

## Calcolo Scientifico: II lezione

sviluppo di software

efficiente

per le operazioni di base  
del calcolo matriciale

1

## Efficienza del software...

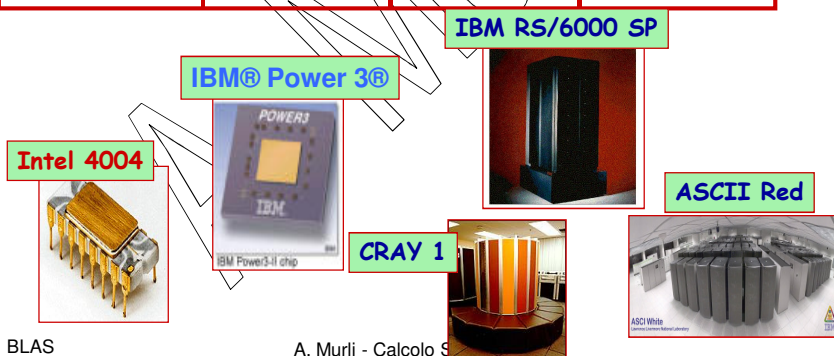
$$T_s \cong k \cdot \mu \cdot T(N)$$

Tecnologia  
Hardware

Algoritmo  
Software

## Evoluzione della tecnologia hardware

	1970	1995	2010
<b>Velocità del proc</b>	20 M flops	6 G Flops	10 Tera flops
<b>Memoria RAM</b>	200 K words	50 M words	100 G words
<b>Hard disk</b>	200 M words	500 G words	1 Peta words



## Evoluzione della tecnologia software

	1945	1955	1965	1975	1985
<b>Tempo</b>	$2 \times 10^6$ anni	20 anni	1 giorno	12 ore	0.2 sec
<b>Memoria</b>	800 M words	5 M words	300 K words	170 K words	50 K words

costo computazionale necessario alla risoluzione  
di un problema differenziale 3D  
su uno stesso calcolatore

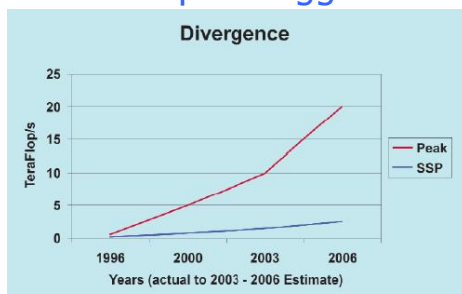
BLAS

A. Murli - Calcolo Scientifico

4

## MA .....

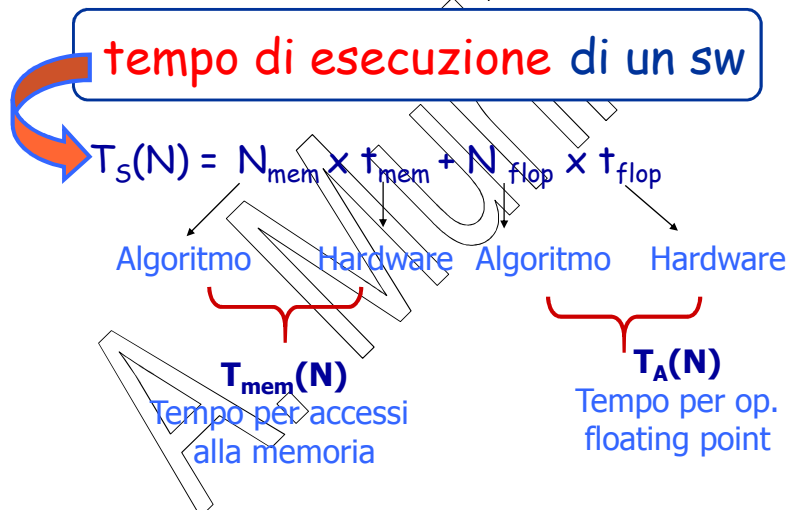
Sono sempre maggiori le difficoltà nell'ottenere le massime prestazioni



**Obiettivo:** ridurre il gap tra la *peak* performance (la massima prestazione del calcolatore) e la *sustained* performance (prestazione massima del software su quel calcolatore) (SSP)

5

## Nel dettaglio...



## Quando si raggiunge la massima prestazione ?

$$T_S(N) = T_A(N) \Rightarrow \frac{T_S(N)}{T_A(N)} = 1$$

OVVERO quando il tempo di esecuzione del software  $T_S(N)$   
 è uguale al tempo  $T_A(N)$  dell'algoritmo

## MA in generale.....

$$\frac{T_S(N)}{T_A(N)} = \frac{N_{\text{mem}} \times t_{\text{mem}} + N_{\text{flop}} \times t_{\text{flop}}}{N_{\text{flop}} \times t_{\text{flop}}} = 1 + \frac{N_{\text{mem}} \times t_{\text{mem}}}{N_{\text{flop}} \times t_{\text{flop}}}$$

avviene che

$$T_S(N) > T_A(N)$$

per migliorare l'efficienza dobbiamo  
 rendere "piccolo" il rapporto

Per migliorare l'efficienza...

$$\frac{T_S(N)}{T_A(N)} = 1 + \frac{N_{\text{mem}} \times t_{\text{mem}}}{N_{\text{flop}} \times t_{\text{flop}}}$$

### Problema 1

Ridurre il **numero** di accessi alla memoria

### Problema 2

Ridurre il **tempo** di accesso alla memoria

BLAS

A. Murli - Calcolo Scientifico

9

## Problema 1

$$\frac{T_S(N)}{T_A(N)} = 1 + \frac{N_{\text{mem}} \times t_{\text{mem}}}{N_{\text{flop}} \times t_{\text{flop}}}$$

$$q = \frac{N_{\text{mem}}}{N_{\text{flop}}} = \frac{\text{\# accessi in memoria}}{\text{\# operazioni floating - point}}$$

Parametro di valutazione del traffico "parassita"

E' possibile ridurre il **traffico parassita** q ?

BLAS

A. Murli - Calcolo Scientifico

10

Quanto vale il parametro  $q$   
(misura del traffico parassita)  
per le operazioni di base di algebra lineare

?

Aggiornamento di un vettore

Prodotto scalare

Prodotto Matrice vettore

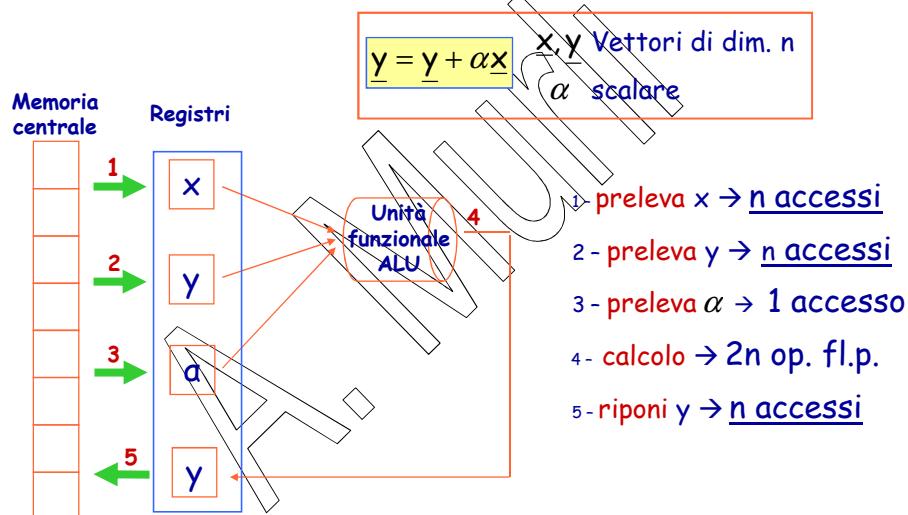
Prodotto Matrice -Matrice

BLAS

A. Murli - Calcolo Scientifico

11

Esempio1: esecuzione di una aggiornamento di un vettore (saxpy)



