

*D R A F T*

Document for a Standard Message-Passing Interface

Message Passing Interface Forum

October 31, 2022

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

# 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 (e.g., communicators, etc.) as needed to restore the full capability of MPI to carry out 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 eventually 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. If an application needs comprehensive 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, see Section 16.2.2.

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 must be portable across MPI implementations (including those which do not provide fault tolerance). If an MPI implementation does not provide fault tolerance features, it may exhibit undefined behavior after a process failure at any MPI process. Fault tolerant applications may determine if the implementation supports fault tolerance by querying the predefined attribute `MPI_FT` on `MPI_COMM_WORLD` (see 9.1.2.)

*Advice to users.* The MPI standard does not specify transparent process recovery

1 upon MPI process failure. In particular, restoring the lost dataset, spawning spare  
2 processes or taking other recovery actions are the responsibility of the user.

3 Many of the operations and semantics described in this chapter are applicable only  
4 when the MPI application has replaced the default error handler  
5 `MPI_ERRORS_ARE_FATAL` on the communicators and windows it uses. (*End of advice*  
6 *to users.*)  
7

## 8 9 16.2 Failure Notification

10 All MPI procedures must eventually return (possibly by raising one of the process failure  
11 error classes listed in Section 16.4) even if they involve a failed MPI process.

12  
13 *Advice to implementors.* As long as an implementation can successfully complete  
14 some operations during a completing procedure, it may choose to delay raising an  
15 error. Another valid implementation might choose to raise an error as quickly as  
16 possible. (*End of advice to implementors.*)  
17

18 When an operation raises a process failure error it may not satisfy its specification (like  
19 any other error, see 9.4). Note that the remainder of this chapter defines operations that  
20 maintain full specification semantic after raising a fault tolerance error; such exceptions will  
21 be explicitly stated.

22 Operations do not raise fault tolerance errors during their initialization, initiation, or  
23 starting stages. The corresponding completion stage raises a fault tolerance error when  
24 appropriate.  
25

26 *Advice to users.* Persistent collective communication initiation procedures are them-  
27 selves collective operations, and may therefore raise a process failure error of their  
28 own. (*End of advice to users.*)  
29

### 30 16.2.1 Revoked Communicator

31 A communicator is **revoked** at a given MPI process either when `MPI_COMM_REVOKE` is  
32 locally called on it, or when any MPI operation on the communicator raises an error of class  
33 `MPI_ERR_REVOKED` at that process. Once a communicator has been revoked at an MPI  
34 process, all subsequent non-local operations on that communicator are considered local and  
35 must complete by raising an error of class `MPI_ERR_REVOKED` at that MPI process.  
36

37 *Advice to implementors.* An implementation may detect that a communicator should  
38 become revoked when completing procedures can still complete some operations suc-  
39 cessfully. The implementation may delay raising the `MPI_ERR_REVOKED` error to a  
40 later operation, but it cannot delay revoking the communicator indefinitely and should  
41 favor raising the error promptly. (*End of advice to implementors.*)  
42

### 43 16.2.2 Error Reporting Range

44 By default, process failure errors are raised only during communication operations in which  
45 a failed process is involved, but users can control the range of MPI processes whose failure  
46 cause MPI operations to raise errors on their communicators by setting the following values  
47 to the info key "mpi\_error\_range" on their communicators:  
48

"operation" A process fault tolerance error shall be raised only if the operation involves a failed MPI process (e.g., a point-to-point message between two non-failed MPI processes shall not raise a process fault tolerance error). An MPI process is considered involved in an operation (for the purpose of this definition) if its failure may prevent it from calling a matching procedure for the operation, that is, an MPI process is involved in a communication if any of the following is true: the process is in the group(s) over which the operation is collective; the process is a destination or a specified or a matched source in a point-to-point communication; the process belongs to the source group in an MPI\_ANY\_SOURCE receive operation; or the process is a specified target in a remote memory operation.

"group" The failure of any MPI process in the group of the communicator causes the communicator to become revoked (as if MPI\_COMM\_REVOKE had been called on the communicator).

"universe" The failure of any MPI process in the MPI universe causes the communicator to become revoked (as if MPI\_COMM\_REVOKE had been called on the communicator).

When the info key "mpi\_error\_range" is not set, it is equivalent to having the key set to "operation".

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

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

- MPI\_ERR\_REVOKED indicates that the communicator is revoked.
- MPI\_ERR\_PROC\_FAILED\_PENDING indicates, for a nonblocking or persistent 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. The request remains active.
- 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 representing a point-to-point communication, it is completed.

