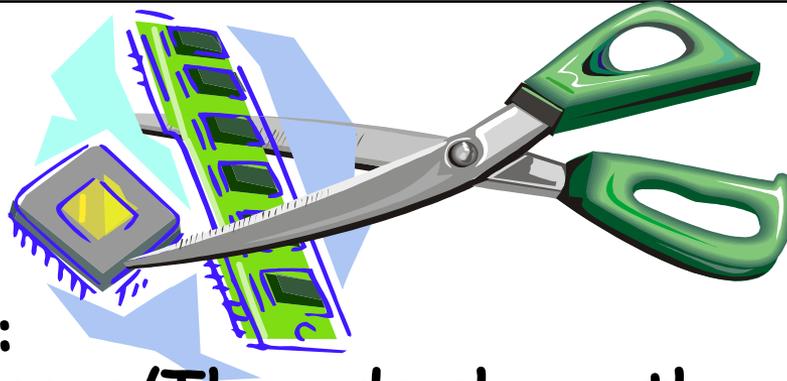# Main memory and virtual memory

Adapted by Tiziano Villa from lecture notes by Prof. John Kubiatowicz (UC Berkeley)
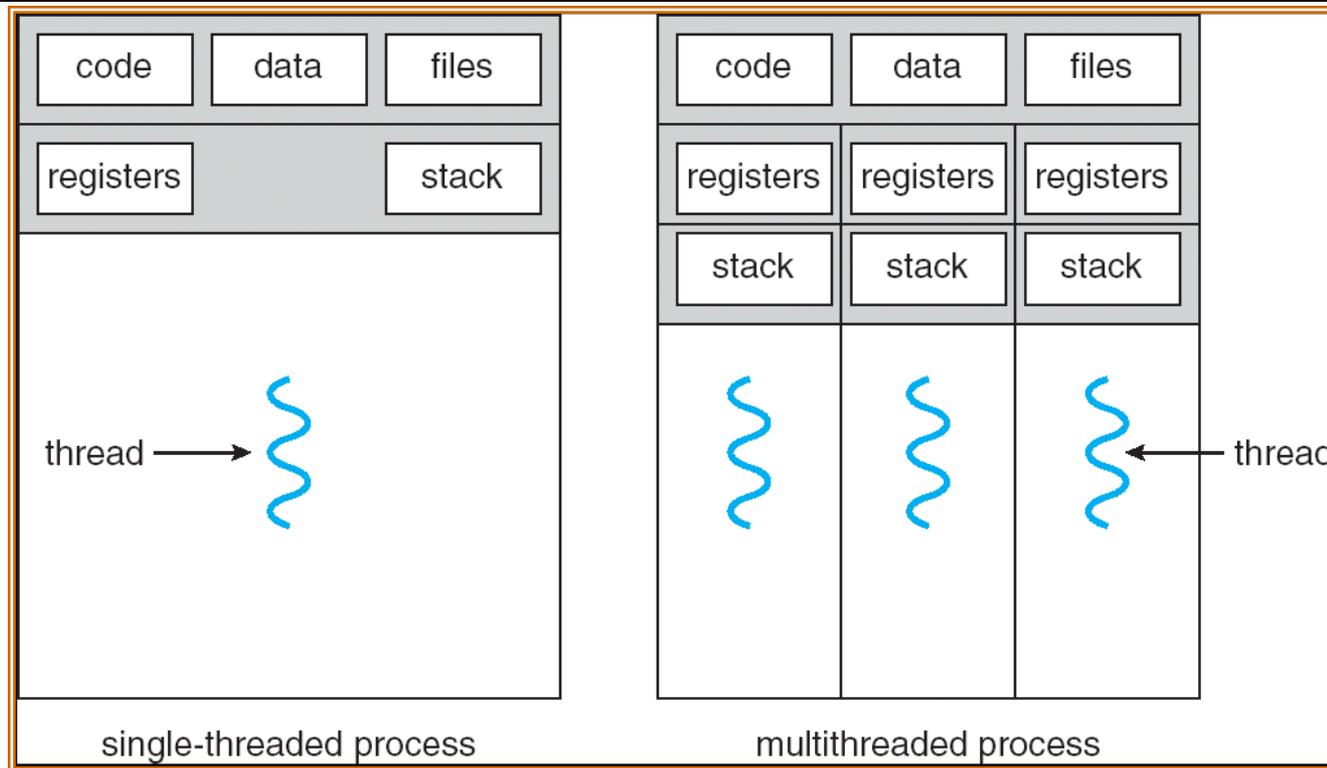
Main memory
and
Address Translation

# Virtualizing Resources

- ## Physical Reality: Different Processes/Threads share the same hardware
  - Need to multiplex CPU (Just finished: scheduling)
  - Need to multiplex use of Memory (Today)
  - Need to multiplex disk and devices (later in term)
- ## Why worry about memory sharing?
  - The complete working state of a process and/or kernel is defined by its data in memory (and registers)
  - Consequently, cannot just let different threads of control use the same memory
    - » Physics: two different pieces of data cannot occupy the same locations in memory
  - Probably don't want different threads to even have access to each other's memory (protection)

# Recall: Single and Multithreaded Processes



| code | data | files |
| registers | | stack |
| thread → |

single-threaded process

| code | data | files |
| registers | registers | registers |
| stack | stack | stack |
| | | ← thread |

multithreaded process

- **Threads encapsulate concurrency**
    - "Active" component of a process
- **Address spaces encapsulate protection**
    - Keeps buggy program from trashing the system
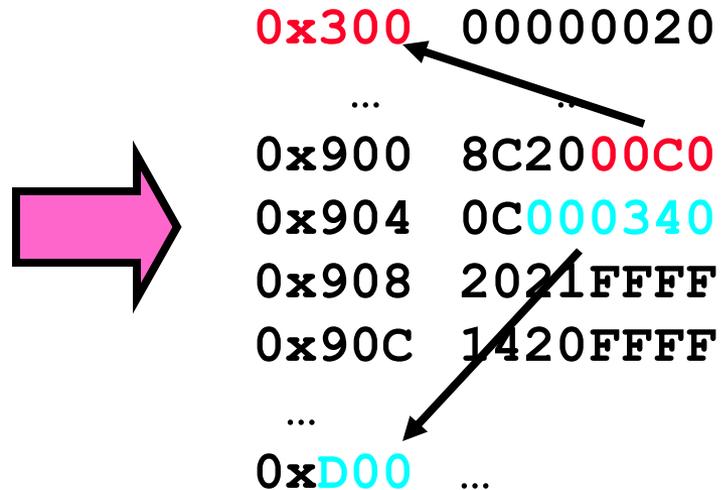    - "Passive" component of a process

# Important Aspects of Memory Multiplexing

- **Controlled overlap:**
  - Separate state of threads should not collide in physical memory.  Obviously, unexpected overlap causes chaos!
  - Conversely, would like the ability to overlap when desired (for communication)
- **Translation:**
  - Ability to translate accesses from one address space (virtual) to a different one (physical)
  - When translation exists, processor uses virtual addresses, physical memory uses physical addresses
  - Side effects:
    » Can be used to avoid overlap
    » Can be used to give uniform view of memory to programs
- **Protection:**
  - Prevent access to private memory of other processes
    » Different pages of memory can be given special behavior (Read Only, Invisible to user programs, etc).
    » Kernel data protected from User programs
    » Programs protected from themselves

# Binding of Instructions and Data to Memory

- **Binding of instructions and data to addresses:**
  - **Choose addresses for instructions and data from the standpoint of the processor**

```
data1: dw    32                        0x300  00000020
             …                         …      …
start: lw    r1,0(data1)               0x900  8C2000C0
       jal   checkit                   0x904  0C000340
loop:  addi  r1, r1, -1                0x908  2021FFFF
       bnz   r1, r0, loop              0x90C  1420FFFF
             …                         …
checkit: …                             0xD00  …
```
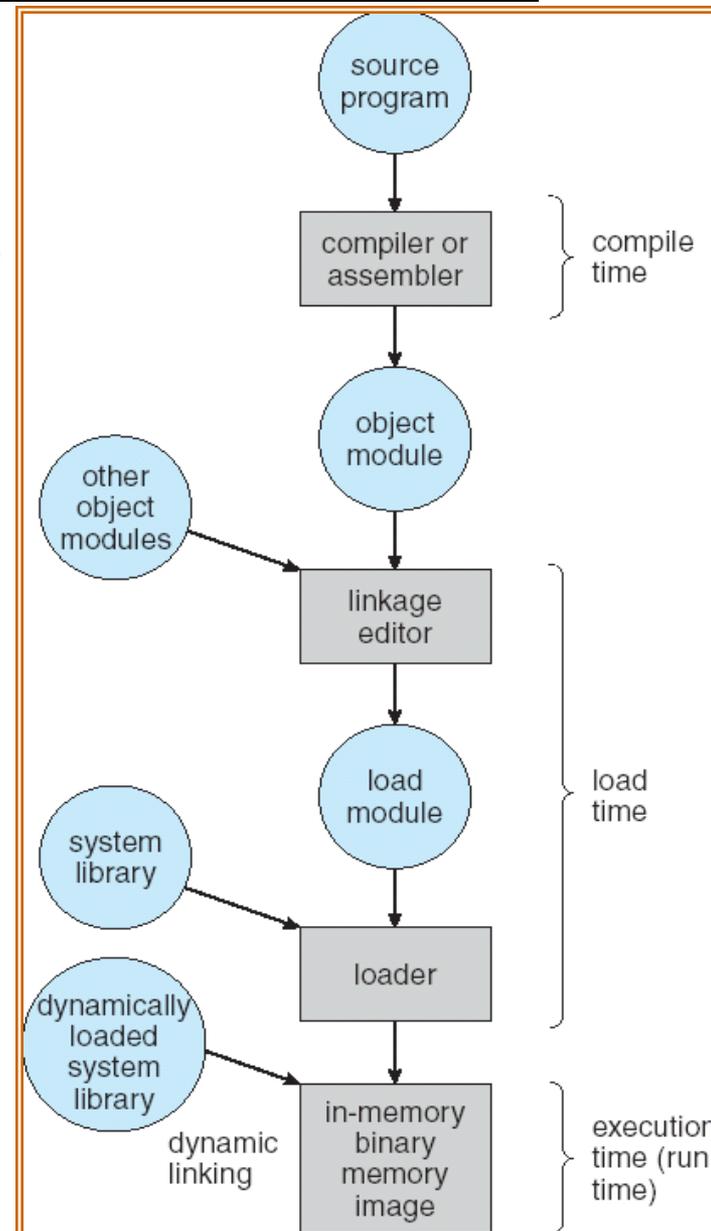
  - **Could we place** `data1`, `start`, **and/or** `checkit` **at different addresses?**
    - » Yes
    - » When? Compile time/Load time/Execution time
  - **Related: which physical memory locations hold particular instructions or data?**
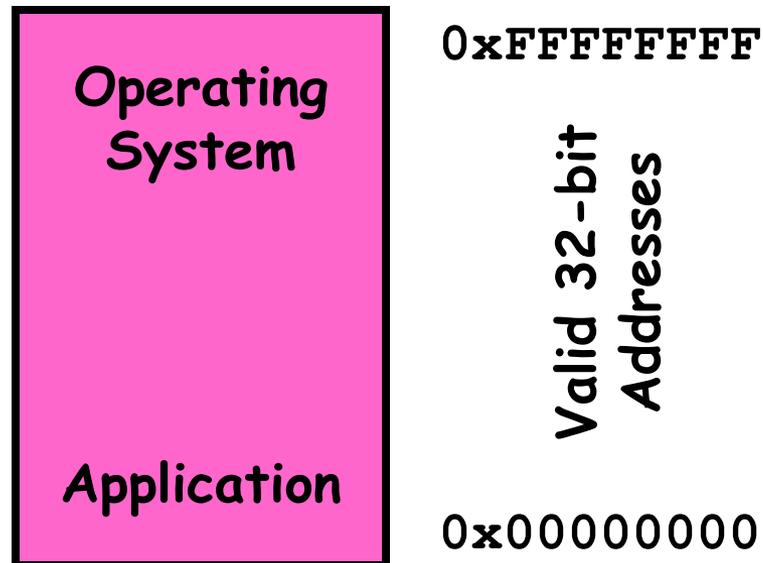
# Multi-step Processing of a Program for Execution

- **Preparation of a program for execution involves components at:**
  - Compile time (i.e. "gcc")
  - Link/Load time (unix "ld" does link)
  - Execution time (e.g. dynamic libs)
- **Addresses can be bound to final values anywhere in this path**
  - Depends on hardware support
  - Also depends on operating system
- **Dynamic Libraries**
  - Linking postponed until execution
  - Small piece of code, *stub*, used to locate the appropriate memory-resident library routine
  - Stub replaces itself with the address of the routine, and executes routine

# Uniprogramming

- **Uniprogramming (no Translation or Protection)**
  - Application always runs at same place in physical memory since only one application at a time
  - Application can access any physical address

**Operating System**

**Application**

0xFFFFFFFF

Valid 32-bit Addresses

0x00000000

  - Application given illusion of dedicated machine by giving it reality of a dedicated machine
- **Of course, this doesn't help us with multithreading**

# Multiprogramming (First Version)

- **Multiprogramming without Translation or Protection**
  - **Must somehow prevent address overlap between threads**

| Operating System | 0xFFFFFFFF |
|---|---|
| | |
| Application2 | 0x00020000 |
| | |
| Application1 | 0x00000000 |

  - **Trick: Use Loader/Linker: Adjust addresses while program loaded into memory (loads, stores, jumps)**
    - » Everything adjusted to memory location of program
    - » Translation done by a linker-loader
    - » Was pretty common in early days
- **With this solution, no protection: bugs in any program can cause other programs to crash or even the OS**

# Multiprogramming (Version with Protection)

- ## Can we protect programs from each other without translation?

| | | |
|---|---|---|
| **Operating System** | 0xFFFFFFFF | |
| | | Limit=+0x10000 |
| **Application2** | | Base=0x20000 |
| | 0x00020000 | |
| **Application1** | | |
| | 0x00000000 | |

- – Yes: use two special registers *Base* and *Limit* to prevent user from straying outside designated area
  - » If user tries to access an illegal address, cause an error
- – During switch, kernel loads new base/limit from TCB
  - » User not allowed to change base/limit registers

# Segmentation with Base and Limit registers



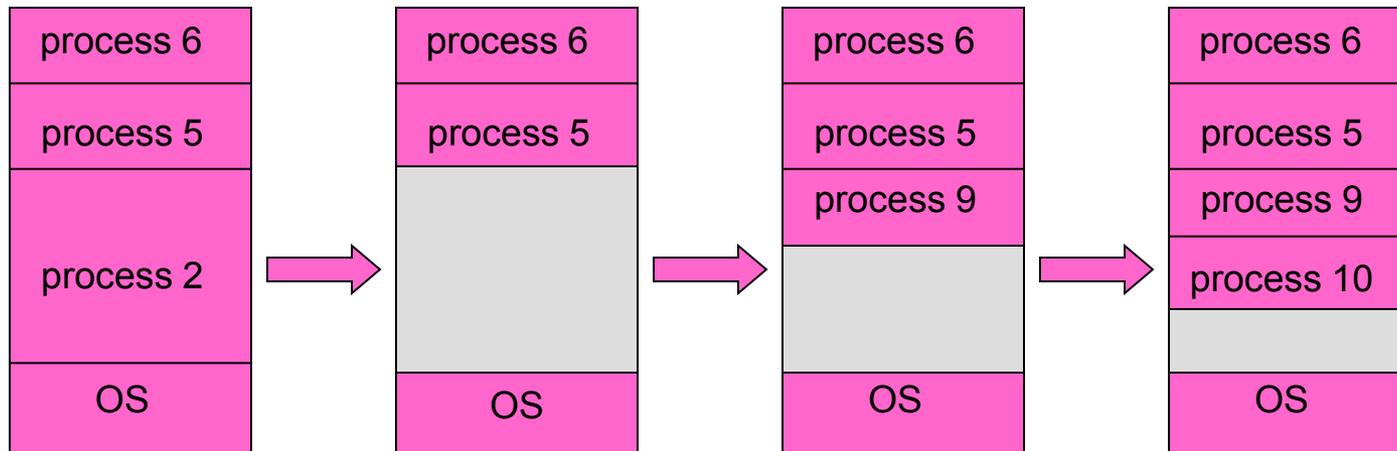- **Could use base/limit for dynamic address translation** (often called "segmentation"):
  - Alter address of every load/store by adding "base"
  - User allowed to read/write within segment
    - » Accesses are relative to segment so don't have to be relocated when program moved to different segment
  - User may have multiple segments available (e.g x86)
    - » Loads and stores include segment ID in opcode:
      x86 Example: `mov [es:bx],ax`.
    - » Operating system moves around segment base pointers as necessary

# Issues with simple segmentation method

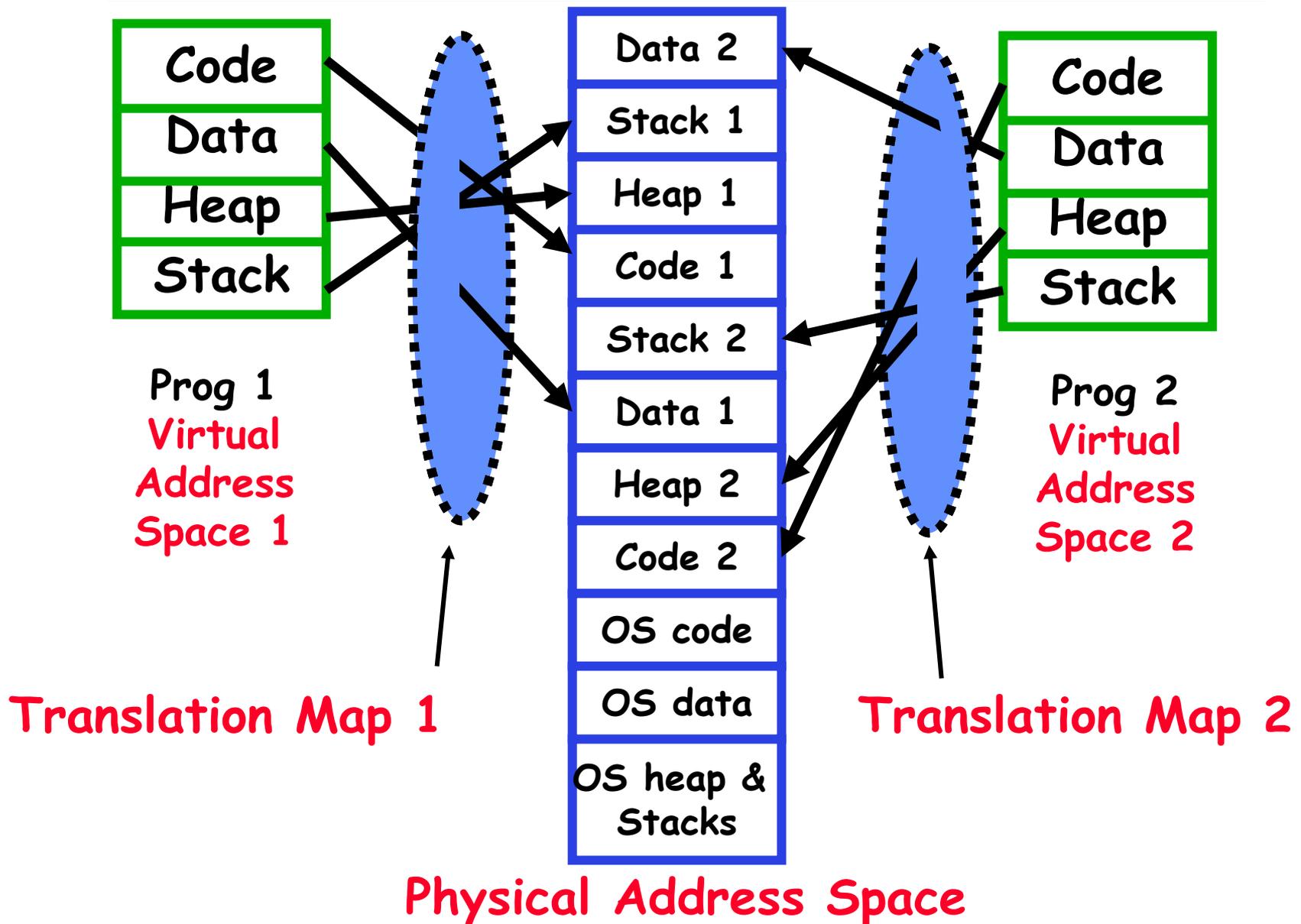| process 6 | | process 6 | | process 6 | | process 6 |
|---|---|---|---|---|---|---|
| process 5 | → | process 5 | → | process 5 | → | process 5 |
| process 2 | | | | process 9 | | process 9 |
| | | | | | | process 10 |
| OS | | OS | | OS | | OS |

- **Fragmentation problem**
  - **Not every process is the same size**
  - **Over time, memory space becomes fragmented**
- **Hard to do inter-process sharing**
  - **Want to share code segments when possible**
  - **Want to share memory between processes**
  - **Helped by by providing multiple segments per process**
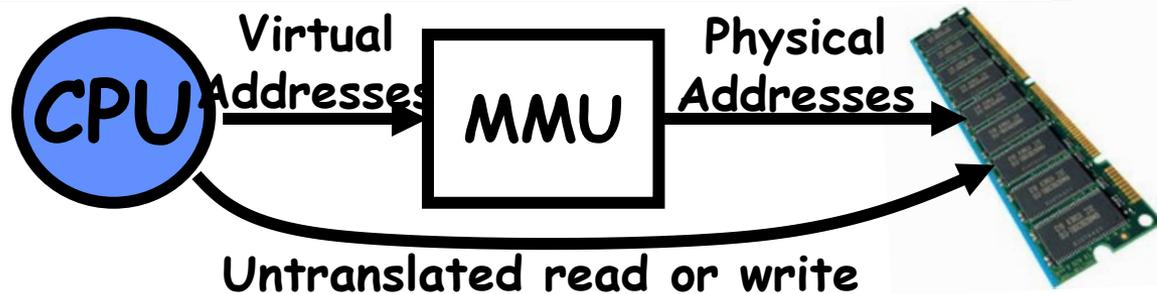- **Need enough physical memory for every process**

# Multiprogramming (Translation and Protection version 2)

- Problem: Run multiple applications in such a way that they are protected from one another
- Goals:
  - Isolate processes and kernel from one another
  - Allow flexible translation that:
    - » Doesn't lead to fragmentation
    - » Allows easy sharing between processes
    - » Allows only part of process to be resident in physical memory
- (Some of the required) Hardware Mechanisms:
  - General Address Translation
    - » Flexible: Can fit physical chunks of memory into arbitrary places in users address space
    - » Not limited to small number of segments
    - » Think of this as providing a large number (thousands) of fixed-sized segments (called "pages")
  - Dual Mode Operation
    - » Protection base involving kernel/user distinction

# Example of General Address Translation



Code
Data
Heap
Stack

Prog 1
Virtual
Address
Space 1

Data 2
Stack 1
Heap 1
Code 1
Stack 2
Data 1
Heap 2
Code 2
OS code
OS data
OS heap &
Stacks

