#### **CODES-TRACER TUTORIAL**

**ENABLING HPC DESIGN SPACE EXPLORATION VIA DISCRETE-EVENT SIMULATION** 

#### **NIKHIL JAIN**

Compute Architect NVIDIA

#### **NEIL MCGLOHON**

PhD Student Researcher Department of Computer Science Rensselaer Polytechnic Institute

Adapted from N. Jain and M. Mubarak's slides for HOTI2017

### LINKS

- Slides available for download at:
  - <u>http://nikhil-jain.github.io/codes-hoti.pptx</u>
- QoS Slides available for download at:
  - <u>https://nmcglo.com/public-talks/df-dfp-qos-isc19.pdf</u>
- CODES Repo and Wiki:
  - https://github.com/codes-org/codes
- Tracer Repo:
  - https://github.com/LLNL/TraceR
- NERSC DOE Design Forward DUMPI Traces:
  - https://portal.nersc.gov/project/CAL/designforward.htm

### **TUTORIAL OUTLINE**

- Session I: Introduction to CODES and TraceR, and their use for simulations
- Session II: Replaying HPC workloads, simulating MPI operations and collectives
- Session III: Storage simulations, Misc

#### **PROJECT CONTRIBUTORS**

- Argonne National Laboratory, Mathematics and Computer Science
- Lawrence Livermore National Laboratory, Center for Applied Scientific Computing
- Nvidia, Compute Arch
- Rensselaer Polytechnic Institute, Computer Science
- University of Arizona, Computer Science
- UC Davis, Computer Science
- University of Maryland, Computer Science

### SESSION I: INTRODUCTION, CASE STUDY, AND INTERCONNECT MODELS

# **GOALS OF SESSION I**

- Introduction
- Case study
- Features
- HPC Interconnect models
- Configuring interconnect simulations
- Interpreting simulation output

#### **INTRODUCTION & BACKGROUND**

# **BUILDING BLOCKS OF HPC INTERCONNECTS**

- HPC performance is sensitive to the interconnect latency
- Interconnects are custom built to achieve maximum throughput
- Basic Components:
  - Network topology: Arrangement of network nodes and links to maximize performance.
  - Routing: The protocol using which a packet selects its next destination.
     Depends on the underlying network topology.
  - Flow control: Allocating network links and buffers to packets as they traverse through the network. Usually independent of the underlying topology and routing protocols.

#### WHY HPC INTERCONNECT SIMULATION?

- HPC interconnects are complex!
- Need to look at several factors that influence network performance:
  - Interconnect Topology
  - Routing
  - Quality of Service
  - Job placement policies
  - Inter-job and intra-job interference
  - Storage placement on network
- Improved performance is contingent on an optimal combination of the above factors

#### **ROSS: PARALLEL DISCRETE-EVENT SIMULATOR**

- Discrete event simulation (DES): a computer model for a system where changes in the state of the system occur at *discrete points* in simulation time
- Parallel DES allows execution of simulation on a parallel platform
- Rensselaer Optimistic Simulator System (ROSS) provides PDES capability for CODES
  - Optimistically schedules events. Rollback realized via reverse computation
  - Logical processes (LPs) model state of the system



Simulation of a 24.5K node system (16 endpoints/node)

#### **CODES: MODELING & SIMULATION FRAMEWORK**

- Accelerate HPC system co-design by providing a detailed simulation of HPC interconnects, storage, workloads and surrounding environment
- Couple best-of-breed parallel discrete event simulation with experts in interconnects and storage architecture design
- Incrementally develop HPC simulation capability, validating approach and components along the way
- Complement experimentation on real systems



Figure: CODES Architecture Diagram

#### **TRACER: REPLAYING MPI TRACES**

- OTF2 traces for MPI
- BigSim traces for Charm++, AMPI
- Default and user-defined job placement and task mapping
- MPI point-to-point semantics and protocols
- Inbuilt collectives: tree based bcast, reduce, allreduce, and barrier; message size based algorithms for alltoall and allgather
- Simulation time scaling

#### **FEATURES: HPC SIMULATIONS**

- Packet-level simulations of HPC interconnect topologies
- Trace-driven analysis (DUMPI+OTF2), synthetic workloads
- Multiple jobs can be replayed on the network
- Different job placement schemes can be used
- Multiple ranks mapped to network nodes can be used
- Detailed statistics generation
- MPI collective operations can be simulated
- General purpose storage model that uses concurrent, pipelined RDMA read/write requests (simulating burst buffers/SSD)



#### **HPC INTERCONNECT MODELS**

# STRUCTURE OF MODERN INTERCONNECTTOPOLOGIESSlim Fly Network



Dragonfly Class Networks

- Groups with All-To-All Global Connections
- Flexible Routing
- Strong resilience to inter-job interference



An 8 Group Slim Fly

#### **Torus Network**

- Homogeneous structure
- Dedicated routers per compute node



A 3-ary 3-d Torus network

A 7 Group Dragonfly



#### Fat tree Network

- Static Routing
  - Job placement impacts inter-job interference

#### A 3-level pruned Fat Tree with 64 compute nodes

# **OVERVIEW OF NETWORK MODELS**

- Multiple network models are supported 1D/2D Dragonfly, Megafly, Slim Fly, Fat Tree, Express Mesh, Arbitrary graphs, Torus...
- Abstraction layer 'model-net' sits on top of network models
  - Breaks messages into packets
  - Offers FIFO, round robin and priority queues
- To try different networks, simply switch the network configuration files!
- Storage models, MPI simulation and workload replay layers are independent of the underlying networks

#### SIMPLE NET NETWORK MODEL

- A latency/bandwidth model where message is directly sent from source to destination
- Uses infinite size queuing
- Easy setup—uses a startup delay and link bandwidth for configuration
- Mostly for debugging/testing purposes-- Can be used as a starting point when replaying MPI traces
- It can be used as a baseline network model with no contention and no routing

#### **CONFIGURING SIMPLE-NET LOGP MODEL**



Configuration file can be found in codes/tests/conf/modelnet-test.conf

# **RUNNING A SIMPLE-NET LOGP MODEL**

- ./tests/modelnet-test --sync=1 -- tests/conf/modelnettest.conf
- A simple test in which a simulated MPI rank sends message to the next rank, which replies back
- Continues until a certain number of messages is reached

# **DRAGONFLY NETWORK MODEL**

- Multiple forms of routing are supported: minimal, adaptive, non-minimal and progressive adaptive
- Packet based simulation with credit based flow control
- Uses multiple virtual channels for deadlock prevention



# **CONFIGURING DRAGONFLY NETWORK MODEL**



- For simulating multiple MPI processes per node → nw-lp=num-procs \* number of network nodes
- Self messages → messages sent to the same network node
- Overhead for sending self message can be configured

