

Supercomputer Design Through Simulation

Cray User Group (CUG) Meeting
Lugano, Switzerland

Rolf Riesen

`rolf@cs.sandia.gov`

Sandia National Laboratories

May 9, 2006

Talk Overview

- 1 Goal
- 2 Design
- 3 Usage
- 4 Validation
- 5 Experiments
- 6 Comparison to Other Work
- 7 Future Work
- 8 Summary



Talk
Overview

Goal

Design

Usage

Validation

Experiments

Related Work

Future Work

Summary

End

Section Outline

- 1 Goal
- 2 Design
- 3 Usage
- 4 Validation
- 5 Experiments
- 6 Comparison to Other Work
- 7 Future Work
- 8 Summary



Talk
Overview

Goal

Design

Usage

Validation

Experiments

Related Work

Future Work

Summary

End

Goal

- Simulate a supercomputer; e.g., Red Storm, using federated discrete event simulators
- With enough fidelity to make future purchase and design decisions concerning things like:
 - CPU choice
 - Memory size and speed
 - Network interface
 - Topology
 - Application behavior
 - Research directions
 - etc.
- Created initial prototype with promising attributes
- This talk describes simulator
- Collective results on Thursday



Talk
Overview

Goal

Design

Usage

Validation

Experiments

Related Work

Future Work

Summary

End

Section Outline

1 Goal

2 Design

- Node (Application)
- MPI Wrapper Library
- MPI Communicators
- Virtual Time
- Linking

3 Usage

4 Validation

5 Experiments



Talk
Overview

Goal

Design

Node (Application)
MPI Wrappers
Communicators
Virtual Time
Linking

Usage

Validation

Experiments

Related Work

Future Work

Summary

End

Node (Application)

Hybrid simulator:

- App runs regularly and uses MPI to exchange data
- Each MPI send and receive generates an event to the network simulator
- Sim generates rcv events that are matched by clients
- Algorithm determines when and how to update virtual time on each node
- Use MPI wrappers and profiling interface
- Current network simulator uses simple model:

$$\Delta = \frac{s}{B} + L$$

Δ network delay
 s message size

B network bandwidth
 L network latency



Talk
Overview

Goal

Design

Node (Application)

MPI Wrappers

Communicators

Virtual Time

Linking

Usage

Validation

Experiments

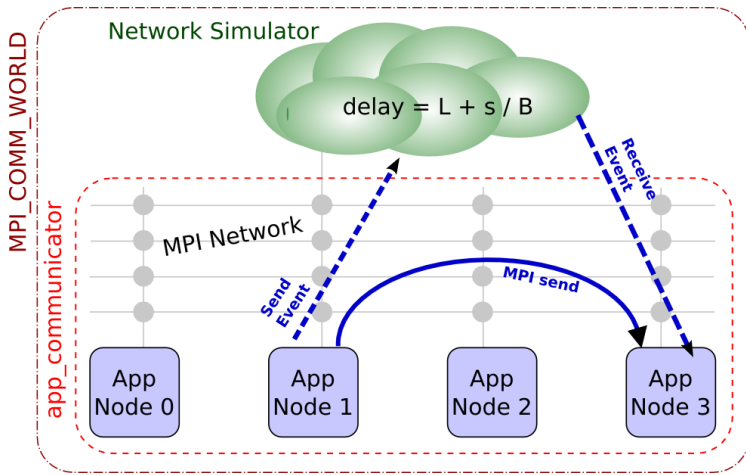
Related Work

Future Work

Summary

End

MPI Wrapper Library



Talk
Overview

Goal

Design

Node (Application)

MPI Wrappers

Communicators

Virtual Time

Linking

Usage

Validation

Experiments

Related Work

Future Work

Summary

End

MPI Wrapper Library

```
int MPI_Send(void *data,
             int len,
             MPI_Datatype dt,
             int dest,
             int tag,
             MPI_Comm comm)
{
    t_x = get_vtime();

    // Send the MPI message
    rc = PMPI_Send(data, len, dt, dest, tag, comm);

    // Send event to simulator
    event_send(t_x, len, dt, dest, tag);

    return rc;
}
```



Talk
Overview

Goal

Design

Node (Application)

