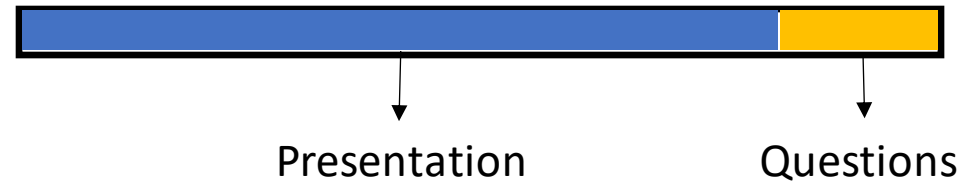


Machine Learning Tools

Bazyli Polednia 2021

Plan for today

Total: 1h



Deep Learning Boom in recent years

1. Hardware and software

- Introduction of TPUs and more powerful GPUs
- Parallelization of computations
- Programming interfaces e.g. Nvidia CUDA
- Deep Learning libraries and toolkits e.g. Tensorflow, PyTorch, Keras



2. Datasets and benchmarks

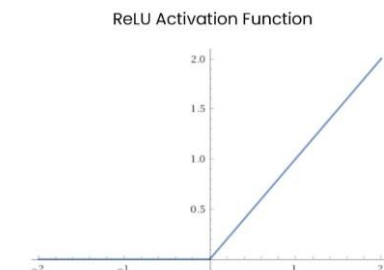
- „Big Data”
- Exponential growth of hardware storage
- Public datasets e.g. ImageNet, Kaggle competitions

IMAGENET

kaggle™

3. Algorithmic advances

- New activation functions
- New optimizers
- Batch optimization
- ...



Deep Learning Boom in recent years

1. Hardware and software

- Introduction of TPUs and more powerful GPUs
- Parallelization of computations
- Programming interfaces e.g. Nvidia CUDA
- Deep Learning libraries and toolkits e.g. Tensorflow, PyTorch, Keras



2. Datasets and benchmarks

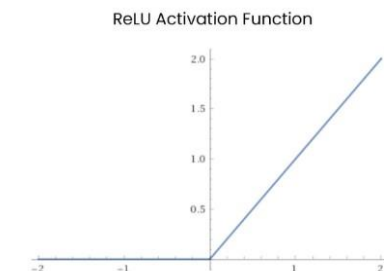
- „Big Data”
- Exponential growth of hardware storage
- Public datasets e.g. ImageNet, Kaggle competitions

IMAGENET

kaggle™

3. Algorithmic advances

- New activation functions
- New optimizers
- Batch optimization
- ...



CPU vs. GPU

CPU vs. GPU

CPU

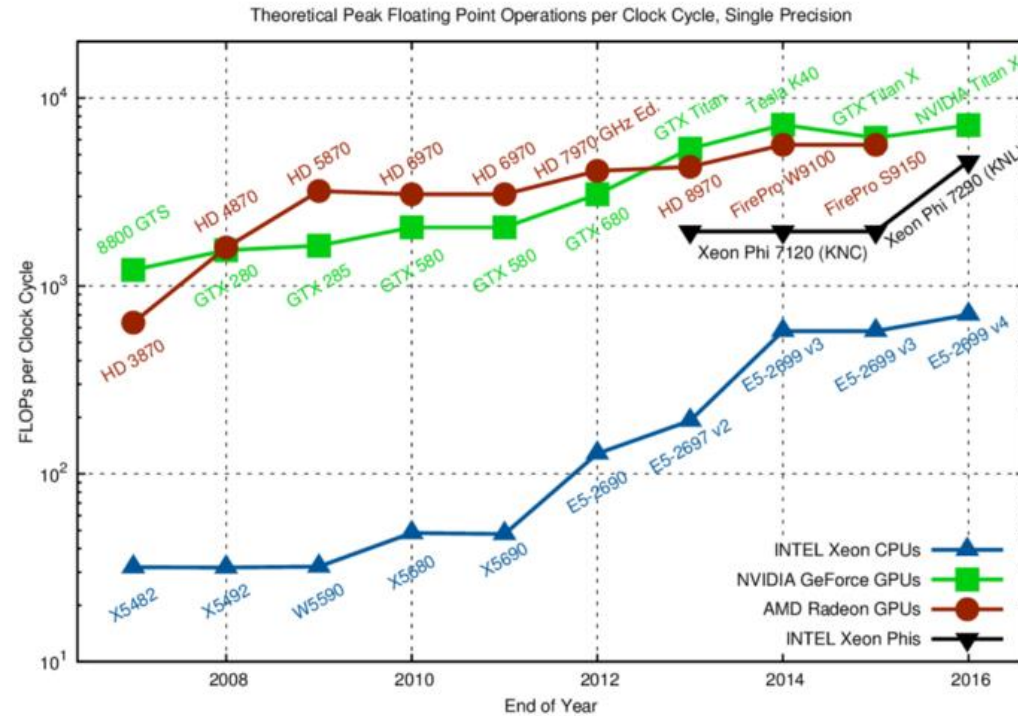
- Few complex cores
- Single-thread performance optimization
- Transistors dedicated to complex ILP
- Small die surface for integer and floating point units

