

CS194-24  
Advanced Operating Systems  
Structures and Implementation  
Lecture 10

Virtual Memory, TLBs

February 26<sup>th</sup>, 2014  
Prof. John Kubiawicz  
<http://inst.eecs.berkeley.edu/~cs194-24>

Goals for Today

- Virtual Memory
- TLBs

Interactive is important!  
Ask Questions!

Note: Some slides and/or pictures in the following are adapted from slides ©2013

2/26/14

Kubiawicz CS194-24 ©UCB Fall 2014

Lec 10.2

Recall: Definition of Monitor

- **Monitor**: a lock and zero or more condition variables for managing concurrent access to shared data
  - Use of Monitors is a programming paradigm
- **Lock**: provides mutual exclusion to shared data:
  - Always acquire before accessing shared data structure
  - Always release after finishing with shared data
- **Condition Variable**: a queue of threads waiting for something *inside* a critical section
  - Key idea: allow sleeping inside critical section by atomically releasing lock at time we go to sleep
  - Contrast to semaphores: Can't wait inside critical section

2/26/14

Kubiawicz CS194-24 ©UCB Fall 2014

Lec 10.3

Recall: Programming with Monitors

- Monitors represent the logic of the program
  - Wait if necessary
  - Signal when change something so any waiting threads can proceed
- Basic structure of monitor-based program:

```
lock
while (need to wait) {
    condvar.wait();
}
unlock
```

} Check and/or update state variables  
Wait if necessary

do something so no need to wait

```
lock

condvar.signal();

unlock
```

} Check and/or update state variables

2/26/14

Kubiawicz CS194-24 ©UCB Fall 2014

Lec 10.4

## pThreads Monitors (Mutex + Condition Vars)

- To create a mutex:

```
pthread_mutex_t amutex = PTHREAD_MUTEX_INITIALIZER;
pthread_mutex_init(&amutex, NULL);
```

- To use it:

```
int pthread_mutex_lock(amutex);
int pthread_mutex_unlock(amutex);
```

- To create condition variables:

```
pthread_cond_t mycondvar = PTHREAD_COND_INITIALIZER;
pthread_cond_init(&mycondvar, NULL);
```

- To use them:

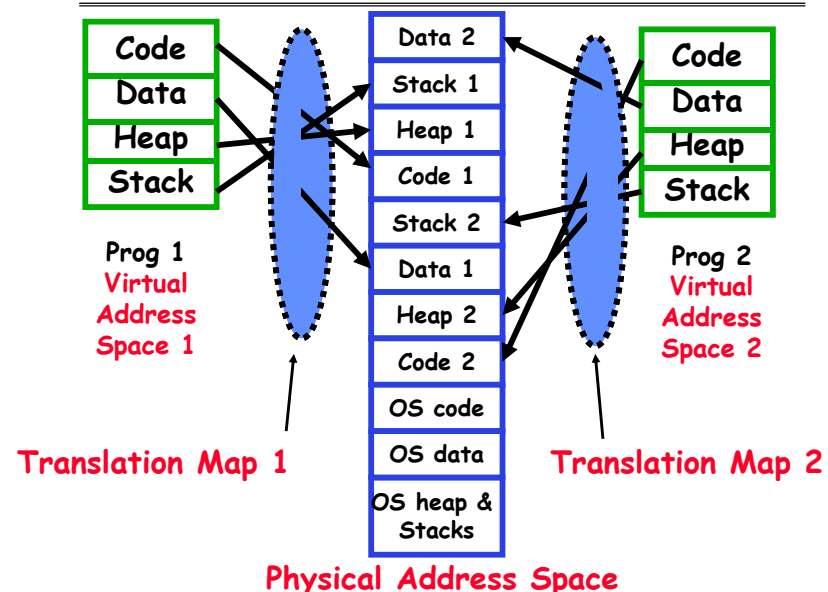
```
pthread_cond_wait(mycondvar, amutex);
pthread_cond_signal(mycondvar);
pthread_cond_broadcast(mycondvar);
```

2/26/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 10.5

## Recall: Example of Address Translation

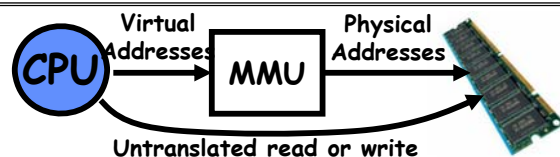


2/26/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 10.6

## Two Views of Memory with Dynamic Translation



- Two views of memory:

- View from the CPU (what program sees, virtual memory)
- View from 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

