

Data Science Survival Skills

Hardware relevance in data science

Agenda

- Why you should even bother
- CPUs
- GPUs
- TPUs
- Hardware accelerators

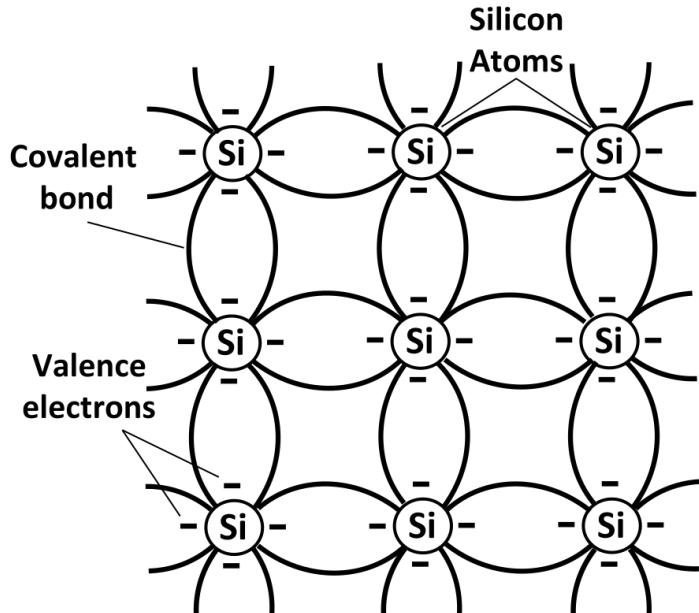
A computer from the inside

BACK

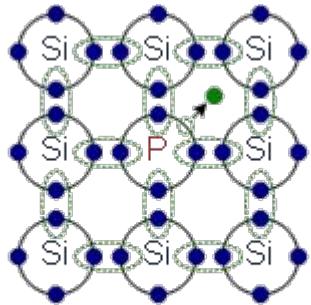
FRONT



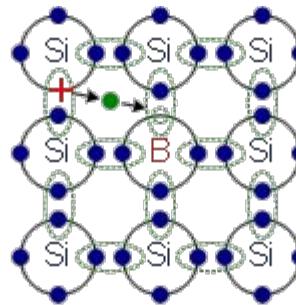
A computer is based on chemistry and physics!



Doping

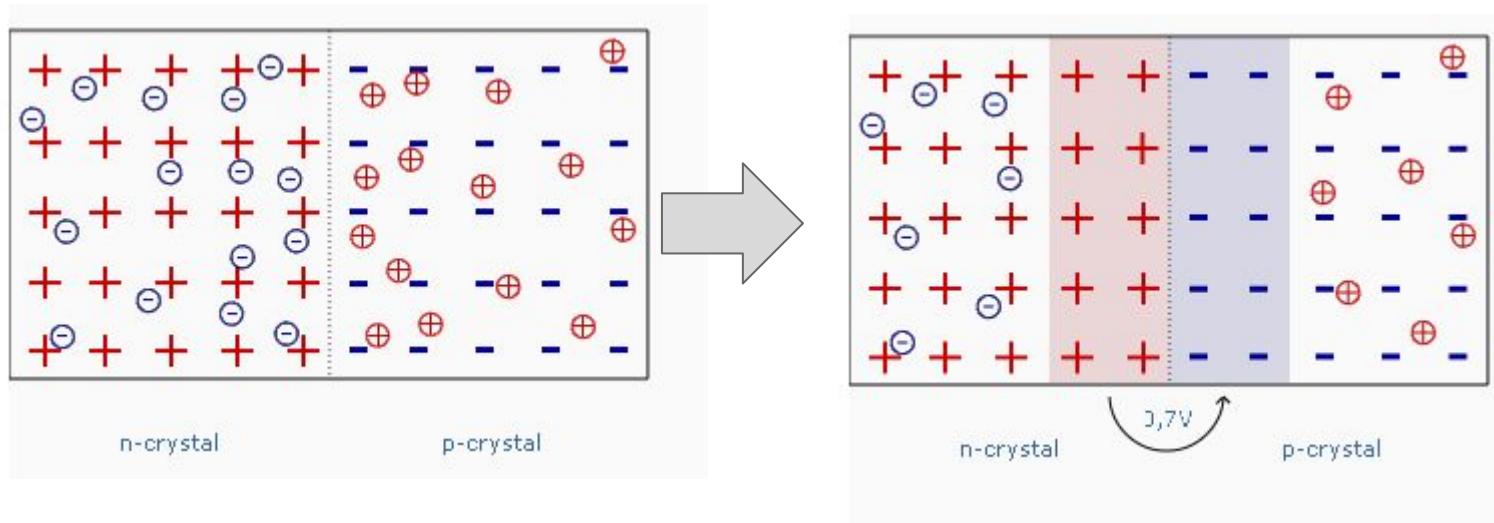


The phosphorus atom donates 1st fifth valence electron. It acts as a free charge carrier.



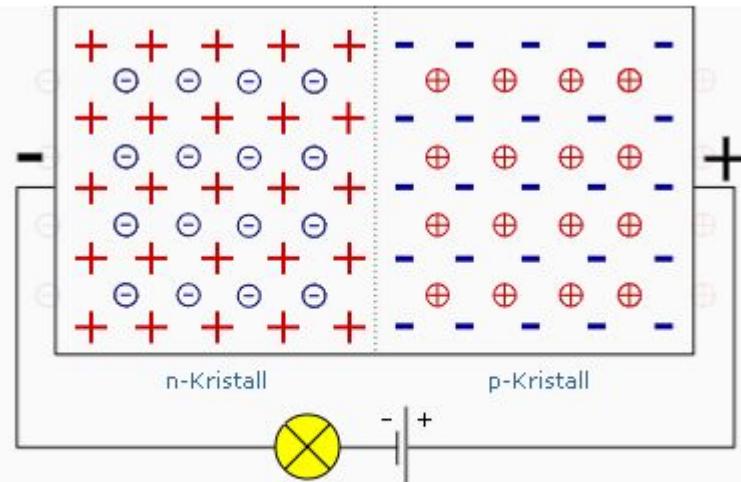
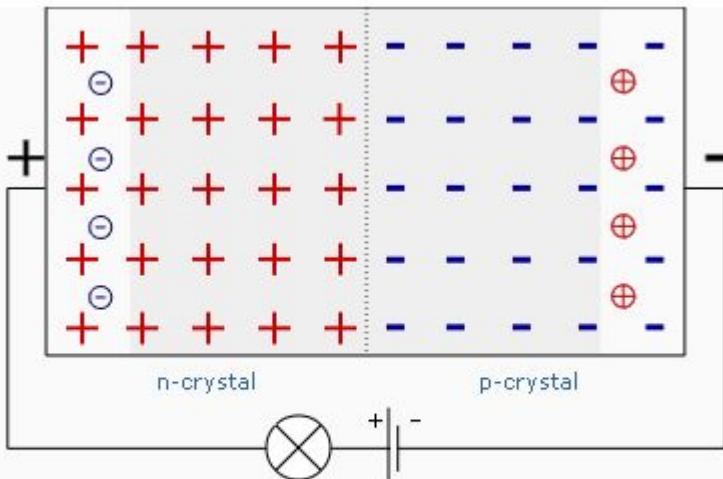
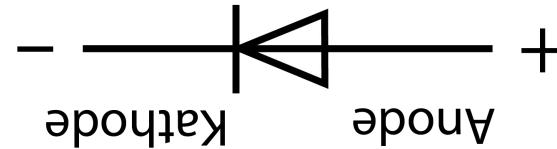
The free place on the boron atom is filled with an electron. Therefore a new hole ("defect electron") is generated. This holes move in the opposite direction to the electrons

A diode

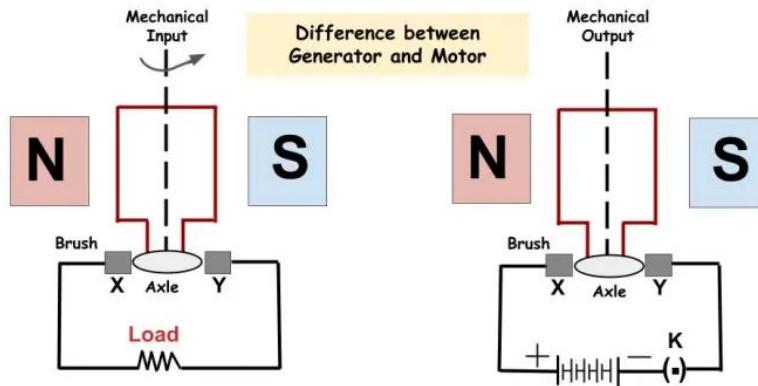


A diode

Durchlassrichtung
←

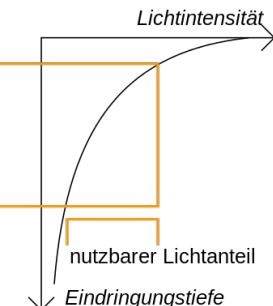
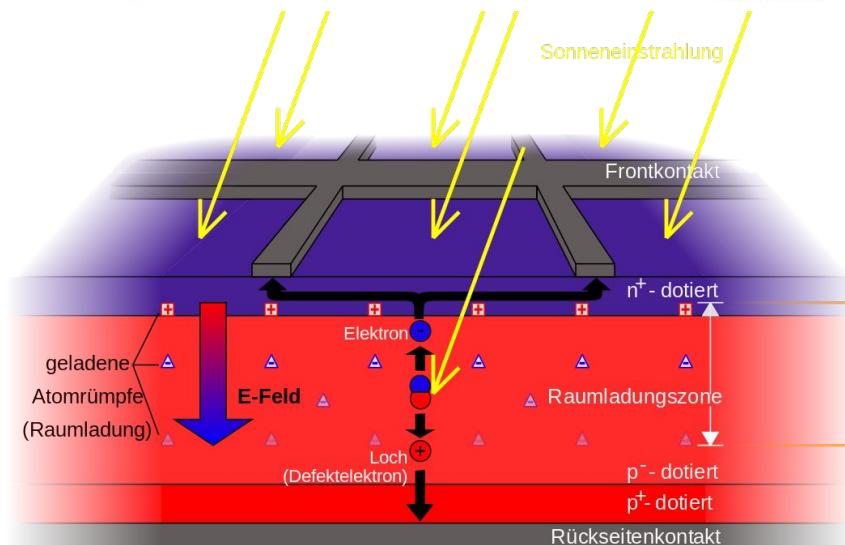
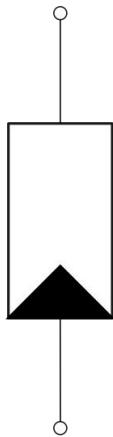


