CMSC 313 COMPUTER ORGANIZATION & ASSEMBLY LANGUAGE PROGRAMMING

Lecture 27, Fall 2014

CACHING

- Why: bridge speed difference between CPU and RAM
- Modern RAM allows blocks of memory to be read quickly
- Principle of locality: temporal and spatial

During each memory access :

- CPU checks if memory location is already in cache
- Found = cache hit:
 - read from or write to cache
- Not Found = cache miss:
 - Fetch entire memory block of location into cache

CACHE MAPPING SCHEMES

Direct Mapping:

- Each memory block mapped to 1 cache block
- Many memory blocks for each cache block
- Use *tag* to check if block in cache is the one needed

Fully Associative Mapping:

- Each memory block can be placed in any cache block
- Associative memory finds cache block with *tag*

Set Associative Mapping:

• Hybrid of direct mapping and fully associative mapping

DIRECT MAPPING

- The purpose of cache memory is to speed up accesses by storing recently used data closer to the CPU, instead of storing it in main memory.
- Although cache is much smaller than main memory, its access time is a fraction of that of main memory.
- Unlike main memory, which is accessed by address, cache is typically accessed by content; hence, it is often called *content addressable memory*.
- Because of this, a single large cache memory isn't always desirable-- it takes longer to search.

- The "content" that is addressed in content addressable cache memory is a subset of the bits of a main memory address called a *field*.
 - Many blocks of main memory map to a single block of cache. A *tag* field in the cache block distinguishes one cached memory block from another.
 - A *valid bit* indicates whether the cache block is being used.
 - An offset field points to the desired data in the block.



- The simplest cache mapping scheme is *direct mapped cache*.
- In a direct mapped cache consisting of N blocks of cache, block X of main memory maps to cache block Y = X mod N.
- Thus, if we have 10 blocks of cache, block 7 of cache may hold blocks 7, 17, 27, 37, . . . of main memory.

The next slide illustrates this mapping.

 With direct mapped cache consisting of *N* blocks of cache, block *X* of main memory maps to cache block *Y* = *X* mod *N*.



- EXAMPLE 6.1 Consider a word-addressable main memory consisting of four blocks, and a cache with two blocks, where each block is 4 words.
- This means Block 0 and 2 of main memory map to Block 0 of cache, and Blocks 1 and 3 of main memory map to Block 1 of cache.
- Using the tag, block, and offset fields, we can see how main memory maps to cache as follows.

- <u>EXAMPLE 6.1</u> Consider a word-addressable main memory consisting of four blocks, and a cache with two blocks, where each block is 4 words.
 - First, we need to determine the address format for mapping.
 Each block is 4 words, so the offset field must contain 2 bits;
 there are 2 blocks in cache, so the block field must contain 1 bit;
 this leaves 1 bit for the tag (as a main memory address has 4 bits because there are a total of 2⁴=16 words).



• EXAMPLE 6.1 Cont'd

- Suppose we need to access main memory address 3₁₆ (0011 in binary). If we partition 0011 using the address format from Figure a, we get Figure b.
- Thus, the main memory address
 0011 maps to cache block 0.
- Figure c shows this mapping, along with the tag that is also stored with the data.

The next slide illustrates another mapping.



© 2012 Jones & Bartiett Learning, LLC www.jblearning.com



- EXAMPLE 6.2 Assume a byte-addressable memory • consists of 2¹⁴ bytes, cache has 16 blocks, and each - The number of memory blocks are: $\frac{2^{14}}{2^3} = 2^{11}$ block has 8 bytes.

 - Each main memory address requires 14 bits. Of this 14-bit address field, the rightmost 3 bits reflect the offset field
 - We need 4 bits to select a specific block in cache, so the block field consists of the middle 4 bits.
 - The remaining 7 bits make up the tag field.