Configuration file can be found in codes/src/networkworkloads/dragonfly-custom

#### **CONFIGURING DRAGONFLY NETWORK MODEL**

#### Params

packet size in the network packet\_size="4096"; modelnet\_order=( "dragonfly\_custom","dragonfly\_custom\_router" ); # scheduler options modelnet\_scheduler="fcfs"; chunk size in the network (when chunk size = packet size, packets will not be divided into chunks chunk\_size="4096"; # modelnet scheduler="round-robin": ser of routers within each group # this is dictated by the dragonfly configuration file num\_router\_rows="6"; # number of router columns num\_router\_cols="16"; # number of groups in the retwork m\_aroups="25"; vtes for local virtual channels local\_vc\_size="8192"; size in bytes for global virtual channel global\_vc\_size="16384"; ize in bytes for compute node virtual channels cn vc size="8192": for local channels local bandwidth="5.25": bandwidth in GiB/s for global channels global\_bandwidth="4.69"; bandwidth in GiB/s for compute node-router channels cn\_bandwidth="16.0"; ROSS message size message\_size="592"; of compute nodes connected to router, dictated by dragonfly config num\_cns\_per\_router="4"; umber of global channe per router num global channels="4"; config file for intra-group connections intra-group-connections="...(src/network-workloads/conf/dragonfly-custom/intra-9K-custom"; network config file for inter-o connections /src/network-workloads/conf/dragonfly-custom/inter-9K-custom"; inter-group-connections=". routing="prog-adaptive":

Router arrangement within a group. Should match the input network configuration

Buffer size of virtual channels can be configured

Number of compute nodes per router is configurable

Network configuration files – can be custom generated (see scripts/gencray-topo/README.txt).

#### RUNNING A DRAGONFLY NETWORK SIMULATION

- Download the traces:
  - wget https://portal.nersc.gov/project/CAL/doe-miniappsmpi-traces/AMG/df\_AMG\_n1728\_dumpi.tar.gz
- Run the simulation:
  - ./src/network-workloads/model-net-mpi-replay --sync=1 -disable\_compute=1 --workload\_type="dumpi" -workload\_file=df\_AMG\_n1728\_dumpi/dumpi-2014.03.03.14.55.50- --num\_net\_traces=1728 --../src/network-workloads/conf/dragonfly-custom/modelnettest-dragonfly-edison.conf

# FAT TREE NETWORK MODEL

- Can simulate two and three level fat tree networks
- Width of the tree (number of pods) can also be configured
- Two forms of routing are supported:
  - static: uses destination-based look-up tables
  - adaptive: selects least congested output port
- Packet based simulation with credit based flow control



#### FAT TREE NETWORK CONFIGURATIONS



- Tapering can be used to connect more nodes to the leaf switches
  - Reduces the bandwidth, switches and links at higher level
- To get higher bandwidth, nodes can connect to multiple ports (multi-rail) in one or more plane (multi-plane)
  - These configurations can also be tapered to reduce switches, links at higher levels
- Model supports configurations for multiple rails, multiple plane and tapering

#### **CONFIGURING FAT TREE NETWORK MODEL**

#### LPGROUPS

```
MODELNET_GRP
```

```
repetitions="198";
nw-lp="144";
modelnet_fattree="18";
fattree_switch="3";
```

*Nw-lp is a simulated MPI process* 

A simulated fat tree network node

Three simulated fat tree network switches (one in each level of the network)

#### **CONFIGURING FAT TREE NETWORK MODEL**

#### PARAMS

packet\_size="4096"; message\_size="624"; chunk\_size="4096"; modelnet\_scheduler="fcfs"; modelnet\_order=( "fattree" ft\_type="0"; num\_levels="3"; switch\_count="198"; switch\_radix="36"; vc\_size="65536": cn\_vc\_size="65536"; link\_bandwidth="12.5"; cn\_bandwidth="12.5" routing="static"; routing\_folder="/Fat-Tree/summit" dot\_file="summit-3564"; dump\_topo="0";

Switch arrangement should match the input network configuration

Switch radix can be configured

Static routing requires precomputed destination routing tables

https://github.com/codesorg/codes/wiki/codes-

fattree#enabling-static-routing

#### SLIM FLY NETWORK MODEL



- Topology of interconnected router groups built with MMS graphs
- The max network diameter is always 2
- Packet based simulation with credit based flow control
- Multiple forms of routing are supported:
  - minimal: 2 virtual channels
  - non-minimal: 4 virtual channels
  - adaptive: 4 virtual channels



Fig. Slim Fly with q=5

#### **CONFIGURING SLIM FLY NETWORK MODEL**

#### PARAMS

packet\_size="4096"; chunk\_size="4096"; message\_size="592"; modelnet\_order=( "slimfly" ); modelnet\_scheduler="fcfs"; num\_routers="13"; num\_terminals="9"; global\_channels="13"; local\_channels="6": generator\_set\_X=("1","10","9","12","3","4"); generator\_set\_X\_prime=("6","8","2","7","5","11") local\_vc\_size="25600", global\_vc\_size="25600"; cn\_vc\_size="25600"; local\_bandwidth="12.5"; global\_bandwidth="12.5"; cn\_bandwidth="12.5"; routing="minimal"; num\_vcs="4";

Router arrangement within a group.

Generator sets are set of indices used to calculate connections between routers in the same subgraph. They must be precomputed.

# HYPER-X, EXPRESS MESH, AND TORUS

- Express Mesh: low-diameter densely connected grids
  - Allows for specifying connection gap
  - Gap = 1 -> HyperX
- Torus: based on a n-dimensional k-ary network topology
  - Number of torus dimensions and length of each dimension can be configured
  - Supports dimension order routing
- Uses bubble escape virtual channel for deadlock prevention



#### **ARBITRARY GRAPHS**

- Can also input arbitrarily connected graphs
  - Defined using DOT format
- Static routing is required
  - Generated using OpenSM, courtesy Jens Domke

#### **INTERPRETING SIMULATION OUTPUT**

### **INTERPRETING SIMULATION OUTPUT**

Total GVT Computations Total All Reduce Calls Average <u>Reduction / GVT</u>

Total bytes sent 13584368 recvd 13584368 max runtime 449332.124035 ns avg runtime 443706.882419 max comm time 449332.124035 avg comm time 443706.882419 max send time 5142770.436275 avg send time 2779472.247926 max recv time 4149449.596308 avg recv time 2335071.940672 max recv time 4129449.596308 avg recv time 2335071.940672 max recv time 432820.362362 avg wait time 430457.043672 .P-IO: writing output to dragonfly-simple-33408-1499374633/dragonfly-router-traffic dragonfly-simple-33488-1499374633/dragonfly-router-traffic dragonfly-simple-33488-1499374633/dragonfly-musg-stats dragonfly-simple-33488-1499374633/model-net-category-all dvagonfly-simple-33488-1499374633/model-net-category-test dragonfly-simple-33488-1499374633/model-net-category-test 0 nan

Application level statistics e.g. time spent in overall execution, communication, wait operations, amount of data transferred etc.

