# Fine-grained Automated Failure Management on Extreme-Scale GPU Systems

2025 Feb 20

Balazs Gerofi <balazs.gerofi@intel.com>

Multicore World 2025, Christchurch, New Zealand

intel.

# Legal Disclaimer

Notice: This document contains information on products in the design phase of development. The information here is subject to change without notice. Do not finalize a design with this information.

Intel technologies may require enabled hardware, software or service activation.

You may not use or facilitate the use of this document in connection with any infringement or other legal analysis concerning Intel products described herein. You agree to grant Intel a non-exclusive, royalty-free license to any patent claim thereafter drafted which includes subject matter disclosed herein.

No license (express or implied, by estoppel or otherwise) to any intellectual property rights is granted by this document.

The products described may contain design defects or errors known as errata which may cause the product to deviate from published specifications. Current characterized errata are available on request.

Performance results are based on testing as of dates shown in configurations and may not reflect all publicly available updates. See backup for configuration details. No product or component can be absolutely secure.

Performance varies by use, configuration, and other factors. Learn more on the Performance Index site.

Your costs and results may vary.

"Conflict-free" refers to products, suppliers, supply chains, smelters, and refiners that, based on our due diligence, do not contain or source tantalum, tin, tungsten or gold (referred to as "conflict minerals" by the U.S. Securities and Exchange Commission) that directly or indirectly finance or benefit armed groups in the Democratic Republic of the Congo or adjoining countries.

All product plans and roadmaps are subject to change without notice.

Code names are used by Intel to identify products, technologies, or services that are in development and not publicly available. These are not "commercial" names and not intended to function as trademarks.

Intel disclaims all express and implied warranties, including without limitation, the implied warranties of merchantability, fitness for a particular purpose, and non-infringement, as well as any warranty arising from course of performance, course of dealing, or usage in trade.

Altering clock frequency or voltage may void any product warranties and reduce stability, security, performance, and life of the processor and other components. Check with system and component manufacturers for details.

Results have been estimated or simulated.

Intel does not control or audit third-party data. You should consult other sources to evaluate accuracy.

Copies of documents which have an order number and are referenced in this document may be obtained by calling 1-800-548-4725 or visiting the Intel Resource and Document Center.
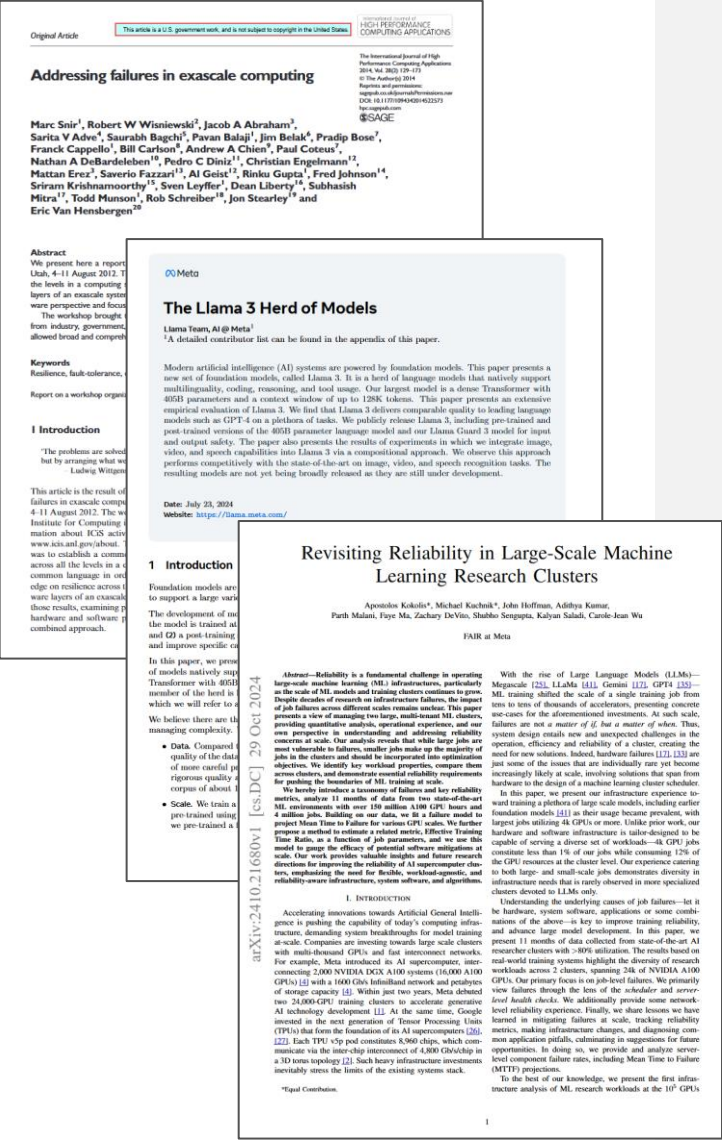
© 2025 Intel Corporation. Intel, the Intel logo, OpenVINO and the OpenVINO logo, and other Intel marks are trademarks of Intel Corporation or its subsidiaries. Other names and brands may be claimed as the property of others.
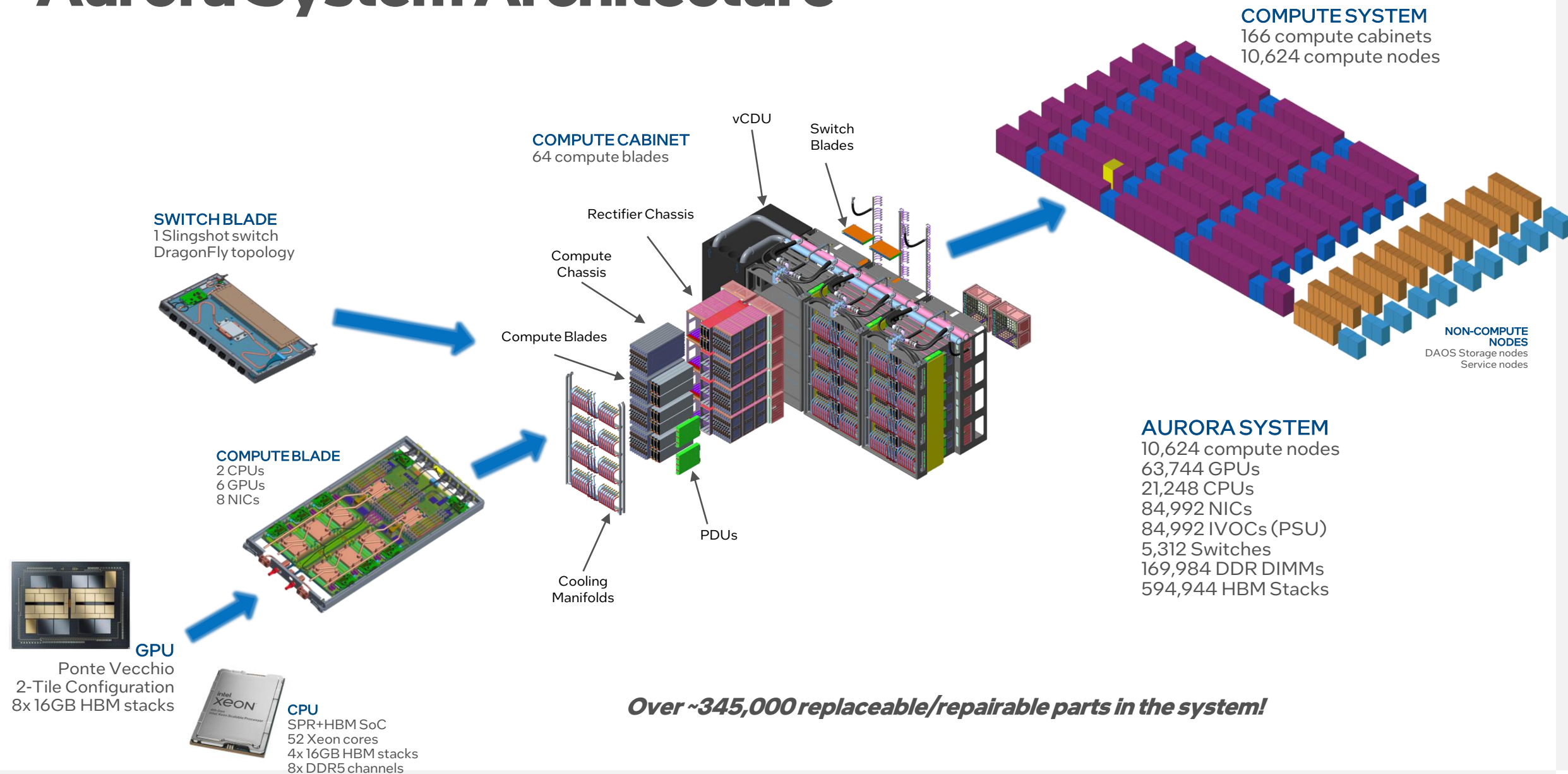
intel.

