# *Distributed (Parallel) Simulation with ns-3*

WNS3 2015 Tutorial, Castelldefels (Barcelona), Spain

https://www.nsnam.org/wiki/AnnualTraining2015

May 12, 2015

Peter D. Barnes, Jr (LLNL)

Ken Renard (ARL)

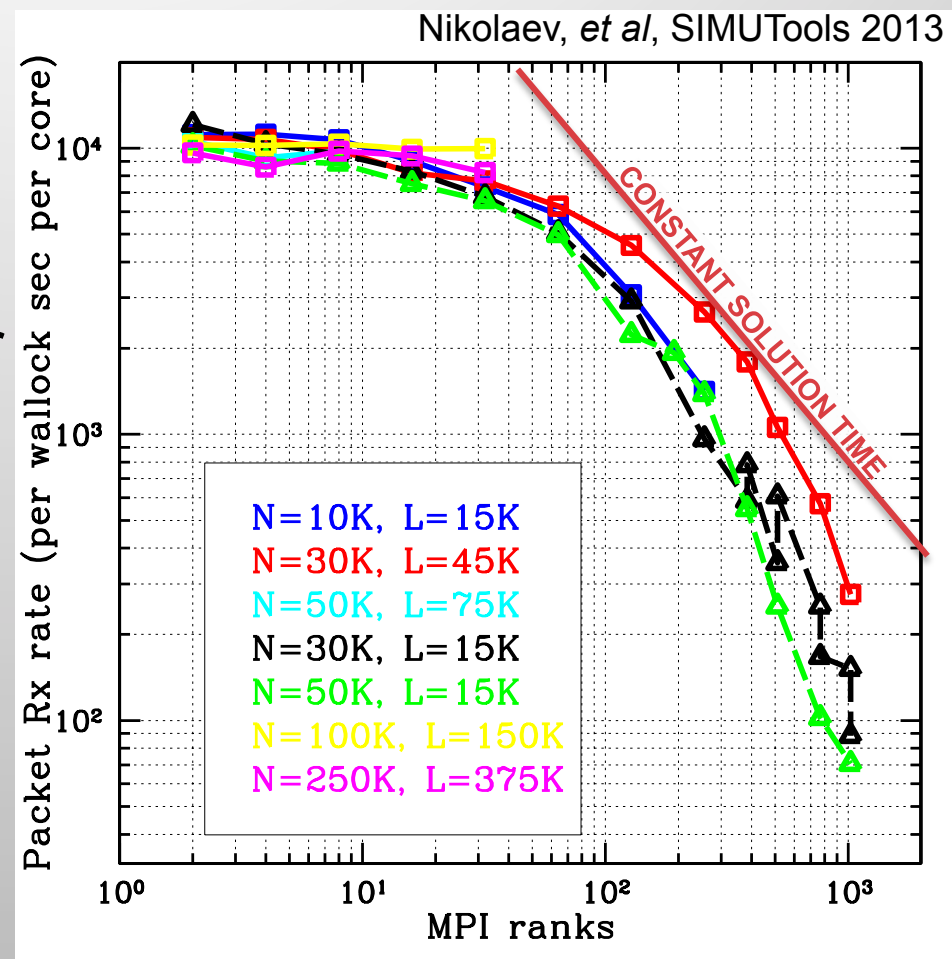# Why should you care about distributed (parallel) simulation?

- ## Faster execution
  - Measure $\sim 10^4$ packet receives/wall clock second/core


- ## Large models, too big for one compute node

- ## Heavy-weight nodes
  - DCE applications

  - Virtual machines

  - Core routers with large forwarding tables

# Motivation for High Performance, Scalable Network Simulation

- Reduce simulation run-time for large, complex network simulations
  - Complex models require more CPU cycles and memory
    - MANETs, robust radio devices
    - More realistic application-layer models and traffic loading
    - Load balancing among CPUs
  - Potential to enable real-time performance for NS-3 emulation

- Enable larger simulated networks
  - Distribute memory footprint to reduce swap usage
  - Potential to reduce impact of $N^2$ problems such as global routing

- Allows network researchers to run multiple simulations and collect significant data

# ns-3 Execution Scaling

- ## $10^{\sim 4}$ packets/core/sec
  - Independent of model size
  - 100 cores is 100x faster than 1 core

Nikolaev, *et al*, SIMUTools 2013

**Packet Rx rate (per wallock sec per core)** vs **MPI ranks**

CONSTANT SOLUTION TIME

N=10K, L=15K
N=30K, L=45K
N=50K, L=75K
N=30K, L=15K
N=50K, L=15K
N=100K, L=150K
N=250K, L=375K

.ıllNS-3

4

# How many hardware threads do you have?

**This laptop**

- MacBook Pro, mid 2009

- Intel Core 2 Duo:
  2 hardware threads

**My other computer**

- BlueGene/Q
  - Currently TOP500 #3
  - 8M hardware threads

# Outline

- Motivation
- Intro to MPI
- Parallel Discrete Event Si
- PDES in ns-3
- Constructing models
- Example
- Error Conditions
- Performance
- Future Capabilities

**Big topics. Focus on the concepts and terms needed to understand //ns-3**

| |
|---|
| Motivation |
| MPI |
| PDES |
| ns-3 |
| Models |
| Example |
| Errors |
| Performance |
| Future |

# Parallelizing ns-3 Models Is Straightforward, But…



https://xkcd.com/1205

# And the Reality…
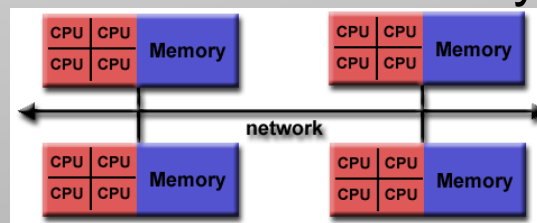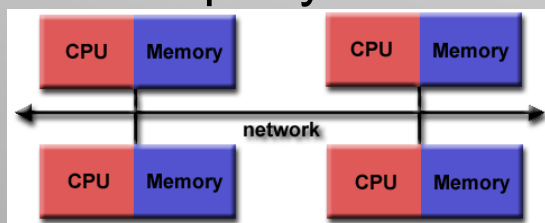


https://xkcd.com/1319

# MPI Topics

- Core features

- API specification

- Ranks and communicators

- Point-to-point messages

- Collectives

- Getting and using MPI

- Examples
  - Hello World
  - Simple messaging
  - Ghost cell pattern

# Message Passing Interface (MPI) Features 1

- *De facto* standard programming model for parallel scientific codes (but see Charm++ for an alternative)

- Basic functionality is sending messages (data) between processes, which are called *ranks*

- Core features
  - API specification
    - ~Language independent (FORTRAN, C, C++, Python, Java,…)
    - Supports (doesn't preclude) high performance and scalability
  - Point-to-point (src,dst) messaging, as well as collectives
    - Broadcast, reduction (compute min value), …
  - Works equally well on distributed and shared memory



https://computing.llnl.gov/tutorials/mpi/

# MPI Features 2

- API specification, implementation up to "vendor"
  - Vary in performance, runtime launch, …
  - Architecture-specific libraries can target specialized hardware
    — High-speed interconnects:  Infiniband, PAMI, …
    — Specialized network topologies:  Fat-tree, Dragonfly, 5-d torus
    — Specialized network interfaces: low-latency, high-throughput
    — Multi-path routing
  - Multiple implementations can coexist (but not interoperate)
    — OpenMPI, MPICH, IBM, …
    — Language-specific: mpi4py, mpiJava, …

# MPI Concepts
# Ranks and Communicators

- ## Processes are called *ranks*

- ## Communicator
  - ### Group of ranks, numbered [0,R) within the group
  - ### Enables separating messages by purpose
  - ### Initial default communicator is MPI_COMM_WORLD
  - ### Several functions for creating communicators, to support specific topologies

ns-3 Rank 0

ns-3 Rank 1

ns-3 Rank N

ns-3 Communicator

Service Rank 0

Service Rank 1

Service Rank M

Service Communicator

# MPI Concepts
# Point-to-Point Messages

Send a message to a specific rank in a communicator

- `MPI_Send(data, data_length, data_type,`
  `            destination, tag, communicator)`

  — `Data/data_length`    Message contents

  — `data_type`           MPI-defined data types
                          Or a custom data type (*i.e*, a `struct`)

  — `destination`         Rank Number

  — `tag`                 Application tag to distinguish types

  — `communicator`

- Matching `MPI_Recv()`

- Blocking and non-blocking versions

- Various optimized calls, for managing memory

- Wait for a message, test for new messages

# MPI Concepts 3
# Collective Communications

Higher level patterns involving more than 2 ranks

- ## Synchronization
  - `MPI_Barrier`

- ## Data movement
  - `MPI_Bcast`
    `MPI_Scatter`
    `MPI_Gather`
    `MPI_Allgather`

- ## Reductions
  - `MPI_Reduce`
    `MPI_Allreduce`

- ## Combinations
  - `MPI_Reduce_scatter`
    `MPI_Alltoall`
    `MPI_Scan`

broadcast

scatter

gather

reduction

**...ılIıNS-3**

# MPI Full API

## Environment Management

| | | | | | |
|---|---|---|---|---|---|
| MPI_Abort | MPI_Errorhandler_create | MPI_Errorhandler_free | MPI_Errorhandler_get | MPI_Errorhandler_set | MPI_Error_class |
| MPI_Error_string | MPI_Finalize | MPI_Get_processor_name | MPI_Get_version | MPI_Init | MPI_Initialized |
| MPI_Wtick | MPI_Wtime | | | | |

## Point-to-Point Communication

| | | | | | |
|---|---|---|---|---|---|
| MPI_Bsend | MPI_Bsend_init | MPI_Buffer_attach | MPI_Buffer_detach | MPI_Cancel | MPI_Get_count |
| MPI_Get_elements | MPI_Isend | MPI_Iprobe | MPI_Irecv | MPI_Irsend | MPI_Isend |
| MPI_Issend | MPI_Probe | MPI_Recv | MPI_Recv_init | MPI_Request_free | MPI_Rsend |
| MPI_Rsend_init | MPI_Ssend | MPI_Ssend_init | MPI_Start | MPI_Startall | MPI_Test |
| MPI_Test_cancelled | MPI_Testall | MPI_Testany | MPI_Testsome | MPI_Wait | MPI_Waitall |
| MPI_Waitany | MPI_Waitsome | | | | |

## Collective Communication

| | | | | | |
|---|---|---|---|---|---|
| MPI_Allgather | MPI_Allgatherv | MPI_Allreduce | MPI_Alltoall | MPI_Alltoallv | MPI_Barrier |
| MPI_Bcast | MPI_Gather | MPI_Gatherv | MPI_Op_create | MPI_Op_free | MPI_Reduce |
| MPI_Reduce_scatter | MPI_Scan | MPI_Scatter | MPI_Scatterv | | |

## Process Group

| | | | | | |
|---|---|---|---|---|---|
| MPI_Group_compare | MPI_Group_difference | MPI_Group_excl | MPI_Group_free | MPI_Group_incl | MPI_Group_intersection |
| MPI_Group_range_excl | MPI_Group_range_incl | MPI_Group_rank | MPI_Group_size | MPI_Group_translate_ran | MPI_Group_union |

