

CS252
Graduate Computer Architecture
Lecture 4

Caches and Memory Systems

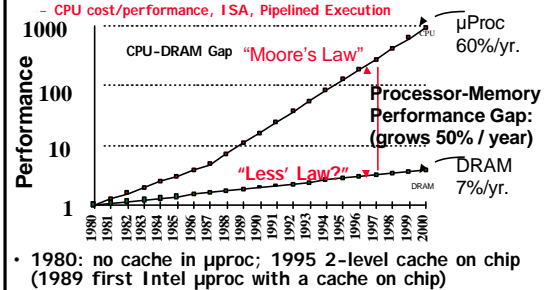
January 26, 2001
Prof. John Kubiatowicz

1/26/01

CS252/Kubiatowicz
Lec. 4.1

Review: Who Cares About the Memory Hierarchy?

• Processor Only Thus Far in Course:



1/26/01

CS252/Kubiatowicz
Lec. 4.2

Review: Cache performance

• Miss-oriented Approach to Memory Access:

$$CPUtime = IC \cdot \frac{CPI}{Execution} + \frac{MemAccess}{Inst} \cdot MissRate \cdot MissPenalty \cdot CycleTime$$

$$CPUtime = IC \cdot \frac{CPI}{Execution} + \frac{MemMisses}{Inst} \cdot MissPenalty \cdot CycleTime$$

- $CPI_{Execution}$ includes ALU and Memory instructions

• Separating out Memory component entirely

- AMAT = Average Memory Access Time

- CPI_{ALUOps} does not include memory instructions

$$CPUtime = IC \cdot \frac{AluOps}{Inst} \cdot CPI_{AluOps} + \frac{MemAccess}{Inst} \cdot AMAT \cdot CycleTime$$

$$AMAT = HitTime + MissRate \cdot MissPenalty$$

$$= (HitTime_{Inst} + MissRate_{Inst} \cdot MissPenalty_{Inst}) + (HitTime_{Data} + MissRate_{Data} \cdot MissPenalty_{Data})$$

1/26/01

CS252/Kubiatowicz
Lec. 4.3

Review: Reducing Misses

• Classifying Misses: 3 Cs

- **Compulsory**—Misses in even an Infinite Cache
- **Capacity**—Misses in Fully Associative Size X Cache
- **Conflict**—Misses in N-way Associative, Size X Cache

• More recent, 4th "C":

- **Coherence** - Misses caused by cache coherence.

1/26/01

CS252/Kubiatowicz
Lec. 4.4

Review: Miss Rate Reduction

$$AMAT = HitTime + MissRate \cdot MissPenalty$$

• 3 Cs: Compulsory, Capacity, Conflict

1. Reduce Misses via Larger Block Size
2. Reduce Misses via Higher Associativity
3. Reducing Misses via Victim Cache
4. Reducing Misses via Pseudo-Associativity
5. Reducing Misses by HW Prefetching Instr, Data
6. Reducing Misses by SW Prefetching Data
7. Reducing Misses by Compiler Optimizations

• Prefetching comes in two flavors:

- Binding prefetch: Requests load directly into register.
 - » Must be correct address and register!
- Non-Binding prefetch: Load into cache.
 - » Can be incorrect. Frees HW/SW to guess!

1/26/01

CS252/Kubiatowicz
Lec. 4.5

Improving Cache Performance Continued

1. Reduce the miss rate,
2. Reduce the miss penalty, or
3. Reduce the time to hit in the cache.

$$AMAT = HitTime + MissRate \cdot MissPenalty$$

1/26/01

CS252/Kubiatowicz
Lec. 4.6

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).

1/26/01
CS252/Kubitowicz Lec 4.7

Write Policy: Write-Through vs Write-Back

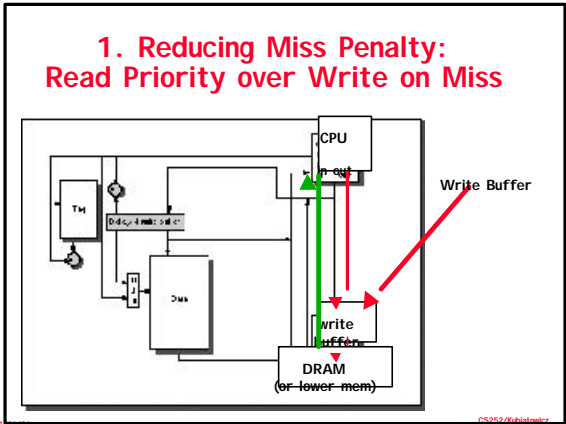
- Write-through:** all writes update cache and underlying memory/cache
 - Can always discard cached data - most up-to-date data is in memory
 - Cache control bit: only a *valid* bit
- Write-back:** all writes simply update cache
 - Can't just discard cached data - may have to write it back to memory
 - Cache control bits: both *valid* and *dirty* bits
- Other Advantages:**
 - Write-through:
 - memory (or other processors) always have latest data
 - Simpler management of cache
 - Write-back:
 - much lower bandwidth, since data often overwritten multiple times
 - Better tolerance to long-latency memory?