BLAS

A. Murli - Calcolo Scientifico

12

## Calcoliamo q....

SAXPY



$$q = \frac{m}{f} = \frac{3n+1}{2n} \approx \frac{3}{2}$$

- 1 - preleva  $x \rightarrow n$  accessi
- 2 - preleva  $y \rightarrow n$  accessi
- 3 - preleva  $\alpha \rightarrow 1$  accesso
- 4 - calcolo  $\rightarrow 2n$  op. fl.p.
- 5 - riponi  $y \rightarrow n$  accessi

# accessi > # operazioni f.p.

BLAS

A. Murli - Calcolo Scientifico

13

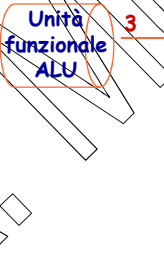
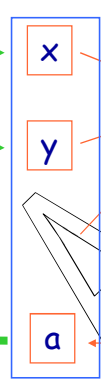
## Esempio2: esecuzione di un prodotto scalare (dot)

$$a = x^T y$$

$x, y$  Vettori di dim.  $n$   
 $\alpha$  scalare

Memoria centrale

Registri



- 1 - preleva  $x \rightarrow n$  accessi
- 2 - preleva  $y \rightarrow n$  accessi
- 3 - calcolo  $\rightarrow 2n$  op. fl.p.
- 4 riponi  $\alpha \rightarrow 1$  accesso

BLAS

A. Murli - Calcolo Scientifico

14

## Calcoliamo q....

DOT



$$q = \frac{m}{f} = \frac{2n+1}{2n} \approx 1$$

1 - preleva  $x \rightarrow n$  accessi

2 - preleva  $y \rightarrow n$  accessi

3 - calcolo  $\rightarrow 2n$  op. fl.p.

4 riponi  $\alpha \rightarrow 1$  accesso

# accessi = # operazioni f.p.

BLAS

A. Murli - Calcolo Scientifico

15

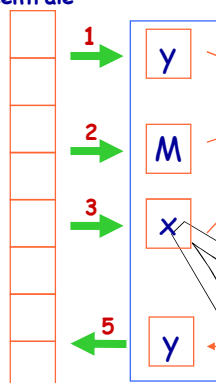
## Esempio3: esecuzione di un prodotto matrice vettore (gemmv)

$$\underline{y} = \underline{y} + M \underline{x}$$

$\underline{x}, \underline{y}$  Vettori di dim.  $n$   
 $M$  Matrice di dim.  $n$

Memoria centrale

Registri



Unità funzionale  
ALU

- ♦ Preleva dalla memoria  $\Rightarrow n$  accessi
- ♦ Preleva dalla memoria  $\Rightarrow n^2$  accessi
- ♦ Preleva dalla memoria  $\Rightarrow n$  accessi
- ♦ Esegui  $\Rightarrow 2n^2$  flop
- ♦ Riponi in memoria  $\Rightarrow n$  accessi

BLAS

A. Murli - Calcolo Scientifico

16



## Calcoliamo q....

GEMMV

$$q_{\text{gemv}} = \frac{3n + n^2}{2n^2} = \frac{1}{2} + \frac{3}{2n}$$

- ♦ Preleva dalla memoria  $\Rightarrow n$  accessi
- ♦ Preleva dalla memoria  $\Rightarrow n^2$  accessi
- ♦ Preleva dalla memoria  $\Rightarrow n$  accessi
- ♦ Esegui  $\Rightarrow 2n^2$  flop
- ♦ Riponi in memoria  $\Rightarrow n$  accessi

# accessi > # operazioni f.p.

BLAS

A. Murli - Calcolo Scientifico

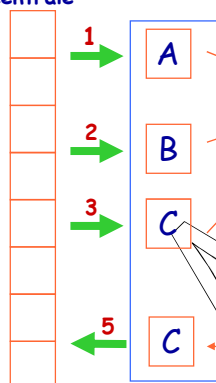
17

## Esempio 4: esecuzione di un prodotto matrice matrice (gemm)

$$C = C + AB \quad \begin{array}{l} A, B \text{ Matrici di dim. } n \\ C \text{ Matrice di dim. } n \end{array}$$

Memoria centrale

Registri



Unità funzionale ALU

- ♦ Preleva dalla memoria  $\Rightarrow n^2$  accessi
- ♦ Preleva dalla memoria  $\Rightarrow n^2$  accessi
- ♦ Preleva dalla memoria  $\Rightarrow n^2$  accessi
- ♦ Esegui  $\Rightarrow 2n^3$  flop
- ♦ Riponi in memoria  $\Rightarrow n^2$  accessi

BLAS

A. Murli - Calcolo Scientifico

18

## Calcoliamo q....

GEMM



$$q_{\text{gemm}} = \frac{3n^2 + (n)^2}{2(n)^3} = \frac{2}{n}$$

- ♦ Preleva dalla memoria  $\Rightarrow n^2$  accessi
- ♦ Preleva dalla memoria  $\Rightarrow n^2$  accessi
- ♦ Preleva dalla memoria  $\Rightarrow n^2$  accessi
- ♦ Esegui  $\Rightarrow 2n^3$  flop
- ♦ Riponi in memoria  $\Rightarrow n^2$  accessi

# accessi < # operazioni f.p.

BLAS

A. Murli - Calcolo Scientifico

19

## In conclusione....

$$q_{\text{saxpy}} \approx \frac{3}{2}$$

# accessi > # operazioni f.p.

$$q_{\text{dot}} \approx 1$$

# accessi = # operazioni f.p.

$$q_{\text{gemv}} \approx \frac{1}{2}$$

# accessi < # operazioni f.p.

$$q_{\text{gemm}} \approx \frac{2}{n}$$

# accessi << # operazioni f.p.

Il traffico parassita  $q$  è  
una costante che  
**DIPENDE** dal nucleo  
computazionale di base

(un prodotto tra matrici  
e' piu' conveniente in termini  
di traffico parassita )!!

A. Murli - Calcolo Scientifico

20

In sintesi per migliorare l'efficienza...

$$\frac{T_S(N)}{T_A(N)} = 1 + \frac{N_{\text{mem}} \times t_{\text{mem}}}{N_{\text{flop}} \times t_{\text{flop}}} = 1 + q \frac{t_{\text{mem}}}{t_{\text{flop}}}$$

Il **numero** di accessi alla memoria  $q$  è costante e dipende dal problema

**BISOGNA**  
Ridurre il **tempo** di accesso alla memoria

E' possibile ridurre il tempo di accesso alla "memoria" ?

BLAS

A. Murli - Calcolo Scientifico

21

Risposta

$$\frac{T_S(N)}{T_A(N)} = 1 + q \frac{t_{\text{mem}}}{t_{\text{flop}}}$$

**numero** di accessi alla memoria  $q$  costante

Ridurre il **tempo** di accesso alla memoria

*Principio di Località dei dati*

Utilizzando la memoria "gerarchica"

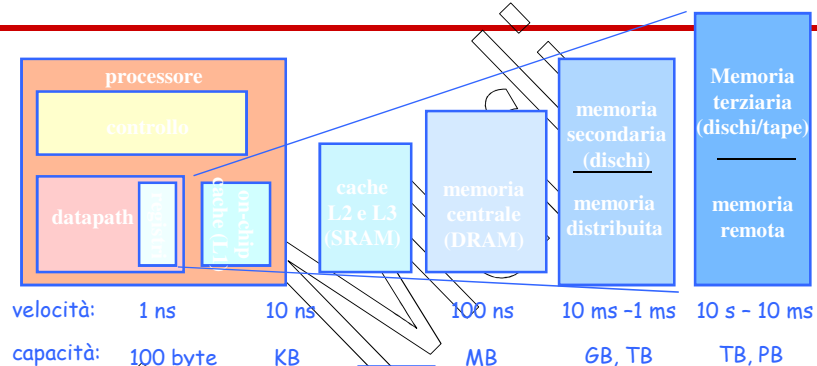
Rivisitazione degli algoritmi al fine di...