### 16.2.4 Fault Tolerance Errors in Collective Operations

A collective operation raises errors of the following classes to notify users that the communication could not complete successfully because of the failure of at least one involved MPI process:

- MPI\_ERR\_REVOKED indicates that the communicator is revoked.
- MPI\_ERR\_PROC\_FAILED indicates that the failure prevents the operation from following its failure-free specification. If there is a request representing a collective communication, it is completed.

1       *Advice to users.*

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

7       (*End of advice to users.*)

9       *Advice to users.*

10       Note that communicator creation functions (e.g., MPI\_COMM\_DUP or  
11       MPI\_COMM\_SPLIT) are collective operations. As such, if a failure happened during  
12       the procedure, an error might be raised at some MPI processes while others succeed  
13       and obtain a new communicator handle. Although it is valid to communicate be-  
14       tween MPI processes that succeeded in creating the new communicator handle, the  
15       user is responsible for ensuring a consistent view of the communicator creation, if  
16       needed. A conservative solution is to check the global outcome of the communicator  
17       creation function with MPI\_COMM\_AGREE (defined in Section 16.3.1), as illustrated  
18       in Example 16.1. (*End of advice to users.*)

19       After an MPI process failure, MPI\_COMM\_FREE (as with all other collective operations)  
20       may not complete successfully at all MPI processes. For any MPI process that receives the  
21       return code MPI\_SUCCESS, the behavior is defined in Section 7.4.3. If an MPI process  
22       raises a process failure error (classes MPI\_ERR\_PROC\_FAILED or MPI\_ERR\_REVOKED), the  
23       communicator handle `comm` is set to MPI\_COMM\_NULL; however, the implementation makes  
24       no guarantee about the success or failure of the MPI\_COMM\_FREE operation, locally or  
25       remotely.  
26       remotely.

27         
28       *Advice to users.* Users are encouraged to call MPI\_COMM\_FREE on communicators  
29       they do not wish to use anymore, even when they contain failed MPI processes. Al-  
30       though the operation may raise a fault tolerance error and not synchronize properly,  
31       this gives a high quality implementation an opportunity to release local resources and  
32       memory consumed by the object. (*End of advice to users.*)  
33       

## 34       Error Uniformity

35       As noted above, by default, collective communication do not enforce uniformity in error  
36       raising across MPI processes. Despite the performance advantages that non-uniformity  
37       offers, a common usage pattern in applications is to transform non-uniform error raising  
38       into a uniform behavior across all MPI processes of the group.

39       Users can set collective operations on a communicator to enforce uniform error raising  
40       by setting the following values in the info key "mpi\_error\_uniform" on the communicator:  
41       

42         
43       "local" Process fault tolerance errors are raised to indicate that an MPI process failure  
44       prevents from guaranteeing the specified behavior at this process for the collective  
45       communication operation (e.g., the output buffer contains invalid data). Other MPI  
46       processes may have satisfied their specification (e.g., the output buffer is valid at that  
47       process) and may have returned MPI\_SUCCESS.  
48

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

"construct" Process fault tolerance errors are raised to indicate that an MPI process failure prevents from guaranteeing the specified behavior at any process for the collective MPI communication context constructor/destructor operation (e.g., MPI\_COMM\_DUP could not create a new communicator at any process in the group of the communicator). Non-synchronizing context constructor/destructor operations become synchronizing. Collective communication that do not create or free a communication context are not impacted.

When the info key "mpi\_error\_uniform" is not set on the communicator, it is equivalent to having the key set to "local". The info key "mpi\_error\_uniform" must be set to the same value at all MPI processes for a given communicator.

### 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,

1 and instead delay raising a fault tolerance error to a later communication involving  
2 this process. (*End of advice to implementors.*)

3  
4 *Advice to users.* To determine how many new MPI processes have effectively been  
5 spawned, the normal semantics for hard and soft spawn applies: if the requested  
6 number of processes is unavailable for a hard spawn, an error of class `MPI_ERR_SPAWN`  
7 is raised (possibly leaving MPI in an undefined state), and an appropriate error code  
8 is set in the `array_of_errcodes` parameter. Note however that an implementation may  
9 report that it has spawned the requested number of MPI processes even when some  
10 MPI processes have failed before exiting `MPI_INIT`. This condition can be detected  
11 by communicating over the created intercommunicator with these processes. (*End of*  
12 *advice to users.*)

