

# MPI Course

Victor Eijkhout

2025 TACC APPI

# 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

1. The SPMD Model *link*
2. Collectives *link*
3. Point-to-point *link*
4. Derived datatypes *link*
5. Communicators *link*
6. MPI I/O *link*
7. One-sided communication *link*
8. Big data *link*
9. Advanced collectives *link*
10. Shared memory *link*
11. Process management *link*
12. Process topologies *link*
13. Trace and performance *link*





# Basics

# Supercomputer clusters

# Cluster setup

Typical cluster:

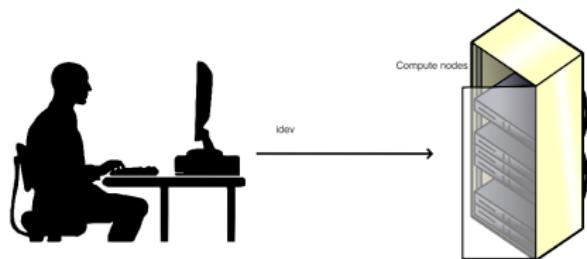
- Login nodes, where you ssh into; usually shared with 100 (or so) other people. You don't run your parallel program there!
- Compute nodes: where your job is run. They are often exclusive to you: no other users getting in the way of your program.

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



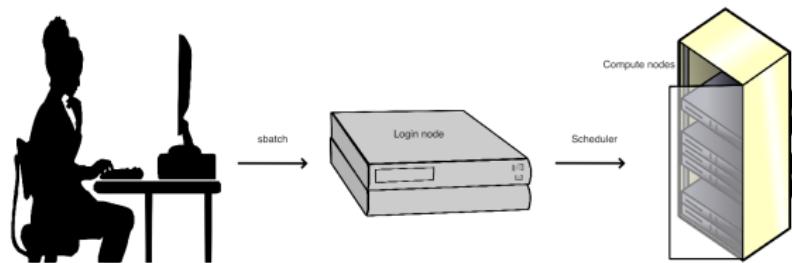
# Interactive run

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



# 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, ...



# Exercise 1

- Connect to your favorite cluster
- Start an interactive session with `idev`;  
what is the hostname? how many users are logged in?
- Run: `ibrun hostname`  
also `ibrun -n 3 hostname`
- Create a job script that will run on 2 nodes;  
again let it run the `hostname` command.

# The SPMD model

# 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,*

# Practicalities



# Lab setup

- Clone the repository  
`https://github.com/VictorEijkhout/TheArtOfHPC_vol2_parallelprogramming`
- Directory: `exercises-mpi-c` or `cxx` or `f` or `f08` or `p` or `mpl`
- 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.
- Type `make exercisename` to compile it
- Run with `ibrun` or `mpexec` (see above)



# Python

Python: setup once per session

```
module load python3
```

No compilation needed. Run:

```
ibrun python3 yourprogram.py
```

(on a compute node!)

# Compiling

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 cetera
```

(For CMake see slide 409.)

# Running

Running your program at TACC:

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

the number of processes is determined by SLURM.

Outside TACC:

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

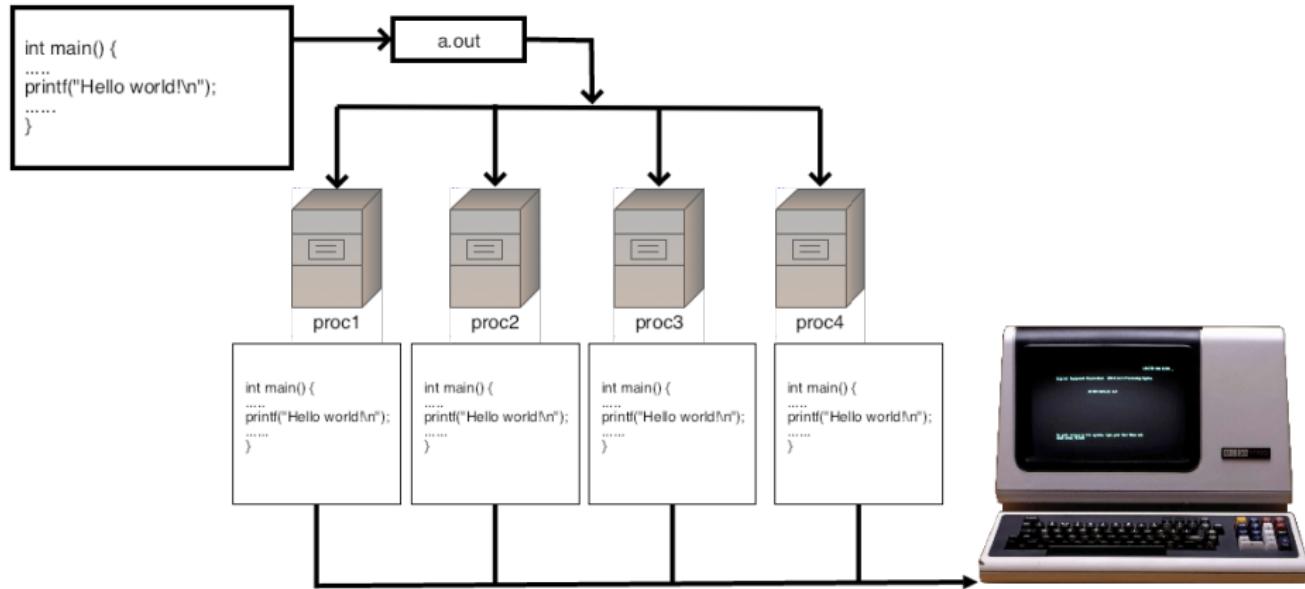


## Exercise 2 (hello)

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

Do: `ibrun python3 hello.py` to execute.

# In a picture



## The MPI worldview: SPMD

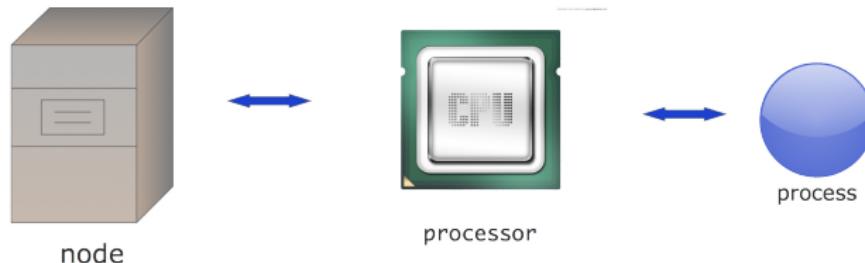
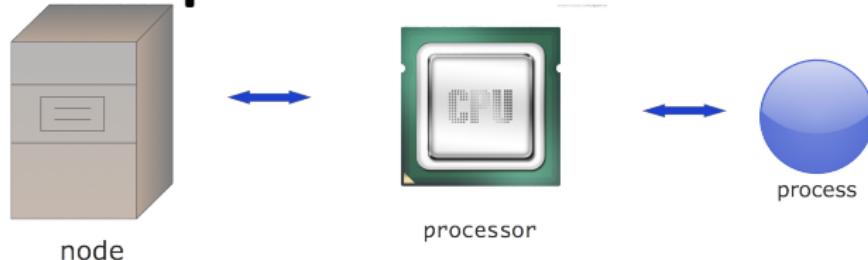
# SPMD

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.

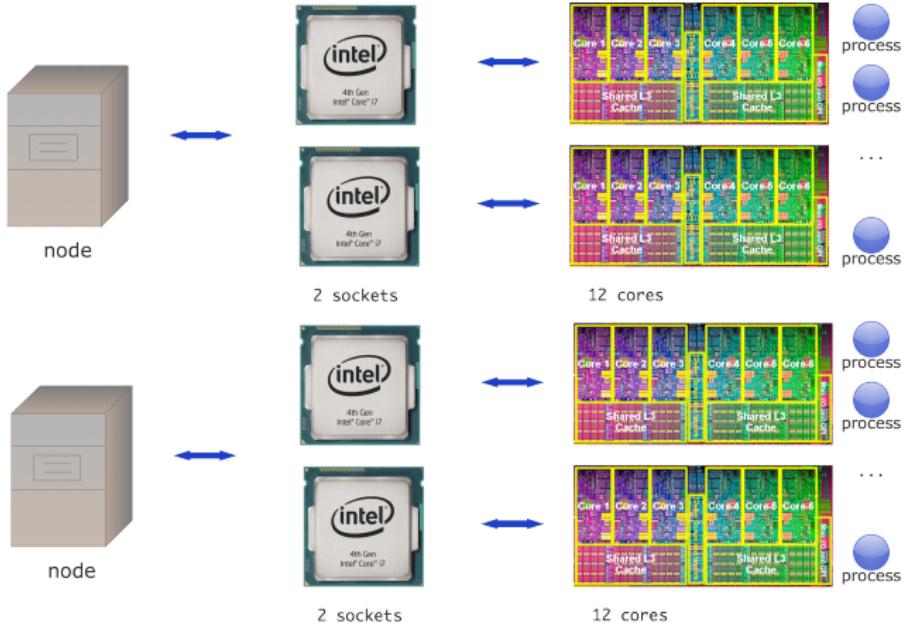
# Computers when MPI was designed



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

⇒ process model, no data sharing.

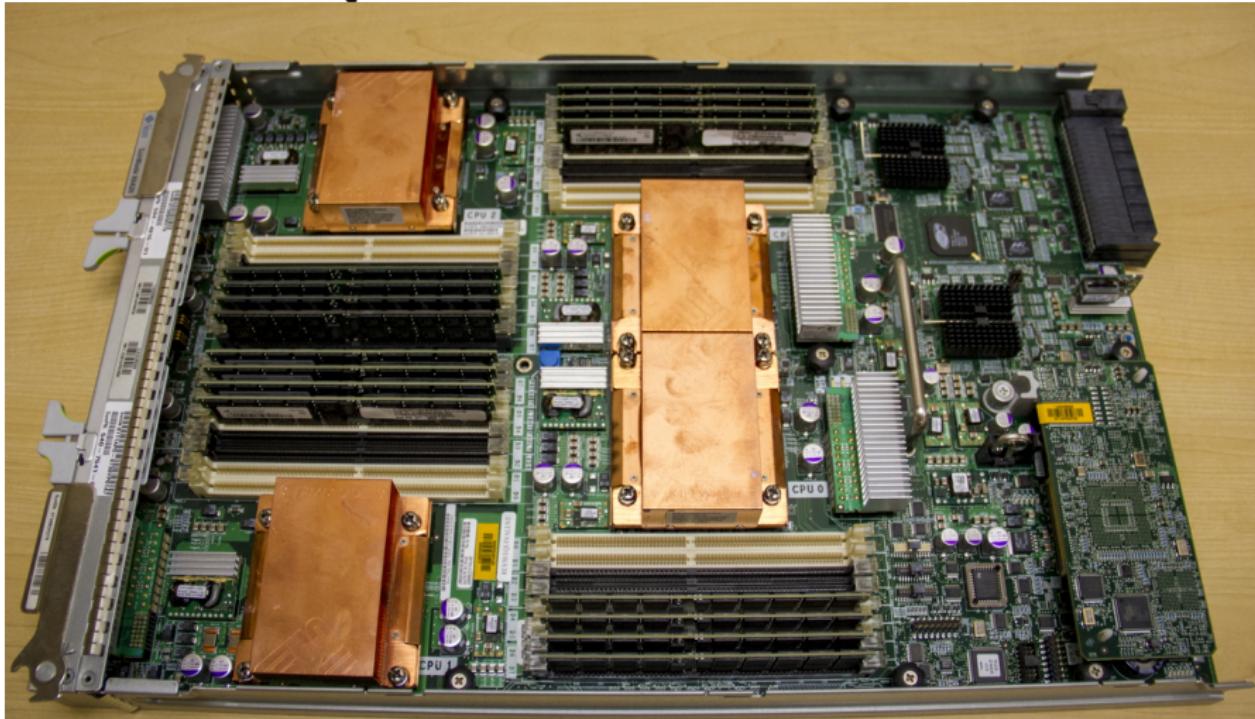
# Pure MPI



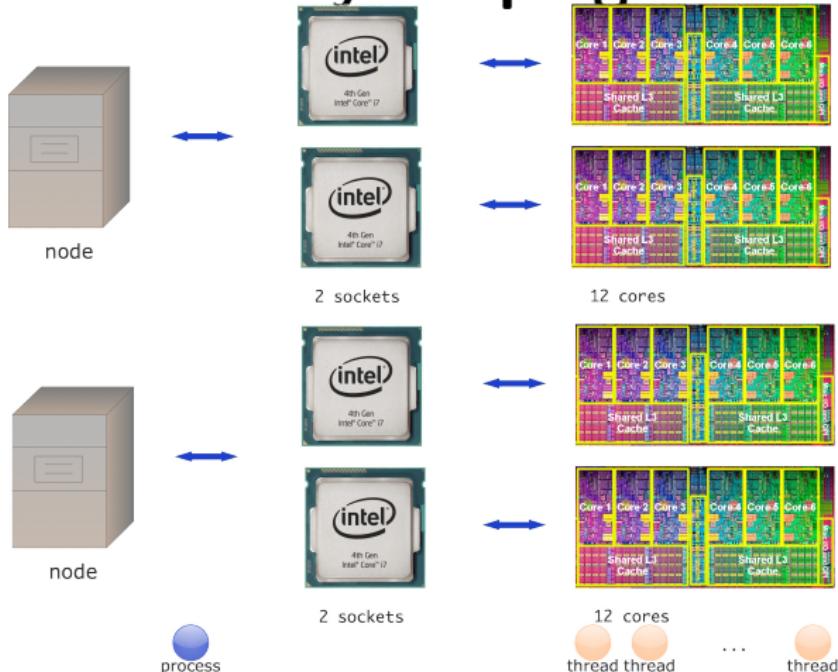
A node has multiple sockets, each with multiple cores.

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

# Quad socket node



# Hybrid programming



Hybrid programming puts a process per node or per socket;  
further parallelism comes from threading.

# 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.

# 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. . .

# 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](http://mpich.org)

# We start learning MPI!



# 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, . . .

# Python init/finalize

Done by the import / at end of the program.

```
from mpi4py import MPI
```

# Exercise 3 (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.

## About library calls and bindings

# 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

# MPI headers: C

You need an include file:

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

This defines all routines and constants.

# Python bindings

- Not part of the standard:  
private project by Lisandro 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



# Ranks

# Process identification

- Processes 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.



# World communicator

For now, the communicator will be *MPI\_COMM\_WORLD*.

C:

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

F:

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

P:

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

MPL:

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



# MPI\_Comm\_size

```
Python:  
MPI.Comm.Get_size(self)
```

## MPI\_Comm\_rank

```
Python:  
MPI.Comm.Get_rank(self)
```

# 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:

```
1 MPI_Comm_size( comm,&nprocs );
```



# 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.)



## Exercise 4 (commrank)

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

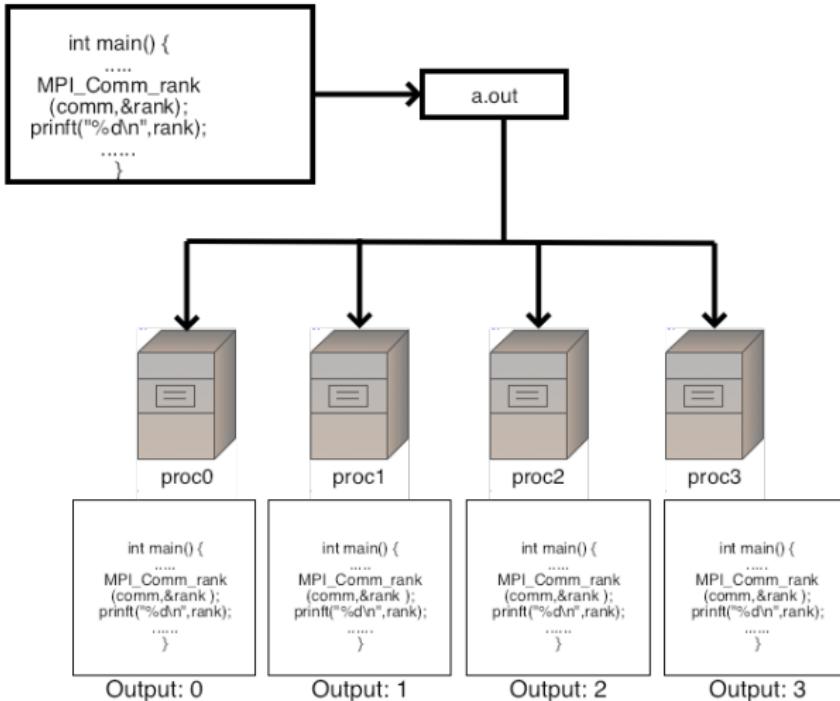
Hello from process 2 out of 5!

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 5 (commrank)

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

# Illustration



# About errors

MPI routines invoke an error handler (slide 412)

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.
- MPL: throws exceptions
- 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.



# MPI\_Get\_processor\_name

```
Python:  
MPI.Get_processor_name()
```

# Exercise 6

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/xxx` and do `make procname`, then `ibrun procname`

# 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);
```

# 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



THE UNIVERSITY OF  
**TEXAS**  
AT AUSTIN



# Processor name: Python

Processes (can) run on physically distinct locations.

```
1 ## procname.py
2 from mpi4py import MPI
3 procname = MPI.Get_processor_name()
```

# Your first useful parallel program

# 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 7 (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 8

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.

# Collectives

# 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*

# 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.

# Collectives

Gathering and spreading information:

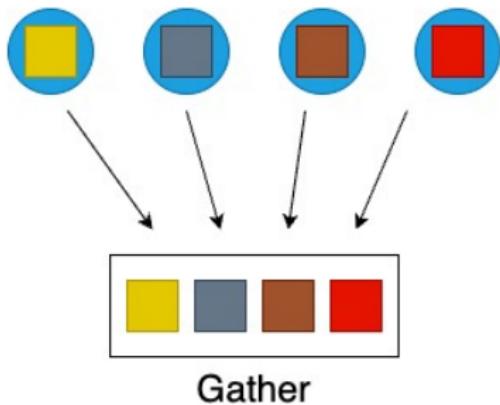
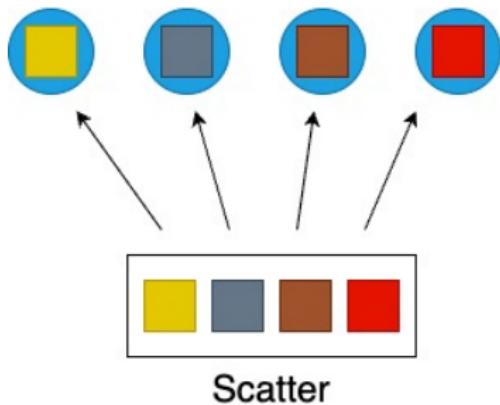
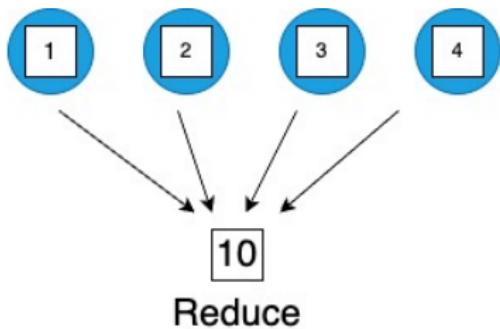
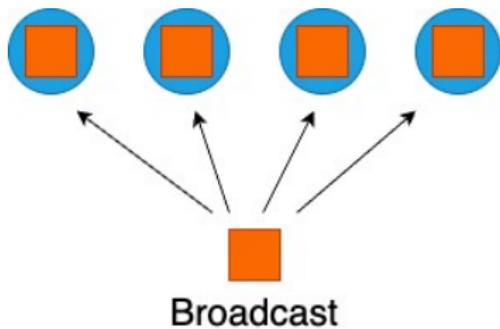
- Every process has data, you want to bring it together;
- One process has data, you want to spread it around.

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.





# Exercise 9

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.



# Allreduce: reduce-to-all

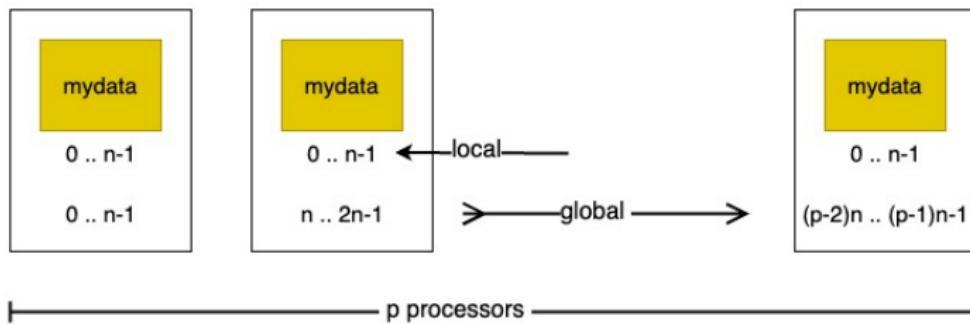
Exercise 2 above contains a common case:  
do a reduction, but everyone needs the result.

- `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.

# Motivation for allreduce

Example: normalizing a vector

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



# Structure of allreduce

- 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?

# 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 process stores just one  $x_i$  value.

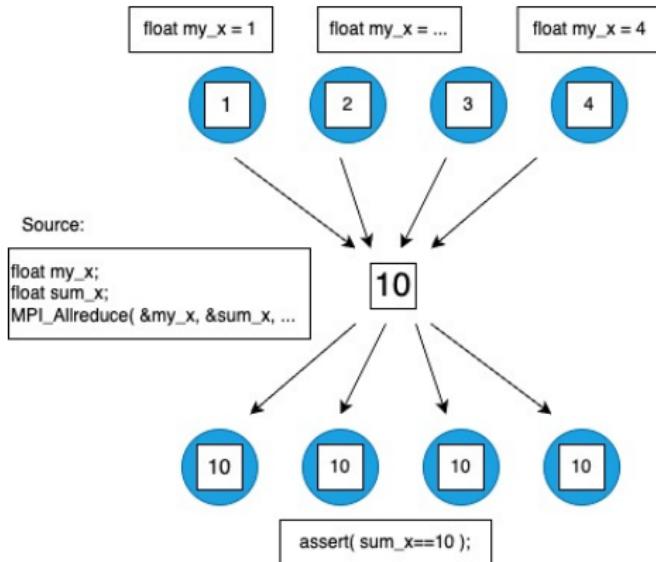
How do we compute this?

1. The calculation of the average  $\mu$  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 processes.
3.  $\sum_i (x_i - \mu)$  is another sum of distributed data, so we need another reduction operation. Might as well use allreduce.



# Conceptual picture

Recall SPMD: every process has the input and output variable



(What actually happens is a different story!)

# 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.

# MPI\_Allreduce

```
Python native:  
recvobj = MPI.Comm.allreduce(self, sendobj, op=SUM)  
Python numpy:  
MPI.Comm.Allreduce(self, sendbuf, recvbuf, Op op=SUM)
```

## Exercise 10 (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.



# Buffers

# 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.



# 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.
- MPL: `long int` and `size_t` supported for layouts.

```
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);
```



# 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



# Exercise 11

Extend the previous exercise to letting each process have an array.

## Collective basics

# Elementary datatypes

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

A bunch more.



