# MPIBlib: MPI Benchmark library

Version 1.2.0 (Revision 308)

Generated by Doxygen 1.7.1

Wed Jun 8 2011 18:06:02

# Contents

# 1 Introduction

Accurate estimation of the execution time of MPI communication operations plays an important role in optimization of parallel applications. A priori information about the performance of each MPI operation allows a software developer to design a parallel application in such a way that it will have maximum performance. This data can also be useful for tuning collective communication operations and for the evaluation of different available implementations. The choice of collective algorithms becomes even more important in heterogeneous environments.

A typical MPI benchmarking suite uses only one timing method to estimate the execution time of the MPI communications. The method provides a certain accuracy and efficiency. The efficiency of the timing method is particularly important in self-adaptable parallel applications using runtime benchmarking of communication operations to optimize their performance on the executing platform. In this case, less accurate results can be acceptable in favor of a rapid response from the benchmark. We design a new MPI benchmarking suite called MPIBlib [1] that provides a variety of timing methods. This suite supports both fast measurement of collective operations and exhaustive benchmarking.

In addition to general timing methods that are universally applicable to all communication operations, MPIBlib includes methods that can only be used for measurement of one or more specific operations. Where applicable, these operation-specific methods work faster than their universal counterparts and can be used as their time-efficient alternatives.

Most of the MPI benchmarking suites are designed in the form of a standalone executable program that takes the parameters of communication experiments and produce a lot of output data for further analysis. As such, they cannot be integrated easily and efficiently into application-level software. Therefore, there is a need for a benchmarking library that can be used in parallel applications or programming systems for communication performance modeling and tuning communication operations. MPIBlib is such a library that can be linked to other applications and used at runtime.

## 1.1 Authors

Alexey Lastovetsky, Vladimir Rychkov, Maureen O'Flynn, Kiril Dichev

Heterogeneous Computing Laboratory

School of Computer Science and Informatics, University College Dublin

Belfield, Dublin 4, Ireland

http://hcl.ucd.ie

{alexey.lastovetsky, vladimir.rychkov}@ucd.ie

# References

[1] A. Lastovetsky, V. Rychkov, and M. OFlynn. MPIBlib: Benchmarking MPI communications for parallel computing on homogeneous and heterogeneous clusters. In A. Lastovetsky, T. Kechadi, and J. Dongarra, editors, *Recent Advances in Parallel Virtual Machine and Message Passing Interface: 15th European PVM/MPI User's Group Meeting*, volume 5205, pages 227–238, Dublin, Ireland, September 7-10 2008. Springer-Verlag Berlin Heidelberg. 1, 11

[2] B.R. de Supinski and N.T. Karonis. Accurately measuring MPI broadcasts in a computational grid. In *HPDC'99: Proceedings of the 8th IEEE International Symposium on High Performance Distributed Computing*, pages 29–37, Washington, DC, USA, 1999. IEEE Computer Society. 10

[3] J.L. Traff. Hierarchical gather/scatter algorithms with graceful degradation. In *Proceedings of IPDPS'04*, pages 80–89, 2004. 9, 31, 58, 59

[4] T. Worsch, R Reussner, and W. Augustin. On benchmarking collective MPI operations. In *Proceedings of the 9th European PVM/MPI Users' Group Meeting on Recent Advances in Parallel Virtual Machine and Message Passing Interface*, pages 271–279, London, UK, 2002. Springer-Verlag. 13

# 2 Installation

```
Installation
============

Required software:
1. any C/C++ and MPI (MPICH-1 does not support shared libraries)
2. GSL (GNU Scientific Library)
3. Boost (The Boost C++ libraries: Graph)

Optional software:
1. Gnuplot (An Interactive Plotting Program) -
   optional (for performance diagrams)
2. Graphviz (Graph Visualization Software: dot) -
   optional (for tree visualization)

GSL
If GSL is installed in a non-default directory
$ export LD_LIBRARY_PATH=DIR/lib:$LD_LIBRARY_PATH

Boost
1. Boost should be configured with at least the Graph library
   (default: all)
$ ./configure --prefix=DIR --with-libraries=graph
2. Default installation:
 - DIR/include/boost_version/boost
```

```
  - DIR/lib/libboost_library_versions.*
Create symbolic links:
$ cd DIR/include; ln -s boost_version/boost
$ cd DIR/lib; ln -s libboost_[library]_[version].[a/so] libboost_[library].[a/so]

$ export LD_LIBRARY_PATH=DIR/lib:$LD_LIBRARY_PATH

For users
---------

Download the latest package from http://hcl.ucd.ie/project/mpiblib

$ tar -zxvf mpiblib-X.X.X.tar.gz
$ cd mpiblib-X.X.X
$ ./configure
$ make all install

Configuration
-------------

Packages:
  --with-gsl-dir=DIR      GNU Scientific Library directory
  --with-boost-dir=DIR    The Boost C++ libraries directory

Check configure options:
$ ./configure -h

For developers
--------------

Required software:
1. Subversion
2. GNU autotools
3. Doxygen, Graphviz and any TeX - optional (for reference manual)

$ svn co https://hcl.ucd.ie/repos/CPM/trunk/MPIBlib
$ cd MPIBlib
$ autoreconf -i
$ ./configure --enable-debug
$ make all

To create a package:
$ svn log -v > ChangeLog
$ make dist
```

# 3 The software design

MPIBlib is implemented in C/C++ on top of MPI. The package consists of libraries, tools and tests. The libraries implement point-to-point and collective benchmarks and different algorithms of point-to-point and collective communication operations. The tools verify algorithms and perform benchmarks. Their output includes the results of measurements and communication trees (for tree-based algorithms of collective operations). The results of measurements can be visualized by the gnuplot utility; MPIBlib provides the basic gnuplot scripts. The communication trees built in the tree-based collective algorithms can be visualized by the dot utility (a part of Graphviz).

```
┌─────────────────────┐  ┌─────────────────┐  ┌──────────────────────┐
│ Benchmarks library: │  │ P2P library:    │  │ Collectives libraries:│
│ libmpib.a           │  │ libmpib_p2p.a   │  │ libmpib_coll.so      │
└─────────────────────┘  └─────────────────┘  └──────────────────────┘
```



# 4  Benchmarks library

A static library `libmpib.a` implements point-to-point and collective benchmarks. Main library modules (the modules marked grey can be extended, providing the benchmarking of user-defined point-to-point and collective operations):

- Measurement: base data structures and functions (depends on the GNU Scientific Library)

- Operation-specific benchmarks (now includes only benchmark for bcast)

- Point-to-point benchmarks (seqential/parallel point-to-point benchmark)

- Containers for point-to-point communication operations (data structures describing point-to-point communications to be measured, which are used as an argument of the point-to-point benchmark function)

- Collective benchmarks (collective benchmarks based on root, max and global timing methods)

- Containers for collective communication operations (data structures describing collective communications to be measured, which are used as an argument of collective benchmark functions)

- Definitions of MPI communication operations (definitions of MPI communication operations)

## 4.1 Tools

- P2P benchmark

- Collective benchmark

- Generator of multiplying factors for benchmarking scatterv/gatherv

Command-line arguments are described in Options for executables.

## 4.2 Tests

- Scatter-gather-based p2p benchmarks

- Point-to-point communication by fragmenting in a hybrid MPI/OpenMP approach

- Eager protocol point-to-point benchmark

## 4.3 P2P benchmark

Performs adaptive p2p benchmark between a pair of processors 0-1. Arguments are described in Options for executables.

**Example:**

```
$ mpirun -n 2 p2p -m 1024 -M 2049 > p2p.out
```

will benchmark the p2p communication between 2 processes for message sizes of 1024 and 2048 bytes

**p2p.plot** draws the graph of the execution time (sec) against message size (kb).

- input: p2p.out
- output: p2p.eps (0-1 with error bars)

Using the gnuplot script:

```
$ gnuplot p2p.plot
```

## 4.4 Collective benchmark

Performs regular universal or operation-specific collective benchmark. Basic arguments are described in Options for executables.

**Example:**

```
$ mpirun -n 16 collective -O MPIB_Scatter_binomial -l libmpib_coll.so -m 1024 -M
      2049 > collective.out
```

will benchmark the default binomial tree scatter operation with message sizes per process of 1024 and 2048

**collective.plot** draws the graph of the execution time (sec) against message size (kb) with error bars.

- input: collective.out
- output: collective.eps

Using the gnuplot script:

```
$ gnuplot collective.plot
```

## 4.5 Generator of multiplying factors for benchmarking scatterv/gatherv

A standalone tool for generating factors which can be used as input file for benchmarking scatterv/gatherv in Collective benchmark. The factors are normalized so that the average factor is 1. Arguments can be:

- -r for random factors
- -c for factors based on a small GSL matrix multiplication at every process

**Example:**

```
mpirun -n 2 generate_factors -c > factors
```

## 4.6 Scatter-gather-based p2p benchmarks

Program for testing the correctness and benchmarking of the modified scatter-gather-based p2p. The correctness is checked by verifying the same data is received that is sent. The benchmarks are done by taking the maximum time for each operation to complete from sender or receiver, and taking the minimum from a number of iterations with this approach

To run:

```
cd <install-dir>/bin
mpirun -n 2 <src-dir>/MPIBlib/tests/sg-based-p2p
```

## 4.7 Point-to-point communication by fragmenting in a hybrid MPI/OpenMP approach

This standalone program uses a hybrid approach to send a message in fragments from MPI process 0 to MPI process 1. We create a number of OpenMP threads on each MPI process. Each process then uses all OpenMP threads to submit equal fragments of the message. Each fragment is transferred by a different OpenMP thread using MPI_Send/MPI_Recv.

The program can only run properly if the MPI library supports MPI_THREAD_MULTIPLE

```
$ mpirun -n 2 ./fragment-omp-for
```

## 4.8 Eager protocol point-to-point benchmark

Point-to-point benchmarks using MPI_Rsend/MPI_Recv pairs. It usually performs better than MPI_-Send/MPI_Recv on high-latency links, probably because it uses the eager protocol.

The maximum time measured at sender or receiver is taken, and after a number of iterations the minimum time is taken.

# 5 P2P library

A static library `libmpib_p2p.a` implements different algorithms of p2p communications.

- The API for using the modified p2p communication is available here: Wrappers for p2p operations.

- The spawned process which is always used for the modified p2p communication is Spawned process for scatter-gather point-to-point communication

## 5.1 Spawned process for scatter-gather point-to-point communication

This program is launched by a master process performing scatter-gather based point-to-point communication. It can be launched a number of times on each host. The program listens for incoming control messages from any process in a loop:

- if a control message has `info[0] != 1` , then the program terminates

- if a control message has `info[0] == 1` , then the process receives a chunk of a message from a sender process and forwards the message to a receiver process

# 6 Collectives library

A shared library `libmpib_coll.so` implements different algorithms of collectives

- Basic algorithms of MPI collective operations (mostly linear)

- Tree-based algorithms of MPI collective operations (depends on the Boost C++ libraries)

The command-line arguments:

- **verbose** verbose mode (default: no)

- **sgv** scatterv/gatherv mode

    - *0* propagated on duplicated communicator (default)

    - *1* broadcasted counts

    - *2* propagated with tags: not safe

    - *3* everywhere-defined counts: not standard

## 6.1 Tools

- Test for MPI collective operations

## 6.2 Examples

- Substitute for native MPI collectives

- Test for tree builders

- Test for the tree-based scatterv/gatherv

## 6.3 Test for MPI collective operations

Performs tests for the implementations of MPI collective operations with different root

## 6.4 Substitute for native MPI collectives

Substitute for native MPI collective operations in parallel applications

- add `-I$(MPIBlib-dir)/include` to CPPFLAGS

- add `-L$(MPIBlib-dir)/lib -lmpib_coll` to LDFLAGS

- replace `#include <mpi.h>` by `#include "mpib_coll.h"`

- define macro substitutes for native MPI collective operations `#define MPI_Scatter MPIB_-Scatter_flat`

```
//#include <mpi.h>
#include "mpib_coll.h"
#define MPI_Scatter MPIB_Scatter_flat_nb
#define MPI_Gather MPIB_Gather_flat_nb

#include <malloc.h>

int main(int argc, char** argv)
{
    MPI_Init(&argc, &argv);
    int rank;
    MPI_Comm_rank(MPI_COMM_WORLD, &rank);
    int size;
    MPI_Comm_size(MPI_COMM_WORLD, &size);
    int M = 1024;
    char* buffer = rank == 0 ? (char*)malloc(sizeof(char) * M * size) : (char*)ma
      lloc(sizeof(char) * M);
```

```
    MPI_Scatter(buffer, M, MPI_CHAR, buffer, M, MPI_CHAR, 0, MPI_COMM_WORLD);
    MPI_Gather(buffer, M, MPI_CHAR, buffer, M, MPI_CHAR, 0, MPI_COMM_WORLD);
    free(buffer);
    MPI_Finalize();
    return 0;
}
```

## 6.5 Test for tree builders

Usage:

- include a header file with your builder class

- use your builder class in the main function

```
#include <getopt.h>

#include <boost/graph/graphviz.hpp>

#include "collectives/brsg_tree_builders.hpp"
using namespace MPIB::BRSG;

int main(int argc, char** argv) {
    int n = 8;
    int c;
    while ((c = getopt(argc, argv, "hn:")) != -1) {
        switch (c) {
            case 'h':
                cerr <<
"usage: trees [options]\n\
-h        help\n\
-n I      number of nodes in the tree\n";
                return 0;
                break;
            case 'n':
                n = atoi(optarg);
                break;
        }
    }
    Graph graph;
    Vertex r, u, v;
    Binomial_builder builder;
    builder.build(n, 0, 0, 0, graph, r, u, v);
    cout << "#Tree\n";
    write_graphviz(cout, graph);
    return 0;
}
```

## 6.6 Test for the tree-based scatterv/gatherv

Test for tree-based scatterv/gatherv. Based on [3].

