Memory Hierarchy (IV): Programming Techniques to Cache Performance & Basic Pipelined Processor Design

Hung-Wei Tseng

Performance gap between Processor/Memory





Which of the following schemes can help Athlon 64?

- How many of the following schemes mentioned in "improving direct-mapped cache performance by the addition of a small fully-associative cache and prefetch buffers" would help AMD Phenom II for the code in the previous slide?
 - Missing cache help improving conflict misses
 - Victim cache help improving conflict misses
 - ③ Prefetch
 - can potentially hurt
 - Stream buffer only help improving compulsory misses
 - A. 0 B. 1
 - C. 2



E. 4

Multibanks & non-blocking caches





return block 0xDEAEBE



assume the bus between L1/L2 only allows a quarter of the cache block go through it

Write buffer!



assume the bus between L1/L2 only allows a quarter of the cache block go through it

Summary of Optimizations

- Hardware
 - Prefetch compulsory miss
 - Write buffer miss penalty
 - Bank/pipeline miss penalty
 - Critical word first and early restart miss panelty



The result of sizeof(struct student)

• Consider the following data structure: struct student { int id; double *homework; int participation; double midterm; double average; };



64-bit

What's the output of

printf("%lu\n",sizeof(struct student))?

- A. 20
- B. 28
- C. 32
- D. 36





Memory addressing/alignment

- Almost every popular ISA architecture uses "byte-addressing" to access memory locations
- Instructions generally work faster when the given memory address is aligned
 - Aligned if an instruction accesses an object of size n at address X, the access is **aligned** if $X \mod n = 0$.
 - Some architecture/processor does not support aligned access at all
 - Therefore, compilers only allocate objects on "aligned" address
- Compiler optimization cannot help!



Team scores



4.5	5.5	4.5



- Programmer's optimizations for cache performance
- Basic Pipelined Processor Design

Programming and memory performance

Array of structures or structure of arrays

		Array of objects							object
<pre>struct grades { int id; double *homework; double average; };</pre>					<pre>struct grades { int *id; double **homework; double *average; };</pre>				
								ID	ID
ID	*hom	ework	average	ID	*homework	average		homework	homework
								average	average
average of homew	<pre>ge of each mework for(i=0;i<homework_items; (double)total_number_students;="" +="gradesheet[j].homework[i];" =="" for(j="0;j<total_number_students;j++)" gradesheet[total_number_students].homework[i]="" i++)="" pre="" {="" }<=""></homework_items;></pre>		<pre>for(i = { grades for(j { gradeshe } gr total_nu }</pre>	0;i < homewo heet.homewor = 0; j <tota adesheet.hom et.homework[adesheet.hom mber_student</tota 	rk_items; i+ k[i][total_n l_number_stu ework[i][tot i][j]; ework[i][tot s;				





Poll close in 1:30

What data structure is performing better

	Array of objects	object
	<pre>struct grades { int id; double *homework; double average; };</pre>	<pre>struct grades { int *id; double **homework; double *average; };</pre>
average of each homework	<pre>for(i=0;i<homework_items; (double)total_number_students;="" +="gradesheet[j].homework[i];" =="" for(j="0;j<total_number_students;j++)" gradesheet[total_number_students].homework[i]="" i++)="" pre="" {="" }<=""></homework_items;></pre>	<pre>for(i = 0;i < homework_items; i+ { gradesheet.homework[i][total_n for(j = 0; j <total_number_stu gradesheet.homework[i][j];="" gradesheet.homework[i][tot="" pre="" total_number_students;="" {="" }="" }<=""></total_number_stu></pre>

