

Vienna



# Deep Learning Meetup

## Deep Learning Hardware Overview: What and where to buy or rent



Jan Schlueter  
JKU Linz



René Donner  
contextflow



Thomas Lidy  
Musimap

# Outline

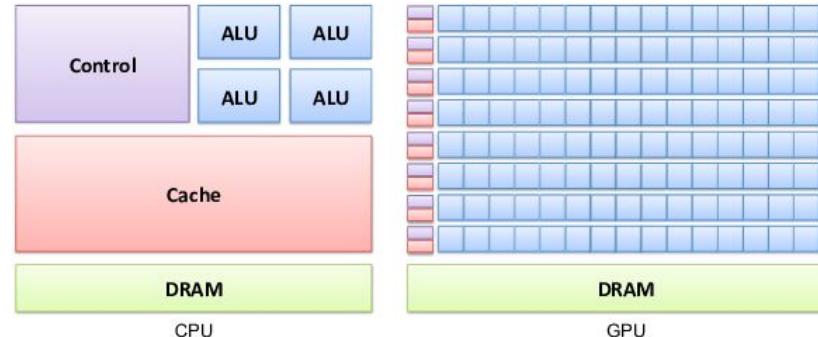
- NVIDIA GPUs
- Intel/AMD GPU-ready Servers
- AMD GPUs
- GPU Cloud Providers

# NVIDIA GPUs

# NVIDIA GPUs

## Why GPUs?

- Deep learning: lots of linear algebra (matrix-matrix products, convolutions)
  - Embarrassingly parallel
- GPUs designed to run *same computation* on different pieces of data
  - multi-core CPUs designed to run *different computations* in parallel
  - CPUs have SIMD (single instruction, multiple data) only on smaller scale
- Larson/McAllister 2001: “Fast Matrix Multiplies using Graphics Hardware”



Lounis et al., 2015

<https://doi.org/10.1145/582034.582089>

# NVIDIA GPUs

## Why NVIDIA?

- saw the potential of GPGPUs (general-purpose graphical processing units)
- started selling GPUs *without graphics output* in 2007
- invested a lot in software to make it accessible: CUDA, cuBLAS, cuDNN, ...
- academic hardware grant program to donate GPUs to universities
- led to wide adoption in academia, and academic software libraries

# NVIDIA GPUs

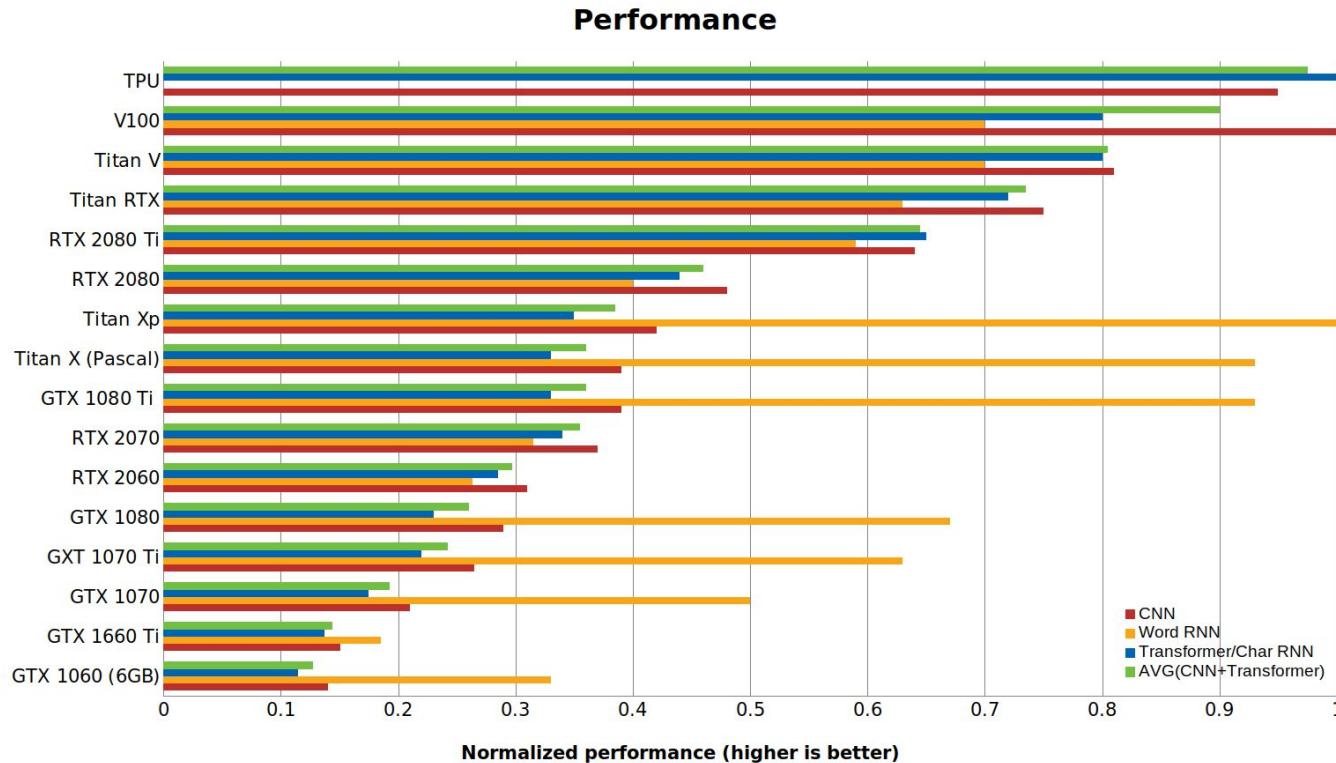
## Main criteria

- Number of CUDA cores (32-bit and/or 64-bit precision)
- Number of Tensor cores (32-bit and 16-bit precision, possibly int16/8/4)
- Memory size
- Memory bandwidth

## Other criteria

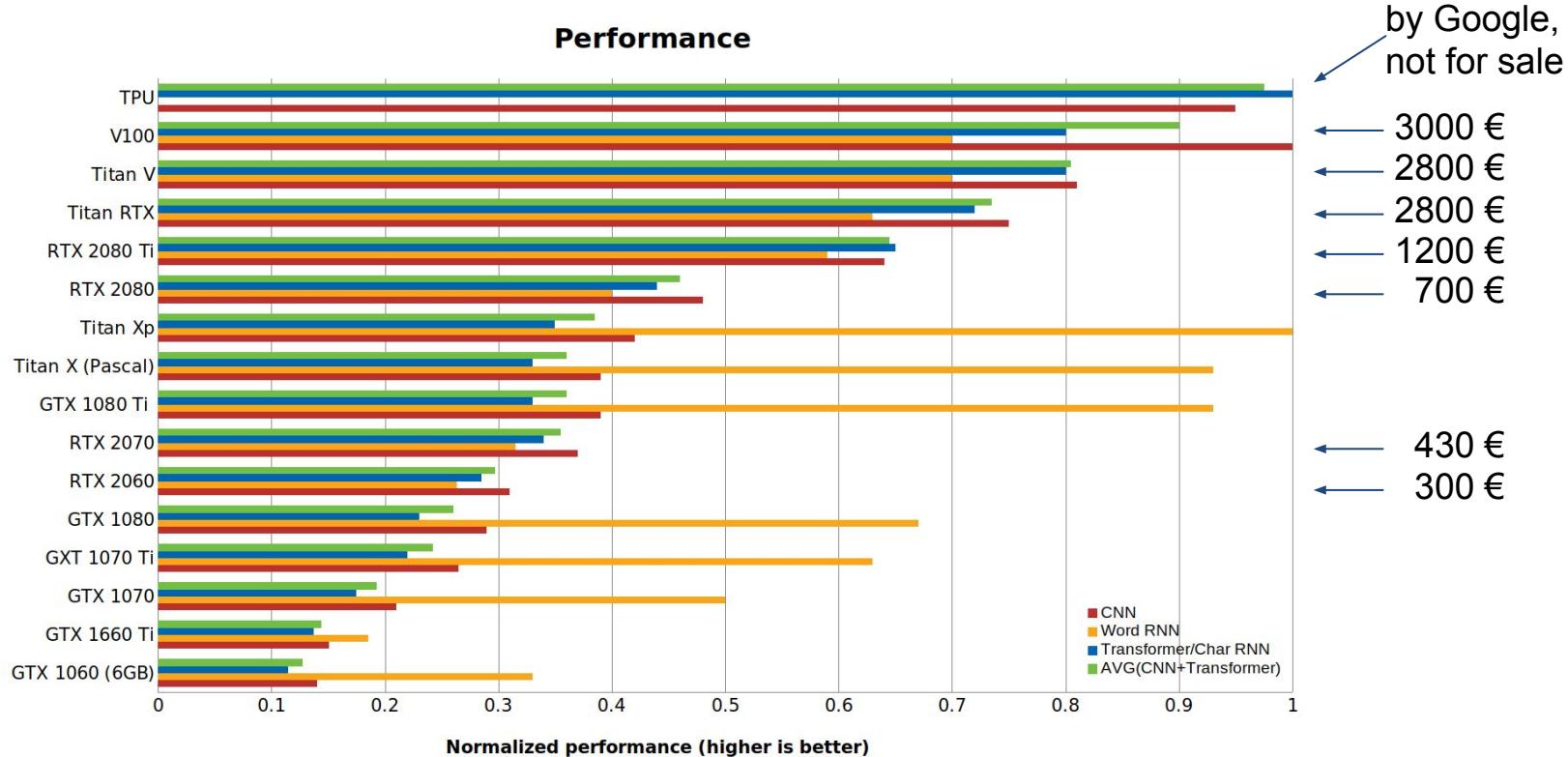
- 64-bit performance
- power draw / cooling requirements

