

Improving the Efficiency of Deadlock Detection in MPI Programs through Trace Compression (Full Paper with Proofs)

Yu Huang, Tao Wang, Zihui Yin, Eric Mercer, Benjamin Ogles

Abstract—This paper presents a static deadlock analysis for single-path MPI programs. Deadlock is when processes are blocked indefinitely by a circular communication dependency. A single path program is one that does not decode messages for control flow. The analysis records a program execution in the form of a trace and then determines from that trace whether there exists any feasible deadlocking schedules. The primary contribution is the combining of identical consecutive sends or receives into single macro actions. This simplified trace is analyzed for potential deadlock cycles. An abstract machine identifies infeasible cycles, and those not identified by the machine are encoded as satisfiability problems for an SMT solver to resolve. The action combination reduces the complexity of identifying and filtering cycles before needing the costly SMT solver. This paper shows the effectiveness of the action combination in experiments on a benchmark suite comparing to traces without action combination and other state-of-the-art deadlock analyses.

Index Terms—Static Analysis, MPI, SMT



1 INTRODUCTION

The *message passing interface* (MPI) is the *de facto* standard for communication and synchronization in high performance distributed programs. MPI programs contain a finite set of *endpoints* that send and receive messages concurrently. A *single-path* MPI program is one where the order of actions issued by each endpoint is deterministic for a given input. A common error in MPI programs, referred to as *deadlock*, occurs when one or more endpoints block indefinitely due to a circular communication dependency. This paper presents a static analysis to discover deadlocks in single-path MPI programs.

This new static deadlock analysis is an optimization to the predictive analysis in [1], and as such, this paper follows that presentation making clear where it deviates from the original work. The work in [1] statically detects deadlock in *concurrent trace programs* (CTPs). A CTP is an observed execution of an MPI program that records the sequences of send, receive, and blocking actions on each endpoint.

The static deadlock analysis in [1] examines the CTP to identify potential *deadlock cycles* that might exist in other execution orders that are allowed by the MPI semantics. A deadlock cycle is a sequence of actions that form a circular communication dependency and thus never complete in the runtime. The analysis identifies such candidate deadlock cycles by constructing and then searching a dependency graph from the CTP. Each candidate deadlock cycle is then used as a target state in an abstract machine that runs the CTP with

abstract semantics. The machine takes linear time and space to identify those candidates that may exist in some feasible execution schedule in the original MPI programs. Potential feasible deadlock cycles are then proved to be reachable, or not, in the original single-path MPI program with a novel encoding as a satisfiability problem that can be dispatched to an SMT solver.

The work in [1] includes a proof that the static deadlock analysis is sound: meaning that if no deadlock candidates are found, then no deadlock cycle exists in the original single-path MPI program. The work includes experiments in a benchmark set. The results of those experiments show that the static deadlock analysis is effective in only generating a small set of candidates, if any, that need to be verified with a satisfiability problem. It further shows that the static deadlock analysis is significantly faster, and more efficient, than existing state of the art deadlock analysis tools. The experiments include a set of CTPs for which the static deadlock analysis completes while all the other tools timeout.

Although [1] is more efficient at deadlock detection than other state-of-art tools, it suffers from state explosion in the cycle analysis on the dependency graph from the CTP. The dependency graph constructed from a CTP creates bidirectional edges between send and receive pairs with common endpoints to represent viable pairings for message communication. *Wildcard* receives multiply such edges as these accept messages from any sending source endpoint to the destination endpoint. The final dependency graph thus over-approximates what is feasible in any MPI runtime and is thereby sound in the cycle analysis.

The time complexity of cycle analysis is $O((n+e)(c+1))$ where n is the number of nodes in the graph, e is the number of edges, and c is the number of cycles. Every added edge increases the complexity of the algorithm and more so since adding one edge can create many new cycles. This work reduces the cost of cycle detection by reducing n , e , and c in

- Yu Huang, Tao Wang and Zihui Yin are with the School of Information and Engineering, Southwestern University of Finance and Economics, Chengdu, China. Email: {yuhuang, ccnuwt}@swufe.edu.cn, yinzihui98@gmail.com.
- Eric Mercer and Benjamin Ogles are with the Computer Science Department, Brigham Young University, Provo, USA. Email: egm@cs.byu.edu, benjaminogles@gmail.com.

the dependency graph.

Common to the CTPs in the benchmark suite is a communication pattern where one endpoint sends consecutive messages to another endpoint, or one endpoint receives consecutive messages from another endpoint. The work in this paper makes the observation that these consecutive groups of actions can be collapsed into a single macro action in each endpoint where they occur. With that, the original CTP can thus be reduced to a new CTP over these macro actions. A dependency graph is defined over CTPs with macro actions that inherently has fewer nodes, edges, and thereby cycles when compared to dependency graphs from the original CTP. Cycle detection in the reduced dependency graph is shown to be sound with respect to finding deadlock cycles and significantly more efficient in finding those cycles when compared to its unsimplified counterpart.

The macro actions necessitated changes in the abstract machine that detects infeasible, or unreachable, cycles. Those changes are somewhat direct and expected but merit proper documentation as given here. The proof that the abstract machine is sound is updated accordingly and stated in the full version of paper [2] because of page limit. The time and space complexity of the abstract machine does not change with the addition of macro actions.

The use of macro actions to reduce the CTP complicates the encoding of the deadlock cycle as a satisfiability problem to prove its feasibility in the original CTP. Deadlock cycles in the reduced CTP include all sends and receives in the macro actions even though only a prefix of those may be needed to create the deadlock in the original program. New in this paper is an improved encoding of the deadlock cycle that allows prefix matching of macro actions. The new encoding additionally reduces the number of variables the solver must resolve to prove or disprove the deadlock leading to runtime improvements in the backend solver.

The reduced CTP over macro actions, with the new dependency graph and cycle detection, the new abstract machine with macro actions, the new satisfiability encoding allowing for prefix matching, and proofs of soundness stated in [2] that constitute the contributions in this research are over those in [1]. The effectiveness of this improved static deadlock analysis is shown in the benchmarks from [1] and in new benchmarks where the new analysis completes and the old analysis does not.

A threat to validity in this research and its results is that the benchmark suite is not representative of MPI program paradigms in general or that the restriction to single-path programs makes it hard to say anything about arbitrary input. There is no easy way to argue either side of these critiques. The authors have sought to include benchmark programs from all relevant related work and other programs from public repositories. The efficacy of the macro actions depends on size and frequency of consecutive identical actions. When they exist, the new analysis is more efficient, when they do not exist, the analysis is still more efficient with the new satisfiability encoding than existing state-of-art. Single-path programs represent behavior on a given input. Multiple inputs would need to be considered to cover all branching behavior in an MPI program which is a difficult problem in and of itself.

The formalism and notation from [1] is used for back-

```

a      ::= nb | bb
nb     ::= (s i p src dst n)
        | (r i p src dst n)
bb     ::= (w i p i)
        | (b i p g)
src    ::= p | *
dst    ::= p

```

Fig. 1. Types of Actions $a \in \mathcal{A}$.

ground and context here. As an aid to the reader, each section opens with an explanation of what is new. Additionally, quotation marks, "...", delineate text that is included verbatim with only minor edits from [1].

2 CONCURRENT TRACE PROGRAMS

We formalize an observed MPI program execution as a concurrent trace program (CTP). A CTP is a set of communication actions $A \subset \mathcal{A}$ where \mathcal{A} is the set of all possible MPI actions. The syntax and structure of actions in \mathcal{A} is shown in Fig. 1. It differs from [1] in the addition of counters for the send and receive actions to use in the trace compression.

Each action has a unique identifier i to index its position in the owning endpoint $p \in P(A)$ where $P(A)$ is the set of all endpoints in the CTP A . For a set of actions X and for each endpoint $p \in P(X)$, $X_p \subseteq X$ is the projection of the actions in X onto the endpoint p .

Non-blocking actions are used to asynchronously send (s) and receive (r) a number of messages between endpoints. Such actions can be singleton, one message, or a macro with n indicating the number of consecutively issued actions with the same source and destination endpoints. The source (src) and destination (dst) endpoints for a message are indicated by endpoint identifiers in the body of these actions where src can be $*$ to allow for *wildcard* receive. A wildcard receive accepts messages from any source endpoint.

Messages in an MPI program are held in buffers allocated by the user program and copied from source to destination buffers by the MPI runtime. The message buffers are not part of the CTP because the CTP is only the communication observed on a single-path of the original MPI program on some input. The analysis in this paper assumes that that observed communication is deterministic meaning that the same messages are sent and received whenever that same input is given. In this way the analysis considers different schedules on that communication to identify schedules that lead to deadlock.

Blocking actions are used to halt an endpoint until some condition becomes true, possibly synchronizing two or more endpoints. The wait (w) action blocks until the message buffer of a non-blocking action is available. Even though message buffers are not modeled in the CTP, the wait actions are needed to indicate the point where its corresponding actions are complete vis-a-vis the issuing endpoint. The barrier (b) action blocks until each endpoint in a group (indicated with a unique identifier g) has reached the same barrier. Blocking send and receive actions are not included in \mathcal{A} because they are accurately modeled by placing a wait action directly after the non-blocking action.

p_0	p_1	p_2
$s_1(p_1, 1)$	$r_0(*, 1)$	
	$w_2(r_0)$	
	$r_4(p_2, 1)$	$s_3(p_1, 1)$
$w_5(s_1)$		$w_6(s_3)$
	$w_8(r_4)$	$s_7(p_1, 1)$
	$r_9(p_2, 1)$	
$s_{10}(p_2, 1)$		$w_{11}(s_7)$
	$w_{13}(r_9)$	$r_{12}(p_0, 1)$
	$r_{14}(*, 1)$	
		$w_{15}(r_{12})$
$w_{17}(s_{10})$	$w_{18}(r_{14})$	$s_{16}(p_1, 1)$
		$w_{19}(s_{16})$
$b_{20}(0)$	$b_{21}(0)$	$b_{22}(0)$

Fig. 2. Example CTP with hidden deadlock.

For simplicity in the presentation of the analysis, we assume that the last action in each endpoint $p \in P(A)$ is a barrier action ($b \ i_p \ p \ g$) for some group $g \in G(A)$ where $G(A)$ is the set of barrier groups in A . This assumption does not affect the programs we analyze as we can always extend the observed CTP with these actions. The added barrier action at the end of each endpoint is necessary to model all potential deadlocks in the dependency graph including those arising from wildcard receives.

The CTP can be constructed by the original send and receive actions with $n = 1$, or it can be further compressed by combining these actions into new send and receive actions with $n > 1$. The rules of action combination is given in Section 5. The analysis in this paper is more efficient and scalable for a compressed CTP that combines consecutive identical actions, but both the uncompressed and compressed CTPs can be correctly analyzed.

The disjoint sets of send, receive, wait, and barrier actions in A are denoted respectively by $S(A)$, $R(A)$, $W(A)$ and $B(A)$. For an action $a \in A$, $id(a)$ and $pid(a)$ denote the action and endpoint identifiers of a .

For an action $a \in S(A) \cup R(A)$, $src(a)$ and $dst(a)$ denote the source and destination endpoint of a . $src(a) = pid(a)$ and $dst(a) = pid(a)$ for sends and receives respectively. The notation $n(a)$ denotes the number of actions that are combined into a . For an action $a \in W(A)$, $req(a)$ is the non-blocking action that a waits for. For an action $a \in B(A)$, $g = grp(a)$ is the group identifier for a and $B_g(A)$ is the set of actions in that group.

