

## EE 457 Unit 7c

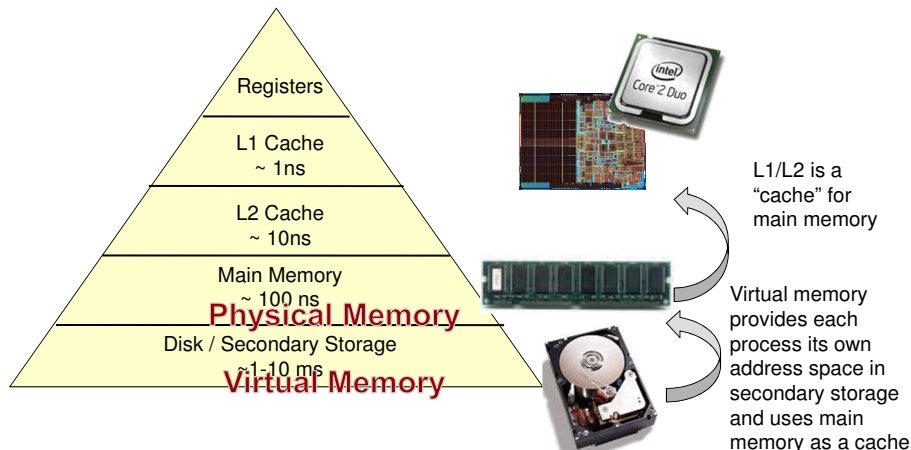
### Virtual Memory

## Virtual Memory Concept

- A mechanism for hiding the details of how much physical memory exists and how it's being shared
- Allows the OS to
  - Efficiently share the physical memory between several running programs/processes and provide protection against each other
  - Remove the need of the programmer to know how much memory is physically present and/or give the illusion of more or less physical memory than is present
- Use MM as a cache for multiple programs and their data as they run using secondary storage (hard-drive) as the home location

## Memory Hierarchy & Caching

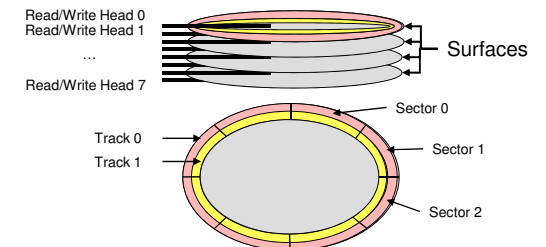
- Lower levels act as a cache for upper levels



## Virtual Memory Motivation

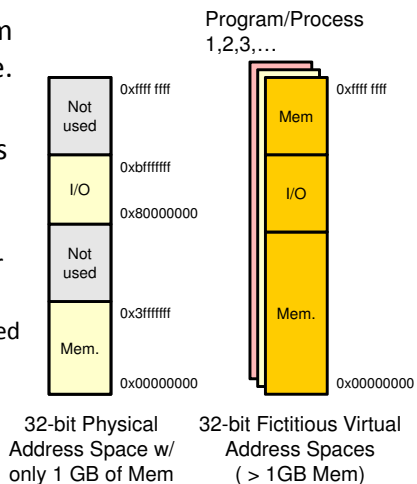
- Virtual memory is largely discussed in operating systems courses
  - We will focus on HW support for VM
- Magnetic hard drive consists of
  - Double sided surfaces/platters (with R/W head)
  - Each platter is divided into concentric tracks of small sectors that each store several thousand bits

• Seek Time: Time needed to position the read-head above the proper track	3-12 ms
• Rotational delay: Time needed to bring the right sector under the read-head	5-6 ms
• Transfer Time:	0.1 ms
• Disk Controller Overhead:	+ 2.0 ms
	~20 ms



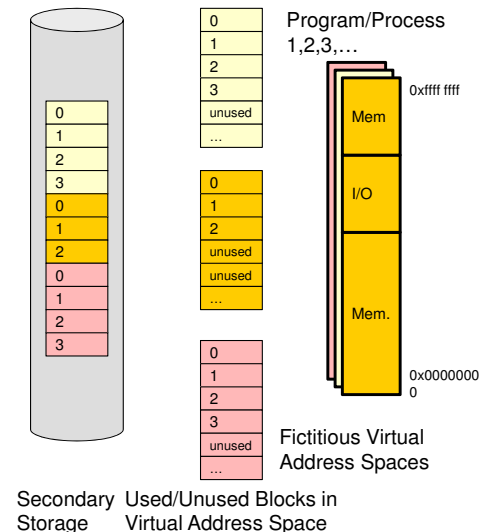
# Address Spaces

- Physical address spaces corresponds to the actual system address bus width and range (i.e. main memory and I/O)
- Each process/program runs in its own private "virtual" address space
  - Virtual address space can be larger (or smaller) than physical memory
  - Virtual address spaces are protected from each other