13  
14 *Advice to implementors.* `MPI_COMM_JOIN` does not require any supplementary  
15 semantics. When the remote MPI process on the fd socket has failed, the operation  
16 succeeds and sets `intercomm` to `MPI_COMM_NULL`. (*End of advice to implementors.*)

17  
18 After an MPI process failure, `MPI_COMM_DISCONNECT` (as with all other collective  
19 operations) may not complete successfully at all MPI processes. For any process that re-  
20 ceives the return code `MPI_SUCCESS`, the behavior is defined in 11.10.4. If an MPI process  
21 raises a fault tolerance error (classes `MPI_ERR_PROC_FAILED` or `MPI_ERR_REVOKED`), the  
22 communicator handle `comm` is set to `MPI_COMM_NULL`; however, the implementation makes  
23 no guarantee about the success or failure of the `MPI_COMM_DISCONNECT` operation, lo-  
24 cally or remotely.

25  
26 *Advice to users.* Users are encouraged to call `MPI_COMM_DISCONNECT` on com-  
27 municators they do not wish to use anymore, even when they contain failed MPI  
28 processes. Although the operation may raise a fault tolerance error and not synchro-  
29 nize properly, this gives a high quality implementation an opportunity to release local  
30 resources and memory consumed by the object. (*End of advice to users.*)

### 31 16.2.5 Fault Tolerance Errors in One-Sided Communication

32  
33 When an operation on a window raises a fault tolerance error, the state of all data held  
34 in memory exposed by that window becomes undefined at all MPI processes for which  
35 a one-sided communication operation could have modified local data in that window (a  
36 target in a remote write, or accumulate operation, or an origin in a remote read operation),  
37 and the operation completion has not been semantically guaranteed (as an example by a  
38 successful synchronization between the origin and the target, after the origin had issued an  
39 `MPI_WIN_FLUSH`).

40  
41 *Advice to users.* Assessing if a particular portion of the exposed memory remains  
42 correct is the responsibility of the user. Note that in passive target mode, when an  
43 error is raised at the origin, the target memory data may become undefined before a  
44 synchronization raises an error at the target.

45  
46 The exposed memory data becomes undefined for all uses, not only the window in  
47 which the error was raised. Any overlapping windows or uses involving shared memory  
48 also read undefined data (even if they do not call MPI procedures). (*End of advice to*  
*users.*)



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

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

*Advice to implementors.* Non-synchronizing operations (`MPI_WIN_FLUSH_LOCAL`, `MPI_WIN_FLUSH_LOCAL_ALL`) are not required to raise a fault tolerance error. (*End 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 Fault Tolerance Errors in File I/O operations

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.

1 If an MPI process failure prevents a file operation from completing, an MPI error of  
 2 class `MPI_ERR_PROC_FAILED` is raised. Once an MPI implementation has raised an error of  
 3 class `MPI_ERR_PROC_FAILED`, the state of the file pointers involved in the operation that  
 4 raised the error is *undefined*.

5  
 6 *Advice to users.* Since collective I/O operations may not synchronize with other MPI  
 7 processes, process failures may not be reported during a collective I/O operation.  
 8 Users are encouraged to use `MPI_COMM_AGREE` on a communicator containing the  
 9 same group as the file handle when they need to deduce the completion status of  
 10 collective operations on file handles and maintain a consistent view of file pointers.  
 11 The file pointer can be reset by using `MPI_FILE_SEEK` with the `MPI_SEEK_SET` update  
 12 mode. (*End of advice to users.*)

13  
 14 After an MPI process failure, `MPI_FILE_CLOSE` (as with all other collective operations)  
 15 may not complete successfully at all MPI processes. For any MPI process that receives the  
 16 return code `MPI_SUCCESS`, the behavior is defined in Section 14.2.2. If an MPI process  
 17 raises a process failure error (classes `MPI_ERR_PROC_FAILED` or `MPI_ERR_REVOKED`), the  
 18 file handle `fh` is set to `MPI_FILE_NULL`; however, the implementation makes no guarantee  
 19 about the success or failure of the `MPI_FILE_CLOSE` operation, locally or remotely.

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

## 27 16.3 Failure Mitigation Procedures

### 29 16.3.1 Communicator Procedures

30  
 31 Process failure notification is not global in MPI. MPI processes that do not call operations  
 32 involving a failed MPI process are possibly never notified of its failure (see Section 16.2). If  
 33 a notification must be propagated, MPI provides a function to revoke a communicator at  
 34 all members.