# 7 Module Documentation

## 7.1 Operation-specific benchmarks

**Functions**

- void MPIB_bcast_timer_init (MPI_Comm comm, int parallel, MPIB_precision precision)

- void MPIB_measure_bcast (MPIB_Bcast bcast, MPI_Comm comm, int root, int M, int max_reps, MPIB_result ∗result)

### 7.1.1    Detailed Description

This module provides operation-specific benchmarks.

### 7.1.2    Function Documentation

#### 7.1.2.1    void MPIB_bcast_timer_init ( MPI_Comm *comm,* int *parallel,* MPIB_precision *precision* )

Performs the p2p benchmark with empty message with given precision and sets up an internal global variable, which is used by MPIB_measure_bcast. Called by MPIB_measure_bcast. Can be called directly before measurements.

**Attention**

As the global variable is an array, the memory should be freed by

```
MPIB_bcast_timer_init(MPI_COMM_NULL, (MPIB_precision){0, 0, 0})
```

**Parameters**

*comm* communicator

*parallel* several non-overlapped point-to-point communications at the same time if non-zero

*precision* measurement precision

#### 7.1.2.2    void MPIB_measure_bcast ( MPIB_Bcast *bcast,* MPI_Comm *comm,* int *root,* int *M,* int *max_reps,* MPIB_result ∗ *result* )

Measures the execution time of bcast between root and the rest of processes. Based on [2]. Reuses already obtained empty-roundtrip times if the previous initialization was performed on the same communicator. Otherwise, initializes the bcast timer by calling MPIB_bcast_timer_init.

In the loop over processes, except root:

In the loop over repetitions:

- bcast at all processes except root and current process in the loop

- bcast and empty recv at root and measures the execution time

- bcast and empty send at current process in the loop

Having substructed the half of empty-roundtrip times, finds maximum.

Broadcasts the result.

**Note**

No double barriers as no need in synchronization.
No precision as we cannot interrupt the process of measurement by broadcasting the error value after each repetition.

---

**Parameters**

> ***bcast***  bcast implementation
>
> ***comm***  communicator
>
> ***root***  root process
>
> ***M***  message size
>
> ***max_reps***  maximum number of repetitions
>
> ***result***  measurement result

## 7.2    Collective benchmarks

**Classes**

- struct MPIB_coll_container

**Typedefs**

- typedef int(∗ MPIB_measure_coll )(MPIB_coll_container ∗container, MPI_Comm comm, int root, int M, MPIB_precision precision, MPIB_result ∗result)
- typedef void(∗ MPIB_measure_coll_msgset )(MPIB_coll_container ∗container, MPI_Comm comm, int root, MPIB_msgset msgset, MPIB_precision precision, int ∗count, MPIB_result ∗∗results)

**Functions**

- int MPIB_measure_max (MPIB_coll_container ∗container, MPI_Comm comm, int root, int M, MPIB_precision precision, MPIB_result ∗result)
- void MPIB_root_timer_init (MPI_Comm comm, int reps)
- int MPIB_measure_root (MPIB_coll_container ∗container, MPI_Comm comm, int root, int M, MPIB_precision precision, MPIB_result ∗result)
- void MPIB_global_timer_init (MPI_Comm comm, int parallel, int reps)
- int MPIB_measure_global (MPIB_coll_container ∗container, MPI_Comm comm, int root, int M, MPIB_precision precision, MPIB_result ∗result)

### 7.2.1    Detailed Description

This module provides collective benchmarks based on root, max and global timing methods (see [1]).

### 7.2.2    Typedef Documentation

#### 7.2.2.1    typedef int(∗ MPIB_measure_coll)(MPIB_coll_container ∗container, MPI_Comm comm, int root, int M, MPIB_precision precision, MPIB_result ∗result)

Collective benchmark. Measures the execution time of collective operation for a given message size.

**Parameters**

> ***container***  communication operaion container
>
> ***comm***  communicator
>
> ***root***  root process

*M* message size

*precision* measurement precision

*result* measurement result

**7.2.2.2 typedef void(∗ MPIB_measure_coll_msgset)(MPIB_coll_container ∗container, MPI_Comm comm, int root, MPIB_msgset msgset, MPIB_precision precision, int ∗count, MPIB_result ∗∗results)**

Collective benchmark for multiple message sizes. Measures the execution time of collective operation for different message sizes.

**Parameters**

*container* communication operaion container

*comm* communicator

*root* root process

*msgset* message sizes

*precision* measurement precision

*count* the number of measurements performed (significant only at the root processor)

*results* array of measurement results (significant only at the root processor, allocated by this function and must be deallocated by user)

### 7.2.3 Function Documentation

#### 7.2.3.1 int MPIB_measure_max ( MPIB_coll_container ∗ *container,* MPI_Comm *comm,* int *root,* int *M,* MPIB_precision *precision,* MPIB_result ∗ *result* )

Measures the execution time of collective operation at all processes and finds a maximum. In the loop over repetitions:

- Synchronizes the processes by double barrier.

- Measures the execution time of collective operation at all processes.

- Finds maximum by allreduce.

- Performs statistical analysis.

#### 7.2.3.2 void MPIB_root_timer_init ( MPI_Comm *comm,* int *reps* )

Measures the average execution time of barrier at all processes in the communicator and sets up an internal global variable, which is used by MPIB_measure_root. Called by MPIB_measure_root. Can be called directly before measurements.

**Parameters**

*comm* communicator

*reps* number of repetitions (not precision - we suppose no pipeline effect with barrier)

### 7.2.3.3 int MPIB_measure_root ( MPIB_coll_container ∗ *container,* MPI_Comm *comm,* int *root,* int *M,* MPIB_precision *precision,* MPIB_result ∗ *result* )

Measures the execution time of collective operation at the root process. Reuses already obtained barrier time if the previous initialization was performed over the same communicator. Otherwise, initializes the root timer by calling MPIB_root_timer_init.

In the loop over repetitions:

- Synchronizes the processes by double barrier.

- Measures the execution time of collective operation and barrier confirmation at the root process.

- Subtracts the execution time of barrier.

- Performs statistical analysis at root.

Broadcasts the result.

### 7.2.3.4 void MPIB_global_timer_init ( MPI_Comm *comm,* int *parallel,* int *reps* )

Measures the offsets between local clocks of of all processes in the communicator and sets up an internal global variable, which is used by MPIB_measure_global. If MPI_WTIME_IS_GLOBAL (MPI global timer) is defined, the offsets are set to zero. Otherwise, the offsets are measured. Called by MPIB_measure_global. Can be called directly before measurements.

**Attention**

As the global variable is an array, the memory should be freed by

```
MPIB_global_timer_init(MPI_COMM_NULL, 0)
```

**Parameters**

*comm* communicator

*parallel* several non-overlapped point-to-point communications at the same time if non-zero

*reps* number of repetitions (not precision because series of ping-pongs are required - we cannot interrupt them by sending/receiving the result of the statistical estimation)

### 7.2.3.5 int MPIB_measure_global ( MPIB_coll_container ∗ *container,* MPI_Comm *comm,* int *root,* int *M,* MPIB_precision *precision,* MPIB_result ∗ *result* )

Measures the execution time of collective operation between processes using global time. Based on [4]. Reuses already obtained offsets between local clocks if the previous initialization was performed on the same communicator. Otherwise, initializes the global timer by calling MPIB_global_timer_init.

In the loop over repetitions:

- Synchronizes the processes by double barrier.

- Measures the moment of start at the root and the moment of finish at the rest of processes.

- Having substructed the offset, finds maximum by reducing to the root.

- Performs statistical analysis at root.

Broadcasts the result.

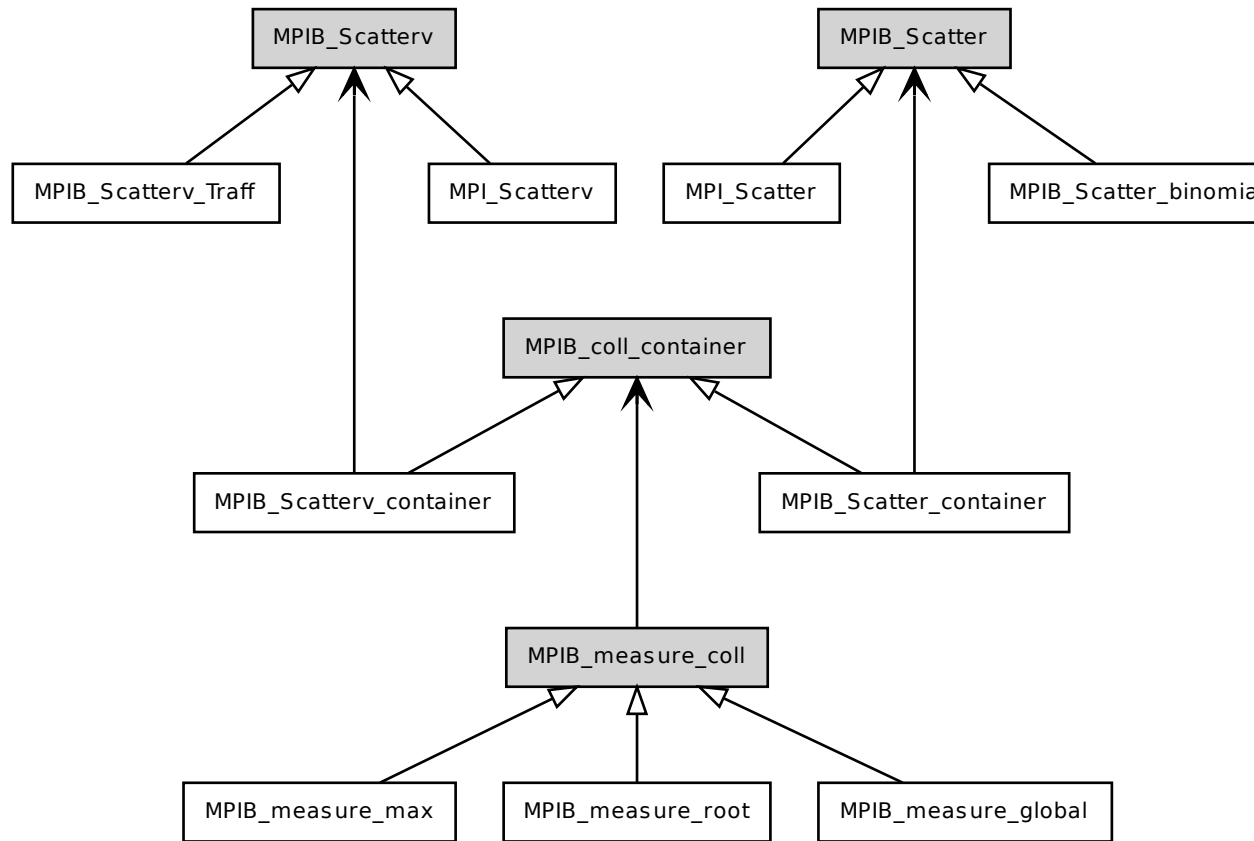## 7.3 Containers for collective communication operations

**Classes**

- class MPIB_buffer_container
- class MPIB_SG_container
- class MPIB_Scatter_container
- class MPIB_Gather_container
- class MPIB_Bcast_container
- class MPIB_Reduce_container
- class MPIB_Comm_dup_free_container
- class MPIB_SGv_container
- class MPIB_Scatterv_container
- class MPIB_Gatherv_container

**Functions**

- void MPIB_coll_container_free (MPIB_coll_container *container)
- MPIB_coll_container * MPIB_Scatter_container_alloc (MPIB_Scatter scatter)
- MPIB_coll_container * MPIB_Gather_container_alloc (MPIB_Gather gather)
- MPIB_coll_container * MPIB_Bcast_container_alloc (MPIB_Bcast bcast)
- MPIB_coll_container * MPIB_Reduce_container_alloc (MPIB_Reduce reduce)
- MPIB_coll_container * MPIB_Comm_dup_free_container_alloc ()
- MPIB_coll_container * MPIB_Scatterv_container_alloc (MPIB_Scatterv scatterv, const double *factors)
- MPIB_coll_container * MPIB_Gatherv_container_alloc (MPIB_Gatherv gatherv, const double *factors)

### 7.3.1 Detailed Description

Data structures describing collective communications to be measured, which are used as an argument of collective benchmark functions. Collective container allocators usually have an argument that points to the implementation of the collective communication operation, for example MPIB_Scatter_container_alloc has an argument `MPIB_Scatter`, pointer to a scatter implementation. This provides three-level extension:

The same container can be used in different benchmarks. Container can call different algorithms of the same collective operation.

