

Case Study – Matrix Multiplication





- › Serve as an example of design exploration of matrix multiplication
 - › While examples are for a processor with cache, they are equally valid for an FPGA with external memory
-

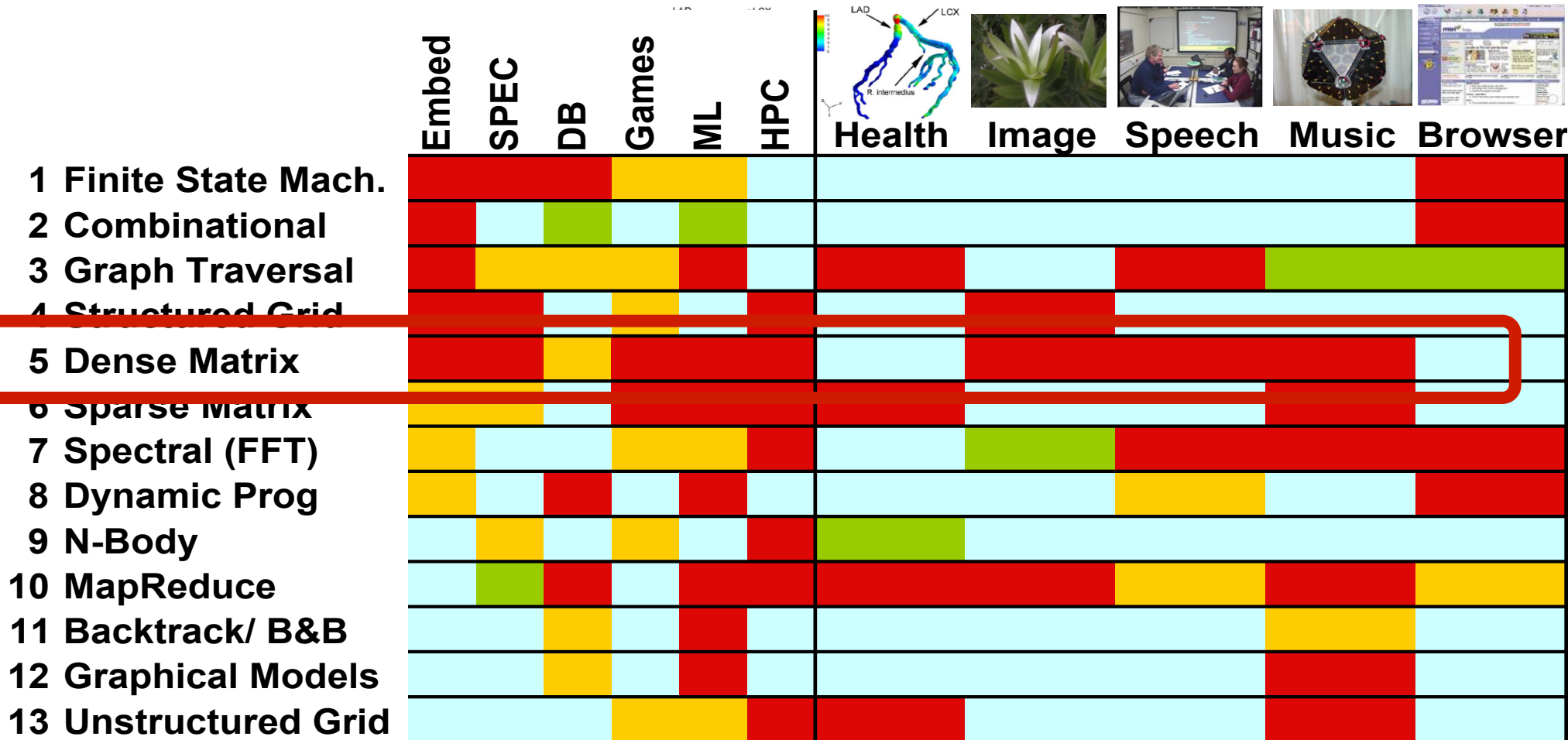


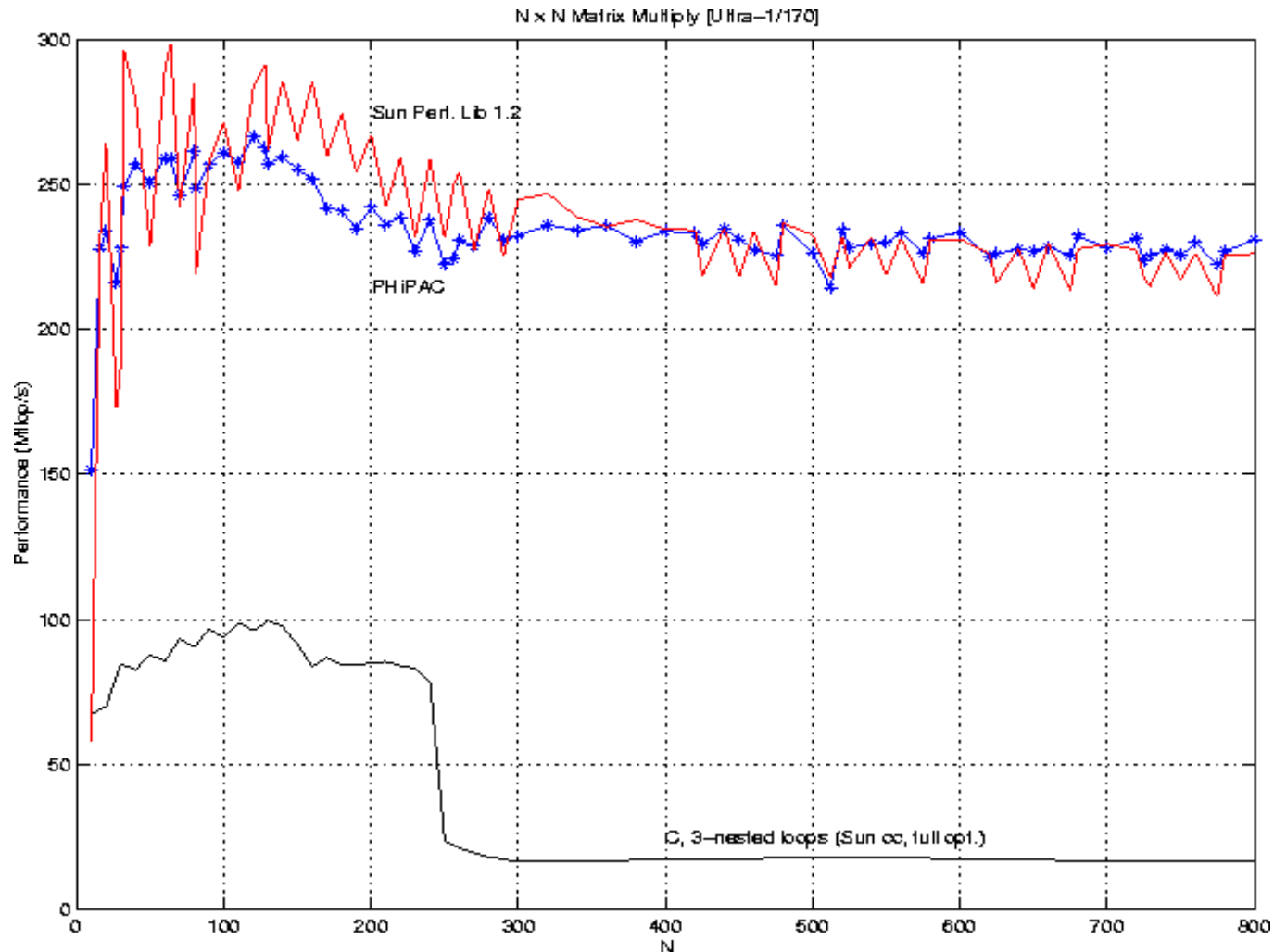
- › Performance Modeling
 - › Matrix-Vector Multiply (Warmup)
 - › Matrix Multiply Cache Optimizations
-