#### Enabling Ip-io-dir generates detailed network statistics files

Average number of hops traversed 1.709869 average chunk latence 0.925252 us maximum chunk latency 9.312357 us avg message size 812.563110 bytes finished messages 16820 finished chunks 65012 ADAPTIVE ROUTING STATS: 00012 chunks routed minimally 0 chunks routed non-minimally completed packets 65012

Total packets generated 39722 finished 39722

- Average and maximum times are reported for all the application runs
- Network statistics (hops traversed, latency, routing etc.) are reported for the entire network
- Detailed statistics for each MPI rank, network node, router, port are generated using Ip-io-dir option
- --Ip-io-dir=my-dir can be used to enable statistics generation (Each Ip writes it statistics to a summary file)

# STATISTICS REPORTED BY LP-IO

- Dragonfly-msg-stats:
  - number of hops, packet latency, packets sent/received, link saturation time reported for each network node
- Dragonfly-router-stats
  - link saturation time each router port
- Dragonfly-router-traffic
  - Traffic sent for each router port
- Fat tree and slim fly networks have similar statistics files.
- Mpi-replay-stats (generated for any network model):
  - bytes sent/received per MPI process
  - time spent in communication per MPI process
  - Number of sends and receives per MPI process

#### CASE STUDY

#### FIT FLY: INTERCONNECT INNOVATION THROUGH PARALLEL SIMULATION
#### CODES is designed to help find answers to the "What If..." type questions in HPC Interconnection research

### **MULTI-PLANE NETWORKS**

- Hand in hand with multi-rail networks
  - Additional links for packet injection
- Additional independent planes of routers
  - Often sharing terminals across planes





#### Single Plane

**Dual Plane** 

## SLIM FLY NETWORK TOPOLOGY

- Routers in network are organized into groups
- Each Router
  - Some degree of Local connectivity
  - Some degree of Global connectivity
  - Some degree of Terminal connectivity
- Guaranteed Diameter-2
- Groups are divided into two subgraphs
  - No global connections between two groups within same subgraph
- Connections are determined via non-trivial generation method
  - Makes it challenging to physically build



#### FIT FLY

- Multi-Planar Slim Fly Network
- Planes share single set of terminals
- Each plane follows same Slim Fly network generation method
- Terminal to Router mapping is alternating mirrored on each new plane
  - Increase path diversity





### **MORE RAILS = MORE THROUGHPUT**

- Fit Fly Based on previously validated Slim Fly Model
- Additional planes bring additional throughput
- Observed expected increase in throughput with synthetic uniform random traffic



#### **EXPERIMENTS OVERVIEW**

#### Experiment Set 1 Cross Network



#### **Experiment Set 2** Equalized Bandwidth



#### **NETWORK COMPARISON: AMG 1728**



#### Lower is better

#### **NETWORK COMPARISON: MG 1000**



#### Lower is better

## **DISCUSSION: NETWORK COMPARISON**

- Slim Fly performed well against state of the art Dragonfly and Megafly networks
  - Possible future exascale networks
- Fit Fly showed great resilience to high levels of interference traffic
  - Beat Slim Fly by an order of magnitude

#### Slim Fly and Fit Fly networks show great promise

- Low-diameter-high-path-diversity

#### EQUALIZED BANDWIDTH 12.5GIB/S – AMG 1728



#### Lower is better

#### EQUALIZED BANDWIDTH 25GIB/S – AMG 1728



#### Lower is better

# DISCUSSION: EQUALIZED BANDWIDTH

- Equalizing the aggregate bandwidth across networks slightly reduced the advantage that Fit Fly had
  - Fit Fly still pulled ahead
    - Greater interference resilience
- Additional planes of routers give less chance for any two packets to interact
  - Less interference
  - Less buffer wait time
  - Increased Application Performance

 More planes of cheaper routers may be a better option to single-planehigh-bandwidth networks

# **CASE STUDY: CONCLUSION**

- Fit Fly well outperforms state-of-the-art interconnects
- More routers and planes = Less Interference
- More routers = More Cost
  - But cheaper routers and links could be used



 CODES provides a strong environment for answering "What If..." questions and fostering future innovation in the field of HPC interconnection networks

#### **RUNNING INTERCONNECT SIMULATIONS**

Checkout the exercises at the wiki link:

https://github.com/codes-org/codes/wiki/quick-start-interconnects

# SESSION II: APPLICATION SIMULATION, WORKLOADS, MPI

### FOUR STEPS TO SIMULATIONS

- 1. Prototype system design
  - Discussed in the previous session
  - Set up using network parameters
- 2. Workload selection
  - Depends on the use case
  - Application traces
  - Synthetic patterns
  - Skeletons
- 3. Workload creation
- 4. Execution



#### General framework for replaying traces on HPC interconnect simulation



**CODES** specific framework for replaying traces on HPC interconnect simulations

### WORKLOADS

#### Synthetic Workloads:

- Follow specific communication pattern and a constant injection rate
- Often used to stress the network topology to identify best and worst case performance
- Examples include uniform random, all to all, bisection pairing, bit permutation
- Don't require simulation of MPI operations

#### HPC Application Traces:

- Useful for network performance prediction of production HPC applications
- Trace size can be large for long running or communication intensive applications
- Potential to capture computation-communication interplay
- Require accurate simulation of MPI operations
- Simulation results can be complex to analyze

#### Intel SWM Online Workloads:

- Accurate workload representations
- Decoupled from original application
- Portable to arbitrary simulation environments
- Generates traffic on-the-fly

### **DUMPI MPI TRACE LIBRARY**

- Provides trace collection and replay tools for MPI based applications
- Trace collection is simple link the MPI application with libdumpi
- Trace can be replayed using libundumpi utility
- Libundumpi provides callbacks you can use when MPI operations are replayed
- Preserves the causality order of MPI operations
- Captures detailed statistics for each MPI operation call

### **CAPTURING APPLICATION TRACES WITH DUMPI**

- Repository can be cloned at:
  - git clone https://github.com/sstsimulator/sst-dumpi.git
- Configure and build using any MPI compiler
- Make sure to use '—enable-libdumpi' when configuring
- Once installed, simply add '-L\$(DUMPI\_INSTALL) -Idumpi' in your application
- DUMPI traces will be generated automatically with each application run
- Naming convention: dumpi-yyyy.dd.mm.hh.mm.ss-MPI-RANK-ID.bin
- More information can be found at: <u>https://github.com/sstsimulator/sst-dumpi</u>
- HPC application traces in DUMPI format: <u>https://portal.nersc.gov/project/CAL/designforward.htm</u>

