

Getting started with MPI

Parallel Programming with the Message Passing Interface

Carsten Kutzner, 12.10.2006



compiling & running • most useful commands • parallelization concepts • performance monitoring • MPI resources

What is MPI?

- ▶ a library, not a language
 - ▶ specifies names, parameters and results of functions / subroutines to be called from C, C++, and Fortran programs
 - ▶ programs that use MPI are compiled with ordinary compilers (gcc, icc) and linked with the MPI library
- ▶ MPI standard defined in 1994
 - ▶ MPI-1 core functionality
- ▶ MPI-2 add-ons 1995-97
 - ▶ remote memory operations
 - ▶ parallel I/O
 - ▶ dynamic process management
- ▶ MPI-3 will address future needs
 - ▶ interoperability among MPI implementations?
 - ▶ non-blocking collectives?
 - ▶ ...



Why use MPI?

- ▶ high performance
 - ▶ today used from 1 – 131072 processors
- ▶ flexible
- ▶ portable

	MPI	Cluster OpenMP	OpenMP	threads
shared memory	yes	yes	yes	yes
distributed memory	yes	yes	not possible	not possible
performance	++	+/o	+	+
ease of use	o	+/o	+	o

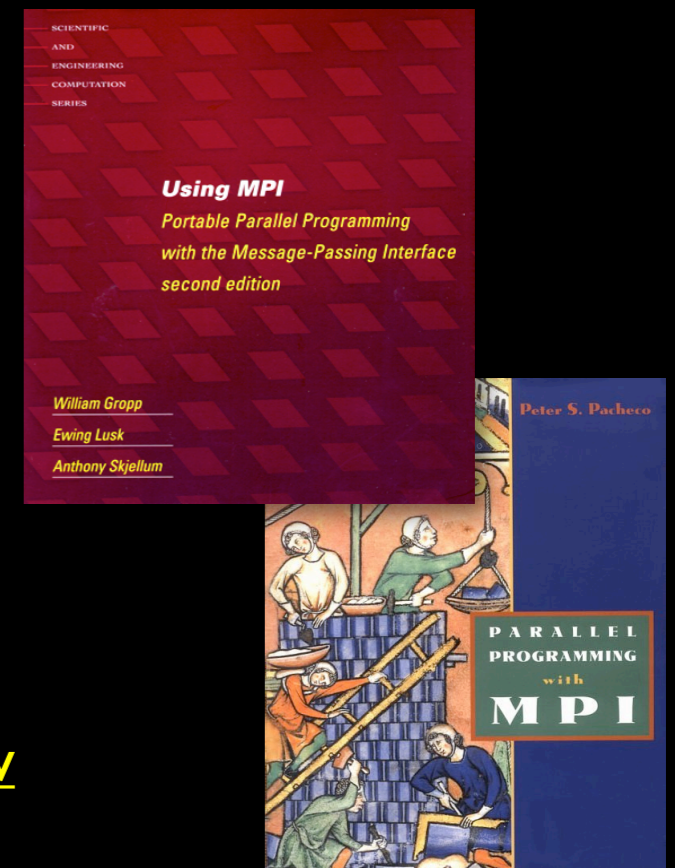
MPI Resources

▶ Books

- ▶ Pacheco 1997: *Parallel Programming with MPI*
- ▶ Gropp, Lusk, Skjellum 1999: *Using MPI*
- ▶ Gropp, Lusk, Thakur 1999: *Using MPI-2*

▶ On the web

- ▶ the MPI forum: www-unix.mcs.anl.gov/mpi/ supplies the MPI standards as PS, PDF
- ▶ web pages for MPI and MPE: www.mcs.anl.gov/mpi/www (multiple mirrors)
- ▶ MPICH (MPI CHameleon): www-unix.mcs.anl.gov/mpi/mpich2
- ▶ LAM/MPI (Local Area Multicomputer): www.lam-mpi.org
- ▶ OpenMPI: www.open-mpi.org
- ▶ en.wikipedia.org/wiki/Message_Passing_Interface
- ▶ www.linux-mag.com/extreme/
- ▶ newsgroup comp.parallel.mpi



Getting started with LAM

▶ edit your .bashrc

▶ on your workstation:

```
export PATH=/usr/local/lam-7.0.4/bin:$PATH  
export LAMHOME=/usr/local/lam-7.0.4
```

▶ on the clusters (e.g. coral4) export

```
/usr/local/Cluster-Apps/lam-7.1.1-64/
```

▶ compile ...

▶ mpicc hellompi.c -o hellompi.x

▶ mpicc – a compiler wrapper around gcc/icc/*cc that automatically links your code to the mpi library

▶ mpif77, mpif90 – for fortran

▶ mpic++ / mpiCC – for c++

▶ ... and run

▶ lamboot

▶ mpirun -np 2 ./hellompi.x

Getting started with LAM

- ▶ on multiple nodes a hostfile is needed:
 - ▶ `lamboot -v bhost`

```
node01 cpu=2  
node27 cpu=2  
node07 cpu=1  
node09 cpu=1
```

- ▶ tests can be done with any number of processes even on your workstation!
 - ▶ `lamboot`
 - ▶ `mpirun -np 17 ./hellompi.x`

Getting started with LAM

- ▶ you can also use our network of workstations:
 - ▶ you will need a homedir on each workstation with a .bashrc reading

```
export PATH=/usr/local/lam-7.0.4/bin:$PATH
export LAMHOME=/usr/local/lam-7.0.4
```
 - ▶ you must be able to login via ssh without password:
execute `ssh-keygen` (if you not have already) on your machine
`cd .ssh`
`scp id_dsa.pub you@remotehost:~/.ssh/authorized_keys`
 - ▶ `export LAMRSH="ssh"`
`lamboot -v bhost`

wes	cpu=2
platypus	cpu=2
tivoli	cpu=2
- ```
ckutzne@wes:~> lamnodes
n0 wes.mpibpc.gwdg.de:2:origin,this_node
n1 platypus.mpibpc.gwdg.de:2:
n2 tivoli.mpibpc.gwdg.de:2:
```
- ▶ the executable has to reside on all of the workstations in the same directory, e.g. `/netmount/coral4/you/`

# Do's and don'ts

- ▶ never mix LAM versions!!!
  - ▶ LAM-A-compiled program might run on LAM-B, but nothing is guaranteed!
- ▶ recompile when you move your program from A to B

| Cluster    | LAM v.  | arch   |
|------------|---------|--------|
| tivoli     | 7.0.4   | 32 bit |
| orca 1-2   | 7.0.6   | 32 bit |
| coral 1-4  | 7.1.1   | 64 bit |
| beluga 1-3 | 7.1.1   | 64 bit |
| kea 1-3    | 7.1.1-2 | 64 bit |
| slomo      | 7.1.1   | 32 bit |





