

Matrix-vector Multiplication

Michael Quinn

Parallel Programming with MPI and OpenMP

Material do autor

objetivos

- explorar operações MPI
- mostrar formas diferentes de dividir um mesmo programa sequencial

Sequential Algorithm

2	1	0	4
3	2	1	1
4	3	1	2
3	0	2	0

 ×

1
3
4
1

 =

9
14
19
11

matrix decomposition

2	1	0	4
3	2	1	1
4	3	1	2
3	0	2	0

2	1	0	4
3	2	1	1
4	3	1	2
3	0	2	0

2	1	0	4
3	2	1	1
4	3	1	2
3	0	2	0

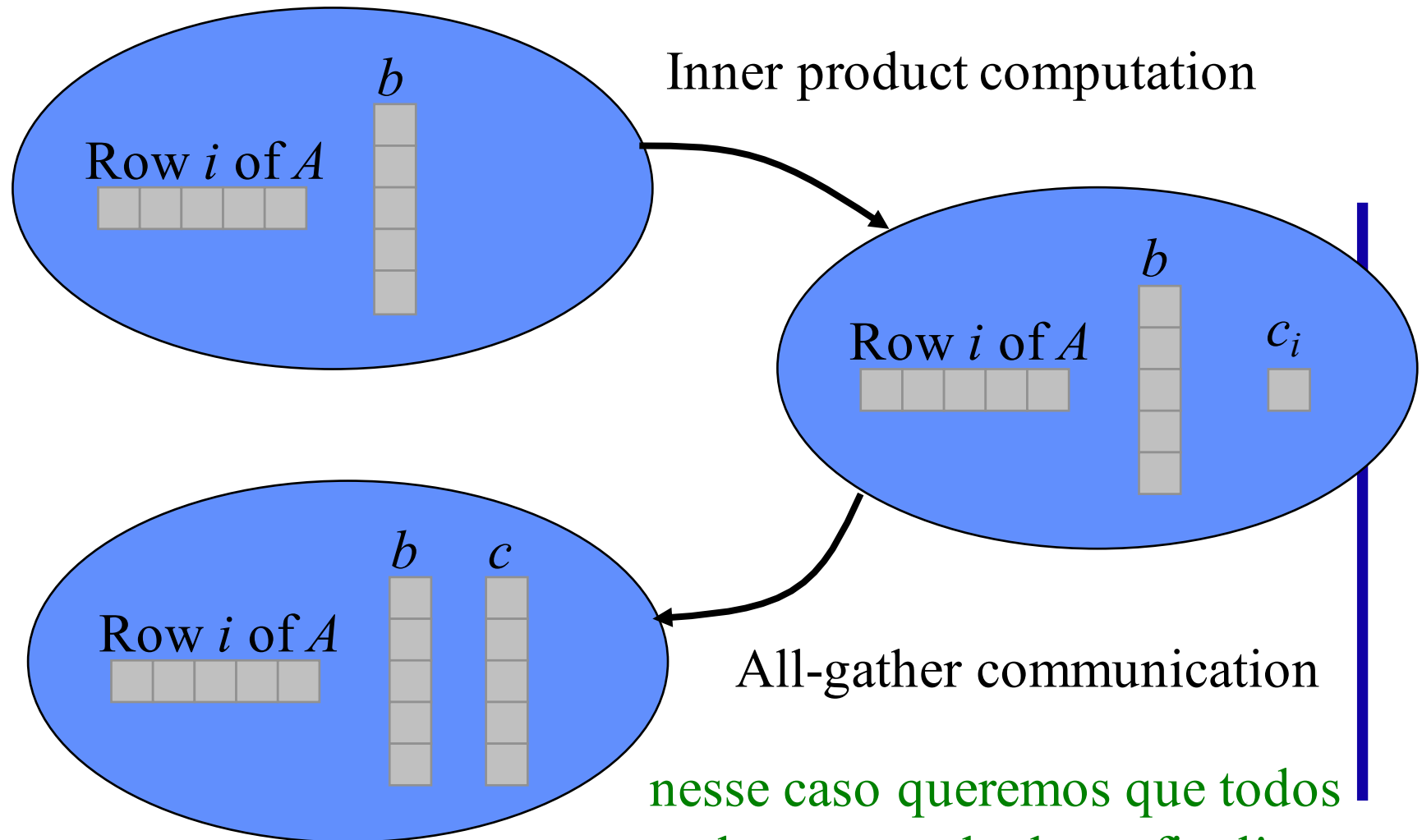
Storing Vectors

- Divide vector elements among processes
- Replicate vector elements
- Vector replication acceptable because vectors have only n elements, versus n^2 elements in matrices

Rowwise Block Striped Matrix

- Partitioning through domain decomposition
- Primitive task associated with
 - Row of matrix
 - Entire vector

Phases of Parallel Algorithm - rowwise



nesse caso queremos que todos tenham o resultado ao final!

Agglomeration and Mapping

- Static number of tasks
- Regular communication pattern (all-gather)
- Computation time per task is constant
- Strategy:
 - Agglomerate groups of rows
 - Create one task per MPI process

com matrizes esparsas poderíamos ter desbalanceamento!

