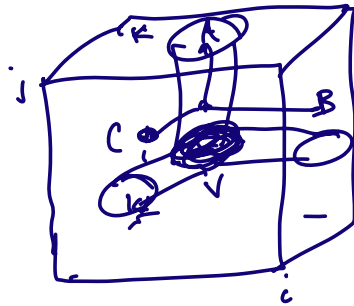


Welcome to Ma221! Lecture 10, Sep 15

Optimize Matmul: minimize comm  
between main mem + cache:

Do as many flops as possible given  
 $M$  words of data in cache



$$C(i,j) = A(i,k)B(k,j)$$

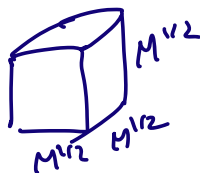
$$\# \text{ iterations} = |V| \leq \sqrt{|V_A| \cdot |V_B| \cdot |V_C|}$$

$$|V_A| + |V_B| + |V_C| = M$$

$$|V| \approx M^{3/2}$$

What "shape" of  $V$  attains this bound?

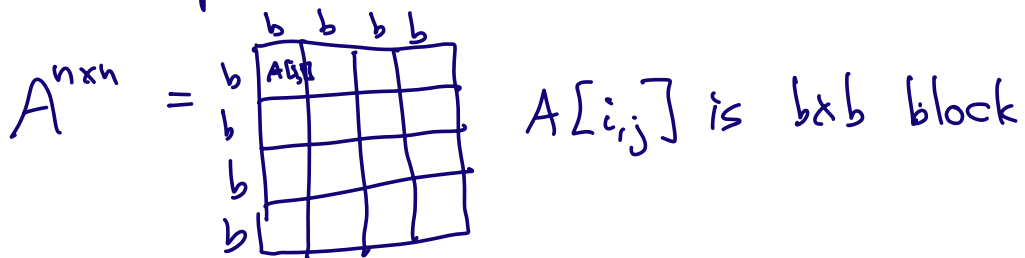
a cube:  $M^{1/2} \times M^{1/2} \times M^{1/2}$



Algorithm: break  $A, B, C$  into square  
submatrices that we can fit in Cache

read 3 blocks into cache  
 Do local mat mul, update C block  
 no data movement

repeat



for  $i = 1$  to  $n/b$   
 for  $j = 1$  to  $n/b$

read  $C[i,j]$  into cache, ...  $b^2$  words

for  $k = 1$  to  $n/b$

read  $A[i,k], B[k,j]$  into cache.

...  $2b^2$  words

$$C[i,j] = C[i,j] + A[i,k] \cdot B[k,j]$$

$b \times b$  matmul, 3 more loops  
 no data movement

end for

write  $C[i,j]$  to main mem --  $b^2$  words

end for

end for

$$\begin{aligned} \text{Total words moved} &= \overset{\text{read } C}{\frac{n}{b} \cdot \frac{n}{b} \cdot b^2} + \overset{\text{read } A, B}{\left(\frac{n}{b}\right)^3} 2b^2 + \overset{\text{write}}{\left(\frac{n}{b}\right)^2} b^2 \\ &= 2n^2 + 2\frac{n^3}{b} \end{aligned}$$

$$b \approx \sqrt{\frac{n}{3}}$$

What about rectangular case?  $m \times n \times k$   
What if some dimension  $< \sqrt{M}$ ?

Ex:  $y = Ax$  GEMV

Optimal Lower Bound =  $\Omega\left(\max\left(\frac{m \cdot n \cdot k}{\sqrt{M}}, \text{size input, size output}\right)\right)$

Idea of using Loomis-Whitney works for all dense LA, also extends to any algorithm that looks like nested loops accessing arrays "any" subscripts like  $i, i+j, i-2j+k \dots$

Can derive lower bound:

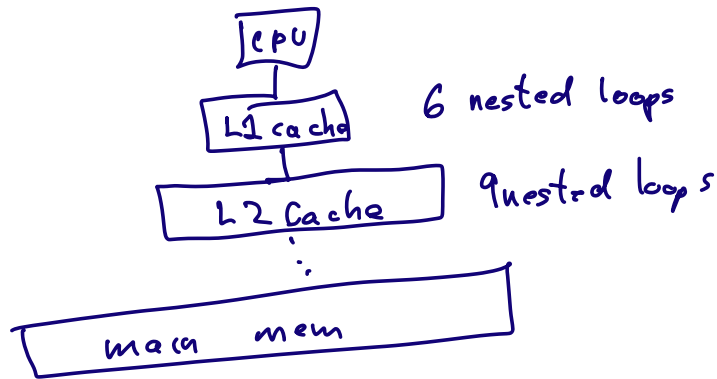
$$\# \text{ words moved} = \Omega\left(\frac{\# \text{ loop iterations}}{M^e}\right)$$

$e$  depends on details of alg, uses generalization of Loomis-Whitney called Brascamp-Lieb (extension by Christ, Tao et al)

still open questions...

One problem: need to know  $M$

What if there is more than one cache?