A solar cell

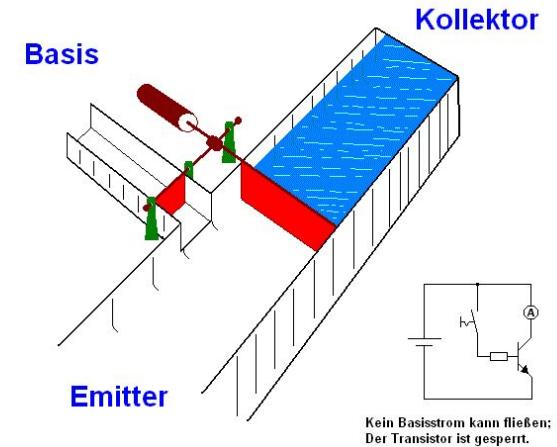
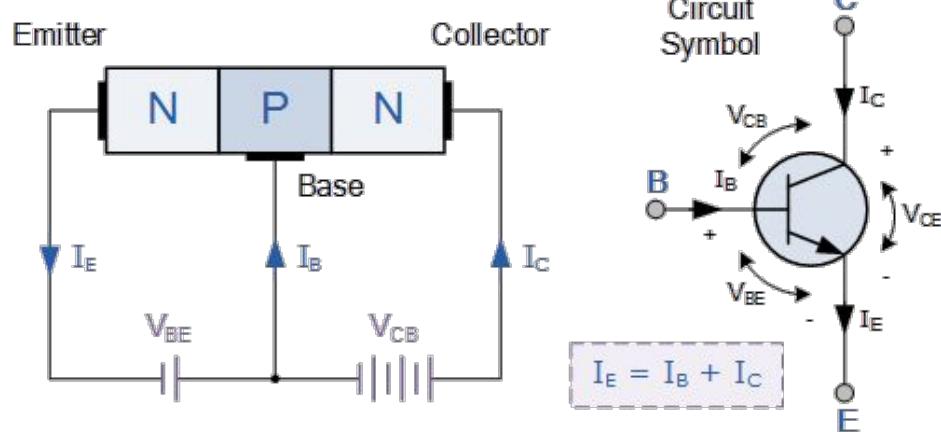
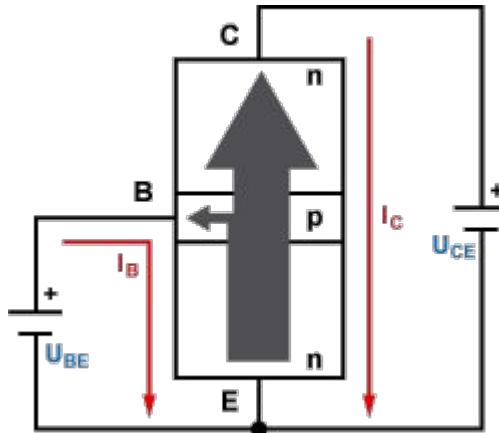
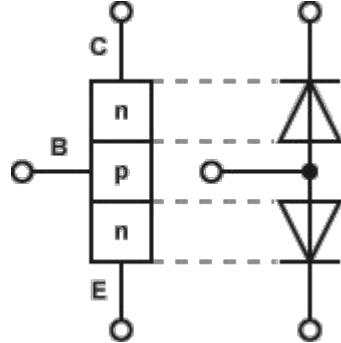


Generator

Motor DewWool.com

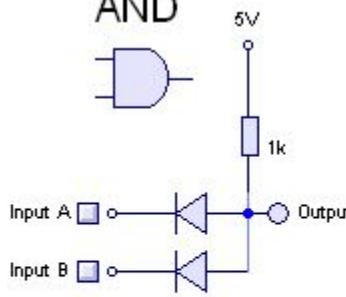


A transistor

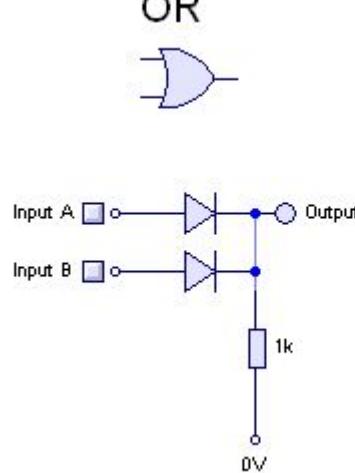


Transistor can perform computations

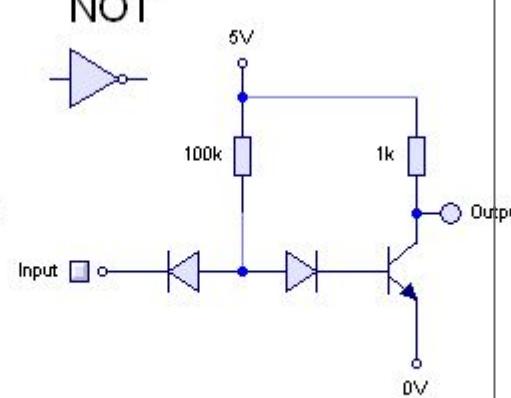
AND



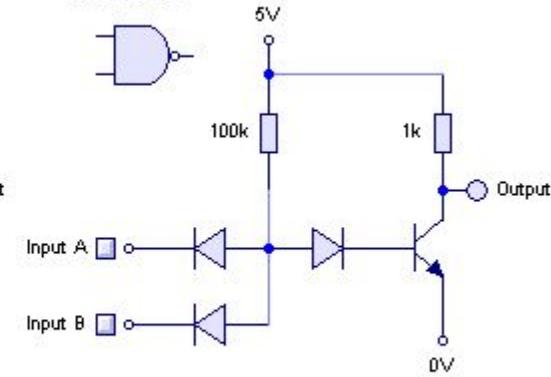
OR



NOT

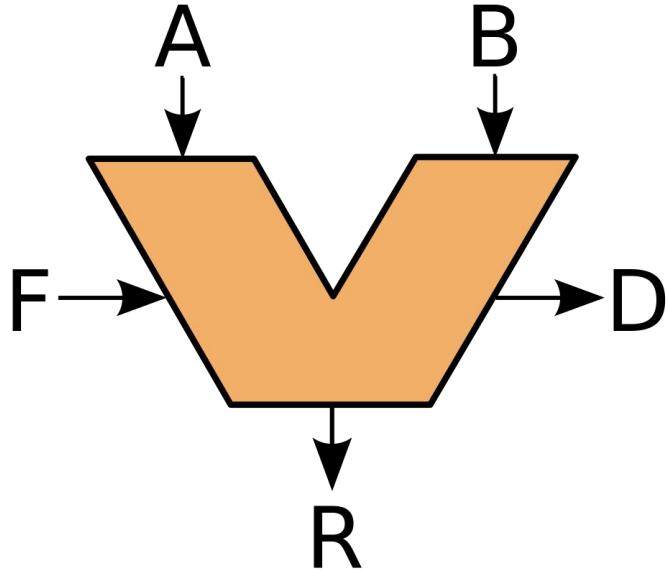
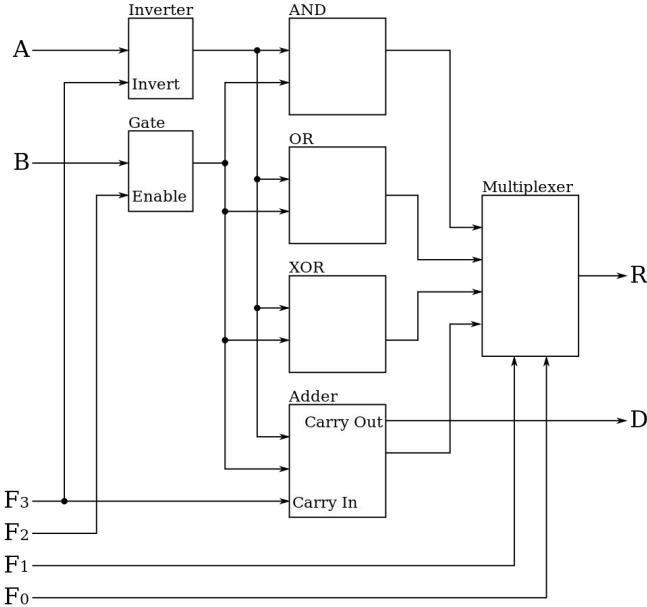


NAND

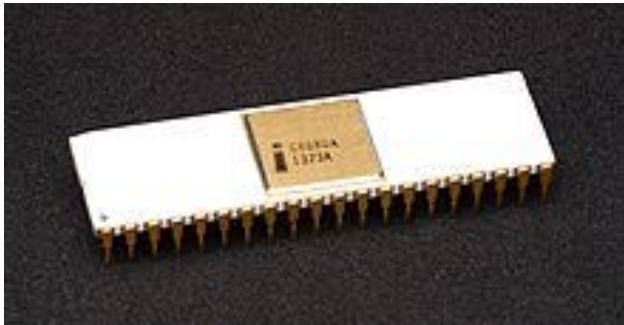


NAND is logic complete, i.e. can be used to gain every other function!

