# Dense Matrix Algorithms

**KUMAR Chapter 8**

Jan Lemeire

Parallel Systems lab

December 16th 2011

Vrije Universiteit Brussel

**1. Matrix-vector Multiplication**

**2. Matrix-matrix Multiplication**

# Utility of Matrix Algorithms

Applied in several numerical and non-numerical contexts:

- 3D image calculations
- Solving (linear) equations
- Simulations of physical systems

E.g.: makes the basis of LINPACK, a software library for performing numerical linear algebra, by using the BLAS (Basic Linear Algebra Subprograms) libraries for performing basic vector and matrix operations

# Dense versus Sparse Matrices

**Dense matrices**: have no or few known zero entries

**Sparse matrices**: are populated primarily with zeros

- – often appear in science or engineering when solving partial differential equations
- – easily compressed
- – very large sparse matrices are impossible to manipulate with the standard algorithms, due to memory limitations
- – special versions of the algorithms are necessary and are more efficient
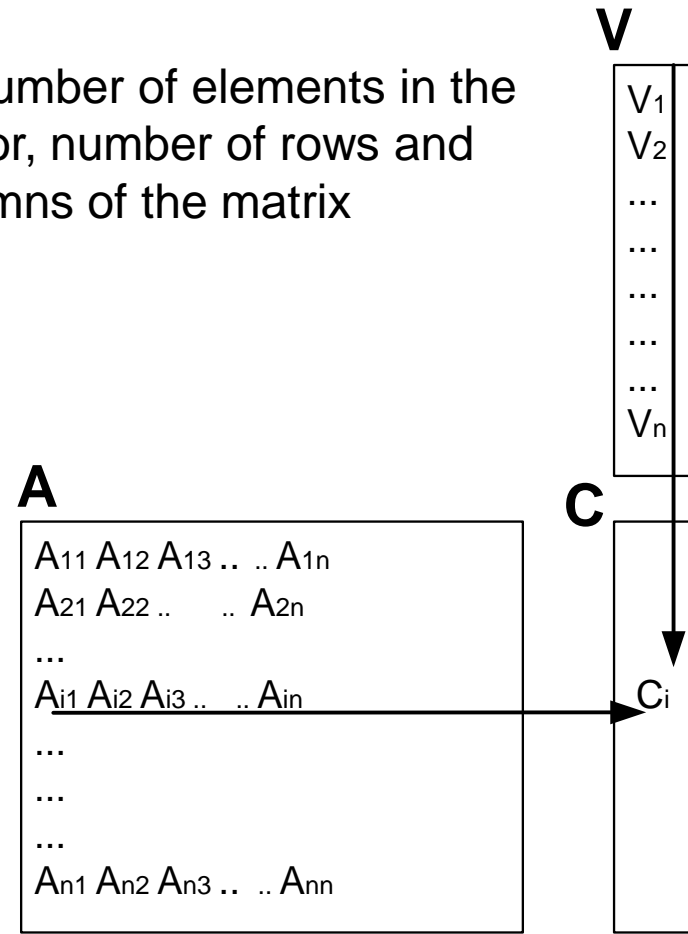
# Matrix – Vector Multiplication

$$C = A \times V$$
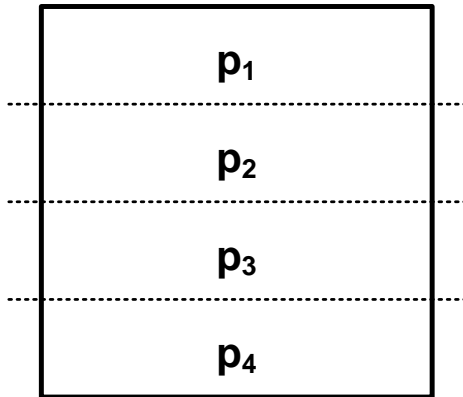
$$C_i = \sum_{k=1}^{n} A_{ik}.V_k \quad (i:1..n)$$

$$T_{seq} = \delta_{mm}.n^2$$

$n$ : number of elements in the vector, number of rows and columns of the matrix

**V**
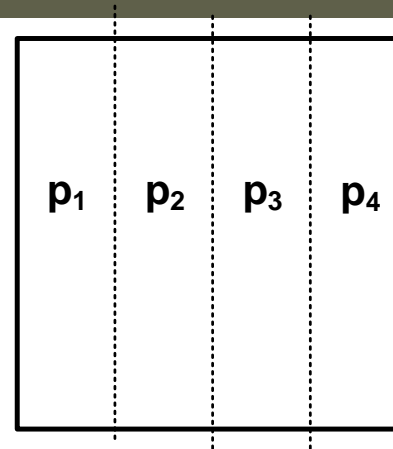
| |
|---|
| $V_1$ |
| $V_2$ |
| ... |
| ... |
| ... |
| ... |
| ... |
| $V_n$ |

```
for (i=0; i<n; i++){
   C[i]=0;
   for (k=0; k<n; k++){
      C[i]+=A[i, k]*V[k];
   }
}
```

