

# Collective Contracts for Message-Passing Parallel Programs

Anonymous

No Institute Given

**Abstract.** Procedure contracts are a well-known approach for specifying programs in a modular way. We investigate a new contract theory for collective procedures in parallel message-passing programs. As in the sequential setting, one can verify that a procedure  $f$  conforms to its contract using only the contracts, and not the implementations, of the collective procedures called by  $f$ . We apply this approach to C programs that use the Message Passing Interface (MPI), introducing a new contract language that extends the ANSI/ISO C Specification Language. We present contracts for the standard MPI collective functions, as well as many user-defined collective functions. A prototype verification system has been implemented using the CIVL model checker for checking contract satisfaction within small bounds on the number of processes.

**Keywords:** contract, message-passing, MPI, verification, collective

## 1 Introduction

*Procedure contracts* [27, 44, 45] are a well-known way to decompose program verification. In this approach, each procedure  $f$  is specified independently with pre- and postconditions or other invariants. To verify  $f$ , one needs only the contracts, not the implementations, of the procedures called by  $f$ .

Contract languages have been developed for many programming languages. These include the *Java Modeling Language* (JML) [37] for Java and the *ANSI C Specification Language* (ACSL) [10] for C. A number of tools have been developed which (partially) automate the process of verifying that a procedure satisfies its contract; an example for C is Frama-C [18] with the WP plugin [9].

In this paper, we explore a procedure contract system for message-passing parallel programs, specifically for programs that use the Message-Passing Interface (MPI) [43], the de facto standard for high performance computing.

Our contracts apply to *collective-style* procedures in these programs. These are procedures  $f$  called by all processes and that are *communication-closed*: any message issued by a send statement in  $f$  is received by a receive statement in  $f$ , and vice-versa. The processes executing  $f$  coordinate in order to accomplish a coherent change in the global state. Examples include all of the standard blocking MPI collective functions [43, Chapter 5], but also many user-defined procedures, such as a procedure to exchange ghost cells in a stencil computation. These procedures are typically specified informally by describing the effect they

produce when called by all processes, rather than the effect of an individual process. They should be formally specified and verified in the same way.

Developers often construct applications by composing collective procedures. As examples, consider the Monte Carlo particle transport code OpenMC [51] (over 24K lines of C++/MPI code) and module `parcsr_ls` in the algebraic multigrid solver AMG [60] (over 35K lines of C/MPI code). Through manual inspection, we confirmed that every function in these codes that involves MPI communication is collective-style.

We begin in Section 2 with a toy message-passing language, so the syntax, semantics, and theoretical results can be stated and proved precisely. The main result is a theorem that justifies a method for verifying a collective procedure using only the contracts of the collective procedures called, as in the sequential case.

Section 3 describes changes needed to apply this system to C/MPI programs. We handle a significant subset of MPI that does not include `MPI_ANY_SOURCE` (“wildcard”) receives. This means program behavior is largely independent of interleaving [53]. There are enough issues to deal with, such as MPI datatypes, input nondeterminism, and nontermination, that we feel it best to leave wildcards for a sequel. A prototype verification system for such programs, using the CIVL model checker, is described and evaluated in Section 4. Related work is discussed in Section 5. In Section 6, we wrap up with a discussion of the advantages and limitations of our system, and work that remains.

In summary, this paper makes the following contributions: (1) a contract theory for collective message-passing procedures, with mathematically precise syntax and semantics, (2) a theorem justifying a method for verifying that a collective procedure conforms to its contract, (3) a contract language for a large subset of MPI, based on the theory but also dealing with additional intricacies of MPI, and (4) a prototype verification tool for checking that collective MPI procedures conform to their contracts.

## 2 A Theory of Collective Contracts

### 2.1 Language

We describe a simple message-passing language MINIMP with syntax in Figure 1. There is one datatype: integers; 0 is interpreted as *false* and any non-zero integer as *true*. A program consists of global variable declarations followed by (mutually recursive) procedure definitions. Global variables may start with arbitrary values. Each procedure takes a sequence of formal parameters. The procedure body consists of local variable declarations followed by a sequence of statements. Local variables are initially 0. Assignment, branch, loop, call, and compound statements have the usual semantics. Operations have the usual meaning and always return some value—even if the second argument of division is 0, e.g. Operators with ‘\’, described below, occur only in the optional contract.

A procedure is executed by specifying a positive integer  $n$ , the number of processes. Each process executes its own “copy” of the code; there is no shared

```

program ::= ( int x ; ) * procdef +
procdef ::= contract? void f ( ( int x ( , int x ) * ) ? ) { ( int x ; ) * s * }
s ∈ stmt ::= x = e ; | f ( ( e ( , e ) * ) ? ) ; | if ( e ) s ( else s ) ? | while ( e ) s
           | { s * } | send e to e ; | recv x from e ;
e ∈ expr ::= c | x | nprocs | pid | ⊖ e | e ⊙ e | \on(e, e) | \old(e)
contract ::= /*@ requires e; ensures e; assigns (x, (x) * ) ? ;
           waitsfor { e | int x ; e } ; */
c ∈ ℤ      x, f ∈ ID      ⊖ ∈ { -, ! }      ⊙ ∈ { +, -, *, /, %, ==, <, <=, &&, || }

```

Fig. 1. MINIMP syntax

memory. Each process has a unique ID number in  $PID = \{0, \dots, n-1\}$ . A process can obtain its ID using the primitive `pid`; it can obtain  $n$  using `nprocs`.

The command “`send data to dest`” sends the value of `data` to the process with ID `dest`. There is one FIFO message buffer for each ordered pair of processes  $p \rightarrow q$  and the effect of `send` is to enqueue the message on the buffer for which  $p$  is the ID of the sender and  $q$  is `dest`. The buffers are unbounded, so `send` never blocks. Command “`recv buf from source`” removes the oldest buffered message originating from `source` and stores it in variable `buf`; this command blocks until a message becomes available. A `dest` or `source` not in  $PID$  results in a no-op.

A procedure  $f$  with a contract is a *collective procedure*. The contract encodes a claim about executions of  $f$ : if  $f$  is called collectively (by all processes), in such a way that the precondition (specified in the `requires` clause) holds, then all of the following hold for each process  $p$ :  $p$  will eventually return;  $p$ ’s postcondition (specified in the `ensures` clause) will hold at the post-state; all variables not listed in  $p$ ’s `assigns` clause will have their pre-state values at the post-state; and if  $q$  is in  $p$ ’s `waitsfor` set then  $p$  will not return before  $q$  enters the call. These notions will be made precise below.

Global variables and the formal parameters of the procedure are the only variables that may occur free in a contract; only globals may occur in the `assigns` clause. A postcondition may use `\old(e)` to refer to the value of expression  $e$  in the pre-state; `\old` may not occur in this  $e$ . Pre- and postconditions can use `\on(e, i)` to refer to the value of  $e$  on process  $i$ . These constructs allow contracts to relate the state of different processes, and the state before and after the call.

*Example 1.* The program of Figure 2 has two procedures, both collective. Procedure `g` accepts an argument `k` and sends its value for global variable `x` to its right neighbor, in a cyclic ordering. It then receives into local variable `y` from its left neighbor  $q$ , adds `k` to the received value, and stores the result in `x`. The contract for `g` states that when  $p$  exits (returns), the value of `x` on  $p$  is the sum of  $k$  and the original value of `x` on  $q$ . It also declares  $p$  cannot exit until  $q$  has entered. Procedure `f` calls `g nprocs` times. Its contract requires that all processes call `f` with the same value for `k`. It ensures that upon return, the value of `x` is the sum of its original value and the product of `nprocs` and `k`. It also declares that no process can exit until every process has entered.

```

int x;
/*@ requires 1; ensures x == \on(\old(x), (pid+nprocs-1)%nprocs) + k;
    assigns x; waitsfor { j | int j; j == (pid+nprocs-1)%nprocs }; */
void g(int k) {
    int y;
    send x to (pid+1)%nprocs;
    recv y from (pid+nprocs-1)%nprocs;
    x = y+k;
}
/*@ requires k == \on(k,0); ensures x == \old(x) + nprocs*k;
    assigns x; waitsfor { j | int j; 0<=j && j<nprocs }; */
void f(int k) { int i; i = 0; while (i<nprocs) { g(k); i = i+1; } }

```

Fig. 2. cyc: a MINIMP program

## 2.2 Semantics

Semantics for procedural programs are well-known (e.g., [2]), so we will only summarize the standard aspects of the MINIMP semantics. Fix a program  $P$  and an integer  $n \geq 1$  for the remainder of this section. Each procedure in  $P$  may be represented as a *program graph*, which is a directed graph in which nodes correspond to locations in the procedure body. Each program graph has a designated start node. An edge is labeled by either an expression  $\phi$  (a *guard*) or one of the following kinds of statements: *assignment*, *call*, *return*, *send* or *receive*. An edge labeled *return* is added to the end of each program graph, and leads to the terminal node, which has no outgoing edges.

A *process state* comprises an assignment of values to global variables and a call stack. Each entry in the stack specifies a procedure  $f$ , the values of the local variables (including formal parameters) for  $f$ , and the program counter, which is a node in  $f$ 's program graph. A *state* specifies a process state for each process, as well as the state of channel  $p \rightarrow q$  for all  $p, q \in \text{PID}$ . The channel state is a finite sequence of integers, the buffered messages sent from  $p$  to  $q$ .

An *action* is a pair  $t = \langle e, p \rangle$ , where  $e$  is an edge  $u \xrightarrow{\alpha} v$  in a program graph and  $p \in \text{PID}$ . Action  $t$  is *enabled* at state  $s$  if the program counter of the top entry of  $p$ 's call stack in  $s$  is  $u$  and one of the following holds:  $\alpha$  is a guard  $\phi$  and  $\phi$  evaluates to *true* in  $s$ ;  $\alpha$  is an assignment, call, return, or send; or  $\alpha$  is a receive with source  $q$  and channel  $q \rightarrow p$  is nonempty in  $s$ . The execution of an enabled action from  $s$  results in a new state  $s'$  in the natural way. In particular, execution of a call pushes a new entry onto the stack of the calling process; execution of a return pops the stack and, if the resulting stack is not empty, moves the caller to the location just after the call. The triple  $s \xrightarrow{t} s'$  is a *transition*.

Let  $f$  be a procedure and  $s_0$  a state with empty channels, and in which each process has one entry on its stack, the program counter of which is the start location for  $f$ . An  $n$ -process *execution*  $\zeta$  of  $f$  is a finite or infinite chain of transitions  $s_0 \xrightarrow{t_1} s_1 \xrightarrow{t_2} \dots$ . The *length* of  $\zeta$ , denoted  $\text{len}(\zeta)$ , is the number of transitions in  $\zeta$ . An execution must be *fair*: if a process  $p$  becomes enabled at

some point in an infinite execution, then eventually  $p$  will execute. Note that, once  $p$  becomes enabled, it will remain enabled until it executes, as no process other than  $p$  can remove a buffered message with destination  $p$ .

A process  $p$  *terminates* in  $\zeta$  if for some  $i$ , the stack for  $p$  is empty in  $s_i$ . We say  $\zeta$  *terminates* if  $p$  terminates in  $\zeta$  for all  $p \in \text{PID}$ . The execution *deadlocks* if it is finite, does not terminate, and ends in a state with no enabled action.

It is often convenient to add a “driver” to  $P$  when reasoning about executions of a collective procedure  $f$ . Say  $f$  takes  $m$  formal parameters. Form a program  $P^f$  by adding fresh global variables  $x_1, \dots, x_m$  to  $P$ , and adding a procedure

$$\text{void main() } \{ f(x_1, \dots, x_m); \}.$$

By “execution of  $P^f$ ,” we mean an execution of **main** in this new program.

### 2.3 Collective Correctness

The goal of this section is to define what it means for a procedure to conform to its contract. This notion comprises several conditions on the invocation of collective procedures and their interaction with communication operations. Some of the conditions are generic and others are specified by the contract clauses.

Fix a program  $P$  and integer  $n \geq 1$ . Let  $\mathcal{C}$  be the set of names of collective procedures of  $P$ . Let  $\zeta$  be an execution  $s_0 \xrightarrow{t_1} s_1 \xrightarrow{t_2} \dots$  of a procedure in  $P$ . For  $i \in 1..\text{len}(\zeta)$ , let  $\zeta^i$  denote the prefix of  $\zeta$  of length  $i$ , i.e., the execution  $s_0 \xrightarrow{t_1} \dots \xrightarrow{t_i} s_i$ .

The first correctness condition for  $\zeta$  is *collective consistency*. To define this concept, consider strings over the alphabet consisting of symbols of the form  $\mathbf{e}^f$  and  $\mathbf{x}^f$ , for  $f \in \mathcal{C}$ . Given an action  $t$  and  $p \in \text{PID}$ , define string  $T_p(t)$  as follows:

- if  $t$  is a call by  $p$  to some  $f \in \mathcal{C}$ ,  $T_p(t) = \mathbf{e}^f$  ( $t$  is called an *enter* action)
- if  $t$  is a return by  $p$  from some  $f \in \mathcal{C}$ ,  $T_p(t) = \mathbf{x}^f$  ( $t$  is called an *exit* action)
- otherwise,  $T_p(t)$  is the empty string.

Now let  $T_p(\zeta)$  be the concatenation  $T_p(t_1)T_p(t_2)\dots$ . Hence  $T_p(\zeta)$  records the sequence of collective actions—enter or exit actions—taken by  $p$ .

**Definition 1.** An execution  $\zeta$  is *collective consistent* if there is some  $p \in \text{PID}$  such that for all  $q \in \text{PID}$ ,  $T_q(\zeta)$  equals or is a prefix of  $T_p(\zeta)$ . We say  $\zeta$  commits a consistency violation at step  $i$  if  $\zeta^{i-1}$  is collective consistent but  $\zeta^i$  is not.

For the rest of this section, assume  $\zeta$  is collective consistent.

The sequence of actions performed by  $p$  in  $\zeta$  is divided into segments whose boundaries are the collective actions of  $p$ . More precisely, given  $i \in 0..\text{len}(\zeta)$  and  $p \in \text{PID}$ , define  $k = \text{seg}_p(\zeta, i)$  to be the number of collective actions of  $p$  in  $t_1, \dots, t_i$ . We say  $p$  is in *segment  $k$  at state  $i$* .

We now turn to the issue of evaluation of pre- and postconditions. Let  $f$  be a collective procedure in  $P$  with precondition  $\text{pre}(f)$  and postcondition  $\text{post}(f)$ . Let  $V_f$  be the union of the set of formal parameters of  $f$  and the global variables

of  $P$ . As noted above, these are the only variables that may occur free in  $\text{pre}(f)$  and  $\text{post}(f)$ . An  $f$ -valuation is a function  $\alpha: \text{PID} \rightarrow (V_f \rightarrow \mathbb{Z})$ . For each process,  $\alpha$  specifies a value for each free variable that may occur in  $\text{pre}(f)$  or  $\text{post}(f)$ .

For any expression  $e$  that may occur as a sub-expression of  $\text{pre}(f)$ , and  $p \in \text{PID}$ , define  $\llbracket e \rrbracket_{\alpha,p} \in \mathbb{Z}$  as follows:

$$\begin{aligned} \llbracket c \rrbracket_{\alpha,p} &= c & \llbracket \ominus e \rrbracket_{\alpha,p} &= \ominus \llbracket e \rrbracket_{\alpha,p} \\ \llbracket x \rrbracket_{\alpha,p} &= \alpha(p)(x) & \llbracket e_1 \odot e_2 \rrbracket_{\alpha,p} &= \llbracket e_1 \rrbracket_{\alpha,p} \odot \llbracket e_2 \rrbracket_{\alpha,p} \\ \llbracket \text{nprocs} \rrbracket_{\alpha,p} &= n & \llbracket \backslash \text{on}(e_1, e_2) \rrbracket_{\alpha,p} &= \llbracket e_1 \rrbracket_{\alpha,q}, \text{ where } q = \llbracket e_2 \rrbracket_{\alpha,p}. \\ \llbracket \text{pid} \rrbracket_{\alpha,p} &= p \end{aligned}$$

This is the result of evaluating  $e$  in process  $p$ . Note how  $\backslash \text{on}$  shifts the evaluation context from process  $p$  to the process specified by  $e_2$ , allowing the precondition to refer to the value of an expression on another process.

Evaluation of an expression involving  $\backslash \text{old}$ , which may occur only in  $\text{post}(f)$ , requires a second  $f$ -valuation  $\beta$  specifying values in the pre-state. The definition of  $\llbracket \cdot \rrbracket_{\alpha,\beta,p}$  repeats the rules above, replacing each subscript “ $\alpha$ ” with “ $\alpha, \beta$ ”, and adds one rule:

$$\llbracket \backslash \text{old}(e) \rrbracket_{\alpha,\beta,p} = \llbracket e \rrbracket_{\beta,p}.$$

Say  $1 \leq i \leq \text{len}(\zeta)$  and  $t_i$  is an  $\mathbf{e}^f$  action in process  $p$ . Let  $r = \text{seg}_p(\zeta, i)$  and

$$Q = \{q \in \text{PID} \mid \text{seg}_q(\zeta, i) \geq r\}, \quad \alpha': Q \rightarrow (V_f \rightarrow \mathbb{Z}),$$

where  $\alpha'(q)(v)$  is the value of  $v$  on process  $q$  in state  $s_{j(q)}$ , and  $j(q)$  is the unique integer in  $1..i$  such that  $t_{j(q)}$  is the  $r$ -th collective action of  $q$  in  $\zeta$ . (As  $\zeta$  is collective consistent,  $t_{j(q)}$  is also an  $\mathbf{e}^f$  action.) In other words,  $\alpha'$  uses the values of process  $q$ 's variables just after  $q$  entered the call. Now,  $\alpha'$  is not an  $f$ -valuation unless  $Q = \text{PID}$ . Nevertheless, we can ask whether  $\alpha'$  can be extended to an  $f$ -valuation  $\alpha$  such that  $\llbracket \text{pre}(f) \rrbracket_{\alpha,q}$  holds for all  $q \in \text{PID}$ . If no such  $\alpha$  exists, we say a *precondition violation* occurs at step  $i$ .

*Example 2.* Consider program `cyc` of Figure 2. Suppose process 1 calls  $f(1)$  and process 2 calls  $f(2)$ . Then a precondition violation of  $f$  occurs with the second call, because there is no value that can be assigned to  $k$  on process 0 for which  $1 = \backslash \text{on}(k, 0)$  and  $2 = \backslash \text{on}(k, 0)$  both hold.

