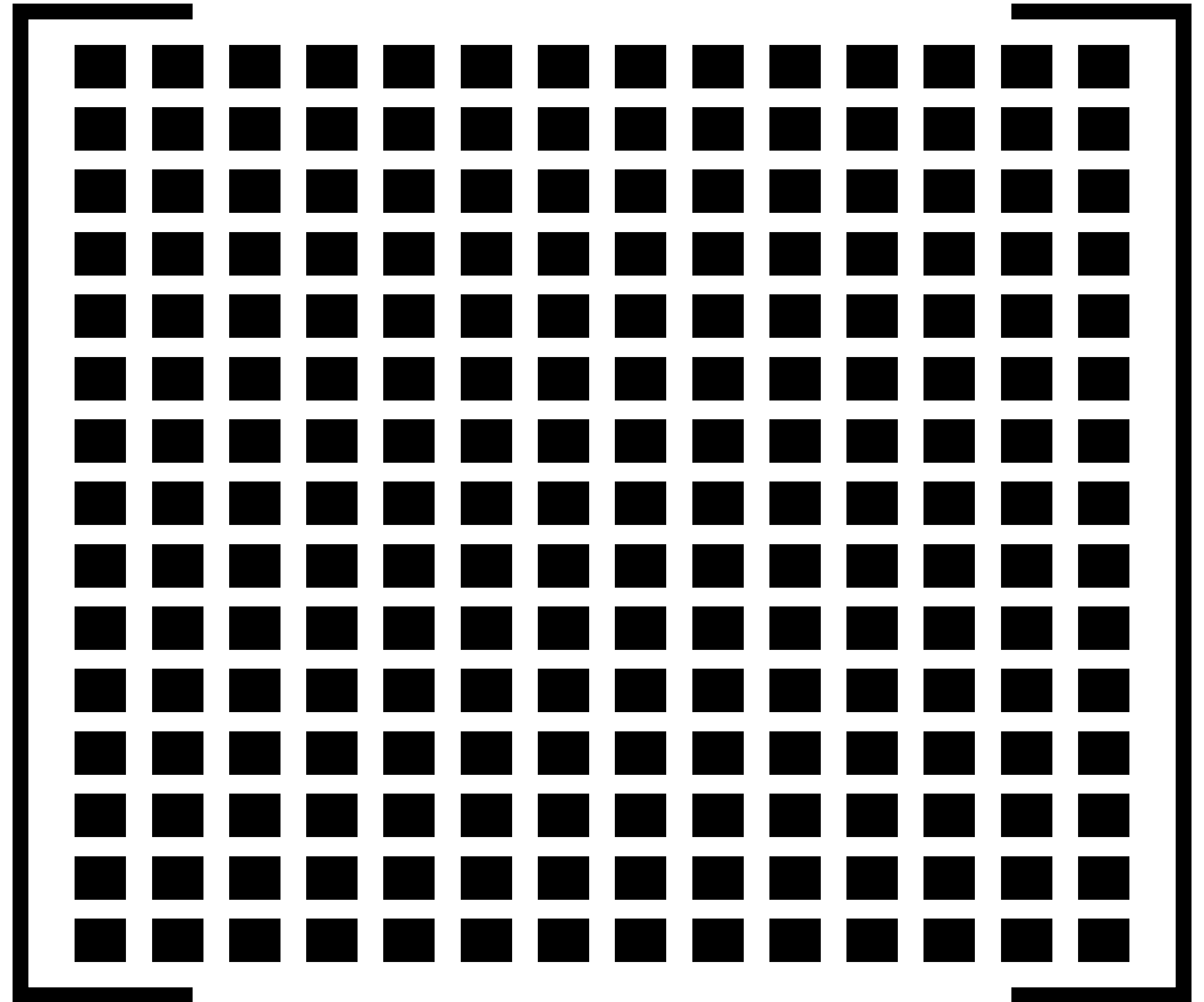# Introduction to Parallel Processing

Lecture 5: MPI Matrix-Matrix Multiplication

Professor Amanda Bienz
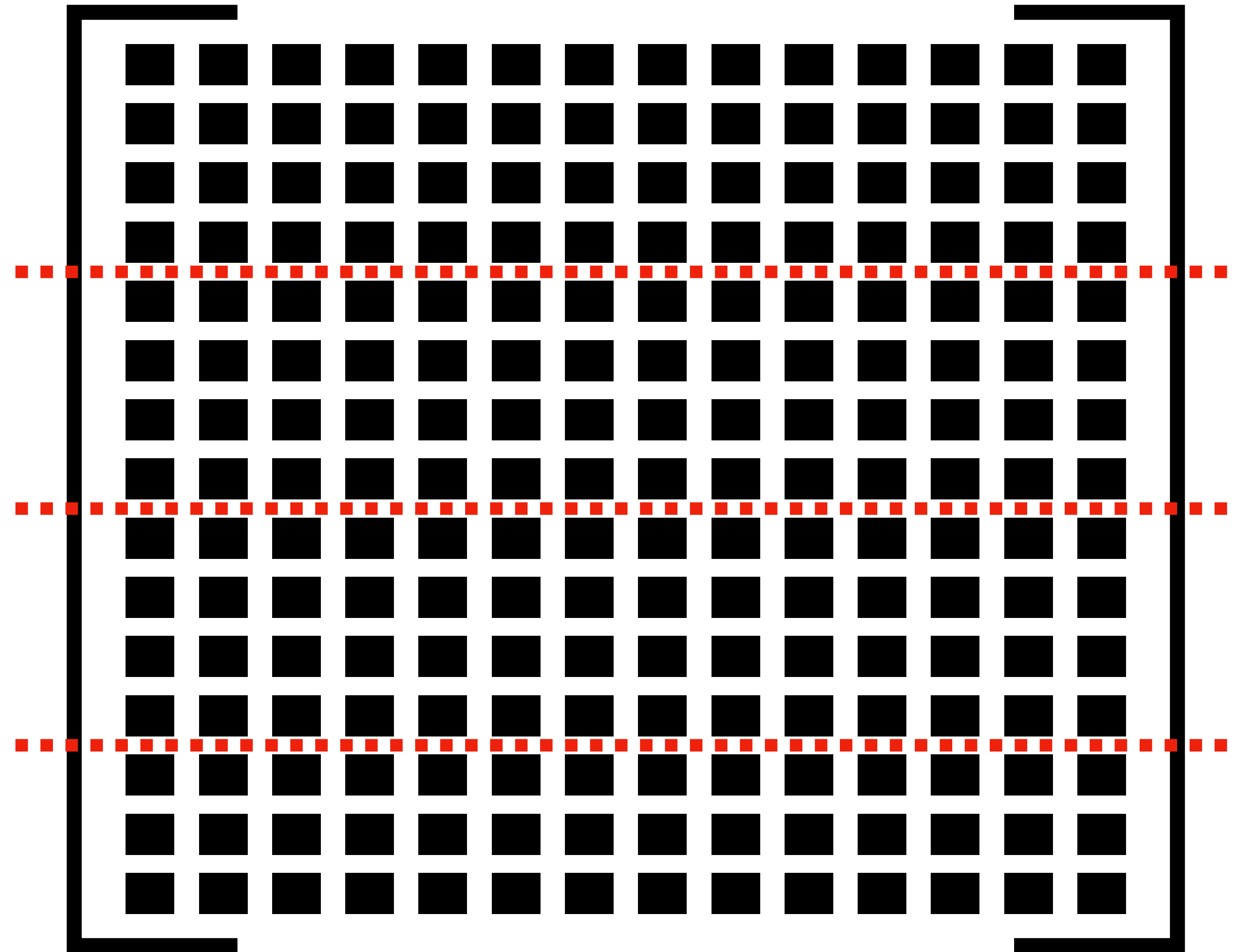
# Parallel Matrix Partition

