

Short notes on computational
Linear Algebra.

Fabio AFFINITO

SCAI dept.

CINECA

BASICS ALGORITHMS AND NOTATION ①

$$A \in \mathbb{R}^{m \times n} \Leftrightarrow A = (a_{ij}) = \begin{bmatrix} a_{11} & \dots & a_{1n} \\ \vdots & & \vdots \\ a_{m1} & & a_{mn} \end{bmatrix} \quad a_{ij} \in \mathbb{R}$$

MATRIX OPERATIONS :

- Transposition : $C = A^T \Leftrightarrow c_{ij} = a_{ji}$
- Addition : $C = A + B \Leftrightarrow c_{ij} = a_{ij} + b_{ij}$
- Scalar-Matrix multiplication : $C = \alpha A \Leftrightarrow c_{ij} = \alpha a_{ij}$
- Matrix-Matrix multiplication : $C = AB \Leftrightarrow c_{ij} = \sum_{k=1}^r a_{ik} b_{kj}$

VECTOR NOTATION

$$x \in \mathbb{R}^n \Leftrightarrow x = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \quad x_i \in \mathbb{R}$$

②

VECTOR OPERATIONS

- Dot product: $c = x^T y \Leftrightarrow \sum_{i=1}^N x_i y_i$
- Vector multiply: $z = x * y$
 $z_i = x_i y_i$ (Hadamard product)

• SAXPY

$$y = ax + y \quad \Leftrightarrow y_i = ax_i + y_i$$

... COMPUTATIONALLY SPEAKING ...

Dot product:

$$c = 0$$

for $i = 1:n$

$$c = c + x(i)y(i)$$

end

$$O(n)$$

SAXPY

for $i = 1:n$

$$y(i) = x(i)a + y(i)$$

end

$$O(n)$$

3

MATRIX - VECTOR

$$y = Ax + y \Rightarrow y_i = \sum_{j=1}^n a_{ij} x_j + y_i$$

MULTIPLICATION (GAXPY)

generalized

GAXPY (ROW VERSION)

$$A \in \mathbb{R}^{m \times n} \quad x \in \mathbb{R}^n \\ y \in \mathbb{R}^m$$

$$y \mapsto y = Ax + y$$

for $i = 1:m$

for $j = 1:n$

$$y(i) = A(i,j) \times (j) + y(i)$$

end

end

We can look at Ax as a combination of A 's columns

$$\begin{bmatrix} 1 \\ 3 \\ 5 \end{bmatrix} \begin{bmatrix} 2 \\ 4 \\ 6 \end{bmatrix} = 7 \begin{bmatrix} 1 \\ 3 \\ 5 \end{bmatrix} + 8 \begin{bmatrix} 2 \\ 4 \\ 6 \end{bmatrix}$$

GAXPY (COLUMN VERSION)

for $j = 1:n$

for $i = 1:m$

$$y(i) = A(i,j) \times (j) + y(i)$$

end

end

SAXPY!

4

MATRIX - MATRIX MULTIPLICATION

• Dot product formulation

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \begin{bmatrix} 5 & 6 \\ 7 & 8 \end{bmatrix} = \begin{bmatrix} 1 \cdot 5 + 2 \cdot 7 & 1 \cdot 6 + 2 \cdot 8 \\ 3 \cdot 5 + 4 \cdot 7 & 3 \cdot 6 + 4 \cdot 8 \end{bmatrix}$$

• Saxxy version (each column is lin. comb of columns of A)

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \begin{bmatrix} 5 & 6 \\ 7 & 8 \end{bmatrix} = \begin{bmatrix} 5 \\ 7 \end{bmatrix} + 7 \begin{bmatrix} 2 \\ 4 \end{bmatrix}; 6 \begin{bmatrix} 1 \\ 3 \end{bmatrix} + 8 \begin{bmatrix} 2 \\ 4 \end{bmatrix}$$

• Outer product version

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \begin{bmatrix} 5 & 6 \\ 7 & 8 \end{bmatrix} = \begin{bmatrix} 1 \\ 3 \end{bmatrix} \begin{bmatrix} 5 & 6 \end{bmatrix} + \begin{bmatrix} 2 \\ 4 \end{bmatrix} \begin{bmatrix} 7 & 8 \end{bmatrix}$$

5

MATRIX MULTIPLICATION : ijk VARIANT

$$C = A \cdot B + C$$

for $i = 1:m$

for $j = 1:n$

for $k = 1:p$

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

end

end

end

} 3 loops : $3! = 6$
variations !