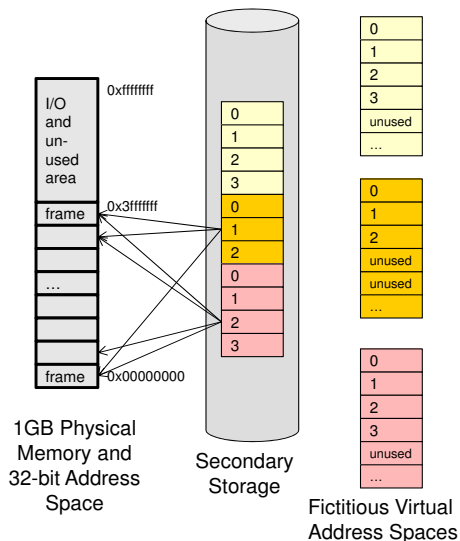
# Virtual Address Spaces

- Virtual address spaces are broken into blocks called "pages"
- Depending on the program, much of the virtual address space will not be used
- All used pages are "housed" in secondary storage (hard drive)



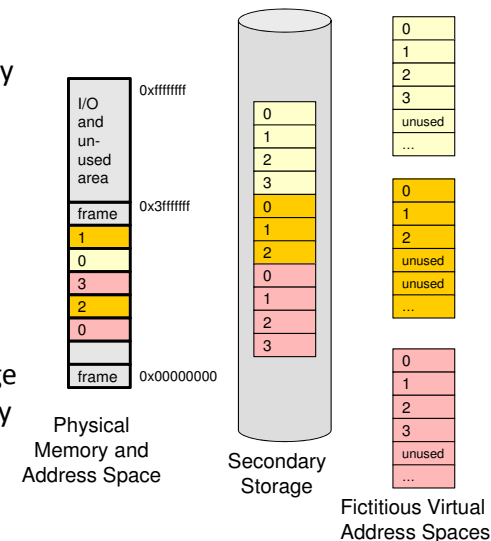
# Physical Address Space

- Physical memory is broken into page-size blocks called "page frames"
- Multiple programs are run concurrently and their pages (code & data) need to reside in physical memory
- Physical memory acts as a cache for pages from the secondary storage as each program executes



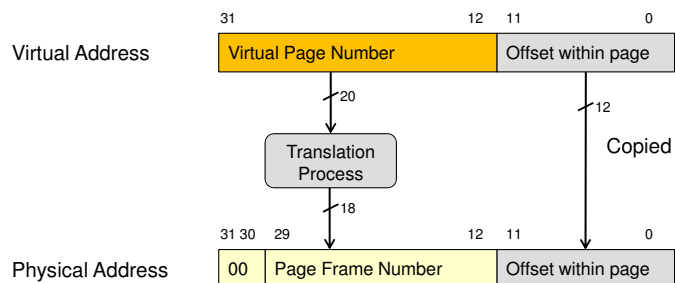
# Physical Memory Usage

- HW & the OS will translate the virtual addresses used by the program to the physical address where that page resides
- If an attempt is made to access a page that is not in physical memory, a "page fault exception" is declared and the OS brings in the page to physical memory (possibly evicting another page)



# Page Size and Address Translation

- Usually pages are several KB in size to amortize the large access time
- Example: 32-bit virtual & physical address, 1 GB physical memory, 4 KB pages
- Virtual page number to physical page frame translation performed by HW unit = MMU (Mem. Management Unit)

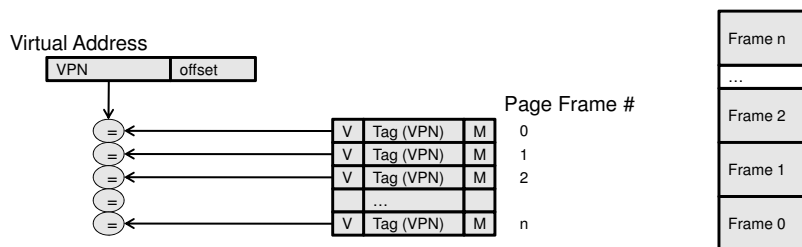


# VM Design Implications

- SLOW secondary storage access on page faults (10 ms)
  - Implies page size should be fairly large (i.e. once we've taken the time to find data on disk, make it worthwhile by accessing a reasonably large amount of data)
  - Implies the placement of pages in main memory should be fully associative to reduce conflicts and maximize page hit rates
  - Implies a "page fault" is going to take so much time to even access the data that we can handle them in software (via an exception) rather than using HW like typical cache misses
  - Implies eviction algorithms like LRU can be used since reducing page miss rates will pay off greatly
  - Implies write-back (write-through would be too expensive)