Code
Data
Heap
Stack

Prog 2
Virtual
Address
Space 2

**Translation Map 1**

**Translation Map 2**

**Physical Address Space**

# Two Views of Memory



CPU — Virtual Addresses → MMU — Physical Addresses → (memory)
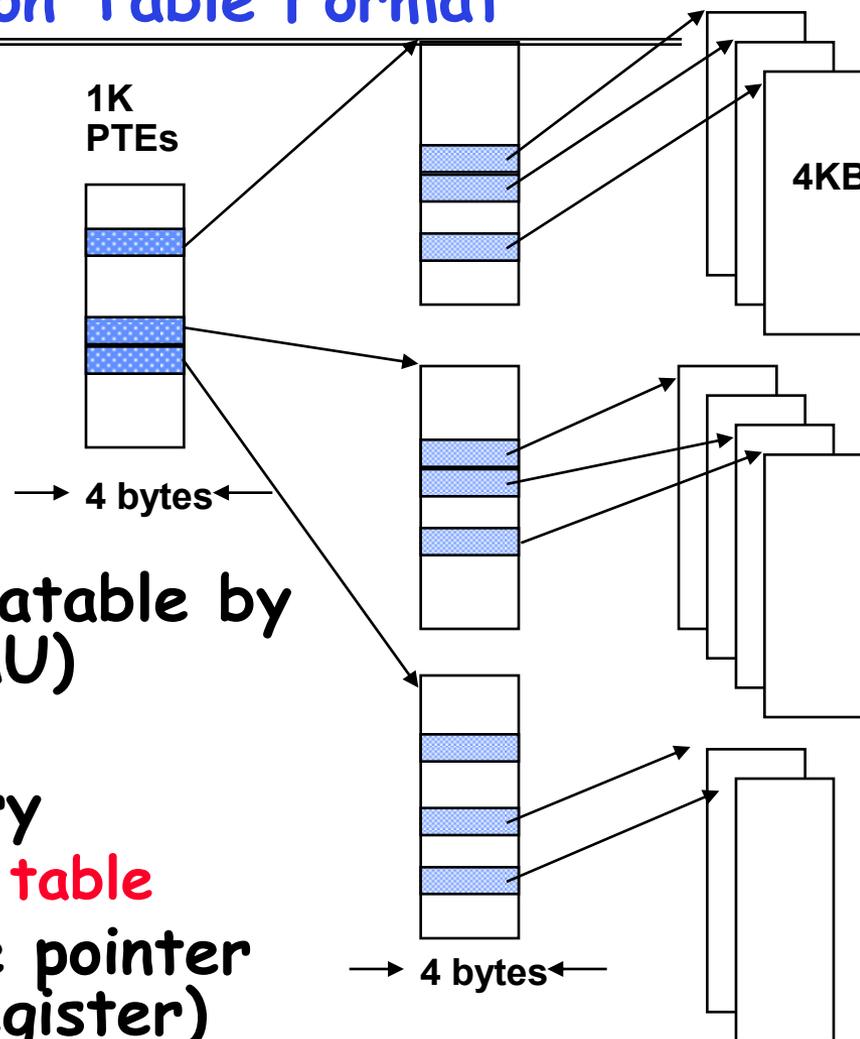
Untranslated read or write

- **Recall: Address Space:**
  - – **All the addresses and state a process can touch**
  - – **Each process and kernel has different address space**
- **Consequently: two views of memory:**
  - – **View from the CPU (what program sees, virtual memory)**
  - – **View fom memory (physical memory)**
  - – **Translation box converts between the two views**
- **Translation helps to implement protection**
  - – **If task A cannot even gain access to task B's data, no way for A to adversely affect B**
- **With translation, every program can be linked/loaded into same region of user address space**
  - – **Overlap avoided through translation, not relocation**

# Example of Translation Table Format

**Two-level Page Tables**
**32-bit address:**

| 10 | 10 | 12 |
|---|---|---|
| P1 index | P2 index | page offset |

**1K PTEs**

→ 4 bytes ←

4KB

→ 4 bytes ←

- **Page: a unit of memory translatable by memory management unit (MMU)**
  - **Typically 1K – 8K**
- **Page table structure in memory**
  - **Each user has different page table**
- **Address Space switch: change pointer to base of table (hardware register)**
  - **Hardware traverses page table (for many architectures)**
  - **MIPS uses software to traverse table**

# Address Translation Schemes

- Segmentation
- Paging
- Multi-level translation
- Paged page tables
- Inverted page tables

# More Flexible Segmentation

**Logical View**

- subroutine
- stack
- symbol table
- *Sqrt*
- main program

logical address

1

2

3

4

user view of memory space

1
4
2
3

physical memory space

- **Logical View: multiple separate segments**
  - Typical: Code, Data, Stack
  - Others: memory sharing, etc
- **Each segment is given region of contiguous memory**
  - Has a base and limit
  - Can reside anywhere in physical memory

# Implementation of Multi-Segment Model



- **Segment map resides in processor**
  - Segment number mapped into base/limit pair
  - Base added to offset to generate physical address
  - Error check catches offset out of range
- **As many chunks of physical memory as entries**
  - Segment addressed by portion of virtual address
  - However, could be included in instruction instead:
    » x86 Example: `mov [es:bx],ax.`
- **What is "V/N"?**
  - Can mark segments as invalid; requires check as well

# Intel x86 Special Registers

## 80386 Special Registers

**Segment registers**

| | |
|---|---|
| Code Seg. | Data Seg. |
| 15 CS 0 | 15 DS 0 |
| Stack Seg. | Extra Seg. |
| 15 SS 0 | 15 ES 0 |
| Extra Seg. | Extra. Seg |
| 15 FS 0 | 15 GS 0 |

| X | N T | IO PL | O F | D F | I F | T F | S F | Z F | X | A F | X | P F | X | C F |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 15 | 14 | 13 | 12 | 11 | 10 | 9 | 8 | 7 | 6 | 5 | 4 | 3 | 2 | 1 | 0 |

| P G | | | E T | T S | T S | M P | P E | CR0 |
|---|---|---|---|---|---|---|---|---|
| 31 30 | | | 5 | 4 | 3 | 2 | 1 | 0 |

| Unused | CR1 |
|---|---|
| 31 | 0 Flags |

| Page Fault Linear Address | CR2 |
|---|---|
| 31 | 0 |

| Page Directory Base Register | Not Used | CR3 |
|---|---|---|
| 31 | 7 | 0 |

PG=Paging Enable
ET=Emulation Type
TS=Task Switched
EM=Emulate Coprocessor
MP=Math coprocessor present
PE=Protected Mode enable

X=Reserved
NT=Nested Task
IOPL=I/O Privilege Level
OF=Overflow Flag
DF=Direction Flag
IF=Interrupt Flag
TF=Trap Flag
SF=Sign Flag
ZF=Zero Flag
AF=Auxiliary Flag
PF=Parity Flag
CF=Carry Flag

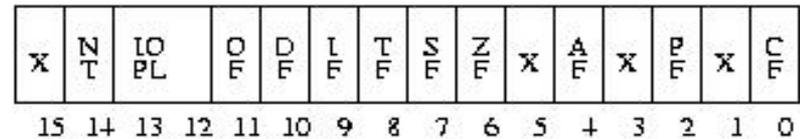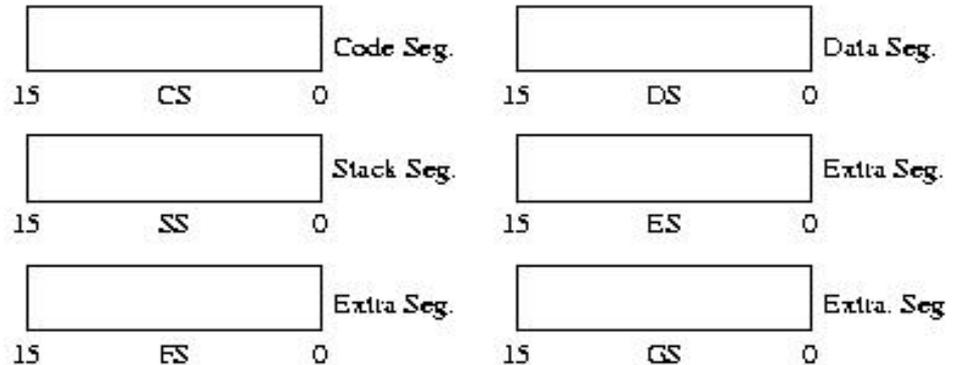| 15 | | | 1 | 0 |
|---|---|---|---|---|
| Index | T I | RPL | |

RPL = Requestor Privilege Level
TI = Table Indicator
    (0 = GDT, 1 = LDT)
Index = Index into table

Protected Mode segment selector

**Typical Segment Register
Current Priority is RPL
Of Code Segment (CS)**

# Example: Four Segments (16 bit addresses)

| Seg | Offset |
|-----|--------|

15 14 13                                          0

**Virtual Address Format**

| Seg ID # | Base | Limit |
|----------|--------|--------|
| 0 (code) | 0x4000 | 0x0800 |
| 1 (data) | 0x4800 | 0x1400 |
| 2 (shared) | 0xF000 | 0x1000 |
| 3 (stack) | 0x0000 | 0x3000 |

0x0000

0x4000

0x8000

0xC000

**Virtual
Address Space**

0x0000

0x4000
0x4800

0x5C00

0xF000

Might
be shared

Space for
Other Apps

Shared with
Other Apps

**Physical
Address Space**

# Example of segment translation

```
0x240   main:     la $a0, varx
0x244             jal strlen
   …                  …
0x360   strlen:   li   $v0, 0   ;count
0x364   loop:     lb   $t0, ($a0)
0x368             beq  $r0,$t1, done
   …                  …
0x4050  varx      dw   0x314159
```

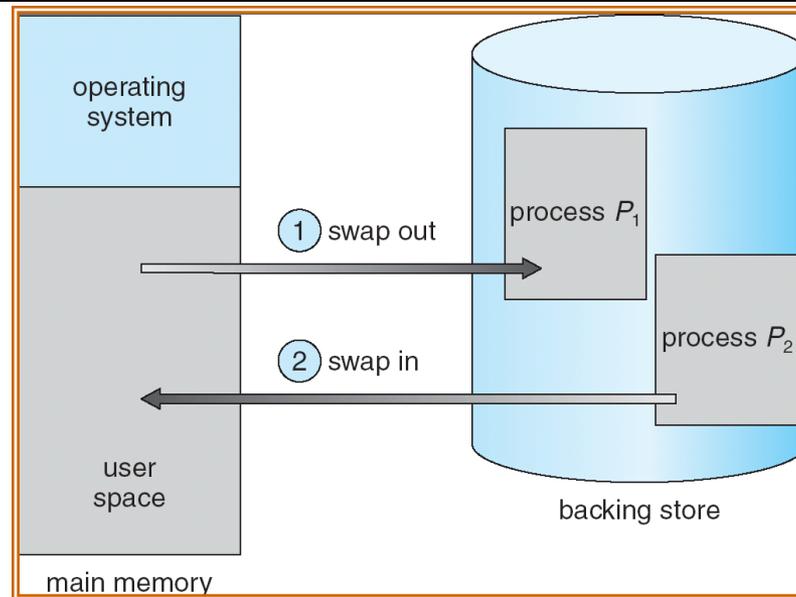| Seg ID # | Base | Limit |
|---|---|---|
| 0 (code) | 0x4000 | 0x0800 |
| 1 (data) | 0x4800 | 0x1400 |
| 2 (shared) | 0xF000 | 0x1000 |
| 3 (stack) | 0x0000 | 0x3000 |

Let's simulate a bit of this code to see what happens (PC=0x240):

1. Fetch 0x240. Virtual segment #? 0; Offset? 0x240
   Physical address? Base=0x4000, so physical addr=0x4240
   Fetch instruction at 0x4240. Get "la $a0, varx"
   Move 0x4050 → $a0, Move PC+4→PC

2. Fetch 0x244. Translated to Physical=0x4244. Get "jal strlen"
   Move 0x0248 → $ra (return address!), Move 0x0360 → PC

3. Fetch 0x360. Translated to Physical=0x4360. Get "li $v0,0"
   Move 0x0000 → $v0, Move PC+4→PC

4. Fetch 0x364. Translated to Physical=0x4364. Get "lb $t0,($a0)"
   Since $a0 is 0x4050, try to load byte from 0x4050
   Translate 0x4050. Virtual segment #? 1; Offset? 0x50
   Physical address? Base=0x4800, Physical addr = 0x4850,
   Load Byte from 0x4850→$t0, Move PC+4→PC

# Observations about Segmentation

- **Virtual address space has holes**
  - Segmentation efficient for sparse address spaces
  - A correct program should never address gaps (except as mentioned in moment)
    - » If it does, trap to kernel and dump core
- **When it is OK to address outside valid range:**
  - This is how the stack and heap are allowed to grow
  - For instance, stack takes fault, system automatically increases size of stack
- **Need protection mode in segment table**
  - For example, code segment would be read-only
  - Data and stack would be read-write (stores allowed)
  - Shared segment could be read-only or read-write
- **What must be saved/restored on context switch?**
  - Segment table stored in CPU, not in memory (small)
  - Might store all of processes memory onto disk when switched (called "swapping")

# Schematic View of Swapping



- **Extreme form of Context Switch: Swapping**
  - **In order to make room for next process, some or all of the previous process is moved to disk**
    - » **Likely need to send out complete segments**
  - **This greatly increases the cost of context-switching**
- **Desirable alternative?**
  - **Some way to keep only active portions of a process in memory at any one time**
  - **Need finer granularity control over physical memory**

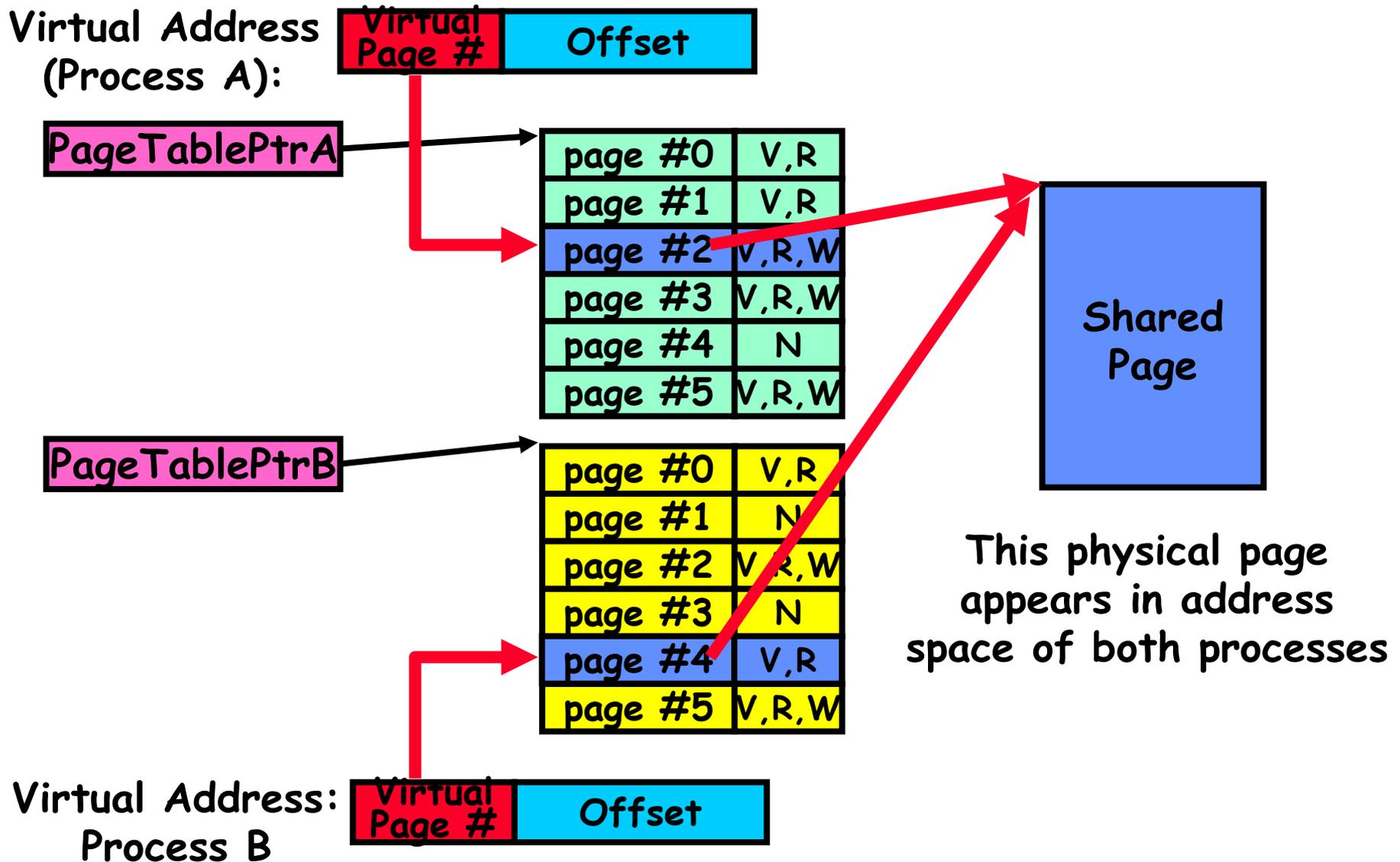# Paging: Physical Memory in Fixed Size Chunks

- **Problems with segmentation?**
  - Must fit variable-sized chunks into physical memory
  - May move processes multiple times to fit everything
  - Limited options for swapping to disk
- **Fragmentation: wasted space**
  - **External:** free gaps between allocated chunks
  - **Internal:** don't need all memory within allocated chunks
- **Solution to fragmentation from segments?**
  - Allocate physical memory in fixed size chunks ("pages")
  - Every chunk of physical memory is equivalent
    - » Can use simple vector of bits to handle allocation:
      00110001110001101 ... 110010
    - » Each bit represents page of physical memory
      $1 \Rightarrow$ allocated, $0 \Rightarrow$ free
- **Should pages be as big as our previous segments?**
  - No: Can lead to lots of internal fragmentation
    - » Typically have small pages (1K-16K)
  - Consequently: need multiple pages/segment

# How to Implement Paging?

**Virtual Address:** | Virtual Page # | Offset |

PageTablePtr → | page #0 | V,R |
| page #1 | V,R |
| page #2 | V,R,W |
| page #3 | V,R,W |
| page #4 | N |
| page #5 | V,R,W |

PageTableSize → **>**

**Access Error**

Physical Page # | Offset

**Physical Address**

Check Perm

**Access Error**

- **Page Table (One per process)**
  - **Resides in physical memory**
  - **Contains physical page and permission for each virtual page**
    - » Permissions include: Valid bits, Read, Write, etc
- **Virtual address mapping**
  - **Offset from Virtual address copied to Physical Address**
    - » Example: 10 bit offset $\Rightarrow$ 1024-byte pages
  - **Virtual page # is all remaining bits**
    - » Example for 32-bits: 32-10 = 22 bits, i.e. 4 million entries
    - » Physical page # copied from table into physical address
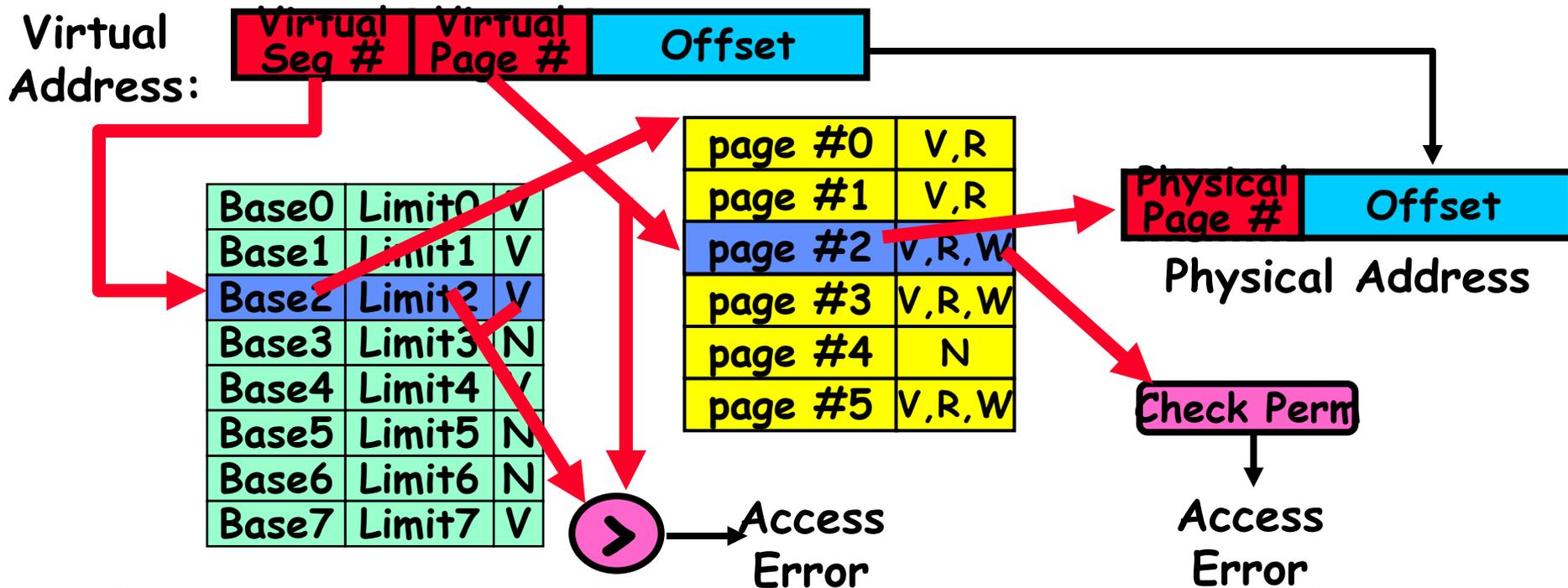  - **Check Page Table bounds and permissions**

# What about Sharing?

**Virtual Address (Process A):**

| Virtual Page # | Offset |
|---|---|

**PageTablePtrA**

| page #0 | V,R |
|---|---|
| page #1 | V,R |
| page #2 | V,R,W |
| page #3 | V,R,W |
| page #4 | N |
| page #5 | V,R,W |

**PageTablePtrB**

| page #0 | V,R |
|---|---|
| page #1 | N |
| page #2 | V,R,W |
| page #3 | N |
| page #4 | V,R |
| page #5 | V,R,W |

