# CS252
## Graduate Computer Architecture
## Lecture 7

### Cache Design (continued)

Feb 12, 2002
Prof. David Culler

---

# How to Improve Cache Performance?

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

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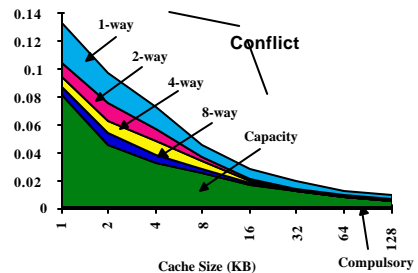3. Reduce the time to hit in the cache.

---

# Where to misses come from?

- **Classifying Misses: 3 Cs**
  - *Compulsory*—The first access to a block is not in the cache, so the block must be brought into the cache. Also called *cold start misses* or *first reference misses*.
    *(Misses in even an Infinite Cache)*
  - *Capacity*—If the cache cannot contain all the blocks needed during execution of a program, capacity misses will occur due to blocks being discarded and later retrieved.
    *(Misses in Fully Associative Size X Cache)*
  - *Conflict*—If block-placement strategy is set associative or direct mapped, conflict misses (in addition to compulsory & capacity misses) will occur because a block can be discarded and later retrieved if too many blocks map to its set. Also called *collision misses* or *interference misses*.
    *(Misses in N-way Associative, Size X Cache)*
- **4th "C":**
  - *Coherence* - Misses caused by cache coherence.

---

# 3Cs Absolute Miss Rate (SPEC92)

---

# Reducing Misses by <u>Hardware</u> Prefetching of Instructions & Data

- **E.g., Instruction Prefetching**
  - Alpha 21064 fetches 2 blocks on a miss
  - Extra block placed in "<u>stream buffer</u>"
  - 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**

---

# Reducing Misses by <u>Software</u> 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
- **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. Faults?
- **Issuing Prefetch Instructions takes time**
  - Is cost of prefetch issues < savings in reduced misses?
  - Higher superscalar reduces difficulty of issue bandwidth

## 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

## Merging Arrays Example

```
/* Before: 2 sequential arrays */
int val[SIZE];
int key[SIZE];

/* After: 1 array of stuctures */
struct merge {
  int val;
  int key;
};
struct merge merged_array[SIZE];
```

Reducing conflicts between val & key; improve spatial locality

## Loop Interchange Example

```
/* Before */
for (k = 0; k < 100; k = k+1)
  for (j = 0; j < 100; j = j+1)
    for (i = 0; i < 5000; i = i+1)
      x[i][j] = 2 * x[i][j];
/* After */
for (k = 0; k < 100; k = k+1)
  for (i = 0; i < 5000; i = i+1)
    for (j = 0; j < 100; j = j+1)
      x[i][j] = 2 * x[i][j];
```

Sequential accesses instead of striding through memory every 100 words; improved spatial locality

## Loop Fusion Example

```
/* Before */
for (i = 0; i < N; i = i+1)
  for (j = 0; j < N; j = j+1)
    a[i][j] = 1/b[i][j] * c[i][j];
for (i = 0; i < N; i = i+1)
  for (j = 0; j < N; j = j+1)
    d[i][j] = a[i][j] + c[i][j];
/* After */
for (i = 0; i < N; i = i+1)
  for (j = 0; j < N; j = j+1)
  {   a[i][j] = 1/b[i][j] * c[i][j];
      d[i][j] = a[i][j] + c[i][j];}
```

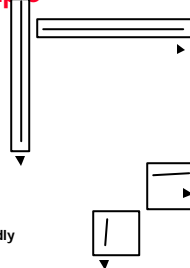2 misses per access to a & c vs. one miss per access; improve spatial locality

## Blocking Example

```
/* Before */
for (i = 0; i < N; i = i+1)
  for (j = 0; j < N; j = j+1)
    {r = 0;
     for (k = 0; k < N; k = k+1){
       r = r + y[i][k]*z[k][j];};
     x[i][j] = r;
    };
```

