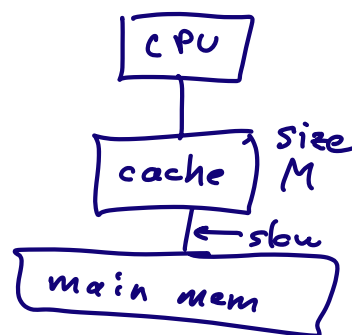


Welcome to Ma 221! Lec 7, Fall 24

Last time: considered matmul on



lower bound on # words moved
between cache and main mem

$$= \Omega\left(\frac{n^3}{\sqrt{M}}\right) = \Omega\left(\frac{\# \text{ flops}}{\sqrt{M}}\right)$$

enough for linear algebra, which looks

like 3 nested loops with $C(i,j) = C(i,j) + A(i,k) \cdot B(k,j)$

used Loomis-Whitney

Generalizes to any algorithm that

has nested loops (any number)

any number of arrays

any subscripts $A(i), B(i, i+j)$

$C(i+2j-3k, \dots)$

Can Get a lower bound and optimal alg

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

e depends on details of subscripts, need generalization of Loomis-Whitney to get e

called Hölder-Brascamp-Lieb inequality

Last time: optimal matmul, minimized communication by "blocking", breaking up matrices into $b \times b$ submatrices
 $b \sim \sqrt{M/3}$

Problems: Need to know about hardware:
depends on M , multiple levels of cache
may need 6 nested loops, or 9, or 12, ...

Goal: optimal matmul independent of HW

use recursion:

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

... RMM = Recursive Mat Mul

... Simplicity: assume $A^{n \times n}, B^{n \times n}, n = 2^m$

if $n = 1$

$$C = A \cdot B$$

else ...

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

$$C_{11} = \text{RMM}(A_{11}, B_{11}) + \text{RMM}(A_{12}, B_{21})$$

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

end: f

correct by induction

Cost analysis

$$A(n) = \# \text{ arithmetic ops for size } n \\ = 8 \cdot A\left(\frac{n}{2}\right) + 4 \left(\frac{n}{2}\right)^2, \quad A(1) = 1$$

solve for $n = 2^m$

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

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

geometric sum

answer is $2n^3 - n^2$ same as usual matmul

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

Base case when all 3 matrices fit in cache

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

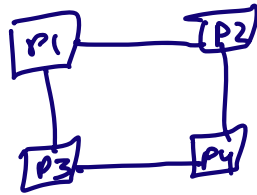
$$b = \sqrt{M/3}$$

Still geometric sum, $W(n) = O\left(\frac{n^3}{\sqrt{M}}\right)$

"Cache Oblivious": works for

any cache size, any # levels of cache,
works for much of linear algebra

Extension to Parallel case



want to minimize communication over network

each processor does $\frac{1}{P}$ of flops
"load balanced"

each processor stores $\frac{1}{P}$ of all data
"memory balanced"

"fast memory" = local to each proc

"slow memory" = all rest of memory on other procs

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

Surprise: can use more memory to replicate data

→ denominator bigger

→ lower bound lower : attainable!

Going faster than n^3 flops

Strassen (1967) matmul possible
in $O(n^{\log_2 7})$ flops $\sim O(n^{2.81})$

Trick: recursion with 7 calls, not 8

$$A = \begin{bmatrix} A_{11}^{n/2} & A_{12}^{n/2} \\ A_{21} & A_{22} \end{bmatrix}^{n/2}$$

$$\left. \begin{array}{l} P_1 = (A_{12} - A_{22}) \cdot (B_{21} + B_{22}) \\ P_2 = (A_{11} + A_{22}) \cdot (B_{11} + B_{22}) \\ \vdots \\ P_7 = \dots \end{array} \right\} \begin{array}{l} 7 \text{ recursive} \\ \text{calls} \end{array}$$

$$\left. \begin{array}{l} C_{11} = P_1 + P_2 - P_4 + P_6 \\ C_{12} = P_4 + P_5 \\ \vdots \end{array} \right\} \begin{array}{l} \text{total of} \\ 18 \text{ add/subs} \\ \text{of } (\frac{n}{2}) \times (\frac{n}{2}) \text{ matrices} \end{array}$$

