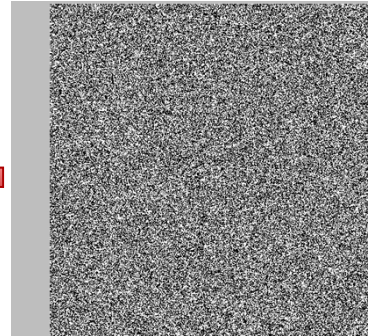


Exceptional service in the national interest



Continuous
Monitoring of HPC
platforms and the
applications run
upon them



Chaotic, random
applications,
workflows, and
dependencies

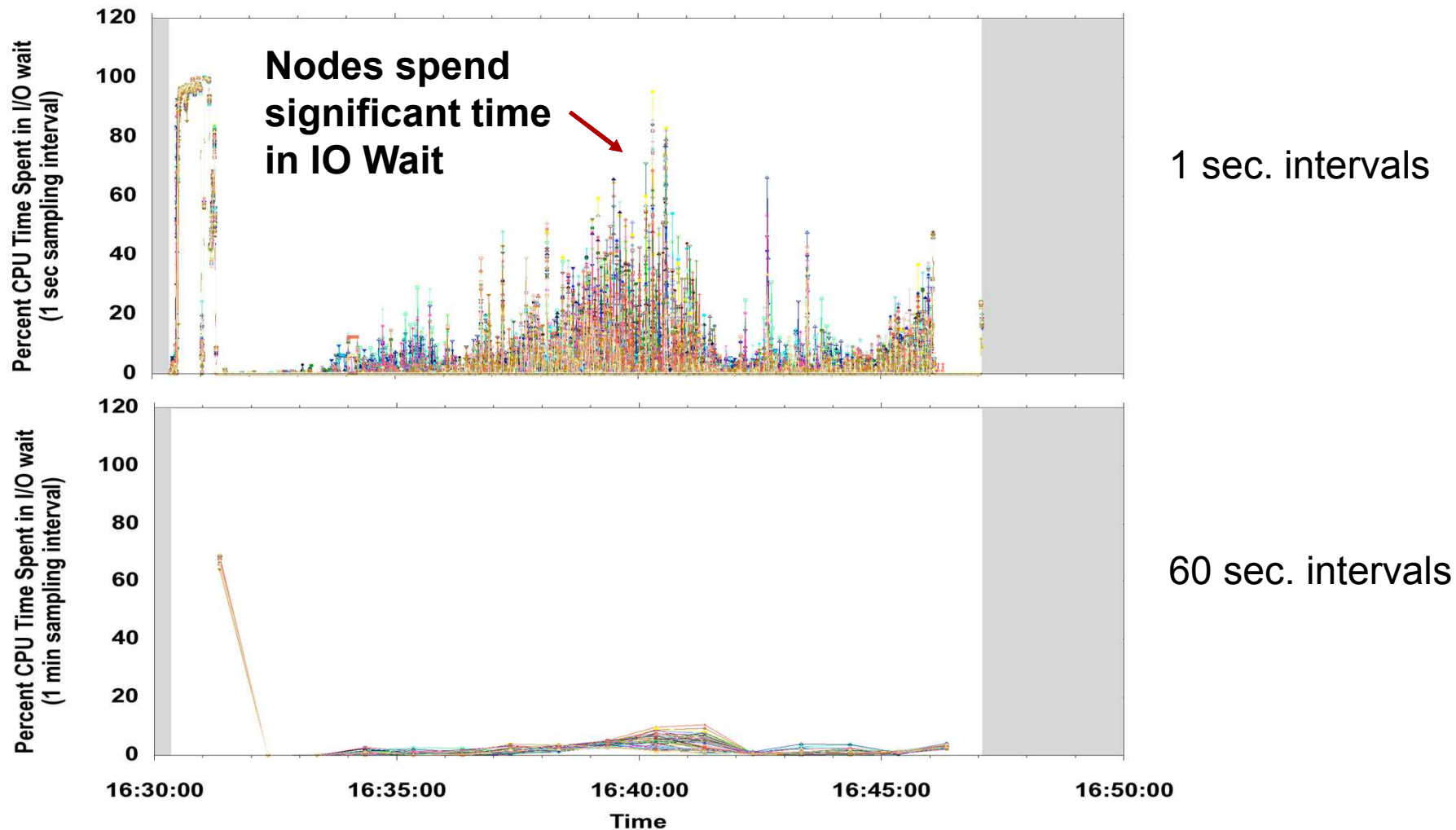
Defining Metrics to Distill Large-Scale HPC Platform and Application Performance Data into Actionable Quantities – Resource Contention of File System and Aries Interconnect

(alphabetical) Anthony Agelastos, James Brandt, Ann Gentile, Justin Lamb, Kevin Ruggirello, Joel Stevenson

