May 4: Caches





Caches

Memory & data

Integers & floats

Machine code & C

Roadmap

C:

Java:



2

How does execution time grow with SIZE?

```
int array[SIZE];
```

```
int sum = 0;
```

```
for (int i = 0; i < 200000; i++) {
    for (int j = 0; j < SIZE; j++) {
        sum += array[j];
     }
</pre>
```



Actual Data



Making memory accesses fast!

- Cache basics
- Principle of locality
- Memory hierarchies
- Cache organization
- Program optimizations that consider caches

Problem: Processor-Memory Bottleneck



Problem: lots of waiting on memory

cycle: single machine step (fixed-time)

Problem: Processor-Memory Bottleneck



Solution: caches

cycle: single machine step (fixed-time)



 A hidden storage space for provisions, weapons, and/or <u>treasures</u>



- Computer memory with short access time used for the storage of frequently or recently used instructions or data (i-cache and dcache)
 - More generally:

used to optimize data transfers between any system elements with different characteristics (network interface cache, I/O cache, etc.)

General Cache Mechanics



General Cache Concepts: Hit



Data in block b is needed

Block b is in cache: Hit!

General Cache Concepts: Miss



Data in block b is needed

Block b is not in cache: Miss!

Block b is fetched from memory

Block b is stored in cache

- Placement policy: determines where b goes
- Replacement policy: determines which block gets evicted (victim)

Locality: Programs tend to use data and instructions with addresses near or equal to those they have used recently

Locality: Programs tend to use data and instructions with addresses near or equal to those they have used recently

Temporal locality:

- Recently referenced items are *likely* to be referenced again in the near future
- Why is this important?





Locality: Programs tend to use data and instructions with addresses near or equal to those they have used recently

Temporal locality:

 Recently referenced items are *likely* to be referenced again in the near future



Spatial locality?

Locality: Programs tend to use data and instructions with addresses near or equal to those they have used recently

Temporal locality:

 Recently referenced items are *likely* to be referenced again in the near future

Spatial locality:

- Items with nearby addresses *tend* to be referenced close together in time
- How do caches take advantage of this?





Example: Any Locality?

```
sum = 0;
for (i = 0; i < n; i++) {
   sum += a[i];
}
return sum;
```

Example: Any Locality?

```
sum = 0;
for (i = 0; i < n; i++) {
   sum += a[i];
}
return sum;
```

Data:

- <u>Temporal</u>: <u>sum</u> referenced in each iteration
- Spatial: arraya[] accessed in stride-1 pattern
- Instructions:
 - <u>Temporal</u>: cycle through loop repeatedly
 - **Spatial**: reference instructions in sequence

Being able to assess the locality of code is a crucial skill for a programmer

```
int sum_array_rows(int a[M][N])
{
    int i, j, sum = 0;
    for (i = 0; i < M; i++)
        for (j = 0; j < N; j++)
            sum += a[i][j];
    return sum;
}</pre>
```

a[0][0]	a[0][1]	a[0][2]	a[0][3]
a[1][0]	a[1][1]	a[1][2]	a[1][3]
a[2][0]	a[2][1]	a[2][2]	a[2][3]



stride-1

76 is just one possible starting address of array a

```
int sum_array_cols(int a[M][N])
{
    int i, j, sum = 0;
    for (j = 0; j < N; j++)
        for (i = 0; i < M; i++)
            sum += a[i][j];
    return sum;
}</pre>
```

a[0][0]	a[0][1]	a[0][2]	a[0][3]
a[1][0]	a[1][1]	a[1][2]	a[1][3]
a[2][0]	a[2][1]	a[2][2]	a[2][3]



Layout in Memory





Locality Example #3 int sum_array_3d(int a[M][N][L]) {

```
int i, j, k, sum = 0;
```



What is wrong with this code?

• How can it be fixed?





•••

Make memory fast again.

- Cache basics
- Principle of locality
- Memory hierarchies
- Cache organization
- Program optimizations that consider caches

Cost of Cache Misses

- Huge difference between a hit and a miss
 - Could be 100x, if just L1 and main memory

Would you believe 99% hits is twice as good as 97%?