MPI Wrappers

Communicators

Virtual Time

Linking

Usage

Validation

Experiments

Related Work

Future Work

Summary

End

MPI Wrapper Library

```
int MPI_Recv(void *data, int len, MPI_Datatype dt, int src,
             int tag, MPI_Comm comm, MPI_Status *stat)
{
    t1 = get_vtime();

    // Receive the MPI message
    rc = PMPI_Recv(data, len, dt, src, tag, comm, stat);

    // Wait for the matching event
    event_wait(&tx, &Δ, stat->MPI_TAG, stat->MPI_SOURCE);

    if (tx + Δ > t1)
        t3 = tx + Δ;
    else
        t3 = t1;
    set_vtime(t3); // Adjust virtual time
    return rc;
}
```



Talk
Overview

Goal

Design

Node (Application)

MPI Wrappers

Communicators

Virtual Time

Linking

Usage

Validation

Experiments

Related Work

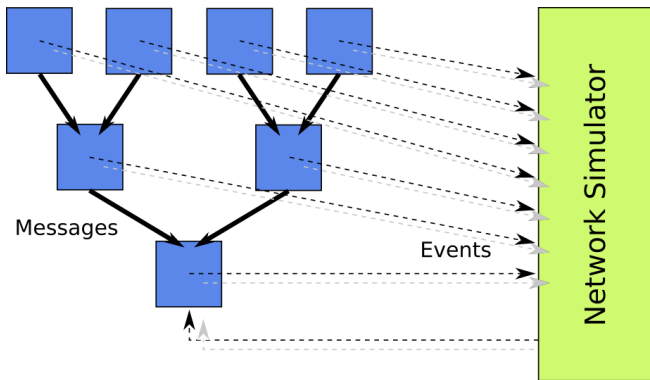
Future Work

Summary

End

MPI Wrapper Library

Event traffic for collectives



Talk
Overview

Goal

Design

Node (Application)

MPI Wrappers

Communicators

Virtual Time

Linking

Usage

Validation

Experiments

Related Work

Future Work

Summary

End

MPI Communicators

- Simulator framework sets up communicator for application nodes only
- **MPI_COMM_WORLD** covers application and simulator
- Wrappers swap **MPI_COMM_WORLD** with internal communicator when application calls MPI
- Application never sees real **MPI_COMM_WORLD**



Talk
Overview

Goal

Design

Node (Application)

MPI Wrappers

Communicators

Virtual Time

Linking

Usage

Validation

Experiments

Related Work

Future Work

Summary

End

Virtual Time

```
if ( $t_x + \Delta > t_1$ )  
     $t_3 = t_x + \Delta$ ;  
else  
     $t_3 = t_1$ ;  
set_vtime( $t_3$ );
```

- If message was sent earlier than we started looking for it, we have to assume it was already here
 - Just “erase” the time we spent actually receiving it
- If message arrived after we started waiting for it, use the virtual send time + Δ to set local virtual clock



Talk
Overview

Goal

Design

Node (Application)

MPI Wrappers

Communicators

Virtual Time

Linking

Usage

Validation

Experiments

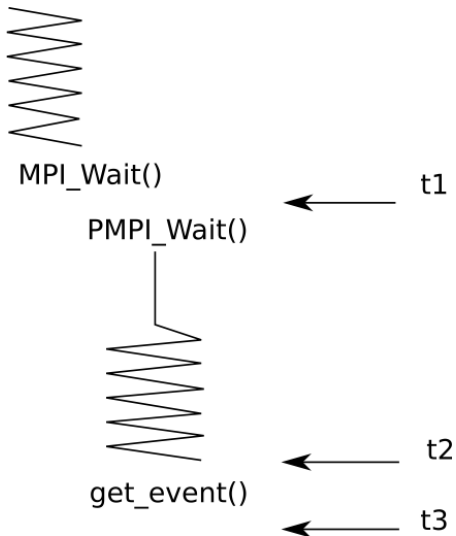
Related Work

Future Work

Summary

End

Virtual Time



Talk
Overview

Goal

Design

Node (Application)

MPI Wrappers

Communicators

Virtual Time

Linking

Usage

Validation

Experiments

Related Work

Future Work

Summary

End



- Currently need to rename **main()** to **main_node()**
- Should not be necessary when we use **MPI_Init()**
- In Fortran programs **program** has to be changed to **subroutine main_node**
- **No Changes to application are necessary!**