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

$$\downarrow$$
$$A(n) = O(n^{\log_2 7}) \quad w = \log_2 7$$

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

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

Thm (2015) Extends to all

"Strassen-like" algs

Lastest record for smallest w

Thm (2024, Williams, Xu, Xu, Zhou) $w \approx 2.371552$

not practical, Strassen is practical

Thm (2008) All linear algebra can be done in $O(n^w)$ ops and "stably"

Error Analysis for Strassen:

Recall usual n^3 bound (Q1.10)

$$|fl(A \cdot B) - A \cdot B| \leq n \cdot \epsilon \cdot |A| \cdot |B|$$

Strassen:

$$\|fl(A \cdot B) - A \cdot B\| = O(\epsilon) \|A\| \cdot \|B\|$$

Gauss's trick for complex matmul

$$(A + iB) \cdot (C + iD)$$

$$T_1 = A \cdot C$$

$$T_2 = B \cdot D$$

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

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

cost : 3 matmuls + 5 add/sub

vs : 4 matmuls + 2 add/sub

cost = $\frac{3}{4}$ of usual matmul

Gaussian Elimination

Additional Goals, on top of
avoiding communication, $O(n^4)$ flops

: Backward stability: exact solution
of $(A+E)\hat{x} = b+f$

$$\frac{\|E\|}{\|A\|} = O(\epsilon), \quad \frac{\|f\|}{\|b\|} = O(\epsilon)$$

: Exploit math. structure of A

A : symmetric, positive definite

"sparse" = depends on $\ll n^2$ parameters

so could have many 0 entries,
or be dense and depend on $\ll n^2$
parameters

Eg: Vandermonde Matrix

$$V_{ij} = x_i^{j-1} \quad \text{given } [x_1, \dots, x_n]$$

Ex: multiply $V \cdot c = b$

= polynomial evaluation

$$b_i = \sum_j c_j x_i^{j-1}$$

solving $Vc = b$ for c

polynomial interpolation

much cheaper than $O(n^3)$

Seek Matrix Factorization

$A =$ product of simpler matrices

$$\text{SVD: } A = U \Sigma V^T \\ = \text{orthog.} \cdot \text{diag.} \cdot \text{orthog.}$$

$$\text{Gauss \& Lim } A = P \cdot L \cdot U$$

$P =$ permutation

$L =$ lower triangular
("unit", $L_{ii} = 1$)

$U =$ upper triangular

$$\text{Least Squares } A = QR \\ = \text{orthog.} \cdot \text{upper triang.}$$

$$\text{Eigenproblems } A = Q \cdot T \cdot Q^T \\ Q \text{ orthog.} \\ T \text{ triang.}$$

Def: Permutation Matrix
: identity matrix with permuted rows

Facts: let P, P_1, P_2 be perms.

P has exactly one 1 in each row
and each column

$P \cdot X = X$ with permuted rows

$X \cdot P = X$ with permuted columns

$P_1 \circ P_2 = \text{permutation}$

$P^{-1} = P^T$ i.e. P is orthogonal

proof: $(P^T P)_{ii} = 1 \Rightarrow P^T P = I$

$\det(P) = \pm 1$

storing and multiplying by P is cheap,
(store indices of $1s$, copy rows)

Thm: (LU decomposition)

Given any $m \times n$ full rank A , $m \geq n$

$\exists m \times m$ perm P

$m \times n$ unit lower triangular $L = \begin{bmatrix} 1 & & & \\ & \ddots & & \\ & & \ddots & \\ & & & 1 \end{bmatrix}$

$n \times n$ nonsingular U

s.t. $A = P \cdot L \cdot U$

to solve $Ax = b$, A square

(1) Factor $A = P \cdot L \cdot U$

(expensive part: # flops = $\frac{2}{3}n^3 + O(n^2)$)

(2) Solve $P \cdot L \cdot U \cdot x = b$ for $L \cdot U \cdot x = P^T \cdot b$
by permuting b , cost = $O(n)$

(3) Solve $L \cdot U \cdot x = P^T \cdot b$ for $U \cdot x = L^{-1} \cdot P^T \cdot b$
by forward substitution
cost = n^2

