

Story of Two GPUs: Characterizing the Resilience of Hopper H100 and Ampere A100 GPUs

Shengkun Cui*
UIUC
Urbana, IL, USA

Archit Patke*
UIUC
Urbana, IL, USA

Hung Nguyen
UIUC
Urbana, IL, USA

Aditya Ranjan
UIUC
Urbana, IL, USA

Ziheng Chen
UIUC
Urbana, IL, USA

Phuong Cao
UIUC
Urbana, IL, USA

Brett Bode
NCSA, UIUC
Urbana, IL, USA

Gregory Bauer
NCSA, UIUC
Urbana, IL, USA

Catello Di Martino
Nokia Bell Labs
Sao Paulo, Brazil

Saurabh Jha
IBM Research
Yorktown Heights, NY, USA

Chandra Narayanaswami
IBM Research
Yorktown Heights, NY, USA

Daby Sow
IBM Research
Yorktown Heights, NY, USA

Zbigniew T. Kalbarczyk
UIUC
Urbana, IL, USA

Ravishankar K. Iyer
UIUC
Urbana, IL, USA

Abstract

This study characterizes GPU resilience in *Delta*¹, a large-scale AI system that consists of 1,056 A100 and H100 GPUs, with over 1,300 petaflops of peak throughput. We used 2.5 years of operational data (11.7 million GPU hours) on GPU errors. Our major findings include: (i) H100 GPU memory resilience is worse than A100 GPU memory, with 3.2x lower per-GPU MTBE for memory errors, (ii) The GPU memory error-recovery mechanisms on H100 GPUs are insufficient to handle the increased memory capacity, (iii) H100 GPUs demonstrate significantly improved GPU hardware resilience over A100 GPUs with respect to critical hardware components, (iv) GPU errors on both A100 and H100 GPUs frequently result in job failures due to the lack of robust recovery mechanisms at the application level, and (v) We project the impact of GPU node availability on larger-scales and find that significant overprovisioning of 5% is necessary to handle GPU failures.

Keywords

Large-scale AI/HPC System; Reliability Evaluation and Analysis; GPU Resilience; Application Impacts.

1 Introduction

Large-scale HPC systems are important not only for scientific workloads [47] but also for data analytics [1] and machine learning (ML) [21]. The main components of these systems are specialized accelerators, such as GPUs, that enable acceleration of computations, such as ML training [2, 51], ML inference [25], and simulations [39, 46].

¹*Delta* is an HPC system operated by the National Center for Supercomputing Applications (NCSA) at the University of Illinois Urbana-Champaign.

*Equal Contribution.

This paper studies the resilience of A100 GPUs and compares the result with H100 GPUs of the GH200 Grace Hopper Superchips², together with their associated memory: 40 GB HBM2e on each A100 GPU, and 96 GB HBM3 on each H100 GPU, respectively. The A100 GPU nodes and H100 GPUs (GH200 Superchip) nodes are operated as two independent systems sharing a storage cluster running the Lustre file system, allowing us to study and compare the two systems (refer to as *Delta*). The workflow on *Delta* involves users from universities nationwide and presents a spectrum of HPC and ML workloads. The study uses 2.5 years of data on critical GPU errors collected across the stated GPUs, encompassing 9.6 million GPU hours of A100 GPUs and 2.1 million GPU hours of H100, a combined 11.7 million GPU hours.

This study assesses (i) the resilience of GPU hardware and memory components; (ii) the error propagation paths in GPU memory, GPU hardware, and NVLink interconnect; and (iii) the impact of the observed GPU errors on user jobs. Figure 1 shows an example error propagation path for an uncorrectable double-bit memory error in H100 GPU that caused user job failure. The complete recovery process required node draining and GPU reset, which took 19 hours following the error detection.

Our major findings include:

(i) H100 shows 3.2× lower per-GPU mean time between errors (MTBE) compared to A100 for uncorrectable ECC memory errors. The per-GB MTBE of the H100’s HBM3 memory is 24% lower (~ 8.5M hours) than the A100’s HBM2e memory (~ 11.3M hours). We conjecture that the reduction in memory resilience stems from H100’s higher memory capacity.

(ii) The GPU memory error-recovery mechanisms on A100 and H100 GPUs (e.g., memory row remapping, error containment) [35]

²GH200 Superchip, hereafter.

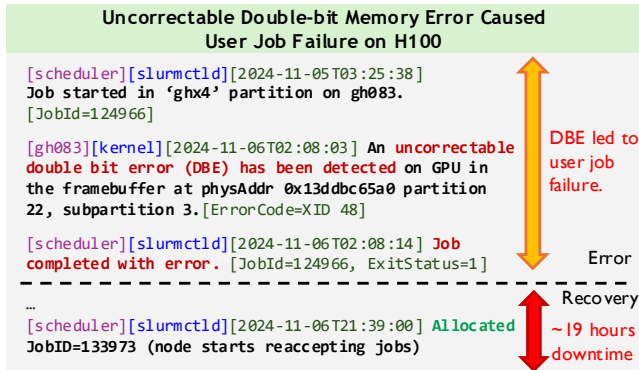


Figure 1: A double-bit memory error (XID 48) occurred, and it is uncorrectable by the SECCDED ECC HBM3 memory. Due to this double-bit error, the user job scheduled on that GPU failed, as reflected in the scheduler logs. Subsequently, this uncorrectable memory error requires a node draining and reset to complete the row remapping recovery action. The total recovery process for this incident took 19 hours, during which the node was unavailable for accepting new jobs. This incident shows that a GPU error can lead to user job failures and significantly impact node availability.

improve GPU memory resilience and reduce service interruption. We observed that these mechanisms mitigate (e.g., using memory row remapping) 92% of uncorrectable ECC memory errors on H100 GPUs. However, the memory error-recovery mechanism is insufficient to handle the increase in memory capacity and the corresponding increase in row remapping events on the H100 GPUs.

(iii) H100 GPUs demonstrate significantly improved hardware resilience over A100 GPUs, with respect to critical components such as GSP³, NVLink, and PMU SPI⁴, which were major sources of job failures in A100 systems. We attribute this to driver-level enhancements and tighter integration [32, 34], which contribute to improved resilience. Specifically, comparing H100 and A100 GPU hardware, we observed (a) a significant reduction of GSP errors on H100 (only 3 cases in our measurement period) and (b) the elimination of PMU SPI error propagations, which on A100 GPUs can lead to MMU errors 88% of times with 90% leading to user job failures, (c) no NVLink errors on H100 GPUs during the measurement period.

(iv) GPU errors on both A100 and H100 GPUs frequently result in job failures due to the lack of robust recovery mechanisms at the application level. Except for MMU and NVLink errors, other GPU errors cannot be handled by application-level mechanisms, resulting in close to 100% job failure rate. The underlying cause of job failures differs by GPU type: hardware errors are the predominant cause in A100 GPUs, whereas memory errors are the primary cause in H100 GPUs. Overall, we find that application-based recovery strategies are largely ineffective; hence, there is a compelling need to improve resilience at the GPU memory and hardware level.

(v) The overall availability per-GPU node is approximately 99.4% for A100 GPUs and 99.3% for H100 GPUs, corresponding

³A GPU system processor (GSP) is an onboard co-processor that offloads driver tasks from the CPU for latency and performance improvement.

⁴PMU on an NVIDIA’s GPU regulates the frequency, voltage, and power of the GPU based on various factors such as temperature and power cap. SPI stands for serial peripheral interface, which serves as the communication channel between peripheral hardware.

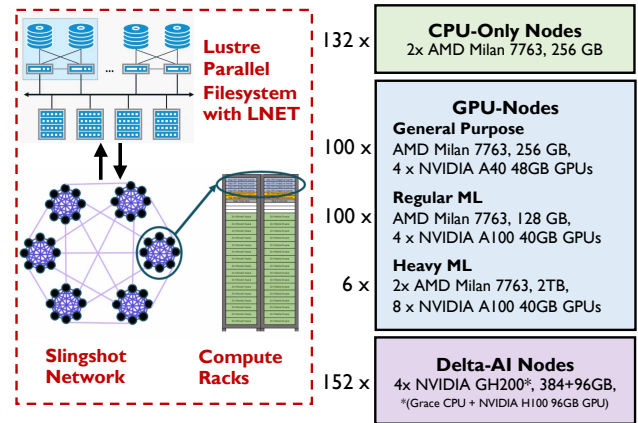


Figure 2: System architecture and specifications of *Delta*. This study focuses on the H100 and A100 GPU nodes.

to a downtime between 9–10 minutes per day. We projected the impact of this measured availability on larger scales via emulation. For example, to maintain 99.9% availability at the job level, overprovisioning of 5% would be necessary. While at first glance, such overprovisioning would appear to be a small cost, for the above example, it would cost over \$1 million per month. If GPU node availability were improved to 99.9%, the required overprovisioning would reduce by 2.5x.

Putting the paper in perspective. Previous studies on characterizing GPU resilience in large-scale systems [10, 11, 17, 18, 30, 31, 36, 37, 49, 50] focus on GPU memory errors in older GPU generations (Tesla, Kepler, and Volta) that lack the latest resilience (e.g., row remapping, error containment, NVLink CRC-retry) and performance features (e.g., GPU System Processor) introduced in NVIDIA Ampere-generation GPUs. A recent study from Meta [24] characterizes cluster-level resilience for two large-scale machine learning clusters equipped with A100 GPUs. Our paper provides a deeper understanding of GPU errors and failures of two recent generations of GPUs, A100 and H100, and their impact on a broad set of HPC/ML applications. To the best of our knowledge, this is the first study on A100 and H100 GPU errors in HPC/ML systems.

