# Dell Wyse Datacenter

# Accelerated Graphics Reference Architecture for

**CITRIX**

XenDesktop

5/27/2014
Phase 6
Version 1.4

# Contents

# 1 Introduction

## 1.1 Purpose

This document describes:

1. Dell Wyse Datacenter for Citrix XenDesktop Reference Architecture specifically as it pertains to applicable pass-through and shared graphic acceleration options on NVIDIA GRID boards.

This document addresses the design, configuration and implementation considerations for the key components of the architecture required to deliver graphically-accelerated virtual desktops via XenDesktop 7 on VMware vSphere 5 and XenServer 6. Please refer to the Dell Wyse Datacenter (DWD) Reference Architecture for Citrix XenDesktop for more details on the broader solution stack and design considerations: Link

## 1.2 Scope

Relative to delivering the virtual desktop environment, the objectives of this document are to:

- Define the hardware requirements to support the design.
- Define the design constraints which are relevant to the design.
- Define relevant risks, issues, assumptions and concessions – referencing existing ones where possible.
- Provide a breakdown of the design into key elements such that the reader receives an incremental or modular explanation of the design.

## 1.3 What's New in This Release

- vGPU: Shared and pass-through graphics using NVIDIA Grid cards
- Updated Dell Wyse Cloud Clients

# 2 Solution Architecture Overview

## 2.1 Introduction

The Dell Wyse Datacenter Solution leverages a core set of hardware and software components consisting of 4 primary layers:

- Networking Layer
- Compute Server Layer
- Management Server Layer
- Storage Layer
- End Point Layer

These components have been integrated and tested to provide the optimal balance of high performance and lowest cost per user. Additionally, the Wyse Datacenter Solution includes an approved extended list of optional components in the same categories. These components give IT departments the flexibility to custom tailor the solution for environments with unique VDI feature, scale or performance needs. The Wyse Datacenter stack is designed to be a cost effective starting point for IT departments looking to migrate to a fully virtualized desktop environment slowly. This approach allows you to grow the investment and commitment as needed or as your IT staff becomes more comfortable with VDI technologies.

## 2.2  Physical Architecture Overview

The core Dell Wyse Datacenter architecture consists of two models: Local Tier1 and Shared Tier1. "Tier 1" in the Wyse Solutions Engineering context defines from which disk source the VDI sessions execute. Local Tier1 includes rack servers only while Shared Tier 1 can include rack or blade servers due to the usage of shared Tier 1 storage. Tier 2 storage is present in both solution architectures and, while having a reduced performance requirement, is utilized for user profile/data and Management VM execution. Management VM execution occurs using Tier 2 storage for all solution models. Dell Wyse Datacenter is a 100% virtualized solution architecture.

**Local Tier 1**

| MGMT Server | | Compute Server | | |
|---|---|---|---|---|
| CPU | RAM | CPU | **VDI Disk** | RAM |

Mgmt VMs

VDI VMs

| **Mgmt Disk** | User Data |
|---|---|

T2 Shared Storage

In the Shared Tier 1 solution model, an additional high-performance shared storage array is added to handle the execution of the VDI sessions. All Compute and Management layer hosts in this model are diskless.

**Shared Tier 1**

| MGMT Server | | Compute Server | |
|---|---|---|---|
| CPU | RAM | CPU | RAM |

Mgmt VMs

VDI VMs

| **Mgmt Disk** | User Data |
|---|---|

| **VDI Disk** |
|---|

T2 Shared Storage

T1 Shared Storage

## 2.3 Dell Wyse Datacenter – Solution Layers

Only a single high performance Dell Force10 48-port switch is required to get started in the Network layer. This switch will host all solution traffic consisting of 1Gb iSCSI and LAN sources for smaller stacks. For 500 users or above we recommend that LAN and iSCSI traffic be separated into discrete switching fabrics. Additional switches are added and stacked as required to provide High Availability for the Network layer.



The Compute layer consists of the server resources responsible for hosting the XenDesktop user sessions, hosted via VMware vSphere for local or shared tier 1 solution models (local tier 1 pictured below).



VDI management components are dedicated to their own layer so as to not negatively impact the user sessions running in the compute layer. This physical separation of resources provides clean, linear, and predictable scaling without the need to reconfigure or move resources within the solution as you grow. The Management layer will host all the VMs necessary to support the VDI infrastructure.



The Storage layer consists of options provided by EqualLogic for iSCSI and Compellent arrays for Fiber Channel to suit your Tier 1 and Tier 2 scaling and capacity needs.



## 2.4 Local Tier 1

### 2.4.1 Local Tier 1 – Rack (iSCSI – EQL)

The Local Tier 1 solution model provides a scalable rack-based configuration that hosts user VDI sessions on local disk in the Compute layer. vSphere or XenServer based solutions are available and scale based on the chosen hypervisor.

## 2.4.1.1 Local Tier 1 – Network Architecture (iSCSI)

In the local tier 1 architecture, a single Force10 switch is shared among all network connections for both Management and Compute, up to 500 users. For 500 or more users Dell Wyse Solutions Engineering recommends separating the network fabrics to isolate iSCSI and LAN traffic as well as making each switch stack redundant. Only the Management servers connect to iSCSI storage in this model. All Top of Rack (ToR) traffic has been designed to be layer 2/ switched locally, with all layer 3/ routable VLANs trunked from a core or distribution switch. The following diagrams illustrate the logical data flow in relation to the core switch.

## 2.4.1.2 Local Tier 1 Cabling (Rack – HA)



## 2.5 Shared Tier 1

### 2.5.1 Shared Tier 1 – Rack (iSCSI – EQL)

For POCs or small deployments, Tier1 and Tier2 are combined on a single 6210XS storage array. Above 500 users, a separate array is used for Tier 2.

#### 2.5.1.1 Shared Tier 1 Rack – Network Architecture (iSCSI)

In the Shared Tier 1 architecture for rack servers, both Management and Compute servers connect to shared storage in this model. All ToR traffic has designed to be layer 2/ switched locally, with all layer 3/ routable VLANs routed through a core or distribution switch. The following diagrams illustrate the server NIC to ToR switch connections, vSwitch assignments, as well as logical VLAN flow in relation to the core switch.



#### 2.5.1.2 Shared Tier 1 Cabling – (Rack – EQL)



### 2.5.2 Shared Tier 1 – Rack (FC – CML)

Utilizing Compellent storage for Shared Tier 1 provides a fiber channel solution where Tier 1 and Tier 2 are functionally combined in a single array. Tier 2 functions (user data + Mgmt VMs) are removed from the array if the customer has another solution in place. Doing this will net an additional resource capability per Compellent array for Tier 1 user desktop sessions. Scaling this

solution is very linear by predictably adding Compellent arrays for every 2000 Standard users, on average.



### 2.5.2.1 Shared Tier 1 Rack – Network Architecture (FC)

In the Shared Tier 1 architecture for rack servers using fiber channel, a separate switching infrastructure is required for FC. Management and Compute servers will both connect to shared storage using FC. Both Management and Compute servers connect to all network VLANs in this model. All ToR traffic has designed to be layer 2/ switched locally, with all layer 3/ routable VLANs routed through a core or distribution switch. The following diagrams illustrate the server NIC to ToR switch connections, vSwitch assignments, as well as logical VLAN flow in relation to the core switch.

### 2.5.2.2 Shared Tier 1 Cabling (Rack – CML)

# 3  Hardware Components

## 3.1  Servers

### 3.1.1  PowerEdge R720

The rack server platform for the DWD solution is the best-in-class Dell PowerEdge R720 (12G). This dual socket CPU platform runs the fastest Intel Xeon E5-2600 family of processors, can host up to 768GB RAM, and supports up to 16 2.5" SAS disks. The graphics-enabled Dell PowerEdge R720 offers uncompromising performance and scalability in a 2U form factor. For more information, please visit: Link



### 3.1.2  Dell Precision R7610

The rack workstation platform for the DWD solution is the top-of-the-line Dell Precision R7610. Add a Dell Precision R7610 rack Workstation to your data center for exceptional performance and security. The Dell Precision R7610 keeps your valuable data in your data center at all times and allows secure access only through soft- or zero-client devices. This rack Workstation solution also reduces the potential liability of lost or damaged physical assets in the field. For more information, please visit: Link



## 3.2  GPUs

### 3.2.1  NVIDIA Grid K1 and K2 Boards

NVIDIA GRID™ technology offers the ability to offload graphics processing from the CPU to the GPU in virtualized environments, allowing the data center manager to deliver true PC graphics-rich experiences to more users for the first time. NVIDIA's Kepler™-based GRID K1 and K2 boards are specifically designed to enable rich graphics in virtualized environments.

**GPU Virtualization**

GRID boards allow hardware virtualization of the GPU. This means multiple users can share a single GPU, improving

user density while providing true PC performance and compatibility.

**Low-Latency Remote Display**

NVIDIA's patented low-latency remote display technology greatly improves the user experience by reducing the lag that users feel when interacting with their virtual machine. With this technology, the virtual desktop screen is pushed directly to the remoting protocol.

**Maximum User Density**

NVIDIA GRID boards have an optimized multi-GPU design that helps to maximize user density. GRID K1 boards, which include four Kepler-based GPUs and 16GB of memory, are designed to host the maximum number of concurrent users. GRID K2 boards, which include two higher end Kepler GPUs and 8GB of memory, deliver maximum density for users of graphics-intensive applications.

| | Grid K1 | Grid K2 |
| --- | --- | --- |
| Number of GPUs | 4 x entry Kepler GPUs (GK107) | 2 x high-end Kepler GPUs (GK104) |
| Total NVIDIA CUDA cores | 768 (192 per GPU) | 3072 (1536 per GPU) |
| Core Clock | 850 MHz | 745 MHz |
| Total memory size | 16 GB DDR3 | 8 GB GDDR5 |
| Max power | 130 W | 225 W |
| Form Factors | Dual slot (4.4" x 10.5") | Dual slot (4.4" x 10.5") |
| Aux power | 6-pin connector | 8-pin connector |
| PCIe | x16 (Gen3) | x16 (Gen3) |
| Cooling solution | Passive | Passive/ Active |

For more information on NVIDIA Grid, please visit: Link

## 3.3  Dell Wyse Cloud Clients



The following Dell Wyse clients will deliver a superior Citrix user experience and are the recommended choices for this solution.

### 3.3.1  Windows Embedded 7 – Z90Q7

The Dell Wyse Z90Q7 is a super high-performance Windows Embedded Standard 7 thin client for virtual desktop environments. Featuring a quad-core AMD processor and an integrated graphics engine that significantly boost performance, the Z90Q7 achieves exceptional speed and power for the most demanding VDI and embedded Windows applications, rotational 3D graphics, 3D simulation and modeling, unified communications, and multi-screen HD multimedia. Take a unit from box to productivity in minutes. Just select the desired configuration and the Z90Q7 does the rest automatically—no need to reboot. Scale to tens of thousands of endpoints with Dell Wyse

WDM software or leverage your existing Microsoft System Center Configuration Manager platform. The Z90Q7 is the thin client for power users who need workstation-class performance on their desktop or within a desktop virtualization environment. For more information, please visit: Link

### 3.3.2 Windows Embedded 8 – Z90Q8

Dell Wyse Z90Q8 is a super high-performance Windows Embedded 8 Standard thin client for virtual desktop environments. Featuring a quad-core AMD processor, the Z90Q8 offers a vibrant Windows 8 experience and achieves exceptional speed and power for the most demanding embedded Windows applications, rich 3D graphics and HD multimedia. And you can scale to tens of thousands of Z90Q8 endpoints with Dell Wyse Device Manager (WDM) software, or leverage your existing Microsoft System Center Configuration Manager platform. With single-touch or multi-touch capable displays, the Z90Q8 adds the ease of an intuitive touch user experience. The Z90Q8 is an ideal thin client for offering a high-performance Windows 8 experience with the most demanding mix of virtual desktop or cloud applications. For more information please visit: Link

### 3.3.3 Suse Linux – Z50D

Designed for power users, the new Dell Wyse Z50D is the highest performing thin client on the market. Highly secure and ultra-powerful, the Z50D combines Dell Wyse-enhanced SUSE Linux Enterprise with a dual-core AMD 1.65 GHz processor and a revolutionary unified engine for an unprecedented user experience. The Z50D eliminates performance constraints for high-end, processing-intensive applications like computer-aided design, multimedia, HD video and 3D modeling. Scalable enterprise-wide management provides simple deployment, patching and updates. Take a unit from box to productivity in minutes with auto configuration. Delivering unmatched processing speed and power, security and display performance, it's no wonder no other thin client can compare. For more information, please visit: Link

### 3.3.4 Wyse Xenith 2

Establishing a new price/performance standard for zero clients for Citrix, the new Dell Wyse Xenith 2 provides an exceptional user experience at a highly affordable price for Citrix XenDesktop and XENAPP environments. With zero attack surface, the ultra-secure Xenith 2 offers network-borne viruses and malware zero target for attacks. Xenith 2 boots up in just seconds and delivers exceptional performance for Citrix XenDesktop and XENAPP users while offering usability and management features found in premium Dell Wyse cloud client devices. Xenith 2 delivers outstanding performance based on its system-on-chip (SoC) design optimized with its Dell Wyse zero architecture and a built-in media processor delivers smooth multimedia, bi-directional audio and Flash playback. Flexible mounting options let you position Xenith 2 vertically or horizontally on your desk, on the wall or behind your display. Using about 7 Watts of power in full operation, the Xenith 2 creates very little heat for a greener, more comfortable working environment. For more information, please visit: Link

### 3.3.5 Xenith Pro 2

Dell Wyse Xenith Pro 2 is the next-generation zero client for Citrix HDX and Citrix XenDesktop, delivering ultimate performance, security and simplicity. With a powerful dual core AMD G-series processor, Xenith Pro 2 is faster than competing devices. This additional computing horsepower allows dazzling HD multimedia delivery without overtaxing your server or network. Scalable enterprise-wide management provides simple deployment, patching and updates—your Citrix XenDesktop server configures it out-of-the-box to your preferences for plug-and-play speed and ease of use. Completely virus and malware immune, the Xenith Pro 2 draws under 9 watts of power in full operation—that's less than any PC on the planet. For more information please visit: Link

# 4   Solution Architecture for XenDesktop 7

## 4.1   XenDesktop with HDX 3D Pro

XenDesktop with HDX 3D Pro is a desktop and app virtualization solution that supports high-end designers and engineers of 3D professional graphics applications and provides cost-effective support to viewers and editors of 3D data. With XenDesktop, you can deliver a persistent user experience and leverage other virtualization benefits such as single-image management and improved data security.