GPU

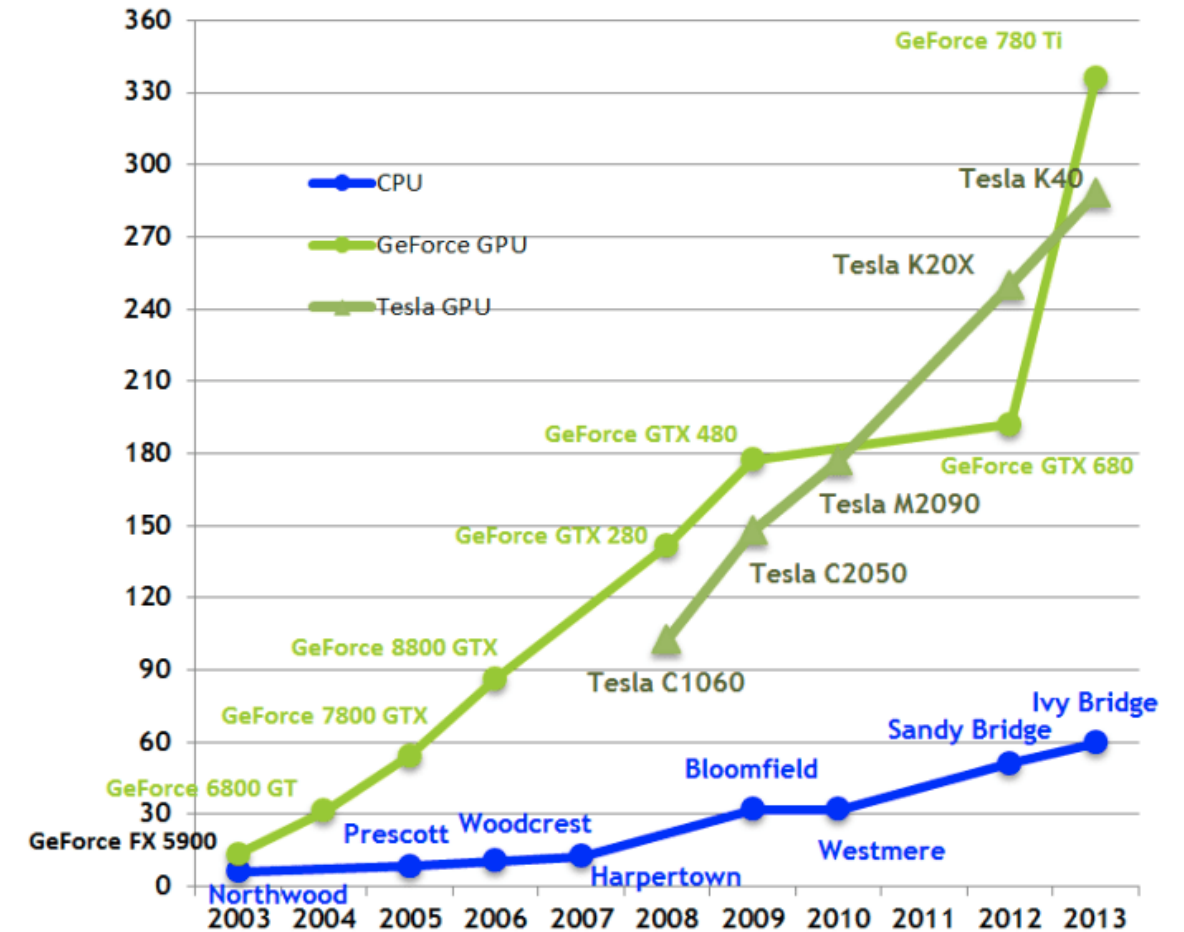
- Many simpler cores
- Many concurrent hardware threads
- Maximization of float point throughput
- Large die surface for integer and floating point units

