MPI Basics	Where to find help
Goal: Overview of MPI How to compile and run MPI programs Main MPI commands - examples See many other online resources available	 Many tutorials and other documents including the MPI standard available online The openMPI latest documentation https://www.open-mpi.org/doc/current/ The original MPI site: https://www.mcs.anl.gov/research/projects/mpi/index.html The LLNL link: https://computing.llnl.gov/tutorials/mpi/
9-1 Getting started: hello world from each process	9-2 mpi 9-2
<pre>#include <stdio.h> #include <mpi.h> int main(int argc, char *argv[]) { int nPEs, rank; MPI_Init(&argc, &argv); MPI_Comm_size(MPI_COMM_WORLD, &nPEs); MPI_Comm_rank(MPI_COMM_WORLD, &rank); printf("Hello from process %d out of %d \n", rank,nPEs); MPI_Finalize(); }</mpi.h></stdio.h></pre>	 #include <mpi.h> is mandatory: contains MPI definitions, types, etc,</mpi.h> MPI_Init(.,.) and MPI_Finalize() start and end MP respectively MPI_Comm_size gives the number of processes MPI_Comm_rank gives the id (proc. number) of this process Note: all non-MPI functions/ commands are local: e.g., print will run on each PE.

9-3

mpi

9-4

- mpi

Communicators

- ► A communicator defines a group of processes which communicate with each other.
- ▶ MPI_COMM_WORLD is the (default) communicator: an internal structure which includes all processes that are allowed to communicate with each other within this run.
- > Ranks, sizes, are all relative to a communicator.
- Can define other communicators than MPI_COMM_WORLD by command MPI_Comm.
- > Can work with communicators using various commands:

MPI_Comm_group, MPI_Group_excl, MPI_Comm_create,

The "phixx cluster"

- A small cluster of Xeon processors connected via a network
- Machines in the cluster: phi01, phi02, ···, phi08
- Login to any of these..
- System is dedicated to teaching parallel programming in CSE

Compiling, linking, and running a program

On the phixx cluster - a recent version of openmpi was installed (v. 4.0.3) - it is loaded by default.

9-5

- Instructions will be provided on how to load required modules, compiling, etc
- [One note: It is important that *all nodes* run the same version!]
- ► For compiling: MPI provides mpicc for C and mpif77 (fortran 77), mpif90 (fortran 90 when installed) for fortran.
- > For the "hello world" example:

mpicc -o hello.ex hello.c

9-7

For large programs use makefiles

Running a program

> in MPICH, and openMPI a program is run with mpirun:

mpirun -np 2 hello.ex

9-6

A few options

help: print mpi options

mpirun -h[elp]

> use filename as host-file:

mpirun -np 4 -machinefile < filename > executable.ex

Note you will need a machinefile on the phixx cluster. [a default one will be made available].
9-8

9-8

Remember: hardest aspect of parallel programming is: Debugging	Sends and Receives	
 Explore the use of MCA parameters (Modular Component Archi- tecture) for debugging (and more) in 	Simplest types of communication yet there are several options For example:	
https://www.open-mpi.org/faq/?category=debugging	* Synchronous send	
More information will be added as needed.	* Blocking send / blocking receive	
	* Non-blocking send / non-blocking receive	
	* Buffered send	
9-9 — mpi	9-10 — mpi	
9-9	9-10	

Blocking send will complete (return) only when the buffered data has been sent or saved and it is safe to free or reuse send buffer

Synchronus blocking send: Handshaking takes place between receiving and sending processes before actual send.

Asynchronus blocking send: data buffered until eventual delivery to receiving process

Blocking Receives complete (return) after the data has arrived to receiving process

Non-blocking sends and receives are similar - they return almost immediately. No waiting for communication to complete

9-11

► Non-blocking communications are useful when overlapping computation with communication

▶ Non-blocking operations "request" the MPI library to perform the operation when possible. Cannot predict when this takes place

➤ The application buffer (variable space, i.e., array to be sent) should not be changed until it is known that the requested non-blocking operation was performed

9-12

> Can use "wait" routines and MPI_Iprobe for this (see later)

Sends and Receives: Blocking send

MPI_Send(start,count,datatype,dest,tag,comm)

- **start** is a pointer to first entry of data to be sent
- **count** is the length of array
- > datatype one of MPI's datatypes: MPI_INT , MPI_FLOAT , MPI_DOUBLE , etc,...
- dest is the destination process
- comm is the communicator
- **tag** is a tag assigned to message so it is recognized by a matching receive.

There should be an associated receive. General form:

MPI_Recv(start,count,datatyp,source,tag,comm,status)

- > Start, count, datatype, comm, have the same meaning as in send
- source is the rank of sending process
- tag is the same as the tag used in the send
- Can also use: MPI_ANY_TAG and / or MPI_ANY_SOURCE
- status is a struct of type MPI_Status
- ► The source, tag, and count of the message actually received can be retrieved from status.

9-13	– mpi	9-14	— mpi
9-13			9-14

<code>MPI_TAG</code> , <code>MPI_SOURCE</code> <code>useful</code> in <code>case</code> when <code>MPI_ANY_TAG</code> and <code>/or</code> <code>MPI_ANY_SOURCE</code> <code>used</code> in the receive

MPI_Get_count will tell us how much data of a particular type was received [details shortly]

9-15

What is in the MPI_status struct? MPI_SOURCE - id of processor sending the message MPI_TAG - the message tag MPI_ERROR - error status MPI_LENGTH (Not accessible) MPI_COMM - communicator.

