

The Communication Semantics of the Message Passing Interface

*Robert Palmer, Ganesh Gopalakrishnan, and
Robert M. Kirby*

UUCS-06-012

School of Computing
University of Utah
Salt Lake City, UT 84112 USA

31 October 2006

Abstract

The Message Passing Interface (MPI) standard is a natural language document that describes a software library for interprocess communication. Automatic reasoning about the reactive nature of programs communicating via MPI libraries is not possible without also analyzing the library being used. Many distributed programs that use MPI are relatively brief compared to the libraries that implement MPI. A formal specification of the communication semantics of the MPI standard (i) enables modular automatic reasoning of MPI based parallel programs independent of the library implementation, (ii) provides a mathematically precise declaration of the natural language intent of the MPI specification, (iii) enables mathematical reasoning about libraries that implement the standard, and (iv) allows for reasoning about the standard itself. We have created such a specification of the point to point operations and present it in this report. We also discuss some preliminary efforts to accomplish (i) above.

Disclaimer: While the semantics have been proof-read once, the actual semantics document is continually evolving. We are developing a tool—MPIC—that can be used to verify programs against this semantic specification. When the MPIC tool is released there will be a new version. The MPIC tool will also have an accompanying technical report.

Although every effort has been made to correctly model the intent of the MPI 1.1 specification, we make no claim regarding the correctness of the model contained herein. Please notify the authors if a discrepancy is found.

The Communication Semantics of the Message Passing Interface *

Robert Palmer Ganesh Gopalakrishnan Robert M. Kirby

31 October 2006

Abstract

The Message Passing Interface (MPI) standard is a natural language document that describes a software library for interprocess communication. Automatic reasoning about the reactive nature of programs communicating via MPI libraries is not possible without also analyzing the library being used. Many distributed programs that use MPI are relatively brief compared to the libraries that implement MPI. A formal specification of the communication semantics of the MPI standard (i) enables modular automatic reasoning of MPI based parallel programs independent of the library implementation, (ii) provides a mathematically precise declaration of the natural language intent of the MPI specification, (iii) enables mathematical reasoning about libraries that implement the standard, and (iv) allows for reasoning about the standard itself. We have created such a specification of the point to point operations and present it in this report. We also discuss some preliminary efforts to accomplish (i) above.

1 Introduction

Standards documents are one of the powerful tools for developing portable, reusable, and correct implementations of complex systems. In almost all cases, they are initially created as semi-formal documents, often containing gaping holes and potentially ambiguous statements. Over time, thanks to the experience gained from

*Supported in part by NSF award CNS-0509379 and a grant from Microsoft Corporation.

the widespread use of systems built according to the standard, they evolve to become much more rigorous and coherent. Yet, without a *mathematical* (formal) description, they still leave much room for misinterpretation – often unfortunately coinciding with the increased scale of design and deployment of systems built according to the standard.

The IEEE Floating Point standard [6] is a resounding success story in this area. It was initially conceived as a standard that helped minimize the danger of non-portable floating point implementations. As a fortunate side effect of the infamous Intel Pentium division bug, it now has incarnation in various higher order logic specifications (e.g., [5]), and routinely finds applications in *formal proofs* of modern microprocessor floating point hardware circuits. We strongly believe that the MPI communication standard – one of the most widely used in high performance computing – has the vast potential of being solidified in a similar fashion.

MPI is already a success story in the area of software library standardization, in that a collection of primitives that support message passing based communication for high performance computing has been widely adopted. Unfortunately, the MPI standard [10] uses natural language descriptions, as well as examples to communicate definitions, semantics, and other important details. Experience shows that this can lead to errors that result from unstated assumptions, ambiguities, and unclear causal dependencies. Some of the recent additions to MPI, such as *one-sided communication constructs*, are so tricky to understand that even simple algorithms using them have been shown to be incorrect (e.g., [9, 12]). These errors can be progressively eliminated by relying on *formal* (mathematical) descriptions of MPI, and employing modern formal verification techniques such as model checking [3].

1.1 Related Work

Significant inroads have been made in formalizing the MPI standard. The earliest work we are aware of is that of Georgelin et. al. [4] where the authors create a LOTOS description of MPI_BSEND, MPI_SSEND, MPI_RSEND, and MPI_RECV along with the collective operations MPI_BROADCAST, MPI_GATHER, and MPI_SCATTER. The MPI_ANY_SOURCE and MPI_ANY_TAG wild-cards are modeled. The above operations are modeled using the channel primitives of the LOTOS description language. They then apply the model to verification of some MPI programs.

More recent work by Siegel and Avrunin includes:

- In [13, 15] the authors create a mathematically precise model of `MPI_SEND`, `MPI_RECV`, `MPI_SENDRECV`, and `MPI_BARRIER`. The `MPI_ANY_SOURCE` and `MPI_ANY_TAG` wild-cards are expressly disallowed due to the additional non-determinism that are introduced with their use. The process universe is a fully connected graph where edges in the graph are a pair of FIFO channels. Synchronous communications are modeled as an interleaved execution of two processes where a send is followed immediately by the corresponding receive in an execution trace. A number of theorems and their proofs are presented with regards to synchronous communication in the proposed model.
- In [14] the authors model a simple MPI based 2D diffusion simulation and verify the model using both SPIN and INCA. Models of `MPI_SEND`, `MPI_RECV`, `MPI_BARRIER`, and `MPI_SENDRECV` are used in connection with the diffusion simulation. Some of the results from [13] are also presented.
- In [16] the authors verify the output of a distributed numerical program using a model checker and a sequential version of the same program.

Other work in modeling MPI related programs includes [12] where model checking is applied to a program using the one-sided locking routines of MPI 2.

We have also modeled `MPI_SEND`, `MPI_RECV` and `MPI_BARRIER` for use with both SPIN [2] and Zing [11].

1.2 Motivations

With several existing models of MPI one may naturally ask, why have another? To answer this question, consider the following points.

1. There are 35 operations related to point to point communication described in chapter 3 of the MPI 1.1 standard [10]. The most aggressive modeling effort we are aware of contains four (4) of these operations.
2. There is no tool based reasoning support for any mathematically rigorous model of MPI. The only rigorous model we are aware of is the one described in [13]. Models created in languages such as SPIN and Zing are

dependent upon the model checker implementation for a formal description of the language semantics.¹

3. Many of the errors that are seen in MPI programs are derived not in the use of the blocking sends, rather in the translation from the use of these simple send primitives to the more aggressive counterparts in the ongoing effort to optimize. None of the existing models have representations of these more aggressive operations.
4. A mathematically precise representation of a larger subset of the MPI operations is necessary to create an industrially useful tool for reasoning about programs that communicate using MPI libraries.
5. To serve as a specification, the MPI standard does not mandate implementation details beyond the function signatures and the existence of some symbols. No mention is made of how messages are to be transmitted from one process to another. To make a sufficiently complete model of any MPI program, these details must be filled in. Existing models make no distinction between what is specified in the standard and what is added to support tool based reasoning.

A mathematically precise specification the MPI standard can serve, not only to reason about programs that employ the MPI libraries, it can be used to reason about the various MPI library implementations and the standard itself.

1.3 A Formal Model

Our formal model of the MPI specification is expressed using the Temporal Logic of Actions [7, 1]. TLA is a formal logic containing standard ZF set theory, an action operator that induces a transition relation, and some limited temporal logic. It's semantics are well understood by mathematicians and computer scientists, independent of any verification tool. Implementation details can be abstracted using set theoretic operations in combination with the action operator.

A program that uses MPI can be specified as a formula in TLA representing a distributed computation. The MPI operations are modeled as operators on variables in the formula. Comments accompany individual logical clauses referencing

¹While this may not necessarily be the case for LOTOS or INCA, any language developed as input for a verification tool that is undergoing active development could deviate from the previously published semantics.

the page and line number of the MPI standard that requires the given clause where possible.

The action operator of TLA induces a transition relation on logical formulas. Variables that are *primed* (i.e., foo') indicate the value of that variable in the next state of the system. Every transition specifies the values of all variables in the next state. Transitions are total functions from valuations of variables to valuations of variables. New operators can be defined for any finite arity as a combination of the existing operators, user defined variables and constants, and parameters to the operator.

1.3.1 Transition Granularity

Using a TLA operator to represent an MPI operation implies that such operations will require exactly one transition to complete. Our model assumes that only one MPI operation will be applied in a given transition. Although 28 of the operations related to point to point communication are modeled in this way, we could see no way to model some of the point to point operations using only one transition. In particular, `MPI_SEND`, `MPI_SSEND`, `MPI_BSEND`, `MPI_RSEND`, `MPI_RECV`, `MPI_SENDRECV`, and `MPI_SENDRECV_REPLACE` could not be modeled in a single transition. The reason for this is quite simple: Each of these operations writes some variable *and then* waits for some other variable to be written by another process. The specification of which explicitly requires at least two (2) transitions. We will demonstrate in this paper how to model the remaining operations as a sequential composition of the provided operators.

To model operations requiring more than one transition we adopt the same convention described in [13] noting that the sequential composition of `MPI_ISEND` and `MPI_WAIT` on a single process is semantically equivalent to `MPI_SEND` and that a transition sequence with a `MPI_SEND` followed immediately by the corresponding `MPI_RECV` can be considered synchronous. Although it is possible to apply sequential and operator composition in such a way that all 35 operations can be derived using a minimal subset of the MPI operations, to facilitate cross referencing the MPI standard, all 28 single transition operations in our model are modeled separately.

1.4 What is and is not modeled

Before diving into all the details, it is important to note that not everything in MPI 1.1 is present in the model. In particular, the following are either not present, are

limited in their current modeling, or are currently only placeholders:

- Data. Data, such as arrays of floating point values, objects, etc., could be modeled using TLA. It is, however, not necessary in most cases to retain the actual data of the distributed simulation to verify reactive properties of nodes participating in the distributed simulation. Therefore we allow a placeholder for data such that it can be included when necessary.
- Data manipulation operations. There are many operations specified by MPI to pack and manipulate data. These are not currently modeled, but could be if there were sufficient interest.
- Operations on communicators and topologies. These are modeled to a limited extent to enable point to point communications on intra-communicators. We currently model the operations shown in Figure 2 in addition to the point to point operations of chapter 3 of MPI 1.1 shown in Figure 1. Such operations on communicators and topologies should be a strait forward extension of this work.
- Implementation details. To the greatest extent possible we have avoided asserting implementation details that might constrain an implementation. One obvious ramification of this omission is that modeling return codes of MPI operations is completely eliminated (see Pg 11 of [10]).
- Transient buffering of messages created by the standard mode send (MPI_SEND, MPI_ISEND, MPI_SEND_INIT). We require the system to either eventually buffer these send requests or to never buffer them. It is not clear at the time of this writing how to model the situation where a buffer may be available for some but not all of the program execution.

Our model includes the point to point operations shown in figure 1 as single transition TLA operators. The argument order and meaning is as specified in the MPI standard for each operator, adding the pid of the process that is applying the operator as the last argument.

2 Conventions

In TLA, the whitespace in the document is significant. Sequences of logical conjuncts can become quite large and are therefore formatted as bulleted lists, the bullet being the logical and (\wedge) or the logical or (\vee) operator.

MPI_GET_COUNT	MPI_REQUEST_FREE	MPI_TEST_CANCELED
MPI_BUFFER_ATTACH	MPI_WAITANY	MPI_SEND_INIT
MPI_BUFFER_DETACH	MPI_TESTANY	MPI_BSEND_INIT
MPI_ISEND	MPI_WAITALL	MPI_SSEND_INIT
MPI_IBSEND	MPI_TESTALL	MPI_RSEND_INIT
MPI_ISSEND	MPI_WAIT SOME	MPI_RECV_INIT
MPI_IRSEND	MPI_TEST SOME	MPI_START
MPI_IRecv	MPI_IPROBE	MPI_STARTALL
MPI_WAIT	MPI_PROBE	
MPI_TEST	MPI_CANCEL	

Figure 1: Point to point operations included in the TLA specification.

MPI_BARRIER	MPI_GROUP_SIZE	MPI_GROUP_RANK
MPI_COMM_SIZE	MPI_COMM_RANK	MPI_COMM_COMPARE
MPI_INIT	MPI_FINALIZE	MPI_INITIALIZED
MPI_ABORT		

Figure 2: Additional MPI operations modeled to enable tool based reasoning on MPI based parallel programs.

Our modeling is influenced by the desire to model SPMD style programs in connection with the TLA MPI specification. As such, all program variables are assumed to be arrays of variables (one for each process in the computation).

When specifying the next state of a variable, it is necessary to completely specify that next state. As an example suppose variable $rank \in [0..(N-1) \rightarrow 0..(N-1)]$ is the rank variable declared by a process. When a process calls `MPI_COMM_RANK` the rank variable would be passed into the function. Our model of `MPI_COMM_RANK` requires that the $rank$ be passed to the operator, not $rank[pid]$. As such, we assume that any parameter that might be written by an operator is an array $a : 0..(N-1) \rightarrow \alpha$ where only the i^{th} element (i.e., $a[i]$) is ever accessed by the applying process.

When using the action operator, the value of all variables in the next state must be specified. The specification of the MPI operations includes sometimes many `UNCHANGED` commands which are short hand for $f' = f$. The MPI operators completely specify all of the MPI variables. In addition, those user variables that may be changed by application of the operator are also either updated or marked as `UNCHANGED`.

Comments of the form $n.m$ indicate the corresponding page (n) and line (m) numbers that require the particular feature. All comments are enclosed in shaded regions.

3 Data Structures

This section presents the elements of the model that are introduced to mathematically specify the constructs of MPI. Appendix A contains the entire model. We will refer to it throughout the remainder of the presentation.

3.1 Constants

Symbols that are defined in the MPI standard are modeled as constant values. We have included the subset of symbols that are necessary for the point to point communications on intra-communicators.

In addition to these symbols, we introduce four (4) additional constants. These constant values are useful to (i) make an instantiated program model finite, and (ii) to provide some information that is implicitly available to the MPI system. The additional constants are:

- N . The number of processes in the distributed computation.

- `MAX_COMM`. The maximum number of communicators.
- `TYPES`. The set of strings representing user specified types.
- `TAGS`. The set of integers representing user specified tags.
- `SEND_IS_BUFFERED`. A flag to indicate whether a send can be buffered by the MPI system.

3.2 Variables

Variables are functions. Functions need not have homogeneous domains or ranges. The elements of the domains or ranges need not be numbers (they could be other functions, or strings, or values). The variables in the model are `group`, `communicator`, `requests`, `initialized`, `bufsize`, `message_buffer`, and `collective`.

Functions therefore model data structures such as records, arrays, and sequences. Functions can represent a sequence in that elements can be modified, added, or deleted in the range or domain using the action operator. For example, a sequence $\langle 3, 2, 1 \rangle$ can be modeled by the function

$$s(x) = \begin{cases} 3 & \text{if } x = 1, \\ 2 & \text{if } x = 2, \\ 1 & \text{if } x = 3, \text{ and} \\ \text{undefined} & \text{otherwise} \end{cases}$$

If we wish to append $\langle 4, 5 \rangle$ to this sequence we would let $\langle 3, 2, 1 \rangle \circ \langle 4, 5 \rangle = \langle 3, 2, 1, 4, 5 \rangle$ as

$$s'(x) = \begin{cases} 4 & \text{if } x = 4 \\ 5 & \text{if } x = 5 \\ s(x) & \text{otherwise} \end{cases}$$

As a shorthand we write $x = a..b$ for $x = \{y \in \mathbb{N} : a \leq y \leq b\}$. If x is a set, TLA denotes $SUBSET x$ to be the power-set or the set of all possible subsets of x .

3.2.1 Groups and Communicators:

A group is a set of integers representing process IDs $members \in SUBSET(0..(\mathbb{N}-1))$ and the *size* of *members*, $size = |members|$. If *foo* is a Group then

$foo.members$ is the set of pids in the group and $foo.size$ is the number of elements in $foo.members$ (i.e., $foo[members] \in SUBSET(0..N-1)$ and $foo[size] = |foo(members)|$).

A ranking function and inverse ranking function are maps $ranking : 0..(N-1) \rightarrow 0..(N-1)$, $invranking : 0..(N-1) \rightarrow 0..(N-1)$ such that $\forall k \in Dom(ranking) : \exists n \in 0..(N-1) : ranking[k] = n \wedge invranking[n] = k \wedge \forall m \in Dom(ranking) : ranking[k] = ranking[m] \Rightarrow k = m$. A ranking and inverse ranking function are associated with each group.

A communication universe is record containing a group handle $group$ and a collective context handle $collective$. Groups and Communicators are referenced by handles on processes. Thus the mapping from handles to group or communicator records may be different on any process.

3.2.2 The collective context

Each communicator has a collective context associated with it. The collective context is not directly accessible to the user program, only through the handle in the associated communicator.

Our model currently includes `MPI_BARRIER` as two transitions: `MPI_BARRIER_INIT` and `MPI_BARRIER_WAIT`. The collective context is a record having

<i>participants</i>	$\rightarrow x \in SUBSET(0..(N-1))$
<i>root</i>	$\rightarrow 0..(N-1)$
<i>type</i>	$\rightarrow \{ "barrier" \}$
<i>state</i>	$\rightarrow \{ "in", "out", "vacant" \}$

All processes in the communicator's group must participate in the collective communication. Collective operations operate under a simple state machine. When no process is in the communication the state is "vacant" and the participants set is empty. As processes enter the operation their pid is added to the participants set and the first process changes the state from "vacant" to "in" and sets the type of the communication to "barrier". Processes are only allowed to enter the communication when the state is "in". `MPI_BARRIER_INIT` performs the addition of a process to the participant set when the state is "vacant" or "in" and the process is not represented in the set of participants; blocking the process applying this operator otherwise.

When all processes in the group are in the participant set then the state of the operation changes from "in" to "out" and processes are allowed to exit. `MPI_BARRIER_WAIT`

blocks the calling process until the state is “out”, removes the process applying the operator from the participant set, and sets the state to “vacant” if the process is the last to leave the communication.

Additional collective operations can be implemented by adding additional collective message types to the range of *collective.type* and appropriate checks on the parameters that are passed to the operators.

