



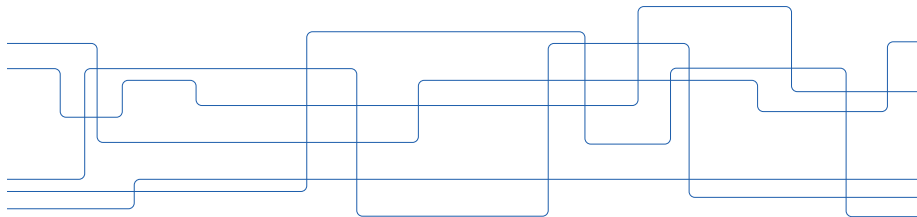
MPI: Part III

AQTIVATE Training Workshop I

Dirk Pleiter

CST | EECS | KTH

December 2023





Overview

Communicators

Derived Data Types

One-Sided Communication

Hybrid Parallelisation



Content

Communicators

Derived Data Types

One-Sided Communication

Hybrid Parallelisation

Communicator Creation: Duplication

- ▶ Create a new communicator based on an existing communicator

```
1 int MPI_Comm_dup(  
2     MPI_Comm comm,      /* Existing communicator */  
3     MPI_Comm *newcomm  /* New communicator */  
4 );
```

- ▶ The new communicator
 - ▶ Refers to the same set of processes and leaves rank ordering and other associated values unchanged
 - ▶ Defines a new context for message communication (used, e.g., by parallel libraries)

Communicator Creation: Partition (1/2)

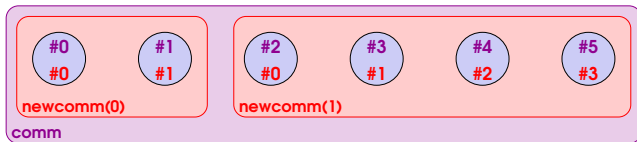
- ▶ MPI allows to partition the processes in disjoint subgroups in a flexible manner using

```
1  int MPI_Comm_split(  
2    MPI_Comm comm,      /* Existing communicator */  
3    int color ,         /* Control of subset assignment */  
4    int key ,           /* Control of rank assignment */  
5    MPI_Comm *newcomm  /* New communicator */  
6  );
```

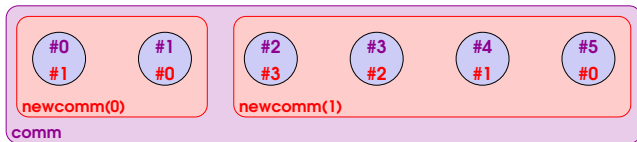
- ▶ For each value of color a different subgroup (and communicator) is created
- ▶ Within each subgroup, the processes are ranked in the order defined by key
 - ▶ The rank of the old group is used by tie breaking

Communicator Creation: Partition (2/2)

- ▶ Example #1: $\text{color} = (\text{rank} < 2) ? 0 : 1; \text{key} = 0$



- ▶ Example #2: $\text{color} = (\text{rank} < 2) ? 0 : 1; \text{key} = -\text{rank}$



Inter-Communicators (1/2)

- ▶ Communications between different subgroups (partitions) can be established using **inter-communicators**, e.g. using

```
1 int MPI_Intercomm_create(  
2     MPI_Comm local_comm,      /* Local communicator */  
3     int local_leader,         /* Rank of local leader */  
4     MPI_Comm peer_comm,      /* Peer communicator */  
5     int remote_leader,       /* Rank of remote leader */  
6     int tag,                 /* Tag */  
7     MPI_Comm *newintercomm /* Inter-communicator */  
8 );
```

- ▶ `local_comm` and `local_leader` identify lead process within local subgroup
- ▶ `peer_comm` must be a group to which both the local and remote leader process belong to
- ▶ The remote leader is identified within the `peer_comm` through `peer_leader`

Inter-Communicators (2/2)

- ▶ The syntax of point-to-point and collective communication is the same for both inter- and intra-communication
 - ▶ A target process is addressed by its rank in the remote group