- Considering your workload would like to calculate the average score of one of the homework for all students, which data structure would deliver better performance?
 - A. Array of objects
 - B. Object of arrays

of arrays

+)

umber_students] = 0.0; dents; j++)

al_number_students] +=

al_number_students] /=

Poll close in 1:30

What data structure is performing bet

	Array of objects	object
	<pre>struct grades { int id; double *homework; double average; };</pre>	<pre>struct grades { int *id; double **homework; double *average; };</pre>
average of each homework	<pre>for(i=0;i<homework_items; (double)total_number_students;="" +="gradesheet[j].homework[i];" =="" for(j="0;j<total_number_students;j++)" gradesheet[total_number_students].homework[i]="" i++)="" pre="" {="" }<=""></homework_items;></pre>	<pre>for(i = 0;i < homework_items; i+ { gradesheet.homework[i][total_n for(j = 0; j <total_number_stu gradesheet.homework[i][j];="" gradesheet.homework[i][tot="" pre="" total_number_students;="" {="" }="" }<=""></total_number_stu></pre>

- Considering your workload would like to calculate the average score of one of the homework for all students, which data structure would deliver better performance?
 - A. Array of objects
 - B. Object of arrays



of arrays

+)

umber_students] = 0.0; dents; j++)

cal_number_students] +=

cal_number_students] /=

What data structure is performing better

	Array of objects	object
	<pre>struct grades { int id; double *homework; double average; };</pre>	<pre>struct grades { int *id; double **homework; double *average; };</pre>
average of each homework	<pre>for(i=0;i<homework_items; (double)total_number_students;="" +="gradesheet[j].homework[i];" =="" for(j="0;j<total_number_students;j++)" gradesheet[total_number_students].homework[i]="" i++)="" pre="" {="" }<=""></homework_items;></pre>	<pre>for(i = 0;i < homework_items; i+ { gradesheet.homework[i][total_n for(j = 0; j <total_number_stu gradesheet.homework[i][j];="" gradesheet.homework[i][tot="" pre="" total_number_students;="" {="" }="" }<=""></total_number_stu></pre>

- Considering your workload would like to calculate the average score of one of the homework for all students, which data structure would deliver better performance? What if we want to calculate average scores for each student?
 - A. Array of objects

B. Object of arrays

of arrays

+)

umber_students] = 0.0; dents; j++)

al_number_students] +=

al_number_students] /=

Column-store or row-store

• If you're designing an in-memory database system, will you be using

Rowld	Empld	Lastname	Firstname	Salary
1	10	Smith	Joe	40000
2	12	Jones	Mary	50000
3	11	Johnson	Cathy	44000
4	22	Jones	Bob	55000

• column-store — stores data tables column by column

10:001,12:002,11:003,22:004; Smith:001, Jones:002, Johnson:003, Jones:004 select Lastname, Firstname from table Joe:001, Mary:002, Cathy:003, Bob:004; 40000:001,50000:002,44000:003,55000:004;

row-store — stores data tables row by row

001:10,Smith,Joe,40000; 002:12, Jones, Mary, 50000; 003:11, Johnson, Cathy, 44000; 004:22, Jones, Bob, 55000;



if the most frequently used query looks like -

Loop interchange/fission/fusion

Loop interchange

```
for(i = 0; i < ARRAY_SIZE; i++)
{
    for(j = 0; j < ARRAY_SIZE; j++)
    {
        c[i][j] = a[i][j]+b[i][j];
    }
}</pre>
```



$O(n^2)$	Complexity
Same	Instruction Count?
Same	Clock Rate
Better	CPI

for(j = 0; j < ARRAY_SIZE; j++) { for(i = 0; i < ARRAY_SIZE; i++) { c[i][j] = a[i][j]+b[i][j]; } }</pre>

 $O(n^2)$







AMD Phenom II

- D-L1 Cache configuration of AMD Phenom II
 - Size 64KB, 2-way set associativity, 64B block, LRU policy, write-allocate, write-back, and assuming 64-bit address.

```
int a[16384], b[16384], c[16384];
/* c = 0x10000, a = 0x20000, b = 0x30000 */
for(i = 0; i < 512; i++) {
    c[i] = a[i] + b[i];
   //load a, b, and then store to c
}
```

What's the data cache miss rate for this code?