# Address Translation Issues

- A virtual page with 20-bit VPN can be sitting anywhere in the  $256K = 2^{18}$  page frames in physical memory
  - TAG = 20 + 1 = 21 bits
- This is impractical
- Instead, most systems implement full associativity using a look-up table = **PAGE TABLE**



# Analogy for Page Tables

- Suppose we want to build a caller-ID mechanism for your contacts on your cell phone
  - Let us assume 1000 contacts represented by a 3-digit integer (0-999) by the cell phone (this ID can be used to look up their names)
  - We want to use a simple Look-Up Table (LUT) to translate phone numbers to contact ID's, how shall we organize/index our LUT

① LUT indexed w/ contact ID	② Sorted LUT indexed w/ used phone #'s	③ LUT indexed w/ all possible phone #'s
000 213-745-9823 001 626-454-9985 002 818-329-1980 ... 999 323-823-7104	213-730-2198 436 213-745-9823 000 323-823-7104 999 ... 818-329-1980 002	000-000-0000 null .. 213-745-9823 000 ... 999-999-9999 null

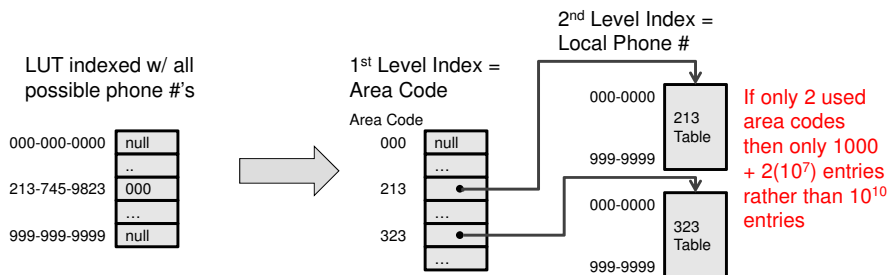
Doesn't Work  
We are given phone # and need to translate to ID (1000 accesses)

Could Work  
Since its in sorted order we could use a binary search ( $\log_2 1000$  accesses)

Could Work  
Easy to index & find but LARGE (1 access)

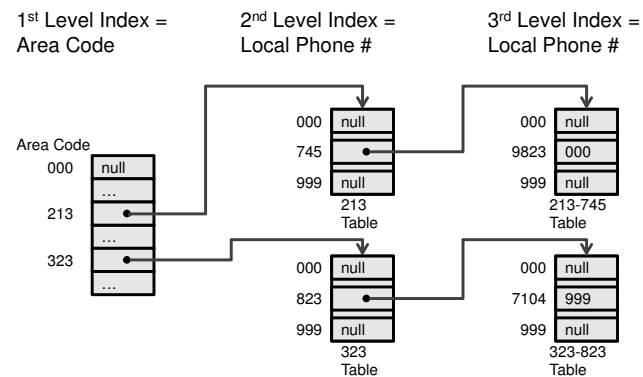
# Analogy for Page Tables

- Can we use the table indexed using all possible phone numbers (because it only requires 1 access) but somehow reduce the size especially since much of it is unused?
- Do you have friends from every area code? Likely contacts are clustered in only a few area codes.
- Use a 2-level organization
  - 1<sup>st</sup> level LUT is indexed on area code and contains pointers to 2<sup>nd</sup> level tables
  - 2<sup>nd</sup> level LUT's indexed on local phone numbers and contains contact ID entries



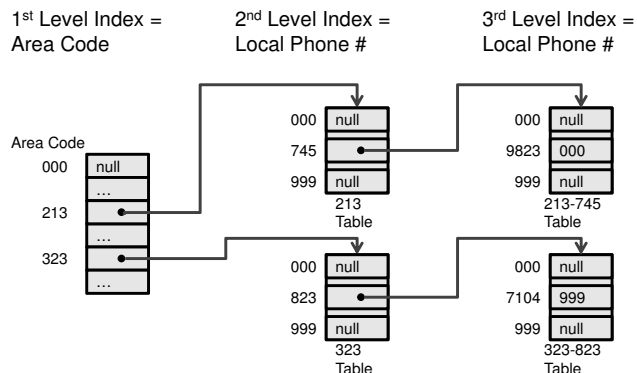
# Analogy for Page Tables

- Could extend to 3 levels if desired
  - 1<sup>st</sup> Level = Area code and pointers to 2<sup>nd</sup> level tables
  - 2<sup>nd</sup> Level = First 3-digits of local phone and pointers to 3<sup>rd</sup> level tables
  - 3<sup>rd</sup> Level = Contact ID's