Arithmetic Logic Unit (ALU)



A modern CPU (central processing unit)



Intel's 8 bit processor (8080, 1974),
6000 Transistors, 6 um per transistor

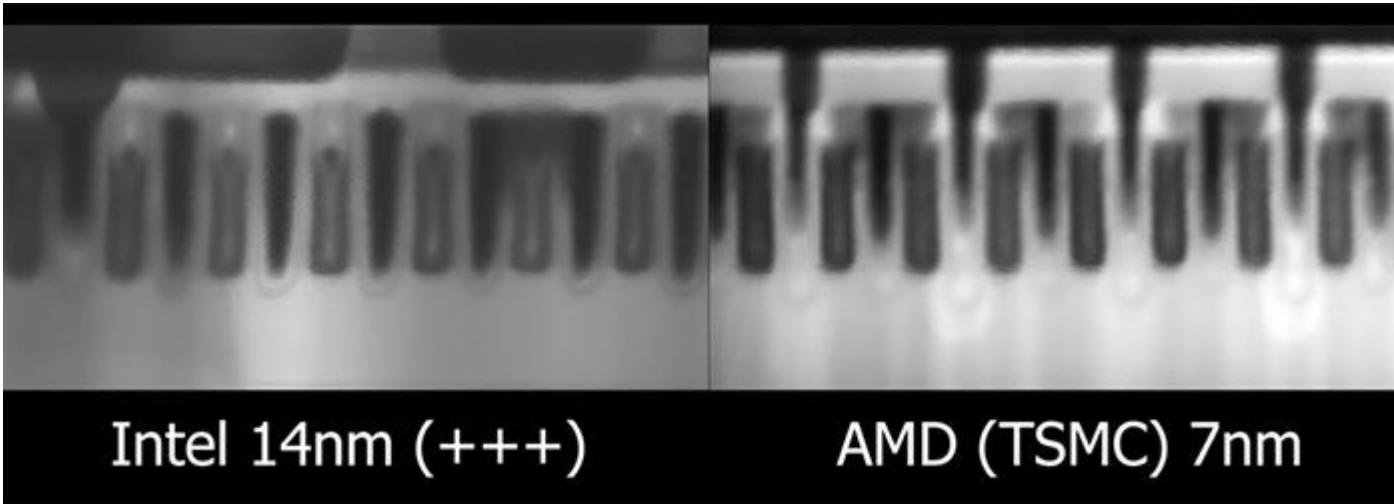


60 MHz
0.8 um (transistor)
3.1M transistors
1993



AMD Athlon
First consumer CPU > 1 GHz
End 90s, beginning 2000s
0.25 um - 0.13 um

Computations using transistors



Nowadays CPU

Desktop Ryzen 9 [Bearbeiten | Quelltext bearbeiten]

AMD



Lagernd
CUSTOMERS' CHOICE

AMD Ryzen 7 5800X 8x 3.80GHz So.AM4 WOF
€ 359,-*
inkl. 19% USt + Versandkosten
★ Mindstar Highlight ★



Lagernd
Artnr.: 74531
inkl. 19% USt + Versandkosten
über 7.080 verkauft



Lagernd
CUSTOMERS' CHOICE
inkl. 19% USt + Versandkosten
über 23.490 verkauft



über 7.080 verkauft



über 7.080 verkauft



über 7.080 verkauft

| Modellnummer | Ordering Part Number | | Mikroarchitektur | Fab | Soc-Kel | Kerne/Threads | Cache | | | Prozessortakt / Maximaltakt [ann 2] | Speicher-Controller [ann 1] | TDP | Einführungsdatum |
|-------------------|----------------------|------------------|------------------|------|---------|---------------|-------|-------|---------------|-------------------------------------|-----------------------------|----------------|------------------|
| | Tray | PIB | | | | | L1 | L2 | L3 | | | | |
| Ryzen 9 3900 | 100-000000070 | - | Zen 2 | AM4 | 12 / 24 | 768 KB | 6 MB | 64 MB | 3,1 / 4,3 GHz | Dual-Channel DDR4-3200 | 65 W | 24. Sep. 2019 | |
| Ryzen 9 Pro 3900X | 100-000000072 | - | | | | | | | | | | | |
| Ryzen 9 3900X | 100-000000023 | 100-100000023BOX | | | | | | | | | | | |
| Ryzen 9 3900XT | 100-00000277 | 1000000277WOF | | TSMC | 7FF | 16 / 32 | 1 MB | 8 MB | 3,5 / 4,7 GHz | 105 W | 7. Juli 2021 | | |
| Ryzen 9 3950X | 100-00000051 | 100-100000051WOF | | | | | | | | | | | |
| Ryzen 9 5900X | 100-00000062 | - | Zen 3 | AM4 | 12 / 24 | 768 KB | 6 MB | 64 MB | 3,0 / 4,7 GHz | Dual-Channel DDR4-3200 | 65 W | 13. April 2021 | |
| Ryzen 9 5900X | 100-00000061 | 100-100000061WOF | | | | | | | | | | | |
| | | | | | | | | | | | | | |



Lagernd
Artnr.: 75251

Intel Core i9 11900K 8x 3.50GHz So.1200 WOF

€ 519,-*

inkl. 19% USt + Versandkosten

über 570 verkauft



Lagernd
Artnr.: 74143

Intel Core i7 10700K 8x 3.80GHz So.1200 WOF

€ 322,-*

inkl. 19% USt + Versandkosten

über 7.640 verkauft



Lagernd
Artnr.: 75252

Intel Core i5 11400F 6x 2.70GHz So.1200 WOF

€ 158,97*

inkl. 19% USt + Versandkosten

über 3.000 verkauft



Die Modelle der Rocket Lake-Generation (komplett)

| Modell | Release | Cores (Threads) | PCIe Lanes | Takt | | TDP | L3-Cache | Bemerkungen |
|-----------------|---------|-----------------|------------|---------|---------|-------|----------|-------------|
| | | | | Base | Turbo | | | |
| Core™ i9-11900K | Mrz. 21 | 8 (16) | 20 | 3,5 GHz | 5,2 GHz | 125 W | 16 MB | B |
| Core™ i7-11700K | | 8 (16) | | 3,6 GHz | 5,0 GHz | | 16 MB | |

€ 427,76*

inkl. 19% USt + Versandkosten

über 2.900 verkauft

€ 248,97*

inkl. 19% USt + Versandkosten

über 1.200 verkauft

CPU

| | | | | | | | | | | | | | |
|-------------------------------|---------------|--------------|-------|-------------|----------|---------|--------|--------|---------------|-------------------------------|-------------------------------|---------------|------------------|
| Ryzen Threadripper 3960X | - | 100000010WOF | Zen 2 | TSMC 7FF | sTRX4 | 24 / 48 | 1,5 MB | 12 MB | 128 MB | 3,8 / 4,5 GHz | Quad- Channel DDR4-3200 | 280 W | 25. Nov. 2019 |
| Ryzen Threadripper 3970X | - | 100000011WOF | | | 32 / 64 | 2 MB | 16 MB | MB | 3,7 / 4,5 GHz | | | | |
| Ryzen Threadripper 3990X | | 100000163WOF | | | 64 / 128 | 4 MB | 32 MB | | 2,9 / 4,3 GHz | | | | |
| Ryzen Threadripper Pro 3945WX | 100-000000168 | - (OEM only) | | WRX80 | 12 / 24 | 768 KB | 6 MB | 64 MB | 4,0 / 4,3 GHz | Octo- Channel DDR4-3200 | | 14. Juli 2020 | |
| Ryzen Threadripper Pro 3955WX | 100-000000167 | - (OEM only) | | | 16 / 32 | 1 MB | 8 MB | 64 MB | 3,9 / 4,3 GHz | | | | |
| Ryzen Threadripper Pro 3975WX | 100-000000086 | - (OEM only) | | | 32 / 64 | 2 MB | 16 MB | 128 MB | 3,5 / 4,2 GHz | | | | |
| Ryzen Threadripper Pro 3995WX | 100-000000087 | - (OEM only) | | | 64 / 128 | - MB | 32 MB | 256 MB | 2,7 / 4,2 GHz | | | | |



AMD Ryzen Threadripper PRO 3995WX, 64C/128T, 2.70-4.20GHz, boxed ohne Kühler (100-100000087WOF)

★★★★★ jetzt bewerten!

| | |
|----------------------|-------------------------------------|
| Sockel | sWRX8 (LGA) |
| Codename | Castle Peak |
| Grafik | nein |
| TDP | 280W |
| Kerne | 64 |
| Threads | 128 |
| Basistakt | 2.70GHz |
| Turbotakt | 4.20GHz |
| SMT | ja |
| Speichercontroller | Octa Channel PC4-25600L (DDR4-3200) |
| ECC-Unterstützung | ja |
| Freier Multiplikator | nein |

▼ Alle Produkteigenschaften anzeigen

Info beim Hersteller

Aktueller Preisbereich

€ 5079,00 bis € 7392,29

Preisentwicklung 1W 1M 3M 6M 1J



Preisentwicklung öffnen

Preisalarm setzen

Zur Wunschliste hinzufügen

Zur Vergleichsliste hinzufügen

Feedback senden

Server CPUs



Verfügbar
Artnr: 8972690

AMD Epyc 7F32 8x 3.70GHz So.SP3
TRAY

€ 2.109,47*

inkl. 19% USt + [Versandkosten](#)



Verfügbar
Artnr: 8978613

AMD Epyc 7F52 16x 3.50GHz
So.SP3 TRAY

€ 3.426,95*

inkl. 19% USt + [Versandkosten](#)

über 5 verkauft



Verfügbar
Artnr: 8984108

AMD Epyc 7H12 64x 2.60GHz
So.SP3 TRAY

€ 5.900,86*

inkl. 19% USt + [Versandkosten](#)



Verfügbar
Artnr: 8847113

AMD Epyc 7401P 24x 2.00GHz
So.SP3 TRAY

€ 1.064,86*

inkl. 19% USt + [Versandkosten](#)

über 10 verkauft



Verfügbar
Artnr: 9029943

AMD EPYC MILAN 48-CORE 7643
2.3GHZ

€ 5.421,50*

inkl. 19% USt + [Versandkosten](#)

CPU - what for?!

- Mobile vs. desktop
- Frequency (clock rate, how fast it runs)
- Thermal power design (TPD) in Watts -> how much power your PC consumes/heat it produces
- L2/L3 cache (fast cache for operations, up to several MBs)
- On board GPU (only for checking emails...)
- General purpose!

