

Semantics Driven Dynamic Partial-order Reduction of MPI-based Parallel Programs

Robert Palmer

Intel Validation Research Labs, Hillsboro, OR
(work done at the Univ of Utah as PhD student)

Ganesh Gopalakrishnan

Robert M. Kirby

School of Computing
University of Utah

Supported by:

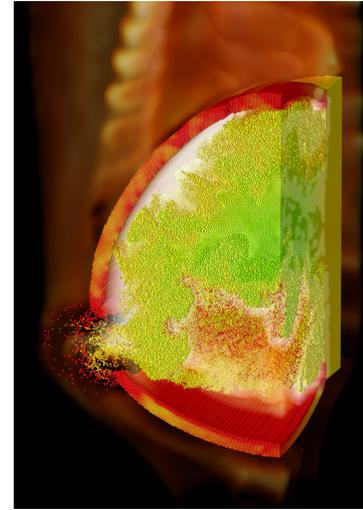
Microsoft HPC Institutes

NSF CNS 0509379

MPI is the de-facto standard for programming cluster machines



(BlueGene/L - Image courtesy of IBM / LLNL)



(Image courtesy of Steve Parker, CSAFE, Utah)

Our focus: Eliminate Concurrency Bugs from HPC Programs !

An Inconvenient Truth: Bugs → More CO₂, Bad Numbers !

So many ways to eliminate MPI bugs ...

- **Inspection**

- Difficult to carry out on MPI programs (low level notation)

- **Simulation Based**

- Run given program with manually selected inputs
- Can give poor coverage in practice

- **Simulation with runtime heuristics to find bugs**

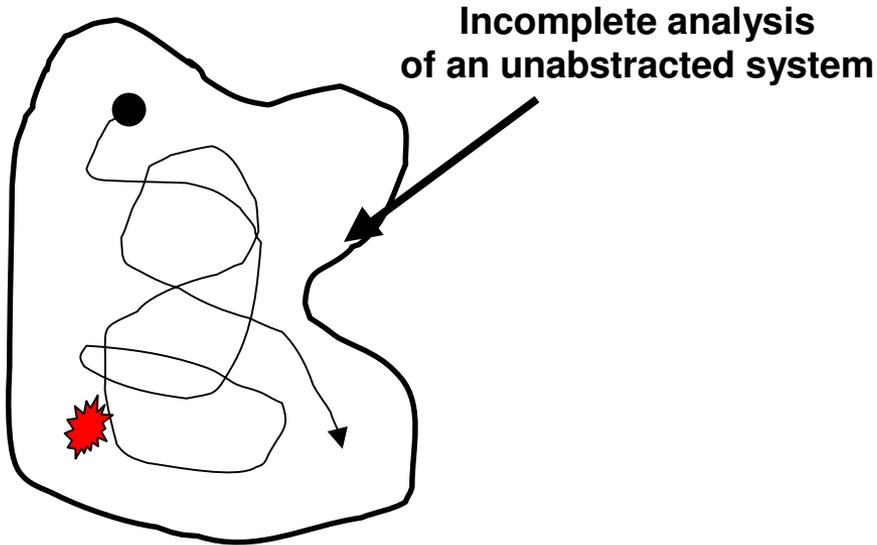
- Marmot: Timeout based deadlocks, random executions
- Intel Trace Collector: Similar checks with data checking
- TotalView: Better trace viewing – still no “model checking”(?)
- We don't know if any formal coverage metrics are offered

- **Model Checking Based**

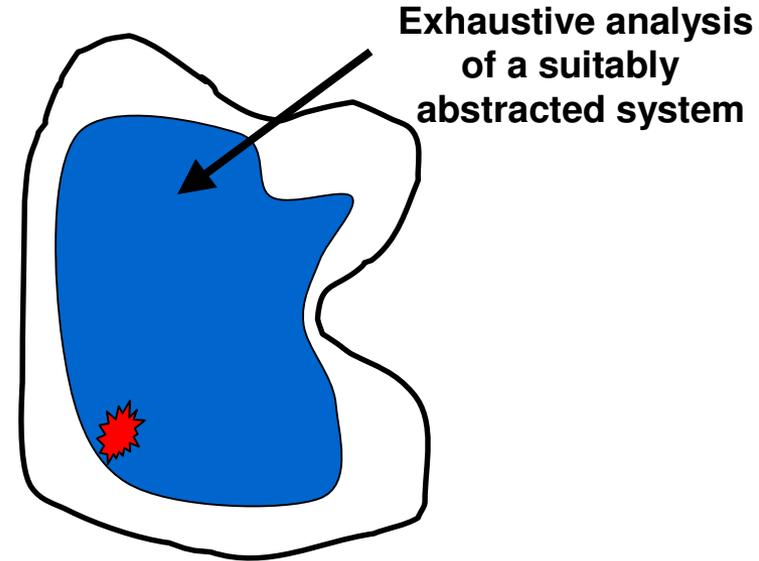
- Being widely used in practice
- Can provide superior debugging for reactive bugs
- Has made considerable strides in abstraction (data, control)

Our Core Technique: Model Checking

“Ad-hoc Testing”



“Model Checking”



Why model checking works in practice:

- * It applies *Exhaustive Analysis*, as opposed to *Incomplete Analysis*
- * It relies on *Abstraction* (both manual, and automated)

Exhaustive analysis of suitably abstracted systems helps catch more bugs than incomplete analysis of unabstracted systems

[Rushby, SRI International]

Model Checking Approaches for MPI

- **MC Based On “Golden” Semantics of MPI**

Limited Subsets of MPI / C Translated to TLA+ (FMICS 2007)

Limited C Front-End with Slicing using Microsoft Phoenix

- **Hand Modeling / Automated Verif. in Executable Lower Level Formal Notations**

Modeling / Verif in Promela (Siegel, Avrunin, et.al. – several papers)

Non-Blocking MPI Operations in Promela + C (Siegel)

