
Memory Hierarchies and Optimizations: Case Study in Matrix Multiplication

Kathy Yelick

yelick@cs.berkeley.edu

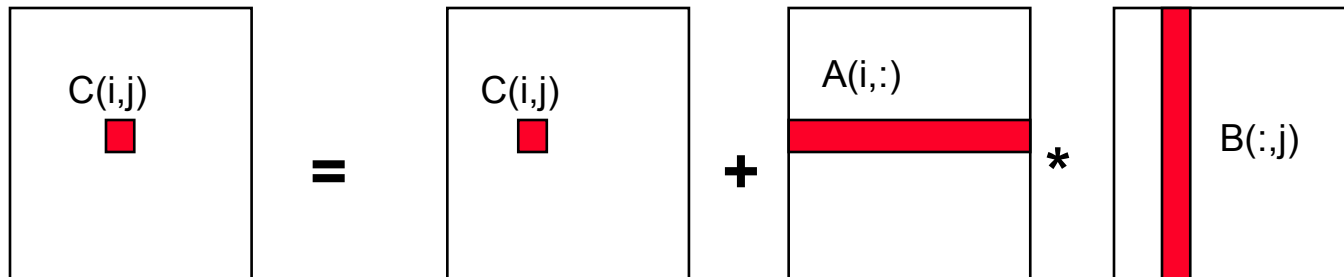
www.cs.berkeley.edu/~yelick/cs194f07

Naïve Matrix Multiply

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

Algorithm has $2*n^3 = O(n^3)$ Flops and operates on $3*n^2$ words of memory

Reuse quotient ($q = \text{flops/word}$) in the algorithm is *potentially* as large as
 $2*n^3 / 3*n^2 = O(n)$



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

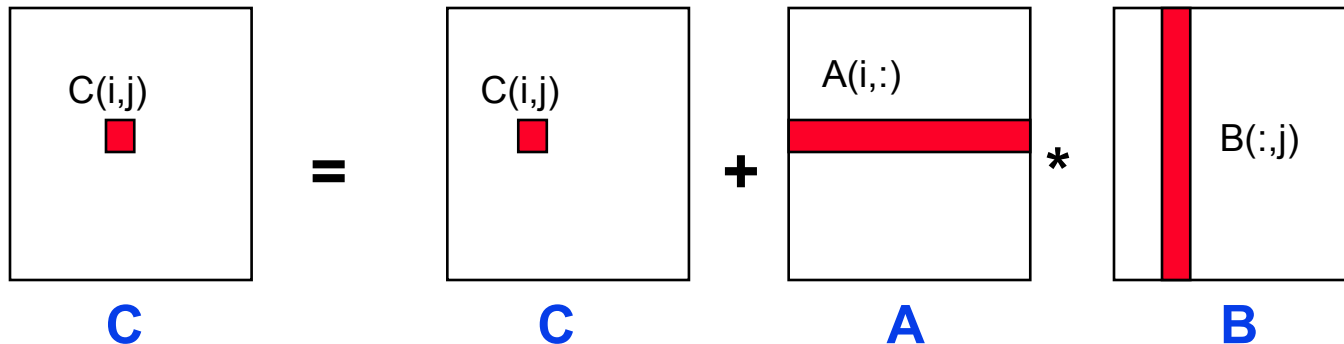
{read $C(i,j)$ into fast memory}

{read column j of B into fast memory}

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}

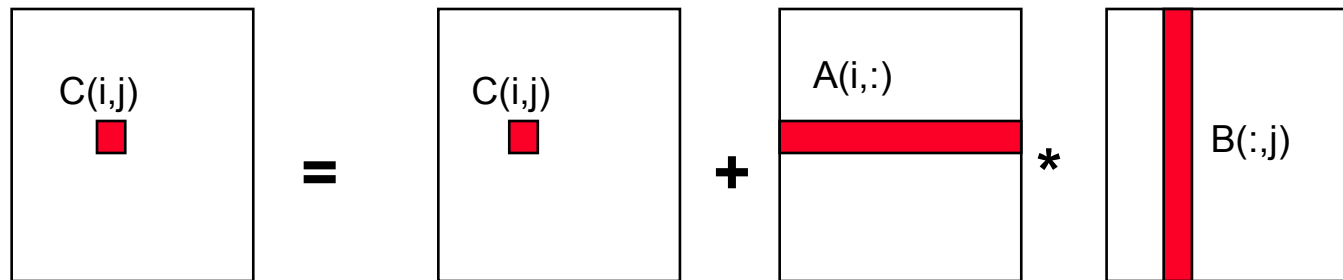


Naïve Matrix Multiply

- Number of slow memory references on unblocked matrix multiply

$$\begin{aligned} m &= n^3 && \text{to read each column of B } n \text{ times} \\ &+ n^2 && \text{to read each row of A once} \\ &+ 2n^2 && \text{to read and write each element of C once} \\ &= n^3 + 3n^2 \end{aligned}$$

- So the re-use quotient, $q = f / m = 2n^3 / (n^3 + 3n^2)$
 ~ 2 for large n
- This is no better than matrix-vector multiply
- And is far from the “best possible” which is $2/3 \cdot n$ for large n
- And this doesn't take into account cache lines



Blocked (Tiled) Matrix Multiply

Consider A,B,C to be N-by-N matrices of b-by-b subblocks where $b = n / N$ is called the **block size**

for $i = 1$ to N

for $j = 1$ to N

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

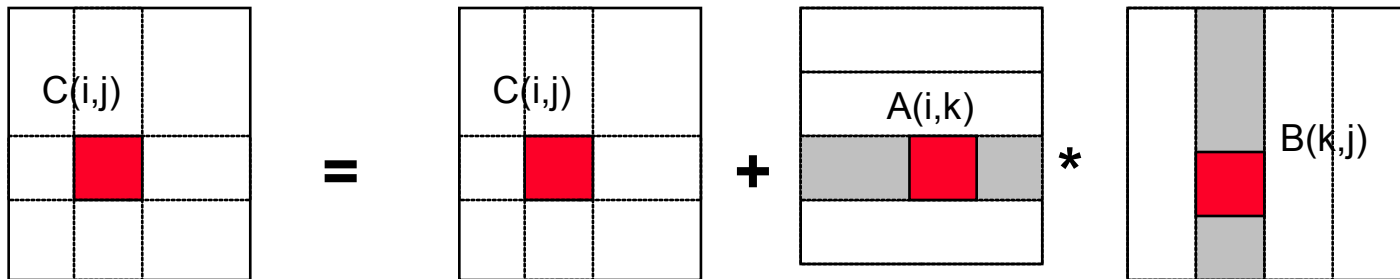
for $k = 1$ to N

{read block $A(i,k)$ into fast memory}

{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}



Blocked (Tiled) Matrix Multiply

- Recall:

m is amount memory traffic between slow and fast memory

matrix has $n \times n$ elements, also $N \times N$ blocks each of size $b \times b$ ($N = n/b$)

f is number of floating point operations, $2n^3$ for this problem

$q = f / m$ is our measure of algorithm efficiency in the memory system

$$\begin{aligned} m &= N * n^2 && \text{read each block of B } N^3 \text{ times } (N^3 * b^2 = N^3 * (n/N)^2 = N * n^2) \\ &+ N * n^2 && \text{read each block of A } N^3 \text{ times} \\ &+ 2n^2 && \text{read and write each block of C once} \\ &= (2N + 2) * n^2 \end{aligned}$$

- Re-use quotient is: $q = f / m = 2n^3 / ((2N + 2) * n^2)$
 $\approx n / N = b$ for large n
- We can improve performance by increasing the blocksize b
- Can be much faster than matrix-vector multiply ($q=2$)

Using Analysis to Understand Machines

The blocked algorithm has computational intensity $q \approx b$

- The larger the block size, the more efficient our algorithm will be
- Limit: All three blocks from A,B,C must fit in fast memory (cache), so we cannot make these blocks arbitrarily large
- Assume your fast memory has size M_{fast}

$$3b^2 \leq M_{\text{fast}}, \text{ so } q \approx b \leq \sqrt{M_{\text{fast}}/3}$$

- To build a machine to run matrix multiply at 1/2 peak arithmetic speed of the machine, we need a fast memory of size

$$M_{\text{fast}} \geq 3b^2 \approx 3q^2 = 3(t_m/t_f)^2$$

- This size is reasonable for L1 cache, but not for register sets
- Note: analysis assumes it is possible to schedule the instructions perfectly

	t_m/t_f	required KB
Ultra 2i	24.8	14.8
Ultra 3	14	4.7
Pentium 3	6.25	0.9
Pentium3M	10	2.4
Power3	8.75	1.8
Power4	15	5.4
Itanium1	36	31.1
Itanium2	5.5	0.7

Limits to Optimizing Matrix Multiply

- The blocked algorithm changes the order in which values are accumulated into each $C[i,j]$ by applying associativity
 - Get slightly different answers from naïve code, because of roundoff - OK
- The previous analysis showed that the blocked algorithm has computational intensity:

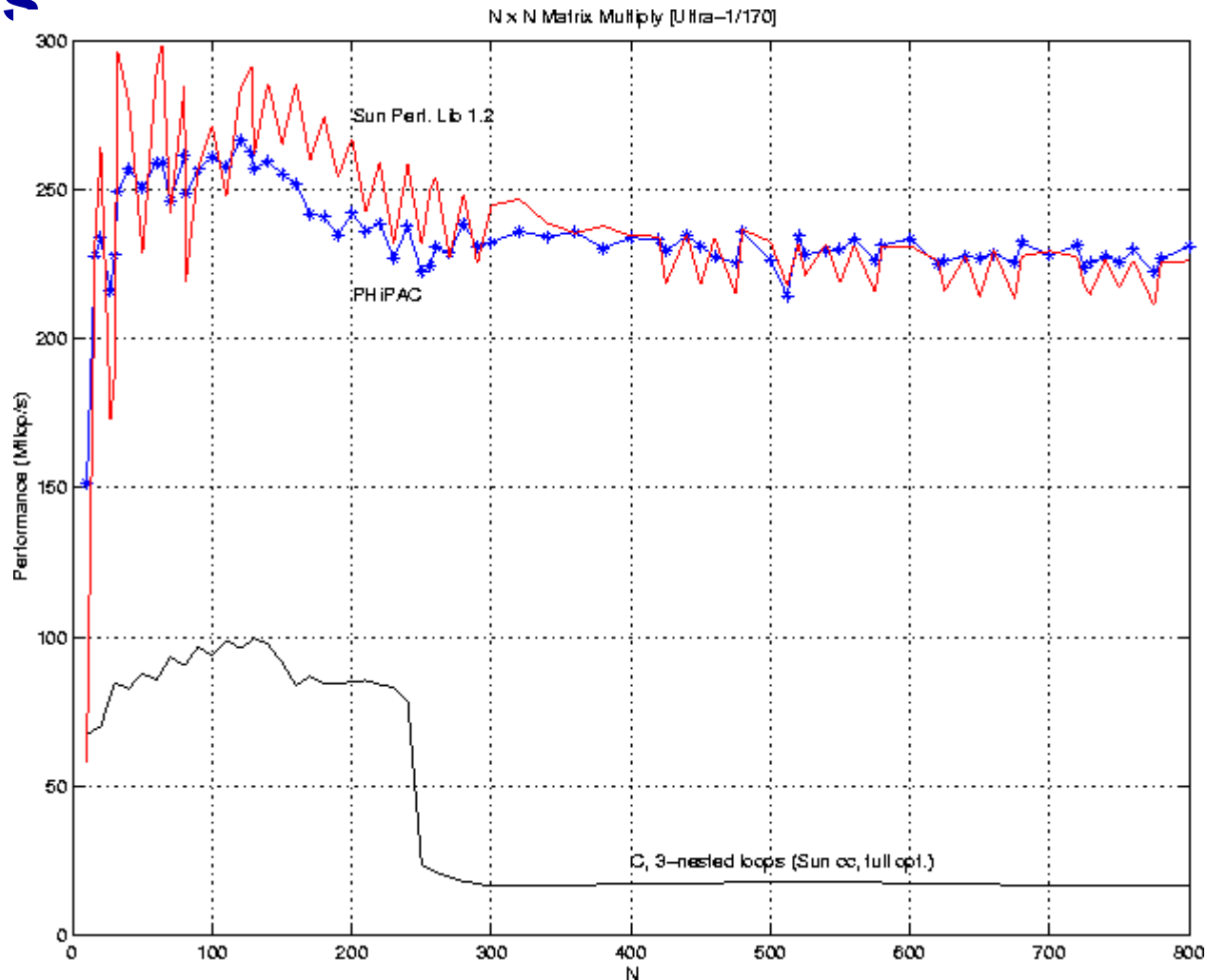
$$q \sim b \leq \sqrt{M_{\text{fast}}/3}$$

- Aside (for those who have taken CS170 or equivalent)
- There is a *lower bound* result that says we cannot do any better than this (using only associativity)
- Theorem (Hong & Kung, 1981): Any reorganization of this algorithm (that uses only associativity) is limited to $q = O(\sqrt{M_{\text{fast}}})$
- What if more levels of memory hierarchy?

Tiling for Multiple Levels

- Multiple levels: pages, caches, registers in machines
- Each level could have it's only 3-nested loops:
 for i, for j, for k {matmul blocks that fit in L2 cache}
 for ii, for jj, for kk {matmul blocks that fit in L1}
 for iii, for jjj, for kkk {matmul blocks that fit in registers}
- But in practice don't use so many levels and fully unroll code for register "tiles"
 - E.g., if these are 2x2 matrix multiplies, can write out all 4 statements without loops
 $c[0,0] = c[0,0] + a[0,0] * b[0,0] + a[0,1]*b[1,0]$
 $c[1,0] =$
 $c[0,1] =$
 $c[1,1] =$
 Many possible code variations; see Mark's notes on code generators in HW page

Matrix-multiply, optimized several ways



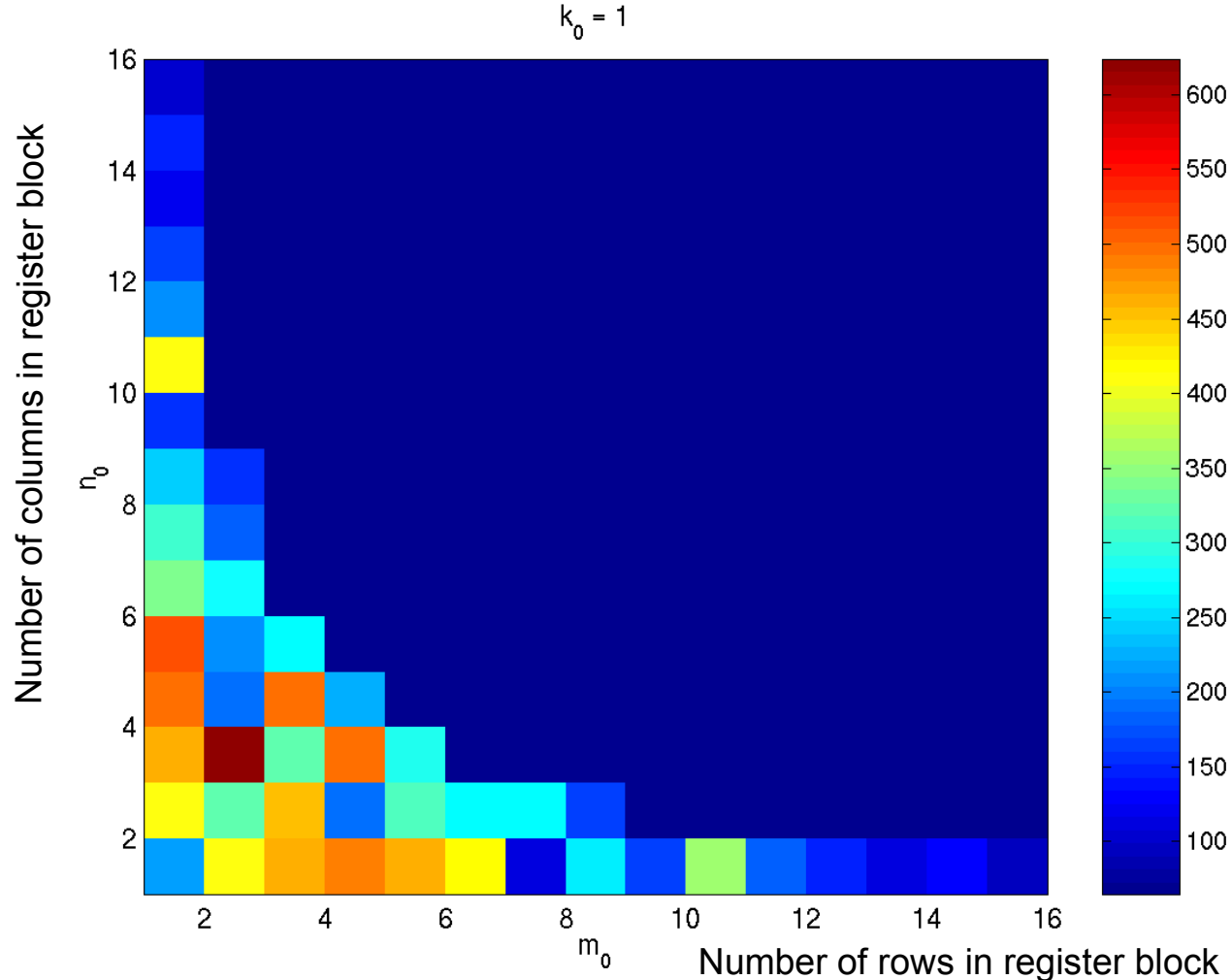
Speed of n-by-n matrix multiply on Sun Ultra-1/170, peak = 330 MFlops

Search Over Block Sizes

Strategies for choose block sizes:

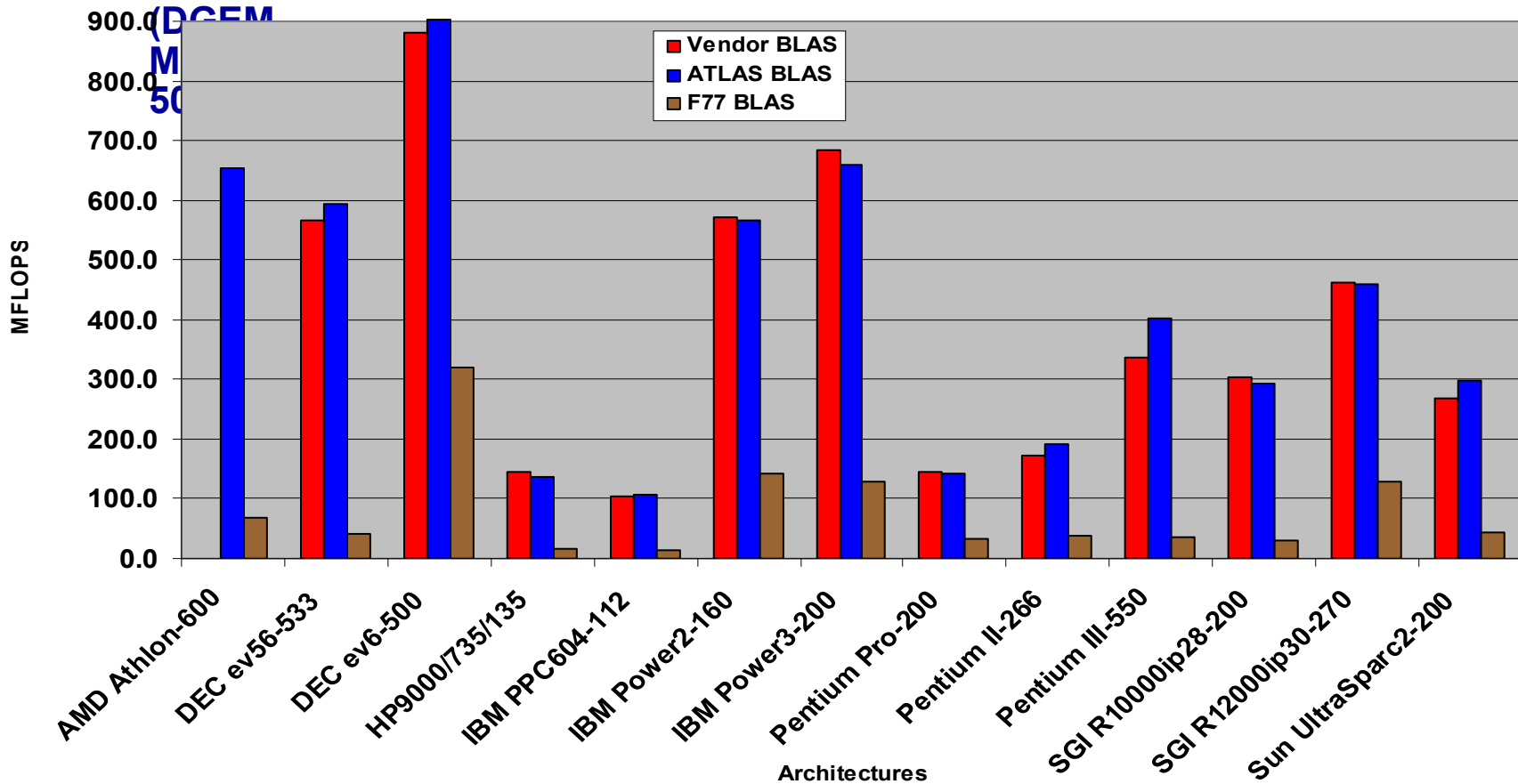
- Find out hardware parameters:
 - Read vendor manual for register #, cache sizes (not part of ISA), page sizes (OS config)
 - Measure yourself using memory benchmark from last time
 - But in practice these don't always work well; memory systems are complex as we saw
- Manually try many variations ☹
- Write “autotuner” to search over “design space” of possible implementations
 - Atlas – incorporated into Matlab
 - PhiPAC – original dense linear algebra tuning project from Berkeley; came from homework assignment like HW2

What the Search Space Looks Like



A 2-D slice of a 3-D register-tile search space. The dark blue region was pruned.
(Platform: Sun Ultra-Ili, 333 MHz, 667 Mflop/s peak, Sun cc v5.0 compiler)

Source: Jack Dongarra

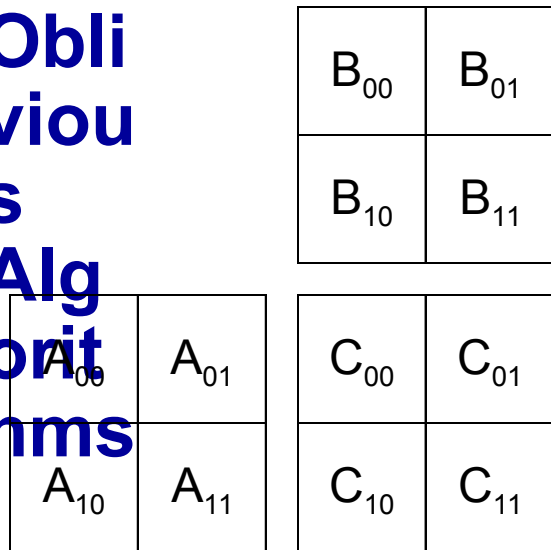


- ATLAS is faster than all other portable BLAS implementations and it is comparable with machine-specific libraries provided by the vendor.

Recursion: Cache Oblivious Algorithms

- The tiled algorithm requires finding a good block size
- Cache Oblivious Algorithms offer an alternative
 - Treat $n \times n$ matrix multiply set of smaller problems
 - Eventually, these will fit in cache
- The idea of cache-oblivious algorithms is use for other problems than matrix multiply. The general idea is:
 - Think of recursive formulation
 - If the subproblems use smaller data sets and some reuse within that data set, then a recursive order may improve performance

Cac he- Obli viou s Alg orit hms



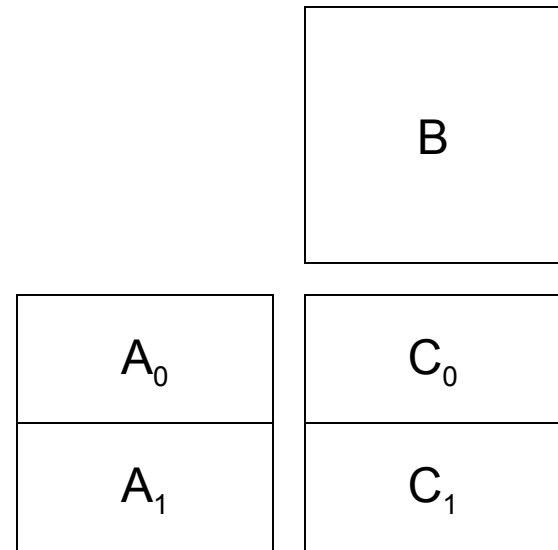
$$C_{00} = A_{00} * B_{00} + A_{01} * B_{10}$$

$$C_{01} = A_{01} * B_{11} + A_{00} * B_{01}$$

$$C_{11} = A_{11} * B_{01} + A_{10} * B_{01}$$

$$C_{10} = A_{10} * B_{00} + A_{11} * B_{10}$$

- Divide all dimensions (AD)
- 8-way recursive tree down to 1x1 blocks
 - Gray-code order promotes reuse
- Bilardi, et al.



$$C_0 = A_0 * B$$

$$C_1 = A_1 * B$$

- Divide largest dimension (LD)
- Two-way recursive tree down to 1x1 blocks
- Frigo, Leiserson, et al.

Cac

he- Obli viou

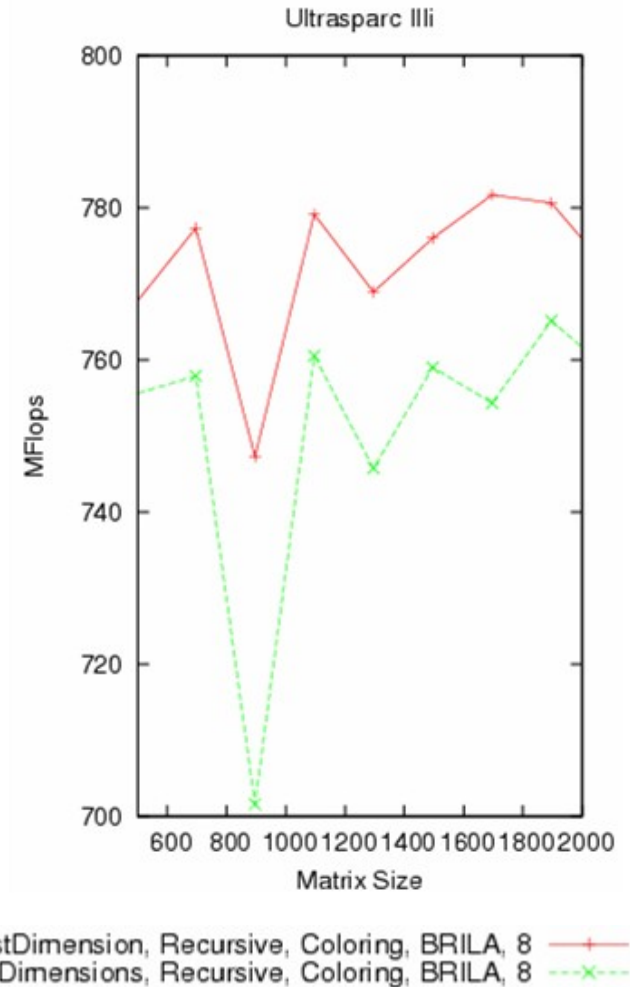
s:

Dis

cus sion

- Block sizes
 - Generated dynamically at each level in the recursive call tree
- Our experience
 - Performance is similar
 - Use AD for the rest of the talk

Mflop/s (up is good)