Process Topologies

- ▶ Communicators can be used to define a **virtual network topology** that defines the neighbourhood relations between the processes within the given group
- ▶ Benefits of such virtual topologies:
 - ▶ Allow to express communication patterns in a more natural way
 - ▶ Provide middleware with hints about target communication patterns
 - ▶ MPI may use this information to select a better mapping of processes to underlying **physical network topology** (if user allows for rank re-ordering)

Process Topologies: Cartesian (1/2)

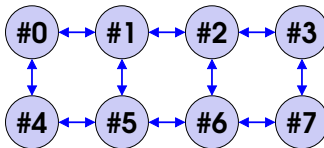
- With a Cartesian communicator processes are assumed to be arranged on a one or higher-dimensional mesh or torus

```

1 int MPI_Cart_create(
2     MPI_Comm comm,           /* Existing communicator */
3     int ndims,               /* Number of dimensions of grid */
4     const int dims[],        /* Grid size */
5     const int periods[],     /* True (false) if grid is (not) periodic */
6     int reorder,             /* Allow to reorder or not */
7     MPI_Comm *comm_cart     /* New communicator */
8 );

```

- Example with 8 processes and $\text{ndims} = 2$, $\text{dims} = [4, 2]$, $\text{periods} = [0, 0]$ and $\text{reorder} = 0$:



Process Topologies: Cartesian (2/2)

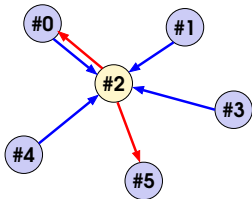
- ▶ MPI provides a (convenience) function to compute source and destination to communication with the grid:

```
1 int MPI_Cart_shift(  
2     MPI_Comm comm_cart, /* Cartesian communicator */  
3     int direction,      /* Communication direction */  
4     int disp,           /* Displacement (sign determines direction) */  
5     int *source,        /* Returned source rank */  
6     int *dest           /* Returned destination rank */  
7 );
```

- ▶ The returned source and dest ranks can be directly used as input for MPI_Sendrecv

Process Topologies: Graph (1/2)

- ▶ For graph-based algorithms, MPI allows defining communicator with each process specifying each of its incoming and outgoing (adjacent) edges in the logical communication graph
- ▶ Weights may be provided to influence the rank reordering
 - ▶ `MPI_UNWEIGHTED` may be used in case no weights are known
- ▶ Example:
 - ▶ `indegree = 4`
`sources = [0, 1, 3, 4]`
 - ▶ `outdegree = 2`
`sources = [0, 5]`



Process Topologies: Graph (2/2)

► Syntax:

```
1 int MPI_Dist_graph_create_adjacent(  
2     MPI_Comm comm,           /* Existing communicator */  
3     int indegree,            /* Size of sources/sourceweights arrays */  
4     const int sources[],      /* List of ranks where called is a  
5                               possible destinations */  
6     const int sourceweights[], /* Weight of the edges */  
7     int outdegree,           /* Size of destination/destinationweights  
8                               arrays */  
9     const int destinations[], /* List of possible destinations */  
10    const int destweights[],  /* Weight of the edges */  
11    MPI_Info info,            /* Information object */  
12    int reorder,              /* Allow to reorder or not */  
13    MPI_Comm *comm_dist_graph /* New communicator */  
14 );
```



Content

Communicators

Derived Data Types

One-Sided Communication

Hybrid Parallelisation

Creating Derived Data Types

- ▶ A **derived data** type (also called general data type) is an opaque object that specifies two things:
 - ▶ A sequence of basic data types
 - ▶ A sequence of integer (byte) displacements
- ▶ Creating derived data types involves the following steps:
 - ▶ Definition of derived data types through type constructor functions
 - ▶ Commit of the new type
 - ▶ Release of all types