### 7.3.2    Function Documentation

#### 7.3.2.1    void MPIB_coll_container_free ( MPIB_coll_container ∗ *container* )

Frees collective container

#### 7.3.2.2    MPIB_coll_container∗ MPIB_Scatter_container_alloc ( MPIB_Scatter *scatter* )

Allocates Scatter container

#### 7.3.2.3    MPIB_coll_container∗ MPIB_Gather_container_alloc ( MPIB_Gather *gather* )

Allocates Gather container

**7.3.2.4 MPIB_coll_container∗ MPIB_Bcast_container_alloc ( MPIB_Bcast *bcast* )**

Allocates Bcast container

**7.3.2.5 MPIB_coll_container∗ MPIB_Reduce_container_alloc ( MPIB_Reduce *reduce* )**

Allocates Reduce container

**7.3.2.6 MPIB_coll_container∗ MPIB_Comm_dup_free_container_alloc ( )**

Allocates MPI_Comm_dup-MPI_Comm_free container

**7.3.2.7 MPIB_coll_container∗ MPIB_Scatterv_container_alloc ( MPIB_Scatterv *scatterv,* const double ∗ *factors* )**

Allocates Scatterv container

**7.3.2.8 MPIB_coll_container∗ MPIB_Gatherv_container_alloc ( MPIB_Gatherv *gatherv,* const double ∗ *factors* )**

Allocates Gatherv container

## 7.4 Definitions of MPI communication operations

**Typedefs**

- typedef int(∗ MPIB_Scatter )(void ∗sendbuf, int sendcount, MPI_Datatype sendtype, void ∗recvbuf, int recvcount, MPI_Datatype recvtype, int root, MPI_Comm comm)
- typedef int(∗ MPIB_Gather )(void ∗sendbuf, int sendcount, MPI_Datatype sendtype, void ∗recvbuf, int recvcount, MPI_Datatype recvtype, int root, MPI_Comm comm)
- typedef int(∗ MPIB_Scatterv )(void ∗sendbuf, int ∗sendcounts, int ∗displs, MPI_Datatype sendtype, void ∗recvbuf, int recvcount, MPI_Datatype recvtype, int root, MPI_Comm comm)
- typedef int(∗ MPIB_Gatherv )(void ∗sendbuf, int sendcount, MPI_Datatype sendtype, void ∗recvbuf, int ∗recvcounts, int ∗displs, MPI_Datatype recvtype, int root, MPI_Comm comm)
- typedef int(∗ MPIB_Bcast )(void ∗buffer, int count, MPI_Datatype datatype, int root, MPI_Comm comm)
- typedef int(∗ MPIB_Reduce )(void ∗sendbuf, void ∗recvbuf, int count, MPI_Datatype datatype, MPI_Op op, int root, MPI_Comm comm)

### 7.4.1 Detailed Description

This module provides the definitions of MPI communication operations.

### 7.4.2 Typedef Documentation

**7.4.2.1 typedef int(∗ MPIB_Scatter)(void ∗sendbuf, int sendcount, MPI_Datatype sendtype, void ∗recvbuf, int recvcount, MPI_Datatype recvtype, int root, MPI_Comm comm)**

Scatter typedef

**7.4.2.2 typedef int(∗ MPIB_Gather)(void ∗sendbuf, int sendcount, MPI_Datatype sendtype, void ∗recvbuf, int recvcount, MPI_Datatype recvtype, int root, MPI_Comm comm)**

Gather typedef

**7.4.2.3 typedef int(∗ MPIB_Scatterv)(void ∗sendbuf, int ∗sendcounts, int ∗displs, MPI_Datatype sendtype, void ∗recvbuf, int recvcount, MPI_Datatype recvtype, int root, MPI_Comm comm)**

Scatterv typedef

**7.4.2.4 typedef int(∗ MPIB_Gatherv)(void ∗sendbuf, int sendcount, MPI_Datatype sendtype, void ∗recvbuf, int ∗recvcounts, int ∗displs, MPI_Datatype recvtype, int root, MPI_Comm comm)**

Gatherv typedef

**7.4.2.5 typedef int(∗ MPIB_Bcast)(void ∗buffer, int count, MPI_Datatype datatype, int root, MPI_Comm comm)**

Bcast typedef

**7.4.2.6 typedef int(∗ MPIB_Reduce)(void ∗sendbuf, void ∗recvbuf, int count, MPI_Datatype datatype, MPI_Op op, int root, MPI_Comm comm)**

Reduce typedef

## 7.5 Options for executables

### 7.5.1 Detailed Description

Since getopt cannot be reused, this module provides an interface for getopt, which can be used as follows:

```
X x;
Y y;
MPIB_getopt_x_default(&x);
MPIB_getopt_y_default(&y);
if (root == 0)
{
    char options[256] = "";
    strcat(options, MPIB_getopt_x_options());
    strcat(options, MPIB_getopt_y_options());
    while ((c = getopt(argc, argv, options)) != -1)
    {
        MPIB_getopt_x_optarg(c, &x);
        MPIB_getopt_y_optarg(c, &y);
    }
}
MPIB_getopt_x_bcast(&x, 0, MPI_COMM_WORLD);
MPIB_getopt_y_bcast(&y, 0, MPI_COMM_WORLD);
```

Options are divided into sets, which are managed by the following functions:

- **MPIB_getopt_X_default** - fills the parameters by default values

- **MPIB_getopt_X_options** - returns a string of the getopt options

- **MPIB_getopt_X_optarg** - fills the parameters by the getopt argument

- **MPIB_getopt_X_bcast** - broadcasts the parameters (parallel) where X stands for a name of the set.

Typical options for executables are as follows:

- **-h** help

- **-v** verbose

- **-l** *S* collectives shared library (required: a full path or relative to LD_LIBRARY_PATH)

- **-o** *S* suboptions (comma separated options for the shared library: `subopt1,subopt2=value2`)

- **-O** *S* collective operation defined in the shared library (required: the name must contain Bcast, Scatter, etc)

- **-t** *S* timing: max, root, global (default: max)

- **-m** *I* minimum message size (default: 0)

- **-M** *I* maximum message size (default: 204800)

- **-S** *I* stride between message sizes (default: 1024)

- **-d** *D* maximum relative difference: $0 < D < 1$ (default: 0.1)

- **-s** *I* minimum stride between message sizes (default: 64)

- **-n** *I* maximum number of message sizes (default: 100)

- **-p** 0/1 parallel p2p benchmarking (default: 1)

- **-r** *I* minimum number of repetitions (default: 5)

- **-R** *I* maximum number of repetitions (default: 100)

- **-c** *D* confidence level: $0 < D < 1$ (default: 0.95)

- **-e** *D* relative error: $0 < D < 1$ (default: 0.025)

- **-p** *I* parallel mode or not: 0 / 1 (default: 0)

- **-P** *I* type of P2P experiments: 0 is default, 1 is scatter/gather based P2P: 0 / 1 (default: 0). Note: If using -P 1 you have to be in the <install-dir>/bin directory

where:

- *S* - string

- *I* - integer

- *D* - double

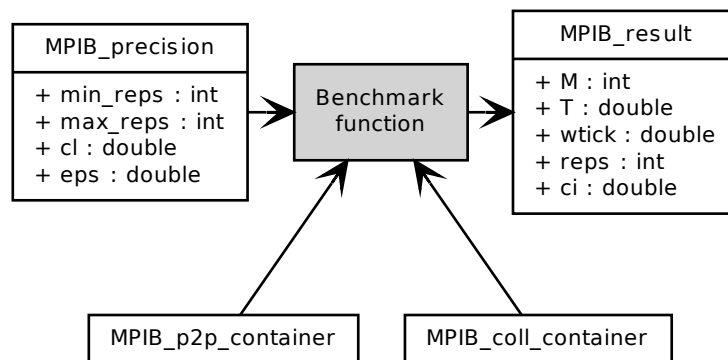## 7.6 Measurement: base data structures and functions

**Classes**

- struct MPIB_result
- struct MPIB_msgset
- struct MPIB_precision

---

**Functions**

- void MPIB_Comm (MPI_Comm comm, MPI_Comm ∗newcomm)
- double MPIB_diff (MPIB_result result, MPIB_result results[2])
- void MPIB_max_wtick (MPI_Comm comm, double ∗wtick)
- double MPIB_ci (double cl, int reps, double ∗T)

### 7.6.1    Detailed Description

This module provides base data structures and functions for measurement:
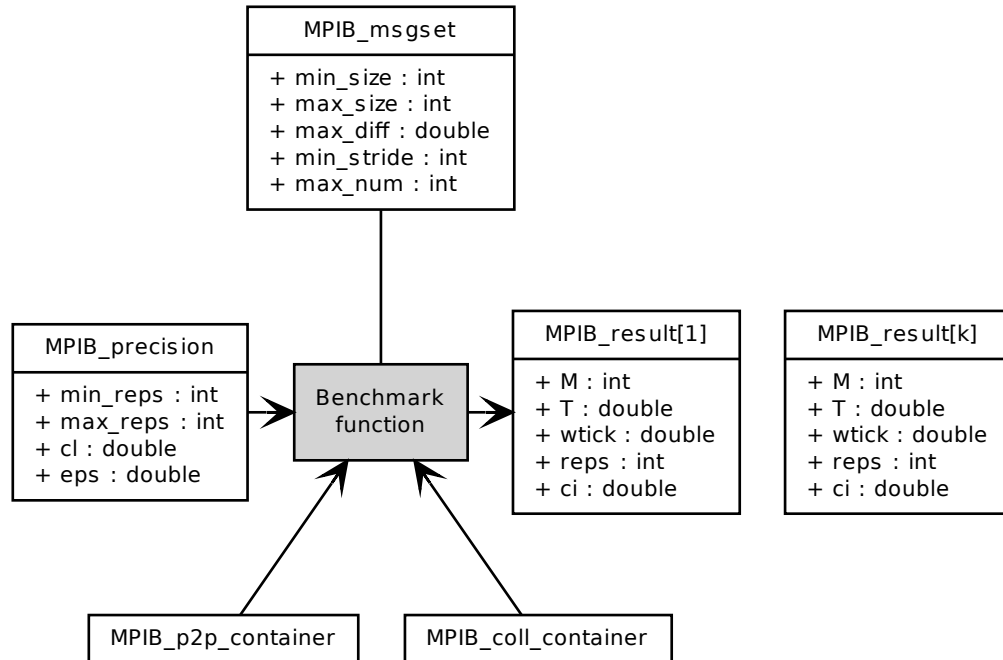


Benchmark input:

- The MPIB_precision data structure defines statistical precision of measurement. Statistical analysis is implemented with help of the GNU Scientific Library.

- The MPIB_p2p_container and MPIB_coll_container data structures encapsulate the communication operation to be measured.

Benchmark output:

- The MPIB_result data structure represents result of measurement and its reliability.

