

Composable Asynchronous Communication Graphs and Streams in MPI

Quincey Koziol, Amazon, qkoziol@amazon.com
Ryan Grant, Queens' University, ryan.grant@queensu.ca

The MPI standard has historically employed a “weak” progress model for overlapping communication, I/O, and computation with ‘nonblocking’ API calls. Weak progress means that the MPI library must be entered for nonblocking operations to be guaranteed of making progress toward completion. A “strong” progress model would mean that nonblocking operations would progress asynchronously toward completion, without requiring that the MPI library be entered for that progress to occur.

This proposal describes extensions to MPI that allow applications to request strong progress from an MPI implementation, and to optionally perform ‘true’ asynchronous operations.

Comment [QK1]: Is there a definitive description of weak and strong progress to refer to?

Add ‘strong’ progress to MPI Init, or session init

Revision History

Version Number	Date	Comments
v1-9	Mar-Sept, 2022	Hand-written lab notebook circulated informally amongst collaborators
v10	Sept. 14, 2022	Shared with MPI Collectives/Persistence Working Group
v11	Sept. 30, 2022	Shared at the MPI Forum.
v12-14	Oct. 17, 2022	Shared with HPE and NVIDIA.
v15	Oct. 18, 2022	Weekly discussions w/Ryan Grant.

Contents

List of Figures	5
List of Listings	6
1. Introduction	7
2. Background: MPI Today	7
3. Motivation / Etc.	8
4. Overview	8
4.1. Deferred Operations	8
4.1.1. Data Dependencies For Deferred Operations	8
4.2. Aggregating Deferred Operations: MPI Graphs and Streams	9
4.2.1. Errors When Executing Deferred Operations Asynchronously	10
4.3. Strong Progress	10
5. Approach	10
6. New API Routines	11
6.1. Deferred MPI Operations	11
6.2. Graph Management	12
6.2.1. Graph Variables	18
6.3. Stream Management	21
6.4. Memory Operations	23
6.5. Synchronization Operations	23
6.6. Local MPI Data Types	24
7. Use Cases	24
7.1. Use Case #1: Communication and Computation Overlap	24
7.2. Use Case #2: Asynchronously Allocate and Receive Unknown Message	26
7.3. Use Case #3: Prefetch compressed data from file and broadcast to other ranks	27
7.4. Use Case #4: Offload Collective I/O For Writing Checkpoint File	30
7.5. Use Case #5: Data Dependent Asynchronous Communication	32
References	37
A. Appendix: An Abstract Data Movement Machine	37
A.1. Data Containers	37
A.2. Container Operations	37
A.2.1. Abstract Copy Operations	37
A.2.2. Abstract Reduce Operations	38
A.2.3. Abstract Comparison Operations	38
A.2.4. Abstract Synchronization Operations	38
A.2.5. Constants	38
A.3. Graphs	39
A.3.1. Graph Task Execution Unit	39
A.3.2. Graph Dependency Unit	39
A.3.3. Graph Construction Rules	44
A.3.4. Graph Variables	44

A.4. Streams	45
A.4.1. Stream Task Execution Unit	46
A.4.2. Stream Construction Rules	46
A.4.3. Stream Variables	46
A.5. Scopes	49
B. Appendix: Comparison to proposed Continuations and MPIX_Stream extensions	50
C. Appendix: What this proposal is <i>not</i>	51
D. Appendix: Details for Asynchronous Operations Option #2	52

List of Figures

1.	Task dependencies	40
2.	General form of unconditional organizational graph dependencies	40
3.	Common forms of unconditional organizational graph dependencies	41
4.	Threshold organizational graph dependency	42
5.	Branch organizational graph dependency	42
6.	Control flow operation	43
7.	Example graph for control flow destinations	43
8.	Connecting multiple task dependencies	44
9.	Connecting multiple organizational dependencies	45
10.	Combining task and organizational dependencies	46
11.	Combining tasks, organizational dependencies, and control flow operations to create loop in graph	47
12.	Combining tasks, organizational dependencies, and sub-graphs	47
13.	Tasks in a stream	48
14.	Tasks and sub-graphs in a stream	48

List of Listings

1.	Deferred operations added to graphs	9
2.	Deferred operations added to streams	10
3.	Data-dependent communication operations	25
4.	Asynchronously allocate space for and receive message of unknown length	27
5.	Asynchronously prefetch a compressed file on rank 0 and broadcast the decompressed buffer to other ranks	30
6.	Asynchronously write checkpoint file with collective I/O	32
7.	Data Dependent Asynchronous Operations	35

1. Introduction

The goal of this RFC is to present a proposal for improving performance of MPI applications in 3 primary ways: reduce synchronization penalties for collective operations, fully overlap communication, I/O, and computation, and enable more MPI operations to execute concurrently with application computation.

Expand on the goals

2. Background: MPI Today

MPI currently provides two kinds of communication and I/O operations: blocking and nonblocking (which are also sometimes called “immediate”). Blocking operations do not return until the operation is complete (or an error occurs). Nonblocking operations return immediately and provide a “request” object (`MPI_Request`) to the caller, which can use the request to check for the operation’s completion (or error). Initially, nonblocking operations appear to meet the goals for strong progress with asynchronous operation, but there are many restrictions and caveats to them, as shown below.

To begin with, a simple set of blocking MPI operations might look like this:

<u>Rank 0</u>	<u>Rank 1</u>
<code>MPI_Send(...)</code>	<code>MPI_Recv(...)</code>
<code>MPI_Send(...)</code>	<code>MPI_Recv(...)</code>
<code>...</code>	<code>...</code>
<code>MPI_Send(...)</code>	<code>MPI_Recv(...)</code>

Clearly, there’s no way to overlap these communication operations with any computation, as all the MPI operations are blocking.

Attempting to use nonblocking MPI operations might look something like this:

<u>Rank 0</u>	<u>Rank 1</u>
<code>MPI_Isend(..., &req[0])</code>	<code>MPI_Irecv(..., &req[0])</code>
<code>MPI_Isend(..., &req[1])</code>	<code>MPI_Irecv(..., &req[1])</code>
<code>...</code>	<code>...</code>
<code>MPI_Isend(..., &req[n])</code>	<code>MPI_Irecv(..., &req[n])</code>
<code>[Compute]</code>	<code>[Compute]</code>
<code>MPI_Waitall(..., n+1, req)</code>	<code>MPI_Waitall(..., n+1, req)</code>

However, using nonblocking operations has many drawbacks:

- a) Nonblocking operations are not guaranteed to make progress unless the MPI library is entered
- b) An application must track individual nonblocking operations with a request object for each one
- c) There is no way to indicate dependencies between nonblocking operations
- d) There is no way to invoke a user operation for asynchronous execution¹
- e) Applications and middleware libraries can’t easily nest asynchronous operations on or around the MPI API

¹Generalized requests are not a solution that meets the requirements outlined here.

- f) A result from one asynchronous operation can't be used as a "future" value for another operation
- g) There are no asynchronous operations for memory allocation or release, as well as local memory movement
- h) Many file operations are not supported

The asynchronous operations described in this document correct the drawbacks of nonblocking operations as well as put new capabilities into the hands of application developers.

3. Motivation / Etc.

4. Overview

Building a robust set of extensions to add "true" asynchronous operations to MPI can achieve the goal of improving application performance by enabling full overlap with computation, which can reduce the costs for communication and I/O to nearly zero.

Additional benefits of achieving the primary performance goals in an elegant and well-designed way are: an improved 'user experience' for developers using asynchronous (currently 'nonblocking') operations, hiding the latency of operations, exposing more opportunities for optimizing performance of MPI operations, enabling offload of more operations to networking hardware, and enabling applications to build powerful data movement orchestration operations.

Include benchmarks that show the benefit of overlapping communication and I/O with computation

4.1. Deferred Operations

This document describes two mechanisms for executing asynchronous MPI operations: defining an aggregated set of operations (a "graph") to execute later and executing operations immediately in an ordered manner (a "stream"). Each of those mechanisms relies on a new concept: "deferred" MPI operations. Deferred operations are very similar to persistent operations, but are designed to cover all kinds of MPI operations, unlike persistent operations, which focus on communication operations.^{2 3}

Executing deferred MPI operations in a "fully" asynchronous manner, which are guaranteed to complete without re-entering the MPI implementation, requires strong progress, described in section 4.3.

Listings 1 and 2 are pseudocode examples that shows the approach, with graphs and streams.

Implementing deferred operations in this way allows for a single API definition for each operation, which returns a `MPI_Token` object that can be added to a graph or enqueued on a stream.

4.1.1. Data Dependencies For Deferred Operations

Deferred operations may have data dependencies on values produced in earlier operations, as well as produce values that later operations may wish to consume. MPI operations have two kinds of data dependencies: strong dependencies, which involve *handles* of objects, and weak dependencies, which involve the *contents* of objects. For example, a strong dependency is created by the `MPI_File` file handle produced from a deferred call to `MPI_File_open` that is used in a later deferred call to `MPI_File_get_size`. A weak

²Deferred operations can cover `MPI_File_open`, `MPI_Comm_create`, etc. along with `MPI_Send` and related operations that are covered by the set of persistent operations.

³Nonblocking operations will not work, because they are allowed to start execution immediately and are not compatible with the approach described here.


```

// Info keys control graph execution behavior
MPIX_Graph_create(session, info, &graph);
...
// Define deferred operations and add to graph
MPIX_File_open_def(..., &token);
MPIX_Graph_add(graph, &token, <dependency info>);
...
MPIX_File_read_def(..., &token);
MPIX_Graph_add(graph, &token, <dependency info>);
...
MPIX_File_close_def(..., &token);
MPIX_Graph_add(graph, &token, <dependency info>);
...
MPIX_Bcast_def(..., &token);
MPIX_Graph_add(graph, &token, <dependency info>);
...
// Create deferred operation token for graph
MPIX_Graph_def(graph, info, &token);
...
// Create an [inactive] request for the graph token
MPIX_Execute_init(token, info, MPIX_OFFLOAD_CUDA, &cuda_stream, &req);
...
// Start execution of graph
MPIX_Start(req);
...
<compute or other overlap w/graph execution>
...
// Conclude graph's operations
MPI_Wait(&req, &status);

```

Listing 1. – Deferred operations added to graphs

data dependency is demonstrated by the data in a buffer from a deferred call to `MPI_Recv` being used as the source buffer for a deferred call to `MPI_Bcast` that has a dependency on the call to `MPI_Recv`.

4.2. Aggregating Deferred Operations: MPI Graphs and Streams

Graphs and streams are the core objects that “hold” deferred operations. Deferred operations can be added to a graph, with optional dependencies on other asynchronous operations. Dependencies between deferred operations can describe a directed graph and deferred operations without dependencies can execute concurrently. Deferred operations added to a stream execute in FIFO order, without requiring explicit dependencies, but also without the possible concurrent execution that is possible with graphs.⁴

MPI graphs have an additional supporting object: variables. MPI graphs allow the composition of Turing-complete data movement kernels using graph variables to control the execution of if/else blocks and loops. Graph variables are also used to parameterize graphs, allowing for their re-use with new inputs.

⁴For full details on graph and stream construction rules, see sections A.3.3 and A.4.2 respectively.

```

// Info keys control stream execution behavior
MPIX_Stream_create(session, info, MPIX_OFFLOAD_CUDA, &cuda_stream, &stream);
...
// Pause stream from processing operations that are added
MPIX_Stream_pause(stream);
...
// Define deferred operations and enqueue into stream
MPIX_File_open_def(..., &token);
MPIX_Stream_enqueue(stream, &token);
...
MPIX_File_read_def(..., &token);
MPIX_Stream_enqueue(stream, &token);
...
MPIX_File_close_def(..., &token);
MPIX_Stream_enqueue(stream, &token);
...
MPIX_Bcast_def(..., &token);
MPIX_Stream_enqueue(stream, &token);
...
// Start processing all the enqueued operations
MPIX_Stream_resume(stream);
...
// Ensure that all operations currently on stream have finished
MPIX_Stream_sync(stream); // Optional
...

