# Introduction to the Parallelization with MPI
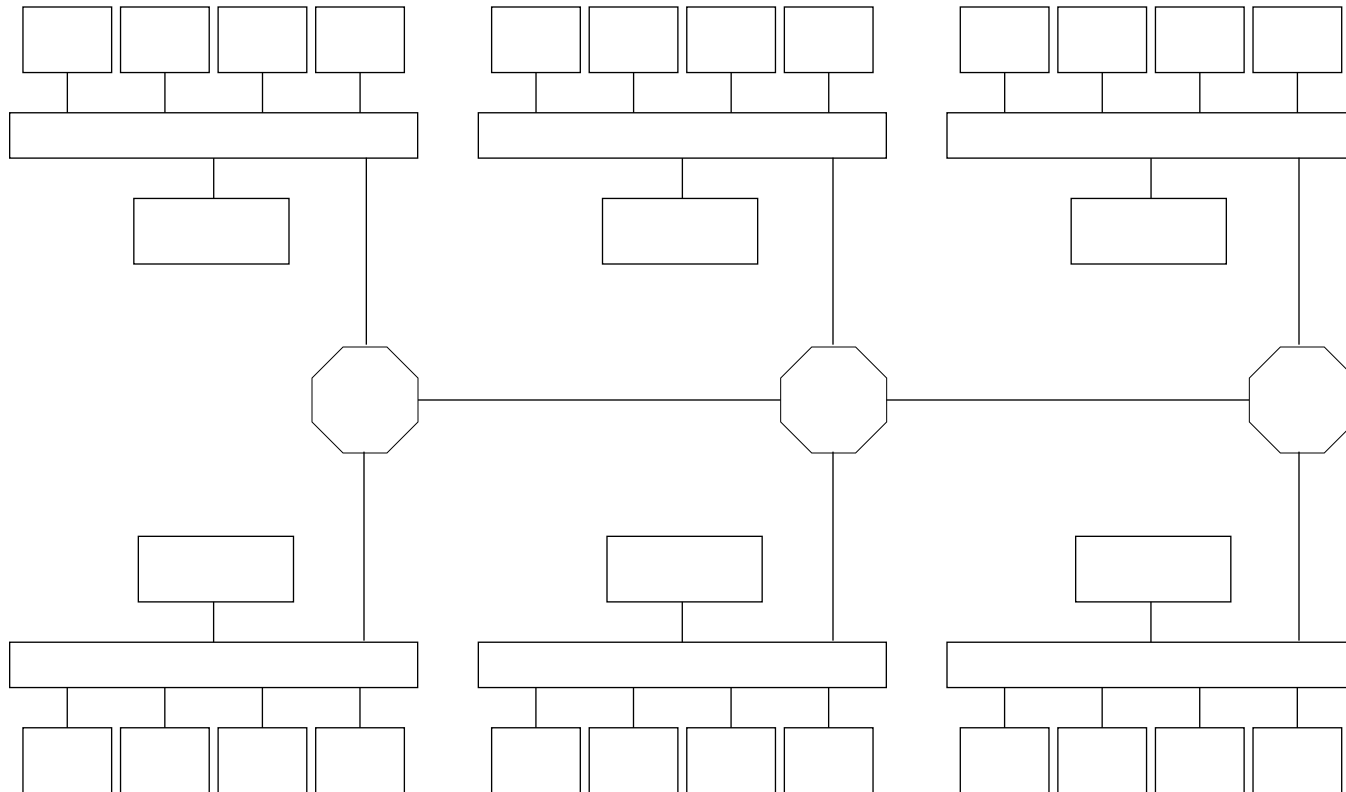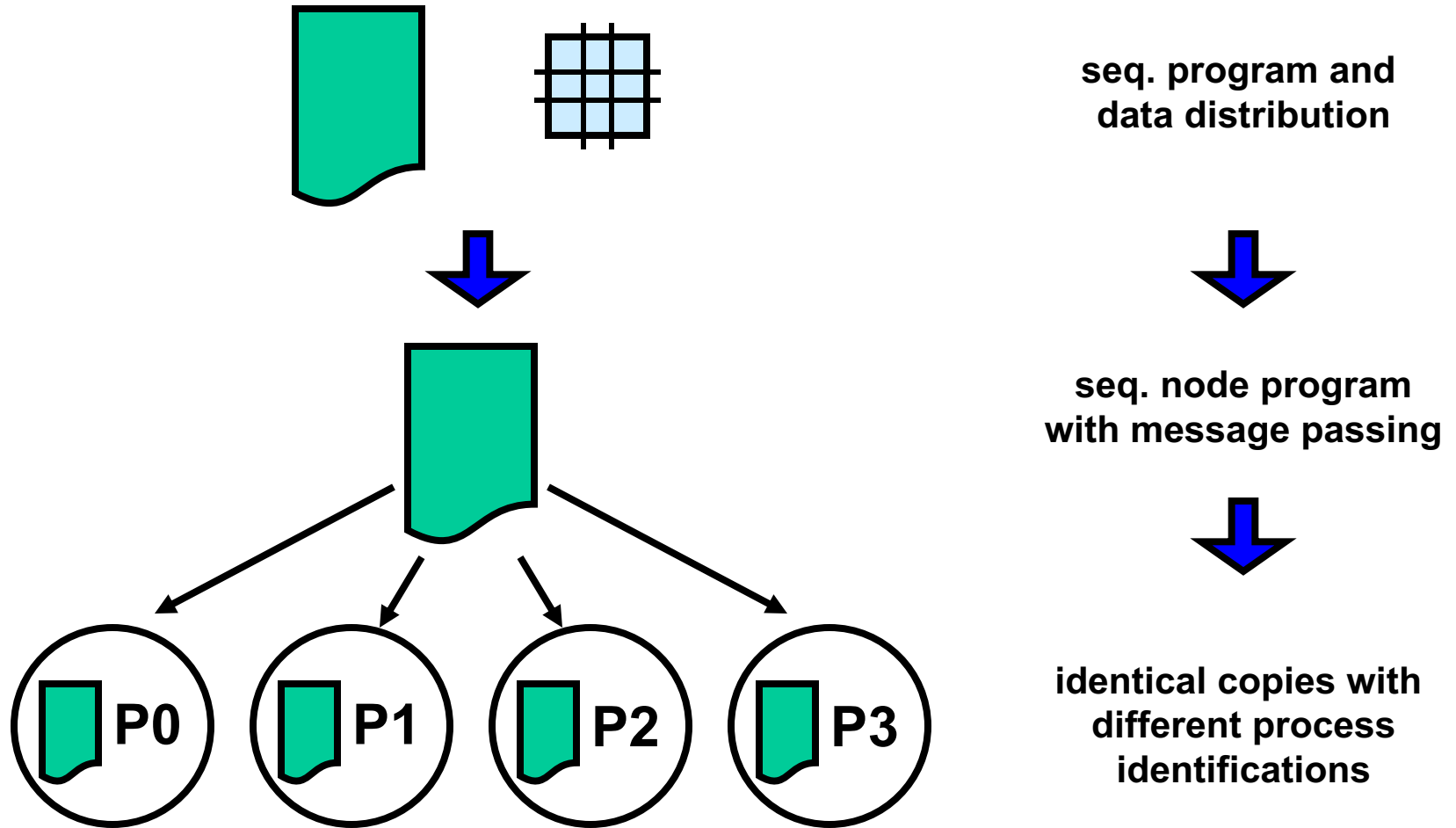
Michael Gerndt