3.2.3 Requests

The set of requests represent the point to point contexts of all communicators. Messages are paired only if they have the same communicator handle (which in our model are unique across space and time).

A message is represented by the envelope that includes all information needed to pair and transmit point to point communication operations. We model messages as a record (i.e., a function having character strings as elements of the domain) as follows:

<i>data</i>	→ <code>Buffers</code>
<i>src</i>	→ $0..(N - 1) \cup \{\text{MPI_ANY_SOURCE}\}$
<i>dest</i>	→ $0..(N - 1)$
<i>msgtag</i>	→ <code>TAGS</code> \cup <code>{MPI_ANY_TAG}</code>
<i>dtype</i>	→ <code>TYPES</code> \cup <code>{MPI_TYPES}</code>
<i>num</i>	→ \mathbb{N}
<i>universe</i>	→ $0..(\text{MAX_COMM} - 1)$
<i>state</i>	→ <code>{“send”, “recv”}</code>

Where `Buffers` is a placeholder for future inclusion of data in a model.

A request is the bookkeeping information needed to manage messages within a process. Request objects are required to be opaque to the user process and are therefore represented by a function `requests : $\mathbb{N} \rightarrow Request$` where the set *Request* is the set of all possible request objects. The request handle is the element of the domain of the `requests` function which returns the associated request object.

With *Seq*(\mathbb{N}) as the set of all sequences of natural numbers, we model request

objects as records as follows:

<i>error</i>	$\rightarrow \mathbb{N}$
<i>active</i>	$\rightarrow \{TRUE, FALSE\}$
<i>transmitted</i>	$\rightarrow \{TRUE, FALSE\}$
<i>buffered</i>	$\rightarrow \{TRUE, FALSE\}$
<i>started</i>	$\rightarrow \{TRUE, FALSE\}$
<i>canceled</i>	$\rightarrow \{TRUE, FALSE\}$
<i>deallocated</i>	$\rightarrow \{TRUE, FALSE\}$
<i>ctype</i>	$\rightarrow \{ "send", "bsend", "ssend", "rsend", "recv" \}$
<i>persist</i>	$\rightarrow \{TRUE, FALSE\}$
<i>match</i>	$\rightarrow Seq(\mathbb{N})$
<i>message</i>	$\rightarrow Messages$

A new request is appended to the requests function as described above. Each request record is accessible by the user process through its associated handle until that record is marked as deallocated either by successful application of a message completion operator such as `MPI_WAIT` or `MPI_REQUEST_FREE`. The handles are set to `MPI_REQUEST_NULL` at this time and become unaccessible to the user process.

3.2.4 Message buffers and buffer size

Users may wish to provide buffer space to the MPI system and allow the MPI system to manage that buffer space. Calls to `MPI_BSEND`, `MPI_IBSEND`, and `MPI_BSEND_INIT` use this buffer that is specified through `MPI_BUFFER_ATTACH`.

Not modeling data, the buffers are represented by a counting semaphore to track resource availability. Only one buffer can be attached to the MPI system for a process at a time. We approximate the use of the buffer space as follows. The user specifies how many messages can be stored in the buffer by the call to `MPI_BUFFER_ATTACH`. When a message is activated one buffer slot is consumed until the message is transmitted or canceled. Accordingly, `MPI_BUFFER_DETACH` blocks the process applying the operator until all buffered messages have either been transmitted or canceled.

$$\begin{aligned}
 message_buffer &: 0..(N - 1) \rightarrow \mathbb{N} \\
 buf_size &: 0..(N - 1) \rightarrow \mathbb{N}
 \end{aligned}$$

The `message_buffer` variable is a function that represents the counting semaphore for each process. The `bufsize` variable is a function that represents the maximum values for each of the associated `message_buffer` variables.

3.3 Statuses

MPI operations return information to the user program in two ways. The first is the return value of a function. We do not model this. The second way information is returned to the user program is via the status object.

We model a status as a record with members as follows:

<i>state</i>	$\rightarrow \{ \text{"defined"}, \text{"undefined"}, \text{"empty"} \}$
<i>MPI_SOURCE</i>	$\rightarrow 0..(N - 1) \cup \{ \text{MPI_PROC_NULL}, \text{MPI_ANY_SOURCE} \}$
<i>MPI_TAG</i>	$\rightarrow \text{TAGS} \cup \{ \text{MPI_ANY_TAG} \}$
<i>MPI_ERROR</i>	$\rightarrow \mathbb{N}$
<i>count</i>	$\rightarrow \mathbb{N}$
<i>canceled</i>	$\rightarrow \{ \text{TRUE}, \text{FALSE} \}$

4 Collective Communications

5 Closing the model for use with model checking

Many things are left and specified by MPI. Among these are details on how messages are communicated between processes. So far we've introduced the request, status, and communicator records. Using the temporal logic of actions we now have sufficient structure in our model to specify the pairing, buffering, transmitting, and completing of messages.

5.1 Completing messages

5.1.1 Envelope matching

Envelopes match according to the operator `Match` shown in appendix A.

5.1.2 Pairing messages

Messages are paired together as a send request and a receive request. Program order must be observed on both the send and receive process when matching two

requests. To enforce this policy, the operator that performs message pairing specifies the earliest active message in the sequence that has not been canceled, transmitted, or paired previously. Operator Pair contains the logic of this operation. Pairs of messages that have been started, not matched, not canceled, not transmitted, and where one message is a send and the other message is a receive can be paired. In addition we require messages to have matching envelopes.

5.1.3 Transmitting messages

Once messages are paired appropriately they may complete in any order. Thus it is not enough to model communication as the pairing of messages. The Transmit operator contains the logic involved in passing data from one process to another. Only messages that have been started, have not been canceled, have not previously been transmitted, and have been previously paired can be transmitted. The request is updated to reflect that the corresponding message has been transmitted.

5.1.4 Buffering messages

Message buffering can happen under two circumstances. The first is when the user specifically requests MPI to buffer the outgoing messages using commands such as `MPI_IBSEND`. These messages may be buffered at any time after the message is started and before the message has transmitted. The operator `Buffer_bsend` contains the logic to mark requests when messages have been buffered appropriately. Thereby allowing the sending process to continue when the corresponding message completion operator is applied.

When using `MPI_SEND` this system may choose to buffer the outgoing message. We allow this to happen at any time after the message is posted up until the message is transmitted or canceled. However it may also be the case that the MPI system will never buffer such a message. The operator `Buffer_send` performs this operation.

It is possible, from the user's perspective, for the message to be buffered and transmitted before the user program regains control. For this reason we allow the message to be buffered up to the point where the message is actually transmitted or canceled by the user.

6 Modeling MPI Programs

There are many ways to model programs in TLA. The +CAL tool makes it significantly easier to take this step [8]. We will describe a similar modeling paradigm that suits our needs.

In modeling MPI programs in connection with the TLA MPI specification, we assume for simplicity that all programs are written in the SPMD style. Although this is not required, it is required that all variables be declared as arrays as described in section 2.

It is also convenient to assume that all programs make only MPI function calls, although adding procedure calls is a relatively trivial extension. Closing the environment and making available other standard system procedures is an important area of research but is beyond the scope of this work.

6.1 Sequential execution

Let PC be an array $[0..(N - 1) \rightarrow Labels]$ such that each process $i \in 0..(N - 1)$ in the distributed computation has a program counter represented by $PC[i]$. The transition relation of a sequential program can be specified as a disjunct of conjuncts where each conjunct has (i) a current PC guard, (ii) the specified next PC after executing the conjunct, and (iii) an action associated with the current PC that modifies the state – perhaps only the PC itself.

All control statements can be modeled using the explicit PC and an IF construct provided by TLA.

6.2 Multiple step MPI procedures

As mentioned before, when using a multi-step MPI operator these can be compiled into some sequence of single-step operators. We present possible solutions for the seven contained in MPI that are not present in our TLA model.

The MPI operations can be modeled using a sequence of transitions with *proc* being the pid of the process, “in”² being the starting PC of the call to MPI_SEND, “intermediate” being the middle PC, and “out” being the return PC, and the variable $req \in [0..(N - 1) \rightarrow Request]$. We also consider the status variable $stat : [0..(N - 1) \rightarrow Status]$ as defined in the appendix.

²Strings are valid PC values in TLA. Recall that the PC is a function whose domain is the set of pids and the range is in this case a string.

6.2.1 MPI_SEND

Applying MPI_SEND in a program having sequential execution can be implemented follows:

$$\begin{aligned} & \vee \wedge pc[proc] = \text{"in"} \\ & \wedge pc = [pc \text{ EXCEPT } ![proc] = \text{"intermediate"}] \\ & \wedge \text{MPI_Isend}(\text{buf}, \text{count}, \text{datatype}, \text{dest}, \text{tag}, \text{com}, \text{req}, \text{proc}) \\ & \vee \wedge pc[proc] = \text{"intermediate"} \\ & \wedge pc' = [pc \text{ EXCEPT } ![proc] = \text{"out"}] \\ & \wedge \text{MPI_Wait}(\text{req}, \text{stat}, \text{proc}) \end{aligned}$$

6.2.2 MPI_BSEND

Applying MPI_BSEND is as follows:

$$\begin{aligned} & \vee \wedge pc[proc] = \text{"in"} \\ & \wedge pc = [pc \text{ EXCEPT } ![proc] = \text{"intermediate"}] \\ & \wedge \text{MPI_Bsend}(\text{buf}, \text{count}, \text{datatype}, \text{dest}, \text{tag}, \text{com}, \text{req}, \text{proc}) \\ & \vee \wedge pc[proc] = \text{"intermediate"} \\ & \wedge pc' = [pc \text{ EXCEPT } ![proc] = \text{"out"}] \\ & \wedge \text{MPI_Wait}(\text{req}, \text{stat}, \text{proc}) \end{aligned}$$

The restrictions on attaching buffers and managing the buffer space are identical.

6.2.3 MPI_SSEND

Applying MPI_SSEND is as follows:

$$\begin{aligned} & \vee \wedge pc[proc] = \text{"in"} \\ & \wedge pc = [pc \text{ EXCEPT } ![proc] = \text{"intermediate"}] \\ & \wedge \text{MPI_Ssend}(\text{buf}, \text{count}, \text{datatype}, \text{dest}, \text{tag}, \text{com}, \text{req}, \text{proc}) \\ & \vee \wedge pc[proc] = \text{"intermediate"} \\ & \wedge pc' = [pc \text{ EXCEPT } ![proc] = \text{"out"}] \\ & \wedge \text{MPI_Wait}(\text{req}, \text{stat}, \text{proc}) \end{aligned}$$

6.2.4 MPI_RSEND

Applying MPI_RSEND is as follows:

$$\begin{aligned} & \vee \wedge pc[proc] = \text{"in"} \\ & \wedge pc = [pc \text{ EXCEPT } ![proc] = \text{"intermediate"}] \\ & \wedge \text{MPI_Irsend}(\text{buf}, \text{count}, \text{datatype}, \text{dest}, \text{tag}, \text{com}, \text{req}, \text{proc}) \\ & \vee \wedge pc[proc] = \text{"intermediate"} \\ & \wedge pc' = [pc \text{ EXCEPT } ![proc] = \text{"out"}] \\ & \wedge \text{MPI_Wait}(\text{req}, \text{stat}, \text{proc}) \end{aligned}$$

6.2.5 MPI_RECV

Applying MPI_RECV is as follows:

$$\begin{aligned} & \vee \wedge pc[proc] = \text{"in"} \\ & \wedge pc = [pc \text{ EXCEPT } ![proc] = \text{"intermediate"}] \\ & \wedge \text{MPI_Irecv}(\text{buf}, \text{count}, \text{datatype}, \text{source}, \text{tag}, \text{com}, \text{req}, \text{proc}) \\ & \vee \wedge pc[proc] = \text{"intermediate"} \\ & \wedge pc' = [pc \text{ EXCEPT } ![proc] = \text{"out"}] \\ & \wedge \text{MPI_Wait}(\text{req}, \text{stat}, \text{proc}) \end{aligned}$$

6.2.6 MPI_SENDRECV

Overloading *req* and *stat* to be arrays of records appropriately, MPI_SENDRECV could be implemented as follows:

$$\begin{aligned} &\vee \wedge pc[proc] = \text{"in"} \\ &\wedge pc' = [pc \text{ EXCEPT } ![proc] = \text{"intermediate_recv"} \\ &\wedge \text{MPI_Isend}(\text{sendbuf}, \text{sendcount}, \text{sendtype}, \text{dest}, \text{sendtag}, \text{com}, \text{req1}, \text{proc}) \\ &\vee \wedge pc[proc] = \text{"intermediate_recv"} \\ &\wedge pc' = [pc \text{ EXCEPT } ![proc] = \text{"wait"} \\ &\wedge \text{MPI_Irecv}(\text{recvbuf}, \text{recvcount}, \text{recvtype}, \text{source}, \text{recvtag}, \text{com}, \text{req2}, \text{proc}) \\ &\vee \wedge pc[proc] = \text{"in"} \\ &\wedge pc' = [pc \text{ EXCEPT } ![proc] = \text{"intermediate_send"} \\ &\wedge \text{MPI_Irecv}(\text{recvbuf}, \text{recvcount}, \text{recvtype}, \text{source}, \text{recvtag}, \text{com}, \text{req2}, \text{proc}) \\ &\vee \wedge pc[proc] = \text{"intermediate_send"} \\ &\wedge pc' = [pc \text{ EXCEPT } ![proc] = \text{"wait"} \\ &\wedge \text{MPI_Isend}(\text{sendbuf}, \text{sendcount}, \text{sendtype}, \text{dest}, \text{sendtag}, \text{com}, \text{req1}, \text{proc}) \\ &\vee \wedge pc = \text{"wait"} \\ &\wedge pc' = [pc \text{ EXCEPT } ![proc] = \text{"out"} \\ &\wedge \text{MPI_Waitall}(2, [\text{req EXCEPT } ![proc] = [0 \mapsto \text{req1}[proc], 1 \mapsto \text{req2}[proc]]], \text{stat}, \text{proc}) \end{aligned}$$

6.2.7 MPI_SENDRECV_REPLACE

In addition to overloading *req* and *stat* to be arrays of records appropriately we add a temporary variable for receiving the results. MPI_SENDRECV_REPLACE

could be implemented as follows:

$$\begin{aligned}
& \vee \wedge pc[proc] = \text{"in"} \\
& \wedge pc' = [pc \text{ EXCEPT } ![proc] = \text{"intermediate_recv"} \\
& \wedge \text{MPI_Isend}(buf, sendcount, sendtype, dest, sendtag, com, req1, proc) \\
& \vee \wedge pc[proc] = \text{"intermediate_recv"} \\
& \wedge pc' = [pc \text{ EXCEPT } ![proc] = \text{"wait"} \\
& \wedge \text{MPI_Irecv}(tempbuf, recvcount, recvtype, source, recvtag, com, req2, proc) \\
& \vee \wedge pc[proc] = \text{"in"} \\
& \wedge pc' = [pc \text{ EXCEPT } ![proc] = \text{"intermediate_send"} \\
& \wedge \text{MPI_Irecv}(tempbuf, recvcount, recvtype, source, recvtag, com, req2, proc) \\
& \vee \wedge pc[proc] = \text{"intermediate_send"} \\
& \wedge pc' = [pc \text{ EXCEPT } ![proc] = \text{"wait"} \\
& \wedge \text{MPI_Isend}(sendbuf, sendcount, sendtype, dest, sendtag, com, req1, proc) \\
& \vee \wedge pc = \text{"wait"} \\
& \wedge pc' = [pc \text{ EXCEPT } ![proc] = \text{"copy"} \\
& \wedge \text{MPI_Waitall}(2, [req \text{ EXCEPT } ![proc] = [0 \mapsto req1[proc], 1 \mapsto req2[proc]]], stat, proc) \\
& \vee \wedge pc = \text{"copy"} \\
& \wedge pc' = [pc \text{ EXCEPT } ![proc] = \text{"out"} \\
& \wedge sendbuf' = [buf \text{ EXCEPT } ![proc] = temp[proc]]
\end{aligned}$$

6.3 An example

An example program is included in Appendix A. This program exercises the immediate mode synchronous send, along with the immediate mode receive. Processes are conceptually placed in a ring. Even ranked processes send to the neighbor with higher rank (mod ring size), synchronize on the barrier, and then receive from the neighbor having lower rank (again mod ring size). Odd ranked processes receive from the neighbor with lower rank, synchronize on the barrier and then send to the neighbor with higher rank.

The program is represented as a disjunct of conjuncts similar in style to Section 6.2. This operator has one parameter which is the process id of the process that is currently executing—therein we model the SPMD style where every process executes the same program image.

The next state relation for the entire system is the initial state of the model *Init* and henceforth () the *Next* relation that performs either a *Pair*, *Transmit*, *Buffer*, or *Proc* move for some pid at any step.

7 Conclusions

The TLA model of MPI in connection with this paper describes the reactive behavior of all 35 point to point communication operations from chapter 3 of the MPI 1.1 standard.

We have closed the model for model checking single threaded programs that communicate via MPI point to point operations. We have provided the additional MPI operations necessary to initialize, determine the rank of a process, the size of a communicator's group, and exit according to the MPI standard.

References

- [1] Martín Abadi, Leslie Lamport, and Stephan Merz. A TLA solution to the RPC-Memory specification problem. In M. Broy, S. Merz, and K. Spies, editors, *Formal System Specification: The RPC-Memory Specification Case Study*, volume 1169 of *Lecture Notes in Computer Science*, pages 21–66. Springer-Verlag, Berlin, 1996.
- [2] Steven Barrus, Ganesh Gopalakrishnan, Robert M. Kirby, and Robert Palmer. Verification of MPI programs using SPIN. Technical Report UUCS-04-008, The University of Utah, 2004.
- [3] E.M. Clarke, O. Grumberg, and D. Peled. *Model Checking*. MIT Press, 1999.
- [4] Philippe Georgelin, Laurence Pierre, and Tin Nguyen. A formal specification of the MPI primitives and communication mechanisms. Technical report, LIM, 1999.
- [5] John Harrison. Formal verification of square root algorithms. *Formal Methods in System Design*, 22(2):143–154, March 2003. Guest Editors: Ganesh Gopalakrishnan and Warren Hunt, Jr.
- [6] IEEE standard for radix-independent floating-point arithmetic, ANSI/IEEE Std 854-1987.
- [7] Leslie Lamport. *Specifying Systems: The TLA+ Language and Tools for Hardware and Software Engineers*. Addison Wesley Professional, 2002.