A. 6.25% C = ABSB. 56.25% 64KB = 2 * 64 * SC. 66.67% S = 512offset = lg(64) = 6 bits 68.75% D. index = lg(512) = 9 bits E. 100% tag = 64 - lg(512) - lg(64) = 49 bits

Poll close in 1:30

What if the code look like this?

- D-L1 Cache configuration of AMD Phenom II
 - Size 64KB, 2-way set associativity, 64B block, LRU policy, write-allocate, write-back, and assuming 64-bit address.

```
int a[16384], b[16384], c[16384];
/* c = 0x10000, a = 0x20000, b = 0x30000 */
for(i = 0; i < 512; i++)
    c[i] = a[i]; //load a and then store to c
for(i = 0; i < 512; i++)
    c[i] += b[i]; //load b, load c, add, and then store to c
```

What's the data cache miss rate for this code?

- A. 6.25%
- B. 56.25%
- C. 66.67%
- 68.75% D.
- E. 100%



Poll close in 1:30

What if the code look like this?

- D-L1 Cache configuration of AMD Phenom II
 - Size 64KB, 2-way set associativity, 64B block, LRU policy, write-allocate, write-back, and assuming 64-bit address.

```
int a[16384], b[16384], c[16384];
/* c = 0x10000, a = 0x20000, b = 0x30000 */
for(i = 0; i < 512; i++)
    c[i] = a[i]; //load a and then store to c
for(i = 0; i < 512; i++)
    c[i] += b[i]; //load b, load c, add, and then store to c
```

What's the data cache miss rate for this code?

- A. 6.25%
- 56.25% B.
- C. 66.67%
- 68.75% D.
- E. 100%

Loop fission







What if the code look like this?

- D-L1 Cache configuration of AMD Phenom II
 - Size 64KB, 2-way set associativity, 64B block, LRU policy, write-allocate, write-back, and assuming 64-bit address.

```
int a[16384], b[16384], c[16384];
/* c = 0x10000, a = 0x20000, b = 0x30000 */
for(i = 0; i < 512; i++)
    c[i] = a[i]; //load a and then store to c
for(i = 0; i < 512; i++)
    c[i] += b[i]; //load b, load c, add, and then store to c
```

What's the data cache miss rate for this code?



E. 100%



Poll close in 1:30

What if the code look like this? — intel



for (i = 0; i < N; i = i+1) ${}$ }

- D-L1 Cache configuration of Intel Core i7
 - Size 32KB, 8-way set associativity, 64B block, LRU policy, writeallocate, write-back, and assuming 64-bit address.

Which version of code will perform better?

- A. Version A
- B. Version B
- C. They're about the same



- 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];

Poll close in 1:30

What if the code look like this? — in



for (i = 0; i < N; i = i+1)}

- D-L1 Cache configuration of Intel Core i7
 - Size 32KB, 8-way set associativity, 64B block, LRU policy, writeallocate, write-back, and assuming 64-bit address.

Which version of code will perform better?

- A. Version A
- B. Version B
- C. They're about the same



- 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];

Loop Fusion

```
/* 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

What if the code look like this? — intel

- D-L1 Cache configuration of Intel Core i7
 - Size 32KB, 8-way set associativity, 64B block, LRU policy, writeallocate, write-back, and assuming 64-bit address.

Which version of code will perform better?

Version A A.

Version B

C. They're about the same



- i < N; i = i+1) = 0; j < N; j = j+1)
- **i][j]** = 1/b[i][j] * **c[i][j]**; i][j] = a[i][j] + c[i][j];

Loop fusion

Blocking

Case study: Matrix Multiplication