- In summary, direct mapped cache maps main memory blocks in a modular fashion to cache blocks. The mapping depends on:
- The number of bits in the main memory address (how many addresses exist in main memory)
- The number of blocks are in cache (which determines the size of the block field)
- How many addresses (either bytes or words) are in a block (which determines the size of the offset field)

FULLY ASSOCIATIVE MAPPING

- Suppose instead of placing memory blocks in specific cache locations based on memory address, we could allow a block to go anywhere in cache.
- In this way, cache would have to fill up before any blocks are evicted.
- This is how *fully associative* cache works.
- A memory address is partitioned into only two fields: the tag and the word.

 Suppose, as before, we have 14-bit memory addresses and a cache with 16 blocks, each block of size 8. The field format of a memory reference is: 11 bits 3 bits



- When the cache is searched, all tags are searched in parallel to retrieve the data quickly.
- This requires special, costly hardware.

SET ASSOCIATIVE MAPPING

- Set associative cache combines the ideas of direct mapped cache and fully associative cache.
- An *N*-way set associative cache mapping is like direct mapped cache in that a memory reference maps to a particular location in cache.
- Unlike direct mapped cache, a memory reference maps to a set of several cache blocks, similar to the way in which fully associative cache works.
- Instead of mapping anywhere in the entire cache, a memory reference can map only to the subset of cache slots.

- The number of cache blocks per set in set associative cache varies according to overall system design.
 - For example, a 2-way set associative cache can be conceptualized as shown in the schematic below.
 - Each set contains two different memory blocks.





- In set associative cache mapping, a memory reference is divided into three fields: tag, set, and offset.
- As with direct-mapped cache, the offset field chooses the word within the cache block, and the tag field uniquely identifies the memory address.
- The set field determines the set to which the memory block maps.

- <u>EXAMPLE 6.5</u> Suppose we are using 2-way set associative mapping with a word-addressable main memory of 2¹⁴ words and a cache with 16 blocks, where each block contains 8 words.
 - Cache has a total of 16 blocks, and each set has 2 blocks, then there are 8 sets in cache.
 - Thus, the set field is 3 bits, the offset field is 3 bits, and the tag field is 8 bits.



CACHING POLICIES

- Cache replacement policy
 - For fully associative and set associative mapping
 - Which cache block gets kicked out?
 - Some schemes: first-in first-out, least recently used, ...
- Cache write policy
 - Write through: always write to main memory
 - Write back: write to main memory when replaced

CACHE PERFORMANCE

- The performance of hierarchical memory is measured by its *effective access time* (EAT).
- EAT is a weighted average that takes into account the hit ratio and relative access times of successive levels of memory.
- The EAT for a two-level memory is given by: EAT = $H \times Access_{C} + (1-H) \times Access_{MM}$.

where H is the cache hit rate and $Access_C$ and $Access_{MM}$ are the access times for cache and main memory, respectively.

- For example, consider a system with a main memory access time of 200ns supported by a cache having a 10ns access time and a hit rate of 99%.
- Suppose access to cache and main memory occurs concurrently. (The accesses overlap.)
- The EAT is:

0.99(10ns) + 0.01(200ns) = 9.9ns + 2ns = 11ns.

- For example, consider a system with a main memory access time of 200ns supported by a cache having a 10ns access time and a hit rate of 99%.
- If the accesses do not overlap, the EAT is:

0.99(10ns) + 0.01(10ns + 200ns)

= 9.9ns + 2.01ns = 12ns.

• This equation for determining the effective access time can be extended to any number of memory levels, as we will see in later sections.

VIRTUAL MEMORY

MEMORY PROBLEMS

Not enough memory

- Many processes ran simultaneously
- Large applications, but most code is unused (MS Word)

Fragmentation

- Processes need contiguous blocks of memory
- Total amount of free memory is sufficient, but largest block of contiguous memory is too small

Unprotected memory

- Many processes ran simultaneously
- "Bad" processes can overwrite other processes' memory

- Cache memory enhances performance by providing faster memory access speed.
- Virtual memory enhances performance by providing greater memory capacity, without the expense of adding main memory.
- Instead, a portion of a disk drive serves as an extension of main memory.
- If a system uses paging, virtual memory is partitioned into individually managed *pages*, that are written to *(or paged to)* disk when they are not immediately needed.