2/26/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 10.7

## 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

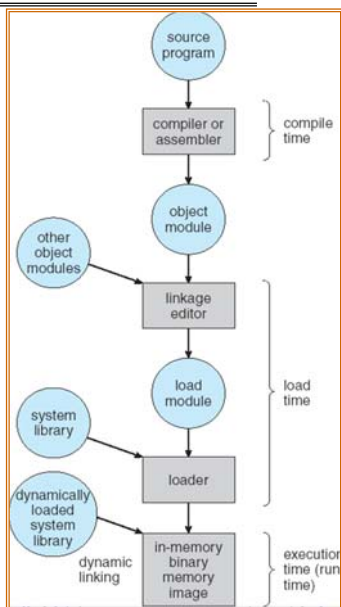
2/26/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 10.8

## 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



2/26/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 10.9

## 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
        ...
start: lw r1,0(data1)
        jal checkit
loop:  addi r1, r1, -1
        bnz r1, r0, loop
        ...
checkit: ...
    
```
- Memory addresses shown on the right:
- ```

0x300 00000020
...
0x900 8C2000C0
0x904 0C000340
0x908 2021FFFF
0x90C 1420FFFF
...
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?

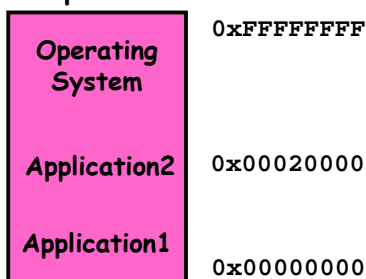
2/26/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 10.10

## Multiprogramming (First Version)

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



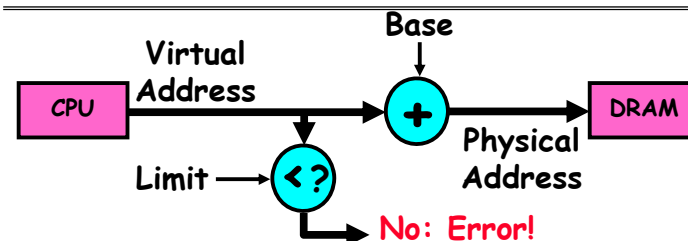
- Trick: Use Loader/Linker: Adjust addresses while program loaded into memory (loads, stores, jumps)
  - » Everything adjusted to memory location of program
  - » Translation done by a linker-loader
  - » Was pretty common in early days
- With this solution, no protection: bugs in any program can cause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

2/26/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 10.11

## Simple Segmentation: Base and Limit registers (CRAY-1)



- Could use base/limit for dynamic address translation (often called "segmentation"):
  - Alter address of every load/store by adding "base"
  - Generate error if address bigger than limit
- This gives program the illusion that it is running on its own dedicated machine, with memory starting at 0
  - Program gets continuous region of memory
  - Addresses within program do not have to be relocated when program placed in different region of DRAM
- User may have multiple segments available (e.g x86)
  - Loads and stores include segment ID in opcode: x86 Example: `mov [es:bx],ax.`

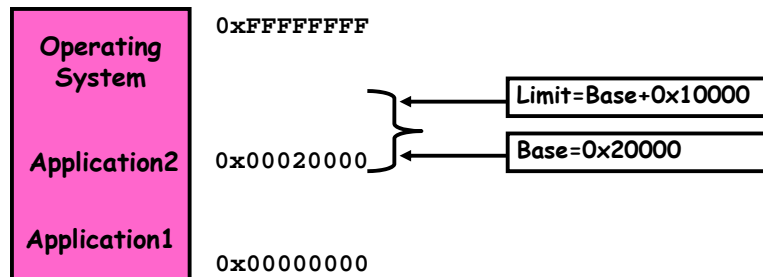
2/26/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 10.12

## Multiprogramming (Version with Protection)

- Can we protect programs from each other without translation?



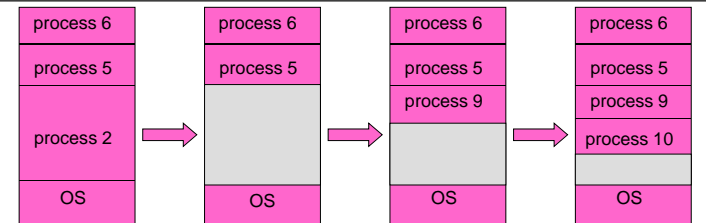
- Yes: use two special registers *Base* and *Limit* to prevent user from straying outside designated area
  - » If user tries to access an illegal address, cause an error
- During switch, kernel loads new base/limit from TCB
  - » User not allowed to change base/limit registers

2/26/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 10.13

## Issues with simple segmentation method



- 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 providing multiple segments per process
- Need enough physical memory for every process

