

## MPI Runtime Error Detection with MUST

---

Joachim Protze  
IT Center RWTH Aachen University  
Feb 2022

---

## How many issues can you spot in this tiny example?

---

```
#include <mpi.h>
#include <stdio.h>

int main (int argc, char** argv)
{
    int rank, size, buf[8];

    MPI_Comm_rank (MPI_COMM_WORLD, &rank);
    MPI_Comm_size (MPI_COMM_WORLD, &size);

    MPI_Datatype type;
    MPI_Type_contiguous (2, MPI_INTEGER, &type);

    MPI_Recv (buf, 2, MPI_INT, size - rank, 123, MPI_COMM_WORLD, MPI_STATUS_IGNORE);
    MPI_Send (buf, 2, type, size - rank, 123, MPI_COMM_WORLD);

    printf ("Hello, I am rank %d of %d.\n", rank, size);

    return 0;
}
```

At least 8 issues in this code example!

# Content

---

- Motivation
- **MPI usage errors**
- Examples: Common MPI usage errors
  - Including MUST's error descriptions
- MUST usage
- Hands-on

## Motivation

---

- MPI programming is error prone
- Portability errors  
(just on some systems, just for some runs)
- Bugs may manifest as:
  - Crash
  - Application hanging
  - Finishes
- Questions:
  - Why crash/hang?
  - Is my result correct?
  - Will my code also give correct results on another system?
- Tools help to pin-point these bugs



## Common MPI Error Classes

- Common syntactic errors:
  - Incorrect arguments
  - Resource usage
  - Lost/Dropped Requests
  - Buffer usage
  - Type-matching
  - Deadlocks

Tool to use:  
**MUST,**  
**Static analysis tool,**  
**(Debugger)**

- Semantic errors that are correct in terms of MPI standard, but do not match the programmer's intention:
  - Displacement/Size/Count errors

Tool to use:  
**Debugger**

# Content

---

- Motivation
- MPI usage errors
- **Examples: Common MPI usage errors**
  - **Including MUST's error descriptions**
- MUST usage
- Hands-on

## Already fixed missing MPI\_Init/Finalize:

---

```
#include <mpi.h>
#include <stdio.h>

int main (int argc, char** argv)
{
    int rank, size, buf[8];

    MPI_Init (&argc, &argv);
    MPI_Comm_rank (MPI_COMM_WORLD, &rank);
    MPI_Comm_size (MPI_COMM_WORLD, &size);

    MPI_Datatype type;
    MPI_Type_contiguous (2, MPI_INTEGER, &type);

    MPI_Recv (buf, 2, MPI_INT, size - rank - 1, 123, MPI_COMM_WORLD, MPI_STATUS_IGNORE);

    MPI_Send (buf, 2, type, size - rank - 1, 123, MPI_COMM_WORLD);

    printf ("Hello, I am rank %d of %d.\n", rank, size);

    MPI_Finalize ();

    return 0;
}
```

# Must detects deadlocks

MUST Output, starting date: Fri Mar 24 11:59:41 20...

Rank(s)	Type	Message
	Error	The application issued a set of MPI calls that can cause a deadlock! A graphical representation of this situation is available in a <a href="#">detailed de...</a>

Details:

Message	From	References
The application issued a set of MPI calls that can cause a deadlock! A graphical representation of this situation is available in a <a href="#">detailed deadlock view (MUST_Output-files/MUST_Deadlock.html)</a> . References 1-2 list the involved calls (limited to the first 5 calls, further calls may be involved). The application still runs if the deadlock manifested (e.g. caused a hang on this MPI implementation) you can attach to the involved ranks with a debugger or abort the application (if necessary).		References of a representative process: reference 1 rank 0: <b>MPI_Recv</b> (1st occurrence) called from: #0 main@example.c:15  reference 2 rank 3: <b>MPI_Recv</b> (1st occurrence) called from: #0 main@example.c:15

Who?

What?

Where?

Details

Click for graphical representation of the detected deadlock situation.



Message	
<p>The application issued a set of MPI calls that can cause a deadlock! The graphs below show details on this situation. This includes a wait-for graph that shows active wait-for dependencies between the processes that cause the deadlock. Note that this process set only includes processes that cause the deadlock and no further processes. A legend details the wait-for graph components in addition, while a parallel call stack view summarizes the locations of the MPI calls that cause the deadlock. Below these graphs, a message queue graph shows active and unmatched point-to-point communications. This graph only includes operations that could have been intended to match a point-to-point operation that is relevant to the deadlock situation. Finally, a parallel call stack shows the locations of any operation in the parallel call stack. The leaves of this call stack graph show the components of the message queue graph that they span. The application still runs, if the deadlock manifested (e.g. caused a hang on this MPI implementation) you can attach to the involved ranks with a debugger or abort the application (if necessary).</p>	
Active Communicators	
Comm:	A
MPI COMM WORLD	
Wait-for Graph	Legend
<p>Rank 0 waits for rank 1 and vv.</p>	<p>Active MPI Call</p> <p>Sub Operation</p> <p>A → A waits for B and C → B</p> <p>A → A waits for B and C → C</p> <p>A → A waits for B or C → B</p> <p>A → A waits for B or C → C</p>
Call Stack	
<p>main@/rwthfs/rz/cluster/home/pj416018/must-example/VI-HPS/example.c:15</p> <p>Ranks:</p> <p>MPI_Recv</p> <p>Simple call stack for this example.</p>	
Active and Relevant Point-to-Point Messages: Overview	
Message queue	
Active and Relevant Point-to-Point Messages: Callstack-view	
stack	

## Fix1: use asynchronous receive

```
#include <mpi.h>
#include <stdio.h>

int main (int argc, char** argv)
{
    int rank, size, buf[8];

    MPI_Init (&argc, &argv);
    MPI_Comm_rank (MPI_COMM_WORLD, &rank);
    MPI_Comm_size (MPI_COMM_WORLD, &size);

    MPI_Datatype type;
    MPI_Type_contiguous (2, MPI_INTEGER, &type);

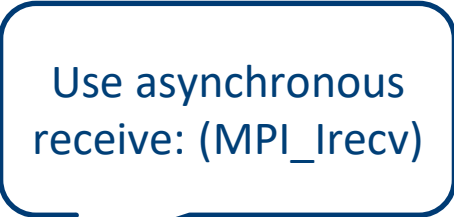
    MPI_Request request;
    MPI_Irecv (buf, 2, MPI_INT, size - rank - 1, 123, MPI_COMM_WORLD, &request);

    MPI_Send (buf, 2, type, size - rank - 1, 123, MPI_COMM_WORLD);

    printf ("Hello, I am rank %d of %d.\n", rank, size);

    MPI_Finalize ();

    return 0;
}
```



Use asynchronous receive: (MPI\_Irecv)