**Shared Page**

**This physical page appears in address space of both processes**

**Virtual Address: Process B**

| Virtual Page # | Offset |
|---|---|

# Simple Page Table Discussion

- **What needs to be switched on a context switch?**
  - Page table pointer and limit
- **Simple Page Table Analysis**
  - Pros
    - » Simple memory allocation
    - » Easy to Share
  - Con: What if address space is sparse?
    - » E.g. on UNIX, code starts at 0, stack starts at $(2^{31}-1)$.
    - » With 1K pages, need 4 million page table entries!
  - Con: What if table really big?
    - » Not all pages used all the time $\Rightarrow$ would be nice to have working set of page table in memory
- **How about combining paging and segmentation?**
  - Segments with pages inside them?
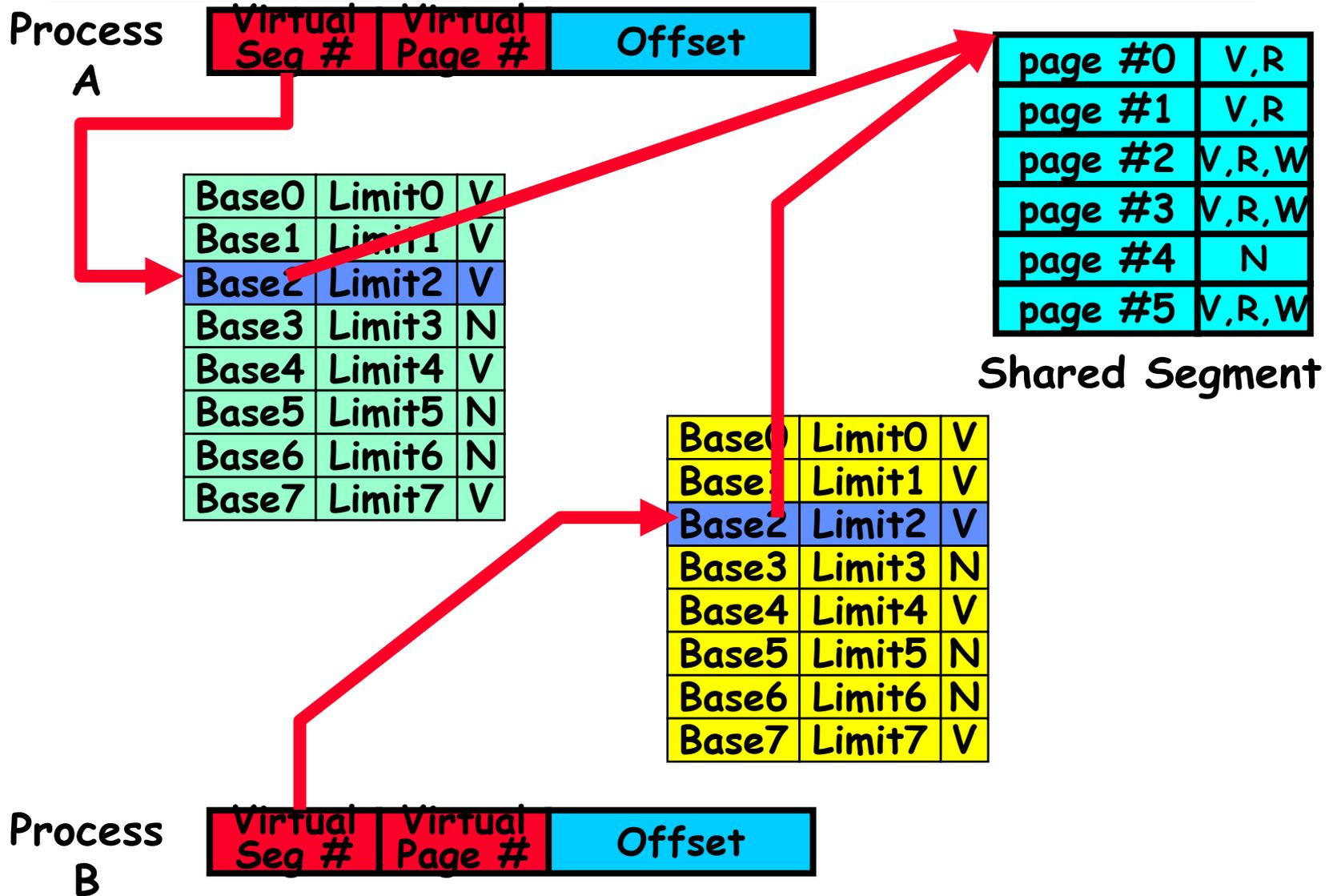  - Need some sort of multi-level translation

# Multi-level Translation: Segments + Pages

- **What about a tree of tables?**
  - Lowest level page table⇒memory still allocated with bitmap
  - Higher levels often segmented
- **Could have any number of levels. Example (top segment):**

Virtual Address:

| Virtual Seg # | Virtual Page # | Offset |
|---|---|---|

| Base0 | Limit0 | V |
|---|---|---|
| Base1 | Limit1 | V |
| Base2 | Limit2 | V |
| Base3 | Limit3 | N |
| Base4 | Limit4 | V |
| Base5 | Limit5 | N |
| Base6 | Limit6 | N |
| Base7 | Limit7 | V |

| page #0 | V,R |
|---|---|
| page #1 | V,R |
| page #2 | V,R,W |
| page #3 | V,R,W |
| page #4 | N |
| page #5 | V,R,W |

| Physical Page # | Offset |
|---|---|

**Physical Address**

**Check Perm**
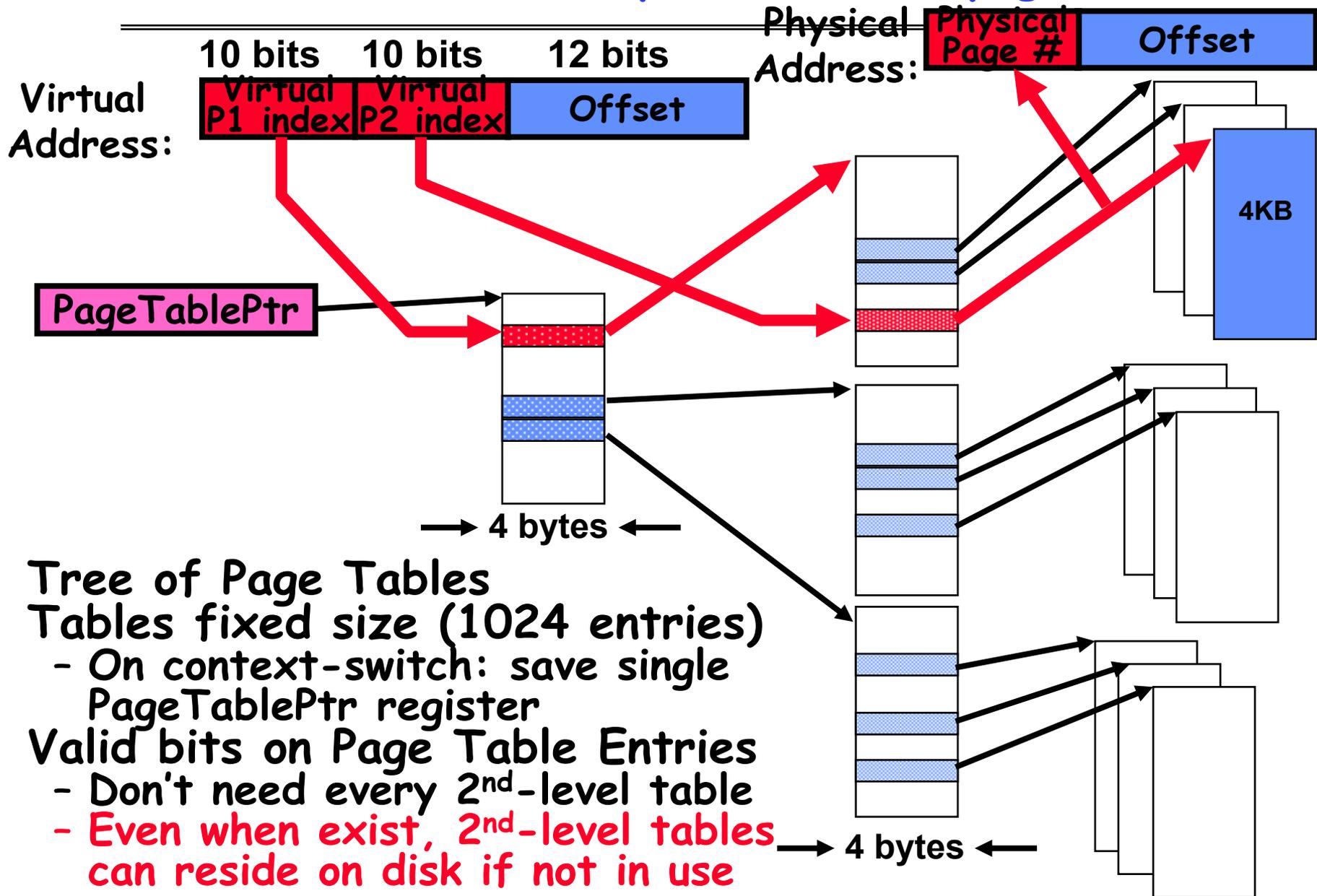
**>** → Access Error

Access Error

- **What must be saved/restored on context switch?**
  - Contents of top-level segment registers (for this example)
  - Pointer to top-level table (page table)

# What about Sharing (Complete Segment)?

Process A

| Virtual Seg # | Virtual Page # | Offset |
|---|---|---|

| Base0 | Limit0 | V |
|---|---|---|
| Base1 | Limit1 | V |
| Base2 | Limit2 | V |
| Base3 | Limit3 | N |
| Base4 | Limit4 | V |
| Base5 | Limit5 | N |
| Base6 | Limit6 | N |
| Base7 | Limit7 | V |

| page #0 | V,R |
|---|---|
| page #1 | V,R |
| page #2 | V,R,W |
| page #3 | V,R,W |
| page #4 | N |
| page #5 | V,R,W |

**Shared Segment**

| Base0 | Limit0 | V |
|---|---|---|
| Base1 | Limit1 | V |
| Base2 | Limit2 | V |
| Base3 | Limit3 | N |
| Base4 | Limit4 | V |
| Base5 | Limit5 | N |
| Base6 | Limit6 | N |
| Base7 | Limit7 | V |

Process B

| Virtual Seg # | Virtual Page # | Offset |
|---|---|---|

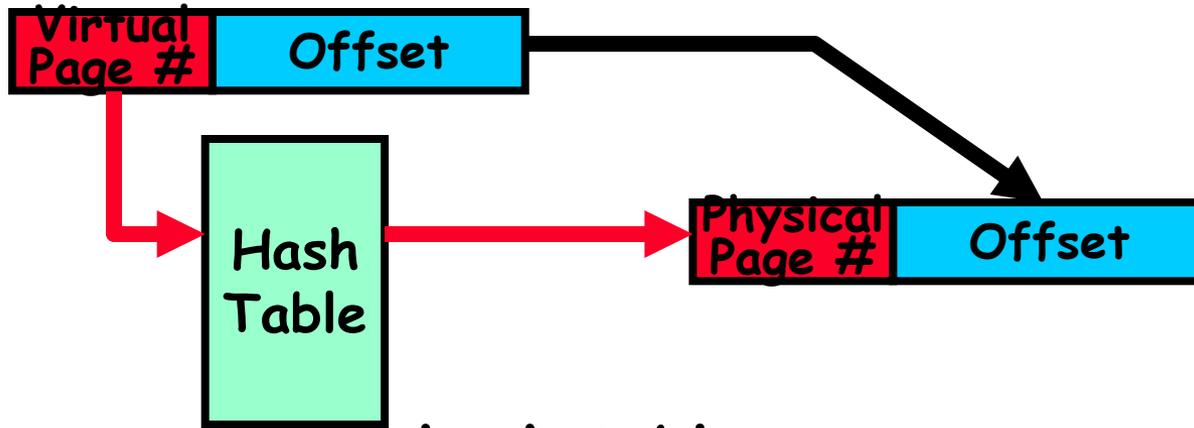# Another common example: two-level page table



- **Tree of Page Tables**
- **Tables fixed size (1024 entries)**
  - On context-switch: save single PageTablePtr register
- **Valid bits on Page Table Entries**
  - Don't need every 2nd-level table
  - Even when exist, 2nd-level tables can reside on disk if not in use

# Multi-level Translation Analysis

- Pros:
  - Only need to allocate as many page table entries as we need for application
    - » In other wards, sparse address spaces are easy
  - Easy memory allocation
  - Easy Sharing
    - » Share at segment or page level (need additional reference counting)
- Cons:
  - One pointer per page (typically 4K – 16K pages today)
  - Page tables need to be contiguous
    - » However, previous example keeps tables to exactly one page in size
  - Two (or more, if >2 levels) lookups per reference
    - » Seems very expensive!

# Inverted Page Table

- **With all previous examples ("Forward Page Tables")**
  - Size of page table is at least as large as amount of virtual memory allocated to processes
  - Physical memory may be much less
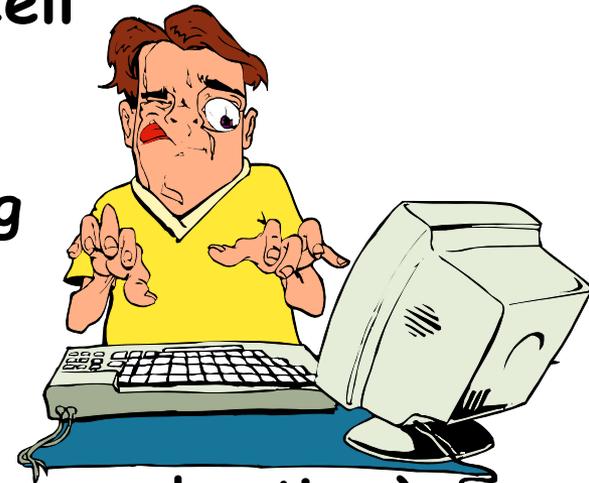    - » Much of process space may be out on disk or not in use



- **Answer: use a hash table**
  - Called an "Inverted Page Table"
  - Size is independent of virtual address space
  - Directly related to amount of physical memory
  - Very attractive option for 64-bit address spaces
- **Cons: Complexity of managing hash changes**
  - Often in hardware!

# Dual-Mode Operation

- **Can Application Modify its own translation tables?**
  - If it could, could get access to all of physical memory
  - Has to be restricted somehow
- **To Assist with Protection, <span style="color:red">Hardware</span> provides at least two modes (Dual-Mode Operation):**
  - "Kernel" mode (or "supervisor" or "protected")
  - "User" mode (Normal program mode)
  - Mode set with bits in special control register only accessible in kernel-mode
- **Intel processor actually has four "rings" of protection:**
  - PL (Priviledge Level) from 0 – 3
    - » PL0 has full access, PL3 has least
  - Privilege Level set in code segment descriptor (CS)
  - Mirrored "IOPL" bits in condition register gives permission to programs to use the I/O instructions
  - Typical OS kernels on Intel processors only use PL0 ("user") and PL3 ("kernel")

# For Protection, Lock User-Programs in Asylum

- **Idea: Lock user programs in padded cell with no exit or sharp objects**
  - Cannot change mode to kernel mode
  - User cannot modify page table mapping
  - Limited access to memory: cannot adversely effect other processes
    - » Side-effect: Limited access to memory-mapped I/O operations (I/O that occurs by reading/writing memory locations)
  - Limited access to interrupt controller
  - What else needs to be protected?
- **A couple of issues**
  - How to share CPU between kernel and user programs?
    - » Kinda like both the inmates and the warden in asylum are the same person.  How do you manage this???
  - How do programs interact?
  - How does one switch between kernel and user modes?
    - » OS $\rightarrow$ user (kernel $\rightarrow$ user mode): getting into cell
    - » User$\rightarrow$ OS (user $\rightarrow$ kernel mode): getting out of cell

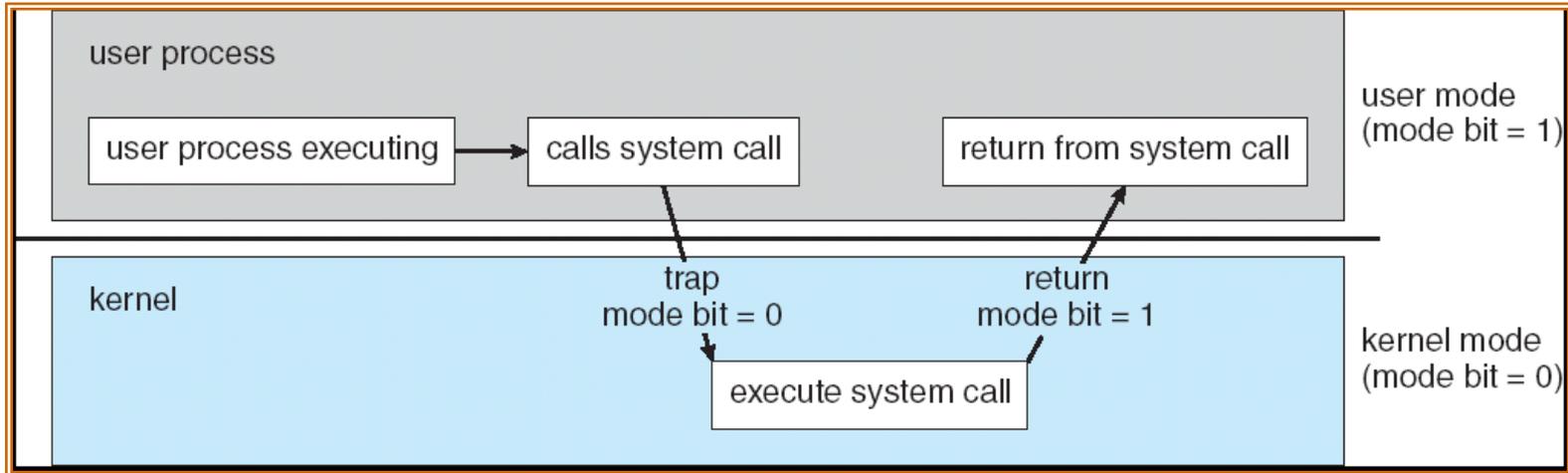# How to get from Kernel→User

- **What does the kernel do to create a new user process?**
    - **Allocate and initialize address-space control block**
    - **Read program off disk and store in memory**
    - **Allocate and initialize translation table**
        - » **Point at code in memory so program can execute**
        - » **Possibly point at statically initialized data**
    - **Run Program:**
        - » **Set machine registers**
        - » **Set hardware pointer to translation table**
        - » **Set processor status word for user mode**
        - » **Jump to start of program**
- **How does kernel switch between processes?**
    - **Same saving/restoring of registers as before**
    - **Save/restore PSL (hardware pointer to translation table)**

# User→Kernel (System Call)

- **Can't let inmate (user) get out of padded cell on own**
  - Would defeat purpose of protection!
  - So, how does the user program get back into kernel?



- **System call: Voluntary procedure call into kernel**
  - Hardware for controlled User→Kernel transition
  - Can any kernel routine be called?
    - » No! Only specific ones.
  - System call ID encoded into system call instruction
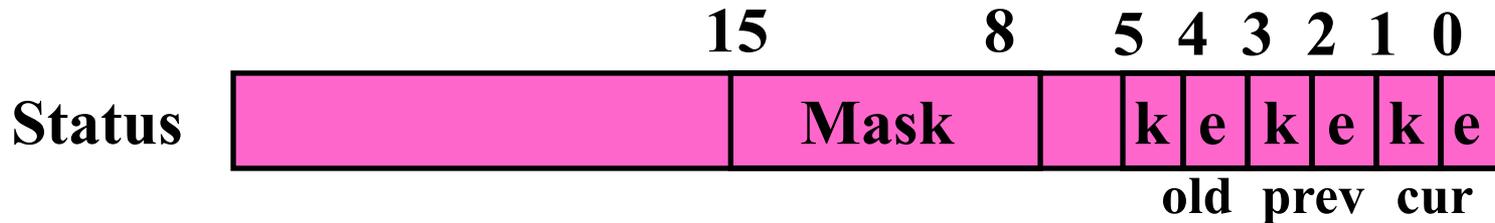    - » Index forces well-defined interface with kernel

# System Call Continued

- **What are some system calls?**
  - I/O: open, close, read, write, lseek
  - Files: delete, mkdir, rmdir, truncate, chown, chgrp, ..
  - Process: fork, exit, wait (like join)
  - Network: socket create, set options
- **Are system calls constant across operating systems?**
  - Not entirely, but there are lots of commonalities
  - Also some standardization attempts (POSIX)
- **What happens at beginning of system call?**
    - » On entry to kernel, sets system to kernel mode
    - » Handler address fetched from table/Handler started
- **System Call argument passing:**
  - In registers (not very much can be passed)
  - Write into user memory, kernel copies into kernel mem
    - » User addresses must be translated!w
    - » Kernel has different view of memory than user
  - Every Argument must be explicitly checked!

# User→Kernel (Exceptions: Traps and Interrupts)

- A system call instruction causes a synchronous exception (or "trap")
  - In fact, often called a software "trap" instruction
- Other sources of *Synchronous Exceptions:*
  - Divide by zero, Illegal instruction, Bus error (bad address, e.g. unaligned access)
  - Segmentation Fault (address out of range)
  - Page Fault (for illusion of infinite-sized memory)
- Interrupts are *Asynchronous Exceptions*
  - Examples: timer, disk ready, network, etc….
  - Interrupts can be disabled, traps cannot!
- On system call, exception, or interrupt:
  - Hardware enters kernel mode with interrupts disabled
  - Saves PC, then jumps to appropriate handler in kernel
  - For some processors (x86), processor also saves registers, changes stack, etc.
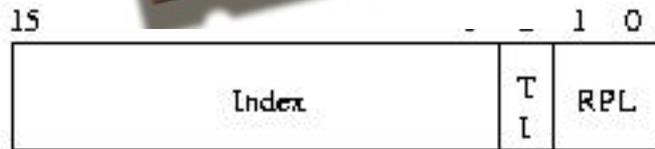- Actual handler typically saves registers, other CPU state, and switches to kernel stack

# Additions to MIPS ISA to support Exceptions?

- **Exception state is kept in "Coprocessor 0"**
  - **Use mfc0 read contents of these registers:**
    - » <span style="color:red">BadVAddr (register 8):</span> contains memory address at which memory reference error occurred
    - » <span style="color:red">Status (register 12):</span> interrupt mask and enable bits
    - » <span style="color:red">Cause (register 13):</span> the cause of the exception
    - » <span style="color:red">EPC (register 14):</span> address of the affected instruction

|  | **15** | **8** | **5** | **4** | **3** | **2** | **1** | **0** |
|---|---|---|---|---|---|---|---|---|
| **Status** | | **Mask** | | **k** | **e** | **k** | **e** | **k** | **e** |

                                                              old   prev   cur

- **Status Register fields:**
  - **Mask: Interrupt enable**
    - » 1 bit for each of 5 hardware and 3 software interrupts
  - **k = kernel/user:       0⇒kernel mode**
  - **e = interrupt enable: 0⇒interrupts disabled**
  - <span style="color:red">**Exception⇒6 LSB shifted left 2 bits, setting 2 LSB to 0:**</span>
    - » <span style="color:red">run in kernel mode with interrupts disabled</span>

# Intel x86 Special Registers

## 80386 Special Registers

Segment registers

| | | |
|---|---|---|
| Code Seg. | | Data Seg. |
| 15 CS 0 | | 15 DS 0 |
| Stack Seg. | | Extra Seg. |
| 15 SS 0 | | 15 ES 0 |
| Extra Seg. | | Extra. Seg |
| 15 ES 0 | | 15 GS 0 |

| X | N T | IO PL | O F | D F | I F | T F | S F | Z F | X | A F | X | P F | X | C F |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 15 | 14 | 13 | 12 | 11 | 10 | 9 | 8 | 7 | 6 | 5 | 4 | 3 | 2 | 1 0 |

| P G | | E T | T S | T S | M P | P E | CR0 |
|---|---|---|---|---|---|---|---|
| 31 30 | | 5 | 4 | 3 | 2 | 1 0 | |

| Unused | CR1 |
|---|---|
| 31 | 0 Flags |

| Page Fault Linear Address | CR2 |
|---|---|
| 31 | 0 |

| Page Directory Base Register | Not Used | CR3 |
|---|---|---|
| 31 | 7 | 0 |

PG=Paging Enable
ET=Emulation Type
TS=Task Switched
EM=Emulate Coprocessor
MP=Math coprocessor present
PE=Protected Mode enable

X=Reserved
NT=Nested Task
IOPL=I/O Privilege Level
OF=Overflow Flag
DF=Direction Flag
IF=Interrupt Flag
TF=Trap Flag
SF=Sign Flag
ZF=Zero Flag
AF=Auxiliary Flag
PF=Parity Flag
CF=Carry Flag

| 15 | | 1 0 |
|---|---|---|
| Index | T I | RPL |

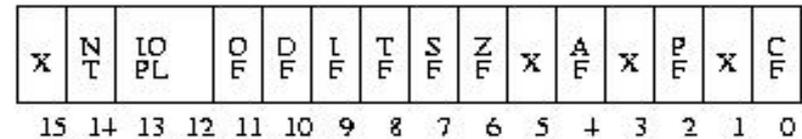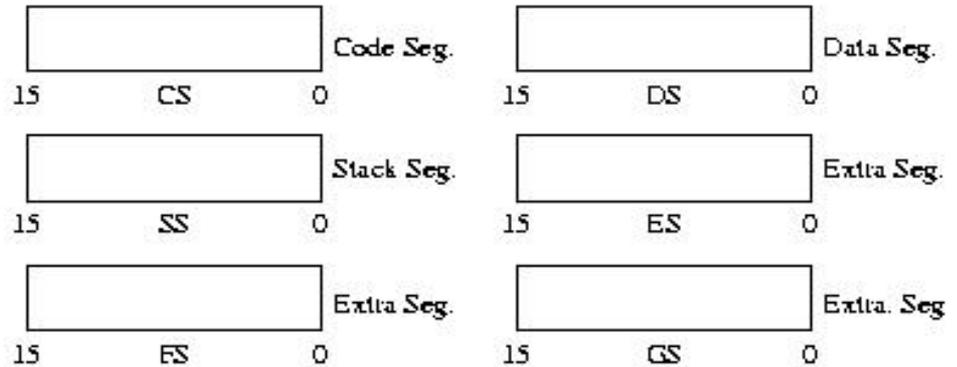RPL = Requestor Privilege Level
TI = Table Indicator
   (0 = GDT, 1 = LDT)
Index = Index into table

Protected Mode segment selector

**Typical Segment Register Current Priority is RPL Of Code Segment (CS)**

# Communication

- **Now that we have isolated processes, how can they communicate?**
  - **Shared memory: common mapping to physical page**
    - » As long as place objects in shared memory address range, threads from each process can communicate
    - » Note that processes A and B can talk to shared memory through different addresses
    - » In some sense, this violates the whole notion of protection that we have been developing
  - **If address spaces don't share memory, all inter-address space communication must go through kernel (via system calls)**
    - » Byte stream producer/consumer (put/get): Example, communicate through pipes connecting `stdin/stdout`
    - » Message passing (send/receive): Will explain later how you can use this to build remote procedure call (RPC) abstraction so that you can have one program make procedure calls to another
    - » File System (read/write): File system is shared state!

