



WWW.TACC.UTEXAS.EDU



Tutorial on MPI programming Basic material

Victor Eijkhout eijkhout@tacc.utexas.edu
SDS PCSE 2023

Materials

Textbooks and repositories:

<https://theartofhpc.com>

Justification

The MPI library is the main tool for parallel programming on a large scale.
This course introduces the main concepts through lecturing and exercises.

Table of Contents

- The SPMD model 6
- Collectives 64
- Point-to-point communication 134

Basics

Part I

The SPMD model

1. Overview

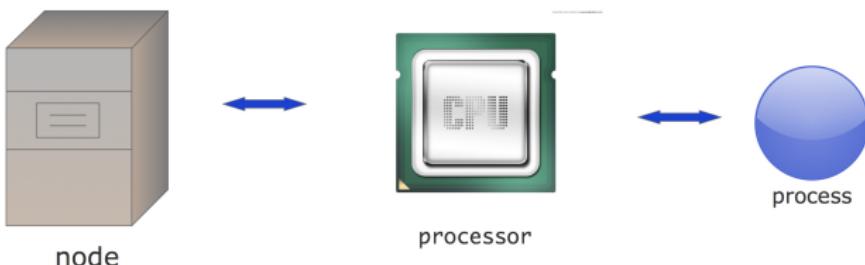
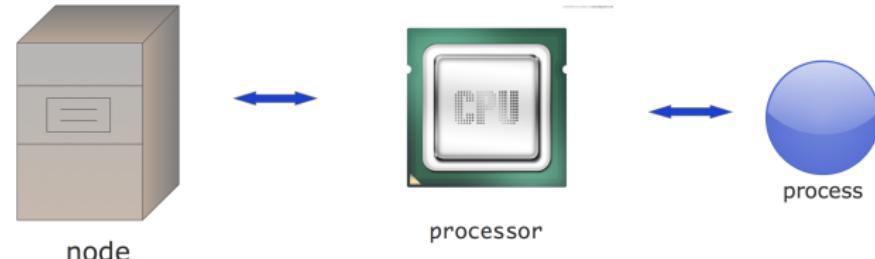
In this section you will learn how to think about parallelism in MPI.

Commands learned:

- `MPI_Init`, `MPI_Finalize`,
- `MPI_Comm_size`, `MPI_Comm_rank`
- `MPI_Get_processor_name`,

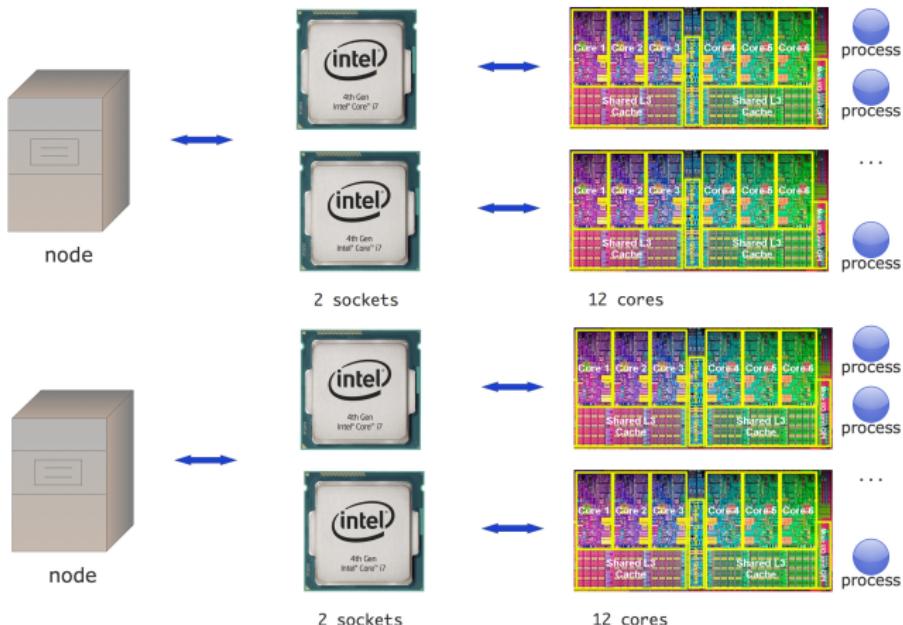
The MPI worldview: SPMD

2. Computers when MPI was designed



One processor and one process per node;
all communication goes through the network.

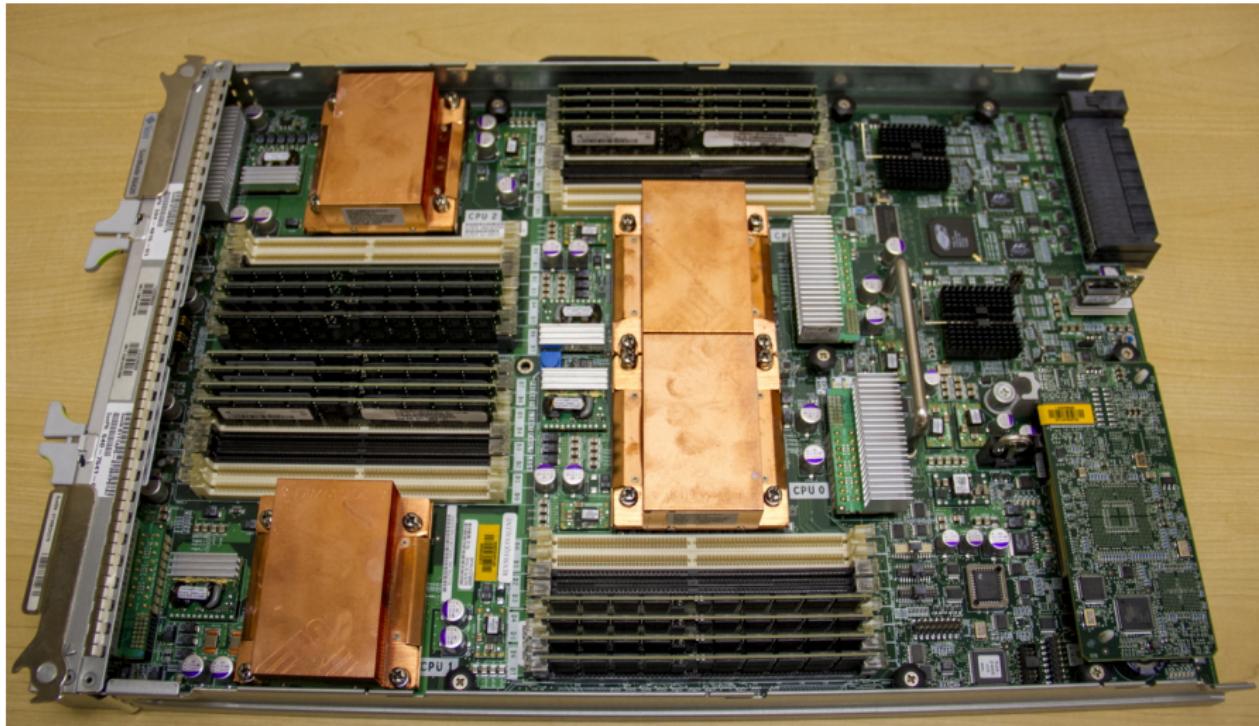
3. Pure MPI



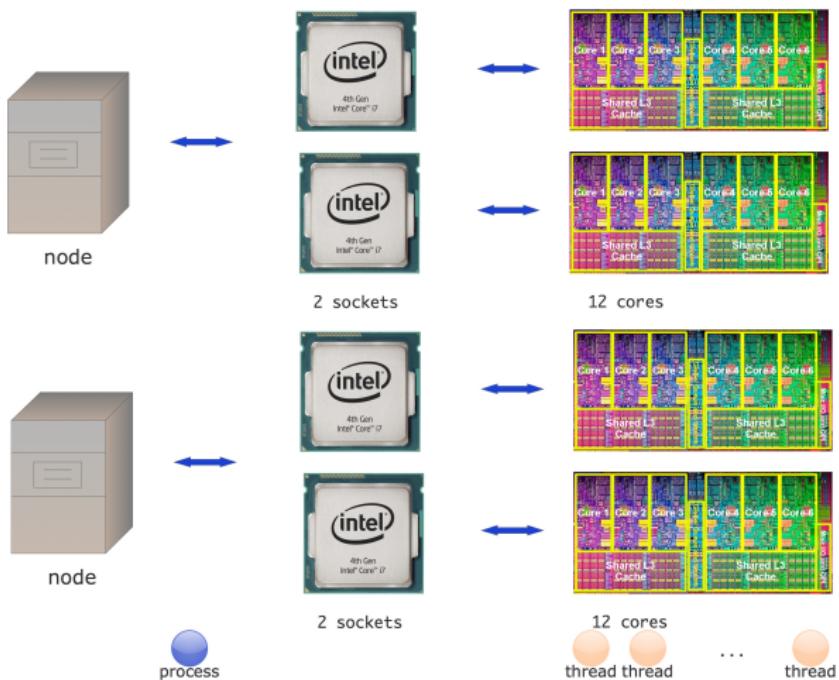
A node has multiple sockets, each with multiple cores.

Pure MPI puts a process on each core: pretend shared memory doesn't exist.

4. Quad socket node



5. Hybrid programming



Hybrid programming puts a process per node or per socket;
further parallelism comes from threading.
Not in this course...

6. Terminology

'Processor' is ambiguous: is that a chip or one independent instruction processing unit?

- Socket: the processor chip
- Processor: we don't use that word
- Core: one instruction-stream processing unit
- Process: preferred terminology in talking about MPI.

7. SPMD

The basic model of MPI is
'Single Program Multiple Data':
each process is an instance of the same program.

Symmetry: There is no 'master process', all processes are equal, start and end at the same time.

Communication calls do not see the cluster structure:
data sending/receiving is the same for all neighbors.

Let's experience the SPMD model

8. Lab setup

- Open a terminal window on a TACC cluster.
- Type `idev -N 2 -n 10 -t 2:0:0` which gives you an interactive session of 2 nodes, 10 cores, for the next 2 hours.
- (After this course, for serious work, you would write a batch script. The idev sessions are strictly limited in time and resources.)
- See external handouts for reservations, project IDs, and location of training materials.
- Next slide for how to make and run exercises.

9. How to make exercises

- Directory: `exercises-mpi-c` or `cxx` or `f` or `f08` or `p` or `mpl`
- If a slide has a `(exercisename)` over it, there will be a template program `exercisename.c` (or `F90` or `py`).
- Type `make exercisename` to compile it
- Run with `ibrun` or `mpexec` (see above)
- Python: setup once per session

```
module load python3
```

No compilation needed. Run:

```
ibrun python3 yourprogram.py
```

- Add an exercise of your own to the `makefile`: add the name to the `EXERCISES`

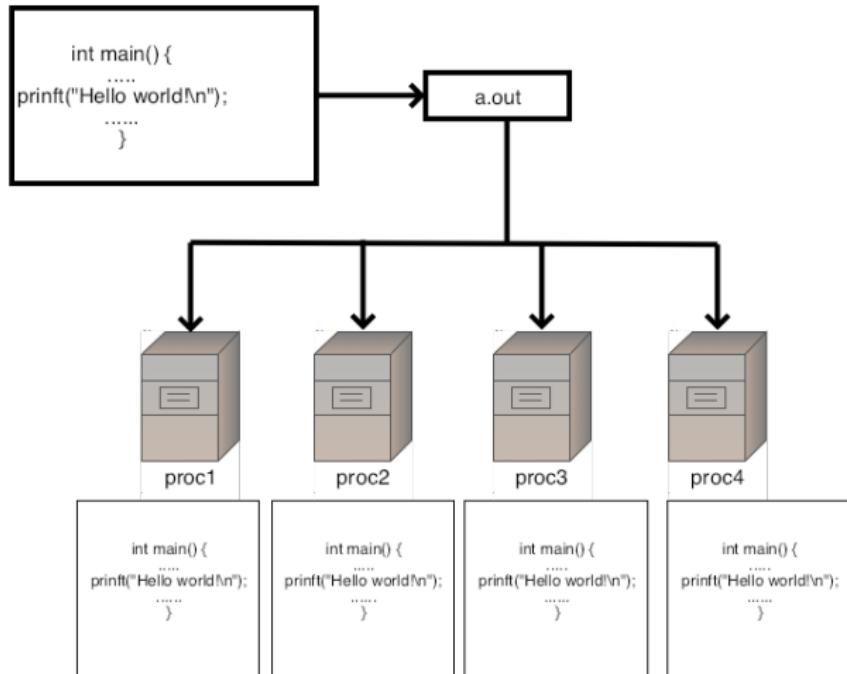
Exercise 1 (hello)

Write a ‘hello world’ program, without any MPI in it, and run it in parallel with `mpiexec` or your local equivalent. Explain the output.

(On TACC machines such as stampede, use `ibrun`, no processor count.)

(In the directories `exercises-mpi-xxx` you can do `make hello` to compile, then `ibrun hello` to execute.)

10. In a picture



Practicalities

11. Compiling and running

MPI compilers are usually called `mpicc`, `mpif90`, `mpicxx`.