# MUST detects errors in transfer buffer sizes / types

Rank(s)	Type		
2(28793)	<b>Error</b>	A receive operation uses a (datatype, count) pair that can not hold the data transferred by the send it matches! The first element of the send...	
Details:			
		Message	References
		<p>A receive operation uses a (datatype, count) pair that can not hold the data transferred by the send it matches! The first element of the send that did not fit into the receive operation is at (contiguous)[0](MPI_INTEGER) in the send type (consult the MUST manual for a detailed description of datatype positions). The send operation was started at reference 1, the receive operation was started at reference 2. (Information on communicator: MPI_COMM_WORLD) (Information on send of count 2 with type:Datatype created at reference 3 is for Fortran, based on the following type(s): { MPI_INTEGER}) (Information on receive of count 2 with type:MPI_INT)</p>	<p>Representative location:  <b>MPI_Send</b> (1st occurrence) called from:                      #0 main@example-fix1.c:18</p> <p>reference 2 rank 1: <b>MPI_Irecv</b> (1st occurrence) called from:                      #0 main@example-fix1.c:16</p> <p>reference 3 rank 2:  <b>MPI_Type_contiguous</b> (1st occurrence) called from:                      #0 main@example-fix1.c:13</p>
1(28792)	<b>Error</b>	A receive operation uses a (datatype, count) pair that can not hold the data transferred by the send it matches! The first element of the send...	
0-3	<b>Error</b>	Argument 3 (datatype) is not committed for transfer, call MPI_Type_commit before using the type for transfer!(Information on datatypeData...	
2(28793)	<b>Error</b>	The memory regions to be transferred by this send operation overlap with regions spanned by a pending non-blocking receive operation!(In...	
1(28792)	<b>Error</b>	The memory regions to be transferred by this send operation overlap with regions spanned by a pending non-blocking receive operation!(In...	
3(28795)	<b>Error</b>	The memory regions to be transferred by this send operation overlap with regions spanned by a pending non-blocking receive operation!(In...	
3(28795)	<b>Error</b>	A receive operation uses a (datatype, count) pair that can not hold the data transferred by the send it matches! The first element of the send...	
0(28794)	<b>Error</b>	The memory regions to be transferred by this send operation overlap with regions spanned by a pending non-blocking receive operation!(In...	
0(28794)	<b>Error</b>	A receive operation uses a (datatype, count) pair that can not hold the data transferred by the send it matches! The first element of the send...	

Size of sent message larger than receive buffer

All detected errors are collapsed for overview - click to expand

## Fix2: use same message size for send and receive

```
#include <mpi.h>
#include <stdio.h>

int main (int argc, char** argv)
{
    int rank, size, buf[8];

    MPI_Init (&argc, &argv);
    MPI_Comm_rank (MPI_COMM_WORLD, &rank);
    MPI_Comm_size (MPI_COMM_WORLD, &size);

    MPI_Datatype type;
    MPI_Type_contiguous (2, MPI_INTEGER, &type);
    MPI_Type_commit (&type);

    MPI_Request request;
    MPI_Irecv (buf, 2, MPI_INT, size - rank - 1, 123, MPI_COMM_WORLD, &request);

    MPI_Send (buf, 1, type, size - rank - 1, 123, MPI_COMM_WORLD);

    printf ("Hello, I am rank %d of %d.\n", rank, size);

    MPI_Finalize ();

    return 0;
}
```

Reduce the  
message size

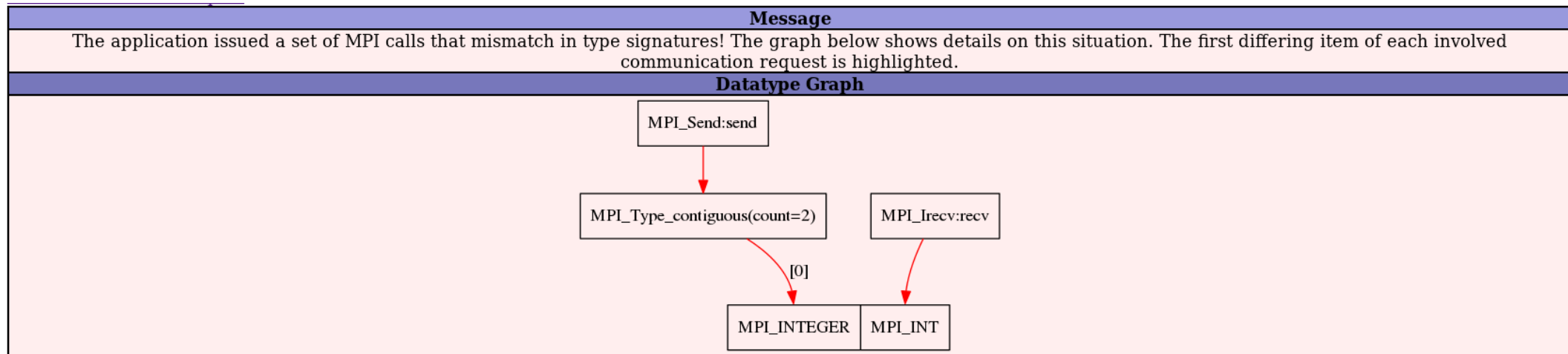
# MUST detects errors in handling datatypes

