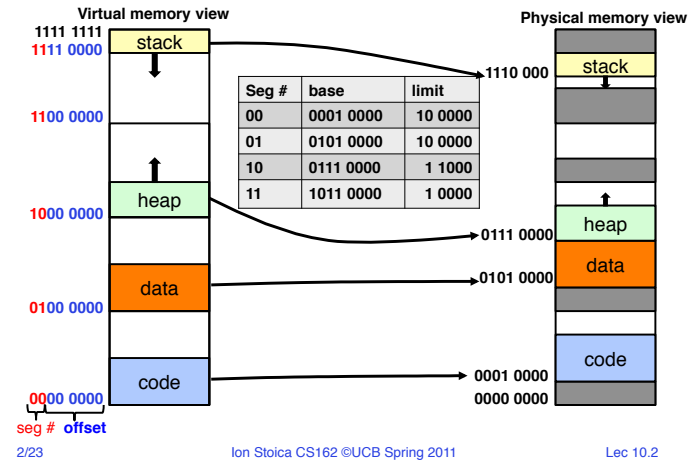


CS162 Operating Systems and Systems Programming Lecture 10

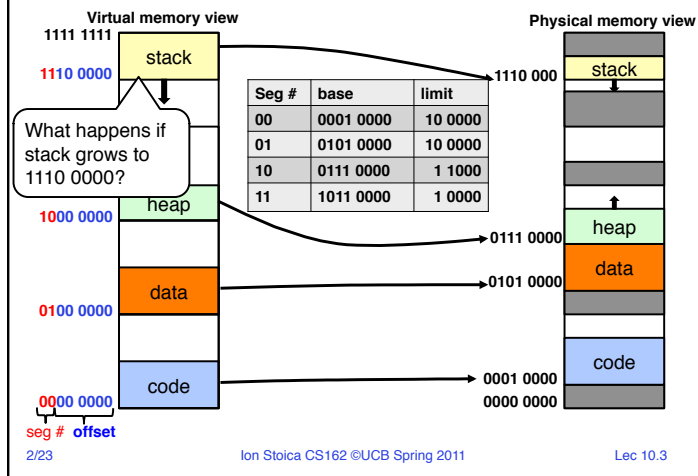
Caches and TLBs

February 23, 2011
Ion Stoica
<http://inst.eecs.berkeley.edu/~cs162>

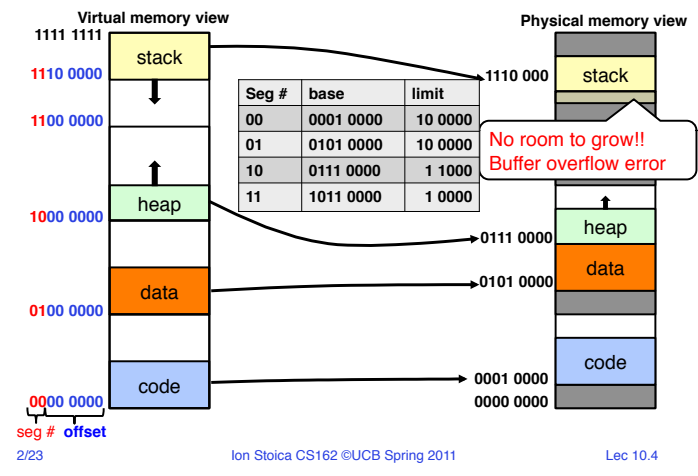
Review: Address Segmentation



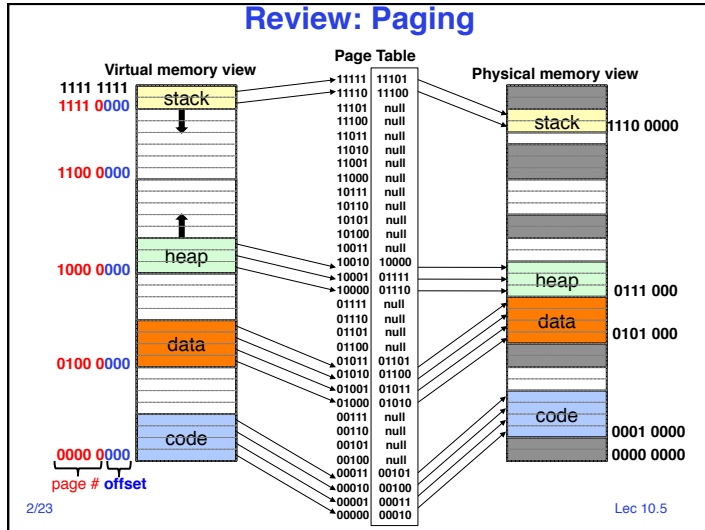
Review: Address Segmentation



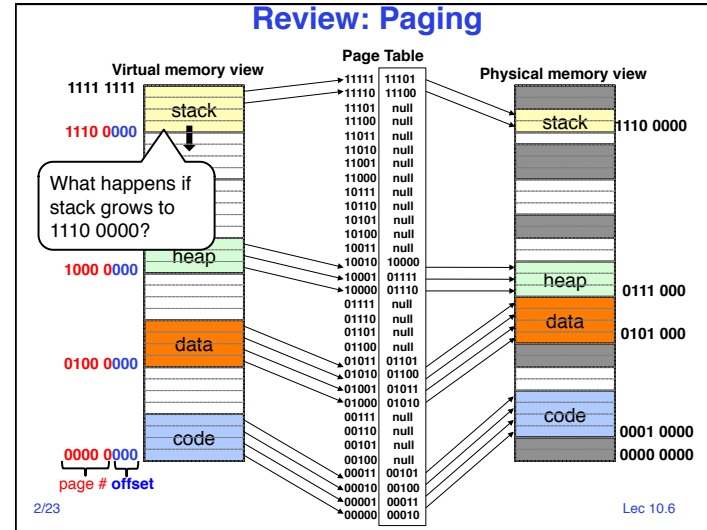
Review: Address Segmentation



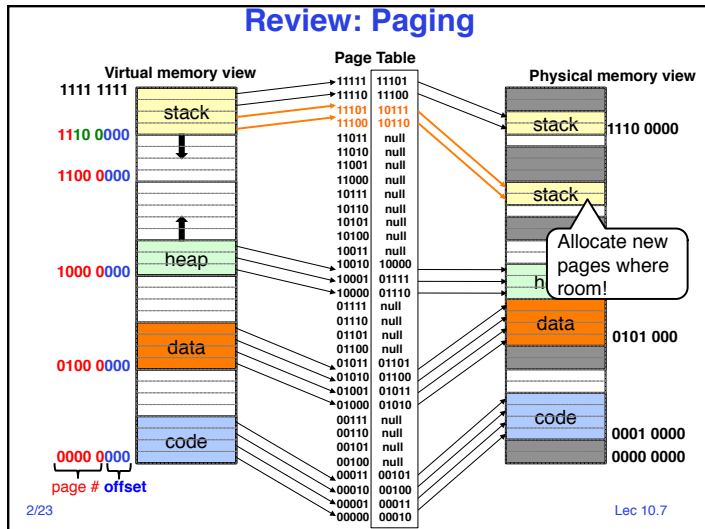