**A**

$A_{11}\ A_{12}\ A_{13}\ ..\ ..\ A_{1n}$
$A_{21}\ A_{22}\ ..\ \quad ..\ A_{2n}$
...
$A_{i1}\ A_{i2}\ A_{i3}\ ..\ \quad ..\ A_{in}$
...
...
...
$A_{n1}\ A_{n2}\ A_{n3}\ ..\ ..\ A_{nn}$

**C**

$C_i$

# Matrix/Vector Partitioning?

| | | | |
|---|---|---|---|
| **p₁** | | | |
| **p₂** | | | |
| **p₃** | | | |
| **p₄** | | | |

**Row-wise block striping**

| p₁ | p₂ | p₃ | p₄ |
|----|----|----|----|

**Column-wise block striping**

| p₁ | p₂ | p₃ |
|----|----|----|
| p₄ | p₅ | p₆ |
| p₇ | p₈ | p₉ |

**Checkerboard partitioning**

| |
|---|
| **p₁** |
| **p₂** |
| **p₃** |
| **p₄** |

**Vector partitioning**

# Data & results distributed

o In parallel applications, data remains distributed while operations (such as vector-matrix operations) are performed on them.

o In the following we assume that the data is already distributed among the processors
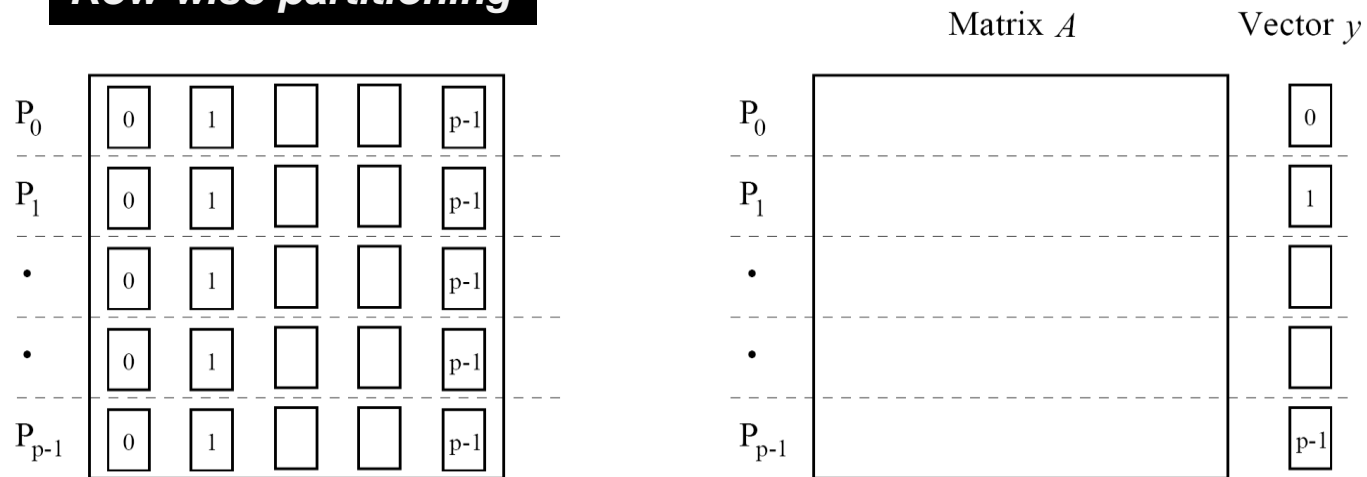
# Parallel M x V Multi-plication Version 1

## (n=p)

**Matrix $A$** — **Vector $x$**

**(a)** Initial partitioning of the matrix and the starting vector $x$

**Row-wise partitioning**

**Processes**

**(b)** Distribution of the full vector among all the processes by all-to-all broadcast

**(c)** Entire vector distributed to each process after the broadcast

**Matrix $A$** — **Vector $y$**

**(d)** Final distribution of the matrix and the result vector $y$

**Figure 8.1** Multiplication of an $n \times n$ matrix with an $n \times 1$ vector using rowwise block 1-D partitioning. For the one-row-per-process case, $p = n$.

# Parallel M x V Multiplication Version 1    p < n

n/p rows of matrix and n/p elements of vector per processor

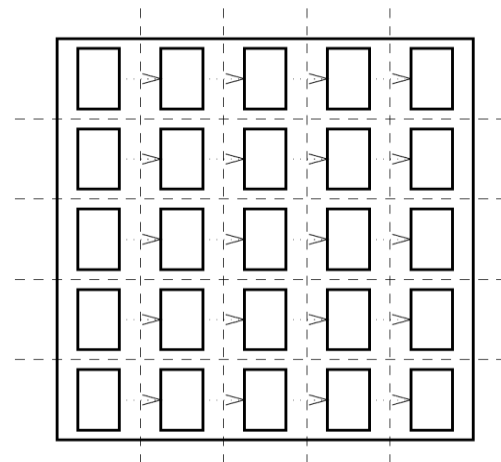# Parallel M x V Multi-plication Version 2