These are not separate compilers, but scripts around the regular C/Fortran compiler. You can use all the usual flags.

```
$ mpicc -show  
icc -I/intel/include/stuff -L/intel/lib/stuff -Wwarnings # et
```

Running your program at TACC:

```
#SBATCH -N 4  
#SBATCH -n 200  
ibrun yourprog
```

the number of processes is determined by SLURM. General case of running code

```
mpiexec -n 4 hostfile ... yourprogram  
mpirun -np 4 hostfile ... yourprogram
```

12. CMake for C

```
1  cmake_minimum_required( VERSION 3.12 )
2  project( ${PROJECT_NAME} VERSION 1.0 )
3
4  # https://cmake.org/cmake/help/latest/module/FindMPI.html
5  find_package( MPI )
6
7  add_executable( ${PROJECT_NAME} ${PROJECT_NAME}.c )
8  target_include_directories(
9      ${PROJECT_NAME} PUBLIC
10     ${MPI_C_INCLUDE_DIRS} ${CMAKE_CURRENT_SOURCE_DIR} )
11 target_link_libraries(
12     ${PROJECT_NAME} PUBLIC
13     ${MPI_C_LIBRARIES} )
14
15 install( TARGETS ${PROJECT_NAME} DESTINATION . )
```

13. CMake for C++

```
1  cmake_minimum_required( VERSION 3.12 )
2  project( ${PROJECT_NAME} VERSION 1.0 )
3
4  # https://cmake.org/cmake/help/latest/module/FindMPI.html
5  find_package( MPI )
6
7  add_executable( ${PROJECT_NAME} ${PROJECT_NAME}.cxx )
8  target_compile_features( ${PROJECT_NAME} PRIVATE cxx_std_20 )
9  target_include_directories(
10      ${PROJECT_NAME} PUBLIC
11      ${MPI_CXX_INCLUDE_DIRS} ${CMAKE_CURRENT_SOURCE_DIR} )
12 target_link_libraries(
13      ${PROJECT_NAME} PUBLIC
14      ${MPI_CXX_LIBRARIES} )
15
16 install( TARGETS ${PROJECT_NAME} DESTINATION . )
```

14. CMake for F90

```
1  cmake_minimum_required( VERSION 3.12 )
2  project( ${PROJECT_NAME} VERSION 1.0 )
3
4  enable_language(Fortran)
5
6  # https://cmake.org/cmake/help/latest/module/FindMPI.html
7  find_package( MPI )
8
9  add_executable( ${PROJECT_NAME} ${PROJECT_NAME}.F90 )
10 target_include_directories(
11     ${PROJECT_NAME} PUBLIC
12     ${MPI_INCLUDE_DIRS} ${CMAKE_CURRENT_SOURCE_DIR} )
13 target_link_directories(
14     ${PROJECT_NAME} PUBLIC
15     ${MPI_LIBRARY_DIRS} )
16 target_link_libraries(
17     ${PROJECT_NAME} PUBLIC
18     ${MPI_Fortran_LIBRARIES} )
19
20 install( TARGETS ${PROJECT_NAME} DESTINATION . )
```

15. CMake for F2008

```
1 cmake_minimum_required( VERSION 3.12 )
2 project( ${PROJECT_NAME} VERSION 1.0 )
3
4 enable_language(Fortran)
5
6 # https://cmake.org/cmake/help/latest/module/FindMPI.html
7 find_package( MPI )
8
9 if( MPI_Fortran_HAVE_F08_MODULE )
10 else()
11     message( FATAL_ERROR "No f08 module for this MPI" )
12 endif()
13
14 add_executable( ${PROJECT_NAME} ${PROJECT_NAME}.F90 )
15 target_include_directories(
16     ${PROJECT_NAME} PUBLIC
17     ${MPI_Fortran_INCLUDE_DIRS} ${CMAKE_CURRENT_SOURCE_DIR} )
18 target_link_directories(
19     ${PROJECT_NAME} PUBLIC
20     ${MPI_LIBRARY_DIRS} )
21 target_link_libraries(
22     ${PROJECT_NAME} PUBLIC
```

16. Do I need a supercomputer?

- With `mpiexec` and such, you start a bunch of processes that execute your MPI program.
- Does that mean that you need a cluster or a big multicore?
- No! You can start a large number of MPI processes, even on your laptop. The OS will use 'time slicing'.
- Of course it will not be very efficient...

17. Installing your own MPI

It is convenient to do MPI development on your laptop/desktop.

- Use a package manager

Apple: brew or macports

Linux: yum, aptget, ...

Windows: I'll have to get back to you on that

- ... or download and compile from source mpich.org

18. Cluster setup

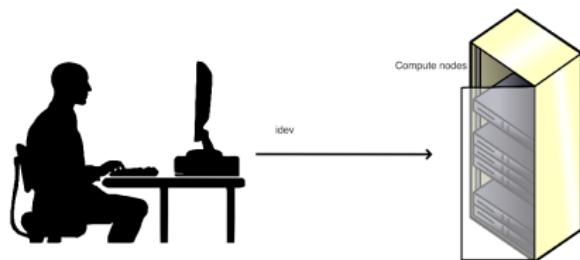
Typical cluster:

- Login nodes, where you ssh into; usually shared with 100 (or so) other people. You don't run your parallel program there!
- Compute nodes: where your job is run. They are often exclusive to you: no other users getting in the way of your program.

Hostfile: the description of where your job runs. Usually generated by a *job scheduler*.

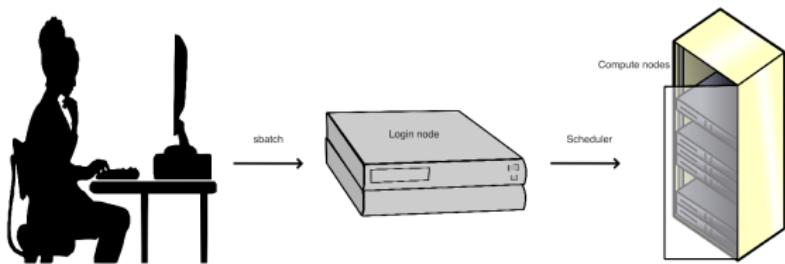
19. Interactive run

- Do not run your programs on a login node.
- Acquire compute nodes with `idev`
(other systems: `qsub -I`)
- Caveat: only small short jobs; nodes may not be available.



20. Batch run

- Submit batch job with `sbatch` or `qsub`
- Your job will be executed ... Real Soon Now.
- See userguide for details about queues, sizes, runtimes, ...



We start learning MPI!

21. MPI Init / Finalize

These calls need to be around the MPI part of your code:

```
1 MPI_Init(&argc,&argv); // zeros allowed  
2 // your code  
3 MPI_Finalize();
```

This is not a 'parallel region':
only internal library initialization:
allocate buffers, discover network, ...

21. Init Finalize, Fortran

No possibility for commandline arguments:

```
1 call MPI_Init()      ! F08 style
2 ! your code
3 call MPI_Finalize()
4
5 call MPI_Init(ierr) ! F90 style
6 ! your code
7 call MPI_Finalize(ierr)
```

21. MPL init/finalize

No explicit init/finalize:

- init is done by the first command that needs it
- finalize in some destructor.

21. Python init/finalize

Done by the import / at end of the program.

Exercise 2 (hello)

Add the commands `MPI_Init` and `MPI_Finalize` to your code. Put three different print statements in your code: one before the init, one between init and finalize, and one after the finalize. Again explain the output.

Run your program on a large scale, using a batch job. Where does the output go? Experiment with

```
MY_MPIRUN_OPTIONS="-prepend-rank" ibrun yourprogram
```

Python, MPL: include magic line from slide 53.

About library calls and bindings

22. Bindings

The standard defines interfaces to MPI from C and Fortran.
These look very similar; sometimes we will only show the C variant.

MPI can also be used from C++ and Python

23. MPI headers: C

You need an include file:

```
#include "mpi.h"
```

This defines all routines and constants.

23. MPI headers: Fortran

- In the standard
- Come in two flavors, Fortran90 vs Fortran2008
you will find many examples online in old style
- We teach you modern style.

You need an include file:

```
! Modern Fortran2008
use mpi_f08
! Legacy Fortran90
use mpi
! Deprecated
#include "mpif.h"
```

23. C++ bindings, MPL

MPI-1 had C++ bindings, by MPI-2 they were deprecated, in MPI-3 they have been removed.

- Easy solution: use the C bindings unaltered.
This is done in the `cxx` exercise directory; Ugly: very un-OO.
- There are private projects for C++ bindings.
In particular MPL: <https://github.com/rabauke/mp1> (Exercises in `mp1` directory.)
Very modern OO, Header-only
Not a full MPI implementation: (I/O and one-sided mostly missing)

In your program:

```
1      #include <mpl/mpl.hpp>
2
```

Compiling:

```
mpicxx -o prog sources.cxx -I${TACC_}
```

23. Python bindings

- Not part of the standard:
private project by Lisando Dalcin
Download <https://github.com/mpi4py/mpi4py>
Docs: <https://mpi4py.readthedocs.io/>
- Comes in two variants:
'pythonic' vs efficient

You need an include file:

```
from mpi4py import MPI
```

You need a python with MPI support
at TACC: module load python3

24. About errors

MPI routines invoke an error handler (slide ??)

default action: abort

Every routine is defined as returning integer error code

- In C: function result.

```
1 ierr = MPI_Init(0,0);  
2 if (ierr!=MPI_SUCCESS) /* do something */
```

But really: can often be ignored; is ignored in this course.

```
1 MPI_Init(0,0);
```

- In Fortran: as optional (F08 only) parameter.
- In Python: throwing exception.

There's not a lot you can do with an error code:

very hard to recover from errors in parallel.

By default code bombs with (hopefully informative) message.

Ranks

25. Process identification

- Processors are organized in ‘communicators’.
- For now only the ‘world’ communicator
- Each process has a unique ‘rank’ wrt the communicator.

```
1 int MPI_Comm_size( MPI_Comm comm, int *nprocs )  
2 int MPI_Comm_rank( MPI_Comm comm, int *procno )
```

Lowest number is always zero.

This is a logical view of parallelism: mapping to physical processors/cores is invisible here.

Exercise 3 (commrank)

Write a program where each process prints out a message reporting its number, and how many processes there are:

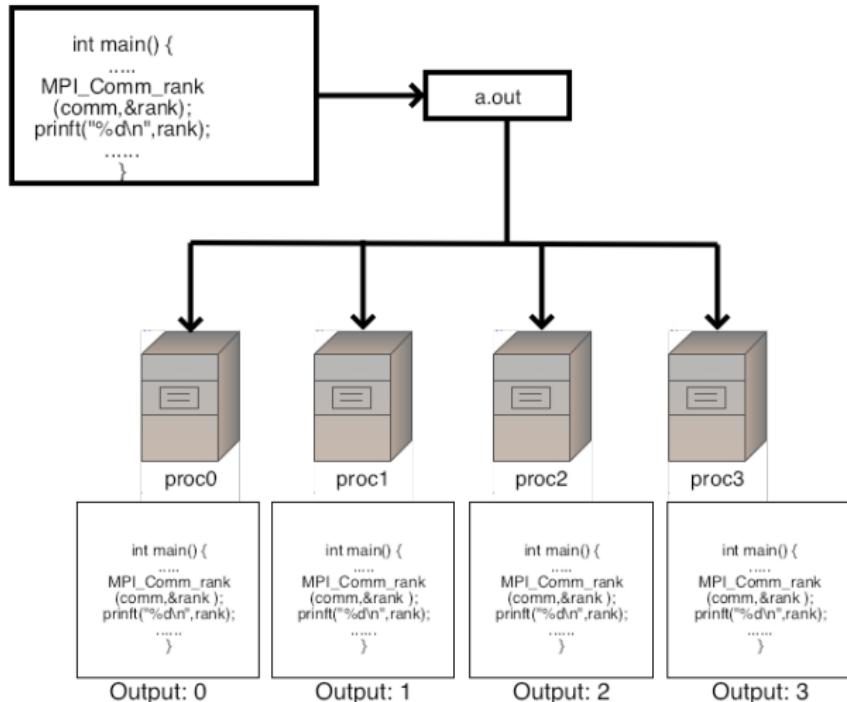
Hello from process 2 out of 5!

Write a second version of this program, where each process opens a unique file and writes to it. *On some clusters this may not be advisable if you have large numbers of processors, since it can overload the file system.*

Exercise 4 (commrank)

Write a program where only the process with number zero reports on how many processes there are in total.

26. Illustration



27. About routine signatures: C/C++

Signature:

```
1 int MPI_Comm_size(MPI_Comm comm, int *nprocs)
```

Use:

```
1 MPI_Comm comm = MPI_COMM_WORLD;
2 int nprocs;
3 int errorcode;
4 errorcode = MPI_Comm_size( comm,&nprocs );
```

(but forget about that error code most of the time)

27. About routine signatures: Fortran2008

Signature

```
1 MPI_Comm_size(comm, size, ierror)
2 Type(MPI_Comm), INTENT(IN) :: comm
3 INTEGER, INTENT(OUT) :: size
4 INTEGER, OPTIONAL, INTENT(OUT) :: ierror
```

Use:

```
1 Type(MPI_Comm) :: comm = MPI_COMM_WORLD
2 integer :: size
3 CALL MPI_Comm_size( comm, size ) ! F2008 style
```

- final parameter optional.
- MPI_... types are Type.

27b. About routine signatures: Fortran90

Signature

```
1 MPI_Comm_size(comm, size, ierror)
2 Integer, Intent(in) :: comm
3 Integer, Intent(out) :: ierror
```

Use:

```
1 Integer :: comm = MPI_COMM_WORLD
2 Integer :: size,ierr
3 CALL MPI_Comm_size( comm, size, ierr ) ! F90 style
```

- Final parameter always error parameter. Do not forget!
- MPI_... types are INTEGER.

27. About routine signatures: Python

Signature:

```
1 # object method
2 MPI.Comm.Send(self, buf, int dest, int tag=0)
3 # class method
4 MPI.Request.Waitall(type cls, requests, statuses=None)
```

Use:

```
1 from mpi4py import MPI
2 comm = MPI.COMM_WORLD
3 comm.Send(sendbuf, dest=other)
4 MPI.Request.Waitall(requests)
```

Note: most functions are methods of the *MPI.Comm* class.
(Sometimes of *MPI*, sometimes other.)

28. Communicators

For now, the communicator will be `MPI_COMM_WORLD`.

C:

```
1     MPI_Comm comm = MPI_COMM_WORLD;  
2
```

F:

```
1     Type(MPI_Comm) :: comm = MPI_COMM_WORLD  
2
```

P:

```
1     from mpi4py import MPI  
2     comm = MPI.COMM_WORLD  
3
```

MPL:

```
1     const mpl::communicator &comm_world =  
2         mpl::environment::comm_world();  
3
```

29. Have you been paying attention?

T/F?

- ① In C, the result of `MPI_Comm_rank` is a number from zero to number-of-processes-minus-one, inclusive.
- ② In Fortran, the result of `MPI_Comm_rank` is a number from one to number-of-processes, inclusive.

