

D R A F T

Document for a Standard Message-Passing Interface

Message Passing Interface Forum

November 8, 2022

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

Chapter 16

Completion Continuations

Some applications may need to handle large numbers of requests or require fast reaction to the completion or cancellation of an MPI operation. The reaction to the completion of an operation can be expressed as a **continuation**. A continuation is a callback function provided by the application that is invoked by MPI once completion of the operation is detected.

Continuations are *attached* to either a single operation request or a set of operation requests and *registered* with a continuation request. A continuation request is a persistent request that has to be initialized and freed by the application and can be used to test or wait for its completion. A continuation request completes once the callbacks of all continuations previously registered with it have completed execution. After initialization or completion, a continuation request has to be started to enable the execution of continuations. However, continuations can be registered with a valid continuation request at any time.

Continuation requests themselves may have a continuation attached, which will be invoked once the continuation request is complete. Attaching a continuation to a non-persistent request returns ownership of that request to MPI, i.e., the request may subsequently not be used to test or wait for the completion of the respective operation. The ownership of persistent requests is returned to the application at the start of the execution of the continuation callback. The outcome of attaching more than one continuation to a request is undefined.

Execution of continuation callbacks can occur on any application thread calling into MPI or be restricted to any thread testing or waiting on the associated continuation request (see Section 16.3.3). When attaching a continuation to an operation a status object may be provided, which will be filled before the continuation is invoked. The application may pass `MPI_STATUS[ES]_IGNORE` in lieu of a status object or status objects.

MPI procedures may be called from within the continuation callback. Applications strive for short callback functions so as not to unnecessarily prolong the execution time of the MPI procedure from within which the callback is executed. Specifically, while not explicitly prohibited the use of blocking MPI procedures inside the continuation callback is discouraged.

Example 16.1 uses continuations to implement a scheme for offloading work to other processes. By using continuations, the requests needed to track the corresponding send and receive operations do not have to be managed in application space. Instead, progressing the continuation request is sufficient to react to the completion of both operations.

Example 16.1. Library functions to offload work to other processes and receive the result. A continuation is attached to the send and receive operation, which will be executed once both operations are completed.

```
#include <mpi.h>
```

```

1  /* work descriptor */
2  struct work_item {
3      MPI_Request reqs[2];
4      MPI_Status stats[2];
5      work_t *msg;
6      int size;
7      int reply;
8  };
9
10 MPI_Request cont_request;
11 MPI_Comm comm;
12
13 void init()
14 {
15     MPI_Continue_init(0, 0, MPI_INFO_NULL, &cont_request);
16     MPI_Start(&cont_request);
17     MPI_Comm_dup(MPI_COMM_WORLD, &comm);
18 }
19
20 void fini()
21 {
22     MPI_Request_free(&cont_request);
23     MPI_Comm_free(&comm);
24 }
25
26 /* callback invoked by MPI when send and receive operations are complete */
27 int complete_cb(int rc, void *user_data)
28 {
29     struct work_item *wd = (struct work_item *)user_data;
30     int source = wd->stats[1].MPI_SOURCE;
31     mark_completed(wd->msg, wd->reply, source);
32     free(wd);
33     return MPI_SUCCESS;
34 }
35
36 /* send work to a target process and post a receive for the result */
37 void send_work(work_t *msg, int size)
38 {
39     struct work_item *wd = malloc(sizeof(struct work_item));
40     wd->msg = msg;
41     wd->size = size;
42     int target = next_target();
43     /* send the message */
44     MPI_Isend(msg, size, MPI_BYTE, target, comm, &wd->reqs[0]);
45     /* receive the reply */
46     MPI_Irecv(&wd->reply, 1, MPI_INT, target, comm, &wd->reqs[1]);
47     /* attach a continuation to both requests */
48     MPI_Continueall(2, wd->reqs,
49                    &complete_cb, wd,
50                    /* flags = */0, wd->stats, &cont_request);
51 }
52
53 /* progress outstanding communication and continuations */

```

```

void progress()
{
    int flag;
    MPI_Test(&flag, &cont_request, MPI_STATUS_IGNORE);
    if (flag) {
        MPI_Start(&cont_request);
    }
}

