# COMP 364 / 464
# High Performance Computing

## **Collection Commumcation:**

### Message Passing Interface (MPI)

# MPI Collective Communications

- Every process *within the communicator* MUST call the routine
  - All calls are blocking
  - A process may return when participation is complete
  - May or may not synchronize (implementation dependent)
- You can build your own using p2p methods but MPI has done most of this for us (send/recv are the building blocks)
- Send and Receive sizes must match
  - mapping may vary
- Basic calls have a *root (or destination)* — "all" versions don't

# MPI Collective Communications
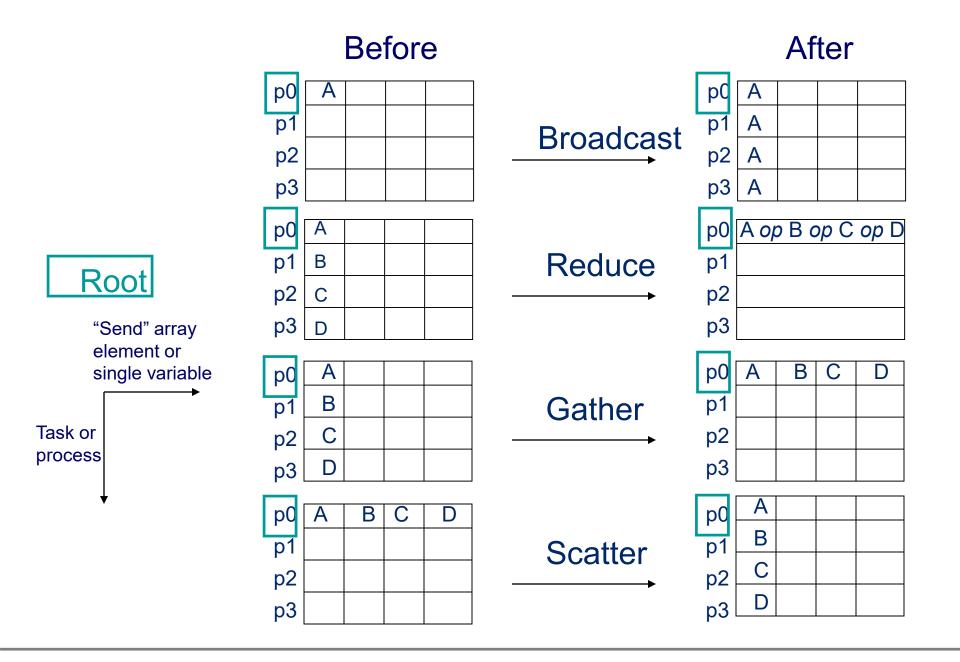
- Involves a group of processes.
- Basic Routines

  Broadcast—      `MPI_Bcast()`

  Reduce—      `MPI_Reduce()`

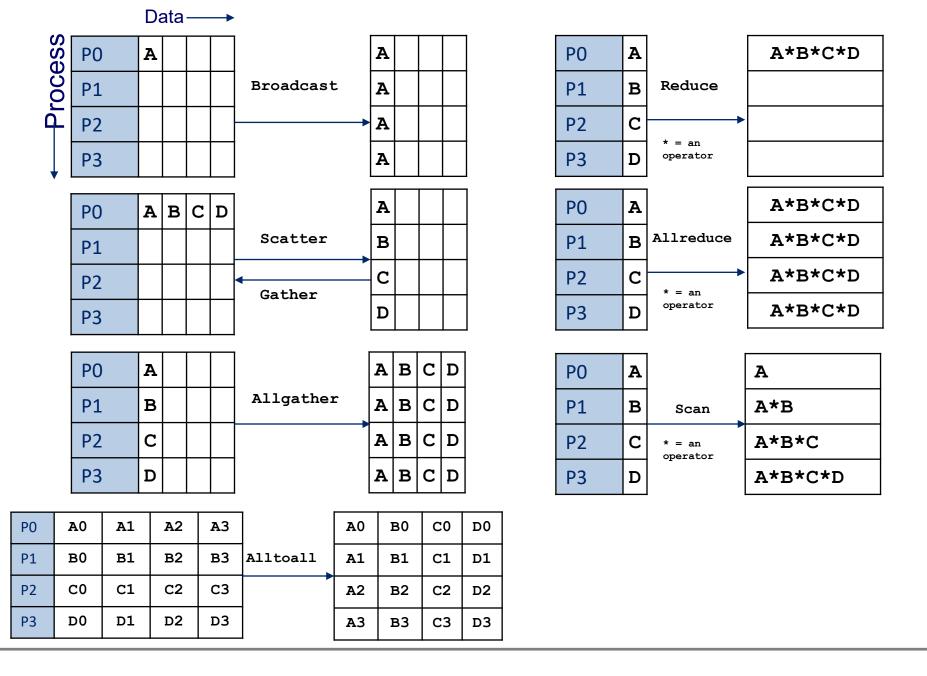  Gather/Scatter— `MPI_Gather()MPI_Scatter()...`

