

*D R A F T*

Document for a Standard Message-Passing Interface

Message Passing Interface Forum

January 24, 2022

This work was supported in part by NSF and ARPA under NSF contract CDA-9115428 and Esprit under project HPC Standards (21111).

This is the result of a LaTeX run of a draft of a single chapter of the MPIF Final Report document.

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48

# Chapter 16

## Process Fault Tolerance

### 16.1 Introduction

In distributed systems with numerous or complex components, a serious risk is that a component fault manifests as a process failure that disrupts the normal execution of a long running application. A process failure is a common outcome for many hardware, network, or software faults that cause a process to crash; it can be more formally defined as a fail-stop failure: the affected MPI process unexpectedly and permanently stops communicating. This chapter introduces MPI features that support the development of applications, libraries, and programming languages that can tolerate MPI process failures. The primary goal is to specify error classes and interfaces that permit users to continue simple MPI communication (e.g., some point-to-point patterns) after failures have impacted the execution and rebuild MPI objects (communicators, files, etc.) as needed to restore the full capability of MPI to carry out application elaborate communication operations (like collective communications), or dynamic process operations (allowing for spawning replacement processes). This specification does not include mechanisms to restore the application data lost due to process failures. The literature is rich with diverse fault tolerance techniques that the users may employ at their discretion, including checkpoint-restart, algorithmic dataset recovery, and continuation ignoring failed MPI processes. All these fault tolerance approaches benefit from, and often require, the definitions and interfaces specified in this chapter in order to resume communicating after a failure.

The expected behavior of MPI in the case of an MPI process failure is defined by the following statements: any MPI operation that involves a failed process must not block indefinitely but either succeed or raise an MPI error (see Section 16.2); an MPI operation that does not involve a failed process will complete normally, unless interrupted by the user through provided functionality. By default, errors indicate only the local impact of the failure on an operation, and make no guarantee that other processes have also been notified of the same failure; asynchronous failure propagation is not guaranteed or required, and users must exercise caution when determining the set of processes where a failure has been detected and raised an error. If an application needs global knowledge of failures, it can use the interfaces defined in Section 16.3 to explicitly propagate the notification of locally detected failures, or set communicators in specific modes that enforce such propagation.

Some usage patterns on reliable machines do not require fault tolerance. An MPI implementation that does not tolerate process failures must never raise a *fault tolerance error* (as listed in Section 16.4). Applications using the interfaces defined in this chapter

1 must be portable across MPI implementations (including those which do not provide fault  
2 tolerance, but in this case the interfaces may exhibit undefined behavior after a process  
3 failure at any MPI process.) Fault tolerant applications may determine if the implementation  
4 supports fault tolerance by querying the predefined attribute `MPI_FT` on `MPI_COMM_WORLD`  
5 (see 9.1.2.)

6  
7 *Advice to users.* The MPI standard does not specify transparent process recovery  
8 upon MPI process failure. In particular, restoring the lost dataset, spawning spare  
9 processes or taking other recovery actions are the responsibility of the user.

10 Many of the operations and semantics described in this chapter are applicable only  
11 when the MPI application has replaced the default error handler  
12 `MPI_ERRORS_ARE_FATAL` on the communicators and windows it uses. (*End of advice*  
13 *to users.*)

## 14 16.2 Failure Notification

15  
16  
17 When an operation raises a fault tolerance error it may not satisfy its specification (like  
18 any other error, see 9.4). Note that the remainder of this chapter defines operations that  
19 maintain full specification semantic after raising a fault tolerance error; such exceptions will  
20 be explicitly stated. A list of fault tolerance errors is provided in Section 16.4.

21  
22 Nonblocking operations do not raise fault tolerance errors during creation or initiation.  
23 The corresponding completion call raises a fault tolerance error when appropriate.

24 An operation involving a failed MPI process must always complete in a finite amount  
25 of time (possibly by raising one of the process failure error classes listed in Section 16.4).

26 An MPI process is considered involved in a communication (for the purpose of this  
27 chapter) if its failure may prevent the successful update of user-visible state (e.g., output  
28 buffers, synchronizations, etc.) by its operations. More formally, an MPI process is involved  
29 in a communication if any of the following is true:

- 30 • The process is in the group over which the operation is collective.
- 31
- 32 • The process is a destination or a specified or matched source in a point-to-point  
33 communication.
- 34
- 35 • The operation is an `MPI_ANY_SOURCE` receive operation and the process belongs to  
36 the source group.
- 37
- 38 • The process is a specified target in a remote memory operation.

39  
40 By default, process failure errors are raised only during communication operations  
41 in which a failed process is involved, but the range of processes whose failure may cause  
42 operations to raise an error is user controllable.

### 43 16.2.1 Error Reporting Range

44  
45 Users can control the range of processes whose failure cause MPI operations to raise errors  
46 on their communicators by setting the following values to the info key "mpi\_error\_range" on  
47 their communicators:  
48

"operation" If an operation on the communicator does not involve a failed MPI process (such as a point-to-point message between two non-failed MPI processes), it must not raise a fault tolerance error. This is the default if the info key is not set.

"group" The failure of any MPI process in the group of the communicator cause the communication context to become revoked. This causes communication operations to raise a fault tolerance error of class MPI\_ERR\_REVOKED, even in operations that involve only non-failed processes.

"global" The failure of any MPI process in the MPI universe cause communication context to become revoked. This causes communication operations to raise a fault tolerance error of class MPI\_ERR\_REVOKED, even in operations that involve only non-failed processes.

*Advice to implementors.* As long as an implementation can complete operations, it may choose to delay raising an error. Another valid implementation might choose to raise an error as quickly as possible. (*End of advice to implementors.*)

### 16.2.2 Fault Tolerance Errors in Point-to-Point Communication

An MPI implementation raises errors of the following classes in order to notify users that a point-to-point communication operation could not complete successfully because of the failure of at least one involved MPI process:

- MPI\_ERR\_PROC\_FAILED\_PENDING indicates, for a nonblocking communication, that the communication is a receive operation from MPI\_ANY\_SOURCE and no send operation has matched, yet a potential sending MPI process has failed. Neither the operation nor the request identifying the operation is completed.
- In all other cases, the operation raises an error of class MPI\_ERR\_PROC\_FAILED to indicate that the failure prevents the operation from following its failure-free specification. If there is a request identifying a point-to-point communication, it is completed. Communication involving the failed MPI process, initiated on this communicator after the error raised, must also raise an error of class MPI\_ERR\_PROC\_FAILED.

### 16.2.3 Fault Tolerance Errors in Collective Communication

When a collective operation cannot be completed because of the failure of an involved MPI process, the collective operation raises an error of class MPI\_ERR\_PROC\_FAILED.

*Advice to users.*

Depending on how the collective operation is implemented and when an MPI process failure occurs, some participating MPI processes may raise an error while other MPI processes return successfully from the same collective operation. For example, in MPI\_BCAST, the root process may succeed before a failed MPI process disrupts the operation, resulting in some other processes raising an error.