2 Background

This section provides information on (i) *Delta* specification and overall GPU utilization, (ii) critical GPU error categories used in this study, and (iii) GPU error management and recovery.

2.1 *Delta* Specifications and Utilization

***Delta* Specifications.** Figure 2 shows the layout of *Delta*, which consists of a cluster of 106 4- and 8-way A100 40 GB GPU nodes with 448 A100 GPUs, and in addition, it has a second cluster of 152 4-way GH200 nodes with 608 H100 96 GB GPUs, a total of 1,056 A100 and H100 GPUs. The two clusters share a common storage cluster running Lustre Filesystem. Our study focuses on the resilience of A100 GPUs and H100 GPUs (the GPU of GH200), as they are optimized for AI/ML workloads, exhibit the highest utilization, and incorporate the latest resilience features.

***Delta* GPU Utilization.** *Delta*’s NVIDIA A100 GPUs are frequently scheduled and utilized, with an average GPU utilization

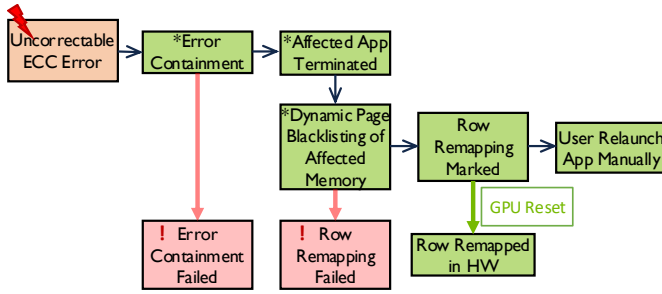


Figure 3: NVIDIA memory error recovery process for A100 and H100 GPUs.

of 51% during the operational period. NVIDIA H100 GPUs show a slightly lower average utilization (41%) than A100 GPUs.

2.2 NVIDIA GPU Error Categories

NVIDIA GPU errors are reported as XID errors. In this study, we selected a subset of XID errors that are described as common and high-impact by NVIDIA’s Developer Manuals [33, 35], NVIDIA Developer Forums and Blogs, and *Delta* site reliability engineers (SREs). We primarily collected errors and their associated recovery events. The selected XID errors/events indicate GPU issues that often cannot be resolved without SRE’s interventions (e.g., node service and GPU replacement). The selected XIDs and their corresponding GPU errors are described in Section 4, Table 1. We categorize the selected GPU errors into three categories: (i) GPU hardware, which includes all onboard hardware except for GPU memory and NVLink interconnect, (ii) NVLink interconnect, and (iii) GPU memory. Note that General GPU Software Error (XID 13) and Reset Channel Verification Error (XID 43) are usually caused by user jobs and do not impact the health of the GPU [33]; we excluded those errors from our study.

GPU Hardware Errors. The critical GPU hardware errors we studied include MMU⁵ errors, GPU Fallen Off the Bus errors, GSP errors, and PMU communication errors. We do not consider other GPU hardware errors in our study. GPU hardware errors can lead to user job failures, GPU halt, and data corruption. Among those errors, GPU Fallen Off the Bus and GSP errors lead to GPU failures, and manual GPU resets or node reboots are required to recover from the error [33]. *Delta* SREs monitor GSP errors closely to ensure timely recovery to maintain GPU availability.

GPU Interconnect (NVLink) Errors. GPU-GPU NVLink errors are caused by faulty GPU hardware, connectors, or improper connector installation during system integration, and can lead to GPU unavailability and user job failures. NVLink errors impede data transfer between GPUs and reduce computational throughput. A GPU reset or node reboot is required to clear NVLink errors [33].

GPU Memory Errors. GPU memory errors included in this study are double-bit errors (DBEs) and consecutive single-bit errors (SBEs)⁶. Individual SBEs are not logged, as they are automatically corrected by ECC. DBEs and consecutive SBEs are considered uncorrectable ECC memory errors by the NVIDIA driver, and they trigger downstream error-recovery mechanisms [33, 35], which are introduced in Section 2.1. Failures in these mechanisms can lead to

⁵The memory management unit (MMU) provides essential memory I/O functionalities.

⁶Consecutive SBEs are multiple single-bit error (SBE) occurrences at the same memory location.

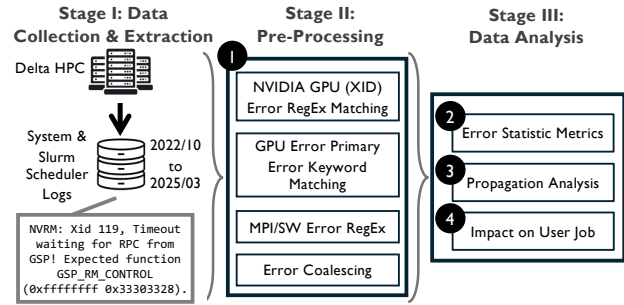


Figure 4: Overview of our data collection, processing, and analysis pipeline.

GPU/node failures and require GPU or node reboots to recover [33]. *Delta* SREs continuously monitor uncorrectable ECC memory errors and error-recovery failures to ensure timely replacement of faulty GPUs.

2.3 NVIDIA GPU Error Management

Here, we provide an overview of the resilience architecture of NVIDIA A100 and H100 GPUs.

GPU Memory. Figure 3 shows the uncorrectable ECC memory error-recovery process [35] for A100 and H100 in more detail. The primary mechanism for mitigating uncorrectable ECC memory errors for A100 and H100 GPUs is *row-remapping*, wherein the faulty memory row is replaced with a spare row, and a row-remapping event (RRE) is logged. The actual row remapping happens at the next GPU reset (e.g., during node reboot or maintenance). If there are no spare memory rows, a row remapping failure (RRF) is indicated [33, 35].

A100 and H100 GPUs support online recovery mechanisms such as *error containment* and *dynamic page offlining* [33, 35] for mitigating uncorrectable ECC memory errors with minimal node interruption. The dynamic page offlining marks the faulty memory page as unusable without requiring a GPU reset to maintain availability. The error containment procedure terminates user processes using the faulty memory address to prevent error propagation to other applications. Successful error containment is logged as a Contained Memory Error, whereas an unsuccessful error containment is logged as an Uncontained Memory Error. Failure in a row-remapping or error containment can cause a GPU failure that requires a GPU reset or node reboot. *Delta* SREs monitor row-remapping failures and replace GPUs that repeatedly emit such errors.

GPU Hardware. While the GPU caches and memory are SECEDED protected, information on failure-recovery mechanisms on GPU hardware, including peripheral hardware such as GSP, PMU, or SPI communication channels, is limited.

GPU Interconnect (NVLink). NVLink employs Cyclic Redundancy Checks (CRCs) for error detection to ensure the integrity of flow control digits and data. Upon encountering a CRC checksum error, NVLink retries packet transmissions from the last-known good packet.

Algorithm 1: Error Coalescing and Persistence Analysis.

```

Input :Error logs with timestamps  $E = \{(e_1, t_1), \dots, (e_n, t_n)\}$ , regex
         patterns  $\mathcal{R} = \{r_1, r_2, \dots\}$ , time window  $\Delta t$ 
Output: Coalesced errors with persistence duration  $E'$ 
 $E' \leftarrow \emptyset$  // Initialize output set
foreach pattern  $r \in \mathcal{R}$  do
     $E_r \leftarrow \{(e_i, t_i) \in E \mid e_i \text{ matches } r\}$  // Filter errors with regex
     $i \leftarrow 1$  // Loop through errors in a matched group
    while  $i \leq |E_r|$  do
         $(e_{\text{first}}, t_{\text{start}}, t_{\text{latest}}) \leftarrow (e_i, t_i), t_i$ 
        // Loop through later errors within the matched group
        while  $i + 1 \leq |E_r|$  do
             $(e_{\text{next}}, t_{\text{next}}) \leftarrow (e_{i+1}, t_{i+1})$ 
            // Error has identical message and is close in time
            if  $e_{\text{next}} = e_{\text{first}}$  and  $t_{\text{next}} - t_{\text{latest}} \leq \Delta t$  then
                 $t_{\text{latest}} \leftarrow t_{\text{next}}$  // Discard latest error
                 $i \leftarrow i + 1$ 
            else
                break
            end if
        // Store coalesced error and persistence duration
        Add  $(e_{\text{first}}, t_{\text{start}}, t_{\text{latest}} - t_{\text{start}})$  to  $E'$ 
         $i \leftarrow i + 1$  // Move to the next unprocessed error
    end while
end foreach
return  $E'$ 

```

3 Methodology

3.1 Data Sources

Our analysis was performed on data collected from *Delta* over its operational period: (i) 895 days from October 2022 to March 2025 for A100 GPUs and (ii) 146 days from October 2024 to March 2025 for H100 GPUs, covering 11.7 million GPU hours. This section describes data sources for *Stage I*: data collection and extraction in the pipeline in Figure 4.

System logs. System logs collected from all compute nodes capture events across system components. We created a set of regular expression (RegEX) patterns and used it to extract GPU error-recovery log entries by referring to NVIDIA XID messages [33] from the system logs (Figure 4, item (1)). The GPU error logs were our major sources of error and recovery information.

Slurm scheduler database. *Delta* uses the Slurm Workload Manager [58] (“Slurm scheduler” hereafter) for scheduling user jobs. The Slurm scheduler database keeps track of user job scheduling events, including the start and end times, the scheduled nodes, resource usage, job status, exit codes, and the `srun` command line. We used the Slurm database for user job failure characterization.

NVIDIA DCGM database. *Delta* uses NVIDIA Data Center GPU Manager (DCGM) to collect metrics and status data from all GPUs with a one-minute granularity. We used DCGM metric data to characterize GPU utilization.

