

# Operating Systems 412

---

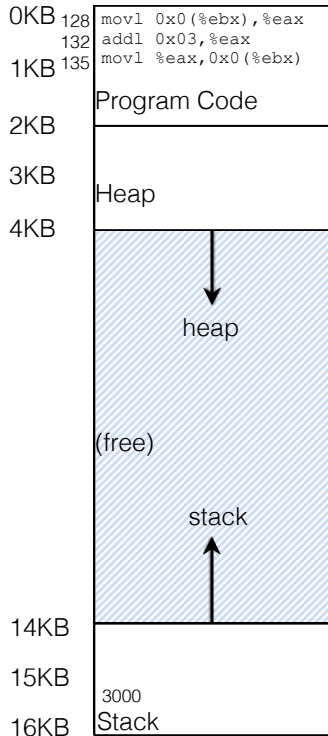
Pete Keleher

*Memory*

## Memory

- 13 - Address Spaces
- 15 - Address Translation
- 16 - Segmentation
- 17 - Free Space Management
- 18 - Paging
- 19 - Translation Lookaside Buffers
- 20 - Advanced Paging
- 21 - Swapping
- 22 - Swapping Policy

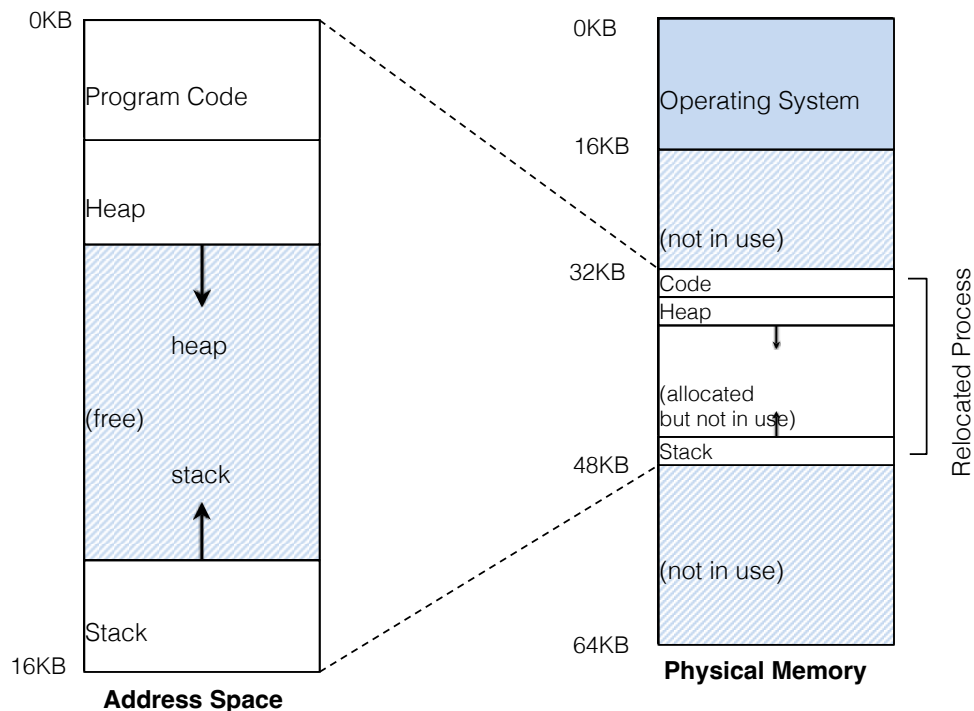
# Example: Address Translation



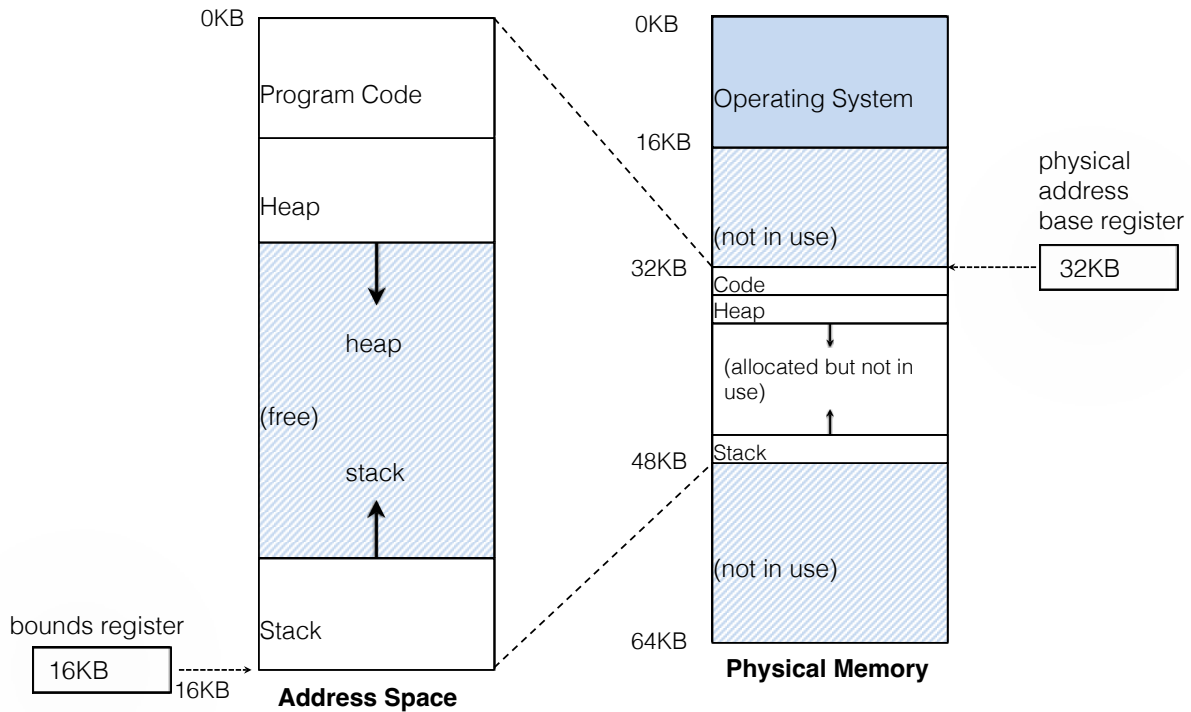
- Fetch instruction at address 128
- Execute instruction (load from address 15KB)
- Fetch instruction at address 132
- Execute instruction (no memory reference)
- Fetch the instruction at address 135
- Execute instruction (store to address 15 KB)

*But not all programs can be at location 0*

# A Single Relocated Process



# One Approach *base and bounds*



104

## Dynamic(Hardware base) Relocation

- OS decides where in physical memory a process is loaded.
  - Set the **base** register:  
 $\text{physical address} = \text{virtual address} + \text{base}$
  - Virtual addresses must **not be greater than bound** or **negative**:  
 $0 \leq \text{virtual address} < \text{bound}$

