

# *The Message Passing Interface: MPI 3.1 and Plans for MPI 4.0*

**Martin Schulz**

LLNL / CASC

Chair of the MPI Forum

MPI Forum BOF @ SC14



<http://www.mpi-forum.org/>

This work was performed under the auspices of the U.S. Department of Energy by Lawrence Livermore National Laboratory under Contract DE-AC52-07NA27344.



# Overview

- **Current State of MPI**
  - Features in MPI 3.0
  - Implementation status
- **Timeline for MPI 3.1**
- **Initiatives for MPI 4.0**
  - Fault tolerance
  - Support for hybrid programming
  - Stream/channel communication
- **How to contribute to the MPI Forum**

**Let's keep this interactive – Please feel free to ask questions!**

## Where We Are

- **MPI 3.0 ratified in September 2012**
  - Available at <http://www.mpi-forum.org/>
  - 852 pages, 430 functions
  - Adaption in most MPIs progressing fast



# MPI 3.0 has 430 Functions

MPI_ABORT	MPI_ERRHANDLER_GET	MPI_GROUP_DIFFERENCE	MPI_QUERY_THREAD	MPI_TYPE_DELETE_ATTR
MPI_ACCUMULATE	MPI_ERRHANDLER_SET	MPI_GROUP_EXCL	<b>MPI_RACCOMULATE</b>	MPI_TYPE_DUP
MPI_ADD_ERROR_CLASS	MPI_ERROR_CLASS	MPI_GROUP_F2C	MPI_RECV	MPI_TYPE_DUP_FN
MPI_ADD_ERROR_CODE	MPI_ERROR_STRING	MPI_GROUP_FREE	MPI_RECV_INIT	MPI_TYPE_EXTENT
MPI_ADD_ERROR_STRING	MPI_EXSCAN	MPI_GROUP_INCL	MPI_REDUCE	MPI_TYPE_F2C
MPI_ADDRESS	MPI_FETCH_SYNC_REG	MPI_GROUP_INTERSECTION	<b>MPI_REDUCE_LOCAL</b>	MPI_TYPE_FREE
MPI_ALLGATHER	MPI_FETCH_AND_OP	MPI_GROUP_RANGE_EXCL	MPI_REDUCE_SCATTER	MPI_TYPE_FREE_KEYVAL
MPI_ALLGATHERV	MPI_FILE_C2F	MPI_GROUP_RANGE_INCL	<b>MPI_REDUCE_SCATTER_BLOCK</b>	MPI_TYPE_GET_ATTR
MPI_ALLOC_MEM	MPI_FILE_CALL_ERRHANDLER	MPI_GROUP_RANK	MPI_REGISTER_DATAREP	MPI_TYPE_GET_CONTENTS
MPI_ALLOC_OPEN_CPTR	MPI_FILE_CLOSE	MPI_GROUP_SIZE	MPI_REQUEST_C2F	MPI_TYPE_GET_ENVELOPE
MPI_ALLOC_VIEW	MPI_FILE_CREATE_ERRHANDLER	MPI_GROUP_TRANSLATE_RANKS	MPI_REQUEST_C2C	MPI_TYPE_GET_EXTENT