Review: Paging



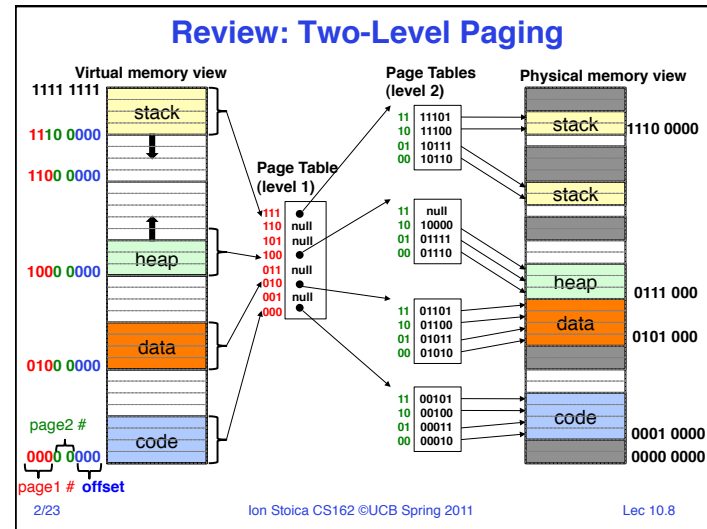
Review: Paging

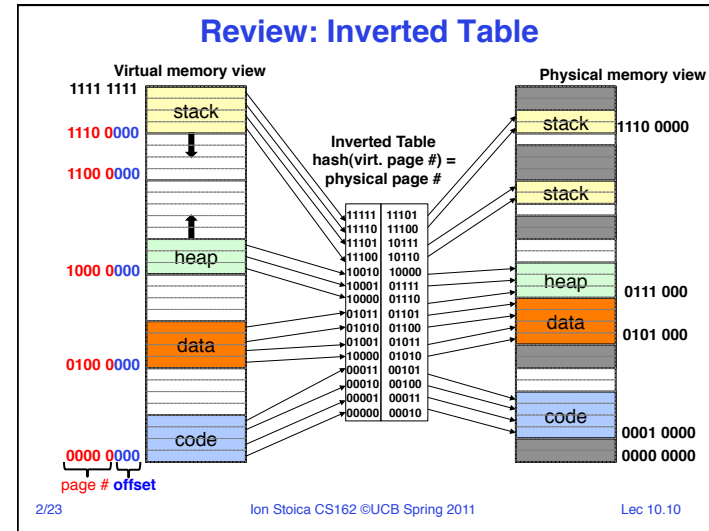
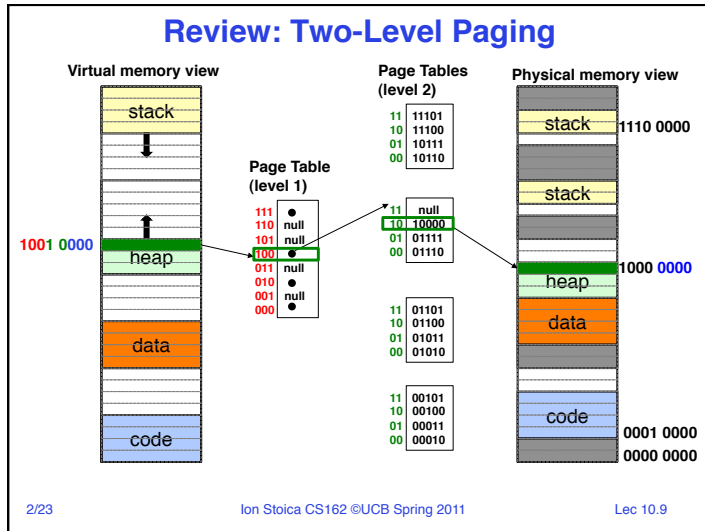


Review: Paging



Review: Two-Level Paging

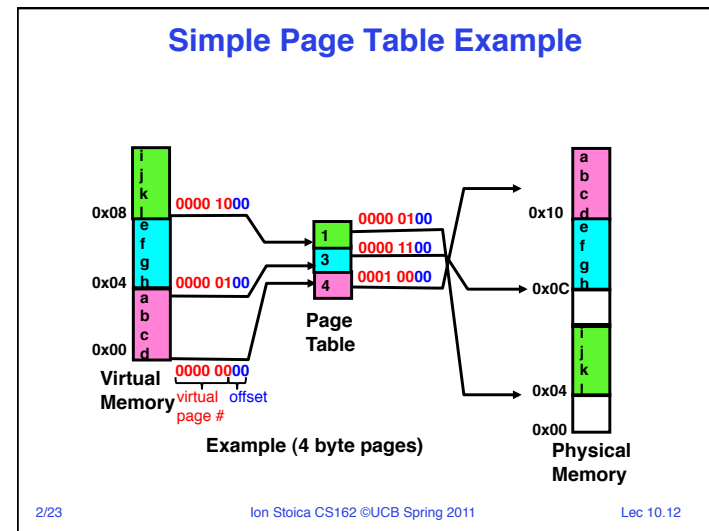




Address Translation Comparison

	Advantages	Disadvantages
Segmentation	Fast context switching: Segment mapping maintained by CPU	External fragmentation
Paging (single-level page)	No external fragmentation	<ul style="list-style-type: none"> • Large size: Table size ~ virtual memory • Internal fragmentation
Paged segmentation	<ul style="list-style-type: none"> • No external fragmentation 	<ul style="list-style-type: none"> • Multiple memory references per page access
Two-level pages	• Table size ~ memory used by program	• Internal fragmentation
Inverted Table		Hash function more complex

2/23 Ion Stoica CS162 ©UCB Spring 2011 Lec 10.11



How is the translation accomplished?



