# Data Fusion and Statistical Analysis: Piercing the Darkness of the Black Box

Jim Brandt, Frank Chen, Vince De Sapio, Ann Gentile, Jackson Mayo, Phillippe Pebay, Diana Roe, David Thompson, Matthew Wong

With support from:

Marcus Epperson and Jerry Smith
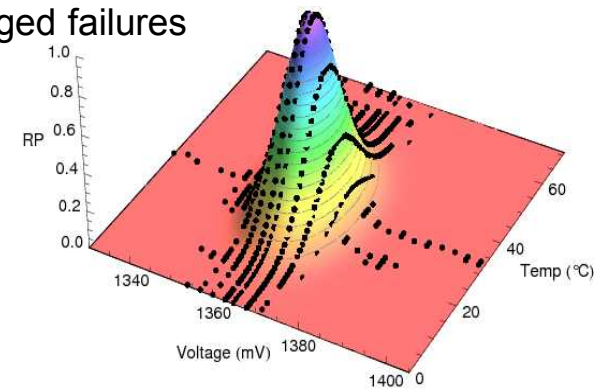
ASC™

# Black Box Defn.
www.thefreedictionary.com

**a.** A device or theoretical construct with known or specified performance characteristics but unknown or unspecified constituents and means of operation.

**b.** Something that is mysterious, especially as to function.

# Motivation

- Large amount of time spent/(wasted?) in checkpointing
- Black Box approach to failures in platforms implies scaling wall
  - Best we can do is adjust checkpointing freq based on the measured system average MTTI
- Assumption -- If we understood failure mechanisms
  - Instrument to get precursor warnings where possible
  - System and Application interfaces that would facilitate such recovery
  - Know where to invest more up front for more robust systems
    - Power supply redundancy cost/fault tolerance tradeoff
    - Memory socket cost/MTTF tradeoff

# Resilience Fault Prediction Strategy

- Discover predictors, accuracy, time windows, and coverage with respect to all non-recoverable faults
  - Scalable data collection
    - HW related metrics
      - Limited by current instrumentation
      - Discovery can help drive future system instrumentation
    - System related metrics
      - RM databases, log files, troubleshooting notes, etc.
      - Work with System Administrators to capture as much as possible
    - Not available
      - Human errors, power grid outages, etc.
  - Scalable data analysis
    - Definition of analysis methods that make sense given the data and time scales
      - Currently: correlate low probability behaviors with logged failures
    - Efficient data exploration tools
      - UI
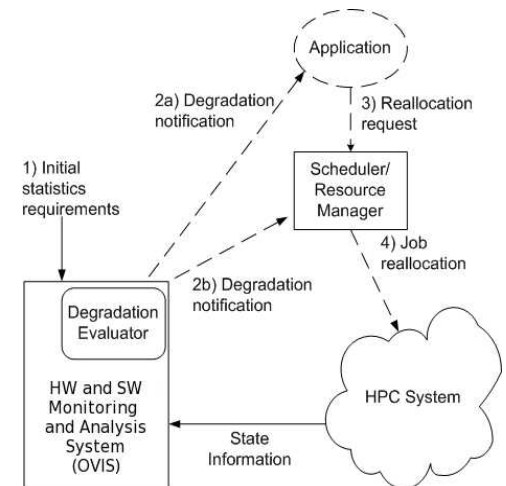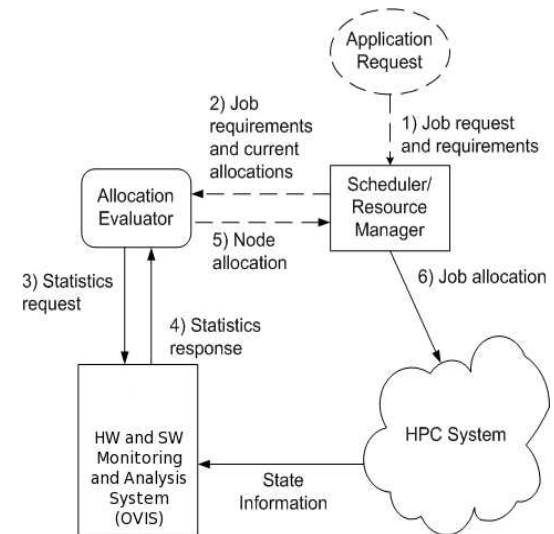      - Visualization
- Quantify Prediction Effectiveness