## (n²=p)



(a) Initial data distribution and communication steps to align the vector along the diagonal
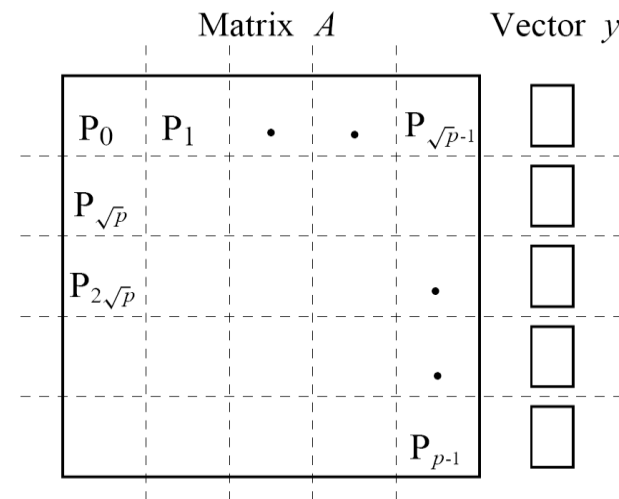
(b) One-to-all broadcast of portions of the vector along process columns

**checkerboard partitioning**

(c) All-to-one reduction of partial results

(d) Final distribution of the result vector

**Figure 8.2** Matrix-vector multiplication with block 2-D partitioning. For the one-element-per-process case, $p = n^2$ if the matrix size is $n \times n$.

# Parallel M x V Multiplication Version 2    $p < n^2$

$(n/\sqrt{p})$ x $(n/\sqrt{p})$ blocks of matrix and $n/\sqrt{p}$ elements of vector per processor

## 1. Matrix-vector Multiplication

## 2. Matrix-matrix Multiplication

# Matrix Multiplication

$$C = A \times B$$

$$C_{ij} = \sum_{k=1}^{n} A_{ik}.B_{kj} \quad (i, j : 1..n)$$

$$T_s = \delta_{mm}.n^3$$

```
for (i=0; i<n; i++){
    for (j=0; j<n; j++){
        C[i,j]=0;
        for (k=0; k<n; k++){
            C[i,j]+=A[i, k]*B[k,j];
        }
    }
}
```

**B**

```
B11 B12 .. B1j .. .. B1n
B21 B22 .. B2j    .. B2n
...
...
...
...
...
Bn1 Bn2 .. Bnj  .. Bnn
```

**A**

```
A11 A12 A13 .. .. A1n
A21 A22 ..       .. A2n
...
Ai1 Ai2 Ai3 ..    .. Ain
...
...
...
An1 An2 An3 ..  .. Ann
```

**C**

$C_{ij}$

# MxM: one-step version

- B is sent to all processors
- Computation in 1 step
-  Same amount of communication



computation step

# Alternate shift-compute version



- Algorithm alternates p computation and communication steps
- Computation step: each processor multiplies its A submatrix with its B submatrix, resulting in a submatrix of C. The black circles indicate the step in which each submatrix is computed.
- After multiplication: processor sends it B submatrix to next processor and receives one from the preceding processor. The communication forms a **circular shift operation**.

# **Parallel MxM**: Execution Profile



# **Speedup=2.55 Efficiency = 85%**

# Theoretical Analysis of Matrix Multiplication

◆ Computation time

$$T^i_{work} = \frac{n^3}{p}.\delta_{mm}$$

◆ Communication time (slave)