If  $t_i$  is an  $\mathbf{x}^f$  action, define  $Q$  and  $j(q)$  as above; for any  $q \in Q$ ,  $t_{j(q)}$  is also an  $\mathbf{x}^f$  action. Let  $\alpha'(q)(v)$  be the value of  $v$  in  $q$  at state  $s_{j(q)-1}$ , i.e., just before  $q$  exits. Define  $k(q) \in 1..j(q)-1$  so that  $t_{k(q)}$  is the  $\mathbf{e}^f$  action in  $q$  corresponding to  $t_{j(q)}$ , i.e.,  $t_{k(q)}$  is the call that led to the return  $t_{j(q)}$ . Define  $\beta': Q \rightarrow (V_f \rightarrow \mathbb{Z})$  so that  $\beta'(q)(v)$  is the value of  $v$  on  $q$  in state  $s_{k(q)}$ , i.e., in the pre-state. A *postcondition violation* occurs if it is not the case that there are extensions of  $\alpha'$  and  $\beta'$  to  $f$ -valuations  $\alpha$  and  $\beta$  such that  $\llbracket \text{post}(f) \rrbracket_{\alpha,\beta,q}$  holds for all  $q \in \text{PID}$ .

We now explain the *waitsfor* contract clause. Assume again that  $t_i$  is an  $\mathbf{x}^f$  action in process  $p$ , and that  $k$  is the index of the corresponding  $\mathbf{e}^f$  action in  $p$ . The expression in the *waitsfor* clause is evaluated at the pre-state  $s_k$  to yield a

set  $W \subseteq \text{PID}$ . A *waitsfor violation* occurs at step  $i$  if there is some  $q \in W$  such that  $\text{seg}_q(\zeta, i) < \text{seg}_p(\zeta, k)$ , i.e.,  $p$  exits a collective call before  $q$  has entered it.

We can now encapsulate all the ways something may go wrong with collective procedures and their contracts:

**Definition 2.** Let  $P$  be a program,  $\zeta = s_0 \xrightarrow{t_1} s_1 \cdots$  an execution of a procedure in  $P$ , and  $i \in 1..\text{len}(\zeta)$ . Let  $p$  be the process of  $t_i$  and  $r = \text{seg}_p(\zeta, i)$ . We say  $\zeta$  commits a collective error at step  $i$  if any of the following occur at step  $i$ :

1. a consistency, precondition, postcondition, or waitsfor violation,
2. an assigns violation:  $t_i$  is an exit action and the value of a variable not in  $p$ 's assigns set differs from its pre-state value,
3. a segment boundary violation:  $t_i$  is a receive of a message sent from a process  $q$  at  $t_j$  ( $j < i$ ) and  $\text{seg}_q(\zeta, j) > r$ ; or  $t_i$  is a send to  $q$  and  $\text{seg}_q(\zeta, i) > r$ , or
4. an unreceived message violation:  $t_i$  is a collective action and there is an unreceived message sent to  $p$  from  $q$  at  $t_j$  ( $j < i$ ), and  $\text{seg}_q(\zeta, j) = r - 1$ .

The last two conditions imply that a message that crosses segment boundaries is erroneous. In particular, if an execution terminates without collective errors, every message sent within a segment is received within that same segment.

**Definition 3.** An execution of a procedure is correct if it is finite, does not deadlock, and has no collective errors.

We can now define what it means for a procedure to conform to its contract. Let  $f$  be a collective procedure in  $P$ . By a  $\text{pre}(f)$ -state, we mean a state of  $P^f$  with empty call stacks, empty channels, and an assignment to the global variables satisfying the precondition of  $f$  for all processes.

**Definition 4.** A collective procedure  $f$  conforms (to its contract) if all executions of  $P^f$  from  $\text{pre}(f)$ -states are correct.

Note that any maximal non-deadlocking finite execution terminates. So a conforming procedure will always terminate if invoked from a  $\text{pre}(f)$ -state, i.e., ours is a “total” (not “partial”) notion of correctness in the Hoare logic sense.

## 2.4 Simulation

In the sequential theory, one may verify properties of a procedure  $f$  using only the contracts of the procedures called by  $f$ . We now generalize that approach for collective procedures. We will assume from now on that  $P$  has no “collective recursion.” That is, in the call graph for  $P$ —the graph with nodes the procedures of  $P$  and an edge from  $f$  to  $g$  if the body of  $f$  contains a call to  $g$ —there is no cycle that includes a collective procedure. This simplifies reasoning about termination.

If  $f, g \in \mathcal{C}$ , we say  $f$  *uses*  $g$  if there is a path of positive length in the call graph from  $f$  to  $g$  on which any node other than the first or last is not in  $\mathcal{C}$ .

Given  $f \in \mathcal{C}$ , we construct a program  $\bar{P}^f$  which abstracts away the implementation details of each collective procedure  $g$  used by  $f$ , replacing the body of  $g$  with a stub that simulates  $g$ 's contract. The stub consists of two new statements. The first may be represented with pseudocode

```
havoc(assigns(g)); assume(post(g));
```

This nondeterministic statement assigns arbitrary values to the variables specified in the *assigns* clause of  $g$ 's contract, as long as those values do not commit a postcondition violation for  $g$ . The second statement may be represented

```
wait(\old(waitsfor(g)));
```

and blocks the calling process  $p$  until all processes in  $p$ 's wait set (evaluated in  $p$ 's pre-state) reach this statement. This ensures the stub will obey  $g$ 's *waitsfor* contract clause. Now  $\overline{P}^f$  is a program with the same set of collective procedure names, and same contracts, as  $P^f$ . A *simulation* of  $f$  is an execution of  $\overline{P}^f$ .

**Theorem 1** *Let  $P$  be a program with no collective recursion. Let  $f$  be a collective procedure in  $P$  and assume all collective procedures used by  $f$  conform. If all simulations of  $f$  from a  $\text{pre}(f)$ -state are correct then  $f$  conforms.*

Theorem 1 is the basis for the contract-checking tool described in Section 4.2. The tool consumes a C/MPI program annotated with procedure contracts. The user specifies a single procedure  $f$  and the tool constructs a CIVL-C program that simulates  $f$  by replacing the collective procedures called by  $f$  with stubs derived from their contracts. It then uses symbolic execution and model checking techniques to verify that all simulations of  $f$  behave correctly. By Theorem 1, one can conclude that  $f$  conforms. A detailed proof of Theorem 1 is given in [3]. Here we summarize the main ideas of the proof. We assume henceforth that  $P$  is a collective recursion-free program.

Two actions from different processes commute as long as the second does not receive a message sent by the first. Two executions are *equivalent* if one can be obtained from the other by a finite number of transpositions of commuting adjacent transitions. We first observe that equivalence preserves most violations:

**Lemma 1** *Let  $\zeta$  and  $\eta$  be equivalent executions of a procedure  $f$  in  $P$ . Then*

1.  $\zeta$  commits a consistency, precondition, postcondition, assigns, segment boundary, or unreceived message violation iff  $\eta$  commits such a violation.
2.  $\zeta$  deadlocks iff  $\eta$  deadlocks.
3.  $\zeta$  is finite iff  $\eta$  is finite.

If  $\zeta$  commits a collective error when control is not inside a collective call made by  $f$  (i.e., when  $f$  is the only collective function on the call stack), we say the error is *observable*. If the error is not observable, it is *internal*. We say  $\zeta$  is *observably correct* if it is finite, does not deadlock, and is free of observable collective errors.

We are interested in observable errors because those are the kind that will be visible in a simulation, i.e., when each collective function  $g$  called by  $f$  is replaced with a stub that mimics  $g$ 's contract.

When  $\zeta$  has no observable collective error, it can be shown that a collective call to  $g$  made within  $\zeta$  can be *extracted* to yield an execution of  $g$ . The idea



behind the proof is to transpose adjacent transitions in  $\zeta$  until all of the actions inside the call to  $g$  form a contiguous subsequence of  $\zeta$ . The resulting execution  $\xi$  is equivalent to  $\zeta$ . Using Lemma 1, it can be shown that  $\xi$  is also observably correct and the segment involving the call to  $g$  can be excised to yield an execution of  $g$ . The next step is to show that extraction preserves internal errors:

**Lemma 2** *Assume  $\zeta$  is an observably correct execution of collective procedure  $f$  in  $P$ . Let  $g_1, g_2, \dots$  be the sequence of collective procedures called from  $f$ . If a transition in region  $r$  of  $\zeta$  commits an internal collective error then the execution of  $P^{g_r}$  extracted from region  $r$  of  $\zeta$  is incorrect.*

A corollary of Lemma 2 may be summarized as “conforming + observably correct = correct”. More precisely,

**Lemma 3** *Let  $f$  be a collective procedure of  $P$ . Assume all collective procedures used by  $f$  conform. Let  $\zeta$  be an execution of  $P^f$ . Then  $\zeta$  is correct if and only if  $\zeta$  is observably correct.*

To see this, suppose  $\zeta$  is observably correct but commits an internal collective error. Let  $r$  be the region of the transition committing the first internal collective error of  $\zeta$ . Let  $g$  be the associated collective procedure used by  $f$ , and  $\chi$  the execution of  $P^g$  extracted from region  $r$  of  $\zeta$ . By Lemma 2,  $\chi$  is incorrect, contradicting the assumption that  $g$  conforms.

Next we show that observable errors will be picked up by some simulation. The following is proved using extraction and Lemma 3:

**Lemma 4** *Suppose  $f$  is a collective procedure of  $P$ , all collective procedures used by  $f$  conform, and  $\zeta$  is an execution of  $P^f$ . If  $\zeta$  has an observable collective error or ends in deadlock then there exists an incorrect simulation of  $f$ .*

Since infinite executions are also considered erroneous, we must ensure they are detected by simulation:

**Lemma 5** *Suppose  $f$  is a collective procedure of  $P$ , and all collective procedures used by  $f$  conform. If  $\zeta$  is an infinite execution of  $P^f$  with no observable collective error then there exists an incorrect simulation of  $f$ .*

Finally, we prove Theorem 1. Assume  $f$  is a collective procedure in  $P$  and all collective procedures used by  $f$  conform. Suppose  $f$  does not conform; we must show there is an incorrect simulation of  $f$ . As  $f$  does not conform, there is an incorrect execution  $\zeta$  of  $P^f$  from a  $pre(f)$ -state. By Lemma 3,  $\zeta$  is not observably correct. If  $\zeta$  is finite or commits an observable collective error, Lemma 4 implies an incorrect simulation exists. Otherwise, Lemma 5 implies such a simulation exists. This completes the proof.

### 3 Collective Contracts for C/MPI

In Section 3.1, we summarize the salient aspects of C/MPI needed for a contract system. Section 3.2 describes the overall grammar of MPI contracts and summarizes the syntax and semantics of each new contract primitive.

### 3.1 Background from MPI

In the toy language of Section 2, every collective procedure was invoked by all processes. In MPI, a collective procedure is invoked by all processes in a *communicator*, an abstraction representing an ordered set of processes and an isolated communication universe.<sup>1</sup> Programs may use multiple communicators. The *size* of a communicator is the number of processes. Each process has a unique *rank* in the communicator, an ID number in  $0..size - 1$ .

In Section 2, a receive always selects the oldest message in a channel. In MPI, a point-to-point send operation specifies a *tag*, an integer attached to the “message envelope.” A receive can specify a tag, in which case the oldest message in the channel with that tag is removed, or the receive can use `MPI_ANY_TAG`, in which case the oldest message is. MPI collective functions do not use tags.

MPI communication operations use *communication buffers*. A buffer *b* is specified by a pointer *p*, *datatype* *d* (an object of type `MPI_Datatype`), and nonnegative integer *count*. There are constants of type `MPI_Datatype` corresponding to the C basic types: `MPI_INT`, `MPI_DOUBLE`, etc. MPI provides functions to build aggregate datatypes. Each datatype specifies a *type map*: a sequence of ordered pairs  $(t, m)$  where *t* is a basic type and *m* is an integer displacement in bytes. A type map is *nonoverlapping* if the memory regions specified by distinct entries in the type map do not intersect. A receive operation requires a nonoverlapping type map; no such requirement applies to sends. For example, the type map  $\{(\text{int}, 0), (\text{double}, 8)\}$ , together with *p*, specifies an `int` at *p* and a `double` at  $(\text{char}*)p+8$ . As long as  $\text{sizeof}(\text{int}) \leq 8$ , this type map is nonoverlapping.

The *extent* of *d* is the distance from its lowest to its highest byte, including possible padding bytes at the end needed for alignment; the precise definition is given in the MPI Standard. The type map of *b* is defined to be the concatenation of  $T_0, \dots, T_{count-1}$ , where  $T_i$  is the type map obtained by adding  $i * \text{extent}(d)$  to the displacements of the entries in the type map of *d*. For example, if *count* is 2,  $\text{sizeof}(\text{double}) = 8$  and `ints` and `doubles` are aligned at multiples of 8 bytes, the buffer type map in the example above is

$$\{(\text{int}, 0), (\text{double}, 8), (\text{int}, 16), (\text{double}, 24)\}.$$

A message is created by reading memory specified by the send buffer, yielding a sequence of basic values. The message has a *type signature*—the sequence of basic types obtained by projecting the type map onto the first component. The receive operation consumes a message and writes the values into memory according to the receive buffer’s type map. Behavior is undefined if the send and receive buffers do not have the same type signature.

### 3.2 Contract structure

We now describe the syntax and semantics for C/MPI function contracts. A contract may specify either an MPI collective function, or a user-defined collective

<sup>1</sup> We consider only *intra-communicators* in this paper.

```

function-contract ::= requires-clause* terminates-clause* decreases-clause?
                    simple-clause* comm-clause* named-behavior* completeness-clause*
                    collective-contract*
simple-clause ::= assigns-clause | ensures-clause | allocation-clause | abrupt-clause
named-behavior ::= behavior id : assumes-clause* requires-clause*
                 simple-clause* comm-clause*
comm-clause ::= mpi uses term (, term)* ;
collective-contract ::= mpi collective(term):requires-clause* simple-clause*
                     waitsfor-clause* mpi-named-behavior* completeness-clause*
mpi-named-behavior ::= behavior id : assumes-clause* requires-clause*
                    simple-clause* waitsfor-clause*

```

**Fig. 3.** Grammar for ACSL function contracts, extended for MPI.

function. A user function may be implemented using one or more communicators, point-to-point operations, and MPI collectives.

The top level grammar is given in Figure 3. A function contract begins with a sequence of distinct behaviors, each with an assumption that specifies when that behavior is active. Clauses in the global contract scope preceding the first named behavior are thought of as comprising a single behavior with a unique name and assumption *true*. The behaviors may be followed by **disjoint behaviors** and **complete behaviors** clauses, which encode claims that the assumptions are pairwise disjoint, and their disjunction is equivalent to *true*, respectively. All of this is standard ACSL, and we refer to it as the *sequential part* of the contract.

A new kind of clause, the *comm-clause*, may occur in the sequential part. A comm-clause begins “**mpi uses**” and is followed by a list of terms of type `MPI_Comm`. Such a clause specifies a guarantee that no communication will take place on a communicator *not* in the list. When multiple comm-clauses occur within a behavior, it is as if the lists were appended into one.

Collective contracts appear after the sequential part. A collective contract begins “**mpi collective**” and names a communicator *c* which provides the context for the contract; *c* must occur in a comm-clause from the sequential part. A collective contract on *c* encodes the claim that the function conforms to its contract (Definition 4) with the adjustment that all of the collective errors defined in Definition 2 are interpreted with respect to *c* only.

A collective contract may comprise multiple behaviors. As with the sequential part, clauses occurring in the collective contract before the first named behavior are considered to comprise a behavior with a unique name and assumption *true*.

**Type signatures.** The new logic type `mpi_sig_t` represents MPI type signatures. Its domain consists of all finite sequences of basic C types. As with all ACSL types, equality is defined and `==` and `!=` can be used on two such values in a logic specification. If *t* is a term of integer type and *s* is a term of type `mpi_sig_t`, then *t\*s* is a term of type `mpi_sig_t`. If the value of *t* is *n* and *n* ≥ 0, then *t\*s* denotes the result of concatenating the sequence of *s* *n* times.

**Operations on datatypes.** Two logic functions and one predicate are defined:

```

int \mpi_extent(MPI_Datatype datatype);
mpi_sig_t \mpi_sig(MPI_Datatype datatype);
\mpi_nonoverlapping(MPI_Datatype datatype);

```

The first returns the extent (in bytes) of a datatype. The second returns the type signature of the datatype. The predicate holds iff the type map of the datatype is nonoverlapping, a requirement for any communication buffer that receives data.

**Value sequences.** The domain of type `mpi_seq_t` consists of all finite sequences of pairs  $(t, v)$ , where  $t$  is a basic C type and  $v$  is a value of type  $t$ . Such a sequence represents the values stored in a communication buffer or message. Similar to the case with type signatures, we define multiplication of an integer with a value of type `mpi_seq_t` to be repeated concatenation.

**Communication buffers.** Type `mpi_buf_t` is a struct with fields `base` (of type `void*`), `count` (`int`), and `datatype` (`MPI_Datatype`). A value of this type specifies an MPI communication buffer and is created with the logic function

```
mpi_buf_t \mpi_buf(void * base, int count, MPI_Datatype datatype);
```

The ACSL predicate `\valid` is extended to accept arguments of type `mpi_buf_t` and indicates that the entire extent of the buffer is allocated memory; predicate `\valid_read` is extended similarly.

**Buffer arithmetic.** An integer and a buffer can be added or multiplied. Both operations are commutative. These are defined by

```

n * \mpi_buf(p, m, dt) == \mpi_buf(p, n * m, dt)
n + \mpi_buf(p, m, dt) == \mpi_buf((char*)p + n*\mpi_extent(dt), m, dt)

```

Multiplication corresponds to multiplying the size of a buffer by  $n$ . It is meaningful only when both  $n$  and  $m$  are nonnegative. Addition corresponds to shifting a buffer by  $n$  units, where a unit is the extent of the datatype `dt`. It is meaningful for any integer  $n$ .

**Buffer dereferencing.** The dereference operator `*` may take an `mpi_buf_t`  $b$  as an argument. The result is the value sequence (of type `mpi_seq_t`) obtained by reading the sequence of values from the buffer specified by  $b$ .