# Approaches

Seek "real time" functional component (Memory bank, core, communication bus, power supply sub-systems, etc.) level understanding of health

- ## Statistical analysis
  - Learn statistical characteristics of components that lead to both normal and abnormal operation of a system or platform

- ## Data Fusion
  - Look for correlations in disparate data to facilitate understanding of component, system and application level behavioral interaction, operation, and failure
    - Could be outlier or smack in the middle of normal operational regions -- TBD

# Ultimate Goal

- Health based allocation given parameters of application run

- Preemptive task/state migration based on failure prediction

- Targeted checkpoint frequency based on health (e.g. resource with higher chance of failure might save state frequently but this would be much faster than saving the aggregate. Have to address how much larger the memory footprint gets for healthy resources having to maintain state in case of rollback)
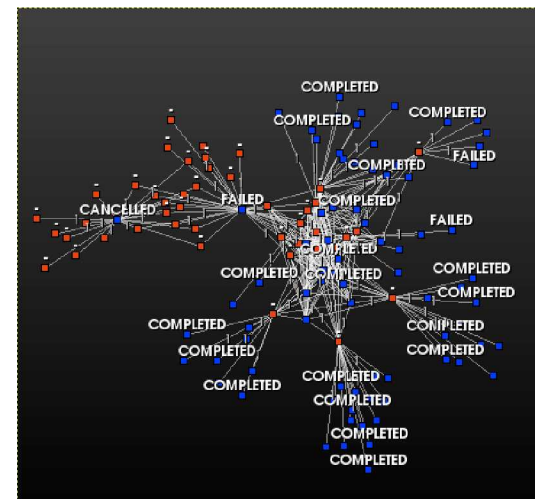
# How does this differ?

- Many system level studies (Black Box approach) give MTTI over the system that can lead to extrapolated MTTI for a pool size – Same studies show this approach to checkpoint frequency calculation doesn't scale

# Statistical Approaches

- Descriptive
  - Numeric metrics such as mem err rates, cpu utlization, voltage, PS ripple
- Correlation
  - Looking at numeric ensemble behaviors in run-time environment
- Time Series
  - Taking into account temporal displacements
  - Temporal ensemble behaviors relating to actual codes

