

# Message Passing Interface (MPI); Collective Communication

EECS 221: Intro to High-Performance Computing

*Aparna Chandramowlishwaran*

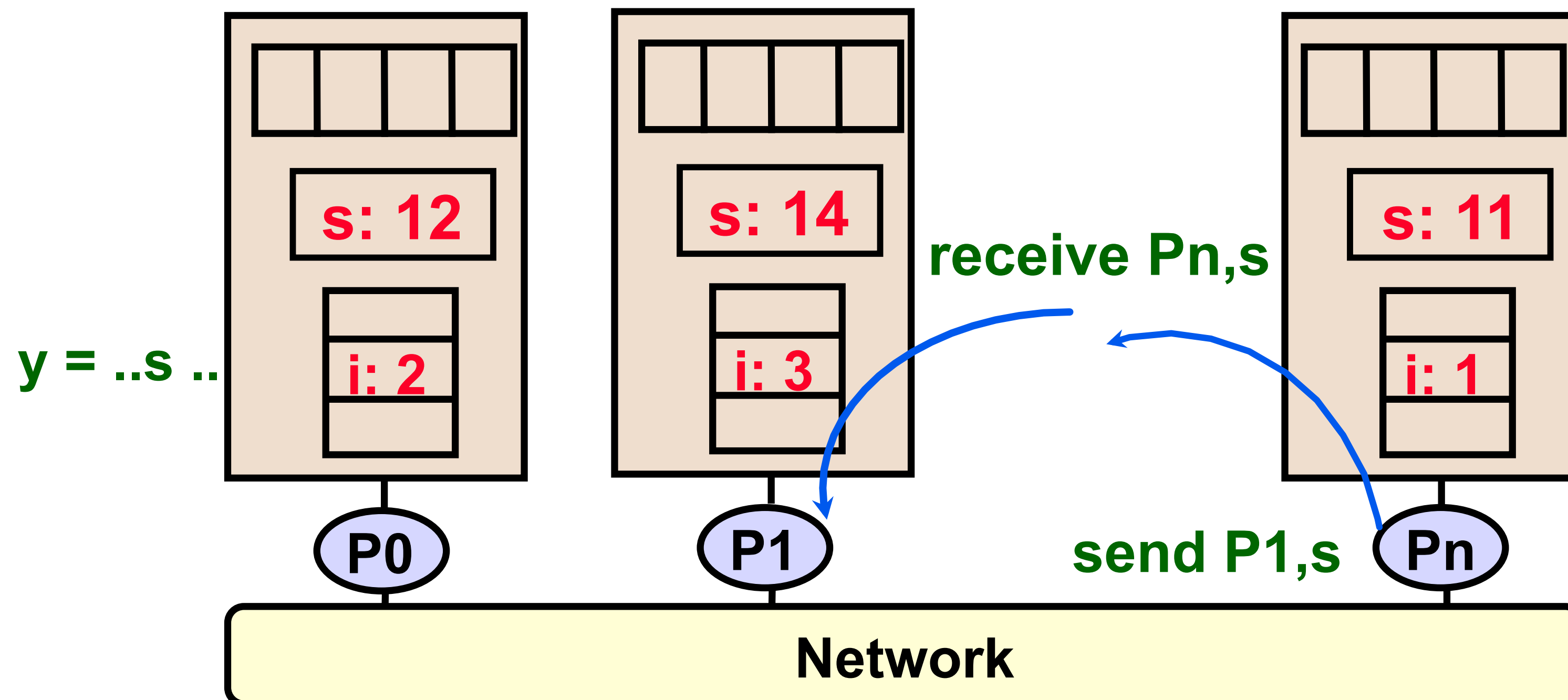
April 27, 2017

# Administrivia

- ▶ Homework 2: **Hands-on lab** on Tuesday May 2<sup>nd</sup> (next class)
- ▶ Readings posted on the class website

# Message Passing Model

- ▶ Program consists of a collection of **named** processes
- ▶ Processes communicate by **explicit send/receive messages**
- ▶ Coordination is implicit in every communication
- ▶ **MPI** has become the de facto standard



# Message Passing Interface

- ▶ All communication, synchronization require subroutine calls
  - ▶ No shared variables
- ▶ Communication primitives
  - ▶ Pairwise, or **point-to-point**: send & receive ([non]blocking, [a]synchronous)
  - ▶ **Collectives**
    - ▶ **Move data**: Broadcast, Scatter/gather
    - ▶ **Compute and move**: sum, product, max, prefix sum, etc,.
- ▶ **Barrier** synchronization
  - ▶ No locks because there are no shared variables to protect
- ▶ **See**: <http://www.mpi-forum.org> and <http://www.mcs.anl.gov/research/projects/mpi>

# MPI environment

- ▶ Two important questions arise in a parallel program
  - ▶ How many processes are participating in this computation?
  - ▶ Which one am I?

# MPI environment

- ▶ Two important questions arise in a parallel program
  - ▶ How many processes are participating in this computation?  
`MPI_Comm_size` returns the number of processors, ***size***
  - ▶ Which one am I? `MPI_Comm_rank` reports the ***rank***, a number between **0** and **size-1**, identifying the calling process

# “Hello World” example in MPI

```
int main(int argc, char *argv[])
{
    int rank, size;

    printf ("I am %d of %d\n", rank, size); // Execute in parallel

    return 0;
}
```

# “Hello World” example in MPI

```
#include <mpi.h>
```

```
int main(int argc, char *argv[])  
{  
    int rank, size;
```

```
    printf (“I am %d of %d\n”, rank, size);
```

```
    return 0;
```

```
}
```



# “Hello World” example in MPI

```
#include <mpi.h>

int main(int argc, char *argv[])
{
    int rank, size;
    MPI_Init (&argc, &argv);

    printf ("I am %d of %d\n", rank, size);
    MPI_Finalize ();
    return 0;
}
```