- [8] Leslie Lamport. A +CAL user's manual. <http://research.microsoft.com/users/lamport/tla/p-manual.pdf>, April 2006.
- [9] Glenn R. Luecke, Silvia Spanoyannis, and Marina Kraeva. The performance and scalability of SHMEM and MPI-2 one-sided routines on a SGI Origin 2000 and a cray T3E-600. *Concurrency and Computation: Practice and Experience*, 16(10):1037–1060, 2004.
- [10] MPI: A Message-Passing Interface Standard. <http://www.mpi-forum.org/docs/mpi-11-html/mpi-report.html>.
- [11] Robert Palmer, Steven Barrus, Yu Yang, Ganesh Gopalakrishnan, and Robert M. Kirby. Gauss: A framework for verifying scientific computing software. In *SoftMC: Workshop on Software Model Checking*, number 953 in ENTCS, August 2005.
- [12] Salman Pervez. Byte-range-locks using mpi-one-sided communication: a report. Technical report, The University of Utah, 2006. Submitted.
- [13] Stephen F. Siegel and George Avrunin. Analysis of mpi programs. Technical Report UM-CS-2003-036, Department of Computer Science, University of Massachusetts Amherst, 2003.
- [14] Stephen F. Siegel and George S. Avrunin. Verification of mpi-based software for scientific computation. In *Proceedings of the 11th International SPIN Workshop on Model Checking Software*, volume 2989 of LNCS, pages 286–303, Barcelona, April 2004. Springer.
- [15] Stephen F. Siegel and George S. Avrunin. Modeling wildcard-free MPI programs for verification. In *ACM SIGPLAN Symposium on Principles and Practices of Parallel Programming*, pages 95–106, Chicago, June 2005.
- [16] Stephen F Siegel, Anastasia Mironova, and George S Avrunin nad Lori A Clarke. Using model checking with symbolic execution to verify parallel numerical programs. In Lori Pollock and Mauro Pezz, editors, *Proceedings of the ACM SIGSOFT International Symposium on Software Testing and Analysis*, pages 157–168, Portland, ME, July 2006.

A The Full Specification

The formal *MPI* library specification.

Robert Palmer

The University of *Utah*

School of Computing

Some notes : – Need to split the buffer rule into rules - one for user specified buffering and one for system provided buffering. – don't really know how

- Need to add deallocation of requests to the model as in *mpi_wait*.

- Need to add more semantics.

- Need to cause buffers to be freed appropriately when a message is sent.

- Need to add a return code to indicate success or error and error handling.

- Need to fix the buffering of standard mode sends such that they might block forever.

EXTENDS *Naturals, TLC, Sequences, FiniteSets*

Constants are given values in the configuration file that accompanies this document: *mpi_base.cfg*

CONSTANTS

<i>N</i> ,	The number of processes in the computation.
<i>MAX_COMM</i> ,	The highest allowed handle value for a communicator. This is not in the standard but makes our model finite.
<i>MAX_GROUP</i> ,	The highest allowed handle value for a group.
<i>TYPES</i> ,	The set of user defined types.
<i>TAGS</i> ,	The set of user defined tags.
<i>SEND_IS_BUFFERED</i> ,	A flag to indicate whether sends are to be buffered.
<i>RANK_ORDERINGS_SIGNIFICANT</i> ,	a flag to indicate whether all possible ranking orders should be considered in verification
<i>MPI_COMM_WORLD</i> ,	The handle for <i>MPI_COMM_WORLD</i> .
<i>MPI_ANY_SOURCE</i> ,	The wildcard source rank.
<i>MPI_ANY_TAG</i> ,	The wildcard tag value.
<i>MPI_PROC_NULL</i> ,	Section 3.11 Null Processes
<i>MPI_REQUEST_NULL</i> ,	A special handle value for requests.

	Set this to 0 in the configuration file and make the initial values of the requests occupied to avoid an array out-of-bounds error.
<i>MPI_SUCCESS</i> ,	The return value of a successful call to an <i>MPI</i> procedure.
<i>MPI_IDENT</i> ,	5.4: Two communicator handles refer to the same communicator.
<i>MPI_CONGRUENT</i> ,	The communicator handles are different; communicators differ only in context.
<i>MPI_SIMILAR</i> ,	The communicator handles are different; communicators have the same group, however both context and ranking differ.
<i>MPI_UNEQUAL</i> ,	The communicator handles are different; communicators have different groups, contexts, and rankings.
<i>MPI_UNDEFINED</i> ,	A special rank returned to a process that is not a member of the queried communicator.
<i>MPI_INT</i> ,	<i>MPI</i> defined <i>datatype</i> for integers
<i>MPI_FLOAT</i> ,	<i>MPI</i> defined <i>datatype</i> for floating point numbers
<i>UB</i>	The upper bound on the tag range 19.27 – 19.31
<i>MPI_GROUP_EMPTY</i> \ *	The empty group

Variables represent the state of the *MPI* system at any given time. None of these state elements are specified by the standard. However they are useful to describe what is specified. In particular mention is made of handles that reference opaque objects. The communicator and requests arrays are such opaque objects that are referenced by integer handles that in our model are unique across both space and time (*i.e.*, the same value is used for *MPI_COMM_WORLD* on all processes for the entire execution etc.).

VARIABLES

<i>communicator</i> ,	An array of communication universe objects.
<i>bufsize</i> ,	The size of the user attached <i>message_buffer</i> .
<i>message_buffer</i> ,	The user attached buffer.
<i>requests</i> ,	A array of message requests lists, one per process. Although we do model the allocation of request objects by adding a structure to a list of requests, we are not modeling the freeing of requests more than setting the associated handle to <i>MPI_REQUEST_NULL</i> .
<i>initialized</i> ,	An array of flags that indicate whether <i>MPI_Init</i>

	has been called by a given process.
<i>collective</i> ,	The collective contexts for all communicators
<i>group</i> ,	The array of groups
<i>Memory</i>	A model of memory for individual processes.

Type invariant

Memory is considered a *program_var*

$mpi_vars \triangleq \langle group, communicator, bufsize, message_buffer, requests, initialized, collective \rangle$

$Messages \triangleq [src : (0 .. (N - 1)) \cup \{MPI_ANY_SOURCE\}, 21.24 - 21.25$
 $dest : (0 .. (N - 1)), 19.39$
 $msgtag : 0 .. UB \cup \{MPI_ANY_TAG\}, 19.28$
 $dtype : TYPES \cup \{MPI_FLOAT, MPI_INT\},$
 $numelements : Nat,$
 $universe : (MPI_COMM_WORLD .. (MPI_COMM_WORLD + MAX_COMM)),$
 $state : \{“send”, “recv”\},$
 $addr : Nat]$

$Message_types \triangleq \{“send”, “bsend”, “ssend”, “rsend”, “recv”\}$

$Collective_types \triangleq \{“barrier”\}$

$Collective_states \triangleq \{“in”, “out”, “vacant”\}$

$Request \triangleq [error : Nat,$
 $active : BOOLEAN ,$
 $transmitted : BOOLEAN ,$
 $buffered : BOOLEAN ,$
 $started : BOOLEAN ,$
 $cancelled : BOOLEAN ,$
 $deallocated : BOOLEAN ,$
 $ctype : Message_types,$
 $persist : BOOLEAN ,$
 $match : Seq(Nat),$
 $message : Messages]$

$Requests \triangleq [(0 .. (N - 1)) \rightarrow Seq(Request)]$

22.1 – 22.8
 $Statuses \triangleq [state : \{“defined”, “undefined”, “empty”\},$
 $MPI_SOURCE : (0 .. (N - 1)) \cup \{MPI_PROC_NULL, MPI_ANY_SOURCE\},$
 $MPI_TAG : TAGS \cup \{MPI_ANY_TAG\},$
 $MPI_ERROR : Nat,$

```

count      : Nat,
cancelled  : BOOLEAN ]

```

Refactored into *Memory* to allow for a uniform treatment of the model of memory and to facilitate modelling using pointer arithmetic for member accesses.

21.45 – 21.48

Status variables are explicitly allocated by the user. Therefore they are present in the *Memory* of individual processes. We will use a simple offset mechanism to return the individual member addresses within *Memory*.

22.1 – 22.8

```

Status_Cancelled(base)  $\triangleq$  base
Status_Count(base)     $\triangleq$  base + 1
Status_Source(base)    $\triangleq$  base + 2
Status_Tag(base)       $\triangleq$  base + 3
Status_Err(base)       $\triangleq$  base + 4

```

Initialized \triangleq $[0 .. (N - 1) \rightarrow \{\text{"initialized"}, \text{"uninitialized"}, \text{"finalized"}\}]$

MessageBuffers \triangleq $[0 .. (N - 1) \rightarrow \text{Nat}]$

BufferSizes \triangleq $[0 .. (N - 1) \rightarrow \text{Nat}]$

Groups can be different on different processes.

```

Group  $\triangleq$   $[0 .. (N - 1) \rightarrow$  138.37 – 138.38
   $[$ MPI_COMM_WORLD .. (MPI_COMM_WORLD + MAX_GROUP)  $\rightarrow$ 
     $[$ members : SUBSET (0 .. (N - 1)),
      size   : 0 .. N,
      ranking :  $[0 .. (N - 1) \rightarrow 0 .. (N - 1)]$ ,
      invranking :  $[0 .. (N - 1) \rightarrow 0 .. (N - 1)]$ ]]]

```

```

Communicator  $\triangleq$   $[0 .. (N - 1) \rightarrow$ 
   $[$ MPI_COMM_WORLD .. (MPI_COMM_WORLD + MAX_COMM)  $\rightarrow$ 
     $[$ group : MPI_COMM_WORLD .. (MPI_COMM_WORLD + MAX_GROUP),
      collective : MPI_COMM_WORLD .. (MPI_COMM_WORLD + MAX_COMM)]]]

```

```

Collective  $\triangleq$   $[$ (MPI_COMM_WORLD .. (MPI_COMM_WORLD + MAX_COMM))  $\rightarrow$ 
   $[$ participants : SUBSET (0 .. (N - 1)),
    root : (0 .. (N - 1)),
    type : Collective_types,
    state : Collective_states]]

```

Comm_inv \triangleq *communicator* \in *Communicator*

Buff_inv \triangleq *bufsize* \in *BufferSizes*

Msg_buf_inv \triangleq *message_buffer* \in *MessageBuffers*

Initialized_inv \triangleq *initialized* \in *Initialized*

Request_inv \triangleq *requests* \in *Requests*

$Col_inv \triangleq collective \in Collective$

$group_inv \triangleq group \in Group$

$MPI_Type_Invariant \triangleq$

$\wedge communicator \in Communicator$
 $\wedge bufsize \in BufferSizes$
 $\wedge message_buffer \in MessageBuffers$
 $\wedge initialized \in Initialized$
 $\wedge requests \in Requests$
 $\wedge collective \in Collective$

$Make_request$ is a rule to simplify the expressions that create a new request object. Section 3.7.1

$Make_request(err, act, com, sta, buf, cty, per, mat, can, mes) \triangleq$

$error$	$\mapsto err,$	The error code associated with this request
$active$	$\mapsto act,$	The message was initiated
$transmitted$	$\mapsto com,$	Data was transmitted by this message
$started$	$\mapsto sta,$	Start this request
$buffered$	$\mapsto buf,$	The data was copied from the input address
$cancelled$	$\mapsto can,$	Whether the request was <i>cancelled</i>
$deallocated$	$\mapsto FALSE,$	A new request is created in an allocated state
$ctype$	$\mapsto cty,$	The type of message (send, <i>b</i> send, <i>r</i> send, or <i>s</i> send)
$persist$	$\mapsto per,$	Whether the request is a persistent communication
$match$	$\mapsto mat,$	The matching < process,handle >
$message$	$\mapsto mes]$	The message envelope associated with this request

The initial values for the *MPI* specification state variables. These are not specified by the standard, however these initial values make the TLA+ representation complete such that it can be verified using *TLC*.

$MPI_Specification_Init \triangleq$

$\wedge requests = [i \in (0 .. (N - 1)) \mapsto$ Create an instance of *MPI_REQUEST_NULL*
 $\langle Make_request(0, FALSE, FALSE, FALSE, FALSE,$ for each process.
 $“send”, FALSE, \langle \rangle, TRUE,$
 $[src \mapsto 0,$
 $dest \mapsto 0,$
 $msgtag \mapsto MPI_ANY_TAG,$
 $dtype \mapsto 0,$
 $numelements \mapsto 0,$
 $universe \mapsto MPI_COMM_WORLD,$
 $state \mapsto “send”,$
 $addr \mapsto 0]]]$

$\wedge bufsize = [i \in (0 .. (N - 1)) \mapsto 0]$ Each process starts with no user attached buffer.
 $\wedge message_buffer = [i \in (0 .. (N - 1)) \mapsto 0]$ Each process starts with no messages buffered.
 $\wedge initialized = [i \in (0 .. (N - 1)) \mapsto “uninitialized”]$ Each process starts uninitialized.

\wedge *communicator* = $[a \in 0 \dots (N-1) \mapsto [i \in \text{MPI_COMM_WORLD} \dots (\text{MPI_COMM_WORLD} + \text{MAX_COMM}) \mapsto$
 \wedge *collective* = $[i \in (\text{MPI_COMM_WORLD} \dots (\text{MPI_COMM_WORLD} + \text{MAX_COMM})) \mapsto$
 $[$ *participants* $\mapsto \{\},$
 $\textit{root} \mapsto 0,$
 $\textit{type} \mapsto \text{"barrier"},$
 $\textit{state} \mapsto \text{"vacant"}]$

$\wedge \vee \wedge \neg \text{RANK_ORDERINGS_SIGNIFICANT} \setminus * \text{ In this case, choose an arbitrary ordering}$

\wedge CHOOSE $f \in [0 \dots (N-1) \rightarrow 0 \dots (N-1)] : \setminus * \text{12.41 - 12.42 order is not specified.}$

CHOOSE $\textit{finv} \in [0 \dots (N-1) \rightarrow 0 \dots (N-1)] : \setminus * \text{The inverse of } f$

$\forall k \in \text{DOMAIN } f :$

$\exists n \in 0 \dots (N-1) :$

$\wedge f[k] = n$

$\wedge \textit{finv}[n] = k$

$\wedge \forall m \in \text{DOMAIN } f : f[k] = f[m] \Rightarrow k = m$

\wedge *group* = $[a \in 0 \dots (N-1) \mapsto [i \in (\text{MPI_COMM_WORLD} \dots$
 $((\text{MPI_COMM_WORLD} + \text{MAX_GROUP})) \mapsto$

IF $i = \text{MPI_COMM_WORLD}$

THEN

$[$ *members* $\mapsto \{x \in 0 \dots (N-1) : \text{TRUE}\},$

$\textit{size} \mapsto N,$

$\textit{ranking} \mapsto f,$

$\textit{invranking} \mapsto \textit{finv}]$

ELSE

$[$ *members* $\mapsto \{\},$

$\textit{size} \mapsto 0,$

$\textit{ranking} \mapsto [j \in 0 \dots (N-1) \mapsto 0],$

$\textit{invranking} \mapsto [j \in 0 \dots (N-1) \mapsto 0]]]$

$\vee \wedge \text{RANK_ORDERINGS_SIGNIFICANT} \setminus * \text{ in this case, try all orderings}$

$\wedge \exists f \in [0 \dots (N-1) \rightarrow 0 \dots (N-1)] : \text{12.41 - 12.42 order is not specified.}$

$\exists \textit{finv} \in [0 \dots (N-1) \rightarrow 0 \dots (N-1)] : \text{The inverse of } f$

$\forall k \in \text{DOMAIN } f :$

$\exists n \in 0 \dots (N-1) :$

$\wedge f[k] = n$

$\wedge \textit{finv}[n] = k$

$\wedge \forall m \in \text{DOMAIN } f :$

$\wedge f[k] = f[m] \Rightarrow k = m$

\wedge *group* = $[a \in 0 \dots (N-1) \mapsto$

$[i \in (\text{MPI_COMM_WORLD} \dots ((\text{MPI_COMM_WORLD} + \text{MAX_GROUP})) \mapsto$

IF $i = \text{MPI_COMM_WORLD}$

THEN

$[$ *members* $\mapsto \{x \in 0 \dots (N-1) : \text{TRUE}\},$

$\textit{size} \mapsto N,$

$\textit{ranking} \mapsto f,$

$\textit{invranking} \mapsto \textit{finv}]$

ELSE

$[$ *members* $\mapsto \{\},$

$\textit{size} \mapsto 0,$

$$\begin{aligned} \textit{ranking} &\mapsto [j \in 0 \dots (N - 1) \mapsto 0], \\ \textit{invranking} &\mapsto [j \in 0 \dots (N - 1) \mapsto 0]]] \end{aligned}$$

A correct *MPI* program is one in which all messages that are posted are eventually transmitted or *cancelled*. A message that is posted but never transmitted is in error. It seems that a message that is transmitted but never completed locally may also be in error... I should check on this.

Messages_sent_are_received_and_completed \triangleq
 $\forall i \in (0 \dots (N - 1)) :$
 $\forall m \in (1 \dots \textit{Len}(\textit{requests}[i])) :$
 LET $r \triangleq \textit{requests}[i][m]$ IN
 $r.\textit{active} \rightsquigarrow$
 $\wedge \vee r.\textit{transmitted}$
 $\vee r.\textit{cancelled}$
 $\wedge \neg r.\textit{active}$

There is some issue with regards to where the UNCHANGED identifiers should be living. I am using the following protocol:

1. Rules that have parameters that might be changed will declare the UNCHANGED value appropriately inside the rule for those parameters.
2. Variables that are passed as parameters to rules must be declared as UNCHANGED appropriately outside the rule unless the parameter might be modified by the rule when the rule is used.
3. Constants (such as a literal number, 0 for example) or CONSTANT values need not be declared as UNCHANGED.
4. *MPI* based rules always indicate the UNCHANGED terms for *MPI* state variables. Program models also indicate UNCHANGED for *MPI* variables only when no *MPI* rule is fired in that transition.

