

CS152
Computer Architecture and Engineering
Lecture 18

Memory and Caches

April 7, 1999

John Kubiatowicz (<http://cs.berkeley.edu/~kubitron>)

lecture slides: <http://www-inst.eecs.berkeley.edu/~cs152/>

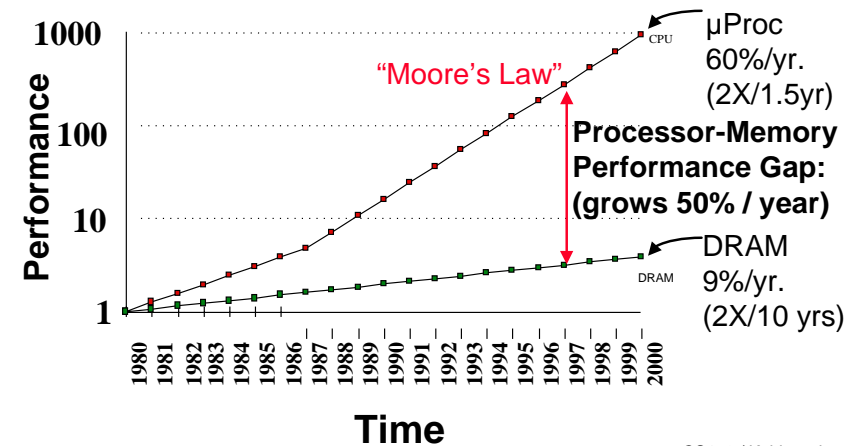
4/7/99

©UCB Spring 1999

CS152 / Kubiatowicz
Lec18.1

Recap: Who Cares About the Memory Hierarchy?

Processor-DRAM Memory Gap (latency)



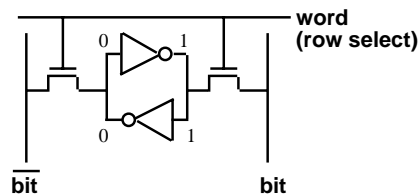
4/7/99

©UCB Spring 1999

CS152 / Kubiatowicz
Lec18.2

Recap: Static RAM Cell

6-Transistor SRAM Cell

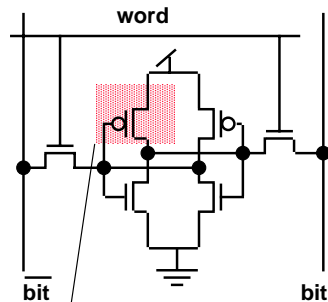


◦ Write:

1. Drive bit lines (bit=1, bit=0)
- 2.. Select row

◦ Read:

1. Precharge bit and bit to Vdd
- 2.. Select row
3. Cell pulls one line low
4. Sense amp on column detects difference between bit and bit



replaced with pullup
to save area

4/7/99

©UCB Spring 1999

CS152 / Kubiatowicz
Lec18.3

Recap: 1-Transistor Memory Cell (DRAM)

◦ Write:

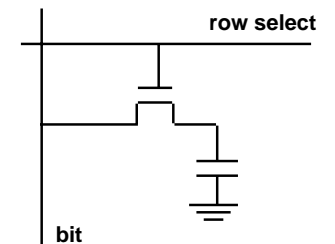
1. Drive bit line
- 2.. Select row

◦ Read:

1. Precharge bit line to Vdd
- 2.. Select row
3. Cell and bit line share charges
 - Very small voltage changes on the bit line
4. Sense (fancy sense amp)
 - Can detect changes of ~1 million electrons
5. Write: restore the value

◦ Refresh

1. Just do a dummy read to every cell.



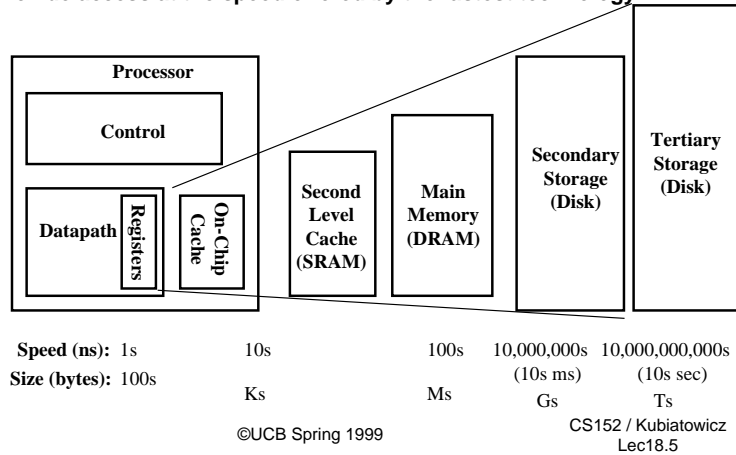
4/7/99

©UCB Spring 1999

CS152 / Kubiatowicz
Lec18.4

Recap: Memory Hierarchy of a Modern Computer System

- **By taking advantage of the principle of locality:**
 - Present the user with as much memory as is available in the cheapest technology.
 - Provide access at the speed offered by the fastest technology.



4/7/99

Recap: Memory Systems

- **Two Different Types of Locality:**
 - **Temporal Locality (Locality in Time):** If an item is referenced, it will tend to be referenced again soon.
 - **Spatial Locality (Locality in Space):** If an item is referenced, items whose addresses are close by tend to be referenced soon.
- **By taking advantage of the principle of locality:**
 - Present the user with as much memory as is available in the cheapest technology.
 - Provide access at the speed offered by the fastest technology.
- **DRAM is slow but cheap and dense:**
 - Good choice for presenting the user with a BIG memory system
- **SRAM is fast but expensive and not very dense:**
 - Good choice for providing the user FAST access time.

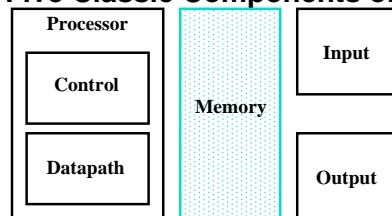
4/7/99

©UCB Spring 1999

CS152 / Kubiawicz
Lec18.6

The Big Picture: Where are We Now?

- **The Five Classic Components of a Computer**



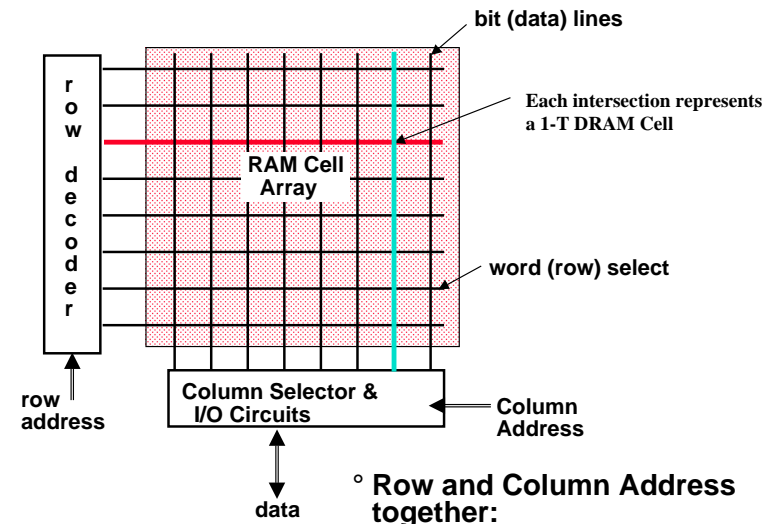
- **Today's Topics:**
 - Recap last lecture
 - Continue discussion of DRAM
 - Cache Review
 - Advanced Cache
 - Virtual Memory
 - Protection
 - TLB

4/7/99

©UCB Spring 1999

CS152 / Kubiawicz
Lec18.7

Classical DRAM Organization (square)