The term `*b` used in an `assigns` clause specifies that any of the memory locations associated to  $b$  may be modified; these are the bytes in the range  $p + m$  to  $p + m + \text{sizeof}(t) - 1$ , for some entry  $(t, m)$  in the type map of  $b$ .

The ACSL predicate `\separated` takes a comma-separated list of expressions, each of which denotes a set of memory locations. It holds if those sets are pairwise disjoint. We extend the syntax to allow expressions of type `mpi_buf_t` in the list; these expressions represent sets of memory locations as above.

**Terms.** The grammar for ACSL *terms* is extended:

```
term ::= \mpi_comm_rank | \mpi_comm_size | \mpi_on(term, term)
```

The term `\mpi_comm_size` is a constant, the number of processes in the communicator; `\mpi_comm_rank` is the rank of “this” process. In the term `\mpi_on(t, r)`,  $r$  must have integer type and is the rank of a process in the communicator. Term  $t$  is evaluated in the state of the process of rank  $r$ . For convenience, we define a macro `\mpi_agree(x)` which expands to `x==\mpi_on(x,0)`. This is used to say the value of  $x$  is the same on all processes.

**Reduction.** A predicate for reductions is defined:

```
\mpi_reduce(mpi_seq_t out, integer lo, integer hi,
            MPI_Op op, (integer)->mpi_seq_t in);
```

The predicate holds iff the value sequence `out` on this process is a point-wise reduction, using operator `op`, of the  $hi - lo$  value sequences `in(lo)`, `in(lo + 1)`, ..., `in(hi - 1)`. Note `in` is a function from `integer` to `mpi_seq_t`. We say *a* reduction, and not *the* reduction, because `op` may not be strictly commutative and associative (e.g., floating-point addition).

## 4 Evaluation

In this section we describe a prototype tool we developed for MPI collective contract verification, and experiments applying it to various example codes. All experimental artifacts, including the tool source code, are available online [3].

### 4.1 Collective Contract Examples

The first part of our evaluation involved writing contracts for a variety of collective functions. We started with the 17 MPI blocking collective functions specified in [43, Chapter 5]. These represent the most commonly used message-passing patterns, such as broadcast, scatter, gather, transpose, and reduce (fold). The MPI Standard is a precisely written natural language document, similar to the C Standard. We scrutinized each sentence in the description of each function and checked that it was reflected accurately in the contract.

Figure 4 shows the contract for the MPI collective function `MPI_Allreduce`. This function “combines the elements provided in the input buffer of each process... using the operator `op`” and “the result is returned to all processes” [43]. This guarantee is reflected in line 13. “The ‘in place’ option ... is specified by passing the value `MPI_IN_PLACE` to the argument `sendbuf` at all processes. In this case, the input data is taken at each process from the receive buffer, where it will be replaced by the output data.” This option is represented using two behaviors. These are just a few examples of the tight mapping between the natural language and the contract.

The only ambiguity we could not resolve concerned synchronization. The Standard is clear that collective operations may or may not impose barriers. It is less clear on whether certain forms of synchronization are implied by the semantics of the operation. For example, many users assume that a non-root

```

1 #define SBUF \mpi_buf(sbuf, count, dt)
2 #define RBUF \mpi_buf(rbuf, count, dt)
3 /*@ mpi uses comm; mpi collective(comm):
4   requires \valid(RBUF) && \mpi_nonoverlapping(dt);
5   requires \mpi_agree(count) && \mpi_agree(dt) && \mpi_agree(op) && count >= 0;
6   requires \separated(RBUF, { SBUF | int i; sbuf != MPI_IN_PLACE });
7   assigns *RBUF;
8   ensures \mpi_agree(*RBUF);
9   waitsfor { i | int i; 0 <= i < \mpi_comm_size };
10  behavior not_in_place:
11    assumes sbuf != MPI_IN_PLACE;
12    requires \mpi_agree(sbuf != MPI_IN_PLACE) && \valid_read(SBUF);
13    ensures \mpi_reduce(*RBUF, 0, \mpi_comm_size, op, \lambda integer t; \mpi_on(*SBUF, t));
14  behavior in_place:
15    assumes sbuf == MPI_IN_PLACE;
16    requires \mpi_agree(sbuf == MPI_IN_PLACE);
17    ensures
18      \mpi_reduce(*RBUF, 0, \mpi_comm_size, op, \lambda integer t; \mpi_on(\old(*RBUF), t));
19  disjoint behaviors; complete behaviors; */
20 int MPI_Allreduce(const void *sbuf, void *rbuf, int count, MPI_Datatype dt, MPI_Op op,
21                  MPI_Comm comm);

```

Fig. 4. The contract of the `MPI_Allreduce` function.

process must wait for the root in a broadcast, or that all-reduce necessarily entails a barrier. But these operations could be implemented with no synchronization when *count* is 0. (Similarly, a process executing all-reduce with *logical and* could return immediately if its contribution is *false*.) This issue has also been discussed in the MPI Forum [17]. Our `MPI_Allreduce` contract declares, on line 9, that barrier synchronization will occur, but other choices could also be encoded.

In addition to the MPI collectives, we wrote contracts for a selection of user-defined collectives from the literature, including:

1. `exchange`: “ghost cell exchange” in 1d-diffusion solver [56]
2. `diff1dIter`: computes one time step in 1d-diffusion [56]
3. `dotProd`: parallel dot-product procedure from Hypre [23]
4. `matmat`: matrix multiplication using a block-striped decomposition [50]
5. `oddEvenIter`: odd-even parallel sorting algorithm [30, 40].

We also implemented `cyc` of Figure 2 in MPI with contracts.

Fig. 5 shows the contract and the implementation for `dotProd`. The functions `hypre_MPI*` are simple wrappers for the corresponding MPI functions. The input vectors are block distributed. Each process gets its blocks and computes their inner product. The results are summed across processes with an all-reduce. The contract uses the ACSL `\sum` function to express the local result on a process (line 3) as well as the global result (line 13). Thus the contract is only valid if a real number model of arithmetic is used. This is a convenient and commonly-used assumption when specifying numerical code. We could instead use our predicate `\mpi_reduce` for a contract that holds in the floating-point model.

## 4.2 Bounded verification of collective contracts

For the second part of our evaluation, we developed a prototype tool for verifying that C/MPI collective procedures conform to their contracts. We used

```

1 #define hypre_ParVectorComm(vector) ((vector) -> comm)
2 #define PAR_SIZE x->local_vector->size * x->local_vector->num_vectors
3 #define LOCAL_RESULT \sum(0, PAR_SIZE-1, \lambda int t; \
4 x->local_vector->data[t] * y->local_vector->data[t])
5 /*@ requires \valid_read(x) && \valid_read(x->local_vector);
6     requires \valid_read(y) && \valid_read(y->local_vector);
7     requires \valid_read(x->local_vector->data + (0..PAR_SIZE-1));
8     requires \valid_read(y->local_vector->data + (0..PAR_SIZE-1));
9     requires x->local_vector->size > 0 && x->local_vector->num_vectors > 0;
10    mpi uses hypre_ParVectorComm(x);
11    mpi collective(hypre_ParVectorComm(x)):
12        assigns \nothing;
13        ensures \result == \sum(0, \mpi_comm_size-1,
14                               \lambda integer k; \mpi_on(LOCAL_RESULT, k));
15        waitsfor {i | int i; 0 <= i < \mpi_comm_size}; */
16 HYPRE_Real hypre_ParVectorInnerProd(hypre_ParVector *x, hypre_ParVector *y) {
17     MPI_Comm comm = hypre_ParVectorComm(x);
18     hypre_Vector *my_x = hypre_ParVectorLocalVector(x);
19     hypre_Vector *my_y = hypre_ParVectorLocalVector(y);
20     HYPRE_Real result = 0.0;
21     HYPRE_Real local_result = hypre_SeqVectorInnerProd(my_x, my_y);
22     hypre_MPI_Allreduce(&local_result, &result, 1, hypre_MPI_REAL,
23                        hypre_MPI_SUM, comm);
24     return result;
25 }

```

**Fig. 5.** The parallel `dotProd` function from Hypre [23], with contract.

CIVL, a symbolic execution and model checking framework [55] written in Java, because it provides a flexible intermediate verification language and it already has strong support for concurrency and MPI [42]. We created a branch of CIVL and modified the Java code in several ways, which we summarize here.

We modified the front-end to accept contracts in our extended version of ACSL. This required expanding the grammar, adding new kinds of AST nodes, and updating the analysis passes. Our prototype can therefore parse and perform basic semantic checks on contracts.

We then added several new primitives to the intermediate language to support the formal concepts described in Section 2. For example, in order to evaluate pre- and postconditions using `\mpi_on` expressions, we added a type for *collective state*, with operations to take a “snapshot” of a process state and merge snapshots into a program state, in order to check collective conditions.

Finally, we implemented a *transformer*, which consumes a C/MPI program annotated with contracts and the name of the function  $f$  to be verified. It generates a program similar to  $\overline{P^f}$  (Section 2.4). This program has a driver that initializes the global variables and arguments for  $f$  to arbitrary values constrained only by  $f$ ’s precondition, using CIVL’s `$assume` statement. The body of a collective function  $g$  used by  $f$  is replaced by code of the form

```

wait(waitsfor(g)); $assert(precondition); $havoc(assigns(g));
wait(waitsfor(g)); $assume(postcondition);

```

where `wait` is implemented using CIVL primitive `$when`, which blocks until a condition holds. When the CIVL verifier is applied to this program, it explores all simulations of  $f$ , verifying they terminate and are free of collective errors. By Thm. 1, the verifier can prove, for a bounded number of processes,  $f$  conforms.

function	states	prover	time(s)	function	states	prover	time(s)
<b>g</b> (cyc)	3,577	6	7	<b>allgather</b>	14,636	356	88
<b>allreduceDR</b>	7,406	15	10	<b>reduce</b>	118,388	54	110
<b>f</b> (cyc)	7,898	7	34	<b>scatter</b>	126,436	306	179
<b>oddEvenIter</b>	14,216	91	19	<b>gather</b>	126,834	259	183
<b>bcast</b>	29,311	80	36	<b>matmat</b>	8,360	275	216
<b>allreduce</b>	14,204	64	39	<b>reduceScatterNC</b>	264,473	261	594
<b>dotProd</b>	4,690	102	51	<b>reduceScatter</b>	211,799	501	1,233
<b>diffIdIter</b>	4,777	130	87	<b>exchange</b>	901,264	9,659	1,244

**Fig. 6.** Verification performance for  $nprocs \leq 5$ .

Our prototype has several limitations. It assumes no wildcard is used in the program. It does not check *assigns violation* for the verifying function. It assumes all communication uses standard mode blocking point-to-point functions and blocking MPI collective functions. Nevertheless, it can successfully verify a number of examples with nontrivial bounds on the number of processes.

For the experiment, we found implementations for several of the MPI collective functions. Some of these are straightforward; e.g., the implementation of `MPI_Allreduce` consists of calls to `MPI_Reduce` followed by a call to `MPI_Bcast`. Two of these implementations are more advanced: `allreduceDR` implements `MPI_Allreduce` using a double recursive algorithm; `reduceScatterNC` implements `MPI_Reduce_scatter` using an algorithm optimized for noncommutative reduction operations [12].

We applied our prototype to these collective implementations, using the contracts described in Section 4.1. We also applied it to the 5 user-defined collectives listed there. We were able to verify these contracts for up to 5 processes (no other input was bounded), using a Mac Mini with an M1 chip and 16GB memory. Results are given in Figure 6. For each problem, we give the number of states saved by CIVL, the number of calls to an external theorem prover (CVC4 [8] or Z3 [47]), and the verification time in seconds.

The times range from 3 seconds to 9 minutes. In general, time increases with the number of states and prover calls. Exceptions to this pattern occur when prover queries are very complex and the prover times out—two seconds in our case. For example, `matmat`, whose queries involve integer multiplications and uninterpreted functions, times out often. It is slower than most of the test cases despite a smaller state space.

Comparing `reduceScatter` with `reduceScatterNC`, it is noteworthy that verifying the simple implementation costs more than the advanced version. This is because the simple implementation re-uses verified collective functions. Reasoning about the contracts of those functions may involve expensive prover calls.

For `exchange`, nearly one million states are saved though its implementation involves only two MPI point-to-point calls. This is due to the generality of its contract. A process communicates with its left and right “neighbors” in this function. The contract assumes that the neighbors of a process can be any two



processes—as long as each pair of processes agree on whether they are neighbors. Hence there is combinatorial explosion generating the initial states.

For each example, we made erroneous versions and confirmed that CIVL reports a violation or “unknown” result.

## 5 Related Work

The ideas underlying code contracts originate in the work of Floyd on formal semantics [26], the proof system of Hoare [29], the specification system Larch [27], and Meyer’s work on Eiffel [44, 45]. Contract systems have been developed for many other languages, including Java [25, 32, 37], Ada [5], C# [7], and C [10, 18].

Verification condition generation (VCG) [6, 20, 25, 38] and symbolic execution [35, 36, 49] are two techniques used to verify that code conforms to a contract. *Extended static checking* is an influential VCG approach for Java [25, 32, 38]. Frama-C’s WP plugin [9, 18] is a VCG tool for ACSL-annotated C programs, based on the Why platform [24]. The Kiasan symbolic execution platform [19] has been applied to both JML and Spark contracts [11].

Several contract systems have been developed for shared memory concurrency. The VCC verifier [15, 16, 46] takes a contract approach, based on object invariants in addition to pre- and postconditions, to shared-memory concurrent C programs. VeriFast is a deductive verifier for multithreaded C and Java programs [31]. Its contract language is based on concurrent separation logic [14]. These systems focus on issues, such as ownership and permission, that differ from those that arise in distributed computing.

For distributed concurrency, type-theoretic approaches based on *session types* [48, 52, 57] are used to describe communication protocols; various techniques verify an implementation conforms to a protocol. ParTypes [39] applies this approach to C/MPI programs using a user-written protocol that specifies the sequence of messages transmitted in an execution. Conformance guarantees deadlock-freedom for an arbitrary number of processes. However, ParTypes protocols cannot specify programs with wildcards or functional correctness, and they serve a different purpose than our contracts. Our goal is to provide a public contract for a collective procedure—the messages transmitted are an implementation detail that should remain “hidden” to the extent possible.

There are a number of correctness tools for MPI programs, including the dynamic model checkers ISP [58] and DAMPI [59], the static analysis tool MPI-Checker [22], and the dynamic analysis tool MUST [28]. These check for certain pre-defined classes of defects, such as deadlocks and incorrectly typed receive statements; they are not used to specify or verify functional correctness.

Ashcroft introduced the idea of verifying parallel programs by showing every atomic action preserves a global invariant [4]. This approach is applied to a simple message-passing program in [41] using Frama-C+WP and ghost variables to represent channels. The contracts are quite complicated; they are also a bespoke solution for a specific problem, rather than a general language. However, the approach applies to non-collective as well as collective procedures.

A parallel program may also be specified by a functionally equivalent sequential version [54]. This works for whole programs which consume input and produce output, but it seems less applicable to individual collective procedures.

*Assume-guarantee reasoning* [1, 21, 33, 34] is another approach that decomposes along process boundaries. This is orthogonal to our approach, which decomposes along procedure boundaries.

## 6 Discussion

We have summarized a theory of contracts for collective procedures in a toy message-passing language. We have shown how this theory can be realized for C programs that use MPI using a prototype contract-checking tool. The approach is applicable to programs that use standard-mode blocking point-to-point operations, blocking MPI collective functions, multiple communicators, user-defined datatypes, pointers, pointer arithmetic, and dynamically allocated memory. We have used it to fully specify all of the MPI blocking collective functions, and several nontrivial user-defined collective functions.

MPI’s nonblocking operations are probably the most important and widely-used feature of MPI not addressed here. In fact, there is no problem specifying a collective procedure that uses nonblocking operations, as long as the procedure completes all of those operations before returning. For such procedures, the nonblocking operations are another implementation detail that need not be mentioned in the public interface. However, some programs may use one procedure to post nonblocking operations, and another procedure to complete them; this is in fact the approach taken by the new MPI “nonblocking collective” functions [43, Sec. 5.12]. The new “neighborhood collectives” [43, Sec. 7.6] may also require new abstractions and contract primitives.

Our theory assumes no use of `MPI_ANY_SOURCE` “wildcard” receives. It is easy to construct counterexamples to Theorem 1 for programs that use wildcards. New conceptual elements will be required to ensure a collective procedure implemented with wildcards will always behave as expected.

Our prototype tool for verifying conformance to a contract uses symbolic execution and bounded model checking techniques. It demonstrates the feasibility of this approach, but can only “verify” with small bounds placed on the number of processes. It would be interesting to see if the verification condition generation (VCG) approach can be applied to our contracts, so that they could be verified without such bounds. This would require a kind of Hoare calculus for message-passing parallel programs, and/or a method for specifying and verifying a global invariant.

One could also ask for runtime verification of collective contracts. This is an interesting problem, as the assertions relate the state of multiple processes, so checking them would require communication.

