HITACHI
Inspire the Next

# VMware Horizon View with Graphics Acceleration on Hitachi Unified Compute Platform HC Horizon Fast Track Certified

## Reference Architecture Guide

By Jose Perez

February 2019

# Feedback

Hitachi Vantara welcomes your feedback. Please share your thoughts by sending an email message to SolutionLab@HitachiVantara.com. To assist the routing of this message, use the paper number in the subject and the title of this white paper in the text.

## Revision History

| Revision | Changes | Date |
|---|---|---|
| MK-SL-128-00 | Initial release | February 12, 2019 |

# Table of Contents

# VMware Horizon View with Graphics Acceleration on Hitachi Unified Compute Platform HC Horizon Fast Track Certified

## Reference Architecture Guide

This document describes the reference architecture for the implementation of VMware Horizon View with Graphics Acceleration on Hitachi Unified Compute Platform HC (UCP HC).

This reference architecture combines UCP HC All-Flash appliance with VMware vSAN technology, VMware Horizon View, and NVIDIA Virtual GPU technology to enable high performance and graphics capabilities for graphics-intensive applications.
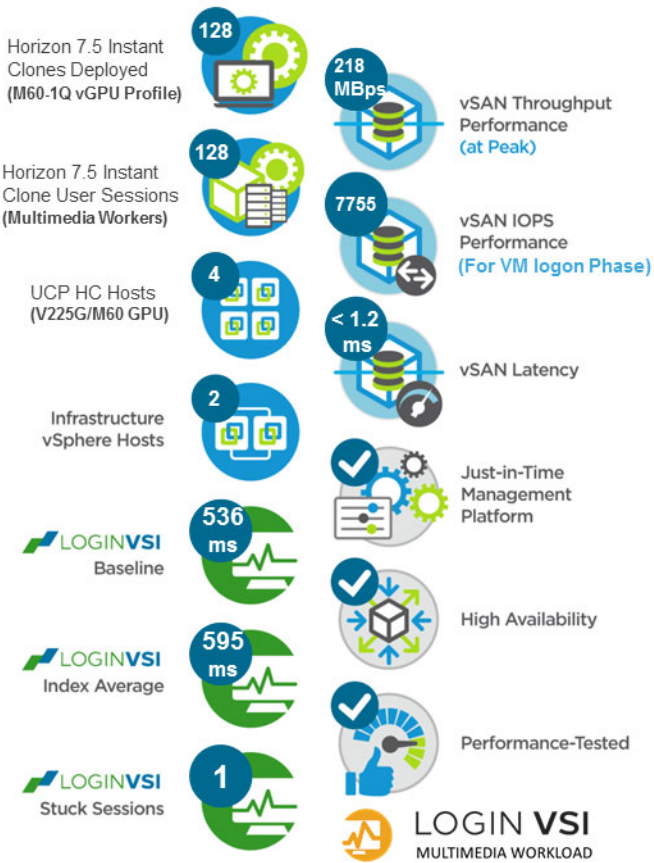
Graphics Processing Units (GPUs) have changed the dynamics and landscape of Virtual Desktop Infrastructure. By employing GPU technology, the strain caused by the increased load can now be offloaded from the CPU into the GPUs, improving performance, user experience, and increasing virtual desktop density per host.

Extensive testing was executed to evaluate the performance and graphics capabilities of UCP HC with 3 different models of Tesla GPU cards, although the hardware used in this guide configured with 2 GPUs per server. This solution supports up to 4 GPUs per server, expanding the density of vGPU instances that can be deployed per compute node.

This paper includes the results of a series of tests generated with Login VSI, the industry standard benchmarking tool for VDI workloads. For these tests, representative vGPU profiles were selected and only Power Worker and Multimedia Worker user workloads were tested.

Some of the results for one of the Multimedia Worker user workload tests are summarized in Figure 1.

**Figure 1**



---

**Note** — These practices were developed in a lab environment. Many things affect production environments beyond prediction or duplication in a lab environment. Follow recommended practice by conducting proof-of-concept testing for acceptable results before implementing this solution in your production environment. Test the implementation in a non-production, isolated test environment that otherwise matches your production environment.

---

## Solution Overview

Hitachi Unified Compute Platform HC V225G Series is a high-performance appliance to run I/O and graphics intensive applications. This UCP HC system combines VMware vSAN, the Intel Xeon Scalable processor, All-Flash storage, and NVIDIA GPU technology to enable high performance and graphics capabilities.

The combination of Hitachi Advanced Server DS225 family, NVIDIA GPU, and GRID technology, enables not only high graphics performance, but also support for high number of users since a single DS225 server can support up to 4 × Tesla M60/P40 GPUs or up to 3 × Tesla M10 GPUs. The number of users per server is determined by the combination of the type and number of GPU cards, the vGPU profile, and licensing.

Considering this number of GPU cards per server and the vGPU profile, a single DS225 server can support up to 192 vGPU VMs when using M10 GPUs, up to 128 vGPU VMs when using M60 GPUs, and up to 96 vGPU VMs when using P40 GPUs.
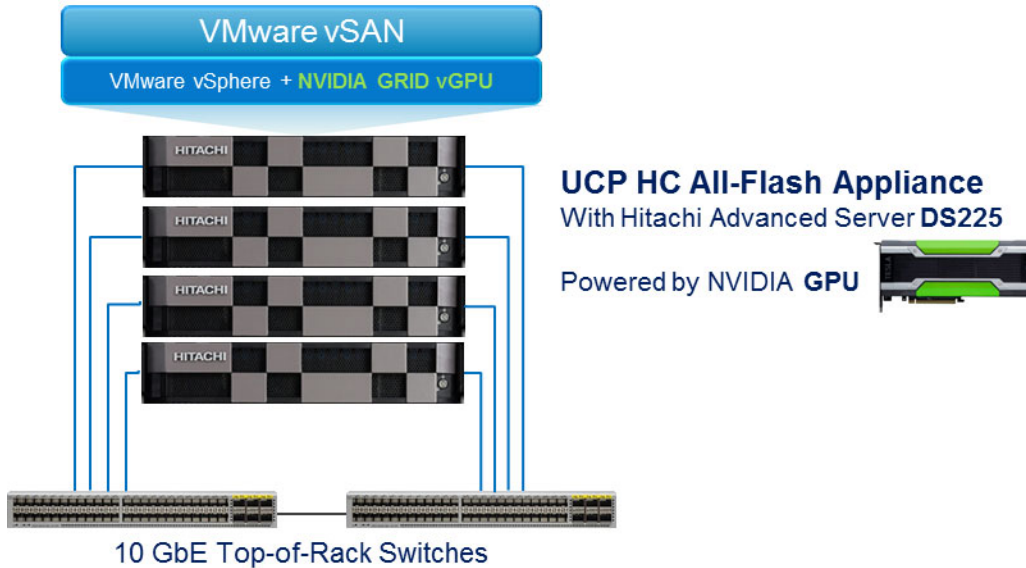
UCP HC together with VMware Horizon and NVIDIA GRID vGPU technology provides a superior user experience for graphics-heavy VDI users.

Hitachi Unified Compute Platform HC providing a scalable, building block style approach to deploying infrastructure for VDI, offering the enterprise predictable cost, and delivering a high-performing desktop experience to end users.

### Logical Architecture

Figure 2 shows the high-level infrastructure for this solution.

**Figure 2**



## Product Features

The following information describes the hardware and software features used in this reference architecture.

### Unified Compute Platform HC

Hitachi Unified Compute Platform HC (UCP HC) solutions combine compute, storage, and virtualization into a simple scalable and easy-to-deploy and manage hyperconverged infrastructure powered by software-defined storage VMware vSAN and Hitachi software to extend the agility and simplicity of the UCP family.

Ensure powerful performance and availability for traditional and cloud-native application deployment. UCP HC solutions deliver a reliable platform for business-critical applications, databases, virtual desktop infrastructure (VDI), DevOps, containers, and remote and branch office (ROBO) deployments, among other use cases.

UCP HC V225G is part of the UCP HC family that supports NVIDIA Tesla graphics processing unit (GPU) technology.

### VMware vSphere

VMware vSphere is a virtualization platform that provides a datacenter infrastructure. It helps you get the best performance, availability, and efficiency from your infrastructure and applications. Virtualize applications with confidence using consistent management.

VMware vSphere has the following components:

- **VMware vSphere ESXi** — A hypervisor that loads directly on a physical server. ESXi provides a robust, high-performance virtualization layer that abstracts server hardware resources and makes them shareable by multiple virtual machines.

- **vCenter Server** — Management of the vSphere environment through a single user interface. With vCenter, there are features available such as vMotion, Storage vMotion, Storage Distributed Resource Scheduler, High Availability, and Fault Tolerance.

## VMware vSAN

VMware vSAN is the industry-leading software that powers hyperconverged Infrastructure solutions.

VMware hyperconverged software combines compute, networking, storage, and management resources into a hyperconverged infrastructure appliance to create a simple, easy to deploy, all-in-one solution.

## VMware Horizon

VMware Horizon transforms static desktops into secure, virtual workspaces that can be delivered on demand. Provision virtual or remote desktops and applications through a single VDI platform to streamline management and easily entitle end users.

Dynamically allocate resources with virtual storage, virtual compute, and virtual networking to simplify management and drive down costs. With Horizon, reduce day-to-day operations costs with a single platform that allows you to extend virtualization from the datacenter to your devices.
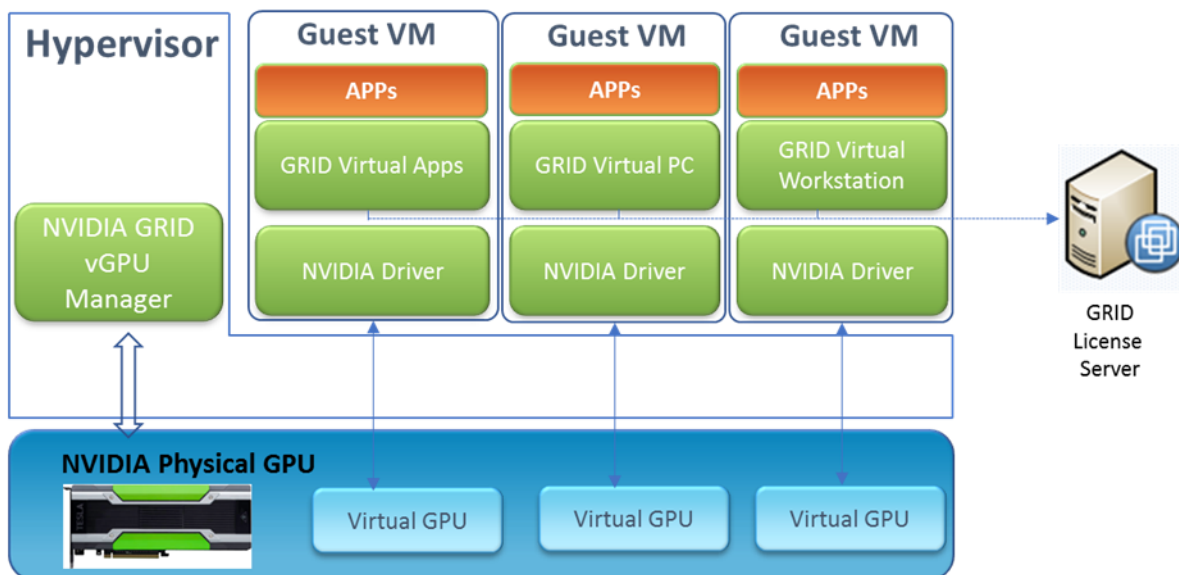
## NVIDIA GRID GPU

Nvidia GRID Virtual GPU (vGPU) is a graphics virtualization software technology that enables multiple virtual machines (VMs) to have simultaneous access to NVIDIA GPU technology, using the same graphics drivers that are deployed on non-virtualized systems. This provides VMs unparalleled graphics performance and scalability by sharing a GPU among multiple workloads/users.