1/26/01
CS252/Kubitowicz Lec 4.8

Write Policy 2: Write Allocate vs Non-Allocate (What happens on write-miss)

- Write allocate:** allocate new cache line in cache
 - Usually means that you have to do a "read miss" to fill in rest of the cache-line!
 - Alternative: per/word valid bits
- Write non-allocate (or "write-around"):**
 - Simply send write data through to underlying memory/cache - don't allocate new cache line!

1/26/01
CS252/Kubitowicz Lec 4.9



1. Reducing Miss Penalty: Read Priority over Write on Miss

- Write-through with write buffers offer RAW conflicts with main memory reads on cache misses**
 - If simply wait for write buffer to empty, might increase read miss penalty (old MIPS 1000 by 50%)
 - Check write buffer contents before read; if no conflicts, let the memory access continue
- Write-back also want buffer to hold misplaced blocks**
 - Read miss replacing dirty block
 - Normal: Write dirty block to memory, and then do the read
 - Instead copy the dirty block to a write buffer, then do the read, and then do the write
 - CPU stall less since restarts as soon as do read

1/26/01
CS252/Kubitowicz Lec 4.11

2. Reduce Miss Penalty: Early Restart and Critical Word First

- Don't wait for full block to be loaded before restarting CPU
 - Early restart**—As soon as the requested word of the block arrives, send it to the CPU and let the CPU continue execution
 - Critical Word First**—Request the missed word first from memory and send it to the CPU as soon as it arrives; let the CPU continue execution while filling the rest of the words in the block. Also called *wrapped fetch* and *requested word first*
- Generally useful only in large blocks,
- Spatial locality a problem; tend to want next sequential word, so not clear if benefit by early restart

1/26/01
CS252/Kubitowicz Lec 4.12

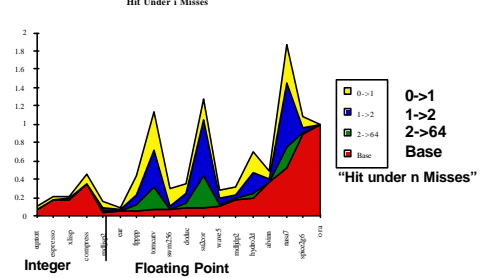
3. Reduce Miss Penalty: Non-blocking Caches to reduce stalls on misses

- **Non-blocking cache** or **lookup-free cache** allow data cache to continue to supply cache hits during a miss
 - requires F/E bits on registers or out-of-order execution
 - requires multi-bank memories
- **"hit under miss"** reduces the effective miss penalty by working during miss vs. ignoring CPU requests
- **"hit under multiple miss"** or **"miss under miss"** may further lower the effective miss penalty by overlapping multiple misses
 - Significantly increases the complexity of the cache controller as there can be multiple outstanding memory accesses
 - Requires **multiple memory banks** (otherwise cannot support)
 - Pentium Pro allows 4 outstanding memory misses

1/26/01

CS252/Kabiatowicz Lec. 4.13

Value of Hit Under Miss for SPEC (Normalized to blocking cache)



FP programs on average: AMAT= 0.68 -> 0.52 -> 0.34 -> 0.26
 Int programs on average: AMAT= 0.24 -> 0.20 -> 0.19 -> 0.19
 8 KB Data Cache, Direct Mapped, 32B block, 16 cycle miss