Conventions on parameters.

1. Parameters that are set (*i.e.*, OUT or *INOUT*) are all arrays from $0 \dots (N - 1)$ with one instance of each object for each process in the model.
2. All other parameters (*i.e.*, IN) are the single instance of the variable value being passed, or are constant.

These rules perform the communication or “matching” of messages that is necessary to complete the *MPI* communication infrastructure. They are in no way specified in the standard, except that messages are spoken of as being transmitted from one process to another and matching.

$\alpha \rightarrow \beta \rightarrow \text{BOOLEAN}$
 No change in state

$$\begin{aligned}
\text{Match}(a, b) &\triangleq \\
&\wedge \text{Assert}((a.\text{state} = \text{"recv"} \wedge b.\text{state} = \text{"send"}) \vee \\
&\quad (a.\text{state} = \text{"send"} \wedge b.\text{state} = \text{"recv"}), \\
&\quad \text{"Error: Match attempted with two send or receives."}) \\
&\wedge (a.\text{src} = b.\text{src} \vee a.\text{src} = \text{MPI_ANY_SOURCE} \vee b.\text{src} = \text{MPI_ANY_SOURCE}) \quad 21.14 - 21.15 \\
&\wedge a.\text{dest} = b.\text{dest} \\
&\wedge a.\text{dtype} = b.\text{dtype} \quad 23.17, 23.24 - 23.27 \\
&\wedge (a.\text{msgtag} = b.\text{msgtag} \vee a.\text{msgtag} = \text{MPI_ANY_TAG} \vee b.\text{msgtag} = \text{MPI_ANY_TAG}) \quad 21.15 - 21.16 \\
&\wedge a.\text{universe} = b.\text{universe} \quad 19.34 - 19.37 \\
&\wedge \neg a.\text{src} = \text{MPI_PROC_NULL} \quad 60.48 - 61.1 \\
&\wedge \neg a.\text{dest} = \text{MPI_PROC_NULL} \\
&\wedge \neg b.\text{src} = \text{MPI_PROC_NULL} \\
&\wedge \neg b.\text{dest} = \text{MPI_PROC_NULL} \\
&21.13 - 21.14 \text{ count need not be matched in point to point messages.}
\end{aligned}$$

Messages match in program order pairwise between processes, however they may complete in a nondeterministic order on both the sender and receiver. This tends to imply that *Communicate* should in fact be two rules. And it also seems to imply that completion of a message can happen on one side and then on the other also in a non-deterministic way. Therefore *Transmit* should complete only one side of the communication.

Pairs messages together such that they result in a communication eventually.

$$\begin{aligned}
\text{Pair} &\triangleq \\
&\wedge \exists i \in 0 \dots (N - 1) : \\
&\quad \exists j \in 0 \dots (N - 1) : \\
&\quad \quad \exists m \in 1 \dots \text{Len}(\text{requests}[i]) : \\
&\quad \quad \quad \exists n \in 1 \dots \text{Len}(\text{requests}[j]) : \\
&\quad \quad \quad \text{LET } a \triangleq \text{requests}[i][m] \text{ IN} \\
&\quad \quad \quad \text{LET } b \triangleq \text{requests}[j][n] \text{ IN} \\
&\quad \quad \quad \wedge a.\text{started} \\
&\quad \quad \quad \wedge b.\text{started} \\
&\quad \quad \quad \wedge \neg a.\text{cancelled} \\
&\quad \quad \quad \wedge \neg b.\text{cancelled} \\
&\quad \quad \quad \wedge \neg a.\text{transmitted} \\
&\quad \quad \quad \wedge \neg b.\text{transmitted} \\
&\quad \quad \quad \wedge \vee \wedge a.\text{message.state} = \text{"send"} \\
&\quad \quad \quad \quad \wedge b.\text{message.state} = \text{"recv"} \\
&\quad \quad \quad \quad \vee \wedge a.\text{message.state} = \text{"recv"} \\
&\quad \quad \quad \quad \quad \wedge b.\text{message.state} = \text{"send"} \\
&\quad \quad \quad \wedge a.\text{match} = \langle \rangle \\
&\quad \quad \quad \wedge b.\text{match} = \langle \rangle \\
&\quad \quad \quad \wedge \text{Match}(a.\text{message}, b.\text{message}) \\
&\quad \quad \wedge \forall r \in 1 \dots \text{Len}(\text{requests}[i]) : \quad \text{This conjunct enforces the fifo} \\
&\quad \quad \quad \forall s \in 1 \dots \text{Len}(\text{requests}[j]) : \\
&\quad \quad \quad \quad \text{LET } c \triangleq \text{requests}[i][r] \text{ IN} \\
&\quad \quad \quad \quad \text{LET } d \triangleq \text{requests}[j][s] \text{ IN}
\end{aligned}$$

$$\begin{aligned}
& \wedge \text{IF } m.ctype = \text{"b\text{send"}} \\
& \quad \text{THEN} \\
& \quad \quad message_buffer' = [message_buffer \text{ EXCEPT } ![i] = @ - 1] \\
& \quad \text{ELSE} \\
& \quad \quad \text{UNCHANGED } \langle message_buffer \rangle \\
& \wedge \text{UNCHANGED } \langle group, communicator, bufsize, initialized, collective \rangle
\end{aligned}$$

The specification indicates that messages are buffered in an asynchronous manner. The rule *Buffer* is not part of the standard but necessary to allow buffering to complete asynchronously.

$$\begin{aligned}
Buffer & \triangleq \\
& \vee \wedge \exists i \in (0 \dots (N - 1)) : \\
& \quad \exists m \in 1 \dots Len(requests[i]) : \\
& \quad \quad \wedge requests[i][m].started \\
& \quad \quad \wedge requests[i][m].active \\
& \quad \quad \wedge \neg requests[i][m].buffered \\
& \quad \quad \wedge \neg requests[i][m].cancelled \\
& \quad \quad \wedge \neg requests[i][m].transmitted \\
& \quad \quad \wedge \vee \wedge requests[i][m].ctype = \text{"b\text{send}} \quad \text{Buffering is provided explicitly by the user.} \\
& \quad \quad \wedge requests' = \\
& \quad \quad \quad [requests \text{ EXCEPT } ![i] = \\
& \quad \quad \quad \quad [@ \text{ EXCEPT } ![m] = \\
& \quad \quad \quad \quad \quad [@ \text{ EXCEPT } !.buffered = TRUE]]] \\
& \quad \quad \vee \wedge requests[i][m].ctype = \text{"send"} \quad \text{Buffering may be provided by the system.} \\
& \quad \quad \wedge \vee requests' = \\
& \quad \quad \quad [requests \text{ EXCEPT } ![i] = \\
& \quad \quad \quad \quad [@ \text{ EXCEPT } ![m] = \\
& \quad \quad \quad \quad \quad [@ \text{ EXCEPT } !.buffered = TRUE]]] \\
& \quad \quad \vee \text{UNCHANGED } requests \\
& \wedge \text{UNCHANGED } \langle group, communicator, bufsize, message_buffer, initialized, collective \rangle
\end{aligned}$$

General Comments:

1.19.23 – 19.24 The message source is provided in the envelope implicitly. Operators in our model must be passed this information as a parameter. As such we extend the argument list to include *proc*, being the unique identity of the applying process.

Section 3.2 Blocking Send and Receive Operations

Section 3.2.1 Blocking send

Can these really be done in a single transition? I am thinking that it is not possible under an interleaving semantics. In particular, either the send must be two transitions or the receive must be two transitions, it cannot be the case that they are both only one transition.

$MPI_Send(buf, count, datatype, dest, tag, comm, proc) \triangleq$
 $MPI_Isend; MPI_Wait$

Section 3.2.4 Blocking receive If receive is modeled using only one transition, it is just a combination of the MPI_Irecv and Communicate rules.

$MPI_Recv(buf, count, datatype, source, tag, comm, status) \triangleq$
 $MPI_Irecv; MPI_Wait$

Section 3.2.5 Return status

Returns in count the number of data elements in the message represented by status.

$MPI_Get_count(status, datatype, count, return, proc) \triangleq$ 22.24 – 22.37
 $\wedge Assert(Memory[proc][Status_Cancelled(status)] = FALSE,$ 54.47
 “Error: count is undefined on a status from a cancelled message.”)
 $\wedge Assert(initialized[proc] = \text{“initialized”},$ 200.10 – 200.12
 “Error: MPI_Get_count called before process was initialized.”)
 $\wedge Memory' = [Memory \text{ EXCEPT } ![proc] = [@ \text{ EXCEPT } ![count] = Memory[proc][Status_Count(status)]]]$
 $\wedge UNCHANGED \langle group, communicator, bufsize, message_buffer, requests, initialized, collective \rangle$

Section 3.4 Communication Modes

Notes: These, like the above blocking communications really should be modeled using two transitions. In this way, the interleaving semantics is able to schedule another process to complete the communications.

$MPI_Bsend(buf, count, datatype, dest, tag, comm, proc) \triangleq$

$MPI_Ssend(buf, count, datatype, dest, tag, comm, proc) \triangleq$

$MPI_Rsend(buf, count, datatype, dest, tag, comm, proc) \triangleq$

Section 3.6 Buffer allocation and usage

We ignore the buffer argument as data is abstracted away in our model. Buffering is modeled as a counting semaphore, keeping track of the resources available but not exactly which resources are used or what is done with those resources.

Return value is unspecified.

$MPI_Buffer_attach(buffer, size, return, proc) \triangleq$ < 34.17 – 34.33 >
 $\wedge Assert(initialized[proc] = \text{“initialized”},$ 200.10 – 200.12
 “Error: MPI_Buffer_attach called with proc not in initialized state.”)
 $\wedge Assert(bufsize[proc] = 0,$ 34.32
 “Error: MPI_Buffer_attach called when processes buffer is non-zero.”)

$\wedge \text{bufsize}' = [\text{bufsize} \text{ EXCEPT } ![\text{proc}] = \text{size}[\text{proc}]]$ < /34.17 – 34.33 >
 $\wedge \text{UNCHANGED} \langle \text{group}, \text{communicator}, \text{message_buffer}, \text{requests}, \text{initialized}, \text{collective} \rangle$
 $\wedge \text{UNCHANGED } \text{Memory}$

Again we ignore the *buffer_addr* argument as we are abstracting data.

The standard does not indicate what the result is when there is no buffer currently attached.

$\text{MPI_Buffer_detach}(\text{buffer_addr}, \text{size}, \text{return}, \text{proc}) \triangleq$ < 34.36 – 35.2 >
 $\wedge \text{Assert}(\text{initialized}[\text{proc}] = \text{"initialized"},$ 200.10 – 200.12
 "Error: MPI_Buffer_detach called with proc not in initialized state.")
 $\wedge \text{Assert}(\text{bufsize}[\text{proc}] \neq 0,$
 "Error: MPI_Buffer_detach called when no buffer is currently associated with this process.")
 $\wedge \text{bufsize}' = [\text{bufsize} \text{ EXCEPT } ![\text{proc}] = 0]$ 34.46
 $\wedge \forall j \in 1 \dots \text{Len}(\text{requests}[\text{proc}]) :$ 34.47
 $\text{requests}[\text{proc}][j].\text{ctype} = \text{"bsend"} \Rightarrow \text{requests}[\text{proc}][j].\text{transmitted}$
 < /34.36 – 35.2 >
 $\wedge \text{Memory} = [\text{Memory} \text{ EXCEPT } ![\text{proc}] = [\text{@} \text{ EXCEPT } ![\text{size}] = \text{bufsize}[\text{proc}]]]$ 34.47
 $\wedge \text{UNCHANGED} \langle \text{group}, \text{communicator}, \text{message_buffer}, \text{requests}, \text{initialized}, \text{collective} \rangle$

Section 3.7.2 Communication initiation

Notes: I am not sure how to model this construct. The main problem lies in the nondeterministic buffering scheme that the standard refers to. For a correct program one must expect no buffering, however is it possible to write a program in such a way as to require synchronous handshakes?

Start a non-blocking standard send. 38.17 – 38.35, 58.13 – 58.18

$\text{MPI_Isend}(\text{buf}, \text{count}, \text{datatype}, \text{dest}, \text{tag}, \text{comm}, \text{request}, \text{return}, \text{proc}) \triangleq$
 $\wedge \text{Assert}(\text{initialized}[\text{proc}] = \text{"initialized"},$ 200.10 – 200.12
 "Error: MPI_Isend called with proc not in initialized state.")
 $\wedge \text{Assert}(\text{proc} \in \text{group}[\text{proc}][\text{communicator}[\text{proc}][\text{comm}].\text{group}].\text{members},$
 "Error: MPI_Isend called on a communicator which this process is not a member of.")
 $\wedge \text{LET } \text{msg} \triangleq$
 $[\text{addr} \quad \mapsto \text{buf},$
 $\text{src} \quad \mapsto \text{group}[\text{proc}][\text{communicator}[\text{proc}][\text{comm}].\text{group}].\text{ranking}[\text{proc}],$
 $\text{dest} \quad \mapsto \text{dest},$
 $\text{msgtag} \quad \mapsto \text{tag},$
 $\text{dtype} \quad \mapsto \text{datatype},$
 $\text{numelements} \mapsto \text{count},$
 $\text{universe} \quad \mapsto \text{comm},$
 $\text{state} \quad \mapsto \text{"send"}]$
 IN
 $\text{requests}' = [\text{requests} \text{ EXCEPT } ![\text{proc}] =$ 40.40, 35.37 – 35.39
 $\text{@} \circ \langle \text{Make_request}(0, \text{TRUE}, \text{FALSE}, \text{TRUE}, \text{TRUE}, \text{"send"}, \text{FALSE}, \langle \rangle, \text{FALSE}, \text{msg}) \rangle]$
 $\wedge \text{Memory}' = [\text{Memory} \text{ EXCEPT } ![\text{proc}] = [\text{@} \text{ EXCEPT } ![\text{request}] = \text{Len}(\text{requests}[\text{proc}]) + 1]]$ 40.41

\wedge UNCHANGED $\langle group, communicator, bufsize, message_buffer, initialized, collective \rangle$

Set up a non-blocking buffered send. 39.1 – 39.19, 58.13 – 58.18

$MPI_Ibsend(buf, count, datatype, dest, tag, comm, request, return, proc) \triangleq$
 \wedge Assert($initialized[proc] = \text{"initialized"}$, 200.10 – 200.12
 “Error: MPI_Ibsend called with proc not in initialized state.”)
 \wedge Assert($message_buffer[proc] < bufsize[proc]$, 28.6, 35.34 – 35.35
 “Error: MPI_Ibsend called when insufficient buffering was available.”)
 \wedge Assert($proc \in group[proc][communicator[proc][comm].group].members$,
 “Error: MPI_Ibsend called on a communicator which this process is not a member of.”)
 \wedge LET $msg \triangleq$
 $[addr \quad \mapsto buf,$
 $src \quad \mapsto group[proc][communicator[proc][comm].group].ranking[proc],$
 $dest \quad \mapsto dest,$
 $msgtag \quad \mapsto tag,$
 $dtype \quad \mapsto datatype,$
 $numelements \mapsto count,$
 $universe \quad \mapsto comm,$
 $state \quad \mapsto \text{"send"}]$
 IN
 $requests' = [requests \text{ EXCEPT } ![proc] =$ 40.40
 $\quad @ \circ \langle Make_request(0, TRUE, FALSE, TRUE, TRUE, \text{"bsend"}, FALSE, \langle \rangle, FALSE, msg) \rangle]$
 $\wedge Memory' = [Memory \text{ EXCEPT } ![proc] = [@ \text{ EXCEPT } ![request] = Len(requests[proc]) + 1]]$ 40.41
 $\wedge message_buffer' = [message_buffer \text{ EXCEPT } ![proc] = @ + 1]$ 28.6 Consume necessary buffer space
 \wedge UNCHANGED $\langle group, communicator, bufsize, initialized, collective \rangle$

Tested

Set up a non-blocking synchronous send. 39.21 – 39.39, 58.13 – 58.18

$MPI_Issend(buf, count, datatype, dest, tag, comm, request, return, proc) \triangleq$
 \wedge Assert($initialized[proc] = \text{"initialized"}$, 200.10 – 200.12
 “Error: MPI_Issend called with proc not in initialized state.”)
 \wedge Assert($proc \in group[proc][communicator[proc][comm].group].members$,
 “Error: MPI_Issend called on a communicator which this process is not a member of.”)
 \wedge LET $msg \triangleq$
 $[addr \quad \mapsto buf,$
 $src \quad \mapsto group[proc][communicator[proc][comm].group].ranking[proc],$
 $dest \quad \mapsto dest,$
 $msgtag \quad \mapsto tag,$
 $dtype \quad \mapsto datatype,$
 $numelements \mapsto count,$
 $universe \quad \mapsto comm,$
 $state \quad \mapsto \text{"send"}]$
 IN
 $requests' = [requests \text{ EXCEPT } ![proc] =$ 40.40
 $\quad @ \circ \langle Make_request(0, TRUE, FALSE, TRUE, FALSE, \text{"ssend"}, FALSE, \langle \rangle, FALSE, msg) \rangle]$

\wedge $Memory' = [Memory \text{ EXCEPT } ![proc] = [@ \text{ EXCEPT } ![request] = Len(requests[proc]) + 1]]$ 40.41
 \wedge UNCHANGED $\langle group, communicator, message_buffer, bufsize, initialized, collective \rangle$

