# The Distributed Application Debugger (DAD)

Michael Q. Jones & Matt B. Pedersen
University of Nevada Las Vegas

# Introduction

- The Distributed Application Debugger is a debugging tool for parallel programs
- Targets the MPI platform
- Runs remotley even on private networks
- Has record and replay features.
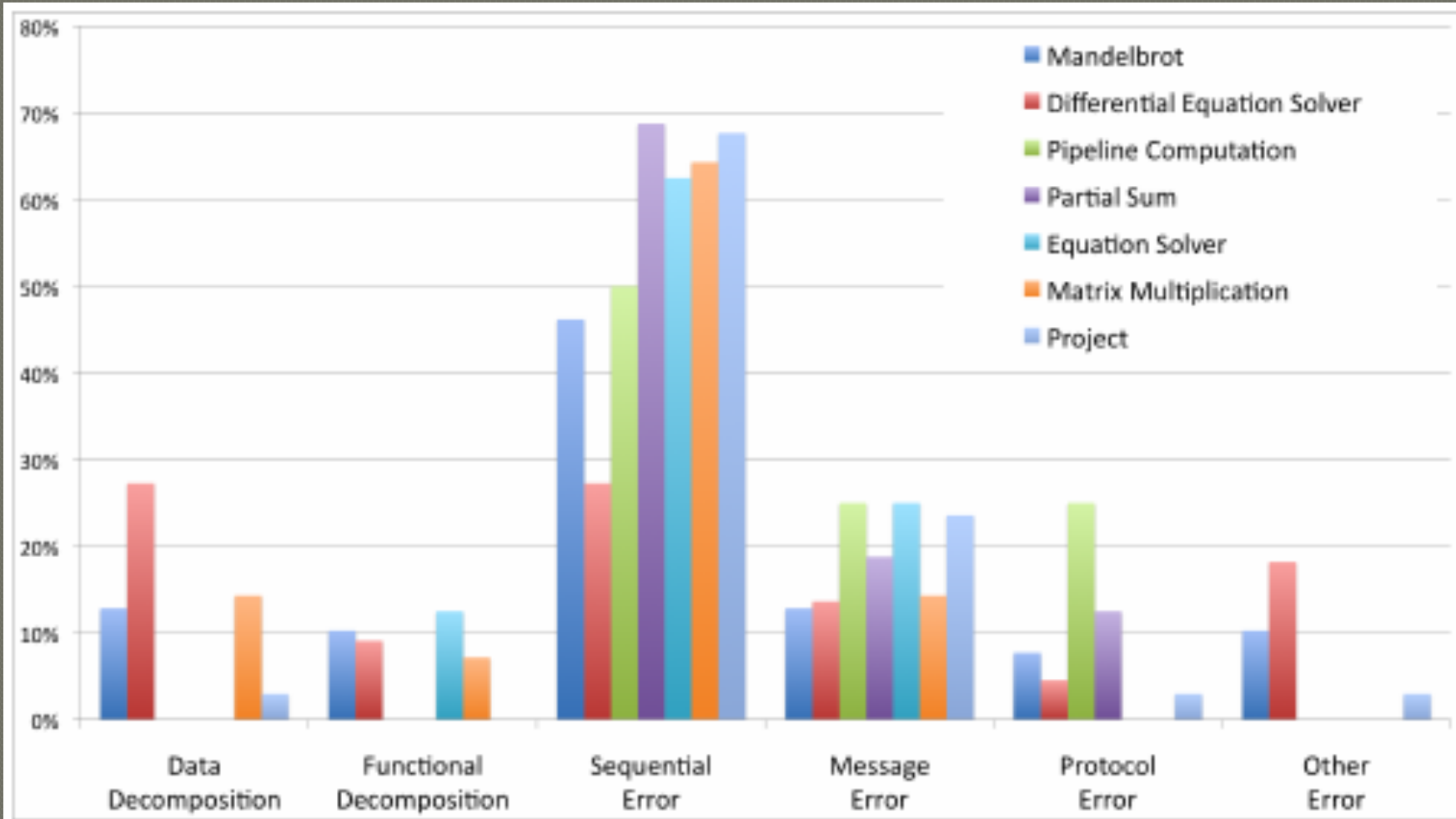- Integrates GDB