Executive Summary

- **Objective:** Formally define the needed metrics/methods that distill vast quantities of HPC monitoring data to a minimum set of actionable and interpretable quantities to characterize **resource contention** that can be used by **application developers, system administrators, production analysts, and HPC platform designers**
- **Methodologies Used**
 - 1 whole-system, high-fidelity monitoring service (e.g., LDMS)
 - 1 mini-app with known network interconnect profiles (e.g., IMB)
 - 1 mini-app with known file system behavior (e.g., IOR)
 - 1 production-relevant test case (e.g., CTH)
- **Results**
 - **Actionable Metric Model (AMM)** produces good estimations of time spent in contention and overall impact of contention for single-node cases

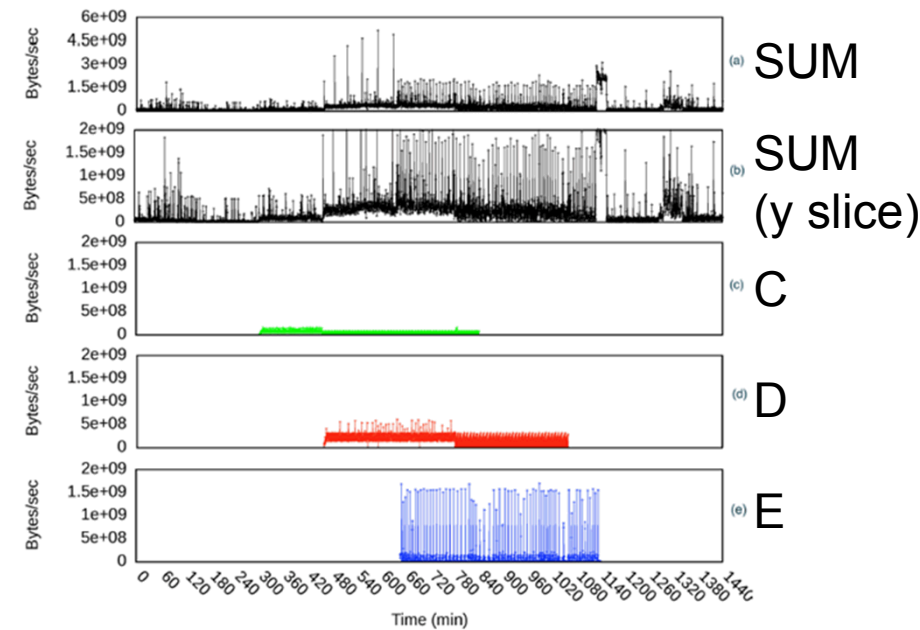
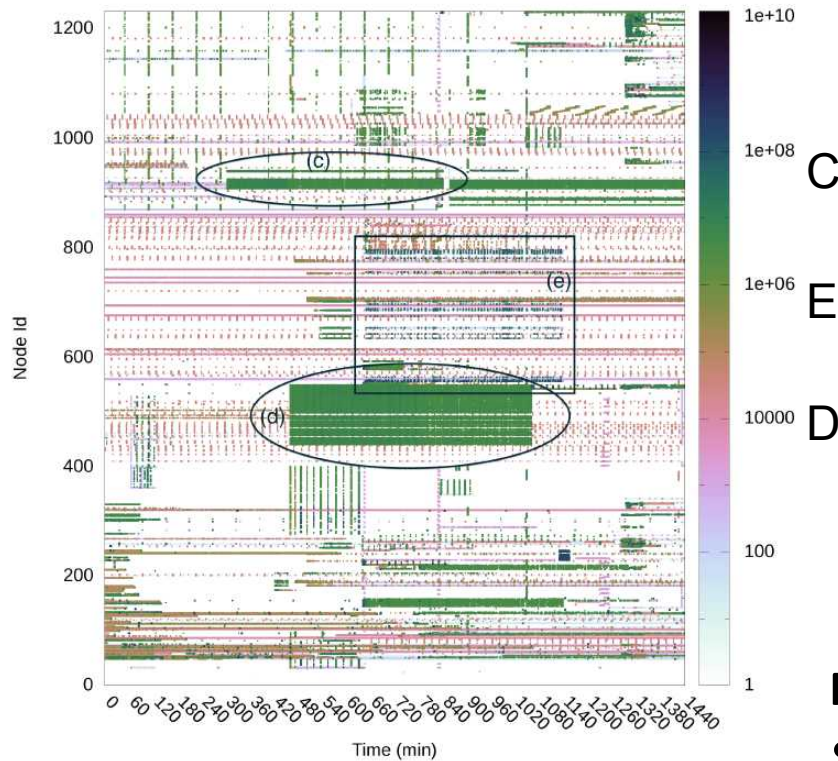
High-frequency Sampling