Set up a non-blocking ready send. 40.1 – 40.19, 58.13 – 58.18

$MPI_Irsend(buf, count, datatype, dest, tag, comm, request, return, proc) \triangleq$
 \wedge $Assert(initialized[proc] = \text{"initialized"},$ 200.10 – 200.12
 "Error: MPI_Irsend called with proc not in initialized state.")
 \wedge $Assert(proc \in group[proc][communicator[proc][comm].group].members,$
 "Error: MPI_Irsend called on a communicator which this process is not a member of.")
 \wedge $Assert(\exists k \in (1 \dots Len(requests[dest])) :$ 37.6 – 37.8
 $\wedge requests[dest][k].active$
 $\wedge \neg requests[dest][k].transmitted$
 $\wedge \neg requests[dest][k].cancelled$
 $\wedge Match(requests[proc][request].message, requests[dest][k].message),$
 "Error: MPI_Start tried to start a rsend request when no matching message exists.")
 $\wedge requests' = [requests \text{ EXCEPT } ![proc] =$ 40.40
 LET $msg \triangleq$
 $[addr \quad \mapsto buf,$
 $src \quad \mapsto group[proc][communicator[proc][comm].group].ranking[proc],$
 $dest \quad \mapsto dest,$
 $msgtag \quad \mapsto tag,$
 $dtype \quad \mapsto datatype,$
 $numelements \mapsto count,$
 $universe \quad \mapsto comm,$
 $state \quad \mapsto \text{"send"}]$
 IN
 $@ \circ \langle Make_request(0, TRUE, FALSE, TRUE, FALSE, \text{"rsend"}, FALSE, \langle \rangle, FALSE, msg) \rangle$
 \wedge $Memory' = [Memory \text{ EXCEPT } ![proc] = [@ \text{ EXCEPT } ![request] = Len(requests[proc]) + 1]]$ 40.41
 \wedge UNCHANGED $\langle group, communicator, bufsize, message_buffer, requests, initialized, collective \rangle$

Set up a non-blocking receive. 40.21 – 40.39, 58.13 – 58.18

$MPI_Irecv(buf, count, datatype, source, tag, comm, request, return, proc) \triangleq$
 \wedge $Assert(initialized[proc] = \text{"initialized"},$ 200.10 – 200.12
 "Error: MPI_Irecv called with proc not in initialized state.")
 \wedge $Assert(proc \in group[proc][communicator[proc][comm].group].members,$
 "Error: MPI_Irecv called on a communicator which this process is not a member of.")
 \wedge LET $msg \triangleq$
 $[addr \quad \mapsto buf,$
 $src \quad \mapsto source,$
 $dest \quad \mapsto group[proc][communicator[proc][comm].group].ranking[proc],$
 $msgtag \quad \mapsto tag,$
 $dtype \quad \mapsto datatype,$
 $numelements \mapsto count,$
 $universe \quad \mapsto comm,$

$state \mapsto \text{"recv"}$
 IN
 $requests' = [requests \text{ EXCEPT } ![proc] = \text{40.40}$
 $\quad @ \circ \langle Make_request(0, \text{TRUE}, \text{FALSE}, \text{TRUE}, \text{FALSE}, \text{"recv"}, \text{FALSE}, \langle \rangle, \text{FALSE}, msg) \rangle$
 $\wedge Memory' = [Memory \text{ EXCEPT } ![proc] = [@ \text{ EXCEPT } ![request] = Len(requests[proc]) + 1]] \text{40.41}$
 $\wedge \text{UNCHANGED } \langle group, communicator, message_buffer, bufsize, initialized, collective \rangle$

Section 3.7.3 Communication Completion

Would if...then...else be a better, more readable form here? Maybe not because we need to block.

Wait for request to complete. Return information about the message in status. 41.23 – 42.6

No specification on what the status value is when a send is posted with *MPI_PROC_NULL*

Specifies next state for status and request

$MPI_Wait(request, status, return, proc) \triangleq$

LET $r \triangleq requests[proc][Memory[proc][request]]$ IN

$\wedge Assert(initialized[proc] = \text{"initialized"}, \text{200.10} - \text{200.12}$

$\quad \text{"Error: MPI.Wait called with proc not in initialized state."})$

$\wedge \vee \wedge Memory[proc][request] \neq MPI_REQUEST_NULL \text{41.32} - \text{41.39}$ The request handle is not the null handle

$\wedge r.active$

The request is active.

$\wedge \vee \wedge r.message.src \neq MPI_PROC_NULL$

The message source is not null

$\wedge r.message.dest \neq MPI_PROC_NULL$

The message destination is not null

41.32 – Blocks until complete

$\wedge \vee r.transmitted$

The communication actually happened or

$\vee r.cancelled$

the communication got *cancelled* by the user program or

$\vee r.buffered$

the communication got buffered either into explicit user provided

buffer space or into system provided buffer space (if regular send is used).

A status object for a completed communication.

$\wedge Memory' =$

$[Memory \text{ EXCEPT } ![proc] = \text{41.36}$

$\quad [@ \text{ EXCEPT } ![Status_Cancelled(status)] = r.cancelled \wedge \neg r.transmitted, \text{54.46}$

$\quad ![Status_Count(status)] = r.message.numelements,$

$\quad ![Status_Source(status)] = r.message.src,$

$\quad ![Status_Tag(status)] = r.message.msgtag,$

$\quad ![Status_Err(status)] = r.error,$

$\quad ![request] = \text{IF } r.persist \text{ THEN } @ \text{ ELSE } MPI_REQUEST_NULL]] \text{41.32} - \text{41.35,}$

$\vee \wedge \vee r.message.src = MPI_PROC_NULL$

The source or destination was actually

$\vee r.message.dest = MPI_PROC_NULL$

the null process

$\wedge Memory' = [Memory \text{ EXCEPT } ![proc] = \text{41.36}$

$\quad [@ \text{ EXCEPT } ![Status_Cancelled(status)] = r.cancelled,$

$\quad ![Status_Count(status)] = 0,$

$\quad ![Status_Source(status)] = MPI_PROC_NULL,$

$\quad ![Status_Tag(status)] = MPI_ANY_TAG,$

$$\begin{aligned}
& \text{!}[flag] = \text{TRUE}, \\
& \text{!}[request] = \text{IF } r.\text{persist} \text{ THEN } @ \text{ ELSE } \text{MPI_REQUEST_NULL}] \\
\wedge \text{requests}' = & \\
& [\text{requests EXCEPT !}[proc] = \\
& \quad [@ \text{ EXCEPT !}[Memory[proc][request]] = \\
& \quad \quad [@ \text{ EXCEPT !.active} = \text{FALSE}]]] \quad 42.22 - 42.23, 58.34 \text{ Not modeling deallocation} \\
\text{ELSE} & \quad 42.23 - 42.24 \\
& \wedge \text{Memory}' = [\text{Memory EXCEPT !}[proc] = [@ \text{ EXCEPT !}[flag] = \text{FALSE}]] \quad \text{status is undefined 4} \\
& \wedge \text{UNCHANGED } \langle \text{requests} \rangle \\
\vee \wedge \vee r.\text{message.src} = \text{MPI_PROC_NULL} & \quad \text{The source or destination were actually} \\
\vee r.\text{message.dest} = \text{MPI_PROC_NULL} & \quad \text{the null process 42.29 - 42.31} \\
\wedge \text{Memory}' = [\text{Memory EXCEPT !}[proc] = \\
& \quad [@ \text{ EXCEPT !}[Status_Cancelled(status)] = \text{FALSE}, \\
& \quad \quad \text{!}[Status_Count(status)] = 0, \\
& \quad \quad \text{!}[Status_Source(status)] = \text{MPI_PROC_NULL}, \\
& \quad \quad \text{!}[Status_Tag(status)] = \text{MPI_ANY_TAG}, \\
& \quad \quad \text{!}[Status_Err(status)] = 0, \\
& \quad \quad \text{!}[flag] = \text{TRUE}, \\
& \quad \quad \text{!}[request] = \text{IF } r.\text{persist} \text{ THEN } @ \text{ ELSE } \text{MPI_REQUEST_NULL}] \\
\wedge \text{requests}' = & \\
& [\text{requests EXCEPT !}[proc] = \\
& \quad [@ \text{ EXCEPT !}[Memory[proc][request]] = \\
& \quad \quad [@ \text{ EXCEPT !.active} = \text{FALSE}]]] \quad 42.22 - 42.23, 58.34 \text{ Not modeling deallocation} \\
\vee \wedge \vee \neg r.\text{active} & \quad \text{The request is not active or the request} \\
\vee \text{Memory}[proc][request] = \text{MPI_REQUEST_NULL} & \quad \text{handle is null 42.29 - 42.31} \\
\wedge \text{Memory}' = [\text{Memory EXCEPT !}[proc] = \\
& \quad [@ \text{ EXCEPT !}[Status_Cancelled(status)] = \text{FALSE}, \\
& \quad \quad \text{!}[Status_Count(status)] = 0, \\
& \quad \quad \text{!}[Status_Source(status)] = \text{MPI_ANY_SOURCE}, \\
& \quad \quad \text{!}[Status_Tag(status)] = \text{MPI_ANY_TAG}, \\
& \quad \quad \text{!}[Status_Err(status)] = 0, \\
& \quad \quad \text{!}[flag] = \text{TRUE}]] \\
& \wedge \text{UNCHANGED } \langle \text{requests} \rangle \\
& \wedge \text{UNCHANGED } \langle \text{group, communicator, bufsize, message_buffer, initialized, collective} \rangle \\
& \text{Frees the request specified.} \\
& \text{Modifies request.} \\
\text{MPI_Request_free}(request, return, proc) \triangleq & \\
& \wedge \text{Assert}(\text{initialized}[proc] = \text{"initialized"}, \quad 200.10 - 200.12 \\
& \quad \text{"MPI_Request_free called with proc not in initialized state."}) \\
& \wedge \text{Assert}(\neg \text{requests}[proc][\text{Memory}[proc][request]].\text{active}, \quad 43.37 - 43.39 \\
& \quad \text{"MPI_Request_free called with an inactive request."}) \\
& \wedge \text{Memory}' = [\text{Memory EXCEPT !}[proc] = [@ \text{ EXCEPT !}[request] = \text{MPI_REQUEST_NULL}]] \quad 43.20 \text{ Not } \\
& \wedge \text{UNCHANGED } \langle \text{group, communicator, bufsize, message_buffer, requests, initialized, collective} \rangle
\end{aligned}$$

Section 3.7.5 Multiple Completions

Wait for one of the requests referenced in *array_of_requests* to complete.

Specifies next state for index and status

$MPI_Waitany(count, array_of_requests, index, status, return, proc) \triangleq$

LET $r(v) \triangleq requests[proc][Memory[proc][array_of_requests + v]]$ IN

$\wedge Assert(initialized[proc] = \text{"initialized"}, 200.10 - 200.12$

$\text{"Error: MPI_Waitany called with proc not in initialized state."})$

$\wedge \vee \exists i \in (0 .. (count - 1)) :$ 45.44 – 45.46 Blocks—chooses arbitrarily one that satisfies the following:

$\wedge Memory[proc][array_of_requests + i] \neq MPI_REQUEST_NULL$ The request handle is not the null

$\wedge r(i).active$ The request is active.

$\wedge \vee \wedge r(i).message.src \neq MPI_PROC_NULL$ The message source is not null

$\wedge r(i).message.dest \neq MPI_PROC_NULL$ The message destination is not null

$\wedge \vee r(i).transmitted$ The communication actually happened or

$\vee r(i).cancelled$ the communication got *cancelled* by the user program or

$\vee r(i).buffered$ the communication got buffered either into explicit user provided buffer space or into system provided buffer space (if regular send is used).

$\wedge Memory' = [Memory \text{ EXCEPT } ![proc] = 45.46 - 45.47$

$[@ \text{ EXCEPT}$

$![Status_Source(status)] = r(i).message.src, 45.47 - 45.48$

$![Status_Tag(status)] = r(i).message.msgtag,$

$![Status_Err(status)] = r(i).error,$

$![Status_Count(status)] = r(i).message.numelements,$

$![Status_Cancelled(status)] = r(i).cancelled \wedge \neg r(i).transmitted, 54.46$

$![array_of_requests + i] = \text{IF } r(i).persist \text{ THEN } @ \text{ ELSE } MPI_REQUEST_NULL, 46.1 -$

$![index] = i]] 45.46$

$\vee \wedge \vee r(i).message.src = MPI_PROC_NULL$ The source or destination was actually

$\vee r(i).message.dest = MPI_PROC_NULL$ the null process

$\wedge Memory' = [Memory \text{ EXCEPT } ![proc] = [@ \text{ EXCEPT}$

$![Status_Source(status)] = MPI_PROC_NULL,$

$![Status_Tag(status)] = MPI_ANY_TAG,$

$![Status_Err(status)] = 0,$

$![Status_Count(status)] = 0,$

$![Status_Cancelled(status)] = r(i).cancelled,$

$![array_of_requests + i] = \text{IF } r(i).persist \text{ THEN } @ \text{ ELSE } MPI_REQUEST_NULL, 46.2 -$

$![index] = i]] 45.46$

$\wedge requests' = [requests \text{ EXCEPT } ![proc] = 46.1, 58.34$

$[@ \text{ EXCEPT } ![Memory[proc][array_of_requests + i]] =$

$[@ \text{ EXCEPT } !.active = FALSE]]]$

$\vee \forall i \in (0 .. (count - 1)) :$ 46.3 – 46.4

$\wedge \vee \neg r(i).active$ The request is not active or the request

$\vee Memory[proc][array_of_requests + i] = MPI_REQUEST_NULL$ handle is null

$\wedge Memory' = [Memory \text{ EXCEPT } ![proc] = 46.5$

[@ EXCEPT
 ! $[\text{Status_Source}(\text{status})] = \text{MPI_ANY_SOURCE}$,
 ! $[\text{Status_Tag}(\text{status})] = \text{MPI_ANY_TAG}$,
 ! $[\text{Status_Err}(\text{status})] = 0$,
 ! $[\text{Status_Count}(\text{status})] = 0$,
 ! $[\text{Status_Cancelled}(\text{status})] = \text{FALSE}$,
 ! $[\text{index}] = \text{MPI_UNDEFINED}]$ 46.5
 \wedge UNCHANGED $\langle \text{requests} \rangle$
 \wedge UNCHANGED $\langle \text{group}, \text{communicator}, \text{bufsize}, \text{message_buffer}, \text{initialized}, \text{collective} \rangle$

Test whether one of the requests referenced in array_of_requests has completed.
 $\text{MPI_Testany}(\text{count}, \text{array_of_requests}, \text{index}, \text{flag}, \text{status}, \text{return}, \text{proc}) \triangleq$
 LET $r(v) \triangleq \text{requests}[\text{proc}][\text{Memory}[\text{proc}][\text{array_of_requests} + v]]$ IN
 $\wedge \text{Assert}(\text{initialized}[\text{proc}] = \text{"initialized"})$, 200.10 – 200.12
 “Error: MPI_Testany called with proc not in initialized state.”
 $\wedge \vee \exists i \in (0 \dots (\text{count} - 1))$: 46.28 – 46.29
 $\wedge \text{array_of_requests}[\text{proc}][i] \neq \text{MPI_REQUEST_NULL}$ The request handle is not the null handle.
 $\wedge r(i).active$ The request is active.
 $\wedge \vee \wedge r(i).message.src \neq \text{MPI_PROC_NULL}$ The message source is not null
 $\wedge r(i).message.dest \neq \text{MPI_PROC_NULL}$ The message destination is not null
 \wedge IF $\vee r(i).transmitted$ The communication actually happened or
 $\vee r(i).cancelled$ the communication got *cancelled* by the user program or
 $\vee r(i).buffered$ the communication got buffered either into explicit user provided
 buffer space or into system provided buffer space (if regular send is used).
 THEN
 $\wedge \text{Memory}' = [\text{Memory} \text{ EXCEPT } ![\text{proc}] = [@ \text{ EXCEPT}$
 ! $[\text{flag}] = \text{TRUE}$, 46.29
 ! $[\text{Status_Source}(\text{status})] = r(i).message.src$, 46.30
 ! $[\text{Status_Tag}(\text{status})] = r(i).message.msgtag$,
 ! $[\text{Status_Err}(\text{status})] = r(i).error$,
 ! $[\text{Status_Count}(\text{status})] = r(i).message.numelements$,
 ! $[\text{Status_Cancelled}(\text{status})] = r(i).cancelled \wedge \neg r(i).transmitted$, 54.46
 ! $[\text{index}] = i$, 46.29
 ! $[\text{array_of_requests} + i] = \text{IF } r(i).persist \text{ THEN } @ \text{ ELSE } \text{MPI_REQUEST_NULL}]$ 46.30 – 46.31, 58.34
 $\wedge \text{requests}' =$
 $[\text{requests} \text{ EXCEPT } ![\text{proc}] =$
 $[@ \text{ EXCEPT } ![\text{Memory}[\text{proc}][\text{array_of_requests} + i]] =$
 $[@ \text{ EXCEPT } !.active = \text{FALSE}]]]$
 ELSE
 $\wedge \text{Memory}' = [\text{Memory} \text{ EXCEPT } ![\text{proc}] = [@ \text{ EXCEPT}$ status is explicitly undefined.
 ! $[\text{flag}] = \text{FALSE}$, 46.33
 ! $[\text{index}] = \text{MPI_UNDEFINED}]$ 46.33 – 46.34
 \wedge UNCHANGED $\langle \text{requests} \rangle$
 $\vee \wedge \vee r(i).message.src = \text{MPI_PROC_NULL}$ The source or destination were actually

$$\begin{aligned}
& \vee r(i).message.dest = MPI_PROC_NULL \quad \text{the null process 61.3 – 61.4} \\
\wedge Memory' = [Memory \text{ EXCEPT } ![proc] = [@ \text{ EXCEPT} \\
& \quad ![flag] = \text{TRUE}, \quad 46.29 \\
& \quad ![Status_Source(status)] = MPI_PROC_NULL, \\
& \quad ![Status_Tag(status)] = MPI_ANY_TAG, \\
& \quad ![Status_Err(status)] = 0, \\
& \quad ![Status_Count(status)] = 0, \\
& \quad ![Status_Cancelled(status)] = r(i).cancelled, \\
& \quad ![index] = i, \quad 46.29 \\
& \quad ![array_of_requests + i] = \text{IF } r(i).persist \text{ THEN } @ \text{ ELSE } MPI_REQUEST_NULL]] \quad 46.31 - \\
\wedge requests' = \quad 46.31 - 46.32, 58.34 \\
& \quad [requests \text{ EXCEPT } ![proc] = \\
& \quad \quad [@ \text{ EXCEPT } ![Memory[proc][array_of_requests + i]] = \\
& \quad \quad \quad [@ \text{ EXCEPT } !.active = \text{FALSE}]]] \\
\vee \forall i \in (0 \dots (count - 1)) : \quad 46.35 - 46.37 \\
& \quad \wedge \vee \neg r(i).active \quad \text{The request is not active or the request} \\
& \quad \quad \vee array_of_requests[proc][i] = MPI_REQUEST_NULL \quad \text{handle is null} \\
\wedge Memory' = [Memory \text{ EXCEPT } ![proc] = [@ \text{ EXCEPT} \\
& \quad ![flag] = \text{TRUE}, \quad 46.36 \\
& \quad ![Status_Source(status)] = MPI_ANY_SOURCE, \quad 46.36 \\
& \quad ![Status_Tag(status)] = MPI_ANY_TAG, \\
& \quad ![Status_Err(status)] = 0, \\
& \quad ![Status_Count(status)] = 0, \\
& \quad ![Status_Cancelled(status)] = \text{FALSE}, \\
& \quad ![index] = MPI_UNDEFINED]] \quad 46.36 \\
& \quad \wedge \text{UNCHANGED } \langle requests \rangle \\
\wedge \text{UNCHANGED } \langle group, communicator, bufsize, message_buffer, initialized, collective \rangle
\end{aligned}$$

A long version of *MPI_Waitall* – includes the line by line reference.

