

# MPI Global-Restart Fault Tolerance Specification

Version 0.2.1

Unofficial, for comment only

Ignacio Laguna and Giorgis Georgakoudis  
ilaguna@llnl.gov, georgakoudis1@llnl.gov

Lawrence Livermore National Laboratory

December 7, 2021



# Chapter 1

## Global-Restart Fault Tolerance

### 1.1 Introduction

The traditional method to handle process failures in large-scale scientific applications is periodic, global synchronous checkpoint/restart (CPR). When a process failure occurs in a bulk synchronous MPI program, it quickly propagates to other processes so re-starting the application from a previously-saved checkpoint is a simple solution to recover from failures.

A large number of MPI applications already use some form of global synchronous CPR. The goal of global-restart fault tolerance is to provide an easy-to-use interface to improve the efficiency of CPR in bulk synchronous applications by reducing as much as possible the recovery time when failure occurs.

In this chapter, we refer to the global-restart fault tolerance model and interface as the **Reinit** (i.e., re-initialization) model and interface, respectively.

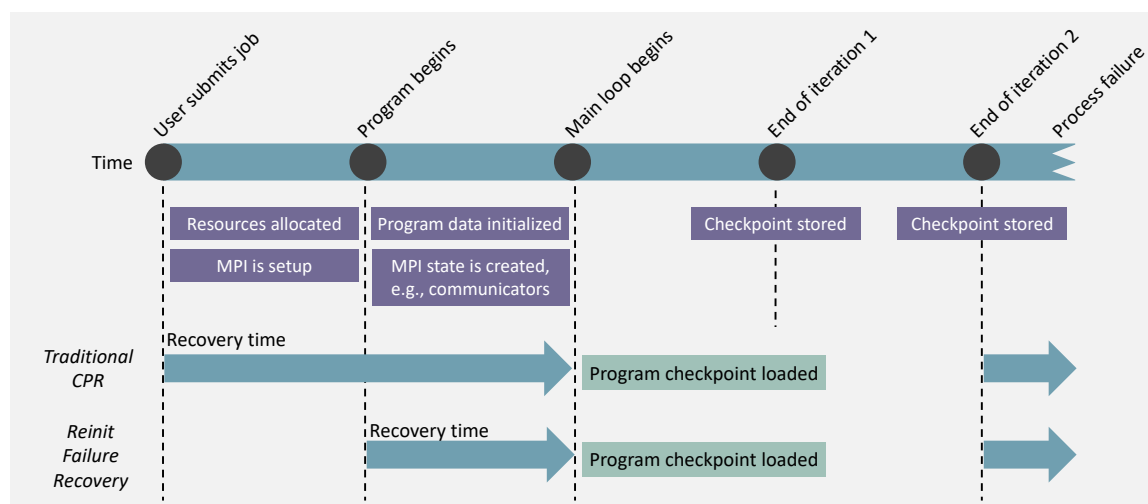


Figure 1.1: The global-restart fault tolerance model (Reinit) provides a mechanisms to reduce the recovery time for bulk synchronous applications that use periodic synchronous checkpoint/restart.

## 1.2 Fault Model

The Reinit model provides a pre-defined fault-tolerance mechanism to survive **MPI process failures**. We use the definition of process failures used in Section 2.8, i.e., a process failure occurs when an MPI process unexpectedly and permanently stops communicating (e.g., a software or hardware crash results in an MPI process terminating unexpectedly). In the rest of the chapter, when we refer to *failures* we mean *MPI process failures*.

The Reinit model assumes that the application's data will be recovered after a failure. The application can use different mechanisms to recover its data, for example, reloading a checkpoint that was saved before the failure occurred or re-generating the data.

## 1.3 Reinit MPI Interface

The Reinit interface for global-restart fault tolerance is composed of two MPI functions: `MPI_REINIT` and `MPI_TEST_FAILURE`. This section describes the syntax of these MPI functions.

### MPI\_Reinit

```
int MPI_Reinit(resilient_fn, void *data)
```

IN `resilient_fn` user defined procedure (function)

IN `data` pointer to user defined data

The user-defined procedure should be in C, a function of type `MPI_Reinit_function` which is defined as:

```
typedef MPI_Reinit_fn void (*)(void *data);
```

The first argument is a user defined procedure, `resilient_fn`, which is called by the `MPI_Reinit` procedure. The second argument is a pointer to user defined data. This pointer is passed as an argument to the user defined procedure, `resilient_fn`, when the procedure is called. A valid MPI program must contain at most one call to the `MPI_Reinit` procedure. Calling `MPI_Reinit` more than one time results in undefined behavior.

