

# Numerical Simulation and Scientific Computing II

## Lecture 1: Introduction and Distributed Parallel Computing I



Josef Weinbub, Paul Manstetten, Luiz  
Aguinsky, Heinz Pettermann, Marius  
Schasching, Jesús Carrete Montana,  
Francesco Zonta, Kevin Sturm

Institute for Microelectronics  
TU Wien



[nssc@iue.tuwien.ac.at](mailto:nssc@iue.tuwien.ac.at)

# Outline

---

- **Introduction to the Lecture**
- **Distributed Parallel Computing I**
  - Primer
  - Overview
  - Process Model and Language Bindings
  - Messages and Point-to-Point Communication
  - Examples
- Quiz

# Team Presentation

Core Team (Organization, MPI, ODE, FVM)

- **Josef Weinbub**
- **Paul Manstetten**
- **Luiz Felipe Aguirsky**  
*Institute for Microelectronics*



Extended Interdisciplinary Team

- **Francesco Zonta** (Fluid Dynamics)  
*Institute of Fluid Mechanics and Heat Transfer*
- **Heinz Pettermann** and **Marius Schasching** (Finite Elements)  
*Institute of Lightweight Design and Structural Biomechanics*
- **Jesús Carrete Montana** (Molecular Dynamics)  
*Institute of Materials Chemistry*
- **Kevin Sturm** (Partial Differential Equations)  
*Institute of Analysis and Scientific Computing*



[nssc@iue.tuwien.ac.at](mailto:nssc@iue.tuwien.ac.at)

- Acceptance to course via TISS: Check today <https://tiss.tuwien.ac.at/course/courseAnnouncement.xhtml?dswid=9250&dsrid=235&courseNumber=360243&courseSemester=2021S>
- News will be sent via TISS' notification feature: Regularly monitor your TISS News (at least weekly)!

360.243 Numerical Simulation and Scientific Computing II 2021S ▾  
2021S, VU, 3.0h, 6.0EC Announcement | Processing | Communication | Event coordination

Description **News** Course registration Feedback

- To TUWEL online course
- To course forum

**TUWEL**

**Properties**

- Semester hours: 3.0
- Credits: 6.0
- Type: VU Lecture and Exercise
- Format: Distance Learning

**Learning outcomes**

After successful completion of the course, students are able to select and apply fundamental methods of scientific computing and to judge the challenges regarding computing time and implementation effort. Furthermore, the students have solution skills


- Course material (slides, handouts, zoom link) provided via TUWEL course:  
<https://tuwel.tuwien.ac.at/course/view.php?idnumber=360243-2021S>
- TUWEL course also “reachable” via TISS course, watch out for the TUWEL icon:

360.243 Numerical Simulation and Scientific Computing II 2021S ▾

2021S, VU, 3.0h, 6.0EC Announcement | Processing | Communication | Event coordination

Description News Course registration Feedback

- To TUWEL online course
- To course forum



Properties

- Semester hours: 3.0
- Credits: 6.0
- Type: VU Lecture and Exercise
- Format: Distance Learning

Learning outcomes

After successful completion of the course, students are able to select and apply fundamental methods of scientific computing and to judge the challenges regarding computing time and implementation effort. Furthermore, the students have solution skills

# General Goals

---

- **Introduction to advanced methods of CSE**
  - **Distributed Parallel Computing**
  - **Methods for Ordinary Differential Equations**
  - **Classification and Analyses of Partial Differential Equations**
  - **Finite Volume Method**
  - **Finite Element Method**
  - **Fluid Dynamics**
  - **Molecular Dynamics**

# Course Calendar

March 4	Lecture
March 11	Lecture with Exercise Handout
March 18	Lecture with Exercise Support
March 25	Lecture with Exercise Support
<i>April 13</i>	<i>Exercise 1 Submission</i>
April 15	Lecture
April 22	Lecture with Exercise Handout
April 29	Lecture with Exercise Support
<i>May 4</i>	<i>Exercise 2 Submission</i>
May 6	Lecture with Exercise Handout
May 12	Exercise Support
May 20	Lecture with Exercise Handout and Support
<i>May 26</i>	<i>Exercise 3 Submission</i>
May 27	Lecture with Exercise Support
June 2	Exercise Support
<i>June 10</i>	<i>Exercise 4 Submission</i>
June 30	Written (Online) Exam

- **Every slot starts at 14:00 sharp!**
- **Will put lecture slides/exercise material on TUWEL right before slot.**
- ***Exercise submissions until 8:00!***
- **Schedule also in TISS!**

# Necessary Background

- **C++ and Python!** (see TISS)
  - If not experienced: Checkout online tutorials asap!
- **Not a formal requirement**  
yet the expected level of expertise:  
*Numerical Simulation and Scientific Computing I*  
*360.242*



# Exercise Rules and Course Grade

- 4 mandatory exercises over the whole lecture
- Groups of 2-3 students must be formed; student-picked.
- **Code copied? 0 points for both groups!**
- Each exercise will be graded separately from 0-10 points
- Access for the final exam: **sum of the points  $\geq 28.0$**   
(same procedure as with NSSC I)
- Course grade: 1/4<sup>th</sup> Exercises, 3/4<sup>th</sup> Exam
- **Submission deadlines are hard deadlines**
- Submissions: per email to [nssc@iue.tuwien.ac.at](mailto:nssc@iue.tuwien.ac.at) as a single compressed file  
(More information with each exercise)

# Rules - Quizzes

---

- In the end of each lecture, you will receive 5 questions
- 3 questions reviewing the current lecture
- 2 questions preparing for the next
- Discussion in the beginning of next class
- Participation is voluntary but encouraged
- These questions might help you learn for the exam!

# Computer Resources

---

- **It is expected that you have access to a modern, x86-based, personal computer (laptop, home desktop, etc.)**
- **Linux (virtual machine, dual boot or full OS) is required**
  - **More information on the development environment for a particular exercise will be made available with the exercise handout**

# Outline

---

- Introduction to the Lecture
- **Distributed Parallel Computing I**
  - **Primer**
  - Overview
  - Process Model and Language Bindings
  - Messages and Point-to-Point Communication
  - Examples
- Quiz

# Acknowledgments

---

Thanks to Rolf Rabenseifner's (HLRS)  
2014 Parallel Programming Workshop Lecture Material on:  
**Introduction to the Message Passing Interface (MPI)**

# Sources

- High Performance Computing Center Stuttgart (HLRS)  
Online Courses  
<https://www.hlr.de/about-us/media-publications/teaching-training-material/>
- YouTube -- search for “Introduction to MPI”, e.g.,  
<https://youtu.be/RoQJNx5npF4> -- Part I of III

# Additional Courses at TU Wien

## **ECTS Courses** (count as free electives)

- 057.020 (Winter term)

### **VSC-School I Courses in High Performance Computing**

- 057.021 (Summer term)

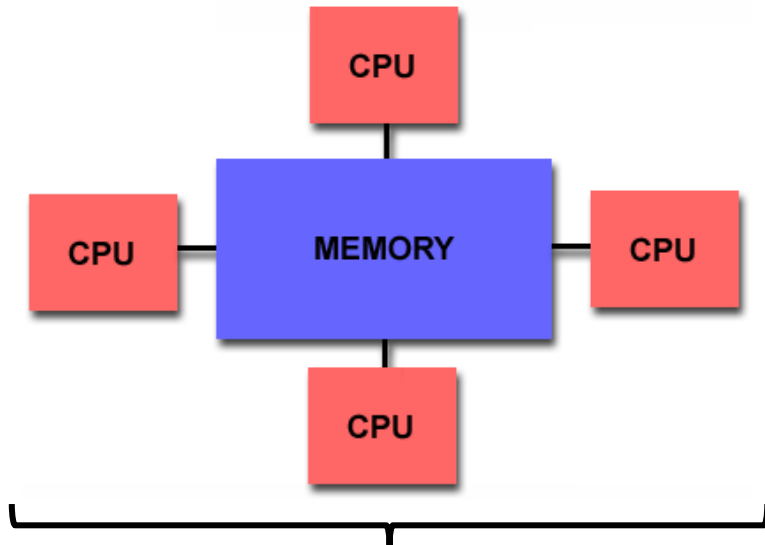
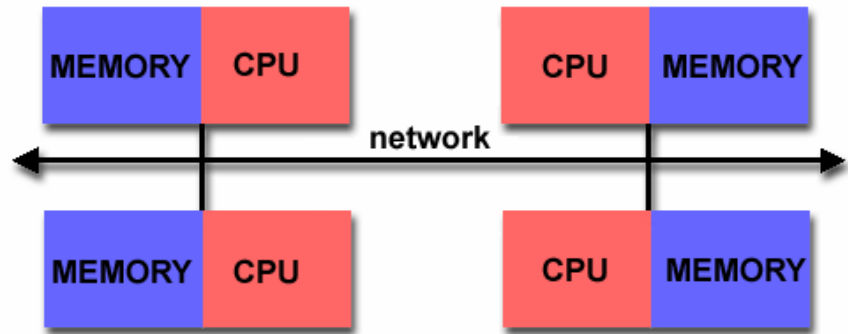
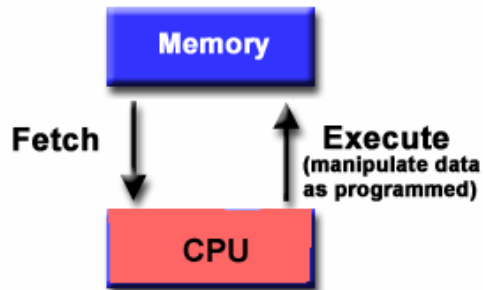
### **VSC-School II Courses in High Performance Computing**