Use HDX 3D technologies with:

- Computer-aided design, manufacturing, and engineering (CAD/CAM/CAE) applications

- Geographical information system (GIS) software

- Picture Archiving Communication System (PACS) workstations for medical imaging

- Latest OpenGL, DirectX, CUDA and CL versions supported

- Latest NVIDIA Grid cards

- Shared or dedicated GPUs or a mix of both

### 4.1.1   Virtual Shared Graphics Acceleration (vSGA)

vSGA provides the ability for multiple VMs to share physical hardware GPUs for 3D acceleration. vSGA supports a maximum of 512MB of virtual memory assignable to each VM, half of which comes directly from the physical GPU. The limit of the number of VMs that can share a GPU is determined by the amount of available GPU RAM as well as other traditional mitigating factors in VDI, such as host CPU and RAM. GPU hardware resources are reserved using a first-come, first-served basis during VM power on.

## 4.1.2 Virtual Dedicated Graphics Acceleration (vDGA)

VMware virtual dedicated graphics acceleration (vDGA), also known as pass through graphics support, refers to the technology of mapping a virtual desktop directly to a physical GPU on a high-end display adapter such as an NVIDIA GRID K1 or K2 card. These adapters, when used in vDGA mode, enable the user of a virtual desktop to run high-end, graphics-intensive applications such as CAD or other graphics editing and authoring applications. This solution offers a greater density of high-end graphics users per server when compared to the traditional one-to-one model of a dedicated graphics workstation.



## 4.1.3 Virtual GPU (vGPU)

NVIDIA GRID™ vGPU™ brings the full benefit of NVIDIA hardware-accelerated graphics to virtualized solutions. With GRID vGPU technology, the graphics commands of each virtual machine are passed directly to the GPU, without translation by the hypervisor. This allows the GPU hardware to be time-sliced to deliver the ultimate in shared virtualized graphics performance. Multiple vGPU profiles are available.

vGPU Manager enables up to eight users to share each physical GPU, assigning the graphics resources of the available GPUs to virtual machines in a balanced approach. Each NVIDIA GRID K1 card has up to four GPUs, allowing up to 32 users to share a single card.

| NVIDIA GRID Graphics Board | Virtual GPU Profile | Application Certifications | Graphics Memory | Max Displays Per User | Max Resolution Per Display | Max Users Per Graphics Board | Use Case |
|---|---|---|---|---|---|---|---|
| GRID K2 | K260Q | ✔ | 2,048 MB | 4 | 2560x1600 | 4 | Designer/Power User |
| | K240Q | ✔ | 1,024 MB | 2 | 2560x1600 | 8 | Designer/Power User |
| | K200 | | 256 MB | 2 | 1900x1200 | 16 | Knowledge Worker |
| GRID K1 | K140Q | ✔ | 1,024 MB | 2 | 2560x1600 | 16 | Power User |
| | K100 | | 256 MB | 2 | 1900x1200 | 32 | Knowledge Worker |

For more information please visit: Link

### 4.1.4 Virtualized Graphics Comparison

| Name | Description | Considerations |
|---|---|---|
| **SVGA** (Super Video Graphics Array) | VMware WDDM (Windows Display Driver Model) 1.1-compliant driver | Software rendering for 2D and 3D. Live migration is supported. |
| **vSGA** (Virtual <u>Shared</u> Graphics Acceleration) | Multiple virtual machines leverage physical GPUs installed locally in the ESXi hosts to provide hardware-accelerated 3D graphics to multiple virtual desktops | 512MB max video memory per VM, half from RAM, half from GPU. GPU resources are assigned first come, first serve. vSphere required. |
| **vDGA** (Virtual <u>Dedicated</u> Graphics Acceleration) | Graphics acceleration capability provided by VMware ESXi for high-end workstation graphics where a discrete GPU is needed | More suited to persistent desktop use case. Live vMotion is not supported (cold only). Intel VT-d required. vSphere required. |
| **vGPU** (Virtual <u>Shared</u> or <u>Dedicated</u> Graphics Acceleration) | True hardware graphics acceleration capability provided by Citrix and NVIDIA that enables flexibility and custom assignment of GPU resources | Up to 8 users can share a physical GPU, Windows 8 support currently experimental, XenServer required. |

For more information, please visit: LINK

## 4.2 Graphics Compute Server Infrastructure

### 4.2.1 Local Tier 1 Rack

In the Local Tier 1 model, VDI sessions execute on local storage on each compute server. In this model, the management server hosts access shared storage to support the VDI management VMs. Because of this, the compute and management servers are configured with different add-on NICs to support their pertinent network fabric connection requirements. Due to the reduced densities on GPU-enabled compute hosts, recommended RAM has been reduced. The management server host does not require local disk to host the management VMs.

| Local Tier 1 GFX Compute Host – PowerEdge R720 | | |
|---|---|---|
| 2 x Intel Xeon E5-2680v2 Processor (2.8Ghz) | | 2 x Intel Xeon E5-2680v2 Processor (2.8Ghz) |
| 96GB Memory (6 x 16GB DIMMs @ 1600Mhz) | | 96GB Memory (6 x 16GB DIMMs @ 1600Mhz) |
| VMware **vSphere** on internal 2GB Dual SD | | Citrix **XenServer** on 12 x 300GB 15K SAS disks |
| 10 x 300GB SAS 6Gbps 15k Disks (VDI) | | PERC H710 Integrated RAID Controller – RAID10 |
| PERC H710 Integrated RAID Controller – RAID10 | **OR** | Broadcom 5720  1Gb QP NDC (LAN) |
| Broadcom 5720  1Gb QP NDC (LAN) | | Broadcom 5720 1Gb DP NIC (LAN) |
| Broadcom 5720 1Gb DP NIC (LAN) | | iDRAC7 Enterprise w/ vFlash, 8GB SD |
| iDRAC7 Enterprise w/ vFlash, 8GB SD | | 2 x 1100W PSUs |
| 2 x 1100W PSUs | | |

| Local Tier 1 Management Host – PowerEdge R720 | | |
|---|---|---|
| 2 x Intel Xeon E5-2670v2 Processor (2.5Ghz) | | 2 x Intel Xeon E5-2670v2 Processor (2.5Ghz) |
| 96GB Memory (6 x 16GB DIMMs @ 1600Mhz) | | 96GB Memory (6 x 16GB DIMMs @ 1600Mhz) |
| VMware **vSphere** on internal 2GB Dual SD | | Citrix **XenServer** on 2 x 300GB 15K SAS disks |
| Broadcom 5720  1Gb QP NDC (LAN/iSCSI) | **OR** | Broadcom 5720  1Gb QP NDC (LAN/iSCSI) |
| Broadcom 5719 1Gb QP NIC (LAN/iSCSI) | | Broadcom 5719 1Gb QP NIC (LAN/iSCSI) |
| iDRAC7 Enterprise w/ vFlash, 8GB SD | | iDRAC7 Enterprise w/ vFlash, 8GB SD |
| 2 x 750W PSUs | | 2 x 750W PSUs |

### 4.2.2 Local Tier 1 Rack Workstation

The Dell Precision R7610 rack workstation is another compute host option for the Local Tier 1 solution that allows for up to three NVIDIA K2A cards to be installed per host. The R720 is also the recommended management host platform for environments implementing the R7610, as noted in the previous section.

| Local Tier 1 GFX Workstation – Precision R7610 |
|---|
| 2 x Intel Xeon E5-2697v2 Processor (2.7Ghz) |
| 128GB Memory (6 x 16GB DIMMs @ 1600Mhz) |
| Citrix **XenServer** on 4 x 300GB 15K SAS disks |
| LSI Mega RAID Controller – RAID10 |
| Up to 3 x NVIDIA Grid K2A boards |
| Intel 82579 1Gb DP LOM (LAN) |
| Intel I350 1Gb DP NIC (LAN) |
| 2 x 1400W PSUs |

### 4.2.3 Shared Tier 1 Rack (iSCSI)

In the Shared Tier 1 model, VDI sessions execute on shared storage so there is no need for local disk on each server. To provide server-level network redundancy using the fewest physical NICs possible, both the compute and management servers use a split QP NDC: 2 x 10Gb ports for iSCSI, 2 x 1Gb ports for LAN. 2 additional DP NICs (2 x 1Gb + 2 x 10Gb) provide slot and connection redundancy for both network fabrics. All configuration options otherwise are identical to the Local Tier 1 host.

| Shared Tier 1 GFX Compute Host (iSCSI) – PowerEdge R720 | | |
|---|---|---|
| 2 x Intel Xeon E5-2680v2 Processor (2.8GHz) | **OR** | 2 x Intel Xeon E5-2680v2 Processor (2.8Ghz) |
| 96GB Memory (6 x 16GB DIMMs @ 1600Mhz) | | 96GB Memory (6 x 16GB DIMMs @ 1600Mhz) |
| VMware **vSphere** on internal 2GB Dual SD | | Citrix **XenServer** on 2 x 300GB 15K SAS disks |
| Broadcom 57800  2 x 10Gb SFP+ + 2 x 1Gb NDC | | Broadcom 57800  2 x 10Gb SFP+ + 2 x 1Gb NDC |
| 1 x Broadcom 5720 1Gb DP NIC (LAN) | | 1 x Broadcom 5720 1Gb DP NIC (LAN) |
| 1 x Intel X520 2 x 10Gb SFP+ DP NIC (iSCSI) | | 1 x Intel X520 2 x 10Gb SFP+ DP NIC (iSCSI) |
| iDRAC7 Enterprise w/ vFlash, 8GB SD | | iDRAC7 Enterprise w/ vFlash, 8GB SD |
| 2 x 1100W PSUs | | 2 x 1100W PSUs |

| Shared Tier 1 Management Host (iSCSI) – PowerEdge R720 | | |
|---|---|---|
| 2 x Intel Xeon E5-2670v2 Processor (2.5Ghz) | **OR** | 2 x Intel Xeon E5-2670v2 Processor (2.5Ghz) |
| 96GB Memory (6 x 16GB DIMMs @ 1600Mhz) | | 96GB Memory (6 x 16GB DIMMs @ 1600Mhz) |
| VMware **vSphere** on internal 2GB Dual SD | | Citrix **XenServer** on 2 x 300GB 15K SAS disks |
| Broadcom 57800  2 x 10Gb SFP+ + 2 x 1Gb NDC | | Broadcom 57800  2 x 10Gb SFP+ + 2 x 1Gb NDC |
| 1 x Broadcom 5720 1Gb DP NIC (LAN) | | 1 x Broadcom 5720 1Gb DP NIC (LAN) |
| 1 x Intel X520 2 x 10Gb SFP+ DP NIC (iSCSI) | | 1 x Intel X520 2 x 10Gb SFP+ DP NIC (iSCSI) |
| iDRAC7 Enterprise w/ vFlash, 8GB SD | | iDRAC7 Enterprise w/ vFlash, 8GB SD |
| 2 x 750W PSUs | | 2 x 750W PSUs |

### 4.2.4 Shared Tier 1 Rack (FC)

Fiber Channel is an optional block storage protocol for Compute and Management hosts with Compellent Tier 1 and Tier 2 storage. Aside from the use of FC HBAs to replace the 10Gb NICs used for iSCSI, the rest of the server configurations are the same.

| Shared Tier 1 GFX Compute Host (FC) – PowerEdge R720 | | |
|---|---|---|
| 2 x Intel Xeon E5-2680v2 Processor (2.8GHz) | **OR** | 2 x Intel Xeon E5-2680v2 Processor (2.8Ghz) |
| 96GB Memory (6 x 16GB DIMMs @ 1600Mhz) | | 96GB Memory (6 x 16GB DIMMs @ 1600Mhz) |
| VMware **vSphere** on internal 2GB Dual SD | | Citrix **XenServer** on 2 x 300GB 15K SAS disks |
| Broadcom 5720  1Gb QP NDC (LAN) | | Broadcom 5720  1Gb QP NDC (LAN) |
| Broadcom 5720 1Gb DP NIC (LAN) | | Broadcom 5720 1Gb DP NIC (LAN) |
| 2 x QLogic 2562 8Gb DP FC HBA | | 2 x QLogic 2562 8Gb DP FC HBA |
| iDRAC7 Enterprise w/ vFlash, 8GB SD | | iDRAC7 Enterprise w/ vFlash, 8GB SD |
| 2 x 1100W PSUs | | 2 x 1100W PSUs |

| Shared Tier 1 Management Host (FC) – PowerEdge R720 | | |
|---|:---:|---|
| 2 x Intel Xeon E5-2670v2 Processor (2.5Ghz) | | 2 x Intel Xeon E5-2670v2 Processor (2.5Ghz) |
| 96GB Memory (6 x 16GB DIMMs @ 1600Mhz) | | 96GB Memory (6 x 16GB DIMMs @ 1600Mhz) |
| VMware **vSphere** on internal 2GB Dual SD | | Citrix **XenServer** on 2 x 300GB 15K SAS disks |
| Broadcom 5720  1Gb QP NDC (LAN) | **OR** | Broadcom 5720  1Gb QP NDC (LAN) |
| Broadcom 5720 1Gb DP NIC (LAN) | | Broadcom 5720 1Gb DP NIC (LAN) |
| 2 x QLogic 2562 8Gb DP FC HBA | | 2 x QLogic 2562 8Gb DP FC HBA |
| iDRAC7 Enterprise w/ vFlash, 8GB SD | | iDRAC7 Enterprise w/ vFlash, 8GB SD |
| 2 x 750W PSUs | | 2 x 750W PSUs |

In the above configurations, the R720-based Dell Wyse Datacenter Solution can support the following user counts per server based on the configurations specified.

| Local/ Shared Tier 1 Rack Densities | | |
|---|:---:|:---:|
| Mode | 2 x NVidia  K1 | 2 x NVidia K2 |
| vSGA (Shared) | 18 | 20 |
| vDGA (Pass-Through) | 8 | 4 |
| vGPU (R720) | 64<br>(K100Q) | 8<br>(K260Q) |
| vGPU Pass-Through (R7610) | - | 6 |
| vGPU (R7610) | - | 8 |

# 5 Solution Performance and Testing

## 5.1 Load Generation and Monitoring

### 5.1.1 Login VSI – Login Consultants

Login VSI is the de-facto industry standard tool for testing VDI environments and server-based computing / terminal services environments. It installs a standard collection of desktop application software (e.g. Microsoft Office, Adobe Acrobat Reader etc.) on each VDI desktop; it then uses launcher systems to connect a specified number of users to available desktops within the environment.  Once the user is connected the workload is started via a logon script which starts the test script once the user environment is configured by the login script. Each launcher system can launch connections to a number of 'target' machines (i.e. VDI desktops), with the launchers being managed by a centralized management console, which is used to configure and manage the Login VSI environment.

