# Combining System Characterization and Novel Execution Models to Achieve Scalable Robust Computing

**H. Adalsteinsson, J. Brandt, B. Debusschere, A. Gentile, J. Mayo, P. Pebay, D. Thompson, M. Wong**

**Sandia National Laboratories**

**http://ovis.ca.sandia.gov**

Sandia National Laboratories

# Motivation

- Platform growth in size and complexity
  - ~ 100 cores per CPU
  - Heterogeneous CPUs
  - Multiple CPUs per node
    - Hierarchical access to memory, communications, etc.
  - Thousands of nodes
  - Millions of cores

- Expected system wide decrease in mean time to component failure (~constant per component)

- Many current MPI applications are not robust to component failure and don't take platform non-uniformity into account
  - Hangs if one process fails to reach a barrier (can happen if underlying component failure occurs)
  - Decompositions don't take resources/hierarchies into account → imbalance

# Big Goal

- Infinitely scalable and 100% efficient applications and platforms

# Moderate but attainable Goals

- Understand current and future architectures
  - Failure modes and root causes
  - Associated failure and pre-failure symptoms
    - Detect
    - Predict
  - Monitor and communicate system state and state change (Granularity?)
    - System/Resource manager
    - Application
  - Implications of hierarchical and heterogeneous resources on application performance
    - Propagation of faults

# Moderate but attainable Goals Cont.

- Monitoring and analysis infrastructure that can work hand in hand with the system and application to facilitate high performance and resilience to faults

- Novel execution models that are able to take advantage of new and future platform architectures and leverage advanced monitoring and analysis capabilities

# Current interrelated Projects aligned with goals

- OVIS: A tool for scalable real time monitoring and analysis

- Quantification of ability to predict failures

- Novel execution models to enable scalable, reliable, and self-balancing applications
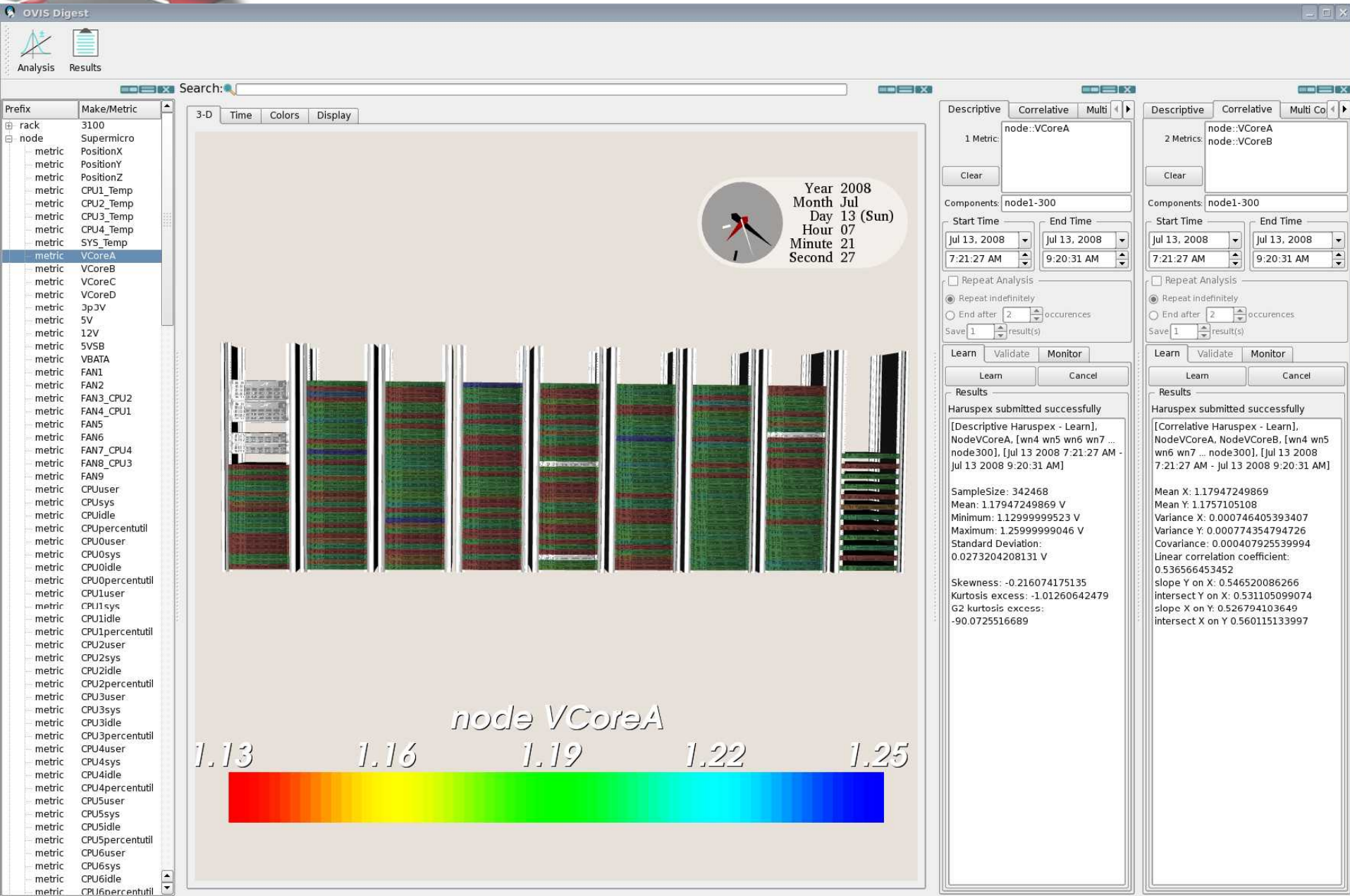
# OVIS Goals