## References

1. Abadi, M., Lamport, L.: Conjoining specifications. *ACM Trans. Program. Lang. Syst.* **17**(3), 507–535 (May 1995). <https://doi.org/10.1145/203095.201069>
2. Alur, R., Bouajjani, A., Esparza, J.: *Model Checking Procedural Programs*, chap. 17, pp. 541–572. Springer, Cham (2018). [https://doi.org/10.1007/978-3-319-10575-8\\_17](https://doi.org/10.1007/978-3-319-10575-8_17)
3. Anonymous: Collective contracts for MPI: Extended report, artifact, and source code. Tech. Rep. 2023-0001, Anonymous Institution (Aug 2023), <https://github.com/Anonymous4Anonymous/MPI-Contract>
4. Ashcroft, E.A.: Proving assertions about parallel programs. *J. Comput. Syst. Sci.* **10**(1), 110–135 (Feb 1975). [https://doi.org/10.1016/S0022-0000\(75\)80018-3](https://doi.org/10.1016/S0022-0000(75)80018-3)
5. Barnes, J.: *High Integrity Software: The SPARK Approach to Safety and Security*. Addison-Wesley, Boston, MA, USA (2003)
6. Barnett, M., Chang, B.Y.E., DeLine, R., Jacobs, B., Leino, K.R.M.: Boogie: A modular reusable verifier for object-oriented programs. In: de Boer, F.S., Bonsangue, M.M., Graf, S., de Roever, W.P. (eds.) *Formal Methods for Components and Objects*, 4th International Symposium, FMCO 2005, Amsterdam, The Netherlands, November 1-4, 2005, Revised Lectures. *Lecture Notes in Computer Science*, vol. 4111, pp. 364–387. Springer, Springer, Berlin, Heidelberg (2005). [https://doi.org/10.1007/11804192\\_17](https://doi.org/10.1007/11804192_17)
7. Barnett, M., Fähndrich, M., Leino, K.R.M., Müller, P., Schulte, W., Venter, H.: Specification and verification: the Spec# experience. *Commun. ACM* **54**(6), 81–91 (2011). <https://doi.org/10.1145/1953122.1953145>
8. Barrett, C., Conway, C.L., Deters, M., Hadarean, L., Jovanović, D., King, T., Reynolds, A., Tinelli, C.: CVC4. In: *International Conference on Computer Aided Verification*. pp. 171–177. Springer, Berlin, Heidelberg (2011), <http://dl.acm.org/citation.cfm?id=2032305.2032319>
9. Baudin, P., Bobot, F., Correnson, L., Dargaye, Z., Blanchard, A.: WP plugin manual: Frama-C 22.0 (titanium) (2020), <https://frama-c.com/download/frama-c-wp-manual.pdf>
10. Baudin, P., Cuoq, P., Filliâtre, J.C., Marché, C., Monate, B., Moy, Y., Prevosto, V.: ACSL: ANSI/ISO C Specification Language, version 1.16. <http://frama-c.com/download/acsl-1.16.pdf> (2012)
11. Belt, J., Hatcliff, J., Robby, Chalin, P., Hardin, D., Deng, X.: Bakar Kiasan: Flexible contract checking for critical systems using symbolic execution. In: Bobaru et al. [13], pp. 58–72. [https://doi.org/10.1007/978-3-642-20398-5\\_6](https://doi.org/10.1007/978-3-642-20398-5_6)
12. Bernaschi, M., Iannello, G., Lauria, M.: Efficient Implementation of Reduce-scatter in MPI. In: *Proceedings of the 10th Euromicro Conference on Parallel, Distributed and Network-based Processing*. pp. 301–308. EUROMICRO-PDP’02, IEEE Computer Society, Washington, DC, USA (2002), <http://dl.acm.org/citation.cfm?id=1895489.1895529>
13. Bobaru, M.G., Havelund, K., Holzmann, G.J., Joshi, R. (eds.): *NASA Formal Methods - Third International Symposium, NFM 2011, Pasadena, CA, USA, April 18-20, 2011. Proceedings, LNCS*, vol. 6617. Springer (2011). <https://doi.org/10.1007/978-3-642-20398-5>
14. Brookes, S.: A semantics for concurrent separation logic. *Theoretical Computer Science* **375**(1), 227–270 (2007). <https://doi.org/10.1016/j.tcs.2006.12.034>, festschrift for John C. Reynolds’s 70th birthday

15. Cohen, E., Dahlweid, M., Hillebrand, M., Leinenbach, D., Moskal, M., Santen, T., Schulte, W., Tobies, S.: VCC: A practical system for verifying concurrent C. In: Berghofer, S., Nipkow, T., Urban, C., Wenzel, M. (eds.) Theorem Proving in Higher Order Logics, LNCS, vol. 5674, pp. 23–42. Springer, Berlin, Heidelberg (2009). [https://doi.org/10.1007/978-3-642-03359-9\\_2](https://doi.org/10.1007/978-3-642-03359-9_2)
16. Cohen, E., Moskal, M., Schulte, W., Tobies, S.: Local verification of global invariants in concurrent programs. In: Touili, T., Cook, B., Jackson, P. (eds.) Computer Aided Verification, 22nd International Conference, CAV 2010, Edinburgh, UK, July 15–19, 2010. Proceedings. LNCS, vol. 6174, pp. 480–494. Springer, Berlin, Heidelberg (2010). [https://doi.org/10.1007/978-3-642-14295-6\\_42](https://doi.org/10.1007/978-3-642-14295-6_42)
17. Community, M.: Collective Synchronization. <https://github.com/mpi-forum/mpi-issues/issues/257> (2020), accessed Aug. 13, 2021
18. Cuoq, P., Kirchner, F., Kosmatov, N., Prevosto, V., Signoles, J., Yakobowski, B.: Frama-C — a software analysis perspective. In: Eleftherakis, G., Hinchey, M., Holcombe, M. (eds.) Software Engineering and Formal Methods - 10th International Conference, SEFM 2012, Thessaloniki, Greece, October 1–5, 2012. Proceedings. LNCS, vol. 7504, pp. 233–247. Springer-Verlag, Berlin, Heidelberg (2012). [https://doi.org/10.1007/978-3-642-33826-7\\_16](https://doi.org/10.1007/978-3-642-33826-7_16)
19. Deng, X., Lee, J., Robby: Bogor/Kiasan: A k-bounded symbolic execution for checking strong heap properties of open systems. In: 21st IEEE/ACM International Conference on Automated Software Engineering (ASE 2006), 18–22 September 2006, Tokyo, Japan. pp. 157–166. IEEE Computer Society, USA (2006). <https://doi.org/10.1109/ASE.2006.26>
20. Detlefs, D.L., Leino, K.R.M., Nelson, G., Saxe, J.B.: Extended static checking. Research Report 159, COMPAQ Systems Research Center (Dec 1998), <https://www.hpl.hp.com/techreports/Compaq-DEC/SRC-RR-159.pdf>
21. Dingel, J.: Computer-assisted assume/guarantee reasoning with VeriSoft. In: Proceedings of the 25th International Conference on Software Engineering. pp. 138–148. ICSE ’03, IEEE Computer Society, Washington, DC, USA (2003). <https://doi.org/10.1109/ICSE.2003.1201195>
22. Droste, A., Kuhn, M., Ludwig, T.: MPI-Checker: Static analysis for MPI. In: Proceedings of the Second Workshop on the LLVM Compiler Infrastructure in HPC. pp. 3:1–3:10. LLVM ’15, ACM, New York (2015). <https://doi.org/10.1145/2833157.2833159>
23. Falgout, R.D., Yang, U.M.: *hypre*: A Library of High Performance Preconditioners. In: Sloot, P.M.A., Hoekstra, A.G., Tan, C.J.K., Dongarra, J.J. (eds.) Computational Science — ICCS 2002. pp. 632–641. Springer Berlin Heidelberg, Berlin, Heidelberg (2002), [https://doi.org/10.1007/3-540-47789-6\\_66](https://doi.org/10.1007/3-540-47789-6_66)
24. Filliâtre, J.C., Paskevich, A.: Why3: Where programs meet provers. In: Felleisen, M., Gardner, P. (eds.) Proceedings of the 22nd European Conference on Programming Languages and Systems (ESOP’13). pp. 125–128. Springer-Verlag, Berlin, Heidelberg (2013). [https://doi.org/10.1007/978-3-642-37036-6\\_8](https://doi.org/10.1007/978-3-642-37036-6_8)
25. Flanagan, C., Leino, K.R.M., Lillibridge, M., Nelson, G., Saxe, J.B., Stata, R.: Extended static checking for Java. In: Knoop, J., Hendren, L.J. (eds.) Proceedings of the 2002 ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI), Berlin, Germany, June 17–19, 2002. pp. 234–245. Association for Computing Machinery, New York, NY, USA (2002). <https://doi.org/10.1145/512529.512558>
26. Floyd, R.W.: Assigning meanings to programs. Mathematical Aspects of Computer Science **19**, 19–32 (1967)

27. Guttag, J.V., Horning, J.J., Wing, J.M.: The Larch family of specification languages. *IEEE Software* **2**(5), 24–36 (1985). <https://doi.org/10.1109/MS.1985.231756>
28. Hilbrich, T., Protze, J., Schulz, M., de Supinski, B.R., Müller, M.S.: MPI runtime error detection with MUST: Advances in deadlock detection. In: Hollingsworth, J.K. (ed.) *International Conference on High Performance Computing Networking, Storage and Analysis, SC '12*, Salt Lake City, UT, USA - November 11–15, 2012. pp. 30:1–30:11. IEEE Computer Society Press, Los Alamitos, CA, USA (2012). <https://doi.org/10.1109/SC.2012.79>
29. Hoare, C.A.R.: An axiomatic basis for computer programming. *Communications of the ACM* **12**(10), 576–580 (1969). <https://doi.org/10.1145/363235.363259>
30. Huisman, M., Monahan, R., Müller, P., Mostowski, W., Ulbrich, M.: *VerifyThis 2017: A Program Verification Competition*. Tech. Rep. Karlsruhe Reports in Informatics 2017,10, Karlsruhe Institute of Technology, Faculty of Informatics (2017). <https://doi.org/10.5445/IR/1000077160>
31. Jacobs, B., Piessens, F.: Expressive modular fine-grained concurrency specification. In: *Proceedings of the 38th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*. pp. 271–282. POPL '11, Association for Computing Machinery, New York, NY, USA (2011). <https://doi.org/10.1145/1926385.1926417>
32. James, P.R., Chalin, P.: Faster and more complete extended static checking for the Java Modeling Language. *Journal of Automated Reasoning* **44**, 145–174 (2010). <https://doi.org/10.1007/s10817-009-9134-9>
33. Jones, C.B.: Tentative steps toward a development method for interfering programs. *ACM Trans. Program. Lang. Syst.* **5**(4), 596–619 (Oct 1983). <https://doi.org/10.1145/69575.69577>
34. Jones, C.B.: Specification and design of (parallel) programs. In: Mason, R.E.A. (ed.) *Information Processing 83, Proceedings of the IFIP 9th World Computer Congress, Paris, France, September 19–23, 1983*. pp. 321–332. North-Holland/IFIP, Newcastle University (1983)
35. Khurshid, S., Păsăreanu, C.S., Visser, W.: Generalized symbolic execution for model checking and testing. In: Garavel, H., Hatcliff, J. (eds.) *Tools and Algorithms for the Construction and Analysis of Systems, 9th International Conference, TACAS 2003, Held as Part of the Joint European Conferences on Theory and Practice of Software, ETAPS 2003, Warsaw, Poland, April 7–11, 2003, Proceedings*. LNCS, vol. 2619, pp. 553–568. Springer-Verlag, Berlin, Heidelberg (2003). [https://doi.org/10.1007/3-540-36577-X\\_40](https://doi.org/10.1007/3-540-36577-X_40)
36. King, J.C.: Symbolic execution and program testing. *Communications of the ACM* **19**(7), 385–394 (1976). <https://doi.org/10.1145/360248.360252>
37. Leavens, G.T., Baker, A.L., Ruby, C.: Preliminary design of JML: a behavioral interface specification language for Java. *SIGSOFT Softw. Eng. Notes* **31**(3), 1–38 (2006). <https://doi.org/10.1145/1127878.1127884>
38. Leino, K.R.M.: Extended static checking: A ten-year perspective. In: Wilhelm, R. (ed.) *Informatics - 10 Years Back. 10 Years Ahead*. LNCS, vol. 2000, pp. 157–175. Springer, Berlin, Heidelberg (2001). [https://doi.org/10.1007/3-540-44577-3\\_11](https://doi.org/10.1007/3-540-44577-3_11)
39. López, H.A., Marques, E.R.B., Martins, F., Ng, N., Santos, C., Vasconcelos, V.T., Yoshida, N.: Protocol-based verification of message-passing parallel programs. In: Aldrich, J., Eugster, P. (eds.) *Proceedings of the 2015 ACM SIGPLAN International Conference on Object-Oriented Programming, Systems, Languages, and Applications, OOPSLA 2015, part of SPLASH 2015, Pittsburgh, PA,*

- USA, October 25–30, 2015. pp. 280–298. ACM, New York, NY, USA (2015). <https://doi.org/10.1145/2814270.2814302>
40. Luo, Z., Siegel, S.F.: Symbolic execution and deductive verification approaches to VerifyThis 2017 challenges. In: Margaria, T., Steffen, B. (eds.) *Leveraging Applications of Formal Methods, Verification and Validation (ISoLA 2018)*, Proceedings, Part II: Verification. LNCS, vol. 11245, pp. 160–178. Springer, International Publishing (2018)
  41. Luo, Z., Siegel, S.F.: Towards deductive verification of message-passing parallel programs. In: Laguna, I., Rubio-González, C. (eds.) *2018 IEEE/ACM 2nd International Workshop on Software Correctness for HPC Applications (Correctness)*. pp. 59–68. IEEE, USA (2018). <https://doi.org/10.1109/Correctness.2018.00012>
  42. Luo, Z., Zheng, M., Siegel, S.F.: Verification of MPI programs using CIVL. In: *Proceedings of the 24th European MPI Users’ Group Meeting*. pp. 6:1–6:11. EuroMPI ’17, ACM, New York, NY, USA (2017). <https://doi.org/10.1145/3127024.3127032>
  43. Message-Passing Interface Forum: MPI: A Message-Passing Interface standard, version 3.1 (Jun 2015), <https://www.mpi-forum.org/docs/mpi-3.1/mpi31-report.pdf>
  44. Meyer, B.: Applying “Design by Contract”. *IEEE Computer* **25**(10), 40–51 (1992). <https://doi.org/10.1109/2.161279>
  45. Meyer, B., Nerson, J.M., Matsuo, M.: EIFFEL: Object-oriented design for software engineering. In: Nichols, H.K., Simpson, D. (eds.) *ESEC ’87, 1st European Software Engineering Conference*, Strasbourg, France, September 9–11, 1987, Proceedings. LNCS, vol. 289, pp. 221–229. Springer, Berlin, Heidelberg (1987). <https://doi.org/10.1007/BFb0022115>
  46. Moskal, M.: Verifying functional correctness of C programs with VCC. In: Bobaru et al. [13], pp. 56–57. [https://doi.org/10.1007/978-3-642-20398-5\\_5](https://doi.org/10.1007/978-3-642-20398-5_5)
  47. de Moura, L., Bjørner, N.: Z3: An Efficient SMT Solver. In: Ramakrishnan, C.R., Rehof, J. (eds.) *Tools and Algorithms for the Construction and Analysis of Systems*. pp. 337–340. Springer Berlin Heidelberg, Berlin, Heidelberg (2008)
  48. Ng, N., Yoshida, N., Honda, K.: Multiparty Session C: Safe parallel programming with message optimisation. In: Furia, C.A., Nanz, S. (eds.) *Objects, Models, Components, Patterns*. LNCS, vol. 7304, pp. 202–218. Springer, Berlin, Heidelberg (2012). [https://doi.org/10.1007/978-3-642-30561-0\\_15](https://doi.org/10.1007/978-3-642-30561-0_15)
  49. Păsăreanu, C., Visser, W.: A survey of new trends in symbolic execution for software testing and analysis. *Int. J. Softw. Tools Technol. Transf.* **11**(4), 339–353 (Oct 2009). <https://doi.org/10.1007/s10009-009-0118-1>
  50. Quinn, M.: *Parallel Programming in C with MPI and OpenMP*. McGraw-Hill, (2004)
  51. Romano, P.K., Horelik, N.E., Herman, B.R., Nelson, A.G., Forget, B., Smith, K.: OpenMC: A state-of-the-art Monte Carlo code for research and development. *Annals of Nuclear Energy* **82**, 90–97 (2015). <https://doi.org/10.1016/j.anucene.2014.07.048>
  52. Scalas, A., Yoshida, N., Benussi, E.: Verifying message-passing programs with dependent behavioural types. In: *Proceedings of the 40th ACM SIGPLAN Conference on Programming Language Design and Implementation*. pp. 502–516. PLDI 2019, Association for Computing Machinery, New York, NY, USA (2019). <https://doi.org/10.1145/3314221.3322484>
  53. Siegel, S.F., Avrunin, G.S.: Modeling wildcard-free MPI programs for verification. In: *Proceedings of the Tenth ACM SIGPLAN Symposium on Principles and Prac-*

- tice of Parallel Programming. pp. 95–106. PPOPP '05, Association for Computing Machinery, New York, NY, USA (2005). <https://doi.org/10.1145/1065944.1065957>
54. Siegel, S.F., Mironova, A., Avrunin, G.S., Clarke, L.A.: Combining symbolic execution with model checking to verify parallel numerical programs. *ACM Transactions on Software Engineering and Methodology* **17**(2), Article 10, 1–34 (2008). <https://doi.org/10.1145/1348250.1348256>
  55. Siegel, S.F., Zheng, M., Luo, Z., Zirkel, T.K., Marianiello, A.V., Edenhofner, J.G., Dwyer, M.B., Rogers, M.S.: CIVL: The Concurrency Intermediate Verification Language. In: *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*. pp. 61:1–61:12. SC '15, ACM, New York (2015), <http://doi.acm.org/10.1145/2807591.2807635>
  56. Siegel, S.F., Zirkel, T.K.: FEVS: A Functional Equivalence Verification Suite for high performance scientific computing. *Mathematics in Computer Science* **5**(4), 427–435 (2011)
  57. Takeuchi, K., Honda, K., Kubo, M.: An interaction-based language and its typing system. In: Halatsis, C., Maritsas, D., Philokyprou, G., Theodoridis, S. (eds.) *PARLE'94 Parallel Architectures and Languages Europe*. LNCS, vol. 817, pp. 398–413. Springer, Berlin, Heidelberg (1994). [https://doi.org/10.1007/3-540-58184-7\\_118](https://doi.org/10.1007/3-540-58184-7_118)
  58. Vakkalanka, S., Gopalakrishnan, G., Kirby, R.M.: Dynamic verification of MPI programs with reductions in presence of split operations and relaxed orderings. In: Gupta, A., Malik, S. (eds.) *Computer Aided Verification, 20th International Conference, CAV 2008, Princeton, NJ, USA, July 7–14, 2008, Proceedings*. LNCS, vol. 5123, pp. 66–79. Springer, Berlin, Heidelberg (2008). [https://doi.org/10.1007/978-3-540-70545-1\\_9](https://doi.org/10.1007/978-3-540-70545-1_9)
  59. Vo, A., Aananthakrishnan, S., Gopalakrishnan, G., Supinski, B.R.d., Schulz, M., Bronevetsky, G.: A scalable and distributed dynamic formal verifier for MPI programs. In: *Proceedings of the 2010 ACM/IEEE International Conference for High Performance Computing, Networking, Storage and Analysis*. pp. 1–10. SC '10, IEEE Computer Society, Washington, DC, USA (2010). <https://doi.org/10.1109/SC.2010.7>
  60. Yang, U., Falgout, R., Park, J.: Algebraic Multigrid Benchmark, Version 00 (Aug 2017), <https://www.osti.gov/servlets/purl/1389816>

## A Proof of Theorem 1

### A.1 Preliminaries

We first establish some notation and vocabulary that will be used throughout this appendix.

