

Module 2 Part 2

High Performance Application Optimization on CPUs

Carnegie Mellon University
18-645

Jike Chong
Ian Lane
With slides from Yevgen Voronenko

18-645 – How to Write Fast Code? 1



What we discussed last time:

Fast Platforms

- Multicore platforms
- Manycore platforms
- Cloud platforms

+

Good Techniques

- Data structures
- Algorithms
- Software Architecture

- Highlighted the difference between multicore and manycore platforms
- **Exposing** concurrency in k-means, **Exploiting** parallelism by **exploring** mappings from application to platform
- OpenMP: **An abstraction for multicore parallel programming**
- The **choice and responsibilities** we all have to influence the industry!
- This lecture:
 - Optimization techniques on a modern CPU
 - Mapping of a simple problem of matrix-matrix multiplication to a complex CPU

18-645 – How to Write Fast Code? Carnegie Mellon University (c) 2012-2014 2

Answers you should know after this...

- Why does naïve matrix-multiply does not achieve peak performance on the CPU?
- What are the different data layouts for matrices?
- Is blocking sufficient?
- What can be learned from this for other computations?

Outline

- CPU achievable peak performance
- Matrix-multiplication discussion
- Optimization techniques

Peak Single-precision Performance

- Intel Sandy Bridge 3.5 Ghz, 4 cores * 2 sockets

[GFLOPS]	1 core	2 cores	4 cores	8 cores
None, scalar	7	14	28	56
SSE	28	56	112	224
AVX	56	112	224	448

- NVIDIA Tesla C2xxx
1030 GFLOPs
- NVIDIA Quadro 5000
722 GFLOPs

Most CPU applications achieve a fraction of this

In addition ...

- Deep memory hierarchy
 - Registers
 - L1 cache
 - L2 cache
 - TLB (Translation Lookaside Buffer)
- Other forms of parallelism
 - Pipelining
 - Instruction-level parallelism (multiple functional units)
- “Free operations” are not free