- **Row and Column Address together:**

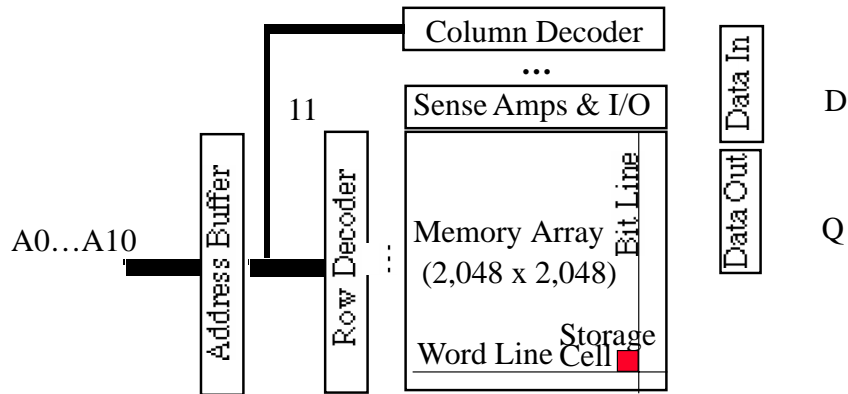
- Select 1 bit at a time

4/7/99

©UCB Spring 1999

CS152 / Kubiawicz
Lec18.8

DRAM logical organization (4 Mbit)



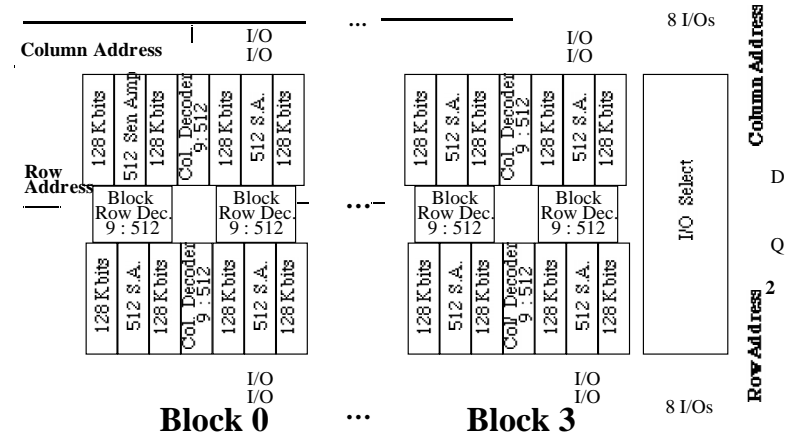
- Square root of bits per RAS/CAS

4/7/99

©UCB Spring 1999

CS152 / Kubiawicz
Lec18.9

DRAM physical organization (4 Mbit)

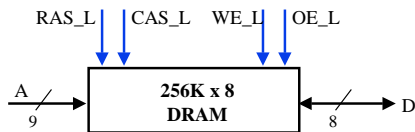


4/7/99

©UCB Spring 1999

CS152 / Kubiawicz
Lec18.10

Logic Diagram of a Typical DRAM



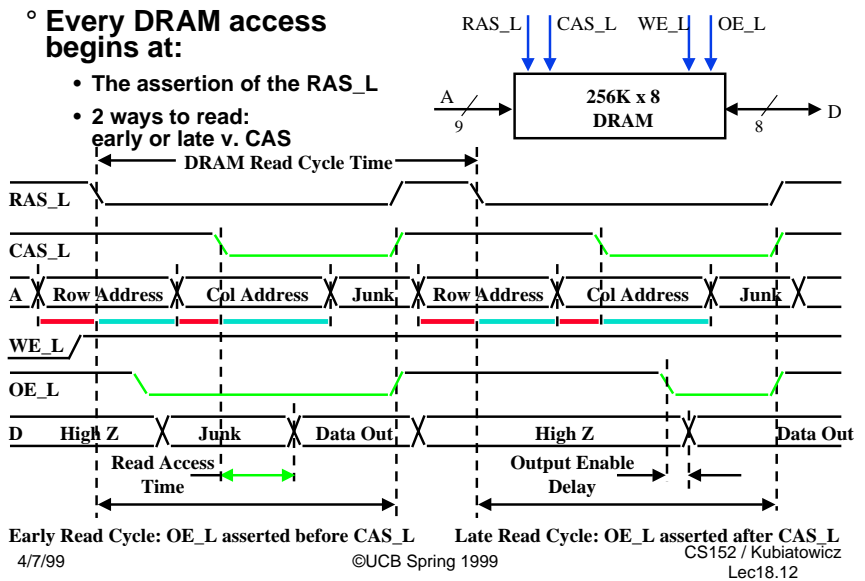
- Control Signals (RAS_L, CAS_L, WE_L, OE_L) are all active low
- Din and Dout are combined (D):
 - WE_L is asserted (Low), OE_L is disasserted (High)
 - D serves as the data input pin
 - WE_L is disasserted (High), OE_L is asserted (Low)
 - D is the data output pin
- Row and column addresses share the same pins (A)
 - RAS_L goes low: Pins A are latched in as row address
 - CAS_L goes low: Pins A are latched in as column address
 - RAS/CAS edge-sensitive

4/7/99

©UCB Spring 1999

CS152 / Kubiawicz
Lec18.11

DRAM Read Timing



4/7/99

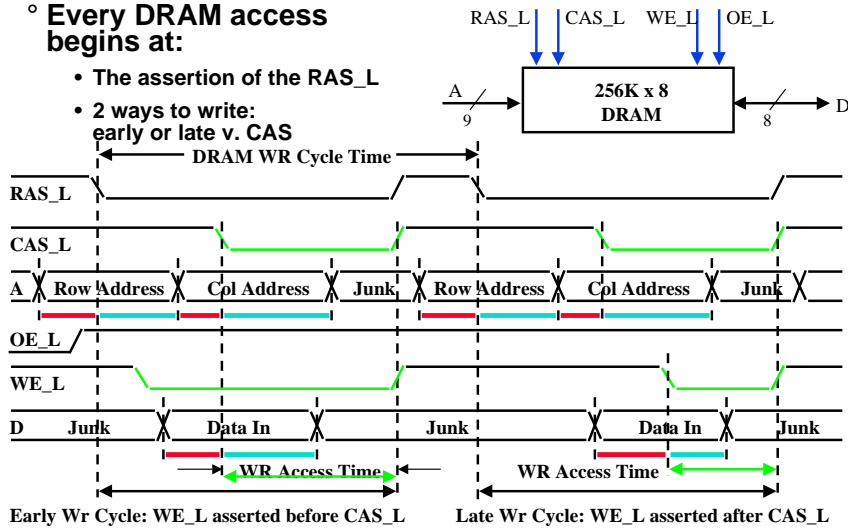
©UCB Spring 1999

CS152 / Kubiawicz
Lec18.12

DRAM Write Timing

Every DRAM access begins at:

- The assertion of the RAS_L
- 2 ways to write: early or late v. CAS

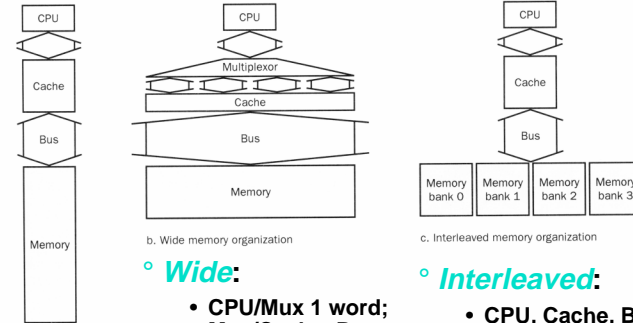


4/7/99

©UCB Spring 1999

CS152 / Kubiawicz
Lec18.13

Main Memory Performance



Simple:

- CPU, Cache, Bus, Memory same width (32 bits)

Wide:

- CPU/Mux 1 word; Mux/Cache, Bus, Memory N words (Alpha: 64 bits & 256 bits)

Interleaved:

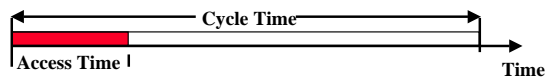
- CPU, Cache, Bus 1 word; Memory N Modules (4 Modules); example is word interleaved

4/7/99

©UCB Spring 1999

CS152 / Kubiawicz
Lec18.14

Main Memory Performance



DRAM (Read/Write) Cycle Time >> DRAM (Read/Write) Access Time

- 2:1; why?