```

16.1 Continuation Requests

`MPI_CONTINUE_INIT(flags, max_poll, info, cont_req)`

IN	flags	flags (integer)
IN	max_poll	maximum number of continuations to execute, or 0 for no limit (non-negative integer)
IN	info	info argument (handle)
OUT	cont_req	continuation request (handle)

C binding

```

int MPI_Continue_init(int flags, int max_poll, MPI_Info info,
                    MPI_Request *cont_req)

```

Fortran 2008 binding

```

MPI_Continue_init(flags, max_poll, info, cont_req, ierror)
    INTEGER, INTENT(IN) :: flags, max_poll
    TYPE(MPI_Info), INTENT(IN) :: info
    TYPE(MPI_Request), INTENT(OUT) :: cont_req
    INTEGER, OPTIONAL, INTENT(OUT) :: ierror

```

Fortran binding

```

MPI_CONTINUE_INIT(FLAGS, MAX_POLL, INFO, CONT_REQ, IERROR)
    INTEGER FLAGS, MAX_POLL, INFO, CONT_REQ, IERROR

```

A call to this procedure creates a new continuation request in `cont_req`. The `flags` argument is used to control aspects of the continuation request, with the following predefined flag:

MPI_CONT_POLL_ONLY Marks the continuation request as poll-only, i.e., continuations registered with this continuation request are only executed when a thread tests or wait for the completion of the continuation request.

Additional flags may be added in the future.

The `max_poll` argument controls the maximum number of continuations to execute when testing for the completion of this continuation request. If multiple continuation requests are tested for completion, the maximum number of continuations executed is the sum of all of their limits. A value of 0 signals that all available continuations may be executed when testing or waiting for the completion of the continuation request.

The `info` argument is used to further control the aspects of the continuation request. Predefined info keys are described in Section 16.1.1.

A call to a test or wait procedure indicating completion of a continuation request does not free the request. A continuation request is released through a call to `MPI_REQUEST_FREE`, which marks the request for deallocation. The request will be deallocated once all registered continuations have been executed. Continuation requests may not be canceled.

16.1.1 Predefined Info Keys

The execution context of continuations registered with a continuation request can be controlled using the following info keys:

"mpi_continue_thread" This key may be set to one of the following two values: "application" and "any". The "application" value indicates that continuations may only be executed by threads controlled by the application, i.e., any application thread that calls into MPI. This is the default. The value "any" indicates that continuations may be executed by *any* thread, including MPI-internal progress threads if available. This key has no effect on implementations that do not use an internal progress thread.

Rationale. Some applications may rely on thread-local data being initialized outside of the continuation or use callbacks that are not thread-safe, in which case the use of "any" would lead to correctness issues. (*End of rationale.*)

"mpi_continue_async_signal_safe" If the value is set to "true", the application provides a hint to the implementation that the continuations are async-signal safe and thus may be invoked from within a signal handler. This limits the capabilities of the callback, excluding calls back into the MPI library and other unsafe operations. The default is "false".

16.2 Callback Function Signature

```
typedef int MPI_Continue_cb_function(int error_code, void *user_data);
```

ABSTRACT INTERFACE

```
  SUBROUTINE MPI_Continue_cb_function(error_code, user_data, ierror)
```

```
    INTEGER :: error_code
```

```
    INTEGER(KIND=MPI_ADDRESS_KIND) :: user_data
```

```
    INTEGER, OPTIONAL :: ierror
```

```
  SUBROUTINE MPI_CONTINUE_CB_FUNCTION(ERROR_CODE, USER_DATA, IERROR)
```

```
    INTEGER ERROR_CODE, IERROR
```

```
    INTEGER(KIND=MPI_ADDRESS_KIND) USER_DATA
```

The continuation callback function has the type `MPI_Continue_cb_function` and upon invocation is passed an error code signalling the state of the associated operations as well as the pointer to additional data provided when attaching a continuation to the operation or set of operations. Unless `MPI_CONT_INVOKE_FAILED` is specified when attaching a continuation (see Section 16.3.3), the error code passed into the continuation will always be `MPI_SUCCESS`. The continuation callback should return `MPI_SUCCESS` if the continuation

executed successfully or an error code otherwise. In the latter case, the continuation will be marked as failed. The handling of failed continuations is explained in Section 16.4.