## **Non-ECTS Trainings** (don't count as free electives)

- Node-Level Performance Engineering
- OpenMP
- MPI
- Deep-Learning und GPU programming (OpenACC)
- Hybrid-programming MPI+X

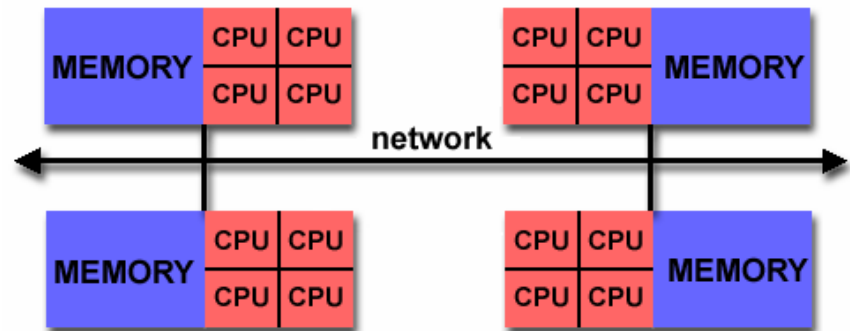
<http://typo3.vsc.ac.at/research/vsc-research-center/vsc-school-seminar/>

# Parallel Computer Architectures



**Easy to use, but not scalable**

Source: Lyle N. Long, PSU



**Difficult to use, but scalable**

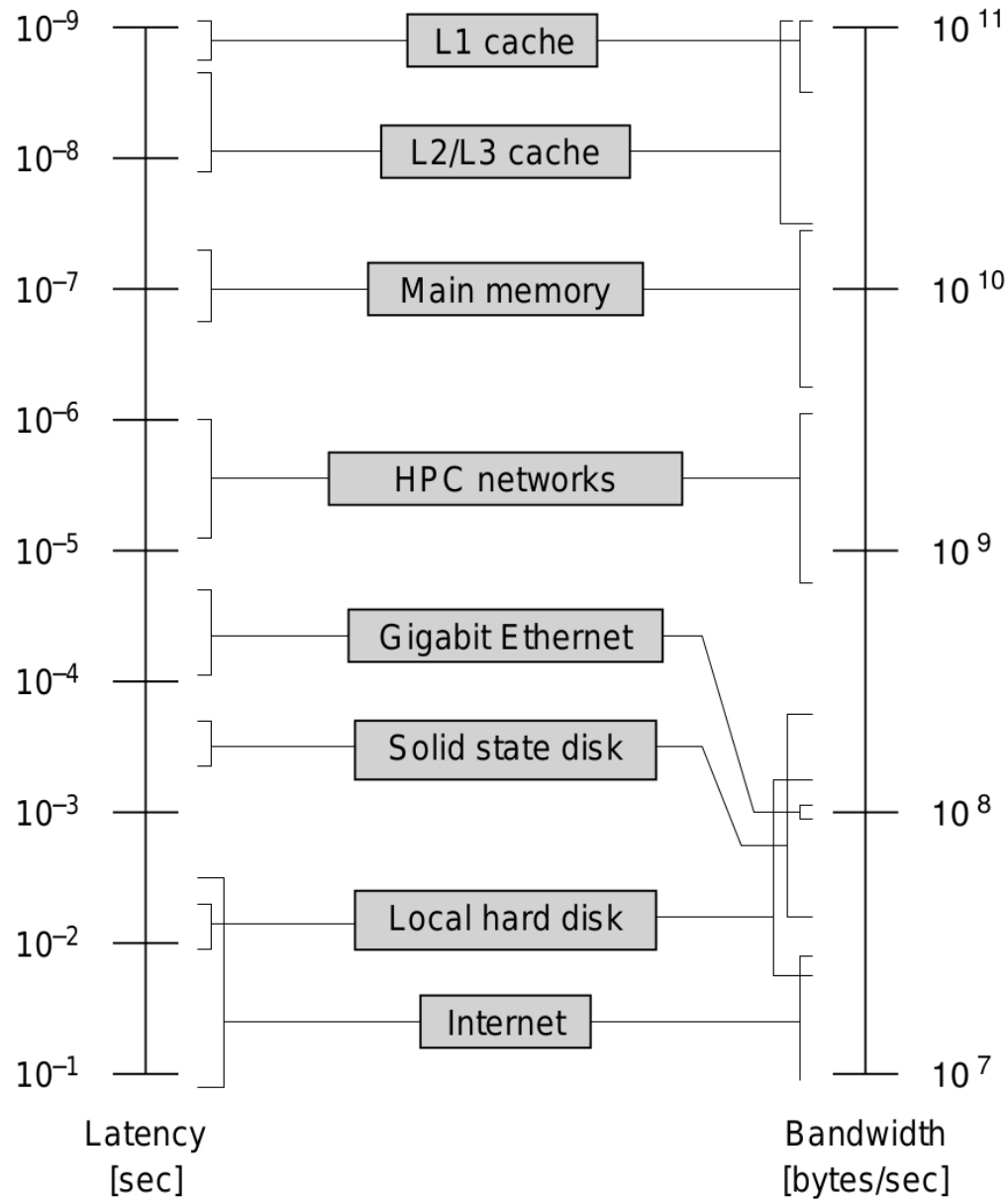


# Distributed Parallel Computing Clusters



**IBM Blue Gene/P: Intrepid  
Argonne National Laboratory**

# Latency / Bandwidth Hierarchy

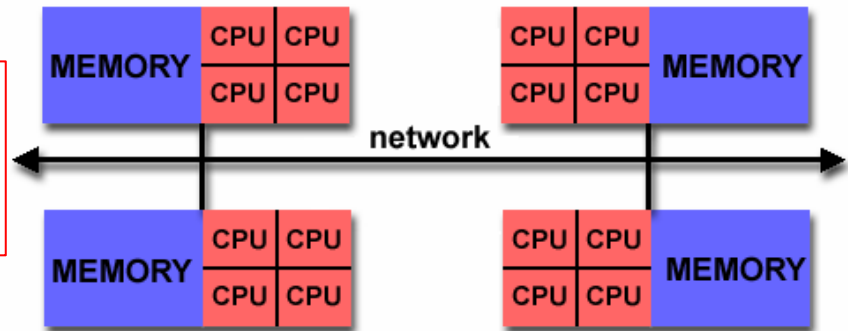


# Processes and Threads

- **Process**

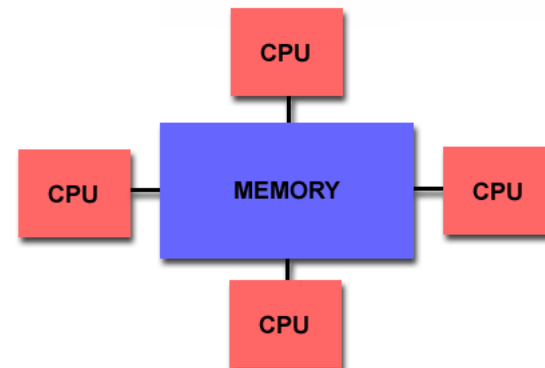
- A process is a program in execution
- Represented by a process control block: all necessary information to run a program

Data exchange happens explicitly via sending messages between processes!



- **Thread**

- Unit of CPU utilization
- A single process can contain multiple threads
- Threads belonging to the same process can share resources



# Outline

---

- Introduction to the Lecture
- Distributed Parallel Computing I
  - Primer
  - **Overview**
  - Process Model and Language Bindings
  - Messages and Point-to-Point Communication
  - Examples
- Quiz

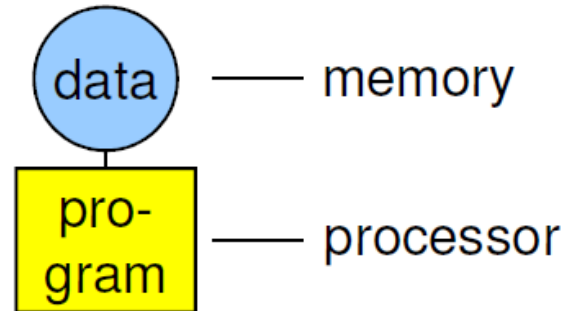
# Message Passing Interface (MPI)

```
#include "mpi.h"
#include <stdio.h>
int main( int argc, char *argv[] )
{
    int myid, numprocs;
    MPI_Init(&argc,&argv);
    MPI_Comm_size(MPI_COMM_WORLD,&numprocs);
    MPI_Comm_rank(MPI_COMM_WORLD,&myid);
    printf("I am %d out of %d\n", myid, numprocs);
    MPI_Finalize();
    return 0;
}
```

