# vMotion for NVIDIA GRID vGPU Virtual Machines: Case Study of vMotion Using MLaaS

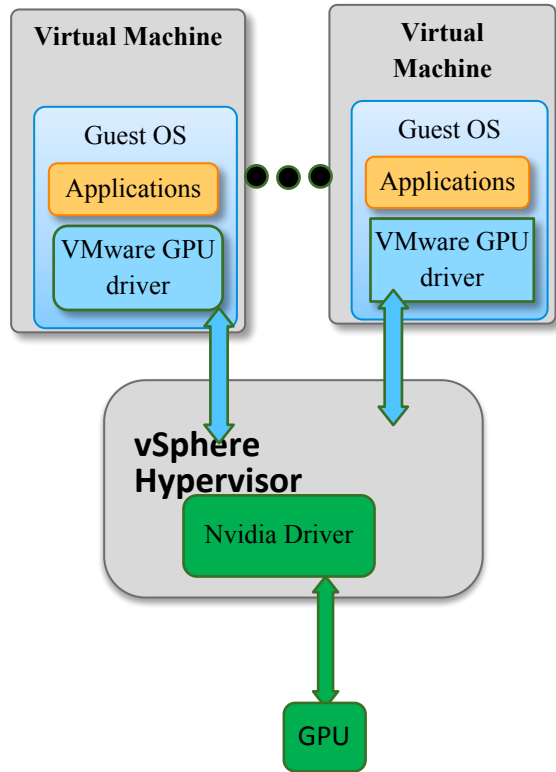Hari Sivaraman, Dimitrios Skarlatos
Lan Vu, Uday Kurkure

**GTC 2019**
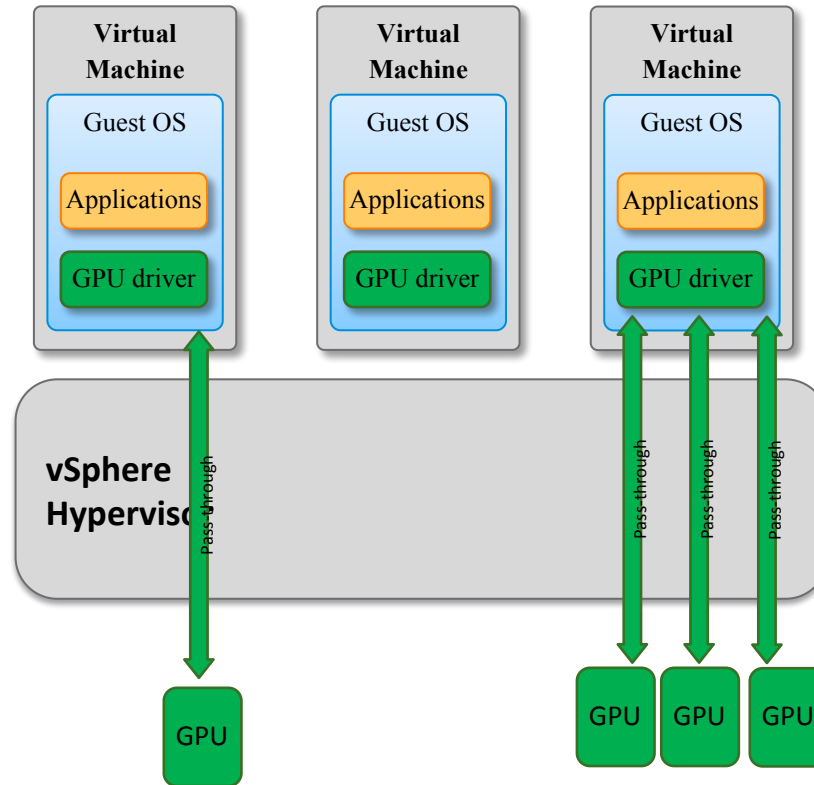
# vMotion for NVIDIA GRID vGPU - Agenda

- GPUs in vSphere.

- vMotion for vGPU Architecture.

- Performance of vMotion for vGPU.

- MLaaS – a case study for vMotion performance.

- Conclusions and future work.
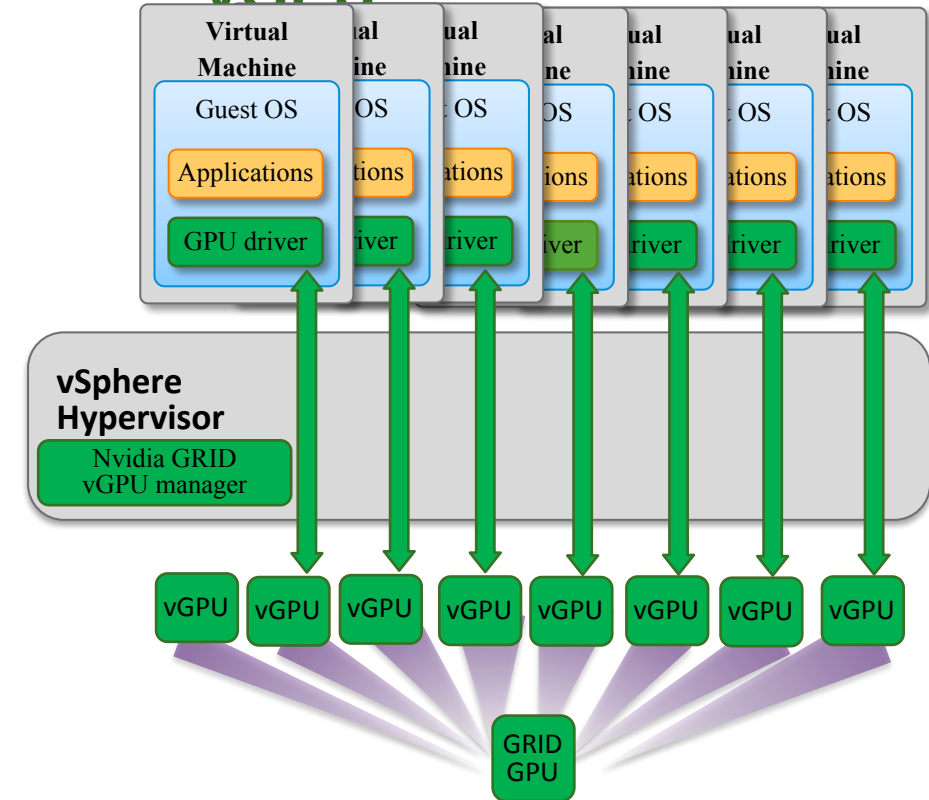
# vMotion for NVIDIA GRID vGPU – GPUs in vSphere



## vSGA

**Virtual Machine**

Guest OS
- Applications
- VMware GPU driver

**Virtual Machine**

Guest OS
- Applications
- VMware GPU driver

**vSphere Hypervisor**
- Nvidia Driver

GPU

☑ ☑ vMotion
☑ Sharing

## VMware DirectPath I/O

**Virtual Machine**

Guest OS
- Applications
- GPU driver

**Virtual Machine**

Guest OS
- Applications
- GPU driver

**Virtual Machine**

Guest OS
- Applications
- GPU driver

**vSphere Hypervisor**

Pass-through

GPU

Pass-through | Pass-through | Pass-through

GPU | GPU | GPU

✗ vMotion
✗ Sharing

## Nvidia GRID vGPU

**Virtual Machine**

Guest OS
- Applications
- GPU driver

Virtual Machine | Virtual Machine | Virtual Machine | Virtual Machine | Virtual Machine | Virtual Machine

OS | OS | OS | OS | OS | OS
- ations | ations | ations | ations | ations | ations
- river | river | river | river | river | river

**vSphere Hypervisor**
- Nvidia GRID vGPU manager

vGPU | vGPU | vGPU | vGPU | vGPU | vGPU | vGPU | vGPU

GRID GPU

☑ vMotion
☑ Sharing
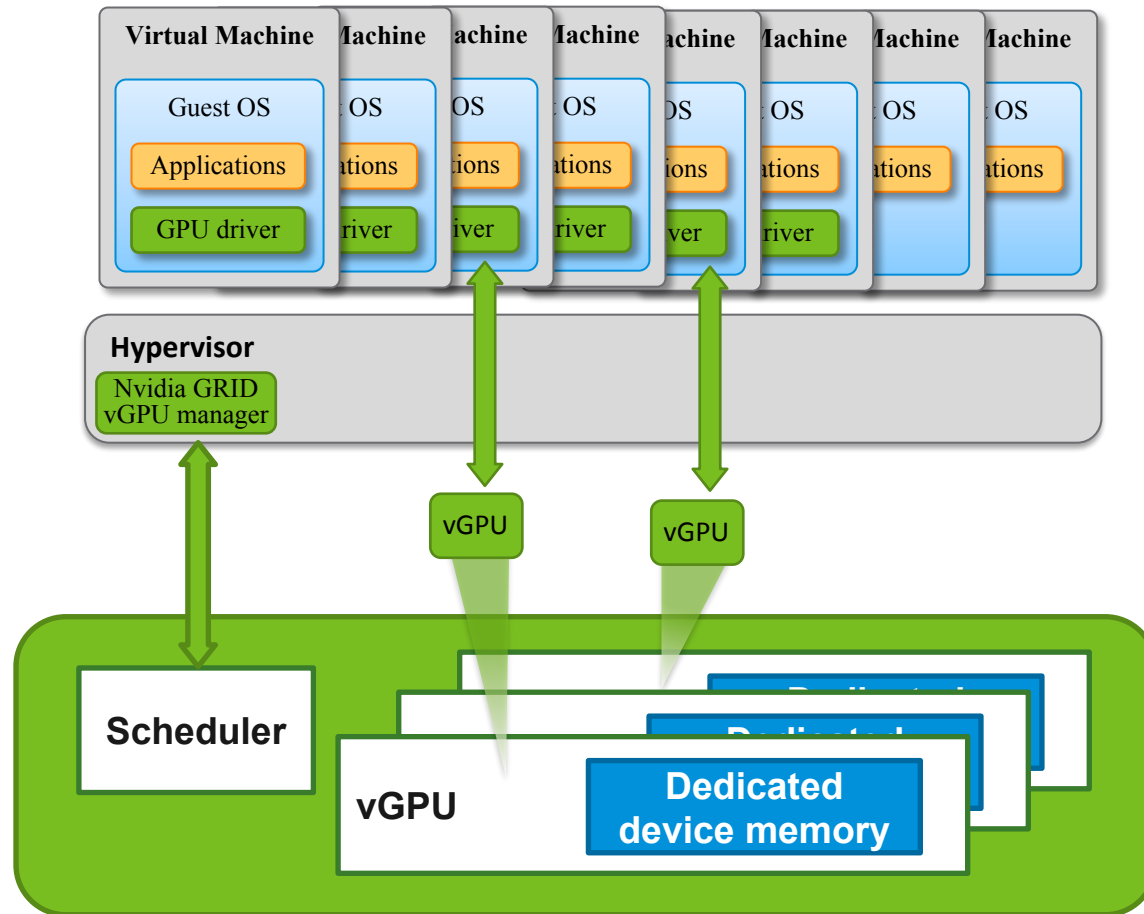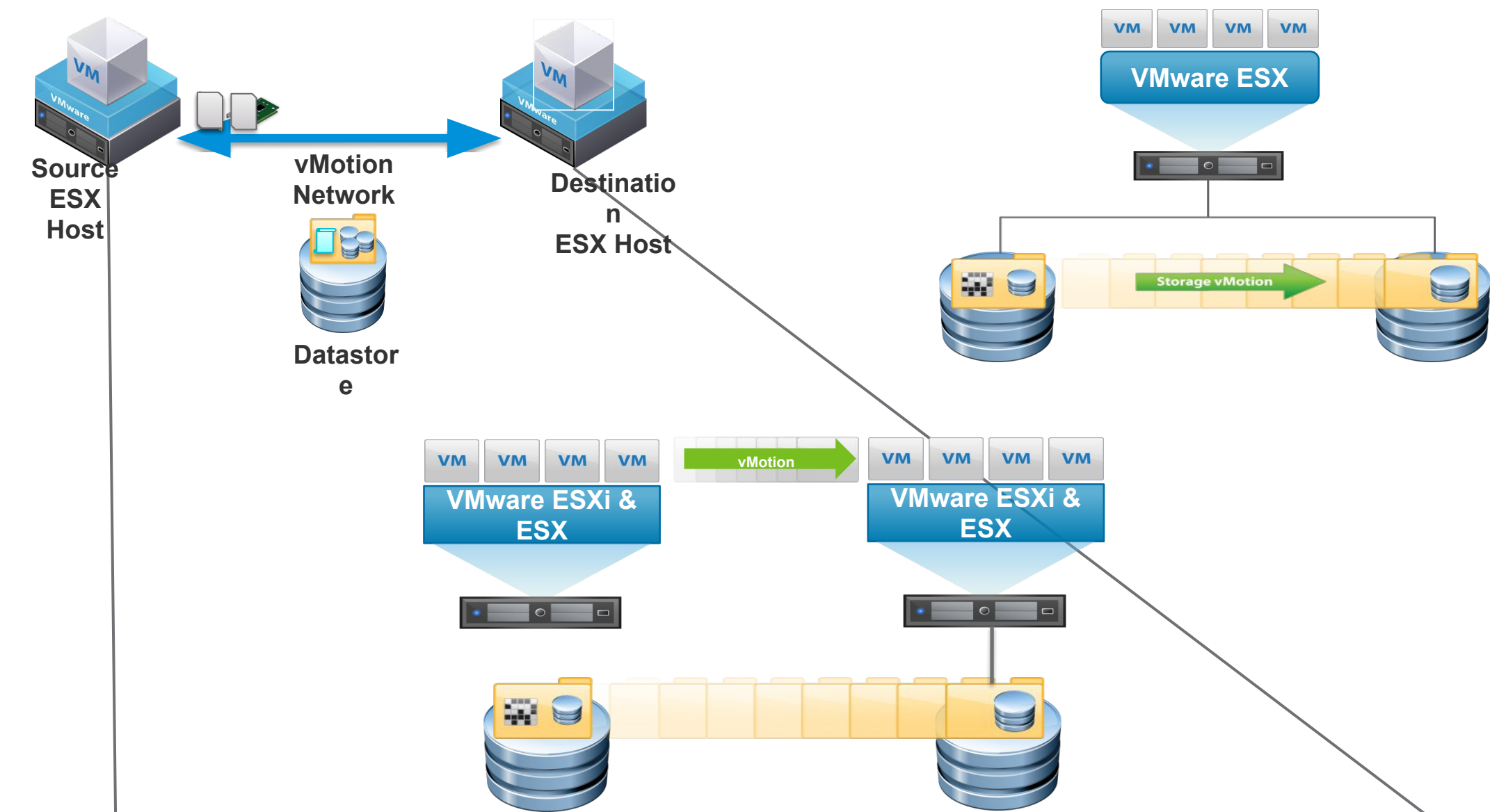
# vMotion for NVIDIA GRID vGPU – vGPU
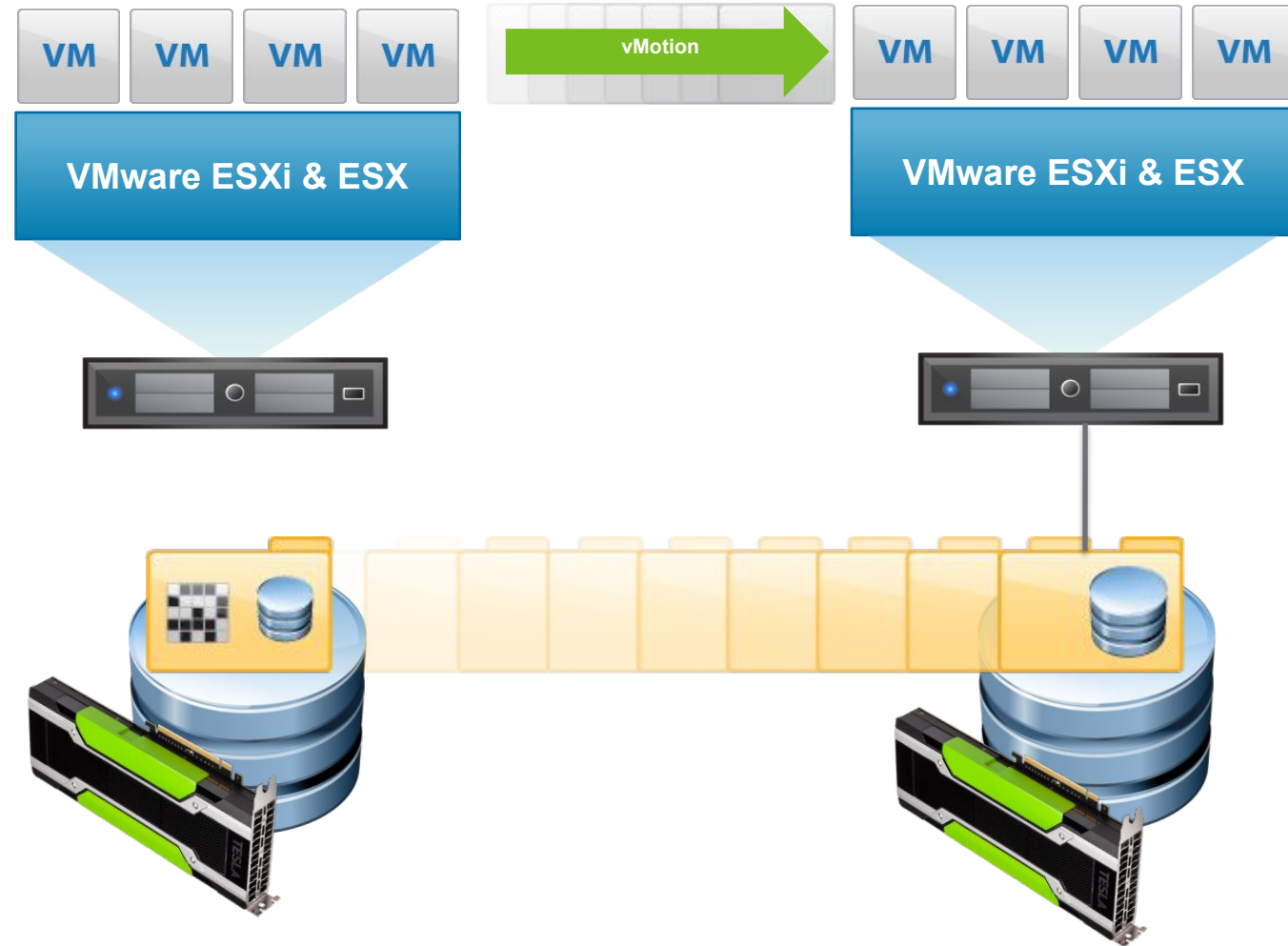
## Nvidia GRID vGPU