```
for(i = 0; i < ARRAY_SIZE; i++) {</pre>
                                        Algorithm class tells you it's O(n<sup>3</sup>)
 for(j = 0; j < ARRAY_SIZE; j++) {</pre>
   for(k = 0; k < ARRAY_SIZE; k++) {</pre>
                                          If n=1024, it takes about 1 sec
     c[i][j] += a[i][k]*b[k][j];
   }
 }
                                   How long is it take when n=2048?
}
```



Matrix Multiplication



- If each dimension of your matrix is 2048
 - Each row takes 2048*8 bytes = 16KB
 - The L1 \$ of intel Core i7 is 32KB, 8-way, 64-byte blocked
 - You can only hold at most 2 rows/columns of each matrix!
 - You need the same row when j increase!

Very likely a miss if array is large



b

Block algorithm for matrix multiplication



You only need to hold these sub-matrices in your cache

for(kk = k; kk < k+(ARRAY_SIZE/n); kk++)</pre>



How do you know it's better?

- Discover the cache miss rate
 - valgrind --tool=cachegrind cmd
 - cachegrind is a tool profiling the cache performance
 - Performance counter
 - Intel[®] Performance Counter Monitor http://www.intel.com/software/pcm/



Matrix Transpose

```
// Transpose matrix b into b_t
                                                                        b_t[i][j] += b[j][i];
                                                                    }
for(i = 0; i < ARRAY_SIZE; i+=(ARRAY_SIZE/n)) {</pre>
                                                                  }
  for(j = 0; j < ARRAY_SIZE; j+=(ARRAY_SIZE/n)) {</pre>
    for(k = 0; k < ARRAY_SIZE; k+=(ARRAY_SIZE/n)) {</pre>
         for(ii = i; ii < i+(ARRAY_SIZE/n); ii++)</pre>
           for(jj = j; jj < j+(ARRAY_SIZE/n); jj++)</pre>
             for(kk = k; kk < k+(ARRAY_SIZE/n); kk++)</pre>
               c[ii][jj] += a[ii][kk]*b[kk][jj];
```

for(i = 0; i < ARRAY_SIZE; i+=(ARRAY_SIZE/n)) {</pre> for(j = 0; j < ARRAY_SIZE; j+=(ARRAY_SIZE/n)) {</pre>

for(i = 0; i < ARRAY_SIZE; i+=(ARRAY_SIZE/n)) {</pre> for(j = 0; j < ARRAY_SIZE; j+=(ARRAY_SIZE/n)) {</pre> for(k = 0; k < ARRAY_SIZE; k+=(ARRAY_SIZE/n)) {</pre> for(ii = i; ii < i+(ARRAY_SIZE/n); ii++)</pre> for(jj = j; jj < j+(ARRAY_SIZE/n); jj++)</pre> for(kk = k; kk < k+(ARRAY_SIZE/n); kk++)</pre> // Compute on b_t c[ii][jj] += a[ii][kk]*b_t[jj][kk];

Poll close in 1:30

Block

What kind(s) of misses can matrix transpose remove?

• By transposing a matrix, the performance of matrix multiplication can be further improved. What kind(s) of cache misses does matrix transpose help to remove?

```
for(i = 0; i < ARRAY_SIZE; i+=(ARRAY_SIZE/n)) {</pre>
  for(j = 0; j < ARRAY_SIZE; j+=(ARRAY_SIZE/n)) {</pre>
    for(k = 0; k < ARRAY_SIZE; k+=(ARRAY_SIZE/n)) {</pre>
      for(ii = i; ii < i+(ARRAY_SIZE/n); ii++)</pre>
         for(jj = j; jj < j+(ARRAY_SIZE/n); jj++)</pre>
           for(kk = k; kk < k+(ARRAY_SIZE/n); kk++)</pre>
             c[ii][jj] += a[ii][kk]*b[kk][jj];
}
```

- A. Compulsory miss
- B. Capacity miss
- C. Conflict miss
- D. Capacity & conflict miss
- E. Compulsory & conflict miss



for(i = 0; i < ARRAY_SIZE; i+=(ARRAY_SIZE/n)) {</pre> for(j = 0; j < ARRAY_SIZE; j+=(ARRAY_SIZE/n)) {</pre>

for(i = 0; i < ARRAY_SIZE; i+=(ARRAY_SIZE/n)) {</pre> for(j = 0; j < ARRAY_SIZE; j+=(ARRAY_SIZE/n)) {</pre> for(k = 0; k < ARRAY_SIZE; k+=(ARRAY_SIZE/n)) {</pre> for(ii = i; ii < i+(ARRAY_SIZE/n); ii++)</pre> for(jj = j; jj < j+(ARRAY_SIZE/n); jj++)</pre> for(kk = k; kk < k+(ARRAY_SIZE/n); kk++</pre> // Compute on b_t c[ii][jj] += a[ii][kk]*b_t[jj][kk];

Poll close in 1:30

Block

What kind(s) of misses can matrix transpose ren

• By transposing a matrix, the performance of matrix multiplication can be further improved. What kind(s) of cache misses does matrix transpose help to remove?

