Memory Hierarchy

Computer Organization and Assembly Languages Yung-Yu Chuang 2007/01/08

with slides by CMU15-213

Reference



 Chapter 6 from "Computer System: A Programmer's Perspective"





- Final project demo time on 1/24 or 1/31. The mechanism for signup will be mailed to the mailing list soon.
- Mail your report to TA before your demo time. The length of report depends on your project type. It can be html, pdf, doc, ppt...
- TA evaluation today

Computer system model



• We assume memory is a linear array which holds both instruction and data, and CPU can access memory in a constant time.





SRAM vs DRAM



	Tran. per bit	Access time	Needs refresh	? Cost	Applications
SRAM	4 or 6	1X	No	100X	cache memories
DRAM	1	10X	Yes	1X	Main memories, frame buffers

The gap widens between DRAM, disk, and CPU speeds.



Memory hierarchies



- Some fundamental and enduring properties of hardware and software:
 - Fast storage technologies cost more per byte, have less capacity, and require more power (heat!).
 - The gap between CPU and main memory speed is widening.
 - Well-written programs tend to exhibit good locality.
- They suggest an approach for organizing memory and storage systems known as a memory hierarchy.





Why does it work?



- Most programs tend to access the storage at any particular level more frequently than the storage at the lower level.
- Locality: tend to access the same set of data items over and over again or tend to access sets of nearby data items.

Why learning it?



- A programmer needs to understand this because the memory hierarchy has a big impact on performance.
- You can optimize your program so that its data is more frequently stored in the higher level of the hierarchy.
- For example, the difference of running time for matrix multiplication could up to a factor of 6 even if the same amount of arithmetic instructions are performed.

Locality



- Principle of Locality: programs tend to reuse data and instructions near those they have used recently, or that were recently referenced themselves.
 - Temporal locality: recently referenced items are likely to be referenced in the near future.
 - Spatial locality: items with nearby addresses tend to be referenced close together in time.
- In general, programs with good locality run faster then programs with poor locality
- Locality is the reason why cache and virtual memory are designed in architecture and operating system. Another example is web browser caches recently visited webpages.

Locality example



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

- Data
 - Reference array elements in succession (stride-1 reference pattern): **Spatial locality**
 - Reference sum each iteration: Temporal locality
- Instructions
 - Reference instructions in sequence: Spatial locality
 - Cycle through loop repeatedly: Temporal locality

Locality example



• Being able to look at code and get a qualitative sense of its locality is important. Does this function have good locality?





Locality example • Does this function have good locality? funt i, j, sum = 0; for (j = 0; j < N; j++) for (i = 0; i < M; i++) sum += a[i][j]; return sum; } stride-N reference pattern</pre>

Cache memories



- Cache memories are small, fast SRAM-based memories managed automatically in hardware.
- CPU looks first for data in L1, then in L2, then in main memory.
- Typical system structure:



Caching in a memory hierarchy



Type of cache misses



- Cold (compulsory) miss: occurs because the cache is empty.
- Capacity miss: occurs when the active cache blocks (working set) is larger than the cache.
- Conflict miss
 - Most caches limit blocks at level k+1 to a small subset of the block positions at level k, e.g. block i at level k+1 must be placed in block (i mod 4) at level k.
 - Conflict misses occur when the level k cache is large enough, but multiple data objects all map to the same level k block, e.g. Referencing blocks 0, 8, 0, 8, 0, 8, ... would miss every time.





 Program needs object d, which is stored in some block b.