16.3 Attaching Continuations

16.3.1 Attaching to a Single Request

`MPI_CONTINUE(op_request, cb, cb_data, flags, status, cont_request)`

INOUT	<code>op_request</code>	operation request (handle)
IN	<code>cb</code>	callback to be invoked once the operation is complete (function)
IN	<code>cb_data</code>	pointer to a user-controlled buffer
IN	<code>flags</code>	flags controlling aspects of the continuation (integer)
IN	<code>status</code>	status object (array of status)
IN	<code>cont_request</code>	continuation request (handle)

C binding

```
int MPI_Continue(MPI_Request *op_request, MPI_Continue_cb_function cb,
                void *cb_data, int flags, MPI_Status *status,
                MPI_Request cont_request)
```

Fortran 2008 binding

```
MPI_Continue(op_request, cb, cb_data, flags, status, cont_request, ierror)
  TYPE(MPI_Request), INTENT(INOUT) :: op_request
  MPI_Continue_cb_function, INTENT(IN) :: cb
  INTEGER(KIND=MPI_ADDRESS_KIND), INTENT(IN) :: cb_data
  INTEGER, INTENT(IN) :: flags
  TYPE(MPI_Status), INTENT(IN), ASYNCHRONOUS :: status(1)
  TYPE(MPI_Request), INTENT(IN) :: cont_request
  INTEGER, OPTIONAL, INTENT(OUT) :: ierror
```

Fortran binding

```
MPI_CONTINUE(OP_REQUEST, CB, CB_DATA, FLAGS, STATUS, CONT_REQUEST, IERROR)
  INTEGER OP_REQUEST, FLAGS, STATUS(MPI_STATUS_SIZE, 1), CONT_REQUEST, IERROR
  MPI_Continue_cb_function CB
  INTEGER(KIND=MPI_ADDRESS_KIND) CB_DATA
```

This function attaches a continuation to the operation represented by the request `op_request` and registers it with the continuation request `cont_request`. The callback function `cb` will be invoked after the MPI implementation finds the operation to be complete. Upon invocation, `cb_data` pointer will be passed to the callback function. The request `cont_request` must be a continuation request created through a call to `MPI_CONTINUE_INIT`.

The `flags` argument is either 0 or an OR-combination of one or more of the flags outlined in Section 16.3.3.

1 If the operation represented by `op_request` is complete at the time of the call to
 2 `MPI_CONTINUE`, the implementation may invoke the callback function `cb` immediately,
 3 unless the `MPI_CONT_DEFER_COMPLETE` flag is provided.

4 Unless `MPI_STATUS_IGNORE` is passed as `status`, the status object pointed to by `status`
 5 must be accessible until the continuation callback is invoked.

6 The memory location holding the request must be accessible until the continuation
 7 callback is invoked and the request will be set to `MPI_REQUEST_NULL` before the invocation.
 8 If the `MPI_CONT_REQBUF_VOLATILE` flag was provided the request buffer may be reused
 9 after the call to `MPI_CONTINUE` returns.

10 A persistent operation request will not be set to `MPI_REQUEST_NULL`. Upon invocation
 11 of the attached continuation, the persistent request will be inactive and may be started
 12 inside the continuation callback. The outcome of attaching more than one continuation
 13 to an operation is undefined. After a continuation has been attached, persistent operation
 14 requests may not be tested or waited on until the execution of the continuation callback has
 15 begun. A persistent operation request with an attached continuation may be canceled (if
 16 allowed for the type of operation), which will make the continuation eligible for execution.
 17 Whether or not a request has been canceled can be queried from the associated status
 18 object.