36 `MPI_COMM_REVOKE(comm)`

37     IN           comm                           communicator (handle)

#### 40 **C binding**

41 `int MPI_Comm_revoke(MPI_Comm comm)`

#### 42 **Fortran 2008 binding**

43 `MPI_Comm_revoke(comm, ierror)`

44     TYPE(MPI\_Comm), INTENT(IN) :: comm

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

#### 47 **Fortran binding**

48 `MPI_COMM_REVOKE(COMM, IERROR)`

INTEGER COMM, IERROR

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

`MPI_COMM_IS_REVOKED(comm, flag)`

|     |                   |   |
|-----|-------------------|---|
| IN  | <code>comm</code> | communicator (handle)                         |
| OUT | <code>flag</code> | true if the communicator is revoked (logical) |

### C binding

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

### Fortran 2008 binding

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

### Fortran binding

```
MPI_COMM_IS_REVOKED(COMM, FLAG, IERROR)
  INTEGER COMM, IERROR
  LOGICAL FLAG
```

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.

*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.*)

`MPI_COMM_SHRINK(comm, newcomm)`

|     |                      |                       |
|-----|----------------------|-----------------------|
| IN  | <code>comm</code>    | communicator (handle) |
| OUT | <code>newcomm</code> | communicator (handle) |

### C binding

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

### Fortran 2008 binding

```
MPI_Comm_shrink(comm, newcomm, ierror)
  TYPE(MPI_Comm), INTENT(IN) :: comm
```

```

1     TYPE(MPI_Comm), INTENT(OUT) :: newcomm
2     INTEGER, OPTIONAL, INTENT(OUT) :: ierror

```

### Fortran binding

```

4 MPI_COMM_SHRINK(COMM, NEWCOMM, IERROR)
5     INTEGER COMM, NEWCOMM, IERROR

```

7 This collective operation creates a new intra- or intercommunicator  
8 `newcomm` from the intra- or intercommunicator `comm`, respectively, by excluding the group  
9 of failed MPI processes as agreed upon during the operation. The groups of  
10 `newcomm` must include every MPI process that returns from `MPI_COMM_SHRINK`, and it  
11 must exclude every MPI process whose failure caused an operation on `comm` to raise an MPI  
12 error of class `MPI_ERR_PROC_FAILED` or `MPI_ERR_PROC_FAILED_PENDING` at a member of  
13 the groups of `newcomm`, before that member initiated `MPI_COMM_SHRINK`. This procedure  
14 is semantically equivalent to an `MPI_COMM_SPLIT` operation that would succeed despite  
15 failures, where members of the groups of `newcomm` participate with the same color and a  
16 key equal to their rank in `comm`.

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

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

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

```

28 MPI_COMM_ISHRINK(comm, newcomm, request)

```

```

29     IN      comm      communicator (handle)
30     OUT     newcomm   communicator (handle)
31     OUT     request   communication request (handle)

```

### C binding

```

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

```

### Fortran 2008 binding

```

37 MPI_Comm_ishrink(comm, newcomm, request, ierror)
38     TYPE(MPI_Comm), INTENT(IN) :: comm
39     TYPE(MPI_Comm), INTENT(OUT), ASYNCHRONOUS :: newcomm
40     TYPE(MPI_Request), INTENT(OUT) :: request
41     INTEGER, OPTIONAL, INTENT(OUT) :: ierror

```

### Fortran binding

```

44 MPI_COMM_ISHRINK(COMM, NEWCOMM, REQUEST, IERROR)
45     INTEGER COMM, NEWCOMM, REQUEST, IERROR

```

47 `MPI_COMM_ISHRINK` is a nonblocking variant of `MPI_COMM_SHRINK`. With the  
48 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)

|     |           |                                    |
|-----|-----------|------------------------------------|
| IN  | comm      | communicator (handle)              |
| OUT | failedgrp | group of failed processes (handle) |

### 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 procedure returns the group failedgrp of MPI processes from the communicator comm that are known to have failed by the calling MPI process. Any MPI process whose failure caused a procedure on comm to raise a process failure error at this MPI process must appear in that group. The failedgrp can be empty, that is, equal to MPI\_GROUP\_EMPTY.

When an MPI process is a member of failedgrp, it is also a member, with the same rank, in any group produced by a subsequent call to MPI\_COMM\_GET\_FAILED on comm at the same calling process.