Consider:

Cache hit time of 1 cycle Miss penalty of 100 cycles

cycle = single fixed-time machine step

Cost of Cache Misses

- Huge difference between a hit and a miss
 - Could be 100x, if just L1 and main memory
- Would you believe 99% hits is twice as good as 97%?
 - Consider: Cache hit time of 1 cycle Miss penalty of 100 cycles

cycle = single fixed-time machine step

Average access time: check the cache every time

97% hits: 1 cycle + 0.03 * 100 cycles = 4 cycles

99% hits: 1 cycle + 0.01 * 100 cycles = 2 cycles

This is why "miss rate" is used instead of "hit rate"

Cache Performance Metrics



Miss Rate

Fraction of memory references not found in cache (misses / accesses)
 = 1 - hit rate

Caches

- Typical numbers (in percentages):
 - 3% 10% for L1
 - Can be quite small (e.g., < 1%) for L2, depending on size, etc.

Hit Time

- Time to deliver a line in the cache to the processor
 - Includes time to determine whether the line is in the cache
- Typical hit times: 4 clock cycles for L1; 10 clock cycles for L2

Miss Penalty

- Additional time required because of a miss
- Typically 50 200 cycles for missing in L2 & going to main memory (Trend: increasing!)

Spring 2016

Can we have more than one cache?

Why would we want to do that?

Memory Hierarchies

- Some fundamental and enduring properties of hardware and software systems:
 - Faster storage technologies almost always cost more per byte and have
 lower capacity
 - The gaps between memory technology speeds are widening
 - True for: registers \leftrightarrow cache, cache \leftrightarrow DRAM, DRAM \leftrightarrow disk, etc.
 - Well-written programs tend to exhibit good locality
- These properties complement each other beautifully
- They suggest an approach for organizing memory and storage systems known as a <u>memory hierarchy</u>

May 6

An Example Memory Hierarchy



An Example Memory Hierarchy



An Example Memory Hierarchy



Memory Hierarchies

Fundamental idea of a memory hierarchy:

 For each level k, the faster, smaller device at level k serves as a cache for the larger, slower device at level k+1.

Why do memory hierarchies work?

- Because of locality, programs tend to access the data at level k more often than they access the data at level k+1.
- Thus, the storage at level k+1 can be slower, and thus larger and cheaper per bit.

Big Idea: The memory hierarchy creates a large pool of storage that costs as much as the cheap storage near the bottom, but that serves data to programs at the rate of the fast storage near the top.

Intel Core i7 Cache Hierarchy

Processor package



Block size: 64 bytes for all caches.

L1 i-cache and d-cache: 32 KB, 8-way, Access: 4 cycles

L2 unified cache: 256 KB, 8-way, Access: 11 cycles

L3 unified cache:

8 MB, 16-way, Access: 30-40 cycles

Making memory accesses fast!

- Cache basics
- Principle of locality
- Memory hierarchies
 - Cache organization
 - Direct-mapped (sets; index + tag)
 - Associativity (ways)
 - Replacement policy
 - Handling writes
 - **Program optimizations that consider caches**
Cache Organization

Where should data go in the cache?

- We need a mapping from memory addresses to specific locations in the cache to make checking the cache for an address **fast**
 - Otherwise each memory access requires "searching the entire cache" (slow!)
- What is a data structure that provides fast lookup?

Aside: Hash Tables for Fast Lookup



Aside: Hash Tables for Fast Lookup



Where should we put data in the cache?



Where should we put data in the cache?



Use tags to record which location is cached



What's a cache block? (or cache line)



A puzzle.

- What can you infer from this:
- Cache starts *empty*
- Access (addr: hit/miss) stream:

(12: miss), (13: hit), (14: miss)
block size >= 2 bytes

Direct mapped cache

Direct mapped:

- Each memory address can be mapped to exactly one index in the cache.
- Easy to find an address!
 Cheap to implement.
- But...
- What happens if a program uses addresses
 2, 6, 2, 6, 2, ...?

2 and 6 *conflict,* rest of cache is <u>unused</u>

Memory **Address** 0000 0001 0010 0011 0100 Index Tag 0101 00 0110 01 0111 \mathfrak{O} 10 1000 11 1001

Caches

1010

1011 1100 1101

1110 1111

Associativity

- What if we could store data in any place in the cache?
 - Minimize conflicts!
 - More complicated hardware, consumes more *power*, slower.





