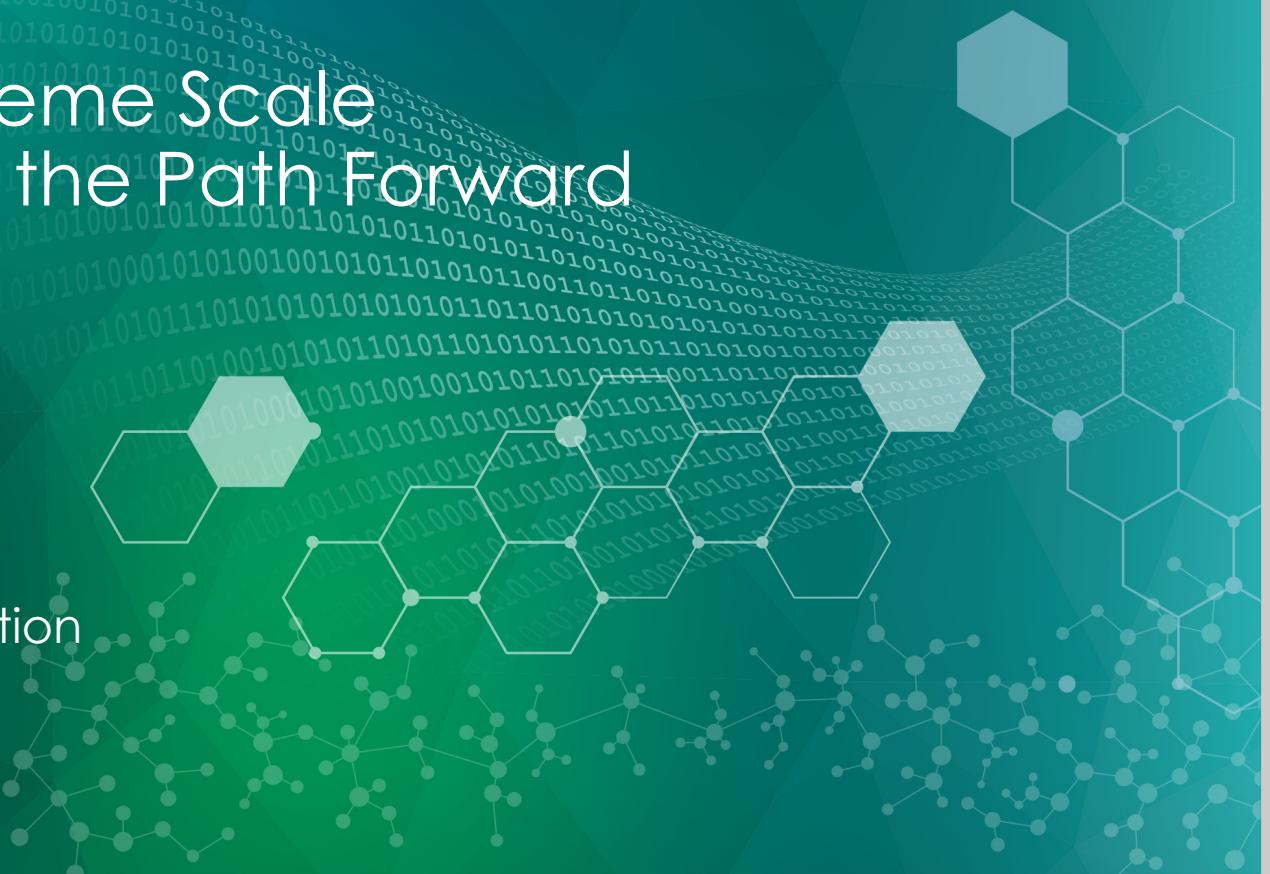


The Resilience Problem in Extreme Scale Computing: Experiences and the Path Forward

Christian Engelmann, Ph.D.

Senior Scientist & Group Leader
Intelligent Systems and Facilities Group
Advanced Computing Systems Research Section
Computer Science and Mathematics Division
Oak Ridge National Laboratory



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

Motivation

- Resilience in extreme-scale supercomputers is an optimization problem between the key design and deployment cost factors:
 - Performance, resilience, and power consumption
- The challenge is to build a reliable system within a given cost budget that achieves the expected performance.
- This requires fully understanding the resilience problem and offering efficient resilience mitigation technologies.
 - What is the fault model of such systems?
 - What is the impact of faults on applications?
 - How can mitigation in hard-/software help and at what cost?

Characterizing Supercomputer Faults, Errors and Failures

Novel Ideas:

- Applies a unified taxonomy for supercomputer faults, errors and failures
 - Understanding resilience is a data analytics problem, requiring fusion and analysis of different logs and system health data

Impact:

- Develops an understanding of observed and inferred supercomputer reliability conditions
 - Extrapolates this knowledge to future systems
 - Enables the systematic improvement of resilience in extreme-scale systems
 - Keeps applications running to a correct solution in a timely and efficient manner in spite of frequent faults, errors, and failures

Accomplishments:

- Analyzed 1.2 billion node hours of logs from the Jaguar, Titan, and Eos systems at OLCF
 - Developed tools for analyzing logs and creating a fault, error and failure catalog
 - Created novel modeling techniques to characterize temporal and spatial failure behavior

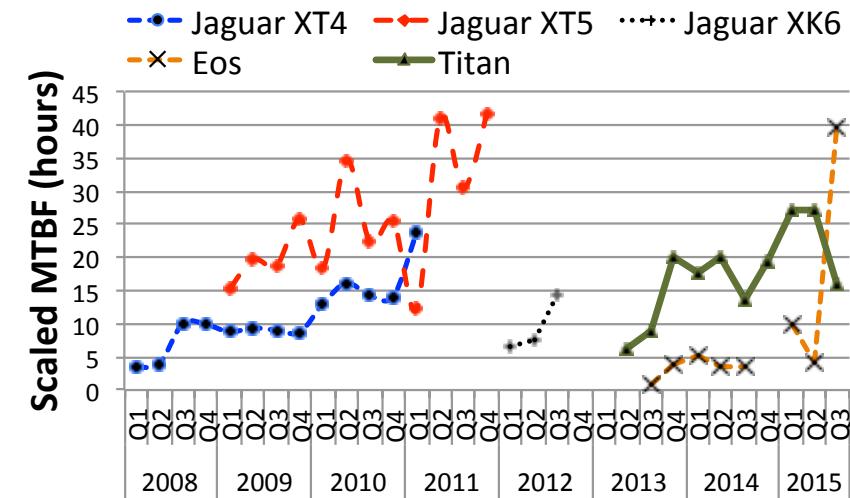
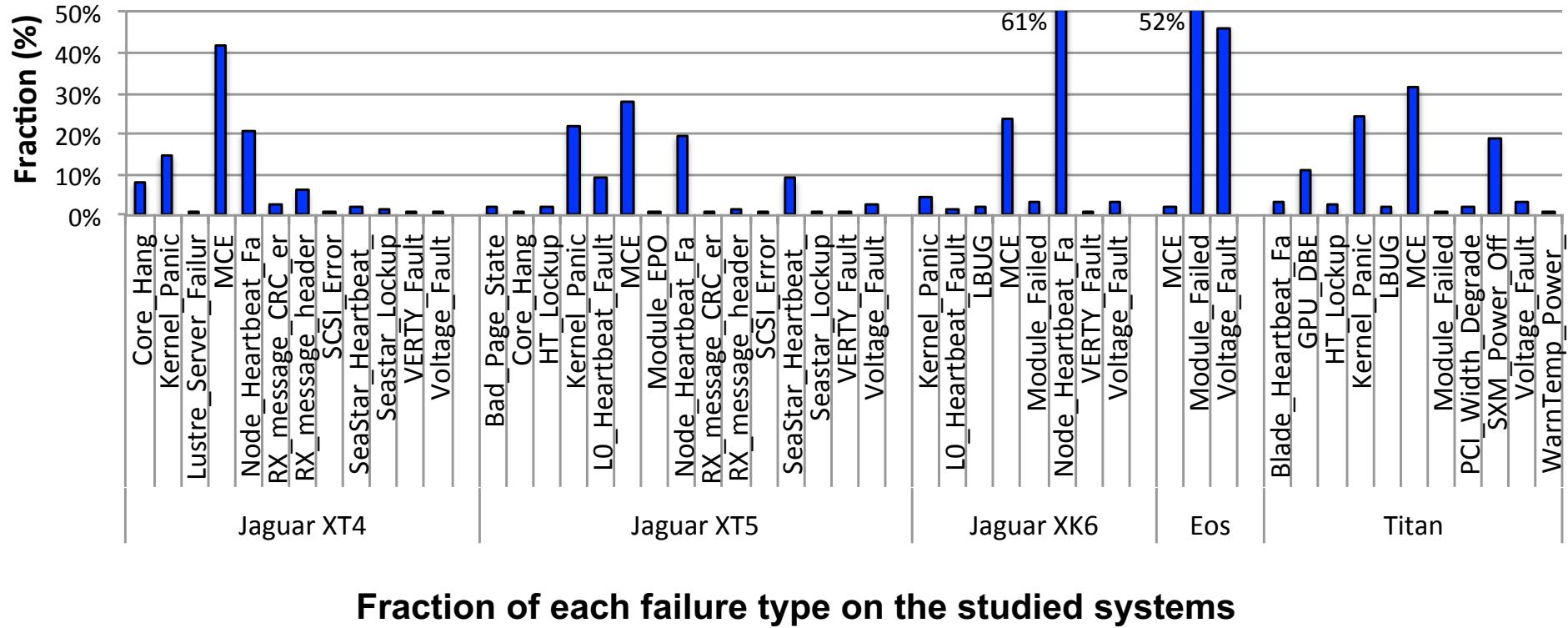