# Distributed Memory Architektur

# Single Program Multiple Data (SPMD)

seq. program and
data distribution

seq. node program
with message passing

P0   P1   P2   P3

identical copies with
different process
identifications

# Program Parallelization

1. **Adaptation of array declarations**
   - Local size of distributed arrays covers only the part of the data structure assigned to the process.

2. **Index transformation**
   - Global indices are mapped into a tupel of node number and of a local index.

3. **Work distribution**
   - Computations are executed by the process owning the assigned variable.

4. **Communication**
   - Accesses to array elements of other processes have to be implemented by message passing.

# Scope of the Message Passing Interface

- MPI 1.2
  - Point-to-Point communication
  - Collective communication
  - Communicators
  - Process topologies
  - User-defined data types
  - Operations and properties of the execution environment
  - Profiling interface
- MPI 2.0
  - Dynamic process creation
  - One-sided communication
  - Parallel IO
- MPI 3.0 (2012)
  - Nonblocking collectives
  - Additional one-sided communication operations
- MPI 4.0 Draft (2020)
- https://www.mpi-forum.org

# Core Routines

- MPI 1.2 has 129 functions
- It is possible to write real programs with only six functions:
  - MPI_Init
  - MPI_Finalize
  - MPI_Comm_size
  - MPI_Comm_rank
  - MPI_Send
  - MPI_Recv

# MPI_Init

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

        IN argc, argv:   *arguments*
        return:         *MPI_SUCCESS or error codes*

- **This routine has to be called by each MPI process before any other MPI routine is executed**

- **Fortran interface**
  - MPI_INIT (integer ierror)
  - The name is written in capital letters and the error code is returned via an additional argument.