## Communicators

| | | | | | |
|---|---|---|---|---|---|
| MPI_Comm_compare | MPI_Comm_create | MPI_Comm_dup | MPI_Comm_free | MPI_Comm_group | MPI_Comm_rank |
| MPI_Comm_remote_group | MPI_Comm_remote_size | MPI_Comm_size | MPI_Comm_split | MPI_Comm_test_inter | MPI_Intercomm_create |
| MPI_Intercomm_merge | | | | | |

## Derived Types

| | | | | | |
|---|---|---|---|---|---|
| MPI_Type_commit | MPI_Type_contiguous | MPI_Type_extent | MPI_Type_free | MPI_Type_hindexed | MPI_Type_hvector |
| MPI_Type_indexed | MPI_Type_lb | MPI_Type_size | MPI_Type_struct | MPI_Type_ub | MPI_Type_vector |

## Virtual Topology

| | | | | | |
|---|---|---|---|---|---|
| MPI_Cart_coords | MPI_Cart_create | MPI_Cart_get | MPI_Cart_map | MPI_Cart_rank | MPI_Cart_shift |
| MPI_Cart_sub | MPI_Cartdim_get | MPI_Dims_create | MPI_Graph_create | MPI_Graph_get | MPI_Graph_map |
| MPI_Graph_neighbors | MPI_Graph_neighbors_count | MPI_Graphdims_get | MPI_Topo_test | | |

## Miscellaneous

| | | | | | |
|---|---|---|---|---|---|
| MPI_Address | MPI_Attr_delete | MPI_Attr_get | MPI_Attr_put | MPI_Keyval_create | MPI_Keyval_free |
| MPI_Pack | MPI_Pack_size | MPI_Pcontrol | MPI_Unpack | | |

# Getting and Using MPI

- Check your package manager

- Only the API defined
  - Tool names and configuration vary
  - OpenMPI commands used here for illustration

- Building your code
  - Typically a compiler wrapper script, to ensure correct includes and libraries: `$ mpicc …`
  - Often hidden inside your build system (`Makefile`, `wscript`)

- Multiple executables
  - Possible to run different executables on different ranks
  - But job launch commands depend on package, so not portable
  - Typically build everything into one executable, select functions based on rank id at runtime

# Getting and Using MPI

- Where to run ranks?
  - Single computer: typically defaults to all hardware threads
  - Ad hoc cluster: `--hostfile` node names and max number of ranks
  - HPC cluster: typically via a batch job system, which selects physical nodes and launches jobs as a shell script
  - "Overcommitment": running more ranks than cores or hardware threads

- Launching
  `$ mpirun –n <nranks> <executable>`
  - Need to launch on each host

**Example OpenMPI Hosts File**

```
# This is an example hostfile.  Comments begin with #
#
# The following node is a single processor machine:
foo.example.com

# The following node is a dual-processor machine:
bar.example.com slots=2

# The following node is a quad-processor machine,
# over-subscribing disallowed
yow.example.com slots=4 max-slots=4
```

**Example OpenMPI Build and Run**

```
$ mpicxx -o hello hello.cc
$ mpirun -np 4 ./hello
Hello World from rank 3 of 4 (35986)
Hello World from rank 0 of 4 (35983)
Hello World from rank 1 of 4 (35984)
Hello World from rank 2 of 4 (35985)
```

# Parallel Hello World

- ## Typical Structure
  1. Include header
  2. Initialize MPI with command-line args
  3. Get world size, my rank index
  4. Parallel code
     - Send messages, synchronize…
  5. Clean shutdown

  6.  Build and Launch

### hello.cc

```
#include <mpi.h>

int
main (int argc, char **argv)
{
  int size, rank, rc;

  rc = MPI_Init (&argc, &argv);
  if (rc != MPI_SUCCESS)
    MPI_Abort(MPI_COMM_WORLD, rc);

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

  printf ("Hello World from rank %d of %d
          (%d)\n", rank, size, getpid ());

  MPI_Finalize();
}
```

**Example OpenMPI Build and Run**

```
$ mpicxx -o hello hello.cc
$ mpirun -n 4 ./hello
Hello World from rank 3 of 4 (35986)
Hello World from rank 0 of 4 (35983)
Hello World from rank 1 of 4 (35984)
Hello World from rank 2 of 4 (35985)
```

.ıılllNS-3

# Simple Messaging Example

Rank0 --"Hello"--> Rank1

■ Typical Structure

1. Check #ranks

2. Message data buffers

3. Each rank runs different code

4. Sends paired with Recv

5. Send and Recv data lengths, types match

```
$ mpicxx -o send1 send1.cc
$ mpirun -np 4 ./send1
Rank 1 received message "Hello" (5) from rank 0 tag 0.
```

**Parallel Body of send1.cc**

**①**
```
if (size < 2) {
    printf ("Need two ranks\n");
    MPI_Abort(MPI_COMM_WORLD, 0);
}
```

**②**
```
char *msg = (char *)"Hello";
int msg_len = strlen(msg);
char in_msg[msg_len + 1];  // leave space to add null
```
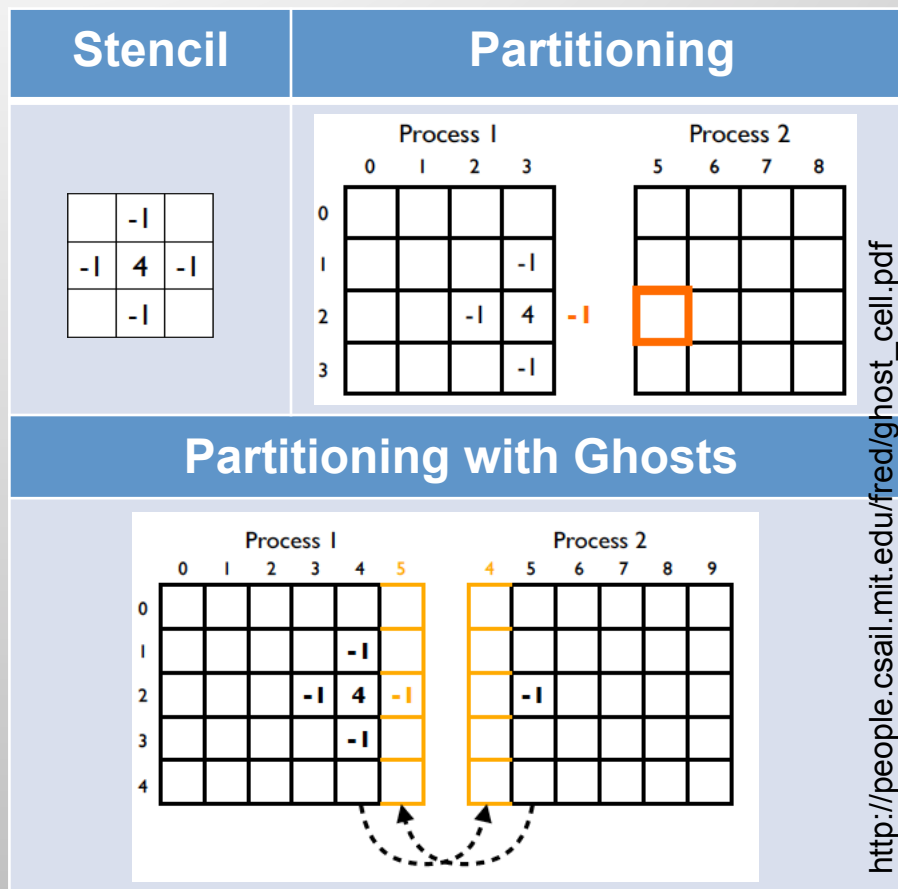
**③**
```
if (rank == 0) {
    int dest = 1;
    rc = MPI_Send (msg, msg_len, MPI_CHAR, dest,
                   0, MPI_COMM_WORLD);
}
```

**④** **⑤**
```
if (rank == 1) {
    int count = 0;
    MPI_Status stat;

    rc = MPI_Recv (&in_msg, msg_len, MPI_CHAR,
                   MPI_ANY_SOURCE, 0, MPI_COMM_WORLD,
                   &stat);
    in_msg[msg_len] = (char) 0;
    MPI_Get_count (&stat, MPI_CHAR, &count);
    printf("Rank %d received message \"%s\" (%d) "
           "from rank %d tag %d.\n",
           rank, in_msg, count,
           stat.MPI_SOURCE, stat.MPI_TAG);
}
```

# Ghost Cell Design Pattern

- Decomposition
  - Need neighbors' data: stencil
  - Some neighbors are remote

- Solution:
  - Ghosts replicate data
  - Two-phase execution
    - Exchange neighbor data
      *Communication*
    - Compute local update
      *Computation*

| Stencil | Partitioning |
|---|---|
| | |



**Partitioning with Ghosts**



http://people.csail.mit.edu/fred/ghost_cell.pdf

**Maximize *computation/communication*. Overlap computation and communication.**

# PDES Topics

- **Discrete Event Simulation**
  - Mathematical paradigm and time control
  - State and time evolution
  - Event scheduling
  - Time consuming processes

- **Parallel DES**
  - Logical processors
  - Causality
  - Granted-time synchronization
  - Lookahead
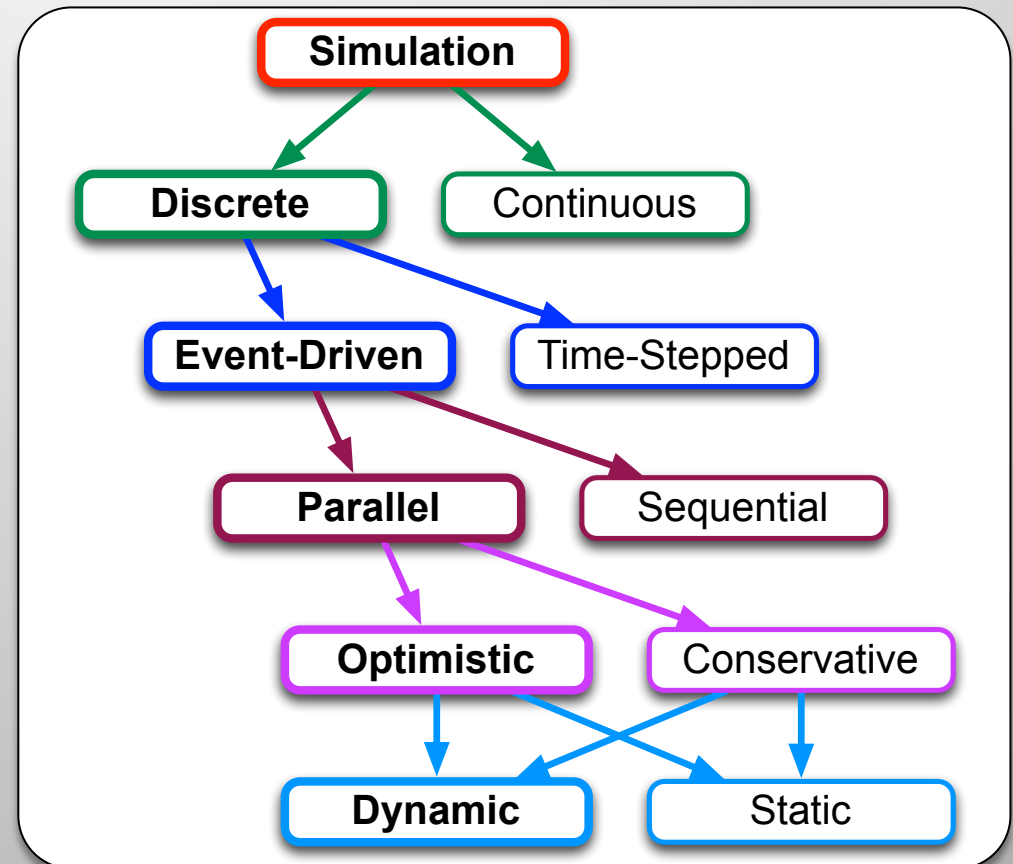  - Null-message synchronization

# Classification of Simulation Techniques

Mathematical Paradigm

Static *vs.* dynamic time control

Parallelism?

