

# SUMMA

Scalable Universal Matrix Multiplication Algorithm

# Naïve matrix multiply

```
For i = 0 to n
  For j = 0 to n
    For k = 0 to n
      C[i,j] += A[i,k]*B[k,j]
```

Calculates  $n^2$  dot products (inner products)

$C[i,j] = A[i,:]*B[:,j]$

# Naïve matrix multiply

What happens if we switch the order?

```
For k = 0 to n
```

```
  For i = 0 to n
```

```
    For j = 0 to n
```

```
      C[i,j] += A[i,k]*B[k,j]
```



# Naïve matrix multiply

What happens if we switch the order?

For  $k = 0$  to  $n$

For  $i = 0$  to  $n$

For  $j = 0$  to  $n$

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

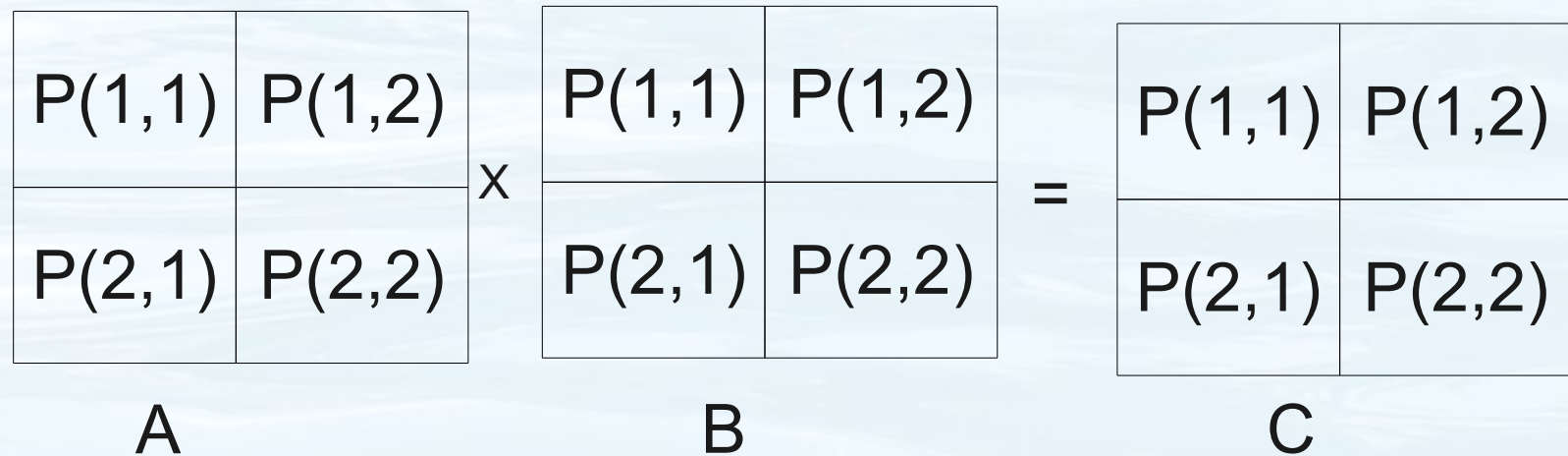
Calculates  $n$  outer products

# Outer product

$$\begin{bmatrix} a \\ b \\ c \end{bmatrix} * [ e \quad f \quad g ] = \begin{bmatrix} a*e & a*f & a*g \\ b*e & b*f & b*g \\ c*e & c*f & c*g \end{bmatrix}$$

# Introducing SUMMA

- Processors arranged in grid,  $P(i,j)$
- Example, 2x2 grid:





# Introducing SUMMA

- Let  $A = m \times k$ ,  $B = k \times n \Rightarrow C = m \times n$
- Let each process do  $k$  outer products

How do we handle the communication?

# SUMMA

- For each  $k$  (between 0 and  $n-1$ ),
  - Owner of partial row  $K$  broadcasts that row along its process column
  - Owner of partial column  $K$  broadcasts that column along its process row

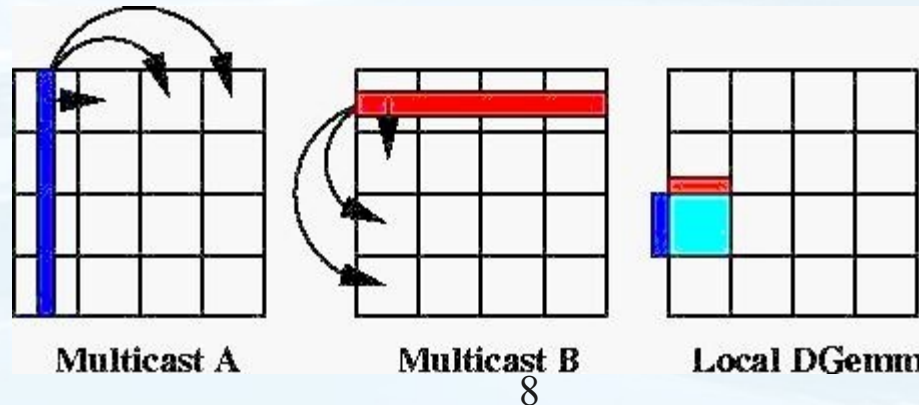


Image credit:  
Stephen J. Fink



# SUMMA

Complete algorithm. On each process  $P(i,j)$ :