*Sequences.* For a finite or infinite sequence  $\eta$ , write  $|\eta|$  for the length of  $\eta$ . The elements of  $\eta$  are numbered from 1. For sequences  $\eta$  and  $\xi$ ,  $\eta \leq \xi$  denotes that  $\eta$  is a prefix of, or equal to,  $\xi$ .

*Associated transitions.* Let  $P$  be a program and  $\zeta = (s_0 \xrightarrow{t_1} s_1 \xrightarrow{t_2} \dots)$  an execution of a procedure in  $P$ . Suppose  $\zeta$  is collective consistent. Choose  $p \in \text{PID}$



such that  $T_q(\zeta) \leq T_p(\zeta)$  for all  $q \in \text{PID}$ . The *number of collective events in  $\zeta$* , denoted  $\text{ncol}(\zeta)$ , is  $|T_p(\zeta)|$ . If  $1 \leq k \leq \text{ncol}(\zeta)$ , the  *$k$ -th collective event of  $\zeta$*  is the  $k$ -th symbol of  $T_p(\zeta)$ . By collective consistency, these definitions are independent of the choice of  $p$ .

For any integer  $k \geq 1$ , let

$$Q(\zeta, k) = \{p \in \text{PID} \mid |T_p(\zeta)| \geq k\}.$$

A process  $p$  is in  $Q(\zeta, k)$  if  $p$  performs at least  $k$  collective actions in  $\zeta$ .

Let  $k \in 1..\text{ncol}(\zeta)$ . Then  $Q(\zeta, k) \neq \emptyset$ . For each  $p \in Q(\zeta, k)$ , there is a unique  $i \in 1..\text{len}(\zeta)$  such that  $t_i$  is the  $k$ -th collective action of  $p$  in  $\zeta$ . We say the transition  $s_{i-1} \xrightarrow{t_i} s_i$  is *associated to the  $k$ -th collective event of  $\zeta$* .

## A.2 Commuting Actions

The key to the proof is the ability to transpose commuting adjacent actions from different processes. We say transitions  $t_1$  and  $t_2$  are *adjoining* if the destination state of  $t_1$  is the source state of  $t_2$ . The basic fact is captured by the following:

**Lemma 1.** *Let  $P$  be a program. Suppose  $s_0 \xrightarrow{t_1} s_1 \xrightarrow{t_2} s_2$  are adjoining transitions in  $P$ . Assume  $t_1$  and  $t_2$  are in different processes and it is not the case that  $t_1$  is a send of a message received by  $t_2$ . Then there exists a state  $s'_1$  of  $P$  such that  $s_0 \xrightarrow{t_2} s'_1 \xrightarrow{t_1} s_2$  are adjoining transitions.*

*Proof.* An assignment, call, or return updates the state of the invoking process only. A send or receive updates the states of the invoking process (the program counter) and a single message channel. An enabled receive in one process can never be disabled by an action from another process. Furthermore, the variables in a guard are local to one process (as there are no shared variables), so the truth value of a guard is not impacted by an action from another process. Therefore an assignment, call, return, or guard in one process commutes with any action from another process.

A send from one process commutes with a send from another process as they access separate channels. Likewise, a send from one process  $p$  commutes with a receive from another process  $q$  if the destination of the send statement is not  $q$  or the source in the receive statement is not  $p$ .

So assume one of the statements is a send in process  $p$  to process  $q$ , and the other is a receive in process  $q$  from  $p$ . There are two cases:

1. There is at least one buffered message in the  $p \rightarrow q$  channel in state  $s_0$ . In this case, the send enqueues a new message onto the queue and the receive dequeues the oldest message in the queue, and it is clear these two actions commute.
2. The  $p \rightarrow q$  channel is empty in  $s_0$ . In this case,  $t_1$  must be the send, and  $t_2$  therefore receives the message enqueued by  $t_1$ , contradicting the hypothesis of this lemma.



**Definition 5.** Let  $P$  be a program and  $f$  a procedure in  $P$ . Let  $\sim$  be the smallest equivalence relation on the set of executions of  $f$  such that  $\zeta \sim \eta$  if  $\eta$  is obtained from  $\zeta$  by transposing two adjacent transitions satisfying the hypotheses of Lemma 1. We say  $\zeta$  is equivalent to  $\eta$  if  $\zeta \sim \eta$ .

Hence, if  $\zeta \sim \eta$  then  $\eta$  can be obtained from  $\zeta$  by performing a finite sequence of transpositions of commuting adjacent transitions.

**Lemma 2.** Let  $P$  be a program and  $f$  a procedure in  $P$ . Suppose  $\zeta$  and  $\eta$  are executions of  $f$  and  $\zeta \sim \eta$ . Then all of the following hold:

1.  $\zeta$  commits a consistency, precondition, postcondition, assigns, segment boundary, or unreceived message violation iff  $\eta$  commits a violation of one of those kinds.
2.  $\zeta$  deadlocks iff  $\eta$  deadlocks.
3.  $\zeta$  is finite iff  $\eta$  is finite.

In the remainder of this section we sketch the proof of Lemma 2. For the most part, the proof is straightforward, but tedious, since each kind of error must be considered in turn.

*Consistency.* The transitions involved in each transposition come from different processes. Therefore  $T_p(\zeta) = T_p(\eta)$  for all  $p \in \text{PID}$ . It follows that  $\zeta$  is collective consistent iff  $\eta$  is collective consistent.

Hence if  $\zeta$  has a consistency violation, we are done. So assume  $\zeta$  (and therefore  $\eta$ ) is collective consistent.

*Precondition.* Suppose  $1 \leq r \leq \text{ncol}(\zeta)$  and the  $r$ -th collective event of  $\zeta$  is  $e^g$  where  $g \in \mathcal{C}$ . For  $q \in Q(\zeta, r)$ , let  $\sigma_{r,q}$  be the process state of process  $q$  just after executing the  $r$ -th collective action in  $q$ . Define  $\alpha': Q \rightarrow (V_f \rightarrow \mathbb{Z})$  by setting  $\alpha'(q)(v)$  to be the value of  $v$  in  $\sigma_{r,q}$ . Hence  $\alpha'$  is the partial  $g$ -valuation used to determine whether a precondition violation occurs when the last process in  $Q(\zeta, r)$  to execute its  $r$ -th collective transition executes that transition.

We claim that a precondition violation involving the  $r$ -th collective event occurs in  $\zeta$  iff there is no satisfying extension of  $\alpha'$ , i.e., an extension to a  $g$ -valuation  $\alpha$  such that  $\llbracket \text{pre}(g) \rrbracket_{\alpha, p}$  holds for all  $p \in \text{PID}$ . Indeed, if there is no such extension, then a precondition violation occurred with the last process in  $Q$  to execute its associated transition. Conversely, if a precondition violation occurs at some (possibly earlier) associated transition, then the set of processes used to determine a violation is some subset  $Q' \subseteq Q$ , the partial valuation used is  $\alpha'|_{Q'}$ , and there is no satisfying extension of  $\alpha'|_{Q'}$ . Since any extension of  $\alpha'$  is an extension of  $\alpha'|_{Q'}$ , it follows there is no satisfying extension of  $\alpha'$ , proving the claim.

This means the question of whether a precondition violation for collective event  $r$  occurs in  $\zeta$  depends only on the set of process states

$$\{\sigma_{r,q} \mid q \in Q(\zeta, r)\}.$$

This set is invariant under transposition of adjacent commuting transitions: while such a transposition may change the order in which processes execute their associated enter transitions, they do not change the final set of such processes or their local states at the time they enter. Hence a precondition violation occurs in  $\zeta$  iff such a violation occurs in  $\eta$ .

*Postcondition.* The argument is the same as that for preconditions, but uses the final sets of process states in the prestate and poststate for the call.

*Assigns.* For  $p \in \text{PID}$  we may project the execution  $\zeta$  onto  $p$  to obtain a “local execution” of process  $p$ , i.e., a sequence of process states and actions in  $p$ . This projection is invariant under transposition of commuting transitions. Since an assigns violation in process  $p$  depends only upon the sequence of process states and actions from  $p$ ,  $\zeta$  has an assigns violation iff  $\eta$  has an assigns violation.

*Segment boundary.* The first way a segment boundary violation could occur in  $\zeta$  is that process  $p$  receives a message in segment  $a$  that was sent by process  $q$  in segment  $b$ , with  $b > a$ . A transposition does not change the segment number in which the send or receive occurs, and since a transposition will not be applied if the first transition is a send and the second a receive of the message sent, this kind of violation will be preserved by a transposition.

The second kind of segment boundary violation occurs when process  $p$  in segment  $a$  sends to  $q$ , when  $q$  is in segment  $b$  and  $b > a$ . The only way this can be affected by transposition is if the first transition is a collective action of  $q$  and the second transition is the send by  $p$ . After transposing, the send takes place when  $p$  is at segment  $a$  and  $q$  is at segment  $b - 1$ . If this is not also a segment boundary violation, then  $a = b - 1$ , i.e., the send takes place when  $p$  and  $q$  are at the same segment. But then the collective action of  $q$  results in an unreceived message violation.

*Unreceived message.* Suppose there is an unreceived message for  $p$  sent in segment  $r - 1$  when  $p$  executes its  $r$ -th collective action. The only way a transposition could change this condition is if the first transition is the send that enqueues the message, and the second transition is the collective action by  $p$ . After transposing, there is no unreceived message when  $p$  executes its  $r$ -th collective action. However, the send now commits a segment boundary violation: the send to  $p$  is executed when  $p$  has a higher segment number than that of the sending process.

*Deadlock and finiteness.* Deadlock is clearly preserved since  $\zeta$  and  $\eta$  have the same final state. Transposition does not change the length of an execution, so in particular  $\zeta$  is finite iff  $\eta$  is finite.

This completes the proof sketch for Lemma 2.

*The issue with waitsfor.* The one kind of error not necessarily preserved by equivalence is a *waitsfor* violation. If all of the following hold:

- the first transition is an  $e^g$  in process  $p$  from state  $s_{i-1}$ ,

- the second transition is an  $x^g$  in process  $q$  from state  $s_i$ ,
- $\text{seg}_p(\zeta, i-1) = \text{seg}_q(\zeta, i)$ , and
- $p$  is in the wait set for  $q$ ,

then transposing results in a *waitsfor* violation where none existed before. However, if any of those conditions does not hold, the transposition cannot introduce a *waitsfor* violation.

### A.3 Observable Correctness

Let  $P$  be a program with no collective recursion. Let  $\mathcal{C}$  be the set of collective procedure names from  $P$ . Let  $f \in \mathcal{C}$  and assume all collective procedures used by  $f$  conform. Let  $\tilde{\mathcal{C}} \subseteq \mathcal{C}$  consist of  $f$  and all collective procedures used by  $f$ . We call the elements of  $\tilde{\mathcal{C}}$  the *observable collective procedures*.

Let  $\tilde{P}^f$  be the program that is the same as  $P^f$  except contracts are removed from all procedures in  $\mathcal{C} \setminus \tilde{\mathcal{C}}$ , so the collective procedures for  $\tilde{P}^f$  are precisely the observable collective procedures. An *observable collective action* is a collective action involving  $\tilde{\mathcal{C}}$ , i.e., an  $e^g$  or  $x^g$  action for  $g \in \tilde{\mathcal{C}}$ . Note that any execution  $\zeta$  of  $P^f$  is also an execution of  $\tilde{P}^f$ , but the notions of *segment* and *collective error* (and all the specific kinds of collective errors) depend on whether one considers  $\zeta$  to be an execution of  $P^f$  or of  $\tilde{P}^f$ .

**Definition 6.** Let  $\zeta$  be an execution of  $P^f$ . Let  $\tilde{\zeta}$  be the same execution, considered as an execution of  $\tilde{P}^f$ .

- We say  $\zeta$  commits an observable collective error at step  $i$  if  $\tilde{\zeta}$  commits a collective error at step  $i$ .
- If  $\zeta$  commits a collective error at step  $i$  but  $\tilde{\zeta}$  does not, then  $\zeta$  commits an internal collective error at step  $i$ .
- We say  $\zeta$  is observably correct if  $\tilde{\zeta}$  is correct.

Let  $\zeta = (s_0 \xrightarrow{t_1} s_1 \xrightarrow{t_2} \dots)$  be an execution of  $P^f$ .

It is not hard to see that if  $\zeta$  is free of collective errors then so is  $\tilde{\zeta}$ . Intuitively, the segment decomposition of  $\zeta$  refines that of  $\tilde{\zeta}$ , so, for example, a message that crosses a segment boundary in  $\tilde{\zeta}$  also crosses a segment boundary in  $\zeta$ . A pre- or postcondition, assigns, *waitsfor*, or unreceived message violation involving some  $g \in \tilde{\mathcal{C}}$  is also a violation in  $\zeta$ , as  $g \in \mathcal{C}$ .

Clearly,  $\zeta$  is finite if and only if  $\tilde{\zeta}$  is finite, and  $\zeta$  deadlocks if and only if  $\tilde{\zeta}$  deadlocks. Hence if  $\zeta$  is correct, then  $\zeta$  is observably correct.

At any state in  $\zeta$ , the number of observable collective procedures on the call stack of a process is at most 2. In fact, for any  $p \in \text{PID}$ , if  $T_p(\tilde{\zeta})$  is finite then it is a prefix of a sequence of the form

$$e^f e^{g_1} x^{g_1} \dots e^{g_m} x^{g_m} x^f,$$

where  $g_1, \dots, g_m$  are (not necessarily distinct) collective procedures used by  $f$ . If  $T_p(\tilde{\zeta})$  is infinite, it has the form

$$e^f e^{g_1} x^{g_1} e^{g_2} x^{g_2} \dots$$

If  $s$  is a state occurring on  $\zeta$ , we say  $p$  is *at an internal state at  $s$*  if there are exactly two observable collective procedures on the call stack for  $p$  in  $s$ , otherwise  $p$  is *at an external state at  $s$* .

Let  $i \in 1..\text{len}(\zeta)$ , so  $s_{i-1} \xrightarrow{t_i} s_i$  is a transition in  $\zeta$ . Let  $p$  be the process of  $t_i$ . Suppose  $r \geq 1$ . We say  $t_i$  is *in region  $r$*  of  $\zeta$  if  $\text{seg}_p(\tilde{\zeta}, i-1) = 2r$  or  $\text{seg}_p(\tilde{\zeta}, i) = 2r$ . In other words, all transitions in  $p$  from  $e^{gr}$  to  $x^{gr}$ , inclusive, belong to region  $r$ . We say  $t_i$  is *strictly in region  $r$*  if  $\text{seg}_p(\tilde{\zeta}, i-1) = 2r$  and  $\text{seg}_p(\tilde{\zeta}, i) = 2r$ , i.e.,  $t_i$  is in region  $r$  but is not the initial  $e^{gr}$  action nor the final  $x^{gr}$  action (if the final action occurs in  $\zeta$ ). We say  $t_i$  is *post-region  $r$*  if  $\text{seg}_p(\tilde{\zeta}, i-1) > 2r$ . We say  $t_i$  is *pre-region  $r$*  if  $\text{seg}_p(\tilde{\zeta}, i) < 2r$ .

#### A.4 Extraction of executions

Fix a collective-recursion-free program  $P$ . Let  $f$  be a collective procedure of  $P$ , and  $\zeta$  a finite execution of  $P^f$  with no observable collective error.

Suppose  $r \geq 1$  and there is some action in  $\zeta$  in region  $r$ . Let  $g = g_r$ . We now describe how one can extract an execution of  $P^g$  from region  $r$  of  $\zeta$ . It will be shown that this extracted execution starts from a  $\text{pre}(g)$ -state.

First, suppose

$$s_{i-1} \xrightarrow{t_i} s_i \xrightarrow{t_{i+1}} s_{i+1}$$

are two consecutive transitions in  $\zeta$ ,  $t_i$  and  $t_{i+1}$  are in two different processes,  $t_i$  is post-region  $r$ , and  $t_{i+1}$  is in region  $r$ . Let  $p$  be the process of  $t_i$  and  $q$  that of  $t_{i+1}$ . We have

$$\text{seg}_p(\tilde{\zeta}, i) \geq \text{seg}_p(\tilde{\zeta}, i-1) > 2r \geq \text{seg}_q(\tilde{\zeta}, i+1).$$

It follows that  $t_i$  cannot be a send of a message that is received by  $t_{i+1}$ , else a segment boundary violation occurs in  $\tilde{\zeta}$  at step  $i+1$ , contradicting the assumption that  $\tilde{\zeta}$  is collective error-free. By Lemma 1, we may transpose  $t_i$  and  $t_{i+1}$  to yield an execution  $\xi$  equivalent to  $\zeta$ .

We claim that  $\xi$  is also free of observable collective errors. By Lemma 2, we only need to show that  $\tilde{\xi}$  (i.e.,  $\xi$  considered as an execution of  $\tilde{P}^f$ ) has no *waitsfor* violation. According to the comments following the proof of Lemma 2, such a violation can be introduced only if  $\text{seg}_p(\tilde{\zeta}, i-1) = \text{seg}_q(\tilde{\zeta}, i)$ . However, in the transformation above,  $\text{seg}_p(\tilde{\zeta}, i-1) > \text{seg}_q(\tilde{\zeta}, i+1) \geq \text{seg}_q(\tilde{\zeta}, i)$ .

Similarly, if  $t_i$  is in region  $r$  and  $t_{i+1}$  is pre-region  $r$ , it can be seen that the two transitions commute, and the resulting execution is free of observable collective errors.

As  $\zeta$  is finite, by a finite number of transpositions of these two kinds, one can transform  $\zeta$  to an equivalent execution, free of observable collective errors, which is the concatenation  $\eta \circ \theta \circ \phi$ , where all transitions in  $\eta$  are pre-region  $r$ , all transitions in  $\theta$  are in region  $r$ , and all transitions in  $\phi$  are post-region  $r$ . Let  $s$  be the final state of  $\eta$ , which is also the initial state of  $\theta$ . Note that all channels are empty in  $s$ : otherwise, an observable unreceived message violation would occur.