- GPU Memory is statically shared

- CUDA cores are time-shared

- GPU memory per VM is called vGPU Profile

- **For example**: P40-1q profile for P40 GPU
  - vGPU has 1GB of device memory
  - **24** vGPUs per **1** physical P40

# vMotion for NVIDIA GRID vGPU – Types of vMotion
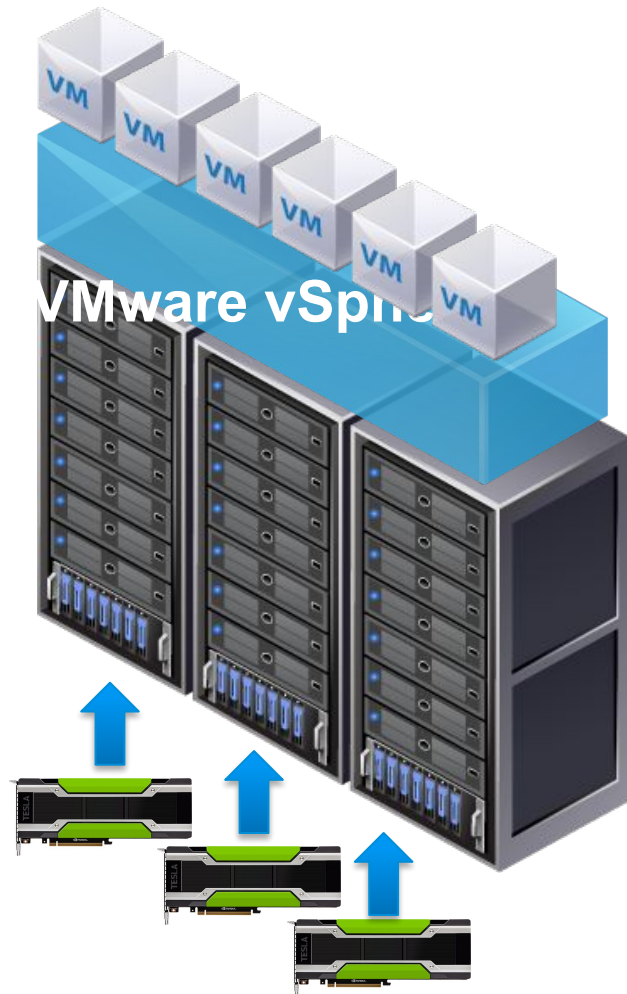
# vMotion for NVIDIA GRID vGPU – vMotion



**1** pre-copy memory pages

**2** Stun the VM

**3** Checkpoint devices

**4** Xfer device checkpoint data (includes vGPU memory data)

**5** Power on VM & xfer pages from main memory

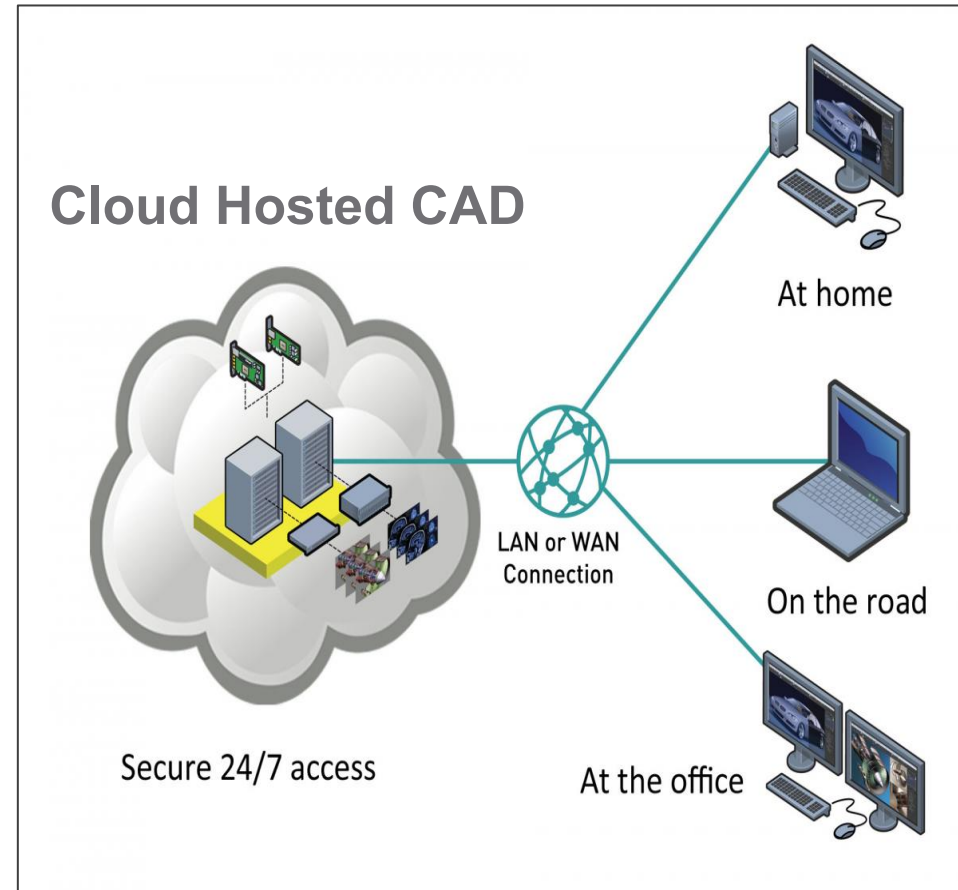# vMotion for NVIDIA GRID vGPU - Agenda

- GPUs in vSphere.

- vMotion for vGPU Architecture.

- **Performance of vMotion for vGPU.**

- **MLaaS – a case study for vMotion performance.**

- **Conclusions and future work.**

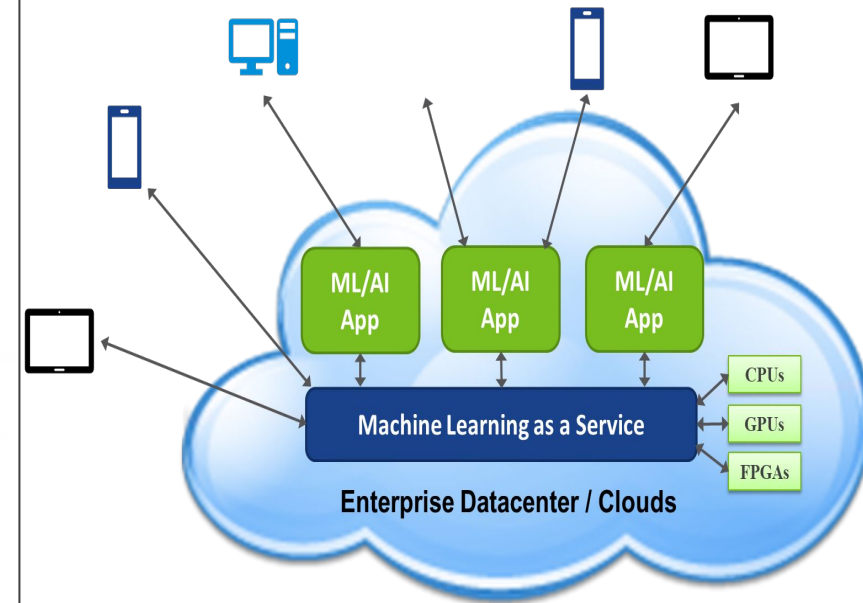# vMotion for NVIDIA GRID vGPU - Workloads