- Main memory and virtual memory are divided into equal sized pages.
- The entire address space required by a process need not be in memory at once. Some parts can be on disk, while others are in main memory.
- Further, the pages allocated to a process do not need to be stored contiguously-- either on disk or in memory.
- In this way, only the needed pages are in memory at any time, the unnecessary pages are in slower disk storage.

- Information concerning the location of each page, whether on disk or in memory, is maintained in a data structure called a *page table* (shown below).
- There is one page table for each active process.



- When a process generates a virtual address, the operating system translates it into a physical memory address.
- To accomplish this, the virtual address is divided into two fields: A *page* field, and an *offset* field.
- The page field determines the page location of the address, and the offset indicates the location of the address within the page.
- The logical page number is translated into a physical page frame through a lookup in the page table.

- If the valid bit is zero in the page table entry for the logical address, this means that the page is not in memory and must be fetched from disk.
 - This is a page fault.
 - If necessary, a page is evicted from memory and is replaced by the page retrieved from disk, and the valid bit is set to 1.
- If the valid bit is 1, the virtual page number is replaced by the physical frame number.
- The data is then accessed by adding the offset to the physical frame number.

 As an example, suppose a system has a virtual address space of 8K and a physical address space of 4K, and the system uses byte addressing.

- We have $2^{13}/2^{10} = 2^3$ virtual pages.

- A virtual address has 13 bits (8K = 2¹³) with 3 bits for the page field and 10 for the offset, because the page size is 1024.
- A physical memory address requires 12 bits, the first two bits for the page frame and the trailing 10 bits the offset.



- Suppose we have the page table shown below.
- What happens when CPU generates address 5459_{10} = 10101010011₂ = 1553 ₁₆?



				Ade	dresses		
Page		Ba	se	10		Bas	se 16
0	:	0		1023	0		3FF
1	:	1024		2047	400		7ff
2	:	2048		3071	800		BFF
3	:	3072		4095	C00		FFF
4	:	4096		5119	1000		13FF
5	:	5120		6143	1400		17FF
6	:	6144		7167	1800		1BFF
7	:	7168		8191	1000		1 F F F

• What happens when CPU generates address 5459_{10} = 10101010011₂ = 1553 ₁₆?



The high-order 3 bits of the virtual address, 101 (5_{10}) , provide the page number in the page table.

• The address 1010101010011_2 is converted to physical address $010101010011_2 = 1363_{16}$ because the page field 101 is replaced by frame number 01 through a lookup in the page table.



 What happens when the CPU generates address 100000000100₂?



- We said earlier that effective access time (EAT) takes all levels of memory into consideration.
- Thus, virtual memory is also a factor in the calculation, and we also have to consider page table access time.
- Suppose a main memory access takes 200ns, the page fault rate is 1%, and it takes 10ms to load a page from disk. We have:

EAT = 0.99(200ns + 200ns) + 0.01(10ms)

= 0.99 x 400ns + 0.01 x 10,000ns

= 396ns + 100ns = 496ns

- Even if we had no page faults, the EAT would be 400ns because memory is always read twice: First to access the page table, and second to load the page from memory.
- Because page tables are read constantly, it makes sense to keep them in a special cache called a *translation look-aside buffer* (TLB).
- TLBs are a special associative cache that stores the mapping of virtual pages to physical pages.

The next slide shows address lookup steps when a TLB is involved.

TLB lookup process

1. Extract the page number from the virtual address.

2. Extract the offset from the virtual address.

CPU

3. Search for the virtual page number in the TLB.

4. If the (virtual page #, page frame #) pair is found in the TLB, add the offset to the physical frame number and access the memory location.

5. If there is a TLB miss, go to the page table to get the necessary frame number.

If the page is in memory, use the corresponding frame number and add the offset to yield the physical address.

