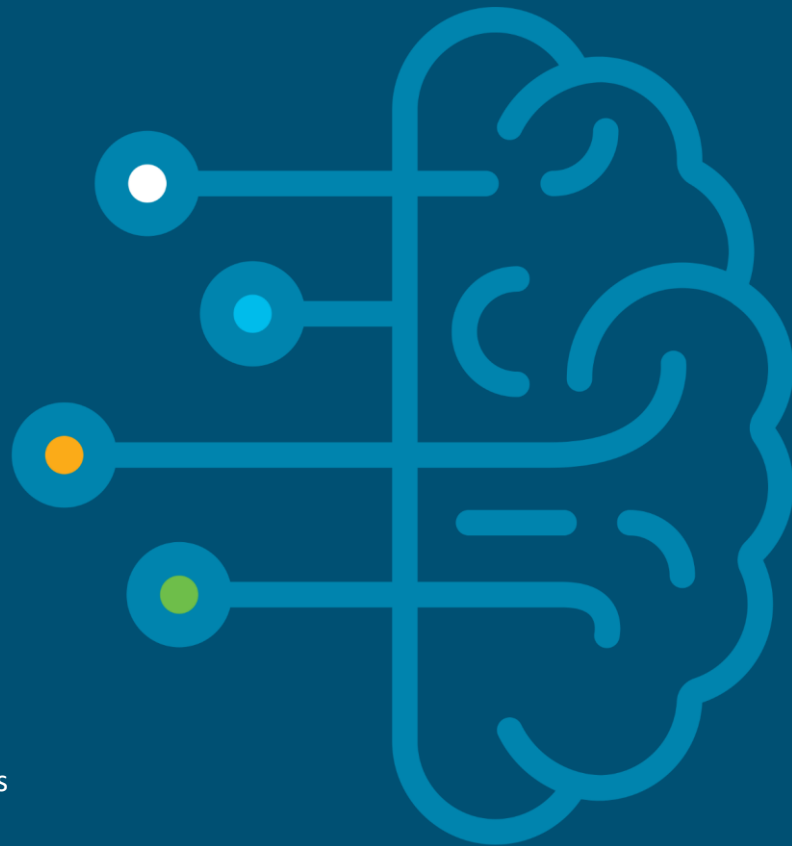# S9881 Using Industry Standard Benchmark Tools to Size Graphics Accelerated Applications

Mike Brennan, Product Manager, Virtual Client Computing and Graphics
Vadim Lebedev, Technical Marketing Engineer
21 March 2019

# What we will cover

- Let's talk about the basics on virtualizing pro graphics apps

- How do you measure performance?

- Key NVIDIA cards

- Sample benchmark performance

- Server/GPU performance

- Where do I start with sizing?

- Cisco lineup

- Key takeaways

- Q&A

# Lets talk basics

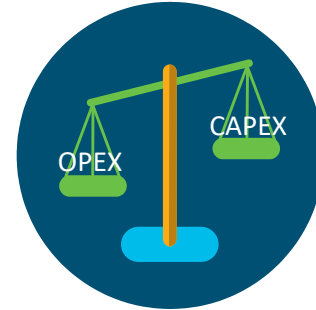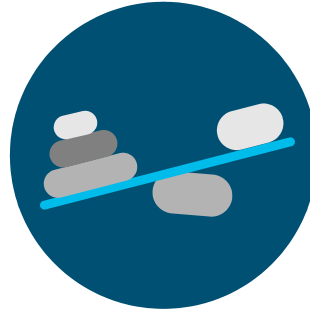# Can you virtualize Catia and SolidWorks?

Yes you can!

Cisco has 14 hardware/ software combinations certified

Dassault VDI Certifications

# Why replace physical graphics workstations?

# Key performance requirements for Virtual Workstations

User requirements

CPU and memory performance

Software requirements

Graphics card oversubscription

Display resolution

FPS
Frame rate

Monitors

Multi-user graphics card scheduling engine

# User requirements

User role

Concurrent applications open

Complexity of graphic application

Collaboration requirements

Working hours

# User roles

## Light user type

- Primarily read only – documentation, project managers
- Small subsets of entire entity

## Medium user type

- Read only and design
- Small and medium sub-assemblies

## Heavy user type

- Design and render
- Large sub-assemblies and full model

# Software requirements

Dassault minimum requirements

Dassault support for virtualization
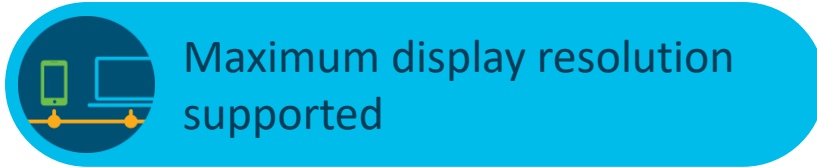
Dassault hardware qualification

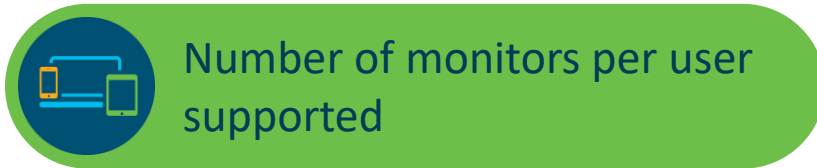Dassault delivery partners

Dassault support model

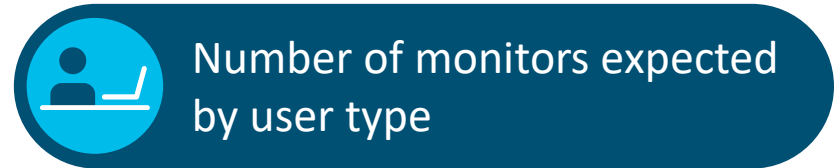# Display resolution and monitors

Maximum display resolution supported

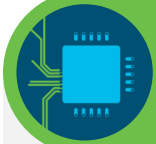Display resolution expected by user type

Number of monitors per user supported
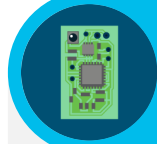
Number of monitors expected by user type

# CPU and memory performance

## CPU selection criteria

- Total frequency in MHz/CPU and server

- Core count

- Planned user count

## Memory selection criteria

- Balance

- Frequency

# Graphics card oversubscription

## User count per graphics card

- Fixed at GPU frame buffer divided by vGPU profile
    - For an NVIDIA P4 card
    - For a 2Q profile: 8GB frame buffer/2GB frame buffer per user = 4 Users per card.

## GPU oversubscription

- NVIDIA concept
- Based on scheduler chosen
- For the T4 card, light user could get more than 12.5% of GPU resources

# Frame rates

FPS

## The great equalizer for performance

- For computer video displays,

  - frame rate = #frames or images displayed per second

- For a given application

  - Provides a mechanism to compare systems performance

  - Describes a mechanism by which system requirements can be stated

## Virtual Graphics Workstation insights

- Frame rate can be controlled – or not

- Frame rate can be set in the NVIDIA and Desktop Broker software policy

# Multi-user graphics cards scheduling engines

## NVIDIA supports 3 models

**Best effort (default)**

- User gets GPU resources based on current availability

- At any given point in time a user MIGHT get more than his fair share of GPU

**Fixed share**

- Each user gets the same dedicated performance at all times

**Equal share**

- Each VM gets and equal share of the GPU resources

# How do you measure performance?

# Performance measurement

## Industry graphic benchmark examples

- SPECviewperf 13
- PassMark Software
- Unigen Heaven, Valley, etc
- Others

SPECviewperf 13 supports nine Virtual Professional Graphics Applications

SPECviewperf 13 provides a composite benchmark score across all nine applications

SPECviewperf 13 provides capability to score individual applications

SPECviewperf 13 provides ability to measure performance across various graphic card, CPU, memory, scheduling and frame rate scenarios