# Closing thought: Protection without Hardware

- **Does protection require hardware support for translation and dual-mode behavior?**
  - No: Normally use hardware, but anything you can do in hardware can also do in software (possibly expensive)
- **Protection via Strong Typing**
  - Restrict programming language so that you can't express program that would trash another program
  - Loader needs to make sure that program produced by valid compiler or all bets are off
  - Example languages: LISP, Ada, Modula-3 and Java
- **Protection via software fault isolation:**
  - Language independent approach: have compiler generate object code that provably can't step out of bounds
    - » Compiler puts in checks for every "dangerous" operation (loads, stores, etc). Again, need special loader.
    - » Alternative, compiler generates "proof" that code cannot do certain things (Proof Carrying Code)
  - Or: use virtual machine to guarantee safe behavior (loads and stores recompiled on fly to check bounds)

# Summary (1/2)

- **Memory is a resource that must be shared**
  - Controlled Overlap: only shared when appropriate
  - Translation: Change Virtual Addresses into Physical Addresses
  - Protection: Prevent unauthorized Sharing of resources
- **Simple Protection through Segmentation**
  - Base+limit registers restrict memory accessible to user
  - Can be used to translate as well
- **Full translation of addresses through Memory Management Unit (MMU)**
  - Every Access translated through page table
  - Changing of page tables only available to user
- **Dual-Mode**
  - Kernel/User distinction: User restricted
  - User$\rightarrow$Kernel: System calls, Traps, or Interrupts
  - Inter-process communication: shared memory, or through kernel (system calls)
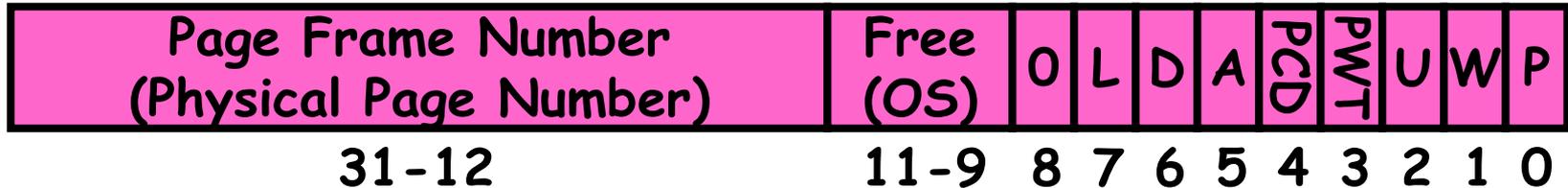
# Summary (2/2)

- **Segment Mapping**
  - **Segment registers within processor**
  - **Segment ID associated with each access**
    - » Often comes from portion of virtual address
    - » Can come from bits in instruction instead (x86)
  - **Each segment contains base and limit information**
    - » Offset (rest of address) adjusted by adding base
- **Page Tables**
  - **Memory divided into fixed-sized chunks of memory**
  - **Virtual page number from virtual address mapped through page table to physical page number**
  - **Offset of virtual address same as physical address**
  - **Large page tables can be placed into virtual memory**
- **Multi-Level Tables**
  - **Virtual address mapped to series of tables**
  - **Permit sparse population of address space**
- **Inverted page table**
  - **Size of page table related to physical memory size**

# Caching and Virtual Memory

# What is in a PTE?

- **What is in a Page Table Entry (or PTE)?**
  - Pointer to next-level page table or to actual page
  - Permission bits: valid, read-only, read-write, write-only
- **Example: Intel x86 architecture PTE:**
  - Address same format previous slide (10, 10, 12-bit offset)
  - Intermediate page tables called "Directories"

| Page Frame Number (Physical Page Number) | Free (OS) | 0 | L | D | A | PCD | PWT | U | W | P |
|---|---|---|---|---|---|---|---|---|---|---|
| 31–12 | 11–9 | 8 | 7 | 6 | 5 | 4 | 3 | 2 | 1 | 0 |

P: Present (same as "valid" bit in other architectures)

W: Writeable

U: User accessible

PWT: Page write transparent: external cache write-through

PCD: Page cache disabled (page cannot be cached)

A: Accessed: page has been accessed recently

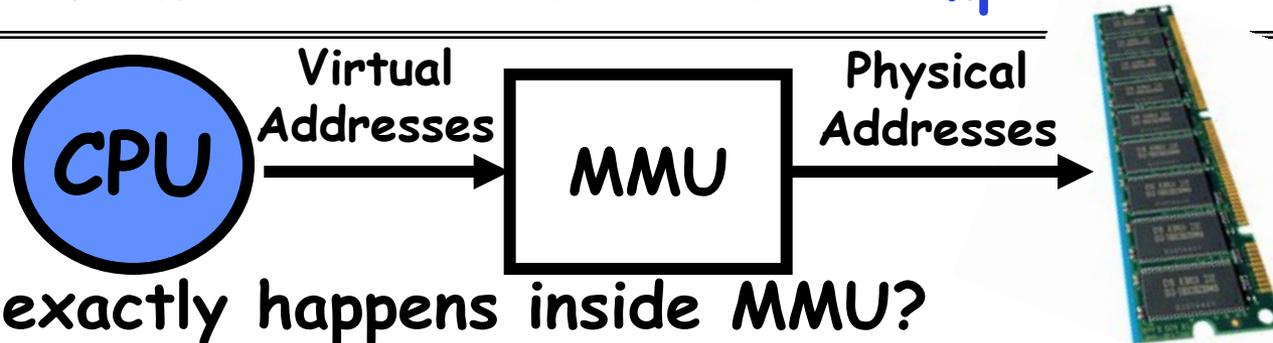D: Dirty (PTE only): page has been modified recently

L: L=1⇒4MB page (directory only).
Bottom 22 bits of virtual address serve as offset
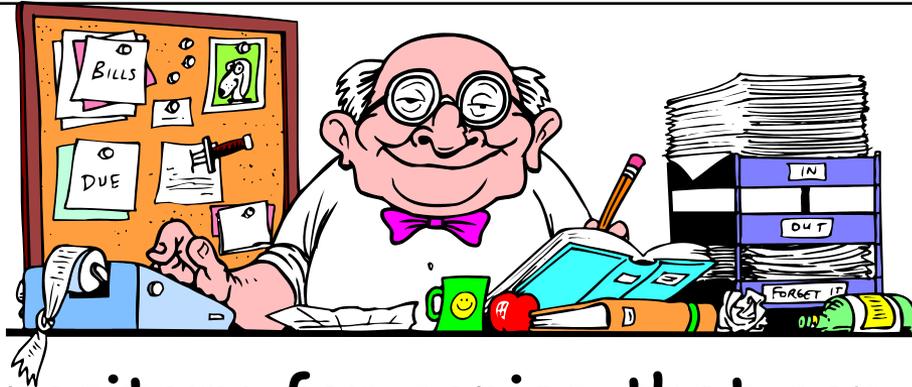
# Examples of how to use a PTE

- **How do we use the PTE?**
  - Invalid PTE can imply different things:
    - » Region of address space is actually invalid or
    - » Page/directory is just somewhere else than memory
  - Validity checked first
    - » OS can use other (say) 31 bits for location info
- **Usage Example: Demand Paging**
  - Keep only active pages in memory
  - Place others on disk and mark their PTEs invalid
- **Usage Example: Copy on Write**
  - UNIX fork gives *copy* of parent address space to child
    - » Address spaces disconnected after child created
  - How to do this cheaply?
    - » Make copy of parent's page tables (point at same memory)
    - » Mark entries in both sets of page tables as read-only
    - » Page fault on write creates two copies
- **Usage Example: Zero Fill On Demand**
  - New data pages must carry no information (say be zeroed)
  - Mark PTEs as invalid; page fault on use gets zeroed page
  - Often, OS creates zeroed pages in background

# How is the translation accomplished?



**CPU** → Virtual Addresses → **MMU** → Physical Addresses →

- **What, exactly happens inside MMU?**
- **One possibility: Hardware Tree Traversal**
  - For each virtual address, takes page table base pointer and traverses the page table in hardware
  - Generates a "Page Fault" if it encounters invalid PTE
    » Fault handler will decide what to do
    » More on this next lecture
  - Pros: Relatively fast (but still many memory accesses!)
  - Cons: Inflexible, Complex hardware
- **Another possibility: Software**
  - Each traversal done in software
  - Pros: Very flexible
  - Cons: Every translation must invoke Fault!
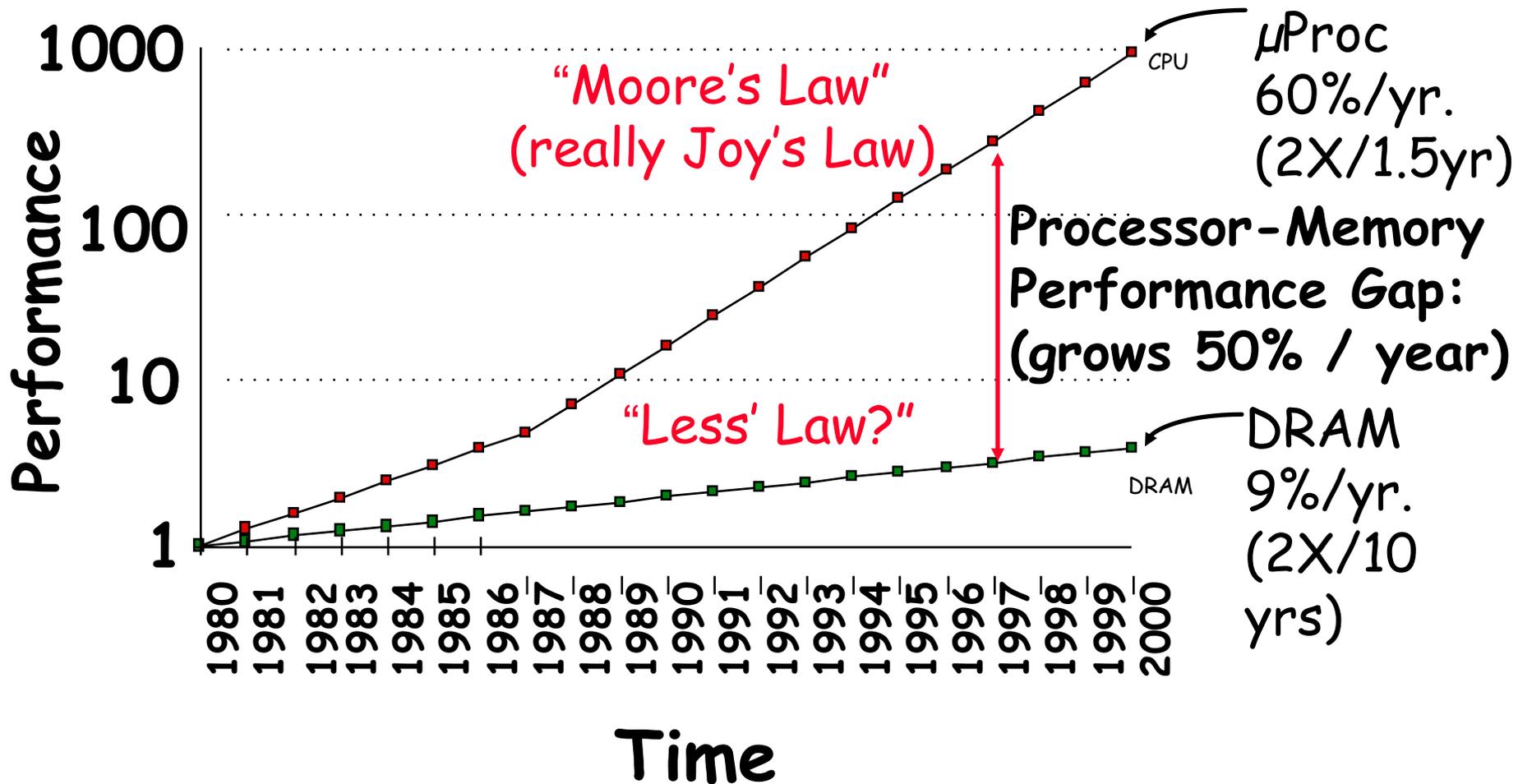- **In fact, need way to cache translations for either case!**
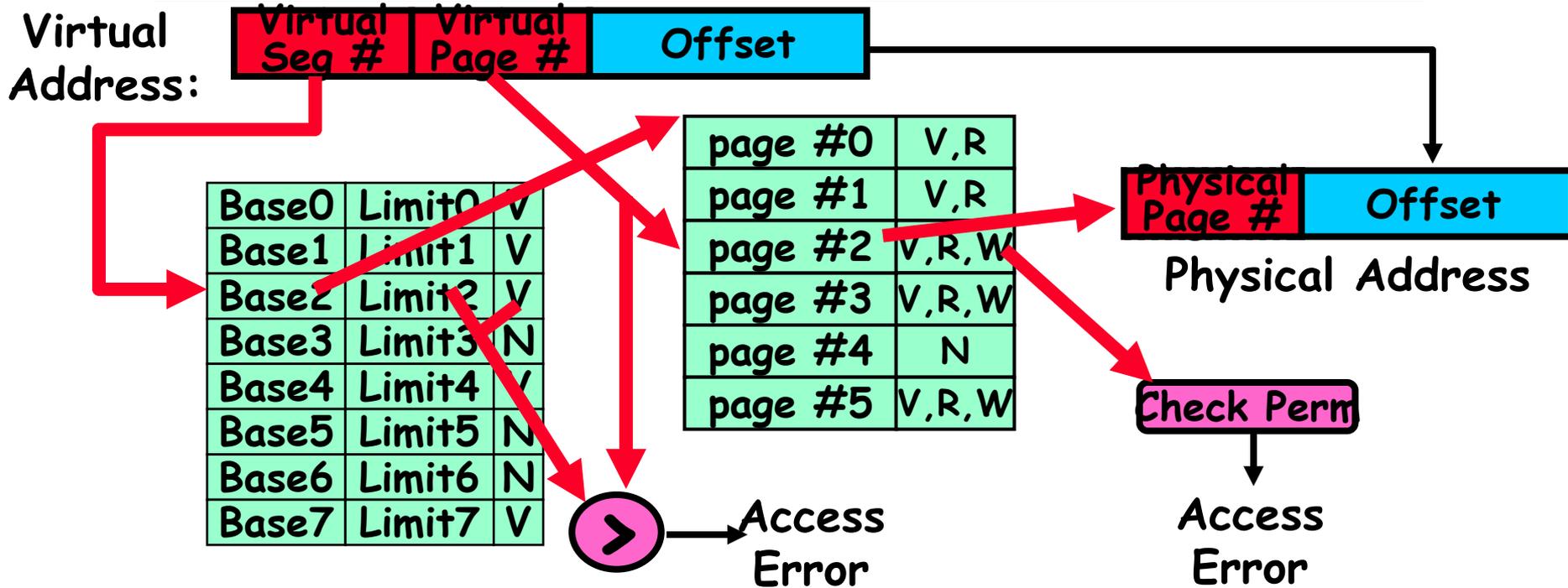
# Caching Concept

- **Cache**: a repository for copies that can be accessed more quickly than the original
  - Make frequent case fast and infrequent case less dominant
- Caching underlies many of the techniques that are used today to make computers fast
  - Can cache: memory locations, address translations, pages, file blocks, file names, network routes, etc...
- Only good if:
  - Frequent case frequent enough and
  - Infrequent case not too expensive
- Important measure: Average Access time =
    (Hit Rate x **Hit Time**) + (Miss Rate x **Miss Time**)

# Why Bother with Caching?
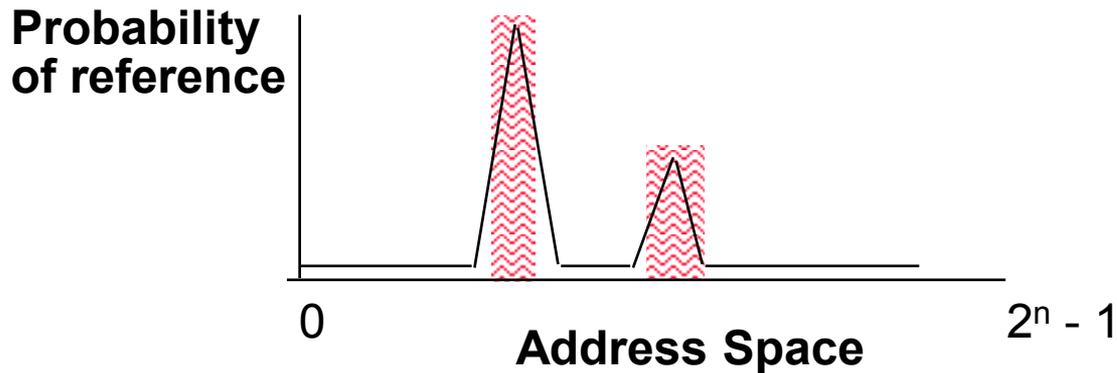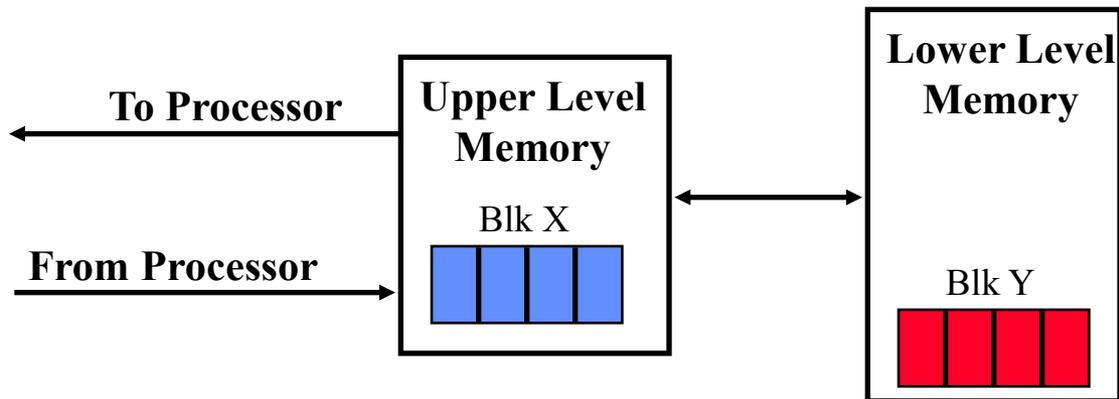
## Processor-DRAM Memory Gap (latency)



A graph titled "Processor-DRAM Memory Gap (latency)". The vertical axis is labeled "Performance" with a logarithmic scale from 1, 10, 100, to 1000. The horizontal axis is labeled "Time" with years from 1980 to 2000.