# **GENERATING OTF2 TRACES (1/2)**

- New Open Trace Format version 2 is supported by several tools
- ScoreP Scalable Performance Measurement Infrastructure for Parallel Codes
- Tool suite with several libraries and helper tools
  - <u>http://www.vi-hps.org/projects/score-p/</u>
- Inside ScoreP source directory
  - CC=mpicc CFLAGS="-O2" CXX=mpicxx CXXFLAGS="-O2" FC=mpif90 F77=mpif77 ./configure --without-gui --prefix=<SCOREP\_INSTALL>
  - make && make install
  - Make sure ScoreP installation's bin directory is in PATH
- Simple case: change the application linker to LD = scorep --user --nocompiler --noopenmp --nocuda --noopenacc -noopencl --nomemory <your\_linker>

# **GENERATING OTF2 TRACES (2/2)**

- Before running, set the following environment variables: *export SCOREP\_ENABLE\_TRACING=1 export SCOREP\_ENABLE\_PROFILING=0 export SCOREP\_MPI\_ENABLE\_GROUPS=ENV,P2P,COLL,XNONBLOCK*
- Turning tracing on/off: make sure these calls are synchronized
  - #include <scorep/SCOREP\_User.h>
  - SCOREP\_RECORDING\_ON(); start recoding
  - SCOREP\_RECORDING\_OFF(); stop recording
- During compilation, add flags: -I\$SCOREP\_INSTALL/include -I\$SCOREP\_INSTALL/include/scorep -DSCOREP\_USER\_ENABLE
- Trace target options export SCOREP\_TOTAL\_MEMORY=256M export SCOREP\_EXPERIMENT\_DIRECTORY=/p/lscratchd/<username>/...

#### **TRACING OUTPUT**

- scorep-\* directory generated with following content: scorep.cfg traces traces.def traces.otf2
- scorep.cfg is human readable; can be used to verify if the environment is correctly generated
- traces.otf2 is a binary meta-file
- traces is a directory that contains the details
- Use otf2-print utility in ScoreP bin to view the traces: otf2-print –L 0 traces.otf2

#### **INFORMATION CAPTURED IN A TYPICAL TRACE** (E.G. IN DUMPI, OTF2)

Time stamp, t (rounded off)	Operation type	Operation data (only critical information is highlighted)
t = 10	MPI_Bcast	root, size of bcast, communicator
t = 10.5	MPI_Irecv	source, tag, communicator, req ID
t = 10.51	user_computation	optional region name - "boundary updates"
t = 12.51	MPI_Isend	dest, tag, communicator, req ID
t = 12.53	user_computation	optional region name - "core updates"
t = 22.53	MPI_Waitall	req IDs
t = 25	MPI_Barrier	communicator

# EXAMPLE TO SHOW THE EFFECT OF REPLAYING TRACES

Original Time stamps	Original duration	New Time stamps	New duration	Operation type
10	0.5	10	0.2	MPI_Bcast
10.5	0.01	10.2	0.01	MPI_Irecv
10.51	2	10.21	2	user_computation
12.51	0.02	12.21	0.02	MPI_lsend
12.53	10	12.23	10	user_computation
22.53	2.47	22.23	0.03	MPI_Waitall
25	1	22.26	1.7	MPI_Barrier

### **DUMPI VS OTF2**

- Most of the information in the trace format is the same
- Different w.r.t. capturing of dynamically determined events: e.g. MPI\_Waitany
- DUMPI: stores all the information passed to the MPI call
  - Simulation decides which request to fulfill: accurate resolution for target systems
  - If the control flow of the program can change significantly due to the ordering of operations, simulations are not entirely correct
- OTF2: stores only the information that is used (e.g. which request was satisfied)
  - Accurately mimics the control flow of the trace run
  - But does not accurately represent execution for the target system
- Artifact of leveraging existing tools not originally intended for PDES!

### INTEL SWM WORKLOADS

- Open Source version hosted at <u>https://github.com/codes-org/SWM-workloads</u>
- Built separately: a CODES-SWM interface has been developed
- Includes several workloads including LAMMPS, Nekbone, Nearest Neighbor, HACC, MILC, Incast, Point-to-Point
- Each workload is configured by its own JSON configuration file
  - Specifies size
- More CODES use information:
  - <u>https://github.com/codes-</u> org/codes/wiki/online-workloads

1	{		
2		"jobs":	{
3			"name": "StandaloneSWM",
4			"app": "dll",
5			"dll_path": "apps/dll/lammps.so",
6			"size": 2048,
7			"time": 0,
8			"cfg":
9			{
10			"num_x_replicas": 3,
11			"num_y_replicas": 3,
12			"num_z_replicas": 3,
13			"num_time_steps": 30,
14			"req_vc" : 0,
15			"resp_vc" : 1,
16			"router_freq" : 800e6,
17			"cpu_freq" : 4e9,
18			"cpu_sim_speedup" : 1e6
19			}
20		}	
21	}		

#### SIMULATING MPI

### **MPI SIMULATION**

- Matching semantics and standard has to be followed for a correct simulation
  So obviously done
- Eager Rendezvous protocol
  - Cutoff can be specified in the config
- Library call overheads handled using a constant cost
- Collectives:
  - OTF2 based simulations implements them internally
  - DUMPI based simulations use Cortex

### TRANSLATING MPI CALLS USING CORTEX

- Internally most MPI implementations support collectives by translating into point to point
- Cortex comes with a set of translation functions to convert collectives into point to point using MPICH algorithms
- When linked with DUMPI and CODES, Cortex translates MPI collectives into point to point sends/receives (simulated by CODES)
- Cortex can also be used to implement your own translation functions (e.g. collective algorithms)
- Cortex tutorial is available at : <u>https://xgitlab.cels.anl.gov/mdorier/dumpi-cortex/wikis/home</u>

#### **CODES, CORTEX AND DUMPI INTERACTION**



### **MPI TRANSLATION WITH CORTEX**

- To enable collective translation, install Cortex and reconfigure CODES with Cortex
- Cortex available for download: git clone <u>https://xgitlab.cels.anl.gov/mdorier/dumpi-cortex.git</u>
- cmake .. -G "Unix Makefiles" -DMPICH\_FORWARD:BOOL=TRUE -DCMAKE\_INSTALL\_PREFIX=\$HOME/CODES/install/cortex -DDUMPI\_ROOT=\$HOME/CODES/install/dumpi
- See instructions at: <u>https://xgitlab.cels.anl.gov/codes/codes/wikis/codes-cortex-install</u>
- Use –with-cortex=/path/to/cortex/install option