(*End of advice to users.*)

*Advice to users.*

Note that communicator creation functions (e.g., MPI\_COMM\_DUP or MPI\_COMM\_SPLIT) are collective operations. As such, if a failure happened during

1 the call, an error might be raised at some MPI processes while others succeed and  
 2 obtain a new communicator handle. Although it is valid to communicate between  
 3 MPI processes that succeeded in creating the new communicator handle, the user is  
 4 responsible for ensuring a consistent view of the communicator creation, if needed.  
 5 A conservative solution is to check the global outcome of the communicator creation  
 6 function with `MPI_COMM_AGREE` (defined in Section 16.3.1), as illustrated in Ex-  
 7 ample 16.1. (*End of advice to users.*)  
 8

9 After an MPI process failure, `MPI_COMM_FREE` (as with all other collective opera-  
 10 tions) may not complete successfully at all processes. For any MPI process that receives  
 11 the return code `MPI_SUCCESS`, the behavior is defined in Section 7.4.3. If an MPI process  
 12 raises a process failure error (classes `MPI_ERR_PROC_FAILED` or `MPI_ERR_REVOKED`), the  
 13 communicator handle `comm` is set to `MPI_COMM_NULL`; however, the implementation makes  
 14 no guarantee about the success or failure of the `MPI_COMM_FREE` operation, locally or  
 15 remotely.  
 16

17 *Advice to users.* Users are encouraged to call `MPI_COMM_FREE` on communicators  
 18 they do not wish to use anymore, even when they contain failed MPI processes. Al-  
 19 though the operation may raise a fault tolerance error and not synchronize properly,  
 20 this gives a high quality implementation an opportunity to release local resources and  
 21 memory consumed by the object. (*End of advice to users.*)  
 22

## 23 Error Uniformity

24 As noted above, by default, collective communication do not enforce uniformity in error  
 25 raising accross processes. Despite the performance advantages that non-uniformity offer,  
 26 a common usage pattern in applications is to transform non-uniform error raising into a  
 27 uniform behavior accross all processes of the group (as illustrated in Example 16.1).  
 28

29 Users can set collective operations on a communicator to enforce uniform error raising  
 30 by setting the following values in the info key "mpi\_error\_uniform" on the communicator:

31 "local" Process fault tolerance errors are raised to indicate that an MPI process failure  
 32 prevents from guaranteeing the specified behavior at the local process for the collective  
 33 communication operation (e.g., the output buffer contains invalid data). Other pro-  
 34 cesses may have locally satisfied their specification (e.g., the output buffer is valid at  
 35 that process) and may have returned `MPI_SUCCESS`.  
 36

37 "coll" Process fault tolerance errors are raised to indicate that an MPI process failure  
 38 prevents from guaranteeing the specified behavior at any process for the collective  
 39 communication operation. Non-synchronizing collective communication become syn-  
 40 chronizing.  
 41

42 "create" Process fault tolerance errors are raised to indicate that an MPI process failure  
 43 prevents from guaranteeing the specified behavior at any process for the MPI com-  
 44 munication context creation/destruction collective operation (e.g., `MPI_COMM_DUP`  
 45 could not create a new communicator at any process in the group of the communi-  
 46 cator). Non-synchronizing context creation/destruction collective operations become  
 47 synchronizing. Collective communication that do not create or free a communication  
 48 context are not impacted.

### 16.2.4 Dynamic Process Management

*Rationale.* As with communicator creation functions, if a failure happens during a dynamic process management operation, an error might be raised at some MPI processes while others succeed and obtain a new valid communicator. For most communicator creation functions, users can validate the success of the operation by communicating on a pre-existing communicator spanning over the same group of processes (in the worst case, from `MPI_COMM_WORLD`). This is however not always possible for dynamic process management operations, since these operations can create a new intercommunicator between previously disconnected MPI processes. The following additional failure case semantics allow for users to validate, on the created intercommunicator itself, the success of the dynamic process management operation. (*End of rationale.*)

If the MPI implementation raises a fault tolerance error at the root process in `MPI_COMM_ACCEPT` or `MPI_COMM_CONNECT`, the corresponding operation must also raise a fault tolerance error at its root process.

*Advice to users.* The root process of an operation can succeed when a fault tolerance error is raised at some other non-root process. (*End of advice to users.*)

When using the intercommunicator returned from `MPI_COMM_SPAWN`, `MPI_COMM_SPAWN_MULTIPLE`, or `MPI_COMM_GET_PARENT`, a communication for which the root process of the spawn operation is the source or the destination must not deadlock. When the root process raises a fault tolerance error from a spawn operation, no MPI processes are spawned.

*Advice to implementors.* An implementation is allowed to abort a spawned MPI process during `MPI_INIT` when it cannot setup an intercommunicator with the root process of the spawn operation because of a process failure.

An implementation may report it spawned all the requested MPI processes even when a process created from `MPI_COMM_SPAWN` or `MPI_COMM_SPAWN_MULTIPLE` failed, and instead delay raising a fault tolerance error to a later communication involving this process. (*End of advice to implementors.*)

*Advice to users.* To determine how many new MPI processes have effectively been spawned, the normal semantics for hard and soft spawn applies: if the requested number of processes is unavailable for a hard spawn, an error of class `MPI_ERR_SPAWN` is raised (possibly leaving MPI in an undefined state), and an appropriate error code is set in the `array_of_errcodes` parameter. Note however that an implementation may report that it has spawned the requested number of MPI processes even when some MPI processes have failed before exiting `MPI_INIT`. This condition can be detected by communicating over the created intercommunicator with these processes. (*End of advice to users.*)

*Advice to implementors.* `MPI_COMM_JOIN` does not require any supplementary semantics. When the remote MPI process on the fd socket has failed, the operation succeeds and sets `intercomm` to `MPI_COMM_NULL`. (*End of advice to implementors.*)

1 After an MPI process failure, `MPI_COMM_DISCONNECT` (as with all other collective  
2 operations) may not complete successfully at all MPI processes. For any process that receives  
3 the return code `MPI_SUCCESS`, the behavior is defined in 11.10.4. If an MPI process raises a  
4 fault tolerance error (classes `MPI_ERR_PROC_FAILED` or `MPI_ERR_REVOKED`), the commu-  
5 nicator handle `comm` is set to `MPI_COMM_NULL`; however, the implementation makes no  
6 guarantee about the success or failure of the `MPI_COMM_DISCONNECT` operation, locally  
7 or remotely.

8  
9 *Advice to users.* Users are encouraged to call `MPI_COMM_DISCONNECT` on com-  
10 municators they do not wish to use anymore, even when they contain failed MPI  
11 processes. Although the operation may raise a fault tolerance error and not synchro-  
12 nize properly, this gives a high quality implementation an opportunity to release local  
13 resources and memory consumed by the object. (*End of advice to users.*)

## 14 16.2.5 One-Sided Communication

15  
16 When an operation on a window raises a fault tolerance error, the state of all data held  
17 in memory exposed by that window becomes undefined at all MPI processes for which  
18 a one-sided communication operation could have modified local data in that window (a  
19 target in a remote write, or accumulate operation, or an origin in a remote read operation),  
20 and the operation completion has not been semantically guaranteed (as an example by a  
21 successful synchronization between the origin and the target, after the origin had issued an  
22 `MPI_WIN_FLUSH`).

23  
24 *Advice to users.* Assessing if a particular portion of the exposed memory remains  
25 correct is the responsibility of the user. Note that in passive target mode, when an  
26 error is raised at the origin, the target memory data may become undefined before a  
27 synchronization raises an error at the target.