The purpose of the user defined `resilient_fn` procedure is to specify a *rollback location*, i.e., a program location to resume execution after a process failure occurs. Depending on the error handler being used, upon the detection of a process failure, MPI will cause the execution of the program to resume at the `resilient_fn` procedure synchronously or asynchronously (see the Error Handling section for more details).

After the `resilient_fn` procedure is re-executed due to failure recovery, the only valid communication objects are the communicators `MPI_COMM_WORLD`, `MPI_COMM_SELF`, `MPI_COMM_NULL`.

*Advice to users.* MPI objects that are created before `MPI_Reinit` is called will not be valid when the `resilient_fn` procedure is re-executed due to a failure. (*End of advice to users.*)

Calling the `MPI_Reinit` procedure sets the `resilient_fn` procedure to be a rollback location and makes this rollback location active. After activating the rollback location, `MPI_Reinit` calls the `resilient_fn` procedure. After the `MPI_Reinit` procedure returns, the rollback location becomes inactive. If a failure occurs during an inactive rollback location,

MPI cannot resume execution at the rollback location, and as a result cannot recover from failures using the Reinit model.

*Advice to users.* To able to survive most of the process failures that can occur during the execution of the program, most calls to MPI and computation should be executed before MPI\_Reinit returns. (*End of advice to users.*)

An MPI process must invoke MPI\_FINALIZE only after MPI\_Reinit returns.

#### MPI\_Test\_failure

```
int MPI_Test_failure()
```

The MPI\_Test\_failure procedure causes the program to resume execution at the rollback point that was activated by MPI\_Reinit when two conditions occur: (1) the MPI\_ERRORS\_REINIT\_SYNC handler is associated with MPI\_COMM\_WORLD, and (2) a failure has been detected before MPI\_Test\_failure is called.

If no failures were detected before MPI\_Test\_failure is called, the return code value is MPI\_SUCCESS and the procedure performs no operations. If on the other hand failures are detected before the procedure is called, the procedure does not return and it immediately resumes execution at the rollback point.

## 1.4 Error Handling

MPI provides two predefined error handlers that can be used to handle failures using the Reinit model. These error handlers are intended to be used to handle failures when the World Model is used to initialize MPI. The Reinit error handlers have no effect when the Sessions Model is used.

Unlike other predefined error handlers, such as MPI\_ERRORS\_ARE\_FATAL, that can be associated to communicator, window, file, and session objects, the Reinit error handlers must be associated only to the predefined MPI\_COMM\_WORLD communicator in the World Model. Associating the Reinit error handlers to window, file, session objects, or communicators other than MPI\_COMM\_WORLD is undefined.

*Rationale.* Associating the Reinit error handler to MPI\_COMM\_SELF would have no effect if a failure occurs because the process that contains MPI\_COMM\_SELF failed and the error handler cannot be called. Since a process failure during the handling of MPI objects, such as windows, files and sessions eventually manifest itself as a process failure in MPI\_COMM\_WORLD, it makes sense to associate a Reinit error handler to MPI\_COMM\_WORLD only. (*End of rationale.*)

The following Reinit error handlers are available in MPI:

- **MPI\_ERRORS\_REINIT\_ASYNC:** The handler is called by MPI immediately after a process failure is detected. The handler, when called, causes the execution of the program to resume at (or jump back to) the active rollback location that was activated by MPI\_Reinit.

- **MPI\_ERRORS\_REINIT\_SYNC**: The handler has two effects. The first effect is that it enables the `MPI_Test_failure` function to cause the execution of the program to resume at (or jump back to) the active rollback location when `MPI_Test_failure` is called. The second effect is that it returns the error code to the user.

Using the `MPI_ERRORS_REINIT_ASYNC` handler causes MPI to resume execution of the program when an error is detected whether or not the error is detected during a call to MPI. On the other hand, using the `MPI_ERRORS_REINIT_SYNC` handler causes MPI to resume execution only after `MPI_Test_failure` function is called if an error was detected.

### 1.4.1 Association of Error Handlers

The Reinit error handlers must be associated to `MPI_COMM_WORLD` before the `MPI-Reinit` procedure is called. Calling `MPI_Reinit` before associating any of the Reinit error handlers produces undefined behavior.

After a Reinit error handler has been associated to `MPI_COMM_WORLD`, it is invalid to associate a different Reinit error handler to `MPI_COMM_WORLD`.