3.2 Data Processing Pipeline

This section focuses on *Stage II and III* of the data processing pipeline in Figure 4, which pre-processes the raw logs, compute error counts, mean times between errors (MTBE), error propagation, and impact on user jobs.

Error Coalescing Analysis. The error coalescing step in Figure 4, item (1) filters out duplicated errors. While most errors are logged as isolated events, there are frequent periods during which the same error is logged repeatedly in close succession, i.e., there are error bursts. During these bursts, the system continually detects

and attempts to recover from errors, which could lead to either system recovery or failure. To prevent over-counting, we assume that identical error logs within a short time interval (Δt) from the same GPU are caused by the same issue. Thus, the error coalescing step counts only the first occurrence by combining identical error log lines from the same GPU within a predefined time interval (Δt) into a single error (see Algorithm 1). The remaining analyses in this paper were conducted on errors after the coalescing. We set $\Delta t = 5$ as decreases in Δ results in more duplicated logs, while further increase in Δ results in negligible changes in coalesced log count.

Error Statistic Metrics. Using the coalesced error logs as input, the pipeline computes standard error statistics metrics such as the count and the MTBE as in [36] (Figure 4, item (2)). For errors that are not directly countable from XID logs, we estimate their occurrence by correlating related XIDs within the same time interval [33, 35, 36]. In particular, the number of uncorrectable ECC memory errors can be inferred by summing up the number of RREs and RRFs, as they are mutually exclusive outcomes of an uncorrectable ECC memory error recovery event. Subsequently, the number of consecutive SBEs can be obtained by subtracting the number of DBEs from the number of uncorrectable ECC memory errors.

In addition, we computed system-wide MTBE and derived per-node MTBE by normalizing the error count using the number of GPU nodes in *Delta*. The per-node MTBE indicates the operational hours a single *Delta* GPU node can function before encountering an error. For GPU memory resilience characterization, we additionally derived per-GPU MTBE by normalizing the error count using the number of GPU in a node and per-GB MTBE by normalizing the error count with per GPU memory capacity in GB. The per-GPU MTBE reflects the operational hours a single GPU can function before encountering an error.

Error Propagation Analysis. We performed error propagation analysis (Figure 4, item (3)) to capture how errors propagate within a GPU and across different GPUs while measuring the propagation time. The propagation probability from GPU error e_1 to e_2 is defined:

$$P(e_2|e_1) = \frac{\#e_2}{\text{Total } \#e_1}, \quad t_{e_2} - t_{e_1} \leq \Delta t.$$

A propagation path is created if e_2 occurs immediately after e_1 within a predefined time window Δt . If there is no succeeding error e_2 after e_1 within Δt , then e_1 is a terminal error that does not propagate. For intra-GPU⁷ propagation, we require errors e_1 and e_2 to be on the same GPU device, whereas for the inter-GPU propagation, e_1 and e_2 are from two distinct GPUs on the same node. We recorded the time difference between the initial (e_1) and subsequent errors (e_2), referred to as the *propagation time*, for each propagation event. A shorter propagation time suggests a higher correlation between e_1 and e_2 . We applied the same Δt selection criteria as the error coalescing analysis and selected the $\Delta t = 5$ seconds in error propagation analysis.

User Job Impact Analysis. The user job impact analysis step (Figure 4, item (4)) associates GPU errors with failed user jobs to characterize the impact of GPU errors on user jobs. Section 5 provides detail on this analysis.

⁷GPU devices are identified by their node ID and PCI Express bus address.

Table 1: *Delta* NVIDIA Ampere A100 and Hopper H100 (GPU of GH200 Superchip) GPU resilience statistics.

Event Code	Abbr.	Category	Description	Recovery Action	Count		MTBE (hrs)			
					A100	H100	System-wide		Per Node	
							A100	H100	A100	H100
XID 31	MMU Error	Hardware	GPU memory management unit (MMU) error.	MMU error due to invalid memory access or driver/hardware bugs.	8,863	1,737	2.4	2	257	307
XID 48	DBE	Memory	Double bit ECC memory error (DBE).	Triggers RRE; GPU reset or node reboot is needed to clear this error.	1	17	-	206	-	31,330
-	Consecutive SBEs	Memory	Consecutive single-bit ECC memory errors (SBEs).	Triggers RRE; GPU reset or node reboot is needed to clear this error.	33	7	651	501	68,996	76,087
-	Uncorrectable ECC memory Errors	Memory	Consecutive SBEs or a DBE.	Triggers RRE; GPU reset or node reboot is needed to clear this error.	34	24	632	146	66,967	22,192
XID 63	RRE	Memory	Row remapping event triggered by uncorrectable MBE: one DBE or two SBEs at the same memory address.	GPU reset needed for row remapping.	34	16	632	219	66,967	33,288
XID 64	RRF	Memory	Row remapping failure of a row remapping event.	GPU reset is needed to clear this error.	0	8	-	438	-	66,576
XID 74	NVLink Error	Inter-connect	NVLink error indicating connection issues between GPUs via NVLink interconnection.	GPU reset or SRE intervention required.	1,922	-	11	-	1,185	-
XID 79	GPU Fallen Off the Bus Error	Hardware	GPU has fallen off the system bus and is not reachable, typically because of driver or hardware issues.	GPU reset or SRE intervention required.	10	-	2,148	-	227,668	-
XID 94	Contained Memory Error	Memory	Uncorrectable contained ECC error, indicating successful containment.	Not specified.	13	14	1,652	250	175,144	38,043
XID 95	Uncontained Memory Error	Memory	Uncontained memory error, indicating failure in containment.	GPU reset or SRE intervention required.	11	19	1,953	184	206,989	28,032
XID 119/120	GSP Error	Hardware	NVIDIA GPU Systems Processor (GSP) error.	GPU reset or SRE intervention required.	3,857	3	6	1,168	590	177,536
XID 122/123	PMU SPI Error	Hardware	PMU SPI read/write failure, indicating failed communication with the PMU.	Not specified.	77	-	279.0	-	29,570	-

*NVIDIA A100/H100 GPU supports page retirement and up to 512 row remappings (RRE); previous generations support only 64 page retirements (no row remapping support).

*Row remapping, contained memory error, and uncontained memory error are new resilience features introduced starting with the NVIDIA Ampere architecture for uncorrectable ECC memory error management.

*Per-node MTBE in hours is derived by multiplying system MTBE by the number of GPU nodes of the corresponding GPU type. The number of nodes and number of GPUs per node type are specified in Figure 2.

*All XID events presented, except for row remapping events, are errors. However, for simplicity, we treat all XID events as errors in this paper.

*Uncorrectable ECC and Consecutive SBEs are inferred from the corresponding recovery event RRE and RRF (see Section 3.2).

4 Characterizing GPU Resilience

This section characterizes and compares the resilience of *Delta*'s NVIDIA A100 (Ampere) and H100 (Hopper) GPUs. Specifically, we discuss error statistics and error propagation of GPU errors in three categories: (i) GPU memory, (ii) GPU hardware, and (iii) NVLink interconnect, as described in Section 2.2 and Table 1. These errors are critical because they interrupt user jobs and lead to unplanned node downtime, as we show in Section 5. We first highlight error statistics and key findings from our analysis and then focus on error propagation for each of the three error categories on *Delta*'s workload. Note that we do not directly compare *Delta* with Blue Waters [11], Titan [30], or Summit [36] because of the significant generational gap from the GPUs used in those machines.

4.1 Error Statistics and Result Highlights

Table 1 summarizes the selected critical GPU error statistics for A100 GPUs and H100 GPUs during the operational period, including error count, system-wide mean time between errors (MTBE) and per-node MTBE. As described in Section 2.2, these selected errors are critical XID errors that can lead to user job interruption

and node downtime, requiring manual SRE intervention for recovery. In the operational period, a total of 14,821 critical errors (listed in Table 1) were recorded for *Delta*'s A100 GPUs, with a system-wide MTBE of 1.4 hours and node MTBE of 154 hours. *Delta*'s H100 GPUs encountered 1,821 GPU errors, with a system-wide MTBE of 1.9 hours and node MTBE of 292 hours, higher than A100 nodes. Next, we use the results from Table 1 to assess the resilience characteristics of GPU memory, GPU hardware, and NVLink interconnect in greater detail.

GPU Memory. Recall that, each A100 GPU incorporates 40 GB of HBM2e memory, and each H100 GPU incorporates 96 GB of HBM3 memory. H100's shows 3.2× lower per GPU MTBE for uncorrectable ECC memory errors compared to A100's, despite comparable per-GB MTBE. This reduced resilience likely stems from H100's higher memory density, which trades resilience for capacity and performance. The observed lower per node MTBE of RREs and higher RRF counts on H100 GPUs further suggest that current resilience features are insufficient for increased memory capacity. As a uncorrectable memory error can cause a multi-GPU

job to fail, enabling uninterrupted execution in such cases at scale remains an open challenge.

(i) *Delta*'s H100 GPUs exhibit lower memory resilience than the A100 GPUs. Specifically, the per GPU MTBE is 3.2× lower on H100 GPUs at 88,768 hours versus 283,271 hours for A100 GPUs. Based on the per GPU MTBE, we calculated the per GB MTBE to be 8,521,728 hours⁸ for H100's HBM3 memory vs. 11,330,826 hours for A100's HBM2e memory, a 24% reduction.

We attribute the decrease in resilience is primarily due to the higher memory capacity (96 GB vs. 40 GB, a 2.4× increase), which increases the chances of bit flips. We additionally hypothesize that H100 memory resilience is worse due to (a) a lower signaling voltage that increases susceptibility to bit flips [56] and (b) an increased number of stacks that make heat dissipation challenging and degrade the resilience of memory modules, of the HBM3 memory.