Cac

he-

- Usually Iterative

scjo

Nested loops

us

- Implementation of blocking

Algo

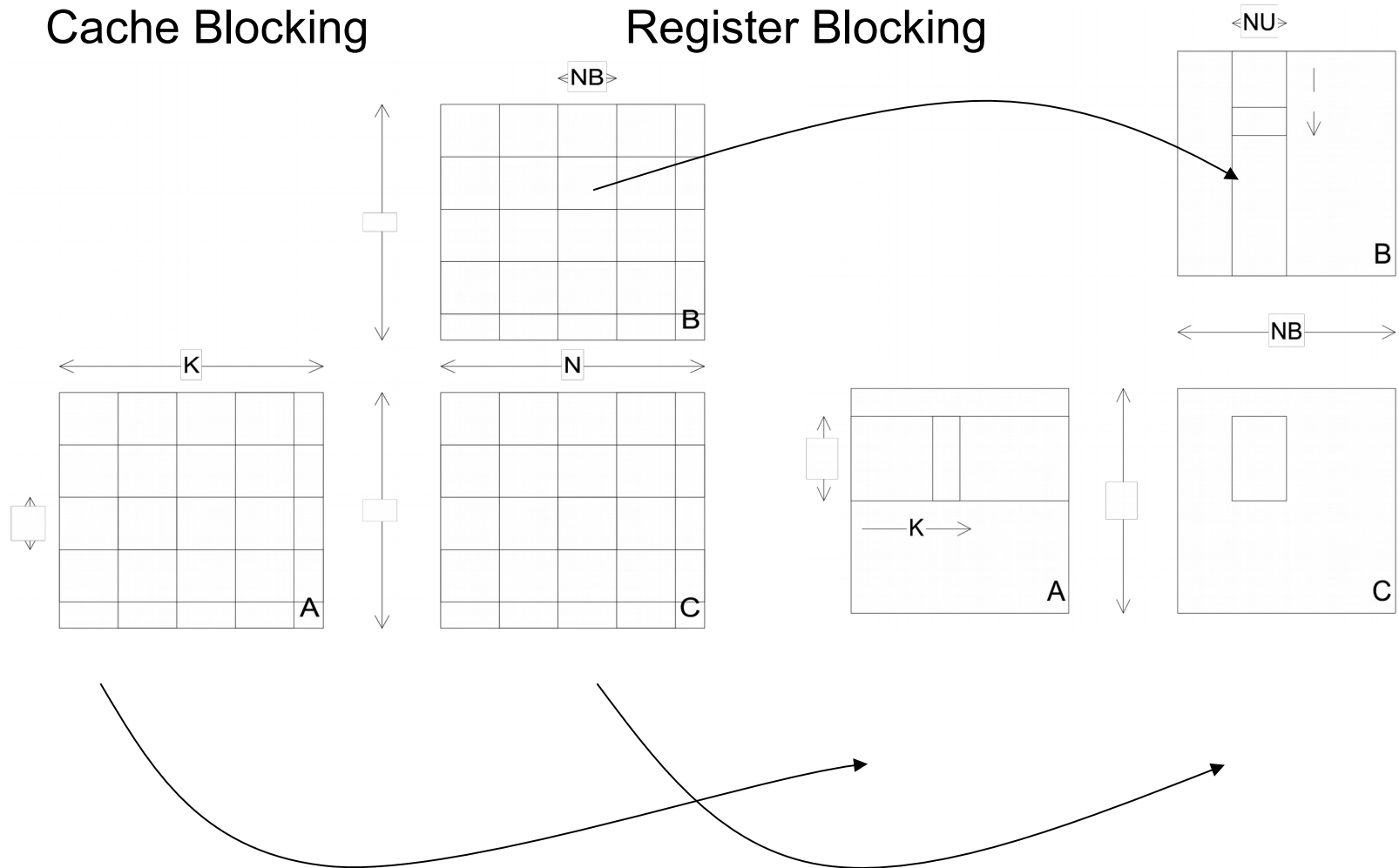
Cache blocking achieved by Loop Tiling

rith

ms

- Register blocking also requires Loop Unrolling

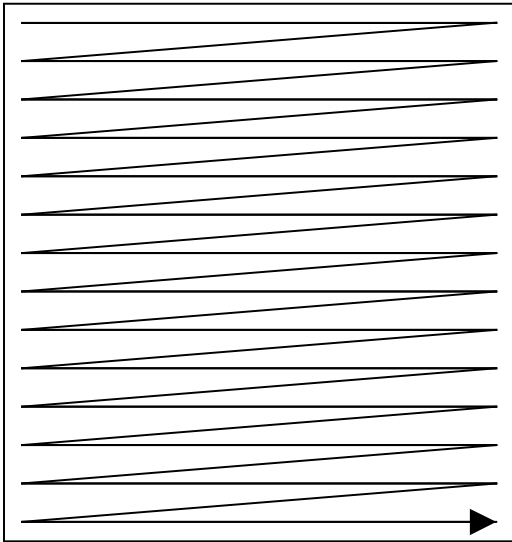
Structure of Tiled (Cache Conscious) Code



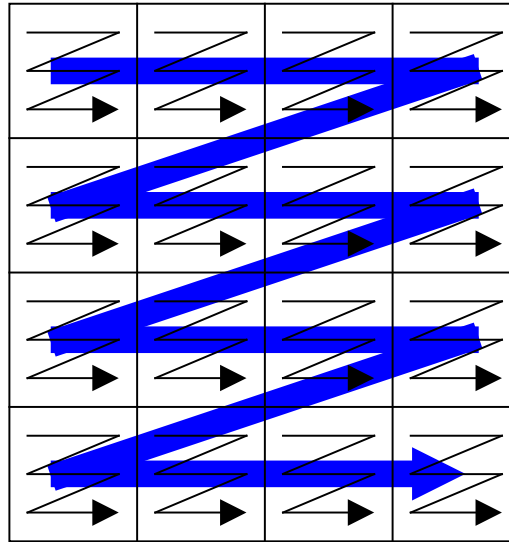
Data

- Fit control structure better
- Improve
 - Spatial locality
 - Streaming, prefetching

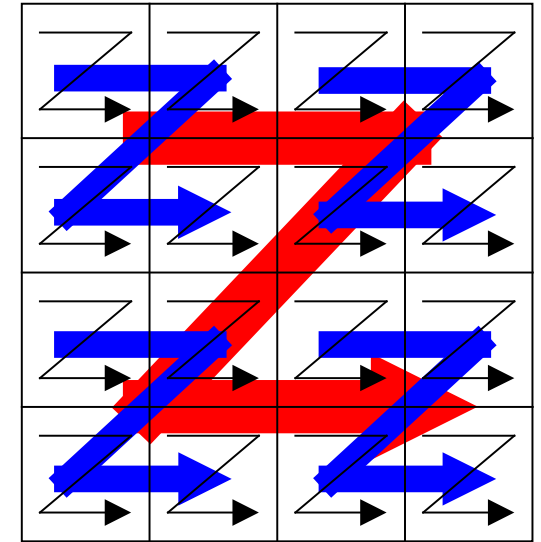
Row-major



Row-Block-Row

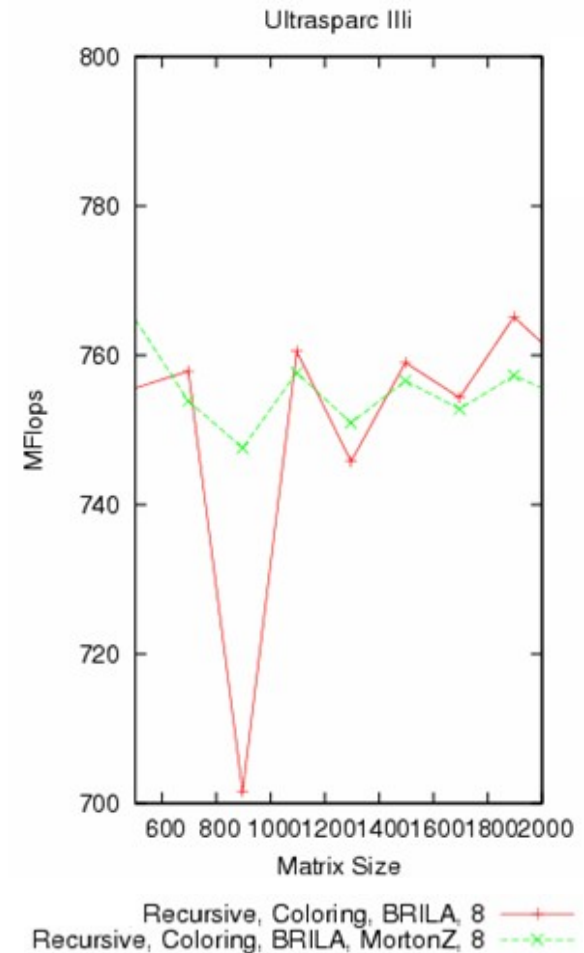


Morton-Z



Data

- Morton-Z
 - Matches recursive control structure better than RBR
 - Suggests better performance for CO
 - More complicated to implement
 - In our experience payoff is small or even negative
- Use RBR for the rest of the talk



Ultr

aSP

AR

C

Mr

- Peak performance
 - 2 GFlops
- Memory hierarchy
 - Registers: 32
 - L1 data cache: 64KB, 4-way
 - L2 data cache: 1MB, 4-way
- Compilers
 - FORTRAN: SUN F95 7.1
 - C: SUN C 5.5