# “Hello World” example in MPI

```
#include <mpi.h>

int main(int argc, char *argv[])
{
    int rank, size;
    MPI_Init (&argc, &argv);
    MPI_Comm_rank (MPI_COMM_WORLD, &rank);
    MPI_Comm_size (MPI_COMM_WORLD, &size);
    printf ("I am %d of %d\n", rank, size);
    MPI_Finalize ();
    return 0;
}
```

# “Hello World” example in MPI

```
#include <mpi.h>
```

```
int main(int argc, char *argv[])
```

```
{
```

```
    int rank, size;
```

```
    MPI_Init (&argc, &argv);
```

```
    MPI_Comm_rank (MPI_COMM_WORLD, &rank);
```

```
    MPI_Comm_size (MPI_COMM_WORLD, &size);
```

```
    printf (“I am %d of %d\n”, rank, size);
```

```
    MPI_Finalize ();
```

```
    return 0;
```

```
}
```

*Compiling:*

```
$ mpicc -o hello hello.c
```

# “Hello World” example in MPI

```
#include <mpi.h>
```

```
int main(int argc, char *argv[])
```

```
{
```

```
    int rank, size;
```

```
    MPI_Init (&argc, &argv);
```

```
    MPI_Comm_rank (MPI_COMM_WORLD, &rank);
```

```
    MPI_Comm_size (MPI_COMM_WORLD, &size);
```

```
    printf (“I am %d of %d\n”, rank, size);
```

```
    MPI_Finalize ();
```

```
    return 0;
```

```
}
```