Each variant involves the same amount of floating point operations but accesses the A, B, C differently

$$\lfloor a_m^T \rfloor$$

$$B = [b_1, \dots, b_n] \quad b_k \in \mathbb{R}^p$$

and then Algorithm 1.1.6 has this interpretation:

```

for i = 1:m
  for j = 1:n
    cij = aiTbj + cij
  end
end

```

Note that the "mission" of the j -loop is to compute update. To emphasize this we could write

```

for i = 1:m
  ciT = aiTB + ciT
end

```

$$C = \begin{bmatrix} c_1^T \\ \vdots \\ c_m^T \end{bmatrix}$$

where c_i is a row partitioning of C . To say the same thing we write

```

for i = 1:m
  C(i,:) = A(i,:)B + C(i,:)
end

```

Either way we see that the inner two loops of the row-oriented gaxpy operation.

is the jki variant. Each of the six possibilities (ijk , jik , ikj , jkj , kij , kji) features an inner loop operation (dot product or saxpy) and has its own pattern of data flow. For example, in the ijk variant, the inner loop over k overwrites a dot product that requires access to a row of A and a column of B . The jki variant involves a saxpy that requires access to a column of C and a column of A . These attributes are summarized in Table 1.1.1 along with an interpretation of what is going on when the middle and inner loop are considered together. Each variant involves the same amount of floating

Loop Order	Inner Loop	Middle Loop	Inner Loop Data Access
ijk	dot	vector \times matrix	A by row, B by column
jik	dot	matrix \times vector	A by row, B by column
ikj	saxpy	row gaxpy	B by row, C by row
jki	saxpy	column gaxpy	A by column, C by column
kij	saxpy	row outer product	B by row, C by row
kji	saxpy	column outer product	A by column, C by column

TABLE 1.1.1. Matrix Multiplication: Loop Orderings and Properties

point arithmetic, but accesses the A , B , and C data differently.

1.1.12 A Dot Product Formulation

The usual matrix multiplication procedure regards AB as an array of dot products to be computed one at a time in left-to-right, top-to-bottom order.

7

MATRIX MULTIPLICATION : DOT PRODUCT VERSION

$$C = AB + C$$

for $i = 1:m$

for $j = 1:n$

$$C(i,j) = A(i,1)C(i,j) + C(i,j)$$

end

for $j = 1:n$

$$C_{ij} = a_i^T b_j + c_{ij}$$

end

end

Partitioned matrix notation:

$$A(k, :) = [a_{k1} \dots a_{kn}]$$

$$A(:, k) = \begin{bmatrix} a_{1k} \\ \vdots \\ a_{mk} \end{bmatrix}$$

for $i = 1:m$

$$C_i^T = a_i^T B + c_i^T$$

end

IT'S A GAXPY

8

MATRIX MULTIPLICATION ; SAXPY VERSION

$$A = [a_1 \dots a_p] \quad a_j \in \mathbb{R}^m$$

$$C = [c_1 \dots c_n] \quad c_j \in \mathbb{R}^m$$

$$c_j = \sum_{k=1}^p b_{kj} a_k + c_j \quad j=1:n$$

(see slide 4)

for $j=1:n$

for $k=1:p$

$$C(:, j) = A(:, k) B(k, j) + C(:, j)$$

end
end

}

for $j=1:n$

$$C(:, j) = A \cdot B(:, j) + C(:, j)$$

end
SAXPY

9

MATRIX MULTIPLICATION : OUTER PRODUCT

$$A \in \mathbb{R}^{m \times p}$$

$$B \in \mathbb{R}^{p \times n}$$

$$C \in \mathbb{R}^{m \times n}$$

for $k = 1:p$

$$\Rightarrow C = A(:, k)B(k, :) + C$$

end

for $k = 1:p$

for $j = 1:n$

for $i = 1:m$

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

end
end
end

$$C = a_k b_k^T + C$$

EXPLOIT THE STRUCTURE !