```

Listing 2. – Deferred operations added to streams

4.2.1. Errors When Executing Deferred Operations Asynchronously

There can be many asynchronous operations in an aggregation object⁵, possibly executing concurrently. Additionally, data dependencies on information sent from other MPI ranks at runtime are possible. This complex and unpredictable environment argues against investing in mechanisms to resume/restart from failures.

Therefore, if an error occurs when executing a deferred operation in an aggregation object, the aggregation object will permanently stop scheduling new operations for execution. Existing operations that are executing will be allowed to complete, but no further operations will be started.

When an error occurs, only three actions on the aggregation object are possible: use the “introspection” API routines on the object (mainly to query about the failed operation, although any introspection operation can be called), duplicate the object (for possible possible “restart from the beginning”), and free the object⁶.

4.3. Strong Progress

5. Approach

⁵An MPI graph or stream

⁶You can autopsy or clone it before you bury it, but it's dead.

Describe enabling strong progress in MPI sessions, similar to enabling threading

Add de-

6. New API Routines

6.1. Deferred MPI Operations

TODO: Describe memory operation APIs...

```
MPIX_ZZZ_DEF(..., token)
-      ...      existing parameters for operation
OUT    token    deferred operation token (handle)
```

`MPIX_ZZZ_DEF` corresponds to a deferred version of the `MPI_ZZZ` operation. The `token` object may be added to a graph or stream.

```
MPIX_OP_DEF(op, extra_state, token)
IN      op      function pointer to user callback to invoke
IN      extra_state  pointer to extra state for user callback
OUT     token   deferred operation token (handle)
```

`MPIX_OP_DEF` creates a deferred invocation of a user callback.⁷ User callbacks invoked through this mechanism are designed to be “support” routines, not intensive computations. The `token` object may be added to a graph or stream.

NOTE: The definition of the user callback would probably look like:

```
void (*op)(MPIX_Graph *graph, void *extra_state)
```

The `op` callback receives a pointer to the `MPIX_Graph` that it was invoked from, as well as the `extra_state` specified in the `MPIX_OP_DEF` call. The `graph` parameter may be used for calls to `MPIX_GRAPH_VAR_GET` and `MPIX_GRAPH_VAR_SET`, to access the graph’s variables.

NOTE: This is designed as an “escape hatch” for user operations that should execute in a graph or stream, but aren’t covered by the current deferred operations in the MPI API.

```
MPIX_EXECUTE(token, info, loc, loc_info)
IN      token   deferred operation token (handle)
IN      info    info object (handle)
IN      loc_type  offload type (handle)
IN      loc_info  pointer to offload location info (choice)
```

`MPIX_EXECUTE` executes the deferred operation specified in `token`. The deferred token object is not consumed by this operation, allowing it to be used for multiple execute operation. Hints may be provided by the `info` argument. The execution offload type of the operation is specified by `loc_type`, with additional information provided in `loc_info`. Valid values for `loc_type` are given in <side document>, but implementations must at least support `MPIX_LOC_INPLACE`, for local execution without offload.

⁷MPI generalized requests are similar, but are not able to be invoked from a graph or stream.

Discuss deferred operations, how they are similar / different to persistent and non-blocking operations.

Mention that deferred persistent, nonblocking, and partitioned operations are possible.

Create table of important deferred operations, including file I/O, and alloc / free mem.

Possibly other parameters?

There are potentially lots of sharp edges here, so we should add more cautionary notes and user guidance.

CUDA, ROCM.

.....

```
MPIX_EXECUTE_INIT(token, info, loc_type, loc_info, req)
  IN      token      deferred operation token (handle)
  IN      info       info object (handle)
  IN      loc_type   offload type (handle)
  IN      loc_info   pointer to offload location info (choice)
  OUT     req        request object (handle)
```

MPIX_EXECUTE_INIT is a persistent version of MPIX_EXECUTE, returning an inactive request object in req. The request must be passed to a function in the MPI_Start family of operations to become active.

.....

```
MPIX_TOKEN_FREE(token)
  INOUT   token      deferred operation token (handle)
```

MPIX_TOKEN_FREE releases a deferred operation and sets the value of token to MPIX_TOKEN_NULL. Freeing a token has no effect on deferred operations or requests which have used the token already.

6.2. Graph Management

TODO: Describe graph management APIs...

Probably
need to
introduce
operation
IDs

```
MPIX_GRAPH_CREATE(session, info, graph)
  IN      session    session (handle)
  IN      info       info object (handle)
  OUT     graph      graph object (handle)
```

MPIX_GRAPH_CREATE creates a new graph object in the specified session. An info object is provided in order to support optimization hints and other information that may be nonstandard.

.....

```
MPIX_GRAPH_DUP(graph, info, newgraph)
  IN      graph      graph (handle)
  IN      info       info object (handle)
  OUT     newgraph   copy of graph (handle)
```

MPIX_GRAPH_DUP duplicates the existing graph graph. Hints provided by the argument info are associated with the output graph newgraph.

.....

```
MPIX_GRAPH_ADD(graph, token, op_id, dep_op_id)
  INOUT  graph      graph (handle)
  INOUT  token      deferred operation token (handle)
  OUT    op_id      operation ID (non-negative integer)
  IN     dep_op_id  operation ID (non-negative integer)
```

`MPIX_GRAPH_ADD` adds the deferred operation `token` to the graph `graph`. The graph takes ownership of the deferred token object, freeing it and setting `token` to `MPIX_TOKEN_NULL`.

Optionally an operation ID `op_id` can be returned to the application to refer to this operation in the graph, to create dependencies on it. `NULL` may be passed to indicate that no operation ID should be returned.

`dep_op_id` may also optionally be provided, to create a dependency on an earlier operation already added to the graph. `MPIX_DEP_NONE` may be passed to indicate that the added operation has no parent dependency.

NOTE: Operation IDs are only valid for creating dependencies between operations in the same graph. Using an operation ID from one graph in another graph is not supported and may cause undefined behavior.

.....

```
MPIX_GRAPH_JOIN(graph, count, array_of_dep_op_ids, op_id)
  INOUT  graph      graph (handle)
  IN     count      array size (non-negative integer)
  IN     array_of_dep_op_ids
                        array of dependent operation IDs
                        (array of non-negative integers)
  OUT    op_id      operation ID (non-negative integer)
```

`MPIX_GRAPH_JOIN` provides the capability to have many-parent to one-child relationships in a graph. `MPIX_GRAPH_JOIN` produces an operation ID `op_id` that can be used as a dependent operation ID for further operations in the graph that depend on all the operations in `array_of_dep_op_ids` completing before the operation can execute. All of the operation IDs in `array_of_dep_op_ids` must be from operations in `graph`.

.....

```
MPIX_GRAPH_IF_ELSE(graph, var, true_graph, false_graph,
                    op_id, dep_op_id)
    INOUT    graph      graph (handle)
    IN       var_id     variable identifier (non-negative integer)
    INOUT    true_graph  graph (handle)
    INOUT    false_graph graph (handle)
    OUT     op_id       operation ID (non-negative integer)
    IN      dep_op_id   operation ID (non-negative integer)
```

`MPIX_GRAPH_IF_ELSE` adds a conditional dependency to `graph`. If the value of the graph variable `var_id` is non-zero, the `true_graph` graph will be the child of the `MPIX_GRAPH_IF_ELSE` dependency and if `var` is zero, the `false_graph` graph will be the child of the dependency.

The `MPIX_GRAPH_IF_ELSE` operation *copies* both the `true_graph` and the `false_graph`, allowing them to be re-used for other purposes (and eventually released with `MPIX_GRAPH_FREE`). `MPIX_GRAPH_NULL` may be passed for either the true or false graph handles, indicating that there is no child dependency for that condition.

Comment [QK2]: Or increments the refcount on?

Optionally an operation ID `op_id` can be returned to the application to refer to this operation in the graph, to create dependencies on it. `dep_op_id` may also optionally be provided, to create a dependency on an earlier operation already added to the graph.

NOTE: Examining the graph variable within the `MPIX_GRAPH_IF_ELSE` operation could have weak data dependencies on other concurrently executing operations in the graph and create race conditions when the graph is executed. Application developers can use other organizational dependencies as preludes to `MPIX_GRAPH_IF_ELSE` to eliminate such race conditions.

.....

```

MPIX_GRAPH_WHILE(graph, var_id, sub_graph, op_id, dep_op_id)
  INOUT  graph      graph (handle)
  IN     var_id     graph variable (handle)
  INOUT  sub_graph  graph (handle)
  OUT    op_id      operation ID (non-negative integer)
  IN     dep_op_id  operation ID (non-negative integer)

```

MPIX_GRAPH_WHILE adds a loop to `graph`. If the value of the graph variable `var_id` is non-zero, the `sub_graph` graph will be executed.

The MPIX_GRAPH_WHILE operation *copies* the `sub_graph`, allowing it to be re-used for other purposes (and eventually released with MPIX_GRAPH_FREE).

Comment [QK3]: Or increments the refcount on?

Optionally an operation ID `op_id` can be returned to the application to refer to this operation in the graph, to create dependencies on it. `dep_op_id` may also optionally be provided, to create a dependency on an earlier operation already added to the graph.

NOTE: Examining the graph variable within the MPIX_GRAPH_WHILE operation could have weak data dependencies on other concurrently executing operations in the graph and create race conditions when the graph is executed. Application developers can use other organizational dependencies as preludes to MPIX_GRAPH_WHILE to eliminate such race conditions.

.....

```

MPIX_GRAPH_GOTO(graph, dest_op_id, dep_op_id)
  INOUT  graph      graph (handle)
  IN     dest_op_id  operation ID (non-negative integer)
  IN     dep_op_id  operation ID (non-negative integer)

```

MPIX_GRAPH_GOTO transfers the flow of control within `graph` to another operation in the graph, `dest_op_id`.

Optionally an operation ID, `dep_op_id`, may be provided, to create a dependency on an earlier operation already added to the graph.

NOTE: It is possible to create > 1 MPIX_GRAPH_GOTO operation as a dependency on `dep_op_id`, starting concurrent execution of the `dest_op_id` for each MPIX_GRAPH_GOTO operation.

.....

```
MPIX_GRAPH_SET_DEP_OP(graph, op_id, dep_op_id)
  INOUT  graph      graph (handle)
  IN     op_id      operation ID (non-negative integer)
  IN     dep_op_id  operation ID (non-negative integer)
```

`MPIX_GRAPH_SET_DEP_OP` sets or modifies the `dest_op_id` for `op_id` in `graph`. The value 0 may be passed for `dest_op_id` to remove the dependent operation for `op_id`. `dep_op_id` may not be set to the same value as `op_id`.

NOTE: Application developers should take care not to create dependency loops.

.....

```
MPIX_GRAPH_NOP(graph, op_id, dep_op_id)
  INOUT  graph      graph (handle)
  IN     op_id      operation ID (non-negative integer)
  IN     dep_op_id  operation ID (non-negative integer)
```

`MPIX_GRAPH_NOP` creates a “no op” operation in `graph`.

Optionally an operation ID `op_id` can be returned to the application to refer to this operation in the graph, to create dependencies on it. `dep_op_id` may also optionally be provided, to create a dependency on an earlier operation already added to the graph.

NOTE: `MPIX_GRAPH_NOP` and `MPIX_GRAPH_SET_DEP_OP` may be useful when creating and updating ‘placeholder’ operations as destinations when constructing complex control flows for graphs.

.....

```
MPIX_GRAPH_SET_PARAM(graph, var_id, value / alias_var_id)
  INOUT  graph      graph (handle)
  IN     var_id      parameter variable identifier (non-negative integer)
  IN     value or    pointer to parameter value (choice) or alias variable
         alias_var_id identifier (non-negative integer)
```

`MPIX_GRAPH_SET_PARAM` is overloaded to accept either a variable ID or a pointer to a value. If a variable ID is passed, the graph parameter variable, `var_id`, must be a reference parameter variable and if a pointer to a value is passed, the graph parameter variable must be a value parameter (see [MPIX_GRAPH_VAR_CREATE](#) for more information). Graph parameter variables can be set or changed at any time before or after `MPIX_GRAPH_DEF` is called.