One core vs multiple cores

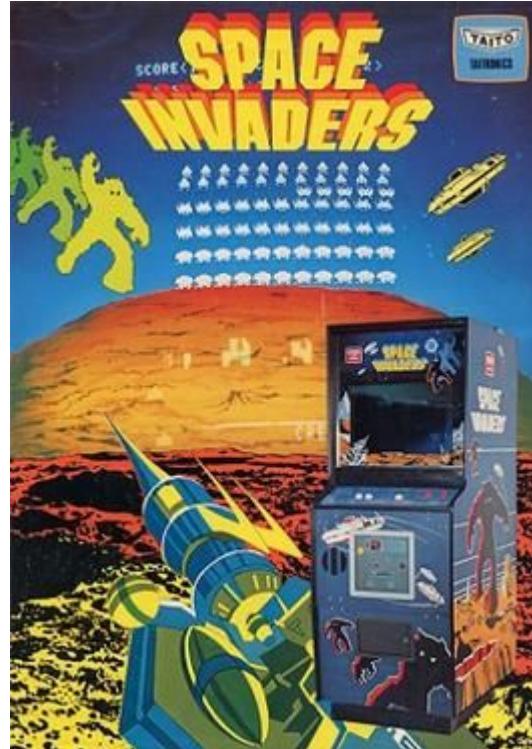
- Single core performance on Intel mostly better than on AMD
 - most code you write
- Multiple core performance (AMD > Intel)
 - multithreading (simultaneous runs of programs without stopping the other one)
 - multiprocessing (computations distributed across cores)
 - ⇒ Later lectures!
- In general, in a scientific environment it depends ON YOUR TASK how many cores you need.

Fluid dynamics: lots of cores → cpu cluster

Deep learning: lots of ???

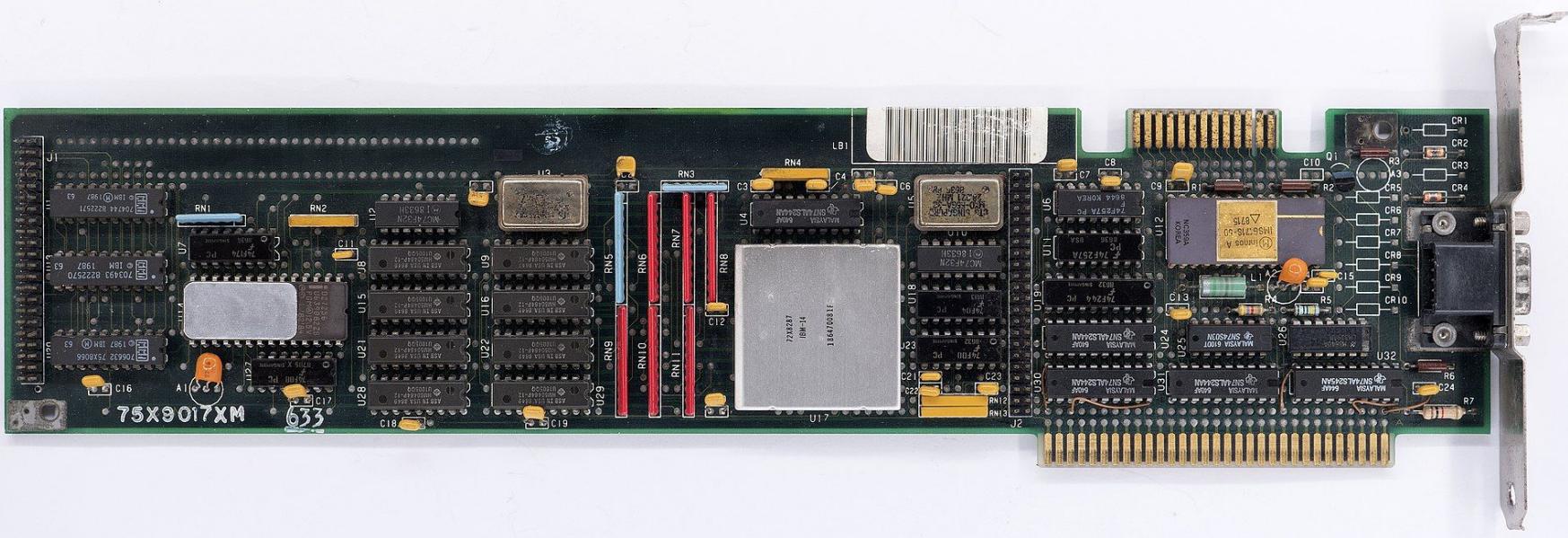
The GPU

Graphics Processing Unit
→ started in the 1970s to play video games



IBM's VGA (Video Graphics Array)

Introduced 1987 → VGA resolution is 640x480



GPUs these days

GIGABYTE™



Lagernd
Artnr: 70637

8GB MSI GeForce RTX 3070
GAMING Z TRIO LHR DDR6 (Retail)

€ 999,-*

inkl. 19% USt, [Gratisversand](#)

über 580 verkauft



Lagernd
Artnr: 74407

24GB MSI GeForce RTX 3090
GAMING X TRIO Aktiv PCIe 4.0 x16

€ 2.499,-*

inkl. 19% USt + [Versandkosten](#)

über 700 verkauft



Lagernd
Artnr: 70334

8GB Gigabyte GeForce RTX 3070
GAMING OC 8G 2.0 LHR

€ 949,-*

inkl. 19% USt, [Gratisversand](#)

★ Mindstar Highlight ★



Lagernd
Artnr: 70988

12GB KFA2 GeForce RTX 3080 Ti
Hall Of Fame 1-Click OC GDDR6X

€ 1.949,-*

inkl. 19% USt + [Versandkosten](#)

über 5 verkauft



Lagernd
Artnr: 70127

8GB Gigabyte GeForce RTX 3070 Ti
AORUS Master Aktiv PCIe 4.0 x16)

€ 1.129,-*

inkl. 19% USt + [Versandkosten](#)

über 280 verkauft



Lagernd



Lagernd



Lagernd



Lagernd



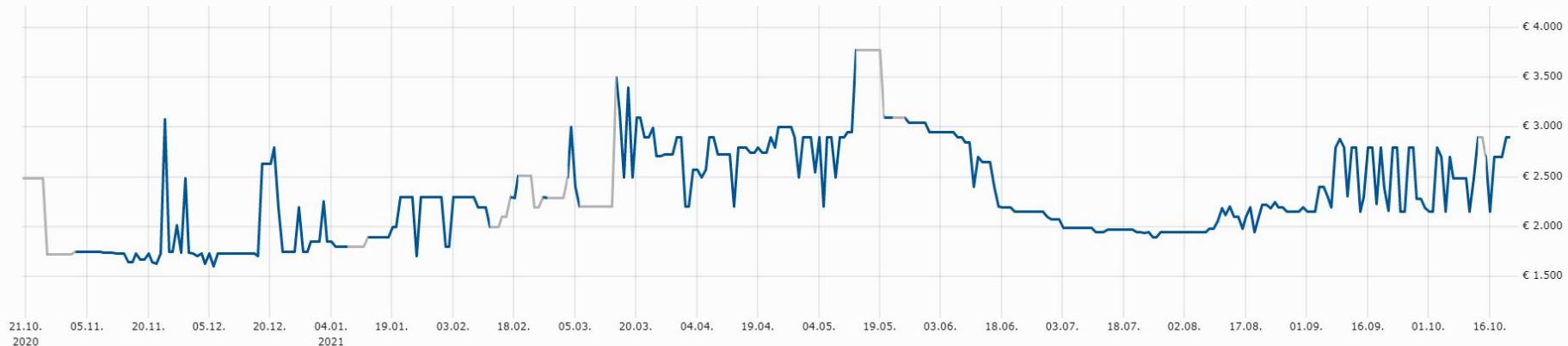
Lagernd

GPU price

Preisentwicklung

Preisentwicklung von Händlern aus DE bis 21.10.2021 (letzter Preis: € 2.899,00)
Keine Angebote in diesem Zeitraum

▼ Optionen Preisalarm setzen
Zeitraum: 1M 3M 6M **1J** Max



Hinweis: Ausfälle von Händler-Websites und andere technische Probleme können zu ungewöhnlichen Preisschwankungen führen!

RTX 3090 UVP: 1500 EUR

Bitcoin etc mining



**ICELAND'S
BITCOIN
SUPER-MINE**



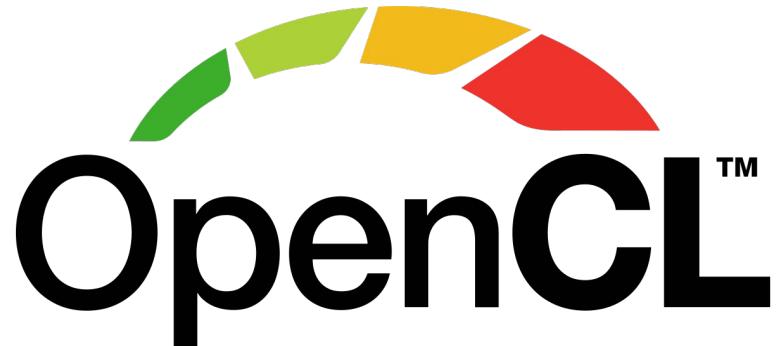
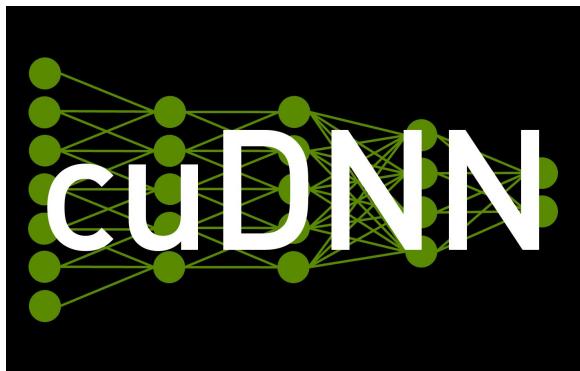
GPU terminology

Usage specific GPU[edit]

Most GPUs are designed for a specific usage, real-time 3D graphics or other mass calculations:

1. Gaming
 - GeForce GTX, RTX
 - Nvidia Titan
 - Radeon HD, R5, R7, R9, RX, Vega and Navi series
 - Radeon VII
2. Cloud Gaming
 - Nvidia GRID
 - Radeon Sky
3. Workstation
 - Nvidia Quadro
 - Nvidia RTX
 - AMD FirePro
 - AMD Radeon Pro
4. Cloud Workstation
 - Nvidia Tesla
 - AMD FireStream
5. Artificial Intelligence training and Cloud
 - Nvidia Tesla
 - AMD Radeon Instinct
6. Automated/Driverless car
 - Nvidia Drive PX

AMD vs NVIDIA?!