19 A continuation may be attached to a continuation request (i.e., the request `op_request`
 20 itself may be a continuation request), in which case the continuation is invoked once all
 21 continuations registered with the continuation request have completed. It is erroneous to
 22 pass the same continuation request as both the `op_request` and `cont_request` argument. No
 23 new continuations may be registered with a continuation request from the point when a
 24 continuation has been attached to it until execution of the continuation callback has begun.
 25

26 *Rationale.* The above restriction prohibiting new continuations to be registered with
 27 a continuation request that has a continuation attached to it is meant to prevent cyclic
 28 dependencies between continuation requests. Otherwise, a deadlock is imminent if the
 29 completion of one continuation request is dependent on the completion of the other.
 30 (*End of rationale.*)
 31

32 16.3.2 Attaching to Multiple Requests

33
 34
 35 `MPI_CONTINUEALL(count, array_of_op_requests, cb, cb_data, flags, array_of_statuses,`
 36 `cont_request)`

37			
38	IN	<code>count</code>	list length (non-negative integer)
39	INOUT	<code>array_of_op_requests</code>	array of requests (array of handles)
40	IN	<code>cb</code>	callback to be invoked once the operation is complete
41			(function)
42	IN	<code>cb_data</code>	pointer to a user-controlled buffer
43	IN	<code>flags</code>	flags controlling aspects of the continuation (integer)
44	IN	<code>array_of_statuses</code>	array of status objects (array of status)
45	IN	<code>array_of_statuses</code>	array of status objects (array of status)
46	IN	<code>cont_request</code>	continuation request (handle)
47			
48			

C binding

```

int MPI_Continueall(int count, MPI_Request array_of_op_requests[],
                   MPI_Continue_cb_function cb, void *cb_data, int flags,
                   MPI_Status array_of_statuses[], MPI_Request cont_request)

```

Fortran 2008 binding

```

MPI_Continueall(count, array_of_op_requests, cb, cb_data, flags,
               array_of_statuses, cont_request, ierror)
INTEGER, INTENT(IN) :: count, flags
TYPE(MPI_Request), INTENT(INOUT) :: array_of_op_requests(count)
MPI_Continue_cb_function, INTENT(IN) :: cb
INTEGER(KIND=MPI_ADDRESS_KIND), INTENT(IN) :: cb_data
TYPE(MPI_Status), INTENT(IN), ASYNCHRONOUS :: array_of_statuses(*)
TYPE(MPI_Request), INTENT(IN) :: cont_request
INTEGER, OPTIONAL, INTENT(OUT) :: ierror

```

Fortran binding

```

MPI_CONTINUEALL(COUNT, ARRAY_OF_OP_REQUESTS, CB, CB_DATA, FLAGS,
               ARRAY_OF_STATUSES, CONT_REQUEST, IERROR)
INTEGER COUNT, ARRAY_OF_OP_REQUESTS(*), FLAGS,
               ARRAY_OF_STATUSES(MPI_STATUS_SIZE, *), CONT_REQUEST, IERROR
MPI_Continue_cb_function CB
INTEGER(KIND=MPI_ADDRESS_KIND) CB_DATA

```

Similar to `MPI_CONTINUE`, this function is used to attach a continuation callback to a set of operation requests. The continuation callback will be invoked once all `count` operations in the list `array_of_op_requests` have completed. If `MPI_STATUSES_IGNORE` is not passed for `array_of_statuses`, the array should be of length `count` and the statuses will be set before the continuation is invoked. Unless `MPI_STATUSES_IGNORE` is provided, the memory containing the statuses must remain accessible until the continuation is invoked. Unless the `MPI_CONT_REQBUF_VOLATILE` flag is set, the memory containing the requests must remain accessible until the continuation is invoked. The rules regarding persistent and non-persistent requests described for `MPI_CONTINUE` also apply here.

16.3.3 Flags For Attaching Continuations

MPI_CONT_DEFER_COMPLETE Do not execute the continuation immediately even if the operation (or all operations) is complete already while the continuation is attached. The continuation will instead be executed at a later point in time.

MPI_CONT_REQBUF_VOLATILE The request buffer is volatile and should not be accessed after the return from the function call. If used, MPI cannot handle errors occurring in any of the involved operations. The behavior is undefined if an error occurs on an associated operation and a non-aborting error handler is installed.

Rationale. In case MPI discovers a fault on an operation it may be required to access to the request objects to properly handle the error and mark the continuation associated with that operation as failed. Thus, this flag should not be used in conjunction with error handlers that do not abort. However, applications that do not wish to handle errors may use this flag to simplify the handling of requests. (*End of rationale.*)

MPI_CONT_PERSISTENT The continuation is marked as persistent. This flag is only allowed in conjunction with persistent requests. If provided, the continuation will remain attached to the persistent operations and be invoked once all requests have completed. The persistent requests have to be freed to remove the persistent continuation. See Section 16.3.4 for a discussion of persistent continuations.

