

## Efficient Verification Solutions for Message Passing Systems

Subodh Sharma and Ganesh Gopalakrishnan

*School of Computing, University of Utah, Salt Lake City, UT*

*Email: {svs,ganesh}@cs.utah.edu*

**Abstract**—We examine the problem of automatically and efficiently verifying the absence of communication related bugs in message passing systems, specifically in programs written using Message Passing Interface (MPI) API. A typical debugging or testing tool will fail to achieve this goal because they do not provide any guarantee of coverage of non-deterministic communication matches in a message passing program. While dynamic verification tools do provide such a guarantee, they are quickly rendered useless when an interleaving explosion is witnessed. The general problem is difficult to solve, though we propose that specialized techniques can be developed that can work on top of dynamic verification schedulers thus making them more efficient.

In this work, we provide point solutions to deal with the interleaving explosion. Specifically, we present algorithms that accomplish the following tasks: (i) identifying irrelevant message passing operations (Barriers) in MPI programs that add to the verification complexity and degrade application's performance and (ii) reducing substantially the *relevant* set of interleavings using symmetry patterns; that needs to be explored for the detection of *refusal deadlocks* in MPI programs.

**Keywords**-MPI; Dynamic Verification; Partial Order Reduction

### I. INTRODUCTION

Message passing systems are ubiquitous with their presence stretching across computing domains - starting from high performance computing all the way to embedded/heterogeneous multicore computing. Most of today's supercomputers and clusters are running programs that are written using MPI [1]. Although use of libraries like MPI provide significant performance boost, they are also a source of worry. The extensive use of non-deterministic communication primitives (e.g., MPI receives with `MPI_ANY_SRC`), are the cause of hard to reproduce bugs (*Hiesenburgs*).

While there have been significant advances in the debugging and testing methods [2]–[4], they lack the fine control necessary to explore different schedules

that arise due to the non-determinism in communication matching.

Model checking tools can provide guaranteed coverage (MPI-SPIN [5]), however constructing models for these codes is often laborious and bug-prone. Recently created *formal dynamic verifiers* such as ISP [6], [7] and DAMPI [8] take an approach that integrates the best features of testing tools (ability to run directly on user applications) and model checking (coverage guarantees). They run the MPI program under the control of a verification scheduler, guarantee to detect all potential matches for non-deterministic (wildcard) receives, and explore each of these matches in different runs of the program. Thus, *they exhaustively explore and ensure full coverage of non-determinism*.

With such unrestrained coverage of non-determinism, the dynamic verifiers like ISP and DAMPI will have to grapple with the inevitable *interleaving explosion*. For instance, consider an MPI program with  $n + 1$  processes where each of the  $n$  processes sends a message to the  $(n + 1)^{th}$  process. The  $(n + 1)^{th}$  process posts  $n$  wildcard receive calls (say, in a loop). One can easily observe that even in such a simple setting, there will be  $n!$  execution schedules. This is clearly unacceptable: all dynamic verifiers must, ideally, be equipped with approaches to detect when such exhaustive explorations are unnecessary, and then avoid them.

**Problem Statement:** The problem of pruning the interleaving space with multiple identical processes with *symmetric reasoning* is formally undecidable [9]. We do not attempt to provide another approximating solution to this problem. However, we do propose a specialized dynamic analysis method that will substantially reduce the number of interleavings while detecting *communication deadlocks*. We have implemented our initial work in this problem domain on top ISP. Our method is an improvement of ISP's algorithm called POE (partial order avoiding elusive interleavings). We

call our method as MSPOE (macroscopic POE).

As a part of a larger effort of interleaving space reduction, we also propose another specialized algorithm to identify *irrelevant* synchronizing MPI operations and remove them. Irrelevant MPI operations are those operations which when removed, do not alter the communication structure of the program. There are two-fold reasons for such operations to be removed: (i) they unnecessarily increase the number of MPI operations in the program, thus indirectly affecting the number of schedules a program can have; (ii) they increase the application running time and consequently increasing the verification time of these applications.

Rest of the paper is presented in the following fashion: Section II will briefly present our proposed solution for detecting *functionally irrelevant barriers* (FIB) in an MPI program with results followed by Section III where we briefly discuss our MSPOE algorithm. Finally, we conclude with remarks on future work in Section IV.

## II. FUNCTIONALLY IRRELEVANT BARRIERS

The importance of detecting irrelevant barriers comes from a number of perspectives. Many MPI users are known to employ collective barriers for “good measure”; they are unsure whether it is necessary. The authors of [10] present an example where a barrier was considered irrelevant and was proposed to be removed. A year later, they realized that its removal introduces a non-benign race condition. In [11], it is shown that barriers can consume a significant fraction of the total application time. Of course, users wanting to control performance by avoiding network or I/O contention may insert collective barriers. In this case, they are employing functionally irrelevant barriers for controlling the non-functional aspects of their program.