DRAM (Read/Write) Cycle Time :

- How frequent can you initiate an access?
- Analogy: A little kid can only ask his father for money on Saturday

DRAM (Read/Write) Access Time:

- How quickly will you get what you want once you initiate an access?
- Analogy: As soon as he asks, his father will give him the money

DRAM Bandwidth Limitation analogy:

- What happens if he runs out of money on Wednesday?

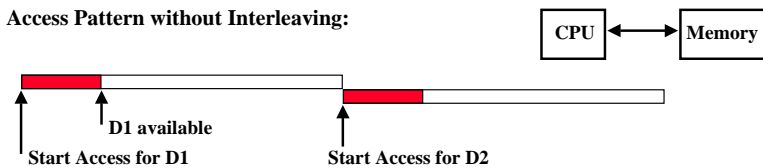
4/7/99

©UCB Spring 1999

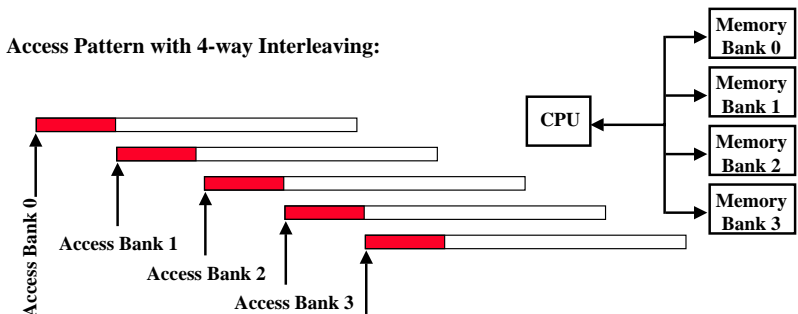
CS152 / Kubiawicz
Lec18.15

Increasing Bandwidth - Interleaving

Access Pattern without Interleaving:



Access Pattern with 4-way Interleaving:



We can Access Bank 0 again

4/7/99

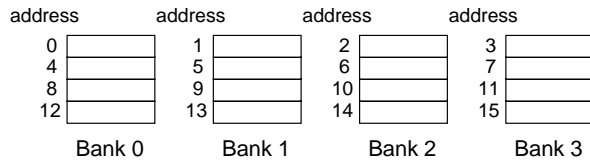
©UCB Spring 1999

CS152 / Kubiawicz
Lec18.16

Main Memory Performance

Timing model

- 1 to send address,
- 4 for access time, 10 cycle time, 1 to send data
- Cache Block is 4 words
- **Simple M.P.** = $4 \times (1+10+1) = 48$
- **Wide M.P.** = $1 + 10 + 1 = 12$
- **Interleaved M.P.** = $1+10+1 + 3 = 15$



4/7/99

©UCB Spring 1999

CS152 / Kubiawicz
Lec18.17

Independent Memory Banks

How many banks?

number banks \geq number clocks to access word in bank

- For sequential accesses, otherwise will return to original bank before it has next word ready
- Increasing DRAM \Rightarrow fewer chips \Rightarrow harder to have banks
 - Growth bits/chip DRAM : 50%-60%/yr
 - Nathan Myrvoid M/S: mature software growth (33%/yr for NT) - growth MB/\$ of DRAM (25%-30%/yr)

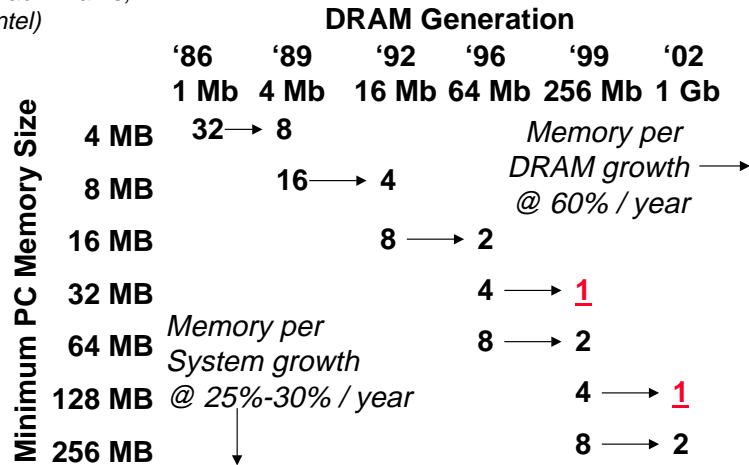
4/7/99

©UCB Spring 1999

CS152 / Kubiawicz
Lec18.18

Fewer DRAMs/System over Time

(from Pete MacWilliams, Intel)



4/7/99

©UCB Spring 1999

CS152 / Kubiawicz
Lec18.19

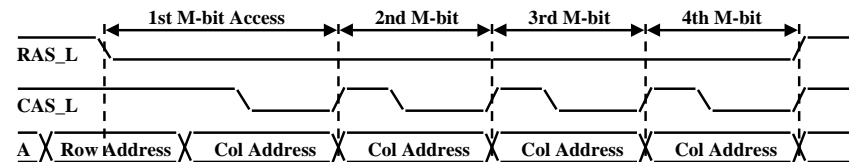
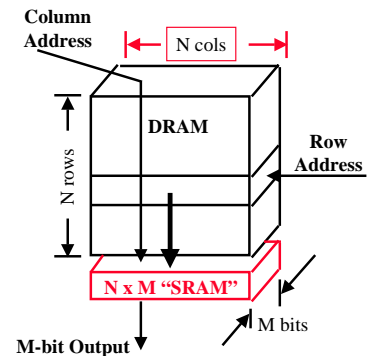
Fast Page Mode Operation

Regular DRAM Organization:

- N rows x N column x M-bit
- Read & Write M-bit at a time
- Each M-bit access requires a RAS / CAS cycle

Fast Page Mode DRAM

- N x M "SRAM" to save a row
- After a row is read into the register
 - Only CAS is needed to access other M-bit blocks on that row
 - RAS_L remains asserted while CAS_L is toggled



4/7/99

©UCB Spring 1999

CS152 / Kubiawicz
Lec18.20

Key DRAM Timing Parameters

- t_{RAC} : minimum time from RAS line falling to the valid data output.
 - Quoted as the speed of a DRAM
 - A fast 4Mb DRAM $t_{RAC} = 60$ ns
- t_{RC} : minimum time from the start of one row access to the start of the next.
 - $t_{RC} = 110$ ns for a 4Mbit DRAM with a t_{RAC} of 60 ns
- t_{CAC} : minimum time from CAS line falling to valid data output.
 - 15 ns for a 4Mbit DRAM with a t_{RAC} of 60 ns
- t_{PC} : minimum time from the start of one column access to the start of the next.
 - 35 ns for a 4Mbit DRAM with a t_{RAC} of 60 ns

4/7/99

©UCB Spring 1999

CS152 / Kubiawicz
Lec18.21

DRAMs over Time

1st Gen. Sample	DRAM Generation					
	'84	'87	'90	'93	'96	'99
Memory Size	1 Mb	4 Mb	16 Mb	64 Mb	256 Mb	1 Gb
Die Size (mm ²)	55	85	130	200	300	450
Memory Area (mm ²)	30	47	72	110	165	250
Memory Cell Area (μm ²)	28.84	11.1	4.26	1.64	0.61	0.23

(from Kazuhiro Sakashita, Mitsubishi)

4/7/99

©UCB Spring 1999

CS152 / Kubiawicz
Lec18.22

DRAM History

- DRAMs: capacity +60%/yr, cost -30%/yr
 - 2.5X cells/area, 1.5X die size in -3 years
- '97 DRAM fab line costs \$1B to \$2B
 - DRAM only: density, leakage v. speed
- Rely on increasing no. of computers & memory per computer (60% market)
 - SIMM or DIMM is replaceable unit
=> computers use any generation DRAM
- Commodity, second source industry
=> high volume, low profit, conservative
 - Little organization innovation in 20 years
page mode, EDO, Synch DRAM
- Order of importance: 1) Cost/bit 1a) Capacity
 - RAMBUS: 10X BW, +30% cost => little impact