# NVIDIA GPUs

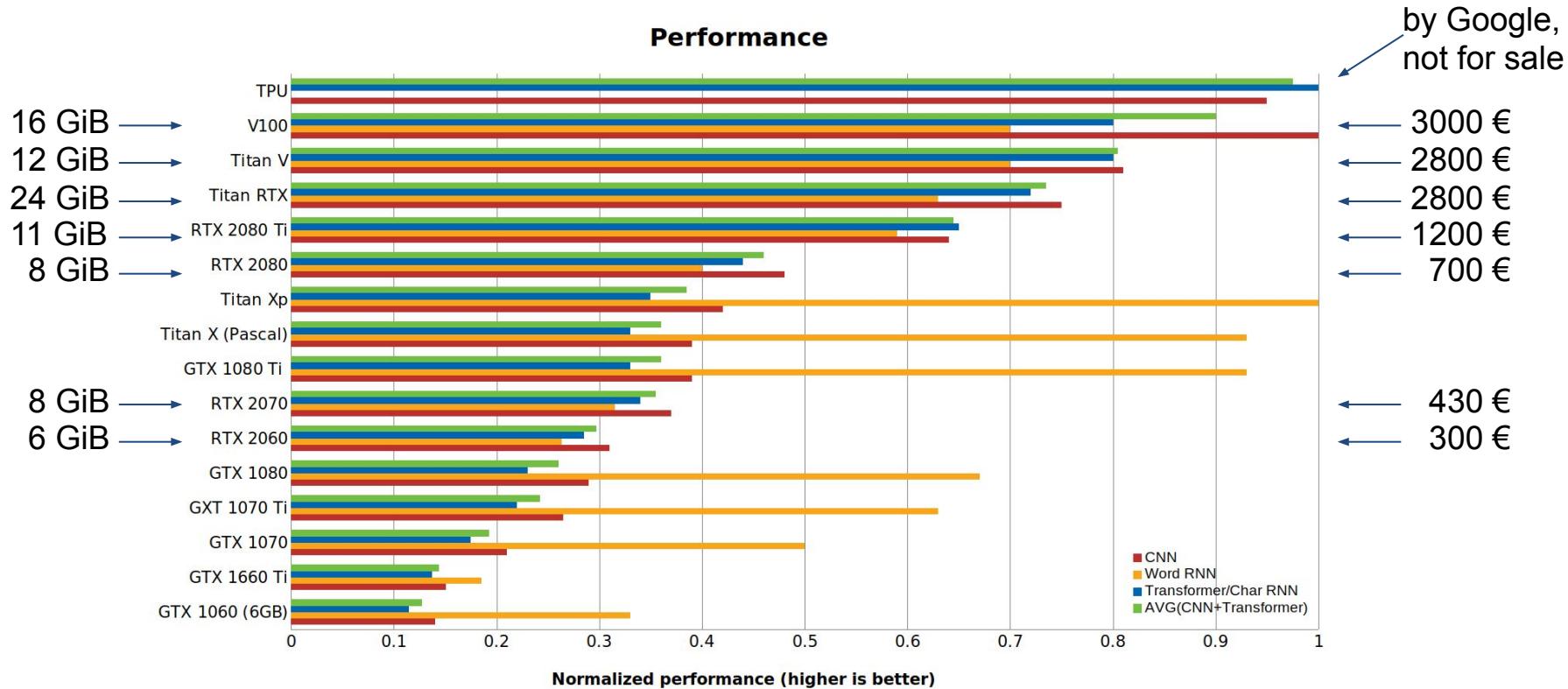


<https://timdettmers.com/2019/04/03/which-gpu-for-deep-learning/>

# NVIDIA GPUs

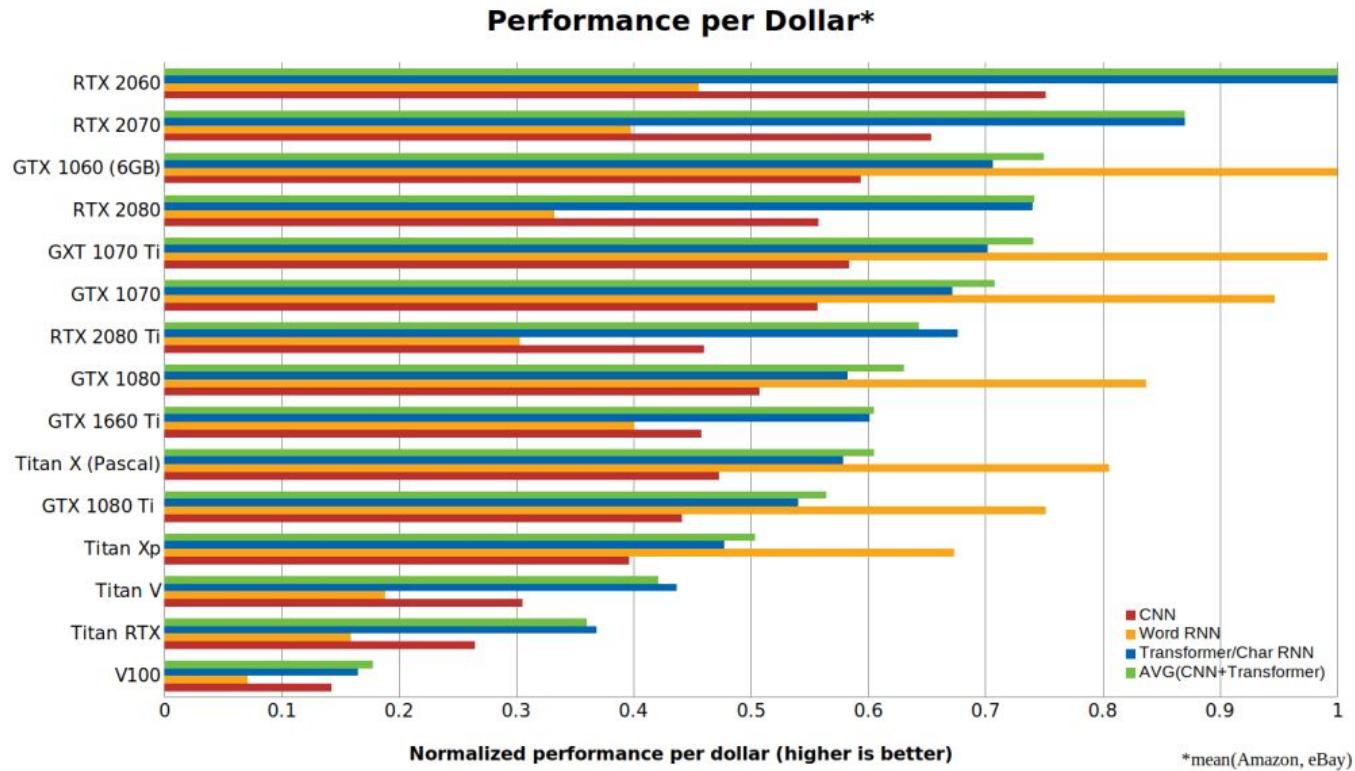


# NVIDIA GPUs



<https://timdettmers.com/2019/04/03/which-gpu-for-deep-learning/>

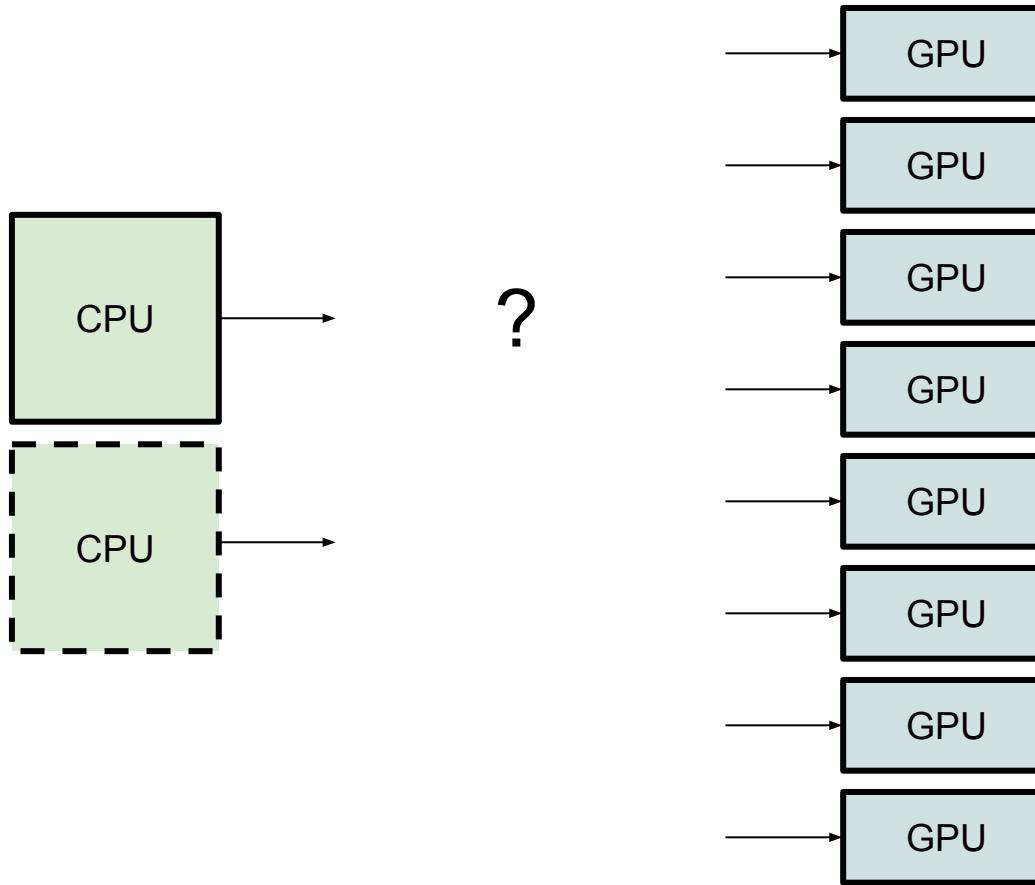
# NVIDIA GPUs



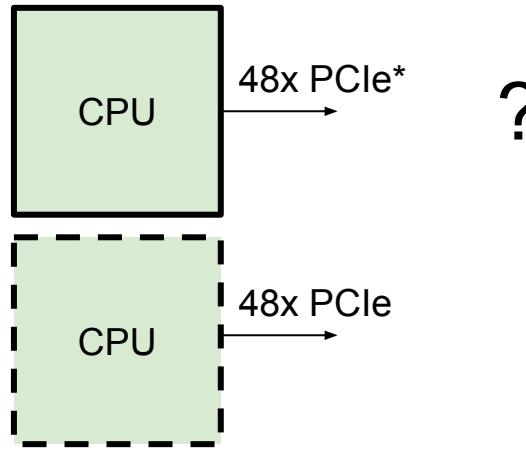
<https://timdettmers.com/2019/04/03/which-gpu-for-deep-learning/>

