William Stallings Computer Organization and Architecture

Chapter 4 Internal Memory

Characteristics

₭ Location
CPU, Internal, External

₭ Unit of transfer
☑Word on bus, block, cluster

₭ Access method☑ Direct, Random, Associative, Sequential ₭ Performance
Access, Cycle, Transfer time

- Corganization
 Physical arrangement of bits into words

Access Methods (1)

#Sequential

 Start at the beginning and read through in order
 Access time depends on location of data and previous location

⊡e.g. tape

#Direct

Individual blocks have unique address

☑Access is by jumping to vicinity plus sequential search

☑Access time depends on location and previous location

⊡e.g. disk



Memory Hierarchy

₩Registers ⊠In CPU

- H Internal or Main memory
 ⊡May include one or more levels of cache
 ⊡"RAM"

Performance

#Access time

☐Time between presenting the address and getting the valid data

#Memory Cycle time

☐Time may be required for the memory to "recover" before next access

⊡Cycle time is access + recovery

#Transfer Rate

□ Rate at which data can be moved

Physical Characteristics

₩Decay ₩Volatility

₩Erasable

₩Power consumption

The Bottom Line

₭ How much?
△Capacity
₭ How fast?
△Time is money

How expensive?

ℋ Tradeoffs among all of these ▷ E.g. Faster = More expensive, More = Less cost (per bit) but slower

☐ Solution : Memory Hierarchy

Hierarchy List

Registers

₩L1 Cache

₩ L2 Cache

Hain memory

₭ Disk cache₭ Disk

₩ Optical

. Ж Таре

₭ As one goes down the hierarchy

⊡ Decreasing cost per bit

☐ Increasing capacity

 \boxdot Increasing access time

Decreasing frequency of access of the memory by the processor – locality of reference

So you want fast?

#It is possible to build a computer which uses only static RAM (see later)

₭This would be very fast

HThis would need no cache

⊡How can you cache cache?

HThis would cost a very large amount



₭ Temporal Locality

☑ Programs tend to reference the same memory locations at a future point in time

Due to loops and iteration, programs spending a lot of time in one section of code

ℜ Spatial Locality

□ Programs tend to reference memory locations that are near other recently-referenced memory locations

☐ Due to the way contiguous memory is referenced, e.g. an array or the instructions that make up a program

Eccality of reference does not always hold, but it usually holds
 Holds



Semiconductor Memory

₩RAM

☐Misnamed as all semiconductor memory is random access

⊡Read/Write

⊠Volatile

⊡Temporary storage

⊡Two main types: **Static** or **Dynamic**

Dynamic RAM

Bits stored as charge in capacitors

₭ Charges leak

X Need refreshing even when powered

Simpler construction

¥ Smaller per bit

₭ Less expensive

Need refresh circuits (every few milliseconds)

Slower

ℜ Main memory

Static RAM

₭ Bits stored as on/off switches via flip-flops

- ₭ No charges to leak
- **#** No refreshing needed when powered
- ₭ More complex construction
- **¥** Larger per bit
- ₿ More expensive
- ₭ Does not need refresh circuits
- ₩ Faster
- ₩ Cache

Read Only Memory (ROM)

#Permanent storage
#Microprogramming
#Library subroutines
#Systems programs (BIOS)
#Function tables









Refreshing

#Refresh circuit included on chip
#Disable chip
#Count through rows
#Read & Write back
#Takes time
#Slows down apparent performance























Mapping Function

₩e'll use the following configuration example
Cache of 64KByte
Cache line / Block size is 4 bytes
⊠i.e. cache is 16,385 (2¹⁴) lines of 4 bytes
Main memory of 16MBytes
24 bit address
(2²⁴=16M)
16Mbytes / 4bytes-per-block → 4 MB of Memory Blocks
Somehow we have to map the 4Mb of blocks in memory onto the 16K of lines in the cache. Multiple blocks will have to map to the same line in the cache!