Most benchmark functions have a counterpart for measurement of execution time with a set of message sizes:

The counterpart has an extra input argument:

- The MPIB_msgset data structure defines the message sizes for which the measurements are to be performed. Message sizes are selected either regularly or adaptively.

The output of the counterpart is an array of MPIB_result.

### 7.6.2    Function Documentation

#### 7.6.2.1    void MPIB_Comm ( MPI_Comm *comm*, MPI_Comm ∗ *newcomm* )

Creates a copy of communicator that includes one process per processor.

#### 7.6.2.2    double MPIB_diff ( MPIB_result *result*, MPIB_result *results[2]* )

Compares the result of measurement with the linear model based on the results of two previous measurements. Required for adaptive selection of message sizes.

#### 7.6.2.3    void MPIB_max_wtick ( MPI_Comm *comm*, double ∗ *wtick* )

Returns a resolution of MPI_Wtime, maximum in the communicator. Result can be used to check the execution time measured at several processors (MPIB_measure_max, MPIB_measure_global): $T_{coll} < wtick_{max}$.

**Parameters**

> *comm*  MPI communicator
>
> *wtick*  a maximum resolution

### 7.6.2.4    double MPIB_ci ( double *cl,* int *reps,* double ∗ *T* )

Returns a confidence interval that contains the average execution time with a certain probability: $Pr(|\bar{T} - \mu| < ci) = cl$.

**Note**

> If communication operations in a series are isolated from each other, we can assume that the execution times form an independent sample from a normally distributed population, and use t distribution to estimate confidence interval.

**Parameters**

> *cl*  confidence level
>
> *reps*  number of measurements (should be $> 1$)
>
> *T*  array of reps measurement results

**Returns**

> confidence interval

## 7.7    Point-to-point benchmarks

**Classes**

- struct MPIB_p2p_container

**Defines**

- #define MPIB_C2(n) (n) ∗ ((n) - 1) / 2
- #define MPIB_IJ2INDEX(n, i, j) (2 ∗ (n) - ((i) < (j) ? (i) : (j)) - 1) ∗ (((i) < (j) ? (i) : (j))) / 2 + (((i) < (j) ? (j) : (i))) - (((i) < (j) ? (i) : (j))) - 1

**Functions**

- void MPIB_measure_p2p (MPIB_p2p_container ∗container, MPI_Comm comm, int measure, int mirror, int M, MPIB_precision precision, MPIB_result ∗result)
- void MPIB_measure_p2p_msgset (MPIB_p2p_container ∗container, MPI_Comm comm, int measure, int mirror, MPIB_msgset msgset, MPIB_precision precision, int ∗count, MPIB_result ∗∗results)
- void MPIB_measure_allp2p (MPIB_p2p_container ∗container, MPI_Comm comm, int parallel, int M, MPIB_precision precision, MPIB_result ∗results)

### 7.7.1    Detailed Description

This module provides the point-to-point benchmarks.

---

### 7.7.2    Define Documentation

#### 7.7.2.1    #define MPIB_C2(  *n*  ) (n) $*$ ((n) - 1) / 2

$$C_n^2$$

#### 7.7.2.2    #define MPIB_IJ2INDEX(  *n,    i,    j*  ) (2 $*$ (n) - ((i) < (j) ? (i) : (j))) - 1) $*$ (((i) < (j) ? (i) : (j))) / 2 + (((i) < (j) ? (j) : (i))) - (((i) < (j) ? (i) : (j))) - 1

For a symmetric square matrix stored in the array of $C_n^2$ elements, returns the index of the $(i, j)$ element, $i \neq j < n$: $\dfrac{(n-1) + (n-I)}{2} I + (J - I - 1), I = min(i, j), J = max(i, j)$

### 7.7.3    Function Documentation

#### 7.7.3.1    void MPIB_measure_p2p ( MPIB_p2p_container $*$ *container,* MPI_Comm *comm,* int *measure,* int *mirror,* int *M,* MPIB_precision *precision,* MPIB_result $*$ *result* )

Point-to-point benchmark. Estimates the execution time of the point-to-point communications between a pair of processors in the MPI communicator. Performs series of communication experiments to obtain reliable results.

**Parameters**

> *container*  communication operation container
>
> *comm*  communicator, number of nodes should be $\geq 2$
>
> *measure*  measure processor
>
> *mirror*  mirror processor
>
> *M*  message size
>
> *precision*  measurement precision
>
> *result*  measurement result (significant only at the measure processor)

#### 7.7.3.2    void MPIB_measure_p2p_msgset ( MPIB_p2p_container $*$ *container,* MPI_Comm *comm,* int *measure,* int *mirror,* MPIB_msgset *msgset,* MPIB_precision *precision,* int $*$ *count,* MPIB_result $**$ *results* )

Point-to-point benchmark. Estimates the execution time of the point-to-point communication between a pair of processors in the MPI communicator for different message sizes. Performs series of communication experiments to obtain reliable results.

**Parameters**

> *container*  communication operation container
>
> *comm*  communicator, number of nodes should be $\geq 2$
>
> *measure*  measure processor
>
> *mirror*  mirror processor
>
> *msgset*  message sizes
>
> *precision*  measurement precision
>
> *count*  the number of measurements performed (significant only at the measure processor)
>
> *results*  array of measurement results (significant only at the measure processor, allocated by this function and must be deallocated by user)

### 7.7.3.3  void MPIB_measure_allp2p ( MPIB_p2p_container ∗ *container,* MPI_Comm *comm,* int *parallel,* int *M,* MPIB_precision *precision,* MPIB_result ∗ *results* )

Point-to-point benchmark. Estimates the execution time of the point-to-point communications between all pairs of processors in the MPI communicator. Performs series of communication experiments to obtain reliable results.

**Parameters**

> *container*  communication operation container
>
> *comm*  communicator, number of nodes should be $\geq 2$
>
> *parallel*  several non-overlapped point-to-point communications at the same time if non-zero
>
> *M*  message size
>
> *precision*  measurement precision
>
> *results*  array of $C_n^2$ measurement results

## 7.8  Containers for point-to-point communication operations

**Functions**

- MPIB_p2p_container ∗ MPIB_roundtrip_container_alloc ()
- void MPIB_roundtrip_container_free (MPIB_p2p_container ∗container)

### 7.8.1  Detailed Description

Data structures describing point-to-point communications to be measured, which are used as an argument of the point-to-point benchmark function.

### 7.8.2  Function Documentation

#### 7.8.2.1  MPIB_p2p_container∗ MPIB_roundtrip_container_alloc (  )

Allocates a roundtrip container

#### 7.8.2.2  void MPIB_roundtrip_container_free ( MPIB_p2p_container ∗ *container* )

Frees the roundtrip container

## 7.9  API for shared libraries

**Defines**

- #define MPIB_COLLECTIVES_INITIALIZE "collectives_initialize"
- #define MPIB_COLLECTIVES_FINALIZE "collectives_finalize"

**Typedefs**

- typedef int(∗ MPIB_collectives_initialize )(MPI_Comm comm, char ∗subopts)
- typedef int(∗ MPIB_collectives_finalize )(MPI_Comm comm)

### 7.9.1   Define Documentation

#### 7.9.1.1   #define MPIB_COLLECTIVES_INITIALIZE "collectives_initialize"

Name of the initialization function in the collectives shared library

#### 7.9.1.2   #define MPIB_COLLECTIVES_FINALIZE "collectives_finalize"

Name of the finalization function in the collectives shared library

### 7.9.2   Typedef Documentation

#### 7.9.2.1   typedef int(∗ MPIB_collectives_initialize)(MPI_Comm comm, char ∗subopts)

Type of functions that parse suboptions and allocate global variables at all processors. Implement the function `collectives_initialize` of this type in the collectives shared library (optional).

**Parameters**

>   *comm*   MPI communicator
>   *subopts*   suboptions

#### 7.9.2.2   typedef int(∗ MPIB_collectives_finalize)(MPI_Comm comm)

Type of functions that free global variables at all processors. Implement the function `collectives_finalize` of this type in the collectives shared library (optional).

**Parameters**

>   *comm*   MPI communicator

## 7.10   Wrappers for p2p operations

**Functions**

- int p2p_init (MPI_Comm comm, int proc_num_spawned)
- int p2p_finalize ()
- int MPIB_Send (void ∗buf, int count, MPI_Datatype datatype, int dest, int tag, MPI_Comm comm)
- int MPIB_Isend (void ∗buf, int count, MPI_Datatype datatype, int dest, int tag, MPI_Comm comm, MPI_Request ∗request)
- int MPIB_Recv (void ∗buf, int count, MPI_Datatype datatype, int source, int tag, MPI_Comm comm, MPI_Status ∗status)
- int MPIB_Irecv (void ∗buf, int count, MPI_Datatype datatype, int source, int tag, MPI_Comm comm, MPI_Request ∗request)
- int MPIB_Waitall (int count, MPI_Request ∗array_of_requests, MPI_Status ∗array_of_statuses)

### 7.10.1  Function Documentation

#### 7.10.1.1  int p2p_init ( MPI_Comm *comm,* int *proc_num_spawned* )

This init routine is called by all processes in the original MPI communicator. Process 0 spawns a number of additional MPI processes everywhere (see Spawned process for scatter-gather point-to-point communication) for the scatter-gather-based point-to-point communication. After the routine, the global communicator intracomm contains all processes (spawning and spawned)

#### 7.10.1.2  int p2p_finalize (  )

This finalize routine is used by process 0 to send a terminate signal to all spawned processes

#### 7.10.1.3  int MPIB_Send ( void ∗ *buf,* int *count,* MPI_Datatype *datatype,* int *dest,* int *tag,* MPI_Comm *comm* )

Wrapper send function for standard or modified operation

#### 7.10.1.4  int MPIB_Isend ( void ∗ *buf,* int *count,* MPI_Datatype *datatype,* int *dest,* int *tag,* MPI_Comm *comm,* MPI_Request ∗ *request* )

Wrapper non-blocking send function for standard or modified operation

#### 7.10.1.5  int MPIB_Recv ( void ∗ *buf,* int *count,* MPI_Datatype *datatype,* int *source,* int *tag,* MPI_Comm *comm,* MPI_Status ∗ *status* )

Wrapper receive function for standard or modified operation

#### 7.10.1.6  int MPIB_Irecv ( void ∗ *buf,* int *count,* MPI_Datatype *datatype,* int *source,* int *tag,* MPI_Comm *comm,* MPI_Request ∗ *request* )

Wrapper non-blocking receive function for standard or modified operation

#### 7.10.1.7  int MPIB_Waitall ( int *count,* MPI_Request ∗ *array_of_requests,* MPI_Status ∗ *array_of_statuses* )

Wrapper waitall function for standard or modified operation

## 7.11  Basic algorithms of MPI collective operations

**Functions**

- int MPIB_Scatter_flat_nb (void ∗sendbuf, int sendcount, MPI_Datatype sendtype, void ∗recvbuf, int recvcount, MPI_Datatype recvtype, int root, MPI_Comm comm)
- int MPIB_Gather_flat_nb (void ∗sendbuf, int sendcount, MPI_Datatype sendtype, void ∗recvbuf, int recvcount, MPI_Datatype recvtype, int root, MPI_Comm comm)
- int MPIB_Scatter_flat_rsend (void ∗sendbuf, int sendcount, MPI_Datatype sendtype, void ∗recvbuf, int recvcount, MPI_Datatype recvtype, int root, MPI_Comm comm)

- int MPIB_Gather_flat_rsend (void ∗sendbuf, int sendcount, MPI_Datatype sendtype, void ∗recvbuf, int recvcount, MPI_Datatype recvtype, int root, MPI_Comm comm)
- int MPIB_Gatherv_flat_sync (void ∗sendbuf, int sendcount, MPI_Datatype sendtype, void ∗recvbuf, int ∗recvcounts, int ∗displs, MPI_Datatype recvtype, int root, MPI_Comm comm)
- int MPIB_Scatterv_flat (void ∗sendbuf, int ∗sendcounts, int ∗displs, MPI_Datatype sendtype, void ∗recvbuf, int recvcount, MPI_Datatype recvtype, int root, MPI_Comm comm)
- int MPIB_Gatherv_flat (void ∗sendbuf, int sendcount, MPI_Datatype sendtype, void ∗recvbuf, int ∗recvcounts, int ∗displs, MPI_Datatype recvtype, int root, MPI_Comm comm)
- int MPIB_Scatterv_sorted_flat_asc (void ∗sendbuf, int ∗sendcounts, int ∗displs, MPI_Datatype sendtype, void ∗recvbuf, int recvcount, MPI_Datatype recvtype, int root, MPI_Comm comm)
- int MPIB_Gatherv_sorted_flat_asc (void ∗sendbuf, int sendcount, MPI_Datatype sendtype, void ∗recvbuf, int ∗recvcounts, int ∗displs, MPI_Datatype recvtype, int root, MPI_Comm comm)
- int MPIB_Scatterv_sorted_flat_dsc (void ∗sendbuf, int ∗sendcounts, int ∗displs, MPI_Datatype sendtype, void ∗recvbuf, int recvcount, MPI_Datatype recvtype, int root, MPI_Comm comm)
- int MPIB_Gatherv_sorted_flat_dsc (void ∗sendbuf, int sendcount, MPI_Datatype sendtype, void ∗recvbuf, int ∗recvcounts, int ∗displs, MPI_Datatype recvtype, int root, MPI_Comm comm)

### 7.11.1    Detailed Description

This module provides basic, mostly flat-tree, algorithms of MPI collective operations.

