



NVIDIA® VIRTUAL COMPUTE SERVER (VCS) POWER THE MOST COMPUTE-INTENSIVE WORKLOADS WITH VIRTUAL GPU_s

TRANSFORMING VIRTUALIZED COMPUTE

As the number of servers grow across the data center, IT admins expect to manage them with standard server virtualization platforms from VMware, Red Hat, Nutanix, and Citrix. According to Gartner, "hypervisor-based server virtualization is now mature, with 80% to 90% of server workloads running in a virtual machine (VM) for most midsize to large enterprises." However, this traditional data center infrastructure using hypervisor-based virtualization has been limited to CPU-only servers, with VDI as an exception. As a result, GPU accelerated servers running AI, deep learning, data science, and high-performance computing (HPC) workloads are often isolated in other servers in the data center, limiting utilization, flexibility, and manageability.

[NVIDIA® Virtual Compute Server \(vCS\)](#) enables the benefits of hypervisor-based server virtualization for GPU-accelerated servers. Data center admins are now able to power any compute-intensive workload with GPUs in a virtual machine (VM).

vCS software virtualizes NVIDIA GPUs to accelerate large workloads, including more than 600 GPU accelerated applications for AI, deep learning, and HPC. With GPU sharing, multiple VMs can be powered by a single GPU, maximizing utilization and affordability, or a single VM can be powered by multiple virtual GPUs, making even the most intensive workloads possible. And with support for nearly all major hypervisor virtualization platforms, data center admins can use the same management tools for their GPU-accelerated servers as they do for the rest of their data center.

LICENSED FOR COMPUTE

Unlike [NVIDIA® GRID® vPC/vApps](#) and [Quadro® Virtual Data Center Workstation](#) (Quadro vDWS), vCS is not tied to a user with a display. It's licensed per GPU as a 1-year subscription with NVIDIA enterprise support included. This allows a number of compute workloads in multiple VMs to be run on a single GPU, maximizing utilization of resources and ROI.

OPTIMIZED FOR CONTAINERS WITH NGC SOFTWARE

vCS supports [NVIDIA NGC](#) GPU-optimized software for deep learning, machine learning, and HPC. NGC software includes containers for the top AI and data science software, tuned, tested, and optimized by NVIDIA, as well as fully-tested containers for HPC applications and data analytics.

NGC also offers pre-trained models for a variety of common AI tasks that are optimized for NVIDIA [Tensor Core](#) GPUs and includes instructions and scripts for creating deep learning models with sample performance and accuracy metrics. This allows data scientists, developers, and researchers to reduce deployment times and project complexity so that they can focus on building solutions, gathering insights, and delivering business value.

FEATURES

- > **GPU Performance** - Access the most powerful GPUs in a virtualized environment.
- > **Management and Monitoring** - Streamline data center manageability by leveraging hypervisor-based tools.
- > **Live Migration** - Live migrate GPU-accelerated VMs without disruption, easing maintenance and upgrades.
- > **Maximize Utilization** - Increase utilization and productivity with both GPU sharing and aggregation of multiple GPUs.
- > **Security** - Extend the benefits of server virtualization to GPU workloads.
- > **Multi-Tenant** - Isolate workloads and securely support multiple users.
- > **Rapid Deployment** - Leverage GPU-optimized NGC containers for AI, data science, and HPC.
- > **Reliability** - Prevent against data corruption with error-correcting code (ECC) and dynamic page retirement.
- > **Enterprise Software Support** - Enterprise Software Support - Get support with NVIDIA Enterprise and NVIDIA NGC Support Services.

NVIDIA vCS FEATURES LIST

Configuration and Deployment		Data Center Management	
GPU Sharing (fractional)	✓	Host-, Guest-, and Application-Level Monitoring	✓
GPU Aggregation (Multi-vGPU)	✓	Live Migration	✓
Peer-to-Peer over NVLink	✓	Support	
ECC & Dynamic Page Retirement	✓	NVIDIA Direct Enterprise-Level Technical Support	✓
Linux OS Support	✓	Maintenance Releases, Defect Resolutions, and Security Patches ²	✓
Windows OS Support	X	NGC Support Services ³	✓
NVIDIA Compute Driver	✓		
NVIDIA Graphics Driver	X		
NVIDIA Quadro Driver	X		
Quality-of-Service Scheduling	✓		

vCS PROFILES

Maximum Frame Buffer Supported	48GB
Minimum Frame Buffer Supported	4GB
Maximum Multi-Tenancy	8:1
Available Profiles	4C, 5C ⁸ , 6C, 8C, 10C ⁸ , 12C, 16C, 20C ⁸ , 24C ⁴ , 32C ⁵ , 40C ⁸ 48C ⁶

RECOMMENDED GPUs FOR vCS

	NVIDIA A100 Tensor Core GPU	NVIDIA V100S Tensor Core GPU	NVIDIA Quadro RTX™ 8000 GPU	NVIDIA Quadro RTX 6000 GPU	NVIDIA T4 Tensor Core GPU
Memory	40 GB HBM2	32 GB HBM2	48 GB HBM2	24 GB HBM2	16 GB HBM2
Peak FP32	19.5 TFLOPS	16.4 TFLOPS	14.9 TFLOPS	14.9 TFLOPS	8.1 TFLOPS
Peak FP64	9.7 TFLOPS	8.2 TFLOPS	-	-	-
NVLink: Number of GPUs per VM	Up to 8	Up to 8	2	2	-
ECC and Page Retirement	✓	✓	✓	✓	✓
Multi-vGPU per VM ⁷	Up to 16	Up to 16	Up to 16	Up to 16	Up to 16

ADDITIONAL SUPPORTED GPUs

NVIDIA® P40, P100, and P6 for blade form factor.

¹ Gartner, [Market Guide for Server Virtualization](#), April 24, 2019. ID G00350674.

² Available with an active Support, Updates, and Maintenance (SUMS) contract.

³ Not included with vCS license, but available separately through [NVIDIA NGC Support Service partners](#).

⁴ 24C profile available with Quadro RTX 6000 and RTX 8000.

⁵ 32C profile available with NVIDIA V100.

⁶ 48C profile supported with Quadro RTX 8000, V100 and V100S Tensor Core GPUs.

⁷ Number of multi-GPUs supported may vary by hypervisor.

⁸ Support for these profiles will be available in the upcoming vGPU summer 2020 release for use with the NVIDIA A100 Tensor Core GPU.