TRIANGULAR MATRIX MULTIPLICATION

$$C = \begin{bmatrix} a_{11}b_{11} & a_{11}b_{12} + a_{12}b_{22} & a_{11}b_{13} + a_{12}b_{23} + a_{13}b_{33} \\ 0 & a_{22}b_{22} & a_{22}b_{23} + a_{23}b_{33} \\ 0 & 0 & a_{33}b_{33} \end{bmatrix}$$

$$C = A \cdot B = \begin{pmatrix} \text{shaded triangle} \end{pmatrix} \cdot \begin{pmatrix} \text{shaded triangle} \end{pmatrix}$$

$$a_{ik} b_{kj} = 0$$

if $k < i$
or $j < k$

$$c_{ij} = \sum_{k=i}^j a_{ik} b_{kj}$$

$c = 0$
 for $i = 1:n$
 for $j = i:n$
 for $k = i:j$
 $c(i,j) = A(i,k)B(k,j) + c(i,j)$
 end x 3

11

HOW MUCH DID WE SAVE? COUNT THE FLOPS

FLOP = Floating Point Operation

Dot product (sexpy) = $2n$ flops
Gexpy (outer product) = $m \times n$ flops
 $C = A \cdot B + C = 2 \times m \times n \times p$ flops

$A, B \in \mathbb{R}^{n \times n}$

$$\sum_{i=1}^n \sum_{j=1}^n 2(j-i+1) =$$

$$= \sum_{i=1}^n \sum_{j=i+1}^n 2j \approx$$

$$\approx \frac{2(n-i+1)^2}{2} = \sum_{i=1}^n i^2 = \frac{n^3}{3}$$

$C = 0$

for $i = 1:n$

for $j = i:n$

for $k = i:j$

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

end

end

end

$$\sum_{i=1}^n p \approx \frac{9^2}{2}$$

$$\sum_{i=1}^n p^{2n} \approx \frac{9^3}{3}$$

BLOCK MATRICES

We can partition a $m \times n$ matrix in blocks: $A = \begin{bmatrix} A_{11} & \dots & A_{1r} & \dots & A_{1q} \\ \vdots & & \vdots & & \vdots \\ A_{q1} & \dots & A_{qr} & \dots & A_{qj} \end{bmatrix}$ $\begin{matrix} m_1 & & n_r & & m_q \end{matrix}$

where $m_1 + \dots + m_q = m$
 $n_1 + \dots + n_r = n$

$A_{\alpha\beta}$ is a block of dimension $m_\alpha \times n_\beta$

$A_{\alpha\beta}$ is a $q \times r$ block matrix

$$A = \begin{bmatrix} A_{11} & \dots & A_{1s} & \dots & A_{1s} \\ \vdots & & \vdots & & \vdots \\ A_{q1} & \dots & A_{qs} & \dots & A_{qs} \end{bmatrix} \begin{matrix} m_1 & & m_1 & & m_q \end{matrix}; \quad B = \begin{bmatrix} B_{11} & \dots & B_{1r} \\ \vdots & & \vdots \\ B_{s1} & \dots & B_{sr} \end{bmatrix} \begin{matrix} p_1 & & n_r & & p_s \end{matrix}$$

TH: If

$$C = \begin{bmatrix} C_{11} & \dots & C_{1r} \\ \vdots & & \vdots \\ C_{q1} & \dots & C_{qr} \end{bmatrix} \begin{matrix} m_1 & & n_r & & m_q \end{matrix}$$

\Rightarrow

$$C_{\alpha\beta} = \sum_{\gamma=1}^s A_{\alpha\gamma} B_{\gamma\beta}$$

$\alpha = 1:q$
 $\beta = 1:r$
 (proof: cfr Globb-VL Pg. 27)