General purpose, but 30% slower on
NVIDIA GPUs...

What do GPUs do?

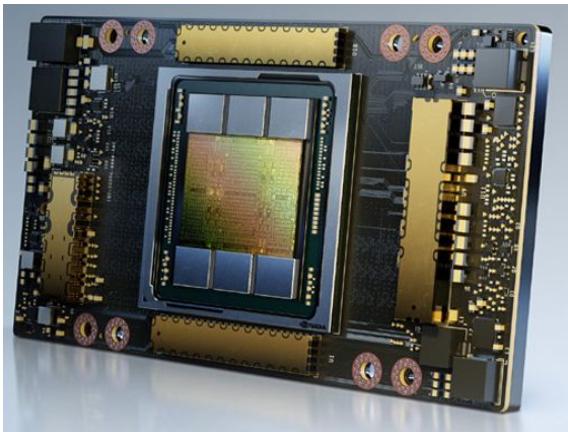
- Highly parallel processing
- Dedicated ICs for video encoding and decoding
- Ray tracing/rendering
- Matrix multiplications!

GPU specs (example RTX 3090)

| Graphics Processor | |
|--------------------|------------------------|
| GPU Name: | GA102 |
| GPU Variant: | GA102-300-A1 |
| Architecture: | Ampere |
| Foundry: | Samsung |
| Process Size: | 8 nm |
| Transistors: | 28,300 million |
| Die Size: | 628 mm ² |

| Theoretical Performance | Memory | Render Config |
|--|-----------------------------|---------------------------|
| Pixel Rate: 189.8 GPixel/s | Memory Size: 24 GB | Shading Units: 10496 |
| Texture Rate: 556.0 GTexel/s | Memory Type: GDDR6X | TMUs: 328 |
| FP16 (half) performance: 35.58 TFLOPS (1:1) | Memory Bus: 384 bit | ROPs: 112 |
| FP32 (float) performance: 35.58 TFLOPS | Bandwidth: 936.2 GB/s | SM Count: 82 |
| FP64 (double) performance: 556.0 GFLOPS (1:64) | Graphics Features | Tensor Cores: 328 |
| Graphics Card | DirectX: 12 Ultimate (12_2) | RT Cores: 82 |
| Release Date: Sep 1st, 2020 | OpenGL: 4.6 | L1 Cache: 128 KB (per SM) |
| Availability: Sep 24th, 2020 | OpenCL: 3.0 | L2 Cache: 6 MB |
| Generation: GeForce 30 | Vulkan: 1.2 | |
| Predecessor: GeForce 20 | CUDA: 8.6 | |
| Production: Active | Shader Model: 6.6 | |
| Launch Price: 1,499 USD | | |
| Current Price: Amazon / Newegg | | |
| Bus Interface: PCIe 4.0 x16 | | |
| Reviews: 56 in our database | | |

The real power horses



OPTIMIZED SOFTWARE AND SERVICES FOR ENTERPRISE



EVERY DEEP LEARNING FRAMEWORK

mxnet

PYTORCH



NVIDIA A100 TENSOR CORE GPU SPECIFICATIONS (SXM4 AND PCIe FORM FACTORS)

| | A100 40GB PCIe | A100 80GB PCIe | A100 40GB SXM | A100 80GB SXM |
|--------------------------------|--------------------|--------------------------|--------------------|---------------------|
| FP64 | | | 9.7 TFLOPS | |
| FP64 Tensor Core | | | 19.5 TFLOPS | |
| FP32 | | | 19.5 TFLOPS | |
| Tensor Float 32 (TF32) | | 156 TFLOPS 312 TFLOPS* | | |
| BFLOAT16 Tensor Core | | 312 TFLOPS 624 TFLOPS* | | |
| FP16 Tensor Core | | 312 TFLOPS 624 TFLOPS* | | |
| INT8 Tensor Core | | 624 TOPS 1248 TOPS* | | |
| GPU Memory | 40GB HBM2 | 80GB HBM2e | 40GB HBM2 | 80GB HBM2e |
| GPU Memory Bandwidth | 1,555GB/s | 1,935GB/s | 1,555GB/s | 2,039GB/s |
| Max Thermal Design Power (TDP) | 250W | 300W | 400W | 400W |
| Multi-Instance GPU | Up to 7 MIGs @ 5GB | Up to 7 MIGs @ 10GB | Up to 7 MIGs @ 5GB | Up to 7 MIGs @ 10GB |
| Form Factor | PCIe | | SXM | |

A100



Preisentwicklung

Preisentwicklungsdaten
Keine Angebote

EANnummer 74392 EAN 3536403378035 SKU TCSA100M-PB



40GB PNY A100 PCIe, HBM2 (TCSA100M-PB)

Verfügbar
(Max. Bestellmenge: 2)

[0 Bewertungen](#)

nur € 10.420,36*

* inkl. 19% USt zzgl. [Versandkosten / Lieferbeschränkungen](#)

+ € 3,90* [Geschenkverpackung](#)



1

In den Warenkorb

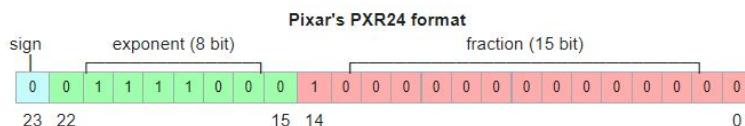
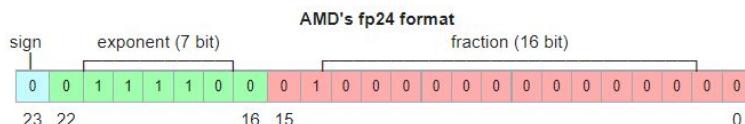
Auf den Wunschzettel

- PaketShop-Lieferung möglich
- Versandkosten sparen
- ⓘ Feedback zum Artikel
- 🖨 Artikelinfos drucken

[Preisalarm setzen](#)

3M 6M 1J Max

Floating point precision



INT8

Table 2: ResNet-50 ImageNet validation set accuracy per math type

| Math type | Multiply-add type | top-1 acc (%) | top-5 acc (%) |
|---------------------|-------------------|---------------|---------------|
| float32 | FMA | 76.130 | 92.862 |
| (8, 1, 5, 5, 7) log | ELMA | -0.90 | -0.20 |
| (7, 1) posit | EMA | -4.63 | -2.28 |
| (8, 0) posit | EMA | -76.03 | -92.36 |
| (8, 1) posit | EMA | -0.87 | -0.19 |
| (8, 2) posit | EMA | -2.20 | -0.85 |
| (9, 1) posit | EMA | -0.30 | -0.09 |
| Jacob et al. [15]: | | | |
| float32 | FMA | 76.400 | n/a |
| int8/32 | MAC | -1.50 | n/a |
| Migacz [23]: | | | |
| float32 | FMA | 73.230 | 91.180 |
| int8/32 | MAC | -0.20 | -0.03 |

Rethinking floating point for deep learning

Mixed precision

A: All models are suitable for AMP, although the speed-up may vary from model to model. The following table provides some examples of the speed-up for different models:
Table 3. Speed-ups FP32 and Mixed Precision models.

| Model Script ¹ | Framework | Data Set | FP32 Accuracy | Mixed Precision Accuracy | FP32 Throughput | Mixed Precision Throughput | Speed-up |
|--|------------|-------------------------|-----------------|--------------------------|-------------------------------|-------------------------------|-----------|
| BERT Q&A ² | TensorFlow | SQuAD | 90.83 Top 1% | 90.99 Top 1% | 66.65 sentences/sec | 129.16 sentences/sec | 1.94 |
| SSD w/RN50 ¹ | TensorFlow | COCO 2017 | 0.268 mAP | 0.269 mAP | 569 images/sec | 752 images/sec | 1.32 |
| GNMT ³ | PyTorch | WMT16 English to German | 24.16 BLEU | 24.22 BLEU | 314,831 tokens/sec | 738,521 tokens/sec | 2.35 |
| Neural Collaborative Filter ¹ | PyTorch | MovieLens 20M | 0.959 HR | 0.960 HR | 55,004,590 samples/sec | 99,332,230 samples/sec | 1.81 |
| U-Net Industrial ¹ | TensorFlow | DAGM 2007 | 0.965-0.988 | 0.960-0.988 | 445 images/sec | 491 images/sec | 1.10 |
| ResNet-50 v1.5 ¹ | MXNet | ImageNet | 76.67 Top 1% | 76.49 Top 1% | 2,957 images/sec | 10,263 images/sec | 3.47 |
| Tacotron 2 / WaveGlow 1.0 ¹ | PyTorch | LJ Speech Dataset | 0.3629/-6.1087 | 0.3645/-6.0258 | 10,843 tok/s 257,687 smp/s | 12,742 tok/s 500,375 smp/s | 1.18/1.94 |

GPU memory

Short: the more, the better (larger models, bigger batch sizes, ...)

How can I fit +24GB models into 10GB memory?

It is a bit contradictory that I just said if you want to train big models, you need lots of memory, but we have been struggling with big models a lot since the onslaught of BERT and solutions exists to train 24 GB models in 10 GB memory. If you do not have the money or what to avoid cooling/power issues of the RTX 3090, you can get RTX 3080 and just accept that you need do some extra programming by adding memory-saving techniques. There are enough techniques to make it work, and they are becoming more and more commonplace.