### 5.1.2 Liquidware Labs Stratusphere UX

Stratusphere UX was used during each test run to gather data relating to User Experience and desktop performance.  Data was gathered at the Host and Virtual Machine layers and reported back to a central server (Stratusphere Hub).  The hub was then used to create a series of "Comma Separated Values" (.csv) reports which have then been used to generate graphs and summary tables of key information.  In addition the Stratusphere Hub generates a magic quadrate style scatter plot showing the Machine and IO experience of the sessions.  The Stratusphere hub was deployed onto the core network therefore its monitoring did not impact the servers being tested. This core network represents an existing customer environment and also includes the following services:

- Active Directory
- DNS
- DHCP
- Anti-Virus

Stratusphere UX calculates the User Experience by monitoring key metrics within the Virtual Desktop environment, the metrics and their thresholds are shown in the following screen shot:

**Machine Experience Indicators**

| | Weight (%) | Good | | | Fair | | | Poor | |
|---|---|---|---|---|---|---|---|---|---|
| **Login Delay :** Time it takes to login (sec.) ? | 20 | 0 | <= | 15 | <= | 60 | <= | unbounded |
| **Application Load Time :** Avg. startup time for applications (sec.) ? | 20 | 0 | <= | 10 | <= | 30 | <= | unbounded |
| **CPU Queue Length :** Length of CPU queue at inspection time ? | 20 | 0 | <= | 3 | <= | 6 | <= | unbounded |
| **Page Faults :** Number of page faults during inspection interval ? | 20 | 0 | <= | 2,000 | <= | 10,000 | <= | unbounded |
| **Non-Responding Applications :** Number of unresponsive applications at inspection time ? | 20 | 0 | <= | 2 | <= | 3 | <= | unbounded |

**I/O Experience Indicators**

| | Weight (%) | Good | | | Fair | | | Poor | |
|---|---|---|---|---|---|---|---|---|---|
| **Disk Load :** Avg. disk IO per second ? | 25 | 0 | <= | 25 | <= | 75 | <= | unbounded |
| **Disk Queue Length :** Avg. length of disk queue(s) ? | 25 | 0 | <= | 1 | <= | 3 | <= | unbounded |
| **Network Latency :** Avg. network roundtrip time (ms) ? | 25 | 0 | <= | 150 | <= | 300 | <= | unbounded |
| **Failed Connections :** Number of outgoing connection attempts that failed ? | 25 | 0 | <= | 5 | <= | 15 | <= | unbounded |

### 5.1.3 VMware vCenter

VMware vCenter has been used for VMware vSphere-based solutions to gather key data (CPU, Memory and Network usage) from each of the desktop hosts during each test run. This data was exported to .csv files for each host and then consolidated to show data from all hosts. While the report does not include specific performance metrics for the Management host servers, these servers were monitored during testing and were seen to be performing at an expected performance level.

## 5.2 Testing and Validation

### 5.2.1 Testing Process

The purpose of the single server testing is to validate the architectural assumptions made around the server stack. Each user load is tested against 4 runs. A pilot run to validate that the infrastructure is functioning and valid data can be captured and 3 subsequent runs allowing correlation of data. Summary of the test results will be listed out in the below mentioned tabular format.

At different stages of the testing the testing team will complete some manual "User Experience" Testing while the environment is under load. This will involve a team member logging into a session during the run and completing tasks similar to the User Workload description. While this experience will be subjective, it will help provide a better understanding of the end user experience of the desktop sessions, particularly under high load, and ensure that the data gathered is reliable.

Login VSI has two modes for launching user's sessions:

- Parallel
  - Sessions are launched from multiple launcher hosts in a round robin fashion; this mode is recommended by Login Consultants when running tests against multiple host servers. In parallel mode the VSI console is configured to launch a number of sessions over a specified time period (specified in seconds)
- Sequential
  - Sessions are launched from each launcher host in sequence, sessions are only started from a second host once all sessions have been launched on the first host and this is repeated for each launcher host. Sequential launching is recommended by Login Consultants when testing a single desktop host server. The VSI console is configure to launch a specified number of session at a specified interval specified in seconds

All test runs which involved the 6 desktop hosts were conducted using the Login VSI "Parallel Launch" mode, all sessions were launched over an hour to try and represent the typical 9am logon storm. Once the last user session has connected, the sessions are left to run for 15 minutes prior to the sessions being instructed to logout at the end of the current task sequence, this allows every user to complete a minimum of two task sequences within the run before logging out. The single server test runs were configured to launch user sessions every 60 seconds, as with the full bundle test runs sessions were left to run for 15 minutes after the last user connected prior to the sessions being instructed to log out.

### 5.2.2 Dell Wyse Datacenter Profiles

The table shown below presents the profiles used during PAAC on Dell Wyse Datacenter solutions. These profiles have been carefully selected to provide the optimal level of resources for common use cases.

| Profile Name | Number of vCPUs per Virtual Desktop | Nominal RAM (GB) per Virtual Desktop | Use Case |
|---|---|---|---|
| Standard | 1 | 2 | Task Worker |
| Enhanced | 2 | 3 | Knowledge Worker |
| Professional | 2 | 4 | Power User |
| Shared Graphics | 2 + Shared GPU | 3 | Knowledge Worker with high graphics consumption requirements |
| Pass-through Graphics | 4 + Pass-through GPU | 32 | Workstation type user e.g. producing complex 3D models. |

### 5.2.3 Dell Wyse Datacenter Workloads

Load-testing on each of the profiles described in the above table is carried out using an appropriate workload that is representative of the relevant use case. In the case of the non-graphics workloads, these workloads are Login VSI workloads and in the case of graphics workloads, these are specially designed workloads that stress the VDI environment to a level that is appropriate for the relevant use case. This information is summarized in the table below.

| Profile Name | Workload | OS Image |
|---|---|---|
| Standard | Login VSI Light | Shared |
| Enhanced | Login VSI Medium | Shared |
| Professional | Login VSI Heavy | Shared + Profile Virtualization |
| Shared Graphics | Fishbowl / eFigures | Shared + Profile Virtualization |
| Pass-through Graphics | eFigures / AutoCAD - SPEC Viewperf | Persistent |

As noted in the table above, further information for each of the workloads is given below. It is noted that for Login VSI testing, the following login and boot paradigm is used:

- For single-server / single-host testing (typically carried out to determine the virtual desktop capacity of a specific physical server), users are logged in every 30 seconds.

- For multi-host / full solution testing, users are logged in over a period of 1-hour, to replicate the normal login storm in an enterprise environment.

- All desktops are pre-booted in advance of logins commencing.

For all testing, all virtual desktops run an industry-standard anti-virus solution (currently McAfee VirusScan Enterprise) in order to fully replicate a customer environment.

### 5.2.3.1 Workloads Running on Shared Graphics Profile

Graphics hardware vendors (e.g. Nvidia) typically market a number of graphics cards that are targeted at different markets. Consequently, it is necessary to provide two shared graphics workloads – one for mid-range cards and the other for high-end cards. These workloads are described in more detail below. It is noted that technologies such as the Citrix / Nvidia vGPU technology (where the Nvidia drivers reside on the virtual desktop, giving shared-level density with near pass-through functionality) mean that in some cases, the higher-end workloads, traditionally used for pass-through GPU PAAC, may be more appropriate for shared GPU PAAC. Such scenarios will explicitly state the workload used.

**Mid-Range Shared Graphics Workload**

The mid-range shared graphics workload is a modified Login VSI medium workload with 60 seconds of graphics-intensive activity (Microsoft Fishbowl at http://ie.microsoft.com/testdrive/performance/fishbowl/) added to each loop.

**High-End Shared Graphics Workload**

The high-end shared graphics workload consists of one desktop running Heaven Benchmark and n-1 desktops running eFigures Advanced Animation activity where n = per-host virtual desktop density being tested at any specific time.

### 5.2.3.2 Workloads Running on Pass-through Graphics Profile

Similarly for pass-through graphics, two workloads have been defined in order to align with graphics cards of differing capabilities.

**Mid-Range Pass-through Graphics Workload**

The mid-range pass-through graphics workload consists of one desktop running Heaven Benchmark and n-1 desktops running eFigures Advanced Animation activity where n = per-host virtual desktop density being tested at any specific time.

**High-End Pass-through Graphics Workload**

One desktop running Viewperf benchmark; n-1 desktops running AutoCAD auto-rotate activity where n = per host virtual desktop density being tested at any specific time.

## 5.2.4 vDGA Results

These results describe validation efforts undertaken on vSphere using NVidia K1 and K2 Grid cards to study their behavior when used with graphic-intensive applications. Broadly speaking, users of graphics-intensive applications in a virtual desktop environment can be subdivided into 2 categories, as discussed below:

- "Premium Plus" VDI users are users who may be consuming relatively high-end graphics through relatively high frame-rate applications such as Google Earth, graphics-rich HTML5 pages etc. and also reviewing electrical, mechanical CAD drawings etc. *[The term 'Premium Plus' is used to distinguish this user type from the existing 'Premium' user type that is used in current Dell Wyse Solutions Engineering-Wyse PAAC and Sizing Activities].*

- Workstation users, as the name implies, are users who would typically have used high-end physical workstations (e.g. Dell Precision); typical activities carried out by these users would include 3D modelling for the oil and gas industry, involving a large amount of resource –intensive activities such as model rotation etc.

The validation effort described in this document is for Workstation users. In the cases described here the Grid cards are operating in pass-through mode, i.e. a VM has direct access to a GPU on the K1 or K2 card. The GPU cannot be shared with other VMs on the Hypervisor. No other VMs are provisioned on the server except those involved in the Workstation user testing. It is assumed that the server will be dedicated to the number of workstation users that can be supported by the appropriate Grid card in pass-through mode.

All of the following results were gathered with either two K1 cards or two K2 cards installed in the server. The K1 card allows 4 GPU pass-through sessions per card (**8 per Server**), while the K2 card allows 2 pass-through sessions per card (**4 per server**).

## Graphics-Specific Performance Analysis Results

**Viewperf Benchmark**

SPECviewperf is a widely used benchmark in the workstation domain for benchmarking graphics performance and on this basis it has been chosen as the appropriate benchmark to use for the pass-through use case. It was used for both K1 and K2 testing. The tests run in SPECviewperf are defined by viewsets.  The viewsets chosen for the K1 and K2 pass-through testing are CATIA, EnSight, Maya, Pro/ENGINEER, SolidWorks, Teamcenter Visualization Mockup, Siemens NX and Lightwave. For a description of these viewsets please see http://www.spec.org/gwpg/gpc.static/vp11info.html.
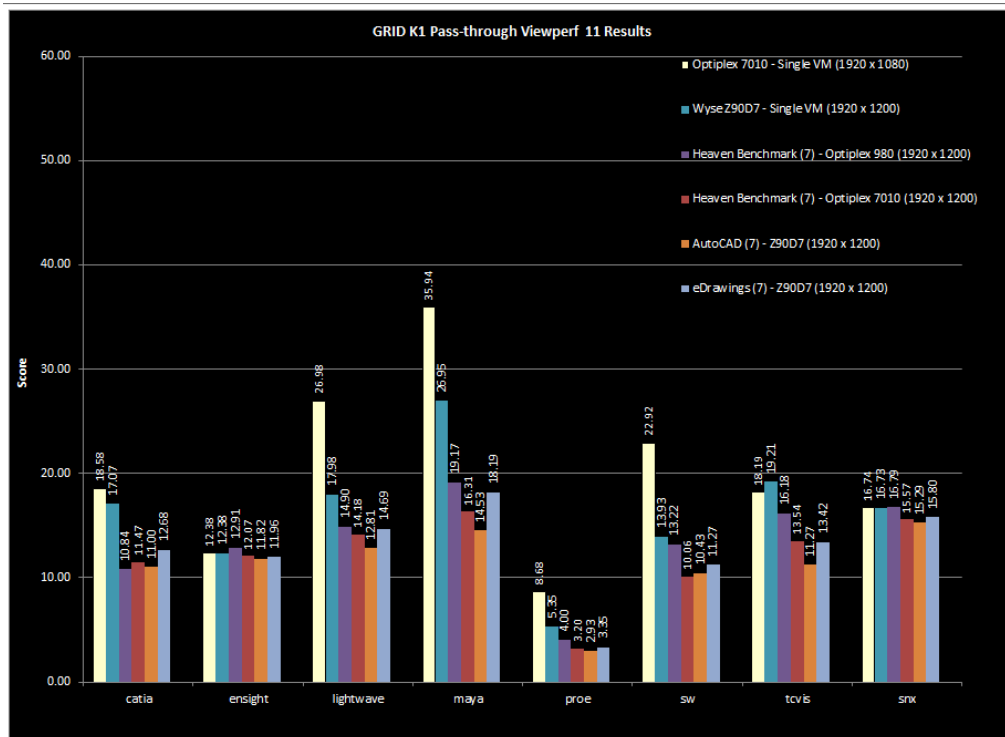
The SPECviewperf workloads were run against both the Dell OptiPlex 7010 and the Dell Wyse Z90D7.

A number of tests were run against K1 grid card as follows
- Viewperf test run on a single VM against both endpoints with no companion tests running
- Viewperf test run against 2 OptiPlex endpoints with Heaven companion workload on remaining VMs
- Viewperf test run against Wyse Z90D7 with Solidworks eDrawings companion workload on remaining VMs
- Viewperf test run against Wyse Z90D7 endpoint with AutoCAD companion workload on remaining VMs

The SPECviewperf results recorded are presented in the following graph:



**GRID K1 Pass-through Viewperf 11 Results**

A number of tests were run against K2 grid card as follows:

- Viewperf test run on a single VM against both endpoints with no companion tests running. An additional test was run against OptiPlex 980

- Viewperf test run against OptiPlex 7010 endpoint with Heaven companion workload on remaining VMs

- Viewperf test run against Wyse Z90D7 and OptiPlex 7010 endpoints with Solidworks eDrawings companion workload on remaining VMs

- Viewperf test run against Wyse Z90D7 and OptiPlex 7010 endpoints with AutoCAD companion workload on remaining VMs

The SPECviewperf results recorded are presented in the following graph:



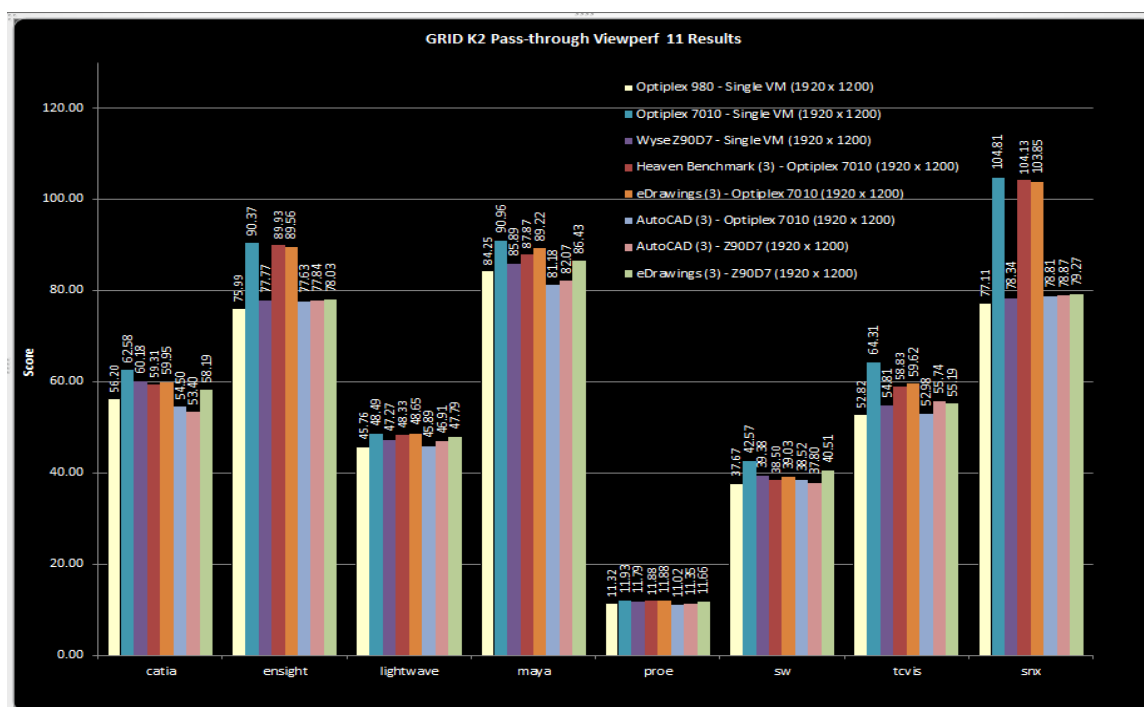SPECviewperf reports results in frames per second. The most obvious result is that the K2 card gives a consistently higher score than the K1 card. This is to be expected as the NVidia GRID K2 card is positioned a high-end graphics card with performance similar to the Quadro K5000 card while the NVidia K1 is positioned as a mid-range graphics card equivalent to a Quadro K600.