2/26/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 10.14

## Administrivia

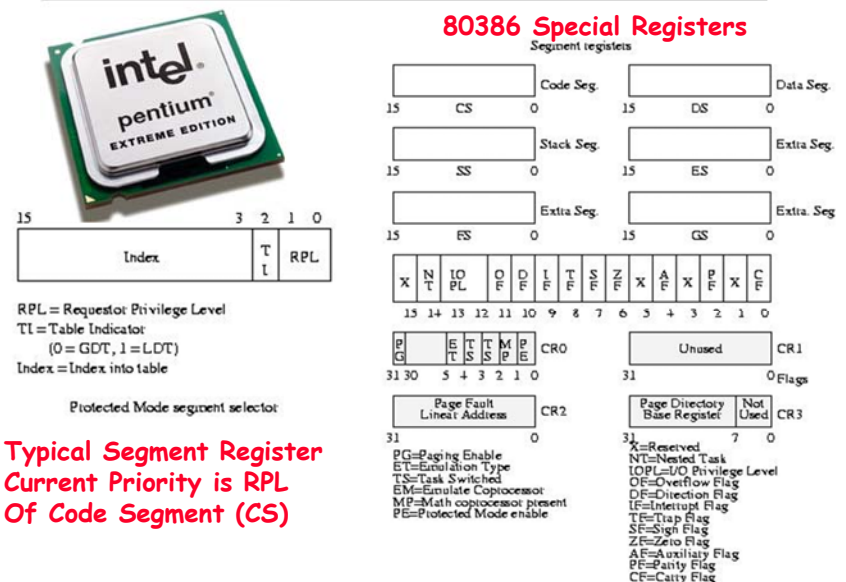
- Recall: Midterm I: 2½ weeks from today!
  - Wednesday 3/12
  - Intention is a 1.5 hour exam over 3 hours
  - No class on day of exam!
- Midterm Timing:
  - Probably 4:00-7:00PM in 320 Soda Hall
  - Does this work? (Still confirming room)
- Topics: everything up to the previous Monday
  - OS Structure, BDD, Process support, Synchronization, Memory Management, File systems
- Lab 1 Due!
  - Code checkin tomorrow at 9pm
  - Final project report due Friday @Noon
  - Final partner evaluations due Friday by midnight

2/26/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 10.15

## Intel x86 Special Registers



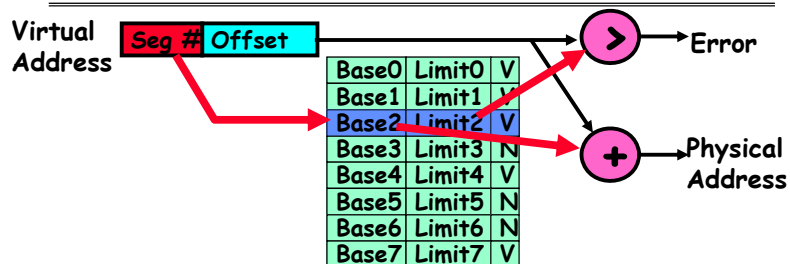
Typical Segment Register  
Current Priority is RPL  
Of Code Segment (CS)

2/26/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 10.16

## Implementation of Multi-Segment Model



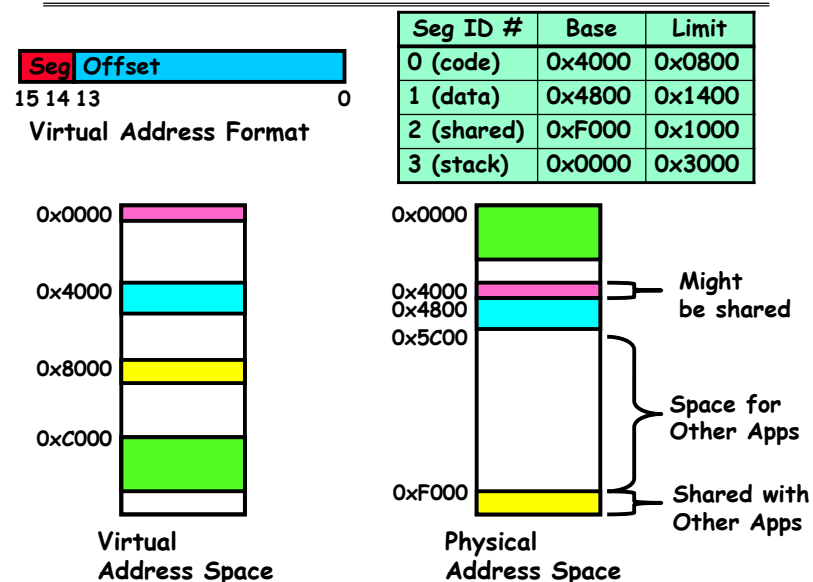
- Segment map resides in processor
  - Segment number mapped into base/limit pair
  - Base added to offset to generate physical address
  - Error check catches offset out of range
- As many chunks of physical memory as entries
  - Segment addressed by portion of virtual address
  - However, could be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