$$T^i_{comm} = (p+1).t_s + (1+\frac{2}{p}).n^2.t_w$$
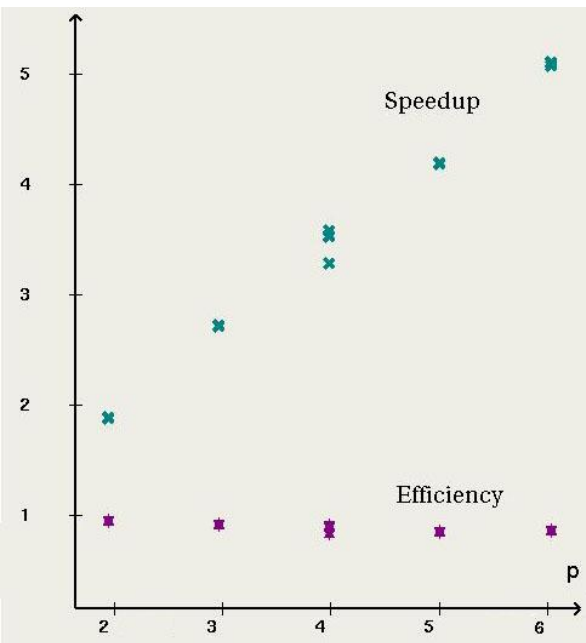
➡ Explains parameter dependence of performance:
  ✦ Total overhead = $O(n^2.p)$
  ✦ Computation = $O(n^3/p)$        ➡ Overhead ratio=$O(p/n)$
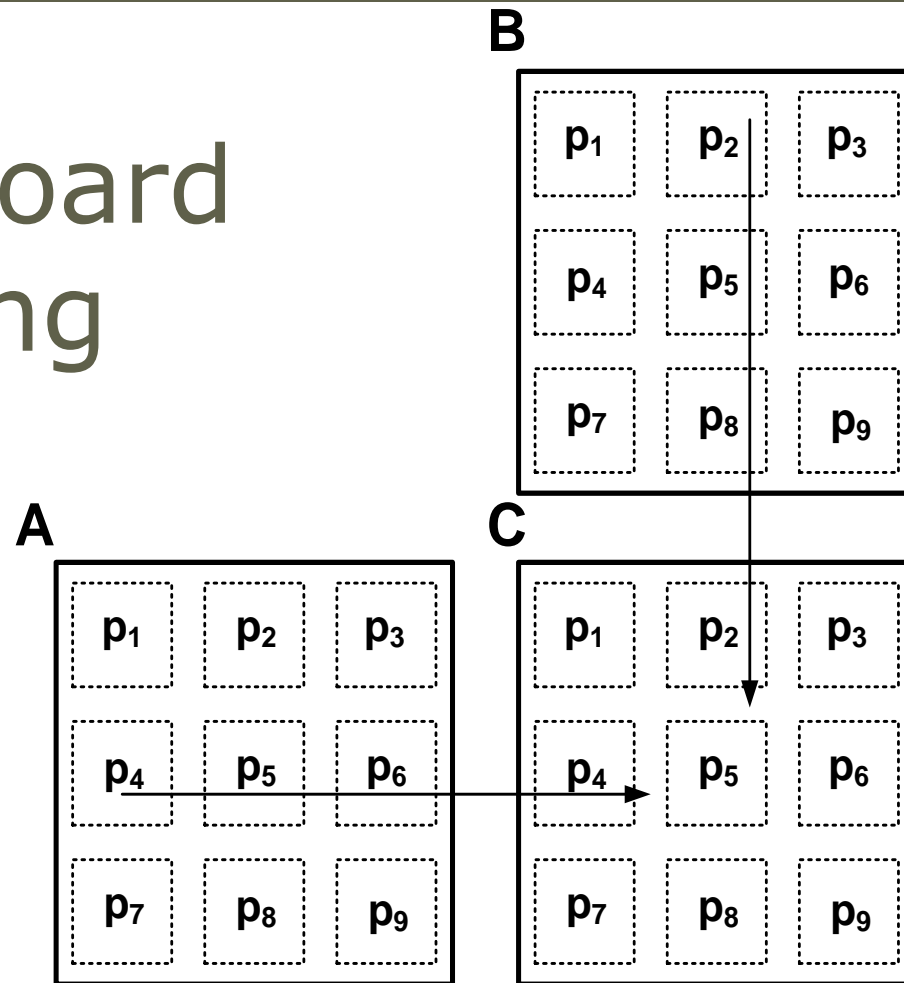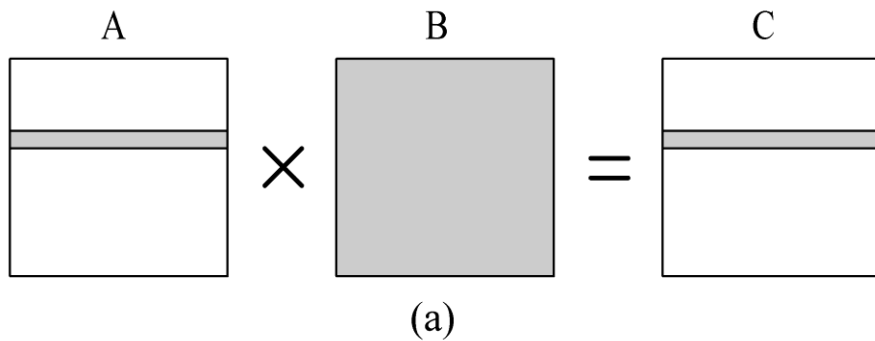  ✦ Ratio computation/communication = $2n/p$