# **Intel/AMD GPU-ready servers**

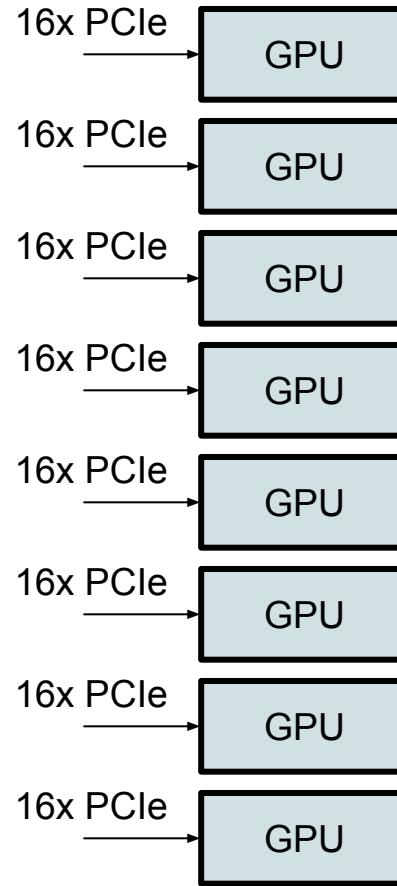
# Intel/AMD GPU-ready servers



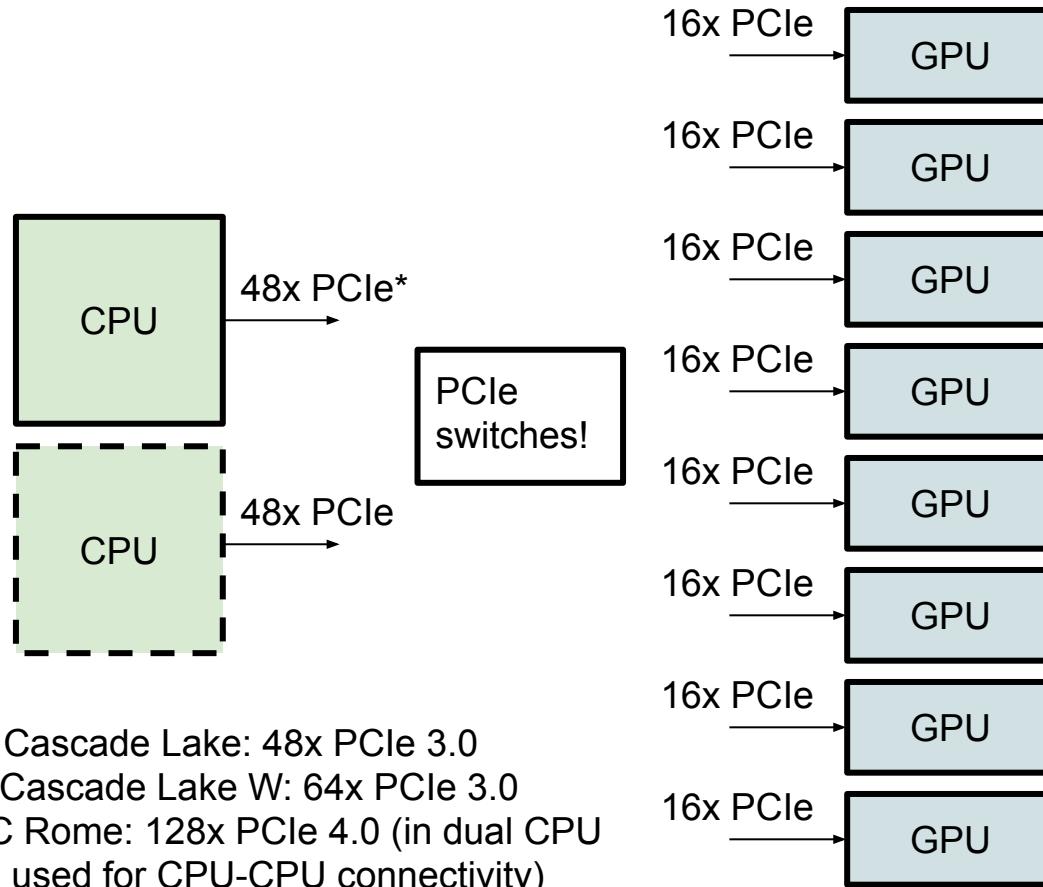
# Intel/AMD GPU-ready servers



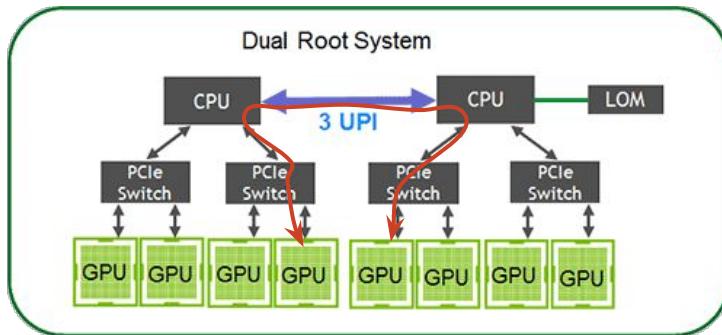
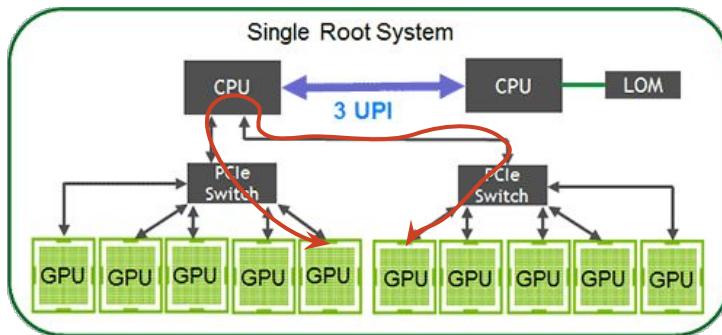
\* Intel Xeon Cascade Lake: 48x PCIe 3.0  
Intel Xeon Cascade Lake W: 64x PCIe 3.0  
AMD EPYC Rome: 128x PCIe 4.0 (in dual CPU system: 64 used for CPU-CPU connectivity)



# Intel/AMD GPU-ready servers

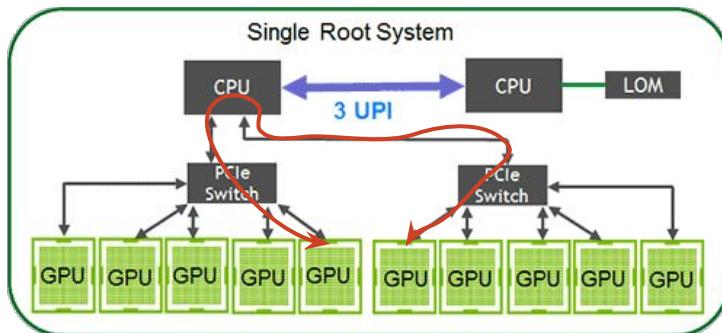


# Single root complex vs. dual root complex

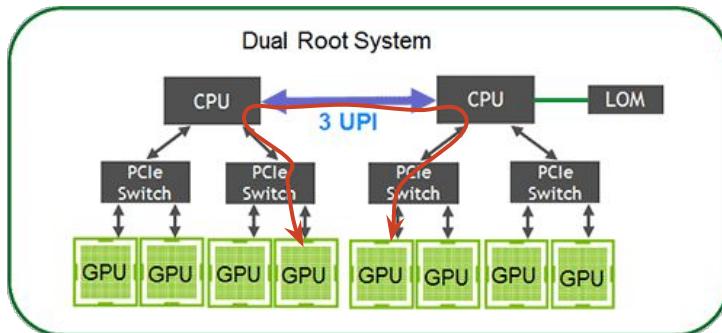


<https://www.supermicro.com/products/system/4U/4029/PCIe-Root-Architecture.cfm>

# Single root complex vs. dual root complex

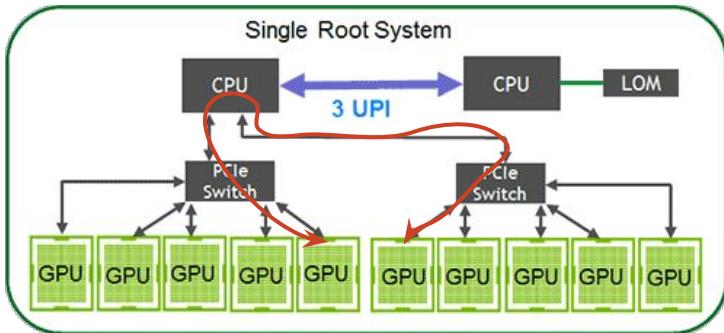


high GPU-GPU  
bandwidth

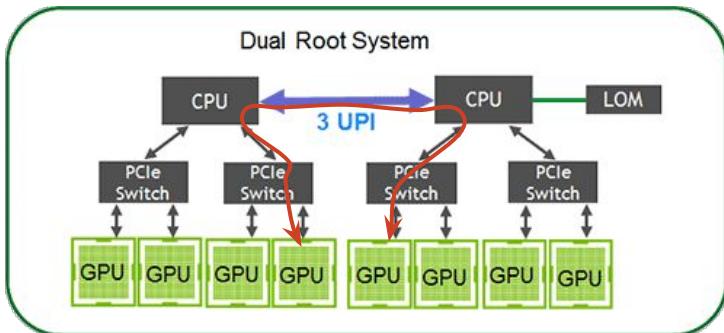


high GPU-CPU  
bandwidth

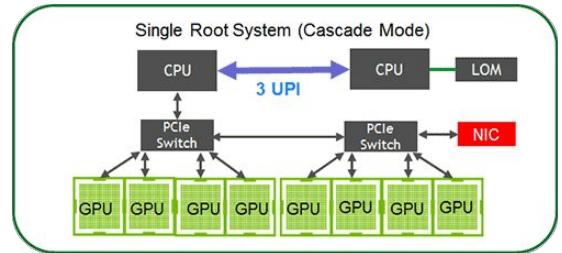
# Single root complex vs. dual root complex