# SPECviewperf 13 has the following minimum requirements:

✔ Microsoft Windows 10 64-bit RS3 or later VM

✔ OpenGL 4.0

✔ Direct X12 support

✔ 8GB of installed system memory

✔ 80GB available disk space

✔ 1920x1080 screen resolution for submissions published on the SPEC website

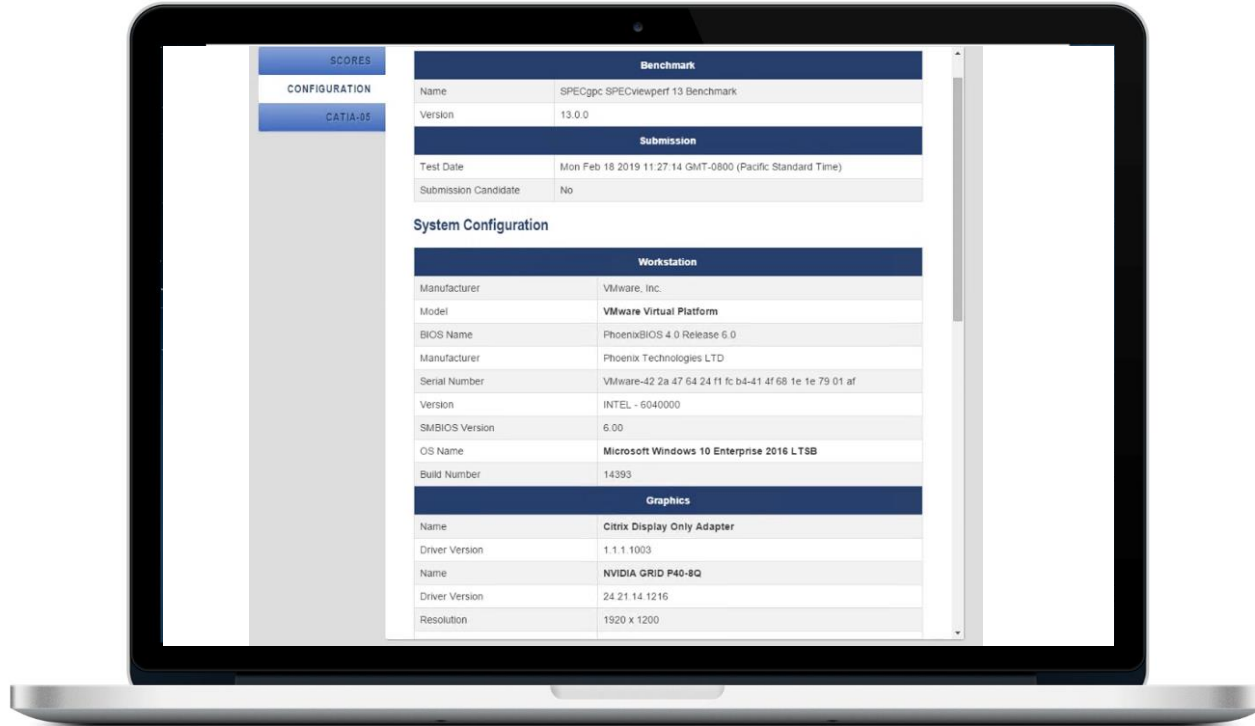# Applications driving large TAM, Verticals

# SPECviewperf 13 Test Console

# SPECviewperf 13 Test Results - Composite

# SPECviewperf 13 Test Results - Configuration

# SPECviewperf 13 Test Results – Viewset Catia

# NVIDIA Tesla T4 and P40

# Tesla T4 Key Specifications



| GPU Architecture | **NVIDIA Turing** |
|---|---|
| NVIDIA Turing Tensor Cores | **320** |
| NVIDIA CUDA® Cores | **2,560** |
| RT Cores | **40** |
| Giga Rays/second | **5** |
| Memory Size | **16 GB GDDR6** |
| Memory BW | **Up to 320 GB/s** |
| vGPU Profiles | **1 GB, 2 GB, 4 GB, 8 GB, 16 GB** |
| Form Factor | **PCIe 3.0 single slot (half height & length)** |
| Power | **70W** |
| Thermal | **Passive** |

# Tesla P6 Key Specifications



| GPU | 1 NVIDIA Pascal GPU |
|---|---|
| **CUDA Cores** | 2,048 |
| **Memory Size** | 16 GB GDDR5 |
| **H.264 1080p30 streams** | 24 |
| **Max vGPU instances** | 16 (1 GB Profile) |
| **vGPU Profiles** | 1 GB, 2 GB, 4 GB, 8 GB, 16 GB |
| **Form Factor** | MXM (blade servers) |
| **Power** | 90 W (70 W opt) |
| **Thermal** | Bare Board |

# Tesla P40 Key Specifications

| | |
|---|---|
| **GPU** | 1 NVIDIA Pascal GPU |
| **CUDA Cores** | 3,840 |
| **Memory Size** | 24 GB GDDR5 |
| **H.264 1080p30 streams** | 24 |
| **Max vGPU instances** | 24 (1 GB Profile) |
| **vGPU Profiles** | 1 GB, 2 GB, 3 GB, 4 GB, 6 GB, 8 GB, 12 GB, 24 GB |
| **Form Factor** | PCIe 3.0 Dual Slot (rack servers) |
| **Power** | 250 W |
| **Thermal** | Passive |

# NVIDIA TESLA GPUs
## Recommended for Virtualization

| | V100 | P40 | T4 | P4 | M60 | M10 | P6 |
|---|---|---|---|---|---|---|---|
| GPUs / Board (Architecture) | 1 (Volta) | 1 (Pascal) | 1 (Turing) | 1 (Pascal) | 2 (Maxwell) | 4 (Maxwell) | 1 (Pascal) |
| CUDA Cores | 5,120 | 3,840 | 2,560 | 2,560 | 4,096 (2,048 per GPU) | 2,560 (640 per GPU) | 2,048 |
| Memory Size | 32 GB/16 GB HBM2 | 24 GB GDDR5 | 16 GB GDDR6 | 8 GB GDDR5 | 16 GB GDDR5 (8 GB per GPU) | 32 GB GDDR5 (8 GB per GPU) | 16 GB GDDR5 |
| vGPU Profiles | 1 GB, 2 GB, 4 GB, 8 GB, 16 GB, 32 GB | 1 GB, 2 GB, 3 GB, 4 GB, 6 GB, 8 GB, 12 GB, 24 GB | 1 GB, 2 GB, 4 GB, 8 GB, 16 GB | 1 GB, 2 GB, 4 GB, 8 GB | 0.5 GB, 1 GB, 2 GB, 4 GB, 8 GB | 0.5 GB, 1 GB, 2 GB, 4 GB, 8 GB | 1 GB, 2 GB, 4 GB, 8 GB, 16 GB |
| Form Factor | PCIe 3.0 Dual Slot & SXM2 (rack servers) | PCIe 3.0 Dual Slot (rack servers) | PCIe 3.0 Single Slot (rack servers) | PCIe 3.0 Single Slot (rack servers) | PCIe 3.0 Dual Slot (rack servers) | PCIe 3.0 Dual Slot (rack servers) | MXM (blade servers) |
| Power | 250W/300W | 250W | 70W | 75W | 300W (225W opt) | 225W | 90W |
| Thermal | passive | passive | passive | passive | active/passive | passive | bare board |