1/26/01

CS252/Kabiatowicz Lec. 4.14

4. Second level cache

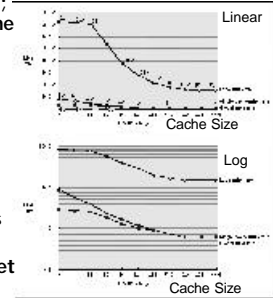
- **L2 Equations**
 $AMAT = Hit\ Time_{L1} + Miss\ Rate_{L1} \times Miss\ Penalty_{L1}$
 $Miss\ Penalty_{L1} = Hit\ Time_{L2} + Miss\ Rate_{L2} \times Miss\ Penalty_{L2}$
 $AMAT = Hit\ Time_{L1} + Miss\ Rate_{L1} \times (Hit\ Time_{L2} + Miss\ Rate_{L2} \times Miss\ Penalty_{L2})$
- **Definitions:**
 - **Local miss rate**— misses in this cache divided by the total number of memory accesses *to this cache* ($Miss\ rate_{L2}$)
 - **Global miss rate**—misses in this cache divided by the total number of memory accesses *generated by the CPU* ($Miss\ Rate_{L1} \times Miss\ Rate_{L2}$)
 - Global Miss Rate is what matters

1/26/01

CS252/Kabiatowicz Lec. 4.15

Comparing Local and Global Miss Rates

- 32 KByte 1st level cache; Increasing 2nd level cache
- Global miss rate close to single level cache rate provided $L2 \gg L1$
- Don't use local miss rate
- L2 not tied to CPU clock cycle!
- Cost & A.M.A.T.
- Generally Fast Hit Times and fewer misses
- Since hits are few, target miss reduction



1/26/01

CS252/Kabiatowicz Lec. 4.16

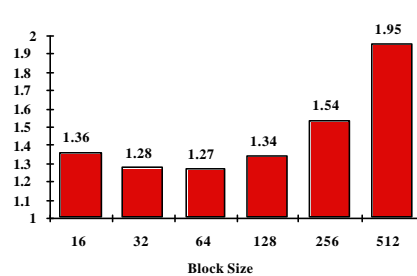
Reducing Misses: Which apply to L2 Cache?

- **Reducing Miss Rate**
 1. Reduce Misses via Larger Block Size
 2. Reduce Conflict Misses via Higher Associativity
 3. Reducing Conflict Misses via Victim Cache
 4. Reducing Conflict Misses via Pseudo-Associativity
 5. Reducing Misses by HW Prefetching Instr., Data
 6. Reducing Misses by SW Prefetching Data
 7. Reducing Capacity/Conf. Misses by Compiler Optimizations

1/26/01

CS252/Kabiatowicz Lec. 4.17

L2 cache block size & A.M.A.T.



• 32KB L1, 8 byte path to memory

1/26/01

CS252/Kabiatowicz Lec. 4.18

Reducing Miss Penalty Summary

$$AMAT = HitTime + MissRate \cdot MissPenalty$$

- Four techniques
 - Read priority over write on miss
 - Early Restart and Critical Word First on miss
 - Non-blocking Caches (Hit under Miss, Miss under Miss)
 - Second Level Cache
- Can be applied recursively to Multilevel Caches
 - Danger is that time to DRAM will grow with multiple levels in between
 - First attempts at L2 caches can make things worse, since increased worst case is worse

1/26/01

CS252/Kubitowicz Lec. 4.19

Main Memory Background

Performance of Main Memory:

- **Latency:** Cache Miss Penalty
 - > **Access Time:** time between request and word arrives
 - > **Cycle Time:** time between requests
- **Bandwidth:** I/O & Large Block Miss Penalty (L2)

Main Memory is **DRAM:** Dynamic Random Access Memory

- Dynamic since needs to be **refreshed** periodically (8 ms, 1% time)
- Addresses divided into 2 halves (Memory as a 2D matrix):
 - > **RAS** or **Row Access Strobe**
 - > **CAS** or **Column Access Strobe**

Cache uses **SRAM:** Static Random Access Memory

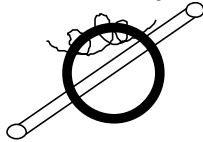
- No refresh (6 transistors/bit vs. 1 transistor)
- Size:** DRAM/SRAM - 4-8,
- Cost/Cycle time:** SRAM/DRAM - 8-16

1/26/01

CS252/Kubitowicz Lec. 4.20

Main Memory Deep Background

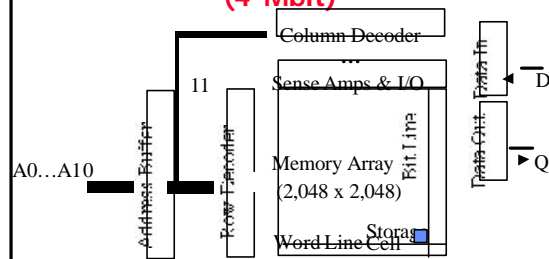
- "Out-of-Core", "In-Core", "Core Dump"?
- "Core memory"?
- Non-volatile, magnetic
- Lost to 4 Kbit DRAM (today using 64Kbit DRAM)
- Access time 750 ns, cycle time 1500-3000 ns