high GPU-GPU bandwidth

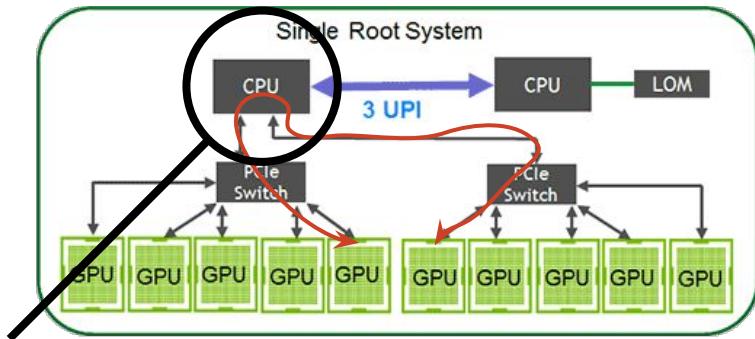


high GPU-CPU bandwidth

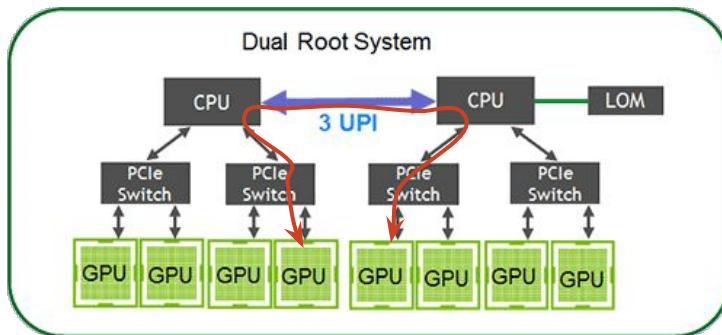


high GPU-GPU bandwidth, but only 16 lanes between CPU and all GPUs

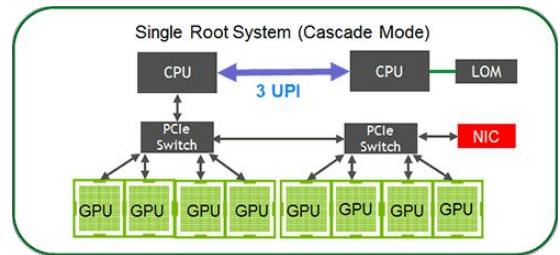
# Single root complex vs. dual root complex



high GPU-GPU bandwidth



high GPU-CPU bandwidth



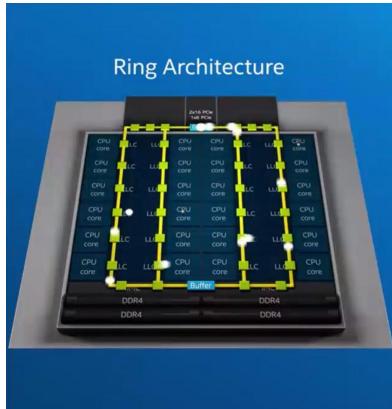
high GPU-GPU bandwidth, but only 16 lanes between CPU and all GPUs

# Single root complex vs. dual root complex

## Intel Broadwell:

Ring architecture

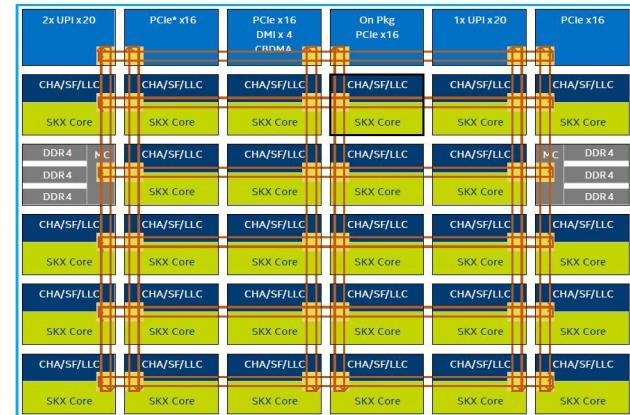
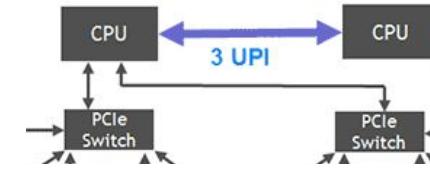
Two PCIe controllers in single root



## Intel Cascade Lake:

Mesh architecture

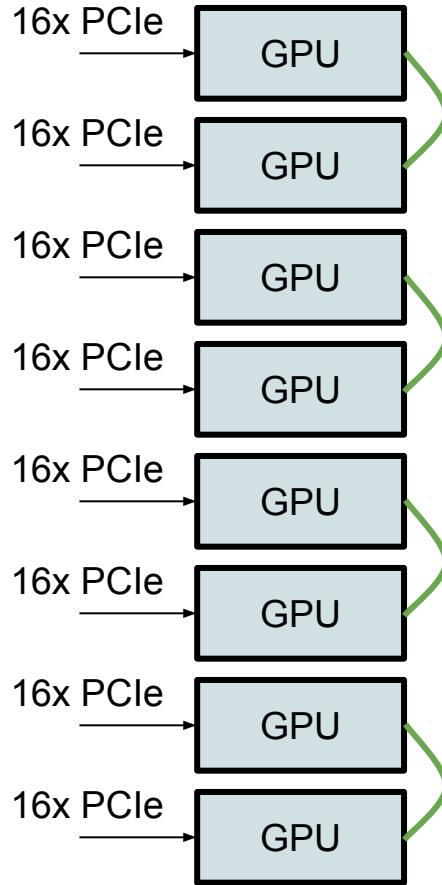
Two PCIe controllers on different hops



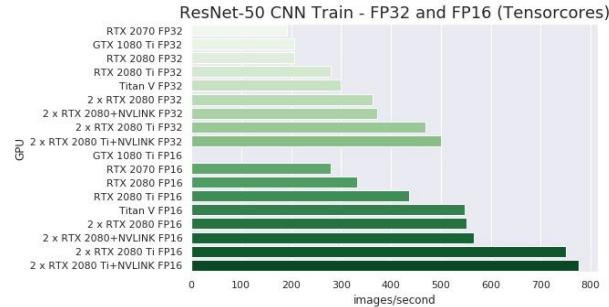
CHA – Caching and Home Agent ; SF – Snoop Filter; LLC – Last Level Cache;  
SKY Core – Skylake Server Core; UPI – Intel® UltraPath Interconnect

<https://www.servethehome.com/how-intel-skylake-sp-changes-impacted-single-root-pcie-due-to-deep-learning-servers/>

# NVlink



- High-bandwidth interconnect between pairs of GPUs
- 3-4% performance improvement in dual-GPU ResNet training



<https://www.pugetsystems.com/labs/hpc/RTX-2080Ti-with-NVLINK---TensorFlow-Performance-Includes-Comparison-with-GTX-1080Ti-RTX-2070-2080-2080Ti-and-Titan-V-1267/>

# Server sizes

## 1U (rack height unit)

- up to 4 GPUs
- hard to cool
- only for inference, e.g. NVIDIA T4
- some vendors will sell it with RTX 2080 Ti, not recommended!



## 2U

- up to 8 GPUs
- still hard to cool, need to choose carefully



## 4U

- up to 10 GPUs
- has space for NVlink

# Server recommendations

## Gigabyte G291-Z20

- 2U server for 8 GPUs
- AMD EPYC Rome
- cool enough if only putting 4 GPUs
- ~10k€ all in all



## Supermicro 4029-TRT2

- 4U server for 10 GPUs
- Intel Cascade Lake
- need special top cover for RTX cards and NVlink (makes it 4.5U)
- ~20k€ all in all



## Other options

- smaller Supermicro servers: not cooled well enough
- Tyan, Asus, Asrock: more expensive
- Zstor GS-P206G: custom 2U 6GPU design with good cooling, but expensive

# Possible vendors

- server-bau.de
- deltacomputer.de
- sysgen.de
- bressner.de
- thomas-krenn.com
- zstor.de

# AMD GPUs

# AMD GPUs

- Good value / money
- Focussed on gaming so far
- CUDA -> ROCm ecosystem

# AMD GPUs

- Vega 20 chip
- Navi 10/14 chips

	Score	GPU	Base/Boost	Memory	Power	Buy	
Nvidia Titan RTX	100	TU102	1350/1770 MHz	24GB GDDR6	280W	<b>US\$2,449</b>	<a href="#">VIEW DEAL AMAZON</a>
Nvidia GeForce RTX 2080 Ti	98.4	TU102	1350/1635 MHz	11GB GDDR6	260W	<b>US\$1,186.94</b>	<a href="#">VIEW DEAL AMAZON</a>
Nvidia GeForce RTX 2080 Super	98.2	TU104	1650/1815 MHz	8GB GDDR6	250W	<b>US\$696.16</b>	<a href="#">VIEW DEAL AMAZON</a>
Nvidia GeForce RTX 2080	96.1	TU104	1515/1800 MHz	8GB GDDR6	225W	<b>US\$1,102.64</b>	<a href="#">VIEW DEAL AMAZON</a>
Nvidia Titan X	96.0	GP102	1405/1480 MHz	12GB GDDR5X	250W	<a href="#">CHECK AMAZON</a>	
Nvidia GeForce GTX 1080 Ti	96.0	GP102	1480/1582 MHz	11GB GDDR5X	250W	<b>US\$769.99</b>	<a href="#">VIEW DEAL AMAZON</a>
AMD Radeon RX 5700 XT	95.8	Navi 10	1605/1905 MHz	8GB GDDR6	225W	<b>US\$389.99</b>	<a href="#">VIEW DEAL AMAZON</a>
Nvidia GeForce RTX 2070 Super	94.1	TU104	1605/1770 MHz	8GB GDDR6	215W	<b>US\$599.42</b>	<a href="#">VIEW DEAL AMAZON</a>
AMD Radeon VII	92.4	Vega 20	1400/1750 MHz	16GB HBM2	300W	<b>US\$599.99</b>	<a href="#">VIEW DEAL AMAZON</a>

# AMD ROCm software stack

- Replaces CUDA
- cuDNN -> MIOpen



The screenshot shows the official ROCm website. At the top, there's a dark header with the ROCm logo on the left and the text "ROCM, a New Era in Open GPU Computing" in large white font. Below the header is a sub-header "Platform for GPU-Enabled HPC and Ultrascale Computing". A navigation bar follows, containing links for Overview, Getting Started, Documentation, ROCm Ecosystem, Deep Learning, Tutorials, and Community. There's also a "Contribute" button. The main content area features a large heading "Welcome to the ROCm Platform". Below it is a paragraph about the platform's purpose and philosophy. A note at the bottom states "ROCM is built for scale; it supports multi-GPU computing in and out of server-node communication through RDMA. It also simplifies the stack when the driver directly incorporates RDMA peer access support." Finally, a URL "https://rocm.github.io/" is provided.