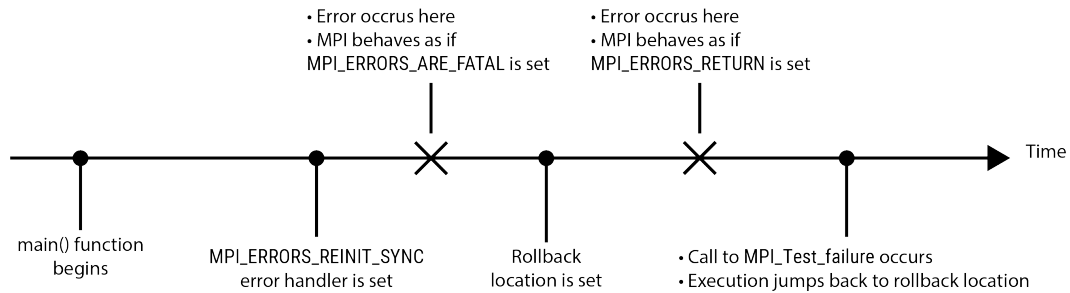


Figure 1.2: Different error scenarios for the `MPI_ERRORS_REINIT_SYNC` error handler.

### 1.4.2 Behavior for Specific Error Conditions

If an error occurs and one of the Reinit error handlers has been set but there is no active Reinit rollback location, MPI will behave as if the `MPI_ERRORS_RETURN` error handler is set (see Figure 1.2).

Errors can occur between the moment the `MPI_ERRORS_REINIT_SYNC` handler is set and the `MPI_Test_failure` function is called. If an error occurs in such period of time, MPI behaves as if the `MPI_ERRORS_RETURN` handler is set.

## 1.5 Tools

The Reinit interface supports the use of MPI tools. The following must be taken into consideration when writing MPI tools:

- The Reinit interface assumes that, when a process failure occurs, data may be lost. If a tool requires data that can be lost due to failures, the tool must implement a mechanism to recover such data, for example, reloading a checkpoint.

- An MPI implementation should provide a performance variable of type `MPI_T_PVAR_CLASS_COUNTER` that reflects the number of times the MPI process has been reinitialized due to failures. The variable has a value of zero initially and it is incremented every time the program resumes execution at the rollback location.
- The performance variables that are provided by an MPI implementation are not reset when execution resumes at the rollback location. Tools are responsible for presenting information about performance variables to users after taking into account failures.

## 1.6 Failures During Device Code Execution

MPI applications may execute code in hardware devices, such as GPUs, which can suffer from failures. In general, it may not be possible to stop the execution of device code. When MPI causes the program execution to resume at a rollback location when a device code region is being executed, this device code region may not be terminated automatically. The `MPIERRORS_REINIT_ASYNC` handler along with the `MPI_Test_failure` function can be used to enable the program execution to be resumed only when device code is not being executed.

## 1.7 Examples

We present a few examples of how to use the Reinit interface with synchronous and asynchronous error handlers.

**Example 1.1** Using Reinit with asynchronous error handling to recover from process failures

```
typedef struct {
    int argc;
    char **argv;
} data_t;

void resilient_function(void *arg)
{
    data_t *data = (data_t *)arg;
    // Cleanup library, if needed
    cleanup_library_state();
    // Resume computation from checkpoint
    // or initialize application data
    if( load_checkpoint() )
        printf("Resume from checkpoint\n");
    else
        init_app_data(data->argc, data->argv);
    bool done = false;
    while(!done) {
        done = compute();
        store_checkpoint();
    }
}
```

```

1      }
2  }
3
4  int main(int argc, char *argv[])
5  {
6      // Initialize user defined data type
7      data_t data = { argc, argv };
8
9      MPI_Init(argc, argv);
10     MPI_Comm_set_errhandler(MPI_COMM_WORLD, MPI_ERRORS_REINIT_ASYNC);
11     // MPI_Reinit sets the rollback location
12     // to resilient_function and calls it.
13     // In asynchronous error handling, the program
14     // will go to the rollback location as soon a
15     // failure is detected
16     MPI_Reinit(&data, resilient_function);
17     MPI_Finalize();
18
19     return 0;
20 }

```

**Example 1.2** Using Reinit with synchronous error handling to recover from process failures

```

26 void resilient_function(void *arg)
27 {
28     data_t *data = (data_t *)arg;
29     // Cleanup library, if needed
30     cleanup_library_state();
31     // Resume computation from checkpoint
32     // or initialize application data
33     if( load_checkpoint() )
34         printf("Resume from checkpoint\n");
35     else
36         init_app_data(data->argc, data->argv);
37     bool done = false;
38     while(!done) {
39         done = compute();
40         MPI_Test_failure();
41         store_checkpoint();
42         // Calling MPI_Test_failure will go to the
43         // rollback location, that is resilient_function,
44         // in case of a failure
45         // MPI + computation
46         compute();
47         MPI_Test_failure();
48         // MPI + computation

```



```
        compute();  
        MPI_Test_failure();  
    }  
}
```

## 1.8 To-Do List

While the Reinit specification is almost complete, a few aspects may require further clarification, which include the following:

1. Do we need to define FORTRAN bindings?
2. Can we allow multiple and consecutive Reinit points (blocks)?

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