# Why a new tool?

- Results from survey of students learning parallel programming concluded 3 things:
  - 1. Sequential errors are still frequent
  - 2. Message errors are time consuming
  - 3. Print statements are still used for debugging

# Survey

- Survey results categorized according to the domains of multilevel debugging
  - Sequential errors
  - Message errors
  - Protocol errors
- In addition to
  - Data decomposition errors
  - Functional decomposition errors

# Survey Results

# Survey Results

| | Data Decomp. | Functional Decomp. | Sequential Error | Message Error | Protocol Error | Other |
|---|---|---|---|---|---|---|
| Average Time | 19.9 | 68.1 | 24.3 | 61.4 | 50.0 | 30.6 |
| Total Time Spent | 278 | 545 | 1,846 | 1,536 | 451 | 545 |
| # Errors | 14 | 8 | 76 | 25 | 9 | 8 |
| Total time Spent in % | 5.67% | 11.12% | 37.67% | 31.34% | 9.20% | 5.0% |

# The Components

- ## The Client
  - The GUI interacting with the programmer
- ## The Call Center
  - A central messaging hub (running on the cluster) for
    - Routing messages from the MPI program to The Client
    - Routing commands from The Client to the MPI program
- ## Bridges
  - A relay application for passing data between The Client and The Call Center, when The Call Center is not directly accessible (cluster behind firewall)
- ## The Runtime
  - A libraries with wrapper code for the MPI functions (talks to The Call Center)