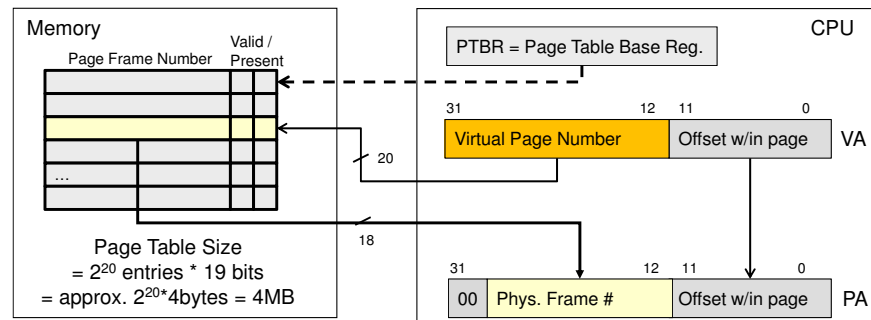
# Analogy for Page Tables

- If we add a friend from area code 408 we would have to add a second and third level table for just this entry



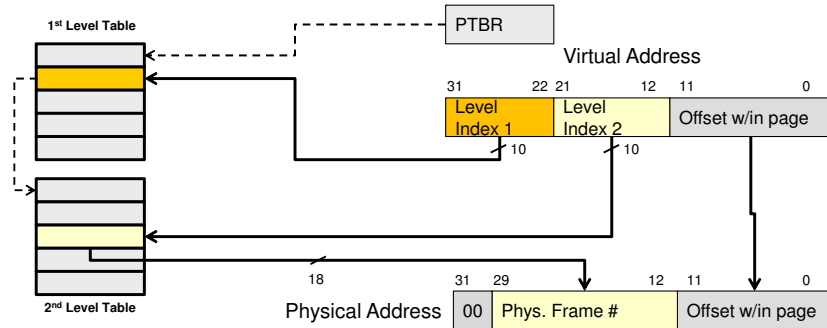
# Page Tables

- Page table is built by the OS and maintained in the physical memory at some chosen place by the OS
  - Allows virtual memory page placement to be fully associative in physical memory
  - One page table per process and indexed on virtual address
  - PTBR is a processor register pointing to the start address of currently executing process' page table

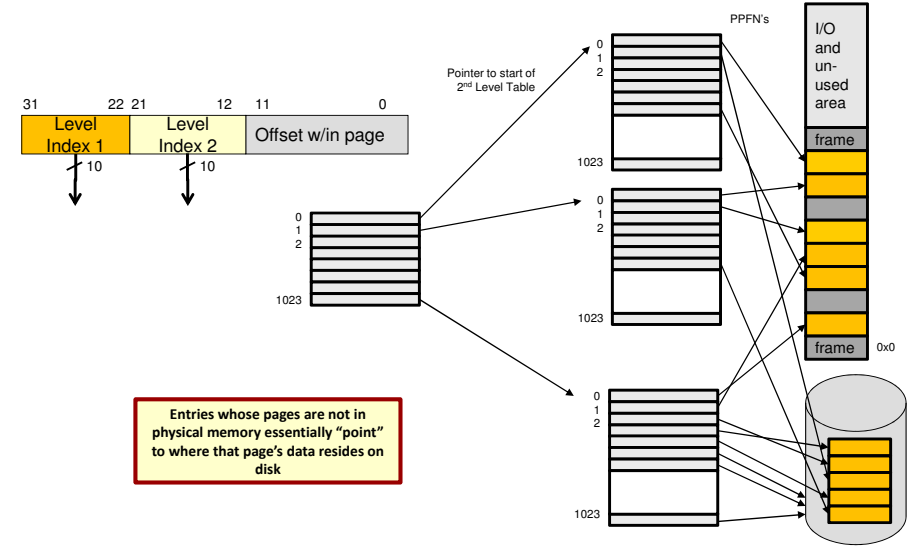


# Multi-Level Page Table

- VPN is broken into fields to index each level of the multi-level page table

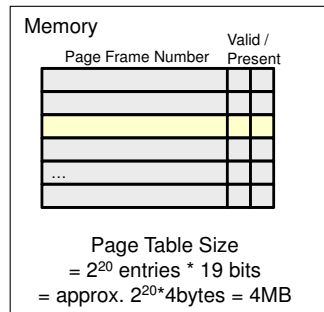


# Another View



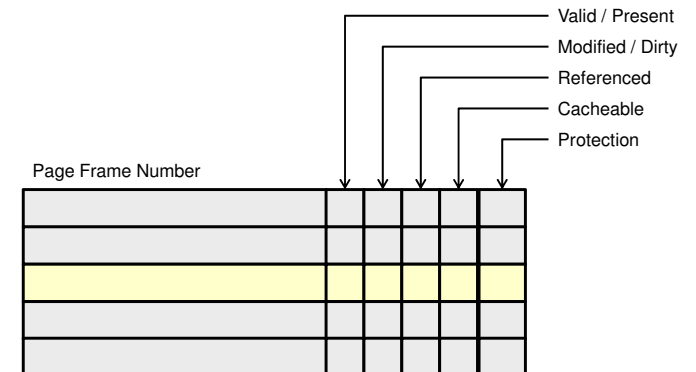
# To Tag or Not?