MPI_ALLOC_VIEW	MPI_FILE_DELETE	MPI_GROUP_UNION	MPI_REQUEST_FREE	<b>MPI_TYPE_GET_EXTENT_X</b>
MPI_ALLTOALL	MPI_FILE_F2C	<b>MPI_IALLGATHER</b>	MPI_REQUEST_GET_STATUS	MPI_TYPE_GET_NAME
MPI_ALLTOALLV	MPI_FILE_GET_AMODE	<b>MPI_IALLGATHERV</b>	<b>MPI_RGET</b>	MPI_TYPE_GET_TRUE_EXTENT
MPI_ALLTOALLW	MPI_FILE_GET_ATOMICTY	<b>MPI_IALLREDUCE</b>	<b>MPI_RGET_ACCUMULATE</b>	<b>MPI_TYPE_GET_TRUE_EXTENT_X</b>
MPI_ATTR_DELETE	MPI_FILE_GET_BYTE_OFFSET	<b>MPI_IALLTOALL</b>	MPI_SEND	MPI_TYPE_INDEXED
MPI_ATTR_GET	MPI_FILE_GET_ERRHANDLER	<b>MPI_IALLTOALLV</b>	MPI_SEND	MPI_TYPE_HVECTOR
MPI_ATTR_PUT	MPI_FILE_GET_GROUP	<b>MPI_IALLTOALLW</b>	MPI_SEND_INIT	MPI_TYPE_INDEXED
MPI_BARRIER	MPI_FILE_GET_INFO	<b>MPI_IBARRIES</b>	MPI_SCAN	MPI_TYPE_LB
MPI_BCAST	MPI_FILE_GET_POSITION	<b>MPI_IBCAST</b>	MPI_SCATTER	MPI_TYPE_MATCH_SIZE
MPI_BSEND	MPI_FILE_GET_POSITION_SHARED	<b>MPI_IBSEND</b>	MPI_SCATTERV	MPI_TYPE_NULL_COPY_FN
MPI_BSEND_INIT	MPI_FILE_GET_SIZE	<b>MPI_IXSCAN</b>	MPI_SEND	MPI_TYPE_NULL_DELETE_FN
MPI_BUFFER_ATTACH	MPI_FILE_GET_TYPE_EXTENT	<b>MPI_IGATHER</b>	MPI_SEND_INIT	MPI_TYPE_SET_ATTR
MPI_BUFFER_DETACH	MPI_FILE_GET_VIEW	<b>MPI_IGATHERV</b>	MPI_SENDREC	MPI_TYPE_SET_NAME
MPI_CANCEL	MPI_FILE_IREAD	<b>MPI_IGATHERV</b>	MPI_SENDREC_REPLACE	MPI_TYPE_SIZE
MPI_CART_COORDS	MPI_FILE_IWRITE	<b>MPI_IRECV</b>	MPI_SIZEOF	<b>MPI_TYPE_SIZE_X</b>
MPI_CART_CREATE	MPI_FILE_READ	<b>MPI_INEIGHBOR_ALLGATHER</b>	MPI_SSEND	MPI_TYPE_STRUCT
MPI_CART_GET	MPI_FILE_READ_ALL	<b>MPI_INEIGHBOR_ALLGATHERV</b>	MPI_SSEND_INIT	MPI_TYPE_UB
MPI_CART_MAP	MPI_FILE_READ_BEGIN	<b>MPI_INEIGHBOR_ALLTOALL</b>	MPI_START	MPI_TYPE_VECTOR
MPI_CART_RANK	MPI_FILE_READ_END	<b>MPI_INEIGHBOR_ALLTOALLV</b>	MPI_STARTALL	MPI_UNPACK
MPI_CART_SHIFT	MPI_FILE_READ_AT	<b>MPI_INEIGHBOR_ALLTOALLW</b>	MPI_STATUS_C2F	MPI_UNPACK_EXTERNAL
MPI_CART_SUB	MPI_FILE_READ_ORDERED	MPI_INFO_C2F	<b>MPI_STATUS_C2E08</b>	MPI_UNPUBLISH_NAME
MPI_CARTDIM_GET	MPI_FILE_READ_ORDERED_BEGIN	MPI_INFO_CREATE	<b>MPI_STATUS_F082C</b>	MPI_WAIT