1/26/01

CS252/Kubitowicz Lec. 4.21

DRAM logical organization (4 Mbit)



- Square root of bits per RAS/CAS

1/26/01

CS252/Kubitowicz Lec. 4.22

4 Key DRAM Timing Parameters

- **t_{RAC}**: minimum time from RAS line falling to the valid data output.
 - Quoted as the speed of a DRAM when buy
 - A typical 4Mb DRAM t_{RAC} = 60 ns
 - Speed of DRAM since on purchase sheet?
- **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

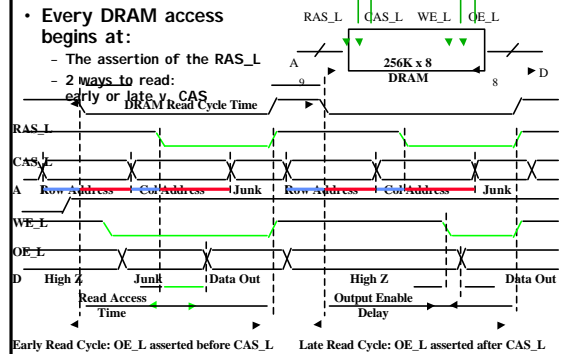
1/26/01

CS252/Kubitowicz Lec. 4.23

DRAM Read Timing

- Every DRAM access begins at:

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



1/26/01

CS252/Kubitowicz Lec. 4.24

DRAM Performance

- A 60 ns (t_{RAC}) DRAM can
 - perform a row access only every 110 ns (t_{RC})
 - perform column access (t_{CAC}) in 15 ns, but time between column accesses is at least 35 ns (t_{PC}).
 - » In practice, external address delays and turning around buses make it 40 to 50 ns
- These times do not include the time to drive the addresses off the microprocessor nor the memory controller overhead!

1/26/01

CS252/Kubitowicz Lec. 4.29

DRAM History

- DRAMs: capacity +60%/yr, cost -30%/yr
 - 2.5X cells/area, 1.5X die size in -3 years
- '98 DRAM fab line costs \$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
- Order of importance: 1) Cost/bit 2) Capacity
 - First RAMBUS: 10X BW, +30% cost => little impact

1/26/01

CS252/Kubitowicz Lec. 4.28

DRAM Future: 1 Gbit DRAM (ISSCC '96; production '02?)

	Mitsubishi	Samsung
• Blocks	512 x 2 Mbit	1024 x 1 Mbit
• Clock	200 MHz	250 MHz
• Data Pins	64	16
• Die Size	24 x 24 mm	31 x 21 mm
- Sizes will be much smaller in production		
• Metal Layers	3	4
• Technology	0.15 micron	0.16 micron

1/26/01

CS252/Kubitowicz Lec. 4.27

Fast Memory Systems: DRAM specific

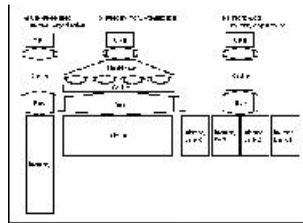
- Multiple CAS accesses: several names (page mode)
 - **Extended Data Out (EDO)**: 30% faster in page mode
- New DRAMs to address gap; what will they cost, will they survive?
 - **RAMBUS**: startup company; reinvent DRAM interface
 - » Each Chip a module vs. slice of memory
 - » Short bus between CPU and chips
 - » Does own refresh
 - » Variable amount of data returned
 - » 1 byte / 2 ns (500 MB/s per chip)
 - » 20% increase in DRAM area
 - **Synchronous DRAM**: 2 banks on chip, a clock signal to DRAM, transfer synchronous to system clock (66 - 150 MHz)
 - Intel claims RAMBUS Direct (16 b wide) is future PC memory?
 - » Possibly not true! Intel to drop RAMBUS?
- Niche memory or main memory?
 - e.g., Video RAM for frame buffers, DRAM + fast serial output

1/26/01

CS252/Kubitowicz Lec. 4.26

Main Memory Organizations

- **Simple**:
 - CPU, Cache, Bus, Memory same width (32 or 64 bits)
- **Wide**:
 - CPU/Mux 1 word; Mux/Cache, Bus, Memory N words (Alpha: 64 bits & 256 bits; UltraSPARC 512)
- **Interleaved**:
 - CPU, Cache, Bus 1 word; Memory N Modules (4 Modules): example is **word interleaved**