# **Agenda**

- Motivation
- Background
  - Aurora Overview
  - Failures in Large Scale Systems
- Problem Statement
- StabilityDB Architectural Overview
- Failure Strike Policy
- Failure Management Automation
- Results
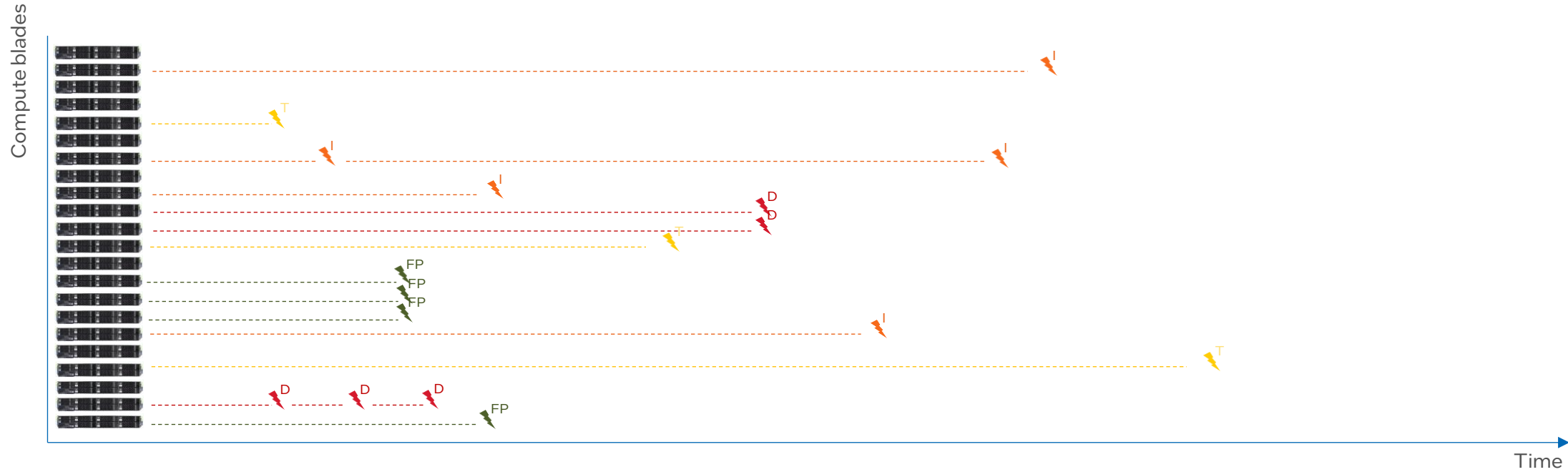- Conclusion

intel.