# Parameter Dependence of Matrix Multiplication
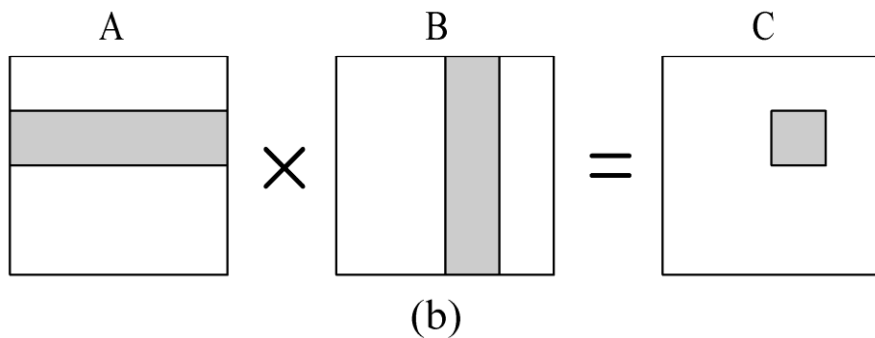


*n*: work size, here: matrix size

# V3: Cannon's algorithm

Checkerboard partitioning

# Elements of A and B Needed to Compute a Process's Portion of C
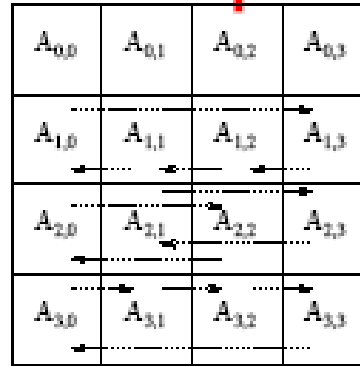
A × B = C
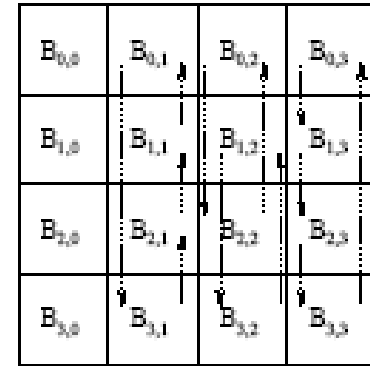
(a)

Algorithm 1 & 2

A × B = C

(b)

Cannon's Algorithm

Why faster? Ratio perimeter/surface is minimal for a square!
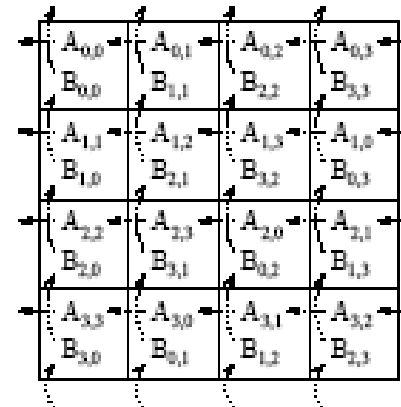
# Cannon's parallel MxM

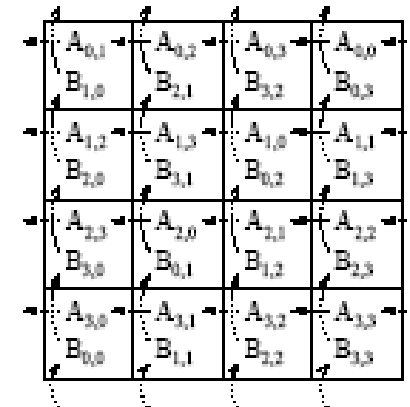## Communication steps on 16 processes
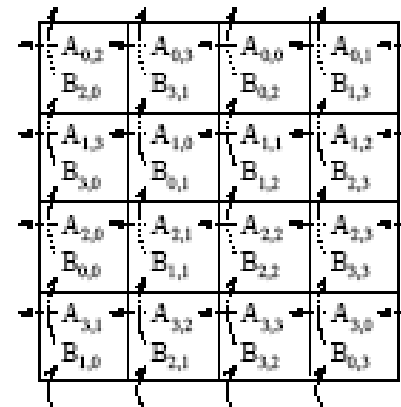


(a) Initial alignment of A
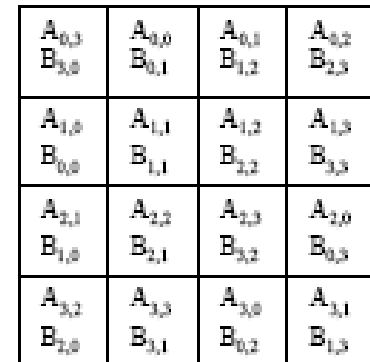
(b) Initial alignment of B

(c) A and B after initial alignment

(d) Submatrix locations after first shift
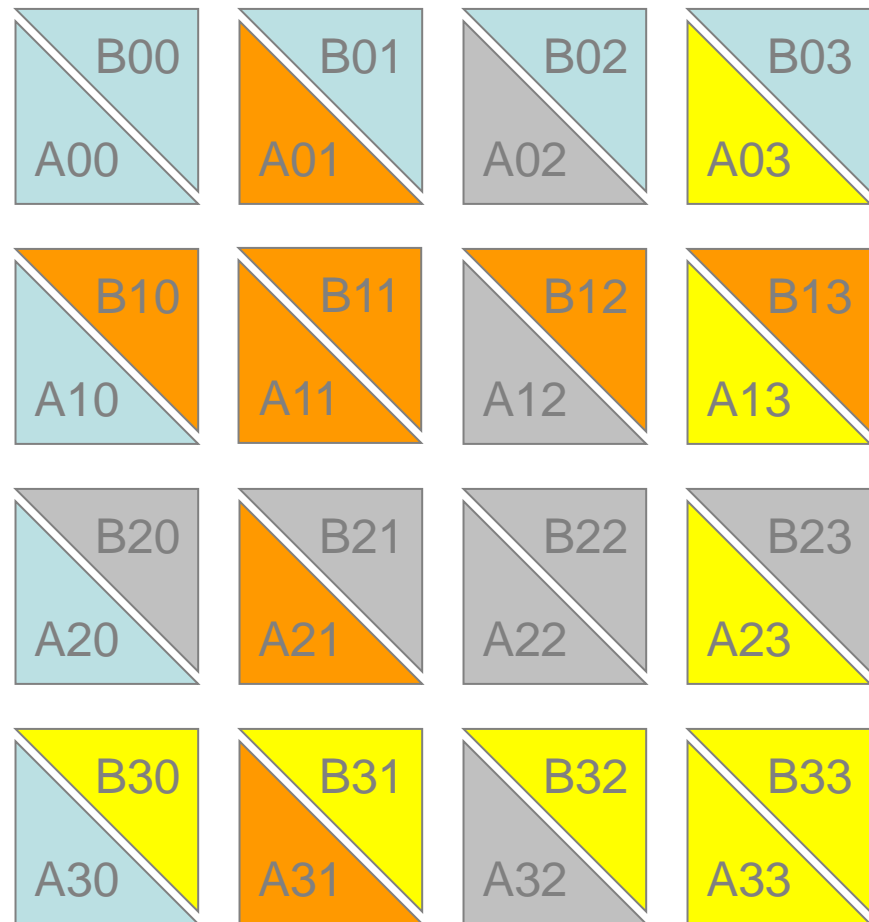
(e) Submatrix locations after second shift