Specifies the next state for *array_of_requests* and *array_of_statuses*.

$$\begin{aligned}
& MPI_Waitall(count, array_of_requests, array_of_statuses, return, proc) \triangleq \\
& \text{LET } r(v) \triangleq requests[proc][Memory[proc][array_of_requests + v]] \text{ IN} \\
& \quad \wedge \text{Assert}(initialized[proc] = \text{"initialized"}, \quad 200.10 - 200.12 \\
& \quad \quad \text{"Error: MPI_Waitall called with proc not in initialized state."}) \\
& \quad \wedge \forall i \in (0 \dots (count - 1)) : \quad 47.18 \\
& \quad \quad \wedge Memory[proc][array_of_requests + i] \neq MPI_REQUEST_NULL \quad \text{The request handle is not the null} \\
& \quad \quad \wedge r(i).active \quad \text{The request is active.} \\
& \quad \quad \wedge \vee \wedge r(i).message.src \neq MPI_PROC_NULL \quad \text{The message source is not null} \\
& \quad \quad \quad \wedge r(i).message.dest \neq MPI_PROC_NULL \quad \text{The message destination is not null} \\
& \quad \quad \quad \wedge \vee r(i).transmitted \quad \text{The communication actually happened or} \\
& \quad \quad \quad \quad \vee r(i).cancelled \quad \text{the communication got } cancelled \text{ by the user program or} \\
& \quad \quad \quad \quad \vee r(i).buffered \quad \text{the communication got buffered either into explicit user provided} \\
& \quad \quad \quad \quad \quad \text{buffer space or into system provided buffer space (if regular send is used).} \\
& \quad \quad \vee \wedge \vee r(i).message.src = MPI_PROC_NULL \quad \text{The source or destination was actually}
\end{aligned}$$

```

     $\forall r(i).message.dest = MPI\_PROC\_NULL$  the null process
 $\wedge$   $array\_of\_requests' =$ 
     $[array\_of\_requests$  EXCEPT  $![proc] =$  47.22 – 47.23
     $[j \in 0 .. (count - 1) \mapsto$ 
    IF  $r(j).persist$ 
    THEN
     $array\_of\_requests[proc][j]$ 
    ELSE
     $MPI\_REQUEST\_NULL]$ 
 $\wedge$   $array\_of\_statuses' = [array\_of\_statuses$  EXCEPT  $![proc] =$  47.18 – 47.21
 $[j \in (0 .. count - 1) \mapsto$ 
    IF  $\forall r(j).message.src = MPI\_PROC\_NULL$ 
     $\forall r(j).message.dest = MPI\_PROC\_NULL$ 
    THEN 61.3 – 61.4
     $[state \mapsto$  "defined", A status object for a communication
     $MPI\_SOURCE \mapsto MPI\_PROC\_NULL,$  with a null process
     $MPI\_TAG \mapsto MPI\_ANY\_TAG,$ 
     $MPI\_ERROR \mapsto 0,$ 
     $count \mapsto 0,$ 
     $cancelled \mapsto r(j).cancelled]$ 
    ELSE
     $[state \mapsto$  "defined", A status object for a completed communication.
     $MPI\_SOURCE \mapsto r(j).message.src,$ 
     $MPI\_TAG \mapsto r(j).message.msgtag,$ 
     $MPI\_ERROR \mapsto r(j).error,$ 
     $count \mapsto r(j).message.numelements,$ 
     $cancelled \mapsto r(j).cancelled \wedge \neg r(i).transmitted]]$  54.46
 $\wedge$   $requests' = [requests$  EXCEPT  $![proc] =$  47.22, 58.34 Not modeling deallocation
 $[j \in 1 .. Len(@) \mapsto$ 
    IF  $\exists k \in 0 .. (count - 1) : j = array\_of\_requests[proc][k]$ 
    THEN
     $[requests[proc][j]$  EXCEPT  $!.active = FALSE]$ 
    ELSE
     $requests[proc][j]]]$ 
 $\forall \forall i \in 0 .. (count - 1) :$  47.23 – 47.24
 $\wedge \forall Memory[proc][array\_of\_requests + i] = MPI\_REQUEST\_NULL$  The request handle is null or
 $\forall \neg r(i).active$  not active
 $\wedge$   $Memory' = [Memory$  EXCEPT  $![proc] =$ 
 $[j \in 1 .. Len(Memory[proc]) \mapsto$ 
    IF  $j \in array\_of\_statuses .. (array\_of\_statuses + ((count * 5) - 1))$ 
    THEN
    IF  $(j - array\_of\_statuses)\%5 = 0$ 
    THEN FALSE
    ELSE
    IF  $(j - array\_of\_statuses)\%5 = 1$ 

```

```

THEN 0
ELSE
  IF (j - array_of_statuses)%5 = 2
  THEN MPI_ANY_SOURCE
  ELSE
    IF (j - array_of_statuses)%5 = 3
    THEN MPI_ANY_TAG
    ELSE
      IF (j - array_of_statuses)%5 = 4
      THEN 0
      ELSE Assert(FALSE, "Internal Error: Cannot have any other cases.")
    ELSE Memory[proc][j]]
  ∧ UNCHANGED ⟨array_of_requests, requests⟩
∧ UNCHANGED ⟨group, communicator, bufsize, message_buffer, requests, initialized, collective⟩

```

Test whether all requests referenced in *array_of_requests* have completed.

MPI_Testall(*count*, *array_of_requests*, *flag*, *array_of_statuses*, *return*, *proc*) \triangleq

LET $r(v) \triangleq$ requests[proc][array_of_requests[proc][v]] IN

∧ Assert(initialized[proc] = "initialized", 200.10 – 200.12)

"Error: MPI_Testall called with proc not in initialized state.")

∧ IF $\forall i \in (0 .. (count - 1))$: 48.15

$\vee \wedge$ array_of_requests[proc][i] \neq MPI_REQUEST_NULL The request handle is not the null handle.

\wedge r(i).active The request is active.

$\wedge \vee \wedge$ r(i).message.src \neq MPI_PROC_NULL The message source is not null

\wedge r(i).message.dest \neq MPI_PROC_NULL The message destination is not null

$\wedge \vee$ r(i).transmitted The communication actually happened or

\vee r(i).cancelled the communication got cancelled by the user program or

\vee r(i).buffered the communication got buffered either into explicit user provided

buffer space or into system provided buffer space (if regular send is used).

$\vee \wedge \vee$ r(i).message.src = MPI_PROC_NULL The source or destination were actually

\vee r(i).message.dest = MPI_PROC_NULL the null process

$\vee \forall i \in (0 .. (count - 1))$: 48.16

\vee array_of_requests[proc][i] = MPI_REQUEST_NULL

$\vee \neg$ r(i).active

THEN

∧ array_of_statuses' = [array_of_statuses EXCEPT ![proc] =

[i \in (0 .. (count - 1)) \mapsto

IF \vee r(i).message.src = MPI_PROC_NULL

\vee r(i).message.dest = MPI_PROC_NULL

THEN

[state \mapsto "defined", 61.3 – 61.4

MPI_SOURCE \mapsto MPI_PROC_NULL,

MPI_TAG \mapsto MPI_ANY_TAG,

MPI_ERROR \mapsto 0,

A status object for a communication
with a null process

```

    count          ↦ 0,
    cancelled      ↦ FALSE]
ELSE
IF  $\forall$  array_of_requests[proc][i] = MPI_REQUEST_NULL 48.21
 $\vee \neg r(i).active$ 
THEN
[state          ↦ "empty",           The resultant empty status.
 MPI_SOURCE ↦ MPI_ANY_SOURCE,
 MPI_TAG     ↦ MPI_ANY_TAG,
 MPI_ERROR   ↦ 0,
 count     ↦ 0,
 cancelled ↦ FALSE]
ELSE 48.17 – 48.18
[state          ↦ "defined",         A status object for a completed communication.
 MPI_SOURCE ↦ r(i).message.src,
 MPI_TAG     ↦ r(i).message.msgtag,
 MPI_ERROR   ↦ r(i).error,
 count     ↦ r(i).message.numelements,
 cancelled ↦ r(i).cancelled  $\wedge \neg r(i).transmitted$ ] 54.46
 $\wedge$  requests' = 48.18 – 48.19, 58.34 Not modeling deallocation
[requests EXCEPT ![proc] =
 [i  $\in$  1 .. Len(@)  $\mapsto$ 
 IF  $\exists j \in 0 \dots (count - 1) : array\_of\_requests[proc][j] = i$ 
 THEN
 [requests[proc][i] EXCEPT !.active = FALSE]
 ELSE
 requests[proc][i]]
 $\wedge$  array_of_requests' = [array_of_requests EXCEPT ![proc] =
 [i  $\in$  0 .. (count - 1)  $\mapsto$ 
 IF r(i).persist
 THEN
 array_of_requests[proc][i] 58.34 – 58.35
 ELSE
 MPI_REQUEST_NULL] 48.19 – 48.21
 $\wedge$  flag' = [flag EXCEPT ![proc] = TRUE] 48.15
ELSE
 $\wedge$  flag' = [flag EXCEPT ![proc] = FALSE] 48.21 – 48.22
 $\wedge$  array_of_statuses' = [array_of_statuses EXCEPT ![proc] =
 [i  $\in$  0 .. (count - 1)  $\mapsto$ 
 [array_of_statuses[proc][i] EXCEPT !.state = "undefined"]]
 $\wedge$  UNCHANGED (array_of_requests, requests)
 $\wedge$  UNCHANGED (group, communicator, bufsize, message_buffer, initialized, collective)

```

Wait for some subset of the requests referenced in *array_of_requests* to complete.
The ordering of *array_of_indices* or *array_of_statuses* is not specified.
Not modeling the possibility of arbitrary ordering of the *array_of_indices* or *array_of_statuses*.

MPI_Waitsome(*incount*, *array_of_requests*, *outcount*,
array_of_indices, *array_of_statuses*, *return*, *proc*) \triangleq
LET $r(v) \triangleq$ *requests*[*proc*][*array_of_requests*[*proc*][*v*]] IN
LET $msgs \triangleq$
{ $x \in (0 .. (incount - 1))$: The messages that have completed in the *array_of_requests*
 \wedge *array_of_requests*[*proc*][*x*] \neq *MPI_REQUEST_NULL* The request handle is not the null handle.
 \wedge *r*(*x*).*active* The request is active.
 \wedge \vee *r*(*x*).*transmitted* The communication actually happened or
 \vee *r*(*x*).*cancelled* the communication got *cancelled* by the user program or
 \vee *r*(*x*).*buffered* } the communication got buffered either into explicit user provided
buffer space or into system provided buffer space (if regular send is used).
IN
 \wedge *Assert*(*initialized*[*proc*] = "initialized", 200.10 – 200.12
"Error: MPI_Waitsome called with proc not in initialized state.")
 \wedge \vee \wedge *Cardinality*(*msgs*) > 0 48.45
 \wedge *outcount*' = [*outcount* EXCEPT ![*proc*] = *Cardinality*(*msgs*)] 48.46
 \wedge \exists *seq* \in *Seq*(*msgs*) : from *FiniteSets.tla* module!
 \wedge \forall *s* \in *msgs* :
 \exists *n* \in 1 .. *Len*(*seq*) :
 \wedge *seq*[*n*] = *s*
 \wedge \forall *m* \in 1 .. *Len*(*seq*) : *seq*[*n*] = *seq*[*m*] \Rightarrow *m* = *n*
 \wedge *array_of_indices*' = [*array_of_indices* EXCEPT ![*proc*] =
[*i* \in 0 .. (*incount* - 1) \mapsto
IF *i* < *Len*(*seq*)
THEN *seq*[*i* + 1]
ELSE *array_of_indices*[*proc*][*i*]]]
 \wedge *array_of_statuses*' = [*array_of_statuses* EXCEPT ![*proc*] =
[*i* \in 0 .. (*incount* - 1) \mapsto
IF *i* < *Len*(*seq*)
THEN
[*state* \mapsto "defined", A status object for a completed communication.
MPI_SOURCE \mapsto *r*(*seq*[*i* + 1]).*message.src*,
MPI_TAG \mapsto *r*(*seq*[*i* + 1]).*message.msgtag*,
MPI_ERROR \mapsto *r*(*seq*[*i* + 1]).*error*,
count \mapsto *r*(*seq*[*i* + 1]).*message.numelements*,
cancelled \mapsto *r*(*seq*[*i* + 1]).*cancelled* \wedge \neg *r*(*seq*[*i* + 1]).*transmitted*] 54.46
ELSE
array_of_statuses[*proc*][*i*]]]
 \wedge *requests*' = [*requests* EXCEPT ![*proc*] =
[*i* \in 1 .. *Len*(*requests*[*proc*]) \mapsto
IF \exists *m* \in *msgs* : *i* = *array_of_requests*[*proc*][*m*]]

```

    THEN [r(i) EXCEPT !.active = FALSE]
    ELSE r(i)]
  ∧ array_of_requests' = [array_of_requests EXCEPT ![proc] =
    [i ∈ 0 .. (incount - 1) ↦ 49.2 - 49.4
    IF ∧ ∃ m ∈ msgs : i = array_of_requests[proc][m]
    ∧ r(i).persist
    THEN
      array_of_requests[proc][i]
    ELSE
      MPI_REQUEST_NULL]]
  ∨ ∧ ∀ i ∈ (0 .. (incount - 1)) : 49.5
    ∨ array_of_requests[proc][i] = MPI_REQUEST_NULL
    ∨ ¬requests[proc][array_of_requests[proc][i]].active
    ∧ outcount' = [outcount EXCEPT ![proc] = MPI_UNDEFINED] 49.5 - 49.6
    ∧ UNCHANGED ⟨array_of_indices, array_of_statuses, requests, array_of_requests⟩
  ∧ UNCHANGED ⟨group, communicator, bufsz, message_buffer, initialized, collective⟩

```

Test for some subset of the requests referenced in the *array_of_requests* to complete.

Defined in terms of *MPI_Waitsome*.

```

MPI_Testsome(incount, array_of_requests, outcount,
  array_of_indices, array_of_statuses, return, proc) ≜
  LET r(v) ≜ requests[proc][array_of_requests[proc][v]] IN
  LET msgs ≜
    {x ∈ (0 .. (incount - 1)) : The messages that have completed in the array_of_requests
    ∧ array_of_requests[proc][x] ≠ MPI_REQUEST_NULL The request handle is not the null handle.
    ∧ r(x).active The request is active.
    ∧ ∨ r(x).transmitted The communication actually happened or
    ∨ r(x).cancelled the communication got cancelled by the user program or
    ∨ r(x).buffered} the communication got buffered either into explicit user
    buffer space or into system provided buffer space (if regular send is used).
  IN
  ∧ Assert(initialized[proc] = "initialized", 200.10 - 200.12
    "Error: MPI_Testsome called with proc not in initialized state.")
  ∧ ∨ ∃ i ∈ (0 .. (incount - 1)) : 49.35 - 49.36, 49.5
    ∧ array_of_requests[proc][i] ≠ MPI_REQUEST_NULL
    ∧ r(i).active
    ∧ IF Cardinality(msgs) > 0 number of completed messages
    THEN
      ∧ outcount' = [outcount EXCEPT ![proc] = Cardinality(msgs)] 48.46
      ∧ ∃ seq ∈ Seq(msgs) : from FiniteSets.tla module!
        ∧ ∃ s ∈ msgs :
          ∃ n ∈ 1 .. Len(seq) :
            ∧ seq[n] = s 48.47 - 49.2
            ∧ ∃ m ∈ 1 .. Len(seq) : seq[n] = seq[m] ⇒ m = n

```



```

 $\wedge$  array_of_indices' = [array_of_indices EXCEPT ![proc] =
  [j  $\in$  0 .. (incount - 1)  $\mapsto$ 
    IF j < Len(seq)
      THEN seq[j + 1]
      ELSE array_of_indices[proc][j]]]
 $\wedge$  array_of_statuses' = [array_of_statuses EXCEPT ![proc] =
  [j  $\in$  0 .. (incount - 1)  $\mapsto$ 
    IF j < Len(seq)
      THEN
        [state  $\mapsto$  "defined", A status object for a completed communication.
          MPI_SOURCE  $\mapsto$  r(seq[j + 1]).message.src,
          MPI_TAG  $\mapsto$  r(seq[j + 1]).message.msgtag,
          MPI_ERROR  $\mapsto$  r(seq[j + 1]).error,
          count  $\mapsto$  r(seq[j + 1]).message.numelements,
          cancelled  $\mapsto$  r(seq[j + 1]).cancelled  $\wedge$   $\neg$ r(seq[j + 1]).transmitted] 54.46
        ELSE
          array_of_statuses[proc][j]]]
   $\wedge$  requests' = [requests EXCEPT ![proc] =
    [j  $\in$  1 .. Len(requests[proc])  $\mapsto$ 
      IF  $\exists m \in$  msgs : j = array_of_requests[proc][m]
        THEN [r(j) EXCEPT !.active = FALSE] 58.34
        ELSE r(j)]]
   $\wedge$  array_of_requests' = [array_of_requests EXCEPT ![proc] =
    [j  $\in$  0 .. (incount - 1)  $\mapsto$  49.2 - 49.4
      IF  $\wedge \exists m \in$  msgs : j = array_of_requests[proc][m]
         $\wedge$  r(j).persist 49.2 - 49.4
      THEN
        array_of_requests[proc][j] 58.34 - 58.35
      ELSE
        MPI_REQUEST_NULL]]]
  ELSE 49.35
     $\wedge$  outcount' = [outcount EXCEPT ![proc] = 0]
     $\wedge$  UNCHANGED  $\langle$ array_of_indices, array_of_statuses, requests, array_of_requests $\rangle$ 
   $\vee$   $\wedge \forall i \in$  (0 .. (incount - 1)) : 49.5
     $\vee$  array_of_requests[proc][i] = MPI_REQUEST_NULL
     $\vee$   $\neg$ requests[proc][array_of_requests[proc][i]].active
     $\wedge$  outcount' = [outcount EXCEPT ![proc] = MPI_UNDEFINED] 49.5 - 49.6, 49.36
   $\wedge$  UNCHANGED  $\langle$ group, communicator, bufsize, message_buffer, requests, initialized, collective $\rangle$ 