Here just a list of common techniques:

- FP16/BF16 training ([apex](#))
- [Gradient checkpointing](#) (only store some of the activations and recompute them in the backward pass)
- GPU-to-CPU [Memory Swapping](#) (swap layers not needed to the CPU; swap them back in just-in-time for backprop)
- [Model Parallelism](#) (each GPU holds a part of each layer; supported by fairseq)
- Pipeline parallelism (each GPU holds a couple of layers of the network)
- [ZeRO](#) parallelism (each GPU holds partial layers)
- [3D parallelism](#) (Model + pipeline + ZeRO)
- CPU Optimizer state (store and update Adam/Momentum on the CPU while the next GPU forward pass is happening)

The best GPU for Deep LEarning

<https://timdettmers.com/2020/09/07/which-gpu-for-deep-learning/>

TL;DR advice

Best GPU overall: RTX 3080 and RTX 3090.

GPUs to avoid (as an individual): Any Tesla card; any Quadro card; any Founders Edition card; Titan RTX, Titan V, Titan XP.

Cost-efficient but expensive: RTX 3080.

Cost-efficient and cheaper: RTX 3070, RTX 2060 Super

I have little money: Buy used cards. Hierarchy: RTX 2070 (\$400), RTX 2060 (\$300), GTX 1070 (\$220), GTX 1070 Ti (\$230), GTX 1650 Super (\$190), GTX 980 Ti (6GB \$150).

I have almost no money: There are a lot of startups that promo their clouds: Use free cloud credits and switch companies accounts until you can afford a GPU.

I do Kaggle: RTX 3070.

I am a competitive computer vision, pretraining, or machine translation researcher: 4x RTX 3090. Wait until working builds with good cooling, and enough power are confirmed (I will update this blog post).

I am an NLP researcher: If you do not work on machine translation, language modeling, or pretraining of any kind, an RTX 3080 will be sufficient and cost-effective.

What do I have in the lab?

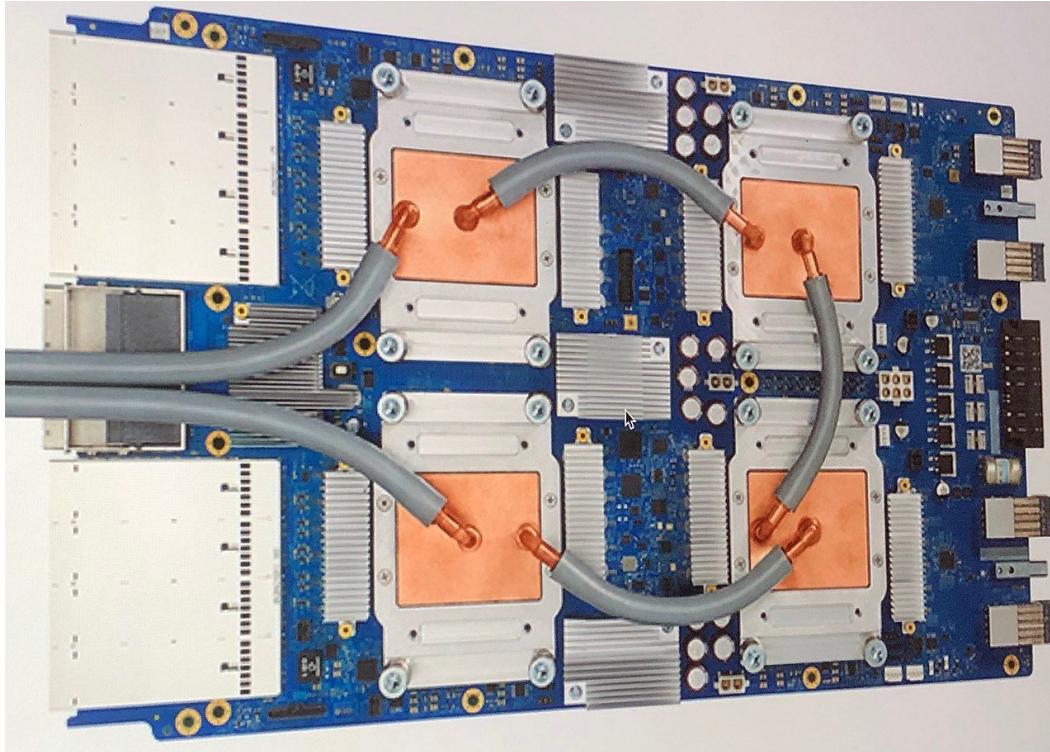
- RTX 2080 Ti (1x)
- Titan RTX
- RTX 3090 (3x)



Application Specific ICs (integrated circuits) ASIC



TPUs



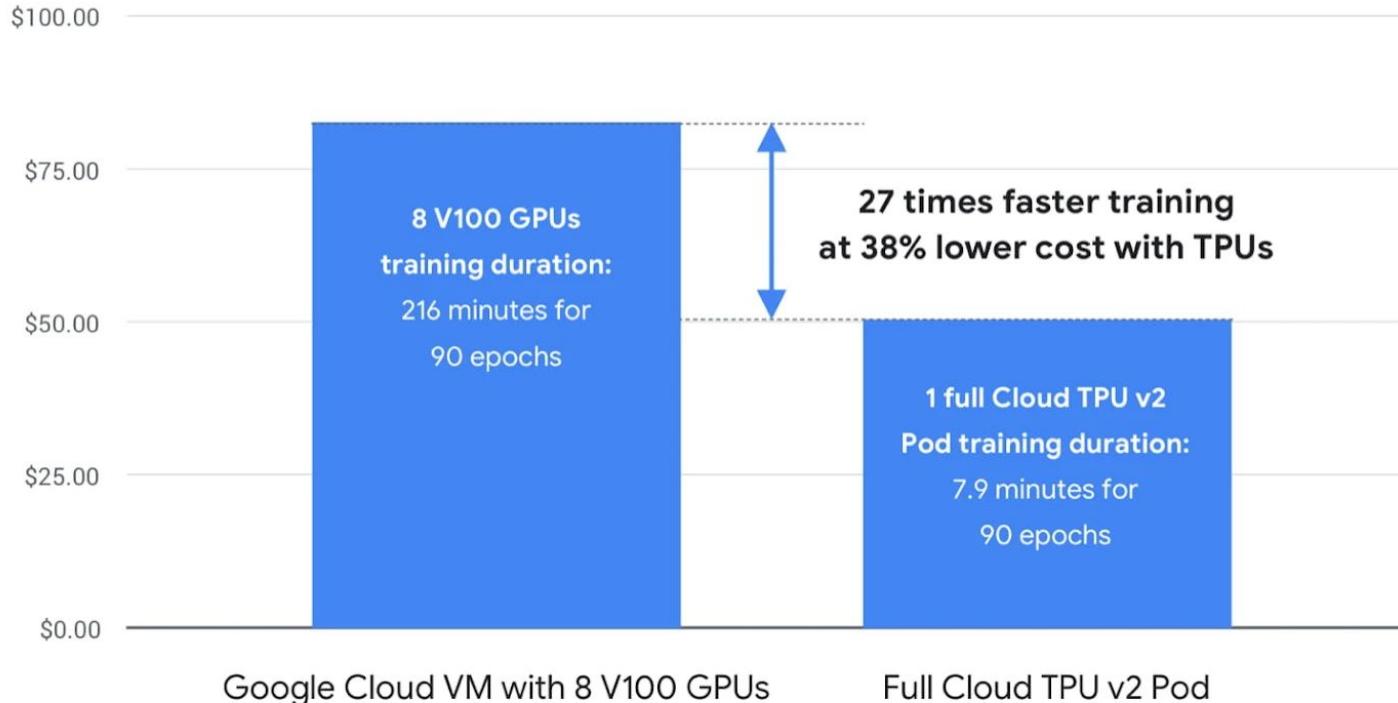
Products^[13] [edit]

| | TPUv1 | TPUv2 | TPUv3 | TPUv4 ^[14] | Edge v1 |
|-----------------------------|----------|----------|----------|-----------------------|---------|
| Date Introduced | 2016 | 2017 | 2018 | 2021 | 2018 |
| Process Node | 28 nm | 16 nm | 16 nm | 7 nm | |
| Die Size (mm ²) | 331 | < 625 | < 700 | < 400 | |
| On chip memory (MiB) | 28 | 32 | 32 | 144 | |
| Clock Speed (MHz) | 700 | 700 | 940 | 1050 | |
| Memory (GB) | 8GB DDR3 | 16GB HBM | 32GB HBM | 8GB | |
| TDP(W) | 75 | 280 | 450 | 175 | 2 |
| TOPS | 23 | 45 | 90 | ? | 4 |



TPUs are faster (in some regard)

ResNet-50 Training Cost Comparison



TPUs are in “nodes”

Cloud TPUs are available in:

US Europe

US EURO

Version

Cloud T

Cloud T

v3-256

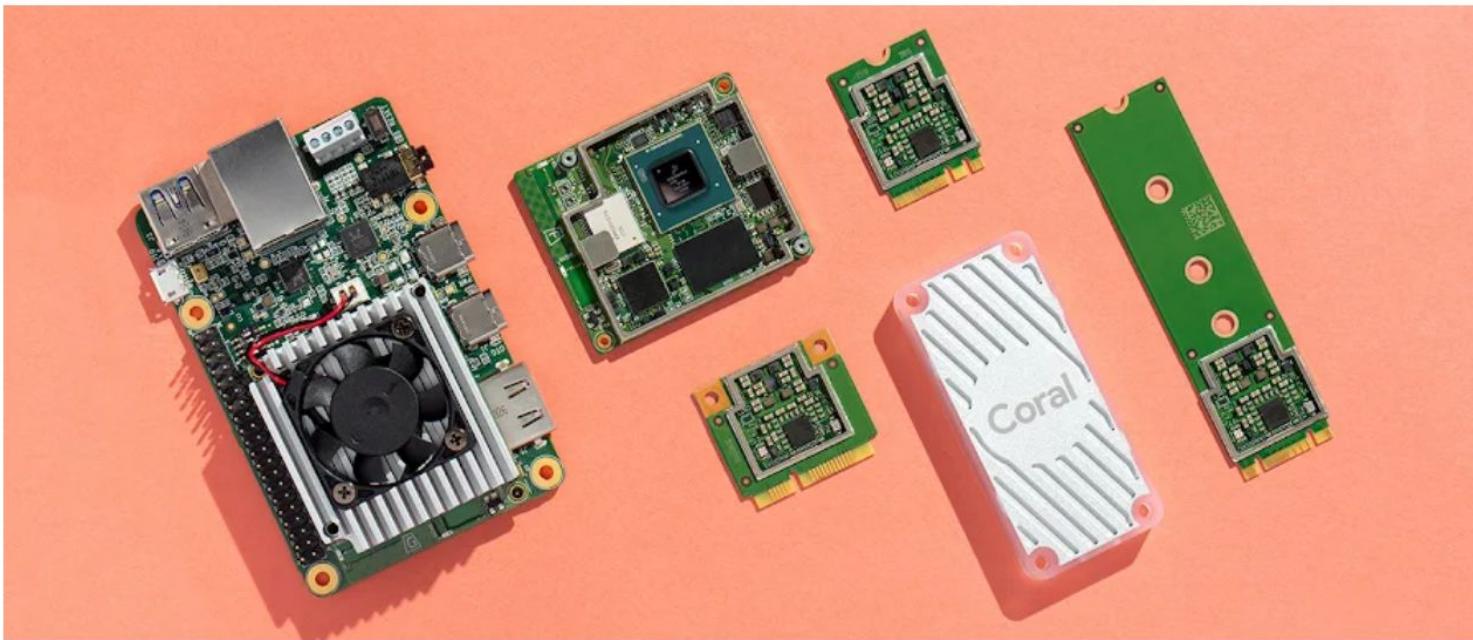
v3-512

v3-1024

v3-2048

| Cloud TPU v2 Pod | Evaluation Price / hr | 1-yr Commitment Price (37% discount) |
|--------------------|-----------------------|---|
| 32-core Pod slice | \$24 USD | \$132,451 USD |
| 128-core Pod slice | \$96 USD | \$529,805 USD |
| 256-core Pod slice | \$192 USD | \$1,059,610 USD |
| 512-core Pod slice | \$384 USD | \$2,119,219 USD |
| Cloud TPU v3 Pod | Evaluation Price / hr | 1-yr Commitment Price (37% discount) |
| 32-core Pod slice | \$32 USD | \$176,601 USD |

The Edge TPU

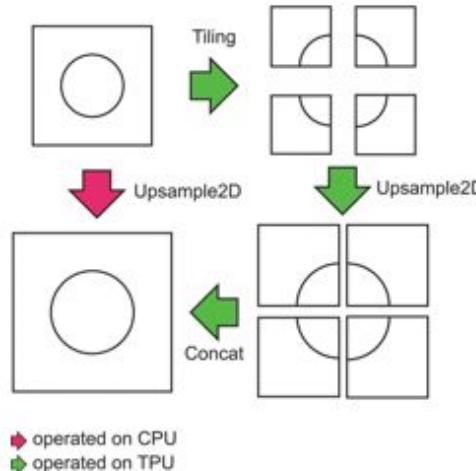
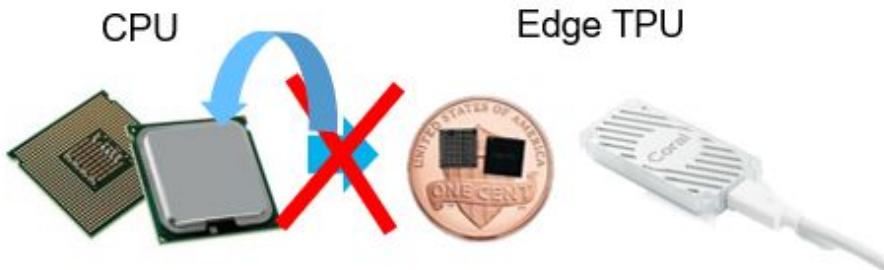
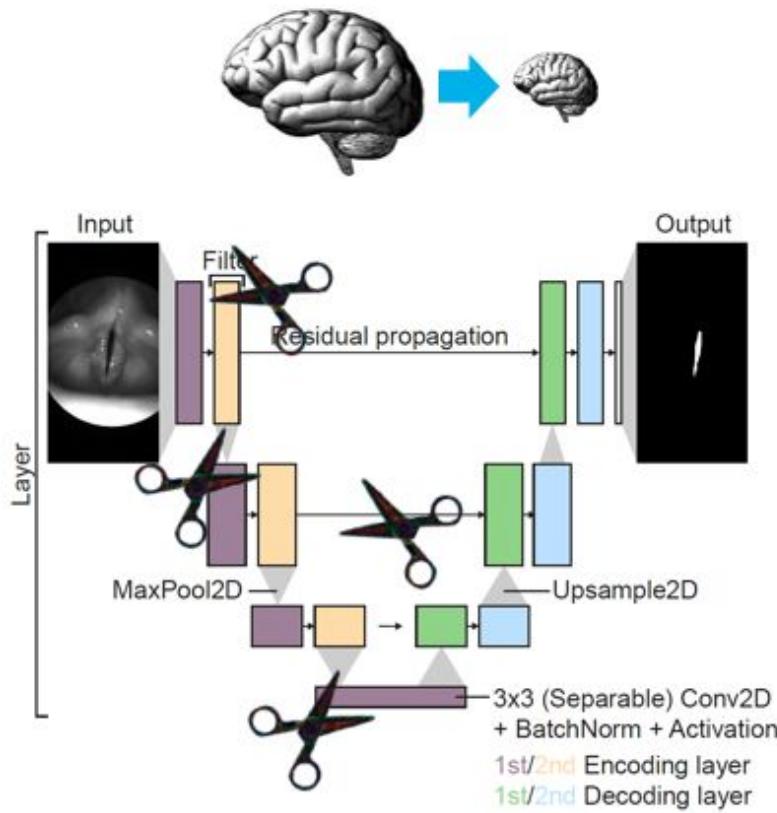


Inference boost

Table 1. Time per inference, in milliseconds (ms)

| Model architecture | Desktop CPU ¹ | Desktop CPU ¹ + USB Accelerator (USB 3.0) <i>with Edge TPU</i> | Embedded CPU ² | Dev Board ³ <i>with Edge TPU</i> |
|----------------------------------|--------------------------|---|---------------------------|--|
| Unet Mv2 (128x128) | 27.7 | 3.3 | 190.7 | 5.7 |
| DeepLab V3 (513x513) | 394 | 52 | 1139 | 241 |
| DenseNet (224x224) | 380 | 20 | 1032 | 25 |
| Inception v1 (224x224) | 90 | 3.4 | 392 | 4.1 |
| Inception v4 (299x299) | 700 | 85 | 3157 | 102 |
| Inception-ResNet V2 (299x299) | 753 | 57 | 2852 | 69 |
| MobileNet v1 (224x224) | 53 | 2.4 | 164 | 2.4 |
| MobileNet v2 (224x224) | 51 | 2.6 | 122 | 2.6 |
| MobileNet v1 SSD (224x224) | 109 | 6.5 | 353 | 11 |
| MobileNet v2 SSD (224x224) | 106 | 7.2 | 282 | 14 |
| ResNet-50 V1 (299x299) | 484 | 49 | 1765 | 56 |

Semantic segmentation on Edge TPUs



More AI hardware accelerators



Nvidia Jetson platform
(GPU based)
→ Jetson SDK



VPU (Vision Processing unit) → OpenVINO

FPGAs

Field-programmable gate arrays



Circuit can be programmed

Amazing for dedicated tasks

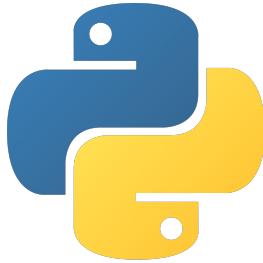
Deep learning:

Xilinx Alveo U50/U250...

Summary

- CPUs
 - Quick prototyping that requires maximum flexibility
 - Simple models that do not take long to train
 - Small models with small effective batch sizes
 - Models that are dominated by [custom TensorFlow operations written in C++](#)
 - Models that are limited by available I/O or the networking bandwidth of the host system
- GPUs
 - Models for which source does not exist or is too onerous to change
 - Models with a significant number of custom TensorFlow operations that must run at least partially on CPUs
 - Models with TensorFlow ops that are not available on Cloud TPU (see the list of [available TensorFlow ops](#))
 - Medium-to-large models with larger effective batch sizes
- TPUs
 - Models dominated by matrix computations
 - Models with no custom TensorFlow operations inside the main training loop
 - Models that train for weeks or months
 - Larger and very large models with very large effective batch sizes

How do I use CPU/GPU/TPUs?!?!



Programming on CPUs and GPUs

Array programming with NumPy

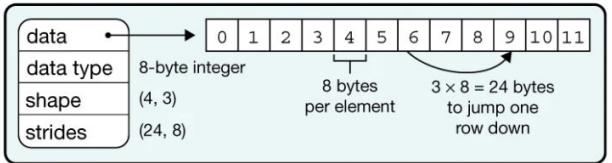
Charles R. Harris, K. Jarrod Millman, Travis E. Oliphant

Nature 585, 357–362 (2020) | Cite this article

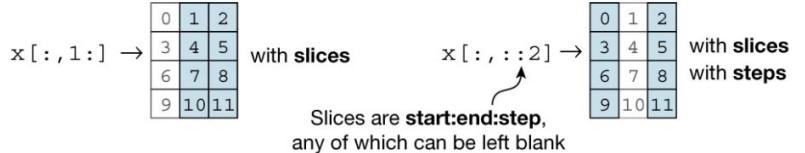
238k Accesses | 1019 Citations | 1992 Altmetric | Metrics

a Data structure

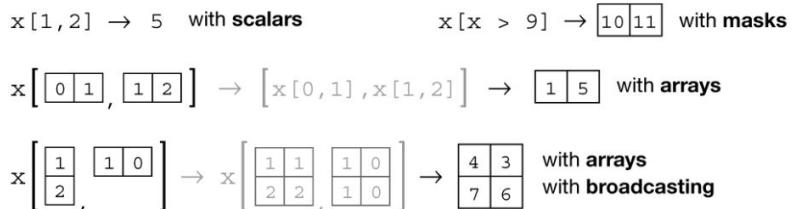
| | | |
|---|----|----|
| 0 | 1 | 2 |
| 3 | 4 | 5 |
| 6 | 7 | 8 |
| 9 | 10 | 11 |