*Compiling:*

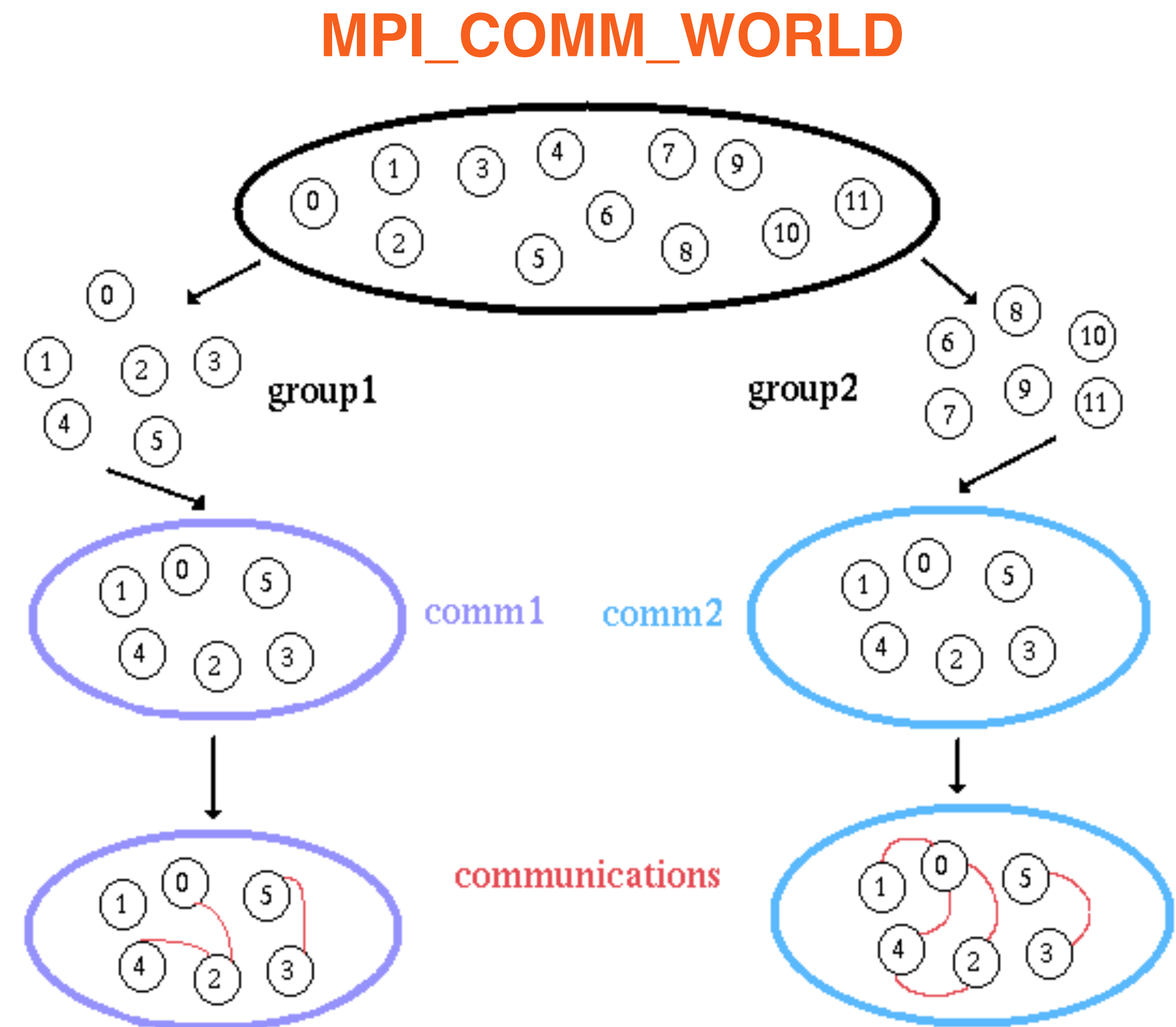
```
$ mpicc -o hello hello.c
```

*Run:*

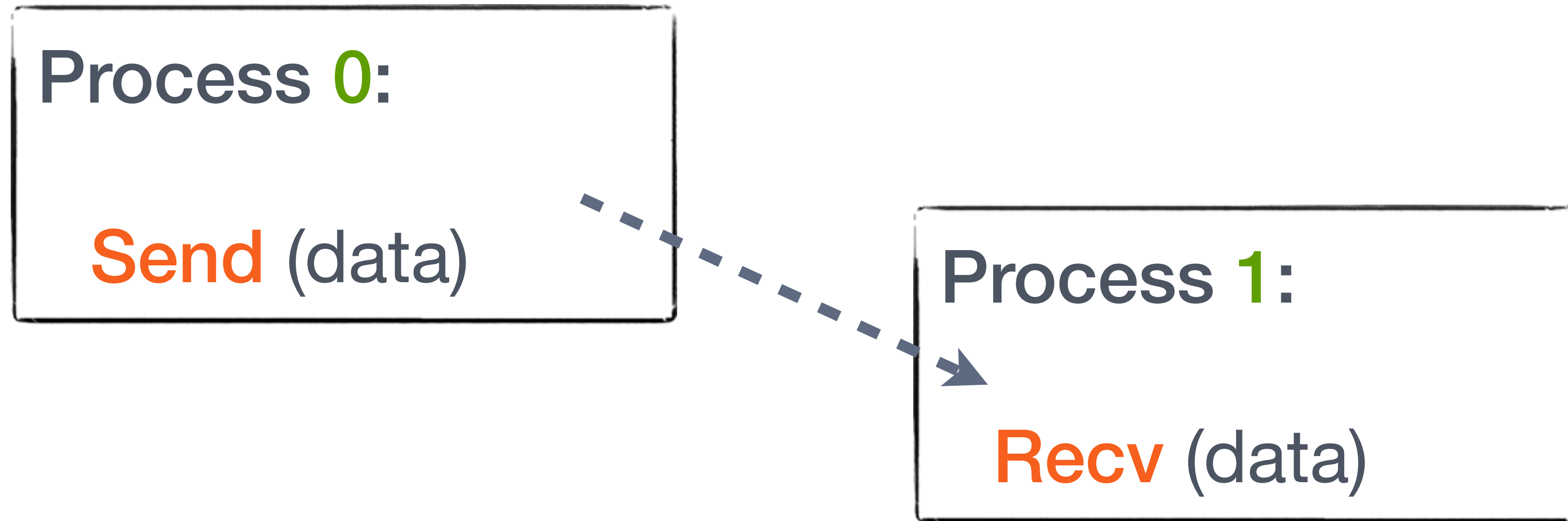
```
$ mpirun -np 4 ./hello
```

# Basic concept: Communicators

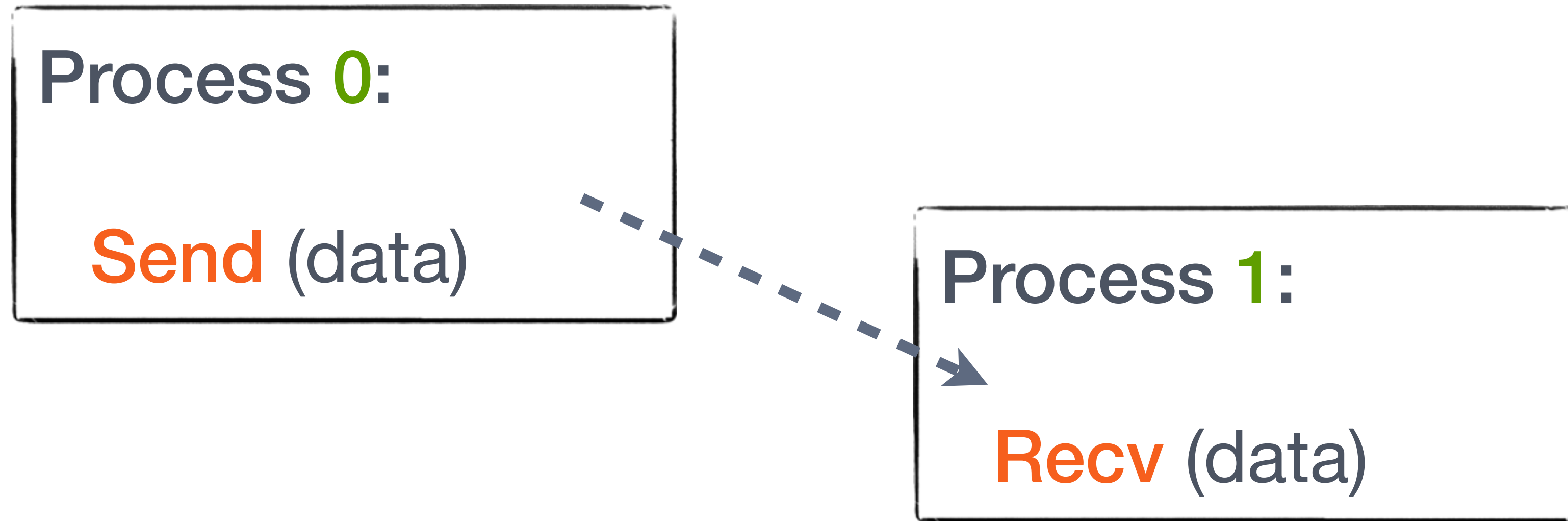
- ▶ **Group**: collection of processors
- ▶ Each message is sent in a **context**, and must be received in the same context
- ▶ **Communicator** = group + context
- ▶ **Rank**: process ID in its communicator
- ▶ **MPI\_COMM\_WORLD** = Group consisting of all processes
- ▶ **MPI\_ANY\_RANK** = Wildcard rank (receive only)