(ii) Uncorrectable ECC memory error-recovery mechanisms (e.g., row remapping and error containment, see Section 2.1) improve GPU memory resilience and reduce service interruptions [35]. The error recovery mechanisms mitigate uncorrectable memory errors with a probability of 0.92 on H100 according to error propagation analysis in Section 4.2. However, these mechanisms may not scale well with the increase in GPU memory capacity in H100. For example, the available spare rows for row remapping are capped at the same 512 rows [35], which is not proportional to the 2.4× increase in memory capacity. Such insufficiency can be evident from the significantly lower per node MTBE of RRE on the H100 GPUs than the A100 GPUs, indicating more frequent recovery events. In addition, we observed 8 RRFs on H100 GPUs during the early operational period, which indicates memory recovery failure due to exhaustion of spared memory rows. We have yet to observe an RRF on the A100 GPUs during the much longer operational period.

GPU Hardware and NVLink Interconnect. GPU hardware such as GSP, PMU SPI, and NVLink are critical components from a resilience perspective for A100 GPUs, leading to node downtime and job failures. In this context, H100 GPUs demonstrate significant improvements in GPU hardware resilience, especially in critical components such as the GSP, PMU SPI, and NVLink. We believe that these improvements are likely due to the tightly integrated heterogeneous CPU-GPU architecture of GH200 [34] and driver-level enhancements [32]. We present a detail analysis below.

(i) Among the GPU hardware components on A100 GPUs, GSP, intended as a performance enhancement, is the most vulnerable due to its lack of robust detection and recovery. Our error propagation analysis (Section 4.2) shows that over 99% of GSP errors put the GPU in an error state and lead to job failure if encountered. A GSP error requires a GPU reset or a node reboot to recover, introducing significant overheads. In addition, PMU SPI communication errors propagate downstream and cause MMU errors with a probability of 0.88, leading to MMU errors, which then, in turn, result in a job failure over 90% of the time. Although PMU SPI communication errors are high-impact errors, they are not highlighted in NVIDIA's developer's manual [33].

⁸We calculated the per GB MTBE by multiplying the per node MTBE with the total memory of all GPUs in GB for that node. For example, in a 4-way H100 GPU node with 96 GB of memory, the per GB MTBE = 22,192 × 96 × 4 = 8,521,728 hours.

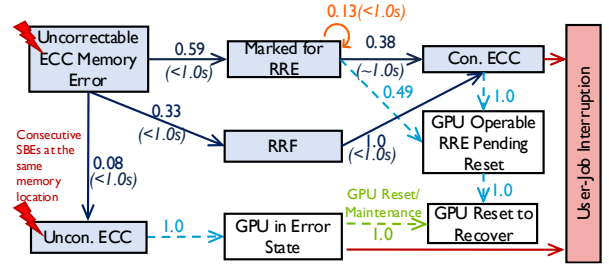


Figure 5: Intra-GPU uncorrectable memory error recovery paths in H100 GPUs. Numbers on the edges show propagation probability. The precise sub-second timing information is not available in the H100 nodes' system logs. All propagation time for H100 GPUs memory errors were within one second.

(ii) Despite error detection mechanisms such as CRC and recovery mechanisms such as message retransmitting, NVLink GPU interconnect errors are frequent (1,922 total NVLink errors) on A100 GPUs, with a node MTBE of 1,185 hours (system-wide MTBE of 11 hours); 42% (801) of the NVLink errors affected two or more GPUs. Although NVIDIA [33] indicates that a GPU reset or a node reboot is required to clear NVLink, *Delta* SREs suggested that NVLink errors were largely benign and did not always lead to job failures. Indeed, the user job impact analysis shows that an NVLink error has a 54% chance of leading to a job failure (see Section 5). When an NVLink error is observed but does not affect a user job, it may be because NVLink is primarily used for communication rather than computation in many jobs or because multiple NVLink errors within the same job have a consolidated impact, resulting in high occurrence rates but minimal application disruption.

(iii) H100 GPUs have substantially improved GPU hardware resilience over A100 GPUs: only 3 GSP errors were observed during the measurement period. Moreover, the number of GPU Fallen off the bus errors, NVLink errors⁹, and PMU SPI errors were not observed in H100 GPUs. We attribute this improvement to (a) the better packaging and tight of CPU and GPU into a single heterogeneous compute module, which significantly reduces integration errors and enhances the resilience of the CPU-GPU complex, and (b) the NVIDIA driver upgrades appear to have improved GSP and NVLink stability (also observed by the *Delta* SREs).

GPU Failure and Replacement. During the measurement period (895 days) of A100 GPUs, four A100 GPUs were replaced due to GPU's failure to boot-up. In comparison, during the measurement period (146 days) of the H100 GPUs, two GPUs were replaced due to uncorrectable memory errors and row remapping failures. This provides an additional support for the observed result that memory is a resilience weak link on H100 GPUs, while the A100 GPUs are more susceptible to errors from GPU hardware components.

4.2 GPU Error Propagation

This section describes results on both *intra-GPU* and *inter-GPU* error propagation during the operational period for *Delta*'s A100 and H100 GPUs. Understanding GPU error propagation reveals resilience weak links between GPU components. We break down the error-recovery propagation into three categories for the GPU errors listed in Table 1: (i) GPU memory, (ii) GPU hardware, and (iii) NVLink

⁹We tested the NVLink to ensure that the NVLinks were enabled.

interconnect, and estimate the propagation probabilities from GPU error logs (see Section 3.2). The propagation paths in Figures 5 to 7 are highlighted with lightning signs that indicate the beginning of each path. If two succeeding errors occur in close successions, i.e., a short propagation time, it suggests causality.

For GPU memory errors, we observed propagation primarily on H100 GPUs, which coincides with our findings on worsened memory resilience in Section 4.1. For GPU hardware and NVLink interconnect errors, we observed error propagation paths primarily on A100 GPUs because of the improved hardware resilience of H100 as discussed in Section 4.1.

4.2.1 GPU Memory-Error Propagation. Intra-GPU uncorrectable memory error recovery paths are shown in Figure 5 for H100 GPUs during the operational period. Memory error recovery path for A100 GPUs were a subset shown in Figure 5. Hence, we only show the memory propagation path for H100 GPUs.

Successful Error Recovery and Containment. As shown in Figure 5, row remapping recovery (RRE) triggered by an uncorrectable ECC memory error has a success rate of 0.59 on H100 GPUs. For the 33% of the row remapping events that fail (RRF), the GPU still contains the uncorrectable memory error because only the affected user jobs were terminated, and only the faulty page was offlined. Overall, considering both uncorrectable memory error recovery paths (RRE and error containment after an RRF), the impact of uncorrectable memory errors was alleviated 92% of the time on H100, while the GPU can remained operable. Such uninterrupted operations were not achievable on previous-generation GPUs (e.g., Kepler in [11, 49] and Volta in [36]), as an uncorrectable memory error would immediately cause user job interruption and put GPU in an error state, necessitating a GPU reset to recover [33].

Unsuccessful Error Containment. The above uncorrectable ECC memory error recovery or containment process can fail, resulting in uncontained memory errors (Section 2.1). The error containment process failed 8% of the times on H100 GPUs during the operational period (see Figure 5). Moreover, as informed by the *Delta* SREs, uncontained memory errors can be highly bursty and persistent and may spam the console logs, consuming useful compute cycles and leaving the GPU inoperable. We did not observed highly bursty and persistent uncontained memory error in A100 or H100 GPUs during the operational period.

Overall, the propagation analysis suggests that the new memory error recovery mechanisms on H100 GPUs alleviated impact of uncorrectable memory errors 92% of the time. That said, *Delta*'s H100 GPUs exhibited significantly more memory error recovery propagations than A100 GPUs, highlighting the worsened memory resilience compared to A100s. In addition, the highly bursty and persisting nature of uncontained memory errors can lead to node/system operation disruptions, as we learned from *Delta* SREs.

4.2.2 GPU Hardware-Error Propagation. This section primarily focuses on A100 because the H100 hardware (these components) experienced almost no error during the operational period. Figure 6 shows error propagation across GPU hardware components in A100 GPUs during the operational period. We found three dominant GPU hardware error propagation paths, originating in (i) GSP (GPU System Processor) errors, (ii) PMU (Performance Management Unit) SPI errors, and (iii) GPU Fallen Off the Bus errors. We omit hardware

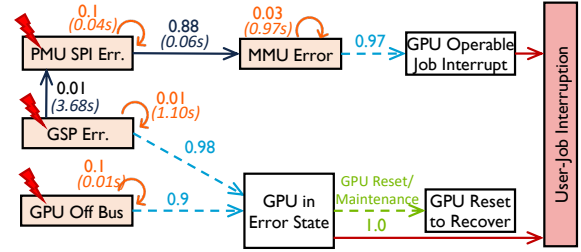


Figure 6: Intra-GPU hardware error propagation probabilities in A100 GPUs. Numbers on the edges show propagation probabilities and average propagation time in seconds.

error propagation graphs for H100 GPUs as we only observed three GSP errors and their propagation paths are similar to that of the A100 GPUs.

GSP-related Errors. Error propagations that originate in GSP-related errors are the most prominent among GPU hardware errors (see Table 1) on A100 GPUs. A GSP error arises when the GSP fails to respond to the remote procedure calls from the GPU driver. Figure 6 shows that, with a probability of 0.99, GSP errors lead to the recurrence of the same error or put the GPU in an inoperable state. The remaining 0.01 (15 cases) of GSP errors caused PMU SPI communication errors (see the follow-up description) that led to user job failure, as depicted in Section 5. Our analysis additionally shows that 99% of GSP errors appeared in isolation without a preceding error.