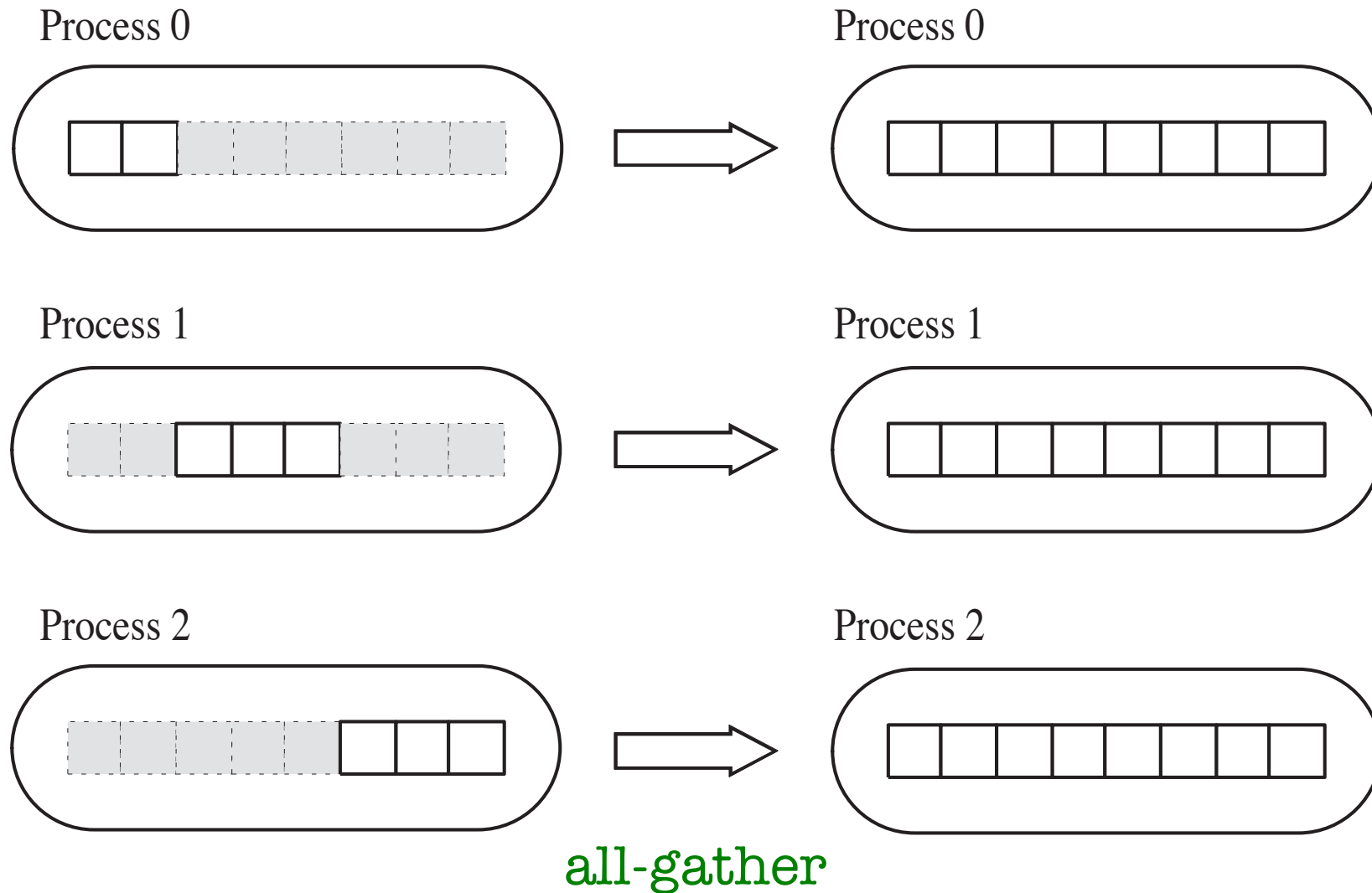
Read matrix/distribute rows



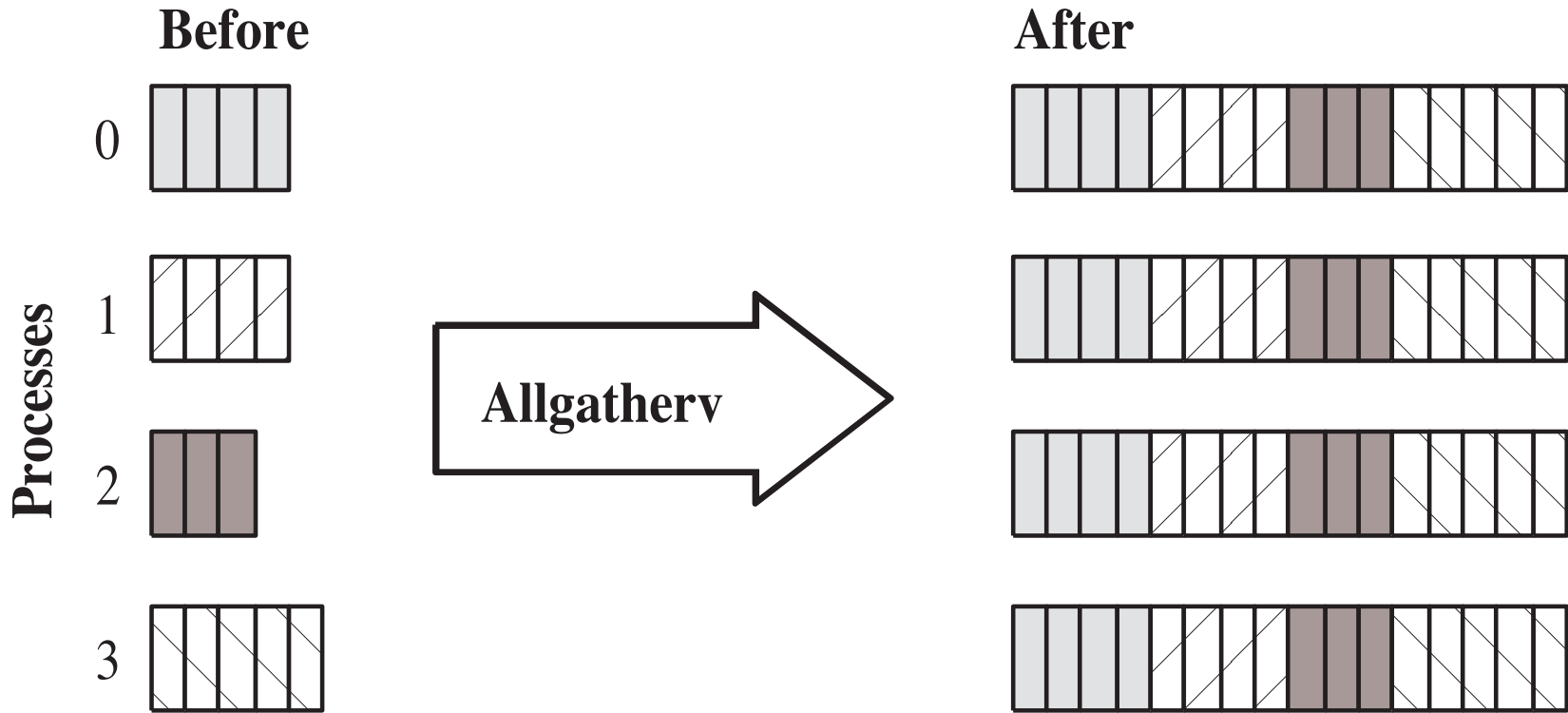
Function `read_replicated_vector`

- Process $p-1$
 - Opens file
 - Reads vector length
- Broadcast vector length (root process = $p-1$)
- Allocate space for vector
- Process $p-1$ reads vector, closes file
- Broadcast vector