Cache hit

- Program finds b in the cache at level k. E.g., block 14.
- Cache miss
 - b is not at level k, so level k cache must fetch it from level k+1.
 E.g., block 12.
 - If level k cache is full, then some current block must be replaced (evicted). Which one is the "victim"?
 - Placement policy: where can the new block go? E.g., b mod 4
 - Replacement policy: which block should be evicted? E.g., LRU

General organization of a cache







Accessing direct-mapped caches



- Line matching and word selection
 - Line matching: Find a valid line in the selected set with a matching tag
 - Word selection: Then extract the word =1? (1) The valid bit must be set



Direct-mapped cache simulation

t=1	s=2	b=1	M=16 S=4 se
X	XX	X	Addre

oyte addresses, B=2 bytes/block,
ets, E=1 entry/set

ess trace (reads): 0 [0000₂], miss 1 [0001₂], hit

- 7 [0111₂], miss 8 [1000₂], miss
- 0 [0000₂] miss

V	tag	data
1	0	M[0-1]

•	[• -]
0	M[6-7]
	0

Accessing direct-mapped caches



- Line matching and word selection
 - Line matching: Find a valid line in the selected set with a matching tag
 - Word selection: Then extract the word



What's wrong with direct-mapped?



float dotprod(float x[8], y[8]) {

float sum=0.0;

for (int i=0; i<8; i++)</pre>

sum+= x[i]*y[i];

return sum;

two sets block size=16 bytes

element	address	set	element	address	set
x[0]	0	0	y[0]	32	0
x[1]	4	0	y[1]	36	0
x[2]	8	0	y[2]	40	0
x[3]	12	0	y[3]	44	0
x[4]	16	1	y[4]	48	1
x[5]	20	1	y[5]	52	1
x[6]	24	1	y[6]	56	1
x[7]	28	1	y[7]	60	1

Solution? padding



float dotprod(float x[12], y[8]) {

float sum=0.0;
for (int i=0; i<8; i++)
 sum+= x[i]*y[i];</pre>

return sum;

ι.						
ſ	element	address	set	element	Address	set
	x[0]	0	0	y[0]	48	1
	x[1]	4	0	y[1]	52	1
	x[2]	8	0	y[2]	56	1
	x[3]	12	0	y[3]	60	1
	x[4]	16	1	y[4]	64	0
	x[5]	20	1	y[5]	68	0
	x[6]	24	1	y[6]	72	0
	x[7]	28	1	y[7]	76	0

Accessing set associative caches

Set selection



Set associative caches



Characterized by more than one line per set



Accessing set associative caches



 Line matching and word selection - must compare the tag in each valid line in the selected set. =1? (1) The valid bit must be set 1001 selected set (i): 0110 $\mathbf{w}_0 \mid \mathbf{w}_1 \mid \mathbf{w}_2 \mid$ Wa (2) The tag bits in one of the cache lines If (1) and (2), then cache hit = 2 must match the tag bits in the address t bits s bits b bits 0110 100 m-1 set index block offset⁰ tag

Accessing set associative caches



- Line matching and word selection
 - Word selection is the same as in a direct mapped cache



Why use middle bits as index?

4-line Cache	l	High-Order Bit Indexing	, M	Aiddle-Order Bit Indexing	r
00	0000		0000		
01	0001		0001		
10	0010		0010		
11	0011		0011		
	0100		0100		
High-order bit	0101		0101		
indexing	0110		0110		
adjacent memory lines	0111		0111		
would man to same	<u>10</u> 00		10 <u>00</u>		
sacha antry	<u>10</u> 01		10 <u>01</u>		
	1010		1010		
- poor use of spatial	1011		1011		
ocality	1100		1100		
	1101		1101		
	1110		1110		
	1111		$11\overline{11}$		

2-Way associative cache simulation



- - Write-through
 - Write-back (need a dirty bit)
- What to do on a replacement?
 - Depends on whether it is write through or write back

Multi-level caches



 Options: separate data and instruction caches, or a unified cache



Intel Pentium III cache hierarchy





The memory mountain



- Read throughput: 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.

```
void test(int elems, int stride) {
    int i, result = 0;
    volatile int sink;
    for (i = 0; i < elems; i += stride)
        result += data[i];
    /* So compiler doesn't optimize away the loop */
    sink = result;</pre>
```

The memory mountain



```
/* Run test(elems, stride) and return read
    throughput (MB/s) */
double run(int size, int stride, double Mhz)
{
    double cycles;
    int elems = size / sizeof(int);
    test(elems, stride); /* warm up the cache */
    /* call test(elems,stride) */
    cycles = fcyc2(test, elems, stride, 0);
    /* convert cycles to MB/s */
    return (size / stride) / (cycles / Mhz);
}
```



A slope of spatial locality



Slice through memory mountain (size=256KB)
 shows cache block size.