- Graph based
  - Exploration of both numeric and text (log) data for telltale relationships using clustering techniques
- Incorporate both parametric and failure data to produce models that target the prediction of those failures
  - Looking for numeric metric correlations with failure that aren't tied to anomalous behavior
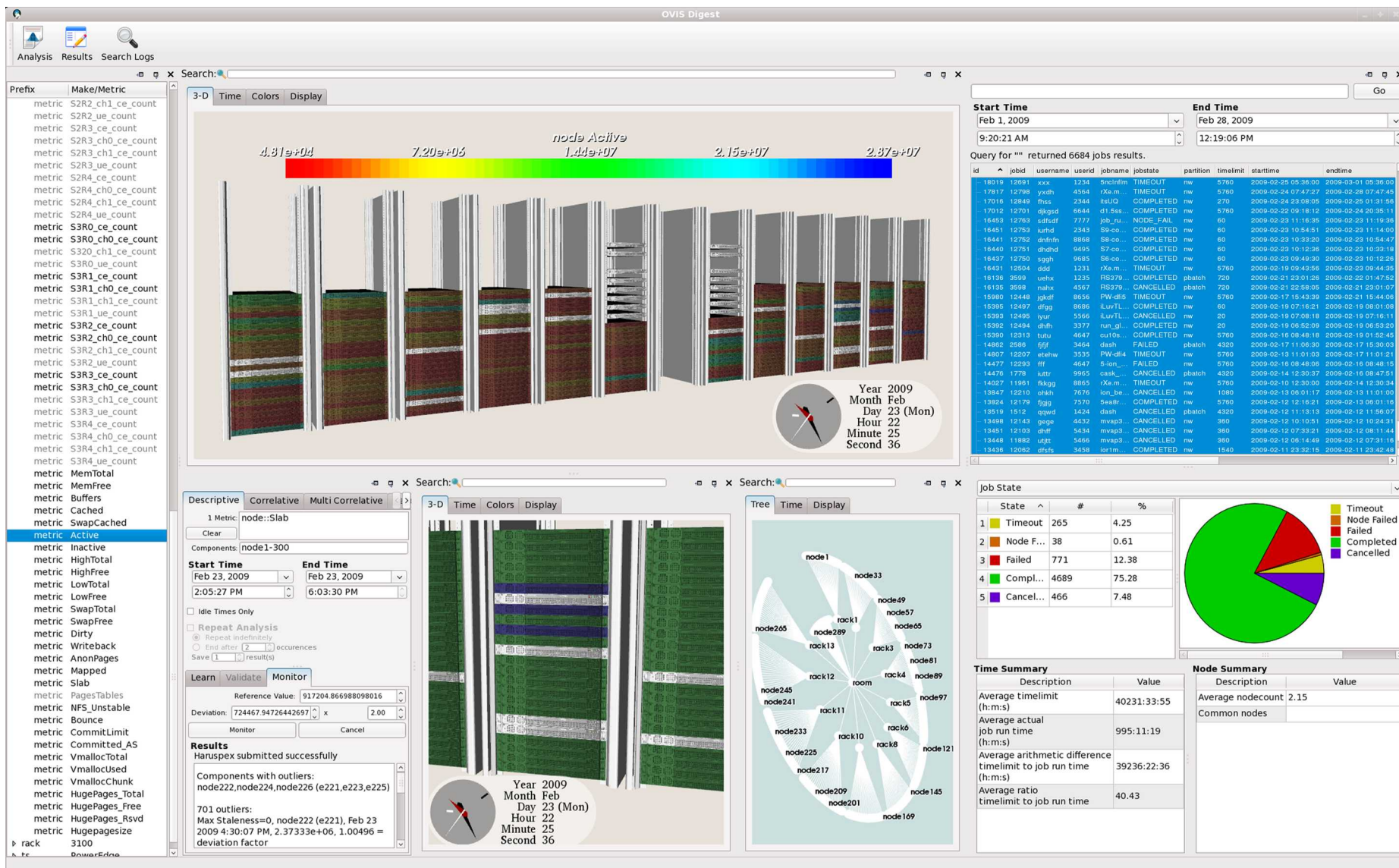  - Treats failure as a generic event without subtypes

# Data Fusion

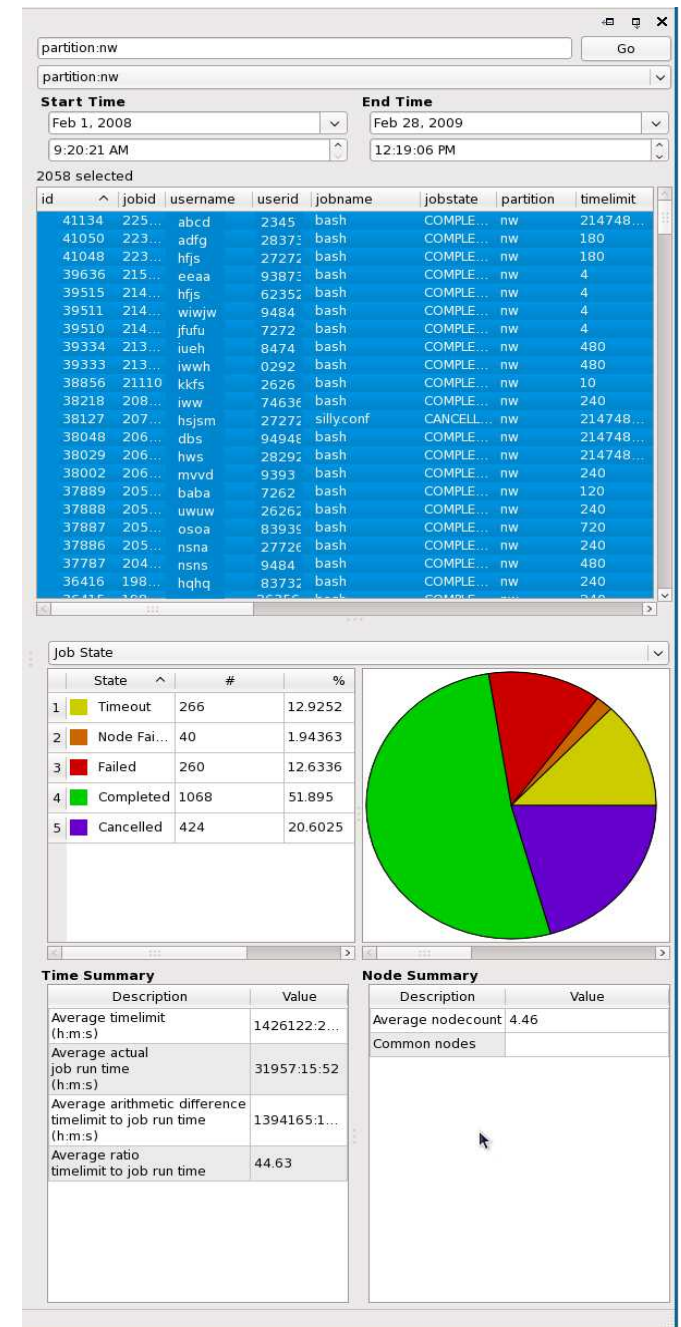Cockpit approach with interacting capabilities