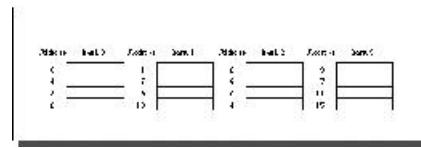


1/26/01

CS252/Kubitowicz Lec. 4.29

Main Memory Performance

- Timing model (word size is 32 bits)
 - 1 to send address,
 - 6 access time, 1 to send data
 - Cache Block is 4 words
- **Simple M.P.** = $4 \times (1+6+1) = 32$
- **Wide M.P.** = $1 + 6 + 1 = 8$
- **Interleaved M.P.** = $1 + 6 + 4 \times 1 = 11$

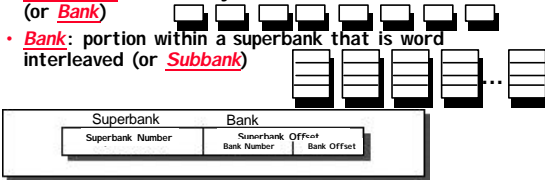


1/26/01

CS252/Kubitowicz Lec. 4.30

Independent Memory Banks

- Memory banks for independent accesses vs. faster sequential accesses
 - Multiprocessor
 - I/O
 - CPU with Hit under n Misses, Non-blocking Cache
- Superbank:** all memory active on one block transfer (or **Bank**)
- Bank:** portion within a superbank that is word interleaved (or **Subbank**)



1/26/01

CS252/Kubitowicz Lec. 4.31

Independent Memory Banks

- How many banks?
 - number banks \neq number clocks to access word in bank
 - For sequential accesses, otherwise will return to original bank before it has next word ready
 - (like in vector case)
- Increasing DRAM => fewer chips => harder to have banks

1/26/01

CS252/Kubitowicz Lec. 4.32

Avoiding Bank Conflicts

- Lots of banks


```
int x[256][512];
for (j = 0; j < 512; j = j+1)
    for (i = 0; i < 256; i = i+1)
        x[i][j] = 2 * x[i][j];
```
- Even with 128 banks, since 512 is multiple of 128, conflict on word accesses
- SW: loop interchange or declaring array not power of 2 ("array padding")
- HW: Prime number of banks
 - bank number = address mod number of banks
 - address within bank = address / number of words in bank
 - modulo & divide per memory access with prime no. banks?
 - address within bank = address mod number words in bank
 - bank number? easy if 2^n words per bank

1/26/01

CS252/Kubitowicz Lec. 4.33

Fast Bank Number

- Chinese Remainder Theorem**
 - As long as two sets of integers a_i and b_i follow these rules

$$b_i = x \text{ mod } a_i, 0 \leq b_i < a_i, 0 \leq x < a_0 \times a_1 \times a_2 \times \dots$$
 - and that a_i and a_j are co-prime if $i \neq j$, then the integer x has only one solution (unambiguous mapping):
 - bank number = b_0 , number of banks = a_0 (= 3 in example)
 - address within bank = b_1 , number of words in bank = a_1 (= 8 in example)
 - N word address 0 to N-1, prime no. banks, words power of 2

