

Toward Lightweight, Actionable Analytical Tools Based on Statistical Learning for Efficient System Operations

Devesh Tiwari

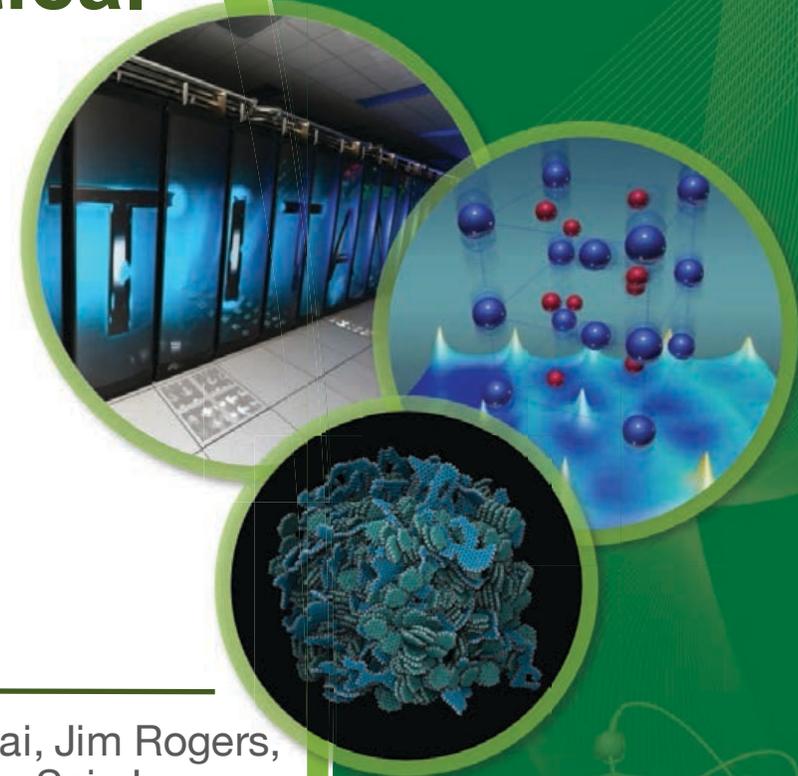
Oak Ridge National Laboratory

tiwari@ornl.gov

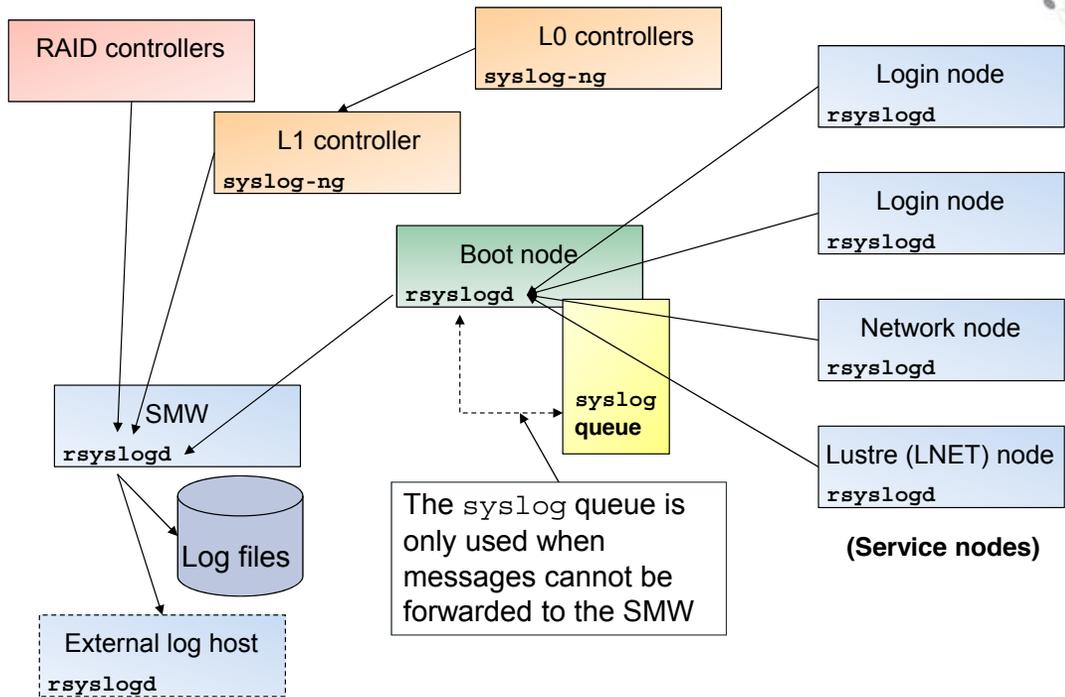
Thanks to Saurabh Gupta, Christian Engelmann, Sudharshan Vazhkudai, Jim Rogers, Bin Nie, Evgenia Smirni, Franck Cappello, Rinku Gupta, Sheng Di, Marc Snir, Leo Bautista-Gomez, Swann Perarnau, Nathan Debardeleben, Paolo Rech, Don Maxwell, Changhee Jung, Martin Schulz, Ignacio Laguna, and many more smart folks!

This work used the resources of the Oak Ridge Leadership Computing Facility, located in the National Center for Computational Sciences at the Oak Ridge National Laboratory, which is managed by UT Battelle, LLC for the U.S. Department of Energy, under the contract No. DEAC05-00OR22725.

ORNL is managed by UT-Battelle
for the US Department of Energy



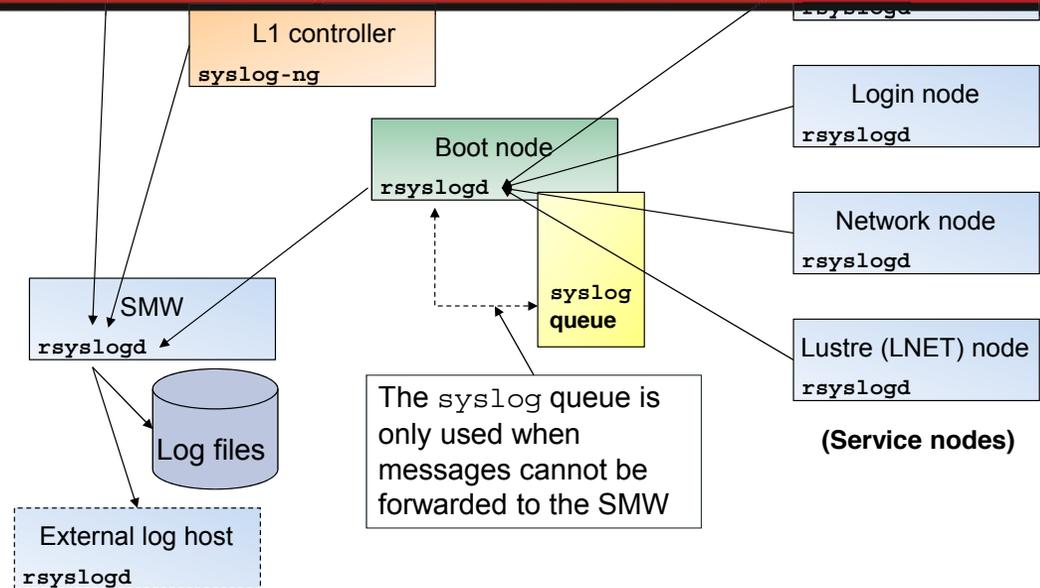
System Generated Data for ModSim Efforts