- Scalable exploration of measurable attributes via both visual and quantitative analysis
- Run time anomaly detection
  - Hypothesize that statistical outliers can be an indication of impending failure
- Component health evaluation
  - Relative
  - Degradation
- Inference based root cause analysis
- Probability based failure prediction
- Interaction with
  - Resource management
  - Applications
  - System administrators

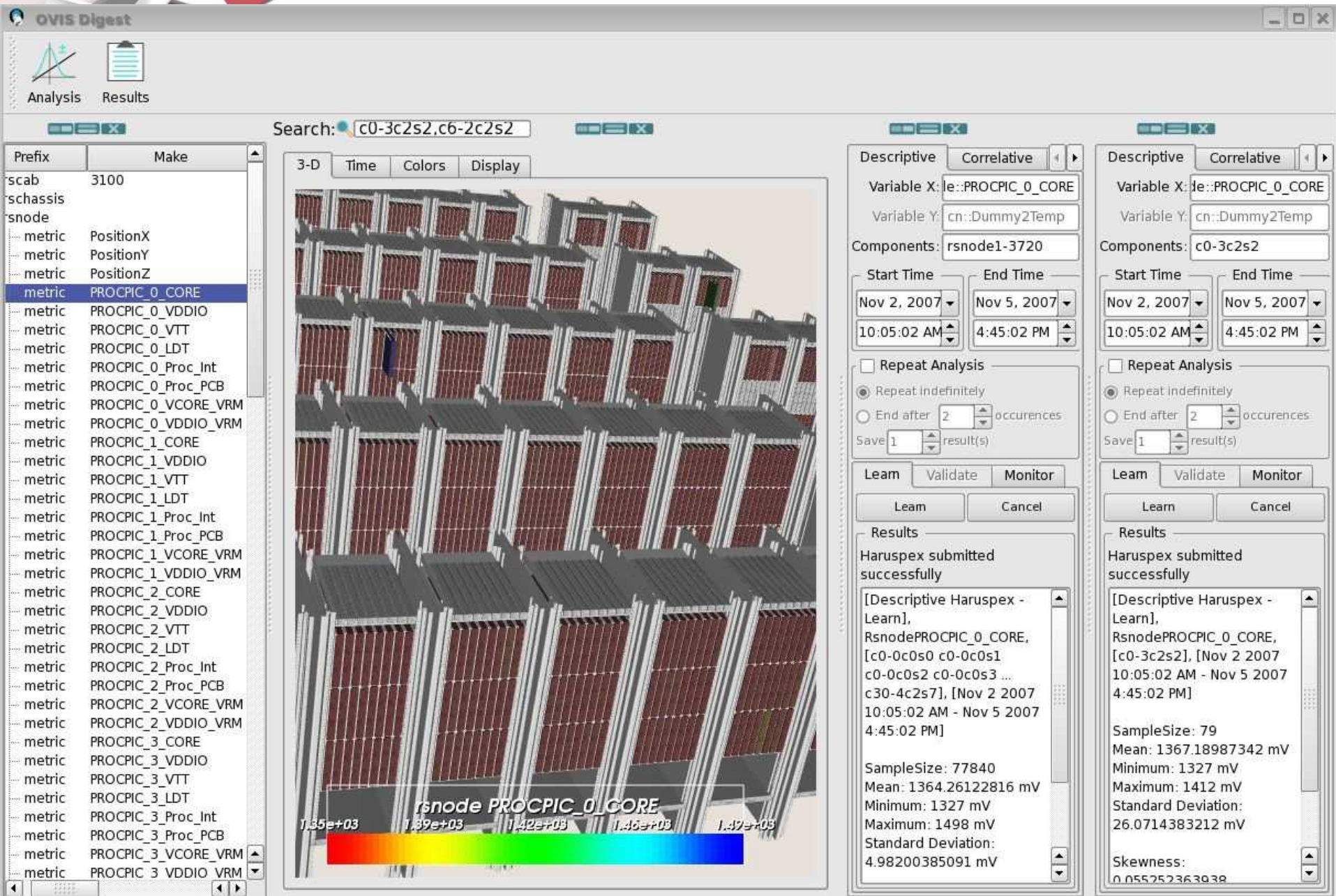Sandia National Laboratories

# Features of current release (2.0)

- Visual analysis
  - Realistic rendering of physical system
  - Components colored by relative magnitude of measurements
  - Play live feeds or historic data to see temporal variation
- Analysis engines
  - Descriptive statistics
  - Multi-variate correlative statistics
  - Bayesian inference
- Learn and Monitor modes of operation