Limited Modeling in LOTOS (Pierre et.al. – in the 90’s)

- **Modeling in MPI / C – Automatic Model Extraction**

Limited Conversion to Zing (Palmer et.al. – SoftMC 05)

Limited Conversion to MPIC-IR (Palmer et.al. – FMICS 07)

- **Direct Model Checking of Promela / C programs**

Pervez et.al. using PMPI Instrumentation – EuroPVM / MPI

Demo of One-Sided + a Few MPI Ops (Pervez et.al, EuroPVM / MPI 07)

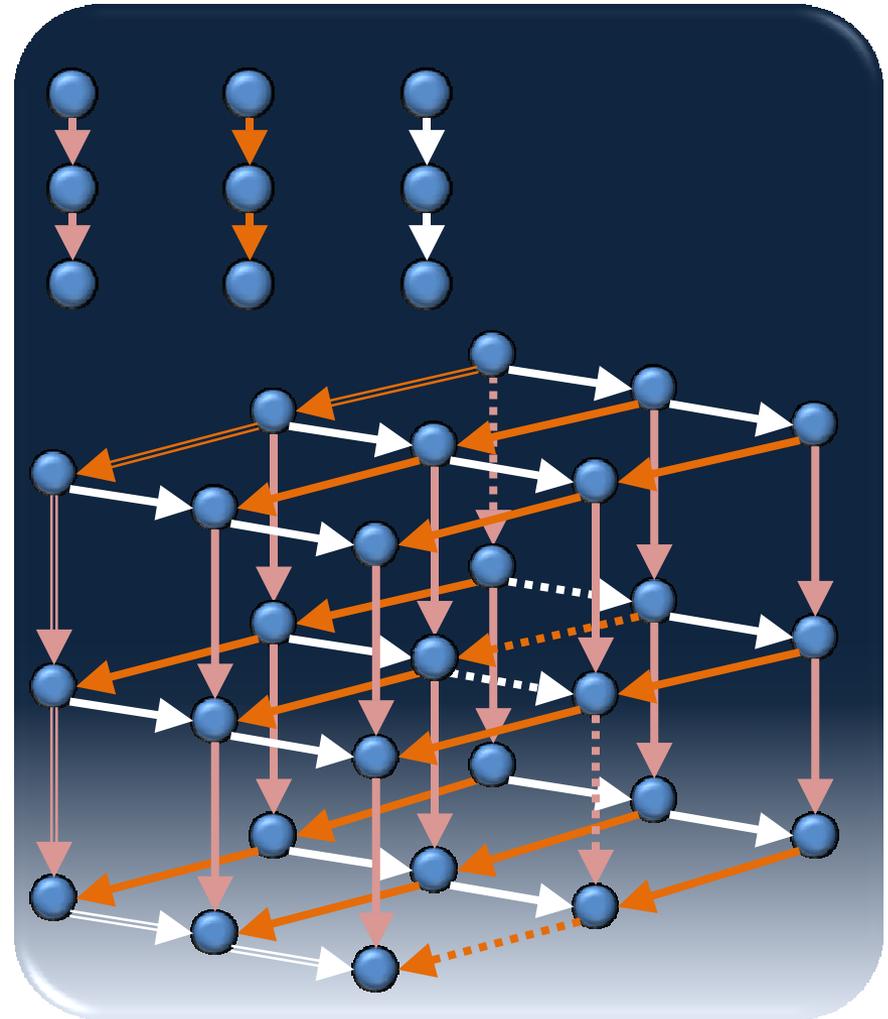
Model Checking Approaches for MPI

- 1. MC Based On “Golden” Semantics of MPI**
 - 1. Limited* Subsets of **MPI / C** Translated to **TLA+** (FMICS 2007)
 - 2. Limited* C Front-End with Slicing using **Microsoft Phoenix**
- 2. Hand Modeling / Automated Verif. in Executable Lower Level Formal Notations**
 1. Modeling / Verif in **Promela** (Siegel, Avrunin, et.al. – several papers)
 2. Non-Blocking MPI Operations in **Promela + C** (Siegel)
 3. *Limited* Modeling in **LOTOS** (Pierre et.al. – in the 90’s)
- 3. Modeling in MPI / C – Automatic Model Extraction**
 1. *Limited* Conversion to **Zing** (Palmer et.al. – SoftMC 05)
 2. *Limited* Conversion to **MPIC-IR** (Palmer et.al. – FMICS 07)
- 4. Direct Model Checking of Promela / C programs**
 1. Pervez et.al. using **PMPI Instrumentation** – EuroPVM / MPI
 2. Demo of **One-Sided + a Few MPI Ops** (Pervez et.al, EuroPVM / MPI 07)