105

# Address Translation *base and bounds*

**128 : movl 0x0(%ebx), %eax**

- **Fetch** instruction at address 128

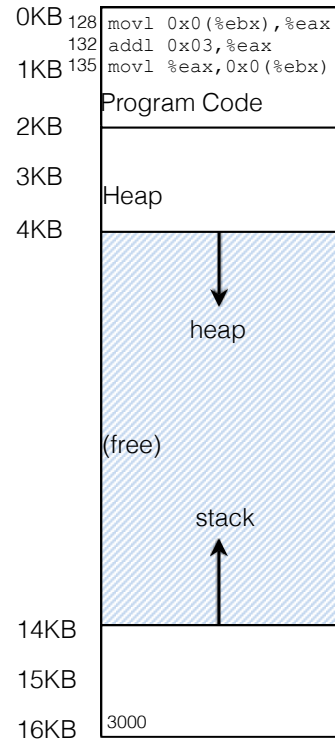
$$32896 = 128 + 32KB(base)$$

- **Execute** this instruction

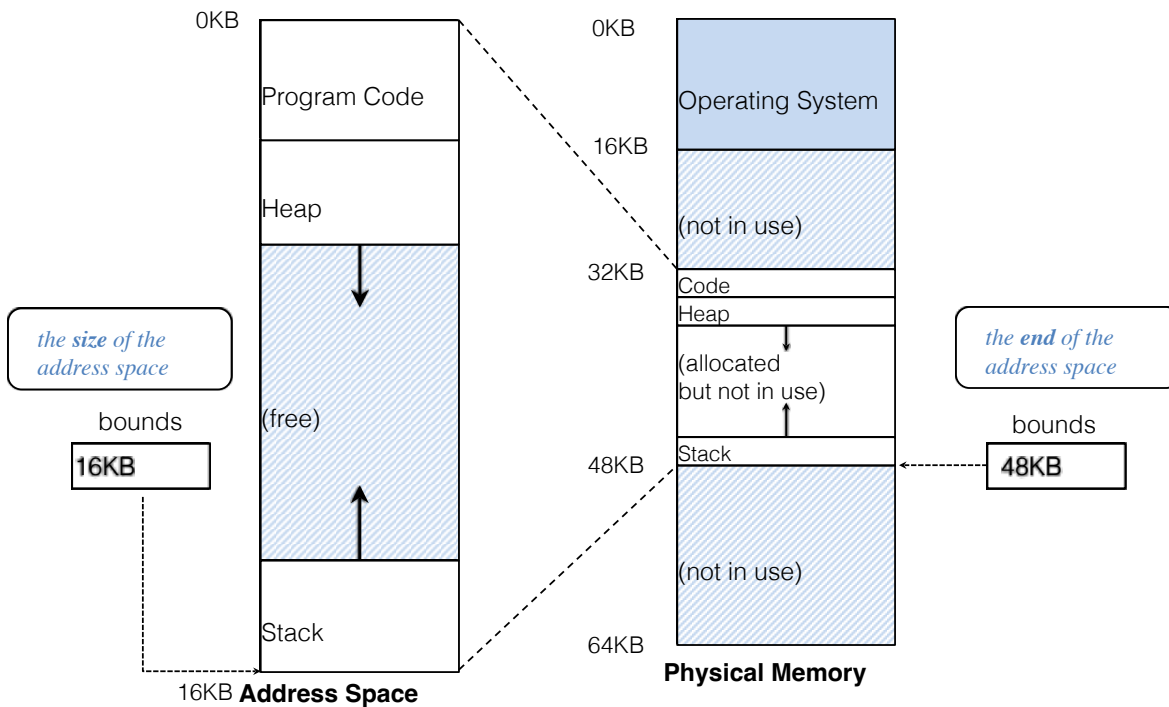
- Load from address 15KB

$$47KB = 15KB + 32KB(base)$$

$$physical\ address = offset + base$$



## Two ways to Use the Bounds Register



## Hardware Requirements *for base and bounds*

- Privileged mode
  - user processes should not execute privileged operations
- Base and bounds registers
  - per CPU / core
- Translate virtual addresses, and check bounds
  - instructions
  - data
- Privileged instructions to update base/bound registers
- Ability to raise exceptions
  - out-of-bound accesses

108

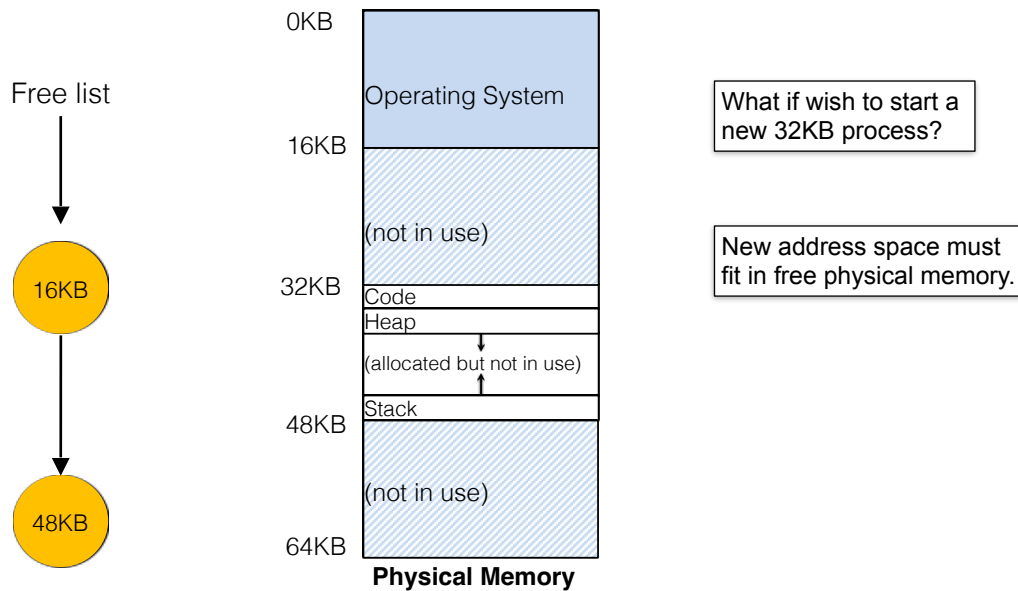
## OS Requirements *for base and bounds*

- OS must intervene at three critical junctures:
  - When a process starts running:
    - find space for address space in physical memory
  - When a process is terminated:
    - reclaims the memory for use
  - When context switch occurs:
    - Save and store the base-and-bounds pair

109

# OS Issues: *single chunk of memory*

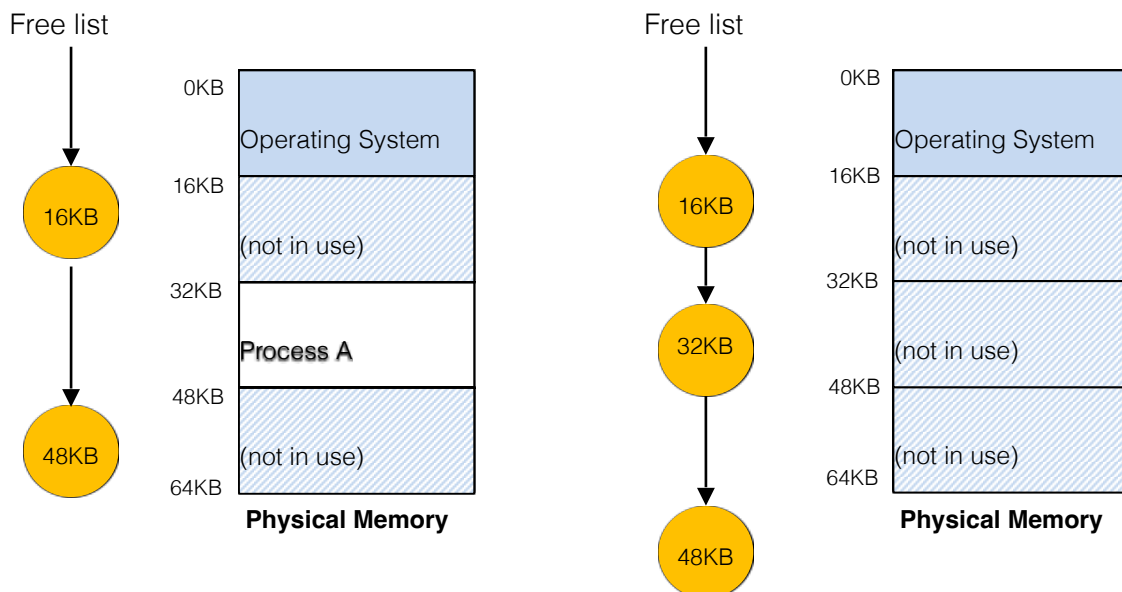
- The OS must find a room for a new address space.
  - free list : a list of unused ranges of physical memory