The “extracted” execution of  $P^g$  is formed by modifying  $\theta$ . Let  $\chi = (s'_0 \xrightarrow{t'_1} s_1 \xrightarrow{t'_2} \dots)$  be a copy of  $\theta$  which we will modify in a sequence of steps. First, remove the portion of the call stack below the call to  $g$  in every state occurring in  $\chi$ . After this modification, every process has an empty call stack in the initial state  $s'_0$ . Then, insert an entry on the bottom of each stack in every state corresponding to the *main* function in  $P^g$ .

Let  $V_g$  be the global variables of  $P^g$ , which consists of the global variables of  $P$  together with a set of fresh variables corresponding to the formal parameters of  $g$ . A state of  $P^g$  must specify, for each  $p \in \text{PID}$ , the value of each global variable of  $P^g$  on process  $p$ , i.e., a  $g$ -valuation  $\alpha: \text{PID} \rightarrow (V_g \rightarrow \mathbb{Z})$ .

Let  $Q$  be the set of processes that have at least one transition in region  $r$  of  $\zeta$ ; the processes of  $Q$  are at a location just before a call to  $g$  in  $s$ . For  $q \in Q$  and formal parameter  $x$ , let  $\alpha(q)(x)$  be the result of evaluating the corresponding actual argument in the call on process  $q$  at  $s$ . For  $v$  a global of  $P$ , let  $\alpha(q)(v)$  be the value of  $v$  in the process state for  $q$  at  $s$ .

At this point, we have defined a partial  $g$ -valuation  $Q \rightarrow (V_g \rightarrow \mathbb{Z})$ , which is the valuation used to check the precondition of  $g$  when the last process in  $Q$  to enter region  $r$  of  $\zeta$  enters. Since  $\zeta$  has no observable collective error, no precondition violation occurs at that step, i.e., the partial  $g$ -valuation can be extended to a full  $g$ -valuation  $\alpha: \text{PID} \rightarrow (V_g \rightarrow \mathbb{Z})$  for which  $\llbracket \text{pre}(g) \rrbracket_{\alpha, p}$  holds for all  $p \in \text{PID}$ . Use valuation  $\alpha$  to specify the value of all global variables in  $s'_0$ ; note that the resulting  $s'_0$  is a  $\text{pre}(g)$ -state.

The processes not in  $Q$  never execute in  $\chi$ , so those process states remain the same at every state in  $\chi$ . In fact the values assigned to the variables of processes not in  $Q$  are irrelevant as they are never read nor modified and therefore have no impact on the extracted execution  $\chi$ . They were only needed to ensure that the initial state of the extracted execution satisfies the definition of a  $\text{pre}(g)$ -state.

We summarize the discussion above in the following:

**Lemma 3.** *Let  $P$  be a collective-recursion-free program,  $f$  a collective procedure of  $P$ , and  $\zeta$  a finite execution of  $f$  with no observable collective error. Let  $r \geq 1$  and  $Q$  the set of processes that have at least one transition in region  $r$  of  $\zeta$ . Assume  $Q \neq \emptyset$ . Let  $g$  be the observable collective procedure associated to region  $r$  of  $\zeta$ . For  $p \in Q$ , let  $\sigma_p \xrightarrow{\tau_R} \sigma'_p$  be the first transition in  $\zeta$  from process  $p$  in region  $r$ . Then there is an execution  $\chi$  of  $P^g$  with initial state  $s'_0$  such that*

1.  $s'_0$  is a  $\text{pre}(g)$ -state,
2. for  $p \in Q$ , the value assigned to global variable  $v$  of  $P$  in process  $p$  of  $s'_0$  is the value assigned to  $v$  by process  $p$  in  $\sigma_p$ ,
3. for  $p \in Q$ , the value assigned to formal parameter  $x$  of  $g$  in process  $p$  of  $s'_0$  is the value assigned to  $x$  by process  $p$  in  $\sigma'_p$ , and
4. the sequence of actions in  $\chi$  is the sequence of region  $r$  actions in  $\zeta$ , in the same order.

We next show that internal errors in  $\zeta$  are preserved in an extracted execution.

**Lemma 4.** *Assume  $\zeta$  is an observably correct execution of collective procedure  $f$  in a collective-recursion-free program  $P$ . Let  $g_1, g_2, \dots$  be the sequence of collective procedures called from  $f$ . If a transition in region  $r$  of  $\zeta$  commits an internal collective error then the execution of  $P^{g_r}$  extracted from region  $r$  of  $\zeta$  is incorrect.*

*Proof.* Write  $\zeta = (s_0 \xrightarrow{t_1} s_1 \xrightarrow{t_2} \dots)$ . Let  $g = g_r$ . Let  $\chi$  be the execution of  $P^g$  extracted from  $\zeta$ . By Lemma 3,  $\chi$  preserves the order of all region  $r$  transitions in  $\zeta$ .

Suppose  $\zeta$  commits a consistency violation at step  $i$ , and  $t_i$  is in region  $r$  of  $\zeta$ . We will see that  $\chi$  commits a consistency violation. Let  $p$  be the process of  $t_i$ .

Suppose  $t_i$  is the first action of  $p$  in region  $r$  of  $\zeta$ , i.e., the  $e^g$  action entering the region. The string  $T_p(\zeta^i)$  has the following form:

$$T_p(\zeta^i) = e^f e^{g_1} \dots x^{g_1} \dots e^{g_{r-1}} \dots x^{g_{r-1}} e^g.$$

For a consistency violation to occur at step  $i$ , there must be some other process  $q$  with

$$T_q(\zeta^i) = e^f e^{g_1} \dots x^{g_1} \dots e^{g_{r-1}} \dots x^{g_{r-1}} e^h \dots,$$

where  $h$  is an observable collective procedure used by  $f$  and  $h \neq g$ . But this is an observable consistency violation, contradicting the assumption that  $\zeta$  is observably correct.

So  $t_i$  must occur after the  $e^g$  transition in  $p$ , and there is some process  $q$ , strings  $\mu$  and  $\nu$ , and distinct collective action symbols  $c$  and  $c'$  such that

$$\begin{aligned} T_p(\zeta) &= \mu e^g \nu c \\ T_q(\zeta) &= \mu e^g \nu c' \dots \end{aligned}$$

Both the  $c$  and  $c'$  transitions occur in region  $r$  of  $\zeta$ . But then

$$\begin{aligned} T_p(\chi) &= e^g \nu c \\ T_q(\chi) &= e^g \nu c' \dots, \end{aligned}$$

and  $\chi$  has a consistency violation, as required.

So assume that  $\zeta$  is collective consistent. It follows that  $\chi$  is collective consistent: each  $T_p(\chi)$  is the substring of  $T_p(\zeta)$  starting with the  $e^g$  and ending with the corresponding  $x^g$  (if the  $x^g$  occurs in  $T_p(\zeta)$ ). Now we consider the other kinds of violations.

If a precondition or postcondition violation occurs at transition  $i$  in  $\zeta$  in region  $r$ , then the same violation occurs at the corresponding transition in  $\chi$ : the process states are exactly the same except for the additional stack entries on the bottom of the stacks in  $\zeta$ , which have no effect on the valuations used to evaluate pre- or postconditions.

An *assigns* or *waitsfor* violation in  $\zeta$  in region  $r$  depends only on entrance and exit actions made by processes in region  $r$ , and the process states immediately preceding or following those actions. Since the order of the region  $r$  actions in  $\zeta$

is exactly the same as the order of the actions in  $\chi$ , these violations carry over to  $\chi$ .

The first kind of segment boundary violation in region  $r$  of  $\zeta$  occurs when a process  $p$  in segment  $a$  of region  $r$  receives a message that was sent by a process  $q$  in segment  $b$ , where  $a < b$ . If the send occurred post-region  $r$ , then the violation is observable, contradicting the assumption that  $\zeta$  is observably correct. So the send must occur in region  $r$ , and cross a segment boundary within region  $r$ , and therefore the exact same violation occurs in  $\chi$ . A similar argument applies to the second kind of segment boundary violation.

If  $t_i$  is an exit action by process  $p$  when there is an unreceived message for  $p$  sent from a previous segment of  $\zeta$ , then the same happens at the corresponding transition in  $\chi$ .

An immediate corollary is:

**Lemma 5.** *Let  $P$  be a collective-recursion-free program and  $f$  a collective procedure of  $P$ . Assume all collective procedures used by  $f$  conform. Let  $\zeta$  be an execution of  $P^f$ . Then  $\zeta$  is correct if and only if  $\zeta$  is observably correct.*

*Proof.* If  $\zeta$  is correct then  $\zeta$  is observably correct by the comments following Definition 6.

So assume  $\zeta$  is observably correct but not correct; we will arrive at a contradiction. It follows from Definition 6 that  $\zeta$  has an internal collective error. Let  $r$  be the region of the transition committing the first internal collective error of  $\zeta$ , and  $g$  the associated collective procedure used by  $f$ . Let  $\chi$  be the execution of  $P^g$  extracted from region  $r$  of  $\zeta$ . By Lemma 3, the initial state of  $\chi$  is a  $pre(g)$ -state. By Lemma 4,  $\chi$  is incorrect. But this contradicts the assumption that  $g$  conforms.

The following is an application of the notion of extracted execution. We use the concept to show that calling a conforming collective procedure cannot lead to error, assuming the call did not commit a precondition violation.

**Lemma 6.** *Suppose  $P$  is a collective-recursion-free program,  $f$  is a collective procedure of  $P$ , and all collective procedures used by  $f$  conform. Let  $\zeta = (s_0 \xrightarrow{t_1} s_1 \xrightarrow{t_2} \dots)$  be an execution of  $P^f$ . Suppose  $1 \leq i < \text{len}(\zeta)$ ,  $\zeta^i$  is observably correct,  $p$  is the process of  $t_{i+1}$ , and  $p$  is at an internal state at  $s_i$ . Then  $\zeta$  does not commit a collective error at step  $i + 1$ .*

*Proof.* Say  $t_{i+1}$  is in region  $r$  of  $\zeta$ , inside a call to the collective procedure  $g$  used by  $f$ . Let  $\chi$  be the execution of  $P^g$  extracted from region  $r$  of  $\zeta^i$ . Let  $s'_0$  be the initial state of  $\chi$ . By Lemma 3,  $s'_0$  is a  $pre(g)$ -state.

*Case 1:* action  $t = t_{i+1}$  is enabled at the final state  $s'_f$  of  $\chi$ . Let  $\chi'$  be the execution of  $\chi$  extended by the single action  $t$ . Since  $g$  conforms and  $\chi'$  is an execution of  $P^g$  from a  $pre(g)$ -state,  $\chi'$  is correct.

If  $t$  commits a consistency, precondition, postcondition, assigns, *waitfor*, or unreceived message violation in  $\zeta$ , the same violation is committed by  $t$  in  $\chi'$ ,

contradicting the fact that  $\chi'$  is correct. Suppose  $t$  commits a segment boundary violation in  $\zeta$ . If the violation is of the second kind, i.e.,  $p$  sends a message to a process  $q$  in a later segment of  $\zeta$ , then it is also a segment violation of the second kind in  $\chi'$ : even if, in  $\chi'$ ,  $q$  has terminated,  $q$  is at least one segment after  $p$  when  $t$  executes. Again, this contradicts the fact that  $\chi'$  is correct.

If the segment boundary violation is of the first kind then  $t$  is a receive by  $p$  of a message from  $q$  when  $q$  is at a higher segment in  $\zeta$ . As we are assuming  $t$  is enabled at  $s'_f$ , the send must also occur in  $\chi'$ , and  $t$  commits a segment boundary violation in  $\chi'$ .

*Case 2:  $t$  is not enabled at  $s'_f$ .* Then  $t$  is a receive statement in  $p$  with source  $q$ , and the channel  $p \rightarrow q$  is empty in  $s'_f$ . Since  $t$  is enabled at the final state of  $\zeta^i$ , the matching send in  $q$  must occur in  $\zeta^i$ , and since that transition does not occur in  $\chi$ , it must occur in  $\zeta$  after  $q$  exits region  $r$ . Hence at  $s'_f$ ,  $q$  has terminated. But this means that no extension of  $\chi'$  can terminate:  $p$  is permanently blocked at the receive with empty channel. Hence any maximal extension of  $\chi'$  is incorrect: either infinite, or finite ending in deadlock. This contradicts the assumption that  $g$  conforms.

## A.5 Existence of simulations

Assume  $P$  is a collective-recursion-free program with collective procedure  $f$ , and  $\zeta$  is an execution of  $P^f$  with no observable collective error. Then one can form a simulation  $\eta$  (i.e., an execution of  $\overline{P^f}$ ) from  $\zeta$  by essentially deleting all transitions within regions and replacing each with a single *havoc* followed by *wait* statement. We now make this precise.

Say  $\zeta = (s_0 \xrightarrow{t_1} s_1 \xrightarrow{t_2} \dots)$ . For each  $i \geq 0$ , we will define a finite simulation  $\eta_i$ , such that  $\eta_i$  is a prefix of or equal to  $\eta_{i+1}$ . We will let  $\eta$  be the simulation which is the limit of the  $\eta_i$ . We will also establish the following invariant: for all  $i \geq 0$ , the final state  $s'_i$  of  $\eta_i$  is related to  $s_i$  as follows: for any  $p, q \in \text{PID}$ , the state of the  $p \rightarrow q$  channel in  $s'_i$  is obtained by deleting from the corresponding channel in  $s_i$  any messages sent from a process at an internal state. Moreover, if  $p$  is at an external state at step  $i$  of  $\zeta$ , the the process of state of  $p$  at  $s_i$  is identical to the process state of  $p$  at  $s'_i$ .

For  $i = 0$ ,  $\eta_0$  is the execution of length 0 starting at the initial state of  $\overline{P^f}$  corresponding to  $s_0$ .

Assume  $i \geq 0$  and  $\eta_i$  has been defined. Let  $p$  be the process of  $t_{i+1}$ . If  $p$  is at an internal state at both  $s_i$  and  $s_{i+1}$  (i.e.,  $t_{i+1}$  is strictly within some region), let  $\eta_{i+1} = \eta_i$ . The invariant on channels is maintained because (i) if  $t_{i+1}$  is a send, the message sent is from an internal region, and (ii) if  $t_{i+1}$  is a receive then the message received was sent within the region, as  $\zeta$  is free of collective errors.

If  $p$  is at an external state at  $s_i$ ,  $\eta_{i+1}$  is obtained by appending  $t_{i+1}$  and the resulting state to  $\eta_i$  (this includes the case where  $t_{i+1}$  enters some region). If  $p$  is at an internal state at  $s_i$  and at an external state at  $s_{i+1}$ , i.e.,  $t_{i+1}$  exits a region, proceed as follows: let  $g$  be the collective procedure being exited. Then  $\eta_{i+1}$  is obtained from  $\eta_i$  by adding the following three transitions to  $\eta_i$ :



1. A *havoc* transition which assigns to the variables in  $g$ 's *assigns* list the values they have in  $s_i$ . The fact that  $\zeta$  has no collective error guarantees that these values will not cause a postcondition violation and are therefore allowed in a simulation.
2. A *wait* transition on  $g$ 's wait set. The latter must be enabled in  $\eta$ , otherwise  $\zeta$  would have an observable *waitsfor* violation.
3. The  $x^g$  transition.

After the exit, the process state for the executing process  $p$  in  $s'_{i+1}$  is identical to the process state of  $p$  in  $s_{i+1}$ , maintaining the claimed invariant.

Recall that infinite executions are considered erroneous. Here we show that certain infinite executions will be detected by simulations.

**Lemma 7.** *Suppose  $P$  is a collective-recursion-free program,  $f$  is a collective procedure of  $P$ , and all collective procedures used by  $f$  conform. If  $\zeta$  is an infinite execution of  $P^f$  with no observable collective error then there exists an incorrect simulation of  $f$ .*

*Proof.* Suppose  $\zeta$  has only a finite number of regions. Since  $\zeta$  is infinite, there is some region  $r$  with an infinite number of transitions from  $\zeta$ . Let  $g$  be the collective procedure associated to region  $r$ . Let  $m$  be the index of the last transition in  $\zeta$  to enter region  $r$ . For each  $i \geq m$ , let  $\chi_i$  be the execution of  $P^g$  extracted from region  $r$  of  $\zeta^i$ . By Lemma 3, each  $\chi_i$  starts from the same initial state. Moreover the action sequence of  $\chi_{i+1}$  is either the same sequence as that of  $\chi_i$ , or extends that sequence by one action. It follows that for each  $i$ ,  $\chi_i$  is a (not necessarily strict) prefix of  $\chi_{i+1}$ , and, for an infinite number of  $i$ ,  $\chi_i$  is a strict prefix of  $\chi_{i+1}$ . Let  $\chi$  be the limit of the  $\chi_i$ . Then  $\chi$  is an infinite execution of  $P^g$  from a  $\text{pre}(g)$ -state, contradicting the assumption that  $g$  conforms.

So the regions in  $\zeta$  increase without bound. That means there are an infinite number of transitions from external states, so the simulation formed from  $\zeta$  is infinite.

The following shows that collective errors and deadlocks are detected by simulations:

**Lemma 8.** *Suppose  $P$  is a collective-recursion-free program,  $f$  is a collective procedure of  $P$ , all collective procedures used by  $f$  conform, and  $\zeta$  is an execution of  $P^f$ . If  $\zeta$  has an observable collective error or ends in deadlock then there exists an incorrect simulation of  $f$ .*

*Proof.* Let  $i \in 0..\text{len}(\zeta)$  be the maximal integer for which  $\zeta^i$  has no observable collective error. The final state of  $\zeta^i$  is  $s_i$ . Let  $\eta$  be the simulation generated from  $\zeta^i$ . Let  $s'$  be the final state of  $\eta$ .

*Case 1:  $s_i$  is not deadlocked.* Then  $t = t_{i+1}$  commits an observable collective error in  $\zeta$ . Let  $p$  be the process of  $t_{i+1}$ . By Lemma 6,  $p$  is at an external state at  $s_i$ . Hence  $t$  is also an enabled transition in the simulation state  $s'$ .

If  $t$  commits an observable consistency violation, then it must be an  $e^g$  action for some procedure  $g$  used by  $f$ . (It cannot be an  $x^g$  action, since those are executed from an internal state.) Then action  $t$  causes the same collective error from  $s'$ .

Likewise, if  $t$  commits an observable postcondition, assigns, *waitsfor*, or unreceived message violation, then  $t$  must be an  $x^f$  action, and the same action is enabled at  $s'$ , and commits the same collective error.

