

GPU Lifetimes on Titan Supercomputer: Survival Analysis and Reliability

George Ostrouchov*, Don Maxwell ⁺,
Rizwan A. Ashraf ^{*}, Christian Engelmann ^{*},
Mallikarjun Shankar ⁺, James H. Rogers ⁺

^{*} Computer Science and Mathematics Division

⁺ National Center for Computational Sciences Division

Oak Ridge National Laboratory

Effects of Architecture and Scheduling on GPU Reliability

The Science

- Over 100,000 collective years of GPU operation
- Advanced statistical methods typical in biomedicine
- Data and analysis codes made publicly available for reproducibility

Findings

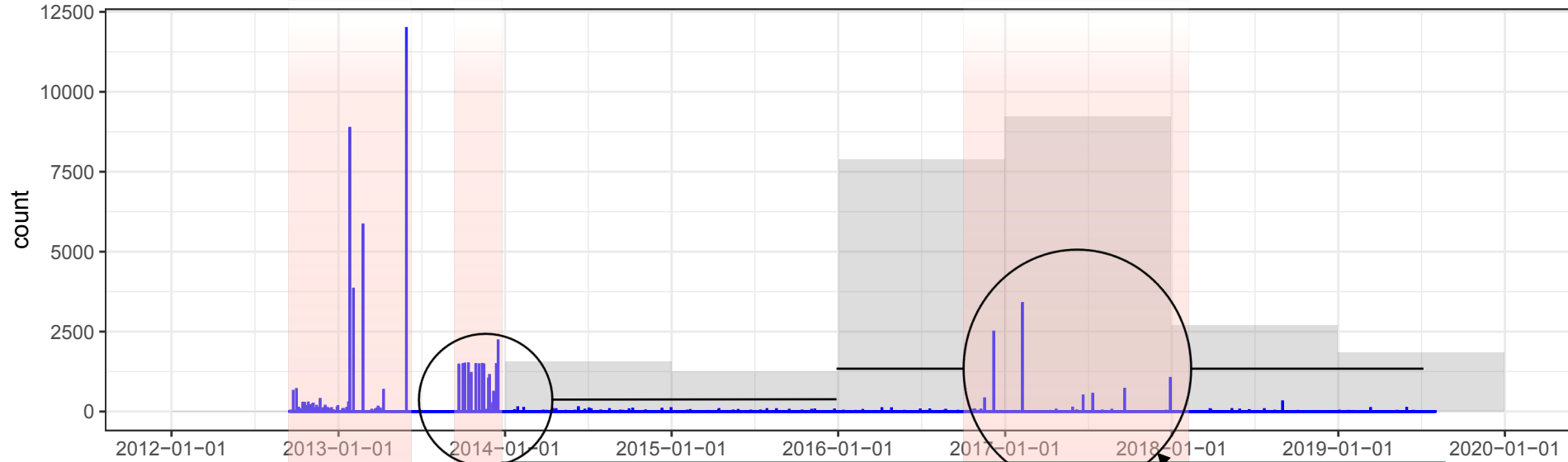
- Failure mode provided strong reliability signal
- Correlations explainable by cooling air transport
- Correlations explainable by job scheduling

Overview

- GPU-focused history of Titan
- Unexpected mid-life failures
- Titan GPU architecture
- Data collection and curation
- Traditional MTBF analysis
- Statistical survival analysis
- Conclusions

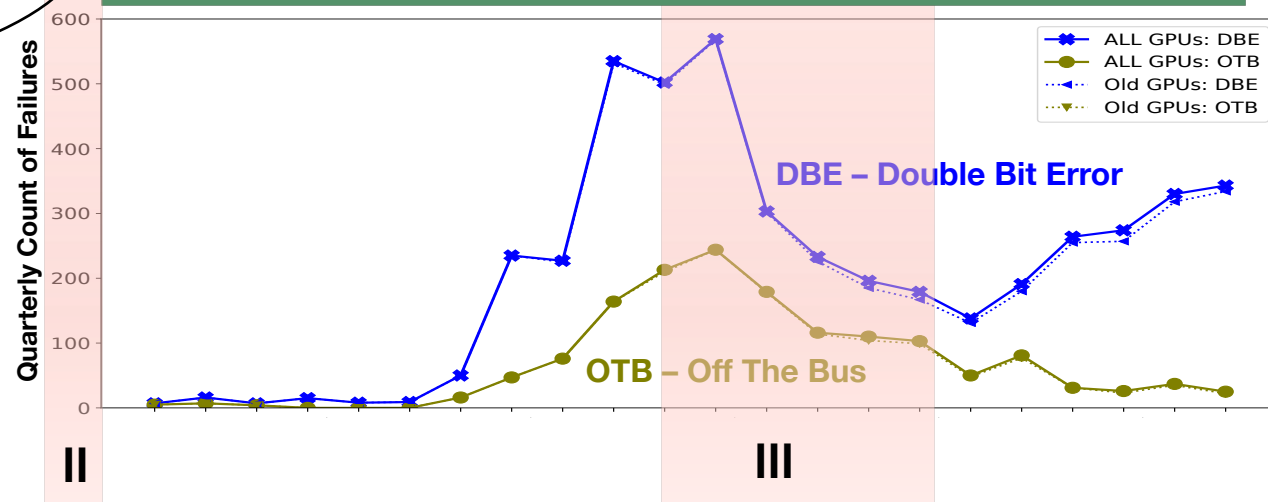
Three Rework Cycles and Years of Stable Operation

GPU swaps detected at inventories (narrow blue) and yearly sum totals for 2014 and later (wide gray)