# MPI commands & keywords

MPIO\_Request\_c2f  
MPIO\_Request\_f2c  
MPIO\_Test  
MPIO\_Wait  
MPI\_Abort  
MPI\_Address  
MPI\_Allgather  
MPI\_Allgatherv  
MPI\_Allreduce  
MPI\_Alltoall  
MPI\_Alltoallv  
MPI\_Attr\_delete  
MPI\_Attr\_get  
MPI\_Attr\_put  
MPI\_Barrier  
MPI\_Bcast  
MPI\_Bsend  
MPI\_Bsend\_init  
MPI\_Buffer\_attach  
MPI\_Buffer\_detach  
MPI\_CHAR  
MPI\_Cancel  
MPI\_Cart\_coords  
MPI\_Cart\_create  
MPI\_Cart\_get  
MPI\_Cart\_map  
MPI\_Cart\_rank  
MPI\_Cart\_shift  
MPI\_Cart\_sub  
MPI\_Cartdim\_get  
MPI\_Comm\_compare  
MPI\_Comm\_create  
MPI\_Comm\_dup  
MPI\_Comm\_free  
MPI\_Comm\_get\_name  
MPI\_Comm\_group  
MPI\_Comm\_rank  
MPI\_Comm\_remote\_group  
MPI\_Comm\_remote\_size  
MPI\_Comm\_set\_name  
MPI\_Comm\_size  
MPI\_Comm\_split  
MPI\_Comm\_test\_inter  
MPI\_DUP\_FN  
MPI\_Dims\_create  
MPI\_Errhandler\_create  
MPI\_Errhandler\_free  
MPI\_Errhandler\_get  
MPI\_Errhandler\_set  
MPI\_Error\_class  
MPI\_Error\_string  
MPI\_File\_c2f  
MPI\_File\_close  
MPI\_File\_delete  
MPI\_File\_f2c  
MPI\_File\_get\_amode  
MPI\_File\_get\_atomicity  
MPI\_File\_get\_byte\_offset  
MPI\_File\_get\_errhandler  
MPI\_File\_get\_group  
MPI\_File\_get\_info  
MPI\_File\_get\_position  
MPI\_File\_get\_position\_shared  
MPI\_File\_get\_size  
MPI\_File\_get\_type\_extent  
MPI\_File\_get\_view  
MPI\_File\_iread  
MPI\_File\_iread\_at  
MPI\_File\_iread\_shared  
MPI\_File\_iwrite  
MPI\_File\_iwrite\_at  
MPI\_File\_iwrite\_shared  
MPI\_File\_open  
MPI\_File\_open  
MPI\_File\_preallocate  
MPI\_File\_read  
MPI\_File\_read\_all  
MPI\_File\_read\_all\_begin  
MPI\_File\_read\_all\_end  
MPI\_File\_read\_at  
MPI\_File\_read\_at\_all  
MPI\_File\_read\_at\_all\_begin  
MPI\_File\_read\_at\_all\_end  
MPI\_File\_read\_ordered  
MPI\_File\_read\_ordered\_begin  
MPI\_File\_read\_ordered\_end  
MPI\_File\_read\_shared  
MPI\_File\_seek  
MPI\_File\_seek\_shared  
MPI\_File\_set\_atomicity  
MPI\_File\_set\_errhandler  
MPI\_File\_set\_info  
MPI\_File\_set\_size  
MPI\_File\_set\_view  
MPI\_File\_sync  
MPI\_File\_write  
MPI\_File\_write\_all  
MPI\_File\_write\_all\_begin  
MPI\_File\_write\_all\_end  
MPI\_File\_write\_at  
MPI\_File\_write\_at\_all  
MPI\_File\_write\_at\_all\_begin  
MPI\_File\_write\_at\_all\_end  
MPI\_File\_write\_ordered  
MPI\_File\_write\_ordered\_begin  
MPI\_File\_write\_ordered\_end  
MPI\_File\_write\_shared  
MPI\_Finalize  
MPI\_Finalized  
MPI\_Gather  
MPI\_Gatherv  
MPI\_Get\_count  
MPI\_Get\_elements  
MPI\_Get\_processor\_name  
MPI\_Get\_version  
MPI\_Graph\_create  
MPI\_Graph\_get  
MPI\_Graph\_map  
MPI\_Graph\_neighbors  
MPI\_Graph\_neighbors\_count  
MPI\_Graphdims\_get  
MPI\_Group\_compare  
MPI\_Group\_difference  
MPI\_Group\_excl  
MPI\_Group\_free  
MPI\_Group\_incl  
MPI\_Group\_intersection  
MPI\_Group\_range\_excl  
MPI\_Group\_range\_incl  
MPI\_Group\_rank  
MPI\_Group\_size  
MPI\_Group\_translate\_ranks  
MPI\_Group\_union  
MPI\_Ibsend  
MPI\_Info\_c2f  
MPI\_Info\_create  
MPI\_Info\_delete  
MPI\_Info\_dup  
MPI\_Info\_f2c  
MPI\_Info\_free  
MPI\_Info\_get  
MPI\_Info\_get\_nkeys  
MPI\_Info\_get\_nthkey  
MPI\_Info\_get\_valuelen  
MPI\_Info\_set  
MPI\_Init  
MPI\_Init\_thread  
MPI\_Init\_thread  
MPI\_Initialized  
MPI\_Initialize  
MPI\_Initialize  
MPI\_Intercomm\_create  
MPI\_Intercomm\_merge  
MPI\_Iprobe  
MPI\_Irecv  
MPI\_Irsend  
MPI\_Isend  
MPI\_Issend  
MPI\_Keyval\_create  
MPI\_Keyval\_free  
MPI\_NULL\_COPY\_FN  
MPI\_NULL\_DELETE\_FN  
MPI\_Op\_create  
MPI\_Op\_free  
MPI\_Pack  
MPI\_Pack\_size  
MPI\_Pcontrol  
MPI\_Probe  
MPI\_Recv  
MPI\_Recv\_init  
MPI\_Reduce  
MPI\_Reduce\_scatter  
MPI\_Request\_c2f  
MPI\_Request\_free  
MPI\_Rsend  
MPI\_Rsend\_init  
MPI\_Scan  
MPI\_Scatter  
MPI\_Scatterv  
MPI\_Send  
MPI\_Send\_init  
MPI\_Sendrecv  
MPI\_Sendrecv\_replace  
MPI\_Ssend  
MPI\_Ssend\_init  
MPI\_Start  
MPI\_Startall  
MPI\_Status\_c2f  
MPI\_Status\_set\_cancelled  
MPI\_Status\_set\_elements  
MPI\_Test  
MPI\_Test\_cancelled  
MPI\_Testall  
MPI\_Testany  
MPI\_Testsome  
MPI\_Topology  
MPI\_Type\_commit  
MPI\_Type\_contiguous  
MPI\_Type\_create\_darray  
MPI\_Type\_create\_indexed\_block  
MPI\_Type\_create\_subarray  
MPI\_Type\_extent  
MPI\_Type\_free  
MPI\_Type\_get\_contents  
MPI\_Type\_get\_envelope  
MPI\_Type\_hindexed  
MPI\_Type\_hvector  
MPI\_Type\_indexed  
MPI\_Type\_lb  
MPI\_Type\_size  
MPI\_Type\_struct  
MPI\_Type\_ub  
MPI\_Type\_vector  
MPI\_Unpack  
MPI\_Wait  
MPI\_Waitall  
MPI\_Waitany  
MPI\_Waitsome  
MPI\_Wtick  
MPI\_Wtime

MPI\_Something

# What you need to start

- ▶ **Minimum subset to start with**
  - ▶ `MPI_Init` – Initialize MPI
  - ▶ `MPI_Comm_size` – How many processes are there?
  - ▶ `MPI_Comm_rank` – My process number
  - ▶ `MPI_Send` – Send a message
  - ▶ `MPI_Recv` – Receive a message
  - ▶ `MPI_Finalize` – Close MPI universe
- ▶ **Collective communication (powerful)**
  - ▶ `MPI_Bcast` – Broadcast a variable to all processes
  - ▶ `MPI_Reduce` – Add up a variable across all processes
  - ▶ `MPI_Alltoall` – Complete communication across all processes
  - ▶ `MPI_Barrier` – Synchronize all processes
  - ▶ ...

# Initializing & quitting

Fortran:

```
PROGRAM hello

INCLUDE 'mpif.h'
INTEGER err

CALL MPI_INIT(err)
PRINT *, "Hello world!"
CALL MPI_FINALIZE(err)

END
```

- ▶ a header file (mpi.h / mpif.h) has to be included which contains the MPI definitions and function prototypes
- ▶ MPI routines return an error code indicating whether or not they run successfully.

```
if (err == MPI_SUCCESS)
{
 ... /* routine ran correctly */
}
else
{
}
```

- ▶ output

```
> mpirun -np 4 ./helloworld.x
Hello world!
Hello world!
Hello world!
Hello world!
```

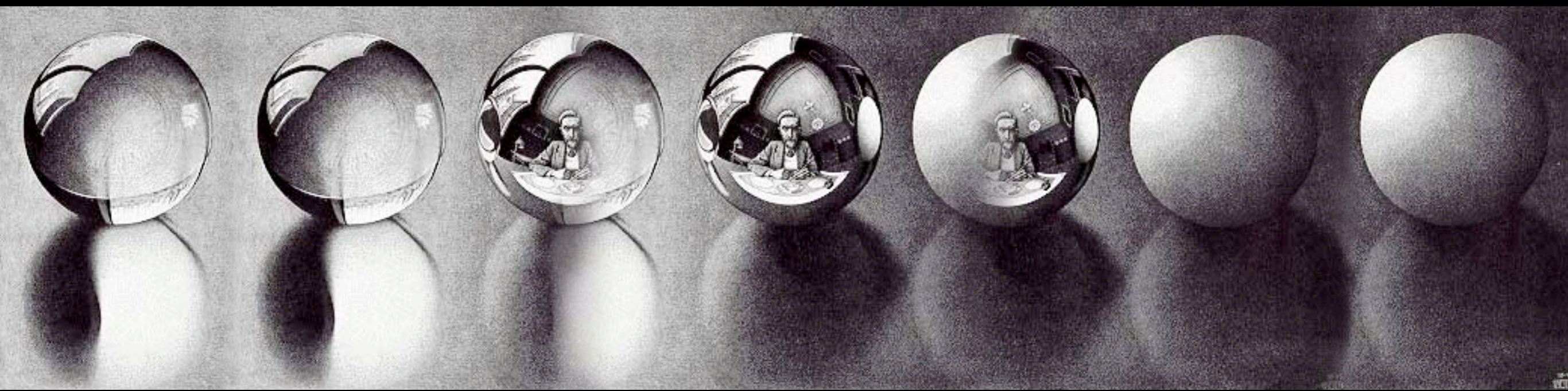
C:

```
#include <stdio.h>
#include <mpi.h>

main (int argc, char *argv[])
{
 int err;
 err = MPI_Init(&argc, &argv);
 printf("Hello world!\n");
 err = MPI_Finalize();
}
```

# Try to think parallel

- ▶ on `mpirun -np N` each of the  $N$  processes runs one copy of your code
- ▶ each variable is duplicated  $N$  times and may have different values on the different processes
- ▶ if you want to check the value of variables with `printf`, always output the rank with each print statement!