110

# OS Issues: Process Termination

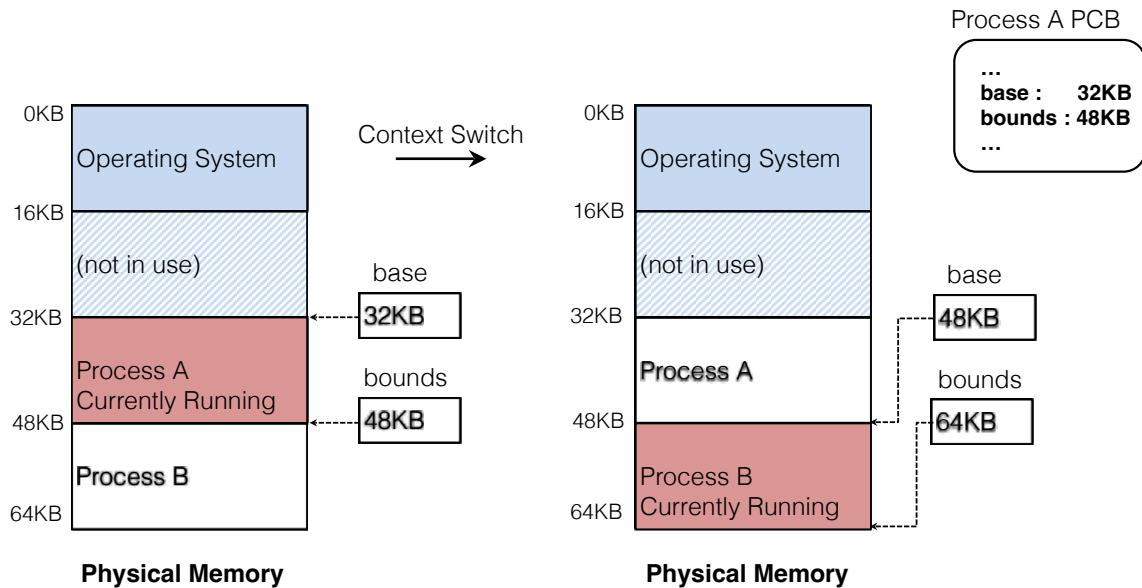
- OS must put memory back on the free list.



111

# OS Issues: Context Switches

- OS must save and restore base-and-bounds pair.
  - In process structure or process control block (PCB)

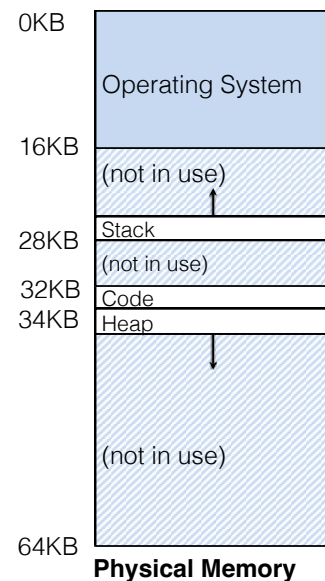


112

## Segmentation *another approach*

- Segment is a *contiguous* portion of the address space:
  - code, stack, heap, ...
- Can be placed anywhere in *contiguous* physical memory.
  - Basically base and bounds *per segment*

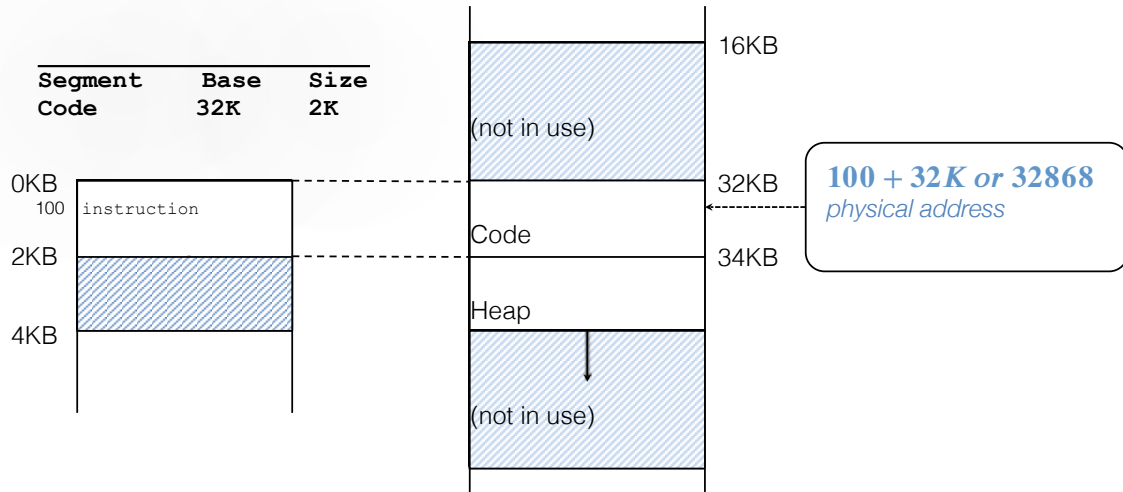
Segment	Base	Bound
Code	32K	2K
Heap	34K	2K
Stack	28K	2K



113

# Address Translation for Segments

- The *offset* of virtual address 100 is: 100
  - The code segment **starts at virtual address 0** in address space.

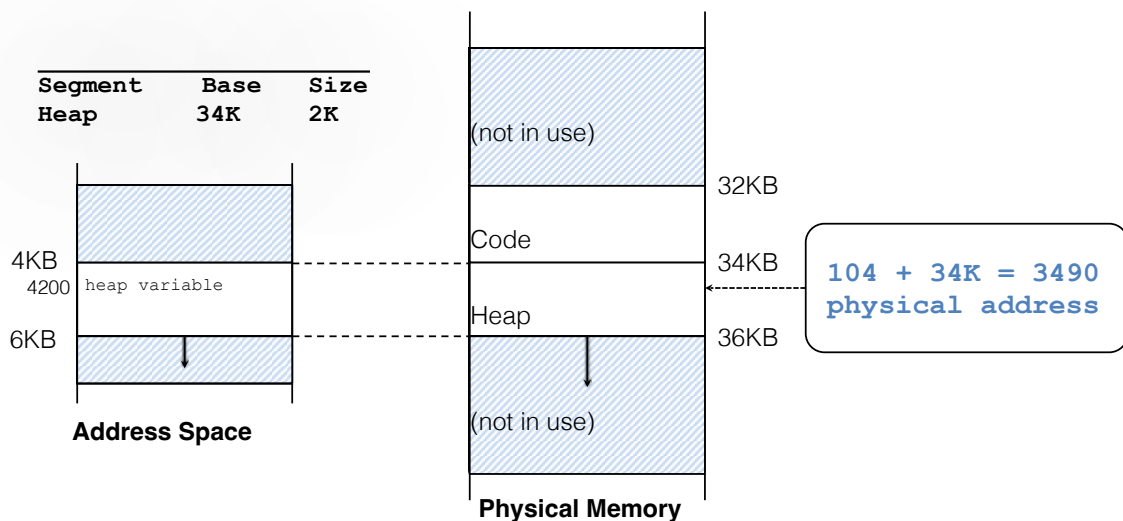


$$\text{physical address} = \text{offset} + \text{base}$$

114

# Address Translation for Segments

- The *offset* of virtual address 4200 is 104.
  - The heap segment **starts at virtual address 4096** in address space.