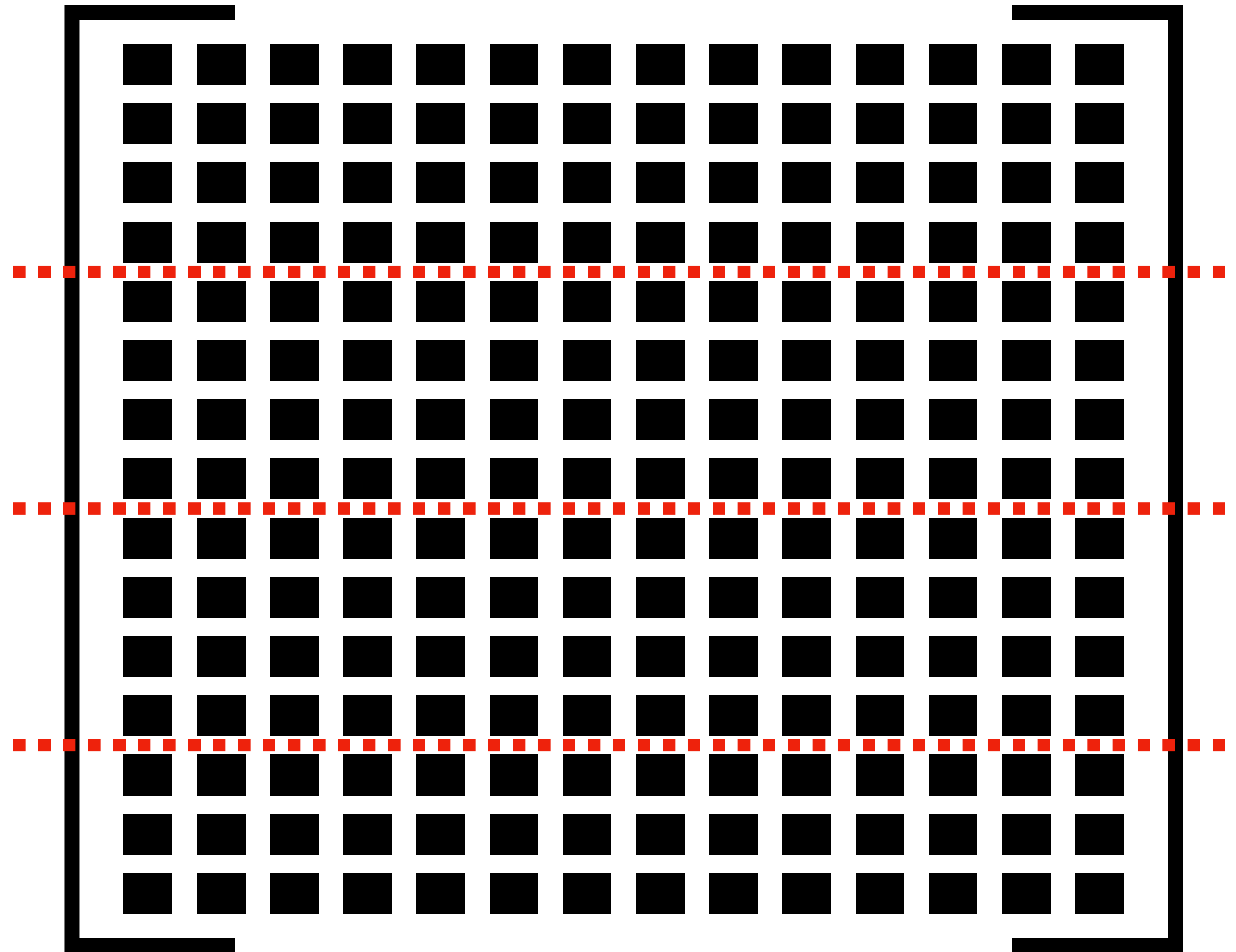- How do we partition this matrix, so that a part of the matrix is on each process?

# Parallel Matrix Partition

- How do we partition this matrix, so that a part of the matrix is on each process?