30. Processor name

Processes (can) run on physically distinct locations.

```
1 // procname.c
2 int name_length = MPI_MAX_PROCESSOR_NAME;
3 char proc_name[name_length];
4 MPI_Get_processor_name(proc_name,&name_length);
5 printf("Process %d/%d is running on node <<%s>>\n",
6 procid,nprocs,proc_name);
```

31. In a picture

Four processes on two nodes (`idev -N 2 -n 4`)

```
Program:  
number <- MPI_Comm_rank  
  
name <- MPI_Get_processor_name
```

```
Program:  
number <- MPI_Comm_rank  
0  
name <- MPI_Get_processor_name  
c111.tacc.utexas.edu
```

```
Program:  
number <- MPI_Comm_rank  
1  
name <- MPI_Get_processor_name  
c111.tacc.utexas.edu
```

```
Program:  
number <- MPI_Comm_rank  
2  
name <- MPI_Get_processor_name  
c222.tacc.utexas.edu
```

```
Program:  
number <- MPI_Comm_rank  
3  
name <- MPI_Get_processor_name  
c222.tacc.utexas.edu
```

c111.tacc.utexas.edu

c222.tacc.utexas.edu

MPI_Get_processor_name

Name	Param name	Explanation	C type	F type	inc
MPI_Get_processor_name (
name		A unique specifier for the actual (as opposed to virtual) node.	char*	CHARACTER	
resultlen		Length (in printable characters) of the result returned in name	int*	INTEGER	
)					

MPL: MPI_Get_processor_name

Missing MPL proto: mpi_get_processor_name

Exercise 5

Use the command `MPI_Get_processor_name`. Confirm that you are able to run a program that uses two different nodes.

TACC nodes have a hostname cRRR-CNN, where RRR is the rack number, C is the chassis number in the rack, and NN is the node number within the chassis. Communication is faster inside a rack than between racks!

Go to `examples/mpi/c` and do `make procname`, then `ibrun procname`

Your first useful parallel program

32. Functional Parallelism

Parallelism by letting each process do a different thing.

Example: divide up a search space.

Each process knows its rank, so it can find its part of the search space.

Exercise 6 (prime)

Is the number $N = 2,000,000,111$ prime? Let each process test a disjoint set of integers, and print out any factor they find. You don't have to test all integers $< N$: any factor is at most $\sqrt{N} \approx 45,200$.

(Hint: $i\%0$ probably gives a runtime error.)

Can you find more than one solution?

Exercise 7

Allocate on each process an array:

```
1 int my_ints[10];
```

and fill it so that process 0 has the integers $0 \dots 9$, process 1 has $10 \dots 19$, et cetera.

It may be hard to print the output in a non-messy way.

Part II

Collectives

33. Overview

In this section you will learn ‘collective’ operations, that combine information from all processes.

Commands learned:

- `MPI_Bcast`, `MPI_Reduce`, `MPI_Gather`, `MPI_Scatter`
- `MPI_All_...` variants, `MPI_....v` variants
- `MPI_Barrier`, `MPI_Alltoall`, `MPI_Scan`

34. Technically

Routines can be ‘collective on a communicator’:

- They involve a communicator;
- if one process calls that routine, every process in that communicator needs to call it
- Mostly about combining data, but also opening shared files, declaring ‘windows’ for one-sided communication.

Concepts

35. Collectives

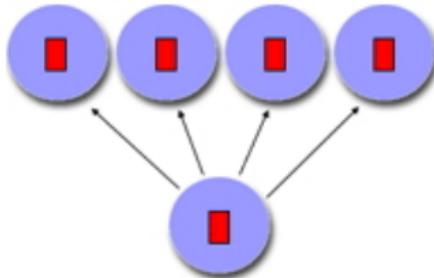
Gathering and spreading information:

- Every process has data, you want to bring it together;
- One process has data, you want to spread it around.

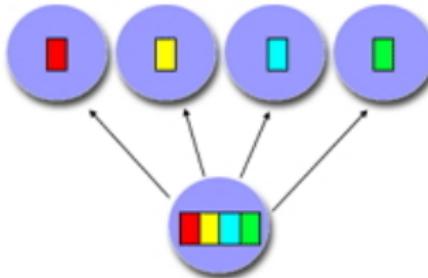
Root process: the one doing the collecting or disseminating.

Basic cases:

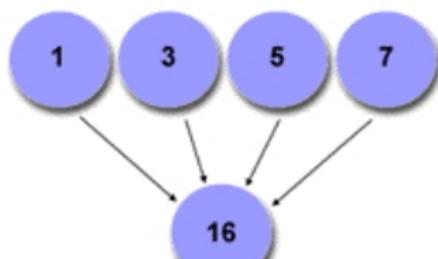
- Collect data: gather.
- Collect data and compute some overall value (sum, max): reduction.
- Send the same data to everyone: broadcast.
- Send individual data to each process: scatter.



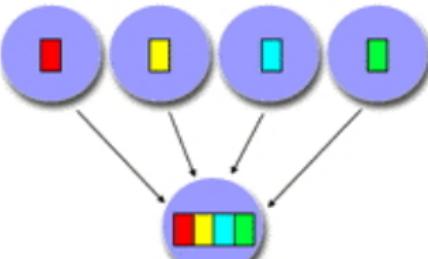
broadcast



scatter



reduction



gather

Exercise 8

How would you realize the following scenarios with MPI collectives?

- ① Let each process compute a random number. You want to print the maximum of these numbers to your screen.
- ② Each process computes a random number again. Now you want to scale these numbers by their maximum.
- ③ Let each process compute a random number. You want to print on what processor the maximum value is computed.

Think about time and space complexity of your suggestions.

36. Allreduce: reduce-to-all

Exercise 2 above contains a common case:
do a reduction, but everyone needs the result.

- `MPI_Allreduce` does the same as:
`MPI_Reduce` (reduction) followed by `MPI_Bcast` (broadcast)
- Same running time as either, half of reduce-followed-by-broadcast
(no proof given here)
- Common use case, symmetrical expression.

37. Allreduce syntax

```
1 int MPI_Allreduce(  
2     const void* sendbuf,  
3     void* recvbuf, int count, MPI_Datatype datatype,  
4     MPI_Op op, MPI_Comm comm)
```

- All processes have send and recv buffer
- (No root argument)
- *count* is number of items in the buffer: 1 for scalar.
 > 1: pointwise application of the reduction operator
- *MPI_Datatype* is *MPI_INT*, *MPI_FLOAT*, *MPI_REAL8* et cetera.
- *MPI_Op* is *MPI_SUM*, *MPI_MAX* et cetera.

Exercise 9 (randommax)

Let each process compute a random number, and compute the sum of these numbers using the `MPI_Allreduce` routine.

$$\xi = \sum_i x_i$$

Each process then scales its value by this sum.

$$x'_i \leftarrow x_i / \xi$$

Compute the sum of the scaled numbers

$$\xi' = \sum_i x'_i$$

and check that it is 1.

38. Motivation for allreduce

Example: normalizing a vector

$$y \leftarrow x / \|x\|$$

- Vectors x, y are distributed: every process has certain elements
- The norm calculation is an all-reduce: every process gets same value
- Every process scales its part of the vector.
- Question: what kind of reduction do you use for an inf-norm?
One-norm? Two-norm?

39. Another Allreduce

Standard deviation:

$$\sigma = \sqrt{\frac{1}{N} \sum_i^N (x_i - \mu)^2} \quad \text{where} \quad \mu = \frac{\sum_i^N x_i}{N}$$

and assume that every processor stores just one x_i value.

How do we compute this?

- ① The calculation of the average μ is a reduction.
- ② Every process needs to compute $x_i - \mu$ for its value x_i , so use allreduce operation, which does the reduction and leaves the result on all processors.
- ③ $\sum_i(x_i - \mu)$ is another sum of distributed data, so we need another reduction operation. Might as well use allreduce.

MPI_Allreduce

Name	Param name	Explanation	C type	F type	inc
MPI_Allreduce (
sendbuf		starting address of send buffer	const void*	TYPE(*), DIMENSION(..)	IN
recvbuf		starting address of receive buffer	void*	TYPE(*), DIMENSION(..)	OUT
count		number of elements in send buffer	int	INTEGER	IN
datatype		datatype of elements of send buffer	MPI_Datatype	TYPE(MPI_Datatype)	IN
op		operation	MPI_Op	TYPE(MPI_Op)	IN
comm		communicator	MPI_Comm	TYPE(MPI_Comm)	IN
)					

Basic collectives

40. Elementary datatypes

C	Fortran	meaning
MPI_CHAR	MPI_CHARACTER	only for text
MPI_SHORT	MPI_BYTE	8 bits
MPI_INT	MPI_INTEGER	like the C/F types
MPI_FLOAT	MPI_REAL	
MPI_DOUBLE	MPI_DOUBLE_PRECISION MPI_COMPLEX MPI_LOGICAL	
unsigned	extensions	
		MPI_Aint MPI_Offset

A bunch more.

41. MPI operators

MPI_Op	description
MPI_MAX	maximum
MPI_MIN	minimum
MPI_SUM	sum
MPI_PROD	product
MPI_LAND	logical and
MPI_BAND	bitwise and
MPI_LOR	logical or
MPI_BOR	bitwise or
MPI_LXOR	logical xor
MPI_BXOR	bitwise xor
MPI_MAXLOC	location of max
MPI_MINLOC	location of min

A couple more.

42. Buffers in C

General principle: buffer argument is address in memory of the data.

- Buffer is void pointer:
- write `&x` or `(void*)&x` for scalar
- write `x` or `(void*)x` for array

42. Buffers in Fortran

General principle: buffer is address in memory of the data.

- Fortran always passes by reference:
- write x for scalar
- write x for array

43. Buffers in C++

- Scalars same as in C.
- Use of `std::vector` or `std::array`:

```
1  vector<float> xx(25);
2  MPI_Send( xx.data(),25,MPI_FLOAT, .... );
3  MPI_Send( &xx[0],25,MPI_FLOAT, .... );
4  MPI_Send( &xx.front(),25,MPI_FLOAT, .... );
```

- Can not send from iterator / let recv determine size/capacity.

44. Large buffers

As of MPI-4 a buffer can be longer than 2^{31} elements.

- Use `MPI_Count` for count
- In C: use `MPI_Reduce_c`
- in Fortran: polymorphism means no change to the call.

```
1 MPI_Count buffersize = 1000;  
2 double *indata,*outdata;  
3 indata = (double*) malloc( buffersize*sizeof(double) );  
4 outdata = (double*) malloc( buffersize*sizeof(double) );  
5 MPI_Allreduce_c(indata,outdata,buffersize,  
6                   MPI_DOUBLE,MPI_SUM,MPI_COMM_WORLD);
```

44. Buffers in Python

For many routines there are two variants:

- lowercase: can send Python objects;
output is *return* result

```
result = comm.recv(...)
```

this uses pickle: slow.

- uppercase: communicates numpy objects;
input and output are function argument.

```
result = np.empty(.....)
```

```
comm.Recv(result, ...)
```

basically wrapper around C code: fast

45. Inner product calculation

Given vectors x, y :

$$x^t y = \sum_{i=0}^{N-1} x_i y_i$$

Start out with distributed vectors x, y ,
assume same distribution.

Proposed solution:

`MPI_Gather` or `MPI_Allgather` and calculate locally.

Comments?

46. Inner product calculation another way

What are (at least two) problems with:

```
1 double local_prod[localsize],global_inprod[localsize];  
2 for (i=0; i<localsize; i++)  
3     local_prod[i] = x[i]*y[i];  
4 MPI_Allreduce( &local_prod, &global_inprod,  
5                 localsize,MPI_DOUBLE,MPI_SUM,comm )
```

47. Inner product calculation: the right way

Compute local part, then collect local sums.

```
1 local_inprod = 0;  
2 for (i=0; i<localsize; i++)  
3     local_inprod += x[i]*y[i];  
4 MPI_Allreduce( &local_inprod, &global_inprod,  
5                 1,MPI_DOUBLE,MPI_SUM,comm )
```

48. Reduction to single process

Regular reduce: great for printing out summary information at the end of your job.

49. Reduction to root

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

- Buffers: *sendbuf*, *recvbuf* are ordinary variables/arrays.
- Every process has data in its *sendbuf*,
Root combines it in *recvbuf* (ignored on non-root processes).
- *count* is number of items in the buffer: 1 for scalar.
- *MPI_Op* is *MPI_SUM*, *MPI_MAX* et cetera.

50. Broadcast

```
1 int MPI_Bcast(  
2     void *buffer, int count, MPI_Datatype datatype,  
3     int root, MPI_Comm comm )
```

- All processes call with the same argument list
- *root* is the rank of the process doing the broadcast
- Each process allocates buffer space;
root explicitly fills in values,
all others receive values through broadcast call.
- Datatype is `MPI_FLOAT`, `MPI_INT` et cetera, different between C/Fortran.
- *comm* is usually `MPI_COMM_WORLD`

51. Gauss-Jordan elimination

<https://youtu.be/aQYuwatlWME>

Exercise 10 (jordan)

The *Gauss-Jordan algorithm* for solving a linear system with a matrix A (or computing its inverse) runs as follows:

for pivot $k = 1, \dots, n$

 let the vector of scalings $\ell_i^{(k)} = A_{ik}/A_{kk}$

 for row $r \neq k$

 for column $c = 1, \dots, n$

$$A_{rc} \leftarrow A_{rc} - \ell_r^{(k)} A_{kc}$$

where we ignore the update of the righthand side, or the formation of the inverse.

Let a matrix be distributed with each process storing one column. Implement the Gauss-Jordan algorithm as a series of broadcasts: in iteration k process k computes and broadcasts the scaling vector $\{\ell_i^{(k)}\}_i$. Replicate the right-hand side on all processors.

Exercise (optional) 11

Bonus exercise: can you extend your program to have multiple columns per processor?

Scan

52. Scan

Scan or ‘parallel prefix’: reduction with partial results

- Useful for indexing operations:
- Each process has an array of n_p elements;
- My first element has global number $\sum_{q < p} n_q$.
- Two variants: `MPI_Scan` inclusive, and `MPI_Exscan` exclusive.