Synchronization Style

Load Distribution

```
                    Simulation

         Discrete            Continuous

    Event-Driven        Time-Stepped

         Parallel            Sequential

    Optimistic          Conservative

         Dynamic             Static
```

# Mathematical Paradigm

| Aspect | Discrete | Continuous |
|---|---|---|
| Form of model | Discrete systems<br>Automata, agents, particle systems, stochastic processes, *etc.* | Ordinary or partial differential equations |
| Time, space, state | Continuous or discrete | All continuous |
| State changes | Discontinuous in time<br>Constant between state changes | Continuous in time<br>Occasional discontinuities<br>Piecewise differentiable |
| Mathematical tools | Probability and statistics | Numerical analysis |

- Discrete simulation is natural when there are no underlying physical equations
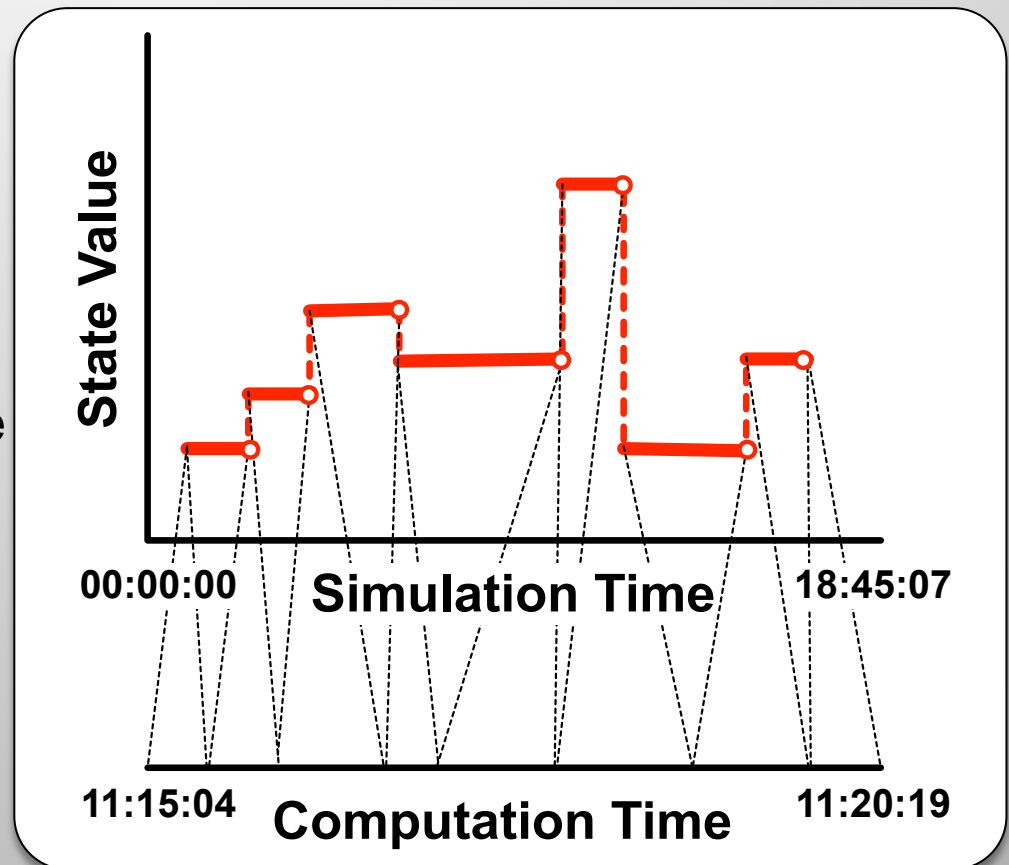
# Time Control

| Aspect | **Event-Driven** | Time-Stepped |
|---|---|---|
| Event times | Dynamically computed | Statically chosen |
| Time resolution | (Ideally) Floating point time<br>Zero lower limit on resolution: inherently multi-scale | (Usually) Integer time<br>Nonzero lower limit on resolution |
| Event distribution in space and time | Sparse and irregular | Dense and regular |
| Appropriate for | Irregular, asynchronous and/or multi-scale models | Spatially *and* temporally regular models |

- Event-driven execution imposes no timescale
  - Supports simulation with wide dynamic range in natural time scales and/or long quiescent periods
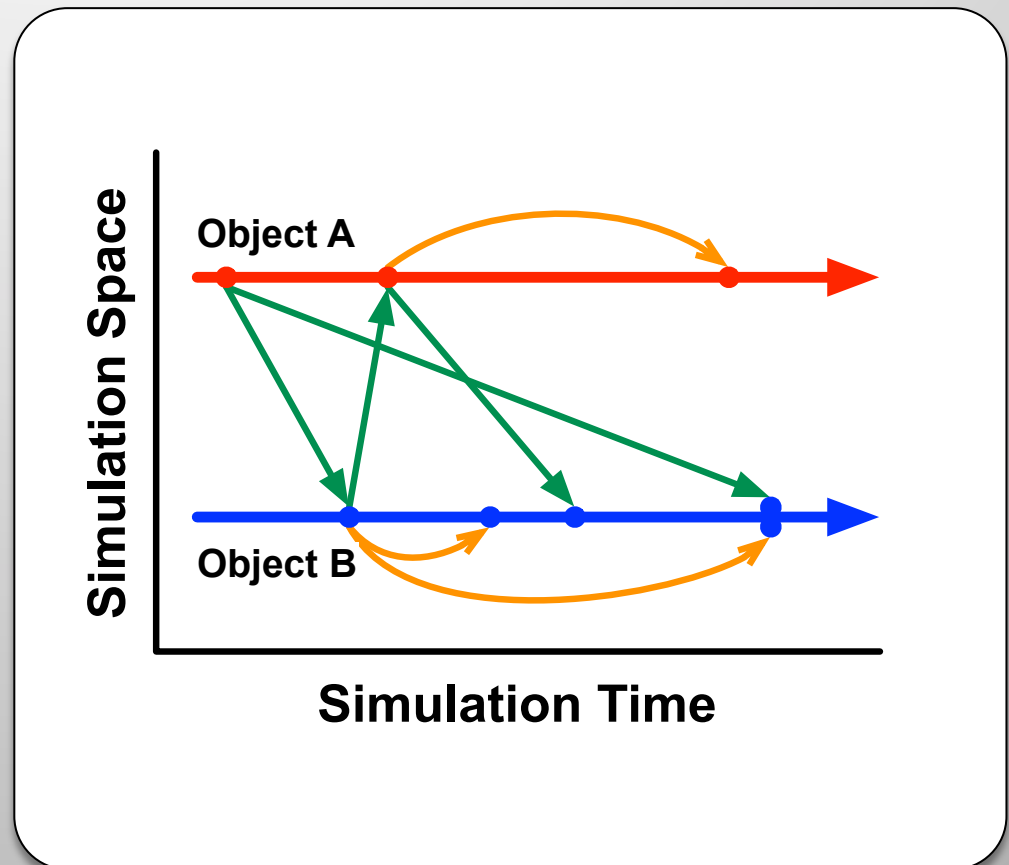
.ıllnS-3

# DES State and Time Evolution

- State values change discontinuously in simulation time
  - Constant between state changes
  - Time interval between state changes is not fixed

- Computation (real world time) required to compute new state value
  - Computation occurs at a fixed value of simulation time

- Rate of model evolution not fixed
  - Faster or slower than real time
  - (Best effort) real time, to interoperate with external real systems

# DES Event Scheduling

- Objects communicate by sending *messages* = schedule events
  - *Event* is a function call, to be executed $\Delta t$ in the future–*no backwards arrows*
  - Typically an event schedules one or more future events

- All event types allowed
  - Event to self
  - Event sends multiple messages
  - Event sends no messages
  - Events can tie
  - Non-FIFO scheduling



