# CSC367 Parallel computing

# Lecture 3: Single Processor Machines-Performance Model

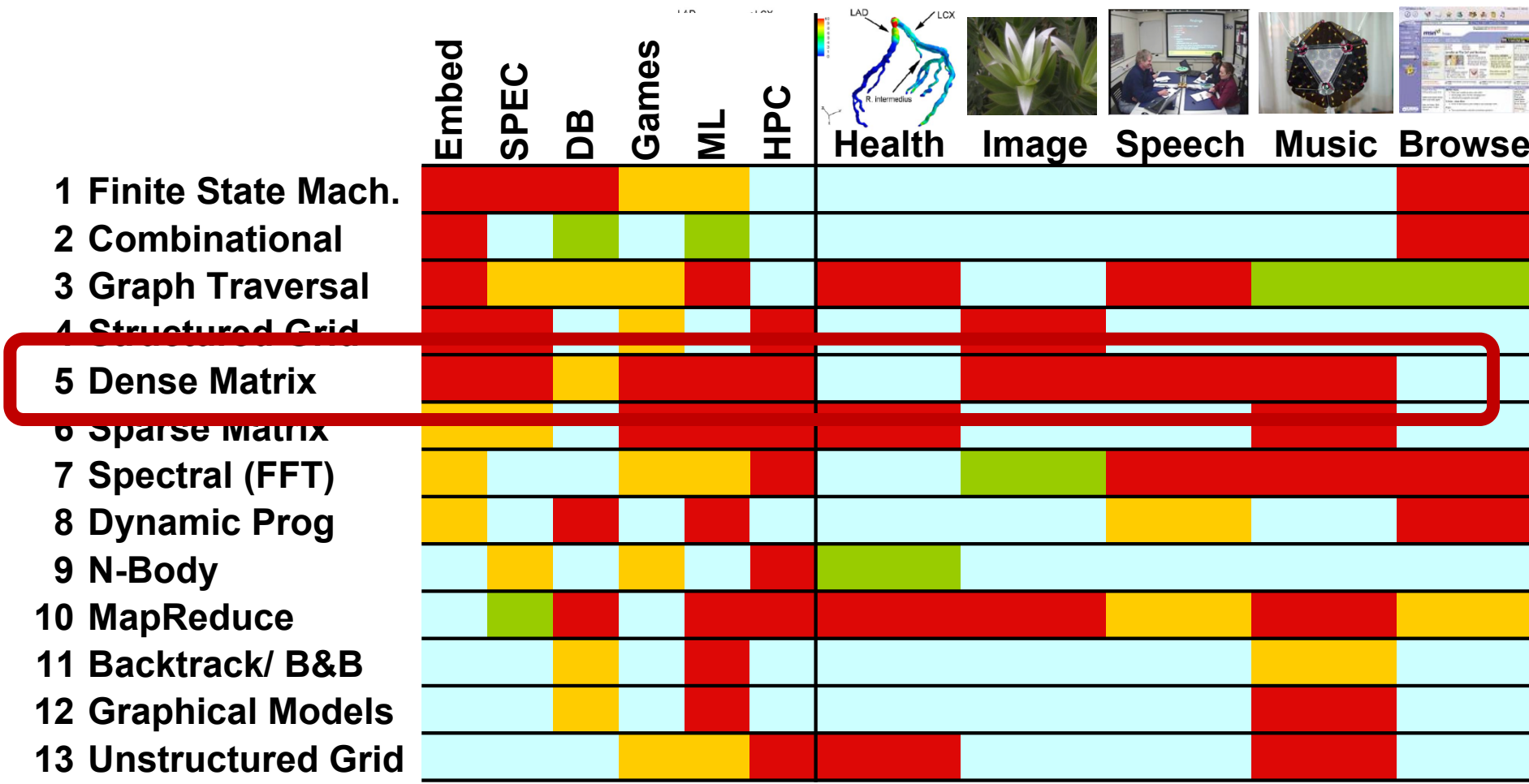Maryam Mehri Dehnavi
mmehride@cs.toronto.edu

1

# Outline

A performance model for Matrix Multiplication

- Use of performance models to understand performance
- Attainable lower bounds on communication
- Simple cache model
- Warm-up: Matrix-vector multiplication
- Naïve vs optimized Matrix-Matrix Multiply
  - Minimizing data movement
  - Beating $O(n^3)$ operations
- BLAS routines