Block-to-replicated Transformation - allgather



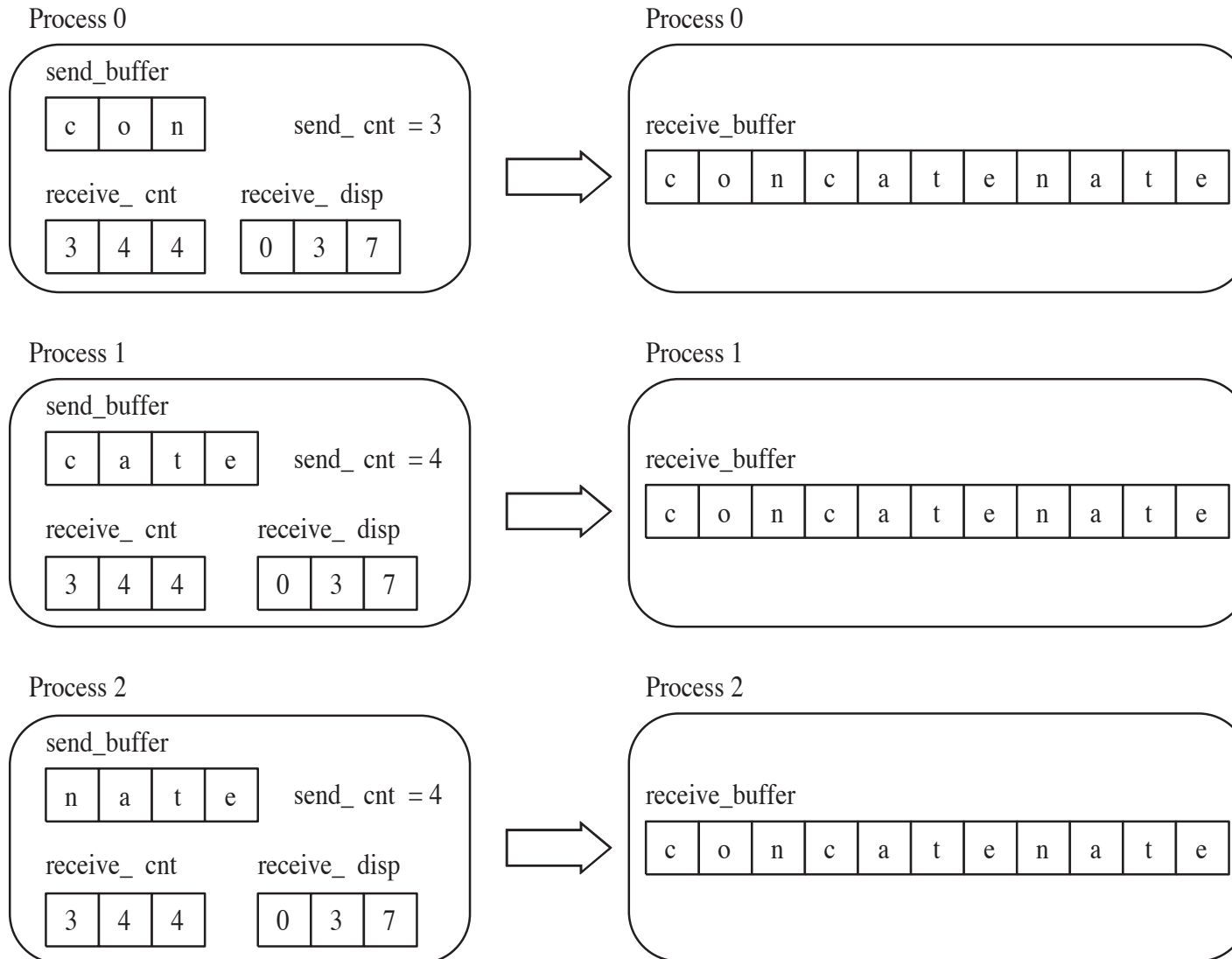
MPI_Allgatherv



MPI_Allgatherv

```
int MPI_Allgatherv (  
    void                *send_buffer,  
    int                 send_cnt, #itens q cada processo  
                                contribui (um inteiro)  
    MPI_Datatype        send_type,  
    void                *receive_buffer,  
    int                 *receive_cnt, #itens q espera de cada  
                                    processo (um array)  
    int                 *receive_disp, índice onde colocar  
                                    dados de cada  
    MPI_Datatype        receive_type, processo  
    MPI_Comm            communicator)
```