# Finding out who I am

```
#include <stdio.h>
#include <mpi.h>

int gmx_setup(int *nnodes)
{
 int resultlen;
 int mpi_num_nodes;
 int mpi_my_rank;
 char mpi_hostname[MPI_MAX_PROCESSOR_NAME];

 MPI_Comm_size(MPI_COMM_WORLD, &mpi_num_nodes);
 MPI_Comm_rank(MPI_COMM_WORLD, &mpi_my_rank);
 MPI_Get_processor_name(mpi_hostname, &resultlen);

 fprintf(stderr, "NNODES=%d, MYRANK=%d, HOSTNAME=%s\n",
 mpi_num_nodes, mpi_my_rank, mpi_hostname);

 *nnodes=mpi_num_nodes;

 return mpi_my_rank;
}
```

```
int main(int argc, char *argv[])
{
 int nnodes, nodeid;

 MPI_Init(&argc, &argv);
 nodeid = gmx_setup(&nnodes);
 MPI_Finalize();
 return 0;
}
```

```
> mpirun -np 6 ./gmxsetup.x
NNODES=6, MYRANK=2, HOSTNAME=node02
NNODES=6, MYRANK=0, HOSTNAME=node01
NNODES=6, MYRANK=1, HOSTNAME=node01
NNODES=6, MYRANK=4, HOSTNAME=node12
NNODES=6, MYRANK=3, HOSTNAME=node02
NNODES=6, MYRANK=5, HOSTNAME=node12
>
```

# Sending and receiving data

- ▶ int MPI\_Send(void \*buf, int count, MPI\_Datatype dtype, int dest, int tag, MPI\_Comm comm);
  - ▶ **body** – read `count` elements of datatype `dtype` from memory address `buf`.
  - ▶ **envelope:** – send this message to the process with rank `dest` in `comm` and label it `tag`.
  - ▶ **error code**
- ▶ int MPI\_Recv(void \*buf, int count, MPI\_Datatype dtype, int source, int tag, MPI\_Comm comm, MPI\_Status \*status);
  - ▶ **body** – write `count` elements of datatype `dtype` to memory address `buf`.
  - ▶ **envelope** – only accept a message tagged `tag` from process with rank `source` in communicator `comm`.
  - ▶ **error code**

# Sending and receiving data

```
/* simple send and receive */
#include <stdio.h>
#include <mpi.h>

main (int argc, char **argv)
{
 int myrank;
 MPI_Status status;
 double a[100];

 MPI_Init(&argc, &argv); /* Initialize MPI */
 MPI_Comm_rank(MPI_COMM_WORLD, &myrank); /* Get rank */

 if(myrank == 0) /* Send a message */
 MPI_Send(a, 100, MPI_DOUBLE, 1, 17, MPI_COMM_WORLD);
 else if(myrank == 1) /* Receive a message */
 MPI_Recv(a, 100, MPI_DOUBLE, 0, 17, MPI_COMM_WORLD,
 &status);

 MPI_Finalize(); /* Terminate MPI */
}
```

# MPI datatypes

| <b>MPI datatype</b> | <b>C type</b>      |
|---------------------|--------------------|
| MPI_CHAR            | signed char        |
| MPI_SHORT           | signed short int   |
| MPI_INT             | signed int         |
| MPI_LONG            | signed long int    |
| MPI_UNSIGNED_CHAR   | unsigned char      |
| MPI_UNSIGNED_SHORT  | unsigned short int |
| MPI_UNSIGNED        | unsigned int       |
| MPI_UNSIGNED_LONG   | unsigned long int  |
| MPI_FLOAT           | float              |
| MPI_DOUBLE          | double             |
| MPI_LONG_DOUBLE     | long double        |
| MPI_BYTE            | (none)             |



# Sending stuff to everyone

$p_0 \rightarrow p_1 \rightarrow p_2 \rightarrow \dots \rightarrow p_{(N-2)} \rightarrow p_{(N-1)}$

```
int main(int argc, char *argv[])
{
 int nnodes, nodeid, right, left;
 float test=0.0;
 MPI_Status status;

 MPI_Init(&argc, &argv);
 nodeid = gmx_setup(&nnodes);
 gmx_left_right(nnodes, nodeid, &left, &right);

 if (nodeid == 0)
 test = 3.1415;

 if (nodeid != 0)
 MPI_Recv(&test, 1, MPI_FLOAT, left, 0, MPI_COMM_WORLD,
 &status);
 if (nodeid != (nnodes-1))
 MPI_Send(&test, 1, MPI_FLOAT, right, 0, MPI_COMM_WORLD);

 MPI_Finalize();
}
```

```
void
gmx_left_right(int nnodes, int nodeid,
int *left, int *right)
{
 *left = (nnodes+nodeid-1) % nnodes;
 *right = (nodeid +1) % nnodes;
}
```

# Printf output from all procs

```
ckutzne@coral4:~/mpi-techtea> mpirun -np 5 ./gmxsetup.x
NNODES=5, MYRANK=0, HOSTNAME=coral4
NNODES=5, MYRANK=1, HOSTNAME=coral4
Nodeid: 0, value of test 3.141500
Nodeid: 0, value of test 3.141500
Nodeid: 1, value of test 0.000000
Nodeid: 1, value of test 3.141500
NNODES=5, MYRANK=2, HOSTNAME=coral4
NNODES=5, MYRANK=3, HOSTNAME=coral4
NNODES=5, MYRANK=4, HOSTNAME=coral4
Nodeid: 2, value of test 0.000000
Nodeid: 2, value of test 3.141500
Nodeid: 3, value of test 0.000000
Nodeid: 3, value of test 3.141500
Nodeid: 4, value of test 0.000000
Nodeid: 4, value of test 3.141500
```

```
int main(int argc, char *argv[])
{
 int nnodes, nodeid, right, left;
 float test=0.0;
 MPI_Status status;

 MPI_Init(&argc, &argv);
 nodeid = gmx_setup(&nnodes);
 gmx_left_right(nnodes, nodeid, &left, &right);

 if (nodeid == 0)
 test = 3.1415;

 printf("Nodeid: %d, value of test: %d\n", nodeid, test);

 if (nodeid != 0)
 MPI_Recv(&test, 1, MPI_FLOAT, left, 0, MPI_COMM_WORLD,
 &status);
 if (nodeid != (nnodes-1))
 MPI_Send(&test, 1, MPI_FLOAT, right, 0, MPI_COMM_WORLD);

 printf("Nodeid: %d, value of test: %d\n", nodeid, test);
 MPI_Finalize();
}
```

# Printf output from all procs

```
ckutzne@coral4:~/mpi-techtea> mpirun -np 5 ./gmxsetup.x
```

```
NNODES=5, MYRANK=0, HOSTNAME=coral4
NNODES=5, MYRANK=2, HOSTNAME=coral4
NNODES=5, MYRANK=1, HOSTNAME=coral4
NNODES=5, MYRANK=3, HOSTNAME=coral4
NNODES=5, MYRANK=4, HOSTNAME=coral4
Nodeid: 0, value of test 3.141500
Nodeid: 2, value of test 0.000000
Nodeid: 1, value of test 0.000000
Nodeid: 4, value of test 0.000000
Nodeid: 3, value of test 0.000000
Nodeid: 0, value of test 3.141500
Nodeid: 2, value of test 3.141500
Nodeid: 1, value of test 3.141500
Nodeid: 3, value of test 3.141500
Nodeid: 4, value of test 3.141500
```

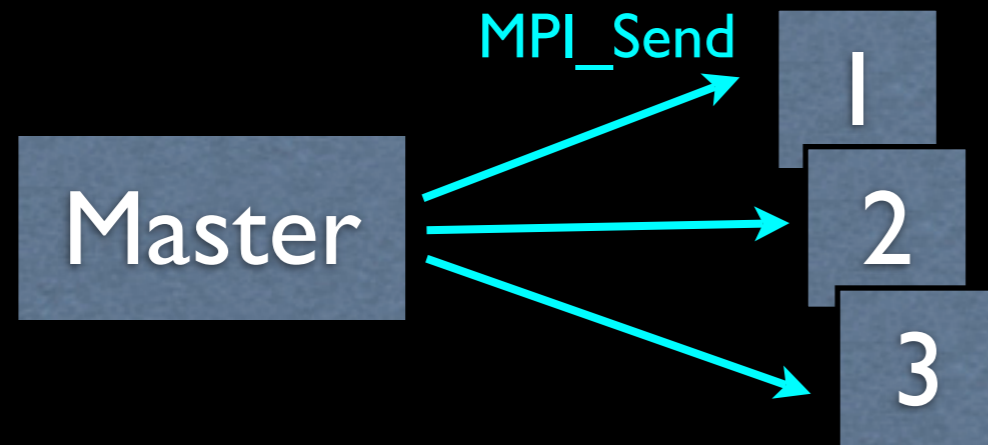
```
int main(int argc, char *argv[])
{
 int nnodes, nodeid, right, left;
 float test=0.0;
 MPI_Status status;

 MPI_Init(&argc, &argv);
 nodeid = gmx_setup(&nnodes);
 gmx_left_right(nnodes, nodeid, &left, &right);

 if (nodeid == 0)
 test = 3.1415;
 MPI_Barrier(MPI_COMM_WORLD);
 printf("Nodeid: %d, value of test: %d\n", nodeid, test);

 if (nodeid != 0)
 MPI_Recv(&test, 1, MPI_FLOAT, left, 0, MPI_COMM_WORLD,
 &status);
 if (nodeid != (nnodes-1))
 MPI_Send(&test, 1, MPI_FLOAT, right, 0, MPI_COMM_WORLD);
 MPI_Barrier(MPI_COMM_WORLD);
 printf("Nodeid: %d, value of test: %d\n", nodeid, test);
 MPI_Finalize();
}
```

