CS 201

Writing Cache-Friendly Code

Gerson Robboy Portland State University

An Example Memory Hierarchy



Why Cache-Friendly Code is Important

Cache type	Size of item (bytes)	Latency (cpu cycles)
Registers	4 bytes	0
L1 Cache	32 bytes	1
L2 Cache	32 bytes	10
Main Memory	4-KB pages	100
Disk		millions

On ia32 processor, with few registers, even local

variables are likely to spill to memory.

We want them in cache!

Just what does a cache do?

The cache stores memory in units or *cache lines*

- Fixed length chunks, hardware dependent
- For our example, let's say cache lines are 32 bytes
- Aligned on a cache-line (32 byte) boundary

When the CPU accesses a memory address (store or load), the cache line containing that address is pulled into the cache

Examples

Suppose a certain processor has a 32-byte cache line size.

You access address 0x3a40. What addresses are pulled into the cache?

You access address 0x3a94. What addresses are pulled into the cache?

Next you access 0x3a48. What happens?

You access 4 32-bit words sequentially, from 0x8000 to 0x801c

How many cache misses and how many cache hits?

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.

Locality Example:

- 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

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



Claim: Being able to look at code and get a qualitative sense of its locality is a key skill for a professional programmer.

Question: Does this function have good locality?

Spatial, temporal, both, or neither?

```
int sumarrayrows(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>
```

Locality Example

Question: Does this function have good locality?

Spatial, temporal, both, or neither?

```
int sumarraycols(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>
```

Locality Example

Question: Can you permute the loops so that the function scans the 3-d array a [] with a stride-1 reference pattern (and thus has good spatial locality)?

```
int sumarray3d(int a[M][N][N])
{
    int i, j, k, sum = 0;
    for (i = 0; i < M; i++)
        for (j = 0; j < N; j++)
            for (k = 0; k < N; k++)
                sum += a[k][i][j];
    return sum
}</pre>
```

Why does traversing a matrix with stride 1 give you good spatial locality?

Why do strides other than 1 give you bad spatial locality?

Writing Cache Friendly Code

Repeated references to variables are good (temporal locality)

Stride-1 reference patterns are good (spatial locality)

Examples:

cold cache, 4-byte words, 8-word cache blocks

```
int sumarrayrows(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>
```

```
Miss rate = 1/8 = 12.5\%
```

```
int sumarraycols(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>
```

Miss rate = 100%

15-213, F'02

- 11 -

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

```
/* The test function */
void test(int elems, int stride) {
    int i, result = 0;
   volatile int sink;
    for (i = 0; i < elems; i += stride)</pre>
        result += data[i];
    sink = result; /* So compiler doesn't optimize away the loop */
}
/* 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);
                                             /* warm up the cache */
    test(elems, stride);
    cycles = fcyc2(test, elems, stride, 0); /* call test(elems, stride) */
    return (size / stride) / (cycles / Mhz); /* convert cycles to MB/s */
}
```

Memory Mountain Main Routine

```
/* mountain.c - Generate the memory mountain. */
#define MINBYTES (1 << 10) /* Working set size ranges from 1 KB */</pre>
#define MAXBYTES (1 << 23) /* ... up to 8 MB */
#define MAXSTRIDE 16
                       /* Strides range from 1 to 16 */
#define MAXELEMS MAXBYTES/sizeof(int)
int main()
{
   int stride; /* Stride (in array elements) */
   double Mhz;
                 /* Clock frequency */
   init data(data, MAXELEMS); /* Initialize each element in data to 1 */
   Mhz = mhz(0);
                          /* Estimate the clock frequency */
   for (size = MAXBYTES; size >= MINBYTES; size >>= 1) {
      for (stride = 1; stride <= MAXSTRIDE; stride++)</pre>
          printf("%.1f\t", run(size, stride, Mhz));
      printf("\n");
   }
   exit(0);
```

- 14 -

The Memory Mountain



- 15 -

Ridges of Temporal Locality

- 16 -

Slice through the memory mountain with stride=1

 illuminates read throughputs of different caches and memory



A Slope of Spatial Locality

Slice through memory mountain with size=256KB

shows cache block size.