- **Row Wise Partition : Each process holds contiguous chunk of rows of the matrix**
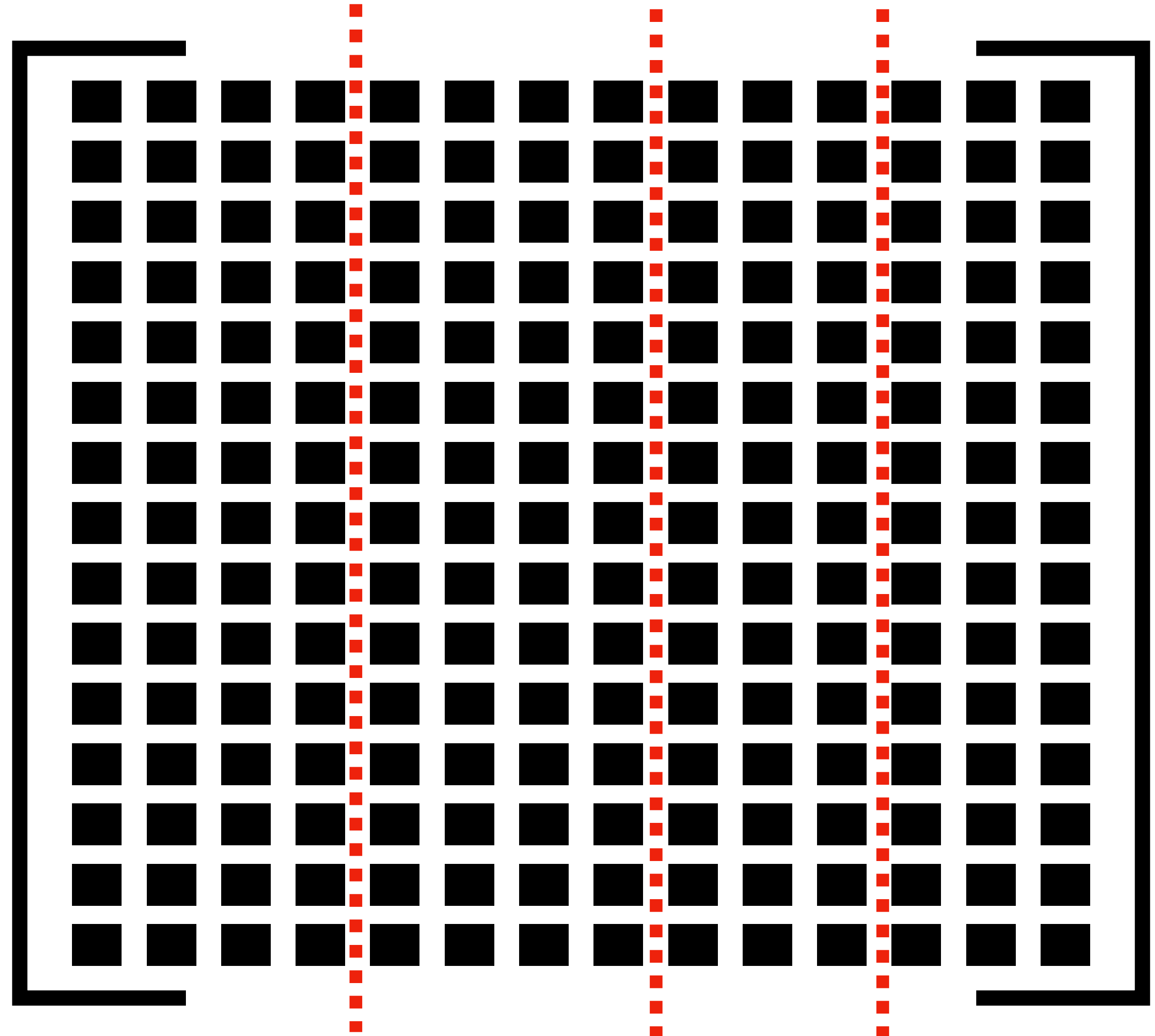
# Parallel Matrix Partition

- How do we partition this matrix, so that a part of the matrix is on each process?

- **Problem : What is the smallest partition that we can make?**

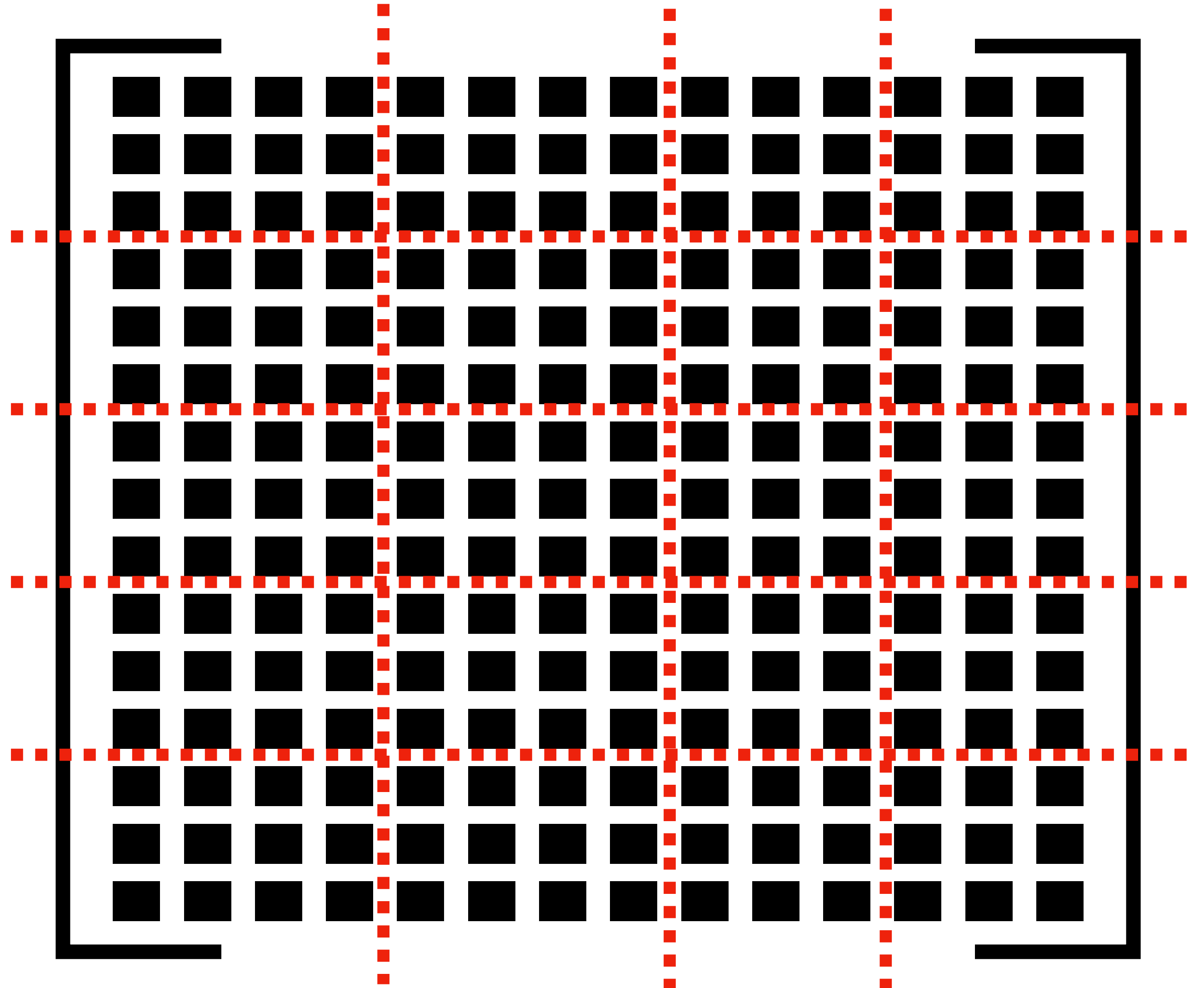- **For an nxn matrix, smallest row-wise partition holds n values**

# Parallel Matrix Partition

- How do we partition this matrix, so that a part of the matrix is on each process?