MPI_CONT_INVOKE_FAILED The continuation is invoked even if an error is detected for one or more of the associated operations. In that case, the error code for the operation (in the case of `MPI_CONTINUE`) or `MPI_ERR_IN_STATUS` is passed as the error code to the continuation callback. If the continuation subsequently returns `MPI_SUCCESS` the continuation will not be marked as failed.

16.3.4 Persistent Continuations

By default, continuations that are attached to persistent requests are not themselves persistent, i.e., a continuation attached to a persistent request is removed from that request once the callback has executed. The next time the request completes there will be no continuation to execute, unless a new continuation has been attached.

By specifying the flag `MPI_CONT_PERSISTENT`, the continuation can be marked as persistent. In that case, the continuation will remain attached to the persistent request after it completed and has been started again. The persistent request has to be freed to remove the continuation. In the case of `MPI_CONTINUEALL`, all persistent requests have to be started and subsequently complete before the continuation can be invoked again by MPI. The outcome is undefined if only a subset of the requests associated with a continuation are started again. This may lead to a deadlock since the continuation callback is never invoked again and the continuation request to which the continuation is registered remains active.

It is erroneous to use this flag with requests that are not persistent.

16.4 Error Handling

A continuation may fail for two reasons. First, any of the operations the continuation is attached to fails. This will cause the continuation to be marked as failed and the corresponding error and MPI object of the operation will be associated with the continuation. The status for the request of the failed continuation will contain the error code for that operation. By default, the continuation will not be executed if one or more of its operations have failed.

However, if `MPI_CONT_INVOKE_FAILED` has been specified then the continuation callback will be invoked and the error code of the first failed operation will be passed as its first argument. If the application is able to handle that error inside the continuation callback, it may subsequently return `MPI_SUCCESS` from the continuation callback to prevent the error from being propagated outside of the continuation. If, however, a continuation returns any error other than `MPI_SUCCESS` the continuation will be marked as failed and the returned error as well as `MPI_COMM_SELF` will be associated with the continuation.

The test or wait on a continuation request with a registered failed continuation shall return for that request the error of the first continuation detected as failed. The appropriate error handler will be invoked on the error and MPI object associated with that continuation.

If that error handler does not abort the application, the set of failed continuations can be queried from the continuation request using `MPI_CONTINUE_GET_FAILED`. A continu-

ation request that contains failed continuations cannot be restarted until all failed continuations have been queried.

`MPI_CONTINUE_GET_FAILED(cont_request, count, cb_data)`

IN	<code>cont_request</code>	continuation request (handle)
INOUT	<code>count</code>	list length (non-negative integer)
OUT	<code>cb_data</code>	address of a buffer of <code>count</code> pointer to callback data (choice)

C binding

```
int MPI_Continue_get_failed(MPI_Request cont_request, int *count,
                           void *cb_data)
```

Fortran 2008 binding

```
MPI_Continue_get_failed(cont_request, count, cb_data, ierror)
  TYPE(MPI_Request), INTENT(IN) :: cont_request
  INTEGER, INTENT(INOUT) :: count
  TYPE(*), DIMENSION(..) :: cb_data
  INTEGER, OPTIONAL, INTENT(OUT) :: ierror
```

Fortran binding

```
MPI_CONTINUE_GET_FAILED(CONT_REQUEST, COUNT, CB_DATA, IERROR)
  INTEGER CONT_REQUEST, COUNT, IERROR
  <type> CB_DATA(*)
```

This function returns at most `count` pointers to callback data of failed continuations in the buffer pointed to by `cb_data`. The argument `cb_data` should be an array pointers to callback data of length at least `count`.

Rationale. The use of a formal parameter `cb_data` of type `void*` (rather than `void**`) avoids the messy type casting that would be needed if the callback data pointer are declared with a type other than `void*`. (*End of rationale.*)

Upon return, `count` will be set to the actual number of callback data pointers stored in the first `count` positions of the `cb_data` array. If the value of `count` upon return is the same as its input value then there may be more failed continuations to query. Conversely, if `count` upon return is smaller than the input value then all failed continuations have been queried. The callback data pointer for any given continuation shall not be returned twice, unless $N > 1$ continuations share the same pointer, in which case that pointer value shall be returned N times.