– 17 –

Matrix Multiplication Example

Major Cache Effects to Consider

- Total cache size
 - Exploit temporal locality and keep the working set small (e.g., by using blocking)
- 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 4 64-bit words)
- Matrix dimension (N) is very large
 - Approximate 1/N as 0.0
- Cache is not 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:

- accesses successive elements
- if block size (B) > 4 bytes, exploit spatial locality
 - compulsory miss rate = 4 bytes / B

Stepping through rows in one column:

- for (i = 0; i < n; i++)
 sum += a[i][0];</pre>
- accesses distant elements
- no spatial locality!
 - compulsory miss rate = 1 (i.e. 100%)

Matrix Multiplication (ijk)





Misses per Inner Loop Iteration:

A	<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
    }
}</pre>
```



Misses per	Inner Loop	Iteration:
<u>A</u>	<u>B</u>	<u>C</u>
0.25	1.0	0.0

Matrix Multiplication (kij)





Misses per	Inner Loo	p Iteration:
Α	<u>B</u>	<u>C</u>
0.0	0.25	0.25

Matrix Multiplication (ikj)





Misses per	Inner Loo	p Iteration:
Α	<u>B</u>	<u>C</u>
0.0	0.25	0.25

Matrix Multiplication (jki)



Matrix Multiplication (kji)



Summary of Matrix Multiplication

 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
for (i=0; i <n; i++)="" td="" {<=""><td>for (k=0; k<n; k++)="" td="" {<=""><td>for (j=0; j<n; j++)="" td="" {<=""></n;></td></n;></td></n;>	for (k=0; k <n; k++)="" td="" {<=""><td>for (j=0; j<n; j++)="" td="" {<=""></n;></td></n;>	for (j=0; j <n; j++)="" td="" {<=""></n;>
for (j=0; j <n; j++)="" td="" {<=""><td>for (i=0; i<n; i++)="" td="" {<=""><td>for (k=0; k<n; k++)="" td="" {<=""></n;></td></n;></td></n;>	for (i=0; i <n; i++)="" td="" {<=""><td>for (k=0; k<n; k++)="" td="" {<=""></n;></td></n;>	for (k=0; k <n; k++)="" td="" {<=""></n;>
sum = 0.0;	r = a[i][k];	r = b[k][j];
for (k=0; k <n; k++)<="" td=""><td>for (j=0; j<n; j++)<="" td=""><td>for (i=0; i<n; i++)<="" td=""></n;></td></n;></td></n;>	for (j=0; j <n; j++)<="" td=""><td>for (i=0; i<n; i++)<="" td=""></n;></td></n;>	for (i=0; i <n; i++)<="" td=""></n;>
sum += a[i][k] * b[k][j];	c[i][j] += r * b[k][j];	c[i][j] += a[i][k] * r;
c[i][j] = sum;	}	}
}	}	}
}		

Pentium Matrix Multiply Performance

Miss rates are helpful but not perfect predictors.

• Code scheduling matters, too.



Improving Temporal Locality by Blocking

Example: Blocked matrix multiplication

- "block" (in this context) does not mean "cache block".
- Instead, it mean a sub-block within the matrix.
- Example: N = 8; sub-block size = 4

$$\begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} X \begin{bmatrix} B_{11} & B_{12} \\ B_{21} & B_{22} \end{bmatrix} = \begin{bmatrix} C_{11} & C_{12} \\ C_{21} & C_{22} \end{bmatrix}$$

<u>Key idea:</u> Sub-blocks (i.e., 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 (bijk)

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

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*



Let's try to see what this does

```
for (jj=0; jj<n; jj+=bsize) {// for each bsize block</pre>
  //skip zeroing C for now
  for (kk=0; kk<n; kk+=bsize) {</pre>
    for (i=0; i<n; i++) \{ // for each row of A
      // for each column of the block of B
      for (j=jj; j < min(jj+bsize,n); j++) {</pre>
        sum = 0.0
        // For each element of the sliver of A/column of B
        for (k=kk; k < min(kk+bsize,n); k++) {
          sum += a[i][k] * b[k][j];
        }
        c[i][j] += sum;
      }
    }
  }
                               — kk → — ii →
}
                                      kk
                               l i
                                          B
                                                   С
                                 Α
```

15-213, F'02

So here's the point of blocking

- Use a block size smaller than the size of the CPU cache
- The row sliver and the block in B are re-used many times in a row.
- They are in cache after the first time they are used.
- Then go on to another small block, get it in the cache.
- If you do it in the right order, you multiply all the horizontal slivers in A times one block in B, before going on to another block in B.



Pentium Blocked Matrix Multiply Performance

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

relatively insensitive to array size.



Concluding Observations

Programmer can optimize for cache performance

- How data structures are organized
- How data are accessed
 - 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)