- **Similarly, could look at column-wise partition, but we run into the same issue**

# Parallel Matrix Partition
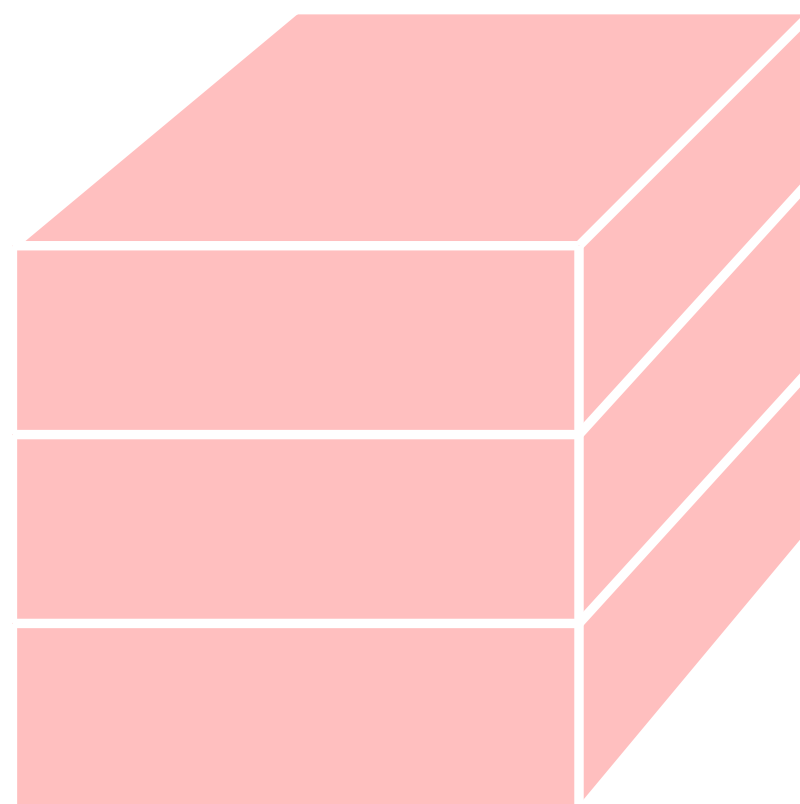
- How do we partition this matrix, so that a part of the matrix is on each process?

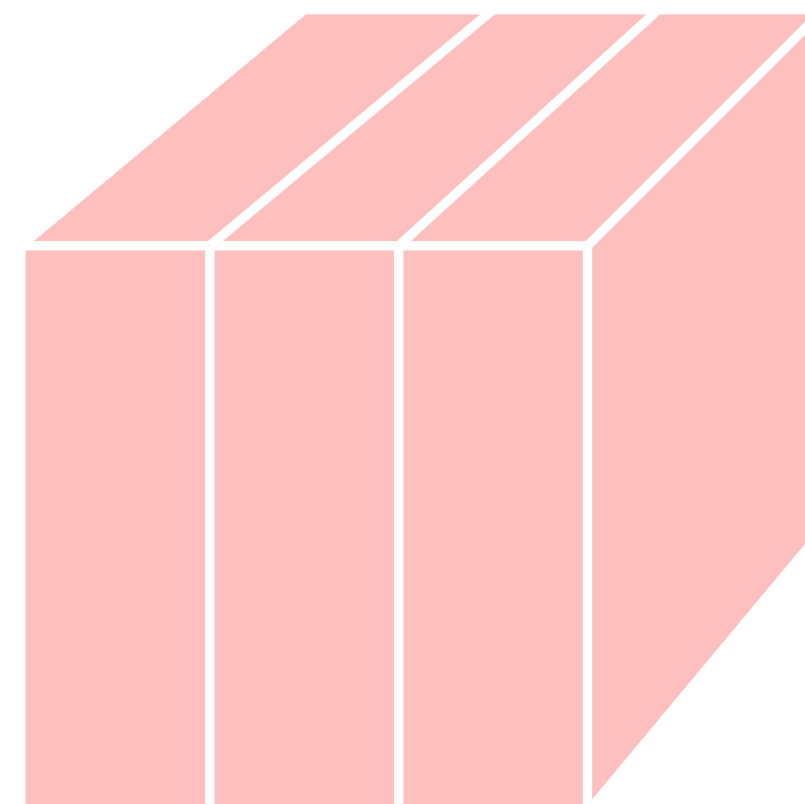- **Solution : 2-dimensional partitions**

# Parallel Matrix Partition

- Could have a 3D partition

- Partition by operations rather than elements of matrix

- Sometimes a good idea for GPUs, but we typically stick with 2D in HPC due to communication overhead

1-D row       1-D col       2-D       3-D

# Matrix Multiplication

# Simplest Implementation

# Simplest Implementation

- Every process sends entire rows of A and columns of B to every process that needs them

$$
\begin{bmatrix}
C0 & C1 & C2 & C3 \\
C4 & C5 & C6 & C7 \\
C8 & C9 & C10 & C11 \\
C12 & C13 & C14 & C15
\end{bmatrix}
=
\begin{bmatrix}
A3 \\
A7 \\
A11 \\
A15
\end{bmatrix}
\begin{bmatrix}
B0 & B1 & B2 & B3 \\
B4 & B5 & B6 & B7 \\
B8 & B9 & B10 & B11 \\
B12 & B13 & B14 & B15
\end{bmatrix}
$$

# Simplest Implementation

- Every process sends entire rows of A and columns of B to every process that needs them