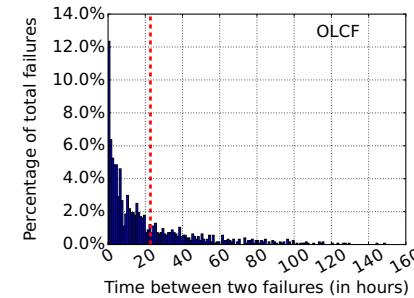
Figure: Each system goes through phases of high and low stability due to continuous efforts of system administrators to improve overall system reliability

Characterizing Supercomputer Faults, Errors and Failures

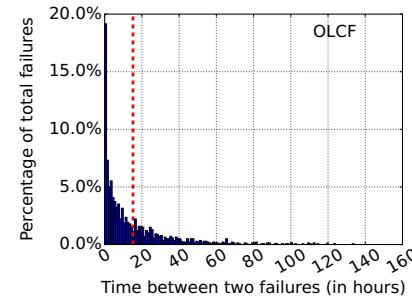


Fraction of each failure type on the studied systems

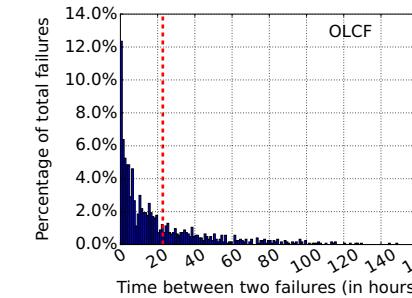
Characterizing Supercomputer Faults, Errors and Failures



(a) Jaguar XT4

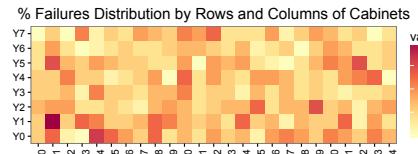


(b) Titan

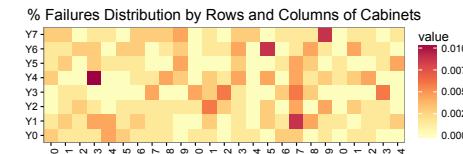


(c) Eos

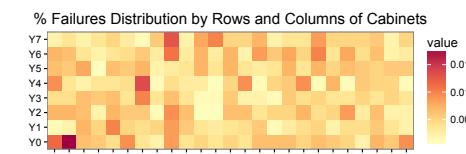
Failure inter-arrival time for 3 studied systems (MTBF as red vertical line)



(a) Jaguar XT5

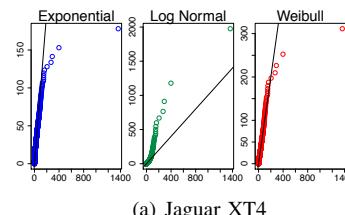


(b) Jaguar XK6

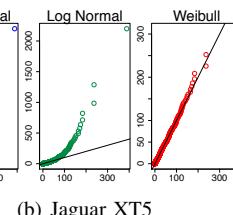


(c) Titan

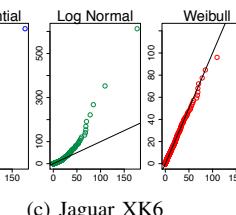
Spatial distribution of failures among cabinets for 3 studied systems