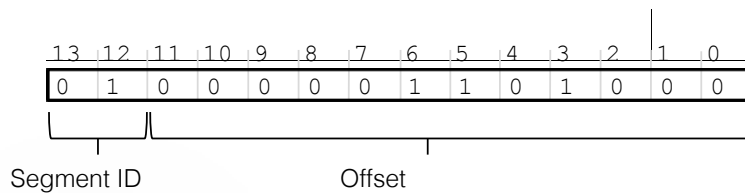
$$\text{physical address} = \text{offset} + \text{base}$$

115



# Segment Descriptors

- Explicit approach
  - Chop up the address space into segments based on the **top few bits** of virtual address.
- Example: virtual address 4200 (01000001101000)



Segment	bits
Code	00
Heap	01
Stack	10
-	11

116

# Segment Descriptors

- Bits
  - SEG\_MASK = 0x3000 (110000000000000)
  - SEG\_SHIFT = 12
  - OFFSET\_MASK = 0xFFF (001111111111111)

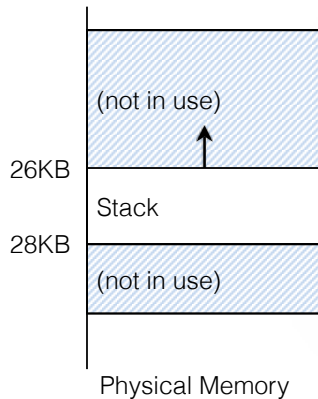
```
1 // get top 2 bits of 14-bit VA
2 Segment = (VirtualAddress & SEG_MASK) >> SEG_SHIFT

3 // now get offset
4 Offset = VirtualAddress & OFFSET_MASK
5 if (Offset >= Bounds[Segment])
6     RaiseException(PROTECTION_FAULT)
7 else
8     PhysAddr = Base[Segment] + Offset
9     Register = AccessMemory(PhysAddr)
```

117

# Referring to Stack Segment

- Stack grows *backward*.
- *Extra hardware support needed*.
  - The hardware checks which way the segment grows.
  - 1: positive direction, 0: negative direction



Segment Register(with Negative-Growth Support)

Segment	Base	Size	Grows Positive?
Code	32K	2K	1
Heap	34K	2K	1
Stack	28K	2K	0

$physical\ address = base + offset - sizeof(stack)$

*two's complement arithmetic*

118

# Support for Sharing

- Segments can be shared between address spaces
  - Code sharing still used
- Need hardware support in form of *protection bits*.
  - Bits indicate read, write and execute permissions.

Segment Register Values(with Protection)

Segment	Base	Size	Grows Positive?	Protection
Code	32K	2K	1	Read-Execute
Heap	34K	2K	1	Read-Write
Stack	28K	2K	0	Read-Write

119

# Fine-Grained and Coarse-Grained

- Coarse-Grained is small number of segments
  - e.g., code, heap, stack.
- Fine-Grained segmentation allows more flexibility
  - Hardware-supported segment tables

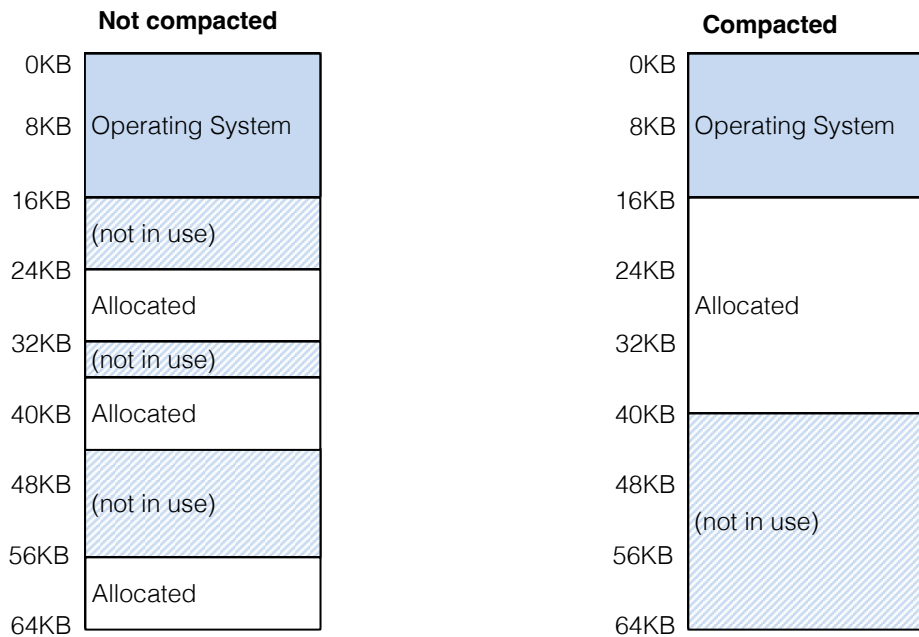
120

# OS support: Fragmentation

- External Fragmentation:
  - Distinct runs of free space in physical memory
  - Might be **24KB free**, but **not in one contiguous** segment.
  - The OS **cannot** immediately satisfy the **20KB request**.
- Compaction: consolidating segments in physical memory.
  - Compaction is **costly**.
    - **Stop** running process.
    - **Copy** data to somewhere.
    - **Change** segment register value.

121

# Memory Compaction



122

## GeekOS

- [segmented memory addresses](#)
  - 16-bit “segment selector”, 32-bit offset
  - segment selector has:
    - 1 bit: GDT or LDT
    - 13 bits: index into GDT or LDT
    - 2 bits: protection level of segment
  - segment descriptor (from table) has:
    - linear *base* physical address of segment: 32 bits
    - limit (size) of segment: 20 bits
    - descriptor privilege level (dpl): 2 bits
    - type of segment (data, code, system, tss, gate): 4 bits
    - present (in-memory): 1 bit
    - etc.

123

# GeekOS

- GDT
  - entries point to kernel segments, optionally user segments
  - entry 0 (null selector) is not used to access memory
  - `gdt r` register points to the GDT
- LDT similar, but
  - points to segments of a single process
  - entry 0 can be used
  - any number of LDTs can be in memory
  - `ldt r` register points (via GDT) to currently used LDT

124

# Virtual Memory

- 14 - Memory API
- 13 - Address Spaces
- 15 - Address Translation
- 16 - Segmentation
- 17 - Free Space Management
- 18 - Paging
- 19 - Translation Lookaside Buffers
- 20 - Advanced Paging
- 21 - Swapping
- 22 - Swapping Policy