# Motivation

- Failures in leadership class accelerated HPC and AI systems have become the norm rather than the exception
  - This was anticipated a decade ago in the HPC community...

- As systems continue to scale in size, the frequency of failures on the entire system is expected to increase

- Tightly coupled parallel workloads (e.g., HPC modeling/simulation and AI training/fine-tuning) are highly sensitive to failures
  - A single "*lemon node*" can ruin the run

- To ensure efficient deployment and operation, automated failure management is essential

intel.

# Aurora System Architecture

**COMPUTE SYSTEM**
166 compute cabinets
10,624 compute nodes

vCDU

Switch Blades

**COMPUTE CABINET**
64 compute blades

Rectifier Chassis

**SWITCH BLADE**
1 Slingshot switch
DragonFly topology

Compute Chassis

Compute Blades

**NON-COMPUTE NODES**
DAOS Storage nodes
Service nodes

**COMPUTE BLADE**
2 CPUs
6 GPUs
8 NICs

PDUs

**AURORA SYSTEM**
10,624 compute nodes
63,744 GPUs
21,248 CPUs
84,992 NICs
84,992 IVOCs (PSU)
5,312 Switches
169,984 DDR DIMMs
594,944 HBM Stacks

Cooling Manifolds

**GPU**
Ponte Vecchio
2-Tile Configuration
8x 16GB HBM stacks

**CPU**
SPR+HBM SoC
52 Xeon cores
4x 16GB HBM stacks
8x DDR5 channels

intel xeon

*Over ~345,000 replaceable/repairable parts in the system!*

# Hardware Failures in Large Scale Systems



(D) Failures due to design bugs which fail on every instance of component <X> under specific conditions (e.g., MERT accumulator overflow bug in PVC)

(I) Failures due to intermittent faults

▪  Aged/degraded components, systematic issues due to marginalities (could also be deficiencies in design, validation, etc.), manufacturing defects

(T) Transient/random failures due to particle/EM radiation/cosmic rays

▪  These will create a "background" noise of failures that should be distributed evenly across the system

(FP) Software, networking, or correlated faults and other errors that could create failures that look like hardware component issues

# Key Observations and Problem Statement

- How to distinguish *intermittent* and *transient* failures?

- How to handle first strikes?
    - Replacing every component on first strikes is *impractical*

- Probability of an intermittent error occurring on the same component twice is extremely small
    - Indication of defects?

- Need to understand *reoccurrence rates* and the *statistical properties of durations between strikes*
    - These are specific to failure modes

- Failure history needs to be captured *in context*
    - Firmware/software versions, external conditions have impacts

- Need for automated failure categorization

- Automated failure servicing/management?

intel.