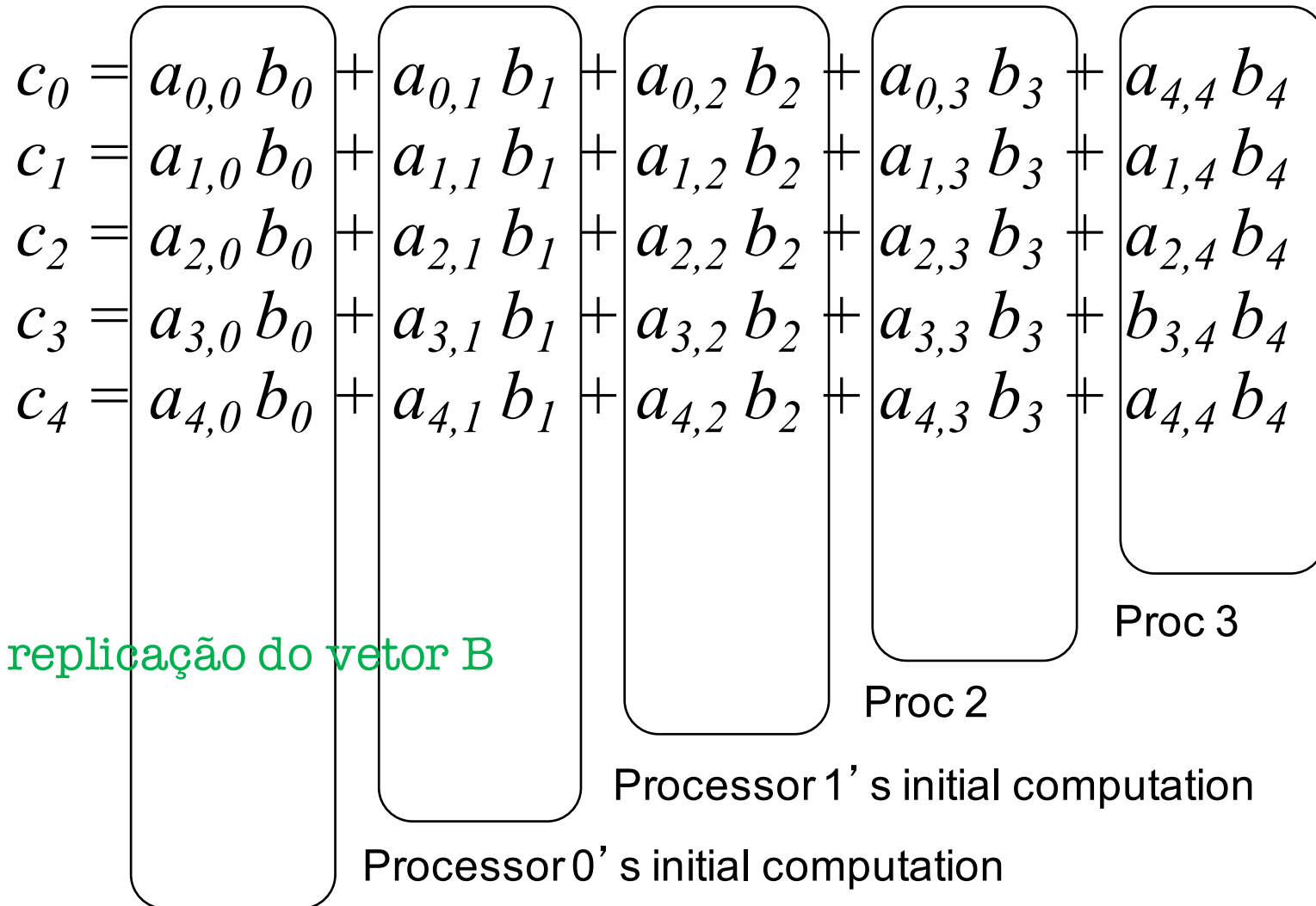
MPI_Allgatherv in Action



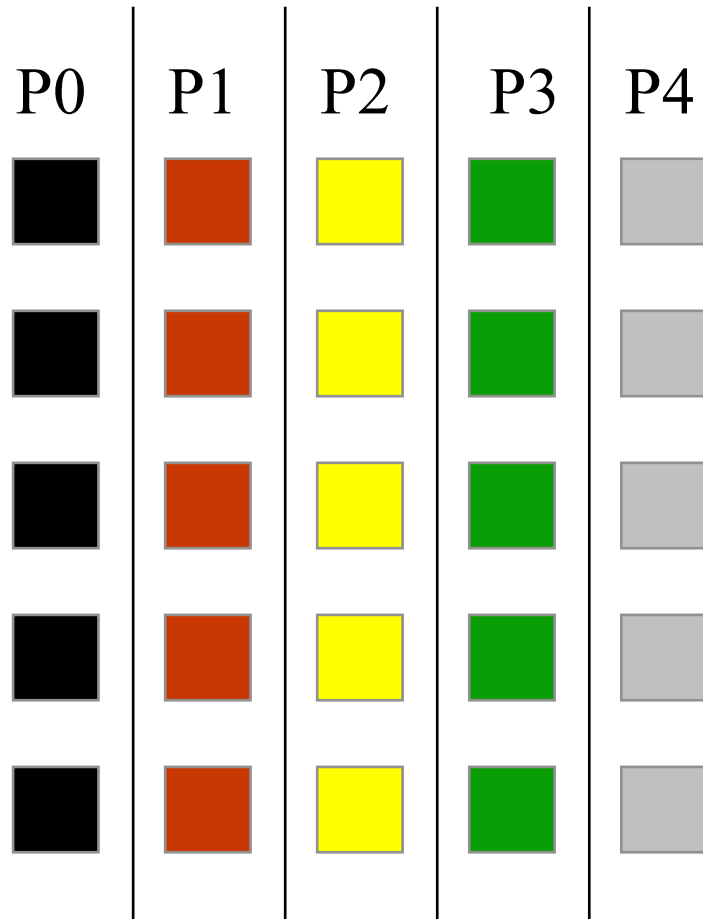
Columnwise Block Striped Matrix

- Partitioning through domain decomposition
- Task associated with
 - Column of matrix
 - Vector element