Rank(s)	Type	Message						
2(17250)	Error	A send and a receive operation use datatypes that do not match! Mismatch occurs at (contiguous)[0](MPI_INTEGER) in the send type and a...						
Details:								
		<table border="1"> <thead> <tr> <th>Message</th> <th>From</th> <th>References</th> </tr> </thead> <tbody> <tr> <td>A send and a receive operation use datatypes that do not match! Mismatch occurs at (contiguous)[0](MPI_INTEGER) in the send type and at (MPI_INT) in the receive type (consult the MUST manual for a detailed description of datatype positions). A graphical representation of this situation is available in a <a href="#">detailed type mismatch view (MUST_Output-files/MUST_Typemismatch_74088185856002.html)</a>. The send operation was started at reference 1, the receive operation was started at reference 2. (Information on communicator: MPI_COMM_WORLD) (Information on send of count 1 with type:Datatype created at reference 3 is for Fortran, based on the following type(s): { MPI_INTEGER}) (Information on receive of count 2 with type:MPI_INT)</td> <td>Representative location: <b>MPI_Send</b> (1st occurrence) called from: #0 main@example-fix2.c:18</td> <td>References of a representative process:  reference 1 rank 2: <b>MPI_Send</b> (1st occurrence) called from: #0 main@example-fix2.c:18  reference 2 rank 1: <b>MPI_Irecv</b> (1st occurrence) called from: #0 main@example-fix2.c:16  reference 3 rank 2: <b>MPI_Type_contiguous</b> (1st occurrence) called from: #0 main@example-fix2.c:13</td> </tr> </tbody> </table>	Message	From	References	A send and a receive operation use datatypes that do not match! Mismatch occurs at (contiguous)[0](MPI_INTEGER) in the send type and at (MPI_INT) in the receive type (consult the MUST manual for a detailed description of datatype positions). A graphical representation of this situation is available in a <a href="#">detailed type mismatch view (MUST_Output-files/MUST_Typemismatch_74088185856002.html)</a> . The send operation was started at reference 1, the receive operation was started at reference 2. (Information on communicator: MPI_COMM_WORLD) (Information on send of count 1 with type:Datatype created at reference 3 is for Fortran, based on the following type(s): { MPI_INTEGER}) (Information on receive of count 2 with type:MPI_INT)	Representative location: <b>MPI_Send</b> (1st occurrence) called from: #0 main@example-fix2.c:18	References of a representative process:  reference 1 rank 2: <b>MPI_Send</b> (1st occurrence) called from: #0 main@example-fix2.c:18  reference 2 rank 1: <b>MPI_Irecv</b> (1st occurrence) called from: #0 main@example-fix2.c:16  reference 3 rank 2: <b>MPI_Type_contiguous</b> (1st occurrence) called from: #0 main@example-fix2.c:13
Message	From	References						
A send and a receive operation use datatypes that do not match! Mismatch occurs at (contiguous)[0](MPI_INTEGER) in the send type and at (MPI_INT) in the receive type (consult the MUST manual for a detailed description of datatype positions). A graphical representation of this situation is available in a <a href="#">detailed type mismatch view (MUST_Output-files/MUST_Typemismatch_74088185856002.html)</a> . The send operation was started at reference 1, the receive operation was started at reference 2. (Information on communicator: MPI_COMM_WORLD) (Information on send of count 1 with type:Datatype created at reference 3 is for Fortran, based on the following type(s): { MPI_INTEGER}) (Information on receive of count 2 with type:MPI_INT)	Representative location: <b>MPI_Send</b> (1st occurrence) called from: #0 main@example-fix2.c:18	References of a representative process:  reference 1 rank 2: <b>MPI_Send</b> (1st occurrence) called from: #0 main@example-fix2.c:18  reference 2 rank 1: <b>MPI_Irecv</b> (1st occurrence) called from: #0 main@example-fix2.c:16  reference 3 rank 2: <b>MPI_Type_contiguous</b> (1st occurrence) called from: #0 main@example-fix2.c:13						
0(17249)	Error	that do not m						
1(17248)	Error	that do not m						
3(17251)	Error	that do not						
0-3	Error	transfer, call						
Details:								
		<table border="1"> <thead> <tr> <th>Message</th> <th>From</th> <th>References</th> </tr> </thead> <tbody> <tr> <td>Argument 3 (datatype) is not committed for transfer, call MPI_Type_commit before using the type for transfer! (Information on datatypeDatatype created at reference 1 is for Fortran, based on the following type(s): { MPI_INTEGER})</td> <td>Representative location: <b>MPI_Send</b> (1st occurrence) called from: #0 main@example-fix2.c:18</td> <td>References of a representative process:  reference 1 rank 2: <b>MPI_Type_contiguous</b> (1st occurrence) called from: #0 main@example-fix2.c:13</td> </tr> </tbody> </table>	Message	From	References	Argument 3 (datatype) is not committed for transfer, call MPI_Type_commit before using the type for transfer! (Information on datatypeDatatype created at reference 1 is for Fortran, based on the following type(s): { MPI_INTEGER})	Representative location: <b>MPI_Send</b> (1st occurrence) called from: #0 main@example-fix2.c:18	References of a representative process:  reference 1 rank 2: <b>MPI_Type_contiguous</b> (1st occurrence) called from: #0 main@example-fix2.c:13
Message	From	References						
Argument 3 (datatype) is not committed for transfer, call MPI_Type_commit before using the type for transfer! (Information on datatypeDatatype created at reference 1 is for Fortran, based on the following type(s): { MPI_INTEGER})	Representative location: <b>MPI_Send</b> (1st occurrence) called from: #0 main@example-fix2.c:18	References of a representative process:  reference 1 rank 2: <b>MPI_Type_contiguous</b> (1st occurrence) called from: #0 main@example-fix2.c:13						

Use of Fortran type in C, datatype mismatch between sender and receiver

Use of uncommitted datatype: type

## MUST detects errors in handling datatypes



Graphical representation of the type mismatch

## Fix3: use MPI\_Type\_commit

## Fix4: use C integer type

```
#include <mpi.h>
#include <stdio.h>

int main (int argc, char** argv)
{
    int rank, size, buf[8];

    MPI_Init (&argc, &argv);
    MPI_Comm_rank (MPI_COMM_WORLD, &rank);
    MPI_Comm_size (MPI_COMM_WORLD, &size);

    MPI_Datatype type;
    MPI_Type_contiguous (2, MPI_INT, &type);
    MPI_Type_commit (&type);

    MPI_Request request;
    MPI_Irecv (buf, 2, MPI_INT, size - rank - 1, 123, MPI_COMM_WORLD, &request);

    MPI_Send (buf, 2, type, size - rank - 1, 123, MPI_COMM_WORLD);

    printf ("Hello, I am rank %d of %d.\n", rank, size);

    MPI_Finalize ();

    return 0;
}
```

Use the integer datatype intended for usage in C

Commit the datatype before usage

# MUST detects data races in asynchronous communication

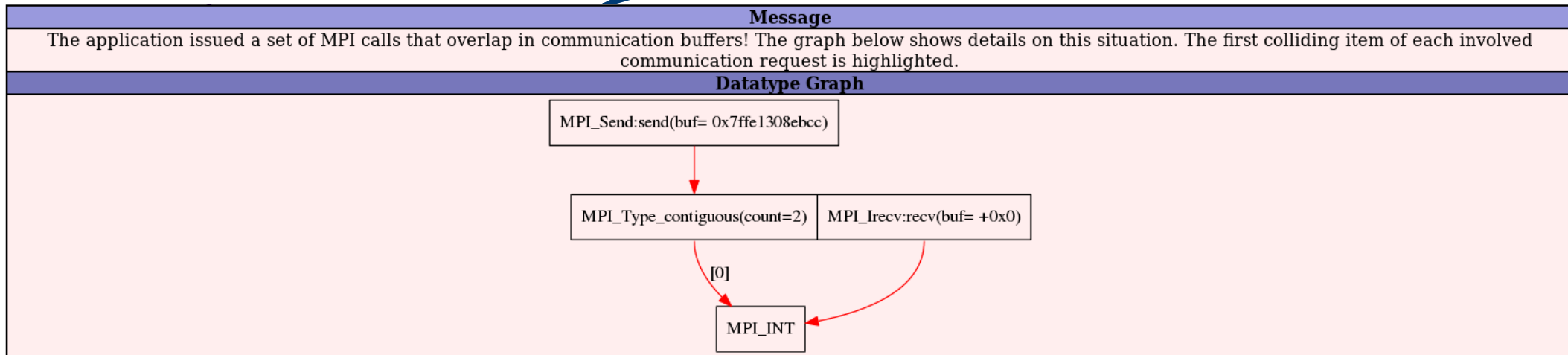
Data race between send and asynchronous receive operation