```
for(i = 0; i < ARRAY_SIZE; i+=(ARRAY_SIZE/n)) {</pre>
  for(j = 0; j < ARRAY_SIZE; j+=(ARRAY_SIZE/n)) {</pre>
    for(k = 0; k < ARRAY_SIZE; k+=(ARRAY_SIZE/n)) {</pre>
      for(ii = i; ii < i+(ARRAY_SIZE/n); ii++)</pre>
         for(jj = j; jj < j+(ARRAY_SIZE/n); jj++)</pre>
           for(kk = k; kk < k+(ARRAY_SIZE/n); kk++)</pre>
             c[ii][jj] += a[ii][kk]*b[kk][jj];
}
```

- A. Compulsory miss
- B. Capacity miss
- C. Conflict miss
- D. Capacity & conflict miss
- E. Compulsory & conflict miss





for(i = 0; i < ARRAY_SIZE; i+=(ARRAY_SIZE/n)) {</pre> for(j = 0; j < ARRAY_SIZE; j+=(ARRAY_SIZE/n)) {</pre>

for(i = 0; i < ARRAY_SIZE; i+=(ARRAY_SIZE/n)) {</pre> for(j = 0; j < ARRAY_SIZE; j+=(ARRAY_SIZE/n)) {</pre> for(k = 0; k < ARRAY_SIZE; k+=(ARRAY_SIZE/n)) {</pre> for(ii = i; ii < i+(ARRAY_SIZE/n); ii++)</pre> for(jj = j; jj < j+(ARRAY_SIZE/n); jj++)</pre> for(kk = k; kk < k+(ARRAY_SIZE/n); kk++</pre> // Compute on b_t c[ii][jj] += a[ii][kk]*b_t[jj][kk];

What kind(s) of misses can matrix transpose remove?

• By transposing a matrix, the performance of matrix multiplication can be further improved. What kind(s) of cache misses does matrix transpose help to remove?



- A. Compulsory miss
- B. Capacity miss
- C. Conflict miss
- D. Capacity & conflict miss
- E. Compulsory & conflict miss



for(i = 0; i < ARRAY_SIZE; i+=(ARRAY_SIZE/n)) {</pre> for(j = 0; j < ARRAY_SIZE; j+=(ARRAY_SIZE/n)) {</pre>

for(i = 0; i < ARRAY_SIZE; i+=(ARRAY_SIZE/n)) {</pre> for(j = 0; j < ARRAY_SIZE; j+=(ARRAY_SIZE/n)) {</pre> for(k = 0; k < ARRAY_SIZE; k+=(ARRAY_SIZE/n)) {</pre> for(ii = i; ii < i+(ARRAY_SIZE/n); ii++)</pre> for(jj = j; jj < j+(ARRAY_SIZE/n); jj++)</pre> for(kk = k; kk < k+(ARRAY_SIZE/n); kk++</pre> // Compute on b_t c[ii][jj] += a[ii][kk]*b_t[jj][kk];

Summary of Optimizations

- Software
 - Data layout capacity miss, conflict miss, compulsory miss
 - Blocking capacity miss
 - Transpose conflict miss
 - Loop fission conflict miss when \$ has limited way associativity
 - Loop fusion capacity miss when \$ has enough way associativity
 - Loop interchange conflict/capacity miss
- Hardware
 - Prefetch compulsory miss
 - Write buffer miss penalty
 - Bank/pipeline miss penalty
 - Critical word first and early restart miss panelty





Which version is faster?