BLAS

A. Murli - Calcolo Scientifico

22

## ...Utilizzare la localita' dei dati ...



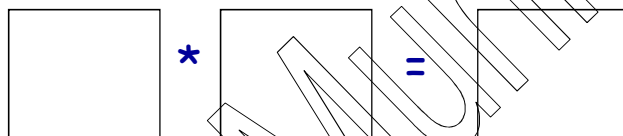
**OVVERO... progettare sw che  
riusano i dati nei  
"Livelli alti della memoria"**

BLAS

23

## Un caso studio il prodotto di due matrici:

$$A * B = C$$



```

for ...
  for ...
    for ...
       $c(i,j) = c(i,j) + a(i,k) * b(k,j)$ 
    endfor
  endfor
endfor

```

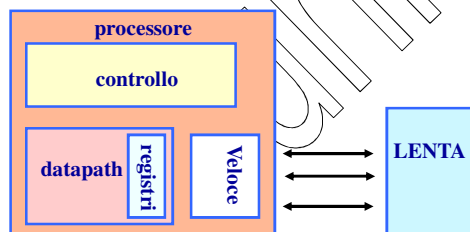
BLAS

A. Murli - Calcolo Scientifico

24

## Alcune Ipotesi:

1. Il Calcolatore in considerazione ha soli 2 livelli nella gerarchia di memoria ("veloce" e "lenta")



2. Tutti i dati all'inizio sono presenti nella memoria lenta

3. Nel modello di prestazione si assume che:

$$\frac{T_S(N)}{T_A(N)} = 1 + q \frac{t_{\text{mem}}}{t_{\text{flop}}} \rightarrow 10^{-1} \text{ sec}$$

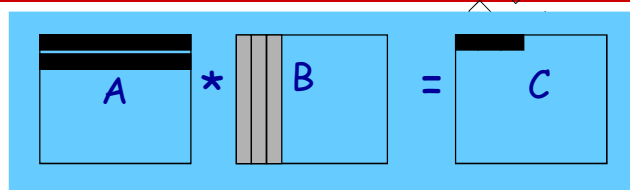
BLAS

A. Murli - Calcolo Scientifico

25

Esempio: A,B,C matrici quadrate di dimensione 10

### Strategia I



Memoria lenta

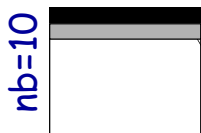
♦ Il primo elemento di C e' ottenuto da un prodotto scalare (saxpy)

(trasferimento memoria Lenta-Veloce 1 riga di A ed 1 di B)

nb=10

♦ Il secondo elemento di C e' ottenuto da un prodotto scalare (saxpy)

(trasferimento memoria Lenta-Veloce 1 riga di A ed 1 di B)



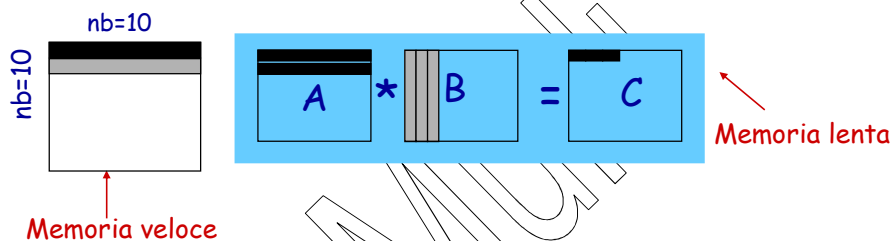
Memoria veloce

BLAS

A. Murli - Calcolo Scientifico

26

## Strategia I (cont)



La prima riga di C è ottenuta mediante 10 saxpy

(trasferimento memoria Lenta-Veloce 10 riga di A ed 10 di B)

Gli elementi di C sono ottenuti mediante 100 saxpy

(trasferimento memoria Lenta-Veloce 10 riga di A ed 10 di B 10 volte!)

BLAS

A. Murli - Calcolo Scientifico

27

## Domanda

### Strategia I:

Ogni elemento di C ottenuto facendo un  
prodotto scalare (saxpy)...

... quanto costa tale strategia in termini di  
**ACCESSI alla memoria ?**

BLAS

A. Murli - Calcolo Scientifico

28

## Risposta

Consideriamo gli accessi alla memoria

$$\frac{T_S(N)}{T_A(N)} = 1 + q \frac{t_{\text{mem}}}{t_{\text{flop}}}$$

Sono necessarie 100 saxpy per ottenere C dunque...

$$100 \times q_{\text{saxpy}} \times \frac{t_{\text{mem}}}{t_{\text{flop}}} = 100 \times 1 \times 10^{-1} = 10 \text{ sec}$$

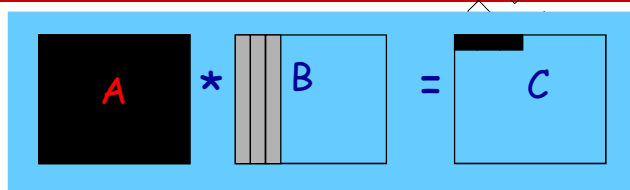
La strategia I richiede 10 sec solo per gli accessi alla memoria!

BLAS

29

Esempio: A,B,C matrici quadrate di dimensione 10

### Strategia II



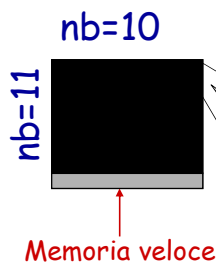
Memoria lenta

♦ La prima riga di C è ottenuta mediante un prodotto matrice vettore (GEMMV)

(trasferimento memoria Lenta-Veloce di 1 sola riga di B e tutta A)

♦ La seconda riga di C è ottenuta mediante un prodotto matrice vettore (GEMMV)

(trasferimento memoria Lenta-Veloce di 1 sola riga di B)



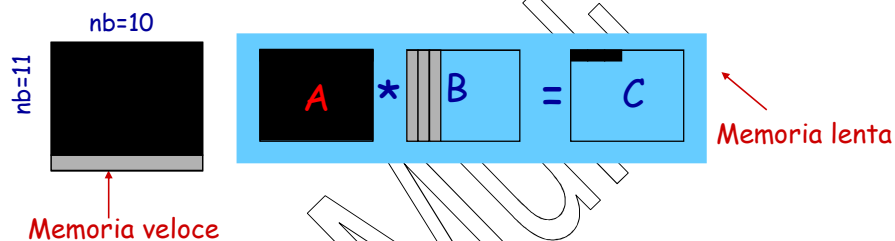
Memoria veloce

BLAS

A. Murli - Calcolo Scientifico

30

## Strategia II (cont)



Le righe di  $C$  sono ottenute mediante 1 GEMMV

(trasferimento memoria Lenta-Veloce 1 elemento di  $B$  ad ogni passo e di tutta  $A$  una volta sola)

Gli elementi di  $C$  sono ottenuti mediante 10 GEMMV

BLAS

A. Murli - Calcolo Scientifico

31

## Domanda

### Strategia II:

Ogni riga di  $C$  è ottenuta mediante un prodotto matrice-vettore (GEMMV)...  
... Ma quanto costa tale strategia in termini di

**ACCESSI alla memoria ?**

BLAS

A. Murli - Calcolo Scientifico

32



## Risposta

Consideriamo gli accessi alla memoria

$$\frac{T_S(N)}{T_A(N)} = 1 + q \frac{t_{\text{mem}}}{t_{\text{flop}}}$$

Sono necessarie 10 GEMMV per ottenere C dunque...

$$10 \times q_{\text{gemmv}} \times \frac{t_{\text{mem}}}{t_{\text{flop}}} = 10 \times \frac{1}{2} \times 10^{-1} = 0.5 \text{ sec}$$

$\approx 1/2$        $10^{-1}$

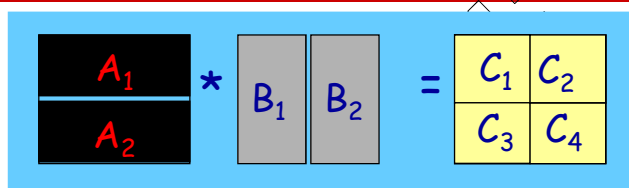
La strategia II richiede 0.5 sec solo per gli accessi alla memoria!

