

CSE 351 Spring 2019

Instructor:	Teaching Assis	tants:	
Ruth Anderso	on Gavin Cai Britt Henderson Sophie Tian Casey Xing	Jack Eggleston n Richard Jiang Connie Wang Chin Yeoh	John Feltrup Jack Skalitzky Sam Wolfson
WHAT'S THIS? HUH?	"I ALWAYS THOUGHT "THE "WAS A HUGE, AMORPHOUS	HOW? YOU'RE ON SHOULD THE CORD BE A CABLE MODEM. STRETCHED ACROSS	WHAT IF SOMEONE TRIPS ON IT? / WHO WOULD WANT TO DO THAT?



http://xkcd.com/908/

Administrivia

- Lab 3, due TONIGHT, Wednesday (5/15)
- Homework 4 , due Wed (5/22) (Structs, Caches)
- Lab 4, Coming soon!
 - Cache parameter puzzles and code optimizations



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

Memory



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. 3

mov OxFACE, F







mov OxFACE, F



Step 1: Bring **F** into cache





mov OxFACE, F



Step 2: Write OxFACE to cache only and set dirty bit

Memory



mov 0xFACE, F mov 0xFEED, F





Memory



















Peer Instruction Question

- Which of the following cache statements is FALSE?
 - Vote at <u>http://pollev.com/rea</u>
 - A. We can reduce compulsory misses by decreasing our block size
 - **B.** We can reduce conflict misses by increasing associativity
 - C. A write-back cache will save time for code with good temporal locality on writes
 - **D.** A write-through cache will always match data with the memory hierarchy level below it
 - E. We're lost...

Optimizations for the Memory Hierarchy

- Write code that has locality!
 - Spatial: access data contiguously
 - <u>Temporal</u>: make sure access to the same data is not too far apart in time
- How can you achieve locality?
 - Adjust memory accesses in *code* (software) to improve miss rate (MR)
 - Requires knowledge of *both* how caches work as well as your system's parameters
 - Proper choice of algorithm
 - Loop transformations

Example: Matrix Multiplication



Matrices in Memory

- How do cache blocks fit into this scheme?
 - Row major matrix in memory:



Naïve Matrix Multiply





Cache Miss Analysis (Naïve)



- Scenario Parameters:
 - Square matrix (n × n), elements are doubles
 - Cache block size K = 64 B = 8 doubles
 - Cache size C << n (much smaller than n)



Cache Miss Analysis (Naïve)



- Scenario Parameters:
 - Square matrix (n × n), elements are doubles
 - Cache block size K = 64 B = 8 doubles
 - Cache size C << n (much smaller than n)



Cache Miss Analysis (Naïve)



- Scenario Parameters:
 - Square matrix (n × n), elements are doubles
 - Cache block size K = 64 B = 8 doubles
 - Cache size C << n (much smaller than n)



Linear Algebra to the Rescue (1)



- Can get the same result of a matrix multiplication by splitting the matrices into smaller submatrices (matrix "blocks")
- For example, multiply two 4×4 matrices:

$$A = \begin{bmatrix} a_{11} & a_{12} & a_{13} & a_{14} \\ a_{21} & a_{22} & a_{23} & a_{24} \\ a_{31} & a_{32} & a_{33} & a_{34} \\ a_{41} & a_{42} & a_{43} & a_{44} \end{bmatrix} = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}, \text{ with } B \text{ defined similarly.}$$
$$AB = \begin{bmatrix} (A_{11}B_{11} + A_{12}B_{21}) & (A_{11}B_{12} + A_{12}B_{22}) \\ (A_{21}B_{11} + A_{22}B_{21}) & (A_{21}B_{12} + A_{22}B_{22}) \end{bmatrix}$$

Linear Algebra to the Rescue (2)





Matrices of size $n \times n$, split into 4 blocks of size r (n=4r)

$$C_{22} = A_{21}B_{12} + A_{22}B_{22} + A_{23}B_{32} + A_{24}B_{42} = \sum_{k} A_{2k}^{*}B_{k2}$$

- Multiplication operates on small "block" matrices
 - Choose size so that they fit in the cache!
 - This technique called "cache blocking"

Blocked Matrix Multiply

Blocked version of the naïve algorithm:

```
# move by rxr BLOCKS now
for (i = 0; i < n; i += r)
for (j = 0; j < n; j += r)
for (k = 0; k < n; k += r)
    # block matrix multiplication
    for (ib = i; ib < i+r; ib++)
       for (jb = j; jb < j+r; jb++)
       for (kb = k; kb < k+r; kb++)
            c[ib*n+jb] += a[ib*n+kb]*b[kb*n+jb];
```

r = block matrix size (assume r divides n evenly)

Cache Miss Analysis (Blocked)

- Scenario Parameters:
 - Cache block size K = 64 B = 8 doubles
 - Cache size C << n (much smaller than n)
 - Three blocks \blacksquare ($r \times r$) fit into cache: $3r^2 < C$

 r^2 elements per block, 8 per cache block

- Sech block iteration:
 - $r^{2}/8$ misses per block

$$2n/r \times r^2/8 = nr/4$$







n/r blocks in row and column

Cache Miss Analysis (Blocked)

- Scenario Parameters:
 - Cache block size K = 64 B = 8 doubles
 - Cache size $C \ll n$ (much smaller than n)
 - Three blocks \blacksquare ($r \times r$) fit into cache: $3r^2 < C$

 r^2 elements per block, 8 per cache block

- Sech block iteration:
 - $r^{2}/8$ misses per block

$$2n/r \times r^2/8 = nr/4$$

n/r blocks in row and column

 Afterwards in cache (schematic)





X



n/r blocks

Ignoring

matrix

Cache Miss Analysis (Blocked)

- Scenario Parameters:
 - Cache block size K = 64 B = 8 doubles
 - Cache size C << n (much smaller than n)
 - Three blocks \blacksquare ($r \times r$) fit into cache: $3r^2 < C$

 r^2 elements per block, 8 per cache block

- Sech block iteration:
 - $r^2/8$ misses per block

$$2n/r \times r^2/8 = nr/4$$





n/r blocks in row and column

- Total misses:
 - $nr/4 \times (n/r)2 = \frac{n^3}{(4r)}$

Matrix Multiply Visualization



Cache-Friendly Code

- 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 size, cache block size, associativity, 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



Learning About Your Machine

& Linux:

- lscpu
- Is /sys/devices/system/cpu/cpu0/cache/index0/
 - <u>Ex</u>: cat /sys/devices/system/cpu/cpu0/cache/index*/size

Windows:

- wmic memcache get <query> (all values in KB)
- <u>Ex</u>: wmic memcache get MaxCacheSize
- Modern processor specs: <u>http://www.7-cpu.com/</u>