28 The exposed memory data becomes undefined for all uses, not only the window in  
29 which the error was raised. Any overlapping windows or uses involving shared memory  
30 also read undefined data (even if they do not involve MPI calls). (*End of advice to*  
31 *users.*)

32  
33 *Advice to implementors.* A high quality implementation should limit the scope of the  
34 exposed memory that becomes undefined (for example, only the memory addresses  
35 and ranges that have been targeted by a remote write, or accumulate, or have been  
36 an origin in a remote read). In that case, we encourage implementations to document  
37 the provided behavior, and to expose the availability of this feature at runtime, as  
38 an example by caching an implementation specific attribute on the window. (*End of*  
39 *advice to implementors.*)

40 Non-synchronizing one-sided communication operations (as an example `MPI_GET`,  
41 `MPI_PUT`) whose outputs are undefined, due to an MPI process failure, are not required to  
42 raise a fault tolerance error. However, if a communication cannot complete correctly due  
43 to process failures, the synchronization operation must raise a fault tolerance error at least  
44 at the origin.

45  
46 *Advice to implementors.* Non-synchronizing operations (`MPI_WIN_FLUSH_LOCAL`,  
47 `MPI_WIN_FLUSH_LOCAL_ALL`) are not required to raise a fault tolerance error. (*End*  
48 *of advice to implementors.*)



*Advice to users.* As with collective operations over MPI communicators, active target one-sided synchronization operations may raise a fault tolerance error at some MPI process while the corresponding operation returned MPI\_SUCCESS at some other MPI process. (*End of advice to users.*)

Passive target synchronization operations may raise a process failure error when any MPI process in the window has failed (even when the target specified in the argument of the passive target synchronization has not failed).

*Rationale.* An implementation of passive target synchronization may need to communicate with non-target MPI processes in the window, as an example, a previous owner of an access epoch on the target window. (*End of rationale.*)

After an MPI process failure, MPI\_WIN\_FREE (as with all other collective operations) may not complete successfully at all MPI processes. For any process that receives the return code MPI\_SUCCESS, the behavior is defined in Section 12.2.5. If a process raises a process failure error (classes MPI\_ERR\_PROC\_FAILED or MPI\_ERR\_REVOKED), the window handle win is set to MPI\_WIN\_NULL; however, the implementation makes no guarantee about the success or failure of the MPI\_WIN\_FREE operation, locally or remotely.

*Advice to users.* Users are encouraged to call MPI\_WIN\_FREE on windows they do not wish to use anymore, even when they contain failed MPI processes. Although the operation may raise a fault tolerance error and not synchronize properly, this gives a high quality implementation an opportunity to release local resources and memory consumed by the object. Before calling MPI\_WIN\_FREE, it may be required to call MPI\_WIN\_REVOKE to close an epoch that couldn't be completed as a consequence of a process failure (see Section 16.3.2). (*End of advice to users.*)

### 16.2.6 I/O

This section defines the behavior of I/O operations when MPI process failures prevent their successful completion. I/O backend failure error classes and their consequences are defined in Section 14.7.

If an MPI process failure prevents a file operation from completing, an MPI error of class MPI\_ERR\_PROC\_FAILED is raised. Once an MPI implementation has raised an error of class MPI\_ERR\_PROC\_FAILED, the state of the file pointers involved in the operation that raised the error is *undefined*.

*Advice to users.* Since collective I/O operations may not synchronize with other MPI processes, process failures may not be reported during a collective I/O operation. Users are encouraged to use MPI\_COMM\_AGREE on a communicator containing the same group as the file handle when they need to deduce the completion status of collective operations on file handles and maintain a consistent view of file pointers. The file pointer can be reset by using MPI\_FILE\_SEEK with the MPI\_SEEK\_SET update mode. (*End of advice to users.*)

After an MPI process failure, MPI\_FILE\_CLOSE (as with all other collective operations) may not complete successfully at all MPI processes. For any MPI process that receives the return code MPI\_SUCCESS, the behavior is defined in Section 14.2.2. If an MPI process

1 raises a process failure error (classes `MPI_ERR_PROC_FAILED` or `MPI_ERR_REVOKED`), the  
 2 file handle `fh` is set to `MPI_FILE_NULL`; however, the implementation makes no guarantee  
 3 about the success or failure of the `MPI_FILE_CLOSE` operation, locally or remotely.

4  
 5 *Advice to users.* Users are encouraged to call `MPI_FILE_CLOSE` on files they do  
 6 not wish to use anymore, even when they contain failed MPI processes. Although the  
 7 operation may raise a fault tolerance error and not synchronize properly, this gives  
 8 a high quality implementation an opportunity to release local resources and memory  
 9 consumed by the object. (*End of advice to users.*)

## 11 16.3 Failure Mitigation Functions

### 13 16.3.1 Communicator Functions

14  
 15 Process failure notification is not global in MPI. MPI processes that do not call operations  
 16 involving a failed MPI process are possibly never notified of its failure (see Section 16.2). If  
 17 a notification must be propagated, MPI provides a function to revoke a communicator at  
 18 all members.

19  
 20 `MPI_COMM_REVOKE(comm)`

21  
 22 IN comm communicator (handle)

#### 23 **C binding**

24  
 25 `int MPI_Comm_revoke(MPI_Comm comm)`

#### 26 **Fortran 2008 binding**

27  
 28 `MPI_Comm_revoke(comm, ierror)`  
 29 TYPE(MPI\_Comm), INTENT(IN) :: comm  
 30 INTEGER, OPTIONAL, INTENT(OUT) :: ierror

#### 31 **Fortran binding**

32  
 33 `MPI_COMM_REVOKE(COMM, IERROR)`  
 34 INTEGER COMM, IERROR

35  
 36 This function notifies all MPI processes in the groups (local and remote) associated  
 37 with the communicator `comm` that this communicator is revoked. The revocation of a  
 38 communicator by any MPI process completes non-local MPI operations on `comm` at all MPI  
 39 processes by raising an error of class `MPI_ERR_REVOKED` (with the exception of  
 40 `MPI_COMM_SHRINK`, `MPI_COMM_AGREE`, and `MPI_COMM_IAGREE`). This function is  
 41 not collective and therefore does not have a matching call on remote MPI processes. All non-  
 42 failed MPI processes belonging to `comm` will be notified of the revocation despite failures.

43  
 44 A communicator is revoked at a given MPI process either when  
 45 `MPI_COMM_REVOKE` is locally called on it, or when any MPI operation on `comm` raises an  
 46 error of class `MPI_ERR_REVOKED` at that process. Once a communicator has been revoked  
 47 at an MPI process, all subsequent non-local operations on that communicator (with the  
 48 same exceptions as above), are considered local and must complete by raising an error of  
 class `MPI_ERR_REVOKED` at that MPI process.