The FIB algorithm [12] can help these users by checking that these barriers are indeed functionally irrelevant.

**Motivating Example:** Consider the example shown below.

```
P0: Irecv(*, &handle); Barrier; Wait(&handle);
P1: Isend(to P0); Barrier; ...rest of P1...
P2: ...some code... Barrier; Isend(to P0);
```

Notice *Irecv* is the asynchronous MPI receive and *Wait* call is a blocking call to check the successful completion of the corresponding non-blocking request.

Observe that there is no *Happens-Before* (HB) ordering between *Isends* of *P1* and *P2*. The *Irecv* and *Isend* from *P0* and *P1* can exist past barriers. Thus, removing the barriers will not create any new communication matchings for the *Irecv* and therefore they are irrelevant.

Now consider an alternative example in which the *Wait* in *P0* is moved to be *before* its *Barrier*. Now, the collective barrier becomes **relevant**. This is because there would be a HB edge from *Wait* to *Barrier* as shown in Figure 1. Hence, *Barrier* cannot be crossed until the *Irecv* finishes. Therefore the *Isend* from *P2* cannot issue, and *Irecv* must finish based on the *Isend* from *P1*.

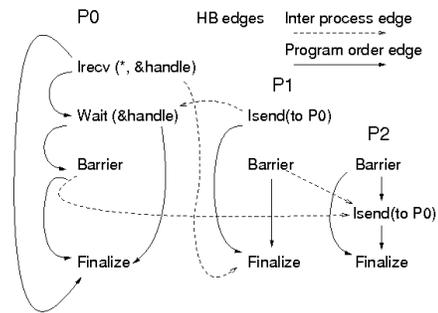


Figure 1. Example in Section II with Happens-Before edges

**Results:** The algorithm is implemented on top of ISP. Our web page [13] provides detailed results regarding FIB. To briefly summarize :

- **Monte-Carlo:** The code of Monte-Carlo, did not have any barrier calls. To acid-test our implementation, we deliberately inserted an irrelevant collective barrier, which our implementation flagged as such. The run times of the Fib algorithm are as follows: with 5 processes, it explored 24 interleavings in 1.52 seconds. The overhead due to FIB analysis was marginal (< 1%).
- **2D Diffusion** This code had 2 irrelevant barriers which were caught by the tool. In fact, this example does not employ wildcard receives, and so all its barriers are irrelevant.
- **Umpire test suite:** We ran our tool successfully on all the 69 tests that came along with Umpire tool [14]. Of the 36 tests that had barriers, all were flagged as **irrelevant**, with negligible runtimes.

## III. SYMMETRIC MSPOE

Our method MSPOE [15] is implemented by augmenting the ISP tool and its POE. We first let POE

Benchmark	# of procs	Deadlocks?	Interleavings		Time(sec)	
			POE	MSPOE	POE	MSPOE
Matrix Multiply	4	No	18	1	2.93	.16
Red-2D-Diffusion	4	No	90	1	29.4	.33
2D-Diffusion	4	No	>5000	1	>3600	3.34
Monte Carlo	6	No	120	1	24.86	.005
Integrate	4	No	81	31	19.46	7.4
Madre	4	No	>8000	1	>3600	.77

Table I

compute the potential send matches for MPI non-deterministic receives as it currently does. The execution history, following the non-deterministic receive, is then examined by MSPOE. It chooses to include only some of these sends (called *relevant sends*) to match this non-deterministic receive for later explorations. These sends are the ones considered relevant to cause refusal deadlocks. In effect, instead of exploring all executions MSPOE explores *representative executions* sufficient to reveal refusal deadlock.

**Observation:** For an MPI program that does not decode data and has a refusal deadlock, it must either have an unequal number of sends and receives in some execution path, or must satisfy the following conditions: (i) it employs a process posting a wildcard receive and a specific receive and (ii) a previous wildcard receive consumes a send that was meant for the later occurring specific receive, thus orphaning the specific receive. MSPOE exploits this observation and computes relevant sends based on the occurrence of specific receives.

**Motivating Example:** In the example shown below there is a deadlock introduced by the use of the deterministic receive call.

P0:  $S_{0,1}$ (to P4); | P1:  $S_{1,1}$ (to P4);  
P2:  $S_{2,1}$ (to P4); | P3:  $S_{3,1}$ (to P4);  
P4:  $R_{4,1}$ (\*);  $R_{4,2}$ (from P3);  $R_{4,3}$ (\*);  $R_{4,4}$ (\*);

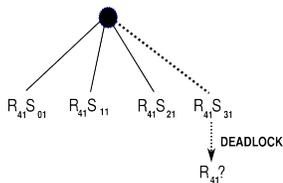


Figure 2. Possibilities after first  $R(*)$  match

Figure 2 shows that if  $R_{4,1}$  were to match  $S_{3,1}$  (rightmost transition from the initial node), the subsequent deterministic call ( $R_{4,2}$ ) will be orphaned, thus creating a refusal deadlock. ISP would explore all the matches starting from leftmost choice shown in Figure 2 and then moving right with every new run, generating four interleavings before finding the deadlock. MSPOE will, on the other hand, choose  $S_{3,1}$  as the next relevant send to explore after any initial run. This guarantees that the deadlock will be detected in at most two interleavings.