### 7.11.2    Function Documentation

#### 7.11.2.1    int MPIB_Scatter_flat_nb ( void ∗ *sendbuf,* int *sendcount,* MPI_Datatype *sendtype,* void ∗ *recvbuf,* int *recvcount,* MPI_Datatype *recvtype,* int *root,* MPI_Comm *comm* )

Flat-tree scatter using non-blocking standard or modified p2p

#### 7.11.2.2    int MPIB_Gather_flat_nb ( void ∗ *sendbuf,* int *sendcount,* MPI_Datatype *sendtype,* void ∗ *recvbuf,* int *recvcount,* MPI_Datatype *recvtype,* int *root,* MPI_Comm *comm* )

Flat-tree gather using standard or modified p2p

#### 7.11.2.3    int MPIB_Scatter_flat_rsend ( void ∗ *sendbuf,* int *sendcount,* MPI_Datatype *sendtype,* void ∗ *recvbuf,* int *recvcount,* MPI_Datatype *recvtype,* int *root,* MPI_Comm *comm* )

Flat-tree scatter using MPI_Rsend => eager protocol

#### 7.11.2.4    int MPIB_Gather_flat_rsend ( void ∗ *sendbuf,* int *sendcount,* MPI_Datatype *sendtype,* void ∗ *recvbuf,* int *recvcount,* MPI_Datatype *recvtype,* int *root,* MPI_Comm *comm* )

Flat-tree gather using MPI_Rsend => eager protocol

#### 7.11.2.5    int MPIB_Gatherv_flat_sync ( void ∗ *sendbuf,* int *sendcount,* MPI_Datatype *sendtype,* void ∗ *recvbuf,* int ∗ *recvcounts,* int ∗ *displs,* MPI_Datatype *recvtype,* int *root,* MPI_Comm *comm* )

Flat-tree gatherv with sync, taken from OMPI trunk

**7.11.2.6   int MPIB_Scatterv_flat ( void ∗ *sendbuf,* int ∗ *sendcounts,* int ∗ *displs,* MPI_Datatype *sendtype,* void ∗ *recvbuf,* int *recvcount,* MPI_Datatype *recvtype,* int *root,* MPI_Comm *comm* )**

Flat-tree scatterv

**7.11.2.7   int MPIB_Gatherv_flat ( void ∗ *sendbuf,* int *sendcount,* MPI_Datatype *sendtype,* void ∗ *recvbuf,* int ∗ *recvcounts,* int ∗ *displs,* MPI_Datatype *recvtype,* int *root,* MPI_Comm *comm* )**

Flat-tree gatherv

**7.11.2.8   int MPIB_Scatterv_sorted_flat_asc ( void ∗ *sendbuf,* int ∗ *sendcounts,* int ∗ *displs,* MPI_Datatype *sendtype,* void ∗ *recvbuf,* int *recvcount,* MPI_Datatype *recvtype,* int *root,* MPI_Comm *comm* )**

Flat-tree scatterv with sendcounts sorted in ascending order

**7.11.2.9   int MPIB_Gatherv_sorted_flat_asc ( void ∗ *sendbuf,* int *sendcount,* MPI_Datatype *sendtype,* void ∗ *recvbuf,* int ∗ *recvcounts,* int ∗ *displs,* MPI_Datatype *recvtype,* int *root,* MPI_Comm *comm* )**

Flat-tree gatherv with recvcounts sorted in ascending order

**7.11.2.10   int MPIB_Scatterv_sorted_flat_dsc ( void ∗ *sendbuf,* int ∗ *sendcounts,* int ∗ *displs,* MPI_Datatype *sendtype,* void ∗ *recvbuf,* int *recvcount,* MPI_Datatype *recvtype,* int *root,* MPI_Comm *comm* )**

Flat-tree scatterv with sendcounts sorted in descending order

**7.11.2.11   int MPIB_Gatherv_sorted_flat_dsc ( void ∗ *sendbuf,* int *sendcount,* MPI_Datatype *sendtype,* void ∗ *recvbuf,* int ∗ *recvcounts,* int ∗ *displs,* MPI_Datatype *recvtype,* int *root,* MPI_Comm *comm* )**

Flat-tree gatherv with recvcounts sorted in descending order

## 7.12   Tree-based algorithms of MPI collective operations

**Functions**

- int MPIB_Bcast_binomial (void ∗buffer, int count, MPI_Datatype datatype, int root, MPI_Comm comm)
- int MPIB_Reduce_binomial (void ∗sendbuf, void ∗recvbuf, int count, MPI_Datatype datatype, MPI_Op op, int root, MPI_Comm comm)
- int MPIB_Scatter_binomial (void ∗sendbuf, int sendcount, MPI_Datatype sendtype, void ∗recvbuf, int recvcount, MPI_Datatype recvtype, int root, MPI_Comm comm)
- int MPIB_Gather_binomial (void ∗sendbuf, int sendcount, MPI_Datatype sendtype, void ∗recvbuf, int recvcount, MPI_Datatype recvtype, int root, MPI_Comm comm)

- int MPIB_Scatterv_binomial (void ∗sendbuf, int ∗sendcounts, int ∗displs, MPI_Datatype sendtype, void ∗recvbuf, int recvcount, MPI_Datatype recvtype, int root, MPI_Comm comm)
- int MPIB_Gatherv_binomial (void ∗sendbuf, int sendcount, MPI_Datatype sendtype, void ∗recvbuf, int ∗recvcounts, int ∗displs, MPI_Datatype recvtype, int root, MPI_Comm comm)
- int MPIB_Scatterv_sorted_binomial_asc (void ∗sendbuf, int ∗sendcounts, int ∗displs, MPI_Datatype sendtype, void ∗recvbuf, int recvcount, MPI_Datatype recvtype, int root, MPI_Comm comm)
- int MPIB_Gatherv_sorted_binomial_asc (void ∗sendbuf, int sendcount, MPI_Datatype sendtype, void ∗recvbuf, int ∗recvcounts, int ∗displs, MPI_Datatype recvtype, int root, MPI_Comm comm)
- int MPIB_Scatterv_sorted_binomial_dsc (void ∗sendbuf, int ∗sendcounts, int ∗displs, MPI_Datatype sendtype, void ∗recvbuf, int recvcount, MPI_Datatype recvtype, int root, MPI_Comm comm)
- int MPIB_Gatherv_sorted_binomial_dsc (void ∗sendbuf, int sendcount, MPI_Datatype sendtype, void ∗recvbuf, int ∗recvcounts, int ∗displs, MPI_Datatype recvtype, int root, MPI_Comm comm)
- int MPIB_Scatterv_Traff (void ∗sendbuf, int ∗sendcounts, int ∗displs, MPI_Datatype sendtype, void ∗recvbuf, int recvcount, MPI_Datatype recvtype, int root, MPI_Comm comm)
- int MPIB_Gatherv_Traff (void ∗sendbuf, int sendcount, MPI_Datatype sendtype, void ∗recvbuf, int ∗recvcounts, int ∗displs, MPI_Datatype recvtype, int root, MPI_Comm comm)

### 7.12.1    Detailed Description

This module provides tree-based algorithms of MPI collective operations. For rapid implementation, the Graph Boost C++ library is used. The sources of tree algorithms of collectives must be compiled by C++.

A tree-based algorithm of a collective operation X consists of the following components:

- **Communication tree** MPIB::comm_tree

- **Function template for the tree-based algorithm** MPIB_X_tree_algorithm (described in Function templates for tree-based algorithms of MPI collective operations)

- **Communication tree builder** Y_builder, a C++ class that builds a communication tree. (described in MPIB::BRSG and MPIB::SGv)

The tree-based implementation looks as follows:

```
extern "C" int MPIB_X_Y(standard args) {
    return MPIB_X_tree_algorithm(Y_builder(), standard args);
}
```

For example, MPIB_Scatter_binomial, a binomial scatter.

This approach provides flexibility:

- Tree-based algorithm of the collective operation is a basis for implementation of algorithms with communication trees of different shapes (for example, binomial and binary algorithms of scatter use the general tree-based algorithm of scatter).

- Communication tree is a universal interface between tree-based algorithms and tree builders. It can be reused for different communication operations (for example, scatter/gather and bcast/reduce use the same simple communication tree).

- Tree and communication tree builders can be reused per tree shape (for example, the same binomial communication tree builder is used in binomial scatter/gather and bcast/reduce).

### 7.12.2    Function Documentation

#### 7.12.2.1    int MPIB_Bcast_binomial ( void ∗ *buffer,* int *count,* MPI_Datatype *datatype,* int *root,* MPI_Comm *comm* )

Binomial Bcast

#### 7.12.2.2    int MPIB_Reduce_binomial ( void ∗ *sendbuf,* void ∗ *recvbuf,* int *count,* MPI_Datatype *datatype,* MPI_Op *op,* int *root,* MPI_Comm *comm* )

Binomial Reduce

**7.12.2.3    int MPIB_Scatter_binomial ( void ∗ *sendbuf,* int *sendcount,* MPI_Datatype *sendtype,* void ∗ *recvbuf,* int *recvcount,* MPI_Datatype *recvtype,* int *root,* MPI_Comm *comm* )**

Binomial Scatter

**7.12.2.4    int MPIB_Gather_binomial ( void ∗ *sendbuf,* int *sendcount,* MPI_Datatype *sendtype,* void ∗ *recvbuf,* int *recvcount,* MPI_Datatype *recvtype,* int *root,* MPI_Comm *comm* )**

Binomial Gather

**7.12.2.5    int MPIB_Scatterv_binomial ( void ∗ *sendbuf,* int ∗ *sendcounts,* int ∗ *displs,* MPI_Datatype *sendtype,* void ∗ *recvbuf,* int *recvcount,* MPI_Datatype *recvtype,* int *root,* MPI_Comm *comm* )**

Binomial scatterv

**7.12.2.6    int MPIB_Gatherv_binomial ( void ∗ *sendbuf,* int *sendcount,* MPI_Datatype *sendtype,* void ∗ *recvbuf,* int ∗ *recvcounts,* int ∗ *displs,* MPI_Datatype *recvtype,* int *root,* MPI_Comm *comm* )**

Binomial gatherv.

**7.12.2.7    int MPIB_Scatterv_sorted_binomial_asc ( void ∗ *sendbuf,* int ∗ *sendcounts,* int ∗ *displs,* MPI_Datatype *sendtype,* void ∗ *recvbuf,* int *recvcount,* MPI_Datatype *recvtype,* int *root,* MPI_Comm *comm* )**

Binomial scatterv with children sorted by counts in ascending order.

**7.12.2.8    int MPIB_Gatherv_sorted_binomial_asc ( void ∗ *sendbuf,* int *sendcount,* MPI_Datatype *sendtype,* void ∗ *recvbuf,* int ∗ *recvcounts,* int ∗ *displs,* MPI_Datatype *recvtype,* int *root,* MPI_Comm *comm* )**

Binomial gatherv with children sorted by counts in ascending order.

**7.12.2.9    int MPIB_Scatterv_sorted_binomial_dsc ( void ∗ *sendbuf,* int ∗ *sendcounts,* int ∗ *displs,* MPI_Datatype *sendtype,* void ∗ *recvbuf,* int *recvcount,* MPI_Datatype *recvtype,* int *root,* MPI_Comm *comm* )**

Binomial scatterv with children sorted by counts in descending order.