VMware vSphere

Cloud Hosted CAD

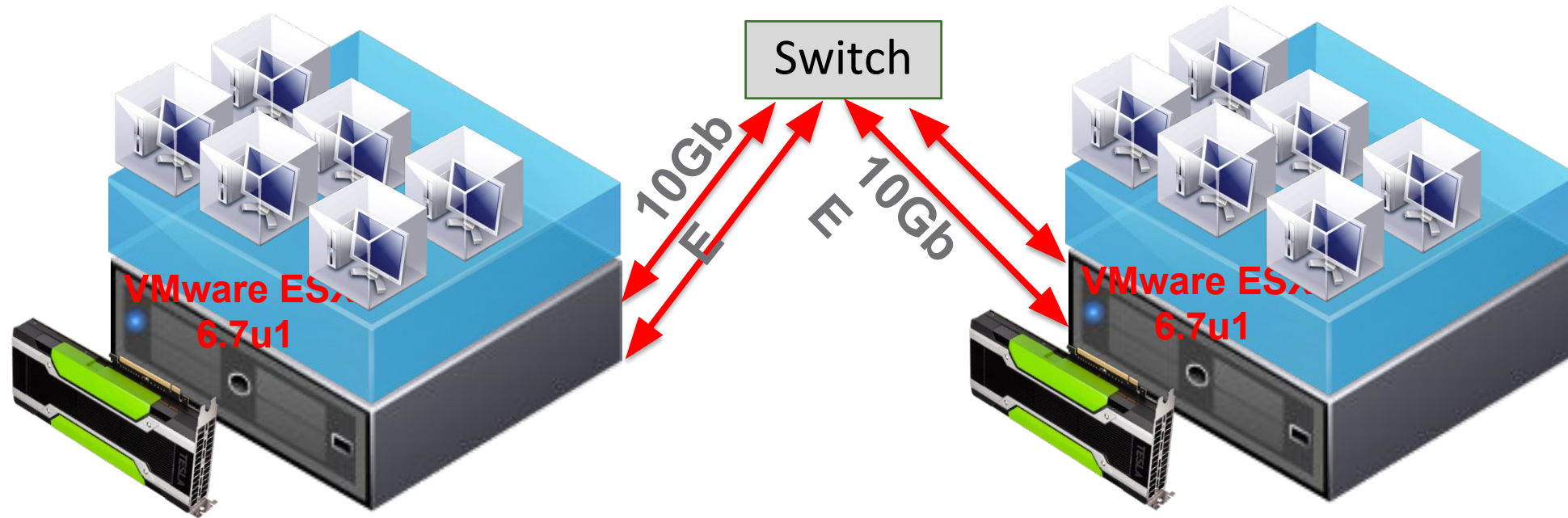At home

LAN or WAN
Connection

On the road

Secure 24/7 access

At the office

ML/AI App

ML/AI App

ML/AI App

Machine Learning as a Service

CPUs

GPUs

FPGAs

Enterprise Datacenter / Clouds

**VDI**

**Cloud Hosted CAD**

**MLaaS**

# vMotion for NVIDIA GRID vGPU – Test-bed



Switch

10Gb E

10Gb E

VMware ESX 6.7u1

VMware ESX 6.7u1

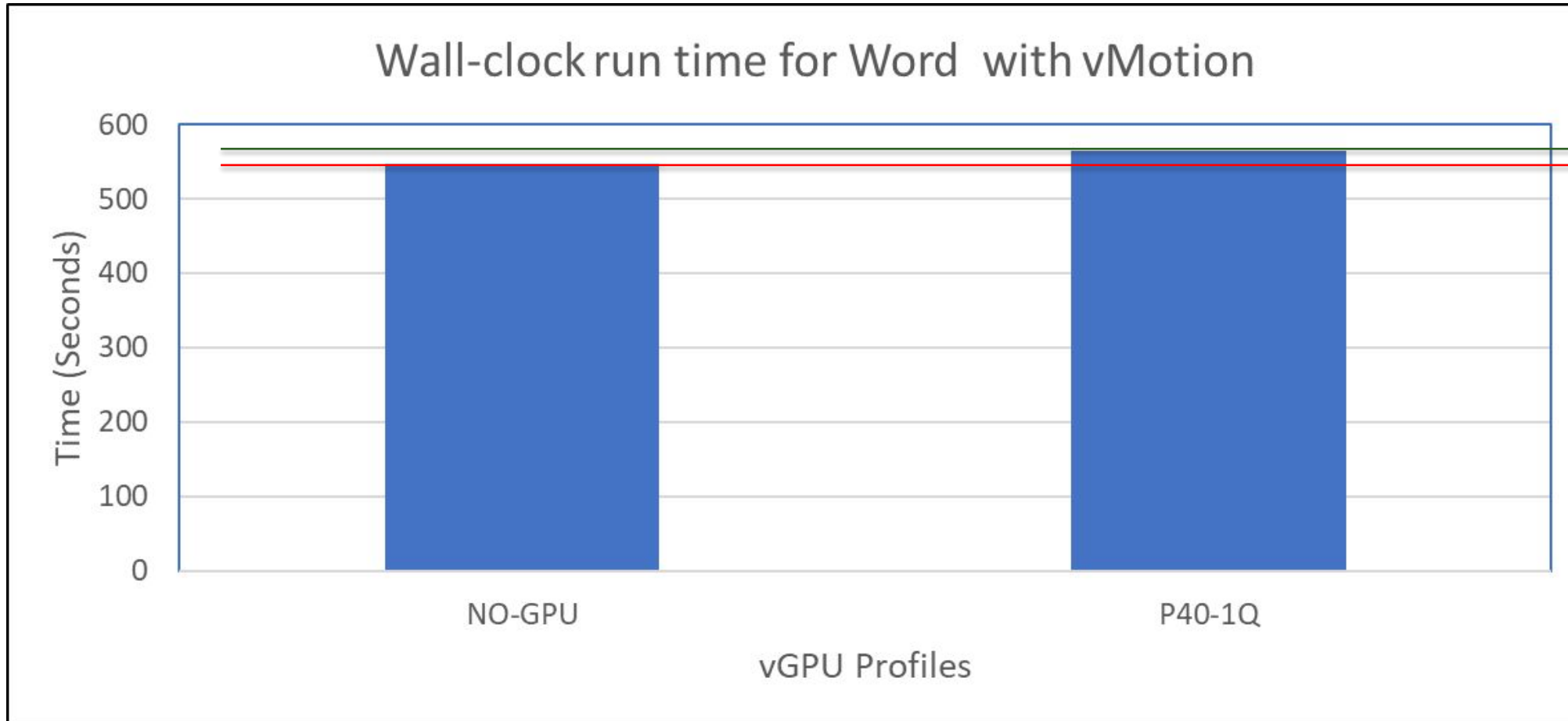**Dell R730 – Intel Broadwell CPUs + 1 x NVidia GRID P40**
40 cores (2 x 20-core socket) E5-2698 v4
768 GB RAM

**Dell R730 – Intel Broadwell CPUs + 1 x NVidia GRID P40**
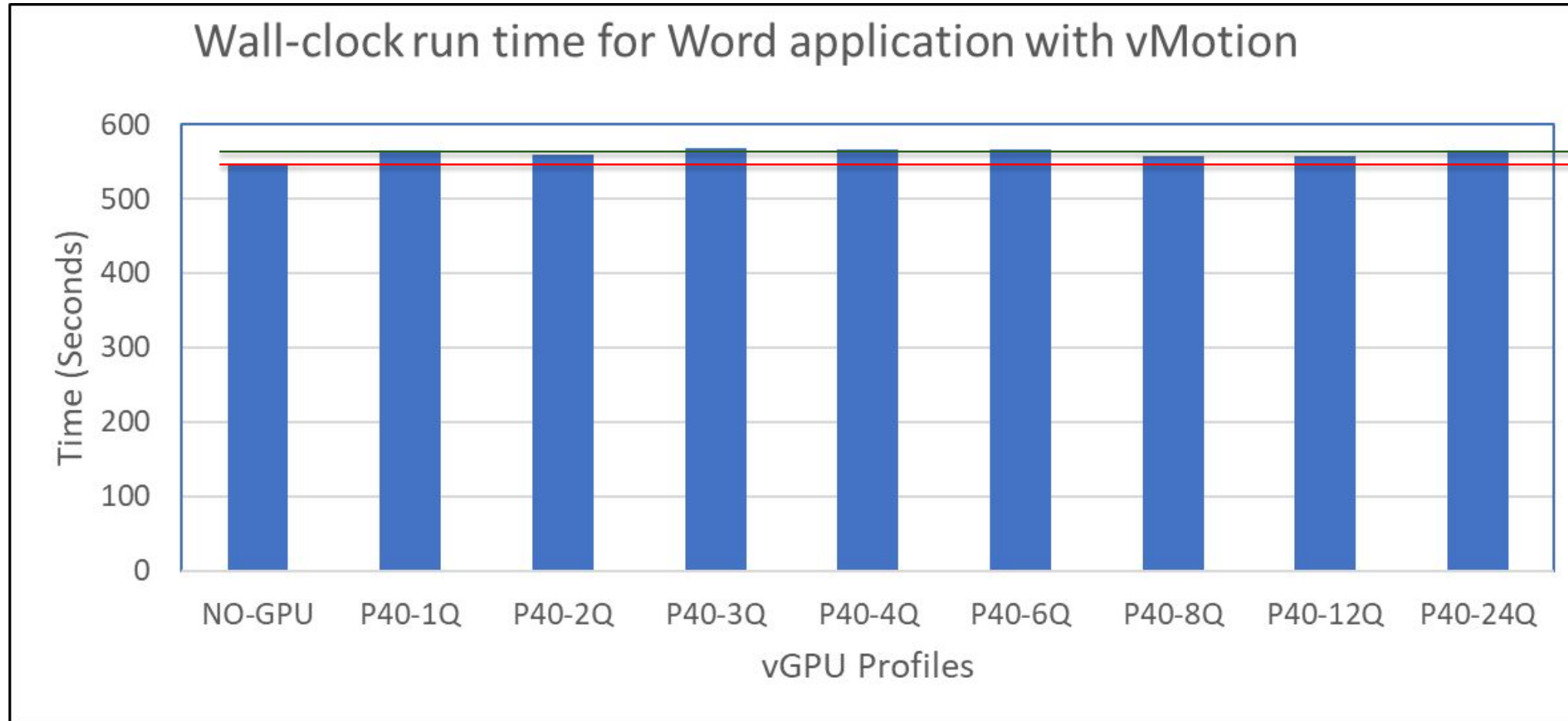40 cores (2 x 20-core socket) E5-2698 v4
768 GB RAM

- **ESX**: 6.7u1   **Nvidia Driver**:  410.68
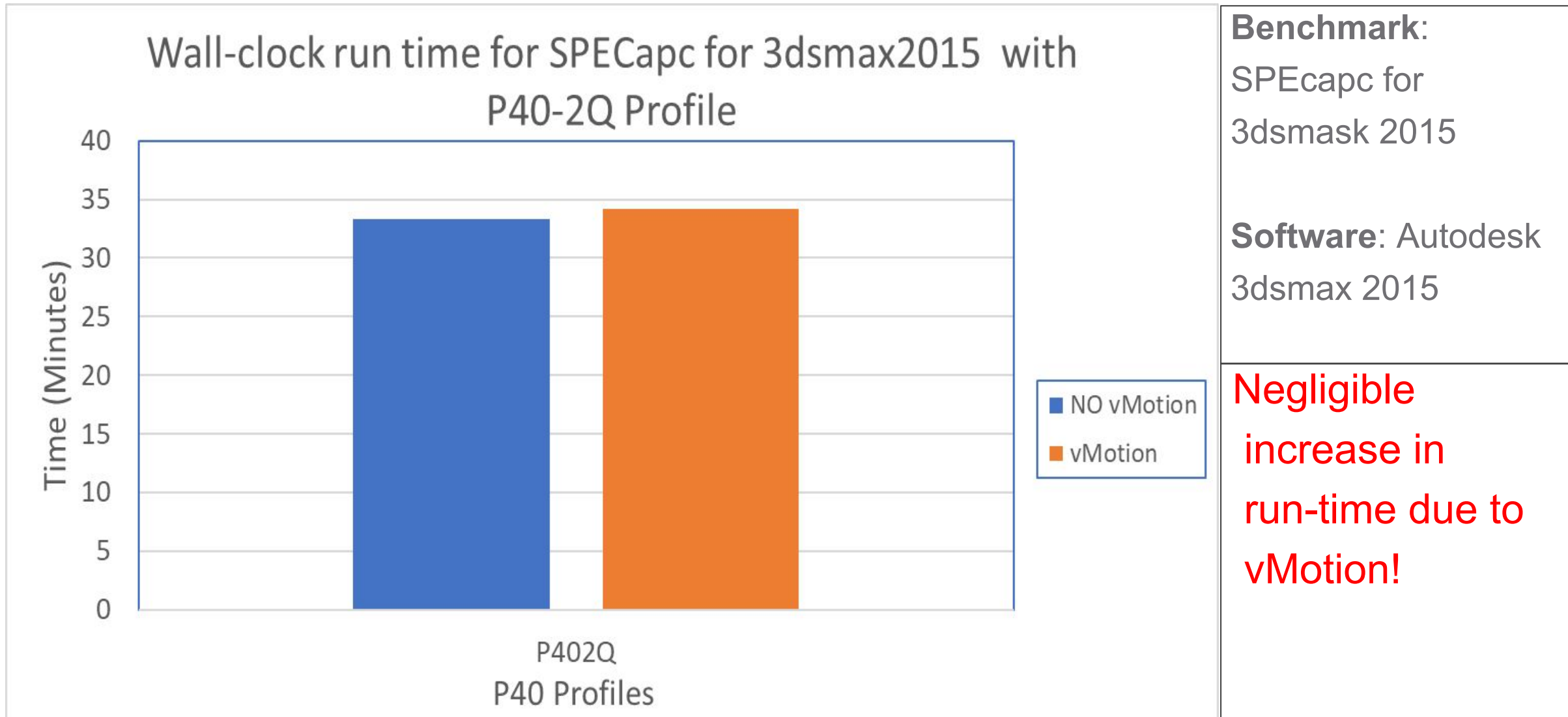
# vMotion for NVIDIA GRID vGPU – Performance of Word



Wall-clock run time for Word with vMotion

Increase in vMotion time due to vGPU is just marginally more than measurement noise.