# Fox's Algorithm

| | | |
|---|---|---|
| A0 B0 | A1 B1 | A2 B2 |
| A3 B3 | A4 B4 | A5 B5 |
| A6 B6 | A7 B7 | A8 B8 |

- The processes are laid out in a grid, so each process p has a row I and column J in the process grid
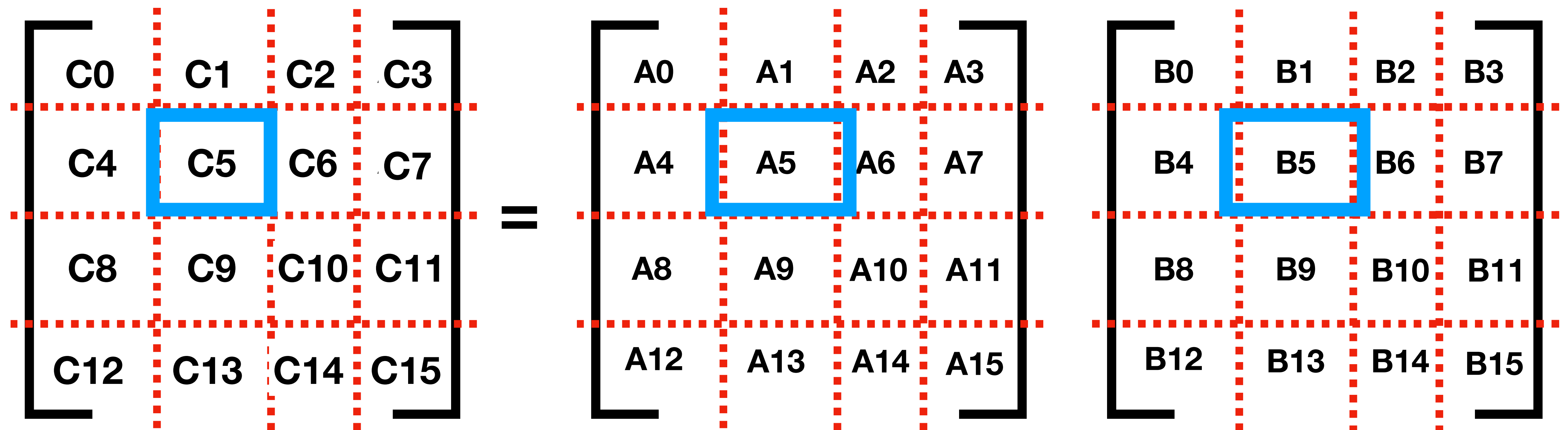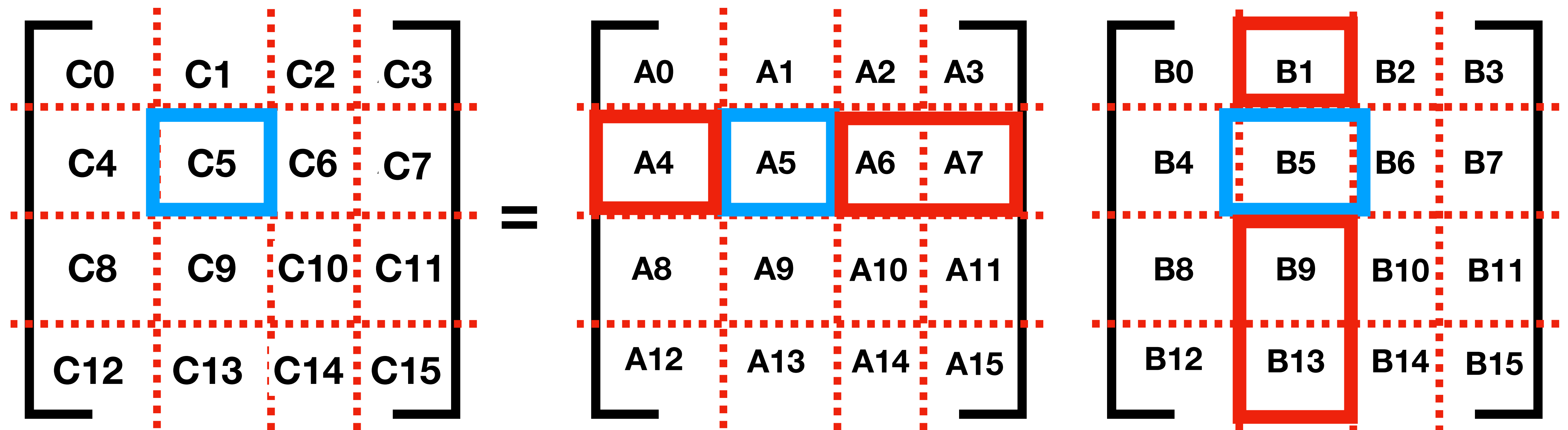
- Gather all of A among all processes in each process row I

# Fox's Algorithm



| A0 B0 | A1 B1 | A2 B2 |
| A3 B3 | A4 B4 | A5 B5 |
| A6 B6 | A7 B7 | A8 B8 |

- Multiple A by local portion of B, then rotate B along process column J

  - Process in row I, column J sends local portion of B to process in row I+1, column J

  - Process in row I, column J receives new portion of B from process in row I-1, column J

# Fox's Algorithm

| A0 B0 | A1 B1 | A2 B2 |
|---|---|---|
| A3 B3 | A4 B4 | **A5 B5** |
| A6 B6 | A7 B7 | A8 B8 |

- Need to multiply A0 by :

  - B0, B1, B2

- Multiply A1 by :

  - B1, B4, B7

- C5 = [A3*B2 + A4*B5 + A5*B8]

# Fox's Algorithm

| | | |
|---|---|---|
| **A0:2 B0** | **A0:2 B1** | **A0:2 B2** |
| **A3:5 B3** | **A3:5 B4** | **A3:5 B5** |
| **A6:8 B6** | **A6:8 B7** | **A6:8 B8** |

- Need to multiply A0 by :

  - B0, B1, B2

- Multiply A1 by :

  - B1, B4, B7

- C5 = [A3*B2 + A4*B5 + A5*B8]

# Fox's Algorithm



- Need to multiply A0 by :

  - B0, B1, B2

- Multiply A1 by :

  - B1, B4, B7

- C5 = [A3*B2 + A4*B5 + A5*B8]

# Fox's Algorithm

| | | |
|---|---|---|
| **A0:2 B3** | **A0:2 B4** | **A0:2 B5** |
| **A3:5 B6** | **A3:5 B7** | **A3:5 B8** |
| **A6:8 B0** | **A6:8 B1** | **A6:8 B2** |