GSP errors can be caused by either GSP firmware bugs [32] or demanding workload. For example, *Delta* SREs observed that these errors were highly correlated with heavy ML benchmarks, and they suggested that GSP errors are high-impact errors whose recovery requires manual node draining and reboots. Our propagation analysis confirm that the GSP is a single point of failure on both A100 GPUs in part because of their spontaneous nature and high downstream impact (e.g., GPU hangs) on the GPU.

PMU SPI Errors. Communication errors with the performance management unit (PMU) over the Serial Peripheral Interface (SPI), known as *PMU SPI errors*, can cause performance management issues (e.g., inability to change the core frequency). We observed that such errors could lead to MMU errors with a probability of 0.88 (see Figure 6), ultimately leading to user job failures. The majority of the rest of the PMU SPI errors resulted in another PMU SPI error in close succession, leading to persisting error patterns. We observed 77 occurrences of PMU SPI errors (see Table 1) with a 0.98 probability of leading to user job failures (see Table 2) on A100s.

GPU Fallen Off the Bus. GPU Fallen Off the Bus errors were logged when the GPU driver could not reach the GPU over the system bus. This error is an integration error often caused by a loose GPU-motherboard connection or contact failure because of thermal cycles [49]. Over 99% of the errors of this type lead to similar errors in close successions and eventually put the GPU into an error state.

Our *GPU hardware error propagation* analysis suggests that the error detection and recovery of GSP, PMU, and the communication interfaces (e.g., PMU SPI) need to be improved via duplications and error-detection and correction mechanisms to prevent single

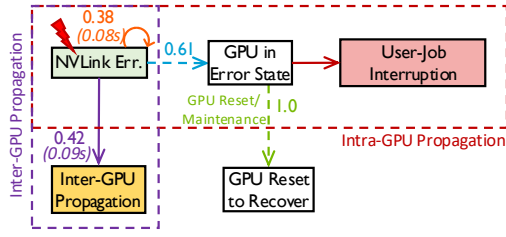


Figure 7: NVLink intra-GPU and inter-GPU error propagation in A100 GPUs. Numbers on the edges show propagation probabilities and average propagation time in seconds.

points of failure, as evident in A100 GPUs. In fact, AWS recommends disabling GSP for stability over performance benefits [4]. By improving GSP hardware and driver software combined with a tightly integrating CPU-GPU architecture, the H100 GPUs in the GH200 Superchips significantly improve GPU hardware resilience over the A100 GPUs. Notably, apart from three GSP errors, which follow the same propagation paths as GSP errors in A100 GPUs, we observed no other hardware error propagation paths on H100 GPUs.

4.2.3 NVLink Interconnect-Error Propagation. NVLink is an GPU-to-GPU interconnect for communication and data exchanges. An NVLink error can impact a single or multiple GPUs on the same node, possibly rendering the entire multi-GPU compute pool unavailable (see Figure 7). We observed both kinds of propagation in our error logs on A100 GPUs during the operational period on A100 GPUs.

NVLink Inter- and Intra-GPU Propagation. An NVLink error occurs when one or more NVLinks experience an error. Our analysis showed that of the 1,922 NVLink errors, 42% propagated to connected GPUs; of those errors, 17% involved three or more GPUs of the same compute node. The rest of the 1,121 NVLink errors did not propagate across GPUs. NVLink errors, like GPU hardware errors, happen in isolation without preceding errors. On the same GPU device, an NVLink error either leads to another NVLink error soon after (with a probability of 0.38) or potentially leaves the faulty GPU in an error state (with a probability of 0.61).

Although *Delta* SREs reported that most NVLink errors were benign, we found that the probability of leading to user job failure when encountered is 54% (43 cases) during the operational period on A100 GPUs. Moreover, in two incidents, GPU resets are needed to recover from critical NVLink errors, leading to over 2 hours of node interruption, and, as with GPU hardware errors, we found no preceding hardware errors before NVLink errors, making them less predictable than memory-related errors.

We did not observe NVLink-related errors on H100 GPUs during the operational period despite observing related non-error NVLink events (e.g., link initialization). We additionally conducted NVLink tests on an *Delta*-AI node to confirm that NVLinks are enabled and fully functional. We conjecture that the improvement is due to both NVLink hardware and driver upgrades and potential changes in NVLink error logging mechanisms.

Incident 1: GSP Error Caused User Job Failure on A100	Incident 2: Contained Uncorrectable Error Caused User Job Failure on H100
<pre>[scheduler][slurmctld][2022-10-11T06:59:19] Job started in 'gpuA100x4' partition on gpus080. [JobId=976760] [gpus080][kernel][2022-10-11T08:23:52] Timeout waiting for RPC from GSP. [ErrorCode=XID 119] [scheduler][slurmctld][2022-10-11T08:23:52] Job completed with error. [JobId=976760, ExitStatus=1]</pre>	<pre>[scheduler][slurmctld][2024-11-11T07:41:33] Job started in 'ghx4' partition on gh083. [JobId=134883] [gh083][kernel][2024-11-11T20:23:22] Contained ECC: SM (0x1). RST: No, D-RST: No [ErrorCode=XID 94] [scheduler][slurmctld][2024-11-11T20:23:29] Job completed with error. [JobId=134883, ExitStatus=1]</pre>
Error	Error
<pre>[scheduler][slurmctld][2022-10-11T09:08:24] Node gpus080 state set to draining. [scheduler][slurmctld][2022-10-12T08:21:36] Node gpus080 rebooted.</pre>	<pre>[gh083][slurmctld][2024-11-11T20:23:55] clean shutdown with force=no halt=no reboot=yes [scheduler][slurmctld][2024-11-11T20:25:34] Node gh083 rebooted.</pre>
Recovery	Recovery

Figure 8: Incidents in which GPU errors led to job failure.

4.3 Examples of Error Propagation to User Jobs

In this section, we show how measured errors relate to NVIDIA’s measured DCGM GPU utilization metric, a step towards relating GPU memory and hardware errors with their impact on the user jobs. Figure 8 presents two example incidents.

Incident 1: A GSP error stalled GPU control functions and rendered the GPU inoperable on an A100 GPU. Consequently, the user job scheduled on that GPU failed. The GSP error required the draining of the node and a full node reboot to recover, which led to the draining of all pending user jobs on that node. From the beginning of the node drain to the completion of the node reboot, the total recovery time for this incident was 23 node hours (09:08 AM to 08:21 AM the next day), during which the node was unavailable. Figure 9a shows that the corresponding GPU utilization dropped quickly following the error incident due to job failure. This incident shows that a GPU error can significantly interrupt user jobs and node availability.

Incident 2: In this incident, the H100 GPU running the user job experienced an uncorrectable ECC memory error, which led to an error containment event. The error containment event contains the uncorrectable memory error by terminating the user’s job. Subsequently, a node reboot is triggered to recover from this error, resulting in a two minute node downtime. Figure 9b shows that GPU utilization dropped as the user job was terminated. This incident shows that although an uncorrectable memory error might be contained, it still results in user job termination and requires node reboot to fully recover, resulting in unexpected node downtime.

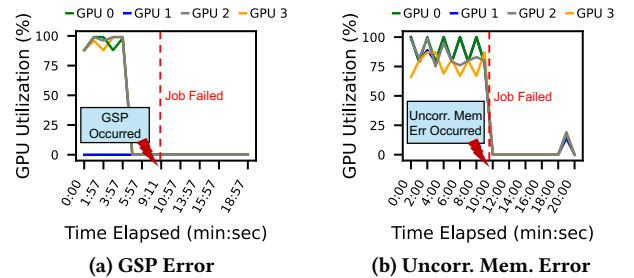


Figure 9: GPU utilization during GPU error incidents.

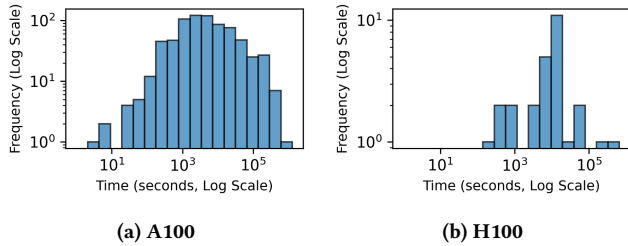


Figure 10: Node unavailability time after GPU failures.

5 Propagation of Errors to Jobs

This section provides an in-depth analysis of job-level resilience and associated GPU downtime.

A GPU error can lead to:

- (1) *Job Failure*: A GPU error may not be handled by jobs, either because the error itself is not contained or because there is a lack of appropriate error-handling mechanisms. For example, GSP errors lead to GPU failures and require node reboots. To recover from GPU failures, jobs need to be re-executed from the beginning or rolled back to the closest checkpoint.
- (2) *GPU and Node Downtime*: Downtime can occur when GPUs need to be reset or replaced by the operator, and no jobs can be scheduled on the GPU and the corresponding node in the interim.

5.1 Result Highlights

Table 2 show the overall impact of hardware errors on applications. In summary, we observe that:

(i) Except for MMU and NVLink errors, no other GPU errors are handled by jobs, thus resulting in their failure. Depending on the GPU type, the underlying cause of job failures differs. In the case of A100 GPUs, hardware errors predominantly lead to job failures, whereas memory errors are the primary cause of H100 GPUs. Therefore, there is a need to improve the resilience of underlying hardware to minimize such failures. While checkpointing is an option, checkpointing routines have a high overhead of up to 40% [29, 53, 54], including management, storage, and restore.

