

Faults, Errors and Failures in Extreme-Scale Supercomputers

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



Saurabh Gupta, Devesh Tiwari, Tirthak Patel, and Christian Engelmann. **Reliability of HPC systems:** Large-term Measurement, Analysis, and Implications. SC'17. DOI 10.1145/3126908.3126937.

Characterizing Supercomputer Faults, Errors and Failures



Fraction of each failure type on the studied systems



Characterizing Supercomputer Faults, Errors and Failures





Spatial distribution of failures among cabinets for 3 studied systems



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



SOAK RIDGE National Laboratory

G. Ostrouchov, D. Maxwell, R. Ashraf, C. Engelmann, Mallikarjun Shankar, and James Rogers. GPU Lifetimes on Titan Supercomputer: Survival Analysis and Reliability. SC'20. DOI 10.1109/SC41405.2020.00045.

Root Cause: Non-ASR Components on SXM GPU







Silver-sulfide corrosion "Flowers-of-Sulfur"



GPU Failures and Replacements in ORNL's Titan

8



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

Machine Room Layout: GPU Locators

C##-#C#S#N# alo a b g o d i e t e n e t





= cabinet in column 17 - row 4



Cabinet Mechanical Packaging: Locating a GPU



GPU Life Data Built Incrementally from Two Sources



CAK RIDGE National Laboratory

GPU Life Visualization: Serial Number View

Critical for:

• Understanding data

SN

- Defining GPU Life
- Data processing verification







GPU Life Visualization: Location View

c24-7c2s5n0 c24-2c1s5n1 c24-1c1s2n2 c24-0c1s3n2 c23-5c1s7n0 c23-0c2s2n0 c22-7c1s1n0 5c0s2n'

c22=0c2s7n1 22-0c1s5n

4c0s6n1

18-2c2s7n0

c17-2c0s2n3

c17-0c1s6n0 c16-4c0s2n1 c16-1c0s0n3 7c2s0n3

14-3c2s7n0 c14-3c2s3n0 c14-1c2s0n1

c13-6c1s2n3 c13-6c1s2n3 c13-0c0s2n2 c12-7c0s4n2 c12-6c2s7n0 c12-2c2s6n3 c12-2c0s4n0 c12-1c0s2n2

c11-1c0s6n3 c10-7c0s1n0

c10-4c2s4n1

c10-2c2s0n2

c9-6c2s5n1 c9-3c2s7n2

9-0c1s0n0

7c0s0n3 7-3c0s4n1 7-1c2s4n2 -1c2s3n1 c6-7c2s4n2 -6-6c2s2n2 6c0s6n1

ocation

Critical for:

- Understanding data
- Defining GPU Life
- Data processing ٠ verification

5-5c0s7n3 5c0s0n0 4c1s0n2 5-4c0s4n2 4-4c1s3n2 c4-4c1s2n2 c4-1c2s5n3 24-0c1s5n3 c3-0c0s0n2 c2-4c0s6n3 c1-5c0s7n1 c1-5c0s2n0 c1-0c0s5n0 c0-5c1s3n3c0-3c2s6n2 c0-2c1s0n0 2014-01-01



Produced in R via ggplot2 and lubridate packages

CAK RIDGE National Laboratory 13

Traditional Reliability in HPC is Focused on MTBF



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

CAK RIDGE National Laboratory

14



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, $S(t) = Pr\{T \ge t\} = 1 F(t)$, 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



CAK RIDGE National Laboratory

16

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 $H_k(t) = H_0(t)e^{\sum_{i=1}^{n}\beta_i x_i}$
 - 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





18

Cooling Architecture and Scheduling Affect Reliability



Reordered columns based on Torus network



19

- Reordered column by fill-in scheduling torus coordinate
- Hot room spot due to other servers and Atlas file system
- Higher cages ran hotter
- Node 0 and 1 ran hotter







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



Short-term Future Research and Development Needs

- Portable system/center monitoring and analysis solutions
 - Collecting the right metrics
 - Proper identification of faults (online)
 - Fast and accurate root cause analysis (online and offline)
 - Using advanced statistical techniques and ML
 - There is some ongoing work at the facilities, but it is disconnected from recent research
- Low-overhead software mitigation techniques (beyond global checkpoint/restart)
 - OS/R and programming model runtime resilience features
 - Resilience for workflows
 - There is some ongoing work in fault tolerant programming models, but it is underfunded and community adoption is low



Other Future Research and Development Needs (1/2)

• Smart systems and facilities

- Autonomous resource management that considers the system/facility state and the involved trade-offs
- Automatic adaptation of systems and facilities in real-time to emerging reliability issues using AI
- Machine-in-the-loop operational intelligence (OODA loop to improve productivity and lower costs)
- Resilience in federated/distributed/complex computing environments
 - Instruments/laboratories using edge and center computing for science feedback on experiments
 - Real-time and urgent computing that has specific resilience needs
- Understanding the resilience problem in non von Neumann architectures
 - E.g., neuromorphic computing



Other Future Research and Development Needs (2/2)

• Resilience by design

- Design space exploration that considers resilience in addition to performance and power/energy
- Performance/energy/resilience co-design
- Programming for resilience (higher-level abstractions and programming models)
- Resilient algorithms and probabilistic/approximate computing
 - Algorithm-based fault tolerance
 - Coded computing
 - Naturally resilient algorithms
- End-to-end resilience (integrity of data and computation)
- Intersection between cyber security and resilience



Questions?

