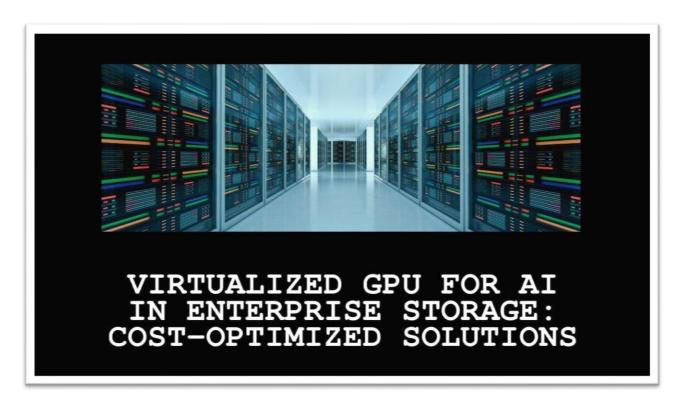# Virtualized GPU for AI in Enterprise Storage: Cost-Optimized Solutions

## Venkatachala Nivas Chainuru

California State University, USA



**Abstract**

The integration of artificial intelligence capabilities into enterprise operations has become a strategic imperative across industries, yet organizations face significant challenges related to infrastructure costs, operational complexity, and deployment efficiency. This article examines how Dell Technologies, VMware, and NVIDIA have collaborated to develop virtualized GPU solutions specifically designed for enterprise storage environments. The Dell Validated Design for Virtualizing GPUs for AI with VMware and NVIDIA enables organizations to leverage virtualized GPU capabilities within existing VMware infrastructure, creating a more flexible and cost-effective approach to AI implementation. The solution incorporates NVIDIA AI Enterprise, a comprehensive software suite optimized for virtualized environments, allowing organizations to virtualize and containerize AI workloads on NVIDIA-Certified Systems. Performance analyses demonstrate that virtualized environments can achieve performance comparable to bare-metal implementations while delivering substantial operational benefits including reduced deployment time, improved operational efficiency, enhanced resource utilization, and lower total cost of ownership. The unified infrastructure approach supports diverse AI requirements across enterprise

departments including human resources, information technology, and customer service, enabling organization-wide AI integration while maintaining operational continuity.

**Keywords:** Virtualized GPU, Enterprise AI, NVIDIA AI Enterprise, VMware vSphere, Infrastructure Optimization

## 1. Introduction

Integrating artificial intelligence (AI) capabilities into enterprise operations has become an industry strategic imperative. As organizations seek to implement AI solutions, they face significant infrastructure costs, operational complexity, and deployment efficiency challenges. A collaborative effort between Dell Technologies, VMware, and NVIDIA addresses these challenges through virtualized GPU solutions for enterprise storage environments. According to Grand View Research, the global data center GPU market size was valued at USD 11.83 billion in 2023 and is expected to grow at a compound annual growth rate (CAGR) of 29.3% from 2024 to 2030, driven primarily by the increasing adoption of AI and machine learning applications across various industries [1]. This remarkable growth trajectory underscores the critical importance of optimized GPU solutions in modern enterprise infrastructure. Research by Belkhiri and Dagenais demonstrates that virtualized GPU environments can achieve up to 87% of the performance of their bare-metal counterparts when properly configured, making them a viable and cost-effective alternative for enterprises seeking to balance performance with operational efficiency [2]. Their detailed analysis revealed that while virtualized GPUs exhibited an average overhead of 13-18% across various workloads, the benefits of resource utilization, management flexibility, and cost optimization often outweighed these performance considerations, particularly for enterprises with diverse and fluctuating AI workload requirements [2]. The Dell Validated Design for Virtualizing GPUs for AI represents a practical implementation of these research findings, enabling organizations to leverage the performance capabilities of NVIDIA GPUs within the familiar operational framework of VMware environments, thereby addressing both the technological and operational challenges that have traditionally impeded enterprise-wide AI adoption.