4/7/99

©UCB Spring 1999

CS152 / Kubiawicz
Lec18.23

DRAM v. Desktop Microprocessors Cultures

Standards	pinout, package, refresh rate, capacity, ...	binary compatibility, IEEE 754, I/O bus
Sources	Multiple	Single
Figures of Merit	1) capacity, 1a) \$/bit 2) BW, 3) latency	1) SPEC speed 2) cost
Improve Rate/year	1) 60%, 1a) 25%, 2) 20%, 3) 7%	1) 60%, 2) little change

4/7/99

©UCB Spring 1999

CS152 / Kubiawicz
Lec18.24

Administrative Issues

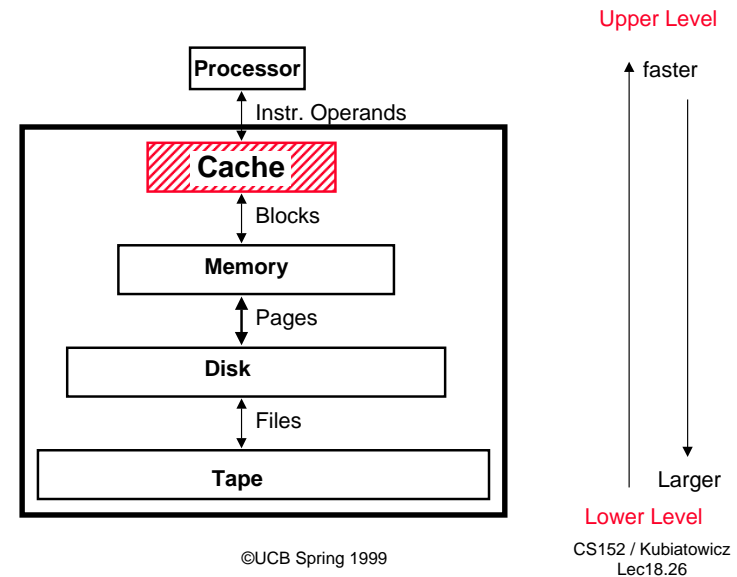
- Due tonight: breakdown of lab 6
- Continue reading Chapter 7 of your book (Memory Hierarchy)
- Second midterm coming up (Wed, April 21)
 - Microcoding/implementation of complex instructions
 - Pipelining
 - Hazards, branches, forwarding, CPI calculations
 - (may include something on dynamic scheduling)
 - Memory Hierarchy
 - Possibly something on I/O (see where we get in lectures)

4/7/99

©UCB Spring 1999

CS152 / Kubiawicz
Lec18.25

Recall: Levels of the Memory Hierarchy



4/7/99

©UCB Spring 1999

CS152 / Kubiawicz
Lec18.26

Cache performance equations:

- Time = IC x CT x (ideal CPI + memory stalls/inst)
- memory stalls/instruction =

$$\text{Average accesses/inst} \times \text{Miss Rate} \times \text{Miss Penalty} =$$

$$(\text{Average IFETCH/inst} \times \text{MissRate}_{\text{Inst}} \times \text{Miss Penalty}_{\text{Inst}}) +$$

$$(\text{Average Data/inst} \times \text{MissRate}_{\text{Data}} \times \text{Miss Penalty}_{\text{Data}})$$
- Assumes that ideal CPI includes Hit Times.
- Average Memory Access time =

$$\text{Hit Time} + (\text{Miss Rate} \times \text{Miss Penalty})$$

4/7/99

©UCB Spring 1999

CS152 / Kubiawicz
Lec18.27

Impact on Performance

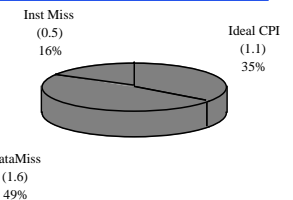
- Suppose a processor executes at
 - Clock Rate = 200 MHz (5 ns per cycle)
 - CPI = 1.1
 - 50% arith/logic, 30% ld/st, 20% control
- Suppose that 10% of memory operations get 50 cycle miss penalty
- Suppose that 1% of instructions get same miss penalty
- CPI = ideal CPI + average stalls per instruction

$$1.1(\text{cycles/ins}) +$$

$$[0.30 (\text{DataMops/ins}) \times 0.10 (\text{miss/DataMop}) \times 50 (\text{cycle/miss})] +$$

$$[1 (\text{InstMop/ins}) \times 0.01 (\text{miss/InstMop}) \times 50 (\text{cycle/miss})]$$

$$= (1.1 + 1.5 + .5) \text{ cycle/ins} = 3.1$$
- 58% of the time the proc is stalled waiting for memory!

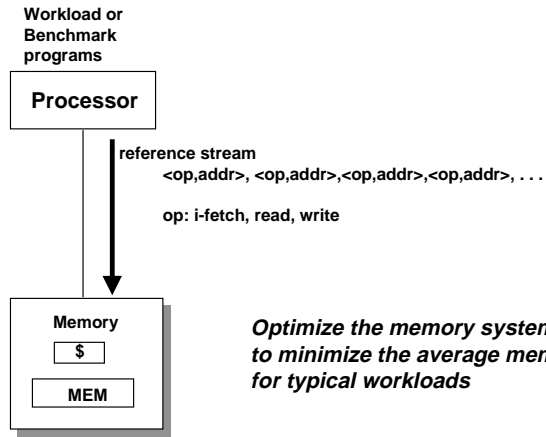


4/7/99

©UCB Spring 1999

CS152 / Kubiawicz
Lec18.28

The Art of Memory System Design



Optimize the memory system organization to minimize the average memory access time for typical workloads

4/7/99

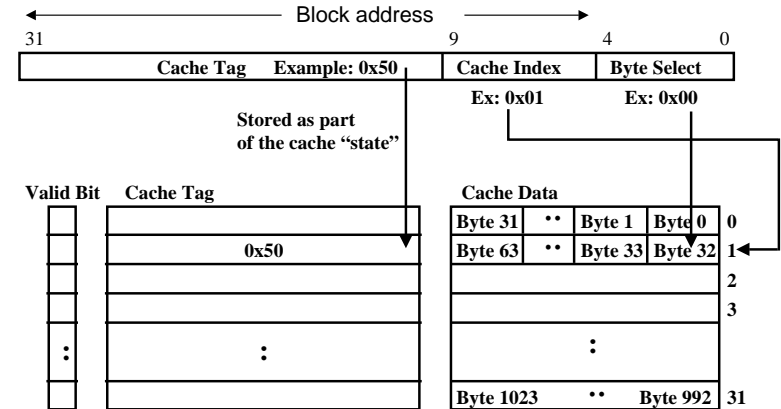
©UCB Spring 1999

CS152 / Kubiawicz
Lec18.29

Example: 1 KB Direct Mapped Cache with 32 B Blocks

For a 2^N byte cache:

- The uppermost $(32 - N)$ bits are always the Cache Tag
- The lowest N bits are the Byte Select (Block Size = 2^N)



4/7/99

©UCB Spring 1999

CS152 / Kubiawicz
Lec18.30

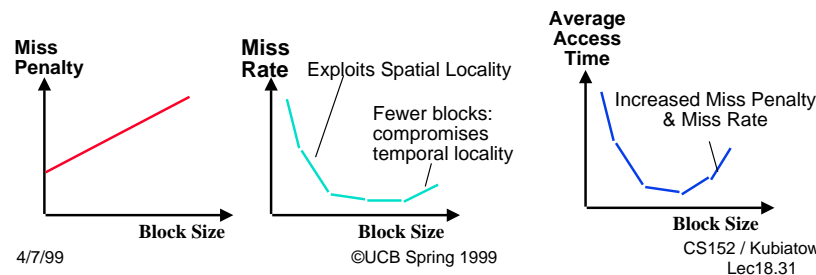
Block Size Tradeoff

In general, larger block size take advantage of spatial locality **BUT**:

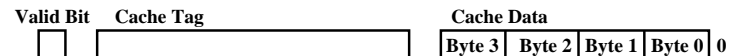
- Larger block size means larger miss penalty:
 - Takes longer time to fill up the block
- If block size is too big relative to cache size, miss rate will go up
 - Too few cache blocks

In general, Average Access Time:

$$= \text{Hit Time} \times (1 - \text{Miss Rate}) + \text{Miss Penalty} \times \text{Miss Rate}$$