# What do commercial and CSE applications have in common?

## Motif/Dwarf: Common Computational Methods
### (Red Hot → Blue Cool)

| | Embed | SPEC | DB | Games | ML | HPC | Health | Image | Speech | Music | Browser |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 Finite State Mach. | | | | | | | | | | | |
| 2 Combinational | | | | | | | | | | | |
| 3 Graph Traversal | | | | | | | | | | | |
| 4 Structured Grid | | | | | | | | | | | |
| 5 Dense Matrix | | | | | | | | | | | |
| 6 Sparse Matrix | | | | | | | | | | | |
| 7 Spectral (FFT) | | | | | | | | | | | |
| 8 Dynamic Prog | | | | | | | | | | | |
| 9 N-Body | | | | | | | | | | | |
| 10 MapReduce | | | | | | | | | | | |
| 11 Backtrack/ B&B | | | | | | | | | | | |
| 12 Graphical Models | | | | | | | | | | | |
| 13 Unstructured Grid | | | | | | | | | | | |

# Note on Matrix Storage

- A matrix is a 2-D array of elements, but memory addresses are "1-D"

- Conventions for matrix layout
    - by column, or "column major" (Fortran default); A(i,j) at A+i+j*n
    - by row, or "row major" (C default) A(i,j) at A+i*n+j
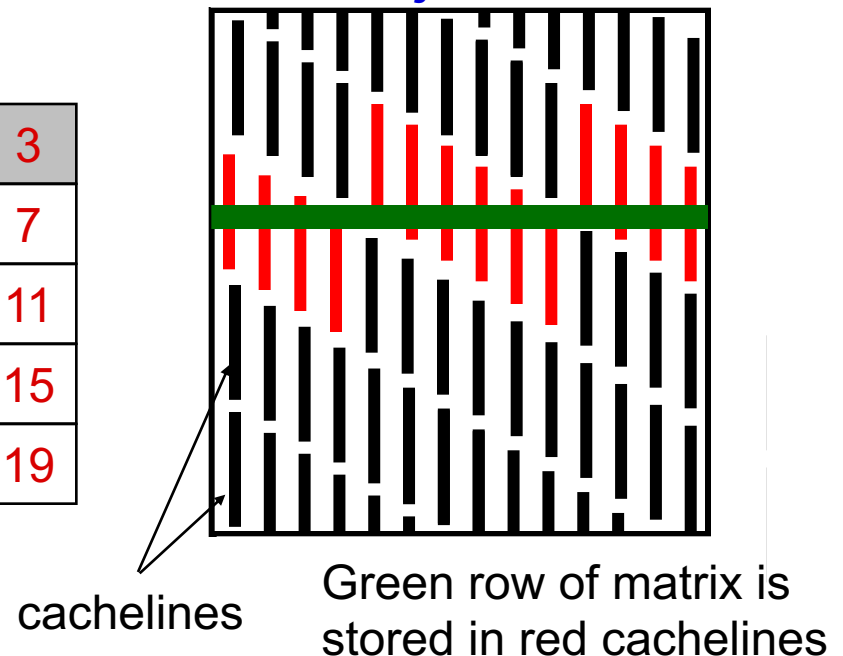    - Recursive

**Column major**

| | | | |
|---|---|---|---|
| 0 | 5 | 10 | 15 |
| 1 | 6 | 11 | 16 |
| 2 | 7 | 12 | 17 |
| 3 | 8 | 13 | 18 |
| 4 | 9 | 14 | 19 |

**Row major**

| | | | |
|---|---|---|---|
| 0 | 1 | 2 | 3 |
| 4 | 5 | 6 | 7 |
| 8 | 9 | 10 | 11 |
| 12 | 13 | 14 | 15 |
| 16 | 17 | 18 | 19 |

**Column major matrix in memory**



cachelines

Green row of matrix is stored in red cachelines

- Column major (for now)

# Using a Simple Model of Memory to Optimize

- Assume just 2 levels in the hierarchy, fast and slow
- All data initially in slow memory
  - $m$ = number of memory elements (words) moved between fast and slow memory
  - $t_m$ = time per slow memory operation
  - $f$ = number of arithmetic operations
  - $t_f$ = time per arithmetic operation $<< t_m$
  - $q = f / m$  average number of flops per slow memory access

*Computational Intensity:* **Key to algorithm efficiency**

- Minimum possible time = $f * t_f$ when all data in fast memory
- Actual time
  - $f * t_f + m * t_m = f * t_f * (1 + t_m/t_f * 1/q)$

*Machine Balance:* **Key to machine efficiency**

- Larger $q$ means time closer to minimum $f * t_f$
  - $q \geq t_m/t_f$  needed to get at least half of peak speed
  - Speed is inverse of time so peak speed is $1/(f * t_f)$

Assume an element is a float and is equal to word!