Goal: optimal matmul independent of HW

Use recursion

function  $C = \text{RMM}(A, B)$

... RMM = Recursive Matmul

... Simple:  $A^{n \times n}, B^{n \times n}, C^{n \times n}, n = 2^m$

$$\dots C = \begin{bmatrix} C_{11} & C_{12} \\ C_{21} & C_{22} \end{bmatrix}_{n/2}^{n/2}$$

... ditto for  $A, B$

$$\text{if } n=1, C=A \cdot B, \text{ else } C_{11} = \text{RMM}(A_{11}, B_{11}) + \text{RMM}(A_{12}, B_{21})$$

... ditto for  $C_{12}, C_{21}, C_{22}$

Correct by induction

Cost analysis:

$$A(n) = \# \text{ arithmetic ops}$$

$$= 8 A\left(\frac{n}{2}\right) + 4 \left(\frac{n}{2}\right)^2, \quad A(1) = 1$$

Solve:  $n = 2^m$

$$\frac{a(m)}{8^m} = \frac{A(2^m)}{8^m} = \frac{8 \cdot a(m-1)}{8^m} + \frac{2^{2m}}{8^m}$$

$$b(m)$$

$$b(m) = b(m-1) + \frac{1}{2^m} \quad \dots \text{geometric sum}$$

$$\text{get } A(n) = 2n^3 + \text{lower order term}$$

$$\begin{aligned} W(n) &= \# \text{ words moved} \\ &= 8W\left(\frac{n}{2}\right) + 4 \cdot 3 \left(\frac{n}{2}\right)^2 \end{aligned}$$

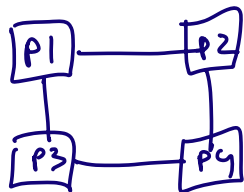
recurrence stops when matrices fit in cache

$$3 \cdot 2^m = M \quad W(b) = 3b^2 \text{ if } 3b^2 = M$$

$$\text{St. of geometric sum } W(n) = O\left(\frac{n^3}{\sqrt{M}}\right)$$

"Cache oblivious algorithm" works for any cache size  $M$ , any number of levels of cache

Lower bound on comm extends to parallel matmul



$P$  processors  
each stores  $\frac{1}{p}$  of all data  
each does  $\frac{1}{p}$  of all work

"fast memory" = local on processor

"slow memory" = memory on other procs.

$$\text{Same lower bound} = \Omega\left(\frac{\# \text{ flops per proc}}{\sqrt{\text{mem per proc}}}\right)$$

$$= \Omega\left(\frac{n^3/p}{\sqrt{3n^2/p}}\right) = \Omega\left(\frac{n^2}{\sqrt{p}}\right)$$

attainable by "SUMMA"  
 can be beaten by using more memory  
 and replicating data (see CS267)

---

Going faster than  $O(n^3)$  flops  
 Strassen (1967): matmul possible in  $O(n^{\log_2 7})$   
 operations  $\approx O(n^{2.81})$

trick: RMM using 7 recursive calls, not 8

$$A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}$$

$$\left. \begin{array}{l} P1 = (A_{12} - A_{22}) \cdot (B_{21} + B_{22}) \\ P2 = (A_{11} + A_{22}) \cdot (B_{11} + B_{22}) \\ \vdots \\ P7 = \end{array} \right\} \begin{array}{l} 7 \text{ recursive} \\ \text{calls} \end{array}$$

$$\left. \begin{array}{l} C_{11} = P1 + P2 - P4 + P6 \\ \vdots \\ C_{22} = \dots \end{array} \right\} \begin{array}{l} + 18 \text{ matrix} \\ \text{additions} \end{array}$$

$$A(n) = 7A\left(\frac{n}{2}\right) + 18\left(\frac{n}{2}\right)^2 = O(n^{\log_2 7})$$

$$W(n) = O\left(\frac{n^w}{M^{w/2-1}}\right) \quad w = \log_2 7$$

Thm (2010)  $W(n)$  attains lower bound

Thm (2015) Extends lower bound to all  
 "Strassen-like" algorithms

What is smallest  $w$ ? July 2023

$$w = 2.371552 \text{ (Williams + Xu + Xu + Zhou)}$$

(not practical)

Nature article last year: using ML  
to discover new matmul algs

Thm (2008): All linear algebra can be  
done in  $O(n^w)$

Error Analysis for Strassen

Classical:  $\|A \cdot B - A \cdot B\| \leq n \epsilon \|A\| \|B\|$

Strassen-like:  $\|A \cdot B - A \cdot B\| = O(\epsilon) \|A\| \|B\|$

Gauss's Trick for  $n$  matmul  
complex

$$(A + iB) \cdot (C + iD) = (T_1 - T_2) + i(T_3 - T_1 - T_2)$$

usual: 4 mat muls of real matrices

$$T_1 = A \cdot C \quad T_2 = B \cdot D$$

$$T_3 = (A+B) \cdot (C+D)$$

cost: 3 real mat muls + 5 adds

classical: 4 real " + 2 "

cost =  $\frac{3}{4}$  the classical formula