3 COMPRESSION EXAMPLE

The example in this section is different from that in [1]. It also uses a more clear and compact visualization of illustra-

p_0	p_1	p_2
$s_1(p_1, 1)$	$r_0(*, 1)$	
	$w_2(r_0)$	
	$r_4(p_2, 2)$	$s_3(p_1, 2)$
$w_5(s_1)$		$w_{11}(s_3)$
$s_{10}(p_2, 1)$		$r_{12}(p_0, 1)$
	$w_{13}(r_4)$	
	$r_{14}(*, 1)$	
		$w_{15}(r_{12})$
$w_{17}(s_{10})$	$w_{18}(r_{14})$	$s_{16}(p_1, 1)$
		$w_{19}(s_{16})$
$b_{20}(0)$	$b_{21}(0)$	$b_{22}(0)$

Fig. 3. Compressed CTP from Fig. 2.

$$\begin{aligned}
s_1 &\rightarrow w_5 \rightarrow s_{10} \rightarrow r_{12} \rightarrow w_{15} \\
&\rightarrow s_{16} \rightarrow r_4 \rightarrow w_{13} \rightarrow r_{14} \rightarrow s_1
\end{aligned}$$

Fig. 4. Cycle for deadlock in Fig. 3.

tive CTP programs as seen in Fig. 2. Specifically, the action IDs are subscripts on the action types while the issuing endpoint of all actions, source endpoint of a send action, and the destination endpoint of a receive actions are all conveyed visually through columniation of the endpoints. The vertical spacing in the figure represents the observed total order of actions in the CTP.

In the example shown in Fig. 2, the wildcard receive r_0 can match with the sends s_1 or s_3 . In the observed execution, r_0 matches with s_1 , leaving the receives r_4 and r_9 to match with the sends s_3 and s_7 respectively. If however the message race is resolved by matching r_0 with s_3 , then a deadlock occurs.

The extended analysis in this paper reduces the example CTP by combining matching actions within a process. The reduced CTP in Fig. 3 replaces the four actions in p_1 with two actions, $r_4(p_2, 2)$ and $w_8(r_4)$ and the four actions in

p_0	p_1	p_2
	$r_0(*, 1)$	$s_3(p_1, 2)$
	$w_2(r_0)$	
	$r_4(p_2, 2)$	
$s_1(p_1, 1)$		$w_{11}(s_3)$
$w_5(s_1)$		$r_{12}(p_0, 1)$
	$w_{13}(r_4)$	
		$w_{15}(r_{12})$

Fig. 5. Witness execution for deadlock in Fig. 4.

p_2 with two actions, $s_3(p_1, 2)$ and $w_{11}(s_3)$. It then builds a dependency graph for the reduced CTP to find deadlock candidates.

The graph for Fig. 3 contains 22 nodes and 69 edges, versus the 26 nodes and 97 edges for that in Fig. 2. Our algorithm detects three deadlock candidates and filters away two that are provably infeasible, while there are eight candidates identified in the original CTP. The fewer candidates is a result of the action combination. For example, the reduced CTP combines s_3 and s_7 into one action and that one action appears in an infeasible deadlock candidate involving b_{21} , whereas the original CTP would generate two infeasible candidates with the same blocking action since it considers both send actions separately.

After filtering infeasible candidates, our analysis gives the cycle in Fig. 4 as a potential deadlock. A deadlock candidate itself is a characterization of the deadlock involving only the blocking actions on the cycle: $\{w_5, w_{13}, w_{15}\}$. This candidate is encoded as an SMT problem to see if it exists in some real execution of the system.

The encoding for the candidate asks the SMT solver to find a state where w_5, w_{13} and w_{15} are issued and where s_1, r_4 and r_{12} cannot be matched. Fig. 5 shows the witness execution from the satisfying assignment discovered by the solver. Here the send s_3 is issued, and since there are two sends from s_3 , it matches the single receive of r_0 and one of the two receives from r_4 . These can only match on messages from p_2 . That leaves an outstanding receive from r_4 still pending so the process is blocked on w_{13} waiting for that future send. The send s_1 is issued and p_1 is then blocked on w_5 . The receive r_{12} is issued and then blocked on w_{15} . No more actions can issue, and none of the pending send and receives in the runtime match. The CTP is deadlock.

4 SEMANTICS OF CONCURRENT TRACE PROGRAMS

The semantics in [1] are extended here to define the behavior of the macro send and receive actions with the added counter in the syntax. The changes are localized in two of the semantic rules, **Issue-Send-Receive** and **Match-Send-Recv**, discussed in Section 4.1. The following background is verbatim from the prior work.

“The semantics of A is given by a finite state machine $\mathcal{F}(A) = \langle Q, q_0, \rightarrow \rangle$ where $Q \subseteq 2^A \times 2^A \times 2^A$ is the set of states, $q_0 = \langle A, \emptyset, \emptyset \rangle$ is the start state and $\rightarrow \subseteq Q \times Q$ is the transition relation. In a state $q = \langle A, I, M \rangle$, A is the set of actions in the CTP, $I \subseteq A$ is the set of actions that have been issued to the runtime, and $M \subseteq A$ is the set of actions that have been *matched*.

For an endpoint $p \in P(A)$, the actions owned by p are always issued sequentially. This constraint can be captured as a partial order over A if we assume that action identifiers were assigned in ascending order while observing the original execution.

Definition 1 (Endpoint order). Endpoint order is a partial order (A, \leq_{po}^A) where

$$\forall a, a' \in A : a \leq_{po}^A a' \iff pid(a) = pid(a') \wedge id(a) \leq id(a')$$

Issue-Send-Receive Transition

$$\frac{a \in S(A) \cup R(A) \setminus I \quad \langle \leq_{po}^{-1} [a] \subseteq I \quad (\langle \leq_{po}^{-1} [a] \cap (W(A) \cup B(A))) \subseteq M \quad a' = a[0] \rangle}{\langle A, I, M \rangle \rightarrow \langle A, I \cup \{a\}, M \cup \{a'\} \rangle}$$

Issue-Other Transition

$$\frac{a \in B(A) \cup W(A) \setminus I \quad \langle \leq_{po}^{-1} [a] \subseteq I \quad (\langle \leq_{po}^{-1} [a] \cap (W(A) \cup B(A))) \subseteq M \rangle}{\langle A, I, M \rangle \rightarrow \langle A, I \cup \{a\}, M \rangle}$$

Match-Send-Recv Transition

$$\frac{a, a' \in I \setminus M \quad a \in S(A) \quad a' \in R(A) \quad dst(a) = dst(a') \quad src(a') \in \{src(a), *\} \quad \langle \leq_{mo}^{-1} [a] \cup \leq_{mo}^{-1} [a'] \subseteq M \quad a_m, a'_m \in M \quad a_m[0] = a[0] \quad a'_m[0] = a'[0] \rangle}{\langle A, I, M \rangle \rightarrow \langle A, I, M \setminus \{a_m, a'_m\} \cup \{a_m[n(a_m) + 1], a'_m[n(a'_m) + 1]\} \rangle}$$

Match-Wait Transition

$$\frac{a \in I \setminus M \quad a \in W(A) \quad \langle \leq_{mo}^{-1} [a] \subseteq M \rangle}{\langle A, I, M \rangle \rightarrow \langle A, I, M \cup \{a\} \rangle}$$

Match-Barrier Transition

$$\frac{g \in G(A) \quad \forall a \in B_g(A), a \in I \wedge \langle \leq_{mo}^{-1} [a] \subseteq M \rangle}{\langle A, I, M \rangle \rightarrow \langle A, I, M \cup B_g(A) \rangle}$$

Fig. 6. Transition Rules for Concrete Semantics.

For Definition 1 and similarly for any other partial order defined as follows, we use \leq_{po}^A to mean the partial order with reflexivity removed and we omit A from the notation when it is clear from context.

Send and receive actions are matched together by the runtime when their source and destination endpoints are compatible. Wait actions do not match other actions but are instead “matched” with themselves after their associated message request is matched. Finally, barrier actions match the other actions in their group when they have all been issued.

Some actions in the same endpoint may be matched in an order different from how they were issued while others must be matched in the order they were issued. This constraint is captured by two more partial orders.

Definition 2 (Queue order). Queue order is a partial order (A, \leq_{qo}^A) where for all actions $a, a' \in A$, $a \leq_{qo}^A a'$ if and only if $a \leq_{po}^A a'$ and one of the following is true:

- 1) $\{a, a'\} \subseteq S(A) \wedge dst(a) = dst(a')$
- 2) $\{a, a'\} \subseteq R(A) \wedge src(a) \in \{src(a'), *\}$

Definition 2 defines a first-in-first-out (FIFO) ordering over messages communicated on the same endpoint. With one exception, this order fully supports the “non-overtaking” property of ordered messages as defined in the MPI standard. The exception occurs when a deterministic

receive is followed by a wildcard receive in the same endpoint.

In this case, FIFO ordering over the two actions is enforced only if they can both match the same send action. This condition is schedule-dependent as its value changes depending on whether a send action has been issued that can match the first action when the second action is issued. As in [3], we leave such receive actions unordered.

Definition 3 (Match order). Match order is a partial order (A, \leq_{mo}^A) where for all actions $a, a' \in A$, $a \leq_{mo}^A a'$ if and only if $a \leq_{po}^A a'$ and one of the following is true:

- 1) $a \leq_{go}^A a'$
- 2) $a \in W(A) \cup B(A)$
- 3) $a \in S(A) \cup R(A) \wedge a' \in W(A) \wedge a = req(a')$

Definition 3 ensures that (1) queue order is preserved when messages are matched, (2) blocking actions are matched before subsequent actions in the same endpoint and (3) message requests are matched before their associated wait actions. With match order defined, we can define the transition relation \rightarrow shown in Fig. 6. Given a relation Q over a set X and an element $x \in X$, $Q^{-1}[x] = \{y \in X : yQx\}$ is the preimage of x under Q .

4.1 Semantic Rules

The issue of a new action a is completed by the **Issue-Send-Receive** or **Issue-Other** transitions when all of the actions preceding a in the same endpoint have been issued and all of the blocking actions preceding a in the same endpoint have been matched. Each transition updates the state by putting a into the set I. If a is a send or receive action, the **Issue-Send-Receive** transition additionally creates an action a' , which is identical to a except that n in the tuple is initialized to 0. The action a' itself has no meaning to the CTP execution, but when it is updated in the set M by incrementing the number n iteratively, the action a can be matched gradually. To help presenting the step of matching send and receive actions, let the action $a[n]$ be defined as follows.

$$a[n] = \begin{cases} (s \ id(a) \ pid(a) \ src(a) \ dst(a) \ n) & a \in S(A) \\ (r \ id(a) \ pid(a) \ src(a) \ dst(a) \ n) & a \in R(A) \end{cases}$$

The **Match-Send-Recv** transition completes type compatible matches for the actions $a \in R(A)$ and $a' \in S(A)$ when their match order dependencies have been satisfied. Additionally, the associated actions a_m and a'_m are updated in the set M by incrementing the number n in each of their bodies by one at a time, indicating that a send and a receive combined in a and a' respectively are matched. Eventually, the action a is completely matched if the number n in a_m reaches $n(a)$ after all the **Match-Send-Recv** transitions that a takes part in are completed. In other words, the structure of action a_m is the same with a at this moment. Note matching a' needs similar iterations.

A complete match is made by the **Match-Wait** transition. The rule forces executions to follow a *rendevous protocol*. Rendezvous protocol synchronizes the sender and receiver by blocking both until the message transfer is completed.