# The Setup

Home

Firewall

Login Server

Cluster Login Server

Cluster

Login from Home to Cluster not Directly possible

# The Setup

Home

Firewall

Login Server

Cluster Login Server

Client

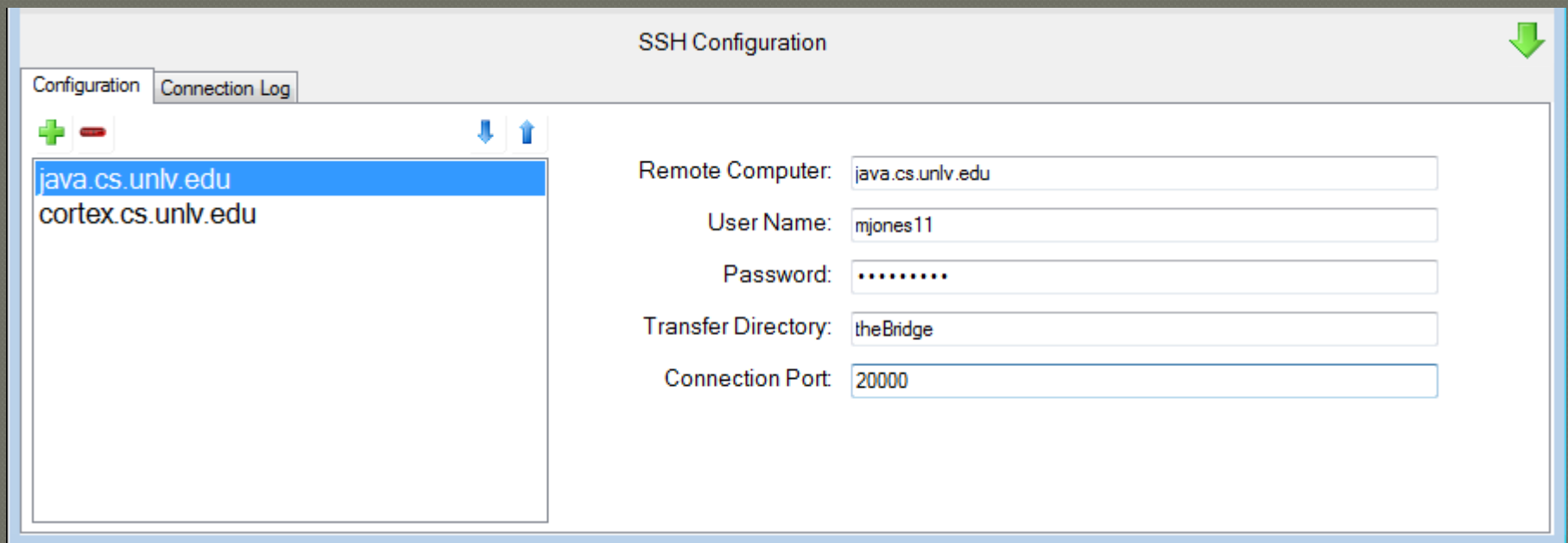Bridge

Bridge

Cluster

Call Center

MPI

- Client runs at home
- Bridges on the servers in between home and the cluster
- Call Center on the cluster
- MPI processes on the cluster

# The Distributed Application Debugger

# Client Side

- The user provides a connection path and credentials on all machines

# Client Side

- The user provides a connection path and credentials on all machines
- The system initiates SSH connections to each configured computer and launches a Bridge or The Call Center.
- Each component then connects to each other via TCP.

# Connection Initiated

# MPI Cluster Side

- Include a special mpi.h header file
- MPI calls are caught by wrapper functions
- Upon start up, each node creates a callback connection to The Call Center
- Data passed to MPI functions is sent back.