Control Structures

- Iterative: triply nested loop
- Recursive: down to 1 x 1 x 1

Outer Control Structure

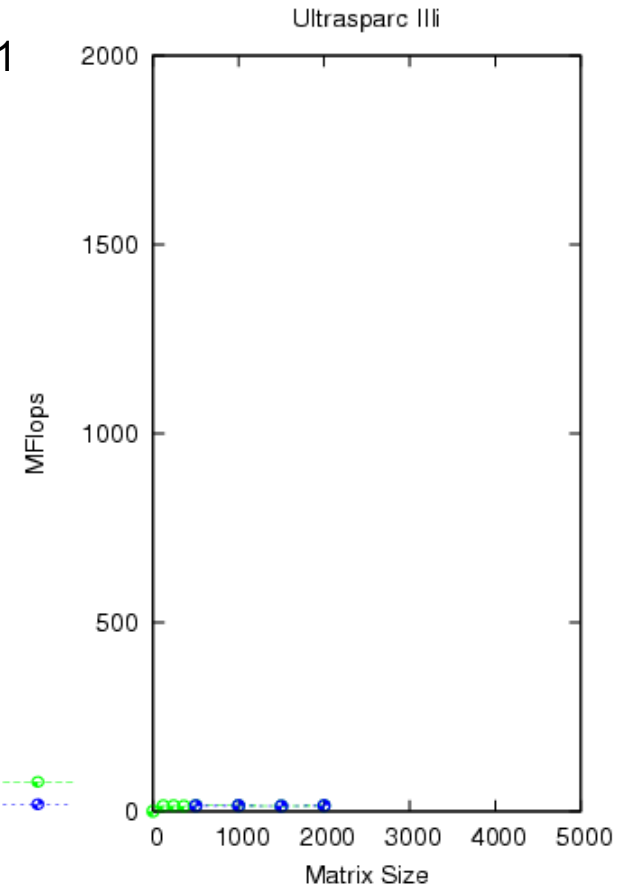
Iterative

Recursive

Inner Control Structure

Statement

Iterative, Statement, None, None, Compiler, 1
Recursive, Recursive, Micro, None, Compiler, 1



Control Structures

Outer Control Structure

Iterative Recursive

Inner Control Structure

Statement

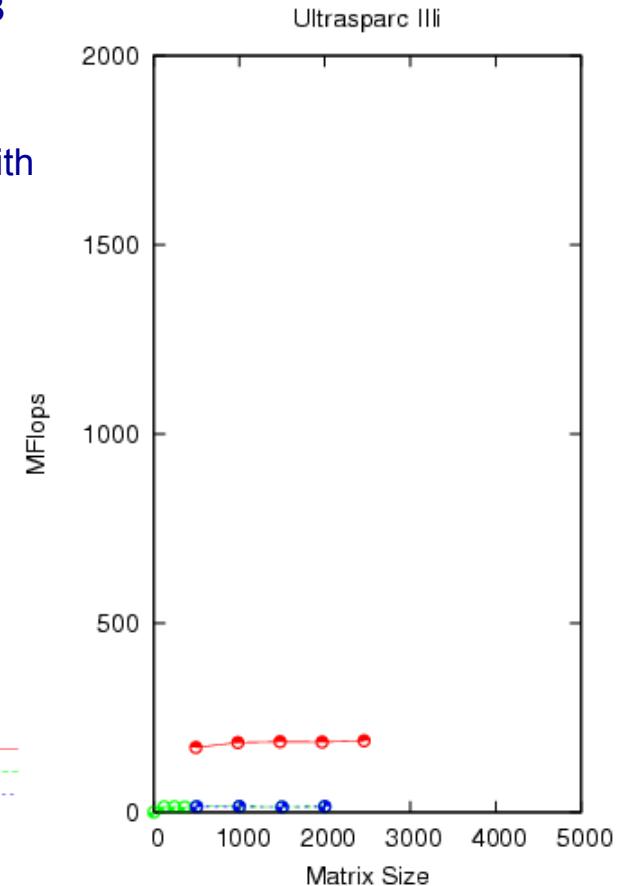
Recursive

Micro-Kernel

None /
Compiler

- Recursion down to NB
 - Unfold completely below NB to get a basic block
- Micro-Kernel:
 - The basic block compiled with native compiler
- Best performance for NB = 12
- Compiler unable to use registers
- Unfolding reduces control overhead
 - limited by I-cache

Recursive, Recursive, Micro, None, Compiler, 12
 Iterative, Statement, None, None, Compiler, 1
 Recursive, Recursive, Micro, None, Compiler, 1



Control Structures

Outer Control Structure

Iterative

Recursive

Inner Control Structure

Statement

Recursive

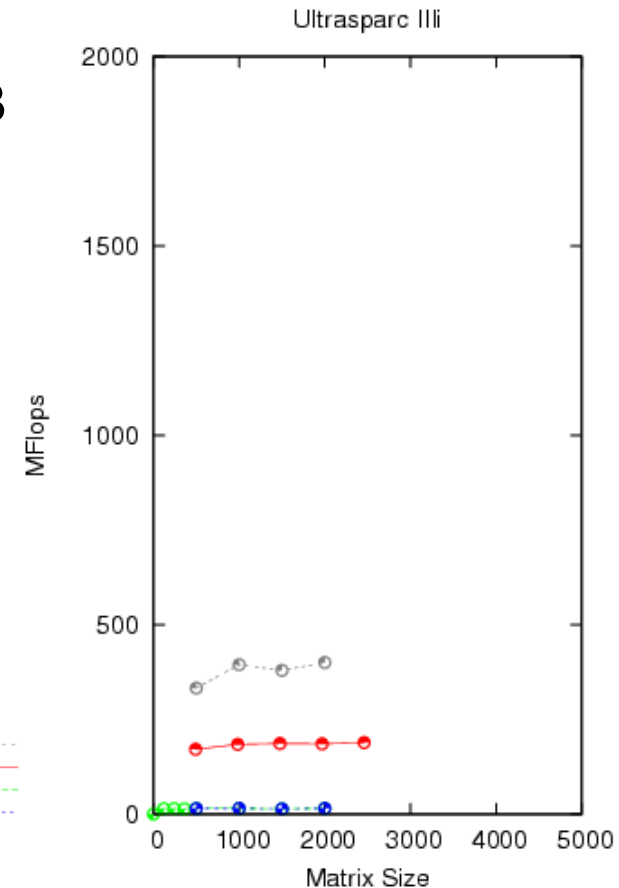
Micro-Kernel

None / Compiler

Scalarized / Compiler

- Recursion down to NB
 - Unfold completely below NB to get a basic block
- Micro-Kernel
 - Scalarize all array references in the basic block
 - Compile with native compiler

Recursive, Recursive, Micro, Scalarized, Compiler, 4
 Recursive, Recursive, Micro, None, Compiler, 12
 Iterative, Statement, None, None, Compiler, 1
 Recursive, Recursive, Micro, None, Compiler, 1



Control Structures

Outer Control Structures

Iterative

Recursive

Inner Control Structure

Statement

Recursive

Micro-Kernel

None /
Compiler

Scalarized /
Compiler

Belady /
BRILA

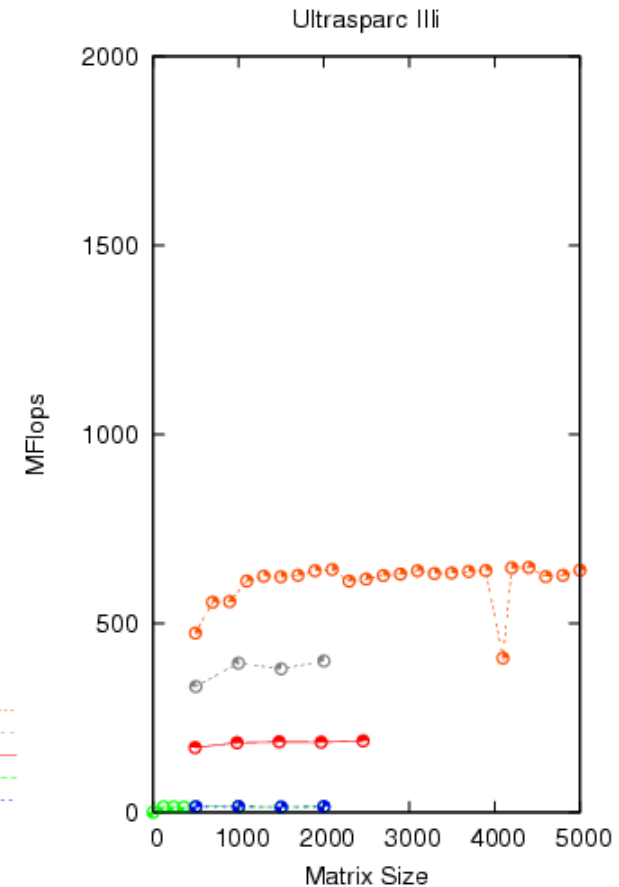
Recursion down to NB

- Unfold completely below NB to get a basic block

Micro-Kernel

- Perform Belady's register allocation on the basic block
- Schedule using BRILA compiler

Recursive, Recursive, Micro, Belady, BRILA, 8
Recursive, Recursive, Micro, Scalarized, Compiler, 4
Recursive, Recursive, Micro, None, Compiler, 12
Iterative, Statement, None, None, Compiler, 1
Recursive, Recursive, Micro, None, Compiler, 1