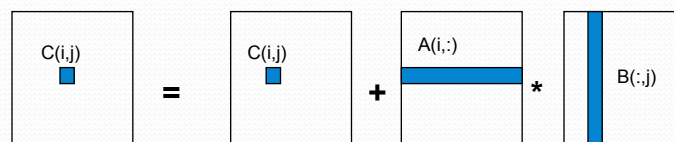
Outline

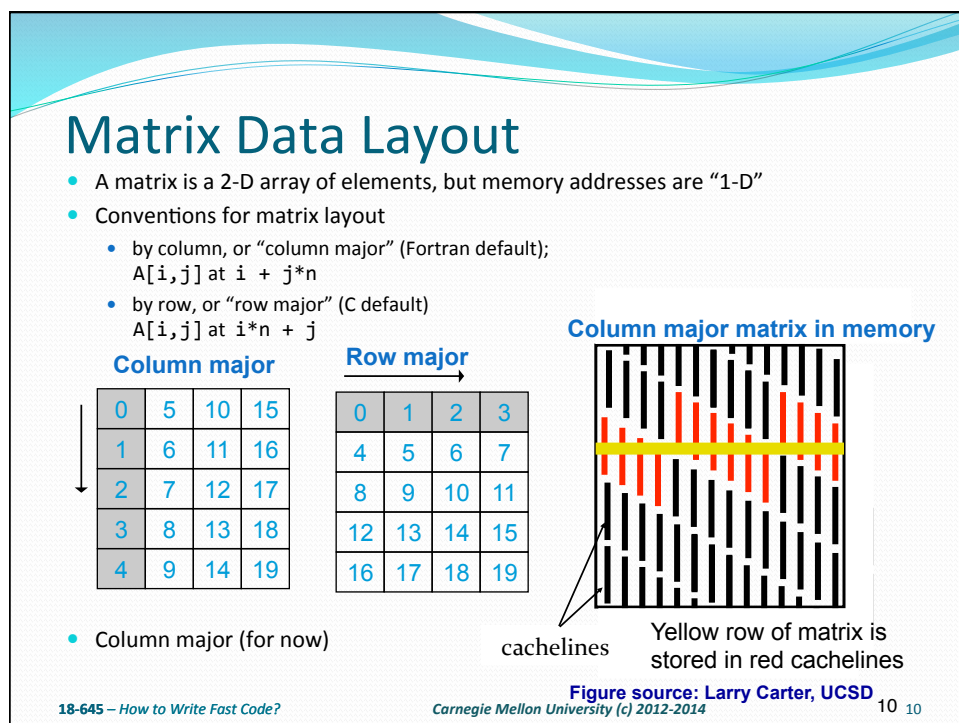
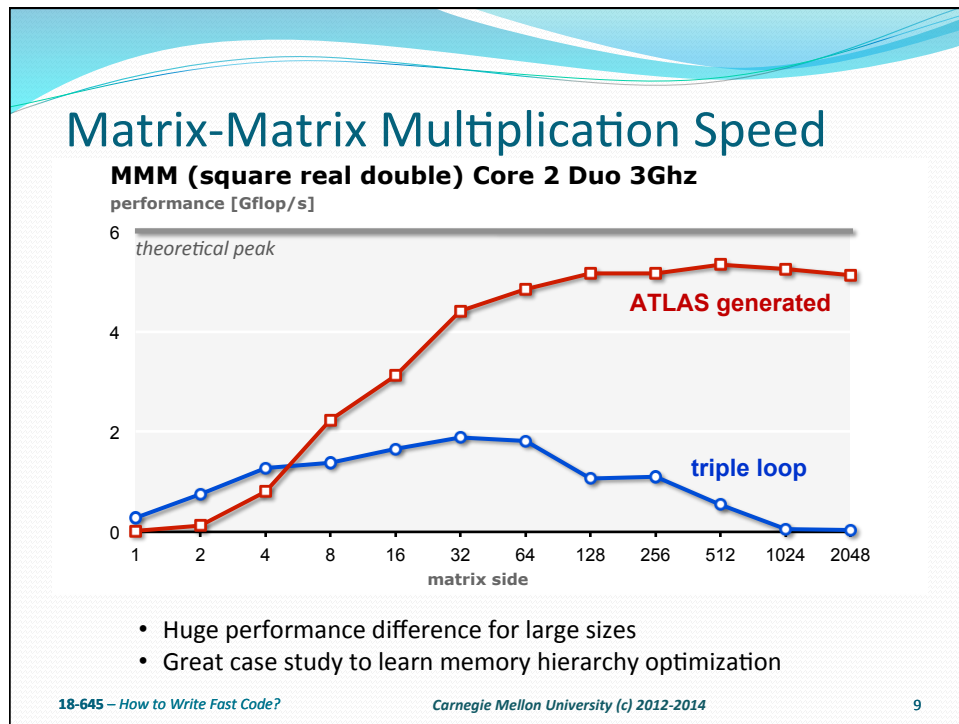
- CPU achievable peak performance
- Matrix-multiplication discussion
- Optimization techniques

Naïve Matrix Multiply

{implements $C = C + A * B$: Affine Transformation}

```
for i = 1 to n
  for j = 1 to n
    for k = 1 to n
      C(i,j) = C(i,j) + A(i,k) * B(k,j)
```



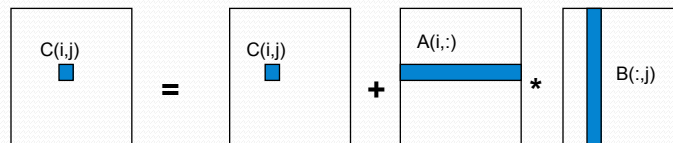


Naïve Matrix Multiply

```

{implements C = C + A*B}
for i = 1 to n
  {read row i of A into fast memory}
  for j = 1 to n
    {read C(i,j) into fast memory}
    {read column j of B into fast memory} // O(n^3) slow loads
    for k = 1 to n
      C(i,j) = C(i,j) + A(i,k) * B(k,j)
    {write C(i,j) back to slow memory}

```



18-645 – How to Write Fast Code?