# Python datatypes

- Elementary types not needed: type can be deduced from the Numpy buffer
- Buffer / count / datatype triples can be used in exceptional circumstances.

# Reduction operators

MPI type	meaning	applies to
MPI.Op		
<i>MPI_MAX</i>	<code>MPI.MAX</code>	maximum
<i>MPI_MIN</i>	<code>MPI.MIN</code>	minimum
<i>MPI_SUM</i>	<code>MPI.SUM</code>	sum
<i>MPI_PROD</i>	<code>MPI.PROC</code>	product
<i>MPI_REPLACE</i>	<code>MPI.REPLACE</code>	overwrite
<i>MPI_NO_OP</i>	<code>MPI.OP'NULL</code>	no change
<i>MPI_BAND</i>	<code>MPI.LAND</code>	logical and
<i>MPI_BOR</i>	<code>MPI.LOR</code>	logical or
<i>MPI_BXOR</i>	<code>MPI.LXOR</code>	logical xor
<i>MPI_BAND</i>	<code>MPI.BAND</code>	bitwise and
<i>MPI_BOR</i>	<code>MPI.BOR</code>	bitwise or
<i>MPI_BXOR</i>	<code>MPI.BXOR</code>	bitwise xor
<i>MPI_MAXLOC</i>	<code>MPI.MAXLOC</code>	max value and location
<i>MPI_MINLOC</i>	<code>MPI.MINLOC</code>	min value and location
		<i>MPI_DOUBLE_INT</i> and such



# Reduction to single process

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

# 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.



# In-place operations

```
1 // allreduceinplace.c
2 for (int irand=0; irand<nrandoms; irand++)
3     myrandoms[irand] = (float) rand()/(float)RAND_MAX;
4 // add all the random variables together
5 MPI_Allreduce(MPI_IN_PLACE,myrandoms,
6                 nrandoms,MPI_FLOAT,MPI_SUM,comm);
```



# More in-place operations

```
1 if (procno==root)
2   MPI_Reduce(MPI_IN_PLACE,myrandoms,
3             nrandoms,MPI_FLOAT,MPI_SUM,root,comm);
4 else
5   MPI_Reduce(myrandoms,MPI_IN_PLACE,
6             nrandoms,MPI_FLOAT,MPI_SUM,root,comm);
```

or

```
1 float *sendbuf,*recvbuf;
2 if (procno==root) {
3   sendbuf = MPI_IN_PLACE; recvbuf = myrandoms;
4 } else {
5   sendbuf = myrandoms; recvbuf = MPI_IN_PLACE;
6 }
7 MPI_Reduce(sendbuf,recvbuf,
8            nrandoms,MPI_FLOAT,MPI_SUM,root,comm);
```



# 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*



# Gauss-Jordan elimination

<https://youtu.be/aQYuwatlWME>

# MPI\_Bcast

```
Python native:  
rbuf = MPI.Comm.bcast(self, obj=None, int root=0)  
Python numpy:  
MPI.Comm.Bcast(self, buf, int root=0)
```

## Exercise 12 (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) 13

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

# Scan



# 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:  $\text{MPI\_Scan}$  inclusive, and  $\text{MPI\_Exscan}$  exclusive.

# In vs Exclusive

process :      0                  1                  2                   $\cdots$                    $p - 1$

data :       $x_0$                    $x_1$                    $x_2$                    $\cdots$                    $x_{p-1}$

inclusive :       $x_0$                    $x_0 \oplus x_1$        $x_0 \oplus x_1 \oplus x_2$        $\cdots$        $\oplus_{i=0}^{p-1} x_i$

exclusive :    unchanged       $x_0$                    $x_0 \oplus x_1$        $\cdots$        $\oplus_{i=0}^{p-2} x_i$

# MPI\_Scan

Python:

```
res = Intracomm.scan( sendobj=None,recvobj=None,op=MPI.SUM)
res = Intracomm.exscan( sendobj=None,recvobj=None,op=MPI.SUM)
```



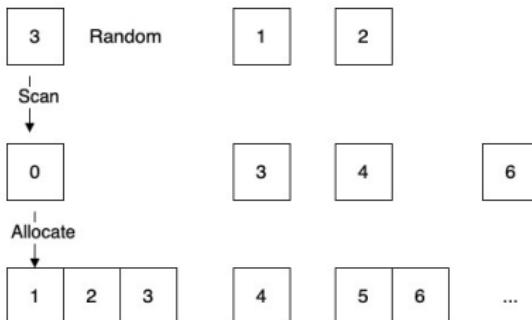
# MPI\_Exscan

# Exercise 14 (scangather)

- Let each process compute a random value  $n_{\text{local}}$ , 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 14.)



## **Gather/Scatter, Barrier, and others**

# MPI\_Gather

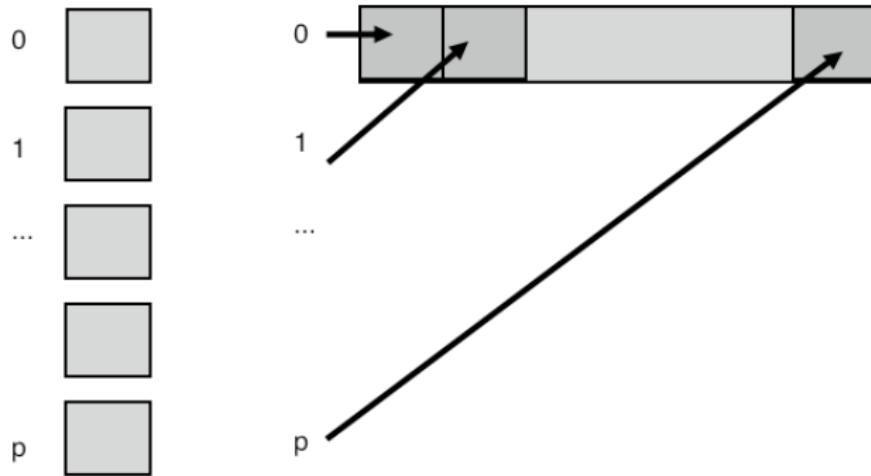
```
Python:  
MPI.Comm.Gather  
    (self, sendbuf, recvbuf, int root=0)
```

# MPI\_Scatter

# 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.

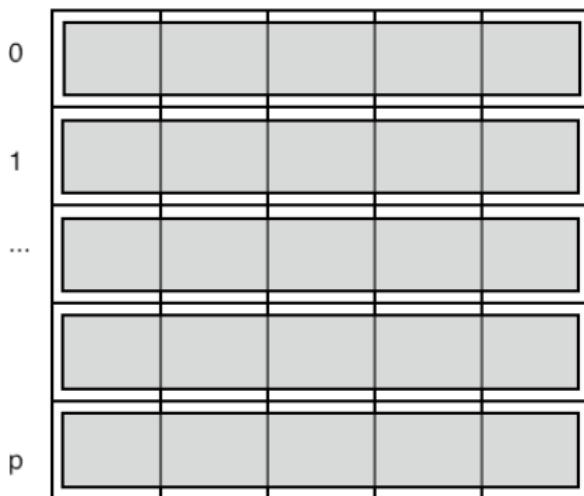
# Gather pictured



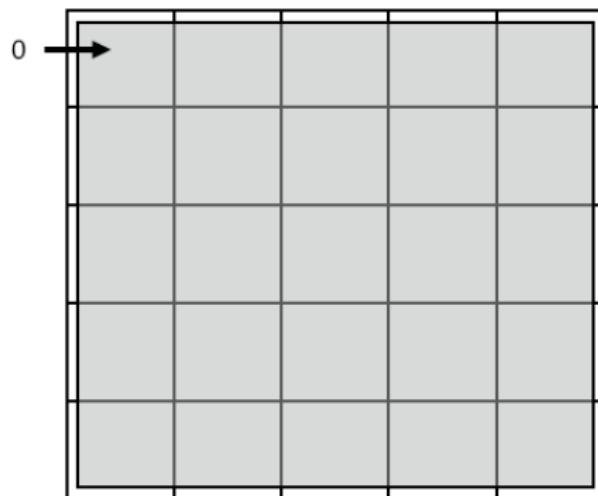
# Popular application of gather

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

distributed matrix



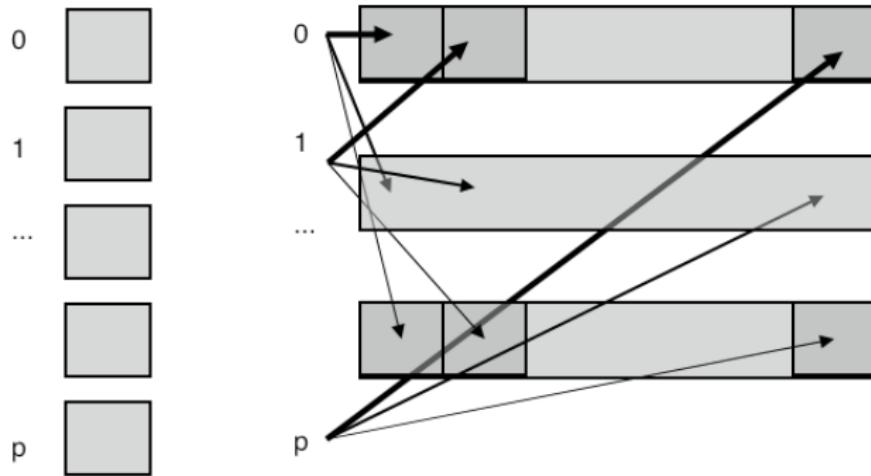
gathered matrix



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

# MPI\_Allgather

# Allgather pictured



# V-type collectives

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

# MPI\_Gatherv

```
Python:  
Gatherv(self, sendbuf, [recvbuf,counts], int root=0)
```

# Exercise 15 (scangather)

Take the code from exercise 14 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

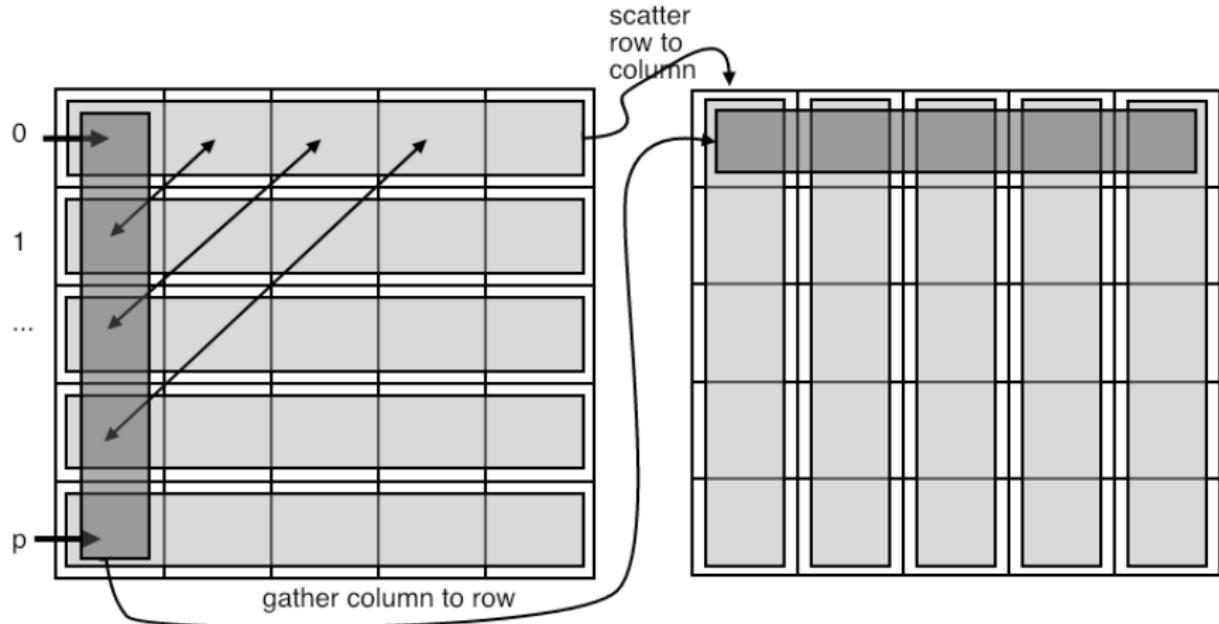
An *MPI\_Scatter* call puts the same data on each process

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

# All-to-all

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

# Data transposition



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

# All-to-allv

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

# Radix sort

Sort 4 numbers on two processes:

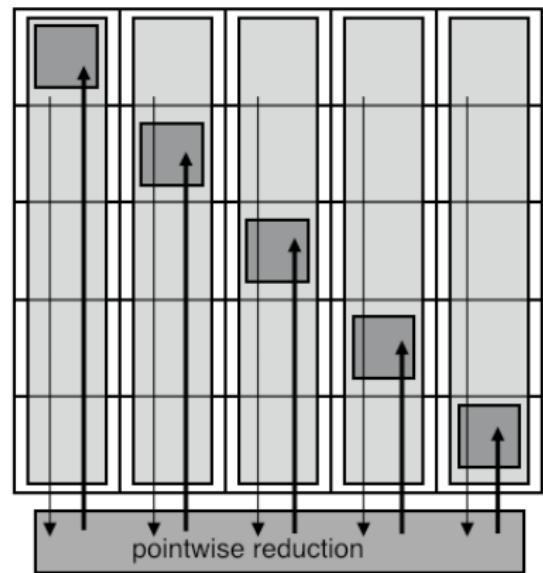
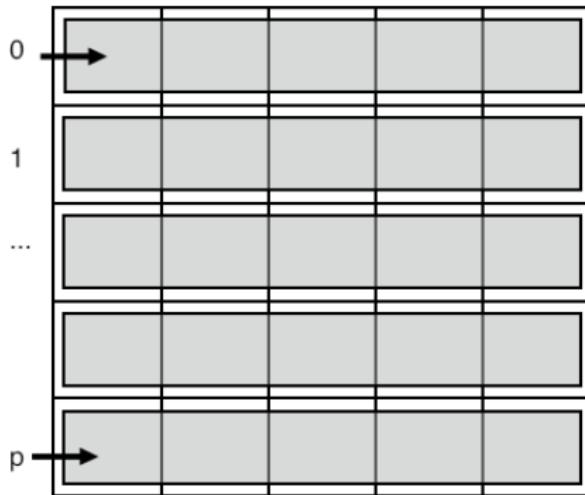
array binary	proc0		proc1	
	2	5	7	1
010	101	111	001	
stage 1				
last digit	0	1	1	1
(this serves as bin number)				
sorted	010		101	111 001
stage 2				
next digit	1		0 1	0
(this serves as bin number)				
sorted	101	001	010	111
stage 3				
next digit	1	0	0 1	
(this serves as bin number)				
sorted	001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

# 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.

# 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



# 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

# 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*

# MPI\_Op\_create

```
Python:  
MPI.Op.create(cls,function,bool commute=False)
```

# Example

Smallest nonzero:

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

# Review 2

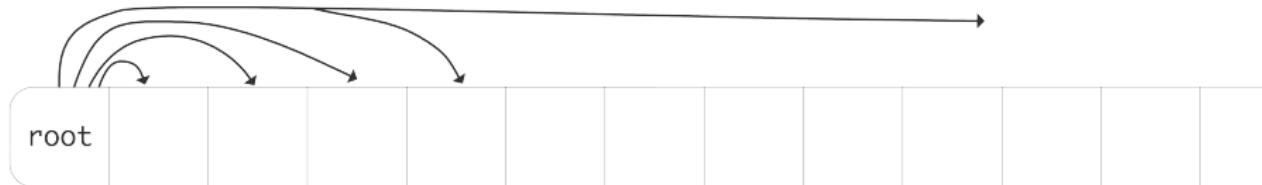
The  $\|\cdot\|_2$  norm (sum of squares) needs a custom operator.

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

## **Performance of collectives**

# Naive realization of collectives

Broadcast:



Single message:

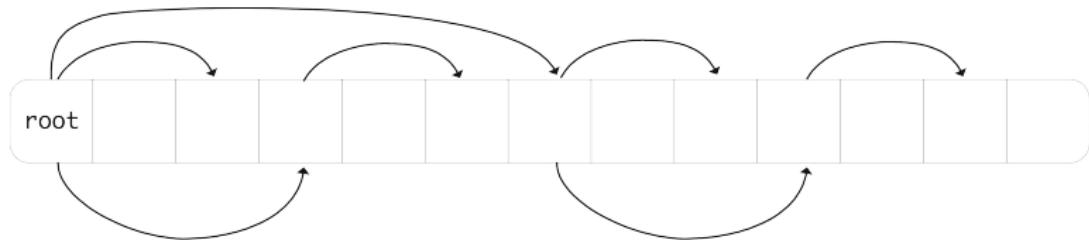
$$\alpha = \text{message startup} \approx 10^{-6} s, \quad \beta = \text{time per word} \approx 10^{-9} s$$

- Time for message of  $n$  words:

$$\alpha + \beta n$$

- Time for collective? Can you improve on that?

# Better implementation of collective



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

# Implementation of Reduce

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

# 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**

# 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 16 (onenorm)

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

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

# Point-to-point communication

# 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**

# 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.

# Ping-pong

A sends to B, B sends back to A

What is the code for A? For B?

# 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 }
```



# Sample send and recv calls

```
1 double x[10],y[10];
2 MPI_Send( x,10,MPI_DOUBLE, tgt,0,comm );
3 MPI_Status status;
4 MPI_Recv( y,10,MPI_DOUBLE, src,0,comm,&status );
```

# MPI\_Send

```
Python:  
MPI.Comm.send(self, obj, int dest, int tag=0)  
Python numpy:  
MPI.Comm.Send(self, buf, int dest, int tag=0)
```

# MPI\_Recv