It is also noticeable that the addition workload when companion tests are running affects the results for the K1 more than the K2.

For the K2 card these results appear in line with some of the results obtained on workstations submitted to the SPECviewperf website results. Please see a summary of submitted results for various workloads across multiple workstations (including 2 Dell models) at http://www.spec.org/gwpg/gpc.data/vp11/summary.html.

**Results for Dell Precision Workstation R5500 with NVIDIA Quadro 6000:**

| Company/Product | Catia | Ensight | Lightwave | Maya | ProE | Snx | SW | Tcvis |
|---|---|---|---|---|---|---|---|---|
| Dell Precision Workstation R5500 with NVIDIA Quadro 6000 | 50.68 | 58.23 | 59.33 | 108.77 | 9.95 | 60.14 | 58.27 | 47.8 |

These results were submitted by Dell

Improved results may be expected against workstations with higher spec. processors using the K5000 video card. It appears that the K2 GRID card is capable of competing against high end workstations but the K1 GRID card should be reserved for medium use cases.

In addition to the Viewperf metrics recorded the Frame rate per second was recorded from Citrix XenDesktop as described below.

*It is first necessary to enable debug logging using the HDX3D cmd line tool using the command shown below:*

*"C:\Program Files\Citrix\ICAService\HDX3DConfigCmdLineX64.exe" debug_logging 1"*

Next, start DebugView – a copy comes with the HDX 3d Pro health check tool and this should be used.  Once this tool has been started, it is possible to filter on "FPS" and save this frame-rate information to a log file.

An interesting result here is that on the K1 tests the companion workload maintained a reasonably high frame rate even though the Viewperf frame rate was low.
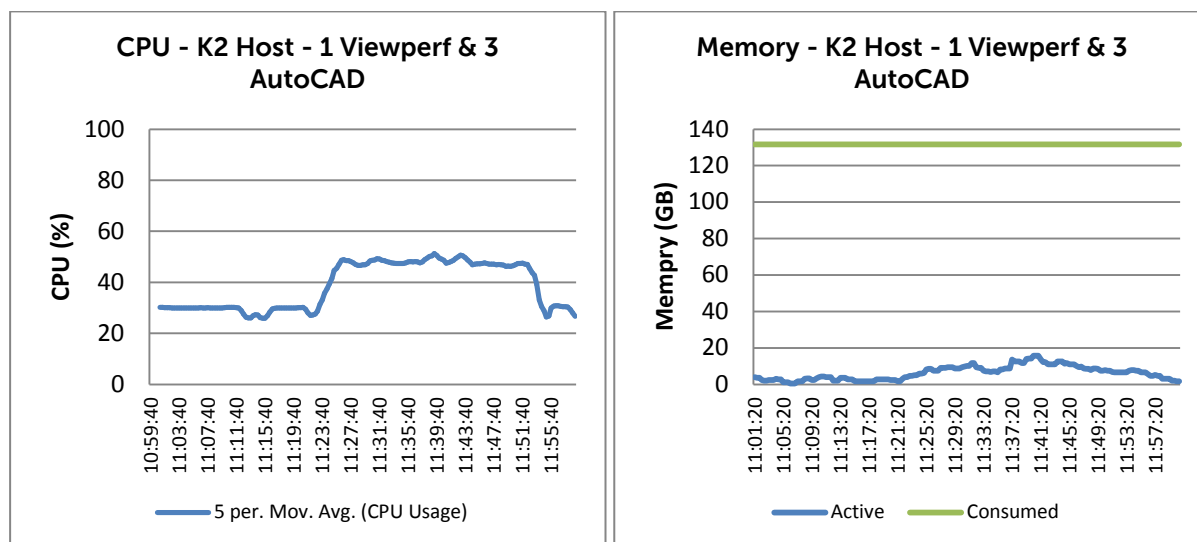
In addition to specific graphics results obtained results were gathered for the host performance during various workloads. This was done to ensure that the host was not stressed and the recorded results fall within the thresholds below.
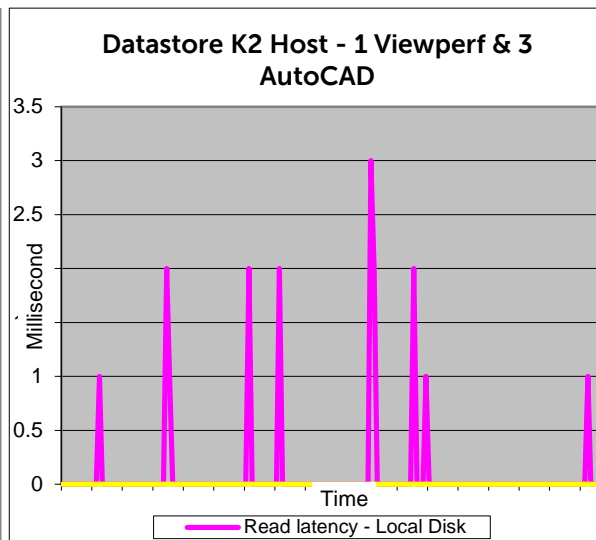
| Resource Utilization Parameter | Threshold |
|---|---|
| Compute Host CPU | 85% |
| Computer Host Memory | 85% |
| Network Throughput | 85% |
| Tier 1 Storage Latency | 20ms |

Results were gathered using the VSphere Performance Reporting tool.  The results shown below were gathered when running SPECviewperf tests on the K2 card. These represent the most load seen on the R720 host during the validation.

**AUTOCAD & Viewperf Workload**

The following graphs are gathered when the system is running the SPECviewperf test in conjunction with 3 companion AUTCAD tests.

**Network - K2 Host - 1 Viewperf & 3 AutoCAD**



**Datastore K2 Host - 1 Viewperf & 3 AutoCAD**

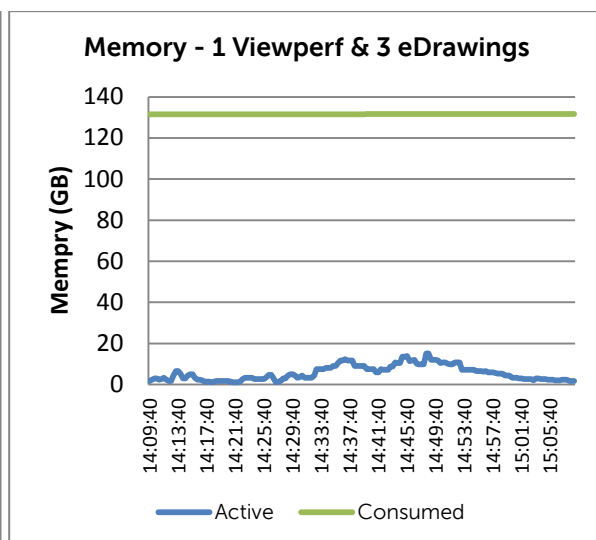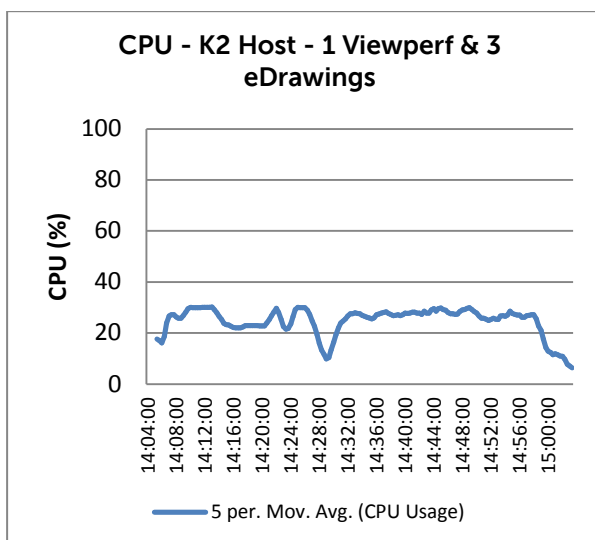## eDrawings & Viewperf Workload

The following graphs are gathered when the system is running the SPECviewperf test in conjunction with 3 companion eDrawings workloads.



**CPU - K2 Host - 1 Viewperf & 3 eDrawings**



**Memory - 1 Viewperf & 3 eDrawings**

**Network - 1 Viewperf & 3 eDrawings**

*(chart: KBps vs. time from 14:11:00 to 15:07:00, legend: 5 per. Mov. Avg. (Total Usage))*

**Heaven Benchmark – No Viewperf**

In the following case the Heaven workload was running on all 4 VMs:

**CPU - K2 Host - Heaven Benchmark**

*(chart: CPU (%) vs. time from 15:27:00 to 16:23:00, legend: 5 per. Mov. Avg. (CPU Usage))*

**Memory - K2 Host - Heaven Benchmark**

*(chart: Mempry (GB) vs. time from 15:28:40 to 16:24:40, legend: Active, Consumed)*

**Network - K2 Host - Heaven Benchmark**

It can be seen that in all cases the performance of the host remains well within the thresholds defined. Maximum CPU is approximately 50% and network usage does not go above 4 Mbps. Active and consumed memory are both well within the capabilities of the host and no ballooning is evident.

Although only one set of results are included for the Datastore this shows that the local disk latency did not go above 1 millisecond.  Similar results are observed in the raw data for other tests.

**Conclusion**

The results from the pass-through tests indicate that the R720 host is capable of handling the various graphics workloads. It should be noted that the specification for the processor is reduced from the normal 2.9GHz Dell Wyse Solutions Engineering spec. to 2.6 GHz because of power/thermal limitations.

The K2 card earns its place as the superior offering over the K1, as expected, and it seems to match the performance of many of the workstation results on the SPECviewperf website. However, there are probably newer higher specification workstations that show superior performance. The K1 cards perform well for the mid-size market segment it is positioned for.

### 5.2.5  vSGA Results

This section describes validation efforts undertaken on vSphere using NVIDIA Grid K1 cards to study their behavior when used in shared mode for virtual desktops with graphics-intensive applications.

Broadly speaking, users of graphics-intensive applications in a virtual desktop environment can be subdivided into 2 categories, as discussed below:

- "Premium Plus" VDI users are users who may be consuming relatively high-end graphics through relatively high frame-rate applications such as Google Earth, graphics-rich HTML5 pages etc. and also reviewing electrical, mechanical CAD drawings etc. *[The term 'Premium Plus' is used to distinguish this user type from the existing 'Premium' user type that is used in current Dell Wyse Solutions Engineering-Wyse PAAC and Sizing Activities].*

- Workstation users, as the name implies, are users who would typically have used high-end physical workstations (e.g. Dell Precision); typical activities carried out by these users

would include 3D modelling for the oil and gas industry, involving a large amount of resource –intensive activities such as model rotation etc.
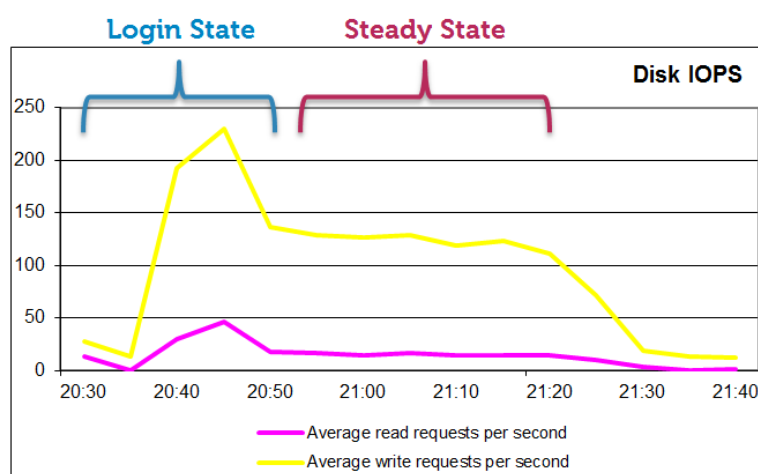
The validation effort described in this document is for Workstation users. In the cases described here the Grid cards are operating in shard mode (vSGA). VMware vSGA provides the ability for multiple virtual machines to leverage physical GPUs installed locally in the ESXi hosts to provide hardware-accelerated 3D graphics. i.e. a VM does not have direct access to a GPU on the K1 cards but multiple VM share a GPU through the hypervisor. No other VMs are provisioned on the server except those involved in the Workstation user testing. It is assumed that the server will be dedicated to the number of workstation users that can be supported by the K1 Grid cards in shared mode. All of the following results were gathered with two K1 cards installed in the server.