*Virtual GPU Profiles*

Each of the GPU cards implement multiple physical GPUs: Tesla M10 have 4 GPUs, Tesla M60 has 2 GPUs, and Tesla P40 has 1 GPU. And each physical GPU supports several types of virtual GPUs. These virtual GPU types, also called virtual GPU profiles, have a fixed amount of graphics memory (frame buffer), number of supported display heads, and maximum resolution, and they are targeted at different classes of workloads. These profiles can be grouped into three different series:

- **Q-Series** virtual GPU type: targeted at designers and power users

- **B-Series** virtual GPU type: targeted at power users

- **A-Series** virtual GPU type: targeted at virtual applications users

**Figure 3**



## Networking

Cisco Nexus 9000 Series Switches are designed to tune datacenter networks into modern, automated environments that deliver exceptional application performance and effective IT operations. Cisco Nexus 9000 Series Switches provide the foundation for the Cisco Application Centric Infrastructure (ACI) and deliver savings in capital expenditures (CapEx) and operating expenses (OpEx) to achieve an agile IT environment. As you plan your Cisco ACI implementation, Cisco Nexus 9000 Series Switches can be deployed in standalone mode.

## Login VSI

Login VSI successfully predicts, validates, and manages the performance of virtualized desktop environments.

Login VSI makes it easy to load test, benchmark, and plan capacity to improve end user experience and productivity for even the most complex virtualized desktop environments. Login VSI tests performance using virtual users, so your real users benefit from consistently great performance.

With Login VSI, you can gain performance insights that enable you to:

- Predict the performance impact of necessary updates and upgrades

- Know the maximum user capacity of your current infrastructure

- Understand the end users' perspective on performance

## Solution Components

This section describes the Hardware and Software components used to validate this reference architecture.

### Hardware Components

Figure 4 shows the Hitachi Advanced Server DS225 with support for up to four GPU cards.

**Figure 4**

Table 1 describes the details of the hardware configuration of the UCP HC with the GP-GPU solution used to run the different test cases of this reference architecture.

TABLE 1. HARDWARE COMPONENTS

| Vendor | Hardware | Detail Description | Version | Quantity |
|---|---|---|---|---|
| Hitachi Vantara | Hitachi Advanced Server DS225 | UCP HC V225G – 4-Node Cluster<br><br>DS225 Server Specifications:<br><br>■ 2 × Intel Xeon Gold 6154 CPU @ 3.0Ghz<br>■ Memory: 512 GB DDR4<br>■ 1 × Intel X527-DA4 10 GbE dual-port SFP+<br>■ 1 × 1GbE RJ45 port for OOB<br>■ 1 × LSI SAS 9305-16i PCIe Controller<br>■ vSAN Cache: 2 × Intel S4600 960 GB<br>■ vSAN Capacity: 4 × Intel S4500 1.92 TB<br><br>■ NVIDIA GPU Cards:<br>　■ Test Case 1: 2 × M10<br>　■ Test Case 2: 2 × M60<br>　■ Test Case 3: 2 × P40 | BIOS Version: 3A10.H8<br><br>BMC: 4.23.06 | 1 |
| Cisco | Top-of-Rack 10 GbE switch | Cisco Nexus 93180YC-EX NX-OS Version | Firmware:<br><br>7.0(3)I4(6) | 2 |

## Software Components

Table 2 describes the software components of the UCP HC solution used to run the different test cases of this reference architecture.

TABLE 2. SOFTWARE COMPONENTS

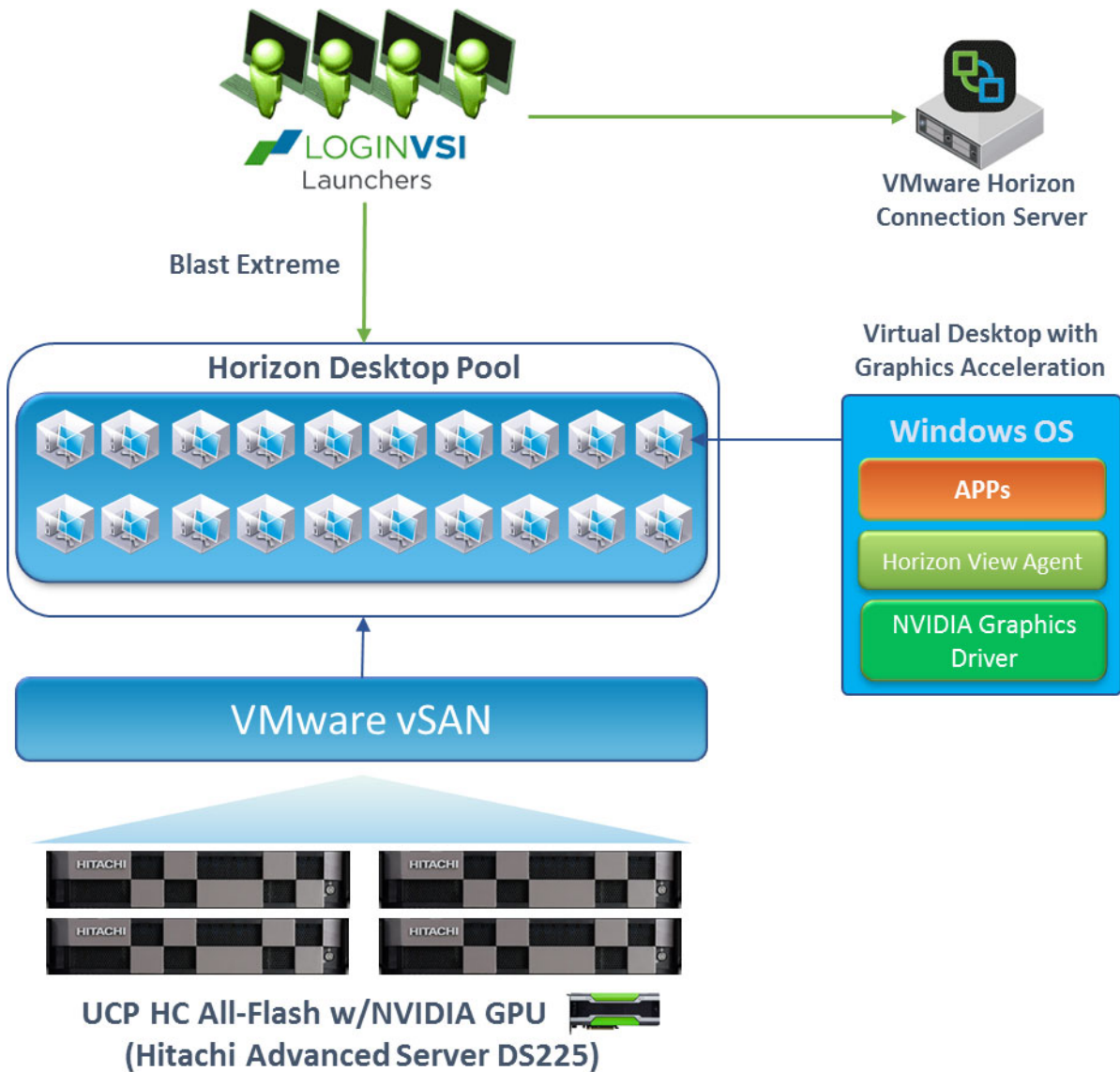| Software | Version | Function |
|---|---|---|
| VMware vCenter server | 6.7, Build 9232942 | Management console |
| VMware ESXi | 6.7, Build 9214924 | Hypervisor |
| NVIDIA Virtual GPU Manager | 6.2 | GPU Virtualization |

TABLE 2. SOFTWARE COMPONENTS (CONTINUED)

| Software | Version | Function |
|---|---|---|
| Microsoft Windows Server 2012 | Datacenter, R2 | Operating System |
| Microsoft Windows Server 2016 | Datacenter | Operating System |
| Microsoft SQL Server | 2014 SP1 | Database server |
| Microsoft Windows 10 | Enterprise Edition, SP1 | Operating System |
| VMware Horizon View | 7.5 | Management console |
| VMware Horizon Client | 4.8 | |
| Login VSI | 4.1.32.1 | Benchmarking tool |

## Solution Design

The infrastructure servers for VMware Horizon and Login VSI used for this solution were placed on separate infrastructure clusters with dedicated resources. The UCP HC cluster was dedicated for Horizon desktops only. Figure 5 shows the key hardware and software components used for the reference architecture tests.

**Figure 5**

## UCP HC Configuration

### *Storage Configuration*

The vSAN cluster was based on UCP HC All-Flash with GPU. Each host has 2 disk groups, each disk group with 1 cache and 2 capacity SSDs as shown in Figure 6. While the cluster was configured three different ways to run the test cases with different GPUs, the vSAN cluster configuration remained the same across all the tests.

**Figure 6**



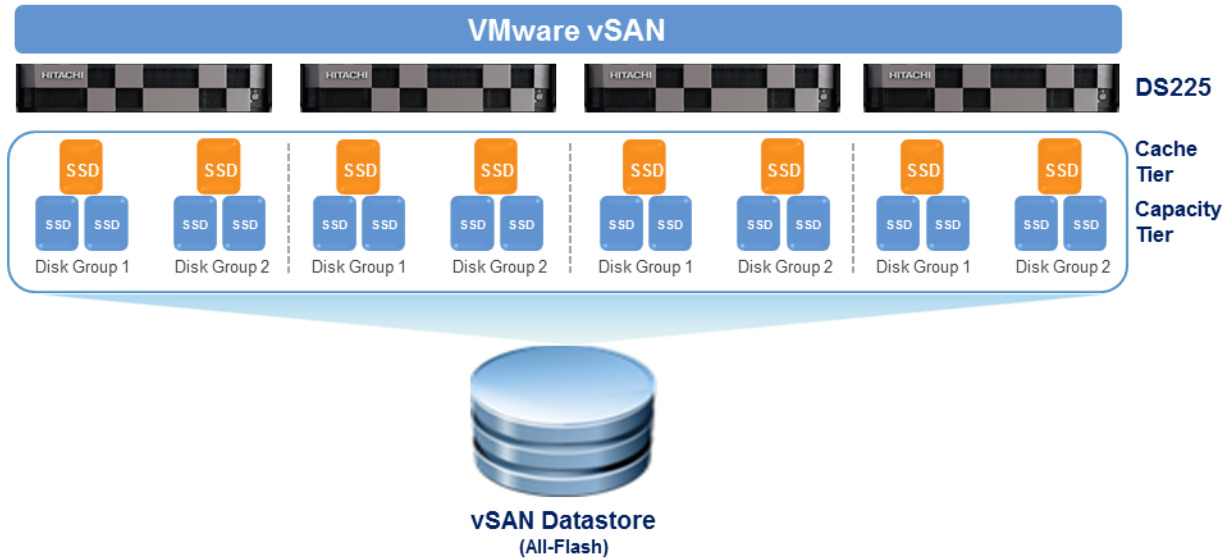Table 3 provides additional details about the vSAN cluster configuration.

TABLE 3. VSAN STORAGE CONFIGURATION

| UCP HC All Flash with GP-GPU | |
|---|---|
| Number of Nodes | 4 |
| vSAN Storage | 2 × Intel SSD DC S4600 Series (960 GB, 2.5"), for cache |
| | 4 × Intel SSD DC S4500 Series (1.92 TB, 2.5"), for capacity |
| Storage Policy | Built-in Storage Policy for VMware vSAN |
| Deduplication/compression | Disabled |

The desktop pools were created with Horizon's built-in storage policy for VMware vSAN. Figure 7 shows the set of storage policies that are configured automatically when using the built-in storage policy option from Horizon.

**Figure 7**



## Network Configuration

The UCP HC cluster was configured with vSphere Distributed Switch and each host with dual port 10 GbE network interfaces. Distributed port groups for management traffic, vSphere vMotion, vSAN and Virtual machine traffic were created as shown in Figure 8. Also, the jumbo frames option was enabled on the vSAN, vMotion vmkernel ports, and network devices.

**Figure 8**



## Graphics Acceleration Configuration

From a GPU perspective, UCP HC based on DS225 supports NVIDIA Tesla M10, Tesla M60, and Tesla P40 GPU devices. The number of GPU cards in a system varies depending on the GPU model (max 4 × M60 or 4 × P40 or 3 × M10). All ESXi hosts in the UCP HC cluster were configured with the latest NVIDIA Virtual GPU Manager software supported by VMware Horizon. The master image used for desktop pools was configured with the corresponding NVIDIA Graphics driver as well.

Table 4 shows the supported GPU cards, for additional details about these GPUs see "Appendix A: Specifications of Supported GPU Cards on Hitachi Advanced Server DS225" on page 37.