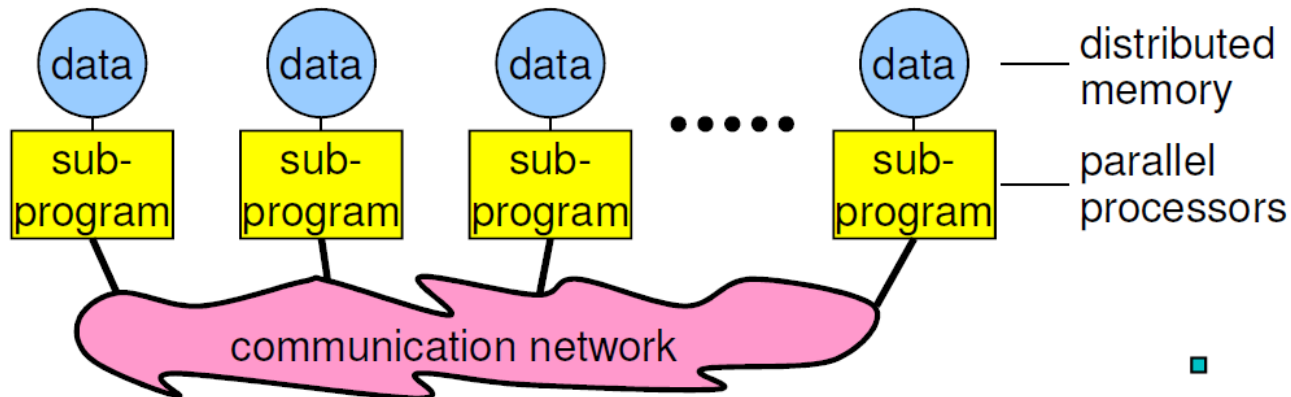
```
> mpiexec -n 4 ./hello
I am 3 out of 4
I am 1 out of 4
I am 0 out of 4
I am 2 out of 4
```

# Message Passing Programming Paradigm

- Sequential Programming Paradigm

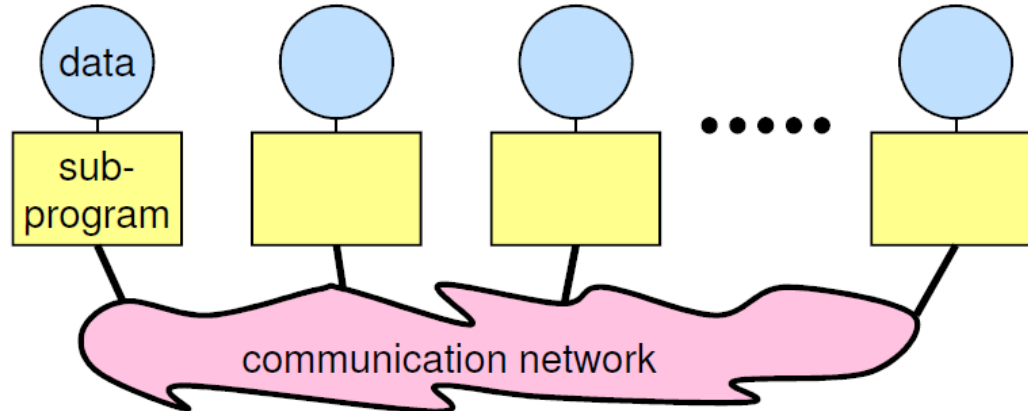


- Message Passing Programming Paradigm



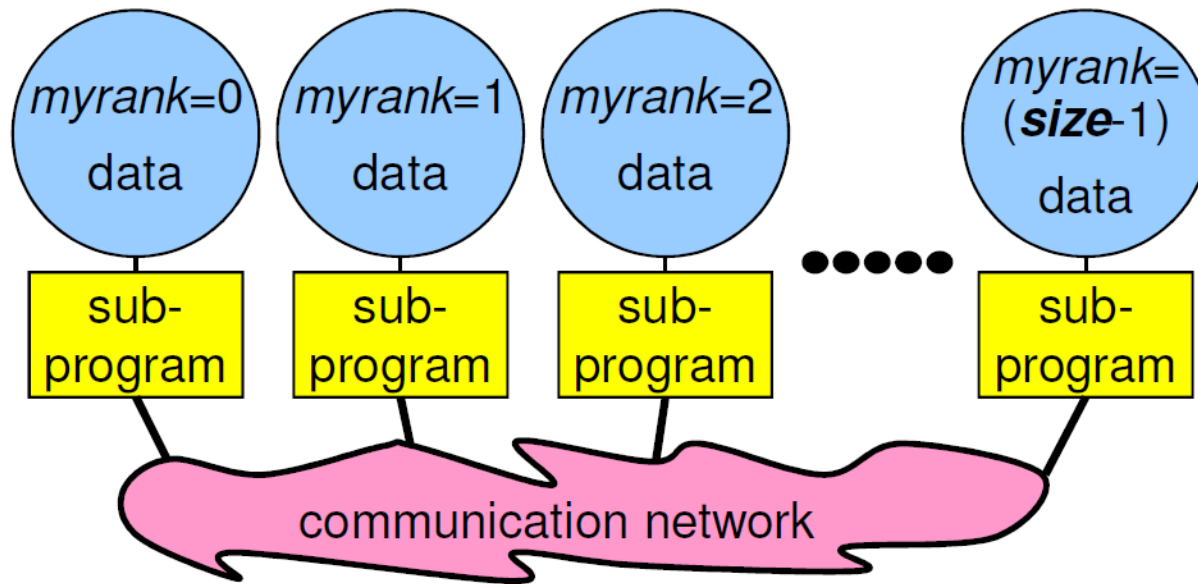
# Message Passing Programming Paradigm

- Each processor in a message passing program runs a *sub-program*:
  - written in a conventional sequential language, e.g., C or Fortran,
  - typically the same on each processor,
  - the variables of each sub-program have
    - the same name
    - but different locations (distributed memory) and different data!
    - i.e., all variables are private
- communicate via special send & receive routines (*message passing*)



# Data and Work Distribution

- the value of *myrank* is returned by special library routine
- the system of *size* processes is started by special MPI initialization
- program (*mpirun* or *mpiexec*)
- all distribution decisions are based on *myrank*
- i.e., which process works on which data





# What is SPMD?

- **Single Program, Multiple Data**
- **Same (sub-)program runs on each processor**
- **MPI allows also MPMD, i.e., Multiple Program, ...**
- **but some vendors may be restricted to SPMD**
- **MPMD can be emulated with SPMD**

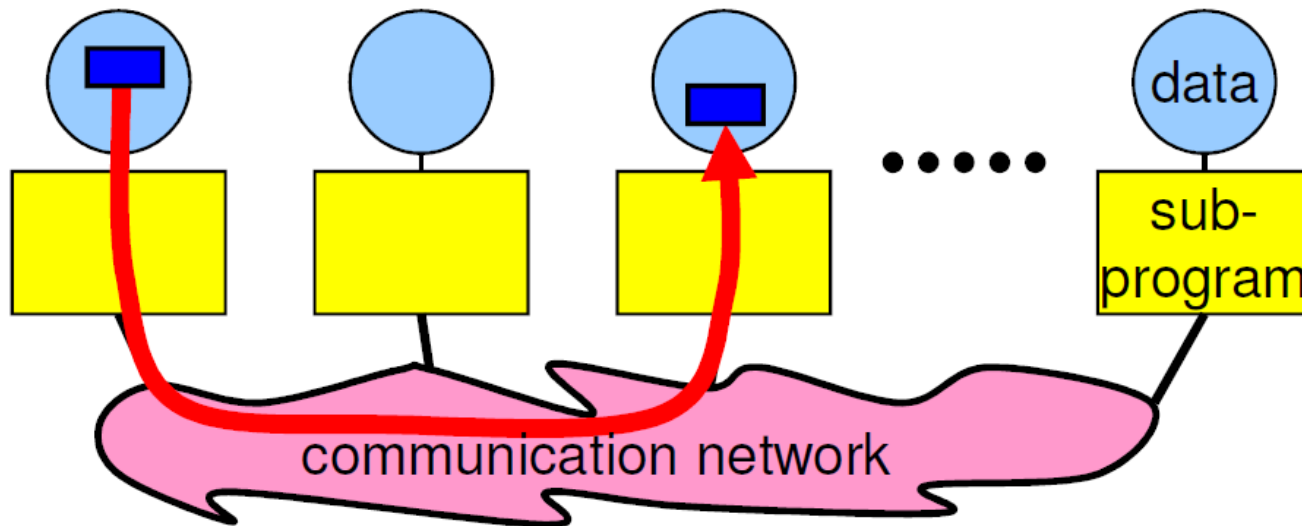
# Emulation of Multiple Program (MPMD)

```
main(int argc, char **argv)
{
    if (myrank < .... /* process should run the ocean model */)
    {
        ocean( /* arguments */ );
    }
    else{
        weather( /* arguments */ );
    }
}
```

# Messages

- Messages are packets of data moving between sub-programs
- Necessary information for the message passing system:
  - sending process – receiving process
  - source location – destination location
  - source data type – destination data type
  - source data size – destination buffer size

} the ranks  
}



- **A sub-program needs to be connected to a message passing system**
- **A message passing system is similar to:**
  - mail box
  - phone line
  - fax machine
  - etc.
- **MPI:**
  - sub-program must be linked with an MPI library
  - sub-program must use include file of this MPI library
  - the total program (i.e., all sub-programs of the program) must be started with the MPI startup tool

# Addressing

- **Messages need to have addresses to be sent to.**
- **Addresses are similar to:**
  - mail addresses
  - phone number
  - fax number
  - etc.
- **MPI: addresses are ranks of the MPI processes (sub-programs)**