## Test Results Summary

This validation was performed for XenDesktop 7 virtual machines running on an R720 host, 96GB of RAM and dual 2.5 GHz processors. Validation was performed using Dell Wyse Solutions Engineering standard testing methodology using LoginVSI load generation tool for VDI benchmarking that simulates production user workloads. The "multimedia workload" profile in LoginVSI was used for testing the graphics capabilities. In addition, the Fishbowl HTML5 test was also used. As a result of this testing, the optimal supported configuration proved to be:

- Maximum of **18** XenDesktop 7 virtual machines per physical R720 host
- Each XenDesktop 7 virtual machine configured with four 2x vCPUs
- Each XenDesktop 7 virtual machine configured with 4GB RAM

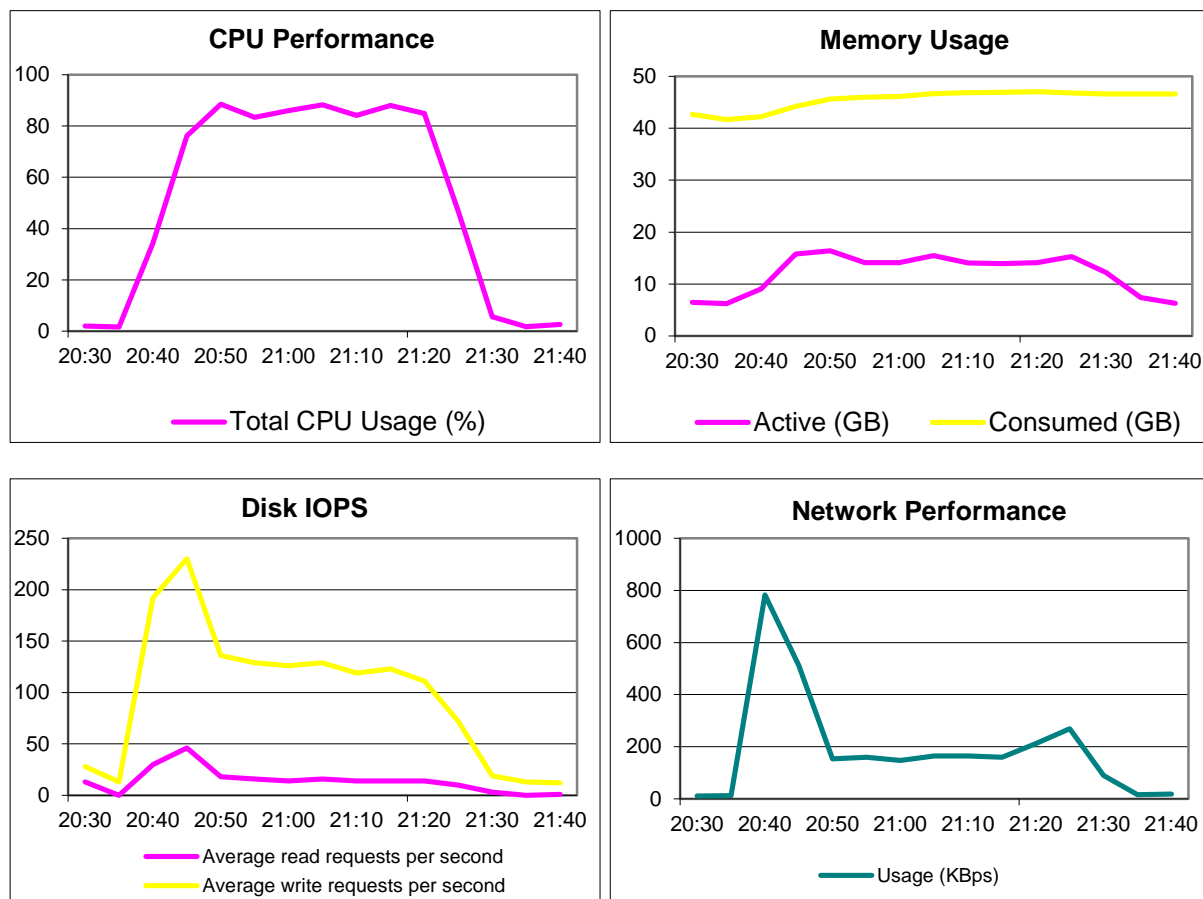| Workload | Server Density | CPU | Login State IOPS | Login State IOPS per User | Steady State IOPS | Steady State IOPS per User |
|---|---|---|---|---|---|---|
| Multimedia Workload & Fishbowl HTML5 Test | 18 | 83.3% - 88.5% | 276 | 15.3 | 145 | 8 |



This virtual server-based configuration was able to sustain a user density of 18 sessions per physical host server depending on the type of user workload chosen as well as the file access and scanning configuration of the anti-virus tools.

## Server Resource Performance (18 Users)

The graphs below show the CPU, Consumed Memory, Local disk IOPS, Disk IOPs, Network and VDI UX scatter plot results from this validation. The CPU usage for this test reached 88.5% thus confirming that the server with this configuration can support up to 18 Standard users. Memory was configured 4GB per VM. The following graphs show CPU, memory, local disk IOPS, network and VDI UX scatter plot results.



Stratusphere UX indicates all desktops had a good user experience. All desktops were in the upper right corner of the upper right quadrant indicating that 18 standard provisioned users can be supported with good performance.

### 5.2.6 vGPU Results

#### 5.2.6.1 Dell PowerEdge R720

These results describe validation efforts undertaken on XenServer using NVidia K1 and K2 Grid cards to study their behavior when used with graphic-intensive applications. NVIDIA Grid K1 and K2 cards contain multiple GPUs on board. Each physical GPU (pGPU) can host several different types of virtual GPUs (vGPU). These vGPU types are preset profiles with a fixed amount of frame-buffer, fixed number of supported display heads and maximum resolutions, and fixed number of CUDA cores designed for different classes of workload. Each vGPU profile is suitable for different use case types. This validation effort concentrates on K140Q vGPU profile on the Grid K1 card, and the K260Q vGPU profile on the K2 card.

To validate these two vGPU profiles, the following two VM user profiles have been used.

- "Professional Shared Graphics" users are those who may be consuming relatively high-end graphics through relatively high frame-rate applications such as editing electrical, mechanical CAD drawings.
    - **vGPU Profile:** GRID K260Q
    - **VM Profile:** 2 vCPU, 4GB RAM with Windows 7 SP1 x64

- "Shared Graphics" VDI users are those who may be performing moderately intensive graphics work, e.g. viewing 3D drawings.
    - **vGPU Profile:** GRID K140Q
    - **VM Profile:** 2 vCPU, 3GB RAM with Windows 7 SP1 x64

The validation of these use cases is performed with 2x Grid K1 cards, and also with 2x Grid K2 cards. There is a finite limit on the number of vGPUs that can be created on a given card. This validation is performed with the maximum number of supported vGPUs for both, K140Q and K260Q, vGPU profiles.

| NVIDIA Grid Card | NVIDIA vGPU Profile Under Test | No. of Cards Per Host Server | No. of Physical GPUs Per Card | No. of vGPUs per pGPU | No. of vGPUs per card | No. of vGPUs per Server |
|---|---|---|---|---|---|---|
| K1 | K140Q | 2 | 4 | 4 | 16 | 32 |
| K2 | K260Q | 2 | 2 | 2 | 4 | 8 |

To test with the maximum number of supported vGPUs, many VMs are created and each VM is assigned a vGPU. No other VMs are provisioned on the server except those involved in this graphics testing.

## NVIDIA K140Q vGPU Profile

To validate the K140Q vGPU profile, two Grid K1 cards are used to support total of 32 VMs. The workload is simulated using eDrawings Viewer and a fixed frames per second (30 FPS) video. eDrawings Viewer is run on 31 out of the 32 VMs, with the video run on the 32nd.

**Solidworks eDrawings Viewer**
eDrawings Viewer is used to view 3D models. The goal of this activity is to simulate viewing 3D drawings. To simulate this activity, a 3D model is "played" continuously on all 31 out of the 32 VMs.

The continuous playing of the 3D model rotates the 3D object and zooms in and out periodically.
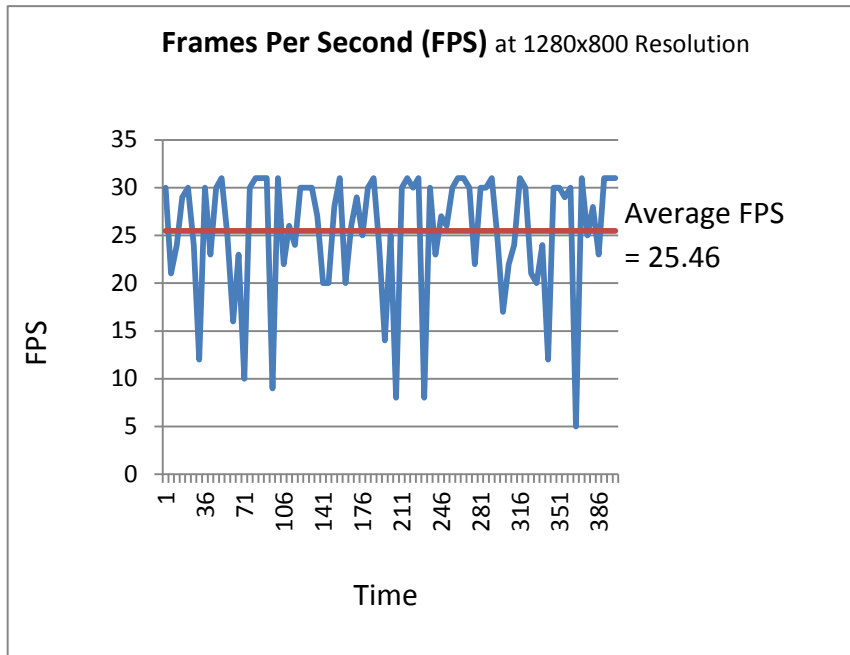
**Fixed FPS Video**
While considering the end-user experience in terms of perceived video/graphics smoothness, a useful domain to use for assessment is broadcast video. NTSC (US analog television system) is transmitted at 30 FPS, while PAL (widely used in Europe) is transmitted at 25 FPS. To simulate the user experience, a movie clip at a fixed 30 FPS has been created and used for this test.

**End-User Experience & Resolution**
While running the workloads as mentioned above, we paid close attention to the end-user experience with a high-performance **Dell Wyse Z90Q7 client**. One of the observations with the K140Q vGPU profile was that the user experience degraded at the higher resolutions, even when the GPU resources were not completely utilized. We observed that the best user experience for all 32 VMs running the workload as mentioned above occurs at resolution 1280x800 or lower. Considering this, all of the testing on K140Q vGPU profile was performed at resolution 1280x800. It should be noted that if the number of VMs are reduced or if a different set of workloads are used, the optimal resolution may be higher.
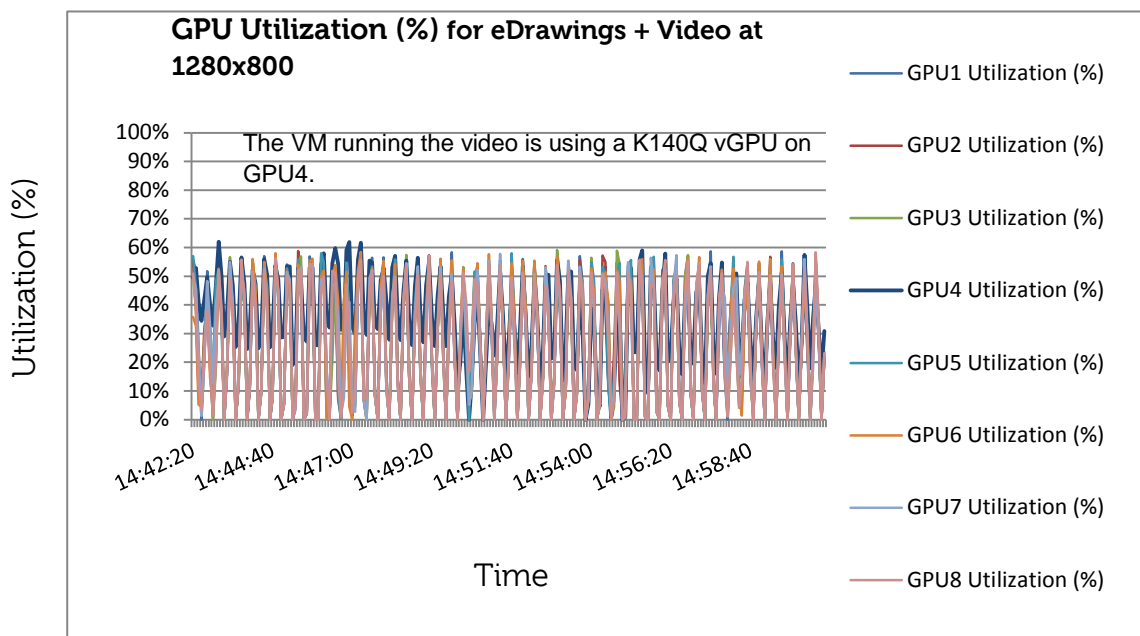
User experience for eDrawings Viewer: The user experience was good. The visual clarity and the colors were great. The rotation and the zooming in/out of the 3D objects were fluid and smooth. As a part of the effort to understand the user experience, we also measured the FPS at which the video is rendered on the end-point (Dell Wyse Z90Q7 client). This is the frames per second that HDX 3D Pro would allow to be rendered at the end-device. The FPS is measured using HDX Monitor provided by Citrix. As shown below, the video is rendered with an average of 25.46 FPS, which is considered to be good.
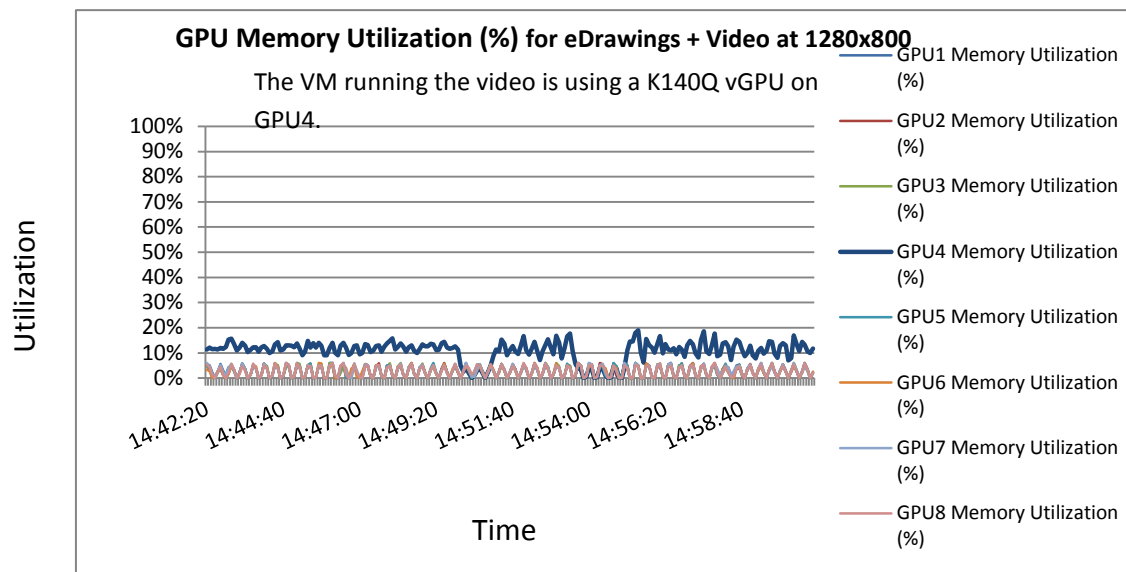
**Frames Per Second (FPS)** at 1280x800 Resolution



Average FPS = 25.46

## GPU Utilization

While the eDrawings viewer and the video were running, the physical GPU utilization and the GPU video memory utilization was continuously monitored. The goal was to ensure that the GPU cores were not saturated. As you can see, the GPU utilization is below 60% for all GPUs and the video memory utilization is below 20% for all GPUs. This represents the most load seen on the R720 host during the validation for K140Q vGPU profile.