TABLE 4. SUPPORTED GPU CARDS ON HITACHI ADVANCED SERVER DS225

| Specification | Tesla M10 | Tesla M60 | Tesla P40 |
|---|---|---|---|
| GPUs per card | 4 × mid-range NVIDIA Maxwell GPUs | 2 × high-end NVIDIA Maxwell GPUs | 1 × high-end NVIDIA Pascal GPUs |
| CUDA cores | 2560 (640 per GPU) | 4096 (2048 per GPU) | 3840 |
| Memory size | 32GB (8GB per GPU) | 16GB (8GB per GPU) | 24GB |

**NVIDIA GRID License Server**

GRID vGPU is a license feature, when booted on a supported GPU without a license, users are warned each time a vGPU tries to obtain a license. A license server must be configured to provide licenses to a vGPU. For this test, the license server was configured on a Windows Server 2012 R2.

TABLE 5. LICENSE SERVER DETAILS

| Server Name | vCPU | Memory | OS Disk Size | Operating System |
|---|---|---|---|---|
| NVIDIA GRID License Server | 2 | 4 GB | 40 GB | Windows Server 2012 R2 |

See the GRID Licensing User Guide for more details about the different NVIDIA vGPU software license editions, other supported operating systems that can be used for the license server, and how to enable and use them on the supported GPUs.

**Horizon Infrastructure**

Table 6 shows the VMware Horizon View Components.

TABLE 6. VMWARE HORIZON VIEW COMPONENTS

| Server Name | vCPU | Memory | Disk Size | Disk Type | Operating System |
|---|---|---|---|---|---|
| View Connection Server | 4 | 16 GB | 40 GB | Eager Zeroed Thick | Windows Server 2012 R2 |
| Domain Controller | 4 | 8 GB | 40 GB | Eager Zeroed Thick | Windows Server 2012 R2 |
| Database Server | 4 | 16 GB | 40 GB (operating system) 60 GB (data) | Eager Zeroed Thick | Windows Server 2012 R2 Microsoft SQL Server 2014 SP1 |

The domain controller was deployed to support user authentication and domain services for the VMware Horizon infrastructure and Login VSI.

- The number of Horizon virtual desktops deployed is based on the number of GPU cards per server, and the NVIDIA vGPU Profile being used, see the Test Cases section for more details.

## Deploy the Solution

Deploying this solution requires doing the following procedures.

### *Deploy VMware Horizon Infrastructure Servers*

This process consists of the installation and configuration of the following core Horizon View components:

1. **View Connection Server**: The most important VMware View component is the View Connection Server. The View Connection Server is a connection broker that is responsible for authenticating clients and connecting them to the appropriate virtual desktop.

2. **View Administrator:** The View Administrator or management console is a web component for deploying and managing virtual desktops.

3. **View Client:** Establishes a connection from physical devices to a View Connection Server. View Client is installed on the user's devices (such as thin clients, zero clients, mobile devices, laptops, desktops, or any other devices supported by View Client).

4. **View agent:** The application installed on virtual desktops that allows VMware View to manage access from clients.

## UCP HC with GP-GPU Configuration Workflow

Configuration of the ESXi hosts and Virtual Machines varies based on the type of graphic acceleration used (vSGA, vDGA, or vGPU). There are four major areas to consider when configuring GPUs, some of the settings required for vGPU are covered in the following sections:

1. Prepare the host with a graphics driver

    - Verify that the NVIDIA GRID package is installed on the ESXi servers.

    - Configure ESXi Graphics Settings for vGPU

2. Prepare a GRID License Server.

    - Install and configure the NVIDIA GRID license server software in a server with a supported operating system for the graphics-enabled VMs to obtain a license and be entitled for vGPU.

3. Prepare the virtual machine with supported OS and graphics driver

    Consider the following when preparing a Virtual Machine to be used with VMware Horizon:

    - Create a Virtual Machine and install the Operating System.

    - Apply Virtual Machine settings depending on the graphic acceleration type being used (vSGA, vDGA  or vGPU). For vGPU, add a virtual PCI device, select NVIDIA GRID vGPU, and apply the required GPU profile.

    - For Guest Operating Systems, apply best practices and the OS optimization tool.

    - Install VMware Tools and Horizon View Agent.

    - Install the GPU drivers in the Guest OS. This driver should match the version of the driver installed on the ESXi host.

    - Install Microsoft Office and additional software required by Login VSI.

    - Add the GRID vGPU license.

    - Join the VM to the domain.

    - Shut down and take a snapshot of the VM.

4. Configure Horizon for 3D rendering

    - Configuring a 3D desktop pool is the same as a normal pool until you reach the desktop pool settings section.

    - When adding a desktop pool for vGPU, on the **Remote Display Protocol** section, select **NVIDIA GRID vGPU** for **3D Renderer**.

For detailed settings and best practices of how to prepare Virtual Machines or Configure Horizon Pools, follow VMware and NVIDIA documentation based on the type of graphics acceleration:

- [Deploying Hardware-Accelerated Graphics with VMware Horizon](#)

- [NVIDIA GRID vGPU Deployment Guide for VMware Horizon](#)

- [NVIDIA GRID vGPU Profile Sizing for Windows 10](#)

*Deploy Login VSI for Load Generation*

The launchers and Login VSI environment are configured and managed by a centralized management console. Login VSI Power Workload (with Pro Library) and Multimedia Workload were used:

- **Power Worker Workload**

  The Power Worker workload is a very intensive workload that puts maximum stress on the system. Many applications on this workload use larger files and higher resolution media that are provided by the extra content of the Login VSI Pro Library.

- **Multimedia Worker Workload**

  The Multimedia Worker workload places a more severe demand on the environment and is designed to stress the CPU and GPU.

## Testing and Results

This section describes the tests, tools, and configuration used to validate VMware Horizon Desktops on the UCP HC with Graphics Acceleration.

### Test Methodology

*Load Generation Tool*

Login VSI was used to generate heavy workload on the desktops. Login VSI launchers were each configured to initiate up to 16 sessions with the VMware Blast protocol to the VMware Horizon View Connection Server to simulate end-to-end execution of the entire VMware Horizon View infrastructure stack.

Login VSI Power Worker and Multimedia Worker profiles were used. Login VSI Pro Library (for Power Workers) and VSI Multimedia workload (for Multimedia Workers) packages were installed to test UCP HC with the graphics processing unit (GPU).

For the test cases with Multimedia Worker user workloads, the desktops were configured with access to the internet to enable the testing of the Google Earth application.

For all the tests we used the VMware Blast Extreme display protocol for the desktop sessions.

Table 7 lists the applications that Login VSI exercised during workload testing for both workload profiles.

TABLE 7. WORKLOAD APPLICATION DEFINITIONS

| Workload Type/ Applications | Power Workload (with Pro Library) | Multimedia Workload (GPU/CPU-intensive operations) |
|---|---|---|
| Applications Exercised | ▪ Adobe Acrobat<br>▪ Freemind/Java<br>▪ Microsoft Internet Explorer<br>▪ Microsoft Excel<br>▪ Microsoft Outlook<br>▪ Microsoft PowerPoint<br>▪ Microsoft Word<br>▪ Photo Viewer<br>▪ Simulated application installs<br>▪ 7-Zip Compression<br>▪ Microsoft Windows Media Player<br>▪ Login VSI Pro Library | ▪ Adobe Acrobat<br>▪ Google Chrome<br>▪ Google Earth<br>▪ Microsoft Excel<br>▪ HTML 5 3D Spinning Balls<br>▪ Microsoft Internet Explorer<br>▪ MP3<br>▪ Microsoft Outlook<br>▪ Microsoft PowerPoint<br>▪ Microsoft Word<br>▪ Streaming video |

*Desktop VM Configuration*

Table 8 shows the configuration profiles of the Windows 10 VM templates used for the Power and Multimedia worker workloads. Each image was optimized with the VMware OS Optimization tool in accordance with [Creating an Optimized Windows Image for a VMware Horizon Virtual Desktop](#) and in conformance with the Login VSI test standards.

TABLE 8. PROFILES FOR WINDOWS 10 VMS

| VM Profile | VM OS | vCPU | VM Memory (configured) | VM Memory (reserved) | Resolution |
|---|---|---|---|---|---|
| Power Worker (with Pro Library) | Windows 10 Enterprise 64-bit | 4 | 4 GB | 4 GB | 1920 × 1080 |
| Multimedia Worker (with Login VSI Multimedia Workload) | Windows 10 Enterprise 64-bit | 4 | 8 GB | 8 GB | 1920 × 1080 |

## Test Cases

The series of tests that were run for this reference architecture are designed to capture the performance capabilities and user experience (based on Login VSI workloads) of UCP HC with Graphics Acceleration. Table 9 shows a summary of the test cases, type of GPU used, and total of vGPU VMs tested in each case.

TABLE 9. TEST CASES

| Test Case | VM Profile/Login VSI Workload | Number of GPUs per Host | NVIDIA vGPU Profile | Density per host | Total VMs in the Cluster |
|---|---|---|---|---|---|
| Test Case 1 | Power Worker (with Pro Library) | 2 × M10 | M10-1Q | 64<br><br>(32 per physical GPU) | 256 |
| Test Case 2 | Multimedia Worker (with Login VSI Multimedia Workload) | 2 × M60 | M60-1Q | 32<br><br>(16 per physical GPU) | 128 |
| Test Case 3 | Multimedia Worker (with Login VSI Multimedia Workload) | 2 × P40 | P40-2Q | 24<br><br>(12 per physical GPU) | 96 |

**Test Case 1: Power Worker Desktops with NVIDIA GPU M10**

For this test the four hosts in the cluster were configured with 2 × NVIDIA M10 GPUs each. The Horizon desktop pool was created with 256 instant clones floating desktops running on a four-node vSAN cluster.

The virtual machine template used for the instant clones was configured for a Power Worker user workload as defined in Table 8, "Profiles for Windows 10 VMs," on page 16. The virtual machine was also configured with memory reserved as mandated when using NVIDIA GRID vGPUs.

Each host in the cluster ran exactly 64 desktops (32 per physical GPU). This number of VMs is based on the NVIDIA GRID vGPU profile. NVIDIA GRID vGPU profile "**M10-1Q**" was used. See "Appendix B: GRID vGPU Profiles" on page 38 for more details about the GRID vGPU profiles.

**Figure 9**



**Test Case 2: Multimedia Worker Desktop with NVIDIA GPU M60**

For this test the four hosts in the cluster were configured with 2 × NVIDIA M60 GPUs each. The Horizon desktop pool was created with 128 instant clones floating desktops running on a four-node vSAN cluster.

The virtual machine template used for the instant clones was configured for a Multimedia Worker user workload as defined in Table 8, "Profiles for Windows 10 VMs," on page 16. The virtual machine was also configured with memory reserved as mandated when using NVIDIA GRID vGPUs.

Each host in the cluster ran exactly 32 desktops (16 per physical GPU). This number of VMs is based on the NVIDIA GRID vGPU profile. The NVIDIA GRID vGPU profile "M60-1Q" was used. See "Appendix B: GRID vGPU Profiles" on page 38 for more details about the GRID vGPU profiles.

**Figure 10**

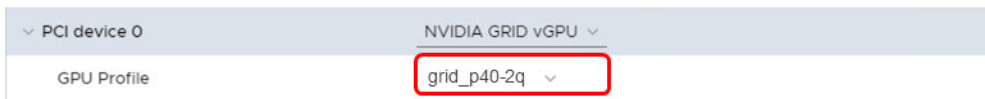| | |
|---|---|
| ∨ PCI device 0 | NVIDIA GRID vGPU ∨ |
| GPU Profile | grid_m60-1q ∨ |

**Test Case 3: Multimedia Worker Desktop with NVIDIA GPU P40**

For this test the four hosts in the cluster were configured with 2 × NVIDIA P40 GPUs each. The Horizon desktop pool was created with 96 instant clones floating desktops running on a four-node vSAN cluster.

The virtual machine template used for the instant clones was configured for a Multimedia Worker user workload as defined on Table 8, "Profiles for Windows 10 VMs," on page 16. The virtual machine was also configured with memory reserved as mandated when using NVIDIA GRID vGPUs.