# Agenda

- Motivation
- Background
    - Aurora Overview
    - Failures in Large Scale Systems
- Problem Statement
- StabilityDB Architectural Overview
- Failure Strike Policy
- Failure Management Automation
- Results
- Conclusion
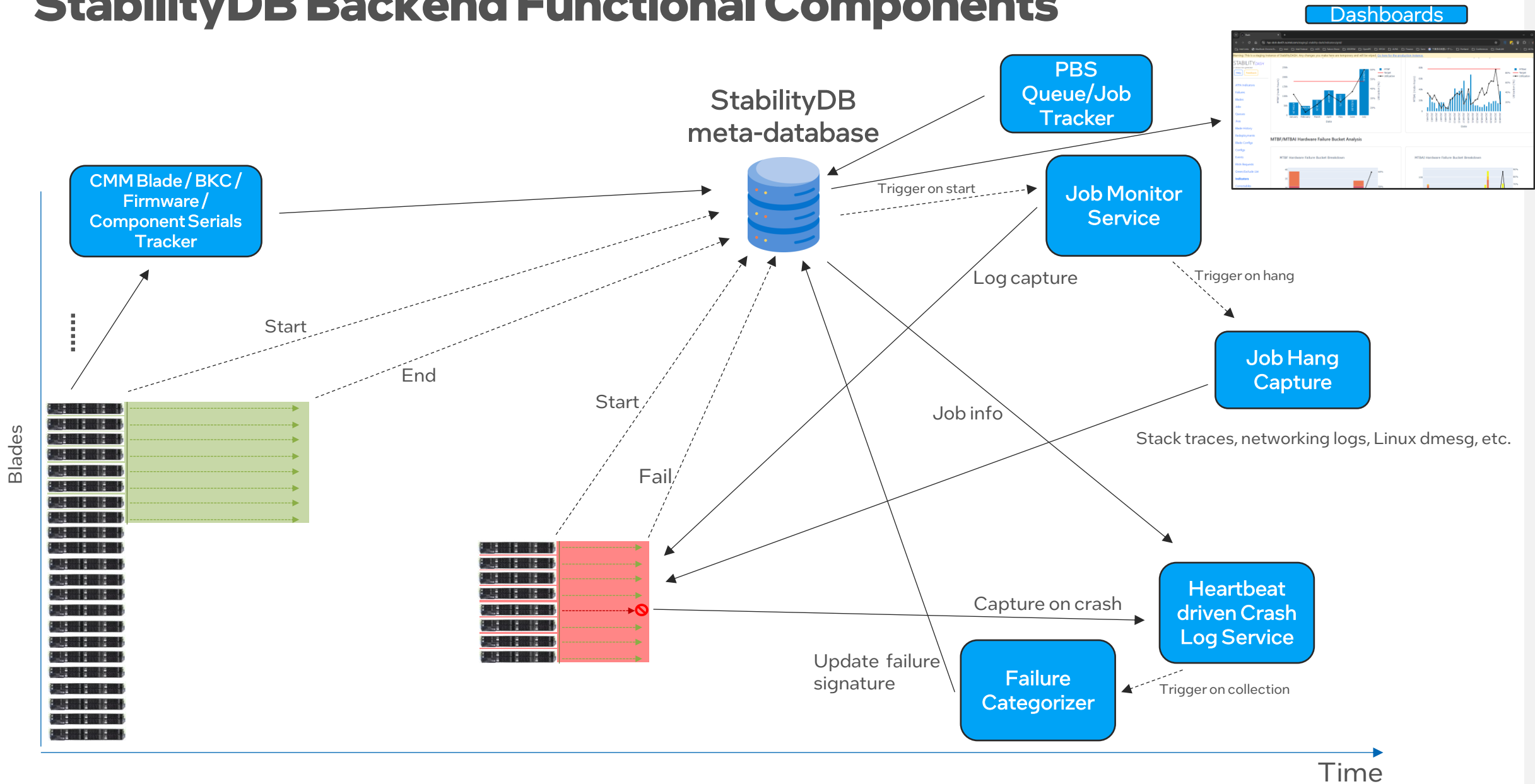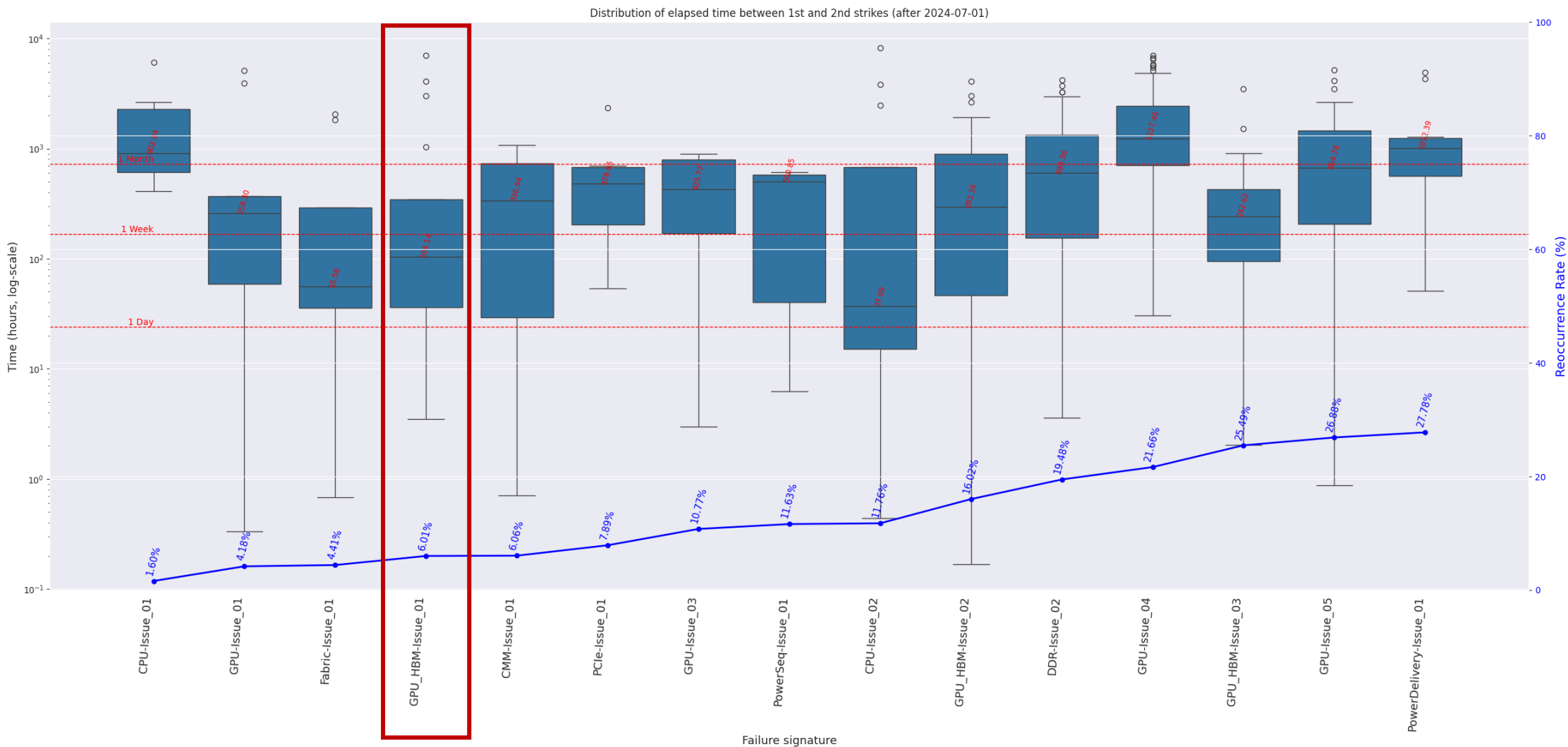
intel.

# StabilityDB Infrastructure Architectural Overview



Bundled failure log capturing

Automated triage

Failure rates and indicators

Real time system view

Inventory and node metadata

Relational DB

StabilityDB is an integrated timeseries based meta-database on compute nodes in the system, history of:
- Failures
- Inventory/component serials
- Versioning (firmware/software/etc.)
- Workloads
- Resource manager state, e.g., queue changes, etc.