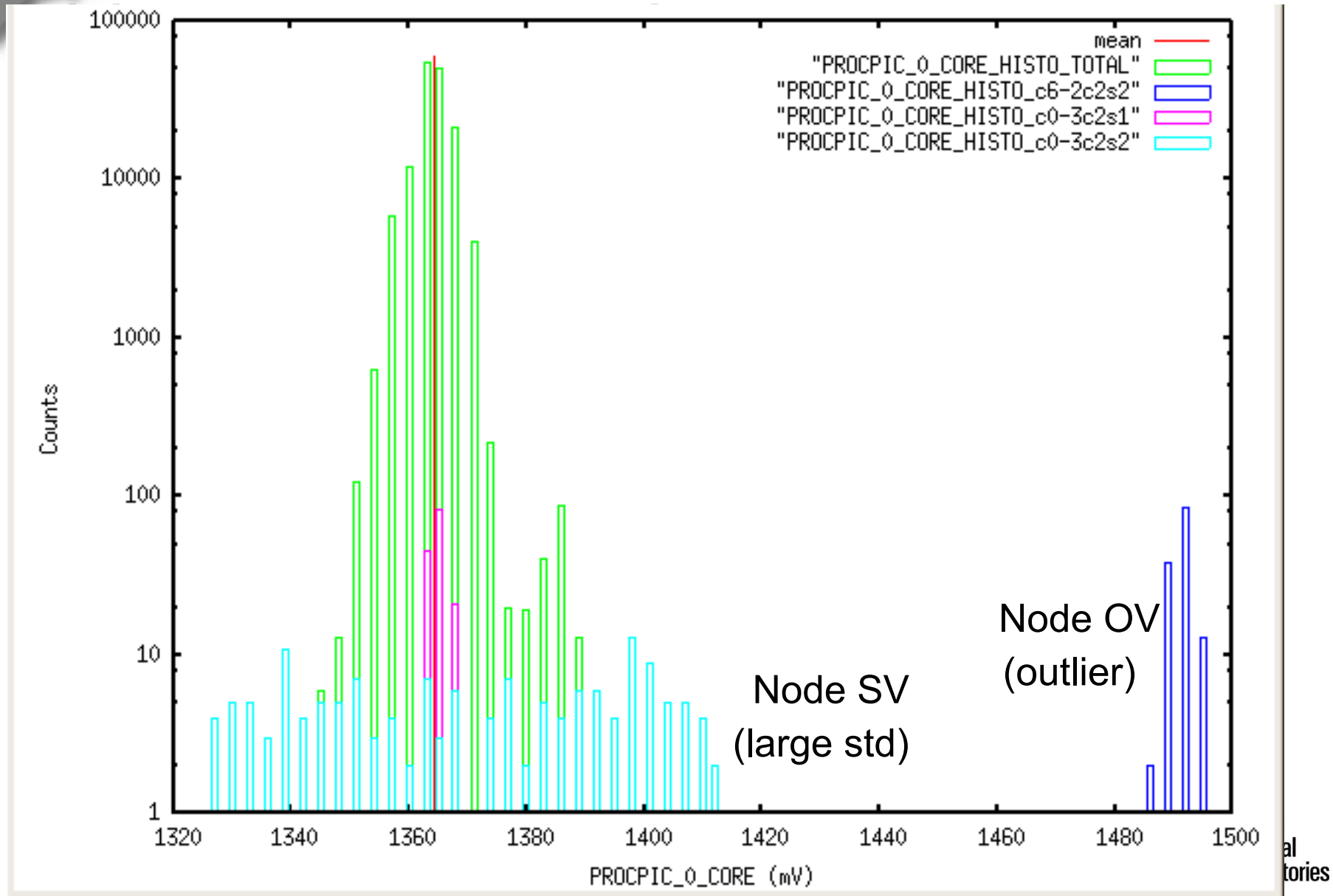Sandia National Laboratories

# OVIS (TLCC)

# OVIS (Red