CACHE MEMORY CS 045

Computer Organization and Architecture

Prof. Donald J. Patterson Adapted from Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

MEMORY HIERARCHIES



EXAMPLES OF CACHING IN THE MEMORY HIERARCHY

Cache Type	What is Cached?	Where is it Cached?	Latency (cycles)	Managed By
Registers	4-8 bytes words	CPU core	0	Compiler
TLB	Address translations	On-Chip TLB	0	Hardware MMU
L1 cache	64-byte blocks	On-Chip L1	4	Hardware
L2 cache	64-byte blocks	On-Chip L2	10	Hardware
Virtual Memory	4-KB pages	Main memory	100	Hardware + OS
Buffer cache	Parts of files	Main memory	100	OS
Disk cache	Disk sectors	Disk controller	100,000	Disk firmware
Network buffer cache	Parts of files	Local disk	10,000,000	NFS client
Browser cache	Web pages	Local disk	10,000,000	Web browser
Web cache	Web pages	Remote server disks	1,000,000,000	Web proxy server

CACHE MEMORY

CACHE MEMORY ORGANIZATION AND OPERATION

- PERFORMANCE IMPACT OF CACHES
 - THE MEMORY MOUNTAIN
 - REARRANGING LOOPS TO IMPROVE
 SPATIAL LOCALITY
 - USING BLOCKING TO IMPROVE
 TEMPORAL LOCALITY

GENERAL CACHE CONCEPT





CACHE MEMORIES

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 cache
- Typical system structure:



GENERAL CACHE ORGANIZATION (S, E, B)



CACHE READ



• Locate set

• Check if any line in set

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





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

DIRECT-MAPPED CACHE SIMULATION



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

Address trace (reads, one byte per read):







E-WAY SET ASSOCIATIVE CACHE (E = 2)

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



Address of short int:

E-WAY SET ASSOCIATIVE CACHE (E = 2)

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



block offset



Address of short int:

E-WAY SET ASSOCIATIVE CACHE (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), ...

2-WAY SET ASSOCIATIVE CACHE SIMULATION



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

Address trace (reads, one byte per read):





WHAT ABOUT WRITES?

Multiple copies of data exist:

L1, L2, L3, 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

INTEL CORE 17 CACHE HIERARCHY

Processor package



WCPKNEEL

Cache Breaks



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:
 - 4 clock cycle for L1
 - 10 clock cycles for L2

Miss Penalty

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

LET'S 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



CACHE MEMORY

CACHE MEMORY ORGANIZATION AND
 OPERATION

PERFORMANCE IMPACT OF CACHES

- THE MEMORY MOUNTAIN
- REARRANGING LOOPS TO IMPROVE
 SPATIAL LOCALITY
- USING BLOCKING TO IMPROVE
 TEMPORAL LOCALITY

CACHE MEMORY

- CACHE MEMORY ORGANIZATION AND
 OPERATION
- PERFORMANCE IMPACT OF CACHES
 THE MEMORY MOUNTAIN
 REARRANGING LOOPS TO IMPROVE SPATIAL LOCALITY
 USING BLOCKING TO IMPROVE TEMPORAL LOCALITY

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

```
long data[MAXELEMS]; /* Global array to traverse */
/* test - Iterate over first "elems" elements of
          array "data" with stride of "stride", using
*
          using 4x4 loop unrolling.
 *
*/
int test(int elems, int stride) {
    long i, sx2=stride*2, sx3=stride*3, sx4=stride*4;
    long acc0 = 0, acc1 = 0, acc2 = 0, acc3 = 0;
    long length = elems, limit = length - sx4;
    /* Combine 4 elements at a time */
    for (i = 0; i < limit; i += sx4) {</pre>
        acc0 = acc0 + data[i];
        acc1 = acc1 + data[i+stride];
        acc2 = acc2 + data[i+sx2];
        acc3 = acc3 + data[i+sx3];
    }
    /* Finish any remaining elements */
    for (; i < length; i++) {</pre>
        acc0 = acc0 + data[i];
    }
    return ((acc0 + acc1) + (acc2 + acc3));
```

Call test() with many combinations of elems and stride.

For each elems and stride:

1. Call test()
once to warm up
the caches.

2. Call test()
again and measure
the read
throughput(MB/s)



CACHE MEMORY

- CACHE MEMORY ORGANIZATION AND
 OPERATION
- PERFORMANCE IMPACT OF CACHES
 THE MEMORY MOUNTAIN
 - REARRANGING LOOPS TO IMPROVE
 SPATIAL LOCALITY
 - USING BLOCKING TO IMPROVE
 TEMPORAL LOCALITY

MATRIX MULTIPLICATION EXAMPLE

Description:

- Multiply N x N matrices
- Matrix elements are doubles (8 bytes)
- O(N³) total operations
- N reads per source element
- N values summed per destination
 - but may be able to hold in register

/* 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;
 }
 matmult/mm.c</pre>

MISS RATE ANALYSIS FOR MATRIX MULTIPLY

Assume:

- Block size = 32B (big enough for four doubles)
- 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



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:

sum += a[0][i];

- accesses successive elements
- if block size (B) > sizeof(a_{ii}) bytes, exploit spatial locality
 - miss rate = sizeof(a_{ii}) / 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%)

MATRIX MULTIPLICATION (IJK)



Misses per inner loop iteration:

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

MATRIX MULTIPLICATION (JIK)



Inner loop:



Misses per inr	<u>ner loop ite</u>	eration:
<u>A</u>	<u>B</u>	<u>C</u>
0.25	1.0	0.0

MATRIX MULTIPLICATION (KIJ)



Inner loop:



Misses per in	<u>ner loop ite</u>	eration:
<u>A</u>	B	<u>C</u>
0.0	0.25	0.25



MATRIX MULTIPLICATION (IKJ)



Misses per in	<u>ner loop ite</u>	eration:
<u>A</u>	B	<u>C</u>
0.0	0.25	0.25



MATRIX MULTIPLICATION (JKI)



Misses per in	<u>ner loop ite</u>	eration:
A	<u>B</u>	<u>C</u>
1.0	0.0	1.0



MATRIX MULTIPLICATION (KJI)



Misses per ini	ner loop ite	eration:
A	<u>B</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;
    }
}</pre>
```

```
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];
  }
}</pre>
```

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**

CORE 17 MATRIX MULTIPLY PERFORMANCE



CACHE MEMORY

- CACHE MEMORY 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







CACHE MISS ANALYSIS

Assume:

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



CACHE MISS ANALYSIS

Assume:

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



BLOCKING SUMMARY

- No blocking: (9/8) * n³
- Blocking: 1/(4B) * n³

Suggest 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 elements used O(n) times!
- But program has to be written properly



CACHE SUMMARY

Cache memories can have significant performance impact

You can write your programs to exploit this!

- Focus on the inner loops, where bulk of computations and memory accesses occur.
- Try to maximize spatial locality by reading data objects with sequentially with stride 1.
- Try to maximize temporal locality by using a data object as often as possible once it's read from memory.