(4) Solve $U \cdot x = L^{-1} \cdot P^T \cdot b$ for $x = U^{-1} \cdot L^{-1} \cdot P^T \cdot b$
 by back substitution, cost = n^2

Given another b' can solve $Ax' = b'$
 in just $O(n^2)$ flops

Note: Should never compute A^{-1} and
 multiply by it because

(1) 3x more expensive in dense case
 (much worse in sparse case,
 up to $O(n^2)$ slower)

(2) not numerically stable

Proof of $A = P \cdot L \cdot U$ (Gauss Elim)

If A full rank \Rightarrow first col nonzero

$\Rightarrow \exists$ perm P s.t. $(PA)(1,1) \neq 0$

$$PA = \begin{array}{c|c} \begin{matrix} 1 & n-1 \\ \hline \end{matrix} \begin{matrix} A_{11} & A_{12} \\ \hline \end{matrix} \\ \begin{matrix} m-1 \\ \hline \end{matrix} \begin{matrix} A_{21} & A_{22} \\ \hline \end{matrix} \end{array} = \begin{array}{c|c} \begin{matrix} 1 & m-1 \\ \hline \end{matrix} \begin{matrix} 1 & 0 \\ \hline \end{matrix} \\ \begin{matrix} m-1 \\ \hline \end{matrix} \begin{matrix} \frac{A_{21}}{A_{11}} & I \\ \hline \end{matrix} \end{array} \cdot \begin{array}{c|c} \begin{matrix} 1 & n-1 \\ \hline \end{matrix} \begin{matrix} A_{11} & A_{12} \\ \hline \end{matrix} \\ \begin{matrix} m-1 \\ \hline \end{matrix} \begin{matrix} 0 & A_{22} - \frac{A_{21}A_{12}}{A_{11}} \\ \hline \end{matrix} \end{array}$$

$S =$ Schur
 Complement

A full rank $\Rightarrow PA$ full rank

$\Rightarrow S$ full rank: otherwise $\exists x \neq 0 : Sx = 0$

$\Rightarrow A \begin{bmatrix} -A_{21}x / A_{11} \\ x \end{bmatrix} = 0 \Rightarrow A$ not full rank
 contradiction

Simpler in square case:

$$\begin{aligned} 0 \neq \det(A) &= \pm \det(PA) = \pm \det(\text{1st factor}) \\ &\quad \cdot \det(\text{2nd factor}) \\ &= \pm \cdot 1 \cdot A_{11} \cdot \det(S) \end{aligned}$$

$$\Rightarrow \det(S) \neq 0$$

Apply induction to get $S = P' \cdot L' \cdot U'$

$$PA = \left[\begin{array}{c|c} 1 & 0 \\ \hline A_{21} & I \\ A_k & \end{array} \right] \cdot \left[\begin{array}{c|c} A_{11} & A_{12} \\ \hline 0 & \underbrace{P' \cdot L' \cdot U'}_S \end{array} \right]$$

$$= \left[\begin{array}{c|c} 1 & 0 \\ \hline A_{21} & P' \cdot L' \\ A_{11} & \end{array} \right] \cdot \left[\begin{array}{c|c} A_{11} & A_{12} \\ \hline 0 & U' \end{array} \right]$$

$$= \left[\begin{array}{c|c} 1 & 0 \\ \hline 0 & P' \end{array} \right] \cdot \left[\begin{array}{c|c} 1 & 0 \\ \hline P'^T A_{21} & L' \\ A_{11} & \end{array} \right] \cdot \left[\begin{array}{c|c} A_{11} & A_{12} \\ \hline 0 & U' \end{array} \right]$$

$$A = \underbrace{P^T \left[\begin{array}{c|c} 1 & 0 \\ \hline 0 & P' \end{array} \right]}_{\text{perm}} \cdot \underbrace{\left[\begin{array}{c|c} 1 & 0 \\ \hline P'^T A_{21} & L' \\ A_{11} & \end{array} \right]}_{\text{unit lower triang}} \cdot \underbrace{\left[\begin{array}{c|c} A_{11} & A_{12} \\ \hline 0 & U' \end{array} \right]}_{\text{upper triang}} \quad Q \in O$$