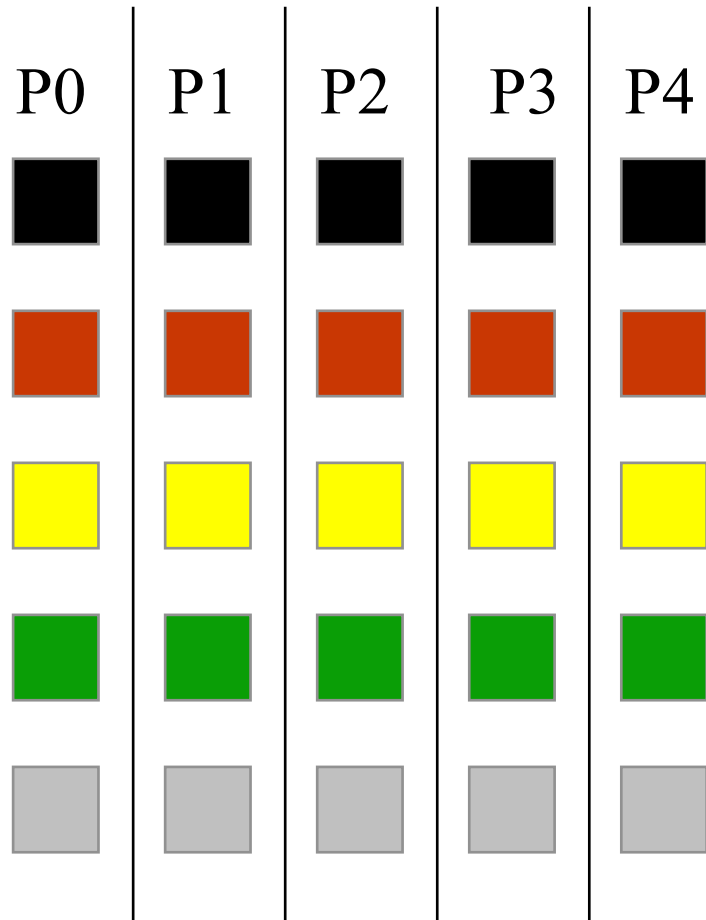
Matrix-Vector Multiplication



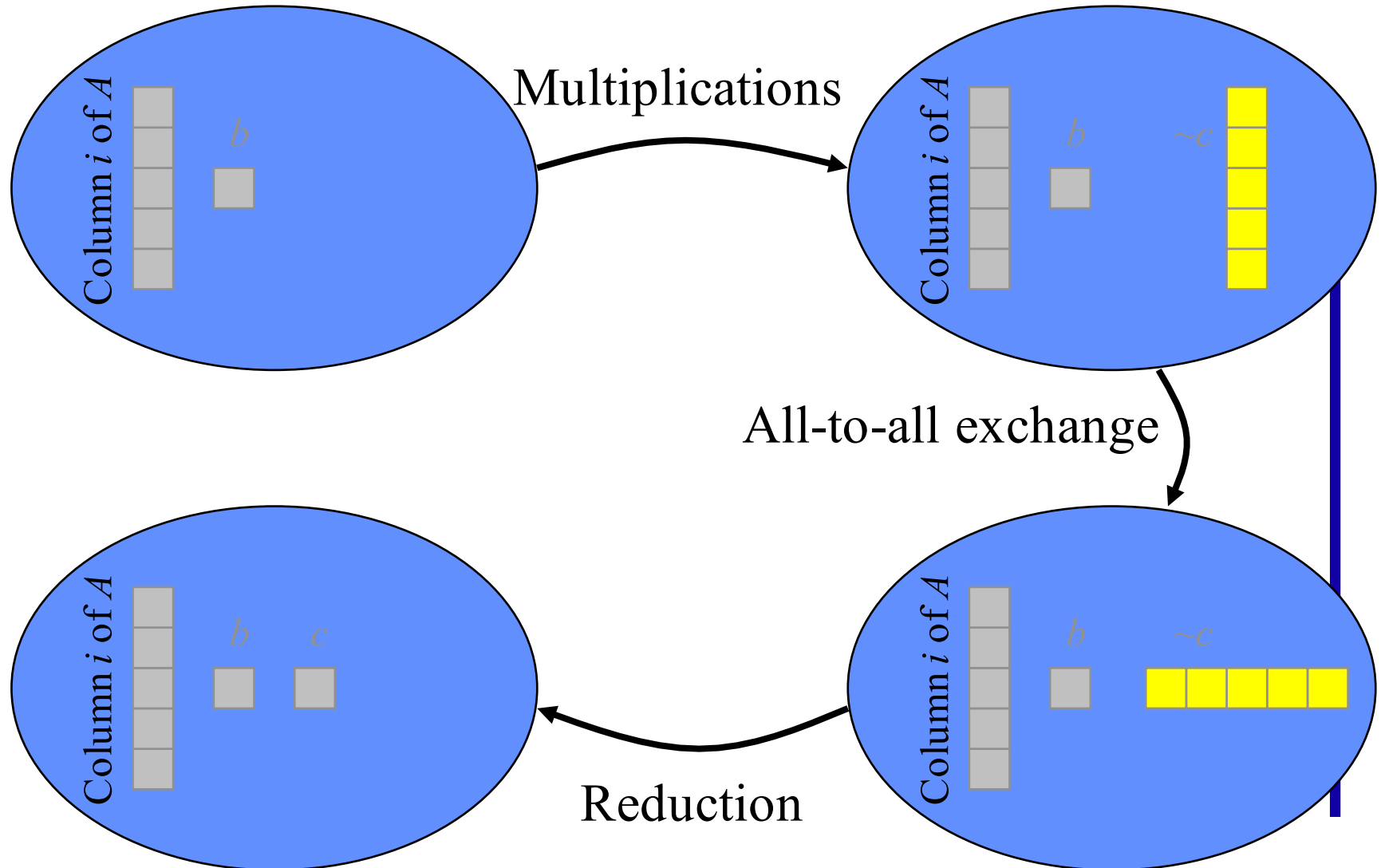
All-to-all Exchange (before)



All-to-all Exchange (after)