Each host in the cluster ran exactly 24 desktops (12 per physical GPU). This number of VMs is based on the NVIDIA GRID vGPU profile. The NVIDIA GRID vGPU profile "P40-2Q" was used. See "Appendix B: GRID vGPU Profiles" on page 38 for more details about the GRID vGPU profiles.

**Figure 11**

| | |
|---|---|
| ∨ PCI device 0 | NVIDIA GRID vGPU ∨ |
| GPU Profile | grid_p40-2q ∨ |

## Test Case 1: Power Worker Desktops with NVIDIA GPU M10 (M10-1Q)

*Workload Testing*

Login VSI was configured on the Management Console to launch a Power Worker workload profile. The test was executed with the following configuration:

- The test was executed with 16 launchers with 16 sessions per launcher.

- Login VSI was configured to stagger logins of all 256 users with a single login occurring every 7.5 seconds. This ensured the 256 users logged into their desktops and began steady state workloads within a 32-minute period. The steady workload ran during a 20-minute period, then logged off all the sessions.

- The connection resolution used for the sessions was 1920 × 1080.
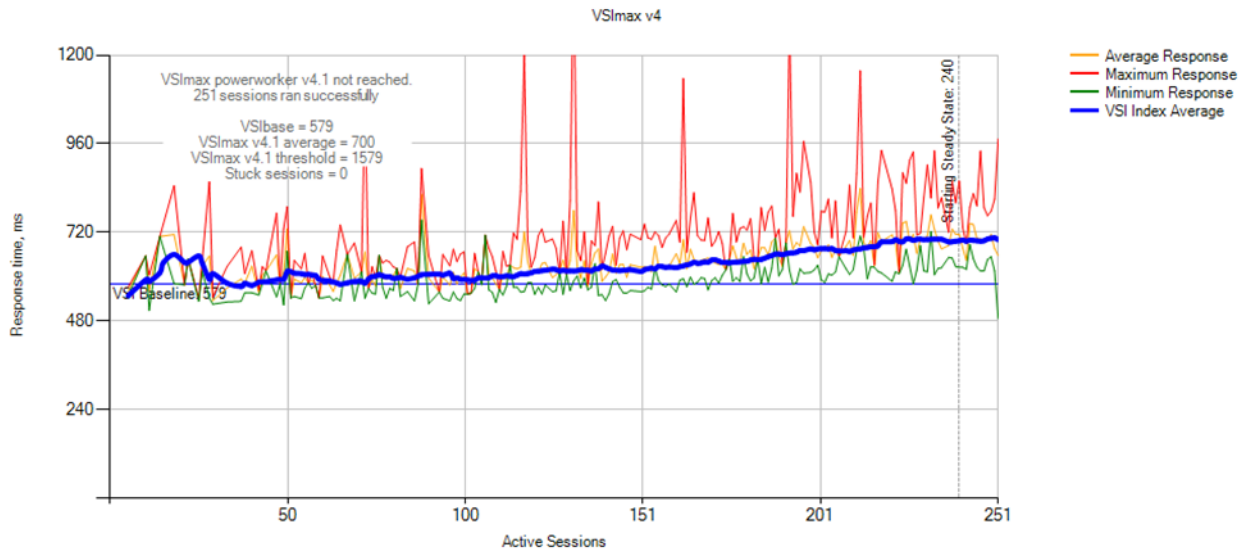
*Login VSI - Test Results*

The Login VSI Max user experience score shown in Figure 9 for the test with GPU M10 and Power Worker user workloads was not reached.

- The test completed with 251 (out of 256) successful power worker sessions.

- With this number of sessions, the maximum capacity VSImax (v4.1) was not reached and the Login VSI baseline performance score was **579**.

This indicates that the number of power workers tested did not put any strain on the UCP HC system resources.
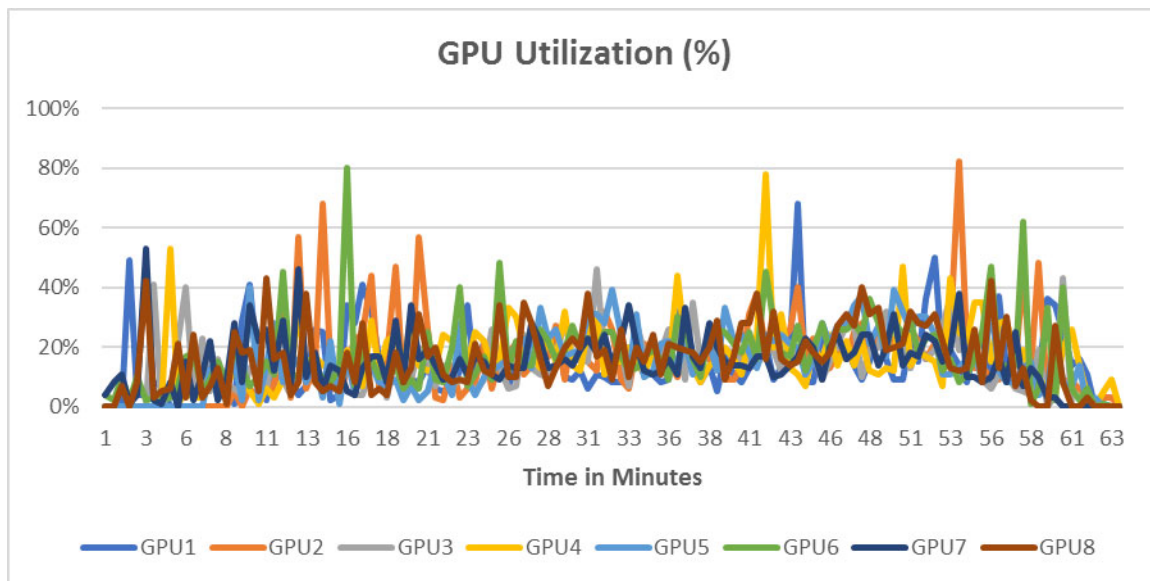
**Figure 12**



VSImax v4

VSImax powerworker v4.1 not reached.
251 sessions ran successfully

VSIbase = 579
VSImax v4.1 average = 700
VSImax v4.1 threshold = 1579
Stuck sessions = 0

VSI Baseline 579

Legend:
— Average Response
— Maximum Response
— Minimum Response
— VSI Index Average

Staring Steady State: 240

Y-axis: Response time, ms
X-axis: Active Sessions

*Compute Infrastructure - Performance Results*

**GPU Utilization**

The ESXi hosts in the cluster were populated with 2 × NVIDIA M10 GPU cards. Figure 13 shows the GPU utilization representative of 2 GPU cards (each M10 GPU card had 4 × GPUs) in one of the ESXi hosts during the test period.

- With a few peaks of 68% and 80% of GPU utilization for very short periods of time, the GPU utilization was below 50% during most of the test.

- This indicates that even while running the maximum of 256 vGPU VMs (based on M10-1Q GRID vGPU profile) with the Login VSI Power Worker workload, the GPUs still have some resources for additional workloads.
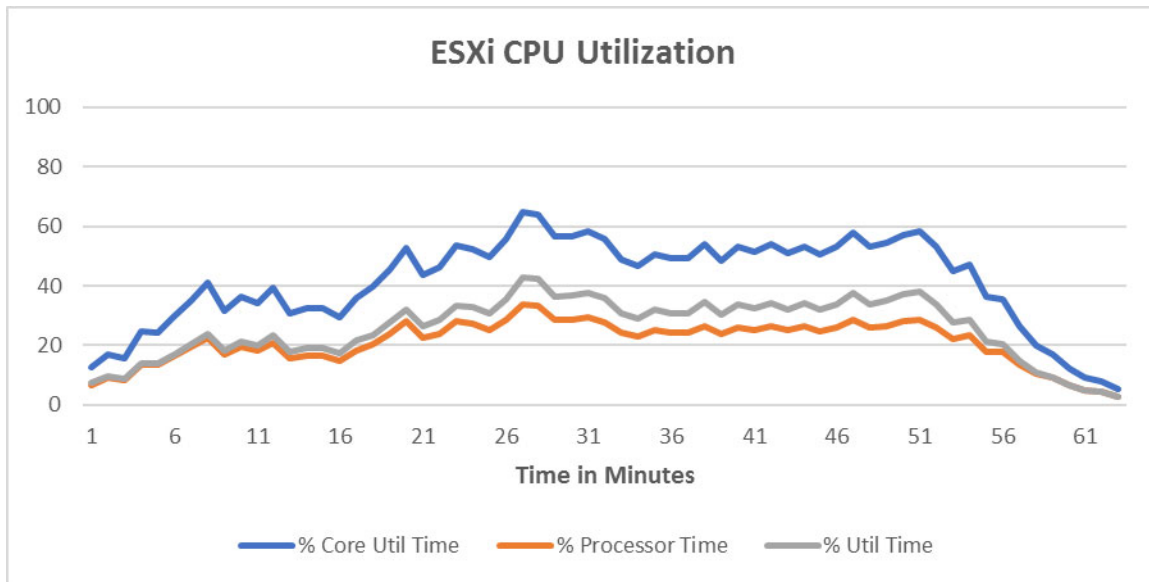
**Figure 13**



**Hypervisor CPU Performance**

Figure 14 illustrates the physical CPU utilization representative of one of the hosts in the cluster running Power Worker user workload during the login, steady state, and logoff operations.

- There are three operations that occur during this test executed by Login VSI:
  - The logon of all power users was completed in 32 minutes.
  - The steady state workload ran for the next 20 minutes.
  - It took an additional 8 minutes for logoff operations
- The performance metrics show the following:
  - Percent core utilization does not rise above 64% during login storm.
  - Percent core utilization does not rise above 58% during steady state.
  - This shows that there is still headroom from the CPU perspective on the UCP HC cluster to support bursts in workloads while maintaining acceptable end user performance.

A key aspect to note with these results is the fact that even when running multiple power workers with intensive workloads, the CPU usage does not increase much, and the this is because the video processing is being offloaded to the GPU.

**Figure 14**



ESXi CPU Utilization chart with Y-axis showing values 0 to 100 and X-axis labeled "Time in Minutes" ranging from 1 to 61. Legend shows: % Core Util Time (blue), % Processor Time (orange), % Util Time (gray).
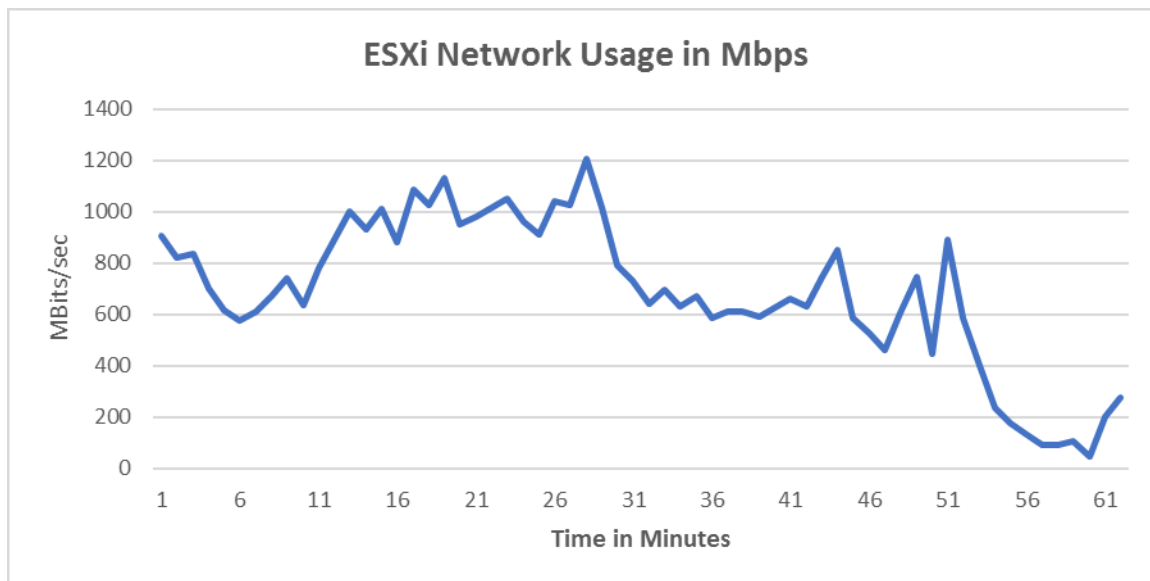
**Hypervisor Network Performance**

Figure 15 shows the network performance representative of one of the hosts in the cluster running a Power Worker user workload during the login, steady state, and logoff operations.