125

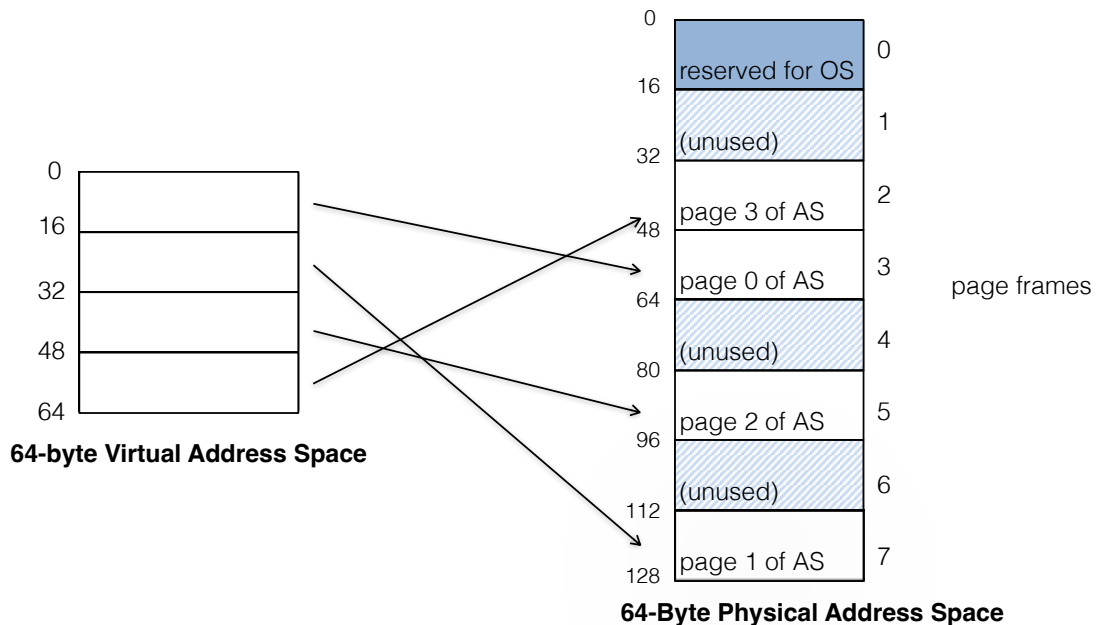
# Paging

- Paging splits address space into fixed-size *pages*.
  - vs segmentation: variable size of logical segments
- Physical memory holding a page is the *page frame*
- Per-process page tables
  - translate virtual address to physical address.
- Flexibility:
  - No assumptions on how heap and stack grow or are used
- Simplicity: ease of free-space management
  - All pages and page frames are the same size
  - Free lists are easy...

126

# Paging Example

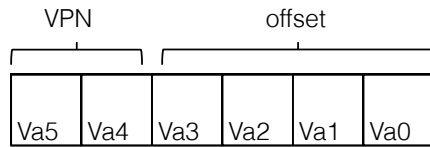
- 128-byte physical memory with eight 16-byte page frames
- 64-byte address space with 16-byte pages



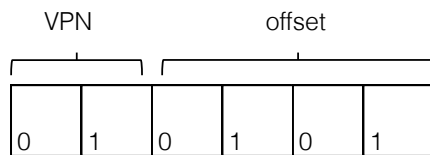
127

# Address Translation

- Two components in the virtual address
  - VPN: virtual page number
  - Offset: offset within the page



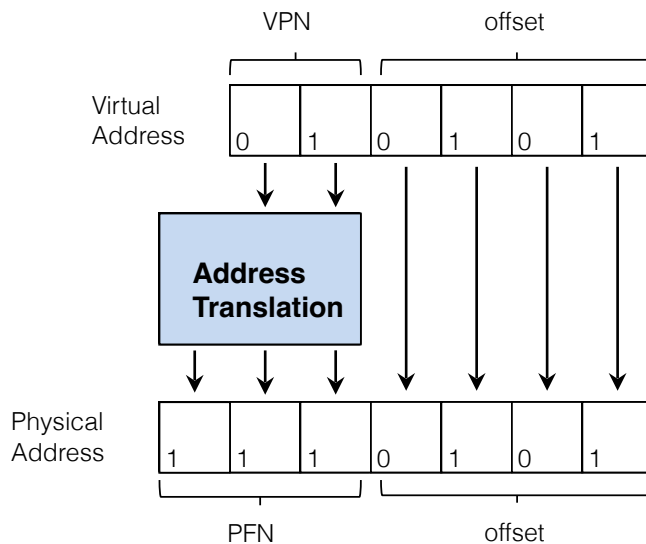
- Example: virtual address 21 in 64-byte address space



128

# Example: Address Translation

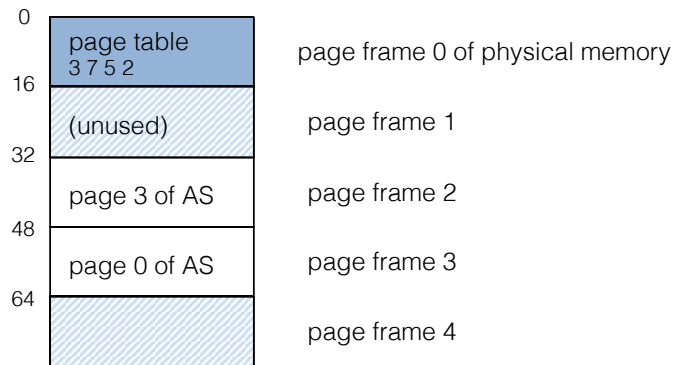
- The virtual address 21 in 64-byte address space



129

# Where Are Page Tables Stored?

- Page tables can be large...
  - 32-bit address space with 4-KB pages, 20 bits for VPN
    - assume entry is 4 bytes:
    - page table size is  $2^{20} * 4 = 2^{22} = 4\text{MB}$  of space
- Page tables for each process are stored in memory...



130

# What Is In The Page Table?

- A page table is just a **data structure** that is used to map the virtual address to physical address.
  - Simplest form: a linear page table, an array
- The OS/hardware accesses a page-table entry by indexing into the array by virtual page-number
- Common bits:
  - Valid Bit: whether the particular translation is valid.
  - Protection Bit: read, write, execute
  - Present Bit: in physical memory or swapped out
  - Dirty Bit: page modified since it brought into memory
  - Reference Bit (Accessed Bit): page has been accessed

131





# Accessing Memory With Paging

```
1 // Extract the VPN from the virtual address
2 VPN = (VirtualAddress & VPN_MASK) >> SHIFT
3
4 // Form the address of the page-table entry (PTE)
5 PTEAddr = PTBR + (VPN * sizeof(PTE))
6
7 // Fetch the PTE
8 PTE = AccessMemory(PTEAddr)
9
10 // Check if process can access the page
11 if (PTE.Valid == False)
12     RaiseException(SEGMENTATION_FAULT)
13 else if (CanAccess(PTE.ProtectBits) == False)
14     RaiseException(PROTECTION_FAULT)
15 else
16     // Access is OK: form physical address and fetch
17     offset = VirtualAddress & OFFSET_MASK
18     PhysAddr = (PTE.PFN << PFN_SHIFT) | offset
19     Register = AccessMemory(PhysAddr)
```

134

## A Memory Trace

- Example: A Simple Memory Access

```
int array[1000];
...
for (i = 0; i < 1000; i++)
    array[i] = 0;
```