Primary insight was "essential data in context", with goals of:
- Easy analysis
- Data driven guidance for bring-up/debug

- MTBF/MTBAI analysis
- Failure stats/breakdowns with workload association
- Failure stats with component/FW versioning association
- Inventory/part replacement tracking and projections
- Automated failure mitigation

# Aurora: Interaction with Cluster Management/Telemetry



Aurora Monitoring

StabilityDB meta-database

Argonne Management Node(s) — Leader Node(s) — Compute Node(s)

HPCM-internal Storage
- ELK DB (ElasticSearch)
- Timescale DB

Kafka Bus

Logstash
Conserver
LDMS Data Aggregator — IB — LDMS PVC sampler
HPCM Sensor Aggregator — OOB — BMC

DAI/DS NRE
PBS Pro

ALCF Storage node(s)
Long-term storage
RAS / Telemetry Inventory / Job
ALCF-Managed DB (PostgreSQL)

[]

- Plugs into standard telemetry/event data streams:
  - E.g., Kafka and RabbitMQ
- Interacts with cluster management software components (e.g., HPE HPCM node management)
- Interacts with batch scheduler (e.g., PBS)

# StabilityDB Backend Functional Components



Dashboards

StabilityDB meta-database

PBS Queue/Job Tracker

CMM Blade / BKC / Firmware / Component Serials Tracker

Job Monitor Service

Trigger on start

Trigger on hang

Job Hang Capture

Log capture

Start

End

Blades

Start

Job info

Stack traces, networking logs, Linux dmesg, etc.

Fail

Capture on crash

Heartbeat driven Crash Log Service

Update failure signature

Failure Categorizer

Trigger on collection

Time

# Agenda

- Motivation
- Background
  - Aurora Overview
  - Failures in Large Scale Systems
- Problem Statement
- StabilityDB Architectural Overview
- **Failure Strike Policy**
- **Failure Management Automation**
- **Results**
- **Conclusion**

intel.

# Distribution of Elapsed Time between 1ˢᵗ and 2ⁿᵈ Strikes



Distribution of elapsed time between 1st and 2nd strikes (after 2024-07-01)

- Low reoccurrence rate
- Higher reoccurrence rate but long elapsed time between strikes

intel.

# Fine-Grained Multi-Strike Policies

| Updated_ts | Sig | Live | Criteria | Strike_1 | Strike_2 | Strike_3 | Strike_4 | Strike_5 |
|---|---|---|---|---|---|---|---|---|
| 2024-12-08 13:02:41 | fatal_7_multiple_quads | 1 | single chip | HOST_COLD_RESET,REDEPLOY_STABILITYDB... | RMA_BTK | CONTACT_INTEL_REP | CONTACT_INTEL_REP | CONTACT_INTEL_REP |
| 2024-12-08 13:02:42 | fatal_7_others | 1 | single chip | HOST_COLD_RESET,REDEPLOY_STABILITYDB... | RMA_BTK | CONTACT_INTEL_REP | CONTACT_INTEL_REP | CONTACT_INTEL_REP |
| 2024-12-08 13:02:42 | fpu_fatal_singlesubslice_error | 1 | single source | DSS_SWAP,MONITOR | DSS_SWAP,MONITOR | CONTACT_INTEL_REP | CONTACT_INTEL_REP | CONTACT_INTEL_REP |
| 2024-12-08 13:02:42 | general_multibank_correctable_error | 1 | single chip | HOST_COLD_RESET,REDEPLOY_STABILITYDB... | REDEPLOY_STABILITYDB,MONITOR | REDEPLOY_STABILITYDB,MONITOR | REDEPLOY_STABILITYDB,MONITOR | REDEPLOY_STABILITYDB,MONITOR |
| 2024-12-08 13:02:42 | general_multibank_error | 0 | single chip | HOST_COLD_RESET,REDEPLOY_STABILITYDB... | GT IFR | RMA_BTK | CONTACT_INTEL_REP | CONTACT_INTEL_REP |
| 2024-12-08 13:02:42 | general_multisubslice_correctable_error | 1 | single chip | HOST_COLD_RESET,REDEPLOY_STABILITYDB... | REDEPLOY_STABILITYDB,MONITOR | REDEPLOY_STABILITYDB,MONITOR | REDEPLOY_STABILITYDB,MONITOR | REDEPLOY_STABILITYDB,MONITOR |
| 2024-12-08 13:02:42 | general_multisubslice_error | 1 | single chip | CONTACT_INTEL_REP | CONTACT_INTEL_REP | CONTACT_INTEL_REP | CONTACT_INTEL_REP | CONTACT_INTEL_REP |
| 2024-12-08 13:02:42 | general_singlesubslice_error | 0 | single source | DSS_SWAP,MONITOR | DSS_SWAP,MONITOR | CONTACT_INTEL_REP | CONTACT_INTEL_REP | CONTACT_INTEL_REP |
| 2024-12-08 13:02:42 | grf_fatal_singlesubslice_error | 1 | single source | DSS_SWAP,MONITOR | DSS_SWAP,MONITOR | CONTACT_INTEL_REP | CONTACT_INTEL_REP | CONTACT_INTEL_REP |
| 2024-12-08 13:02:42 | grf_singlesubslice_correctable_error | 1 | single source | DSS_SWAP,MONITOR | DSS_SWAP,MONITOR | CONTACT_INTEL_REP | CONTACT_INTEL_REP | CONTACT_INTEL_REP |
| 2024-12-08 13:02:42 | GPU_HBM-Issue_01 | 1 | single source | HBM_IFR,PVC_COLD_RESET,MONITOR | RMA_BTK | CONTACT_INTEL_REP | CONTACT_INTEL_REP | CONTACT_INTEL_REP |
| 2024-12-08 13:02:42 | hbm_controller_multi_errors | 0 | single chip | HBM_IFR,PVC_COLD_RESET,MONITOR | RMA_BTK | CONTACT_INTEL_REP | CONTACT_INTEL_REP | CONTACT_INTEL_REP |
| 2024-12-08 13:02:42 | hbm_data_errors | 1 | single source | HBM_IFR,MONITOR | CONTACT_INTEL_REP | CONTACT_INTEL_REP | CONTACT_INTEL_REP | CONTACT_INTEL_REP |
| 2024-12-08 13:02:42 | hbm_training_errors | 1 | single source | HBM_IFR,PVC_COLD_RESET,MONITOR | RMA_BTK | CONTACT_INTEL_REP | CONTACT_INTEL_REP | CONTACT_INTEL_REP |