5 COMPRESSION OF CONCURRENT TRACE PROGRAMS

This entirely new section discusses the detailed steps of action combination. Action combination is only adaptable to a set of consecutive sends or a set of consecutive receives in an identical endpoint. More precisely,

Definition 4 (Action Combination). Two actions a and a' such that,

- 1) $pid(a) = pid(a') \wedge src(a) = src(a') \wedge dst(a) = dst(a')$, and
- 2) $id(a) < id(a') \wedge \forall a'' \in A : id(a) < id(a'') \wedge id(a'') < id(a') \Rightarrow a'' \in W(A)$,

can be combined, denoted as $Comb(a, a', A)$, to yield a new action

$$a_c = \begin{cases} (s \ id(a) \ pid(a) \ src(a) \ dst(a) \ (n(a) + n(a'))) & a, a' \in S(A) \\ (r \ id(a) \ pid(a) \ src(a) \ dst(a) \ (n(a) + n(a'))) & a, a' \in R(A) \end{cases}$$

and a new CTP

$$A' = A \cup \{a_c\} \cup \{\mathbf{w}_{(id(a_c)+1)}(pid(a_c))id(a_c)\} \setminus (\{a, a'\} \cup \{w \in W(a) : req(w) = a \vee req(w) = a'\})$$

The combination in Definition 4 works for any pair of two sends or two receives if their endpoints match and they are consecutively ordered in an identical endpoint. The new action a_c and its associated wait are created to replace the original actions and their waits, and consequently updates the CTP.

Given the basic steps of action combination in Definition 4, a CTP can be compressed by the operator \xrightarrow{Comb} .

Definition 5 (Compressing a CTP with Action Combination). The operator \xrightarrow{Comb} is a relation on $2^A \times 2^A$ such that $\xrightarrow{Comb} = \{< A \subset \mathcal{A}, A' \subset \mathcal{A} > : \exists a, a' \in A, \exists a_c \in A', \{a_c, A'\} = Comb(a, a', A)\}$, where $\xrightarrow{Comb^*}$ denotes the transitive closure of \xrightarrow{Comb} .

This paper uses $A \xrightarrow{Comb} A'$ to represent $< A, A' > \in \xrightarrow{Comb}$. Additionally, any CTP A can be applied for an utmost compression in Definition 6 by multiple iterations of action combination.

Definition 6 (Utmost Compression). For any CTP A , there exists an utmost compressed CTP A_f such that

$$A \xrightarrow{Comb^*} A_f \wedge \nexists A'_f \subseteq \mathcal{A}, A_f \xrightarrow{Comb} A'_f.$$

The CTP A_f is constructed by combining the actions to the utmost extent, which stands for the largest simplification for the execution. The analysis in this paper intends to compress every CTP observed from the execution, A_0 to the CTP A_f , which directly optimizes the process of deadlock detection.

6 DEADLOCK

This section is unchanged from [1] and restated verbatim for convenience as it formalizes the problem statement for deadlock detection in A . The problem statement is defined over a generic transition relation $\delta \subseteq \mathcal{Q} \times \mathcal{Q}$ rather than that in the previous section because the statement is reused in

the the proof for the abstract transition relation defined in the full version of paper [2].

“Let $\Sigma_\delta^q \subseteq \mathcal{Q}$ denote the reachable states of A from the state q with respect to a transition relation δ :

$$\Sigma_\delta^q = \{q' \in \mathcal{Q} : (q, q') \in \delta^*\}$$

where δ^* denotes the transitive closure of δ .

Definition 7 (Deadlock). A state $q = \langle A, I, M \rangle$ is deadlocked with respect to a transition relation δ if there are no enabled transitions and there are actions left to be issued or matched:

$$Dead_\delta(q) \iff (I \neq A \vee M \neq A) \wedge (\forall q' \in \mathcal{Q}, (q, q') \notin \delta)$$

The deadlock discovery problem, given in Definition 8, asks whether $\mathcal{F}(A)$ can reach a deadlocked state. This problem is NP-Complete and can be directly encoded as a propositional formula [3].

Definition 8 (Deadlock discovery problem).

$$\exists q \in \Sigma_\delta^{q_0}, Dead_\delta(q)$$

The search for an arbitrary feasible deadlock state can be extremely expensive for many programs. We can give the search a kind of head start by finding a simple way to characterize the types of states that may deadlock. A convenient way to describe a state is by its *control point*. The control point of a state is simply the set of last issued actions from each endpoint:

$$Ctrl(\langle A, I, M \rangle) = \{a \in I : \forall a' \in I_{pid(a)}, a' \leq_{po} a\}.$$

If we only provide the last issued action for a subset of endpoints, we obtain a partial control point that describes the collection of states which include it as a subset of their control points. Definition 9 augments the problem statement in Definition 8 to ask for a deadlock state that matches a partial program point D (also called a deadlock candidate)."

Definition 9 (Constrained deadlock discovery problem).

$$\exists q \in \Sigma_\delta^{q_0}, D \subseteq Ctrl(q) \wedge Dead_\delta(q)$$

7 DEPENDENCY GRAPH

This section defines how to compute a sound set of deadlock candidates $\mathbb{D}(A)$ for A . The dependency graph construction uses the same set of rules as those in [1]. New here is the definition of a deadlock path with macro actions, the set of potential matches for macro actions, and how that set is over-approximate to build the dependency graph. The algorithm for finding potential deadlock cycles is changed accordingly.

"The soundness property according to Definition 9 is formally stated in the following theorem.

Theorem 1 (Deadlock candidates sound). For all states $q \in \Sigma_\delta^{q_0}$, If $Dead_\delta(q)$, then $\exists d \in \mathbb{D}(A), d \subseteq Ctrl(q)$.

We generate $\mathbb{D}(A)$ by detecting cycles in a graph (N, E) where $N = A \cup \{\perp_p : p \in P(A)\}$ is the set of nodes and $E : N \times N$ is the set of edges. The node \perp_p is used to explicitly represent the end of endpoint p in the graph. An edge $(a, a') \in E$ represents a potential communication dependency of a' on a in some execution of A . A proof that

our technique satisfies Theorem 1, along with proofs for the other theorems stated in this paper, is given in the appendix of the full paper [2].

Before presenting the dependency graph, we describe how one of its cycles can represent a deadlock. This not only motivates the rules for adding edges to the graph, but also leads to a precise understanding of the type of cycle the analysis must report and the types it can ignore. This is important because the graphs we build may contain a huge number of cycles which are expensive to enumerate. The more we can ignore, the more efficient our analysis will be.

Fix a deadlock state $q \in \Sigma_\delta^{q_0}$ with control point $D = Ctrl(q)$. For each endpoint $p \in P(A)$, there is an action $a \in D$ with $p = pid(a)$. We will define a few new terms that allow us to talk about why a is blocking p from progressing in the state q .

First, we call a the *deadlock action* for p . Next, let a' be the earliest action in p that is issued in q but not matched with $a' \leq_{mo} a$. We call a' the *orphaned action* for p . Finally let a'' be an action that is not issued in q but would allow p to progress if it is matched with a' . We call a'' the *parent action*. If the parent action does not exist, it is represented in the graph by a \perp node (discussed more below). We will use the following definition to translate these concepts to the context of a path of edges in E ."

Definition 10 (Deadlock path). Let $a_0 \rightarrow a_1 \rightarrow \dots \rightarrow a_{m-1} \rightarrow a_m$ be a path of edges in E . This path is a *deadlock path* for the endpoint $p \in P(A)$ if

- 1) $pid(a_0) \neq p$ and
- 2) $pid(a_i) = p$ for all $i \in \{1 \dots m\}$ and
- 3) if $m > 1$, $a_i \in W(A) \cup B(A)$ for some $i \in \{1 \dots m-1\}$ and
- 4) if $m = 1$, $a_1 \in R(A)$ such that $n(a_1) > 1$ and $src(a_1) = *$.

Definition 10 specifies two typical instances of deadlock path. The first instance is constrained by the conditions (1) – (3), where a_0 and a_1 are interpreted as parent and orphan actions of the endpoint p for a deadlock path $a_0 \rightarrow a_1 \rightarrow \dots \rightarrow a_{m-1} \rightarrow a_m$ for p . The edge connecting them represents the possibility that a_1 may depend on a_0 being issued and available to match to complete in some execution.

Specially, if $n(a_0) > 1$ or $n(a_1) > 1$, the edge implies more than one (probably a huge number of) possible dependencies between the actions combined in a_0 and a_1 . The integration of these dependencies to a single edge is helpful to highly reduce the number of detected cycles, while still preserves all possible circular dependency between the actions.

The earliest blocking action issued by p and contained in the path is interpreted as the deadlock action of p . The path requires that the deadlock action not to be the last action in the path. This is because a deadlock cycle is constructed by composing the deadlock paths of two or more endpoints. In other words a_m in the deadlock path of p will be the parent action in a different deadlock path for some endpoint $p' \neq p$.

The second instance of deadlock path is enforced by the conditions (1), (2) and (4), where a_1 is the only action in p . Unlike a common understanding of deadlock path, the action a_1 here forms two feasible connections $a_0 \rightarrow a_1$ and

$a_1 \rightarrow a_2$ without inserting any deadlock action. The action a_2 is the orphaned action in some other deadlock path. In this case, the action a_1 can only be a wildcard receive with $n(a_1) > 1$. The actual deadlock action is the very next wait w for a_1 in $pid(a_1)$. w must exist according to action combination. To let the deadlock path contribute to a feasible deadlock cycle, a_1 has to be either unmatched or prefix matched, meaning not all actions in a_1 are matched, leaving some suffix action orphaned in the associated deadlock cycle. Note that a_1 can not be a send action even if $n(a_1) > 1$, as it can not connect with two different endpoints.

Definition 11 (Deadlock cycle). A cycle of edges in E is a *deadlock cycle* if it can be constructed from a set of deadlock paths, with each endpoint contributing at most one path, and the orphaned action of each path must not be a compatible match for the orphaned action of any path in the cycle.

For example, the deadlock cycle in Fig. 4 includes three deadlock paths:

- 1) $r_{14} \rightarrow s_1 \rightarrow w_5 \rightarrow s_{10}$,
- 2) $s_{10} \rightarrow r_{12} \rightarrow w_{15} \rightarrow s_{16}$, and
- 3) $s_{16} \rightarrow r_4 \rightarrow w_{13} \rightarrow r_{14}$,

each of which is contributed by a separate endpoint. The set of deadlock actions $\{w_5, w_{15}, w_{13}\}$ is extracted as a deadlock candidate.

Deadlock cycles are detected by constructing the aforementioned graph and then enumerating the cycles in the graph. The nodes, N , being the actions, A , in the CTP with a special \perp_p node for each endpoint. The edges between nodes are defined by actions that can match, a send with a receive for example, and the match order \leq_{mo}^A .

Definition 12 (Potential matches). For an action $a \in A$, $\mathbb{M}(a)$ denotes the set of actions where each can be matched with a in at least one transition between reachable states of A . More precisely, if $M(q)$ denotes the matched set in the state q , then

$$\begin{aligned} \mathbb{M}(a) = & \{a' : \exists q, q' \in \Sigma_{\rightarrow}^{q_0}, q \rightarrow q' \\ & \wedge \{a[m], a'[m']\} \subseteq M(q) \\ & \wedge \{a[m+1], a'[m'+1]\} \subseteq M(q') \\ & \wedge 0 \leq m < n(a) \wedge 0 \leq m' < n(a')\} \setminus \{a\}. \end{aligned}$$

Computing precisely \mathbb{M} is as hard as the deadlock discovery problem itself. Potential matches are thus over approximated to ensure that the whole analysis is sound. The algorithm for approximating potential matches in [4] is extended here to work with macro send and receive actions. The extension generates a match set, \mathbb{M}_0 , over the original, uncompressed actions, A_0 , from the CTP. It then iteratively creates \mathbb{M} with a mapping from the original actions in \mathbb{M}_0 .