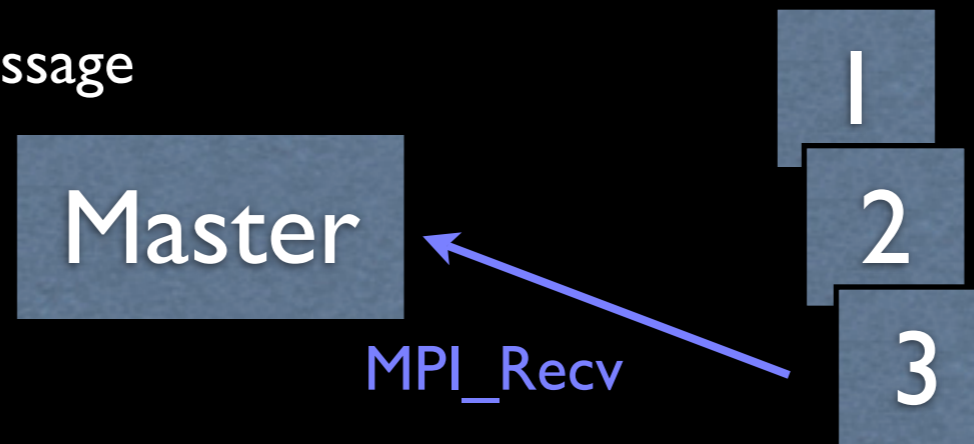
# Master-slave scheme



```
MPI_Recv(data, ..., MPI_ANY_TAG, ..., status)
if (status->MPI_TAG == 1)
{
 do_work;
 MPI_Send(results back to master);
}
```

Master waits for any message

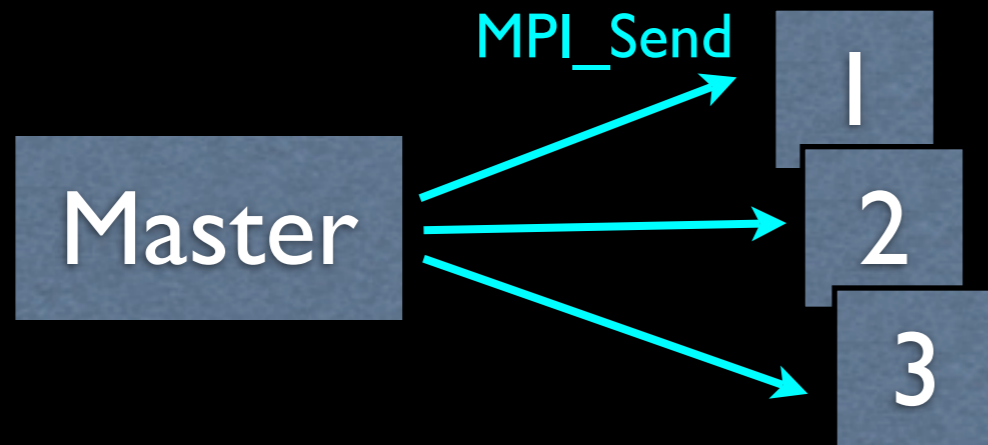
```
MPI_Recv(...,
MPI_ANY_SOURCE,
MPI_ANY_TAG, ...,
status)
```



accumulates the data ...

- ▶ **status** needed when using wildcards
  - ▶ status -> MPI\_SOURCE
  - ▶ status -> MPI\_TAG

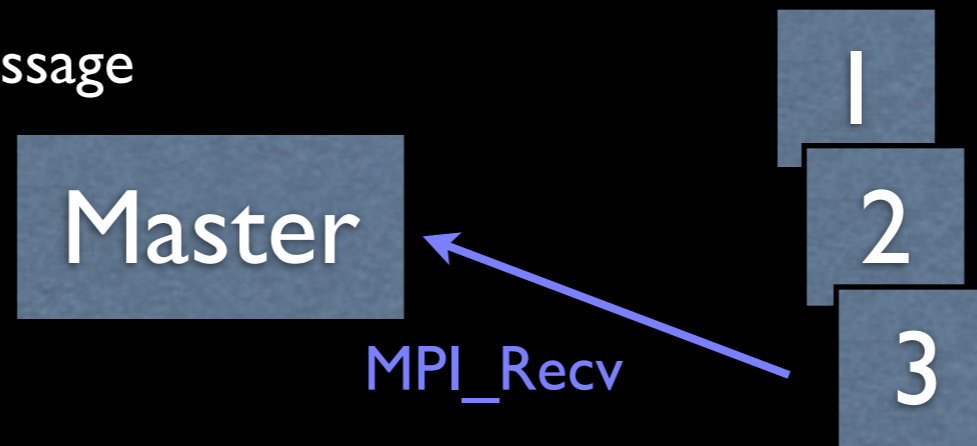
# Master-slave scheme



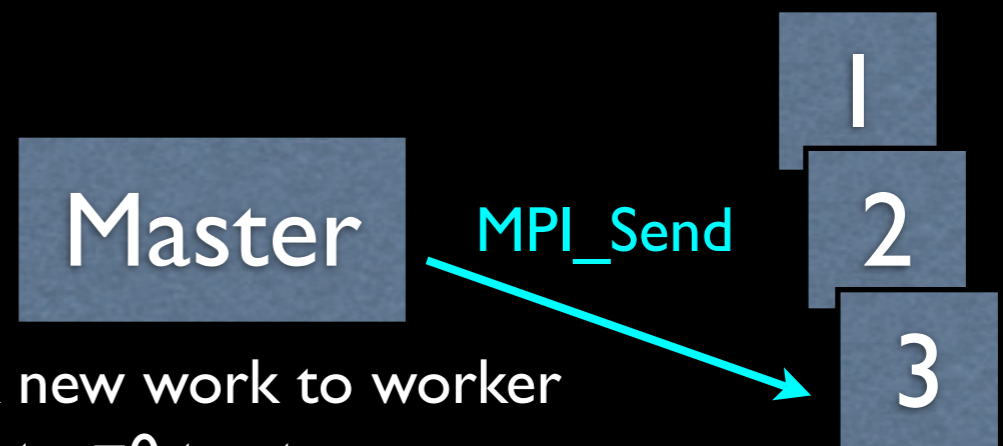
```
MPI_Recv(data, ..., MPI_ANY_TAG, ..., status)
if (status->MPI_TAG == 1)
{
 do_work;
 MPI_Send(results back to master);
}
```

Master waits for any message

```
MPI_Recv(...,
MPI_ANY_SOURCE,
MPI_ANY_TAG, ...,
status)
```



accumulates the data ...



... and sends back new work to worker  
or a message with tag=0 to stop

# Pitfall: deadlock

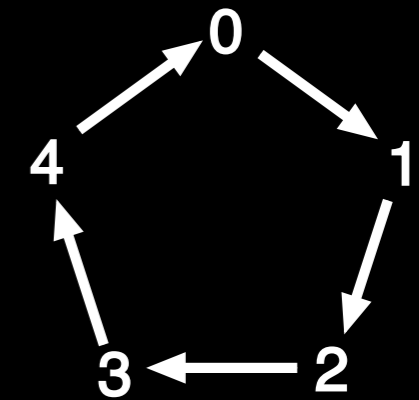
- ▶ MPI\_Send and MPI\_Recv block, so deadlock may occur

| Time | Process 0       | Process 1       |
|------|-----------------|-----------------|
| 1    | MPI_Send to 1   | MPI_Send to 0   |
| 2    | MPI_Recv from 1 | MPI_Recv from 0 |

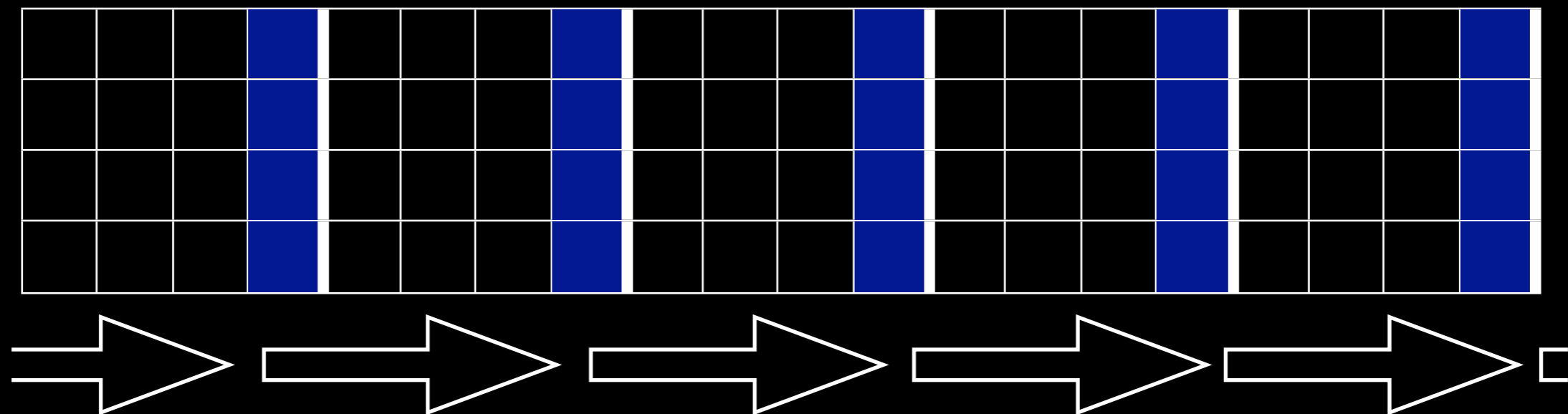
- ▶ if the messages cannot be buffered, the order becomes relevant!