Derived Data Types: Contiguous Type

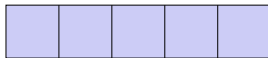
- Constructor allows replication of a data type into contiguous locations:

```
1 int MPI_Type_contiguous(  
2   int count,           /* Replication count */  
3   MPI_Datatype type_exist, /* Existing data type */  
4   MPI_Datatype *type_new /* New data type */  
5 );
```

type_exist



type_new



count = 5

Derived Data Types: Vector Type

- `MPI_Type_vector` is a more general constructor that allows replication of a data type into locations that consist of equally spaced blocks:

```

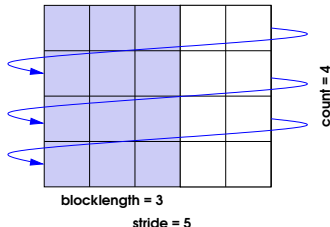
1 int MPI_Type_vector(
2   int count,           /* Number of blocks */
3   int blocklength,     /* Number of elements in each block */
4   int stride,          /* Number of elements between start of each block */
5   MPI_Datatype type_exist, /* Existing data type */
6   MPI_Datatype *type_new /* New data type */
7 );

```

`type_exist`



`type_new`



Derived Data Types: Subarray Type

- To describe an n-dimensional sub-array of an n-dimensional array use

```
1 int MPI_Type_create_subarray(
2   int ndims,
3   const int size[],
4   const int subsize[],
5   const int start[],
6   int order,
7   MPI_Datatype type_exist,
8   MPI_Datatype *type_new
9 );
```

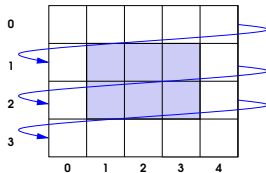
- Example:

- `ndims = 2`
- `size[] = [5, 4]`
- `subsize[] = [3, 2]`
- `start[] = [1, 1]`
- `order = MPI_ORDER_FORTRAN`

type_exist



type_new



Derived Data Types: Indexed Type (1/2)

- Even more general: Replicate an existing data type into a sequence of blocks, where each block can contain a different number of copies and have a different displacement:

```

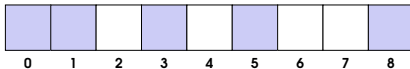
1  int MPI_Type_indexed(
2  int count,           /* Array length */
3  const int blocklength[], /* Number of elements per block */
4  const int displacement[], /* Displacement for each block */
5  MPI_Datatype type_exist, /* Existing data type */
6  MPI_Datatype *type_new  /* New data type */
7  );
    
```

- Example: $\text{count} = 4$, $\text{blocklength}[] = [2, 1, 1, 1]$,
 $\text{displacement}[] = [0, 3, 5, 8]$

type_exist



type_new



Derived Data Types: Indexed Type (2/2)

- ▶ Slightly less general: Assume all blocks to be of the same length:

```
1 int MPI_Type_create_indexed_block(
2   int count,
3   int blocklength,
4   const int displacement[],
5   MPI_Datatype type_exist,
6   MPI_Datatype *newtype
7 );
```

- ▶ Example: $\text{count} = 5$, $\text{blocklength} = 1$,
 $\text{displacement}[] = [0, 1, 3, 5, 8]$

type_exist



type_new



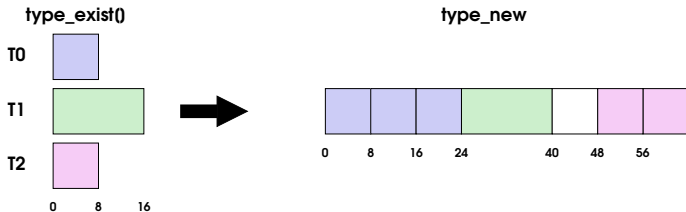
Derived Data Types: Structured Type

- Finally, the most general constructor:

```

1  int MPI_Type_create_struct(
2  int count,                               /* Array length */
3  const int blocklength[],                 /* List of number of elements per block */
4  const MPI_Aint displacement[],          /* List of displacements in Bytes */
5  const MPI_Datatype type[],              /* List of data types */
6  MPI_Datatype *type_new                  /* New data type */
7  );
    
```

- Example: $\text{count} = 3$, $\text{blocklength} = [3, 1, 2]$,
 $\text{displacement} = [0, 24, 48]$, $\text{type} = [T0, T1, T2]$,