GPUs are optimized for ***processing multiple computations*** simultaneously.

CPU vs. GPU



Theoretical GB/s

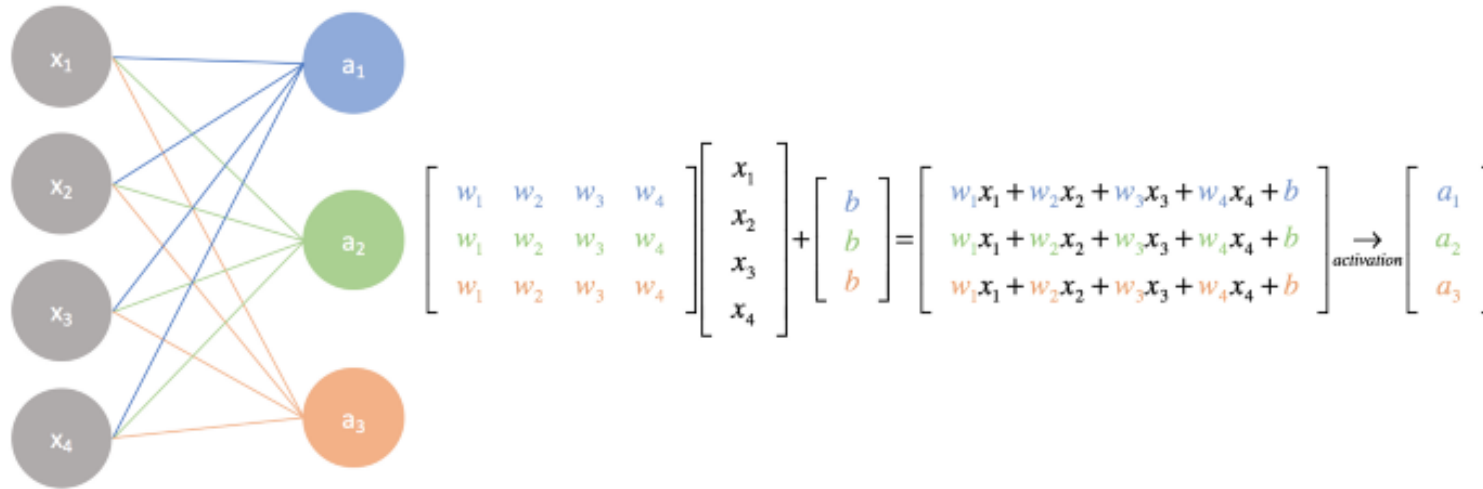


Why are GPUs used in DL?