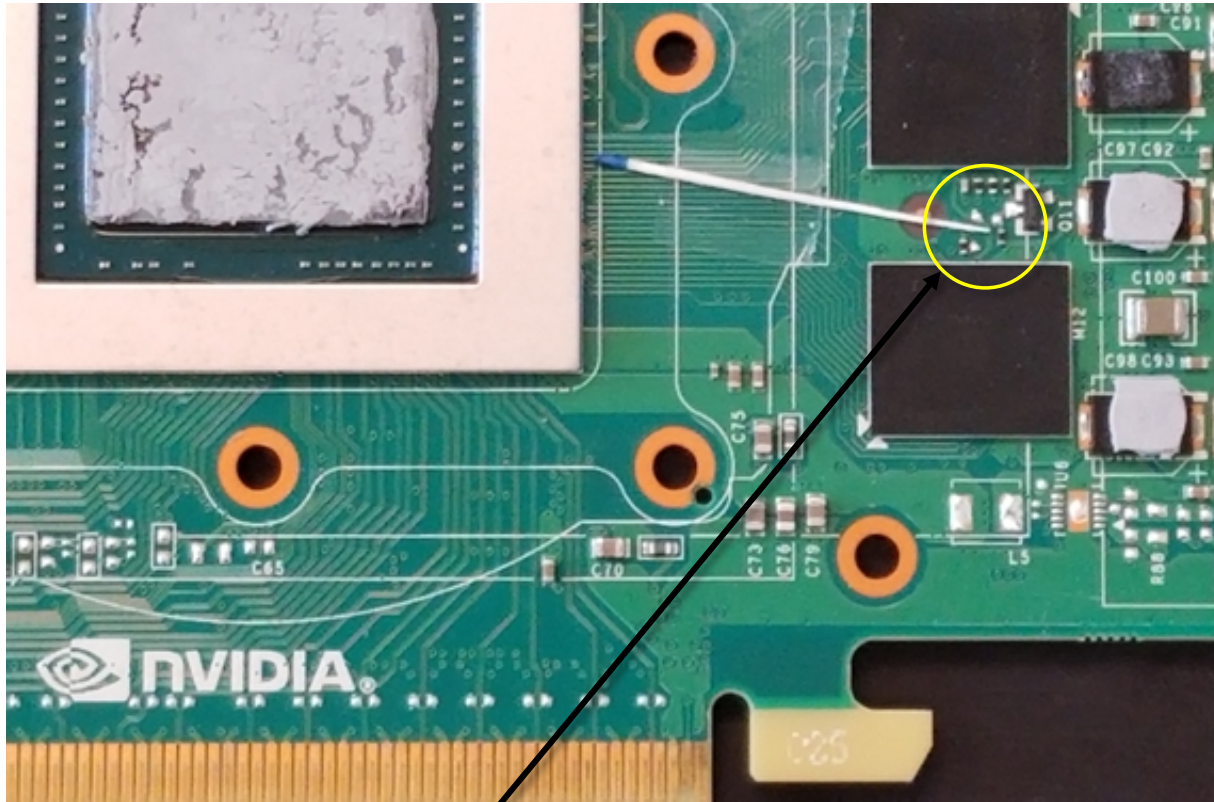


“Old” GPUs

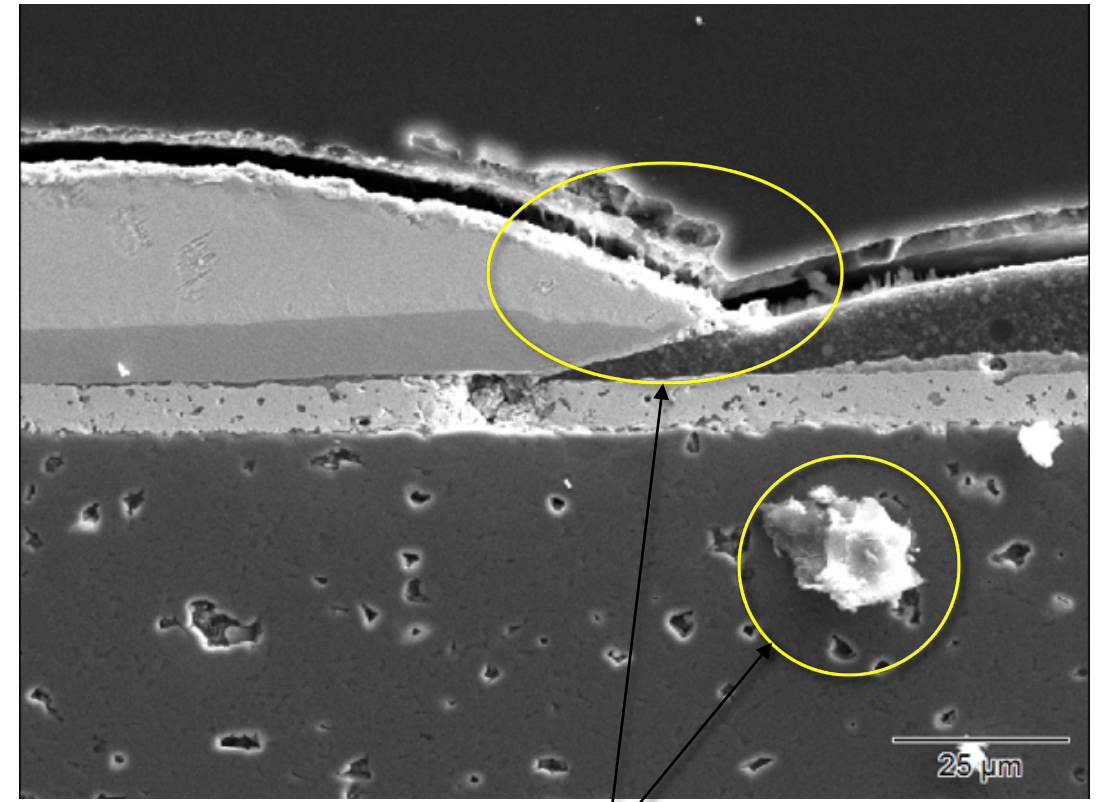
“New” GPUs



Root Cause: Non-ASR Components on SXM GPU



NVIDIA SXM – Location of a non-ASR

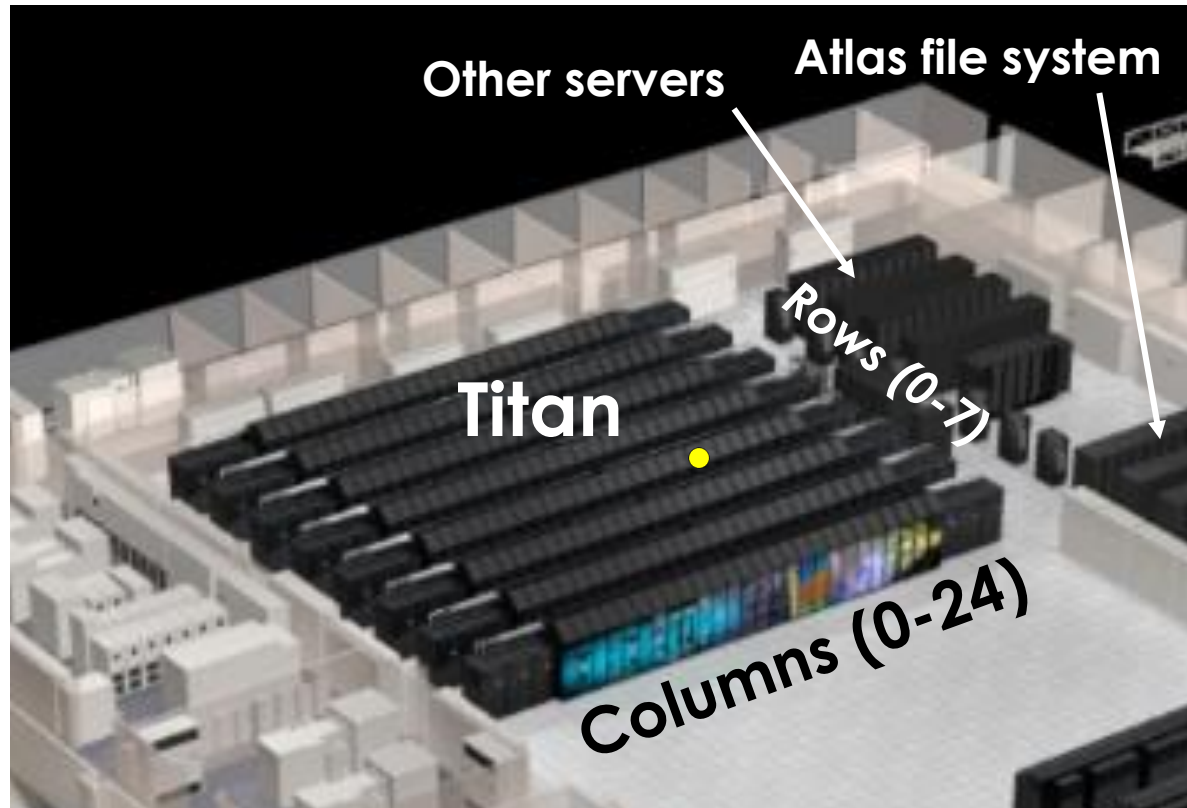


*Silver-sulfide corrosion
"Flowers-of-Sulfur"*