Bank Number:	Seq. Interleaved			Modulo Interleaved		
	0	1	2	0	1	2
Address	0	1	2	0	1	2
within Bank:	0	3	6	0	4	8
	1	4	7	1	5	11
	2	5	8	2	6	12
	3	6	9	3	7	13
	4	7	10	4	8	14
	5	8	11	5	9	15
	6	9	12	6	10	16
	7	10	13	7	11	17
	8	11	14	8	12	18
	9	12	15	9	13	19
	10	13	16	10	14	20
	11	14	17	11	15	21
	12	15	18	12	16	22
	13	16	19	13	17	23
	14	17	20	14	18	24
	15	18	21	15	19	25
	16	19	22	16	20	26
	17	20	23	17	21	27
	18	21	24	18	22	28
	19	22	25	19	23	29
	20	23	26	20	24	30
	21	24	27	21	25	31
	22	25	28	22	26	32
	23	26	29	23	27	33
	24	27	30	24	28	34
	25	28	31	25	29	35
	26	29	32	26	30	36
	27	30	33	27	31	37
	28	31	34	28	32	38
	29	32	35	29	33	39
	30	33	36	30	34	40
	31	34	37	31	35	41
	32	35	38	32	36	42
	33	36	39	33	37	43
	34	37	40	34	38	44
	35	38	41	35	39	45
	36	39	42	36	40	46
	37	40	43	37	41	47
	38	41	44	38	42	48
	39	42	45	39	43	49
	40	43	46	40	44	50
	41	44	47	41	45	51
	42	45	48	42	46	52
	43	46	49	43	47	53
	44	47	50	44	48	54
	45	48	51	45	49	55
	46	49	52	46	50	56
	47	50	53	47	51	57
	48	51	54	48	52	58
	49	52	55	49	53	59
	50	53	56	50	54	60
	51	54	57	51	55	61
	52	55	58	52	56	62
	53	56	59	53	57	63
	54	57	60	54	58	64
	55	58	61	55	59	65
	56	59	62	56	60	66
	57	60	63	57	61	67
	58	61	64	58	62	68
	59	62	65	59	63	69
	60	63	66	60	64	70
	61	64	67	61	65	71
	62	65	68	62	66	72
	63	66	69	63	67	73
	64	67	70	64	68	74
	65	68	71	65	69	75
	66	69	72	66	70	76
	67	70	73	67	71	77
	68	71	74	68	72	78
	69	72	75	69	73	79
	70	73	76	70	74	80
	71	74	77	71	75	81
	72	75	78	72	76	82
	73	76	79	73	77	83
	74	77	80	74	78	84
	75	78	81	75	79	85
	76	79	82	76	80	86
	77	80	83	77	81	87
	78	81	84	78	82	88
	79	82	85	79	83	89
	80	83	86	80	84	90
	81	84	87	81	85	91
	82	85	88	82	86	92
	83	86	89	83	87	93
	84	87	90	84	88	94
	85	88	91	85	89	95
	86	89	92	86	90	96
	87	90	93	87	91	97
	88	91	94	88	92	98
	89	92	95	89	93	99
	90	93	96	90	94	100
	91	94	97	91	95	101
	92	95	98	92	96	102
	93	96	99	93	97	103
	94	97	100	94	98	104
	95	98	101	95	99	105
	96	99	102	96	100	106
	97	100	103	97	101	107
	98	101	104	98	102	108
	99	102	105	99	103	109
	100	103	106	100	104	110
	101	104	107	101	105	111
	102	105	108	102	106	112
	103	106	109	103	107	113
	104	107	110	104	108	114
	105	108	111	105	109	115
	106	109	112	106	110	116
	107	110	113	107	111	117
	108	111	114	108	112	118
	109	112	115	109	113	119
	110	113	116	110	114	120
	111	114	117	111	115	121
	112	115	118	112	116	122
	113	116	119	113	117	123
	114	117	120	114	118	124
	115	118	121	115	119	125
	116	119	122	116	120	126
	117	120	123	117	121	127
	118	121	124	118	122	128
	119	122	125	119	123	129
	120	123	126	120	124	130
	121	124	127	121	125	131
	122	125	128	122	126	132
	123	126	129	123	127	133
	124	127	130	124	128	134
	125	128	131	125	129	135
	126	129	132	126	130	136
	127	130	133	127	131	137
	128	131	134	128	132	138
	129	132	135	129	133	139
	130	133	136	130	134	140
	131	134	137	131	135	141
	132	135	138	132	136	142
	133	136	139	133	137	143
	134	137	140	134	138	144
	135	138	141	135	139	145
	136	139	142	136	140	146
	137	140	143	137	141	147
	138	141	144	138	142	148
	139	142	145	139	143	149
	140	143	146	140	144	150
	141	144	147	141	145	151
	142	145	148	142	146	152
	143	146	149	143	147	153
	144	147	150	144	148	154
	145	148	151	145	149	155
	146	149	152	146	150	156
	147	150	153	147	151	157
	148	151	154	148	152	158
	149	152	155	149	153	159
	150	153	156	150	154	160
	151	154	157	151	155	161
	152	155	158	152	156	162
	153	156	159	153	157	163
	154	157	160	154	158	164
	155	158	161	155	159	165
	156	159	162	156	160	166
	157	160	163	157	161	167
	158	161	164	158	162	168
	159	162	165	159	163	169
	160	163	166	160	164	170
	161	164	167	161	165	171
	162	165	168	162	166	172
	163	166	169	163	167	173
	164	167	170	164	168	174
	165	168	171	165	169	175
	166	169	172	166	170	176
	167	170	173	167	171	177
	168	171	174	168	172	178
	169	172	175	169	173	179