- Fully associative caches needed to store TAGs to check if the block is present.
- Do we need to store tags with the PPFN in the page table?
- Consider a book, assuming we start numbering pages at 1, do we need to print the page number along with the page content?
  - No since every page exists we can just count 1,2,3,...
  - Since we have an entry in the Page Table for every Virtual Page Number, we don't need to tag our entries



# Handling Page Faults

- Valid bit (1 = desired page in memory / 0 = page not present / page fault)
- Referenced = To implement pseudo-LRU replacement
- Protection: Read/Write/eXecute



## Page Fault Steps

- HW will...
  - Record the offending address, the EPC, and cause (page fault)
- SW will...
  - Pick an empty frame or select a page to evict
  - Writeback the evicted page if it has been modified
    - May block process while waiting and yield processor
  - Bring in the desired page and update the page table
    - May block process while waiting and yield processor
  - Restart the offending instruction

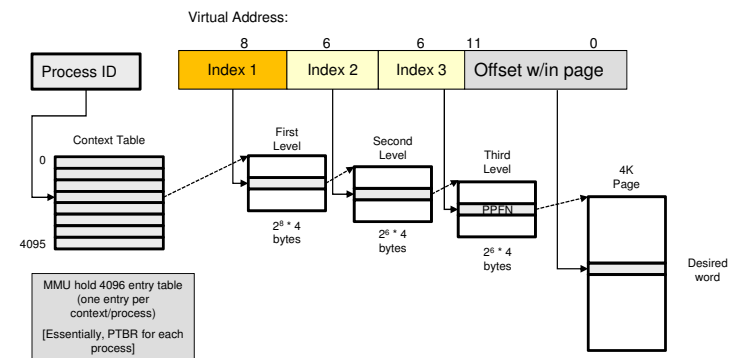
## Page Replacement Policies

- Possible algorithms: LRU, FIFO, Random
- Since page misses are so costly (slow) we can afford to spend sometime keeping statistics to implement LRU
- Implementing exact LRU would require updating statistics every access (using some kind of timestamp). This is too much to do in HW and we don't want to use SW when we have hits
- HW will implement simple mechanism that allows SW to implement a pseudo-LRU algorithm
  - HW will set the "Referenced" bit when a page is used
  - At certain intervals, SW will use these reference bits to keep statistics on which pages have been used in that interval and then clear the reference bits
  - On replacement, these statistics can be used to find the pseudo-LRU page

## Cache & VM Comparison

	Cache	Virtual Memory
<b>Block Size</b>	16-64B	4 KB – 64 MB
<b>Mapping Schemes</b>	Direct or Set Associative	Fully Associative
<b>Miss handling and replacement</b>	HW	SW
<b>Replacement Policy</b>	Full LRU if low associativity / Random is also used	Pseudo-LRU can be implemented

## SPARC VM Implementation



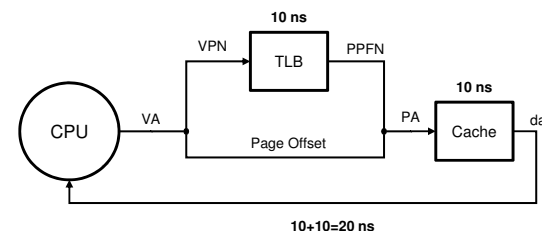
How many accesses to memory does it take to get the desired word that corresponds to the given virtual address? Would that change for a 1- or 2- level table?