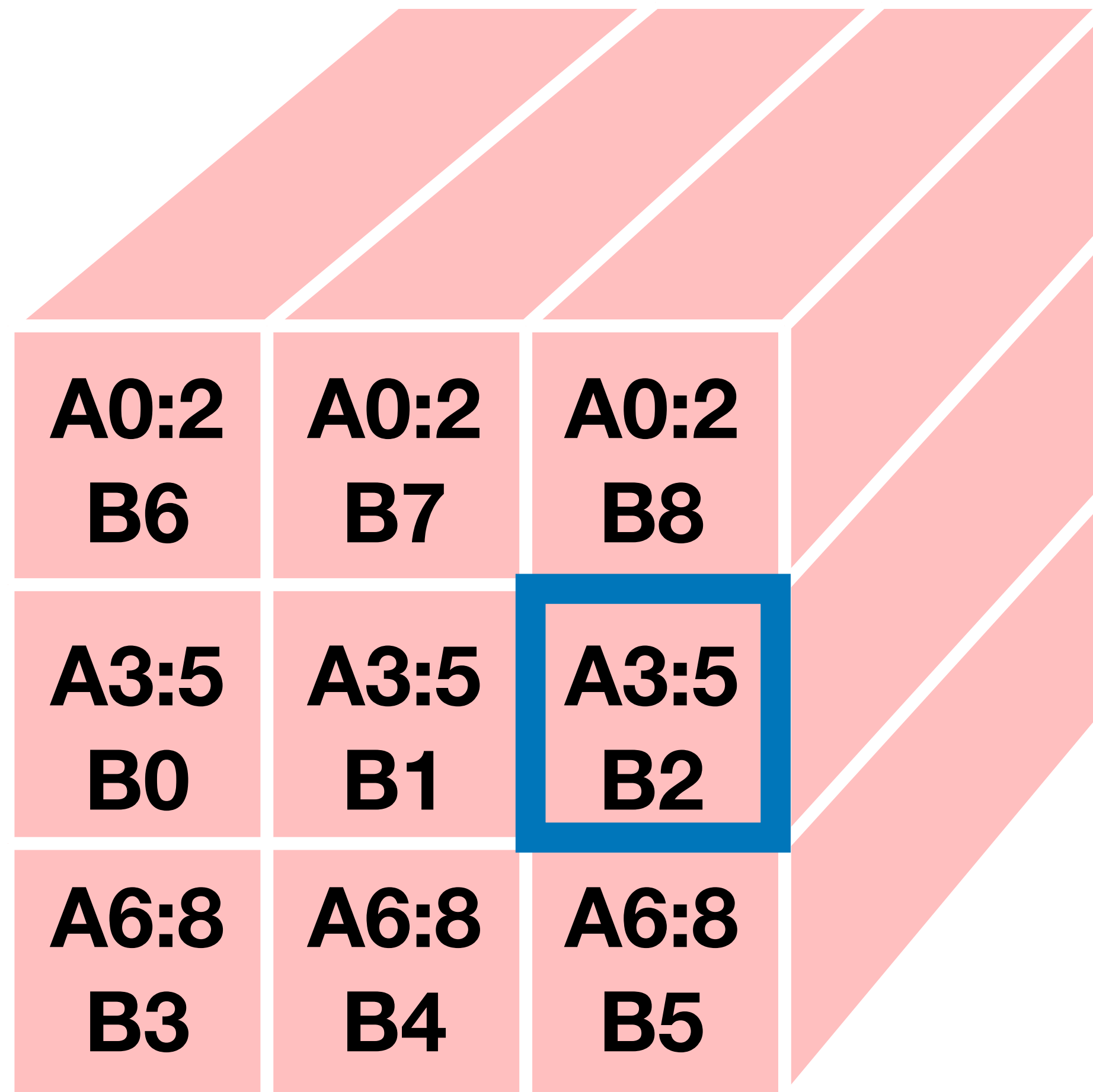- Need to multiply A0 by :
  - B0, B1, B2
- Multiply A1 by :
  - B1, B4, B7
- C5 = [A3*B2 + A4*B5 + A5*B8]

# Fox's Algorithm

- Hold all needed values of A, but only one partition of B at a time

$$
\begin{bmatrix}
C0 & C1 & C2 & C3 \\
C4 & C5 & C6 & C7 \\
C8 & C9 & C10 & C11 \\
C12 & C13 & C14 & C15
\end{bmatrix}
=
\begin{bmatrix}
& & A3 & \\
& & A7 & \\
& & A11 & \\
& & A15 &
\end{bmatrix}
\begin{bmatrix}
B0 & B1 & B2 & B3 \\
B4 & B5 & B6 & B7 \\
B8 & B9 & B10 & B11 \\
B12 & B13 & B14 & B15
\end{bmatrix}
$$

# Fox's Algorithm

- Hold all needed values of A, but only one partition of B at a time

# Fox's Algorithm

- Hold all needed values of A, but only one partition of B at a time

# Fox's Algorithm

- Hold all needed values of A, but only one partition of B at a time

# Cannon's Algorithm

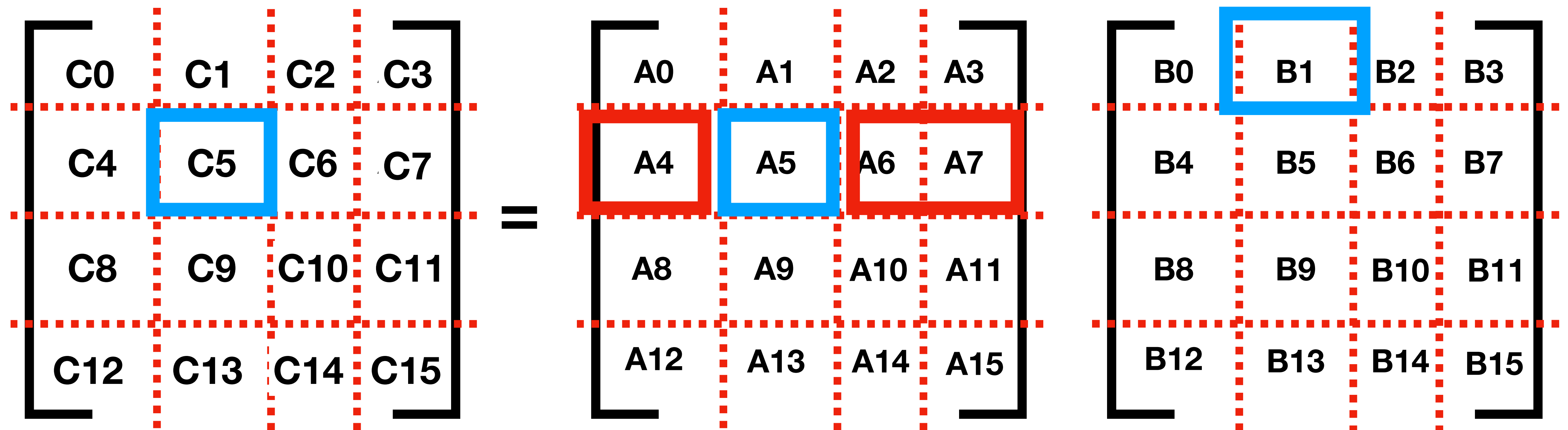|  | Col 0 | Col 1 | Col 2 |
|---|---|---|---|
| **Row 0** | **A0** **B0** | **A1** **B1** | **A2** **B2** |
| **Row 1** | **A3** **B3** | **A4** **B4** | **A5** **B5** |
| **Row 2** | **A6** **B6** | **A7** **B7** | **A8** **B8** |