That mapping, $\phi : A \rightarrow 2^{A_0 \cup A}$, is defined such that for an action a , if it is a send or receive action, it returns the set of actions in A_0 that are combined by a ; otherwise, it returns $\{a\}$. Given \mathbb{M}_0 , computed by [4], and ϕ , then $\mathbb{M}(a) = \{a' \in A : \exists b, b' \in A_0, b' \in \mathbb{M}_0(b) \wedge b \in \phi(a) \wedge b' \in \phi(a')\}$. If two actions matched in \mathbb{M}_0 , then the actions that cover those also match in \mathbb{M} .

Definition 13 (Edges). $\forall (a, a') \in N \times N$, $(a, a') \in E$ if and only if any one of the following hold

- 1) $a \in A \wedge a' = \perp_{pid(a)}$
- 2) $a, a' \in A \wedge a <_{mo} a'$
- 3) $a, a' \in A \wedge a' \in \mathbb{M}(a)$
- 4) $a = \perp_{src(a')} \wedge a' \in R(A) \wedge \exists a_r \in R(A)(a_r <_{po} a' \wedge src(a_r) = * \wedge src(a') \neq *)$
- 5) $a' \in S(A) \wedge \exists a_r \in R(A)(a = \perp_{dst(a_r)} \wedge dst(a_r) = dst(a') \wedge src(a_r) = *)$

Rule one connects every action to the end of its owning endpoint. Rule two encodes the match order from Section 4. It connects deadlock actions to those in the tail of a deadlock path, and it ensures that wait actions (often a deadlock action) have an incoming edge from their associated message request (often an orphaned action).

The rest of the rules ensure that parent actions are connected to orphaned actions in any scenario. The most obvious occurs in rule three when the actions can form a match. The nature of \mathbb{M} adds an edge in each direction to reflect that either a send or a receive can be the parent action in a deadlock.

Rules four and five encode *message starvation*. Message starvation is when a wildcard receive is matched with a send action that leaves some subsequent deterministic receive without any potential future matches. The parent action does not exist for this type of deadlock as the starved action is waiting for something that can never be issued. The added \perp node at the end of every endpoint takes the place of the parent action in these cases.

As an aside, edges are not added from \perp nodes to wildcard receive actions. Deadlocks arising from starved wildcard receives are deterministic; they manifest on any schedule of the program for the given input.

Definition 13 results in graphs with fewer nodes and edges for the CTP compressed by action combination. The number of cycles, therefore, can be noticeably reduced as any edge connecting two sends or receives with $n > 1$ actually implies more than one possible matches for the original actions, which could largely expand the graph if these actions were not combined. That said, many spurious, infeasible, cycles exist in the graph. For example, the third rule alone creates trivial cycles between matching actions.

Johnson's algorithm for enumerating the elementary cycles of a directed graph [5] is modified to only enumerate cycles that match Definition 11. The algorithm enumerates the strongly connected components in the graph finding the cycles in each one. The order in which it considers strongly connected components is set such that the component containing the next unvisited least node on some total order of \leq_{po} is followed. Within that connected component, s , is the least node, and it is where the cycle search begins. This ordering ensures that s is always the orphaned action of the first deadlock path visited by the algorithm.

From here, the algorithm considers edges that can extend the current path being considered. Algorithm 1 is a boolean function that determines when a deadlock path can or cannot be extended given an edge (v, w) and the starting node s . The reference stack is the depth first search stack of visited nodes. The `all_count` and `block_count` functions respectively return the number of all actions and the number of blocking actions in a given endpoint that have been visited by the current stack. The `can_match` function

Algorithm 1 Determine whether the edge (v, w) can possibly reach a deadlock cycle starting at s

```

1: procedure DEADLOCKEDGE( $v, w, s$ )
2:   if  $pid(v) = pid(w)$  then
3:     return true
4:   if  $all\_count(stack, pid(v)) > 1 \wedge$ 
      $block\_count(stack, pid(v)) = 0$  then
5:     return false
6:   if  $block\_count(stack, pid(v)) = 1 \wedge v \in W(A) \cup$ 
      $B(A)$  then
7:     return false
8:   if  $all\_count(stack, pid(v)) = 1 \wedge \neg(v \in R(A) \wedge$ 
      $src(v) = * \wedge n(v) > 1)$  then
9:     return false
10:  if  $\exists a \in orphaned(stack), can\_match(a, w)$  then
11:    return false
12:  if  $w \neq s$  then
13:    return  $all\_count(stack, pid(w)) = 0$ 
14:  return false

```

determines whether two actions may form a compatible match based on their types and endpoints. Finally, the orphaned function returns the orphaned action from each deadlock path in the current stack.

Lines 2-3 of Algorithm 1 extend the tail of the current path along the same endpoint (rule one in Definition 10). If the path cannot be extended along the same endpoint, then v and w are the parent and orphaned actions for a new deadlock path. Lines 4-7 ensure that the current deadlock path contains a deadlock action that is not also a parent action for w (rules two and three in Definition 10). Lines 8-9 ensure that v has to be a wildcard receive with $n(v) > 1$ when the current deadlock path contains exactly one action for the endpoint $pid(v)$ (rule four in Definition 10). Lines 10-11 check that the cycle cannot unwind by matching some orphaned action (Definition 11). Lines 12-13 ensure that each endpoint only contributes one path to the cycle (also Definition 11). The final case occurs when w is equal to s and a cycle is formed.

As a final optimization not shown in Algorithm 1, duplicate candidates are not reported. The optimization tracks how actions are orphaned. If the tail of a path can reach an orphaned action in another endpoint from two different parents, then only the cycle from one such parent is needed.

8 ABSTRACT SEMANTICS OF CONCURRENT TRACE PROGRAMS

This section extends [1] to update the abstract machine semantics for compressed send and receive actions. The changes are made for the issue and match transitions.

"The abstract machine $\hat{F}(A) = \langle \mathcal{Q}, q_0, \rightarrow_{abs} \rangle$ augments the semantics of CTPs to efficiently filter away infeasible deadlock candidates. This filtering is an important stage in the analysis because it can drastically reduce the number of calls to the SMT solver. The abstract transition relation \rightarrow_{abs} is shown in Fig. 7 with the barrier transition omitted as it is unchanged from \rightarrow . In this transition relation we create

Issue-Send Transition

$$\begin{array}{c}
 a \in S(A) \setminus I \quad \leq_{po}^{-1} [a] \subseteq I \\
 (\leq_{po}^{-1} [a] \cap (W(A) \cup B(A))) \subseteq M \\
 s = (s \text{ (id}(a) - 1) \text{ pid}(a) * \text{dst}(a) \text{ n}(a)) \quad s' = s[0] \\
 a' = a[0] \\
 \hline
 \langle A, I, M \rangle \rightarrow_{abs} \langle A \cup \{s\}, I \cup \{s, a\}, M \cup \{s', a'\} \rangle
 \end{array}$$

Issue-Receive Transition

$$\begin{array}{c}
 a \in R(A) \setminus I \quad \leq_{po}^{-1} [a] \subseteq I \\
 (\leq_{po}^{-1} [a] \cap (W(A) \cup B(A))) \subseteq M \quad a' = a[0] \\
 \hline
 \langle A, I, M \rangle \rightarrow \langle A, I \cup \{a\}, M \cup \{a'\} \rangle
 \end{array}$$

Issue-Other Transition

$$\begin{array}{c}
 a \in B(A) \cup W(A) \setminus I \quad \leq_{po}^{-1} [a] \subseteq I \\
 (\leq_{po}^{-1} [a] \cap (W(A) \cup B(A))) \subseteq M \\
 \hline
 \langle A, I, M \rangle \rightarrow_{abs} \langle A, I \cup \{a\}, M \rangle
 \end{array}$$

Match-Send-Recv Transition

$$\begin{array}{c}
 a, a' \in I \setminus M \quad a \in S(A) \quad a' \in R(A) \quad \text{dst}(a) = \text{dst}(a') \\
 \text{src}(a') = \text{src}(a) \quad (\leq_{mo}^{-1} [a] \cup \leq_{mo}^{-1} [a']) \subseteq M \\
 a_m, a'_m \in M \quad a_m[0] = a[0] \quad a'_m[0] = a'[0] \\
 \hline
 \langle A, I, M \rangle \rightarrow \\
 \langle A, I, M \setminus \{a_m, a'_m\} \cup \{a_m[n(a_m) + 1], a'_m[n(a'_m) + 1]\} \rangle
 \end{array}$$

Match-Wait Transition

$$\begin{array}{c}
 a \in I \setminus M \quad a \in W(A) \\
 s = (s \text{ (id}(req(a)) - 1) \text{ pid}(req(a)) * \text{dst}(req(a)) \text{ n}(req(a))) \\
 \exists a_s, s' \in M, (a_s[0] = req(a)[0]) \wedge (s'[0] = \\
 s[0]) \wedge (n(req(a)) = n(a_s) + n(s')) \\
 \hline
 \langle A, I, M \rangle \rightarrow_{abs} \langle A, I, M \cup \{a\} \rangle
 \end{array}$$

Fig. 7. Transition Rules for Abstract Semantics.

a dedicated *wildcard endpoint* for each source process. This eliminates the possibility of message starvation."

The **Issue-Send** transition generates a fresh *wildcard send* for the new endpoint and issue it alongside the original send action. Similar to the concrete semantics, the copies of both *wildcard send* and original send actions are added to M for matching. A wait on the send action is allowed to match in the **Match-Wait** transition if the action has been completely matched, indicating the accumulation of the n numbers for the two copies of sends in M reaches that of the original send. The wildcard send is only allowed to match wildcard receive actions issued by the destination process and the original send action is only allowed to match deterministic receive actions.

"We filter a deadlock candidate D by deriving a CTP A_D that contains the actions in D , the actions process ordered before actions in D and all of the actions in other processes:

$$A_D = \bigcup_{a \in D} \leq_{po}^{-1} [a] \cup \bigcup_{p \notin P(D)} A_p$$

We then attempt to execute $\hat{F}(A_D)$ to determine whether

$$\exists q \in \Sigma_{\rightarrow abs}^{q_0}, D \subseteq Ctrl(q) \wedge Dead_{\rightarrow abs}(q)$$

Note that this is just Definition 9 with the abstract transition relation substituted in. If the abstract execution is able to issue every action in D , then the candidate may represent a real deadlock and it is added to the set of candidates to be encoded as an SMT formula. Otherwise, the candidate is infeasible and is discarded.

The abstract machine is sound if it never discards a reachable control point. Let the set of reachable control points from a state q and a transition relation δ be \mathbb{C}_δ^q .

$$\mathbb{C}_\delta^q = \{Ctrl(q') : q' \in \Sigma_\delta^q\}$$

Theorem 2 states that the reachable control points of the abstract machine subsumes the reachable control points of the concrete machine. Theorem 3 states that if the abstract machine cannot issue every action in the deadlock candidate in one execution, it will not be able to issue them all in any execution. Together these theorems prove that the candidate D can be filtered away in a single execution of the abstract machine when it fails to issue all of the actions in D .

Theorem 2 (Abstract candidate simulation). Let $q \in \Sigma_{\rightarrow}^{q_0}$ and $q' \in \Sigma_{\rightarrow abs}^{q_0}$ be a concrete and abstract state reachable from the start state q_0 . If $Ctrl(q) = Ctrl(q')$, then for all control points $D \in \mathbb{C}_\delta^q$, it follows that $D \in \mathbb{C}_{\rightarrow abs}^{q'}$.

Theorem 3 (Abstract deadlock deterministic). Let $q \in \Sigma_{\rightarrow abs}^{q_0}$ be a reachable abstract state with $Dead_{\rightarrow abs}(q)$. If $D \not\subseteq Ctrl(q)$, then for all $q' \in \Sigma_{\rightarrow abs}^{q_0}$, $D \not\subseteq Ctrl(q')$.