THIS PAPER : Explain new DPOR Idea Underlying **3.2, 4.2**

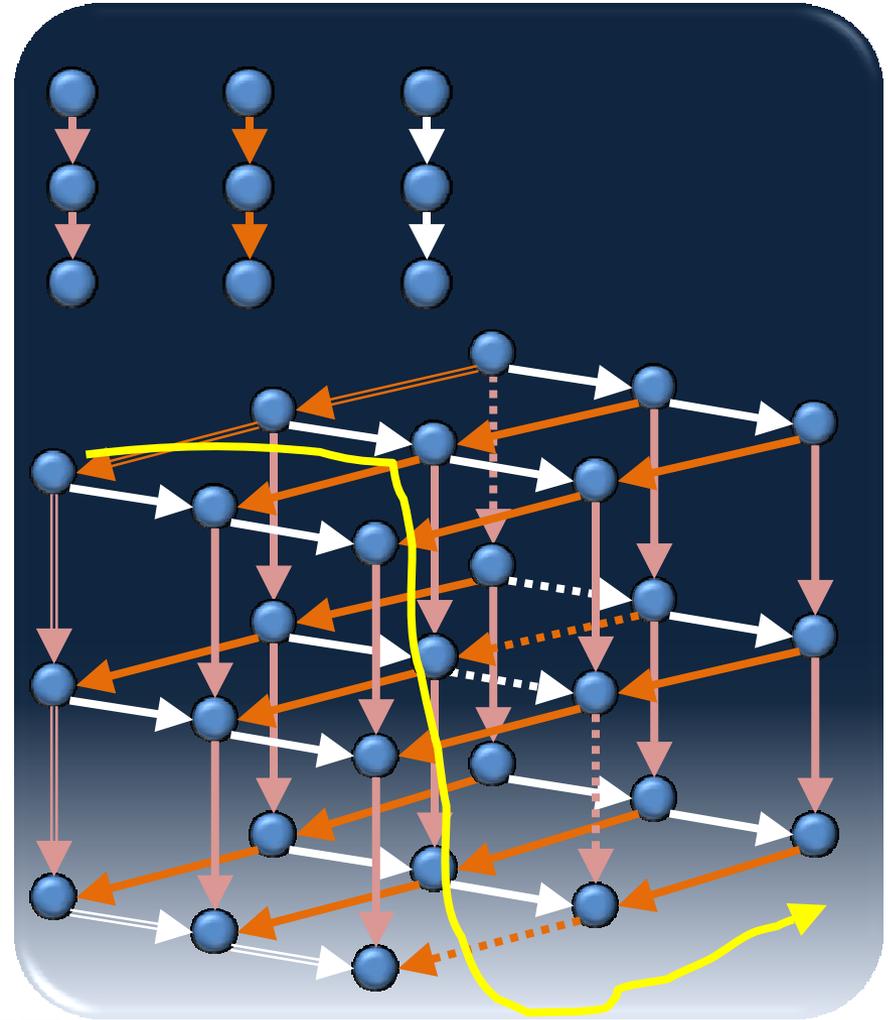
The Importance of Partial Order Reduction During Model Checking

- With 3 processes, the size of an interleaved state space is $p^s=27$
- Partial-order reduction explores representative sequences from each equivalence class
- Delays the execution of independent transitions

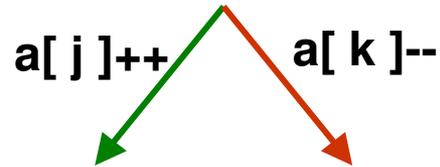


The Importance of Partial Order Reduction for Model Checking

- With 3 processes, the size of an interleaved state space is $p^s=27$
- Partial-order reduction explores representative sequences from each equivalence class
- Delays the execution of independent transitions
- In this example, it is possible to “get away” with 7 states (one interleaving)



Static POR Won't Always Do *(Flanagan and Godefroid, POPL 05)*



- **Action Dependence Determines COMMUTABILITY**
(POR theory is really detailed; it is more than commutability, but let's pretend it is ...)
- **Depends on $j == k$, in this example**
- **Can be very difficult to determine statically**
- **Can determine dynamically**

Similar Situation Arises with Wildcards...

Proc P:

Proc Q:

Proc R:

Send(to Q)

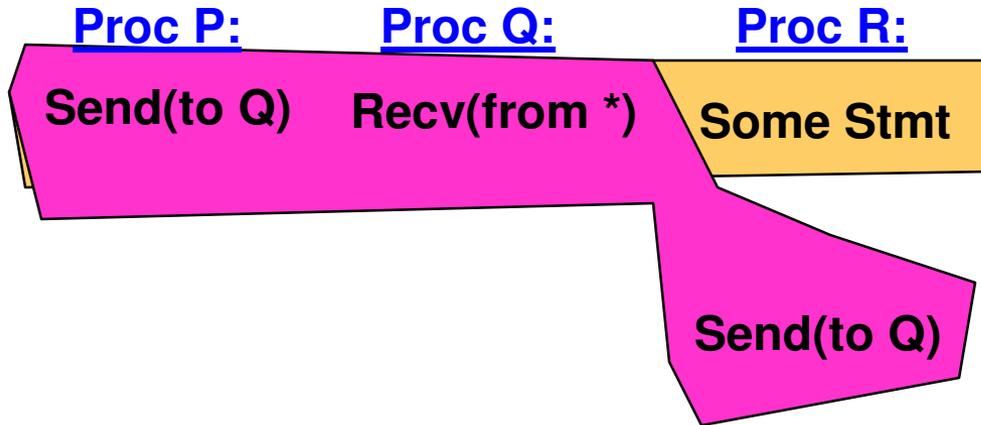
Recv(from *)

Some Stmt

Send(to Q)

- Dependencies may not be fully known, **JUST** by looking at enabled actions
- So **Conservative Assumptions** to be made (as in Urgent Algorithm)
- If not, Dependencies may be **Overlooked**
- The same problem exists with other “dynamic situations”
 - e.g. MPI_Cancel

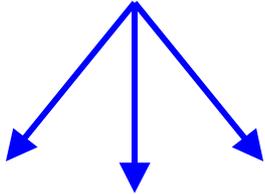
POR in the presence of Wildcards...



- Illustration of a Missed Dependency that would have been detected, had Proc R been scheduled first...

DPOR Exploits Knowledge of “Future” to Compute Dependencies More Accurately

Ample determined using “local” criteria



Nearest Dependent Transition Looking Back

{ **BT** }, { **Done** }

Add **Red** Process to “Backtrack Set”

This builds the “**Ample set**” incrementally based on observed dependencies

Blue is in “Done” set

Current State

Next move of **Red** process

How to define “Dependence” for MPI ?