- Initial Shift:

  - Process in row I, column J sends local portion of A to process in row I, column J-I

  - Process in row I, column J sends local portion of B to process in row I-J, column J

# Cannon's Algorithm

|  | Col 0 | Col 1 | Col 2 |
|---|---|---|---|
| Row 0 | A0 B0 | A1 B4 | A2 B8 |
| Row 1 | A4 B3 | A5 B7 | A3 B2 |
| Row 2 | A8 B6 | A6 B1 | A7 B5 |

- Initial Shift:

  - Process in row I, column J sends local portion of A to process in row I, column J-I

  - Process in row I, column J sends local portion of B to process in row I-J, column J

# Cannon's Algorithm

- After initial shift, rotate in opposite direction:

  - Multiply local portion of A with local portion of B

  - Process [I,J] sends local portion of A to Process[I,J+1]

  - Process [I,J] sends local portion of B to process [I+1,J]

# Cannon's Algorithm



- Need to multiply A0 by :

  - B0, B1, B2

- Multiply A1 by :

  - B1, B4, B7

- C5 = [A3*B2 + A4*B5 + A5*B8]

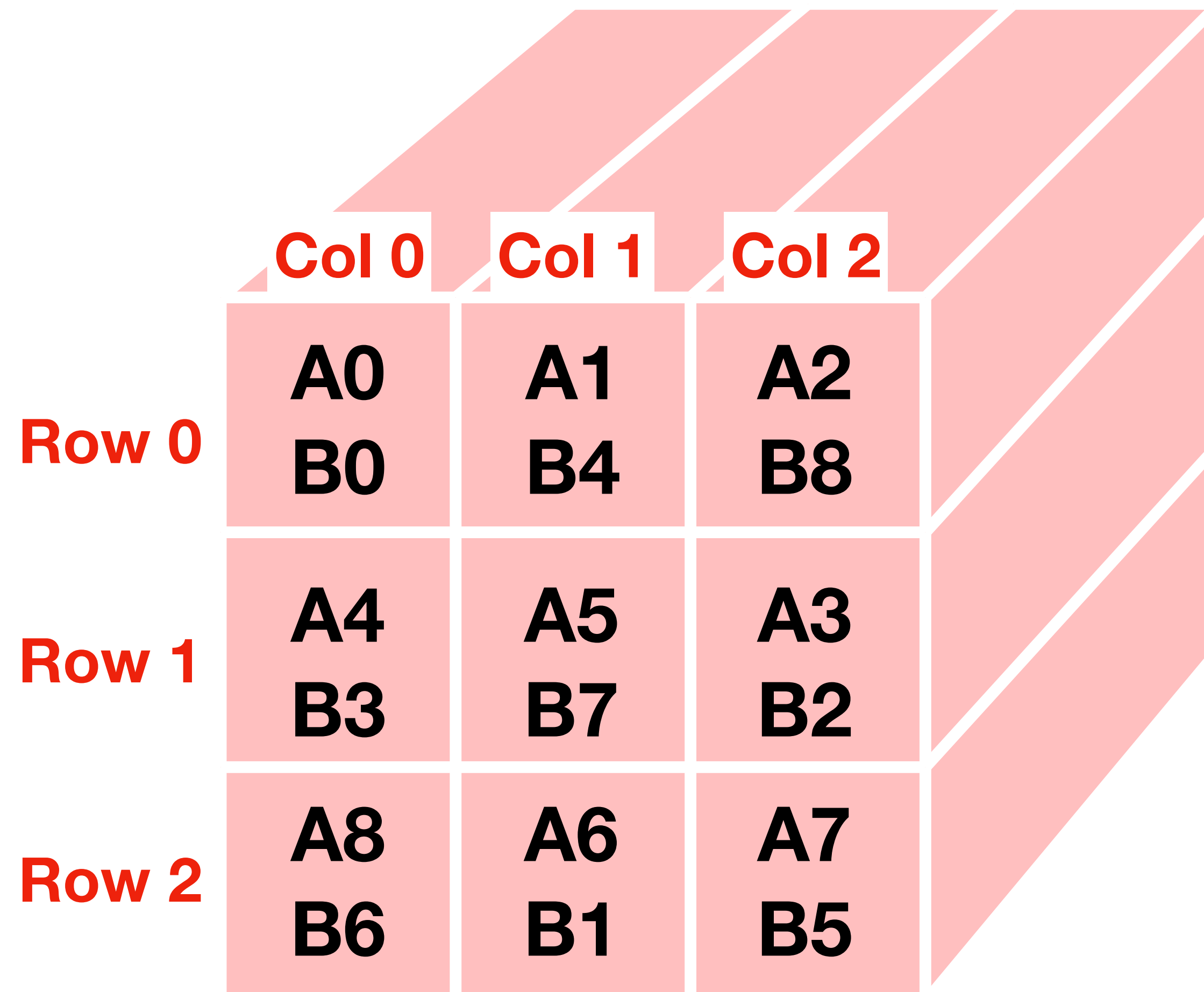# Cannon's Algorithm

| A0 B0 | A1 B4 | A2 B8 |
| A4 B3 | A5 B7 | A3 B2 |
| A8 B6 | A6 B1 | A7 B5 |

- Need to multiply A0 by :

  - B0, B1, B2

- Multiply A1 by :

  - B1, B4, B7

- C5 = [A3*B2 + A4*B5 + A5*B8]

# Cannon's Algorithm

| | | |
|---|---|---|
| A2 B6 | A0 B1 | A1 B5 |
| A3 B0 | A4 B4 | **A5 B8** |
| A7 B3 | A8 B7 | A6 B2 |

- Need to multiply A0 by :

  - B0, B1, B2

- Multiply A1 by :