- Two Inner Loops:
  - Read all NxN elements of z[]
  - Read N elements of 1 row of y[] repeatedly
  - Write N elements of 1 row of x[]
- Capacity Misses a function of N & Cache Size:
  - $2N^3 + N^2$ => (assuming no conflict; otherwise …)
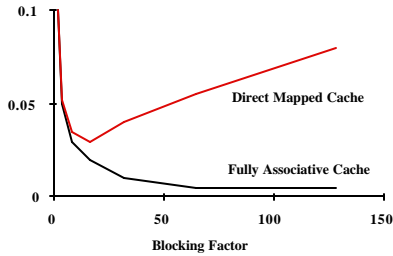- Idea: compute on BxB submatrix that fits

## Blocking Example

```
/* After */
for (jj = 0; jj < N; jj = jj+B)
for (kk = 0; kk < N; kk = kk+B)
for (i = 0; i < N; i = i+1)
  for (j = jj; j < min(jj+B-1,N); j = j+1)
    {r = 0;
     for (k = kk; k < min(kk+B-1,N); k = k+1) {
       r = r + y[i][k]*z[k][j];};
     x[i][j] = x[i][j] + r;
    };
```

- B called *Blocking Factor*
- Capacity Misses from $2N^3 + N^2$ to $N^3/B+2N^2$
- Conflict Misses Too?
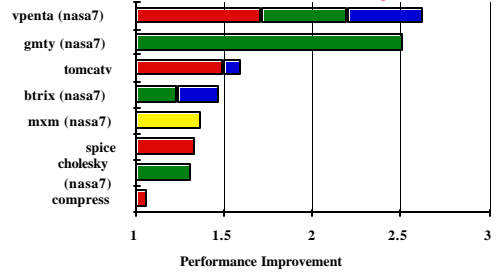
## Reducing Conflict Misses by Blocking



- **Conflict misses in caches not FA vs. Blocking size**
  - Lam et al [1991] a blocking factor of 24 had a fifth the misses vs. 48 despite both fit in cache

---

## Summary of Compiler Optimizations to Reduce Cache Misses (by hand)

---

## Summary: Miss Rate Reduction

$$CPUtime = IC \times \left( CPI_{Execution} + \frac{Memory\ accesses}{Instruction} \times Miss\text{-}rate \times Miss\ penalty \right) \times Clock\ cycle\ time$$

- **3 Cs: Compulsory, Capacity, Conflict**
  - 0. Larger cache
  - 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!

---

## Review: Improving Cache Performance

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

---

## 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?

---

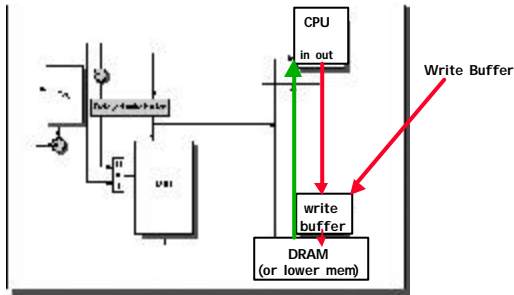## Write Policy 2: Write Allocate vs Non-Allocate (What happens on write-miss)

- **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. Reducing Miss Penalty: Read Priority over Write on Miss



CPU

in out

**Write Buffer**

write buffer

DRAM (or lower mem)

---

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

- **Write-through w/ write buffers => 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 want buffer to hold displaced 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

---

## 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 executio
  - *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 => tend to want next sequential word, so not clear if benefit by early restart**
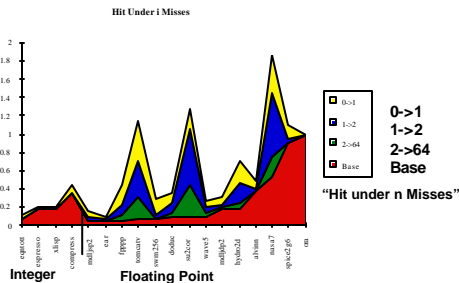
**block**

---

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

- *Non-blocking cache* or *lockup-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 muliple memory banks (otherwise cannot support)
  - Penium Pro allows 4 outstanding memory misses