Control Structures

Outer Control Structures

Iterative

Recursive

Inner Control Structure

Statement

Recursive

Micro-Kernel

None /
Compiler

Scalarized /
Compiler

Belady /
BRILA

Coloring /
BRILA

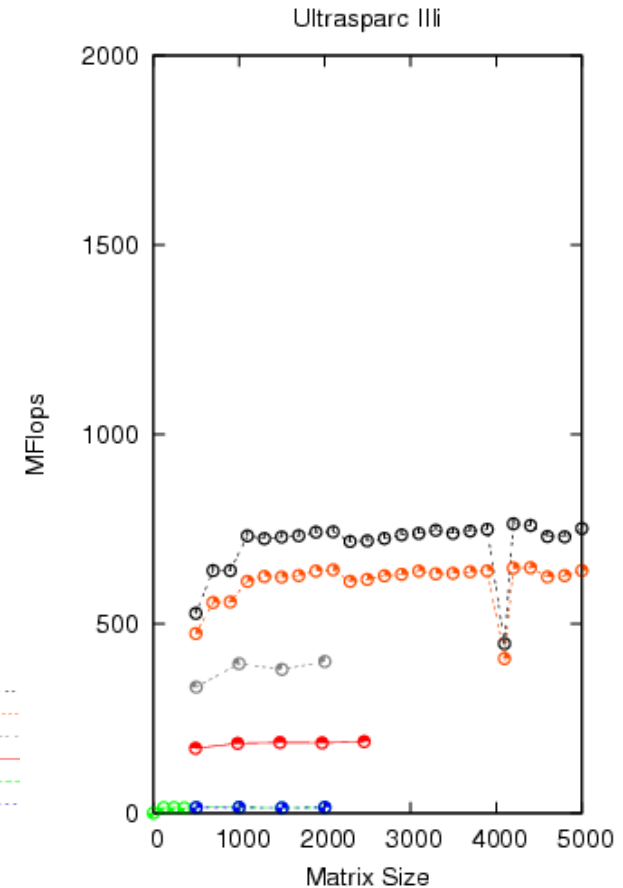
Recursion down to NB

- Unfold completely below NB to get a basic block

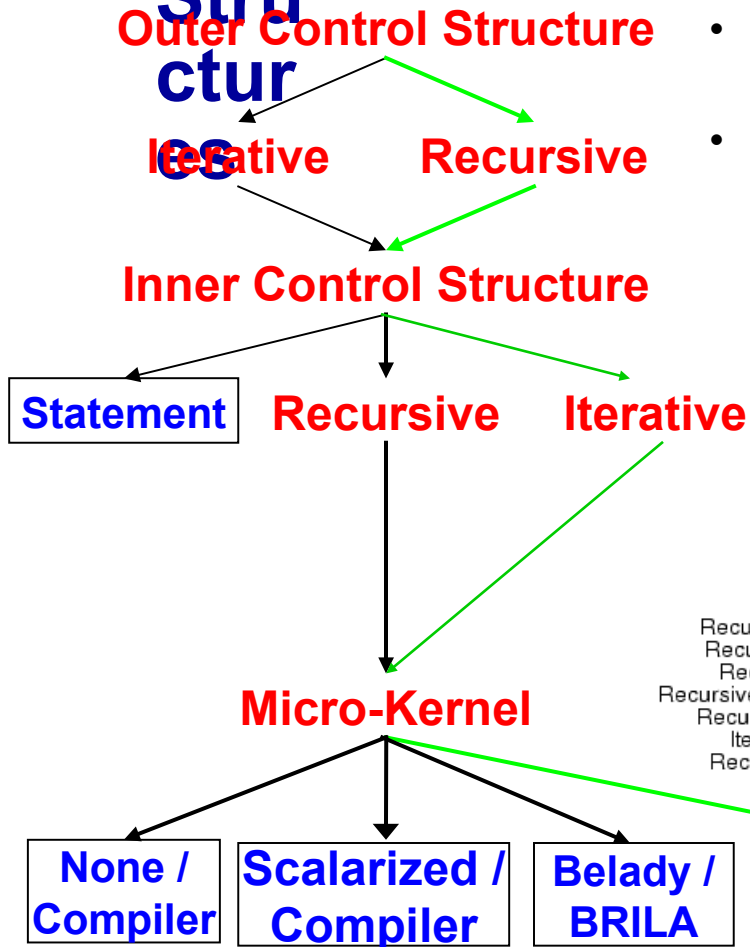
Micro-Kernel

- Construct a preliminary schedule
- Perform Graph Coloring register allocation
- Schedule using BRILA compiler

Recursive, Recursive, Micro, Coloring, BRILA, 8
Recursive, Recursive, Micro, Belady, BRILA, 8
Recursive, Recursive, Micro, Scalarized, Compiler, 4
Recursive, Recursive, Micro, None, Compiler, 12
Iterative, Statement, None, None, Compiler, 1
Recursive, Recursive, Micro, None, Compiler, 1

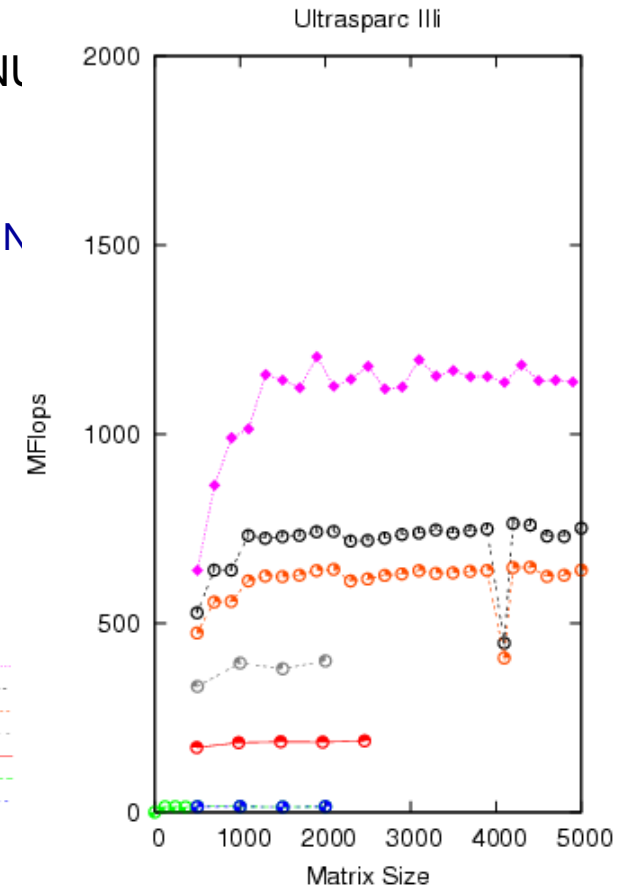


Control Structures

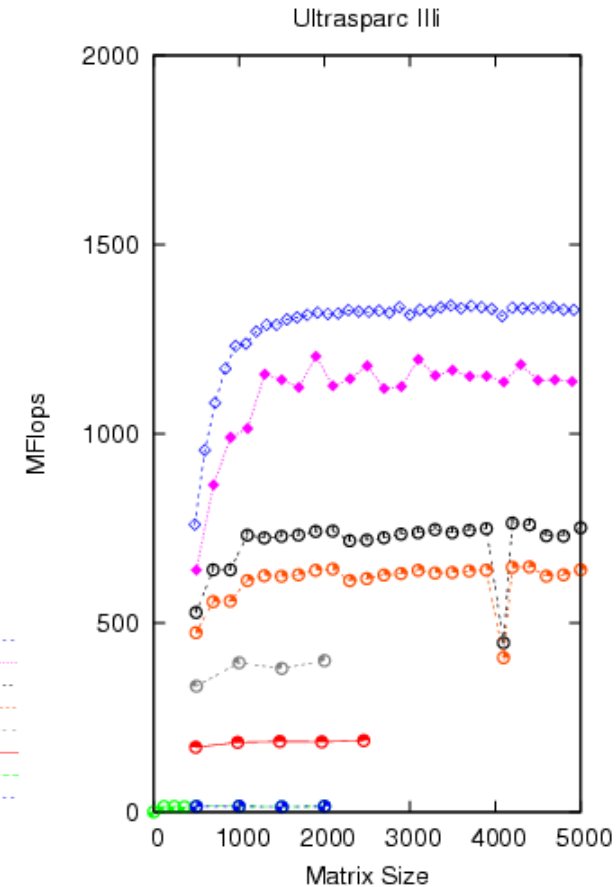
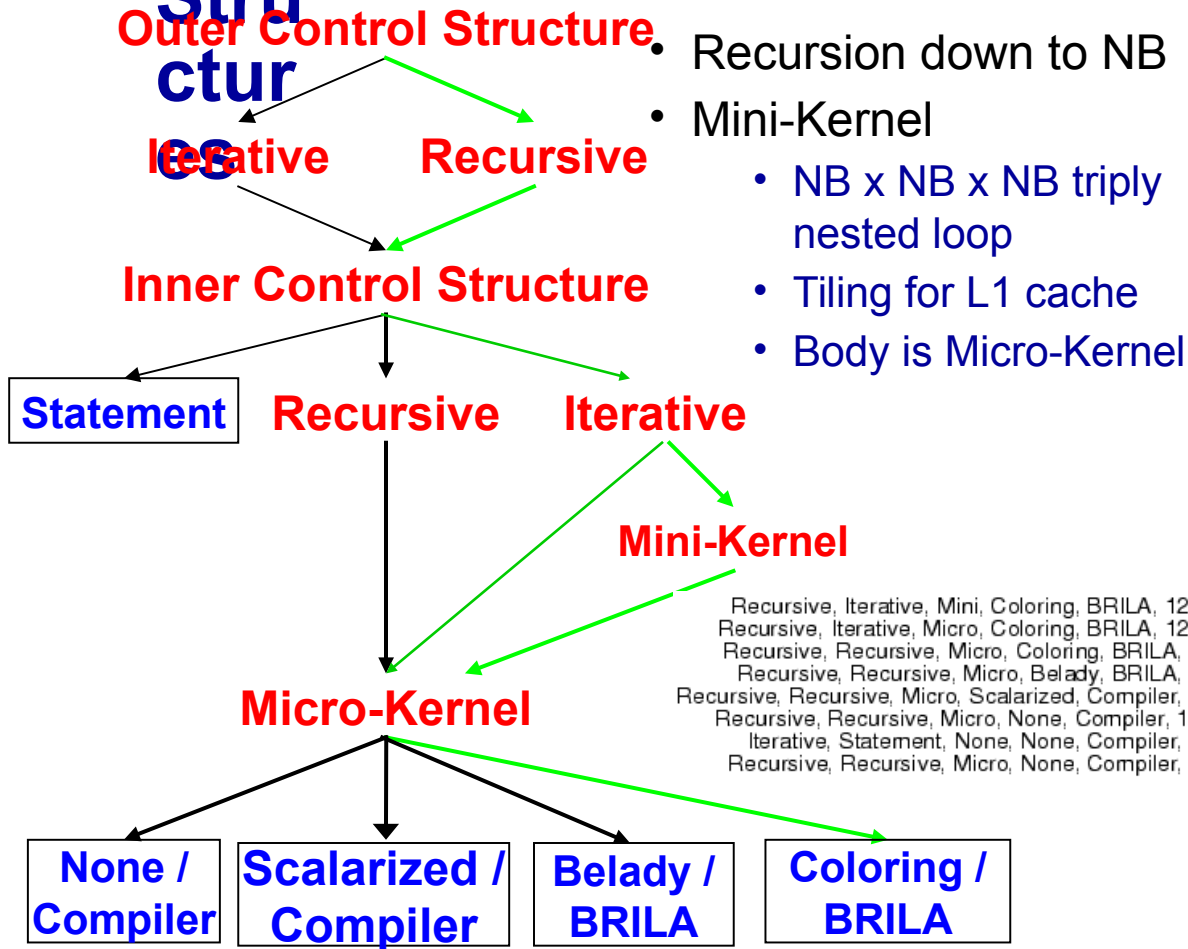