# Reception

---

- **All messages must be received.**

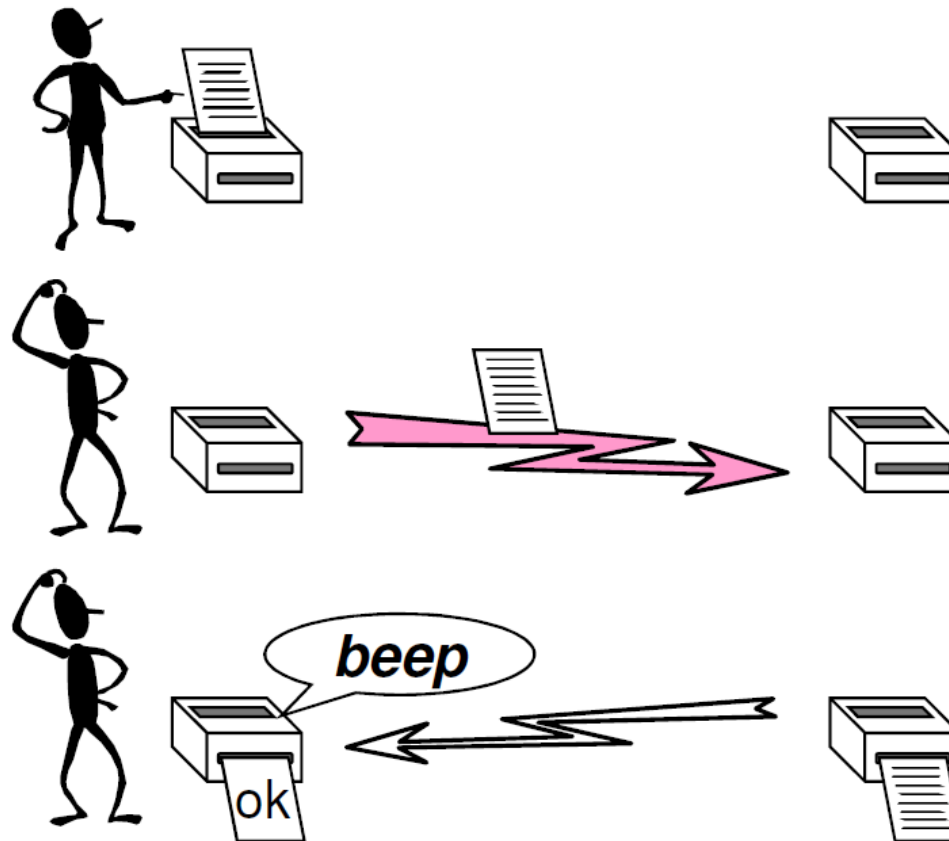
# Point-to-Point Communication

---

- **Simplest form of message passing.**
- **One process sends a message to another.**
- **Different types of point-to-point communication:**
  - **synchronous send**
  - **buffered = asynchronous send**

# Synchronous Sends

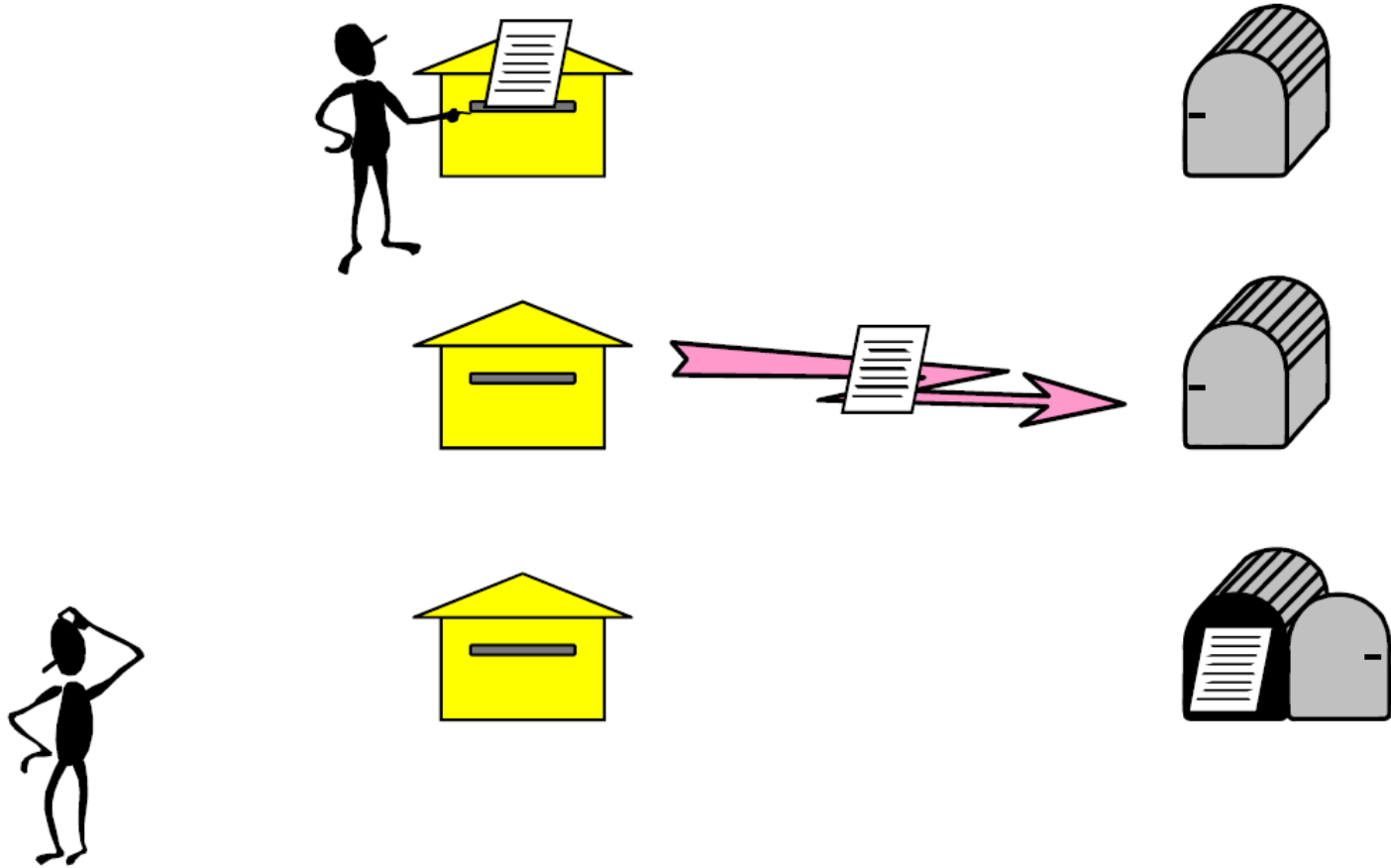
- The sender gets an information that the message is received.
- Analogue to the *beep* or *okay-sheet* of a fax.





# Buffered = Asynchronous Sends

- Only know when the message has left.

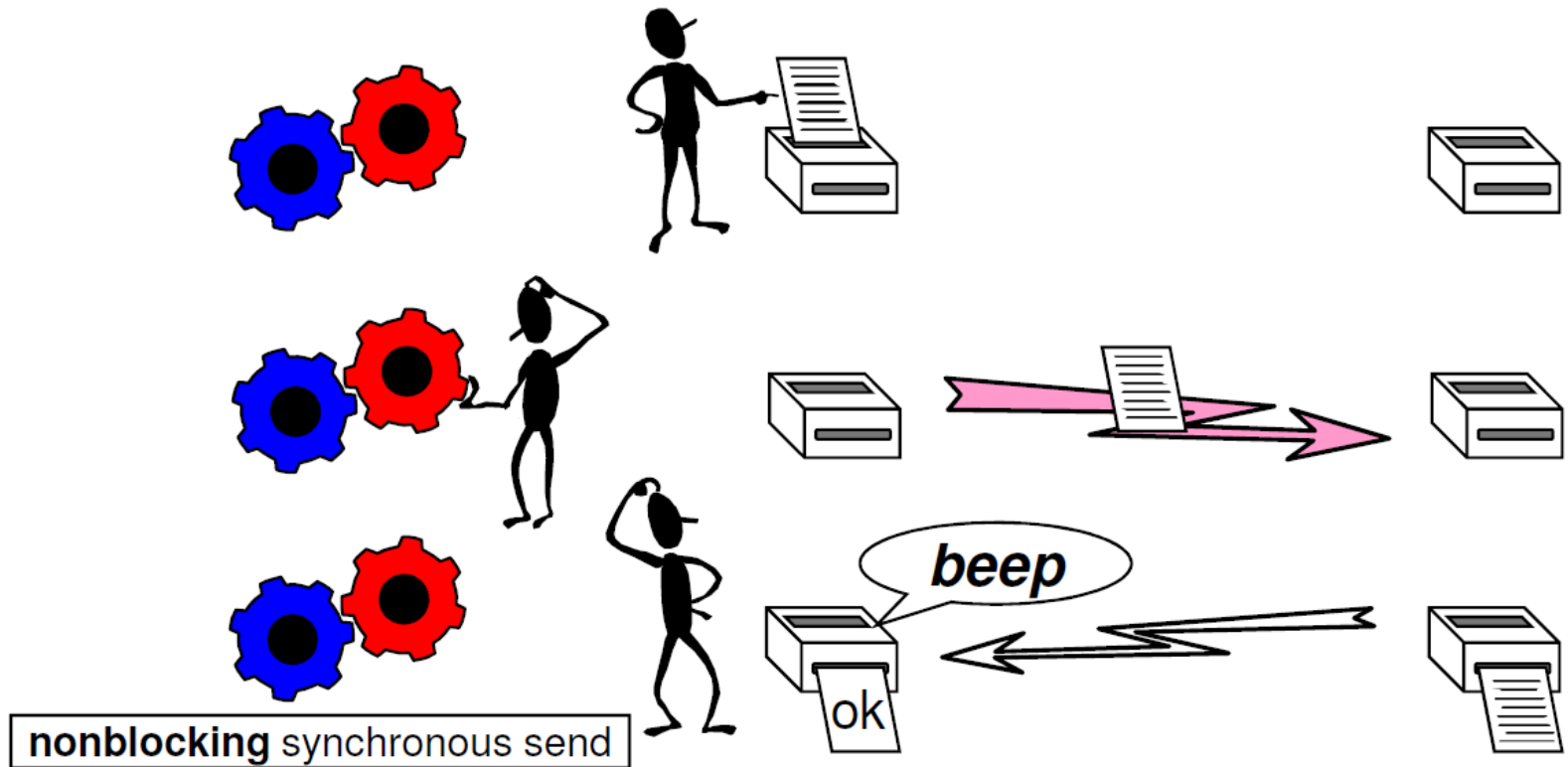


# Blocking Operations

- **Operations are local activities, e.g.,**
  - sending (a message)
  - receiving (a message)
- **Some operations may block until another process acts:**
  - synchronous send operation blocks until receive is posted;
  - receive operation blocks until message was sent.
- **Relates to the completion of an operation.**
- **Blocking subroutine returns only when the operation has completed.**