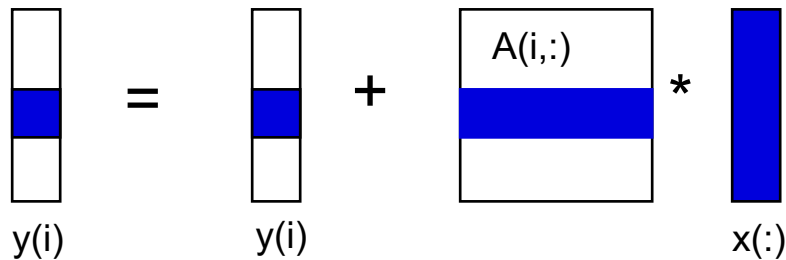# Warm up: Matrix-Vector Multiplication

{implements y = y + A*x}

for i = 1:n

       for j = 1:n

              y(i) = y(i) + A(i,j)*x(j)



y(i)      y(i)                x(:)

- m = ?
- f  = ?
- q  = ?

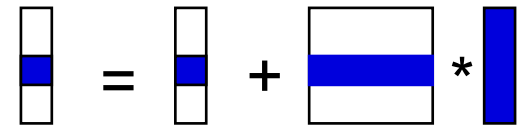# Warm up: Matrix-Vector Multiplication

{read x(1:n) into fast memory} **This line needs n memory references**

{read y(1:n) into fast memory} **This line needs n memory references**

for i = 1:n

    {read row i of A into fast memory} ← **This line needs n memory references but note that there is loop around this line so we do $n^2$ references in total from this line**

    for j = 1:n

        y(i) = y(i) + A(i,j)*x(j)

{write y(1:n) back to slow memory} **This line needs n memory references**

**More explanation: m is computed as follows: n + n + n (n) + n**

- m = number of slow memory refs = $3n + n^2$
- f = number of arithmetic operations = $2n^2$
- q = f / m ≈ 2
- Matrix-vector multiplication limited by slow memory speed

Assume for simplicity vectors x and y, and one row of A, all fit in fast memory

# Modeling Matrix-Vector Multiplication

- Examples of some architectures and their machine balance
- So the computational intensity of 2 in matrix-vector multiply means that we can not get close to half peak of these machines: Matrix-Vector Multiplication is a *memory bound operation!*

| | Clock | Peak | Mem Lat (Min,Max) | | Linesize | t_m/t_f |
|---|---|---|---|---|---|---|
| | MHz | Mflop/s | cycles | | Bytes | |
| Ultra 2i | 333 | 667 | 38 | 66 | 16 | 24.8 |
| Ultra 3 | 900 | 1800 | 28 | 200 | 32 | 14.0 |
| Pentium 3 | 500 | 500 | 25 | 60 | 32 | 6.3 |
| Pentium3M | 800 | 800 | 40 | 60 | 32 | 10.0 |
| Power3 | 375 | 1500 | 35 | 139 | 128 | 8.8 |
| Power4 | 1300 | 5200 | 60 | 10000 | 128 | 15.0 |
| Itanium1 | 800 | 3200 | 36 | 85 | 32 | 36.0 |
| Itanium2 | 900 | 3600 | 11 | 60 | 64 | 5.5 |

*machine balance (q must be at least this for ½ peak speed)*
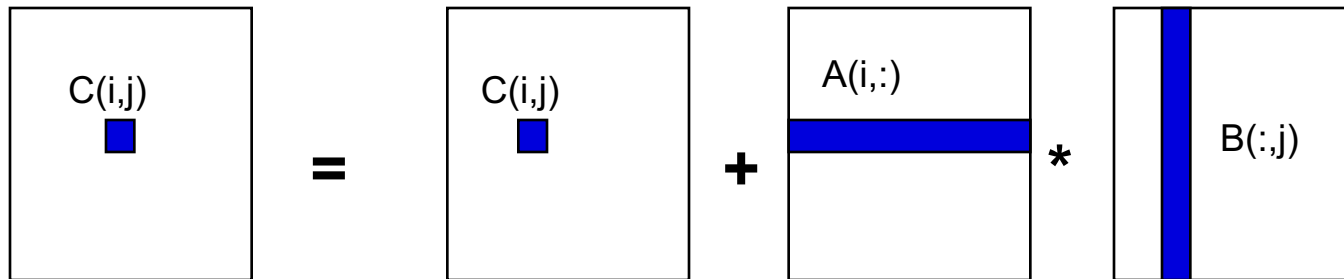
# Naïve Matrix Multiply