  - B1, B4, B7

- C5 = [A3*B2 + A4*B5 + A5*B8]
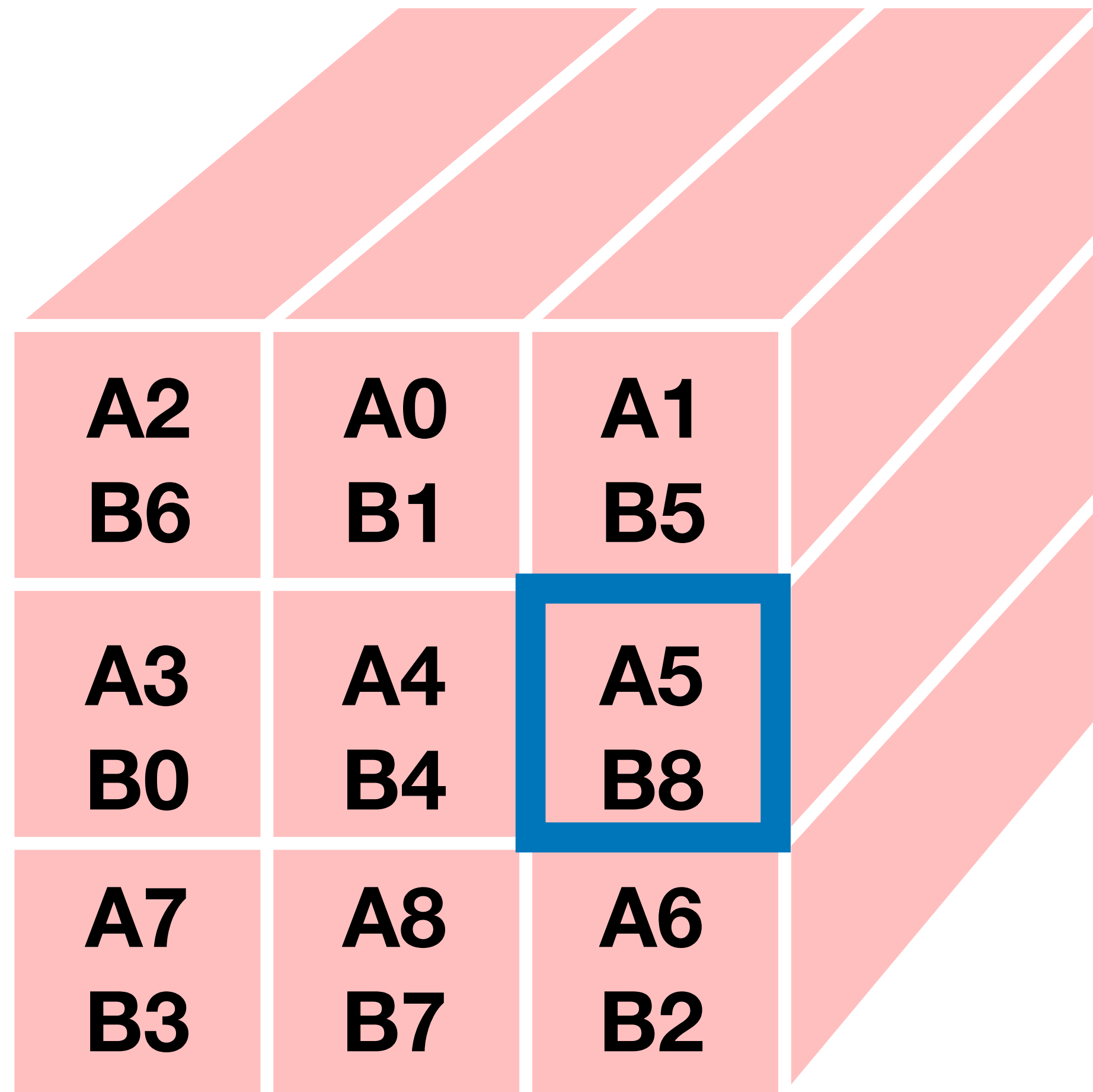
# Cannon's Algorithm



- Need to multiply A0 by :

  - B0, B1, B2

- Multiply A1 by :

  - B1, B4, B7

- C5 = [A3*B2 + A4*B5 + A5*B8]

# Cannon's Algorithm

- Only hold a single partition of A and a single partition of B at any time

$$
\begin{bmatrix}
C0 & C1 & C2 & C3 \\
C4 & C5 & C6 & C7 \\
C8 & C9 & C10 & C11 \\
C12 & C13 & C14 & C15
\end{bmatrix}
=
\begin{bmatrix}
A3 \\
A7 \\
A11 \\
A15
\end{bmatrix}
\begin{bmatrix}
B0 & B1 & B2 & B3 \\
B4 & B5 & B6 & B7 \\
B8 & B9 & B10 & B11 \\
B12 & B13 & B14 & B15
\end{bmatrix}
$$

# Cannon's Algorithm

- First, rotate data

# Cannon's Algorithm

- Then multiply local parts and rotate A to the right and B down

$$
\begin{bmatrix}
C0 & C1 & C2 & C3 \\
C4 & C5 & C6 & C7 \\
C8 & C9 & C10 & C11 \\
C12 & C13 & C14 & C15
\end{bmatrix}
=
\begin{bmatrix}
A3 \\
A7 \\
A11 \\
A15
\end{bmatrix}
\begin{bmatrix}
B0 & B1 & B2 & B3 \\
B4 & B5 & B6 & B7 \\
B8 & B9 & B10 & B11 \\
B12 & B13 & B14 & B15
\end{bmatrix}
$$

# Cannon's Algorithm

- Then multiply local parts and rotate A to the right and B down

$$\begin{bmatrix} C0 & C1 & C2 & C3 \\ C4 & C5 & \boxed{C6} & C7 \\ C8 & C9 & C10 & C11 \\ C12 & C13 & C14 & C15 \end{bmatrix} = \begin{bmatrix} A3 \\ A7 \\ A11 \\ A15 \end{bmatrix} \begin{bmatrix} B0 & B1 & B2 & B3 \\ B4 & B5 & \boxed{B6} & B7 \\ B8 & B9 & B10 & B11 \\ B12 & B13 & B14 & B15 \end{bmatrix}$$

# Cannon's Algorithm

- Then multiply local parts and rotate A to the right and B down

$$
\begin{bmatrix}
C0 & C1 & C2 & C3 \\
C4 & C5 & C6 & C7 \\
C8 & C9 & C10 & C11 \\
C12 & C13 & C14 & C15
\end{bmatrix}
=
\begin{bmatrix}
A3 \\
A7 \\
A11 \\
A15
\end{bmatrix}
\begin{bmatrix}
B0 & B1 & B2 & B3 \\
B4 & B5 & B6 & B7 \\
B8 & B9 & B10 & B11 \\
B12 & B13 & B14 & B15
\end{bmatrix}
$$