---

## Value of Hit Under Miss for SPEC

**Hit Under i Misses**



0->1
1->2
2->64
Base

"Hit under n Misses"

Integer   Floating Point

- 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

---

## 4: Add a second-level cache

- **L2 Equations**

  $AMAT = Hit\ Time_{L1} + Miss\ Rate_{L1}\ x\ Miss\ Penalty_{L1}$

  $Miss\ Penalty_{L1} = Hit\ Time_{L2} + Miss\ Rate_{L2}\ x\ Miss\ Penalty_{L2}$

  $AMAT = Hit\ Time_{L1} +$
  $Miss\ Rate_{L1}\ x\ (Hit\ Time_{L2} + Miss\ Rate_{L2} + Miss\ Penalty_{L2})$

- **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*
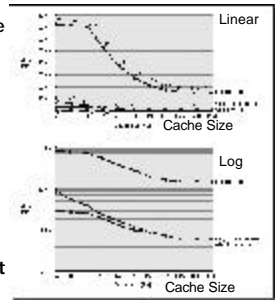
---

## Partner Discussion

What's different in L2 vs L1 Caches?

---

## 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 >> 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**



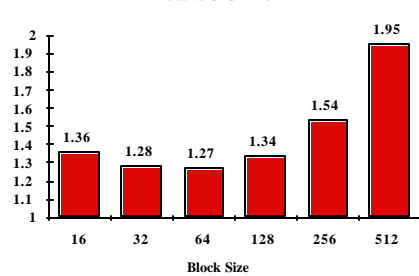Linear — Cache Size

Log — Cache Size

---

## 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**

---

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

Relative CPU Time



| Block Size | Value |
|---|---|
| 16 | 1.36 |
| 32 | 1.28 |
| 64 | 1.27 |
| 128 | 1.34 |
| 256 | 1.54 |
| 512 | 1.95 |

- **32KB L1, 8 byte path to memory**

---

## Reducing Miss Penalty Summary

$$CPUtime = IC \times \left( CPI_{Execution} + \frac{Memory\ accesses}{Instruction} \times \textbf{Miss rate} \times \textbf{Miss\_penalty} \right) \times Clock\ cycle\ time$$
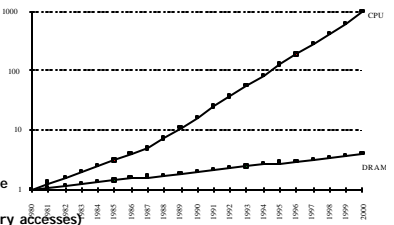
- **Four techniques**
  - **Read priority over write on miss**
  - **Early Restart and Critical Word First on miss**
  - **Non-blocking Caches (Hit under Miss, Miss under Miss)**
  - **Second Level Cache**
- **Can be applied recursively to Multilevel Caches**
  - **Danger is that time to DRAM will grow with multiple levels in between**
  - **First attempts at L2 caches can make things worse, since increased worst case is worse**

---

## What is the Impact of What You've Learned About Caches?

- **1960-1985: Speed = ƒ(no. operations)**
- **1990**
  - **Pipelined Execution & Fast Clock Rate**
  - **Out-of-Order execution**
  - **Superscalar Instruction Issue**
- **1998: Speed = ƒ(non-cached memory accesses)**
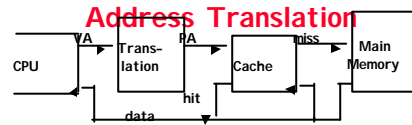- **Superscalar, Out-of-Order machines hide L1 data cache miss (~5 clocks) but not L2 cache miss (~50 clocks)?**



CPU

DRAM

## 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

---

## Address Translation



- Page table is a large data structure in memory
- Two memory accesses for every load, store, or instruction fetch!!!
- Virtually addressed cache?
  - synonym problem
- Cache the address translations?