{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.
- If all the data would fit in fast memory (ideal case!) we would only make $3*n^2$ memory references because you only need to read each matrix once form slow memory and then it sticks around in fast memory!
- q potentially as large as⁂ $2*n^3 / 3*n^2 = O(n)$



Assume for simplicity that data is layed out in slow memory in the order it is being accessed. For example, here A is stored in row-major and B is stored in column-major.

⁂ Replace 3 with a 4 if you need to write C back to slow memory.

# Naïve Matrix Multiply

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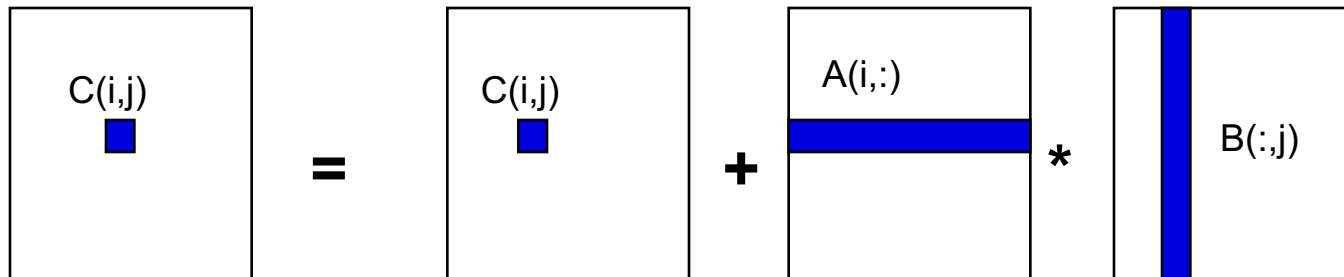  for j = 1 to n
      {read C(i,j) into fast memory}
      {read column j of B into fast memory}
      for k = 1 to n
          C(i,j) = C(i,j) + A(i,k) * B(k,j)
      {write C(i,j) back to slow memory}



m = ?

# 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}

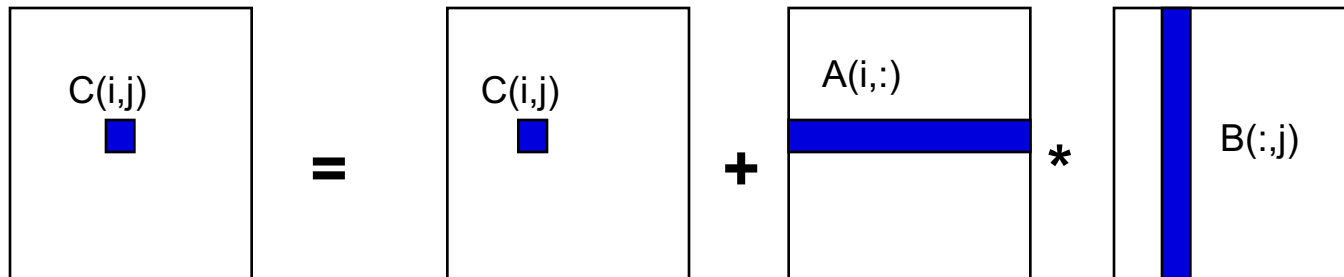     {read column j of B into fast memory}

     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}

**This line needs n memory references but note that there is a loop around this line so we do $n^2$ references in total for this line.**

**This line needs n memory reference but note that there are two loops around this line so we do $n^3$ references in total for this line**

**This line needs 1 memory reference but note that there are two loops around this line so we do $n^2$ references in total for this line**