1. Fast Hit times via Small and Simple Caches

- Why Alpha 21164 has 8KB Instruction and 8KB data cache + 96KB second level cache?
 - Small data cache and clock rate
- Direct Mapped, on chip

1/26/01 CS252/Kubitowicz Lec. 4.43

2. Fast hits by Avoiding Address Translation

1/26/01 CS252/Kubitowicz Lec. 4.43

2. Fast hits by Avoiding Address Translation

- Send virtual address to cache? Called Virtually Addressed Cache or just Virtual Cache vs. Physical Cache
 - Every time process is switched logically must flush the cache; otherwise get false hits
 - » Cost is time to flush + "compulsory" misses from empty cache
 - Dealing with aliases (sometimes called synonyms):
 - Two different virtual addresses map to same physical address
 - I/O must interact with cache, so need virtual address
- Solution to aliases
 - HW guarantees covers index field & direct mapped, they must be unique; called page coloring
- Solution to cache flush
 - Add process identifier tag that identifies process as well as address within process: can't get a hit if wrong process

1/26/01 CS252/Kubitowicz Lec. 4.43

2. Fast Cache Hits by Avoiding Translation: Process ID impact

- Black is uniprocess
- Light Gray is multiprocess when flush cache
- Dark Gray is multiprocess when use Process ID tag
- Y axis: Miss Rates up to 20%
- X axis: Cache size from 2 KB to 1024 KB

1/26/01 CS252/Kubitowicz Lec. 4.43

2. Fast Cache Hits by Avoiding Translation: Index with Physical Portion of Address

- If index is physical part of address, can start tag access in parallel with translation so that can compare to physical tag

31		Page Address		16		11		Page Offset		0	
31		Page address		16		index		Page offset		0	
31		Address tag		16		Index		Block offset		0	
31		Address Tag		Index		Block Offset		0		0	

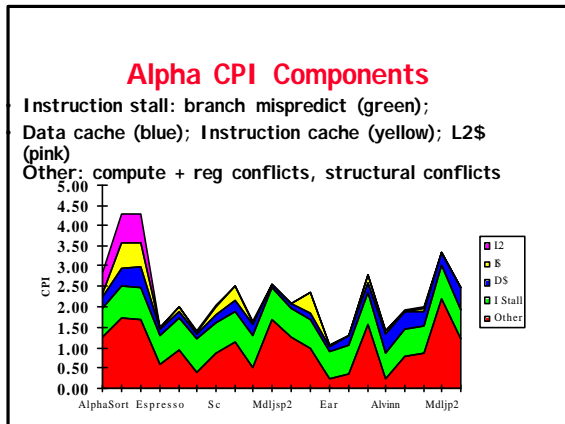
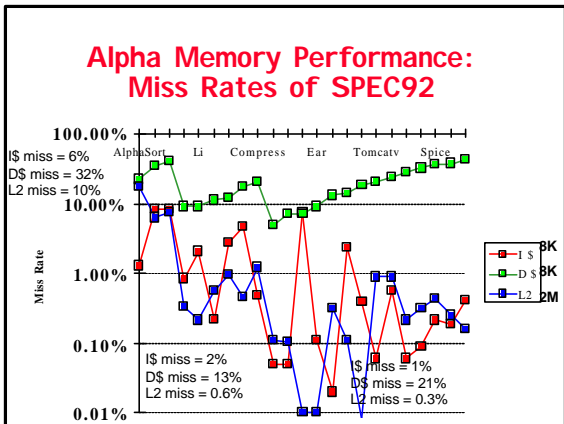
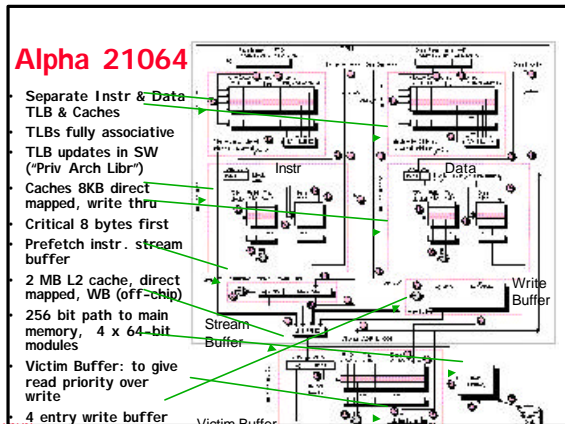
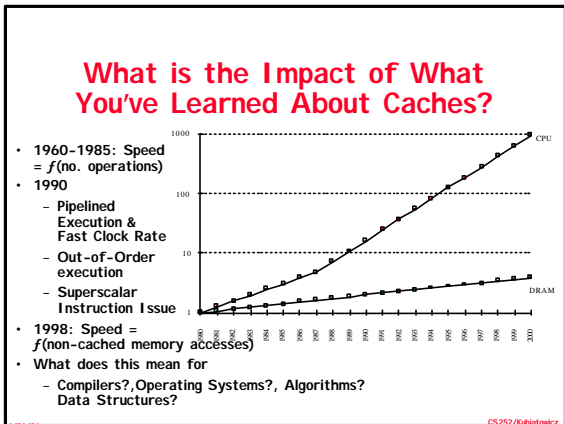
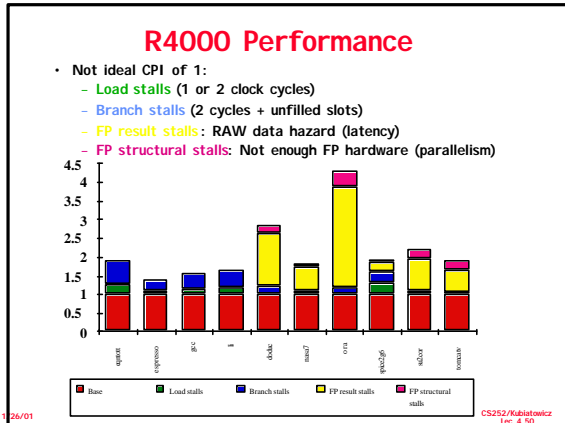
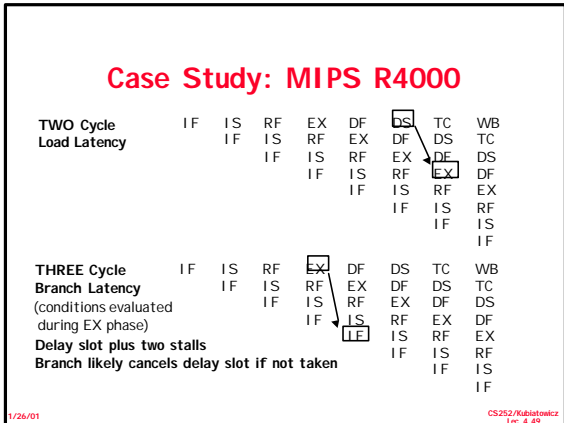
- Limits cache to page size: what if want bigger caches and uses same trick?
 - Higher associativity moves barrier to right
 - Page coloring