> Other members reserved for internal implementation.



0 10
- u - i u

9-19

– mpi

9-20

– mpi

MPI Collective communication routines

Allgather	Allgatherv	Allreduce
Alltoall	Alltoallv	Bcast
Gather	Gatherv	Reduce
ReduceScatter	Scan	Scatter
Scatterv		

➤ The 'v' versions: Scatterv, Gatherv, Allgatherv, Alltoallv, allow variable counts and shifts on the original array for each destination



> You will need to add all the machines into $\sim /.ssh/known_hosts$ file. For this you can ssh once to each machine. This will prevent prompts for confirmations each time you run something.

► The next thing you will need is a hostfile needed by MPI [see examples provided] Here is a sample hostfile:

```
phi01.cselabs.umn.edu slots=16
phi02.cselabs.umn.edu slots=16
phi03.cselabs.umn.edu slots=16
phi04.cselabs.umn.edu slots=16
phi05.cselabs.umn.edu slots=16
phi06.cselabs.umn.edu slots=16
phi07.cselabs.umn.edu slots=16
phi08.cselabs.umn.edu slots=16
```

- Call it hostfile or phi_cluster for example
- > Run with, e.g.,

mpirun -np 32 -hostfile phi_cluster test.ex

9-23

9-23

Appendix: set-up for the phi cluster

The nodes are called phi01, phi02, ..., phi08

You can login to any of them via ssh phi0x.cselabs.umn.edu

Latest version of openMPI loaded by default (no need for modules)

➤ One key requirement is that you should be able to login to any node without requiring a password. Find out how to use ssh-keygen and enable ssh without passwords:

```
ssh-keygen -t rsa -P ""
cat ~/.ssh/id_rsa.pub > ~/.ssh/authorized_keys
```

9-22

Appendix: A few details and an example

Some function calls

Send-Receive Pairs

```
#include <mpi.h>
int MPI_Send(const void *buf, int count,
    MPI_Datatype datatype, int dest,
    int tag, MPI_Comm comm)
int MPI_Recv(void *buf, int count,
    MPI_Datatype datatype, int source,
    int tag, MPI_Comm comm,
    MPI_Status *status)
```

Broadcast

Quoting from the documentation: (open-mpi.org):

MPI_Bcast broadcasts a message from the process with rank root to all processes of the group, itself included. It is called by all members of group using the same arguments for comm, root. On return, the contents of root's communication buffer has been copied to all processes.

▶ Main point: no need to have a corresponding receive - but all nodes in *comm* must issue the same command (the communication is "collective")

► What are MPI Data types?

MPI Datatypes

<i>MPI_CHAR</i>	MPI_C_COMPLEX
<i>MPI_WCHAR</i>	MPI_C_FLOAT_COMPLEX
MPI_SHORT	MPI_C_DOUBLE_COMPLEX
MPI_INT	<i>MPI_C_LONG_DOUBLE_COMPLEX</i>
<i>MPI_LONG</i>	MPI_C_BOOL
MPI_LONG_LONG_INT	MPI_LOGICAL
MPI_LONG_LONG	<i>MPI_C_LONG_DOUBLE_COMPLEX</i>
MPI_SIGNED_CHAR	MPI_INT8_T, MPI_INT16_T
MPI_UNSIGNED_CHAR	MPI_INT32_T
<i>MPI_UNSIGNED_SHORT</i>	MPI_INT64_T
MPI_UNSIGNED_LONG	MPI_UINT8_T, MPI_UINT16_T
<i>MPI_UNSIGNED</i>	MPI_UINT32_T
<i>MPI_FLOAT</i>	MPI_UINT64_T
<i>MPI_DOUBLE</i>	MPI_BYTE
MPI_LONG_DOUBLE	<i>MPI_PACKED</i>

9-26

9-26

Example: Compute values of e^x

▶ No practical value: Only goal is to illustrate a pipelined computation and sends, receives, broadcasts.

9-25

> We will compute exp(x) by its *n*-th Taylor series expansion and use Horner's scheme:

 $e^x \approx 1 + \frac{x}{1} \left(1 + \frac{x}{2} \left(1 + \frac{x}{3} \left(\cdots \left(1 + \frac{x}{n} \right) \cdots \right) \right) \right)$

9-27

> In Python this can be implemented as:

```
x = some-value
n = some value
val = 1.0
for i in range(n,0,-1):
    val = 1.0+(x/i)*val
print(val)
```

General procedure: Think of a linear array. For n = 5:

* Start: proc. n-1 (last one) who sends

- value val = 1 + (x/n) to West processor
- * All other nodes: get val from East; update; and:
- * Send val West *unless* MyId=0 in which case print.
- \blacktriangleright At first assume x is set in each process
- \blacktriangleright Next: Assume x is also sent along with val
- \blacktriangleright Next: Use a broadcast of x instead

> Next: In the above only one process is active at a time. What if I need to compute a bunch of values $e^{x_0}, e^{x_1}, \cdots e^{x_k}$

9-28

9-27

- mpi

– mp