Talk
Overview

Goal

Design

Node (Application)

MPI Wrappers

Communicators

Virtual Time

Linking

Usage

Validation

Experiments

Related Work

Future Work

Summary

End

Section Outline

- 1 Goal
- 2 Design
- 3 Usage**
- 4 Validation
- 5 Experiments
- 6 Comparison to Other Work
- 7 Future Work
- 8 Summary



Talk
Overview

Goal

Design

Usage

Validation

Experiments

Related Work

Future Work

Summary

End



- Two steps:
 - Create point-to-point model
 - Create collective model
- Measure two-node latency curve and write function to model it
- Measure all-to-all performance and write model

Talk
Overview

Goal

Design

Usage

Validation

Experiments

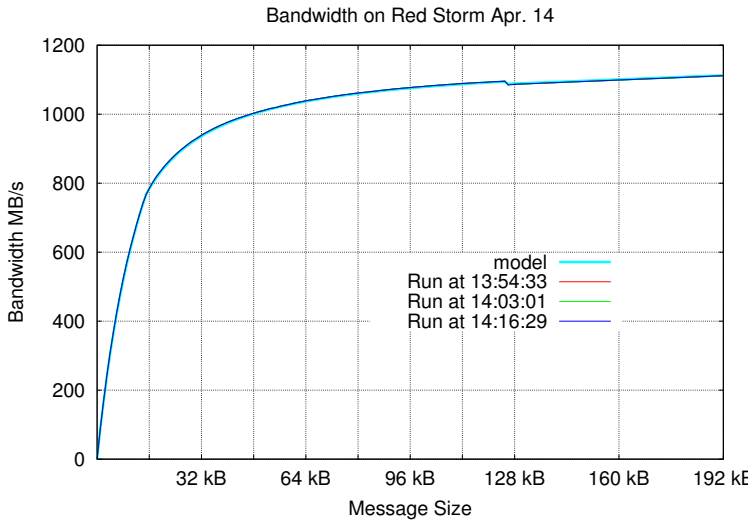
Related Work

Future Work

Summary

End

Usage



Talk
Overview

Goal

Design

Usage

Validation

Experiments

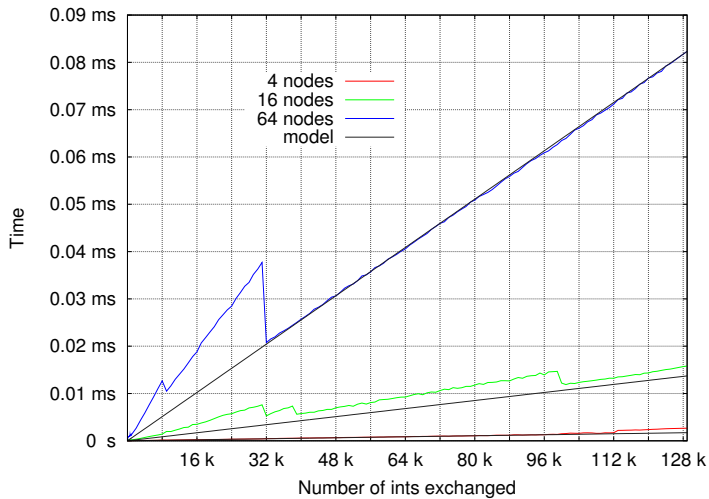
Related Work

Future Work

Summary

End

Usage



Talk
Overview

Goal

Design

Usage

Validation

Experiments

Related Work

Future Work

Summary

End

Section Outline

- 1 Goal
- 2 Design
- 3 Usage
- 4 Validation**
- 5 Experiments
- 6 Comparison to Other Work
- 7 Future Work
- 8 Summary