# MPI_Finalize

int MPI_Finalize ()

- Each process must call MPI_FINALIZE before it exits.
  - Precondition: All pending communication has to be finished.
  - One MPI_FINALIZE returns, no further MPI routines can be executed.
  - MPI_FINALIZE frees any resources.

# MPI_Comm_size

int MPI_Comm_size (MPI_Comm comm, int *size)

    IN comm:     *Communicator*
    OUT size:    *Cardinality of the process group*

- **Communicator**
  - Identifies a process group and defines the communication context. All message tags are unique with respect to a communicator.

- **MPI_COMM_WORLD**
  - This is a predefined standard communicator. Its process group includes all processes of a parallel application.

- **MPI_Comm_size**
  - It returns the number of processes in the process group of the given communicator.

# MPI_Comm_rank

int MPI_Comm_rank (MPI_Comm comm, int *rank)

       IN comm:    *Communicator*
       OUT rank:  *process number of the executing process*

- Process number
  - The process number is a unique identifier within the process group of the communicator.
  - It is the only way to distinguish processes and to implement an SPMD program.
- MPI_Comm_rank returns the process number of the executing process.

# MPI_Send

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

IN buf:      *Address of the send buffer*
IN count:    *Number of data to be sent*
IN dtype:    *Data type*
IN dest:     *Receiver*
IN tag:      *Message tag*
IN comm:     *Communicator*

- MPI_Send
  - Sends the data to the receiver.
  - It is a blocking operation, i.e. it terminates when the send buffer can be reused, either because the message was delivered or the data were copied to a system buffer.

# MPI Data Types

| C | |
|---|---|
| MPI_CHAR | signed char |
| MPI_SHORT | signed short int |
| MPI_INT | signed int |
| MPI_LONG | signed long int |
| MPI_UNSIGNED_CHAR | |
| MPI_UNSIGNED_INT | |
| ... | |
| MPI_FLOAT | float |
| MPI_DOUBLE | double |
| MPI_LONG_DOUBLE | long double |
| MPI_BYTE | |
| MPI_PACKED | |

| FORTRAN | |
|---|---|
| MPI_INTEGER | integer |
| MPI_REAL | real |
| MPI_DOUBLE_PRECISION | double precision |
| MPI_COMPLEX | complex |
| MPI_LOGICAL | logical |
| MPI_CHARACTER | character(1) |
| MPI_BYTE | |
| MPI_PACKED | |

# MPI_Recv

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

OUT buf:      *Address of the receive buffer*
IN count:     *Size of receive buffer*
IN dtype:     *Data type*
IN source:    *Sender*
IN tag:       *Message tag*
IN comm:      *Communicator*
OUT status:   *Status information*

- Properties:
  - It is a blocking operation, i.e. it terminates after the message is available in the receive buffer.
  - The message must not be larger than the receive buffer.
  - The remaining part of the buffer not used for the received message will be unchanged.

# Properties of MPI_Recv

- **Message selection**
  - A message to be received by this function must match
    - the sender
    - the tag
    - the communicator
  - Sender and tag can be specified as wild cards
    - MPI_ANY_SOURCE and MPI_ANY_TAG
  - There is no wild card for the communicator.

- **Status**
  - The data structure MPI_Status includes
    - status(MPI_SOURCE): sender of the message
    - status(MPI_TAG): message tag
    - status(MPI_ERROR): error code
  - The actual length of the received message can be determined via MPI_Get_count.

# Circular Left Shift Application

**mpirun –np 4 shifts <number of positions>**

### Description

- Position 0 of an array with 100 entries is initialized to 1. The array is distributed among all processes in a blockwise fashion.

- A number of circular left shift operations is executed.