# MPI Code Redirected

```c
#include "mpi.h"


int main(int argc, char *argv[]){
  MPI_Init(&argc,&argv);
  MPI_Comm_size(MPI_COMM_WORLD,&numProcs);
  MPI_Comm_rank(MPI_COMM_WORLD,&myId);

  if(myId == 0){
    //This is the master node.
    //Send the whole buffer to the middle index
    MPI_Send(transferBuffer, numSlaves, MPI_INT,
      startingNodeId, numSlaves, MPI_COMM_WORLD);

    //Sync and then distribute the intial values
    MPI_Barrier(MPI_COMM_WORLD);

    //Wait for the result
    MPI_Recv(&finalResult,1, MPI_INT,
      numSlaves, TAG, MPI_COMM_WORLD, &stat);

    //Send result and wait for everyone to get it
    MPI_Send(&finalResult, 1, MPI_INT,
     startingNodeId, TAG, MPI_COMM_WORLD);

    MPI_Barrier(MPI_COMM_WORLD);
  }
  else{
    //Distribute and process the partial sums
```
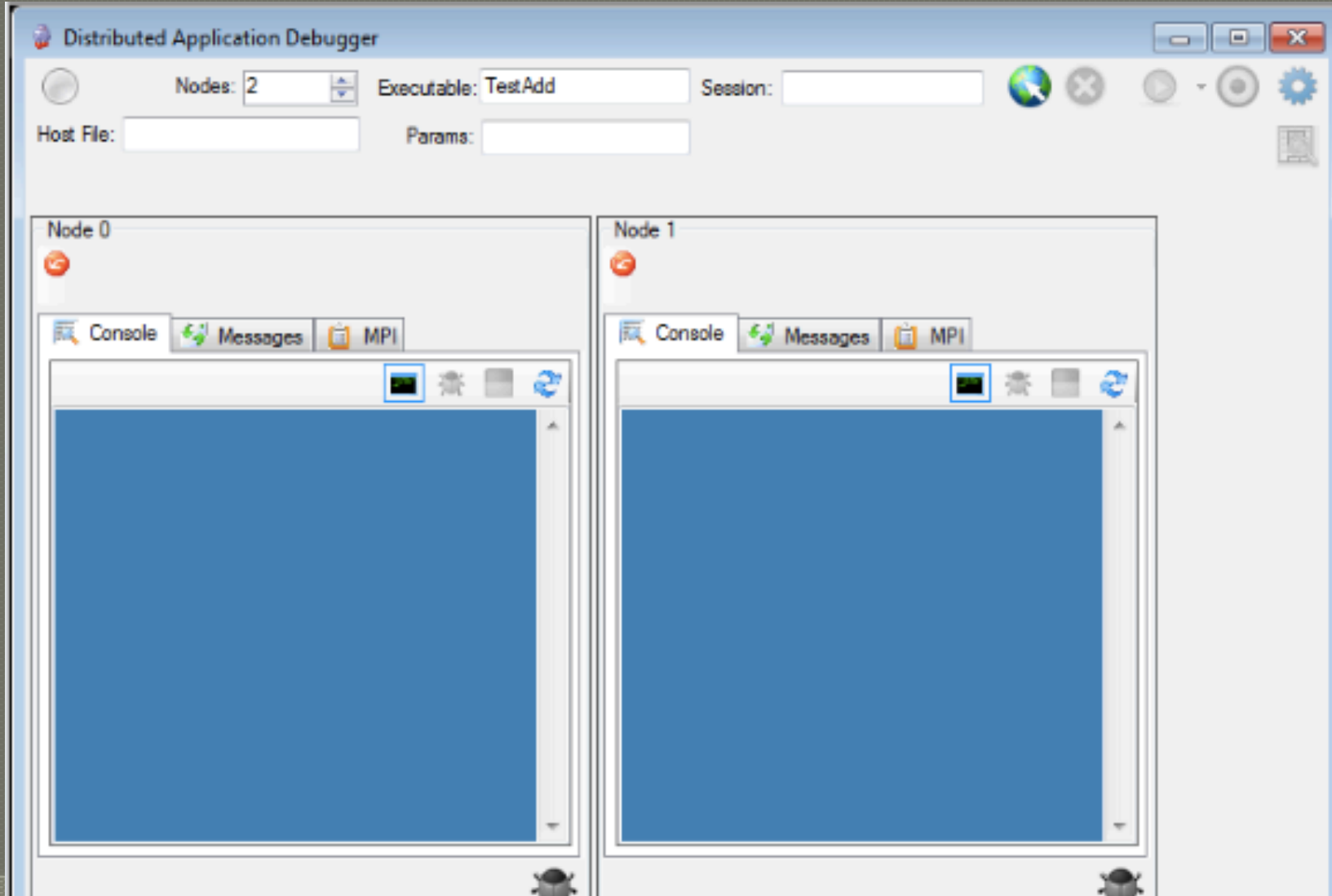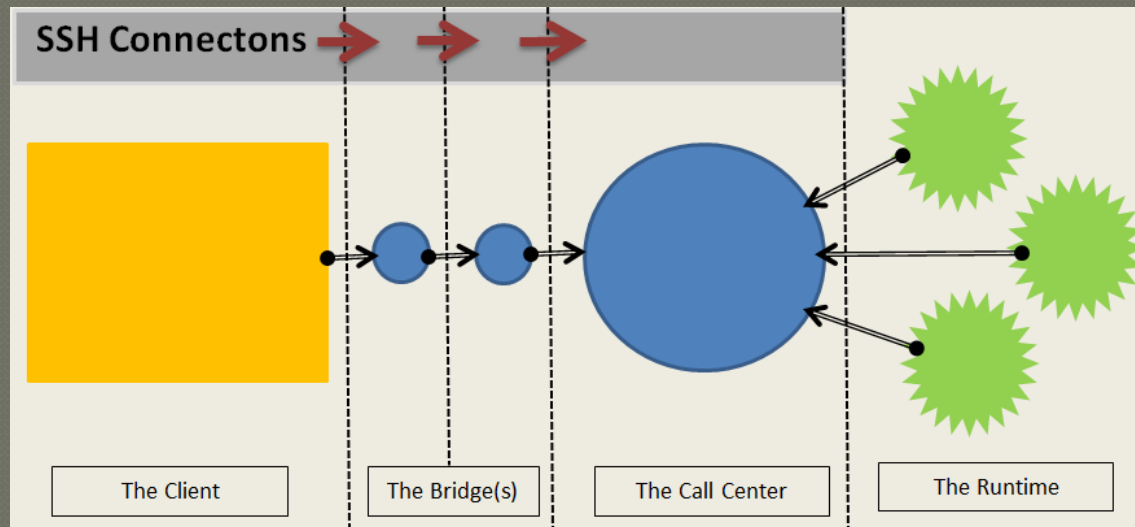
```c
#include <mpi.h>
#include "debug.h"
#include "mpidebug.h"

int main(int argc, char *argv[]){
  _MPI_Init(&argc,&argv);
  _MPI_Comm_size(MPI_COMM_WORLD,&numProcs);
  _MPI_Comm_rank(MPI_COMM_WORLD,&myId);

  if(myId == 0){
    //This is the master node.
    //Send the whole buffer to the middle index
    _MPI_Send(transferBuffer, numSlaves, MPI_INT,
      startingNodeId, numSlaves, MPI_COMM_WORLD);

    //Sync and then distribute the intial values
    _MPI_Barrier(MPI_COMM_WORLD);

    //Wait for the result
    _MPI_Recv(&finalResult,1, MPI_INT,
      numSlaves, TAG, MPI_COMM_WORLD, &stat);

    //Send result and wait for everyone to get it
    _MPI_Send(&finalResult, 1, MPI_INT,
     startingNodeId, TAG, MPI_COMM_WORLD);

    _MPI_Barrier(MPI_COMM_WORLD);
  }
  else{
    //Distribute and process the partial sums
```