Carnegie Mellon University (c) 2012-2014

11

Naïve Matrix Multiply

Number of slow memory references on unblocked matrix multiply

$m = n^3$ to read each column of B n times
 $+ n^2$ to read each row of A once
 $+ 2n^2$ to read and write each element of C once
 $= n^3 + 3n^2$

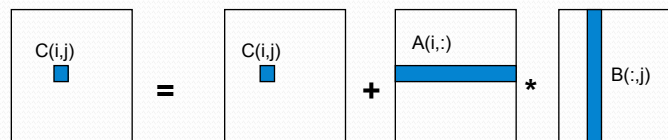
Arithmetic intensity: Total FLOPS/ Total DRAM Bytes

In single precision: $2n^3 / 4*(n^3 + 3n^2)$

≈ 0.5 for large n , no improvement over matrix-vector multiply

Inner two loops are just matrix-vector multiply, of row i of A times B

Similar for any other order of 3 loops

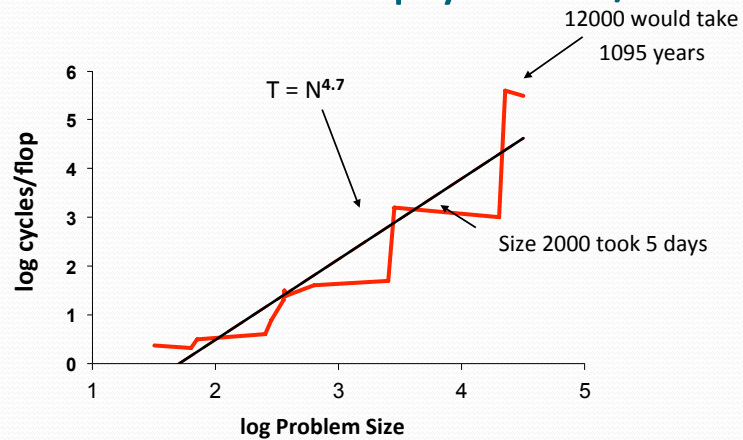


18-645 – How to Write Fast Code?

Carnegie Mellon University (c) 2012-2014

12 12

Naïve Matrix Multiply on RS/6000

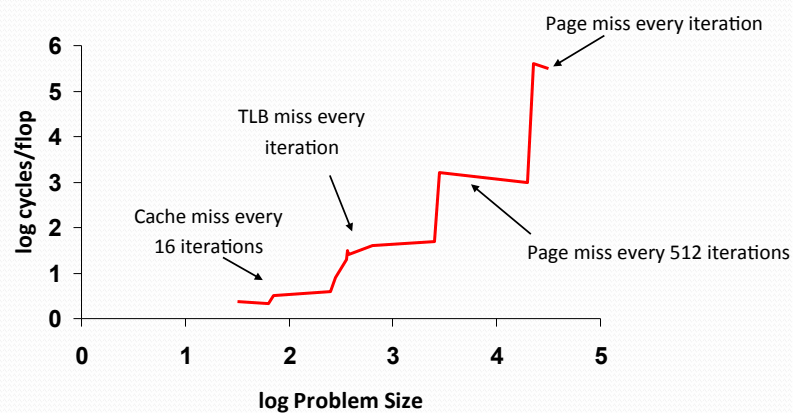


$O(N^3)$ performance would have constant cycles/flop

18-645 – How to Write Fast Code?

Slide source: Larry Carter, UCSD 13

Naïve Matrix Multiply on RS/6000



18-645 – How to Write Fast Code?

Slide source: Larry Carter, UCSD

14

Blocked (Tiled) Matrix Multiply

$A, B, C = N \times N$ matrices of $b \times b$ sub-blocks, $N = n / b$

for $i = 1$ to N

for $j = 1$ to N

{read block $C(i,j)$ into fast memory}

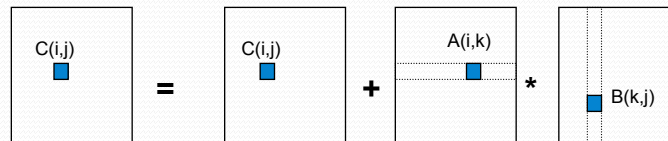
for $k = 1$ to N

{read block $A(i,k)$ into fast memory} // $O(N^3)$ slow loads

{read block $B(k,j)$ into fast memory}

$C(i,j) = C(i,j) + A(i,k) * B(k,j)$ {do a matrix multiply on blocks}

{write block $C(i,j)$ back to slow memory}



18-645 – How to Write Fast Code?