.....

```
MPIX_GRAPH_DEF(graph, info, token)
  INOUT  graph      graph (handle)
  IN     info       info object (handle)
  IN     token      deferred operation token (handle)
```

`MPIX_GRAPH_DEF` binds the parameters for the graph `graph` to a deferred operation token `token`. Further changes to the parameters for the graph do not change the parameters for the initialized graph associated with `token`. The deferred operation token representing the initialized graph may be added to another graph (with `MPIX_GRAPH_ADD`), enqueued in a stream (with `MPIX_STREAM_ENQUEUE`), or executed with `MPIX_EXECUTE` or `MPIX_IEXECUTE`.

NOTE: If any parameter variables are defined but not set for the graph (with `MPIX_GRAPH_SET_PARAM`), `MPIX_GRAPH_DEF` will fail.

.....

```
MPIX_GRAPH_FREE(graph)
  INOUT  graph      graph (handle)
```

`MPIX_GRAPH_FREE` frees `graph` and sets it to `MPIX_GRAPH_NULL`.

6.2.1. Graph Variables

TODO: Describe graph variables and APIs...

```
MPIX_GRAPH_VAR_CREATE(graph, var_id, scope, type)
  INOUT  graph      graph object (handle)
  IN     var_id     variable identifier (non-negative integer)
  IN     scope      scope of variable (handle)
  IN     datatype   datatype of variable (handle)
```

`MPIX_GRAPH_VAR_CREATE` creates a new variable in the specified graph, `graph`, identified with `var_id`. `datatype` must be an atomic MPI datatype⁸. The scope of the variable is one of the following values, with the associated meaning:

- **MPIX_GRAPH_VAR_VAL_PARAM** – a variable that is used as a parameter for the graph, and must be initialized with a value using `MPIX_GRAPH_SET_PARAM` before the graph starts execution.
- **MPIX_GRAPH_VAR_REF_PARAM** – a variable that is an alias for a variable in the application or another graph variable in an outer scope, and must be initialized with `MPIX_GRAPH_SET_PARAM` before the graph starts execution.
- **MPIX_GRAPH_VAR_LOCAL** – a variable that is initialized to “all zeroes” when the graph starts execution, and is not aliased to a variable in an outer scope.

Maybe a better word than ‘scope’?

Variables may be declared at any point during a graph’s definition, but assignments for all value or reference parameter variables defined anywhere in the graph must be provided before `MPIX_GRAPH_DEF` is called on the graph.

Reference parameters are aliases for application variables or graph variables in an outer graph scope. The reference parameter can be thought of as an RMA window on the aliased application or outer graph variable. The value in the aliased variable is not valid for inspection through the aliased variable until the graph completes execution.

Graph variables are valid sources and destinations where memory buffers are specified for **all** deferred MPI data movement operations, such as `MPIX_REDUCE_DEF` and `MPIX_SEND_DEF`, etc. Graph variables can also be used for ‘IN’ parameters to deferred MPI operations and as ‘future’ values from deferred MPI operations with ‘OUT’ parameters, such as `MPI_FILE_OPEN`, etc.

Graph variables may be of a “local” MPI type, allowing a variable to be used as a parameter of the same type for deferred operations. For example, a graph variable may be created with a type of `MPI_TYPE_FILE` and used in all the places where an `MPI_FILE_*` operation accepts (an ‘IN’ parameter) or returns (as an ‘OUT’ parameter) a parameter of type `MPI_File`. Local MPI data types are described in section 6.6.

Should an application wish to set the value of a graph variable during execution of a graph or preserve a variable’s value before the graph completes execution, the `MPIX_COPY_MEM_DEF` operation should be used with a graph variable as the source or destination buffer (or both).

⁸Including `MPI_AINT`, `MPI_COUNT`, and ‘local’ MPI data types like `MPI_FILE`, `MPI_STATUS`, etc.

.....

`MPIX_GRAPH_VAR_SET(graph, var_id, value)`

INOUT	graph	graph (handle)
IN	var_id	parameter variable identifier (non-negative integer)
IN	value	pointer to new value for variable (choice)

`MPIX_GRAPH_VAR_SET` sets the value of variable `var_id` in graph `graph` to the value pointed to by the `value` pointer.

Graph variables set before `MPIX_GRAPH_DEF` is called will be associated with the graph token created by `MPIX_GRAPH_DEF` and later calls to `MPIX_GRAPH_VAR_SET` will not affect the values for the variables associated with that token.

Graph variables set to values in a user callback (with `MPIX_OP_DEF`) are reflected in the values for those variables when the callback completes. Race conditions around graph variables may exist depending on the graph's operations and an application developer should use operation dependencies within the graph to avoid them.

.....

`MPIX_GRAPH_VAR_GET(graph, var_id, value)`

INOUT	graph	graph (handle)
IN	var_id	parameter variable identifier (non-negative integer)
OUT	value	pointer to a location to place the value for variable (choice)

`MPIX_GRAPH_VAR_GET` retrieves the value of variable `var_id` in graph `graph` into the location pointed to by the `value` pointer.

Graph variable values retrieved within a user callback (with `MPIX_OP_DEF`) reflect the value of those variables when `MPIX_GRAPH_VAR_GET` is invoked. Race conditions around graph variables may exist depending on the graph's operations and an application developer should use operation dependencies within the graph to avoid them.

.....

```
MPIX_GRAPH_VAR_REDUCE(graph, in_var1_id, in_var2_id,
                       var_op, out_var_id, op_id, dep_op_id)
  INOUT  graph      graph (handle)
  IN     in_var1_id  variable identifier (non-negative integer)
  IN     in_var2_id  variable identifier (non-negative integer)
  IN     var_op      operation (handle)
  IN     out_var_id  variable identifier (non-negative integer)
  IN     op_id       operation ID (non-negative integer)
  IN     dep_op_id   operation ID (non-negative integer)
```

`MPIX_GRAPH_VAR_REDUCE` is similar to `MPI_REDUCE_LOCAL`, but performs an operation `op` on two graph variables, `in_var1_id` and `in_var2_id` and places the result in `out_var_id`. All three variables must belong to the same graph, `graph`.

`var_op` may be any of the predefined reduction operations defined for `MPI_REDUCE`, with the exception of `MPI_MAXLOC` and `MPI_MINLOC`. Additionally, the following operations are defined for `MPIX_GRAPH_VAR_REDUCE`:

- **MPIX_LESS_THAN_OR_EQUAL** – produces a non-zero or TRUE value when the value of `in_var1_id` is less than or equal to the value of `in_var2_id`, i.e. $in_var1_id \leq in_var2_id$.
- **MPIX_LESS_THAN** – produces a non-zero or TRUE value when the value of `in_var1_id` is less than the value of `in_var2_id`, i.e. $in_var1_id < in_var2_id$.
- **MPIX_EQUAL** – produces a non-zero or TRUE value when the value of `in_var1_id` is equal to the value of `in_var2_id`, i.e. $in_var1_id = in_var2_id$.
- **MPIX_NOT_EQUAL** – produces a non-zero or TRUE value when the value of `in_var1_id` is not equal to the value of `in_var2_id`, i.e. $in_var1_id \neq in_var2_id$.
- **MPIX_GREATER_THAN** – produces a non-zero or TRUE value when the value of `in_var1_id` is greater than the value of `in_var2_id`, i.e. $in_var1_id > in_var2_id$.
- **MPIX_GREATER_THAN_OR_EQUAL** – produces a non-zero or TRUE value when the value of `in_var1_id` is greater than or equal to the value of `in_var2_id`, i.e. $in_var1_id \geq in_var2_id$.

Optionally an operation ID `op_id` can be returned to the application to refer to this operation in the graph, to create dependencies on it. `dep_op_id` may also optionally be provided, to create a dependency on an earlier operation already added to the graph.

6.3. Stream Management

TODO: Describe stream management APIs...

`MPIX_STREAM_CREATE(session, info, loc_type, loc_info, stream)`

IN	session	session (handle)
IN	info	info object (handle)
IN	loc_type	offload type (handle)
IN	loc_info	pointer to offload info (choice)
OUT	stream	stream object (handle)

`MPIX_STREAM_CREATE` creates a new stream object in the specified session. An info object is provided in order to support optimization hints and other information that may be nonstandard. The execution offload type of the operation is specified by `loc_type`, with additional information provided in `loc_info`. Valid values for `loc_type` are given in <side document>, but implementations must at least support `MPIX_LOC_INPLACE`, for local execution without offload.

CUDA,
ROCM,
pthread,
etc.

.....

`MPIX_STREAM_ENQUEUE(stream, token)`

INOUT	stream	stream (handle)
INOUT	token	deferred operation token (handle)

`MPIX_STREAM_ENQUEUE` adds the deferred operation `token` to the stream `stream`. The stream takes ownership of the deferred token object, freeing it and setting `token` to `MPIX_TOKEN_NULL`.

Deferred operations enqueued in a stream are executed in the order they are added to the stream, i.e. ‘FIFO’ ordering. Deferred operation tokens may be individual operations, e.g. from `MPIX_SEND_DEF`, `MPIX_FILE_OPEN_DEF`, etc, or entire graphs, i.e. tokens returned from `MPIX_GRAPH_DEF`.

.....

`MPIX_STREAM_PAUSE(stream)`

INOUT	stream	stream (handle)
-------	--------	-----------------

`MPIX_STREAM_PAUSE` pauses execution of deferred operations on a stream, `stream`. Deferred operations may continue to be enqueued on the stream, but none will begin executing until `MPIX_STREAM_RESUME` is called.

.....

`MPIX_STREAM_IS_PAUSED(stream, paused)`

IN	stream	stream (handle)
OUT	paused	flag indicating stream is paused (logical)

`MPIX_STREAM_IS_PAUSED` queries whether execution of deferred operations on a stream, `stream`, is paused.

.....

`MPIX_STREAM_RESUME (stream)`

 INOUT stream stream (handle)

`MPIX_STREAM_RESUME` resumes execution of deferred operations on a stream, `stream`.

.....

`MPIX_STREAM_SYNC (stream)`

 INOUT stream stream (handle)

`MPIX_STREAM_SYNC` blocks until all enqueued operations on `stream` at the time it is called have finished execution. Deferred operations enqueued (with another thread) after this routine is called, but before it returns, are not included in the list of operations blocked for.

NOTE: If `MPIX_STREAM_SYNC` is called on a paused stream, the stream will resume and complete the deferred operations enqueued when `MPIX_STREAM_SYNC` was called.

.....

`MPIX_STREAM_ISYNC (stream, request)`

 INOUT stream stream (handle)

 OUT request request object (handle)

`MPIX_STREAM_ISYNC` is a nonblocking form of `MPIX_STREAM_SYNC`. The request, `request`, will complete when all the enqueued operations on `stream` at the time `MPIX_STREAM_ISYNC` is called have finished execution. Deferred operations enqueued after this routine is called are not included in the list of operations the request applies to.

NOTE: If `MPIX_STREAM_ISYNC` is called on a paused stream, the stream's execution is **not** resumed. A paused stream must be resumed with either `MPIX_STREAM_RESUME` or `MPIX_STREAM_SYNC`.

.....

`MPIX_STREAM_FREE (stream)`

 INOUT stream stream (handle)

`MPIX_STREAM_FREE` frees `stream` and sets it to `MPIX_STREAM_NULL`.

6.4. Memory Operations

TODO: Describe memory operation APIs...

`MPIX_COPY_MEM(dst, src, size)`

INOUT	<code>dst</code>	pointer to beginning of destination memory segment (choice)
IN	<code>src</code>	pointer to beginning of source memory segment (choice)
IN	<code>size</code>	size of memory segment in bytes (non-negative integer)

`MPIX_COPY_MEM` copies `size` bytes from memory area `src` to memory area `dst`. If `dst` and `src` overlap, behavior is undefined.

NOTE: This is the same behavior as Standard C `memcpy`.

NOTE: `MPIX_COPY_MEM_DEF` is the mechanism for setting / getting values for graph variables.

Add deferred version of this routine also

Maybe we want mem-move also / instead?

6.5. Synchronization Operations

TODO: Describe new synchronization operation APIs...