- Compile and execute

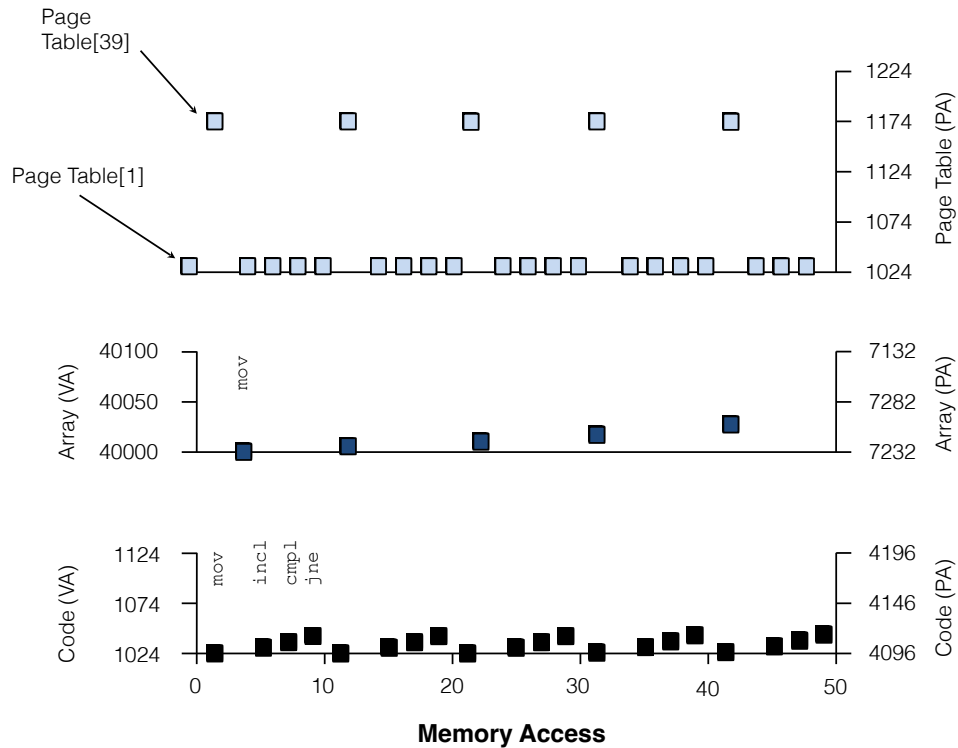
```
prompt> gcc -o array array.c -Wall -o
prompt> ./array
```

- Resulting Assembly code

```
0x1024 movl $0x0, (%edi,%eax,4)
0x1028 incl %eax
0x102c cmpl $0x03e8,%eax
0x1030 jne 0x1024
```

135

# A Virtual(And Physical) Memory Trace



136

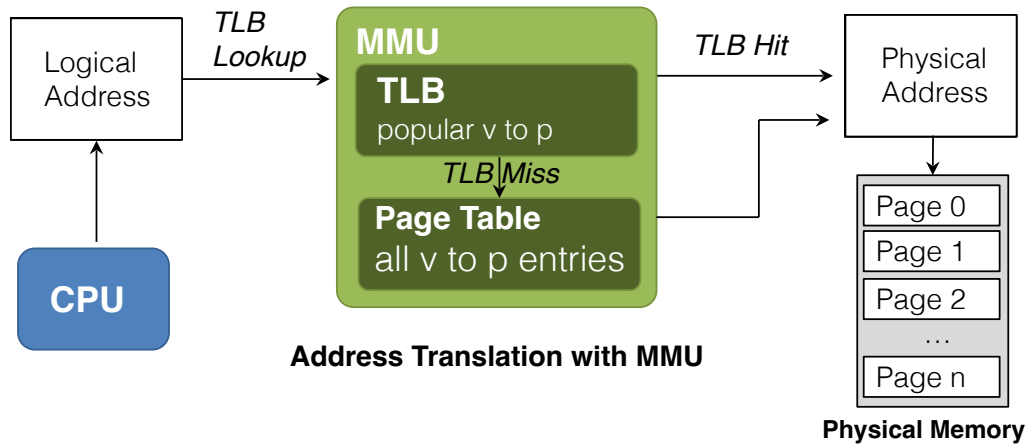
## Virtual Memory

- 13 - Address Spaces
- 14 - Memory API
- 15 - Address Translation
- 16 - Segmentation
- 17 - Free Space Management
- 18 - Paging
- 19 - Translation Lookaside Buffers
- 20 - Advanced Paging
- 21 - Swapping
- 22 - Swapping Policy

137

# TLB

- Part of the chip's memory-management unit (MMU).
- A hardware cache of **popular** virtual-to-physical address translation.



138

## Basic TLB Algorithm

- extract the virtual page number (VPN).
- check for hit in the the TLB
- extract page frame number from relevant TLB entry, form desired physical address, and access memory

139

# Basic TLB Algorithm

```
1  VPN = (VirtualAddress & VPN_MASK) >> SHIFT
2  (Success, TlbEntry) = TLB_Lookup(VPN)
3  if (Success == True)    // TLB Hit
4      if (CanAccess(TlbEntry.ProtectBits) == True)
5          Offset    = VirtualAddress & OFFSET_MASK
6          PhysAddr  = (TlbEntry.PFN << SHIFT) | Offset
7          Register  = AccessMemory(PhysAddr)
8      else
9          RaiseException(PROTECTION_FAULT)
10 else
11     // TLB Miss
12     PTEAddr = PTBR + (VPN * sizeof(PTE))
13     PTE = AccessMemory(PTEAddr)
14     if (PTE.Valid == False)
15         RaiseException(SEGMENTATION_FAULT)
16     else if (CanAccess(PTE.ProtectBits) == False)
17         RaiseException(PROTECTION_FAULT)
18     else
19         TLB_Insert(VPN, PTE.PFN, PTE.ProtectBits)
20     RetryInstruction()
```

Figure 19.1: TLB Control Flow Algorithm

140

## Example: Accessing An Array

- How a TLB can improve its performance.

	OFFSET			
	00	04	08	12
VPN = 00				
VPN = 01				
VPN = 03				
VPN = 04				
VPN = 05				
VPN = 06		a[0]	a[1]	a[2]
VPN = 07	a[3]	a[4]	a[5]	a[6]
VPN = 08	a[7]	a[8]	a[9]	
VPN = 09				
VPN = 10				
VPN = 11				
VPN = 12				
VPN = 13				
VPN = 14				

The TLB improves performance due to **spatial locality**

```
0:     int sum = 0 ;
1:     for( i=0; i<10; i++){
2:         sum+=a[i];
3:     }
```

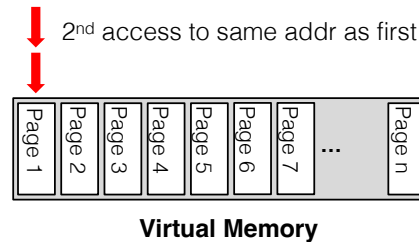
3 TLB misses and 7 hits.  
Thus **TLB hit rate** is 70%.

141

# Locality

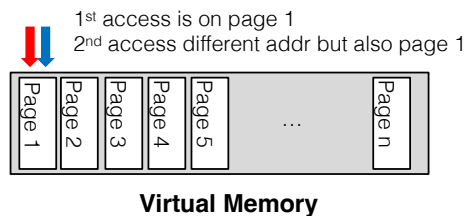
- Temporal Locality

- An instruction or data item that has been recently accessed will likely be re-accessed soon in the future.



- Spatial Locality

- If a program accesses memory at address  $x$ , it will likely soon access memory near  $x$ .



142

# Virtual Memory

- 13 - Address Spaces
- 14 - Memory API
- 15 - Address Translation
- 16 - Segmentation
- 17 - Free Space Management
- 18 - Paging
- 19 - Translation Lookaside Buffers
- 20 - Advanced Paging
- 21 - Swapping
- 22 - Swapping Policy

143

# Who Handles The TLB Miss?

- Hardware handles the TLB miss entirely on CISC processors.
  - The hardware know where the page tables are located
  - ... “walks” the page table, finding the correct entry and extracting the desired translation, update and retry instruction.
  - this is a hardware-managed TLB.
- RISC processors often manage TLBs in software.
  - On a TLB miss, the hardware raises an exception
    - Trap handler is code within the OS that is written with the express purpose of handling TLB misses.