BLAS

33

Esempio: A,B,C matrici quadrate di dimensione 10

### Strategia III



Memoria lenta

nb=10



Memoria veloce

♦  $C_1$  il primo blocco di C e' ottenuto mediante un prodotto matrice matrice (GEMM)

( trasferimento memoria Lenta-Veloce di 1 blocco di A e di B)

♦  $C_2$  il secondo blocco di C e' ottenuto mediante un prodotto matrice matrice (GEMM)

( trasferimento memoria Lenta-Veloce di 1 solo blocco di B)

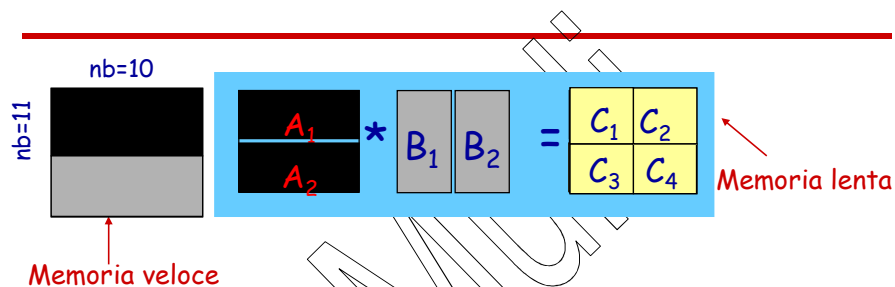
♦ ...

BLAS

A. Murli - Calcolo Scientifico

34

### Strategia III (cont)



Ogni blocco di  $C$  è ottenuto mediante 1 GEMM

Gli elementi di  $C$  sono ottenuti calcolando 4 blocchi ovvero mediante 4 GEMM

BLAS

A. Murli - Calcolo Scientifico

35

### Domanda

#### Strategia III:

Ogni blocco di  $C$  è ottenuto mediante un prodotto matrice-matrice (GEMM)...

... Ma quanto costa tale strategia in termini di

**ACCESSI alla memoria ?**

BLAS

A. Murli - Calcolo Scientifico

36

## Risposta

Consideriamo gli accessi alla memoria

$$\frac{T_S(N)}{T_A(N)} = 1 + q \frac{t_{\text{mem}}}{t_{\text{flop}}}$$

Sono necessarie 4 GEMM per ottenere C dunque...

$$4 \times q_{\text{gemm}} \times \frac{t_{\text{mem}}}{t_{\text{flop}}} = 4 \times \frac{1}{10} \times 10^{-1} = 4 \times 10^{-2} \text{ sec}$$

$\approx 1/10$        $10^{-1}$

La strategia III richiede  $4 \times 10^{-2}$  sec  
per gli accessi alla memoria!

BLAS

37

## In conclusione

Strategia I (saxpy)

richiede 10 sec  
per gli accessi alla memoria!

Strategia II (GEMMV)

richiede 1 sec  
per gli accessi alla memoria!

Strategia III (GEMM)

richiede  $4 \times 10^{-2}$  sec  
per gli accessi alla memoria!

È Conveniente organizzare gli algoritmi con  
operazioni a blocchi  
(nuclei computazionali Matrice - Matrice )

BLAS

A. Manti - Calcolo Scientifico

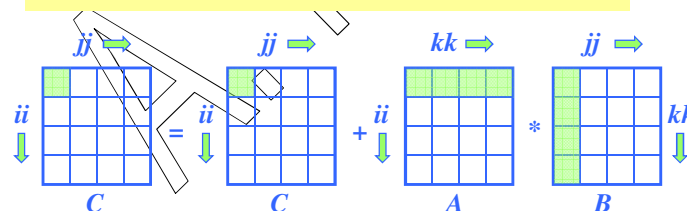
38

## In conclusione

Algoritmi con operazioni a blocchi

Riuso ottimale dei dati nei

"Livelli alti della memoria"



BLAS

A. Murli - Calcolo Scientifico

39

## Quanto descritto è alla base di...

**BLAS:**

Basic Linear Algebra Subroutines:

Libreria di software matematico per  
l'esecuzione di operazioni di base del  
calcolo matriciale che  
ottimizza gli accessi alla memoria

BLAS

A. Murli - Calcolo Scientifico

40

## BLAS

### ◆ BLAS 1:

- Operazioni di base tra vettori

- » Somma
- » Aggiornamento
- » Norma
- » Prodotto scalare
- » .....



### ◆ BLAS 2

- Operazioni di base tra matrice e vettori

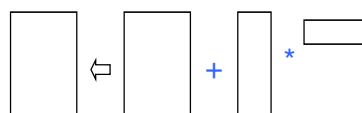
- Prodotto matrice - vettore
- Aggiornamento
- .....



### ◆ BLAS 3

- Operazioni di base tra matrici

- » Prodotto tra matrici
- » Aggiornamento
- » Norma
- » Somma
- » .....



BLAS

A. Murli - Calcolo Scientifico

41

## Confronto BLAS1, BLAS2, BLAS3

### ◆ BLAS 1:

- Ottimizza le operazioni tra vettori (loop unrolling)

### ◆ BLAS 2

- Ottimizza il riutilizzo dei dati che risiedono nei registri (riduce lo spostamento dei dati dalla cache ai registri)

### ◆ BLAS 3

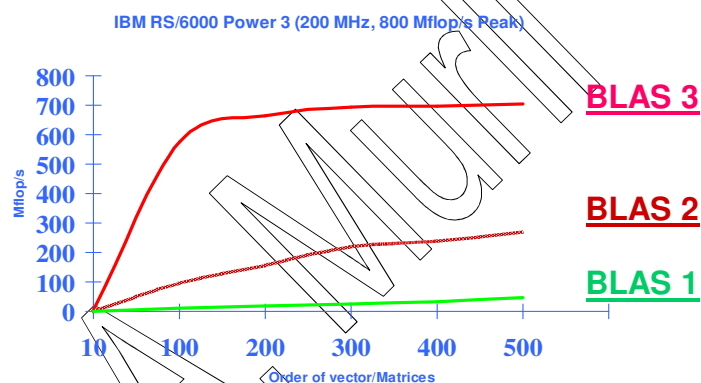
- Ottimizza il riutilizzo dei dati che risiedono nella cache (riduce lo spostamento dei dati dalla memoria centrale alla cache)

BLAS

A. Murli - Calcolo Scientifico

42

## Un confronto tra BLAS1, BLAS2, BLAS3



Solo BLAS 3 riesce a raggiungere  
la peak performance del processore IBM /RS 6000  
perché ottimizza l'uso di **tutti i livelli** della memoria gerarchica

43

“As machines become more powerful, the efficiency of algorithms grows more important, not less.”

Nick Trefethen, 1997

Quanto più potente è l'ambiente di calcolo tanto più diventa complesso e difficile lo sviluppo di software in grado di sfruttarne appieno le capacità.

---

**FINE LEZIONE**

---

## **Esercitazione**

Come sviluppare software efficiente per il  
calcolo matriciale ?

## Case study n. 1: Prodotto matrice x matrice

Calcolo di

$$A \cdot B = C$$

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

```

for _ = 1:n;
  for _ = 1:n;
    for _ = 1:n;
      Ci,j ← Ci,j + Ai,k Bk,j
    end
  end
end

```

Indipendentemente dall' ordine dei cicli  
abbiamo sempre  $O(n^3)$  operazioni f.p.