Direct Mapping Address Structure

#Address is in two parts

- □ Least Significant w bits identify unique word within a cache line
- ☐Most Significant s bits specify one memory block
- □The MSBs are split into a cache line field r and a tag of s-r (most significant)

Direct Mapping Address Structure

V D	Tag s-r	Line or Slot r	Word w				
1 1	8	14	2				
# Given a 24 bit address (to access 16Mb)							
# 2 bit word identifier (4 byte block)							
# 22 bit block identifier							
⊠ 8 bit tag (=22-14)							
[□]	☐ 14 bit slot or line						
% No two blocks in the same line have the same Tag field							
# Check contents of cache by finding line and checking Tag							
# Also need a Valid bit and a Dirty bit							
Valid – Indicates if the slot holds a block belonging to the program being executed							
Dirty – Indicates if a block has been modified while in the cache. Will need to be written back to memory before slot is reused for anot her block							











Associative Mapping Address Structure Word Tag 22 bit 2 bit # 22 bit tag stored with each 32 bit block of data **#** Compare tag field with tag entry in cache to check for hit # Least significant 2 bits of address identify which 8 bit word is required from 32 bit data block ₿e.g. Address: FFFFFC = 1111 1111 1111 1111 1111 1100 ⊠Tag: Left 22 bits, truncate on left: • 11 1111 1111 1111 1111 • 3FFFFF Address: 16339C = 0001 0110 0011 0011 1001 1100 ⊠Tag: Left 22 bits, truncate on left: • 00 0101 1000 1100 1110 0111 • 058CE7









Set Associative Mapping Address Structure





K-Way Set Associative

Content of the second s

#Further increases in the size of the set has little effect other than increased cost of the hardware!

Replacement Algorithms (1) Direct mapping

%No choice**%**Each block only maps to one line**%**Replace that line

Replacement Algorithms (2) Associative & Set Associative

Algorithm must be implemented in hardware (speed)
Least Recently used (LRU)

△ e.g. in 2 way set associative, which of the 2 block is LRU?
△ For each slot, have an extra bit, USE. Set to 1 when accessed, set all others to 0.

△ For more than 2-way set associative, need a time stamp for each slot - expensive

First in first out (FIFO)

△ Replace block that has been in cache longest
△ Easy to implement as a circular buffer

Least frequently used

△ Replace block which has had fewest hits
△ Need a counter to sum number of hits

Random

△ Almost as good as LFU and simple to implement



Write through

- **#** Simplest technique to handle the cache coherency problem All writes go to main memory as well as cache.
- # Multiple CPUs must monitor main memory traffic (snooping) to keep local cache local to its CPU up to date in case another CPU also has a copy of a shared memory location in its cache
- ₭ Simple but Lots of traffic
- **#** Slows down writes
- **#** Other solutions: noncachable memory, hardware to maintain coherency



Updates initially made in cache only

- ₭ Dirty bit for cache slot is cleared when update occurs
- # If block is to be replaced, write to main memory only if dirty bit is set
- **#** Other caches can get out of sync
- H If I/O must access invalidated main memory, one solution is for I/O to go through cache
 ⊡Complex circuitry
- **#** Only ~15% of memory references are writes





Cache Performance Example							
# Sample program executes from memory location 48-95 once. Then it executes from 15-31 in a loop ten times before exiting.							
	Event	Location	Time	Comment			
	1 miss 15 hits 1 miss 15 hits 1 miss 15 hits	48 49-63 64 65-79 80 81 95	2500ns 80ns×15=1200ns 2500ns 80ns×15=1200ns 2500ns 80ns×15=1200ns	Memory block 3 to cache slot 3 Memory block 4 to cache slot 0 Memory block 5 to cache slot 1			
	1 miss 1 miss 1 miss 15 hits 9 hits 144 hits	15 16 17-31 15 16-31	80ns×13=1200ns 2500ns 2500ns 80ns×15=1200ns 80ns×9=720ns 80ns×144=12,240ns	Memory block 0 to cache slot 0 Memory block 1 to cache slot 1 Last nine iterations of loop Last nine iterations of loop			
	Total hits = 213 Total misses = 5						