**GPU Utilization (%) for eDrawings + Video at 1280x800**

The VM running the video is using a K140Q vGPU on GPU4.



- GPU1 Utilization (%)
- GPU2 Utilization (%)
- GPU3 Utilization (%)
- GPU4 Utilization (%)
- GPU5 Utilization (%)
- GPU6 Utilization (%)
- GPU7 Utilization (%)
- GPU8 Utilization (%)

**GPU Memory Utilization (%) for eDrawings + Video at 1280x800**
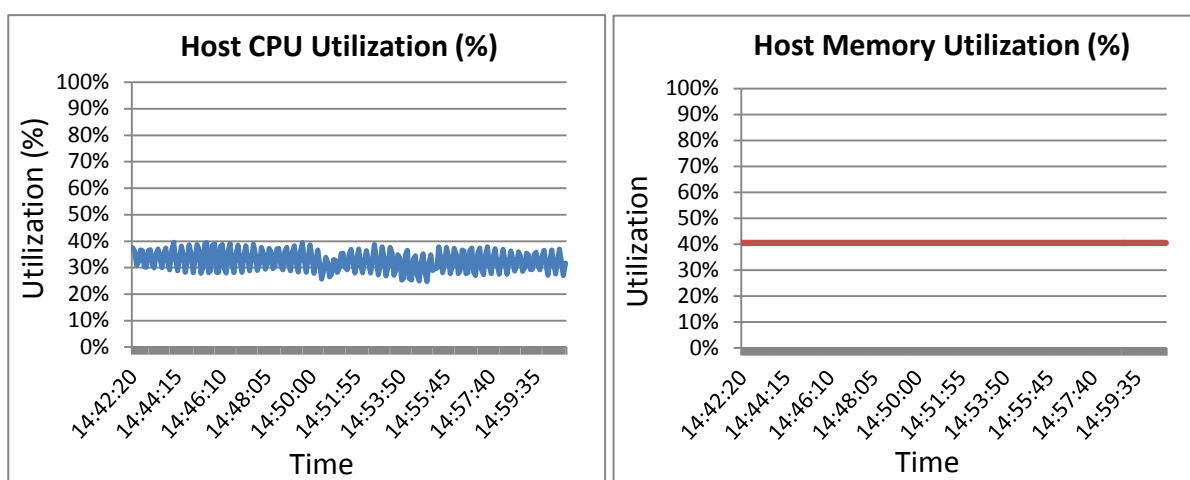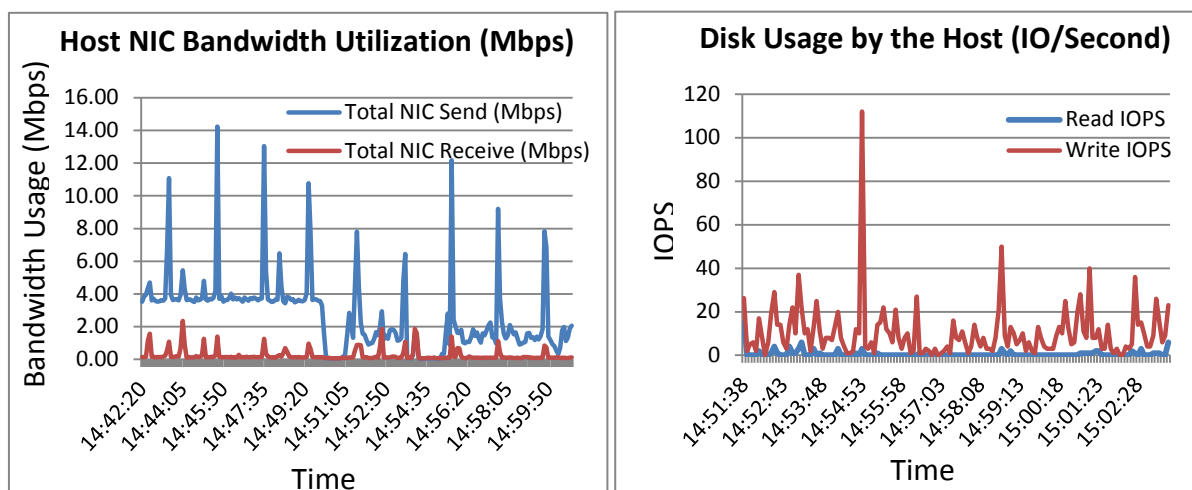
The VM running the video is using a K140Q vGPU on GPU4.

## Host Resource Utilization

During the activities described above, results were gathered for the host performance during various workloads. This was done to ensure that the host was not stressed and the recorded results fall within the thresholds below:

| NVIDIA Grid Card | NVIDIA vGPU Profile Under Test |
|---|---|
| Compute Host CPU | 85% |
| Compute Host Memory | 85% |
| Network Throughput | 85% |

Results were gathered using scripts running on the XenServer host server. The results shown below were gathered when running eDrawings Viewer on 31 VMs and the fixed frame-rate video on one VM. These represent the most load seen on the R720 host during the validation for K140Q vGPU profile.



**Host CPU Utilization (%)**



**Host Memory Utilization (%)**

**Host NIC Bandwidth Utilization (Mbps)**

**Disk Usage by the Host (IO/Second)**

As evident from the plots above, the utilization is well below the desired threshold. The table below shows the peak for each of these plots.

| NVIDIA Grid Card | NVIDIA vGPU Profile Under Test |
|---|---|
| Compute Host CPU | 40% |
| Compute Host Memory | 40% |
| Network Throughput | 1.4% |

**Additional Behavioral Observations**

To further gauge the user experience and the capability of the vGPU profile, user experience was observed while running Auto CAD 2014 on one of the VMs (and eDrawings Viewer on the remaining 31). The user experience for AutoCAD, while running AutoCAD 2014 on one VM and eDrawings Viewer on the other 31 VMs, was as good as the user experience for eDrawings Viewer.

## NVIDIA K260Q vGPU Profile

To validate the K260Q vGPU profile, two Grid K2 cards are used to support total of eight VMs. The workload is simulated using AutoCAD 2014, and SPECviewperf 11 is used to measure performance. AutoCAD 2014 is run on seven out of the eight VMs with SPECviewperf running on the 8[th] VM. AutoCAD running on seven VMs creates load on the physical GPU, creating load situations for the performance measurement using SPECviewperf.

**AutoCAD 2014**
In this validation activity, AutoCAD 2014 is used to view 3D CAD drawings. This viewing activity is simulated on seven VMs using the "continuous orbit" function on a CAD drawing. This creates a graphics intensive environment. The "Sun and Sky" demo file from Autodesk is used for simulation. The object in the drawing is rotated at high speed using the "continuous orbit" function.
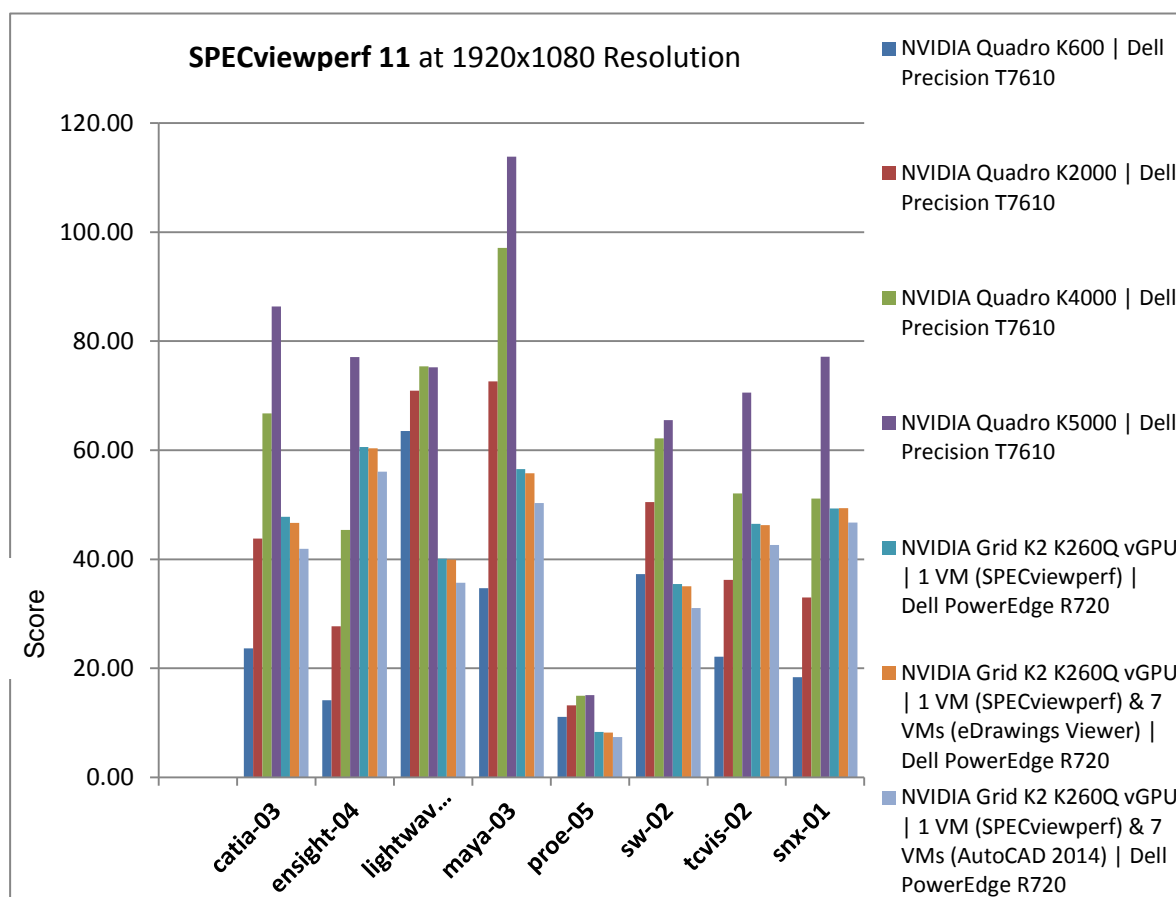
**SPECviewperf 11**

SPECviewperf is a widely used benchmark in the workstation domain for benchmarking graphics performance and on this basis it has been chosen as the appropriate benchmark to use for this high-end graphics use case. The SPECviewperf workloads were run using the Dell Wyse Z90Q7 end point.

Three different tests were run against the K260Q vGPU profile:

- SPECviewperf test run on a single VM against Wyse Z90Q7 with no companion tests running.

- SPECviewperf test run against Wyse Z90Q7 with Solidworks eDrawings Viewer companion workload on the remaining 7 VMs.

- SPECviewperf test run against Wyse Z90Q7 endpoint with AutoCAD 2014 companion workload on the remaining 7 VMs.

The SPECviewperf requirement is to have the minimum resolution of 1920x1080. All three tests, a mentioned above, were performed at that resolution. The SPECviewperf results recorded are presented in the following graph:


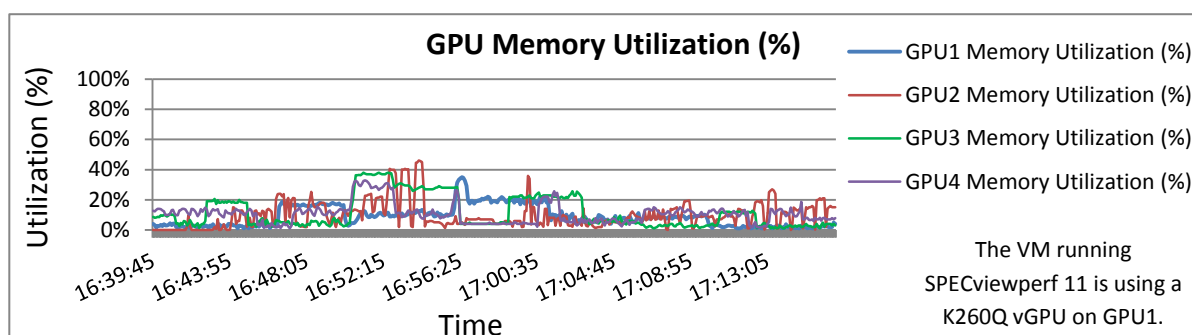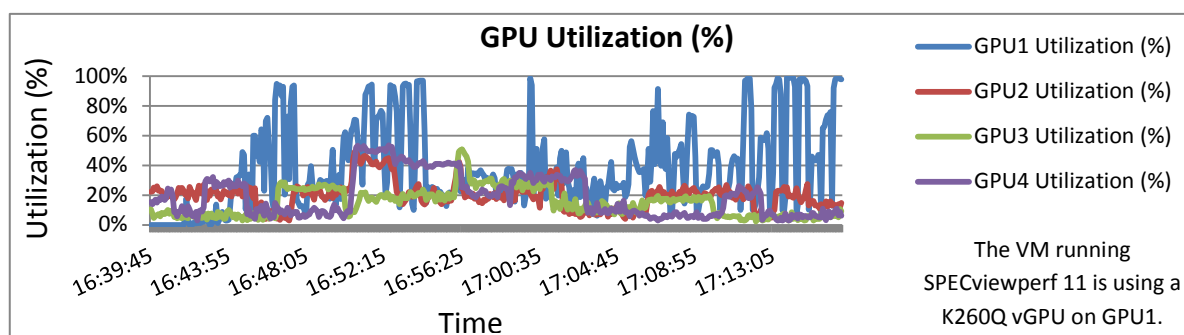
**End-User Experience & Resolution**

While running the workloads as mentioned above, we paid close attention to the end-user experience with a Dell Wyse Z90Q7 end point. One of the observations with the K240Q vGPU profile was that the user experience remained very good at very high resolutions. This was a major difference as compared to the K140Q vGPU profile.

User experience for AutoCAD 2014 was very good with crisp visual clarity and great colors. The rotation and orbiting of the 3D object at very high speed was smooth.

### GPU Utilization

While AutoCAD 2014 and SPECviewperf were running, the physical GPU utilization and the GPU video memory utilization was continuously monitored. The goal was to ensure that the GPU cores were not saturated. As you can see, the GPU utilization goes as high as 100%. This shows that running AutoCAD 2014 and SPECviewperf is very graphics intensive for the K260Q vGPU profile. The video memory utilization is below 45% for all GPUs. This represents the most load seen on the R720 host during the validation for K260Q vGPU profile.



The VM running SPECviewperf 11 is using a K260Q vGPU on GPU1.