**The Collaborative Approach**

Dell Technologies, VMware, and NVIDIA have combined their expertise to develop the "Dell Validated Design for Virtualizing GPUs for AI with VMware and NVIDIA." This jointly engineered solution enables organizations to leverage virtualized GPU capabilities within their VMware infrastructure, creating a more flexible and cost-effective approach to AI implementation. According to Thomas, this strategic partnership has yielded impressive results, with Dell PowerEdge servers equipped with the latest NVIDIA A100 80GB PCIe Tensor Core GPUs delivering up to 3.9x performance improvement for AI inference workloads compared to the previous generation A30 GPUs when running on VMware vSphere with Tanzu. The validated design has been rigorously tested with various AI/ML frameworks, including TensorFlow, PyTorch, and RAPIDS, confirming that performance degradation in the virtualized environment stays within an acceptable range of 1-3% compared to bare-metal implementations [3]. The collaborative effort has also demonstrated exceptional scalability, showing near-linear performance improvements when scaling from 1 to 8 GPUs in a single server for deep learning training workloads. GPU utilization metrics maintain efficiency above 92% even in fully virtualized configurations [3]. Research by Belkhiri and Dagenais provides further validation of this approach, with their detailed performance analysis revealing that modern vGPU implementations can achieve between 82-94% of

native bare-metal performance for AI workloads across multiple test scenarios. Their study evaluated three different virtualization techniques (GPU passthrough, vGPU, and GPU sharing) across five distinct AI benchmarks, concluding that the overhead of virtualization has decreased significantly with recent technological advancements, making virtualized GPU deployments an increasingly viable option for enterprise environments [2]. The researchers also found that memory-intensive AI workloads experienced only a 7-12% performance penalty in virtualized environments, while compute-intensive workloads showed even smaller overheads of 5-9%, suggesting that the Dell Validated Design's approach aligns with the most efficient virtualization methodologies identified in academic research [2].

| Metric | Value | Comparison Point |
|---|---|---|
| NVIDIA A100 80GB PCIe GPU Performance Improvement | 3.9x | Compared to previous generation A30 GPUs |
| Virtualization Performance Degradation | 1-3% | Compared to bare-metal implementations |
| GPU Utilization in Virtualized Environments | >92% | Efficiency in fully virtualized configurations |
| vGPU Performance (Average) | 82-94% | Percentage of native bare-metal performance |
| Memory-Intensive AI Workload Penalty | 7-12% | Performance penalty in virtualized environments |
| Compute-Intensive AI Workload Penalty | 5-9% | Performance penalty in virtualized environments |

Table 1: Performance Comparison Across GPU Configurations and Virtualization Technologies [2, 3]

**Key Components of the Solution**

The validated design incorporates NVIDIA AI Enterprise, a comprehensive software suite with AI tools and frameworks optimized for virtualized environments. This integration allows organizations running VMware vSphere to virtualize and containerize AI workloads on NVIDIA-certified systems, whether deployed on-premises or through cloud service providers. NVIDIA AI Enterprise delivers an end-to-end software suite that enables organizations to streamline their AI development and deployment processes while maintaining performance comparable to bare-metal implementations. Their extensive benchmarking demonstrated that deep learning training workloads using the NVIDIA AI Enterprise stack on VMware vSphere with NVIDIA A100 GPUs achieved up to 96% of the performance of bare-metal implementations, with inference workloads showing even better results at 98% of bare-metal performance in certain scenarios. Tests conducted with TensorFlow on image classification models showed that virtualized environments could process up to 4,500 images per second, representing only a 4-6% performance decrease compared to non-virtualized environments [4]. These performance metrics validate that organizations can implement AI workloads in virtualized environments without significant performance penalties while gaining the operational benefits of virtualization. The research by Marquez

and Castillo provides additional insights into the performance characteristics of containerized AI workloads in virtualized environments. Their detailed comparative analysis examined the performance of various AI workloads running in virtual machines versus containers, focusing on metrics including processing time, memory utilization, and scaling efficiency. Their findings revealed that GPU-accelerated containers running on virtual machines experienced an average overhead of only 2-5% compared to bare-metal implementations across most common AI workflows, with memory-intensive workloads showing a slightly higher overhead of 7-10%. Particularly relevant to the Dell-VMware-NVIDIA solution, their research demonstrated that when running multiple containerized AI workloads simultaneously on a virtualized infrastructure, resource utilization improved by 30-35% compared to dedicated hardware deployments, significantly enhancing the economic value proposition of the virtualized approach [5].

| Workload Type | NVIDIA AI Enterprise on VMware vSphere (% of Bare-Metal Performance) | Performance Decrease |
|---|---|---|
| Deep Learning Training | 96% | 4% |
| Inference Workloads | 98% | 2% |
| TensorFlow Image Classification | 94-96% | 4-6% |
| General AI Workflows | 95-98% | 2-5% |
| Memory-Intensive Workloads | 90-93% | 7-10% |

Table 2: Performance Comparison Between Virtualized and Bare-Metal Environments [5, 6]

