

INTELLIGENCE BRIEFING

Security Command Center

TLP:CLEAR

2026-04-12 18:29 UTC

# GPUBreach: Rowhammer-Style Attack on GDDR6 Memory Enables Host Privilege Escalation

SECURITY ANALYSIS | HIGH | CVSS 7.8

SCC Item ID	SCC-STY-2026-0058
Type	Security Analysis
Severity	HIGH
CVSS Base Score	7.8
Affected Products	Shared GPU infrastructure using GDDR6 memory (NVIDIA GPUs); CUDA-based workloads on multi-tenant systems
Published	2026-04-10
Discovery Source	Gemini

## Executive Summary

According to secondary reporting, security researchers have claimed to demonstrate GPUBreach, a hardware-level attack that exploits bit-flip vulnerabilities in GDDR6 memory to break the assumed isolation between GPU workloads and the host operating system, potentially enabling full CPU-level privilege escalation from an unprivileged CUDA process. As of this report, these technical claims have not been independently verified by NVIDIA or published in peer-reviewed research. The attack would pose direct risk to shared GPU infrastructure, including cloud-based AI training platforms, analytics environments, and multi-tenant research systems, where workload isolation is the primary security boundary. If confirmed, this finding would signal that the rapid expansion of GPU-based compute, particularly in AI and cloud environments, has outpaced the security assumptions those environments were built on, and hardware-layer attacks on accelerator memory would represent an emerging and underexamined threat surface.

## Technical Analysis

GPUBreach, as described in available reporting, would adapt the Rowhammer attack technique - historically associated with DRAM manipulation on CPUs - to GDDR6 memory used in modern NVIDIA GPUs. Rowhammer exploits the physical proximity of memory rows: by repeatedly accessing ('hammering') specific memory locations at high frequency, an attacker induces electrical interference that causes bit-flips in adjacent rows, corrupting values the attacker does not directly address. GPUBreach would apply this principle to the GPU memory subsystem, where an unprivileged CUDA workload serves as the hammering agent. The claimed attack

works in three stages: (1) an attacker with access to a shared GPU environment launches a CUDA workload designed to hammer targeted GDDR6 memory rows; (2) induced bit-flips corrupt memory structures in adjacent rows; (3) the corrupted state is leveraged to escalate from GPU-level execution to host OS privilege. The critical implication would be the boundary crossing: the attack does not stay within the GPU's memory space but propagates upward to affect host CPU execution context. The defensive gap exploited is structural: shared GPU infrastructure commonly treats GPU-to-host isolation as an assumed architectural guarantee rather than an actively enforced security control. Multi-tenant environments, cloud GPU instances, and shared AI training clusters operate on the premise that CUDA workloads are sandboxed. If confirmed, GPUBreach would challenge that premise directly. If confirmed, this would map to MITRE ATT&CK T1068 (Exploitation for Privilege Escalation), T1203 (Exploitation for Client Execution in the context of workload execution environments), and T1548 (Abuse Elevation Control Mechanism). The associated CWEs - CWE-119 (Improper Restriction of Operations within Memory Bounds), CWE-284 (Improper Access Control), and CWE-1220 (Insufficient Granularity of Access Control) - reflect both the memory corruption mechanism and the access control failure that would allow host-level impact. Important verification note: As of this report, no CVE has been publicly assigned, no vendor confirmation from NVIDIA has been identified, and the available sources are secondary press coverage at Tier 3. Peer-reviewed or vendor-confirmed technical details have not been independently verified. Organizations should treat the technical specifics as preliminary and monitor for authoritative disclosure before taking action beyond awareness and architecture review.