- What, exactly happens inside Memory Management Unit (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 invalid Page Table Entry (PTE)
 - » Fault handler will decide what to do (see next)
 - 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!

2/23

Ion Stoica CS162 ©UCB Spring 2011

Lec 10.13

Goals for Today

- Caching
 - Misses
 - Organization
- Translation Look aside Buffers (TLBs)

Note: Some slides and/or pictures in the following are adapted from slides ©2005 Silberschatz, Galvin, and Gagne. Many slides generated from lecture notes by Kubiawicz.

2/23

Ion Stoica CS162 ©UCB Spring 2011

Lec 10.14

Caching Concept



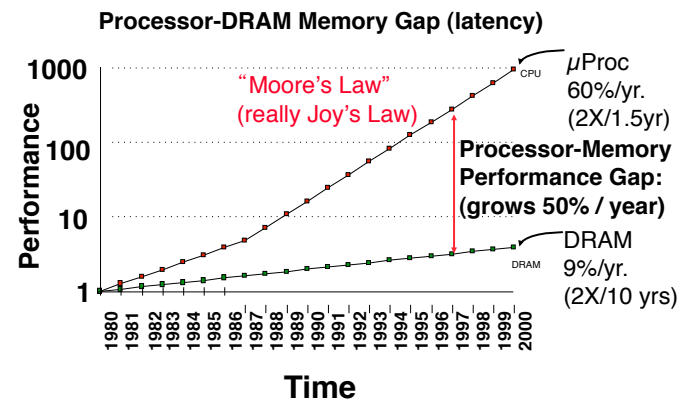
- **Cache**: a repository for copies that can be accessed more quickly than the original
 - Make frequent case fast and infrequent case less dominant
- Caching underlies many of the techniques that are used today to make computers fast
 - Can cache: memory locations, address translations, pages, file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

2/23

Ion Stoica CS162 ©UCB Spring 2011

Lec 10.15

Why Bother with Caching?



2/23

Ion Stoica CS162 ©UCB Spring 2011

Lec 10.16

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 Address: Physical Page #, Offset

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!
- Solution? Cache translations!
 - Translation Cache: TLB ("Translation Lookaside Buffer")

2/23 Ion Stoica CS162 ©UCB Spring 2011 Lec 10.17

Why Does Caching Help? Locality!

- 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

2/23 Ion Stoica CS162 ©UCB Spring 2011 Lec 10.18

Review: Memory Hierarchy

- 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

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

2/23 Ion Stoica CS162 ©UCB Spring 2011 Lec 10.19

Example

- Data in memory, no cache:**

Access time = 100ns
- Data in memory, 10ns cache:**

Average Access time = (Hit Rate x HitTime) + (Miss Rate x MissTime)

 - HitRate + MissRate = 1
 - HitRate = 90% → Average Access Time = 19ns
 - HitRate = 99% → Average Access Time = 10.9ns

2/23 Ion Stoica CS162 ©UCB Spring 2011 Lec 10.20

Review: Sources of Cache Misses

- **Compulsory** (cold start): first reference to a block
 - “Cold” fact of life: not a whole lot you can do about it
 - Note: When running “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 same cache location
 - Solutions: increase cache size, or increase associativity
- **Two others:**
 - **Coherence** (Invalidation): other process (e.g., I/O) updates memory
 - **Policy:** Due to non-optimal replacement policy

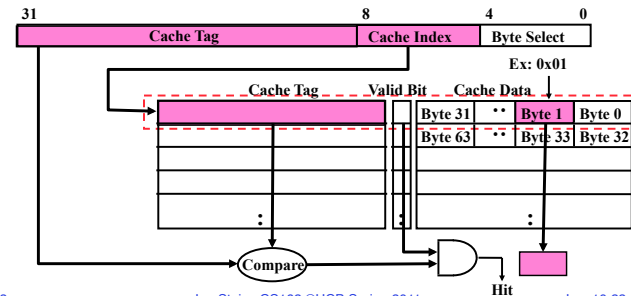
2/23

Ion Stoica CS162 ©UCB Spring 2011

Lec 10.21

Direct Mapped Cache

- Cache index selects a cache block
- “Byte select” selects byte within cache block
 - Example: Block Size=32B blocks
- Cache tag fully identifies the cached data
- Data with same “cache tag” shares the same cache entry
 - Conflict misses