MPI\_COMM\_IS\_REVOKED(comm, flag) 1

IN	comm	communicator (handle) <span style="float: right;">2</span>
OUT	flag	true if the communicator is revoked (logical) <span style="float: right;">3</span>

### C binding 5

```
int MPI_Comm_is_revoked(MPI_Comm comm, int *flag) 6
```

### Fortran 2008 binding 7

```
MPI_Comm_is_revoked(comm, flag, ierror) 8
  TYPE(MPI_Comm), INTENT(IN) :: comm 9
  LOGICAL, INTENT(OUT) :: flag 10
  INTEGER, OPTIONAL, INTENT(OUT) :: ierror 11
```

### Fortran binding 12

```
MPI_COMM_IS_REVOKED(COMM, FLAG, IERROR) 13
  INTEGER COMM, IERROR 14
  LOGICAL FLAG 15
```

Returns flag = true if the communicator associated with the handle comm is revoked at the calling process. It returns flag = false otherwise. The operation is local. 16

*Advice to users.* In a multithreaded application, a thread calling MPI\_COMM\_IS\_REVOKED may return flag = true before the operation that raises the first exception of class MPI\_ERR\_REVOKED has completed in a concurrent thread. (End of advice to users.) 17

MPI\_COMM\_SHRINK(comm, newcomm) 18

IN	comm	communicator (handle) <span style="float: right;">19</span>
OUT	newcomm	communicator (handle) <span style="float: right;">20</span>

### C binding 21

```
int MPI_Comm_shrink(MPI_Comm comm, MPI_Comm *newcomm) 22
```

### Fortran 2008 binding 23

```
MPI_Comm_shrink(comm, newcomm, ierror) 24
  TYPE(MPI_Comm), INTENT(IN) :: comm 25
  TYPE(MPI_Comm), INTENT(OUT) :: newcomm 26
  INTEGER, OPTIONAL, INTENT(OUT) :: ierror 27
```

### Fortran binding 28

```
MPI_COMM_SHRINK(COMM, NEWCOMM, IERROR) 29
  INTEGER COMM, NEWCOMM, IERROR 30
```

This collective operation creates a new intra- or intercommunicator newcomm from the intra- or intercommunicator comm, respectively, by excluding the group of failed MPI processes as agreed upon during the operation. The groups of newcomm must include every MPI process that returns from MPI\_COMM\_SHRINK, and it must exclude 31

every MPI process whose failure caused an operation on `comm` to raise an MPI error of class `MPI_ERR_PROC_FAILED` or `MPI_ERR_PROC_FAILED_PENDING` at a member of the groups of `newcomm`, before that member initiated `MPI_COMM_SHRINK`. This call is semantically equivalent to an `MPI_COMM_SPLIT` operation that would succeed despite failures, where members of the groups of `newcomm` participate with the same color and a key equal to their rank in `comm`.

This function never raises an error of class `MPI_ERR_PROC_FAILED` or `MPI_ERR_REVOKED`. The defined semantics of `MPI_COMM_SHRINK` are maintained when `comm` is revoked, or when the group of `comm` contains failed MPI processes.

*Advice to users.* `MPI_COMM_SHRINK` is a collective operation, even when `comm` is revoked.

The group of `newcomm` may still contain failed MPI processes, whose failure will be detected in subsequent MPI operations. (*End of advice to users.*)

`MPI_COMM_ISHRINK(comm, newcomm, request)`

IN	<code>comm</code>	communicator (handle)
OUT	<code>newcomm</code>	communicator (handle)
OUT	<code>request</code>	communication request (handle)

### C binding

```
int MPI_Comm_ishrink(MPI_Comm comm, MPI_Comm *newcomm,
                    MPI_Request *request)
```

### Fortran 2008 binding

```
MPI_Comm_ishrink(comm, newcomm, request, ierror)
  TYPE(MPI_Comm), INTENT(IN) :: comm
  TYPE(MPI_Comm), INTENT(OUT), ASYNCHRONOUS :: newcomm
  TYPE(MPI_Request), INTENT(OUT) :: request
  INTEGER, OPTIONAL, INTENT(OUT) :: ierror
```

### Fortran binding

```
MPI_COMM_ISHRINK(COMM, NEWCOMM, REQUEST, IERROR)
  INTEGER COMM, NEWCOMM, REQUEST, IERROR
```

`MPI_COMM_ISHRINK` is a nonblocking variant of `MPI_COMM_SHRINK`. With the exception of its nonblocking behavior, the semantics of `MPI_COMM_ISHRINK` are as if `MPI_COMM_SHRINK` was executed at the time `MPI_COMM_ISHRINK` is called. All restrictions and assumptions for nonblocking collective operations (see Section 6.12) apply to `MPI_COMM_ISHRINK` and the returned request.

Note that, as with `MPI_COMM_IDUP` (see Section 7.4.2), it is erroneous to use `newcomm` before `request` has completed.

MPI_COMM_GET_FAILED(comm, failedgrp)			1
IN	comm	communicator (handle)	2
			3
OUT	failedgrp	group of failed processes (handle)	4

**C binding**

```
int MPI_Comm_get_failed(MPI_Comm comm, MPI_Group *failedgrp)
```

**Fortran 2008 binding**

```
MPI_Comm_get_failed(comm, failedgrp, ierror)
  TYPE(MPI_Comm), INTENT(IN) :: comm
  TYPE(MPI_Group), INTENT(OUT) :: failedgrp
  INTEGER, OPTIONAL, INTENT(OUT) :: ierror
```

**Fortran binding**

```
MPI_COMM_GET_FAILED(COMM, FAILEDGRP, IERROR)
  INTEGER COMM, FAILEDGRP, IERROR
```

This local operation returns the group `failedgrp` of processes from the communicator `comm` that are locally known to have failed. The `failedgrp` can be empty, that is, equal to `MPI_GROUP_EMPTY`.

For any two groups obtained from calls to that routine at the same MPI process with the same `comm`, the smallest group is a prefix of the largest group, that is, the same processes have the same ranks in the two groups up to the size of the smallest group.

*Advice to users.* MPI makes no assumption about asynchronous progress of the failure detection. A valid MPI implementation may choose to update the group of locally known failed MPI processes only when it enters a function that must raise a fault tolerance error.

It is possible that only the calling MPI process has detected the reported failure. If global knowledge is necessary, MPI processes detecting failures should use the call `MPI_COMM_REVOKE`. (*End of advice to users.*)

MPI_COMM_ACK_FAILED(comm, nack, nacked)			34
IN	comm	communicator (handle)	35
			36
IN	nack	Maximum number of process failures to acknowledge (integer)	37
			38
OUT	nacked	Number of process failures acknowledged (integer)	39

**C binding**

```
int MPI_Comm_ack_failed(MPI_Comm comm, int nack, int *nacked)
```

**Fortran 2008 binding**

```
MPI_Comm_ack_failed(comm, nack, nacked, ierror)
  TYPE(MPI_Comm), INTENT(IN) :: comm
  INTEGER, INTENT(IN) :: nack
  INTEGER, INTENT(OUT) :: nacked
```

1       INTEGER, OPTIONAL, INTENT(OUT) :: ierror

2       **Fortran binding**

3       MPI\_COMM\_ACK\_FAILED(COMM, NACK, NACKED, IERROR)

4       INTEGER COMM, NACK, NACKED, IERROR

5  
6       This local operation gives the users a way to acknowledge locally notified failures on  
7 *comm*. The operation acknowledges the first *nack* process failures on *comm*, that is, it  
8 acknowledges the failure of members with a rank lower than *nack* in the group that would  
9 be produced by a concurrent call to MPI\_COMM\_GET\_FAILED on the same *comm*.