ROCM, a New Era in Open GPU Computing

Platform for GPU-Enabled HPC and Ultrascale Computing

Overview Getting Started Documentation ROCm Ecosystem Deep Learning Tutorials Community

Contribute

## Welcome to the ROCm Platform

We are excited to present ROCm, the first open-source HPC/Hyperscale-class platform for GPU computing that's also programming-language independent. We are bringing the UNIX philosophy of choice, minimalism and modular software development to GPU computing. The new ROCm foundation lets you choose or even develop tools and a language run time for your application.

**ROCM is built for scale;** it supports multi-GPU computing in and out of server-node communication through RDMA. It also simplifies the stack when the driver directly incorporates RDMA peer access support.

<https://rocm.github.io/>

# AMD ROCm software stack

ROCM, a New Era in Open Computing

Platform for GPU-Enabled HPC and Ultrascale Computing

Overview Getting Started Documentation ROCm Ecosystem Deep Learning Tutorials Contributing

## Deep Learning on ROCm

**TensorFlow:** TensorFlow for ROCm – latest supported official version 1.14.1 and 2.0-beta3 ROCm Community has landed on the official Tensorflow repository.

**MIOpen:** Open-source deep learning library for AMD GPUs – latest supported version 1.7.1

**PyTorch:** PyTorch for ROCm – latest supported version 1.0

## Debugger

- ROCm debugger binary build
- ROCm GDB source code
- ROCm GPU debugger SDK; ROCm debug run time that services GDB

## Profiling Tools

- `rocprof`
- `rocProfiler`
- `rocTracer`

## Math Libraries

- `rocBLAS`
- `rocFFT`
- `Tensile`
- [Eigen C++ Math Library with HIP Based GPU Acceleration](#)
- `clBLAS`
- `clFFT`
- `clSparse`
- `clRNG`

© 2016 AMD Corporation [Disclaimer and Legal Information](#)

<https://rocm.github.io/>

# Benchmarks

▲ bonoboTP 88 days ago | parent | favorite | on: Tensorflow 2.0 AMD Supp

I just ran LambdaLabs' ResNet50 training benchmark (I could do more m

TensorFlow 1.14.4

```
git clone https://github.com/lambdal/lambda-tensorflow-be
python lambda-tensorflow-benchmark/benchmarks/scripts/tf_
```

## FP32

- NVIDIA GTX 1080 Ti: ~215 images/sec
- NVIDIA RTX 2080 Ti: ~300 images/sec
- NVIDIA TITAN RTX: ~320 images/sec
- NVIDIA Tesla V100: ~383 images/sec
- AMD Radeon VII: ~275 images/sec

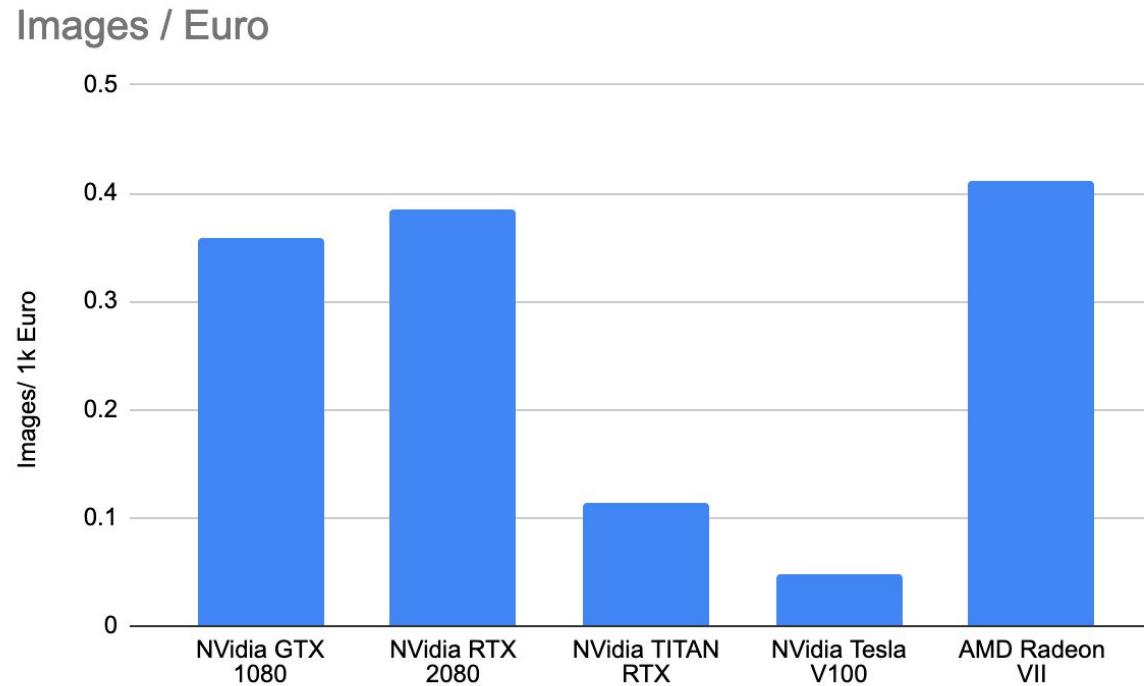
## FP16

- NVIDIA GTX 1080 Ti: ~277 images/sec
- NVIDIA RTX 2080 Ti: ~495 images/sec
- NVIDIA TITAN RTX: ~518 images/sec
- NVIDIA Tesla V100: ~725 images/sec
- AMD Radeon VII: ~373 images/sec



<https://news.ycombinator.com/item?id=21666411>

# Value for money?



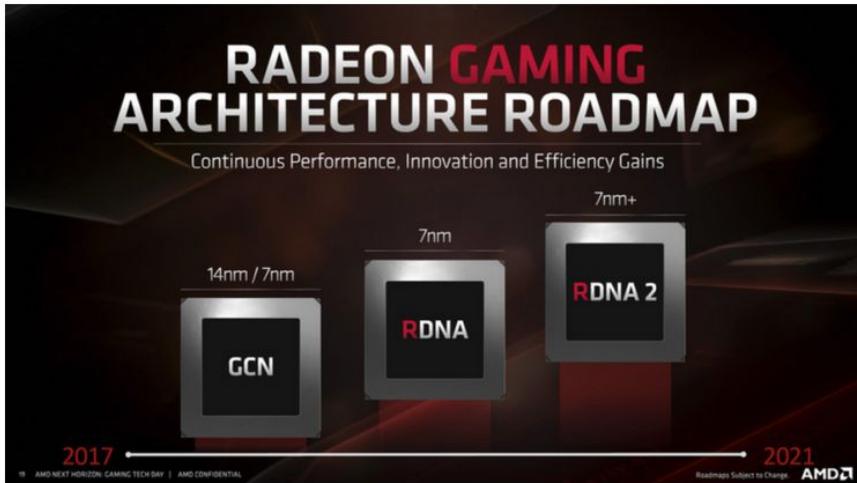
# AMD GPU cloud provider



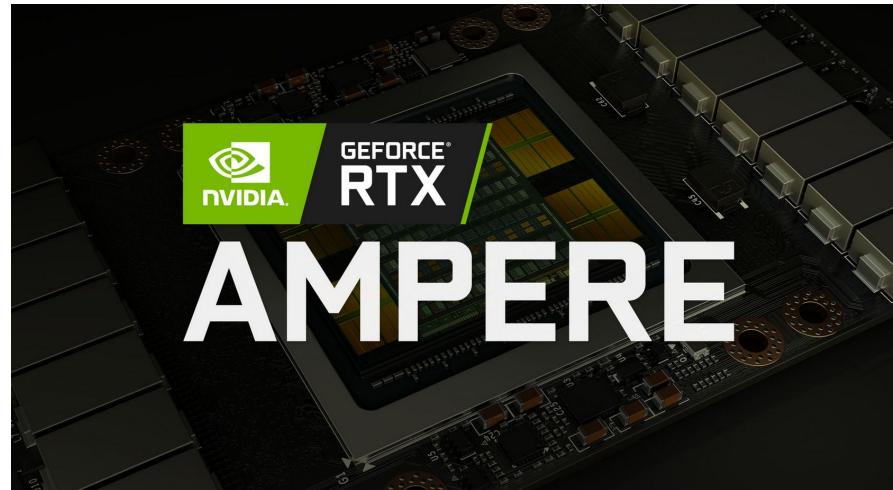
Provider	Plan	TFlops	Cores	GPU Mem	Architecture	Hourly	Monthly
amazon web services	p2.xlarge	4.37	2496	12GB	NVIDIA Kepler	\$0.9000/h	\$648/m
amazon web services	g3.4xlarge	4.82	2048	8GB	NVIDIA Maxwell	\$1.1400/h	\$820/m
GPU EATER	a1.rx580	<b>6.1</b>	<b>2304</b>	<b>8GB</b>	<b>AMD RADEON RX</b>	<b>\$0.3458/h</b>	<b>\$249/m</b>
GPU EATER	a1.vega56	<b>10.5</b>	<b>3584</b>	<b>8GB</b>	<b>AMD RADEON VEGA</b>	<b>\$0.4794/h</b>	<b>\$345/m</b>
GPU EATER	a1.vegafe	<b>13.1</b>	<b>4096</b>	<b>16GB</b>	<b>AMD RADEON VEGA</b>	<b>\$0.6164/h</b>	<b>\$443/m</b>
amazon web services	p3.2xlarge	14.0	5120	16GB	NVIDIA Volta	\$3.0600/h	\$2203/m