2/26/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 10.17

## Example: Four Segments (16 bit addresses)



2/26/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 10.18

## 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*

2/26/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 10.19

## 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")

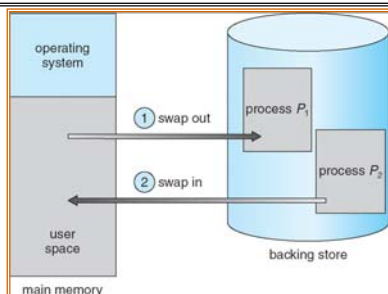
2/26/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 10.20



## 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

2/26/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 10.21

## Paging: Physical Memory in Fixed Size Chunks

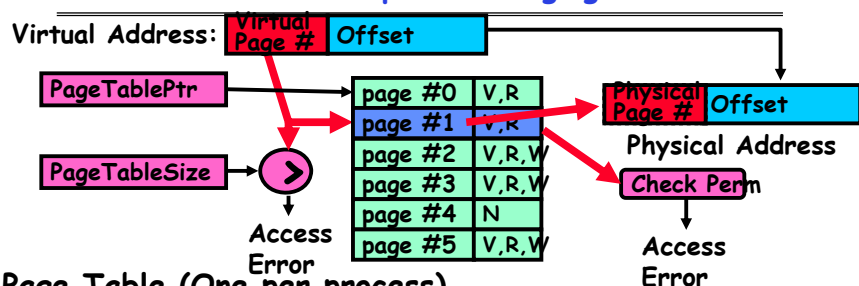
- Problems with segmentation?
  - Must fit variable-sized chunks into physical memory
  - May move processes multiple times to fit everything
  - Limited options for swapping to disk
- **Fragmentation**: wasted space
  - **External**: free gaps between allocated chunks
  - **Internal**: don't need all 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⇒allocated, 0⇒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

2/26/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 10.22

## How to Implement Paging?



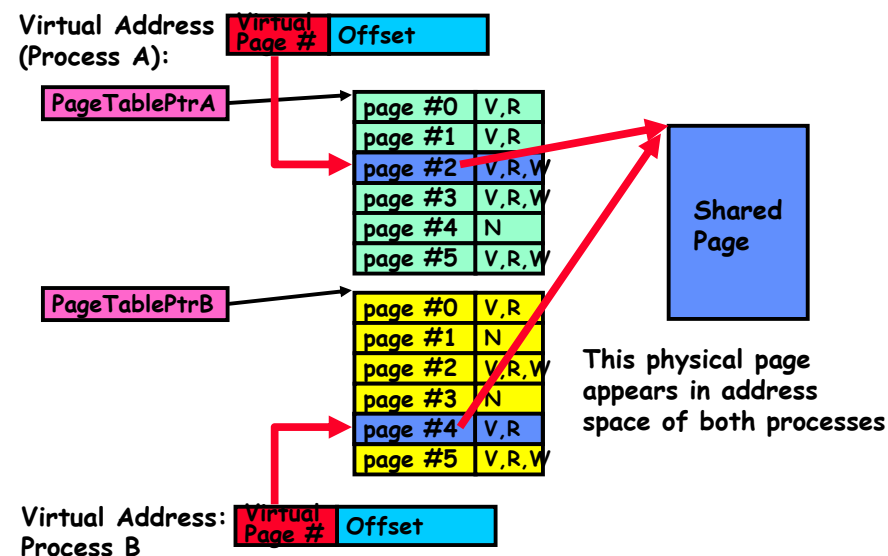
- Page Table (One per process)
  - Resides in physical memory
  - Contains physical page and permission for each virtual page
    - » Permissions include: Valid bits, Read, Write, etc
- Virtual address mapping
  - Offset from Virtual address copied to Physical Address
    - » Example: 10 bit offset ⇒ 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

2/26/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 10.23

## What about Sharing?



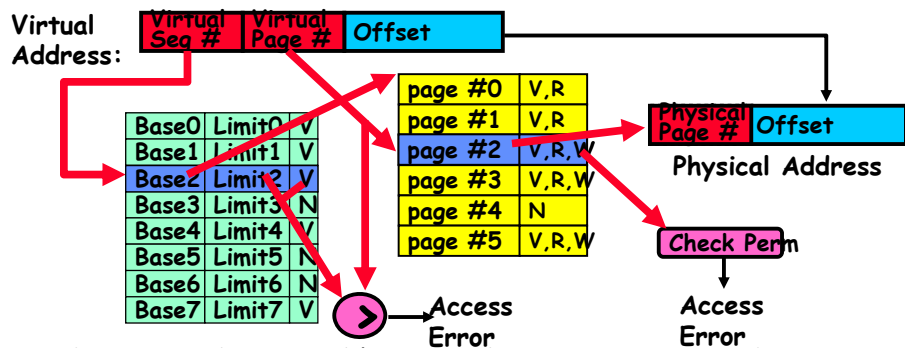
2/26/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 10.24