```
Python native:  
recvbuf = Comm.recv(self, buf=None, int source=ANY_SOURCE, int tag=ANY_TAG,  
                     Status status=None)  
Python numpy:  
Comm.Recv(self, buf, int source=ANY_SOURCE, int tag=ANY_TAG,  
           Status status=None)
```

# 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 17 (pingpong)

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

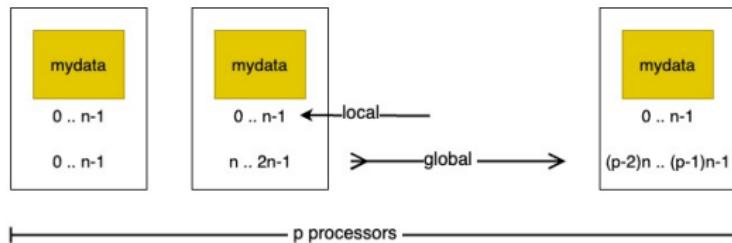
- Run multiple ping-pongs (say a thousand) and put the timer around the loop. The first run may take longer; try to discard it.
- Run your code with the two communicating processes first on the same node, then on different nodes. Do you see a difference?
- Then modify the program to use longer messages.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

For bonus points, can you do a regression to determine  $\alpha, \beta$ ?

# Distributed data

# Distributed data

Distributed array: each process stores disjoint local part



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

# 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) 18

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

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

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



# 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?

# (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**

# 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:  $u_i \equiv u(ih)$ .

# Motivation (continued)

Equations

$$\left\{ \begin{array}{l} -u_{i-1} + 2u_i - u_{i+1} = h^2 f(x_i) \\ 2u_1 - u_2 = h^2 f(x_1) + u_0 \\ 2u_n - u_{n-1} = h^2 f(x_n) + u_{n+1}. \end{array} \right. \quad 1 < i < n$$
$$\begin{pmatrix} 2 & -1 & & \emptyset \\ -1 & 2 & -1 & \\ \emptyset & \ddots & \ddots & \ddots \end{pmatrix} \begin{pmatrix} u_1 \\ u_2 \\ \vdots \end{pmatrix} = \begin{pmatrix} h^2 f_1 + u_0 \\ h^2 f_2 \\ \vdots \end{pmatrix} \quad (1)$$

So we are interested in sparse/banded matrices.

# Matrix vector product

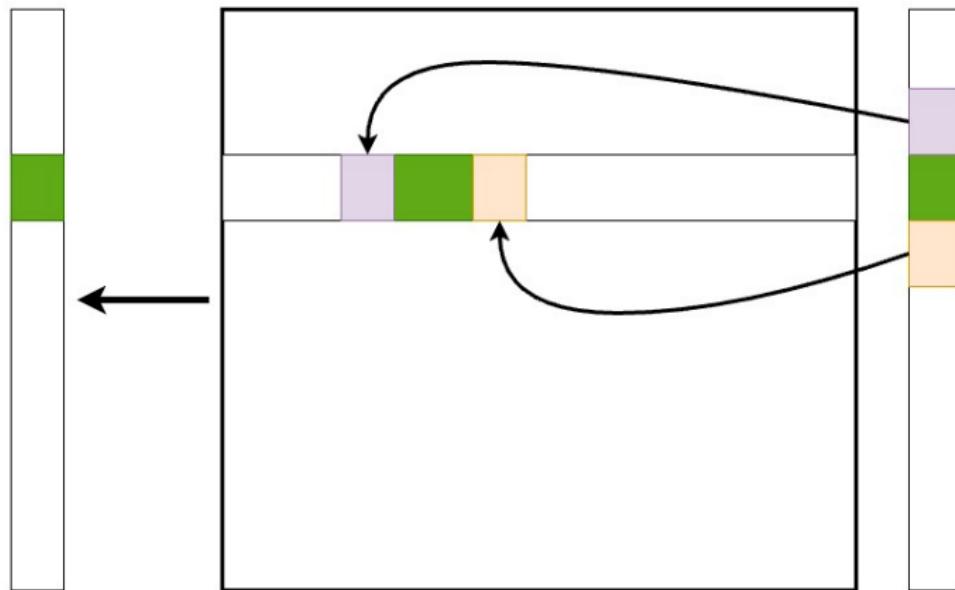
Most common operation: matrix vector product

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

- 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

# Partitioned matrix-vector product

We need a point-to-point mechanism:



each process with ones before/after it.

# Operating on distributed data

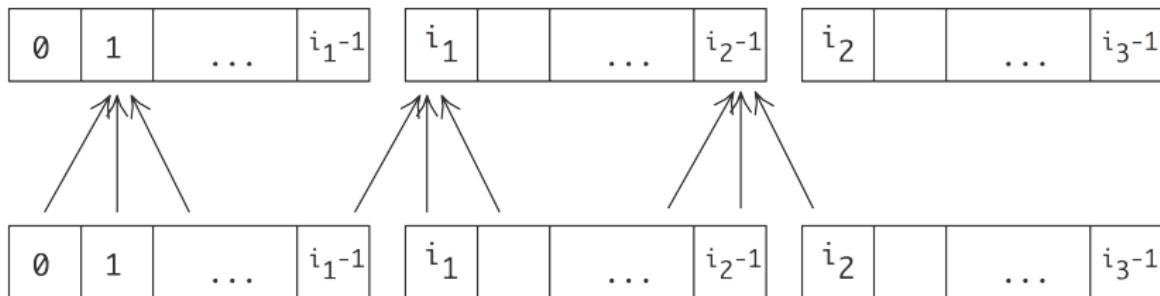
Array of numbers  $x_i : i = 0, \dots, N$

compute

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

'owner computes'

This leads to communication:



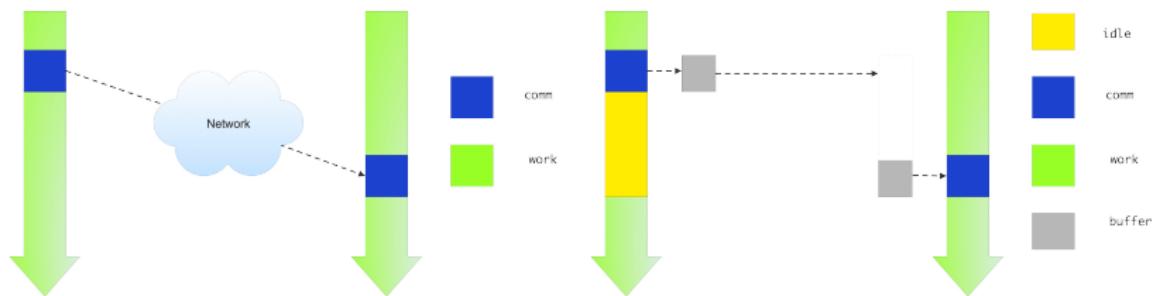
so we need a point-to-point mechanism.

## **Blocking communication**

# 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.

# 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)



# 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 19

(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:

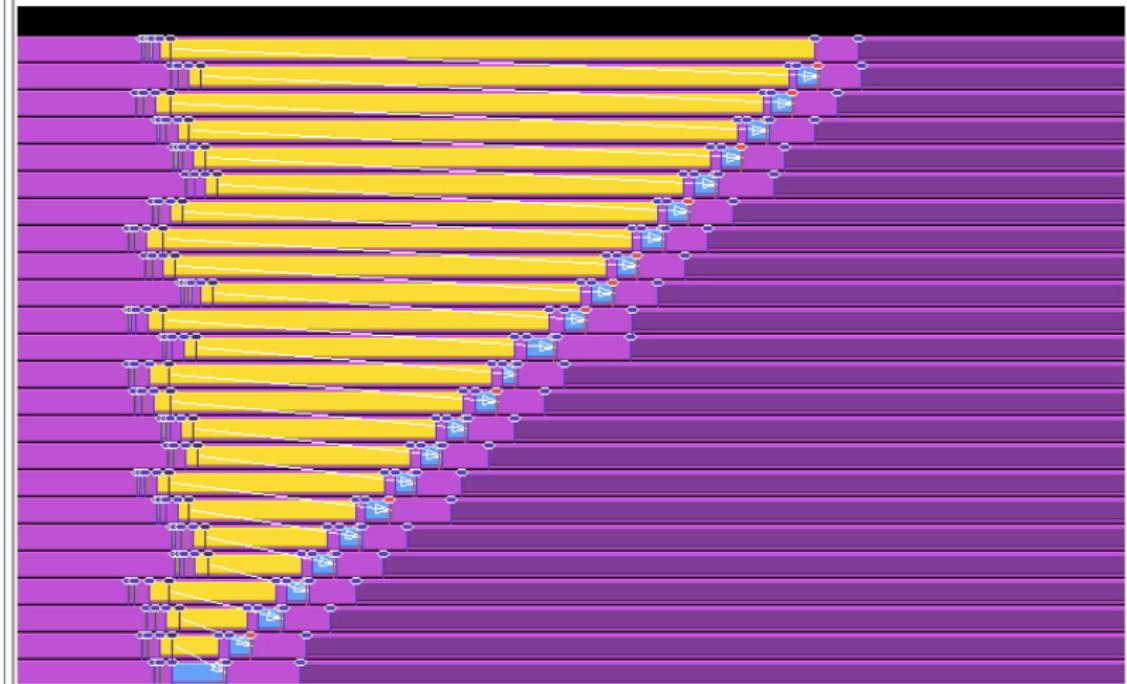
1. Give the paper to your right neighbor;
2. Accept the paper from your left neighbor.

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.

# TAU trace: serialization

TimeLines



## 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 20 (rightsend)

Implement the above algorithm using `MPI_Send` and `MPI_Recv` calls. Run the code, and use TAU to reproduce the trace output of figure 156. If you don't have TAU, can you show this serialization behavior using timings, for instance running it on an increasing number of processes?

# Synchronous send

Synchronous send:

- sender and receiver synchronize
- No ‘eager’ sends
- ⇒ enforced always blocking behavior

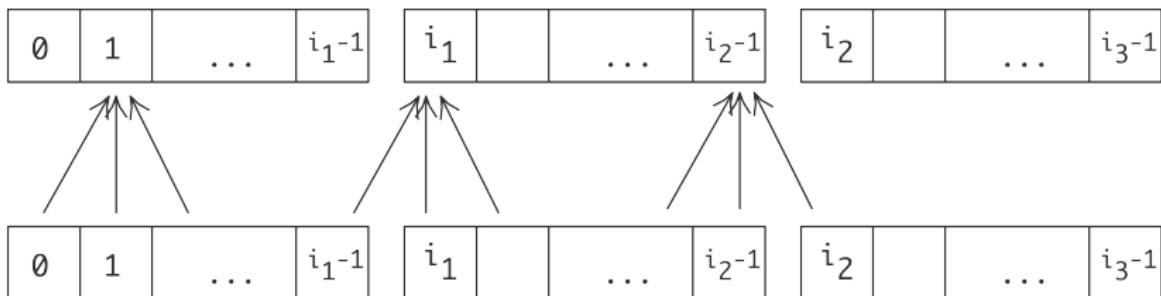
# MPI\_Ssend

## Pairwise exchange

# Operating on distributed data

Take another look:

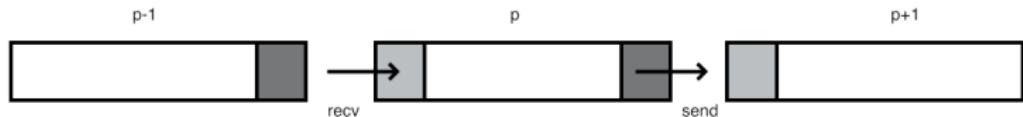
$$y_i = x_{i-1} + x_i + x_{i+1}: i = 1, \dots, N - 1$$



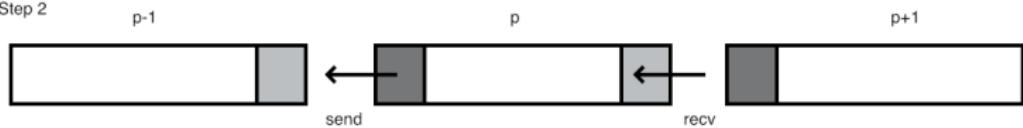
- One-dimensional data and linear process numbering;
- Operation between neighboring indices: communication between neighboring processes.

# 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!

# Sendrecv

Instead of separate send and receive: use

*MPI\_Sendrecv*

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

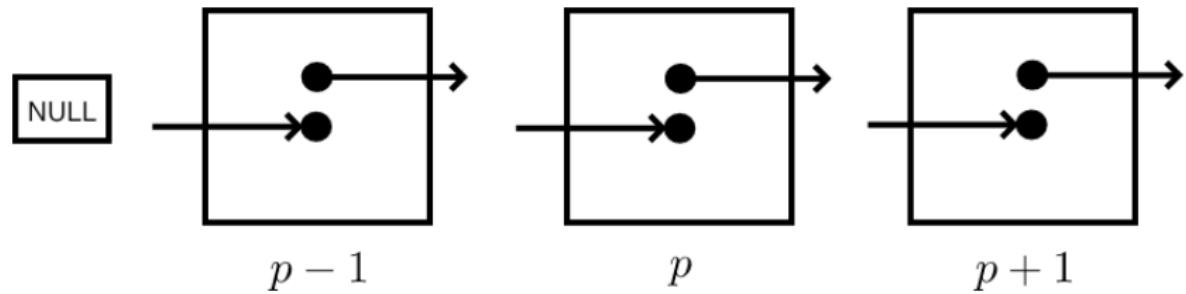
(Also: *MPI\_Sendrecv\_replace* with single buffer.)

# MPI\_Sendrecv

```
Python:  
Sendrecv(self,  
         sendbuf, int dest, int sendtag=0,  
         recvbuf=None, int source=ANY_SOURCE, int recvtag=ANY_TAG,  
         Status status=None)
```

# SPMD picture

What does process  $p$  do?



# 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.)



# 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 else if ( /* i am last */
4   Sendrecv( to=NULL,   from=left );
5 else
6   Sendrecv( to=right, from=left );
```

```
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 );
```

But:

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



# Exercise (optional) 21 (rightsend)

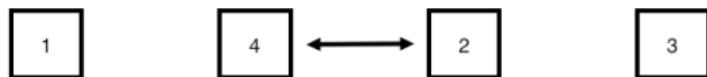
Revisit exercise 19 and solve it using `MPI_Sendrecv`.

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 22 (sendrecv)

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

## Odd-even transposition sort



↔ transpose performed  
↔ no transpose needed

Odd-even transposition sort on 4 elements.

# Exercise (optional) 23

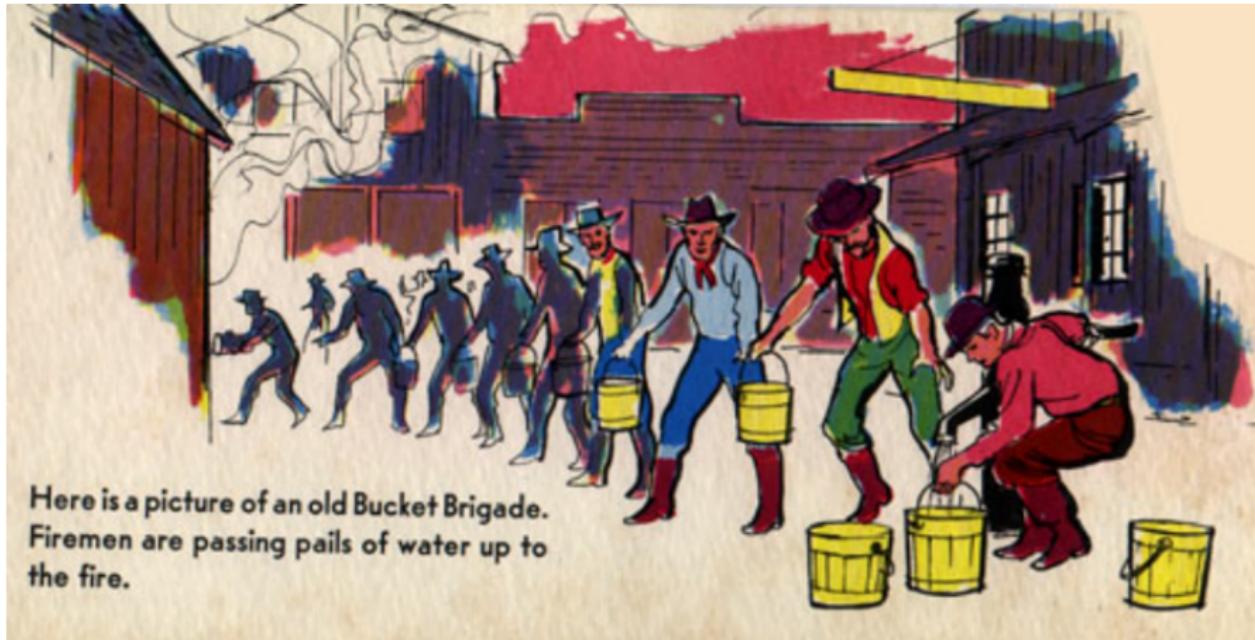
A very simple sorting algorithm is *swap sort* or *odd-even transposition sort*: pairs of processors compare data, and if necessary exchange. The elementary step is called a *compare-and-swap*: in a pair of processors each sends their data to the other; one keeps the minimum values, and the other the 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 171.

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.

# 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 24 (bucketblock)

Take the code of exercise 20 and modify it so that the data from process zero gets propagated to every process. Specifically, compute all partial sums  $\sum_{i=0}^p i^2$ :

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

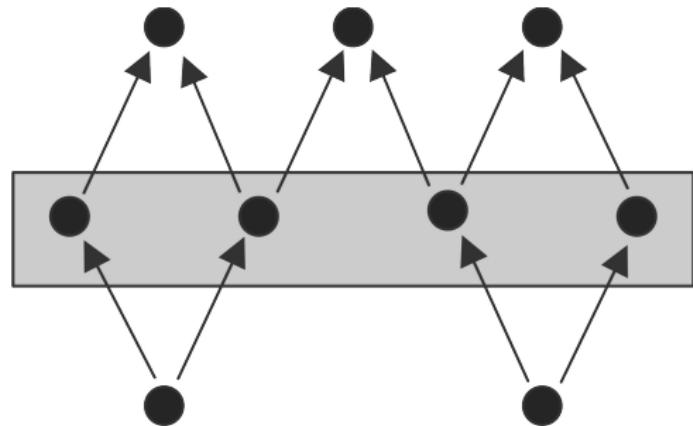
Use `MPI_Send` and `MPI_Recv`; make sure to get the order right.

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**

# Sending with irregular connections

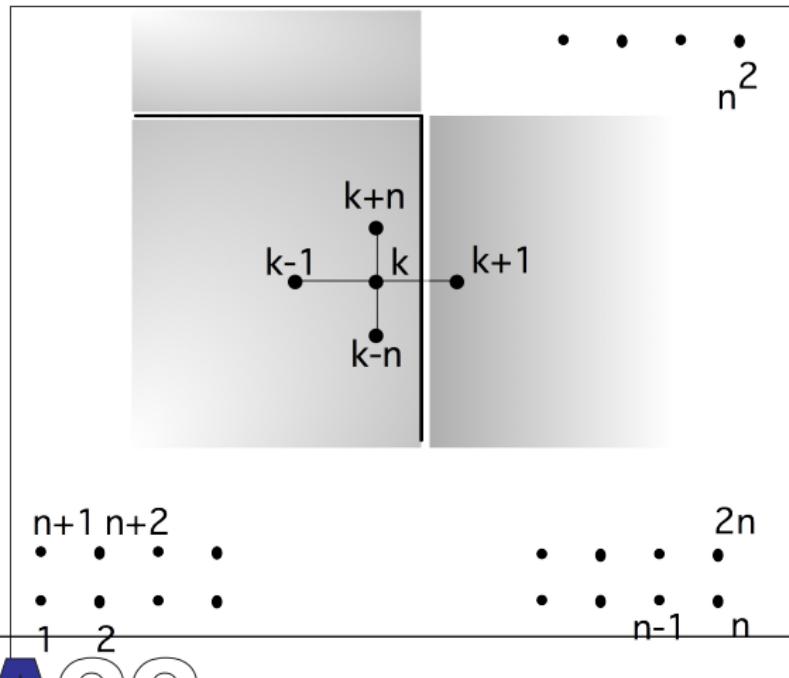
Graph operations:



## **Communicating other than in pairs**

# PDE, 2D case

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

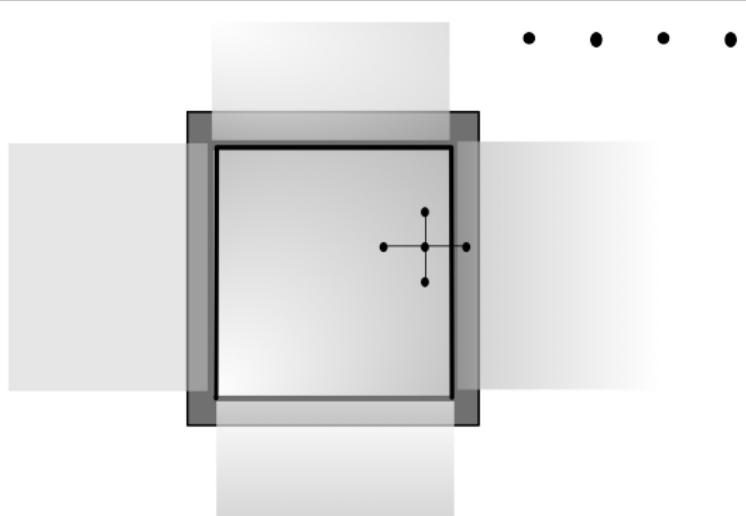


# PDE matrix

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

# Halo region

The halo region of a process, induced by a stencil



$n+1$   $n+2$   
• • • •  
• • • •  
1 2

• • • •  
• • • •  
 $n-1$   $n$   
 $2n$

# 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.

# Non-blocking send/recv

- $\text{MPI\_Isend}$  /  $\text{MPI\_Irecv}$  does not send/receive:
- They declare a buffer.
- The buffer contents are there after a wait call.
- In between the  $\text{MPI\_Isend}$  and  $\text{MPI\_Wait}$  the data may not have been sent.
- In between the  $\text{MPI\_Irecv}$  and  $\text{MPI\_Wait}$  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( ... );
```

# MPI\_Isend

Python:

```
request = MPI.Comm.Isend(self, buf, int dest, int tag=0)
```

# MPI\_Irecv

```
Python native:  
recvbuf = Comm.irecv(self, buf=None, int source=ANY_SOURCE, int tag=ANY_TAG,  
    Request request=None)  
Python numpy:  
Comm.Irecv(self, buf, int source=ANY_SOURCE, int tag=ANY_TAG,  
    Request status=None)
```

# Request waiting

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`

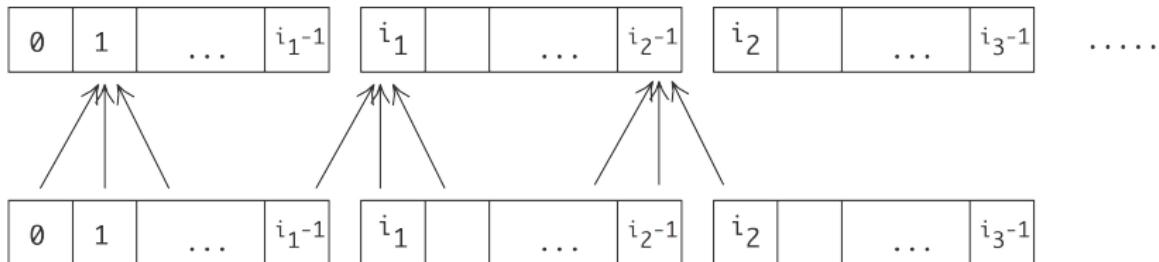
# Exercise 25 (isendirecv)

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 the end points?)



# 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.

# 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(.....)
```

# Pitfalls

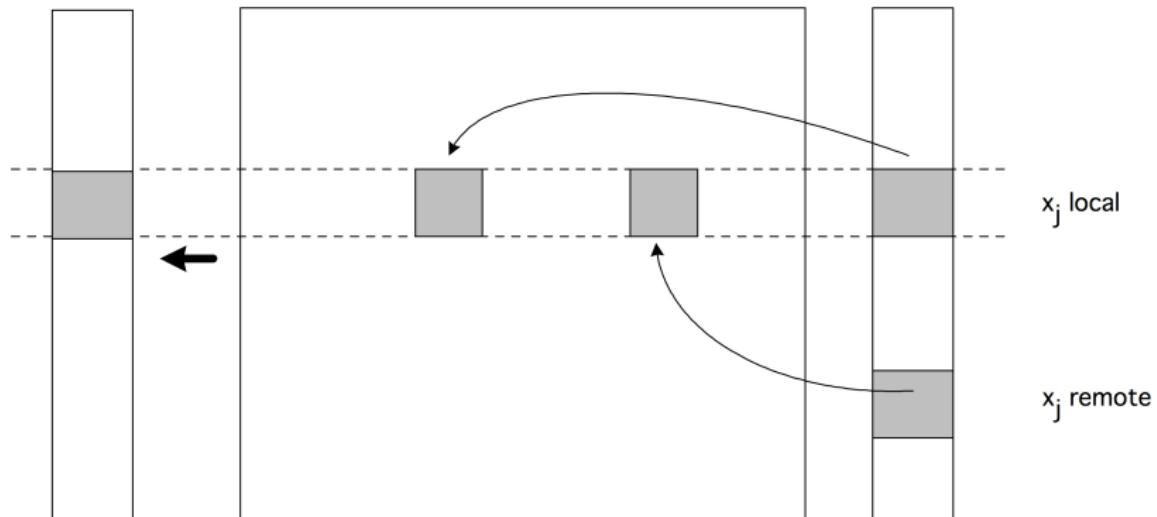
- Strictly one request/wait per `isend`/`irecv`:  
can not use one request for multiple simultaneous `isends`
- Some people argue:  
*Wait for the send is not necessary: if you wait for the receive,  
the message has arrived safely*

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

# Matrices in parallel

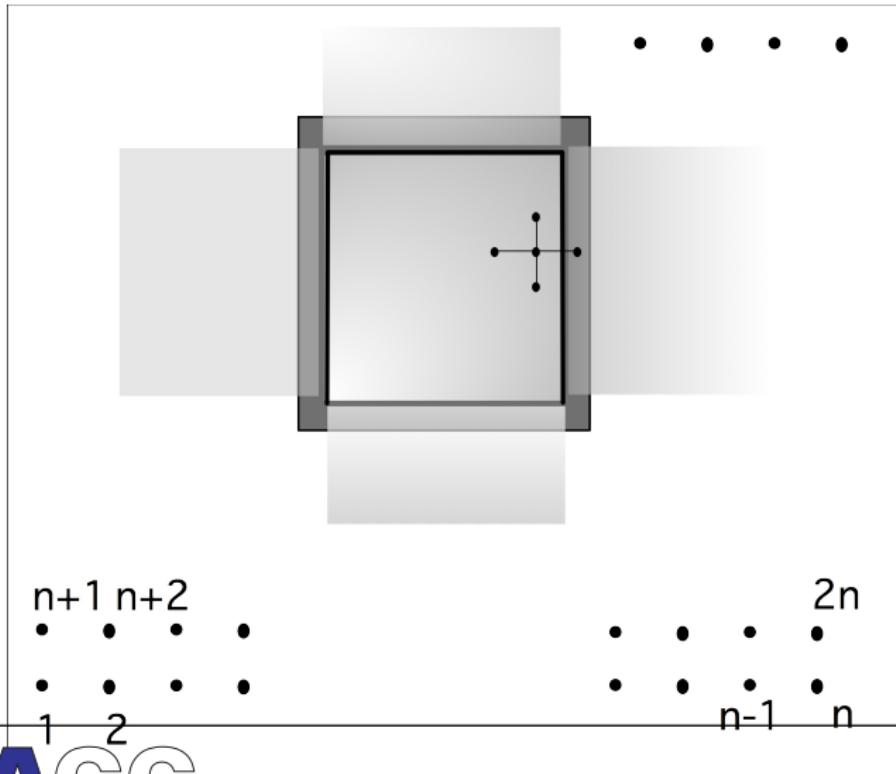
$$y \leftarrow Ax$$

and  $A, x, y$  all distributed:



# Hiding the halo

Interior of a process domain can overlap with halo transfer:



# 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 26 (isendirecvarray)

Take your code of exercise 25 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.

# Mix and match

You can match blocking and non-blocking:

```
1 if ( /* I am Process A */ ) {  
2     MPI_Irecv( /* from B */, &req1 );  
3     MPI_Isend( /* to B */, &req2 );  
4     MPI_Waitall( /* requests 1 and 2 */ );  
5 } else if ( /* I am Process B */ ) {  
6     MPI_Recv( /* from A */ );  
7     MPI_Send( /* to A */ );  
8 }
```



# Test: non-blocking wait

- Post non-blocking receives
- test on the request(s) for incoming messages
- if nothing comes in, do local work

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

Local operation.