Likewise, if  $t$  commits an observable precondition violation in a call to procedure  $g$ , then it commits the same violation at simulation state  $s'$ .

Suppose  $t$  commits a segment boundary violation. Then  $t$  is a send or receive by  $p$  when the matching process  $q$  is in a later segment in  $\zeta$ . Since  $p$  is at an external state at  $s_i$ ,  $q$  is also at a later segment in  $\eta$  at  $s'$ . Hence  $t$  commits a segment boundary violation from  $s'$ .

*Case 2:  $s_i$  is deadlocked.* Then  $\zeta = \zeta^i$  has no observable collective error, there is some process that is not terminated at  $s_i$ , and all non-terminated processes are blocked at receive statements with empty channels at  $s_i$ . Suppose all the non-terminated processes are at external states at  $s_i$ . Then the simulation state  $s'$  is also deadlocked, and therefore incorrect.

So suppose there is some process  $p$  at an internal state at  $s_i$ . Say  $p$  is in a region corresponding to a call to collective procedure  $g$  used by  $f$ . Consider the execution  $\chi$  extracted from this region of  $\zeta$ . By Lemma 3,  $\chi$  starts at a  $pre(g)$ -state. If there were an action enabled at the final state of  $\chi$ , then that same action would be enabled at  $s_i$ ; hence  $\chi$  is also deadlocked. But then  $\chi$  is an incorrect execution of conforming procedure  $g$  from a  $pre(g)$ -state, a contradiction.

Finally, we can prove the main theorem:

*Proof (Proof of Theorem 1).* Assume  $f$  is a collective procedure in  $P$  and all collective procedures used by  $f$  conform. Suppose that  $f$  does not conform; we must show that there is an incorrect simulation of  $f$ . As  $f$  does not conform, there is an incorrect execution  $\zeta$  of  $P^f$  from a  $pre(f)$ -state. By Lemma 5,  $\zeta$  is not observably correct. If  $\zeta$  is finite, or is infinite and commits an observable collective error, Lemma 8 implies an incorrect simulation exists. If  $\zeta$  is infinite and does not commit an observable collective error, Lemma 7 implies an incorrect simulation exists.

## B Contracts of MPI Collective Functions and MPI\_Reduce\_local

### B.1 MPI\_Allgather

```

#define SBUF          \mpi_buf(sbuf, scount, stype)
#define SBUF_OF(id)   \mpi_on(SBUF, (id))
#define RBUF          \mpi_buf(rbuf, rcount * \mpi_comm_size, rtype)
#define RBUF_OF(id)   (\mpi_buf(rbuf, rcount, rtype) + (id)*rcount)
#define SSIG          (\mpi_sig(stype) * scount)
#define SSIG_OF(id)   \mpi_on(SSIG, (id))
#define RSIG          (\mpi_sig(rtype) * rcount)
/*@ mpi uses comm;
    mpi collective(comm):
        requires \mpi_nonoverlapping(rtype) && rcount >= 0 &&
            \valid(RBUF);
        requires \separated(RBUF, {SBUF|int i; sbuf != MPI_IN_PLACE});
        assigns *RBUF;
        waitsfor {i | int i; 0 <= i < \mpi_comm_size};
        behavior not_in_place:
            assumes sbuf != MPI_IN_PLACE;
            requires \mpi_agree(sbuf != MPI_IN_PLACE);
            requires \forall int i; 0 <= i < \mpi_comm_size
                ==> SSIG == \mpi_on(RSIG, i);
            requires scount >= 0 && \valid_read(SBUF);
            ensures \forall int i; 0 <= i < \mpi_comm_size
                ==> *RBUF_OF(i) == \mpi_on(*SBUF, i);
        behavior in_place:
            assumes sbuf == MPI_IN_PLACE;
            requires \mpi_agree(sbuf == MPI_IN_PLACE);
            requires \forall int i; 0 <= i < \mpi_comm_size
                ==> RSIG == \mpi_on(RSIG, i);
            ensures \forall int i; 0 <= i < \mpi_comm_size
                ==> *RBUF_OF(i) == \old(\mpi_on(*RBUF_OF(i), i));
        disjoint behaviors;
        complete behaviors;
    */
int MPI_Allgather(const void *sbuf, int scount, MPI_Datatype stype,
                  void *rbuf, int rcount, MPI_Datatype rtype,
                  MPI_Comm comm);

```

### B.2 MPI\_Allgatherv

```

#define SBUF          \mpi_buf(sbuf, scount, stype)
#define RBUF_OF(i)    (\mpi_buf(rbuf, rcounts[(i)], rtype)+displs[(i)])
#define SSIG_OF(i)    \mpi_on(\mpi_sig(stype) * scount, (i))
#define RSIG_OF(i)    \mpi_sig(rtype) * rcounts[(i)]
/*@ mpi uses comm;
    mpi collective(comm):

```

```

requires \valid_read(rcounts + (0 .. \mpi_comm_size-1));
requires \valid_read(displs + (0 .. \mpi_comm_size-1));
requires \forall int i; 0 <= i < \mpi_comm_size
    ==> rcounts[i] >= 0;
requires \forall int i; 0 <= i < \mpi_comm_size
    ==> \valid(RBUF_OF(i));
requires \mpi_nonoverlapping(rtype);
requires \forall int i, j; 0 <= i < j < \mpi_comm_size
    ==> (displs[i] + rcounts[i] <= displs[j]) ||
        (displs[j] + rcounts[j] <= displs[i]);
requires \separated({RBUF_OF(i) | int i; 0 <= i < \mpi_comm_size},
    {{SBUF | int i; sbuf != MPI_IN_PLACE},
     rcounts + (0 .. \mpi_comm_size-1),
     displs + (0 .. \mpi_comm_size-1)});
assigns {*RBUF_OF(i) | int i; 0 <= i < \mpi_comm_size};
waitsfor {i | int i; 0 <= i < \mpi_comm_size};
behavior not_in_place:
    assumes sbuf != MPI_IN_PLACE;
    requires \mpi_agree(sbuf != MPI_IN_PLACE);
    requires \valid_read(SBUF) && scount >= 0;
    requires \forall int i; 0 <= i < \mpi_comm_size
        ==> RSIG_OF(i) == SSIG_OF(i);
    ensures \forall int i; 0 <= i < \mpi_comm_size
        ==> *RBUF_OF(i) == \mpi_on(*SBUF, i);
behavior in_place:
    assumes sbuf == MPI_IN_PLACE;
    requires \mpi_agree(sbuf == MPI_IN_PLACE);
    requires \forall int i; 0 <= i < \mpi_comm_size
        ==> \mpi_agree(RSIG_OF(i));
    ensures \forall int i; 0 <= i < \mpi_comm_size
        ==> *RBUF_OF(i) == \old(\mpi_on(*RBUF_OF(i), i));
disjoint behaviors;
complete behaviors;
*/
int MPI_Allgatherv(const void *sbuf, int scount, MPI_Datatype stype,
    void *rbuf, const int *rcounts, const int *displs,
    MPI_Datatype rtype, MPI_Comm comm);

```

### B.3 MPI\_Allreduce

```

#define SBUF \mpi_buf(sbuf, count, datatype)
#define RBUF \mpi_buf(rbuf, count, datatype)
#define AGREE(x) \mpi_agree((x))
/*@
    mpi uses comm;
    mpi collective(comm):
        requires \valid(RBUF) && count >= 0 && AGREE(count) &&
            AGREE(datatype) && AGREE(op);
        requires \mpi_nonoverlapping(datatype);
        requires \separated(RBUF, {SBUF | int i; sbuf != MPI_IN_PLACE});

```

```

    assigns *RBUF;
    ensures AGREE(*RBUF);
    waitsfor {i | int i; 0 <= i < \mpi_comm_size};
    behavior not_in_place:
        assumes sbuf != MPI_IN_PLACE;
        requires AGREE(sbuf != MPI_IN_PLACE);
        requires \valid_read(SBUF);
        ensures \mpi_reduce(*RBUF, 0, \mpi_comm_size, op,
                           \lambda integer t; \mpi_on(*SBUF, t));
    behavior in_place:
        assumes sbuf == MPI_IN_PLACE;
        requires AGREE(sbuf == MPI_IN_PLACE);
        ensures \mpi_reduce(*RBUF, 0, \mpi_comm_size, op,
                           \lambda integer t; \mpi_on(\old(*RBUF), t));
    disjoint behaviors;
    complete behaviors;
*/
int MPI_Allreduce(const void *sbuf, void *rbuf, int count,
                  MPI_Datatype datatype, MPI_Op op, MPI_Comm comm);

```

#### B.4 MPI\_Alltoall

```

#define SBUF      \mpi_buf(sbuf, scount * \mpi_comm_size, stype)
#define RBUF      \mpi_buf(rbuf, rcount * \mpi_comm_size, rtype)
#define SBUF_OF(i) (\mpi_buf(sbuf, scount, stype) + (i) * scount)
#define RBUF_OF(i) (\mpi_buf(rbuf, rcount, rtype) + (i) * rcount)
/*@ mpi uses comm;
    mpi collective(comm):
        requires \valid(RBUF) && rcount >= 0 && \mpi_nonoverlapping(rtype);
        requires \separated(RBUF, { SBUF | int i; sbuf != MPI_IN_PLACE });
        assigns *RBUF;
        waitsfor {i | int i; 0 <= i < \mpi_comm_size};
        behavior not_in_place:
            assumes sbuf != MPI_IN_PLACE;
            requires \mpi_agree(sbuf != MPI_IN_PLACE) && \valid_read(SBUF);
            requires scount >= 0 &&
                scount * \mpi_sig(stype) == rcount * \mpi_sig(rtype);
            ensures \forall int i, j;
                0 <= i < \mpi_comm_size && j == \mpi_comm_rank
                    ==> *RBUF_OF(i) == \mpi_on(*SBUF_OF(j), i);
        behavior in_place:
            assumes sbuf == MPI_IN_PLACE;
            requires \mpi_agree(sbuf == MPI_IN_PLACE);
            ensures \forall int i, j;
                0 <= i < \mpi_comm_size && j == \mpi_comm_rank
                    ==> *RBUF_OF(i) == \old(\mpi_on(*RBUF_OF(j), i));
        disjoint behaviors;
        complete behaviors;
*/
int MPI_Alltoall(const void *sbuf, int scount, MPI_Datatype stype,

```

```
void *rbuf, int rcount, MPI_Datatype rtype,
MPI_Comm comm);
```

### B.5 MPI\_Alltoallv

```
#define SBUF_OF(i) (\mpi_buf(sbuf, scounts[(i)], stype) + sdispls[(i)])
#define RBUF_OF(i) (\mpi_buf(rbuf, rcounts[(i)], rtype) + rdispls[(i)])
/*@ mpi uses comm;
mpi collective(comm):
  requires \valid_read(rcounts + (0 .. \mpi_comm_size-1));
  requires \valid_read(rdispls + (0 .. \mpi_comm_size-1));
  requires \forall int i; 0 <= i < \mpi_comm_size
    ==> \valid(RBUF_OF(i));
  requires \forall int i; 0 <= i < \mpi_comm_size
    ==> rcounts[i] >= 0;
  requires \forall int i, j; 0 <= i < j < \mpi_comm_size ==>
    (rdispls[i] + rcounts[i] <= rdispls[j] ||
     (rdispls[j] + rcounts[j] <= rdispls[i]));
  requires \mpi_nonoverlapping(rtype);
  requires \separated({RBUF_OF(i) | int i; 0 <= i < \mpi_comm_size},
    { {SBUF_OF(i), scounts + i, sdispls + i} | int i;
      0 <= i < \mpi_comm_size && sbuf != MPI_IN_PLACE},
      rcounts + (0 .. \mpi_comm_size-1),
      rdispls + (0 .. \mpi_comm_size-1) });
  assigns {*RBUF_OF(i) | int i; 0 <= i < \mpi_comm_size-1};
  waitsfor {i | int i; 0 <= i < \mpi_comm_size};
  behavior not_in_place:
    assumes sbuf != MPI_IN_PLACE;
    requires \mpi_agree(sbuf != MPI_IN_PLACE);
    requires \forall int i; 0 <= i < \mpi_comm_size
      ==> \valid_read(SBUF_OF(i));
    requires \valid_read(scounts + (0 .. \mpi_comm_size-1));
    requires \valid_read(sdispls + (0 .. \mpi_comm_size-1));
    requires \forall int i; 0 <= i < \mpi_comm_size
      ==> scounts[i] >= 0;
    requires \forall int i, j; 0 <= i < \mpi_comm_size &&
      j == \mpi_comm_rank ==>
      \mpi_on(scounts[j] * \mpi_sig(stype), i) ==
      rcounts[i] * \mpi_sig(rtype);
    ensures \forall int i, j; 0 <= i < \mpi_comm_size &&
      j == \mpi_comm_rank ==>
      *RBUF_OF(i) == \mpi_on(*SBUF_OF(j), i);
  behavior in_place:
    assumes sbuf == MPI_IN_PLACE;
    requires \mpi_agree(sbuf == MPI_IN_PLACE);
    requires \forall int i, j; 0 <= i < \mpi_comm_size &&
      j == \mpi_comm_rank ==>
      \mpi_on(rcounts[j] * \mpi_sig(rtype), i) ==
      rcounts[i] * \mpi_sig(rtype);
    ensures \forall int i, j; 0 <= i < \mpi_comm_size &&
```

```

        j == \mpi_comm_rank ==>
            *RBUF_OF(i) == \old(\mpi_on(*RBUF_OF(j), i));
    disjoint behaviors;
    complete behaviors;
*/
int MPI_Alltoallv(const void *sbuf, const int *scounts,
                  const int *sdispls, MPI_Datatype stype, void *rbuf,
                  const int *rcounts, const int *rdispls,
                  MPI_Datatype rtype, MPI_Comm comm);

```

## B.6 MPI\_Alltoallw

```

#define SBUF_OF(i) (\mpi_buf(sbuf, scounts[(i)], stypes[(i)]+sdispls[i])
#define RBUF_OF(i) (\mpi_buf(rbuf, rcounts[(i)], rtypes[(i)]+rdispls[i])
/*@ mpi uses comm;
    mpi collective(comm):
        requires \valid_read(rcounts + (0 .. \mpi_comm_size-1));
        requires \valid_read(rdispls + (0 .. \mpi_comm_size-1));
        requires \valid_read(rtypes + (0 .. \mpi_comm_size-1));
        requires \forall int i; 0 <= i < \mpi_comm_size
            ==> \valid(RBUF_OF(i)) && rcounts[i] >= 0;
        requires \forall int i; 0 <= i < \mpi_comm_size
            ==> \mpi_nonoverlapping(rtypes[i]);
        requires \forall int i, j; 0 <= i < j < \mpi_comm_size
            ==> (rdispls[i] + rcounts[i] <= rdispls[j]) ||
                (rdispls[j] + rcounts[j] <= rdispls[i]);
        requires \separated({RBUF_OF(i) | int i; 0 <= i < \mpi_comm_size},
            { {SBUF_OF(i), scounts + i, sdispls + i, stypes + i} |
              int i; 0 <= i < \mpi_comm_size},
              rcounts + (0 .. \mpi_comm_size-1),
              rdispls + (0 .. \mpi_comm_size-1),
              rtypes + (0 .. \mpi_comm_size-1) });
        assigns { *RBUF_OF(i) | int i; 0 <= i < \mpi_comm_size};
        waitsfor {i | int i; 0 <= i < \mpi_comm_size};
        behavior not_in_place:
            assumes sbuf != MPI_IN_PLACE;
            requires \mpi_agree(sbuf != MPI_IN_PLACE);
            requires \valid_read(scount + (0 .. \mpi_comm_size-1));
            requires \valid_read(sdispls + (0 .. \mpi_comm_size-1));
            requires \valid_read(stypes + (0 .. \mpi_comm_size-1));
            requires \forall int i; 0 <= i < \mpi_comm_size ==>
                \valid_read(SBUF_OF(i)) && scounts[i] >= 0;
            requires \forall int i, j; 0 <= i < \mpi_comm_size &&
                j == \mpi_comm_rank ==>
                    \mpi_on(\mpi_sig(stypes[j]) * scounts[j], i) ==
                    \mpi_sig(rtypes[i]) * rcounts[i];
            ensures \forall int i, j; 0 <= i < \mpi_comm_size &&
                j == \mpi_comm_rank ==>
                    *RBUF_OF(i) == \mpi_on(*SBUF_OF(j), i);
        behavior in_place:

```

```

    assumes sbuf == MPI_IN_PLACE;
    requires \mpi_agree(sbuf == MPI_IN_PLACE);
    requires \forall int i, j; 0 <= i < \mpi_comm_size &&
        j == \mpi_comm_rank ==>
        \mpi_on(\mpi_sig(rtypes[j]) * rcounts[j], i) ==
        \mpi_sig(rtypes[i]) * rcounts[i];
    ensures \forall int i, j; 0 <= i < \mpi_comm_size &&
        j == \mpi_comm_rank ==>
        *RBUF_OF(i) == \old(\mpi_on(*RBUF_OF(j), i));
    disjoint behaviors;
    complete behaviors;
*/
int MPI_Alltoallw(const void *sbuf, const int scounts[],
    const int sdispls[], const MPI_Datatype stypes[],
    void *rbuf, const int rcounts[], const int rdispls[],
    const MPI_Datatype rtypes[], MPI_Comm comm);

```

## B.7 MPI\_Barrier

```

/*@ mpi uses comm;
    mpi collective(comm):
        requires \true;
        assigns \nothing;
        waitsfor {i | int i; 0 <= i < \mpi_comm_size};
*/
int MPI_Barrier(MPI_Comm comm);

```

## B.8 MPI\_Bcast

```

#define BUF \mpi_buf(buf, count, datatype)
#define AGREE(x) \mpi_agree((x))
/*@ mpi uses comm;
    mpi collective(comm):
        requires 0 <= root < \mpi_comm_size && AGREE(root);
        requires AGREE(count * \mpi_sig(datatype)) && 0 <= count;
        requires \mpi_nonoverlapping(datatype);
        ensures AGREE(*BUF);
        behavior root:
            assumes \mpi_comm_rank == root;
            requires \valid_read(BUF);
            assigns \nothing;
        behavior nonroot:
            assumes \mpi_comm_rank != root;
            requires \valid(BUF);
            assigns *BUF;
            waitsfor root;
        complete behaviors;
        disjoint behaviors;
*/