*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. (*End of advice to users.*)

MPI\_COMM\_ACK\_FAILED(comm, num\_to\_ack, num\_acked)

|     |            |   |
|-----|------------|---|
| IN  | comm       | communicator (handle)                                       |
| IN  | num_to_ack | Maximum number of process failures to acknowledge (integer) |
| OUT | num_acked  | Number of process failures acknowledged (integer)           |

### C binding

```
int MPI_Comm_ack_failed(MPI_Comm comm, int num_to_ack, int *num_acked)
```

**Fortran 2008 binding**

```

1 MPI_Comm_ack_failed(comm, num_to_ack, num_acked, ierror)
2     TYPE(MPI_Comm), INTENT(IN) :: comm
3     INTEGER, INTENT(IN) :: num_to_ack
4     INTEGER, INTENT(OUT) :: num_acked
5     INTEGER, OPTIONAL, INTENT(OUT) :: ierror

```

**Fortran binding**

```

8 MPI_COMM_ACK_FAILED(COMM, NUM_TO_ACK, NUM_ACKED, IERROR)
9     INTEGER COMM, NUM_TO_ACK, NUM_ACKED, IERROR

```

This local procedure is used to **acknowledge** locally notified failures on **comm**. The operation acknowledges the first **num\_to\_ack** process failures on **comm**, that is, it acknowledges the failure of members with a rank lower than **num\_to\_ack** in the group that would be produced by a concurrent call to **MPI\_COMM\_GET\_FAILED** on the same **comm**.

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

**num\_acked** can be larger than **num\_to\_ack** when process failures have been acknowledged in a prior call to **MPI\_COMM\_ACK\_FAILED**.

After an MPI process failure is acknowledged on **comm**, unmatched **MPI\_ANY\_SOURCE** receive operations on the same **comm** that would have raised an error of class **MPI\_ERR\_PROC\_FAILED\_PENDING** (see Section 16.2.3) proceed without further raising errors due to this acknowledged failure. Also, **MPI\_COMM\_AGREE** on the same **comm** will not raise an error of class **MPI\_ERR\_PROC\_FAILED** due to this acknowledged failure (according to the specification found later in this section).

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

Calling **MPI\_COMM\_ACK\_FAILED** on a communicator with failed MPI processes has no effect on collective operations (except for **MPI\_COMM\_AGREE**). If a collective operation would raise an error due to the communicator containing a failed process (as defined in Section 16.2.4), it can continue to raise an error even after the failure has been acknowledged. In order to use collective operations between MPI processes of a communicator that contains failed MPI processes, users should create a new communicator by calling **MPI\_COMM\_SHRINK**. (*End of advice to users.*)

```

45 MPI_COMM_AGREE(comm, flag)

```

```

46     IN      comm      communicator (handle)
47     INOUT  flag      bitwise 'AND' of contributed values (integer)

```

**C binding**

```
int MPI_Comm_agree(MPI_Comm comm, int *flag)
```

**Fortran 2008 binding**

```
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.*)

1 MPI\_COMM\_IAGREE(comm, flag, request)

2     IN        comm                   communicator (handle)  
 3  
 4     INOUT    flag                   bitwise ‘AND’ of contributed values (integer)  
 5     OUT       request               communication request (handle)

6  
 7 **C binding**

8 int MPI\_Comm\_iagree(MPI\_Comm comm, int \*flag, MPI\_Request \*request)

9  
 10 **Fortran 2008 binding**

11 MPI\_Comm\_iagree(comm, flag, request, ierror)  
 12     TYPE(MPI\_Comm), INTENT(IN) :: comm  
 13     INTEGER, INTENT(INOUT), ASYNCHRONOUS :: flag  
 14     TYPE(MPI\_Request), INTENT(OUT) :: request  
 15     INTEGER, OPTIONAL, INTENT(OUT) :: ierror

16 **Fortran binding**

17 MPI\_COMM\_IAGREE(COMM, FLAG, REQUEST, IERROR)  
 18     INTEGER COMM, FLAG, REQUEST, IERROR

19  
 20     This function has the same semantics as MPI\_COMM\_AGREE except that it is non-  
 21 blocking.

22  
 23 16.3.2 One-Sided Procedures

24  
 25  
 26 MPI\_WIN\_REVOKE(win)

27     IN        win                   window object (handle)

28  
 29  
 30 **C binding**

31 int MPI\_Win\_revoke(MPI\_Win win)

32 **Fortran 2008 binding**

33 MPI\_Win\_revoke(win, ierror)  
 34     TYPE(MPI\_Win), INTENT(IN) :: win  
 35     INTEGER, OPTIONAL, INTENT(OUT) :: ierror

36 **Fortran binding**

37 MPI\_WIN\_REVOKE(WIN, IERROR)  
 38     INTEGER WIN, IERROR

39  
 40     This function notifies all MPI processes in the group associated with the window win  
 41 that this window is revoked. The revocation of a window by any MPI process completes  
 42 RMA operations on win at all MPI processes and RMA synchronizations on win raise an error  
 43 of class MPI\_ERR\_REVOKED. This function is not collective and therefore does not have a  
 44 matching call on remote MPI processes. All non-failed MPI processes belonging to win will  
 45 be notified of the revocation despite failures.

46     A window is revoked at a given MPI process either when MPI\_WIN\_REVOKE is locally  
 47 called on it, or when any MPI operation on win raises an error of class MPI\_ERR\_REVOKED  
 48



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)

|     |      |   |
|-----|------|---|
| IN  | win  | window object (handle)                  |
| OUT | flag | true if the window is revoked (logical) |

### C binding

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

### Fortran 2008 binding

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

### Fortran binding