#### IN A NUTSHELL: REPLAYING A SINGLE APPLICATION TRACE

./bin/model-net-mpi-replay --sync=1 --disable\_compute=1 -workload\_type="dumpi" --workload\_file=dumpi-2014.03.03.14.55.50- --num\_net\_traces=1728 -- modelnet-testdragonfly-edison.conf

- Runtime options
  - --workload\_type: "dumpi" or "online" for SWM
  - --num\_net\_traces : Number of input network traces
  - --workload\_file: DUMPI trace file
  - Network configuration file: Any of the network files (number of simulated ranks > number of ranks in trace)
  - -- *Ip-io-dir (optional): Generates detailed network counters and statistics*
  - -- *Ip-io-use-suffix (optional): Generates a unique directory per run*
  - --disable\_compute (optional): disable any compute time between MPI events
  - --workload\_conf\_file: for SWM, specifies name of workload to be used
- For running parallel simulations, use mpirun and –sync=3

#### SIMULATING MULTIPLE JOBS ON THE NETWORK

### **REPLAYING MULTIPLE JOBS**

./src/network-workloads/model-net-mpi-replay --sync=1 -disable\_compute=1 --workload\_type="dumpi" -workload\_conf\_file=multiple-workloads.conf -alloc\_fil -modelnet-mpi-test-dragonfly.conf

- For multiple jobs, two of the arguments are different:
  - Workload\_file: Has information on the dumpi/SWM traces for each application
  - Alloc\_file: List of simulated MPI ranks to be assigned to each job

#### EXAMPLE WORKLOAD FILE

216 /path/to/AMG/df\_AMG\_n216\_dumpi/dumpi-2014.03.03.14.55.23-125 /path/to/Multigrid/MultiGrid\_C\_n125\_dumpi/dumpi-2014.03.06.23.48.13-100 /path/to/Crystal\_Router/100/dumpi-2014.04.23.12.12.05-

2048 lammps 2197 nekbone

- Left column: Number of application ranks per job
- Right column: Path and prefix of DUMPI traces for each job or name of SWM
  - Combination SWM and Dumpi slated for future
  - Currently can only combine DUMPI+DUMPI/Synthetic and SWM+SWM/Synthetic
# **EXAMPLE JOB ALLOCATION FILE**

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 5 6 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74

- List of MPI ranks for each job
- There is a job entry per line
- Example at the top shows job placement being done in a linear scheme
- Example below shows job placement in a random fashion (assumes one rank per node).

7022 6303 <u>allocation\_gen - vim - 189×55</u>7 6518 2522 3963 4851 7551 4196 2343 6285 5142 2083 7565 8172 6737 7502 2579 6069 7474 5688 3926 834 3606 5076 7593 8038 3962 7194 2161 6692 4560 116 5981 3311 6857 5552 9 4101 7213 2022 2701 7237 1584 4716 1978 8045 6677 3721 666 2658 4221 7758 904 4581 5872 4549 6771 4490 8015 5013 4883 3696 4955 3196 4783 2410 1570 1548 5092 5457 3958

# GENERATING JOB ALLOCATIONS USING DIFFERENT SCHEMES

- Multiple schemes to map jobs onto the network
  - Randomly selected nodes
  - Contiguous or linear
  - Randomly selected switches (ranks ordered on nodes attached to a switch)
  - Clustered placement
  - ...
- Scripts can be used to generate job allocation files with any of the above schemes
- CODES keeps track of the job ID and provides it in simulation output
- Some python scripts can be found in scripts/allocation\_gen

### SESSION III: STORAGE MODELS AND SYNTHETIC TRAFFIC GENERATION

### **GOALS OF THE SESSION**

- How to do storage placement on networks?
- How to generate background network traffic?
- Using model-net API
- PDES and Networks Internal
- Continue with hands on exercises

#### **STORAGE PLACEMENT ON INTERCONNECTS**

#### **STORAGE PLACEMENT ON HPC SYSTEMS**



\* Image credit: On the role of burst buffers in Leadership class storage system by N. Liu et al. in MSST 2012

### **MODELING BURST BUFFER WITH CODES**

- General purpose model for read and write operations
- Concurrent, pipelined RDMA requests
- Comprises of the following:
  - a storage manager
  - a disk/local storage model
  - A resource tracker
- Placement of storage over the network can be modified using the network config file

### **PROTOCOL FOR WRITE OPERATIONS**



### **USING THE STORAGE MODEL**

- codes\_store\_init\_req (is\_write, priority, obj\_id, xfer\_offset, xfer\_size, codes\_req) → For initializing the request
- codes\_store\_send\_req(codes\_req, dest\_id, sender, network\_id, mapping\_context, ..) → For sending the request
- codes\_store\_send\_req\_rc → For reverse computation
- Repo available at:<u>https://xgitlab.cels.anl.gov/codes/codes-</u> storage-server

# **CONFIGURING STORAGE OVER THE NETWORK**



read overheads = ("20.0");

Number of concurrent requests
Buffer size for each thread
Size of the Memory (RAM)
Storage size (for disk/LSM)
Aggregate memory+storage size

Disk bandwidth/seek configuration

### **CONFIGURING STORAGE OVER THE NETWORK**

#### PGROUPS

```
#dragonfly_router can be in one group only!
DRAGONFLY GRP
    repetitions = "150";
    codes-store="2"; <
    lsm="2"; <
    resource="2";
    test-dummy = "2";
    test-checkpoint-client="60";
    modelnet_simplenet="1";
    modelnet_dragonfly_custom="64";
    modelnet_dragonfly_custom_router="16
EXTERNAL STR GRP
   repetitions="1";
   codes-external-store="1"; 
   modelnet_simplenet="1";
```

Two storage manager entities per 60 clients/compute nodes (Cray Cori configuration)

Local storage model entity (disk).
 One to one correspondence

Dummy nodes are to balance node to router ratio for BB routers

A total of 64 network nodes

*If the data from burst buffer needs to be drained to the external storage entity* 

#### **GENERATING BACKGROUND NETWORK TRAFFIC**

# WHY BACKGROUND TRAFFIC?

- On production HPC systems, a significant fraction of network nodes can be occupied
- How to introduce communication interference if a single application trace is being replayed on the simulation?
- Running multiple traces at a large-scale can be expensive
- One solution is to mix synthetic traffic patterns and HPC application traces

### **EXAMPLE SYNTHETIC PATTERNS**

- Uniform Random: A network node is equally likely to send to any other network node (traffic distributed throughout the network)
- All to All: Each network node communicates with all other network nodes
- Nearest neighbor: A network node communicates with near by network nodes (or the ones that are at minimal number of hops)
- Permutation traffic: Source node sends all traffic to a single destination based on a permutation matrix
- Bisection pairing: Node 0 communicates with Node 'n', node 1 with 'n-1' and so on.