## Action Checklist

1. Step 1: Assess exposure. Determine whether your organization operates shared GPU infrastructure using NVIDIA GPUs with GDDR6 memory, particularly in multi-tenant cloud, AI training, analytics, or research environments where untrusted or third-party CUDA workloads execute alongside privileged processes.
2. Step 2: Review controls. Audit whether GPU-to-host isolation is enforced as an active security control or assumed as an architectural guarantee; verify whether hypervisor or container-level controls (e.g., GPU passthrough configurations, IOMMU enforcement, workload sandboxing) are in place and correctly configured.
3. Step 3: Update threat model. Incorporate hardware-layer memory manipulation attacks on GPU accelerators as an explicit threat vector; add multi-tenant GPU environments to your shared-resource privilege escalation threat scenarios and map against T1068 and T1548.
4. Step 4: Communicate findings. Brief infrastructure and cloud security leads on the GPU isolation assumption gap; frame the risk specifically around AI training clusters, GPU cloud instances, and any environment where external or untrusted workloads share GPU resources with privileged processes.
5. Step 5: Monitor for authoritative disclosure. Track for NVIDIA security advisory, CVE assignment, and peer-reviewed publication. If NVIDIA confirms the attack technique and assigns a CVE, escalate to incident response and patch planning. Until confirmation, continue with architecture review and control validation. Subscribe to NVIDIA Product Security and CISA advisories for follow-up disclosure.

## IR / Forensic Enrichment

Triage Priority

URGENT

<b>Escalation Criteria</b>	Escalate immediately to CISO and infrastructure leadership if nvidia-smi or audit logs reveal an unprivileged CUDA process that has achieved <code>uid=0</code> on a GPU host, if NVIDIA publishes a CVE with a CVSS score $\geq 7.0$ for GDDR6 memory integrity, or if the organization operates AI training infrastructure subject to SOC 2, FedRAMP, or HIPAA where GPU workload co-tenancy with sensitive data processing creates a breach notification exposure.
<b>Recovery Notes</b>	If exploitation is confirmed or suspected, GPU host systems should be rebuilt from known-good images rather than patched in place, as a successful Rowhammer-style escalation to root-level on the host provides an adversary with the capability to persist in firmware or bootloader — standard OS-level remediation is insufficient. Following rebuild, re-validate IOMMU enforcement and NVIDIA driver version against any available NVIDIA advisory before returning nodes to multi-tenant service. Monitor rebuilt hosts for 30 days using <code>auditd</code> rules targeting <code>setuid/setreuid</code> syscalls from CUDA-affiliated processes, and retain those logs for a minimum of 90 days to support any downstream forensic or regulatory inquiry.
<b>Forensic Artifacts</b>	Linux <code>auditd</code> logs for <code>setuid/setreuid/setresuid</code> syscalls ( <code>audit.log</code> , typically <code>/var/log/audit/audit.log</code> ) — a successful GPUBreach privilege escalation from an unprivileged CUDA process to host root would produce a <code>uid</code> transition event here that is anomalous relative to the baseline captured in Step 3   <code>nvidia-smi</code> process accounting output ( <code>'nvidia-smi pmon -s u -d 1'</code> captured over time) — records which CUDA PIDs held GPU memory handles and compute resources at the time of a suspected exploitation event, establishing workload co-tenancy and identifying the source tenant process   Kernel ring buffer ( <code>dmesg</code> ) and <code>/var/log/kern.log</code> — GDDR6 bit-flip events causing memory corruption that crosses isolation boundaries may produce IOMMU fault messages, PCIe AER (Advanced Error Reporting) errors, or NVIDIA driver ECC (Error Correction Code) uncorrectable error events that are logged to the kernel ring buffer   <code>/proc//maps</code> for all CUDA workload processes — captures GPU device memory mappings (entries referencing <code>/dev/nvidia*</code> or <code>/dev/nvidiactl</code> ) for each CUDA process at the time of incident, documenting which processes had direct GPU memory access and the virtual address ranges involved   NVIDIA driver ECC error logs via <code>'nvidia-smi --query-gpu=ecc.errors.uncorrected.volatile.total --format=csv'</code> and <code>'nvidia-smi -q -d ECC'</code> — uncorrectable ECC errors on GDDR6 memory are a hardware-level artifact of successful Rowhammer-style bit-flip activity and represent the closest available forensic indicator of the GPUBreach memory manipulation stage prior to privilege escalation

**Per-Action IR Details**

**Step 1: Assess exposure — determine whether your organization operates shared GPU infrastructure using NVIDIA GPUs with GDDR6 memory, particularly in multi-tenant cloud, AI training, analytics, or research environments where untrusted or third-party CUDA workloads execute alongside privileged processes**

**NIST Phase:** Preparation

**Reference:** NIST 800-61r3 §2 — Preparation: Establishing IR capability requires knowing which assets are exposed; inventory of GDDR6-equipped NVIDIA GPU nodes in shared or multi-tenant configurations is a prerequisite to any detection or response action for GPUBreach