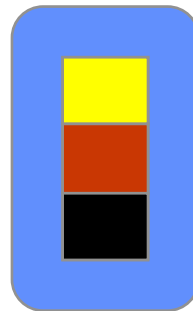
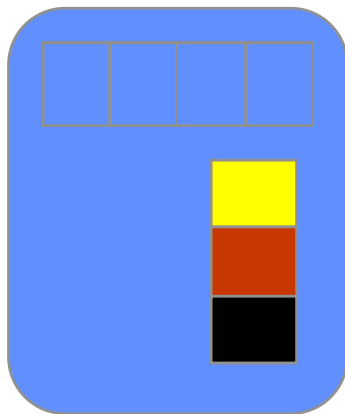
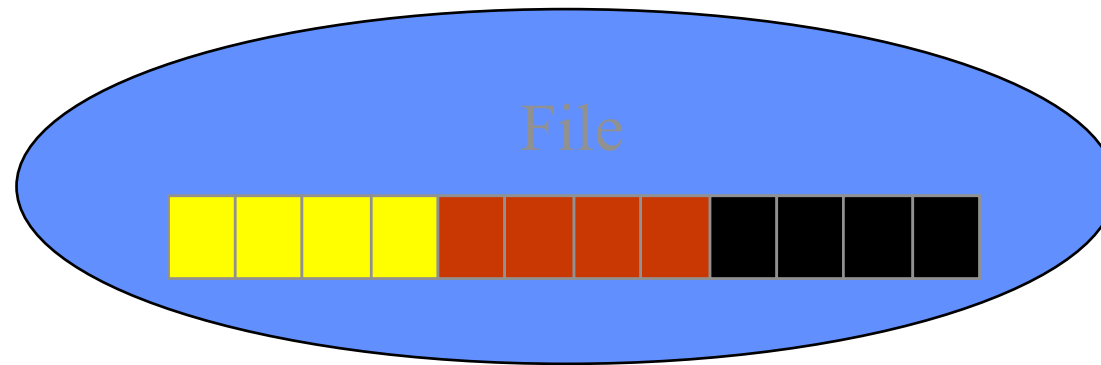
Phases of Parallel Algorithm



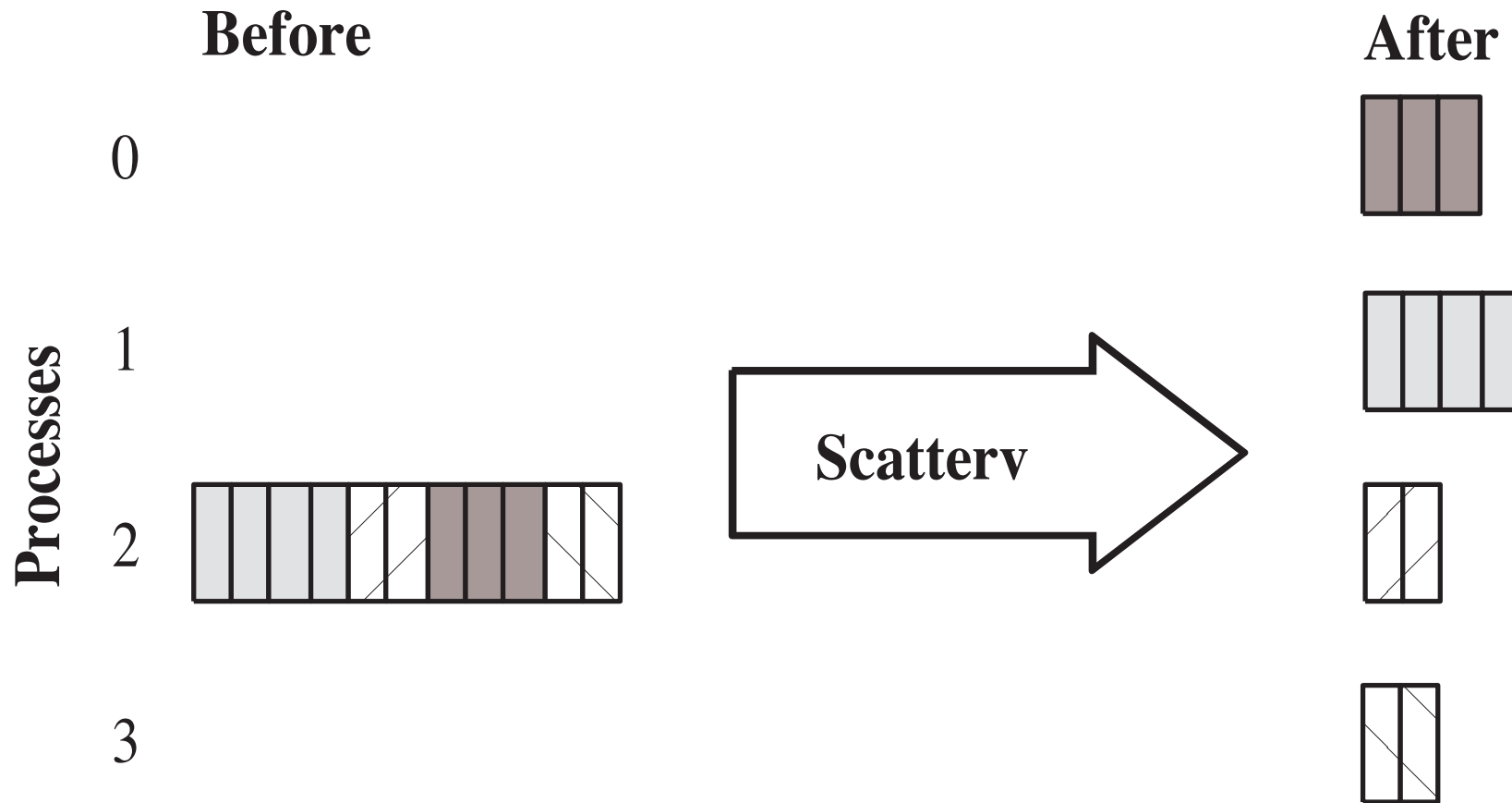
Agglomeration and Mapping

- Static number of tasks
- Regular communication pattern (all-to-all)
- Computation time per task is constant
- Strategy:
 - Agglomerate groups of columns
 - Create one task per MPI process