1/26/01 CS252/Kubitowicz Lec. 4.47

3: Fast Hits by pipelining Cache Case Study: MIPS R4000

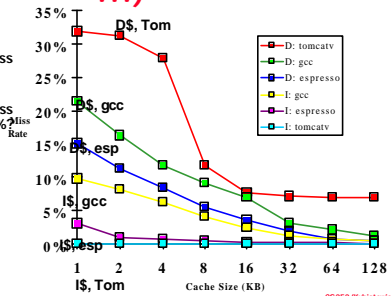
- 8 Stage Pipeline:
 - IF-first half of fetching of instruction; PC selection happens here as well as initiation of instruction cache access.
 - IS-second half of access to instruction cache.
 - RF-instruction decode and register fetch, hazard checking and also instruction cache hit detection.
 - EX-execution, which includes effective address calculation, ALU operation, and branch target computation and condition evaluation.
 - DF-data fetch, first half of access to data cache.
 - DS-second half of access to data cache.
 - TC-tag check, determine whether the data cache access hit.
 - WB-write back for loads and register-register operations.
- What is impact on Load delay?
 - Need 2 instructions between a load and its use!

1/26/01 CS252/Kubitowicz Lec. 4.48



Pitfall: Predicting Cache Performance from Different Prog. (ISA, compiler, ...)

- 4KB Data cache miss rate 8%, 12%, or 28%?
- 1KB Instr cache miss rate 0%, 3%, or 10%?
- Alpha vs. MIPS for 8KB Data \$: 17% vs. 10%
- Why 2X Alpha v. MIPS?



1/26/01

CS252/Kobiatowicz Lec. 4.55

Cache Optimization Summary

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
Priority to Read Misses		+		1
Early Restart & Critical Word 1st		+		2
Non-Blocking Caches		+		3
Second Level Caches		+		2
Better memory system		+		3
Small & Simple Caches	-		+	0
Avoiding Address Translation			+	2
Pipelining Caches			+	2

1/26/01

CS252/Kobiatowicz Lec. 4.56