`MPIX_REMOTE_BARRIER(comm, watched, info)`

IN	<code>comm</code>	communicator (handle)
IN	<code>watched</code>	rank of watched MPI process (integer)
IN	<code>info</code>	info object (handle)

`MPIX_REMOTE_BARRIER` blocks the caller until the `watched` MPI process has called it. For the `watched` process, `MPIX_REMOTE_BARRIER` is local and never blocks. An info object is provided in order to support optimization hints and other information that may be nonstandard.

NOTE: `MPIX_REMOTE_BARRIER_DEF` can be used to create inter-graph dependencies between concurrently executing graphs.

Add deferred version of this routine also

C Typedefs	Local Type Name
MPI_Comm	MPI_LOCAL_TYPE_COMM
MPI_Datatype	MPI_LOCAL_TYPE_DATATYPE
MPI_Errhandler	MPI_LOCAL_TYPE_ERRHANDLER
MPI_File	MPI_LOCAL_TYPE_FILE
MPI_Group	MPI_LOCAL_TYPE_GROUP
MPI_Info	MPI_LOCAL_TYPE_INFO
MPI_Message	MPI_LOCAL_TYPE_MESSAGE
MPI_Op	MPI_LOCAL_TYPE_OP
MPI_Request	MPI_LOCAL_TYPE_REQUEST
MPI_Session	MPI_LOCAL_TYPE_SESSION
MPI_Status	MPI_LOCAL_TYPE_STATUS
const char *	MPI_LOCAL_TYPE_STRING
MPI_Win	MPI_LOCAL_TYPE_WIN

Table 1 – C Typedefs and equivalent Local Data Types

6.6. Local MPI Data Types

TODO: Describe local MPI data types

TODO: Note ‘const char *’ as MPI_LOCAL_TYPE_STRING

7. Use Cases

Use cases that demonstrate the need and applicability of these new capabilities are given below.

7.1. Use Case #1: Communication and Computation Overlap

This [example code](#) demonstrates overlapping compute with communication, including: broadcasting the direction of neighbors to exchange with for the next timestep computation and a exchange of data between neighbors, all of which occurs asynchronously, overlapping the current compute step.

```

1 MPIX_Graph graph;
2 MPIX-Token token;
3 MPI_Request req;
4 enum {BCAST_VAL, DECR_VAL}; // Graph variable IDs
5
6 // Create graph for each timestep
7 MPIX_Graph_create(session, info, &graph);
8
9 // Create a value parameter graph variable to hold the broadcast information
10 MPIX_Graph_var_create(&graph, BCAST_VAL, MPIX_GRAPH_VAR_VAL_PARAM, MPI_INT);
11
12 // Create a local graph variable to hold the constant '-1'

```



```

13 int dec_value = -1;
14 MPIX_Graph_var_create(&graph, DECR_VAL, MPI_INT);
15 MPIX_Graph_var_set(&graph, DECR_VAL, &dec_value);
16
17 // Define deferred bcast w/a graph variable and add to graph
18 // Note: The graph variable ID is passed for the buffer pointer
19 unsigned bcast_op_id = 0;
20 MPIX_Bcast_def(BCAST_VAL, ..., &token);
21 MPIX_Graph_add(&graph, &token, &bcast_op_id, MPIX_DEP_NONE);
22
23 // Define sub-graph for 'while' loop
24 MPIX_Graph_create(session, info, &while_graph);
25 // Exchange with neighbor
26 unsigned sendrecv_op_id = 0;
27 MPIX_Sendrecv_def(..., &token);
28 MPIX_Graph_add(&while_graph, &token, &sendrecv_op_id, MPIX_DEP_NONE);
29 // Decrement loop counter
30 // NOTE: Use of var IDs from outer graph is OK, as variable IDs are resolved
   ↪ when MPIX_Graph_def is called
31 MPIX_Graph_var_reduce(&while_graph, BCAST_VAL, DECR_VAL, MPI_SUM, BCAST_VAL,
   ↪ NULL, sendrecv_op_id);
32
33 // While value of var ID 'BCAST_VAL' is non-zero, execute sub-graph
34 MPIX_Graph_while(&graph, BCAST_VAL, &while_graph, NULL, bcast_op_id);
35
36 // Free sub-graph for while loop
37 MPIX_Graph_free(&while_graph);
38
39 while(<more timesteps>) {
40     // Set parameter (var ID 'BCAST_VAL') for this execution of the graph
41     MPIX_Graph_set_param(&graph, BCAST_VAL, <neighbor exchange value>);
42
43     // Create deferred operation token for parameterized graph
44     MPIX_Graph_def(graph, info, &token);
45
46     // Create an [inactive] request for the graph token
47     MPIX_Execute_init(token, info, MPIX_OFFLOAD_CUDA, &cuda_stream, &req);
48
49     // Release token for parameterized graph
50     MPIX-Token_free(&token);
51
52     // Start execution of graph
53     MPI_Start(&req);
54
55     <compute current timestep>
56
57     // Conclude graph's operations, if it's not finished yet
58     MPI_Wait(&req, &status);
59 } // timestep loop
60
61 // Free graph
62 MPIX_Graph_free(&graph);

```

Listing 3. – Data-dependent communication operations

This use case shows a fairly straightforward overlap of compute and communication. If the # of neighbor exchanges didn't depend on the the information being broadcast, this could be done with nonblocking

operations today. However, the runtime data dependency makes this impossible with today's calls, without waiting for the broadcast to complete before issuing the sendrecv calls.

7.2. Use Case #2: Asynchronously Allocate and Receive Unknown Message

This [example code](#) demonstrates receiving an unknown-sized message in the background, with the buffer allocated asynchronously as well.

Although `MPI_Iprobe` and `MPI_Irecv` provide partial functionality needed for this use case, they can't provide the 'future' values needed (for the message size and allocated buffer), the necessary dependencies between operations, or the asynchronous memory allocation. All of these features must be present to allow this sequence of operations to execute independently of the application.

```

1  MPIX_Graph graph;
2  MPIX-Token token;
3  MPI-Request req;
4  enum {MSG_LEN, MSG_BUF_PTR}; // Graph variable IDs
5
6  // Create graph
7  MPIX_Graph_create(session, info, &graph);
8
9  // Create a reference parameter graph variable to alias the graph variable to
   ↪ an application buffer pointer
10 MPIX_Graph_var_create(&graph, MSG_BUF_PTR, MPIX_GRAPH_VAR_REF_PARAM,
   ↪ MPI_AINT);
11
12 // Create a local graph variable to hold the message size
13 MPIX_Graph_var_create(&graph, MSG_LEN, MPIX_GRAPH_VAR_LOCAL, MPI_COUNT);
14
15 // Receive the buffer size into a graph variable
16 unsigned size_recv_op_id = 0;
17 MPIX_Recv_def(MSG_LEN, 1, MPI_COUNT, MPI_ANY_SOURCE, MPI_ANY_TAG, comm,
   ↪ MPI_STATUS_IGNORE, &token);
18 MPIX_Graph_add(&graph, &token, &size_recv_op_id, MPIX_DEP_NONE);
19
20 // Allocate memory for the message, using the 'future' value for the message
   ↪ size and returning a future for the buffer pointer
21 unsigned buf_alloc_op_id = 0;
22 MPIX_Alloc_mem_def(MSG_LEN, MPI_INFO_NULL, MSG_BUF_PTR, &token);
23 MPIX_Graph_add(&graph, &token, &buf_alloc_op_id, size_recv_op_id);
24
25 // Receive the actual message
26 unsigned msg_recv_op_id = 0;
27 MPIX_Recv_def(MSG_BUF_PTR, MSG_LEN, MPI_BYTE, MPI_ANY_SOURCE, MPI_ANY_TAG,
   ↪ comm, MPI_STATUS_IGNORE, &token);
28 MPIX_Graph_add(&graph, &token, &msg_recv_op_id, buf_alloc_op_id);
29
30 // Alias graph parameter for message buffer to application pointer to buffer
31 void *msg_buf = NULL;
32 MPIX_Graph_set_param(&graph, MSG_BUF_PTR, &msg_buf);
33
34 // Create deferred operation token for parameterized graph
35 MPIX_Graph_def(graph, info, &token);
36
37 // Free graph
38 MPIX_Graph_free(&graph);
39

```

```

40 // Create an [inactive] request for the graph token
41 MPIX_Execute_init(token, info, MPIX_OFFLOAD_PTHREAD, NULL, &req);
42
43 // Release token for parameterized graph
44 MPIX-Token_free(&token);
45
46 // Start execution of graph
47 MPI_Start(&req);
48
49 <compute>
50
51 // Conclude graph's operations, if it's not finished yet
52 MPI_Wait(&req, &status);
53
54 // Use the message buffer
55 <use msg_buf>

```

Listing 4. – Asynchronously allocate space for and receive message of unknown length

7.3. Use Case #3: Prefetch compressed data from file and broadcast to other ranks

This [example code](#) demonstrates asynchronously pre-fetching a compressed file on rank 0, decompressing it, and broadcasting the resulting buffer to the other ranks.

This use case demonstrates the power of offloading operations to a “very smart NIC”, one which can run both communication and I/O operations asynchronously from the application. The code below offloads file I/O, memory allocation, invoking a user callback, and communication, allowing the application to continue to perform computation in the foreground while all of the data movement and preparation operations for the next timestep occur in the background.

```

1  MPIX_Graph graph;
2  MPIX-Token token;
3  MPI_Request req;
4  enum {FILE_NAME, BCAST_COMM, BUF_LEN, BUF_PTR, // Graph parameter variable IDs
5        FILE_HANDLE, FILE_SIZE, COMP_BUF_PTR}; // Graph variable IDs
6
7  // Create graph
8  MPIX_Graph_create(session, info, &graph);
9
10 // Create a reference parameter graph variable to hold the decompressed buffer
11   ↪ size and pointer
12 MPIX_Graph_var_create(&graph, BUF_LEN, MPIX_GRAPH_VAR_REF_PARAM, MPI_COUNT);
13 MPIX_Graph_var_create(&graph, BUF_PTR, MPIX_GRAPH_VAR_REF_PARAM, MPI_AINT);
14
15 // Create a local graph variable to hold the communicator to use for
16   ↪ broadcasting data within the graph
17 MPIX_Graph_var_create(&graph, BCAST_COMM, MPIX_GRAPH_VAR_LOCAL,
18   ↪ MPI_LOCAL_TYPE_COMM);
19
20 // Rank 0 does the file I/O and bcasts the decompressed data
21 if (0 == rank) {
22   // Create a value parameter graph variable for the filename
23   MPIX_Graph_var_create(&graph, FILE_NAME, MPIX_GRAPH_VAR_VAL_PARAM,
24   ↪ MPI_LOCAL_TYPE_STRING);
25
26   // Create a local graph variable to hold the file handle