Also *MPI\_Testall* et cetera.



# Probe for message

Is there a message?

```
1 // probe.c
2 if (procno==receiver) {
3     MPI_Status status;
4     MPI_Probe(sender,0,comm,&status);
5     int count;
6     MPI_Get_count(&status,MPI_FLOAT,&count);
7     float recv_buffer[count];
8     MPI_Recv(recv_buffer,count,MPI_FLOAT, sender,0,comm,MPI_STATUS_IGNORE);
9 } else if (procno==sender) {
10     float buffer[buffer_size];
11     ierr = MPI_Send(buffer,buffer_size,MPI_FLOAT, receiver,0,comm); CHK(ierr);
12 }
```

(Also non-blocking `MPI_Iprobe`.)

These commands force global progress.



# 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) 27**

## **(bucketpipenonblock)**

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

## Exercise 28 (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`

# The wheel of reinvention

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

# More sends and receive

- $\text{MPI_Bsend}$ ,  $\text{MPI_Ibsend}$ : buffered send
- $\text{MPI_Ssend}$ ,  $\text{MPI_Issend}$ : synchronous send
- $\text{MPI_Rsend}$ ,  $\text{MPI_Irsend}$ : 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         MPI_Isend(sbuffers[p],buflen,MPI_INT,p,0,comm,&(requests[ireq++]));
5 for (int p=0; p<nprocs; p++)
6     if (p!=procid)
7         MPI_Recv(rbuffer,buflen,MPI_INT,p,0,comm,MPI_STATUS_IGNORE);
8 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         MPI_Irecv(rbuffers[p],buflen,MPI_INT,p,0,comm,&(requests[ireq++]));
5 MPI_Waitall(nprocs-1,requests,MPI_STATUSES_IGNORE);
6 for (int p=0; p<nprocs; p++)
7     if (p!=procid)
8         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         MPI_Irecv(rbuffers[p],buflen,MPI_INT,p,0,comm,&(requests[ireq++]));
5 for (int p=0; p<nprocs; p++)
6     if (p!=procid)
7         MPI_Send(sbuffer,buflen,MPI_INT,p,0,comm);
8 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

# Intermediate topics

# Justification

MPI basic concepts suffice for many applications. The Intermediate Topics section deals with more complicated data, process groups, file I/O, and the basics of one-sided communication.

# Derived Datatypes

# Overview

In this section you will learn about derived data types.

Commands learned:

- *MPI\_Type\_contiguous / vector / indexed / struct MPI\_Type\_create\_subarray*
- *MPI\_Pack / MPI\_Unpack*
- F90 types

# Motivation: datatypes in MPI

All examples so far:

- contiguous buffer
- elements of single type

We need data structures with gaps, or heterogeneous types.

- Send real or imaginary parts out of complex array.
- Gather/scatter cyclicly.
- Send *struct* or *Type* data.

MPI allows for recursive construction of data types.



# Datatype topics

- Elementary types: built-in.
- Derived types: user-defined.
- Packed data: not really a datatype.

# Elementary datatypes

C/C++	Fortran
<i>MPI_CHAR</i>	<i>MPI_CHARACTER</i>
<i>MPI_UNSIGNED_CHAR</i>	
<i>MPI_SIGNED_CHAR</i>	<i>MPI_LOGICAL</i>
<i>MPI_SHORT</i>	
<i>MPI_UNSIGNED_SHORT</i>	
<i>MPI_INT</i>	<i>MPI_INTEGER</i>
<i>MPI_UNSIGNED</i>	
<i>MPI_LONG</i>	
<i>MPI_UNSIGNED_LONG</i>	
<i>MPI_FLOAT</i>	<i>MPI_REAL</i>
<i>MPI_DOUBLE</i>	<i>MPI_DOUBLE_PRECISION</i>
<i>MPI_LONG_DOUBLE</i>	
	<i>MPI_COMPLEX</i>
	<i>MPI_DOUBLE_COMPLEX</i>



# How to use derived types

Create, commit, use, free:

```
1 MPI_Datatype newtype;
2 MPI_Type_xxx( ... oldtype ... &newtype);
3 MPI_Type_commit ( &newtype );
4
5 // code using the new type
6
7 MPI_Type_free ( &newtype );

1 Type(MPI_Datatype) :: newtype ! F2008
2 Integer           :: newtype ! F90
```

The *oldtype* can be elementary or derived.  
Recursively constructed types.



# Contiguous type

```
1 int MPI_Type_contiguous(  
2     int count, MPI_Datatype old_type, MPI_Datatype *new_type_p)
```



This one is indistinguishable from sending `count` instances of the `old_type`.

# Example: non-contiguous data

Matrix in column storage:

- Columns are contiguous
- Rows are not contiguous

Logical:

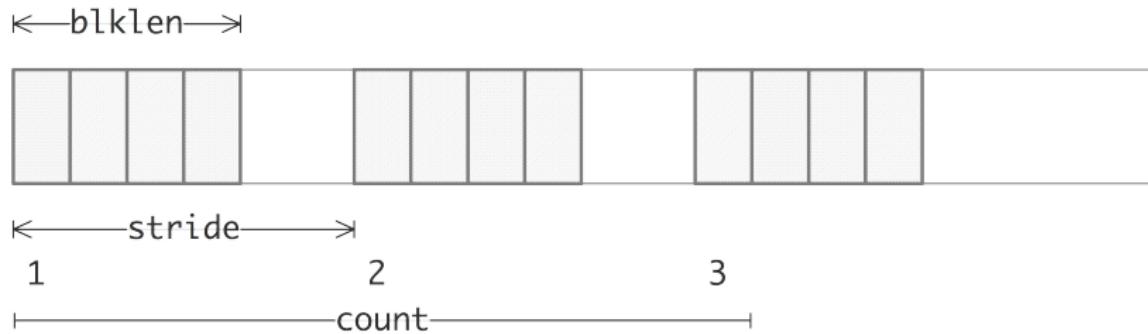
(1,1)	(1,2)	
(2,1)		
(3,1)		

Physical:

(1,1)	(2,1)	(3,1)	...	(1,2)	...
-------	-------	-------	-----	-------	-----

# Vector type

```
1 int MPI_Type_vector(  
2     int count, int blocklength, int stride,  
3     MPI_Datatype old_type, MPI_Datatype *newtype_p  
4 );
```



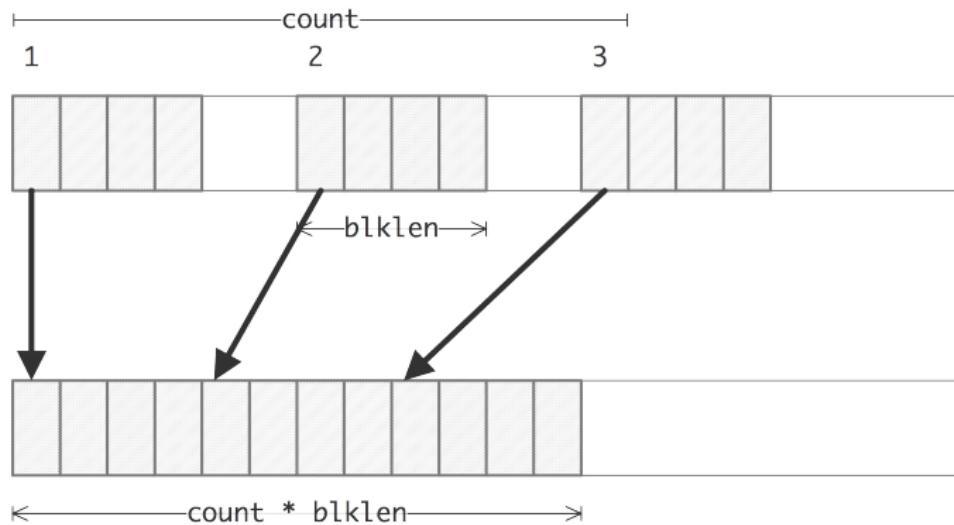
Used to pick a regular subset of elements from an array.

# Different send and receive types

Send and receive type can differ. Example:

Sender type: vector

receiver type: contiguous or elementary



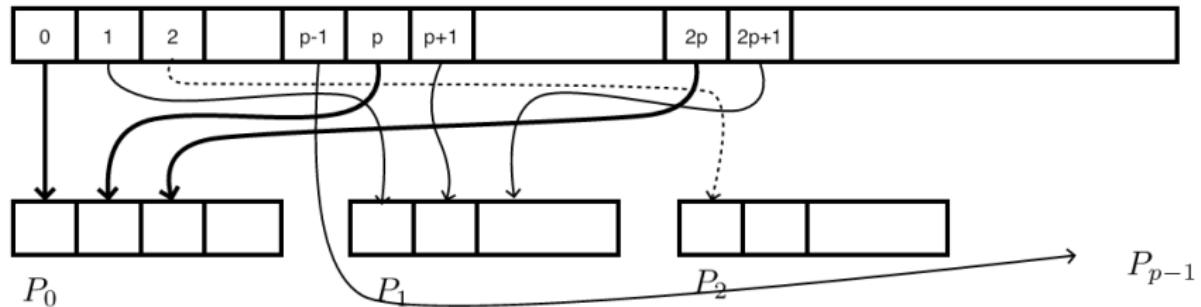
Receiver has no knowledge of the stride of the sender.

# Send vs recv type

```
1 // vector.c
2 source = (double*) malloc(stride*count*sizeof(double));
3 target = (double*) malloc(count*sizeof(double));
4 MPI_Datatype newvectortype;
5 if (procno==sender) {
6     MPI_Type_vector(count,1,stride,MPI_DOUBLE,&newvectortype);
7     MPI_Type_commit(&newvectortype);
8     MPI_Send(source,1,newvectortype,the_other,0,comm);
9     MPI_Type_free(&newvectortype);
10 } else if (procno==receiver) {
11     MPI_Status recv_status;
12     int recv_count;
13     MPI_Recv(target,count,MPI_DOUBLE,the_other,0,comm,
14             &recv_status);
15     MPI_Get_count(&recv_status,MPI_DOUBLE,&recv_count);
16     ASSERT(recv_count==count);
17 }
```



# Illustration of the next exercise



Sending strided data from process zero to all others

## Exercise 29 (stridesend)

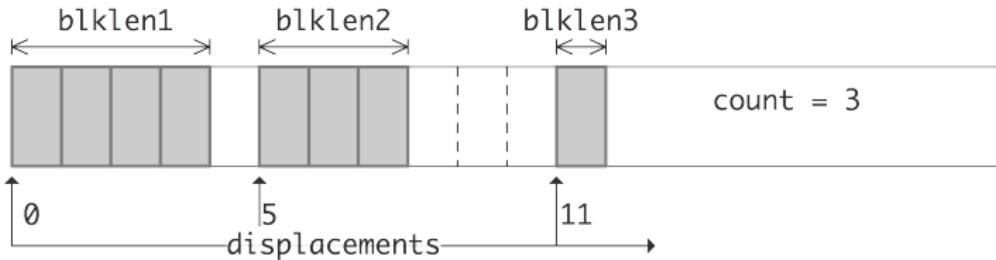
Let processor 0 have an array  $x$  of length  $10P$ , where  $P$  is the number of processors. Elements  $0, P, 2P, \dots, 9P$  should go to processor zero,  $1, P + 1, 2P + 1, \dots$  to processor 1, et cetera.

- Code this as a sequence of send/recv calls, using a vector datatype for the send, and a contiguous buffer for the receive.
- For simplicity, skip the send to/from zero. What is the most elegant solution if you want to include that case?
- For testing, define the array as  $x[i] = i$ .

# Exercise 30

Allocate a matrix on processor zero, using Fortran column-major storage. Using  $P$  sendrecv calls, distribute the rows of this matrix among the processors.

# Indexed type



```
1 int MPI_Type_indexed(
2     int count, int blocklens[], int displacements[],
3     MPI_Datatype old_type, MPI_Datatype *newtype);
```

# Hindexed type

Similar to indexed but using byte offsets:  
explicit memory address.

Example usage scenario: send linked list.

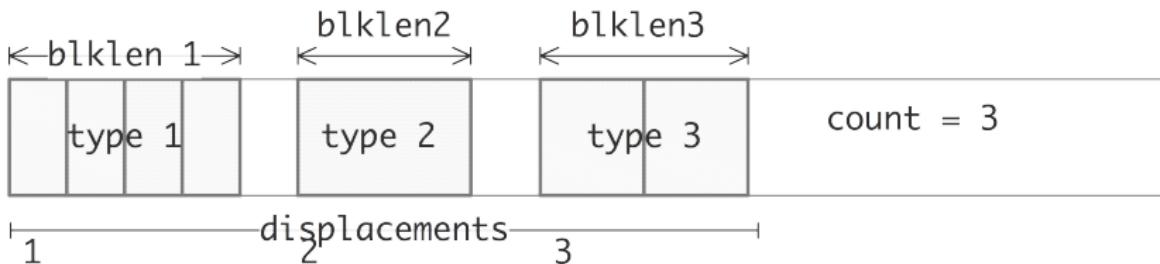
Use `MPI_Get_address` (figure 1)

Figure 1 `MPI_Get_address`

Name	Param name	Explanation	C type	F type
<code>MPI_Get_address</code>				
	<code>location</code>	location in caller memory	<code>const void*</code>	<code>TYPE(*), DIMENSION(.</code>
	<code>address</code>	address of location	<code>MPI_Aint*</code>	<code>INTEGER (KIND=MPI_A</code>
)				

# Heterogeneous: Structure type

```
1 int MPI_Type_create_struct(  
2     int count, int blocklengths[], MPI_Aint displacements[],  
3     MPI_Datatype types[], MPI_Datatype *newtype);
```

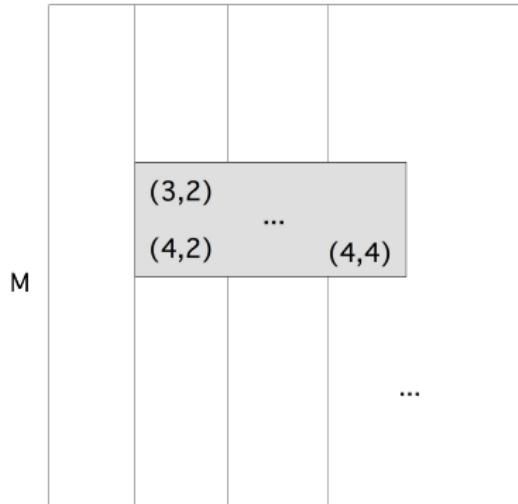


This gets very tedious...

## **Subarray type**

# Submatrix storage

Logical:



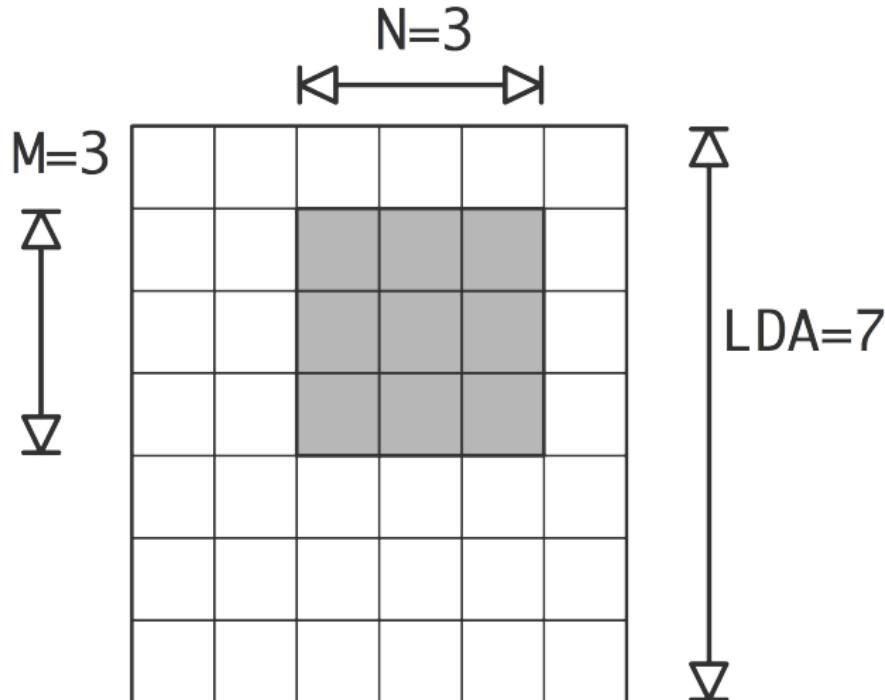
Physical:



- Location of first element
- Stride, blocksize

# BLAS/Lapack storage

Three parameter description:



# Subarray type

- Vector type is convenient for 2D subarrays,
- it gets tedious in higher dimensions.
- Better solution: `MPI_Type_create_subarray`

```
1 MPI_Type_create_subarray(  
2     ndims, array_of_sizes, array_of_subsizes,  
3     array_of_starts, order, oldtype, newtype)
```

Subtle: data does not start at the buffer start

## Exercise 31 (cubegather)

Assume that your number of processors is  $P = Q^3$ , and that each process has an array of identical size. Use `MPI_Type_create_subarray` to gather all data onto a root process. Use a sequence of send and receive calls; `MPI_Gather` does not work here.

If you haven't started `idev` with the right number of processes, use

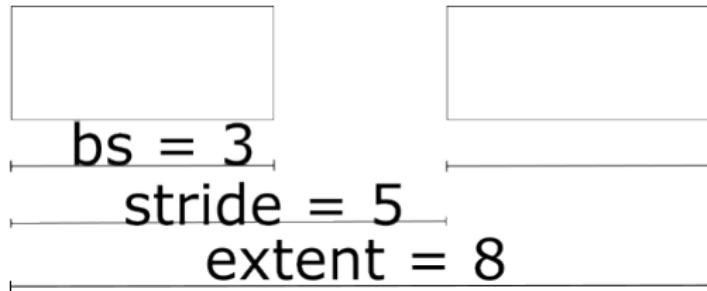
```
ibrun -np 27 cubegather
```

Normally you use `ibrun` without process count argument.

## Extent and resizing

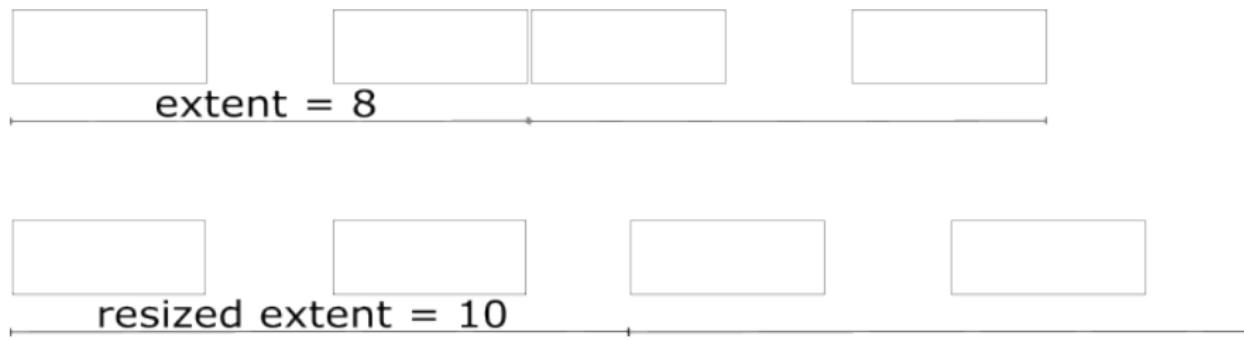
# Extent

Extent: ‘size’ of a type,  
especially useful for derived types.



# Extent resizing: enlarging

Multiple derived types may not be what you intended  
extent resizing makes it artificially larger:



# Extent computation

Use `MPI_Type_get_extent` to query extent  
note: parameters are measured in bytes.

```
1 MPI_Aint lb,asize;
2 MPI_Type_vector(count,bs,stride,MPI_DOUBLE,&newtype);
3 MPI_Type_commit(&newtype);
4 MPI_Type_get_extent(newtype,&lb,&asize);
5 ASSERT( lb==0 );
6 ASSERT( asize==((count-1)*stride+bs)*sizeof(double) );
7 MPI_Type_free(&newtype);
```

# Naive code

Send multiple derived types from

0 1 2 3 4 5 6 7 8 9 10

```
1 // vectorpadsend.c
2 for (int i=0; i<max_elements; i++) sendbuffer[i] = i;
3 MPI_Type_vector(count,blocklength,stride,MPI_INT,&stridetype);
4 MPI_Type_commit(&stridetype);
5 MPI_Send( sendbuffer,ntypes,stridetype, receiver,0, comm );
```

Receive as single block:

```
1 MPI_Recv( recvbuffer,max_elements,MPI_INT, sender,0, comm,&status );
2 int count; MPI_Get_count(&status,MPI_INT,&count);
3 printf("Receive %d elements:",count);
4 for (int i=0; i<count; i++) printf(" %d",recvbuffer[i]);
5 printf("\n");
```

giving an output of:

Receive 6 elements: 0 2 4 5 7 9



# Resizing code

Extend the vector type with padding:

```
1 MPI_Type_get_extent(stridetype,&l,&e);
2 printf("Stride type l=%ld e=%ld\n",l,e);
3 e += ( stride-blocklength ) * sizeof(int);
4 MPI_Type_create_resized(stridetype,l,e,&paddedtype);
5 MPI_Type_get_extent(paddedtype,&l,&e);
6 printf("Padded type l=%ld e=%ld\n",l,e);
7 MPI_Type_commit(&paddedtype);
8 MPI_Send( sendbuffer,ntypes,paddedtype, receiver,0, comm );
```

giving:

Strided type l=0 e=20

Padded type l=0 e=24

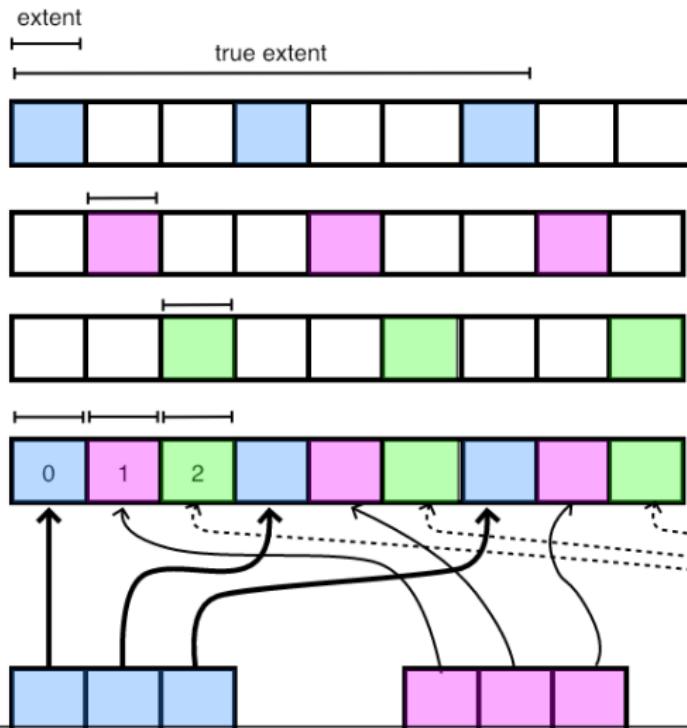
Receive 6 elements: 0 2 4 6 8 10



# MPI\_Type\_create\_resized

# Extent resizing: shrinking

Elements are placed at distance equal to extent:



## Exercise 32 (stridescatter)

Change the stridesend code to use a scatter call, rather than a sequence of sends.