BLOCK MATRIX MULTIPLICATION

Let's suppose A, B, C are all $N \times N$ made of blocks of size ℓ

$$A = (A_{\alpha\beta})$$

$$B = (B_{\gamma\delta})$$

$$C = (C_{\alpha\beta})$$

$$C_{\alpha\beta} = \sum_{\gamma=1}^N A_{\alpha\gamma} B_{\gamma\beta} + C_{\alpha\beta}$$

$$\alpha = 1:N$$

$$\beta = 1:N$$

for $\alpha = 1:N$

$i = (\alpha-1)\ell + 1 : \alpha\ell$

for $\beta = 1:N$

$j = (\beta-1)\ell + 1 : \beta\ell$

for $\gamma = 1:N$

$k = (\gamma-1)\ell + 1 : \gamma\ell$

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

end

end

end

N.B. for $\ell = 1$

$$\alpha \equiv i$$

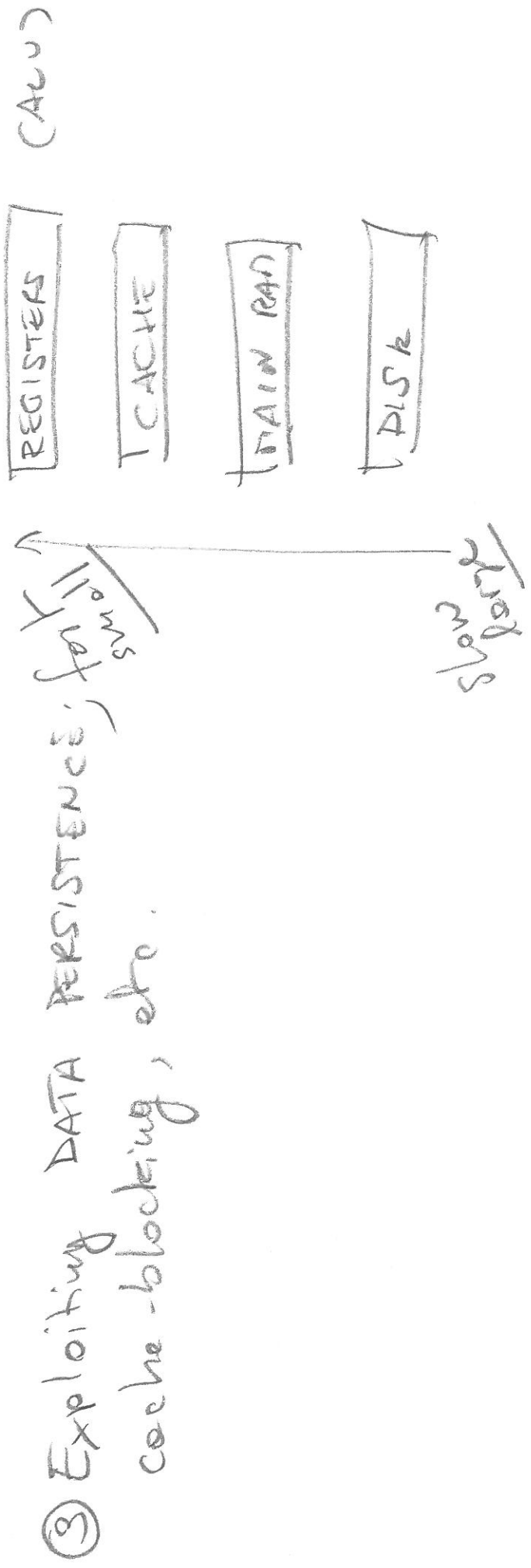
$$\beta \equiv j$$

$$\gamma \equiv k$$

we recover the well know algorithm

Optimization KEY FACTORS

- ① Computational complexity: how many flops?
- ② Access to data: stride, vector units...



slow
fast

BLOCKING AND RE-USE

$$C = A \cdot B + C$$

$A, B, C \in \mathbb{R}^{n \times n}$ reside

in the main memory

$n \gg$ cache dimension \rightarrow

$\#M$

Strategy:

$$B = \begin{bmatrix} B_1 & \dots & B_\ell \\ \vdots & & \vdots \\ \vdots & & \vdots \end{bmatrix}$$

$$C = \begin{bmatrix} C_1 & \dots & C_\ell \\ \vdots & & \vdots \\ \vdots & & \vdots \end{bmatrix}$$

$$n = \ell \ell$$

$$C_\alpha = A B_\alpha + C_\alpha =$$

$$= \sum_{k=1}^n A(:, k) B_\alpha(k, :) + C_\alpha$$

16

$$C_{\alpha} = \sum_{k=1}^n A(:, k) B_{\alpha}(k, :) + C_{\alpha}$$

```
for a = 1:n
  load Ba and Ca in cache
  for k = 1:n
    load A(:, k) in cache
    Ca = A(:, k) Ba(k, :) + Ca
  end
  store Ca in main memory
end
```