# Non-Blocking Operations

- **Non-blocking operation:** returns immediately and allow the sub-program to perform other work.
- **At some later time the sub-program must test or wait for the completion of the non-blocking operation.**



# Non-Blocking Operations (cont'd)

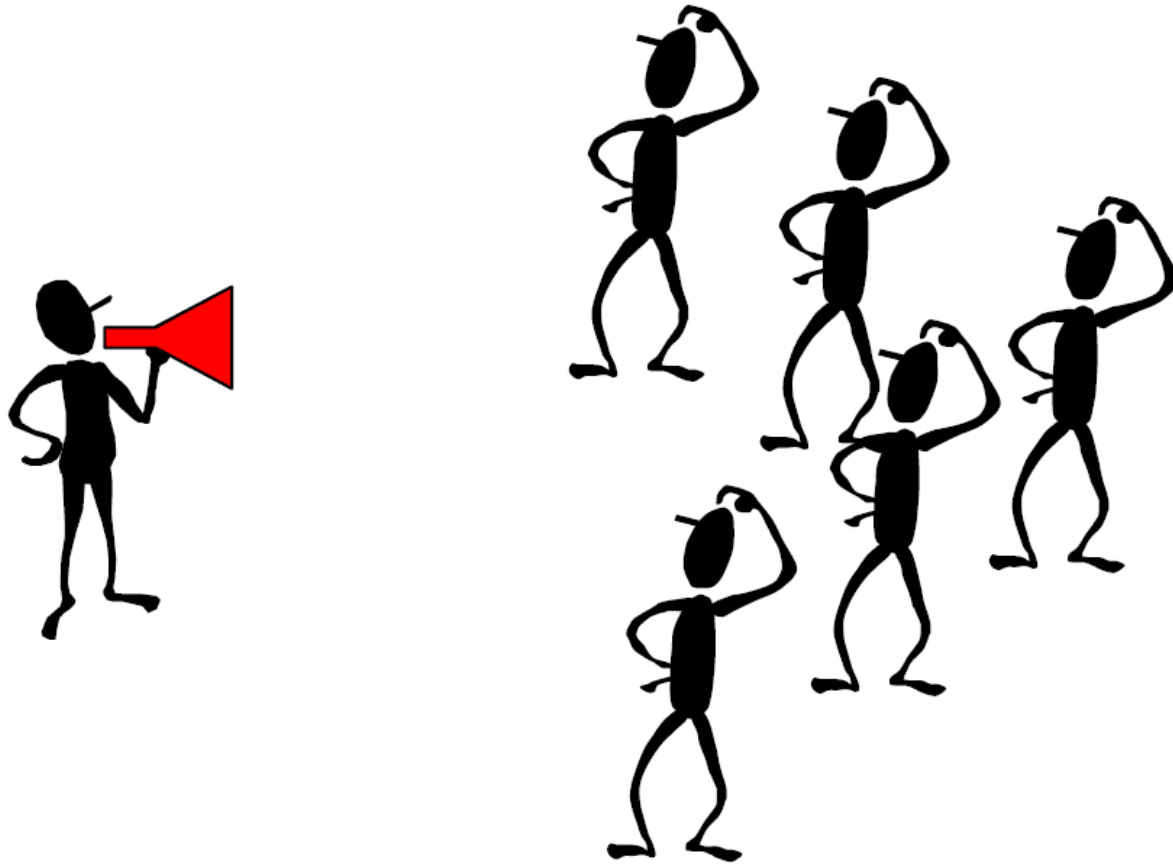
- All non-blocking operations must have matching wait (or test) operations.  
(Some system or application resources can be freed only when the non-blocking operation is completed.)
- A non-blocking operation immediately followed by a matching wait is equivalent to a blocking operation.
- Non-blocking operations are not the same as sequential subroutine calls:
  - the operation may continue while the application executes the next statements!

# Collective Communications

- **Collective communication routines are higher level routines.**
- **Several processes are involved at a time.**
- **May allow optimized internal implementations, e.g., tree based algorithms.**
- **Can be built out of point-to-point communications.**

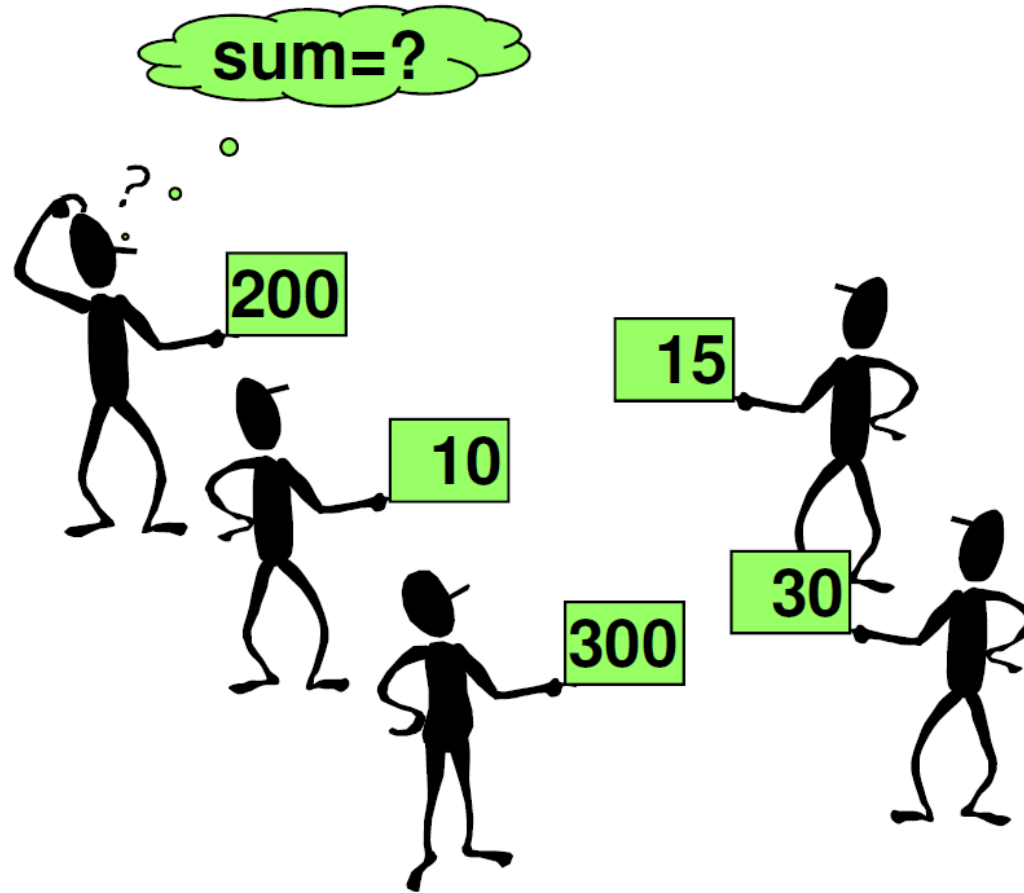
# Broadcast

- A one-to-many communication.