Talk
Overview

Goal

Design

Usage

Validation

Experiments

Related Work

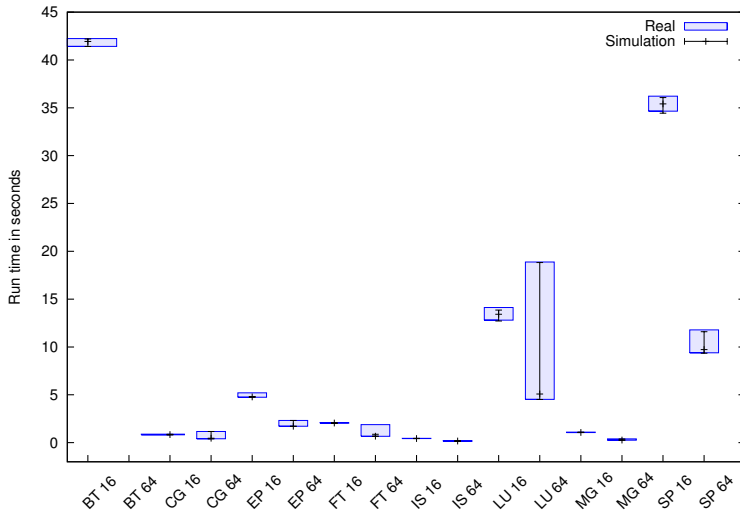
Future Work

Summary

End

Validation

NAS Class A Run Times



Talk
Overview

Goal

Design

Usage

Validation

Experiments

Related Work

Future Work

Summary

End

Section Outline

- 1 Goal
- 2 Design
- 3 Usage
- 4 Validation
- 5 Experiments**
 - Communication Patterns
 - Varying Bandwidth and Latency
 - Zero-Cost Collectives
 - Intrusion-Free MPI Traces

6 Comparison to Other Work



Talk
Overview

Goal

Design

Usage

Validation

Experiments

Patterns

Varying

Collectives

MPI Traces

Related Work

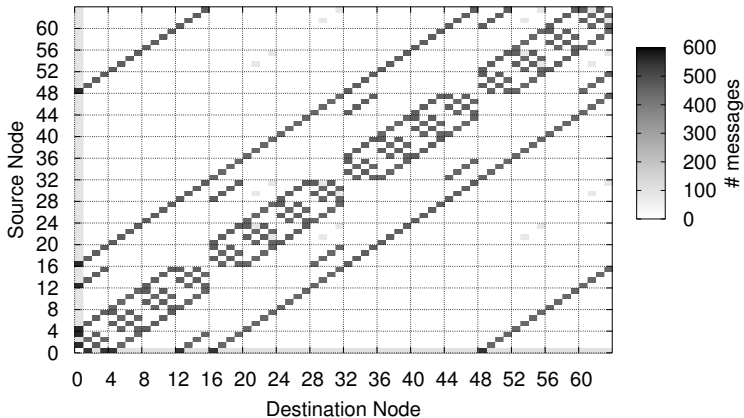
Future Work

Summary

End

Communication Patterns

MG (class B) message density distribution



Talk
Overview

Goal

Design

Usage

Validation

Experiments

Patterns

Varying

Collectives

MPI Traces

Related Work

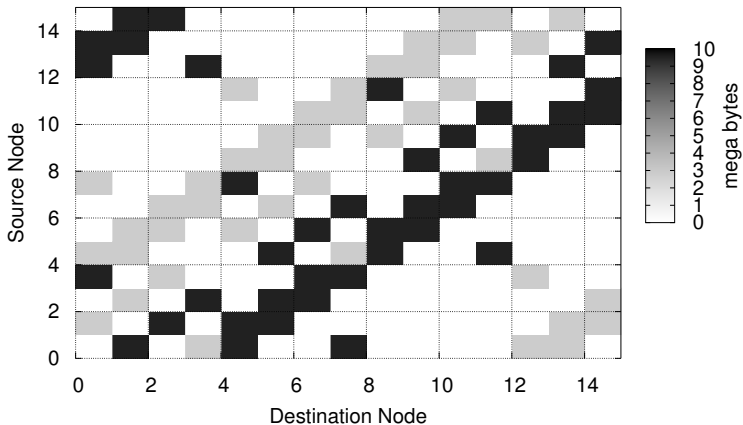
Future Work

Summary

End

Communication Patterns

BT (class A) data density distribution



Talk
Overview

Goal

Design

Usage

Validation

Experiments

Patterns

Varying

Collectives

MPI Traces

Related Work

Future Work

Summary

End

Varying Bandwidth and Latency

- Simulator can change bandwidth and latency independently
- This can be used to evaluate application performance under varying network characteristics
- → predict impact of new network