# Performance Issues

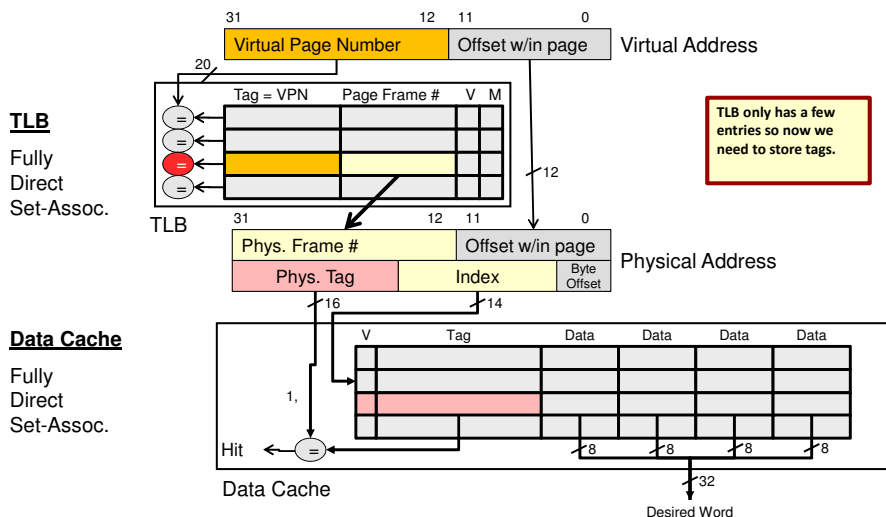
- Let cache hits = 10ns, memory accesses=100ns
- Assume a program makes an access to data located in cache...
  - Without VM, only requires 10ns cache access time
  - With VM, address must first be translated via the page table (recall page table is in memory)
    - If a single-level, one access to the page table (MM) = 100ns
    - If two-levels, two access to the page tables = 200ns
    - If three-levels, three access to the page tables = 300ns
    - Finally, physical address can access cache = 10 ns (if hit)
    - Total time equals 100\*L+10 (where L=# of Level of Page Table)
- Translation is extremely costly as currently implemented!!!

# Translation Lookaside Buffer (TLB)

- Solution: Let's create a cache for translations = Translation Lookaside Buffer (TLB)
- Needs to be small (64-128 entries) so it can be fast, with high degree of associativity (at least 4-way and many times fully associative) to avoid conflicts
  - On hit, the PPFN is produced and concatenated with the offset
  - On miss, a page table walk is needed

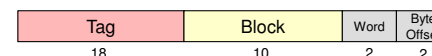


# Translation Lookaside Buffer (TLB)



# TLB Block Size

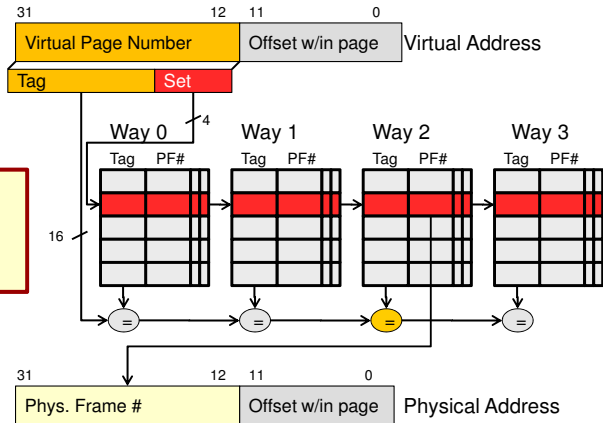
- A block in cache may be
  - 1 word
  - 2 words
  - 4 words
- Consider a direct mapped cache mapping can the word field be 0-bits?



- But an entry in the TLB is 1 value
  - $\log_2(1) = 0$ ...TLB mappings have no word field

# A 4-Way Set Associative TLB

- 64 entry 4-way SA TLB (wet field indexes each “way”)
  - On hit, page frame # supplied quickly w/o page table access



What is the page size? \_\_\_\_\_  
 Tag size? \_\_\_\_\_  
 Comparator Width? \_\_\_\_\_

# Virtual Memory System Examples

Microprocessor	AMD Opteron	P4	PPC 7447a
Virtual Address	48-bit	32- or 48-bit	52-bit
Physical Address	40-bit	36-bit	32- or 36-bit
TLB Entries (I/D/L2 TLB)	L1: 40/40 L2: 512/512	L1: 128/128	L1: 128 / 128
TLB Mapping	L1: Fully L2: 4-way SA	Fully (? 4-way)	2-way set associative
Min. Page Size	4 KB	4 KB	4 KB

Notes: Large VA's include ASID (process ID's) and other segment information  
 Sources: H&P, "CO&D", 3<sup>rd</sup> ed., Freescale.com,

# TLB Issues

- Because of high degree of associativity and limited working set of pages (usually) we can get VERY HIGH hit rates for the TLB
  - Variable page size settable by OS to allow for different working set sizes
  - Example: 64 TLB entries and 4 KB pages = 256KB
- Often times, separate TLB's for instruction and data address translation

# TLB Miss Process

- On a TLB miss, there is some division of work between the hardware (MMU) and OS
- Option 1
  - MMU can perform the TLB search followed by a page table walk if needed
  - If page fault occurs, OS takes over to bring in the page
- Option 2
  - MMU performs TLB Search
  - If TLB miss, OS can perform page table walk and bring in page if necessary
- When we want to remove a page from MM
  - First flush out blocks belong to that page from cache (writing back if necessary)
  - Invalidate tags of those blocks
  - Invalidate TLB entry (if any) corresponding to that page
    - If D=1, set dirty bit in page table
  - If page is dirty, copy page back to the disk
  - Simple way to remember this...
    - If parents (page) leave a party then the children (cache blocks & TLB entries) must leave too