BLAS

A. Murli - Calcolo Scientifico

47

## 6 Varianti della moltiplicazione di matrici

```

for _ = 1:n;
  for _ = 1:n;
    for _ = 1:n;
      Ci,j ← Ci,j + Ai,k Bk,j
    end
  end
end

```

BLAS

A. Murli - Calcolo Scientifico

48

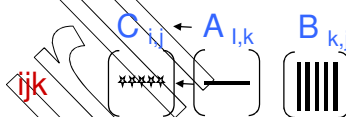


## 6 Varianti della moltiplicazione di matrici

```

for  $\underline{i}$  = 1:n;
  for  $\underline{j}$  = 1:n;
    for  $\underline{k}$  = 1:n;
       $C_{i,j} \leftarrow C_{i,j} + A_{i,k} B_{k,j}$ 
    end
  end
end

```



BLAS

A. Murli - Calcolo Scientifico

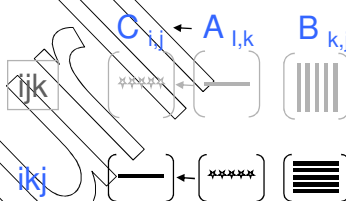
49

## 6 Varianti della moltiplicazione di matrici

```

for  $\underline{i}$  = 1:n;
  for  $\underline{k}$  = 1:n;
    for  $\underline{j}$  = 1:n;
       $C_{i,j} \leftarrow C_{i,j} + A_{i,k} B_{k,j}$ 
    end
  end
end

```



BLAS

A. Murli - Calcolo Scientifico

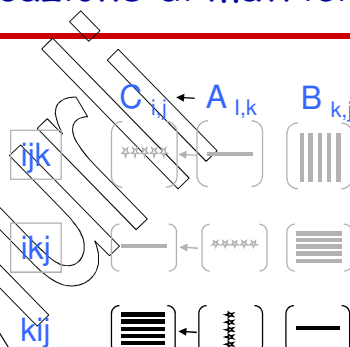
50

## 6 Varianti della moltiplicazione di matrici

```

for  $\underline{k} = 1:n$ ;
  for  $\underline{i} = 1:n$ ;
    for  $\underline{j} = 1:n$ ;
       $C_{i,j} \leftarrow C_{i,j} + A_{i,k} B_{k,j}$ 
    end
  end
end

```



BLAS

A. Murli - Calcolo Scientifico

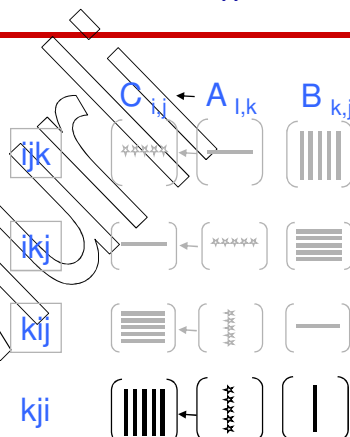
51

## 6 Varianti della moltiplicazione di matrici

```

for  $\underline{k} = 1:n$ ;
  for  $\underline{j} = 1:n$ ;
    for  $\underline{i} = 1:n$ ;
       $C_{i,j} \leftarrow C_{i,j} + A_{i,k} B_{k,j}$ 
    end
  end
end

```



BLAS

A. Murli - Calcolo Scientifico

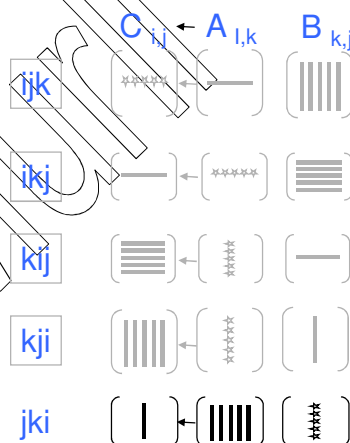
52

## 6 Varianti della moltiplicazione di matrici

```

for j = 1:n;
  for k = 1:n;
    for i = 1:n;
      Ci,j ← Ci,j + Ai,k Bk,j
    end
  end
end

```



BLAS

A. Murli - Calcolo Scientifico

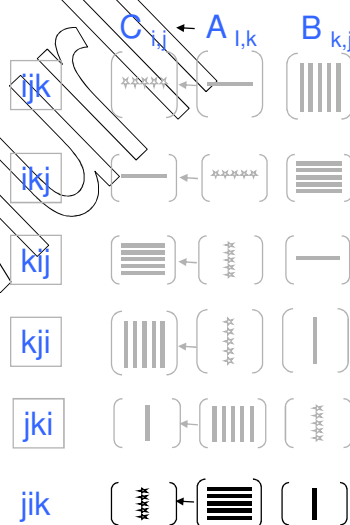
53

## 6 Varianti della moltiplicazione di matrici

```

for j = 1:n;
  for i = 1:n;
    for k = 1:n;
      Ci,j ← Ci,j + Ai,k Bk,j
    end
  end
end

```



BLAS

A. Murli - Calcolo Scientifico

54

Quale versione implementare?

BLAS

A. Murli - Calcolo Scientifico

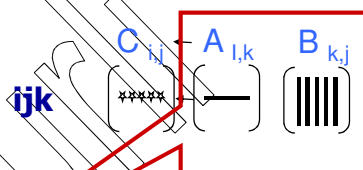
55

## 6 Varianti della moltiplicazione di matrici

```

for  $\underline{i}$  = 1:n;
  for  $\underline{j}$  = 1:n;
    for  $\underline{k}$  = 1:n;
       $C_{i,j} \leftarrow C_{i,j} + A_{i,k} B_{k,j}$ 
    end
  end
end

```



Operazione di base:  
**prodotto scalare**  
 (BLAS1)

della riga  $i$ -ma di  $A$  e  
 della colonna  $j$ -ma di  $B$

BLAS

A. Murli - Calcolo Scientifico

56

## 6 Varianti della moltiplicazione di matrici

for **i** = 1:n;

for **k** = 1:n;

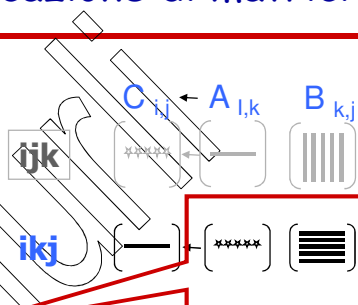
for **j** = 1:n;

$C_{i,j} \leftarrow C_{i,j} + A_{i,k} B_{k,j}$

end

end

end



Operazione di base:

Matrice B x vettore(riga)  $A_{i,k}$   
(aggiornamento per riga)

(BLAS 2)

BLAS

A. Murli - Calcolo Scientifico

57

## 6 Varianti della moltiplicazione di matrici

for **k** = 1:n;

for **i** = 1:n;

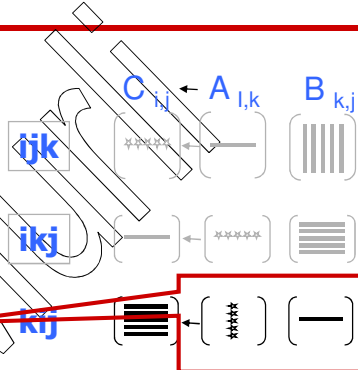
for **j** = 1:n;

$C_{i,j} \leftarrow C_{i,j} + A_{i,k} B_{k,j}$

end

end

end



Operazione di base: saxpy

della riga k -ma di B  
con lo scalare  $A_{i,k}$

(BLAS 1)

BLAS

A. Murli - Calcolo Scientifico

58

## 6 Varianti della moltiplicazione di matrici

```
for  $\underline{k}$  = 1:n;
```

```
  for  $\underline{j}$  = 1:n;
```

```
    for  $\underline{i}$  = 1:n;
```

```
       $C_{i,j} \leftarrow C_{i,j} + A_{i,k} B_{k,j}$ 
```

```
    end
```

```
  end
```

```
end
```

BLAS

A. Murli - Calcolo Scientifico

(BLAS 1)

59

Operazione di base: saxpy  
della colonna  $k$ -ma di  $A$   
con lo scalare  $B_{kj}$

## 6 Varianti della moltiplicazione di matrici

```
for  $\underline{j}$  = 1:n;
```

```
  for  $\underline{k}$  = 1:n;
```

```
    for  $\underline{i}$  = 1:n;
```

```
       $C_{i,j} \leftarrow C_{i,j} + A_{i,k} B_{k,j}$ 
```

```
    end
```

```
  end
```

```
end
```

BLAS

(BLAS 2)

A. Murli - Calcolo Scientifico

60

Operazione di base:  
Matrice  $\times$  vettore  
(aggiornamento per colonna) $_j$

## 6 Varianti della moltiplicazione di matrici

for  $j = 1:n$ ;

for  $i = 1:n$ ;

for  $k = 1:n$ ;

$$C_{i,j} \leftarrow C_{i,j} + A_{i,k} B_{k,j}$$

end

end

end

BLAS

Operazione di base:

prodotto scalare

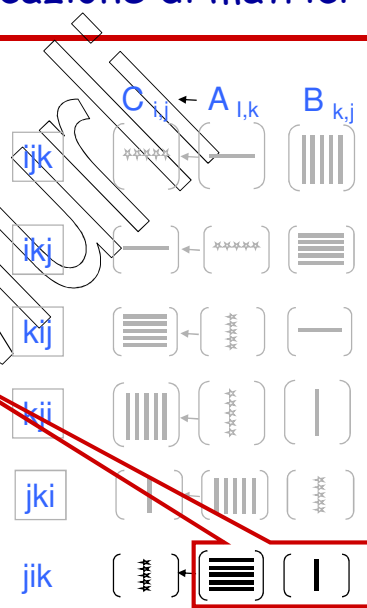
(BLAS1)

della riga  $i$ -ma di  $A$  e

della colonna  $j$ -ma di  $B$

A. Murli - Calcolo Scientifico

61



Ciascuna permutazione degli indici implica  
l'esecuzione di una diversa operazione di base  
tra vettori (BLAS 1) o  
tra vettori e matrici (BLAS 2)