Input layer

Output layer

A simple neural network



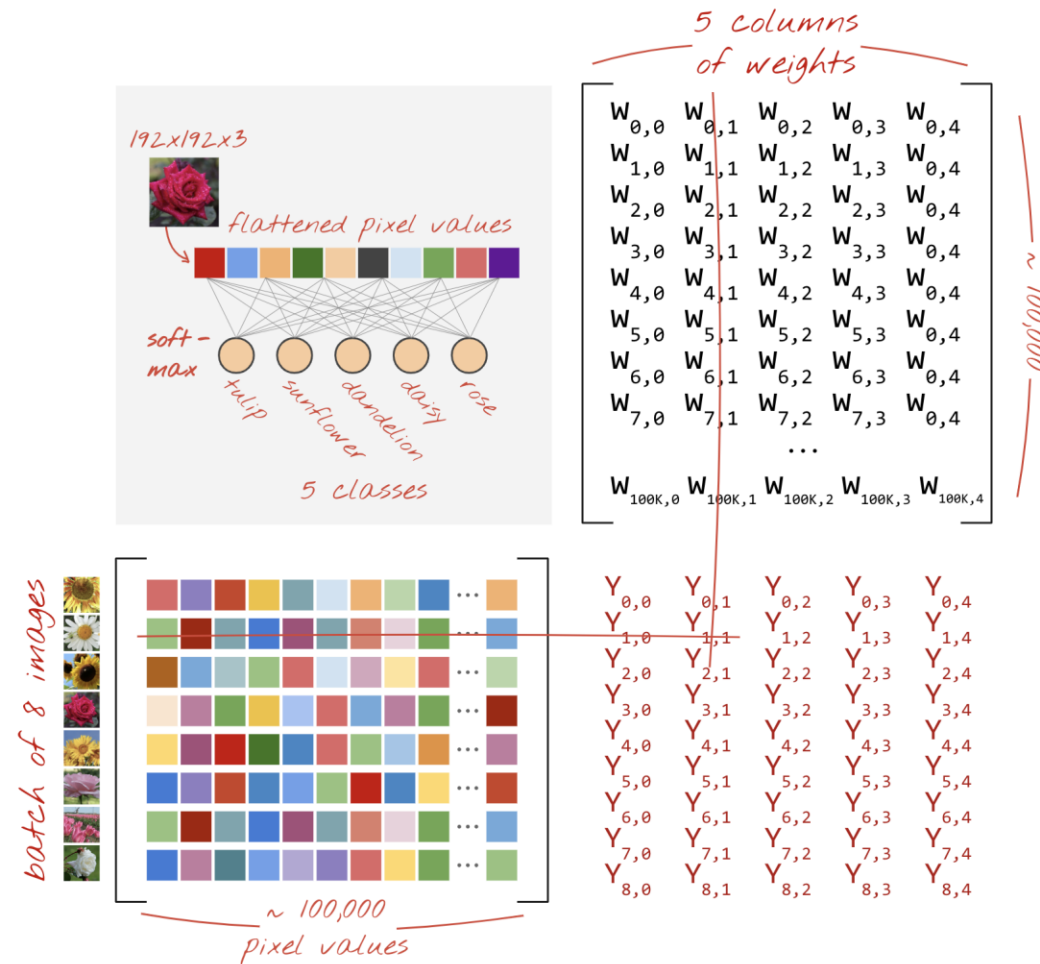
Can we do even better than GPU?

TPU – Tensor Processing Unit

Modern GPUs are organized around programmable "cores", a very flexible architecture that allows them to handle a variety of tasks such as 3D rendering, deep learning, physical simulations, etc.

TPUs on the other hand pair a classic vector processor with a dedicated matrix multiply unit and excel at any task where large matrix multiplications dominate, such as neural networks.