Derived Data Types: Caveats

- ▶ The data type constructors typically assume that objects and arrays are stored linearly in memory
 - ▶ For dynamically allocated multi-dimensional arrays in C or C++ this may not be the case
- ▶ Performance depends on the data types (and the MPI implementation)
 - ▶ Generally expect more general data types to be slower



Content

Communicators

Derived Data Types

One-Sided Communication

Hybrid Parallelisation

MPI One-Sided Communication

- ▶ Send-receive communication assumes that both end-points of the communication take action (**two-sided communication**)
- ▶ MPI since version 2.0 supports also communication based on a get/put semantics with only one end-point being active (**one-sided communication**)
 - ▶ Also called **remote memory access (RMA)**
- ▶ Challenges:
 - ▶ Need for a common address space
 - ▶ Suitable synchronisation to ensure RAW and WAR dependencies to be respected



One-Sided Communication: Basic Operations

- ▶ Create a memory window for remote access
- ▶ Communicate using, e.g.,
 - ▶ MPI_Put: Write to remote memory
 - ▶ MPI_Get: Read from remote memory
- ▶ Synchronise using MPI_Win_fence

Memory Window Creation

- ▶ Specify a window of existing memory that processes exposes to RMA accesses:

```
1 int MPI_Win_create(  
2     void *base,      /* Initial address of window */  
3     MPI_Aint size,    /* Size of window [Bytes] */  
4     int disp_unit,    /* Local unit size for displacements [Bytes] */  
5     MPI_Info info,    /* Info object */  
6     MPI_Comm comm,    /* Communicator */  
7     MPI_Win *win      /* Window handle (out) */  
8 );
```

- ▶ Operation is a collective call, all processes must participate
- ▶ All address computations are scaled by the `disp_unit`
 - ▶ Choose `disp_unit = 1` for no scaling
- ▶ The `MPI_Info` object allows to provide optimisation hints
(not covered here)
 - ▶ `MPI_Info` allows to store an unordered set of (key,value) pairs
 - ▶ May use `info = MPI_INFO_NULL`

Synchronisation

- ▶ All RMA calls related to a particular window can be synchronised using

```
1 int MPI_Win_fence(  
2   int assert, /* Program assertion */  
3   MPI_Win win /* Window handle */  
4 );
```

- ▶ The application may make assertions, which may allow for optimisations but `assert = 0` is always valid
- ▶ `MPI_Win_fence` separates different **access epochs**
 - ▶ Epoch starts with an RMA synchronisation call on `win`
 - ▶ Next, a number of communication calls may happen
 - ▶ Epoch stops with another synchronisation call
- ▶ Other synchronisation mechanisms available (not considered here)

MPI_Put

- ▶ The analogue of a send operation is the following operation that puts data into remote memory:

```
1 int MPI_Put(  
2     void *oaddr,           /* Address of origin buffer */  
3     int ocount,            /* Number of entries in origin buffer */  
4     MPI_Datatype odatatype, /* Data type of entries in origin buffer */  
5     int trank,              /* Target rank */  
6     MPI_Aint tdisp,         /* Displacement from window start to target buffer */  
7     int tcount,            /* Number of entries in target buffer */  
8     MPI_Datatype tdatatype, /* Data type of entries in origin buffer */  
9     MPI_Win win             /* Window handle */  
10 );
```

- ▶ The data is
 - ▶ Read from the local memory at address oaddr
 - ▶ written to the remote memory at address
 $taddr = window_base + tdisp \times disp_unit$
- ▶ Need synchronisation to avoid
 - ▶ Remote data is overwritten before read of old data (WAR)
 - ▶ Remote data is used before new data is written (RAW)

- ▶ To copy data from the target memory to the origin use

```
1 int MPI_Get(  
2     void *oaddr,           /* Address of origin buffer */  
3     int ocount,            /* Number of entries in origin buffer */  
4     MPI_Datatype odatatype, /* Data type of entries in origin buffer */  
5     int trunk,             /* Target rank */  
6     MPI_Aint tdisp,        /* Displacement from window start to target buffer */  
7     int tcount,            /* Number of entries in target buffer */  
8     MPI_Datatype tdatatype, /* Data type of entries in origin buffer */  
9     MPI.Win win            /* Window handle */  
10 );
```

