

Parallel Programming

Prof. **Paolo Bientinesi**

`pauldj@aices.rwth-aachen.de`

WS 16/17



Collective Communication

```
Barrier
Broadcast ↔ Reduce
Scatter ↔ Gather
Allgather ↔ Reduce-scatter
Allreduce
Alltoall
⋮
```

```
int MPI_BCast(...)
```

Before:

| Node ₀ | Node ₁ | Node ₂ | Node ₃ |
|-------------------|-------------------|-------------------|-------------------|
| | | α | |

After:

| Node ₀ | Node ₁ | Node ₂ | Node ₃ |
|-------------------|-------------------|-------------------|-------------------|
| α | α | α | α |

```
int MPI_BCast(...)
```

| Before: | Node ₀ | Node ₁ | Node ₂ | Node ₃ |
|---------|-------------------|-------------------|-------------------|-------------------|
| | | | α | |

| After: | Node ₀ | Node ₁ | Node ₂ | Node ₃ |
|--------|-------------------|-------------------|-------------------|-------------------|
| | α | α | α | α |

```
MPI_Bcast( Buffer, count, type, root, communicator );
```

```
int MPI_Reduce(...)
```

| Before: | Node ₀ | Node ₁ | Node ₂ | Node ₃ |
|---------|-------------------|-------------------|-------------------|-------------------|
| | δ_0 | δ_1 | δ_2 | δ_3 |

| After: | Node ₀ | Node ₁ | Node ₂ | Node ₃ |
|--------|-------------------|-------------------|----------------------------------|-------------------|
| | | | $\bigoplus_{i=0}^{p-1} \delta_i$ | |

`int MPI_Reduce(...)`

| Before: | Node ₀ | Node ₁ | Node ₂ | Node ₃ |
|---------|-------------------|-------------------|-------------------|-------------------|
| | δ_0 | δ_1 | δ_2 | δ_3 |

| After: | Node ₀ | Node ₁ | Node ₂ | Node ₃ |
|--------|-------------------|-------------------|----------------------------------|-------------------|
| | | | $\bigoplus_{i=0}^{p-1} \delta_i$ | |

```
MPI_Reduce( sendBuffer, recvBuffer, count, type,
            operation, root, communicator );
```

`int MPI_Reduce(...)`

| Before: | Node ₀ | Node ₁ | Node ₂ | Node ₃ |
|---------|-------------------|-------------------|-------------------|-------------------|
| | δ_0 | δ_1 | δ_2 | δ_3 |

| After: | Node ₀ | Node ₁ | Node ₂ | Node ₃ |
|--------|-------------------|-------------------|----------------------------------|-------------------|
| | | | $\bigoplus_{i=0}^{p-1} \delta_i$ | |

```
MPI_Reduce( sendBuffer, recvBuffer, count, type,  
            operation, root, communicator );
```

```
MPI_Op_create
```

```
int MPI_Scatter(...)
```

Before:

| Node ₀ | Node ₁ | Node ₂ | Node ₃ |
|-------------------|-------------------|------------------------------|-------------------|
| | | v[0] v[1] v[2] v[3] | |

After:

| Node ₀ | Node ₁ | Node ₂ | Node ₃ |
|-------------------|-------------------|-------------------|-------------------|
| v[0] | v[1] | v[2] | v[3] |


```
int MPI_Scatter(...)
```

| Before: | Node ₀ | Node ₁ | Node ₂ | Node ₃ |
|---------|-------------------|-------------------|-------------------|-------------------|
| | | | v[0] | |
| | | | v[1] | |
| | | | v[2] | |
| | | | v[3] | |

| After: | Node ₀ | Node ₁ | Node ₂ | Node ₃ |
|--------|-------------------|-------------------|-------------------|-------------------|
| | v[0] | | | |
| | | v[1] | | |
| | | | v[2] | |
| | | | | v[3] |

```
MPI_Scatter( sendBuffer, sendCount, sendType,  
             recvBuffer, recvCount, recvType,  
             root, communicator );
```

```
int MPI_Gather(...)
```

| Before: | Node ₀ | Node ₁ | Node ₂ | Node ₃ |
|---------|-------------------|-------------------|-------------------|-------------------|
| | v[0] | v[1] | v[2] | v[3] |

| After: | Node ₀ | Node ₁ | Node ₂ | Node ₃ |
|--------|-------------------|-------------------|------------------------------|-------------------|
| | | | v[0] v[1] v[2] v[3] | |

```
int MPI_Gather(...)
```