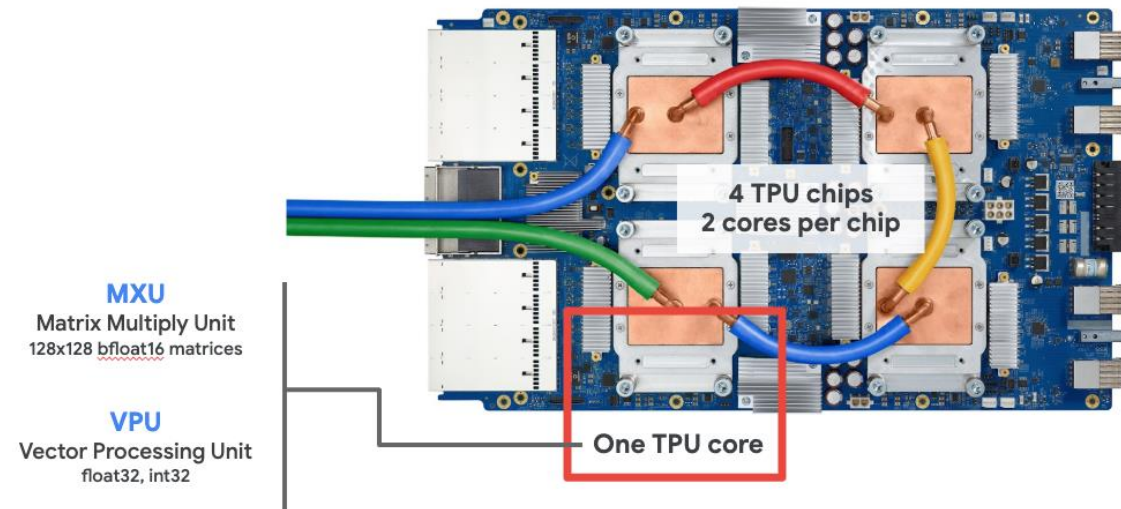
TPU – Tensor Processing Unit



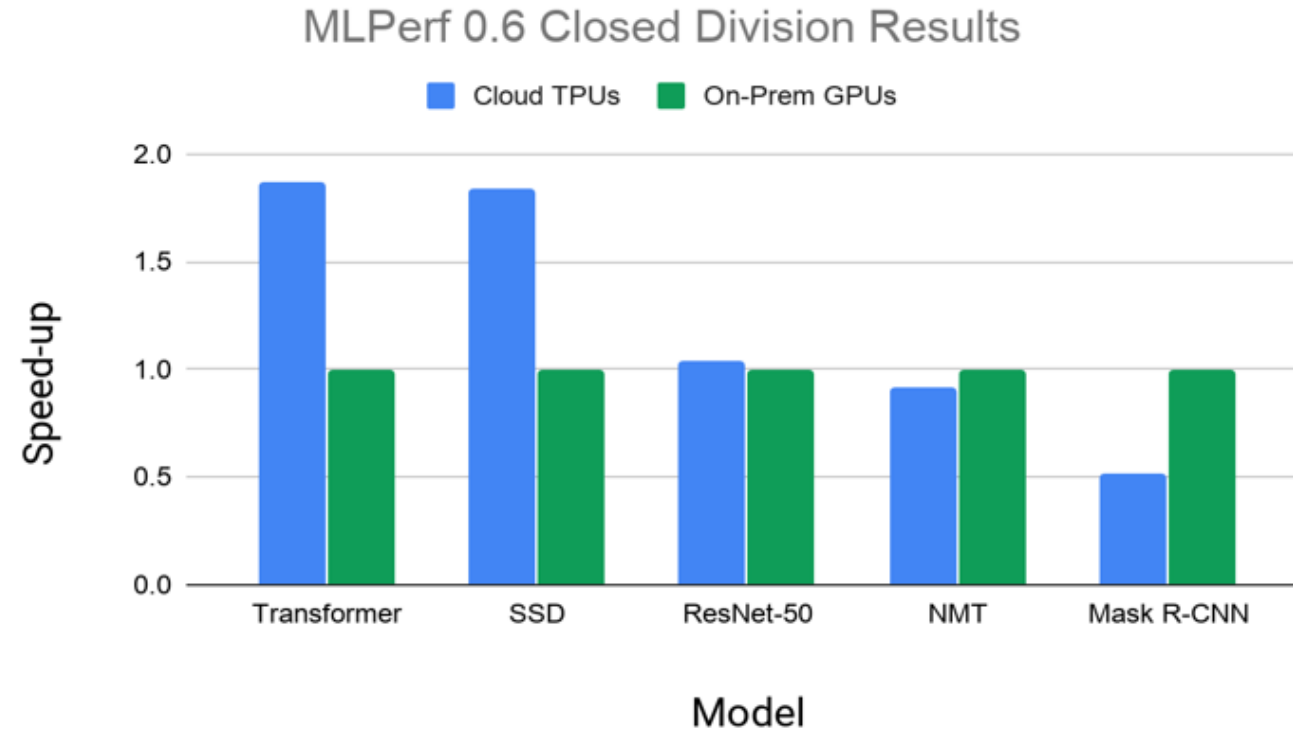
TPU – Tensor Processing Unit

MXU (Matrix Multiply Unit) – responsible for matrix multiplication

VPU (Vector Processing Unit) – handles other DL tasks such as activation functions



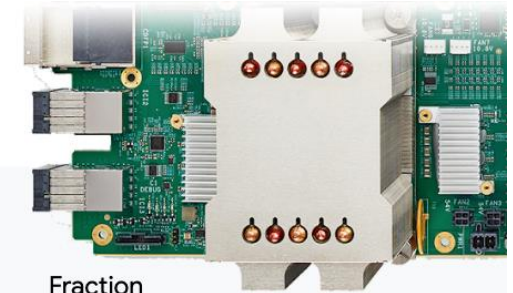
TPU – Tensor Processing Unit



Google Cloud TPU v3 Pod speed-ups over the largest-scale on-premise NVIDIA DGX-2h clusters entered in the MLPerf 0.6 Closed Division. The Cloud TPU Pod submissions use 1024, 1024, 1024, 512, and 128 chips respectively; the NVIDIA DGX-2h clusters use 480, 240, 1536, 256, and 192 chips respectively. ^{1,2}