```
MPI_WIN_IS_REVOKED(WIN, FLAG, IERROR)
  INTEGER WIN, IERROR
  LOGICAL FLAG
```

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.

*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.*)

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

|     |           |                        |
|-----|-----------|------------------------|
| IN  | win       | window object (handle) |
| OUT | failedgrp | (handle)               |

### C binding

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

### Fortran 2008 binding

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

### Fortran binding

```
MPI_WIN_GET_FAILED(WIN, FAILEDGRP, IERROR)
```

1       INTEGER WIN, FAILEDGRP, IERROR

2  
3       This local operation returns the group `failedgrp` of MPI processes from the window  
4 `win` that are locally known to have failed. The `failedgrp` can be empty, that is, equal to  
5 `MPI_GROUP_EMPTY`.

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

10  
11       *Advice to users.* It is possible that only the calling MPI process has detected the  
12 reported failure. If global knowledge is necessary, MPI processes detecting failures  
13 should use the call `MPI_WIN_REVOKE`. (*End of advice to users.*)

### 14 16.3.3 File I/O Procedures

15  
16  
17  
18 `MPI_FILE_REVOKE(fh)`

19       IN        fh                               file (handle)

#### 20 **C binding**

21  
22 `int MPI_File_revoke(MPI_File fh)`

#### 23 **Fortran 2008 binding**

24 `MPI_File_revoke(fh, ierror)`  
25       TYPE(MPI\_File), INTENT(IN) :: fh  
26       INTEGER, OPTIONAL, INTENT(OUT) :: ierror

#### 27 **Fortran binding**

28 `MPI_FILE_REVOKE(FH, IERROR)`  
29       INTEGER FH, IERROR

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

37  
38       A file handle is revoked at a given MPI process either when `MPI_FILE_REVOKE` is lo-  
39 cally called on it, or when any MPI operation on `fh` raises an error of class `MPI_ERR_REVOKED`  
40 at that process. Once a file handle has been revoked at an MPI process, all subsequent non-  
41 local operations on that file handle are considered local and must complete by raising an  
42 error of class `MPI_ERR_REVOKED` at that process.

|                               |      |  |   |
|-------------------------------|------|--|---|
| MPI_FILE_IS_REVOKED(fh, flag) |      |  | 1 |
| IN                            | fh   | file (handle)                                | 2 |
|                               |      |  | 3 |
| OUT                           | flag | true if the file handle is revoked (logical) | 4 |
|                               |      |  | 5 |

**C binding**

```
int MPI_File_is_revoked(MPI_File fh, int *flag)
```

**Fortran 2008 binding**

```
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:

|                             |   |
|-----------------------------|---|
| MPI_ERR_PROC_FAILED         | The operation could not complete because of an MPI process failure (a fail-stop failure).   |
| MPI_ERR_PROC_FAILED_PENDING | 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. |
| MPI_ERR_REVOKED             | 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 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

```

12 int Comm_split_consistent(MPI_Comm parent, int color, int key,
13                          MPI_Comm* child)
14 {
15     rc = MPI_Comm_split(parent, color, key, child);
16     split_ok = (MPI_SUCCESS == rc);
17     MPI_Comm_agree(parent, &split_ok);
18     if(split_ok) {
19         /* All surviving processes have created the "child" comm
20          * It may contain supplementary failures and the first
21          * operation on it may raise an error, but it is a
22          * workable object that will yield well specified outcomes */
23         return MPI_SUCCESS;
24     }
25     else {
26         /* At least one process did not create the child comm properly
27          * if the local process did succeed in creating it, it disposes
28          * of it, as it is a broken, inconsistent object */
29         if(MPI_SUCCESS == rc) {
30             MPI_Comm_free(child);
31         }
32         return MPI_ERR_PROC_FAILED;
33     }
34 }