Querying failed continuations does not change the state of the continuation request, i.e., the continuation request has to be started again to enable the execution of continuations and new continuations may be registered with a continuation request that has failed continuations left to query.

16.5 Continuations and Threads

If the implementation supports `MPI_THREAD_MULTIPLE`, it is safe to register continuations with the same continuation request from within multiple threads. This allows multiple threads within an application to attach continuations to operations without requiring mutual exclusion. However, starting as well as testing and waiting for its completion must be limited to a single thread at a time.

If multiple threads exist within an application, any thread calling into MPI may execute continuation callbacks registered with any continuation request. In order to limit the execution of continuations callbacks to a single thread at a time, the `MPI_CONT_POLL_ONLY` flag may be passed to `MPI_CONTINUE_INIT`. Consequently, the callbacks of continuations registered with this continuation request will only be executed by a thread testing or waiting for the completion of the continuation request.

16.6 Examples

Example 16.2. Using a persistent continuation on a persistent receive request, restarting the request after processing an incoming message. This examples uses `MPI_WAIT` to wait for the completion of all continuations registered with the continuation request and thus to wait for all messages to be processed.

```

22 #include <stdlib.h>
23 #include <mpi.h>
24
25 #define NUM_VARS 1024
26 #define TAG 1001
27
28 static MPI_Request recv_request, cont_request;
29 static volatile int num_recvs = 0;
30 static int world_size;
31
32 int completion_cb(int rc, void *user_data)
33 {
34     process_vars(user_data);
35     num_recvs += 1;
36     if (num_recvs < world_size-1) {
37         MPI_Start(&recv_request);
38         /* the continuation was marked persistent */
39     }
40     return MPI_SUCCESS;
41 }
42
43 int main(int argc, char *argv[])
44 {
45     int rank;
46     MPI_Status status;
47     MPI_Init(&argc, &argv);
48     MPI_Comm comm = MPI_COMM_WORLD;
49     MPI_Comm_size(comm, &world_size);
50     MPI_Comm_rank(comm, &rank);

```

```

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
double *vars = malloc(sizeof(double)*NUM_VARS);
if (rank == 0) {
    MPI_Continue_init(0, 0, MPI_INFO_NULL, &cont_request);
    MPI_Recv_init(vars, NUM_VARS, MPI_DOUBLE, MPI_ANY_SOURCE, TAG,
                 comm, &recv_request);
    MPI_Start(&recv_request);
    MPI_Continue(&recv_request, &completion_cb, vars,
                MPI_CONT_PERSISTENT,
                &status, cont_request);
    /* wait for all messages to be received */
    MPI_Start(&cont_request);
    MPI_Wait(&cont_request, MPI_STATUS_IGNORE);
    MPI_Request_free(&recv_request);
    MPI_Request_free(&cont_request);
} else {
    create_vars(vars);
    MPI_Send(vars, NUM_VARS, MPI_DOUBLE, 0, TAG, comm);
}

free(vars);
MPI_Finalize();
return 0;
}

```

Example 16.3. Using continuations to react to an arbitrary number of messages (sender not shown) in a library and checking for cancellation of the receive request inside the continuation. The progress function should be called periodically by the library's user.

```

25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
#include <mpi.h>

static MPI_Request cont_request = MPI_REQUEST_NULL;
static int cont_init_count = 0;
struct callback_data {
    MPI_Request recv_request;
    MPI_Status status;
    void *buffer;
};

int completion_cb(int rc, void *user_data)
{
    int cancelled;
    /* test whether the receive was cancelled, process otherwise */
    MPI_Test_cancelled(status, &cancelled);
    if (cancelled) {
        MPI_Request_free(&recv_request);
    } else {
        process_msg(user_data);
        MPI_Start(&recv_request);
        /* the continuation was marked persistent */
    }
    return MPI_SUCCESS;
}