- Recursion down to $MU \times N \times KU$
- Micro-Kernel
 - Completely unroll $MU \times N \times KU$ triply nested loop
 - Construct a preliminary schedule
 - Perform Graph Coloring register allocation
 - Schedule using BRILA compiler

Recursive, Iterative, Micro, Coloring, BRILA, 120
 Recursive, Recursive, Micro, Coloring, BRILA, 8
 Recursive, Recursive, Micro, Belady, BRILA, 8
 Recursive, Recursive, Micro, Scalarized, Compiler, 4
 Recursive, Recursive, Micro, None, Compiler, 12
 Iterative, Statement, None, None, Compiler, 1
 Recursive, Recursive, Micro, None, Compiler, 1



Control Structures



Control Structures

Outer Control Structures

Iterative

Recursive

Inner Control Structures

Statement

Recursive

Iterative

Mini-Kernel

Micro-Kernel

None /
Compiler

Scalarized /
Compiler

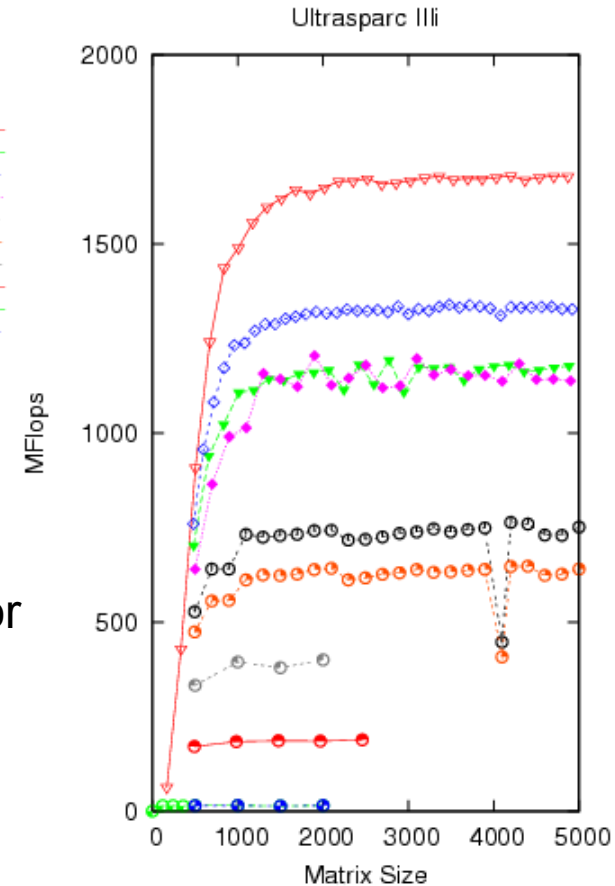
Belady /
BRILA

Coloring /
BRILA

Recursive, Iterative, Mini, ATLAS, Unleashed, 168
Recursive, Iterative, Mini, ATLAS, CGwS, 44
Recursive, Iterative, Mini, Coloring, BRILA, 120
Recursive, Iterative, Micro, Coloring, BRILA, 120
Recursive, Recursive, Micro, Coloring, BRILA, 8
Recursive, Recursive, Micro, Belady, BRILA, 8
Recursive, Recursive, Micro, Scalarized, Compiler, 4
Recursive, Recursive, Micro, None, Compiler, 12
Iterative, Statement, None, None, Compiler, 1
Recursive, Recursive, Micro, None, Compiler, 1

Specialized
code generator
with search

ATLAS CGw/S
ATLAS Unleashed



Control Structures

Outer Control Structure

Iterative

Recursive

Inner Control Structure

Statement

Recursive

Iterative

Mini-Kernel

Micro-Kernel

None /
Compiler

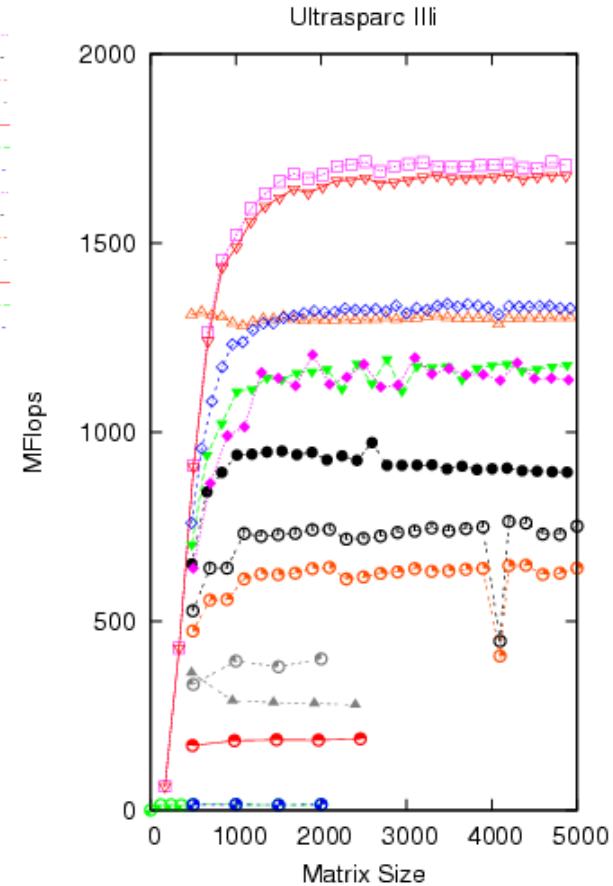
Scalarized /
Compiler

Belady /
BRILA

Coloring /
BRILA

ATLAS CGw/S
ATLAS Unleashed

Iterative, Iterative, Mini, ATLAS, Unleashed, 168
Iterative, Iterative, Mini, ATLAS, CGwS, 44
Iterative, Iterative, Mini, Coloring, BRILA, 120
Iterative, Iterative, Micro, Coloring, BRILA, 120
Recursive, Iterative, Mini, ATLAS, Unleashed, 168
Recursive, Iterative, Mini, ATLAS, CGwS, 44
Recursive, Iterative, Mini, Coloring, BRILA, 120
Recursive, Iterative, Micro, Coloring, BRILA, 120
Recursive, Recursive, Micro, Coloring, BRILA, 8
Recursive, Recursive, Micro, Belady, BRILA, 8
Recursive, Recursive, Micro, Scalarized, Compiler, 4
Recursive, Recursive, Micro, None, Compiler, 12
Iterative, Statement, None, None, Compiler, 1
Recursive, Recursive, Micro, None, Compiler, 1



Experience

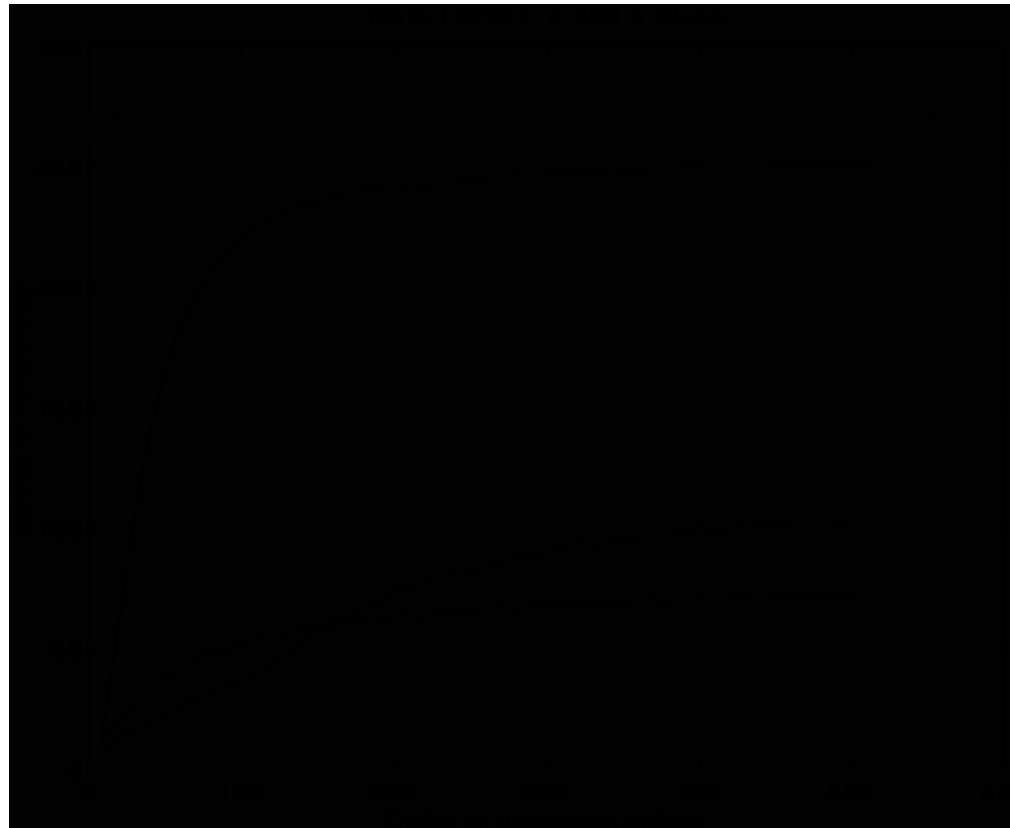
- In practice, need to cut off recursion
- Implementing a high-performance Cache-Oblivious code is not easy
 - Careful attention to micro-kernel and mini-kernel is needed
- Using fully recursive approach with highly optimized recursive micro-kernel, Pingali et al report that they never got more than 2/3 of peak.
- Issues with Cache Oblivious (recursive) approach
 - Recursive Micro-Kernels yield less performance than iterative ones using same scheduling techniques
 - Pre-fetching is needed to compete with best code: not well-understood in the context of CO codes