```

Section 3.8 Probe and Cancel

What happens in the following scenerio: 1: send 2: probe 1: cancel 2: *recv*

Probe for a message. Nonblocking; note the leading IF

$$\begin{aligned}
 & MPI_Iprobe(source, tag, comm, flag, status, return, proc) \triangleq \\
 & \wedge \text{Assert}(initialized[proc] = \text{"initialized"}, \text{200.10} - \text{200.12}) \\
 & \quad \text{"Error: MPI_Testany called with proc not in initialized state."}) \\
 & \wedge \text{IF } \exists i \in (0 \dots (N - 1)) : \text{51.39} - \text{51.41} \\
 & \quad \exists j \in (1 \dots Len(requests[i])) : \\
 & \quad \text{LET } m \triangleq requests[i][j].message \text{ IN} \\
 & \quad \wedge \vee m.src = source \\
 & \quad \quad \vee source = MPI_ANY_SOURCE \\
 & \quad \wedge \vee m.msgtag = tag \\
 & \quad \quad \vee tag = MPI_ANY_TAG \\
 & \quad \wedge m.universe = comm \quad \text{unique across space/time - not required by standard} \\
 & \quad \wedge m.state = \text{"send"} \quad \text{51.41} - \text{51.42} \text{ must match} \\
 & \quad \wedge requests[i][j].active \quad \text{51.41} - \text{51.42} \\
 & \quad \wedge \neg requests[i][j].transmitted \\
 & \quad \wedge \neg requests[i][j].cancelled \\
 & \text{THEN} \\
 & \quad \exists i \in (0 \dots (N - 1)) : \text{51.39} - \text{51.41} \\
 & \quad \exists j \in (1 \dots Len(requests[i])) : \\
 & \quad \text{LET } m \triangleq requests[i][j].message \text{ IN} \\
 & \quad \wedge \vee m.src = source \\
 & \quad \quad \vee source = MPI_ANY_SOURCE \\
 & \quad \wedge \vee m.msgtag = tag \\
 & \quad \quad \vee tag = MPI_ANY_TAG \\
 & \quad \wedge m.universe = comm \quad \text{unique across space/time - not required by standard} \\
 & \quad \wedge m.state = \text{"send"} \quad \text{51.41} - \text{51.42} \text{ must match} \\
 & \quad \wedge requests[i][j].active \quad \text{51.41} - \text{51.42} \\
 & \quad \wedge \neg requests[i][j].transmitted \\
 & \quad \wedge \neg requests[i][j].cancelled \\
 & \quad \wedge \forall k \in (1 \dots Len(requests[i])) : \text{least match} \\
 & \quad \quad \wedge requests[i][k].active \\
 & \quad \quad \wedge \neg requests[i][k].cancelled \\
 & \quad \quad \wedge \neg requests[i][k].transmitted \\
 & \quad \quad \Rightarrow j \leq k \\
 & \wedge Memory' = [Memory \text{ EXCEPT } ![proc] = \\
 & \quad [[loc \in 1 \dots Len(Memory[proc]) \mapsto \text{51.42} \\
 & \quad \text{IF } loc = Status_Cancelled(status) \\
 & \quad \text{THEN FALSE} \\
 & \quad \text{ELSE} \\
 & \quad \text{IF } loc = Status_Count(status) \\
 & \quad \text{THEN } m.numelements \\
 & \quad \text{ELSE} \\
 & \quad \text{IF } loc = Status_Source(status) \\
 & \quad \text{THEN } m.src \\
 & \quad \text{ELSE}
 \end{aligned}$$

```

        IF  $loc = Status\_Tag(status)$ 
        THEN  $m.msgtag$ 
        ELSE
            IF  $loc = Status\_Err(status)$ 
            THEN  $requests[i][j].error$ 
            ELSE  $Memory[proc][loc]$ 
        EXCEPT  $![flag] = TRUE$ ] 51.39
    ELSE
         $\wedge Memory' = [Memory \text{ EXCEPT } ![proc] = [ @ \text{ EXCEPT } ![flag] = FALSE]]$  51.44 Status is undefined
     $\wedge UNCHANGED \langle group, communicator, bufsize, message\_buffer, requests, initialized, collective \rangle$ 

```

Wait on a probe for a message. 52.24 – 52.25

```

 $MPI\_Probe(source, tag, comm, status, return, proc) \triangleq$ 
     $\wedge Assert(initialized[proc] = \text{"initialized"},$  200.10 – 200.12
        "Error: MPI_Testany called with proc not in initialized state.")
     $\wedge \exists i \in (0 .. (N - 1)) :$ 
         $\exists j \in (1 .. Len(requests[i])) :$ 
            LET  $m \triangleq requests[i][j].messageIN$ 
             $\wedge \vee m.src = source$ 
             $\vee source = MPI\_ANY\_SOURCE$ 
             $\wedge \vee m.msgtag = tag$ 
             $\vee tag = MPI\_ANY\_TAG$ 
             $\wedge m.universe = comm$   unique across space/time – not required by standard
             $\wedge m.state = \text{"send"}$ 
             $\wedge requests[i][j].active$ 
             $\wedge \neg requests[i][j].transmitted$ 
             $\wedge \neg requests[i][j].cancelled$ 
             $\wedge \forall k \in (1 .. Len(requests[i])) :$ 
                 $\wedge requests[i][k].active$ 
                 $\wedge \neg requests[i][k].cancelled$ 
                 $\wedge \neg requests[i][k].transmitted$ 
                 $\Rightarrow j \leq k$ 
             $\wedge Memory' = [Memory \text{ EXCEPT } ![proc] =$ 
                 $[loc \in 1 .. Len(Memory[proc]) \mapsto$  51.42
                    IF  $loc = Status\_Cancelled(status)$ 
                    THEN  $requests[i][j].cancelled \wedge \neg requests[i][j].transmitted$ 
                    ELSE
                        IF  $loc = Status\_Count(status)$ 
                        THEN  $m.numelements$ 
                        ELSE
                            IF  $loc = Status\_Source(status)$ 
                            THEN  $m.src$ 
                            ELSE
                                IF  $loc = Status\_Tag(status)$ 

```

```

THEN m.msgtag
ELSE
  IF loc = Status_Err(status)
    THEN requests[i][j].error
    ELSE Memory[proc][loc]
 $\wedge$  UNCHANGED  $\langle$ group, communicator, bufsize, message_buffer, requests, initialized, collective $\rangle$ 

```

Cancel an active request.

What do you do when the request is *MPI_REQUEST_NULL*?

MPI_Cancel(request, return, proc) \triangleq 54.8 – 54.10

```

 $\wedge$  requests' = [requests EXCEPT ![proc] =
  [ @ EXCEPT ![Memory[proc][request]] =
  [ @ EXCEPT !.cancelled = TRUE]]]

```

\wedge UNCHANGED \langle *group, communicator, bufsize, message_buffer, requests, initialized, collective* \rangle

\wedge UNCHANGED \langle *Memory* \rangle

Test whether a request was *cancelled* successfully.

MPI_Test_cancelled(status, flag, return, proc) \triangleq 54.46 – 55.1

```

 $\wedge$  Memory' = [Memory EXCEPT ![proc] = [ @ EXCEPT ![flag] = Memory[proc][Status_Cancelled(status)]]]

```

\wedge UNCHANGED \langle *group, communicator, bufsize, message_buffer, requests, initialized, collective* \rangle

Section 3.9 Persistent communication requests

Create a persistent standard mode send request.

MPI_Send_init(buf, count, datatype, dest, tag, comm, request, return, proc) \triangleq

```

 $\wedge$  Assert(initialized[proc] = "initialized", 200.10 – 200.12

```

“Error: MPI_Send_init called with proc not in initialized state.”)

```

 $\wedge$  requests' = [requests EXCEPT ![proc] = 56.4 – 56.5

```

```

  LET msg  $\triangleq$  [addr  $\mapsto$  buf,
    src  $\mapsto$  group[proc][communicator[proc][comm].group].ranking[proc],
    dest  $\mapsto$  dest,
    msgtag  $\mapsto$  tag,
    dtype  $\mapsto$  datatype,
    numelements  $\mapsto$  count,
    universe  $\mapsto$  comm,
    state  $\mapsto$  “send”]

```

IN

```

  @  $\circ$   $\langle$ Make_request(0, FALSE, FALSE, FALSE, FALSE, "send", TRUE,  $\langle$  $\rangle$ , FALSE, msg) $\rangle$ 

```

57.42 – 57.46

```

 $\wedge$  Memory' = [Memory EXCEPT ![proc] = [ @ EXCEPT ![request] = Len(requests[proc]) + 1]]

```

\wedge UNCHANGED \langle *group, communicator, bufsize, message_buffer, initialized, collective* \rangle

Create a persistent buffered mode send request.

$MPI_Bsend_init(buf, count, datatype, dest, tag, comm, request, return, proc) \triangleq$
 $\wedge \text{Assert}(initialized[proc] = \text{"initialized"}, \text{200.10} - \text{200.12})$
 $\quad \text{"Error: MPI_Bsend_init called with proc not in initialized state."})$
 $\wedge requests' = [requests \text{ EXCEPT } ![proc] = \text{56.26}]$
 $\text{LET } msg \triangleq [addr \quad \mapsto buf,$
 $\quad \quad \quad src \quad \quad \mapsto group[proc][communicator[proc][comm].group].ranking[proc],$
 $\quad \quad \quad dest \quad \quad \mapsto dest,$
 $\quad \quad \quad msgtag \quad \mapsto tag,$
 $\quad \quad \quad dtype \quad \quad \mapsto datatype,$
 $\quad \quad \quad numelements \mapsto count,$
 $\quad \quad \quad universe \mapsto comm,$
 $\quad \quad \quad state \quad \quad \mapsto \text{"send"}]$
 IN
 $\quad @ \circ \langle Make_request(0, \text{FALSE}, \text{FALSE}, \text{FALSE}, \text{FALSE}, \text{"bsend"}, \text{TRUE}, \langle \rangle, \text{FALSE}, msg) \rangle$
 $\quad \quad \quad \text{57.42} - \text{57.46}$
 $\wedge Memory' = [Memory \text{ EXCEPT } ![proc] = [@ \text{ EXCEPT } ![request] = Len(requests[proc]) + 1]]$
 $\wedge \text{UNCHANGED} \langle group, communicator, bufsize, message_buffer, initialized, collective \rangle$

Create a persistent synchronous mode send request.

$MPI_Ssend_init(buf, count, datatype, dest, tag, comm, request, return, proc) \triangleq$
 $\wedge \text{Assert}(initialized[proc] = \text{"initialized"}, \text{200.10} - \text{200.12})$
 $\quad \text{"Error: MPI_Ssend_init called with proc not in initialized state."})$
 $\wedge requests' = [requests \text{ EXCEPT } ![proc] = \text{56.46}]$
 $\text{LET } msg \triangleq [addr \quad \mapsto buf,$
 $\quad \quad \quad src \quad \quad \mapsto group[proc][communicator[proc][comm].group].ranking[proc],$
 $\quad \quad \quad dest \quad \quad \mapsto dest,$
 $\quad \quad \quad msgtag \quad \mapsto tag,$
 $\quad \quad \quad dtype \quad \quad \mapsto datatype,$
 $\quad \quad \quad numelements \mapsto count,$
 $\quad \quad \quad universe \mapsto comm,$
 $\quad \quad \quad state \quad \quad \mapsto \text{"send"}]$
 IN
 $\quad @ \circ \langle Make_request(0, \text{FALSE}, \text{FALSE}, \text{FALSE}, \text{FALSE}, \text{"ssend"}, \text{TRUE}, \langle \rangle, \text{FALSE}, msg) \rangle$
 $\quad \quad \quad \text{57.42} - \text{57.46}$
 $\wedge Memory' = [Memory \text{ EXCEPT } ![proc] = [@ \text{ EXCEPT } ![request] = Len(requests[proc]) + 1]]$
 $\wedge \text{UNCHANGED} \langle group, communicator, bufsize, message_buffer, initialized, collective \rangle$

Create a persistent ready mode send request.

$MPI_Rsend_init(buf, count, datatype, dest, tag, comm, request, return, proc) \triangleq$
 $\wedge \text{Assert}(initialized[proc] = \text{"initialized"}, \text{200.10} - \text{200.12})$
 $\quad \text{"Error: MPI_Rsend_init called with proc not in initialized state."})$
 $\wedge requests' = [requests \text{ EXCEPT } ![proc] = \text{57.18}]$
 $\text{LET } msg \triangleq [addr \quad \mapsto buf,$
 $\quad \quad \quad src \quad \quad \mapsto group[proc][communicator[proc][comm].group].ranking[proc],$

$$\begin{array}{l}
\text{dest} \quad \mapsto \text{dest}, \\
\text{msgtag} \quad \mapsto \text{tag}, \\
\text{dtype} \quad \mapsto \text{datatype}, \\
\text{numelements} \mapsto \text{count}, \\
\text{universe} \quad \mapsto \text{comm}, \\
\text{state} \quad \mapsto \text{"send"}] \\
\text{IN} \\
@ \circ \langle \text{Make_request}(0, \text{FALSE}, \text{FALSE}, \text{FALSE}, \text{FALSE}, \text{"rsend"}, \text{TRUE}, \langle \rangle, \text{FALSE}, \text{msg}) \rangle \\
\text{57.42} - \text{57.46} \\
\wedge \text{Memory}' = [\text{Memory} \text{ EXCEPT } ![\text{proc}] = [@ \text{ EXCEPT } ![\text{request}] = \text{Len}(\text{requests}[\text{proc}]) + 1]] \\
\wedge \text{UNCHANGED} \langle \text{group}, \text{communicator}, \text{bufsize}, \text{message_buffer}, \text{initialized}, \text{collective} \rangle
\end{array}$$

Create a persistent receive request.

$$\begin{array}{l}
\text{MPI_Recv_init}(\text{buf}, \text{count}, \text{datatype}, \text{source}, \text{tag}, \text{comm}, \text{request}, \text{return}, \text{proc}) \triangleq \\
\wedge \text{Assert}(\text{initialized}[\text{proc}] = \text{"initialized"}, \text{200.10} - \text{200.12}) \\
\quad \text{"Error: MPI_Recv_init called with proc not in initialized state."}) \\
\wedge \text{requests}' = [\text{requests} \text{ EXCEPT } ![\text{proc}] = \text{57.39}] \\
\text{LET } \text{msg} \triangleq [\text{addr} \quad \mapsto \text{buf}, \\
\quad \text{src} \quad \mapsto \text{source}, \\
\quad \text{dest} \quad \mapsto \text{group}[\text{proc}][\text{communicator}[\text{proc}][\text{comm}].\text{group}].\text{ranking}[\text{proc}], \\
\quad \text{msgtag} \quad \mapsto \text{tag}, \\
\quad \text{dtype} \quad \mapsto \text{datatype}, \\
\quad \text{numelements} \mapsto \text{count}, \\
\quad \text{universe} \quad \mapsto \text{comm}, \\
\quad \text{state} \quad \mapsto \text{"recv"}] \\
\text{IN} \\
@ \circ \langle \text{Make_request}(0, \text{FALSE}, \text{FALSE}, \text{FALSE}, \text{FALSE}, \text{"recv"}, \text{TRUE}, \langle \rangle, \text{FALSE}, \text{msg}) \rangle \\
\text{57.42} - \text{57.46} \\
\wedge \text{Memory}' = [\text{Memory} \text{ EXCEPT } ![\text{proc}] = [@ \text{ EXCEPT } ![\text{request}] = \text{Len}(\text{requests}[\text{proc}]) + 1]] \\
\wedge \text{UNCHANGED} \langle \text{group}, \text{communicator}, \text{bufsize}, \text{message_buffer}, \text{initialized}, \text{collective} \rangle
\end{array}$$

Start a persistent communication.

What happens when a ready mode send is started and then the receive is *cancelled* before the communication has a chance to transmit?

$$\begin{array}{l}
\text{MPI_Start}(\text{request}, \text{return}, \text{proc}) \triangleq \\
\wedge \text{Assert}(\text{initialized}[\text{proc}] = \text{"initialized"}, \text{200.10} - \text{200.12}) \\
\quad \text{"Error: MPI_Start called with proc not in initialized state."}) \\
\wedge \text{Assert}(\neg \text{requests}[\text{proc}][\text{Memory}[\text{proc}][\text{request}]].\text{active}, \text{58.9}) \\
\quad \text{"Error: MPI_Start tried to start a request that is already active."}) \\
\wedge \text{Assert}(\text{Memory}[\text{proc}][\text{request}] \neq \text{MPI_REQUEST_NULL}, \\
\quad \text{"Error: MPI_Start tried to start a request that is null."}) \\
\wedge \text{Assert}(\text{requests}[\text{proc}][\text{Memory}[\text{proc}][\text{request}]].\text{ctype} = \text{"rsend"} \Rightarrow \text{58.10} - \text{58.11}) \\
\quad \exists j \in (0 \dots (N - 1)) : \\
\quad \quad \exists k \in (1 \dots \text{Len}(\text{requests}[j])) :
\end{array}$$

```

    ∧ requests[j][k].active
    ∧ ¬requests[j][k].transmitted
    ∧ ¬requests[j][k].cancelled
    ∧ Match(requests[proc][Memory[proc][request]].message, requests[j][k].message),
    “Error: MPI_Start tried to start a rsend request when no matching message exists.”)
  ∧ Assert(requests[proc][Memory[proc][request]].ctype = “bsend” ⇒
    message_buffer[proc] < bufsize[proc],
    “Error: MPI_Start tried to start a bsend request when insufficient buffering was available.”)
  ∧ Assert(requests[proc][Memory[proc][request]].persist, 57.44 – 57.45, 58.8
    “Error: MPI_Start tried to start a non-persistent request.”)
  ∧ requests' = [requests EXCEPT ![proc] =
    [ @ EXCEPT ![Memory[proc][request]] =
      [ @ EXCEPT
        !.active = TRUE, 58.9
        !.started = TRUE,
        !.transmitted = FALSE,
        !.cancelled = FALSE]]]
  ∧ IF requests[proc][Memory[proc][request]].ctype = “bsend”
    THEN
      message_buffer' = [message_buffer EXCEPT ![proc] = @ + 1]
    ELSE
      UNCHANGED ⟨message_buffer⟩
  ∧ UNCHANGED ⟨group, communicator, bufsize, initialized, collective⟩
  ∧ UNCHANGED ⟨Memory⟩

```

Start a list of persistent communications.