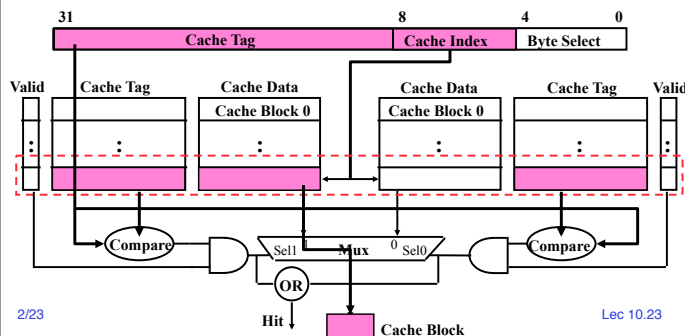
2/23

Ion Stoica CS162 ©UCB Spring 2011

Lec 10.22

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
 - Two tags in the set are compared to input in parallel
 - Data is selected based on the tag result

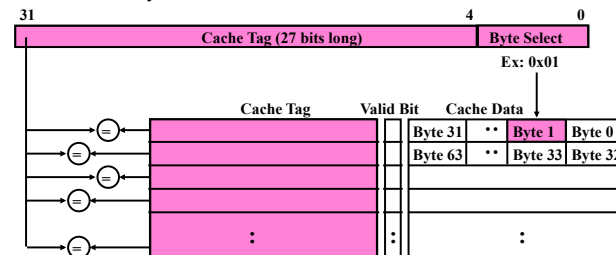


2/23

Lec 10.23

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



2/23

Ion Stoica CS162 ©UCB Spring 2011

Lec 10.24

Administrivia

- Project 1 due:
 - Code: Tuesday, **March 1st**
 - Final document, peer evaluation: Wednesday, **March 2nd**

- Project 2 starts after you are done with Project 1

- Peer evaluation:
 - **Individually** submit a form with a quantitative assessment of the relative contributions of each of your teammates, **excluding yourself**
 - You have 100 points to allocate among your teammates, proportional to how much they contributed to **just this assignment**.

2/23

Ion Stoica CS162 ©UCB Spring 2011

Lec 10.25

5min Break

2/23

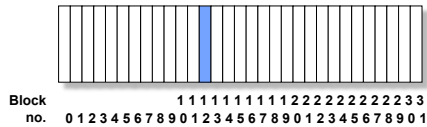
Ion Stoica CS162 ©UCB Spring 2011

Lec 10.26

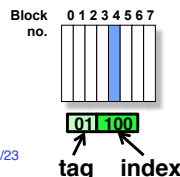
Where does a Block Get Placed in a Cache?

- Example: Block 12 placed in 8 block cache

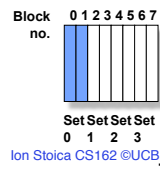
32-Block Address Space:



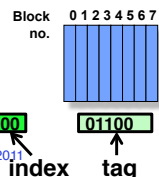
Direct mapped:
block 12 (01100) can go only into block 4 (12 mod 8)



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



Fully associative:
block 12 can go anywhere



2/23

Ion Stoica CS162 ©UCB Spring 2011

Lec 10.27

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)

Size	2-way		4-way		8-way	
	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%

2/23

Ion Stoica CS162 ©UCB Spring 2011

Lec 10.28

What happens on a write?

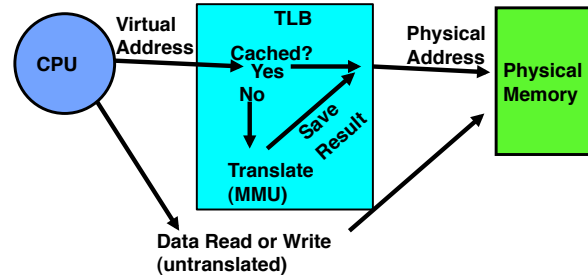
- **Write through:** The information is written both to the block in the cache and to the block in the lower-level memory
- **Write back:** The information is written only to the block in the cache.
 - Modified cache block is written to main memory only when it is replaced
 - Question is block clean 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

2/23

Ion Stoica CS162 ©UCB Spring 2011

Lec 10.29

Caching Applied to Address Translation



- 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

2/23

Ion Stoica CS162 ©UCB Spring 2011

Lec 10.30

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
 - 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

2/23

Ion Stoica CS162 ©UCB Spring 2011

Lec 10.31

What happens on a Context Switch?