```

```

23 MPIX_Graph_var_create(&graph, FILE_HANDLE, MPIX_GRAPH_VAR_LOCAL,
   ↪ MPI_LOCAL_TYPE_FILE);
24
25 // Open the file
26 unsigned file_open_op_id = 0;
27 MPIX_File_open_def(MPI_COMM_SELF, FILE_NAME, MPI_MODE_RDONLY |
   ↪ MPI_MODE_SEQUENTIAL, MPI_INFO_NULL, FILE_HANDLE, &token);
28 MPIX_Graph_add(&graph, &token, &file_open_op_id, MPIX_DEP_NONE);
29
30 // Create a local graph variable to hold the file's size
31 MPIX_Graph_var_create(&graph, FILE_SIZE, MPIX_GRAPH_VAR_LOCAL, MPI_OFFSET);
32
33 // Get the file's size
34 unsigned file_size_op_id = 0;
35 MPIX_File_get_size_def(FILE_HANDLE, FILE_SIZE, &token);
36 MPIX_Graph_add(&graph, &token, &file_size_op_id, file_open_op_id);
37
38 // Create a local graph variable to hold the pointer to the buffer with
   ↪ compressed data
39 MPIX_Graph_var_create(&graph, COMP_BUF_PTR, MPIX_GRAPH_VAR_LOCAL, MPI_AINT);
40
41 // Allocate buffer for compressed data
42 unsigned alloc_comp_buf_op_id = 0;
43 MPIX_Alloc_mem_def(FILE_SIZE, MPI_INFO_NULL, COMP_BUF_PTR, &token);
44 MPIX_Graph_add(&graph, &token, &alloc_comp_buf_op_id, file_size_op_id);
45
46 // Read compressed data
47 unsigned file_read_op_id = 0;
48 MPIX_File_read_def(FILE_HANDLE, COMP_BUF_PTR, FILE_SIZE, MPI_BYTE,
   ↪ MPI_STATUS_IGNORE, &token);
49 MPIX_Graph_add(&graph, &token, &file_read_op_id, alloc_comp_buf_op_id);
50
51 // Close file
52 // NOTE: Concurrent with decompressing data
53 MPIX_File_close_def(FILE_HANDLE, &token);
54 MPIX_Graph_add(&graph, &token, NULL, file_read_op_id);
55
56 // Decompress data with user callback
57 // NOTE: Concurrent with closing the file
58 //
59 // The decompression function would look something like this:
60 // void decompress_fn(MPIX_Graph *graph, void *extra_state)
61 // {
62 //     MPI_Offset compressed_size = 0;
63 //     void *compressed_buf = NULL;
64 //     MPI_Count decompressed_size = 0;
65 //     void *decompressed_buf = NULL;
66 //
67 //     MPIX_Graph_var_get(graph, FILE_SIZE, &compressed_size);
68 //     MPIX_Graph_var_get(graph, COMP_BUF_PTR, &compressed_buf);
69 //     <decompression routine>(*IN:*/compressed_size, /*IN:*/compressed_buf,
   ↪ /*OUT:*/&decompressed_size, /*OUT:*/&decompressed_buf);
70 //     MPIX_Graph_var_set(graph, BUF_LEN, &decompressed_size);
71 //     MPIX_Graph_var_set(graph, BUF_PTR, &decompressed_buf);
72 // }
73 unsigned decomp_op_id = 0;
74 MPIX_Op_def(&decompress_fn, NULL, &token);
75 MPIX_Graph_add(&graph, &token, &decomp_op_id, file_read_op_id);

```

```

76
77 // Bcast the decompressed buffer size, from rank 0
78 unsigned size_bcast_op_id = 0;
79 MPIX_Bcast_def(BUF_LEN, 1, MPI_COUNT, 0, BCAST_COMM, &token);
80 MPIX_Graph_add(&graph, &token, &size_bcast_op_id, decomp_op_id);
81
82 // Bcast the decompressed buffer, from rank 0
83 MPIX_Bcast_def(BUF_PTR, BUF_LEN, MPI_BYTE, 0, BCAST_COMM, &token);
84 MPIX_Graph_add(&graph, &token, NULL, size_bcast_op_id);
85 } // end (0 == rank)
86 else { // (0 != rank)
87 // Receive the decompressed buffer size bcast from rank 0
88 unsigned size_bcast_op_id = 0;
89 MPIX_Bcast_def(BUF_LEN, 1, MPI_COUNT, 0, BCAST_COMM, &token);
90 MPIX_Graph_add(&graph, &token, &size_bcast_op_id, MPIX_DEP_NONE);
91
92 // Allocate buffer for decompressed data
93 unsigned alloc_decomp_buf_op_id = 0;
94 MPIX_Alloc_mem_def(BUF_LEN, MPI_INFO_NULL, BUF_PTR, &token);
95 MPIX_Graph_add(&graph, &token, &alloc_decomp_buf_op_id, size_bcast_op_id);
96
97 // Receive the decompressed buffer bcast from rank 0
98 MPIX_Bcast_def(BUF_PTR, BUF_LEN, MPI_BYTE, 0, BCAST_COMM, &token);
99 MPIX_Graph_add(&graph, &token, NULL, alloc_decomp_buf_op_id);
100 } // end (0 != rank)
101
102 // Set communicator for bcast operations in graph
103 MPI_Comm graph_comm = MPI_COMM_WORLD;
104 MPIX_Graph_var_set(&graph, BCAST_COMM, &graph_comm);
105
106 // Prefetch next timestep's data from disk while current timestep is
107 ↪ computing
108 while(<more timesteps>) {
109 // Rank 0 does the file I/O and bcasts the decompressed data
110 if (0 == rank) {
111     const char *filename = <filename for timestep 'n'>;
112     MPIX_Graph_set_param(&graph, FILE_NAME, &filename);
113 }
114
115 // Alias graph parameters for decompressed data buffer info to application
116 ↪ variables
117 void *buf_ptr = NULL;
118 MPI_Count buf_len = 0;
119 MPIX_Graph_set_param(&graph, BUF_PTR, &buf_ptr);
120 MPIX_Graph_set_param(&graph, BUF_LEN, &buf_len);
121
122 // Create deferred operation token for parameterized graph
123 MPIX_Graph_def(graph, info, &token);
124
125 // Create an [inactive] request for the graph token
126 MPIX_Execute_init(token, info, MPIX_OFFLOAD_NITRO, NULL, &req);
127
128 // Release token for parameterized graph
129 MPIX-Token_free(&token);
130
131 // Start execution of graph
132 MPI_Start(&req);
133

```

```

132 <compute>
133
134 // Conclude graph's operations, if it's not finished yet
135 MPI_Wait(&req, &status);
136
137 <swap data buffers for next timestep>
138 }
139
140 // Free graph
141 MPIX_Graph_free(&graph);

```

Listing 5. – Asynchronously prefetch a compressed file on rank 0 and broadcast the decompressed buffer to other ranks

7.4. Use Case #4: Offload Collective I/O For Writing Checkpoint File

This [example code](#) demonstrates asynchronously writing a checkpoint, using collective I/O.

This use case again shows fully overlapping I/O with compute. In particular, this example includes concurrently executing collective I/O operations, and using the “join” operator to create a single dependency for the file close operation to depend on.

```

1  MPIX_Graph graph;
2  MPIX-Token token;
3  MPI-Request req;
4  enum {FILE_NAME, FILE_COMM, // Graph parameter variable IDs for file
5        BUF1_OFFSET, BUF1_PTR, BUF1_COUNT, BUF1_TYPE, // Graph parameter
6        ↪ variable IDs for buffer #1
7        BUF2_OFFSET, BUF2_PTR, BUF2_COUNT, BUF2_TYPE, // Graph parameter
8        ↪ variable IDs for buffer #2
9        BUF3_OFFSET, BUF3_PTR, BUF3_COUNT, BUF3_TYPE, // Graph parameter
10       ↪ variable IDs for buffer #3
11       FILE_HANDLE}; // Graph local variable IDs for file
12
13 // Create graph
14 MPIX_Graph_create(session, info, &graph);
15
16 // Create a value parameter graph variable for the filename and file's
17 ↪ communicator
18 MPIX_Graph_var_create(&graph, FILE_NAME, MPIX_GRAPH_VAR_VAL_PARAM,
19 ↪ MPI_LOCAL_TYPE_STRING);
20 MPIX_Graph_var_create(&graph, FILE_COMM, MPIX_GRAPH_VAR_VAL_PARAM,
21 ↪ MPI_LOCAL_TYPE_COMM);
22
23 // Create a value parameter graph variables for the three buffers to write
24 MPIX_Graph_var_create(&graph, BUF1_OFFSET, MPIX_GRAPH_VAR_VAL_PARAM,
25 ↪ MPI_OFFSET);
26 MPIX_Graph_var_create(&graph, BUF1_PTR, MPIX_GRAPH_VAR_VAL_PARAM, MPI_AINT);
27 MPIX_Graph_var_create(&graph, BUF1_COUNT, MPIX_GRAPH_VAR_VAL_PARAM,
28 ↪ MPI_COUNT);
29 MPIX_Graph_var_create(&graph, BUF1_TYPE, MPIX_GRAPH_VAR_VAL_PARAM,
30 ↪ MPI_LOCAL_TYPE_DATATYPE);
31
32 MPIX_Graph_var_create(&graph, BUF2_OFFSET, MPIX_GRAPH_VAR_VAL_PARAM,
33 ↪ MPI_OFFSET);
34 MPIX_Graph_var_create(&graph, BUF2_PTR, MPIX_GRAPH_VAR_VAL_PARAM, MPI_AINT);
35 MPIX_Graph_var_create(&graph, BUF2_COUNT, MPIX_GRAPH_VAR_VAL_PARAM,
36 ↪ MPI_COUNT);

```

```

26 MPIX_Graph_var_create(&graph, BUF2_TYPE, MPIX_GRAPH_VAR_VAL_PARAM,
   ↪ MPI_LOCAL_TYPE_DATATYPE);
27
28 MPIX_Graph_var_create(&graph, BUF3_OFFSET, MPIX_GRAPH_VAR_VAL_PARAM,
   ↪ MPI_OFFSET);
29 MPIX_Graph_var_create(&graph, BUF3_PTR, MPIX_GRAPH_VAR_VAL_PARAM, MPI_AINT);
30 MPIX_Graph_var_create(&graph, BUF3_COUNT, MPIX_GRAPH_VAR_VAL_PARAM,
   ↪ MPI_COUNT);
31 MPIX_Graph_var_create(&graph, BUF3_TYPE, MPIX_GRAPH_VAR_VAL_PARAM,
   ↪ MPI_LOCAL_TYPE_DATATYPE);
32
33 // Create a local graph variable to hold the file handle
34 MPIX_Graph_var_create(&graph, FILE_HANDLE, MPIX_GRAPH_VAR_LOCAL,
   ↪ MPI_LOCAL_TYPE_FILE);
35
36 // Create the file
37 unsigned file_create_op_id = 0;
38 MPIX_File_open_def(FILE_COMM, FILE_NAME, MPI_MODE_CREATE | MPI_MODE_RDWR,
   ↪ MPI_INFO_NULL, FILE_HANDLE, &token);
39 MPIX_Graph_add(&graph, &token, &file_create_op_id, MPIX_DEP_NONE);
40
41 // Collectively write all three buffers, concurrently
42 unsigned file_write_buf_op_ids[3] = {0, 0, 0};
43 MPIX_File_write_at_all_c_def(FILE_HANDLE, BUF1_OFFSET, BUF1_PTR, BUF1_COUNT,
   ↪ BUF1_TYPE, MPI_STATUS_IGNORE, &token);
44 MPIX_Graph_add(&graph, &token, &file_write_buf_op_ids[0], file_create_op_id);
45
46 MPIX_File_write_at_all_c_def(FILE_HANDLE, BUF2_OFFSET, BUF2_PTR, BUF2_COUNT,
   ↪ BUF2_TYPE, MPI_STATUS_IGNORE, &token);
47 MPIX_Graph_add(&graph, &token, &file_write_buf_op_ids[1], file_create_op_id);
48
49 MPIX_File_write_at_all_c_def(FILE_HANDLE, BUF3_OFFSET, BUF3_PTR, BUF3_COUNT,
   ↪ BUF3_TYPE, MPI_STATUS_IGNORE, &token);
50 MPIX_Graph_add(&graph, &token, &file_write_buf_op_ids[2], file_create_op_id);
51
52 // Create a 'join' operation ID, for all writes
53 unsigned join_writes_op_id = 0;
54 MPIX_Graph_join(&graph, 3, file_write_buf_op_ids, &join_writes_op_id);
55
56 // Close file
57 // NOTE: Depends on join of all concurrent collective writes
58 MPIX_File_close_def(FILE_HANDLE, &token);
59 MPIX_Graph_add(&graph, &token, NULL, join_writes_op_id);
60
61 // Compute, then write checkpoint, overlapped with next compute
62 boolean checkpoint_io_started = FALSE;
63 while(<more timesteps>) {
64     <compute>
65
66     // Make certain that previous checkpoint has completed
67     if (checkpoint_io_started) {
68         MPI_Wait(&req, &status);
69         checkpoint_io_started = FALSE;
70     }
71
72     // Write checkpoint
73     if (<time to make a checkpoint>) {
74         const char *filename = <checkpoint filename>;