C(i,j) = C(i,j) + A(i,:) * B(:,j)

$$m = n * (n + n * (1 + n + 1)) = n^3 + 3n^2$$

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
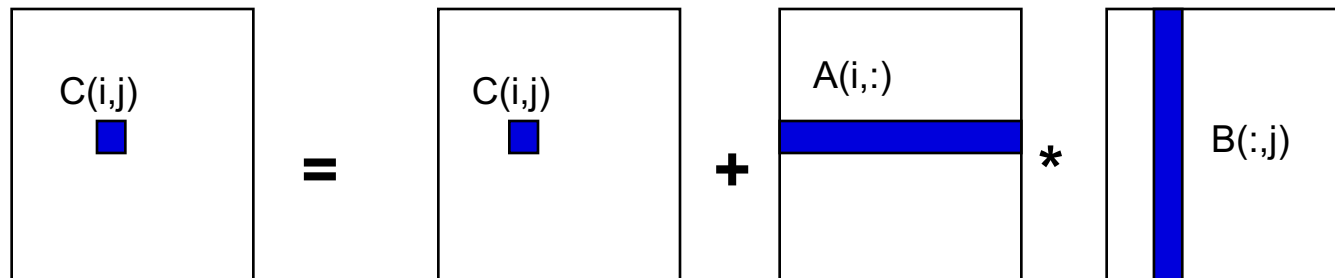
     $= n^3 + 3n^2$

So $q = f \, / \, m = 2n^3 \, / \, (n^3 + 3n^2)$

     $\approx 2$ 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



**Another way: B gets read n times, A gets read once, and C is read/written 2 times so m is n (n²)+ n² + 2 n²**

# 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}
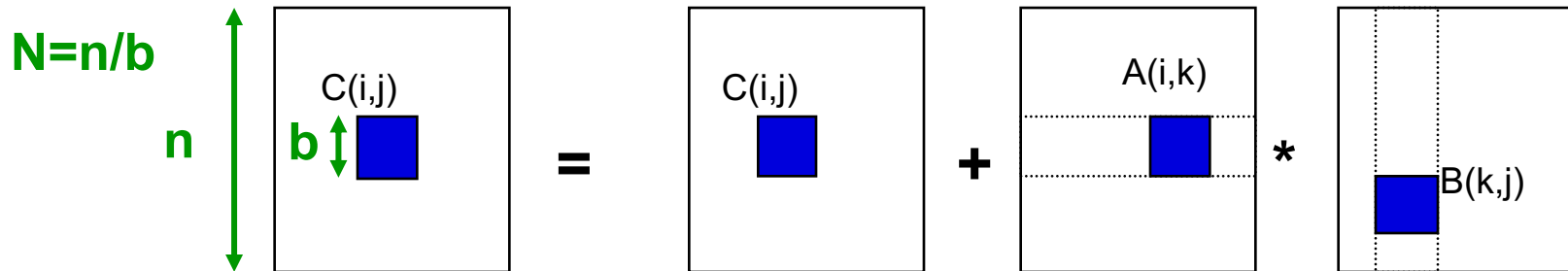                {read block B(k,j) into fast memory}
                C(i,j) = C(i,j) + A(i,k) * B(k,j) {do a matrix multiply on blocks}
            {write block C(i,j) back to slow memory}

**cache does this automatically**

**3 nested loops inside**

**block size = loop bounds**



N=n/b

n

b

C(i,j)    =    C(i,j)    +    A(i,k)    *    B(k,j)

$m = N * N ( (n/N)^2 + N * ((n/N)^2 + (n/N)^2) + (n/N)^2) = (2N+2) n^2$

*Tiling for registers (managed by you/compiler) or caches (hardware)*

13

# Blocked (Tiled) Matrix Multiply

Recall:

    $m$ is amount memory traffic between slow and fast memory

    matrix has nxn elements, and NxN blocks each of size bxb

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

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

So:

$m =$  $N*n^2$     read each block of B  $N^3$ times $(N^3 * b^2 = N^3 * (n/N)^2 = N*n^2)$

    $+ N*n^2$   read each block of A  $N^3$ times

    $+ 2n^2$     read and write each block of C once

    $= (2N + 2) * n^2$

So computational intensity $q = ?$

**Another way: If initially we assume each matrix as N by N elements (the blocks are each one element) then B and A each get read N times, and C is read/written 2 times. Since each element/block is size of $(n/N)^2$ so m is $N^3 (n/N)^2 + N^3 (n/N)^2 + 2 N^2 (n/N)^2$**