**Results:** We implemented the MSPOE algorithm on top of ISP. All experiments were run on Intel Core Duo (2 Ghz) with 3 GB RAM. Table I on page 3 summarizes MSPOE results.

#### IV. CONCLUSIONS AND FUTURE WORK

We have presented our initial proposal to combat non-determinism while verifying message passing systems. Our initial results are inspiring. We obtained dramatic interleaving space reductions using our MSPOE algorithm; in some cases, from million interleavings to one. Our FIB method successfully detected multiple irrelevant barriers in several MPI benchmarks.

**Future Work:** We will extend our MSPOE algorithm to handle cases which have an obvious received data decoding that controls the ensuing communication flow of the program. To handle data decoding, we would require an MPI specific control flow graph (CFG) of the program. In [16], a CFG for MPI programs (*p-cfg*) is presented. However, *p-cfg* can handle only deterministic MPI programs. We intend to work on improving the *p-cfg* to handle non-deterministic MPI operations. Furthermore, we will develop flow-sensitive static analysis methods on top of the improved *p-cfg* to analyze conditional communication patterns.

## REFERENCES

- [1] “Message Passing Forum.” [Online]. Available: <http://www.mpi-forum.org/docs/>
- [2] B. Krammer, K. Bidmon, M. Miller, and M. Resch, “Marmot: An mpi analysis and checking tool,” in *PCST, Algorithms, Architectures and Applications*, ser. Advances in Parallel Computing. North-Holland, 2004, vol. 13, pp. 493 – 500.
- [3] “TotalView Concurrency Tool.” [Online]. Available: <http://www.totalviewtech.com>
- [4] M. Noeth, P. Ratn, F. Mueller, M. Schulz, and B. R. de Supinski, “Scalatrace: Scalable compression and replay of communication traces for high-performance computing,” *J. Parallel Distrib. Comput.*, vol. 69, pp. 696–710, August 2009.
- [5] S. F. Siegel and G. S. Avrunin, “Verification of mpi-based software for scientific computation,” in *Model Checking Software: 11th International SPIN Workshop*. Springer-Verlag, 2004, pp. 286–303.
- [6] S. Vakkalanka, G. Gopalakrishnan, and R. M. Kirby, “Dynamic verification of MPI programs with reductions in presence of split operations and relaxed orderings,” in *Computer Aided Verification (CAV 2008)*, 2008, pp. 66–79.
- [7] A. Vo, S. Vakkalanka, M. DeLisi, G. Gopalakrishnan, R. M. Kirby, and R. Thakur, “Formal verification of practical mpi programs,” in *Proceedings of the 14th ACM SIGPLAN PPOPP*, ser. PPOPP. New York, NY, USA: ACM, 2009, pp. 261–270.
- [8] A. Vo, S. Aananthakrishnan, G. Gopalakrishnan, B. R. de Supinski, M. Schulz, and G. Bronevetsky, “A scalable and distributed dynamic formal verifier for mpi programs,” *SC Conference*, vol. 0, pp. 1–10, 2010.
- [9] E. M. Clarke, O. Grumberg, and D. A. Peled, *Model Checking*. MIT Press, 2000.
- [10] G. S. Avrunin, S. F. Siegel, and A. R. Siegel, “Finite-state verification for high performance computing,” in *Proceedings of SE-HPCS*. New York, NY, USA: ACM, 2005, pp. 68–72.
- [11] R. Rabenseifner, “Automatic profiling of mpi applications with hardware performance counters,” in *Proceedings of the 6th EuroPVM/MPI*. London, UK: Springer-Verlag, 1999, pp. 35–42.
- [12] S. Sharma, S. S. Vakkalanka, G. Gopalakrishnan, R. M. Kirby, R. Thakur, and W. Gropp, “A formal approach to detect functionally irrelevant barriers in mpi programs,” in *PVM/MPI*, 2008, pp. 265–273.
- [13] S. Sharma, “FIBResults,” 2008. [Online]. Available: [http://www.cs.utah.edu/svs/FIB\\_results/](http://www.cs.utah.edu/svs/FIB_results/)
- [14] J. S. Vetter and B. R. de Supinski, “Dynamic software testing of mpi applications with umpire,” in *Proceedings of SC 2000*, ser. Supercomputing ’00. Washington, DC, USA: IEEE Computer Society, 2000.
- [15] S. Sharma and G. Gopalakrishnan, “Scalable analysis for deadlock detection in mpi programs,” *Submitted to ICS*, 2011.
- [16] G. Bronevetsky, “Communication-sensitive static dataflow for parallel message passing applications,” in *CGO*, 2009, pp. 1–12.