3

Cache Memories

15-213: Introduction to Computer Systems 11th Lecture, Feb. 19, 2013

Instructors:

Seth Copen Goldstein, Anthony Rowe, Greg Kesden

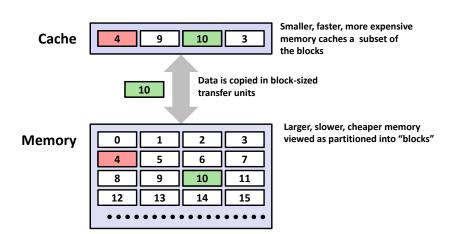
Today

- Cache memory organization and operation
- **■** Performance impact of caches
 - The memory mountain
 - Rearranging loops to improve spatial locality
 - Using blocking to improve temporal locality

Carnegie Mell

Carnegie Mellon

General Cache Concept (Reminder)



Many types of caches

Examples

- Hardware: L1 and L2 CPU caches, TLBs, ...
- Software: virtual memory, FS buffers, web browser caches, ...

Many common design issues

- each cached item has a "tag" (an ID) plus contents
- need a mechanism to efficiently determine whether given item is cached
 - combinations of indices and constraints on valid locations
- on a miss, usually need to pick something to replace with the new item
 - called a "replacement policy"
- on writes, need to either propagate change or mark item as "dirty"
 - · write-through vs. write-back

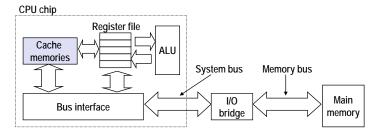
Different solutions for different caches

Lets talk about CPU caches as a concrete example...

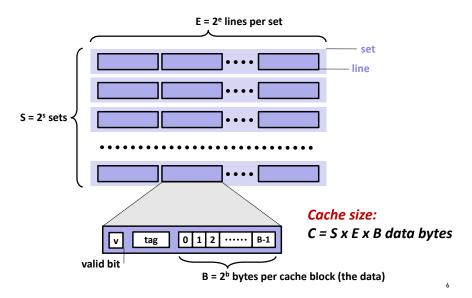
Carnegie Mellon

CPU Cache Memories

- CPU Cache memories are small, fast SRAM-based memories managed automatically in hardware
 - Hold frequently accessed blocks of main memory
- CPU looks first for data in caches (e.g., L1, L2, and L3), then in main memory
- Typical system structure:

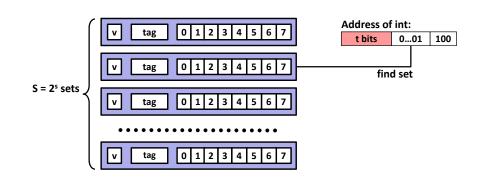


General Cache Organization (S, E, B)

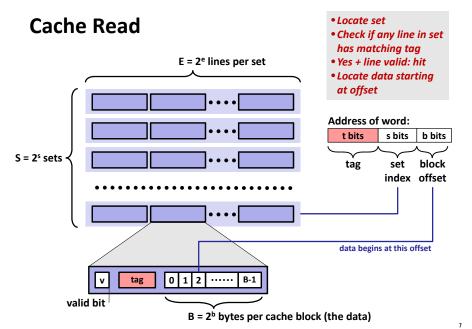


Example: Direct Mapped Cache (E = 1)

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

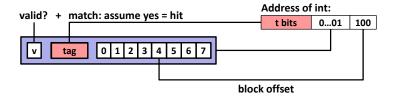


Carnegie Mellon



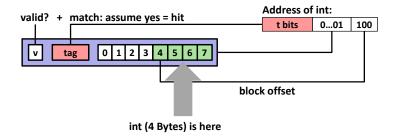
Example: Direct Mapped Cache (E = 1)

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



If tag doesn't match: old line is evicted and replaced

...