- Perform analyses on aggregate of parametric data, Resource Manager (Slurm) data, Syslog, and Console logs

  - Relate user, application (problems here), resources used, application and resource based failures

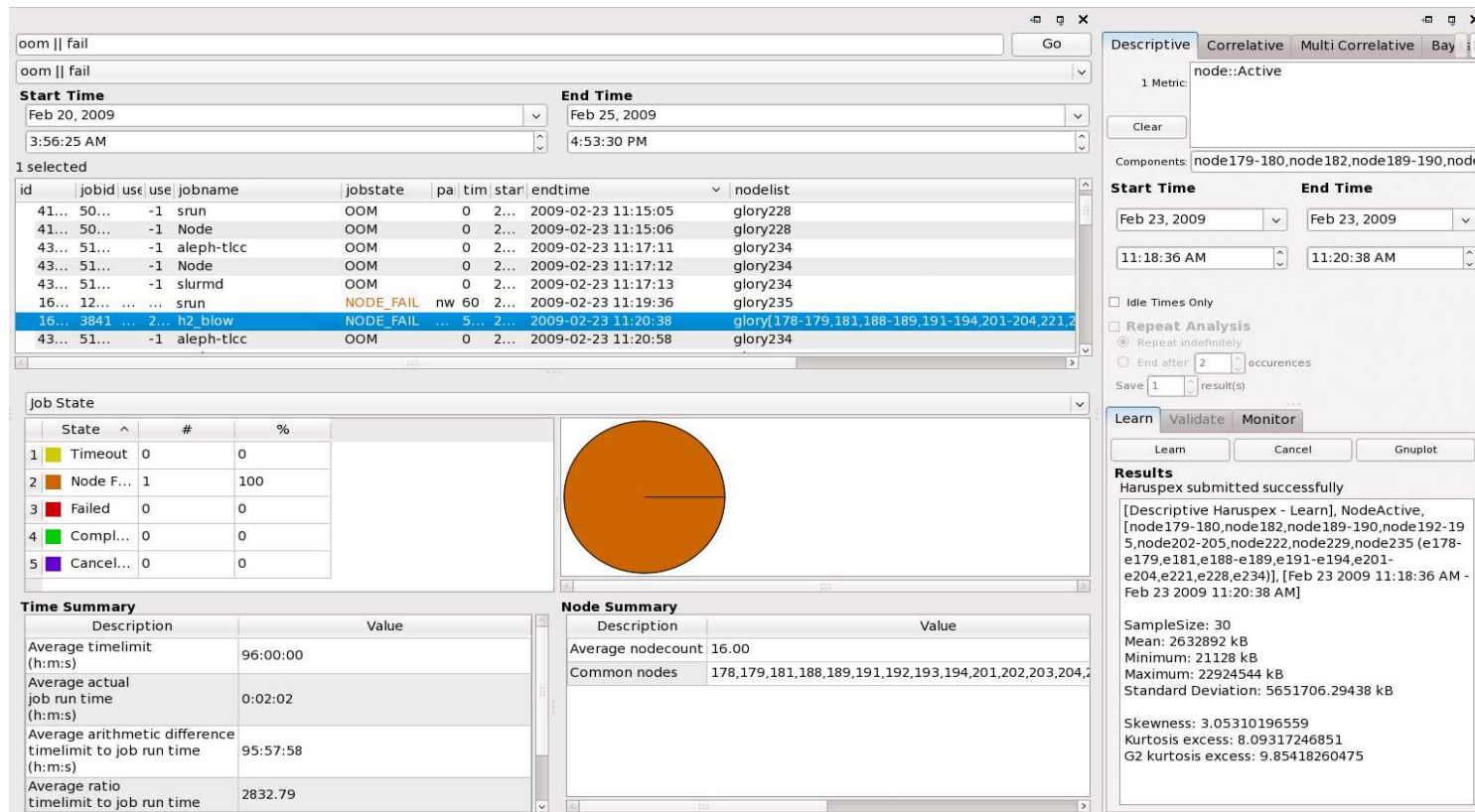# Data Fusion and Interacting Capabilities