```

```

1  }
2
3  /* initialize a receive with a given tag and register a continuation */
4  void init_recv(void *buffer, int num_bytes, int tag)
5  {
6      /* initialize continuation request and start a receive */
7      struct callback_data *cb_data = malloc(sizeof(*cb_data));
8      if (cont_request == MPI_REQUEST_NULL) {
9          MPI_Continue_init(0, 0, MPI_INFO_NULL, &cont_request);
10         MPI_Start(&cont_request);
11         cont_init_count++;
12     }
13     MPI_Recv_init(vars, num_bytes, MPI_BYTE, MPI_ANY_SOURCE, tag,
14                 comm, &cb_data->recv_request);
15     MPI_Start(&cb_data->recv_request);
16     MPI_Continue(&cb_data->recv_request, &completion_cb, buffer,
17                MPI_CONT_PERSISTENT,
18                &cb_data->status, cont_request);
19 }
20
21 void progress()
22 {
23     int flag;
24     /* progress outstanding continuations */
25     MPI_Test(&cont_request, &flag, MPI_STATUS_IGNORE);
26     if (flag) {
27         MPI_Start(&cont_request);
28     }
29 }
30
31 void end_recv()
32 {
33     /* cancel the request and wait for the last continuation to complete */
34     MPI_Cancel(&recv_request);
35     MPI_Wait(&cont_request, MPI_STATUS_IGNORE);
36     --cont_init_count;
37     if (cont_init_count == 0) {
38         MPI_Request_free(&cont_request);
39     }
40 }

```

Example 16.4. Using continuations to handle detached OpenMP tasks communicating through MPI. An additional background thread is needed to ensure progress on outstanding continuations. For both continuations, the flag `MPI_CONT_REQBUF_VOLATILE` is used because the request variable is located at the stack and will go out of scope once the task completes. The `detach` clause on the receive task marks the task as detached, i.e., although the task completes its dependencies will not be fulfilled until the `omp_fulfill_event` OpenMP procedure is called on the event. This mechanism can be used to encapsulate MPI communication in OpenMP tasks without blocking the thread executing the task.

```
#include <stdlib.h>
```

```

1  #include <unistd.h>
2  #include <pthread.h>
3  #include <omp.h>
4  #include <mpi.h>
5
6  #define NUM_VARS 1024
7  #define TAG 1001
8
9  void send_completion_cb(MPI_Status *status, void *user_data)
10 {
11     free(user_data);
12 }
13
14 void recv_completion_cb(MPI_Status *status, void *user_data)
15 {
16     omp_fulfill_event((omp_event_t) user_data);
17 }
18
19 static volatile int need_progress = 1;
20 void* progress_thread(void *arg)
21 {
22     int flag;
23     MPI_Request *cont_request = (MPI_Request*)arg;
24     while (need_progress) {
25         MPI_Test(cont_request, &flag, MPI_STATUS_IGNORE);
26         if (flag) {
27             MPI_Start(&cont_request);
28         }
29         usleep(100);
30     }
31     return NULL;
32 }
33
34 int main(int argc, char *argv[])
35 {
36     int rank, size, provided;
37     MPI_Request op_request, cont_request;
38     omp_event_t event;
39
40     MPI_Init_thread(&argc, &argv, MPI_THREAD_MULTIPLE, &provided);
41     MPI_Comm comm = MPI_COMM_WORLD;
42     MPI_Comm_size(comm, &size);
43     MPI_Comm_rank(comm, &rank);
44     MPI_Continue_init(MPI_INFO_NULL, &cont_request);
45     MPI_Start(&cont_request);
46     /* thread that progresses outstanding continuations */
47     pthread_t thread;
48     pthread_create(&thread, NULL, &progress_thread, &cont_request);
49
50     #pragma omp parallel master
51     {
52         if (rank == 0) {
53             #pragma omp taskloop