- CPU line labeled "μProc 60%/yr. (2X/1.5yr)" with "Moore's Law" (really Joy's Law)"
- DRAM line labeled "DRAM 9%/yr. (2X/10 yrs)" with "Less' Law?"
- "Processor-Memory Performance Gap: (grows 50% / year)"

# Another Major Reason to Deal with Caching

**Virtual Address:**

| Virtual Seg # | Virtual Page # | Offset |
|---|---|---|

| | | |
|---|---|---|
| Base0 | Limit0 | V |
| Base1 | Limit1 | V |
| Base2 | Limit2 | V |
| Base3 | Limit3 | N |
| Base4 | Limit4 | V |
| Base5 | Limit5 | N |
| Base6 | Limit6 | N |
| Base7 | Limit7 | V |

| | |
|---|---|
| page #0 | V,R |
| page #1 | V,R |
| page #2 | V,R,W |
| page #3 | V,R,W |
| page #4 | N |
| page #5 | V,R,W |

| Physical Page # | Offset |
|---|---|

**Physical Address**

> → Access Error

**Check Perm**

Access Error

- **Cannot afford to translate on every access**
  - At least three DRAM accesses per actual DRAM access
  - Or: perhaps I/O if page table partially on disk!
- **Even worse: What if we are using caching to make memory access faster than DRAM access???**
- **Solution? Cache translations!**
  - Translation Cache: TLB ("Translation Lookaside Buffer")

# Why Does Caching Help? Locality!

**Probability of reference**

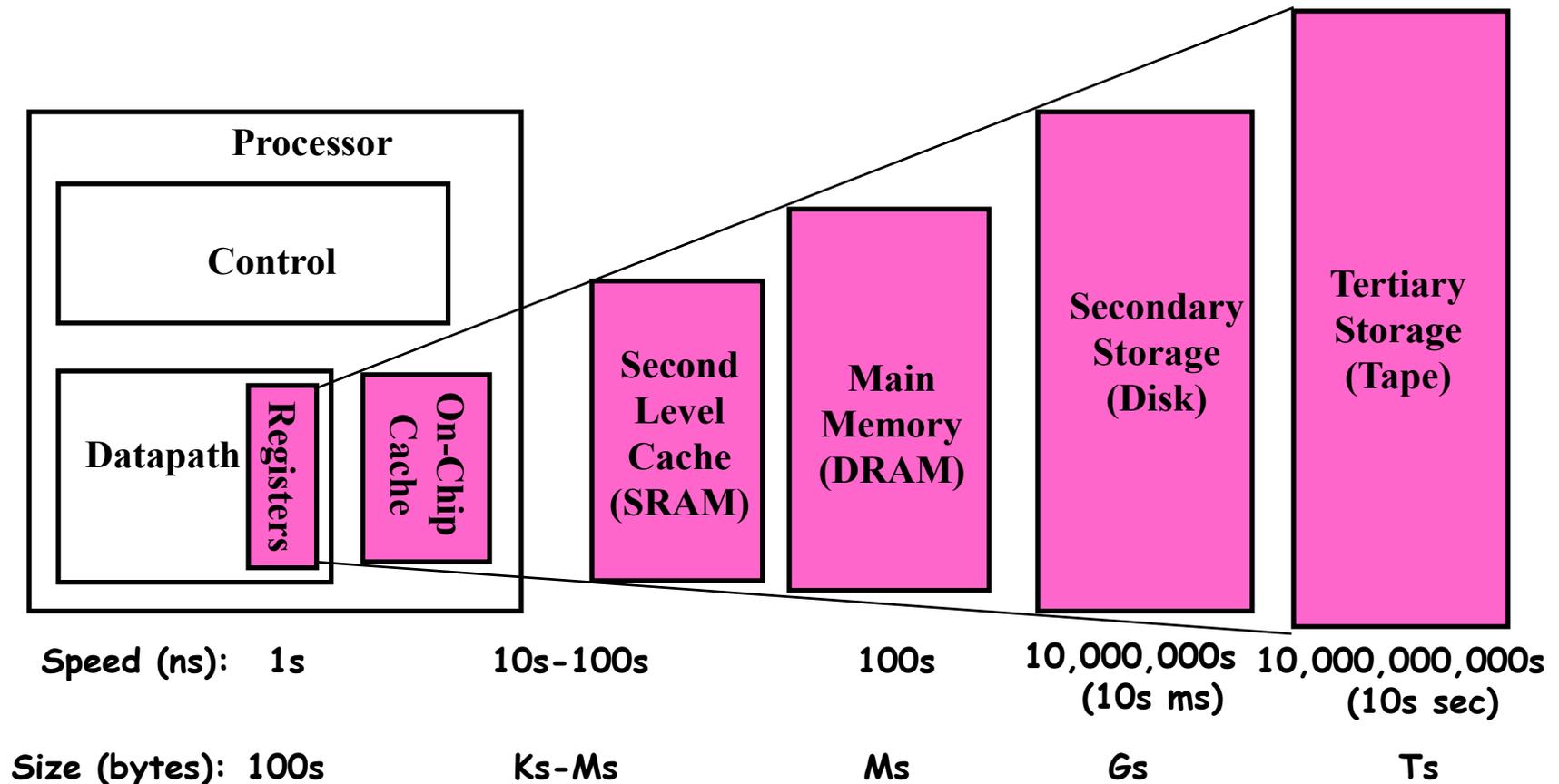**0**        **Address Space**        $2^n - 1$

- **Temporal Locality** (Locality in Time):
    - Keep recently accessed data items closer to processor
- **Spatial Locality** (Locality in Space):
    - Move contiguous blocks to the upper levels

To Processor

From Processor

**Upper Level Memory**

Blk X

**Lower Level Memory**

Blk Y

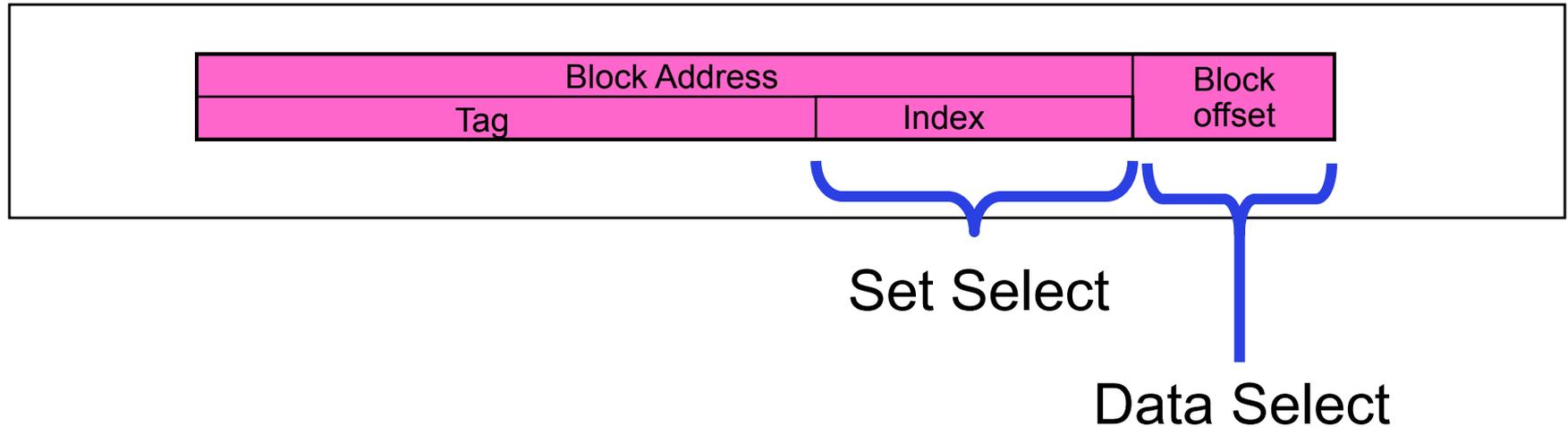# Memory Hierarchy of a Modern Computer System

- Take advantage of the principle of locality to:
  - Present as much memory as in the cheapest technology
  - Provide access at speed offered by the fastest technology

| | Processor | | | Second Level Cache (SRAM) | Main Memory (DRAM) | Secondary Storage (Disk) | Tertiary Storage (Tape) |
|---|---|---|---|---|---|---|---|
| | Control | | | | | | |
| | Datapath | Registers | On-Chip Cache | | | | |
| **Speed (ns):** | 1s | | 10s-100s | | 100s | 10,000,000s (10s ms) | 10,000,000,000s (10s sec) |
| **Size (bytes):** | 100s | | Ks-Ms | | Ms | Gs | Ts |

# A Summary on Sources of Cache Misses

- **Compulsory** (cold start or process migration, first reference): first access to a block
  - "Cold" fact of life: not a whole lot you can do about it
  - Note: If you are going to run "billions" of instruction, Compulsory Misses are insignificant
- **Capacity**:
  - Cache cannot contain all blocks access by the program
  - Solution: increase cache size
- **Conflict** (collision):
  - Multiple memory locations mapped to the same cache location
  - Solution 1: increase cache size
  - Solution 2: increase associativity
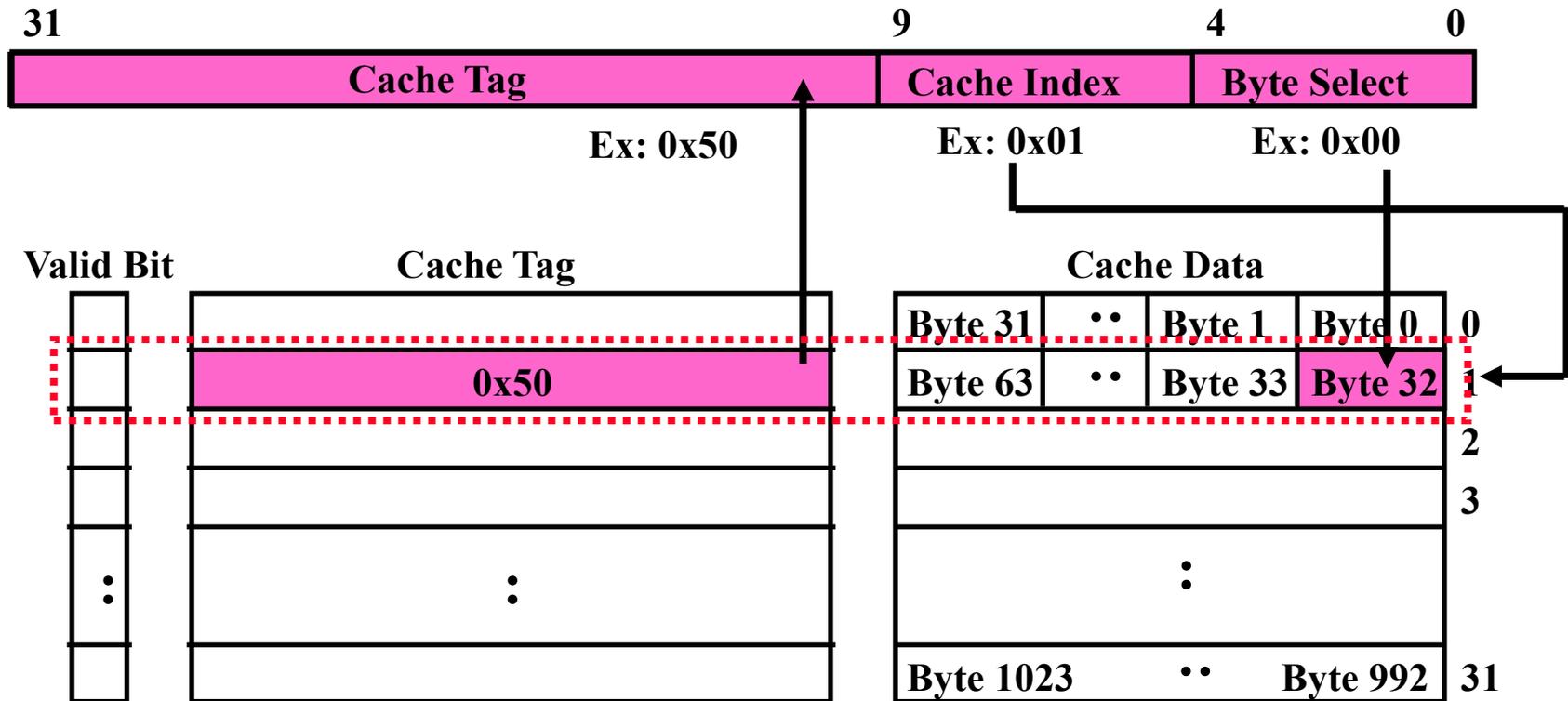- **Coherence** (Invalidation): other process (e.g., I/O) updates memory

# How is a Block found in a Cache?

| Block Address | | Block offset |
|---|---|---|
| Tag | Index | |

Set Select

Data Select

- **Index Used to Lookup Candidates in Cache**
  - **Index identifies the set**
- **Tag used to identify actual copy**
  - **If no candidates match, then declare cache miss**
- **Block is minimum quantum of caching**
  - **Data select field used to select data within block**
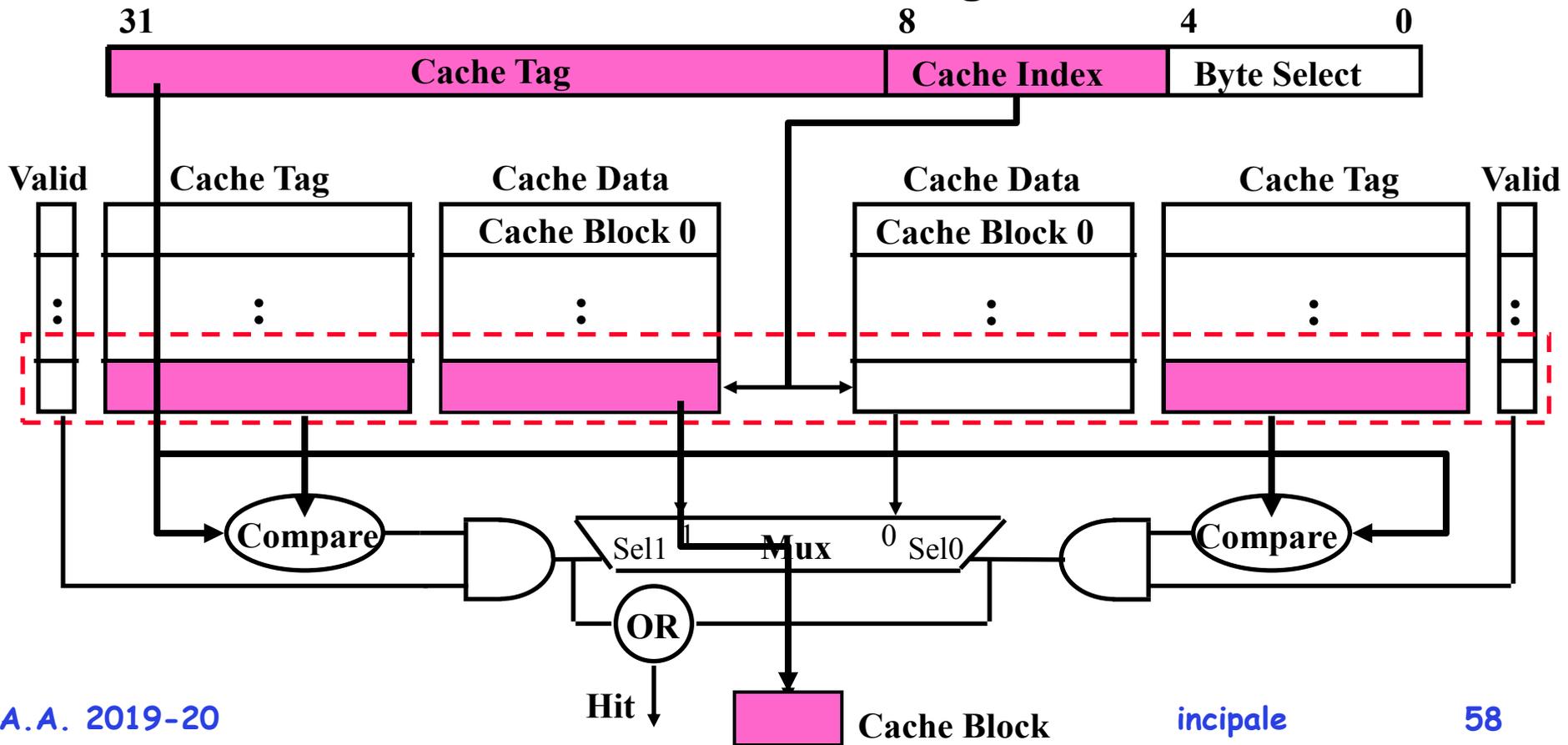  - **Many caching applications don't have data select field**

# Direct Mapped Cache

- **Direct Mapped $2^N$ byte cache**:
  - **The uppermost (32 - N) bits are always the Cache Tag**
  - **The lowest M bits are the Byte Select (Block Size = $2^M$)**
- **Example: 1 KB Direct Mapped Cache with 32 B Blocks**
  - **Index chooses potential block**
  - **Tag checked to verify block**
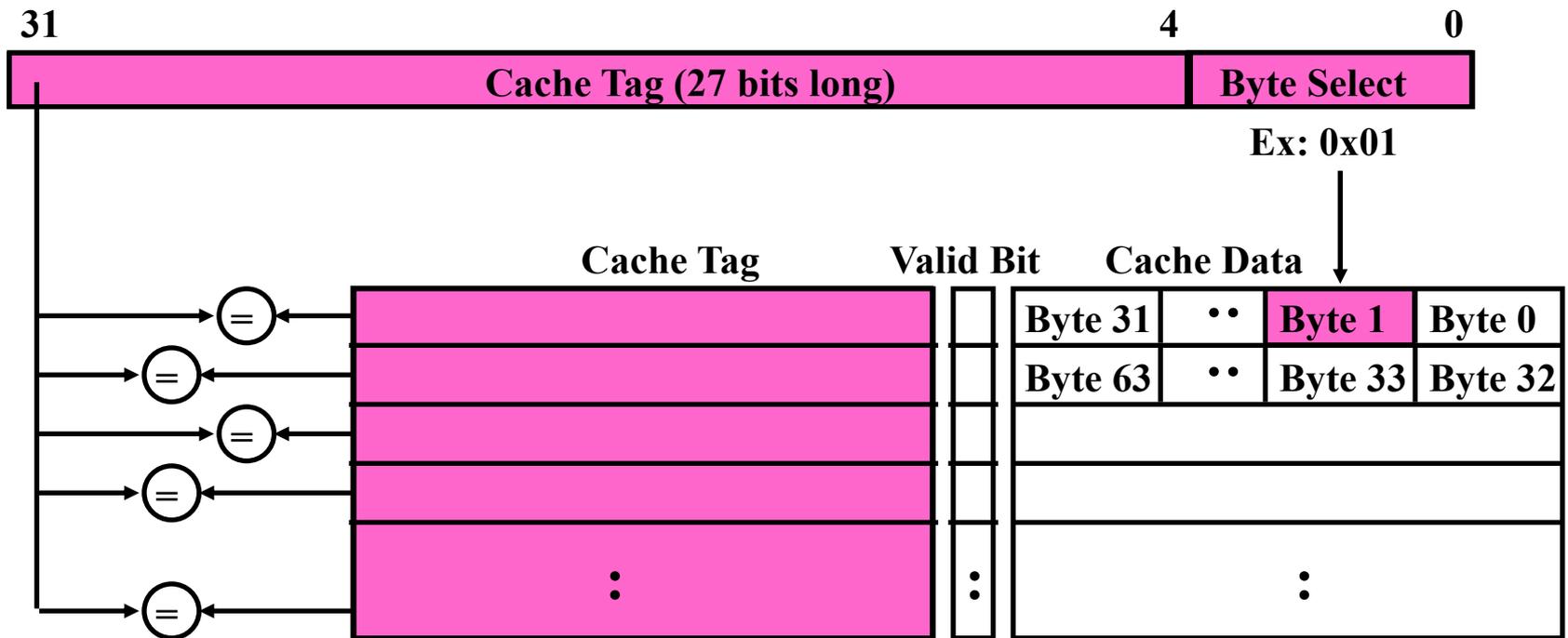  - **Byte select chooses byte within block**

| 31 | | 9 | 4 | 0 |
|---|---|---|---|---|
| | **Cache Tag** | | **Cache Index** | **Byte Select** |

Ex: 0x50          Ex: 0x01          Ex: 0x00

| Valid Bit | Cache Tag | Cache Data | |
|---|---|---|---|
| | | Byte 31 ·· Byte 1 Byte 0 | 0 |
| | **0x50** | Byte 63 ·· Byte 33 **Byte 32** | 1 |
| | | | 2 |
| | | | 3 |
| : | : | : | |
| | | Byte 1023 ·· Byte 992 | 31 |

# Set Associative Cache

- **N-way set associative**: N entries per Cache Index
  - N direct mapped caches operates in parallel
- Example: Two-way set associative cache
  - Cache Index selects a "set" from the cache
  - Two tags in the set are compared to input in parallel
  - Data is selected based on the tag result
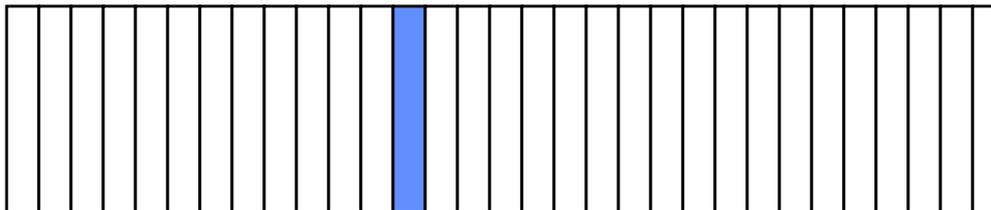
# Fully Associative Cache

- **Fully Associative**: **Every block can hold any line**
  - **Address does not include a cache index**
  - **Compare Cache Tags of all Cache Entries in Parallel**
- **Example: Block Size=32B blocks**
  - **We need N 27-bit comparators**
  - **Still have byte select to choose from within block**

| 31 | 4 | 0 |
|---|---|---|
| **Cache Tag (27 bits long)** | **Byte Select** | |

**Ex: 0x01**

**Cache Tag**  **Valid Bit**  **Cache Data**

| Byte 31 | •• | Byte 1 | Byte 0 |
|---|---|---|---|
| Byte 63 | •• | Byte 33 | Byte 32 |

# Where does a Block Get Placed in a Cache?

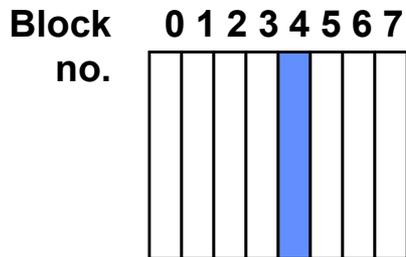- **Example: Block 12 placed in 8 block cache**

**32-Block Address Space:**

Block no.
1 1 1 1 1 1 1 1 1 1 2 2 2 2 2 2 2 2 2 2 3 3
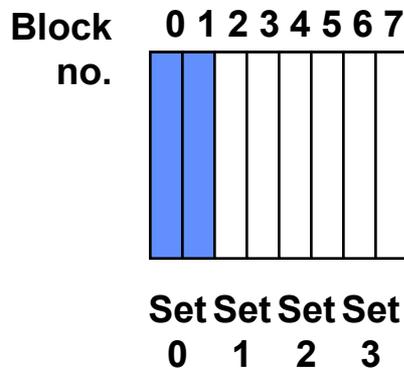0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1

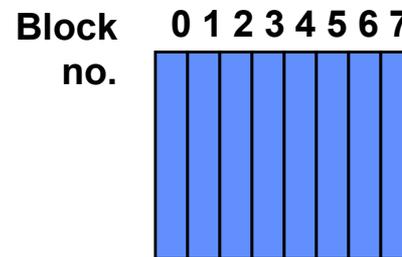**Direct mapped:**
block 12 can go
only into block 4
(12 mod 8)

Block no.    0 1 2 3 4 5 6 7

**Set associative:**
block 12 can go
anywhere in set 0
(12 mod 4)

Block no.    0 1 2 3 4 5 6 7

Set Set Set Set
0    1    2    3

**Fully associative:**
block 12 can go
anywhere

Block no.    0 1 2 3 4 5 6 7

# Review: Which block should be replaced on a miss?

- **Easy for Direct Mapped: Only one possibility**
- **Set Associative or Fully Associative:**
  - Random
  - LRU (Least Recently Used)

| | 2-way | | 4-way | | 8-way | |
|---|---|---|---|---|---|---|
| Size | LRU | Random | LRU | Random | LRU | Random |
| 16 KB | 5.2% | 5.7% | 4.7% | 5.3% | 4.4% | 5.0% |
| 64 KB | 1.9% | 2.0% | 1.5% | 1.7% | 1.4% | 1.5% |
| 256 KB | 1.15% | 1.17% | 1.13% | 1.13% | 1.12% | 1.12% |

# Review: What happens on a write?

- **Write through**: The information is written to both the block in the cache and to the block in the lower-level memory
- **Write back**: The information is written only to the block in the cache.
  - Modified cache block is written to main memory only when it is replaced
  - Question is block clean or dirty?
- Pros and Cons of each?
  - WT:
    - » PRO: read misses cannot result in writes
    - » CON: Processor held up on writes unless writes buffered
  - WB:
    - » PRO: repeated writes not sent to DRAM processor not held up on writes
    - » CON: More complex Read miss may require writeback of dirty data

# Cache performance

- **Miss-oriented Approach to Memory Access:**

$$CPUtime = IC \times \left( CPI_{Execution} + \frac{MemAccess}{Inst} \times MissRate \times MissPenalty \right) \times CycleTime$$