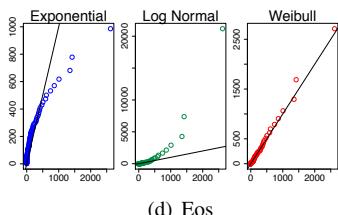
(a) Jaguar XT4



(b) Jaguar XT5



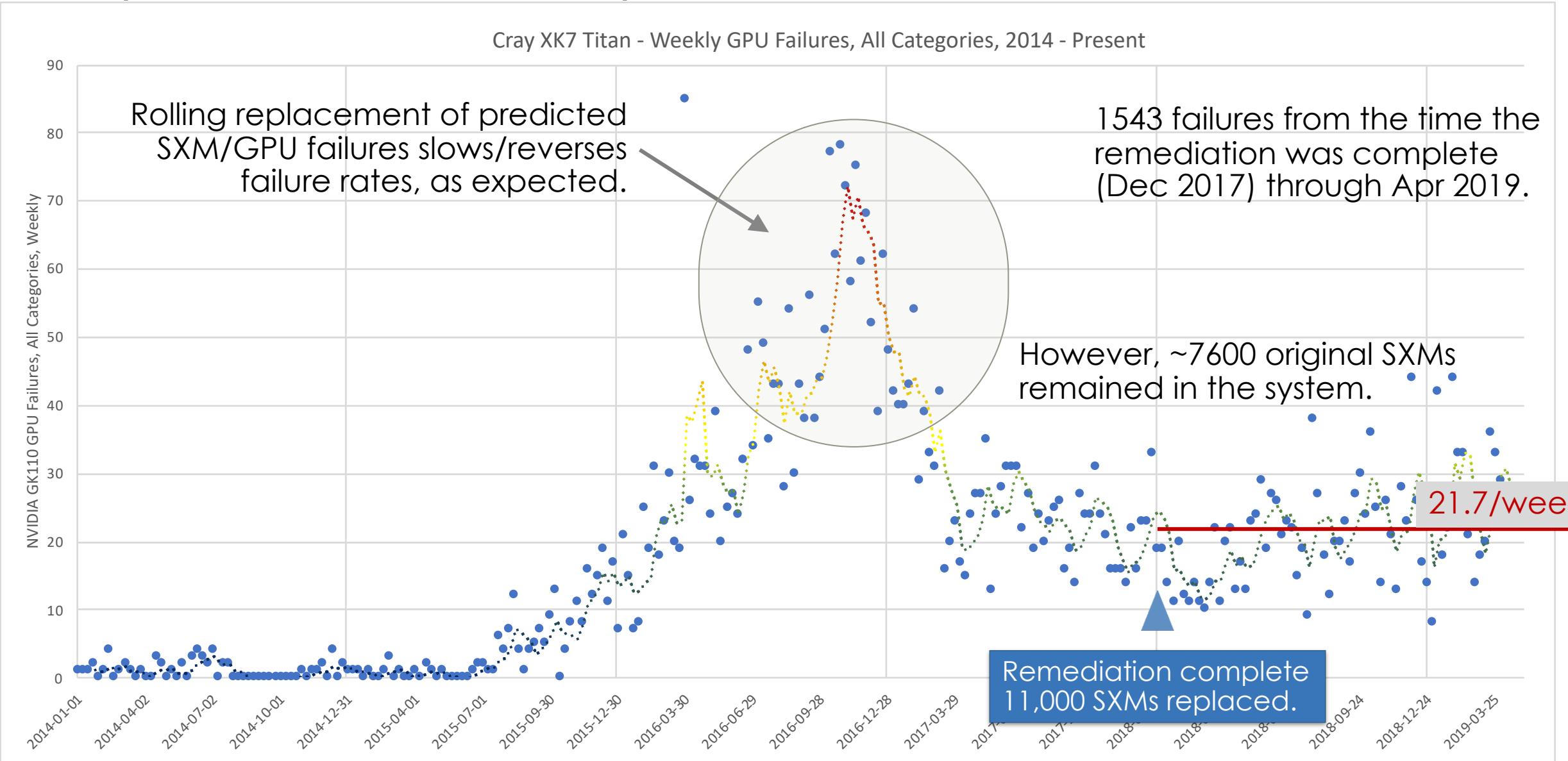
(c) Jaguar XK6



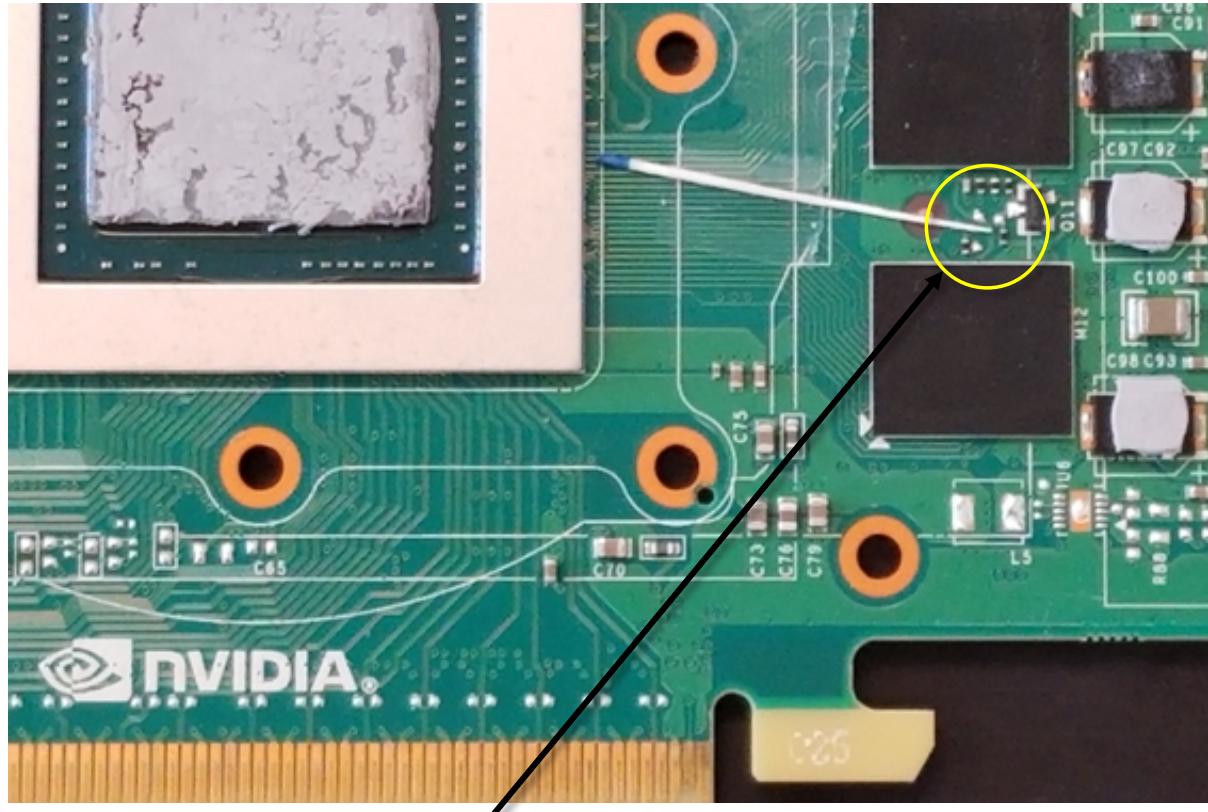
(d) Eos

QQ-plots showing goodness of fit for the failure inter-arrival times for 4 studied systems with different failure probability density functions (Weibull fits best)