MPI_CLOSE_PORT	MPI_FILE_READ_ORDERED_END	MPI_INFO_DELETE	<b>MPI_STATUS_F082F</b>	MPI_WAITALL
MPI_COMM_ACCEPT	MPI_FILE_READ_SHARED	MPI_INFO_DUP	MPI_STATUS_F2C	MPI_WAITANY
MPI_COMM_C2F	MPI_FILE_READ_SHARED_END	MPI_INFO_F2C	<b>MPI_STATUS_F2E08</b>	MPI_WAIT SOME
MPI_COMM_CALL_ERRHANDLER	MPI_FILE_READ_SHARED_END	MPI_INFO_FREE	MPI_STATUS_SET_CANCELLED	<b>MPI_WIN_ALLOC</b>
MPI_COMM_COMPARE	MPI_FILE_READ_SHARED_END	MPI_INFO_GET	MPI_STATUS_SET_ELEMENTS	<b>MPI_WIN_ALLOCATE</b>
MPI_COMM_CONNECT	MPI_FILE_READ_SHARED_END	MPI_INFO_GET_NKEYS	<b>MPI_STATUS_SET_ELEMENTS_X</b>	<b>MPI_WIN_ALLOCATE_CPTR</b>
MPI_COMM_CREATE	MPI_FILE_READ_SHARED_END	MPI_INFO_GET_NTHKEY	<b>MPI_T_CATEGORY_CHANGED</b>	<b>MPI_WIN_ALLOCATE_SHARED</b>
MPI_COMM_CREATE_ERRHANDLER	MPI_FILE_READ_SHARED_END	MPI_INFO_GET_VALUELEN	<b>MPI_T_CATEGORY_GET_CATEGORIES</b>	<b>MPI_WIN_ALLOCATE_SHARED_CPTR</b>
<b>MPI_COMM_CREATE_GROUP</b>	MPI_FILE_READ_SHARED_END	MPI_INFO_SET	<b>MPI_T_CATEGORY_GET_CVARS</b>	MPI_WIN_ATTACH
MPI_COMM_CREATE_KEYVAL	MPI_FILE_READ_SHARED_END	MPI_INIT	<b>MPI_T_CATEGORY_GET_INFO</b>	MPI_WIN_C2F
MPI_COMM_DELETE_ATTR	MPI_FILE_READ_SHARED_END	MPI_INIT_THREAD	<b>MPI_T_CATEGORY_GET_NUM</b>	MPI_WIN_CALL_ERRHANDLER
MPI_COMM_DISCONNECT	MPI_FILE_READ_SHARED_END	MPI_INITIALIZED	<b>MPI_T_CATEGORY_GET_PVARS</b>	MPI_WIN_COMPLETE
MPI_COMM_DUP	MPI_FILE_READ_SHARED_END	MPI_INTERCOMM_CREATE	<b>MPI_T_CVAR_GET_INFO</b>	MPI_WIN_CREATE
MPI_COMM_DUP_FN	MPI_FILE_READ_SHARED_END	MPI_INTERCOMM_MERGE	<b>MPI_T_CVAR_GET_NUM</b>	<b>MPI_WIN_CREATE_DYNAMIC</b>
<b>MPI_COMM_DUP_WITH_INFO</b>	MPI_FILE_READ_SHARED_END	MPI_IPROBE	<b>MPI_T_CVAR_HANDLE_ALLOC</b>	MPI_WIN_CREATE_ERRHANDLER
MPI_COMM_F2C	MPI_FILE_READ_SHARED_END	MPI_IRECV	<b>MPI_T_CVAR_HANDLE_FREE</b>	MPI_WIN_CREATE_KEYVAL
MPI_COMM_FREE	MPI_FILE_READ_SHARED_END	<b>MPI_IREDUCE</b>	<b>MPI_T_CVAR_READ</b>	MPI_WIN_DELETE_ATTR