```



```
int MPI_Bcast(void * buf, int count, MPI_Datatype datatype,
             int root, MPI_Comm comm);
```

### B.9 MPI\_Exscan

```
#define SBUF \mpi_buf(sbuf, count, datatype)
#define RBUF \mpi_buf(rbuf, count, datatype)
#define AGREE(x) \mpi_agree((x))
/*@ mpi uses comm;
    mpi collective(comm):
        requires count >= 0 && AGREE(count) && AGREE(datatype) &&
            AGREE(op);
        requires \mpi_nonoverlapping(datatype);
        requires sbuf != MPI_IN_PLACE ==> \valid_read(SBUF);
        requires sbuf == MPI_IN_PLACE || \mpi_comm_rank > 0
            ==> \valid(RBUF);
        waitsfor {i | int i; 0 <= i < \mpi_comm_rank-1};
        behavior zero:
            assumes \mpi_comm_rank == 0;
            assigns \nothing;
        behavior others:
            assumes \mpi_comm_rank > 0;
            requires \separated(RBUF, {SBUF|int i; sbuf!=MPI_IN_PLACE});
            assigns *RBUF;
            ensures \mpi_reduce(*RBUF, 0, \mpi_comm_rank, op,
                \lambda integer t;
                \mpi_on(sbuf != MPI_IN_PLACE?*SBUF:\old(*RBUF),t));
        disjoint behaviors;
        complete behaviors;
    */
int MPI_Exscan(const void *sbuf, void *rbuf, int count,
              MPI_Datatype datatype, MPI_Op op, MPI_Comm comm);
```

### B.10 MPI\_Gather

```
#define SBUF \mpi_buf(sbuf, scount, stype)
#define SBUF_OF(id) \mpi_on(SBUF, (id))
#define RBUF \mpi_buf(rbuf, rcount * \mpi_comm_size, rtype)
#define RBUF_OF(id) (\mpi_buf(rbuf, rcount, rtype) + (id)*rcount)
#define SSIG (\mpi_sig(stype) * scount)
#define SSIG_OF(id) \mpi_on(SSIG, (id))
#define RSIG (\mpi_sig(rtype) * rcount)
#define AGREE(x) \mpi_agree((x))
/*@
    mpi uses comm;
    mpi collective(comm) :
        requires AGREE(root) && 0 <= root < \mpi_comm_size;
        behavior root:
            assumes \mpi_comm_rank == root;
```

```

requires \valid(RBUF) && \mpi_nonoverlapping(rtype);
requires \forall int id; 0 <= id < \mpi_comm_size
    ==> id != root ==> SSIG_OF(id) == RSIG;
requires sbuf == MPI_IN_PLACE || \valid_read(SBUF);
requires sbuf != MPI_IN_PLACE ==> (SSIG == RSIG);
requires \separated(RBUF, {SBUF|int i; sbuf!=MPI_IN_PLACE});
assigns *RBUF;
ensures sbuf != MPI_IN_PLACE ==> *RBUF_OF(root) == *SBUF;
ensures sbuf == MPI_IN_PLACE ==>
    *RBUF_OF(root) == \old(*RBUF_OF(root));
ensures \forall int id; 0 <= id < \mpi_comm_size ==>
    id != root ==>
        *RBUF_OF(id) == *SBUF_OF(id);
waitfor {i | int i; 0 <= i < \mpi_comm_size};
behavior not_root:
    assumes \mpi_comm_rank != root;
    requires \valid_read(SBUF);
    assigns \nothing;
    waitfor \nothing;
disjoint behaviors;
complete behaviors;
*/
int MPI_Gather(const void* sbuf, int scount, MPI_Datatype stype,
               void* rbuf, int rcount, MPI_Datatype rtype,
               int root, MPI_Comm comm);

```

### B.11 MPI\_Gatherv

```

#define SBUF \mpi_buf(sbuf, scount, stype)
#define RBUF(i) (\mpi_buf(rbuf, rcounts[i], rtype) + displs[i])
#define AGREE(x) \mpi_agree((x))
/*@ mpi uses comm;
    mpi collective(comm):
        requires AGREE(root) && 0 <= root < \mpi_comm_size;
        behavior root:
            assumes \mpi_comm_rank == root;
            requires \mpi_nonoverlapping(rtype);
            requires \valid_read(rcounts + (0 .. \mpi_comm_size-1));
            requires \valid_read(displs + (0 .. \mpi_comm_size-1));
            requires \forall int i; 0 <= i < \mpi_comm_size
                ==> rcounts[i] >= 0;
            requires \forall int i; 0 <= i < \mpi_comm_size
                ==> \valid(RBUF(i));
            requires \forall int i,j; 0 <= i < j < \mpi_comm_size
                ==> (displs[i] + rcounts[i] <= displs[j] ||
                    displs[j] + rcounts[j] <= displs[i]);
            requires \forall int i;
                0 <= i < \mpi_comm_size && i != root
                ==> \mpi_on(scount*\mpi_sig(stype), i) ==

```

```

        rcounts[i]*\mpi_sig(rtype);
requires sbuf != MPI_IN_PLACE ==> \valid_read(SBUF) &&
        scount*\mpi_sig(stype) ==
        rcounts[root]*\mpi_sig(rtype);
requires \separated( {RBUF(i)|int i; 0 <= i < \mpi_comm_size},
        { {SBUF, &scount, &stype} |int i; sbuf!=MPI_IN_PLACE},
        rcounts + (0 .. \mpi_comm_size-1),
        displs + (0 .. \mpi_comm_size-1)} );
assigns { *RBUF(i) | int i; 0 <= i < \mpi_comm_size };
ensures sbuf!=MPI_IN_PLACE ? *RBUF(root) == *SBUF :
        *RBUF(root) == \old(*RBUF(root));
ensures \forall int i; 0 <= i < \mpi_comm_size && i != root
        ==> *RBUF(i) == \mpi_on(*SBUF, i);
waitsfor { i | int i; 0 <= i < \mpi_comm_size};
behavior non_root:
    assumes \mpi_comm_rank != root;
    requires \valid_read(SBUF);
    assigns \nothing;
    waitsfor \nothing;
complete behaviors;
disjoint behaviors;*/
int MPI_Gatherv(const void* sbuf, int scount, MPI_Datatype stype,
        void* rbuf, const int rcounts[], const int displs[],
        MPI_Datatype rtype, int root, MPI_Comm comm);

```

## B.12 MPI\_Reduce

```

#define SBUF \mpi_buf(sbuf, count, datatype)
#define RBUF \mpi_buf(rbuf, count, datatype)
#define AGREE(x) \mpi_agree((x))
/*@ mpi uses comm;
mpi collective(comm):
    requires AGREE(root) && 0 <= root < \mpi_comm_size;
    requires 0 <= count && \mpi_nonoverlapping(datatype);
    requires AGREE(count) && AGREE(datatype) && AGREE(op);
    behavior root:
        assumes \mpi_comm_rank == root;
        requires \valid(RBUF);
        requires sbuf == MPI_IN_PLACE || \valid_read(SBUF);
        requires sbuf != MPI_IN_PLACE ==> \separated(RBUF, SBUF);
        assigns *RBUF;
        ensures \mpi_reduce(*RBUF, 0, \mpi_comm_size, op,
            \lambda integer t; \mpi_on(
                sbuf != MPI_IN_PLACE ? *SBUF : \old(*RBUF),
                t));
        waitsfor {i | int i; 0 <= i < \mpi_comm_size};
    behavior non_root:
        assumes \mpi_comm_rank != root;
        requires \valid_read(SBUF);
        assigns \nothing;

```

```

        waitsfor \nothing;
        disjoint behaviors;
        complete behaviors;
    */
int MPI_Reduce(const void *sbuf, void *rbuf, int count,
               MPI_Datatype datatype, MPI_Op op, int root,
               MPI_Comm comm);

```

### B.13 MPI\_Reduce\_local

```

#define INBUF      \mpi_buf(inbuf, count, datatype)
#define INOUTBUF   \mpi_buf(inoutbuf, count, datatype)
/*@ requires \valid(INBUF) && \valid(INOUTBUF);
    requires count >= 0;
    requires \mpi_nonoverlapping(datatype);
    requires \separated(INOUTBUF, INBUF);
    assigns *INOUTBUF;
    ensures \mpi_reduce(*INOUTBUF, 0, 2, op, \lambda integer t;
                        t == 0 ? *INBUF : \old(*INOUTBUF));
*/
int MPI_Reduce_local(const void *inbuf, void *inoutbuf, int count,
                    MPI_Datatype datatype, MPI_Op op);

```

### B.14 MPI\_Reduce\_scatter

```

#define MY_RANK      \mpi_comm_rank
#define SCOUNT        \sum(0, \mpi_comm_size-1, \lambda int k0; rcounts[k0])
#define COUNTS(i)    (i==0?0:\sum(0, (i-1), \lambda int k1; rcounts[k1]))
#define SBUF          \mpi_buf(sbuf, SCOUNT, datatype)
#define MY_SBUF_BLK(i) \
    (\mpi_buf(sbuf, rcounts[(i)], datatype)+COUNTS(i))
#define MY_RBUF_AS_SBUF_BLK(i) \
    (\mpi_buf(rbuf, rcounts[(i)], datatype)+COUNTS(i))
#define RBUF_AS_SBUF \
    \mpi_buf(rbuf, SCOUNT, datatype)
#define MY_RBUF      \mpi_buf(rbuf, rcounts[MY_RANK], datatype)
#define AGREE(x)      \mpi_agree((x))
/*@ mpi uses comm;
    mpi collective(comm):
        requires \valid_read(rcounts + (0 .. \mpi_comm_size-1));
        requires rcounts[MY_RANK] >= 0 && AGREE(datatype) &&
            AGREE(SCOUNT) && AGREE(op);
        requires \forall int i; 0 <= i < \mpi_comm_size
            ==> AGREE(rcounts[i]);
        requires \mpi_nonoverlapping(datatype);
        requires \separated(MY_RBUF,
            {rcounts + (0 .. \mpi_comm_size-1),
             {SBUF | int i; sbuf != MPI_IN_PLACE }});
        assigns *MY_RBUF;

```

```

ensures \forall int i; i == \mpi_comm_rank ==>
    \mpi_reduce(*MY_RBUF, 0, \mpi_comm_size, op,
    \lambda integer t; \mpi_on(
        sbuf != MPI_IN_PLACE ? *MY_SBUF_BLK(i) :
            \old(*MY_RBUF_AS_SBUF_BLK(i)),
        t));
waitsfor {i | int i; 0 <= i < \mpi_comm_size};
behavior not_in_place:
    assumes sbuf != MPI_IN_PLACE;
    requires \valid_read(SBUF);
    requires \valid(MY_RBUF);
behavior in_place:
    assumes sbuf == MPI_IN_PLACE;
    requires \valid(RBUF_AS_SBUF);
disjoint behaviors;
complete behaviors;
*/
int MPI_Reduce_scatter(const void *sbuf, void *rbuf,
    const int *rcounts, MPI_Datatype datatype,
    MPI_Op op, MPI_Comm comm);

```

### B.15 MPI\_Reduce\_scatter\_block

```

#define SBUF \mpi_buf(sbuf, rcount * \mpi_comm_size, datatype)
#define MY_SBUF_BLK \
    (\mpi_buf(sbuf, rcount, datatype)+\mpi_comm_rank*rcount)
#define MY_RBUF \mpi_buf(rbuf, rcount, datatype)
#define RBUF_AS_SBUF \mpi_buf(rbuf, rcount * \mpi_comm_size, datatype)
#define AGREE(x) \mpi_agree((x))
/*@ mpi uses comm;
    mpi collective(comm):
        requires \valid(MY_RBUF) && rcount >= 0 && AGREE(rcount) &&
            AGREE(datatype) && AGREE(op);
        requires \mpi_nonoverlapping(datatype);
        assigns *MY_RBUF;
        ensures \mpi_reduce(*MY_RBUF, 0, \mpi_comm_size, op,
            \lambda integer t; \mpi_on(
                sbuf != MPI_IN_PLACE ? *MY_SBUF_BLK :
                    \old(*MY_RBUF),
                t));
        waitsfor {i | int i; 0 <= i < \mpi_comm_size};
        behavior not_in_place:
            assumes sbuf != MPI_IN_PLACE;
            requires AGREE(sbuf != MPI_IN_PLACE);
            requires \valid_read(SBUF);
        behavior in_place:
            assumes sbuf == MPI_IN_PLACE;
            requires AGREE(sbuf == MPI_IN_PLACE);
            requires \valid(RBUF_AS_SBUF);
        disjoint behaviors;

```

```

        complete behaviors;
    */
int MPI_Reduce_scatter_block(const void *sbuf, void *rbuf, int rcount,
                           MPI_Datatype datatype, MPI_Op op,
                           MPI_Comm comm);

```

## B.16 MPI\_Scan

```

#define SBUF \mpi_buf(sbuf, count, datatype)
#define RBUF \mpi_buf(rbuf, count, datatype)
#define AGREE(x) \mpi_agree((x))
/*@ mpi uses comm;
    mpi collective(comm):
        requires \valid(RBUF) && count >= 0;
        requires AGREE(count) && AGREE(datatype) && AGREE(op);
        requires \mpi_nonoverlapping(datatype);
        requires sbuf != MPI_IN_PLACE ==> \valid_read(SBUF);
        requires \separated(RBUF, {SBUF | int i; sbuf != MPI_IN_PLACE});
        assigns *RBUF;
        ensures \mpi_reduce(*RBUF, 0, \mpi_comm_rank + 1, op,
                           \lambda integer t; \mpi_on(
                               sbuf != MPI_IN_PLACE ? *SBUF :
                               \old(*RBUF),
                           t));
        waitsfor {i | int i; 0 <= i < \mpi_comm_rank};
    */
int MPI_Scan(const void *sbuf, void *rbuf, int count,
            MPI_Datatype datatype, MPI_Op op, MPI_Comm comm);

```

## B.17 MPI\_Scatter

```

#define SBUF_OF(id) (\mpi_buf(sbuf, scout, stype) + (id) * scout)
#define SBUF \mpi_buf(sbuf, scout * \mpi_comm_size, stype)
#define RBUF \mpi_buf(rbuf, rcount, rtype)
#define SSIG (scout*\mpi_sig(stype))
#define RSIG (rcount*\mpi_sig(rtype))
#define AGREE(x) \mpi_agree((x))
/*@ mpi uses comm;
    mpi collective(comm):
        requires AGREE(root) && 0 <= root < \mpi_comm_size;
        behavior root:
            assumes \mpi_comm_rank == root;
            requires scout >= 0 && \valid_read(SBUF);
            requires rbuf == MPI_IN_PLACE || \valid(RBUF);
            requires rbuf != MPI_IN_PLACE ==> SSIG == RSIG &&
                \mpi_nonoverlapping(rtype);
            requires rbuf != MPI_IN_PLACE ==> \separated(RBUF, SBUF);
            assigns {*RBUF | int i; rbuf != MPI_IN_PLACE};
            ensures rbuf != MPI_IN_PLACE ==> *RBUF == *SBUF_OF(root);

```

```

    waitsfor \nothing;
behavior nonroot:
    assumes \mpi_comm_rank != root;
    requires rcount >= 0 && \valid(RBUF) &&
        \mpi_nonoverlapping(rtype);
    requires RSIG == \mpi_on(SSIG, root);
    assigns *RBUF;
    ensures \forall int i; i == \mpi_comm_rank ==>
        *RBUF == \mpi_on(*SBUF_OF(i), root);
    waitsfor root;
disjoint behaviors;
complete behaviors;
*/
int MPI_Scatter(const void* sbuf, int scount, MPI_Datatype stype,
               void* rbuf, int rcount, MPI_Datatype rtype, int root,
               MPI_Comm comm);

```

### B.18 MPI\_Scatterv

```

#define SBUF_OF(id) (\mpi_buf(sbuf, scounts[(id)], stype)+displs[(id)])
#define RBUF        \mpi_buf(rbuf, rcount, rtype)
#define SSIG_OF(id) (\mpi_sig(stype) * scounts[(id)])
#define RSIG        (\mpi_sig(rtype) * rcount)
#define AGREE(x)     \mpi_agree((x))
/*@ mpi uses comm;
mpi collective(comm):
    requires AGREE(root) && 0 <= root < \mpi_comm_size;
behavior root:
    assumes \mpi_comm_rank == root;
    requires \valid_read(scounts + (0 .. \mpi_comm_size-1));
    requires \valid_read(displs + (0 .. \mpi_comm_size-1));
    requires \forall int i; 0 <= i < \mpi_comm_size
        ==> scounts[i] >= 0;
    requires \forall int i, j; 0 <= i < j < \mpi_comm_size
        ==> (displs[i] + scounts[i] <= displs[j]) ||
            (displs[j] + scounts[j] <= displs[i]);
    requires \forall int i; 0 <= i < \mpi_comm_size &&
        i != root
        ==> \valid_read(SBUF_OF(i));
    requires rbuf != MPI_IN_PLACE
        ==> \valid_read(SBUF_OF(root)) && \valid(RBUF) &&
            SSIG_OF(root) == RSIG &&
            \mpi_nonoverlapping(rtype);
    requires rbuf != MPI_IN_PLACE
        ==> \separated(RBUF,
            {{ SBUF_OF(i) | int i; 0 <= i < \mpi_comm_size },
             scounts + (0 .. \mpi_comm_size-1),
             displs + (0 .. \mpi_comm_size-1)});

    assigns {*RBUF | int i; rbuf != MPI_IN_PLACE};

```

```

        ensures rbuf != MPI_IN_PLACE ==> *RBUF == *SBUF_OF(root);
        waitsfor \nothing;
    behavior nonroot:
        assumes \mpi_comm_rank != root;
        requires rcount >= 0 && \valid(RBUF) &&
            \mpi_nonoverlapping(rtype);
        requires \forall integer i; i == \mpi_comm_rank ==>
            RSIG == \mpi_on(SSIG_OF(i), root);
        assigns *RBUF;
        ensures \forall integer i; i == \mpi_comm_rank ==>
            *RBUF == \mpi_on(*SBUF_OF(i), root);
        waitsfor {root | int i; 0 <= i < \mpi_comm_size};
    disjoint behaviors;
    complete behaviors;
*/
int MPI_Scatterv(const void *sbuf, const int *scounts,
                const int *displs, MPI_Datatype stype, void *rbuf,
                int rcount, MPI_Datatype rtype, int root,
                MPI_Comm comm);

```