”

9 SMT ENCODING

The SMT encoding rules extend those in [1] to allow prefix matching for orphaned actions. These changes are isolated to the definitions and usage of two new set of variables, timestamps and matches for each send or receive action, and the MATCH COUNT and MATCH CORRECT rules.

“If $\hat{F}(A_D)$ is able to issue each action in the candidate D , then it is used to construct an SMT formula F . A satisfying assignment for F can be used to construct a witness execution for the deadlock candidate. If F is unsatisfiable, then there is no feasible deadlock state that contains D as part of its control point.”

To achieve a precise detection, the formula F expands each action $a \in S(A_D) \cup R(A_D)$ so that the match relation of each action in $\phi(a)$ is resolved. This is significant for our SMT solution because we do not know how many actions there are for completely matching a assuming $n(a) > 1$ and whether their match ordering is feasible in a witness execution without precisely encoding the actions in $\phi(a)$.

The expansion includes two steps that are novel to the existing SMT encoding. First, the formula defines two sets of variables: $\{t_a^i : i \in \{1 \dots n(a)\}\}$ and $\{m_a^i : i \in \{1 \dots n(a)\}\}$, where t_a^i and m_a^i respectively represent the timestamps of each action $b_i \in \phi(a)$ and the action that matches b_i . For simplicity, if $n(a) = 1$, the variables t_a and m_a represent any t_a^i and m_a^i respectively. If $a \in R(A_D)$ and $n(a) > 1$, the formula only preserves the variables $t_a^1, t_a^{n(a)}$,

MATCH ORDER	$\bigwedge_{a \in A_D} \bigwedge_{a' \in <_{mo}^{-1}[a]} (c_a \implies c_{a'}) \wedge t_{a'}^{n(a')} < t_a^1$ $\bigwedge_{a \in R(A_D) \wedge n(a) > 1} t_a^1 < t_a^{n(a)}$ $\bigwedge_{a \in S(A_D) \wedge n(a) > 1} \bigwedge_{i \in \{1 \dots n(a)-1\}} t_a^i < t_a^{i+1}$
QUEUE ORDER	$\bigwedge_{a \in S(A_D) \cup R(A_D)} \bigwedge_{a' \in <_{qo}^{-1}[a]} m_{a'}^{n(a')} < m_a^1$ $\bigwedge_{a \in R(A_D) \wedge n(a) > 1} m_a^1 < m_a^{n(a)}$ $\bigwedge_{a \in S(A_D) \wedge n(a) > 1} \bigwedge_{i \in \{1 \dots n(a)-1\}} m_a^i < m_a^{i+1}$
BARRIERS	$\bigwedge_{a \in B(A_D)} \bigwedge_{a' \in B_{grp(a)}(A_D)} t_a = t_{a'}$
MATCH COUNT	$\bigwedge_{a \in S(A_D)} \bigwedge_{b \in \phi(a)} atm(1, \{\mathbb{M}(a', b) : b \in \mathbb{M}_D^*(a')\})$ $\bigwedge_{a \in R(A_D)} atm(n(a), \{\mathbb{M}(a, a') : a' \in \mathbb{M}_D^*(a)\})$
MATCH CORRECT	$\bigwedge_{a \in S(A_D)} (c_a \iff \bigwedge_{b \in \phi(a)} exa(1, \{\mathbb{M}(a', b) : b \in \mathbb{M}_D^*(a')\}))$ $\bigwedge_{a \in R(A_D)} (c_a \iff exa(n(a), \{\mathbb{M}(a, a') : a' \in \mathbb{M}_D^*(a)\}))$
REACH	$\bigwedge_{a \in D} \bigwedge_{a' \in <_{mo}^{-1}[a] \setminus O} c_{a'}$
DEADLOCK	$\bigwedge_{a \in D \cup O} \neg c_a$
NO MATCHES	$\bigwedge_{a \in O} \bigwedge_{a' \in \mathbb{M}_D^*(a) \setminus (O \cup D)} c_{a'}$

Fig. 8. Constraints in the formula F

m_a^1 and $m_a^{n(a)}$, which are sufficient to constrain the scope of matching timestamps. This reduction is significantly helpful for reducing the formula size, while the precision in analysis is also kept, which is achieved by letting the explicit timestamps for a be assigned to the m variables of all matched sends after resolving the encoding.

Second, the formula expands the match set \mathbb{M}_D^* from the set \mathbb{M}_D , which is simply the match set \mathbb{M} applied to the CTP A_D . More precisely, for an action $a \in A_D$,

$$\mathbb{M}_D^*(a) = \bigcup_{a' \in \mathbb{M}_D(a)} \phi(a').$$

For each action $a \in S(A_D) \cup R(A_D)$, the set $\mathbb{M}_D^*(a)$ consists of all send and receive actions combined in the actions in $\mathbb{M}_D(a)$.

The formula uses w_a as another name for t_w where $w \in W(A_D)$ and $a = req(w)$. Additionally, it adds a boolean variable c_a for every action that is true if a must be issued and completed in the witness execution.

“The rules for the encoding are shown in Fig. 8. The MATCH ORDER and QUEUE ORDER constraints preserve the meaning of match order and queue order in the encoding. An action can only complete if its $<_{mo}$ predecessors have completed and the timestamps must reflect that. The timestamps of the first and last combined actions in a receive must be ordered. The timestamps of all combined actions in a send must follow the order by their indices. Additionally, matches must conform to the non-overtaking guarantee of MPI executions. In all cases, constraints are

MATCH ORDER	$(c_{s_1} \implies c_{w_5}) \wedge t_{s_1} < t_{w_5} \bigwedge (c_{r_0} \implies c_{w_2}) \wedge t_{r_0} < t_{w_2} \bigwedge (c_{w_2} \implies c_{r_4}) \wedge t_{w_2} < t_{r_4} \bigwedge (c_{r_4} \implies c_{w_{13}}) \wedge t_{r_4}^2 < t_{w_{13}}$ $\bigwedge (c_{s_3} \implies c_{w_{11}}) \wedge t_{s_3}^2 < t_{w_{11}} \bigwedge (c_{w_{11}} \implies c_{r_{12}}) \wedge t_{w_{11}} < t_{r_{12}} \bigwedge (c_{r_{12}} \implies c_{w_{15}}) \wedge t_{r_{12}} < t_{w_{15}} \bigwedge t_{s_3}^1 < t_{s_3}^2 \bigwedge t_{r_4}^1 < t_{r_4}^2$
QUEUE ORDER	$m_{r_0} < m_{r_4}^1 \bigwedge m_{r_4}^1 < m_{r_4}^2 \bigwedge m_{s_3}^1 < m_{s_3}^2$
MATCH COUNT	$atm(1, \{M(r_0, s_3^1), M(r_4, s_3^1)\}) \wedge atm(1, \{M(r_4, s_3^2)\})$ $\bigwedge atm(1, \{M(r_0, s_1)\}) \bigwedge atm(1, \{M(r_0, s_1), M(r_0, s_3^1)\}) \bigwedge atm(2, \{M(r_4, s_3^1), M(r_4, s_3^2)\})$
MATCH CORRECT	$(c_{s_3} \iff exa(1, \{M(r_0, s_3^1), M(r_4, s_3^1)\}) \wedge exa(1, \{M(r_4, s_3^2)\}))$ $\bigwedge (c_{s_1} \iff exa(1, \{M(r_0, s_1)\})) \bigwedge (c_{r_0} \iff exa(1, \{M(r_0, s_1), M(r_0, s_3^1)\})) \bigwedge (c_{r_4} \iff exa(2, \{M(r_4, s_3^1), M(r_4, s_3^2)\}))$
REACH	$c_{r_0} \bigwedge c_{w_2} \bigwedge c_{s_3} \bigwedge c_{w_{11}}$
DEADLOCK	$\neg c_{s_1} \bigwedge \neg c_{w_5} \bigwedge \neg c_{r_4} \bigwedge \neg c_{w_{13}} \bigwedge \neg c_{r_{12}} \bigwedge \neg c_{w_{15}}$
NO MATCHES	$c_{r_0} \bigwedge c_{s_3}$

Fig. 9. SMT encoding for the witness execution in Fig. 5

omitted when they are obviously redundant with respect to existing constraints and the transitivity of $<$ and $=$ over the integers.

The BARRIERS constraint encodes the inter-process synchronization behavior of barrier actions by asserting that groups complete at the same time. All other timestamps are asserted to be distinct.”

The MATCH COUNT constraint enforces the maximum count of matches for each send or receive action because there are at most $n(a)$ many actions can be matched. The $atm(k, Z)$ and $exa(k, Z)$ constraints are inspired by the encoding in [3], which are true if and only if *at most* and *exactly* k many constraints from the set Z are true, respectively.

The MATCH CORRECT constraint encodes that any send or receive action is completed if and only if all actions in $\phi(a)$ are matched.

The MATCH COUNT and MATCH CORRECT constraints do not exclude the orphaned action or deadlock action in any deadlock path, as each of these actions can be prefix matched meaning some of its combined actions are matched while the rest are not. Completely matching such an action is not allowed in any witness execution, and the MATCH CORRECT constraint reflects that.

Let O denote the set of orphaned actions for the candidate D .

$$O = \{a \in D \cap R(A_D) : src(a) = * \wedge n(a) > 1\} \\ \cup \{a \in A_D : \exists a' \in W(A_D) \cap D, a = req(a')\}$$

In Section 4, a send and receive action were matched by incrementing their n numbers in the matched set, meaning that each time a combined action is matched. This semantic meaning is preserved in the encoding by the $M(r, s_i)$ constraint where $r \in R(A_D)$ and $s_i \in \phi(s)$ for some $s \in S(A_D)$. Let t_s and m_s be the timestamps of s_i and the action that matches s_i , then $M(r, s_i)$ expands to

$$m_r^1 \leq t_s \leq m_r^{n(r)} \wedge t_r^1 \leq m_s \leq t_r^{n(r)} \wedge t_s < w_r \wedge t_r^{n(r)} < w_a.$$

The constraint enforces the timestamp t_s to indicate that the send must complete between the completion time of the first and last matched sends for r . Similarly, the constraint enforces the timestamp m_s meaning the matched receive for s_i must complete between the completion time of the first

and last receives combined in r . Finally, the timestamps of s_i and r are required to precede the timestamps of both wait actions. Note that in the infinite-buffer setting, the wait for the send action does not exist and so one of these constraints is omitted.

“The REACH constraint asserts that every predecessor of the actions in D is completed except the actions in O . In other words, it asserts that the deadlock D is reachable.

The DEADLOCK constraint simply asserts that the deadlock actions in D and the orphaned actions in O are not complete. The NO MATCHES constraint ensures that any issued actions that could untangle the deadlock are complete, thus forcing them to find matches that exclude the deadlock and orphaned actions. Given a satisfying assignment of the variables in F , the witness execution can be constructed from the t timestamp variables while the matches made can be recovered by consulting the m variables.”

Consider the example CTP in Fig. 3. The formula resulting from the witness execution in Fig. 5 is given in Fig. 9. Here the expansion of M constraints are omitted for space limit. The superscript i of an action a_j^i , if existing, represents the index that a_j^i is combined in a_j . To resolve the encoding, the REACH and DEADLOCK constraints enforce the true values of c_{s_3} and c_{r_0} , and the false value of c_{r_4} , which imply the only possible matches in the execution: $\{r_4, s_3^2\}$ and $\{r_0, s_3^1\}$.