ASR = Anti-Sulfur Resistor

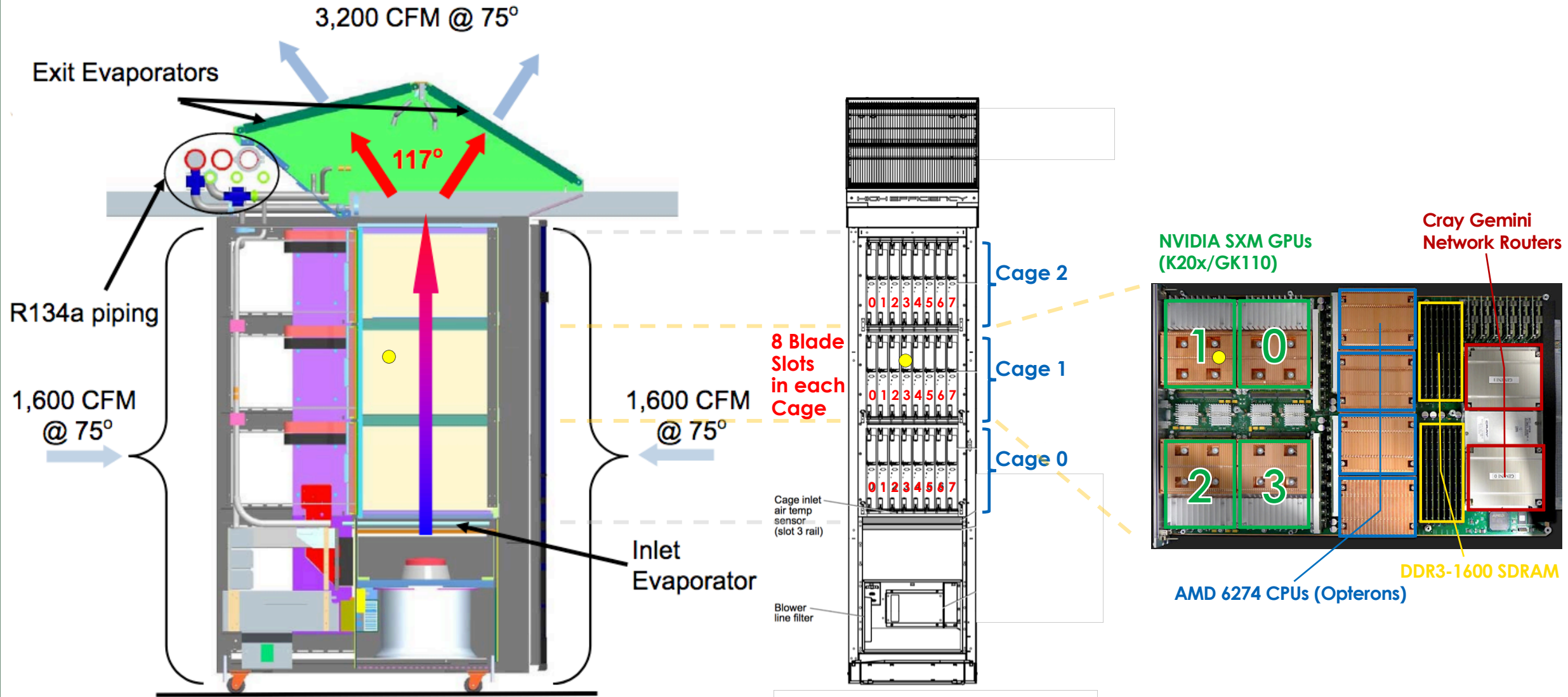
Machine Room Layout: GPU Locators

c###-#c#s#n#
a a l o
b g o d
i e t e
n
e
t



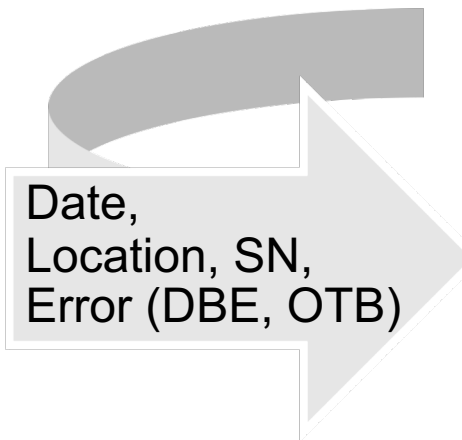
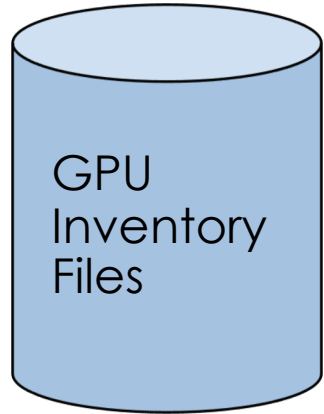
- **c17-4c1s3n1** = cabinet in column **17** – row **4**