| Time | Process 0       | Process 1       |
|------|-----------------|-----------------|
| 1    | MPI_Send to 1   | MPI_Recv from 0 |
| 2    | MPI_Recv from 1 | MPI_Send from 0 |

# Boundary exchange



- ▶ Send and receive with a single call
  - ▶ `int MPI_Sendrecv(void *sbuf, int scount, MPI_Datatype stype, int dest, int stag, void *rbuf, int rcount, MPI_Datatype rtype, int source, int rtag, MPI_Comm comm, MPI_Status *status);`
  - ▶ MPI cares that no deadlock occurs!
  - ▶ Send to right neighbor, receive from left neighbor



# Broadcasting data

- ▶ `int MPI_Bcast(void *buf, int count, MPI_Datatype dtype, int rank, MPI_Comm comm);`
  - ▶ Broadcasts data from process `rank` to all other processes in `comm`
  - ▶ **error code**
- ▶ Pitfall: Matching broadcast with receive

```
#include <mpi.h>
main(int argc, char *argv[])
{
 int rank;
 double param;

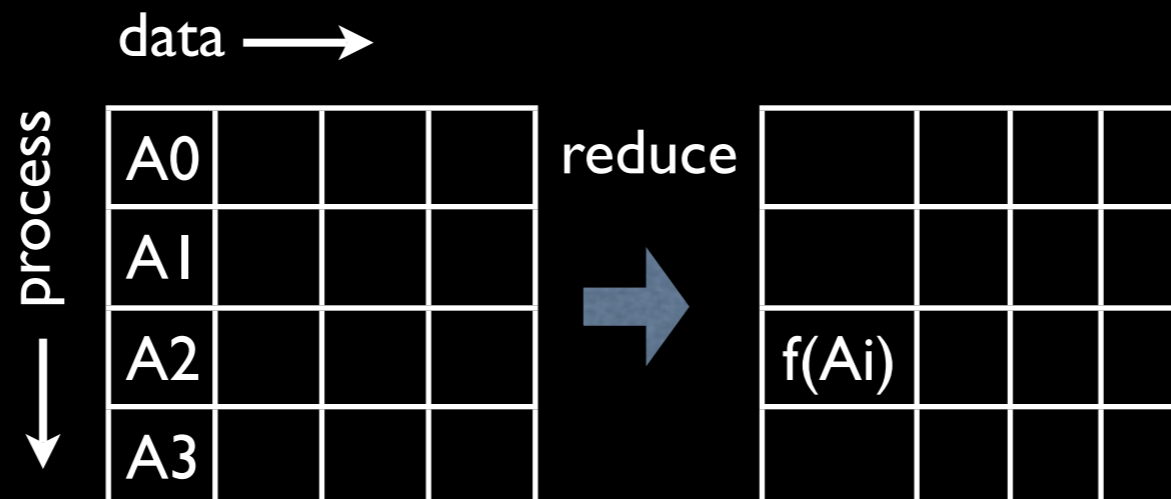
 MPI_Init(&argc, &argv);
 MPI_Comm_rank(MPI_COMM_WORLD, &rank);
 if(rank==5) param=23.0;
 MPI_Bcast(¶m, 1, MPI_DOUBLE, 5, MPI_COMM_WORLD);
 printf("P:%d after broadcast param is %f \n", rank, param);
 MPI_Finalize();
}
```

```
P: 0 after broadcast param is 23.
P: 5 after broadcast param is 23.
P: 2 after broadcast param is 23.
P: 3 after broadcast param is 23.
P: 4 after broadcast param is 23.
P: 1 after broadcast param is 23.
P: 6 after broadcast param is 23.
```



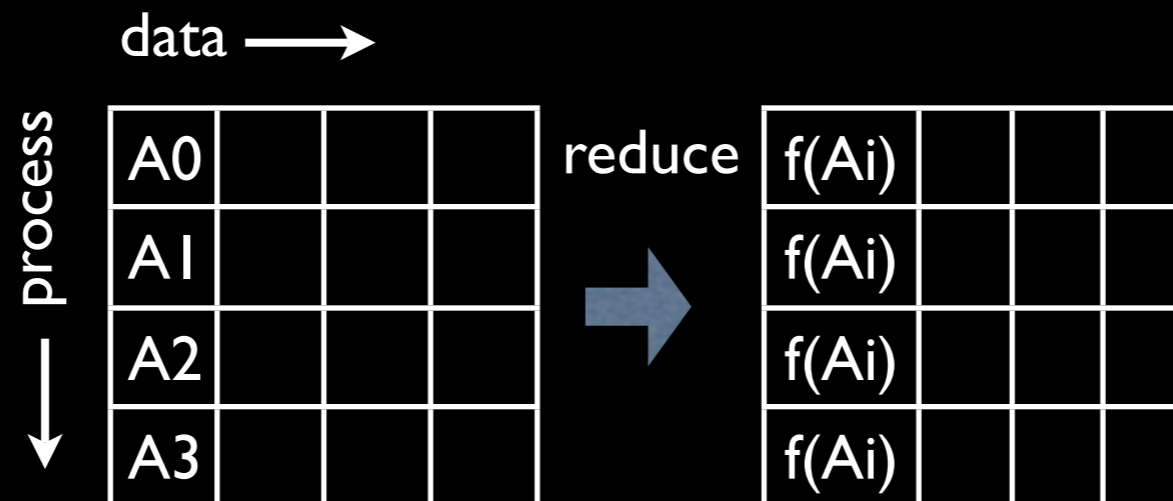
# Combine data from all procs

- ▶ int MPI\_Reduce(void \*operand, void \*result, int count, MPI\_Datatype dtype, MPI\_Op op, int root, MPI\_Comm comm);
  - ▶ Combines the operands using the operator op and stores the result in result on process root
  - ▶ error code



# Combine data from all procs

- ▶ int MPI\_Allreduce(void \*operand, void \*result, int count, MPI\_Datatype dtype, MPI\_Op op, MPI\_Comm comm);
  - ▶ Combines the operands using the operator op and stores the result in result on all processes
  - ▶ error code

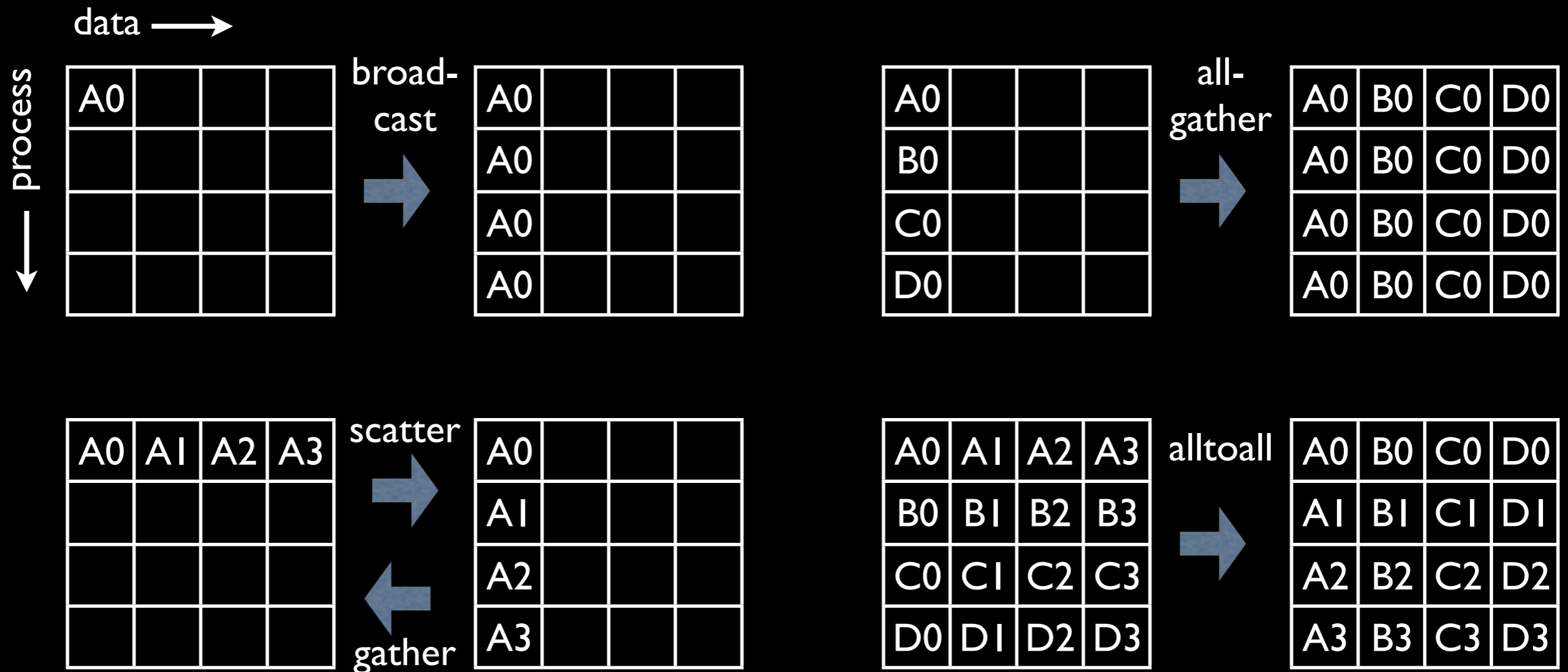


- ▶ Pitfall: Aliasing
  - ▶ attempt to store the result in the same location as the operand
  - ▶ MPI\_Reduce(&operand, &operand, 1, MPI\_FLOAT, MPI\_SUM, 0, comm)