---

## TLBs

A way to speed up translation is to use a special cache of recently used page table entries -- this has many names, but the most frequently used is *Translation Lookaside Buffer* or *TLB*

| Virtual Address | Physical Address | Dirty | Ref | Valid | Access |
|---|---|---|---|---|---|
| | | | | | |

Really just a cache on the page table mappings

TLB access time comparable to cache access time
(much less than main memory access time)

---

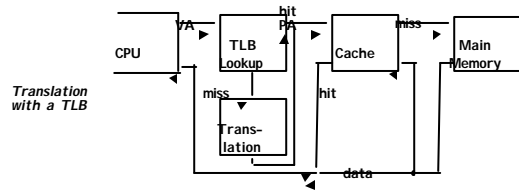## Translation Look-Aside Buffers

Just like any other cache, the TLB can be organized as fully associative, set associative, or direct mapped

TLBs are usually small, typically not more than 128 - 256 entries even on high end machines. This permits fully associative lookup on these machines. Most mid-range machines use small n-way set associative organizations.



*Translation with a TLB*

1/2 t     t     20

---

## 2. Fast hits by Avoiding Address Translation



Conventional Organization

Virtually Addressed Cache Translate only on miss Synonym Problem

Overlap $ access with VA translation: requires $ index to remain invariant across translation

---

## 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



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

---

## 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
    - » Add *process identifier tag* that identifies process as well as address within process: can't get a hit if wrong process
- **Dealing with *aliases* (sometimes called *synonyms*); Two different virtual addresses map to same physical address**
  - solve by fiat: noaliasing! What are the implications?
  - HW antialiasing: guarantees every cache block has unique address
    - » verify on miss (rather than on every hit)
    - » cache set size <= page size ?
    - » what if it gets larger?
  - How can SW simplify the problem? (called *page coloring)*
  - I/O must interact with cache, so need virtual address
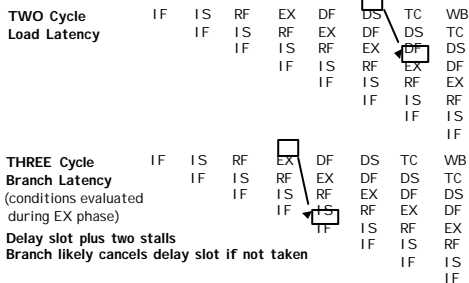
## 3: Fast Hits by pipelining Cache Case Study: MIPS R4000

- **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!**

## Case Study: MIPS R4000

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| **TWO Cycle** | IF | IS | RF | EX | DF | DS | TC | WB |
| **Load Latency** | | IF | IS | RF | EX | DF | DS | TC |
| | | | IF | IS | RF | EX | DF | DS |
| | | | | IF | IS | RF | EX | DF |
| | | | | | IF | IS | RF | EX |
| | | | | | | IF | IS | RF |
| | | | | | | | IF | IS |
| | | | | | | | | IF |

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| **THREE Cycle** | IF | IS | RF | EX | DF | DS | TC | WB |
| **Branch Latency** | | IF | IS | RF | EX | DF | DS | TC |
| (conditions evaluated | | | IF | IS | RF | EX | DF | DS |
| during EX phase) | | | | IF | IS | RF | EX | DF |
| | | | | | IF | IS | RF | EX |
| | | | | | | IF | IS | RF |
| | | | | | | | IF | IS |
| | | | | | | | | IF |

**Delay slot plus two stalls
Branch likely cancels delay slot if not taken**
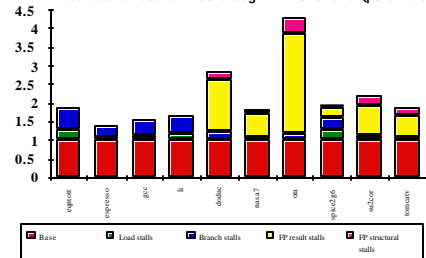
## R4000 Performance

- **Not ideal CPI of 1:**
  - **Load stalls (1 or 2 clock cycles)**
  - **Branch stalls (2 cycles + unfilled slots)**
  - **FP result stalls: RAW data hazard (latency)**
  - **FP structural stalls: Not enough FP hardware (parallelism)**



Legend: Base | Load stalls | Branch stalls | FP result stalls | FP structural stalls

## What is the Impact of What You've Learned About Caches?

- 1960–1985: Speed = ƒ(no. operations)
- 1990
  - Pipelined Execution & Fast Clock Rate
  - Out–of–Order execution
  - Superscalar Instruction Issue
- 1998: Speed = ƒ(non–cached memory accesses)
- What does this mean for
  - Compilers?,Operating Systems?, Algorithms? Data Structures?