**7.12.2.10    int MPIB_Gatherv_sorted_binomial_dsc ( void ∗ *sendbuf,* int *sendcount,* MPI_Datatype *sendtype,* void ∗ *recvbuf,* int ∗ *recvcounts,* int ∗ *displs,* MPI_Datatype *recvtype,* int *root,* MPI_Comm *comm* )**

Binomial gatherv with children sorted by counts in descending order.

**7.12.2.11    int MPIB_Scatterv_Traff ( void ∗ *sendbuf,* int ∗ *sendcounts,* int ∗ *displs,* MPI_Datatype *sendtype,* void ∗ *recvbuf,* int *recvcount,* MPI_Datatype *recvtype,* int *root,* MPI_Comm *comm* )**

<div align="right">Traff algorithm of Scatterv. Based on [3].</div>

**7.12.2.12    int MPIB_Gatherv_Traff ( void ∗ *sendbuf,* int *sendcount,* MPI_Datatype *sendtype,* void ∗ *recvbuf,* int ∗ *recvcounts,* int ∗ *displs,* MPI_Datatype *recvtype,* int *root,* MPI_Comm *comm* )**

<div align="right">Traff algorithm of Gatherv. Based on [3].</div>

## 7.13    Function templates for tree-based algorithms of MPI collective operations

**Functions**

- template<typename Builder >
  int MPIB_Bcast_tree_algorithm (Builder builder, MPIB_child_traverse_order order, void ∗buffer, int count, MPI_Datatype datatype, int root, MPI_Comm comm)
- template<typename Builder >
  int MPIB_Reduce_tree_algorithm (Builder builder, MPIB_child_traverse_order order, void ∗sendbuf, void ∗recvbuf, int count, MPI_Datatype datatype, MPI_Op op, int root, MPI_Comm comm)
- template<typename Builder >
  int MPIB_Scatter_tree_algorithm (Builder builder, MPIB_child_traverse_order order, void ∗sendbuf, int sendcount, MPI_Datatype sendtype, void ∗recvbuf, int recvcount, MPI_Datatype recvtype, int root, MPI_Comm comm)
- template<typename Builder >
  int MPIB_Gather_tree_algorithm (Builder builder, MPIB_child_traverse_order order, void ∗sendbuf, int sendcount, MPI_Datatype sendtype, void ∗recvbuf, int recvcount, MPI_Datatype recvtype, int root, MPI_Comm comm)
- template<typename Builder >
  int MPIB_Scatterv_tree_algorithm (Builder builder, MPIB_child_traverse_order order, void ∗sendbuf, int ∗sendcounts, int ∗displs, MPI_Datatype sendtype, void ∗recvbuf, int recvcount, MPI_-Datatype recvtype, int root, MPI_Comm _comm)
- template<typename Builder >
  int MPIB_Gatherv_tree_algorithm (Builder builder, MPIB_child_traverse_order order, void ∗sendbuf, int sendcount, MPI_Datatype sendtype, void ∗recvbuf, int ∗recvcounts, int ∗displs, MPI_-Datatype recvtype, int root, MPI_Comm _comm)

### 7.13.1    Detailed Description

In addition to the standard arguments of an MPI collective operation, a function template has the communication tree builder and the order arguments:

```
template <typename Builder>
MPIB_X_tree_algorithm(Builder builder, order args, standard args);
```

For example, MPIB_Scatter_tree_algorithm, a base tree algorithm of scatter. In base tree-based algorithm, all point-to-point communications are performed over the communication tree built by the builder in order given by the order argument.

---

Usually, the communication tree is built at all processors independently. If the communication tree can be built only at a designated processor, it must then be sent to other processes along with the data. The Serialization Boost C++ library is used for serialization/deserialization of the communication tree/subtrees in such tree-based algorithms:

```
#include <boost/graph/adj_list_serialize.hpp>
#include <boost/archive/binary_oarchive.hpp>
#include <boost/archive/binary_iarchive.hpp>
#include <sstream>

Graph graph;

if (rank == root) {
    ostringstream oss;
    archive::binary_oarchive ar(oss);
    ar << graph;
    int length = oss.str().length();
    MPI_Send((void*)oss.str().c_str(), length, MPI_CHAR, dest, 1, comm);
}

if (rank == dest) {
    MPI_Status status;
    MPI_Probe(root, 1, comm, &status);
    int length;
    MPI_Get_count(&status, MPI_CHAR, &length);
    buffer = (char*)malloc(sizeof(char) * length);
    MPI_Recv(buffer, length, MPI_CHAR, root, 1, comm, MPI_STATUS_IGNORE);
    istringstream iss(string(buffer, length));
    archive::binary_iarchive ar(iss);
    ar >> graph;
    free(buffer);
}
```

The internal part of the tree-based implementation includes the following auxiliaries united in namespaces in order to avoid duplicates:

- **Tree visitors** traverse communication tree, for example, in order to assemble the data buffer to send or receive:

  ```
  class Visitor {
  public:
      Visitor(args) {...}
      void preorder(Vertex vertex, Tree& tree) {...}
      void inorder(Vertex vertex, Tree& tree) {...}
      void postorder(Vertex vertex, Tree& tree) {...}
  };
  ```

  **Note**

  There may be many pointer or reference arguments in the visitor's constructor because visitors are copied by value.

- **Property writers** print vertex, edge and graph properties during the output of the communication tree:

  ```
  class Vertex_writer {
  public:
      void operator()(std::ostream& out, const Vertex& v) const {
          out << "[label=\"" << ... << "\"]";
      }
  };

  class Edge_writer {
  public:
      void operator()(std::ostream& out, const Edge& e) const {
  ```

```
        out << "[label=\"" << ... << "\"]";
    }
};

class Graph_writer {
public:
    void operator()(std::ostream& out) const {
        out << "graph [...]\n";
        out << "node [...]\n";
        out << "edge [...]\n";
    }
};

write_graphviz(cout, graph, Vertex_writer(), Edge_writer(), Graph_writer());
```

**Note**

Default writers are called when the last three arguments omitted.

### 7.13.2    Function Documentation

#### 7.13.2.1    template<typename Builder > int MPIB_Bcast_tree_algorithm ( Builder *builder,* MPIB_child_traverse_order *order,* void ∗ *buffer,* int *count,* MPI_Datatype *datatype,* int *root,* MPI_Comm *comm* )

Base tree algorithm of bcast

#### 7.13.2.2    template<typename Builder > int MPIB_Reduce_tree_algorithm ( Builder *builder,* MPIB_child_traverse_order *order,* void ∗ *sendbuf,* void ∗ *recvbuf,* int *count,* MPI_Datatype *datatype,* MPI_Op *op,* int *root,* MPI_Comm *comm* )

Base tree algorithm of reduce.

**Note**

Does not perform MPI operations but allocates memory for subproduct. MPI internals should be used. For example, Open MPI:

```
#ifdef HAVE_OPENMPI_OMPI_OP_OP_H
#include <openmpi/ompi/op/op.h>
#endif
...
int MPIB_Reduce_tree_algorithm(...) {
    ...
#ifdef HAVE_OPENMPI_OMPI_OP_OP_H
    ompi_op_reduce(op, buffer, sendbuf, count, datatype);
#endif
    ...
}
```

TODO: implement MPI operation in the reduce tree algorithm

#### 7.13.2.3    template<typename Builder > int MPIB_Scatter_tree_algorithm ( Builder *builder,* MPIB_child_traverse_order *order,* void ∗ *sendbuf,* int *sendcount,* MPI_Datatype *sendtype,* void ∗ *recvbuf,* int *recvcount,* MPI_Datatype *recvtype,* int *root,* MPI_Comm *comm* )

Base tree algorithm of scatter

**7.13.2.4** **template**$<$**typename Builder** $>$ **int MPIB_Gather_tree_algorithm ( Builder** *builder,* **MPIB_child_traverse_order** *order,* **void** $*$ *sendbuf,* **int** *sendcount,* **MPI_Datatype** *sendtype,* **void** $*$ *recvbuf,* **int** *recvcount,* **MPI_Datatype** *recvtype,* **int** *root,* **MPI_Comm** *comm* **)**

Base tree algorithm of gather

**7.13.2.5** **template**$<$**typename Builder** $>$ **int MPIB_Scatterv_tree_algorithm ( Builder** *builder,* **MPIB_child_traverse_order** *order,* **void** $*$ *sendbuf,* **int** $*$ *sendcounts,* **int** $*$ *displs,* **MPI_Datatype** *sendtype,* **void** $*$ *recvbuf,* **int** *recvcount,* **MPI_Datatype** *recvtype,* **int** *root,* **MPI_Comm** *_comm* **)**

Base tree algorithm of scatterv

**7.13.2.6** **template**$<$**typename Builder** $>$ **int MPIB_Gatherv_tree_algorithm ( Builder** *builder,* **MPIB_child_traverse_order** *order,* **void** $*$ *sendbuf,* **int** *sendcount,* **MPI_Datatype** *sendtype,* **void** $*$ *recvbuf,* **int** $*$ *recvcounts,* **int** $*$ *displs,* **MPI_Datatype** *recvtype,* **int** *root,* **MPI_Comm** *_comm* **)**

Base tree algorithm of gatherv

# 8 Namespace Documentation

## 8.1 MPIB::BRSG Namespace Reference

**Classes**

- class Binomial_builder

### 8.1.1 Detailed Description

The communication tree builders for tree-based algorithms of bcast/reduce/scatter/gather must implement a function:

```
void build(int size, int root, int rank, int count,
    Graph& g, Vertex& r, Vertex& u, Vertex& v);
```

**Parameters**

*size* communicator size

*root* root process

*rank* current process

*count* message size in bytes (count $*$ extent)

*g* graph

*r* root vertex (introduced to avoid search)

*u* parent of v (introduced to avoid search)

*v* vertex for the current process (introduced to avoid search)

## 8.2   MPIB::comm_tree Namespace Reference

**Typedefs**

- typedef adjacency_list< listS, listS, directedS, property< vertex_index_t, int >, no_property, no_-property, listS > Graph
- typedef graph_traits< Graph >::vertex_descriptor Vertex
- typedef graph_traits< Graph >::edge_descriptor Edge
- typedef graph_traits< Graph >::vertex_iterator Vertex_iterator
- typedef graph_traits< Graph >::adjacency_iterator Adjacency_iterator
- typedef graph_as_tree< Graph, iterator_property_map< vector< Vertex >::iterator, property_map< Graph, vertex_index_t >::type > > Tree

### 8.2.1   Detailed Description

**Communication tree** is a C++ namespace that contains a set of data structures for tree-based algorithms of collectives. All data types are short aliases for the Boost Graph data structures.

### 8.2.2   Typedef Documentation

#### 8.2.2.1   typedef adjacency_list<listS, listS, directedS, property<vertex_index_t, int>, no_property, no_property, listS> MPIB::comm_tree::Graph

Directed graph (vertex_index = MPI rank)

**Note**

The vertex index property is used for indexing and must starts from zero.
The adjacent_vertices method provides only the parent->child relation. For backward connection, use the graph as tree wrapper Tree.

#### 8.2.2.2   typedef graph_traits<Graph>::vertex_descriptor MPIB::comm_tree::Vertex

Vertex

#### 8.2.2.3   typedef graph_traits<Graph>::edge_descriptor MPIB::comm_tree::Edge

Edge

#### 8.2.2.4   typedef graph_traits<Graph>::vertex_iterator MPIB::comm_tree::Vertex_iterator

Vertex iterator

#### 8.2.2.5   typedef graph_traits<Graph>::adjacency_iterator MPIB::comm_tree::Adjacency_-iterator

Adjacency iterator

**8.2.2.6   typedef graph_as_tree**<**Graph, iterator_property_map**<**vector**<**Vertex**>**::iterator, property_map**<**Graph, vertex_index_t**>**::type**> > **MPIB::comm_tree::Tree**

Graph as tree wrapper. Provides access to the root, parent and children nodes.

**Note**

> We cannot use the `Tree` data structure directly, instead of `Graph`, because of the lack of empty/copy constructors.
>
> Access to parent nodes requires some computation (traversing the tree) - avoid this if possible.