- The results show the total network usage which includes vSAN, management, and Virtual Machine traffic, for both received and transmitted Mbps.

- The busiest period for network traffic was during the first 30 minutes when all the power users were completing logon operations.

- The total network usage reached a peak of 1200 Mbps during the logon phase.

- The total network usage reached a peak of 890 Mbps during the steady state phase

- Considering that each hypervisor has 2 × 10 GbE ports, these results show that network bandwidth is not an issue and there is plenty of bandwidth available for additional workloads.
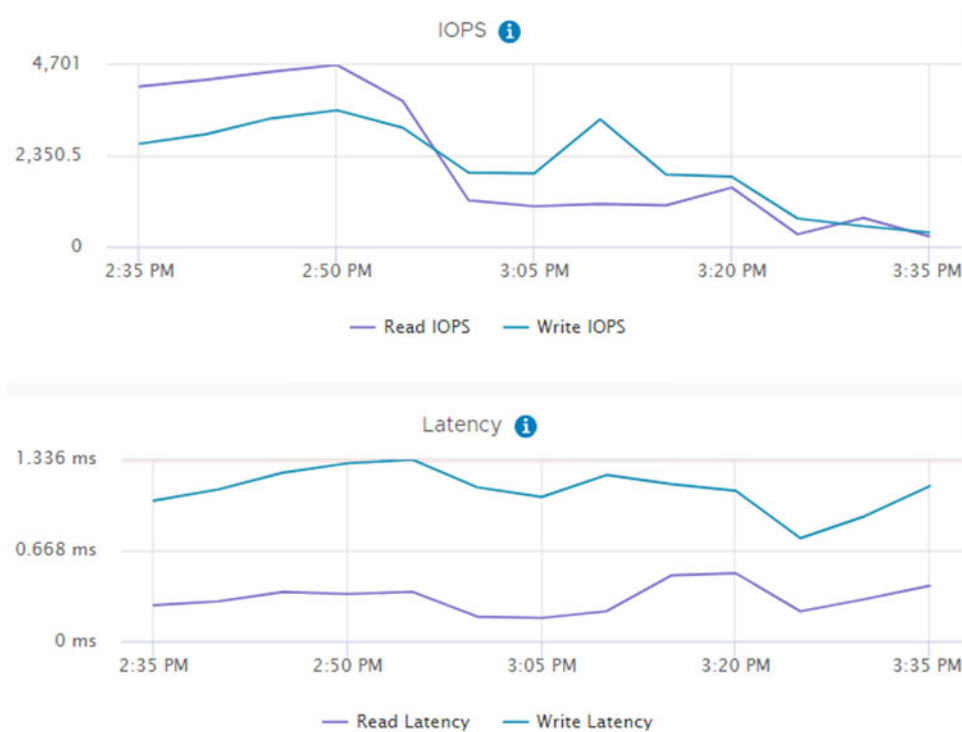
**Figure 15**

**vSAN Storage Performance**

Figure 16 shows IOPS and latency at the cluster level. These graphs are taken directly from the vCenter web client and show the IOPS and latency metrics for the vSAN cluster form the perspective of VM consumption during the logon, steady state, and logoff phases.

- The cluster reached a maximum peak of 4700 Write and 3500 Read IOPS during the logon phase.

- The disk latency during the entire period of the test did not raise above 1.336 ms for writes and 0.363 ms for read operations.

- During the steady state, the disk latency did not rise above 1.223 ms for writes and 0.219 ms for read operations.

- Considering that this is a UCP HC All-Flash vSAN cluster that is not even configured at its full capacity, these low latency results indicate that there is plenty of headroom for additional workloads.

**Figure 16**



**Test Case 2: Multimedia Worker Desktops with NVIDIA GPU M60 (M60-1Q)**

*Workload Testing*

Login VSI was configured on the Management Console to launch a Multimedia Worker workload profile. The test was executed with the following configuration:

- The test was executed with 8 launchers with 16 sessions per launcher.

- Login VSI was configured to stagger logins of all 128 users with a single login occurring every 7.5 seconds. This ensured the 256 users logged into their desktops and began steady state workloads within a 16-minute period. The steady workload ran during a 30-minute period, then logged off all the sessions.

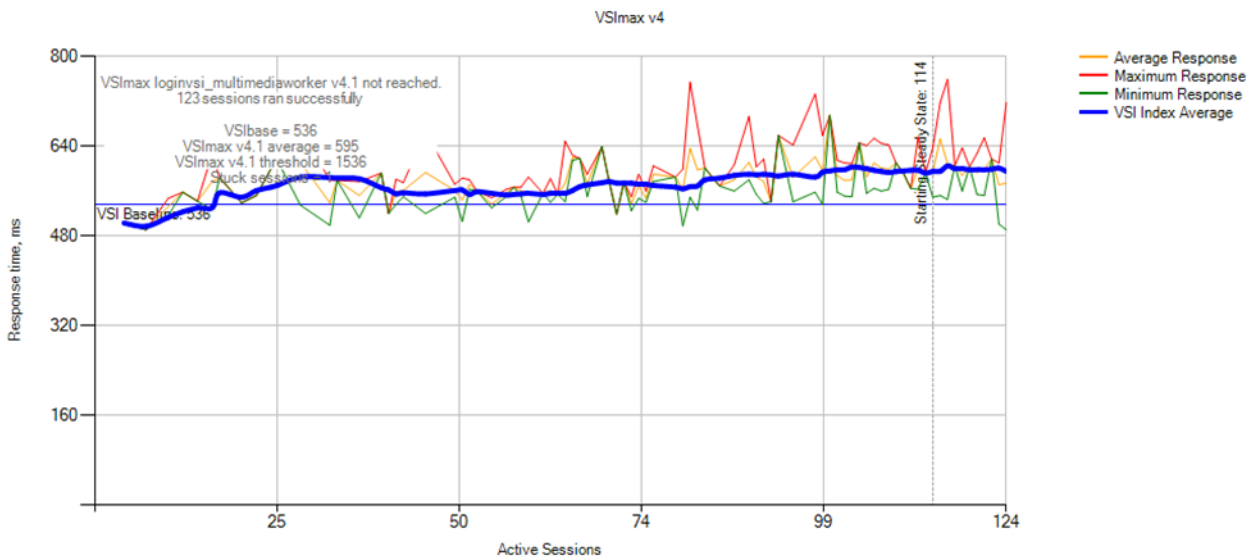- The connection resolution used for the sessions was 1920 × 1080.

## Login VSI - Test Results

The Login VSI Max user experience score shown in Figure 17 for this test with GPU M60 and Multimedia Worker user workloads was not reached.

- The test completed with 123 (out of 128) successful multimedia worker sessions.

- With this number of sessions, the maximum capacity VSImax (v4.1) was not reached and the Login VSI baseline performance score was 536.

This indicates that the number of power workers tested did not put any strain on the UCP HC system resources.
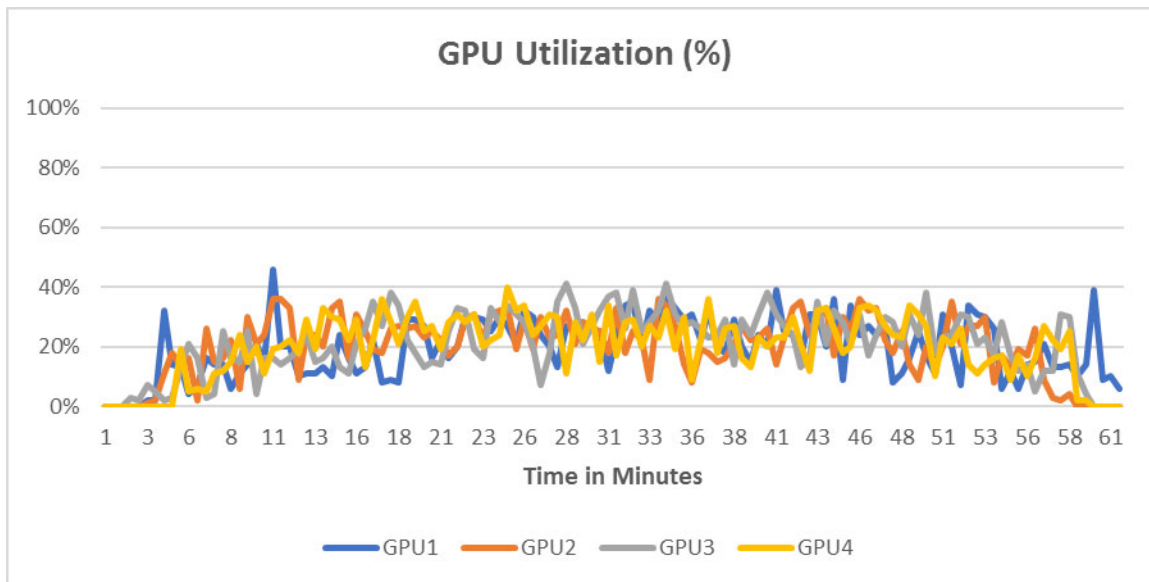
**Figure 17**

*Compute Infrastructure - Performance Results*

**GPU Utilization**

For this test the ESXi hosts in the cluster were populated with 2 × NVIDIA M60 GPU cards. Figure 18 shows the GPU utilization representative of 2 GPU cards (each M60 GPU card has 2 × GPUs) in one of the ESXi hosts during the test period.
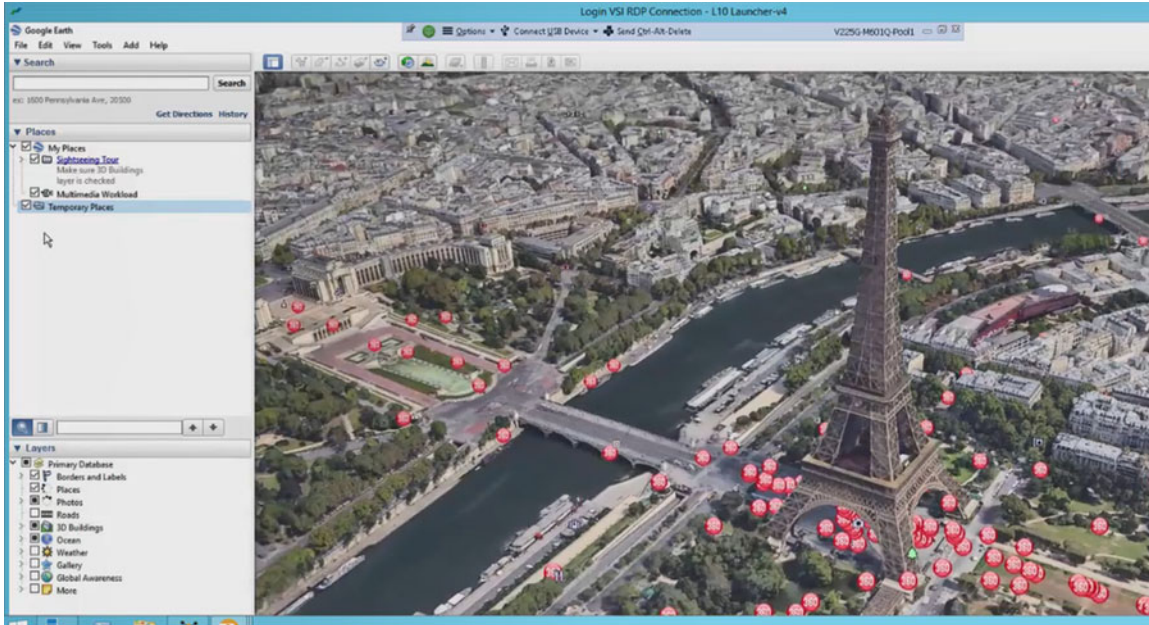
- With only one peak of 46% of GPU utilization during the logon phase, the GPU utilization for all the GPUs was below 40% during steady state phase.

- This indicates that even while running the maximum of 128 vGPU VMs (based on the M60-1Q GRID vGPU profile), the Login VSI Multimedia Workers workload did not drive the GPUs close to their maximum capacity.