**Controls:** NIST RA-3 (Risk Assessment) — assess likelihood and impact of Rowhammer-style GDDR6 bit-flip exploitation against your specific GPU fleet, NIST CM-8 (System Component Inventory) — enumerate all NVIDIA GPU models (A100, H100, RTX 3000/4000 series, and others using GDDR6) across on-premises clusters and cloud-attached GPU instances, CIS 1.1 (Establish and Maintain Detailed Enterprise Asset Inventory) — GPU nodes must be tagged with memory type (GDDR6 vs. HBM2e vs. GDDR5) to scope GPUBreach exposure accurately, CIS 7.1 (Establish and Maintain a Vulnerability Management Process) — hardware-layer vulnerabilities with no CVE require a documented process for tracking researcher disclosures outside NVD

**Compensating:** Run `'nvidia-smi --query-gpu=name,memory.total,driver_version --format=csv'` on all GPU hosts to enumerate NVIDIA GPU models and cross-reference against GDDR6-bearing SKUs (e.g., RTX 30xx, RTX 40xx, A10,

A30). Supplement with 'lshw -class display' on Linux or 'Get-WmiObject Win32\_VideoController' on Windows to confirm memory type where nvidia-smi output is ambiguous. A 2-person team can script this inventory sweep across SSH-accessible nodes using a simple bash loop.

**Evidence:** Before conducting inventory, snapshot existing GPU process tables with 'nvidia-smi pmon -s u' to capture any currently running CUDA processes per GPU — this establishes a baseline of active workload tenancy that will be needed if a later incident investigation must reconstruct who was co-resident on the GPU at time of a suspected exploitation event. Preserve output to timestamped log files.

## **Step 2: Review controls — audit whether GPU-to-host isolation is enforced as an active security control or assumed as an architectural guarantee; verify whether hypervisor or container-level controls (e.g., GPU passthrough configurations, IOMMU enforcement, workload sandboxing) are in place and correctly configured**

**NIST Phase:** Preparation

**Reference:** NIST 800-61r3 §2 — Preparation: Verifying that isolation controls are correctly configured before an incident is a core preparation activity; for GPUBreach, assumed GPU-to-host isolation that is not actively enforced constitutes a capability gap that must be documented in the IR plan

**Controls:** NIST SC-39 (Process Isolation) — confirm that CUDA process isolation between tenants and between GPU workloads and the host kernel is enforced at the hypervisor or OS level, not assumed from NVIDIA driver defaults, NIST CM-6 (Configuration Settings) — verify IOMMU is enabled (Intel VT-d or AMD-Vi) and GPU passthrough is configured to prevent DMA-based escalation paths that GPUBreach's host privilege escalation stage may traverse, NIST CA-2 (Control Assessments) — treat GPU isolation as an untested assumption requiring active verification; document findings as a control deficiency if isolation is architectural rather than enforced, CIS 4.2 (Establish and Maintain a Secure Configuration Process for Network Infrastructure) — extend secure configuration baselines to GPU host configurations, explicitly including IOMMU enforcement settings and NVIDIA MIG (Multi-Instance GPU) partitioning where applicable

**Compensating:** On Linux GPU hosts, verify IOMMU is active with 'dmesg | grep -e IOMMU -e DMAR -e AMD-Vi' and confirm it is not in passthrough mode. Check /proc/cmdline for 'intel\_iommu=on' or 'amd\_iommu=on'. For container-based GPU workloads using nvidia-docker or NVIDIA Container Runtime, verify 'no-new-privileges' and seccomp profiles are applied: 'docker inspect | grep -i seccomp'. For Kubernetes environments, check that GPU resource limits are set per pod and that NVIDIA Device Plugin is not configured with privileged mode enabled unnecessarily.

**Evidence:** Capture IOMMU group assignments before making configuration changes: 'find /sys/kernel/iommu\_groups/-type l' — this documents which PCIe devices (including the GPU) share IOMMU groups with the host, which is directly relevant to whether a Rowhammer-induced bit-flip in GDDR6 could affect host memory mappings. Also collect the current NVIDIA driver version and MIG configuration state via 'nvidia-smi mig -lgip' as evidence of the pre-audit isolation posture.