### SYNTHETIC TRAFFIC IN CODES

```
/* in case of uniform random traffic, send to a random destination. */
if(traffic == UNIFORM)
 b -> c1 = 1;
 local dest = tw rand integer(lp->rng, 0, num nodes -1);
else if(traffic == NEAREST_GROUP)
 local_dest = (local_id + num_nodes_per_grp) % num_nodes;
 //printf("\n LP %ld sending to %ld num nodes %d ", local_id, local_dest, num_nodes);
else if(traffic == NEAREST_NEIGHBOR)
 local dest = (local id + 1) % num nodes;
printf("\n LP %ld sending to %ld num nodes %d ", rep id * 2 + offset, local dest, num nodes);
assert(local dest < num nodes);</pre>
  codes mapping get lp id(group name, lp type name, anno, 1, local dest / num servers per rep, local dest % num servers per rep, &global dest);
global_dest = codes_mapping_get_lpid_from_relative(local_dest, group_name, lp_type_name, NULL, 0);
ns->msg sent count++;
model_net_event(net_id, "test", global_dest, PAYLOAD_SZ, 0.0, sizeof(svr_msg), (const void*)m_remote, sizeof(svr_msg), (const void*)m_local, lp);
```

\* Code snippet from synthetic workload generator

- Typical patterns supported are uniform random and nearest neighbor.
- All to all and stencil patterns have been tested (pending integration)

See src/network-workloads/model-net-synthetic-custom-dfly.c and related files

#### **GENERATING BACKGROUND TRAFFIC WITH CODES**

- Communication based on uniform random traffic
- Kicks off when the main workload starts
- A notification is sent to the background traffic node to stop generating traffic once the main workload finishes
- How to enable synthetic traffic generation?
- Simply add "synthetic" instead of DUMPI trace path in workloads config file

216 synthetic 125 /path/to/Multigrid/Multigrid\_125/dumpi-2014.03.06.23.48.13-

#### PDES AND NETWORK INTERNALS

# **DISCRETE EVENT SIMULATION (DES)**

- Computer model for a system where changes in the state of the system occur at discrete points in simulation time
- In this model, each component of the system being simulated is represented independently via their
- State variables
- Virtual time
- Events scheduled on it and by it



Event scheduling from one component to another progresses and coordinates virtual time across components

Example from slides by Prof Carothers, RPI

#### **IMPLEMENTING DES**



- ROSS lets users define LP (logical processes) on which events can be scheduled with time stamps
- Each LP can have a local state that is accessible and modified only when events are executed on it

# **ROSS LP**

};

```
tw_lptype model_lps[] = {
 {
  (init_f) model_init,
  (event_f) model_event,
  (revent_f) model_event_reverse,
  (final_f) model_final,
  (map_f) model_map,
  sizeof(state)
 },
```

### **EXAMPLE OF AN EVENT FUNCTION**

{

Typical events act based on "type" and "content" of the message

```
static void svr_event(
    svr_state * ns,
   tw_bf * b,
   svr_msg * m,
   tw_lp * lp)
    (void)b;
   switch (m->svr_event_type)
    {
        case REO:
            handle_req_event(ns, m, lp);
            break;
        case ACK:
            handle_ack_event(ns, m, lp);
            break;
        case KICKOFF:
            handle_kickoff_event(ns, m, lp);
            break;
        case LOCAL:
           handle_local_event(ns);
```

```
static void handle_req_event(
   svr_state * ns,
   svr_msg * m,
   tw lp * lp)
ł
   assert(!do_pull);
    svr_msg * m_local = malloc(sizeof(svr_msg));
   svr_msg * m_remote = malloc(sizeof(svr_msg));
   m_local->svr_event_type = LOCAL;
   m_local->src = lp->gid;
   memcpy(m_remote, m_local, sizeof(svr_msg));
   m_remote->svr_event_type = ACK;
   ns->msg_recvd_count++;
   m->ret = model_net_event(net_id, "test", m->src, PAYLOAD_SZ,
                0.0, sizeof(svr_msg), (const void*)m_remote,
                sizeof(svr_msg), (const void*)m_local, lp);
   return;
```

#### ROSS'S LAW OF OPTIMISTIC EXECUTION: FOR EVERY FORWARD ACTION, YOU MUST TELL ROSS HOW TO GO BACKWARDS

```
static void svr_rev_event(
static void svr_event(
                                                  svr_state * ns,
    svr_state * ns,
                                                  tw_bf * b,
    tw_bf * b,
                                                  svr_msg * m,
    svr msq * m,
                                                  tw_lp * lp)
    tw_lp * lp)
                                              ł
{
                                                  (void)b;
    (void)b;
                                                  switch (m->svr_event_type)
   switch (m->svr_event_type)
                                                  {
    ł
                                                      case REO:
        case REO:
                                                          handle_req_rev_event(ns, m, lp);
            handle_req_event(ns, m, lp);
                                                          break:
            break;
                                                      case ACK:
        case ACK:
                                                          handle_ack_rev_event(ns, m, lp);
            handle_ack_event(ns, m, lp);
                                                          break:
            break;
                                                      case KICKOFF:
        case KICKOFF:
                                                          handle_kickoff_rev_event(ns, m, lp);
            handle_kickoff_event(ns, m, lp);
                                                          break;
            break;
                                                      case LOCAL:
        case LOCAL:
                                                          handle_local_rev_event(ns);
           handle_local_event(ns);
                                                          break;
```

```
static void handle_req_rev_event(
                                                          svr_state * ns,
                                                          svr_msg * m,
                                                          tw_lp * lp)
                                                      {
static void handle req event(
                                                          ns->msg_recvd_count--;
    svr_state * ns,
                                                          if (do_pull){
    svr_msq * m,
                                                              model_net_event_rc2(lp, &m->ret);
    tw_lp * lp)
                                                          }
{
                                                          else{
    assert(!do_pull);
                                                              model_net_event_rc2(lp, &m->ret);
    svr_msg * m_local = malloc(sizeof(svr_msg));
                                                          }
    svr_msg * m_remote = malloc(sizeof(svr_msg));
    m_local->svr_event_type = LOCAL;
    m local->src = lp->gid;
    memcpy(m_remote, m_local, sizeof(svr_msg));
    m remote->svr event type = ACK;
    ns->msg_recvd_count++;
    m->ret = model_net_event(net_id, "test", m->src, PAYLOAD_SZ,
                0.0, sizeof(svr_msg), (const void*)m_remote,
                sizeof(svr_msg), (const void*)m_local, lp);
    return;
```

# **APPLICATION SIMULATION IN CODES**



### AVAILABLE MODELS, FEATURES, AND ADDING A NEW NETWORK MODELS

- Available: simple-net model, torus, dragonfly-(custom), fat-tree, slim fly, expressmesh/hyperX
- Typical model consist of NIC (terminals) and switches/routers
- NICs
  - Common code available for within-node, message ordering, etc
  - Plugin code for individual network
- Switch/routers
  - Entirely within a network model
- But, a significant fraction of node is similar for NIC plugin and switch!