# Extent of subarray type

The 'subarray' type:  
data does not start at the start of the type.

`MPI_Type_get_true_extent` (figure 2)

Figure 2 `MPI_Type_get_true_extent`

Name	Param name	Explanation	C type	F type
<code>MPI_Type_get_true_extent</code>				
	<code>MPI_Type_get_true_extent_c</code>			
	<code>datatype</code>	datatype to get information on	<code>MPI_Datatype</code>	<code>TYPE (MPI_Datatype)</code>
	<code>true_lb</code>	true lower bound of datatype	<code>[ MPI_Aint*</code> <code>  MPI_Count*</code>	<code>INTEGER (KIND=MPI_AI</code>
	<code>true_extent</code>	true extent of datatype	<code>[ MPI_Aint*</code> <code>  MPI_Count*</code>	<code>INTEGER (KIND=MPI_AI</code>
	)			

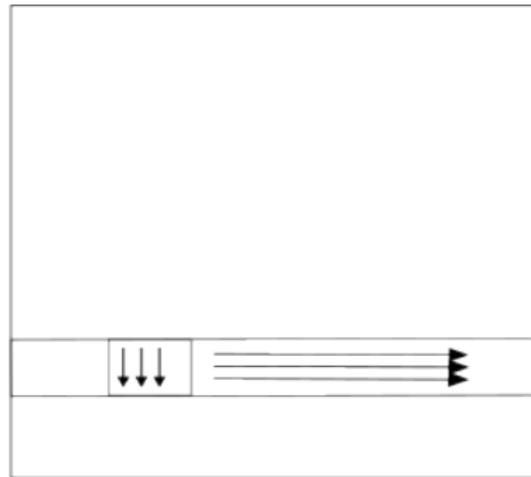
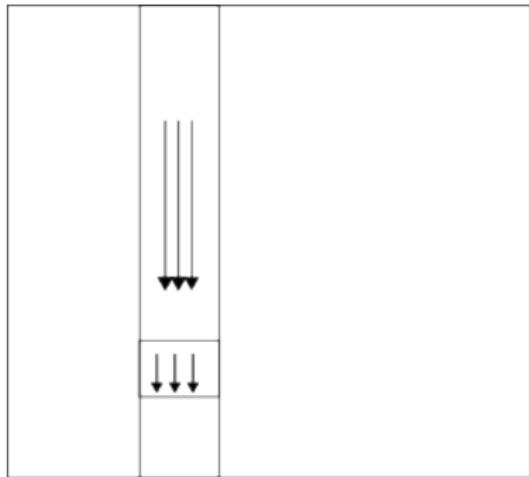
returns non-zero lower bound.



# Transposition

- Basic block of FFT
- Before: Each process stores a block column,
- After: Each process stores a block row

in a picture



# Exercise 33 (transposeblock)

Fill in the missing bits of the skeleton code.

## Packed data

# Packing into buffer

- Construct a buffer with `MPI_Pack`
- Send with `MPI_Send` and such
- Receive
- Unpack buffer elements with `MPI_Unpack`

# MPI\_Pack

# MPI\_Unpack

# Example

```
1 if (procno==sender) {  
2     position = 0;  
3     MPI_Pack(&nseeds,1,MPI_INT,  
4                 buffer,buflen,&position,comm);  
5     for (int i=0; i<nseeds; i++) {  
6         double value = rand()/(double)RAND_MAX;  
7         printf("[%d] pack %e\n",procno,value);  
8         MPI_Pack(&value,1,MPI_DOUBLE,  
9                   buffer,buflen,&position,comm);  
10    }  
11    MPI_Pack(&nseeds,1,MPI_INT,  
12                 buffer,buflen,&position,comm);  
13    MPI_Send(buffer,position,MPI_PACKED,other,0,comm);  
14 } else if (procno==receiver) {  
15     int irecv_value;  
16     double xrecv_value;  
17     MPI_Recv(buffer,buflen,MPI_PACKED,other,0,  
18                 comm,MPI_STATUS_IGNORE);  
19     position = 0;  
20     MPI_Unpack(buffer,buflen,&position,  
21                 &nseeds,1,MPI_INT,comm);  
22     for (int i=0; i<nseeds; i++) {  
23         MPI_Unpack(buffer,buflen,  
24                     &position,&xrecv_value,1,MPI_DOUBLE,comm);  
25         printf("[%d] unpack %e\n",procno,xrecv_value);  
26     }  
27     MPI_Unpack(buffer,buflen,&position,  
28                     &irecv_value,1,MPI_INT,comm);  
29     if (irecv_value==nseeds);  
30 }
```



# Communicator manipulations

# Overview

In this section you will learn about various subcommunicators.

Commands learned:

- *MPI\_Comm\_dup*, discussion of library design
- *MPI\_Comm\_split*
- discussion of groups
- discussion of inter/intra communicators.

# Sub-computations

Simultaneous groups of processes, doing different tasks, but loosely interacting:

- Simulation pipeline: produce input data, run simulation, post-process.
- Climate model: separate groups for air, ocean, land, ice.
- Quicksort: split data in two, run quicksort independently on the halves.
- Process grid: do broadcast in each column.

New communicators are formed recursively from `MPI_COMM_WORLD`.



# Communicator duplication

Simplest new communicator: identical to a previous one.

```
1 int MPI_Comm_dup(MPI_Comm comm, MPI_Comm *newcomm)
```

This is useful for library writers:

```
1 MPI_Isend(...); MPI_Irecv(...);
2 // library call
3 MPI_Waitall(...);
```

- Naively, the library can ‘catch’ the user messages.
- With a duplicate communicator there is no confusion:  
user and library both have their own ‘context’ for their messages.

# Interleaved library and user code

```
1 library my_library(comm);
2 MPI_Isend(&sdata,1,MPI_INT,other,1,comm,&(request[0]));
3 my_library.communication_start();
4 MPI_Irecv(&rdata,1,MPI_INT,other,MPI_ANY_TAG,
5           comm,&(request[1]));
6 MPI_Waitall(2,request,status);
7 my_library.communication_end();
```

# Library internally has messages

```
1 int library::communication_start() {
2     int sdata=6,rdata;
3     MPI_Isend(&sdata,1,MPI_INT,other,2,comm,&(request[0]));
4     MPI_Irecv(&rdata,1,MPI_INT,other,MPI_ANY_TAG,
5                comm,&(request[1]));
6     return 0;
7 }
8
9 int library::communication_end() {
10    MPI_Status status[2];
11    MPI_Waitall(2,request,status);
12    return 0;
13 }
```

# Wrong way of setting up the library

```
1 // commdupwrong.cxx
2 class library {
3 private:
4     MPI_Comm comm;
5     int procno,nprocs,other;
6     MPI_Request request[2];
7 public:
8     library(MPI_Comm incomm) {
9         comm = incomm;
10        MPI_Comm_rank(comm,&procno);
11        other = 1-procno;
12    };
13    int communication_start();
14    int communication_end();
15 };
```

# Right way of setting up the library

```
1 // commdupright.cxx
2 class library {
3 private:
4     MPI_Comm comm;
5     int procno,nprocs,other;
6     MPI_Request request[2];
7 public:
8     library(MPI_Comm incomm) {
9         MPI_Comm_dup(incomm,&comm);
10    MPI_Comm_rank(comm,&procno);
11    other = 1-procno;
12 };
13 ~library() {
14     MPI_Comm_free(&comm);
15 }
16 int communication_start();
17 int communication_end();
18 };
```



# Disjoint splitting

Split a communicator in multiple disjoint others.

Give each process a ‘color’, group processes by color:

```
1 int MPI_Comm_split(MPI_Comm comm, int color, int key,  
2 MPI_Comm *newcomm)
```

(key determines ordering: use rank unless you want special effects)

# MPI\_Comm\_split

Name	Param name	Explanation	C type	F type
MPI_Comm_split (				
comm	communicator		MPI_Comm	TYPE(MPI_Comm)
color	control of subset assignment		int	INTEGER
key	control of rank assignment		int	INTEGER
newcomm	new communicator		MPI_Comm*	TYPE(MPI_Comm)
)				

# Row/column example

Simulate a process grid

create subcommunicator per column (or row)

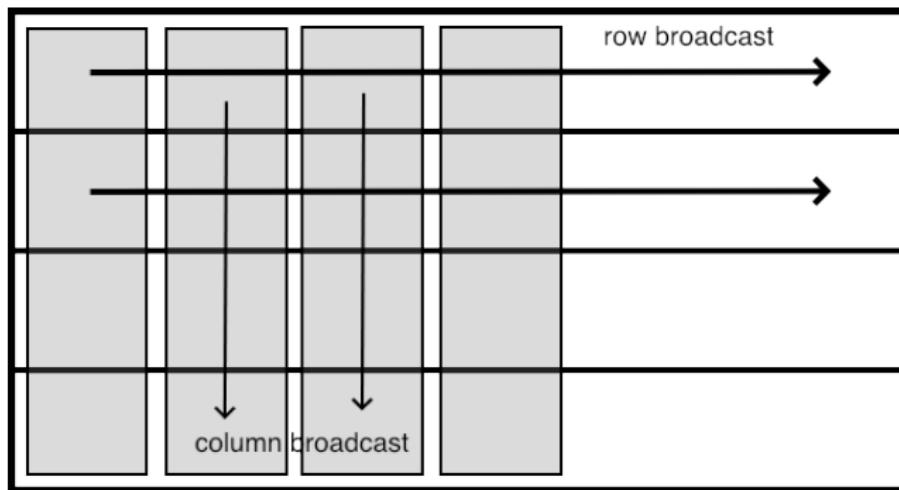
assume processes numbered by rows

```
1 MPI_Comm_rank( MPI_COMM_WORLD, &procno );
2 proc_i = procno % proc_column_length;
3 proc_j = procno / proc_column_length;
4
5 MPI_Comm column_comm;
6 MPI_Comm_split( MPI_COMM_WORLD, proc_j, procno, &column_comm );
7
8 MPI_Bcast( data, ... column_comm );
```

Food for thought: there are many columns, but only one `column_comm` variable. Why?



# Row and column communicators



Row and column broadcasts in subcommunicators

## Exercise 34 (procgrid)

Organize your processes in a grid, and make subcommunicators for the rows and columns. For this compute the row and column number of each process.

In the row and column communicator, compute the rank. For instance, on a  $2 \times 3$  processor grid you should find:

Global ranks:			Ranks in row:			Ranks in column:		
0	1	2	0	1	2	0	0	0
3	4	5	0	1	2	1	1	1

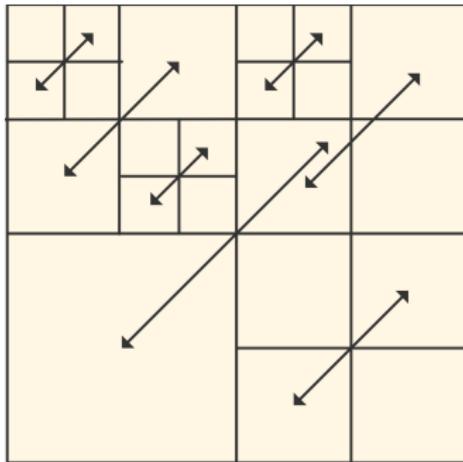
Check that the rank in the row communicator is the column number, and the other way around.

Run your code on different number of processes, for instance a number of rows and columns that is a power of 2, or that is a prime number. This is one occasion where you could use `ibrun -np 9`; normally you would *never* put a processor count on `ibrun`.



# Exercise 35

Implement a recursive algorithm for matrix transposition:



- Swap blocks (1, 2) and (2, 1); then
- Divide the processors into four subcommunicators, and apply this algorithm recursively on each;
- If the communicator has only one process, transpose the matrix in place.

# Splitting by shared memory

- `MPI_Comm_split_type` splits into communicators of same type.
- MPI-3: only `MPI_COMM_TYPE_SHARED` splitting by shared memory.
- MPI-4: `MPI_COMM_TYPE_HW_GUIDED` split using an `info` value from `MPI_Get_hw_resource_types`.

```
1 // commsplittype.c
2 MPI_Info info;
3 MPI_Comm_split_type
4     (MPI_COMM_WORLD,
5      MPI_COMM_TYPE_SHARED,
6      procno,info,&sharedcomm);
7 MPI_Comm_size
8     (sharedcomm,&new_nprocs);
9 MPI_Comm_rank
10    (sharedcomm,&new_procno);
```

# Inter-communicators

- Communicators so far are of *intra-communicator* type.
- Bridge between two communicators: *inter-communicator*.
- Example: communicator with newly spawned processes

# In a picture



Illustration of ranks in an inter-communicator setup

```
1 // intercomm.c
2 MPI_Comm intercomm;
3 MPI_Intercomm_create
4     /* local_comm:          */ split_half_comm,
5     /* local_leader:        */ local_leader_in_inter_comm,
6     /* peer_comm:           */ MPI_COMM_WORLD,
7     /* remote_peer_rank: */ global_rank_of_other_leader,
8     /* tag:                 */ inter_tag,
9     /* newintercomm:        */ &intercomm );
```

# Concepts

- Two local communicators
- The ‘peer’ communicator that contains them
- Leaders in each of them
- An inter-communicator over the leaders.

# Routines

- *MPI\_Intercomm\_create*: create
- *MPI\_Comm\_get\_parent*: the other leader (see process management)
- *MPI\_Comm\_remote\_size*, *MPI\_Comm\_remote\_group*: query the other communicator
- *MPI\_Comm\_test\_inter*: is this an inter or intra?

# More

- Non-disjoint subcommunicators through process groups.
- Process topologies: cartesian and graph.  
There will also be a section about this, later.

# Cartesian topologies

# Cartesian decomposition

Code:

```
1 // cartdims.c
2 int *dimensions = (int*) malloc(dim*
3     sizeof(int));
4 for (int idim=0; idim<dim; idim++)
5     dimensions[idim] = 0;
6 MPI_Dims_create(nprocs, dim, dimensions);
```

Output:

```
1 mpicc -o cartdims
         ↪cartdims.o
2 Cartesian grid size: 3
         ↪dim: 1
3   3
4 Cartesian grid size: 3
         ↪dim: 2
5   3 x 1
6 Cartesian grid size: 4
         ↪dim: 1
7   4
8 Cartesian grid size: 4
         ↪dim: 2
9   2 x 2
10  Cartesian grid size: 4
         ↪dim: 3
11  2 x 2 x 1
12  Cartesian grid size: 12
         ↪dim: 1
13  12
14  Cartesian grid size: 12
         ↪dim: 2
15  4 x 3
16  Cartesian grid size: 12
```



# Create/test Cartesian topology

```
1 MPI_Comm cart_comm;
2 int *periods = (int*) malloc(dim*sizeof(int));
3 for ( int id=0; id<dim; id++ ) periods[id] = 0;
4 MPI_Cart_create
5   ( comm,dim,dimensions,periods,
6     0,&cart_comm );

1 int dim;
2 MPI_Cartdim_get( cart_comm,&dim );
3 int *dimensions = (int*) malloc(dim*sizeof(int));
4 int *periods    = (int*) malloc(dim*sizeof(int));
5 int *coords     = (int*) malloc(dim*sizeof(int));
6 MPI_Cart_get( cart_comm,dim,dimensions,periods,coords );
```

# Rank translation

```
1 // cart.c
2 MPI_Comm comm2d;
3 int periodic[ndim]; periodic[0] = periodic[1] = 0;
4 MPI_Cart_create(comm,ndim,dimensions,periodic,1,&comm2d);
5 if (comm2d==MPI_COMM_NULL) {
6   printf("Process %d not included\n",procno);
7 } else {
8   MPI_Cart_coords(comm2d,procno,ndim,coord_2d);
9   MPI_Cart_rank(comm2d,coord_2d,&rank_2d);
10  printf("I am %d: (%d,%d); originally %d\n",
11        rank_2d,coord_2d[0],coord_2d[1],procno);
```

# Cartesian communication

```
1 // cartcoord.c
2 for ( int id=0; id<dim; id++)
3     periods[id] = id==0 ? 1 : 0;
4 MPI_Cart_create
5   ( comm,dim,dimensions,periods,
6     0,&period_comm );
```

Code:

```
1 int pred,succ;
2 MPI_Cart_shift
3   (period_comm,/* dim: */ 0,/* up: */ 1,
4     &pred,&succ);
5 printf
6   ("periodic dimension 0:\n    src=%d, tgt
     =%d\n",
7   pred,succ);
8 MPI_Cart_shift
9   (period_comm,/* dim: */ 1,/* up: */ 1,
10    &pred,&succ);
11 printf
12   ("non-periodic dimension 1:\n    src=%d,
      tgt=%d\n",
13   pred,succ);
```

Output:

```
1 Grid of size 6 in 3
   ↪dimensions:
2   3 x 2 x 1
3 Shifting process 0.
4 periodic dimension 0:
5   src=4, tgt=2
6 non-periodic dimension
   ↪1:
7   src=-1, tgt=1
```

# Subgrids

Code:

```
1 MPI_Cart_sub( period_comm,remain,&
    hyperplane );
2 if (procno==0) {
3     MPI_Topo_test( hyperplane,&topo_type );
4     MPI_Cartdim_get( hyperplane,&hyperdim )
        ;
5     printf("hyperplane has dimension %d,
        type %d\n",
        hyperdim,topo_type);
6     MPI_Cart_get( hyperplane,dim,dims,
        period,coords );
7     printf(" periodic: ");
8     for (int id=0; id<2; id++)
9         printf("%d,",period[id]);
10    printf("\n");
```

Output:

```
1 Grid of size 6 in 3
    ↪dimensions:
2     3 x 2 x 1
3 hyperplane has
    ↪dimension 2,
    ↪type 2
4     periodic: 1,0,
```



# Exercise 36 (isendirecvcart)

Use Cartesian topology routines to extend exercise 26 to two dimensions.

# MPI File I/O

# Overview

This section discusses parallel I/O. What is the problem with regular I/O in parallel?

Commands learned:

- `MPI_File_open/write/close` and variants
- parallel file pointer routines: `MPI_File_set_view/write_at`

# The trouble with parallel I/O

- Multiple process reads from one file: no problem.
- Multiple writes to one file: big problem.
- Everyone writes to separate file: stress on the file system, and requires post-processing.

# MPI I/O

- Part of MPI since MPI-2
- Joint creation of one file from bunch of processes.
- You could also use hdf5, netcdf, silo ...

# The usual bits

```
1 MPI_File mpifile;
2 MPI_File_open(comm,"blockwrite.dat",
3                 MPI_MODE_CREATE | MPI_MODE_WRONLY,MPI_INFO_NULL,
4                 &mpifile);
5 if (procno==0) {
6     MPI_File_write
7         (mpifile,output_data,nwords,MPI_INT,MPI_STATUS_IGNORE);
8 }
9 MPI_File_close(&mpifile);

1 type(MPI_File) :: mpifile ! F08
2 integer          :: mpifile ! F90
```



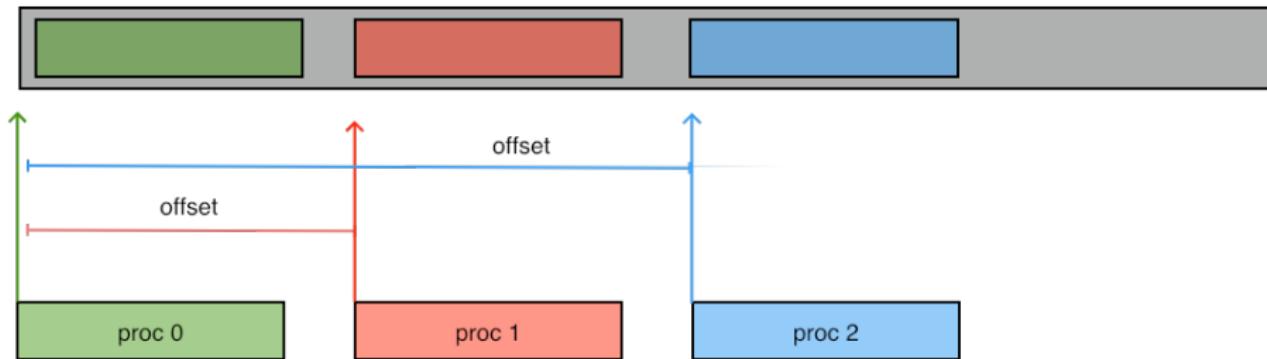
# How do you make it unique for a process?

```
1 MPI_File_write_at  
2     (mpifile,offset,output_data,nwords,  
3      MPI_INT,MPI_STATUS_IGNORE);
```

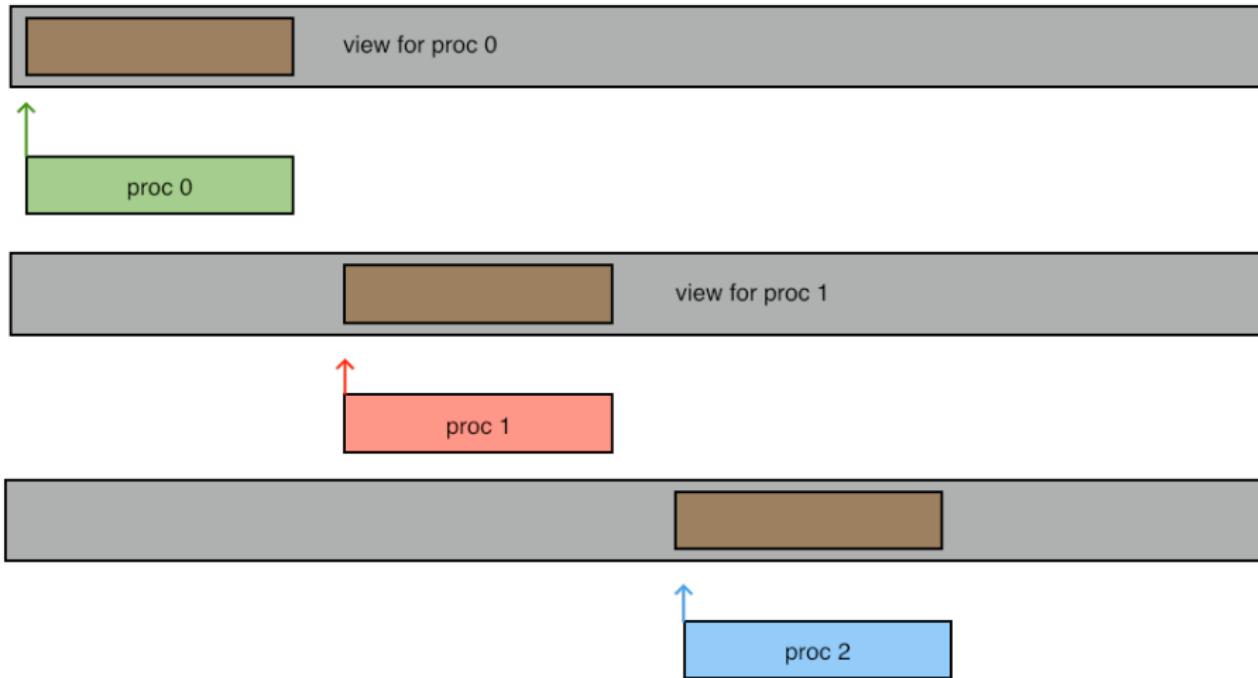
or

```
1 MPI_File_set_view  
2     (mpifile,  
3      offset,datatype,  
4      MPI_INT,"native",MPI_INFO_NULL);  
5 MPI_File_write // no offset, we have a view  
6     (mpifile,output_data,nwords,MPI_INT,MPI_STATUS_IGNORE);
```