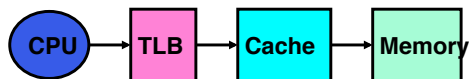
- Need to do something, since TLBs map virtual addresses to physical addresses
 - Address Space just changed, so TLB entries no longer valid!
- Options?
 - Invalidate TLB: simple but might be expensive
 - » What if switching frequently between processes?
 - Include ProcessID in TLB
 - » This is an architectural 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!

2/23

Ion Stoica CS162 ©UCB Spring 2011

Lec 10.32

What TLB organization makes sense?



- Needs to be really fast
 - Critical path of memory access
 - 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

2/23

Ion Stoica CS162 ©UCB Spring 2011

Lec 10.33

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"
- When does TLB lookup occur?
 - Before cache lookup?
 - In parallel with cache lookup?

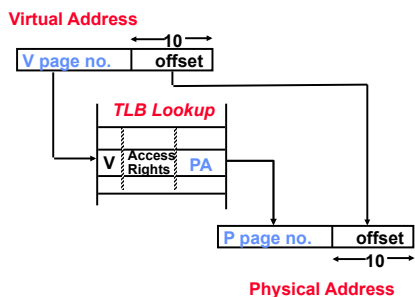
2/23

Ion Stoica CS162 ©UCB Spring 2011

Lec 10.34

Reducing translation time further

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



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

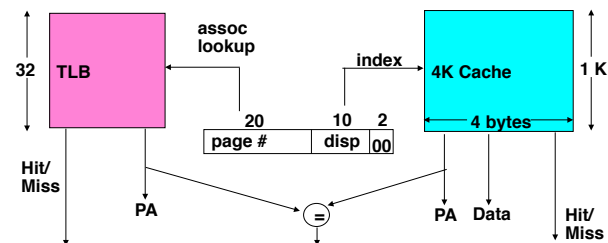
2/23

Ion Stoica CS162 ©UCB Spring 2011

Lec 10.35

Overlapping TLB & Cache Access

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



- 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

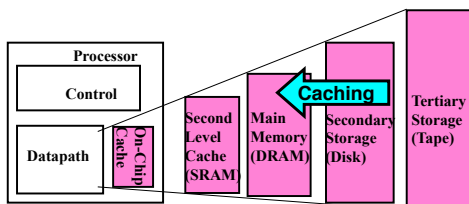
2/23

Ion Stoica CS162 ©UCB Spring 2011

Lec 10.36

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

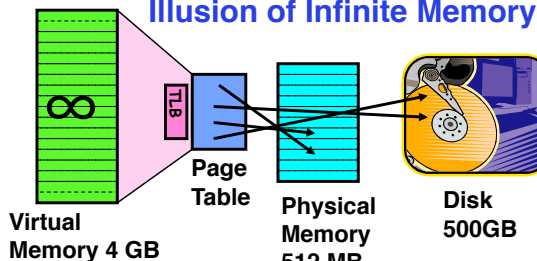


2/23

Ion Stoica CS162 ©UCB Spring 2011

Lec 10.37

Illusion of Infinite Memory



- Disk is larger than physical memory \Rightarrow
 - In-use virtual memory can be bigger than physical memory
 - Combined memory of running processes much larger than physical memory
- 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

2/23

Ion Stoica CS162 ©UCB Spring 2011

Lec 10.38

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 \rightarrow 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 a "dirty" bit (D)!

2/23

Ion Stoica CS162 ©UCB Spring 2011

Lec 10.39

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

Cache

2/23

Ion Stoica CS162 ©UCB Spring 2011

Lec 10.40

Summary (1/2)

- 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

2/23

Ion Stoica CS162 ©UCB Spring 2011

Lec 10.41

Summary (2/2)

- 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

2/23

Ion Stoica CS162 ©UCB Spring 2011

Lec 10.42