10 EXPERIMENTS

The runtime trace for our approach is observed through code instrumentation. The MPICH library is used for the actual runtime [6]. Two CTPs for each benchmark are used in the experiments: the direct translation of an observed MPI execution, A_0 , and the compressed CTP of A_0 by Definition 6, A_f . The translation and compression are largely trivial and proceeds as expected. Anything outside message passing as defined in this paper is ignored. Some care must be taken with collective operations but it is all mechanical. Any deadlock cycle is checked once it is detected. Once a cycle verified to be feasible, the test terminates.

The SMT solver Z3 is used by our approach for validation [7]. The experiments are run on an Intel i7-8700K 6-Core

processor with 24 GB of memory running Ubuntu 18.04 LTS. A time limit of one hour is set for each test.

10.1 Summary of Previous Results

The experiments stated in [1] compare the performance of our approach for analyzing A_0 with three state-of-the-art MPI verifiers MOPPER, a SAT based tool [8], [3], ISP, a dynamic analyzer [9], [10], and Aislinn, another dynamic analyzer [11]. The experiments are launched for seven benchmark programs with various settings of processes and buffering, including *Monte* [12], *Integrate* [13], *Diffusion 2D* [13], *Floyd* [14], *GE* [14], *Heat* [15] and *IS* [16], that exhibit multiple communication patterns with complicated message races.

The direct comparison of time cost to the three tools demonstrates that our approach for analyzing A_0 is much more efficient. Overall, our approach is able to finish most tests under a second. ISP and Aislinn both have unexpectedly high time costs for the tests, and even time out for some tests. MOPPER is much quicker than the other two tools. The optimized MOPPER is even more efficient for large benchmarks, but still suffers from the scalability problem for a few tests (e.g., *Monte* with 16 processes).

In the experiment of this paper, the results are extended to

10.2 Effectiveness of Action Combination

As discussed earlier, the prior analysis may suffer from an unexpected long time cost of cycle searching, and thus is inefficient if the count of processes and actions increases to a certain number, e.g., 128 processes for *Integrate*.

Thus, the experiments in this paper test the effectiveness of action combination by running the complete cycle detection independently on the CTPs A_f and A_0 . For that purpose, the benchmark size is largely expanded from those in [1] so as to differentiate between the analyses with and without action combination. Fig. 10 show the comparison of A_f and A_0 for the number of nodes and edges in the generated graphs and the corresponding time costs of complete cycle detection for some typical benchmarks.

The growths of time costs in Fig. 10 illustrates that for the benchmarks *Integrate*, *Diffusion 2D* and *Floyd*, the CTP A_f has an evident speedup on finishing the cycle detection over the CTP A_0 . Among those benchmarks, *Integrate* has a long sequence of wildcard receives in the first process, which introduces many edges for message race. The action combination largely reduces the count of edges and so highly speeds the cycle detection, especially for the test of 128 processes.

Diffusion 2D has several sequences of wildcard receives in the first process and multiple pairs of consecutive send actions in each process, so the reduction of time cost is evident. Note that since the numbers of nodes and edges for *Diffusion 2D* are extremely large for over 64 processes, the time cost for zero buffer is reasonable. The time cost for infinite buffer, however, is under half a second for all the processes, because the nodes do not form valid deadlock paths, and our filtering algorithm in Algorithm 1 can effectively prune a large number of branches in depth first search.

Floyd has multiple consecutive wildcard receives in each process, whose count is the number of processes. The speedup for cycle detection is more obvious with the growth of #process, as more actions are combined.

Heat also has multiple actions that can be combined, but each combination can only happen between two actions, which does not merge many edges in the graph. Therefore, the improvement of cycle detection is not as evident as the tests of other benchmarks.

Overall, the action combination is able to largely help the analysis in deadlock detection. The level of action combination (how many actions can be combined to a single action) is more important than the count of combination (how many new actions are created) for the analysis to be efficient and scalable.

11 RELATED WORK

Much of the literature on predictive program analysis focuses on detecting as many errors as possible from a single observed execution [17], [18]. As mentioned in Section 1, the analysis presented here is maximally predictive for single-path programs but cannot reason about messages that are issued in a schedule-dependent manner. Instead, this paper focuses on improving the efficiency of an SMT based analysis for single-path programs while leaving extensions to more general programs as future work.

The approach in this paper is inspired by several works. Trace compression is usually important because thread traces can be very large in size. Kini et al. proposed a data race detector that uses trace compression [19]. Unlike our approach that simply collapses the consecutive sends or receives in an identical process, the work represents a trace as a special program written by context-free grammars, which achieves a significant compression.

Joshi et al. proposed a method for multi-threaded Java programs by first detecting potential lock dependency cycles with an imprecise dynamic analyzer and then finding real deadlocks by a random thread scheduler with high probability [20]. The refining strategy of the work inspires the staged approach taken in this paper, but the tool presented in this paper is guaranteed to report a deadlock if it exists in the single-path program.

Sherlock is a tool that uses concolic execution for deadlock detection in Java programs [17]. The key idea is similar to our approach: finding potential deadlocks, and then searching for a feasible schedule that leads to the deadlock. The difference is that Sherlock repeatedly finds alternate schedules through solving constraints that describe new permutations of previously observed schedules rather than leveraging an abstract machine to filter out false deadlocks.

A precise SMT encoding technique was proposed by Huang et al. for verifying properties over MCAPI programs containing message race [21]. The encoding does not require a precise match set and was extended to checking zero buffer incompatibility for MPI programs [22]. This technique is adapted for validating deadlocks in this paper.

The POE approach is a dynamic partial order reduction solution [23] for MPI program verification [10], [9]. The approach was extended to POE_{MSE}, which first uses a precise happens-before relation to find the potential sends that may

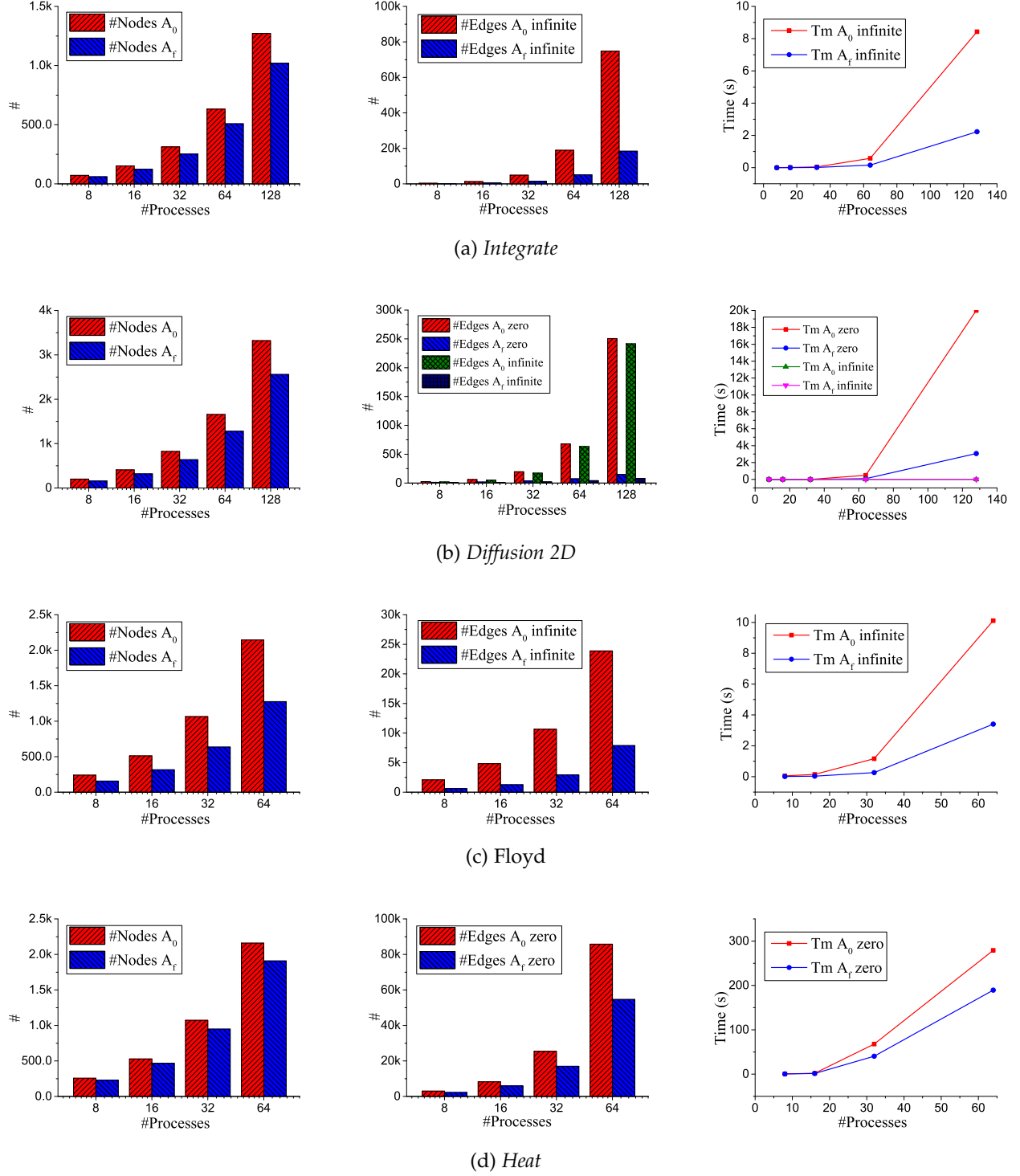


Fig. 10. The #nodes, #edges and time costs of complete cycle detection for the benchmarks (a) *Integrate*, (b) *Diffusion 2D*, (c) *Floyd* and (d) *Heat* under zero and/or infinite buffer settings.

cause different behaviors based on the initial trace, then replays the execution at each potential send with a different choice, i.e. buffering the send instead of matching it [24].

MOPPER is an MPI deadlock detector based on boolean satisfiability encoding [3], [8]. While the solution is precise, the size of the encoding is cubic meaning it only scales to a low degree of message non-determinism.

CIVL is a model checker that uses symbolic execution to verify a number of safety properties of various types of concurrent programs including message passing programs [25], [26], [27]. The tool is outperformed by MOPPER.

An extension to the model checker SPIN [28], is MPI-SPIN that is specific to verifying MPI programs [29], [30]. Since a massive number of states are explored, the work is not scalable.

Böhm et al. provide an approach that aims to find deadlocks for an MPI program under both environments of synchronization and no synchronization [11]. The approach first uses standard partial order reduction to find deadlocks assuming the environment has no synchronization. It then uses an algorithm to search missed deadlocks by enforcing synchronization in the basic operations such as send and collective operations.

MPI-Checker is a static analyzer based on abstract syntax tree of the source code of MPI programs [31]. The tool is able to check many errors in a program, However, it is limited to check deadlocks caused by complicated semantics of communication.

ParTypes is a type-based approach for verifying MPI programs by developing a protocol language for a type system [32]. Since the approach is able to avoid traversing the state space, the analysis is scalable for large programs.

Umpire is an approach of runtime verification for checking multiple MPI errors such as deadlock and resource tracking [33]. The error checking is taken by spawning one manager thread and several outfielder threads in the execution of an MPI program. An extension to Umpire is Marmot [34]. The work uses a centralized server instead of multiple threads for error checking. Another extension to Umpire is MUST [35], [36]. The structure of MUST allows the users to execute the error checking either in an application process itself or in extra processes that are used to offload these analyses. However, just like Umpire and Marmot, the approach is neither sound nor complete for deadlock detection.