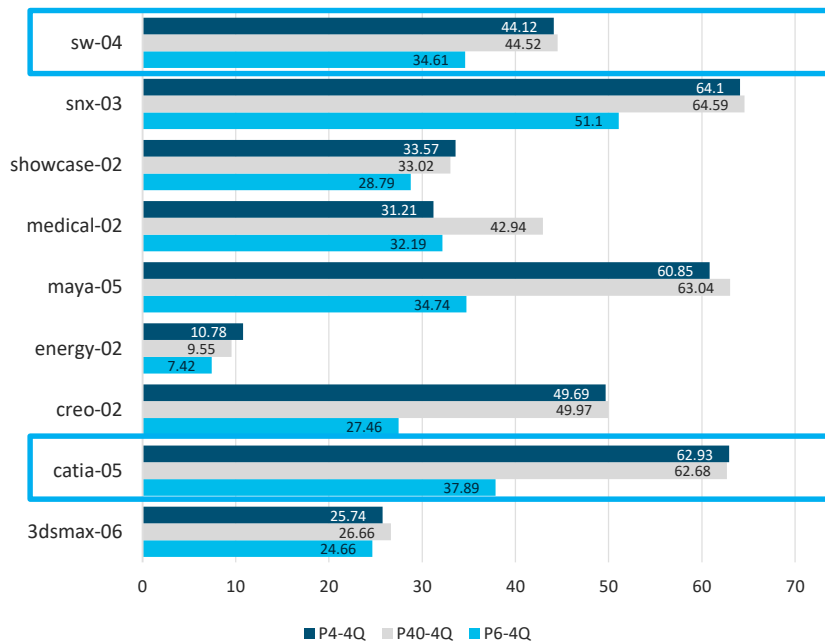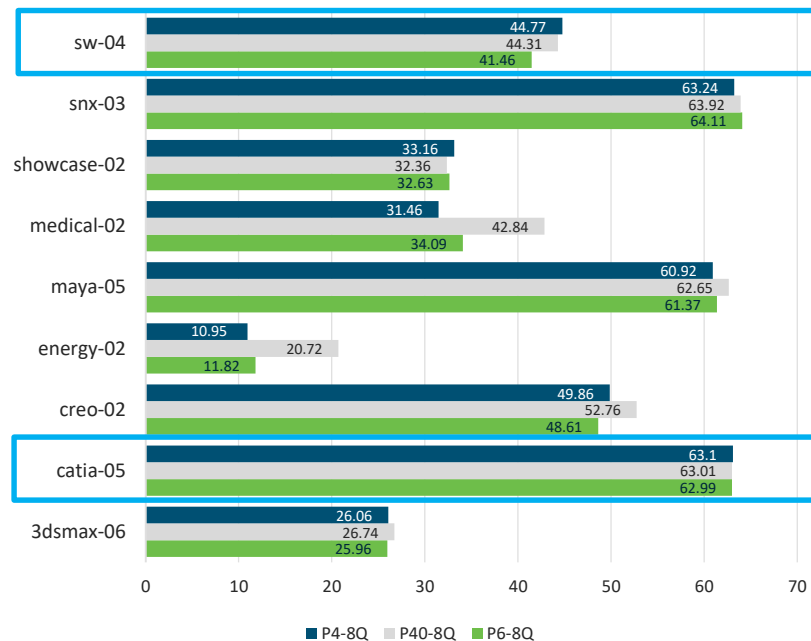| | | PERFORMANCE Optimized | | | DENSITY Optimized | BLADE Optimized |
|---|---|---|---|---|---|---|

# Example Benchmark Insights

# SPECviewperf 13 benchmark results
## Compare three cards, two profiles, 1 VM, best effort, FRL On, Xeon 6140

**1 VM 4Q on XenServer Host**

| Benchmark | P4-4Q | P40-4Q | P6-4Q |
|---|---|---|---|
| sw-04 | 44.12 | 44.52 | 34.61 |
| snx-03 | 64.1 | 64.59 | 51.1 |
| showcase-02 | 33.57 | 33.02 | 28.79 |
| medical-02 | 31.21 | 42.94 | 32.19 |
| maya-05 | 60.85 | 63.04 | 34.74 |
| energy-02 | 10.78 | 9.55 | 7.42 |
| creo-02 | 49.69 | 49.97 | 27.46 |
| catia-05 | 62.93 | 62.68 | 37.89 |
| 3dsmax-06 | 25.74 | 26.66 | 24.66 |

**1 VM 8Q on XenServer Host**

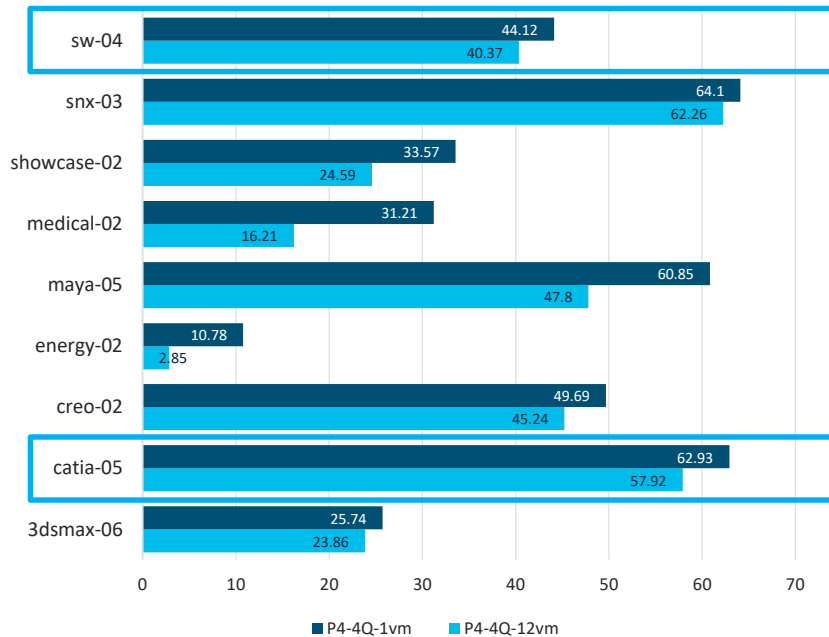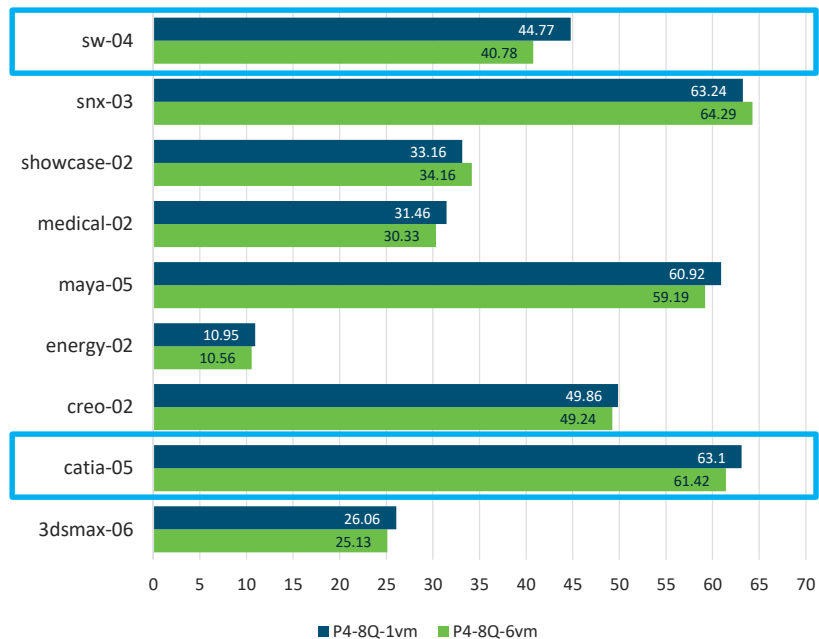| Benchmark | P4-8Q | P40-8Q | P6-8Q |
|---|---|---|---|
| sw-04 | 44.77 | 44.31 | 41.46 |
| snx-03 | 63.24 | 63.92 | 64.11 |
| showcase-02 | 33.16 | 32.36 | 32.63 |
| medical-02 | 31.46 | 42.84 | 34.09 |
| maya-05 | 60.92 | 62.65 | 61.37 |
| energy-02 | 10.95 | 20.72 | 11.82 |
| creo-02 | 49.86 | 52.76 | 48.61 |
| catia-05 | 63.1 | 63.01 | 62.99 |
| 3dsmax-06 | 26.06 | 26.74 | 25.96 |