## **Step 3: Update threat model — incorporate hardware-layer memory manipulation attacks on GPU accelerators as an explicit threat vector; add multi-tenant GPU environments to your shared-resource privilege escalation threat scenarios and map against T1068 and T1548**

**NIST Phase:** Preparation

**Reference:** NIST 800-61r3 §2 — Preparation: Threat modeling that reflects current attack research ensures detection and response playbooks address realistic adversary capabilities; GPUBreach introduces a hardware-layer escalation path absent from most existing GPU threat models

**Controls:** NIST RA-3 (Risk Assessment) — update risk assessment to include hardware-level memory manipulation (Rowhammer-style GDDR6 bit-flip) as a distinct threat scenario with its own likelihood and impact ratings for multi-tenant GPU environments, NIST SI-4 (System Monitoring) — detection strategy must now account for MITRE ATT&CK T1068 (Exploitation for Privilege Escalation) and T1548 (Abuse Elevation Control Mechanism) emanating from CUDA process context, not just traditional user-mode or kernel exploits, NIST IR-4 (Incident Handling) — incident classification criteria should be updated to recognize anomalous privilege transitions from GPU-affiliated processes as a potential GPUBreach exploitation indicator, CIS 7.1 (Establish and Maintain a Vulnerability Management Process) —

threat model updates must feed back into the vuln management process so that future NVIDIA advisories related to GDDR6 memory integrity are triaged against documented risk, not evaluated from scratch

**Compensating:** Add a Sigma rule to your detection stack targeting process creation events where a child process with elevated privileges is spawned from a CUDA-affiliated parent (e.g., processes with GPU memory handles open per `/proc/maps` showing nvidia device mappings). On Linux, monitor for privilege changes using `auditd`: `auditctl -a always,exit -F arch=b64 -S setuid -S setreuid -k gpu_privesc`. Cross-reference with `osquery` query: `'SELECT pid, name, uid, euid FROM processes WHERE euid=0 AND uid!=0;'` run on GPU host nodes on a scheduled basis.

**Evidence:** Before formalizing the updated threat model, pull the current audit log baseline from GPU hosts showing existing CUDA process privilege levels: `cat /proc/status | grep -E "Uid|Gid"` for all running CUDA workload PIDs. This establishes a privilege-state baseline against which future anomalies (euid=0 from a CUDA process) can be compared as a potential GPUBreach exploitation indicator.

#### **Step 4: Communicate findings — brief infrastructure and cloud security leads on the GPU isolation assumption gap; frame the risk specifically around AI training clusters, GPU cloud instances, and any environment where external or untrusted workloads share GPU resources with privileged processes**

**NIST Phase:** Preparation

**Reference:** NIST 800-61r3 §2 — Preparation: Stakeholder communication ensures that infrastructure owners and cloud security leads can act on the GPU isolation gap; without this briefing, teams operating AI training clusters may continue assuming isolation guarantees that GPUBreach has shown are not hardware-enforced

**Controls:** NIST IR-6 (Incident Reporting) — while GPUBreach is currently a research disclosure without confirmed active exploitation, the isolation gap qualifies as a reportable risk condition to organizational leadership and relevant system owners, NIST IR-8 (Incident Response Plan) — update the IR plan to include GPU cloud instances and AI training clusters as high-value targets requiring explicit stakeholder notification procedures if exploitation is later confirmed, NIST SI-5 (Security Alerts, Advisories, and Directives) — distribute the GPUBreach research findings internally as a security advisory to cloud infrastructure teams, explicitly naming GDDR6-equipped NVIDIA GPU environments as the affected scope, CIS 4.6 (Securely Manage Enterprise Assets and Software) — cloud security leads must be informed so that GPU cloud instance configurations (e.g., AWS p4d, Azure NC-series, GCP A100 instances) can be reviewed for workload co-tenancy risks

**Compensating:** Prepare a one-page internal advisory (template: threat summary, affected assets by name from the inventory completed in Step 1, isolation control gap findings from Step 2, and interim mitigations). Distribute via existing security communication channels. For small teams, a documented email thread with read receipts to infrastructure leads satisfies IR-6 notification requirements and creates an auditable communication record. No SIEM required — artifact is the email or ticketing system record.

