

4 November

- 8 classes to go!
- Read 7.3-7.5
- Section 7.5 especially important!
- New Assignment on the web

11/4/2004

Comp 120 Fall 2004

1

Direct-Mapping Example

- With 8 byte BLOCKS, the bottom 3 bits determine the byte in the BLOCK
- With 4 cache BLOCKS, the next 2 bits determine which BLOCK to use

1024d = 10000000000b → line = 00b = 0d
 1000d = 01111101000b → line = 01b = 1d
 1040d = 1000010000b → line = 10b = 2d

Tag	Data	
1024	44	99
1000	17	23
1040	1	4
1016	29	38

Memory	
1000	17
1004	23
1008	11
1012	5
1016	29
1020	38
1024	44
1028	99
1032	97
1036	25
1040	1
1044	4

11/4/2004

Comp 120 Fall 2004

2

Direct Mapping Miss

- What happens when we now ask for address 1008?

1008d = 01111110000b → line = 10b = 2d

but earlier we put 1040d there...

1040d = 1000010000b → line = 10b = 2d

Tag	Data	
1024	44	99
1000	17	23
1008	11	5
1016	29	38

Memory	
1000	17
1004	23
1008	11
1012	5
1016	29
1020	38
1024	44
1028	99
1032	97
1036	25
1040	1
1044	4

11/4/2004

Comp 120 Fall 2004

3

Miss Penalty and Rate

- The *MISS PENALTY* is the time it takes to read the memory if it isn't in the cache
 - 50 to 100 cycles is common.
- The *MISS RATE* is the fraction of accesses which MISS
- The *HIT RATE* is the fraction of accesses which HIT
- $\text{MISS RATE} + \text{HIT RATE} = 1$

Suppose a particular cache has a *MISS PENALTY* of 100 cycles and a *HIT RATE* of 95%. The *CPI* for load is normally 5 but on a miss it is 105. What is the average *CPI* for load?

$$\text{Average CPI} = 10 \quad 5 * 0.95 + 105 * 0.05 = 10$$

Suppose *MISS PENALTY* = 120 cycles?

then *CPI* = 11 (slower memory doesn't hurt much)

11/4/2004

Comp 120 Fall 2004

4

Some Associativity can help

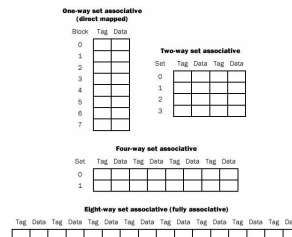
- Direct-Mapped caches are very common but can cause problems...
- SET ASSOCIATIVITY can help.
- Multiple Direct-mapped caches, then compare multiple TAGS
 - 2-way set associative = 2 direct mapped + 2 TAG comparisons
 - 4-way set associative = 4 direct mapped + 4 TAG comparisons
- Now array size == power of 2 doesn't get us in trouble
- But
 - slower
 - less memory in same area
 - maybe direct mapped wins...

11/4/2004

Comp 120 Fall 2004

5

Associative Cache



11/4/2004

Comp 120 Fall 2004

6

What about store?

- What happens in the cache on a store?
 - WRITE BACK CACHE → put it in the cache, write on replacement
 - WRITE THROUGH CACHE → put in cache and in memory
- What happens on store and a MISS?
 - WRITE BACK will fetch the line into cache
 - WRITE THROUGH might just put it in memory

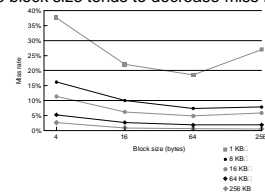
11/4/2004

Comp 120 Fall 2004

7

Cache Block Size and Hit Rate

- Increasing the block size tends to decrease miss rate:



- Use split caches because there is more spatial locality in code:

Program	Block size in words	Instruction miss rate	Data miss rate	Effective combined miss rate
gcc	1	6.1%	2.1%	5.4%
	4	2.0%	1.7%	1.9%
splice	1	1.2%	1.3%	1.2%
	4	0.3%	0.6%	0.4%

11/4/2004

Comp 120 Fall 2004

8

Cache Performance

- Simplified model:

$$\text{execution time} = (\text{execution cycles} + \text{stall cycles}) \times \text{cycle time}$$

$$\text{stall cycles} = \# \text{ of instructions} \times \text{miss ratio} \times \text{miss penalty}$$

- Two ways of improving performance:
 - decreasing the miss ratio
 - decreasing the miss penalty

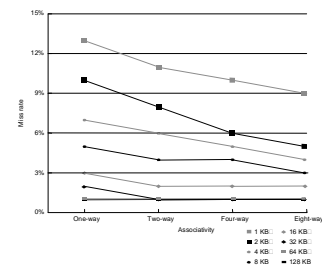
What happens if we increase block size?