### NIC COMMON

- Types of queues:
- fifo
- round-robin
- priority
- Other params
- intra\_bandwidth (10)
- node\_copy\_queues (4)



### **NIC NETWORK SPECIFIC**

```
struct model net method torus method =
ł
   .mn configure = torus_configure,
   .mn register = NULL,
   .model net method packet event = torus packet event,
   .model net method packet event rc = torus packet event rc,
   .model net method recv msg event = NULL,
   .model_net_method_recv_msg_event_rc = NULL,
                                                 tw_lptype torus_lp =
   .mn get lp type = torus get lp type,
   .mn_get_msg_sz = torus_get_msg_sz,
                                                     (init_f) torus_init,
   .mn report stats = torus report stats,
                                                     (pre_run_f) NULL,
   .mn collective call = NULL,
                                                      (event_f) event_handler,
   .mn_collective_call_rc = NULL,
                                                     (revent_f) node_rc_handler,
   .mn sample fn = NULL,
                                                      (commit_f) NULL,
   .mn sample rc fn = NULL,
                                                     (final_f) final,
   .mn sample init fn = NULL,
                                                     (map_f) codes_mapping,
   .mn sample fini fn = NULL
                                                     sizeof(nodes_state),
};
                                                 };
```



#### **CONTRIBUTING:**

- Fork off on the github repository
- Add new features
- Submit a pull request!

#### **THANK-YOU**

#### **ADDITIONAL MATERIAL**

#### **CODES INSTALLATION**

### **INSTALLATION & SETUP 1/2**

#### ROSS INSTALLATION

- Download ROSS repo: git clone https://github.com/ROSS-org/ROSS
- Configure by making a build directory: cd build
- ARCH=x86\_64 CC=mpicc CXX=mpicxx cmake
- -DCMAKE\_INSTALL\_PREFIX=../install ../
- make –j 3 && make install

#### CODES INSTALLATION

- Download CODES repo: git clone https://github.com/codes-org/codes
- ./prepare.sh
- Configure in build directory: cd build

../configure --prefix=/path/to/install CC=mpicc CXX=mpicxx PKG\_CONFIG\_PATH=/path/to/ross/install/lib/pkgconfig

Do both make && make tests

# **INSTALLATION & SETUP 2/2**

#### DUMPI INSTALLATION

- git clone <u>https://github.com/sstsimulator/sst-dumpi</u>
- CFLAGS="-DMPICH\_SUPPRESS\_PROTOTYPES=1 -DHAVE\_PRAGMA\_HP\_SEC\_DEF=1"
- \_./bootstrap.sh
- ./configure --enable-libundumpi CC=mpicc --prefix=\$INSTALL\_PATH
- Use -- with-dumpi=/path/to/dumpi/install option to enable DUMPI with CODES
- OTF2 and BigSim based Tracing
  - ScoreP/OTF2: http://www.vi-hps.org/projects/score-p/
  - Charm++/BigSim:http://charm.cs.illinois.edu/manuals/html/bigsim/manual.html
  - TraceR: https://github.com/LLNL/tracer/
- For installation details and documentation see:
  - <u>https://xgitlab.cels.anl.gov/codes/codes/wikis/home</u>
  - https://xgitlab.cels.anl.gov/codes/codes/wikis/installation

#### MORE ON TRACER
### TRACER – A LAYER FOR CONFIGURABLE REPLAY OF APPLICATION TRACES



Capture application behavior by tracing runs on existing systems

Reproducing the execution: applications' behavior, job placement and mapping, job scheduling, MPI/Charm++, etc.

Simulation of traffic flow on NICs and networks

# DOCUMENTATION

- Distributed with TraceR source code
- README.md getting started
- README.OTF OTF2 installation and usage
- docs/UserWriteUp.txt detailed workflow and usage
- utils/README job placement and task mapping

# **INSTALLING TRACER (1/4)**

- Hosted on github: <u>https://github.com/LLNL/tracer/</u>
- git clone and follow README.md
- Download and install ROSS
  - Last verified commit provided
- Download and install CODES
  - Last verified commit provided

# **INSTALLING TRACER (2/4)**

- Choose a trace format: BigSim or OTF2
- For BigSim, download Charm++
  - git clone http://charm.cs.uiuc.edu/gerrit/charm
- Assuming MPI is available, install two flavors of Charm++
  - For compiling codes for trace generation
    ./build bgampi mpi-linux-x86\_64 bigemulator –O2
  - For compiling TraceR

./build charm++ mpi-linux-x86\_64 bigemulator --with-production

# **INSTALLING TRACER (3/4)**

- For OTF2, download ScoreP
  - <u>http://www.vi-hps.org/projects/score-p/</u>
- Inside ScoreP source directory
  - CC=mpicc CFLAGS="-O2" CXX=mpicxx CXXFLAGS="-O2" FC=mpif90 F77=mpif77 ./configure --without-gui --prefix=<SCOREP\_INSTALL>
  - make && make install
- Make sure ScoreP installation's bin directory is in PATH

# **INSTALLING TRACER (4/4)**

- In tracer/Makefile.common
- Set ROSS to ROSS's installation directory
- Set CODES to CODES's installation directory
- If using BigSim,
  - Set CHARMPATH
  - SELECT\_TRACE = -DTRACER\_BIGSIM\_TRACES=1
- If using OTF2,
  - Make sure ScoreP installation's bin directory is in PATH
  - SELECT\_TRACE = -DTRACER\_OTF\_TRACES=1
- make: generates traceR executable

# **MORE ON GENERATING OTF2 TRACES**

- ScoreP macros can be used to mark special regions
  - SCOREP\_USER\_REGION\_BY\_NAME\_BEGIN( regionname, SCOREP\_USER\_REGION\_TYPE\_COMMON)
     SCOREP\_USER\_REGION\_BY\_NAME\_END(regionname)
- Printing simulation time at locations of interest:
- Region name with prefix TRACER\_WallTime\_<any\_name> prints current time during simulation with tag <any\_name>.