Basic Linear Algebra Subroutines (BLAS)

- Industry standard interface (evolving)
 - www.netlib.org/blas, www.netlib.org/blas/blast--forum
- Vendors, others supply optimized implementations
- History
 - **BLAS1 (1970s):**
 - vector operations: dot product, saxpy ($y = \alpha * x + y$), etc
 - $m = 2 * n$, $f = 2 * n$, $q \sim 1$ or less
 - **BLAS2 (mid 1980s)**
 - matrix-vector operations: matrix vector multiply, etc
 - $m = n^2$, $f = 2 * n^2$, $q \sim 2$, less overhead
 - somewhat faster than BLAS1
 - **BLAS3 (late 1980s)**
 - matrix-matrix operations: matrix matrix multiply, etc
 - $m \leq 3n^2$, $f = O(n^3)$, so $q = f/m$ can possibly be as large as n , so BLAS3 is potentially much faster than BLAS2
- Good algorithms used BLAS3 when possible (LAPACK & ScaLAPACK)
 - See www.netlib.org/{lapack,scalapack}
 - More later in course

BLAS speeds on an IBM RS6000/590

Peak speed = 266 Mflops



Peak
BLAS 3

BLAS 2

BLAS 1

BLAS 3 (n-by-n matrix matrix multiply) vs
BLAS 2 (n-by-n matrix vector multiply) vs
BLAS 1 (saxpy of n vectors)

Strassen's Matrix Multiply

- The traditional algorithm (with or without tiling) has $O(n^3)$ flops
- Strassen discovered an algorithm with asymptotically lower flops
 - $O(n^{2.81})$
- Consider a 2×2 matrix multiply, normally takes 8 multiplies, 4 adds
 - Strassen does it with 7 multiplies and 18 adds

$$\text{Let } M = \begin{pmatrix} m_{11} & m_{12} \\ m_{21} & m_{22} \end{pmatrix} = \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix} \begin{pmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \end{pmatrix}$$

$$\text{Let } p1 = (a_{12} - a_{22}) * (b_{21} + b_{22})$$

$$p5 = a_{11} * (b_{12} - b_{22})$$

$$p2 = (a_{11} + a_{22}) * (b_{11} + b_{22})$$

$$p6 = a_{22} * (b_{21} - b_{11})$$

$$p3 = (a_{11} - a_{21}) * (b_{11} + b_{12})$$

$$p7 = (a_{21} + a_{22}) * b_{11}$$

$$p4 = (a_{11} + a_{12}) * b_{22}$$

$$\text{Then } m_{11} = p1 + p2 - p4 + p6$$

$$m_{12} = p4 + p5$$

$$m_{21} = p6 + p7$$

$$m_{22} = p2 - p3 + p5 - p7$$

Extends to $n \times n$ by divide&conquer

Strassen (continued)

$$\begin{aligned} T(n) &= \text{Cost of multiplying } nxn \text{ matrices} \\ &= 7 * T(n/2) + 18 * (n/2)^2 \\ &= O(n \log_2 7) \\ &= O(n^{2.81}) \end{aligned}$$

- Asymptotically faster
 - Several times faster for large n in practice
 - Cross-over depends on machine
 - Available in several libraries
 - “Tuning Strassen's Matrix Multiplication for Memory Efficiency”, M. S. Thottethodi, S. Chatterjee, and A. Lebeck, in Proceedings of Supercomputing '98
- Caveats
 - Needs more memory than standard algorithm
 - Can be less accurate because of roundoff error

Other Fast Matrix Multiplication Algorithms

- Current world's record is $O(n^{2.376...})$ (Coppersmith & Winograd)
- Why does Hong/Kung theorem not apply?
- Possibility of $O(n^{2+\epsilon})$ algorithm! (Cohn, Umans, Kleinberg, 2003)
- Fast methods (besides Strassen) may need unrealistically large n

Optimizing in Practice

- Tiling for registers
 - loop unrolling, use of named “register” variables
- Tiling for multiple levels of cache and TLB
- Exploiting fine-grained parallelism in processor
 - superscalar; pipelining
- Complicated compiler interactions
- Hard to do by hand (but you’ll try)
- Automatic optimization an active research area
 - BeBOP: bebop.cs.berkeley.edu/
 - PHiPAC: www.icsi.berkeley.edu/~bilmes/hipac
in particular tr-98-035.ps.gz
 - ATLAS: www.netlib.org/atlas

Removing False Dependencies

- Using local variables, reorder operations to remove false dependencies

```
a[i] = b[i] + c;  
a[i+1] = b[i+1] * d;
```

false read-after-write hazard
between a[i] and b[i+1]



```
float f1 = b[i];  
float f2 = b[i+1];  
  
a[i] = f1 + c;  
a[i+1] = f2 * d;
```

With some compilers, you can declare a and b unaliased.

- Done via “restrict pointers,” compiler flag, or pragma)

Exploit Multiple Registers

- Reduce demands on memory bandwidth by pre-loading into local variables

```
while( ... ) {  
    *res++ = filter[0]*signal[0]  
            + filter[1]*signal[1]  
            + filter[2]*signal[2];  
    signal++;  
}
```



```
float f0 = filter[0];  
float f1 = filter[1];  
float f2 = filter[2];  
while( ... ) {  
    *res++ = f0*signal[0]  
            + f1*signal[1]  
            + f2*signal[2];  
    signal++;  
}
```

also: register float f0 = ...;

Example is a convolution

Minimize Pointer Updates

- Replace pointer updates for strided memory addressing with constant array offsets

```
f0 = *r8; r8 += 4;  
f1 = *r8; r8 += 4;  
f2 = *r8; r8 += 4;
```



```
f0 = r8[0];  
f1 = r8[4];  
f2 = r8[8];  
r8 += 12;
```

Pointer vs. array expression costs may differ.

- Some compilers do a better job at analyzing one than the other

Loop Unrolling

- Expose instruction-level parallelism

```
float f0 = filter[0], f1 = filter[1], f2 = filter[2];
float s0 = signal[0], s1 = signal[1], s2 = signal[2];
*res++ = f0*s0 + f1*s1 + f2*s2;
do {
    signal += 3;
    s0 = signal[0];
    res[0] = f0*s1 + f1*s2 + f2*s0;

    s1 = signal[1];
    res[1] = f0*s2 + f1*s0 + f2*s1;

    s2 = signal[2];
    res[2] = f0*s0 + f1*s1 + f2*s2;

    res += 3;
} while( ... );
```

Expose Independent Operations


- Hide instruction latency
 - Use local variables to expose independent operations that can execute in parallel or in a pipelined fashion
 - Balance the instruction mix (what functional units are available?)

```
f1 = f5 * f9;  
f2 = f6 + f10;  
f3 = f7 * f11;  
f4 = f8 + f12;
```

Copy optimization

- Copy input operands or blocks
 - Reduce cache conflicts
 - Constant array offsets for fixed size blocks
 - Expose page-level locality

Original matrix
(numbers are addresses)



0	4	8	12
1	5	9	13
2	6	10	14
3	7	11	15

Reorganized into
2x2 blocks

0	2	8	10
1	3	9	11
4	6	12	13
5	7	14	15

Locality in Other Algorithms

- The performance of any algorithm is limited by q
- In matrix multiply, we increase q by changing computation order
 - increased temporal locality
- For other algorithms and data structures, even hand-transformations are still an open problem
 - sparse matrices (reordering, blocking)
 - Weekly research meetings
 - Bebop.cs.berkeley.edu
 - About to release OSKI – tuning for sparse-matrix-vector multiply
 - trees (B-Trees are for the disk level of the hierarchy)
 - linked lists (some work done here)

Summary

- Performance programming on uniprocessors requires
 - understanding of memory system
 - understanding of fine-grained parallelism in processor
- Simple performance models can aid in understanding
 - Two ratios are key to efficiency (relative to peak)
 - 1.computational intensity of the algorithm:
 - $q = f/m = \# \text{ floating point operations} / \# \text{ slow memory references}$
 - 2.machine balance in the memory system:
 - $t_m/t_f = \text{time for slow memory reference} / \text{time for floating point operation}$
- Want $q > t_m/t_f$ to get half machine peak
- Blocking (tiling) is a basic approach to increase q
 - Techniques apply generally, but the details (e.g., block size) are architecture dependent
 - Similar techniques possible on other data structures / algorithms