(ii) The overall availability per-GPU node is $\sim 99.4\%$ (corresponding to 9 minutes downtime/day) and $\sim 99.3\%$ (corresponding to 10 minutes downtime/day) for A100 and H100 GPUs respectively. This level of unavailability suggests that even if the rest of the infrastructure is highly available, current GPUs may not provide sufficient availability to meet the demands of critical applications that require greater than 3 9’s of availability (downtime of 1.4 mins/day). In addition, we used emulation to project the impact of this availability distribution at increased scales (for both node scale and job duration) and found that significant overprovisioning of 5% would be necessary to handle associated failures (as explained in further detail in Section 5.4).

5.2 Job Statistics

During the characterization period, 1,420,278 user jobs were submitted to GPU nodes with a success rate of 87%. About 74% of user jobs ran on a single GPU; 24% ran on 2–4 GPUs; and only 2% of jobs used five or more GPUs.

Because of a lack of specific information on whether jobs were ML-related, we estimated the percentage of ML user jobs based

on job submission names and the system modules/libraries imported¹⁰. For instance, user jobs with names containing `model` or `train` were likely related to machine learning. We provide detailed statistics on node hours used, durations for both job types, and failure probabilities in Table 3. We find that across all job sizes, the failure probability lies between 6–49%, indicating a lack of adequate recovery mechanisms.

5.3 Job Failure Analysis

To understand job failure patterns, we first separated jobs into “Completed” and “GPU-Failed” depending on their completion status. Based on that job categorization, we analyzed (i) the GPU errors that were most likely lead to a job failure, and (ii) the potential recovery strategies, such as checkpointing and exception handling. *Classifying Job Runs*: We classified jobs based on their exit status and proximity of job failure time to GPU error occurrence time. The job exit status was obtained from the Slurm job scheduler logs (as described in Section 3.1). We marked a job as “GPU-Failed” if a GPU error occurred within a 20-second interval before job failure.¹¹ *Correlating GPU Errors and Job Failures*: We broke down all GPU-Failed jobs by the specific GPU errors that were most likely to lead to job failures. Table 2 provides probabilities of user job failures per GPU error for the A100 and H100 GPUs. As any of the encountered errors may have contributed to a job failure, we consider all GPU errors that occurred within the 20-second interval to be responsible for the failure.

Overall, other than NVLink and MMU errors, all GPU errors, such as GSP RPC timeout and PMU failures, propagate (as discussed in Section 4.2) and cause job failures. Based on previous analysis, we note that hardware errors – such as GSP and PMU errors are dominant in A100 GPUs, whereas memory errors are dominant in H100 GPUs.

NVLink and MMU errors do not necessarily lead to job failures because:

¹⁰Due to privacy restrictions, we could not access the job scripts for use in job classification.

¹¹Note that we do not count “zombie” jobs, i.e. jobs that have failed but not terminated by the Slurm scheduler in the analysis.

XID	GPU Error	# GPU-failed Jobs		# Jobs Encountering Given XID		Job Failure Probability (%)	
		A100	H100	A100	H100	A100	H100
31	MMU err.	3206	93	3543	126	90.48	73.80
74	NVL err.	43	0	80	0	53.75	–
122	SPI PMU RPC failure	40	0	41	0	97.56	–
119	GSP RPC timeout	31	0	31	0	100.00	–
94	Contained ECC	5	5	5	5	100.00	100.00
48	GPU DBE	0	5	0	5	–	100.00
64	Row remapping failed	0	3	0	3	–	100.00
63	Row remapping event	0	2	0	2	–	100.00

Table 2: Distribution of GPU-failed jobs across the different GPU error types for A100 and H100 GPUs. The failure probability is calculated as (# GPU-failed jobs encountering that GPU error) / (# jobs encountering that GPU error). The total number of GPU-failed jobs was 3,359 during the 895-day characterization period.

GPU Count	Count (%)	Elapsed Time (Minutes)		Failed (%)	GPU Hours (k)	
		Mean	P99		ML	Non-ML
1	1,052,993 (74.140%)	134.560	2859.677	129,395 (12.29%)	641.9	1,949.14
2-4	337,637 (23.773%)	159.839	2880.117	40,610 (12.03%)	810.4	2,749.00
5-8	13,907 (0.979%)	231.927	2880.316	4,469 (32.13%)	253.1	290.16
9-32	13,378 (0.942%)	221.636	2880.167	3,045 (22.76%)	307.5	844.40
33-64	1,375 (0.097%)	145.206	2880.017	608 (44.22%)	108.2	149.24
65-128	856 (0.060%)	320.185	2834.413	421 (49.18%)	15.1	442.70
129-256	110 (0.008%)	174.713	2041.312	7 (6.36%)	0.0	57.23
257+	22 (0.002%)	29.128	107.493	5 (22.73%)	0.0	3.44

Table 3: Job distribution, elapsed time statistics (mean, P99), failure count with percentage, and GPU hours divided into ML and non-ML categories for various A100 and H100 GPU configurations.

- (1) For *NVLink errors*, the link or GPU may not be in use by any user jobs (as discussed in Section 4.2)¹².
- (2) For *MMU errors*, there can be application or library-level masking mechanisms. Besides hardware errors, MMU errors can also occur if buggy user code makes illegal memory accesses that cannot be mapped in the virtual-to-physical address space. Such errors can be managed using appropriate application-level exception handlers. Popular libraries and frameworks for machine learning [40–42] have support for handling such exceptions by skipping the associated training iteration, albeit at the cost of model quality.

5.4 Impact of GPU Downtime on Jobs

While significant node hours might be lost because of wasted compute time from failed jobs, additional node hours are also lost because of the time required to recover the impacted GPU node by either resetting it or replacing it entirely. To reset the GPU node, operators typically drain the node, i.e., wait for other jobs running on the node to complete without accepting new jobs and then reboot. After the reboot, if the node successfully passes the health check, the node reset is successful, and new jobs can be scheduled on the GPU node. If the reset is unsuccessful, the node is marked failed until the GPU is additionally tested and physically replaced, if required. To calculate the average system downtime, we estimated the total time when the GPU was unavailable, which primarily included the drain and reboot time. Figure 10 shows the distribution of the unavailable time across the entire characterization duration. Overall, we found that the expected time to service the failed node was 0.88 hours for A100 GPUs and 2.2 hours for H100 GPUs. A total of 5,700 node hours were lost to GPU downtime. Using the node downtime and failure distributions, we can estimate the availability of the GPU node as $\frac{MTTF}{MTTF+MTTR}$ equal to 99.4% and 99.3% for A100 and H100 GPUs respectively¹³.

Projected impact of availability on long-running and large-scale jobs: We provide error and recovery statistics in previous sections; we also attempted to project how those distributions would affect jobs running on a different system. To do so, we built a simulation tool driven by our analysis. The parameters of the simulation tool can be varied based on the scenario under consideration.

¹²Based on our understanding, NVLink and memory errors can occur even when no workload is running [15].

¹³The node MTTF number is estimated from the GPU’s MTBE, for which we conservatively assume that all critical GPU errors lead to node interruption.

Specifically, we simulate the case in which jobs (such as ML training) use the entire set of 608 H100 GPUs and run for a duration of 1 month. These jobs require all GPUs to be operational to make progress, and frequent node failures can lead to resource unavailability and slower job progress. When such failures occur, additional provisioning of GPU resources is necessary to allow the job to resume on alternate nodes while the failed nodes recover.

The simulation uses a discrete time event simulation with node failure probabilities derived from our prior analysis. The recovery time after a failure is dependent on variables such as checkpoint load time and availability of spare GPUs. To account for the variability introduced by these factors, we parameterize recovery time and perform a parameter sweep. For a training job with 608 GPUs and recovery time of 2.2 hours, the required overprovisioning is 5%: i.e., 31 additional GPUs are needed beyond the original 608 to maintain availability of 99.9% at the job level. While at first glance, such overprovisioning would appear to be a small cost, for the above example, it would cost over \$1 million per month for a 1000 node cluster (our analysis is based on AWS H100 GPU rental rates [9]). However, if the recovery time is reduced to 5 minutes, downtime decreases significantly, and the required overprovisioning drops to 2%, a 2.5× reduction. This highlights the criticality of minimizing recovery time to reduce downtime for large and long-running jobs.

In summary, every GPU node in the system has “two nines” of availability. While such availability does not significantly impact small jobs that use recovery mechanisms, large jobs can face significant downtime. Significant overprovisioning of up to 5% is required to eliminate such downtime.

6 Related Work

Existing work has analyzed GPU resilience at the microarchitecture, system, and application levels. This paper extends previous work via comprehensive analyses of GPU error characteristics, failure propagation, and their impact on user jobs.

Microarchitecture-level GPU Resilience. Previous research [20, 44, 52, 57] has primarily focused on the resilience of individual GPUs at the microarchitecture and software levels, e.g., for older generations such as the NVIDIA G80 [6]. However, the earlier work did not evaluate the resiliency of modern GPUs in large-scale HPC settings.

System-level GPU Resilience in HPC Settings. Existing studies have analyzed the resilience of GPUs in HPC systems [7, 10, 16, 19, 24], for example, the NVIDIA Tesla K20X GPUs in various supercomputers [11, 17, 18, 30, 31, 37, 49, 50]. Studies of the Blue Waters, Titan, and Summit supercomputers [11, 24, 30, 31, 36] have examined node failures and GPU error characteristics. However, such work either studied previous generations of GPUs with a focus on GPU memory errors, or focus on cluster-level resilience instead of GPU resilience. Our work complements those efforts by providing resilience insights into the latest generation of GPUs, focusing on a broader range of components.