Courtesy: Rick Slick CUG 2013 Tutorial Slides

System Generated Data for ModSim Efforts

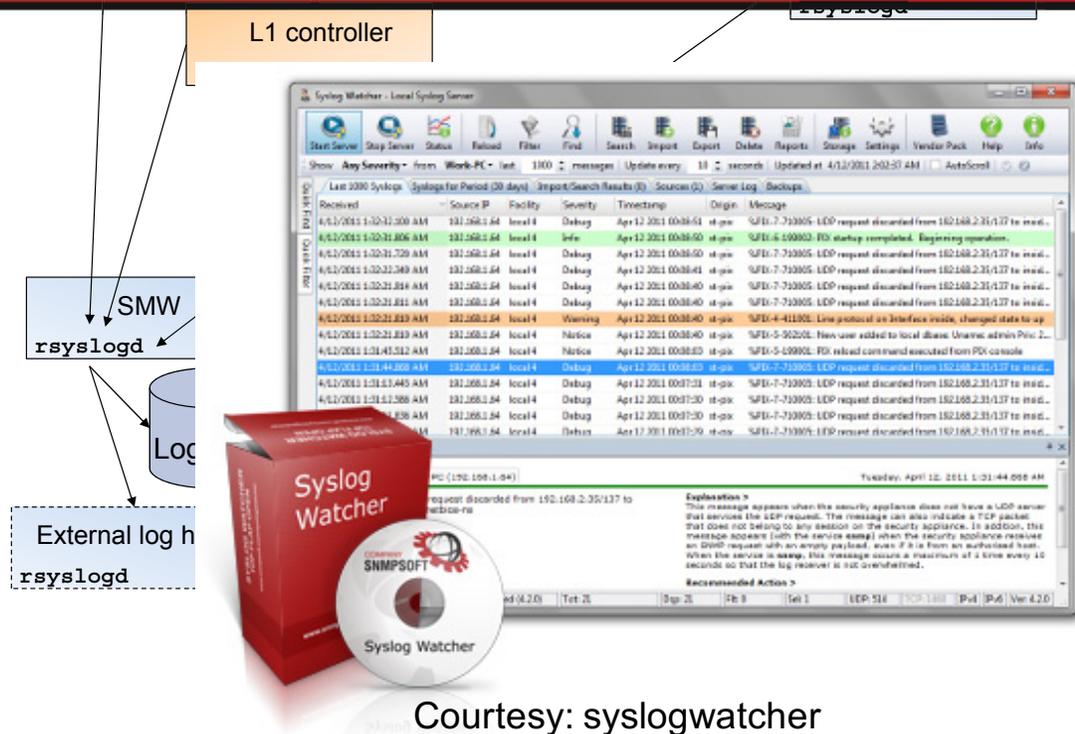
Hard to store and manage



Courtesy: Rick Slick CUG 2013 Tutorial Slides

System Generated Data for ModSim Efforts

Hard to store and manage

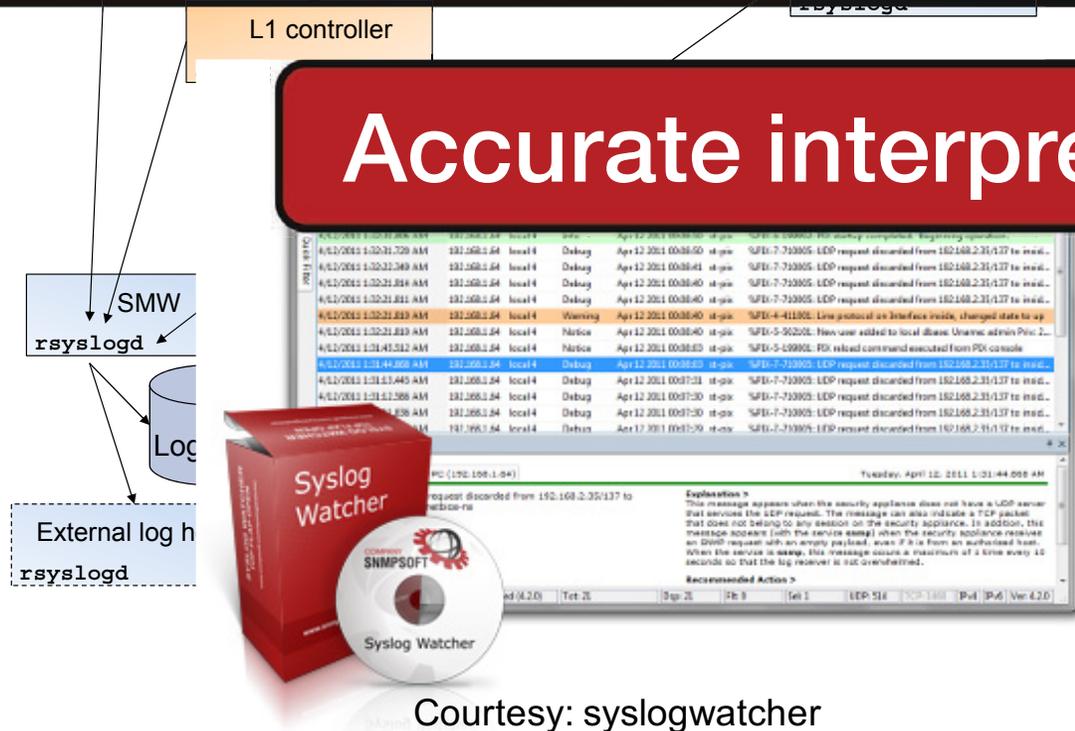


Courtesy: syslogwatcher

System Generated Data for ModSim Efforts

Hard to store and manage

Accurate interpretation hard

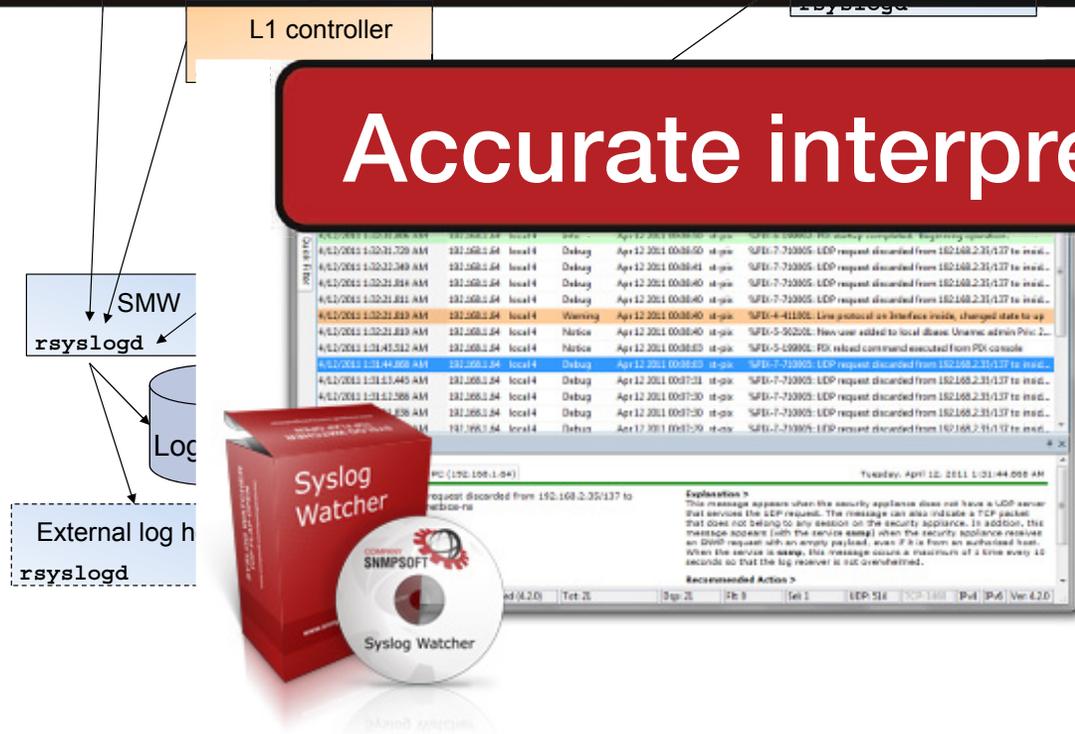


Courtesy: syslogwatcher

System Generated Data for ModSim Efforts

Hard to store and manage

Accurate interpretation hard



Time	IP	Host	Level	Message
Apr 12 2011 00:06:50	192.168.2.31	local4	Debug	%SYS-5-305002: RDP session completed. Beginning operation...
Apr 12 2011 00:06:50	192.168.2.31	local4	Debug	%SYS-7-703005: LDP request discarded from 192.168.2.31(137) to local...
Apr 12 2011 00:06:40	192.168.2.31	local4	Debug	%SYS-7-703005: LDP request discarded from 192.168.2.31(137) to local...
Apr 12 2011 00:06:40	192.168.2.31	local4	Debug	%SYS-7-703005: LDP request discarded from 192.168.2.31(137) to local...
Apr 12 2011 00:06:40	192.168.2.31	local4	Debug	%SYS-7-703005: LDP request discarded from 192.168.2.31(137) to local...