Storm)

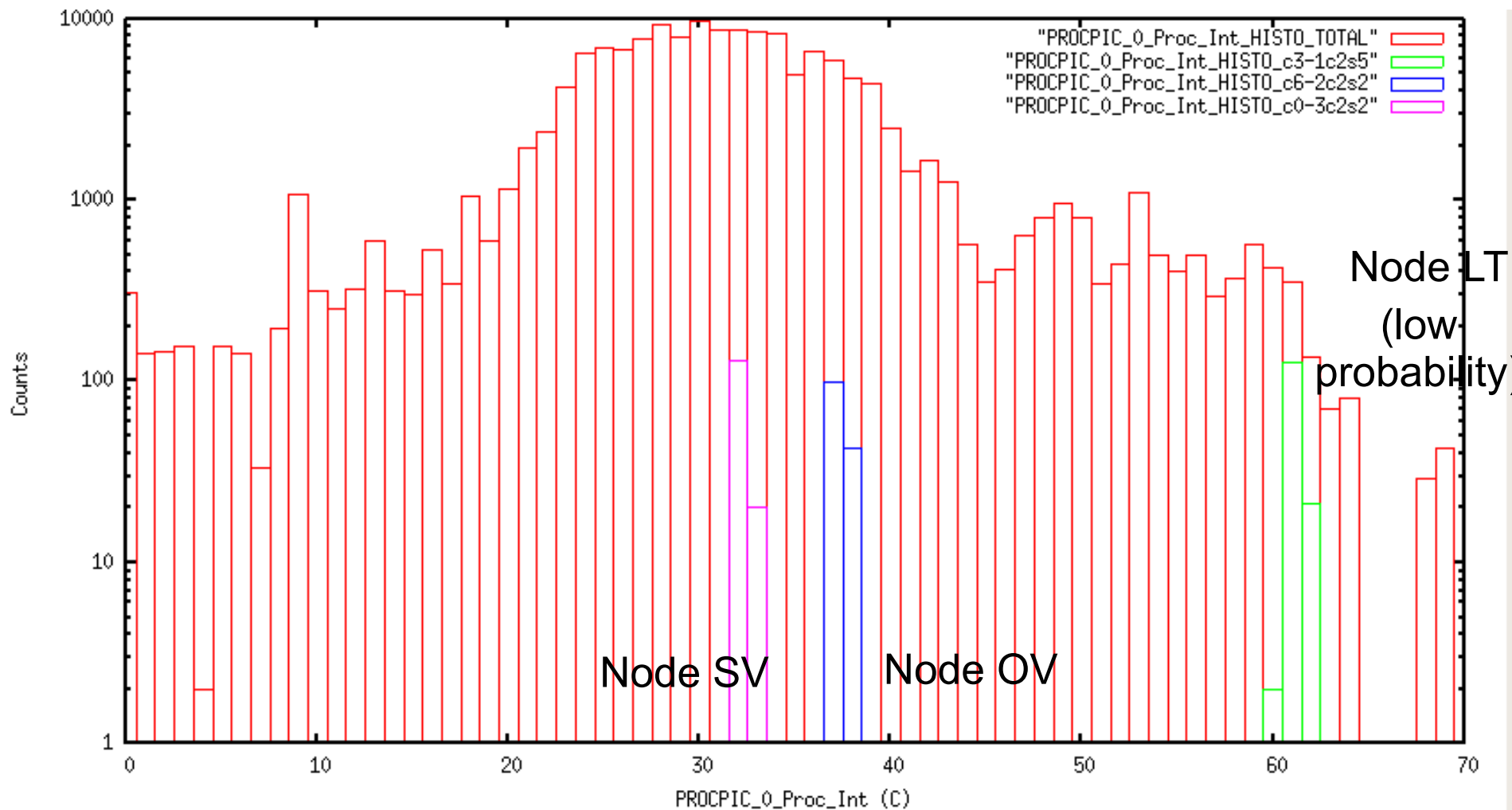# Examples of Interesting Component Related Features