# vMotion for NVIDIA GRID vGPU – Performance of Word



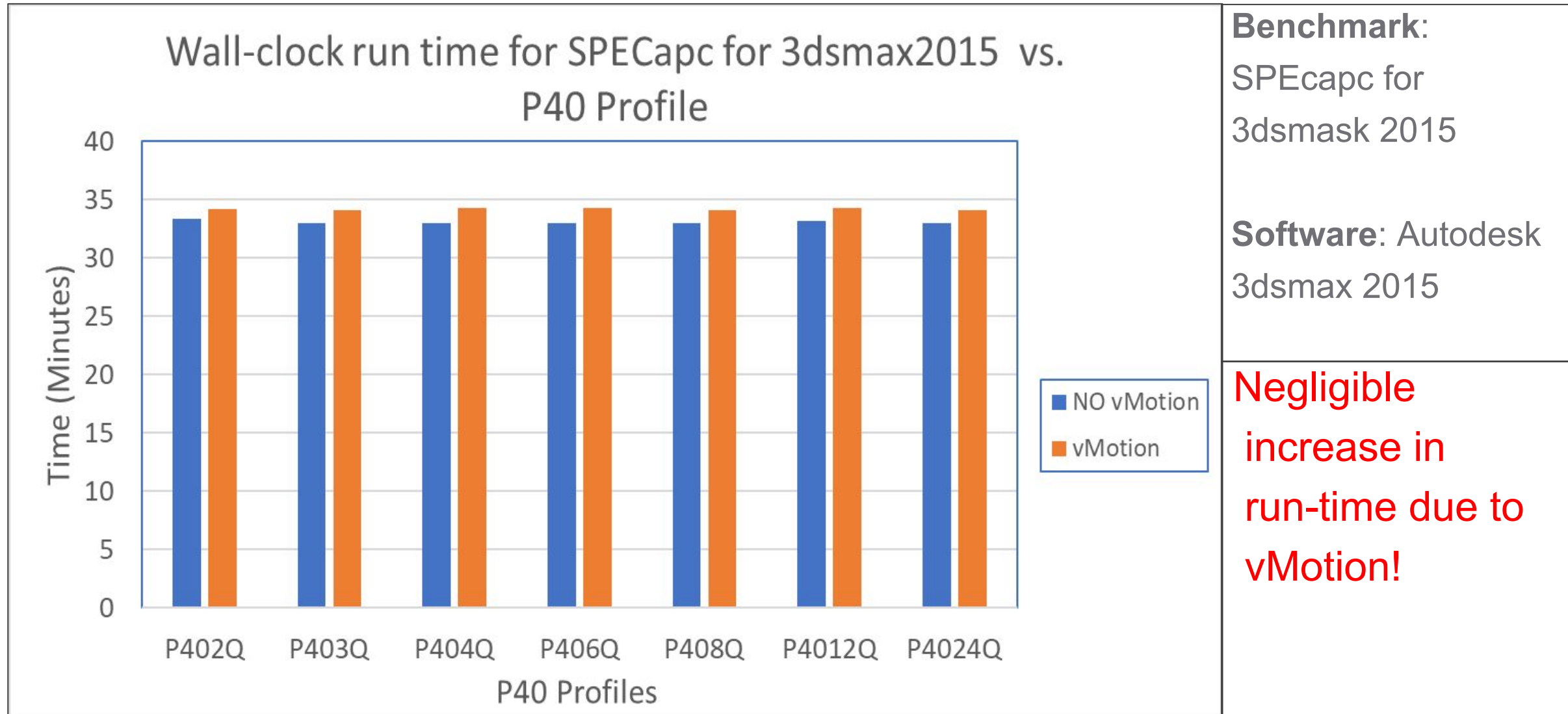Wall-clock run time for Word application with vMotion

Increase in vMotion time due to vGPU is just marginally more than measurement noise.

# vMotion for NVIDIA GRID vGPU – Performance of SPECapc for 3dsmax 2015



Wall-clock run time for SPECapc for 3dsmax2015 with P40-2Q Profile

Time (Minutes) axis: 0, 5, 10, 15, 20, 25, 30, 35, 40

Legend: ■ NO vMotion ■ vMotion

P402Q
P40 Profiles

**Benchmark**: SPEcapc for 3dsmask 2015

**Software**: Autodesk 3dsmax 2015

Negligible increase in run-time due to vMotion!

# vMotion for NVIDIA GRID vGPU – Performance of SPECapc for 3dsmax 2015



**Benchmark**: SPEcapc for 3dsmask 2015

**Software**: Autodesk 3dsmax 2015

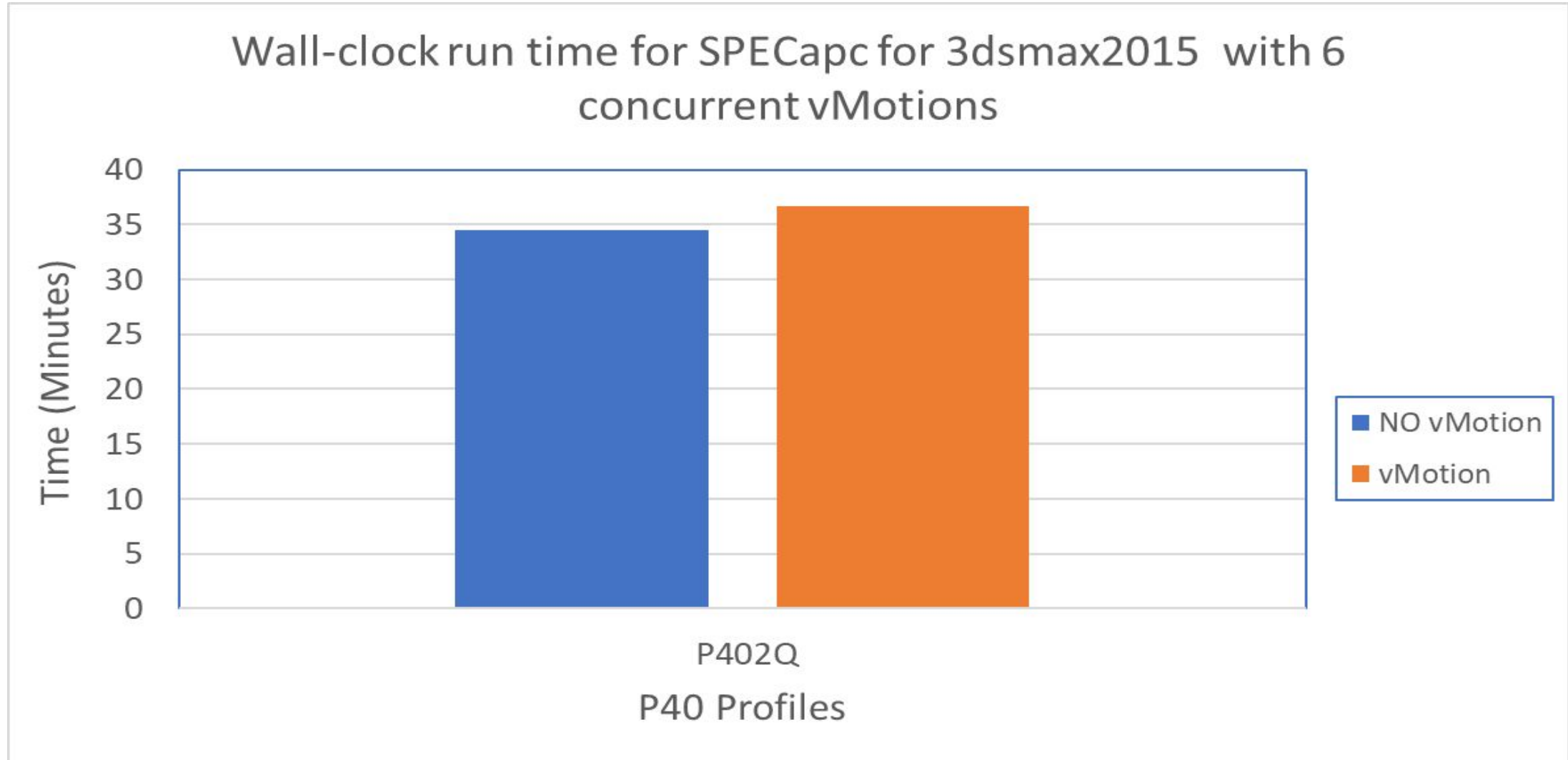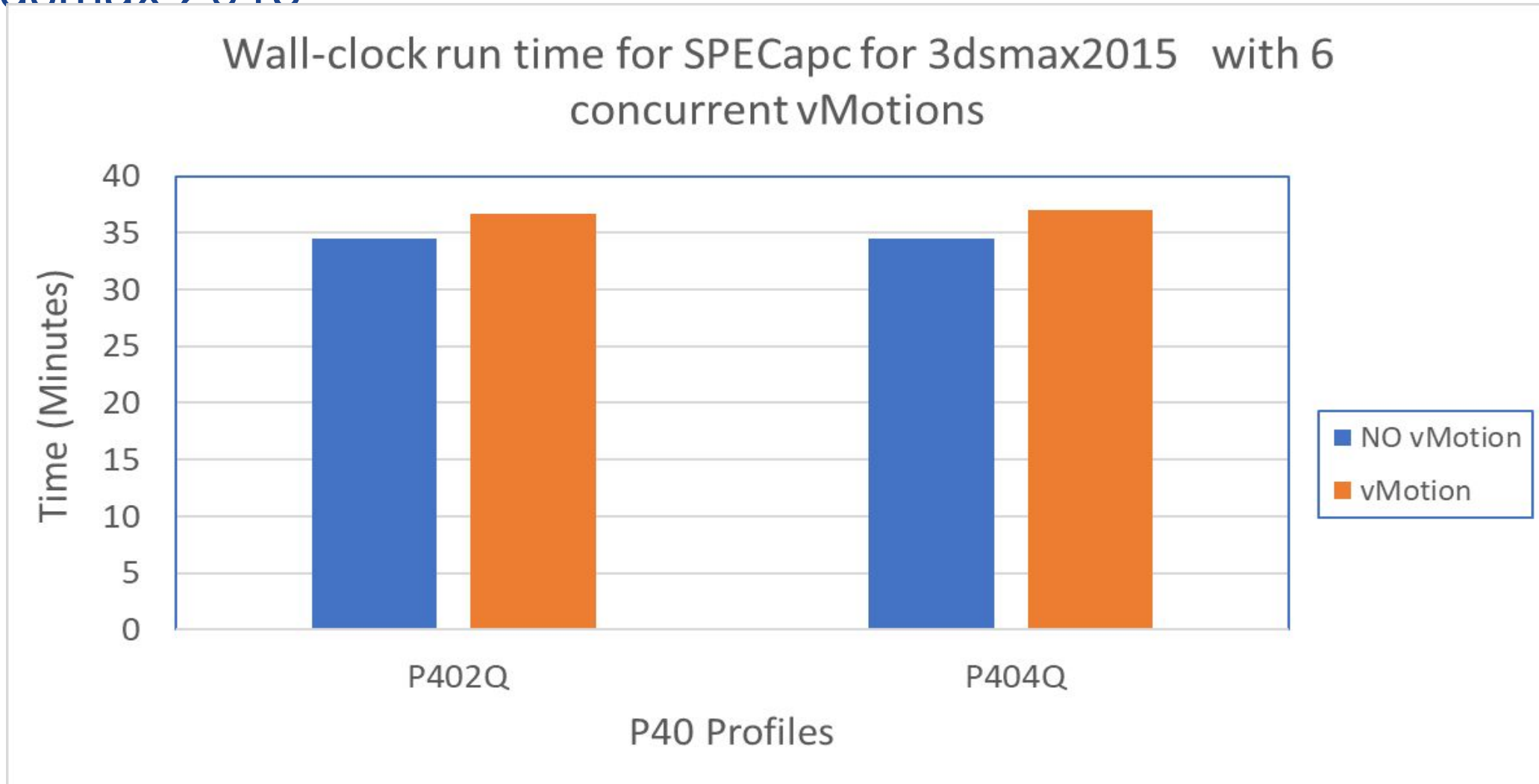Negligible increase in run-time due to vMotion!