Apr 12 2011 00:06:40	192.168.2.31	local4	Warning	%SYS-4-411001: Line protocol on interface inside, changed state to up
Apr 12 2011 00:06:40	192.168.2.31	local4	Notice	%SYS-5-902001: New user added to local database: Username: admin Priv: 2...
Apr 12 2011 00:06:03	192.168.2.31	local4	Notice	%SYS-5-109001: PDK reload command executed from PDK console
Apr 12 2011 00:05:07	192.168.2.31	local4	Debug	%SYS-7-703005: LDP request discarded from 192.168.2.31(137) to local...
Apr 12 2011 00:07:31	192.168.2.31	local4	Debug	%SYS-7-703005: LDP request discarded from 192.168.2.31(137) to local...
Apr 12 2011 00:07:30	192.168.2.31	local4	Debug	%SYS-7-703005: LDP request discarded from 192.168.2.31(137) to local...
Apr 12 2011 00:07:30	192.168.2.31	local4	Debug	%SYS-7-703005: LDP request discarded from 192.168.2.31(137) to local...
Apr 17 2011 00:17:30	192.168.2.31	local4	Notice	%SYS-7-703005: LDP request discarded from 192.168.2.31(137) to local...



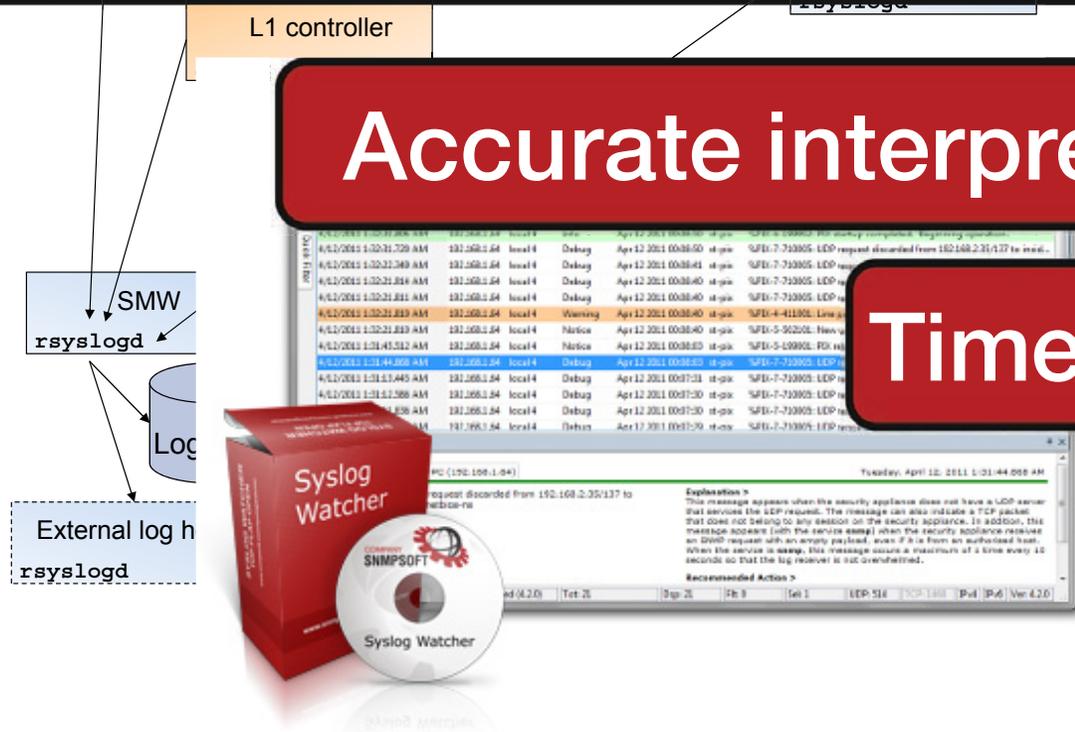
System Generated Data for ModSim Efforts

Hard to store and manage

Accurate interpretation hard

Timely processing and analysis

Insights

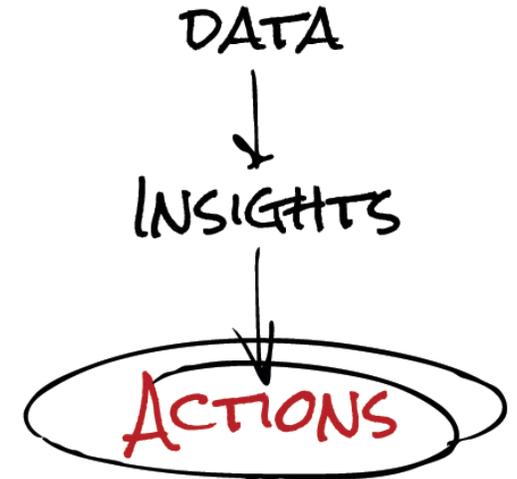
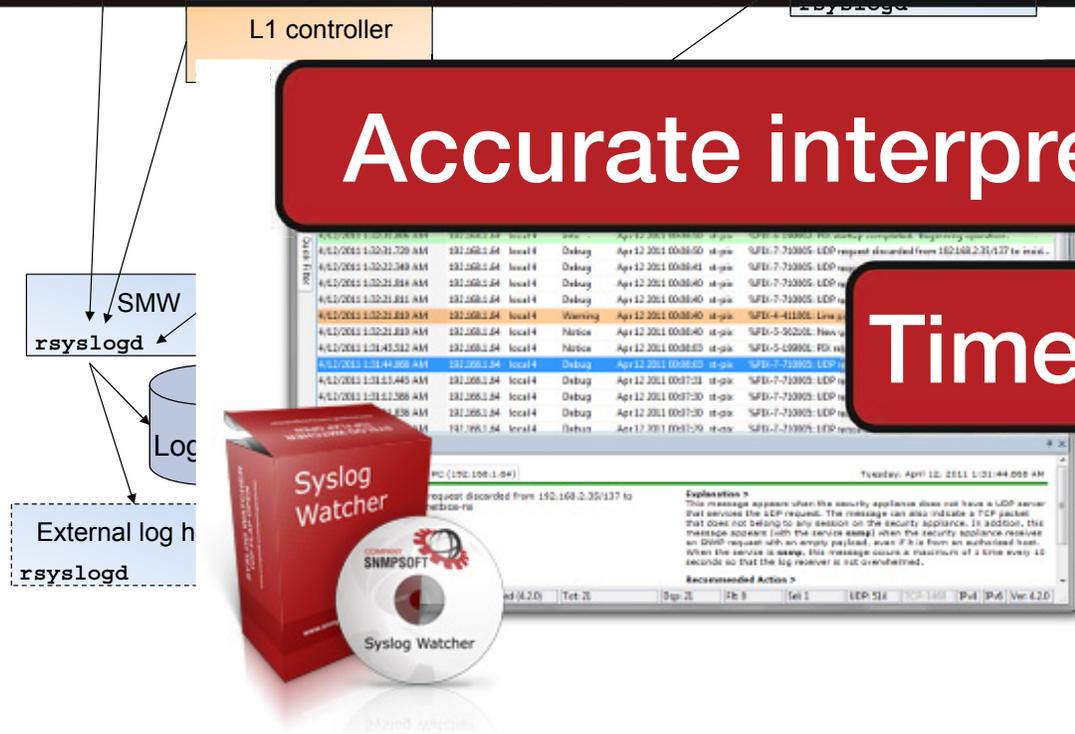