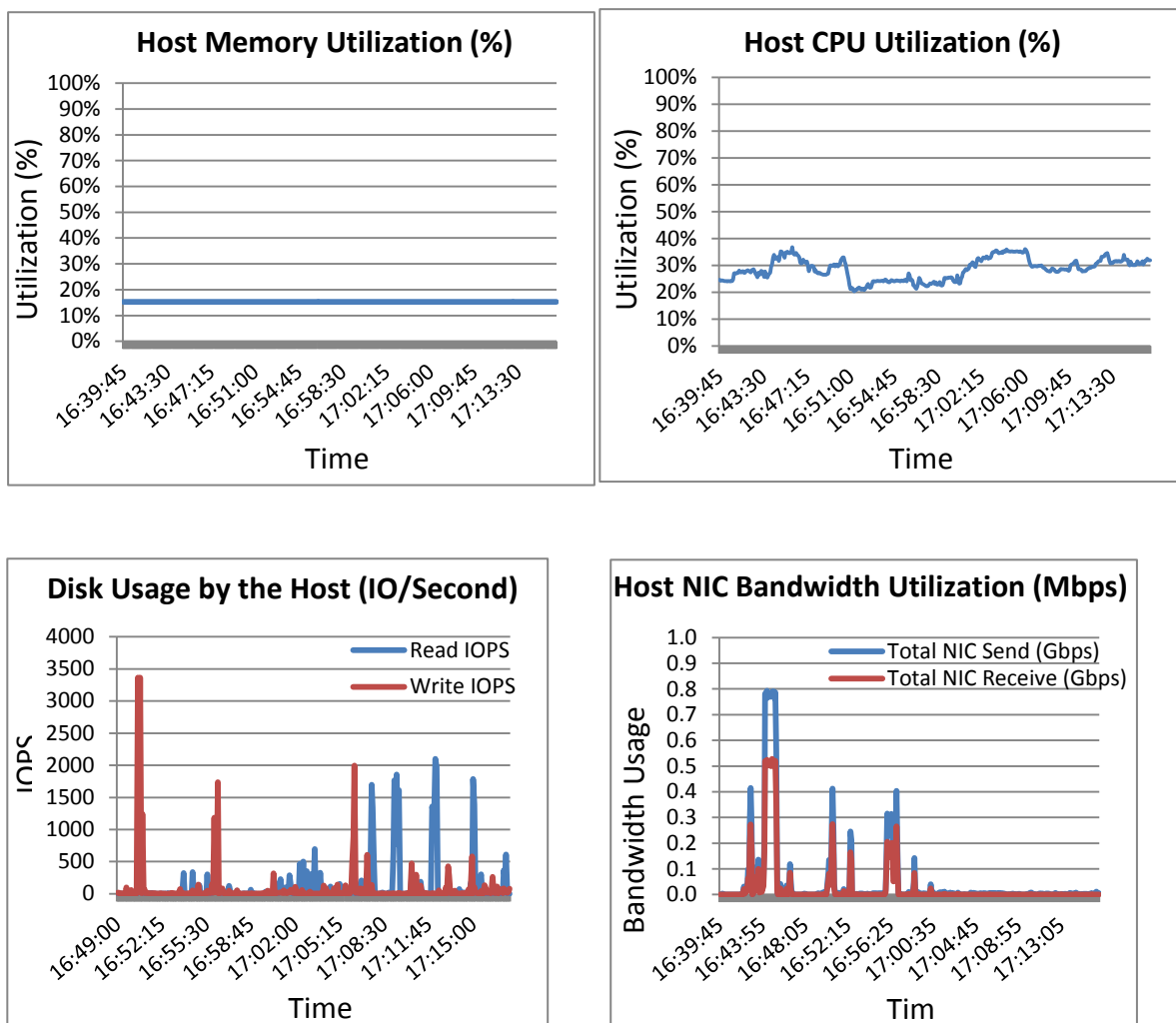The VM running SPECviewperf 11 is using a K260Q vGPU on GPU1.

### Host Resource Utilization

During the activities described above, results were gathered for the host performance during various workloads. This was done to ensure that the host was not stressed and the recorded results fall within the thresholds below:

| NVIDIA Grid Card | NVIDIA vGPU Profile Under Test |
|---|---|
| Compute Host CPU | 85% |
| Compute Host Memory | 85% |
| Network Throughput | 85% |

Results were gathered using scripts running on the XenServer host. The results shown below were gathered when running AutoCAD 2014 on seven VMs and SPECviewperf 11 on one VM. These represent the most load seen on the R720 host during the validation for K260Q vGPU profile.



**Host Memory Utilization (%)**



**Host CPU Utilization (%)**



**Disk Usage by the Host (IO/Second)**



**Host NIC Bandwidth Utilization (Mbps)**

As evident from the plots above, the utilization is well below the desired threshold. The table below shows the peak for each of these plots.

| NVIDIA Grid Card | NVIDIA vGPU Profile Under Test |
|---|---|
| Compute Host CPU | 37% |
| Compute Host Memory | 15% |
| Network Throughput | 80% |

Dell Precision R7610

These results describe validation efforts undertaken on XenServer using NVidia K2a Grid cards to study its behavior when used with 3D graphic-intensive applications. NVIDIA Grid K2a card contains two physical GPUs on the board. Each physical GPU (pGPU) can host several different types of virtual GPUs (vGPU). Alternative, entire physical GPU can also be assigned to a VM. These vGPU types are preset profiles with a fixed amount of frame-buffer, fixed number of supported display heads and maximum resolutions, and fixed number of CUDA cores designed for different classes of workload. Each vGPU profile is suitable for different use case types. This validation effort concentrates on the K260Q vGPU profile on the K2a card; and also on physical GPU passthrough.

The validation activity was undertaken to validate the following configuration of Dell Precision R7610 workstation as the compute host for the XenDesktop VMs.

- CPU: 2x Intel Xeon E5-2697v2 Processor (2.7 GHz)
- Memory: 128 GB (6x 16GB DIMMs @ 1600 MHz)

To validate the vGPU profile and the passthrough mode, the following VM user profile has been used.

- "Professional Shared Graphics" users are those who may be consuming relatively high-end graphics through relatively high frame-rate applications such as editing electrical, mechanical CAD drawings.
    - **vGPU Profile:** GRID K260Q
    - **VM Profile:** 2 vCPU, 4GB RAM with Windows 7 SP1 x64

The validation of this use case is performed with 2x Grid K2a cards. There is a finite limit on the number of vGPUs that can be created on a given card. This validation is performed with the maximum number of supported vGPUs for K260Q vGPU profile, and also with the passthrough mode.

| NVIDIA Grid Card | NVIDIA vGPU Profile Under Test | No. of Cards Per Host Server | No. of Physical GPUs Per Card | No. of Physical GPUs per Server | No. of vGPUs per pGPU | No. of vGPUs per card | No. of vGPUs per Server | No. of supported VMs |
|---|---|---|---|---|---|---|---|---|
| K2a | Passthrough | 3 | 2 | 6 | - | - | - | 6 |
| K2a | K260Q | 2 | 2 | 4 | 2 | 4 | 8 | 8 |

To test with the maximum number of supported vGPUs, an appropriate number of VMs are created on the server and each VM is assigned a vGPU or a physical GPU in pass-through mode. No other VMs are provisioned on the server except those involved in this graphics testing.

## NVIDIA K260Q vGPU Profile

To validate the K260Q vGPU profile, two Grid K2a cards are used to support total of eight VMs. The workload is simulated using AutoCAD 2014, and SPECviewperf 11 is used to measure performance. AutoCAD 2014 is run on seven out of the eight VMs with SPECviewperf running on the 8[th] VM.

AutoCAD running on seven VMs creates load on the physical GPU, creating load situations for the performance measurement using SPECviewperf.

**AutoCAD 2014**
In this validation activity, AutoCAD 2014 is used to view 3D CAD drawings. This viewing activity is simulated on seven VMs using the "continuous orbit" function on a CAD drawing. This creates a graphics intensive environment. The "Sun and Sky" demo file from Autodesk is used for simulation. The object in the drawing is rotated at high speed using the "continuous orbit" function.
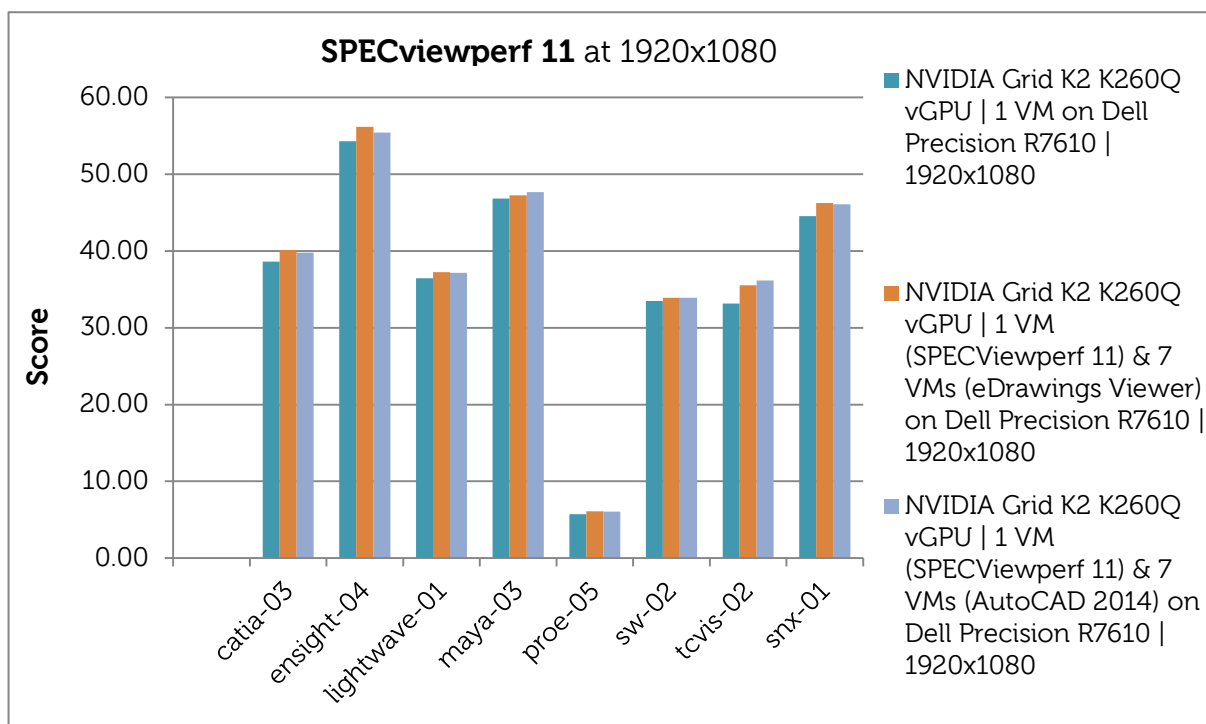
**SPECviewperf 11**
SPECviewperf is a widely used benchmark in the workstation domain for benchmarking graphics performance and on this basis it has been chosen as the appropriate benchmark to use for this high-end graphics use case. The SPECviewperf workloads were run using the Dell Wyse Z90Q7 end point.

Although the focus of the activity is on using AutoCAD 2014 as the graphics intensive workload, the SPECviewperf results were also collected with other workloads and scenarios in order to better understand the workload impact. Three different tests were run against the K260Q vGPU profile:

- SPECviewperf test run on a single VM against Wyse Z90Q7 with no companion tests running.

- SPECviewperf test run against Wyse Z90Q7 with Solidworks eDrawings Viewer companion workload on the remaining 7 VMs.

- SPECviewperf test run against Wyse Z90Q7 endpoint with AutoCAD 2014 companion workload on the remaining 7 VMs.

The SPECviewperf requirement is to have the minimum resolution of 1920x1080. All three tests, a mentioned above, were performed at that resolution. The SPECviewperf results recorded are presented in the following graph.
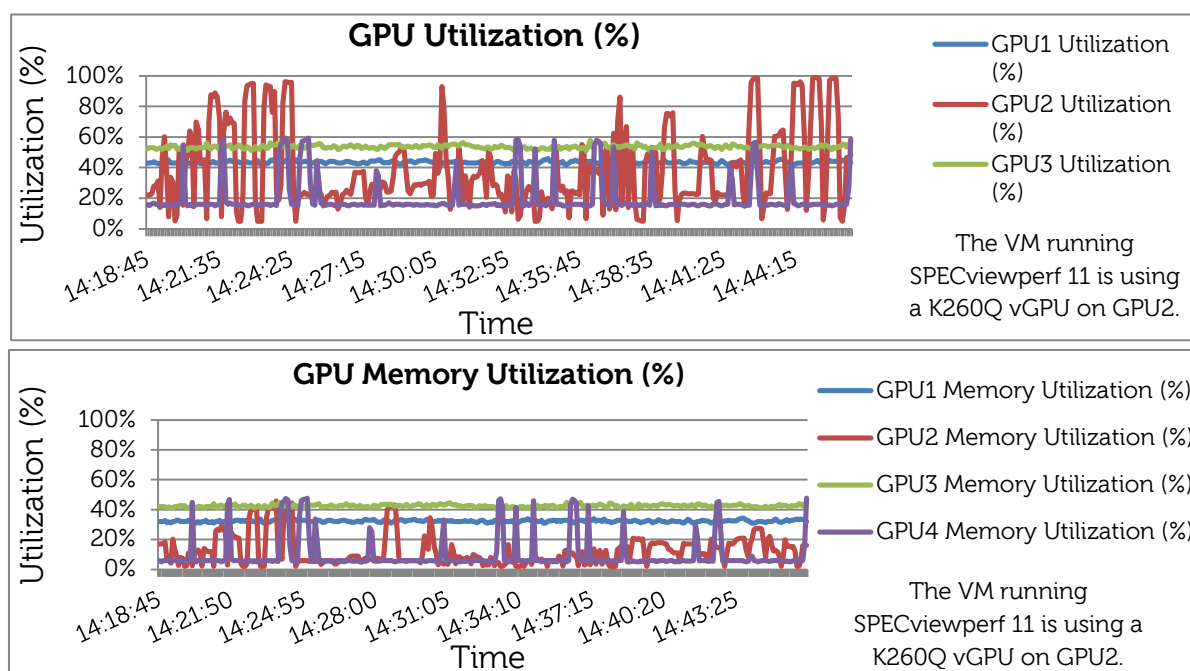
**End-User Experience & Resolution**

While running the workloads as mentioned above, we paid close attention to the end-user experience with a Dell Wyse Z90Q7 end point. One of the observations with the K260Q vGPU profile was that the user experience remained very good at very high resolutions.

User experience for AutoCAD 2014 was very good with crisp visual clarity and great colors. The rotation and orbiting of the 3D object at very high speed was smooth.