## Other Issues

- Property of Inclusion
  - Cache contents are a (subset / superset) of main memory contents
  - Main memory contents are a (subset / superset) of page/swap file on disk
  - TLB contents are a (subset / superset) of page table contents

## Cache, VM, and Main Memory

TLB	VM	Cache	Possible Y/N & Description
Hit	Hit	Hit	Possible: best-case scenario
Hit	Hit	Miss	Possible: TLB hits (hit in VM is implied), then cache miss
Miss	Hit	Hit	TLB misses, then hits in page table, then cache hit
Miss	Hit	Miss	TLB misses, then hits in page table, then cache miss
Miss	Miss	Miss	TLB misses, then page fault, then miss in cache
Hit	Miss	Miss	Impossible: cannot hit in TLB if page not present
Hit	Miss	Hit	Impossible: cannot hit in TLB if page not present
Miss	Miss	Hit	Impossible: data cannot be in cache if page not present

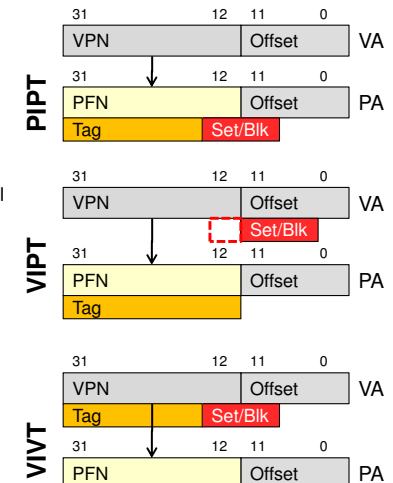
Taken from H & P, "Computer Organization" 3<sup>rd</sup>, Ed.

## Cache Addressing with VM

- Cache review
  - Set or block field indexes LUT holding tags
  - 2 steps to determine hit:
    - Index (lookup) to find tags (using block or set bits)
    - Compare tags to determine hit
    - Sequential connection between indexing and tag comparison
- Rather than waiting for address translation and then performing this two step hit process, can we overlap the translation and portions of the hit sequence?
  - Yes if we choose page size, block size, and set/direct mapping carefully

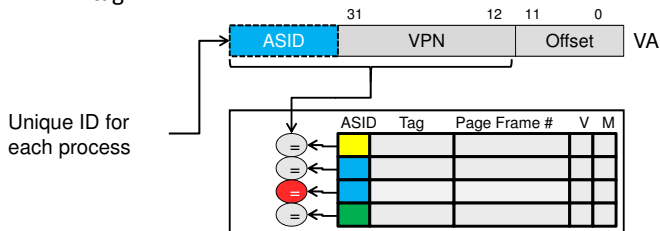
## Cache Index/Tag Options

- Physically indexed, physically tagged (PIPT)
  - Wait for full address translation
  - Then use physical address for both indexing and tag comparison
- Virtually indexed, physically tagged (VIPT)
  - Use portion of the virtual address for indexing then wait for address translation and use physical address for tag comparisons
  - Easiest when index portion of virtual address w/in offset (page size) address bits, otherwise aliasing may occur
- Virtually indexed, virtually tagged (VIVT)
  - Use virtual address for both indexing and tagging...No TLB access unless cache miss
  - Requires invalidation of cache lines on context switch or use of process ID as part of tags



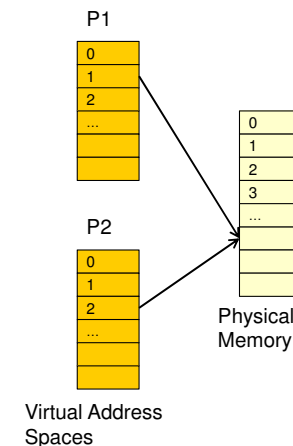
## Multiple Processes

- Recall each process has its own virtual address space, page table, and translations
- How does TLB handle context switch
  - Can choose to only hold translations for current process and thus invalidate all entries on context switch
  - Can hold translations for multiple processes concurrently by concatenating a process or address space ID (PID or ASID) to the VPN tag



## Shared Memory

- In current system, all memory is private to each process
- To share memory between two processes, the OS can allocate an entry in each process' page table to point to the same physical page
- Can use different protection bits for each page table entry (e.g. P1 can be R/W while P2 can be read only)