- No a Priori Definition of when Actions Commute
- MPI Offers MANY API Calls
- So need SYSTEMATIC way to define “Dependence”
- CONTRIBUTION OF THIS PAPER:
 - Define Formal Semantics of MPI
 - Define Commutability Based on Formal Semantics

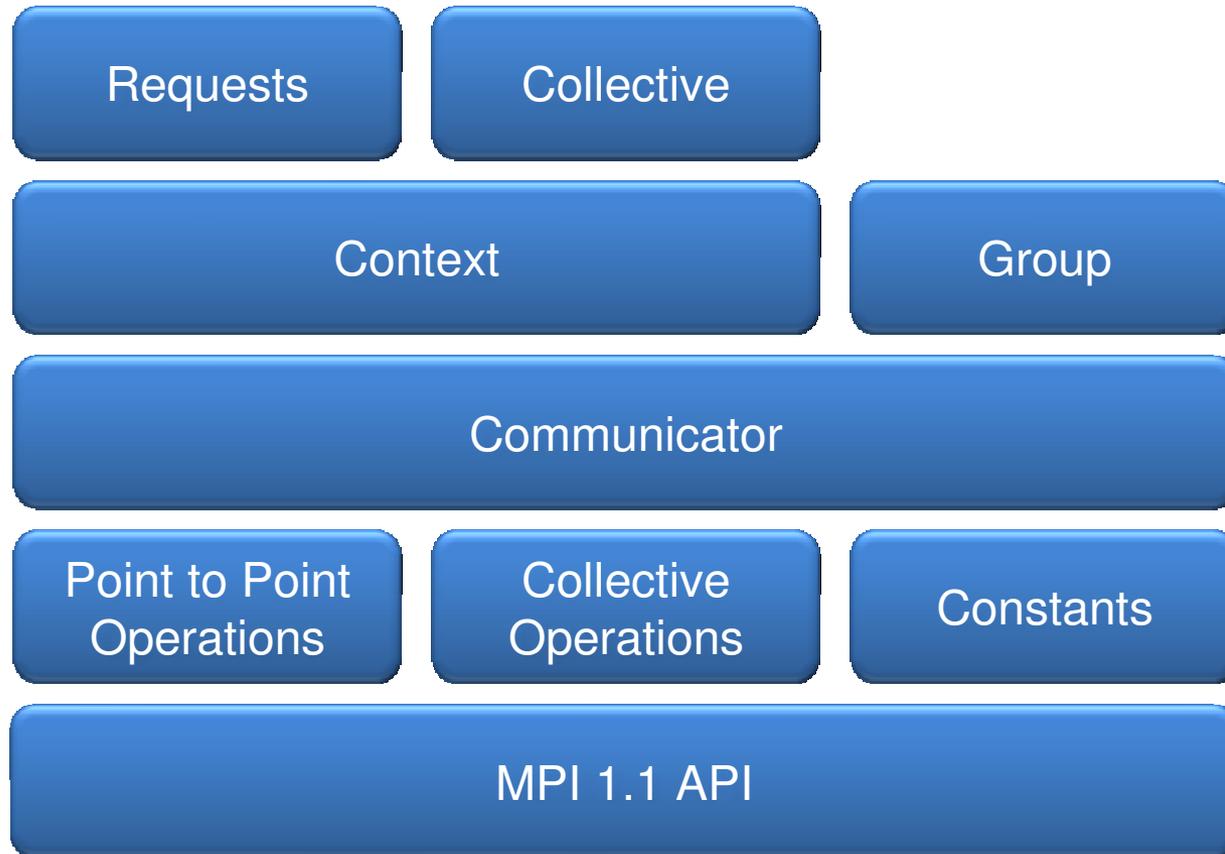
Spec of MPI_Wait (Slide 1/2) – FMICS07

```
1 MPI_Wait(request, status, return, proc) ==
2   LET r == requests[proc][Memory[proc][request]] IN
3   /\ Assert(initialized[proc] = "initialized", \* 200.10-200.12
4     "Error: MPI_Wait called with proc not in initialized state.")
5     \* 41.32-41.39 The request handle is not the null handle.
6   /\ \/ /\ Memory[proc][request] /= MPI_REQUEST_NULL
7     /\ r.localactive \* The request is active locally.
8     /\ \/ /\ r.message.src /= MPI_PROC_NULL \* The message src is not null
9       /\ r.message.dest /= MPI_PROC_NULL \* The message dest is not null
10        \* 41.32 - Blocks until complete
11        /\ \/ r.transmitted \* The message was transmitted or
12          \/ r.canceled \* canceled by the user program or
13          \/ r.buffered \* buffered by the system
14        /\ Memory' =
15          [Memory EXCEPT ![proc] = \* 41.36
16            [@ EXCEPT ![Status_Canceled(status)] =
17              /\ r.canceled
18                /\ \!not r.transmitted, \* 54.46
19                  ![Status_Count(status)] = r.message.numelements,
20                  ![Status_Source(status)] = r.message.src,
21                  ![Status_Tag(status)] = r.message.msgtag,
22                  ![Status_Err(status)] = r.error,
23                  ![request] = \* 41.32-41.35, 58.34-58.35
24                    IF r.persist
25                      THEN @
26                      ELSE MPI_REQUEST_NULL]]
27        \/ /\ \/ r.message.src = MPI_PROC_NULL
28          \/ r.message.dest = MPI_PROC_NULL
29        /\ Memory' = [Memory EXCEPT ![proc] = \* 41.36
30          [@ EXCEPT ![Status_Canceled(status)] = r.canceled,
31            ![Status_Count(status)] = 0,
32            ![Status_Source(status)] = MPI_PROC_NULL,
33            ![Status_Tag(status)] = MPI_ANY_TAG,
34            ![Status_Err(status)] = 0,
35            ![request] = \* 41.32-41.35, 58.34-58.35
36              IF r.persist
37                THEN @
38              ELSE MPI_REQUEST_NULL]]
```