# Predefined reduction ops

| <b>Operator name</b> | <b>Meaning</b>                |
|----------------------|-------------------------------|
| MPI_MAX              | Maximum                       |
| MPI_MIN              | Minimum                       |
| MPI_SUM              | Sum                           |
| MPI_PROD             | Product                       |
| MPI_LAND             | Logical AND                   |
| MPI_BAND             | Bitwise AND                   |
| MPI_LOR              | Logical OR                    |
| MPI_BOR              | Bitwise OR                    |
| MPI_LXOR             | Logical EXCLUSIVE OR          |
| MPI_BXOR             | Bitwise EXCLUSIVE OR          |
| MPI_MAXLOC           | Maximum & location of maximum |
| MPI_MINLOC           | Minimum & location of minimum |

# More collectives



- ▶ **Pitfall: Calling collectives from only some of the processes**

```
if (...) {
 do_something();
} else {
 do_something_else();
 MPI_Barrier(comm);
}
```



# Design considerations

- ▶ “Quick & dirty”
  - ▶ compile your program in parallel, execute with mpicc, decide on the rank, which process operates on which part of the data
  - ▶ advantage & disadvantage: input data is read by every process
  - ▶ results can be written one after another into single output file or to separate files
- ▶ “Master” process reads data and broadcasts to everyone
  - ▶ slightly more work, but faster on more processors
  - ▶ results can be send back to the “master”, who then writes to file
- ▶ Master-slave parallelism  $\Leftrightarrow$  symmetric program?
- ▶ Make performance analysis a part of your development cycle!

# Debugging parallel programs

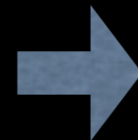
- ▶ serial style will not work!

- ▶ mpicc -g ./buggy.c -o buggy.x

- ▶ ddd mpirun -np 2 buggy.x will try to debug mpirun, not buggy.x

- ▶ what does the trick:

- ▶ keep the program from running away with while statement & barriers



- ▶ run the program:

- mpirun -np 2 buggy.x

- ▶ find out process numbers:

- > ps -C buggy.x

- 3209 buggy.x – MPI rank 0

- 3210 buggy.x

- ▶ attach debugger to process of interest:

- ddd buggy.x 3209 &

```
MPI_Init(&argc, &argv);
MPI_Comm_rank(&nodeid, MPI_COMM_WORLD);

/* wait */
if (nodeid == 0)
{
 int bDebugWait = 1;

 fprintf(stderr, "Wait on proc %d\n",
 nodeid);
 while (bDebugWait)
 ;
}
MPI_Barrier(MPI_COMM_WORLD);
/* end wait */
```

DDD: /netmount/cora4/ckutzne/mpi-techtea/buggy.c

File Edit View Program Commands Status Source Data Help

0: while

Lookup Find>> Break Watch Print Display Plot Hide Rotate Set Undisp

|           |               |                |
|-----------|---------------|----------------|
| 1: nodeid | 2: bDebugWait | 3: status      |
| 1         | 0             | MPI_SOURCE = 0 |
|           |               | MPI_TAG = 0    |
| 4: nnodes |               | MPI_ERROR = 0  |
| 2         |               | st_length = 4  |

```
int main(int argc, char *argv[])
{
 int nnodes, nodeid, rightnode, leftnode;
 float test=0.0;
 MPI_Status status;

 MPI_Init(&argc, &argv);
 nodeid = gmx_setup(&nnodes);
 gmx_left_right(nnodes, nodeid, &leftnode, &rightnode);

 /* wait */
 int bDebugWait = 1;

 if (nodeid == 1)
 {
 fprintf(stderr, "Wait on proc %d —\n", nodeid);
 while (bDebugWait)
 ;
 }

 MPI_Barrier(MPI_COMM_WORLD);
 /* end wait */

 if (nodeid == 0) test = 3.1415;

 MPI_Barrier(MPI_COMM_WORLD);
 printf("Nodeid: %d, value of test %f\n", nodeid, test);

 if (nodeid != 0)
 MPI_Recv(&test, 1, MPI_FLOAT, leftnode, 0, MPI_COMM_WORLD, &status);

 if (nodeid != (nnodes-1))
 MPI_Send(&test, 1, MPI_FLOAT, rightnode, 0, MPI_COMM_WORLD);

 MPI_Barrier(MPI_COMM_WORLD);
 printf("Nodeid: %d, value of test %f\n", nodeid, test);
}
```

DDD

Run

Interrupt

Step Stepi

Next Nexti

Until Finish

Cont Kill

Up Down

Undo Redo

Edit Make

# Performance evaluation

- ▶ Simple time measurements with MPI\_Wtime:
  - ▶ MPI\_Wtick() – timer accuracy
  - ▶ How long does something take on the individual processors? (N timings)  
`t1 = MPI_Wtime();`  
... code to time ...  
`elapsed = MPI_Wtime() - t1;`
  - ▶ How long does it take on all parallel processes? (single timing)  
`MPI_Barrier();`  
`t1 = MPI_Wtime();`  
... code to time ...  
`MPI_Barrier();`  
`elapsed = MPI_Wtime() - t1;`
  - ▶ may be worth to look at the synchronization time!  
(How long does the barrier take?) (N timings, take maximum)  
`t2 = MPI_Wtime();`  
`MPI_Barrier();`  
`synctime = MPI_Wtime() - t2;`





# Multi Processing Environment

## ▶ Instrument your program

- ▶ `include <mpe.h>`
- ▶ `MPE_Init_log();`
- ▶ `MPE_Stop_log(); MPE_Start_log();`
- ▶ `MPE_Describe_state(int start, int end, char *name, char *color);`
- ▶ `MPE_Describe_event(int event, char *name);`
- ▶ `MPE_Log_event(int event, int intdata, char *chardata);`