# Blocked (Tiled) Matrix Multiply

Recall:

    m is amount memory traffic between slow and fast memory

    matrix has nxn elements, and NxN blocks each of size bxb

    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

So:

$m = N*n^2$    read each block of B  $N^3$ times ($N^3 * b^2 = N^3 * (n/N)^2 = N*n^2$)

    $+ N*n^2$   read each block of A  $N^3$ times

    $+ 2n^2$     read and write each block of C once

    $= (2N + 2) * n^2$

So computational intensity $q = f / m = 2n^3 / ((2N + 2) * n^2)$

$$\approx n / N = b \text{ for large } n$$

So we can improve performance by increasing the blocksize b

Can have a much better computational intensity than matrix-vector multiply (q=2)

# Blocked (Tiled) Matrix Multiply

Recall:

   m is amount memory traffic between slow and fast memory

   matrix has nxn elements, and NxN blocks each of size bxb

   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

So:

 m =  N*$n^2$    read each block of B  $N^3$ times ($N^3 * b^2 = N^3 * (n/N)^2 = N*n^2$)

    + N*$n^2$   read each block of A  $N^3$ times

    + $2n^2$     read and write each block of C once

    =  (2N + 2) * $n^2$

**Follow the exact process you did for the non blocked version with the difference that now an element is extended to be a b by b block.**

So computational intensity q = f / m = $2n^3$ / ((2N + 2) * $n^2$)

                    $\approx$ n / N = b  for large n

So we can improve performance by increasing the blocksize b

Can have a much better computational intensity than matrix-vector multiply (q=2)

16

# Blocked (Tiled) Matrix Multiply

Recall:

m is amount memory traffic between slow and fast memory

matrix has nxn elements, and NxN blocks each of size bxb

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

So:

m =  N*$n^2$    read each block of B  $N^3$ times ($N^3$ * $b^2$ = $N^3$ * $(n/N)^2$ = N*$n^2$)

+ N*$n^2$   read each block of A  $N^3$ times

+ 2$n^2$     read and write each block of C once

=  (2N + 2) * $n^2$

**Follow the exact process you did for the non blocked version with the difference that now an element is extended to be a b by b block.**

So computational intensity q = f / m = $2n^3$ / ((2N + 2) * $n^2$)

$\approx$ n / N = b  for large n

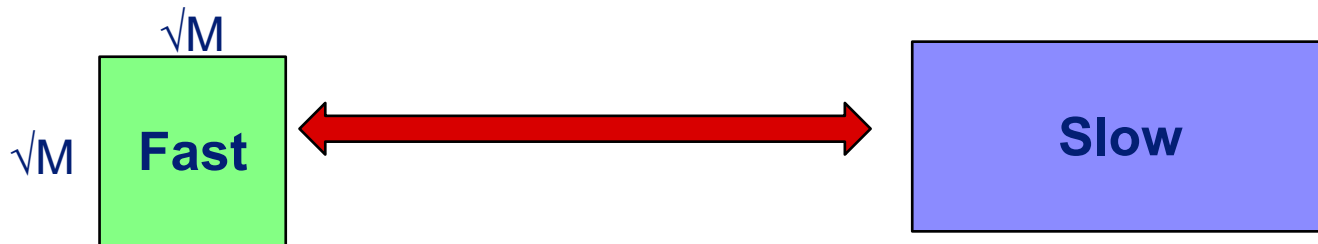**Why not increase b to a very large number?**

So we can improve performance by increasing the blocksize b

Can have a much better computational intensity than matrix-vector multiply (q=2)

# Limits to Optimizing Matrix Multiply

- The tiled matrix multiply analysis assumes that three tiles/blocks fit into fast memory at once.

- If $M_{fast}$ is the size of fast memory then the previous analysis shows that the blocked algorithm has computational intensity:

$$q \approx b \leq (M_{fast}/3)^{1/2}$$

# 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): 15 different operations**

    - **vector operations: dot product, saxpy ($y=\alpha*x+y$), etc**
    - **m=2*n, f=2*n, q = f/m = computational intensity ~1 or less**