Spec of MPI_Wait (Slide 2/2)

```
39     /\ requests' =
40         IF r.match /= << >>
41         THEN
42             [requests EXCEPT ![proc] = \* 58.34
43             [@ EXCEPT
44                 ![Memory[proc][request]] =
45                 IF r.persist
46                 THEN
47                     IF requests[r.match[1]][r.match[2]].localactive
48                     THEN [@ EXCEPT !.localactive = FALSE,
49                         !.globalactive = FALSE]
50                     ELSE [@ EXCEPT !.localactive = FALSE]
51                 ELSE
52                     IF requests[r.match[1]][r.match[2]].localactive
53                     THEN [@ EXCEPT !.localactive = FALSE,
54                         !.globalactive = FALSE,
55                         !.deallocated = TRUE]
56                     ELSE [@ EXCEPT !.localactive = FALSE,
57                         !.deallocated = TRUE]],
58             ![r.match[1]] =
59             [@ EXCEPT ![r.match[2]] =
60                 IF requests[r.match[1]][r.match[2]].localactive
61                 THEN requests[r.match[1]][r.match[2]]
62                 ELSE [@ EXCEPT !.globalactive = FALSE]]]
63         ELSE
64             [requests EXCEPT ![proc] = \* 58.34
65             [@ EXCEPT ![Memory[proc][request]] =
66                 IF r.persist
67                 THEN [@ EXCEPT !.localactive = FALSE]
68                 ELSE [@ EXCEPT !.localactive = FALSE,
69                     !.deallocated = TRUE]]]
70     \/ /\ \/ Memory[proc][request] = MPI_REQUEST_NULL \* 41.40-41.41 The
71     \/ /\ Memory[proc][request] /= MPI_REQUEST_NULL \* request handle is
72     /\ \!not r.localactive \* null or the request is not active
73     /\ Memory' = [Memory EXCEPT ![proc] = \* 41.36
74                 [@ EXCEPT ![Status_Canceled(status)] = FALSE,
75                     ![Status_Count(status)] = 0,
76                     ![Status_Source(status)] = MPI_ANY_SOURCE,
77                     ![Status_Tag(status)] = MPI_ANY_TAG,
78                     ![Status_Err(status)] = 0]]
79     /\ UNCHANGED << requests >>
80 /\ UNCHANGED << group, communicator, bufsize, message_buffer,
81     initialized, collective >>
```

MPI Formal Specification Organization



Example: Challenge posed by a 5-line MPI program...

```
p0: { Irecv(rcvbuf1, from p1);  
      Irecv(rcvbuf2, from p1); ... }
```

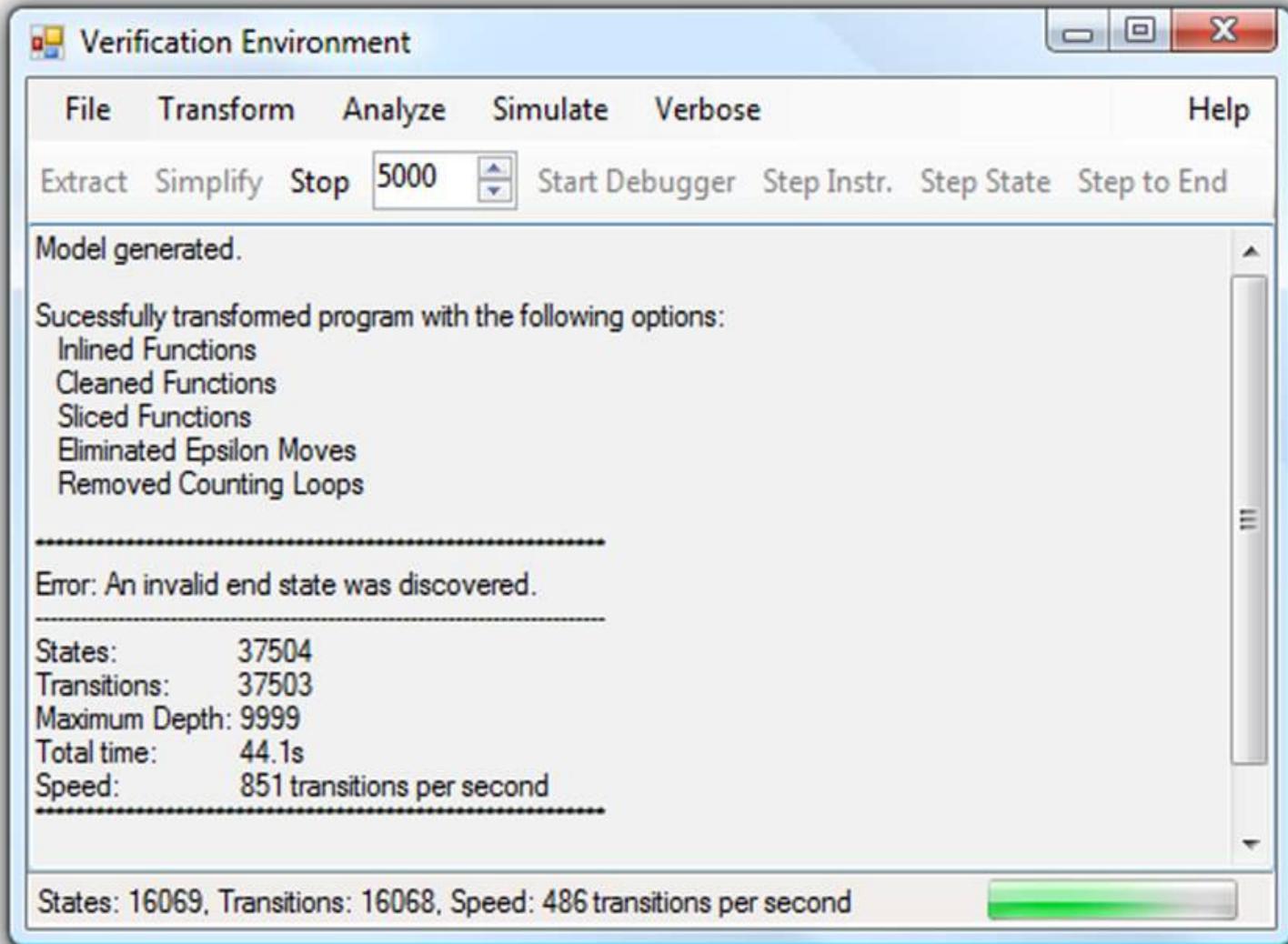
```
p1: { sendbuf1 = 6; sendbuf2 = 7;  
      Isend(sendbuf1, to p0);  
      Isend (sendbuf2, to p0); ... }
```