- › An important kernel in many problems
 - Appears in many linear algebra algorithms
 - Bottleneck for dense linear algebra
 - One of the 7 dwarfs / 13 motifs of parallel computing
 - Closely related to other algorithms, e.g., transitive closure on a graph using Floyd-Warshall
- › Optimization ideas can be used in other problems
- › The best case for optimization payoffs
- › The most-studied algorithm in high performance computing

Motif/Dwarf: Common Computational Methods

(Red Hot → Blue Cool)





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

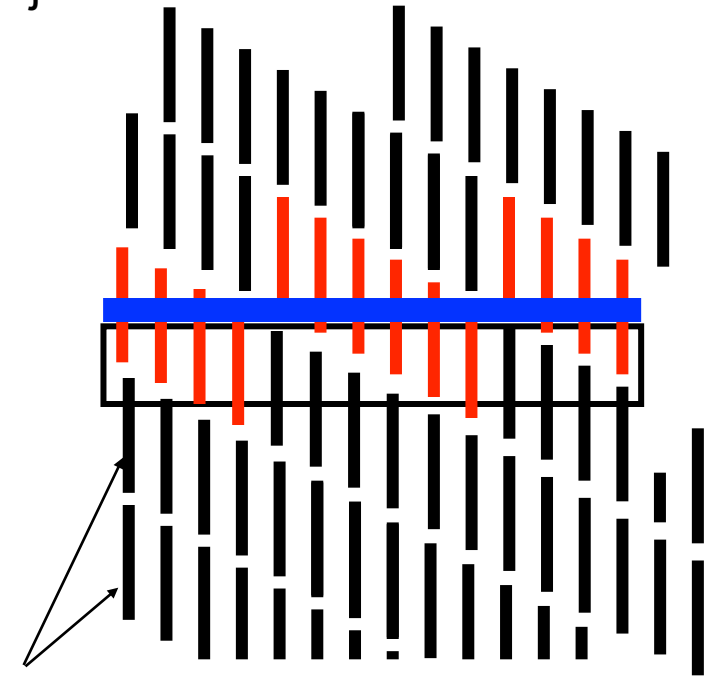
- › A matrix is a 2-D array of elements, but memory addresses are “1-D”
- › Conventions for matrix layout
 - by column, or “column major” (Fortran default); $A(i,j)$ at $A+i*j*n$
 - by row, or “row major” (C default) $A(i,j)$ at $A+i*n+j$
 - recursive (later)

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 (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 $\ll t_m$
 - $q = f / m$ average number of flops per slow memory access
- › Minimum possible time = $f * t_f$ when all data in fast memory
- › Actual time
 - $f * t_f + m * t_m = f * t_f * (1 + \frac{t_m}{t_f} * 1/q)$
- › 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

Computational Intensity: Key to algorithm efficiency

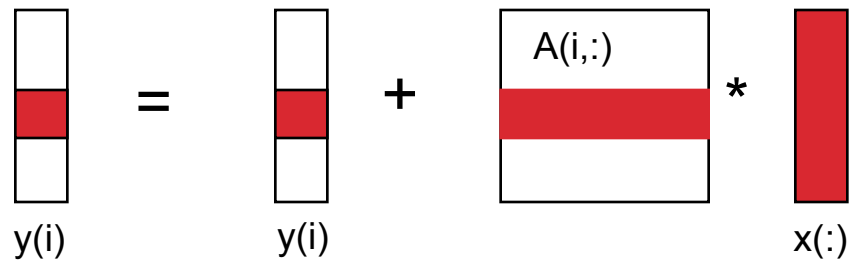
Machine Balance: Key to machine efficiency


```
{implements y = y + A*x}
```

```
for i = 1:n
```

```
    for j = 1:n
```

