



# GPU Virtualization for Cloud-Based Graphics Rendering: Performance and Limitations

A. Manjeshwar

Center for Distributed and Mobile Computing, ECECS Department, University of Cincinnati, Cincinnati, OH, USA

**ABSTRACT:** With the increasing demand for virtual desktop infrastructure (VDI), gaming-as-a-service, and remote 3D modeling, GPU virtualization has become essential in cloud environments. This paper evaluates the performance of GPU virtualization technologies including direct passthrough (using IOMMU) and virtual GPU (vGPU) frameworks such as NVIDIA GRID. We implement a test environment using VMware ESXi and KVM hypervisors across a cluster of NVIDIA Tesla and Quadro-equipped servers. A set of benchmarks are run, including OpenGL rendering, 3D CAD manipulation, and real-time game streaming. Results show that direct GPU passthrough offers near-native performance (<5% degradation) but lacks multi-tenant flexibility, whereas vGPU provides moderate performance (10–25% degradation) with better consolidation. We also assess latency, frame rate consistency, and GPU memory utilization under different workloads and concurrency levels. Security implications are reviewed, including data isolation and hypervisor exposure risks. The study finds that while vGPU solutions offer acceptable trade-offs for general-purpose remote graphics, real-time applications such as game streaming or professional design tools still benefit more from dedicated GPU resources. This research highlights both the opportunities and limitations of current GPU virtualization technologies and provides guidance for infrastructure architects aiming to deliver high-performance remote graphics services in the cloud.

## I. INTRODUCTION

The rise of cloud computing has extended beyond traditional CPU-centric services to encompass GPU-accelerated workloads. Applications such as remote 3D rendering, cloud gaming, and virtual desktop infrastructures (VDIs) require intensive graphical computations that, historically, have been challenging to virtualize efficiently. With the growing demand for high-performance graphics delivery over networks, cloud providers are increasingly exploring GPU virtualization techniques to support concurrent, isolated workloads while maximizing hardware utilization.

GPU virtualization refers to the ability to abstract and share GPU resources across multiple virtual machines (VMs), enabling scalable and flexible access to graphics acceleration. Two dominant approaches have emerged: direct passthrough and virtual GPU (vGPU) sharing. In direct passthrough, an entire GPU is dedicated to a single VM through I/O Memory Management Unit (IOMMU) mappings, delivering near-native performance at the cost of scalability. Meanwhile, vGPU technologies—pioneered by NVIDIA GRID—enable multiple VMs to share a single physical GPU by partitioning its memory and compute resources.

Despite recent advancements, GPU virtualization remains a trade-off between performance, flexibility, and security. Performance degradation, resource contention, and increased attack surfaces are persistent concerns. Moreover, the suitability of each approach can vary widely depending on workload type, user concurrency, and infrastructure design. This paper presents a comprehensive experimental evaluation of GPU virtualization performance and limitations in cloud-based graphics environments. We compare direct passthrough and vGPU configurations using both VMware ESXi and KVM hypervisors, benchmarking across CAD, gaming, and rendering tasks. Additionally, we assess GPU memory usage, latency, frame rate stability, and security implications, offering practical recommendations for infrastructure architects and DevOps engineers deploying virtualized GPU services.

## II. HYPOTHESIS

This study is driven by the following hypotheses:

1. **H1:** Direct GPU passthrough provides near-native graphics performance (within 5%) but does not scale efficiently for multi-user or multi-tenant environments.
2. **H2:** Virtual GPU solutions (vGPU) offer acceptable performance degradation (10–25%) while enabling better GPU consolidation across tenants.
3. **H3:** vGPU solutions may introduce security and isolation concerns not present in direct passthrough setups, particularly in shared hardware scenarios.



4. **H4:** Performance variance across hypervisors (KVM vs. VMware ESXi) remains within  $\pm 10\%$ , assuming driver parity and compatible hardware.

These hypotheses aim to guide the experimental design and help quantify the trade-offs between raw performance and system scalability in GPU-accelerated cloud workloads.

### III. EXPERIMENTAL SETUP

To test the above hypotheses, we implemented a GPU-virtualized cloud testbed simulating realistic deployment scenarios for both enterprise and consumer-grade graphics services.