Extreme Example: single line



Cache Size = 4 bytes

Block Size = 4 bytes

- Only ONE entry in the cache
- If an item is accessed, likely that it will be accessed again soon
 - But it is unlikely that it will be accessed again immediately!!!
 - The next access will likely to be a miss again
 - Continually loading data into the cache but discard (force out) them before they are used again
 - Worst nightmare of a cache designer: **Ping Pong Effect**
- Conflict Misses are misses caused by:
 - Different memory locations mapped to the same cache index
 - Solution 1: make the cache size bigger
 - Solution 2: Multiple entries for the same Cache Index

4/7/99

©UCB Spring 1999

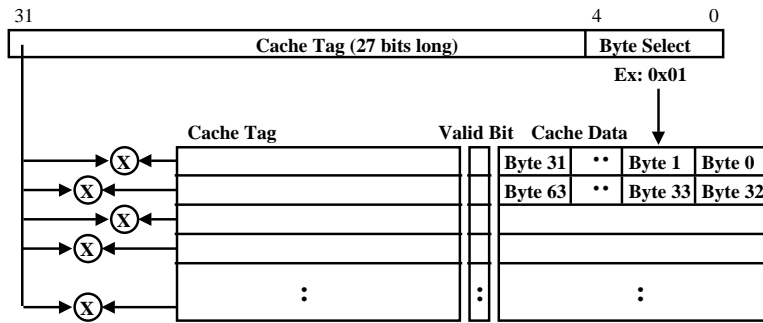
CS152 / Kubiawicz
Lec18.32

Another Extreme Example: Fully Associative

° Fully Associative Cache

- Forget about the Cache Index
- Compare the Cache Tags of all cache entries in parallel
- Example: Block Size = 32 B blocks, we need N 27-bit comparators

° By definition: Conflict Miss = 0 for a fully associative cache



4/7/99

©UCB Spring 1999

CS152 / Kubiawicz
Lec18.33

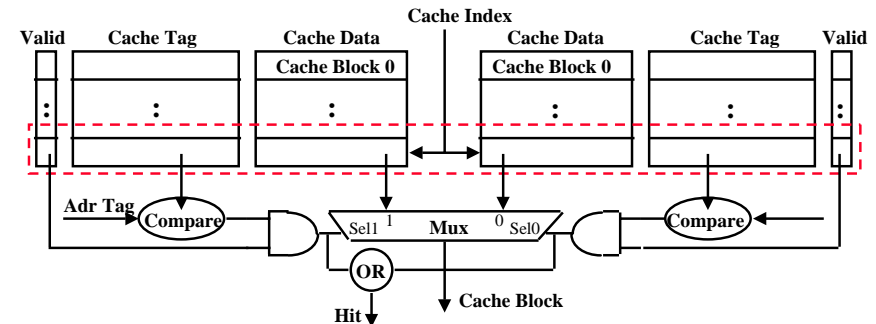
Set Associative Cache

° N-way set associative: N entries for each Cache Index

- N direct mapped caches operates in parallel

° Example: Two-way set associative cache

- Cache Index selects a "set" from the cache
- The two tags in the set are compared to the input in parallel
- Data is selected based on the tag result



4/7/99

©UCB Spring 1999

CS152 / Kubiawicz
Lec18.34

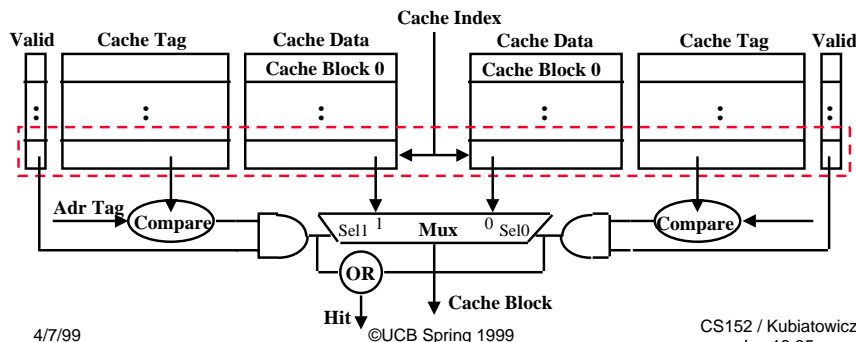
Disadvantage of Set Associative Cache

° N-way Set Associative Cache versus Direct Mapped Cache:

- N comparators vs. 1
- Extra MUX delay for the data
- Data comes **AFTER** Hit/Miss decision and set selection

° In a direct mapped cache, Cache Block is available **BEFORE** Hit/Miss:

- Possible to assume a hit and continue. Recover later if miss.



4/7/99

©UCB Spring 1999

CS152 / Kubiawicz
Lec18.35

A Summary on Sources of Cache Misses

° Compulsory (cold start or process migration, first reference): first access to a block

- "Cold" fact of life: not a whole lot you can do about it
- Note: If you are going to run "billions" of instruction, Compulsory Misses are insignificant

° Conflict (collision):

- Multiple memory locations mapped to the same cache location
- Solution 1: increase cache size
- Solution 2: increase associativity

° Capacity:

- Cache cannot contain all blocks access by the program
- Solution: increase cache size

° Coherence (Invalidation): other process (e.g., I/O) updates memory

4/7/99

©UCB Spring 1999

CS152 / Kubiawicz
Lec18.36

Source of Cache Misses Quiz

Assume constant cost.

	Direct Mapped	N-way Set Associative	Fully Associative
Cache Size: Small, Medium, Big?			
Compulsory Miss:			
Conflict Miss			
Capacity Miss			
Coherence Miss			

Choices: Zero, Low, Medium, High, Same

4/7/99

©UCB Spring 1999

CS152 / Kubiatiowicz
Lec18.37

Sources of Cache Misses Answer

	Direct Mapped	N-way Set Associative	Fully Associative
Cache Size	Big	Medium	Small
Compulsory Miss	Same	Same	Same
Conflict Miss	High	Medium	Zero
Capacity Miss	Low	Medium	High
Coherence Miss	Same	Same	Same

Note:

If you are going to run “billions” of instruction, Compulsory Misses are insignificant.

4/7/99

©UCB Spring 1999

CS152 / Kubiatiowicz
Lec18.38

Four Questions for Caches and Memory Hierarchy

- Q1: Where can a block be placed in the upper level? (*Block placement*)
- Q2: How is a block found if it is in the upper level? (*Block identification*)
- Q3: Which block should be replaced on a miss? (*Block replacement*)
- Q4: What happens on a write? (*Write strategy*)

4/7/99

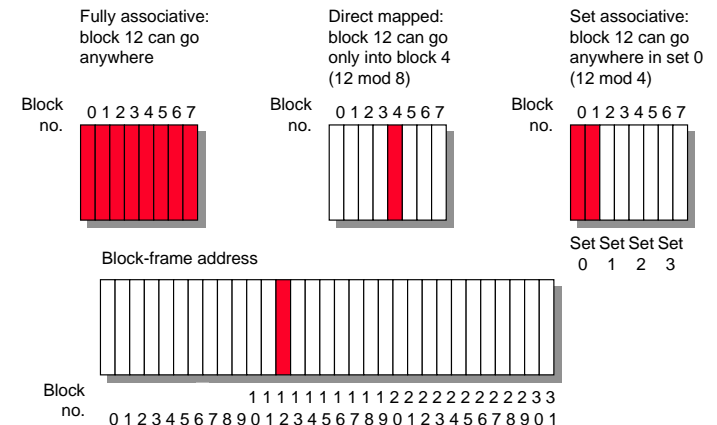
©UCB Spring 1999

CS152 / Kubiatiowicz
Lec18.39

Q1: Where can a block be placed in the upper level?

◦ Block 12 placed in 8 block cache:

- Fully associative, direct mapped, 2-way set associative
- S.A. Mapping = Block Number Modulo Number Sets

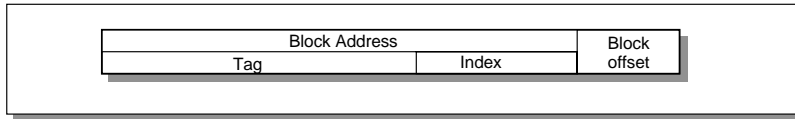


4/7/99

©UCB Spring 1999

CS152 / Kubiatiowicz
Lec18.40

Q2: How is a block found if it is in the upper level?