**Evidence:** Attach the `nvidia-smi` inventory output and IOMMU audit results collected in Steps 1 and 2 to the briefing package as supporting evidence. These artifacts document the specific GPU models, driver versions, and isolation posture in your environment, making the briefing concrete rather than hypothetical and ensuring infrastructure leads can immediately identify which specific clusters require action.

#### **Step 5: Monitor developments — track for NVIDIA security advisory, CVE assignment, peer-reviewed publication, and any patch or firmware mitigation guidance; subscribe to NVIDIA Product Security and CISA advisories for follow-up disclosure**

**NIST Phase:** Post Incident

**Reference:** NIST 800-61r3 §4 — Post-Incident Activity: For a threat with no current CVE or vendor patch, continuous intelligence tracking is the primary ongoing activity; the GPUBreach disclosure lifecycle (research → CVE → vendor advisory → patch) must be monitored so the organization can transition from preparation posture to active remediation when mitigations become available

**Controls:** NIST SI-5 (Security Alerts, Advisories, and Directives) — establish a documented process for receiving and acting on NVIDIA Product Security advisories (<https://www.nvidia.com/en-us/security/>) and CISA Known Exploited Vulnerabilities (KEV) additions related to GPU memory vulnerabilities, NIST IR-5 (Incident Monitoring) — track GPUBreach as an open incident record with status 'monitoring — no patch available'; update the record as NVIDIA advisories, CVE assignments, or peer-reviewed publication details emerge, NIST RA-3 (Risk Assessment) —

re-assess risk rating when a CVE is assigned or when NVIDIA publishes firmware or driver mitigations, as patch availability changes both likelihood (exploitation becomes more tractable post-publication) and impact (mitigation reduces blast radius), CIS 7.2 (Establish and Maintain a Remediation Process) — pre-stage a remediation workflow so that when NVIDIA releases a driver update or firmware patch addressing GDDR6 memory refresh hardening, the organization can execute against the asset inventory from Step 1 without delay

**Compensating:** Set up RSS feed monitoring for NVIDIA Product Security (<https://www.nvidia.com/en-us/security/>) and CISA advisories (<https://www.cisa.gov/news-events/cybersecurity-advisories>) using a free RSS aggregator (e.g., Feedly free tier). Create a NVD CPE watch for NVIDIA GPU driver components using the NVD API (free, no account required): query '<https://services.nvd.nist.gov/rest/json/cves/2.0?keywordSearch=NVIDIA+GPU+memory>' on a weekly cron job and diff output for new CVE entries. A 2-person team can maintain this with under 30 minutes per week of manual review.

**Evidence:** Maintain a timestamped changelog document recording each monitoring check: date, sources reviewed, findings (including 'no update' entries), and any change in risk posture. This log serves as evidence of due diligence for audit purposes and as the trigger document for activating the remediation workflow when NVIDIA releases a patch — the log entry noting patch availability becomes the formal handoff from monitoring to active remediation under NIST SI-2 (Flaw Remediation).

## Detection Guidance

Detection for hardware-level Rowhammer-style attacks is inherently difficult because the malicious activity occurs at the physical memory layer and does not generate conventional application or network log events. At this time, no confirmed IOCs, driver-level signatures, or vendor detection guidance exist for GPUBreach. Recommended detection approach: (1) Baseline GPU memory access patterns for your known workloads using NVIDIA DCGM (GPU metrics) or nvidia-smi telemetry to establish normal behavior; (2) Monitor host-level privilege escalation logs (Windows Security Event ID 4672, Linux audit logs for setuid/setgid execution or unexpected root process spawning) that originate from processes associated with GPU driver stacks; (3) Audit IOMMU and DMA protection configurations to verify that GPU memory regions are correctly bounded and cannot spill over into host memory; (4) In cloud and multi-tenant environments, flag any CUDA workload that requests unusual memory allocation sizes or patterns outside expected model parameters. Given the absence of confirmed IOCs or vendor signatures at this time, prioritize architecture review and anomaly baselining over signature deployment. When NVIDIA releases official advisory or technical details, revisit detection logic to incorporate specific memory access thresholds or driver-level indicators.