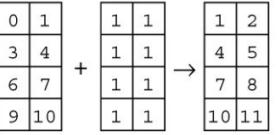
b Indexing (view)



c Indexing (copy)



d Vectorization



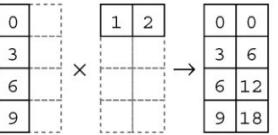
g Example

In [1]: import numpy as np

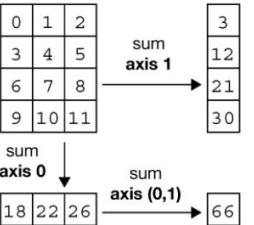
In [2]: $x = \text{np.arange}(12)$

In [3]: $x = x.\text{reshape}(4, 3)$

e Broadcasting



f Reduction



In [4]: x

Out [4]:

```
array([[ 0,  1,  2],
       [ 3,  4,  5],
       [ 6,  7,  8],
       [ 9, 10, 11]])
```

In [5]: `np.mean(x, axis=0)`

Out [5]: `array([4.5, 5.5, 6.5])`

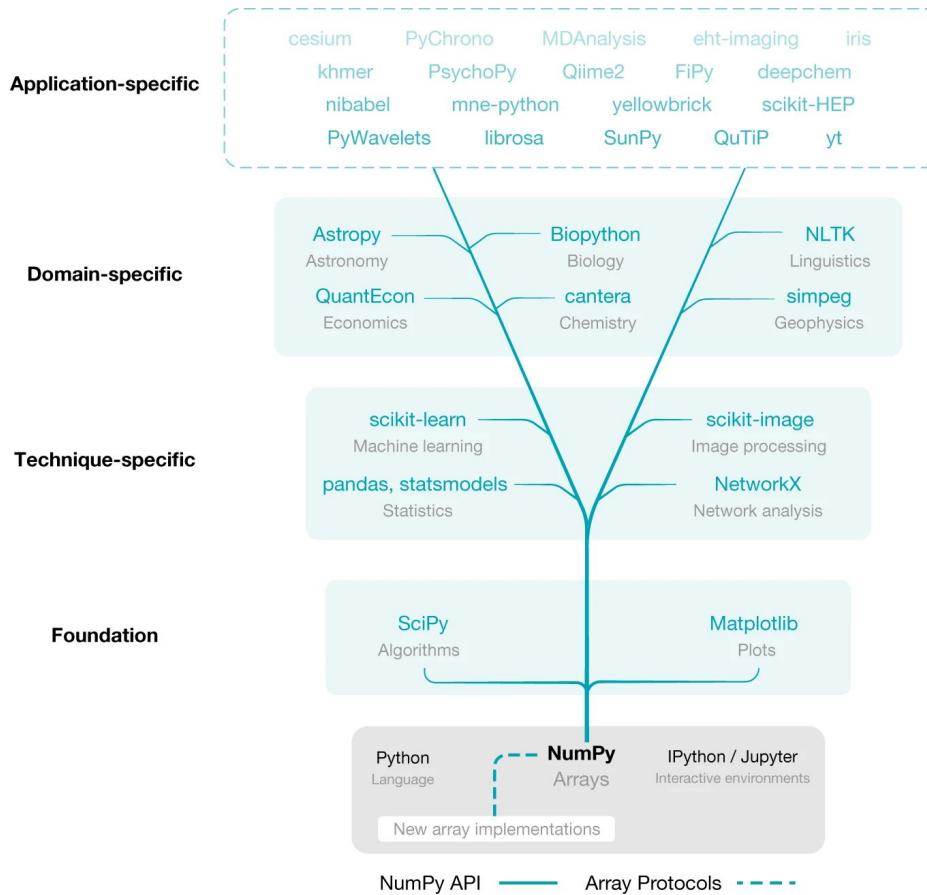
In [6]: $x = x - \text{np.mean}(x, \text{axis}=0)$

In [7]: x

Out [7]:

```
array([[-4.5, -4.5, -4.5],
       [-1.5, -1.5, -1.5],
       [ 1.5,  1.5,  1.5],
       [ 4.5,  4.5,  4.5]])
```

What libraries are using numpy?



... and GPUs?

GPU-Accelerated Computing with Python

NVIDIA's CUDA Python provides a driver and runtime API for existing toolkits and libraries to simplify GPU-based accelerated processing. Python is one of the most popular programming languages for science, engineering, data analytics, and deep learning applications. However, as an interpreted language, it's been considered too slow for high-performance computing.



Numba—a Python compiler from Anaconda that can compile Python code for execution on CUDA®-capable GPUs—provides Python developers with an easy entry into GPU-accelerated computing and for using increasingly sophisticated CUDA code with a minimum of new syntax and jargon. With CUDA Python and Numba, you get the best of both worlds: rapid iterative development with Python combined with the speed of a compiled language targeting both CPUs and NVIDIA GPUs.

Welcome to cuML's documentation!

cuML is a suite of fast, GPU-accelerated machine learning algorithms designed for data science and analytical tasks. Our API mirrors Sklearn's, and we provide practitioners with the easy fit-predict-transform paradigm without ever having to program on a GPU.

As data gets larger, algorithms running on a CPU becomes slow and cumbersome. RAPIDS provides users a streamlined approach where data is initially loaded in the GPU, and compute tasks can be performed on it directly.

cuML is fully open source, and the RAPIDS team welcomes new and seasoned contributors, users and hobbyists! Thank you for your wonderful support!

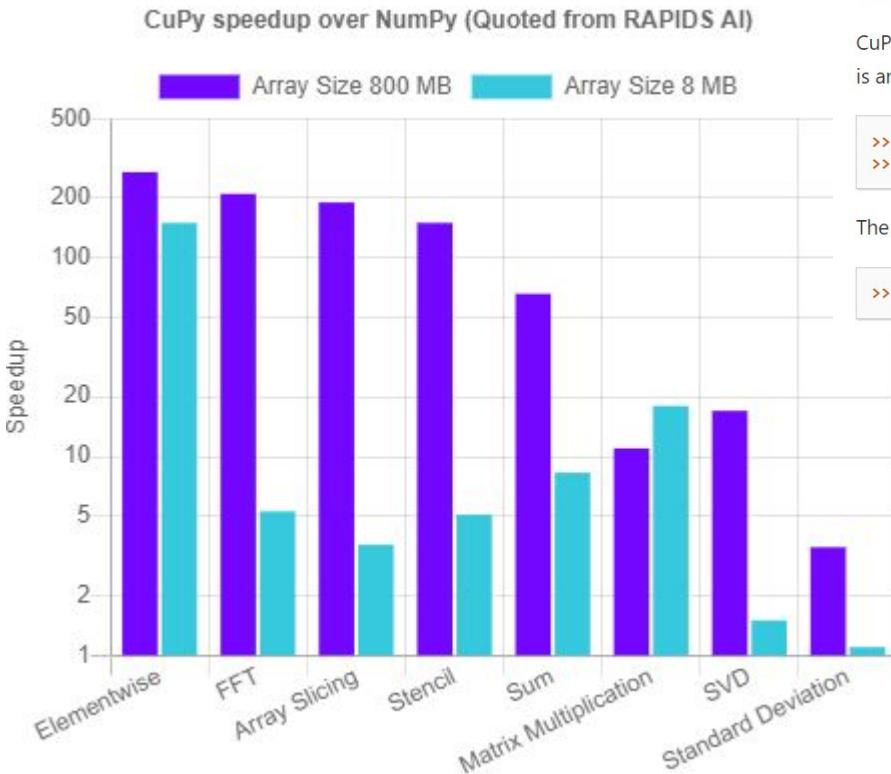
An installation requirement for cuML is that your system must be Linux-like. Support for Windows is possible in the near future.



Welcome to cuDF's documentation!

cuDF is a Python GPU DataFrame library (built on the Apache Arrow columnar memory format) for loading, joining, aggregating, filtering, and otherwise manipulating data. cuDF also provides a pandas-like API that will be familiar to data engineers & data scientists, so they can use it to easily accelerate their workflows without going into the details of CUDA programming.

CuPy



Basics of `cupy.ndarray`

CuPy is a GPU array backend that implements a subset of NumPy interface. In the following code, `cp` is an abbreviation of `cupy`, following the standard convention of abbreviating `numpy` as `np`:

```
>>> import numpy as np  
>>> import cupy as cp
```

The `cupy.ndarray` class is at the core of CuPy and is a replacement class for NumPy's `numpy.ndarray`.

```
>>> x_gpu = cp.array([1, 2, 3])
```

- Routines (SciPy)
 - Discrete Fourier transforms (`cupyx.scipy.fft`)
 - Legacy discrete fourier transforms (`cupyx.scipy.fftpack`)
 - Linear algebra (`cupyx.scipy.linalg`)
 - Multidimensional image processing (`cupyx.scipy.ndimage`)
 - Sparse matrices (`cupyx.scipy.sparse`)
 - Special functions (`cupyx.scipy.special`)
 - Signal processing (`cupyx.scipy.signal`)
 - Statistical functions (`cupyx.scipy.stats`)

The last slide

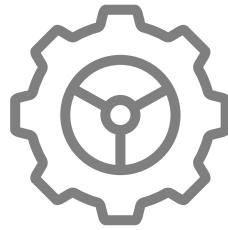
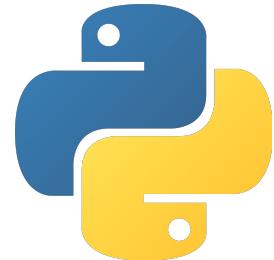
- Comparison CPU and GPU in 1:30



Exercise

Description of the exercise

How to use the GPU on Google Colab and Why? How do the CuPy's arrays work? We are going to discuss these questions.



Description of the exercise

For the very first exercise of the lecture, we are going to have an introduction to arrays. Slicing, indexing, calculations, and types of arrays are some of the topics of this activity.



Description of the exercise

Your tasks?



- Use indexing and slicing to edit an image.
- Compare and contrast CuPy and Numpy for some operations.



CuPy

Numpy