Rank(s)	Type		
0(1605)	<b>Error</b>	The memory regions to be transferred by this	receive operation!(In...
Details:			
Message		From	References
<p>The memory regions to be transferred by this send operation overlap with regions spanned by a pending non-blocking receive operation!</p> <p>(Information on the request associated with the other communication: Point-to-point request activated at reference 1)</p> <p>(Information on the datatype associated with the other communication: MPI_INT)</p> <p>The other communication overlaps with this communication at position:(MPI_INT)</p> <p>(Information on the datatype associated with this communication: Datatype created at reference 2 is for C, committed at reference 3, based on the following type(s): { MPI_INT})</p> <p>This communication overlaps with the other communication at position:(contiguous) [0](MPI_INT)</p> <p>A graphical representation of this situation is available in a <a href="#">detailed overlap view (MUST_Output-files/MUST_Overlap_6893422510080_0.html)</a>.</p>		<p>Representative location: <b>MPI_Send</b> (1st occurrence) called from: #0 main@example-fix3.c:19</p>	<p>References of a representative process:</p> <p>reference 1 rank 0: <b>MPI_Irecv</b> (1st occurrence) called from: #0 main@example-fix3.c:17</p> <p>reference 2 rank 0: <b>MPI_Type_contiguous</b> (1st occurrence) called from: #0 main@example-fix3.c:13</p> <p>reference 3 rank 0: <b>MPI_Type_commit</b> (1st occurrence) called from: #0 main@example-fix3.c:14</p>
3(1610)	<b>Error</b>	The memory regions to be transferred by this send operation overlap with regions spanned by a pending non-blocking receive operation!(In...	
2(1608)	<b>Error</b>	The memory regions to be transferred by this send operation overlap with regions spanned by a pending non-blocking receive operation!(In...	
1(1606)	<b>Error</b>	The memory regions to be transferred by this send operation overlap with regions spanned by a pending non-blocking receive operation!(In...	
0-3	<b>Error</b>	There are 1 datatypes that are not freed when MPI Finalize was issued, a quality application should free all MPI resources before calling ...	
0-3	<b>Error</b>	There are 1 requests that are not freed when MPI Finalize was issued, a quality application should free all MPI resources before calling M...	



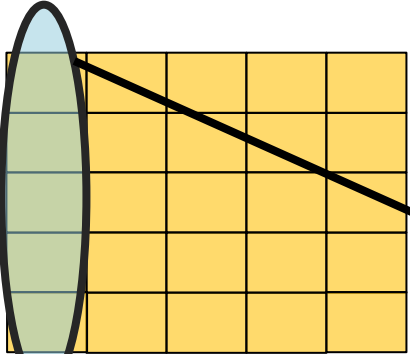
## Graphical representation of the race condition

Graphical representation of the data race location



## Errors with MPI Datatypes - Overview

- Derived datatypes use constructors, example:

- 

2D Field  
(of integers)

```
MPI_Type_vector (  
  NumRows      /*count*/,  
  1            /*blocklength*/,  
  NumColumns   /*stride*/,  
  MPI_INT      /*oldtype*/,  
  &newType);
```

- Errors that involve datatypes can be complex:
  - Need to be detected correctly
  - Need to be visualized

## Errors with MPI Datatypes - Example

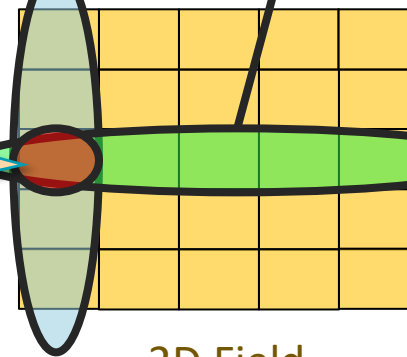
### ■ C Code:

```
...
MPI_Isend(buf, 1 /*count*/, vectortype, target, tag,
          MPI_COMM_WORLD, &request);
MPI_Recv(buf, 1 /*count*/, columntype, target, tag,
         MPI_COMM_WORLD, &status);
MPI_Wait (&request, &status);
...
```

### ■ Memory:

Error: buffer overlap

MPI\_Isend reads, MPI\_Recv writes at the same time



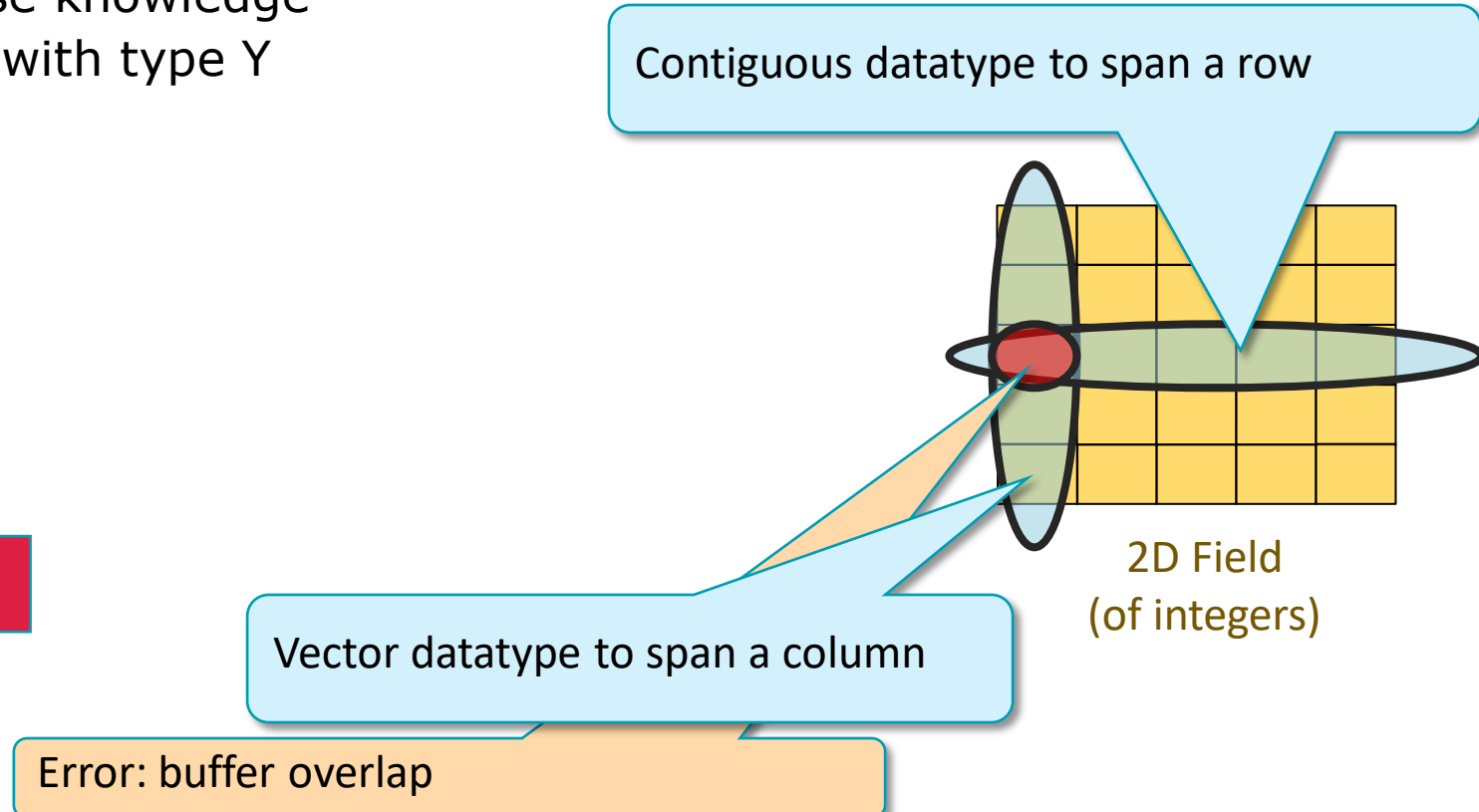
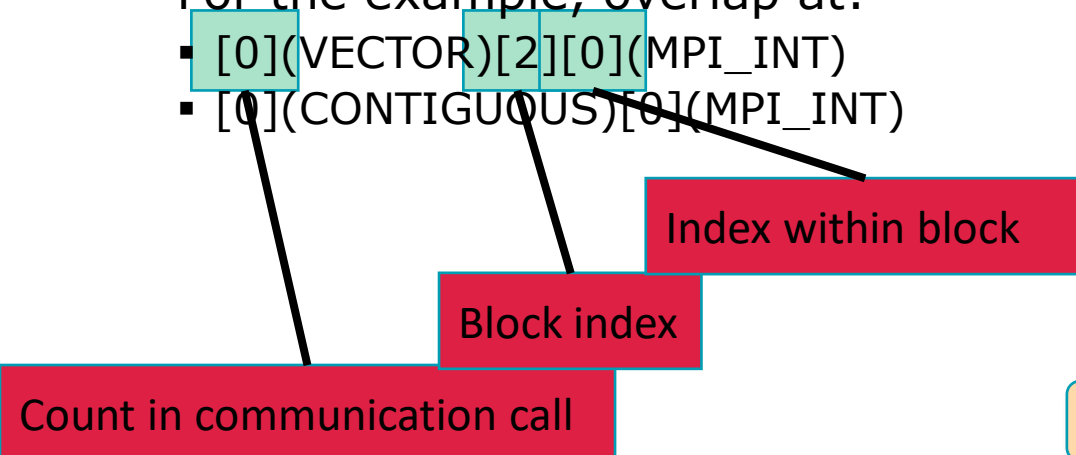
2D Field  
(of integers)

