# Tools for Enabling Automatic Validation of Large-scale Parallel Application Simulations
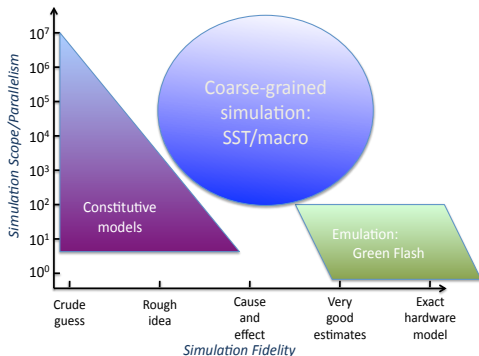
Deli Zhang    Gilbert Hendry    Damian Dechev

University of Central Florida

Sandia National Laboratories

October 2, 2014
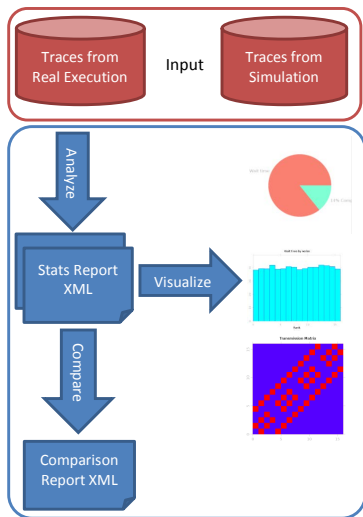
# Hardware/Software Co-design



- Exascale co-design
- Simulation is key

# Hardware Model Validation

- Hardware model is a set of simulation input parameters, e.g., network topology, network bandwidth/latency, node frequency, etc.
- The goal of of simulation validation is to establish the accuracy by quantizing the error between the simulated execution and the execution on the physical machine
- The error of the simulation can be used to guide future tuning process

# Validation Work Flow



- Gather execution traces
- Distill statistical data
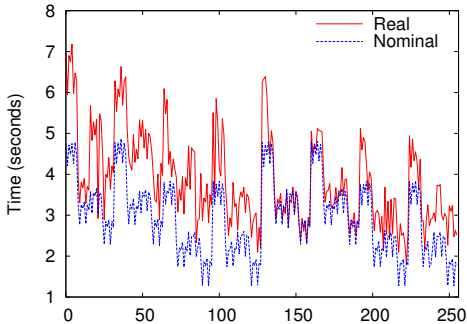- Compute errors through comparison

# Existing Metrics

- Coarse-grained metrics, such as total execution time, lacks fidelity to identify fine-grained execution differences
    - Insensitivity to some parameters
    - Some parameters have adversarial effects
- Detailed traces are not ready for quantitative comparison
    - TAU, Scalasca, Vampir, IPM, mpiP, etc.
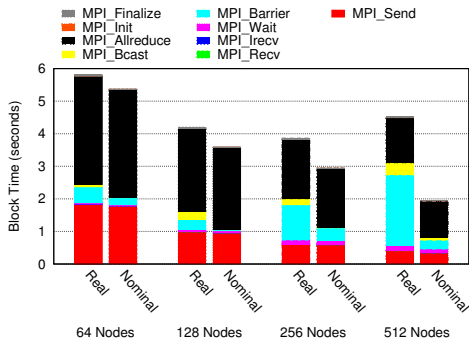
# Proposed Metrics

- Benefits
    - Fine-grained metrics improves validation fidelity
    - Matrix format facilitates quantitative comparison
- Experiment Environment
    - *Hopper* at SNL (a Cray XE6 cluster)
    - Gemini interconnect with two communication paths: fast memory access (FMA) and block transfer engine (BTE)
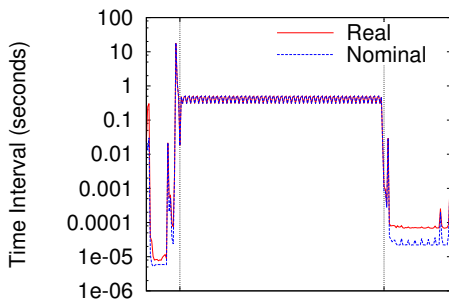    - miniMD and coMD as benchmark

# By-node



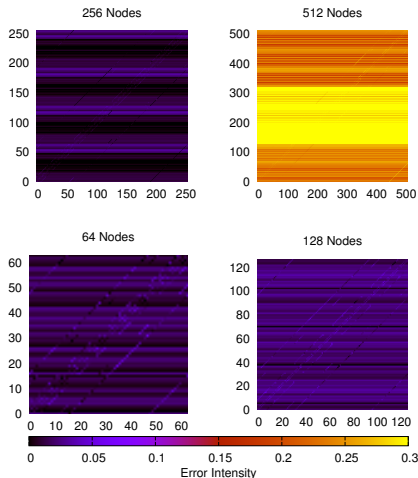- Break down by rank
- $2 \times N$ matrix

# MPI Histogram



- Break down by MPI functions
- $F \times N$ matrix

# Collective Synchronization



- Collective functions as synchronization barriers
- Break down by collective phases
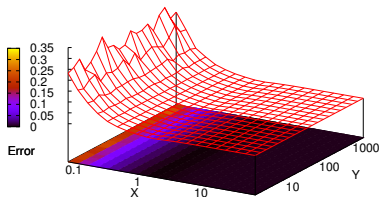- $S \times N$ matrix

# Node-to-node Communication
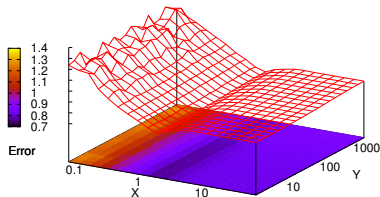


- Pair-wise timing
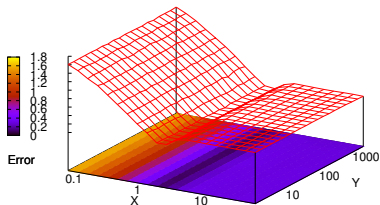- $N \times N$ matrix

# Link Bandwidth and Latency



(a) Total execution time

(b) MPI histogram

The error measured by MPI histogram converges at 2.4Ghz, which is the nominal value

# Link Bandwidth and Latency



(c) Node-to-node timing

(d) Collective Synchronization

The error measured by node-to-node timing and collective synchronization converges at 2.4Ghz.

# Auto-tuning Work Flow



- Search the parameter space for the optimal values