- In-order message delivery (rcvbuf1 == 6)
- Can access the buffers only after a later wait / test
- The second receive may complete before the first
- When Isend (synch.) is posted, all that is guaranteed is that Irecv(rcvbuf1,...) has been posted

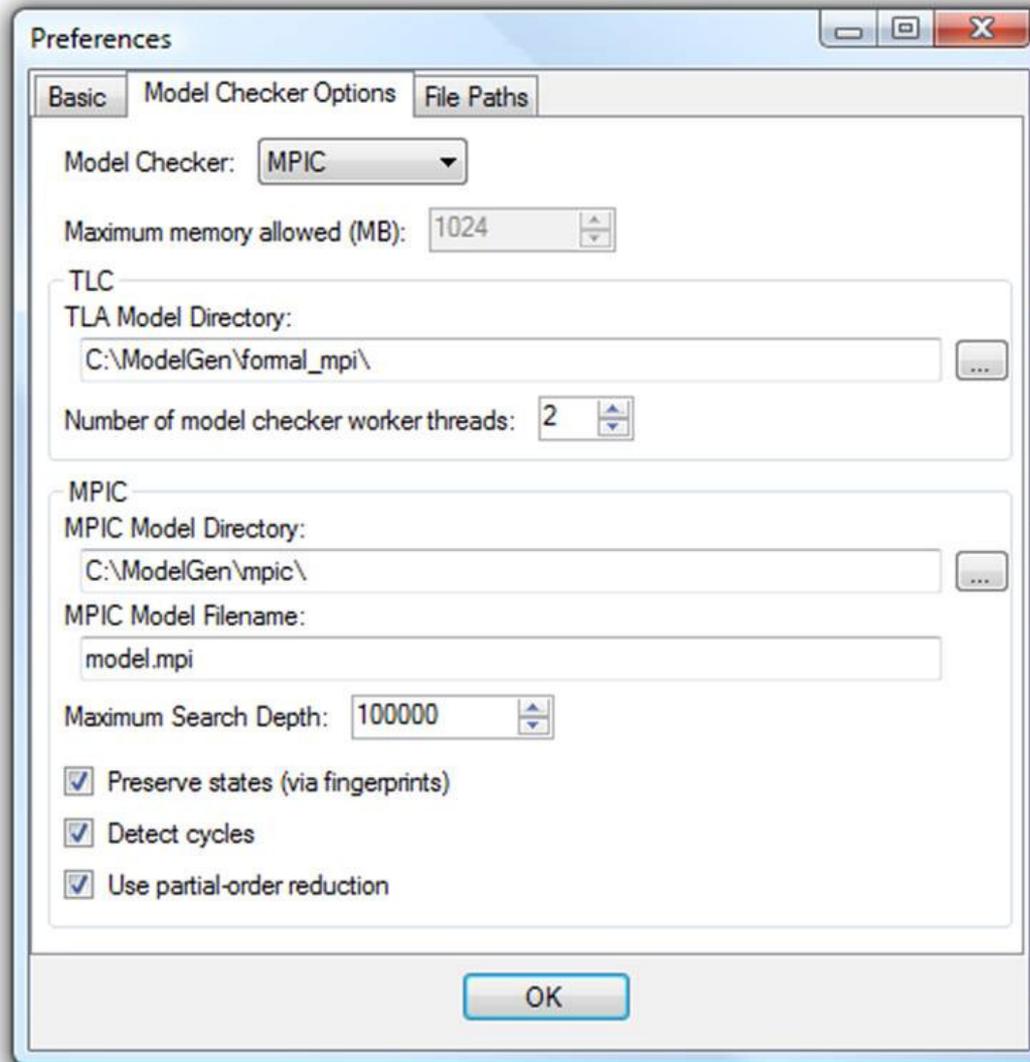
One of our Litmus Tests

```
1 #include "mpi.h"
2
3 int main(int argc, char** argv)
4 {
5     int rank, size, data1, data2, data3, flag;
6     MPI_Request req1, req2, req3;
7     MPI_Status stat;
8     MPI_Init(&argc, &argv);
9     MPI_Comm_rank(MPI_COMM_WORLD, &rank);
10    MPI_Comm_size(MPI_COMM_WORLD, &size);
11    if(rank == 0){
12        data1 = 0;
13        data2 = 0;
14        MPI_Irecv(&data1, 1, MPI_INT, 1,
15                0, MPI_COMM_WORLD, &req1);
16        MPI_Irecv(&data2, 1, MPI_INT, 1,
17                1, MPI_COMM_WORLD, &req2);
18        MPI_Irecv(&data3, 1, MPI_INT, 1,
19                2, MPI_COMM_WORLD, &req3);
20    } else {
21        data1 = 7;
22        data2 = 6;
23        MPI_Issend(&data1, 1, MPI_INT, 0,
24                 1, MPI_COMM_WORLD, &req1);
25    }
26    if(rank == 1){
27        MPI_Wait(&req1, &stat);
28        MPI_Irsend(&data2, 1, MPI_INT, 0,
29                 0, MPI_COMM_WORLD, &req2);
30        MPI_Irsend(&data3, 1, MPI_INT, 0,
31                 2, MPI_COMM_WORLD, &req3);
32    } else {
33        MPI_Wait(&req2, &stat);
34    }
35    if(rank == 0){
36        MPI_Wait(&req1, &stat);
37    } else {
38        MPI_Wait(&req2, &stat);
39    }
40    MPI_Finalize();
41    return 0;
42 }
```