```

```

1   for (int i = 1; i < size; ++i) {
2       double *vars = malloc(sizeof(double)*NUM_VARS);
3       compute_vars_for(vars, i);
4       MPI_Isend(vars, NUM_VARS, MPI_DOUBLE, i, TAG, comm, &op_request);
5       /* attach continuation that frees the buffer once complete */
6       MPI_Continue(&op_request, &send_completion_cb, vars,
7                   MPI_CONT_REQBUF_VOLATILE,
8                   MPI_STATUS_IGNORE, cont_request);
9   }
10  } else {
11  /* task that receives values */
12  double *vars;
13  #pragma omp task depend(out: vars) detach(event)
14  {
15      MPI_Request op_request;
16      vars = malloc(sizeof(double)*NUM_VARS);
17      MPI_Irecv(vars, NUM_VARS, MPI_DOUBLE, 0, TAG, comm, &op_request);
18      MPI_Continue(&op_request, &recv_completion_cb, event, 0,
19                  MPI_CONT_REQBUF_VOLATILE,
20                  MPI_STATUS_IGNORE, cont_request);
21  }
22  /* task processing values, executed once the receiving task's
23  dependencies are released */
24  #pragma omp task depend(in: vars)
25  {
26      compute_vars_from(vars, 0);
27      free(vars);
28  }
29  }
30
31  need_progress = 0;
32  pthread_join(thread, NULL);
33
34  MPI_Request_free(&cont_request);
35  MPI_Finalize();
36  return 0;
37  }

```

Example 16.5. A modified progress function of Example 16.1 querying failed continuations. A failed target is removed from the list of targets and an unsuccessful work item is resubmitted to another process. This requires `MPI_ERRORS_RETURN` to be set as error handler on the communicator `comm`.

```

41  /* progress outstanding communication and continuations */
42  void offload_progress()
43  {
44      int ret, flag;
45      ret = MPI_Test(&flag, &cont_request, MPI_STATUS_IGNORE);
46      if (MPI_SUCCESS != ret) {
47          /* some continuations have failed, query which ones */

```



```
struct work_item *wds[16];
int count = 16;
while (16 == count) {
    MPI_Continue_get_failed(cont_request, &count, wds);
    for (int i = 0; i < count; ++i) {
        /* mark the target as unavailable */
        remove_target(wds[i]->status[1].MPI_SOURCE);
        /* resubmit work to another target */
        send_work(wds[i]->work, wds[i]->size);
        free(wds[i]);
    }
}
/* all failed continuations have been handled,
 * restart the continuation request */
MPI_Start(&cont_request);
} else if (flag) {
    MPI_Start(&cont_request);
}
}
```

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

Index

MPI_COMM_SELF, [8](#)
MPI_CONT_DEFER_COMPLETE, [6](#), [7](#)
MPI_CONT_INVOKE_FAILED, [4](#), [8](#)
MPI_CONT_PERSISTENT, [8](#)
MPI_CONT_POLL_ONLY, [3](#), [10](#)
MPI_CONT_REQBUF_VOLATILE, [6](#), [7](#), [12](#)
MPI_ERR_IN_STATUS, [8](#)
MPI_ERRORS_RETURN, [14](#)
MPI_REQUEST_NULL, [6](#)
MPI_STATUS[ES]_IGNORE, [1](#)
MPI_STATUS_IGNORE, [6](#)
MPI_STATUSES_IGNORE, [7](#)
MPI_SUCCESS, [4](#), [8](#)
MPI_THREAD_MULTIPLE, [10](#)

"mpi_continue_async_signal_safe", [4](#)
"mpi_continue_thread", [4](#)
"any", [4](#)
"application", [4](#)
"false", [4](#)
"true", [4](#)

MPI_CONTINUE, [6-8](#)
MPI_CONTINUE(op_request, cb, cb_data,
 flags, status, cont_request), [5](#)
MPI_Continue_cb_function, [4](#)
MPI_CONTINUE_GET_FAILED, [8](#)
MPI_CONTINUE_GET_FAILED(cont_request,
 count, cb_data), [9](#)
MPI_CONTINUE_INIT, [5](#), [10](#)
MPI_CONTINUE_INIT(flags, max_poll, info,
 cont_req), [3](#)
MPI_CONTINUEALL, [8](#)
MPI_CONTINUEALL(count, array_of_op_requests,
 cb, cb_data, flags, array_of_statuses,
 cont_request), [6](#)
MPI_REQUEST_FREE, [4](#)
MPI_SUCCESS, [8](#)
MPI_WAIT, [10](#)

TERM:continuation, [1](#)

MPI_CONTINUE_CB_FUNCTION, [4](#)
MPI_Continue_cb_function, [4](#)
MPI_Continue_cb_function(
 int error_code, void *user_data), [4](#)