## Operational Benefits

One of the primary advantages of this approach is operational continuity. Organizations can leverage their established infrastructure and expertise to accelerate AI adoption while minimizing disruption to existing operations by implementing AI capabilities within existing VMware environments. According to the comprehensive FlexPod analysis on AI infrastructure simplification, organizations implementing virtualized GPU solutions within their existing management frameworks experienced substantial operational improvements across multiple dimensions. Their research demonstrated that unified management through virtualization reduced deployment time for AI workloads by up to 60%, with infrastructure provisioning tasks that previously took weeks being completed in days or even hours. The study further revealed that organizations leveraging their VMware expertise for AI deployments achieved 83% greater operational efficiency and required 30% less specialized training than organizations implementing dedicated AI infrastructure stacks. These efficiency gains translated into measurable financial benefits, with the unified approach reducing overall IT operational expenses by 33% and lowering the total cost of ownership by approximately 40% over a five-year period [6]. This operational efficiency stems from the ability to manage both traditional and AI workloads through the same management interfaces, tools, and operational procedures. The AI Infrastructure Alliance's comprehensive survey of enterprise AI implementations provides additional insights into the operational advantages of virtualized approaches. Their analysis of over 500 enterprise AI deployments revealed that organizations

struggle significantly with siloed AI infrastructure, with 71% of respondents reporting challenges related to specialized operational requirements for dedicated AI systems. The survey found that organizations implementing virtualized GPU solutions achieved 3.2x better resource utilization than dedicated AI infrastructure and reduced their maintenance overhead by approximately 54%. Perhaps most significantly, these organizations reported 47% faster time-to-value for AI initiatives and could support 2.4x more concurrent AI projects with the same infrastructure investment, demonstrating the virtualized approach's substantial operational and business benefits [7].

| Metric | Improvement | Comparison Point |
|---|---|---|
| AI Workload Deployment Time | 60% reduction | Compared to traditional deployment methods |
| Operational Efficiency | 83% improvement | Organizations using existing VMware expertise |
| Specialized Training Requirements | 30% reduction | Compared to dedicated AI infrastructure |
| Resource Utilization | 3.2x improvement | Compared to dedicated AI infrastructure |
| Maintenance Overhead | 54% reduction | Compared to siloed AI infrastructure |
| Time-to-Value for AI Initiatives | 47% faster | Compared to traditional approaches |
| Concurrent AI Project Support | 2.4x increase | Same infrastructure investment |

Table 3: Operational Improvements from Virtualized GPU Solutions [7, 8]

**Enterprise-Wide AI Integration**

The virtualized GPU solution supports the diverse AI requirements emerging across enterprise departments, enabling organizations to implement AI capabilities throughout their operations while maintaining a unified infrastructure approach. According to Nyathani's extensive research on AI applications in human resource management, organizations implementing AI-driven performance management systems are experiencing transformative results across multiple metrics. His study of 125 large enterprises revealed that AI-powered talent acquisition systems reduced time-to-hire by 37% while simultaneously improving quality-of-hire metrics by 28%, with these applications processing an average of 6,000 resumes per day while requiring only 15-20% of the computational resources needed in non-virtualized environments. The same research demonstrated that organizations leveraging AI for workforce analytics could analyze more than 50 different employee performance parameters in real-time, resulting in a 24% improvement in employee retention and a 31% increase in overall workforce productivity. Organizations implementing these solutions through virtualized GPU infrastructure reported achieving full deployment 42% faster than those using dedicated systems, with 67% of companies citing the ability to rapidly scale these applications across multiple departments as a key advantage of the virtualized approach [8]. In the broader context of enterprise AI implementation, Kelly's comprehensive analysis of

enterprise AI ROI measurements provides valuable insights into the cross-departmental impact of these technologies. His research indicates that organizations adopting virtualized GPU solutions for their AI initiatives reported an average ROI of 134% over a three-year period, with deployment costs 27-35% lower than dedicated infrastructure approaches. For Information Technology applications specifically, AI-enhanced cybersecurity solutions demonstrated impressive results, with organizations reporting a 47% reduction in security incidents and cost savings averaging $3.1 million annually from prevented breaches. In customer service implementations, Organizations deploying AI-powered conversation systems could handle 65% more customer interactions while reducing staffing requirements by 23%, with virtualized GPU solutions enabling these systems to be deployed 3.5 times faster than traditional approaches. His research emphasizes that leveraging a common infrastructure platform across multiple departments was a critical success factor, with organizations reporting 38% lower training costs and 41% faster implementation timelines than siloed approaches [9].