```

```

75 MPI_Comm file_comm = <communicator for file I/O>;
76
77 // Set filename & communicator
78 MPIX_Graph_set_param(&graph, FILE_NAME, &filename);
79 MPIX_Graph_set_param(&graph, FILE_COMM, &file_comm);
80
81 // Set buffer info
82 MPIX_Graph_set_param(&graph, BUF1_OFFSET, &buf1_off[rank]);
83 MPIX_Graph_set_param(&graph, BUF1_PTR, &buf1_ptr[rank]);
84 MPIX_Graph_set_param(&graph, BUF1_COUNT, &buf1_count[rank]);
85 MPIX_Graph_set_param(&graph, BUF1_TYPE, &buf1_type[rank]);
86
87 MPIX_Graph_set_param(&graph, BUF2_OFFSET, &buf2_off[rank]);
88 MPIX_Graph_set_param(&graph, BUF2_PTR, &buf2_ptr[rank]);
89 MPIX_Graph_set_param(&graph, BUF2_COUNT, &buf2_count[rank]);
90 MPIX_Graph_set_param(&graph, BUF2_TYPE, &buf2_type[rank]);
91
92 MPIX_Graph_set_param(&graph, BUF3_OFFSET, &buf3_off[rank]);
93 MPIX_Graph_set_param(&graph, BUF3_PTR, &buf3_ptr[rank]);
94 MPIX_Graph_set_param(&graph, BUF3_COUNT, &buf3_count[rank]);
95 MPIX_Graph_set_param(&graph, BUF3_TYPE, &buf3_type[rank]);
96
97 // Create deferred operation token for parameterized graph
98 MPIX_Graph_def(graph, info, &token);
99
100 // Create an [inactive] request for the graph token
101 MPIX_Execute_init(token, info, MPIX_OFFLOAD_NITRO, NULL, &req);
102
103 // Release token for parameterized graph
104 MPIX-Token_free(&token);
105
106 // Start execution of graph
107 MPI_Start(&req);
108
109 // Indicate that the I/O has started
110 checkpoint_io_started = TRUE;
111
112 <swap data buffers for next timestep>
113 }
114 }
115
116 // Free graph
117 MPIX_Graph_free(&graph);

```

Listing 6. – Asynchronously write checkpoint file with collective I/O

7.5. Use Case #5: Data Dependent Asynchronous Communication

This [example code](#) demonstrates conditionally executing communication operations, depending on data values across the application.

This use case creates a reusable parameterized graph that is invoked to overlap with computation. It uses an asynchronous data dependency to control execution of a sub-graph of asynchronous communication, all of which can be offloaded or executed in the background.

Note that executing this graph is not the same as passing parameters to a subroutine in the application that invokes asynchronous operations. Such a subroutine would need to block on the allreduce operations in order

to compare the results. However, the ‘IfElse’ operation used in the graph is asynchronously executed in the same manner as the asynchronous communication operations and is managed by MPI, not the application. A similar result to the graph execution could be created by running such a subroutine in an application thread, but that would prevent the MPI implementation from seeing the ‘big picture’ of all the planned operations and reduce opportunities for optimizing the operations in the graph.

```

1  MPIX_Graph graph, subgraph;
2  MPIX-Token token;
3  MPI-Request req;
4  enum {GRAPH_COMM, VAR_A, VAR_B, SEND_BUF, RECV_BUF, // Graph parameter
   ↪ variable IDs for conditional communication
5     COMP_VAL}; // Graph parameter variable IDs for comparison
6
7  // Create main graph
8  MPIX_Graph_create(session, info, &graph);
9
10 // Create value parameter graph variables for graph communications
11 MPIX_Graph_var_create(&graph, GRAPH_COMM, MPIX_GRAPH_VAR_VAL_PARAM,
   ↪ MPI_LOCAL_TYPE_COMM);
12
13 // Create value parameter graph variables for controlling data exchange
14 MPIX_Graph_var_create(&graph, VAR_A, MPIX_GRAPH_VAR_VAL_PARAM, MPI_INT);
15 MPIX_Graph_var_create(&graph, VAR_B, MPIX_GRAPH_VAR_VAL_PARAM, MPI_INT);
16
17 // Create a value parameter graph variables for the send & recv buffers
18 MPIX_Graph_var_create(&graph, SEND_BUF, MPIX_GRAPH_VAR_VAL_PARAM, MPI_AINT);
19 MPIX_Graph_var_create(&graph, RECV_BUF, MPIX_GRAPH_VAR_VAL_PARAM, MPI_AINT);
20
21 // Concurrently perform all reduce operations on 'A' and 'B'
22 unsigned allreduce_op_ids[2];
23 MPIX_Allreduce_def(MPI_IN_PLACE, VAR_A, 1, MPI_INT, MPI_MAX, GRAPH_COMM,
   ↪ &token);
24 MPIX_Graph_add(&graph, &token, NULL, &allreduce_op_ids[0]);
25
26 MPIX_Allreduce_def(MPI_IN_PLACE, VAR_B, 1, MPI_INT, MPI_MAX, GRAPH_COMM,
   ↪ &token);
27 MPIX_Graph_add(&graph, &token, NULL, &allreduce_op_ids[1]);
28
29 // Create a 'join' operation ID, for all writes
30 unsigned join_reduces_op_id = 0;
31 MPIX_Graph_join(&graph, 2, allreduce_op_ids, &join_reduces_op_id);
32
33 // Create a reference parameter graph variable for comparison result
34 MPIX_Graph_var_create(&graph, COMP_VAL, MPIX_GRAPH_VAR_REF_PARAM, MPI_INT);
35
36 // Compare 'A' and 'B' values from allreduce operations
37 unsigned compare_op_id = 0;
38 MPIX_Graph_var_reduce(&graph, VAR_A, VAR_B, MPIX_LESS_THAN_OR_EQUAL, COMP_VAL,
   ↪ &compare_op_id, join_reduces_op_id);
39
40 // Create sub-graph, for if/else
41 MPIX_Graph_create(session, info, &subgraph);
42
43 // Sendrecv operation for if/else
44 MPIX_Sendrecv_def(SEND_BUF, 2, MPI_FLOAT, (rank + 1) % num_ranks, 0, RECV_BUF,
   ↪ 2, MPI_FLOAT, (num_ranks + (rank - 1)) % num_ranks, 0, GRAPH_COMM,
   ↪ MPI_STATUS_IGNORE);
45 MPIX_Graph_add(&subgraph, &token, NULL, MPIX_DEP_NONE);

```

```

46
47 // Execute subgraph if 'A' <= 'B'
48 MPIX_Graph_IfElse(&graph, COMP_VAL, &subgraph, MPIX_GRAPH_NULL, NULL,
49 ↪ compare_op_id);
50
51 // Free subgraph
52 MPIX_Graph_free(&subgraph);
53
54 // Set communicator for all graph executions
55 MPI_Comm graph_comm = <communicator>;
56 MPIX_Graph_set_param(&graph, GRAPH_COMM, &graph_comm);
57
58 int comp_val = 0;
59 while(<more timesteps>) {
60 // Check results of previous timestep's graph execution
61 if (<timestep > 0>) {
62 // Make certain that graph has completed
63 MPI_Wait(&req, &status);
64
65 // Check if graph's condition was met
66 if(comp_val) {
67 <use value in recv_buf>
68 comp_val = 0;
69 }
70 }
71
72 // Application variables on each rank, for graph
73 int a = compute_a(timestep);
74 int b = compute_b(timestep);
75 float sendbuf[2], recvbuf[2];
76
77 <init sendbuf & recvbuf>
78
79 // Set parameters for graph on each rank and execute it
80 MPIX_Graph_set_param(&graph, VAR_A, &a);
81 MPIX_Graph_set_param(&graph, VAR_B, &b);
82 MPIX_Graph_set_param(&graph, SEND_BUF, &sendbuf);
83 MPIX_Graph_set_param(&graph, RECV_BUF, &recvbuf);
84 MPIX_Graph_set_param(&graph, COMP_VAL, &comp_val);
85
86 // Create deferred operation token for parameterized graph
87 MPIX_Graph_def(graph, info, &token);
88
89 // Create an [inactive] request for the graph token
90 MPIX_Execute_init(token, info, MPIX_OFFLOAD_PTHREAD, NULL, &req);
91
92 // Release token for parameterized graph
93 MPIX-Token_free(&token);
94
95 // Start execution of graph
96 MPI_Start(&req);
97
98 <compute>
99 }
100
101 // Check results of last timestep's graph execution
102 if (<timestep > 0 >) {

```

```
103 // Make certain that graph has completed
104 MPI_Wait(&req, &status);
105
106 // Check if graph's condition was met
107 if(comp_val)
108     <use value in recv_buf>
109 }
110
111 // Free graph
112 MPIX_Graph_free(&graph);
```

Listing 7. – Data Dependent Asynchronous Operations

A. Appendix: An Abstract Data Movement Machine

If there was a chip that provided “assembly” instructions that could be used to implement the MPI operations described above, what would those instructions be?

A.1. Data Containers

A “container” is an abstraction around a particular sequence of bytes and provides the base object for data operations.

A container is one of four things:

- *Local* memory – memory at this MPI process, which could be attached to a CPU, GPU, FPGA, etc.
- *Remote* memory – memory at another MPI process, which could (also) be attached to a CPU, GPU, FPG, etc.
- *Abstract machine* memory – memory in one of the abstract ‘data movement machines’ being described here. (More details on this type of memory are outlined below, and are called ‘variables’ there)
- *Storage* memory – a file on an SSD, HDD, tape, a cloud object, an RMA window, etc. This form of memory is *passive* – it doesn’t have an MPI process managing it.

The virtual machine must provide mechanisms to: create new, access existing, release access to, and delete these types of containers. (i.e. ‘create’, ‘open’, ‘close’, and ‘delete’ container operations)

Comment [QK4]: Elaborate on Ryan and my discussion today about these types of memory / containers being associated with a process space, except for storage, etc.

A.2. Container Operations

The fundamental operations for containers in the virtual machine are abstractions for distributed “copy data”, “reduce”, “compare”, and “synchronize”. These operations provide the core capabilities necessary to implement MPI operations, as well as other similar data-oriented frameworks, and are described below.

A.2.1. Abstract Copy Operations

The “abstract copy” (“ACOPY”) operation copies bytes from a list of 1+ source locations to a list of 1+ destination locations. A location is defined as a tuple of {container, offset, length}. The offset and length for a location tuple describes a byte range in a container.

The ACOPY operation copies bytes from source containers to destination containers, in the order given in each location list. This copy operation does not require that the lengths of each byte range be equal in the lists, only that the total number of bytes described by the source location list is equal to the total number of bytes described by the destination location list. The ACOPY operation also does not require that byte ranges are non-overlapping, non-repeating, or sorted in any particular order, although that might be imposed by a higher-level abstraction layer (like MPI!).

This ACOPY operation allows for arbitrary scatter-gather operations between any type of data containers. Again, higher-level layers (like MPI) may impose various constraints on the location lists, in order to improve performance or other aspects of execution.

Need a better word than ‘abstract machine’. Maybe “virtual”, “transient”, “temporary”, or “in-tragraph” instead?

A.2.2. Abstract Reduce Operations

Distributed arithmetic operations on data container locations are analogous to the existing MPI ‘reduce’ operation⁹. As such, any type of container (local, remote, abstract machine, or storage) may be used as inputs to or the output from an abstract reduce (“AREDUCE”) operation:

$$output = reduce(op, input_0, input_1, \dots)$$

Valid mathematical operations for the *op* in an AREDUCE operation must be both commutative and associative. Examples include addition, multiplication, logical AND and OR, binary AND and OR, union and intersection on sets, greatest common divisor, least common multiple, minimum, and maximum.

The AREDUCE operation allows for arbitrary reduce operations between locations in any type of data container. Again, higher-level layers (like MPI) may impose constraints on the location of inputs, in order to improve performance or other aspects of execution.

A.2.3. Abstract Comparison Operations

Comparison operations on container locations don’t have a direct analog to an existing MPI operation, but are closest to a distributed consensus operation¹⁰. The abstract comparison (“ACOMPARE”) operation takes a single primary operand (“*primary*”), a comparison operation (“*op*”), a set of *n* secondary operands (“*secondary*₀..*secondary*_{*n*-1}”), and a threshold value (“*threshold*”) as inputs and produces a boolean output value that is true if more than *threshold* of the (*primary*, *op*, *secondary*) comparisons are true and false if the *threshold* value is not reached.

$$output = compare(threshold, primary, op, secondary_0, secondary_1, \dots)$$

Valid comparison operations for the *op* in an ACOMPARE operation must produce a binary true or false value for comparing the *primary* operand to each of the *secondary* operands. Examples include ==, ≠, <, ≤, >, ≥, and possibly others.