10       The operation also sets the value of *nacked* to the current number of acknowledged  
11 process failures in *comm*, that is, a process failure has been acknowledged on *comm* if and  
12 only if the rank of the process is lower than *nacked* in the group that would be produced  
13 by a subsequent call to MPI\_COMM\_GET\_FAILED on the same *comm*.

14       *nacked* can be larger than *nack* when process failures have been acknowledged in a prior  
15 call to MPI\_COMM\_ACK\_FAILED.

16       After an MPI process failure is acknowledged on *comm*, unmatched MPI\_ANY\_SOURCE  
17 receive operations on the same *comm* that would have raised an error of class  
18 MPI\_ERR\_PROC\_FAILED\_PENDING (see Section ??) proceed without further raising errors  
19 due to this acknowledged failure. Also, MPI\_COMM\_AGREE on the same *comm* will not  
20 raise an error of class MPI\_ERR\_PROC\_FAILED due to this acknowledged failure (according  
21 to the specification found later in this section).

22  
23       *Advice to users.* One may query, without side effect, for the number of currently  
24 acknowledged process failures in *comm* by supplying 0 in *nack*. Conversely, one may  
25 unconditionally acknowledge all currently known process failures in  
26 *comm* by supplying the size of the group of *comm* in *nack*. Note that the number of  
27 acknowledged processes, as returned in *nacked*, can be smaller or larger than the value  
28 supplied in *nack*; It is however never larger than the size of the group returned by a  
29 subsequent call to MPI\_COMM\_GET\_FAILED.

30       Calling MPI\_COMM\_ACK\_FAILED on a communicator with failed MPI processes has  
31 no effect on collective operations (except for MPI\_COMM\_AGREE). If a collective  
32 operation would raise an error due to the communicator containing a failed process  
33 (as defined in Section ??), it can continue to raise an error even after the failure  
34 has been acknowledged. In order to use collective operations between MPI processes  
35 of a communicator that contains failed MPI processes, users should create a new  
36 communicator by calling MPI\_COMM\_SHRINK. (*End of advice to users.*)

37  
38  
39       MPI\_COMM\_AGREE(comm, flag)

40       IN        comm                               communicator (handle)

41       INOUT   flag                               bitwise ‘AND’ of contributed values (integer)

42  
43  
44       **C binding**

45       int MPI\_Comm\_agree(MPI\_Comm comm, int \*flag)

46       **Fortran 2008 binding**

47       MPI\_Comm\_agree(comm, flag, ierror)

```

TYPE(MPI_Comm), INTENT(IN) :: comm
INTEGER, INTENT(INOUT) :: flag
INTEGER, OPTIONAL, INTENT(OUT) :: ierror

```

**Fortran binding**

```

MPI_COMM_AGREE(COMM, FLAG, IERROR)
  INTEGER COMM, FLAG, IERROR

```

The purpose of this collective communication is to agree on the integer value `flag` and on the group of failed processes in `comm`.

On completion, all non-failed MPI processes have agreed to set the output integer value of `flag` to the result of a bitwise ‘AND’ operation over the contributed input values of `flag`. If `comm` is an intercommunicator, the value of `flag` is a bitwise ‘AND’ operation over the values contributed by the remote group.

When an MPI process fails before contributing to the operation, the `flag` is computed ignoring its contribution, and `MPI_COMM_AGREE` raises an error of class `MPI_ERR_PROC_FAILED`. However, if all MPI processes have acknowledged this failure prior to the call to `MPI_COMM_AGREE`, using `MPI_COMM_ACK_FAILED`, the error related to this failure is not raised. When an error of class `MPI_ERR_PROC_FAILED` is raised, it is consistently raised at all MPI processes, in both the local and remote groups (if applicable).

After `MPI_COMM_AGREE` raised an error of class `MPI_ERR_PROC_FAILED`, the group produced by a subsequent call to `MPI_COMM_GET_FAILED` on `comm` contains every MPI process that didn’t contribute to the computation of `flag`.

*Advice to users.* Using a combination of `MPI_COMM_ACK_FAILED` and `MPI_COMM_AGREE` as illustrated in Example 16.3, users can propagate and synchronize the knowledge of failures across all MPI processes in `comm`. When `MPI_SUCCESS` is returned locally from `MPI_COMM_AGREE`, the operation has not raised an error of class `MPI_ERR_PROC_FAILED` at any MPI process and thereby returned `MPI_SUCCESS` at all other MPI processes. (*End of advice to users.*)

This function never raises an error of class `MPI_ERR_REVOKED`. The defined semantics of `MPI_COMM_AGREE` are maintained when `comm` is revoked, or when the group of `comm` contains failed MPI processes.

*Advice to users.* `MPI_COMM_AGREE` is a collective operation, even when `comm` is revoked. (*End of advice to users.*)

```

MPI_COMM_IAGREE(comm, flag, request)

```

IN	<code>comm</code>	communicator (handle)
INOUT	<code>flag</code>	bitwise ‘AND’ of contributed values (integer)
OUT	<code>request</code>	communication request (handle)

**C binding**

```

int MPI_Comm_iagree(MPI_Comm comm, int *flag, MPI_Request *request)

```

**Fortran 2008 binding**

```

MPI_Comm_iagree(comm, flag, request, ierror)

```

```

1     TYPE(MPI_Comm), INTENT(IN) :: comm
2     INTEGER, INTENT(INOUT), ASYNCHRONOUS :: flag
3     TYPE(MPI_Request), INTENT(OUT) :: request
4     INTEGER, OPTIONAL, INTENT(OUT) :: ierror

```

### Fortran binding

```

6 MPI_COMM_IAGREE(COMM, FLAG, REQUEST, IERROR)
7     INTEGER COMM, FLAG, REQUEST, IERROR

```

9 This function has the same semantics as MPI\_COMM\_AGREE except that it is non-blocking.

## 16.3.2 One-Sided Functions

### MPI\_WIN\_REVOKE(win)

```

17     IN          win                window object (handle)

```

### C binding

```

20 int MPI_Win_revoke(MPI_Win win)

```

### Fortran 2008 binding

```

23 MPI_Win_revoke(win, ierror)
24     TYPE(MPI_Win), INTENT(IN) :: win
25     INTEGER, OPTIONAL, INTENT(OUT) :: ierror

```

### Fortran binding

```

27 MPI_WIN_REVOKE(WIN, IERROR)
28     INTEGER WIN, IERROR

```

30 This function notifies all MPI processes in the group associated with the window `win` that this window is revoked. The revocation of a window by any MPI process completes RMA operations on `win` at all MPI processes and RMA synchronizations on `win` raise an error of class MPI\_ERR\_REVOKED. This function is not collective and therefore does not have a matching call on remote MPI processes. All non-failed MPI processes belonging to `win` will be notified of the revocation despite failures.

36 A window is revoked at a given MPI process either when MPI\_WIN\_REVOKE is locally called on it, or when any MPI operation on `win` raises an error of class MPI\_ERR\_REVOKED at that process. Once a window has been revoked at an MPI process, all subsequent RMA operations on that window are considered local and RMA synchronizations must complete by raising an error of class MPI\_ERR\_REVOKED at that process. In addition, the current epoch is closed and RMA operations originating from this MPI process are interrupted and completed with undefined outputs.