Application-level GPU Resilience in Data-centers and Deep Learning Workloads. Recent research has focused on understanding GPU power usage [31], GPU component-level failures [23], software-level error handling [14, 27, 28], and the impact of GPU software error propagation on GPGPU applications [3, 26, 43, 55] and emerging GPU workloads such as convolutional neural networks (CNNs) [8,

12], large language models (LLMs) [13, 22], and safety-critical applications [38, 45].

7 Discussion

Justification of Analyzing Errors: Data-driven HPC resilience characterization studies analyze operational data on system/application errors to provide insights into system resilience [10, 11, 17, 18, 30, 31, 36, 37, 49, 50]. While in those studies, the error rate is used as the key metric to quantify resilience, some [5, 48] argue that fault rate is a more appropriate metric.

Errors represent the manifestation of faults and have direct downstream consequences, such as triggering of recovery mechanisms, application interruptions, or system-wide outages (SWOs). While a fault may result in multiple errors, it is the resulting errors that the recovery mechanisms must address to maintain system health. These errors and their recovery process directly impact system health, performance, and availability. Thus, SREs prioritize errors over faults. Hence, like many others who study operational data, we chose to study errors.

Potential errors in NVIDIA driver logging: A potential source of error is the NVIDIA driver’s inconsistent logging of different XID errors. The NVIDIA driver is proprietary software whose source code is not publicly available, making exact specifications of its error-logging mechanism unverifiable. We attempted to minimize the impact of the inconsistency, to the best of our ability, through consistent log data collection and error coalescing.

8 Conclusion

This paper describes the results of a resilience study of *Delta*, which consists of 1,056 NVIDIA A100 and H100 GPUs. The study used up to 2.5 years of operational data on GPU errors collected across those GPUs. We assessed the resilience of GPU components, error propagation paths, and impact on jobs and compared the two generations of GPUs. In future work, we will extend the analysis presented to other accelerators and larger-scale systems, running more complex HPC and ML workloads.

References

- [1] Md Mahbub Alam, Luis Torgo, and Albert Bifet. 2022. A survey on spatio-temporal data analytics systems. *Comput. Surveys* 54, 10s (2022), 1–38.
- [2] Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. 2023. Gemini: A family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805* (2023).
- [3] Abdul Rehman Anwer, Guanpeng Li, Karthik Pattabiraman, Michael Sullivan, Timothy Tsai, and Siva Kumar Sastry Hari. 2020. GPU-Trident: Efficient modeling of error propagation in GPU programs. In *SC20: International Conference for High Performance Computing, Networking, Storage and Analysis*. IEEE, 1–15.
- [4] AWS Official. 2023. *How do I troubleshoot Xid errors on my NVIDIA GPU-accelerated EC2 Linux instance?* <https://repost.aws/knowledge-center/ec2-linux-troubleshoot-xid-errors>
- [5] Majed Valad Beigi, Yi Cao, Sudhanva Gurumurthi, Charles Recchia, Andrew Walton, and Vilas Sridharan. 2023. A Systematic Study of DDR4 DRAM Faults in the Field. In *2023 IEEE International Symposium on High-Performance Computer Architecture (HPCA)*. 991–1002. <https://doi.org/10.1109/HPCA56546.2023.10071066>
- [6] Giani Braga, Marcio M Gonçalves, and José Rodrigo Azambuja. 2023. Software-controlled pipeline parity in GPU architectures for error detection. *Microelectronics Reliability* 148 (2023), 115155.
- [7] Nevin Cini and Gulay Yalcin. 2020. A Methodology for Comparing the Reliability of GPU-Based and CPU-Based HPCs. *ACM Comput. Surv.* 53, 1, Article 22 (Feb. 2020), 33 pages. <https://doi.org/10.1145/3372790>
- [8] Josie E Rodriguez Condia, Juan-David Guerrero-Balaguera, Fernando F Dos Santos, Matteo Sonza Reorda, and Paolo Rech. 2022. A multi-level approach to evaluate the impact of GPU permanent faults on CNN’s reliability. In *2022 IEEE International Test Conference (ITC)*. IEEE, 278–287.

- [9] DataCrunch. 2024. Cloud GPU Pricing Comparison in 2025. <https://datacrunch.io/blog/cloud-gpu-pricing-comparison> Accessed: 2025-04-14.
- [10] Nathan DeBardleben, Sean Blanchard, Laura Monroe, Phil Romero, Daryl Grunau, Craig Idler, and Cornell Wright. 2014. GPU behavior on a large HPC cluster. In *Euro-Par 2013: Parallel Processing Workshops: BigDataCloud, DIHC, FedCL, HeteroPar, HiBB, LSDVE, MHPC, OMHI, PADABS, PROPER, Resilience, ROME, and UCHPC 2013, Aachen, Germany, August 26-27, 2013. Revised Selected Papers 19*. Springer, 680–689.
- [11] Catello Di Martino, Zbigniew Kalbarczyk, Ravishankar K. Iyer, Fabio Baccanico, Joseph Fullop, and William Kramer. 2014. Lessons Learned from the Analysis of System Failures at Petascale: The Case of Blue Waters. In *2014 44th Annual IEEE/IFIP International Conference on Dependable Systems and Networks*. 610–621. <https://doi.org/10.1109/DSN.2014.62>
- [12] Fernando Fernandes dos Santos, Pedro Foletto Pimenta, Caio Lunardi, Lucas Draghetti, Luigi Carro, David Kaeli, and Paolo Rech. 2018. Analyzing and increasing the reliability of convolutional neural networks on GPUs. *IEEE Transactions on Reliability* 68, 2 (2018), 663–677.
- [13] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The Llama 3 herd of models. *arXiv preprint arXiv:2407.21783* (2024).
- [14] Bo Fang, Karthik Pattabiraman, Matei Ripeanu, and Sudhanva Gurumurthi. 2014. GPU-Qin: A methodology for evaluating the error resilience of GPGPU applications. In *2014 IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS)*. IEEE, 221–230.
- [15] NVIDIA Forums. 2024. NVLink error 74 fatal error detected. <https://forums.developer.nvidia.com/t/nvlink-error-74-fatal-error-detected/55582/1> Accessed: 2025-02-24.
- [16] Haryadi S Gunawi, Riza O Suminto, Russell Sears, Casey Golliver, Swaminathan Sundararaman, Xing Lin, Tim Emami, Weiguang Sheng, Nematollah Bidokhti, Caitie McCaffrey, et al. 2018. Fail-slow at scale: Evidence of hardware performance faults in large production systems. *ACM Transactions on Storage (TOS)* 14, 3 (2018), 1–26.
- [17] Saurabh Gupta, Tirthak Patel, Christian Engelmann, and Devesh Tiwari. 2017. Failures in large scale systems: long-term measurement, analysis, and implications. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*. 1–12.
- [18] Saurabh Gupta, Devesh Tiwari, Christopher Jantzi, James Rogers, and Don Maxwell. 2015. Understanding and exploiting spatial properties of system failures on extreme-scale HPC systems. In *2015 45th Annual IEEE/IFIP International Conference on Dependable Systems and Networks*. IEEE, 37–44.
- [19] Imran S Haque and Vijay S Pande. 2010. Hard data on soft errors: A large-scale assessment of real-world error rates in GPGPU. In *2010 10th IEEE/ACM International Conference on Cluster, Cloud and Grid Computing*. IEEE, 691–696.
- [20] Siva Kumar Sastry Hari, Timothy Tsai, Mark Stephenson, Stephen W Keckler, and Joel Emer. 2015. SASSIFI: Evaluating resilience of GPU applications. In *Proceedings of the Workshop on Silicon Errors in Logic-System Effects (SELSE)*.
- [21] Chip Huyen. 2022. *Designing Machine Learning Systems*. O’Reilly Media, Inc.
- [22] Ziheng Jiang, Haibin Lin, Yinmin Zhong, Qi Huang, Yangrui Chen, Zhi Zhang, Yanghua Peng, Xiang Li, Cong Xie, Shibiao Nong, et al. 2024. MegaScale: Scaling large language model training to more than 10,000 GPUs. In *21st USENIX Symposium on Networked Systems Design and Implementation (NSDI 24)*. 745–760.
- [23] Jeageun Jung and Mattan Erez. 2023. Predicting Future-System Reliability with a Component-Level DRAM Fault Model. In *Proceedings of the 56th Annual IEEE/ACM International Symposium on Microarchitecture*. 944–956.
- [24] Apostolos Kokolis, Michael Kuchnik, John Hoffman, Adithya Kumar, Parth Malani, Faye Ma, Zachary DeVito, Shubho Sengupta, Kalyan Saladi, and Carole-Jean Wu. 2025. Revisiting Reliability in Large-Scale Machine Learning Research Clusters. In *2025 IEEE International Symposium on High Performance Computer Architecture (HPCA)*. IEEE Computer Society, Los Alamitos, CA, USA, 1259–1274. <https://doi.org/10.1109/HPCA61900.2025.00096>
- [25] Baolin Li, Yankai Jiang, Vijay Gadepally, and Devesh Tiwari. 2024. LLM inference serving: Survey of recent advances and opportunities. *arXiv preprint arXiv:2407.12391* (2024).
- [26] Guanpeng Li, Karthik Pattabiraman, Chen-Yang Cher, and Pradip Bose. 2016. Understanding Error Propagation in GPGPU Applications. In *SC ’16: Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*. 240–251. <https://doi.org/10.1109/SC.2016.20>
- [27] Xinyu Lian, Sam Ade Jacobs, Lev Kurilenko, Masahiro Tanaka, Stas Bekman, Olatunji Ruwase, and Minjia Zhang. 2025. Universal Checkpointing: A Flexible and Efficient Distributed Checkpointing System for Large-Scale DNN Training with Reconfigurable Parallelism. In *Proceedings of the 2025 USENIX Annual Technical Conference*.
- [28] Naoya Maruyama, Akira Nukada, and Satoshi Matsuoka. 2010. A high-performance fault-tolerant software framework for memory on commodity GPUs. In *2010 IEEE International Symposium on Parallel & Distributed Processing (IPDPS)*. IEEE, 1–12.
- [29] Avinash Maurya, Robert Underwood, M Mustafa Rafique, Franck Cappello, and Bogdan Nicolae. 2024. DataStates-LLM: Lazy asynchronous checkpointing for