12 CONCLUSION

This paper presents a new approach that automatically detects deadlocks in single-path MPI programs after observing a single execution. The actions in the execution, if consecutively sending or consecutively receiving messages in an identical endpoint, can be combined to simplify the presentation of execution, which directly reduces the complication of the mapped dependency graph. The approach leverages a simple characterization of deadlock to efficiently detect deadlock candidates in the dependency graph. An abstract machine is used to quickly disregard many infeasible candidates while the remaining candidates are precisely validated by an SMT solver with an efficient encoding for deadlock. The approach is sound and complete

for deadlock detection in any single-path MPI program on a given input. Experiments show that the new approach with combined actions performs much more efficient for typical benchmarks, comparing to the approach without combination of actions, and the other state-of-the-art MPI verifiers. Future work considers more filtering techniques and extending the approach to support multiple-path MPI programs with more complicated structures.

13 ACKNOWLEDGEMENTS

This work was supported in part by the National Natural Science Foundation of China (NSFC) under Grant 62102325, and in part by the Ministry of Education of Humanities and Social Science Project with Grant 20YJC630146.

REFERENCES

- [1] Y. Huang, B. Ogles, and E. Mercer, "A predictive analysis for detecting deadlock in MPI programs," in *35th IEEE/ACM International Conference on Automated Software Engineering, ASE 2020, Melbourne, Australia, September 21-25, 2020*. IEEE, 2020, pp. 18–28. [Online]. Available: <https://doi.org/10.1145/3324884.3416588>
- [2] The full version of paper. [Online]. Available: <https://github.com/yuhuang/MPIDeadlockExtFullPaper>
- [3] V. Forejt, S. Joshi, D. Kroening, G. Narayanaswamy, and S. Sharma, "Precise predictive analysis for discovering communication deadlocks in MPI programs," *ACM Trans. Program. Lang. Syst.*, vol. 39, no. 4, pp. 15:1–15:27, 2017. [Online]. Available: <http://doi.acm.org/10.1145/3095075>
- [4] Y. Huang, K. Gong, and E. Mercer, "An efficient algorithm for match pair approximation in message passing," *Parallel Computing*, vol. 91, p. 102585, 2020.
- [5] D. B. Johnson, "Finding all the elementary circuits of a directed graph," *SIAM J. Comput.*, vol. 4, no. 1, pp. 77–84, 1975. [Online]. Available: <http://dx.doi.org/10.1137/0204007>
- [6] MPICH, "High-Performance Portable MPI," <http://www.mpich.org>.
- [7] L. de Moura and N. Bjørner, "Z3: An efficient SMT solver," in *TACAS*, vol. 4963. Springer, Heidelberg, 2008, pp. 337–340.
- [8] V. Forejt, D. Kroening, G. Narayanaswamy, and S. Sharma, "Precise predictive analysis for discovering communication deadlocks in MPI programs," in *FM 2014: Formal Methods - 19th International Symposium, Singapore, May 12-16, 2014. Proceedings*, 2014, pp. 263–278. [Online]. Available: http://dx.doi.org/10.1007/978-3-319-06410-9_19
- [9] S. Sharma, G. Gopalakrishnan, and G. Bronevetsky, "A sound reduction of persistent-sets for deadlock detection in MPI applications," in *Formal Methods: Foundations and Applications - 15th Brazilian Symposium, SBMF 2012, Natal, Brazil, September 23-28, 2012. Proceedings*, 2012, pp. 194–209. [Online]. Available: http://dx.doi.org/10.1007/978-3-642-33296-8_15
- [10] S. S. Vakkalanka, S. Sharma, G. Gopalakrishnan, and R. M. Kirby, "ISP: a tool for model checking MPI programs," in *PPOPP*, 2008, pp. 285–286.
- [11] S. Böhm, O. Meca, and P. Jancar, "State-space reduction of non-deterministically synchronizing systems applicable to deadlock detection in MPI," in *FM 2016: Formal Methods - 21st International Symposium, Limassol, Cyprus, November 9-11, 2016. Proceedings*, 2016, pp. 102–118. [Online]. Available: https://doi.org/10.1007/978-3-319-48989-6_7
- [12] J. Burkardt, "MPI Examples." [Online]. Available: http://people.sc.fsu.edu/~jburkardt/cpp_src/mmpi/mmpi.html
- [13] "FEVS benchmark," <http://osl.cis.udel.edu/fevs/index.html>.
- [14] R. Xue, X. Liu, M. Wu, Z. Guo, W. Chen, W. Zheng, Z. Zhang, and G. M. Voelker, "MPIWiz: subgroup reproducible replay of mpi applications," in *Proceedings of the 14th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming, PPOPP 2009, Raleigh, NC, USA, February 14-18, 2009*, 2009, pp. 251–260.
- [15] M. S. Mueller, G. Gopalakrishnan, B. R. de Supinski, D. Lecomber, and T. Hilbrich, "Dealing with MPI Bugs at Scale: Best Practices," in *Automatic Detection, Debugging, and Formal Verification*, 2011.

- [16] D. Bailey, E. Barszcz, B. J.T. B. D.S., C. R.L., D. D., F. R.A., P. Frederickson, L. T.A., R. Schreiber, H. Simon, V. Venkatakrishnan, and W. K., "The nas parallel benchmarks," *International Journal of High Performance Computing Applications*, vol. 5, pp. 63–73, 09 1991.
- [17] M. Eslamimehr and J. Palsberg, "Sherlock: scalable deadlock detection for concurrent programs," in *Proceedings of the 22nd ACM SIGSOFT International Symposium on Foundations of Software Engineering, (FSE-22), Hong Kong, China, November 16 - 22, 2014, 2014*, pp. 353–365. [Online]. Available: <http://doi.acm.org/10.1145/2635868.2635918>
- [18] C. G. Kalhauge and J. Palsberg, "Sound deadlock prediction," *Proceedings of the ACM on Programming Languages*, vol. 2, no. OOPSLA, pp. 1–29, 2018.
- [19] D. Kini, U. Mathur, and M. Viswanathan, "Data race detection on compressed traces," in *Proceedings of the 2018 ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering, ESEC/SIGSOFT FSE 2018, Lake Buena Vista, FL, USA, November 04-09, 2018*, G. T. Leavens, A. Garcia, and C. S. Pasareanu, Eds. ACM, 2018, pp. 26–37. [Online]. Available: <https://doi.org/10.1145/3236024.3236025>
- [20] P. Joshi, C.-S. Park, K. Sen, and M. Naik, "A randomized dynamic program analysis technique for detecting real deadlocks," in *PLDI, 2009*, pp. 110–120.
- [21] Y. Huang, E. Mercer, and J. McCarthy, "Proving MCAPi executions are correct using SMT," in *ASE, 2013*, pp. 26–36.
- [22] Y. Huang and E. Mercer, "Detecting MPI zero buffer incompatibility by SMT encoding," *NFM*, 2015.
- [23] C. Flanagan and P. Godefroid, "Dynamic partial-order reduction for model checking software," in *POPL, 2005*, pp. 110–121.
- [24] S. S. Vakkalanka, A. Vo, G. Gopalakrishnan, and R. M. Kirby, "Precise dynamic analysis for slack elasticity: Adding buffering without adding bugs," in *Recent Advances in the Message Passing Interface - 17th European MPI Users' Group Meeting, EuroMPI 2010, Stuttgart, Germany, September 12-15, 2010. Proceedings, 2010*, pp. 152–159. [Online]. Available: https://doi.org/10.1007/978-3-642-15646-5_16
- [25] S. F. Siegel, M. B. Dwyer, G. Gopalakrishnan, Z. Luo, Z. Rakamaric, R. Thakur, M. Zheng, and T. K. Zirkel, "Civl: The concurrency intermediate verification language," Department of Computer and Information Sciences, University of Delaware, Tech. Rep. UD-CIS-2014/001, 2014.
- [26] M. Zheng, M. S. Rogers, Z. Luo, M. B. Dwyer, and S. F. Siegel, "CIVL: formal verification of parallel programs," in *30th IEEE/ACM International Conference on Automated Software Engineering, ASE 2015, Lincoln, NE, USA, November 9-13, 2015, 2015*, pp. 830–835. [Online]. Available: <https://doi.org/10.1109/ASE.2015.99>
- [27] S. F. Siegel, M. Zheng, Z. Luo, T. K. Zirkel, A. V. Marianello, J. G. Edenhofner, M. B. Dwyer, and M. S. Rogers, "CIVL: the concurrency intermediate verification language," in *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis, SC 2015, Austin, TX, USA, November 15-20, 2015, 2015*, pp. 61:1–61:12. [Online]. Available: <http://doi.acm.org/10.1145/2807591.2807635>
- [28] G. J. Holzmann, "The model checker SPIN," *IEEE Trans. Software Eng.*, vol. 23, no. 5, pp. 279–295, 1997.
- [29] S. F. Siegel, "Model checking nonblocking MPI programs," in *VMCAI, 2007*, pp. 44–58.
- [30] S. Siegel, "Verifying parallel programs with MPI-Spin," in *PVM/MPI, 2007*, pp. 13–14.
- [31] A. Droste, M. Kuhn, and T. Ludwig, "Mpi-checker: static analysis for MPI," in *Proceedings of the Second Workshop on the LLVM Compiler Infrastructure in HPC, LLVM 2015, Austin, Texas, USA, November 15, 2015, 2015*, pp. 3:1–3:10. [Online]. Available: <http://doi.acm.org/10.1145/2833157.2833159>
- [32] H. A. López, E. R. B. Marques, F. Martins, N. Ng, C. Santos, V. T. Vasconcelos, and N. Yoshida, "Protocol-based verification of message-passing parallel programs," in *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, 2015*, pp. 280–298. [Online]. Available: <http://doi.acm.org/10.1145/2814270.2814302>
- [33] J. S. Vetter and B. R. de Supinski, "Dynamic software testing of MPI applications with umpire," in *Proceedings Supercomputing 2000, November 4-10, 2000, Dallas, Texas, USA. IEEE Computer Society, CD-ROM, 2000*, p. 51. [Online]. Available: <http://doi.ieeecomputersociety.org/10.1109/SC.2000.10055>
- [34] B. Krammer, K. Bidmon, M. S. Müller, and M. M. Resch, "MAR-MOT: an MPI analysis and checking tool," in *Parallel Computing: Software Technology, Algorithms, Architectures and Applications, PARCO 2003, Dresden, Germany, 2003*, pp. 493–500.
- [35] T. Hilbrich, M. Schulz, B. R. de Supinski, and M. S. Müller, "MUST: A scalable approach to runtime error detection in MPI programs," in *Tools for High Performance Computing 2009 - Proceedings of the 3rd International Workshop on Parallel Tools for High Performance Computing, September 2009, ZIH, Dresden, 2009*, pp. 53–66. [Online]. Available: http://dx.doi.org/10.1007/978-3-642-11261-4_5
- [36] T. Hilbrich, J. Protze, M. Schulz, B. R. de Supinski, and M. S. Müller, "MPI runtime error detection with MUST: advances in deadlock detection," in *SC Conference on High Performance Computing Networking, Storage and Analysis, SC '12, Salt Lake City, UT, USA - November 11 - 15, 2012, 2012*, p. 30. [Online]. Available: <https://doi.org/10.1109/SC.2012.79>

14 APPENDIX

14.1 Proof of Theorem 1

First we prove a few useful lemmas. Fix a reachable deadlock $q = \langle A, I, M \rangle$ and its full control point $D = Ctrl(q)$.

Lemma 1. [deadlock-actions-blocking]

$$D \subseteq W(A) \cup B(A)$$

Proof: Proof by contradiction. Assume that $\exists a \in D, a \notin W(A) \cup B(A)$.

Then $a \in S(A) \cup R(A)$. Let a' be the very next action after a in $A_{pid(a)}$. a' must exist since a is a send or receive action and the last action in every process is a barrier action.