# Scheduler and Log Search

- Integrated Scheduler and Error Log Searches

- Searchable:
  - e.g., OOM && node 234

- Aggregate statistics

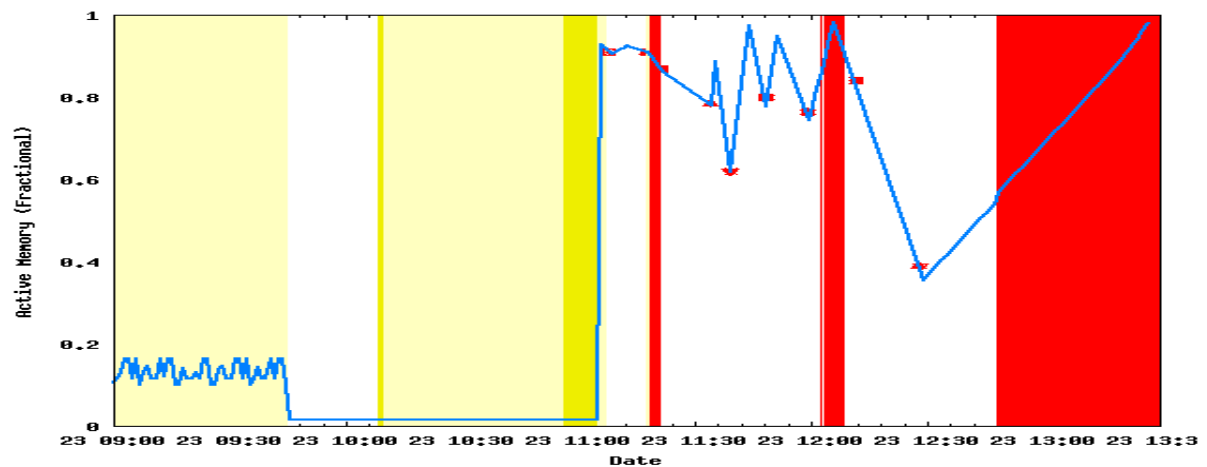- Pie subselection
  - e.g. Distribution of failed jobs by user