- large language models. In *Proceedings of the 33rd International Symposium on High-Performance Parallel and Distributed Computing*. 227–239.
- [30] Bin Nie, Devesh Tiwari, Saurabh Gupta, Evgenia Smirni, and James H Rogers. 2016. A large-scale study of soft-errors on GPUs in the field. In *2016 IEEE International Symposium on High Performance Computer Architecture (HPCA)*. IEEE, 519–530.
- [31] Bin Nie, Ji Xue, Saurabh Gupta, Christian Engelmann, Evgenia Smirni, and Devesh Tiwari. 2017. Characterizing temperature, power, and soft-error behaviors in data center systems: Insights, challenges, and opportunities. In *2017 IEEE 25th International Symposium on Modeling, Analysis, and Simulation of Computer and Telecommunication Systems (MASCOTS)*. IEEE, 22–31.
- [32] NVIDIA. 2023. *NVIDIA Data Center GPU Driver version 525.105.17 (Linux) / 528.89 (Windows)*. https://docs.nvidia.com/datacenter/tesla/pdf/NVIDIA_Data_Center_GPU_Driver_Release_Notes_525_v3.0.pdf
- [33] NVIDIA. 2024. *GPU Deployment and Management*. https://docs.nvidia.com/deploy/pdf/XID_Errors.pdf
- [34] NVIDIA Corporation. 2023. *NVIDIA Grace Hopper Superchip Architecture Whitepaper*. <https://resources.nvidia.com/en-us-grace-cpu/nvidia-grace-hopper>
- [35] NVIDIA Corporation. 2024. *NVIDIA GPU Memory Error Management*. <https://docs.nvidia.com/deploy/a100-gpu-mem-error-mgmt/index.html> Release r555.
- [36] Vladyslav Oles, Anna Schmedding, George Ostrouchov, Woong Shin, Evgenia Smirni, and Christian Engelmann. 2024. Understanding GPU Memory Corruption at Extreme Scale: The Summit Case Study. In *Proceedings of the 38th ACM International Conference on Supercomputing*. 188–200.
- [37] George Ostrouchov, Don Maxwell, Rizwan A Ashraf, Christian Engelmann, Mallikarjun Shankar, and James H Rogers. 2020. GPU lifetimes on Titan supercomputer: Survival analysis and reliability. In *SC20: International Conference for High Performance Computing, Networking, Storage and Analysis*. IEEE, 1–14.
- [38] Jon Perez-Cerrolaza, Jaume Abella, Leonidas Kosmidis, Alejandro J Calderon, Francisco Cazorla, and Jose Luis Flores. 2022. GPU devices for safety-critical systems: A survey. *Comput. Surveys* 55, 7 (2022), 1–37.
- [39] James C Phillips, Rosemary Braun, Wei Wang, James Gumbart, Emad Tajkhorshid, Elizabeth Villa, Christophe Chipot, Robert D Skeel, Laxmikant Kale, and Klaus Schulten. 2005. Scalable molecular dynamics with NAMD. *Journal of Computational Chemistry* 26, 16 (2005), 1781–1802.
- [40] PyTorch Contributors. 2017. *PyTorch Issue #1137: How to Handle Corrupted Data in Dataloader*. <https://github.com/pytorch/pytorch/issues/1137> Accessed: 2024-11-26.
- [41] PyTorch Contributors. 2018. *Discussion: How to Handle Exception in DistributedDataParallel*. <https://discuss.pytorch.org/t/how-to-handle-exception-in-distributeddataparallel/42026> Accessed: 2024-11-26.
- [42] PyTorch Lightning Contributors. 2021. *PyTorch Lightning Discussion: Handling Exceptions in Forward Pass*. <https://github.com/Lightning-AI/pytorch-lightning/discussions/15188> Accessed: 2024-11-26.
- [43] Paolo Rech, Laécio Lima Pilla, Philippe Olivier Alexandre Navaux, and Luigi Carro. 2014. Impact of GPUs parallelism management on safety-critical and HPC applications reliability. In *2014 44th Annual IEEE/IFIP International Conference on Dependable Systems and Networks*. IEEE, 455–466.
- [44] Dimitris Sartzetakis, George Papadimitriou, and Dimitris Gizopoulos. 2022. gpuFI-4: A microarchitecture-level framework for assessing the cross-layer resilience of Nvidia GPUs. In *2022 IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS)*. IEEE, 35–45.
- [45] Deval Shah, Zi Yu Xue, Karthik Pattabiraman, and Tor M Aamodt. 2024. Characterizing and Improving Resilience of Accelerators to Memory Errors in Autonomous Robots. *ACM Transactions on Cyber-Physical Systems* 8, 3 (2024), 1–33.
- [46] Gilad Shainer, Tong Liu, John Michalakes, Jacob Liberman, Jeff Layton, Onur Celebioglu, Scot A Schultz, Joshua Mora, and David Cownie. 2009. Weather research and forecast (WRF) model performance and profiling analysis on advanced multi-core HPC clusters. In *10th LCI International Conference on High-Performance Clustered Computing*.
- [47] David Skinner and William Kramer. 2005. Understanding the causes of performance variability in HPC workloads. In *Proceedings of the IEEE International Workshop/Symposium on Workload Characterization*. IEEE, 137–149.
- [48] Vilas Sridharan, Nathan DeBardeleben, Sean Blanchard, Kurt B. Ferreira, Jon Stearley, John Shalf, and Sudhanva Gurumurthi. 2015. Memory Errors in Modern Systems: The Good, The Bad, and The Ugly. *SIGPLAN Not.* 50, 4 (March 2015), 297–310. <https://doi.org/10.1145/2775054.2694348>
- [49] Devesh Tiwari, Saurabh Gupta, George Gallarno, Jim Rogers, and Don Maxwell. 2015. Reliability lessons learned from GPU experience with the Titan supercomputer at Oak Ridge Leadership Computing Facility. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*. 1–12.
- [50] Devesh Tiwari, Saurabh Gupta, James Rogers, Don Maxwell, Paolo Rech, Sudharshan Vazhkudai, Daniel Oliveira, Dave Londo, Nathan DeBardeleben, Philippe Navaux, et al. 2015. Understanding GPU errors on large-scale HPC systems and the implications for system design and operation. In *2015 IEEE 21st International Symposium on High Performance Computer Architecture (HPCA)*. IEEE, 331–342.
- [51] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971* (2023).
- [52] Alessandro Vallerio, Sotiris Tselonis, Dimitris Gizopoulos, and Stefano Di Carlo. 2018. Multi-faceted microarchitecture level reliability characterization for NVIDIA and AMD GPUs. In *2018 IEEE 36th VLSI Test Symposium (VTS)*. IEEE, 1–6.
- [53] Borui Wan, Mingji Han, Yiyao Sheng, Yanghua Peng, Haibin Lin, Mofan Zhang, Zhichao Lai, Menghan Yu, Junda Zhang, Zuquan Song, et al. 2024. ByteCheckpoint: A Unified Checkpointing System for Large Foundation Model Development. *arXiv preprint arXiv:2407.20143* (2024).
- [54] Long Wang, Karthik Pattabiraman, Zbigniew Kalbarczyk, Ravishankar K Iyer, Lawrence Votta, Christopher Vick, and Alan Wood. 2005. Modeling coordinated checkpointing for large-scale supercomputers. In *2005 International Conference on Dependable Systems and Networks (DSN'05)*. IEEE, 812–821.
- [55] Xiaohui Wei, Hengshan Yue, Shang Gao, Lina Li, Ruyi Zhang, and Jingweijia Tan. 2020. G-SEAP: Analyzing and characterizing soft-error aware approximation in GPGPUs. *Future Generation Computer Systems* 109 (2020), 262–274.
- [56] Ronglong Wu, Shuyue Zhou, Jiahao Lu, Zhirong Shen, Zikang Xu, Jiwu Shu, Kunlin Yang, Feilong Lin, and Yiming Zhang. 2024. Removing Obstacles before Breaking Through the Memory Wall: A Close Look at HBM Errors in the Field. In *2024 USENIX Annual Technical Conference (USENIX ATC 24)*. USENIX Association, Santa Clara, CA, 851–867. <https://www.usenix.org/conference/atc24/presentation/wu-ronglong>
- [57] Lishan Yang, George Papadimitriou, Dimitris Sartzetakis, Adwait Jog, Evgenia Smirni, and Dimitris Gizopoulos. 2024. GPU Reliability Assessment: Insights Across the Abstraction Layers. In *2024 IEEE International Conference on Cluster Computing (CLUSTER)*. IEEE, 1–13.
- [58] Andy B. Yoo, Morris A. Jette, and Mark Grondona. 2003. SLURM: Simple Linux Utility for Resource Management. In *Job Scheduling Strategies for Parallel Processing*, Dror Feitelson, Larry Rudolph, and Uwe Schwiegelshohn (Eds.). Springer Berlin Heidelberg, 44–60.