# vMotion for NVIDIA GRID vGPU – Performance of SPECapc for 3dsmax 2015

**vm**ware®

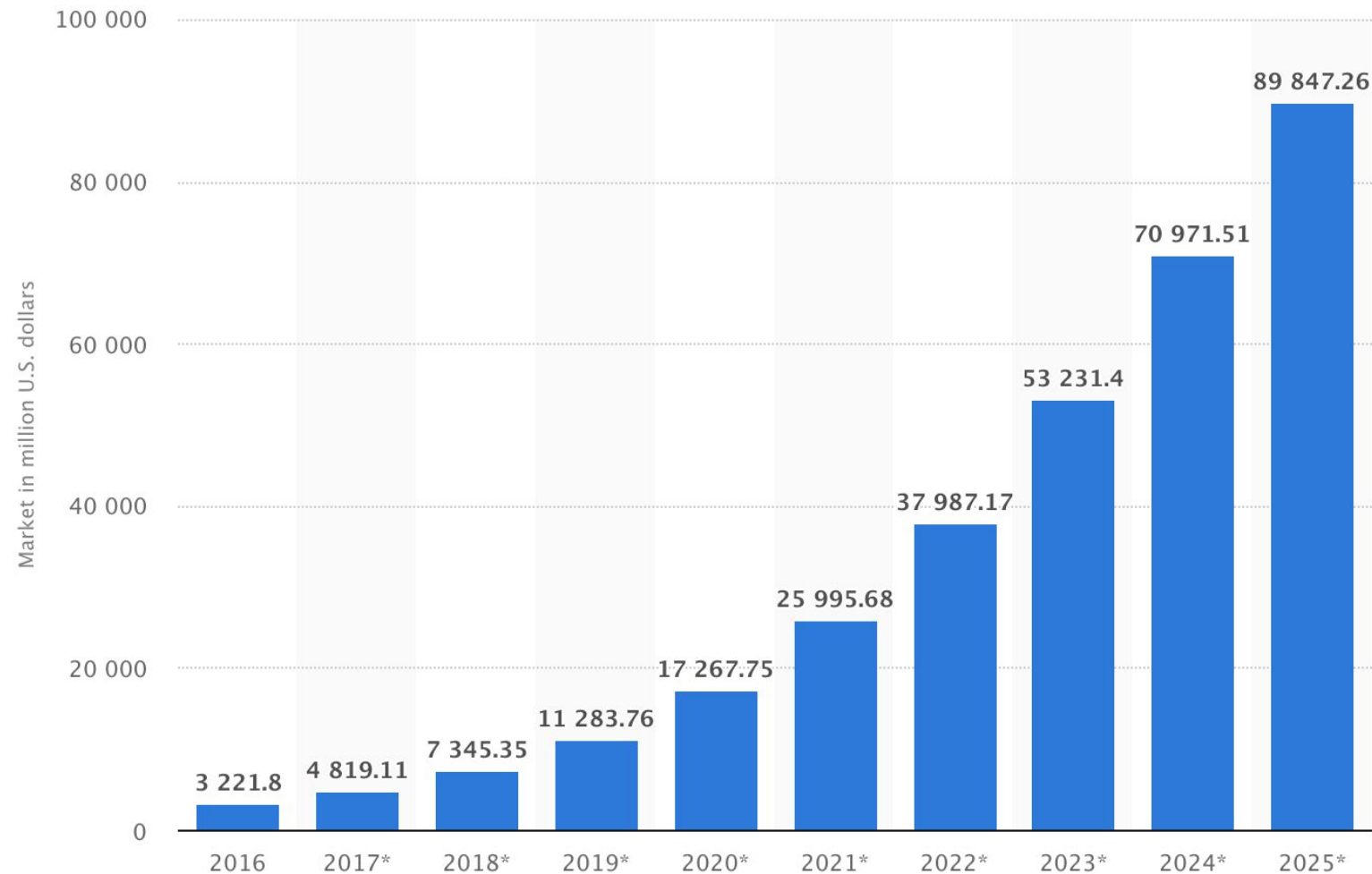# vMotion for NVIDIA GRID vGPU – Performance of SPECapc for 3dsmax 2015

# vMotion for NVIDIA GRID vGPU - Agenda

- GPUs in vSphere.

- vMotion for vGPU Architecture.

- Performance of vMotion for vGPU.

- **MLaaS – a case study for vMotion performance.**
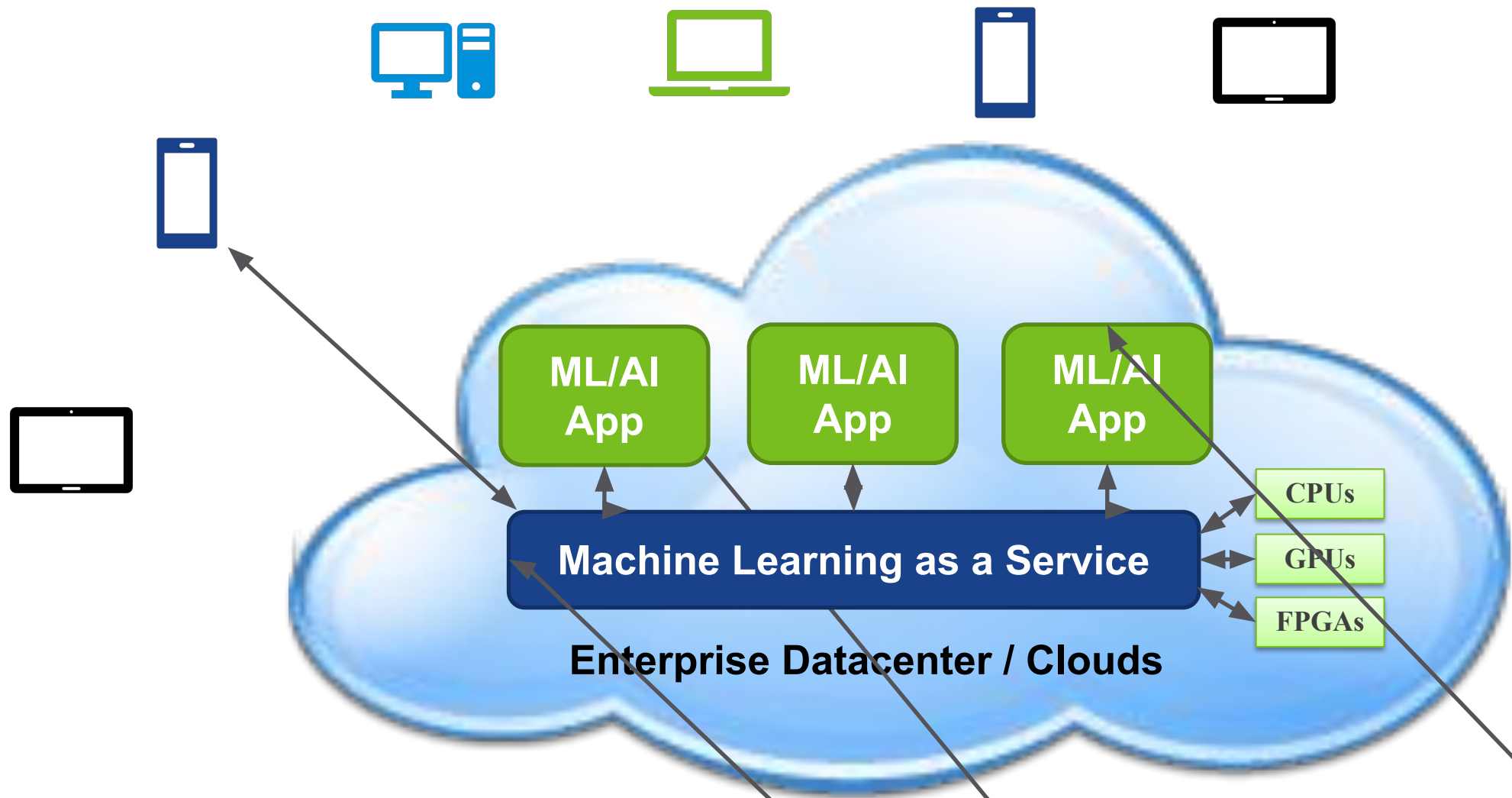
- **Conclusions and future work.**

# Revenues from the Artificial Intelligence (AI) market worldwide from 2016 to 2025
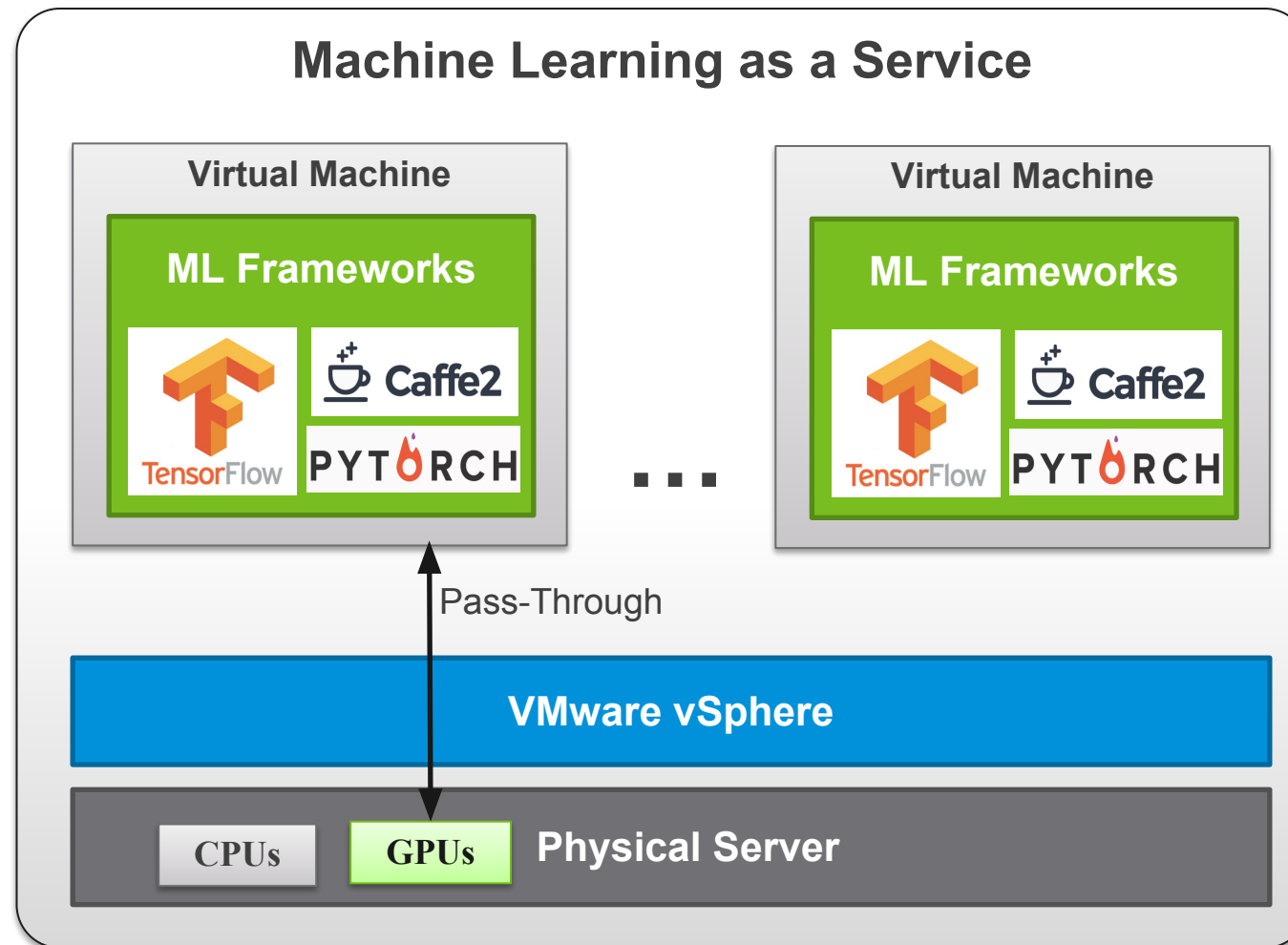


**The largest proportion of revenues come from the ML/AI Enterprise Applications**
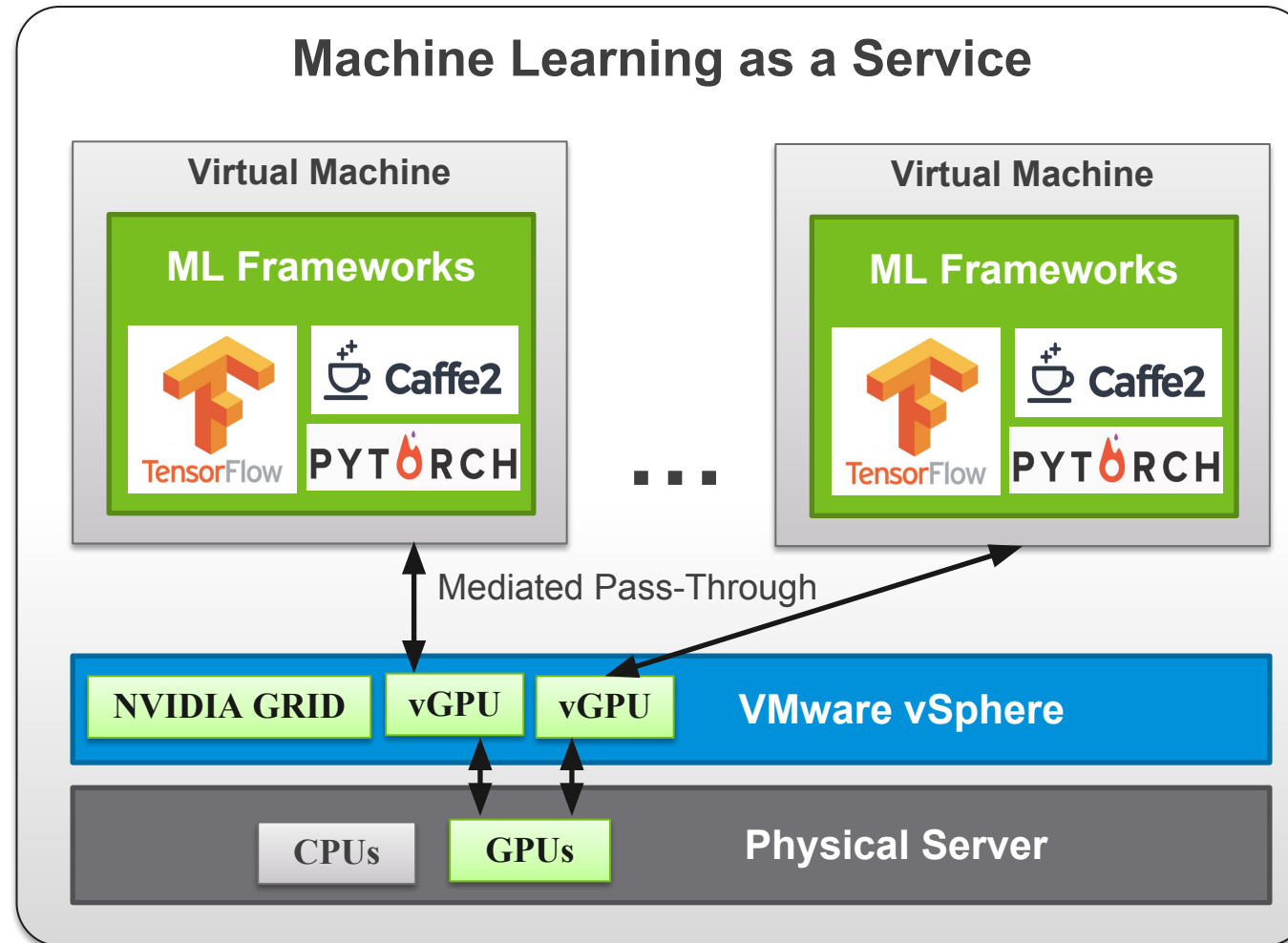
© Statista 2019

©2019 VMware, Inc.