## Indicators of Compromise

Type	Value	Context	Confidence
TOOL	Pending – refer to Security Affairs ( <a href="https://securityaffairs.com/190455">securityaffairs.com/190455</a> ) and The Hacker News ( <a href="https://thehackernews.com/2026/04/new-gpubreach-attack-enables-full-cpu.html">thehackernews.com/2026/04/new-gpubreach-attack-enables-full-cpu.html</a> ) for any published technical indicators	No verifiable IOC values (hashes, domains, IPs, CUDA payload signatures) were present in available secondary sources; if the research team published proof-of-concept code or memory access signatures, they would appear in the original research paper or vendor advisory	LOW

## Framework Mappings

## MITRE-ATTACK

- **T1068** — Exploitation for Privilege Escalation
- **T1203** — Exploitation for Client Execution
- **T1548** — Abuse Elevation Control Mechanism

## NIST-800-53R5

- **AC-6** — Least Privilege
- **SC-7** — Boundary Protection
- **SI-2** — Flaw Remediation
- **SI-3** — Malicious Code Protection
- **SI-4** — System Monitoring
- **CM-6** — Configuration Settings
- **SI-16** — Memory Protection
- **SI-10** — Information Input Validation
- **AC-3** — Access Enforcement

## OWASP-TOP10-2021

- **A03:2021** — Injection
- **A01:2021** — Broken Access Control

## CIS-V8

- **16.10** — Apply Secure Design Principles in Application Architectures
- **6.1** — Establish an Access Granting Process
- **6.2** — Establish an Access Revoking Process
- **5.4** — Restrict Administrator Privileges to Dedicated Administrator Accounts

## SOC2-TSC

- **CC6.1** — The entity implements logical access security software, infrastructure, and architectures over protected information assets
- **CC9.2** — Manages risks associated with vendors and business partners
- **CC6.3** — Authorizes, modifies, or removes access

## HIPAA-SECURITY

- **164.312(a)(1)** — Access Control

## ISO-27001-2022

- **A.5.21** — Managing information security in the ICT supply chain
- **A.5.23** — Information security for use of cloud services

## MITRE ATT&CK Mapping

Technique ID	Technique Name	Tactic
T1068	Exploitation for Privilege Escalation	Privilege-Escalation
T1203	Exploitation for Client Execution	Execution
T1548	Abuse Elevation Control Mechanism	Privilege-Escalation

## Sources

Source	URL	Tier
<b>New GPUBreach Attack Enables Full CPU Privilege Escalation via ...</b>	<a href="https://thehackernews.com/2026/04/new-gpubreach-attack-enables-full...">https://thehackernews.com/2026/04/new-gpubreach-attack-enables-full...</a>	T3
<b>NVIDIA GPUs with GDDR6 memory can be used to take full control ...</b>	<a href="https://www.facebook.com/ethical.hack.group/posts/nvidia-gpus-with-...">https://www.facebook.com/ethical.hack.group/posts/nvidia-gpus-with-...</a>	T3
<b>GPUBreach exploit uses GPU memory bit-flips to achieve full system ...</b>	<a href="https://securityaffairs.com/190455/security/gpubreach-exploit-uses-...">https://securityaffairs.com/190455/security/gpubreach-exploit-uses-...</a>	T3
<b>Powering Your Component Needs - Chip1</b>	<a href="https://chip1.com/insights/article/gpuhammerthreat">https://chip1.com/insights/article/gpuhammerthreat</a>	T3
<b>NVIDIA GPUs with #GDDR6 memory can be used to take full control ...</b>	<a href="https://www.facebook.com/groups/2600net/posts/4521900044699752/">https://www.facebook.com/groups/2600net/posts/4521900044699752/</a>	T3

### DISCLAIMER

This intelligence report is produced by Tech Jacks Solutions Security Command Center (SCC) for informational purposes only. It does not constitute professional security advice, legal counsel, or an incident response engagement. The information herein is derived from publicly available sources and AI-assisted analysis; while every effort is made to ensure accuracy, Tech Jacks Solutions makes no warranties regarding completeness or timeliness. Organizations should conduct their own validation and consult qualified security professionals before taking action based on this report. Tech Jacks Solutions is not liable for any damages resulting from the use of this information.

Generated 2026-04-12 18:29 UTC by TJS Security Command Center