# Power supply voltage

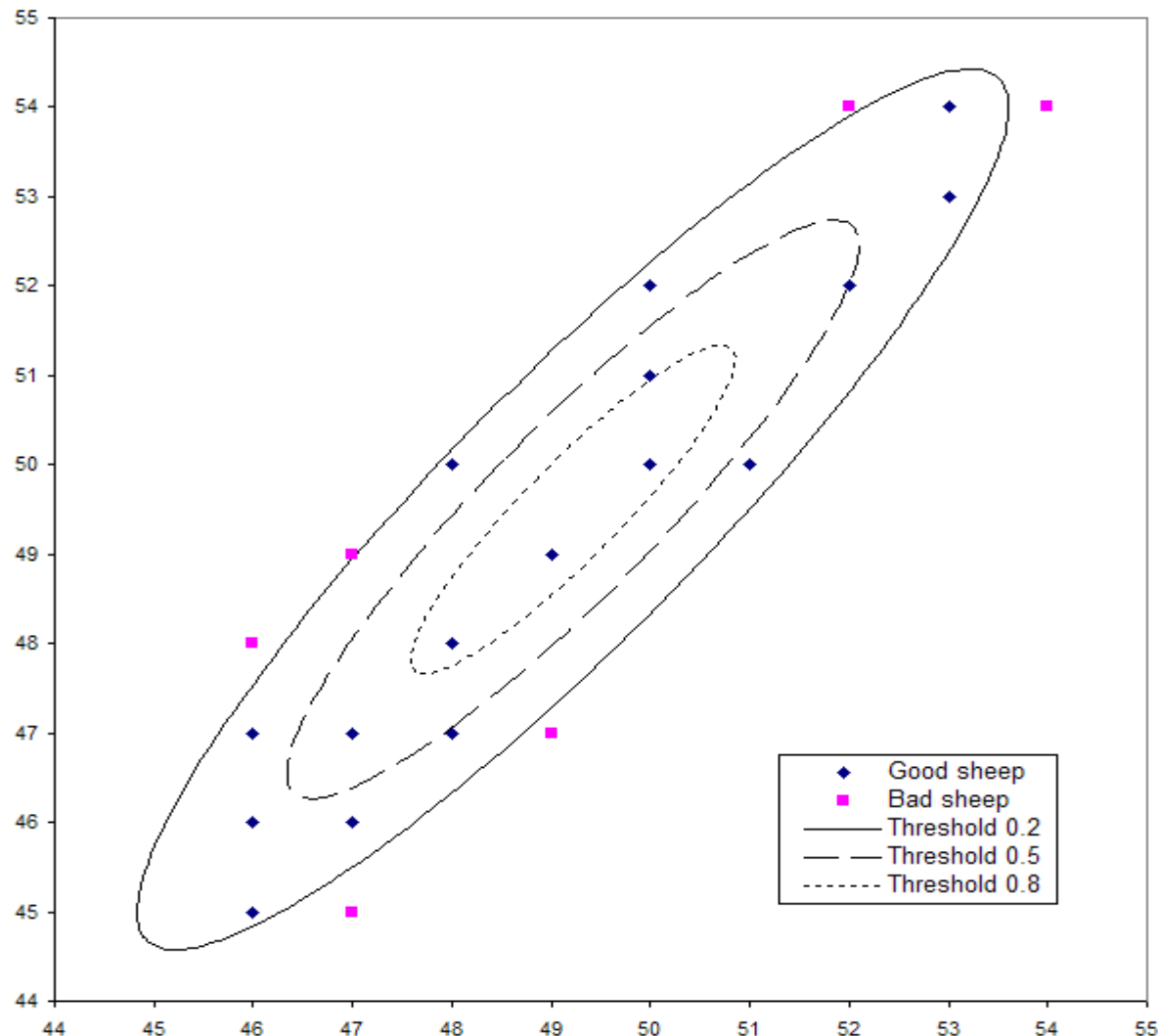# Temperature Histogram (log scale)