Associativity

- What if we could store data in any place in the cache?
 - More complicated hardware, consumes more *power*, slower.
- So we combine the two ideas:
 - Each address maps to exactly one set.
 - Each set can have more than one way.





direct mapped

Now how do I know where data goes?



Now how do I know where data goes?



Example placement in set-associative caches



Example placement in set-associative caches

Caches



Block replacement

- Any empty block in the correct set may be used for storing data.
- If there are no empty blocks, which one should we replace?
 - Obvious for direct-mapped caches, what about set-associative?
 - Caches typically use something close to *least recently used (LRU)* (hardware usually implements "*not most recently used*")



Another puzzle.

What can you infer from this?

Cache starts *empty*

Access (addr: hit/miss) stream:
(10) (11) (12) (12: miss); (10: miss)
(10: miss); (12: miss); (10: miss)
12 is not in the same
12 is not in the same
block as 10
(0) 0(0) 0

direct-mapped cache with <=2 sets

General Cache Organization (S, E, B)



Caches



Example: Direct-Mapped Cache (E = 1)

Direct-mapped: One line per set Assume: cache block size 8 bytes



Example: Direct-Mapped Cache (E = 1)

Direct-mapped: One line per set Assume: cache block size 8 bytes

t 000 000



Example: Direct-Mapped Cache (E = 1)

Direct-mapped: One line per set Assume: cache block size 8 bytes



No match? Then old line gets evicted and replaced



Example (for E = 1)

```
float dotprod(float x[8], float y[8])
{
    float sum = 0;
    int i;
    for (i = 0; i < 8; i++)
        sum += x[i]*y[i];
    return sum;
}</pre>
```

In this example, cache blocks are 16 bytes; 8 sets in cache How many block offset bits? How many set index bits?

Address bits: ttt....t sss bbbb B = 16 = 2^b: b=4 offset bits S = 8 = 2^s: s=3 index bits

0: 000....0 000 0000 128: 000....1 000 0000 160: 000....1 010 0000

if x and y have aligned starting addresses, e.g., &x[0] = 0, &y[0] = 128



if x and y have <u>unaligned</u> starting addresses, e.g., &x[0] = 0, &y[0] = 160



Address of short int:

E-way Set-Associative Cache (Here: E = 2)

E = 2: Two lines per set Assume: cache block size 8 bytes



E-way Set-Associative Cache (Here: E = 2)

E = 2: Two lines per set Assume: cache block size 8 bytes

Address of short int:



block offset

E-way Set-Associative Cache (Here: E = 2)

E = 2: Two lines per set Assume: cache block size 8 bytes

Address of short int:



short int (2 Bytes) is here

No match?

- One line in set is selected for eviction and replacement
- Replacement policies: random, least recently used (LRU), ...

Example (for E = 2)

```
float dotprod(float x[8], float y[8])
{
    float sum = 0;
    int i;
    for (i = 0; i < 8; i++)
        sum += x[i]*y[i];
    return sum;
}</pre>
```

If x and y have aligned starting addresses, e.g. &x[0] = 0, &y[0] = 128, can still fit both because two lines in each set

x[0]	x[1]	x[2]	x[3]	y[0]	y[1]	y[2]	y[3]
x[4]	x[5]	x[6]	x[7]	y[4]	y[5]	y[6]	y[7]

Types of Cache Misses: 3 C's!

Cold (compulsory) miss

Occurs on first access to a block

Conflict miss

- Conflict misses occur when the cache is large enough, but multiple data objects all map to the same slot
 - e.g., referencing blocks 0, 8, 0, 8, ... could miss every time
- direct-mapped caches have more conflict misses than n-way set-associative (where n > 1)

Capacity miss

- Occurs when the set of active cache blocks (the *working set*) is larger than the cache (just won't fit, even if cache was *fully-associative*)
- Note: Fully-associative only has Cold and Capacity misses

What about writes?

Multiple copies of data exist:

L1, L2, possibly L3, main memory

What is the main problem with that?

- Which copies should we update?
- What if it's a hit? Miss?



What about writes?

- Multiple copies of data exist:
 - L1, L2, possibly L3, main memory

What to do on a write-hit?