- Per failure mode fine-grained description
- Policies must be based on statistical information, for that we need a meta-database tracking those metrics
- Defines failure management automation actions

# Automated Failure Management Components
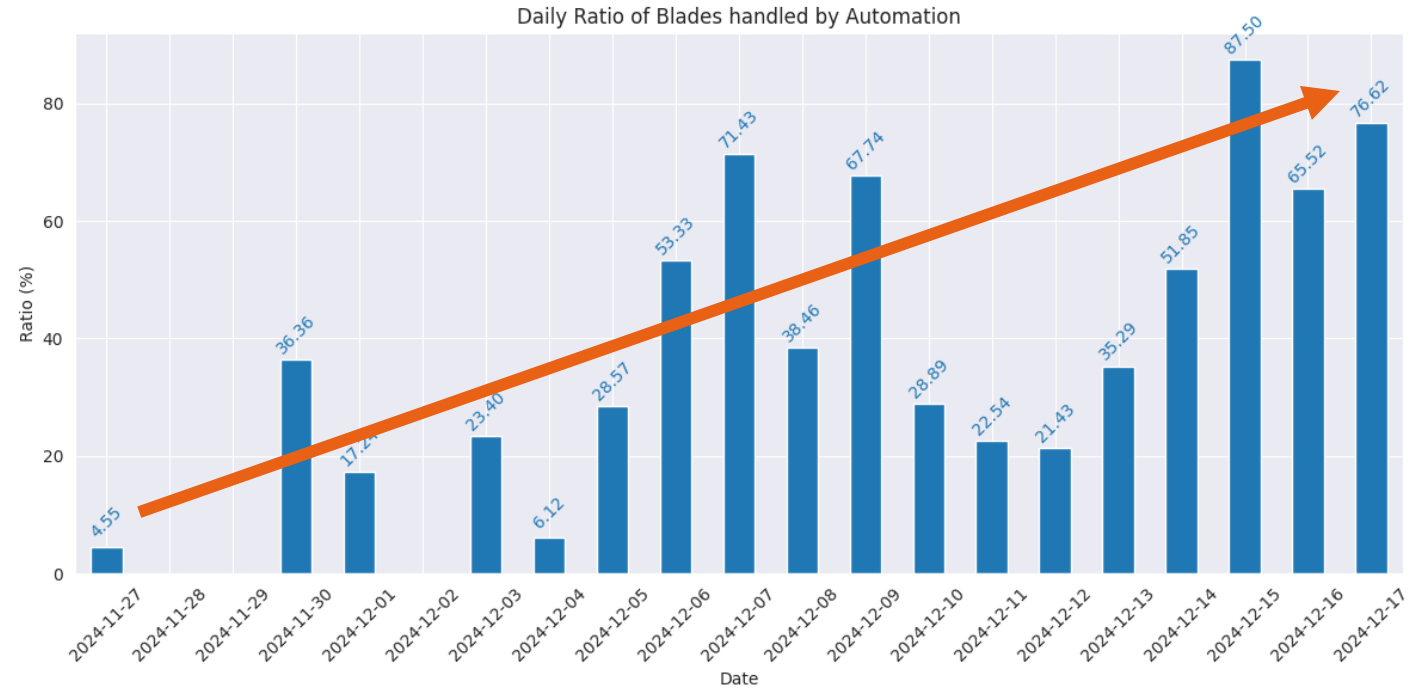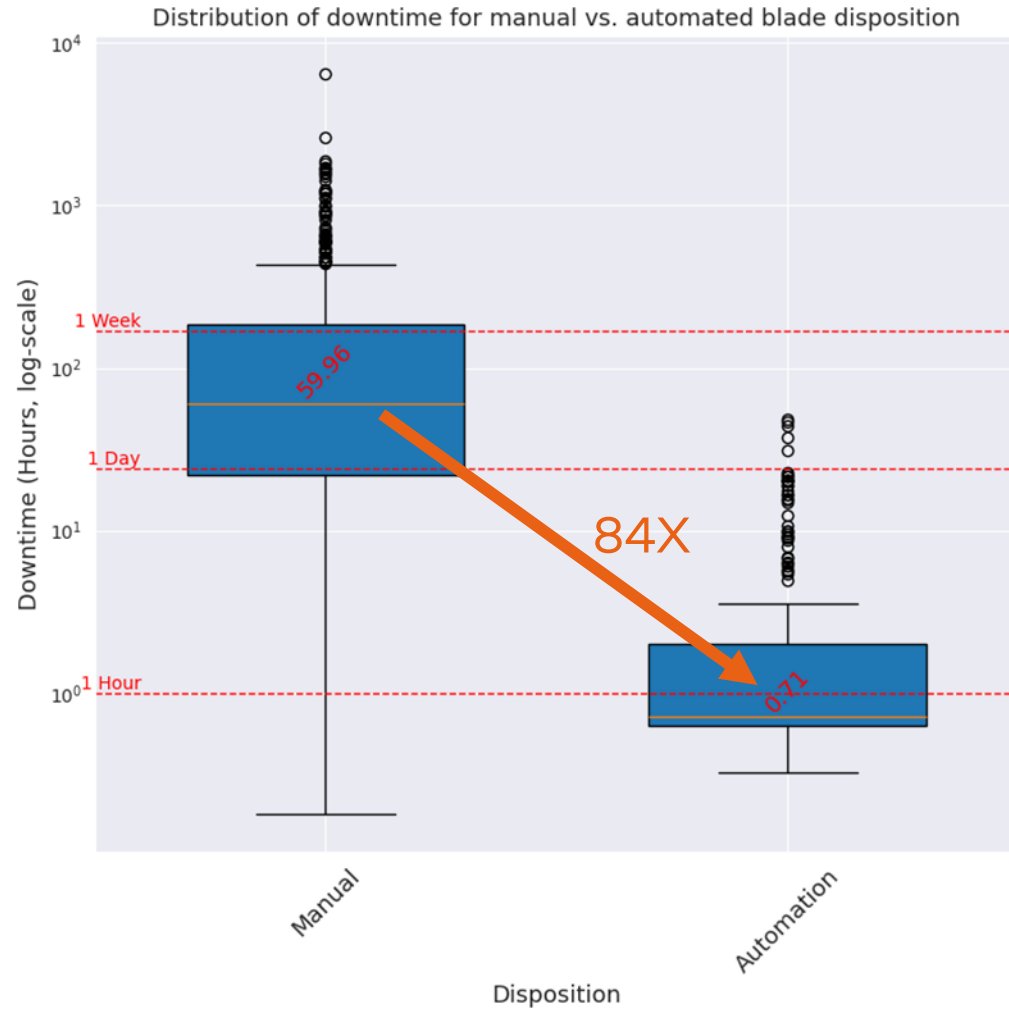