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



```
{read x(1:n) into fast memory}
{read y(1:n) into fast memory}
for i = 1:n
    {read row i of A into fast memory}
    for j = 1:n
        y(i) = y(i) + A(i,j)*x(j)
    {write y(1:n) back to slow memory}
```

- m = number of slow memory refs = $3n + n^2$
- f = number of arithmetic operations = $2n^2$
- $q = f / m \approx 2$

- Matrix-vector multiplication limited by slow memory speed

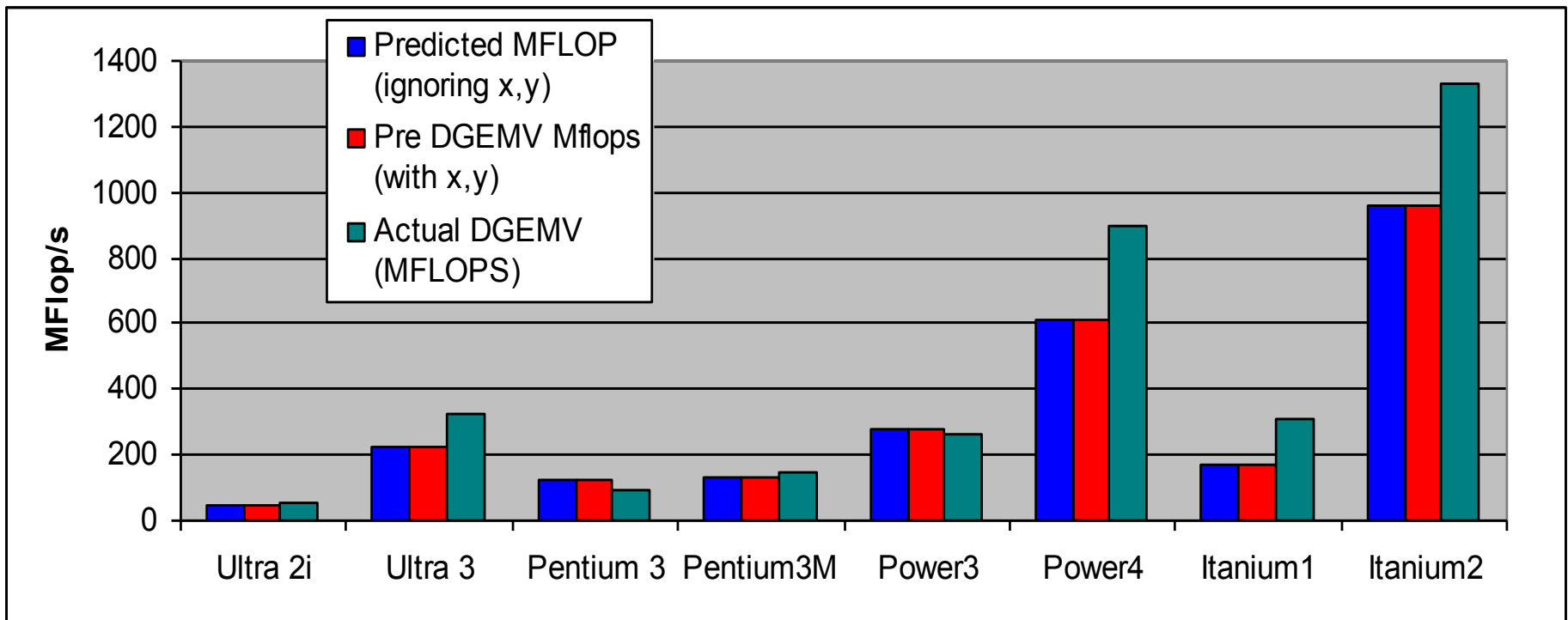
- › Compute time for $n \times n = 1000 \times 1000$ matrix
- › Time
 - $f * t_f + m * t_m = f * t_f * (1 + t_m/t_f * 1/q)$
 - $= 2 * n^2 * t_f * (1 + t_m/t_f * 1/2)$
- › For t_f and t_m , using data from R. Vuduc's PhD (pp 351-3)
 - <http://bebop.cs.berkeley.edu/pubs/vuduc2003-dissertation.pdf>
 - For t_m use minimum-memory-latency / words-per-cache-line