- Write-through: write immediately to memory, all caches in between.
- Write-back: defer write to memory until line is evicted (replaced)
 - Must track which cache lines have been modified ("dirty bit")

What to do on a write-miss?

- Write-allocate("fetch on write"): load into cache, update line in cache.
 - Good if more writes or reads to the location follow
- <u>No-write-allocate("write around")</u>: just write immediately to memory.

Typical caches:

- Write-through + No-write-allocate, occasionally



tag (there is only one set in this tiny cache, so the tag is the entire address!)

Memory T U U ØxBEEF

In this example we are sort of ignoring block offsets. Here a block holds 2 bytes (16 bits, 4 hex digits).

Normally a block would be much bigger and thus there would be multiple items per block. While only one item in that block would be written at a time, the entire line would be brought into cache.

Т

Step 1: Bring T into cacheStep 2: Write ØxFACE to cacheonly and set dirty bit.

mov ØxFEED,T

mov ØxFACE,

Write hit! Write ØxFEED to cache only

mov U,%rax

- 1. Write T back to memory since it's dirty.
- 2. Bring U into the cache so we can copy it into %rax



mov ØxFACE, T



Step 1: Bring T into cache

Memory



mov ØxFACE, T



Step 2: Write ØxFACE to cache only <u>and set</u> <u>dirty bit.</u>

Memory



mov ØxFACE,T mov ØxFEED,T





Memory


Write-back, write-allocate example

mov ØxFACE,T

mov ØxFEED,T mov U,%rax



 Write T back to memory since it is dirty.
 Bring U into the cache so

we can copy it into %rax





Back to the Core i7 to look at ways

Processor package



May 11

Reminders

- Lab 3 due tonight @ 11:59pm
- HW 3 (mostly cache problems) due Friday

Where else is caching used?

Software Caches are More Flexible

Examples

- File system buffer caches, browser caches, etc.
- Content-delivery networks (CDN): cache for the Internet (e.g. Netflix)

Some design differences

- Almost always fully-associative
 - so, no placement restrictions
 - index structures like hash tables are common (for placement)
- More complex replacement policies
 - misses are very expensive when disk or network involved
 - worth thousands of cycles to avoid them
- Not necessarily constrained to single "block" transfers
 - may fetch or write-back in larger units, opportunistically

Optimizations for the Memory Hierarchy

Write code that has locality!

- Spatial: access data contiguously
- Temporal: make sure access to the same data is not too far apart in time

How can you achieve locality?

- Proper choice of algorithm
- Loop transformations

Example: Array of Structs

```
typedef struct {
    char * name;
    int scores[4];
} Student;
/* Array of structs */
Student students[N];
/* Compute average score for assignment 'hw' */
void average_score(int hw) {
    int total = 0;
    for (int s = 0; s < N; s++) {
        total += students[s].scores[hw];
    return total / (double)N;
}</pre>
```

Memory



Example: Array of Structs

capacity:	32 bytes
block size:	16 bytes
sets:	2
ways:	1



```
double average_score(int hw) {
    int total = 0;
    for (int s = 0; s < N; s++) {
        total += students[s].scores[hw];
    return total / (double)N;
}</pre>
```



- 1. students[0].scores[1] *set:0,tag:000* COLD MISS
- 2. students[1].scores[1] set:0,tag:001
 COLD MISS
- 3. students[2].scores[1] *set:1,tag:001* COLD MISS

Example: Array of Structs

Cache Miss Analysis

- Accesses: int (4 bytes)
- *Stride:* sizeof(Student) = 24
- Cache block size: 16 bytes
- No temporal locality, so only **cold** misses
- Miss rate: 100%

```
typedef struct {
    char * name;
    int scores[4];
 Student;
}
/* Array of structs */
Student students[N];
```

```
double average_score(int hw) {
    int total = 0;
    for (int s = 0; s < N; s++) {</pre>
        total += students[s].scores[hw];
    return total / (double)N;
```



Cache

capacity:	32 bytes
block size:	16 bytes
sets:	2
ways:	1

81



Memory



/*	"Struct of arrays" */			
struct {				
<pre>char * names[N];</pre>				
<pre>int scores[4][N];</pre>				
} s	students;			



```
double average_score(int hw) {
    int total = 0;
    for (int s = 0; s < N; s++) {
        total += students[s].scores[hw];
    return total / (double)N;</pre>
```