Failure triage and resource offlining

Historical data analysis and strike policy
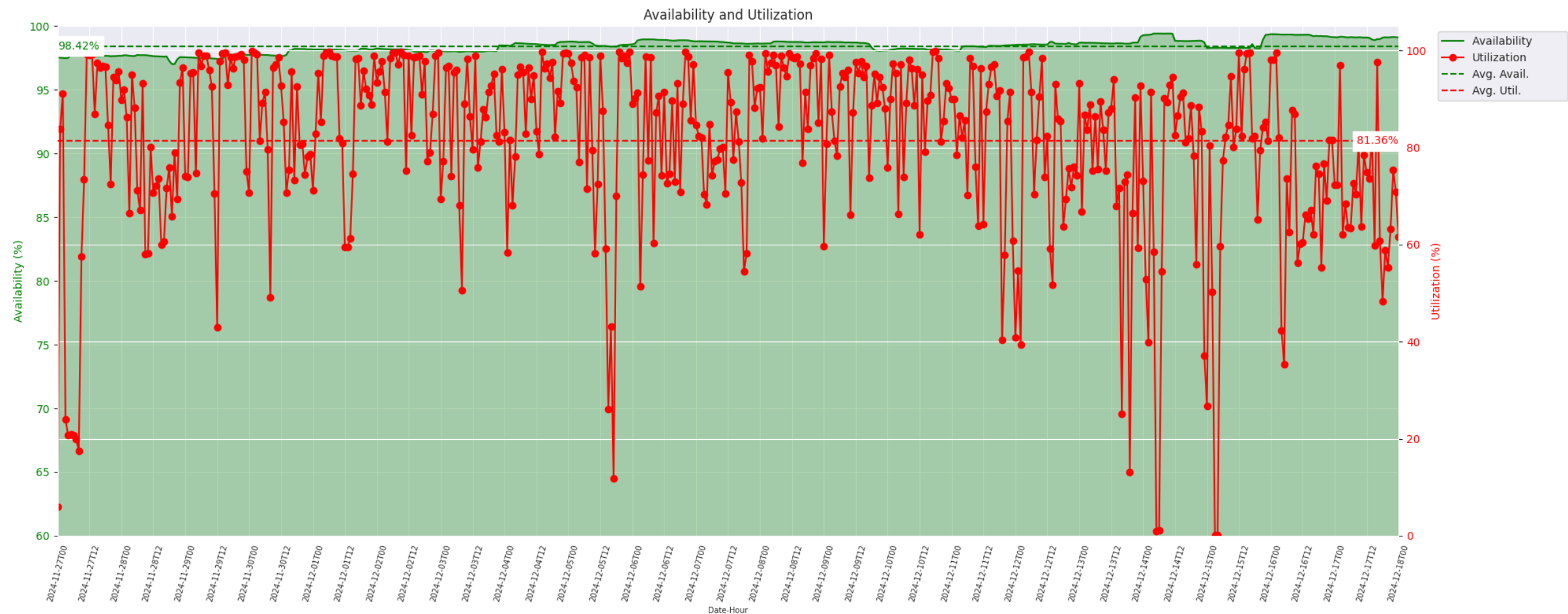
Maintenance actions

# Agenda

- Motivation
- Background
  - Aurora Overview
  - Failures in Large Scale Systems
- Problem Statement
- StabilityDB Architectural Overview
- Failure Strike Policy
- Failure Management Automation
- **Results**
- **Conclusion**

intel.

# Automated Failure Management Results (AT-S)



Distribution of downtime for manual vs. automated blade disposition



Daily Ratio of Blades handled by Automation

- Excludes large correlated events (e.g., Lustre mount failures, rack EPO events, etc.)
- Median time spent on addressing issues is 2 orders of magnitude lower!
- Automation's ratio steadily increasing for blade dispositions

# Availability and Utilization during AT-S



Availability and Utilization

- Requirement of 95% availability in average

# First and Second Strikes Month over Month

First and second strike counts per failure signature per month



- System stabilizing via decreasing 1st and 2nd counts for most failure signatures