## Multi-level Translation: Segments + Pages

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



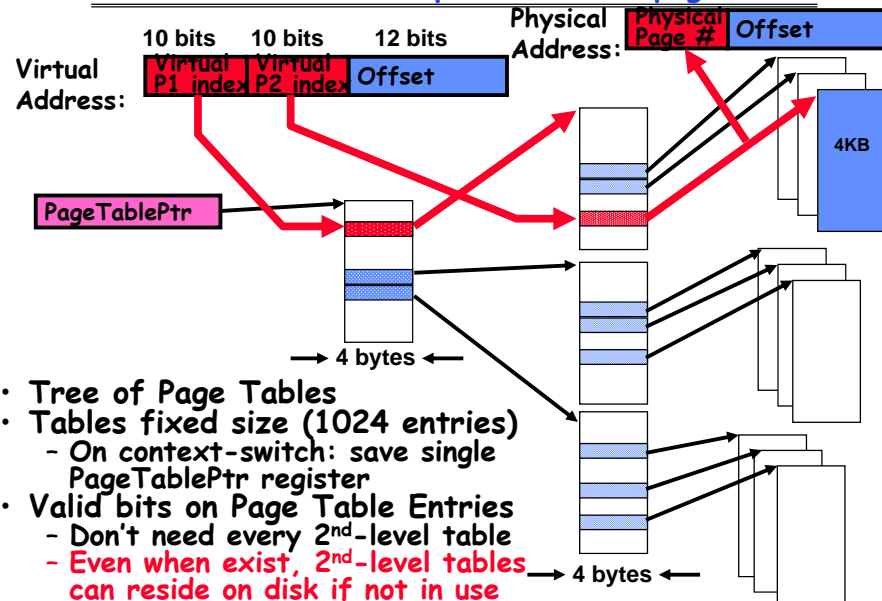
- What must be saved/restored on context switch?
  - Contents of top-level segment registers (for this example)
  - Pointer to top-level table (page table)

2/26/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 10.25

## 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 2<sup>nd</sup>-level table
  - Even when exist, 2<sup>nd</sup>-level tables can reside on disk if not in use

2/26/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 10.26

## Multi-level Translation Analysis

- Pros:
  - Only need to allocate as many page table entries as we need for application
    - » In other words, 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!

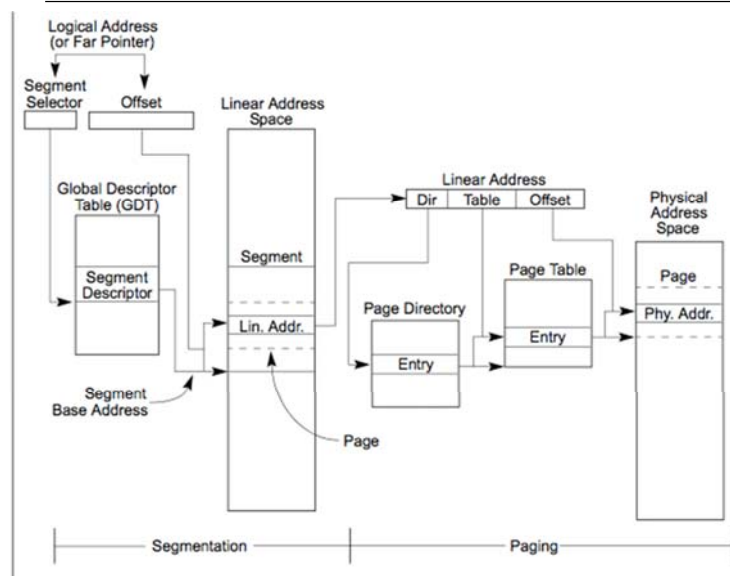
2/26/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 10.27

## Making it real:

### X86 Memory model with segmentation (16/32-bit)



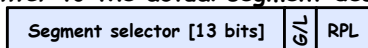
2/26/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 10.28

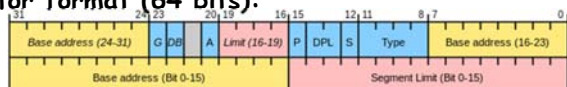
## X86 Segment Descriptors (32-bit Protected Mode)

- Segments are either implicit in the instruction (say for code segments) or actually part of the instruction
  - There are 6 registers: SS, CS, DS, ES, FS, GS
- What is in a segment register?
  - A *pointer* to the actual segment description:



G/L selects between GDT and LDT tables (global vs local descriptor tables)

- Two registers: GDTR and LDTR hold pointers to the global and local descriptor tables in memory
  - Includes length of table (for  $< 2^{13}$  entries)
- Descriptor format (64 bits):