Reading a Block-Column Matrix



MPI_Scatterv



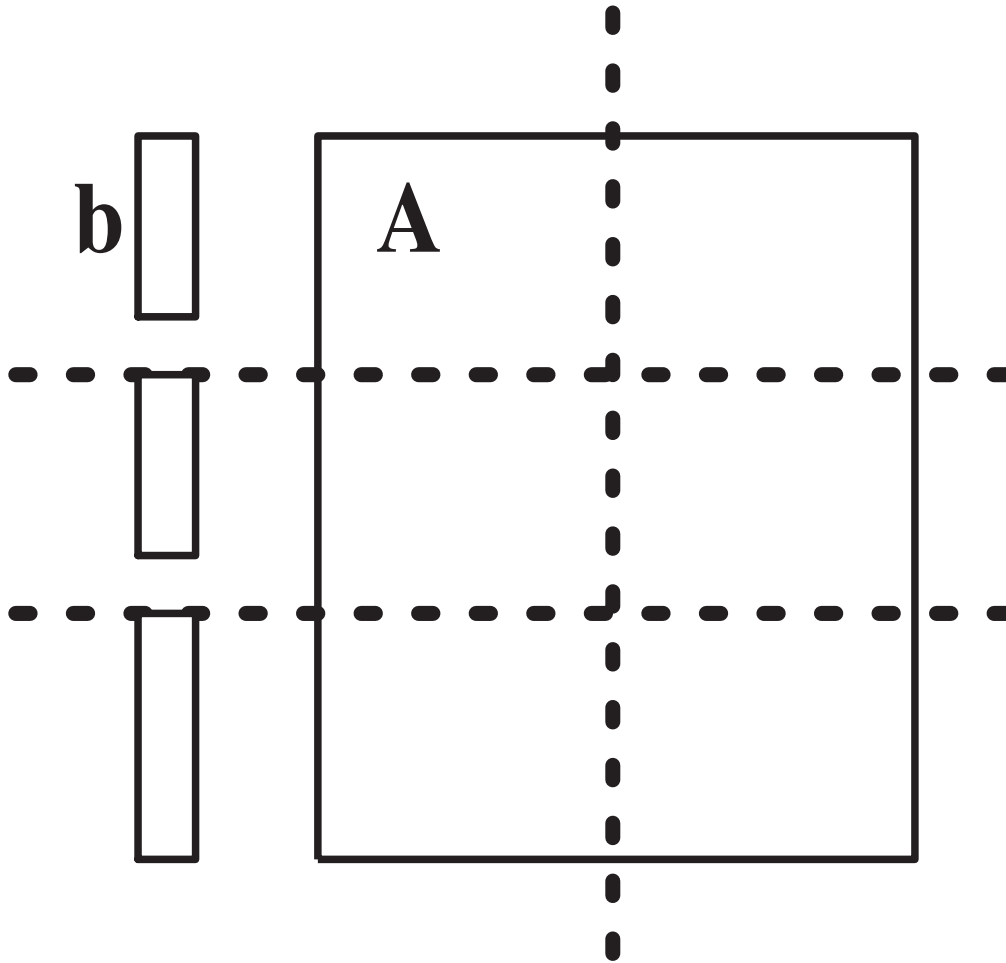
Header for MPI_Scatterv

```
int MPI_Scatterv (  
    void          *send_buffer,  
    int          *send_cnt,  
    int          *send_disp,  
    MPI_Datatype send_type,  
    void          *receive_buffer,  
    int          receive_cnt,  
    MPI_Datatype receive_type,  
    int          root,  
    MPI_Comm     communicator)
```

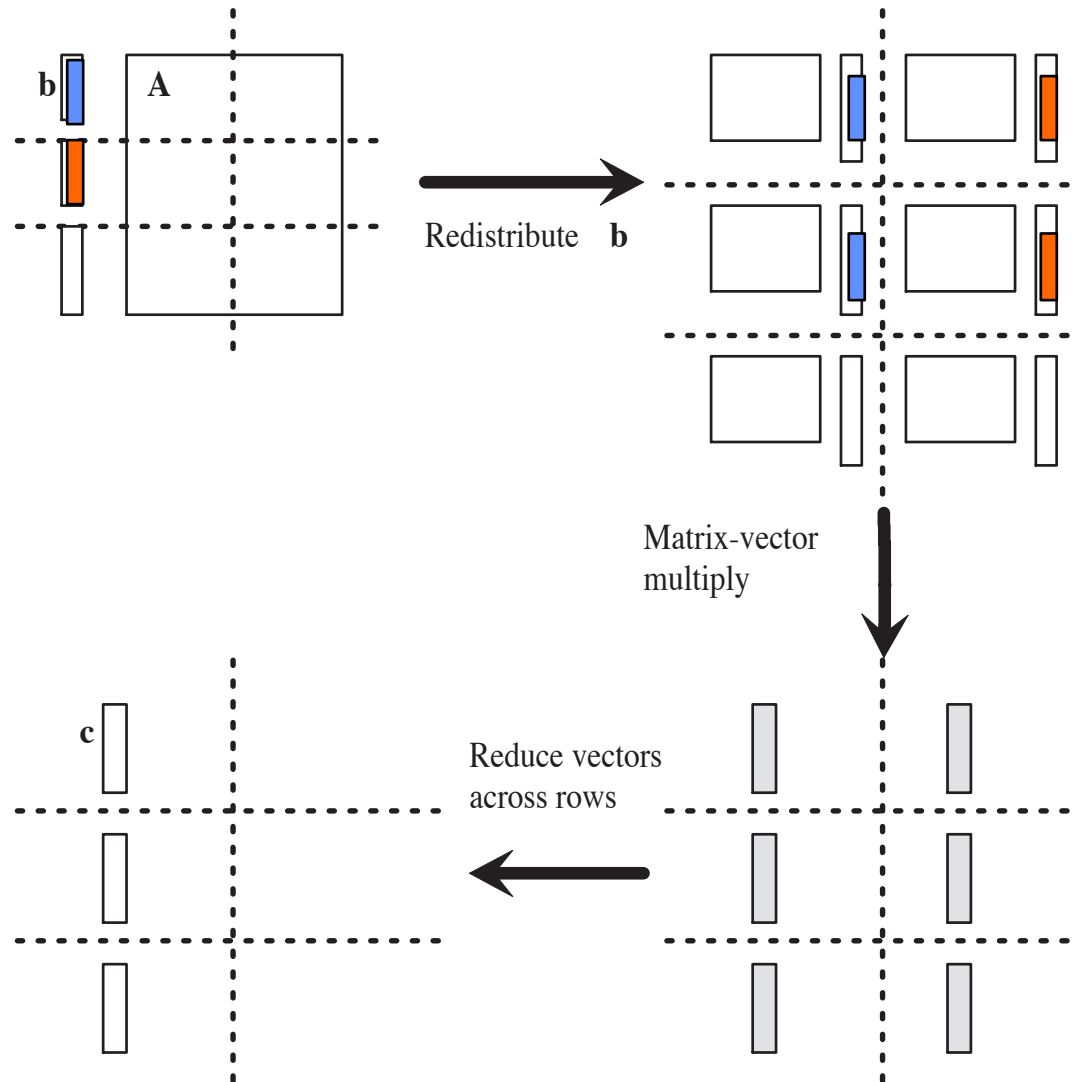
Checkerboard Block Decomposition

- Associate primitive task with each element of the matrix **A**
- Each primitive task performs one multiply
- Agglomerate primitive tasks into rectangular blocks
- Processes form a 2-D grid
- Vector **b** distributed by blocks among processes in first column of grid