(f) Submatrix locations after third shift

# Initially, blocks Need to Be Aligned

Each triangle represents a matrix block

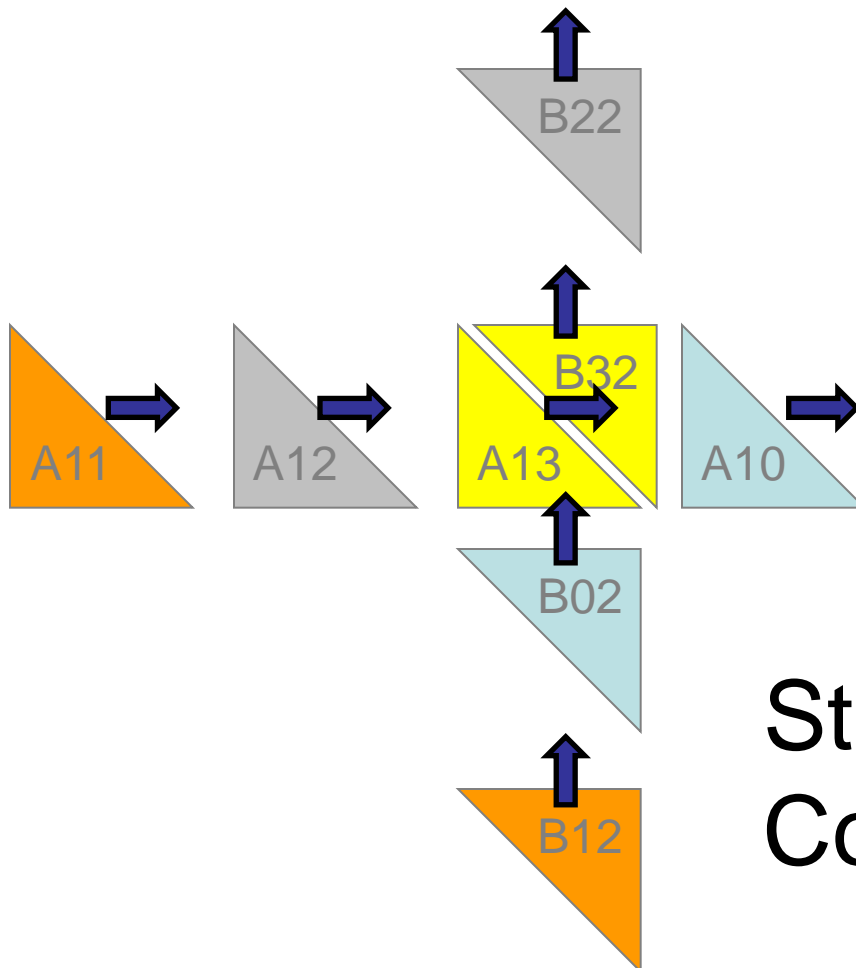*Only same-color triangles should be multiplied*