- **Separating out Memory component entirely**

  – **AMAT = Average Memory Access Time**

$$CPUtime = IC \times \left( CPI_{AluOps} + \frac{MemAccess}{Inst} \times AMAT \right) \times CycleTime$$

$$AMAT = HitRate \times HitTime + MissRate \times MissTime$$
$$= HitTime + MissRate \times MissPenalty$$
$$= Frac_{Inst} \times \left( HitTime_{Inst} + MissRate_{Inst} \times MissPenalty_{Inst} \right) +$$
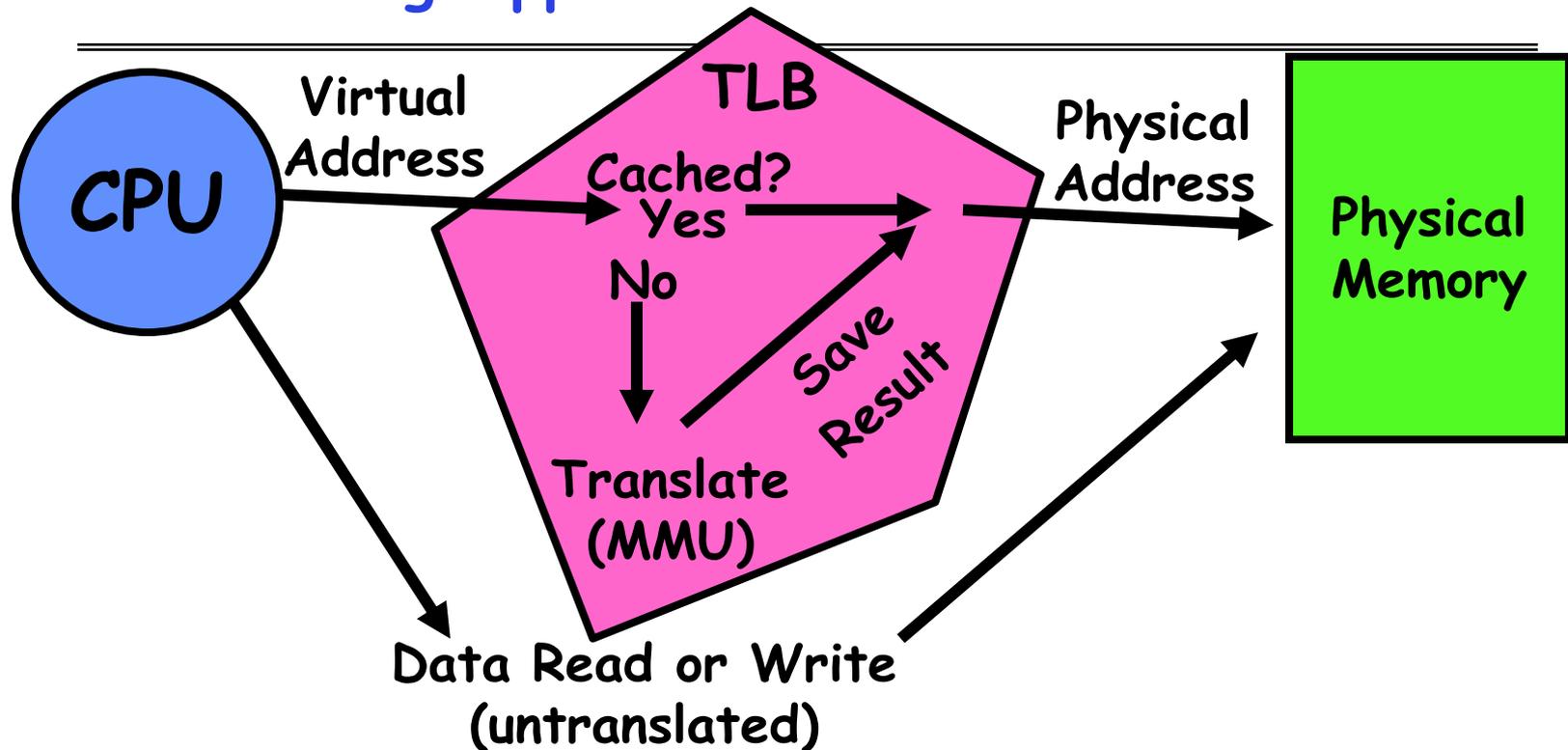$$Frac_{Data} \times \left( HitTime_{Data} + MissRate_{Data} \times MissPenalty_{Data} \right)$$

- **AMAT for Second-Level Cache**

$$AMAT_{1st} = HitTime_{1st} + MissRate_{1st} \times MissPenalty_{1st}$$
$$= HitTime_{1st} + MissRate_{1st} \times AMAT_{2nd}$$
$$= HitTime_{1st} + MissRate_{1st} \times \left( HitTime_{2st} + MissRate_{2st} \times MissPenalty_{2st} \right)$$

# Caching Applied to Address Translation

**CPU** → Virtual Address → **TLB** Cached? Yes → Physical Address → **Physical Memory**

No → Translate (MMU) → Save Result

CPU → Data Read or Write (untranslated) → Physical Memory

- **Question is one of page locality: does it exist?**
  - Instruction accesses spend a lot of time on the same page (since accesses sequential)
  - Stack accesses have definite locality of reference
  - Data accesses have less page locality, but still some…
- **Can we have a TLB hierarchy?**
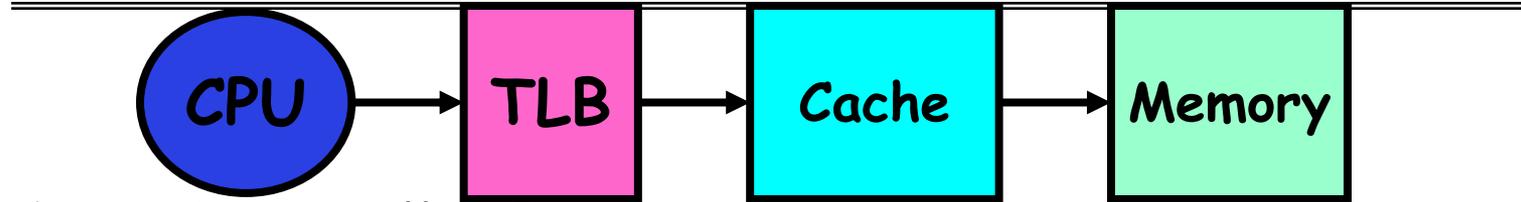  - Sure: multiple levels at different sizes/speeds

# What Actually Happens on a TLB Miss?

- **Hardware traversed page tables:**
  - On TLB miss, hardware in MMU looks at current page table to fill TLB (may walk multiple levels)
    - » If PTE valid, hardware fills TLB and processor never knows
    - » If PTE marked as invalid, causes Page Fault, after which kernel decides what to do afterwards
- **Software traversed Page tables (like MIPS)**
  - On TLB miss, processor receives TLB fault
  - Kernel traverses page table to find PTE
    - » If PTE valid, fills TLB and returns from fault
    - » If PTE marked as invalid, internally calls Page Fault handler
- **Most chip sets provide hardware traversal**
  - Modern operating systems tend to have more TLB faults since they use translation for many things
  - Examples:
    - » shared segments
    - » user-level portions of an operating system

# What happens on a Context Switch?

- **Need to do something, since TLBs map virtual addresses to physical addresses**
  - Address Space just changed, so TLB entries no longer valid!
- **Options?**
  - Invalidate TLB: simple but might be expensive
    - » What if switching frequently between processes?
  - Include ProcessID in TLB
    - » This is an architectural solution: needs hardware
- **What if translation tables change?**
  - For example, to move page from memory to disk or vice versa…
  - Must invalidate TLB entry!
    - » Otherwise, might think that page is still in memory!

# What TLB organization makes sense?

CPU → TLB → Cache → Memory

- **Needs to be really fast**
  - **Critical path of memory access**
    - » In simplest view: before the cache
    - » Thus, this adds to access time (reducing cache speed)
  - **Seems to argue for Direct Mapped or Low Associativity**
- **However, needs to have very few conflicts!**
  - **With TLB, the Miss Time extremely high!**
  - **This argues that cost of Conflict (Miss Time) is much higher than slightly increased cost of access (Hit Time)**
- **Thrashing: continuous conflicts between accesses**
  - **What if use low order bits of page as index into TLB?**
    - » First page of code, data, stack may map to same entry
    - » Need 3-way associativity at least?
  - **What if use high order bits as index?**
    - » TLB mostly unused for small programs

# TLB organization: include protection

- **How big does TLB actually have to be?**
  - Usually small: 128-512 entries
  - Not very big, can support higher associativity
- **TLB usually organized as fully-associative cache**
  - Lookup is by Virtual Address
  - Returns Physical Address + other info
- **What happens when fully-associative is too slow?**
  - Put a small (4-16 entry) direct-mapped cache in front
  - Called a "TLB Slice"
- **Example for MIPS R3000:**

| Virtual Address | Physical Address | Dirty | Ref | Valid | Access | ASID |
|---|---|---|---|---|---|---|
| 0xFA00 | 0x0003 | Y | N | Y | R/W | 34 |
| 0x0040 | 0x0010 | N | Y | Y | R | 0 |
| 0x0041 | 0x0011 | N | Y | Y | R | 0 |

# Example: R3000 pipeline includes TLB "stages"

**MIPS R3000 Pipeline**

| Inst Fetch | | Dcd/ Reg | | ALU / E.A | Memory | Write Reg |
|---|---|---|---|---|---|---|
| TLB | I-Cache | RF | | Operation | | WB |
| | | | | E.A. | TLB | D-Cache |

**TLB**

    **64 entry, on-chip, fully associative, software TLB fault handler**

**Virtual Address Space**

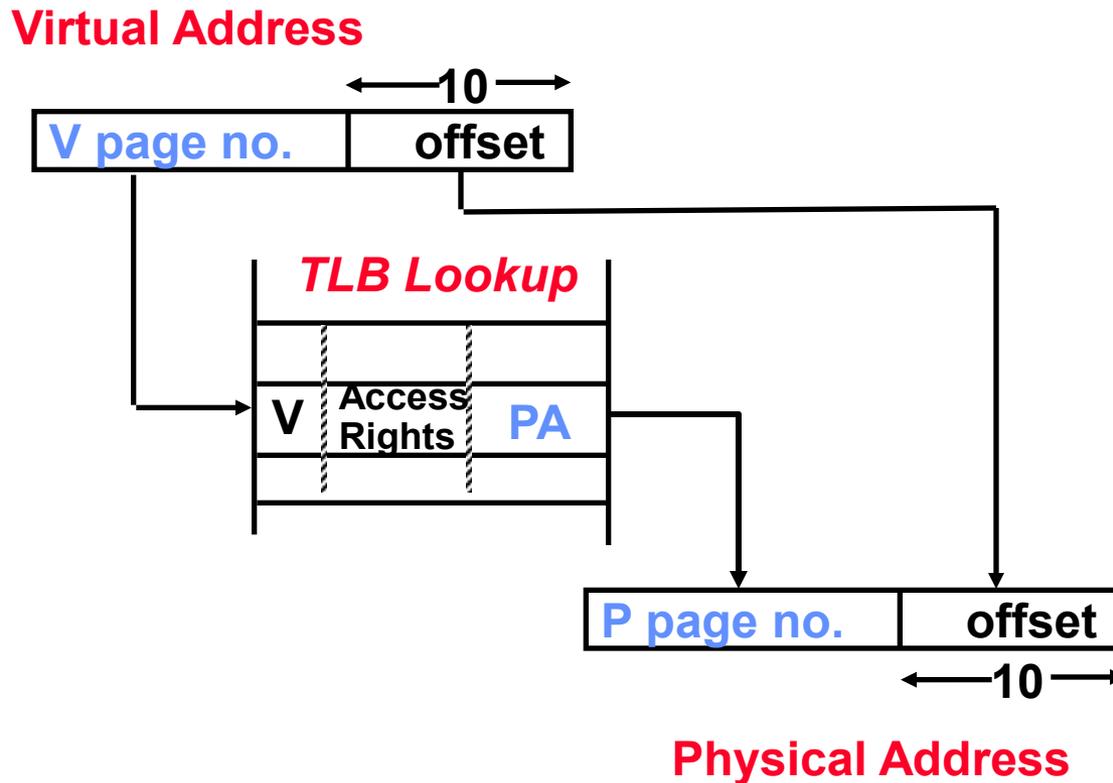| ASID | | V. Page Number | Offset |
|---|---|---|---|
| 6 | | 20 | 12 |

**0xx User segment (caching based on PT/TLB entry)**
**100 Kernel physical space, cached**
**101 Kernel physical space, uncached**
**11x Kernel virtual space**

**Allows context switching among**
**64 user processes without TLB flush**

> Combination
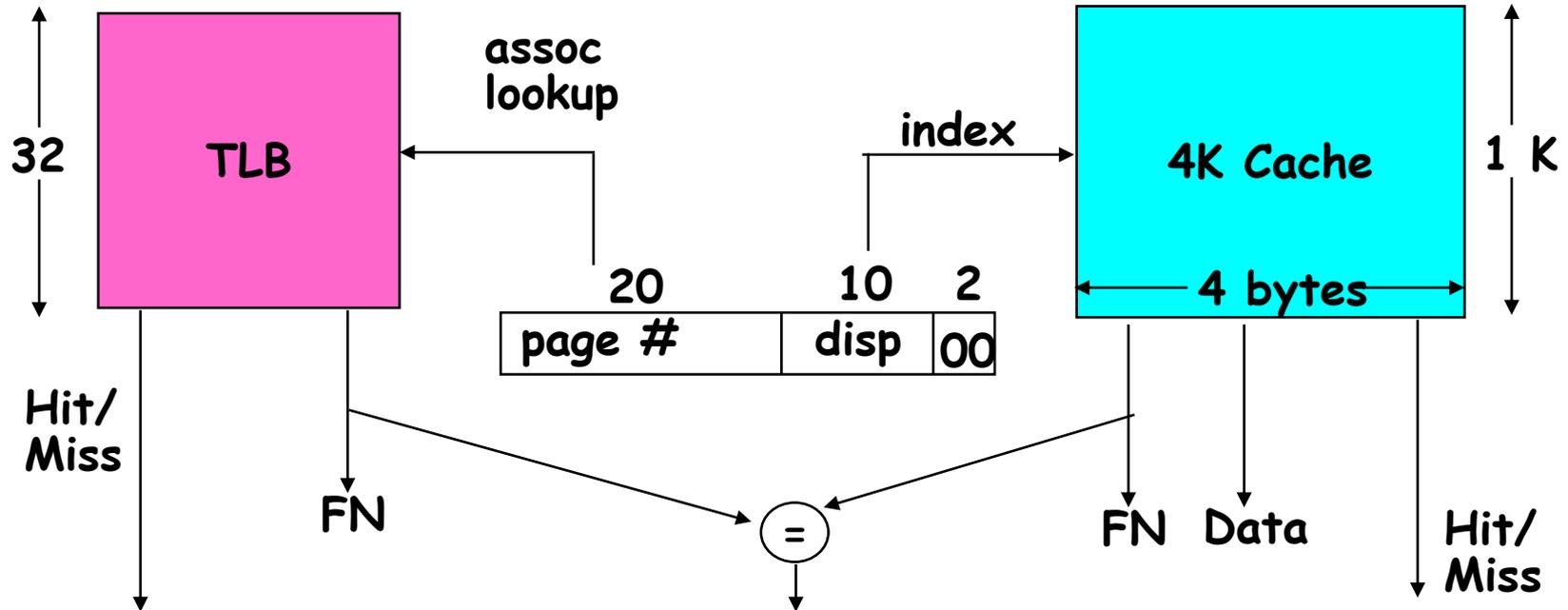> Segments and
> Paging!

# Reducing translation time further

- **As described, TLB lookup is in serial with cache lookup:**

**Virtual Address**

←— **10** —→

| **V page no.** | **offset** |

***TLB Lookup***

| **V** | **Access Rights** | **PA** |

| **P page no.** | **offset** |

←— **10** —→

**Physical Address**

- **Machines with TLBs go one step further: they overlap TLB lookup with cache access.**
  - **Works because offset available early**

# Overlapping TLB & Cache Access

- **Here is how this might work with a 4K cache:**



| | | |
|---|---|---|
| 20 | 10 | 2 |
| page # | disp | 00 |

- **What if cache size is increased to 8KB?**
  - Overlap not complete
  - Need to do something else.  See CS152/252
- **Another option: Virtual Caches**
  - Tags in cache are virtual addresses
  - Translation only happens on cache misses
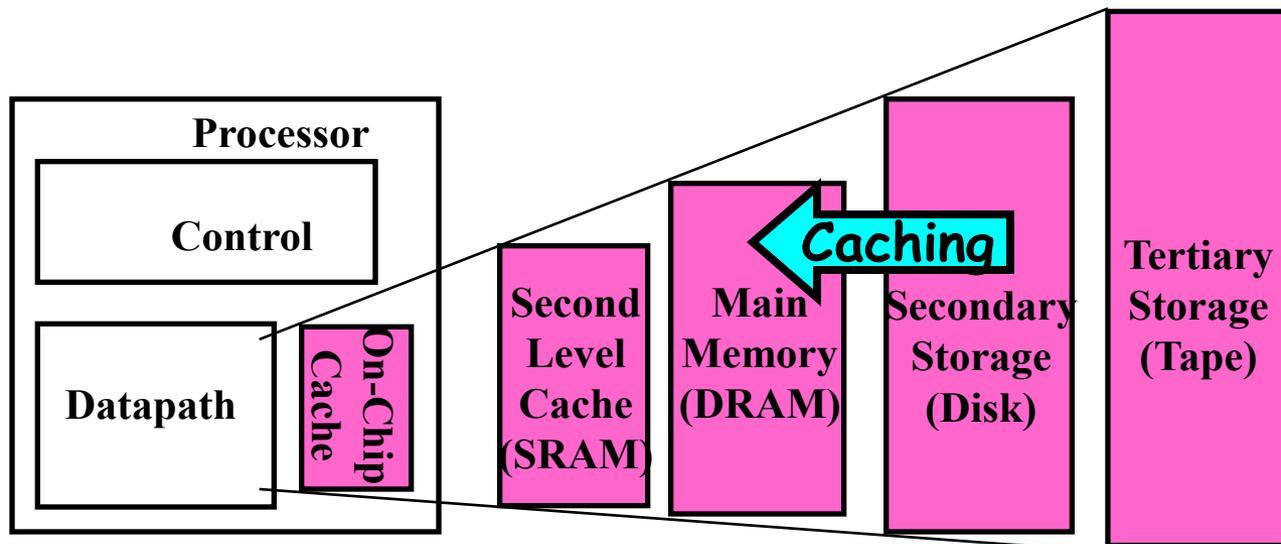
# Summary (1/2): Caching

- **The Principle of Locality:**
  - Program likely to access a relatively small portion of the address space at any instant of time.
    - » **Temporal Locality**: Locality in Time
    - » **Spatial Locality**: Locality in Space
- **Three (+1) Major Categories of Cache Misses:**
  - **Compulsory Misses**: sad facts of life.  Example: cold start misses.
  - **Conflict Misses**: increase cache size and/or associativity
  - **Capacity Misses**: increase cache size
  - **Coherence Misses**: Caused by external processors or I/O devices
- **Cache Organizations:**
  - Direct Mapped: single block per set
  - Set associative: more than one block per set
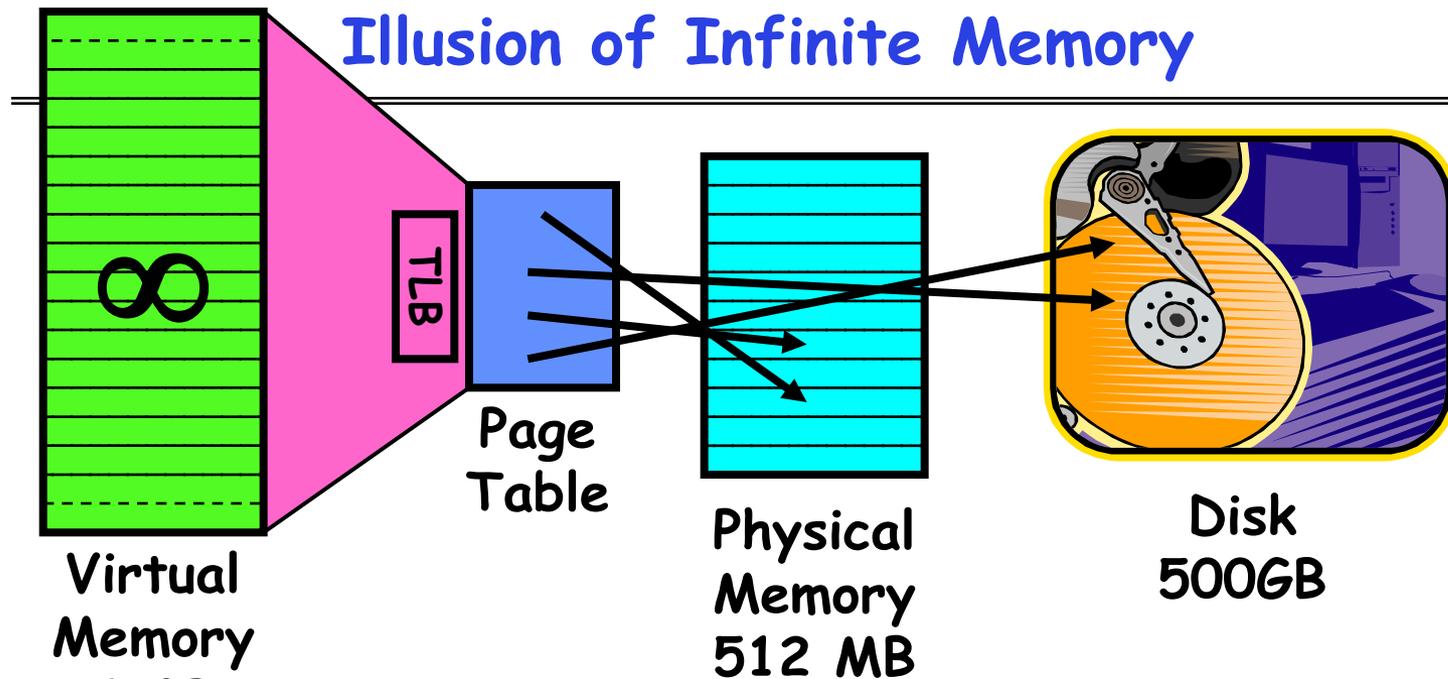  - Fully associative: all entries equivalent

- **PTE: Page Table Entries**
  - Includes physical page number
  - Control info (valid bit, writeable, dirty, user, etc)
- **A cache of translations called a "Translation Lookaside Buffer" (TLB)**
  - Relatively small number of entries (< 512)
  - Fully Associative (Since conflict misses expensive)
  - TLB entries contain PTE and optional process ID
- **On TLB miss, page table must be traversed**
  - If located PTE is invalid, cause Page Fault
- **On context switch/change in page table**
  - TLB entries must be invalidated somehow
- **TLB is logically in front of cache**
  - Thus, needs to be overlapped with cache access to be really fast

# Demand Paging

- **Modern programs require a lot of physical memory**
  - Memory per system growing faster than 25%-30%/year
- **But they don't use all their memory all of the time**
  - 90-10 rule: programs spend 90% of their time in 10% of their code
  - Wasteful to require all of user's code to be in memory
- **Solution: use main memory as cache for disk**

# Illusion of Infinite Memory



**TLB**

**Page Table**

**Virtual Memory 4 GB**

**Physical Memory 512 MB**

**Disk 500GB**

- Disk is larger than physical memory ⇒
  - In-use virtual memory can be bigger than physical memory
  - Combined memory of running processes much larger than physical memory
    » More programs fit into memory, allowing more concurrency
- Principle: Transparent Level of Indirection (page table)
  - Supports flexible placement of physical data
    » Data could be on disk or somewhere across network
  - Variable location of data transparent to user program
    » Performance issue, not correctness issue