Tasks after Agglomeration



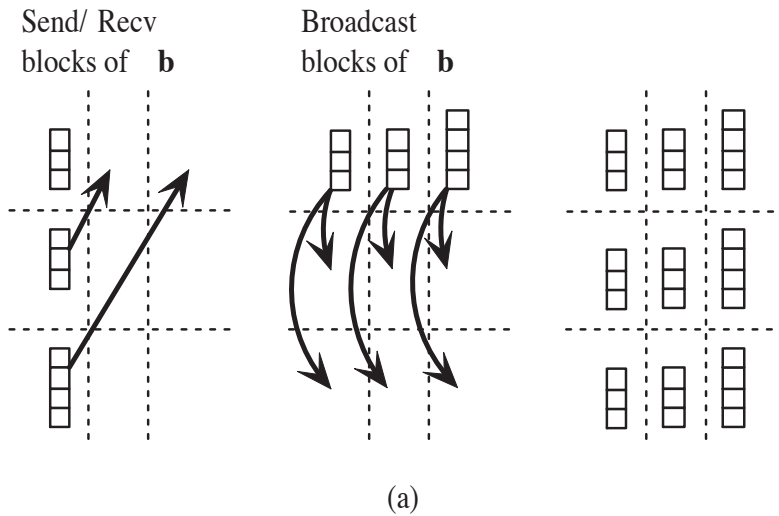
Algorithm's Phases



Redistributing Vector \mathbf{b}

- Step 1: Move \mathbf{b} from processes in first row to processes in first column
 - If p square
 - » First column/first row processes send/receive portions of \mathbf{b}
 - If p not square
 - » Gather \mathbf{b} on process 0, 0
 - » Process 0, 0 broadcasts to first row procs
- Step 2: First row processes scatter \mathbf{b} within columns

Redistributing Vector \mathbf{b}



when p is a square number

Creating Communicators

- Want processes in a virtual 2-D grid
- Create a custom communicator to do this
- Collective communications involve all processes in a communicator
- We need to do broadcasts, reductions among subsets of processes
- We will create communicators for processes in same row or same column

What's in a Communicator?

- Process group
- Context
- Attributes
 - Topology (lets us address processes another way)
 - Others we won't consider

Creating 2-D Virtual Grid of Processes

- MPI_Dims_create
 - Input parameters
 - » Total number of processes in desired grid
 - » Number of grid dimensions
 - Returns number of processes in each dim
- MPI_Cart_create
 - Creates communicator with cartesian topology

MPI_Dims_create

```
int MPI_Dims_create (  
    int nodes,  
        /* Input - Procs in grid */  
  
    int dims,  
        /* Input - Number of dims */  
  
    int *size)  
    /* Input/Output - Size of  
    each grid dimension */
```

- cria um array de dimensões para uso posterior
- size: na entrada pode indicar uma dimensão desejada

```
int p = ...; int size(2) = {0,0};  
MPI_Dims_Create (p, 2, size);
```


MPI_Cart_create

```
int MPI_Cart_create (
    MPI_Comm old_comm, /* Input - old communicator */

    int dims, /* Input - grid dimensions */

    int *size, /* Input - # procs in each dim */

    int *periodic,
        /* Input - periodic[j] is 1 if dimension j
           wraps around; 0 otherwise */

    int reorder,
        /* 1 if process ranks can be reordered */

    MPI_Comm *cart_comm)
    /* Output - new communicator */
```

Using MPI_Dims_create and MPI_Cart_create

```
MPI_Comm cart_comm;
int p;
int periodic[2];
int size[2];
...
size[0] = size[1] = 0;
MPI_Dims_create (p, 2, size);
periodic[0] = periodic[1] = 0;
MPI_Cart_create (MPI_COMM_WORLD, 2, size,
                1, &cart_comm);
```

Useful Grid-related Functions

- `MPI_Cart_rank`
 - Given coordinates of process in Cartesian communicator, returns process rank
- `MPI_Cart_coords`
 - Given rank of process in Cartesian communicator, returns process' coordinates

Header for MPI_Cart_rank

```
int MPI_Cart_rank (  
    MPI_Comm comm,  
    /* In - Communicator */  
    int *coords,  
    /* In - Array containing process'  
        grid location */  
    int *rank)  
    /* Out - Rank of process at  
        specified coords */
```

Header for MPI_Cart_coords

```
int MPI_Cart_coords (
    MPI_Comm comm,
    /* In - Communicator */
    int rank,
    /* In - Rank of process */
    int dims,
    /* In - Dimensions in virtual grid */
    int *coords)
    /* Out - Coordinates of specified
    process in virtual grid */
```

MPI_Comm_split

- Partitions the processes of a communicator into one or more subgroups
- Constructs a communicator for each subgroup
- Allows processes in each subgroup to perform their own collective communications
- Needed for columnwise scatter and rowwise reduce

Header for MPI_Comm_split

```
int MPI_Comm_split (  
    MPI_Comm old_comm,  
    /* In - Existing communicator */  
  
    int partition, /* In - Partition number */  
  
    int new_rank,  
    /* In - Ranking order of processes  
       in new communicator */  
  
    MPI_Comm *new_comm)  
    /* Out - New communicator shared by  
       processes in same partition */
```

Example: Create Communicators for Process Rows

```
MPI_Comm grid_comm; /* 2-D process grid */

MPI_Comm grid_coords[2];
                /* Location of process in grid */

MPI_Comm row_comm;
                /* Processes in same row */

MPI_Comm_split (grid_comm, grid_coords[0],
                grid_coords[1], &row_comm);
```