**Figure 18**

Overall the user experience was excellent with the graphics-intensive applications like Google Earth, graphics heavy HTML5 websites, Microsoft PowerPoint with high resolution transitions, 1080p videos, and navigating across multiple applications. Figure 19 shows a snapshot of one of the Horizon desktops with GRID vGPU (M60-1Q profile) and graphics intensive Google Earth performing smoothly on UCP HC.
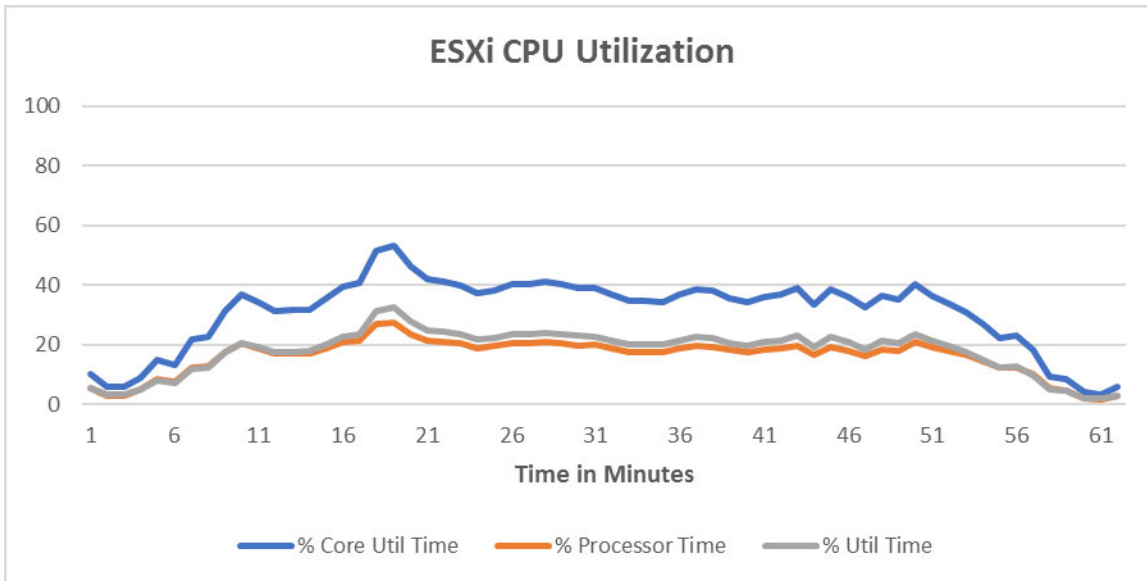
**Figure 19**



**Hypervisor CPU Performance**

Figure 20 illustrates the physical CPU utilization representative of one of the hosts in the cluster running a Multimedia Worker user workload during the login, steady state, and logoff operations.

- There are three operations that occurred during this test executed by Login VSI:
    - The logon of all multimedia users was completed in 16 minutes.
    - The steady state workload ran for the next 30 minutes.
    - It took an additional 8 minutes for logoff operations
- The performance metrics show the following:
    - Percent core utilization did not rise above 53%, and this happened at the beginning of the steady state.
    - Other than this peak utilization, the core utilization did not rise above 41% during the steady state.

This shows that there is still headroom from a CPU perspective on the UCP HC cluster to support bursts in workloads while maintaining acceptable end user performance. A key aspect to note here is that even when running graphics intensive workloads, the CPU does not increase much, because the video processing is being offloaded to the GPU.
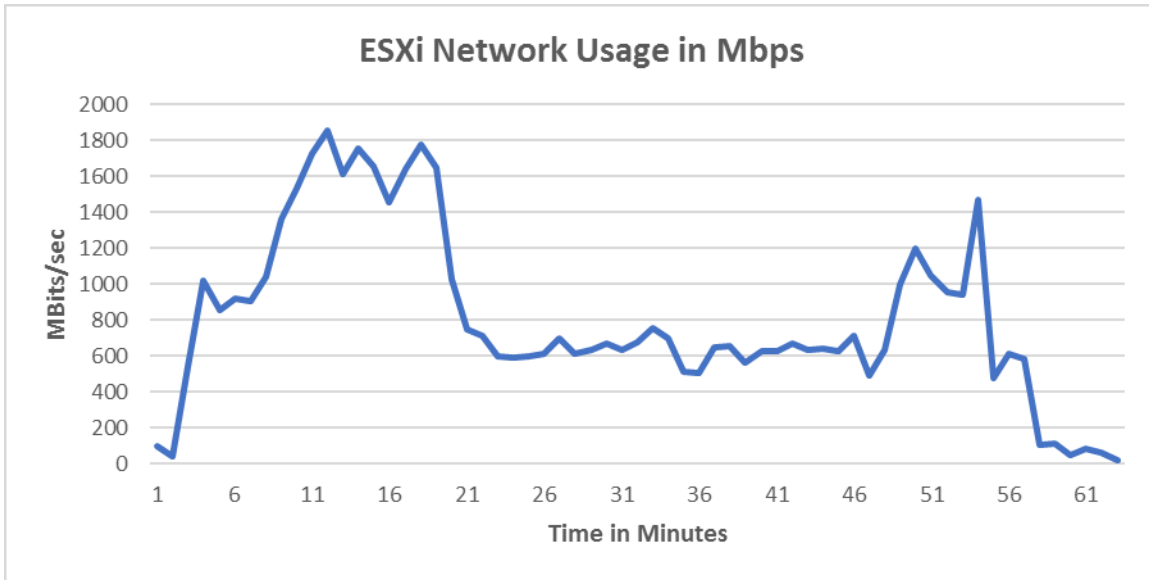
**Figure 20**



**Hypervisor Network Performance**

Figure 21 shows the network performance representative of one of the hosts in the cluster running a Multimedia Worker user workload during the login, steady state, and logoff operations.

- The results show the total network usage which includes vSAN, management, and Virtual Machine traffic, for both received and transmitted Mbps.

- The busiest period for network traffic was during the time that all the multimedia users were completing the logon operations.

- The total network usage reached a peak of 1854 Mbps during the logon phase.

- During most of the steady state period the total network usage was under 750 Mbps.

Considering that each hypervisor has 2 × 10 GbE ports, these results show that network bandwidth is not an issue and there is plenty of bandwidth available for additional workloads.

**Figure 21**



ESXi Network Usage in Mbps

**vSAN Storage Performance**

Figure 22 shows IOPS and latency at the cluster level These graphs are taken directly from vCenter web client and show the IOPS and latency metrics for the vSAN cluster from the perspective of VM consumption during the logon, steady state, and logoff phases.

- The cluster reached a maximum peak of 4755 write and 3000 read IOPS during the logon phase.

- The disk latency during the entire period of the test did not raise above 1.149 ms for writes and 0.650 ms for read operations.

- Considering that this is UCP HC All-Flash vSAN cluster is not even configured at its full capacity, these low latency results indicate that there is plenty of headroom for additional workloads.
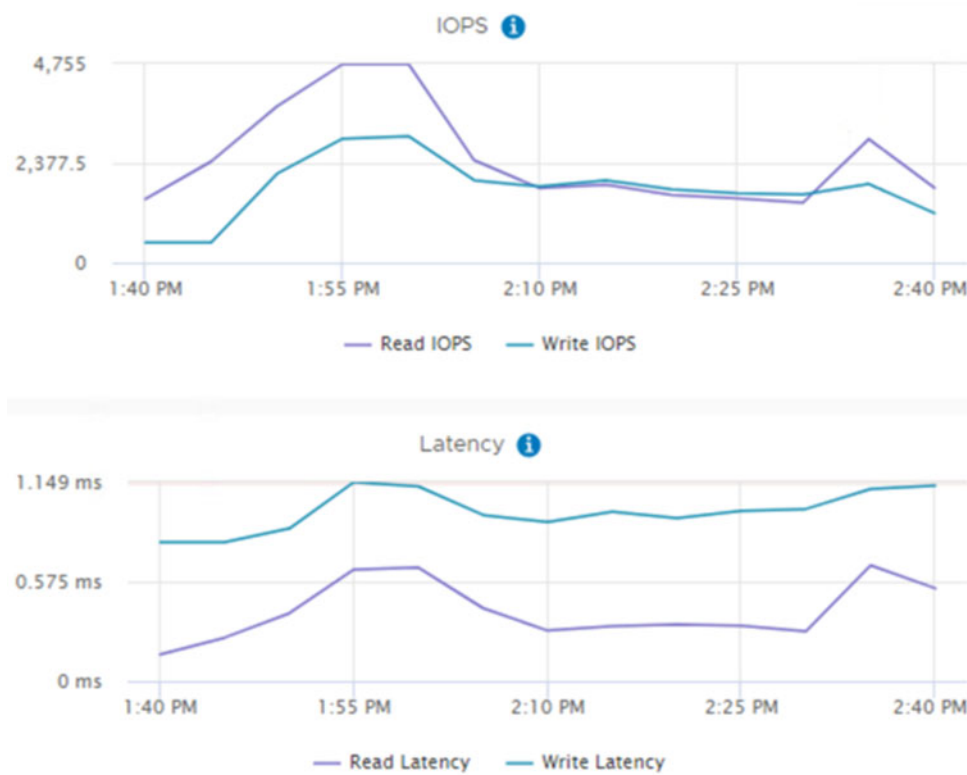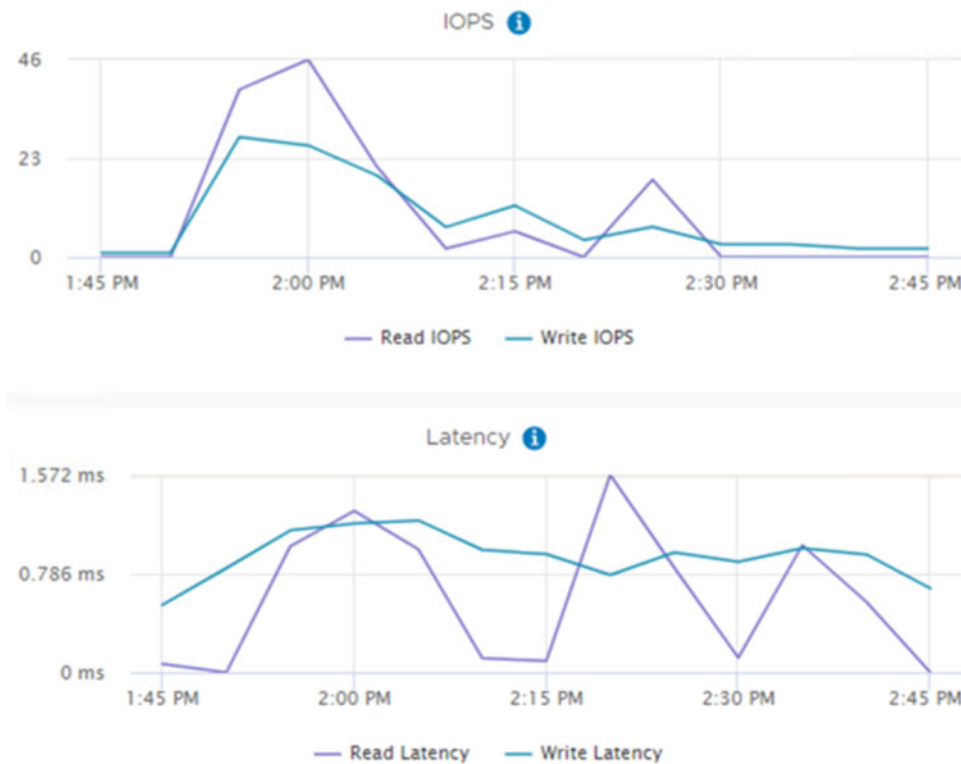
**Figure 22**

Figure 23 shows IOPS and latency representative of one of the desktops in the cluster running a Multimedia Worker user workload. These graphs are taken directly from vCenter web client and show the IOPS and latency from the perspective of VM consumption during the logon, steady state, and logoff phases.

- This VM reached a maximum peak of 46 write and 28 read IOPS during the test.

- The disk latency did not rise above 1.1 ms for writes and 1.6 ms for read operations.

**Figure 23**



## Test Case 3: Multimedia Worker Desktops with NVIDIA GPU P40 (P40-2Q)

*Workload Testing*

Login VSI was configured on the Management Console to launch a Multimedia Worker workload profile. The test was executed with the following configuration:

- The test was executed with 8 launchers with 12 sessions per launcher.