MPI_COMM_FREE_KEYVAL	MPI_FILE_READ_SHARED_END	<b>MPI_IREDUCE_SCATTER</b>	<b>MPI_T_CVAR_WRITE</b>	MPI_WIN_DETACH
MPI_COMM_GET_ATTR	MPI_FILE_READ_SHARED_END	<b>MPI_IREDUCE_SCATTER_BLOCK</b>	<b>MPI_T_ENUM_GET_INFO</b>	MPI_WIN_DUP_FN
MPI_COMM_GET_ERRHANDLER	MPI_FILE_READ_SHARED_END	MPI_IRSEND	<b>MPI_T_ENUM_GET_ITEM</b>	MPI_WIN_F2C
<b>MPI_COMM_GET_INFO</b>	MPI_FILE_READ_SHARED_END	MPI_IS_THREAD_MAIN	<b>MPI_T_FINALIZE</b>	MPI_WIN_FENCE
MPI_COMM_GET_NAME	MPI_FILE_READ_SHARED_END	<b>MPI_ISCAN</b>	<b>MPI_T_INIT_THREAD</b>	<b>MPI_WIN_FLUSH</b>
MPI_COMM_GET_PARENT	MPI_FILE_READ_SHARED_END	<b>MPI_ISCATTER</b>	<b>MPI_T_PVAR_GET_INFO</b>	<b>MPI_WIN_FLUSH_ALL</b>
MPI_COMM_GROUP	MPI_FILE_READ_SHARED_END	<b>MPI_ISCATTERV</b>	<b>MPI_T_PVAR_GET_NUM</b>	<b>MPI_WIN_FLUSH_LOCAL</b>
<b>MPI_COMM_IDUE</b>	MPI_FILE_READ_SHARED_END	MPI_ISEND	<b>MPI_T_PVAR_HANDLE_ALLOC</b>	<b>MPI_WIN_FLUSH_LOCAL_ALL</b>
MPI_COMM_JOIN	MPI_FILE_READ_SHARED_END	MPI_ISSEND	<b>MPI_T_PVAR_HANDLE_FREE</b>	MPI_WIN_FREE
MPI_COMM_KEYVAL_CREATE	MPI_FILE_READ_SHARED_END	MPI_KEYVAL_CREATE	<b>MPI_T_PVAR_READ</b>	MPI_WIN_FREE_KEYVAL
MPI_COMM_NULL_COPY_FN	MPI_FILE_READ_SHARED_END	MPI_KEYVAL_FREE	<b>MPI_T_PVAR_READRESET</b>	MPI_WIN_GET_ATTR
MPI_COMM_NULL_DELETE_FN	MPI_FILE_READ_SHARED_END	<b>MPI_LOCK_ALL</b>	<b>MPI_T_PVAR_RESET</b>	MPI_WIN_GET_ERRHANDLER
MPI_COMM_RANK	MPI_FILE_READ_SHARED_END	MPI_LOOKUP_NAME	<b>MPI_T_PVAR_SESSION_CREATE</b>	MPI_WIN_GET_GROUP
MPI_COMM_REMOTE_GROUP	MPI_FILE_READ_SHARED_END	MPI_MESSAGE_C2F	<b>MPI_T_PVAR_SESSION_FREE</b>	<b>MPI_WIN_GET_INFO</b>
MPI_COMM_REMOTE_SIZE	MPI_FILE_READ_SHARED_END	MPI_MESSAGE_F2C	<b>MPI_T_PVAR_START</b>	MPI_WIN_GET_NAME
MPI_COMM_SET_ERRHANDLER	MPI_FILE_READ_SHARED_END	<b>MPI_MPROBE</b>	<b>MPI_T_PVAR_STOP</b>	MPI_WIN_LOCK
<b>MPI_COMM_SET_INFO</b>	MPI_FILE_READ_SHARED_END	<b>MPI_MRECV</b>	<b>MPI_T_PVAR_WRITE</b>	<b>MPI_WIN_LOCK_ALL</b>
MPI_COMM_SET_NAME	MPI_FILE_READ_SHARED_END	<b>MPI_NEIGHBOR_ALLGATHER</b>	MPI_TEST	MPI_WIN_NULL_COPY_FN