# Basic Linear Algebra Subroutines (BLAS)
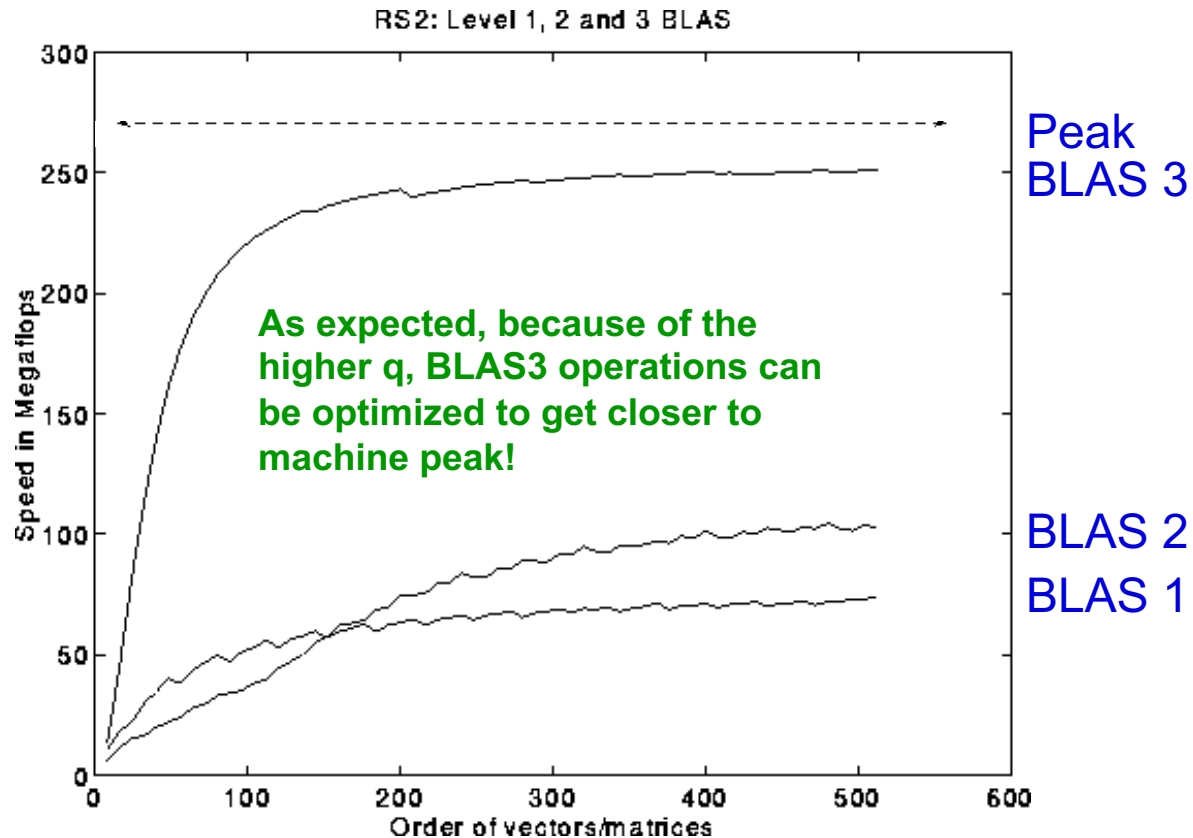
- **Industry standard interface (evolving)**
  - **www.netlib.org/blas,    www.netlib.org/blas/blast--forum**
- **Vendors, others supply optimized implementations**
- **History**
  - **BLAS1 (1970s): 15 different operations**

    - **vector operations: dot product, saxpy (y=$\alpha$*x+y), etc**
    - **m=2*n, f=2*n, q = f/m = computational intensity ~1 or less**
  - **BLAS2 (mid 1980s): 25 different operations**
    - **matrix-vector operations: matrix vector multiply, etc**
    - **m=n^2, f=2*n^2, q~2, less overhead**
    - **somewhat faster than BLAS1**

# 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): 15 different operations**
    - **vector operations: dot product, saxpy (y=$\alpha$*x+y), etc**
    - **m=2*n, f=2*n, q = f/m = computational intensity ~1 or less**
  - **BLAS2 (mid 1980s): 25 different operations**
    - **matrix-vector operations: matrix vector multiply, etc**
    - **m=n^2, f=2*n^2, q~2, less overhead**
    - **somewhat faster than BLAS1**
  - **BLAS3 (late 1980s): 9 different operations (such as matrix-matrix multiply or solving a triangular system or matrix factorization and so on)**
    - **matrix-matrix operations: matrix matrix multiply, etc**
    - **m <= 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 use BLAS3 when possible (LAPACK & ScaLAPACK)**
  - **See `www.netlib.org/{lapack,scalapack}`**