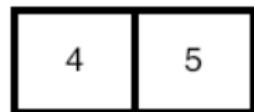
53. In vs Exclusive

process :	0	1	2	...	$p - 1$
data :	x_0	x_1	x_2	...	x_{p-1}
inclusive :	x_0	$x_0 \oplus x_1$	$x_0 \oplus x_1 \oplus x_2$...	$\bigoplus_{i=0}^{p-1} x_i$
exclusive :	unchanged	x_0	$x_0 \oplus x_1$...	$\bigoplus_{i=0}^{p-2} x_i$

MPI_Scan

Name	Param name	Explanation	C type	F type	inc
MPI_Scan (
sendbuf		starting address of send buffer	const void*	TYPE(*), DIMENSION(..)	IN
recvbuf		starting address of receive buffer	void*	TYPE(*), DIMENSION(..)	OUT
count		number of elements in input buffer	int	INTEGER	IN
datatype		datatype of elements of input buffer	MPI_Datatype	TYPE(MPI_Datatype)	IN
op		operation	MPI_Op	TYPE(MPI_Op)	IN
comm		communicator	MPI_Comm	TYPE(MPI_Comm)	IN
)					

54. For the next exercise



Exercise 12 (scangather)

- Let each process compute a random value n_{local} , and allocate an array of that length. Define

$$N = \sum n_{\text{local}}$$

- Fill the array with consecutive integers, so that all local arrays, laid end-to-end, contain the numbers $0 \cdots N - 1$. (See figure 98.)

Gather/Scatter, Barrier, and others

MPI_Gather

Name	Param name	Explanation	C type	F type	in
MPI_Gather (
sendbuf	starting address of send buffer		const void*	TYPE(*), DIMENSION(..)	IN
sendcount	number of elements in send buffer		int	INTEGER	IN
sendtype	datatype of send buffer elements		MPI_Datatype	TYPE(MPI_Datatype)	IN
recvbuf	address of receive buffer		void*	TYPE(*), DIMENSION(..)	OUT
recvcount	number of elements for any single receive		int	INTEGER	IN
recvtype	datatype of recv buffer elements		MPI_Datatype	TYPE(MPI_Datatype)	IN
root	rank of receiving process		int	INTEGER	IN
comm	communicator		MPI_Comm	TYPE(MPI_Comm)	IN
)					

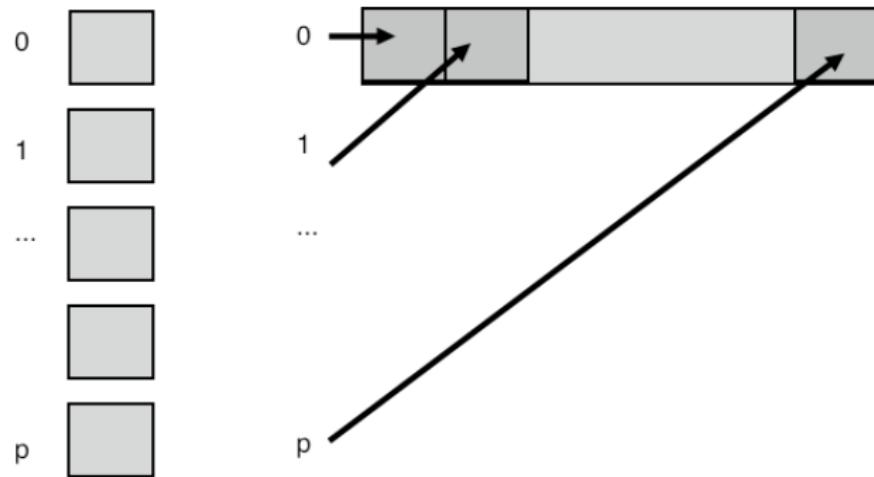
MPI_Scatter

Name	Param name	Explanation	C type	F type	in
MPI_Scatter (
	sendbuf	address of send buffer	const void*	TYPE(*), DIMENSION(..)	IN
	sendcount	number of elements sent to each process	int	INTEGER	IN
	sendtype	datatype of send buffer elements	MPI_Datatype	TYPE(MPI_Datatype)	IN
	recvbuf	address of receive buffer	void*	TYPE(*), DIMENSION(..)	OUT
	recvcount	number of elements in receive buffer	int	INTEGER	IN
	recvtype	datatype of receive buffer elements	MPI_Datatype	TYPE(MPI_Datatype)	IN
	root	rank of sending process	int	INTEGER	IN
	comm	communicator	MPI_Comm	TYPE(MPI_Comm)	IN
)				

55. Gather/Scatter

- Compare buffers to reduce
- Scatter: the `sendcount` / Gather: the `recvcount`:
this is not, as you might expect, the total length of the buffer;
instead, it is the amount of data to/from each process.

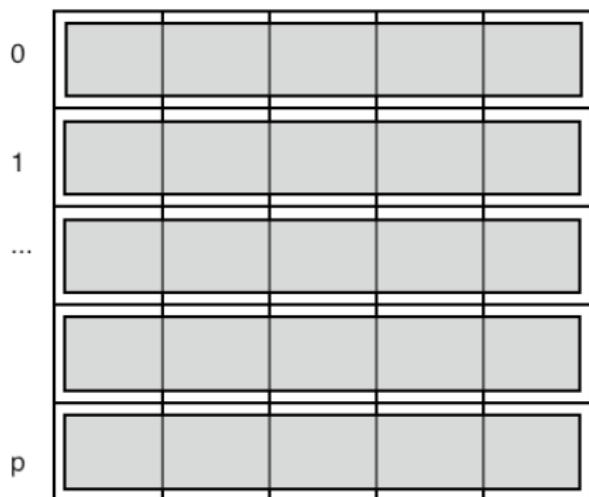
56. Gather pictured



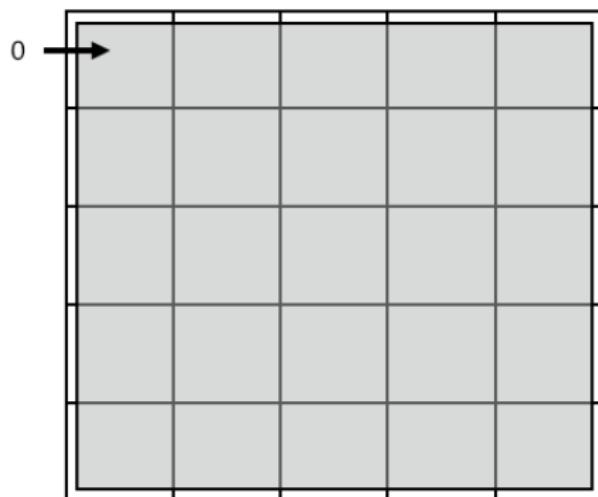
57. Popular application of gather

Matrix is constructed distributed, but needs to be brought to one process:

distributed matrix



gathered matrix

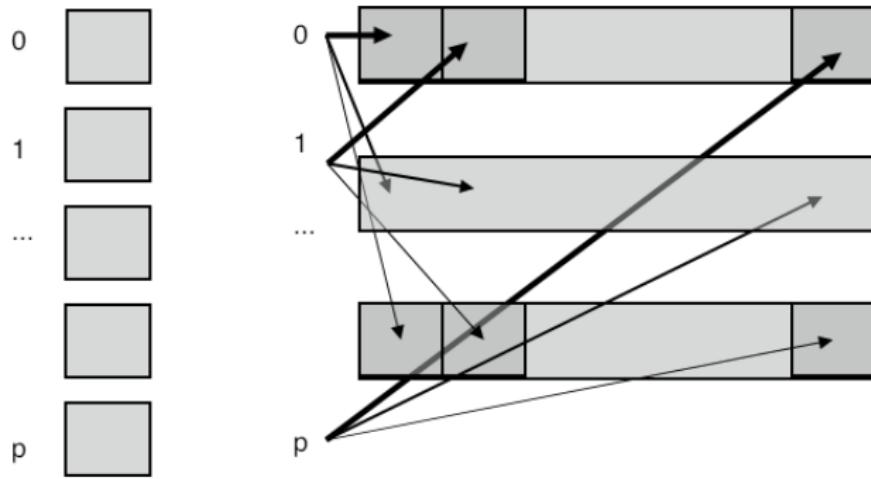


This is not efficient in time or space. Do this only when strictly necessary.
Remember SPMD: try to keep everything symmetrically parallel.

MPI_Allgather

Name	Param name	Explanation	C type	F type	inc
MPI_Allgather (
sendbuf		starting address of send buffer	const void*	TYPE(*), DIMENSION(..)	IN
sendcount		number of elements in send buffer	int	INTEGER	IN
sendtype		datatype of send buffer elements	MPI_Datatype	TYPE(MPI_Datatype)	IN
recvbuf		address of receive buffer	void*	TYPE(*), DIMENSION(..)	OUT
recvcount		number of elements received from any process	int	INTEGER	IN
recvtype		datatype of receive buffer elements	MPI_Datatype	TYPE(MPI_Datatype)	IN
comm		communicator	MPI_Comm	TYPE(MPI_Comm)	IN
)					

58. Allgather pictured



59. V-type collectives

- Gather/scatter but with individual sizes
- Requires displacement in the gather/scatter buffer

MPI_Gatherv

Name	Param name	Explanation	C type	F type	inc
MPI_Gatherv (
	sendbuf	starting address of send buffer	const void*	TYPE(*), DIMENSION(..)	IN
	sendcount	number of elements in send buffer	int	INTEGER	IN
	sendtype	datatype of send buffer elements	MPI_Datatype	TYPE(MPI_Datatype)	IN
	recvbuf	address of receive buffer	void*	TYPE(*), DIMENSION(..)	OUT
	recvcounts	non-negative integer array (of length group size) containing the number of elements that are received from each process	const int[]	INTEGER(*)	IN
	displs	integer array (of length group size). Entry i specifies the displacement relative to recvbuf at which to place the incoming data from process i	const int[]	INTEGER(*)	IN
	recvtype	datatype of recv buffer elements	MPI_Datatype	TYPE(MPI_Datatype)	IN
	root	rank of receiving process	int	INTEGER	IN
	comm	communicator	MPI_Comm	TYPE(MPI_Comm)	IN
)					

Exercise 13 (scangather)

Take the code from exercise 12 and extend it to gather all local buffers onto rank zero. Since the local arrays are of differing lengths, this requires [**MPI_Gatherv**](#).

How do you construct the lengths and displacements arrays?

Review 1

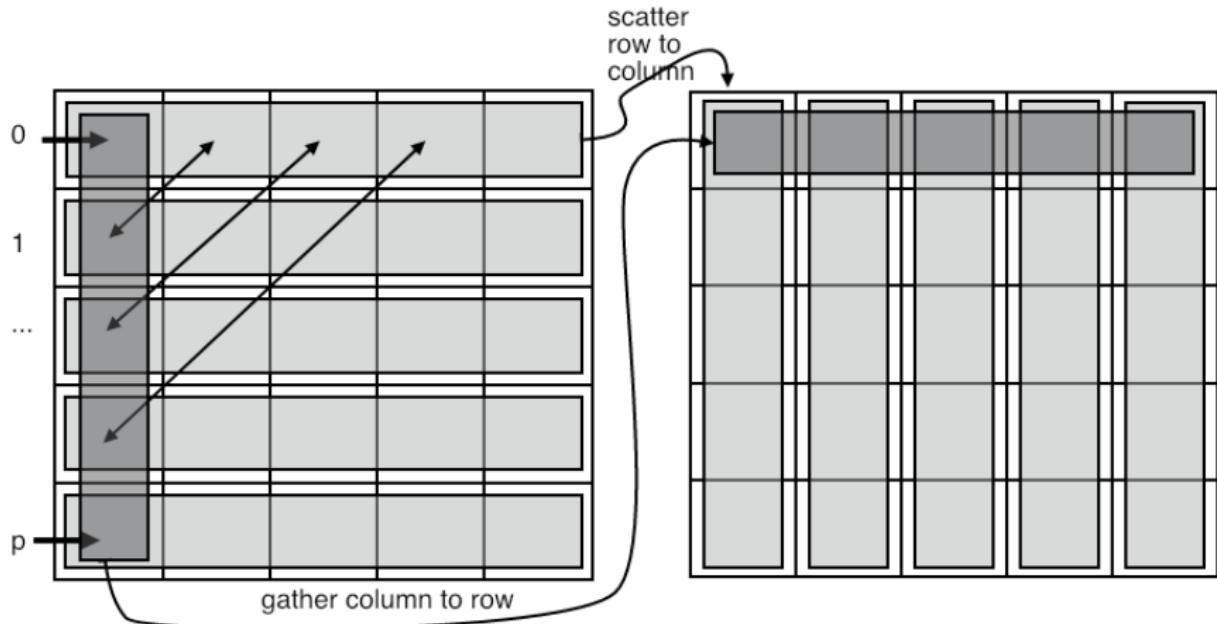
An **MPI_Scatter** call puts the same data on each process

```
/poll "A scatter call puts the same data on each process" "T" "F"
```

60. All-to-all

- Every process does a scatter;
- (equivalently: every process gather)
- each individual data, but amounts are identical
- Example: data transposition in FFT

61. Data transposition



Example: each process knows who to send to,
all-to-all gives information who to receive from

62. All-to-allv

- Every process does a scatter or gather;
- each individual data and individual amounts.
- Example: radix sort by least-significant digit.

63. Radix sort

Sort 4 numbers on two processes:

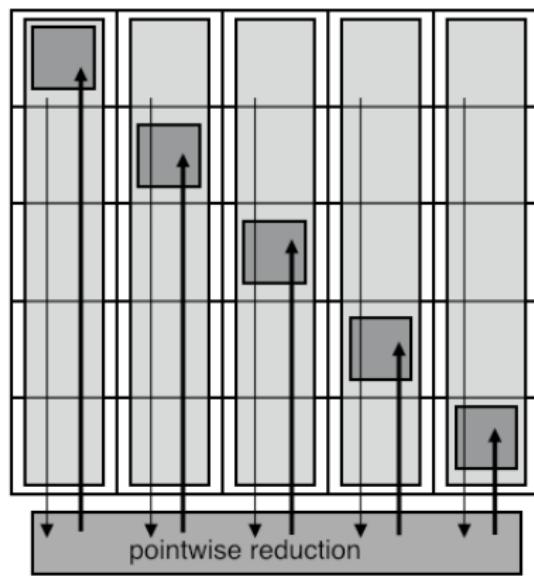
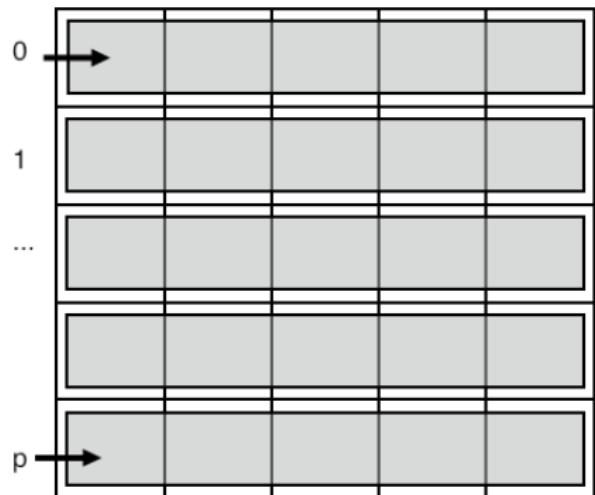
array binary	proc0		proc1	
	2	5	7	1
010	101	111	001	
stage 1				
last digit	0	1	1	1
	(this serves as bin number)			
sorted	010		101	111 001
stage 2				
next digit	1		0	1 0
	(this serves as bin number)			
sorted	101	001	010	111
stage 3				
next digit	1	0	0	1
	(this serves as bin number)			
sorted	001	010	101	111
decimal	1	2	5	7

64. Reduce-scatter

- Pointwise reduction (one element per process) followed by scatter
- Somewhat related to all-to-all: data transpose but reduced information, rather than gathered.
- Applications in both sparse and dense matrix-vector product.