# Basic concepts: Send and receive



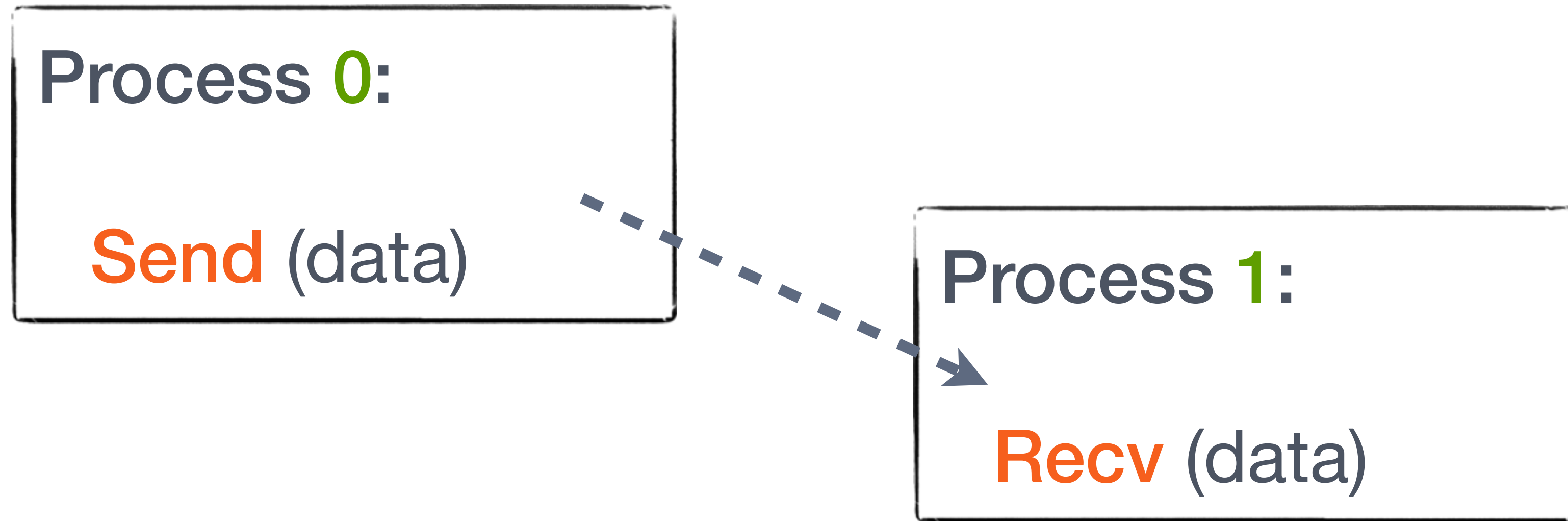
# Basic concepts: Send and receive



- ▶ How to **identify** tasks/processes?
- ▶ How to **describe** “data”?
- ▶ How will **receiver** recognize and screen messages?
- ▶ What does it mean for operations to **complete**?



# Basic concepts: Send and receive



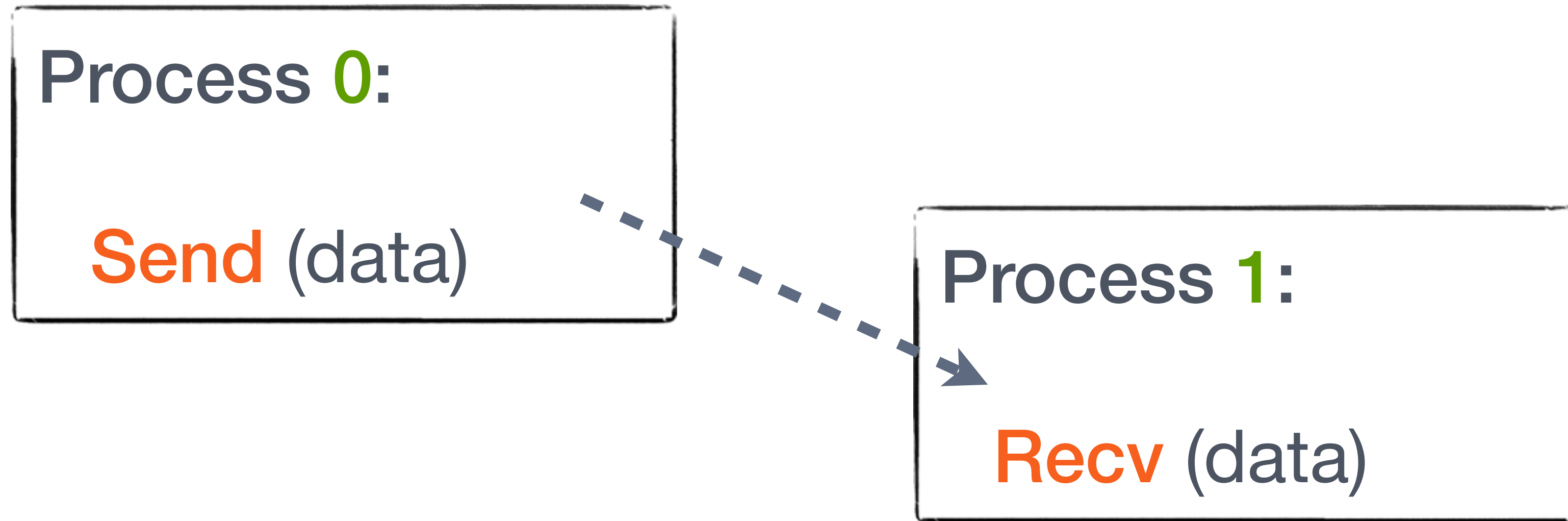
- ▶ How to **identify** tasks/processes?
- ▶ **How to describe “data”?**
- ▶ How will **receiver** recognize and screen messages?
- ▶ What does it mean for operations to **complete**?



# Basic concept: MPI datatypes

- ▶ In MPI call, “data” is described a triple: (**address, count, type**)
- ▶ Datatype = (recursively defined)
  - ▶ “Standard” predefined scalar types: **MPI\_INT, MPI\_DOUBLE, MPI\_CHAR, etc.,**
  - ▶ A contiguous array of datatypes
  - ▶ A strided block of datatypes
  - ▶ An indexed array of blocks of datatypes
  - ▶ An arbitrary structure of datatypes
- ▶ Can construct custom datatypes, in particular ones for subarrays

# Basic concepts: Send and receive



- ▶ How to **identify** tasks/processes?
- ▶ How to **describe** “data”?
- ▶ **How will receiver recognize and screen messages?**
- ▶ What does it mean for operations to **complete**?