Talk
Overview

Goal

Design

Usage

Validation

Experiments

Patterns

Varying

Collectives

MPI Traces

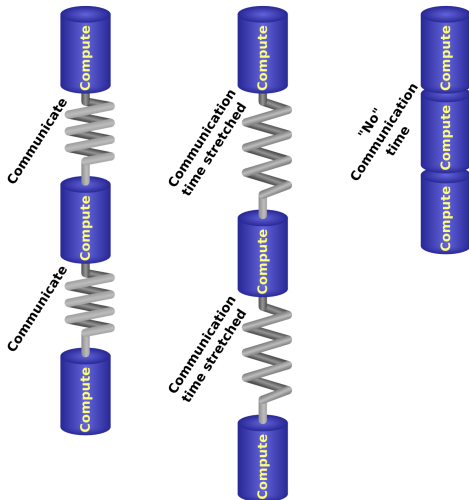
Related Work

Future Work

Summary

End

Varying Bandwidth and Latency



Talk
Overview

Goal

Design

Usage

Validation

Experiments

Patterns

Varying

Collectives

MPI Traces

Related Work

Future Work

Summary

End

Zero-Cost Collectives

- Putting collectives into NIC, building specialized NIC, or optimizing them is interesting
- How much application performance can be gained is not clear
- Simulator can assign $\Delta = 0$ to collectives and leave point-to-point alone
- See talk on Thursday



Talk
Overview

Goal

Design

Usage

Validation

Experiments

Patterns

Varying

Collectives

MPI Traces

Related Work

Future Work

Summary

End

Intrusion-Free MPI Traces

- So far gathered only limited amounts of data
- Simulator can gather, and save to disk, large amount of data
 - Without changing application virtual time



Talk
Overview

Goal

Design

Usage

Validation

Experiments

Patterns

Varying

Collectives

MPI Traces

Related Work

Future Work

Summary

End

Section Outline

- 1 Goal
- 2 Design
- 3 Usage
- 4 Validation
- 5 Experiments
- 6 Comparison to Other Work**
- 7 Future Work
- 8 Summary



Talk
Overview

Goal

Design

Usage

Validation

Experiments

Related Work

Future Work

Summary

End

Comparison to Other Work

- This approach seems to be new
- Combines low-intrusion measurement research with discrete event simulation
- Needs more validation, but seems to be very accurate
- Opens up many different and simple ways of evaluating applications and research directions



Talk
Overview

Goal

Design

Usage

Validation

Experiments

Related Work

Future Work

Summary

End

Comparison to Other Work

- No instrumentation code inserted into app
 - Rename `main()` (`program`) only change to app
- No disturbance of (virtual) runtime of app
 - Independent of amount of data collected.
- No extra memory needed on compute nodes to store trace data
- Language independent (Fortran, Fortran 90 with MPI-2, and C)



Talk
Overview

Goal

Design

Usage

Validation

Experiments

Related Work

Future Work

Summary

End

Section Outline

- 1 Goal
- 2 Design
- 3 Usage
- 4 Validation
- 5 Experiments
- 6 Comparison to Other Work
- 7 Future Work**
- 8 Summary



Talk
Overview

Goal

Design

Usage

Validation

Experiments

Related Work

Future Work

Summary

End

Future Work



Continuing Work

- Need to incorporate more accurate network model
- This will allow simulation of congestion, and evaluation of topology choices, node allocation, etc.
- Move below MPI into NIC for more fine-grained simulation
- Incorporate non-network simulators; CPU and NIC sims

Talk
Overview

Goal

Design

Usage

Validation

Experiments

Related Work

Future Work

Summary

End

Section Outline

- 1 Goal
- 2 Design
- 3 Usage
- 4 Validation
- 5 Experiments
- 6 Comparison to Other Work
- 7 Future Work
- 8 Summary



Talk
Overview

Goal

Design

Usage

Validation

Experiments

Related Work

Future Work

Summary

End

Summary

- Novel tool to collect MPI data
- Language independent
- Only linking with application needed
- Virtual runtime of application is not changed
- Lots of future possibilities



- Talk
- Overview
- Goal
- Design
- Usage
- Validation
- Experiments
- Related Work
- Future Work
- Summary
- End

End

Questions?



Talk
Overview

Goal

Design

Usage

Validation

Experiments

Related Work

Future Work

Summary

End