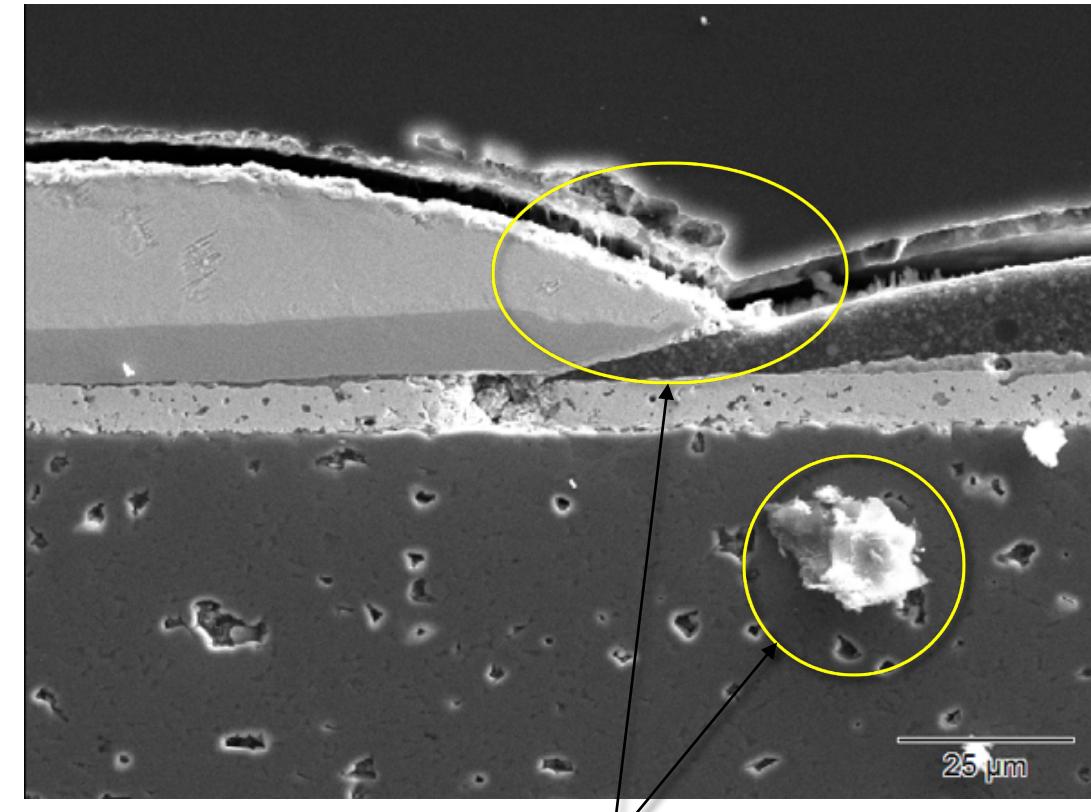
Cray XK7 Titan – Weekly GPU Failures



Root Cause: Non-ASR Components on SXM GPU



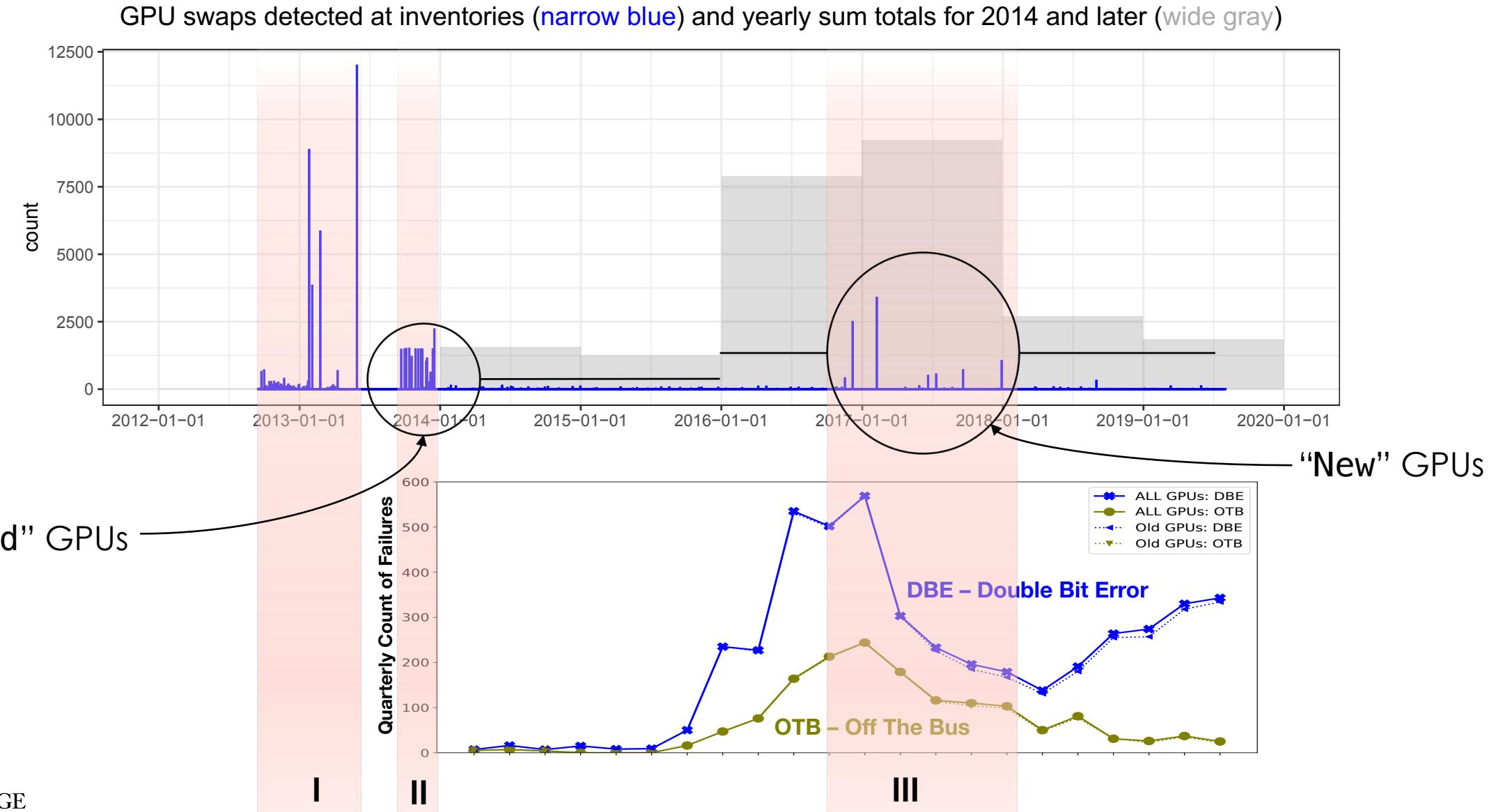
NVIDIA SXM – Location of a non-ASR



Silver-sulfide corrosion
"Flowers-of-Sulfur"

ASR = Anti-Sulfur Resistor

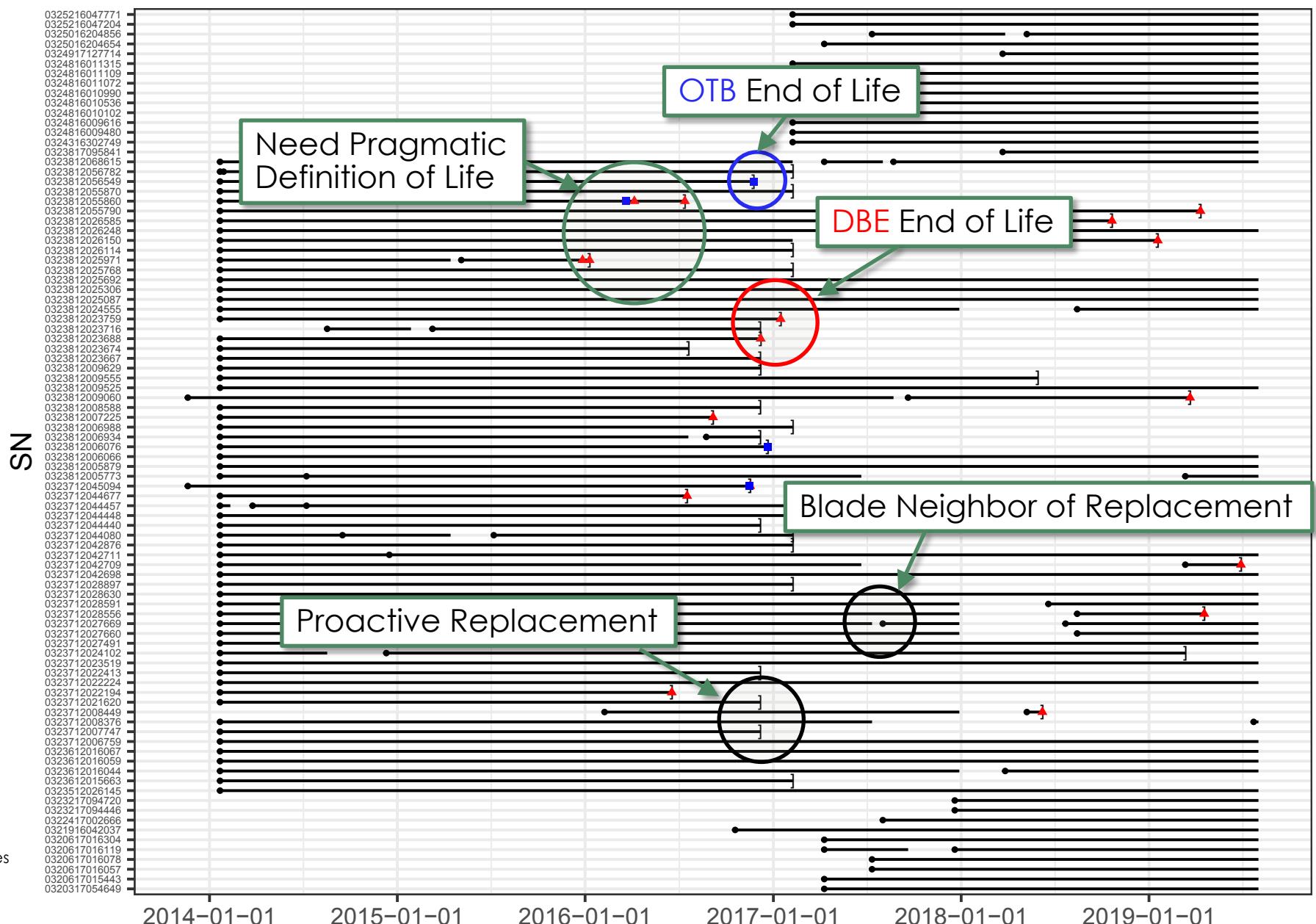
GPU Failures and Replacements in ORNL's Titan



GPU Life Visualization: Serial Number View

Critical for:

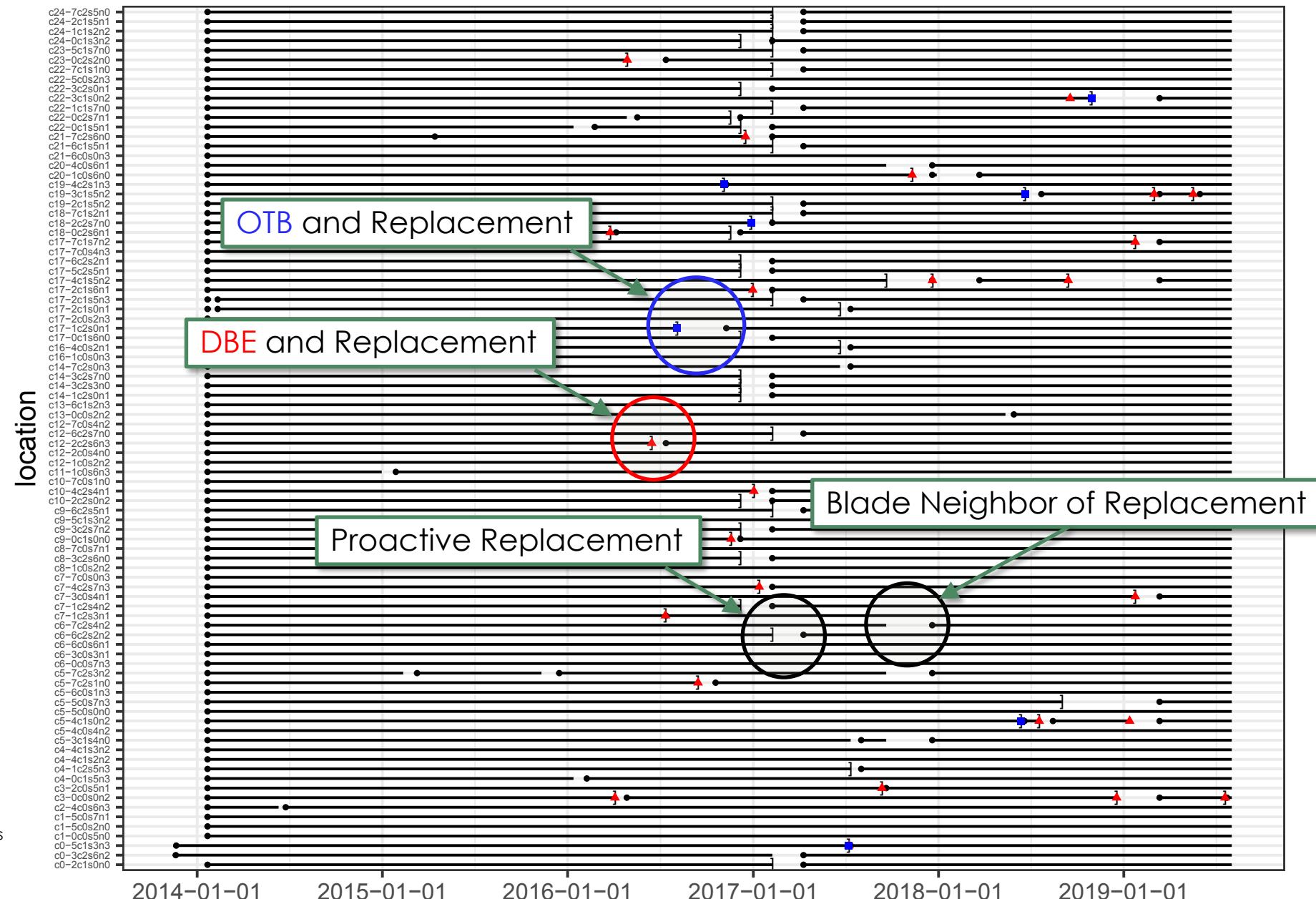
- Understanding data
- Defining GPU Life
- Data processing verification



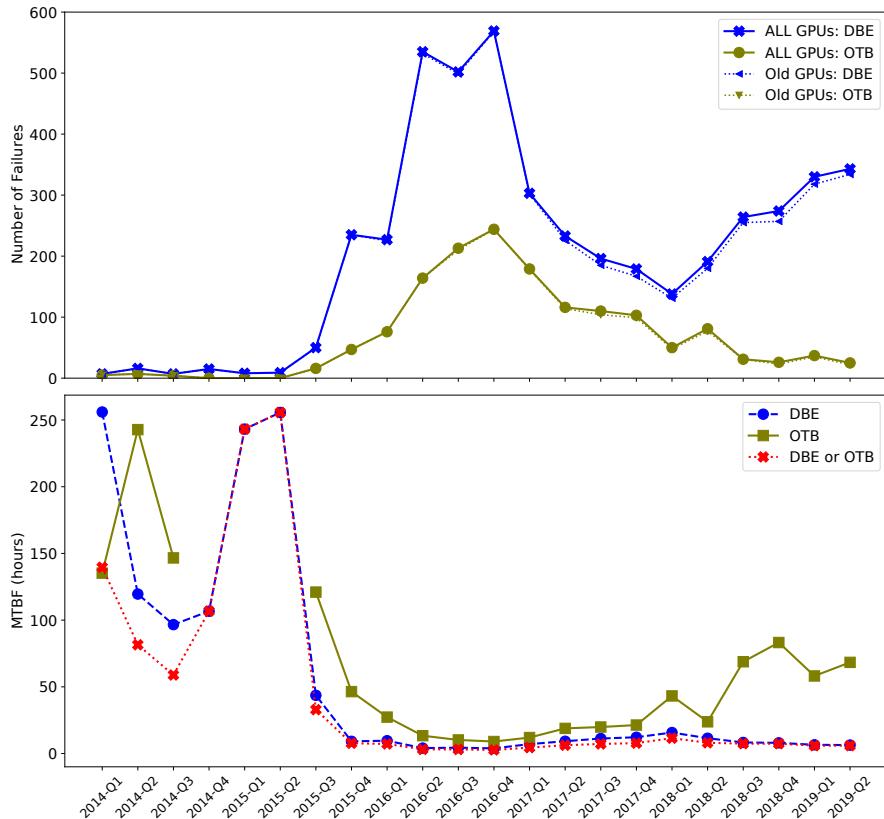
GPU Life Visualization: Location View

Critical for:

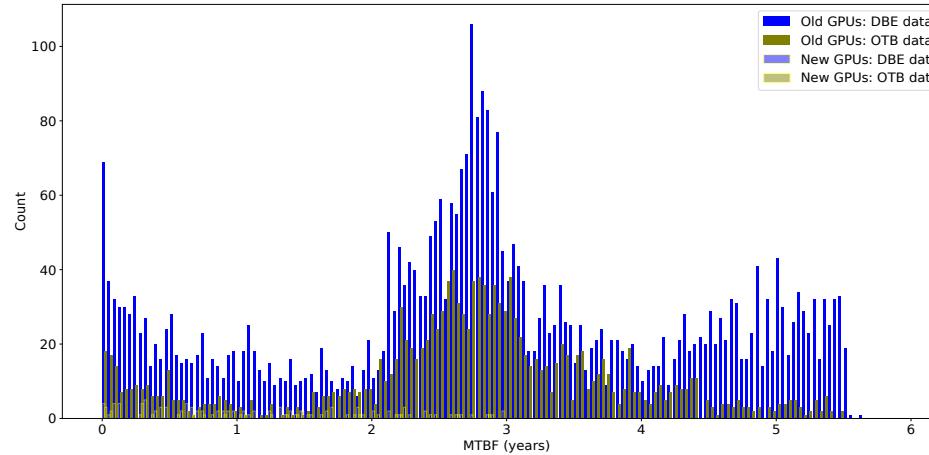
- Understanding data
- Defining GPU Life
- Data processing verification