# Basic concepts: MPI tags

## ► Message tags

- Every message has a user-defined **integer ID** to assist the receiver in identifying the message
- Messages can be screened at the receiver by specifying a tag
- Wildcard: **MPI\_ANY\_TAG**

# Basic concept: Status object

- Opaque structures that contains more information (e.g. size of the message, error)

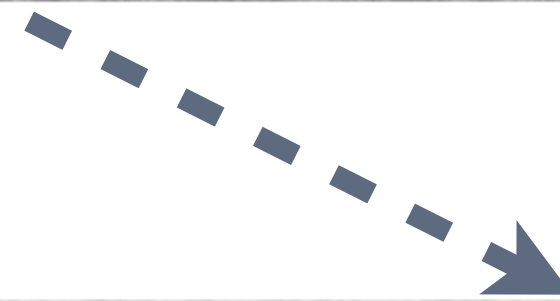
```
int received_tag, received_from, received_count;  
MPI_Status status;  
MPI_Recv (... , MPI_ANY_SOURCE, MPI_ANY_TAG, ..., &status);  
received_tag = status.MPI_TAG;  
received_from = status.MPI_SOURCE;  
MPI_Get_count(status, datatype, &received_count);
```

# Blocking send

**MPI\_Send** (buffer, count, datatype, dest-rank, tag, comm)

Process 0:

```
MPI_Send (A, 10, MPI_INT, 1, tag, comm);
```



Process 1:

```
MPI_Recv (B, 20, MPI_INT, 0, tag, comm, &stat);
```

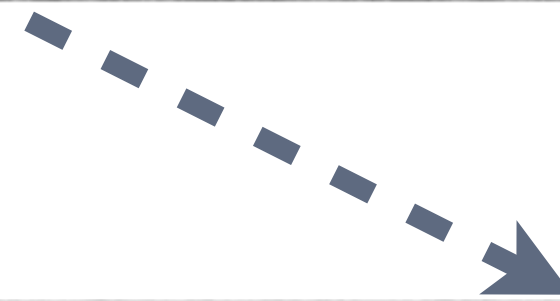
- ▶ Message buffer = (buffer-start, count, datatype)
- ▶ Target = (dest-rank, tag, comm)
- ▶ On return:
  - ▶ Data delivered to the system
  - ▶ Buffer can be reused
  - ▶ Target may not yet have received the message

# Blocking receive

**MPI\_Recv** (buffer, count, datatype, source-rank, tag, comm, status)

Process 0:

```
MPI_Send (A, 10, MPI_INT, 1, tag, comm);
```



Process 1:

```
MPI_Recv (B, 20, MPI_INT, 0, tag, comm, &stat);
```

- ▶ Message buffer = (**buffer-start, count, datatype**) Source = (**source-rank, tag, comm**)
- ▶ Returns when **matching message** (match on source triplet) is received
  - ▶ Source can be a wildcard, **MPI\_ANY\_SOURCE**
  - ▶ Tag can be a wildcard, **MPI\_ANY\_TAG**
  - ▶ Receiving fewer than **count** items is OK, but more is an error
  - ▶ May query **status** for more information (e.g. size of the message, source rank)

# Simple MPI program

```
#include <mpi.h>
int main(int argc, char *argv[])
{
    int rank, buffer;
    MPI_Status status;
    MPI_Init (&argc, &argv);
    MPI_Comm_rank (MPI_COMM_WORLD, &rank);
    if (rank == 0) {
        buffer = 123;
        MPI_Send (&buffer, 1, MPI_INT, 1, 0, MPI_COMM_WORLD);
    }
    else if (rank == 1) {
        MPI_Recv (&buffer, 1, MPI_INT, 0, 0, MPI_COMM_WORLD, &status);
        printf ("Received %d\n", buffer);
    }
    MPI_Finalize ();
    return 0;
}
```



# Beware of deadlock

Process 0:

**Recv** (data  $\leftarrow$  1)

**Send** (data  $\rightarrow$  1)

Process 1:

**Recv** (data  $\leftarrow$  0)

**Send** (data  $\rightarrow$  0)



# Beware of deadlock

Process 0:

**Recv** (data  $\leftarrow$  1)

**Send** (data  $\rightarrow$  1)

Process 1:

**Recv** (data  $\leftarrow$  0)

**Send** (data  $\rightarrow$  0)

Process 0:

**Send** (data  $\rightarrow$  1)

**Recv** (data  $\leftarrow$  1)

Process 1:

**Send** (data  $\rightarrow$  0)

**Recv** (data  $\leftarrow$  0)

# Beware of deadlock

Process 0:

**Recv** (data  $\leftarrow$  1)

**Send** (data  $\rightarrow$  1)

Process 1:

**Recv** (data  $\leftarrow$  0)

**Send** (data  $\rightarrow$  0)

Process 0:

**Send** (data  $\rightarrow$  1)

**Recv** (data  $\leftarrow$  1)

Process 1:

**Send** (data  $\rightarrow$  0)

**Recv** (data  $\leftarrow$  0)

- **Unsafe orderings** of send/recv
- How to avoid?