- Constant rate of CEs on DDR and HBM

# Job Failure Breakdown during AT-S vs. Meta

**Overall Job Failure Breakdown (Flattened)**

*Very similar overall failure profiles!*

This equals to an MTBAI = 3.7 H @ 16K GPUs
Scaling this linearly to Aurora scale:
MTBAI = 55 min @ 64K GPUs

| Component | Category | Interruption Count | % of Interruptions |
|---|---|---|---|
| Faulty GPU | GPU | 148 | 30.1% |
| GPU HBM3 Memory | GPU | 72 | 17.2% |
| Software Bug | Dependency | 54 | 12.9% |
| Network Switch/Cable | Network | 35 | 8.4% |
| Host Maintenance | Unplanned Maintenance | 32 | 7.6% |
| GPU SRAM Memory | GPU | 19 | 4.5% |
| GPU System Processor | GPU | 17 | 4.1% |
| NIC | Host | 7 | 1.7% |
| NCCL Watchdog Timeouts | Unknown | 7 | 1.7% |
| Silent Data Corruption | GPU | 6 | 1.4% |
| GPU Thermal Interface + Sensor | GPU | 6 | 1.4% |
| SSD | Host | 3 | 0.7% |
| Power Supply | Host | 3 | 0.7% |
| Server Chassis | Host | 2 | 0.5% |
| IO Expansion Board | Host | 2 | 0.5% |
| Dependency | Dependency | 2 | 0.5% |
| CPU | Host | 2 | 0.5% |
| System Memory | Host | 2 | 0.5% |

**Table 5 Root-cause categorization of unexpected interruptions during a 54-day period of Llama 3 405B pre-training.** About 78% of unexpected interruptions were attributed to confirmed or suspected hardware issues.

# Supercomputer/Data Center Digital Twins?

- "Digital twins provide living digital models of physical systems that enable **data-driven analysis** and application of artificial intelligence to better manage the datacenter and **drive efficiency** for sustainability." [1]

- "Historically, data center management has been split into silos that each focus on one aspect.. as a result, … different areas can miss the bigger picture.    Digital twins help to **centralize data** from across different areas of concern **into a shared environment**" [2]



FIGURE 5. Data center digital twin software architecture. DB: database.
https://ieeexplore.ieee.org/stamp/stamp.jsp?tp=&arnumber=10687340

- Areas of potential:
  - Design: placing new servers, increasing density, improving thermal performance, etc.
  - Construction: streamlining construction, **reducing waste**, etc.
  - Operations: automating data center processes, **efficient maintenance and repairs**, etc.
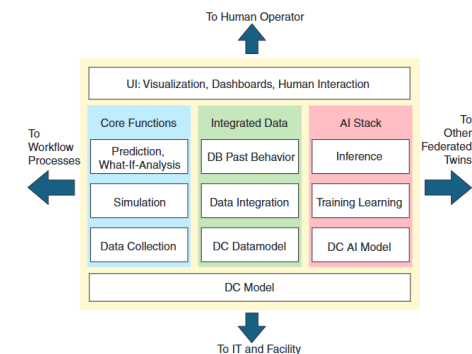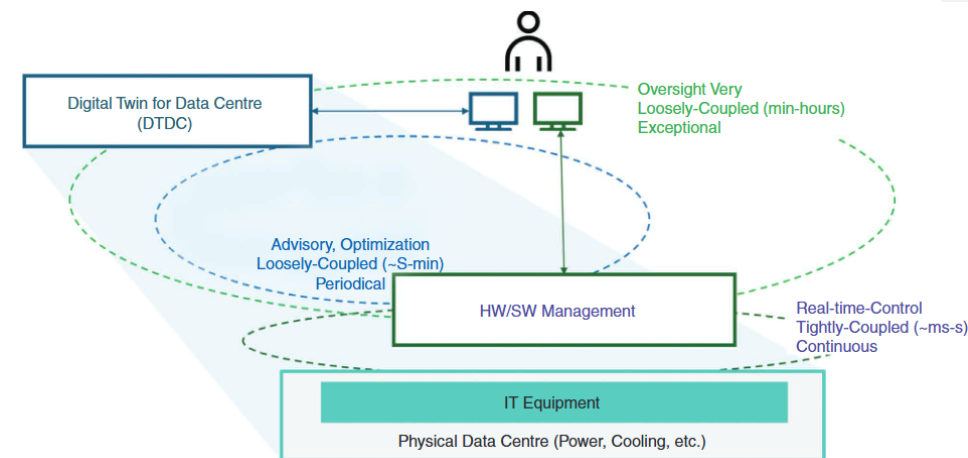  - Planning: ensuring compliance with data twins, understanding material impact, etc.



https://ieeexplore.ieee.org/stamp/stamp.jsp?tp=&arnumber=10687340

[1] J. Athavale et al., "Digital Twins for Data Centers," in Computer, vol. 57, no. 10, pp. 151-158, Oct. 2024, doi: 10.1109/MC.2024.3436945.

[2] https://venturebeat.com/ai/19-ways-digital-twins-improve-data-center-sustainability/

intel.

# Conclusions and Outlook

- As large scale systems grow in size intermittent failures become more prevalent

- Efficient operation requires automated failure management

- Fine-grained multi-strike management policy

- Key is data in context that enables real-time decision making

- Outlook:
  - Predictive (AI?) failure avoidance
  - Continuous fleet scanning
  - Standardized failure reporting across components?

intel.

# Acknowledgments/Contributors

- Yoni Levitt
- Richard Barella
- Erik Adames
- Leobardo Rountree
- Sam Zeltner
- Spurthi Lokeshappa
- Tom Musta
- Damon Millar
- Aravind Balasubramanian
- Patrick Steinbrecher
- Arjun Kripanidhi

- Aakash Patel
- Neha Gupta
- Kevin Canada
- Ky Merril
- Brian Holland
- Sucheta Raghunanda
- Ben Allen (ANL)
- Peter Upton (ANL)
- Doug Waldron (ANL)

# Thank you for your attention!
# Questions?