MPI_COMM_SIZE	MPI_FILE_READ_SHARED_END	<b>MPI_NEIGHBOR_ALLGATHERV</b>	MPI_TEST_CANCELLED	MPI_WIN_NULL_DELETE_FN
MPI_COMM_WORLD	MPI_FILE_READ_SHARED_END	<b>MPI_NEIGHBOR_ALLTOALL</b>	MPI_TESTANY	MPI_WIN_SET
MPI_COMM_WORLD_MULTIPLE	MPI_FILE_READ_SHARED_END	<b>MPI_NEIGHBOR_ALLTOALLV</b>	MPI_TEST SOME	MPI_WIN_SET_ATTR
MPI_COMM_SPLIT	MPI_FILE_READ_SHARED_END	<b>MPI_NEIGHBOR_ALLTOALLW</b>	MPI_TOPO_TEST	MPI_WIN_SET_ERRHANDLER
<b>MPI_COMM_SPLIT_TYPE</b>	MPI_FILE_READ_SHARED_END	MPI_NULL_COPY_FN	MPI_TYPE_C2F	<b>MPI_WIN_SET_INFO</b>
MPI_COMM_TEST_INTER	MPI_FILE_READ_SHARED_END	MPI_NULL_DELETE_FN	MPI_TYPE_COMMIT	MPI_WIN_SET_NAME
<b>MPI_COMPARE_AND_SWAP</b>	MPI_FILE_READ_SHARED_END	MPI_OP_COMMUTATIVE	MPI_TYPE_CONTIGUOUS	MPI_WIN_SHARED_ALLOCATE
MPI_CONVERSION_FN_NULL	MPI_FILE_READ_SHARED_END	MPI_OP_CREATE	MPI_TYPE_DARRAY	MPI_WIN_SHARED_QUERY
MPI_DIMS_CREATE	MPI_FILE_READ_SHARED_END	MPI_OP_F2C	MPI_TYPE_CREATE_F90_COMPLEX	MPI_WIN_SHARED_QUERY_CPTR
<b>MPI_DIST_GRAPH_CREATE</b>	MPI_FILE_READ_SHARED_END	MPI_OP_FREE	MPI_TYPE_CREATE_F90_INTEGER	MPI_WIN_START
<b>MPI_DIST_GRAPH_CREATE_ADJACENT</b>	MPI_FILE_READ_SHARED_END	MPI_OPEN_PORT	MPI_TYPE_CREATE_INDEXED	<b>MPI_WIN_SYNC</b>
<b>MPI_DIST_GRAPH_CREATE_NEIGHBOR</b>	MPI_FILE_READ_SHARED_END	MPI_PACK_EXTERNAL	<b>MPI_TYPE_CREATE_INDEXED_BLOCK</b>	MPI_WIN_UNLOCK
<b>MPI_DIST_GRAPH_CREATE_NEIGHBORS</b>	MPI_FILE_READ_SHARED_END	MPI_PACK_EXTERNAL_SIZE	MPI_TYPE_CREATE_HVECTOR	<b>MPI_WIN_UNLOCK_ALL</b>
<b>MPI_DIST_GRAPH_CREATE_NEIGHBORS_COUNT</b>	MPI_FILE_READ_SHARED_END	MPI_PACK_SIZE	MPI_TYPE_CREATE_INDEXED_BLOCK	MPI_WIN_WAIT
MPI_DUP_FN	MPI_FILE_READ_SHARED_END	MPI_PROBE	MPI_TYPE_CREATE_KEYVAL	MPI_WTICK
MPI_ERRHANDLER_C2F	MPI_FILE_READ_SHARED_END	MPI_PUBLISH_NAME	MPI_TYPE_CREATE_STRUCT	MPI_WTIME
MPI_ERRHANDLER_CREATE	MPI_FILE_READ_SHARED_END	MPI_PUT	MPI_TYPE_CREATE_SUBARRAY	
MPI_ERRHANDLER_F2C				
MPI_ERRHANDLER_FREE				