- G: Granularity of segment (0: 16bit, 1: 4KiB unit)
- DB: Default operand size (0: 16bit, 1: 32bit)
- A: Freely available for use by software
- P: Segment present
- DPL: Descriptor Privilege Level
- S: System Segment (0: System, 1: code or data)
- Type: Code, Data, Segment

2/26/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 10.29

## How are segments used?

- One set of global segments (GDT) for everyone, different set of local segments (LDT) for every process
- In legacy applications (16-bit mode):
  - Segments provide protection for different components of user programs
  - Separate segments for chunks of code, data, stacks
  - Limited to 64K segments
- Modern use in 32-bit Mode:
  - Segments "flattened", i.e. every segment is 4GB in size
  - One exception: Use of GS (or FS) as a pointer to "Thread Local Storage"
    - » A thread can make accesses to TLS like this:
 

```
mov eax, gs(0x0)
```
- Modern use in 64-bit ("long") mode
  - Most segments (SS, CS, DS, ES) have zero base and no length limits
  - Only FS and GS retain their functionality (for use in TLS)

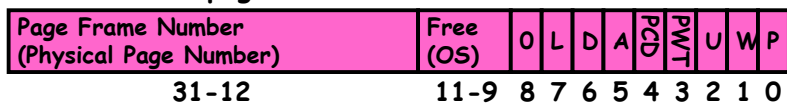
2/26/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 10.30

## What is in a Page Table Entry (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"



- 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

2/26/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 10.31

## Examples of how to exploit 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

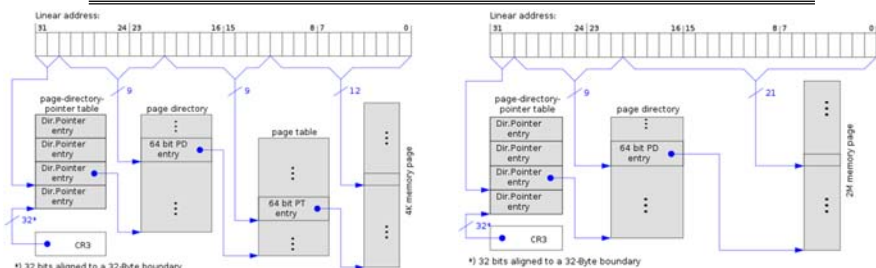
2/26/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 10.32



## Slightly More than 4GB RAM: PAE mode on x86



PAE with 4K pages

PAE with 2MB pages

- Physical Address Extension (PAE)
  - Poor-man's large memory extensions
  - More than 4GB physical memory
  - Every process still can have only 32-bit address space
- 3-Level page table
  - 64-bit PTE format
- How do processes use more than 4GB memory?
  - OS Support for mapping and unmapping physical memory into virtual address space
  - Application Windowing Extensions (AWE)

2/26/14

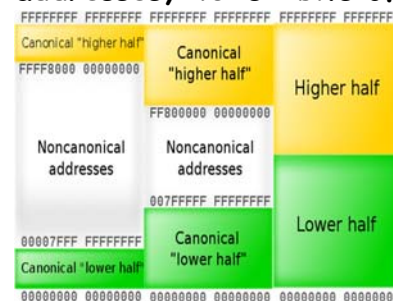
Kubiawicz CS194-24 @UCB Fall 2014

Lec 10.33

## What about 64-bit x86-64 ("long mode")?

- X86 long mode: 64 bit virtual addresses, 40-52 bits of physical memory

- Not all 64-bit virtual addresses translated
- Virtual Addresses must be "canonical": top n bits of must be equal
  - n here might be 48
  - Non-canonical addresses will cause a protection fault



- Using PAE scheme with 64-bit PTE can map 48-bits of virtual memory ( $9 \times 4 + 12 = 48$ )
- As mentioned earlier, segments other than FS/GS disabled in long mode

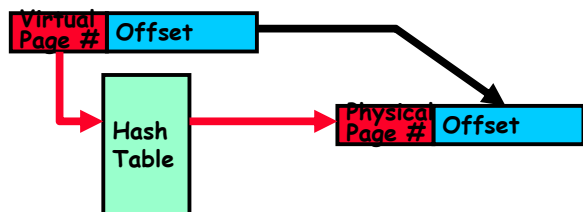
2/26/14

Kubiawicz CS194-24 @UCB Fall 2014

Lec 10.34

## 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



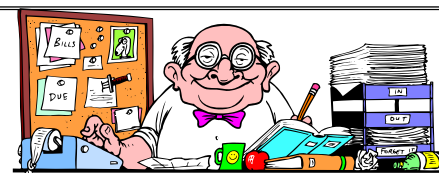
- Answer: use a hash table
  - Called an "Inverted Page Table"
  - Size is independent of virtual address space
  - Directly related to amount of physical memory
  - Very attractive option for 64-bit address spaces
- Cons: Complexity of managing hash changes
  - Often in hardware!