- Direct indexing (using index and block offset), tag compares, or combination
- Increasing associativity shrinks index, expands tag

Q3: Which block should be replaced on a miss?

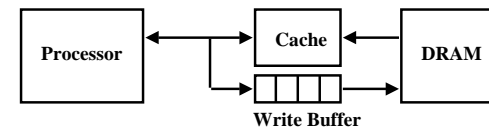
- Easy for Direct Mapped
- Set Associative or Fully Associative:
 - Random
 - LRU (Least Recently Used)

Associativity:	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%

Q4: What happens on a write?

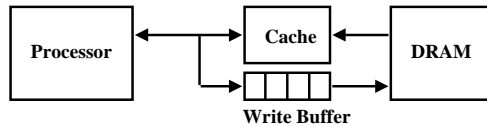
- **Write through**—The information is written to both the block in the cache and to the block in the lower-level memory.
- **Write back**—The information is written only to the block in the cache. The modified cache block is written to main memory only when it is replaced.
 - is block clean or dirty?
- Pros and Cons of each?
 - WT: read misses cannot result in writes
 - WB: no writes of repeated writes
- WT always combined with write buffers so that don't wait for lower level memory

Write Buffer for Write Through



- A Write Buffer is needed between the Cache and Memory
 - Processor: writes data into the cache and the write buffer
 - Memory controller: write contents of the buffer to memory
- Write buffer is just a FIFO:
 - Typical number of entries: 4
 - Works fine if: Store frequency (w.r.t. time) $\ll 1 / \text{DRAM write cycle}$
- Memory system designer's nightmare:
 - Store frequency (w.r.t. time) $> 1 / \text{DRAM write cycle}$
 - Write buffer saturation

Write Buffer Saturation

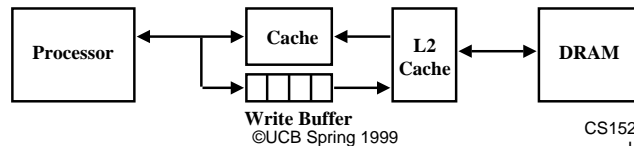


Store frequency (w.r.t. time) > 1 / DRAM write cycle

- If this condition exist for a long period of time (CPU cycle time too quick and/or too many store instructions in a row):
 - Store buffer will overflow no matter how big you make it
 - The CPU Cycle Time <= DRAM Write Cycle Time

Solution for write buffer saturation:

- Use a write back cache
- Install a second level (L2) cache: (does this always work?)



4/7/99

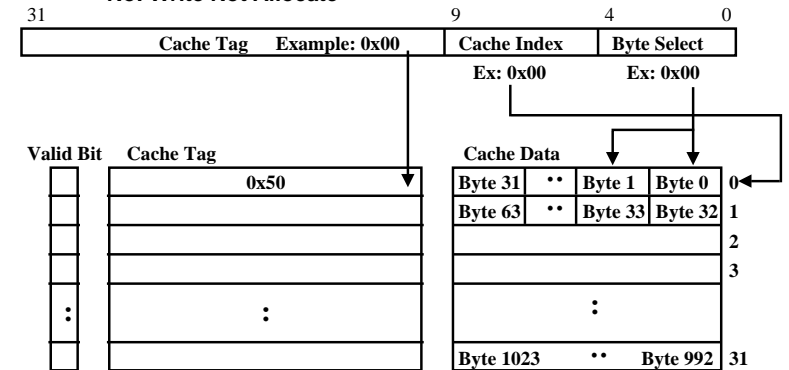
©UCB Spring 1999

CS152 / Kubiawicz
Lec18.45

Write-miss Policy: Write Allocate versus Not Allocate

Assume: a 16-bit write to memory location 0x0 and causes a miss

- Do we read in the block?
 - Yes: Write Allocate
 - No: Write Not Allocate



4/7/99

©UCB Spring 1999

CS152 / Kubiawicz
Lec18.46

Impact of Memory Hierarchy on Algorithms

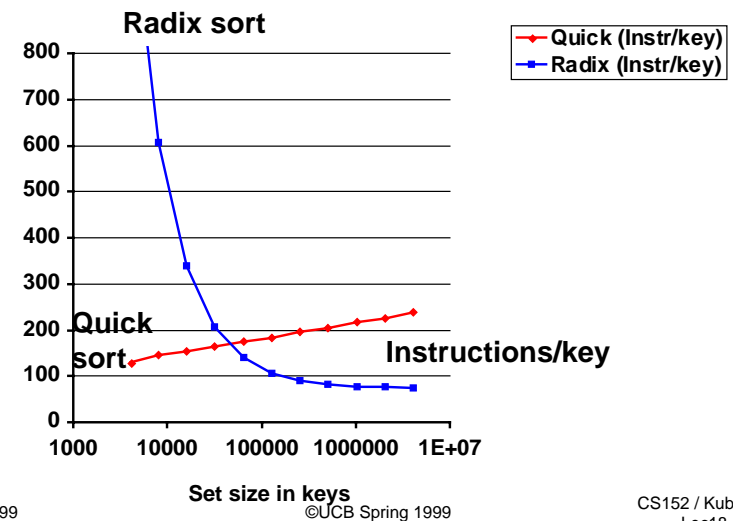
- Today CPU time is a function of (ops, cache misses) vs. just f(ops):
What does this mean to Compilers, Data structures, Algorithms?
- “The Influence of Caches on the Performance of Sorting” by A. LaMarca and R.E. Ladner. *Proceedings of the Eighth Annual ACM-SIAM Symposium on Discrete Algorithms*, January, 1997, 370-379.
- Quicksort: fastest comparison based sorting algorithm when all keys fit in memory
- Radix sort: also called “linear time” sort because for keys of fixed length and fixed radix a constant number of passes over the data is sufficient independent of the number of keys
- For Alphasstation 250, 32 byte blocks, direct mapped L2 2MB cache, 8 byte keys, from 4000 to 4000000

4/7/99

©UCB Spring 1999

CS152 / Kubiawicz
Lec18.47

Quicksort vs. Radix as vary number keys: Instructions

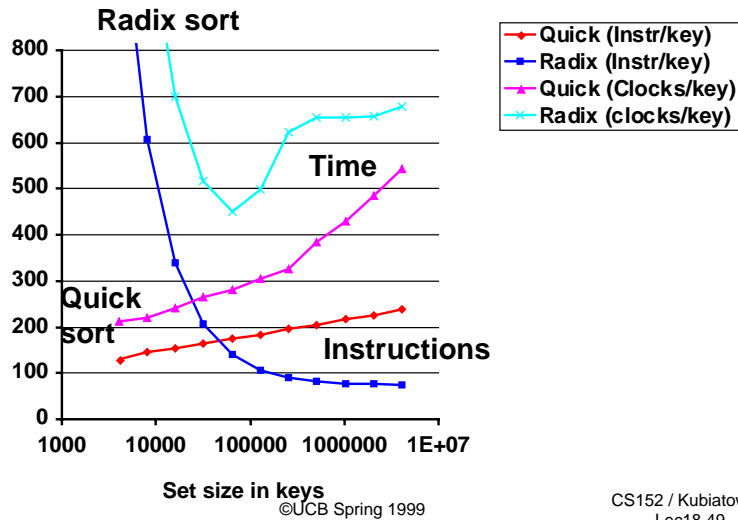


4/7/99

©UCB Spring 1999

CS152 / Kubiawicz
Lec18.48

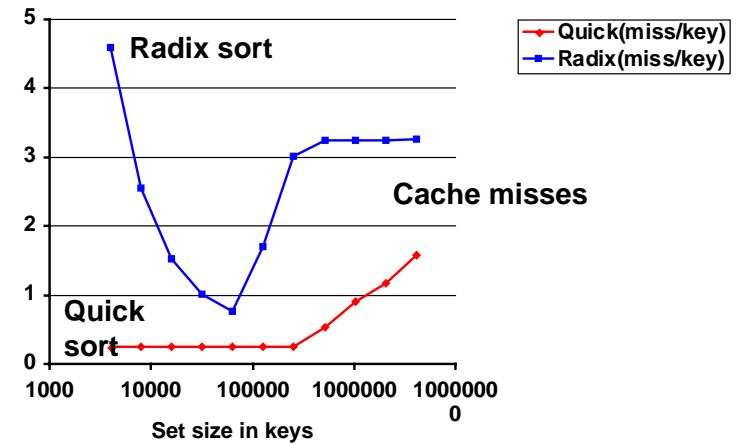
Quicksort vs. Radix as vary number keys: Instrs & Time



4/7/99

CS152 / Kubiatiowicz
Lec18.49

Quicksort vs. Radix as vary number keys: Cache misses



4/7/99

©UCB Spring 1999

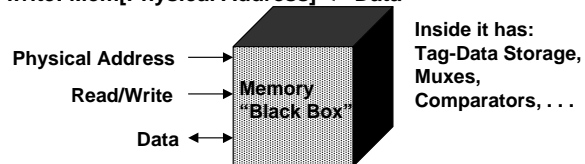
CS152 / Kubiatiowicz
Lec18.50

What is proper approach to fast algorithms?