# The Connected System

# The Connected System

# Session Modes

An MPI session can be run in 3 modes:
- Play
  - Just run like regular MPI
- Record (Record all messages)
  - Record all messages
- Replay
  - Use recorded messages to play back

# Session Modes (Play)

- The Runtime behaves like regular MPI
  - Nothing is saved to disk
  - Nothing is read from disk
  - Messages and parameters ARE sent back to The Client

# Session Modes (Record)

- ## The Runtime
  - Saves messages and parameters to a log file
  - Executes the actual MPI call
  - Saves the result

# Session Modes (Replay)

- The Runtime does not execute any real MPI calls.
  - All data is supplied from log files.
  - No actual communication takes place
  - Guarantees the same run as when the log file was recorded

# Session Mode (Replay – Mixed)

- Mixed mode is special
  - Some processes execute real MPI calls
  - Some replay from log file
    - Sometimes its necessary to execute MPI calls if communicating with someone who is executing real MPI calls; E.g. to avoid buffer overflow
    - Validation is done on real values and log file values

# Debugging Data

- The Runtime sends back 2 debugging messages per MPI command
  - A *PRE* message indicating that an MPI command is about to be executed
  - A *POST* message indicating that an MPI command completed
- Console messages are routed per node to the appropriate window.

# Analyzing Data

- Debugging data gets displayed within the Console, Messages, or MPI tabs

# Analyzing Data

- The Console Tab displays anything that the user's code wrote to **stdout**.

# Analyzing Data

- The Messages Tab displays messages as they come
- Matches Send/ Receive pairs between nodes.
- Messages without a corresponding Send or Receive message get highlighted in red.

# Analyzing Data

- The MPI tab displays all MPI commands
  - in the order they were executed
  - along with their parameters.
- Commands statuses (success, fail, or blocked) are displayed with icons in the Status Column.

Node 0

Host: cortex.cs.unlv.edu  Process Id: 6476

| Console | Messages | MPI |

Count: 13                    Commands (12) ▾

| Id | Command | Line # | Status |
|----|---------|--------|--------|
| 0 | MPI_INIT | 14 | ✓ |
| 1 | MPI_RANK | 15 | ✓ |
| 2 | MPI_SIZE | 16 | ✓ |
| 3 | MPI_RECV | 34 | ✓ |
| 4 | MPI_RECV | 34 | ✓ |
| 5 | MPI_RECV | 34 | ✓ |
| 6 | MPI_RECV | 34 | ✓ |
| 7 | MPI_RECV | 34 | ✓ |
| 8 | MPI_RECV | 34 | ✓ |
| 9 | MPI_RECV | 34 | ✓ |
| 10 | MPI_RECV | 34 | ✓ |
| 11 | MPI_RECV | 34 | ✓ |
| 12 | MPI_RECV | 34 | ? |

Command Details

OriginatorId: 0
CommandId: 12
LineNum: 34
Command: MPI_RECV
Buf:
Count: 1
Datatype: MPI_INT
Src: 1
Tag: 0
Comm: MPI_COMM_WORLD

# Analyzing Data

# Analyzing Data

- Buffer values can be requested and inspected.