# SPECviewperf 13 benchmark results
## Compare one card, two profiles, 1 VM & max VMs, best effort and FRL On, Xeon 6140
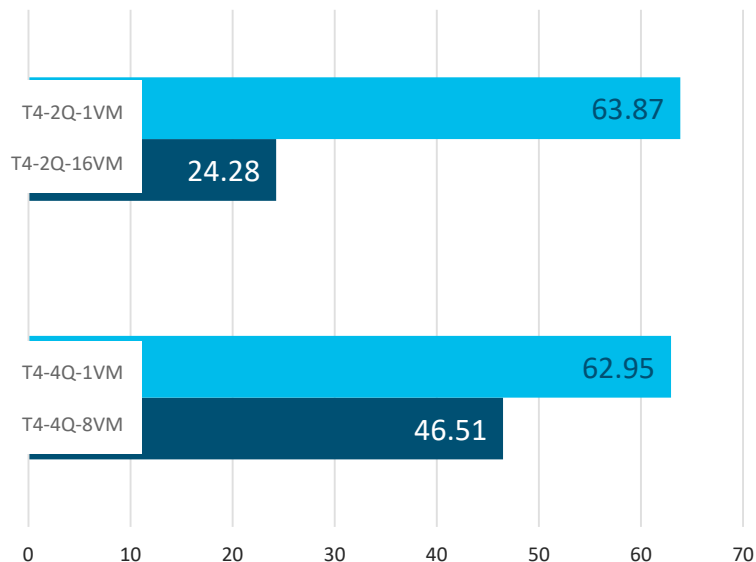
### P-4 4Q profile testing

| Benchmark | P4-4Q-1vm | P4-4Q-12vm |
|---|---|---|
| sw-04 | 44.12 | 40.37 |
| snx-03 | 64.1 | 62.26 |
| showcase-02 | 33.57 | 24.59 |
| medical-02 | 31.21 | 16.21 |
| maya-05 | 60.85 | 47.8 |
| energy-02 | 10.78 | 2.85 |
| creo-02 | 49.69 | 45.24 |
| catia-05 | 62.93 | 57.92 |
| 3dsmax-06 | 25.74 | 23.86 |

■ P4-4Q-1vm ■ P4-4Q-12vm

### P-4 8Q profile testing

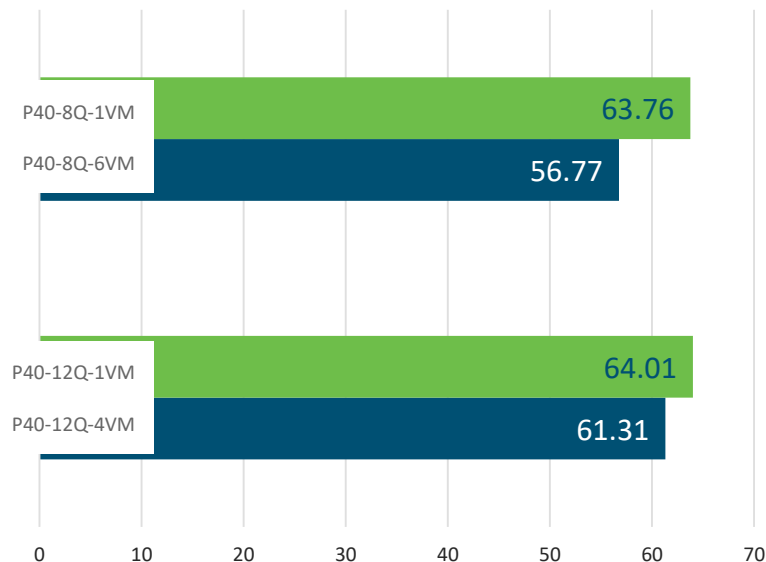| Benchmark | P4-8Q-1vm | P4-8Q-6vm |
|---|---|---|
| sw-04 | 44.77 | 40.78 |
| snx-03 | 63.24 | 64.29 |
| showcase-02 | 33.16 | 34.16 |
| medical-02 | 31.46 | 30.33 |
| maya-05 | 60.92 | 59.19 |
| energy-02 | 10.95 | 10.56 |
| creo-02 | 49.86 | 49.24 |
| catia-05 | 63.1 | 61.42 |
| 3dsmax-06 | 26.06 | 25.13 |

■ P4-8Q-1vm ■ P4-8Q-6vm

# SPECviewperf 13 non-benchmark results - Catia
## Compare two cards, 1 VM & max VMs, best effort and FRL On, Xeon 6140
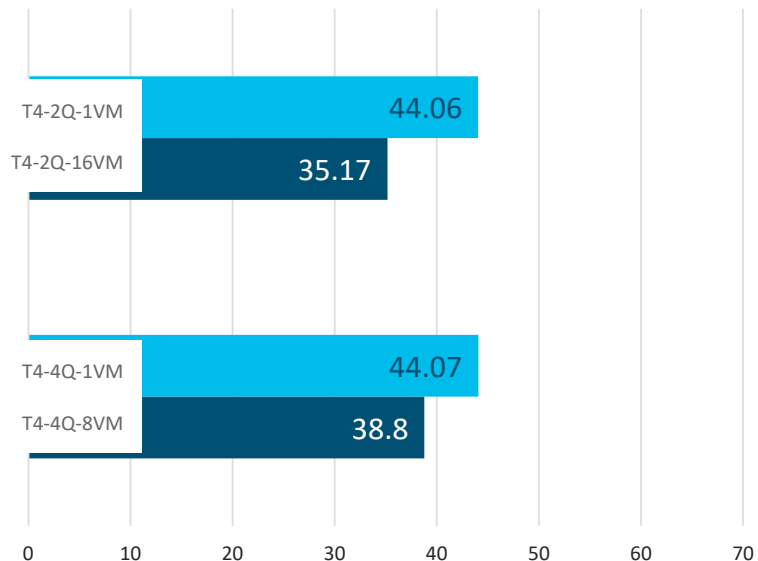
**T-4 2Q, 4Q profile testing**

| Label | Value |
|---|---|
| T4-2Q-1VM | 63.87 |
| T4-2Q-16VM | 24.28 |
| T4-4Q-1VM | 62.95 |
| T4-4Q-8VM | 46.51 |

**P40 8Q, 12Q**

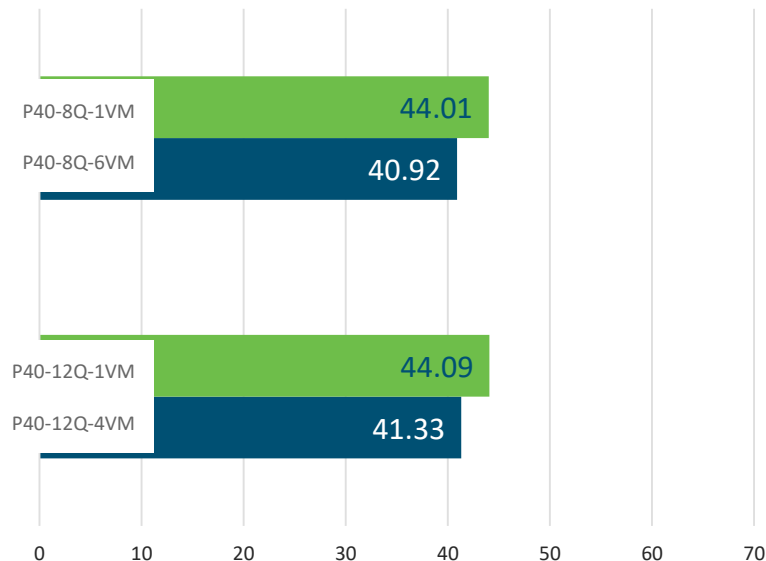| Label | Value |
|---|---|
| P40-8Q-1VM | 63.76 |
| P40-8Q-6VM | 56.77 |
| P40-12Q-1VM | 64.01 |
| P40-12Q-4VM | 61.31 |

# SPECviewperf 13 non-benchmark results -Solidworks
## Compare two cards, 1 VM & max VMs, best effort and FRL On, Xeon 6140
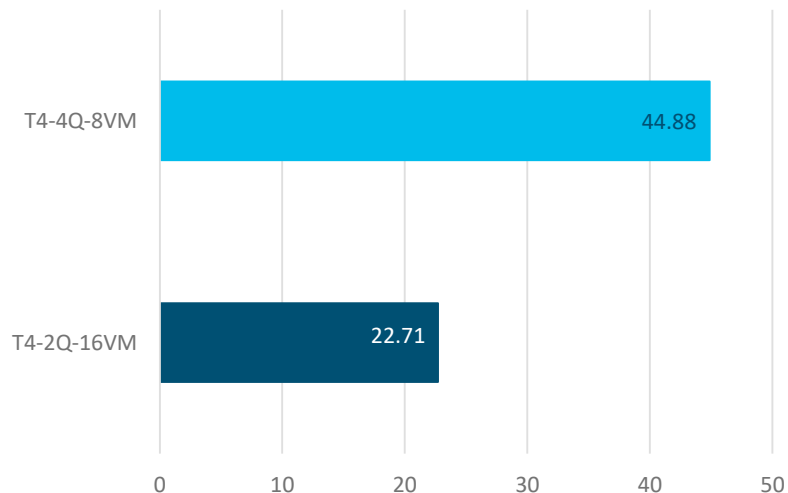