Object A

Object B

**Simulation Space**

**Simulation Time**

# Sequential DES Main Event Loop

```
createInitialObjects();
eventList.insert(initialEvents);            // Priority queue on event time

while ( !(terminationCondition() || eventList.empty()) ) do
{
  event e = eventList.removeMinSimTime(); // Choose next event

  simTime = e.getEventTime();                // Set virtual time and unpack event
  object = e.getEventObject();
  method = e.getMethod();
  args = e.getArgs();

  object.method(args);                       // Invoke the event method
                                             // May change state of object
                                             // May schedule future events
                                             // May create or destroy objects
                                             // May cancel (delete) future events

}

finalize();
```
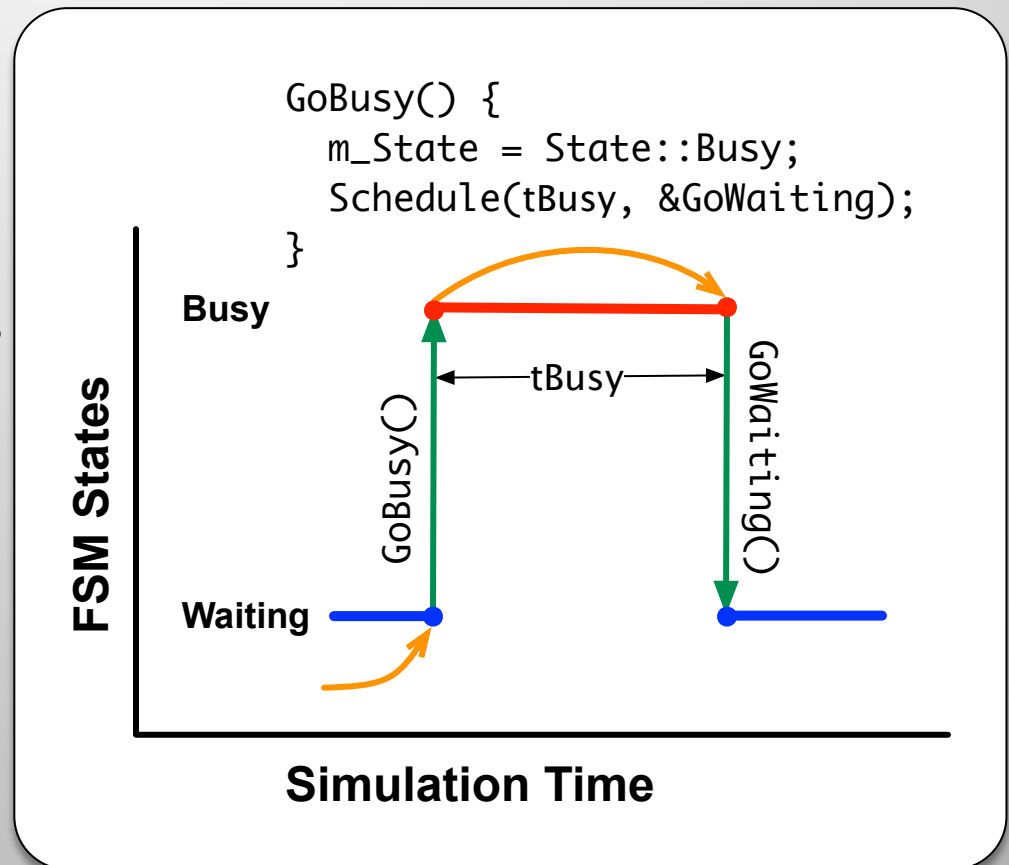
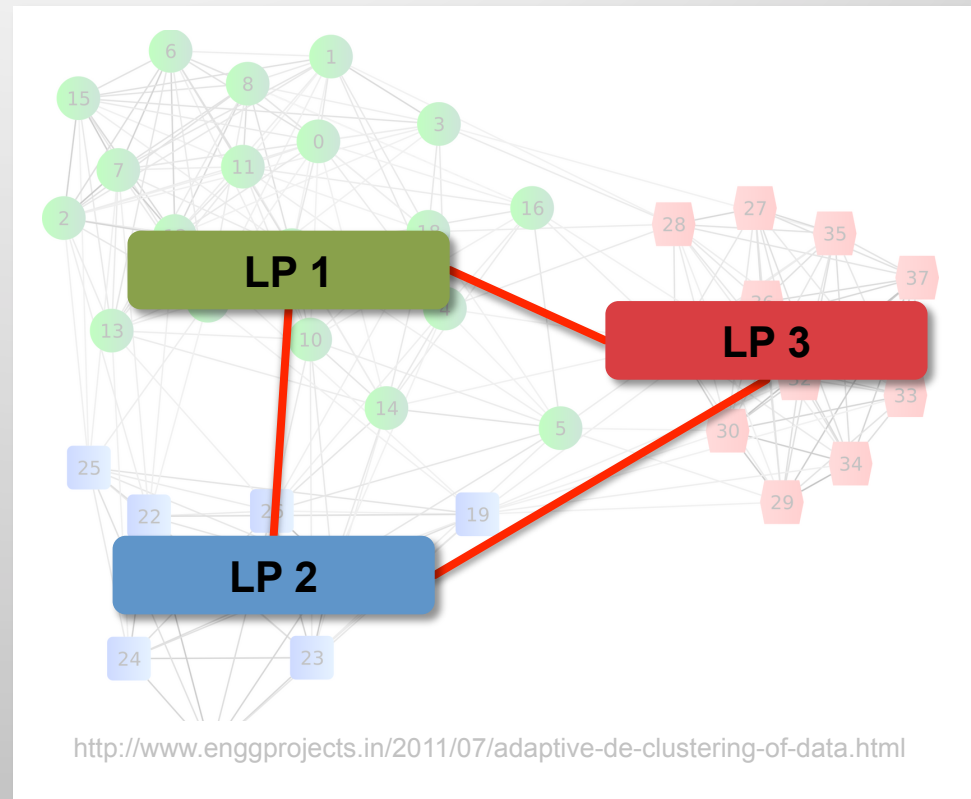**Clean separation between *Simulator* and *Application Model*.**

.....nS-3

# Modeling Time-Consuming Process

- **Model state values change at an instant in simulation time**
  - So how to model time-consuming processes?

- **Finite state machine**
  - Model state is the FSM state
  - Events cause FSM transitions, schedule future transitions

```
GoBusy() {
    m_State = State::Busy;
    Schedule(tBusy, &GoWaiting);
}
```



**FSM States** — **Busy**, **Waiting**, tBusy, GoBusy(), GoWaiting(), **Simulation Time**
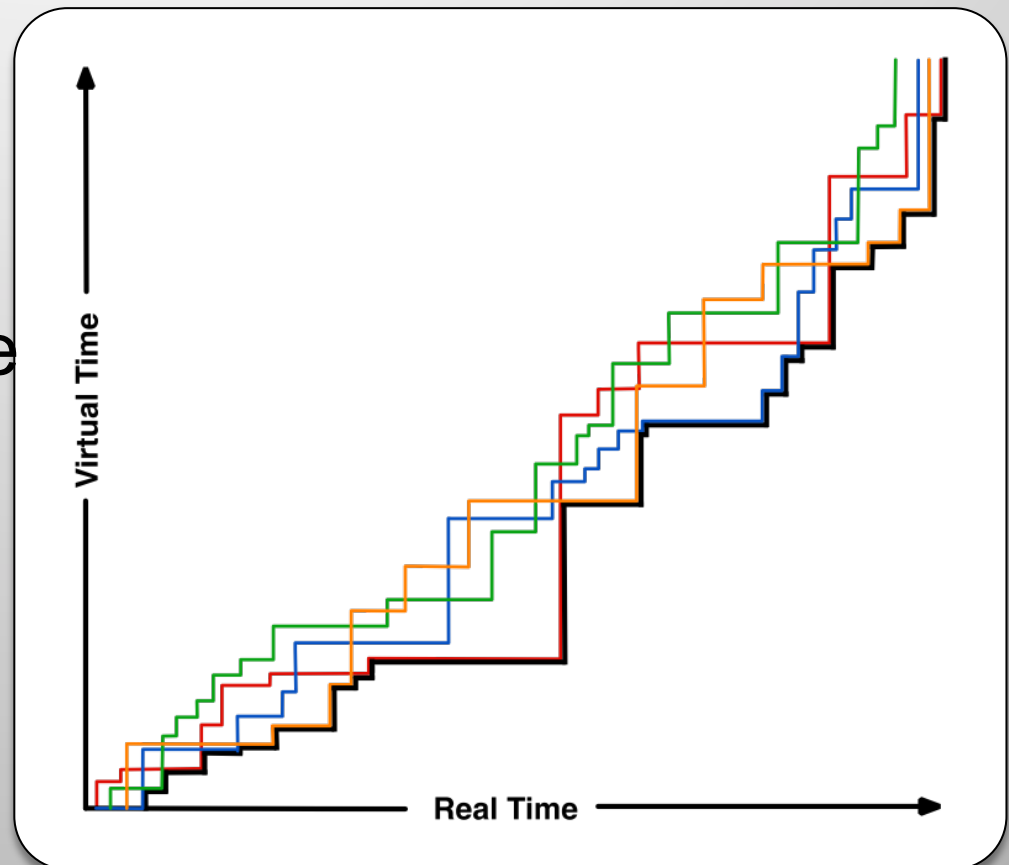
# Parallel Discrete Event Simulation

- **Decompose model into** *Logical Processes*
  - Separate objects and event queues
  - Execute independently
  - Events for other LPs become messages
  - ~ MPI Ranks



http://www.enggprojects.in/2011/07/adaptive-de-clustering-of-data.html

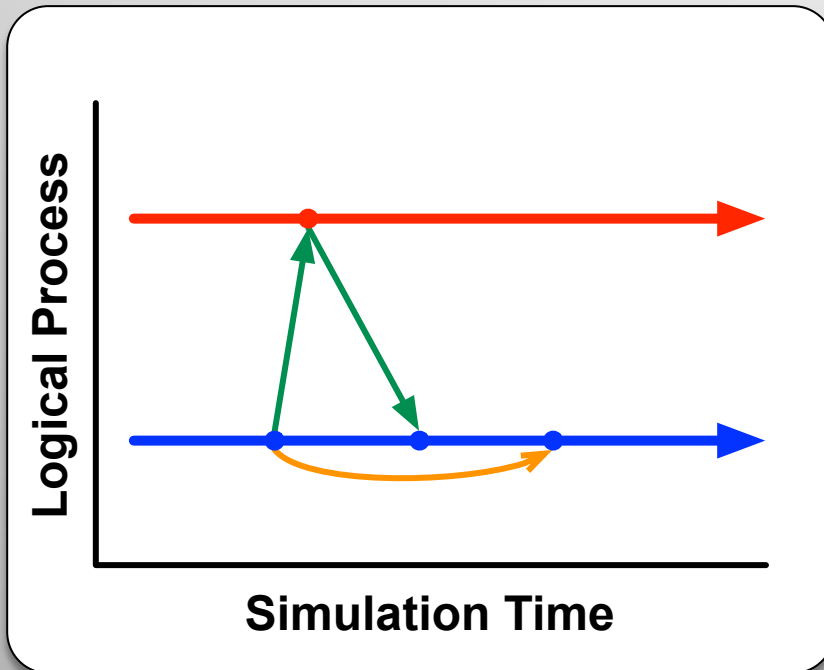**Parallel execution *must* produce exact same results as sequential!**

# PDES Execution: LPs Advance Independently

- **Sometimes ahead in virtual time, sometimes behind**

- **More or less real time per event**

- **Never backwards!**
  - Hallmark of conservative execution (Ask me about optimistic execution ☺)
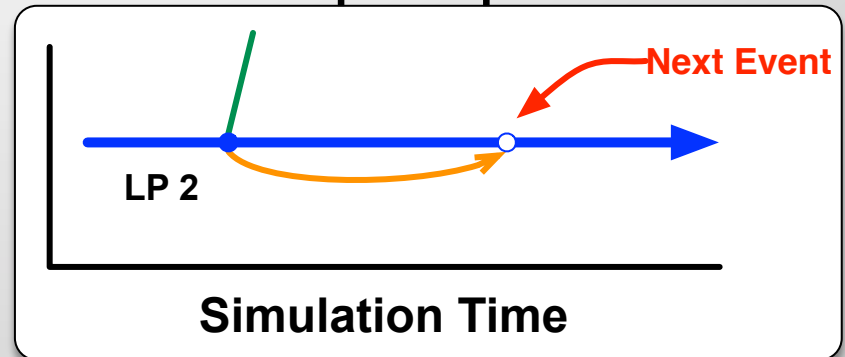


**ıllNS-3**

# Need for Synchronization of LPs: Prevent Causality Violations

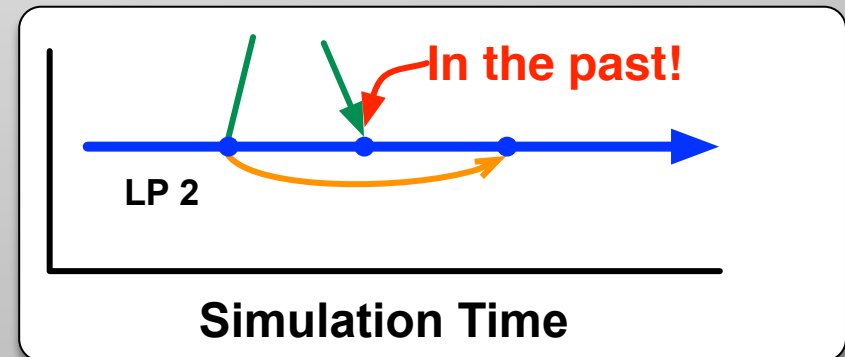- Sequential event sequence

- LP 2's perspective:



- Arrival of LP1 event



**Need to guarantee no messages arrive in the past (for conservative PDES).**

# Granted Time Window Synchronization

- *If* we could guarantee no remote events will arrive before *GrantedTime*
  - All events before *GrantedTime* are safe
  - At *GrantedTime* need to synchronize:
    — Receive and schedule events from other LPs
    — Compute new *GrantedTime*

- Performance
  - Even workload distribution limits cpu idle time
  - Maximize *GrantedTime* to execute more events in parallel between synchronization

# Lookahead and LBTS Provide the Granted Time Guarantee

- Model must provide *Lookahead*
  - Minimum delay for remote events
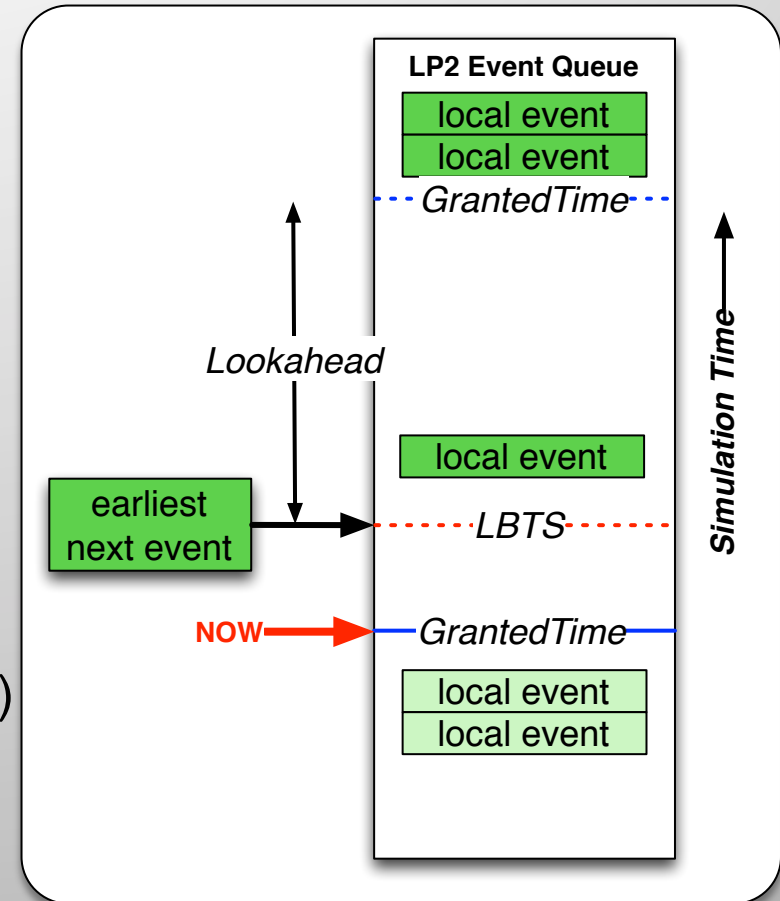    - Example: network channel link latency + transmission time for smallest packet

- Lower Bound Time Stamp (*LBTS*)
  - Min next event time across all LPs

- *GrantedTime = LBTS + Lookahead*

- Synchronization across LPs is expensive
  - Typically barrier (to wait for slowest LP)
  - Plus (at least one) all gather or reduction
  - Each of these is $\log(N_{LP})$ in time

**LP2 Event Queue**

local event
local event
- - - *GrantedTime* - - -

*Lookahead*

local event

earliest next event · · · · *LBTS* · · · ·

NOW → — *GrantedTime* —

local event
local event

*Simulation Time* →

**Finding large *Lookahead* is key to performance**

**....ıllNS-3**

# Null-Message Alternative for Static and Sparse LP Graphs

- *GrantedTime* assumes all LPs can message all other LPs

- But my LP graph is sparse. Why synchronize with everyone?
  - Every message sent could communicate my virtual time
  - Guarantee at least one message every *Lookahead*
  - Send *Null-message* when necessary



**Simulation Time**

**Simulation Time**

34

# To Learn More…

- Much of this material from a short course presented spring 2014
  - David Jefferson, LLNL, co-inventor of Optimistic PDES
  - 15 sessions
    - Sequential DES
    - Ties, LBTS, Lookahead
    - Chandy, Misra, Bryant: YAWNS
    - Deadlock
    - Null Messages
    - Dynamic Object Creation
    - Critical Path
    - Speedup
    - Optimistic DES, TimeWarp
    - Global Virtual Time
    - Commitment
    - Checkpointing
    - Rollback and Reverse Computation
    - Dynamic Load Balancing
    - Mixed Discrete and Continuous
  - Slides and videos publicly available:
    http://pdes-course-2014.ucllnl.org

# Parallel ns-3

- History

- PDES in ns-3

- Mechanics
  - Enabling
  - Running

- PDES Simulators
  - GrantedTime
  - NullMessage

- Parallel Models 1
  - The easy way

- Lookahead

- Under the covers:
  - PointToPointRemoteChannel
  - PointToPointNetDevice

- Parallel Models 2
  - The hard way
  - Limitations

# Parallel ns-3 History

- Initial release in ns-3.8
  - J. Pelkey and G. Riley, "Distributed Simulation with MPI in ns-3," WNS3 2011, Barcelona, Spain.
  - Roots from:
    — Parallel/Distributed ns (pdns)
    — Georgia Tech Network Simulator (GTNetS)

- Publications
  - K. Renard, *et al*, "A Performance and Scalability Evaluation of the ns-3 Distributed Scheduler," SimuTools 2012
  - S. Nikolaev, *et al*, "Performance of Distributed ns-3 Network Simulator," SimuTools 2013
  - WNS3 2015

# PDES in ns-3

## Sequential ns-3

- LP is implicit
  - `ns3::Simulator`

- Event messages
  - Explicit future function calls

    `Schedule (delay, &fn,…)`

- Virtual time discipline

  `DefaultSimulatorImpl`
  `RealtimeSimulatorImpl`
  `VisualSimulatorImpl`

## Parallel ns-3

- Each rank is an LP

- Event messages
  - Local to LP: explicit future function calls
  - Remote: implicit message send

- Virtual time discipline

  `DistributedSimulatorImpl`
  `NullMessageSimulatorImpl`

- Lookahead (later)

# Enabling Parallel ns-3

- Configure with `--enable-mpi`
  - Tries to run `mpic++`
    - — Recognizes OpenMPI and MPICH libraries
  - Defines NS3_MPI and either NS3_OPENMPI or NS3_MPICH

- Followed by usual build

**Configuring ns-3 With MPI**

```
$ ./waf configure --enable-mpi
Setting top to                        : ...
...
---- Summary of optional NS-3 features:
Build profile                  : debug
...
MPI Support                    : enabled
...
'configure' finished successfully (1.295s)

$ ./waf build
...
```

# Running Parallel ns-3 Scripts

- Waf can't distinguish sequential and parallel
  - Need to specify `mpirun` and number of ranks explicitly

## Running Parallel Scripts with `waf` and `mpirun`

```
$ ./waf --run simple-distributed
Waf: Entering directory `build/debug'
Waf: Leaving directory `build/debug'
'build' finished successfully (2.118s)
This simulation requires 2 and only 2 logical processors.
Command ['build/debug/src/mpi/examples/ns3-dev-simple-distributed-debug'] exited with code 1

#  Multiple ranks on a single computer:
$ ./waf --run simple-distributed --command-template="mpirun -np 2 %s"
Waf: Entering directory `build/debug'
Waf: Leaving directory `build/debug'
'build' finished successfully (2.104s)
At time 1.02264s packet sink received 512 bytes from 10.1.1.1 port 49153 total Rx 512 bytes
At time 1.0235s packet sink received 512 bytes from 10.1.2.1 port 49153 total Rx 512 bytes
At time 1.02437s packet sink received 512 bytes from 10.1.3.1 port 49153 total Rx 512 bytes
At time 1.02524s packet sink received 512 bytes from 10.1.4.1 port 49153 total Rx 512 bytes

# Multiple computers:
$ mpirun -np 2 ./waf -run simple-distributed
```

# Switching Between GrantedTime and NullMessage Simulators

- ## Use environment variable

```
$ NS_GLOBAL_VALUE=\
  "SimulatorImplementationType=ns3::NullMessageSimulatorImpl"
\
  ./waf --run ...
```

- ## Use command line:

**Selecting the Parallel Simulator from the Command Line**

```
bool nullmsg = false;
CommandLine cmd;
cmd.AddValue ("nullmsg", "Enable the use of null-message synchronization", nullmsg);
cmd.Parse (argc,argv);
...
if(nullmsg) {
  GlobalValue::Bind ("SimulatorImplementationType",
                     StringValue ("ns3::NullMessageSimulatorImpl"));
} else {
  GlobalValue::Bind ("SimulatorImplementationType",
                     StringValue ("ns3::DistributedSimulatorImpl"));
}
MpiInterface::Enable (&argc, &argv);
```

# Constructing Distributed Models The Easy Way

- **All ranks construct the full topology**
  - All Nodes, NetDevices and Channels
    - — Label Nodes with rank: `Node::Node (uint32_t systemId)`
  - All Internet stacks and addresses
  - Good
    - — Single code for model construction, runs sequential and parallel
    - — Event execution happens in parallel
    - — Enables GOD and NIx-vector routing to work
  - Bad
    - — Memory is used for nodes/stacks/devices that "belong" to other ranks (But come to my talk tomorrow ☺)

- **Install local applications only**
  - Non-local nodes (not on my rank) should not have applications

# Where to Get Lookahead?

- Primarily from link latency

- What about shared channels like CSMA or wireless?
  - Latency can be zero
  - Multiple NetDevices
  - ➢ *Can't span ranks!*

- Only PointToPoint links can cross ranks
  - Global *Lookahead* is smallest cross-rank latency

# Under the Covers: PointToPointHelper::Install

```
bool useNormalChannel = true;
Ptr<PointToPointChannel> channel = 0;

if (MpiInterface::IsEnabled ())  {
  uint32_t currSystemId = MpiInterface::GetSystemId ();
  if (a->GetSystemId () != currSystemId ||
      b->GetSystemId () != currSystemId) {
    useNormalChannel = false;
  }
}
if (useNormalChannel) {
  channel =
    m_channelFactory.Create<Po
} else {
  channel =
    m_remoteChannelFactory.Cre

  Ptr<MpiReceiver> mpiRecA = C
  mpiRecA->SetReceiveCallback
    (MakeCallback (&PointToPoi
  devA->AggregateObject (mpiRe

  // Same for b
}
```



Sequential

Distributed

# Under the Covers: Sending a Packet from PointToPointNetDevice

**PointToPointNetDevice Call Chain**

```
PointToPointNetDevice::Send() {
  TransmitStart() {
    PointToPointRemoteChannel::TransmitStart() {
      MpiInterface::SendPacket();
```

- **MpiInterface::SendPacket()**

  — *Packet data*

  — *Receive time* – Local Now() + Latency + Packet Tx duration

  — Remote SystemId (rank)

  — Remote *NodeId*

  — Remote *InterfaceId*

- Serialize packet and destination data

- Send to remote rank with non-blocking MPI_Isend()

# Under the Covers: Getting a Remote Packet to the `PointToPointNetDevice`

At end of *GrantedTime*, `DistributedSimulatorImpl` calls `GrantedTimeWindowMpiInterface::ReceiveMessages()`