Ridges of temporal locality

- Slice through the memory mountain (stride=1)
 - illuminates read throughputs of different caches and

memory L2 is unified



Matrix multiplication example



held in register

- Major cache effects to consider
 - Total cache size
 - Exploit temporal locality and keep the working set small (e.g., use blocking)

/* ijk */

3

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

sum = 0.0:

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

c[i][j] = sum;

for (k=0; k<n; k++)</pre>

sum += a[i][k] * b[k][j];

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

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 (ijk)





Misses per Inner Loop Iteration:

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

Matrix multiplication (jik)



(i,j)



Matrix multiplication (kij)







Misses per Inner Loop Iteration:				
A	B	<u>C</u>		
0.0	0.25	0.25		

Matrix multiplication (ikj)





Misses per Inner Loop Iteration:

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



Row-wise Row-wise

Matrix multiplication (jki)





Matrix multiplication (kji)







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



Pentium matrix multiply performance

- Miss rates are helpful but not perfect predictors.
- Code scheduling matters, too.



Blocked matrix multiply (bijk)



```
for (jj=0; jj<n; jj+=bsize) {
  for (i=0; i<n; i++)
    for (j=jj; j < min(jj+bsize,n); j++)
        c[i][j] = 0.0;
  for (kk=0; kk<n; kk+=bsize) {
    for (i=0; i<n; i++) {
        for (j=jj; j < min(jj+bsize,n); j++) {
            sum = 0.0
            for (k=kk; k < min(kk+bsize,n); k++) {
                sum += a[i][k] * b[k][j];
            }
            c[i][j] += sum;
            }
        }
    }
}</pre>
```

Improving temporal locality by blocking

- Example: Blocked matrix multiplication
 - Here, "block" does not mean "cache block".
 - Instead, it mean a sub-block within the matrix.
 - Example: N = 8; sub-block size = 4



<u>Key idea:</u> Sub-blocks (i.e., \mathbf{A}_{xy}) can be treated just like scalars.

$$C_{11} = A_{11}B_{11} + A_{12}B_{21} \qquad C_{12} = A_{11}B_{12} + A_{12}B_{22}$$

$$C_{21} = A_{21}B_{11} + A_{22}B_{21} \qquad C_{22} = A_{21}B_{12} + A_{22}B_{22}$$

Blocked matrix multiply analysis



- Innermost loop pair multiplies a 1 X bsize sliver of A by a *bsize X bsize* block of *B* and accumulates into 1 X bsize sliver of C - Loop over *i* steps through *n* row slivers of *A* & *C*, using same B for (i=0; i<n; i++) {</pre> for (j=jj; j < min(jj+bsize,n); j++) {</pre> sum = 0.0for (k=kk; k < min(kk+bsize,n); k++) {</pre> sum += a[i][k] * b[k][j]; c[i][j] += sum; Innermost Loop Pair Update successive row sliver accessed elements of sliver bsize times block reused n times in succession

Blocked matrix multiply performance



 Blocking (bijk and bikj) improves performance by a factor of two over unblocked versions (ijk and jik)



Cache-conscious programming



• make sure that memory is cache-aligned



- Use union and bitfields to reduce size and increase locality
- Split data into hot and cold (list example)







- •Repeated references to variables are good (temporal locality)
- •Stride-1 reference are good (spatial locality)
- •Examples: cold cache, 4-byte words, 4-word cache blocks



```
Miss rate = 1/4 = 25%
```

```
Miss rate =100%
```

Cache-conscious programming



- Prefetching
- Blocked 2D array