- ▶ Need synchronisation to avoid
 - ▶ Remote data is not updated before get is executed (RAW)

One-Sided Communication: Ping (1/2)

```
1 int main(int argc ,char *argv [])
2 {
3     MPI_Init(NULL, NULL);
4
5     int irank;
6     MPI_Comm_rank(MPI_COMM_WORLD, &irank);
7
8     int buf = -1; MPI_Win win;
9     MPI_Win_create(&buf, sizeof(int), sizeof(int),
10                  MPI_INFO_NULL, MPI_COMM_WORLD, &win);
11
12     int src = 1; int dst = 0;
13     MPI_Win_fence(0, win);
14     if (irank == src) {
15         int x = 100 + irank;
16         MPI_Put(&x, 1, MPI_INT, dst, 0, 1, MPI_INT, win);
17         printf("#%d: have put x=%d\n", irank, x);
18     }
19     MPI_Win_fence(0, win);
20
21     printf("#%d: buf=%d\n", irank, buf);
22
23     MPI_Finalize();
24
25     return 0;
26 }
```

What is the expected output?

One-Sided Communication: Ping (2/2)

- ▶ Answer when using 4 processes (order is not deterministic):

#1: have put x=101

#1: buf=-1

#3: buf=-1

#0: buf=101

#2: buf=-1



Content

Communicators

Derived Data Types

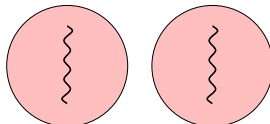
One-Sided Communication

Hybrid Parallelisation

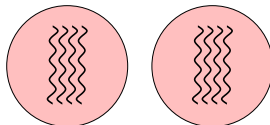
Hybrid Parallelisation Introduction (1/2)

- ▶ MPI parallelism
 - ▶ processes possibly running on different nodes
 - ▶ Separate memory address spaces
- ▶ Thread parallelism (e.g., OpenMP)
 - ▶ Multiple threads running within a single process
 - ▶ Access to a shared memory address space
- ▶ **Hybrid parallelisation** =
MPI + threads

MPI only



MPI + threads



Rank #0

Rank #1

Hybrid Parallelisation Introduction (2/2)

- ▶ All threads within an MPI process share all MPI objects
 - ▶ Example: Communicators
- ▶ Benefits of hybrid parallelisation
 - ▶ Opportunity for more parallelism
 - ▶ Reduced memory footprint
- ▶ Possible need for tuning
 - ▶ Using less MPI ranks per node may reduce overheads
 - ▶ Parallel efficiency reduces for a larger number of threads

Thread-safe MPI: Introduction

- ▶ MPI defines four levels of thread safety:
 - ▶ `MPI_THREAD_FUNNELED`
 - ▶ Multi-threaded, but only the main thread makes MPI calls (the one that calls `MPI_Init_thread`)
 - ▶ `MPI_THREAD_SERIALIZED`
 - ▶ Multi-threaded, but only one thread at a time makes MPI calls
 - ▶ `MPI_THREAD_MULTIPLE`
 - ▶ Multi-threaded and any thread can make MPI calls at any time (with some restrictions)
 - ▶ `MPI_THREAD_SINGLE`
 - ▶ Only one thread exists in the application
- ▶ **Note:** Choosing for a particular level means making commitments to MPI

Thread-safe MPI: MPI_THREAD_SINGLE

- ▶ Only one thread is used
 - ▶ There are no OpenMP parallel regions

```
1 int main(int argc, char ** argv)
2 {
3     double x[100];
4     int rank;
5
6     MPI_Init(&argc, &argv);
7     MPI_Comm_rank(MPI_COMM_WORLD, &rank);
8
9     for (int i = 0; i < 100; i++)
10         f(x[i]);
11
12     /* MPI communication */
13
14     MPI_Finalize();
15
16     return 0;
17 }
```