High-frequency sampling intervals are necessary

Whole-system Views

scratch1: Bytes/sec written over 20 sec interval



Provides insight about:

- system-wide utilization;
- events correlated in time and space;
- identify contention for shared resources;
- understand varying production conditions that can explain performance variations

Whole-system view enables environmental insight

AMM: 1-Node File System

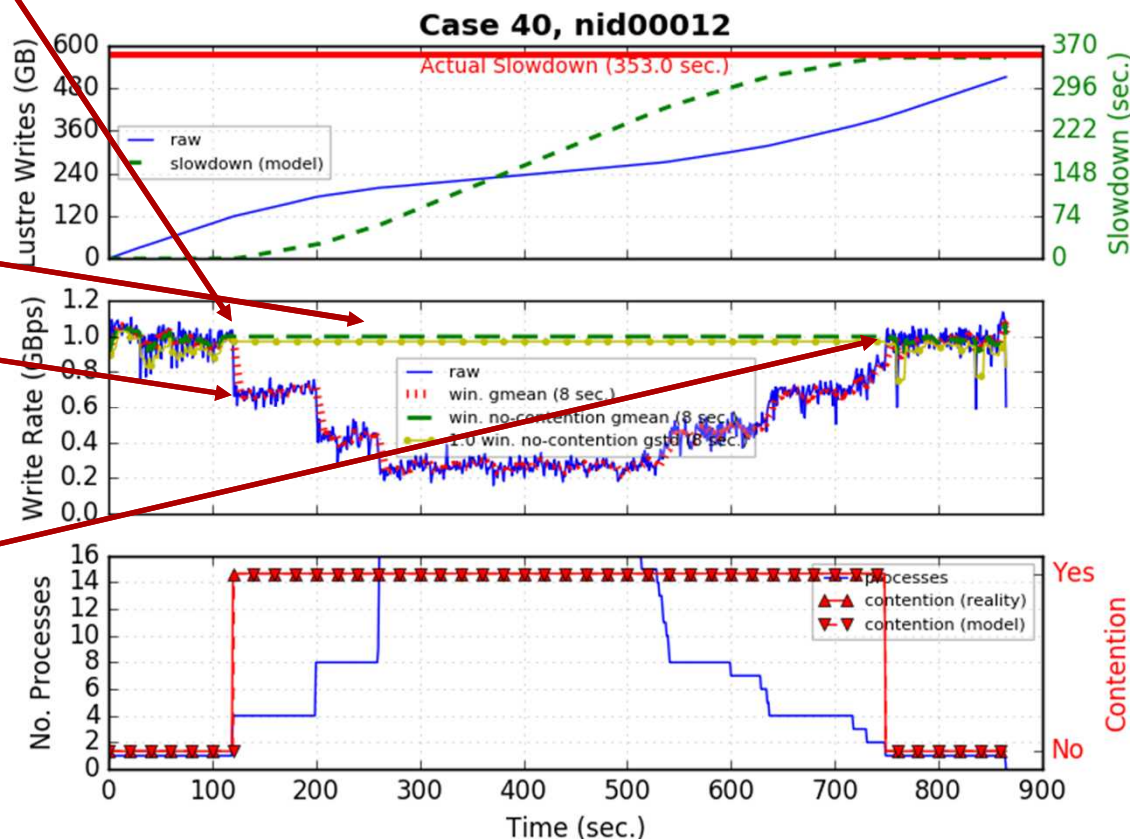
Description: Multiple IORs (1->4->8->16->8->4->1) running on Mutrino with 1 IOR process per node. The long-running, first IOR process is present on NID 00012. With only 4 OSTs on Mutrino's Lustre FS (and with a striping of 4), any more than 1 IOR running in this configuration *will* create contention.

1. Contention occurs

1. This went from 1 IOR to 4

2. Track non-contentious & contentious data

3. Contention ceases when observed trends return to non-contentious values



- **Top sub-figure:** raw Lustre writes in GB (blue, solid line); estimated slowdown model in sec. (green, dashed line)
- **Middle sub-figure:** Lustre write rate in GB/sec (blue, solid line); windowed geometric mean of this rate (red, thin-dashed line); non-contentious portion of this rate (green, dashed line); non-contentious geometric standard deviation (yellow line with dots)
- **Bottom sub-figure:** # of IOR processes running (blue, solid line); modeled (red, dashed line with bottom triangle) and actual (red, solid line with top triangle) contention states

Ongoing Work

1. We are developing improved methods for AMM with Statistical Sciences dept. for multi-dimensional rollups of file-system and interconnect metrics.
2. We are working with Cray to better understand Aries' contention metrics and how to associate them with specific running processes given adaptive routing complexity.
 1. AMM was shown to work on interconnect metrics to estimate time in contention, however other metrics are needed to produce a more accurate slowdown model.

Questions?

- Answers other than 42.