BLAS

A. Murli - Calcolo Scientifico

62

## Case study 2

Risoluzione di

$$Ax = b$$

$$A \in \mathbb{R}^{n \times n}, \quad x, b \in \mathbb{R}^n$$

Mediante

L' algoritmo di  
eliminazione di Gauss

```

for k = 1 to n-1
  for i = k+1 to n
    a(i,k) = a(i,k)/a(k,k)
    for j = k+1 to n
      a(i,j) = a(i,j) - a(i,k) * a(k,j)
    endfor
  endfor
endfor

```

Indipendentemente dall' ordine dei cicli  
abbiamo sempre  $O(n^3)$  operazioni f.p.

BLAS

A. Murli - Calcolo Scientifico

63

## Versione k i j :

Operazione di base?

$$\underline{a}^i = \underline{a}^i - a_{ik} \cdot \underline{a}^k$$

kij



```

for k = 1 to n-1
  for i = k+1 to n
    a(i,k) = a(i,k)/a(k,k)
    for j = k+1 to n
      a(i,j) = a(i,j) - a(i,k) * a(k,j)
    endfor
  endfor
endfor

```

Aggiornamento di un vettore (riga)  
mediante il prodotto di uno scalare  
per un vettore ( $O(n)$  op. f.p.)

BLAS

A. Murli - Calcolo Scientifico

64



## Invertendo i cicli .....Versione (k j i):

Operazione di base?

$$\underline{a}^j = \underline{a}^j - a_{kj} \cdot \underline{a}^k$$

kji



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

```
  for j = k+1 to n
```

```
    for i = k+1 to n
```

```
      a(i,k) = a(i,k)/a(k,k)
```

```
      a(i,j) = a(i,j) - a(i,k) * a(k,j)
```

```
    endfor
```

```
  endfor
```

```
endfor
```

Aggiornamento di un vettore (colonna)  
mediante il prodotto di uno scalare  
per un vettore ( $O(n)$  op. f.p.)

BLAS

A. Murli - Calcolo Scientifico

65

## Invertendo i cicli ....Versione (j k i):

Operazioni di base?

$$\underline{a}^j = \underline{a}^j - M \cdot \underline{a}^j$$

jki



```
for j = ... to ...
```

```
  for k = ... to ...
```

```
    for i = ... to ...
```

```
      a(i,j) = a(i,j) - a(i,k) * a(k,j)
```

```
    endfor
```

```
  endfor
```

```
endfor
```

Aggiornamento di un vettore mediante  
il prodotto di una matrice per un  
vettore ( $O(n^2)$  op. f.p.)

BLAS

A. Murli - Calcolo Scientifico

66

Quale versione dello **stesso** algoritmo  
conviene implementare?

In generale,  
la scelta dipende dall'ambiente di calcolo ...

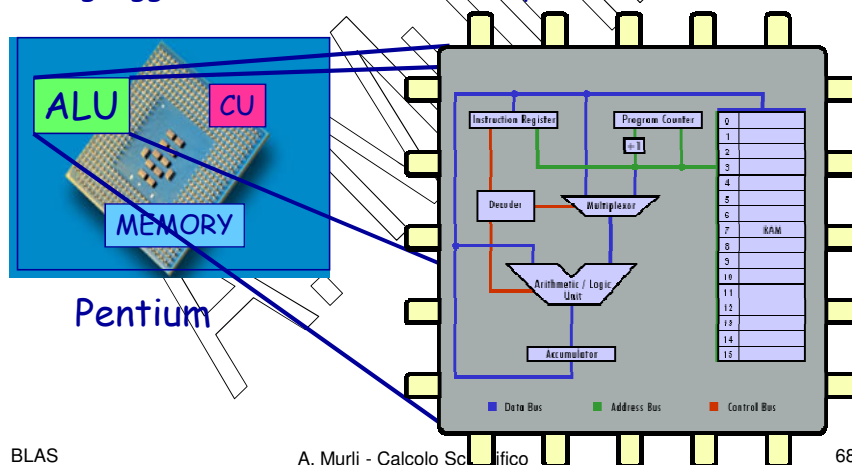
BLAS

A. Murli - Calcolo Scientifico

67

## Ambiente computazionale

- ♦ Architetture monoprocesore
- ♦ Linguaggi ad alto livello (ad esempio Fortran, C)



BLAS

A. Murli - Calcolo Scientifico

68

---

In che modo questo ambiente  
computazionale influenza l'implementazione  
di questi algoritmi?

---

Individuiamo le varianti le cui operazioni di base  
sono  
operazioni tra vettori  
(ad es saxpy, prodotto scalare )

## moltiplicazione di matrici con BLAS 1

```

for _ = 1:n;
  for _ = 1:n;
    for _ = 1:n;
       $C_{i,j} \leftarrow C_{i,j} + A_{i,k} B_{k,j}$ 
    end
  end
end

```

(DOT=prodotto scalare)

BLAS

A. Murli - Calcolo Scientifico

71

## moltiplicazione di matrici

```

for _ = 1:n;
  for _ = 1:n;
    for _ = 1:n;
       $C_{i,j} \leftarrow C_{i,j} + A_{i,k} B_{k,j}$ 
    end
  end
end

```

(saxpy=aggiornamento  
vettore)

BLAS

A. Murli - Calcolo Scientifico

72

## 6 Varianti della moltiplicazione di matrici

```
for _ = 1:n;
```

```
  for _ = 1:n;
```

```
    for _ = 1:n;
```

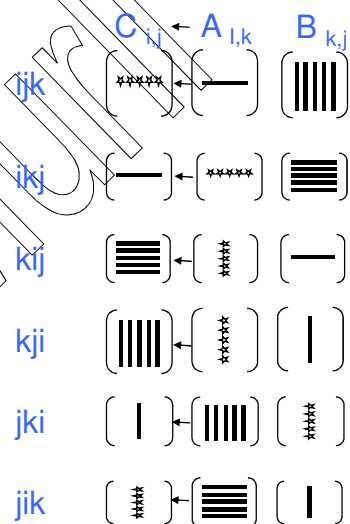
```
       $C_{i,j} \leftarrow C_{i,j} + A_{i,k} B_{k,j}$ 
```

```
    end
```

```
  end
```

```
end
```

**C**  
Fortran  
BLAS



A. Murli - Calcolo Scientifico

73

## 6 Varianti della moltiplicazione di matrici

```
for _ = 1:n;
```

```
  for _ = 1:n;
```

```
    for _ = 1:n;
```

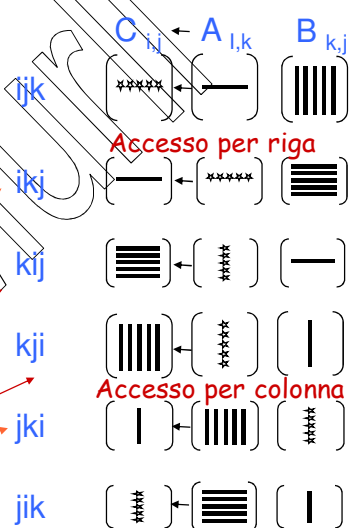
```
       $C_{i,j} \leftarrow C_{i,j} + A_{i,k} B_{k,j}$ 
```

```
    end
```

```
  end
```

```
end
```