## Alpha 21064

- **Separate Instr & Data TLB & Caches**
- **TLBs fully associative**
- **TLB updates in SW ("Priv Arch Libr")**
- **Caches 8KB direct mapped, write thru**
- **Critical 8 bytes first**
- **Prefetch instr. stream buffer**
- **2 MB L2 cache, direct mapped, WB (off–chip)**
- **256 bit path to main memory, 4 x 64–bit modules**
- **Victim Buffer: to give read priority over write**
- **4 entry write buffer between D$ & L2$**



Instr | Data | Write Buffer | Stream Buffer | Victim Buffer

Page 7

## Alpha Memory Performance: Miss Rates of SPEC92

I$ miss = 6%
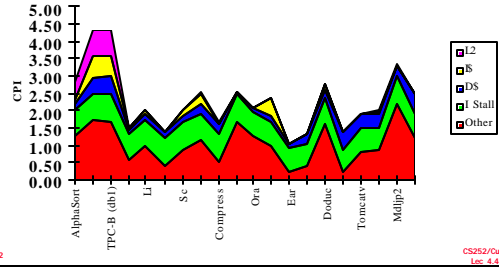D$ miss = 32%
L2 miss = 10%

Miss Rate

100.00%
10.00%
1.00%
0.10%
0.01%

I$ miss = 2%
D$ miss = 13%
L2 miss = 0.6%

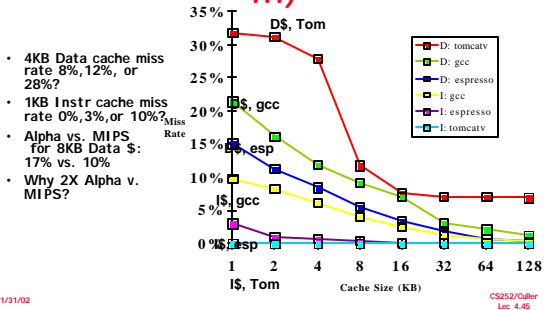I$ miss = 1%
D$ miss = 21%
L2 miss = 0.3%

I $ 8K
D $ 8K
L2 2M

## Alpha CPI Components

- Instruction stall: branch mispredict (green);
- Data cache (blue); Instruction cache (yellow); L2$ (pink)
  Other: compute + reg conflicts, structural conflicts

CPI

5.00
4.50
4.00
3.50
3.00
2.50
2.00
1.50
1.00
0.50
0.00

AlphaSort, TPC-B (db1), Li, Sc, Compress, Ora, Ear, Doduc, Tomcatv, Mdljp2

L2
I$
D$
I Stall
Other

## 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?

Miss Rate

35%
30%
25%
20%
15%
10%
5%
0%

D$, Tom
D$, gcc
D$, esp
I$, gcc
I$, esp
I$, Tom

D: tomcatv
D: gcc
D: espresso
I: gcc
I: espresso
I: tomcatv

1  2  4  8  16  32  64  128
Cache Size (KB)

## Cache Optimization Summary

| | Technique | MR | MP | HT | Complexity |
|---|---|---|---|---|---|
| miss rate | 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 |
| miss penalty | Priority to Read Misses | | + | | 1 |
| | Early Restart & Critical Word 1st | | + | | 2 |
| | Non-Blocking Caches | | + | | 3 |
| | Second Level Caches | | + | | 2 |
| | Better memory system | | + | | 3 |
| hit time | Small & Simple Caches | – | | + | 0 |
| | Avoiding Address Translation | | | + | 2 |
| | Pipelining Caches | | | + | 2 |