2/26/14

Kubiawicz CS194-24 @UCB Fall 2014

Lec 10.35

## Recall: Caching Concept



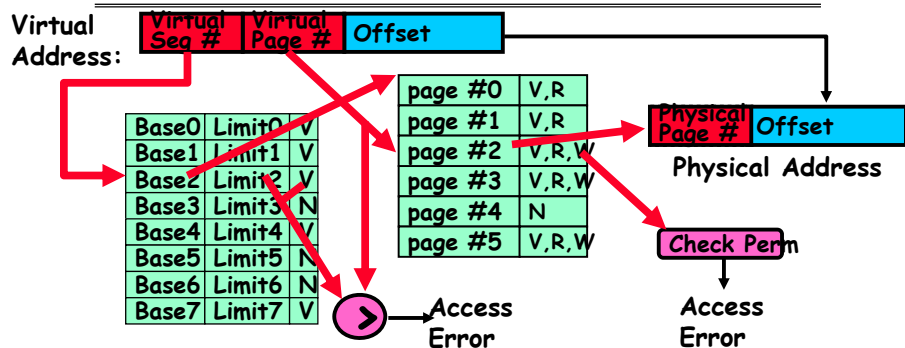
- Cache:** a repository for copies that can be accessed more quickly than the original
  - Make frequent case fast and infrequent case less dominant
- Caching underlies many of the techniques that are used today to make computers fast
  - Can cache: memory locations, address translations, pages, file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

2/26/14

Kubiawicz CS194-24 @UCB Fall 2014

Lec 10.36

## Why does caching matter for Virtual Memory?



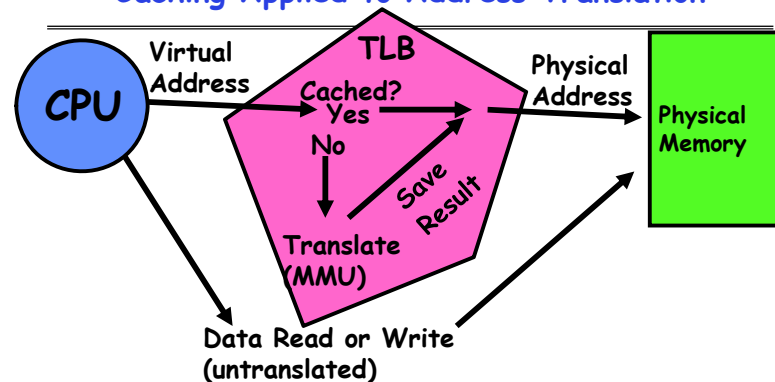
- Cannot afford to translate on every access
  - At least three DRAM accesses per actual DRAM access
  - Or: perhaps I/O if page table partially on disk!
- Even worse: What if we are using caching to make memory access faster than DRAM access???
- Solution? Cache translations!
  - Translation Cache: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

2/26/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 10.37

## 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

2/26/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 10.38

## 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

2/26/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 10.39

## What happens on a Context Switch?

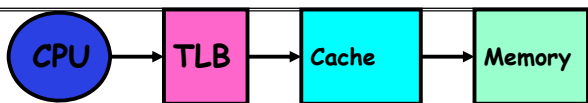
- Need to do something, since TLBs map virtual addresses to physical addresses
  - Address Space just changed, so TLB entries no longer valid!
- Options?
  - Invalidate TLB: simple but might be expensive
    - » What if switching frequently between processes?
  - Include ProcessID in TLB
    - » This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

2/26/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 10.40

## What TLB organization makes sense?



- 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

2/26/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 10.41

## 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

2/26/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 10.42

## Example: R3000 pipeline includes TLB "stages"

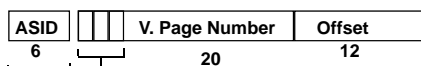
### MIPS R3000 Pipeline

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



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

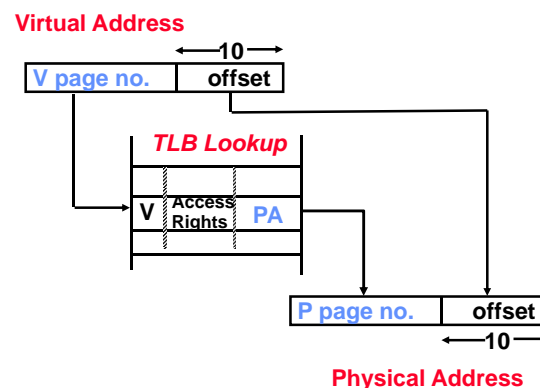
2/26/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 10.43