A tool must:

- Detect the error
- Pinpoint the user to the exact problem

## Errors with MPI Datatypes - Error Positions

- How to point to an error in a derived datatype?
  - Derived types can span wide areas of memory
  - Understanding errors requires precise knowledge
  - E.g., not sufficient: Type X overlaps with type Y
- Example:
- We use path expressions to point to error positions
  - For the example, overlap at:
    - `[0](VECTOR)[2][0](MPI_INT)`
    - `[0](CONTIGUOUS)[0](MPI_INT)`



## Fix5: use independent memory regions

```
#include <mpi.h>
#include <stdio.h>

int main (int argc, char** argv)
{
    int rank, size, buf[8];

    MPI_Init (&argc, &argv);
    MPI_Comm_rank (MPI_COMM_WORLD, &rank);
    MPI_Comm_size (MPI_COMM_WORLD, &size);

    MPI_Datatype type;
    MPI_Type_contiguous (2, MPI_INTEGER, &type);
    MPI_Type_commit (&type);

    MPI_Request request;
    MPI_Irecv (buf, 2, MPI_INT, size - rank - 1, 123, MPI_COMM_WORLD, &request);

    MPI_Send (buf + 4, 1, type, size - rank - 1, 123, MPI_COMM_WORLD);

    printf ("Hello, I am rank %d of %d.\n", rank, size);

    MPI_Finalize ();

    return 0;
}
```

Offset points to  
independent  
memory

# MUST detects leaks of user defined objects

Rank(s)	Type	Message
0-3	<b>Error</b>	There are 1 datatypes that are not freed when MPI Finalize was issued, a quality application should free all MPI resources before calling ...
Details:		
Message	From	References
There are 1 datatypes that are not freed when MPI Finalize was issued, a quality application should free all MPI resources before calling MPI_Finalize. Listing information for these datatypes:  -Datatype 1: Datatype created at reference 1 is for C, committed at reference 2, based on the following type(s): { MPI_INT}	Representative location: <b>MPI_Type_contiguous</b> (1st occurrence) called from: #0 main@example-fix4.c:13	References of a representative process: reference 1 rank 1: <b>MPI_Type_contiguous</b> (1st occurrence) called from: #0 main@example-fix4.c:13  reference 2 rank 1: <b>MPI_Type_commit</b> (1st occurrence) called from: #0 main@example-fix4.c:14
0-3	<b>Error</b>	There are 1 requests that are not freed when MPI Finalize was issued, a quality application should free all MPI resources before calling M...
Details:		
Message	From	References
There are 1 requests that are not freed when MPI Finalize was issued, a quality application should free all MPI resources before calling MPI_Finalize. Listing information for these requests:  -Request 1: Point-to-point request activated at reference 1	Representative location: <b>MPI_Irecv</b> (1st occurrence) called from: #0 main@example-fix4.c:17	References of a representative process: reference 1 rank 1: <b>MPI_Irecv</b> (1st occurrence) called from: #0 main@example-fix4.c:17

- User defined objects include
  - MPI\_Comms (even by MPI\_Comm\_dup)
  - MPI\_Datatypes
  - MPI\_Groups

Unfinished non-blocking receive is resource leak and missing synchronization

Leak of user defined datatype object

## Fix6: Deallocate datatype object

## Fix7: use MPI\_Wait to finish asynchronous communication

```
#include <mpi.h>
#include <stdio.h>

int main (int argc, char** argv)
{
    int rank, size, buf[8];

    MPI_Init (&argc, &argv);
    MPI_Comm_rank (MPI_COMM_WORLD, &rank);
    MPI_Comm_size (MPI_COMM_WORLD, &size);

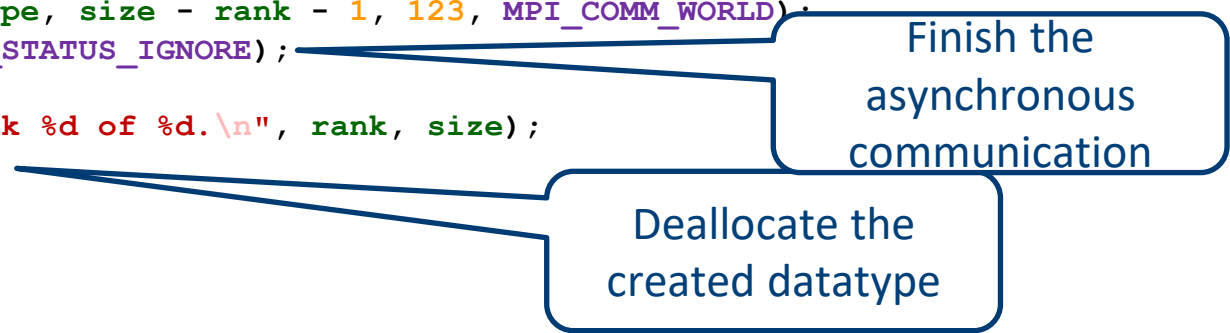
    MPI_Datatype type;
    MPI_Type_contiguous (2, MPI_INT, &type);
    MPI_Type_commit (&type);

    MPI_Request request;
    MPI_Irecv (buf, 2, MPI_INT, size - rank - 1, 123, MPI_COMM_WORLD, &request);

    MPI_Send (buf + 4, 1, type, size - rank - 1, 123, MPI_COMM_WORLD);
    MPI_Wait (&request, MPI_STATUS_IGNORE);

    printf ("Hello, I am rank %d of %d.\n", rank, size);
    MPI_Type_free (&type);

    MPI_Finalize ();
    return 0;
}
```



The diagram consists of two blue-outlined callout boxes with white text. The first callout box, containing the text "Finish the asynchronous communication", has a pointer line that originates from the `MPI_Wait (&request, MPI_STATUS_IGNORE);` line in the code and points to the center of the box. The second callout box, containing the text "Deallocate the created datatype", has a pointer line that originates from the `MPI_Type_free (&type);` line in the code and points to the center of the box.

