

# MPI

- Library intended for distributed, high-performance computing applications
  - The de-facto standard for high-performance computing (HPC)
  - Provides point-to-point message passing via Send and Receive
  - Also provides high-level operations that appear in many HPC apps
    - Some of these are: Barrier, Reduce, Scatter, Gather, All-to-All

# MPI, continued

- Library intended for distributed, high-performance computing applications
  - Programming model is SPMD (single program multiple data), where each process:
    - runs identical code image but operates on different data
      - while each process could do completely different things, no actual MPI program does that
    - occasionally executes global sync operations
  - MPI programmer writes one program
    - the program is then executed on (user-specified) N hosts, according to a user host file
      - if no host file, all MPI processes execute on the same machine
      - details of how program is executed on N hosts handled by MPI implementation

# MPI programs

- Must have two functions
  - MPI\_Init (“start”; note: implicit barrier)
  - MPI\_Finalize (“we’re done”; also implicit barrier)
- Practically, must have two more functions
  - MPI\_Comm\_size (returns total number of MPI processes)
  - MPI\_Comm\_rank (returns caller’s process id)
- The actual computation is placed in between MPI\_Comm\_rank/MPI\_Comm\_size and MPI\_Finalize

# Sending and Receiving in MPI

- All four combinations of blocking/nonblocking send/receive are possible
  - MPI\_Ssend (blocking send)
  - MPI\_Isend (nonblocking send)
  - MPI\_Recv (blocking receive)
  - MPI\_Irecv (nonblocking receive)
  - MPI\_Wait (paired with Isend or Irecv to make sure operation has completed; i.e., it is safe to overwrite [Send] or use [Recv] the data)
  - MPI\_Send: nonblocking if data is small

# Sending and Receiving in MPI

- `MPI_Send` and `MPI_Recv` take as parameters:
  - buffer, which is the data being sent or received
  - number of elements in the buffer
  - type of elements in the buffer
  - destination (`MPI_Send`) or source (`MPI_Recv`)
  - tag---must match other end for send/receive to match
  - “communicator”; always `MPI_COMM_WORLD` in 422
    - communicators can in general be used for non-global barriers
- One extra parameter in `MPI_Recv`, which is the status
  - We will use this in program 3, but it is rarely used

# Example use of MPI\_Isend (with 2 MPI processes)

```
int A[10] = {...}, B[10];  
MPI_Status status; MPI_Request request;  
  
// initialization here (MPI_Init, MPI_Comm_rank, MPI_Comm_size)  
if (myId == 0) {  
    MPI_Isend(A, 10, MPI_INT, 1, 17, MPI_COMM_WORLD, &request);  
    // do anything that doesn't involve writing to A (if A is written here, it's a race condition)  
    MPI_Wait(&request, &status);  
    // now can safely overwrite A  
}  
else {  
    MPI_Recv(B, 10, MPI_INT, 0, 17, MPI_COMM_WORLD, MPI_STATUS_IGNORE);  
}
```

# Example use of MPI\_Irecv (with 2 MPI processes)

```
int A[10] = {...}, B[10];
MPI_Status status; MPI_Request request;

// initialization here (MPI_Init, MPI_Comm_rank, MPI_Comm_size)
if (myId == 0) {
    MPI_Irecv(B, 10, MPI_INT, 1, 17, MPI_COMM_WORLD, &request);
    // do anything that doesn't involve using B (if B is used here, it's a race condition)
    MPI_Wait(&request, &status);
    // now can safely use B
}
else {
    MPI_Send(A, 10, MPI_INT, 0, 17, MPI_COMM_WORLD);
}
```