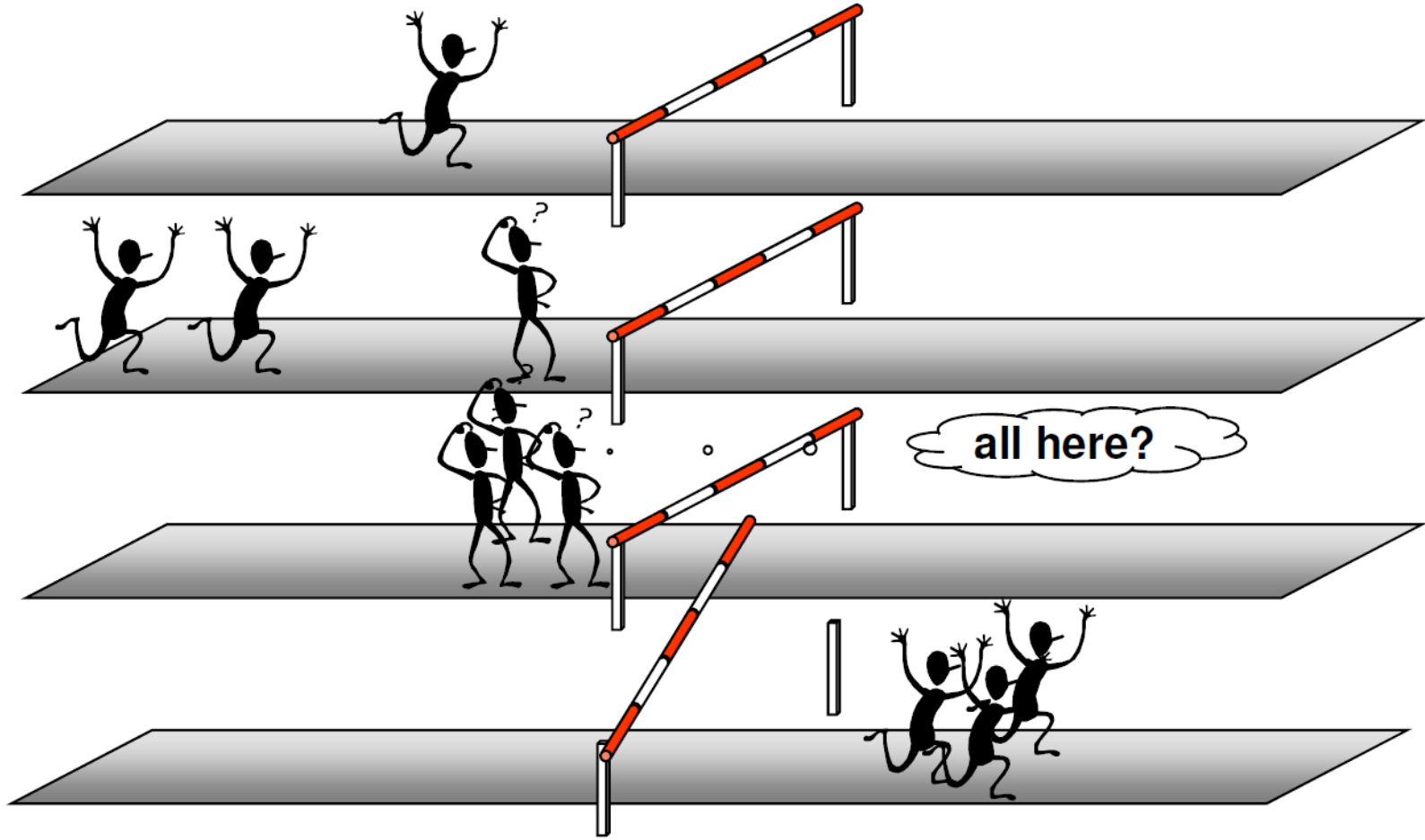
# Reduction Operations

- Combine data from several processes to produce a single result.



# Barriers

- Synchronize processes.





- **MPI is a standard**
  - Each MPI routine is defined
- **MPI Forum: the standardization forum**  
**[www.mpi-forum.org/](http://www.mpi-forum.org/)**
- **MPI-1.0 June 1994**
- **MPI-3.1 June 2015**
- **MPI-4.0 under development**
- **Many MPI libraries ((mostly) adhering to the standard) are available**
  - Open source: OpenMPI, MPICH
  - Native support for C/C++, Fortran  
(this course focuses only on C/C++)

# Outline

---

- Introduction to the Lecture
- Distributed Parallel Computing I
  - Primer
  - Overview
  - **Process Model and Language Bindings**
  - Messages and Point-to-Point Communication
  - Examples
- Quiz

# Header Files and Function Format

- **Header Files**

`#include <mpi.h>`

- **MPI Namespace**

- **MPI\_...** namespace is reserved for MPI constants and routines
- **Application routines and variable names must not begin with MPI\_**

- **Function Format**

`error = MPI_Xxxxxx( parameter, ... );`  
`MPI_Xxxxxx( parameter, ... );`

# Initializing MPI

- **First routine to be called**

```
int MPI_Init( int *argc, char ***argv)
```

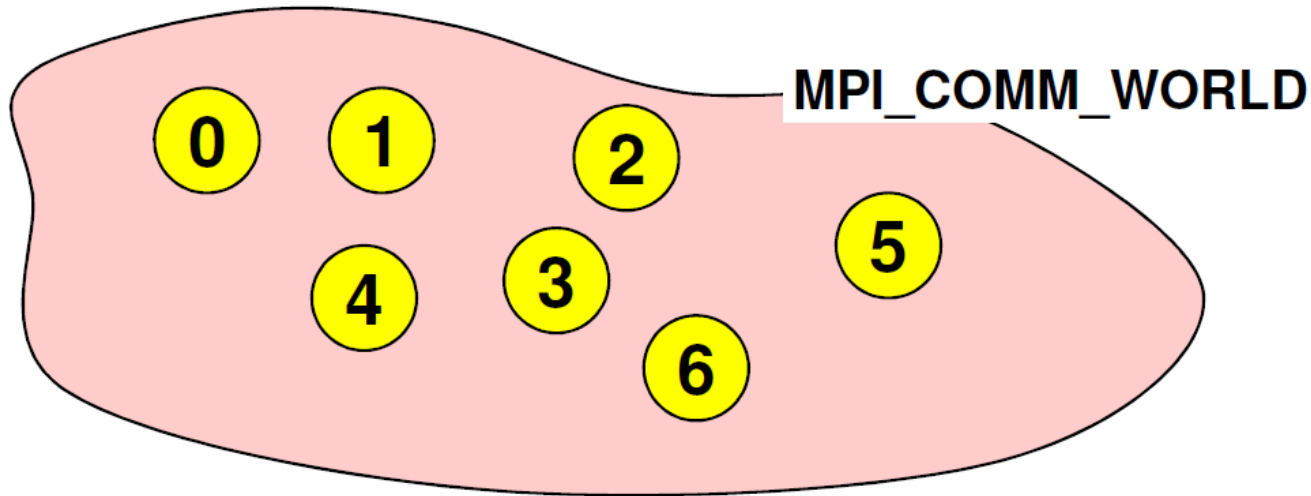
```
#include <mpi.h>
int main(int argc, char **argv)
{
    MPI_Init(&argc, &argv);
    ...
}
```

# Starting the MPI Program

- Start mechanism is implementation dependent
- **mpirun -np *number\_of\_processes* ./executable**  
(most implementations)
- **mpiexec -n *number\_of\_processes* ./executable**  
(with MPI-2 and later)

# Communicator MPI\_COMM\_WORLD

- All processes (= sub-programs) of one MPI program are combined in the communicator MPI\_COMM\_WORLD.
- MPI\_COMM\_WORLD is a predefined handle in mpi.h.
- Each process has its own rank in a communicator:
  - starting with 0
  - ending with (size-1)



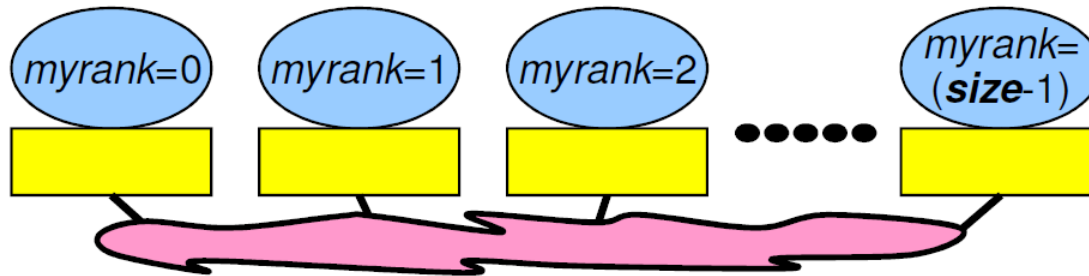
# Handles

- **Handles identify MPI objects.**
- **For the programmer, handles are**
  - **predefined constants in mpi.h**
    - **Example: MPI\_COMM\_WORLD**
    - **Can be used in initialization expressions or assignments.**
    - **The object accessed by the predefined constant handle exists and does not change only between MPI\_Init and MPI\_Finalize.**
  - **values returned by some MPI routines, to be stored in variables, that are defined as**

# Rank

- The rank identifies different processes.
- The rank is the basis for any work and data distribution.

```
int MPI_Comm_rank( MPI_Comm comm, int *rank)
```



```
CALL MPI_COMM_RANK( MPI_COMM_WORLD, myrank, ierror)
```



# Size

- **How many processes are contained within a communicator?**

```
int MPI_Comm_size( MPI_Comm comm, int *size)
```

# Exiting MPI

- **Must be** called last by all processes.
- User must ensure the completion of all pending communications (locally) before calling finalize
- **After MPI\_Finalize:**
  - Further MPI-calls are forbidden
  - Especially re-initialization with MPI\_Init is forbidden
  - May abort all processes except “rank==0” in MPI\_COMM\_WORLD

`int MPI_Finalize()`

# Outline

- Introduction to the Lecture
- Distributed Parallel Computing I
  - Primer
  - Overview
  - Process Model and Language Bindings
  - **Messages and Point-to-Point Communication**
  - Examples
- Quiz

# Messages

- A message contains a number of elements of some particular datatype.
- MPI datatypes:
  - Basic datatype.
  - Derived datatypes.
- Derived datatypes can be built up from basic or derived datatypes.
- Datatype handles are used to describe the type of the data in the memory.