For  $k = 0 \dots n-1$

    Bcast column  $k$  of  $A$  ( $a_i$ ) within row  $i$

    Bcast row  $k$  of  $B$  ( $b_j$ ) within column  $j$

    Do  $C +=$  outer product ( $a_i, b_j$ )

# Communication cost

- What's the cost of communication? Let  $\alpha$  be the startup cost of a message, and  $\beta$  be the bandwidth

# Communication cost

- What's the cost of communication? Let  $\alpha$  be the startup cost of a message, and  $\beta$  be the bandwidth
- Bcast among  $p$  processes takes  $\log(p) * (\alpha + \beta * s)$  time, where  $s$  is the size of the message



# Communication cost

- Bcast among  $p$  processes takes  $\log(p)$   $(\alpha + \beta s)$  time, where  $s$  is the size of the message
- For each  $k$ , there are one Bcast along columns and one Bcast along rows
- Each partial column has size  $m/r$ , each partial row has size  $n/c$

# Communication cost

- Bcast among  $p$  processes takes  $\log(p)$   $(\alpha + \beta s)$  time, where  $s$  is the size of the message
- Total:  $k * (\log(c)(\alpha + \beta * m/r) + \log(r)(\alpha + \beta * n/c))$
- Not very efficient
  - Lots of messages

# Improvements?

- Can we reduce the communication cost somehow?



# Improvements?

- Can we reduce the communication cost somehow?
- Obvious improvement: Instead of broadcasting single rows and columns, do block rows and columns.

# Improvements?

- Obvious improvement: Instead of broadcasting single rows and columns, do block rows and columns.
- Same amount of bytes communicated
- Fewer messages => less overhead
- More efficient ( $O(n^3)$  FLOPS,  $O(n^2)$  loads)

# Pipelined SUMMA

- Another improvement:
- Instead of broadcast one row/column segment to all, pass multiple segments around in a ring
- Each process only communicates with neighbor => no broadcast => fewer messages



# Pipelined SUMMA - Example

- Consider the following matrices A and B (color represents process rank)

a1	a2	a3	b1	b2	b3		e1	e2	e3	f1	f2	f3
a4	a5	a6	b4	b5	b6		e4	e5	e6	f4	f5	f6
a7	a8	a9	b7	b8	b9	x	e7	e8	e9	f7	f8	f9
c1	c2	c3	d1	d2	d3		g1	g2	g3	h1	h2	h3
c4	c5	c6	d4	d5	d6		g4	g5	g6	h4	h5	h6
c7	c8	c9	d7	d8	d9		g7	g8	g9	h7	h8	h9

# Pipelined SUMMA - Example

- K=0: All processes multiply their first (block) column/row

★	★		
<u>a1</u> a2 a3	<u>b1</u> b2 b3	<u>e1</u> <u>e2</u> <u>e3</u>	<u>f1</u> <u>f2</u> <u>f3</u> ★
<u>a4</u> a5 a6	<u>b4</u> b5 b6	e4 e5 e6	f4 f5 f6
<u>a7</u> a8 a9	<u>b7</u> b8 b9	e7 e8 e9	f7 f8 f9
<u>c1</u> c2 c3	<u>d1</u> d2 d3	<u>g1</u> <u>g2</u> <u>g3</u>	<u>h1</u> <u>h2</u> <u>h3</u> ★
<u>c4</u> c5 c6	<u>d4</u> d5 d6	g4 g5 g6	h4 h5 h6
<u>c7</u> c8 c9	<u>d7</u> d8 d9	g7 g8 g9	h7 h8 h9

(for reference)

a1 a2 a3	b1 b2 b3	e1 e2 e3	f1 f2 f3
a4 a5 a6	b4 b5 b6	e4 e5 e6	f4 f5 f6
a7 a8 a9	b7 b8 b9	e7 e8 e9	f7 f8 f9
c1 c2 c3	d1 d2 d3	g1 g2 g3	h1 h2 h3
c4 c5 c6	d4 d5 d6	g4 g5 g6	h4 h5 h6
c7 c8 c9	d7 d8 d9	g7 g8 g9	h7 h8 h9

# Pipelined SUMMA - Example

- K=1: All processes send their first (block) col in A to the right, their first (block) row in B down in a ring pattern