<https://www.gpueater.com/index.html>

# AMD vs NVidia 2020

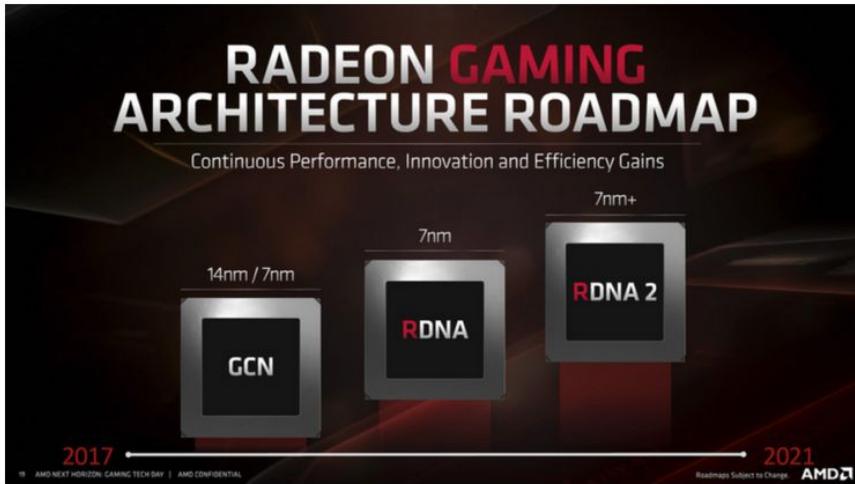


New AMD “Big Navi” chip



New NVidia “Ampere” architecture

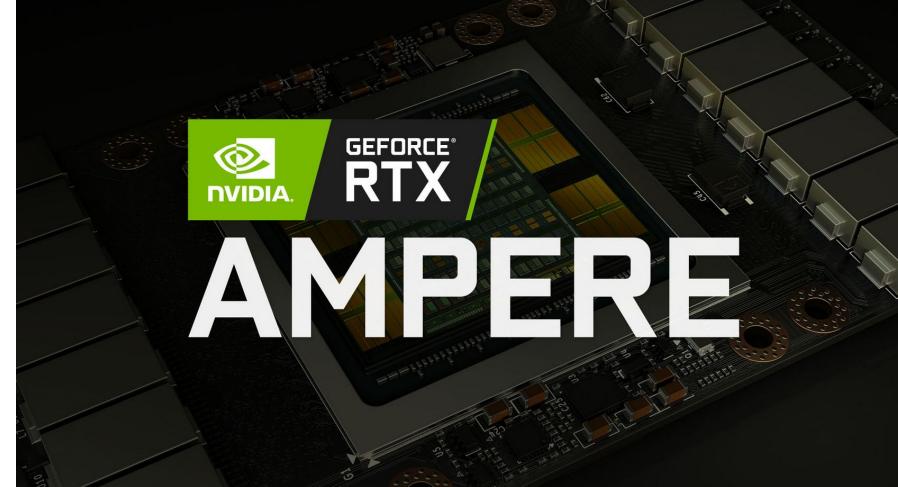
# AMD vs NVidia 2020



New AMD “Big Navi” chip

**AMD Radeon Navi 20 Could Pack 5,120 Cores, 24GB of HBM2: Report**

By Joel Hruska on February 26, 2020 at 7:47 am | 1 Comment



New NVidia “Ampere” architecture

# GPU Cloud Providers

# GPU Cloud Providers: Overview

- Amazon Web Services
- Floydhub
- Linode
- Exoscale
- Paperspace
- LeaderGPU
- Packet
- OVHcloud
- LambdaLabs
- Hetzner
- Vast.ai

## Accelerated Computing:

- P2 Instances: 1, 8 or 16 **NVIDIA K80 GPUs** (12 GB VRAM each)
- P3 Instances: 1, 4 or 8 **NVIDIA Tesla V100 GPUs** (16 GB VRAM each)

Intel Xeon E5-2686 v4 (Broadwell) processors with 2.3 GHz (2.7 GHz turbo)  
(or Xeon P-8175M processors for p3dn.24xlarge)

<https://aws.amazon.com/ec2/instance-types/>



# Amazon Web Services



## Accelerated Computing - P2 Instances (Nvidia K80):

Instance	GPUs	vCPU	Mem (GiB)	GPU		Network Performance	Pricing	
				Memory (GiB)			per hour	per month
p2.xlarge	1	4	61	12		High	\$0,97	\$700
p2.8xlarge	8	32	488	96		10 Gigabit	\$7,78	\$5.599
p2.16xlarge	16	64	732	192		25 Gigabit	\$15,55	\$11.197

Pricing for on-demand in EU-Ireland Region: <https://aws.amazon.com/ec2/pricing/on-demand/>  
(up to 90% reduction using “Spot Instances”)

# Amazon Web Services



## Accelerated Computing - P3 Instances (Nvidia Tesla V100):

Instance	GPUs	vCPU	Mem (GiB)	GPU Mem (GiB)	GPU P2P	Storage (GB)	Dedicated EBS Bandwidth	Networking Performance	Pricing	
									per hour	per month
p3.2xlarge	1	8	61	16	-	EBS-Only	1.5 Gbps	Up to 10 Gigabit	\$3,31	\$2.380
p3.8xlarge	4	32	244	64	NVLink	EBS-Only	7 Gbps	10 Gigabit	\$13,22	\$9.518
p3.16xlarge	8	64	488	128	NVLink	EBS-Only	14 Gbps	25 Gigabit	\$26,44	\$19.037
p3dn.24xlarge	8	96	768	256	NVLink	2 x 900 NVMe SSD	19 Gbps	100 Gigabit	\$33,71	\$24.272

Pricing for on-demand in EU-Ireland Region: <https://aws.amazon.com/ec2/pricing/on-demand/>  
(up to 90% reduction using “Spot Instances”)

# Amazon Web Services



## Amazon SageMaker

fully managed service to build, train, and deploy machine learning models quickly

Label	Build	Train & Tune	Deploy & Manage
<b>Amazon SageMaker Ground Truth</b> Build and manage training data sets		<b>Amazon SageMaker Studio</b> Integrated development environment (IDE) for machine learning	
	<b>Amazon SageMaker Autopilot</b> Automatically build and train models		<b>Amazon SageMaker Model Monitor</b> Automatically detect concept drift
	<b>Amazon SageMaker Notebooks</b> One-click notebooks with elastic compute	<b>Amazon SageMaker Experiments</b> Capture, organize, and search every step	<b>Amazon SageMaker Neo</b> Train once, deploy anywhere
	<b>AWS Marketplace</b> Pre-built algorithms and models	<b>Amazon SageMaker Debugger</b> Debug and profile training runs	<b>Amazon Augmented AI</b> Add human review of model predictions
		<b>Automatic Model Tuning</b> One-click hyperparameter optimization	

<https://aws.amazon.com/sagemaker/>

# FLOYDHUB

Dedicated deep learning machines, ready for your next project.

MACHINE	GPU	CPU	RAM	DISK	PRICE
gpu	Tesla K80 12 GB Memory	4 vCPUs	61 GB	200 GB SSD	\$1.20/hr \$0.00033/sec
gpu2	Tesla V100 16 GB Memory	8 vCPUs	61 GB	200 GB SSD	\$4.20/hr \$0.00117/sec
cpu		2 vCPUs	8 GB	200 GB SSD	\$0.19/hr \$0.00005/sec
cpu2		8 vCPUs	32 GB	200 GB SSD	\$0.48/hr \$0.00013/sec

<https://www.floydhub.com/pricing>

# FLOYDHUB

	<b>Beginner</b>	<b>Data Scientist</b>	<b>Teams</b>
<b>INFRASTRUCTURE</b>			
GPU Access	Tesla K80, Tesla V100	Tesla K80, Tesla V100	Tesla K80, Tesla V100
Machine Priority	Low priority	High priority	High priority
Deep Learning Frameworks	TensorFlow, PyTorch, Keras, <a href="#">others</a>	TensorFlow, PyTorch, Keras, <a href="#">others</a>	TensorFlow, PyTorch, Keras, <a href="#">others</a>
Install Python Packages	✓	✓	✓
Customer Support	Forum / Email	Chat	Chat
Private Projects and Datasets		Unlimited	Unlimited
Public Projects and Datasets	Unlimited	Unlimited	All projects / datasets are private
<b>TRAINING</b>			
Training Duration	1 day	7 days	Custom
Storage	10 GB	100 GB	100 GB
Pricing (per month): <a href="https://www.floydhub.com/pricing">https://www.floydhub.com/pricing</a>	\$0	\$9	\$99



## Everything you need for deep learning

### ✓ Jupyter Notebooks

One-click GPU-powered interactive notebooks.

### ✓ Framework Support

TensorFlow, PyTorch, Keras, and more - fully-configured.

### ✓ Per Second Billing

Pay only for what you use, per second.

### ✓ Version Control

Full environment reproducibility.

### ✓ Parameter Sweeping

Run concurrent jobs for optimization.

### ✓ Public Datasets

Train using data from the open source community.

### ✓ Deploy Trained Models

Create model-serving REST API in a single command.

### ✓ Tensorboard

Visualize and explore your live experiments.