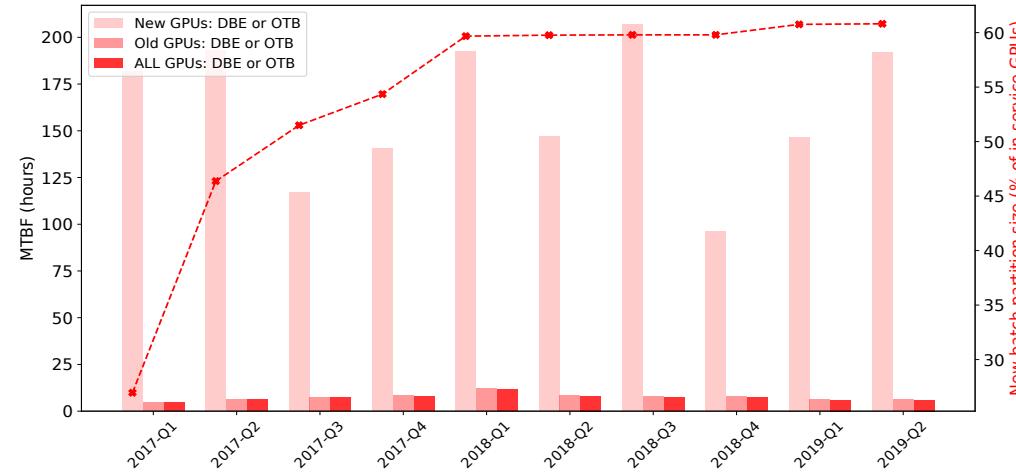
Traditional Reliability in HPC is Focused on MTBF



System-wide Reliability: Quarterly number of failures (top) and MTBF (bottom).



Individual GPU Reliability: MTBF histogram for units that had at least one failure.
Interpret carefully: lacks information from units with no failures!



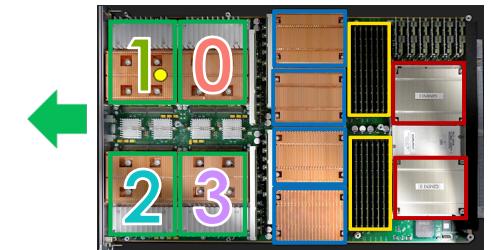
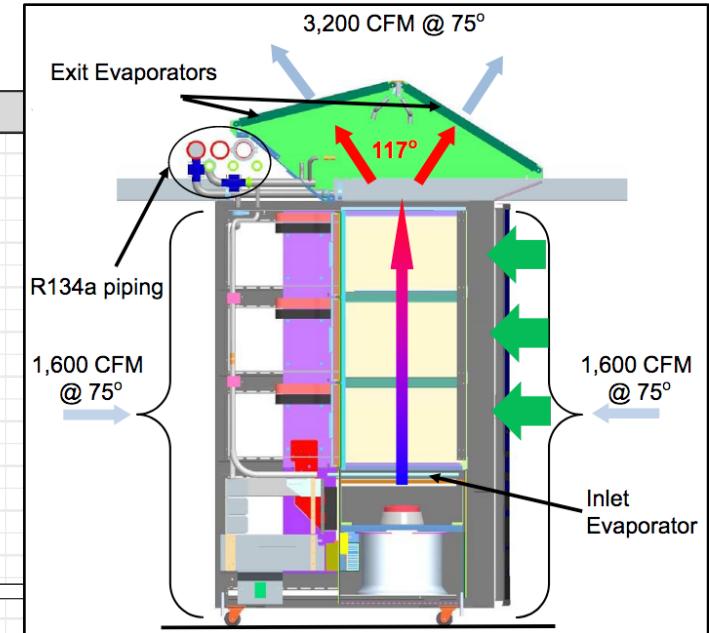
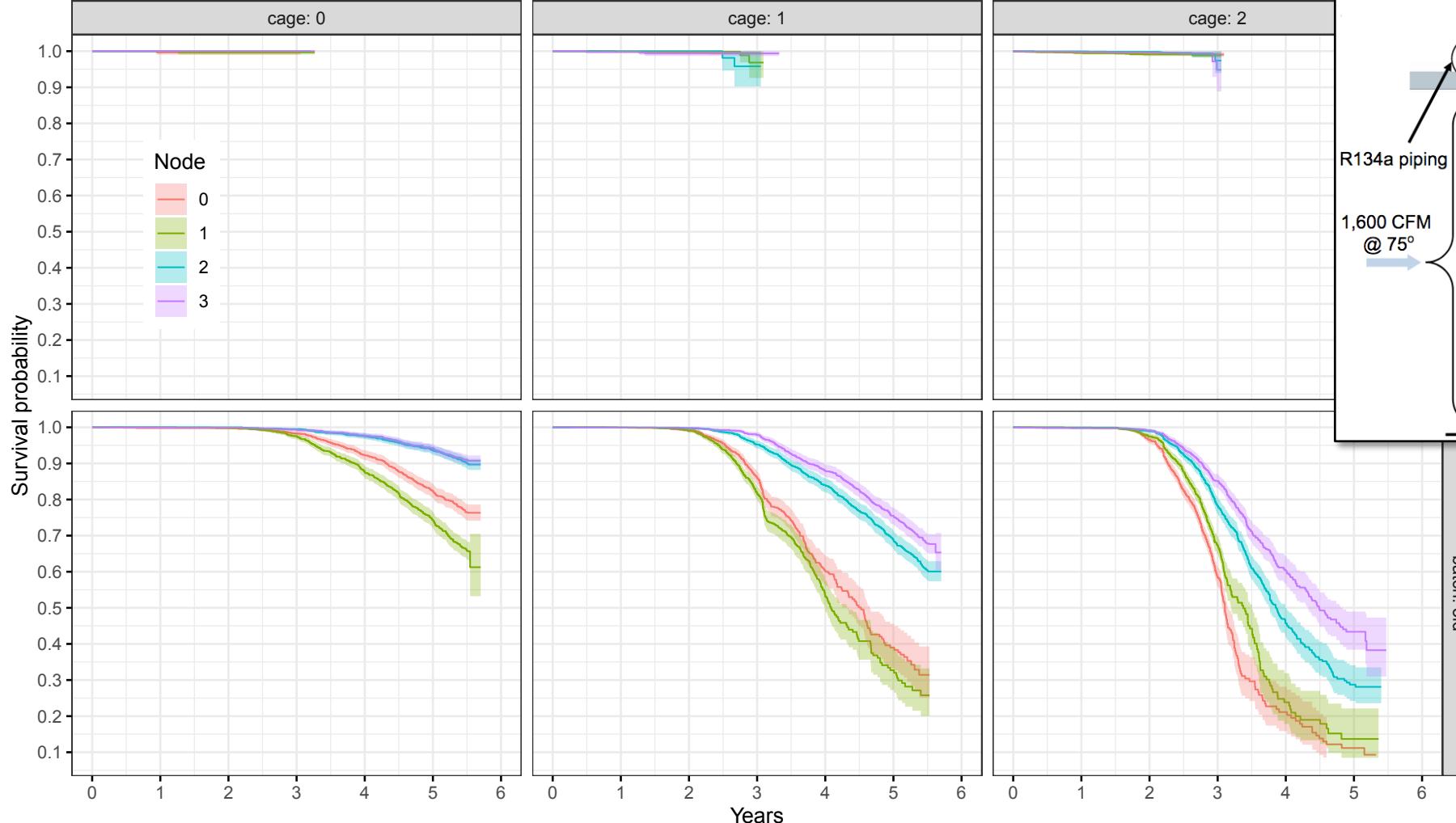
Old-New as Two Partitions: MTBF differs by 12x factor!

Kaplan-Meyer Survival Analysis

- Commonly used in Biostatistics and Biomedical research*
- Nonparametric
 - If T is failure time and $F(t) = \Pr\{T < t\}$ is the cumulative failure distribution function
 - Then the survival probability, $S(t) = \Pr\{T \geq t\} = 1 - F(t)$, is its complement
 - Recursive computation $S(t_2) = \Pr\{\text{survive from } t_1 \text{ to } t_2\} S(t_1)$
- Able to incorporate censoring
- Split population into groups
- Available uncertainty estimate

*E. L. Kaplan and P. Meier, "Nonparametric estimation from incomplete observations," *Journal of the American Statistical Association*, vol. 53, no. 282, pp. 457–481, 1958.

