

Feb 10, 2021

$$m \begin{matrix} \uparrow \\ \boxed{A} \\ \downarrow \end{matrix} \quad n \begin{matrix} \uparrow \\ \boxed{x} \\ \downarrow \end{matrix}$$

① vectors $\alpha x, x+y, x^T y = \sum_{i=1}^n x_i y_i$ } $O(n)$ BLAS 1

② matrix-vector $y = Ax$ $\begin{bmatrix} a_{11} & \dots & a_{1n} \\ \vdots & & \vdots \end{bmatrix} \begin{matrix} \uparrow \\ \boxed{x} \\ \downarrow \end{matrix}$ $y = \sum_{j=1}^n a_{ij} x_j$ } $O(mn)$ BLAS 2
 $y_i = \sum_{j=1}^n a_{ij} x_j$ } $O(n^2)$

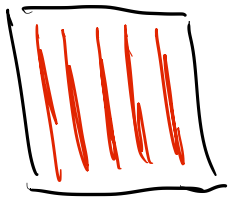
③ matrix-matrix $C = AB$ $c_{ij} = \sum_{k=1}^k a_{ik} b_{kj}$ } $O(mnk)$ BLAS 3
 $O(n^3)$

$$m \begin{matrix} \uparrow \\ \boxed{C} \\ \downarrow \end{matrix} \quad m \begin{matrix} \uparrow \\ \boxed{A} \\ \downarrow \end{matrix} \quad k \begin{matrix} \uparrow \\ \boxed{B} \\ \downarrow \end{matrix} \quad n$$

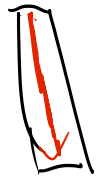
Basic linear algebra subroutines (BLAS)

Intel MKL

cuBLAS



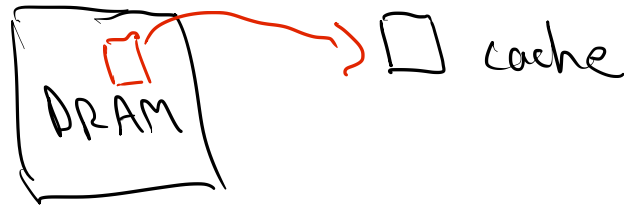
column major



store as contiguous array

Two main costs

- arithmetic (flops)
- fetch data (communication)



~ 100
cycles

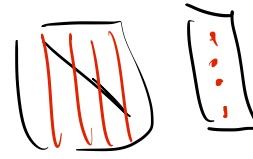
cost (communication) \gg cost (flops) ~ 1 cycle

- spatial locality: faster to access sequentially
- temporal locality: reuse data in cache

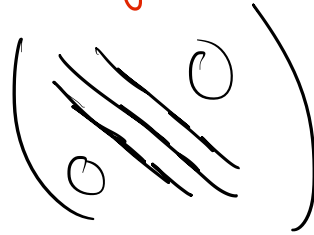
Structured matrices

sparse \Rightarrow many zeros $i=j$

Diagonal = $\begin{pmatrix} & & 0 \\ & & \\ 0 & & \end{pmatrix}$



Tridiagonal



$$I x = x$$

$$I = \begin{pmatrix} 1 & & \\ & \ddots & \\ & & 1 \end{pmatrix}$$