Thread-safe MPI: MPI_THREAD_FUNNELED

- ▶ All MPI calls are made by the master thread
 - ▶ All MPI calls are outside of OpenMP parallel regions

```
1 int main(int argc, char ** argv)
2 {
3     double x[100];
4     int rank, provided;
5
6     MPI_Init_thread(&argc, &argv, MPI_THREAD_FUNNELED, &provided);
7     MPI_Comm_rank(MPI_COMM_WORLD, &rank);
8
9     #pragma omp parallel for
10    for (int i = 0; i < 100; i++)
11        f(x[i]);
12
13    /* MPI communication */
14
15    MPI_Finalize();
16
17    return 0;
18 }
```

Thread-safe MPI: MPI_THREAD_SERIALIZED

- ▶ Only one thread can make MPI calls at a time
 - ▶ Use OpenMP critical clause

```
1 int main(int argc, char ** argv)
2 {
3     double x[100];
4     int rank, provided;
5
6     MPI_Init_thread(&argc, &argv, MPI_THREAD_SERIALIZED, &provided);
7     MPI_Comm_rank(MPI_COMM_WORLD, &rank);
8
9     #pragma omp parallel for
10    for (int i = 0; i < 100; i++) {
11        f(x[i]);
12    }
13    #pragma omp critical
14        /* MPI communication */
15    }
16
17    MPI_Finalize();
18
19    return 0;
20 }
```

Thread-safe MPI: MPI_THREAD_MULTIPLE

- ▶ Any thread can make MPI calls any time, but there are caveats

```
1 int main(int argc, char ** argv)
2 {
3     double x[100];
4     int provided;
5
6     MPI_Init_thread(&argc, &argv, MPI_THREAD_MULTIPLE, &provided);
7     MPI_Comm_rank(MPI_COMM_WORLD, &rank);
8
9     #pragma omp parallel for
10    for (i = 0; i < 100; i++) {
11        f(x[i]);
12
13        /* MPI communication */
14    }
15
16    MPI_Finalize();
17
18    return 0;
19 }
```

MPI_THREAD_MULTIPLE: Caveats

► Ordering

- MPI calls within a thread are executed in order
- MPI calls in different threads are executed in no particular order
- Caveat: User is responsible to prevent different threads making conflicting MPI calls
- Example: Ordering of MPI_Send and MPI_Recv calls

► Blocking

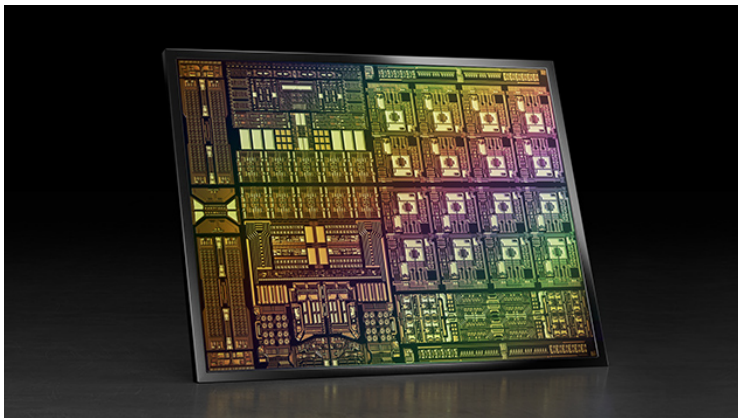
- MPI calls will block only the calling thread
- Caveat: User must ensure that all relevant threads are blocked
- Example: Not all threads calling MPI_Barrier

Thread-safe MPI: Implementation Status

- ▶ All MPI implementations support `MPI_THREAD_SINGLE`
- ▶ `MPI_THREAD_FUNNELED` is typically supported
 - ▶ Situation on Dardel: Current CPE supports this levels
- ▶ `MPI_THREAD_SERIALIZED` and `MPI_THREAD_MULTIPLE` is possibly not supported
 - ▶ Situation on Dardel: Current CPE supports both levels

Finish with Upcoming Networking Technology

NVIDIA BlueField-3 with in-network processing capabilities:



[NVIDIA]