# Beware of deadlock

- Use **safe orderings** of send/recv

Process **0**:

**Send** (data → **1**)

**Recv** (data ← **1**)

Process **1**:

**Recv** (data ← **0**)

**Send** (data → **0**)

- Use **simultaneous** send/recv

Process **0**:

**SendRecv** (data → **1**)

Process **1**:

**SendRecv** (data → **0**)

**MPI\_SendRecv** (Sbuffer, Scount, Stype, dest, Stag, Rbuffer, Rcount, Rtype, src, comm, &stat)

# Beware of deadlock

- **Unsafe orderings** of send/recv

**Process 0:**

**Send** (data → **1**)

**Recv** (data ← **1**)

**Process 1:**

**Send** (data → **0**)

**Recv** (data ← **0**)

- Use **non-blocking** send/recv

**Process 0:**

**Isend** (data1 → **1**)

**Irecv** (data2 ← **1**)

**Waitall**

**Process 1:**

**Isend** (data1 → **0**)

**Irecv** (data2 ← **0**)

**Waitall**

**MPI\_Isend** (buffer, count, datatype, dest-rank, tag, comm, request)

**MPI\_Waitall** (count, request[ ], status[ ])

# Summary: MPI

- ▶ Many parallel programs can be written using six commonly used MPI primitives
  - ▶ **MPI\_Init**
  - ▶ **MPI\_Finalize**
  - ▶ **MPI\_Comm\_rank**
  - ▶ **MPI\_Comm\_size**
  - ▶ **MPI\_Send**
  - ▶ **MPI\_Recv**
- ▶ Non-blocking primitives for correctness and performance

# Alternate approach

- ▶ **Collective communication**

- ▶ Higher-level communication primitives

- ▶ **MPI\_Bcast**: Broadcast data to all processes

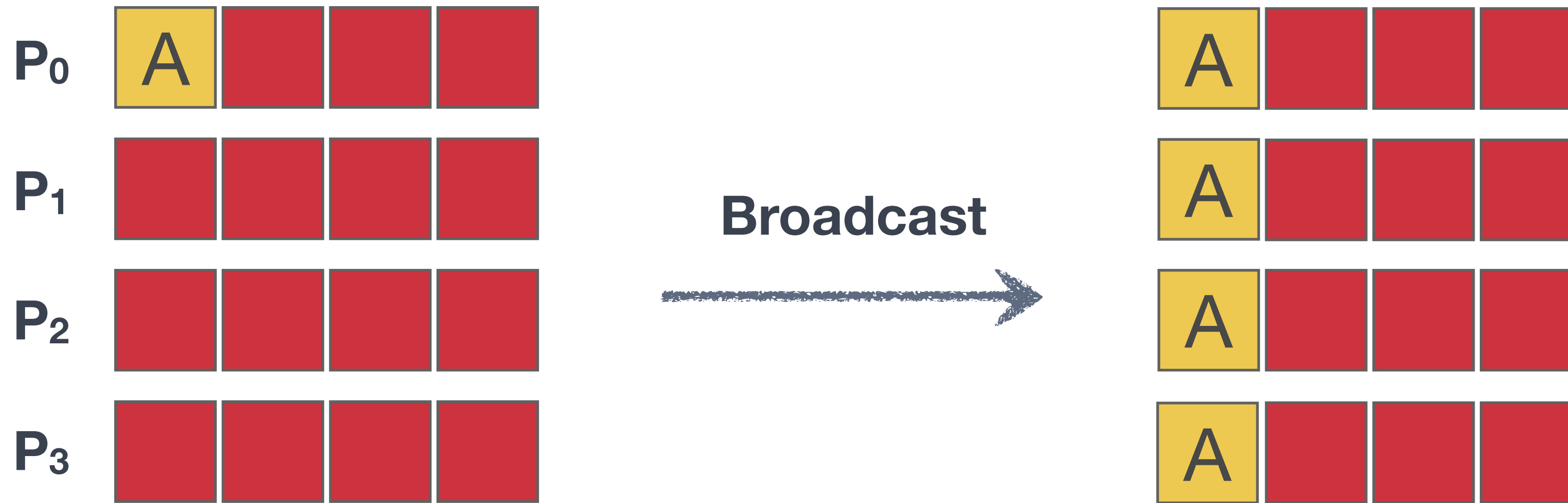
- ▶ **MPI\_Reduce**: Combine data from all processes to one process

- ▶ **MPI\_Barrier**

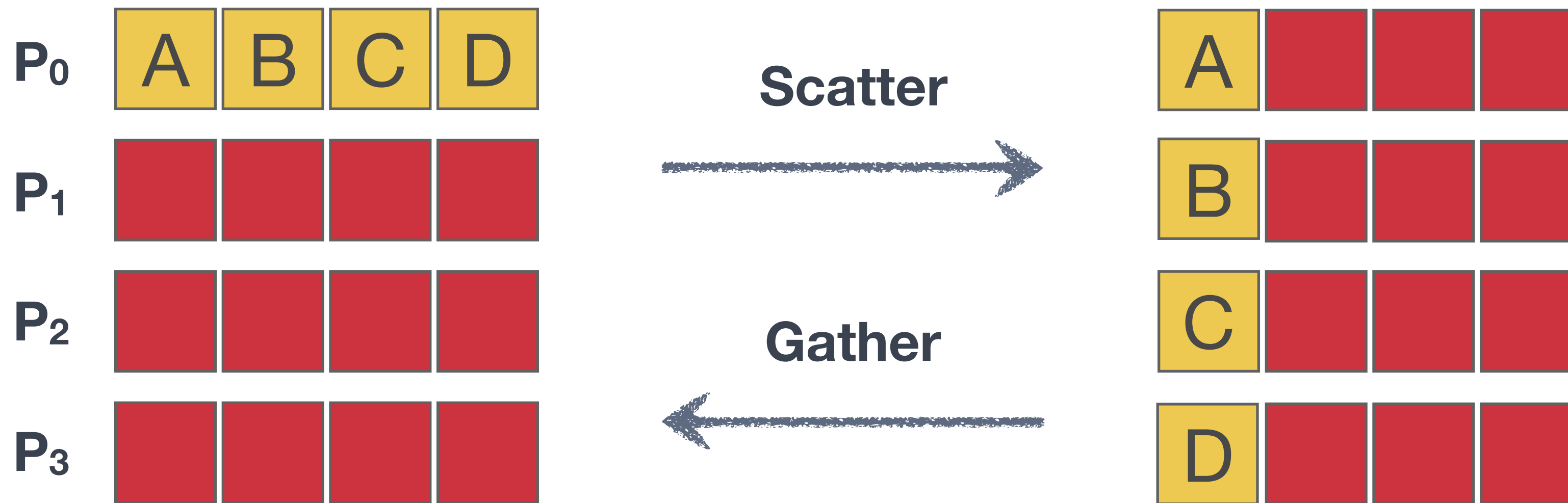
- ▶ Each process executes the **same communication** operation