# Demand Paging is Caching

- **Since Demand Paging is Caching, must ask:**
  - **What is block size?**
    - » 1 page
  - **What is organization of this cache (i.e. direct-mapped, set-associative, fully-associative)?**
    - » Fully associative: arbitrary virtual→physical mapping
  - **How do we find a page in the cache when look for it?**
    - » First check TLB, then page-table traversal
  - **What is page replacement policy? (i.e. LRU, Random…)**
    - » This requires more explanation… (kinda LRU)
  - **What happens on a miss?**
    - » Go to lower level to fill miss (i.e. disk)
  - **What happens on a write? (write-through, write back)**
    - » Definitely write-back.  Need dirty bit!
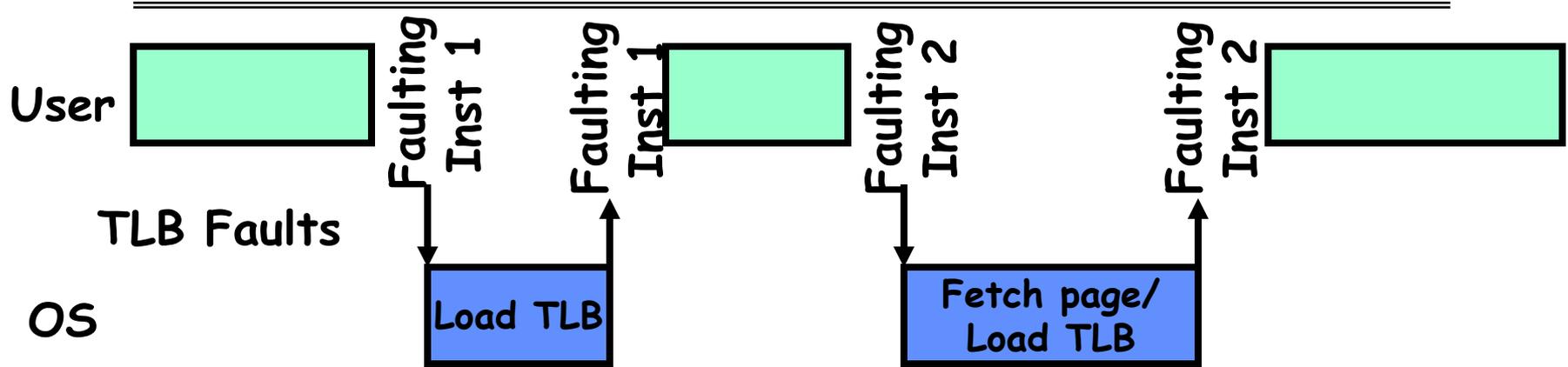
# Demand Paging Mechanisms

- **PTE helps us implement demand paging**
  - Valid $\Rightarrow$ Page in memory, PTE points at physical page
  - Not Valid $\Rightarrow$ Page not in memory; use info in PTE to find it on disk when necessary
- **Suppose user references page with invalid PTE?**
  - Memory Management Unit (MMU) traps to OS
    - » Resulting trap is a "Page Fault"
  - **What does OS do on a Page Fault?:**
    - » **Choose an old page to replace**
    - » **If old page modified ("D=1"), write contents back to disk**
    - » **Change its PTE and any cached TLB to be invalid**
    - » **Load new page into memory from disk**
    - » **Update page table entry, invalidate TLB for new entry**
    - » **Continue thread from original faulting location**
  - TLB for new page will be loaded when thread continued!
  - While pulling pages off disk for one process, OS runs another process from ready queue
    - » Suspended process sits on wait queue

# Software-Loaded TLB

- MIPS/Nachos TLB is loaded by software
  - High TLB hit rate⇒ok to trap to software to fill the TLB, even if slower
  - Simpler hardware and added flexibility: software can maintain translation tables in whatever convenient format
- How can a process run without access to page table?
  - Fast path (TLB hit with valid=1):
    » Translation to physical page done by hardware
  - Slow path (TLB hit with valid=0 or TLB miss)
    » Hardware receives a "TLB Fault"
  - What does OS do on a TLB Fault?
    » Traverse page table to find appropriate PTE
    » If valid=1, load page table entry into TLB, continue thread
    » If valid=0, perform "Page Fault" detailed previously
    » Continue thread
- Everything is transparent to the user process:
  - It doesn't know about paging to/from disk
  - It doesn't even know about software TLB handling

# Transparent Exceptions



- **How to transparently restart faulting instructions?**
  - Could we just skip it?
    » No: need to perform load or store after reconnecting physical page
- **Hardware must help out by saving:**
  - Faulting instruction and partial state
    » Need to know which instruction caused fault
    » Is single PC sufficient to identify faulting position????
  - Processor State: sufficient to restart user thread
    » Save/restore registers, stack, etc
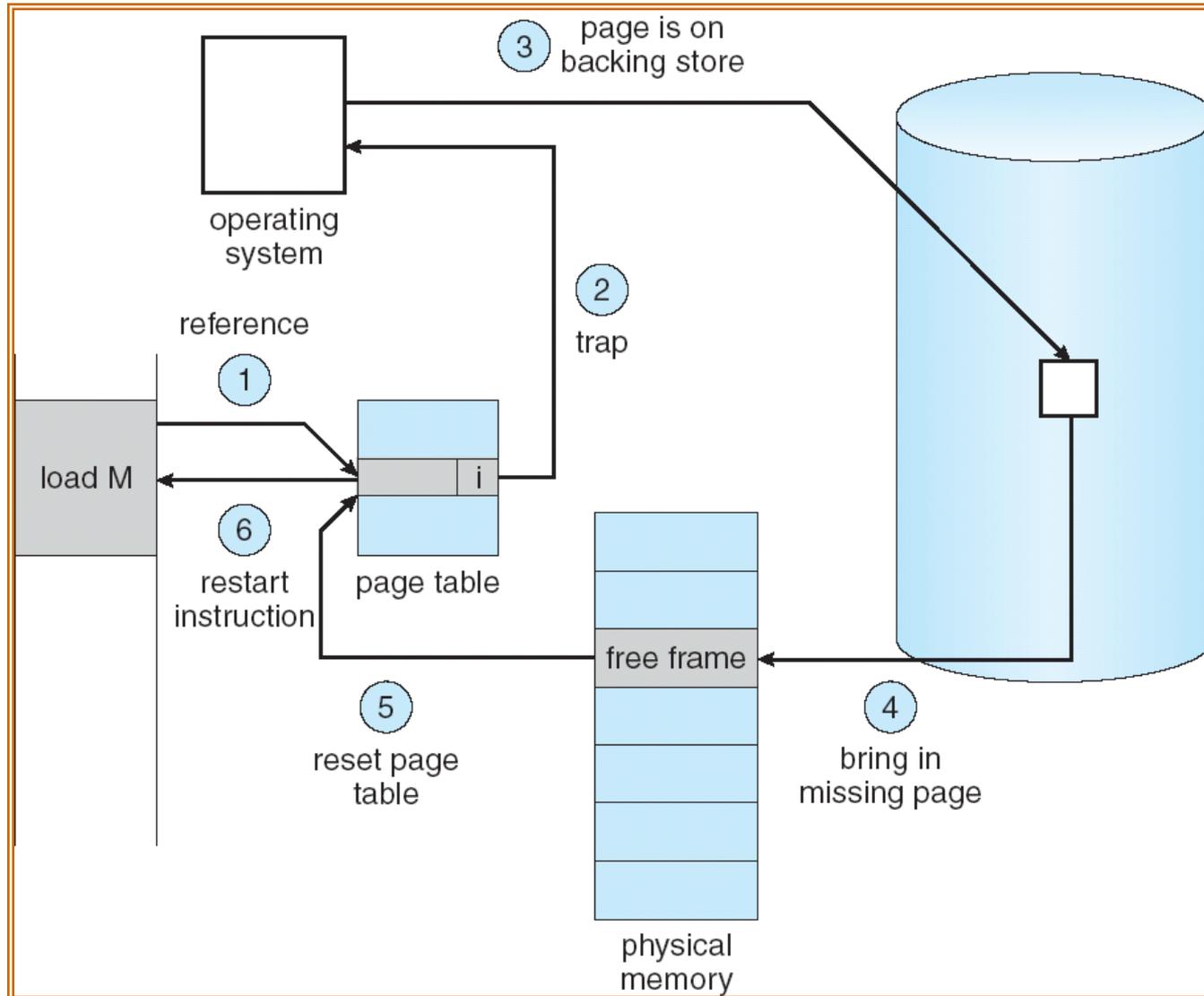- **What if an instruction has side-effects?**

# Consider weird things that can happen

- **What if an instruction has side effects?**
  - Options:
    - » Unwind side-effects (easy to restart)
    - » Finish off side-effects (messy!)
  - Example 1: `mov (sp)+,10`
    - » What if page fault occurs when write to stack pointer?
    - » Did `sp` get incremented before or after the page fault?
  - Example 2: `strcpy (r1), (r2)`
    - » Source and destination overlap: can't unwind in principle!
    - » IBM S/370 and VAX solution: execute twice – once read-only
- **What about "RISC" processors?**
  - For instance delayed branches?
    - » Example:    `bne somewhere`
                `ld r1,(sp)`
    - » Precise exception state consists of two PCs: PC and nPC
  - Delayed exceptions:
    - » Example:    `div r1, r2, r3`
                `ld r1, (sp)`
    - » What if takes many cycles to discover divide by zero, but load has already caused page fault?

# Precise Exceptions

- **Precise ⇒ state of the machine is preserved as if program executed up to the offending instruction**
  - All previous instructions <span style="color:red">completed</span>
  - Offending instruction and all following instructions act <span style="color:red">as if they have not even started</span>
  - Same system code will work on different implementations
  - Difficult in the presence of pipelining, out-of-order execution, ...
  - <span style="color:red">MIPS takes this position</span>
- **Imprecise ⇒ system software has to figure out what is where and put it all back together**
- **Performance goals often lead designers to forsake precise interrupts**
  - system software developers, user, markets etc. usually wish they had not done this
- <span style="color:red">**Modern techniques for out-of-order execution and branch prediction help implement precise interrupts**</span>

③ page is on backing store

operating system

② trap

reference

① 

load M

⑥ restart instruction

page table

⑤ reset page table

free frame

④ bring in missing page

physical memory

# Demand Paging Example

- **Since Demand Paging like caching, can compute average access time! ("Effective Access Time")**
  - EAT = Hit Rate x Hit Time + Miss Rate x Miss Time
  - EAT = Hit Time + Miss Rate x Miss Penalty
- **Example:**
  - Memory access time = 200 nanoseconds
  - Average page-fault service time = 8 milliseconds
  - Suppose p = Probability of miss, 1-p = Probably of hit
  - Then, we can compute EAT as follows:

    $$EAT = 200ns + p \times 8\ ms$$
    $$= 200ns + p \times 8{,}000{,}000ns$$

- **If one access out of 1,000 causes a page fault, then EAT = 8.2 µs:**
  - This is a slowdown by a factor of 40!
- **What if want slowdown by less than 10%?**
  - 200ns x 1.1 > EAT $\Rightarrow$ p < 2.5 x $10^{-6}$
  - This is about 1 page fault in 400000!

# What Factors Lead to Misses?

- **Compulsory Misses:**
  - Pages that have never been paged into memory before
  - How might we remove these misses?
    - » Prefetching: loading them into memory before needed
    - » Need to predict future somehow!  More later.
- **Capacity Misses:**
  - Not enough memory. Must somehow increase size.
  - Can we do this?
    - » One option: Increase amount of DRAM (not quick fix!)
    - » Another option:  If multiple processes in memory: adjust percentage of memory allocated to each one!
- **Conflict Misses:**
  - Technically, conflict misses don't exist in virtual memory, since it is a "fully-associative" cache
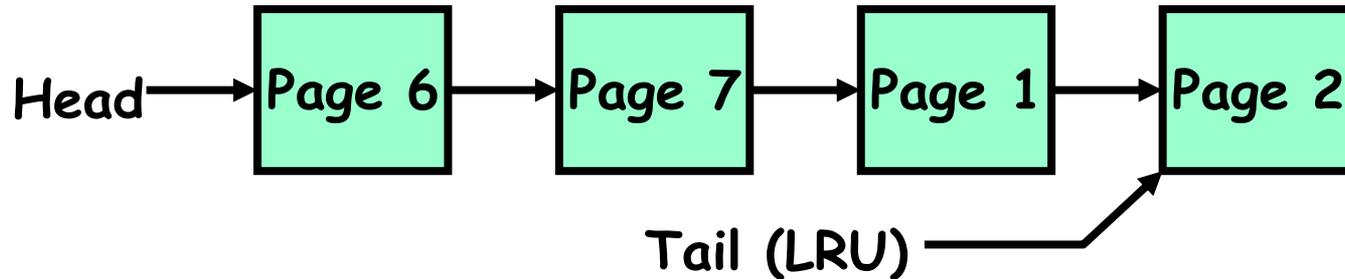- **Policy Misses:**
  - Caused when pages were in memory, but kicked out prematurely because of the replacement policy
  - How to fix? Better replacement policy

# Page Replacement Policies

- **Why do we care about Replacement Policy?**
  - Replacement is an issue with any cache
  - Particularly important with pages
    - » The cost of being wrong is high: must go to disk
    - » Must keep important pages in memory, not toss them out
- **FIFO (First In, First Out)**
  - Throw out oldest page.  Be fair – let every page live in memory for same amount of time.
  - Bad, because throws out heavily used pages instead of infrequently used pages
- **MIN (Minimum):**
  - Replace page that won't be used for the longest time
  - Great, but can't really know future…
  - Makes good comparison case, however
- **RANDOM:**
  - Pick random page for every replacement
  - Typical solution for TLB's.  Simple hardware
  - Pretty unpredictable – makes it hard to make real-time guarantees

# Replacement Policies (Con't)

- **LRU (Least Recently Used):**
  - Replace page that hasn't been used for the longest time
  - Programs have locality, so if something not used for a while, unlikely to be used in the near future.
  - Seems like LRU should be a good approximation to MIN.
- How to implement LRU? Use a list!

Head → Page 6 → Page 7 → Page 1 → Page 2

Tail (LRU) → (points to Page 2)

  - On each use, remove page from list and place at head
  - LRU page is at tail
- Problems with this scheme for paging?
  - Need to know immediately when each page used so that can change position in list...
  - Many instructions for each hardware access
- In practice, people **approximate** LRU (more later)

# Example: FIFO

- **Suppose we have 3 page frames, 4 virtual pages, and following reference stream:**
  - A B C A B D A D B C B
- **Consider FIFO Page replacement:**

| Ref:<br>Page: | A | B | C | A | B | D | A | D | B | C | B |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | A |  |  |  |  | D |  |  |  | C |  |
| 2 |  | B |  |  |  |  | A |  |  |  |  |
| 3 |  |  | C |  |  |  |  |  | B |  |  |

- **FIFO: 7 faults.**
- **When referencing D, replacing A is bad choice, since need A again right away**

# Example: MIN

- **Suppose we have the same reference stream:**
  - A B C A B D A D B C B
- **Consider MIN Page replacement:**

| Ref: | A | B | C | A | B | D | A | D | B | C | B |
|------|---|---|---|---|---|---|---|---|---|---|---|
| Page: | | | | | | | | | | | |
| 1 | A | | | | | | | | | C | |
| 2 | | B | | | | | | | | | |
| 3 | | | C | | | D | | | | | |

  - **MIN: 5 faults**
  - **Where will D be brought in? Look for page not referenced farthest in future.**
- **What will LRU do?**
  - **Same decisions as MIN here, but won't always be true!**

# When will LRU perform badly?

- Consider the following: A B C D A B C D A B C D
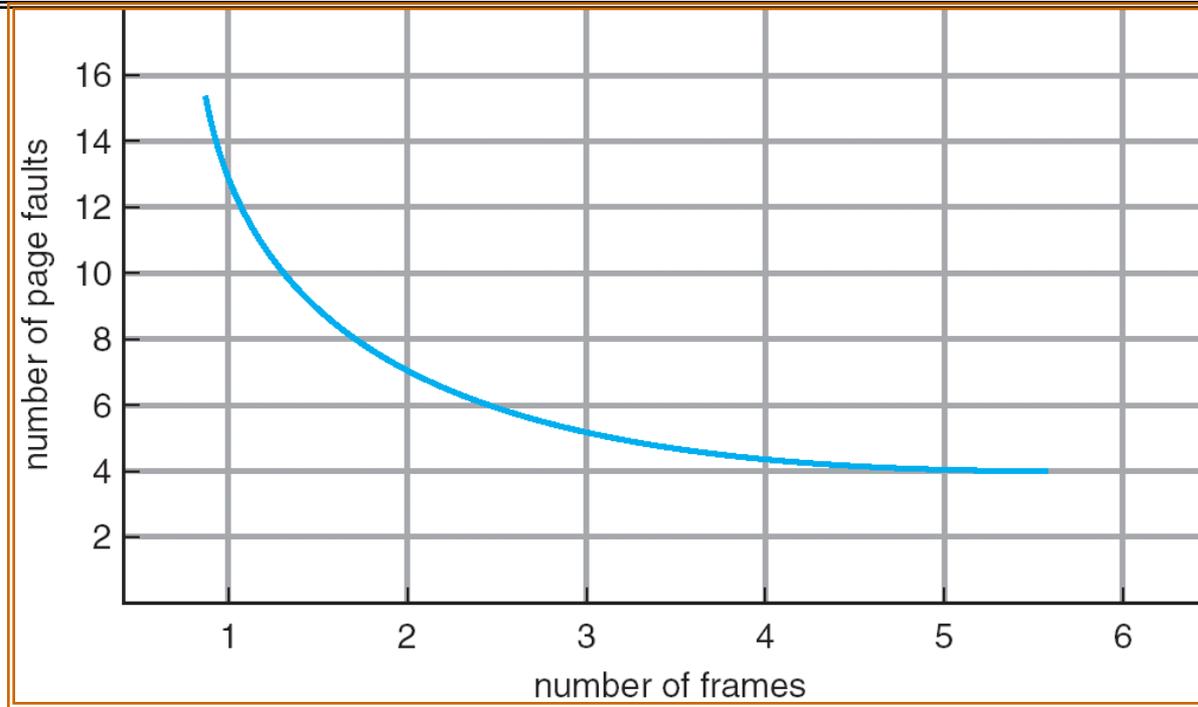- LRU Performs as follows (same as FIFO here):

| Ref: | A | B | C | D | A | B | C | D | A | B | C | D |
|------|---|---|---|---|---|---|---|---|---|---|---|---|
| Page: | | | | | | | | | | | | |
| 1 | A | | | D | | | C | | | B | | |
| 2 | | B | | | A | | | D | | | C | |
| 3 | | | C | | | B | | | A | | | D |

  – Every reference is a page fault!

- MIN Does much better:

| Ref: | A | B | C | D | A | B | C | D | A | B | C | D |
|------|---|---|---|---|---|---|---|---|---|---|---|---|
| Page: | | | | | | | | | | | | |
| 1 | A | | | | | | | | | B | | |
| 2 | | B | | | | | C | | | | | |
| 3 | | | C | D | | | | | | | | |

# Graph of Page Faults Versus The Number of Frames



- **One desirable property: When you add memory the miss rate goes down**
  - **Does this always happen?**
  - **Seems like it should, right?**
- **No: BeLady's anomaly**
  - **Certain replacement algorithms (FIFO) don't have this obvious property!**

# Adding Memory Doesn't Always Help Fault Rate

- **Does adding memory reduce number of page faults?**
  - Yes for LRU and MIN
  - Not necessarily for FIFO!  (Called Belady's anomaly)

| Ref:<br>Page: | A | B | C | D | A | B | E | A | B | C | D | E |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | A |   |   | D |   |   | E |   |   |   |   |   |
| 2 |   | B |   |   | A |   |   |   |   | C |   |   |
| 3 |   |   | C |   |   | B |   |   |   |   | D |   |

| Ref:<br>Page: | A | B | C | D | A | B | E | A | B | C | D | E |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | A |   |   |   |   |   | E |   |   |   | D |   |
| 2 |   | B |   |   |   |   |   | A |   |   |   | E |
| 3 |   |   | C |   |   |   |   |   | B |   |   |   |
| 4 |   |   |   | D |   |   |   |   |   | C |   |   |

- **After adding memory:**
  - With FIFO, contents can be completely different
  - In contrast, with LRU or MIN, contents of memory with X pages are a subset of contents with X+1 Page

# Implementing LRU

- **Perfect:**
  - Timestamp page on each reference
  - Keep list of pages ordered by time of reference
  - Too expensive to implement in reality for many reasons
- **Clock Algorithm: Arrange physical pages in circle with single clock hand**
  - Approximate LRU (approx to approx to MIN)
  - Replace an old page, not the oldest page
- **Details:**
  - Hardware "use" bit per physical page:
    - » Hardware sets use bit on each reference
    - » If use bit isn't set, means not referenced in a long time
    - » Nachos hardware sets use bit in the TLB; you have to copy this back to page table when TLB entry gets replaced
  - On page fault:
    - » Advance clock hand (not real time)
    - » Check use bit: 1→used recently; clear and leave alone
      0→selected candidate for replacement
  - Will always find a page or loop forever?
    - » Even if all use bits set, will eventually loop around⇒FIFO

# Clock Algorithm: Not Recently Used

Set of all pages
in Memory

**Single Clock Hand:**

**Advances only on page fault!**
**Check for pages not used recently**
**Mark pages as not used recently**

- **What if hand moving slowly?**
  - **Good sign or bad sign?**
    - » Not many page faults and/or find page quickly
- **What if hand is moving quickly?**
  - Lots of page faults and/or lots of reference bits set
- **One way to view clock algorithm:**
  - Crude partitioning of pages into two groups: young and old
  - Why not partition into more than 2 groups?