## 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

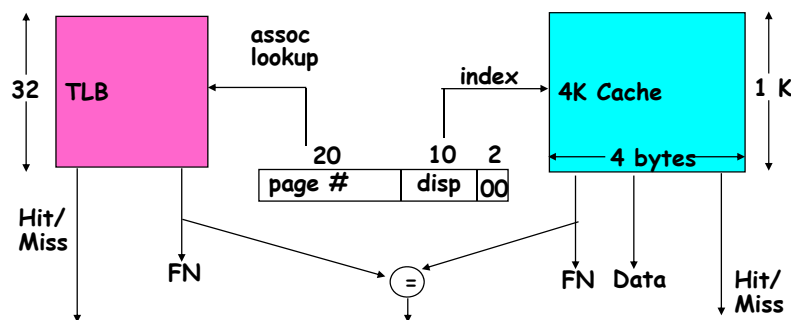
2/26/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 10.44

## Overlapping TLB & Cache Access

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



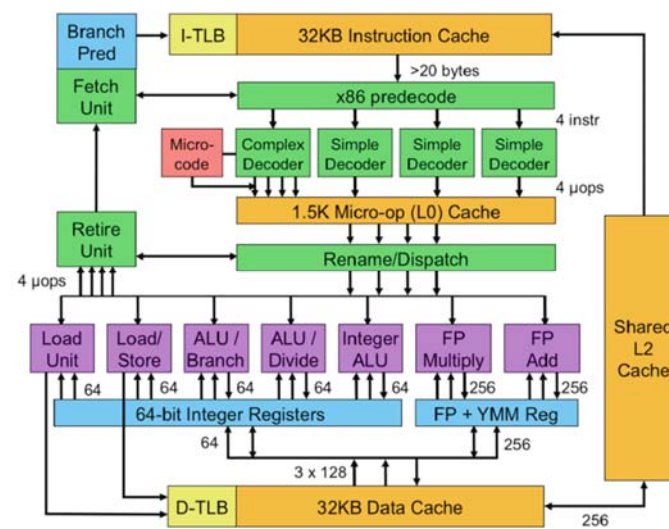
- What if cache size is increased to 8KB?
  - Overlap not complete
  - Need to do something else. See CS152/252
- Another option: Virtual Caches
  - Tags in cache are virtual addresses
  - Translation only happens on cache misses

2/26/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 10.45

## Modern Pipelines: SandyBridge Pipeline

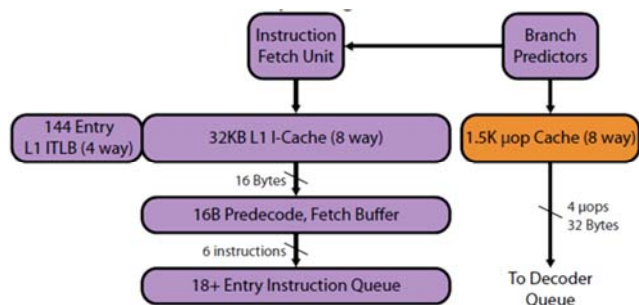


2/26/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 10.46

## Fetch and Branch Prediction



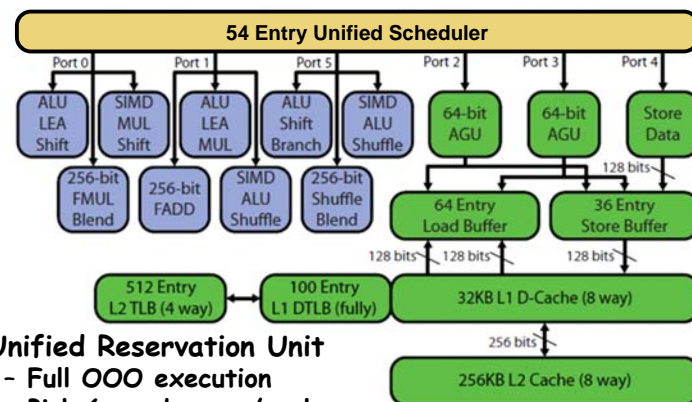
- x86 instructions turned into micro-ops
  - Cached translations are reused
- Branch prediction
  - Not entirely clear, but seems to have some combination of bi-mode, multi-level BTB, stack-based prediction for CALL/RETURN
- Predecoder
  - Finds instruction boundaries
  - Passes at least 6 instructions onto decoding infrastructure

2/26/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 10.47

## Out-of-Order execution: Data TLB (DTLB)



- Unified Reservation Unit
  - Full OOO execution
  - Pick 6 ready μops/cycle
  - Can have two loads or stores/cycle
    - » 2 address generation units (AGUs) + store data
  - Simultaneous 256-bit Multiply and Add
  - Can have 3 regular integer ops/cycle

2/26/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 10.48

## 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
- **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

2/26/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 10.49

## Summary (2/2)

---

- **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

2/26/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 10.50