System Generated Data for ModSim Efforts

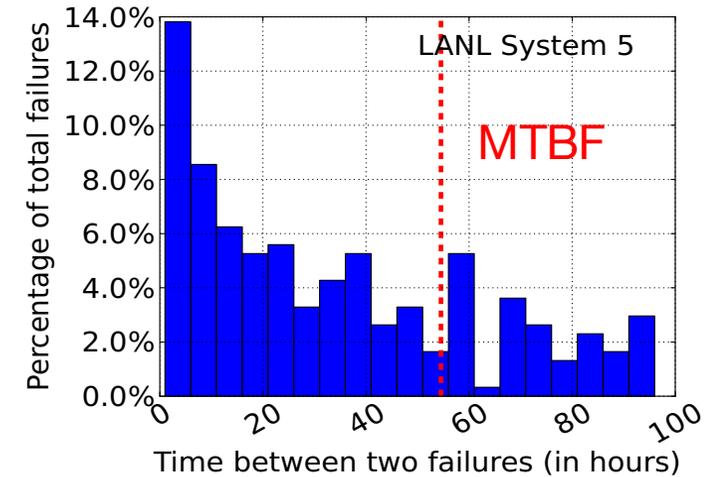
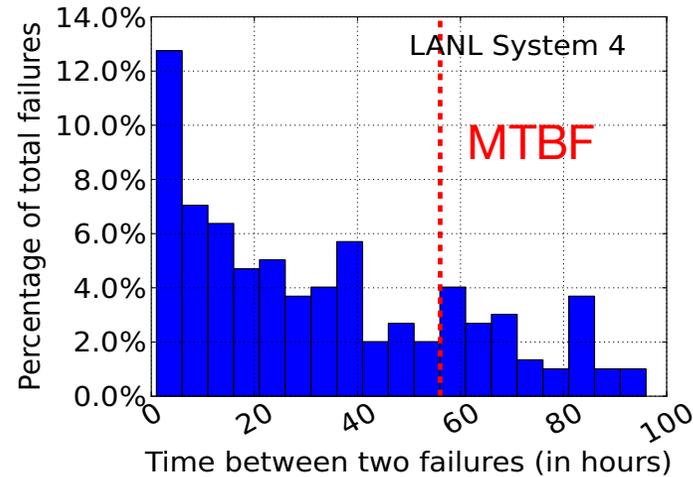
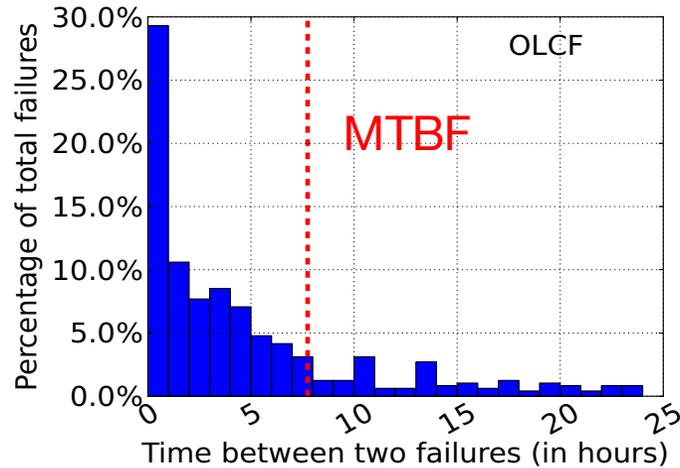
Hard to store and manage

Accurate interpretation hard

Timely processing and analysis

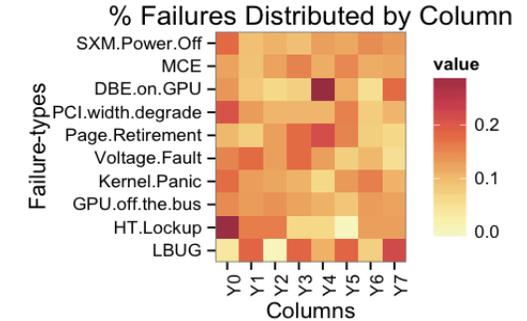
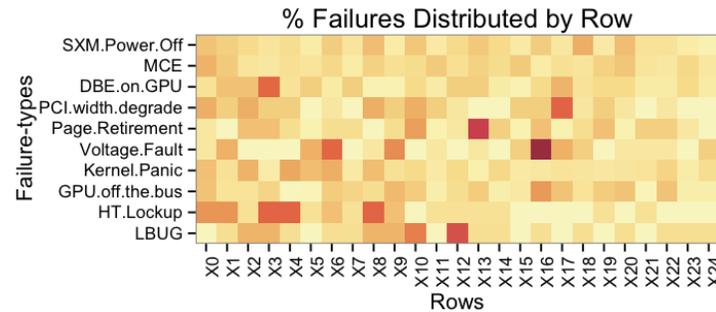
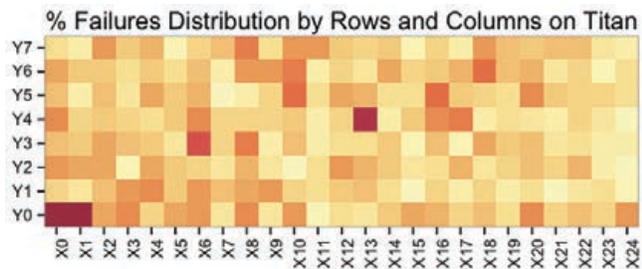
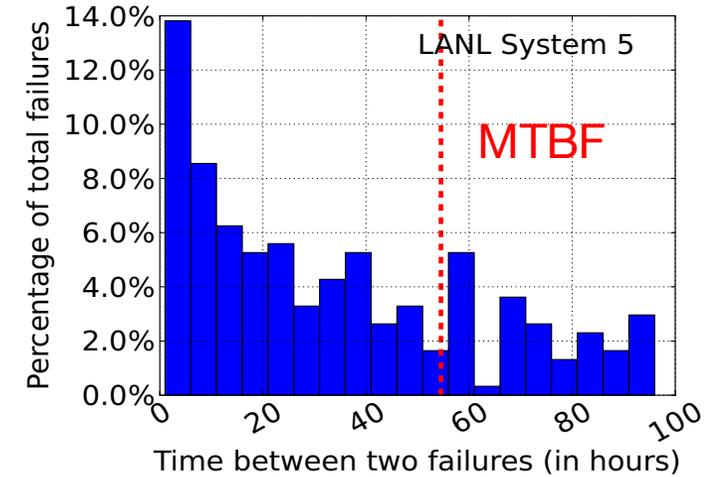
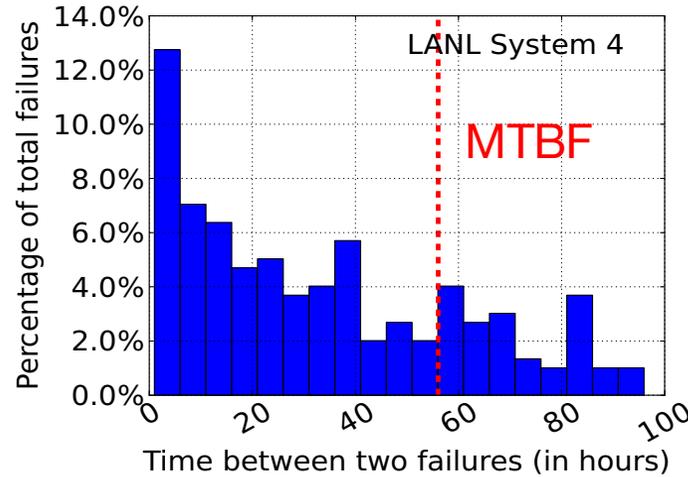
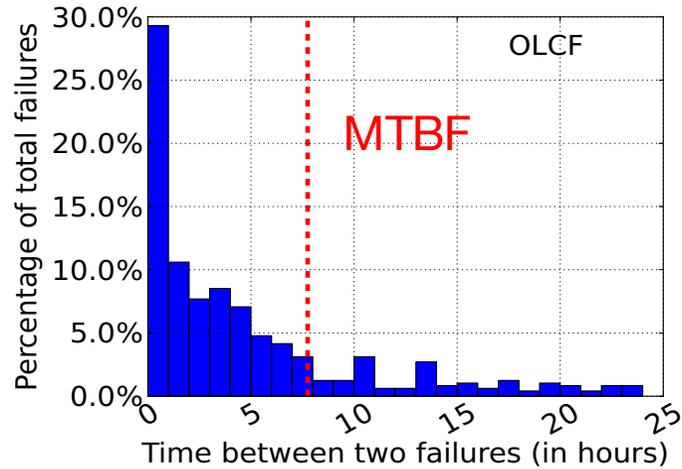


System failures exhibit temporal and spatial locality.