This ACOMPARE operation allows for arbitrary comparison operations between locations in any type of memory container. Again, higher-level layers (like MPI) may impose constraints on the location of inputs, in order to improve performance or other aspects of execution.

A.2.4. Abstract Synchronization Operations

Synchronization operations on container locations are analogous to the existing MPI ‘barrier’ operation or a synchronous send+receive pair. Abstract synchronization (“ASYNCHRONIZE”) operations take two forms:

- *Data* synchronization – Synchronize 2+ processes to be certain that some/all of the processes have the same data in a container location.¹¹
- *Execution* synchronization – Synchronize 2+ processes to be certain that some/all of the processes have/haven’t reached the same point in executing an application.¹²

A.2.5. Constants

Constant (‘literal’) values, such as ‘3’ or ‘2.1’, can be used in the following operations:

- *ACOPY* – a constant may be used as a source location.

⁹Note that some existing MPI operations or reduce operators must be built up from multiple ACOPY or AREDUCE operations. MPI all-reduce and scan operations and the MPI_MINLOC and MPI_MAXLOC reduce operations are examples of these.

¹⁰[https://en.wikipedia.org/wiki/Consensus_\(computer_science\)](https://en.wikipedia.org/wiki/Consensus_(computer_science))

¹¹A pair of synchronous send+receive MPI operations is an example of data synchronization.

¹²The MPI ‘barrier’ operation is an example of execution synchronization.

Need better wording here!

Comment [QK5]: Interesting variations on barrier could include: a “threshold barrier” operations (“at least 6 of *n*_{*i*} processes must have reached this barrier

- *AREDUCE* – a constant may be used as an input.
- *ACOMPARE* – a constant may be used as either a primary or secondary value.

Values for constants used in graph operations are captured during graph definition.

A.3. Graphs

A “virtual task graph” with an execution unit that understands dependencies between 0+ “parent” tasks and 0+ “child” tasks, along with “control flow” operations, and “virtual” memory/registers for the graph execution unit best describes the next components of the “abstract data movement machine”.

A.3.1. Graph Task Execution Unit

The “graph task execution unit” understands the following types of “task operations” as nodes in the graph:

- *Container* operations – *ACOPY*, *AREDUCE*, *ACOMPARE* and *ASYNCHRONIZE* operations on container locations, as well as container create/access/release/delete operations.
- *User* operation – an operation that allows an application callback as an operation in the graph.
- *Subgraph* operation – an entire graph as an “operation” in the current graph.

A.3.2. Graph Dependency Unit

The “graph dependency unit” manages dependencies between graph tasks, i.e. the edges in the graph. There are several types of dependencies, but they fundamentally are a mechanism for ordering execution of graph operations.

Task and Data Dependencies There are two types of dependencies for graph operations: task and data. Tasks dependencies are explicitly created by an application, by specifying parent operations as ‘input’ dependencies for an operation, all of which must complete before the operation can be executed. Data dependencies are indirectly created by an application, by using data or values produced by an earlier graph operation as an input value to another operation.

Task operations (A.3.1) each have a single ‘input’ task dependency (i.e. one parent) and a single ‘output’ task dependency (i.e. one child), as shown in Fig. 1. When their input dependency has been fulfilled (i.e. their parent task has completed), the task may be scheduled for execution. Correspondingly, when the task completes, its output dependency is fulfilled and its child task may be scheduled for execution. More complicated dependencies between graph tasks can be constructed with “organizational” dependencies, described below.

Each task operation may also have 0+ input data dependencies and 0+ output data dependencies. Data dependencies are implicitly created between operations when data or values produced by one operation (i.e. a return value or an “out” parameter) is consumed by another operation (i.e. used as an “in” or “in/out” parameter). There are two kinds of data dependencies:

- *Strong* data dependencies – provide or create access to a container (i.e. a file handle, memory buffer, etc) and can be represented as a “future” value.
- *Weak* data dependencies – access or modification of the contents of container (i.e. *ACOPY*, *AREDUCE*, *ACOMPARE* or *ASYNCHRONIZE* operations on locations within a container).

Strong data dependencies can be reliably tracked between task operations, therefore a dependency between the operation that generates a future value and any operations that use that value doesn’t need to be explicitly defined by an application. However, weak dependencies can have ambiguous ordering that only an application

Disallow ‘self’ as valid “sub-graph op” for now, but consider a mechanism for recursion in the future.

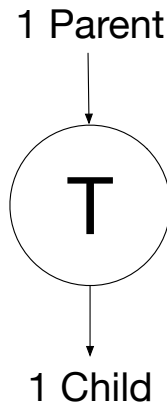


Figure 1 – Task dependencies

programmer understands. Operations that rely on ordered access to or modification of the contents of a container must therefore explicitly define that ordering with task or organizational dependencies between the operations.

Organizational Dependencies “Organizational” dependencies don’t actually perform any work, they just provide “connections” in the graph to indicate and manage ordering of task operations. The graph dependency unit understands the following types of organizational dependencies:

- *Unconditional* dependencies – Unconditional dependencies connect 0+ input dependencies (“parents”) to 0+ output dependencies (“children”) without regard to any other state for the graph. All of the input dependencies for an unconditional organizational operation (if any are defined) must complete before any of the output dependencies may start execution. The general form of an unconditional dependency is represented in Fig. 2.

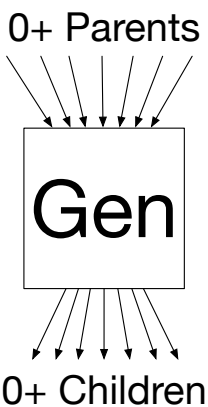


Figure 2 – General form of unconditional organizational graph dependencies

Common specialized forms of unconditional dependencies are root, fan-in, fan-out, and leaf, as shown in Fig. 3.

- *Conditional* dependencies – Conditional dependencies connect 0+ input dependencies (“parents”) to 0+ output dependencies (“children”) and have a condition that must be met to fulfill the output dependency. Two forms of conditional organizational dependencies are:

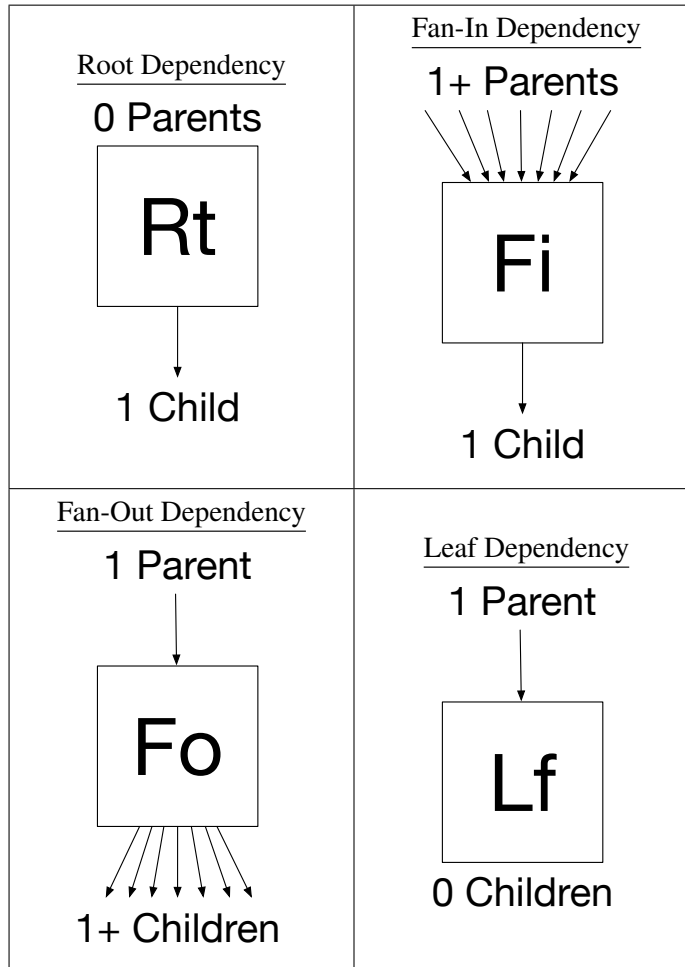
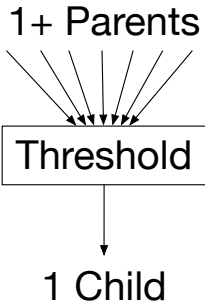


Figure 3 – Common forms of unconditional organizational graph dependencies

- *Consensus* dependencies – Consensus dependencies are fulfilled when a (user-defined) threshold of their input dependencies are fulfilled, providing a “some but not all” or “enough” form of dependency, represented in Fig. 4.



Comment [QK6]: Should the threshold be constant or allowed to depend on a graph variable?

Figure 4 – Threshold organizational graph dependency

- *Branch* dependencies – Branch dependencies connect a different output dependency to an input dependency, depending on the state of a graph variable, as shown in Fig. 5.

Note that a graph variable is examined within the branch dependency, which could have weak data dependencies on other concurrently executing operations in the graph that create race conditions when the graph is executed. Application developers can use other organizational dependencies as preludes to the branch dependency to eliminate such race conditions.

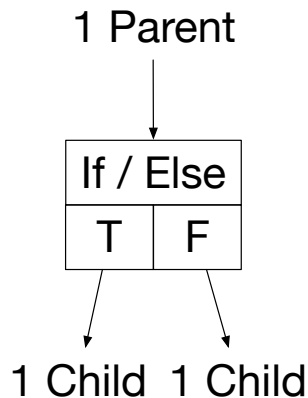


Figure 5 – Branch organizational graph dependency

Graph Control Flow “Control flow” operations also don’t perform any work, they just indicate the next task to execute, as if the destination of the control flow operation was a child dependency of another operation. However, control flow operations don’t create “full” dependencies in the graph – they are child dependencies of their parent task, but not parent dependencies of their destination task. Control flow operations act similar to a ‘goto’ operation in a traditional programming language and are shortened to the ‘AGOTO’ mnemonic in this document and drawn in diagrams as shown in Fig. 6.

The destination of a control flow operation must be a task where all of the destination task’s descendants (children, grandchildren, etc) **only** have parent dependencies that can be traced to the destination task. For example, in Fig. 7, tasks ‘A’, ‘B’, and ‘C’ in the top row are not valid destination tasks, but task ‘D’ in the middle row is, as are tasks ‘E’, ‘F’, and ‘G’ in the bottom two rows.

Infinite loops are possible with AGOTO operations, but can be interrupted by a ‘cancel’ operation on the graph when it is executing.

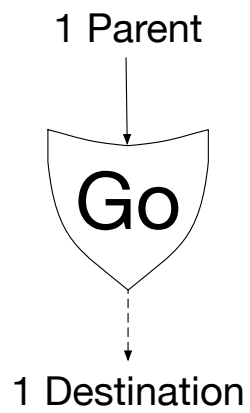


Figure 6 – Control flow operation

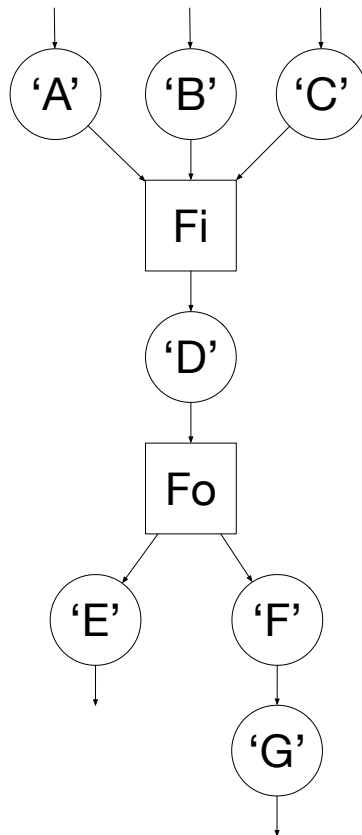


Figure 7 – Example graph for control flow destinations

A.3.3. Graph Construction Rules

Any combination of tasks and organizational dependencies is allowed. Tasks can be connected directly together, as shown in Fig. 8, organizational dependencies can be connected directly together, as shown in Fig. 9, and they can be combined, as shown in Fig. 10.