	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
Itanium 1	800	3200	36	85	32	36.0
Itanium2	900	3600	11	60	64	5.5

machine balance (q must be at least this for 1/2 peak speed)

- › What simplifying assumptions did we make in this analysis?
 - Ignored parallelism in processor between memory and arithmetic within the processor
 - Sometimes drop arithmetic term in this type of analysis
 - Assumed fast memory was large enough to hold three vectors
 - Reasonable if we are talking about any level of cache
 - Not if we are talking about registers (~32 words)
 - Assumed the cost of a fast memory access is 0
 - Reasonable if we are talking about registers
 - Not necessarily if we are talking about cache (1-2 cycles for L1)
 - Memory latency is constant
- › Could simplify even further by ignoring memory operations in X and Y vectors
 - Mflop rate/element = $2 / (2 * t_f + t_m)$

- › How well does the model predict actual performance?
 - Actual DGEMV: Most highly optimized code for the platform
- › Model sufficient to compare across machines
- › But under-predicting on most recent ones due to latency estimate



```
{implements C = C + A*B}
```

```
for i = 1 to n
```

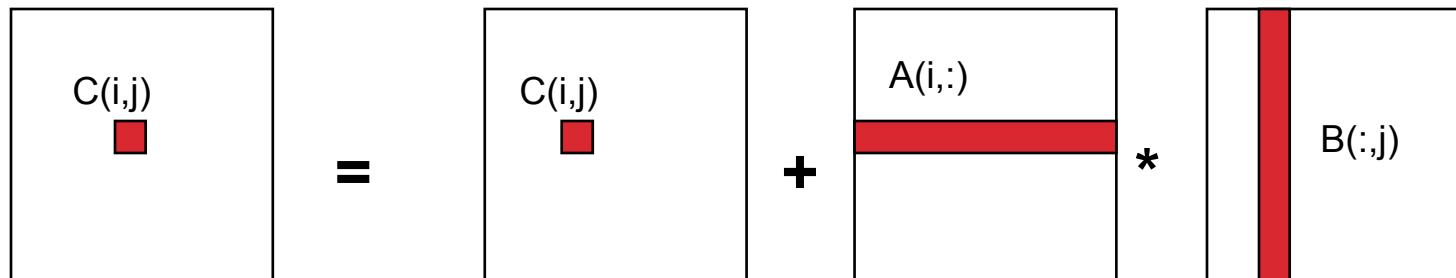
```
  for j = 1 to n
```

```
    for k = 1 to n
```

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

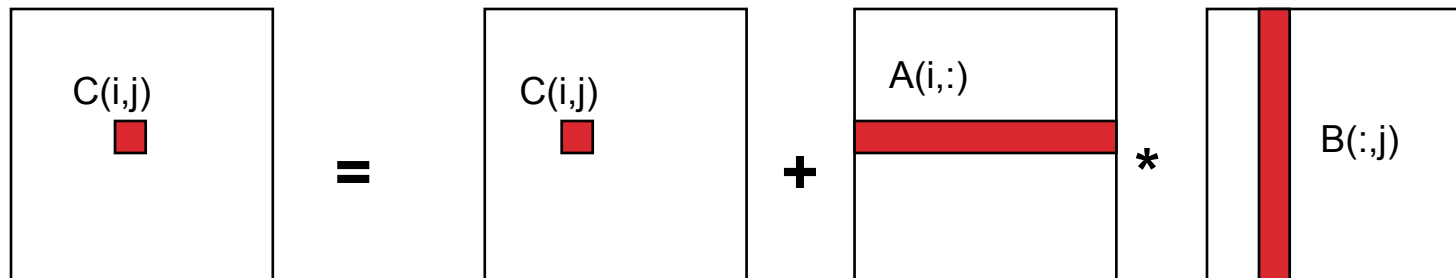
Algorithm has $2*n^3 = O(n^3)$ Flops and operates on
 $3*n^2$ words of memory

q potentially as large as $2*n^3 / 3*n^2 = O(n)$



```

{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}
  