T-4 2Q, 4Q profile testing

| | Value |
|---|---|
| T4-2Q-1VM | 44.06 |
| T4-2Q-16VM | 35.17 |
| T4-4Q-1VM | 44.07 |
| T4-4Q-8VM | 38.8 |

P40 8Q, 12Q

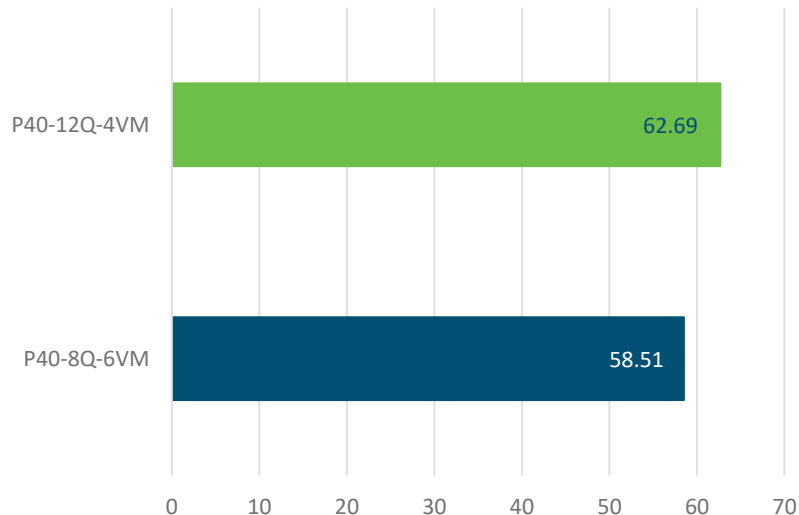| | Value |
|---|---|
| P40-8Q-1VM | 44.01 |
| P40-8Q-6VM | 40.92 |
| P40-12Q-1VM | 44.09 |
| P40-12Q-4VM | 41.33 |

# SPECviewperf 13 non-benchmark results - Catia
## Compare two cards, max VMs, best effort and FRL On, Xeon 6136, 6128



**T-4 2Q, 4Q profile, XEON 6136**

| Configuration | Value |
|---|---|
| T4-4Q-8VM | 44.88 |
| T4-2Q-16VM | 22.71 |

**P40 8Q, 12Q Profile, XEON 6128**

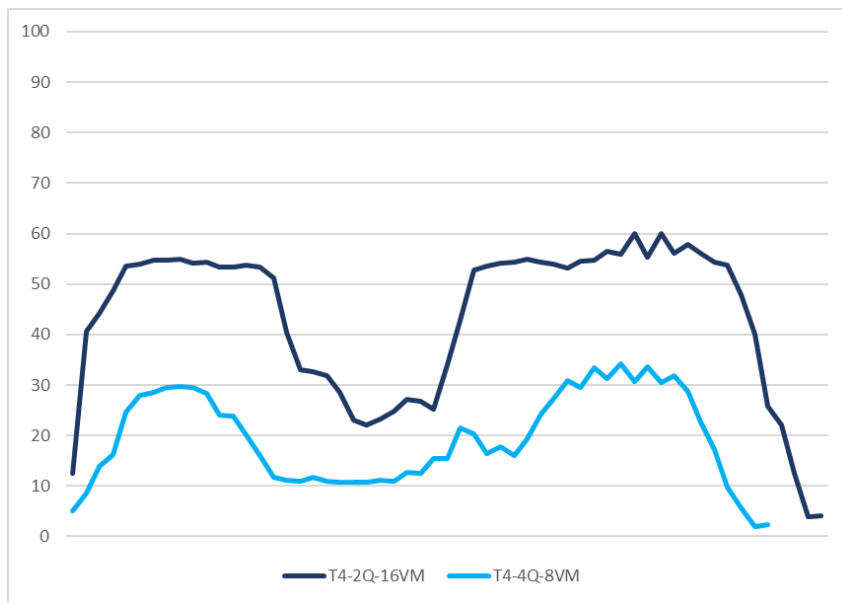| Configuration | Value |
|---|---|
| P40-12Q-4VM | 62.69 |
| P40-8Q-6VM | 58.51 |

# Tying the Benchmarks to CPUs and GPUs

# Intel Scalable Family 6140 and NVIDIA Tesla T4*



Intel Xeon 6140 Utilization

NVIDIA Tesla T4 Utilization
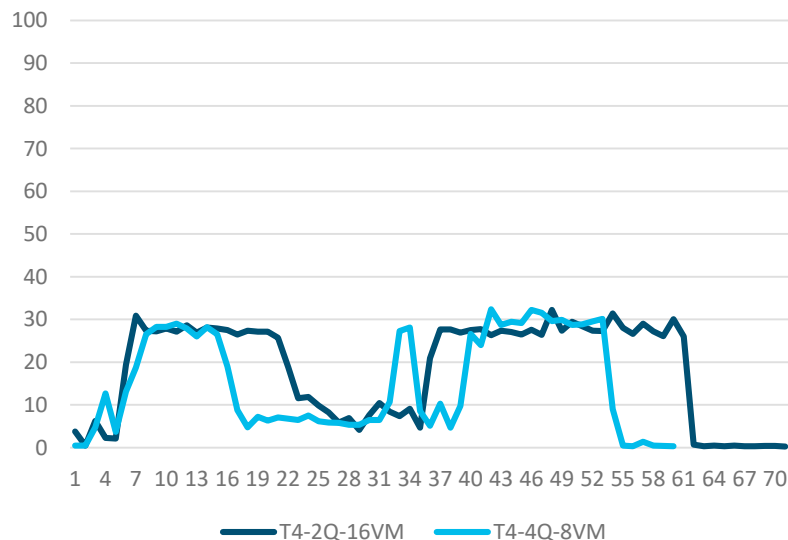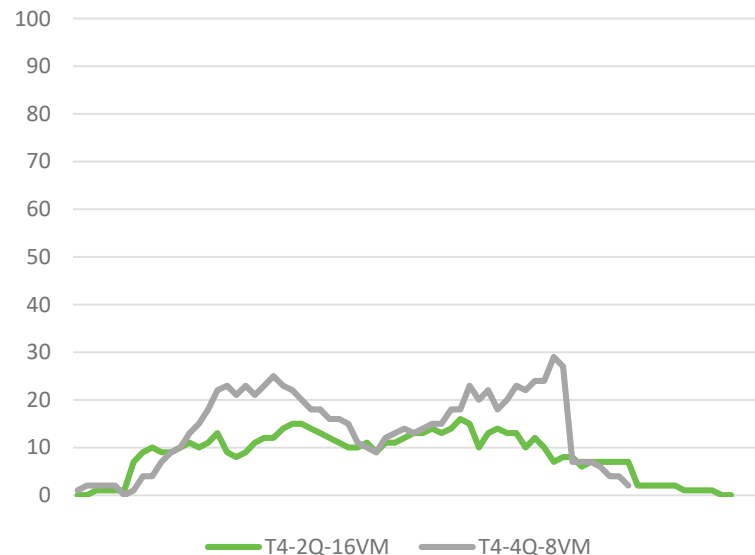
T4-2Q-16VM    T4-4Q-8VM