# ML/AI Enterprise Application Deployment



**ML/AI App**

**ML/AI App**

**ML/AI App**

**Machine Learning as a Service**

CPUs

GPUs

FPGAs

**Enterprise Datacenter / Clouds**

# Example #1 of deploying MLaaS on VMware vSphere



**Machine Learning as a Service**

**Virtual Machine**
**ML Frameworks**
TensorFlow    Caffe2    PYTORCH

**Virtual Machine**
**ML Frameworks**
TensorFlow    Caffe2    PYTORCH

. . .

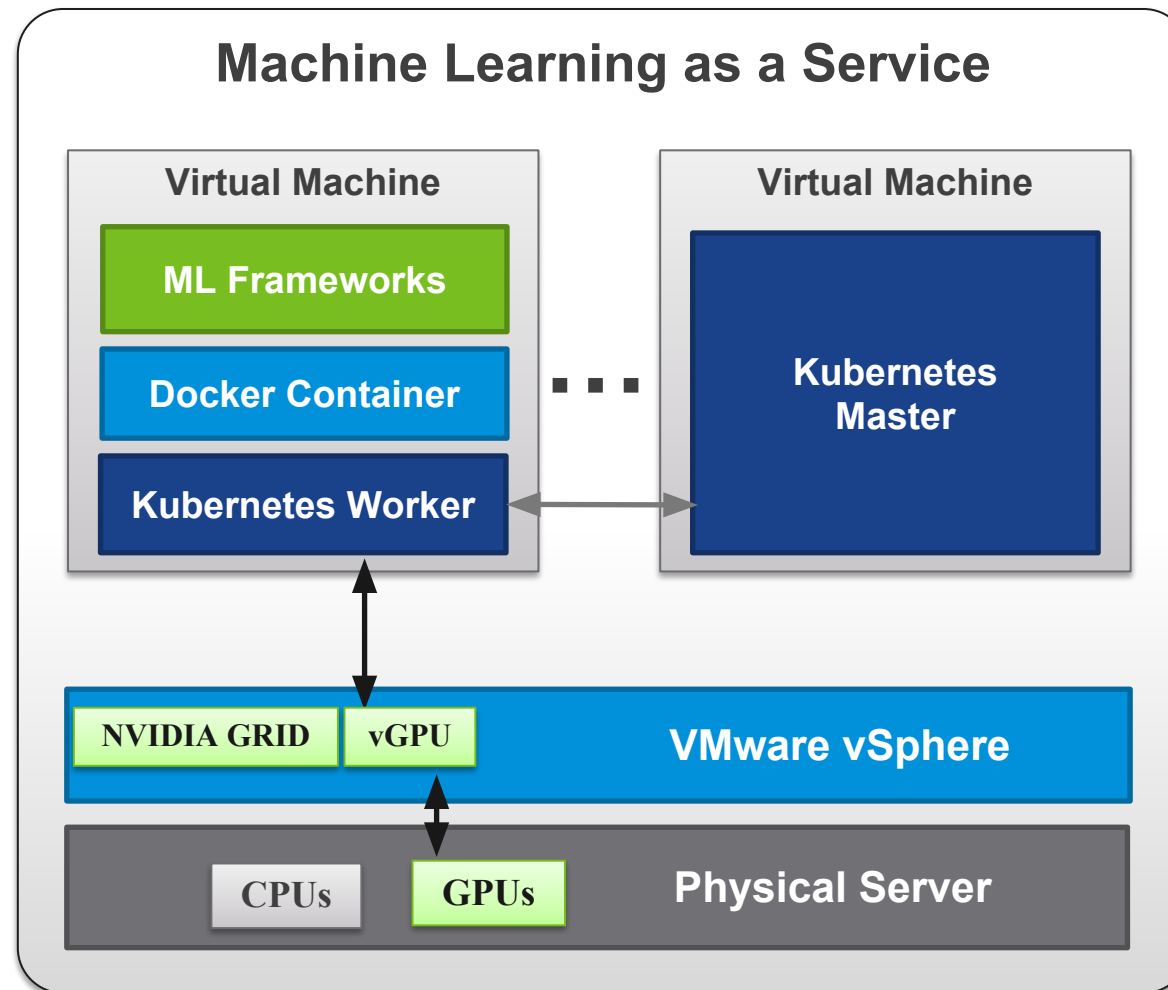Pass-Through

**VMware vSphere**

CPUs    GPUs    **Physical Server**

# Example #2 of deploying MLaaS on VMware vSphere

# Example #3 of deploying MLaaS on VMware vSphere with Container



**Machine Learning as a Service**

**Virtual Machine**
- ML Frameworks
- Docker Container

**Virtual Machine**
- ML Frameworks
- Docker Container

...

NVIDIA GRID | vGPU | vGPU | **VMware vSphere**
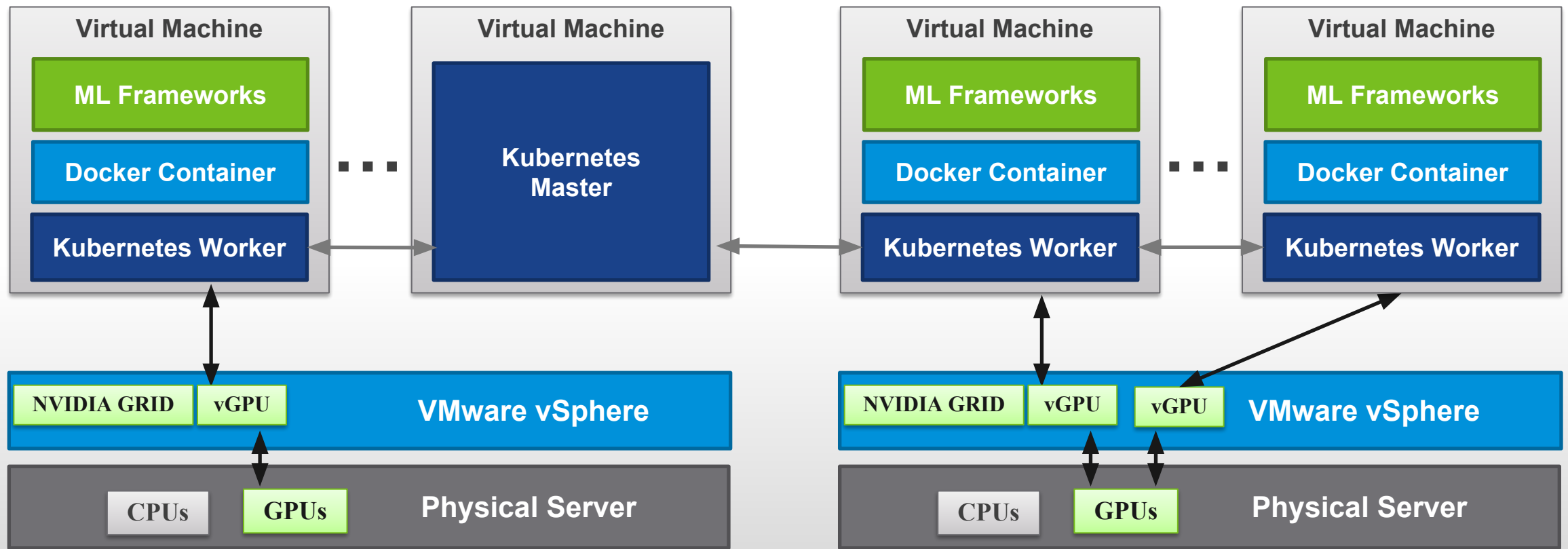
CPUs | GPUs | **Physical Server**

**vm**ware®

# Example #4 of deploying MLaaS on VMware vSphere with Container & Kubernetes

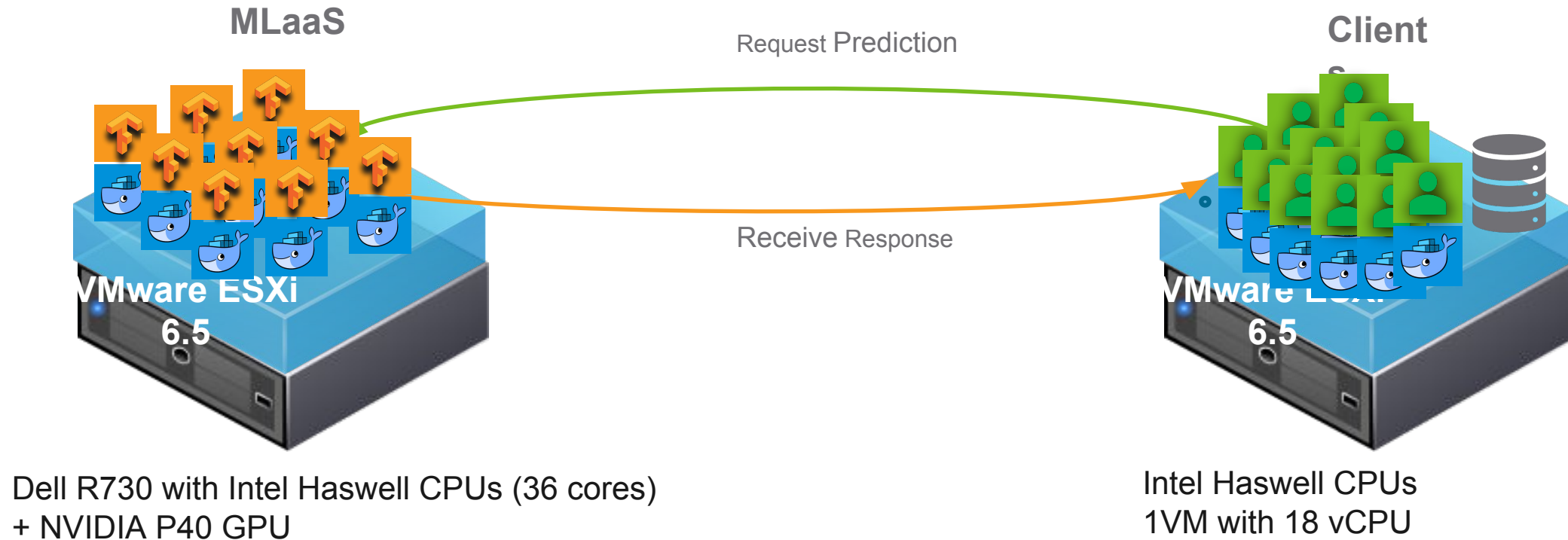# Example #4 of deploying MLaaS on VMware vSphere with Container & Kubernetes



**Machine Learning as a Service**

| Virtual Machine | Virtual Machine | Virtual Machine | Virtual Machine |
|---|---|---|---|
| ML Frameworks | | ML Frameworks | ML Frameworks |
| Docker Container | Kubernetes Master | Docker Container | Docker Container |
| Kubernetes Worker | | Kubernetes Worker | Kubernetes Worker |