# Rearrange Blocks

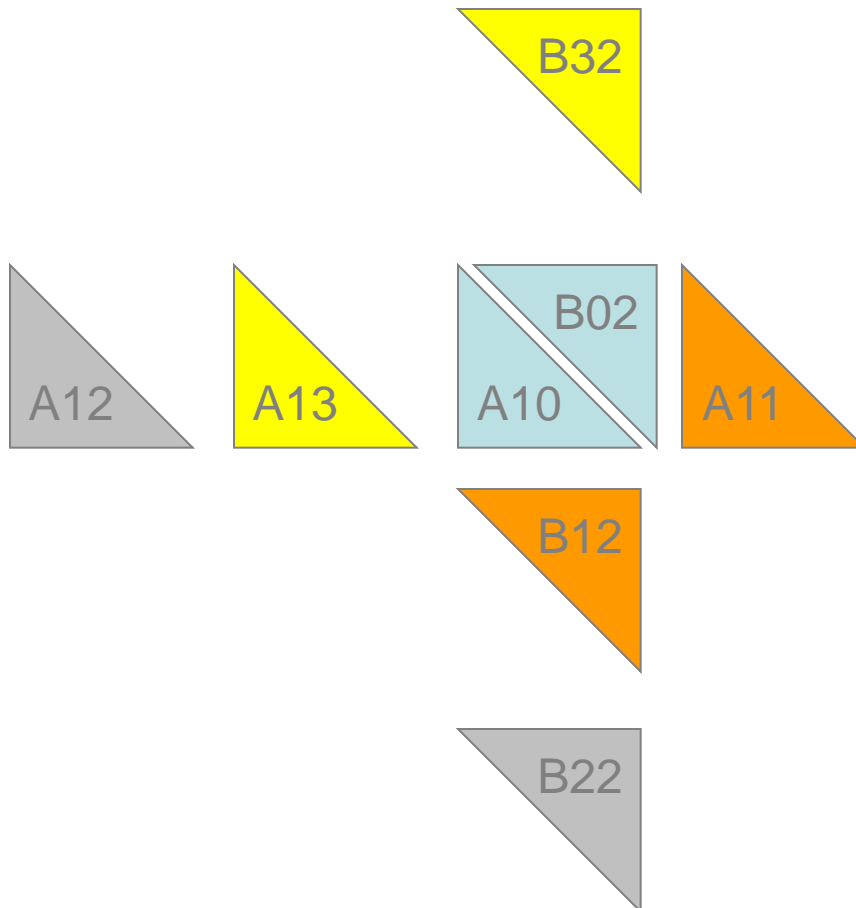| | | | |
|---|---|---|---|
| B00 / A00 | B11 / A01 | B22 / A02 | B33 / A03 |
| B10 / A11 | B21 / A12 | B32 / A13 | B03 / A10 |
| B20 / A22 | B31 / A23 | B02 / A20 | B13 / A21 |
| B30 / A33 | B01 / A30 | B12 / A31 | B23 / A32 |

Block $A_{ij}$ cycles left *i* positions

Block B*ij* cycles up *j* positions

# Consider Process $P_{1,2}$

B22

A11 → A12 → A13 / B32 → A10 →

B02

B12
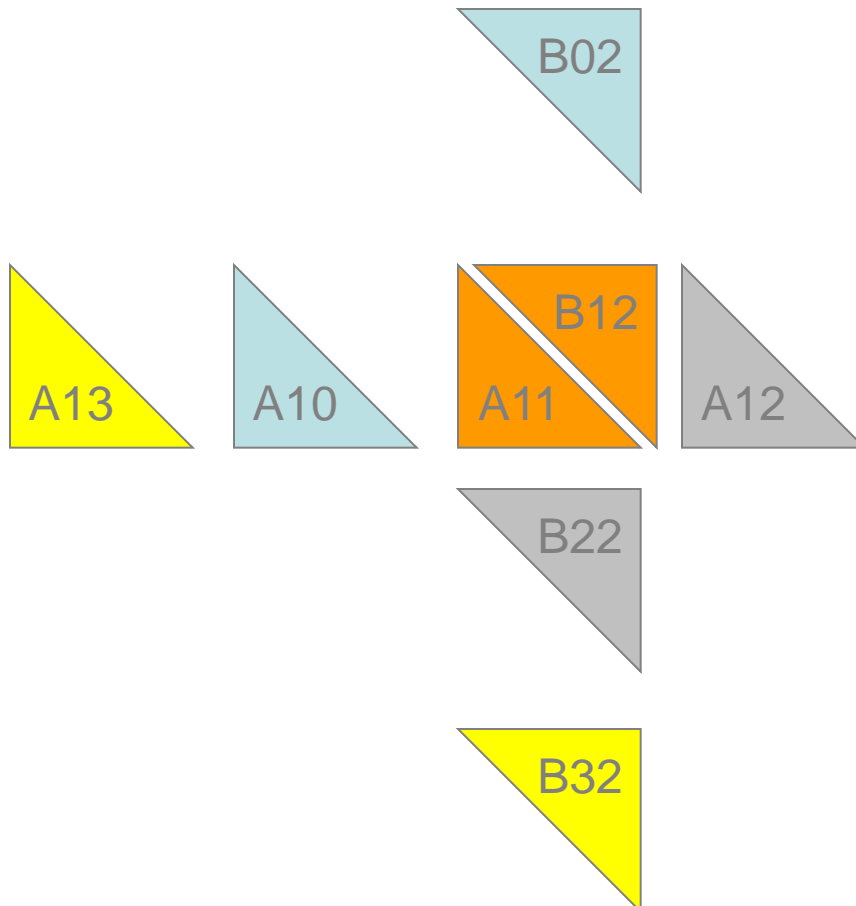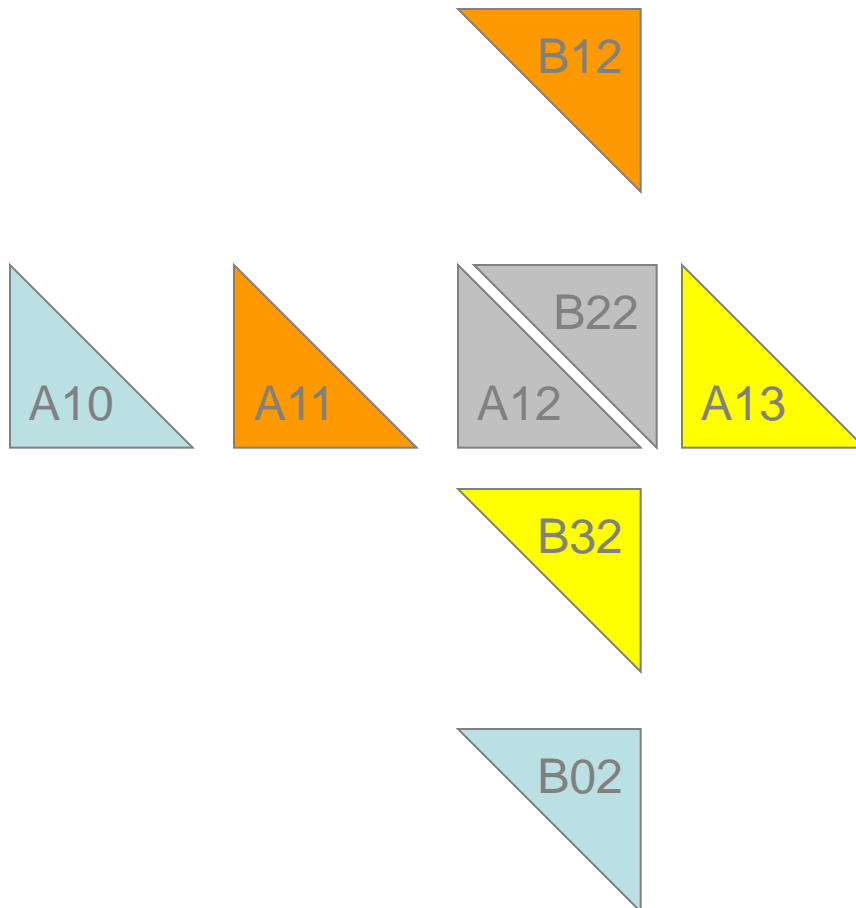
Step 1:
Computation & shift