- Reads all pending MPI messages
  - Deserialize target *Receive time*, *NodeId* and *InterfaceId*
  - Deserialize *packet data*
  - Find Node by *NodeId*
  - Find `NetDevice` on Node with correct InterfaceId
  - Get `MpiReceiver` object aggregated to the `NetDevice`
    - `MpiReceiver` merely holds the correct `NetDevice` Callback
  - Schedule `MpiReceiver::Receive` event at *Receive time*

# Building a Distributed ns-3 Simulation

- **Choose partitioning strategy**
  - Label contiguous regions which can't be partitioned
    - CSMA and wireless
  - Select regions which will share a rank
    - Find large point-to-point latencies for good *Lookahead*
    - Minimize communication between ranks

- **Build topology as normal, assigning Nodes to ranks**
  `CreateObject<Node> (rankId)`

- **Rewrite topology to improve partitioning**
  - CSMA with only 2 nodes
  - Move latency

# Constructing Distributed Models
# The Hard Way

- Use the ghost cell design pattern to save memory
  - Only create local Nodes, Applications, Internet stacks, NetDevices and Channels
  - Plus "ghost" nodes: remote endpoint of PointToPointRemoteChannel

- Requires *manual intervention*
  - Global and NIX routing do not see entire topology
    - Add static, default routes manually.  Hint: IPv6 allows for more "aggregatable" routes
  - Ghost nodes will likely have incorrect remote NodeId, InterfaceId
  - Must align interface identifiers by hand in same fashion

# Limitations of Distributed NS3

- Partitioning is a manual process

- Partitioning is restricted to Point-To-Point links only

  - Partitioning within a wireless network is not supported
    - *Lookahead* is very small and dynamic

- Need full topology in all LPs

  - Exception with careful node ordering, interface numbering, and manual routing

**examples/tutorial/third.cc**

**src/mpi/examples/third-distributed.cc**

**(These have diverged slightly in ns-3-dev.  Differences minimized here.)**

> 1. Include mpi-module.h
> 2. Same topology, split across Point-to-point link

**Left file (examples/tutorial/third.cc):**

```
1  * -*- Mode:C++; c-file-style:"gnu"; indent-tabs-mode:nil; -*- */
2
3
4
5
6
7  * This program is distributed in the hope that it will be useful,
8  * but WITHOUT ANY WARRANTY; without even the implied warranty of
9  * MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE.  See the
10 * GNU General Public License for more details.
11 *
12 * You should have received a copy of the GNU General Public License
13 * along with this program; if not, write to the Free Software
14 * Foundation, Inc., 59 Temple Place, Suite 330, Boston, MA  02111-1307  USA
15 */
16
17 include "ns3/core-module.h"
18 include "ns3/point-to-point-module.h"
19 include "ns3/network-module.h"
20 include "ns3/applications-module.h"
21 include "ns3/wifi-module.h"
22 include "ns3/mobility-module.h"
23 include "ns3/csma-module.h"
24 include "ns3/internet-module.h"
25
26 / Default Network Topology
27 /
28 /    Wifi 10.1.3.0
29 /                  AP
30 / *    *    *    *
31 / |    |    |    |    10.1.1.0
32 / n5   n6   n7   n0 -------------- n1   n2   n3   n4
33 /                  point-to-point  |    |    |    |
34 /                                  ================
35 /                                    LAN 10.1.2.0
36
37 sing namespace ns3;
38
39 S_LOG_COMPONENT_DEFINE ("ThirdScriptExample");
40
41 nt
42 ain (int argc, char *argv[])
43
44   bool verbose = true;
45   uint32_t nCsma = 3;
46   uint32_t nWifi = 3;
47   bool tracing = false;
```

**Right file (src/mpi/examples/third-distributed.cc):**

```
1  * -*- Mode:C++; c-file-style:"gnu"; indent-tabs-mode:nil; -*- */
2
3    program is free software; you can redistribute it and/or modify
4    er the terms of the GNU General Public License version 2 as
5    hed by the Free Software Foundation;
6
7  * This program is distributed in the hope that it will be useful,
8  * but WITHOUT ANY WARRANTY; without even the implied warranty of
9  * MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE.  See the
10 * GNU General Public License for more details.
11 *
12 * You should have received a copy of the GNU General Public License
13 * along with this program; if not, write to the Free Software
14 * Foundation, Inc., 59 Temple Place, Suite 330, Boston, MA  02111-1307  USA
15 */
16
17 include "ns3/core-module.h"
18 include "ns3/point-to-point-module.h"
19 include "ns3/network-module.h"
20 include "ns3/applications-module.h"
21 include "ns3/wifi-module.h"
22 include "ns3/mobility-module.h"
23 include "ns3/csma-module.h"
24 include "ns3/internet-module.h"
25
26 include "ns3/mpi-module.h"
27
28 / Default Network Topology
29 / (same as third.cc from tutorial)
30 / Distributed simulation, split along the p2p link
31 / Number of wifi or csma nodes can be increased up to 250
32 /                         |
33 /               Rank 0    |    Rank 1
34 / -------------------------|-------------------------
35 /    Wifi 10.1.3.0
36 /                  AP
37 / *    *    *    *
38 / |    |    |    |    10.1.1.0
39 / n5   n6   n7   n0 -------------- n1   n2   n3   n4
40 /                  point-to-point  |    |    |    |
41 /                                  ================
42 /                                    LAN 10.1.2.0
43
44 sing namespace ns3;
45
46 S_LOG_COMPONENT_DEFINE ("ThirdExampleDistributed");
47
```

**examples/tutorial/third.cc**    **src/mpi/examples/third-distributed.cc**



1. Different log component name
2. Command line argument to select Null message

Left panel (examples/tutorial/third.cc):

```
25
26
27
28
29
30
31  // |    |    |    |   10.1.1.0
32  // n5   n6   n7   n0 ------------- n1   n2   n3   n4
33  //                  point-to-point |    |    |    |
34  //                                 =====================
35  //                                  LAN 10.1.2.0
36
37  using namespace ns3;
38
39  NS_LOG_COMPONENT_DEFINE ("ThirdScriptExample");
40
41  int
42  main (int argc, char *argv[])
43  {
44    bool verbose = true;
45    uint32_t nCsma = 3;
46    uint32_t nWifi = 3;
47    bool tracing = false;
48
49    CommandLine cmd;
50    cmd.AddValue ("nCsma", "Number of \"extra\" CSMA nodes/devices", nCsma);
51    cmd.AddValue ("nWifi", "Number of wifi STA devices", nWifi);
52    cmd.AddValue ("verbose", "Tell echo applications to log if true", verbose);
53    cmd.AddValue ("tracing", "Enable pcap tracing", tracing);
54
55    cmd.Parse (argc,argv);
56
57    // Check for valid number of csma or wifi nodes
58    // 250 should be enough, otherwise IP addresses
59    // soon become an issue
60    if (nWifi > 250 || nCsma > 250)
61      {
62        std::cout << "Too many wifi or csma nodes, no more than 250 each." << std::
63        return 1;
64      }
65
66    if (verbose)
67      {
68        LogComponentEnable ("UdpEchoClientApplication", LOG_LEVEL_INFO);
69        LogComponentEnable ("UdpEchoServerApplication", LOG_LEVEL_INFO);
70      }
71
```

Right panel (src/mpi/examples/third-distributed.cc):

```
33  //                   Rank 0   |  Rank 1
            --------------|--------------------
                              .3.0
                           AP
             *      *
    //  |     |     |     |   10.1.1.0
39  // n5   n6   n7   n0 ------------- n1   n2   n3   n4
40  //                  point-to-point |    |    |    |
41  //                                 =====================
42  //                                  LAN 10.1.2.0
43
44  using namespace ns3;
45
46  NS_LOG_COMPONENT_DEFINE ("ThirdExampleDistributed");
47
48  int
49  main (int argc, char *argv[])
50  {
51    bool verbose = true;
52    uint32_t nCsma = 3;
53    uint32_t nWifi = 3;
54    bool tracing = false;
55    bool nullmsg = false;
56
57    CommandLine cmd;
58    cmd.AddValue ("nCsma", "Number of \"extra\" CSMA nodes/devices", nCsma);
59    cmd.AddValue ("nWifi", "Number of wifi STA devices", nWifi);
60    cmd.AddValue ("verbose", "Tell echo applications to log if true", verbose);
61    cmd.AddValue ("tracing", "Enable pcap tracing", tracing);
62    cmd.AddValue ("nullmsg", "Enable the use of null-message synchronization",
63
64    cmd.Parse (argc,argv);
65
66    // Check for valid number of csma or wifi nodes
67    // 250 should be enough, otherwise IP addresses
68    // soon become an issue
69    if (nWifi > 250 || nCsma > 250)
70      {
71        std::cout << "Too many wifi or csma nodes, no more than 250 each." << std::
72        return 1;
73      }
74
75    if (verbose)
76      {
77        LogComponentEnable ("UdpEchoClientApplication", LOG_LEVEL_INFO);
78        LogComponentEnable ("UdpEchoServerApplication", LOG_LEVEL_INFO);
79
```

**examples/tutorial/third.cc**

**src/mpi/examples/third-distributed.cc**

```
49   CommandLine cmd;
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66   if (verbose)
67     {
68       LogComponentEnable ("UdpEchoClientApplication", LOG_LEVEL_INFO);
69       LogComponentEnable ("UdpEchoServerApplication", LOG_LEVEL_INFO);
70     }
71
72   NodeContainer p2pNodes;
73   p2pNodes.Create (2);
74
75   PointToPointHelper pointToPoint;
76   pointToPoint.SetDeviceAttribute ("DataRate", StringValue ("5Mbps"));
77   pointToPoint.SetChannelAttribute ("Delay", StringValue ("2ms"));
78
79   NetDeviceContainer p2pDevices;
80   p2pDevices = pointToPoint.Install (p2pNodes);
81
82   NodeContainer csmaNodes;
83   csmaNodes.Add (p2pNodes.Get (1));
84   csmaNodes.Create (nCsma);
85
86   CsmaHelper csma;
87   csma.SetChannelAttribute ("DataRate", StringValue ("100Mbps"));
88   csma.SetChannelAttribute ("Delay", TimeValue (NanoSeconds (6560)));
89
90   NetDeviceContainer csmaDevices;
91   csmaDevices = csma.Install (csmaNodes);
92
93   NodeContainer wifiStaNodes;
94   wifiStaNodes.Create (nWifi);
```

```
80
81   // Sequential fallback values
82   uint32_t systemId = 0;
83   uint32_t systemCount = 1;
84
85   #ifdef NS3_MPI
86
87   // Distributed simulation setup; by default use granted time window algorithm.
88   if(nullmsg)
89     {
90       GlobalValue::Bind ("SimulatorImplementationType",
91                          StringValue ("ns3::NullMessageSimulatorImpl"));
92     }
93   else
94     {
95       GlobalValue::Bind ("SimulatorImplementationType",
96                          StringValue ("ns3::DistributedSimulatorImpl"));
97     }
98
99   MpiInterface::Enable (&argc, &argv);
100
101  systemId = MpiInterface::GetSystemId ();
102  systemCount = MpiInterface::GetSize ();
103
104  // Check for valid distributed parameters.
105  // Must have 2 and only 2 Logical Processors (LPs)
106  if (systemCount != 2)
107    {
108      std::cout << "This simulation requires 2 and only 2 logical processors." <<
109      return 1;
110    }
111
112  #endif // NS3_MPI
113
114  // System id of Wifi side
115  uint32_t systemWifi = 0;
116
117  // System id of CSMA side
118  uint32_t systemCsma = systemCount - 1;
119
120  NodeContainer p2pNodes;
121  Ptr<Node> p2pNode1 = CreateObject<Node> (systemWifi); // Create node with rank
122  Ptr<Node> p2pNode2 = CreateObject<Node> (systemCsma); // Create node with rank
123  p2pNodes.Add (p2pNode1);
124  p2pNodes.Add (p2pNode2);
125
126  PointToPointHelper pointToPoint;
```