## 8.3   MPIB::SG Namespace Reference

**Classes**

- class Assembler
- class Indexer
- class Edge_writer

### 8.3.1   Detailed Description

Auxiliaries for scatter/gather tree algorithms

## 8.4   MPIB::SGv Namespace Reference

**Classes**

- class Assembler
- class Indexer
- class Vertex_writer
- class Edge_writer
- class Binomial_builder
- class Sorted_binomial_builder
- class Traff_builder

**Variables**

- int tag = 0

### 8.4.1   Detailed Description

Auxiliaries for scatterv/gatherv tree algorithms

Communication tree builders for scatterv/gatherv. Must implement a function:

```
void build(int size, int root, int rank, int* counts,
    Graph& g, Vertex& r, Vertex& u, Vertex& v);
```

**Parameters**

> *size*   communicator size

---

    *root*  root process

    *rank*  current process

    *counts*  array of message sizes in bytes (counts ∗ extent) (number of elements = communicator size)

    *g*  graph

    *r*  root vertex (introduced to avoid search)

    *u*  parent of v (introduced to avoid search)

    *v*  vertex for the current process (introduced to avoid search)

**See also**

    Test for the tree-based scatterv/gatherv

### 8.4.2 Variable Documentation

#### 8.4.2.1 int MPIB::SGv::tag = 0

                                     tag for scatterv/gatherv - not safe, incremental on all processors

# 9   Class Documentation

## 9.1   MPIB::SG::Assembler Class Reference

### 9.1.1   Detailed Description

Finds the data to send (scatter) or recv (gather) from the local buffer (rank2index). Default copy constructors are used. Arrays must be preallocated (max = comm size).

The documentation for this class was generated from the following file:

- collectives/sg_tree_algorithms.hpp

## 9.2   MPIB::SGv::Assembler Class Reference

### 9.2.1   Detailed Description

Finds the data to send (scatterv) or recv (gatherv) from the local buffer (rank2index, counts, displs). Default copy constructors are used. Arrays must be preallocated (max = comm size).

The documentation for this class was generated from the following file:

- collectives/sgv_tree_algorithms.hpp

## 9.3   MPIB::BRSG::Binomial_builder Class Reference

### 9.3.1   Detailed Description

Binomial tree builder for bcast/reduce/scatter/gather. The largest subtree on the right. If root <> 0, processes' ranks will be cyclically shifted to 0.

The documentation for this class was generated from the following file:

- collectives/brsg_tree_builders.hpp

## 9.4 MPIB::SGv::Binomial_builder Class Reference

### 9.4.1 Detailed Description

Binomial tree builder for scatterv/gatherv. The largest subtree on the right. If root $<> 0$, processes' procs will be cyclically shifted to 0.

The documentation for this class was generated from the following file:

- collectives/sgv_tree_builders.hpp

## 9.5 MPIB::SGv::Edge_writer Class Reference

### 9.5.1 Detailed Description

Edge writer

The documentation for this class was generated from the following file:

- collectives/sgv_tree_algorithms.hpp

## 9.6 MPIB::SG::Edge_writer Class Reference

### 9.6.1 Detailed Description

Edge writer

The documentation for this class was generated from the following file:

- collectives/sg_tree_algorithms.hpp

## 9.7 MPIB::SG::Indexer Class Reference

### 9.7.1 Detailed Description

Builds the rank2index for the local buffer. Default copy constructors are used.

The documentation for this class was generated from the following file:

- collectives/sg_tree_algorithms.hpp

## 9.8 MPIB::SGv::Indexer Class Reference

**Public Member Functions**

- Indexer (const int ∗_rscounts, int &_index, map< int, int > &_rank2index, int ∗_counts, int ∗_displs, int &_count)

### 9.8.1    Detailed Description

Builds the rank2index, counts and displs for the local buffer. Default copy constructors are used. Arrays must be preallocated (max = comm size).

### 9.8.2    Constructor & Destructor Documentation

**9.8.2.1    MPIB::SGv::Indexer::Indexer ( const int ∗ _rscounts, int & _index, map< int, int > & _rank2index, int ∗ _counts, int ∗ _displs, int & _count ) [inline]**

counts in the buffer to recv/send

The documentation for this class was generated from the following file:

- collectives/sgv_tree_algorithms.hpp

## 9.9    MPIB_Bcast_container Class Reference

Inherits MPIB_buffer_container.

Collaboration diagram for MPIB_Bcast_container:



### 9.9.1    Detailed Description

Bcast container

The documentation for this class was generated from the following file:

- benchmarks/mpib_coll_containers.hpp

## 9.10 MPIB_bcast_timer Struct Reference

Collaboration diagram for MPIB_bcast_timer:



**Public Attributes**

- MPI_Comm comm
- MPIB_result ∗ results

### 9.10.1 Detailed Description

The state of bcast timer

### 9.10.2 Member Data Documentation

#### 9.10.2.1 MPI_Comm MPIB_bcast_timer::comm

The communicator the offset was found for

#### 9.10.2.2 MPIB_result∗ MPIB_bcast_timer::results

The offset between clocks of the current and root processes

The documentation for this struct was generated from the following file:

- benchmarks/mpib_measure_bcast.c

## 9.11 MPIB_buffer_container Class Reference

Inherits MPIB_coll_container.

Inherited by MPIB_Bcast_container, MPIB_SG_container, and MPIB_SGv_container.

---

Collaboration diagram for MPIB_buffer_container:



### 9.11.1   Detailed Description

Base container for Scatter(v)/Gather(v), Bcast. Containes a buffer, which can be of the same (Bcast) or different (Scatter(v)/Gather(v)) sizes at processes.

The documentation for this class was generated from the following file:

- benchmarks/mpib_coll_containers.hpp

## 9.12   MPIB_coll_container Struct Reference

Inherited by MPIB_buffer_container, MPIB_Comm_dup_free_container, and MPIB_Reduce_container.

**Public Attributes**

- const char ∗ operation
- int(∗ initialize )(void ∗_this, MPI_Comm comm, int root, int M)
- int(∗ execute )(void ∗_this, MPI_Comm comm, int root, int M)
- int(∗ finalize )(void ∗_this, MPI_Comm comm, int root)

### 9.12.1   Detailed Description

Container for a collective communication operation to be measured by MPIB_measure_coll (MPIB_-measure_max, MPIB_measure_root, MPIB_measure_global). How to use (example in C):

- Create a data structure with the first field MPIB_coll_container.

```
typedef struct MPIB_Scatter_container {
    MPIB_coll_container base;
    char* buffer;
    MPIB_Scatter scatter;
} MPIB_Scatter_container;
```

- Implement the functions: initialize, execute, finalize, where _this argument can be typecasted to the data structure.

```
void MPIB_Scatter_initialize(void* _this, MPI_Comm comm, int root, int M) {
    MPIB_Scatter_container* container = (MPIB_Scatter_container*)_this;
    int rank;
    MPI_Comm_rank(comm, &rank);
    int size;
    MPI_Comm_size(comm, &size);
    container->buffer = rank == root ?
        (char*)malloc(sizeof(char) * M * size) :
        (char*)malloc(sizeof(char) * M);
}

void MPIB_Scatter_execute(void* _this, MPI_Comm comm, int root, int M) {
    MPIB_Scatter_container* container = (MPIB_Scatter_container*)_this;
    container->scatter(container->base.buffer, M, MPI_CHAR,
        container->base.buffer, M, MPI_CHAR,
        root, comm);
}

void MPIB_Scatter_finalize(void* _this, MPI_Comm comm, int root) {
    MPIB_Scatter_container* container = (MPIB_Scatter_container*)_this;
    free(container->buffer);
}
```

- Implement the functions that allocate and free the data structure.

```
MPIB_Scatter_container* MPIB_Scatter_container_alloc(MPIB_Scatter scatter) {
    MPIB_Scatter_container* container =
        (MPIB_Scatter_container*)malloc(sizeof(MPIB_Scatter_container));
    container->base.operation = "Scatter";
    container->base.initialize = MPIB_Scatter_initialize;
    container->base.execute = MPIB_Scatter_execute;
    container->base.finalize = MPIB_Scatter_finalize;
    container->scatter = scatter;
    return container;
}

void MPIB_Scatter_container_free(void* _this) {
    free(_this);
}
```

    +

In this library, collective containers are implemented in C++.

**Parameters**

    *comm* MPI communicator over which the communication operation will be performed

    *root* root process

    *M* message size

### 9.12.2 Member Data Documentation

#### 9.12.2.1 const char∗ MPIB_coll_container::operation

                                                       Communication operation

#### 9.12.2.2 int(∗ MPIB_coll_container::initialize)(void ∗_this, MPI_Comm comm, int root, int M)

    Initialization of buffers required for the communication operation (in irregular collectives, M can be different at different processors)

**9.12.2.3 int(∗ MPIB_coll_container::execute)(void ∗_this, MPI_Comm comm, int root, int M)**

Communication operation (in irregular collectives, M can be different at different processors)

**9.12.2.4 int(∗ MPIB_coll_container::finalize)(void ∗_this, MPI_Comm comm, int root)**
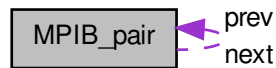
Finalization of buffers required for the communication operation

The documentation for this struct was generated from the following file:

- benchmarks/mpib_coll_benchmarks.h

## 9.13 MPIB_Comm_dup_free_container Class Reference

Inherits MPIB_coll_container.

Collaboration diagram for MPIB_Comm_dup_free_container:



### 9.13.1 Detailed Description

Container for MPI_Comm_dup-MPI_Comm_free

The documentation for this class was generated from the following file:

- benchmarks/mpib_coll_containers.hpp

## 9.14 MPIB_Gather_container Class Reference

Inherits MPIB_SG_container.

Collaboration diagram for MPIB_Gather_container:



### 9.14.1 Detailed Description

Gather container

The documentation for this class was generated from the following file:

- benchmarks/mpib_coll_containers.hpp

## 9.15 MPIB_Gatherv_container Class Reference

Inherits MPIB_SGv_container.

Collaboration diagram for MPIB_Gatherv_container:



### 9.15.1  Detailed Description

Gatherv container

The documentation for this class was generated from the following file:

- benchmarks/mpib_coll_containers.hpp

## 9.16  MPIB_global_timer Struct Reference

**Public Attributes**

- MPI_Comm comm
- double ∗ delta

### 9.16.1  Detailed Description

The state of global timer

### 9.16.2  Member Data Documentation

#### 9.16.2.1  MPI_Comm MPIB_global_timer::comm

The communicator the offset was found for

### 9.16.2.2   double∗ MPIB_global_timer::delta

The offset between clocks of the current and root processes

The documentation for this struct was generated from the following file:

- benchmarks/mpib_measure_global.c

## 9.17   MPIB_msgset Struct Reference

**Public Attributes**

- int min_size
- int max_size
- int stride
- double max_diff
- int min_stride
- int max_num

### 9.17.1   Detailed Description

The message sizes for which the measurements are to be performed. Bounded by min_size and max_size. If stride > 0, message sizes are selected regularly. Otherwise, they are adaptively selected at runtime, based on the max_diff, min_stride and max_num values.

### 9.17.2   Member Data Documentation

#### 9.17.2.1   int MPIB_msgset::min_size

Maximum message size in bytes

#### 9.17.2.2   int MPIB_msgset::max_size

Maximum message size in bytes

#### 9.17.2.3   int MPIB_msgset::stride

Stride in bytes for regular selection of message sizes

#### 9.17.2.4   double MPIB_msgset::max_diff

Maximum relative difference between the result of measurement and the linear model based on the results of two previous measurements that requires further investigation. Must be non-negative, $\leq 1$. Used in adaptive selection of message sizes.

#### 9.17.2.5   int MPIB_msgset::min_stride

Minimum stride between message sizes. Must be positive. Used in adaptive selection of message sizes.

### 9.17.2.6   int MPIB_msgset::max_num

Maximum number of message sizes. Limits the number of different messages sizes the measurement is performed for. Used in adaptive selection of message sizes.
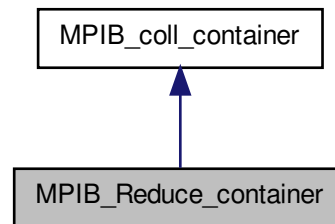
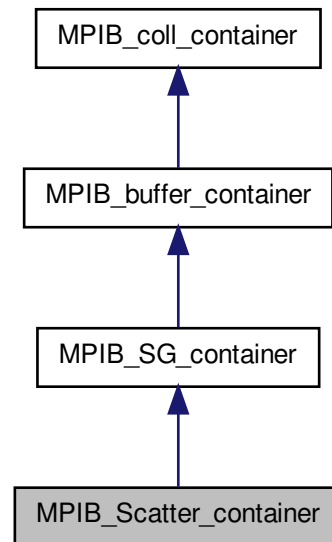The documentation for this struct was generated from the following file:

- benchmarks/mpib_measurement.h

## 9.18   MPIB_p2p_container Struct Reference

Inherited by MPIB_roundtrip_container.