- ▶ MPI provides a rich set of collective operations

- ▶ Presumably optimized/tuned for hardware



**MPI\_Bcast** (buffer, count, datatype, root, comm)



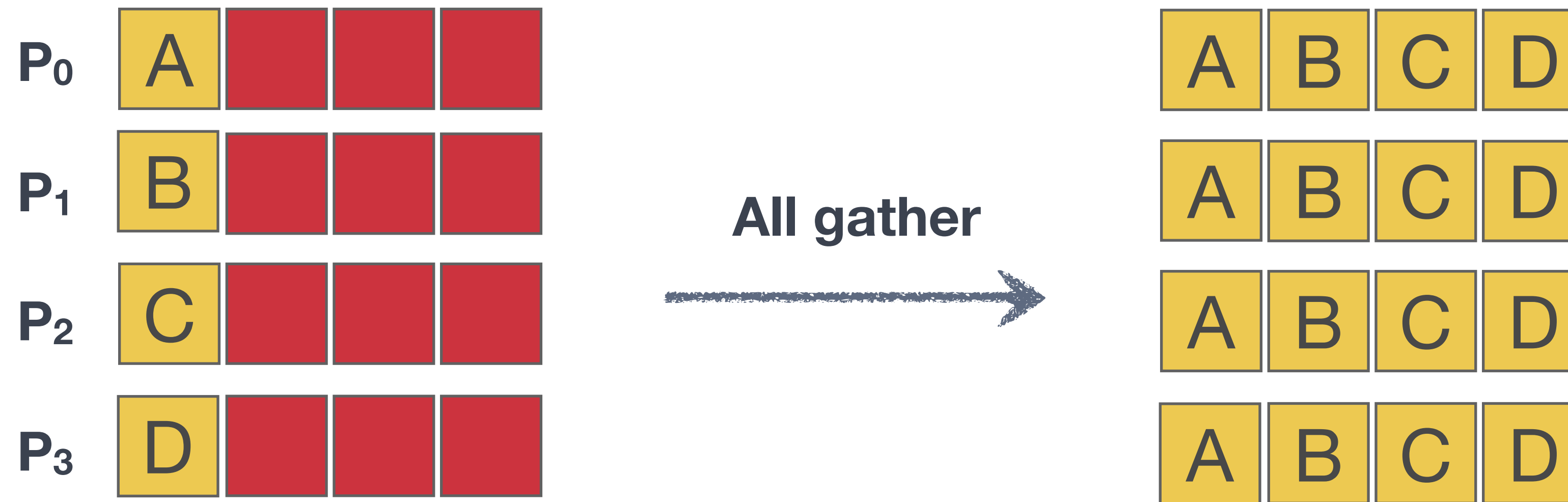
**MPI\_Scatter** (Sbuffer, Scount, Stype, Rbuffer, Rcount, Rtype, root, comm)

**MPI\_Gather** (Sbuffer, Scount, Stype, Rbuffer, Rcount, Rtype, root, comm)

# Case study: Broadcast

- ▶ All collective operations must be called by **all processes** in the communicator
- ▶ Eg,. **MPI\_Bcast** is called by both the root process and all the processes that receive the broadcast
- ▶ Demo example on HPC cluster!