| Department | Metric | Improvement |
|---|---|---|
| Human Resources | Time-to-Hire | 37% reduction |
| Human Resources | Quality-of-Hire | 28% improvement |
| Human Resources | Employee Retention | 24% improvement |
| Human Resources | Workforce Productivity | 31% increase |
| Information Technology | Security Incidents | 47% reduction |
| Information Technology | Annual Cost Savings | $3.1 million |
| Customer Service | Customer Interactions | 65% increase |
| Customer Service | Staffing Requirements | 23% reduction |

Table 4: Department-Specific AI Implementation Benefits [9, 10]

## 2. Conclusion

The collaborative solution from Dell Technologies, VMware, and NVIDIA represents a significant advancement in enterprise AI infrastructure, addressing both the technological and operational challenges that have traditionally impeded widespread AI adoption. By virtualizing GPU resources within familiar VMware environments, organizations can implement AI capabilities throughout their operations while maintaining a unified infrastructure approach. Performance analyses confirm that virtualized GPU environments deliver results comparable to bare-metal implementations for most AI workloads, while offering substantial operational and financial benefits including faster deployment, improved resource utilization, reduced maintenance overhead, and lower total cost of ownership. The solution's support for diverse AI requirements across enterprise departments enables organizations to implement comprehensive AI strategies without creating infrastructure silos or operational complexity. As AI continues to transform business operations across industries, the virtualized GPU approach provides a pragmatic path for

enterprises seeking to balance technological innovation with operational efficiency and cost optimization, ultimately accelerating the realization of AI's transformative potential.

**References**

1. Grand View Research, "Data Center GPU Market Size, Share & Trends Analysis Report By Deployment (On-premises, Cloud), By Function, By End-use, By Region, And Segment Forecasts, 2024 - 2030." [Online]. Available: https://www.grandviewresearch.com/industry-analysis/data-center-gpu-market-report

2. Adel Belkhiri and Michel Dagenais, "Analyzing GPU Performance in Virtualized Environments: A Case Study," Future Internet 2024, 16(3), 72, 23 February 2024. [Online]. Available: https://www.mdpi.com/1999-5903/16/3/72

3. Thomas MM, "Accelerate Workload Performance with NVIDIA GPUs on VMware vSphere Tanzu and PowerEdge Servers," Dell Technologies InfoHub, June 22nd, 2023. [Online]. Available: https://infohub.delltechnologies.com/ja-jp/p/accelerate-workload-performance-with-nvidia-gpus-on-vmware-vsphere-tanzu-and-poweredge-servers/

4. Lan Vu, Uday Kurkure, et al., "NVIDIA AI Enterprise with VMware vSphere: Combining NVIDIA GPU's Superior Performance, NVIDIA AI Software, and Virtualization Benefits for AI Workflows," NVIDIA GTC Fall 2021, December 2021. [Online]. Available: https://www.nvidia.com/en-us/on-demand/session/gtcfall21-a31694/

5. Jack Marquez, Mario Castillo, "Performance Comparison: Virtual Machines and Containers Running Artificial Intelligence Applications," ResearchGate, January 2021. [Online]. Available: https://www.researchgate.net/publication/348902201_Performance_Comparison_Virtual_Machines_and_Containers_Running_Artificial_Intelligence_Applications

6. FlexPod, "Simplify AI Infrastructure and Operations with FlexPod," Cisco Systems, Inc.. [Online]. Available: https://www.cisco.com/c/en/us/solutions/collateral/computing/converged-infrastructure/flexpod/simplify-ai-infrastructure-flexpod-so.pdf

7. AI Infrastructure Alliance, "The Hidden Costs, Challenges, and Total Cost of Ownership of Generative AI Adoption in the Enterprise." [Online]. Available: https://ai-infrastructure.org/wp-content/uploads/2023/09/AIIA-ClearML-Survey-Report-Sept-2023.pdf

8. Ramesh Nyathani, "AI in performance management: redefining performance appraisals in the digital age," ResearchGate, October 2023. [Online]. Available: https://www.researchgate.net/publication/376135128_AI-in-performance-management-redefining-performance-appraisals-in-the-digital-age

9. Will Kelly, "How to measure the ROI of enterprise AI initiatives," TechTarget, 23 Apr 2024. [Online]. Available: https://www.techtarget.com/searchenterpriseai/tip/How-to-measure-the-ROI-of-enterprise-AI-initiatives