# N<sup>th</sup> Chance version of Clock Algorithm

- **N<sup>th</sup> chance algorithm:** Give page N chances
  - OS keeps counter per page: # sweeps
  - On page fault, OS checks use bit:
    - » $1 \Rightarrow$ clear use and also clear counter (used in last sweep)
    - » $0 \Rightarrow$ increment counter; if count=N, replace page
  - Means that clock hand has to sweep by N times without page being used before page is replaced
- How do we pick N?
  - Why pick large N? Better approx to LRU
    - » If N ~ 1K, really good approximation
  - Why pick small N? More efficient
    - » Otherwise might have to look a long way to find free page
- What about dirty pages?
  - Takes extra overhead to replace a dirty page, so give dirty pages an extra chance before replacing?
  - Common approach:
    - » Clean pages, use N=1
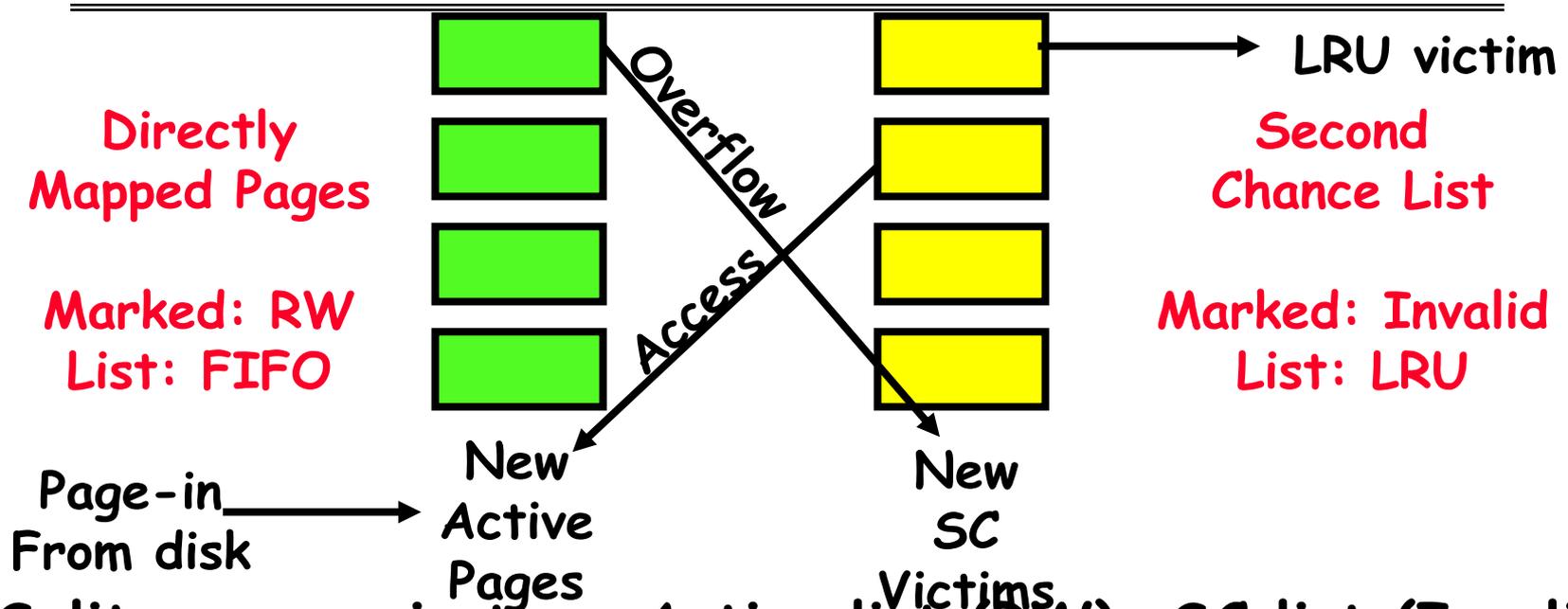    - » Dirty pages, use N=2 (and write back to disk when N=1)

# Clock Algorithms: Details

- **Which bits of a PTE entry are useful to us?**
  - **Use:** Set when page is referenced; cleared by clock algorithm
  - **Modified:** set when page is modified, cleared when page written to disk
  - **Valid:** ok for program to reference this page
  - **Read-only:** ok for program to read page, but not modify
    - » For example for catching modifications to code pages!
- **Do we really need hardware-supported "modified" bit?**
  - No.  Can emulate it (BSD Unix) using read-only bit
    - » Initially, mark all pages as read-only, even data pages
    - » On write, trap to OS. OS sets software "modified" bit, and marks page as read-write.
    - » Whenever page comes back in from disk, mark read-only

# Clock Algorithms Details (continued)

- **Do we really need a hardware-supported "use" bit?**
  - No. Can emulate it similar to above:
    » Mark all pages as invalid, even if in memory
    » On read to invalid page, trap to OS
    » OS sets use bit, and marks page read-only
  - Get modified bit in same way as previous:
    » On write, trap to OS (either invalid or read-only)
    » Set use and modified bits, mark page read-write
  - When clock hand passes by, reset use and modified bits and mark page as invalid again
- **Remember, however, that clock is just an approximation of LRU**
  - Can we do a better approximation, given that we have to take page faults on some reads and writes to collect use information?
  - Need to identify an old page, not oldest page!
  - Answer: second chance list
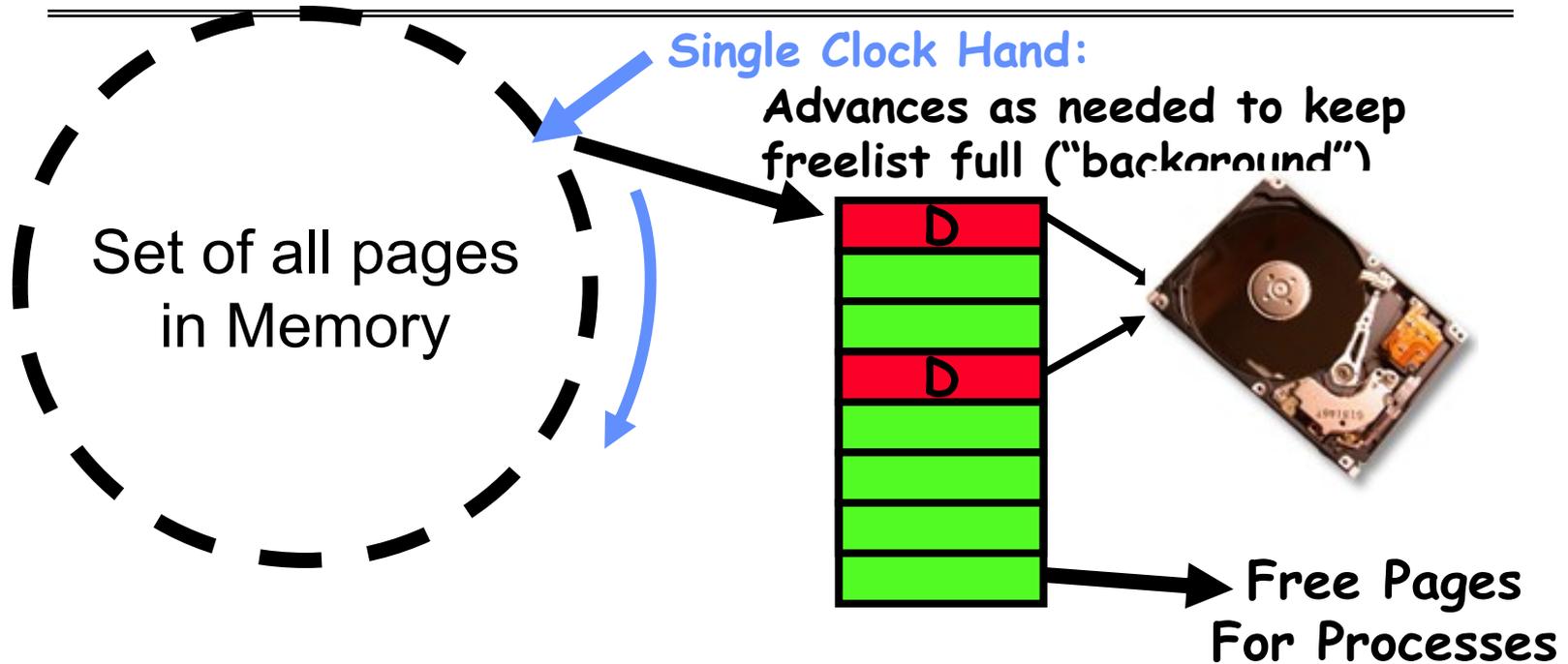
# Second-Chance List Algorithm (VAX/VMS)

**Directly Mapped Pages**

**Marked: RW List: FIFO**

Overflow

Access

**LRU victim**

**Second Chance List**

**Marked: Invalid List: LRU**

Page-in From disk ⟶ New Active Pages

New SC Victims

- **Split memory in two: Active list (RW), SC list (Invalid)**
- **Access pages in Active list at full speed**
- **Otherwise, Page Fault**
  - **Always move overflow page from end of Active list to front of Second-chance list (SC) and mark invalid**
  - **Desired Page On SC List: move to front of Active list, mark RW**
  - **Not on SC list: page in to front of Active list, mark RW; page out LRU victim at end of SC list**

# Second-Chance List Algorithm (con't)

* **How many pages for second chance list?**
  * If 0 $\Rightarrow$ FIFO
  * If all $\Rightarrow$ LRU, but page fault on every page reference
* **Pick intermediate value.  Result is:**
  * Pro: Few disk accesses (page only goes to disk if unused for a long time)
  * Con: Increased overhead trapping to OS (software / hardware tradeoff)
* **With page translation, we can adapt to any kind of access the program makes**
  * Later, we will show how to use page translation / protection to share memory between threads on widely separated machines
* **Question: why didn't VAX include "use" bit?**
  * Strecker (architect) asked OS people, they said they didn't need it, so didn't implement it
  * He later got blamed, but VAX did OK anyway

# Free List

**Single Clock Hand:**

**Advances as needed to keep freelist full ("background")**

Set of all pages in Memory

D

D

Free Pages For Processes

- **Keep set of free pages ready for use in demand paging**
  - **Freelist filled in background by Clock algorithm or other technique ("Pageout demon")**
  - **Dirty pages start copying back to disk when enter list**
- **Like VAX second-chance list**
  - **If page needed before reused, just return to active set**
- **Advantage: Faster for page fault**
  - **Can always use page (or pages) immediately on fault**

# Demand Paging (more details)

- **Does software-loaded TLB need use bit? Two Options:**
  - Hardware sets use bit in TLB; when TLB entry is replaced, software copies use bit back to page table
  - Software manages TLB entries as FIFO list; everything not in TLB is Second-Chance list, managed as strict LRU

- **Core Map**
  - Page tables map virtual page $\rightarrow$ physical page
  - Do we need a reverse mapping (i.e. physical page $\rightarrow$ virtual page)?
    - » Yes. Clock algorithm runs through page frames. If sharing, then multiple virtual-pages per physical page
    - » Can't push page out to disk without invalidating all PTEs

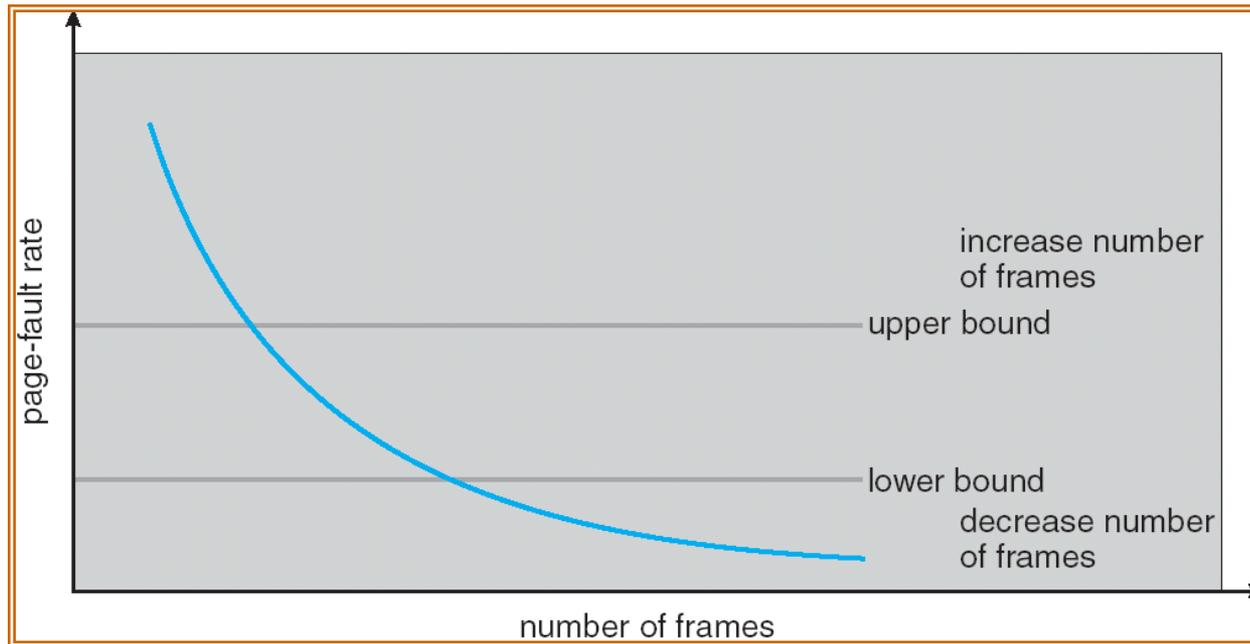# Allocation of Page Frames (Memory Pages)

- **How do we allocate memory among different processes?**
  - Does every process get the same fraction of memory? Different fractions?
  - Should we completely swap some processes out of memory?
- **Each process needs *minimum* number of pages**
  - Want to make sure that all processes <span style="color:red">**that are loaded into memory**</span> can make forward progress
  - Example: IBM 370 – 6 pages to handle SS MOVE instruction:
    - » instruction is 6 bytes, might span 2 pages
    - » 2 pages to handle *from*
    - » 2 pages to handle *to*
- **Possible Replacement Scopes:**

  - <span style="color:red">**Global replacement**</span> – process selects replacement frame from set of all frames; one process can take a frame from another

  - <span style="color:red">**Local replacement**</span> – each process selects from only its own set of allocated frames

# Fixed/Priority Allocation

- **Equal allocation** (Fixed Scheme):
  - Every process gets same amount of memory
  - Example: 100 frames, 5 processes $\Rightarrow$ process gets 20 frames
- **Proportional allocation** (Fixed Scheme)
  - Allocate according to the size of process
  - Computation proceeds as follows:
    $s_i$ = size of process $p_i$ and $S = \Sigma s_i$
    $m$ = total number of frames

    $a_i$ = allocation for $p_i = \dfrac{s_i}{S} \times m$

- **Priority Allocation:**
  - Proportional scheme using priorities rather than size
    » Same type of computation as previous scheme
  - Possible behavior: If process $p_i$ generates a page fault, select for replacement a frame from a process with lower priority number
- Perhaps we should use an adaptive scheme instead???
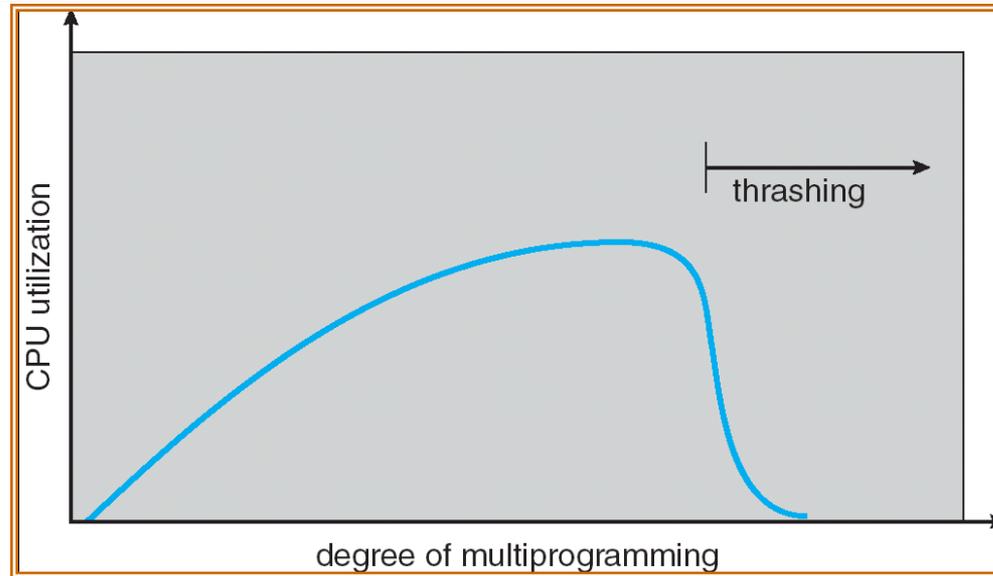  - What if some application just needs more memory?

# Page-Fault Frequency Allocation

- **Can we reduce Capacity misses by dynamically changing the number of pages/application?**



- **Establish "acceptable" page-fault rate**
  - **If actual rate too low, process loses frame**
  - **If actual rate too high, process gains frame**
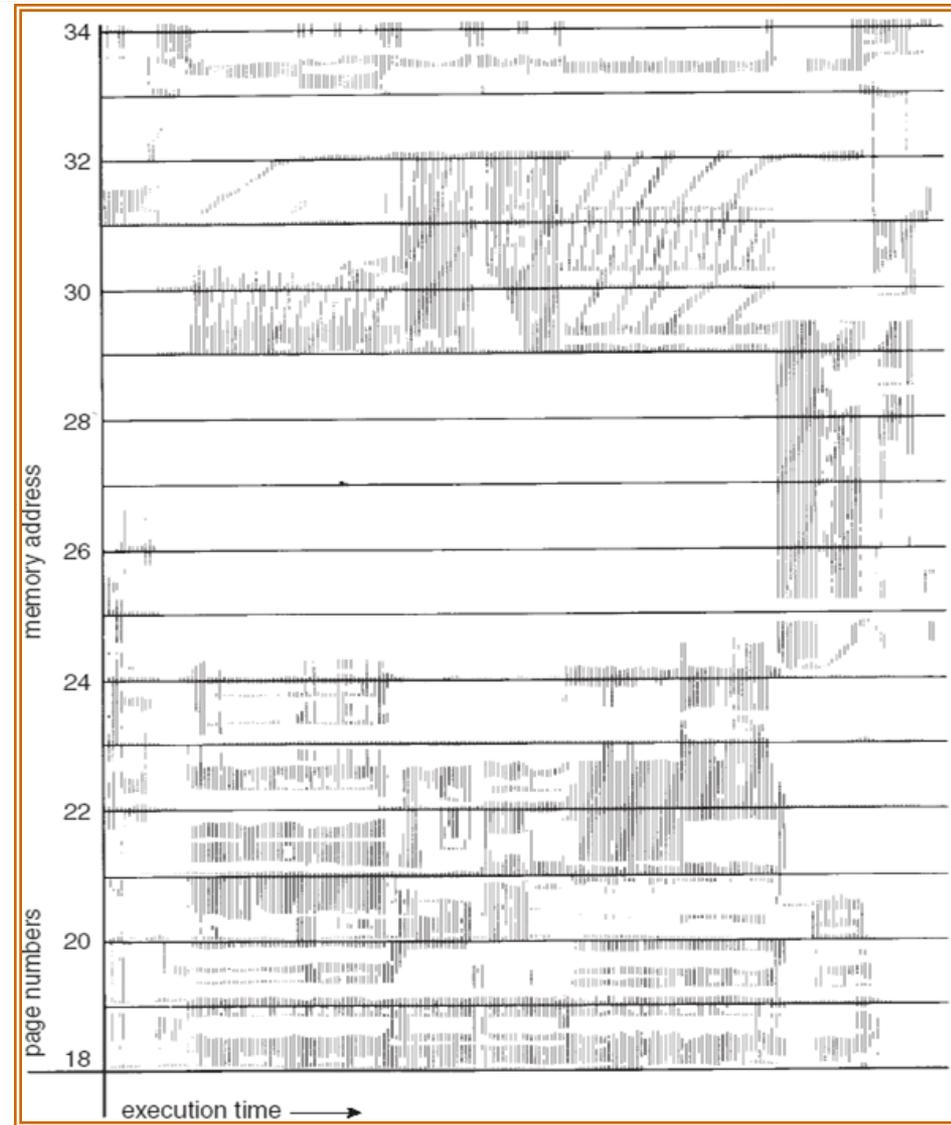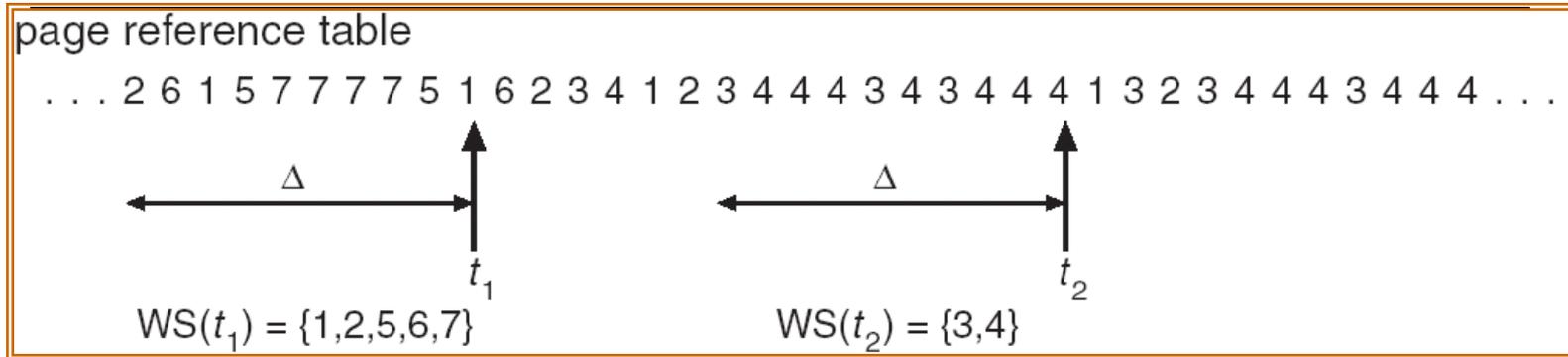- **Question: What if we just don't have enough memory?**

# Thrashing



- **If a process does not have "enough" pages, the page-fault rate is very high. This leads to:**
  - **low CPU utilization**
  - **operating system spends most of its time swapping to disk**
- **Thrashing ≡ a process is busy swapping pages in and out**
- **Questions:**
  - **How do we detect Thrashing?**
  - **What is best response to Thrashing?**

# Locality In A Memory-Reference Pattern

- **Program Memory Access Patterns have temporal and spatial locality**
  - Group of Pages accessed along a given time slice called the "Working Set"
  - Working Set defines minimum number of pages needed for process to behave well

- **Not enough memory for Working Set⇒Thrashing**
  - Better to swap out process?

# Working-Set Model



page reference table

. . . 2 6 1 5 7 7 7 7 5 1 6 2 3 4 1 2 3 4 4 4 3 4 3 4 4 4 1 3 2 3 4 4 4 3 4 4 4 . . .

$\Delta$

$t_1$

$\Delta$

$t_2$

$WS(t_1) = \{1,2,5,6,7\}$    $WS(t_2) = \{3,4\}$

- $\Delta \equiv$ working-set window $\equiv$ fixed number of page references
  - Example:  10,000 instructions
- $WS_i$ (working set of Process $P_i$) = total set of pages referenced in the most recent $\Delta$ (varies in time)
  - if $\Delta$ too small will not encompass entire locality
  - if $\Delta$ too large will encompass several localities
  - if $\Delta = \infty \Rightarrow$ will encompass entire program
- $D = \Sigma |WS_i| \equiv$ total demand frames
- if $D > m \Rightarrow$ Thrashing
  - Policy: if $D > m$, then suspend/swap out processes
  - This can improve overall system behavior by a lot!

# What about Compulsory Misses?

- **Recall that compulsory misses are misses that occur the first time that a page is seen**
  - Pages that are touched for the first time
  - Pages that are touched after process is swapped out/swapped back in
- **Clustering:**
  - On a page-fault, bring in multiple pages "around" the faulting page
  - Since efficiency of disk reads increases with sequential reads, makes sense to read several sequential pages
- **Working Set Tracking:**
  - Use algorithm to try to track working set of application
  - When swapping process back in, swap in working set

# Summary (1/2)

- **TLB is cache on translations**
  - Fully associative to reduce conflicts
  - Can be overlapped with cache access
- **Demand Paging:**
  - Treat memory as cache on disk
  - Cache miss $\Rightarrow$ get page from disk
- **Transparent Level of Indirection**
  - User program is unaware of activities of OS behind scenes
  - Data can be moved without affecting application correctness
- **Software-loaded TLB**
  - Fast Path: handled in hardware (TLB hit with valid=1)
  - Slow Path: Trap to software to scan page table
- **Precise Exception specifies a single instruction for which:**
  - All previous instructions have completed (committed state)
  - No following instructions nor actual instruction have started
- **Replacement policies**
  - FIFO: Place pages on queue, replace page at end
  - MIN: replace page that will be used farthest in future
  - LRU: Replace page that hasn't be used for the longest time

# Summary (2/2)

- **Clock Algorithm: Approximation to LRU**
  - Arrange all pages in circular list
  - Sweep through them, marking as not "in use"
  - If page not "in use" for one pass, than can replace
- **N$^{th}$-chance clock algorithm: Another approx LRU**
  - Give pages multiple passes of clock hand before replacing
- **Second-Chance List algorithm: Yet another approx LRU**
  - Divide pages into two groups, one of which is truly LRU and managed on page faults.
- **Working Set:**
  - Set of pages touched by a process recently
- **Thrashing: a process is busy swapping pages in and out**
  - Process will thrash if working set doesn't fit in memory
  - Need to swap out a process