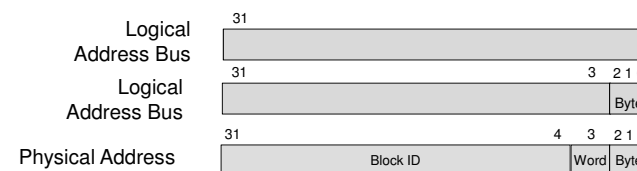


## A Complete VM / Cache Example

- Use the following specification for the following questions
  - 64-bit data, 32-bit virtual/physical address
  - Page Size: 128KB
  - TLB Size: 256 entry 4-way set associative
  - Page Table Org.: 3-levels
    - A 64 entry A-Table (page directory) followed by several 32 entry B-Tables (2<sup>nd</sup> level tables) followed by some number of C-Tables (3<sup>rd</sup> level)
  - Cache Organization
    - Cache Size: 512KB
    - 8-way set associative
    - Block size: 2 words [Word = 64-bits = 8 bytes]

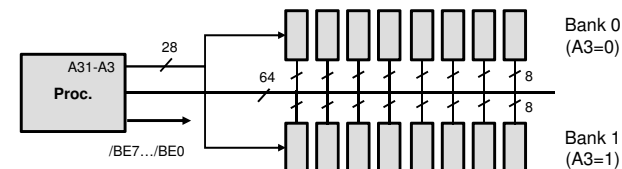
## Address Bus and Interleaving

- Use the following specification for the following questions
  - 64-bit data, 32-bit virtual/physical address
  - Cache Organization: Block size: 2 words [1 Word = 64-bits = 8 bytes]
- How many banks would you suggest for interleaving purposes?



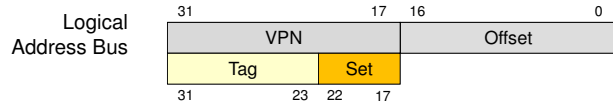
64-bit Data bus = 8 bytes = 8 Byte enables (/BE7.../BE0)

2 Banks so we can quickly get two words to the data cache when a block is transferred



## TLB Mapping

- Use the following specification for the following questions
  - 64-bit data, 32-bit virtual/physical address
  - Page Size: 128KB
  - TLB Size: 256 entry 4-way set associative

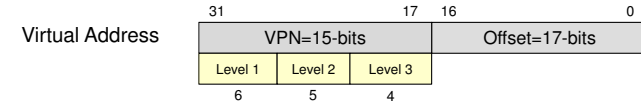


Page Size = 128KB =>  
Offset Field Size = 17-bits

# of TLB Sets: 256 entries / 4 entries per set  
= 64 sets => 6 set bits

## Page Table Mapping

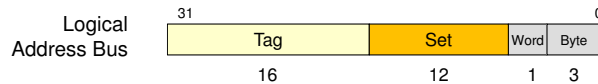
- Use the following specification for the following questions
  - Page Size: 128KB
  - Page Table Org.: 3-levels
    - A 64 entry A-Table (page directory)
    - 32 entry B-Tables (2<sup>nd</sup> level tables)
    - some number of C-Tables (3<sup>rd</sup> level)



VPN = 15-bits  
Level 1 Page Table = 64 (2<sup>6</sup>) entries => 6-bits  
Level 2 Page Table = 32 (2<sup>5</sup>) entries => 5-bits  
Level 3 Page Table = 15-6-5= 4-bits => 16 (2<sup>4</sup>) entries

## Data Cache Design

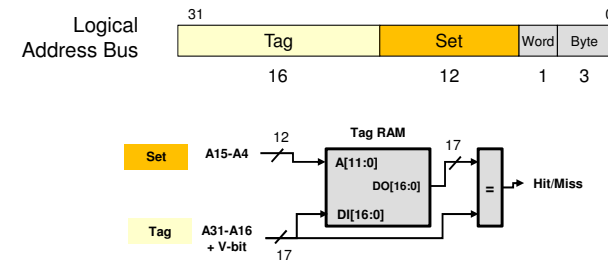
- Use the following specification for the following questions
  - Cache Organization
    - Cache Size: 512KB
    - 8-way set associative
    - Clock size: 2 words [Word = 64-bits = 8 bytes]



64-bit data = 8 (2<sup>3</sup>) bytes per word  
Block Size = 2 (2<sup>1</sup>) words = 1 word bit  
# of Cache blocks = 512KB / 16 bytes per block  
= 2<sup>19</sup> / 2<sup>4</sup> = 2<sup>15</sup>  
# of Sets = 2<sup>15</sup> / 2<sup>3</sup> ways per set = 2<sup>12</sup> sets => 12 set bits  
# of Tag bits = 32 - 12 - 1 - 3 bits = 16-bits

## Data Cache Implementation

- How many comparators and of what size are needed to determine cache hit or miss?
- What is the size of the TAG RAM's?



8 comparators of 17-bits  
Tag RAM Size = 4K x 17