(17)

If M cache size, then we must have =

$$\boxed{2 \cdot nl + n \leq M} \quad A(i, k)$$

B_a, C_a

$T_1 \equiv \#$ of floating point numbers flowing btw cache and RAM

- Every entry of B is loaded into cache once
- Every entry of C is loaded into cache once and stored back in RAM once
- Every entry of A is loaded into cache $N = \frac{n}{l}$ times

$$n^2 + n^2 + n^2 + n^2 \cdot \frac{n}{l} =$$

$$= 3n^2 + \frac{n^3}{l}$$

18

$$\Gamma_1 = 3n^2 + \frac{n^3}{\rho}$$

$$2nl + n \leq M$$

We want to maximize l with the constraint \rightarrow

$$l \approx \frac{1}{2} \left(\frac{M}{n} - 1 \right)$$

$$\Rightarrow \Gamma_1 \approx 3n^2 + \frac{2n^4}{M-n}$$

Using $\Gamma_1 = 3n^2 + \frac{n^3}{\rho}$ we can see 2 extreme cases

1) $l = n$ (large cache, containing the whole matrix)

$$\Rightarrow \Gamma \sim n^2$$

2) $l = 1$ (small cache, containing only some columns...)

$$\Rightarrow \Gamma \sim n^3$$

19

$$C_{\alpha\beta} = \sum_{\gamma=1}^N A_{\alpha\gamma} B_{\gamma\beta}$$

$$\alpha = 1:N \quad \beta = 1:N$$

$$A = (A_{\alpha\beta}), B = (B_{\alpha\beta}) \in \mathbb{R}^{N \times N}$$

for $\alpha = 1:N$

for $\beta = 1:N$

load $C_{\alpha\beta}$ into cache

for $\gamma = 1:N$

load $A_{\alpha\gamma}$ and $B_{\gamma\beta}$ into cache

$$C_{\alpha\beta} = C_{\alpha\beta} + A_{\alpha\gamma} B_{\gamma\beta}$$

end

store $C_{\alpha\beta}$ in main memory

end

end

$$T_2 = 2u^2 + \frac{2u^3}{d}$$

(20)

- each entry in A and B is loaded $N = \frac{N}{\rho}$ times
- each entry of C is loaded once and stored once

$$\Rightarrow T_2 = 2n^2 + 2n^2 \left(\frac{N}{\rho} \right) = \\ = 2n^2 + 2 \frac{n^3}{\rho}$$

Constraint: we want to put a whole block in
the

$$3l^2 \leq M \quad \Rightarrow \quad l \leq \sqrt{\frac{M}{3}}$$

$$T_2 \cong 2n^2 + 2n^3 \sqrt{\frac{3}{M}}$$

$$T_1 = 3n^2 + \frac{2n^4}{H-n}$$

$$T_2 = 2n^2 + 2n^3 \sqrt{\frac{3}{H}}$$

$$\Rightarrow \frac{T_1}{T_2} = \frac{n^2}{n^2}$$

$$\frac{3 + \frac{2n^2}{H-n}}{2 + 2\sqrt{3} \sqrt{\frac{n^2}{H}}}$$

$\frac{\text{matrix size}}{\text{cache size}}$

key quantity is $\frac{n^2}{H}$ = ratio

$$\frac{T_1}{T_2} \approx \frac{n}{\sqrt{3H}}$$

If $\frac{n^2}{H} \gg 1 \Rightarrow$

→ The second strategy is superior from the standpoint of data motion to and from the cache.

References:

- GOLUB, VAN LOAN
Matrix computations
John Hopkins Ed.