# Finally

Rank(s)	Type	Message	
	Information	MUST detected no MPI usage errors nor any suspicious behavior during this application run.	
Details:			
	Message	From	References
	MUST detected no MPI usage errors nor any suspicious behavior during this application run.		

No further error detected

Hopefully this message applies to many applications



# Content

---

- Motivation
- MPI usage errors
- Examples: Common MPI usage errors
  - Including MUST's error descriptions
- **Correctness tools**
- MUST usage
- Hands-on

## Tool Overview - Approaches Techniques

---

- Debuggers:
  - ✓ Helpful to pinpoint any error
  - Finding the root cause may be hard
  - Won't detect sleeping errors
  - E.g.: gdb, TotalView, Allinea DDT
- Static Analysis:
  - Compilers and Source analyzers
  - ✓ Typically: type and expression errors
  - E.g.: MPI-Check
- Model checking:
  - Can find hidden errors
  - Requires a model of your applications
  - State explosion possible
  - E.g.: MPI-Spin

```
MPI_Recv (buf, 5, MPI_INT,  
1,  
123, MPI_COMM_WORLD, &status);
```

“-1” instead of “MPI\_ANY\_SOURCE”

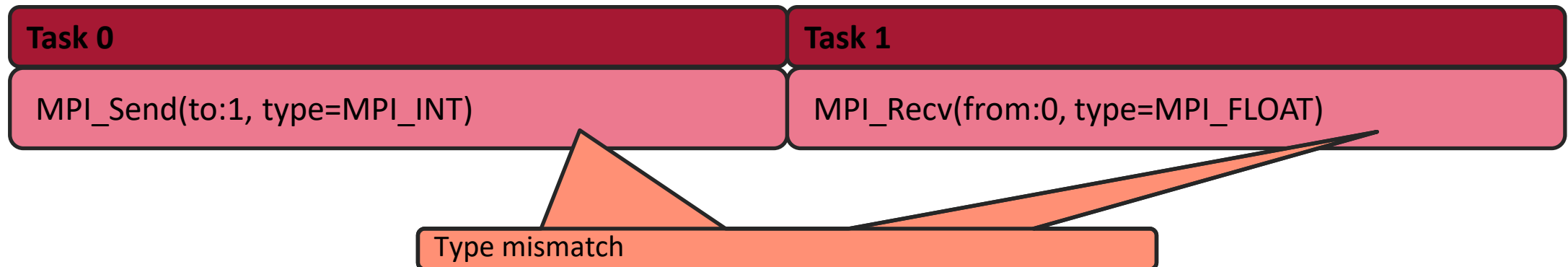
```
if (rank == 1023)  
crash ();
```

Only works with less than 1024 tasks

## Tool Overview - Approaches Techniques (2)

---

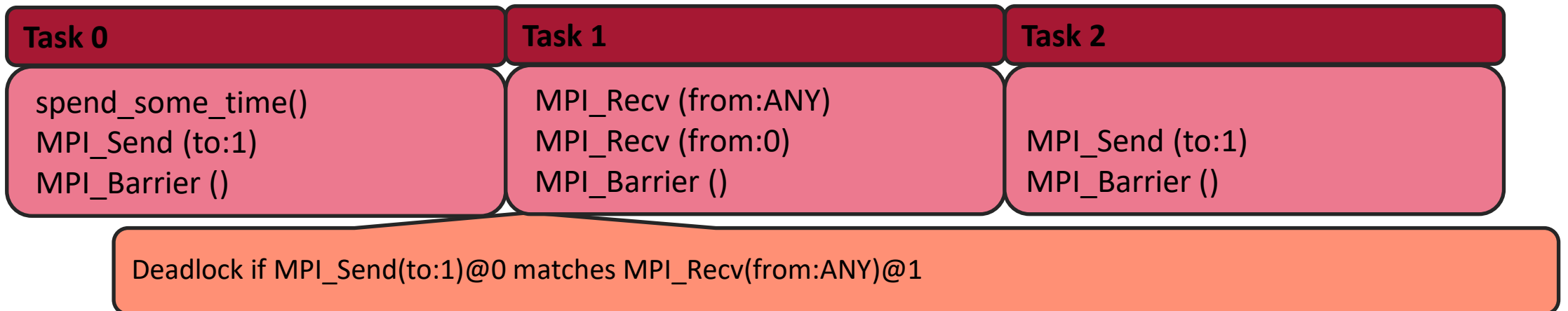
- Runtime error detection:
  - ✓ Inspect MPI calls at runtime
  - Limited to the timely interleaving that is observed
  - Causes overhead during application run
  - E.g.: Intel Trace Analyzer, Umpire, Marmot, MUST



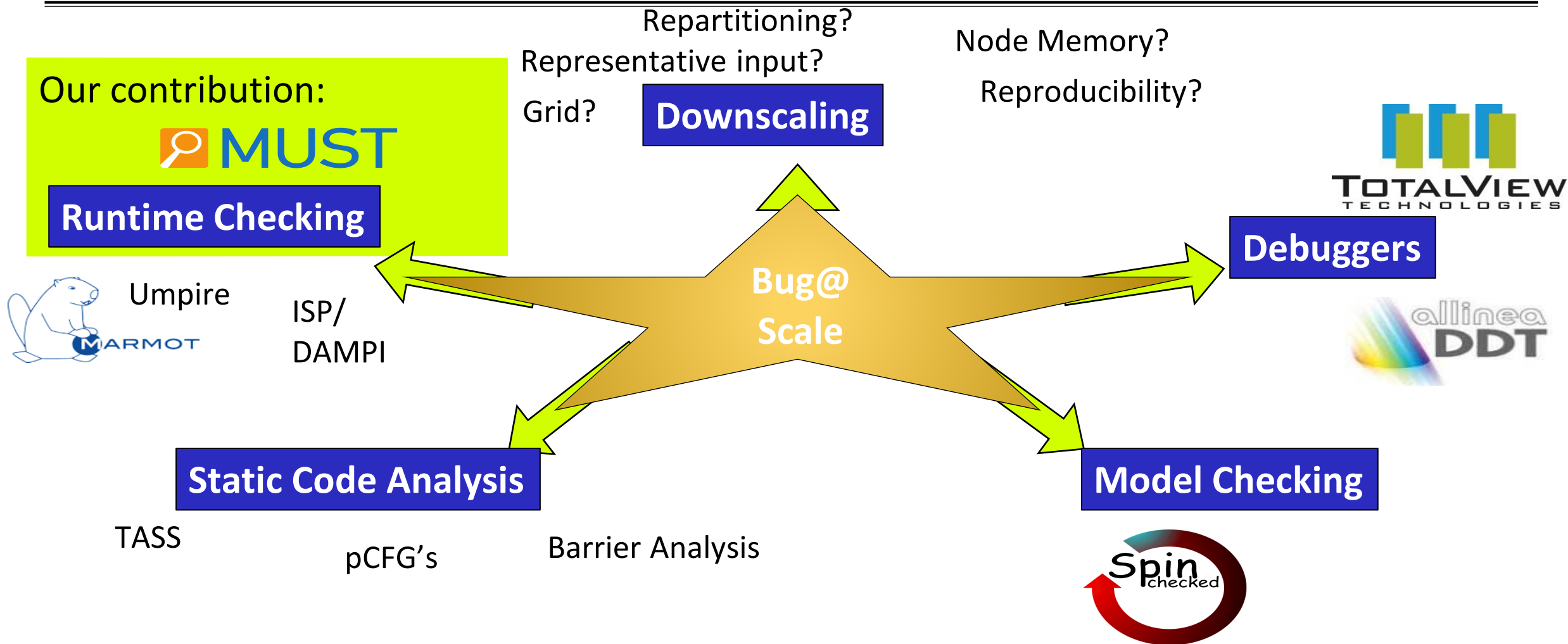
## Tool Overview - Approaches Techniques (3)