NVIDIA GRID | vGPU | VMware vSphere

NVIDIA GRID | vGPU | vGPU | VMware vSphere

CPUs | GPUs | Physical Server

CPUs | GPUs | Physical Server

**vm**ware®

# Experiments of MLaaS on VMware vSphere
## Hardware and Software

MLaaS

Client s

Request Prediction

Receive Response

VMware ESXi
6.5

VMware ESXi
6.5

Dell R730 with Intel Haswell CPUs (36 cores)
+ NVIDIA P40 GPU

Intel Haswell CPUs
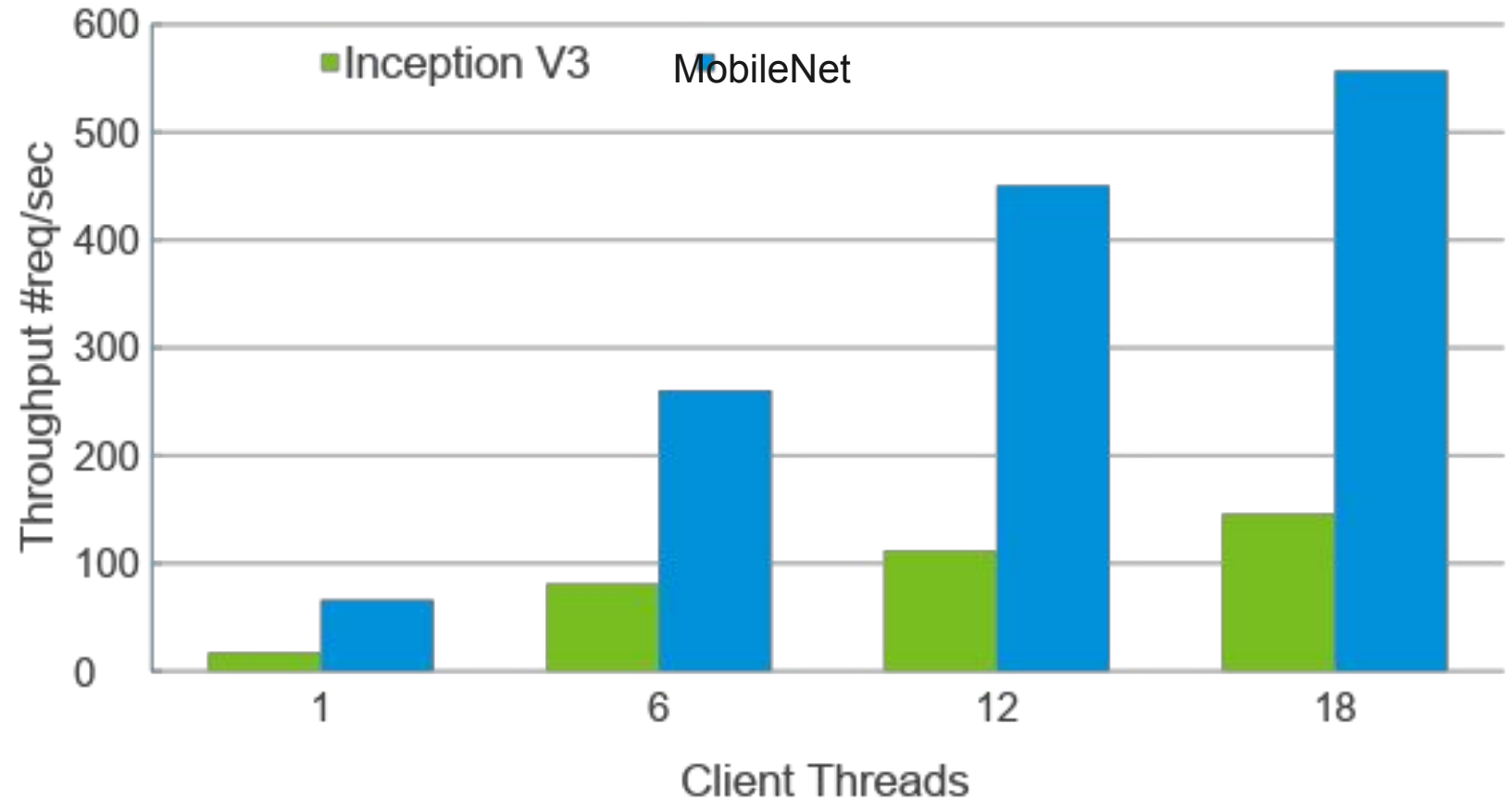1VM with 18 vCPU

**vm**ware®

# Experiment #1: Inference Throughput
Deep Neural Network: Inception V3 vs. MobileNet – Higher is better

**Models:**
*Inception V3*
48 Layers
5000 Million MAC

*MobileNet:*
28 Layers
569 Million MAC

# Experiment #1: Inference Mean Latency

## Deep Neural Network: Inception V3 vs. MobileNet
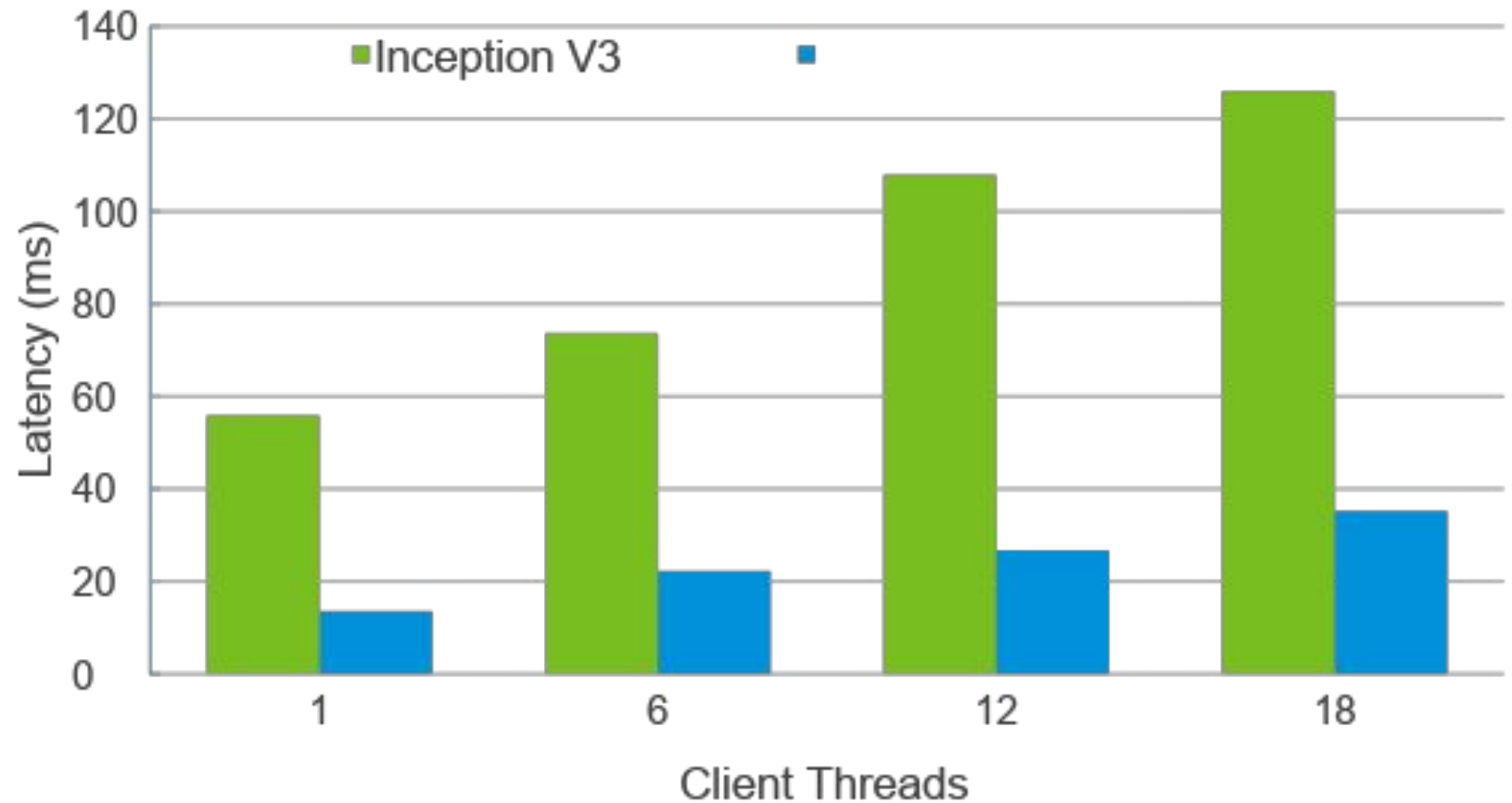
**Models:**
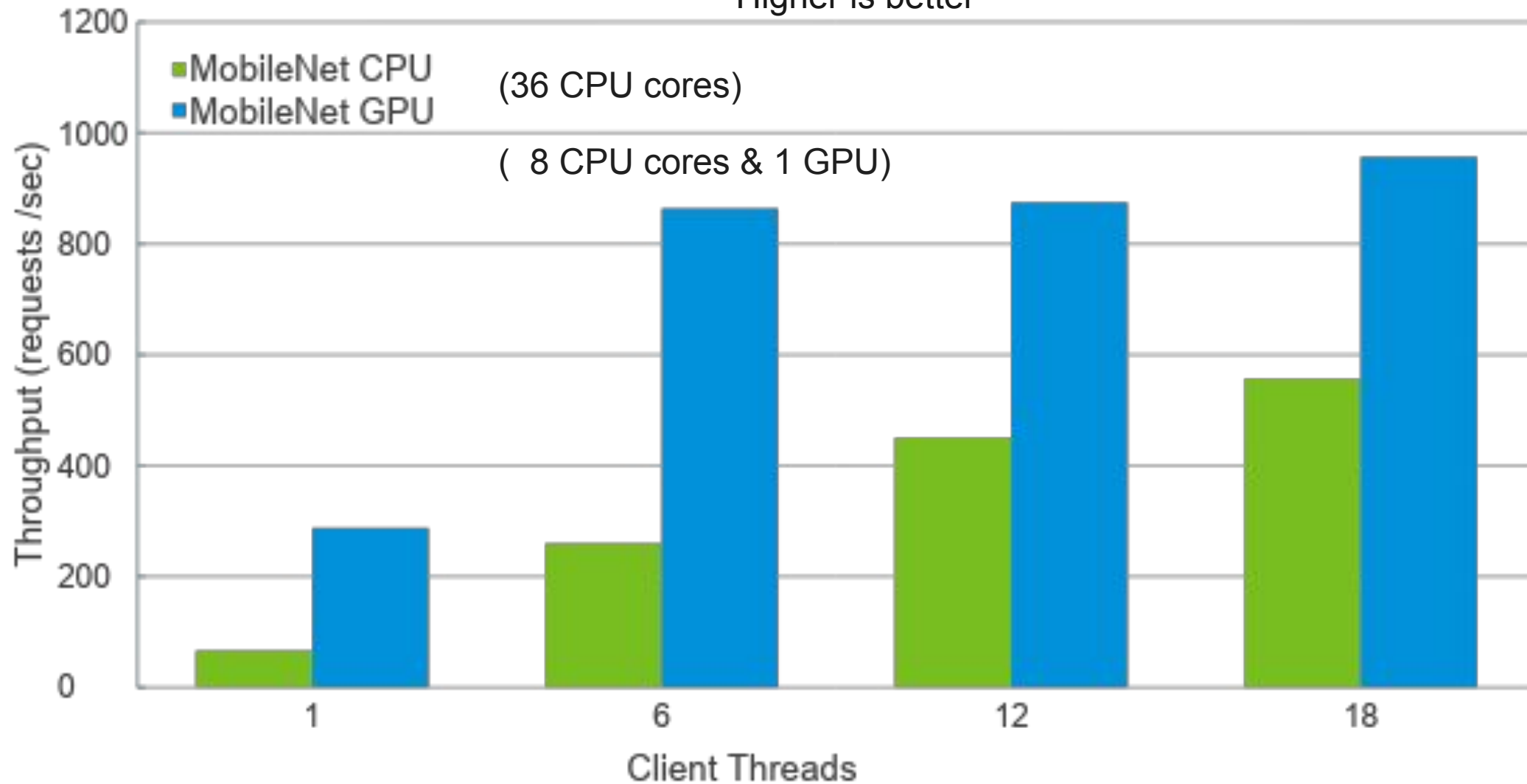*Inception V3*
48 Layers
5000M MAC