1. Condition on NS3_MPI
2. Null message selector
3. Initialize MPI
4. Get rank #, number of ranks
5. Check number of ranks
6. Use symbolic names for each rank
7. Create point-to-point nodes

1
2
3
4
5
6
7

**examples/tutorial/third.cc**

**src/mpi/examples/third-distributed.cc**

Left column (examples/tutorial/third.cc):

```
66    LogComponentEnable ("UdpEchoClientApplication", LOG_LEVEL_INFO);
67
68
69
70
71
72
73    PointToPointHelper pointToPoint;
74    pointToPoint.SetDeviceAttribute ("DataRate", StringValue ("5Mbps"));
75    pointToPoint.SetChannelAttribute ("Delay", StringValue ("2ms"));
76
77    NetDeviceContainer p2pDevices;
78    p2pDevices = pointToPoint.Install (p2pNodes);
79
80    NodeContainer csmaNodes;
81    csmaNodes.Add (p2pNodes.Get (1));
82    csmaNodes.Create (nCsma);
83
84    CsmaHelper csma;
85    csma.SetChannelAttribute ("DataRate", StringValue ("100Mbps"));
86    csma.SetChannelAttribute ("Delay", TimeValue (NanoSeconds (6560)));
87
88    NetDeviceContainer csmaDevices;
89    csmaDevices = csma.Install (csmaNodes);
90
91    NodeContainer wifiStaNodes;
92    wifiStaNodes.Create (nWifi);
93    NodeContainer wifiApNode = p2pNodes.Get (0);
94
95    YansWifiChannelHelper channel = YansWifiChannelHelper::Default ();
96    YansWifiPhyHelper phy = YansWifiPhyHelper::Default ();
97    phy.SetChannel (channel.Create ());
98
99    WifiHelper wifi = WifiHelper::Default ();
100   wifi.SetRemoteStationManager ("ns3::AarfWifiManager");
101
102   NqosWifiMacHelper mac = NqosWifiMacHelper::Default ();
103
104   Ssid ssid = Ssid ("ns-3-ssid");
105   mac.SetType ("ns3::StaWifiMac",
106           "Ssid", SsidValue (ssid),
107           "ActiveProbing", BooleanValue (false));
108
109   NetDeviceContainer staDevices;
110   staDevices = wifi.Install (phy, mac, wifiStaNodes);
111
112   mac.SetType ("ns3::ApWifiMac",
```

Callout (on left column):
1. Create CSMA nodes on one rank
2. Create Wifi nodes on another rank

Right column (src/mpi/examples/third-distributed.cc):

```
119
120   NodeContainer p2pNodes;
121   Ptr<Node> p2pNode1 = CreateObject<Node> (systemWifi); // Create node with rank 0
122   Ptr<Node> p2pNode2 = CreateObject<Node> (systemCsma); // Create node with rank 1
123   p2pNodes.Add (p2pNode1);
124   p2pNodes.Add (p2pNode2);
125
126   PointToPointHelper pointToPoint;
127   pointToPoint.SetDeviceAttribute ("DataRate", StringValue ("5Mbps"));
128   pointToPoint.SetChannelAttribute ("Delay", StringValue ("2ms"));
129
130   NetDeviceContainer p2pDevices;
131   p2pDevices = pointToPoint.Install (p2pNodes);
132
133   NodeContainer csmaNodes;
134   csmaNodes.Add (p2pNodes.Get (1));
135   csmaNodes.Create (nCsma, systemCsma);   // Create csma nodes with rank 1
136
137   CsmaHelper csma;
138   csma.SetChannelAttribute ("DataRate", StringValue ("100Mbps"));
139   csma.SetChannelAttribute ("Delay", TimeValue (NanoSeconds (6560)));
140
141   NetDeviceContainer csmaDevices;
142   csmaDevices = csma.Install (csmaNodes);
143
144   NodeContainer wifiStaNodes;
145   wifiStaNodes.Create (nWifi, systemWifi); // Create wifi nodes with rank 0
146   NodeContainer wifiApNode = p2pNodes.Get (0);
147
148   YansWifiChannelHelper channel = YansWifiChannelHelper::Default ();
149   YansWifiPhyHelper phy = YansWifiPhyHelper::Default ();
150   phy.SetChannel (channel.Create ());
151
152   WifiHelper wifi = WifiHelper::Default ();
153   wifi.SetRemoteStationManager ("ns3::AarfWifiManager");
154
155   NqosWifiMacHelper mac = NqosWifiMacHelper::Default ();
156
157   Ssid ssid = Ssid ("ns-3-ssid");
158   mac.SetType ("ns3::StaWifiMac",
159           "Ssid", SsidValue (ssid),
160           "ActiveProbing", BooleanValue (false));
161
162   NetDeviceContainer staDevices;
163   staDevices = wifi.Install (phy, mac, wifiStaNodes);
164
165   mac.SetType ("ns3::ApWifiMac",
```

ns-3

## examples/tutorial/third.cc

## src/mpi/examples/third-distributed.cc

**1. Install devices, addresses and Internet stack everywhere**
**2. Install applications only on rank-local nodes**

Left panel (examples/tutorial/third.cc):

```
136  stack.Install (csmaNodes);
137
138
139
140
141
142
143
144
145
146
147  Ipv4InterfaceContainer csmaInterfaces;
148  csmaInterfaces = address.Assign (csmaDevices);
149
150  address.SetBase ("10.1.3.0", "255.255.255.0");
151  address.Assign (staDevices);
152  address.Assign (apDevices);
153
154  UdpEchoServerHelper echoServer (9);
155
156  ApplicationContainer serverApps = echoServer.Install (csmaNodes.Get (nCsma));
157  serverApps.Start (Seconds (1.0));
158  serverApps.Stop (Seconds (10.0));
159
160  UdpEchoClientHelper echoClient (csmaInterfaces.GetAddress (nCsma), 9);
161  echoClient.SetAttribute ("MaxPackets", UintegerValue (1));
162  echoClient.SetAttribute ("Interval", TimeValue (Seconds (1.0)));
163  echoClient.SetAttribute ("PacketSize", UintegerValue (1024));
164
165  ApplicationContainer clientApps =
166    echoClient.Install (wifiStaNodes.Get (nWifi - 1));
167  clientApps.Start (Seconds (2.0));
168  clientApps.Stop (Seconds (10.0));
169
170  Ipv4GlobalRoutingHelper::PopulateRoutingTables ();
171
172  Simulator::Stop (Seconds (10.0));
173
174  if (tracing == true)
175    {
176      pointToPoint.EnablePcapAll ("third-distributed-wifi");
177      phy.EnablePcap ("third-distributed-wifi", apDevices.Get (0));
178      csma.EnablePcap ("third-distributed-wifi", csmaDevices.Get (0), true);
179
180      pointToPoint.EnablePcapAll ("third-distributed-csma");
181      phy.EnablePcap ("third-distributed-csma", apDevices.Get (0));
```

Right panel (src/mpi/examples/third-distributed.cc):

```
195  address.SetBase ("10.1.1.0", "255.255.255.0");
196  Ipv4InterfaceContainer p2pInterfaces;
197  p2pInterfaces = address.Assign (p2pDevices);
198
199  address.SetBase ("10.1.2.0", "255.255.255.0");
200  Ipv4InterfaceContainer csmaInterfaces;
201  csmaInterfaces = address.Assign (csmaDevices);
202
203  address.SetBase ("10.1.3.0", "255.255.255.0");
204  address.Assign (staDevices);
205  address.Assign (apDevices);
206
207  // If this rank is systemCsma,
208  // it should contain the server application,
209  // since it is on one of the csma nodes
210  if (systemId == systemCsma)
211    {
212      UdpEchoServerHelper echoServer (9);
213
214      ApplicationContainer serverApps = echoServer.Install (csmaNodes.Get (nCsma))
215      serverApps.Start (Seconds (1.0));
216      serverApps.Stop (Seconds (10.0));
217    }
218
219  // If this rank is systemWifi
220  // it should contain the client application,
221  // since it is on one of the wifi nodes
222  if (systemId == systemWifi)
223    {
224      UdpEchoClientHelper echoClient (csmaInterfaces.GetAddress (nCsma), 9);
225      echoClient.SetAttribute ("MaxPackets", UintegerValue (1));
226      echoClient.SetAttribute ("Interval", TimeValue (Seconds (1.0)));
227      echoClient.SetAttribute ("PacketSize", UintegerValue (1024));
228
229      ApplicationContainer clientApps =
230        echoClient.Install (wifiStaNodes.Get (nWifi - 1));
231      clientApps.Start (Seconds (2.0));
232      clientApps.Stop (Seconds (10.0));
233    }
234
235  Ipv4GlobalRoutingHelper::PopulateRoutingTables ();
236
237  Simulator::Stop (Seconds (10.0));
238
239  if (tracing == true)
240    {
```

**examples/tutorial/third.cc**

**src/mpi/examples/third-distributed.cc**

Left code panel:

```
142   address.SetBase ("10.1.1.0", "255.255.255.0");
143
144
145
146
147
148
149
150
151
152
153
154   UdpEchoServerHelper echoServer (9);
155
156   ApplicationContainer serverApps = echoServer.Install (csmaNodes.Get (nCsma));
157   serverApps.Start (Seconds (1.0));
158   serverApps.Stop (Seconds (10.0));
159
160   UdpEchoClientHelper echoClient (csmaInterfaces.GetAddress (nCsma), 9);
161   echoClient.SetAttribute ("MaxPackets", UintegerValue (1));
162   echoClient.SetAttribute ("Interval", TimeValue (Seconds (1.0)));
163   echoClient.SetAttribute ("PacketSize", UintegerValue (1024));
164
165   ApplicationContainer clientApps =
166     echoClient.Install (wifiStaNodes.Get (nWifi - 1));
167   clientApps.Start (Seconds (2.0));
168   clientApps.Stop (Seconds (10.0));
169
170   Ipv4GlobalRoutingHelper::PopulateRoutingTables ();
171
172   Simulator::Stop (Seconds (10.0));
173
174   if (tracing == true)
175     {
176       pointToPoint.EnablePcapAll ("third-distributed-wifi");
177       phy.EnablePcap ("third-distributed-wifi", apDevices.Get (0));
178       csma.EnablePcap ("third-distributed-wifi", csmaDevices.Get (0), true);
179
180       pointToPoint.EnablePcapAll ("third-distributed-csma");
181       phy.EnablePcap ("third-distributed-csma", apDevices.Get (0));
182       csma.EnablePcap ("third-distributed-csma", csmaDevices.Get (0), true);
183     }
184
185   Simulator::Run ();
186   Simulator::Destroy ();
187   return 0;
188 }
```

