

The Resilience Problem in Extreme Scale Computing

Christian Engelmann Oak Ridge National Laboratory



The Sky is Falling vs. Nothing to See Here

- Nobody buys supercomputers that don't work!
 - High error or failure rates "out of the box" are unacceptable
 - The supercomputer industry does need help with that, though
- However, significant and unexpected reliability issues during operation do happen!
 - Bad solder, dirty power, unexpected early wear-out, etc.
 - See Titan GPU failures as an example (SC'18 paper)
- We need to design the HPC hardware/software ecosystem to be able to deal with high error and failure rates!
 - It's risk mitigation 101
 - There is a cost/benefit trade-off, though





Understanding the Problem is Key

What does fail, why and how? Are our assumptions correct?



Reliability of HPC systems: Large-term Measurement, Analysis, and Implications (1/3)

- Analyzed 1.2 billion node hours of logs from 5 different OLCF supercomputers
- Combined information from different logs and created a consistent log format for analysis
- Used standard and created new methods to model the temporal and spatial behavior of failures
- Analyzed the evolution of temporal and spatial behavior over the years
- Analyzed the correlation of different failure types
- Compared the mean-time between failures of the 5 systems

OAK RIDGE

20082009201020112012201320142015Scale-normalized MTBF of each system
over time (averaged quarterly)

Scale-Normalized MTBF = $\frac{\text{MTBF} \times \text{Num of Nodes in the System}}{\text{Max Number of Nodes across all Systems}}$

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



- Jaguar XT4

(hour

MTBF

Scaled



🗕 Jaguar XT5 🛛 🕂 Jaguar XK6

Reliability of HPC systems: Large-term Measurement, Analysis, and Implications (2/3)



Fraction of each failure type on the studied systems





Reliability of HPC systems: Large-term Measurement, Analysis, and Implications (3/3)





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



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)





If You Can't Measure It, You Can't Improve It.

How do we get from lengthy postmortem analysis to real-time operational intelligence?



LogSCAN Real-Time Processing Architecture



B. H. Park, Y. Hui, S. Boehm, R. Ashraf, C. Engelmann, and C. Layton. **A Big Data Analytics Framework for HPC Log Data: Three Case Studies Using the Titan Supercomputer Log**. HPCMASPA'18.



8

Open slide master to edit

New Metrics to Evaluate HPC System Health



The value of ASI is limited to the range (0, 1). When ASI approaches 1, it represents high sparsity or a time interval in which only a few applications are generating most of system reliability events and vice versa.

Y. Hui, B. H. Park, and C. Engelmann. A Comprehensive Informative Metric for Analyzing HPC System Status using the LogSCAN Platform. FTXS'19.



Real-Time Analysis of System Information Entropy



SIE with Source Type layout (top), SIE with Nodal Map layout (middle), and Total event count (bottom)



Open slide master to edit



Coordinating Multiple Solutions is Key

Why do we abort and restart an entire job when 1 out of 27,648 GPUs has an error? Why don't we just rerun the single failed GPU execution?



Novel Solution: Design Patterns for Resilience

- A design pattern provides a generalizable solution to a recurring problem
- It formalizes a solution with an interface and a behavior specification
- Design patterns do not provide concrete solutions
- They capture the essential elements of solutions, permitting reuse and different implementations
- State patterns provide encapsulation of system state for resilience
- Behavioral patterns provide encapsulation of detection, containment and mitigation techniques for resilience



Anatomy of a Resilience Design Pattern

- A resilience design pattern is defined in an event-driven paradigm
- Instantiation of pattern behaviors may cover combinations of detection, containment and mitigation capabilities
- Enables writing patterns in consistent format to allow readers to quickly understand context and solution





Resilience Design Patterns Classification





Resilience Design Patterns Specification v1.2

- Taxonomy of resilience terms and metrics
- Survey of resilience techniques
- Classification of resilience design patterns
- Catalog of resilience design patterns
 - Uses a pattern language to describe solutions
 - 3 strategy patterns, 5 architectural patterns, 11 structural patterns, and 5 state patterns
- Case studies using the design patterns
- A resilience design spaces framework



S. Hukerikar and C. Engelmann. Resilience Design Patterns: A Structured Approach to Resilience at Extreme Scale (Version 1.2). Tech. Report, ORNL/TM-2017/745, 2017. DOI: 10.2172/1436045



Case Study: Checkpoint Recovery with Rollback





Case Study: Proactive Process Migration





Case Study: Cross-Layer Hardware/Software Hybrid Solution







Resilience by Design and not as an Afterthought

Understanding the cost/benefit trade-off requires a design space exploration process!



Resilience Design Spaces Framework

- Design for resilience can be viewed as a series of refinements
- The design process is defined by 5 design spaces
- Navigating each design space progressively adds more detail to the overall design of the resilience solution
- A single solution may solve more than one resilience problem
- Multiple solutions often solve different resilience problems more efficiently





Design Space Exploration for Resilience

- Vertical and horizontal pattern compositions describe the resilience capabilities of a system
- Pattern coordination leverages beneficial and avoids counterproductive interactions
- Pattern composition optimizes the performance, resilience and power consumption trade-off



21

Case Study: Multiresilient Iterative Linear Solver

- GMRES minimal residual method for solving nonsymmetric linear systems
 - Solve: Ax = b
 - Iterative algorithm
- Resilience patterns provide detection, containment, and mitigation for soft and fail-stop errors
 - Different soft error detectors for inner loop
 - In-memory checkpointing for process failures





Open slide master to edit



Modeling and Simulation for Design Space Exploration (Ongoing Work)

- Model the performance, resilience, and power consumption of an entire system
- Start at compute-node granularity with
 - System component models
 - Resilience design pattern models
 - Application models
- Simulate dynamic interactions between the system, resilience solutions and applications
- Move to finer-grain resolution to include on-node communication, computation and storage



National Laboratory

Resources and Contact

- Catalog: Characterizing Faults, Errors, and Failures in Extreme-Scale Systems
 - <u>https://ornlwiki.atlassian.net/wiki/spaces/CFEFIES</u>
- Resilience Design Patterns: A Structured Approach to Resilience at Extreme Scale
 - <u>https://ornlwiki.atlassian.net/wiki/spaces/RDP</u>
- Christian Engelmann, Oak Ridge National Laboratory
 - <u>engelmannc@ornl.gov</u>