| Before: | Node ₀ | Node ₁ | Node ₂ | Node ₃ |
|---------|-------------------|-------------------|-------------------|-------------------|
| | v[0] | v[1] | v[2] | v[3] |

| After: | Node ₀ | Node ₁ | Node ₂ | Node ₃ |
|--------|-------------------|-------------------|------------------------------|-------------------|
| | | | v[0] v[1] v[2] v[3] | |

```
MPI_Gather( sendBuffer, sendCount, sendType,  
            recvBuffer, recvCount, recvType,  
            root, communicator );
```

```
int MPI_Allgather(...)
```

| | Node ₀ | Node ₁ | Node ₂ | Node ₃ |
|---------|-------------------|-------------------|-------------------|-------------------|
| Before: | v[0] | v[1] | v[2] | v[3] |

| | Node ₀ | Node ₁ | Node ₂ | Node ₃ |
|--------|-------------------|-------------------|-------------------|-------------------|
| After: | v[0] | v[0] | v[0] | v[0] |
| | v[1] | v[1] | v[1] | v[1] |
| | v[2] | v[2] | v[2] | v[2] |
| | v[3] | v[3] | v[3] | v[3] |

```
int MPI_Allgather(...)
```

| | Node ₀ | Node ₁ | Node ₂ | Node ₃ |
|---------|-------------------|-------------------|-------------------|-------------------|
| Before: | v[0] | v[1] | v[2] | v[3] |

| | Node ₀ | Node ₁ | Node ₂ | Node ₃ |
|--------|-------------------|-------------------|-------------------|-------------------|
| After: | v[0] | v[0] | v[0] | v[0] |
| | v[1] | v[1] | v[1] | v[1] |
| | v[2] | v[2] | v[2] | v[2] |
| | v[3] | v[3] | v[3] | v[3] |

```
MPI_Allgather( sendBuffer, sendCount, sendType,  
               recvBuffer, recvCount, recvType,  
               communicator );
```

```
int MPI_Reduce_scatter(...)
```

| | Node ₀ | Node ₁ | Node ₂ | Node ₃ |
|---------|--------------------|--------------------|--------------------|--------------------|
| Before: | v ₀ [0] | v ₁ [0] | v ₂ [0] | v ₃ [0] |
| | v ₀ [1] | v ₁ [1] | v ₂ [1] | v ₃ [1] |
| | v ₀ [2] | v ₁ [2] | v ₂ [2] | v ₃ [2] |
| | v ₀ [3] | v ₁ [3] | v ₂ [3] | v ₃ [3] |

| | Node ₀ | Node ₁ | Node ₂ | Node ₃ |
|--------|---------------------------------------|---------------------------------------|---------------------------------------|---------------------------------------|
| After: | Op v _i [0] _i | | | |
| | | Op v _i [1] _i | | |
| | | | Op v _i [2] _i | |
| | | | | Op v _i [3] _i |

```
int MPI_Reduce_scatter(...)
```

| | Node ₀ | Node ₁ | Node ₂ | Node ₃ |
|---------|--------------------|--------------------|--------------------|--------------------|
| Before: | v ₀ [0] | v ₁ [0] | v ₂ [0] | v ₃ [0] |
| | v ₀ [1] | v ₁ [1] | v ₂ [1] | v ₃ [1] |
| | v ₀ [2] | v ₁ [2] | v ₂ [2] | v ₃ [2] |
| | v ₀ [3] | v ₁ [3] | v ₂ [3] | v ₃ [3] |

| | Node ₀ | Node ₁ | Node ₂ | Node ₃ |
|--------|---------------------------------------|---------------------------------------|---------------------------------------|---------------------------------------|
| After: | Op v _i [0] _i | | | |
| | | Op v _i [1] _i | | |
| | | | Op v _i [2] _i | |
| | | | | Op v _i [3] _i |

```
MPI_Reduce_scatter( sendBuffer,  recvBuffer,
                    recvCounts[], type, operation,
                    communicator );
```

```
int MPI_Allreduce(...)
```

| Before: | Node ₀ | Node ₁ | Node ₂ | Node ₃ |
|---------|-------------------|-------------------|-------------------|-------------------|
| | δ_0 | δ_1 | δ_2 | δ_3 |