- The number is specified via a command line parameter.

| 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

# Shifts: Initialization

```c
#include "mpi.h"

main (int argc,char *argv[]){
int myid, np, ierr, lnbr, rnbr, shifts, i, j;
int *values;
MPI_Status status;

ierr = MPI_Init (&argc, &argv);
if (ierr != MPI_SUCCESS){
  ...
}

MPI_Comm_size(MPI_COMM_WORLD, &np);
MPI_Comm_rank(MPI_COMM_WORLD, &myid);
```
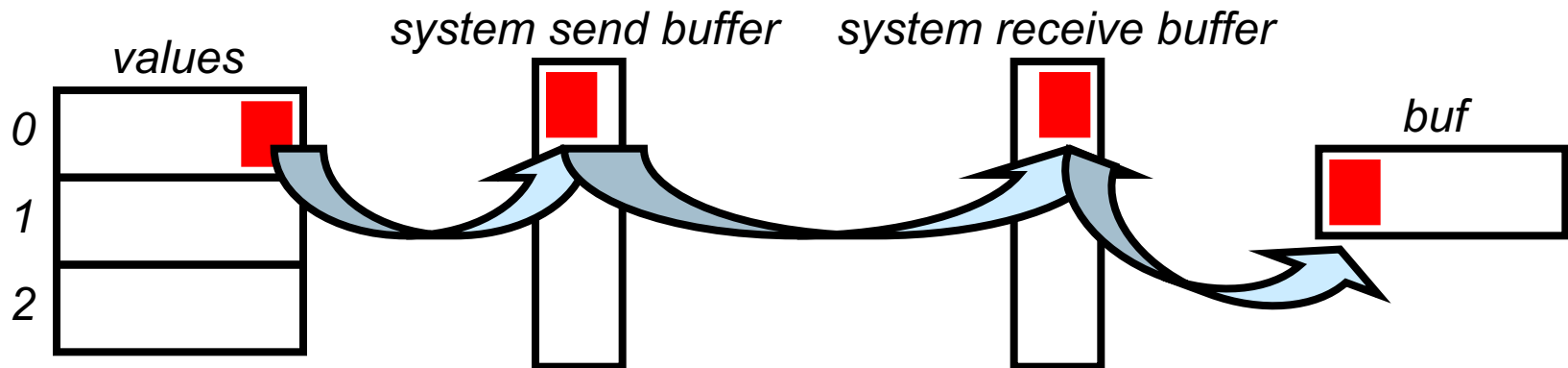
# Shifts: Definition of Neighbors

```
if (myid==0){
  lnbr=np-1; rnbr=myid+1;
}
else if (myid==np-1){
  lnbr=myid-1; rnbr=0;
}
else{
  lnbr=myid-1; rnbr=myid+1;
}

if (myid==0) shifts=atoi(argv[1]);
MPI_Bcast (&shifts, 1, MPI_INT, 0, MPI_COMM_WORLD);

values= (int *) calloc(100/np,sizeof(int));
if (myid==0){
  values[0]=1;
}
```

# Shifts: Shift the array

```
for (i=0;i<shifts;i++){
  int buf;

  MPI_Send(&values[0],1,MPI_INT,lnbr,10,MPI_COMM_WORLD);
  MPI_Recv(&buf, 1, MPI_INT,rnbr,10,
                              MPI_COMM_WORLD, &status);
  for (j=1;j<100/np;j++){
    values[j-1]=values[j];
  }
  values[100/np-1]=buf;
}
```

*system send buffer*     *system receive buffer*

*values*

*buf*

0

1

2

# MPI_Sendrecv

```
int MPI_Sendrecv (void *sendbuf, int sendcount, MPI_Datatype sendtype,
              int dest, int sendtag, void *recvbuf, int recvcount,
              MPI_Datatype recvtype, int source, MPI_Datatype recvtag,
              MPI_Comm comm, MPI_Status *status)
```

- Sendbuf and recvbuf have to be different

- Equivalent to the execution of MPI_Send and MPI_Receive in parallel threads.

- MPI_Sendrecv_replace:
  - Only one buffer.
  - The sent message is replaced by the received message.

# Shifts: Shift the array

```c
for (i=0;i<shifts;i++){
    int buf=values[0];

    for (j=1;j<100/np;j++){
      values[j-1]=values[j];
    }

    MPI_Sendrecv(&buf, 1, MPI_INT, lnbr, 10,
              &values[100/np-1], 1, MPI_INT, rnbr, 10
              MPI_COMM_WORLD, &status);
}
```

# MPI_Wtime

double MPI_Wtime (void)

- Return value represents elapsed wall-clock time since some time in the past measured in seconds.
- The time returned is local to the node that called it.

# Collective Operations

- Properties
  - Must be executed by all processes of the process group.
  - Must be executed in the same sequence.
  - All collective operations are blocking operations.
- MPI provides three classes of collective operations
  - Synchronization
    - Barrier
  - Communication
    - Broadcast
    - gather
    - scatter
  - Reduction
    - Global value returned to one or all processes.
    - Combination with subsequent scatter.
    - Parallel prefix operations