How Do you Design a Cache?

◦ Set of Operations that must be supported

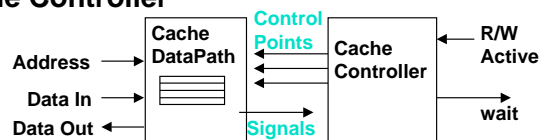
- read: $\text{data} \leftarrow \text{Mem}[\text{Physical Address}]$
- write: $\text{Mem}[\text{Physical Address}] \leftarrow \text{Data}$



◦ Determine the internal register transfers

◦ Design the Datapath

◦ Design the Cache Controller



4/7/99

©UCB Spring 1999

CS152 / Kubiatiowicz
Lec18.51

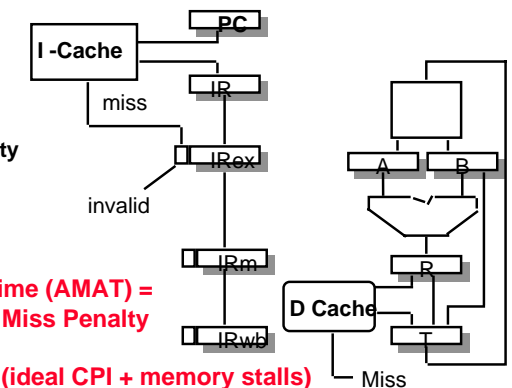
Impact on Cycle Time

Cache Hit Time:

directly tied to clock rate
increases with cache size
increases with associativity

**Average Memory Access time (AMAT) =
Hit Time + Miss Rate x Miss Penalty**

Compute Time = IC x CT x (ideal CPI + memory stalls)



4/7/99

©UCB Spring 1999

CS152 / Kubiatiowicz
Lec18.52

Example: direct map allows miss signal after data

What happens on a Cache miss?

For in-order pipeline, 2 options:

- Freeze pipeline in Mem stage (popular early on: Sparc, R4000)

```
IF ID EX Mem stall stall stall ... stall Mem Wr
IF ID EX stall stall stall ... stall stall Ex Wr
```

- Use Full/Empty bits in registers + MSHR queue

- MSHR = "Miss Status/Handler Registers" (Kroft)
Each entry in this queue keeps track of status of outstanding memory requests to one complete memory line.
 - Per cache-line: keep info about memory address.
 - For each word: register (if any) that is waiting for result.
 - Used to "merge" multiple requests to one memory line
- New load creates MSHR entry and sets destination register to "Empty". Load is "released" from pipeline.
- Attempt to use register before result returns causes instruction to block in decode stage.
- Limited "out-of-order" execution with respect to loads.
Popular with in-order superscalar architectures.

Out-of-order pipelines already have this functionality built in... (load queues, etc).

4/7/99

©UCB Spring 1999

CS152 / Kubiawicz
Lec18.53

Improving Cache Performance: 3 general options

$$\text{Time} = \text{IC} \times \text{CT} \times (\text{ideal CPI} + \text{memory stalls/instruction})$$

memory stalls/instruction =

$$\text{Average memory accesses/inst} \times \text{AMAT} = (\text{Average IFETCH/inst} \times \text{AMAT}_{\text{Inst}}) + (\text{Average DMEM/inst} \times \text{AMAT}_{\text{Data}}) +$$

Average Memory Access time =

$$\text{Hit Time} + (\text{Miss Rate} \times \text{Miss Penalty}) =$$

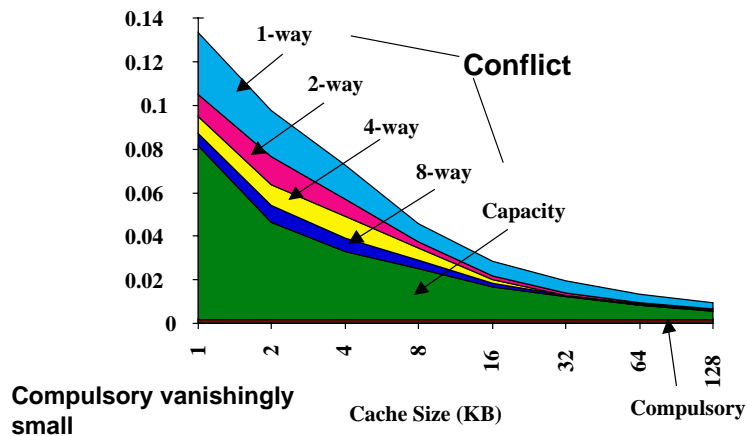
- Reduce the miss rate,
- Reduce the miss penalty, or
- Reduce the time to hit in the cache.

4/7/99

©UCB Spring 1999

CS152 / Kubiawicz
Lec18.54

3Cs Absolute Miss Rate (SPEC92)



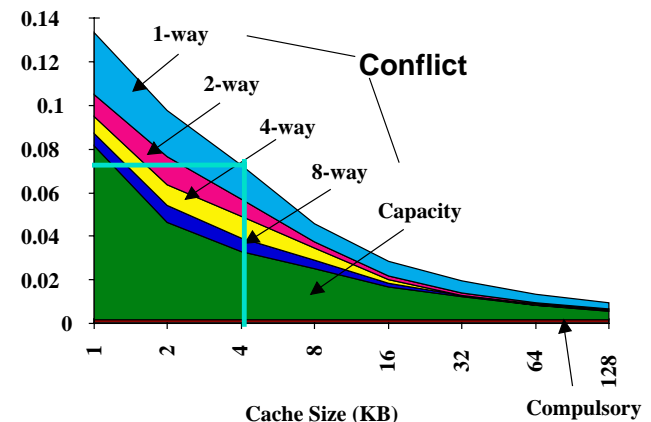
4/7/99

©UCB Spring 1999

CS152 / Kubiawicz
Lec18.55

2:1 Cache Rule

miss rate 1-way associative cache size X
= miss rate 2-way associative cache size X/2