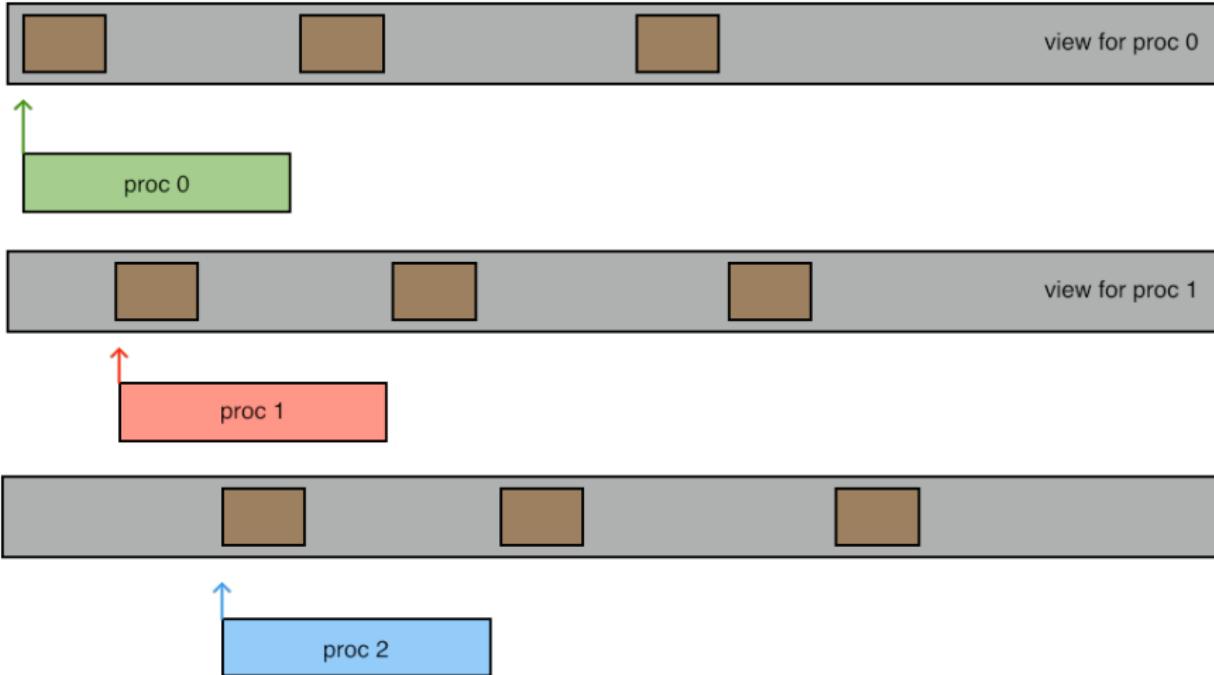
# Write at an offset



# Write to a view



# Write to a view



Made

## Exercise 37 (blockwrite)

The given code works for one writing process. Compute a unique offset for each process (in bytes!) so that all the local arrays are placed in the output file in sequence.

# Exercise 38 (viewwrite)

Solve the previous exercise by using `MPI_File_write` (that is, without offset), but by using `MPI_File_set_view` to specify the location.

## Exercise 39 (scatterwrite)

Now write the local arrays cyclically to the file: with 5 processes and 3 elements per process the file should contain

```
1 4 7 10 13 | 2 5 8 11 14 | 3 6 9 12 15
```

Do this by defining a vector derived type and setting that as the file view.

# One-sided communication

# Overview

This section concerns one-sided operations, which allows 'shared memory' type programming. (Actual shared memory later.)

Commands learned:

- *MPI\_Put, MPI\_Get, MPI\_Accumulate*
- Window commands: *MPI\_Win\_create, MPI\_Win\_allocate*
- Active target synchronization *MPI\_Win\_fence*
- *MPI\_Win\_post/wait/start/complete*
- Passive target synchronization *MPI\_Win\_lock / MPI\_Win\_unlock*
- Atomic operations: *MPI\_Fetch\_and\_op*

## **Basic mechanisms**

# Motivation

With two-sided messaging, you can not just put data on a different process: the other has to expect it and receive it.

- Sparse matrix: it is easy to know what you are receiving, not what you need to send. Usually solved with complicated preprocessing step.
- Neuron simulation: spiking neuron propagates information to neighbors. Uncertain when this happens.
- Other irregular data structures: distributed hash tables.

# Dynamic data

```
1 x = f();  
2 p = hash(x);  
3 MPI_Send( x, /* to: */ p );
```

Problem: how does  $p$  know to post a receive,  
and how does everyone else know not to?

# One-sided concepts

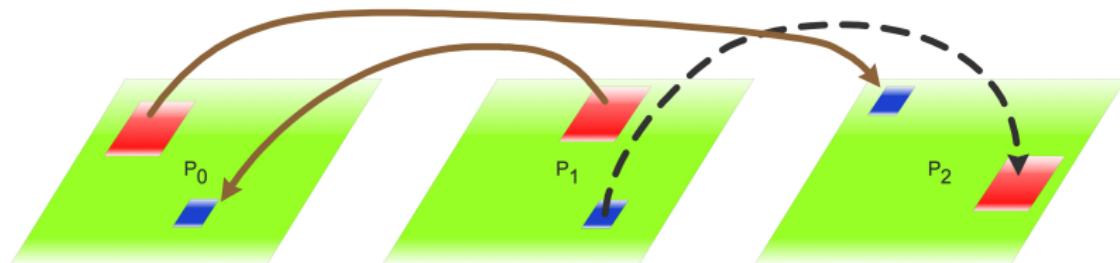


Window

Data

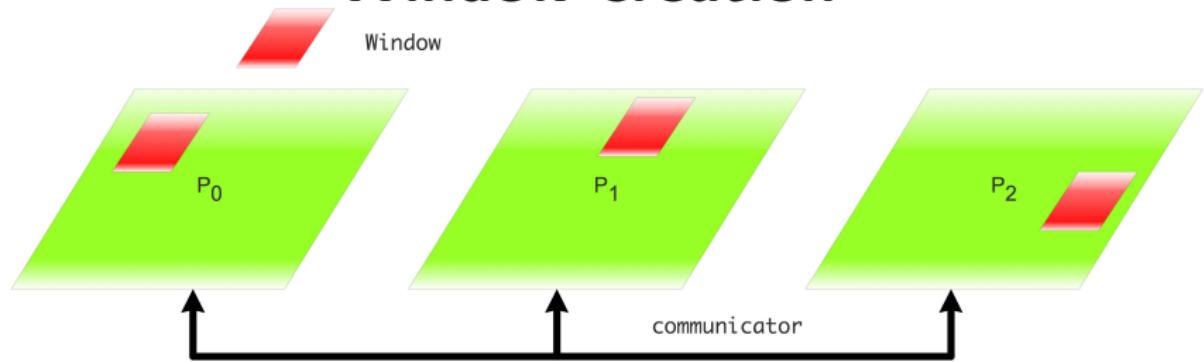
→ Get

→ Put



- A process has a *window* that other processes can access.
- *origin*: process doing a one-sided call  
*target*: process being accessed.
- One-sided calls: *MPI\_Put*, *MPI\_Get*, *MPI\_Accumulate*.
- Various synchronization mechanisms.

# Window creation



```
1 MPI_Win_create (void *base, MPI_Aint size,  
2   int disp_unit, MPI_Info info, MPI_Comm comm, MPI_Win *win)
```

- *size*: in bytes
- *disp\_unit*: `sizeof(type)`

Also call `MPI_Win_free` when done. This is important!

# Window allocation

Instead of passing buffer, let MPI allocate with `MPI_Win_allocate` and return the buffer pointer:

```
1 int MPI_Win_allocate
2   (MPI_Aint size, int disp_unit, MPI_Info info,
3    MPI_Comm comm, void *baseptr, MPI_Win *win)
```

can use dedicated fast memory.

# Active target synchronization

All processes call `MPI_Win_fence`. Epoch is between fences:

```
1 MPI_Win_fence(MPI_MODE_NOPRECEDE, win);
2 if (procno==producer)
3   MPI_Put( /* operands */, win);
4 MPI_Win_fence(MPI_MODE_NOSUCCEED, win);
```

Second fence indicates that one-sided communication is concluded:  
target knows that data has been put.

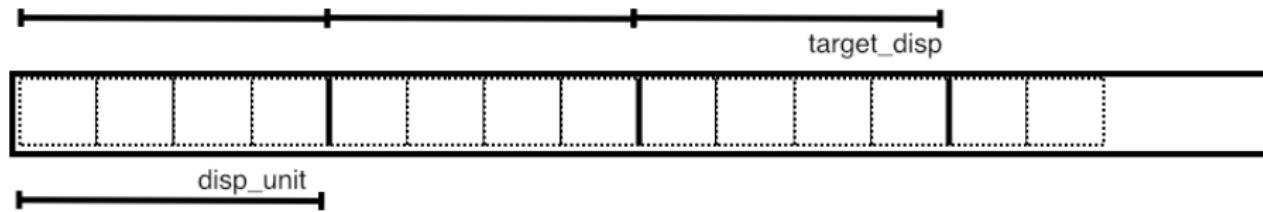
# MPI\_Put

Name	Param name	Explanation	C type	F type
MPI_Put (				
MPI_Put_c (				
origin_addr		initial address of origin buffer	const void*	TYPE(*), DIMENSION(..)
origin_count		number of entries in origin buffer	[ int MPI_Count	INTEGER
origin_datatype		datatype of each entry in origin buffer	MPI_Datatype	TYPE(MPI_Datatype)
target_rank		rank of target	int	INTEGER
target_disp		displacement from start of window to target buffer	MPI_Aint	INTEGER (KIND=MPI_ADDRESS_KIND)
target_count		number of entries in target buffer	[ int MPI_Count	INTEGER
target_datatype		datatype of each entry in target buffer	MPI_Datatype	TYPE(MPI_Datatype)
win		window object used for communication	MPI_Win	TYPE(MPI_Win)
)				

# Location in the window

Location to write:

$$\text{window\_base} + \text{target\_disp} \times \text{disp\_unit}.$$



# Exercise 40 (rightput)

Revisit exercise 19 and solve it using `MPI_Put`.

# Exercise 41 (randomput)

Write code where:

- process 0 computes a random number  $r$
- if  $r < .5$ , zero writes in the window on 1;
- if  $r \geq .5$ , zero writes in the window on 2.

# Exercise (optional) 42 (randomput)

Replace `MPI_Win_create` by `MPI_Win_allocate`.

# Remaining simple routines: Get, Accumulate

- $\text{MPI\_Get}$  is converse of  $\text{MPI\_Put}$ . Like Recv, but no status argument.
- $\text{MPI\_Accumulate}$  is a Put plus a reduction on the result: multiple accumulate calls in one epoch well-defined.  
Can use any predefined  $\text{MPI\_Op}$  (not user-defined) or  $\text{MPI\_REPLACE}$ .

# MPI\_Get

Name	Param name	Explanation	C type	F type
MPI_Get (				
MPI_Get_c (				
origin_addr		initial address of origin buffer	void*	TYPE(*), DIMENSION(..)
origin_count		number of entries in origin buffer	[ int MPI_Count	INTEGER
origin_datatype		datatype of each entry in origin buffer	MPI_Datatype	TYPE(MPI_Datatype)
target_rank		rank of target	int	INTEGER
target_disp		displacement from window start to the beginning of the target buffer	MPI_Aint	INTEGER (KIND=MPI_ADDRESS_KIND)
target_count		number of entries in target buffer	[ int MPI_Count	INTEGER
target_datatype		datatype of each entry in target buffer	MPI_Datatype	TYPE(MPI_Datatype)
win		window object used for communication	MPI_Win	TYPE(MPI_Win)
)				

# MPI\_Accumulate

Name	Param name	Explanation	C type	F type
MPI_Accumulate (				
MPI_Accumulate_c (				
origin_addr		initial address of buffer	const void*	TYPE(*), DIMENSION(..)
origin_count		number of entries in buffer	[ int MPI_Count	INTEGER
origin_datatype		datatype of each entry	MPI_Datatype	TYPE(MPI_Datatype)
target_rank		rank of target	int	INTEGER
target_disp		displacement from start of window to beginning of target buffer	MPI_Aint	INTEGER (KIND=MPI_ADDRESS_KIND)
target_count		number of entries in target buffer	[ int MPI_Count	INTEGER
target_datatype		datatype of each entry in target buffer	MPI_Datatype	TYPE(MPI_Datatype)
op		reduce operation	MPI_Op	TYPE(MPI_Op)
win		window object	MPI_Win	TYPE(MPI_Win)
)				

# **Ordering and synchronization**

# Fence synchronization

Already mentioned active target synchronization:  
the target indicates the start/end of an epoch.

Simplest mechanism: *MPI\_Win\_fence*, collective.

After the closing fence, buffers have been sent / windows have been updated.

# Ordering of operations

Ordering is often undefined:

- No ordering of Get and Put/Accumulate operations
- No ordering of multiple Puts. Use Accumulate.

The following operations are well-defined inside one epoch:

- Instead of multiple Put operations, use Accumulate with `MPI_REPLACE`.
- `MPI_Get_accumulate` with `MPI_NO_OP` is safe.
- Multiple Accumulate operations from one origin are ordered by default.

# Exercise (optional) 43 (countdown)

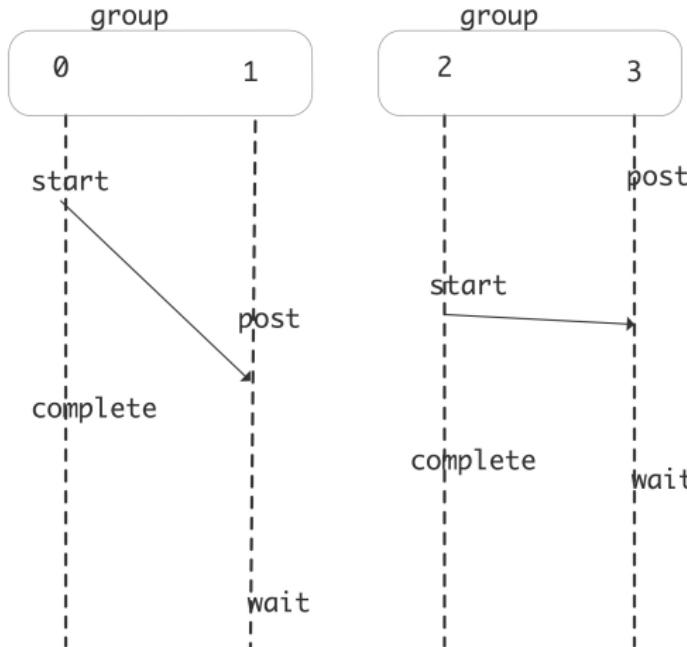
Implement a shared counter:

- One process maintains a counter;
- Iterate: all others at random moments update this counter.
- When the counter is no longer positive, everyone stops iterating.

The problem here is data synchronization: does everyone see the counter the same way?

# A second active synchronization

Use `MPI_Win_post`, `MPI_Win_wait`, `MPI_Win_start`, `MPI_Win_complete` calls



More fine grained than fences.

## **Passive target synchronization**

# Passive target synchronization

Lock a window on the target:

```
1 MPI_Win_lock  
2     (int locktype, int rank, int assert, MPI_Win win)  
3 MPI_Win_unlock  
4     (int rank, MPI_Win win)
```

with types: *MPI\_LOCK\_SHARED MPI\_LOCK\_EXCLUSIVE*

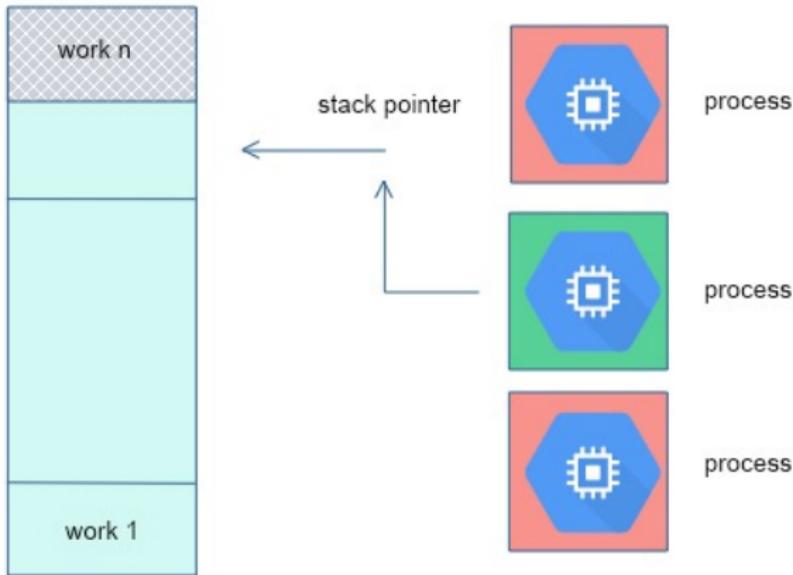
# Justification

MPI-1/2 lacked tools for race condition-free one-sided communication.  
These have been added in MPI-3.

# Emulating shared memory with one-sided communication

- One process stores a table of work descriptors, and a ‘stack pointer’ stating how many there are.
- Each process reads the pointer, reads the corresponding descriptor, and decrements the pointer; and
- A process that has read a descriptor then executes the corresponding task.
- Non-collective behavior: processes only take a descriptor when they are available.

# In a picture



# Simplified model

- One process has a counter, which models the shared memory;
- Each process, if available, reads the counter; and
- ... decrements the counter.
- No actual work: random decision if process is available.

# Shared memory problems: what is a race condition?

Race condition: outward behavior depends on timing/synchronization of low-level events.  
In shared memory associated with shared data.

Example:

```
Init: I=0
process 1: I=I+2
process 2: I=I+3
```

scenario 1.	scenario 2.	scenario 3.
$I = 0$		
read I = 0 local I = 2 write I = 2	read I = 0 local I = 3 write I = 3	read I = 0 local I = 2 write I = 2
		read I = 2 local I = 5 write I = 5
$I = 3$	$I = 2$	$I = 5$

(In MPI, the read/write would be *MPI\_Get / MPI\_Put* calls)



# Case study in shared memory: 1, wrong

```
1 // countdownput.c
2 MPI_Win_fence(0,the_window);
3 int counter_value;
4 MPI_Get( &counter_value,1,MPI_INT,
5         counter_process,0,1,MPI_INT,
6         the_window);
7 MPI_Win_fence(0,the_window);
8 if (i_am_available) {
9     int decrement = -1;
10    counter_value += decrement;
11    MPI_Put
12        ( &counter_value,    1,MPI_INT,
13          counter_process,0,1,MPI_INT,
14          the_window);
15 }
16 MPI_Win_fence(0,the_window);
```

# Discussion

- The multiple `MPI_Put` calls conflict.
- Code is correct if in each iteration there is only one writer.
- Question: In that case, can we take out the middle fence?
- Question: what is wrong with

```
1 MPI_Win_fence(0,the_window);
2 if (i_am_available) {
3   MPI_Get( &counter_value, ... )
4   MPI_Win_fence(0,the_window);
5   MPI_Put( ... )
6 }
7 MPI_Win_fence(0,the_window);
```

?



# Case study in shared memory: 2, hm

```
1 // countdownacc.c
2 MPI_Win_fence(0,the_window);
3 int counter_value;
4 MPI_Get( &counter_value,1,MPI_INT,
5         counter_process,0,1,MPI_INT,
6         the_window);
7 MPI_Win_fence(0,the_window);
8 if (i_am_available) {
9     int decrement = -1;
10    MPI_Accumulate
11        ( &decrement,           1,MPI_INT,
12          counter_process,0,1,MPI_INT,
13          MPI_SUM,
14          the_window);
15 }
16 MPI_Win_fence(0,the_window);
```

# Discussion: need for atomics

- `MPI_Accumulate` is atomic, so no conflicting writes.
- What is the problem?
- Answer: Processes are not reading unique `counter_value` values.
- Conclusion: Read and update need to come together:  
read unique value and immediately update.

Atomic 'get-and-set-with-no-one-coming-in-between':

`MPI_Fetch_and_op` / `MPI_Get_accumulate`.

Former is simple version: scalar only.

# MPI\_Fetch\_and\_op

Name	Param name	Explanation	C type	F type
<code>MPI_Fetch_and_op (</code>				
	<code>origin_addr</code>	initial address of buffer	<code>const void*</code>	<code>TYPE(*), DIMENSION(..)</code>
	<code>result_addr</code>	initial address of result buffer	<code>void*</code>	<code>TYPE(*), DIMENSION(..)</code>
	<code>datatype</code>	datatype of the entry in origin, result, and target buffers	<code>MPI_Datatype</code>	<code>TYPE(MPI_Datatype)</code>
	<code>target_rank</code>	rank of target	<code>int</code>	<code>INTEGER</code>
	<code>target_disp</code>	displacement from start of window to beginning of target buffer	<code>MPI_Aint</code>	<code>INTEGER (KIND=MPI_ADDRESS_KIND)</code>
	<code>op</code>	reduce operation	<code>MPI_Op</code>	<code>TYPE(MPI_Op)</code>
	<code>win</code>	window object	<code>MPI_Win</code>	<code>TYPE(MPI_Win)</code>
	<code>)</code>			

# Case study in shared memory: 3, good

```
1 MPI_Win_fence(0,the_window);
2 int
3     counter_value;
4 if (i_am_available) {
5     int
6         decrement = -1;
7     total_decrement++;
8     MPI_Fetch_and_op
9     ( /* operate with data from origin: */ &decrement,
10      /* retrieve data from target: */ &counter_value,
11      MPI_INT, counter_process, 0, MPI_SUM,
12      the_window);
13 }
14 MPI_Win_fence(0,the_window);
15 if (i_am_available) {
16     my_counter_values[n_my_counter_values++] = counter_value;
17 }
```



# Allowable operators. (Hint!)

MPI type	meaning	applies to
MPI.Op		
<i>MPI_MAX</i>	MPI.MAX	maximum
<i>MPI_MIN</i>	MPI.MIN	minimum
<i>MPI_SUM</i>	MPI.SUM	sum
<i>MPI_PROD</i>	MPI.PROC	product
<i>MPI_REPLACE</i>	MPI.REPLACE	overwrite
<i>MPI_NO_OP</i>	MPI.OP'NULL	no change
<i>MPI_BAND</i>	MPI.LAND	logical and
<i>MPI_BOR</i>	MPI.LOR	logical or
<i>MPI_BXOR</i>	MPI.LXOR	logical xor
<i>MPI_BAND</i>	MPI.BAND	bitwise and
<i>MPI_BOR</i>	MPI.BOR	bitwise or
<i>MPI_BXOR</i>	MPI.BXOR	bitwise xor
<i>MPI_MAXLOC</i>	MPI.MAXLOC	max value and location
<i>MPI_MINLOC</i>	MPI.MINLOC	min value and location
		<i>MPI_DOUBLE_INT</i> and such

# Problem

We are using fences, which are collective.