# Collective Communication in MPI

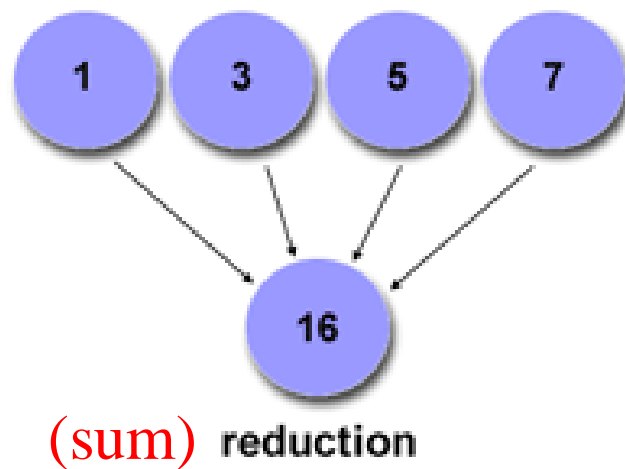
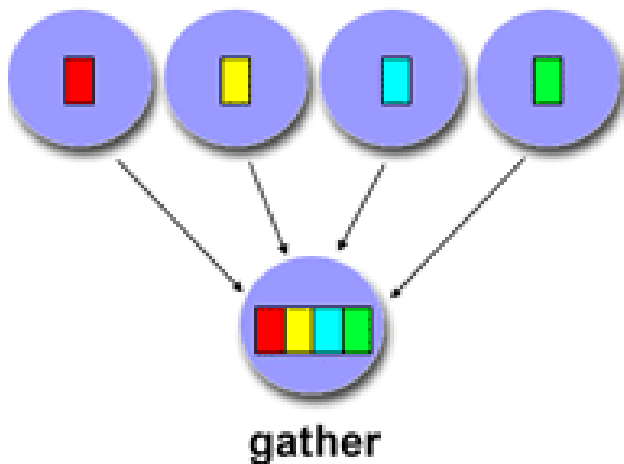
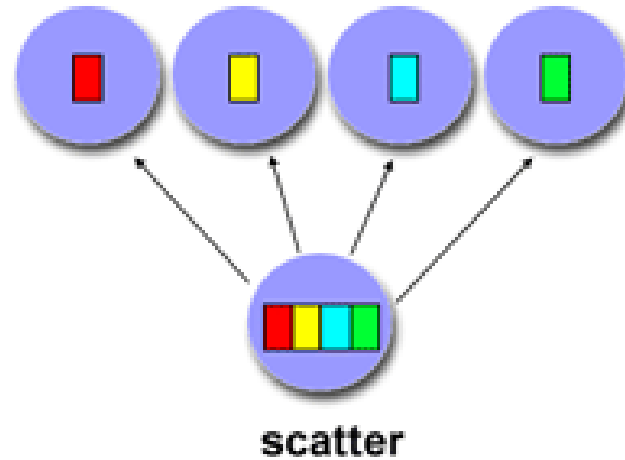
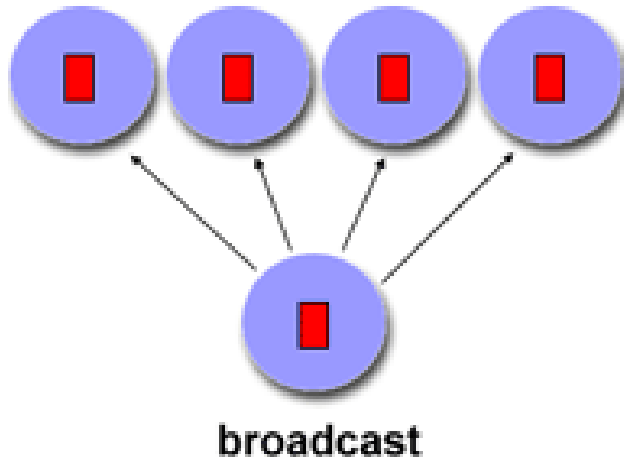
- Collective calls are ones that involve all processes in a communicator (for this class, this means all processes)
  - MPI\_Bcast (usual definition; one sends to many)
  - MPI\_Scatter (given an array on “root” process; send equal-size subarray to each non-root process; awkward because **each** process must specify sendbuf/recvbuf)
  - MPI\_Gather (reverse of MPI\_Scatter)
  - MPI\_Reduce (reduced value ends up at root)
  - MPI\_Allreduce (MPI\_Reduce + dissemination to all)



# Collective Communication in MPI

- Collectives aren't strictly needed
  - I.e., every collective can be implemented with some sequence of sends and receives
- Collectives have advantages, though:
  - Easier for the programmer
  - MPI runtime knows about the operation ahead of time, so it can implement it efficiently
    - Example: MPI\_Allreduce can use a tree of  $\log(P)$  levels or a tree of 1 level