**GPU Utilization**

While AutoCAD 2014 and SPECviewperf were running, the physical GPU utilization and the GPU video memory utilization was continuously monitored. The goal was to ensure that the GPU cores were not saturated. As you can see, the GPU utilization goes as high as 100% but not at all times. This shows that running AutoCAD 2014 and SPECviewperf is graphics intensive for the K260Q vGPU profile. The video memory utilization is below 45% for all GPUs. This represents the most load seen on the R7610 host during the validation of K260Q vGPU profile.
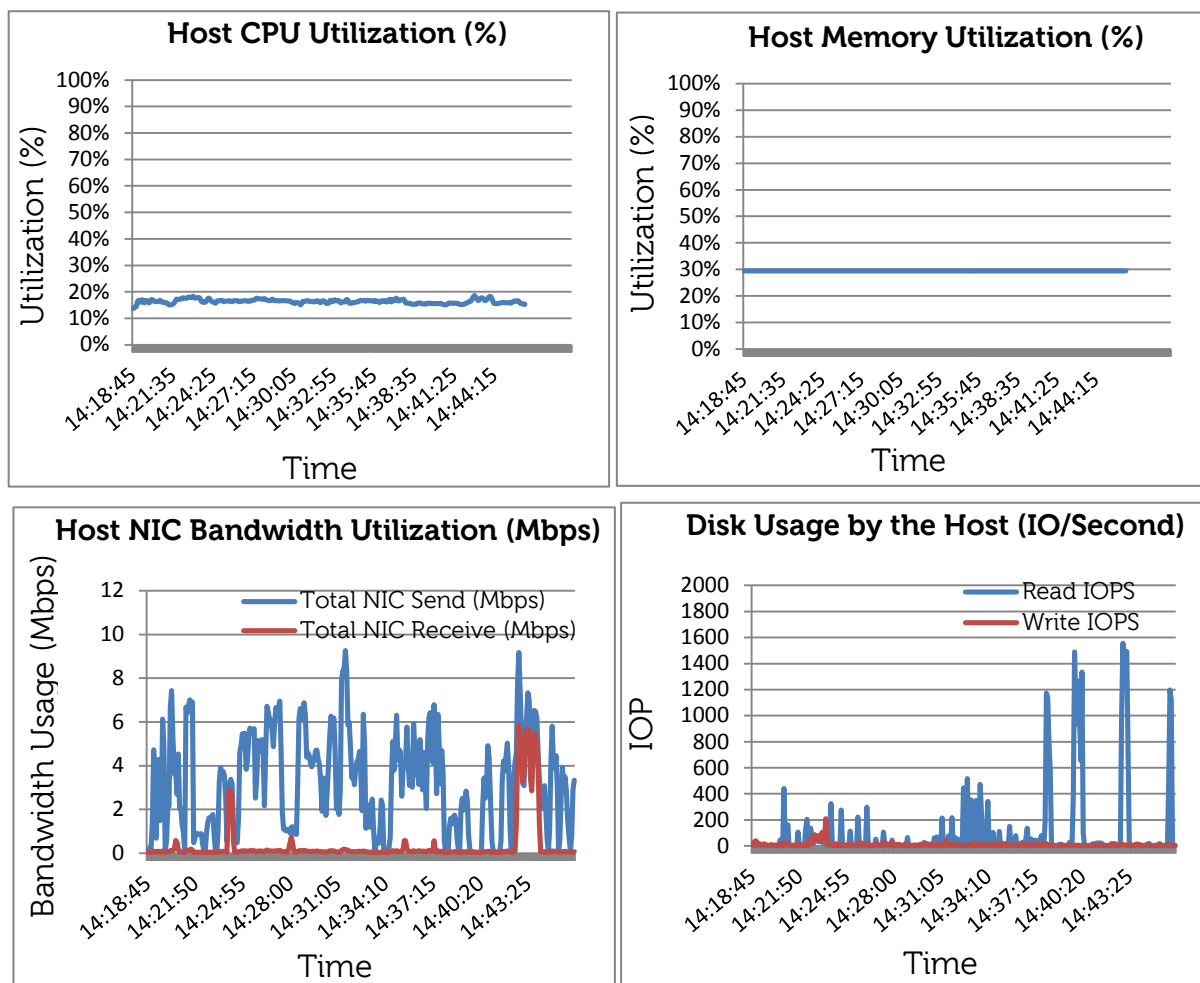


**Host Resource Utilization**

During the activities described above, results were gathered for the host performance during various workloads. This was done to ensure that the host was not stressed and the recorded results fall within the thresholds below:

| NVIDIA Grid Card | NVIDIA vGPU Profile Under Test |
| --- | --- |
| Compute Host CPU | 85% |
| Compute Host Memory | 85% |
| Network Throughput | 85% |

Results were gathered using scripts running on the XenServer host. The results shown below were gathered when running AutoCAD 2014 on seven VMs and SPECviewperf 11 on one VM. These represent the most load seen on the R7610 host during the validation of K260Q vGPU profile.



As evident from the plots above, the utilization is well below the desired threshold. The table below shows the peak for each of these plots:

| NVIDIA Grid Card | NVIDIA vGPU Profile Under Test |
|---|---|
| Compute Host CPU | 19% |
| Compute Host Memory | 29% |
| Network Throughput | 1.4% |

## NVIDIA K2a GPU Passthrough

To validate GPU pass-through, three Grid K2a cards are used to support total of six VMs. The workload is simulated using AutoCAD 2014, and SPECviewperf 11 is used to measure performance. AutoCAD 2014 is run on five out of the six VMs with SPECviewperf running on the 6[th] VM. AutoCAD running on five VMs creates load on the physical GPU, creating load situations for the performance measurement using SPECviewperf.

**AutoCAD 2014**
In this validation activity, AutoCAD 2014 is used to view 3D CAD drawings. This viewing activity is simulated on seven VMs using the "continuous orbit" function on a CAD drawing. This creates a graphics intensive environment. The "Sun and Sky" demo file from Autodesk is used for simulation. The object in the drawing is rotated at high speed using the "continuous orbit" function.
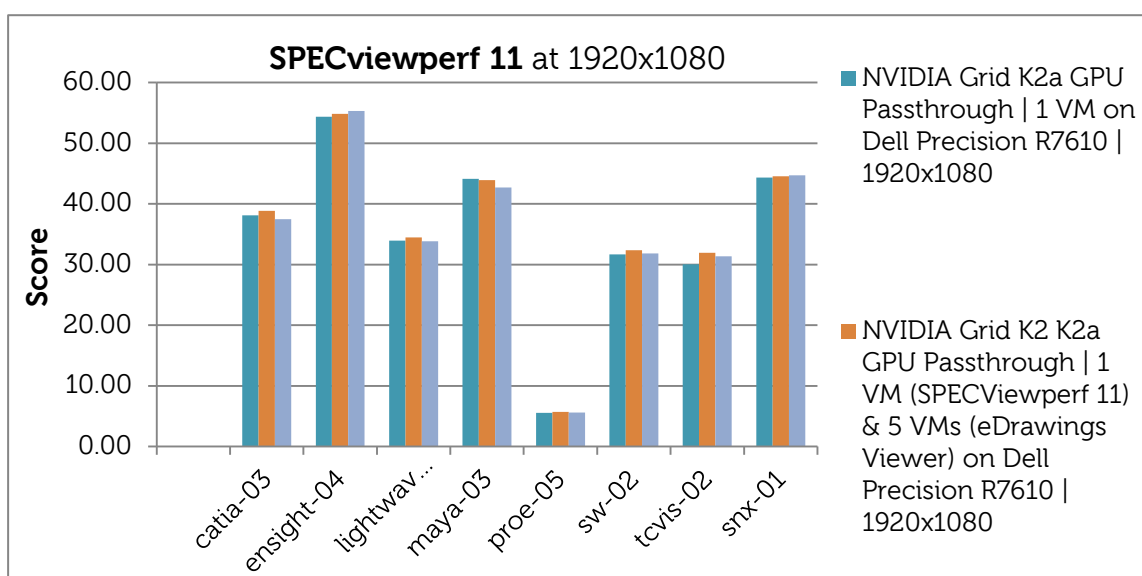
**SPECviewperf 11**
SPECviewperf is a widely used benchmark in the workstation domain for benchmarking graphics performance and on this basis it has been chosen as the appropriate benchmark to use for this high-end graphics use case. The SPECviewperf workloads were run using the Dell Wyse Z90Q7 end point.

Although the focus of the activity is on using AutoCAD 2014 as the graphics intensive workload, the SPECviewperf results were also collected with other workloads and scenarios in order to better understand the workload impact. Three different tests were run against the K260Q vGPU profile:

- SPECviewperf test run on a single VM against Wyse Z90Q7 with no companion tests running.

- SPECviewperf test run against Wyse Z90Q7 with Solidworks eDrawings Viewer companion workload on the remaining 5 VMs.

- SPECviewperf test run against Wyse Z90Q7 endpoint with AutoCAD 2014 companion w orkload on the remaining 5 VMs.

The SPECviewperf requirement is to have the minimum resolution of 1920x1080. All three tests, a mentioned above, were performed at that resolution. The SPECviewperf results recorded are presented in the following graph:

**End-User Experience & Resolution**

While running the workloads as mentioned above, we paid close attention to the end-user experience with a Dell Wyse Z90Q7 end point. One of the observations with the K240Q vGPU profile was that the user experience remained very good at very high resolutions.

User experience for AutoCAD 2014 was very good with crisp visual clarity and great colors. The rotation and orbiting of the 3D object at very high speed was smooth.
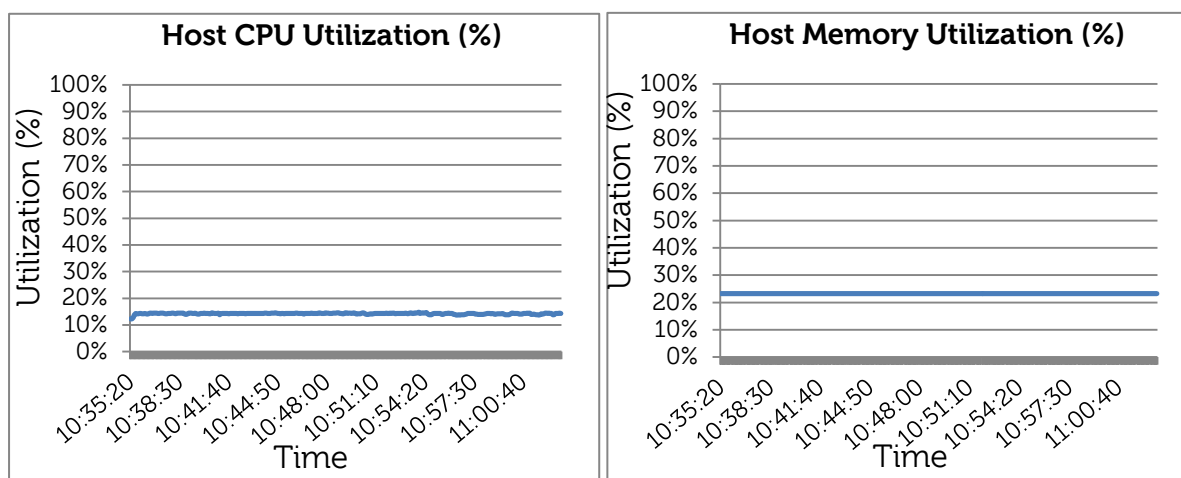
**GPU Utilization**

While AutoCAD 2014 and SPECviewperf were running, the physical GPU utilization and the GPU video memory utilization was continuously monitored. The goal was to ensure that the GPU cores were not saturated. The nvidia-smi commands are not available for passthrough mode, so the detailed plots could not be generated but the observations are noted below. The GPU utilization goes as high as 60% but not at all times. This shows that running AutoCAD 2014 and SPECviewperf is not extremely intensive for GPU passthrough profile. The video memory utilization is below 40% for all GPUs. This represents the most load seen on the R7610 host during the validation GPU pass-through profile.
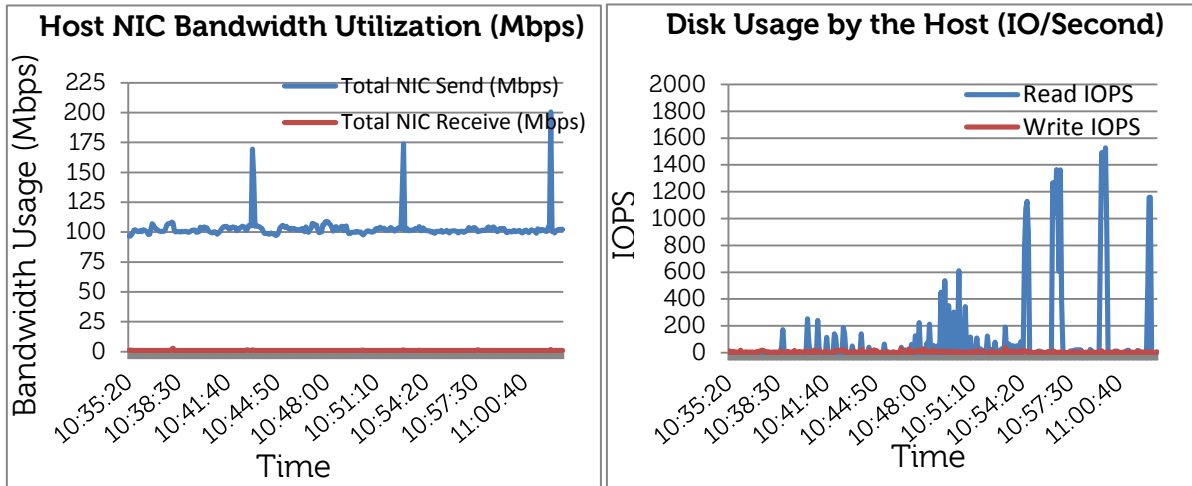
**Host Resource Utilization**

During the activities described above, results were gathered for the host performance during various workloads. This was done to ensure that the host was not stressed and the recorded results fall within the thresholds below:

| NVIDIA Grid Card | NVIDIA vGPU Profile Under Test |
|---|---|
| Compute Host CPU | 85% |
| Compute Host Memory | 85% |
| Network Throughput | 85% |

Results were gathered using scripts running on the XenServer host. The results shown below were gathered when running AutoCAD 2014 on five VMs and SPECviewperf 11 on one VM. These represent the most load seen on the R7610 host during the validation of GPU passthrough profile.

As evident from the plots above, the utilization is well below the desired threshold. The table below shows the peak for each of these plots:

| NVIDIA Grid Card | NVIDIA vGPU Profile Under Test |
|---|---|
| Compute Host CPU | 15% |
| Compute Host Memory | 23% |
| Network Throughput | 20% |

# Acknowledgements

# About the Authors

Peter Fine is the Sr. Principal Solutions Architect for Citrix-based solutions at Dell. Peter has extensive experience and expertise on the broader Microsoft, Citrix and VMware solutions software stacks as well as in enterprise virtualization, storage, networking and enterprise datacenter design.

Rick Biedler is the Solutions Development Manager for Citrix solutions at Dell, managing the development and delivery of Enterprise class Desktop virtualization solutions based on Dell Datacenter components and core virtualization platforms.

Pranav Parekh is a Sr. Solutions engineer at Dell Wyse Solutions Engineering group. Pranav has extensive experience designing desktop virtualization solutions, IaaS private cloud solutions, virtualization solutions, and enterprise class blade servers. Pranav has a master's degree in Electrical & Computer Engineering from the University of Texas at Austin.