T4-2Q-16VM    T4-4Q-8VM

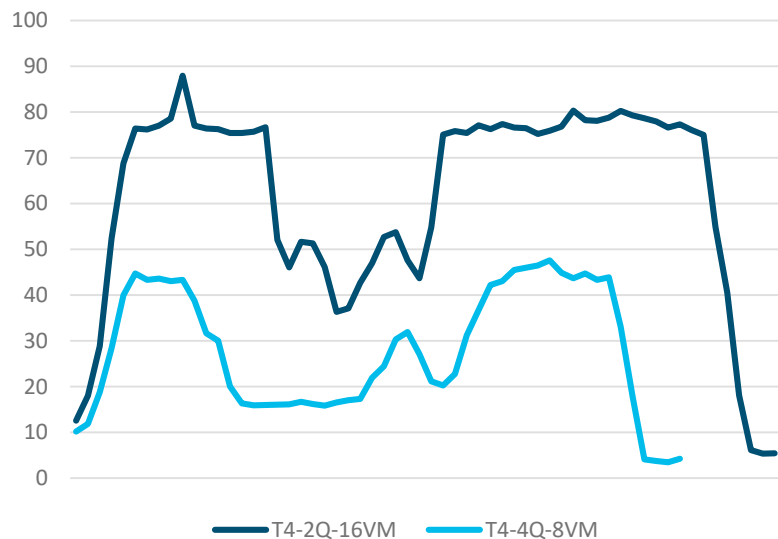* SPECviewperf 13 Catia Test – ESXi Host Data
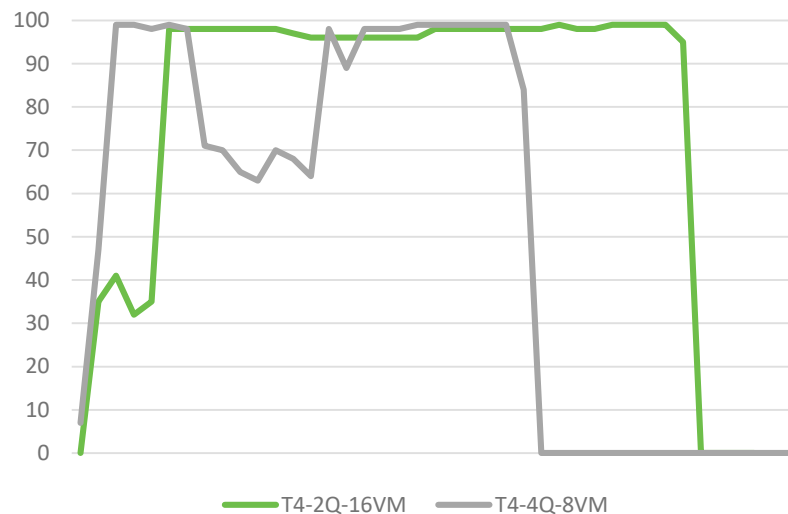
# Windows 10 (1607) VM with Tesla T4 vGPU*



* SPECviewperf 13 Catia Test – Single VM in multiple VM test Perfmon Data with Xeon 6140 host processor

# Intel Scalable Family 6136 and NVIDIA Tesla T4*
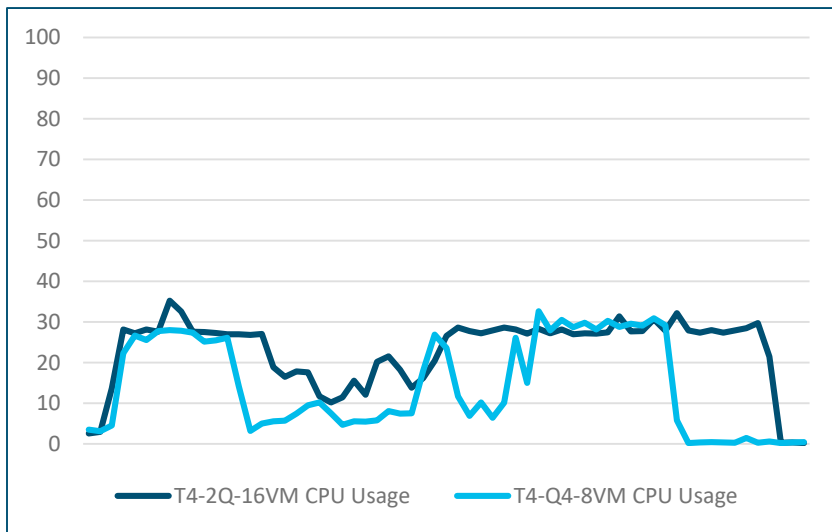
## Intel Xeon 6136 Utilization

Legend: T4-2Q-16VM, T4-4Q-8VM

## NVIDIA Tesla T4 Utilization

Legend: T4-2Q-16VM, T4-4Q-8VM

\* SPECviewperf 13 Catia Test – ESXi Host Data

# Windows 10 (1607) VM with Tesla T4 vGPU*



Perfmon VM Processor Utilization

Perfmon VM GPU Utilization

* SPECviewperf 13 Catia Test – Single VM in multiple VM test Perfmon Data with Xeon 6136 host processor

# Intel Scalable Family 6140 and NVIDIA Tesla P40*

## Intel Xeon 6140 Utilization

Legend: P40-8Q-6VM, P40-12Q-4VM

## NVIDIA Tesla P40 Utilization
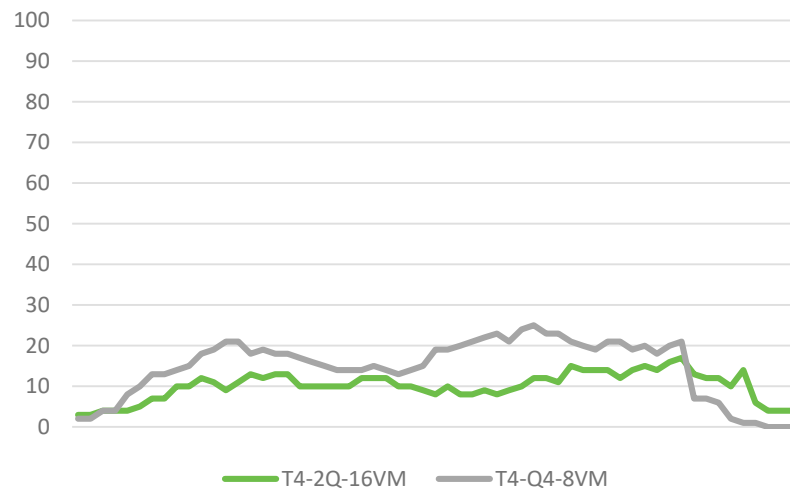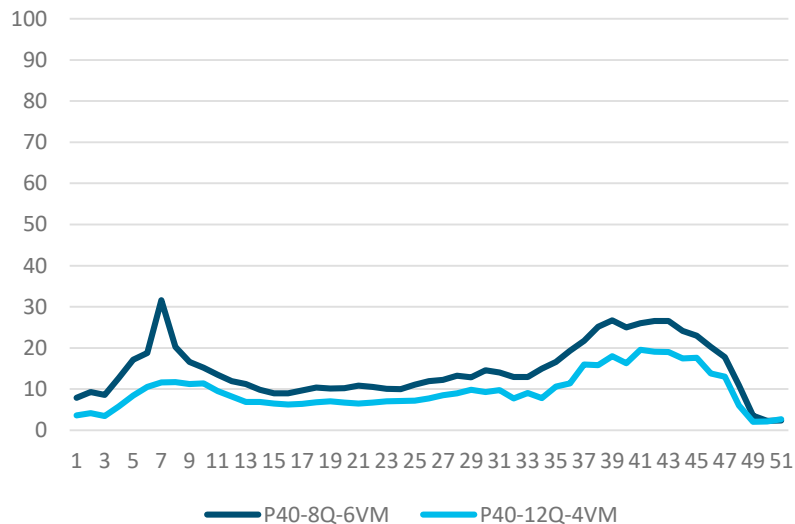
Legend: P40-8Q-6VM, P40-12Q-4VM

* SPECviewperf 13 Catia Test – ESXi Host Data

# Windows 10 (1607) VM with Tesla P40 vGPU*



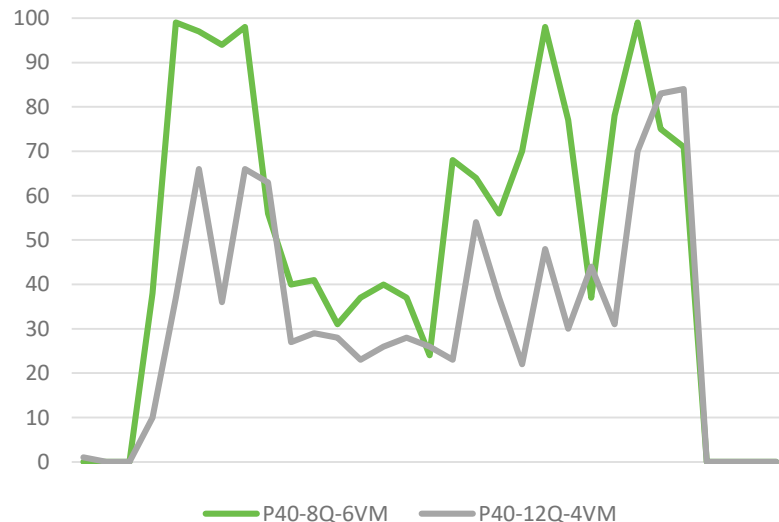Perfmon VM Processor Utilization

Perfmon VM GPU Utilization

P40-8Q-6VM    P40-12Q-4VM

P40-8Q-6VM    P40-12Q-4VM

* SPECviewperf 13 Catia Test – Single VM in multiple VM test Perfmon Data with Xeon 6140 host processor