### Public Attributes

- const char ∗ operation
- void(∗ initialize )(void ∗_this, MPI_Comm comm, int M)
- void(∗ execute_measure )(void ∗_this, MPI_Comm comm, int M, int mirror)
- void(∗ execute_mirror )(void ∗_this, MPI_Comm comm, int M, int measure)
- void(∗ finalize )(void ∗_this, MPI_Comm comm)

### 9.18.1   Detailed Description

Container for a point-to-point communication operation to be measured by MPIB_measure_allp2p. How to use (example in C):

- Create a data structure with the first field MPIB_p2p_container.

```
typedef struct MPIB_roundtrip_container {
    MPIB_p2p_container base;
    char* buffer;
} MPIB_roundtrip_container;
```

- Implement the functions: initialize, execute_measure, execute_mirror, finalize, where _this argument can be typecasted to the data structure.

```
void MPIB_Send_Recv_initialize(void* _this, MPI_Comm comm, int M) {
    MPIB_roundtrip_container* container = (MPIB_roundtrip_container*)_this;
    container->buffer = (char*)malloc(sizeof(char) * M);
}

void MPIB_Send_Recv_execute_measure(void* _this, MPI_Comm comm, int M, int mirror
    ) {
    MPIB_roundtrip_container* container = (MPIB_roundtrip_container*)_this;
    MPI_Send(container->buffer, M, MPI_CHAR, mirror, 0, comm);
    MPI_Recv(container->buffer, M, MPI_CHAR, mirror, 0, comm, MPI_STATUS_IGNORE);

}

void MPIB_Send_Recv_execute_mirror(void* _this, MPI_Comm comm, int M, int measure
    ) {
    MPIB_roundtrip_container* container = (MPIB_roundtrip_container*)_this;
    MPI_Recv(container->buffer, M, MPI_CHAR, measure, 0, comm, MPI_STATUS_IGNORE)
    ;
    MPI_Send(container->buffer, M, MPI_CHAR, measure, 0, comm);
}

void MPIB_Send_Recv_finalize(void* _this, MPI_Comm comm) {
```

```
    MPIB_roundtrip_container* container = (MPIB_roundtrip_container*)_this;
    free(container->buffer);
}
```

- Implement the functions that allocate and free the data structure.

```
MPIB_p2p_container* MPIB_roundtrip_container_alloc() {
    MPIB_p2p_container* container =
        (MPIB_p2p_container*)malloc(sizeof(MPIB_roundtrip_container));
    container->base.operation = "MPI_Send-MPI_Recv";
    container->base.free = MPIB_roundtrip_container_free;
    container->initialize = MPIB_Send_Recv_initialize;
    container->execute_measure = MPIB_Send_Recv_execute_measure;
    container->execute_mirror = MPIB_Send_Recv_execute_mirror;
    container->finalize = MPIB_Send_Recv_finalize;
    return container;
}

void MPIB_Send_Recv_container_free(void* _this) {
    free(_this);
}
```

**Parameters**

> *comm*  MPI communicator over which the communication operation will be performed
>
> *M*  message size
>
> *mirror*  mirror process
>
> *measure*  measure process

### 9.18.2   Member Data Documentation

#### 9.18.2.1   const char∗ MPIB_p2p_container::operation

Communication operation

#### 9.18.2.2   void(∗ MPIB_p2p_container::initialize)(void ∗_this, MPI_Comm comm, int M)

Initializion of buffers required for the communication operation.

#### 9.18.2.3   void(∗ MPIB_p2p_container::execute_measure)(void ∗_this, MPI_Comm comm, int M, int mirror)

Part of communication at the measure side

#### 9.18.2.4   void(∗ MPIB_p2p_container::execute_mirror)(void ∗_this, MPI_Comm comm, int M, int measure)

Part of communication at the mirror side

#### 9.18.2.5   void(∗ MPIB_p2p_container::finalize)(void ∗_this, MPI_Comm comm)

Finalization of buffers required for the communication operation

The documentation for this struct was generated from the following file:

- benchmarks/mpib_p2p_benchmarks.h

## 9.19 MPIB_pair Struct Reference

Collaboration diagram for MPIB_pair:



**Public Attributes**

- int values [2]
- struct MPIB_pair * prev
- struct MPIB_pair * next

### 9.19.1 Detailed Description

List of pairs

### 9.19.2 Member Data Documentation

#### 9.19.2.1 int MPIB_pair::values[2]

Values

#### 9.19.2.2 struct MPIB_pair∗ MPIB_pair::prev

Previous pair

#### 9.19.2.3 struct MPIB_pair∗ MPIB_pair::next

Next pair

The documentation for this struct was generated from the following file:

- benchmarks/mpib_utilities.h

## 9.20    MPIB_pairs Struct Reference

Collaboration diagram for MPIB_pairs:



**Public Attributes**

- MPIB_pair ∗ list
- struct MPIB_pairs ∗ prev
- struct MPIB_pairs ∗ next

### 9.20.1    Detailed Description

List of lists of pairs

### 9.20.2    Member Data Documentation

#### 9.20.2.1    MPIB_pair∗ MPIB_pairs::list

Items

#### 9.20.2.2    struct MPIB_pairs∗ MPIB_pairs::prev

Previous list

#### 9.20.2.3    struct MPIB_pairs∗ MPIB_pairs::next

Next list

The documentation for this struct was generated from the following file:

- benchmarks/mpib_utilities.h

---

## 9.21 MPIB_precision Struct Reference

**Public Attributes**

- int min_reps
- int max_reps
- double cl
- double eps

### 9.21.1 Detailed Description

Precision of measurement. Used as an input argument of benchmark functions. To provide reliable results, the communication experiments in each benchmark are repeated either fixed or variable number of times. This data structure allows the user to control the accuracy and efficiency of benchmarking.

- Assigning to min_reps and max_reps the same values results in the fixed number of repetitions of the communication operation, with the cl and eps arguments being ignored (this allows the user to control the efficiency of benchmarking).

- If min_reps $<$ max_reps, the experiments are repeated until a confidence interval, MPIB_result::ci, found with the confidence level, cl $= Pr(|\bar{T} - \mu| < ci)$, satisfies $\frac{ci}{\bar{T}} <$ eps, or the number of repetitions reaches its maximum, max_reps (this allows the user to control the accuracy of benchmarking).

### 9.21.2 Member Data Documentation

#### 9.21.2.1 int MPIB_precision::min_reps

<div align="right">Minimum number of repetitions</div>

#### 9.21.2.2 int MPIB_precision::max_reps

<div align="right">Maximum number of repetitions</div>

#### 9.21.2.3 double MPIB_precision::cl

<div align="center">Confidence level $\in [0, 1]$: $cl = Pr(|\bar{T} - \mu| < ci) = Pr(\frac{|\bar{T} - \mu|}{\bar{T}} < \epsilon)$.</div>

#### 9.21.2.4 double MPIB_precision::eps

<div align="center">Relative error $\in [0, 1]$: $\frac{|\bar{T} - \mu|}{\bar{T}} < \frac{ci}{\bar{T}} < \epsilon = eps$.</div>

The documentation for this struct was generated from the following file:

- benchmarks/mpib_measurement.h

## 9.22    MPIB_Reduce_container Class Reference

Inherits MPIB_coll_container.

Collaboration diagram for MPIB_Reduce_container:



### 9.22.1    Detailed Description

Reduce container

The documentation for this class was generated from the following file:

- benchmarks/mpib_coll_containers.hpp

## 9.23    MPIB_result Struct Reference

**Public Attributes**

- int M
- double T
- double wtick
- int reps
- double ci

### 9.23.1    Detailed Description

Result of measurement and its reliability. Used as an output argument of benchmark functions.

### 9.23.2    Member Data Documentation

#### 9.23.2.1    int MPIB_result::M

Message size

### 9.23.2.2  double MPIB_result::T

Execution time

### 9.23.2.3  double MPIB_result::wtick

Resolution of MPI_Wtime

### 9.23.2.4  int MPIB_result::reps

Number of repetitions the benchmark has actually taken

### 9.23.2.5  double MPIB_result::ci

Confidence interval, $|\bar{T} - \mu| < ci$. The MPIB_ci function estimates confidence interval, using t distribution.

The documentation for this struct was generated from the following file:

- benchmarks/mpib_measurement.h

## 9.24  MPIB_root_timer Struct Reference

**Public Attributes**

- MPI_Comm comm
- double barrier_time

### 9.24.1  Detailed Description

The state of root timer

### 9.24.2  Member Data Documentation

### 9.24.2.1  MPI_Comm MPIB_root_timer::comm

The communicator the barrier time was found for

### 9.24.2.2  double MPIB_root_timer::barrier_time

The barrier time at the current process

The documentation for this struct was generated from the following file:

- benchmarks/mpib_measure_root.c

## 9.25    MPIB_Scatter_container Class Reference

Inherits MPIB_SG_container.

Collaboration diagram for MPIB_Scatter_container:



### 9.25.1    Detailed Description

Scatter container

The documentation for this class was generated from the following file:

- benchmarks/mpib_coll_containers.hpp

## 9.26    MPIB_Scatterv_container Class Reference

Inherits MPIB_SGv_container.

Collaboration diagram for MPIB_Scatterv_container:



### 9.26.1    Detailed Description

Scatterv container

The documentation for this class was generated from the following file:

- benchmarks/mpib_coll_containers.hpp

## 9.27    MPIB_SG_container Class Reference

Inherits MPIB_buffer_container.

Inherited by MPIB_Gather_container, and MPIB_Scatter_container.

Collaboration diagram for MPIB_SG_container:



### 9.27.1 Detailed Description

Base container for Scatter/Gather

The documentation for this class was generated from the following file:

- benchmarks/mpib_coll_containers.hpp

## 9.28 MPIB_SGv_container Class Reference

Inherits MPIB_buffer_container.

Inherited by MPIB_Gatherv_container, and MPIB_Scatterv_container.

Collaboration diagram for MPIB_SGv_container:



### 9.28.1    Detailed Description

Base container for Scatterv/Gatherv

The documentation for this class was generated from the following file:

- benchmarks/mpib_coll_containers.hpp

## 9.29    MPIB_SGv_count_sorter Struct Reference

Collaboration diagram for MPIB_SGv_count_sorter:

### 9.29.1   Detailed Description

Count sorter

The documentation for this struct was generated from the following file:

- collectives/sgv_flat.c

## 9.30   MPIB_SGv_sorter Struct Reference

**Public Attributes**

- void(∗ sort )(void ∗_this, const int size, const int ∗counts, size_t ∗indices)

### 9.30.1   Detailed Description

Sorter for base sorted flat-tree scatterv/gatherv

### 9.30.2   Member Data Documentation

#### 9.30.2.1   void(∗ MPIB_SGv_sorter::sort)(void ∗_this, const int size, const int ∗counts, size_t ∗indices)

**Parameters**

   _this_   self pointer

   _size_   number of processes

   _counts_   message sizes (in bytes)

   _indices_   sorted ranks (result), must be pre-allocated

The documentation for this struct was generated from the following file:

- collectives/sgv_flat.h

## 9.31   MPIB::SGv::Sorted_binomial_builder Class Reference

### 9.31.1   Detailed Description

Sorted binomial tree builder for scatterv/gatherv. The largest subtree on the right. Processors are sorted by counts in asc/dsc order. Algorithm is described in [3].

The documentation for this class was generated from the following file:

- collectives/sgv_tree_builders.hpp

---

## 9.32    MPIB::SGv::Traff_builder Class Reference

### 9.32.1    Detailed Description

Traff tree builder for scatterv/gatherv. Algorithm is described in [3].

The documentation for this class was generated from the following file:

- collectives/sgv_tree_builders.hpp

## 9.33    MPIB::SGv::Vertex_writer Class Reference

### 9.33.1    Detailed Description

Vertex writer

The documentation for this class was generated from the following file:

- collectives/sgv_tree_algorithms.hpp

# 10    File Documentation

## 10.1    benchmarks/mpib.h File Reference

Include dependency graph for mpib.h:

This graph shows which files directly or indirectly include this file:



### 10.1.1  Detailed Description

For internal use

# Index

Vertex
     MPIB::comm_tree, 35
Vertex_iterator
     MPIB::comm_tree, 35

Wrappers for p2p operations, 24
wtick
     MPIB_result, 53