```



Number of slow memory references on unblocked matrix multiply

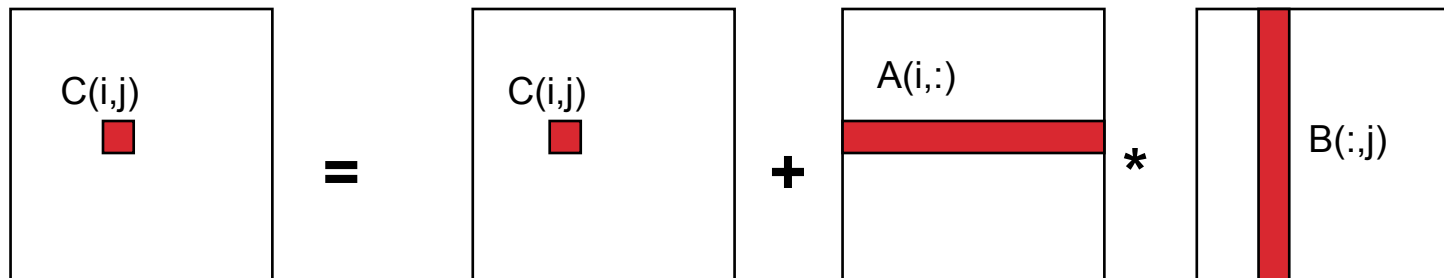
$$\begin{aligned}
 m &= n^3 && \text{to read each column of } B \text{ } n \text{ times} \\
 &+ n^2 && \text{to read each row of } A \text{ once} \\
 &+ 2n^2 && \text{to read and write each element of } C \text{ once} \\
 &= n^3 + 3n^2
 \end{aligned}$$