Cabinet Mechanical Packaging: Locating a GPU



• **c17-4c1s3n1** = cage 1, slot 3, node 1

GPU Life Data Built Incrementally from Two Sources



DOI:10.13139/ORNLNCCS/1657202

GPU Life Data

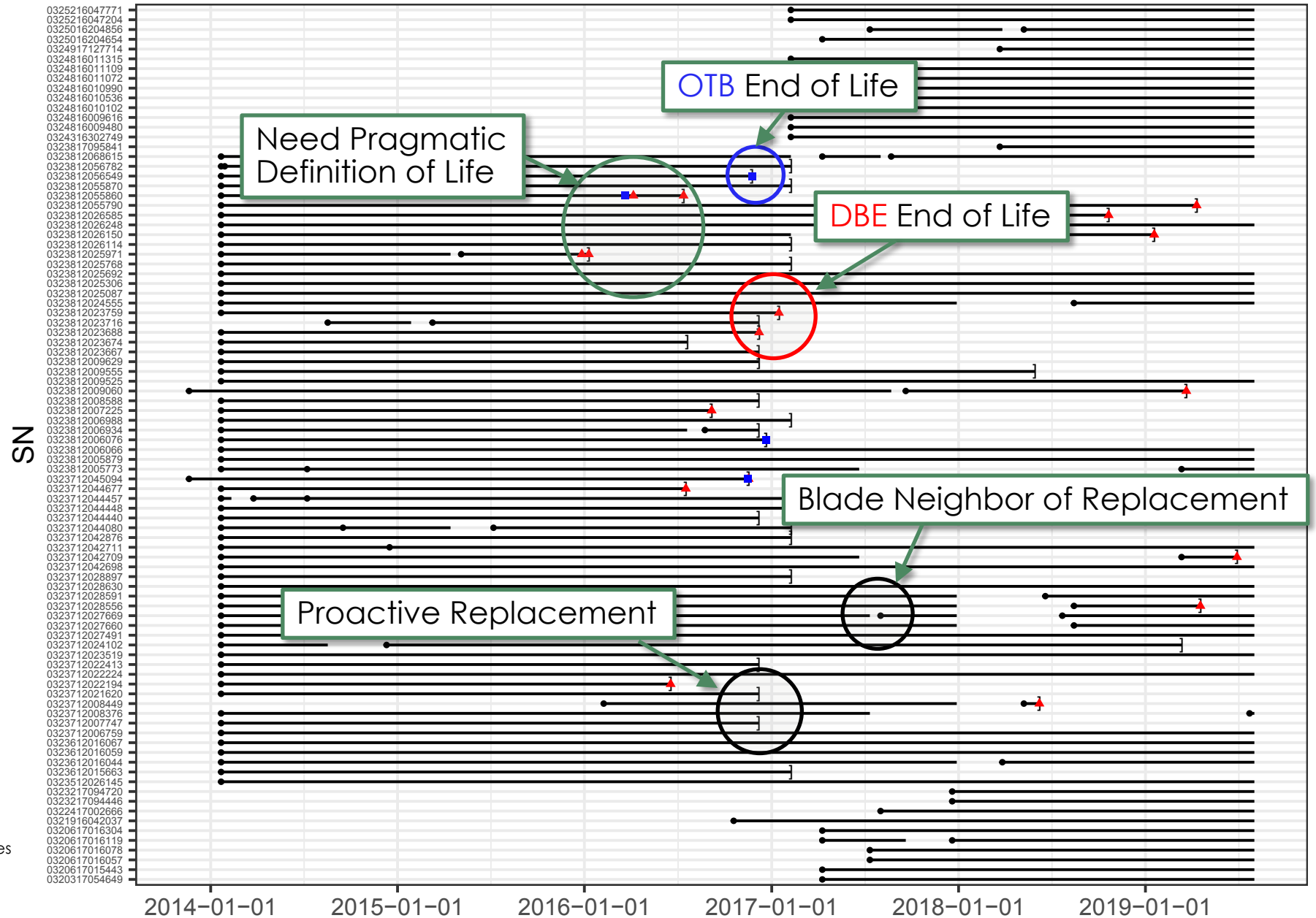
SN	Location	Date In	Date Out
0323812007945	c17-4c1s3n1	09/28/2012 10:29:48	02/02/2013 11:32:29
c20-6c1s5n2	07/23/2019 11:25:33	01/20/2020 18:51:10	
c13-1c1s3n3	01/21/2014 10:28:50	07/11/2017 18:04:25	
c0-1c1s3n3	10/11/2013 15:57:33	10/12/2013 22:09:31	
c21-1c2s5n0	03/19/2013 15:48:11	05/29/2013 11:54:11	
0325216047736	c18-4c1s5n1	04/09/2017 21:36:19	01/20/2020 18:51:10
0323812008856	c5-4c0s7n0	09/30/2012 12:20:00	01/25/2013 15:29:58
c0-6c1s7n2	10/21/2013 14:28:19	10/28/2013 17:52:44	
c3-3c1s5n0	05/29/2013 11:54:11	05/29/2013 11:54:11	
c23-6c1s7n2	01/21/2014 10:28:50	11/02/2018 14:42:34	
DBE	11/02/2018 14:42:34		
.			
.			
.			

Failure type recorded as "Date In"

GPU Life Visualization: Serial Number View

Critical for:

- Understanding data
- Defining GPU Life
- Data processing verification

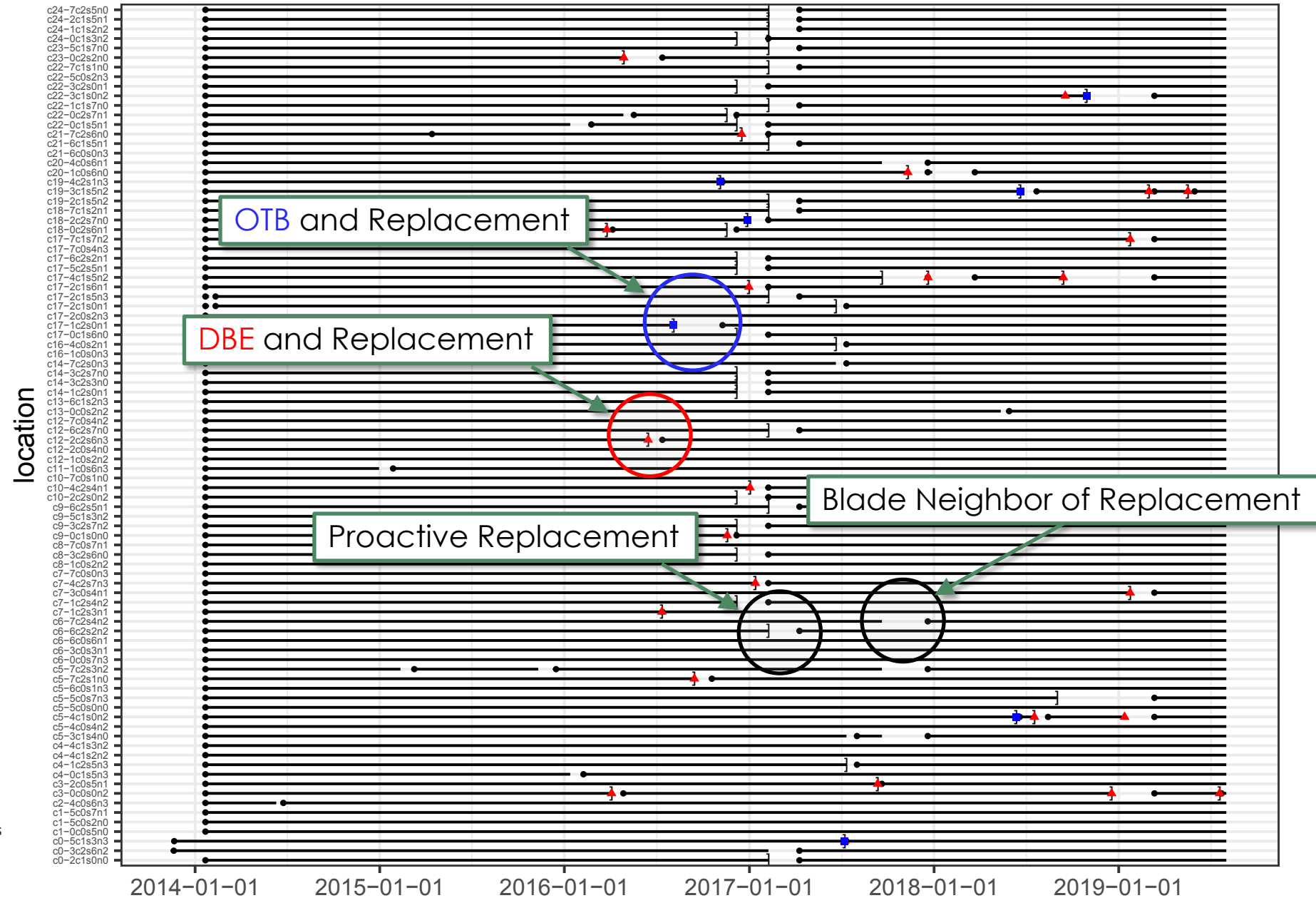


Produced in R via ggpplot2 and lubridate packages

GPU Life Visualization: Location View

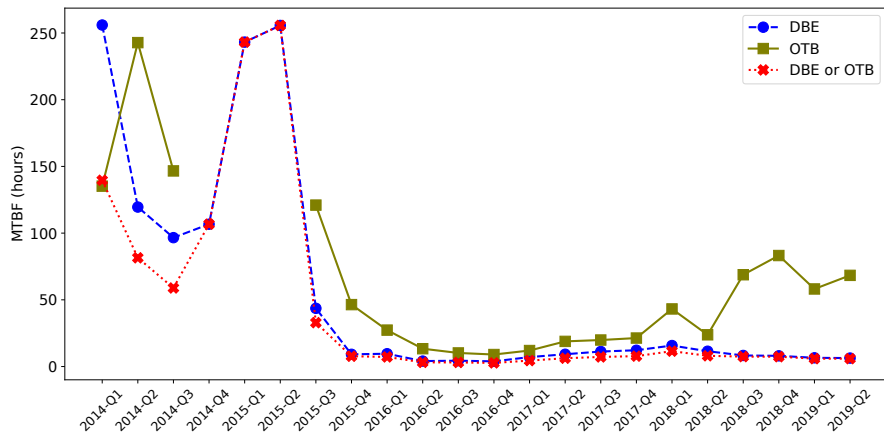
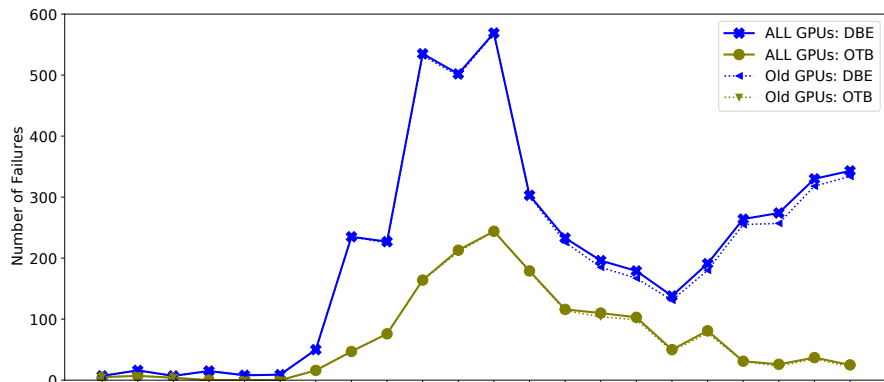
Critical for:

- Understanding data
- Defining GPU Life
- Data processing verification

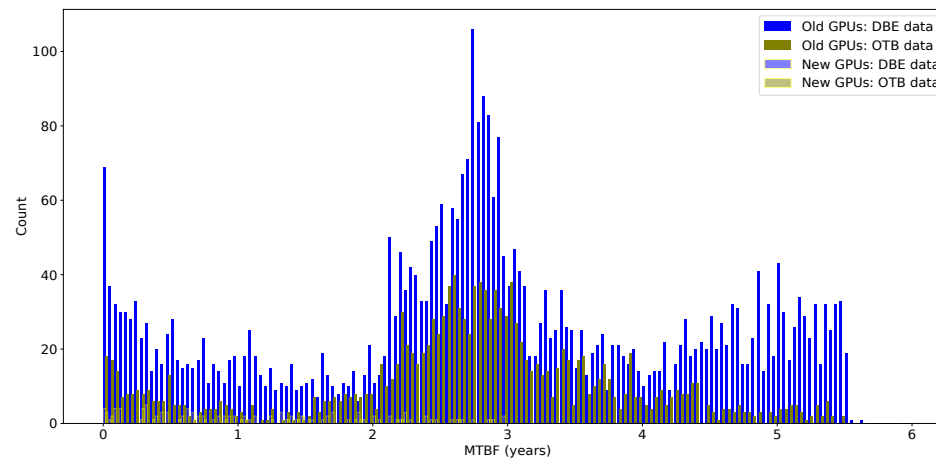


Produced in R via ggpplot2 and lubridate packages

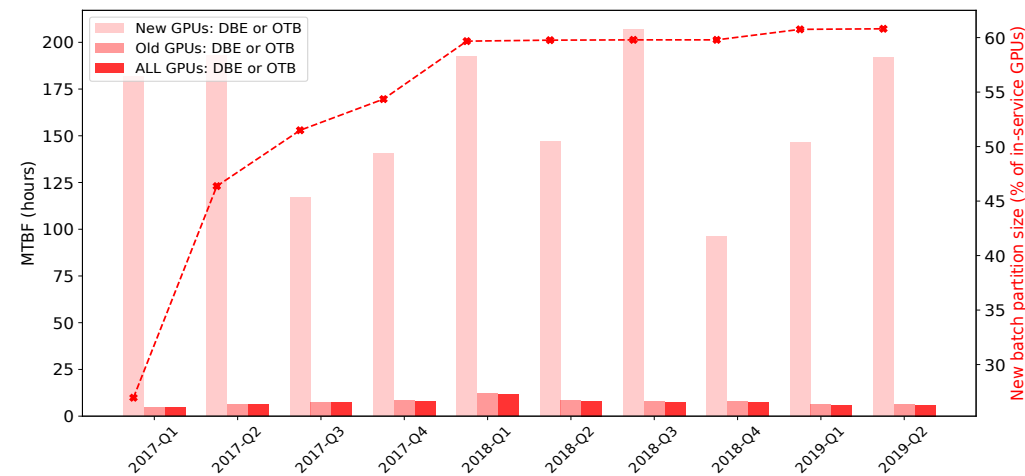
Traditional Reliability in HPC is Focused on MTBF



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!

Survival Analysis Requires a Definition of GPU Lifetime

Lifetime with observed failure:

- Unit experiences at least one DBE or OTB
- Unit is removed from the system
- Lifetime is the cumulative operation time over all its locations

Censored lifetime:

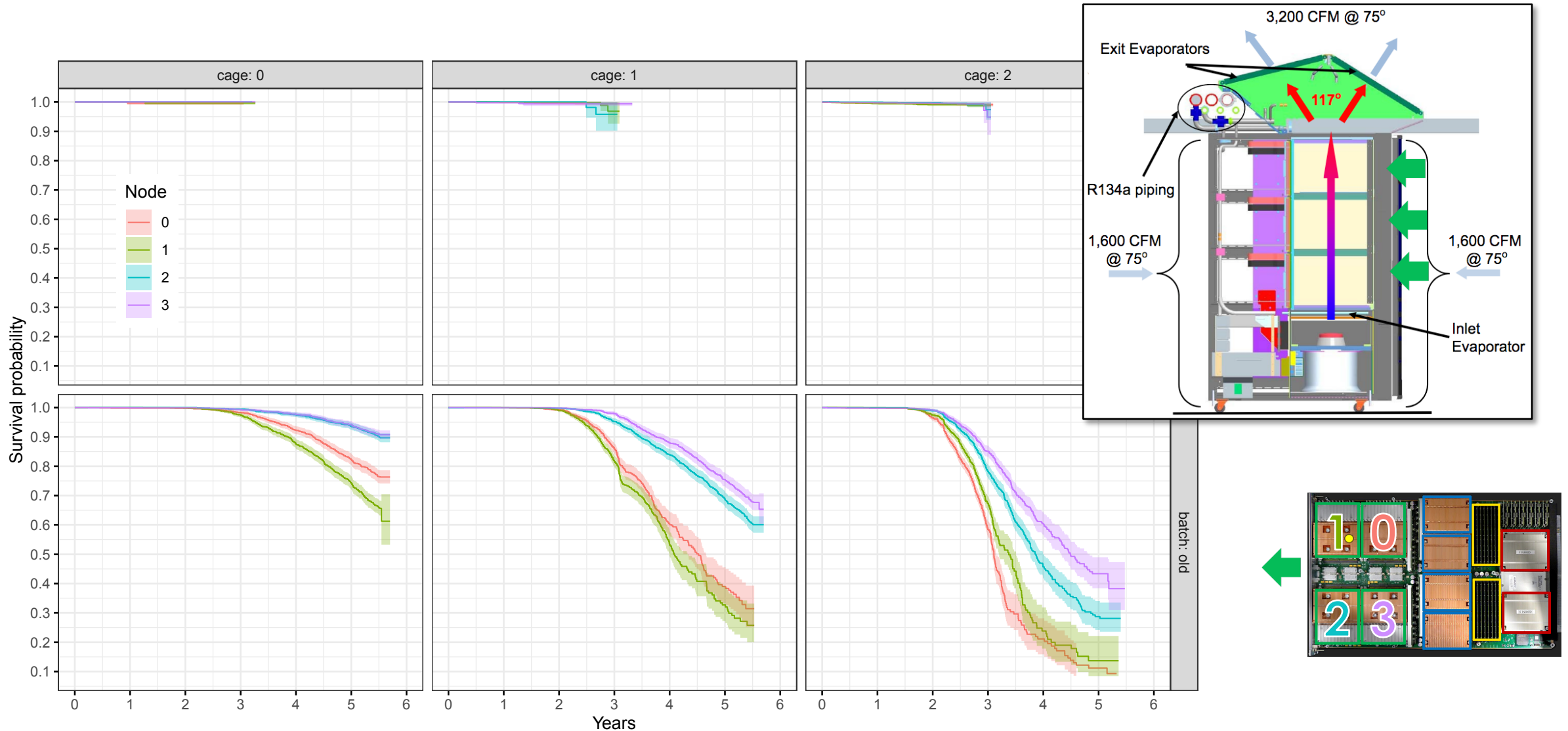
- Unit is removed or is present when system turned off
- If removed, no DBE or OTB events
- Lifetime is the cumulative operation time over all its locations
- Ignore censoring due to observation only at inventories