Carnegie Mellon

Carnegie Mellon

Direct-Mapped Cache Simulation

t=1	s=2	b=1
Х	XX	Х

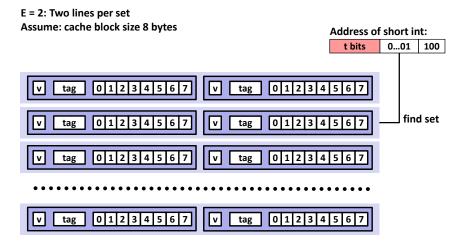
M=16 bytes (4-bit addresses), B=2 bytes/block, S=4 sets, E=1 Blocks/set

Address trace (reads, one byte per read):

trace	(Teaus, one byt	e per rea
0	[0 <u>00</u> 0 ₂],	miss
1	$[0001_{2}],$	hit
7	$[0111_2],$	miss
8	$[1000_2],$	miss
0	[0000]	miss

	V	Tag	Block
Set 0	1	0	M[0-1]
Set 1			
Set 2			
Set 3	1	0	M[6-7]

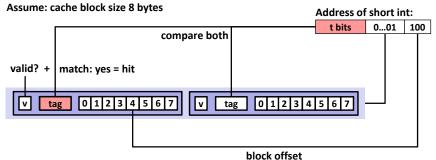
E-way Set Associative Cache (Here: E = 2)



11

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

E = 2: Two lines per set

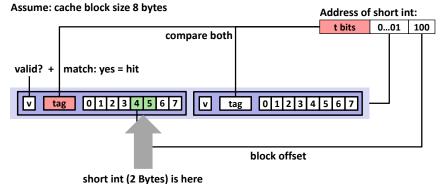


13

15

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

E = 2: Two lines per set



No match:

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

14

Carnegie Mellon

2-Way Set Associative Cache Simulation

t=2 s=1 b=1 xx x x x

M=16 byte addresses, B=2 bytes/block, S=2 sets, E=2 blocks/set

Address trace (reads, one byte per read):

s trace	(reads, one byt	e per rea
0	[00 <u>0</u> 0 ₂],	miss
1	[00 <u>0</u> 1,],	hit
7	[01 <u>1</u> 1,],	miss
8	$[1000_{2}],$	miss
0	[0000,1	hit

	V	Tag	Block
Set 0	1	00	M[0-1]
	1	10	M[8-9]
Set 1	1	01	M[6-7]
	0		

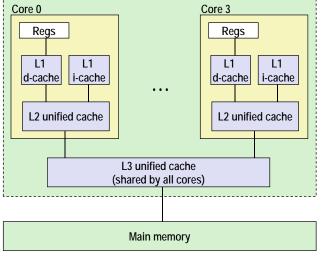
What about writes?

- Multiple copies of data exist:
 - L1, L2, Main Memory, Disk
- What to do on a write-hit?
 - Write-through (write immediately to memory)
 - Write-back (defer write to memory until replacement of line)
 - Need a dirty bit (line different from memory or not)
- What to do on a write-miss?
 - Write-allocate (load into cache, update line in cache)
 - Good if more writes to the location follow
 - No-write-allocate (writes straight to memory, does not load into cache)
- Typical
 - Write-through + No-write-allocate
 - Write-back + Write-allocate

Carnegie Mello

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.

17

Cache Performance Metrics

Miss Rate

- Fraction of memory references not found in cache (misses / accesses)
 = 1 hit rate
- 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 numbers:
 - 1-2 clock cycle for L1
 - 5-20 clock cycles for L2

Miss Penalty

- Additional time required because of a miss
 - typically 50-200 cycles for main memory (Trend: increasing!)

...