```

### 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)

```

47 Comm_failure_allget(MPI_Comm c, MPI_Group * g) {
48

```

```

MPI_Comm s; MPI_Group c_grp, s_grp;

/* Using shrink to create a new communicator, the underlying
 * group is necessarily consistent across all processes, and
 * excludes all processes detected to have failed before the
 * call */
MPI_Comm_shrink(c, &s);
/* Extracting the groups from the communicators */
MPI_Comm_group(c, &c_grp);
MPI_Comm_group(s, &s_grp);
/* s_grp is the group of still alive processes, we want to
 * return the group of failed processes. */
MPI_Group_difference(c_grp, s_grp, g);

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 num_acked;
    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, &num_acked);
        /* 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 == num_acked ) {
        *g = MPI_GROUP_EMPTY;
    }
    else {
        MPI_Comm_get_failed(c, &gf);
        ranges[1] = num_acked - 1;
        MPI_Group_range_incl(gf, 1, ranges, g);
        MPI_Group_free(&gf);
    }
}

```

## 16.5.3 Fault Tolerant Manager/Worker

The example below presents a manager 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 with receive operations using MPI\_ANY\_SOURCE as discussed in Section 16.2.3.

**Example 16.4.** Fault Tolerant Manager Example

```

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

```

```

        if( ec == MPI_ERR_PROC_FAILED_PENDING )
            continue;
    }
}
/* get ready to receive more notifications from workers */
MPI_Irecv(buffer, 1, MPI_INT, MPI_ANY_SOURCE, tag, comm, &req);
}
/* ... cancel request and cleanup ... */
}