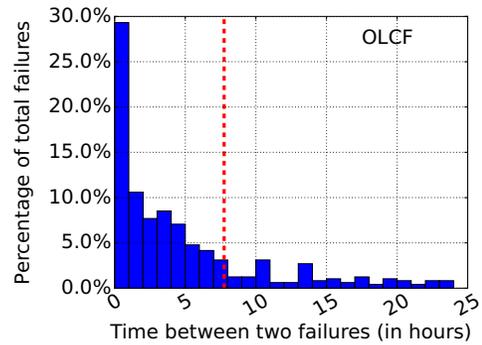
Observation holds true across systems and failure types, consistently across long range of periods.

System failures exhibit temporal and spatial locality.

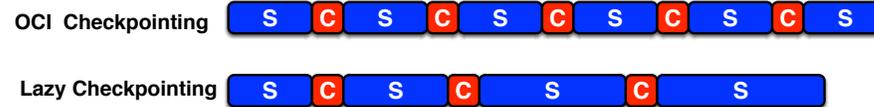
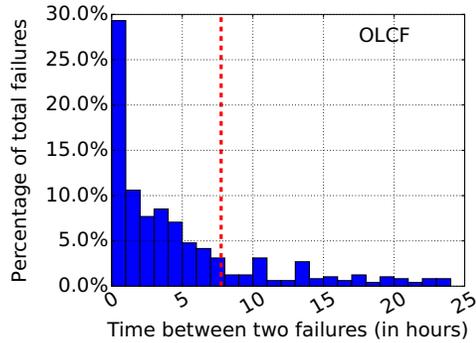


Observation holds true across systems and failure types, consistently across long range of periods.

Lazy Checkpointing Technique



Lazy Checkpointing Technique



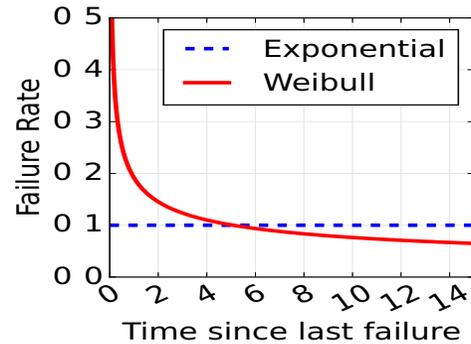
Weibull Distribution

Failure Rate $\frac{k}{\lambda} \left(\frac{t}{\lambda}\right)^{k-1}$

t Time since last failure

λ Scale parameter

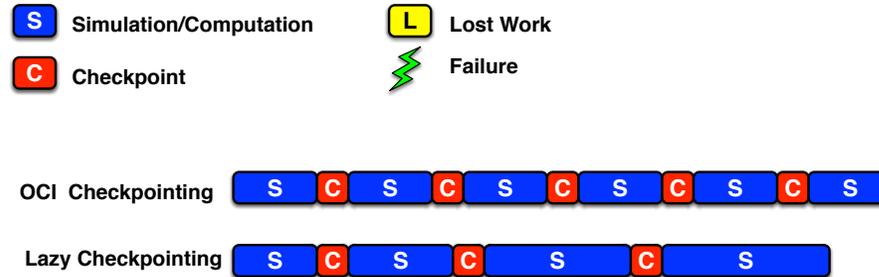
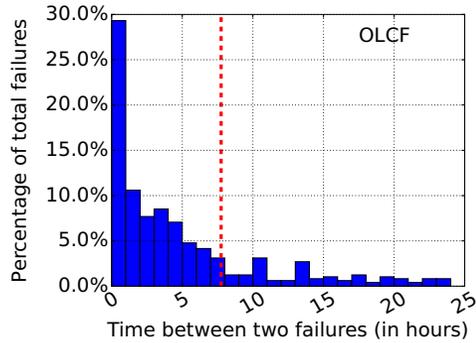
k Shape parameter



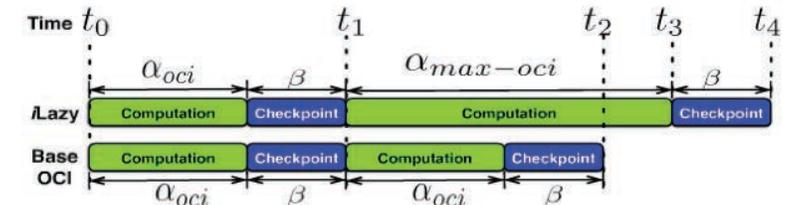
$$\alpha_{oci} = \sqrt{\beta^2 + \frac{\beta\gamma}{\epsilon} + \frac{M\beta}{\epsilon}}$$

$$\alpha_{lazy} = \alpha_{oci} \left(\frac{t}{\alpha_{oci}}\right)^{(1-k)}$$

Lazy Checkpointing Technique



Key is to balance the trade-off between reduction in checkpointing overhead and possible increase in the waste work



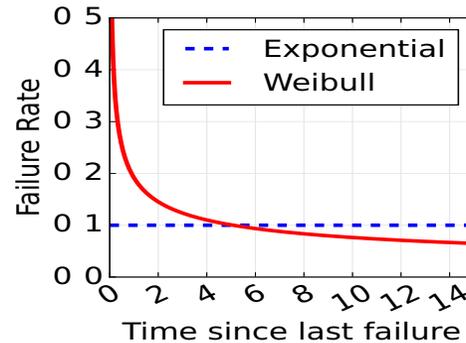
Weibull Distribution

$$\text{Failure Rate} = \frac{k}{\lambda} \left(\frac{t}{\lambda}\right)^{k-1}$$

t Time since last failure

λ Scale parameter

k Shape parameter



$$\alpha_{oci} = \sqrt{\beta^2 + \frac{\beta\gamma}{\epsilon} + \frac{M\beta}{\epsilon}}$$