The Message Passing Interface: MPI 3.1 and Plans for MPI 4.0

Martin Schulz



## Notable Additions to MPI 3.0

- Non-blocking collectives
- Neighborhood collectives
- RMA enhancements
- Shared memory support
- MPI Tool Information Interface
- Non-collective communicator creation
- Fortran 2008 Bindings
- New Datatypes
- Large data counts
- Matched probe



# MPI-3 Impl. as of Novemer 2014 (thanks to Pavan Balaji)

	MPICH	MVAPICH	Open MPI	Cray MPI	Tianhe MPI	Intel MPI	IBM BG/Q MPI <sup>1</sup>	IBM PE MPICH <sup>2</sup>	IBM Platform	SGI MPI	Fujitsu MPI	MS MPI
NB collectives	✓	✓	✓	✓	✓	✓	✓	Q4 '14	✓	✓	✓	
Neighborhood collectives	✓	✓	✓	✓	✓	✓	✓	Q4 '14	Q3 '15	✓	Q2 '15	
RMA	✓	✓	✓	✓	✓	✓	✓	Q4 '14	Q3 '15	✓	Q2 '15	
Shared memory	✓	✓	✓	✓	✓	✓	✓	Q4 '14	Q3 '15	✓	Q2 '15	✓
Tools Interface	✓	✓	✓	✓	✓	✓	✓ <sup>3</sup>	Q4 '14	Q3 '15	✓	Q2 '15	*
Non-collective comm. create	✓	✓	✓	✓	✓	✓	✓	Q4 '14	Q3 '15	✓	Q2 '15	
Fo8 Bindings	✓	✓	✓	Q4 '14	✓	Q4 '14	✓	Q4 '14	Q3 '15	✓	Q2 '15	
New Datatypes	✓	✓	✓	✓	✓	✓	✓	Q4 '14	Q3 '15	✓	Q2 '15	*
Large Counts	✓	✓	✓	✓	✓	✓	✓	Q4 '14	Q3 '15	✓	Q2 '15	*
Matched Probe	✓	✓	✓	✓	✓	✓	✓	Q4 '14	Q3 '15	✓	✓	*

Release dates are estimates and are subject to change at any time.

Empty cells indicate no *publicly announced* plan to implement/support that feature.

<sup>1</sup> Open source, but unsupported

<sup>2</sup> Beta release

<sup>3</sup> No MPI\_T variables exposed

\* Under development

(\*) Platform-specific restrictions might apply for all supported features

# Where We Are

- **MPI 3.0 ratified in September 2012**

- Available at <http://www.mpi-forum.org/>
- 852 pages, 430 functions
- Adaption in most MPIs progressing fast

- **Working towards MPI 3.1**

- Inclusion for errata (mainly RMA, Fortran, MPI\_T)
- Minor updates and additions (address arithmetic and non-block. I/O)
- Currently planned for March 2015

# Where We Are

- **MPI 3.0 ratified in September 2012**
  - Available at <http://www.mpi-forum.org/>
  - 852 pages, 430 functions
  - Adaption in most MPIs progressing fast
- **Working towards MPI 3.1**
  - Inclusion for errata (mainly RMA, Fortran, MPI\_T)
  - Minor updates and additions (address arithmetic and non-block. I/O)
  - Currently planned for March 2015
- **Concurrently discussions for MPI 4.0**
  - Major additions as gating items
  - Schedule tbd. (depends on features)



# Larger Topics under Discussion

- **Fault Tolerance support in MPI**

- Avoid job failure when one node fails
- Suitable fault detection mechanism
- Ability to reason about state of MPI and continue execution

- **Support for Hybrid Programming models**

- Improved thread support
- Easier integration with MPI+X approaches
- Extensions to be able to treat threads as MPI endpoints/processes

- **New point to point mechanisms like Streams/Channels**

- Dedicated point to point connections
- Continued streams of data without the need for individually matching send/recv pairs

# Larger Topics under Discussion

- **Fault Tolerance support in MPI**
  - Part of the Fault Tolerance WG
  - Presentation by Wesley Bland, ANL
- **Support for Hybrid Programming models**
  - Part of the Hybrid WG
  - Presentation by Pavan Balaji, ANL
- **New point to point mechanisms like Streams/Channels**
  - Part of the Point to Point WG
  - Presentation by Daniel Holmes, EPCC