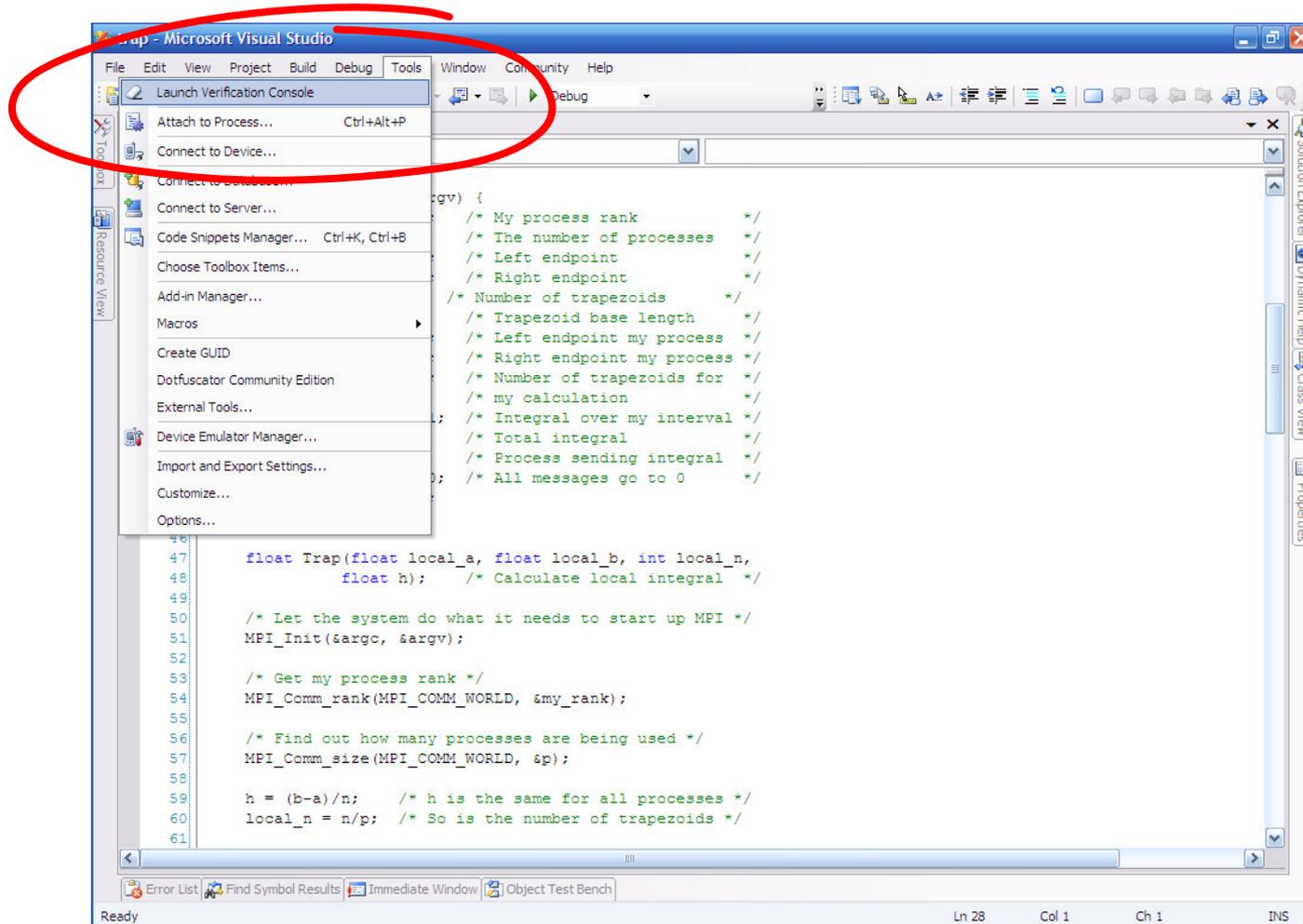
The Histrionics of FV for HPC (1)



The Histrionics of FV for HPC (2)



Error-trace Visualization in VisualStudio



This paper: Simplified Semantics (e.g. as shown by `MPI_Wait`)