- Login VSI was configured to stagger logins of all 96 users with a single login occurring every 7.5 seconds. This ensured the 96 users logged into their desktops and began steady state workloads within a 12-minute period. The steady workload ran during a 30-minute period, then logged off all the sessions.

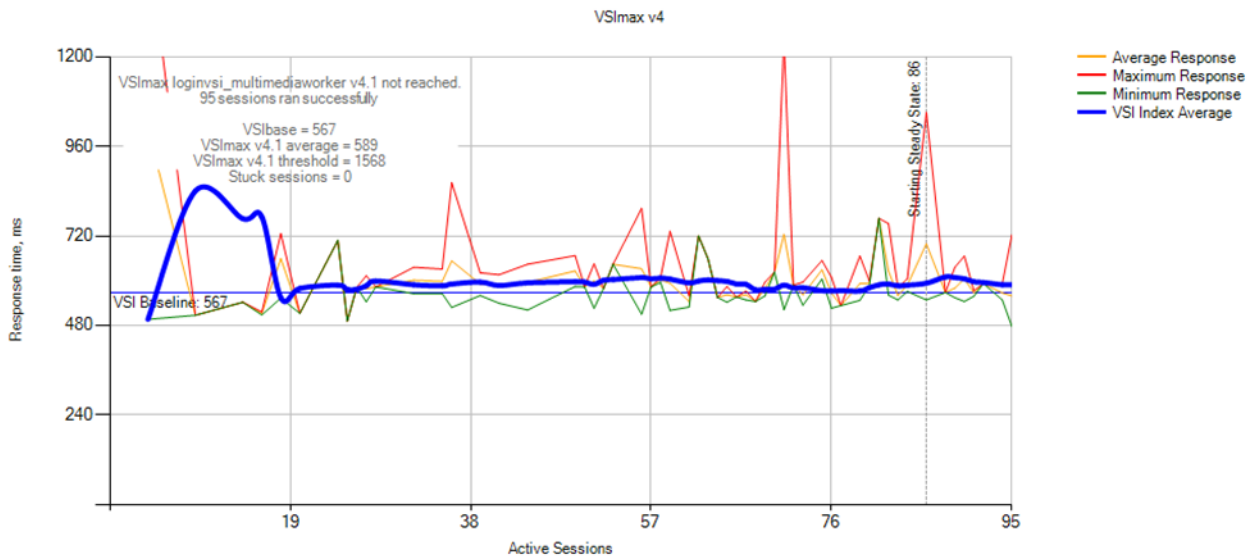- The connection resolution used for the sessions was 1920 × 1080.

*Login VSI - Test Results*

The Login VSI Max user experience score shown in Figure 24 for this test with GPU P40 and Multimedia Worker user workloads was not reached.

- The test completed with 95 (out of 96) successful multimedia worker sessions.

- With this number of sessions, the maximum capacity VSImax (v4.1) was not reached and the Login VSI baseline performance score was **567**.

This indicates that the number of power workers tested did not put any strain on the UCP HC system resources.

**Figure 24**



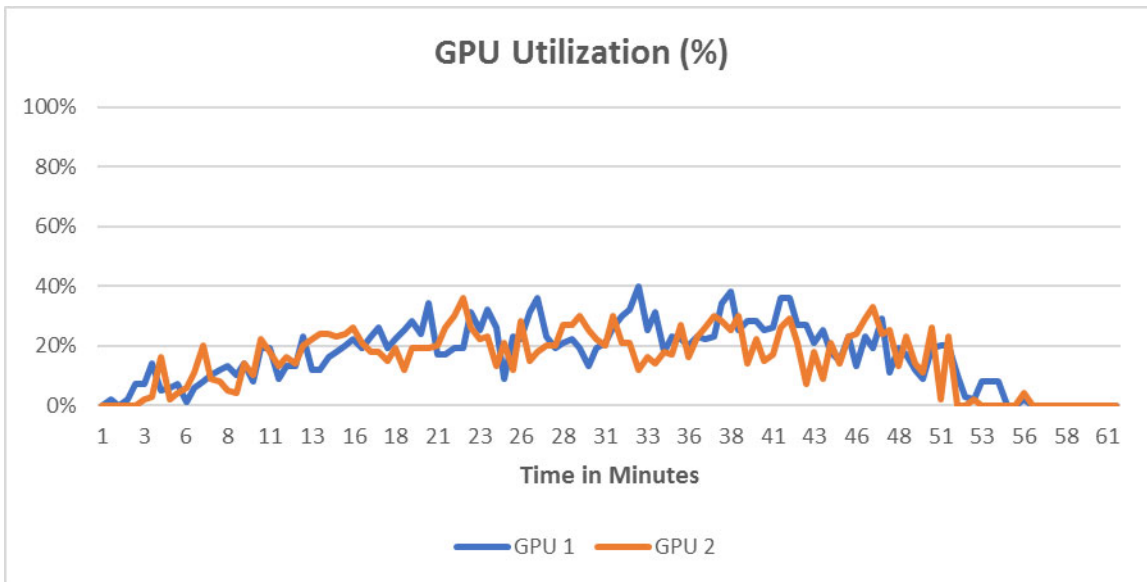*Compute Infrastructure - Performance Results*

**GPU Utilization**

For this test the ESXi hosts in the cluster were populated with 2 × NVIDIA P40 GPU cards. Figure 25 shows the GPU utilization representative of 2 GPU cards in one of the ESXi hosts during the test period.

- With only one peak of 40% during the test, the GPU utilization for all the GPUs was below 36% during the steady state phase.

This indicates that even while running the maximum of 96 vGPU VMs (based on the P40-2Q GRID vGPU profile), the Login VSI Multimedia Workers workload did not drive the GPUs close to even half of their maximum capacity.
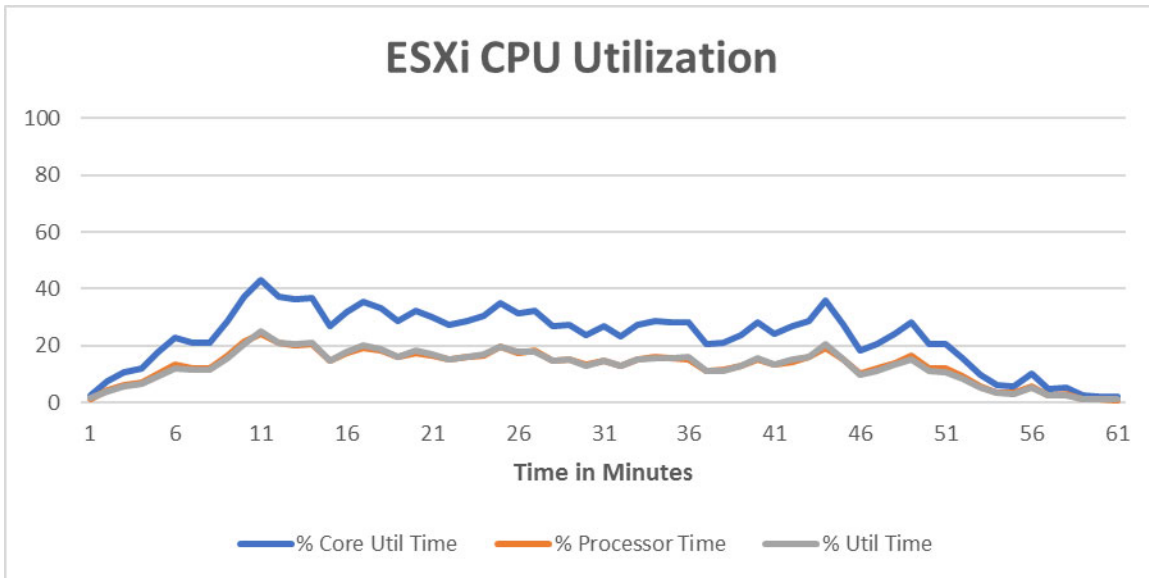
**Figure 25**



**Hypervisor CPU Performance**

Figure 26 illustrates the physical CPU utilization representative of one of the hosts in the cluster running a Multimedia Worker user workload during the login, steady state and logoff operations.

- There are three operations that occurred during this test executed by Login VSI:
    - The logon of all multimedia users was completed in 12 minutes.
    - The steady state workload ran for the next 30 minutes.
    - And there were additional minutes required for logoff operations.
- The performance metrics show the following:
    - Percent core utilization did not rise above 43%, and this happened at the end of the logon phase.
    - The percent core utilization did not rise above 36% during the steady state.

This shows that there is still headroom from a CPU perspective on the UCP HC cluster to support bursts in workloads while maintaining acceptable end user performance. A key aspect to note here is that even when running graphics intensive workloads, the CPU does not increase much, and this is because video processing is being offloaded to the GPU.

**Figure 26**



**Hypervisor Network Performance**

Figure 27 shows the network performance representative of one of the hosts in the cluster running a Multimedia Worker user workload during the login, steady state, and logoff operations.

- The results show the total network usage which includes vSAN, management, and Virtual Machine traffic, for both received and transmitted Mbps.

- The busiest period for network traffic was during the time that all the multimedia users were completing logon operations.

- The total network usage reached a peak of 1641 Mbps during the logon phase.

- During most of the steady state period the total network usage was under 820 Mbps.

Considering that each hypervisor has 2 × 10 GbE ports, these results show that network bandwidth is not an issue and there is plenty of bandwidth available for additional workloads.

**Figure 27**

ESXi Network Usage in Mbps

**vSAN Storage Performance**

Figure 28 shows IOPS and latency at the cluster level. These graphs are taken directly from vCenter web client and show the IOPS and latency metrics for the vSAN cluster from the perspective of VM consumption during the logon, steady state, and logoff phases.

- The cluster reached a maximum peak of 4969 write and 3139 read IOPS during the logon phase.

- The disk latency did not raise above 1.178 ms for writes and 0.472 ms for read operations during the logon and steady state phases.