Carnegie Mellon University (c) 2012-2014

15

Working with a Cache Hierarchy

- Illustration by Vasily Volkov on YouTube
- Naïve Matrix Multiply:
 - http://www.youtube.com/watch?v=j5_JU5rdEi8&NR=1
- Matrix Multiply with cache blocking:
 - <http://www.youtube.com/watch?v=TveIT9Bz6EU&NR=1>
- Multi-level cache blocking:
 - <http://stumptown.cc.gt.atl.ga.us/cse6230-hpcta-fa11/slides/11a-matmul-goto.pdf>

18-645 – How to Write Fast Code?

Carnegie Mellon University (c) 2012-2014

16

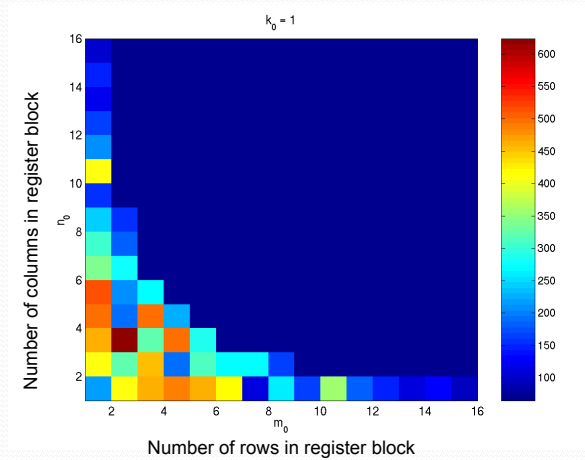
Outline

- CPU achievable peak performance
- Matrix-multiplication discussion
- Optimization techniques

Further Matrix-Multiply Optimizations in Real-world Libraries

- Block size adaptation for appropriate caches
- Register-level blocking
- Copy optimization (data layout)
- Optimizing the mini-matrix-multiply (base case)
- Multi-level blocking
- Multi-level copying

Register-level blocking



A 2-D slice of a 3-D register-tile search space. The dark blue region was pruned.

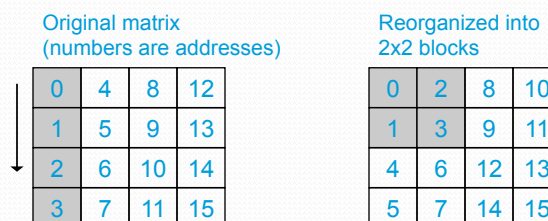
Slide source: Jim Demmel, UCB

18-645 – How to Write Fast Code?

19

Copy optimization

- Copy input operands or blocks
 - Reduce cache conflicts
 - Constant array offsets for fixed size blocks
 - Expose page-level locality
 - Alternative: use different data structures from start (if users willing)



Slide source: Jim Demmel, UCB

18-645 – How to Write Fast Code?

Carnegie Mellon University (c) 2012-2014

20

Optimizing the Base Case

- Base case is usually the guts of the algorithm that computes on data in caches
- Should fit in registers
- Must use SIMD
- **Can be guided by a simple cost model**

Cost Model

- **Algebraic model**
+, -, /, *
sin, cos, etc.
- **Realistic model**
 - +, -, /, *
 - std libs (sin, cos, etc)
 - array[i]
 - A.sum, A.prod
 - (double) x
 - func(...)
 - i < n, compute
 - i < n, test and branch
 - {} unconditional branch

```
int func(int i) {
    return i+2;
}

struct A {
    double sum, prod;
};

A acc(int n, int array[]) {
    A res = {0.0, 1.0};
    for(int i=0; i<n; ++i) {
        A.sum = A.sum +
            (double)func(array[i]);
        A.prod = A.prod *
            (double)func(array[i]);
    }
    return res;
}
```

Cost Model (cont'd)

- Algebraic model

cost = **3 n**

```
int func(int i) {
    return i+2;
}
```

- Realistic model

+ and *
2 int -> double
4 a.x
++i
i<n conditional jump
2 func() calls

```
struct A {
    double sum, prod;
};

A acc(int n) {
    A res = {0.0, 1.0};
    for(int i=0; i<n; ++i) {
        A.sum = A.sum + (double)func(i);
        A.prod = A.prod * (double)func(i);
    }
    return res;
}
```

cost = **19 n**

6x discrepancy in cost

18-645 – How to Write Fast Code?