★			★					
b1	<u>a2</u>	a3	a1	<u>b2</u>	b3	g1	g2	g3
b4	<u>a5</u>	a6	a4	<u>b5</u>	b6	<u>e4</u>	<u>e5</u>	<u>e6</u>
b7	<u>a8</u>	a9	a7	<u>b8</u>	b9	e7	e8	e9
d1	<u>c2</u>	c3	c1	<u>d2</u>	d3	e1	e2	e3
d4	<u>c5</u>	c6	c4	<u>d5</u>	d6	<u>g4</u>	<u>g5</u>	<u>g6</u>
d7	<u>c8</u>	c9	c7	<u>d8</u>	d9	g7	g8	g9

x

★

★

20

(for reference)

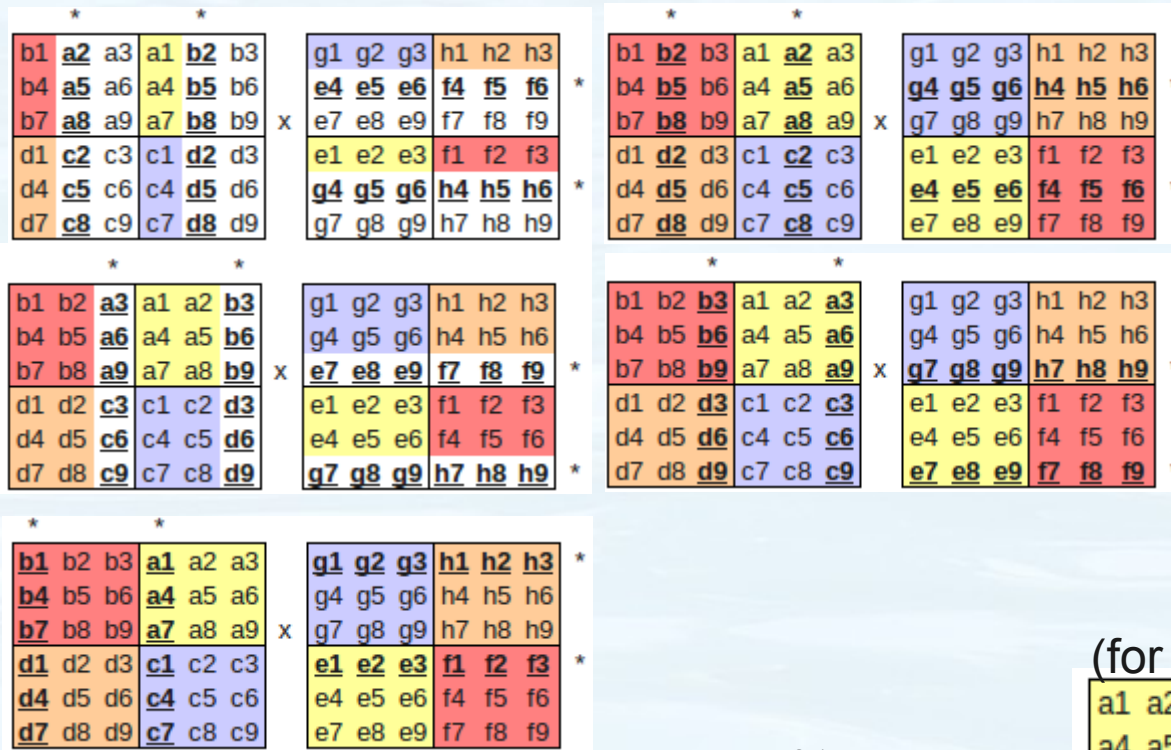
a1	a2	a3	b1	b2	b3	e1	e2	e3	f1	f2	f3
a4	a5	a6	b4	b5	b6	e4	e5	e6	f4	f5	f6
a7	a8	a9	b7	b8	b9	e7	e8	e9	f7	f8	f9
c1	c2	c3	d1	d2	d3	g1	g2	g3	h1	h2	h3
c4	c5	c6	d4	d5	d6	g4	g5	g6	h4	h5	h6
c7	c8	c9	d7	d8	d9	g7	g8	g9	h7	h8	h9

x



# Pipelined SUMMA - Example

- And so on



(for reference)

a1	a2	a3	b1	b2	b3
a4	a5	a6	b4	b5	b6
a7	a8	a9	b7	b8	b9
c1	c2	c3	d1	d2	d3
c4	c5	c6	d4	d5	d6
c7	c8	c9	d7	d8	d9

x

e1	e2	e3	f1	f2	f3
e4	e5	e6	f4	f5	f6
e7	e8	e9	f7	f8	f9
g1	g2	g3	h1	h2	h3
g4	g5	g6	h4	h5	h6
g7	g8	g9	h7	h8	h9