- Considering that this is a UCH HC All-Flash vSAN cluster that is not even configured at its full capacity, these low latency results indicate that there is plenty of headroom for additional workloads.
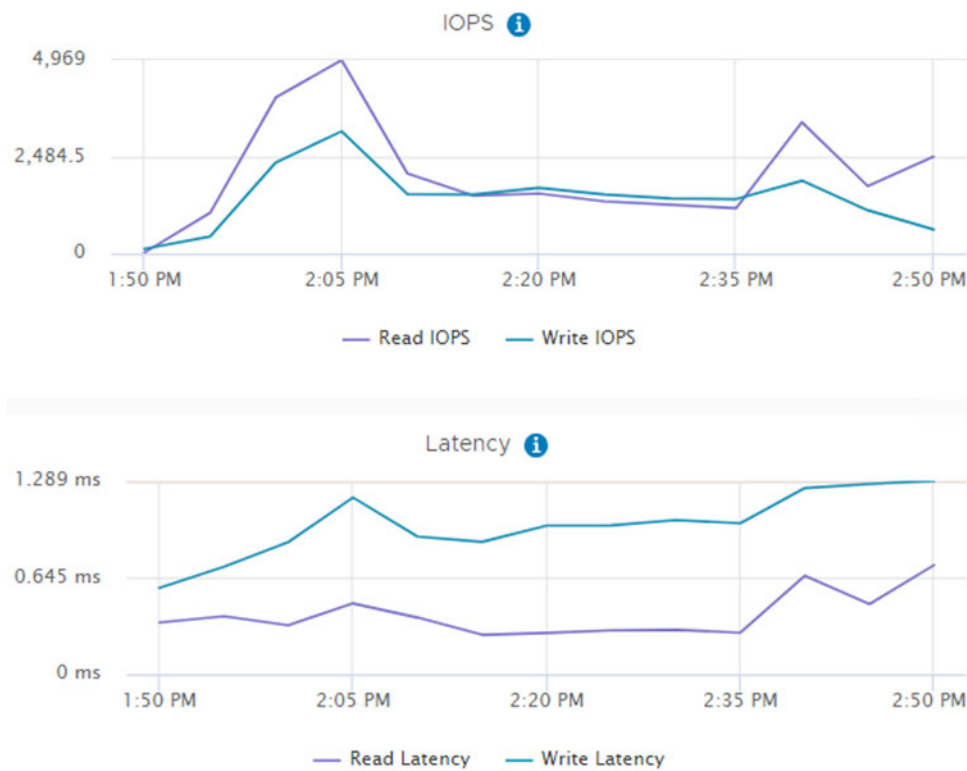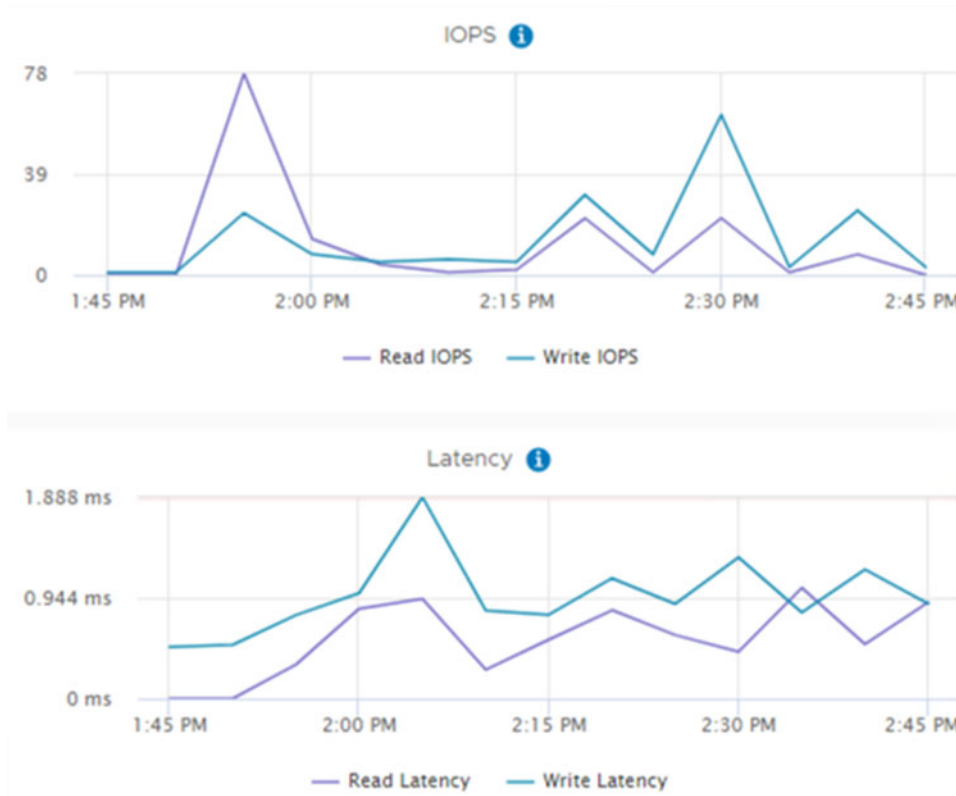
**Figure 28**

Figure 29 shows IOPS and latency representative of one of the desktops in the cluster running a Multimedia Worker user workload. These graphs are taken directly from vCenter web client and show the IOPS and latency from the perspective of VM consumption during the logon, steady state, and logoff phases.

- This VM reached a maximum peak of 78 write and 64 read IOPS during the test.

- The disk latency did not rise above 1.88 ms for writes and 1 ms for read operations.

**Figure 29**



## Conclusion

This architecture provides guidance to organizations implementing VMware Horizon with Graphics Acceleration on a UCP HC infrastructure, and describes tests performed by Hitachi Vantara to validate and measure the performance and graphics capabilities of the recommended solution, including third-party validated performance testing from Login VSI, the industry standard benchmarking tool for virtualized workloads.

Organizations looking for VDI solutions with 3D graphics support can confidently deploy UCP HC and ensure power users, designers, and multimedia users a superior user and graphics experience. Hitachi Vantara's market-leading hyperconverged infrastructure platform helps to deliver the promised benefits of GPU virtualization.

UCP HC for Horizon View with Graphics Acceleration provides:

- Highest density of physical GPU cards in a single server.

- Highest density of vGPU-enabled desktops per node in the hyperconverged infrastructure category.

- Graphics experience that is equivalent to dedicated hardware and delivered with the cost-effectiveness provided by GPU virtualization.

- Simplified deployment with hyperconverged Infrastructure.

- Ability to start small and scale out in affordable increments—from pilot to production.

- Independently validated, unmatched VDI performance for a superb end user experience.

- Enterprise-class data protection and resiliency

# Appendix A: Specifications of Supported GPU Cards on Hitachi Advanced Server DS225

TABLE 10. GPU CARD SPECIFICATIONS

| Specification | | Tesla M10 | Tesla M60 | Tesla P40 |
|---|---|---|---|---|
| GPUs per card | | 4 × mid-range NVIDIA Maxwell GPUs | 2 × high-end NVIDIA Maxwell GPUs | 1 × high-end NVIDIA Pascal GPUs |
| CUDA cores | | 2560 (640 per GPU) | 4096 (2048 per GPU) | 3840 |
| Memory size | | 32GB (8GB per GPU) | 16GB (8GB per GPU) | 24GB |
| Total board power | | 225 W | 300 W | 250 W |
| GPU clock | Idle | 405 MHz | 405 MHz | - |
| | Base | 1033 MHz | 899 MHz | 1303 MHz |
| | Maximum | 1306 MHz | 1178 MHz | 1531 MHz |
| Form factor | | Dual slot 10.5 inch | Dual slot 10.5 inch | Dual slot 10.5 inch |
| PCI Express interface | | x16 (Gen3) | x16 (Gen 3) | x16 (Gen 3) |
| Power connector | | 8-pin connector | 8-pin connector | 8-pin connector |
| Compatibility mode | | Graphics | Graphics (default) and Compute | Compute (default) and Graphics |
| Use case | | User Density-Optimized | Performance-Optimized | NVIDIA Quadro Virtual Data Center Workstation (Quadro vDWS) |

# Appendix B: GRID vGPU Profiles

Table 11, Table 12, and Table 13 show the NVIDIA vGPU profiles supported by different NVIDIA Tesla GPU's, and the maximum vGPU's supported in a single Hitachi Advanced Server DS225.

TABLE 11. TESLA M10 GRID VGPU PROFILES, 4 × PHYSICAL GPUS PER CARD

| vGPU Profile | Graphics Memory (Frame Buffer) | Virtual Display Heads | Maximum Resolution per Display Head | Maximum vGPUs per GPU | Maximum vGPUs per Card | Maximum vGPUs per DS225 Server (**3** Cards) |
|---|---|---|---|---|---|---|
| M10-8Q | 8192 | 4 | 4096×2160 | 1 | 4 | 12 |
| M10-4Q | 4096 | 4 | 4096×2160 | 2 | 8 | 24 |
| M10-2Q | 2048 | 4 | 4096×2160 | 4 | 16 | 48 |
| M10-1Q | 1024 | 2 | 4096×2160 | 8 | 32 | 96 |
| M10-0Q | 512 | 2 | 2560×1600 | 16 | 64 | 192 |
| M10-2B | 2048 | 2 | 4096×2160 | 4 | 16 | 48 |
| M10-1B | 1024 | 4 | 2560×1600 | 8 | 32 | 96 |
| M10-0B | 512 | 2 | 2560×1600 | 16 | 64 | 192 |
| M10-8A | 8192 | 1 | 1280×10241 | 1 | 4 | 12 |
| M10-4A | 4096 | 1 | 1280×10241 | 2 | 8 | 24 |
| M10-2A | 2048 | 1 | 1280×10241 | 4 | 16 | 48 |
| M10-1A | 1024 | 1 | 1280×10241 | 8 | 32 | 96 |

TABLE 12. TESLA M60 GRID VGPU TYPES, 2 × PHYSICAL GPUS PER CARD

| vGPU Profile | Graphics Memory (Frame Buffer) | Virtual Display Heads | Maximum Resolution per Display Head | Maximum vGPUs per GPU | Maximum vGPUs per Card | Maximum vGPUs per DS225 Server (**4** Cards) |
|---|---|---|---|---|---|---|
| M60-8Q | 8192 | 4 | 4096×2160 | 1 | 2 | 8 |
| M60-4Q | 4096 | 4 | 4096×2160 | 2 | 4 | 16 |
| M60-2Q | 2048 | 4 | 4096×2160 | 4 | 8 | 32 |
| M60-1Q | 1024 | 2 | 4096×2160 | 8 | 16 | 64 |

TABLE 12. TESLA M60 GRID VGPU TYPES, 2 × PHYSICAL GPUS PER CARD (CONTINUED)

| vGPU Profile | Graphics Memory (Frame Buffer) | Virtual Display Heads | Maximum Resolution per Display Head | Maximum vGPUs per GPU | Maximum vGPUs per Card | Maximum vGPUs per DS225 Server (**4** Cards) |
|---|---|---|---|---|---|---|
| M60-0Q | 512 | 2 | 2560×1600 | 16 | 32 | 128 |
| M60-2B | 2048 | 2 | 4096×2160 | 4 | 8 | 32 |
| M60-1B | 1024 | 4 | 2560×1600 | 8 | 16 | 64 |
| M60-0B | 512 | 2 | 2560×1600 | 16 | 32 | 128 |
| M60-8A | 8192 | 1 | 1280×10241 | 1 | 2 | 8 |
| M60-4A | 4096 | 1 | 1280×10241 | 2 | 4 | 16 |
| M60-2A | 2048 | 1 | 1280×10241 | 4 | 8 | 32 |
| M60-1A | 1024 | 1 | 1280×10241 | 8 | 16 | 64 |

TABLE 13. TESLA P40 GRID VGPU PROFILES, 1 × PHYSICAL GPU PER CARD

| vGPU Profile | Graphics Memory (Frame Buffer) | Virtual Display Heads | Maximum Resolution per Display Head | Maximum vGPUs per GPU | Maximum vGPUs per Card | Maximum vGPUs per DS225 Server (**4** Cards) |
|---|---|---|---|---|---|---|
| P40-24Q | 24576 | 4 | 4096×2160 | 1 | 1 | 4 |
| P40-12Q | 12288 | 4 | 4096×2160 | 2 | 2 | 8 |
| P40-8Q | 8192 | 4 | 4096×2160 | 3 | 3 | 12 |
| P40-6Q | 6144 | 4 | 4096×2160 | 4 | 4 | 16 |
| P40-4Q | 4096 | 4 | 4096×2160 | 6 | 6 | 24 |
| P40-3Q | 3072 | 4 | 4096×2160 | 8 | 8 | 32 |
| P40-2Q | 2048 | 4 | 4096×2160 | 12 | 12 | 48 |
| P40-1Q | 1024 | 2 | 4096×2160 | 24 | 24 | 96 |
| P40-2B | 2048 | 2 | 4096×2160 | 12 | 12 | 48 |
| P40-1B | 1024 | 4 | 2560×1600 | 24 | 24 | 96 |
| P40-24A | 24576 | 1 | 1280×10241 | 1 | 1 | 4 |
| P40-12A | 12288 | 1 | 1280x10241 | 2 | 2 | 8 |

TABLE 13. TESLA P40 GRID VGPU PROFILES, 1 × PHYSICAL GPU PER CARD (CONTINUED)

| vGPU Profile | Graphics Memory (Frame Buffer) | Virtual Display Heads | Maximum Resolution per Display Head | Maximum vGPUs per GPU | Maximum vGPUs per Card | Maximum vGPUs per DS225 Server (**4** Cards) |
|---|---|---|---|---|---|---|
| P40-8A | 8192 | 1 | 1280×10241 | 3 | 3 | 12 |
| P40-6A | 6144 | 1 | 1280×10241 | 4 | 4 | 16 |
| P40-4A | 4096 | 1 | 1280×10241 | 6 | 6 | 24 |
| P40-3A | 3072 | 1 | 1280×10241 | 8 | 8 | 32 |
| P40-2A | 2048 | 1 | 1280×10241 | 12 | 12 | 48 |
| P40-1A | 1024 | 1 | 1280×10241 | 24 | 24 | 96 |

More details about NVIDIA GRID vGPU profiles can be found on NVIDIA's documentation:

- [Virtual GPU Software User Guide](#).
- [NVIDIA Virtual GPU Technology](#)

**Hitachi Vantara**

Corporate Headquarters
5355 Augustine Drive
Santa Clara, CA 96054 USA
HitachiVantara.com | community.HitachiVantara.com

Contact Information
USA: 1-800-446-0744
Global: 1-858-547-4526
HitachiVantara.com/contact

MK-SL-128-00, February 2019