```

#### 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;
        }
    }
}

```

```
1      gnorm = epsilon + 1.0; /* Force one more iteration */
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
```



# Index

- MPI\_ANY\_SOURCE, [3](#), [12](#), [20](#)
- MPI\_Comm, [8–13](#)
- MPI\_COMM\_NULL, [4](#), [6](#)
- MPI\_COMM\_WORLD, [1](#), [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, [2–4](#), [6–10](#), [13–17](#)
- MPI\_ERR\_SPAWN, [6](#)
- MPI\_ERRORS\_ARE\_FATAL, [2](#)
- MPI\_File, [16](#)
- MPI\_FILE\_NULL, [8](#)
- MPI\_FT, [1](#)
- MPI\_Group, [11](#), [15](#)
- MPI\_GROUP\_EMPTY, [11](#), [16](#)
- MPI\_Request, [13](#)
- MPI\_SEEK\_SET, [8](#)
- MPI\_SUCCESS, [4](#), [6–8](#), [13](#)
- MPI\_Win, [14](#), [15](#)
- MPI\_WIN\_NULL, [7](#)
  
- "mpi\_error\_range", [2](#), [3](#)
- "mpi\_error\_uniform", [4](#), [5](#)
- "coll", [5](#)
- "construct", [5](#)
- "group", [3](#)
- "local", [4](#), [5](#)
- "operation", [3](#)
- "universe", [3](#)
  
- MPI\_ALLREDUCE, [21](#)
- MPI\_BCAST, [4](#)
- MPI\_COMM\_ACCEPT, [5](#)
- MPI\_COMM\_ACK\_FAILED, [12](#), [13](#)
- MPI\_COMM\_ACK\_FAILED(comm, num\_to\_ack, num\_acked), [11](#)
- MPI\_COMM\_AGREE, [4](#), [8](#), [9](#), [12–14](#), [18](#), [21](#)
- MPI\_COMM\_AGREE(comm, flag), [12](#)
- MPI\_COMM\_CONNECT, [5](#)
- MPI\_COMM\_DISCONNECT, [6](#)
- MPI\_COMM\_DUP, [4](#), [5](#)
- MPI\_COMM\_FREE, [4](#)
- MPI\_COMM\_GET\_FAILED, [11–13](#), [18](#)
- MPI\_COMM\_GET\_FAILED(comm, failedgrp), [11](#)
- MPI\_COMM\_GET\_PARENT, [5](#)
- MPI\_COMM\_IAGREE, [9](#)
- MPI\_COMM\_IAGREE(comm, flag, request), [14](#)
- MPI\_COMM\_IDUP, [11](#)
- MPI\_COMM\_IS\_REVOKED, [9](#)
- MPI\_COMM\_IS\_REVOKED(comm, flag), [9](#)
- MPI\_COMM\_ISHRINK, [10](#), [11](#)
- MPI\_COMM\_ISHRINK(comm, newcomm, request), [10](#)
- MPI\_COMM\_JOIN, [6](#)
- MPI\_COMM\_REVOKE, [2](#), [3](#), [21](#)
- MPI\_COMM\_REVOKE(comm), [8](#)
- MPI\_COMM\_SHRINK, [9–12](#), [18](#)
- MPI\_COMM\_SHRINK(comm, newcomm), [9](#)
- MPI\_COMM\_SPAWN, [5](#)
- MPI\_COMM\_SPAWN\_MULTIPLE, [5](#)
- MPI\_COMM\_SPLIT, [4](#), [10](#), [18](#)
- MPI\_FILE\_CLOSE, [8](#)
- MPI\_FILE\_IS\_REVOKED, [17](#)
- MPI\_FILE\_IS\_REVOKED(fh, flag), [17](#)
- MPI\_FILE\_REVOKE, [16](#)
- MPI\_FILE\_REVOKE(fh), [16](#)
- MPI\_FILE\_SEEK, [8](#)
- MPI\_GET, [7](#)
- MPI\_INIT, [5](#), [6](#)
- MPI\_PUT, [7](#)
- MPI\_WIN\_FLUSH, [6](#)
- MPI\_WIN\_FLUSH\_LOCAL, [7](#)
- MPI\_WIN\_FLUSH\_LOCAL\_ALL, [7](#)
- MPI\_WIN\_FREE, [7](#)

1 MPI\_WIN\_GET\_FAILED, [18](#)  
2 MPI\_WIN\_GET\_FAILED(win, failedgrp),  
3 [15](#)  
4 MPI\_WIN\_IS\_REVOKED, [15](#)  
5 MPI\_WIN\_IS\_REVOKED(win, flag), [15](#)  
6 MPI\_WIN\_REVOKE, [7](#), [14](#), [16](#)  
7 MPI\_WIN\_REVOKE(win), [14](#)  
8  
9 NEWEXAMPLES:16.1::C::Fault Tolerant  
10 Comm Split::MPI\_Comm\_split,MPI\_Comm\_agree,MPI\_Comm\_free,  
11 [18](#)  
12 NEWEXAMPLES:16.2::C::Fault Tolerant  
13 Consistent Failed Group (Shrink  
14 variant)::MPI\_Comm\_shrink,MPI\_Group\_difference,MPI\_Comm\_free,  
15 [18](#)  
16 NEWEXAMPLES:16.3::C::Fault Tolerant  
17 Consistent Failed Group (Agree  
18 variant)::MPI\_Comm\_agree,MPI\_Comm\_get\_failed,MPI\_Comm\_ack\_failed,MPI\_Group\_range\_incl,  
19 [19](#)  
20 NEWEXAMPLES:16.4::C::Fault Tolerant  
21 Manager/Worker::MPI\_Comm\_get\_failed,MPI\_Comm\_ack\_failed,  
22 [20](#)  
23 NEWEXAMPLES:16.5::C::Fault Tolerant  
24 Iterative Refinement::MPI\_Comm\_revoke,MPI\_Comm\_agree,MPI\_Comm\_shrink,MPI\_Comm\_free,  
25 [21](#)  
26  
27 TERM:error handling  
28 fault tolerance, [1](#)  
29 ack, [12](#)  
30 agree, [13](#), [14](#)  
31 communicator, [3](#), [8](#)  
32 dynamic process, [5](#)  
33 fault tolerance error, [1](#), [17](#)  
34 I/O, [7](#), [16](#)  
35 inquiry, [1](#)  
36 mitigation, [8](#)  
37 notification, [2](#), [17](#)  
38 one-sided, [6](#), [14](#)  
39 revoke, [14–17](#)  
40 revoked, [2](#), [2](#), [3](#), [9](#)  
41 shrink, [10](#), [10](#)  
42 process failure, [1](#)  
43  
44  
45  
46  
47  
48