144

## TLB Control Flow algorithm (OS Handled)

- The hardware would do the following:

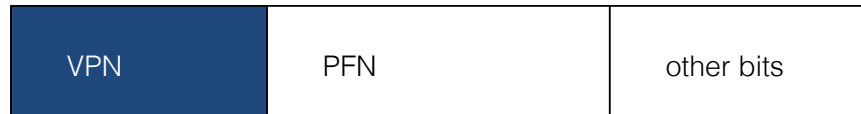
```
1:     VPN = (VirtualAddress & VPN_MASK) >> SHIFT
2:     (Success, TlbEntry) = TLB_Lookup(VPN)
3:     if (Success == True) // TLB Hit
4:         if (CanAccess(TlbEntry.ProtectBits) == True)
5:             Offset = VirtualAddress & OFFSET_MASK
6:             PhysAddr = (TlbEntry.PFN << SHIFT) | Offset
7:             Register = AccessMemory(PhysAddr)
8:         else
9:             RaiseException(PROTECTION_FAULT)
10:    else // TLB Miss
11:        RaiseException(TLB_MISS)
```

- But might be slow, why not just use the hardware approach?

145

# TLB entry

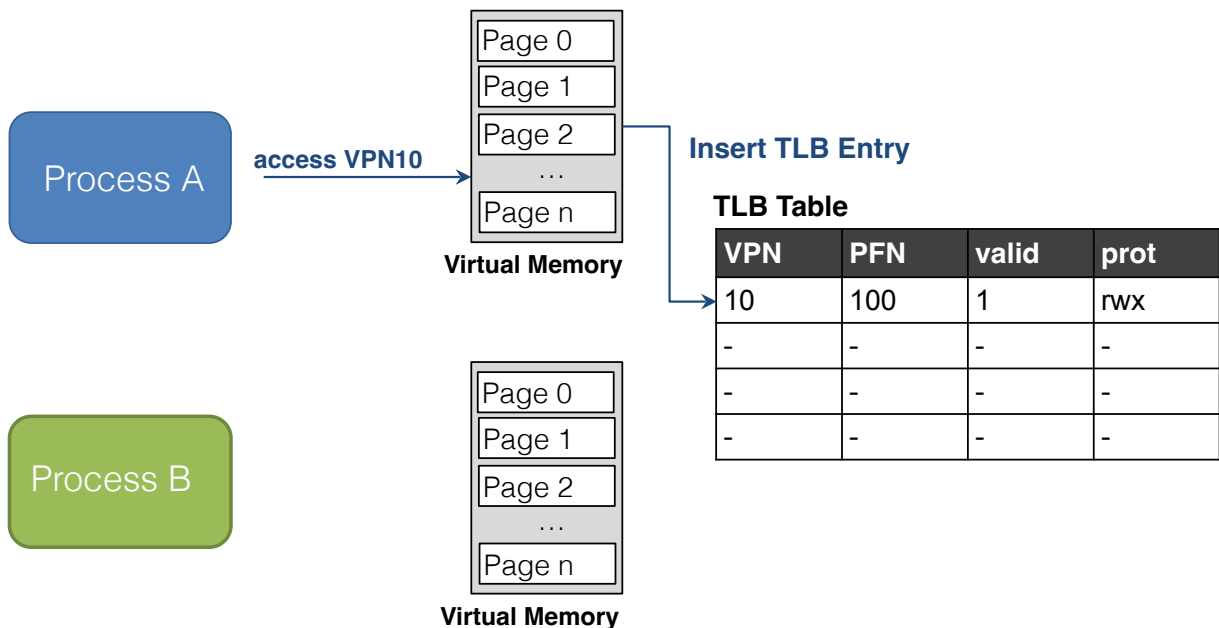
- TLB entries are often *fully associative* (any entry for any mapping)
  - A typical TLB might have 32, 64, or 128 entries.
  - Hardware searches the TLB in parallel to find the translation.
  - other bits: valid, protection, address-space identifier, dirty bit