$$\alpha_{lazy} = \alpha_{oci} \left(\frac{t}{\alpha_{oci}}\right)^{(1-k)}$$

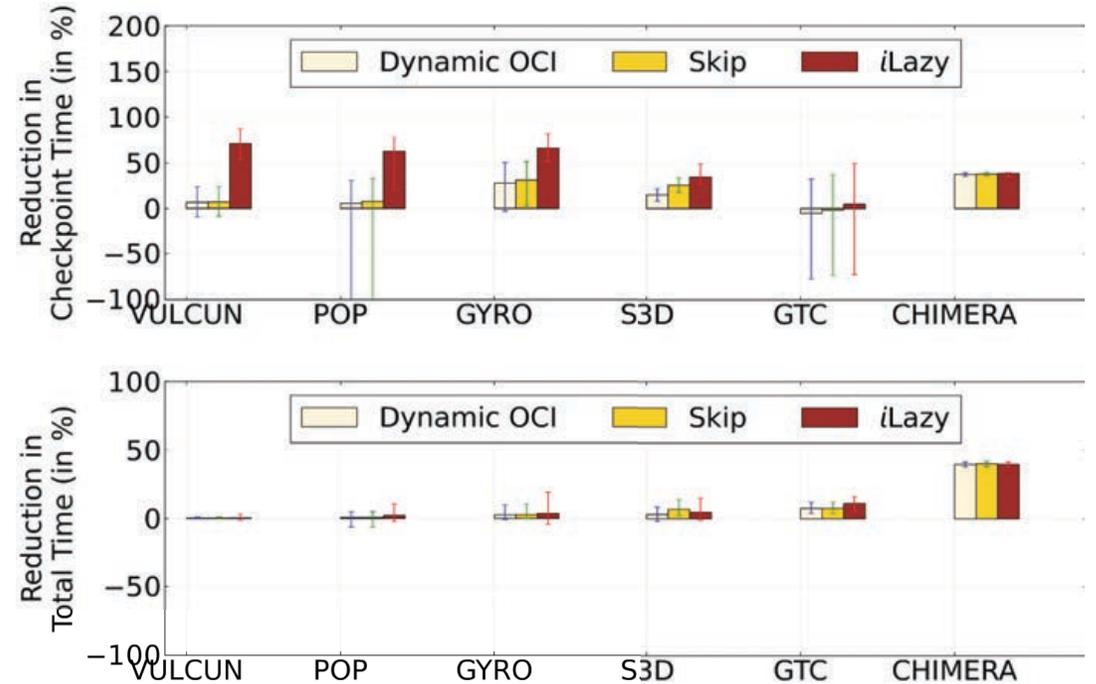
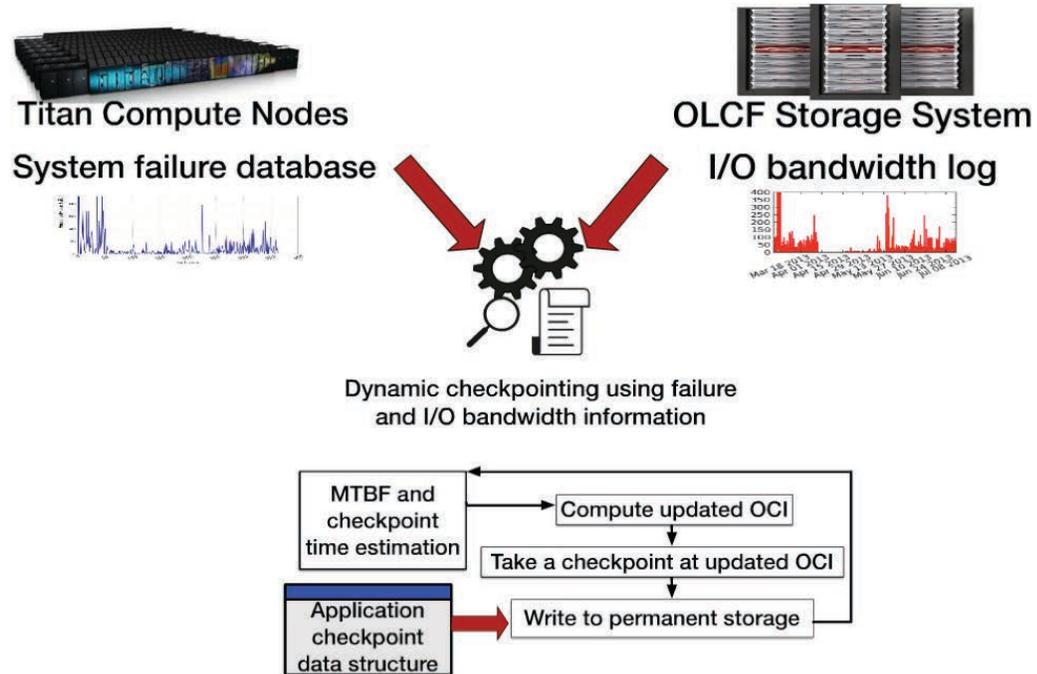
$$\begin{aligned} \text{performance loss} &= (\alpha_{max-oci} - \alpha_{oci}) \left(e^{-\left(\frac{t_2}{\lambda}\right)^k} - e^{-\left(\frac{t_4}{\lambda}\right)^k} \right) \\ &= (\alpha_{max-oci} - \alpha_{oci}) \left(e^{-\left(\frac{2(\alpha_{oci} + \beta)}{\lambda}\right)^k} - e^{-\left(\frac{\alpha_{max-oci} + \alpha_{oci} + 2\beta}{\lambda}\right)^k} \right) \end{aligned}$$

$$\text{performance gain} = \beta e^{-\left(\frac{t_3}{\lambda}\right)^k}$$

$$\begin{aligned} \beta e^{-\left(\frac{\alpha_{max-oci} + \alpha_{oci} + \beta}{\lambda}\right)^k} &= (\alpha_{max-oci} - \alpha_{oci}) e^{-\left(\frac{2(\alpha_{oci} + \beta)}{\lambda}\right)^k} \\ &\quad - (\alpha_{max-oci} - \alpha_{oci}) e^{-\left(\frac{\alpha_{max-oci} + \alpha_{oci} + 2\beta}{\lambda}\right)^k} \end{aligned}$$

Refer to the paper for model validation and simulation results

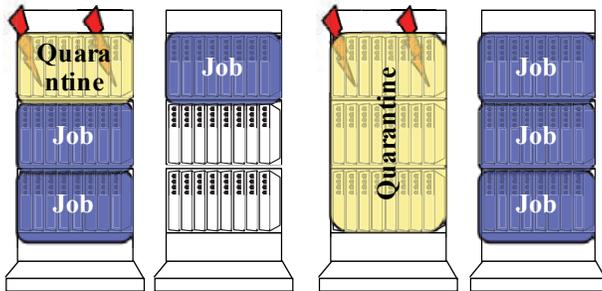
Lazy Checkpointing Technique



Quarantine Job Scheduling Technique

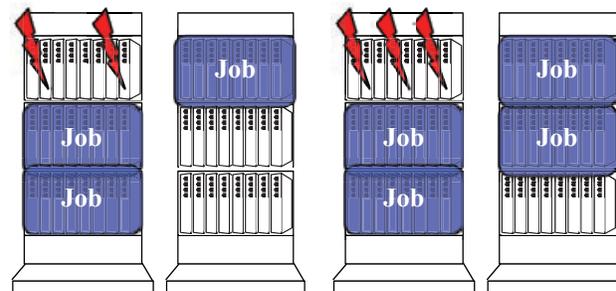
Idea: On job restart or a new job allocation a fraction of compute capacity is not utilized (quarantined)

Quarantine Granularity



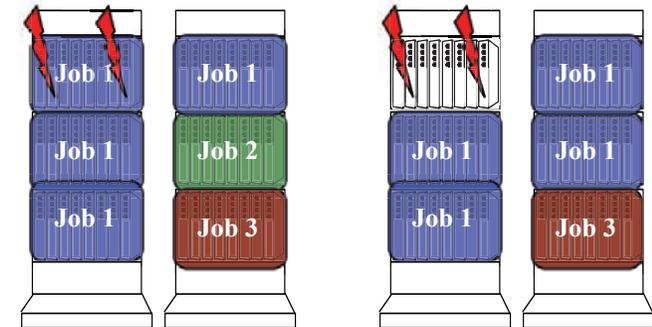
Fraction of avoided system failures versus compute resource waste

Quarantine Time Duration



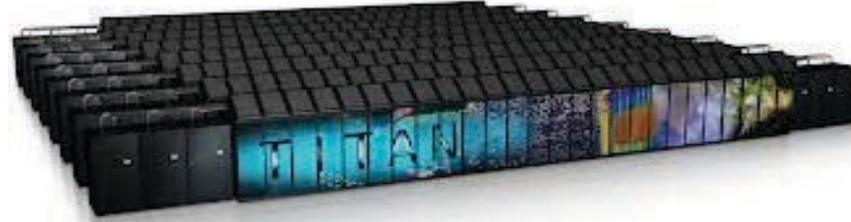
Diminishing returns on the number of avoided failures