Tasks, organizational dependencies, and control flow operations can be combined to make data-dependent loop structures in a graph, as shown in Fig. 11.

Tasks, organizational dependencies, subgraphs, and control flow operations can be combined to make complex nested graphs, as shown in Fig. 12.

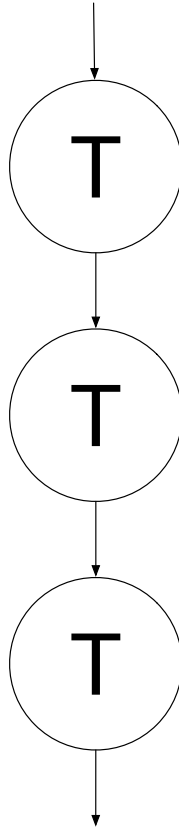


Figure 8 – Connecting multiple task dependencies

A.3.4. Graph Variables

“Graph variables” fill the role of registers and memory for the graph execution virtual machine.

There are three types of graph variables:

- *Value parameter* variable – a variable that is initialized with a value when the graph starts execution.
- *Reference parameter* variable – a variable that is an alias for another container location, set when the graph starts execution.
- *Local* variable – a variable that is initialized to “all zeroes” when the graph starts execution, and is not aliased to a container location.

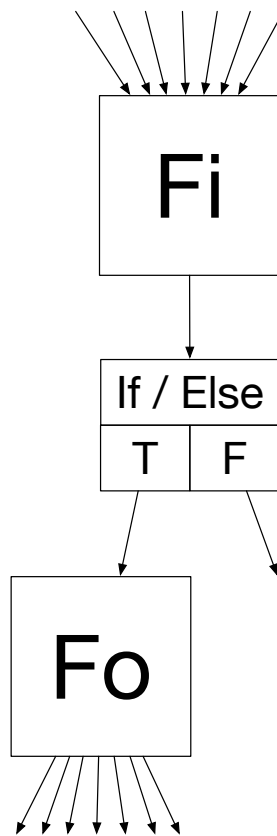


Figure 9 – Connecting multiple organizational dependencies

Variables may be declared at any point during a graph’s definition, but assignments for all value or reference parameter variables defined anywhere in the graph must be provided at the time when the graph starts execution.

The following operations are defined for graph variables:

- *Definition* – declaring a graph variable. Variables are not explicitly deleted, but any resources associated with them are released when a graph terminates.
- *ACOPY* – a graph variable is a valid container location as the source or destination for ACOPY operations.
- *AREDUCE* – a graph variable is valid for both input and output locations for AREDUCE operations.
- *ACOMPARE* – a graph variable is valid for both primary and secondary locations for ACOMPARE operations.
- *ASYNCHRONIZE* – a graph variable is valid for data synchronization ASYNCHRONIZE operations.

‘Set’ or ‘get’ operations on graph variables are ACOPY operations between a graph virtual machine container and another location. If the value of a value parameter or local variable is to be retained after its graph terminates, the value must be ACOPY’d from the graph variable before graph termination.

A.4. Streams

A “virtual task stream” with an execution unit that understands asynchronous execution of an ordered sequence of operations, and “virtual” memory/registers for the stream execution unit best describes the next component of the “abstract data movement machine”.

The MPI API will need to provide a way to provide values/assign references to parameters when a graph is executed.

See if it’s possible to support “string” types.

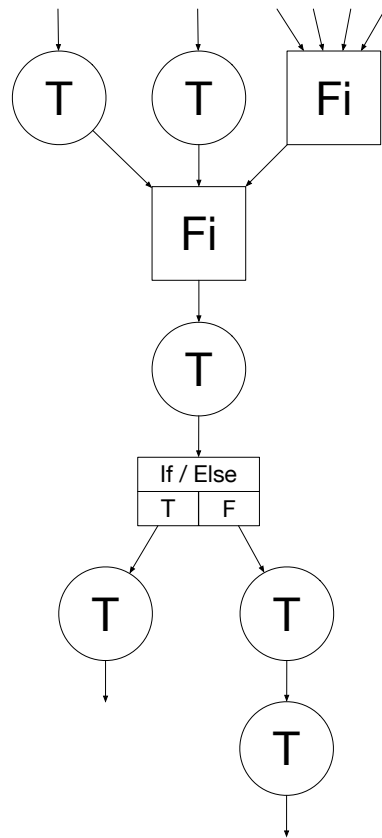


Figure 10 – Combining task and organizational dependencies

A.4.1. Stream Task Execution Unit

Streams operate as FIFO queues to execute operations, similar to a conveyor belt: the first operation in a stream executes to completion and is removed from the stream, then the following operation can be scheduled for execution. The “stream task execution unit” understands the following types of “task operations” as ‘nodes’ in an ordered sequence:

- *Container* operations – ACOPIY, AREDUCE, ACOMPARE and ASYNCHRONIZE operations on container locations, as well as container create/access/release/delete operations.
- *User* operation – an operation that allows an application callback as an operation in the stream.
- *Subgraph* operation – an entire graph as an “operation” in the current stream.

A.4.2. Stream Construction Rules

Any combination of tasks and sub-graphs is allowed. Tasks can be inserted as operations in the stream, as shown in Fig. 13 and sub-graphs can be inserted as stream operations as well, as shown in Fig. 14.

A.4.3. Stream Variables

“Stream variables” fill the role of registers and memory for the stream execution virtual machine and are identical to the graph variables defined in section A.3.4 except that they are defined for a stream, not a graph.

Disallow nesting ‘sub-streams’ for now, but consider in the future.

Comment [QK7]: Not certain it makes sense to have variables for streams.

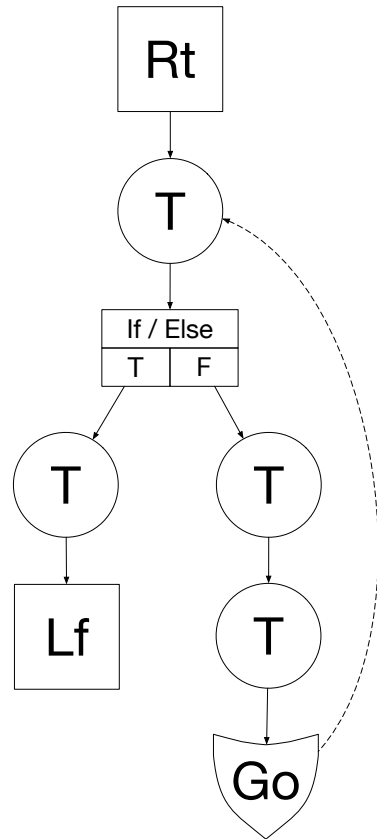


Figure 11 – Combining tasks, organizational dependencies, and control flow operations to create loop in graph

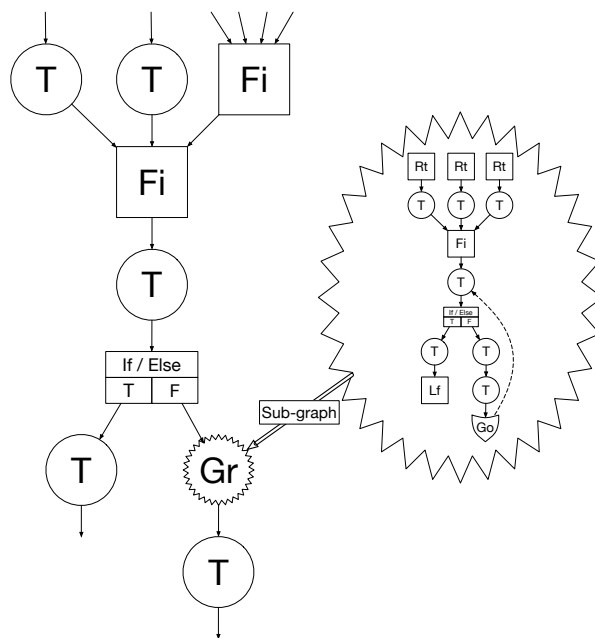


Figure 12 – Combining tasks, organizational dependencies, and sub-graphs

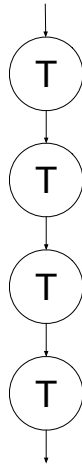


Figure 13 – Tasks in a stream

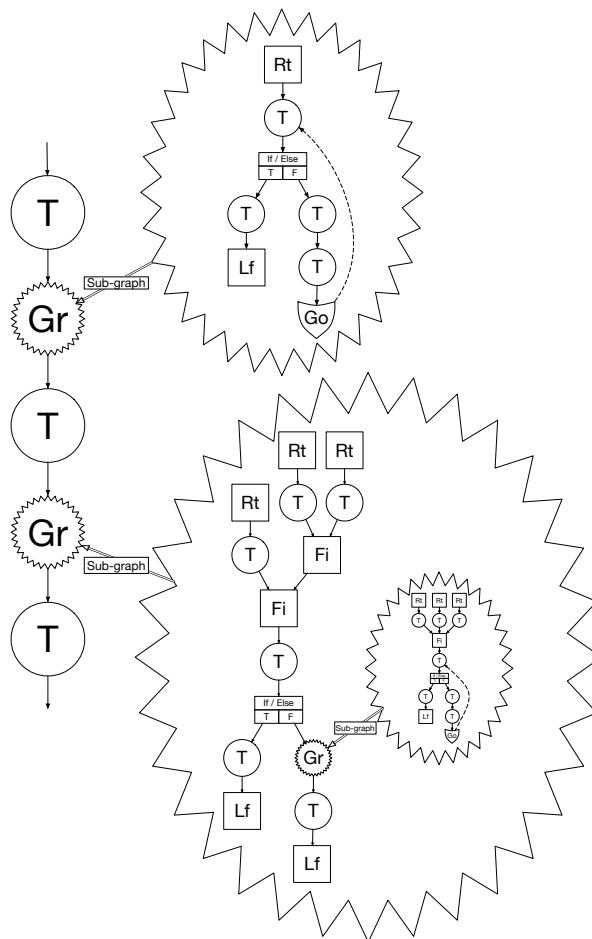


Figure 14 – Tasks and sub-graphs in a stream

A.5. Scopes

B. Appendix: Comparison to proposed Continuations and MPIX_Stream extensions

At first glance, the capabilities described in this document appear similar to the proposed “Continuations” and “MPIX_Stream” extensions..

Continuations allow a user callback to be attached to a nonblocking MPI operation. Then the continuation callback is invoked when the nonblocking operation completes. Although the capabilities described in this document do include an asynchronous user callback, they are not directly comparable with the user callbacks described in this document. A continuation callback could be executed as part of the progress engine, unless they are explicitly excluded and require that the MPI implementation invoke it when the test or wait operation is performed on a request. However, the user callbacks described in this document are standalone asynchronous operations, fully integrated into the flow of other asynchronous communication and I/O operations.

The MPIX_Stream extensions overlaps some of the ideas presented here, in the area of scheduling operations. However, they are mainly focused on enabling communication with compute endpoints on GPUs and similar hardware components. The asynchronous operations described in this document are essentially orthogonal to them, and both extensions can be integrated into the MPI standard without conflict.

Add
bibrefs
for these

C. Appendix: What this proposal is *not*

Although the capabilities described here may operate best with an MPI implementation that executes them with a background thread or on dedicated hardware, that is not a requirement. The capabilities described here will operate correctly with the typical ‘weak’ progress model for MPI, just not as efficiently as if they were executing with ‘strong’ progress. If threads or hardware offload are not used, the ‘asynchronous’ operations described here may be implemented in the same way as the “nonblocking” operations are currently.

The operations described here are also not designed to provide the capabilities of a computationally-focused task execution framework such as CUDA, etc. Although a call to asynchronously execute an application-provided callback function is provided, it is intended as an “escape hatch” mechanism for short-running operations.

Likewise, the asynchronous operations described here are primarily data movement operations that have minimal performance impact on an application’s execution time. A high-quality implementation will schedule asynchronous data movement operations quickly and give up the CPU as soon as possible.

D. Appendix: Details for Asynchronous Operations Option #2

Add details from lab notebook here.