Carnegie Mello

Carnegie Mellon

Lets think about those numbers

- 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
 - Average access 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"

Writing Cache Friendly Code

- Make the common case go fast
 - Focus on the inner loops of the core functions
- Minimize the misses in the inner loops
 - Repeated references to variables are good (temporal locality)
 - Stride-1 reference patterns are good (spatial locality)

Key idea: Our qualitative notion of locality is quantified through our understanding of cache memories

Back to Observations

- Programmer can optimize for cache performance
 - How data structures are organized
 - How data are accessed (examples follow)
 - Nested loop structure
 - Blocking 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)

Today

- Cache organization and operation
- Performance impact of caches
 - The memory mountain
 - Rearranging loops to improve spatial locality
 - Using blocking to improve temporal locality

21

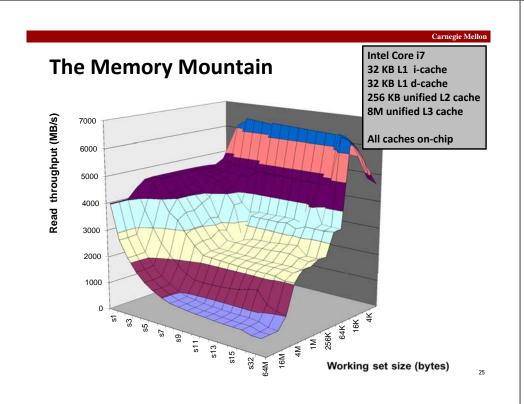
Carnegie Mellon

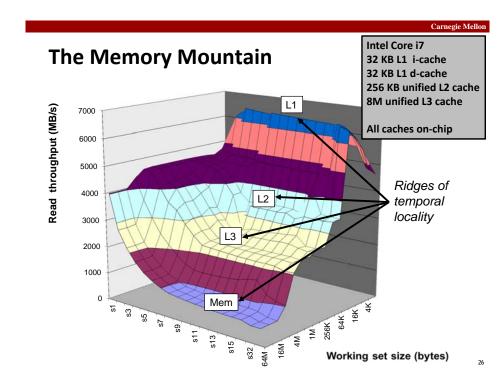
Carnegie Mellon

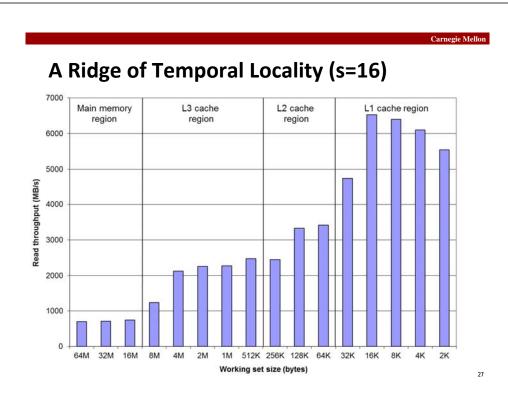
The Memory Mountain

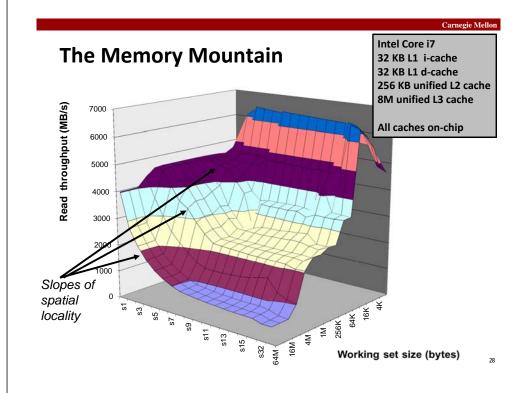
- Read throughput (read bandwidth)
 - Number of bytes read from memory per second (MB/s)
- Memory mountain: Measured read throughput as a function of spatial and temporal locality.
 - Compact way to characterize memory system performance.