Can you start many rsends with only one matching receive posted? –maybe yes

Can you start many bsends with only enough buffering for a subset of the sends? –maybe no

$MPI_Startall(count, array_of_requests, return, proc) \triangleq$

LET $m \triangleq \{x \in (0 .. (count - 1)) : requests[proc][Memory[proc][array_of_requests + x]].ctype = “bsend”\}$

∧ Assert(initialized[proc] = “initialized”, 200.10 – 200.12

“Error: MPI_Startall called with proc not in initialized state.”)

∧ Assert($\forall i \in (0 .. (count - 1)) : \neg requests[proc][array_of_requests[i]].active$,

“Error: MPI_Startall called with some request already active.”)

∧ Assert($\forall i \in (0 .. (count - 1)) : array_of_requests[i] \neq MPI_REQUEST_NULL$,

“Error: MPI_Startall called with some request null.”)

∧ Assert($\forall i \in (0 .. (count - 1)) :$

$requests[proc][array_of_requests[i]].ctype = “rsend” \Rightarrow$ 58.10 – 58.11

$\exists j \in (0 .. (N - 1)) :$

$\exists k \in (1 .. Len(requests[j])) :$

∧ requests[j][k].active

∧ ¬requests[j][k].transmitted

∧ ¬requests[j][k].cancelled

∧ Match(requests[proc][array_of_requests[i]].message, requests[j][k].message),

“Error: MPI_Start tried to start a rsend request when no matching message exists.”)

$$\wedge \text{Assert}(\forall i \in (0 \dots (\text{count} - 1)) : \\ \text{requests}[\text{proc}][\text{array_of_requests}[i]].\text{ctype} = \text{“bsend”} \Rightarrow \\ \text{message_buffer}[\text{proc}] + \text{Cardinality}(m) < \text{bufsize}[\text{proc}], \\ \text{“Error: MPI_Start tried to start a bsend request when insufficient buffering was available.”})$$

$$\wedge \text{Assert}(\forall i \in (0 \dots (\text{count} - 1)) : \\ \text{requests}[\text{proc}][\text{array_of_requests}[i]].\text{persist}, \text{ 57.44 – 57.45, 58.8} \\ \text{“Error: MPI_Start tried to start a non-persistent request.”})$$

$$\wedge \text{requests}' = [\text{requests} \text{ EXCEPT } ![\text{proc}] = \\ [i \in (1 \dots \text{Len}(\text{requests}[\text{proc}])) \mapsto \\ \text{IF } \exists j \in (0 \dots (\text{count} - 1)) : \text{array_of_requests}[j] = i \\ \text{THEN } \text{ 58.9} \\ [\text{requests}[\text{proc}][i] \text{ EXCEPT} \\ \text{!.active} = \text{TRUE}, \\ \text{!.started} = \text{TRUE}, \\ \text{!.transmitted} = \text{FALSE}, \\ \text{!.cancelled} = \text{FALSE}] \\ \text{ELSE} \\ \text{requests}[\text{proc}][i]] \\ \wedge \text{message_buffer}' = [\text{message_buffer} \text{ EXCEPT } ![\text{proc}] = @ + \text{Cardinality}(m)] \\ \wedge \text{UNCHANGED } \langle \text{group}, \text{communicator}, \text{bufsize}, \text{initialized}, \text{collective} \rangle \\ \wedge \text{UNCHANGED } \langle \text{Memory} \rangle$$

Section 3.10 Send-recv

Can this be done with only one transition? I don't think so.

$$\text{MPI_Sendrecv}(\text{sendbuf}, \text{sendcount}, \text{sendtype}, \text{dest}, \\ \text{sendtag}, \text{recvbuf}, \text{recvcount}, \text{recvtype}, \\ \text{source}, \text{recvtag}, \text{comm}, \text{status}) \triangleq$$

Section 4.3 Barrier

$$\text{MPI_Barrier_init}(\text{comm}, \text{return}, \text{proc}) \triangleq \\ \wedge \vee \wedge \text{collective}[\text{communicator}[\text{proc}][\text{comm}].\text{collective}].\text{state} = \text{“vacant”} \\ \wedge \text{collective}' = [\text{collective} \text{ EXCEPT } ![\text{communicator}[\text{proc}][\text{comm}].\text{collective}] = \\ [@ \text{ EXCEPT} \\ \text{!.participants} = @ \cup \{\text{proc}\}, \\ \text{!.type} = \text{“barrier”}, \\ \text{!.state} = \text{“in”}] \\ \vee \wedge \text{collective}[\text{communicator}[\text{proc}][\text{comm}].\text{collective}].\text{state} = \text{“in”} \\ \wedge \text{proc} \notin \text{collective}[\text{communicator}[\text{proc}][\text{comm}].\text{collective}].\text{participants} \\ \wedge \text{collective}' = [\text{collective} \text{ EXCEPT } ![\text{communicator}[\text{proc}][\text{comm}].\text{collective}] = \\ [@ \text{ EXCEPT } \text{!.participants} = @ \cup \{\text{proc}\}]] \\ \wedge \text{UNCHANGED } \langle \text{group}, \text{communicator}, \text{bufsize}, \text{message_buffer}, \text{requests}, \text{initialized} \rangle$$

\wedge UNCHANGED \langle Memory \rangle

$MPI_Barrier_wait(comm, return, proc) \triangleq$
 $\wedge \vee \wedge$ $collective[communicator[proc][comm].collective].participants = group[proc][communicator[proc][comm].collective].participants$
 \wedge $proc \in collective[communicator[proc][comm].collective].participants$
 \wedge $collective[communicator[proc][comm].collective].state = "in"$
 \wedge $collective' = [collective \text{ EXCEPT } ![communicator[proc][comm].collective] =$
 $[@ \text{ EXCEPT}$
 $!.participants = @ \setminus \{proc\},$
 $!.state = "out"]]$
 $\vee \wedge$ $proc \in collective[communicator[proc][comm].collective].participants$
 \wedge $collective[communicator[proc][comm].collective].state = "out"$
 \wedge IF $collective[communicator[proc][comm].collective].participants = \{proc\}$
 THEN
 $collective' = [collective \text{ EXCEPT } ![communicator[proc][comm].collective] =$
 $[@ \text{ EXCEPT}$
 $!.participants = \{\},$
 $!.state = "vacant"]]$
 ELSE
 $collective' = [collective \text{ EXCEPT } ![communicator[proc][comm].collective] =$
 $[@ \text{ EXCEPT } !.participants = @ \setminus \{proc\}]]$
 \wedge UNCHANGED \langle group, communicator, bufsize, message_buffer, requests, initialized \rangle
 \wedge UNCHANGED \langle Memory \rangle

Section 5.3.1 Group Accessors

No text description.

$MPI_Group_size(gr, size, return, proc) \triangleq$
 \wedge Assert($initialized[proc] = "initialized"$, 200.10 – 200.12
 "Error: MPI_Group_size called with proc not in initialized state.")
 \wedge $Memory' = [Memory \text{ EXCEPT } ![proc] = [@ \text{ EXCEPT } ![size] = group[proc][gr].size]$
 \wedge UNCHANGED \langle group, communicator, bufsize, message_buffer, requests, initialized, collective \rangle

No text description.

$MPI_Group_rank(gr, rank, return, proc) \triangleq$
 \wedge Assert($initialized[proc] = "initialized"$, 200.10 – 200.12
 "Error: MPI_Group_rank called with proc not in initialized state.")
 \wedge $Memory' = [Memory \text{ EXCEPT } ![proc] =$
 $[@ \text{ EXCEPT } ![rank] =$
 IF $proc \in group[proc][gr].members$ THEN $group.ranking[proc]$ ELSE $MPI_UNDEFINED]$
 \wedge UNCHANGED \langle group, communicator, bufsize, message_buffer, requests, initialized, collective \rangle

$MPI_Group_translate_ranks(group1, n, ranks1, group2, ranks2, return, proc) \triangleq$

$\wedge \text{Assert}(\text{initialized}[\text{proc}] = \text{"initialized"}, \text{200.10} - \text{200.12})$
 "Error: MPI_Group_translate_ranks called before MPI_Init."
 $\wedge \text{Assert}(\text{group1} \in \text{MPI_COMM_WORLD} \dots (\text{MPI_COMM_WORLD} + \text{MAX_GROUP}),$
 "Error: MPI_Group_translate_ranks called with invalid handle for group1."
 $\wedge \text{Assert}(\text{group2} \in \text{MPI_COMM_WORLD} \dots (\text{MPI_COMM_WORLD} + \text{MAX_GROUP}),$
 "Error: MPI_Group_translate_ranks called with invalid handle for group2."
 $\wedge \text{Assert}(n = \text{Cardinality}(\text{DOMAIN ranks1}), \text{138.3})$
 "Error: MPI_Group_translate_ranks called with invalid n."
 $\wedge \text{Memory}' = [\text{Memory} \text{ EXCEPT } ![\text{proc}] =$
 $[i \in 1 \dots \text{Len}(\text{Memory}[\text{proc}]) \mapsto$
 IF $i \in \text{ranks2} \dots (\text{ranks2} + n)$
 THEN $\text{group}[\text{proc}][\text{group2}].\text{ranking}[\text{group}[\text{proc}][\text{group1}].\text{invranking}[\text{ranks1}[i]]]$
 ELSE $\text{Memory}[\text{proc}][i]$ not quite right as there is no possibility of MPI_UNDEFINED being assigned.
 $\wedge \text{UNCHANGED mpi_vars}$

$\text{MPI_Group_compare}(\text{group1}, \text{group2}, \text{result}, \text{return}, \text{proc}) \triangleq$
 $\wedge \text{Assert}(\text{initialized}[\text{proc}] = \text{"initialized"}, \text{200.10} - \text{200.12})$
 "Error: MPI_Group_compare called before MPI_Init."
 $\wedge \text{result}' = [\text{result} \text{ EXCEPT } ![\text{proc}] =$
 IF $\vee \text{group1} = \text{group2}$ 138.31
 $\vee \wedge \text{group}[\text{proc}][\text{group1}].\text{members} = \text{group}[\text{proc}][\text{group2}].\text{members}$
 $\wedge \text{group}[\text{proc}][\text{group1}].\text{ranking} = \text{group}[\text{proc}][\text{group2}].\text{ranking}$
 THEN
 MPI_IDENT
 ELSE 138.32
 IF $\wedge \text{group}[\text{proc}][\text{group1}].\text{members} = \text{group}[\text{proc}][\text{group2}].\text{members}$
 $\wedge \text{group}[\text{proc}][\text{group1}].\text{ranking} \neq \text{group}[\text{proc}][\text{group2}].\text{ranking}$
 THEN
 MPI_SIMILAR
 ELSE
 MPI_UNEQUAL] 138.33
 $\wedge \text{UNCHANGED mpi_vars}$

Section 5.3.2 Group Constructors

$\text{MPI_Comm_group}(\text{comm}, \text{gr}, \text{return}, \text{proc}) \triangleq$
 $\wedge \text{Assert}(\text{initialized}[\text{proc}] = \text{"initialized"}, \text{200.10} - \text{200.12})$
 "Error: MPI_Comm_group called before MPI_Init."
 $\wedge \text{Memory}' = [\text{Memory} \text{ EXCEPT } ![\text{proc}] = [@ \text{ EXCEPT } ![\text{gr}] = \text{communicator}[\text{proc}][\text{comm}].\text{group}]$ 139.19
 $\wedge \text{UNCHANGED mpi_vars}$

$\text{MPI_Group_union}(\text{group1}, \text{group2}, \text{newgroup}, \text{return}, \text{proc}) \triangleq$
 $\wedge \text{Assert}(\text{initialized}[\text{proc}] = \text{"initialized"}, \text{200.10} - \text{200.12})$
 "Error: MPI_Group_union called before MPI_Init."
 $\wedge \exists i \in 0 \dots (\text{MAX_GROUP} - 1) :$
 LET $\text{newmembers} \triangleq \text{group}[\text{proc}][\text{group1}].\text{members} \cup$

```

      group[proc][group2].members
IN
  ∧ group[proc][i] = MPI_GROUP_EMPTY
  ∧ newgroup' = [newgroup EXCEPT ![proc] = i]
  ∧ group' =
    [group EXCEPT ![proc] =
      [⊙ EXCEPT ![i] =
        [members ↦ newmembers,
          size ↦ Cardinality(newmembers),
          ranking ↦
            [j ∈ 0 .. (Cardinality(newmembers) - 1) ↦
              IF j < group[proc][group1].size
                THEN group[proc][group1].ranking[j]
                ELSE group[proc][group1].ranking[j]]]] \ * incorrect, need to fix
    ]
  ∧ UNCHANGED ⟨communicator, bufsize, message_buffer, requests, initialized, collective⟩

```

Section 5.4.1 Communicator Accessors

```

MPI_Comm_size(comm, size, return, proc) ≜
  ∧ Assert(initialized[proc] = "initialized", 200.10 - 200.12
    "MPI_Comm_size called with proc not in initialized state.")
  ∧ Memory' = [Memory EXCEPT ![proc] = [⊙ EXCEPT ![size] = group[proc][communicator[proc][comm].group]
  ∧ UNCHANGED mpi_vars

```

```

MPI_Comm_rank(comm, rank, return, proc) ≜
  ∧ Assert(initialized[proc] = "initialized", 200.10 - 200.12
    "MPI_Comm_rank called with proc not in initialized state.")
  ∧ Memory' = [Memory EXCEPT ![proc] = [⊙ EXCEPT ![rank] = group[proc][communicator[proc][comm].group]
  ∧ UNCHANGED ⟨group, communicator, bufsize, message_buffer, requests, initialized, collective⟩

```

```

MPI_Comm_compare(comm1, comm2, result, return, proc) ≜
  ∧ Assert(initialized[proc] = "initialized", 200.10 - 200.12
    "MPI_Comm_rank called with proc not in initialized state.")
  ∧ IF comm1 = comm2
    THEN
      result' = MPI_IDENT
    ELSE
      IF ∧ communicator[proc][comm1].group = communicator[proc][comm2].group
        ∧ communicator[proc][comm1].group.ranking = communicator[proc][comm2].group.ranking
        THEN
          result' = MPI_CONGRUENT
        ELSE
          IF communicator[proc][comm1].group = communicator[proc][comm2].group
            THEN
              result' = MPI_SIMILAR
            ELSE
              result' = MPI_UNEQUAL

```

\wedge UNCHANGED $\langle group, communicator, bufsize, message_buffer, requests, initialized, collective \rangle$

Section 7.5 Startup

199.12 – 199.17

Initialize the participation of this process within a distributed computation.

$MPI_Init(argc, argv, return, proc) \triangleq$
 \wedge Assert($initialized[proc] = \text{"uninitialized"}$, 199.12
 "MPI_Init called with proc not in uninitialized state.")
 \wedge $initialized' = [initialized \text{ EXCEPT } ![proc] = \text{"initialized"}]$ 199.13
 \wedge UNCHANGED $\langle Memory \rangle$
 \wedge UNCHANGED $\langle group, communicator, bufsize, message_buffer, requests, collective \rangle$

Finalize the participation of this process within a distributed computation.

Do buffered operations complete when the message is transmitted or buffered?

$MPI_Finalize(return, proc) \triangleq$
 \wedge Assert($initialized[proc] = \text{"initialized"}$, 200.10 – 200.12
 "Error: MPI_Finalize called with proc not in initialized state.")
 \wedge Assert($\forall i \in (1 .. Len(requests[proc])) :$ 199.47
 $\neg requests[proc][i].active,$
 "Error: MPI_Finalize called when some message was still active.")
 \wedge Assert($bufsize[proc] = 0,$
 "Error: MPI_Finalize called before the buffer is detached.")
 \wedge $initialized' = [initialized \text{ EXCEPT } ![proc] = \text{"finalized"}]$ 199.46
 \wedge UNCHANGED $\langle group, communicator, bufsize, message_buffer, requests, collective \rangle$
 \wedge UNCHANGED $\langle Memory \rangle$

Determine whether MPI_Init has been called.

$MPI_Initialized(flag, return, proc) \triangleq$
 \wedge $Memory' = [Memory \text{ EXCEPT } ![proc] = [@ \text{ EXCEPT } ![flag] =$
 IF $initialized[proc] = \text{"initialized"}$ 200.2
 THEN TRUE
 ELSE FALSE]]
 \wedge UNCHANGED $\langle group, communicator, bufsize, message_buffer, requests, initialized, collective \rangle$

"Best effort to clean up"

$MPI_Abort(comm, errorcode, return, proc) \triangleq$
 $\forall p \in (0 .. (N - 1)) :$
 $\forall m \in (1 .. Len(requests[p])) :$
 $\wedge requests[p][m].active$
 $\wedge \neg requests[p][m].transmitted$
 $\Rightarrow requests[p][m]' = [requests[p][m] \text{ EXCEPT } !.cancelled = \text{TRUE}]$
 \wedge UNCHANGED $\langle group, communicator, bufsize, message_buffer, requests, initialized, collective \rangle$