Example: message with 5 integers

2345	654	96574	-12	7676
------	-----	-------	-----	------

# MPI Basic Datatypes

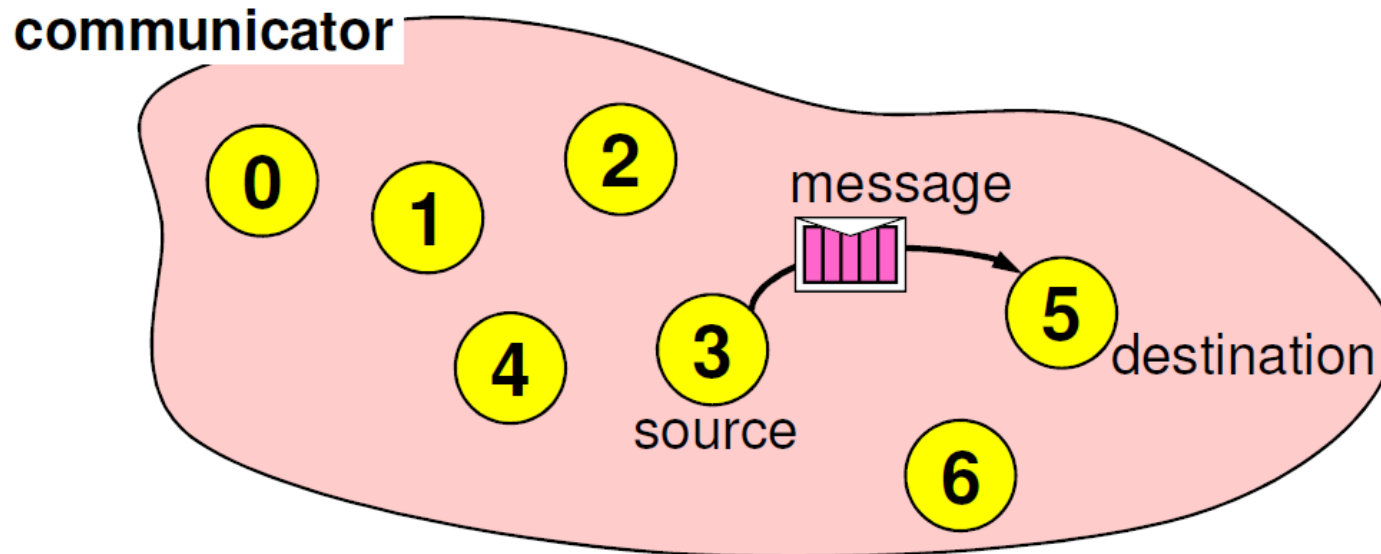
MPI Datatype	C datatype	Remarks
MPI_CHAR	char	Treated as printable character
MPI_SHORT	signed short int	
MPI_INT	signed int	
MPI_LONG	signed long int	
MPI_LONG_LONG	signed long long	
MPI_SIGNED_CHAR	signed char	Treated as integral value
MPI_UNSIGNED_CHAR	unsigned char	Treated as integral value
MPI_UNSIGNED_SHORT	unsigned short int	
MPI_UNSIGNED	unsigned int	
MPI_UNSIGNED_LONG	unsigned long int	
MPI_UNSIGNED_LONG_LONG	unsigned long long	
MPI_FLOAT	float	
MPI_DOUBLE	double	
MPI_LONG_DOUBLE	long double	
MPI_BYTE		
MPI_PACKED		

Further datatypes,  
see, e.g., MPI-3.0,  
Annex A.1

Includes also  
special C++ types,  
e.g., bool,  
see page 666

# Point-to-Point Communication

- Communication between two processes.
- Source process sends message to destination process.
- Communication takes place within a communicator, e.g., `MPI_COMM_WORLD`.
- Processes are identified by their ranks in the communicator.



# Sending a Message

```
int MPI_Send(    void*      buf,  
                int        count,  
                MPI_Datatype datatype,  
                int        dest,  
                int        tag,  
                MPI_Comm   comm)
```

- **buf** is the starting point of the message with **count** elements, each described with **datatype**.
- **dest** is the rank of the destination process within the communicator **comm**.
- **tag** is an additional nonnegative integer piggyback information, additionally transferred with the message.
- The tag can be used by the program to distinguish different types of messages.

# Receiving a Message

```
int MPI_Recv( void* buf,  
              int count,  
              MPI_Datatype datatype,  
              int source,  
              int tag,  
              MPI_Comm comm,  
              MPI_Status* status)
```

- *buf*/*count*/*datatype* describe the receive buffer.
- Receiving the message sent by process with rank *source* in *comm*.
- Envelope information is returned in *status*.
- One can pass `MPI_STATUS_IGNORE` instead of a status argument.
- Output arguments are printed *blue-cursive*.
- Only messages with matching *tag* are received.



# Requirements for Point-to-Point Communications

**For a communication to succeed:**

- **Sender must specify a valid destination rank.**
- **Receiver must specify a valid source rank.**
- **The communicator must be the same.**
- **Tags must match.**
- **Buffer's type must match with the datatype handle (in the send and receive call)**
- **Message datatypes must match.**
- **Receiver's buffer must be large enough.**

# Wildcarding

- Receiver can wildcard.
- To receive from any source:  
source = MPI\_ANY\_SOURCE
- To receive from any tag:  
tag = MPI\_ANY\_TAG
- Actual source and tag are returned in the receiver's *status* parameter.

# Communication Envelope

- Envelope information is returned from MPI\_RECV in *status*.

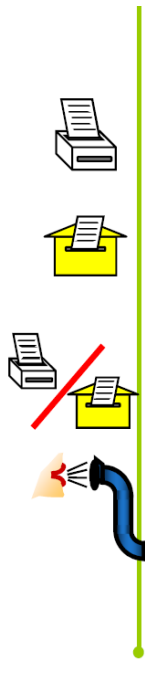
```
MPI_Status status;  
status.MPI_SOURCE  
status.MPI_TAG  
status.MPI_ERROR
```

- Receive Message Count

```
int MPI_Get_count(      MPI_Status*    status,  
                      MPI_Datatype    datatype,  
                      int*            count)
```

# Communication Modes

- **Send communication modes:**
  - synchronous send → **MPI\_SSEND**
  - buffered [asynchronous] send → **MPI\_BSEND**
  - standard send → **MPI\_SEND**
  - Ready send → **MPI\_RSEND**
- **Receiving all modes → **MPI\_RECV****



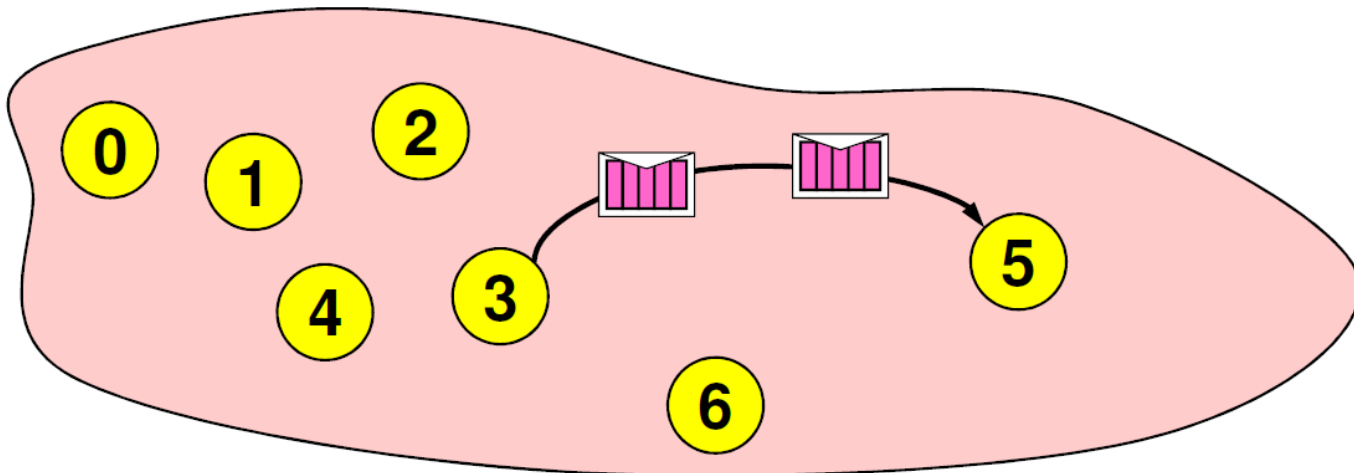
Sender mode	Definition	Notes
Synchronous send <b>MPI_SSEND</b>	Only completes when the receive has started	
Buffered send <b>MPI_BSEND</b>	Always completes (unless an error occurs), irrespective of receiver	needs application-defined buffer to be declared with <b>MPI_BUFFER_ATTACH</b>
Standard send <b>MPI_SEND</b>	Either synchronous or buffered	uses an internal buffer
Ready send <b>MPI_RSEND</b>	May be started <b>only</b> if the matching receive is already posted!	highly dangerous!
Receive <b>MPI_RECV</b>	Completes when a message has arrived	same routine for all communication modes

# Communication Modes Rules

- **Standard send (MPI\_SEND)**
  - minimal transfer time
  - may block due to synchronous mode
- **Synchronous send (MPI\_SSEND)**
  - risk of deadlock
  - risk of waiting → idle time
  - high latency / best bandwidth
- **Buffered send (MPI\_BSEND)**
  - low latency / bad bandwidth
- **Ready send (MPI\_RSEND)**
  - use never, except you have a 200% guarantee that Recv is already called in the current version and all future versions of your code
  - may be the fastest

# Message Order Preservation

- Rule for messages on the same connection,
- i.e., same communicator, source, and destination rank:
- Messages do not overtake each other.
- This is true even for non-synchronous sends.
- If both receives match both messages, then the order is preserved.