- "All" versions

  `MPI_Allreduce()`

  `MPI_Allgather()`

  `MPI_Alltoall()...`

- Others

  `MPI_Barrier ...`

# Before

## Broadcast

| p0 | A | | | |
| p1 | | | | |
| p2 | | | | |
| p3 | | | | |

## Reduce

| p0 | A | | | |
| p1 | B | | | |
| p2 | C | | | |
| p3 | D | | | |

## Gather

| p0 | A | | | |
| p1 | B | | | |
| p2 | C | | | |
| p3 | D | | | |

## Scatter

| p0 | A | B | C | D |
| p1 | | | | |
| p2 | | | | |
| p3 | | | | |

# After

## Broadcast

| p0 | A | | | |
| p1 | A | | | |
| p2 | A | | | |
| p3 | A | | | |

## Reduce

| p0 | A *op* B *op* C *op* D |
| p1 | |
| p2 | |
| p3 | |

## Gather

| p0 | A | B | C | D |
| p1 | | | | |
| p2 | | | | |
| p3 | | | | |

## Scatter

| p0 | A | | | |
| p1 | B | | | |
| p2 | C | | | |
| p3 | D | | | |

Root

"Send" array element or single variable

Task or process

COMP 364/464: High Performance Computing

COMP 364/464: High Performance Computing

# Broadcast Operation: `MPI_Bcast`

- All nodes call `MPI_Bcast`
- One node (`root`) sends a message to all
  - all others receive the message
- All MPI functions have error handling. MPI_SUCCESS = all good.

```
int ierr = MPI_Bcast(&dat, cnt, datatype, root,
   comm);
```

# Reduction Operations

- Used to combine (reduce) partial results from all processors

- Result returned to root processor

- pre-defined or user-defined operations

  - Predefined: associative & commutative (com)
    Order may not be canonical (i.e., rank order)

  - User defined: Must be associative. com or non-com
    "Canonical" evaluation

- Works on a scalar variable or arrays (elemental)

# MPI_Reduce

**`ierr = MPI_Reduce(&sbuf[0], &rbuf[0], count, datatype, operator, root, comm)`**

- Parameters
  - like **`MPI_Bcast`**, a root is specified
  - operation is a type of mathematical operation
- Applies the operator to each element globally
  - send and receive buffers are the same size
- Use **`MPI_Op_create`** for user-defined operation.

# Operations for MPI_Reduce

| | |
|---|---|
| **MPI_PROD** | Product |
| **MPI_SUM** | Sum |
| **MPI_LAND** | **L**ogical and |
| **MPI_LOR** | **L**ogical or |
| **MPI_LXOR** | **L**ogical exclusive or |
| **MPI_BAND** | **B**itwise and |
| **MPI_BOR** | **B**itwise or |
| **MPI_BXOR** | **B**itwise exclusive or |
| **MPI_MAX** | Maximum |
| **MPI_MIN** | Minimum |
| **MPI_MAXLOC** | Maximum value and location |
| **MPI_MINLOC** | Minimum value and location |

# Dot Product of Two Vectors

```c
double a[N], b[N];
…
double localSum=0.0;
for (int i = 0; i < N; ++i)
    localSum += a[i]*b[i];


double globalSum;
const int root = 0;
MPI_Reduce(&localSum, &globalSum, 1, MPI_DOUBLE,
    MPI_SUM, root, MPI_COMM_WORLD);


printf("rank= %d local= %f global= %f\n",
    rank, localSum, globalSum);
```

# MPI_Scatter Syntax

```
ierr = MPI_Scatter(&sbuf[0], scnt, stype, &rbuf[0],
rcnt, rtype, root, comm);
```

- Parameters
  - **sbuf** = array of size np*scnt (np = # of ranks)
  - **scnt** = number of elements sent to each processor
  - **rcnt** = number of element(s) obtained from the root processor
  - **rbuf** = element(s) obtained from the root processor (rcnt in size)

e.g. `MPI_Scatter(S, 1, stype, R, 1, rtype, root, comm)`

              Array       Scalar

COMP 364/464: High Performance Computing