# Scheduler and Log data Drive Analysis



- Time and component relationship of OOM killer and failed job leads to invocation of memory utilization analysis
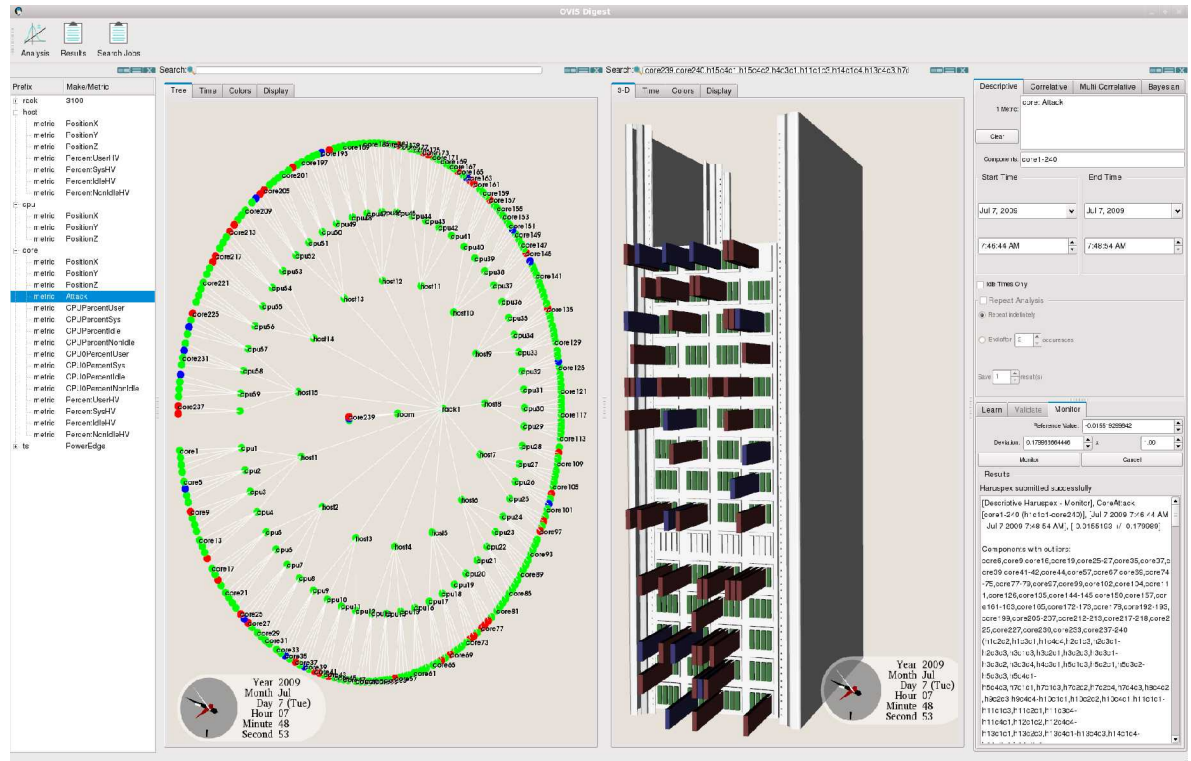
# Data Fusion Visualizations

- Logical and Physical Displays

- Colored by analysis results



- Textual Job data – state and duration (shaded)

- Textual log data (x'es)

- Numerical data (blue)
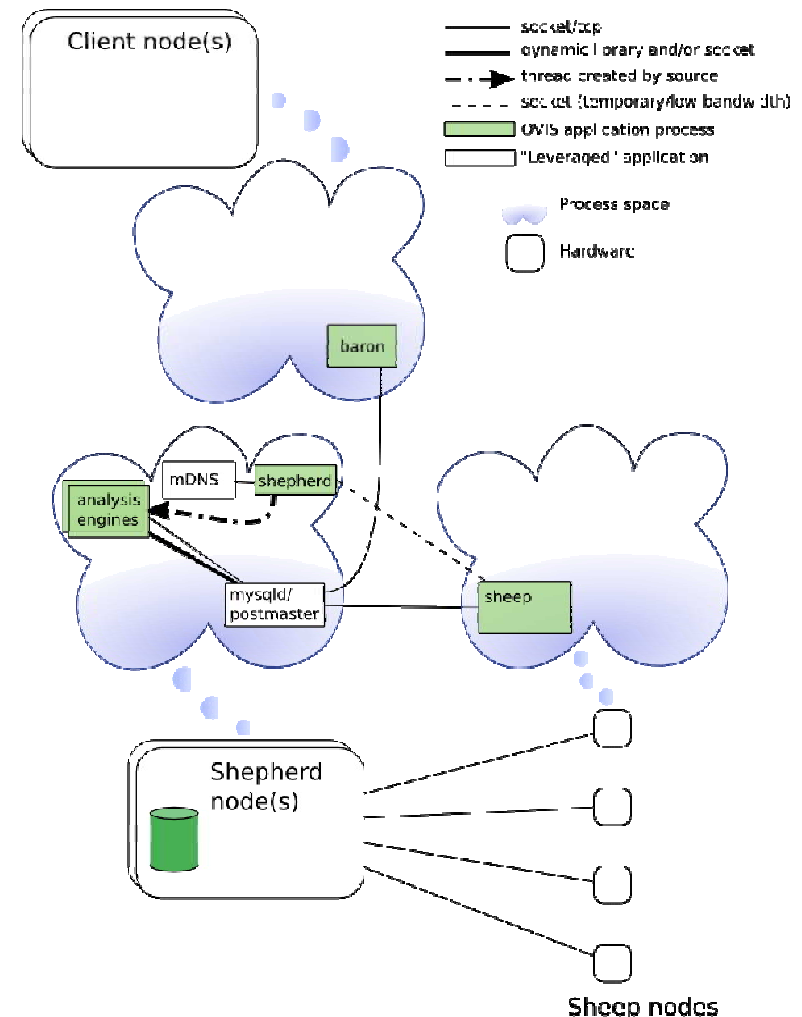
# Data Acquisition