$$\text{So } q = f / m = 2n^3 / (n^3 + 3n^2)$$

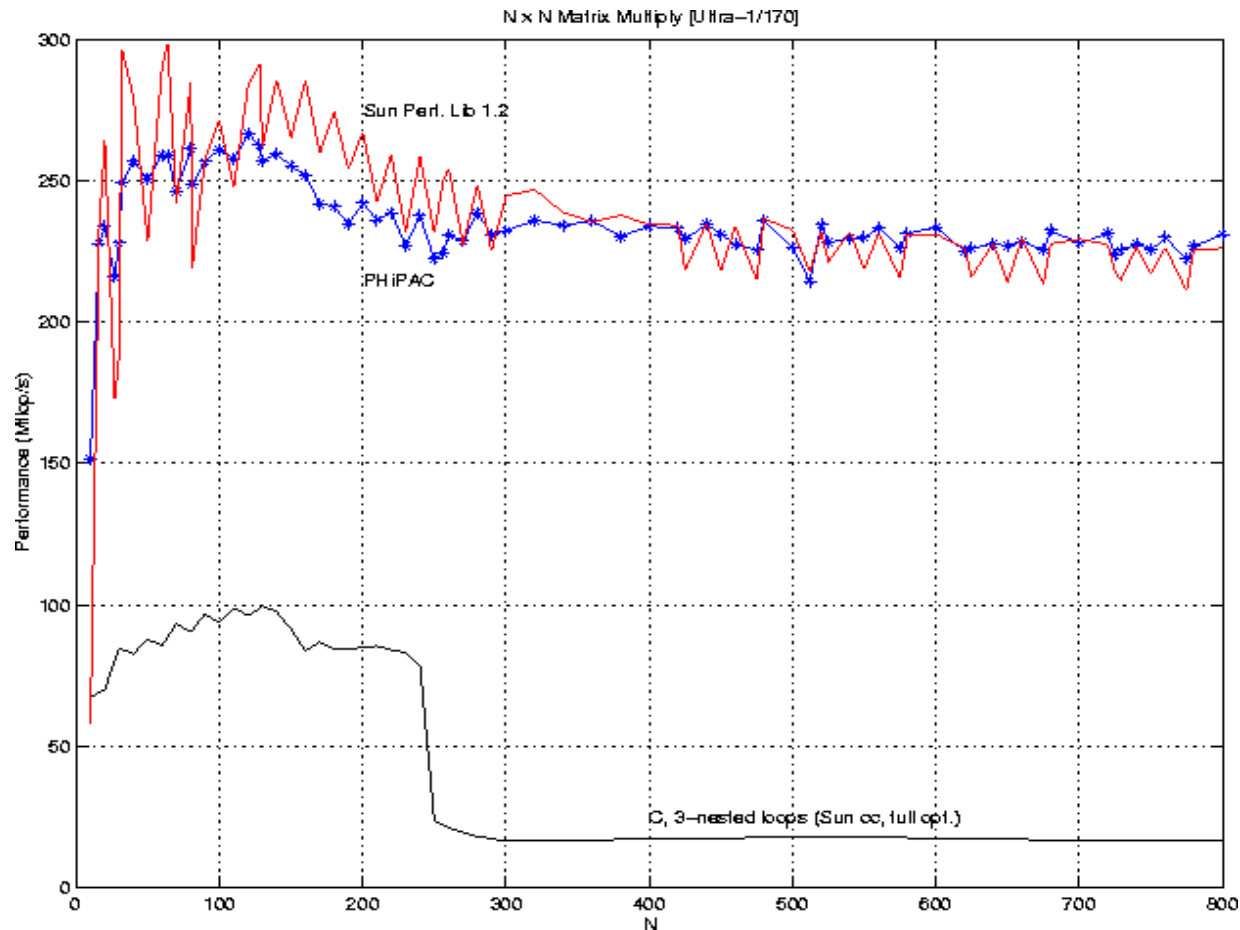
≈ 2 for large n , no improvement over matrix-vector multiply

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

Similar for any other order of 3 loops



Matrix-multiply, optimized several ways



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

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}

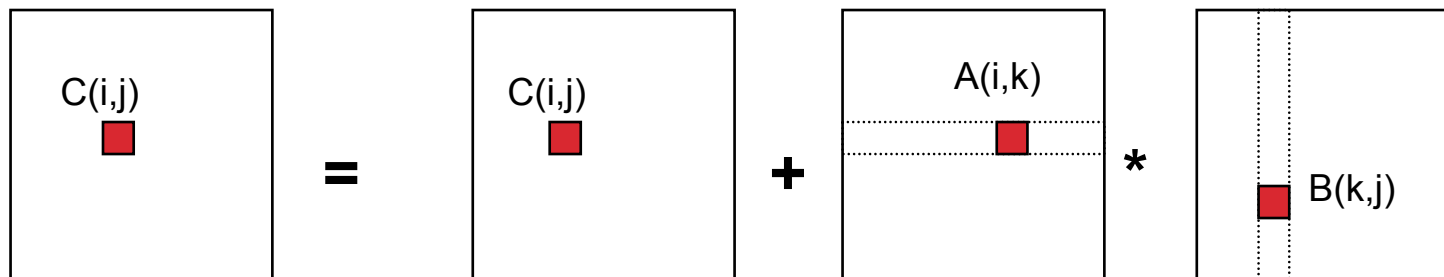
for $k = 1$ to N

{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}



Recall:

m is amount memory traffic between slow and fast memory

matrix has $n \times n$ elements, and $N \times N$ blocks each of size $b \times b$

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

So:

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

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 be much faster than matrix-vector multiply ($q=2$)

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}

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

- 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



- › Described a way to think about computation and memory – computational intensity
 - › Introduced the concept of blocking to increase computational intensity
-

- › Explain in your own words:
 - Computational intensity
- › Do a similar analysis computational intensity analysis for a different algorithm e.g. FFT