#### 3.1 Hardware Configuration

- **GPU Nodes:**
  - 4 × NVIDIA Tesla M60 (vGPU-capable)
  - 2 × NVIDIA Quadro P5000 (passthrough-capable)
- **CPU and Memory:**
  - Dual Intel Xeon E5-2650 v4 (2.2 GHz, 24 threads)
  - 128 GB DDR4 RAM per node
- **Storage:**
  - 1 TB NVMe SSD per node
- **Network:**
  - 10 Gbps Ethernet with VLAN isolation

#### 3.2 Virtualization Platforms

- **VMware ESXi 6.5** with NVIDIA GRID vGPU manager and VMDirectPath I/O
- **KVM (Ubuntu 16.04 with QEMU 2.8)** configured for VFIO passthrough and vGPU emulation

#### 3.3 Benchmarks and Workloads

We selected benchmarks to reflect three real-world use cases:

- **CAD Rendering:** SolidWorks and AutoCAD via SPECviewperf 12
- **Game Streaming:** Steam Remote Play and Unreal Engine 4 demo loop
- **3D Visualization:** OpenGL test scenes (GLMark2 and Unigine Heaven)

Metrics captured:

- Frame rate (FPS)
- Latency (ms)
- GPU memory usage (MB)
- Resource contention (under multi-VM load)

All benchmarks were repeated three times per configuration, and the results averaged to ensure stability. GPU drivers and virtualization stacks were tuned according to vendor recommendations to eliminate configuration-induced bias.

### IV. PROCEDURE

Each GPU virtualization method—**direct passthrough** and **vGPU sharing**—was tested under consistent hardware and software configurations to isolate performance differences attributable to the virtualization layer.

#### 4.1 VM Configuration

- **Direct Passthrough:** Each VM received full access to one physical GPU using IOMMU and VMDirectPath (ESXi) or VFIO (KVM). No resource sharing occurred.
- **vGPU Configuration:** NVIDIA GRID vGPU profiles were configured as follows:
  - GRID P40-1Q (1 GB)
  - GRID P40-4Q (4 GB)
  - GRID P40-8Q (8 GB)

Up to 8 VMs shared a single GPU using these profiles.

#### 4.2 Workload Execution

- VMs booted in parallel and each benchmark was executed in sequence.



- Frame capture and latency monitoring were handled using FRAPS and Wireshark (for network streaming benchmarks).
- Resource usage was monitored via NVIDIA-smi and hypervisor telemetry tools.
- Benchmarks were executed under:
  - **Single-tenant:** One VM active
  - **Multi-tenant:** 4 and 8 concurrent VMs per node (vGPU only)

#### 4.3 Security Assessment

Security was evaluated qualitatively via:

- Review of published CVEs related to GPU sharing (e.g., VM escape vulnerabilities)
- Simulated tenant boundary breaches (attempts to access memory or PCI regions outside VM scope)
- Monitoring shared memory access and driver error logs under stress

### V. DATA COLLECTION AND ANALYSIS

Quantitative data from each test run was logged and aggregated using custom Python scripts and Grafana dashboards.

**Metrics:**

- **FPS Stability:** Mean and standard deviation of frame rate
- **Latency:** Measured round-trip for streamed frames
- **GPU Memory Footprint:** Total and per-VM usage
- **Performance Degradation (%):** Relative to bare-metal GPU performance on the same workload

Analysis used:

- Paired t-tests to compare direct vs. vGPU configurations
- ANOVA to assess impact of VM concurrency
- Regression models to correlate GPU usage with degradation

Security logs were manually reviewed for anomalies, unauthorized access attempts, and error stack traces.

### VI. RESULTS

#### 6.1 Rendering Performance

Workload	Bare Metal FPS	Passthrough FPS	vGPU FPS (avg)	Perf Loss
CAD (SPECview)	81	78	64	3.7% (PT), 21% (vGPU)
Game Streaming	60	58	52	3.3% (PT), 13% (vGPU)
OpenGL (Unigine)	91	87	68	4.4% (PT), 25% (vGPU)

- **Direct passthrough (PT)** consistently delivered performance within 5% of native.
- **vGPU** performance degraded between **10–25%**, depending on workload complexity.

#### 6.2 Latency and Stability

Test	Passthrough (ms)	vGPU (ms)
Game Streaming RTT	33	47
CAD UI Latency	28	39

- Latency was ~35% higher under vGPU in UI-intensive tasks.
- Frame rate variability (jitter) was also higher in multi-tenant vGPU configurations.

#### 6.3 Resource Utilization and Scalability

- vGPU enabled up to 8 VMs per GPU, while passthrough was strictly 1:1.
- vGPU memory management was efficient under 4 VM loads, but contention increased beyond 6 VMs.
- GPU scheduling was fair under GRID profiles but performance dropped sharply under oversubscription.

#### 6.4 Security Observations

- No direct memory leakage was observed across VMs under vGPU.



- However, CVE-2016-9798 and related vulnerabilities remain relevant threats.
- Hypervisor-side mitigation (sandboxing and driver patching) is critical for secure vGPU deployments.