Carnegie Mellon University (c) 2012-2014

23

Using the Cost Model to Optimize

```
int func(int i) { return i+2; }
```

```
struct A { double sum, prod; };
```

```
A acc(int n, int array[]) {
```

```
    A res;
```

```
    double sum=0.0, prod=1.0;
```

Cost = 6n

```
    for(int i=0; i<n; ++i) {
```

```
        f = (double)(i+2);
```

```
        sum = sum + f;
```

```
        prod = prod * f;
```

```
    }
```

```
    return res;
```

```
}
```

Cost Reduction = closing down the gap

18-645 – How to Write Fast Code?

Carnegie Mellon University (c) 2012-2014

24

Using the Cost Model to Optimize II

```
int func(int i) { return i+2; }

struct A { double sum, prod; };

A acc(int n, int array[]) {
    A res;
    double sum=0.0, prod=1.0;

    for(int i=2; i<n+2; ++i) {
        f = (double)(i);
        sum = sum + f;
        prod = prod * f;
    }
    return res;
}
```

Cost = $5n$

Optimizations enable new optimizations

18-645 – How to Write Fast Code?

Carnegie Mellon University (c) 2012-2014

25

Using the Cost Model to Optimize III

```
int func(int i) { return i+2; }

struct A { double sum, prod; };

A acc(int n, int array[]) {
    A res;
    double sum=0.0, prod=1.0;
    #pragma unroll(4)
    for(int i=2; i<n+2; ++i) {
        f = (double)(i);
        sum = sum + f;
        prod = prod * f;
    }
    return res;
}
```

Cost = $3.5n$

Final result: $3.5n$ vs $19n$

18-645 – How to Write Fast Code?

Carnegie Mellon University (c) 2012-2014

26

Cost Reduction: A Compiler Problem

- **Solution:** “human compiler”
- Compilers often fail on complex codes (many assumptions are violated)
- Optimizations
 - Strength reduction (as already shown)
 - Function inlining
 - Loop unrolling
 - Common subexpression elimination
 - Load/store elimination
 - Table lookups
 - Branch elimination

18-645 – How to Write Fast Code?

Carnegie Mellon University (c) 2012-2014

27

More Strength Reduction

```
for i = 1..n
    sum = sum + a[i] / c;
```



```
double c_inv = 1/c;
for i = 1..n
    sum = sum + a[i] * c_inv;
```

Expensive operations:
/, %, sin, cos, log

18-645 – How to Write Fast Code?

Carnegie Mellon University (c) 2012-2014

28

Function Inlining

```
double f(a) { return a/c; }  
...  
for i = 1..n  
    sum = sum + f(a[i])
```



```
for i = 1..n  
    sum = sum + a[i] / c;
```

Now strength reduction can apply

Loop Unrolling

```
for i = 1..n  
    sum = sum + a[i]
```



```
for i = 1:4:n  
    sum = sum + a[i] + a[i+1] + a[i+2] + a[i+3]
```

```
// remaining iterations  
for i = i:n  
    sum = sum + a[i]
```

Often enables other optimizations

Common Sub-expression Elimination

```
for j = 1..m
  for i = 1..n
    sum = sum + cos(j*PI/180)*a[i];
```



```
for j = 1..m
  double c = cos(j*PI/180);
  for i = 1..n
    sum = sum + c*a[i];
```

Table Lookups

```
for j = 1..m
  for i = 1..n
    sum = sum + cos(j*PI/180)*a[i];
```



```
for j = 1..m
  double c = COS_TAB[j];
  for i = 1..n
    sum = sum + c*a[i];
```


Load/Store Elimination Loop Merging

```
for i = 1..n
  sum = sum + f(a[i]);
for i = 1..n
  prod = prod * f(a[i]);
```



```
for i = 1..n
  sum = sum + f(a[i]);
  prod = prod * f(a[i]);
```

One of the most important optimizations!
Almost always enables other optimizations.

18-645 – How to Write Fast Code?

Carnegie Mellon University (c) 2012-2014

33

Branch Elimination

```
for i = 1..n
  if(i != except)
    sum = sum + a[i];
```



```
for i = 1..except-1
  sum = sum + a[i];
for i = except+1..n
  sum = sum + a[i];
```



```
a[except] = 0.0;
for i = n
  sum = sum + a[i];
```

18-645 – How to Write Fast Code?

Carnegie Mellon University (c) 2012-2014

34

Can You Answer These Questions Now?

- Why does naïve matrix-multiply does not achieve peak performance on the CPU?
- Is blocking sufficient?
- What can be learned from this for other computations?