Typical TLB entry

146

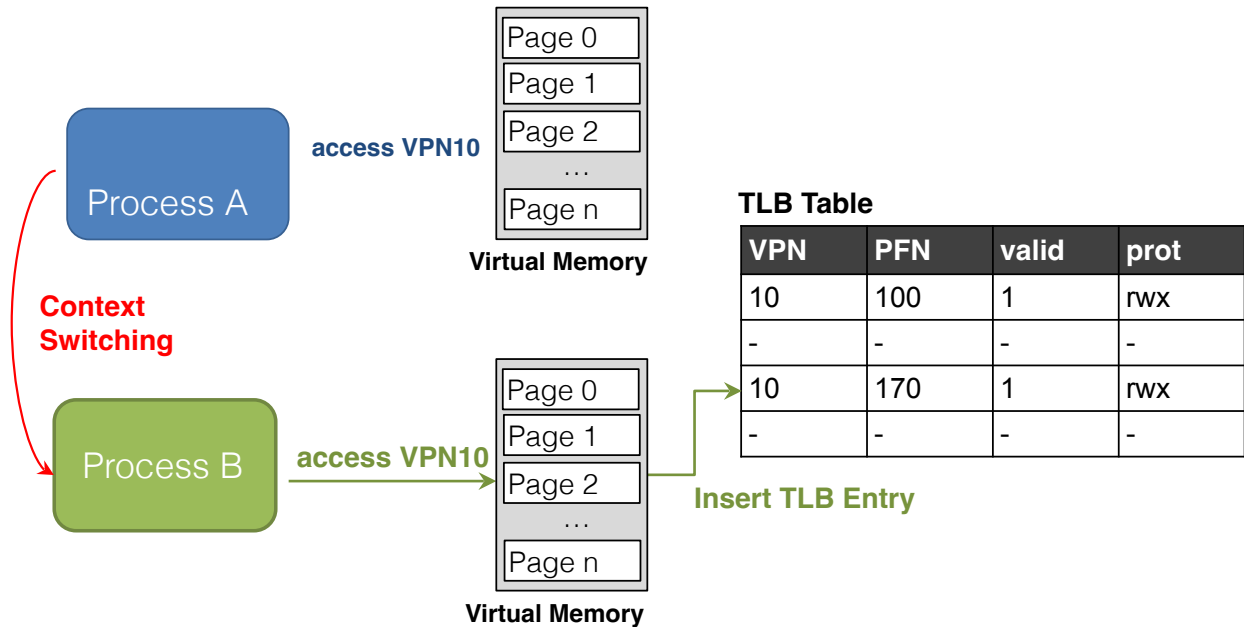
# TLB Issue: Context Switching



147

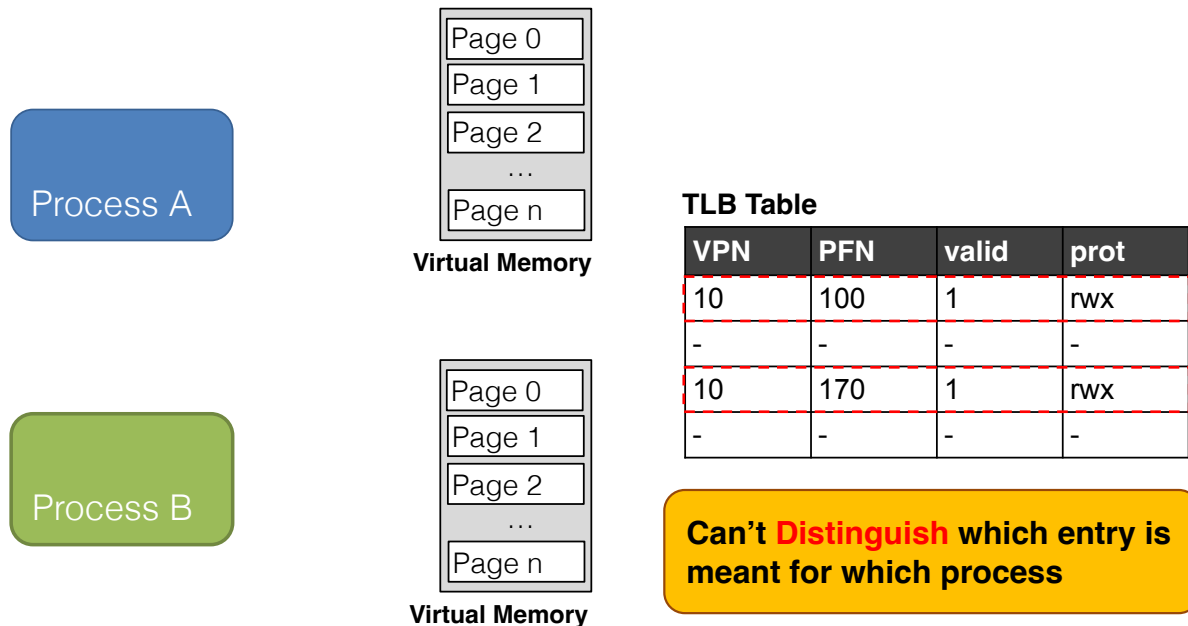


# TLB Issue: Context Switching



148

# TLB Issue: Context Switching

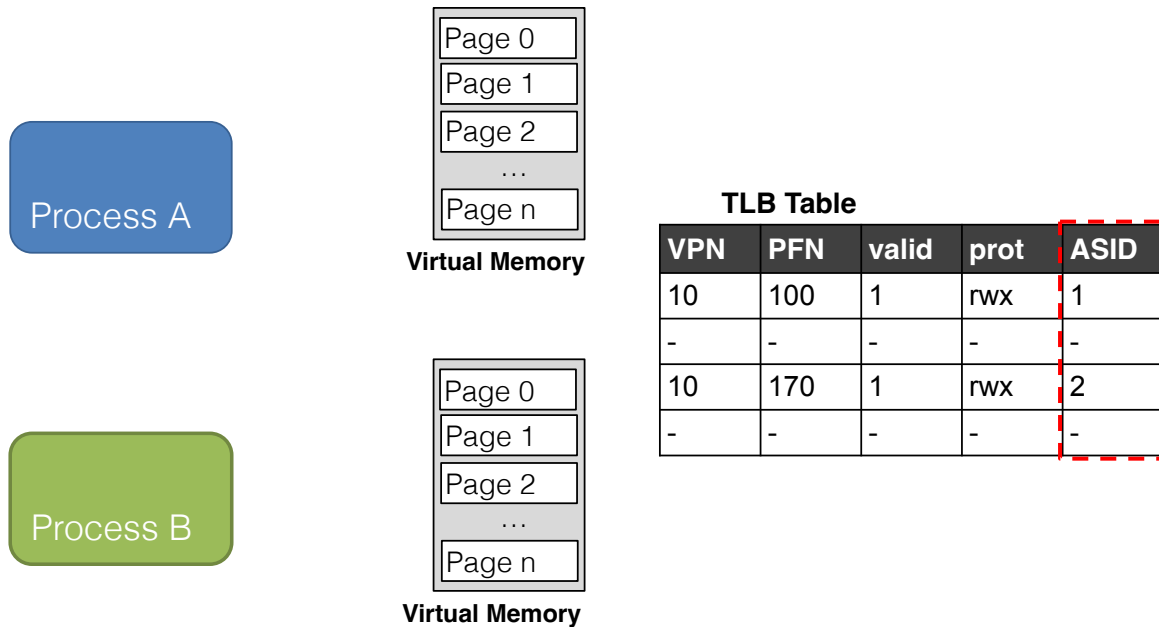


*Could just flush the TLB on every context switch...*

149

# Disambiguating Address Spaces

- Provide an address space identifier (ASID) field in the TLB.



150

## Another Case

- Two processes share a page.
  - Process 1 is sharing physical page 101 with Process 2.
  - P1 maps this page into the 10<sup>th</sup> page of its address space.
  - P2 maps this page to the 50<sup>th</sup> page of its address space.

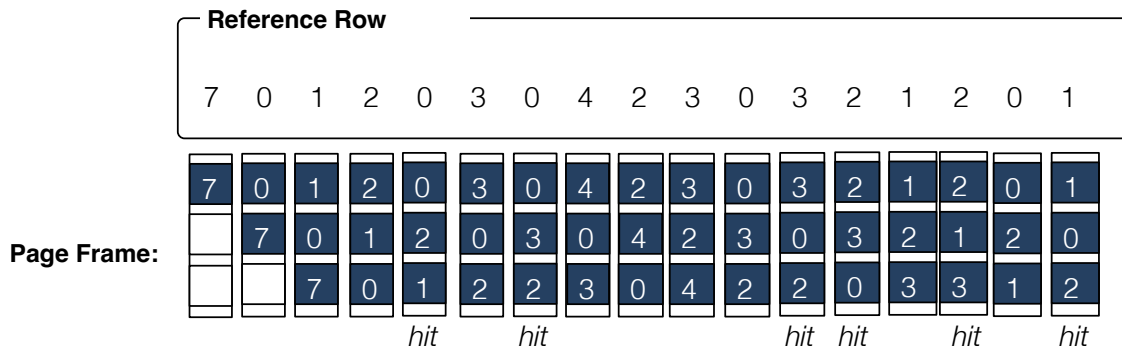
VPN	PFN	valid	prot	ASID
10	101	1	rwx	1
-	-	-	-	-
50	101	1	rwx	2
-	-	-	-	-

Sharing of pages is **useful** as it reduces the number of physical pages in use.

151

# TLB Replacement Policy

- LRU (Least Recently Used)
  - Evict an entry that has not recently been used.
  - Take advantage of *locality* in the memory-reference stream.

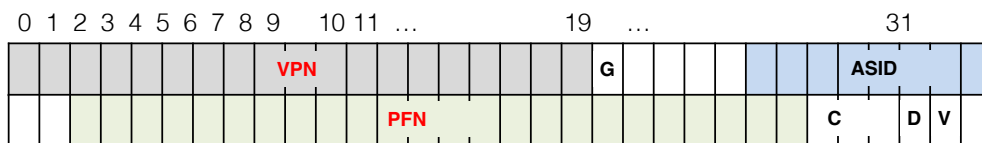


- 6 hits, 11 misses

152

# A Real TLB Entry

64-bit MIPS R4000 TLB entry



Flag	Content
19-bit VPN	The rest reserved for the kernel.
24-bit PFN	Systems can support with up to 64GB of main memory( pages ).
Global bit(G)	Used for pages that are globally-shared among processes.
ASID	OS can use to distinguish between address spaces.
Coherence bit(C)	determine how a page is cached by the hardware.
Dirty bit(D)	marking when the page has been written.
Valid bit(V)	tells the hardware if there is a valid translation present in the entry.

153

# Virtual Memory

- 13 - Address Spaces
- 14 - Memory API
- 15 - Address Translation
- 16 - Segmentation
- 17 - Free Space Management
- 18 - Paging
- 19 - Translation Lookaside Buffers
- 20 - Advanced Paging
- 21 - Swapping
- 22 - Swapping Policy