65. Example: sparse matrix setup

Example: each process knows who to send to,
all-to-all gives information how many messages to expect
reduce-scatter leaves only relevant information



66. Barrier

```
1 int MPI_BARRIER( MPI_Comm comm )
```

- Synchronize processes:
- each process waits at the barrier until all processes have reached the barrier
- **This routine is almost never needed:**
collectives are already a barrier of sorts, two-sided communication is a local synchronization
- One conceivable use: timing

User-defined operators

67. MPI Operators

Define your own reduction operator

- Define operator between partial result and new operand

```
1  typedef void MPI_User_function
2      (void *invec, void *inoutvec, int *len,
3       MPI_Datatype *datatype);
```

- Don't forget to free:

```
1  int MPI_Op_free(MPI_Op *op)
```

- Make your own reduction scheme **`MPI_Reduce_local`**

67. User defined operators, Fortran

```
1 FUNCTION user_function( invec(*), inoutvec(*), length, mpitype)
2 <fortrantlype> :: invec(length), inoutvec(length)
3 INTEGER :: length, mpitype
```

MPI_Op_create

Name	Param name	Explanation	C type	F type	inc
<code>MPI_Op_create (</code>					

<code>user_fn</code>		user defined function	<code>MPI_User_function*</code>	<code>PROCEDURE (MPI_User_function)</code>	IN
<code>commute</code>		true if commutative; false otherwise.	<code>int</code>	<code>LOGICAL</code>	IN
<code>op</code>		operation	<code>MPI_Op*</code>	<code>TYPE(MPI_Op)</code>	OUT
<code>)</code>					

68. Example

Smallest nonzero:

```
1 *(int*)inout = m;  
2 }
```

Review 2

The $\|\cdot\|_2$ norm (sum of squares) needs a custom operator.

```
/poll "The sum of squares norm needs a custom operators" "T" "F"
```

Performance of collectives

69. Naive realization of collectives

Broadcast:



Single message:

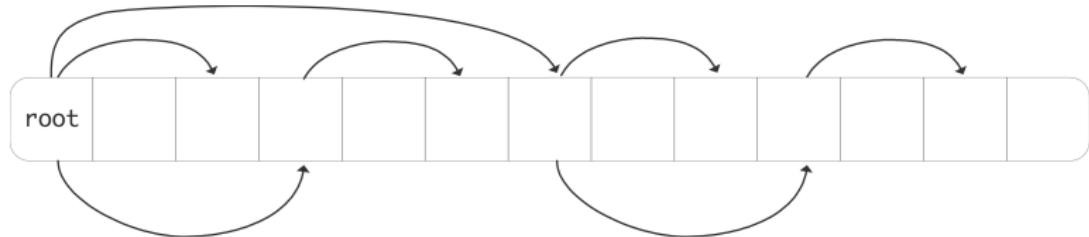
$$\alpha = \text{message startup} \approx 10^{-6} s, \quad \beta = \text{time per word} \approx 10^{-9} s$$

- Time for message of n words:

$$\alpha + \beta n$$

- Time for collective? Can you improve on that?

70. Better implementation of collective



- What is the running time now?
- Can you come up with lower bounds on the α, β terms? Are these achieved here?
- How about the case of really long buffers?

71. Implementation of Reduce

	$t = 1$	$t = 2$	$t = 3$
p_0	$x_0^{(0)}, x_1^{(0)}, x_2^{(0)}, x_3^{(0)}$	$x_0^{(0:1)}, x_1^{(0:1)}, x_2^{(0:1)}, x_3^{(0:1)}$	$x_0^{(0:3)}, x_1^{(0:3)}, x_2^{(0:3)}, x_3^{(0:3)}$
p_1	$x_0^{(1)} \uparrow, x_1^{(1)} \uparrow, x_2^{(1)} \uparrow, x_3^{(1)} \uparrow$		
p_2	$x_0^{(2)}, x_1^{(2)}, x_2^{(2)}, x_3^{(2)}$	$x_0^{(2:3)} \uparrow, x_1^{(2:3)} \uparrow, x_2^{(2:3)} \uparrow, x_3^{(2:3)} \uparrow$	
p_3	$x_0^{(3)} \uparrow, x_1^{(3)} \uparrow, x_2^{(3)} \uparrow, x_3^{(3)} \uparrow$		

72. Implementation of Allreduce

	$t = 1$	$t = 2$	$t = 3$
p_0	$x_0^{(0)} \downarrow, x_1^{(0)} \downarrow, x_2^{(0)} \downarrow, x_3^{(0)} \downarrow$	$x_0^{(0:1)} \downarrow\downarrow, x_1^{(0:1)} \downarrow\downarrow, x_2^{(0:1)} \downarrow\downarrow, x_3^{(0:1)} \downarrow\downarrow$	$x_0^{(0:3)}, x_1^{(0:3)}, x_2^{(0:3)}, x_3^{(0:3)}$
p_1	$x_0^{(1)} \uparrow, x_1^{(1)} \uparrow, x_2^{(1)} \uparrow, x_3^{(1)} \uparrow$	$x_0^{(0:1)} \downarrow\downarrow, x_1^{(0:1)} \downarrow\downarrow, x_2^{(0:1)} \downarrow\downarrow, x_3^{(0:1)} \downarrow\downarrow$	$x_0^{(0:3)}, x_1^{(0:3)}, x_2^{(0:3)}, x_3^{(0:3)}$
p_2	$x_0^{(2)} \downarrow, x_1^{(2)} \downarrow, x_2^{(2)} \downarrow, x_3^{(2)} \downarrow$	$x_0^{(2:3)} \uparrow\uparrow, x_1^{(2:3)} \uparrow\uparrow, x_2^{(2:3)} \uparrow\uparrow, x_3^{(2:3)} \uparrow\uparrow$	$x_0^{(0:3)}, x_1^{(0:3)}, x_2^{(0:3)}, x_3^{(0:3)}$
p_3	$x_0^{(3)} \uparrow, x_1^{(3)} \uparrow, x_2^{(3)} \uparrow, x_3^{(3)} \uparrow$	$x_0^{(2:3)} \uparrow\uparrow, x_1^{(2:3)} \uparrow\uparrow, x_2^{(2:3)} \uparrow\uparrow, x_3^{(2:3)} \uparrow\uparrow$	$x_0^{(0:3)}, x_1^{(0:3)}, x_2^{(0:3)}, x_3^{(0:3)}$

Review 3

True or false: there are collectives that do not communicate data

```
/poll "there are collectives that do not communicate data" "T" "F"
```

Reduction operators

73. User-defined operators

Given a reduction function:

```
1  typedef void user_function
2      (void *invec, void *inoutvec, int *len,
3       MPI_Datatype *datatype);
```

create a new operator:

```
1  MPI_Op rwz;
2  MPI_Op_create(user_function,1,&rwz);
3  MPI_Allreduce(data+procno,&positive_minimum,1,MPI_INT,rwz,comm);
```

Exercise 14 (onenorm)

Write the reduction function to implement the *one-norm* of a vector:

$$\|x\|_1 \equiv \sum_i |x_i|.$$

Part III

Point-to-point communication

74. Overview

This section concerns direct communication between two processes.
Discussion of distributed work, deadlock and other parallel phenomena.

Commands learned:

- `MPI_Send`, `MPI_Recv`, `MPI_Sendrecv`, `MPI_Isend`, `MPI_Irecv`
- `MPI_Wait...`
- Mention of `MPI_Test`, `MPI_Bsend/Ssend/Rsend`.

Point-to-point communication

75. MPI point-to-point mechanism

- Two-sided communication
- Matched send and receive calls
- One process sends to a specific other process
- Other process does a specific receive.

76. Ping-pong

A sends to B, B sends back to A

What is the code for A? For B?

77. Ping-pong in MPI

Remember SPMD:

```
1 if ( /* I am process A */ ) {  
2     MPI_Send( /* to: */ B ..... );  
3     MPI_Recv( /* from: */ B ... );  
4 } else if ( /* I am process B */ ) {  
5     MPI_Recv( /* from: */ A ... );  
6     MPI_Send( /* to: */ A ..... );  
7 }
```

MPI_Send

Name	Param name	Explanation	C type	F type	in/out
MPI_Send (
buf		initial address of send buffer	const void*	TYPE(*), DIMENSION(..)	IN
count		number of elements in send buffer	int	INTEGER	IN
datatype		datatype of each send buffer element	MPI_Datatype	TYPE(MPI_Datatype)	IN
dest		rank of destination	int	INTEGER	IN
tag		message tag	int	INTEGER	IN
comm		communicator	MPI_Comm	TYPE(MPI_Comm)	IN
)					

MPI_Recv

Name	Param name	Explanation	C type	F type	inc
MPI_Recv (
	buf	initial address of receive buffer	void*	TYPE(*), DIMENSION(..)	OUT
	count	number of elements in receive buffer	int	INTEGER	IN
	datatype	datatype of each receive buffer element	MPI_Datatype	TYPE(MPI_Datatype)	IN
	source	rank of source or MPI_ANY_SOURCE	int	INTEGER	IN
	tag	message tag or MPI_ANY_TAG	int	INTEGER	IN
	comm	communicator	MPI_Comm	TYPE(MPI_Comm)	IN
	status	status object	MPI_Status*	TYPE(MPI_Status)	OUT
)					

78. Status object

Use `MPI_STATUS_IGNORE` unless . . .

- Receive call can have various wildcards:
`MPI_ANY_SOURCE, MPI_ANY_TAG`
- Receive buffer size is actually upper bound, not exact
- Use status object to retrieve actual description of the message

```
1 int s = status.MPI_SOURCE;
2 int t = status.MPI_TAG;
3 MPI_Get_count(status,MPI_FLOAT,&c);
```

Exercise 15 (pingpong)

Implement the ping-pong program. Add a timer using `MPI_Wtime`. For the status argument of the receive call, use `MPI_STATUS_IGNORE`.

- Run multiple ping-pongs (say a thousand) and put the timer around the loop. The first run may take longer; try to discard it.
- Run your code with the two communicating processes first on the same node, then on different nodes. Do you see a difference?
- Then modify the program to use longer messages. How does the timing increase with message size?

For bonus points, can you do a regression to determine α, β ?

MPI_Wtime

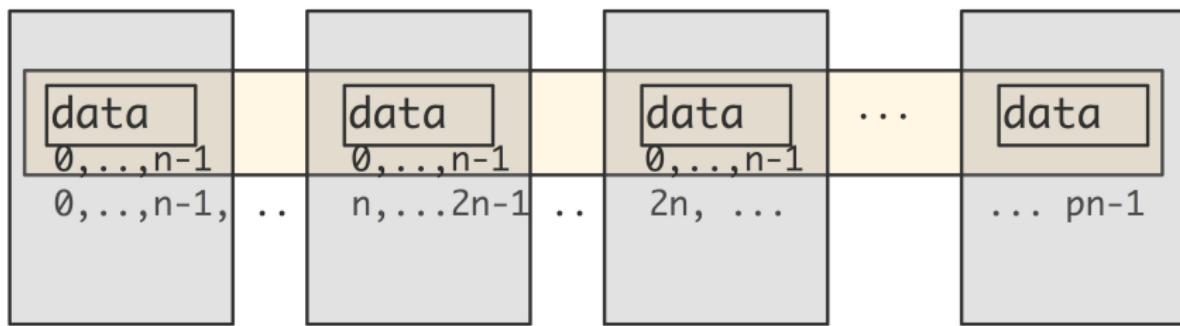
Name	Param name	Explanation	C type	F type	inc
MPI_Wtime ()			

Distributed data

79. Distributed data

Distributed array: each process stores disjoint local part

```
int n;  
double data[n];
```



Local numbering $0, \dots, n_{\text{local}}$;
global numbering is 'in your mind'.

80. Local and global indexing

Every local array starts at 0 (Fortran: 1);
you have to translate that yourself to global numbering:

```
1 int myfirst = .....;
2 for (int ilocal=0; ilocal<nlocal; ilocal++) {
3     int iglobal = myfirst+ilocal;
4     array[ilocal] = f(iglobal);
5 }
```

Exercise (optional) 16

Implement a (very simple-minded) Fourier transform: if f is a function on the interval $[0, 1]$, then the n -th Fourier coefficient is

$$f_n \hat{=} \int_0^1 f(t) e^{-2\pi x} dx$$

which we approximate by

$$f_n \hat{=} \sum_{i=0}^{N-1} f(ih) e^{-in\pi/N}$$

- Make one distributed array for the e^{-inh} coefficients,
- make one distributed array for the $f(ih)$ values
- calculate a couple of coefficients

81. Load balancing

If the distributed array is not perfectly divisible:

```
1 int Nglobal, // is something large
2     Nlocal = Nglobal/nprocs,
3     excess = Nglobal%nprocs;
4 if (procno==nprocs-1)
5     Nlocal += excess;
```

This gives a load balancing problem. Better solution?

82. (for future reference)

Let

$$f(i) = \lfloor iN/p \rfloor$$

and give process i the points $f(i)$ up to $f(i + 1)$.