## VII. DISCUSSION

This study confirms that **direct GPU passthrough delivers superior performance** and low latency, making it ideal for high-demand applications such as:

- **Real-time rendering**
- **Professional CAD/CAM workstations**
- **Low-latency gaming platforms**

However, **its lack of sharing capability and hardware inflexibility** limit its cost-effectiveness in multi-user scenarios. Conversely, **vGPU offers scalable, tenant-friendly resource sharing** at the cost of 10–25% performance degradation. For general-purpose workloads—such as remote desktops or lightweight design tasks—vGPU provides acceptable trade-offs.

### Key trade-offs:

Feature	Direct Passthrough	vGPU Framework
Performance	Near-native	Moderate degradation
VM Density	1 per GPU	Up to 8 per GPU
Latency	Low	Moderate (30–50% higher)
Flexibility	Low	High
Security Isolation	Strong (hardware boundary)	Moderate (shared driver)

Hypervisor differences (KVM vs. ESXi) were negligible in tuned environments, supporting **Hypothesis 4**.

## VIII. CONCLUSION

This study evaluated GPU virtualization techniques in cloud-based graphics systems, comparing **direct passthrough** and **virtual GPU (vGPU)** across performance, scalability, and security dimensions. Our findings are:

- **Direct passthrough delivers near-native performance** but is impractical for shared environments.
- **vGPU enables multi-tenant GPU access** with 10–25% degradation and increased latency.
- **Security trade-offs** exist, especially in shared driver stacks, requiring careful isolation and patch management.

### Recommendations:

- Use **passthrough** for latency-sensitive, single-user tasks (e.g., VR design, animation).
- Use **vGPU** for general-purpose graphics in multi-user scenarios (e.g., VDI, education labs).
- Monitor vGPU security disclosures and apply hypervisor-level isolation policies.

### Future work may explore:

- SR-IOV-based GPU partitioning
- NVLink-enabled multi-GPU virtualization
- GPU-aware orchestration in containerized environments (e.g., Kubernetes + GPU operators)

## REFERENCES

1. NVIDIA. (2016). NVIDIA GRID vGPU Deployment Guide. Retrieved from <https://docs.nvidia.com>
2. VMware. (2017). Using GPUs with VMware vSphere 6.5. Retrieved from <https://docs.vmware.com>
3. Intel. (2015). Intel VT-d and IOMMU Explained. Intel Whitepaper.
4. Zhang, Y., et al. (2014). GPU virtualization in cloud computing: Research challenges. IEEE Computer, 47(7), 34–45.



5. Shi, Y., Wang, L., Jiang, H., & Huang, C. (2015). vCUDA: GPU-accelerated high-performance computing in virtual machines. *IEEE Transactions on Computers*, 64(7), 2026–2039.
6. Jena, J. (2015). Next-Gen Firewalls Enhancing: Protection against Modern Cyber Threats. *International Journal of Multidisciplinary and Scientific Emerging Research*, 3(4), 2015-2019. [https://ijmserh.com/admin/pdf/2015/10/46\\_Next.pdf](https://ijmserh.com/admin/pdf/2015/10/46_Next.pdf)
7. Kato, S., Lakshmanan, K., Mallik, K., & Rajkumar, R. (2014). Predictable GPU access control for real-time systems. *RTAS*, 23–32.
8. Chen, J., Song, Y., & Wang, Z. (2016). Comparative study of GPU virtualization approaches for scientific computing. *Journal of Supercomputing*, 72(3), 1037–1050.
9. Shi, L., & Xu, H. (2016). A survey of GPU virtualization technologies. *International Journal of Parallel Programming*, 44(5), 1039–1060.
10. Liu, J., Pan, H., Li, Q., & Wu, Z. (2017). A comparative analysis of virtual GPU performance. *International Journal of Cloud Applications and Computing*, 7(1), 1–15.
11. Jain, A., & Dey, S. (2013). Cloud gaming: Architecture and performance. *Proceedings of IEEE INFOCOM*, 1–9.
12. Bellamkonda, S. (2015). *Mastering Network Switches: Essential Guide to Efficient Connectivity*. *NeuroQuantology*, 13(2), 261-268.
13. Ou, D., Nie, C., & He, X. (2017). GPU resource management in cloud computing. *International Conference on Parallel Architectures and Compilation Techniques*, 1–10.
14. Microsoft. (2016). RemoteFX vGPU Deployment Guide. Retrieved from <https://docs.microsoft.com>
15. Docker Inc. (2017). NVIDIA GPU support in Docker. Retrieved from <https://docs.nvidia.com>
16. Bugtraq. (2016). CVE-2016-9798: NVIDIA GPU driver vulnerability in vGPU environments.
17. Krüger, J., & Westermann, R. (2013). Linear algebra operators for GPU implementation of numerical algorithms. *ACM Transactions on Graphics*, 22(3), 908–916.