What if a process is still operating on its local work?

Better (but more tricky) solution:  
use passive target synchronization and locks.

# Passive target epoch

```
1 if (rank == 0) {  
2   MPI_Win_lock (MPI_LOCK_EXCLUSIVE, 1, 0, win);  
3   MPI_Put (outbuf, n, MPI_INT, 1, 0, n, MPI_INT, win);  
4   MPI_Win_unlock (1, win);  
5 }
```

No action on the target required!

## Exercise 44 (lockfetch)

Investigate atomic updates using passive target synchronization. Use `MPI_Win_lock` with an exclusive lock, which means that each process only acquires the lock when it absolutely has to.

- All processes but one update a window:

```
1 int one=1;
2 MPI_Fetch_and_op(&one, &readout,
3      MPI_INT, repo, zero_disp, MPI_SUM,
4      the_win);
```

- while the remaining process spins until the others have performed their update.

Use an atomic operation for the latter process to read out the shared value.

Can you replace the exclusive lock with a shared one?



## Exercise 45 (lockfetchshared)

As exercise 44, but now use a shared lock: all processes acquire the lock simultaneously and keep it as long as is needed.

The problem here is that coherence between window buffers and local variables is now not forced by a fence or releasing a lock. Use `MPI_Win_flush_local` to force coherence of a window (on another process) and the local variable from `MPI_Fetch_and_op`.

# Big data communication

# Overview

This section discusses big messages.

Commands learned:

- *MPI\_Send\_c, MPI\_Allreduce\_c, MPI\_Get\_count\_c* etc. (MPI-4)

# The problem with large messages

- There is no problem allocating large buffers:

```
1 size_t bigsize = 1<<33;  
2 double *buffer =  
3     (double*) malloc(bigsize*sizeof(double));
```

- But you can not tell MPI how big the buffer is:

```
1 MPI_Send(buffer,bigsize,MPI_DOUBLE,...) // WRONG
```

because the size argument has to be `int`.

# MPI 3 count type

Count type since MPI 3

C:

```
1 MPI_Count count;
```

Fortran:

```
1 Integer(kind=MPI_COUNT_KIND) :: count
```

Big enough for

- *int*;
- *MPI\_Aint*, used in one-sided (and therefore big enough for *intptr\_t* and *ptrdiff\_t*);
- *MPI\_Offset*, used in file I/O.

However, this type could not be used in MPI-3 to describe send buffers.



# MPI 4 large count routines

C: routines with `_c` suffix

```
1 MPI_Count count;  
2 MPI_Send_c( buff, count, MPI_INT, ... );
```

also `MPI_Reduce_c`, `MPI_Get_c`, ... (some 190 routines in all)

Fortran: polymorphism rules

```
1 Integer(kind=MPI_COUNT_KIND) :: count  
2 call MPI_Send( buff, count, MPI_INTEGER, ... )
```



# Big count example

```
1 // pingpongbig.c
2 assert( sizeof(MPI_Count)>4 );
3 for ( int power=3; power<=10; power++ ) {
4     MPI_Count length=pow(10,power);
5     buffer = (double*)malloc( length*sizeof(double) );
6     MPI_Ssend_c
7         (buffer,length,MPI_DOUBLE,
8          processB,0,comm);
9     MPI_Recv_c
10        (buffer,length,MPI_DOUBLE,
11          processB,0,comm,MPI_STATUS_IGNORE);
```

# Same in F08

```
1 !! pingpongbig.F90
2 integer :: power,countbytes
3 Integer(KIND=MPI_COUNT_KIND) :: length
4 call MPI_Sizeof(length,countbytes,ierr)
5 if (procno==0) &
6     print *, "Bytes in count:",countbytes
7     length = 10**power
8     allocate( senddata(length),recvdata(length) )
9     call MPI_Send(senddata,length,MPI_DOUBLE_PRECISION, &
10                  processB,0, comm)
11    call MPI_Recv(recvdata,length,MPI_DOUBLE_PRECISION, &
12                  processB,0, comm,MPI_STATUS_IGNORE)
```

# MPI\_Send

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

# MPI 4 large count querying

C:

```
1 MPI_Count count;  
2 MPI_Get_count_c( &status,MPI_INT, &count );  
3 MPI_Get_elements_c( &status,MPI_INT, &count );
```

Fortran:

```
1 Integer(kind=MPI_COUNT_KIND) :: count  
2 call MPI_Get_count( status,MPI_INTEGER,count )  
3 call MPI_Get_elements( status,MPI_INTEGER,count )
```

# Compound messages

MPI-3 mechanism, deprecated in MPI-4.1:  
send a number of contiguous types:

```
1 MPI_Datatype blocktype;
2 MPI_Type_contiguous(mediumsize,MPI_FLOAT,&blocktype);
3 MPI_Type_commit(&blocktype);
4 if (procno==sender) {
5   MPI_Send(source,nblocks,blocktype,receiver,0,comm);
```

By composing types you can make a 'big type'. Use

*MPI\_Type\_get\_extent\_x*, *MPI\_Type\_get\_true\_extent\_x*, *MPI\_Get\_elements\_x*  
to query.

```
1 MPI_Count recv_count;
2 MPI_Get_elements_x(&recv_status,MPI_FLOAT,&recv_count);
```



# Advanced (MPI-3/4) topics

# Justification

Recent additions to the MPI standard allow your code to deal with unusual scenarios or very large scale runs.

# Advanced collectives

# Non-blocking collectives

- Collectives are blocking.
- Compare blocking/non-blocking sends:

*MPI\_Send* → *MPI\_Isend*

immediate return of control, produce request object.

- Non-blocking collectives:

*MPI\_Bcast* → *MPI\_Ibcast*

Same:

```
1 MPI_Isomething( <usual arguments>, MPI_Request *req);
```

- Considerations:
  - Calls return immediately;
  - the usual story about buffer reuse
  - Requires *MPI\_Wait...* for completion.
  - Multiple collectives can complete in any order
- Why?
  - Use for overlap communication/computation
  - Imbalance resilience
  - Allows pipelining



# MPI\_Ibcast

Name	Param name	Explanation	C type	F type
MPI_Ibcast (				
MPI_Ibcast_c (				
buffer		starting address of buffer	void*	TYPE(*), DIMENSION(..)
count		number of entries in buffer	[ int MPI_Count	INTEGER
datatype		datatype of buffer	MPI_Datatype	TYPE(MPI_Datatype)
root		rank of broadcast root	int	INTEGER
comm		communicator	MPI_Comm	TYPE(MPI_Comm)
request		communication request	MPI_Request*	TYPE(MPI_Request)
)				

# Overlapping collectives and computation

Independent collective and local operations:

$$y \leftarrow Ax + (x^t x)y$$

```
1 MPI_Iallreduce( .... x ... , &request);  
2 // compute the matrix vector product  
3 MPI_Wait(request);  
4 // do the addition
```

# Simultaneous reductions

Do two reductions (on the same communicator) with different operators simultaneously:

$$\begin{aligned}\alpha &\leftarrow x^t y \\ \beta &\leftarrow \|z\|_\infty\end{aligned}$$

which translates to:

```
1 MPI_Request reqs[2];
2 MPI_Iallreduce
3   ( &local_xy, &global_xy, 1, MPI_DOUBLE, MPI_SUM, comm,
4     &(reqs[0]) );
5 MPI_Iallreduce
6   ( &local_xinf, &global_xin, 1, MPI_DOUBLE, MPI_MAX, comm,
7     &(reqs[1]) );
8 MPI_Waitall(2, reqs, MPI_STATUSES_IGNORE);
```



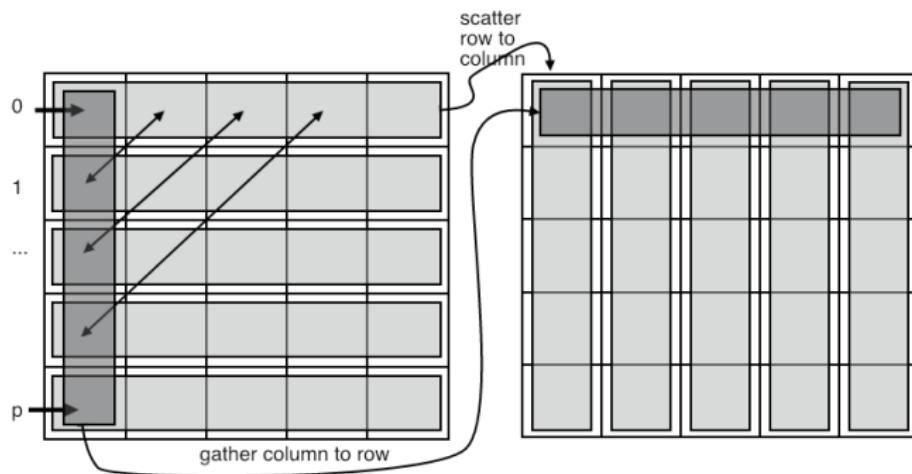
# Matching collectives

Blocking and non-blocking collectives don't match:  
either all processes call the non-blocking or all call the blocking one.  
Thus the following code is incorrect:

```
1 if (rank==root)
2   MPI_Reduce( &x /* ... */ root,comm );
3 else
4   MPI_Ireduce( &x /* ... */ root,comm,&req);
```

This is unlike the point-to-point behavior of non-blocking calls: you can catch a message with *MPI\_Irecv* that was sent with *MPI\_Send*.

# Transpose as gather/scatter



Every process needs to do a scatter or gather.

# Simultaneous collectives

Transpose matrix by scattering all rows simultaneously.  
Each scatter involves all processes, but with a different spanning tree.

```
1 MPI_Request scatter_requests[nprocs];
2 for (int iproc=0; iproc<nprocs; iproc++) {
3     MPI_Iscatter( regular,1,MPI_DOUBLE,
4                   &(transpose[iproc]),1,MPI_DOUBLE,
5                   iproc,comm,scatter_requests+iproc);
6 }
7 MPI_Waitall(nprocs,scatter_requests,MPI_STATUSES_IGNORE);
```

## Persistent collectives

# Persistent collectives (MPI-4)

Similar to persistent send/recv:

```
1 double *buffer;
2 MPI_Allreduce_init( buffer ...., &request );
3 for ( ... ) {
4     // fill buffer
5     MPI_Start( request );
6     // possibly other activities
7     MPI_Wait( &request );
8 }
9 MPI_Request_free( &request );
```

Available for all collectives and neighborhood collectives.

# Example

```
1 // powerpersist1.c
2 double localnorm,globalnorm=1.;
3 MPI_Request reduce_request;
4 MPI_Allreduce_init(
5     ( &localnorm,&globalnorm,1,MPI_DOUBLE,MPI_SUM,
6     comm,MPI_INFO_NULL,&reduce_request);
7 for (int it=0; ; it++) {
8     /*
9      * Matrix vector product
10     */
11    matmult(indata,outdata,buffersize);
12
13 // start computing norm of output vector
14 localnorm = local_12_norm(outdata,buffersize);
15 double old_globalnorm = globalnorm;
16 MPI_Start( &reduce_request );
17
18 // end computing norm of output vector
19 MPI_Wait( &reduce_request,MPI_STATUS_IGNORE );
20 globalnorm = sqrt(globalnorm);
21 // now 'globalnorm' is the L2 norm of 'outdata'
22 scale(outdata,indata,buffersize,1./globalnorm);
23 }
24 MPI_Request_free( &reduce_request );
```

---

Note also the *MPI\_Info* parameter.



# Persistent vs non-blocking

Both request-based.

- Non-blocking is ‘ad hoc’: buffer info not known before the collective call.
- Persistent allows ‘planning ahead’: management of internal buffers and such.

Request handling:

- Non-blocking: wait deallocates the request
- Persistent: wait deactivates the request, still requires `MPI_Request_free`.

## Non-blocking barrier

# Just what is a barrier?

- Barrier is not *time* synchronization but *state* synchronization.
- Test on non-blocking barrier: ‘has everyone reached some state’

# Use case: adaptive refinement

- Some processes decide locally to alter their structure
- ... need to communicate that to neighbors
- Problem: neighbors don't know whether to expect update calls, if at all.
- Solution:
  - send update msgs, if any;
  - then post barrier.
  - Everyone probe for updates, test for barrier.

# Use case: distributed termination detection

- Distributed termination detection (Matocha and Kamp, 1998): draw a global conclusion with local operations
- Everyone posts the barrier when done;
- keeps doing local computation while testing for the barrier to complete

# MPI\_Ibarrier

Name	Param name	Explanation	C type	F type
<code>MPI_Ibarrier (</code> <code>comm</code> <code>request</code> <code>)</code>		communicator communication request	<code>MPI_Comm</code> <code>MPI_Request*</code>	<code>TYPE(MPI_Comm)</code> <code>TYPE(MPI_Request)</code>

# Step 1

Do sends, post barrier.

```
1 // ibarrierprobe.c
2 if (i_do_send) {
3     /*
4      * Pick a random process to send to,
5      * not yourself.
6      */
7     int receiver = rand()%nprocs;
8     MPI_Ssend(&data,1,MPI_FLOAT,receiver,0,comm);
9 }
10 /*
11  * Everyone posts the non-blocking barrier
12  * and gets a request to test/wait for
13 */
14 MPI_Request barrier_request;
15 MPI_Ibarrier(comm,&barrier_request);
```



# Step 2

Poll for barrier and messages

```
1 for ( ; ; step++) {  
2     int barrier_done_flag=0;  
3     MPI_Test(&barrier_request,&barrier_done_flag,  
4               MPI_STATUS_IGNORE);  
5 //stop if you're done!  
6     if (barrier_done_flag) {  
7         break;  
8     } else {  
9 // if you're not done with the barrier:  
10        int flag; MPI_Status status;  
11        MPI_Iprobe  
12            ( MPI_ANY_SOURCE,MPI_ANY_TAG,  
13              comm, &flag, &status );  
14        if (flag) {  
15 // absorb message!
```



# Shared memory

# Shared memory myths

Myth:

*MPI processes use network calls, whereas OpenMP threads access memory directly, therefore OpenMP is more efficient for shared memory.*

Truth:

*MPI implementations use copy operations when possible, whereas OpenMP has thread overhead, and affinity/coherence problems.*

Main problem with MPI on shared memory: data duplication.

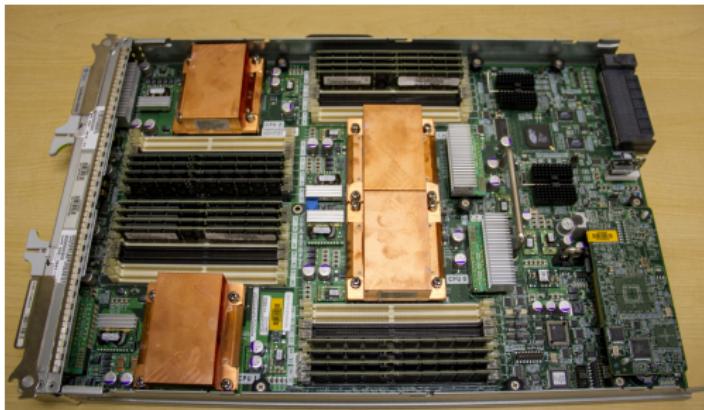
# MPI shared memory

- Shared memory access: two processes can access each other's memory through `double*` (and such) pointers, if they are on the same shared memory.
- Limitation: only window memory.
- Non-use case: remote update. This has all the problems of traditional shared memory (race conditions, consistency).
- Good use case: every process needs access to large read-only dataset  
Example: ray tracing.

# Shared memory treatments in MPI

- MPI uses optimizations for shared memory: copy instead of socket call
- One-sided offers ‘fake shared memory’: you can access another process’ data, but only through function calls.
- MPI-3 shared memory gives you a pointer to another process’ memory,  
*if that process is on the same shared memory.*

# Shared memory per cluster node



- Cluster node has shared memory
- Memory is attached to specific socket
- beware Non-Uniform Memory Access (NUMA) effects



# Shared memory interface

Here is the high level overview; details next.

- Use `MPI_Comm_split_type` to find processes on the same shared memory
- Use `MPI_Win_allocate_shared` to create a window between processes on the same shared memory  
(MPI-4.1: other window creation calls allowed, but the burden is on you!)
- Use `MPI_Win_shared_query` to get pointer to another process' window data.
- You can now use `memcpy` instead of `MPI_Put`.

# Discover shared memory

- `MPI_Comm_split_type` splits into communicators of same type.
- Use type: `MPI_COMM_TYPE_SHARED` splitting by shared memory.

Code:

```
1 // commssplittype.c
2 MPI_Info info;
3 MPI_Comm_split_type
4     (MPI_COMM_WORLD,
5      MPI_COMM_TYPE_SHARED,
6      procno,info,&sharedcomm)
7 ;
8 MPI_Comm_size
9     (sharedcomm,&new_nprocs);
10 MPI_Comm_rank
11     (sharedcomm,&new_procno);
```

Output:

```
1 make[3]: 'commssplittype' is up to
           ↪date.
2 TACC: Starting up job 4356245
3 TACC: Starting parallel tasks...
4 There are 10 ranks total
5 [0] is processor 0 in a shared
           ↪group of 5, running on
           ↪c209-010.frontera.tacc.utexas.edu
6 [5] is processor 0 in a shared
           ↪group of 5, running on
           ↪c209-011.frontera.tacc.utexas.edu
7 TACC: Shutdown complete. Exiting.
```



# Exercise 46 (shareddata)

Write a program that uses `MPI_Comm_split_type` to analyze for a run

1. How many nodes there are;
2. How many processes there are on each node.

If you run this program on an unequal distribution, say 10 processes on 3 nodes, what distribution do you find?

```
1 Nodes: 3; processes: 10
2 TACC: Starting up job 4210429
3 TACC: Starting parallel tasks...
4 There are 3 nodes
5 Node sizes: 4 3 3
6 TACC: Shutdown complete. Exiting.
```

# Allocate shared window

Use `MPI_Win_allocate_shared` to create a window that can be shared;

- Has to be on a communicator on shared memory
- Example: window is one double.

```
1 // sharedbulk.c
2 MPI_Win node_window;
3 MPI_Aint window_size; double *window_data;
4 if (onnode_procid==0)
5     window_size = sizeof(double);
6 else window_size = 0;
7 MPI_Win_allocate_shared
8     ( window_size,sizeof(double),MPI_INFO_NULL,
9     nodecomm,
10     &window_data,&node_window);
```



# Get pointer to other windows

Use `MPI_Win_shared_query`:

```
1 MPI_Aint window_size0; int window_unit; double *win0_addr;  
2 MPI_Win_shared_query  
3   ( node_window,0,  
4     &window_size0,&window_unit, &win0_addr );
```

# MPI\_Win\_shared\_query

Name	Param name	Explanation	C type	F type
MPI_Win_shared_query (				
MPI_Win_shared_query_c (				
win		shared memory window object	MPI_Win	TYPE(MPI_Win)
rank		rank in the group of window win or MPI_PROC_NULL	int	INTEGER
size		size of the window segment	MPI_Aint*	INTEGER (KIND=MPI_ADDRESS_KIND)
disp_unit		local unit size for displacements, in bytes	[ int* MPI_Aint* ]	INTEGER
baseptr		address for load/store access to window segment	void*	TYPE(C_PTR)
)				

# Allocated memory

Memory will be allocated contiguously  
convenient for address arithmetic,  
not for NUMA: set `alloc_shared_noncontig` true in `MPI_Info` object.

Example: each window stores one double. Measure distance in bytes:

Strategy: default behavior of  
shared window allocation

Distance 1 to zero: 8

Distance 2 to zero: 16

Strategy: allow non-contiguous  
shared window allocation

Distance 1 to zero: 4096

Distance 2 to zero: 8192

Question: what is going on here?



# Exciting example: bulk data

- Application: ray tracing:  
large read-only data structure describing the scene
- traditional MPI would duplicate:  
excessive memory demands
- Better: allocate shared data on process 0 of the shared  
communicator
- Everyone else points to this object.

# Split by other hardware types

MPI-4: split by other hardware features through `MPI_COMM_TYPE_HW_GUIDED` and  
`MPI_Get_hw_resource_types`)

# Process management

# Overview

This section discusses processes management; intra communicators.

Commands learned:

- *MPI\_Comm\_spawn, MPI\_UNIVERSE\_SIZE*
- *MPI\_Comm\_get\_parent, MPI\_Comm\_remote\_size*

# Process management

- PVM was a precursor of MPI: could dynamically create new processes.
- It took MPI a while to catch up.
- Use `MPI_Attr_get` to retrieve `MPI_UNIVERSE_SIZE` attribute indicating space for creating more processes outside `MPI_COMM_WORLD`.
- New processes have their own `MPI_COMM_WORLD`.
- Communication between the two communicators: ‘inter communicator’  
(the old type is ‘intra communicator’)

# Space for processes

Probably a machine dependent component.

Suggested standard:

```
mpiexec -n 4 -usize 8 spawn_manager
```

Intel MPI at TACC:

```
MY_MPIRUN_OPTIONS="-usize 8" ibrun -np 4 spawn_manager
```

Discover size of the universe:

```
1 MPI_Attr_get(MPI_COMM_WORLD, MPI_UNIVERSE_SIZE,  
2   (void*)&universe_sizep, &flag);
```



# Manager program

```
1 int universe_size, *universe_size_attr, uflag;
2 MPI_Comm_get_attr
3   (comm_world,MPI_UNIVERSE_SIZE,
4   &universe_size_attr,&uflag);
5 if (uflag) {
6   universe_size = *universe_size_attr;
7 } else {
8   printf("This MPI does not support UNIVERSE_SIZE.\nUsing world size");
9   universe_size = world_n;
10 }
11 int work_n = universe_size - world_n;
12 if (world_p==0) {
13   printf("A universe of size %d leaves room for %d workers\n",
14         universe_size,work_n);
15   printf(.. spawning from %s\n",procname);
16 }
```



# Manager program (cont'd)

```
1 const char *workerprogram = "./spawnapp";
2 MPI_Comm_spawn(workerprogram,MPI_ARGV_NULL,
3                 work_n,MPI_INFO_NULL,
4                 0,comm_world,&comm_inter,NULL);
```

# Worker program

```
1 // spawnworker.c
2 MPI_Comm_size(MPI_COMM_WORLD,&nworkers);
3 MPI_Comm_rank(MPI_COMM_WORLD,&workerno);
4 MPI_Comm_get_parent(&parent);
```

# Were you spawned?

```
1 // spawnapp.c
2 MPI_Comm comm_parent;
3 MPI_Comm_get_parent(&comm_parent);
4 int is_child = (comm_parent!=MPI_COMM_NULL);
5 if (is_child) {
6     int nworkers,workerno;
7     MPI_Comm_size(MPI_COMM_WORLD,&nworkers);
8     MPI_Comm_rank(MPI_COMM_WORLD,&workerno);
9     printf("I detect I am worker %d/%d running on %s\n",
10           workerno,nworkers,procname);
```

# Process topologies

# Overview

This section discusses topologies:

- Cartesian topology
- MPI-1 Graph topology
- MPI-3 Graph topology

Commands learned:

- *MPI\_Dist\_graph\_create*, *MPI\_DIST\_GRAPH*, *MPI\_Dist\_graph\_neighbors\_count*
- *MPI\_Neighbor\_allgather* and such

# Process topologies

- Processes don't communicate at random
- Example: Cartesian grid, each process 4 (or so) neighbors
- Express operations in terms of topology
- Elegance of expression
- MPI can optimize

# Process reordering

- Consecutive process numbering often the best:  
divide array by chunks
- Not optimal for grids or general graphs:
- MPI is allowed to renumber ranks
- Graph topology gives information from which MPI can deduce  
renumbering

# MPI-1 topology

- Cartesian topology
- Graph topology, globally specified.  
Not scalable, do not use!