MPI\_WIN\_IS\_REVOKED(win, flag) 1

	IN	win	window object (handle)	2
				3
	OUT	flag	true if the window is revoked (logical)	4

#### C binding 5

```
int MPI_Win_is_revoked(MPI_Win win, int *flag) 6
```

#### Fortran 2008 binding 7

```
MPI_Win_is_revoked(win, flag, ierror) 8
  TYPE(MPI_Win), INTENT(IN) :: win 9
  LOGICAL, INTENT(OUT) :: flag 10
  INTEGER, OPTIONAL, INTENT(OUT) :: ierror 11
```

#### Fortran binding 12

```
MPI_WIN_IS_REVOKED(WIN, FLAG, IERROR) 13
  INTEGER WIN, IERROR 14
  LOGICAL FLAG 15
```

Returns flag = true if the window associated with the handle win is revoked at the calling process. It returns flag = false otherwise. The operation is local. 16

*Advice to users.* In a multithreaded application, a thread calling MPI\_WIN\_IS\_REVOKED may return flag = true before the operation that raises the first exception of class MPI\_ERR\_REVOKED has completed in a concurrent thread. (*End of advice to users.*) 17

MPI\_WIN\_GET\_FAILED(win, failedgrp) 18

	IN	win	window object (handle)	19
				20
	OUT	failedgrp	(handle)	21

#### C binding 22

```
int MPI_Win_get_failed(MPI_Win win, MPI_Group *failedgrp) 23
```

#### Fortran 2008 binding 24

```
MPI_Win_get_failed(win, failedgrp, ierror) 25
  TYPE(MPI_Win), INTENT(IN) :: win 26
  TYPE(MPI_Group), INTENT(OUT) :: failedgrp 27
  INTEGER, OPTIONAL, INTENT(OUT) :: ierror 28
```

#### Fortran binding 29

```
MPI_WIN_GET_FAILED(WIN, FAILEDGRP, IERROR) 30
  INTEGER WIN, FAILEDGRP, IERROR 31
```

This local operation returns the group failedgrp of MPI processes from the window win that are locally known to have failed. The failedgrp can be empty, that is, equal to MPI\_GROUP\_EMPTY. 32

1 *Advice to users.* MPI makes no assumption about asynchronous progress of the  
 2 failure detection. A valid MPI implementation may choose to update the group of  
 3 locally known failed MPI processes only when it enters a synchronization function and  
 4 must raise a fault tolerance error. (*End of advice to users.*)

5  
 6 *Advice to users.* It is possible that only the calling MPI process has detected the  
 7 reported failure. If global knowledge is necessary, MPI processes detecting failures  
 8 should use the call MPI\_WIN\_REVOKE. (*End of advice to users.*)

### 9 16.3.3 I/O Functions

#### 10 MPI\_FILE\_REVOKE(fh)

11  
 12  
 13  
 14 IN fh file (handle)

#### 15 16 **C binding**

17 int MPI\_File\_revoke(MPI\_File fh)

#### 18 19 **Fortran 2008 binding**

20 MPI\_File\_revoke(fh, ierror)  
 21 TYPE(MPI\_File), INTENT(IN) :: fh  
 22 INTEGER, OPTIONAL, INTENT(OUT) :: ierror

#### 23 24 **Fortran binding**

25 MPI\_FILE\_REVOKE(FH, IERROR)  
 26 INTEGER FH, IERROR

27 This function notifies all MPI processes in the group associated with the file handle  
 28 fh that this file handle is revoked. The revocation of a file handle by any MPI process  
 29 completes non-local MPI operations on fh at all MPI processes by raising an error of class  
 30 MPI\_ERR\_REVOKED. This function is not collective and therefore does not have a matching  
 31 call on remote MPI processes. All non-failed MPI processes belonging to fh will be notified  
 32 of the revocation despite failures.

33 A file handle is revoked at a given MPI process either when MPI\_FILE\_REVOKE is  
 34 locally called on it, or when any MPI operation on fh raises an error of class  
 35 MPI\_ERR\_REVOKED at that process. Once a file handle has been revoked at an MPI pro-  
 36 cess, all subsequent non-local operations on that file handle are considered local and must  
 37 complete by raising an error of class MPI\_ERR\_REVOKED at that process.

#### 38 39 MPI\_FILE\_IS\_REVOKED(fh, flag)

40  
 41 IN fh file (handle)  
 42 OUT flag true if the file handle is revoked (logical)

#### 43 44 **C binding**

45 int MPI\_File\_is\_revoked(MPI\_File fh, int \*flag)

#### 46 47 **Fortran 2008 binding**

48 MPI\_File\_is\_revoked(fh, flag, ierror)

```

TYPE(MPI_File), INTENT(IN) :: fh
LOGICAL, INTENT(OUT) :: flag
INTEGER, OPTIONAL, INTENT(OUT) :: ierror

```

**Fortran binding**

```

MPI_FILE_IS_REVOKED(FH, FLAG, IERROR)
INTEGER FH, IERROR
LOGICAL FLAG

```

Returns `flag = true` if the file handle associated with `fh` is revoked at the calling process. It returns `flag = false` otherwise. The operation is local.

*Advice to users.* In a multithreaded application, a thread calling `MPI_FILE_IS_REVOKED` may return `flag = true` before the operation that raises the first exception of class `MPI_ERR_REVOKED` has completed in a concurrent thread. (*End of advice to users.*)

## 16.4 Fault Tolerance Error Codes and Classes

Among the error classes defined in Section 9.4, the following are **fault tolerance error** classes:

<code>MPI_ERR_PROC_FAILED</code>	The operation could not complete because of an MPI process failure (a fail-stop failure).
<code>MPI_ERR_PROC_FAILED_PENDING</code>	The operation was interrupted by an MPI process failure (a fail-stop failure). The request is still pending and the operation may be completed later.
<code>MPI_ERR_REVOKED</code>	The communication object used in the operation has been revoked.

Table 16.1: Fault tolerance error classes

## 16.5 Examples

### 16.5.1 Safe Communicator Creation

The example below illustrates how a new communicator can be safely created despite disruption by MPI process failures. A child communicator is created with `MPI_COMM_SPLIT`, then the global success of the operation is verified with `MPI_COMM_AGREE`. If any MPI `g` failed to create the child communicator handle, all MPI processes are notified by the value of the integer agreed on. MPI processes that had successfully created the child communicator handle destroy it, as it cannot be used consistently.

**Example 16.1** Fault Tolerant Communicator Split Example

```

int Comm_split_consistent(MPI_Comm parent, int color, int key, MPI_Comm* child)
{

```

```

1  rc = MPI_Comm_split(parent, color, key, child);
2  split_ok = (MPI_SUCCESS == rc);
3  MPI_Comm_agree(parent, &split_ok);
4  if(split_ok) {
5      /* All surviving processes have created the "child" comm
6       * It may contain supplementary failures and the first
7       * operation on it may raise an error, but it is a
8       * workable object that will yield well specified outcomes */
9      return MPI_SUCCESS;
10 }
11 else {
12     /* At least one process did not create the child comm properly
13      * if the local process did succeed in creating it, it disposes
14      * of it, as it is a broken, inconsistent object */
15     if(MPI_SUCCESS == rc) {
16         MPI_Comm_free(child);
17     }
18     return MPI_ERR_PROC_FAILED;
19 }
20 }
21

```

### 16.5.2 Obtaining the consistent group of failed processes

Users can invoke `MPI_COMM_GET_FAILED`, `MPI_WIN_GET_FAILED`, to obtain the group of failed MPI processes, as detected at the local MPI process. However, these operations are local, thereby the invocation of the same function at another MPI process can result in a different group of failed processes being returned.