# Consider Process $P_{1,2}$
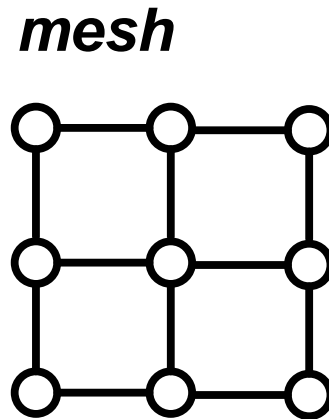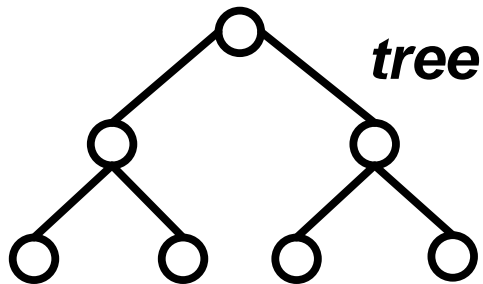


Step 2

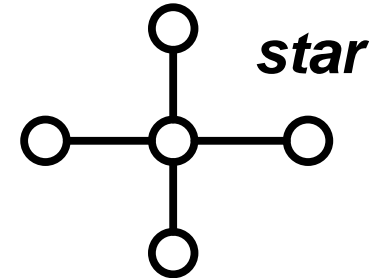# Consider Process $P_{1,2}$



Step 4

# Complexity Analysis

o Algorithm has $\sqrt{p}$ iterations

o During each iteration process multiplies two ($n$ / $\sqrt{p}$ )
  × ($n$ / $\sqrt{p}$ ) matrices: $\Theta(n^3 / p^{3/2})$

o **Computational complexity**: $\Theta(n^3 / p)$  [the same]

o During each iteration process sends and receives
  two blocks of size ($n$ / $\sqrt{p}$ ) × ($n$ / $\sqrt{p}$ )

o **Communication complexity**: $\Theta(n^2 / \sqrt{p})$  [lower!]

# Efficient Interconnection Networks for Cannon's MxM?

*line*

*ring*

*star*

*tree*

*mesh*

*hypercube*

*wraparound mesh*

# Memory need of each processor

*As a function of **n**, **p** and **b** (the number of bytes per element)*

◆ Sequential algorithm:

◆ One-step algorithm:

◆ Alternate algorithm:

◆ Cannon's algorithm:

# Special MPI functions

◆ MxV version 2: Reduce operation

◆ Cannon: Shift operation through a SendRecv_replace call


◆ Submatrix sending:
  ✦ Consist of equally spaced blocks
  ✦ DataType.Vector(int count, int blocklength, int stride, Datatype oldtype)

# MxM on GPU

- Initially, matrices are copied to GPU
  - If they are not still in memory from previous matrix operations, keep pointers in CPU to the data on the GPU
- Every thread computes 1 element of C.
- Not enough memory to put all data in shared memory (16K)
- On one multiprocessor, 1 block of threads computes 1 block of the C matrix
- Iteratively copy A row blocks and B column blocks to shared memory.
- Result: 200x speedup, 50x if compared to quadcore