Memory Mountain Test Function



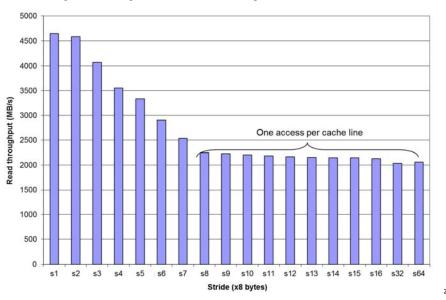


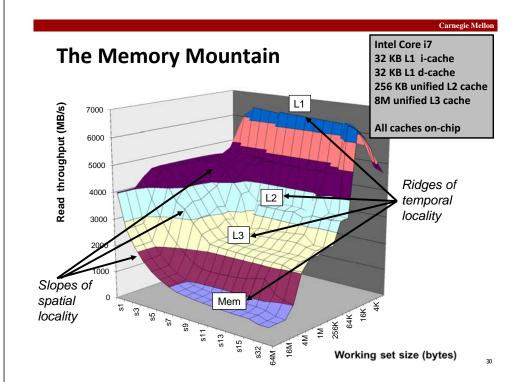






A Slope of Spatial Locality (size=4MB)





Carnegie Mellon

Today

- Cache organization and operation
- **■** Performance impact of caches
 - The memory mountain
 - Rearranging loops to improve spatial locality
 - Using blocking to improve temporal locality

Carnegie Mellon

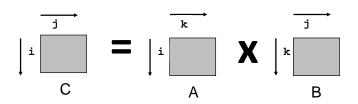
Miss Rate Analysis for Matrix Multiply

Assume:

- Line size = 32B (big enough for four 64-bit words)
- Matrix dimension (N) is very large
 - Approximate 1/N as 0.0
- Cache is not even big enough to hold multiple rows

Analysis Method:

Look at access pattern of inner loop



Matrix Multiplication Example

- Description:
 - Multiply N x N matrices
 - O(N³) total operations
 - N reads per source element
 - N values summed per destination
 - but may be able to hold in register

```
Variable sum
/* ijk */
                      held in register
for (i=0; i<n; i++)
  for (j=0; j<n; j++) {
    sum = 0.0; -
    for (k=0; k<n; k++)
      sum += a[i][k] * b[k][j];
    c[i][j] = sum;
 }
```

Layout of C Arrays in Memory (review)

- C arrays allocated in row-major order
 - each row in contiguous memory locations
- Stepping through columns in one row:

```
for (i = 0; i < N; i++)</pre>
   sum += a[0][i];
```

- accesses successive elements
- if block size (B) > 4 bytes, exploit spatial locality
 - miss rate = 4 bytes / B
- Stepping through rows in one column:

```
for (i = 0; i < n; i++)</pre>
   sum += a[i][0];
```

- accesses distant elements
- no spatial locality!
 - miss rate = 1 (i.e. 100%)

Carnegie Mellon

Carnegie Mello

Matrix Multiplication (ijk)

```
/* ijk */
for (i=0; i<n; i++) {
  for (j=0; j<n; j++) {
    sum = 0.0;
    for (k=0; k< n; k++)
      sum += a[i][k] * b[k][j];
    c[i][j] = sum;
```

```
Inner loop:
                            (<u>i.</u>j)
Row-wise
            Column-
                           Fixed
```

wise

Misses per inner loop iteration:

<u>A</u> <u>B</u> <u>C</u> 0.25 1.0 0.0

Matrix Multiplication (jik)

```
/* jik */
for (j=0; j<n; j++) {
  for (i=0; i<n; i++) {
    sum = 0.0;
    for (k=0; k<n; k++)
      sum += a[i][k] * b[k][j];
    c[i][j] = sum
```

```
Inner loop:
                       (i,j)
Row-wise Column-
                       Fixed
             wise
```

Misses per inner loop iteration:

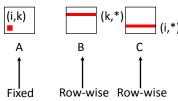
<u>A</u> <u>C</u> 0.25 0.0 1.0

Carnegie Mellor

Matrix Multiplication (kij)

```
/* kij */
for (k=0; k<n; k++) {
  for (i=0; i<n; i++) {
    r = a[i][k];
    for (j=0; j<n; j++)
      c[i][j] += r * b[k][j];
```

Inner loop:

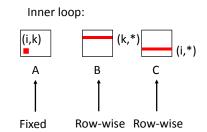


Misses per inner loop iteration:

<u>A</u> 0.0 0.25 0.25

Matrix Multiplication (ikj)

/* ikj */ for (i=0; i<n; i++) { for (k=0; k<n; k++) { r = a[i][k];for (j=0; j<n; j++) c[i][j] += r * b[k][j];



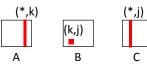
Misses per inner loop iteration:

0.0 0.25 0.25

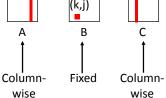
Carnegie Mellon

Matrix Multiplication (jki)

```
/* jki */
for (j=0; j<n; j++) {
  for (k=0; k<n; k++) {
   r = b[k][j];
   for (i=0; i<n; i++)
      c[i][j] += a[i][k] * r;
```



Inner loop:

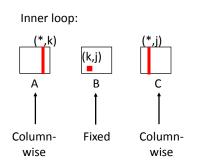


Misses per inner loop iteration:

<u>A</u> <u>B</u> 1.0 0.0 1.0

Matrix Multiplication (kji)

```
/* kji */
for (k=0; k<n; k++) {
  for (j=0; j<n; j++) {
    r = b[k][j];
    for (i=0; i<n; i++)
      c[i][j] += a[i][k] * r;
```



Misses per inner loop iteration:

<u>A</u> <u>C</u> 1.0 0.0 1.0

Summary of Matrix Multiplication

```
for (i=0; i<n; i++) {
  for (j=0; j<n; j++) {
   sum = 0.0;
   for (k=0; k<n; k++)
     sum += a[i][k] * b[k][j];
   c[i][j] = sum;
for (k=0; k<n; k++) {
 for (i=0; i<n; i++) {
 r = a[i][k];
 for (j=0; j<n; j++)
  c[i][j] += r * b[k][j];
for (j=0; j<n; j++) {
 for (k=0; k<n; k++) {
  r = b[k][j];
   for (i=0; i<n; i++)
   c[i][j] += a[i][k] * r;
```

ijk (& jik):

- 2 loads, 0 stores
- misses/iter = 1.25

kij (& ikj):

- 2 loads, 1 store
- misses/iter = **0.5**

jki (& kji):

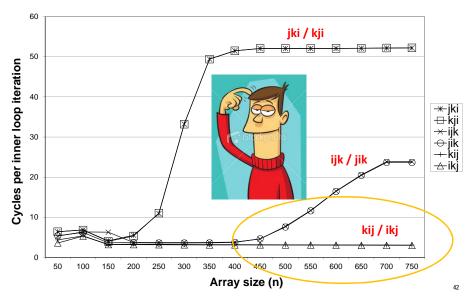
- 2 loads, 1 store
- misses/iter = 2.0

600

41

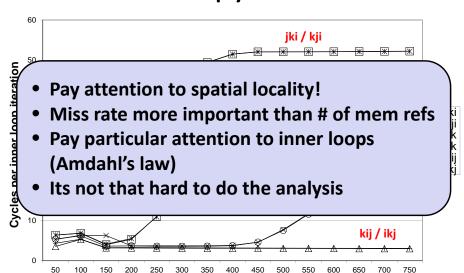
43

Core i7 Matrix Multiply Performance



Carnegie Mellon

Core i7 Matrix Multiply Performance



Array size (n)

Carnegie Mello

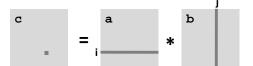
Today

- Cache organization and operation
- **Performance impact of caches**
 - The memory mountain
 - Rearranging loops to improve spatial locality
 - Using blocking to improve temporal locality

Example: Matrix Multiplication

```
c = (double *) calloc(sizeof(double), n*n);

/* Multiply n x n matrices a and b */
void mmm(double *a, double *b, double *c, int n) {
   int i, j, k;
   for (i = 0; i < n; i++)
        for (j = 0; j < n; j++)
        for (k = 0; k < n; k++)
        c[i*n+j] += a[i*n + k]*b[k*n + j];
}</pre>
```