---

- Formal verification:
  - Extension of runtime error detection
  - Explores all relevant interleavings (explore around nondet.)
  - Detects errors that only manifest in some runs
  - Possibly many interleavings to explore
  - E.g.: ISP



# Approaches to Remove Bugs (Selection)



# Content

---

- Motivation
- MPI usage errors
- Examples: Common MPI usage errors
  - Including MUST's error descriptions
- Correctness tools
- **MUST usage**
- Hands-on

## MUST - Basic Usage

---

- Apply MUST as an mpiexec wrapper, that's it:

```
% mpicc source.c -o exe  
% $MPIRUN -n 4 ./exe
```

```
% mpicc -g source.c -o exe  
% mustrun --must:mpiexec $MPIRUN -n 4 ./exe
```

or simply

```
% mustrun -n 4 ./exe
```

- After run: inspect "MUST\_Output.html"
- "mustrun" (default config.) uses an extra process:
  - I.e.: "mustrun -np 4 ..." will use 5 processes
  - Allocate the extra resource in batch jobs!
  - Default configuration tolerates application crash; BUT is slower (details later)

## MUST - Usage on frontend - backend machines

---

- Compile and run using a batch system (only necessary for customized use cases)

```
% mpicc source.c -o exe  
% mpiexec -n 4 ./exe
```

```
% mpicc -g source.c -o exe  
% mustrun -n 4 ./exe --must:mode prepare  
% salloc mustrun -n 4 ./exe --must:mode run
```

- If you see messages about missing dot on the backend, run on frontend:

```
% mustrun --must:dot
```

- Open MUST\_Output.html with a browser



## MUST - At Scale (highly recommended for >10 processes)

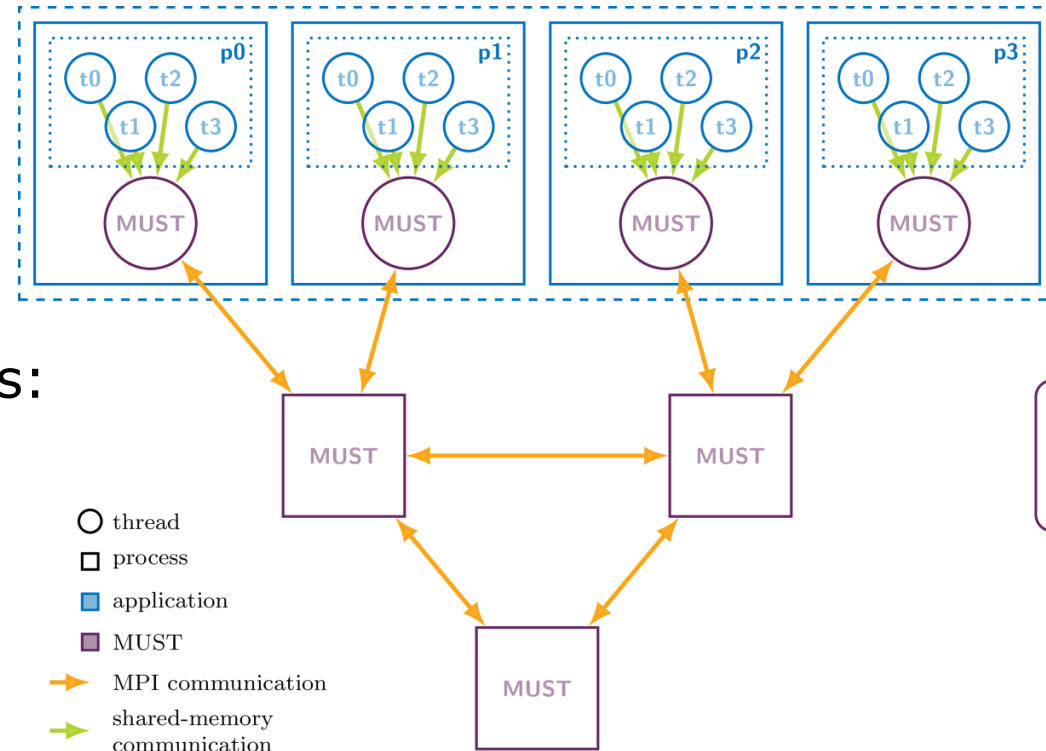
- Provide a branching factor (fan-in) for the tree infrastructure:

```
% mustrun -n 40 ./exe \  
--must:fanin 8
```

- Get info about the number of processes:

```
% mustrun -n 40 ./exe \  
--must:fanin 8 --must:info
```

→ This will give you the number of processes needed with tool attached



### Local analysis

- opaque handle usage
- local type matching
- data race detection
- buffer overlap detection

### Distributed analysis

- distributed deadlock analysis
- point-to-point matching
- p2p type matching
- collective matching

### Centralized analysis

- deadlock graph analysis
- HTML output

## MUST – Execution modes

---

	Application might crash	Application never crashes
Centralized analysis	<div data-bbox="896 454 1648 618"></div> <ul style="list-style-type: none"><li>▪ 1 extra process</li><li>▪ Blocking communication</li></ul>	<div data-bbox="1694 454 2446 618"><code>--must:nocrash</code></div> <ul style="list-style-type: none"><li>▪ 1 extra process</li><li>▪ Non-blocking communication</li></ul>
Distributed analysis	<div data-bbox="896 848 1648 1012"><code>--must:nodesize 8</code></div> <ul style="list-style-type: none"><li>▪ 1 extra process per 7 application processes + tree</li><li>▪ Nodesize must be divisor of ranks sharing memory</li></ul>	<div data-bbox="1694 848 2446 1012"><code>--must:fanin 8</code></div> <ul style="list-style-type: none"><li>▪ 1 extra process per 8 app processes + tree</li></ul>