Kaplan-Meier 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

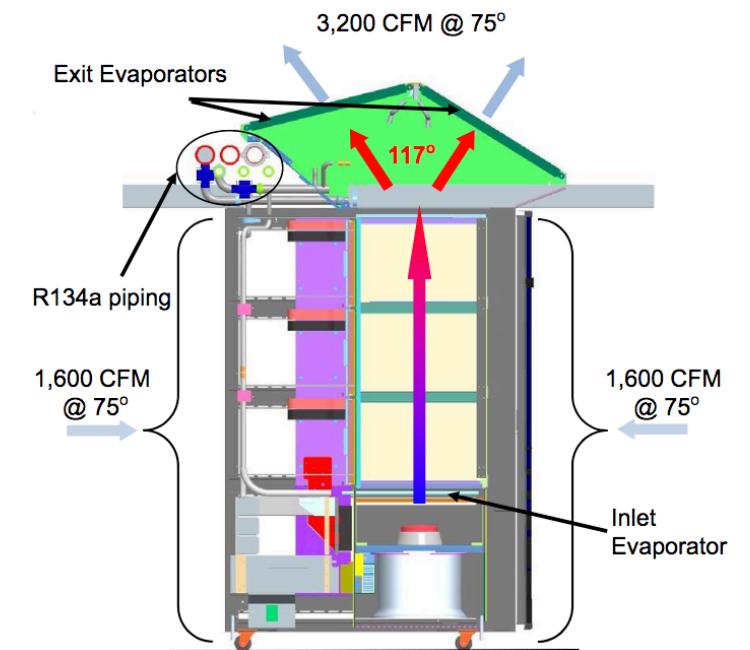
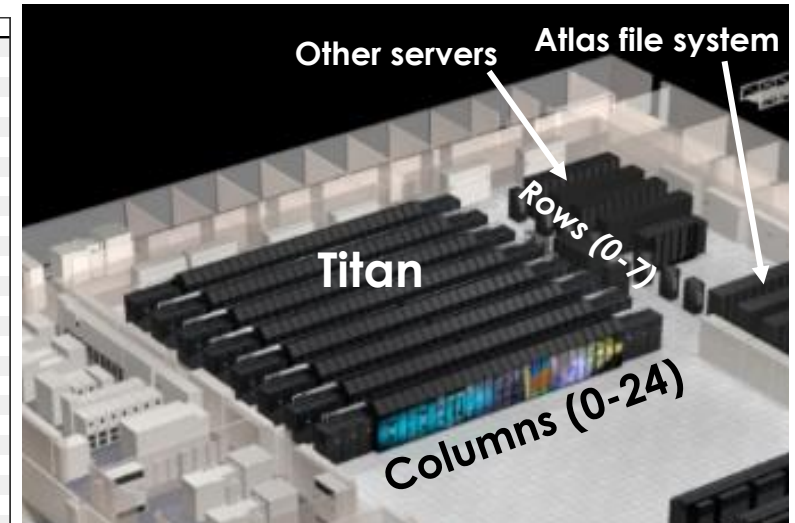
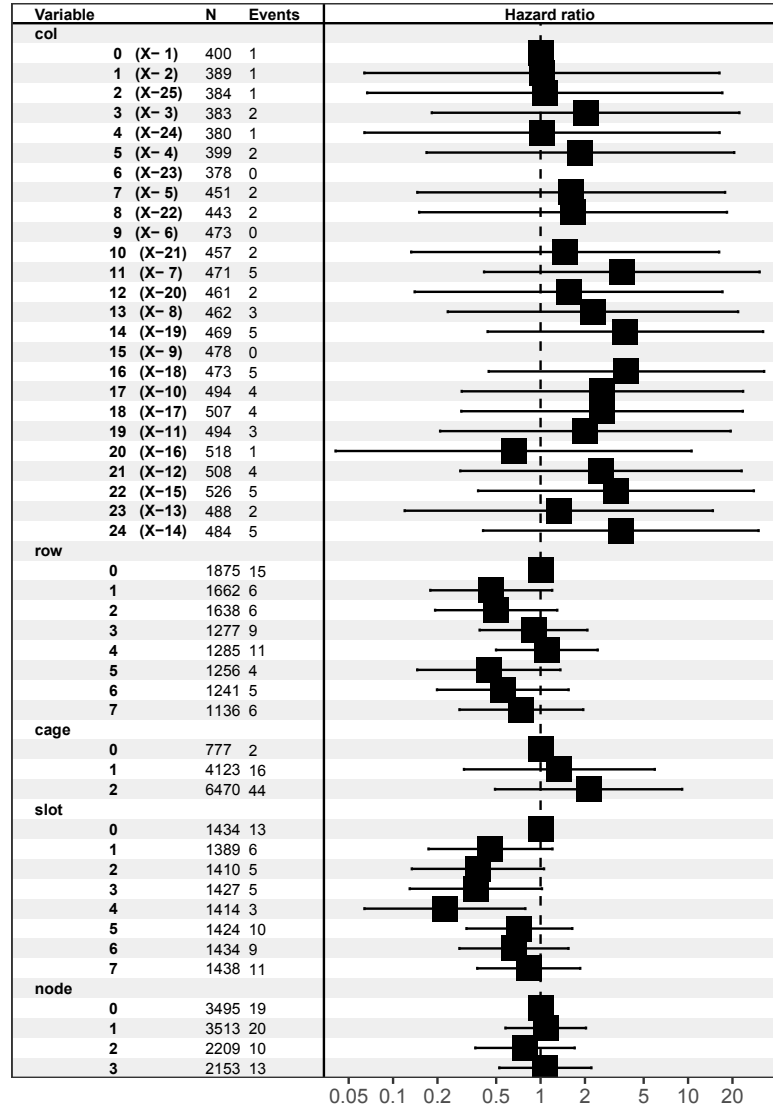
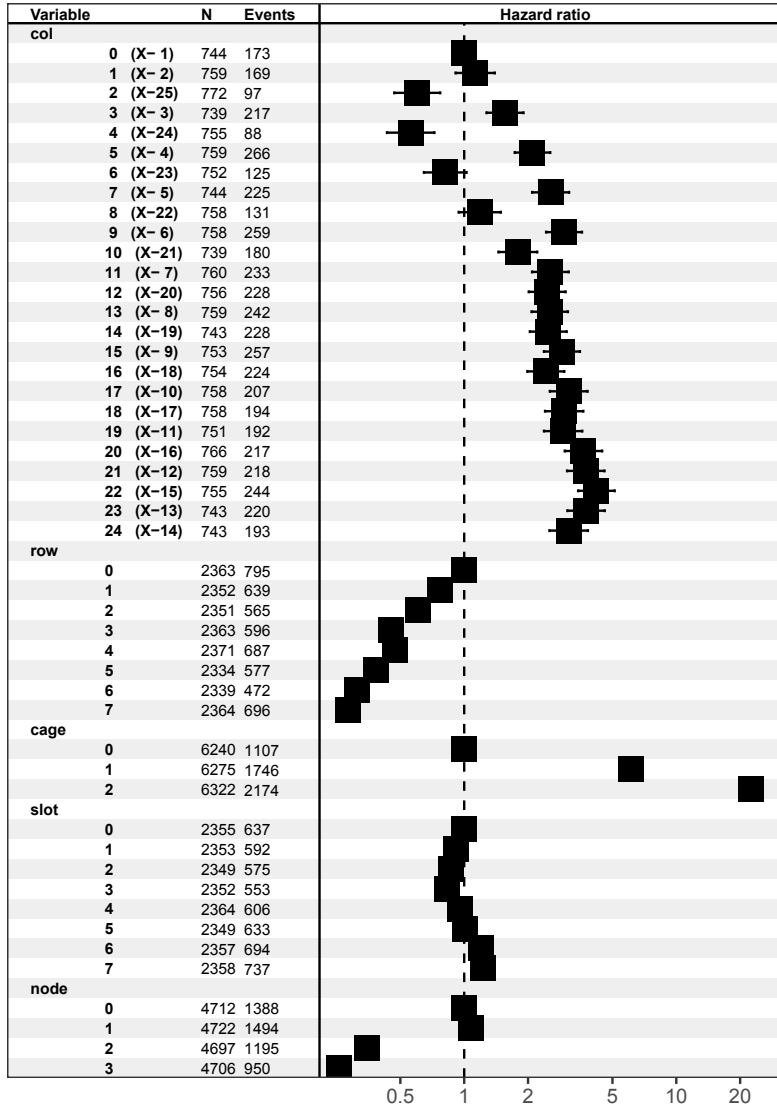


Cox Proportional Hazards Regression Model

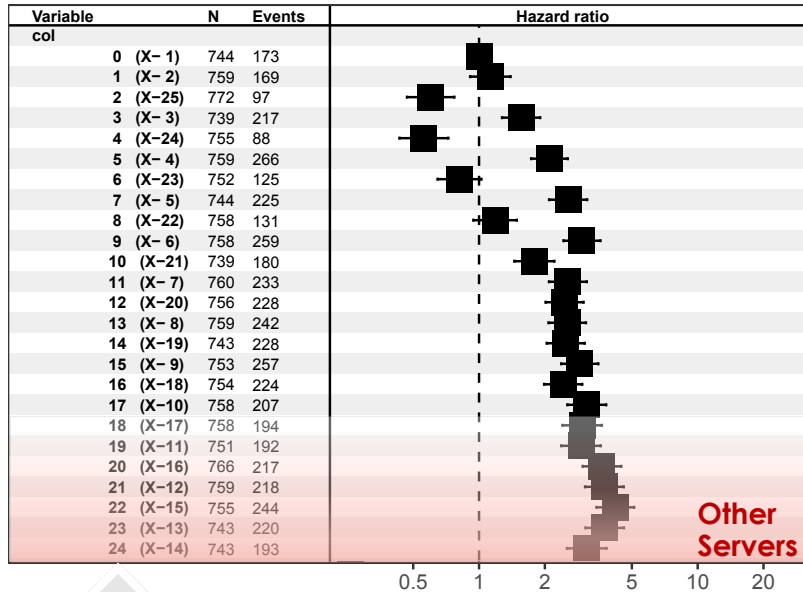
- Commonly used in Biostatistics and Biomedical research*
- Able to adjust for covariate effects
- Each GPU is like a patient, affected by its location (treatment)
- The hazard for patient k is $H_k(t) = H_0(t)e^{\sum_1^n \beta_i x_i}$
 - Base hazard rate, $H_0(t)$, multiplied by a function of covariates (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`.