Cage and Node Effect Explainable by Airflow in Cabinet

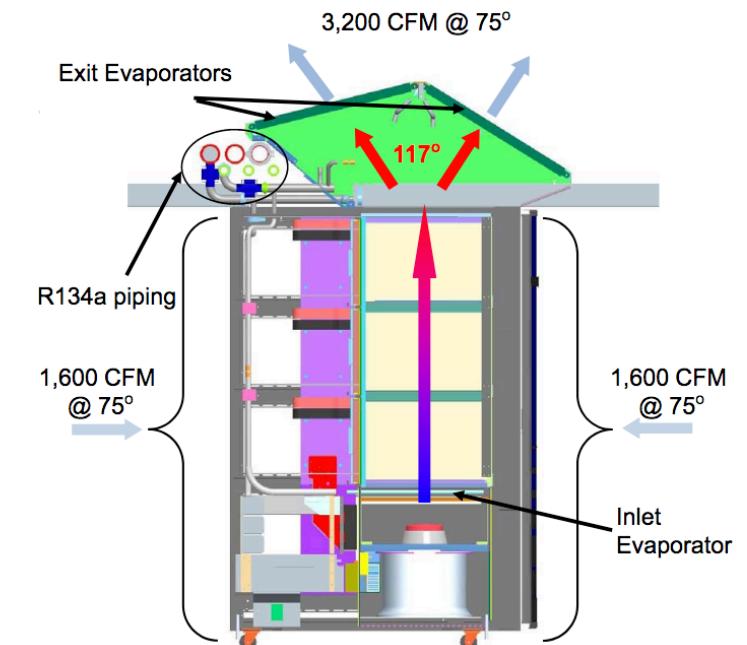
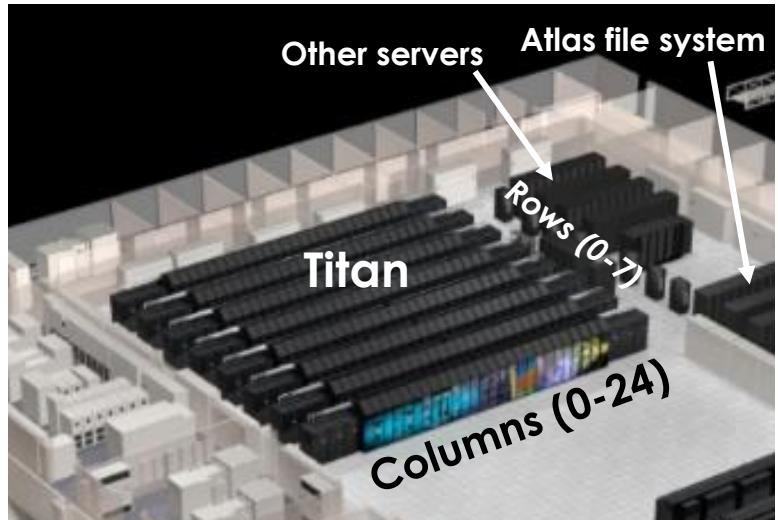
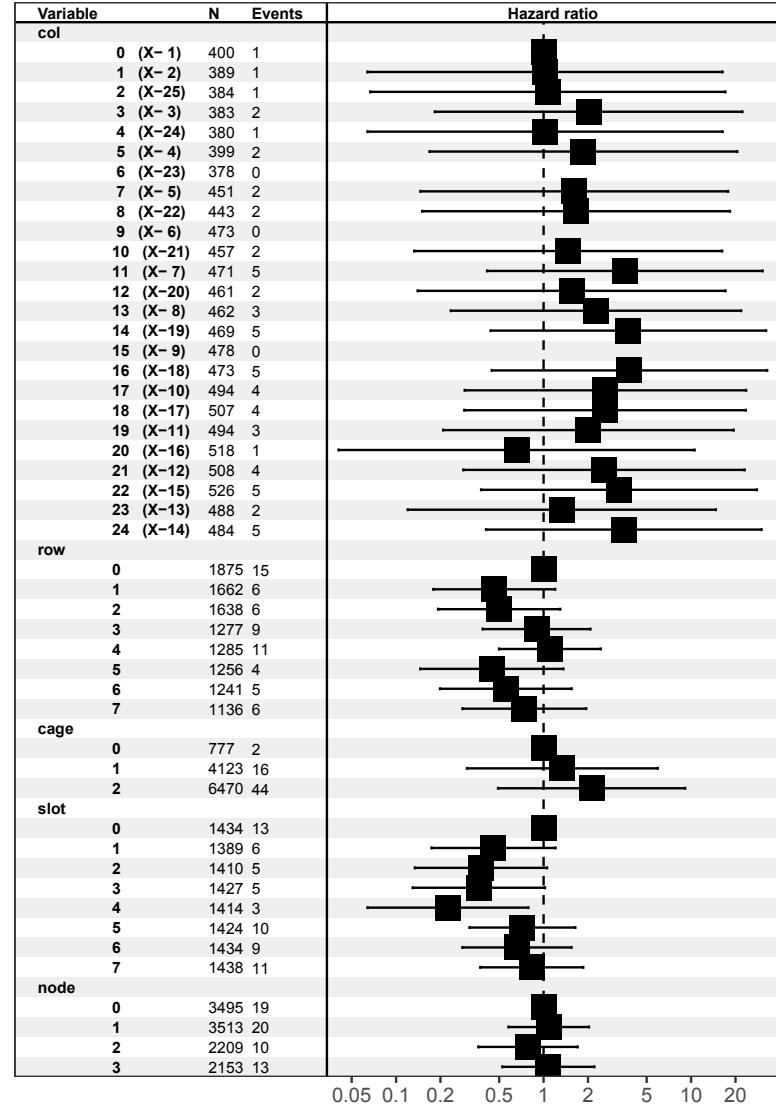
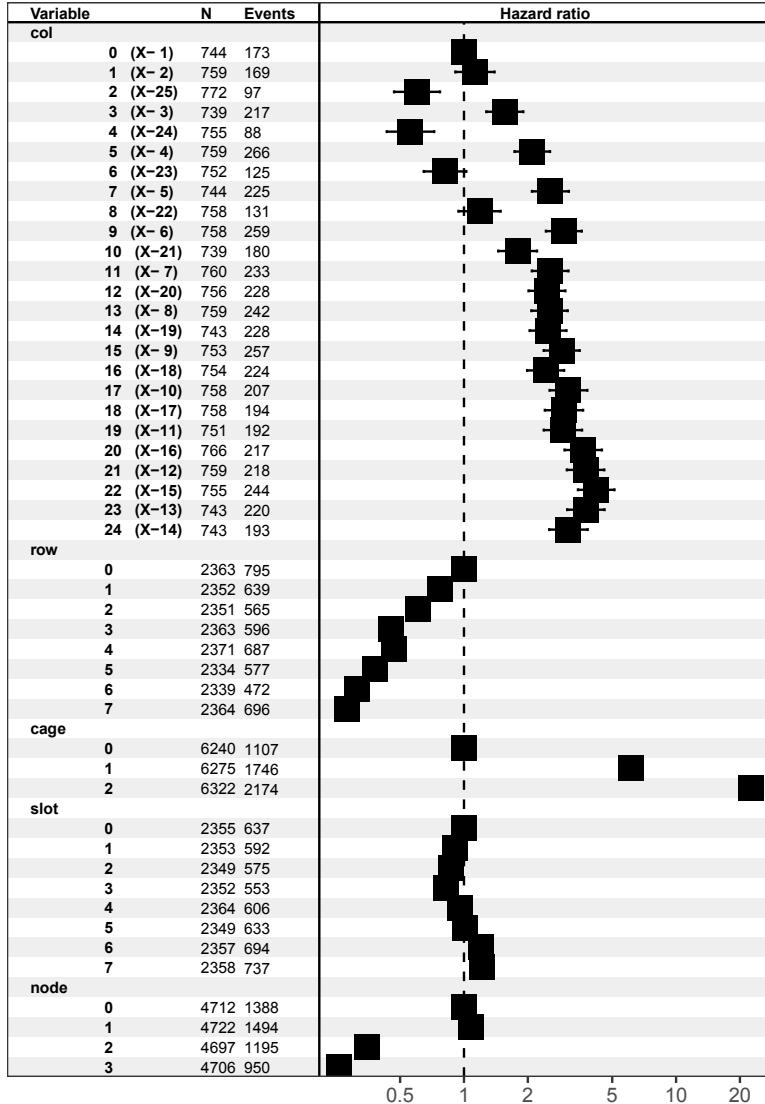