Visualize & measure what is going on in your program

## ▶ Compile with mpecc

- ▶ `mpecc -mpilog mpedemo.c -o mpedemo.x`

## ▶ run

- ▶ `mpirun -np 6 mpedemo.x` generates `mpedemo.x.clog2`

## ▶ view & analyze with jumpshot

# MPE logging

```
int main(int argc, char *argv[])
{
 int nnodes, nodeid, rightnode, leftnode;
 float test;
 MPI_Status status;

 MPI_Init(&argc, &argv);
 MPI_Comm_size(MPI_COMM_WORLD, &nnodes);
 MPI_Comm_rank(MPI_COMM_WORLD, &nodeid);

 /* declare MPE variables & states and initialize MPE */
 int ev_print_start, ev_print_finish;
 ev_print_start = MPE_Log_get_event_number();
 ev_print_finish = MPE_Log_get_event_number();
 MPE_Describe_state(ev_print_start, ev_print_finish, "output", "orange");
 MPE_Init_log();
 /* done MPE setup */

 gmx_left_right(nnodes, nodeid, &leftnode, &rightnode);
 if (nodeid == 0) test = 3.1415;

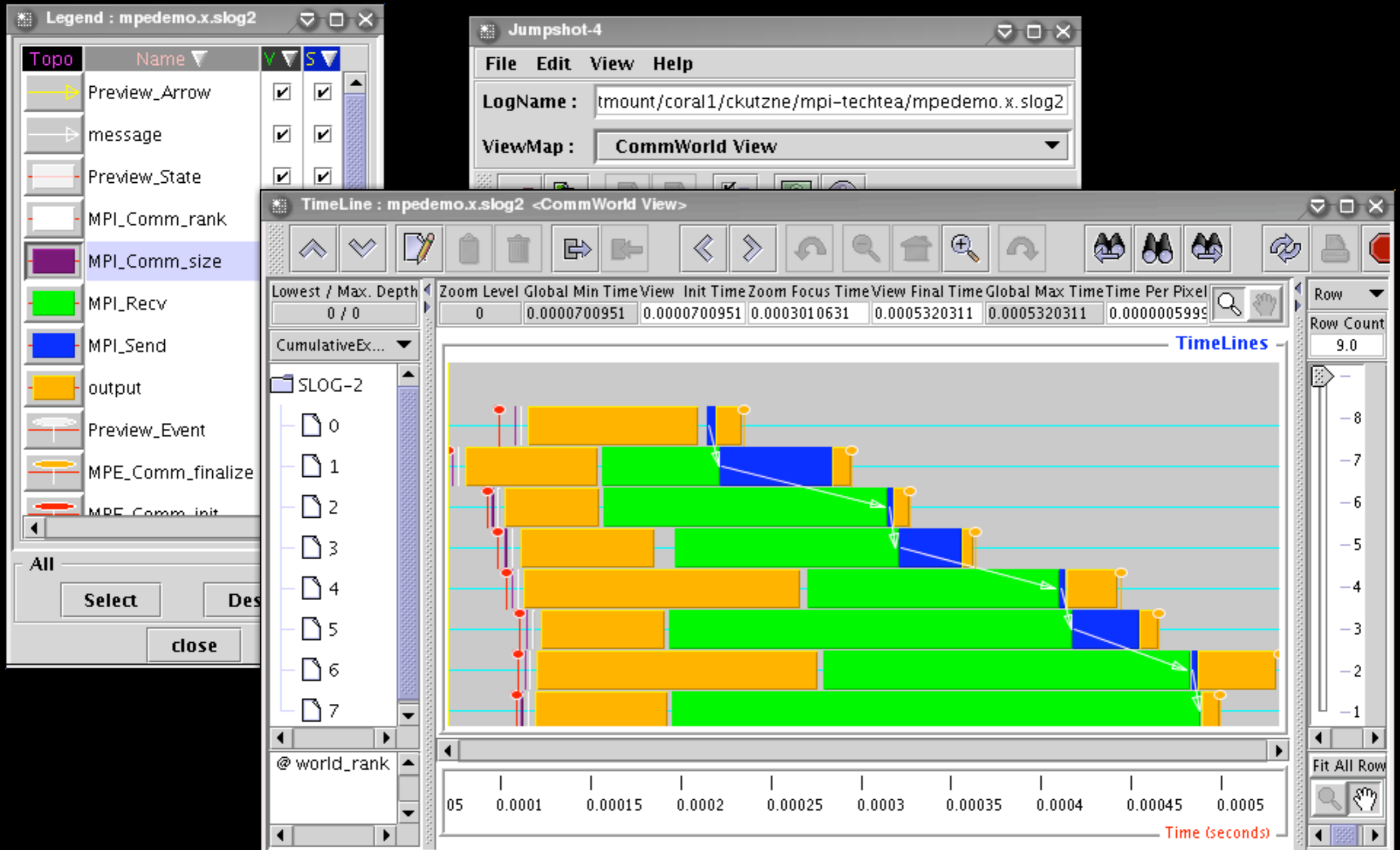
 MPE_Log_event(ev_print_start, 0, "start output");
 printf("Nodeid: %d, value of test %f\n", nodeid, test);
 MPE_Log_event(ev_print_finish, 0, "output finished");

 if (nodeid != 0)
 MPI_Recv(&test, 1, MPI_FLOAT, leftnode, 0, MPI_COMM_WORLD, &status);
 if (nodeid != (nnodes-1))
 MPI_Send(&test, 1, MPI_FLOAT, rightnode, 0, MPI_COMM_WORLD);

 MPE_Log_event(ev_print_start, 0, "start output");
 printf("Nodeid: %d, value of test %f\n", nodeid, test);
 MPE_Log_event(ev_print_finish, 0, "output finished");

 MPI_Finalize();
}
```

# MPE logging



# MPE graphics library

- ▶ **simple graphics interface for parallel output**
  - ▶ shared access by parallel processes to a single X display
  - ▶ each process can individually update the display
- ▶ **usage:**
  - ▶ include `<mpe_graphics.h>`
  - ▶ `MPE_Open/Close_graphics(...)`
  - ▶ `MPE_Draw_point/line/circle/rectangle(...)`
  - ▶ `MPE_Update(...)`
  - ▶ `MPE_Num_colors(...)`
  - ▶ `MPE_Make_color_array(...)`
- ▶ **compile with mpecc**
  - ▶ `mpecc -graphics -lm mpedemo2.c -o mpedemo2.x`

# MPE graphics library

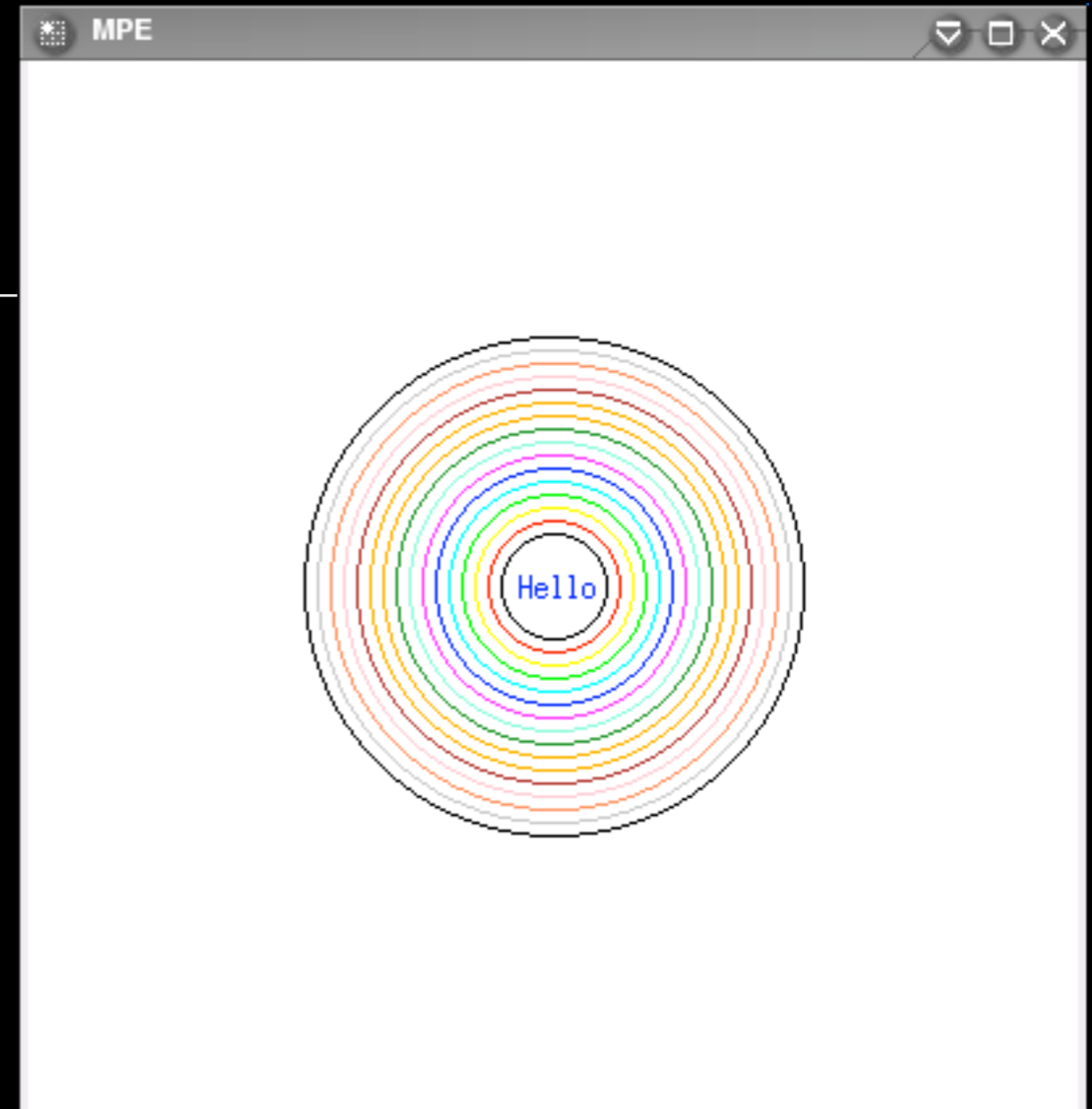
```
#include <stdio.h>
#include <mpi.h>
#include "mpe.h"
#include "mpe_graphics.h"

int main(int argc, char *argv[])
{
 int nnodes, nodeid;
 MPE_XGraph graph;
 char ckey;
 int ierr;
 MPE_Color my_color;

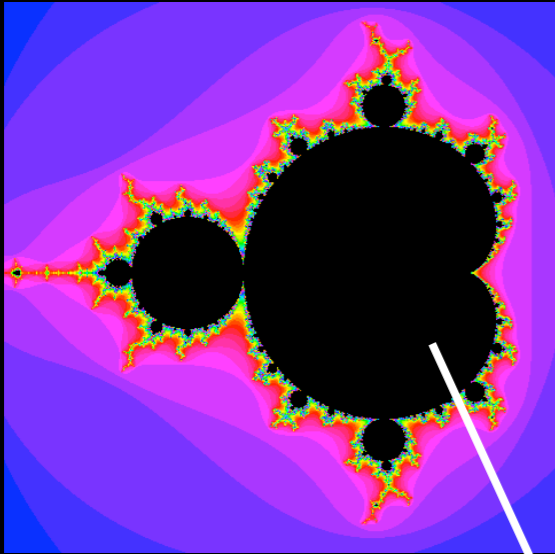
 MPI_Init(&argc, &argv);
 MPI_Comm_size(MPI_COMM_WORLD, &nnodes);
 MPI_Comm_rank(MPI_COMM_WORLD, &nodeid);

 MPE_Open_graphics(&graph, MPI_COMM_WORLD, NULL,
 -1, -1, 400, 400, 0);
 my_color = (MPE_Color) (nodeid + 1);
 if (nodeid == 0)
 ierr = MPE_Draw_string(graph, 187, 205, MPE_BLUE, "Hello");
 ierr = MPE_Draw_circle(graph, 200, 200, 20+nodeid*5, my_color);
 ierr = MPE_Update(graph);

 MPI_Barrier(MPI_COMM_WORLD);
 ierr = MPE_Close_graphics(&graph);
 MPI_Finalize();
}
```