- Impediments
  - Failures are "Rare Events" in small aggregations of machines
  - Data collection mechanisms can cause OS noise which is detrimental in large systems
    - In band (/proc, LM sensors, etc.)
  - Don't know what parameters matter so currently collect as many as possible on as fine a time granularity as possible
    - Data collection scaling issues
      - Storage/retrieval bandwidth
      - Need to store all data for long enough to allow for failure data to be gathered also

# Data Acquisition Cont.

- Out of Band (OOB) collection should have minimal impact
  - Separate processor for acquisition and export
  - Separate network for communication
  - IPMI and SNMP can facilitate out of band collection of a subset of desired parameters
- Many parameters are not available OOB
  - /proc --  CPU, Memory, Network, etc.

# Enabling Architecture for Data Collection and Fusion

- OVIS is a suite of 3 applications – baron, shepherd, sheep – sharing a common distributed database

- **Baron**: A VTK/Qt user interface (ParaView version in development)
  - Data characteristic exploration

- **Shepherd**: service-node program:
  - Advertises DB availability
  - Responds to requests for analyses

- **Sheep**: a service-node or compute-node program:
  - Listens for shepherds
  - Stores measurements to database on shepherd node

# Storage Overview

- OVIS avoids all storage except the database
- Static metadata is stored in an XML file used to initialize a DB
- For scalability, parallel distributed databases are required
  - Analysis requests/results are propagated to all DBs asynchronously
- Some metadata tables are shared between databases
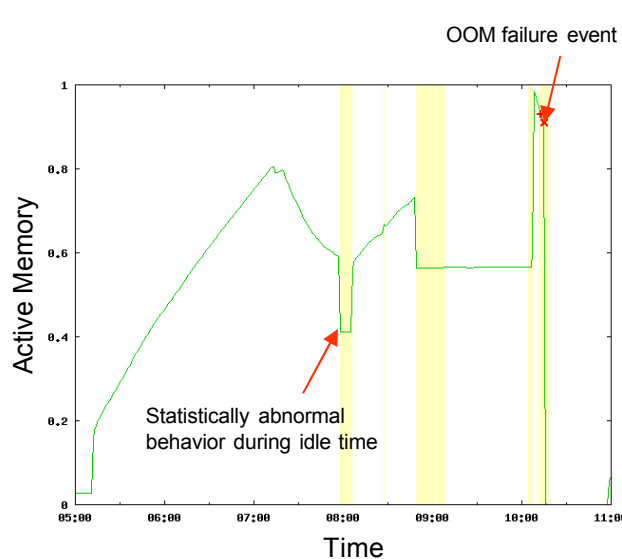- Parallel databases also imply a parallel baron (in development)