By definition D contains the last issued action of every process in $P(A)$. It follows that $a \in I$ and $a' \notin I$. Furthermore, by the definition of deadlock, there is no transition enabled in q that can issue q' .

The definition of \rightarrow requires that every action preceding a in $A_{pid(a)}$ must be issued and every blocking action preceding a must be matched in q . Because $a \in I$ and $a \in S(A) \cup R(A)$, these conditions also hold for a' . Therefore a' can be issued. This contradicts the deadlock of q . \square

Lemma 2. [wildcard-causes-send-starvation] Let $a \in D \cap W(A)$ be a wait action participating in the deadlock at q such that $req(a) \in S(A)$. If there are no receive actions in $(R(A) \setminus I)_{dst(req(a))}$ (no unissued potential matches), then there must exist an issued wildcard receive action $a' \in (R(A) \cap I)_{dst(req(a))}$, $src(a') = *$ in the destination process.

Proof: Proof by contradiction. Assume that there is no receive $a' \in R(A) \cap I$ such that $pid(a') = dst(req(a))$ and $src(a') = *$. Then the number of send and receive match pairs between $src(req(a))$ and $dst(req(a))$ is fixed in every execution of A . Since q is deadlocked and there are no unissued potential matches for $req(a)$, it follows that $\sum_{a_s \in (S(A) \cap I)_{src(req(a))}} n(a_s) > \sum_{a_r \in (R(A) \setminus I)_{dst(req(a))}} n(a_r)$. Therefore, A will deadlock deterministically at a . This contradicts the fact that A was obtained by observing a successful execution of an MPI program. \square

Lemma 3. [wildcard-causes-recv-starvation] Let $a \in D \cap W(A)$ be a wait action participating in the deadlock at q such that $req(a) \in R(A) \wedge src(req(a)) \neq *$. If there are no send actions in $(S(A) \setminus I)_{src(req(a))}$ (no unissued potential matches), then there must exist an issued wildcard receive action $a' \in (R(A) \cap I)_{pid(a)}$, $src(a') = *$ in the destination process.

Proof: Proof by contradiction. Assume that there is no receive $a' \in R(A) \cap I$ such that $a' <_{po} req(a)$ and $src(a') = *$. Then the number of send and receive match pairs between $src(req(a))$ and $dst(req(a))$ is fixed at least until a (where after a there may be more receives but no more sends). Since q is deadlocked and there are no unissued potential matches for $req(a)$, it follows that $\sum_{a_r \in (R(A) \cap I)_{dst(req(a))}} n(a_r) > \sum_{a_s \in (S(A) \setminus I)_{src(req(a))}} n(a_s)$. Therefore, A will deadlock deterministically at a . This contradicts the fact that A was obtained by observing a successful execution of an MPI program. \square

Lemma 4. [parent-action-exists]

$$\forall a \in D, (\exists a' \in A \setminus I, pid(a') \neq pid(a) \wedge (a', a) \in E^*)$$

Proof: Proof by case analysis on the action type of a . By Lemma 1, $a \in W(A) \cup B(A)$.

In the first case, $a \in W(A) \wedge req(a) \in S(A)$. If there is an unissued potentially matching receive action a' in $(A \setminus I)_{dst(req(a))}$, then rule two will ensure that $(a', req(a)) \in E$. If there is no such matching receive, then $a' = \perp_{dst(req(a))}$. By Lemma 2, there must be a wildcard receive in $I_{dst(req(a))}$. Therefore, rule four will again ensure that $(a', req(a)) \in E$. In either case, because $req(a) <_{mo} a$, it follows that $(req(a), a) \in E$ and thus $(a', a) \in E^*$.

In the second case, $a \in W(A) \wedge req(a) \in R(A)$. If there is an unissued potentially matching send action a' in $(A \setminus I)_{src(req(a))}$, then rule two will ensure that $(a', req(a)) \in E$. If there is no such matching send, then $a' = \perp_{src(req(a))}$. By Lemma 3, there must be a wildcard receive in $I_{pid(a)}$. Therefore, rule three will again ensure that $(a', req(a)) \in E$. In either case, because $req(a) <_{mo} a$, it follows that $(req(a), a) \in E$ and thus $(a', a) \in E^*$.

In the final case, $a \in B(A)$. Then $a' \in B(A) \setminus I$ and rule two ensures that $(a', a) \in E$. Such an a' must exist, otherwise a could be completed and q would not be a deadlock.

In the first two cases, $req(a)$ is the orphaned action, in the last case, a itself is the orphaned action. \square

Lemma 5. [actions-preceding-deadlock-actions-connected]

$$\forall a \in D, (\exists a' \in D, a_o \in A_{pid(a)}, a'_o \in A_{pid(a')}, a_o \leq_{po} a \wedge a'_o \leq_{po} a' \wedge (a'_o, a_o) \in E^*)$$

Proof: Direct proof. By Lemma 4, there exists some action $a'' \in A \setminus I$ such that $pid(a'') \neq pid(a)$ and $(a'', a) \in E^*$. Let a_o be an action in $pid(a)$ such that $(a'', a_o), (a_o, a) \in E^*$. By definition of deadlock path, a_o must exist.

By definition, D contains an action $a' \in I_{pid(a'')}$. In the first case of deadlock path, it follows that $a' \leq_{po} a''$. If $a' <_{po} a''$, then $a' <_{mo} a''$ since a' is a blocking action by Lemma 1. In this case, rule one ensures that $(a', a'') \in E^*$. Let a'_o be the action $req(a')$ if $a' \in W(A)$, or equal to a' if $a' \in B(A)$. Then, $a'_o \leq_{mo} a'$ by the definition of match order, and rule one ensures $(a'_o, a') \in E^*$. By the transitivity of E^* , $(a'_o, a_o) \in E^*$.

In the second case of deadlock path, it follows that $a'' \leq_{po} a'$. Let a'_o be equal to a'' . Then, $(a'_o, a_o) \in E^*$ is evident.

If $a' = a'' = \perp_{pid(a')}$, then this is a contradiction since $a' \in I$ by the definition of D and $a'' \notin I$ by Lemma 4. \square

Now we can prove Theorem 1. By Lemma 5, there exist some actions a_o and a'_o preceding the actions a and a' in $pid(a)$ and $pid(a')$ respectively such that $a, a' \in D$, $a \neq a'$ and $(a'_o, a_o) \in E^*$. By the same argument, there exists an action a''_o preceding the action a'' in $pid(a'')$ such that $a'' \in D$, $a' \neq a''$ and $(a''_o, a'_o) \in E^*$. If $a_o = a''_o$, then we have a cycle $C = \{(a_o, a'_o), \dots, (a'_o, a_o)\}$. Otherwise, we can continue applying Lemma 5 until a cycle is created between the actions in D .

Let d be the partial control point or candidate extracted from C . Particularly, the action $a_d \in d$ can be the earliest blocking action in C from process $pid(a_d)$ if there are more

than one action for process $pid(a_d)$ in C . The action a_d can also be the very next wait action for a wildcard receive $req(a_d)$ with $n(req(a_d)) > 1$ if the receive is the only action for process $pid(a_d)$ in C . In the second case, a_d is outside C but can be trivially obtained. The actions that make up edges in C are completely enumerated in Lemmas 4 and 5. C is entirely made up of actions in D , their message requests (which are non-blocking), unissued matches and final barrier actions. Therefore, $d \subseteq D$ holds by construction.

14.2 Proof of Theorem 2

Proof by contradiction. Without loss of generality, assume that no match transitions are enabled from q and q' in \rightarrow and \rightarrow_{abs} (match transitions do not modify the control points). Now assume that some issue transition $q \rightarrow q''$ is enabled such that $Ctrl(q'') \notin \mathbb{C}_{\rightarrow_{abs}}^{q'}$. Then, by the definition of $Ctrl(q'')$, there is some action $a \in A$ that can be issued by \rightarrow but not by \rightarrow_{abs} in the states q and q' respectively. Let $q = \langle A, I, M \rangle$ and $q' = \langle A', I', M' \rangle$.

A few useful facts follow from the definitions of the two transition systems and the assumptions above.

- 1) $A \subseteq A'$ and $A' \setminus A$ consists solely of all wildcard send actions generated by the abstract machine when issuing send actions.
- 2) $I \subseteq I'$ and $I' \setminus I$ consists solely of all wildcard send actions generated by the abstract machine when issuing send actions.
- 3) There exists some earliest action $a' <_{po}^{A'} a$ that is preventing a from being issued by \rightarrow_{abs} in the state q' because it is not issued or matched as required by the issue transition rule.
- 4) a' did not prevent \rightarrow from issuing a in state q and therefore a' must be issued and/or matched as required in that state.

Fact three leads to a few cases. In the first case, $a' \notin I'$, thereby blocking a from being issued in the state q' .

If $a' \notin A$, then it is one of the wildcard send actions generated by the abstract machine. But by the definition of \rightarrow_{abs} , every wildcard send in A' is also in I' which is a contradiction.

So it must be that $a' \in A$ is not a wildcard send action. But then facts two and four imply that $a' \in I'$, again a contradiction.

In the second case, $a' \in W(A') \cup B(A') \wedge a' \notin M'$. If $a' \in B(A')$, then facts two and four imply that $B_{A'}(grp(a')) \subseteq I'$. And because there are no match transitions enabled from state q' in \rightarrow_{abs} , it follows that $B_{A'}(grp(a')) \subseteq M'$ which is a contradiction.

If a' is not a barrier action, then $a' \in W(A')$ must be an incomplete wait action. And because there are no match transitions enabled from state q' in \rightarrow_{abs} , it follows that $req(a')$ has no potential match in I' .

In contrast, fact four implies that $req(a')$ has matched with some action in the state q . Furthermore, fact two guarantees that the same number of send and receive actions from A have been issued between $src(req(a'))$ and $dst(req(a'))$ in the states q and q' . And because the abstract machine issues a wildcard send for every deterministic send, there must exist an enabled match for $req(a')$ in I' . This contradicts a previous assumption.

Because all cases lead to a contradiction, it must be that $Ctrl(q'') \in \mathbb{C}_{\rightarrow_{abs}}^{q'}$. \square

14.3 Proof of Theorem 3

Proof by contradiction. Let $t = q_0 \rightarrow_{abs} q_1 \dots \rightarrow_{abs} q$ be a transition sequence in $\Sigma_{\rightarrow_{abs}}^{q_0}$ that reaches the deadlock state q . Now assume that there exists another state $q' \in \Sigma_{\rightarrow_{abs}}^{q_0}$ such that $Dead_{\rightarrow_{abs}}(q')$ and $D \subseteq Ctrl(q')$.

Then there must be some earliest q_i in t such that $q_i \rightarrow_{abs}^* q'$ and $q_{i+1} \not\rightarrow_{abs}^* q'$. In other words a scheduling choice causes the full control point subsuming D to become reachable. Because the abstract machine isolates wildcard actions on their own endpoints, the only scheduling choice that can affect which blocking actions are able to be matched is the choice of which wildcard send actions are matched with wildcard receive actions.

Therefore the transition $q_i \rightarrow_{abs} q_{i+1}$ must have matched a wildcard send action s with a wildcard receive action r and there must be some other send action s' that can be matched with r along the sequence $q_i \rightarrow_{abs}^* q'_{i+1}$. Furthermore, the matching of s and r must have prevented some blocking action in a process with an action in D to block indefinitely. However the definition of the **Issue-Send** transition in \rightarrow_{abs} , s and s' can be matched with the exact same set of actions and are entirely interchangeable. Therefore, if $q'_{i+1} \rightarrow_{abs}^* q'$, then $q_{i+1} \rightarrow_{abs}^* q'$. This is a contradiction.