Cox Proportional Hazards Regression Model

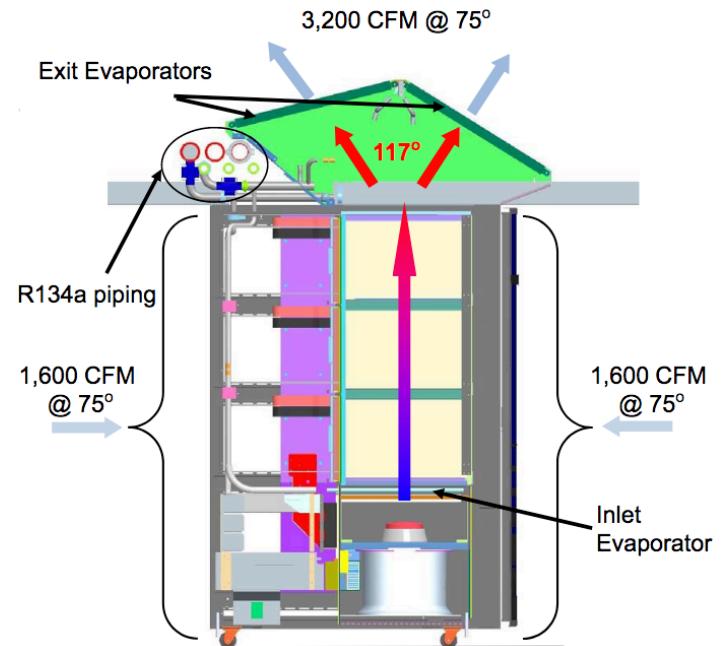
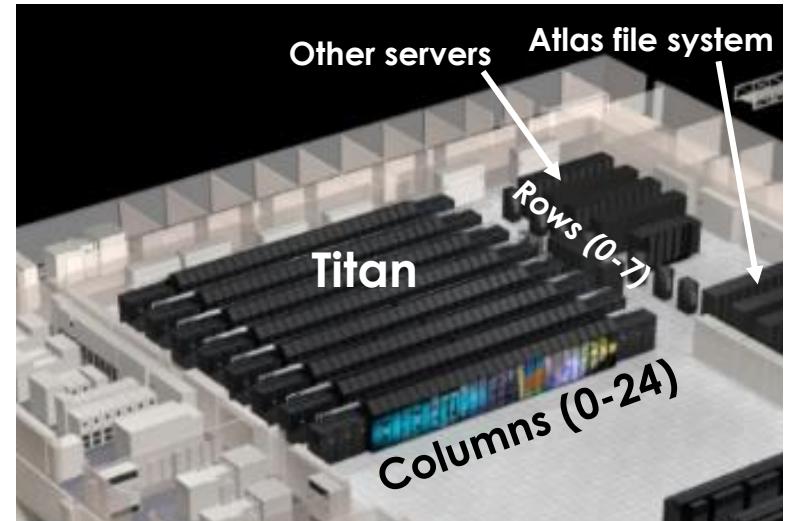
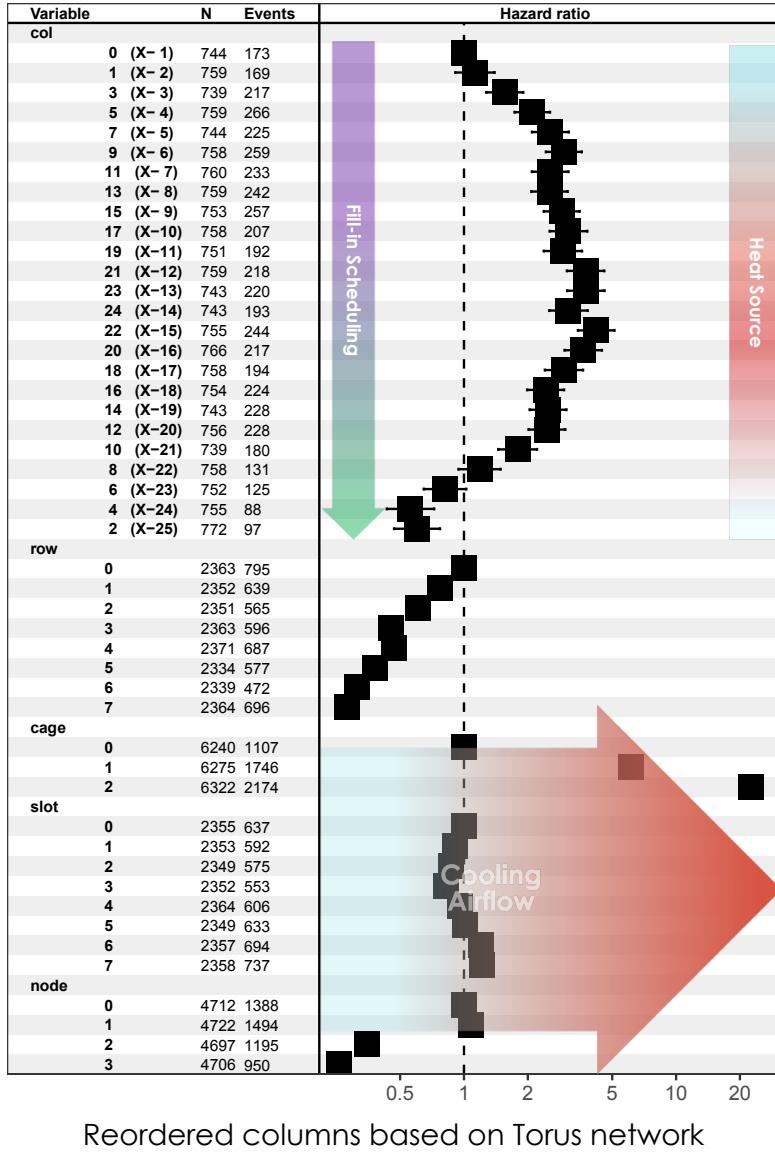
- Commonly used in Biostatistics and Biomedical research*
- Able to adjust for covariate effects
- Each GPU is like a patient, affected by its location (treatment)
- The hazard for patient k is $H_k(t) = H_0(t)e^{\sum_1^n \beta_i x_i}$
 - Base hazard rate, $H_0(t)$, multiplied by a function of covariates (hazard coefficient)
- Semiparametric model
 - Baseline hazard is nonparametric (no functional shape assumption)
 - Hazard coefficient is a parametric function of covariates
- Assumes hazards are proportional

*D. R. Cox, "Regression models and life-tables," *Journal of the Royal Statistical Society. Series B (Methodological)*, vol. 34, no. 2, pp. 187– 220, 1972.
We use R packages `survival` and `survminer`.

Strong Signal in old Batch, Pattern Similar to K-M Analysis



Strong Signal in old Batch, Pattern Similar to K-M Analysis



Future Research and Development Needs

- We need to design the HPC hardware/software ecosystem to be able to deal with high error and failure rates, expected and unexpected!
 - Resilience research and development is, in part, risk mitigation against the unexpected
 - There is always a cost/benefit trade-off that needs to be considered
 - Resilience mitigation mechanisms should be a toolbox with lots of options
- Resilience should be by design and not as an afterthought
 - Resilience is a crosscutting issue that should be considered everywhere (and not only in architecture)
 - After 25 years, MPI is still not fault tolerant, while PVM was fault tolerant 28 years ago in 1993

Questions?