# Attaching GDB
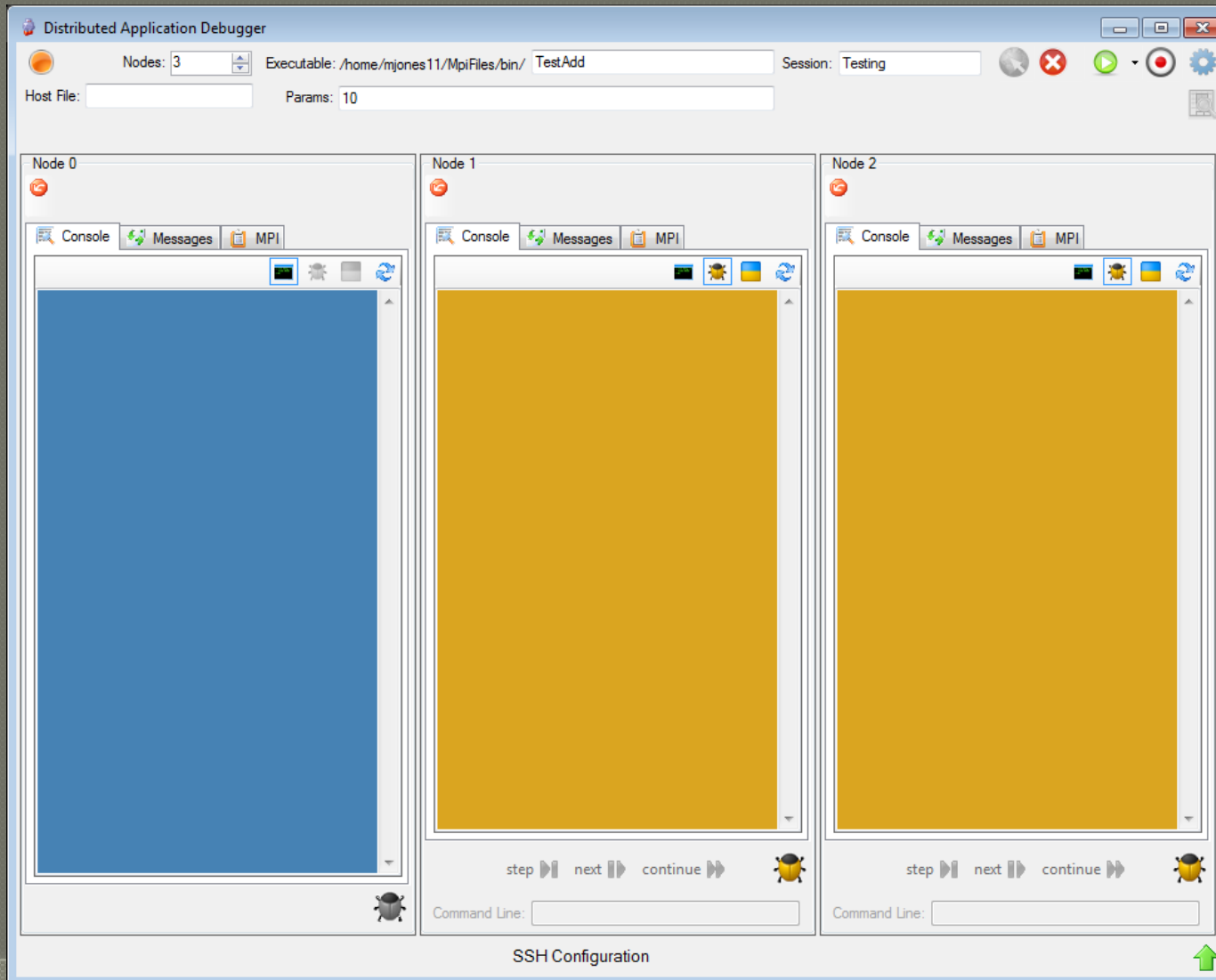
- GDB can be attached to any node and controlled with the GDB Control Panel.

# Attaching GDB

# Source Code

- The source code to The Distributed Application Debugger can be found on GitHub at:


- https://github.com/mjones112000/ DistributedApplicationDebugger

# Questions??