In the following examples, we illustrate two different approaches that permit obtaining the consistent group of failed MPI processes across all MPI processes of a communicator. The first one employs `MPI_COMM_SHRINK` to create a temporary communicator excluding failed MPI processes. The second one employs `MPI_COMM_AGREE` to synchronize the set of acknowledged failures.

#### **Example 16.2** Fault-Tolerant Consistent Group of Failures Example (Shrink variant)

```

35 Comm_failure_allget(MPI_Comm c, MPI_Group * g) {
36     MPI_Comm s; MPI_Group c_grp, s_grp;
37
38     /* Using shrink to create a new communicator, the underlying
39      * group is necessarily consistent across all processes, and excludes
40      * all processes detected to have failed before the call */
41     MPI_Comm_shrink(c, &s);
42     /* Extracting the groups from the communicators */
43     MPI_Comm_group(c, &c_grp);
44     MPI_Comm_group(s, &s_grp);
45     /* s_grp is the group of still alive processes, we want to
46      * return the group of failed processes. */
47     MPI_Group_difference(c_grp, s_grp, g);
48

```

```

MPI_Group_free(&c_grp); MPI_Group_free(&s_grp);
MPI_Comm_free(&s);
}

```

### Example 16.3 Fault-Tolerant Consistent Group of Failures Example (Agree variant)

```

Comm_failure_allget2(MPI_Comm c, MPI_Group * g) {
    int rc; int T=1;
    int size; int nacked;
    MPI_Group gf;
    int ranges[3] = {0, 0, 1};

    MPI_Comm_size(c, &size);

    do {
        /* this routine is not pure: calling MPI_Comm_ack_failed
        * affects the state of the communicator c */
        MPI_Comm_ack_failed(c, size, &nacked);
        /* we simply ignore the T value in this example */
        rc = MPI_Comm_agree(c, &T);
    } while( rc != MPI_SUCCESS );
    /* after this loop, MPI_Comm_agree has returned MPI_SUCCESS at
    * all processes, so all processes have Acknowledged the same set of
    * failures. Let's get that set of failures in the g group. */
    if( 0 == nacked ) {
        *g = MPI_GROUP_EMPTY;
    }
    else {
        MPI_Comm_get_failed(c, &gf);
        ranges[1] = nacked - 1;
        MPI_Group_range_incl(gf, 1, ranges, g);
        MPI_Group_free(&gf);
    }
}

```

### 16.5.3 Fault-Tolerant Master/Worker

The example below presents a master code that handles worker failures by discarding failed worker MPI processes and resubmitting the work to the remaining workers. It demonstrates the different failure cases that may occur when posting receptions from `MPI_ANY_SOURCE` as discussed in the advice to users in Section ??.

### Example 16.4 Fault-Tolerant Master Example

```

int master(void)
{
    MPI_Comm_set_errhandler(comm, MPI_ERRORS_RETURN);
    MPI_Comm_size(comm, &size);
}

```

```

1  MPI_Comm_group(comm, &gcomm);
2
3  /* ... submit the initial work requests ... */
4
5  /* Progress engine: Get answers, send new requests,
6   and handle process failures */
7  MPI_Irecv( buffer, 1, MPI_INT, MPI_ANY_SOURCE, tag, comm, &req );
8  while( (active_workers > 0) && work_available ) {
9      rc = MPI_Wait( &req, &status );
10     if( MPI_SUCCESS == rc ) {
11         /* ... process the answer and update work_available ... */
12     }
13     else {
14         MPI_Error_class(rc, &ec);
15         if( (MPI_ERR_PROC_FAILED == ec) ||
16             (MPI_ERR_PROC_FAILED_PENDING == ec) ) {
17             /* We ack the full size of comm, so we will ack
18              * unconditionally. Variable gsize will contain all
19              * currently known failures. */
20             MPI_Comm_ack_failed(comm, size, &gsize);
21
22             /* ... find the lost work and requeue it ... */
23             MPI_Comm_get_failed(comm, &g);
24             granks = (int*) calloc(active_workers-gsize-1, sizeof(int));
25             cranks = (int*) calloc(active_workers-gsize-1, sizeof(int));
26             for(i = active_workers; i < gsize; i++)
27                 granks[i-active_workers] = i;
28             MPI_Group_translate_ranks(g, gsize, granks, gcomm, cranks);
29             /* iterate over newly failed procs */
30             for(i = active_workers; i < gsize; i++) {
31                 /* resubmit the work */
32             }
33             free(cranks); free(granks);
34             MPI_Group_free(&g);
35
36             active_workers = size - gsize - 1;
37
38             /* no need to repost when the request is still pending */
39             if( ec == MPI_ERR_PROC_FAILED_PENDING )
40                 continue;
41         }
42     }
43     /* get ready to receive more notifications from workers */
44     MPI_Irecv( buffer, 1, MPI_INT, MPI_ANY_SOURCE, tag, comm, &req );
45 }
46 /* ... cancel request and cleanup ... */
47 }
48

```

## 16.5.4 Fault-Tolerant Iterative Refinement

The example below demonstrates a method of fault tolerance for detecting and handling failures. At each iteration, the algorithm checks the return code of the `MPI_ALLREDUCE`. If the return code indicates a process failure for at least one MPI process, the algorithm revokes the communicator, agrees on the presence of failures, and shrinks it to create a new communicator. By calling `MPI_COMM_REVOKE`, the algorithm ensures that all MPI processes will be notified of process failure and enter the `MPI_COMM_AGREE`. If an MPI process fails, the algorithm must complete at least one more iteration to ensure a correct answer.

**Example 16.5** Fault-tolerant iterative refinement with shrink and agreement

```

while( gnorm > epsilon ) {
    /* Add a computation iteration to converge and
       compute local norm in lnorm */
    rc = MPI_Allreduce(&lnorm, &gnorm, 1, MPI_DOUBLE, MPI_MAX, comm);
    MPI_Error_class(rc, &ec);

    if( (MPI_ERR_PROC_FAILED == ec) ||
        (MPI_ERR_REVOKED == ec) ||
        (gnorm <= epsilon) ) {

        /* This process detected a failure, but other processes may have
           * proceeded into the next MPI_Allreduce. Since this process
           * will not match that following MPI_Allreduce, these other
           * processes would be at risk of deadlocking. This process thus
           * calls MPI_Comm_revoke to interrupt other processes and notify
           * them that it has detected a failure and is leaving the
           * failure free execution path to go into recovery. */
        if( MPI_ERR_PROC_FAILED == ec )
            MPI_Comm_revoke(comm);

        /* About to leave: let's be sure that everybody
           received the same information */
        allsucceeded = (rc == MPI_SUCCESS);
        rc = MPI_Comm_agree(comm, &allsucceeded);
        MPI_Error_class(rc, &ec);
        if( ec == MPI_ERR_PROC_FAILED || !allsucceeded ) {
            MPI_Comm_shrink(comm, &comm2);
            MPI_Comm_free(comm); /* Release the revoked communicator */
            comm = comm2;
            gnorm = epsilon + 1.0; /* Force one more iteration */
        }
    }
}