System Utilization vs. Reliability

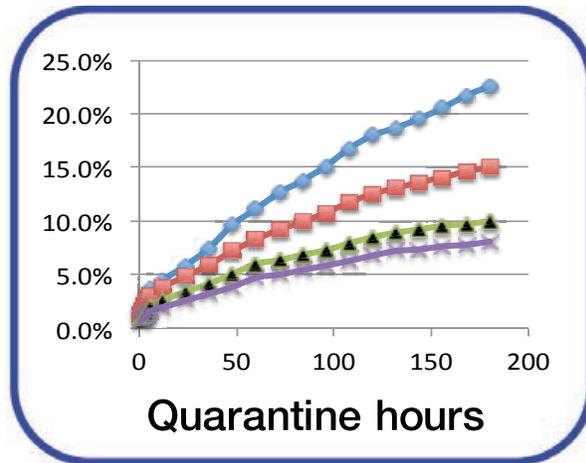


Trading-off lower system utilization for improved reliability

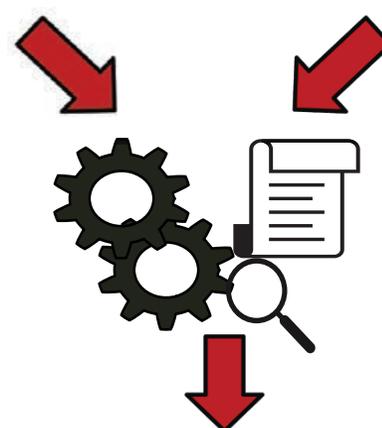
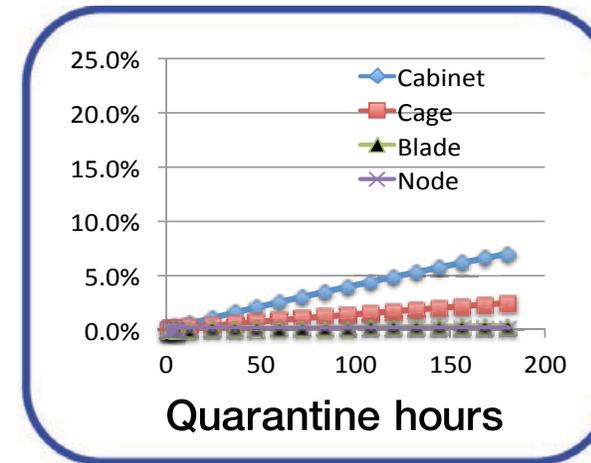
Quarantine Job Scheduling Technique