## GPU Plans

Dedicated Virtual Machines to Speed Up Complex Compute Jobs with RTX 6000 GPUs

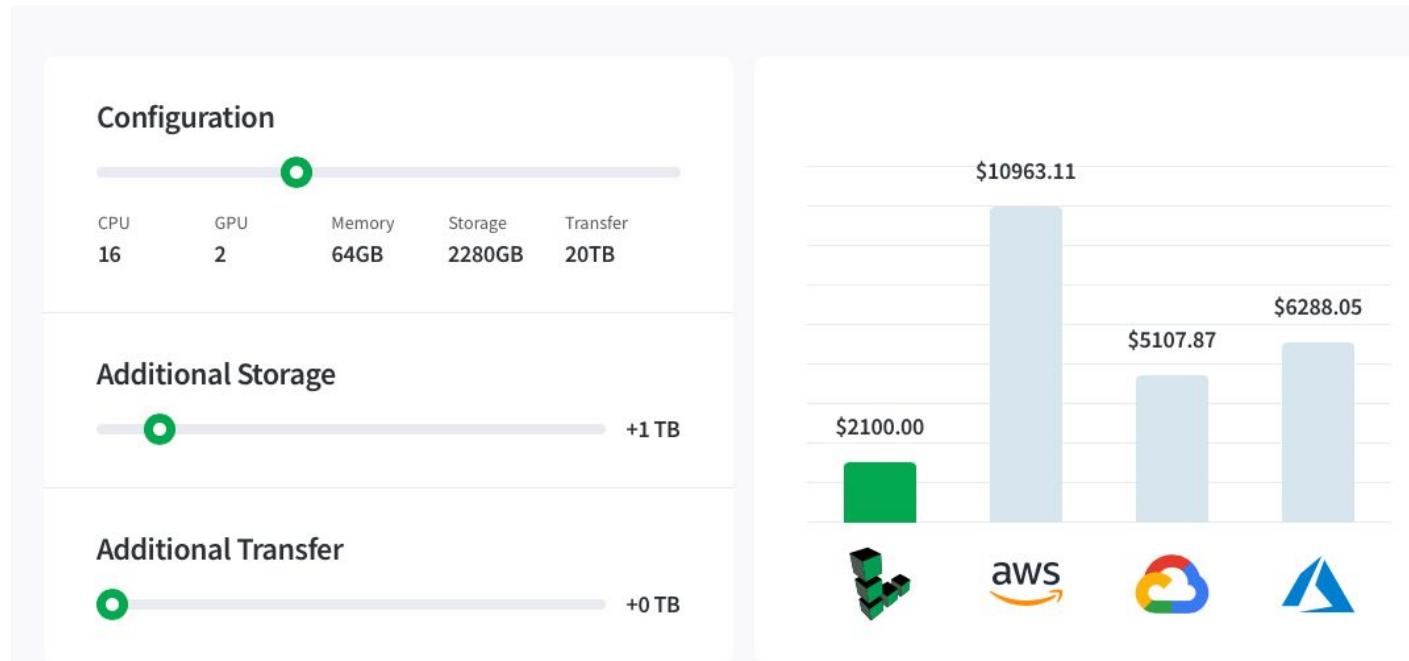
[Learn More >](#)

RAM	CPU	Storage	GPU Cards	Transfer	Network In	Network Out	Price	
<b>32 GB</b>	8 Cores	640 GB SSD	1	16 TB	40 Gbps	10000 Mbps	<b>\$1000 / mo</b>	(\$1.5 / hr)
<b>64 GB</b>	16 Cores	1280 GB SSD	2	20 TB	40 Gbps	10000 Mbps	<b>\$2000 / mo</b>	(\$3.0 / hr)
<b>96 GB</b>	20 Cores	1920 GB SSD	3	20 TB	40 Gbps	10000 Mbps	<b>\$3000 / mo</b>	(\$4.5 / hr)
<b>128 GB</b>	24 Cores	2560 GB SSD	4	20 TB	40 Gbps	10000 Mbps	<b>\$4000 / mo</b>	(\$6.0 / hr)

<https://www.linode.com/products/gpu/>



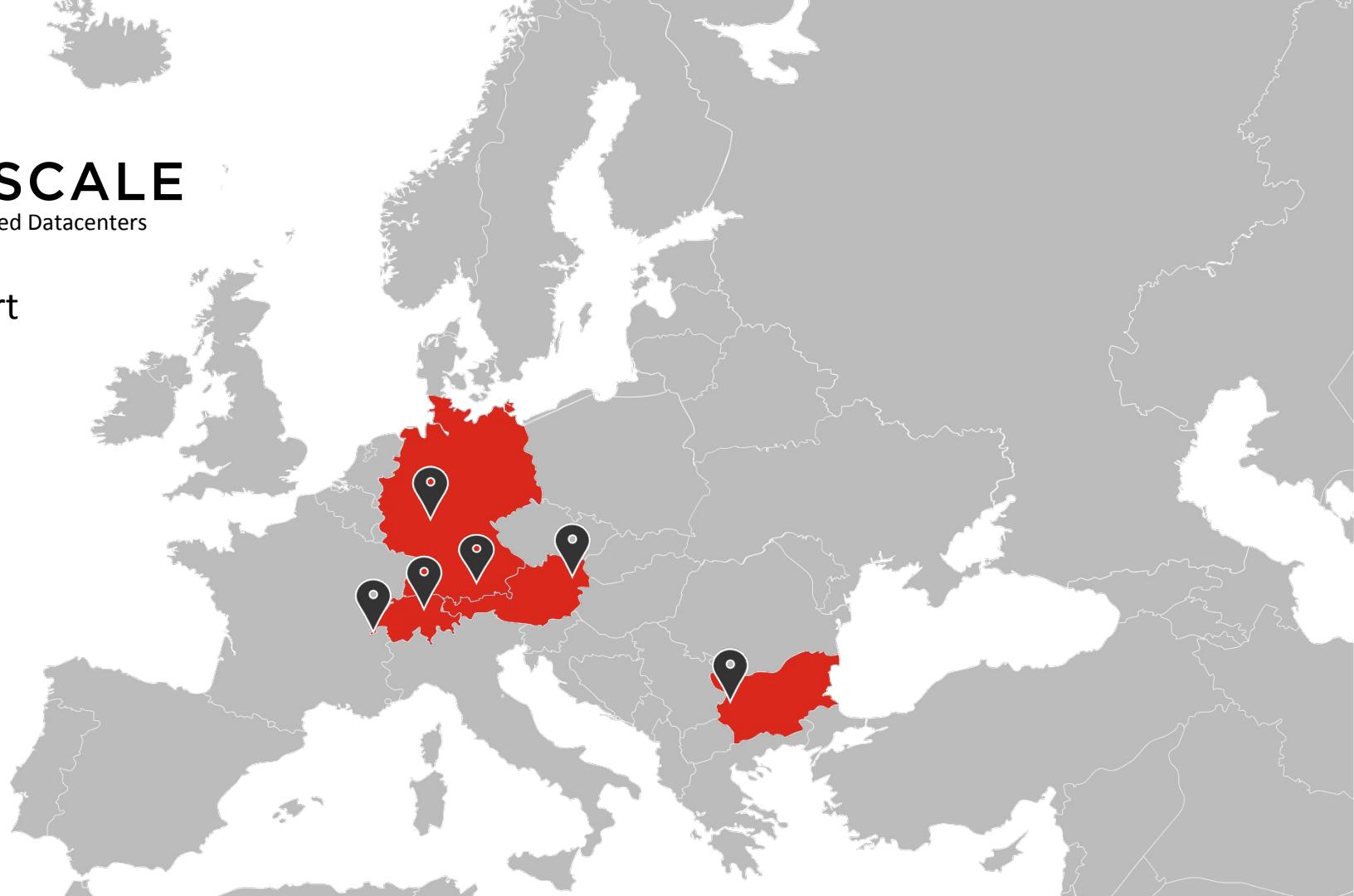
## Pricing Comparison (Example):



<https://www.linode.com>



Frankfurt  
Munich  
Vienna  
Zurich  
Geneva  
Sofia



# NVIDIA Tesla V100



- ▶ From 1 up to 4 **NVIDIA Tesla V100** per GPU2 instance
- ▶ 5,120 CUDA cores per card
- ▶ 640 dedicated Tensor Cores per card
- ▶ 16 GB of dedicated RAM per card
- ▶ Up to 48 CPUs per GPU2 instance
- ▶ Up to 1.6 TB of SSD local storage



	RAM	CPU Cores	GPU Cards	Price (EUR)
Small	56 GB	12 Cores	1 GPU	1.25015/hr
Medium	90 GB	16 Cores	2 GPU	1.83103/hr
Large	120 GB	24 Cores	3 GPU	2.41191/hr
Huge	225 GB	48 Cores	4 GPU	3.01803/hr

# NVIDIA Tesla P100



- ▶ From 1 up to 4 **NVIDIA Tesla P100** per GPU instance
- ▶ 3584 CUDA cores per card
- ▶ 16 GB of dedicated RAM per card
- ▶ Up to 48 CPUs per GPU instance
- ▶ Up to 1.6 TB of SSD local storage



	RAM	CPU Cores	GPU Cards	Price (EUR)
Small	56 GB	12 Cores	1 GPU	1.06200/hr
Medium	90 GB	16 Cores	2 GPU	1.55574/hr
Large	120 GB	24 Cores	3 GPU	2.04949/hr
Huge	225 GB	48 Cores	4 GPU	3.01803/hr

# Real-World Benchmarks



CIFAR10 Tensorflow Example

Instance type	Boot time	Processing Time	Costs / hour	Total Costs
Exoscale V100 GPU Instance	30 sec	10 sec	1.25015 EUR	~0.013 EUR
Exoscale Extra-Large	30 sec	5 mins	0.16164 EUR	~0.014 EUR
Our Office Laptop	Don't ask...	20 mins		Our Sanity

<https://github.com/tensorflow/models/tree/master/tutorials/image/cifar10>

# Why Exoscale?



## European

Built & Managed in Europe



## Transparent Billing

Simple, understandable cost structure



## Performance

Beats most cloud providers



## Compliant

GDPR, ISO27001, & more



## Affordable

30-40% lower costs than big hyperscalers



## Fast Startup Times

New Virtual Machines in under 30 seconds



## Automated

Terraform, Ansible, & more



## Free Internal Bandwidth

Even between zones!



## GPUs

Tesla V100 and P100 GPUs for Machine Learning



## SLA

99.95%



## No Vendor Lock-In

Standardized, vendor-independent systems.



## Dedicated Resources

Get real CPU cores and dedicated resources.

Interested?



Talk to this guy



## Janos Pasztor

Cloud Solution Architect @ Exoscale

[janos.pasztor@a1.digital](mailto:janos.pasztor@a1.digital)





# Paperspace

## Dedicated GPU

GPU+	P5000	P6000	V100
\$ 0.51 / hour + \$5 monthly storage fee	\$ 0.78 / hour + \$5 monthly storage fee	\$ 1.10 / hour + \$5 monthly storage fee	\$ 2.30 / hour + \$5 monthly storage fee
NVIDIA Quadro P4000 with 1792 CUDA cores.	NVIDIA Quadro P5000 with 2560 CUDA cores.	3840 CUDA cores and up to 432 GB/s memory bandwidth.	112 TeraFLOPs VOLTA chipset. The most powerful GPU in the world.
<ul style="list-style-type: none"><li>• 30GB RAM</li><li>• 8 x CPU</li><li>• 8 GB GPU <b>DEDICATED</b></li></ul>	<ul style="list-style-type: none"><li>• 30GB RAM</li><li>• 8 x CPU</li><li>• 16 GB GPU <b>DEDICATED</b></li></ul>	<ul style="list-style-type: none"><li>• 30GB RAM</li><li>• 8 x CPU</li><li>• 24 GB GPU <b>DEDICATED</b></li></ul>	<ul style="list-style-type: none"><li>• 30GB RAM</li><li>• 8 x CPU</li><li>• 16 GB GPU <b>DEDICATED</b></li></ul>

+ Collaboration Tools, Workflow, Templates etc.