11/4/2004

Comp 120 Fall 2004

9

Associative Performance



11/4/2004

Comp 120 Fall 2004

10

Multilevel Caches

- We can reduce the miss penalty with a 2nd level cache
- Add a second level cache:
 - often primary cache is on the same chip as the processor
 - use SRAMs to add another cache above primary memory (DRAM)
 - miss penalty goes down if data is in 2nd level cache
- Example:
 - Base CPI=1.0 on a 500Mhz machine with a 5% miss rate, 200ns DRAM access
 - Adding 2nd level cache with 20ns access time decreases miss rate to 2%
- Using multilevel caches:
 - try and optimize the hit time on the 1st level cache
 - try and optimize the miss rate on the 2nd level cache

11/4/2004

Comp 120 Fall 2004

11

Matrix Multiply

- A VERY common operation in scientific programs
- Multiply a LxM matrix by an MxN matrix to get an LxN matrix result
- This requires L*N inner products each requiring M * and +
- So 2*L*M*N floating point operations
- Definitely a FLOATING POINT INTENSIVE application
- L=M=N=100, 2 Million floating point operations

11/4/2004

Comp 120 Fall 2004

12

Matrix Multiply

```

const int L = 2;
const int M = 3;
const int N = 4;
void mm(double A[L][M], double B[M][N], double C[L][N])
{
    for(int i=0; i<L; i++)
        for(int j=0; j<N; j++) {
            double sum = 0.0;
            for(int k=0; k<M; k++)
                sum = sum + A[i][k] * B[k][j];
            C[i][j] = sum;
        }
}

```

11/4/2004

Comp 120 Fall 2004

13

Matrix Memory Layout

Our memory is a 1D array of bytes
How can we put a 2D thing in a 1D memory?

double A[2][3]:

00	01	02
10	11	12

Row Major

00
01
02
10
11
12

addr = base+(i*3+j)*8

Column Major

00
10
01
11
02
12

addr = base + (i + j*2)*8

11/4/2004

Comp 120 Fall 2004

14

Where does the time go?

The inner loop takes all the time

```

for(int k=0; k<M; k++)
    sum = sum + A[i][k] * B[k][j];

```

```

L1: mul $t1, i, M
    add $t1, $t1, k
    mul $t1, $t1, 8
    add $t1, $t1, A
    l.d $f1, 0($t1)
    mul $t2, k, N
    add $t2, $t2, j
    mul $t2, $t2, 8
    add $t2, $t2, B
    l.d $f2, 0($t2)
    mul.d $f3, $f1, $f2
    add.d $f4, $f4, $f3
    add k, k, 1
    slt $t0, k, M
    bne $t0, $zero, L1

```

11/4/2004

Comp 120 Fall 2004

15

Change Index * to +

The inner loop takes all the time

```

for(int k=0; k<M; k++)
    sum = sum + A[i][k] * B[k][j];

```

```

L1: l.d $f1, 0($t1)
    add $t1, $t1, AColStep
    l.d $f2, 0($t2)
    add $t2, $t2, BRowStep
    mul.d $f3, $f1, $f2
    add.d $f4, $f4, $f3
    add k, k, 1
    slt $t0, k, M
    bne $t0, $zero, L1

```

AColStep = 8
BRowStep = 8 * N

11/4/2004

Comp 120 Fall 2004

16

Eliminate k, use an address instead

The inner loop takes all the time

```

for(int k=0; k<M; k++)
    sum = sum + A[i][k] * B[k][j];

```

```

L1: l.d $f1, 0($t1)
    add $t1, $t1, AColStep
    l.d $f2, 0($t2)
    add $t2, $t2, BRowStep
    mul.d $f3, $f1, $f2
    add.d $f4, $f4, $f3
    bne $t1, LastA, L1

```

11/4/2004

Comp 120 Fall 2004

17

We made it faster

The inner loop takes all the time

```

for(int k=0; k<M; k++)
    sum = sum + A[i][k] * B[k][j];

```

```

L1: l.d $f1, 0($t1)
    add $t1, $t1, AColStep
    l.d $f2, 0($t2)
    add $t2, $t2, BRowStep
    mul.d $f3, $f1, $f2
    add.d $f4, $f4, $f3
    bne $t1, LastA, L1

```

Now this is FAST! Only 7 instructions in the inner loop!

BUT...

When we try it on big matrices it slows way down.

Whas Up?

11/4/2004

Comp 120 Fall 2004

18

Now where is the time?

The inner loop takes all the time

```
for(int k=0; k<M; k++)
    sum = sum + A[i][k] * B[k][j];
```

```
L1: l.d $f1, 0($t1)
    add $t1, $t1, AColStep
    l.d $f2, 0($t2)           lots of time wasted here!
    add $t2, $t2, BRowStep

    mul.d $f3, $f1, $f2
    add.d $f4, $f4, $f3       possibly a little stall right here
    bne $t1, LastA, L1
```

11/4/2004

Comp 120 Fall 2004

19

Why?

The inner loop takes all the time

```
for(int k=0; k<M; k++)
    sum = sum + A[i][k] * B[k][j];
```

```
L1: l.d $f1, 0($t1)           This load usually hits (maybe 3 of 4)
    add $t1, $t1, AColStep
    l.d $f2, 0($t2)           This load always misses!
    add $t2, $t2, BRowStep

    mul.d $f3, $f1, $f2
    add.d $f4, $f4, $f3
    bne $t1, LastA, L1
```

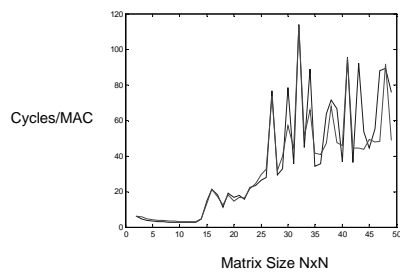
11/4/2004

Comp 120 Fall 2004

20

Matrix Multiply Simulation

Simulation of 2k direct-mapped cache with 32 and 16 byte blocks



11/4/2004

Comp 120 Fall 2004

21

classes to go

7

- Read 7.3-7.5
- Section 7.5 especially important!

11/4/2004

Comp 120 Fall 2004

22