6. If the page is not in main memory, generate a page fault and restart the access when the page fault is complete.



56



VIRTUAL MEMORY IN LINUX

Linux Virtual Memory Space

- Linux reserves 1 Gig memory in the virtual address space
- The size of the Linux kernel significantly affects its performance (swapping is expensive)
- Linux kernel can be customized by including only relevant modules
- Designating kernel space facilitates protection of
- The portion of disk used for paging is called the swap space



Virtual Addressing



Physical address

□ Page faults are costly and take millions of cycles to process (disks are slow)

□ 80386 Page attributes:

- RW: read and write permission
- ➡ US: User mode or kernel mode only access
- ➡ PP: present bit to indicate where the page is

31		12 11	2	1	0
	Address of Page		U S	W R	P P

Page Table

Hardware supported	Page table register	
	Virtual address	
Page table:	31 30 29 28 27 15 14 13 12 11 10	98 3210
★ Resides in main memory	Virtual page number	Page offset
★ One entry per virtual page	20 Valid Physical page number	12
★ No tag is requires since it covers all virtual pages		
★ Point directly to physical pag	e di	
★ Table can be very large	Page table	- L -
★ Operating sys. may maintain one page table per process		
★ A dirty bit is used to track modified pages for copy back	18 If 0 then page is not present in memory	
Indicates whether the virtual page is in main memory or not	29 28 27 ····· 15 14 13 12 11 1 Physical page number	0 9 8 · · · · · 3 2 1 0 Page offset
	Physical address	



Linux 2-Level Page Table



The CR3 register is designated for pointing to the first level page tableThe CR3 is part of the task state that needs to be saved at preemption

31	22	21	12 11		0
	Index into Page Table Table	Index into Page Table		Index into Page	

3.7.1. Linear Address Translation (4-KByte Pages)

Figure 3-12 shows the page directory and page-table hierarchy when mapping linear addresses to 4-KByte pages. The entries in the page directory point to page tables, and the entries in a page table point to pages in physical memory. This paging method can be used to address up to 2^{20} pages, which spans a linear address space of 2^{32} bytes (4 GBytes).



Figure 3-12. Linear Address Translation (4-KByte Pages)

	Page-Directory Entry (4-K	Byt	e Page	T	ab	le)						
31		12	11 9	8	7	6	5	4	3	2	1	0
	Page-Table Base Address		Avail	G	P S	0	A	P C D	P W T	U / S	R / W	Ρ
	Available for system programmer's use Global page (Ignored) — Page size (0 indicates 4 KBytes) — Reserved (set to 0) — Accessed — Cache disabled — Write-through — User/Supervisor — Read/Write — Present —											
	Page-Table Entry (4-I	K By	yte Pag	e)	_	_	_		_	_		_
31		12	2 11 9	8	7	6	5	4	3	2	1	0
	Page Base Address		Avail	G	P A T	D	A	P C D	P W T	U / S	R / W	P
	Available for system programmer's use Global Page Page Table Attribute Index Dirty Accessed Cache Disabled Write-Through User/Supervisor											

Figure 3-14. Format of Page-Directory and Page-Table Entries for 4-KByte Pages and 32-Bit Physical Addresses

Virtual Memory: Problems Solved

Not enough physical memory

- Uses disk space to simulate extra memory
- Pages not being used can be swapped out (how and when you'll learn in CMSC 421 Operating Systems)
- Thrashing: pages constantly written to and retrieved from disk (time to buy more RAM)

• Fragmentation

 Contiguous blocks of virtual memory do not have to map to contiguous sections of real memory

Memory protection

- Seach process has its own page table
- Shared pages are read-only
- User processes cannot alter the page table (must be supervisor)

Memory Protection

- Prevents one process from reading from or writing to memory used by another process
- Privacy in a multiple user environments
- Operating system stability
 - Prevents user processes (applications) from altering memory used by the operating system
 - One application crashing does not cause the entire OS to crash

Virtual Memory: too slow?

• Address translation is done in hardware

In the middle of the fetch execute cycle for:

MOV EAX, [buffer]

the physical address of buffer is computed in hardware.

- Recently computed page locations are cached in the translation lookaside buffer (TLB)
- Page faults are very expensive (millions of cycles)
- Operating systems for personal computers have only recently added memory protection

NEXT TIME

Review