Green callout box:
1. Enable PCAP tracing on local nodes?
2. Close MPI cleanly

Right code panel:

```
223   {
224     UdpEchoClientHelper echoClient (csmaInterfaces.GetAddress (nCsma), 9);
225     echoClient.SetAttribute ("MaxPackets", UintegerValue (1));
226     echoClient.SetAttribute ("Interval", TimeValue (Seconds (1.0)));
227     echoClient.SetAttribute ("PacketSize", UintegerValue (1024));
228
229     ApplicationContainer clientApps =
230       echoClient.Install (wifiStaNodes.Get (nWifi - 1));
231     clientApps.Start (Seconds (2.0));
232     clientApps.Stop (Seconds (10.0));
233   }
234
235   Ipv4GlobalRoutingHelper::PopulateRoutingTables ();
236
237   Simulator::Stop (Seconds (10.0));
238
239   if (tracing == true)
240     {
241       // Depending on the system Id (rank), the pcap information
242       // traced will be different.  For example, the ethernet pcap
243       // will be empty for rank0, since these nodes are placed on
244       // on rank 1.  All ethernet traffic will take place on rank 1.
245       // Similar differences are seen in the p2p and wireless pcaps.
246       if (systemId == systemWifi)
247         {
248           pointToPoint.EnablePcapAll ("third-distributed-wifi");
249           phy.EnablePcap ("third-distributed-wifi", apDevices.Get (0));
250           csma.EnablePcap ("third-distributed-wifi", csmaDevices.Get (0), true);
251         }
252       else // systemCsma
253         {
254           pointToPoint.EnablePcapAll ("third-distributed-csma");
255           phy.EnablePcap ("third-distributed-csma", apDevices.Get (0));
256           csma.EnablePcap ("third-distributed-csma", csmaDevices.Get (0), true);
257         }
258     }
259
260   Simulator::Run ();
261   Simulator::Destroy ();
262
263 #ifdef NS3_MPI
264   // Exit the MPI execution environment
265   MpiInterface::Disable ();
266 #endif
267
268   return 0;
269 }
```

**1** **2**

# Script Output–Identical

```
$ ./waf –run third
```
```
Waf: Entering directory `build/debug'
Waf: Leaving directory `build/debug'
'build' finished successfully (2.152s)
At time 2s client sent 1024 bytes to 10.1.2.4 port 9
At time 2.01796s server received 1024 bytes from 10.1.3.3 port 49153
At time 2.01796s server sent 1024 bytes to 10.1.3.3 port 49153
At time 2.03364s client received 1024 bytes from 10.1.2.4 port 9
```

```
$ ./waf --run third-distributed \
--command-template="mpirun -n 2 %s --tracing"
```
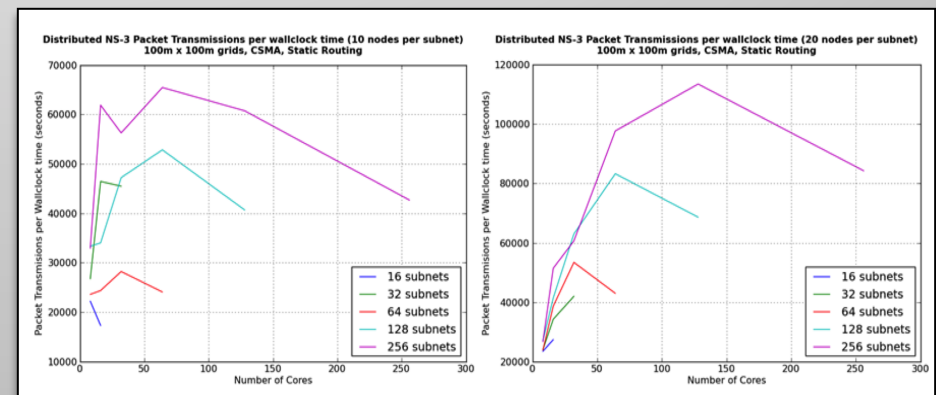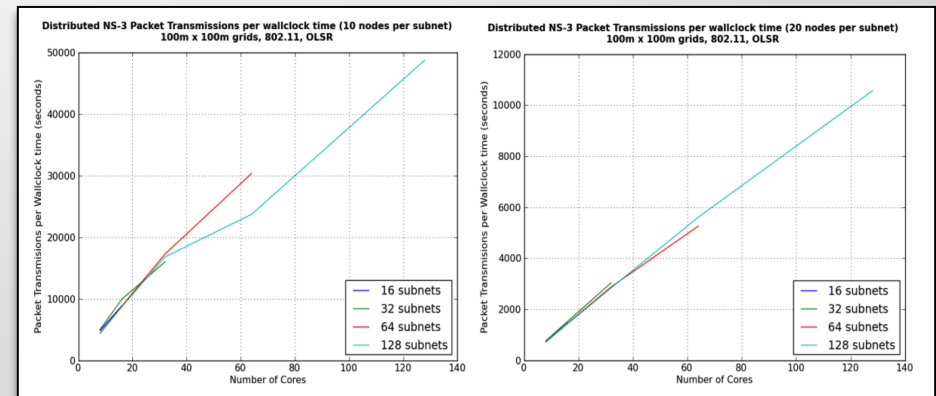```
Waf: Entering directory `build/debug'
Waf: Leaving directory `build/debug'
'build' finished successfully (2.050s)
At time 2s client sent 1024 bytes to 10.1.2.4 port 9
At time 2.01796s server received 1024 bytes from 10.1.3.3 port 49153
At time 2.01796s server sent 1024 bytes to 10.1.3.3 port 49153
At time 2.03364s client received 1024 bytes from 10.1.2.4 port 9
```

**ılıllNS-3**

# Cryptic Error Conditions

- `Can't use distributed simulator without MPI compiled in`

  - Not finding or building with MPI libraries

  ➢Reconfigure NS-3 and rebuild


- `assert failed. cond="pNode && pMpiRec", file=../ src/mpi/model/mpi-interface.cc, line=413`
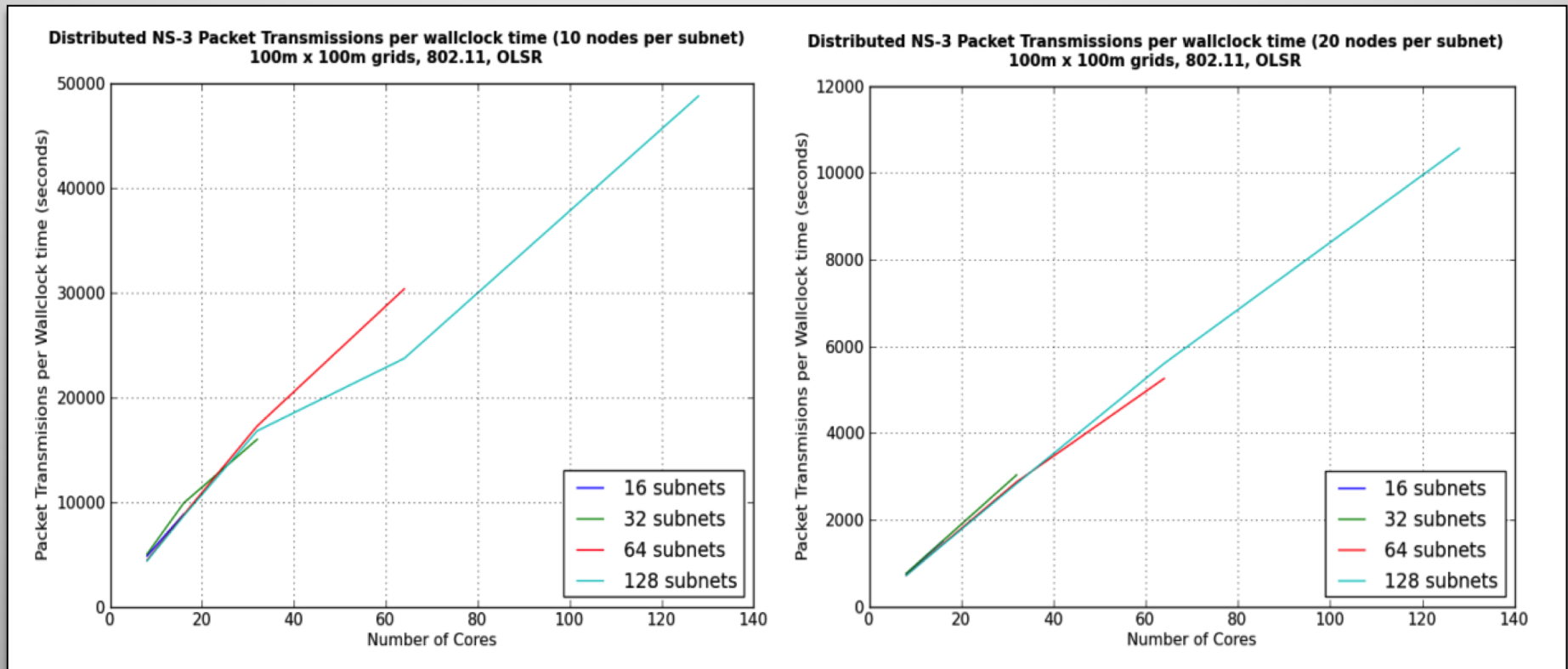
  - Mis-aligned node or interface IDs

# Performance Optimizations

- Larger Lookahead

- Synchronization cost grows exponentially with LP count
  - More work per LP is better
  - Speed gains up to $10^{2\text{-}3}$ ranks, depending on model



- Appropriate performance metric
  - Events/sec can be misleading with varying event cost
  - Packet transmissions (or receives) per wall-clock time

# Parallel Performance with Large Computation Load:  802.11+OLSR

- Linear scaling out to 128 ranks



Distributed NS-3 Packet Transmissions per wallclock time (10 nodes per subnet)
100m x 100m grids, 802.11, OLSR

Distributed NS-3 Packet Transmissions per wallclock time (20 nodes per subnet)
100m x 100m grids, 802.11, OLSR

# Parallel Performance with Small Computation Load:  CSMA+Static

- Performance drops at modest number of ranks

# Work in Progress

- **Automatic memory scaling**
  - Automatic ghost nodes, globally unique node IDs
  - (See my talk tomorrow ☺)

- **Automatic partitioning, ghost alignment**

- **Distributed Real Time**
  - Versus simultaneous real-time emulations:
    — LP-to-LP messaging gives greater *Lookahead* than independent ns-3 instances connected by emulated network devices

- **Scalable default routing**
  - AS-like routing between LPs
  - Scalable replacement for GOD or Nix-vector routing with ghost nodes

# (Mostly) Parallel Partitioning Tools

Serial Code

Parallel Code

model/

**Many** of:
*.xndl
*.xndl.gz
*.xndl.gzip

128 files
123 GB

$ partition -f

model.filelist

8K

$ sector -F model.filelist

model.sectors.xndl.gz

2.7 GB

model.sector.metis

14 GB

$ gpmetis model.sector.metis <nparts>

model.metis.part.<nparts>

1.4 GB

partition-model <first-rank> <last-rank>

**One** per rank of:
*.xndl
*.xndl.gz
*.xndl.gzip

128 files
123 GB

XmlSimulator.cc