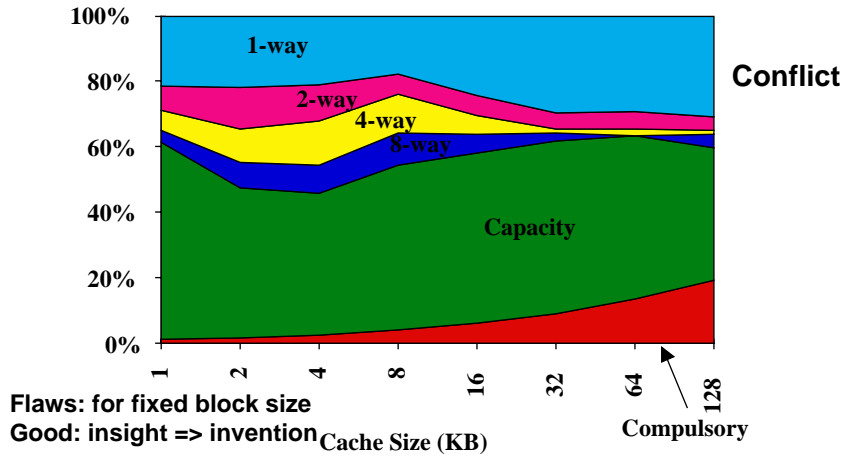


4/7/99

©UCB Spring 1999

CS152 / Kubiawicz
Lec18.56

3Cs Relative Miss Rate

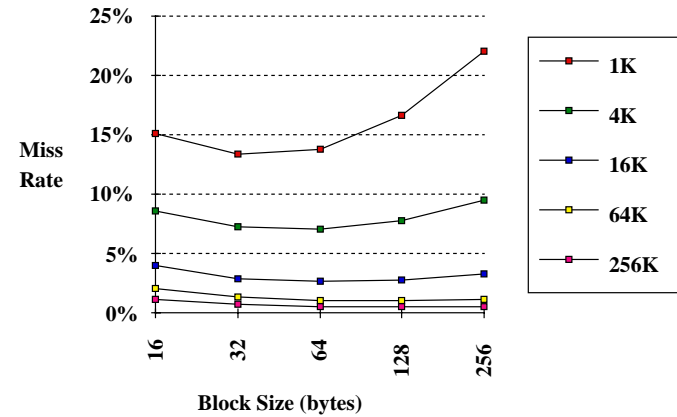


4/7/99

©UCB Spring 1999

CS152 / Kubiawicz
Lec18.57

1. Reduce Misses via Larger Block Size



4/7/99

©UCB Spring 1999

CS152 / Kubiawicz
Lec18.58

2. Reduce Misses via Higher Associativity

- **2:1 Cache Rule:**
 - Miss Rate DM cache size N - Miss Rate 2-way cache size N/2
- **Beware: Execution time is only final measure!**
 - Will Clock Cycle time increase?
 - Hill [1988] suggested hit time for 2-way vs. 1-way external cache +10%, internal + 2%

4/7/99

©UCB Spring 1999

CS152 / Kubiawicz
Lec18.59

Example: Avg. Memory Access Time vs. Miss Rate

- Example: assume CCT = 1.10 for 2-way, 1.12 for 4-way, 1.14 for 8-way vs. CCT direct mapped

Cache Size (KB)	Associativity			
	1-way	2-way	4-way	8-way
1	2.33	2.15	2.07	2.01
2	1.98	1.86	1.76	1.68
4	1.72	1.67	1.61	1.53
8	1.46	1.48	1.47	1.43
16	1.29	1.32	1.32	1.32
32	1.20	1.24	1.25	1.27
64	1.14	1.20	1.21	1.23
128	1.10	1.17	1.18	1.20

(Red means A.M.A.T. not improved by more associativity)

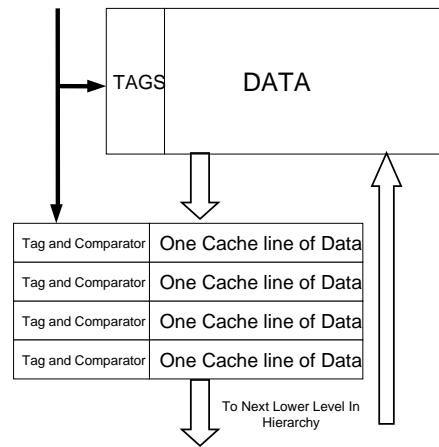
4/7/99

©UCB Spring 1999

CS152 / Kubiawicz
Lec18.60

3. Reducing Misses via a “Victim Cache”

- How to combine fast hit time of direct mapped yet still avoid conflict misses?
- Add buffer to place data discarded from cache
- Jouppi [1990]: 4-entry victim cache removed 20% to 95% of conflicts for a 4 KB direct mapped data cache
- Used in Alpha, HP machines



4/7/99

©UCB Spring 1999

CS152 / Kubiawicz
Lec18.61

4. Reducing Misses via “Pseudo-Associativity”

- How to combine fast hit time of Direct Mapped and have the lower conflict misses of 2-way SA cache?
- Divide cache: on a miss, check other half of cache to see if there, if so have a **pseudo-hit** (slow hit)



- Drawback: CPU pipeline is hard if hit takes 1 or 2 cycles
 - Better for caches not tied directly to processor (L2)
 - Used in MIPS R1000 L2 cache, similar in UltraSPARC

4/7/99

©UCB Spring 1999

CS152 / Kubiawicz
Lec18.62

5. Reducing Misses by Hardware Prefetching

- E.g., Instruction Prefetching
 - Alpha 21064 fetches 2 blocks on a miss
 - Extra block placed in “**stream buffer**”
 - On miss check stream buffer
- Works with data blocks too:
 - Jouppi [1990] 1 data stream buffer got 25% misses from 4KB cache; 4 streams got 43%
 - Palacharla & Kessler [1994] for scientific programs for 8 streams got 50% to 70% of misses from 2 64KB, 4-way set associative caches
- Prefetching relies on having extra memory bandwidth that can be used without penalty

4/7/99

©UCB Spring 1999

CS152 / Kubiawicz
Lec18.63

6. Reducing Misses by Software Prefetching Data

- Data Prefetch
 - Load data into register (HP PA-RISC loads)
 - Cache Prefetch: load into cache (MIPS IV, PowerPC, SPARC v. 9)
 - Special prefetching instructions cannot cause faults; a form of speculative execution
- Issuing Prefetch Instructions takes time
 - Is cost of prefetch issues < savings in reduced misses?
 - Higher superscalar reduces difficulty of issue bandwidth

4/7/99

©UCB Spring 1999

CS152 / Kubiawicz
Lec18.64

7. Reducing Misses by Compiler Optimizations

- McFarling [1989] reduced caches misses by 75% on 8KB direct mapped cache, 4 byte blocks in software
- Instructions
 - Reorder procedures in memory so as to reduce conflict misses
 - Profiling to look at conflicts(using tools they developed)
- Data
 - Merging Arrays:** improve spatial locality by single array of compound elements vs. 2 arrays
 - Loop Interchange:** change nesting of loops to access data in order stored in memory
 - Loop Fusion:** Combine 2 independent loops that have same looping and some variables overlap
 - Blocking:** Improve temporal locality by accessing “blocks” of data repeatedly vs. going down whole columns or rows

4/7/99

©UCB Spring 1999

CS152 / Kubiawicz
Lec18.65

Summary #1 / 3:

- 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 Major Categories of Cache Misses:
 - Compulsory Misses:** sad facts of life. Example: cold start misses.
 - Conflict Misses:** increase cache size and/or associativity. Nightmare Scenario: ping pong effect!
 - Capacity Misses:** increase cache size
 - Coherence Misses: invalidation caused by “external” processors or I/O
- Cache Design Space
 - total size, block size, associativity
 - replacement policy
 - write-hit policy (write-through, write-back)
 - write-miss policy

4/7/99

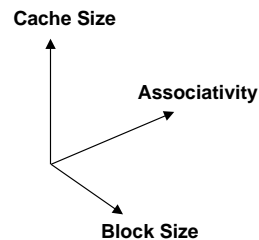
©UCB Spring 1999

CS152 / Kubiawicz
Lec18.66

Summary #2 / 3: The Cache Design Space

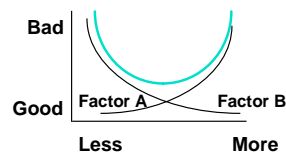
Several interacting dimensions

- cache size
- block size
- associativity
- replacement policy
- write-through vs write-back
- write allocation



The optimal choice is a compromise

- depends on access characteristics
 - workload
 - use (I-cache, D-cache, TLB)
- depends on technology / cost



Simplicity often wins

4/7/99

©UCB Spring 1999

CS152 / Kubiawicz
Lec18.67

Summary #3 / 3: Cache Miss Optimization

Lots of techniques people use to improve the miss rate of caches:

miss rate	Technique	MR	MP	HT	Complexity
	Larger Block Size	+	-	-	0
	Higher Associativity	+			1
	Victim Caches	+			2
	Pseudo-Associative Caches	+			2
	HW Prefetching of Instr/Data	+			2
	Compiler Controlled Prefetching	+			3
	Compiler Reduce Misses	+			0

4/7/99

©UCB Spring 1999

CS152 / Kubiawicz
Lec18.68