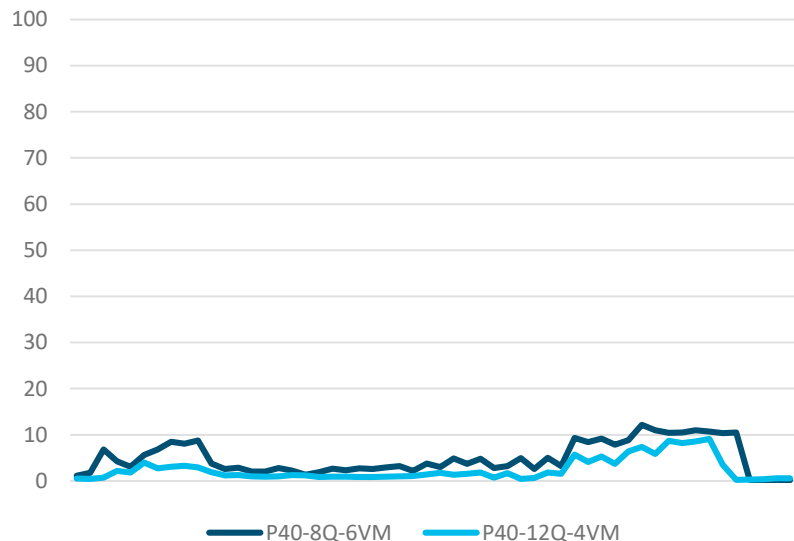# Intel Scalable Family 6128 and NVIDIA Tesla P40*

| Intel Xeon 6128 Utilization | NVIDIA Tesla P40 Utilization |
|---|---|



Left chart legend: P40-8Q-6VM, P40-12Q-4VM
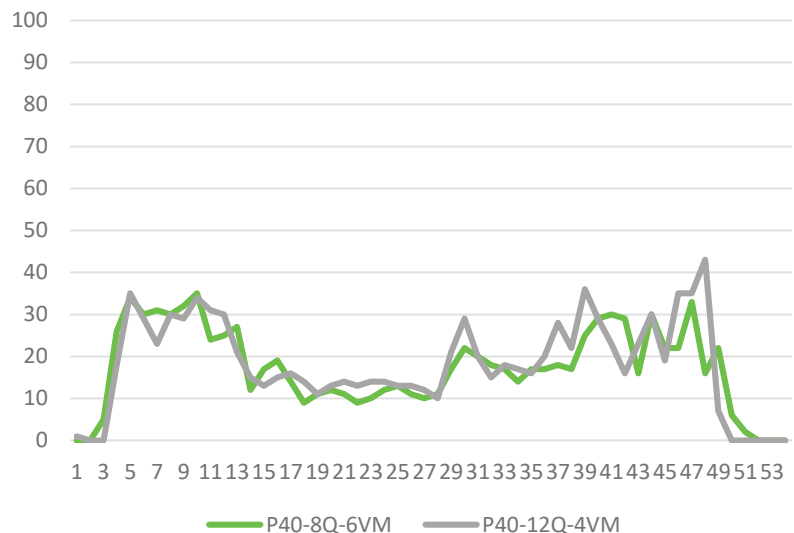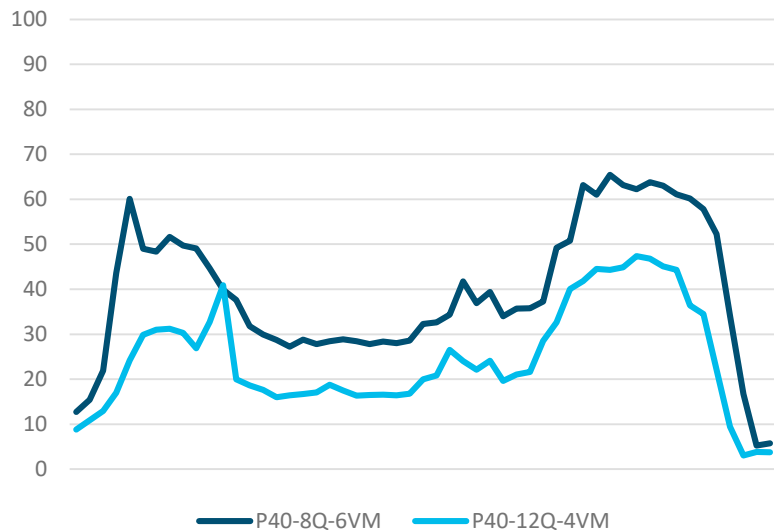Right chart legend: P40-8Q-6VM, P40-12Q-4VM

\* SPECviewperf 13 Catia Test – ESXi Host Data

# Windows 10 (1607) VM with Tesla P40 vGPU*



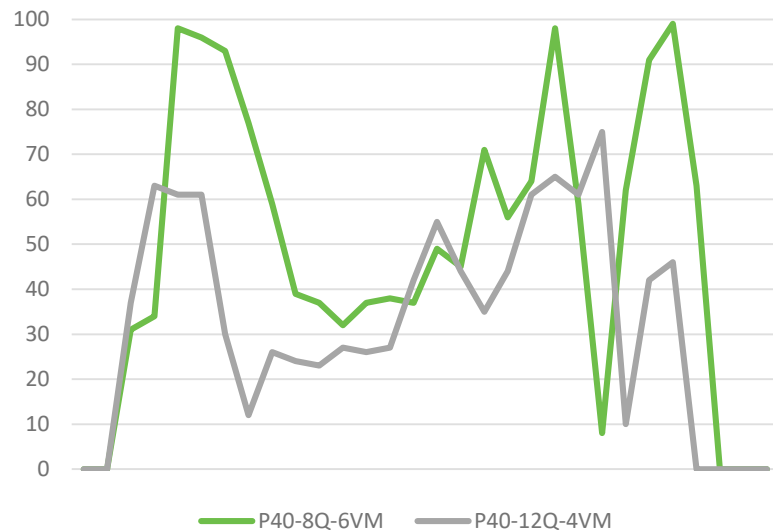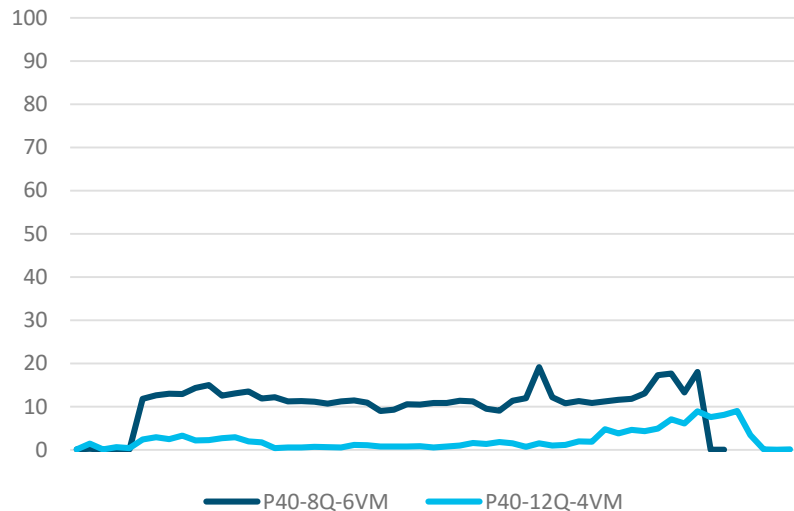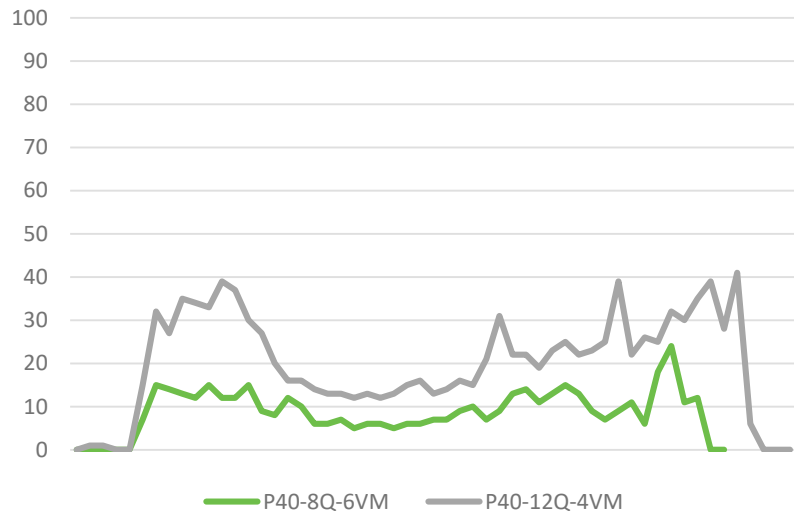Perfmon VM Processor Utilization