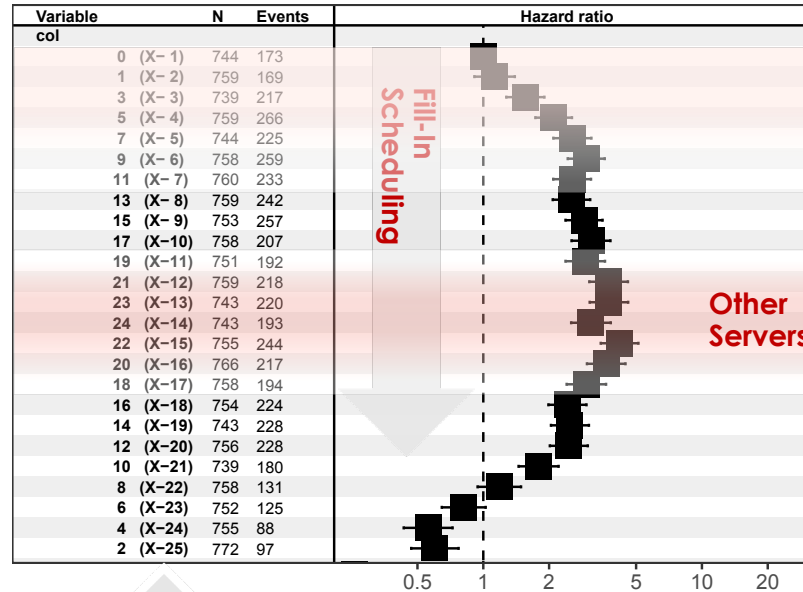
Strong Signal in old Batch, Pattern Similar to K-M Analysis



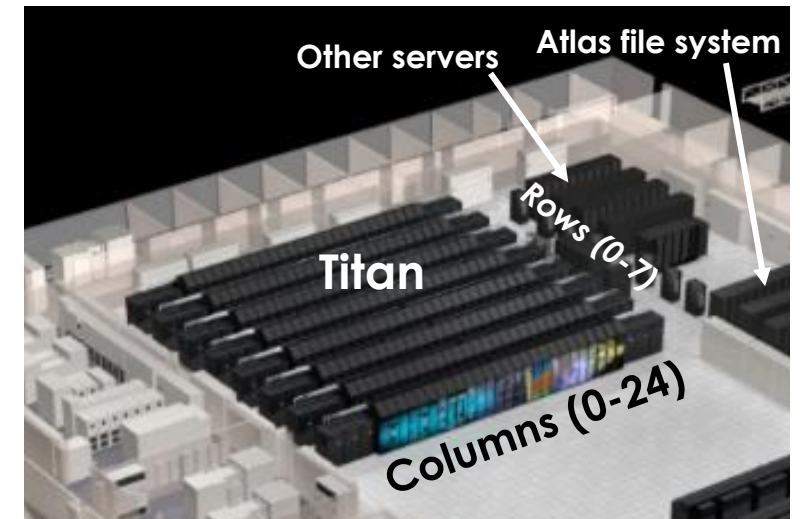
Fill-in Scheduling Effect Explainable via Torus Coordinate



↑
Column Order



↑
Torus Coordinate Order



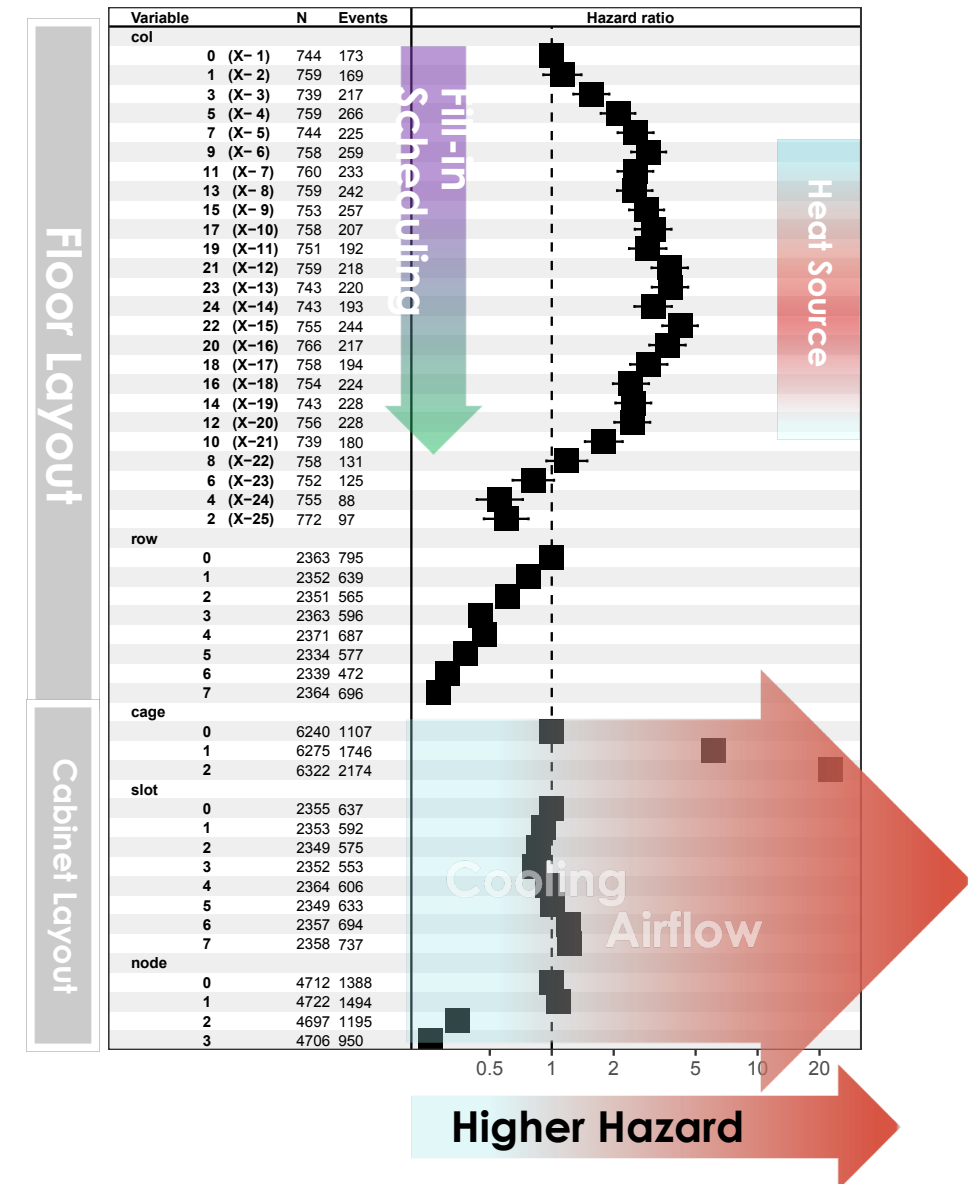
Cooling Architecture and Scheduling Affect Reliability

Conclusions for Future Systems

- Expect the unexpected with available resources
- Collect more detailed sensor data
- Use advanced statistics for configurable causal analysis
 - Provide early warning of reliability issues
 - Inform mitigation strategies

The Science

- 100,000 collective years of GPU life
- Advanced statistical methods typical in biomedicine
- Novel component lifetime visualizations
- Data and analysis codes made publicly available



Thank You!

Acknowledgements:

This research used resources of the Oak Ridge Leadership Computing Facility at the Oak Ridge National Laboratory, which is supported by the Office of Science of the U.S. Department of Energy under Contract No. DE-AC05-00OR22725.

This work was supported by the U.S. Department of Energy, Office of Science, Office of Advanced Scientific Computing Research, Resilience for Extreme Scale Supercomputing Systems Program, with program managers Robinson Pino and Lucy Nowell.