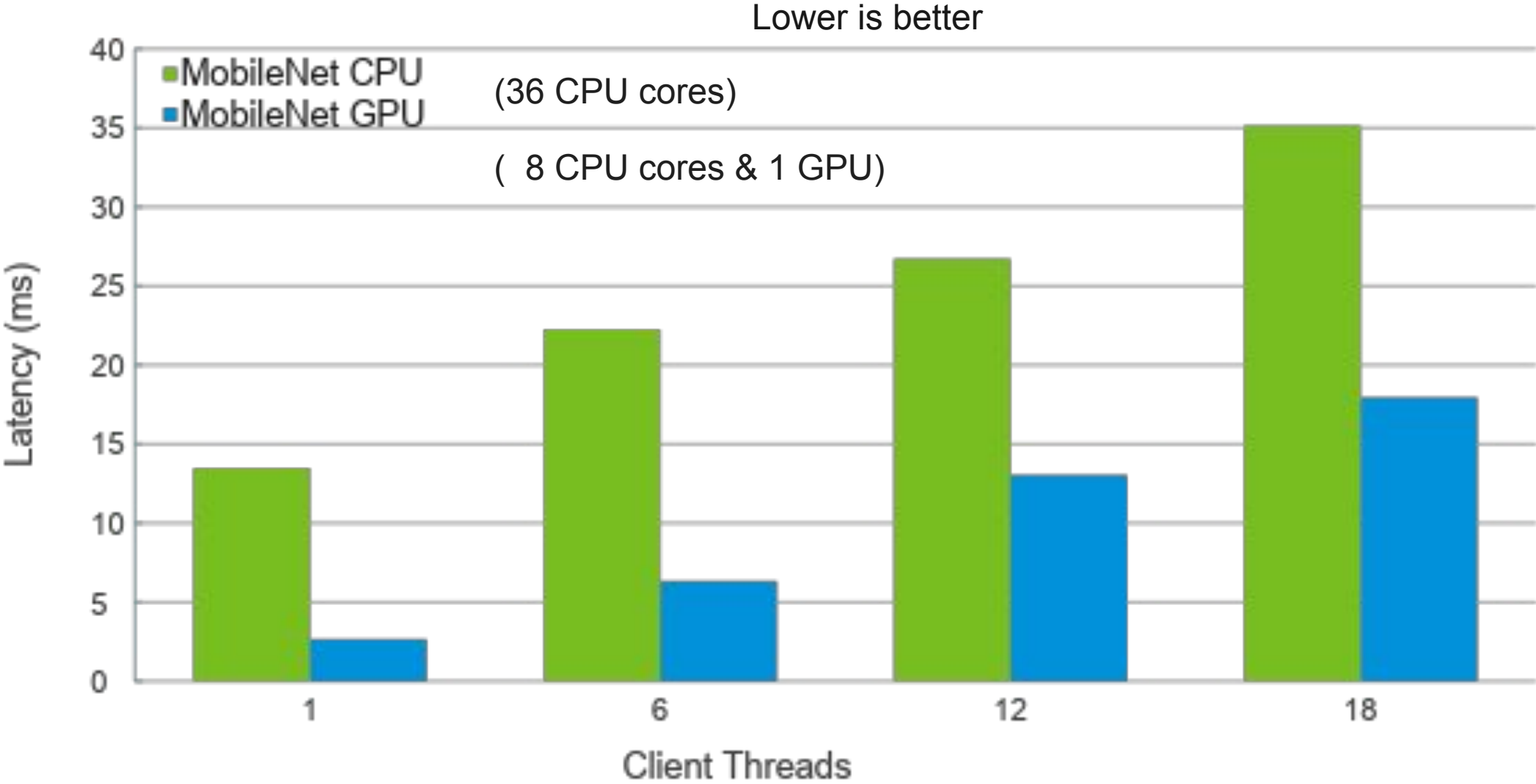*MobileNet:*
28 Layers
569 Million MAC

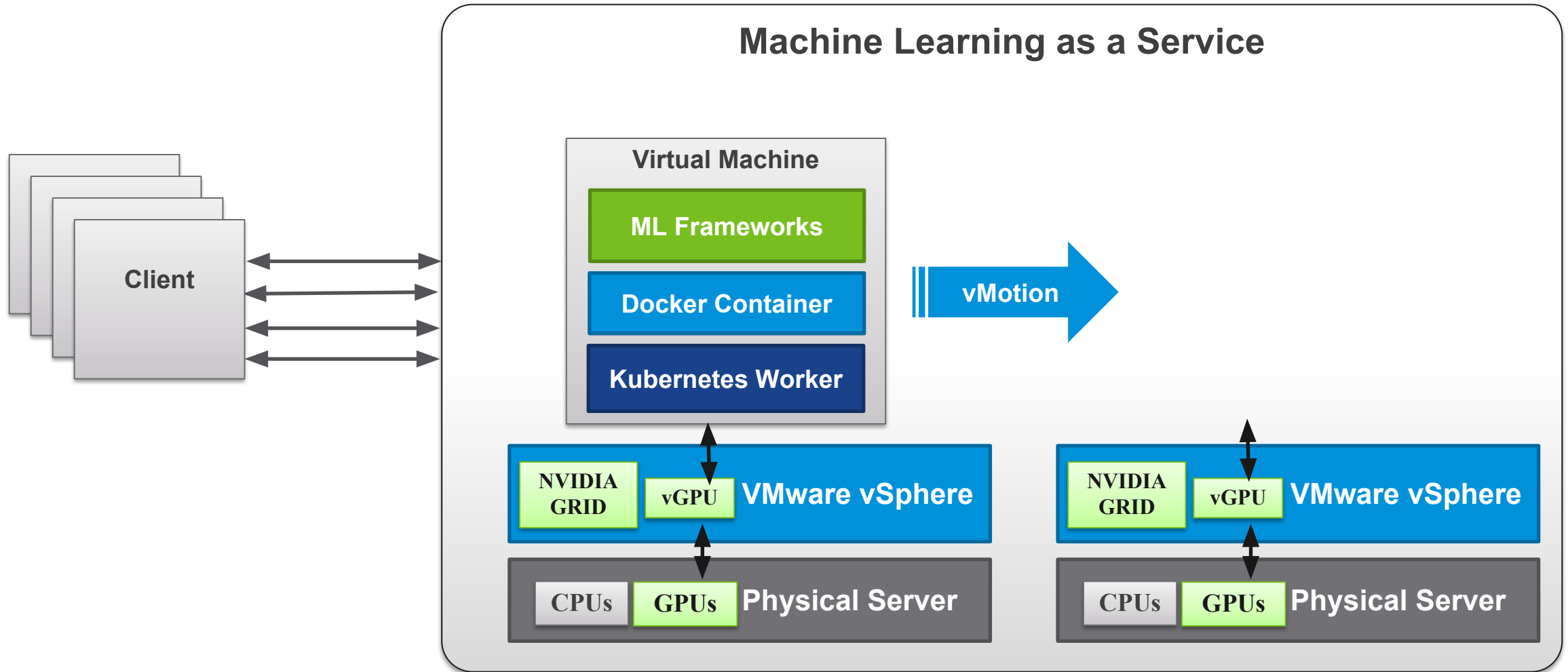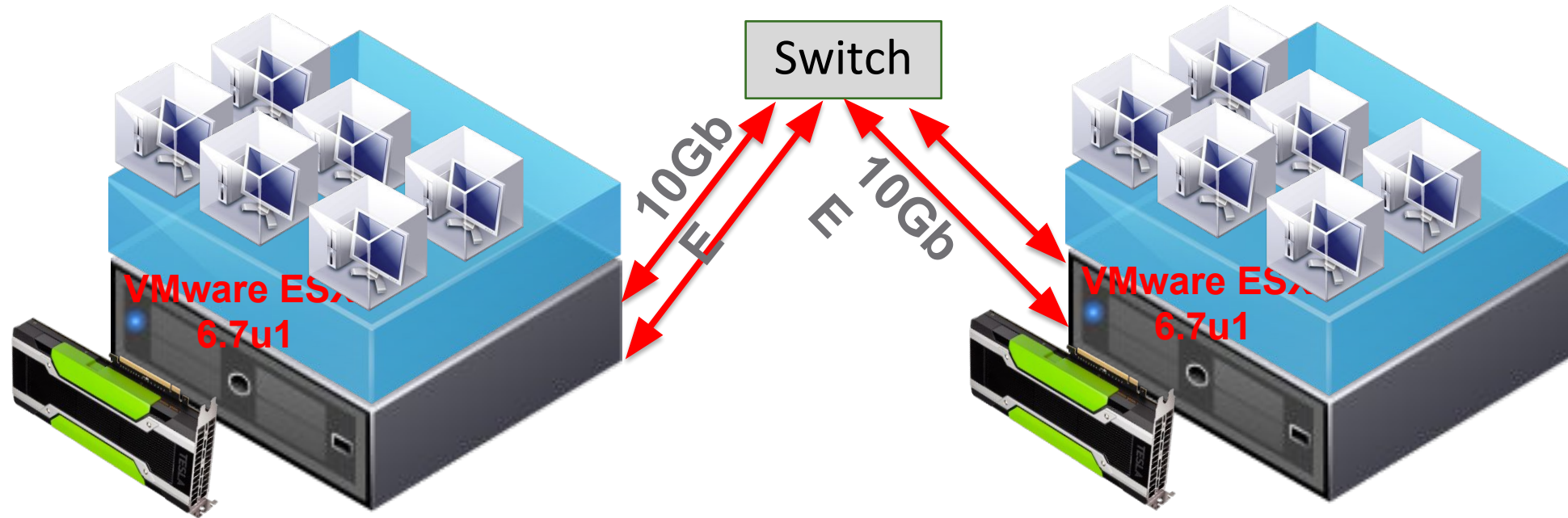# Experiment #2: Inference Throughput



Higher is better

■MobileNet CPU
■MobileNet GPU

(36 CPU cores)

( 8 CPU cores & 1 GPU)

Throughput (requests /sec) vs Client Threads

vmware®

# Experiment #2: Mean Inference Latency



Lower is better

MobileNet CPU (36 CPU cores)
MobileNet GPU ( 8 CPU cores & 1 GPU)

# vMotion for NVIDIA GRID vGPU - MLaaS



**Machine Learning as a Service**

**Client**

**Virtual Machine**
- ML Frameworks
- Docker Container
- Kubernetes Worker

**vMotion**

**NVIDIA GRID** | **vGPU** | **VMware vSphere**
**CPUs** | **GPUs** | **Physical Server**

**NVIDIA GRID** | **vGPU** | **VMware vSphere**
**CPUs** | **GPUs** | **Physical Server**

# vMotion for NVIDIA GRID vGPU – Test-bed



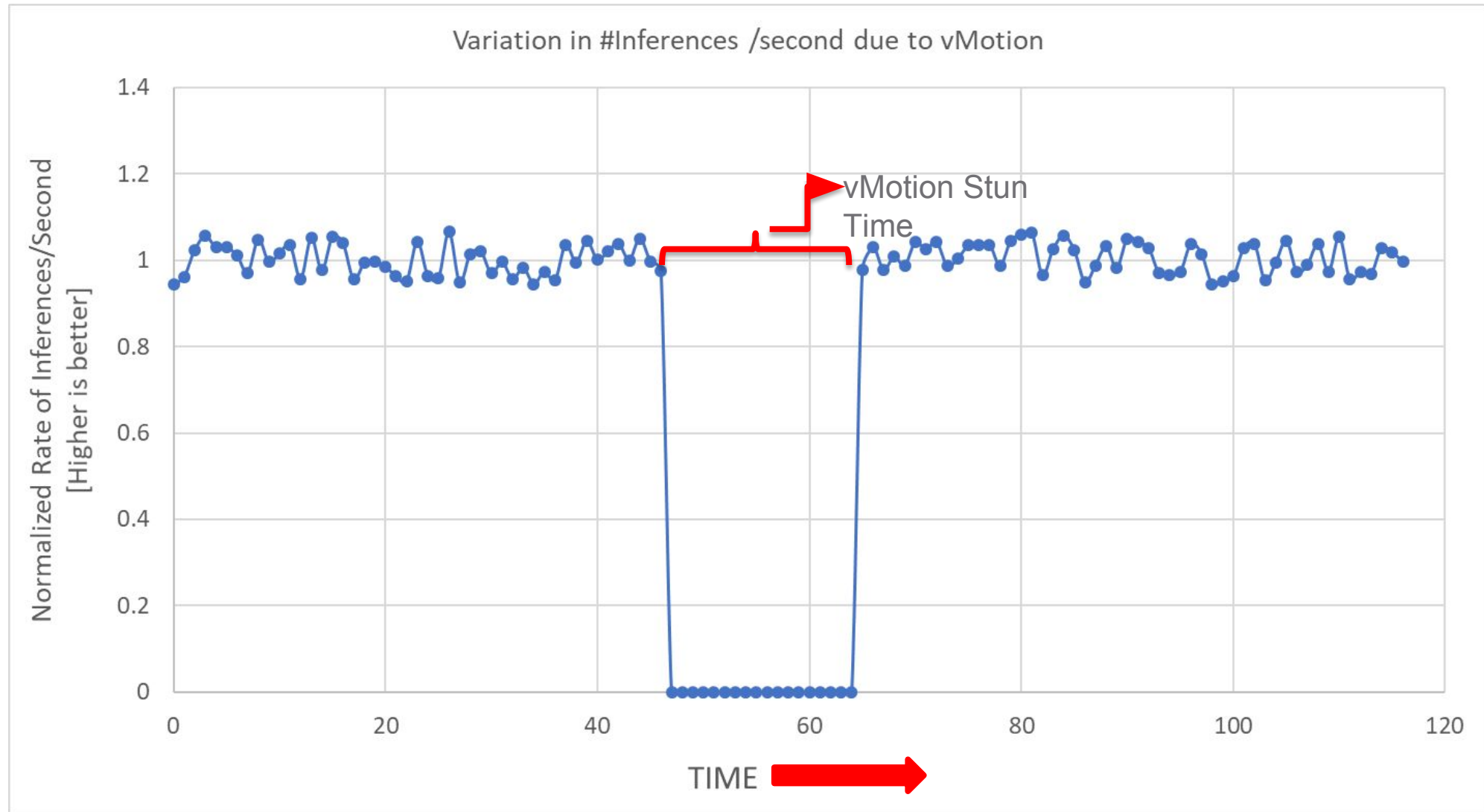**Dell R730 – Intel Broadwell CPUs + 1 x NVidia GRID P40**
40 cores (2 x 20-core socket) E5-2698 v4
768 GB RAM

**Dell R730 – Intel Broadwell CPUs + 1 x NVidia GRID P40**
40 cores (2 x 20-core socket) E5-2698 v4
768 GB RAM

- **ESX**: 6.7u1  **Nvidia Driver**: 410.68

# vMotion for NVIDIA GRID vGPU - MLaaS



Variation in #Inferences /second due to vMotion

# vMotion for Nvidia GRID vGPU: Conclusions and Upcoming Improvements

## Conclusions:

- vMotion for Nvidia GRID vGPU is now available

- The performance impact of vMotion on VDI, CAD and ML applications is negligible or small.

- The performance impact of multiple vMotions running concurrently is small.

## Upcoming Improvements:

- Speedup xfer rate of device checkpoint and vGPU memory data.

- Pre-copy vGPU memory data to reduce stun time to meet or exceed vMotion's standard of 1 second.