# Two Variable Temperature Plot with Contours of Relative Probability

# Challenges

- Long term high frequency (sample every ~seconds) collection of component related data
  - On Whitney ~one million data points per minute with a 5 second sampling interval → ~10GB/day
- Collection of detailed failure data
- Establishing causal relationships between failures and perhaps inter-related behaviors of 10s of variables
- Defining attributes not currently measured that if attainable would significantly contribute to failure prediction
- Quantifying ability to use component characteristic behaviors to predict failures
- Timely feedback to system and application

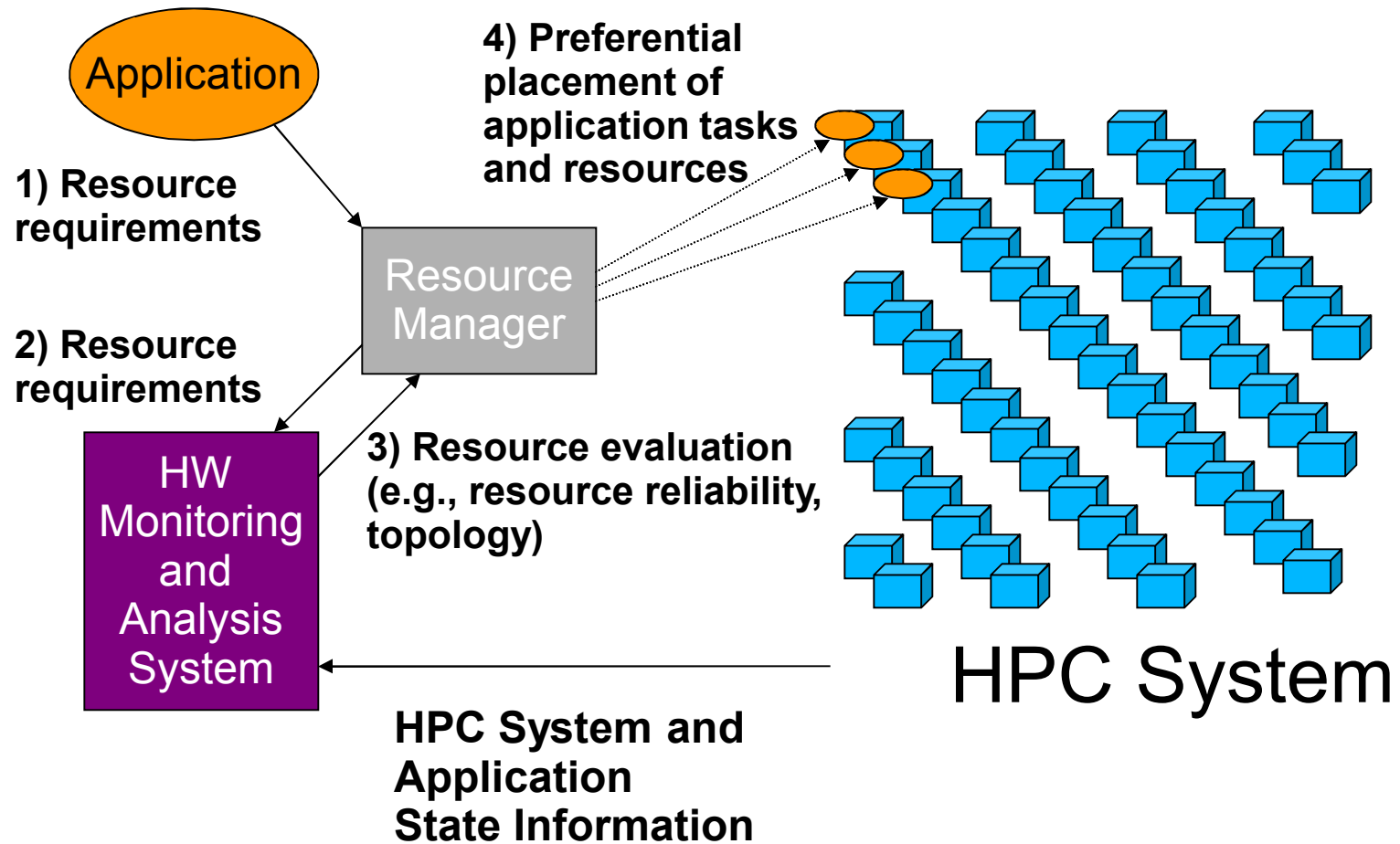Sandia National Laboratories

# Quantify Ability to Predict Failures

- Characterization of multivariate distributions

- Time series analysis

- Classifier for identifying events of interest

- Attribute selection

- Root cause analysis via Inference techniques

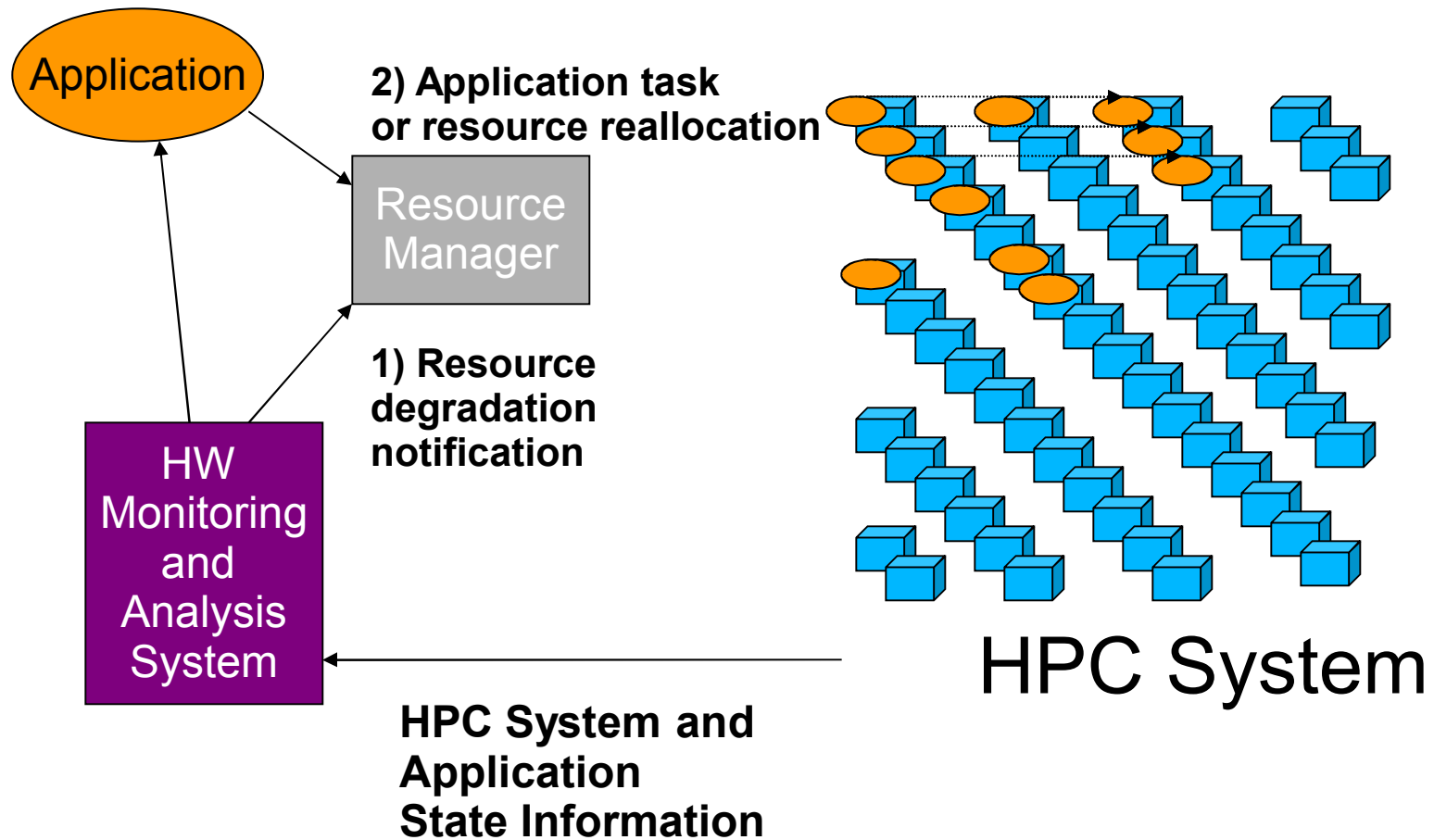- Effectiveness scoring algorithms and metrics