Result:

$$\lfloor N/p \rfloor \leq f(i + 1) - f(i) \leq \lceil N/p \rceil$$

Local information exchange

83. Motivation

Partial differential equations:

$$-\Delta u = -u_{xx}(\bar{x}) - u_{yy}(\bar{x}) = f(\bar{x}) \text{ for } \bar{x} \in \Omega = [0, 1]^2 \text{ with } u(\bar{x}) = u_0 \text{ on } \delta\Omega.$$

Simple case:

$$-u_{xx} = f(x).$$

Finite difference approximation:

$$\frac{2u(x) - u(x+h) - u(x-h)}{h^2} = f(x, u(x), u'(x)) + O(h^2),$$

Finite dimensional: $u_i \equiv u(ih)$.

84. Motivation (continued)

Equations

$$\begin{cases} -u_{i-1} + 2u_i - u_{i+1} &= h^2 f(x_i) \quad 1 < i < n \\ 2u_1 - u_2 &= h^2 f(x_1) + u_0 \\ 2u_n - u_{n-1} &= h^2 f(x_n) + u_{n+1}. \end{cases}$$

$$\begin{pmatrix} 2 & -1 & & \emptyset \\ -1 & 2 & -1 & \\ \emptyset & \ddots & \ddots & \ddots \end{pmatrix} \begin{pmatrix} u_1 \\ u_2 \\ \vdots \end{pmatrix} = \begin{pmatrix} h^2 f_1 + u_0 \\ h^2 f_2 \\ \vdots \end{pmatrix} \quad (1)$$

So we are interested in sparse/banded matrices.

85. Matrix vector product

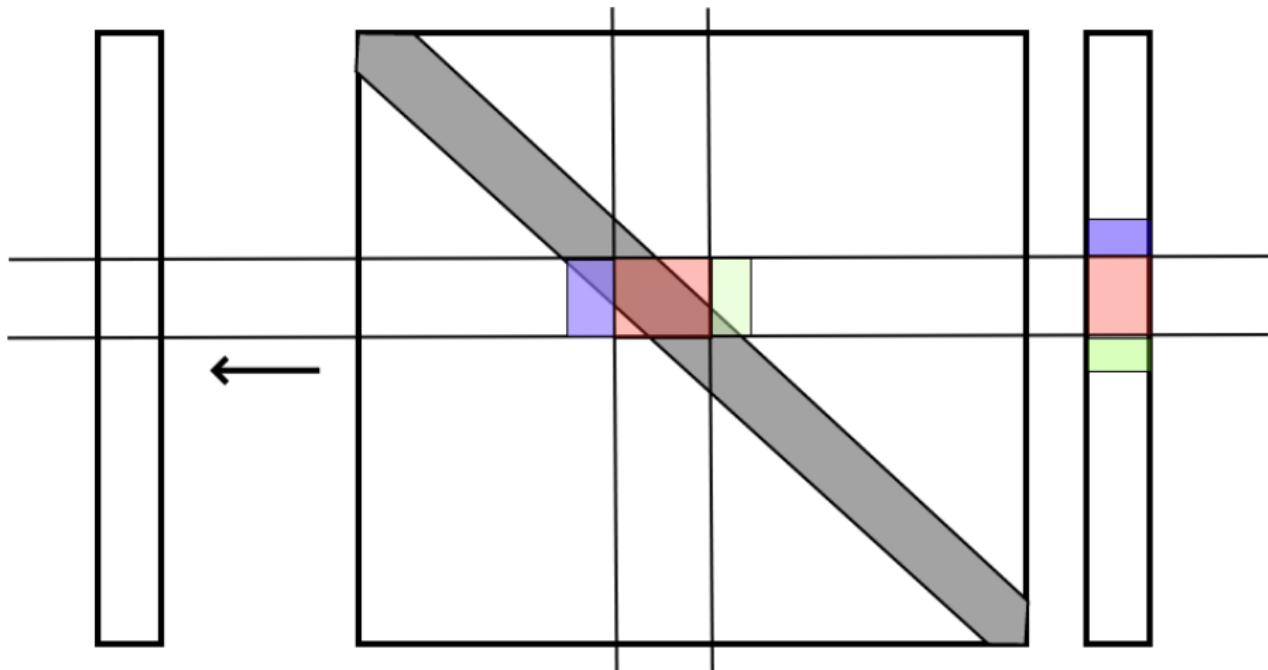
Most common operation: matrix vector product

$$y \leftarrow Ax, \quad A = \begin{pmatrix} 2 & -1 & & \\ -1 & 2 & -1 & \\ & \ddots & \ddots & \ddots \end{pmatrix} \begin{pmatrix} u_1 \\ u_2 \\ \vdots \end{pmatrix}$$

- Component operation: $y_i = 2x_i - x_{i-1} - x_{i+1}$
- Parallel execution: each process has range of i -coordinates
- \Rightarrow segment of vector, block row of matrix

86. Partitioned matrix-vector product

We need a point-to-point mechanism:



each process with ones before/after it.

87. Operating on distributed data

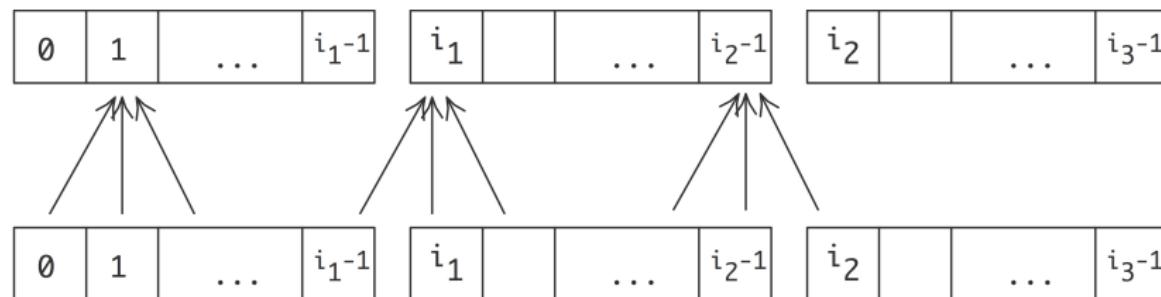
Array of numbers $x_i: i = 0, \dots, N$

compute

$$y_i = -x_{i-1} + 2x_i - x_{i+1}: i = 1, \dots, N-1$$

'owner computes'

This leads to communication:



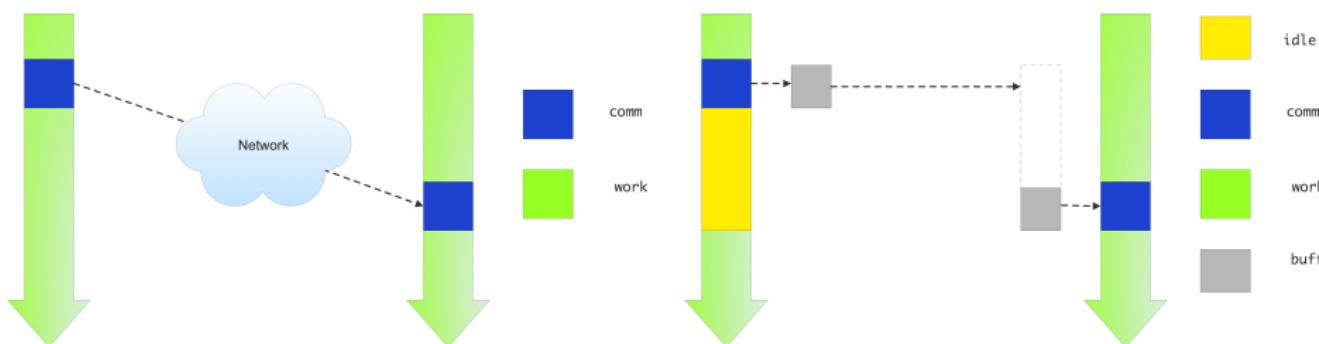
so we need a point-to-point mechanism.

Blocking communication

88. Blocking send/recv

`MPI_Send` and `MPI_Recv` are *blocking* operations:

- The process waits ('blocks') until the operation is concluded.
- A send can not complete until the receive executes.



Ideal vs actual send/recv behaviour.

89. Deadlock

Exchange between two processes:

```
1 other = 1-procno; /* if I am 0, other is 1; and vice versa */
2 receive(source=other);
3 send(target=other);
```

A subtlety.

This code may actually work:

```
1 other = 1-procno; /* if I am 0, other is 1; and vice versa */
2 send(target=other);
3 receive(source=other);
```

Small messages get sent even if there is no corresponding receive.
(Often a system parameter)

90. Protocol

Communication is a ‘rendez-vous’ or ‘hand-shake’ protocol:

- Sender: ‘I have data for you’
- Receiver: ‘I have a buffer ready, send it over’
- Sender: ‘Ok, here it comes’
- Receiver: ‘Got it.’

Small messages bypass this: ‘eager’ send.

Definition of ‘small message’ controlled by environment variables:

`I_MPI_EAGER_THRESHOLD MV2_IBA_EAGER_THRESHOLD`

Exercise 17

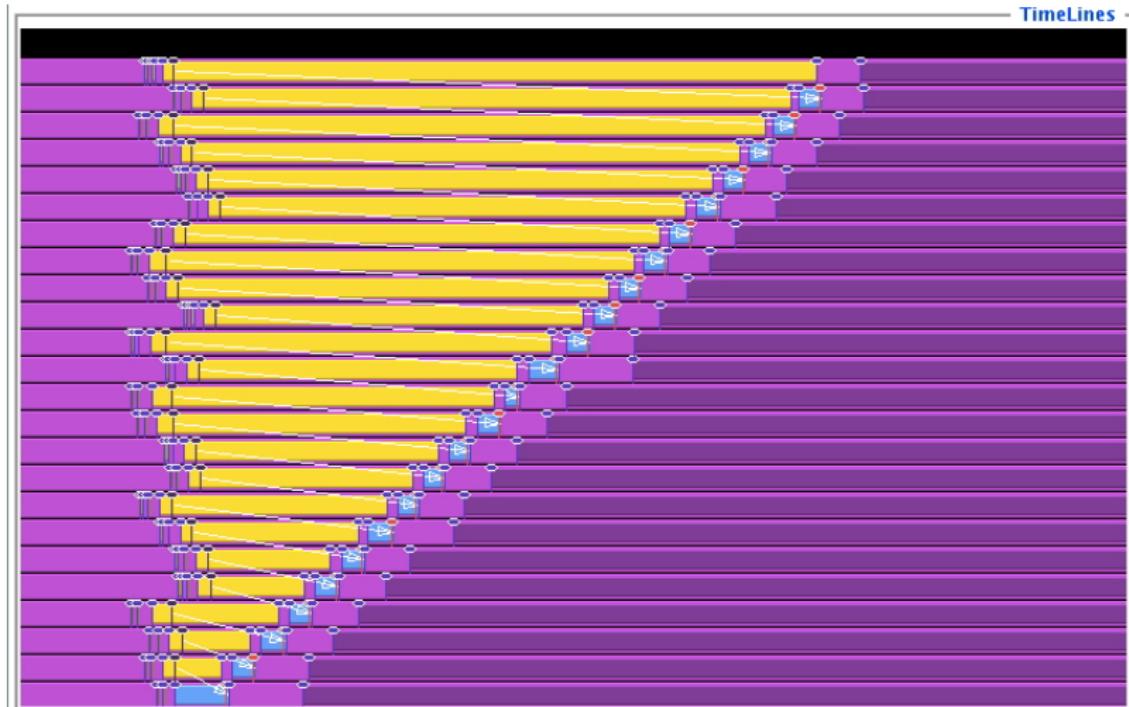
(Classroom exercise) Each student holds a piece of paper in the right hand – keep your left hand behind your back – and we want to execute:

- ① Give the paper to your right neighbor;
- ② Accept the paper from your left neighbor.

Including boundary conditions for first and last process, that becomes the following program:

- ① If you are not the rightmost student, turn to the right and give the paper to your right neighbor.
- ② If you are not the leftmost student, turn to your left and accept the paper from your left neighbor.

91. TAU trace: serialization



92. The problem here...

Here you have a case of a program that computes the right output, just way too slow.

Beware! Blocking sends/receives can be trouble.
(How would you solve this particular case?)

Food for thought: what happens if you flip the send and receive call?

Exercise (optional) 18

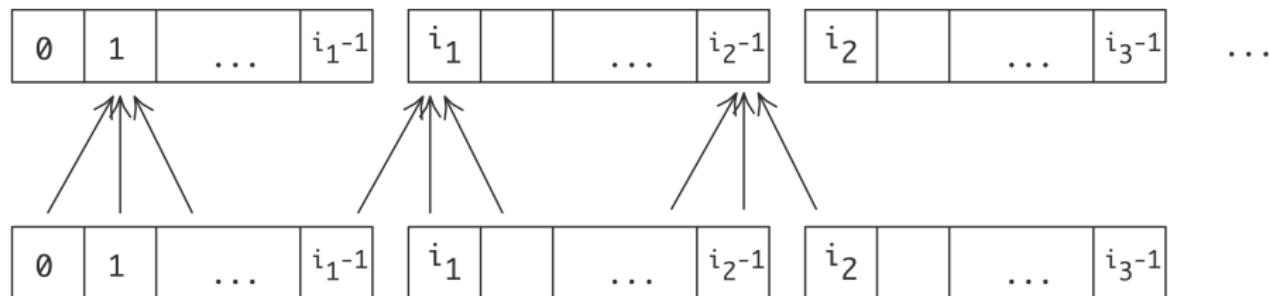
Implement the above algorithm using `MPI_Send` and `MPI_Recv` calls. Run the code, and use TAU to reproduce the trace output of figure 162. If you don't have TAU, can you show this serialization behavior using timings, for instance running it on an increasing number of processes?

Pairwise exchange

93. Operating on distributed data

Take another look:

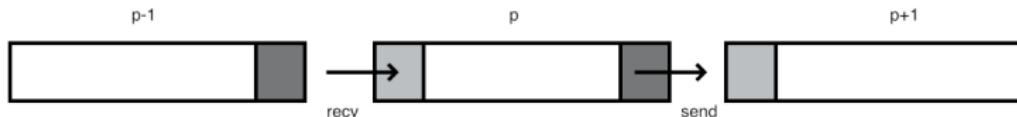
$$y_i = x_{i-1} + x_i + x_{i+1}: i = 1, \dots, N - 1$$



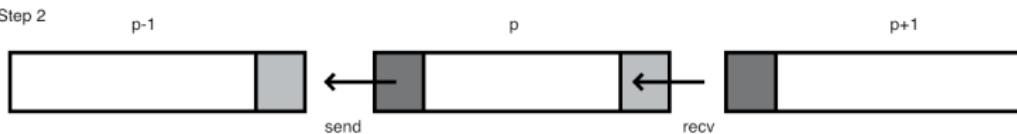
- One-dimensional data and linear process numbering;
- Operation between neighboring indices: communication between neighboring processes.