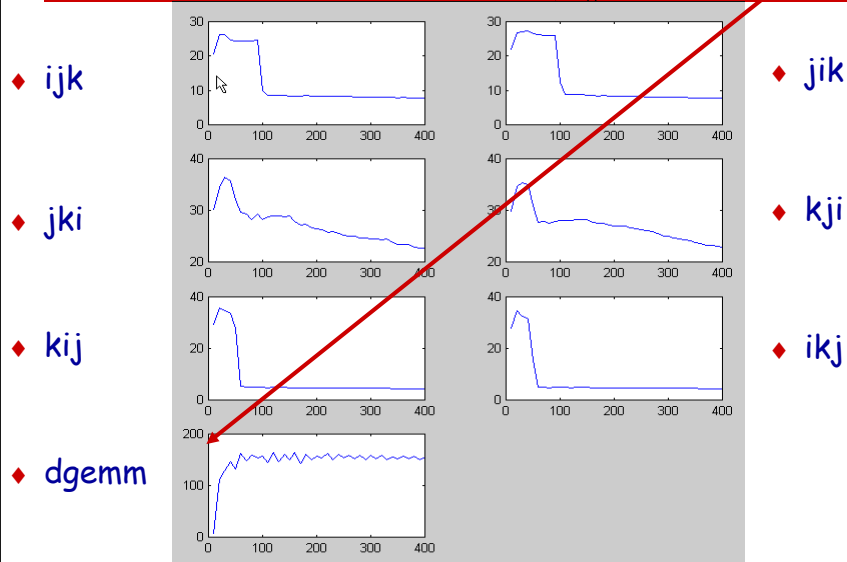
**C**  
Fortran  
BLAS



A. Murli - Calcolo Scientifico

74

## Prestazioni delle 6 varianti di matrice-matrice su un calcolatore SUN Ultra 2 - 200 MHz



BLAS

A. Murli - Calcolo Scientifico

75

## LAPACK

Libreria di software matematico per  
la risoluzione di problemi di algebra  
lineare su architetture a memoria  
gerarchica

Utilizza BLAS 3 come building block

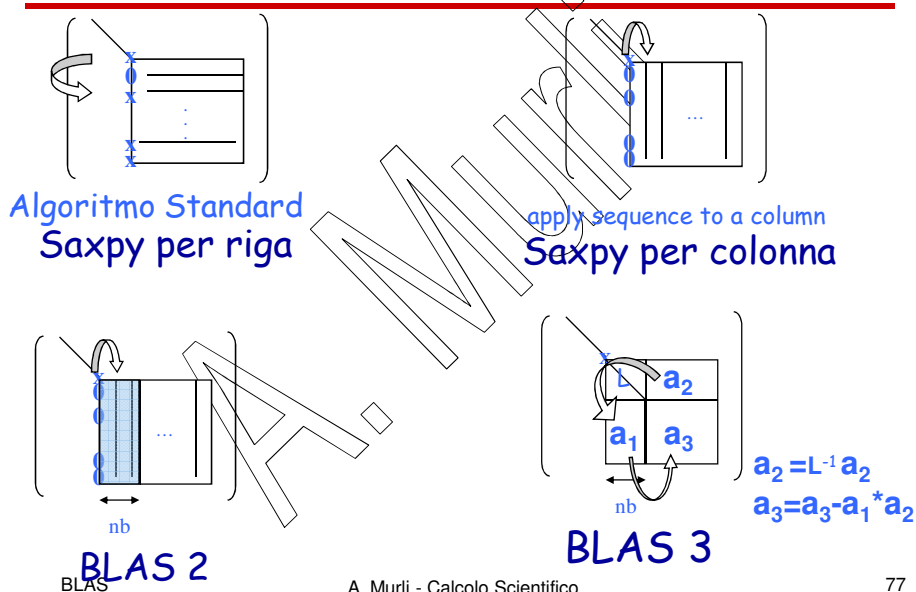
E. Anderson et al., LAPACK User' Guide, SIAM, 1995

BLAS

A. Murli - Calcolo Scientifico

76

## Algoritmo di Gauss con BLAS..



## Alcune problematiche ....

- ♦ Molti parametri da gestire nell'implementazione degli algoritmi (dimensione dei blocchi, permutazione di indici dei loop, profondità dell'unrolling, numero di processi, topologia, ...)
- ♦ L'architettura dei microprocessori diventa sempre più complessa

Necessità di approcci automatici e/o adattativi  
nello sviluppo di software

Differenti caratteristiche capacità della  
cache per differenti calcolatori con  
memoria gerarchica



Sviluppo di implementazioni specializzate di  
BLAS 3

## ATLAS

### Automatically Tuned Linear Algebra Software

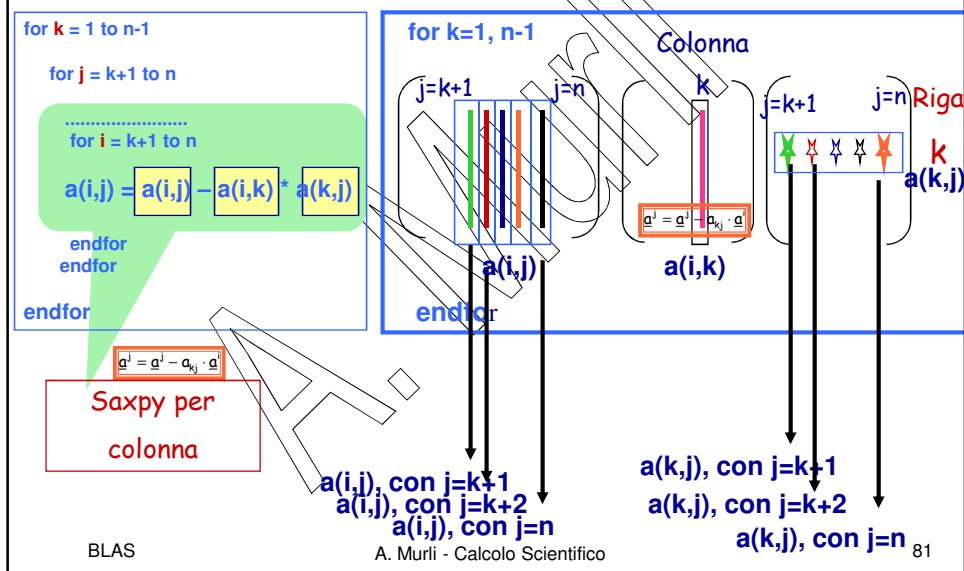
- Versione "self-tuned" di BLAS e di alcune routine di LAPACK
- Kernel fondamentale: GEMM  

$$C = \alpha \text{ op}(A)\text{ op}(B) + \beta \text{ op}(C), \quad \text{op}(X) = X, X^T,$$
- Routine di BLAS 3 implementate in termini di GEMM

Whaley, J. Dongarra, *Automatically Tuned Linear Algebra Software*, Ninth SIAM Conference on Parallel Processing for Scientific Computing, 1999



## Riprendiamo la versione **kji** dell'algoritmo di Gauss



$$\underline{a}^j = \underline{a}^j - a_{kj} \cdot \underline{a}^i$$

Al variare di  $k$ , ogni colonna  $\underline{a}^j$   
viene aggiornata,  
(cioè prelevata e riposta in memoria)

$j$  volte

Uso inefficiente dei registri vettoriali

Mantenere il vettore  $\underline{a}^j$  nei registri  
finché l'aggiornamento non è  
completato

## OVVERO....

Considerare il ciclo "for ... j" come il più esterno

ciò significa far fare all'algorithm

**definitivamente**

tutti gli aggiornamenti sul vettore j

BLAS

A. Murli - Calcolo Scientifico

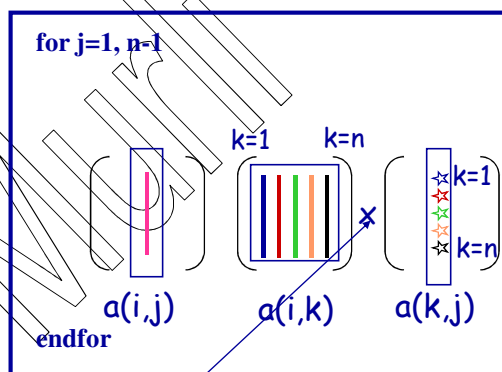
83

## Versione jki: rappresentazione grafica

```

for j = 1 to n
  for k = 1 to n
    .....
    for i = k+1 to n
       $a(i,j) = a(i,j) - a(i,k) * a(k,j)$ 
    endfor
  endfor
endfor