| After: | Node ₀ | Node ₁ | Node ₂ | Node ₃ |
|--------|----------------------------------|----------------------------------|----------------------------------|----------------------------------|
| | $\bigoplus_{i=0}^{p-1} \delta_i$ | $\bigoplus_{i=0}^{p-1} \delta_i$ | $\bigoplus_{i=0}^{p-1} \delta_i$ | $\bigoplus_{i=0}^{p-1} \delta_i$ |


```
int MPI_Allreduce(...)
```

| Before: | Node ₀ | Node ₁ | Node ₂ | Node ₃ |
|---------|-------------------|-------------------|-------------------|-------------------|
| | δ_0 | δ_1 | δ_2 | δ_3 |

| After: | Node ₀ | Node ₁ | Node ₂ | Node ₃ |
|--------|------------------------------|------------------------------|------------------------------|------------------------------|
| | $\prod_{i=0}^{p-1} \delta_i$ | $\prod_{i=0}^{p-1} \delta_i$ | $\prod_{i=0}^{p-1} \delta_i$ | $\prod_{i=0}^{p-1} \delta_i$ |

```
MPI_Allreduce( sendBuffer,  recvBuffer,
               count, type, operation,
               communicator );
```

```
int MPI_Alltoall(...)
```

| | Node ₀ | Node ₁ | Node ₂ | Node ₃ |
|---------|--------------------|--------------------|--------------------|--------------------|
| Before: | v ₀ [0] | v ₁ [0] | v ₂ [0] | v ₃ [0] |
| | v ₀ [1] | v ₁ [1] | v ₂ [1] | v ₃ [1] |
| | v ₀ [2] | v ₁ [2] | v ₂ [2] | v ₃ [2] |
| | v ₀ [3] | v ₁ [3] | v ₂ [3] | v ₃ [3] |

| | Node ₀ | Node ₁ | Node ₂ | Node ₃ |
|--------|--------------------|--------------------|--------------------|--------------------|
| After: | v ₀ [0] | v ₀ [1] | v ₀ [2] | v ₀ [3] |
| | v ₁ [0] | v ₁ [1] | v ₁ [2] | v ₁ [3] |
| | v ₂ [0] | v ₂ [1] | v ₂ [2] | v ₂ [3] |
| | v ₃ [0] | v ₃ [1] | v ₃ [2] | v ₃ [3] |

```
int MPI_Alltoall(...)
```

| | Node ₀ | Node ₁ | Node ₂ | Node ₃ |
|---------|--------------------|--------------------|--------------------|--------------------|
| Before: | v ₀ [0] | v ₁ [0] | v ₂ [0] | v ₃ [0] |
| | v ₀ [1] | v ₁ [1] | v ₂ [1] | v ₃ [1] |
| | v ₀ [2] | v ₁ [2] | v ₂ [2] | v ₃ [2] |
| | v ₀ [3] | v ₁ [3] | v ₂ [3] | v ₃ [3] |

| | Node ₀ | Node ₁ | Node ₂ | Node ₃ |
|--------|--------------------|--------------------|--------------------|--------------------|
| After: | v ₀ [0] | v ₀ [1] | v ₀ [2] | v ₀ [3] |
| | v ₁ [0] | v ₁ [1] | v ₁ [2] | v ₁ [3] |
| | v ₂ [0] | v ₂ [1] | v ₂ [2] | v ₂ [3] |
| | v ₃ [0] | v ₃ [1] | v ₃ [2] | v ₃ [3] |

```
MPI_Alltoall( sendBuffer, sendCount, sendType,  
              recvBuffer, recvCount, recvType,  
              communicator );
```

More Collectives

Variable length

- `MPI_Scatterv`
- `MPI_Gatherv`
- `MPI_Allgatherv`
- `MPI_Alltoallv`

More Collectives

Variable length

- MPI_Scatterv
- MPI_Gatherv
- MPI_Allgatherv
- MPI_Alltoallv

```
MPI_Scatter(  
    sendBuffer, sendCount,  
            sendType,  
    recvBuffer, recvCount,  
            recvType,  
    root, communicator );
```

```
MPI_Scatterv(  
    sendBuffer, sendCounts[],  
    displs[],  sendType,  
    recvBuffer, recvCount,  
            recvType,  
    root, communicator );
```

Even more Collectives

Partial reduction

- `MPI_Scan`

Non-blocking collectives

- `MPI_I*`