# Case Studies

- Red Storm
  - OOB RAS system caused system instability even at a once-per-hour data collection frequency
    - Too course grained
- TLCC
  - Currently collecting on 2SU (288 node, 4K core) system at Sandia Albuquerque (Glory)
    - Hardware related data – once-per-minute
    - Slurm – as it occurs
    - Log file data – before we need it for analysis
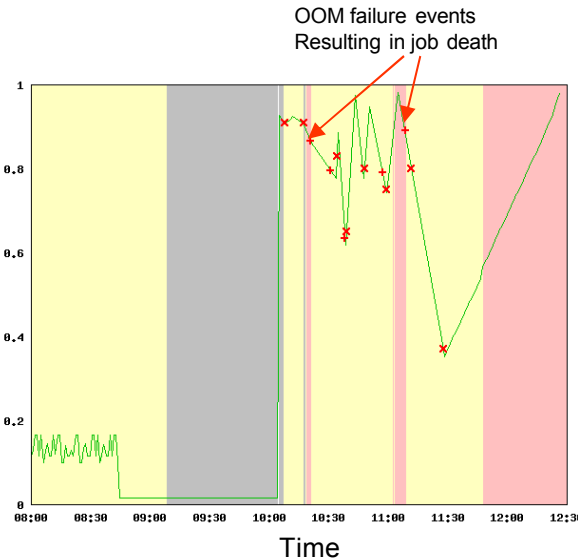
# Glory Failures

- We set out to discover precursors of the top three causes of TLCC failures
    - Power Supply
        - Collect various regulated voltages
            - Problems occur upstream of the regulators (caps, inductors, resistors)
            - Need to work with PS designers to understand what it makes sense to instrument
        - Use only a single PS which then runs closer to spec limits
    - Out of Memory
        - Precursor symptoms discovered but root cause yet unknown
    - Stuck CPU
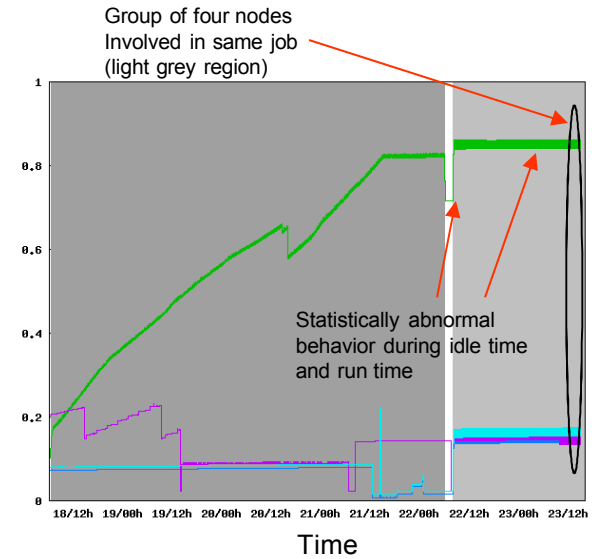        - Still under investigation

# Out of Memory (OOM)



Abnormal Active Memory utilization, detectable during any idle period (yellow shaded region), is seen to persist for hours finally resulting in process being killed due to OOM condition.

Abnormal Active Memory utilization, detectable during any idle period, is seen to result in job failures (pink shaded regions) due to the OOM condition. There are also processes that are killed due to the OOM condition even during idle periods (shaded yellow).

Abnormal Active Memory utilization was detectable both during the idle period (white region) and during the job (right-most grey shaded region). Due to this application's memory requirements the OOM threshold was not reached and this did not result in job death in this case.

*Note that Active Memory is represented as a fraction of available system memory (32GB in this case).*

# Conclusions

- While "black box" approach works for current sized systems it is at best a stop-gap approach
- We are beginning to "pierce the darkness" but still have a long way to go
  - Statistical methods show promise
  - Scalable data collection mechanisms necessary
  - Need more advanced mixed data type analysis methodologies

# Future Work

- Cleansed data sets for other researchers to use

- Deployment on additional and larger platforms

- More advanced analyses

- Tighter integration of disparate data sources

- More UI usability features with respect to plotting and graphing

# Questions?

https://ovis.ca.sandia.gov

# Motivation

- Large amount of time spent/(wasted?) in checkpointing
- Black Box approach to failures in platforms implies scaling wall
  - Best we can do is adjust checkpointing freq based on the measured system average MTTI
- Assumption -- If we understood failure mechanisms
  - Instrument to get precursor warnings where possible
  - Instrument to detect fine granularity failures that could still allow state recovery and hence restart on other resources
  - System and Application interfaces that would facilitate such recovery
  - Know where to invest more up front for higher level of fault tolerance
    - Power supply redundancy cost/fault tolerance tradeoff
    - Memory socket cost/MTTF tradeoff