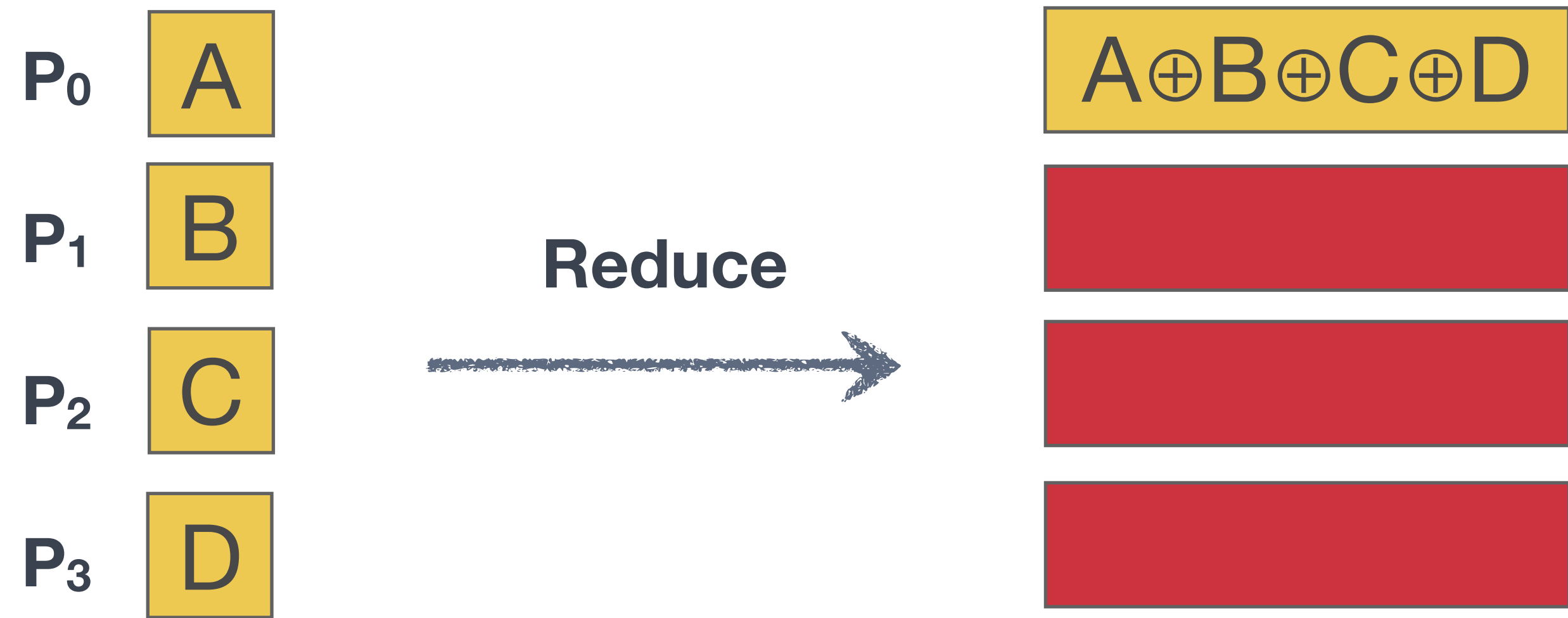




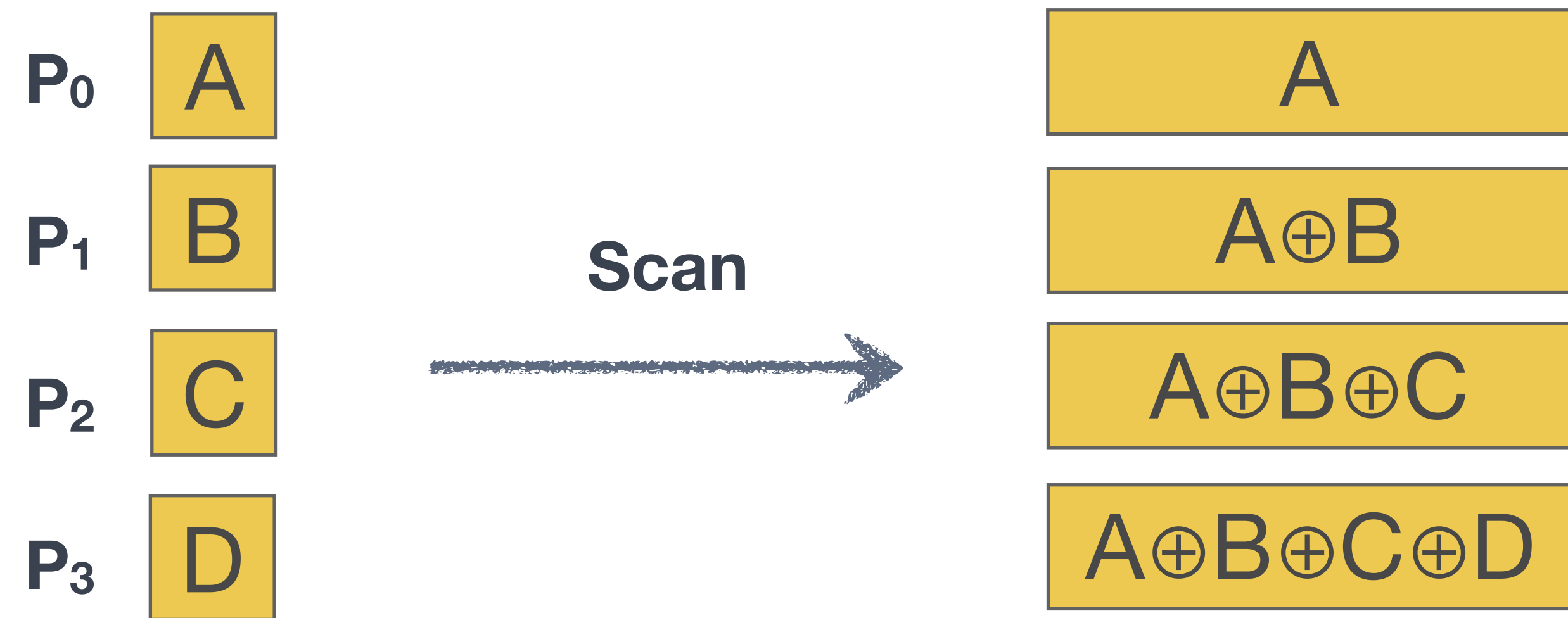
**MPI\_Allgather** (Sbuffer, Scount, Stype, Rbuffer, Rcount, Rtype, comm)



**MPI\_Alltoall** (Sbuffer, Scount, Stype, Rbuffer, Rcount, Rtype, comm)



**MPI\_Reduce** (Sbuffer, Rbuffer, count, datatype, op, root, comm)



**MPI\_Scan** (Sbuffer, Rbuffer, count, datatype, op, comm)

# MPI Built-in Operations

| MPI Reduction Operation |                        |
|-------------------------|------------------------|
| <b>MPI_MAX</b>          | maximum                |
| <b>MPI_MIN</b>          | minimum                |
| <b>MPI_SUM</b>          | sum                    |
| <b>MPI_PROD</b>         | product                |
| <b>MPI_LAND</b>         | logical AND            |
| <b>MPI_BAND</b>         | bit-wise AND           |
| <b>MPI_LOR</b>          | logical OR             |
| <b>MPI_BOR</b>          | bit-wise OR            |
| <b>MPI_LXOR</b>         | logical XOR            |
| <b>MPI_BXOR</b>         | bit-wise XOR           |
| <b>MPI_MAXLOC</b>       | max value and location |
| <b>MPI_MINLOC</b>       | min value and location |

# Other collective routines

- ▶ Many routines: Gather, Gatherv, Allgather, Allgatherv, Reduce\_scatter, Reduce, Allreduce, Scatter, Scatterv, Bcast, Scan
- ▶ **All** = data to all processes
- ▶ **v** = allows data chunks to have variable sizes
- ▶ MPI-2 added **Alltoallw**, **Exscan**, inter-communicator variants for most routines