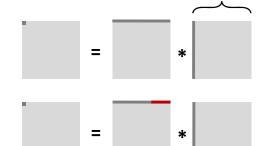


Cache Miss Analysis

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

First iteration:

- n/8 + n = 9n/8 misses
- Afterwards in cache: (schematic)



45

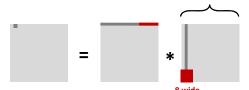
Carnegie Mellon

Cache Miss Analysis

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

Second iteration:

Again: n/8 + n = 9n/8 misses



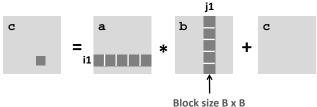
Total misses:

 $9n/8 * n^2 = (9/8) * n^3$

Carnegie Mellon

8 wide

Blocked Matrix Multiplication



n/B blocks

n/B blocks

Cache Miss Analysis

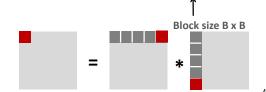
- Assume:
 - Cache block = 8 doubles
 - Cache size C << n (much smaller than n)
 - Three blocks fit into cache: 3B² < C

First (block) iteration:

■ B²/8 misses for each block ■

2n/B * B²/8 = nB/4 (omitting matrix c)

 Afterwards in cache (schematic)



*

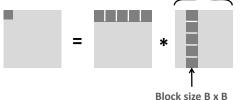
=

Cache Miss Analysis

- Assume:
 - Cache block = 8 doubles
 - Cache size C << n (much smaller than n)
 - Three blocks fit into cache: 3B² < C

Second (block) iteration:

- Same as first iteration
- 2n/B * B²/8 = nB/4



- Total misses:
 - $nB/4 * (n/B)^2 = n^3/(4B)$

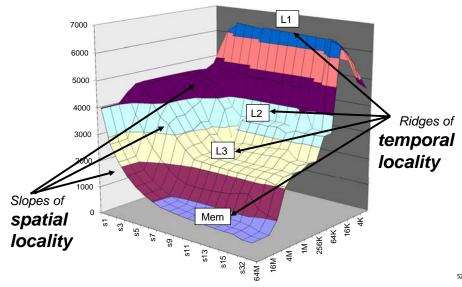
Carnegie Mello

Carnegie Mellon

Summary

- No blocking: (9/8) * n³
- Blocking: 1/(4B) * n³
- Suggest largest possible block size B, but limit 3B² < C!
- Reason for dramatic difference:
 - Matrix multiplication has inherent temporal locality:
 - Input data: 3n², computation 2n³
 - Every array elements used O(n) times!
 - But program has to be written properly

Pay Attention to the Cache!

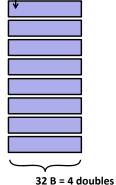


A Higher Level Example

```
int sum_array_rows(double a[16][16])
   int i, j;
   double sum = 0;
   for (i = 0; i < 16; i++)
        for (j = 0; j < 16; j++)
            sum += a[i][j];
   return sum;
```

```
int sum_array_cols(double a[16][16])
   int i, j;
   double sum = 0;
   for (j = 0; i < 16; i++)
       for (i = 0; j < 16; j++)
            sum += a[i][j];
   return sum;
```

Ignore the variables sum, i, j assume: cold (empty) cache, a[0][0] goes here



blackboard

A Higher Level Example

```
int sum_array_rows(double a[16][16])
    int i, j;
   double sum = 0;
   for (i = 0; i < 16; i++)
        for (j = 0; j < 16; j++)
            sum += a[i][j];
   return sum;
```

```
int sum_array_rows(double a[16][16])
    int i, j;
   double sum = 0;
   for (j = 0; i < 16; i++)
        for (i = 0; j < 16; j++)
            sum += a[i][j];
   return sum:
```

Ignore the variables sum, i, j

```
assume: cold (empty) cache,
a[0][0] goes here
      32 B = 4 doubles
```

blackboard