# Outlook Collective Communications

- **Two kinds**
  - Data movement (broadcast, scatter, gather, etc.)
  - Collective computation (min, max, sum, logical OR etc.)
- **Advantages**
  - More convenient
- **Can be used as short cut for an ensemble of P2P operations**
  - More efficient
- **Encapsulate sophisticated algorithms**
- **Implementation can take advantage of the structure of a machine to optimize and increase parallelism in these operations**

# Broadcast

- **Broadcasts a message from the process with rank root to all other processes of the group.**
  - **buf** = starting address of buffer
  - **count** = number of entries in buffer
  - **datatype** = data type of buffer
  - **root** = rank of broadcast root
  - **comm** = communicator

```
int MPI_Bcast( void* buf,  
               int count,  
               MPI_Datatype datatype,  
               int root,  
               MPI_Comm comm)
```



# Reduce

- **Combines the elements in the input buffer of each process using the operation op and returns the combined value in the output buffer of the process with rank root**
  - **sendbuf = address of send buffer**
  - **recvbuf = address of receive buffer**
  - **count = number of elements in send buffer**
  - **datatype = data type of elements of send buffer**
  - **op = reduce operation**
  - **root = rank of root process**
  - **comm = communicator**

```
int MPI_Reduce( void*          sendbuf,  
                void*          recvbuf,  
                int            count,  
                MPI_Datatype    datatype,  
                MPI_Op          op,  
                int            root,  
                MPI_Comm        comm)
```

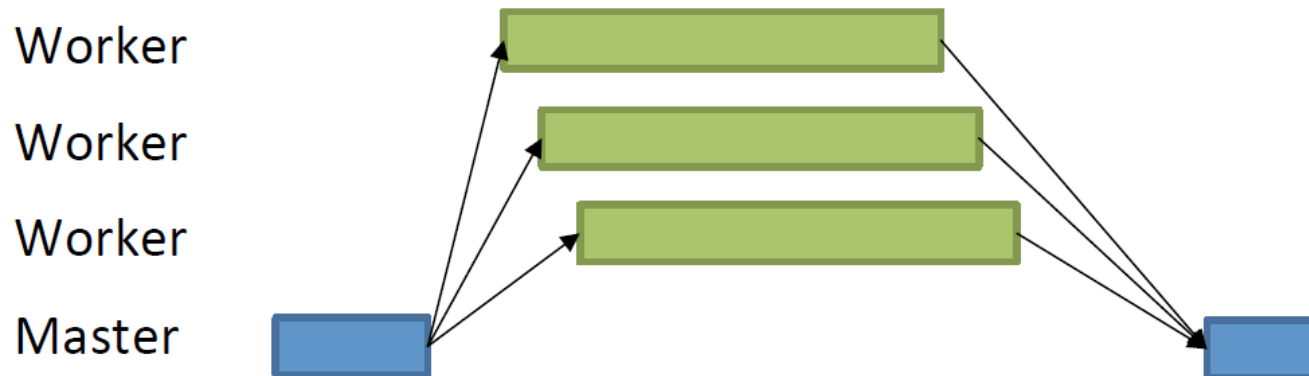
# Outline

---

- Introduction to the Lecture
- Distributed Parallel Computing I
  - Primer
  - Overview
  - Process Model and Language Bindings
  - Messages and Point-to-Point Communication
  - **Example**
- Quiz

# Master Worker

- **Idea: self-scheduling algorithm**
  - Master coordinates processing of tasks by providing input data to workers and collecting results
- **Suitable if**
  - Workers need not communicate with one another
  - Amount of work each worker must perform is difficult to predict
- **Example: matrix-vector multiplication**



# Matrix-Vector Multiplication

$$A \cdot \vec{b} = \vec{c}$$

**Unit of work:**

**dot product of one row of matrix A with vector b**

## **Master**

- **Broadcasts b to each worker**
- **Sends one row to each worker**
- **Loop**
  - **Receives dot product from whichever worker sends one**
  - **Sends next task to that worker**
  - **Termination if all tasks are handed out**

## **Worker**

- **Receives broadcast value of b**
- **while-Loop**
  - **Receives row from A**
  - **Forms dot product**
  - **Returns answer back to master**

# Matrix-Vector Multiplication: Common Part 1

```
#include <stdio.h>
#include <stdlib.h>
#include "mpi.h"
#define MAX_ROWS 1000 // upper limits
#define MAX_COLS 1000 // upper limits
#define MIN(a, b) ((a) > (b) ? (b) : (a))
#define DONE MAX_ROWS+1
int main(int argc, char **argv) {
    double** A;
    double* b, * c, * buffer;
    double ans;
    int myid, master, numprocs, i, j, numsent, sender, done, anstype, row, rows, cols;

    MPI_Status status;
    MPI_Init(&argc, &argv);
    MPI_Comm_rank(MPI_COMM_WORLD, &myid);
    MPI_Comm_size(MPI_COMM_WORLD, &numprocs);

    A = (double **)calloc(MAX_ROWS, sizeof(double *));
    for (i=0; i<MAX_ROWS; i++)
        A[i] = (double *)calloc(MAX_COLS, sizeof(double));
    b = (double *)calloc(MAX_COLS, sizeof(double));
    c = (double *)calloc(MAX_ROWS, sizeof(double));
    buffer = (double *)calloc(MAX_COLS, sizeof(double));
    ...
}
```

# Matrix-Vector Multiplication: Common Part 2

```
...
master = 0;
rows = 2; // for debugging ...
cols = 2; // for debugging ...
if (myid == master) { /* master code: next slides */
} else { /* worker code: next slides */
}
// result c vector finished: considering the test values, all element values should be 3.0
for (i=0; i<MAX_ROWS; i++)
    free(A[i]);
free(A);
free(b);
free(c);
free(buffer);
MPI_Finalize();
return 0;
}
```

# Matrix-Vector Multiplication: Master Part 1

```
...
if (myid == master) {
/* Initialize A and b (arbitrary) */
A[0][0] = 1.0;
A[0][1] = 2.0;
A[1][0] = 1.0;
A[1][1] = 2.0;
b[0] = 1.0;
b[1] = 1.0;

numsent = 0;
/* Send b to each worker process */
MPI_Bcast(b, cols, MPI_DOUBLE, master, MPI_COMM_WORLD);
/* Send a row to each worker process; tag with row number */
for (i = 0; i < MIN(numprocs - 1, rows); i++) {
    MPI_Send(&A[i][0], cols, MPI_DOUBLE, i+1, i, MPI_COMM_WORLD);
    numsent++;
}
...
```

## Master

- **Broadcasts b to each worker**
- **Sends one row to each worker**
- **Loop**
  - **Receives dot product from whichever worker sends one**
  - **Sends next task to that worker**
  - **Termination if all tasks are handed out**

# Matrix-Vector Multiplication: Master Part 2

```
...
for (i = 0; i < rows; i++) {
    MPI_Recv(&ans, 1, MPI_DOUBLE, MPI_ANY_SOURCE, MPI_ANY_TAG,
             MPI_COMM_WORLD, &status);
    sender = status.MPI_SOURCE;
    /* row is tag value */
    anstype = status.MPI_TAG;
    c[anstype] = ans;
    /* send another row */
    if (numsent < rows) {
        MPI_Send(&A[numsent][0], cols,
                 MPI_DOUBLE, sender,
                 numsent, MPI_COMM_WORLD);
        numsent++;
    } else {
        /* Tell sender that there is no more work */
        MPI_Send(MPI_BOTTOM, 0, MPI_DOUBLE, sender, DONE, MPI_COMM_WORLD);
    }
}
/* end of master specific part*/
else { /* start of worker code: next slide */
```

## Master

- Broadcasts **b** to each worker
- Sends one row to each worker
- **Loop**
  - Receives dot product from whichever worker sends one
  - Sends next task to that worker
  - Termination if all tasks are handed out



# Matrix-Vector Multiplication: Worker Part

```
MPI_Bcast(b, cols, MPI_DOUBLE, master, MPI_COMM_WORLD);
/* Skip if more processes than work */
done = myid > rows;
while (!done) {
    MPI_Recv(buffer, cols, MPI_DOUBLE, master, MPI_ANY_TAG,
             MPI_COMM_WORLD, &status);
    done = status.MPI_TAG == DONE;
    if (!done) {
        row = status.MPI_TAG;
        ans = 0.0;
        for (i = 0; i < cols; i++) {
            ans += buffer[i] * b[i];
        }
        MPI_Send(&ans, 1, MPI_DOUBLE, master, row, MPI_COMM_WORLD);
    }
}
/* end of worker code */
```

## Worker

- **Receives broadcast value of b**
- **while-Loop**
  - **Receives row from A**
  - **Forms dot product**
  - **Returns answer back to master**

# Outline

---

- Introduction to the Lecture
- Distributed Parallel Computing I
  - Primer
  - Overview
  - Process Model and Language Bindings
  - Messages and Point-to-Point Communication
  - Examples
- **Quiz**

# Quiz

- **Q1: How is it ensured that a specific message is received by a specific process?**
- **Q2: What is the first and last routine to be called in a MPI program?**
- **Q3: Is “MPI\_Init” executed by one, several or all MPI processes?**
- **Q4: Name typical reduction operations?**
- **Q5: How can a point-to-point communication be made non-blocking? What potential advantage is there?**