System Reliability
Fraction of failures avoided

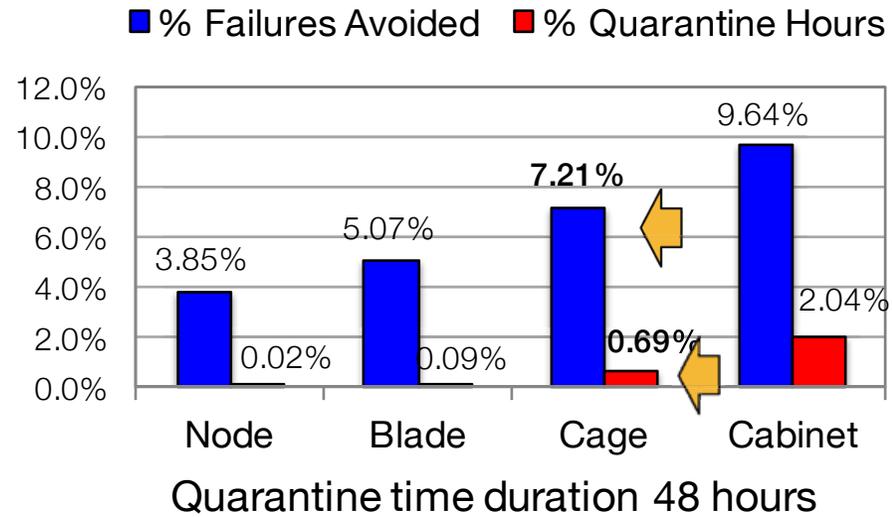


System Utilization
Quarantine node hours

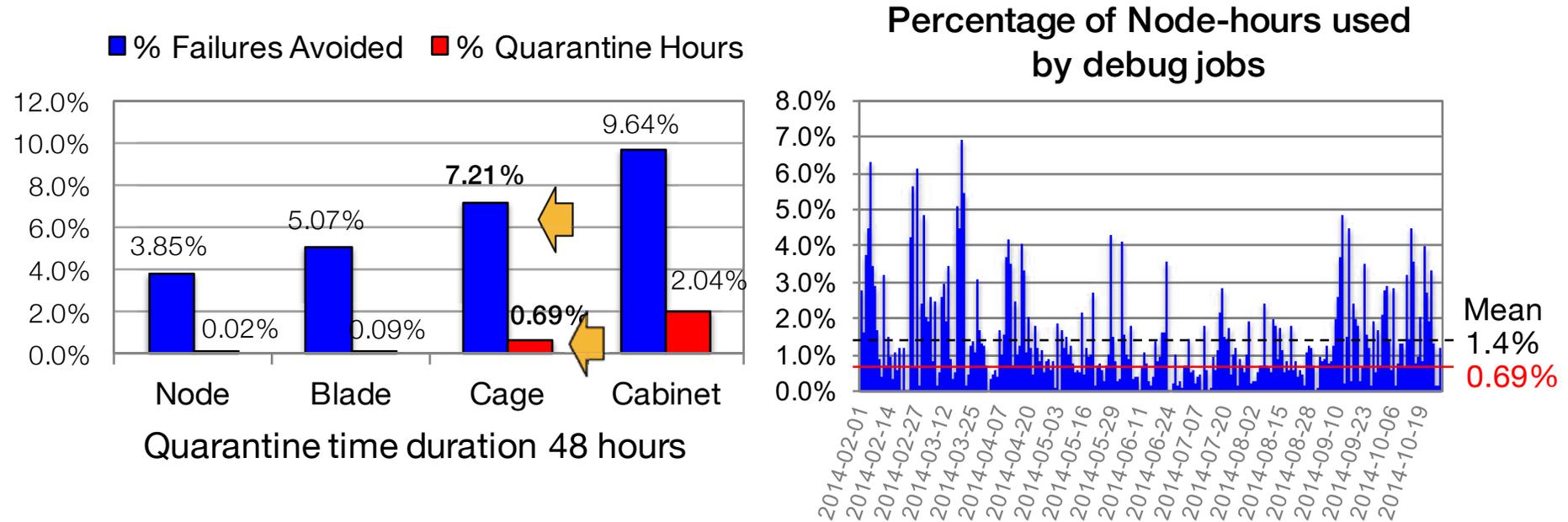


Feedback to the job scheduler

Quarantine Job Scheduling Technique



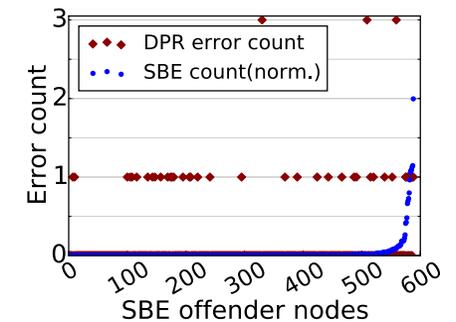
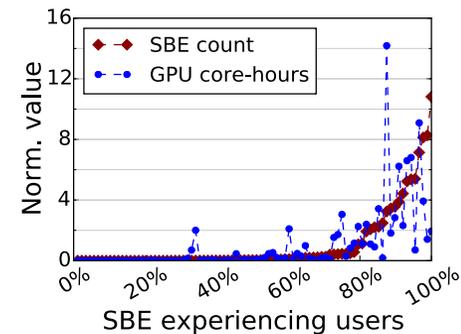
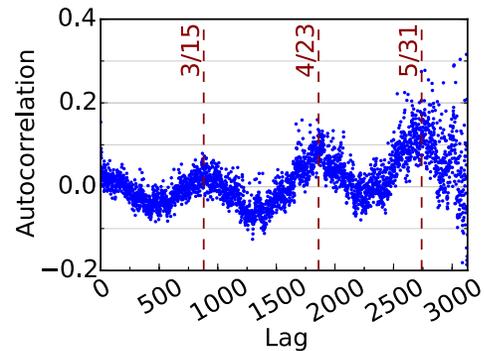
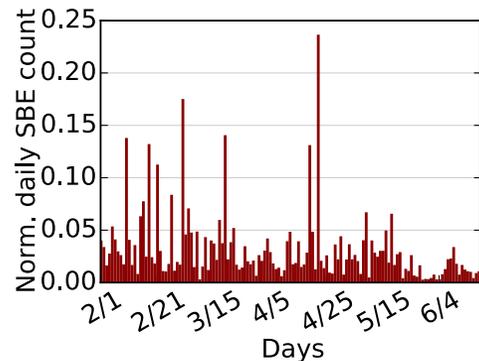
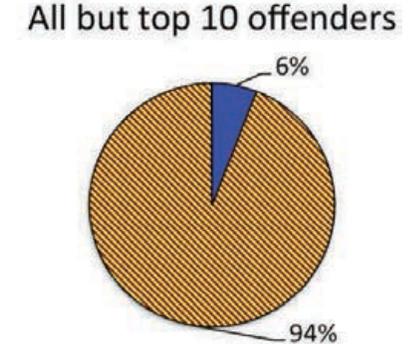
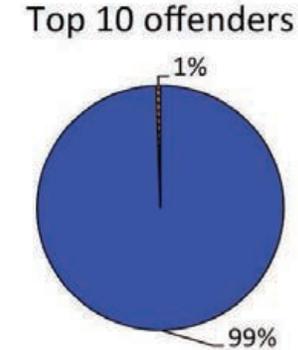
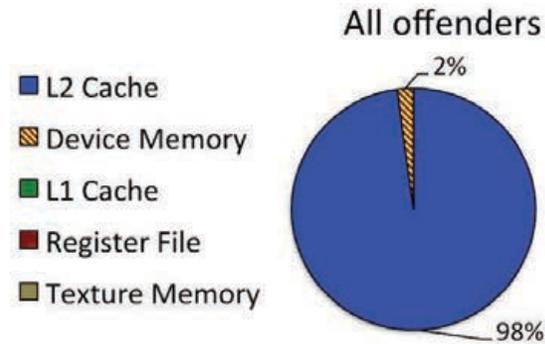
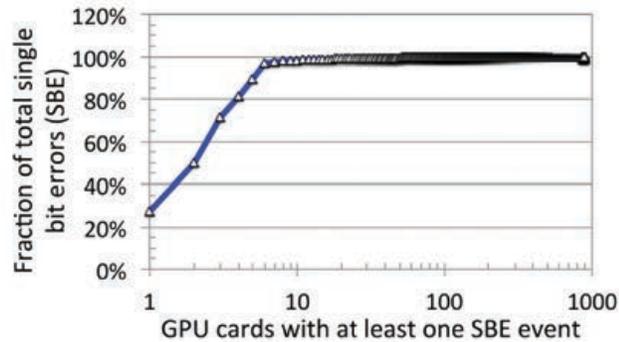
Quarantine Job Scheduling Technique



Significant fraction of failures can be avoided from interrupting production applications

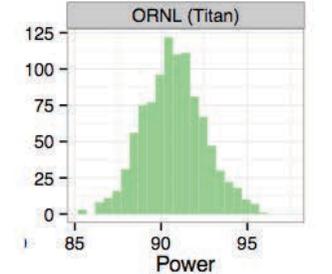
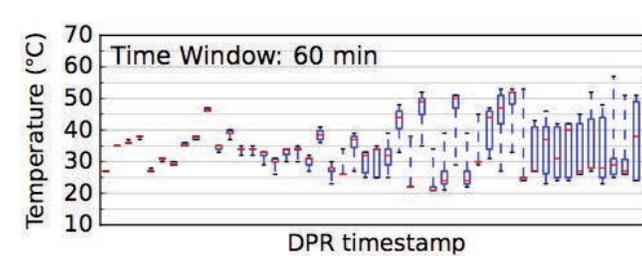
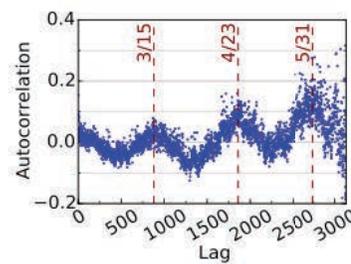
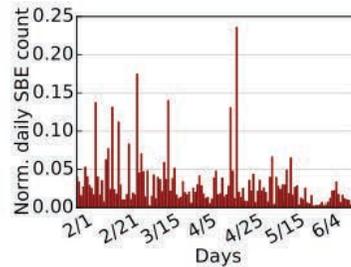
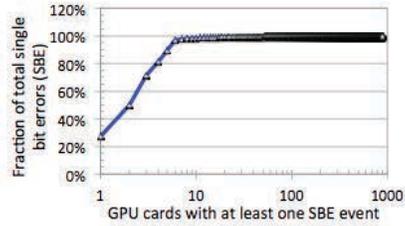
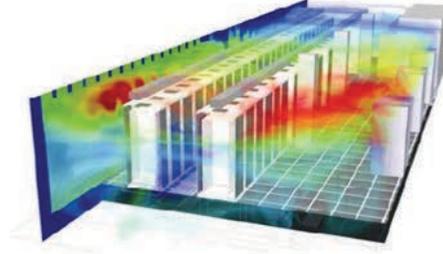
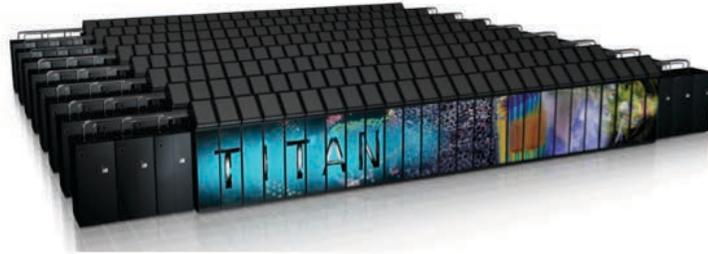
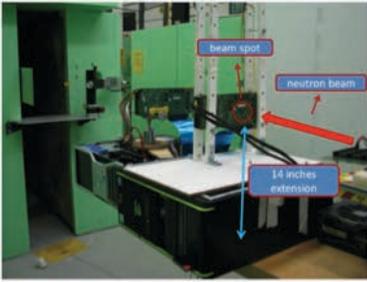
Debug or non-production jobs can be scheduled on quarantine nodes

Feedback Loop for ModSim Community



These insights can potentially change the way we design fault-injection modsim frameworks, operate production machines, and plan for future systems.

Opportunity for Actionable Analytical Tools



Workload and system generated data

Environment and facilities data

How can we (machine learning) fuse all this data to take meaningful, timely, and profitable decisions at-scale?

My Personal View

Future large scale system will have heterogeneity in terms reliability levels, too.

Parts of large systems will go in transient lower reliability, degraded performance, and large performance variability modes.

Traditional “replace and continue” approach will not be sustainable.

We will need theoretically-sound techniques and tools to “dynamically” manage this new kind of heterogeneity.

Traditional “replace and continue” approach will not be sustainable.

We will need theoretically-sound techniques and tools to “dynamically” manage this new kind of heterogeneity.

Denial and blame shifting will continue to work for some time in near future. 😊

Thanks!

Devesh Tiwari

Oak Ridge National Laboratory

tiwari@ornl.gov

devesh.dtiwari@gmail.com