# Self-scheduling master-slave example



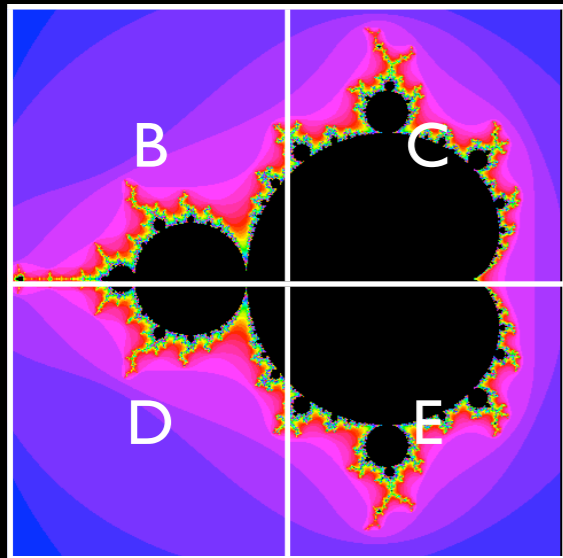
M

- ▶  $z = x + iy$
- ▶ repeat

$$f_c(z) = z^2 + c$$

If  $f$  remains bounded, it is in M,  
otherwise color is a (optional) measure for  
divergence

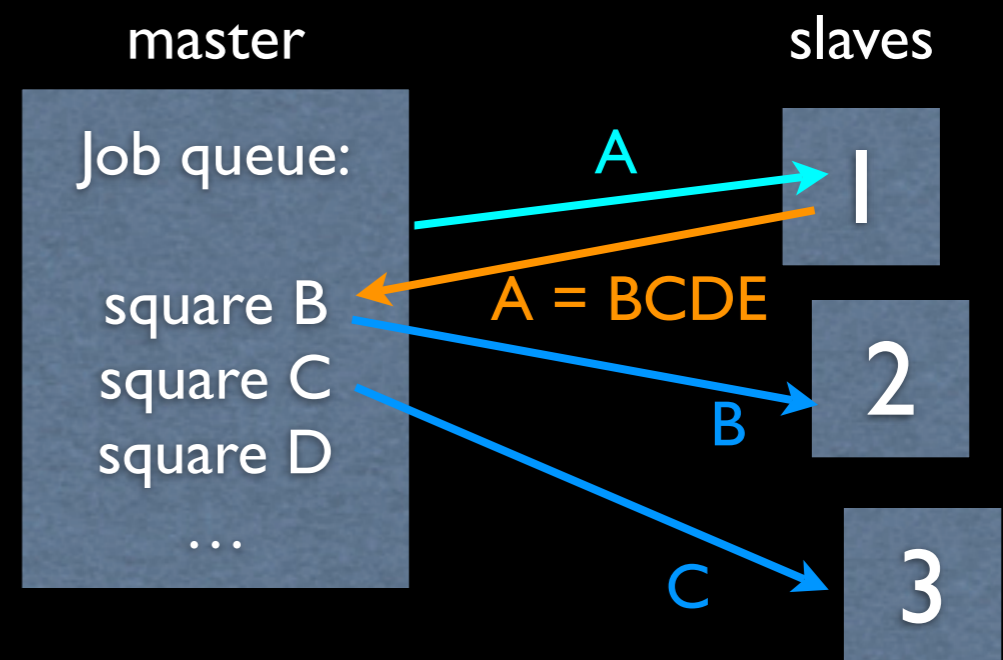
- ▶ “Easy to parallelize”  
since each pixel’s color can be computed  
independently
- ▶ BUT: load imbalances!



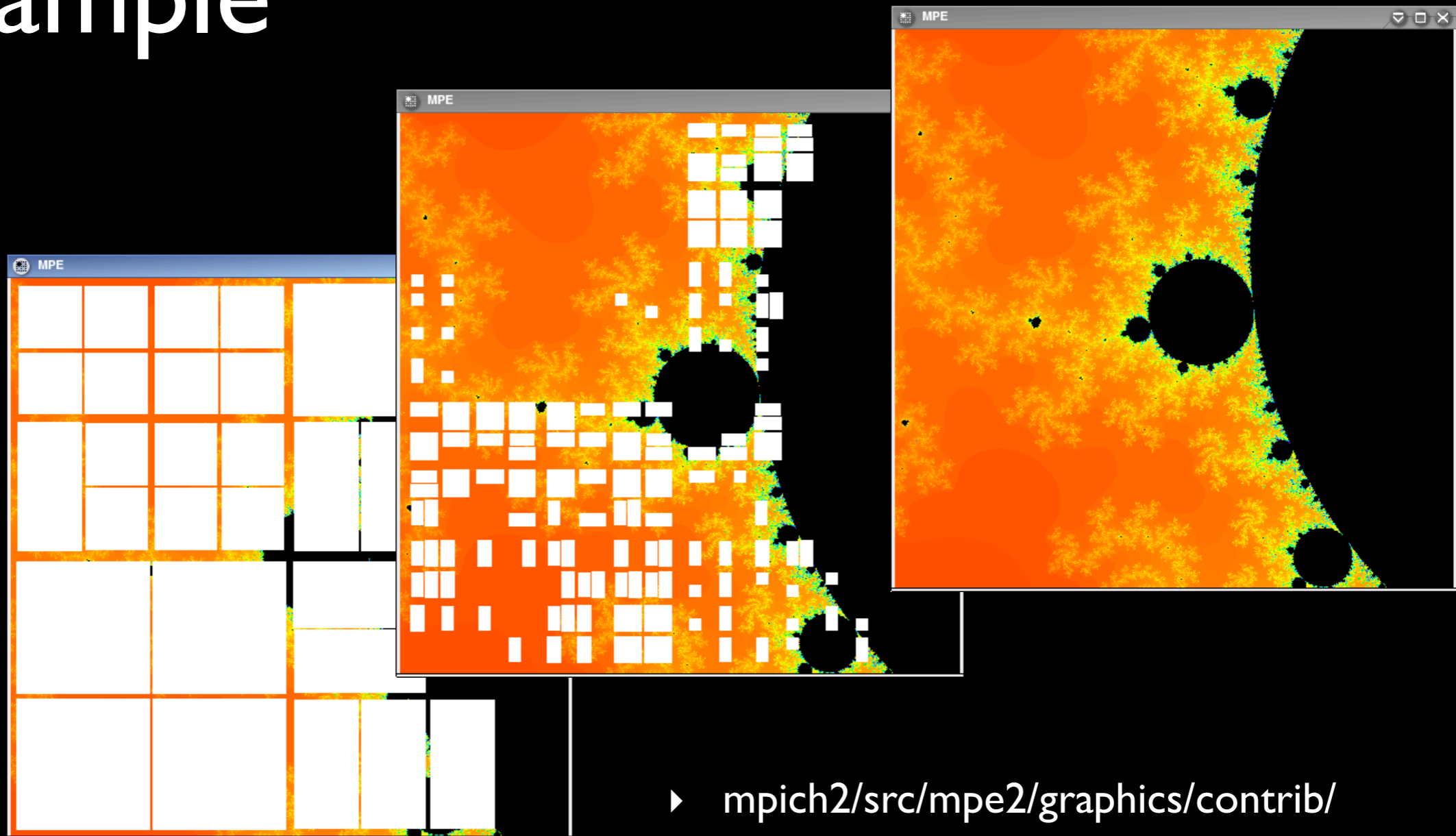
# Self-scheduling master-slave example

## ▶ accelerate computation & balance load:

- ▶ if border of any square is made up of same color, all pixels inside have that color, too
- ▶ start: put whole area (A) into queue
- ▶ master sends queue job to a free slave
  - ▶ slave computes boundary
  - ▶ if same color => fill & display
  - ▶ at first different color => subdivide into 4 subsquares, return them to master, but carry on with boundary. Display boundary when finished.



# Self-scheduling master-slave example



- ▶ `mpich2/src/mpe2/graphics/contrib/`
- ▶ `mpich2-1.0.4/examples`
- ▶ `lam-7.1.2/examples`
- ▶ <http://www-unix.mcs.anl.gov/mpi/usingmpi>