```
for(i=0;i<100000000;i++)</pre>
for(i=0;i<100000000;i++)</pre>
                                            \mathbf{m}
     sum+=data[(i*15) & 131071];
}
                                                }
```

- Both version A and B produces the same output. Without compiler optimization, which version of code would have better performance?
 - A. Version A
 - B. Version B
 - C. They are about the same (less than 5% difference)



sum+=data[((i << 4) - i) & 131071];</pre>



Which version is faster?

```
for(i=0;i<100000000;i++)</pre>
for(i=0;i<100000000;i++)</pre>
                                            \mathbf{m}
     sum+=data[(i*15) & 131071];
}
                                                }
```

- Both version A and B produces the same output. Without compiler optimization, which version of code would have better performance?
 - A. Version A
 - B. Version B
 - C. They are about the same (less than 5% difference)





sum+=data[((i << 4) - i) & 131071];</pre>

Basic Pipelined Processor

Hung-Wei Tseng



Tasks in RISC-V ISA

- Instruction Fetch (IF) fetch the instruction from memory
- Instruction Decode (ID)
 - Decode the instruction for the desired operation and operands
 - Reading source register values
- Execution (**EX**)
 - ALU instructions: Perform ALU operations
 - Conditional Branch: Determine the branch outcome (taken/not taken)
 - Memory instructions: Determine the effective address for data memory access
- Data Memory Access (MEM) Read/write memory
- Write Back (WB) Present ALU result/read value in the target register
- Update PC
 - If the branch is taken set to the branch target address
 - Otherwise advance to the next instruction current PC + 4

Simple implementation w/o branch

- add x1, x2, x3 D IF EX WB
- ld x4, 0(x5)
- sub x6, x7, x8
- sub x9, x10, x11
- sd x1, 0(x12)



IF

D









1







- Different parts of the processor works on different instructions simultaneously
- A clock signal controls and synchronize the beginning and the end of each part of the work
- A pipeline register between different parts of the processor to keep intermediate results necessary for the upcoming work



add x1, x2, x3 ld x4, 0(x5) sub x6, x7, x8 sub x9, x10, x11 sd x1, 0(x12) xor x13, x14, x15 and x16, x17, x18 add x19, x20, x21 sub x22, x23, x24 ld x25, 4(x26) sd x27, 0(x28)

IF

ID	EX	MEM	WB				
IF	ID	EX	MEM	WB		I	
	IF	ID	EX	MEM	WB		
		IF	ID	EX	MEM	WB	
			IF	ID	EX	MEM	
				IF	ID	EX	
					IF	ID	
						IF	
			After this point, we are completing instruction each c				



Which version is faster?



- Both version A and B produces the same output. Without compiler optimization, which version of code would have significantly better performance?
 - A. Version A
 - B. Version B
 - C. They are about the same (less than 10% difference) - Because we have pipelined instructions, the CPI of one instruction doesn't matter as long as we can keep the pipeline busy



sum+=data[((i << 4) - i) & 131071];</pre>



Announcement

- Project is up check the website
- Assignment #3 due next Monday
- Midterm
 - Release next Tuesday (11/10) 12:00am, turn in before next Friday (11/13) 11:59pm
 - You can only open it once and you have to finish a total of 30 questions within 80 minutes.
 - You may open book, but you have to bare the risks of not being able to finish them
- Attendance
 - The attendance throughout the quarter count as one assignment
 - You only need to answer 50% of the Zoom polls to receive full credits
 - Please don't email me for absence we count only 50% to give you flexibility
 - If you just login but never answer questions, you won't receive any.
- Reading Quizzes 2 attempts, average
- Office Hours on Zoom (the office hour link, not the lecture one)
 - Hung-Wei/Prof. Usagi: M 8p-10p (make up for the last week), W 2p-3p
 - Quan Fan: F 1p-3p

1:59pm hin 80 minutes. h them Computer Science & Engineering