# Paperspace

## GPU Instances

A range of GPU types for high-end workloads

New! TPU

advanced AI accelerator chip

K80

\$ 0.25 / hour

P100

\$ 0.59 / hour

V100

\$ 1.15 / hour

V100 x8

\$ 8.43 / hour

TPUv2

\$ 8.42 / hour

- 12GB GDDR5 **DEDICATED**
- 12GB RAM
- 2 vCPU
- 480 GB/s memory bandwidth
- 2,496 CUDA cores

- 16GB GDDR5 **DEDICATED**
- 24GB RAM
- 4 vCPU
- 732 GB/s memory bandwidth
- 3,584 CUDA cores

- 16GB GDDR5 **DEDICATED**
- 30GB RAM
- 8 vCPU
- 900 GB/s memory bandwidth
- 5,120 CUDA cores

- 128GB GDDR5 **DEDICATED**
- 130GB RAM
- 20 vCPU
- 7,200 GB/s memory bandwidth
- 40,960 CUDA cores

- 16GB **DEDICATED**
- 2,400 GB/s memory bandwidth
- 180 TeraFLOPs

<https://www.paperspace.com/pricing>

NEW Forbes covers launch of Gradient Multi-Cloud →

# Production tools for Machine Learning.

From exploration to deployment, Gradient enables individuals and teams to quickly develop, track, and collaborate on Machine Learning models of any size and complexity.

[GET STARTED](#)

[CONTACT SALES](#)





**SALE 20%**

on 8 x 2080 Ti 384 GB

~~€ 2 291,95 month~~ € 1 833,56 month

~~€ 572,99 week~~ € 458,39 week

~~€ 0,11 minute~~ € 0,08 minute

**ORDER NOW**



## 4 x 1080 Ti 128Gb RAM

RECOMMENDED

### GPU speed:



### Server Configuration:

GPU: 4 pcs 1080Ti  
(Each GPU card has:  
3584 CUDA cores  
GeForce GTX1080 Ti)

CPU: 2 x Intel Xeon E5-  
2609v4 1.7 GHz

RAM: 128 GB RAM

SSD: 480 GB SSD

Internal network: 10  
Gbps Port

### Best for:

» Blockchain processing  
» Universal GPU card

### OS:

- CentOS 7
- CentOS 8
- Ubuntu 16.04
- Ubuntu 18.04
- Windows 2016

### Prices:

- ✓ 893.7 € / month
- ✓ 223.43 € / week
- ✓ 0.04 € / minute
- ✓ 0 € - setup fee

(all prices without VAT)

For 1 MONTHS you will  
pay: 893.7 EURO

Available to run

ORDER NOW

## 2 x P100 PCI

### GPU speed:



### Server Configuration:

GPU: 2 pcs P100  
(Each GPU card has:  
3584 CUDA cores)

CPU: 2 x Intel Xeon E5-  
2609v4 1.7 GHz

RAM: 256 GB RAM

SSD: 960 GB SSD

Internal network: 10  
Gbps Port

### Best for:

» Deep Learning  
» AI Training  
» AI Inference  
» HPC

### OS:

- CentOS 7
- CentOS 8
- Ubuntu 16.04
- Ubuntu 18.04
- Windows 2016

### Prices:

- ✓ 1119.4 € / month
- ✓ 279.85 € / week
- ✓ 0.05 € / minute
- ✓ 0 € - setup fee

(all prices without VAT)

For 1 MONTHS you will  
pay: 1119.4 EURO

Available to run

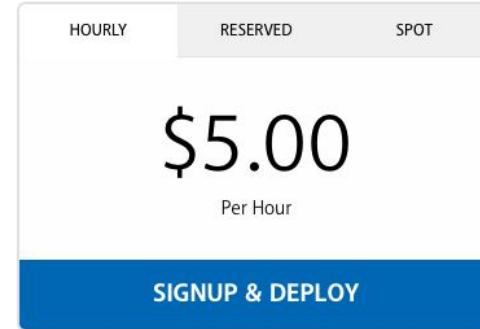
ORDER NOW

+ many more options: <https://www.leadergpu.com>

# packet

We all know that this is the GPU server we would build in our basement...if we had a small powerplant and a lot of cooling down there! Fast cores (thanks to some Intel Xeon 6126 processors), fast RAM, and of course dual Nvidia V100 32GB cards with NVLINK.

Hardware Specs	
CPU	24 Physical Cores @ 2.6 GHz (2 X XEON GOLD 6126)
Memory	192 GB of DDR4 ECC RAM
Storage	2 × 120 GB SSD boot
Storage	2 × 480GB SSD
GPU	2 x Nvidia V100 32GB GPU w/NVLINK



2x Nvidia V100

Currently only this GPU option. <https://www.packet.com/cloud/servers/g2-large/>



## GPU

Experience our most powerful public cloud instances, up to 1,000 times faster than a CPU for parallel processing

[Find out more](#)

Name	Memory	vCore	GPU	Storage	Public network	Private network	Price
t1-45	45 GB	8	Tesla V100 16 GB	400 GB SSD	2 Gbps guaranteed	4 Gbps max.	\$2.661 /hour
t1-90	90 GB	18	2×Tesla V100 16 GB	800 GB SSD	4 Gbps guaranteed	4 Gbps max.	\$5.44 /hour
t1-180	180 GB	32	4×Tesla V100 16 GB	50 GB SSD + 2 TB	10 Gbps	4 Gbps max.	\$10.88 /hour

Up to 1,000 times faster than a CPU on parallel computing, NGC integration makes it easy to use Tensorflow, Caffe, Mxnet, and much more.

<https://www.ovhcloud.com>



Each Lambda gpu.4x instances has four (4) vGPUs with 11 GB of VRAM, 3584 CUDA Cores, and performance equivalent to 4x 1080 Ti GPUs. For exact details see below:

Instance Type	GPUs	GPU Memory (GB)	vCPUs	System Memory (GB)	Instance Speed	On-Demand Hourly Rate	1-yr Reserved Effective Hourly Rate
Lambda gpu.4x	4 GPUs	11 GB / GPU	8 vCPU Cores @ 3.5 GHz	32 GB DDR4	Up to 10 Gbps	\$1.50	\$1.25

- GPU servers optimized for deep learning
- On-demand Pricing: \$0.375 / hour / GPU, Reserved Pricing: \$0.3125 / hour / GPU
- TensorFlow, PyTorch, Keras, and Caffe preinstalled on every server
- Less than half the cost of AWS and Azure

<https://lambdalabs.com/service/lambda-gpu-cloud-faq>

# EX51-SSD-GPU KONFIGURATOR

Startseite > EX51-SSD-GPU > Konfigurator

## EX51-SSD-GPU BASISKONFIGURATION

CPU	Intel® Core™ i7-6700 Quad-Core	inkl. Hyper-Threading-Technologie
RAM	64 GB DDR4	
HARD DRIVES	2x 500 GB SATA SSD	Software-RAID 1
GRAFIKKARTE	GeForce® GTX 1080	

## EX51-SSD-GPU DEDICATED SERVER

### SERVER STANDORT

Server Standort zurücksetzen

## KONFIGURATION

	Anzahl	Preis
Basiskonfiguration in Helsinki	1	111,86 €

Serveranzahl:

1

**111,86 €**

monatlich + einmalig Setup: 0,00 €

Preis inkl. 19 % USt.,

[Verfügbarkeit prüfen >](#)

DERZEIT NICHT LIEFERBAR!

Note: only 1 NVidia GTX 1080 GPU per instance possible

<https://www.hetzner.de/dedicated-rootserver/ex51-ssd-gpu/konfigurator>

# vast.ai

GPU Compute  
“Search Engine”  
for GPU power  
from different  
vendors

## Instance Configuration

Image: tensorflow/tensorflow

Image CUDA version: 10.1  
Offers with incompatible cuda version hidden

Launch Type: Jupyter  
On-start script: Not set

[EDIT IMAGE & CONFIG...](#)

Disk Space To Allocate  
**10.00 GB**

## Filter offers

### Availability

Host Reliability

90%

Max Instance Duration

3 days

Include Unavailable Offers

Include External Offers

Include Unverified Machines

Include Incompatible Machines

### Price

\$/Hour

\$0.00

\$530.00

TFLOPS/\$/Hour

1.00

500.00

Show  interruptible •  on-demand pricing

Auto

474441	838	TB250	↑259.1 Mbps	\$0.182/hr
V	2x GTX 1080 Ti	PCIE 2.0, 1x0.3 GB/s	19.7 DLPerf	Reliability 98.6%

479602	1326	X399 Taichi	↑186.9 Mbps	\$0.387/hr
V	2x RTX 2080 Ti	PCIE 3.0, 8x6.3 GB/s	39.3 DLPerf	Reliability 99.12%

479600	1326	X399 Taichi	↑186.9 Mbps	\$0.767/hr
V	4x RTX 2080 Ti	PCIE 3.0, 8x6.3 GB/s	77.6 DLPerf	Reliability 99.12%

457434	1596	MAXIMUS VIII GE...	↑24.3 Mbps	\$0.113/hr
V	2x GTX 1070	PCIE 3.0, 8x6.1 GB/s	10.8 DLPerf	Reliability 99.79%

479232	1701	X9DRX	↑56.8 Mbps	\$0.202/hr
V	2x GTX 1080 Ti	PCIE 3.0, 8x5.8 GB/s	21.2 DLPerf	Reliability 97.4%

474439	838	TB250	↑259.1 Mbps	\$0.362/hr
V	4x GTX 1080 Ti	PCIE 2.0, 1x0.3 GB/s	33.6 DLPerf	Reliability 98.6%

408336	1896	PRIME Z270	↑9.3 Mbps	\$0.240/hr
V	2x RTX 2070	PCIE 3.0, 8x6.1 GB/s	18.3 DLPerf	Reliability

<https://vast.ai>

# Who provides which GPU types?

Vienna



# Deep Learning Meetup

## Deep Learning Hardware Overview: What and where to buy or rent



Jan Schlueter  
JKU Linz



René Donner  
contextflow



Thomas Lidy  
Musimap