## MUST - Multithreading Support

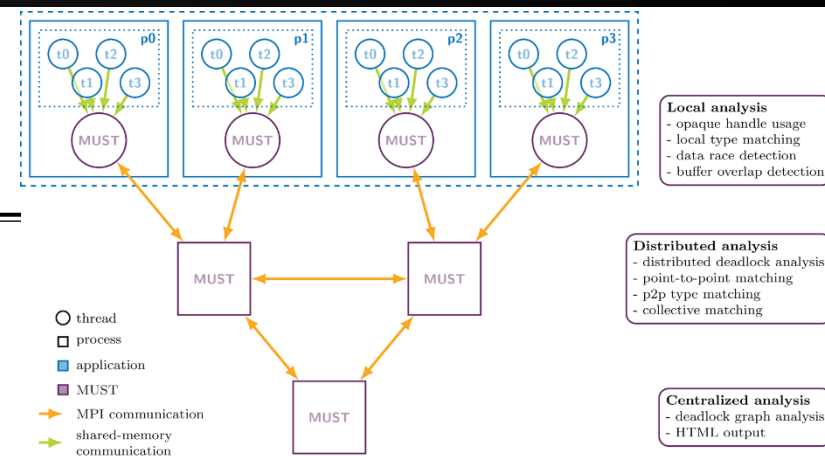
- By default, MUST supports up to `MPI_THREAD_FUNNELED`
- For higher threading levels:

```
% mustrun -n 40 ./exe --must:hybrid
```

- This will raise the required level to `MPI_THREAD_MULTIPLE`!
- Some MPI might need env like: `MPICH_MAX_THREAD_SAFETY=multiple`
- Get info about the resources needed:

```
% mustrun -n 40 ./exe --must:hybrid --must:info
```

→ This will give you the number of processes needed with tool attached



## Stacktraces in MUST

- We use Dyninst or Backward-cpp as an external lib for stacktraces
- Dyninst has officially no support for Intel compiler
  - But: in most cases it's compatible to icc compiled C/C++ applications
- Collecting stack traces can be costly. Select with  
`--must:stacktrace [dyninst|backward|none]`
- Supposed your application has no faults you won't need stacktraces  
 ☺

From
Representative location: <b>MPI_Init_thread</b> (1st occurrence) called from: #0 MAIN_@bt.f:90 #1 main@bt.f:319
Representative location: <b>MPI_Comm_split</b> (1st occurrence) called from: #0 MAIN_@bt.f:90 #1 main@bt.f:319

Rank(s)	Type	Message	From	References
	Information	MUST detected no MPI usage errors nor any suspicious behavior during this application run.		

## MUST modules on juwels

---

---

```
$ module spider MUST/1.7.2
```

```
-----  
MUST: MUST/1.7.2  
-----
```

```
Stages/2022 GCC/11.2.0 OpenMPI/4.1.1  
Stages/2022 GCC/11.2.0 OpenMPI/4.1.2  
Stages/2022 GCC/11.2.0 ParaStationMPI/5.5.0-1  
Stages/2022 Intel/2021.4.0 OpenMPI/4.1.1  
Stages/2022 Intel/2021.4.0 OpenMPI/4.1.2  
Stages/2022 Intel/2021.4.0 ParaStationMPI/5.5.0-1  
Stages/2022 Intel/2021.4.0 ParaStationMPI/5.5.0-1-mt
```

## MUST - Summary

---



- MPI runtime error detection tool
- Open source (BSD license)  
<http://www.itc.rwth-aachen.de/MUST/>
- Wide range of checks, strength areas:
  - Overlaps in communication buffers
  - Errors with derived datatypes
  - Deadlocks
- Largely distributed, able to scale with the application

# Content

---

- Motivation
- MPI usage errors
- Examples: Common MPI usage errors
  - Including MUST's error descriptions
- MUST usage
- **Hands-on**

## Hands-on (part 1)

```
$ source $PROJECT_training2123/setup.sh
$ cd $SCRATCH_training2123/$USER
$ tar zxvf $PROJECT_training2123/examples/must-examples.tgz
```

- Compile with debugflag
- Load must module
- Run interactive debug session or submit job

```
$ cd $SCRATCH_training2123/$USER/must-examples
$ module load MUST/1.7.3
$ mpicc -g example.c -o example
$ mustrun -n 4 ./example
```

- Launch with mustrun
- Abort with Ctrl + C
- Open MUST\_Output.html in browser (e.g., firefox)

```
[MUST] MUST configuration ... centralized checks with fall-back
application crash handling (very slow)
...
[MUST] Executing application:
=====MUST=====
ERROR: MUST detected a deadlock, detailed information is available in
the MUST output file. You should either investigate details with a
debugger or abort, the operation of MUST will stop from now.
=====
^C
[MUST] Execution finished, inspect ".../MUST_Output.html"!
$ firefox MUST_Output.html
```



## Hands-on (part 1)

```
$ source $PROJECT_training2123/setup.sh
$ cd $SCRATCH_training2123/$USER
$ tar zxvf $PROJECT_training2123/examples/must-examples.tgz
```

- Fix the identified issue(s)
- Recompile
- Launch with mustrun
- Choose stdout output if no browser is available
- Abort with Ctrl + C
- Start over

```
$ cd ~/corr-handson
$ vim example.c
$ mpicc -g example.c -o example
$ mustrun --must:output stdout -n 4 ./example
```

Some MPI implementations check for errors and abort.

```
[MUST] MUST configuration ... centralized links with fall-back
application crash handling (very slow)
[MUST] Using prebuilt infrastructure at .../modules/prebuilds/26e361a...
...
[MUST] Executing application:
Invalid datatype, error stack:
MPI_Send(201): MPI_Send(buf=0x7fffe77c58cc, count=2,
dtype=USER<contig>, dest=1, tag=123, comm=0x84000004) failed
MPI_Send(135): Datatype has not been committed
...
Waiting up to 30 seconds for analyses to be finished.
^C
```

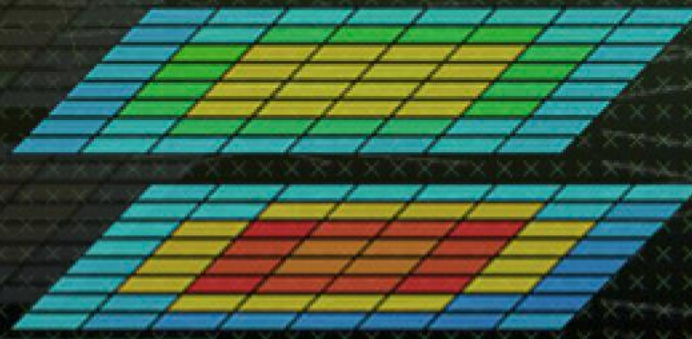
## Hands-on (part 2)

```
$ source $PROJECT_training2123/setup.sh
$ cd $SCRATCH_training2123/$USER
$ tar zxvf $PROJECT_training2123/examples/must-examples.tgz
```

- Apply MUST to tea\_leaf binary
- Launch with mustrun
- Open MUST\_Output.html in browser
- What issues did MUST find for tea\_leaf?

```
$ cd $SCRATCH_training2123/$USER/TeaLeaf_CUDA
$ make COMPILER=PGI
$ mustrun -n 7 ./tea_leaf
```

```
[MUST] MUST configuration ... centralized checks with fall-back
application crash handling (very slow)
[MUST] Using prebuilt infrastructure at ../modules/prebuilds/26e361a...
...
[MUST] Execution finished, inspect ".../MUST_Output.html"!
$ firefox MUST_Output.html
```



# Thank You