```



Prodotto "colonne x colonne"  
di una matrice per un vettore

**GAXPY**

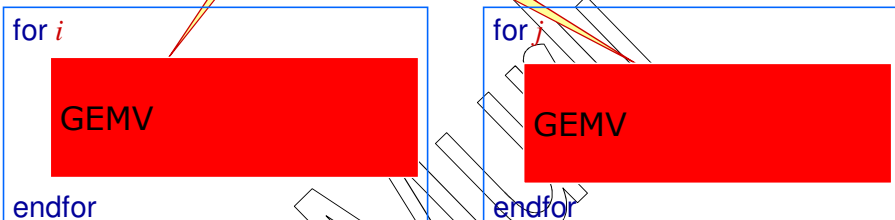
BLAS

A. Murli - Calcolo Scientifico

84



## Versioni **ikj** e **jki** con BLAS 2



Operazione di base:

$$y = y + A x$$

$A$  matrice,  $x, y$  vettori

**GEMV**



$$q = \frac{m}{f} = \frac{n^2 + 3n}{2n^2} \approx \frac{1}{2}$$

BLAS

A. Murli - Calcolo Scientifico

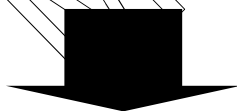
87

Quindi..

**ikj e jki USANO BLAS 2**

La GEMV richiede asintoticamente

**1** accesso alla memoria ogni **2** operazioni fl. point



Gli algoritmi di BLAS 2 sfruttano

la località dei dati che risiedono nella cache

**(riuso dei dati nella cache)**

BLAS

A. Murli - Calcolo Scientifico

88

## Un PROBLEMA .....

La capacità dei registri vettoriali è limitata

Nell'esecuzione della GAXPY se la dimensione del vettore  $y$  e della matrice è maggiore della capacità dei registri, si **prelevano di fatto solo sottovettori di  $y$  e blocchi di  $M$**

"inutile" traffico (parassita) dei dati tra registri e cache

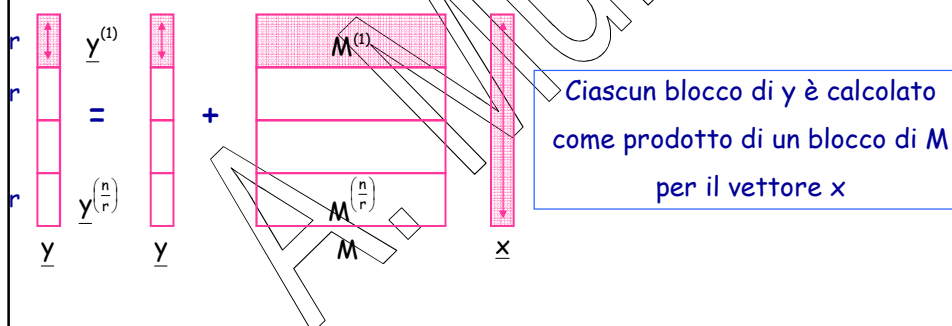
BLAS

A. Murli - Calcolo Scientifico

89

## Strategia

scrivere l'algoritmo in maniera da costringere i blocchi di  $M$  e di  $y$  nei registri **finché** l'aggiornamento su ciascun blocco non sia stato **completato**



BLAS

A. Murli - Calcolo Scientifico

90

## Calcolo di $q$ (Gaxpy)

Per ogni  $i$

- ♦ Preleva un blocco di  $\underline{y}$   $\Rightarrow r$  accessi
- ♦ Preleva un blocco di  $\underline{M}$   $\Rightarrow nr$  accessi
- ♦ Preleva il vettore  $\underline{x}$   $\Rightarrow n$  accessi
- ♦ Esegui  $\underline{y}^{(i)} = \underline{y}^{(i)} + \underline{M}^{(i)} \underline{x} \Rightarrow 2nr$  operazioni fl. p.
- ♦ Riponi il blocco di  $\underline{y}$   $\Rightarrow r$  accessi

$$\tilde{q}_{\text{gaxpy}} = \frac{2r + nr + n}{2nr} \approx \frac{1}{2} + \frac{1}{2r} + \frac{1}{n} \Rightarrow \tilde{q}_{\text{gaxpy}} = \frac{1}{2} + \frac{3}{2n} = q_{\text{gaxpy}} \quad (n = r \text{ caso migliore})$$

BLAS

A. Murli - Calcolo Scientifico

91

## Riprendiamo la versione jki sui blocchi di $y$

```
for j
  for k
    for i
      a(i,j) = a(i,j) - a(i,k) * a(k,j)
    endfor
  endfor
```

La matrice attiva  $M$  viene  
prelevata e riposta nella cache  
più volte al variare di  $j$  e  $k$

IDEA:

Mantenere la matrice  $M$  nella cache finché  
l'aggiornamento non è completato

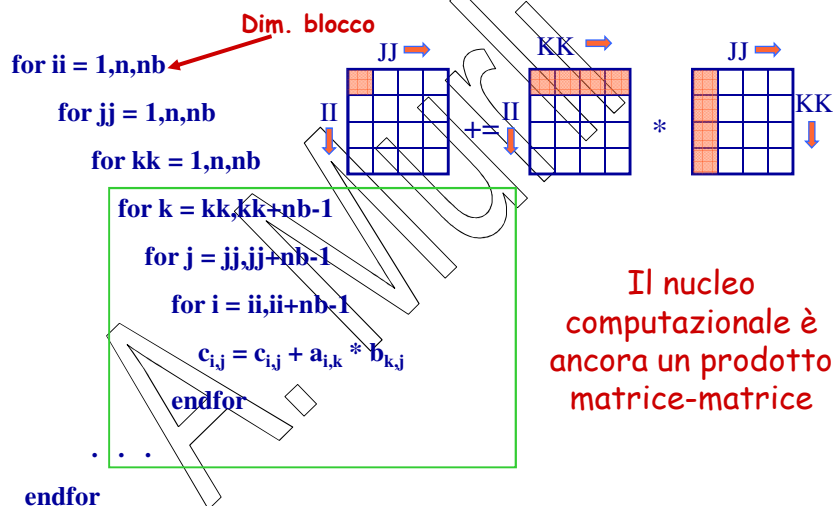
Quale formulazione dell' algoritmo di Gauss consente ciò?

BLAS

A. Murli - Calcolo Scientifico

92

## SGEMM: matrice -matrice a blocchi (BLAS 3)



BLAS

A. Murli - Calcolo Scientifico

93

## Calcolo di q

- ♦ Preleva C  $\Rightarrow n^2$  accessi
- ♦ Preleva A  $\Rightarrow n^2$  accessi
- ♦ Preleva B  $\Rightarrow n^2$  accessi
- ♦ Esegui  $C=C+AB$   $\Rightarrow 2n^3$  operazioni fl. p.
- ♦ Riponi C  $\Rightarrow n^2$  accessi

$$q_{sgemm} = \frac{4n^2}{2n^3} = \frac{2}{n} < \frac{3}{2n} + \frac{1}{2} \approx q_{gaxpy}$$

Ulteriore riduzione del traffico parassita

BLAS

A. Murli - Calcolo Scientifico

94

## Calcolo di q

- ♦ Preleva  $C_{ij}$   $\Rightarrow nb^2$  accessi
- ♦ Preleva  $A_{ik}$  ( $k = 1, \dots, p$ )  $\Rightarrow nb^2$  accessi
- ♦ Preleva  $B_{kj}$  ( $k = 1, \dots, p$ )  $\Rightarrow p(nb)^2$  accessi
- ♦ Esegui  $C_{ij} = C_{ij} + A_{ik} B_{kj}$   $\Rightarrow 2p(nb)^3$  operazioni fl. p.
- ♦ Riponi  $C_{ij}$   $\Rightarrow nb^2$  accessi

$$q_{sgemm} = \frac{3nb^2 + p(nb)^2}{2p(nb)^3} = \frac{p+3}{2p(nb)} = \frac{3}{2pnb} + \frac{1}{2nb} = \frac{2}{n} \quad (\text{caso migliore})$$