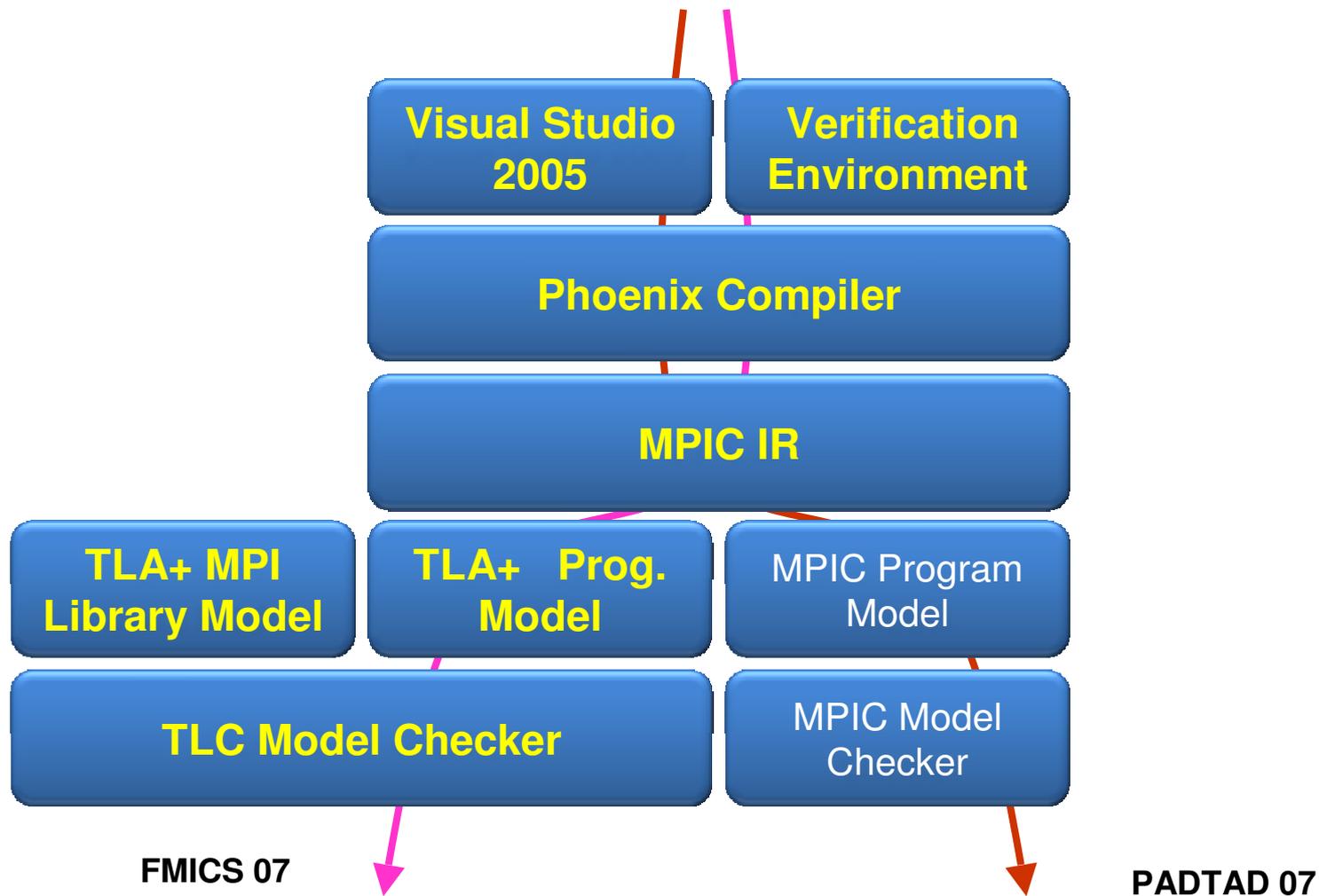
$$\begin{array}{c}
 \Sigma(c, p) \\
 p(i) = (l, g) \\
 \wedge \text{proc}(l(\text{vars}(\text{pc}))) = \text{'wait } e\text{'} \\
 \wedge E[\mathbf{e}, p_i] = 0 \vee \text{Completed}(g(E[\mathbf{e}, p_i])) \\
 \hline
 \Sigma(c, p[\quad i \mapsto (\quad l[\text{vars}(\text{pc}) \mapsto \text{next}(l(\text{vars}(\text{pc}))), \\
 \quad \text{vars}(\mathbf{e}) \mapsto 0]])
 \end{array}$$

$$\begin{array}{c}
 \Sigma(c, p) \wedge \\
 p(i) = (l_i, g_i) \wedge \text{proc}(l_i(\text{vars}(\text{pc}))) = \text{'wait } e\text{'} \\
 \wedge E[\mathbf{e}, p_i] \in \text{Dom}(g_i) \wedge \\
 \exists j : \quad p(j) = (l_j, g_j) \wedge \\
 \quad \exists k : \quad \text{Match}(g_j(k), g_i(E[\mathbf{e}, p_i])) \wedge \\
 \quad \quad \forall m < k : \neg \text{Match}(g_j(m), g_i(E[\mathbf{e}, p_i])) \\
 \hline
 \Sigma(c, p[\quad i \mapsto (\quad l_i[\text{vars}(\text{pc}) \mapsto \text{next}(l_i(\text{vars}(\text{pc}))), \\
 \quad \quad \text{vars}(\mathbf{e}) \mapsto 0] \\
 \quad \quad \quad g_i[E[\mathbf{e}, p_i] \mapsto g_i(E[\mathbf{e}, p_i])[false/true]] \\
 \quad j \mapsto (\quad l_j, g_j[k \mapsto g_j(k)[false/true]]))
 \end{array}$$

Independence Theorems based on Formal Semantics of MPI Subset

- 1. Local actions (Assignment, Goto, Alloc, Assert) are independent of all transitions of other processes.**
- 2. Barrier actions (Barrier_init, Barrier_wait) are independent of all transitions of other processes.**
- 3. Issend and Irecv are independent of all transitions of other processes except Wait and Test.**
- 4. Wait and Test are independent of all transitions of other processes except Issend and Irecv.**

Executable Formal Specification and MPIC Model Checker Integration into VS



A Simple Example:

e.g. mismatched send/rcv causing deadlock

```
/* Add-up integrals calculated by each process */
```

```
if (my_rank == 0) {  
    total = integral;  
    for (source = 0; source < p; source++) {  
        MPI_Recv(&integral, 1, MPI_FLOAT, source,  
                tag, MPI_COMM_WORLD, &status);  
        total = total + integral;  
    }  
} else {  
    MPI_Send(&integral, 1, MPI_FLOAT, dest,  
            tag, MPI_COMM_WORLD);  
}
```

p0:fr 0 p0:fr 1 p0:fr 2
p1:to 0 p2:to 0 p3:to 0

***Partial Demo
of DPOR Tool
for MPIC***

So, the whole story (i.e. Conclusions)...

- **Preliminary Formal Semantics of MPI in Place (50 point-to-point functions)**
- **Can Model-Check this Golden Semantics**
- **About 5 of these 30 have a more rigorous characterization thru Independence Theorems**
- **For MPI Programs using These MPI functions, we have a DPOR based model checker MPIC**
- **Integrated in the VS Framework with MPI-TLC also**

- **Theory Expected to Carry Over into In-Situ Dynamic Partial Order Reduction (model-check without model building – EuroPVM / MPI 2007)**

Questions ?

The verification environment is downloadable from

http://www.cs.utah.edu/formal_verification/mpic

It is at an early stage of development

Answers!

- 1. We are extending it to Collective Operations**
 - lesson learned from de Supinski
- 2. We may perform Formal Testing of MPI Library Implementations based on the Formal Semantics**
- 3. We plan to analyze mixed MPI / Threads**
- 4. That is a very good question – let's talk!**