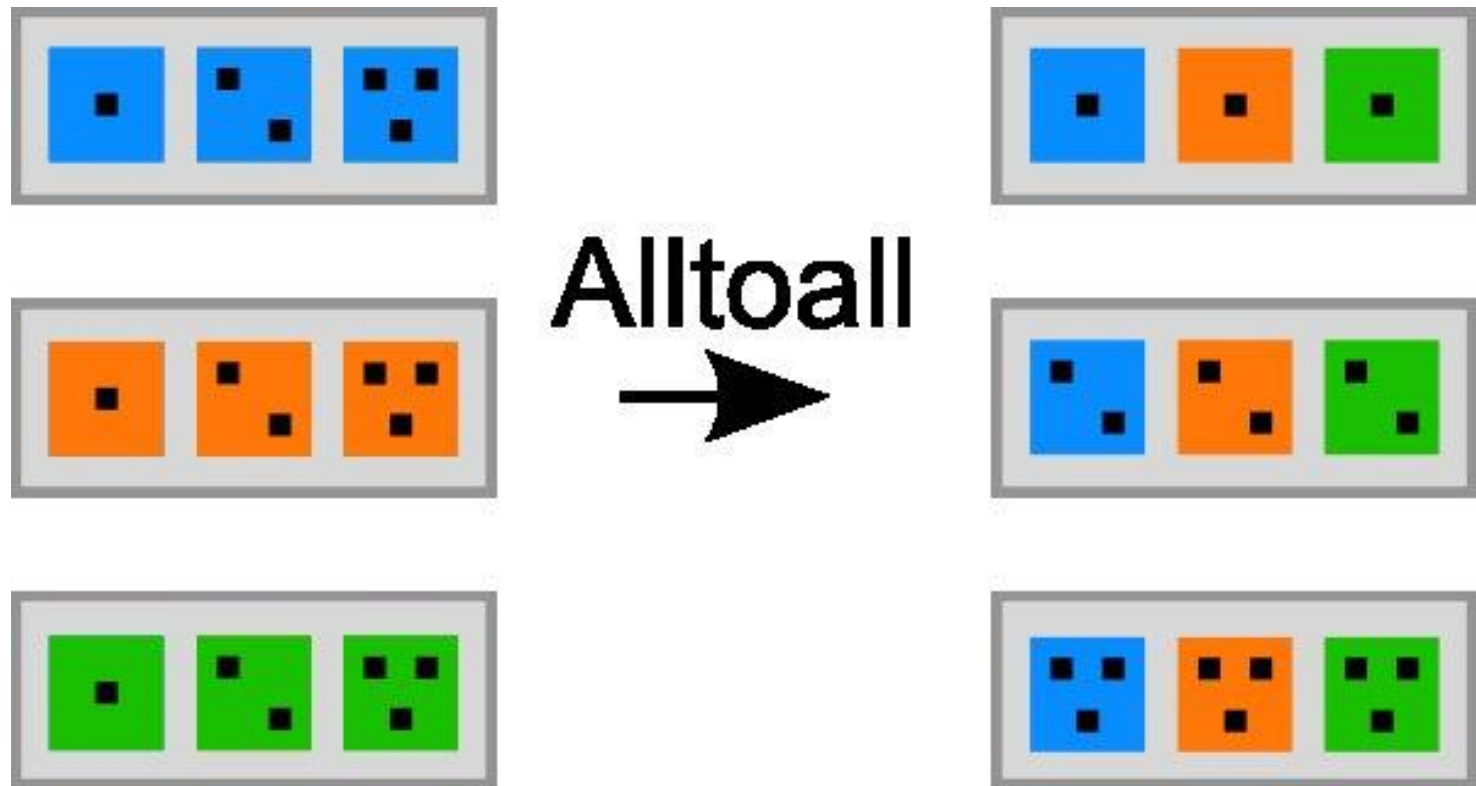
# Collective Communication in MPI



# Collective Communication in MPI

- MPI\_Alltoall (every process sends a unique part of its buffer to each other process)
- MPI\_Alltoallv (generalized MPI\_Alltoall; parts of buffer can have variable size)

# Collective Communication in MPI



# Compiling/Running an MPI Program

- To compile, **you must use** `mpicc`

Example: `mpicc -o mm mpi-mm.c`

- To run, **you must use** `mpirun`. Use the `-n` option to specify number of MPI processes
- Also, `mpirun` arguments come first, then executable, then program command line args

Example: `mpirun -n 4 ./mm 100`

Args to mpirun



Executable

Args to mm program