# **MORE ON GENERATING OTF2 TRACES**

- Simulation time looping:
- Region name TRACER\_Loop can be used to mark beginning and ending of a code loop (currently once)
- In future, region names will be used for
  - Targeted kernel time modifications
  - Targeted message size modifications

## **CODE EXAMPLE FOR TRACING WITH OTF2**

while(iterations < MAX\_ITER) {</pre>

int main(int argc, char **argv) {	if(myRank == 0)	
MPI_Init(&argc, &argv);	SCOREP_USER_REGION_BY_NAME_BEGIN	
SCOREP_RECORDING_OFF();	("TRACER_WallTime_InLoop", SCOREP_USER_REGION_TYPE_COMMON);	
//initializaton code	//kernel and other code	
MPI_Barrier(MPI_COMM_WORLD);		
SCOREP_RECORDING_ON();	}	
SCOREP_USER_REGION_BY_NAME_BEGIN ("TRACER Loop", SCOREP USER REGION TYPE COMMO	SCOREP_USER_REGION_BY_NAME_END ("TRACER_Loop"); DN):	
if(!myRank)	MPI_Barrier(MPI_COMM_WORLD);	
SCOREP_USER_REGION_BY_NAME_BEGIN	endTime = MPI_Wtime();	
("TRACER_WallTime_Total", SCOREP_USER_REGION_TYPE_COMMON);		
startTime = MPI_Wtime();	if(!myRank)	
	SCOREP_USER_REGION_BY_NAME_END ("TRACER_WallTime_Total");	
	SCOREP RECORDING OFF():	

## **RUNNING TRACER**

- A typical run command: mpirun -np 8 ./traceR --sync=3 --nkp=16 --extramem=100000 --max-optlookahead=1000000 --timer-frequency=1000 –lp-io-dir=stats-dir -- torus.conf tracer\_config
- In green, ROSS options
  - --nkp : how many KPs to create per PE = total LPs/<-np>
  - --extramem : how many ROSS messages to allocate = 100K should work for most cases
  - --max-opt-lookahead : optimistic leash = 1 millisecond is a good number

## **RUNNING TRACER**

• A typical run command:

mpirun -np 8 ./traceR --sync=3 --nkp=16 --extramem=100000 --max-optlookahead=1000000 --timer-frequency=1000 –lp-io-dir=stats-dir -- torus.conf tracer\_config

- TraceR-CODES options
  - --timer-frequency : how frequently to print progress of task completion; optional, default = 5000
  - --Ip-io-dir : where to write output stats; optional; code fails if the directory exists to avoid over-writing
  - torus.conf : network config file
  - tracer\_config : TraceR config file

# **TRACER PARAMETER IN NETWORK FILE**

- server in MODELNET\_GRP
  - Number of processes associated with a switch
  - Assigned in a round-robin manner to nodes
- soft\_delay in PARAMS
  - Approximate overhead of making an MPI/runtime call
  - In nanoseconds
- rdma\_delay in PARAMS
  - Overhead of using RDMA call in rzv protocol, in nanoseconds
- eager\_limit in PARAMS
  - Switch over point between eager and rzv protocols, in bytes
- copy\_per\_byte in PARAMS
  - Copy cost for a byte, in nanoseconds per byte

# **TRACER CONFIG FILE (1/2)**

Format:

<global map file> or NA

#jobs

<path to job traces> <task mapping file or NA> <#ranks> <loop iterations>

Example: global\_map.bin 2 traces-64/traces.otf2 job0 64 1 traces-32/traces.otf2 job1 32 1

# **TRACER CONFIG FILE (2/2)**

- At the end of file,
- E <job id> scale\_all <scale factor>
  - Inverse scales computation time by the given factor
  - E.g. : E 0 scale\_all 40
- S <job id> <msg size> <replace by>
  - Change the size of message
  - Under review, to be merged

## SAMPLE OUTPUT

PE0 - LP\_GID:0 : START SIMULATION, TASKS COUNT: 245611, FIRST TASK: 0, RUN TIME TILL NOW=70.000000 s, CURRENT SIM TIME 1.005877

[00:time at task 0/245611 0.000000]

[00:Begin TRACER\_WallTime\_MainLoop 0.000001]

[00:time at task 100/245611 0.000663]

[00:Begin TRACER\_WallTime\_NextTrajec 0.001175]

[ 0 0 : time at task 200/245611 0.003104 ]

. . . .

.... [ 0 0 : time at task 245600/245611 1.485123 ] [ 0 0 : End TRACER\_WallTime\_NextTrajec 1.485264 ] [ 0 0 : End TRACER\_WallTime\_MainLoop 1.485265 ]

#### **CASE STUDIES WITH CODES-TRACER**

#### **Evaluating Methods for Effective Fat-tree Utilization**

- Quantify the effectiveness of multiple network rails in increasing network performance
- Evaluate job placement methodologies and routing choices to maximize the throughput of a multi-rail fat-tree network
- Investigate network performance gains offered by multiple rails in response to increased compute performance of high-density compute nodes
- Improve routing, injection and job allocation design choices of multi-rail fat-tree networks to maximize individual application and system wide performance



#### Fat-tree strong scaling rail performance

#### **Related papers**

- Jain et al. Predicting the Performance Impact of Different Fat-Tree Configurations (To appear in Supercomputing 2017)
- Wolfe et al. Preliminary Performance Analysis of Multi-Rail Fat-tree Networks (CCGRID 2017)

#### VISUALIZING HPC INTERCONNECT PERFORMANCE



Comparing job placement schemes using the CODES simulation framework of the high-radix dragonfly network running the Algebraic Multigrid Solver (AMG), AMR BoxLib and MiniFE applications. The visualizations show the aggregated



twork links.

Application	Ranks	Data	Comm. Pattern
AMG	1728	1.2GB	3D nearest neighbor
AMR Boxlib	1728	2.2GB	Irregular and sparse
MiniFE	1152	147GB	Many-to-many

• K. Li, M. Mubarak, R. Ross et al. "Visual Analytics Techniques for Exploring the Design-space of large-scale high-radix networks", to appear in IEEE Cluster 2017

#### ACCELERATING COLLECTIVE COMMUNICATION ON DRAGONFLY NETWORKS

- Need to design new scalable collective algorithms for high-radix interconnects
- Investigate the parameter space of these algorithms, and the effects of crossapplication communication interference
- Designed/evaluated several collective algorithms for Dragonfly networks



M. Dorier et al. **"Evaluation of Topology-Aware Broadcast Algorithms for Dragonfly Networks** ", IEEE Cluster 2016

127

### **USING MODEL-NET API**

## CONFIGURATION

- Model-net– An abstraction layer on top of network models topology details are specified through the config files
- A valid network configuration file examples can be found in the repo
- Network model must be registered model\_net\_register
- CODES mapping must be setup codes\_mapping\_setup
- Use model-net function calls model\_net\_event(network id, source, destination, message size,...)
- Example of using model-net tests/model-net-test.c