# Novel execution models to enable scalable, reliable, and self-balancing applications

Sandia National Laboratories

# Preferential Resource Allocation

**Application**

**1) Resource requirements**

**2) Resource requirements**

**4) Preferential placement of application tasks and resources**

**Resource Manager**

**HW Monitoring and Analysis System**

**3) Resource evaluation (e.g., resource reliability, topology)**

**HPC System and Application State Information**

**HPC System**

Sandia National Laboratories

# Resource Reallocation Upon Failure Prediction

Application

**2) Application task or resource reallocation**

Resource Manager

**1) Resource degradation notification**

HW Monitoring and Analysis System

**HPC System and Application State Information**

HPC System

# Novel Execution Models

**1) Subset of application tasks complete or require additional resources (Red)**

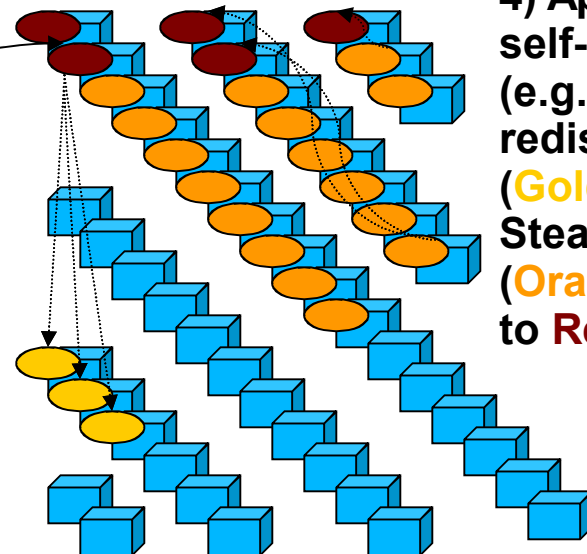**2) Application queries for remaining task (Orange) or system state**

**4) Application self-balances (e.g., redistribution (Gold) or Stealing (Orange tasks to Red)**

**Application**

**3) Analysis system responds with state information**

**HW Monitoring and Analysis System**

## HPC System

**HPC System and Application State Information**

Sandia National Laboratories

# Architecture



- Hybrid model for analysis
- external: memory and state heavy analyses - generate models, generate distributions
- local: model comparisons, small peer group, allow reasonable uncertainty based on limited world view

# Challenges

- Determining criteria for task-to-resource mapping
  - Heterogeneous and hierarchical environment
  - Characterization of relevant information

- Availability of resource state information

- Infrastructure for timely analysis + resource manager + application interaction

Sandia
National
Laboratories

# Future/Ongoing Work

- Quantification of ability to predict failures
  - Hardware numerical data collection
  - Failure modes, symptoms
  - Analysis
  - Missing information identification

- Infrastructure to support up to date resource state knowledge with ability to share that information with both system and application in a useable time frame
  - Small memory footprint agent based
    - Low latency interaction
    - Time vs. global accuracy tradeoff
  - Description of resource state/requirements
    - Total memory, memory bandwidth, compute intensity, etc.
  - Duration of relevancy of historical data

# Questions?

http://ovis.ca.sandia.gov

Demo in ASC booth at SC08

Sandia
National
Laboratories