```

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48

# Bibliography



# Index

- MPI\_ANY\_SOURCE, [2](#), [3](#), [12](#), [19](#)
- MPI\_Comm, [8–13](#)
- MPI\_COMM\_NULL, [4–6](#)
- MPI\_COMM\_WORLD, [2](#), [5](#)
- MPI\_ERR\_PROC\_FAILED, [3](#), [4](#), [6–8](#), [10](#), [12](#), [13](#), [17](#)
- MPI\_ERR\_PROC\_FAILED\_PENDING, [3](#), [10](#), [12](#), [17](#)
- MPI\_ERR\_REVOKED, [3](#), [4](#), [6–10](#), [13–17](#)
- MPI\_ERR\_SPAWN, [5](#)
- MPI\_ERRORS\_ARE\_FATAL, [2](#)
- MPI\_File, [16](#)
- MPI\_FILE\_NULL, [8](#)
- MPI\_FT, [2](#)
- MPI\_Group, [10](#), [15](#)
- MPI\_GROUP\_EMPTY, [11](#), [15](#)
- MPI\_Request, [13](#)
- MPI\_SEEK\_SET, [7](#)
- MPI\_SUCCESS, [4](#), [6](#), [7](#), [13](#)
- MPI\_Win, [14](#), [15](#)
- MPI\_WIN\_NULL, [7](#)
  
- EXAMPLES:Comm\_failure\_allget example, [18](#)
- EXAMPLES:Comm\_failure\_allget2 example, [19](#)
- EXAMPLES:Fault-tolerant iterative refinement with shrink and agreement, [21](#)
- EXAMPLES:Master example, [19](#)
- EXAMPLES:MPI\_COMM\_ACK\_FAILED, [19](#)
- EXAMPLES:MPI\_COMM\_AGREE, [17](#), [19](#), [21](#)
- EXAMPLES:MPI\_COMM\_FREE, [17](#), [18](#), [21](#)
- EXAMPLES:MPI\_COMM\_GET\_FAILED, [19](#)
- EXAMPLES:MPI\_COMM\_GROUP, [18](#)
- EXAMPLES:MPI\_COMM\_REVOKE, [21](#)
- EXAMPLES:MPI\_COMM\_SHRINK, [18](#), [21](#)
- EXAMPLES:MPI\_COMM\_SPLIT, [17](#)
- EXAMPLES:MPI\_GROUP\_DIFFERENCE, [18](#)
  
- EXAMPLES:MPI\_GROUP\_FREE, [18](#)
  
- "mpi\_error\_range", [2](#)
- "mpi\_error\_uniform", [4](#)
- "coll", [4](#)
- "create", [4](#)
- "global", [3](#)
- "group", [3](#)
- "local", [4](#)
- "operation", [3](#)
  
- MPI\_ALLREDUCE, [21](#)
- MPI\_BCAST, [3](#)
- MPI\_COMM\_ACCEPT, [5](#)
- MPI\_COMM\_ACK\_FAILED, [12](#), [13](#)
- MPI\_COMM\_ACK\_FAILED(comm, nack, nacked), [11](#)
- MPI\_COMM\_AGREE, [4](#), [7](#), [8](#), [12–14](#), [17](#), [18](#), [21](#)
- MPI\_COMM\_AGREE(comm, flag), [12](#)
- MPI\_COMM\_CONNECT, [5](#)
- MPI\_COMM\_DISCONNECT, [6](#)
- MPI\_COMM\_DUP, [3](#), [4](#)
- MPI\_COMM\_FREE, [4](#)
- MPI\_COMM\_GET\_FAILED, [12](#), [13](#), [18](#)
- MPI\_COMM\_GET\_FAILED(comm, failedgrp), [11](#)
- MPI\_COMM\_GET\_PARENT, [5](#)
- MPI\_COMM\_IAGREE, [8](#)
- MPI\_COMM\_IAGREE(comm, flag, request), [13](#)
- MPI\_COMM\_IDUP, [10](#)
- MPI\_COMM\_IS\_REVOKED, [9](#)
- MPI\_COMM\_IS\_REVOKED(comm, flag), [9](#)
- MPI\_COMM\_ISHRINK, [10](#)
- MPI\_COMM\_ISHRINK(comm, newcomm, request), [10](#)
- MPI\_COMM\_JOIN, [5](#)
- MPI\_COMM\_REVOKE, [8](#), [11](#), [21](#)
- MPI\_COMM\_REVOKE(comm), [8](#)

1 MPI\_COMM\_SHRINK, [8–10](#), [12](#), [18](#)  
2 MPI\_COMM\_SHRINK(comm, newcomm), [9](#)  
3 MPI\_COMM\_SPAWN, [5](#)  
4 MPI\_COMM\_SPAWN\_MULTIPLE, [5](#)  
5 MPI\_COMM\_SPLIT, [3](#), [10](#), [17](#)  
6 MPI\_FILE\_CLOSE, [7](#), [8](#)  
7 MPI\_FILE\_IS\_REVOKED, [17](#)  
8 MPI\_FILE\_IS\_REVOKED(fh, flag), [16](#)  
9 MPI\_FILE\_REVOKE, [16](#)  
10 MPI\_FILE\_REVOKE(fh), [16](#)  
11 MPI\_FILE\_SEEK, [7](#)  
12 MPI\_GET, [6](#)  
13 MPI\_INIT, [5](#)  
14 MPI\_PUT, [6](#)  
15 MPI\_WIN\_FLUSH, [6](#)  
16 MPI\_WIN\_FLUSH\_LOCAL, [6](#)  
17 MPI\_WIN\_FLUSH\_LOCAL\_ALL, [6](#)  
18 MPI\_WIN\_FREE, [7](#)  
19 MPI\_WIN\_GET\_FAILED, [18](#)  
20 MPI\_WIN\_GET\_FAILED(win, failedgrp), [15](#)  
21 MPI\_WIN\_IS\_REVOKED, [15](#)  
22 MPI\_WIN\_IS\_REVOKED(win, flag), [15](#)  
23 MPI\_WIN\_REVOKE, [7](#), [14](#), [16](#)  
24 MPI\_WIN\_REVOKE(win), [14](#)  
25  
26 TERM:error handling  
27     fault tolerance, [1](#)  
28         ack, [12](#)  
29         agree, [13](#), [14](#)  
30         communicator, [3](#), [8](#)  
31         dynamic process, [5](#)  
32         fault tolerance error, [1](#), [17](#)  
33         I/O, [7](#), [16](#)  
34         inquiry, [2](#)  
35         mitigation, [8](#)  
36         notification, [2](#), [17](#)  
37         one-sided, [6](#), [14](#)  
38         revoke, [8](#), [9](#), [14–17](#)  
39         shrink, [9](#), [10](#)  
40     process failure, [1](#)  
41  
42  
43  
44  
45  
46  
47  
48