# BLAS speeds on an IBM RS6000/590

Peak speed = 266 Mflops

RS2: Level 1, 2 and 3 BLAS



Peak
BLAS 3

As expected, because of the higher q, BLAS3 operations can be optimized to get closer to machine peak!
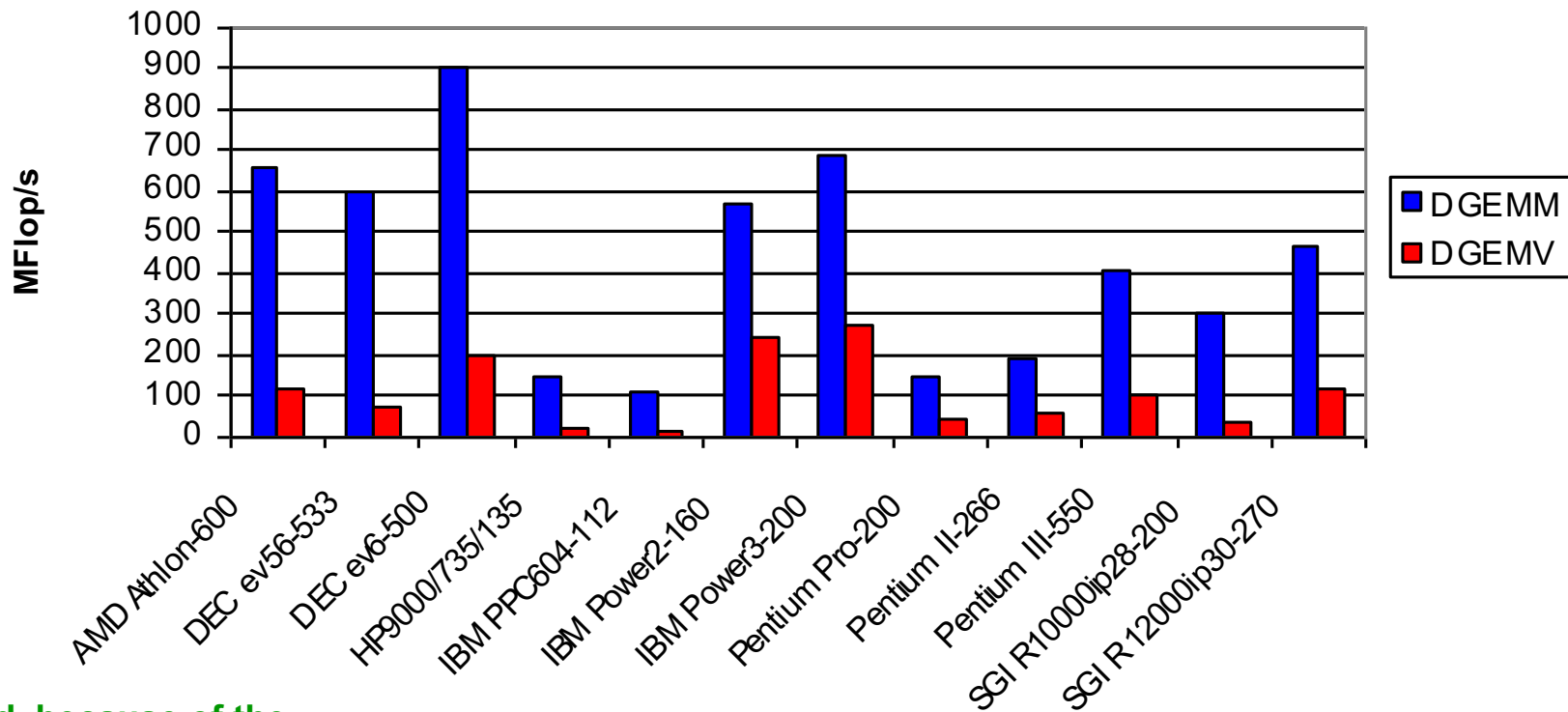
BLAS 2
BLAS 1

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

# Dense Linear Algebra: BLAS2 vs. BLAS3

- BLAS2 and BLAS3 have very different computational intensity, and therefore different performance

**BLAS3 (MatrixMatrix) vs. BLAS2 (MatrixVector)**



**As expected, because of the higher q, BLAS3 operations can be optimized to get closer to machine peak!**

Data source: Jack Dongarra