- 1. students.scores[1][0] *set:1,tag:100* COLD MISS
- 2. students.scores[1][1] set:1,tag:101
 HIT
- 3. students.scores[1][2] *set:1,tag:101* HIT



```
/* "Struct of arrays" */
struct {
    char * names[N];
    int scores[4][N];
} students;
```



Cache Miss Analysis

- Accesses: int (4 bytes)
- *Stride:* 4 bytes
- Cache block size: 16 bytes
- 16 bytes/block / 4 bytes/int = 4 ints/block
- Miss rate:
 - I COLD MISS / block
 - 4 ints/block 1 MISS/block = 3 HIT/block
 - 75% HIT rate

<u>Cache</u>			
capacit	:y:	32 bytes	
block s	ize:	16 bytes	
sets:		2	
ways:		1	

```
/* "Struct of arrays" */
struct {
    char * names[N];
    int scores[4][N];
} students;
```



But don't forget about other uses of this struct/array...

```
double student_grade(int s) {
    int total = 0;
    for (int hw = 0; hw < 4; hw++) {
        total += students.scores[s][hw];
    return total / (double)4;
}</pre>
```





What would the MISS rate be for this?



Extra example (skipping in class)

Example: Matrix Multiplication





memory access pattern?

Extra example **Cache Miss Analysis** spatial locality: (skipping in class) chunks of 8 items in a row in same cache line Assume: each item in column in Matrix elements are doubles different cache line Cache block = 64 bytes = 8 doubles Cache size C << n (much smaller than n, **not** left-shifted by n) n/8 misses n **First iteration:** n/8 + n = 9n/8 misses (omitting matrix c) * Afterwards in cache: (schematic) *

8 doubles wide

n misses

Cache Miss Analysis

Assume:

- Matrix elements are doubles
- Cache block = 64 bytes = 8 doubles
- Cache size C << n (much smaller than n)

Other iterations:

 Again: n/8 + n = 9n/8 misses (omitting matrix c)



Total misses:

9n/8 * n² = (9/8) * n³

once per element

Caches



Spring 2016

Extra example

(skipping in class)

Blocked Matrix Multiplication



Cache Miss Analysis

Assume:

- Cache block = 64 bytes = 8 doubles
- Cache size C << n (much smaller than n)
- Three blocks fit into cache: 3B² < C</p>



Extra example

(skipping in class)

Caches

Cache Miss Analysis

- Cache block = 64 bytes = 8 doubles
- Cache size C << n (much smaller than n)
- Three blocks fit into cache: 3B² < C</p>



- Same as first iteration
- $2n/B * B^2/8 = nB/4$



Total misses:

• $nB/4 * (n/B)^2 = n^3/(4B)$

Extra example (skipping in class)

Caches

Extra example (skipping in class)

Summary

- No blocking: (9/8) * n³
- Blocking: 1/(4B) * n³
- If B = 8 difference is 4 * 8 * 9 / 8 = 36x
- If B = 16 difference is 4 * 16 * 9 / 8 = 72x
- Suggests largest possible block size B, but limit 3B² < C!</p>

Reason for dramatic difference:

- Matrix multiplication has inherent temporal locality:
 - Input data: 3n², computation 2n³
 - Every array element used O(n) times!
- But program has to be written properly

Cache-Friendly Code

Programmer can optimize for cache performance

- How data structures are organized
- How data are accessed
 - Nested loop structure
 - Blocking (previous example) is a general technique

All systems favor "cache-friendly code"

- Getting absolute optimum performance is very platform specific
 - Cache sizes, line sizes, associativities, etc.
- Can get most of the advantage with generic code
 - Keep working set reasonably small (temporal locality)
 - Use small strides (spatial locality)
 - Focus on inner loop code

Intel Core i7 Cache Hierarchy

Processor package



L1 i-cache and d-cache: 32 KB, 8-way, Access: 4 cycles

L2 unified cache: 256 KB, 8-way, Access: 11 cycles

L3 unified cache: 8 MB, 16-way, Access: 30-40 cycles

Block size: 64 bytes for all caches.

Main memory