Perfmon VM GPU Utilization

P40-8Q-6VM    P40-12Q-4VM

P40-8Q-6VM    P40-12Q-4VM

\* SPECviewperf 13 Catia Test – Single VM in multiple VM test Perfmon Data with Xeon 6128 host processor

# Sizing for Dassault Apps

# Dassault Systemes 3DEXPERIENCE UCS C240 M5 Rack Server starting points*

| User type | Equivalent performance | Users/ server | vCPU/ user | Memory/ user | Server CPU | Server memory | NVIDIA GPU | Quadro profile | Storage type | Network |
|-----------|------------------------|---------------|------------|--------------|------------|---------------|------------|----------------|--------------|---------|
| Light | Quadro P1000 | 32 | 4 | 12-16 | Intel Xeon 6136 | 768 | Tesla T4 (4) | T4-2Q | Flash | 10Gb+ |
| Medium | Quadro P2000 | 16 | 4-6 | 16-32 | Intel Xeon 6134 | 768 | Tesla T4 (4) | T4-4Q | Flash | 10Gb+ |
| Heavy | Quadro P5000 | 4-6 | 8-12 | 96+ | Intel Xeon 6128 | 768 | Tesla P40 (2) | P40-8Q P40-12Q | Flash | 10Gb+ |

*The recommendations above reflect starting points. Customers should perform PoCs to determine optimal configurations for their specific environments. Cisco can help.

# Dassault Systemes 3DEXPERIENCE UCS B200 M5 Blade Server Rack Dense starting points*

| User type | Equivalent performance | Users/ server | vCPU/ user | Memory/ user | Server CPU | Server memory | NVIDIA GPU | Quadro profile | Storage type | Network |
|---|---|---|---|---|---|---|---|---|---|---|
| Light | Quadro P1000 | 12 | 4 | 12-16 | Intel Xeon 6128 | 192 | Tesla P6 (2) | P6-2Q | Flash | 10Gb+ |
| Medium | Quadro P2000 | 6 | 4-6 | 16-32 | Intel Xeon 6128 | 192 | Tesla P6 (2) | P6-4Q | Flash | 10Gb+ |
| Heavy | Quadro P5000 | 2-4 | 8-12 | 96+ | Intel Xeon 6128 | 192 | Tesla P6 (2) | P6-8Q P6-16Q | Flash | 10Gb+ |

*The recommendations above reflect starting points. Customers should perform PoCs to determine optimal configurations for their specific environments. Cisco can help.

# The Cisco Lineup

# Cisco graphics accelerated Data Center with NVIDIA

## Racks

### C220 M5



2x NVIDIA P4, T4

### C240 M5



2x NVIDIA V100,P100,P40
2x M10,M60,
6x P4, 4x T4

### C480 M5



6x NVIDIA V100,P100, P40,
M60

3x M10

### C480 ML M5



8x NVIDIA V100 32GB
NvLINK Interconnect

## Blades

### B200 M5



2x NVIDIA P6 GPU/blade
Up to 16x per chassis

### B480 M5



4x P6 GPU/blade,
Up to 16x per chassis

## Hyperconverged

### HyperFlex 240C M5



2x NVIDIA V100, P40

2x M10, M60

6x P4

# Key takeaways

# Keep these things in mind
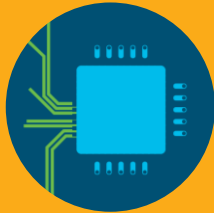
## Understanding the different types of users

## There are three key GPU settings:
- GPU scheduler
- NVIDIA profile selection
- Frame rate control
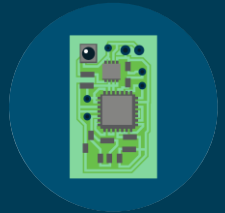  - NVIDIA Tesla card
  - Desktop Broker

## CPU selection is critical
- CPU and GPU work synergistically
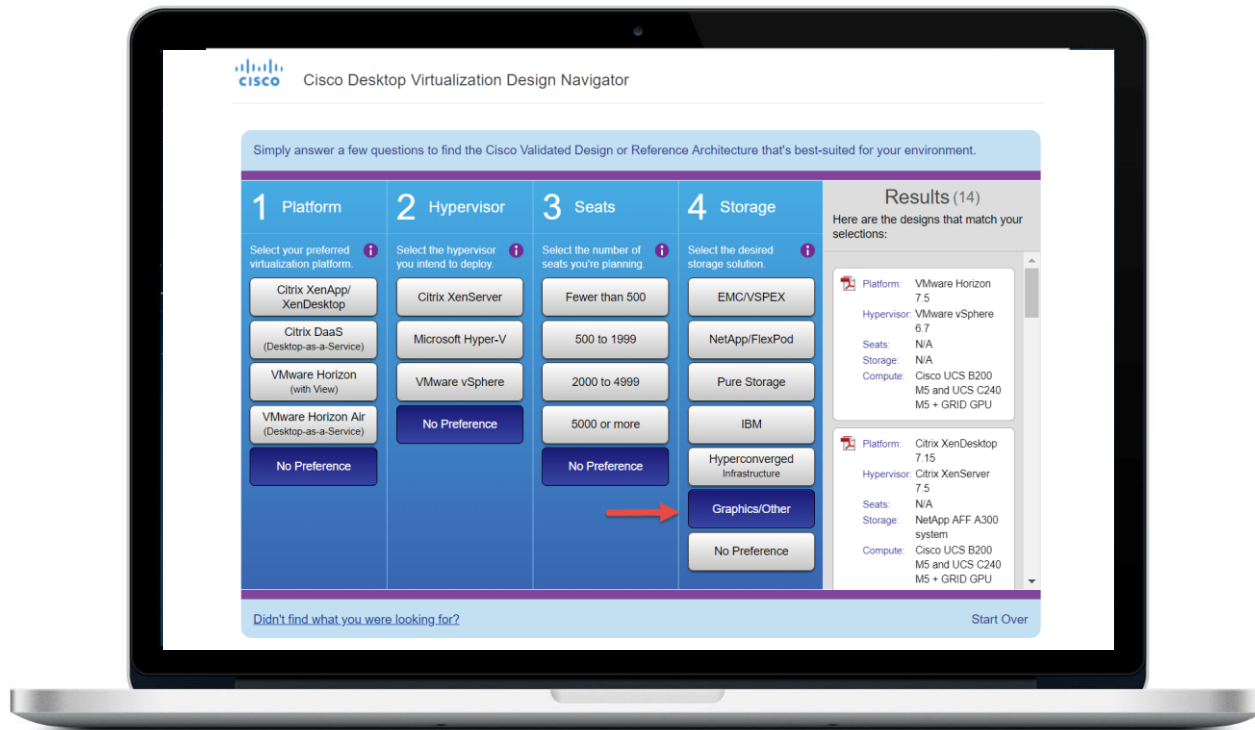- High frequency
- CPU core count

## High frequency memory

# Resources

# VCC Design Navigator
## Your source for VDI content

Q & A