94. Two steps

Step 1



Step 2



First do all the data movement to the right, later to the left.

- Each process does a send and receive
- So everyone does the send, then the receive? We just saw the problem with that.
- Better solution coming up!

95. Sendrecv

Instead of separate send and receive: use

`MPI_Sendrecv`

Combined calling sequence of send and receive;
execute such that no deadlock or sequentialization.

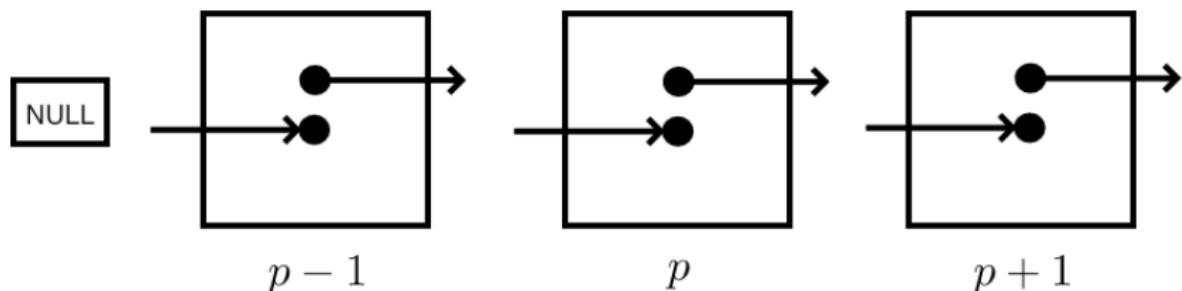
(Also: `MPI_Sendrecv_replace` with single buffer.)

MPI_Sendrecv

Name	Param name	Explanation	C type	F type	inc
MPI_Sendrecv (
sendbuf		initial address of send buffer	const void*	TYPE(*), DIMENSION(..)	IN
sendcount		number of elements in send buffer	int	INTEGER	IN
sendtype		type of elements in send buffer	MPI_Datatype	TYPE(MPI_Datatype)	IN
dest		rank of destination	int	INTEGER	IN
sendtag		send tag	int	INTEGER	IN
recvbuf		initial address of receive buffer	void*	TYPE(*), DIMENSION(..)	OUT
recvcount		number of elements in receive buffer	int	INTEGER	IN
recvtype		type of elements receive buffer element	MPI_Datatype	TYPE(MPI_Datatype)	IN
source		rank of source or MPI_ANY_SOURCE	int	INTEGER	IN
recvtag		receive tag or MPI_ANY_TAG	int	INTEGER	IN
comm		communicator	MPI_Comm	TYPE(MPI_Comm)	IN
status		status object	MPI_Status*	TYPE(MPI_Status)	OUT
)					

96. SPMD picture

What does process p do?



97. Sendrecv with incomplete pairs

```
1 MPI_Comm_rank( .... &procno );
2 if ( /* I am not the first process */ )
3     predecessor = procno-1;
4 else
5     predecessor = MPI_PROC_NULL;
6
7 if ( /* I am not the last process */ )
8     successor = procno+1;
9 else
10    successor = MPI_PROC_NULL;
11
12 sendrecv(from=predecessor,to=successor);
```

(Receive from `MPI_PROC_NULL` succeeds without altering the receive buffer.)

98. A point of programming style

The previous slide had:

- a conditional for computing the sender and receiver rank;
- a single Sendrecv call.

Also possible:

```
1 if ( /* i am first */ )
2   Sendrecv( to=right, from=NULL
3             ↣);
4 else if ( /* i am last */ )
5   Sendrecv( to=NULL,   from=left
6             ↣);
7 else
8   Sendrecv( to=right, from=left
9             ↣);
```

```
1 if ( /* i am first */ )
2   Send( to=right );
3 else if ( /* i am last */ )
4   Recv( from=left );
5 else
6   Sendrecv( to=right, from=left
7             ↣);
```

But:

Code duplication is error-prone, also
chance of deadlock by missing a case

Exercise (optional) 19 (rightsend)

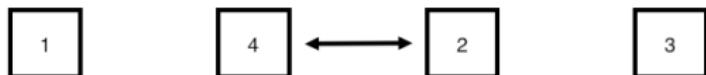
Revisit exercise 17 and solve it using [MPI_Sendrecv](#).

If you have TAU installed, make a trace. Does it look different from the serialized send/recv code? If you don't have TAU, run your code with different numbers of processes and show that the runtime is essentially constant.

Exercise 20 (sendrecv)

Implement the above three-point combination scheme using `MPI_Sendrecv`; every processor only has a single number to send to its neighbor.

99. Odd-even transposition sort



↔ transpose performed
↔ no transpose needed

Odd-even transposition sort on 4 elements.

Exercise (optional) 21

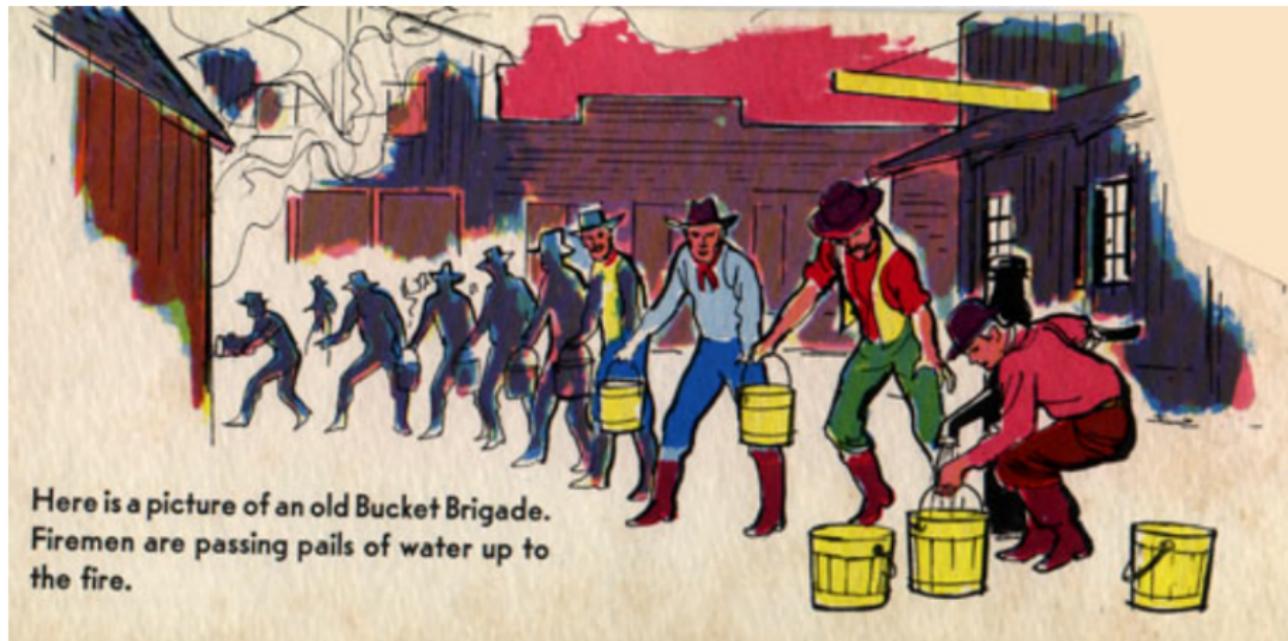
A very simple sorting algorithm is *swap sort* or *odd-even transposition sort*: pairs of processors compare data, and if necessary exchange. The elementary step is called a *compare-and-swap*: in a pair of processors each sends their data to the other; one keeps the minimum values, and the other the maximum. For simplicity, in this exercise we give each processor just a single number.

The transposition sort algorithm is split in even and odd stages, where in the even stage processors $2i$ and $2i + 1$ compare and swap data, and in the odd stage processors $2i + 1$ and $2i + 2$ compare and swap. You need to repeat this $P/2$ times, where P is the number of processors; see figure 175.

Implement this algorithm using `MPI_Sendrecv`. (Use `MPI_PROC_NULL` for the edge cases if needed.) Use a gather call to print the global state of the distributed array at the beginning and end of the sorting process.

100. Bucket brigade

Sometimes you really want to pass information from one process to the next: 'bucket brigade'



Here is a picture of an old Bucket Brigade.
Firemen are passing pails of water up to
the fire.

Exercise 22 (bucketblock)

Take the code of exercise 18 and modify it so that the data from process zero gets propagated to every process. Specifically, compute all partial sums $\sum_{i=0}^p i^2$:

$$\begin{cases} x_0 = 1 & \text{on process zero} \\ x_p = x_{p-1} + (p+1)^2 & \text{on process } p \end{cases}$$

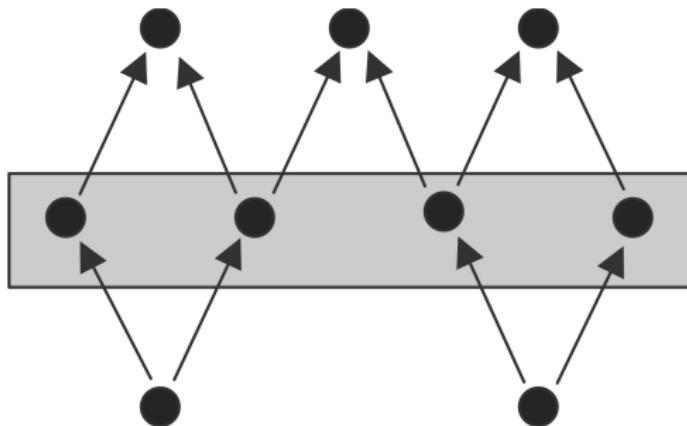
Use `MPI_Send` and `MPI_Recv`; make sure to get the order right.

Food for thought: all quantities involved here are integers. Is it a good idea to use the integer datatype here?

Irregular exchanges: non-blocking communication

101. Sending with irregular connections

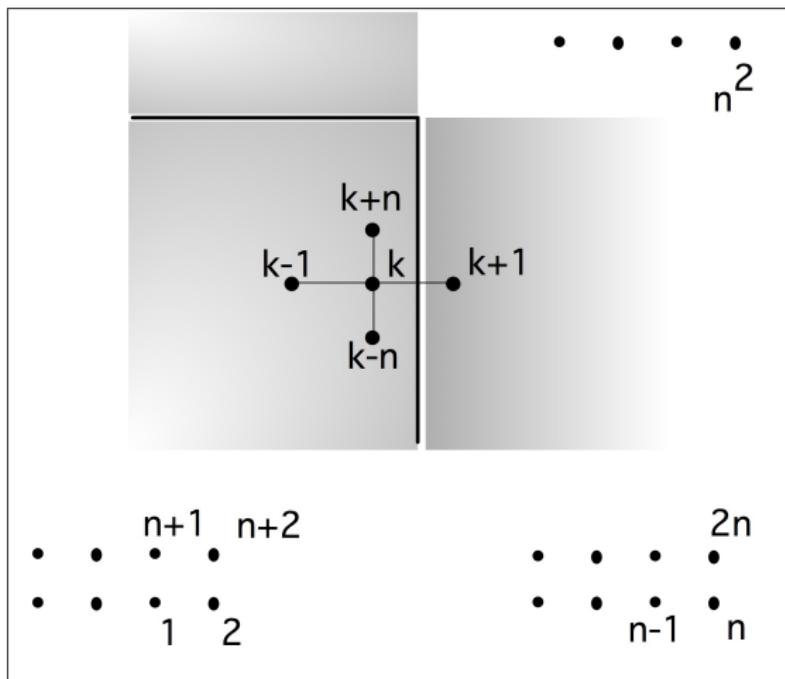
Graph operations:



Communicating other than in pairs

102. PDE, 2D case

A difference stencil applied to a two-dimensional square domain, distributed over processors. A cross-processor connection is indicated \Rightarrow complicated to express pairwise

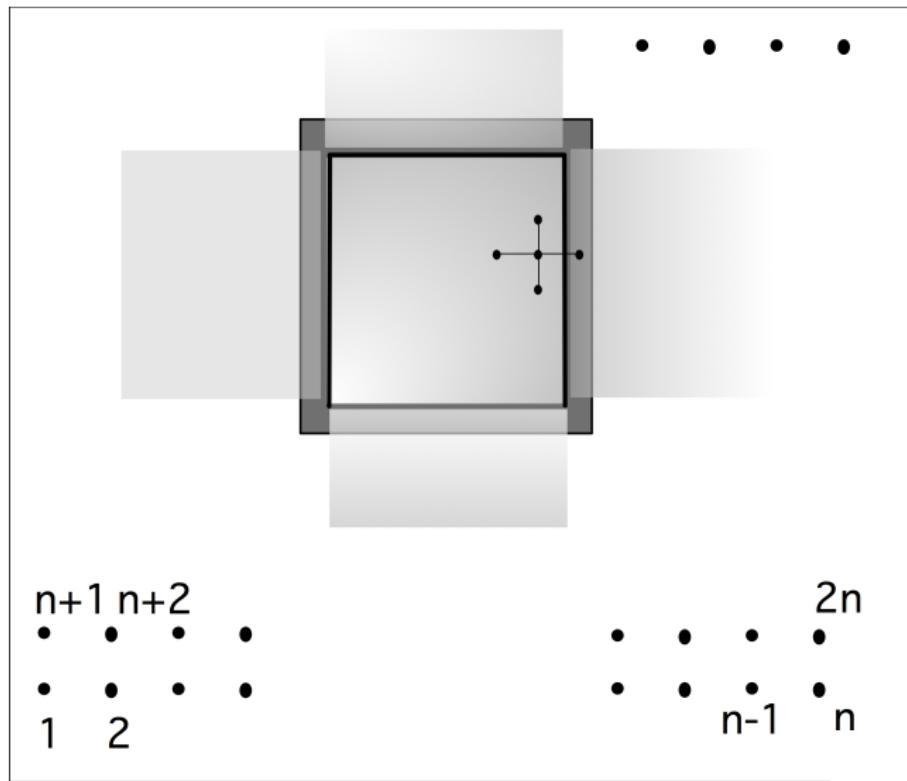


103. PDE matrix

$$A = \left(\begin{array}{ccc|c|cc|c} 4 & -1 & & \emptyset & -1 & -1 & \emptyset \\ -1 & 4 & -1 & & & & \\ \ddots & \ddots & \ddots & & & & \\ & \ddots & \ddots & -1 & & & \\ \emptyset & & -1 & 4 & \emptyset & & -1 \\ \hline -1 & & & \emptyset & 4 & -1 & -1 \\ & -1 & & & -1 & 4 & -1 \\ & & \uparrow & & k-1 & k & k+1 & \uparrow \\ k-n & & & & & & & k+n \\ \hline & & -1 & & & -1 & 4 & \\ \end{array} \right)$$

104. Halo region

The halo region of a process, induced by a stencil



105. How do you approach this?

- It is very hard to figure out a send/receive sequence that does not deadlock or serialize
- Even if you manage that, you may have process idle time.