# MPI-3 topology

- Graph topologies locally specified: scalable!  
Limit cases: each process specifies its own connectivity one process specifies whole graph.
- Neighborhood collectives:  
expression close to the algorithm.

# Graph topologies

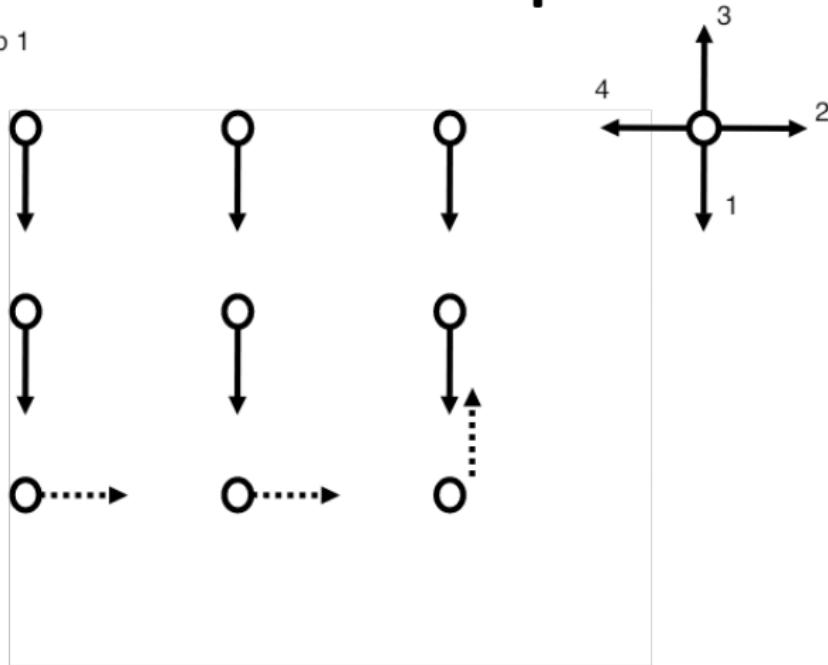
# Example: 5-point stencil

Neighbor exchange, spelled out:

- Each process communicates down/right/up/left
- Send and receive at the same time.
- Can optimally be done in four steps

# Step 1

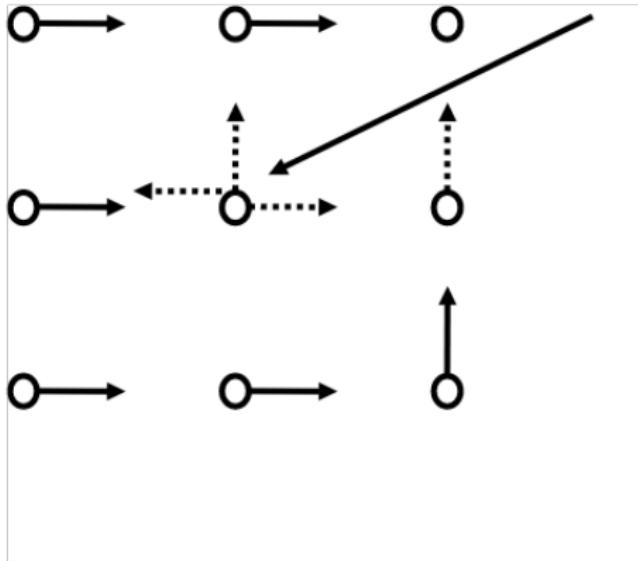
Step 1



Step 2

## Step 2

blocked

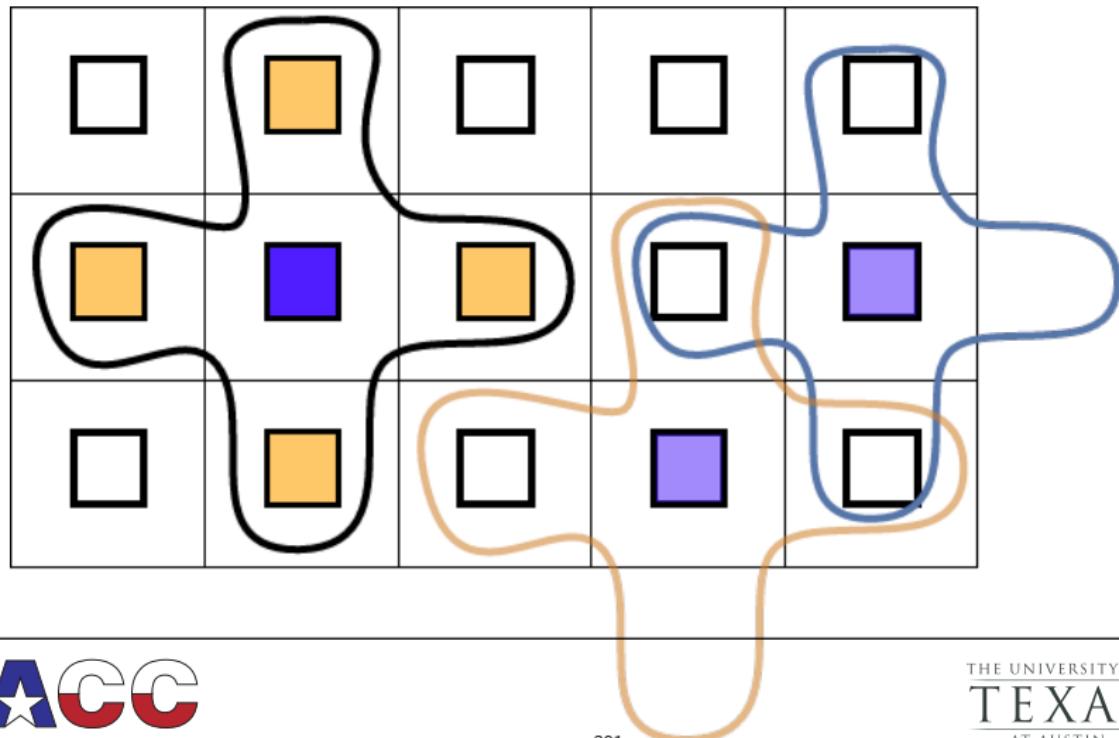


The middle node is blocked because all its targets are already receiving  
or a channel is occupied:  
one missed turn

# Neighborhood collective

This is really a 'local gather':  
each node does a gather from its neighbors in whatever order.

*MPI\_Neighbor\_allgather*



# Why neighborhood collectives?

- Using `MPI_Isend` / `MPI_Irecv` is like spelling out a collective, imposes order;
- Collectives can use pipelining as opposed to sending a whole buffer;
- Collectives can use spanning trees as opposed to direct connections.

# Create graph topology

```
1 int MPI_Dist_graph_create
2   (MPI_Comm comm_old, int nsources, const int sources[],
3    const int degrees[], const int destinations[],
4    const int weights[], MPI_Info info, int reorder,
5    MPI_Comm *comm_dist_graph)
```

- *nsources* how many source nodes described? (Usually 1)
- *sources* the processes being described (Usually *MPI\_Comm\_rank* value)
- *degrees* how many processes to send to
- *destinations* their ranks
- *weights*: usually set to *MPI\_UNWEIGHTED*.
- *info*: *MPI\_INFO\_NULL* will do
- *reorder*: 1 if dynamically reorder processes

# Neighborhood collectives

```
1 int MPI_Neighbor_allgather
2   (const void *sendbuf, int sendcount, MPI_Datatype sendtype,
3    void *recvbuf, int recvcount, MPI_Datatype recvtype,
4    MPI_Comm comm)
```

Like an ordinary `MPI_Allgather`, but  
the receive buffer has a length enough for `degree` messages  
(instead of comm size).

# Neighbor querying

After `MPI_Neighbor_allgather` data in the buffer is *not* in normal rank order.

- `MPI_Dist_graph_neighbors_count` gives actual number of neighbors.  
(Why do you need this?)
- `MPI_Dist_graph_neighbors` lists neighbor numbers.

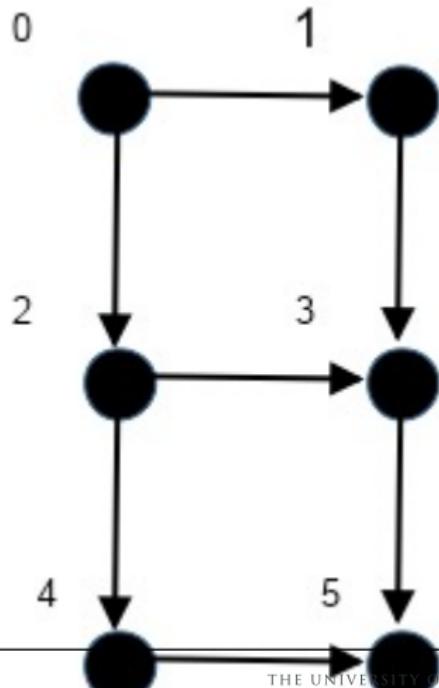
# MPI\_Dist\_graph\_neighbors\_count

# MPI\_Dist\_graph\_neighbors

# Example: Systolic graph

Code:

```
1 // graph.c
2 for ( int i=0; i<=1; i++ ) {
3     int neighb_i = proc_i+i;
4     if (neighb_i<0 || neighb_i>=idim)
5         continue;
6     int j = 1-i;
7     int neighb_j = proc_j+j;
8     if (neighb_j<0 || neighb_j>=jdim)
9         continue;
10    destinations[ degree++ ] =
11        PROC(neighb_i,neighb_j,idim,jdim);
12 }
13 MPI_Dist_graph_create
14 (comm,
15 /* I specify just one proc: me */ 1,
16 &procno,&degree,destinations,weights,
17 MPI_INFO_NULL,0,
18 &comm2d
19 );
```



# Output

Code:

```
1 int indegree,outdegree,  
2   weighted;  
3 MPI_Dist_graph_neighbors_count  
  
4   (comm2d,  
5     &indegree,&outdegree,  
6     &weighted);  
7 int  
8   my_ij[2] = {proci,procj},  
9   other_ij[4][2];  
10 MPI_Neighbor_allgather  
11   ( my_ij,2,MPI_INT,  
12     other_ij,2,MPI_INT,  
13       comm2d );
```

Output:

```
1 [ 0 = (0,0)] has 2  
    ↪outbound: 1, 2,  
2   0 inbound:  
3 [ 1 = (0,1)] has 1  
    ↪outbound: 3,  
4   1 inbound: (0,0)=0  
5 [ 2 = (1,0)] has 2  
    ↪outbound: 3, 4,  
6   1 inbound: (0,0)=0  
7 [ 3 = (1,1)] has 1  
    ↪outbound: 5,  
8   2 inbound: (0,1)=1  
    ↪(1,0)=2  
9 [ 4 = (2,0)] has 1  
    ↪outbound: 5,  
10  1 inbound: (1,0)=2  
11 [ 5 = (2,1)] has 0  
    ↪outbound:  
12  2 inbound: (1,1)=3  
    ↪(2,0)=4
```



# Query

Explicit query of neighbor process ranks.

Code:

```
1 int indegree,outdegree,  
2     weighted;  
3 MPI_Dist_graph_neighbors_count  
4     (comm2d,  
5      &indegree,&outdegree,  
6      &weighted);  
7 int  
8     my_ij[2] = {proci,procj},  
9     other_ij[4][2];  
10 MPI_Neighbor_allgather  
11   ( my_ij,2,MPI_INT,  
12     other_ij,2,MPI_INT,  
13     comm2d );
```

Output:

```
1    0 inbound:  
2    1 inbound: 0  
3    1 inbound: 0  
4    2 inbound: 1 2  
5    1 inbound: 2  
6    2 inbound: 4 3
```

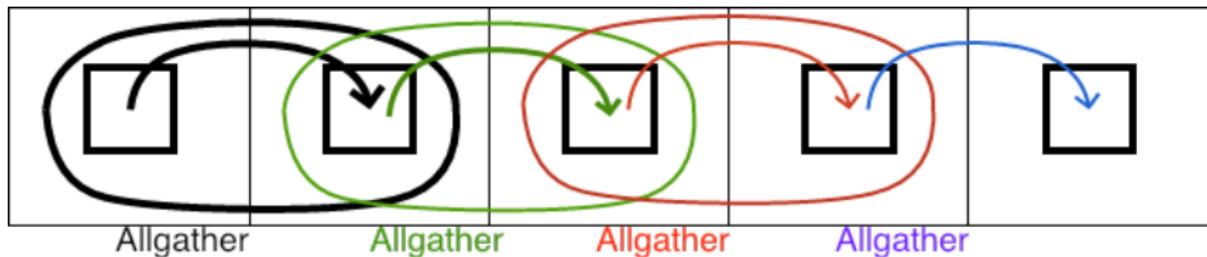
# Exercise 47 (rightgraph)

▶ Earlier rightsend exercise

Revisit exercise 19 and solve it using `MPI_Dist_graph_create`. Use figure 402 for inspiration.

Use a degree value of 1.

# Inspiring picture for the previous exercise



Solving the right-send exercise with neighborhood collectives

# Hints for the previous exercise

Two approaches:

1. Declare just one source: the previous process. Do this! Or:
2. Declare two sources: the previous and yourself. In that case bear in mind slide 395.

# More graph collectives

- Heterogeneous: `MPI_Neighbor_alltoallw`.
- Non-blocking: `MPI_Ineighbor_allgather` and such
- Persistent: `MPI_Neighbor_allgather_init`, `MPI_Neighbor_allgatherv_init`.

# Other

# Tracing, performance, and such

# Overview

We briefly touch on peripheral issues issues to MPI.

# CMake

# cmake

- . MPI is discoverable by cmake.

# 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 . )
```



# Errors

# Built-in handlers

Default: global termination.

```
1 MPI_Comm_set_errhandler(MPI_COMM_WORLD,MPI_ERRORS_ARE_FATAL);
```

**MPI-4:** *Only terminate on communicator: MPI\_ERRORS\_ABORT.*

Local handling: *MPI\_ERRORS\_RETURN*:

# Handlers on specific classes

Associate error handler with communicator:

*MPI\_Comm\_set\_errhandler* *MPI\_Comm\_get\_errhandler*

Other:

- *MPI\_File\_set\_errhandler*, *MPI\_File\_call\_errhandler*,

**MPI-4:** *MPI\_Session\_set\_errhandler*, *MPI\_Session\_call\_errhandler*,

*MPI\_Win\_set\_errhandler*, *MPI\_Win\_call\_errhandler*.

# Handling errors

```
1 char errtxt[MPI_MAX_ERROR_STRING];
2 int err = status.MPI_ERROR;
3 int len=MPI_MAX_ERROR_STRING;
4 MPI_Error_string(err,errtxt,&len);
5 printf("Waitall error: %d %s\n",err,errtxt);
```

# Define new errors

```
1 int nonzero_code;
2 MPI_Add_error_code(nonzero_class,&nonzero_code);
3 MPI_Add_error_string(nonzero_code,"Attempting to send zero buffer");
```

## **Performance measurement**

# Timers

MPI has a *wall clock* timer: `MPI_Wtime` which gives the number of seconds from a certain point in the past.

The timer has a resolution of `MPI_Wtick`

Timers can be global

```
1 int *v,flag;
2 MPI_Attr_get( comm, MPI_WTIME_IS_GLOBAL, &v, &flag );
3 if (mytid==0) printf("Time synchronized? %d->%d\n",flag,*v);
```

but probably aren't.

# Example

```
1 // pingpong.c
2 if (procno==processA) {
3     t = MPI_Wtime();
4     for (int n=0; n<NEXPERIMENTS; n++) {
5         MPI_Send(send,1,MPI_DOUBLE,
6         MPI_Recv(recv,1,MPI_DOUBLE,
7     }
8     t = MPI_Wtime()-t; t /= NEXPERIMENTS;
```

# Global timing

Processes don't start/end simultaneously. What does a timing result mean overall? Take average or maximum?

Alternative:

```
1      MPI_Barrier(comm)
2      t = MPI_Wtime();
3      // something happens here
4      MPI_Barrier(comm)
5      t = MPI_Wtime()-t;
```



# Profiling

See other lecture: MPIP, TAU, et cetera.

# Your own profiling interface

Every routine `MPI_Something` calls a routine `PMPI_Something` that does the actual work. You can now write your `MPI_...` routine which calls `PMPI_...`, and inserting your own profiling calls.

```
main() {  
    MPI_Send ( buffer,ct,tp, ... );  
}
```

call

```
int MPI_Send ( buffer,ct,tp, ... ) {  
    PMPI_Send( buffer,ct,tp, ... );
```

# Programming for performance

# Eager limit

- Optimization for small messages: bypass rendez-vous protocol (slide 154)
- Cross-over point: ‘Eager limit’.
- Force efficient messages by increasing the eager limit.
- Beware: decreasing payoff for large messages, and
- Beware: buffers for eager send eat into your available memory.

# Eager limit setting

- For Intel MPI: I\_MPI\_EAGER\_THRESHOLD
- mvapich2: MV2\_IBA\_EAGER\_THRESHOLD
- OpenMPI: OpenMPI the --mca options *btl\_openib\_eager\_limit* and *btl\_openib\_rndv\_eager\_limit*.

# Blocking versus non-blocking

- Non-blocking sends  $\text{MPI_Isend}$  /  $\text{MPI_Irecv}$  can be more efficient than blocking
- Also: allow overlap computation/communication (latency hiding)
- However: can usually not be considered a replacement.

# Progress

MPI is not magically active in the background, so latency hiding is not automatic. Same for passive target synchronization and non-blocking barrier completion.

- Dedicated communications processor or thread.  
This is implementation dependent; for instance, Intel MPI:  
`I_MPI_ASYNC_PROGRESS_...` variables.
- Force progress by occasional calls to a polling routine such as  
`MPI_Iprobe`.

# Persistent sends

If a communication between the same pair of processes, involving the same buffer, happens regularly, it is possible to set up a *persistent communication*.

- `MPI_Send_init`
- `MPI_Recv_init`
- `MPI_Start`

# Buffering

- MPI has internal buffers: copying costs performance
- Use your own buffer:
  - *MPI\_Buffer\_attach*
  - *MPI\_Bsend*
- Copying is also a problem for derived datatypes.

# Graph topology and neighborhood collectives

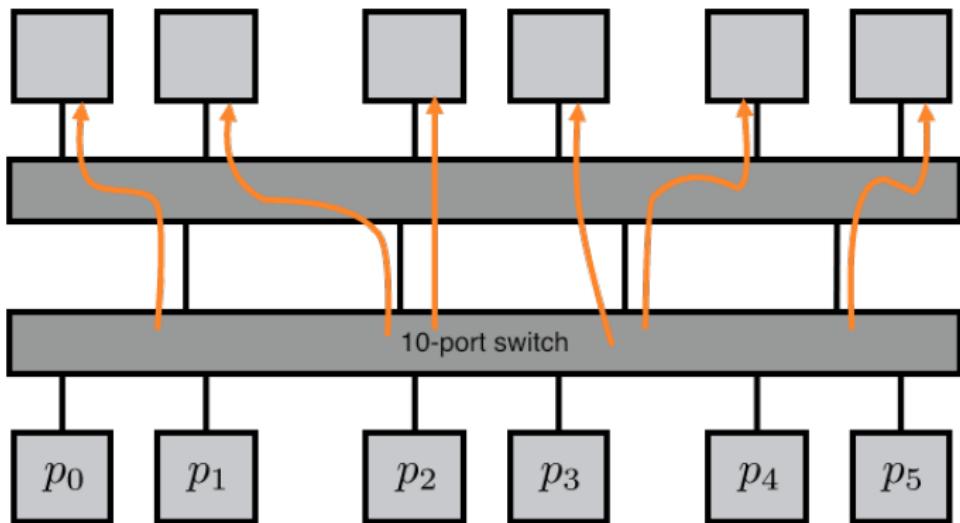
- Mapping problem to architecture sometimes not trivial
- Load balancers: *ParMetis, Zoltan*
- Graph topologies: *MPI\_Dist\_graph\_adjacent*: allowed to reorder ranks for proximity
- Neighborhood collectives allow MPI to schedule optimally.
  - *MPI\_Neighbor\_allgather* (and *MPI\_Neighbor\_allgather\_v*)
  - *MPI\_Neighbor\_alltoall*

# Network issues

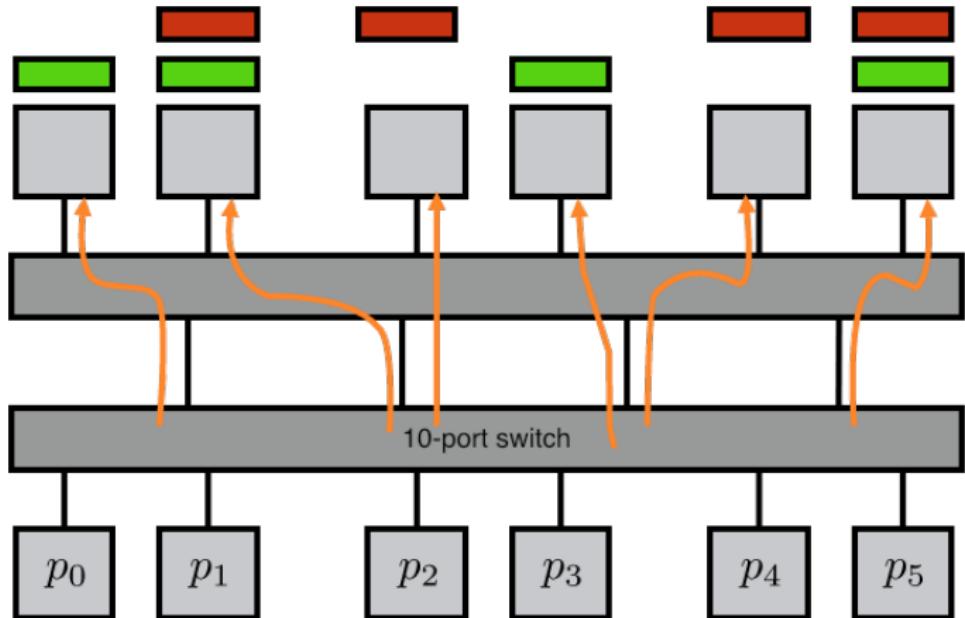
Network contention means that

- Your messages can collide with other jobs
- messages within your job can collide

# Output routing



# Contention



# Offloading and onloading

- Network cards can offer assistance
- Mellanox: off-loading  
limited repertoire of scenarios where it helps
- Intel disagrees: on-loading
- Either way, investigate the capabilities of your network.