Instead:

- Declare 'this data needs to be sent' or 'these messages are expected', and
- then wait for them collectively.

106. Non-blocking send/recv

- `MPI_Isend` / `MPI_Irecv` does not send/receive:
- They declare a buffer.
- The buffer contents are there after a wait call.
- In between the `MPI_Isend` and `MPI_Wait` the data may not have been sent.
- In between the `MPI_Irecv` and `MPI_Wait` the data may not have arrived.

```
1 // start non-blocking communication
2 MPI_Isend( ... ); MPI_Irecv( ... );
3 // wait for the Isend/Irecv calls to finish in any order
4 MPI_Wait( ... );
```

107. Syntax

Very much like blocking MPI_Send/MPI_Recv:

```
1 int MPI_Isend(void *buf,
2     int count, MPI_Datatype datatype, int dest, int tag,
3     MPI_Comm comm, MPI_Request *request)
4 int MPI_Irecv(void *buf,
5     int count, MPI_Datatype datatype, int source, int tag,
6     MPI_Comm comm, MPI_Request *request)
```

Basic wait:

```
1 MPI_Wait( MPI_Request*, MPI_Status* );
```

Most common way of waiting for completion:

```
1 int MPI_Waitall(int count, MPI_Request array_of_requests[],
2     MPI_Status array_of_statuses[])
```

- ignore status: MPI_STATUSES_IGNORE
- also MPI_Wait, MPI_Waitany, MPI_Waitsome

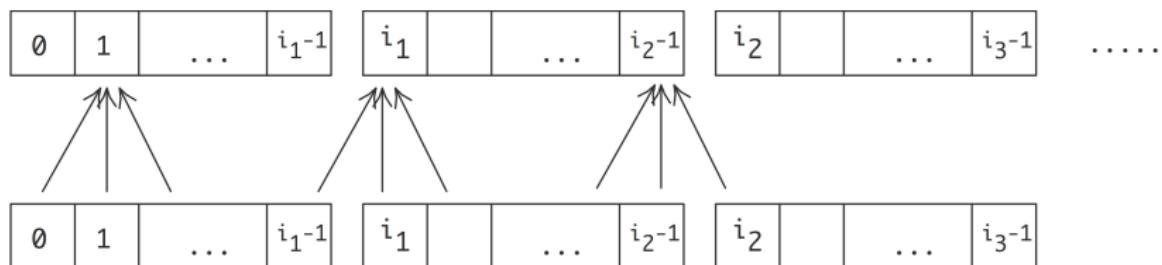
Exercise 23 (isendrecv)

Now use nonblocking send/receive routines to implement the three-point averaging operation

$$y_i = (x_{i-1} + x_i + x_{i+1})/3: i = 1, \dots, N - 1$$

on a distributed array. There are two approaches to the first and last process:

- ① you can use `MPI_PROC_NULL` for the 'missing' communications;
- ② you can skip these communications altogether, but now you have to count the requests carefully.



(Can you think of a different way of handling th

108. Comparison

- Obvious: blocking vs non-blocking behaviour.
- Buffer reuse: when a blocking call returns, the buffer is safe for reuse or free;
- A buffer in a non-blocking call can only be reused/freed after the wait call.

109. Buffer use in blocking/non-blocking case

Blocking:

```
1 double *buffer;
2 // allocate the buffer
3 for ( ... p ... ) {
4     buffer = // fill in the data
5     MPI_Send( buffer, ... /* to: */ p );
```

Non-blocking:

```
1 double **buffers;
2 // allocate the buffers
3 for ( ... p ... ) {
4     buffers[p] = // fill in the data
5     MPI_Isend( buffers[p], ... /* to: */ p );
6     MPI_Waitsomething(.....)
```

110. Pitfalls

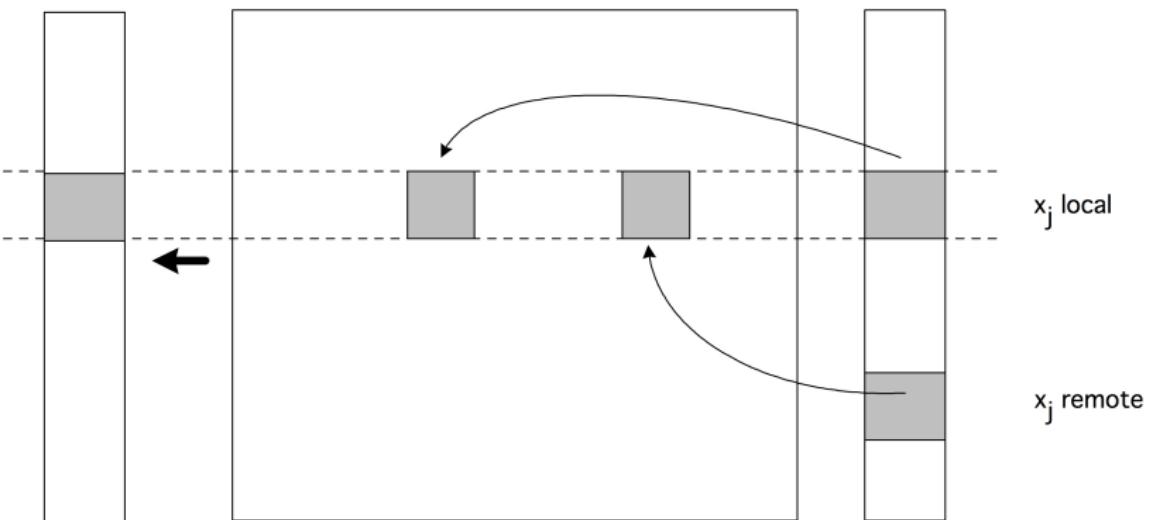
- Strictly one request/wait per `isend`/`irecv`:
can not use one request for multiple simultaneous `isends`
- Some people argue:
Wait for the send is not necessary: if you wait for the receive, the message has arrived safely

This leads to memory leaks! The `wait` call deallocates the request object.

111. Matrices in parallel

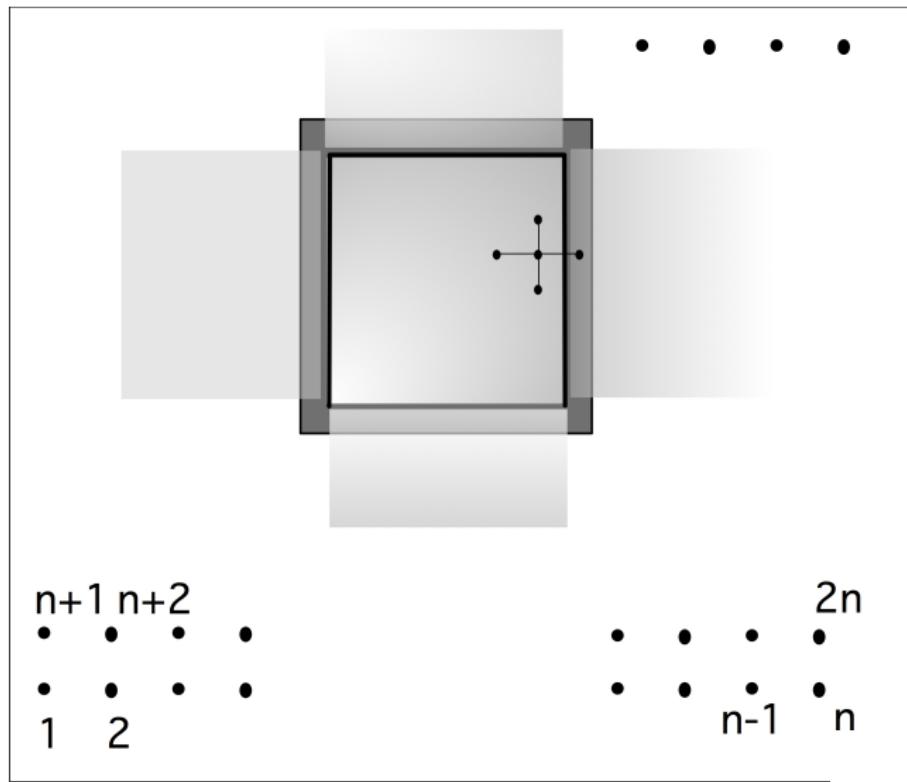
$$y \leftarrow Ax$$

and A, x, y all distributed:



112. Hiding the halo

Interior of a process domain can overlap with halo transfer:



113. Latency hiding

Other motivation for non-blocking calls:

overlap of computation and communication, provided hardware support.

Also known as 'latency hiding'.

Example: three-point combination operation (see above):

- ① Start communication for edge points,
- ② Do local operations while communication goes on,
- ③ Wait for edge points from neighbor processes
- ④ Incorporate incoming data.

Exercise 24 (isendirecvarray)

Take your code of exercise 23 and modify it to use latency hiding.
Operations that can be performed without needing data from neighbors
should be performed in between the `MPI_Isend` / `MPI_Irecv` calls and the
corresponding `MPI_Wait` calls.

Write your code so that it can achieve latency hiding.

114. Test: non-blocking wait

- Post non-blocking receives
- test for incoming messages
- if nothing comes in, do local work

```
1 while (1) {  
2     MPI_Test( /* from: */ MPI_ANY_SOURCE, &flag );  
3     if (flag)  
4         // do something with incoming message  
5     else  
6         // do local work  
7 }
```

115. The Pipeline Pattern

- Remember the bucket brigade: data propagating through processes
- If you have many buckets being passed: pipeline
- This is very parallel: only filling and draining the pipeline is not completely parallel
- Application to long-vector broadcast: pipelining gives overlap

Exercise (optional) 25 (bucketpipenonblock)

Implement a pipelined broadcast for long vectors:
use non-blocking communication to send the vector in parts.

Exercise 26 (setdiff)

Create two distributed arrays of positive integers. Take the set difference of the two: the first array needs to be transformed to remove from it those numbers that are in the second array.

How could you solve this with an `MPI_Allgather` call? Why is it not a good idea to do so? Solve this exercise instead with a circular bucket brigade algorithm.

Consider: `MPI_Send` and `MPI_Recv` vs `MPI_Sendrecv` vs `MPI_Sendrecv_replace` vs `MPI_Isend` and `MPI_Irecv`

116. The wheel of reinvention

The circular bucket brigade is the idea behind the 'Horovod' library, which is the key to efficient parallel Deep Learning.

117. More sends and receive

- `MPI_Bsend`, `MPI_Ibsend`: buffered send
- `MPI_Ssend`, `MPI_Issend`: synchronous send
- `MPI_Rsend`, `MPI_Irsend`: ready send
- Persistent communication: repeated instance of same proc/data description.

MPI-4:

- *Partitioned sends.*

too obscure to go into.

Review 4

Does this code deadlock?

```
1  for (int p=0; p<nprocs; p++)
2      if (p!=procid)
3          MPI_Send(sbuffer,buflen,MPI_INT,p,0,comm);
4  for (int p=0; p<nprocs; p++)
5      if (p!=procid)
6          MPI_Recv(rbuffer,buflen,MPI_INT,p,0,comm,MPI_STATUS_IGNORE);

/poll "This code deadlocks" "Yes" "No" "Maybe"
```

Review 5

Does this code deadlock?

```
1 int ireq = 0;
2 for (int p=0; p<nprocs; p++)
3     if (p!=procid)
4
5         →MPI_Isend(sbuffers[p], buflen, MPI_INT, p, 0, comm, &(requests[ireq++]))
6     for (int p=0; p<nprocs; p++)
7         if (p!=procid)
8             MPI_Recv(rbuffer, buflen, MPI_INT, p, 0, comm, MPI_STATUS_IGNORE);
9 MPI_Waitall(nprocs-1, requests, MPI_STATUSES_IGNORE);

/poll "This code deadlocks" "Yes" "No" "Maybe"
```

Review 6

Does this code deadlock?

```
1 int ireq = 0;
2 for (int p=0; p<nprocs; p++)
3     if (p!=procid)
4
5         →MPI_Irecv(rbuffers[p], buflen, MPI_INT, p, 0, comm, &(requests[ireq++]));
6 MPI_Waitall(nprocs-1, requests, MPI_STATUSES_IGNORE);
7 for (int p=0; p<nprocs; p++)
8     if (p!=procid)
9         MPI_Send(sbuffer, buflen, MPI_INT, p, 0, comm);

/poll "This code deadlocks" "Yes" "No" "Maybe"
```

Review 7

Does this code deadlock?

```
1 int ireq = 0;
2 for (int p=0; p<nprocs; p++)
3     if (p!=procid)
4
5         →MPI_Irecv(rbuffers[p], buflen, MPI_INT, p, 0, comm, &(requests[ireq++]))
6     for (int p=0; p<nprocs; p++)
7         if (p!=procid)
8             MPI_Send(sbuffer, buflen, MPI_INT, p, 0, comm);
9 MPI_Waitall(nprocs-1, requests, MPI_STATUSES_IGNORE);

/poll "This code deadlocks" "Yes" "No" "Maybe"
```

Where to go from here...

- Derived data types: send strided/irregular/inhomogeneous data
- Sub-communicators: work with subsets of `MPI_COMM_WORLD`
- I/O: efficient file operations
- One-sided communication: ‘just’ put/get the data somewhere
- Process management
- Non-blocking collectives
- Graph topology and neighborhood collectives
- Shared memory