bfloat16

bfloat16



bfloat16

range: $\sim 1e^{-38}$ to $\sim 3e^{38}$



float32

range: $\sim 1e^{-38}$ to $\sim 3e^{38}$



float16

range: $\sim 5.9e^{-8}$ to $6.5e^4$



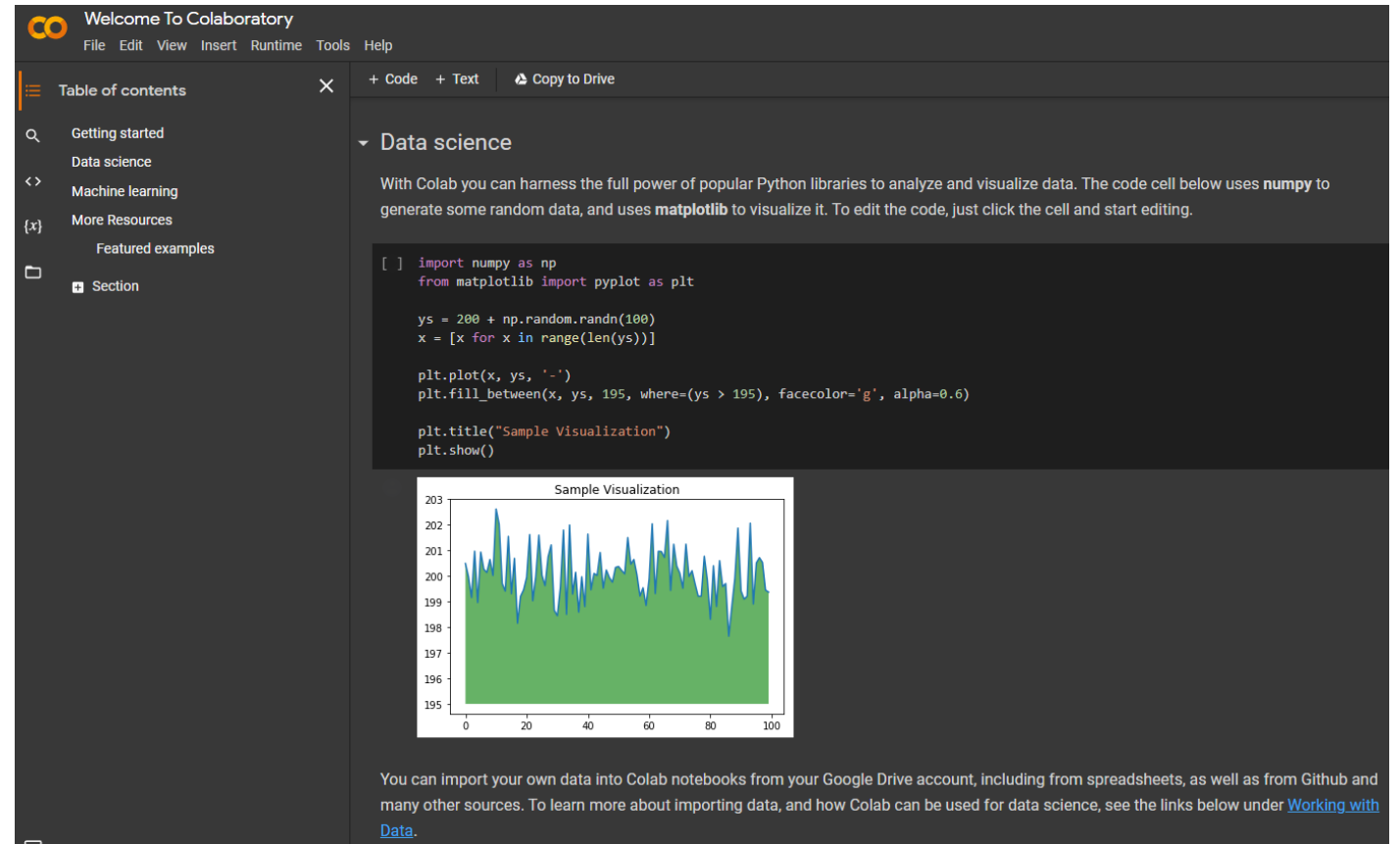
$$(-1)^{b_{31}} \times 2^{(b_{30}b_{29} \dots b_{23})_2 - 127} \times (1.b_{22}b_{21} \dots b_0)_2,$$

Data in Deep Learning is already noisy due to sampling and measurement imprecision

Google Colab

Colab allows you to write and execute Python in your browser, with:

- Zero configuration required
- Free access to GPUs and TPUs
- Easy sharing



The screenshot displays the Google Colaboratory web interface. On the left is a sidebar with a 'Table of contents' and a search bar. The main area is titled 'Welcome To Colaboratory' and contains a 'Data science' section. Below this, a code cell is shown with Python code that generates random data and visualizes it using a line plot with a green shaded area. The output of the code cell is a line plot titled 'Sample Visualization' showing a noisy signal fluctuating between approximately 195 and 203 over 100 iterations.

Table of contents

- Getting started
- Data science
- Machine learning
- More Resources
- Featured examples
- Section

Welcome To Colaboratory

File Edit View Insert Runtime Tools Help

+ Code + Text Copy to Drive

Data science

With Colab you can harness the full power of popular Python libraries to analyze and visualize data. The code cell below uses **numpy** to generate some random data, and uses **matplotlib** to visualize it. To edit the code, just click the cell and start editing.

```
[ ] import numpy as np
    from matplotlib import pyplot as plt

    ys = 200 + np.random.randn(100)
    x = [x for x in range(len(ys))]

    plt.plot(x, ys, '-')
    plt.fill_between(x, ys, 195, where=(ys > 195), facecolor='g', alpha=0.6)

    plt.title("Sample Visualization")
    plt.show()
```

Sample Visualization

You can import your own data into Colab notebooks from your Google Drive account, including from spreadsheets, as well as from Github and many other sources. To learn more about importing data, and how Colab can be used for data science, see the links below under [Working with Data](#).

Hardware Accelerators

- Quick prototype and deployment
- Plug and play
- Support for OpenVINO toolkit
- Pre-trained models
- Model Optimizer
- Inference Engine

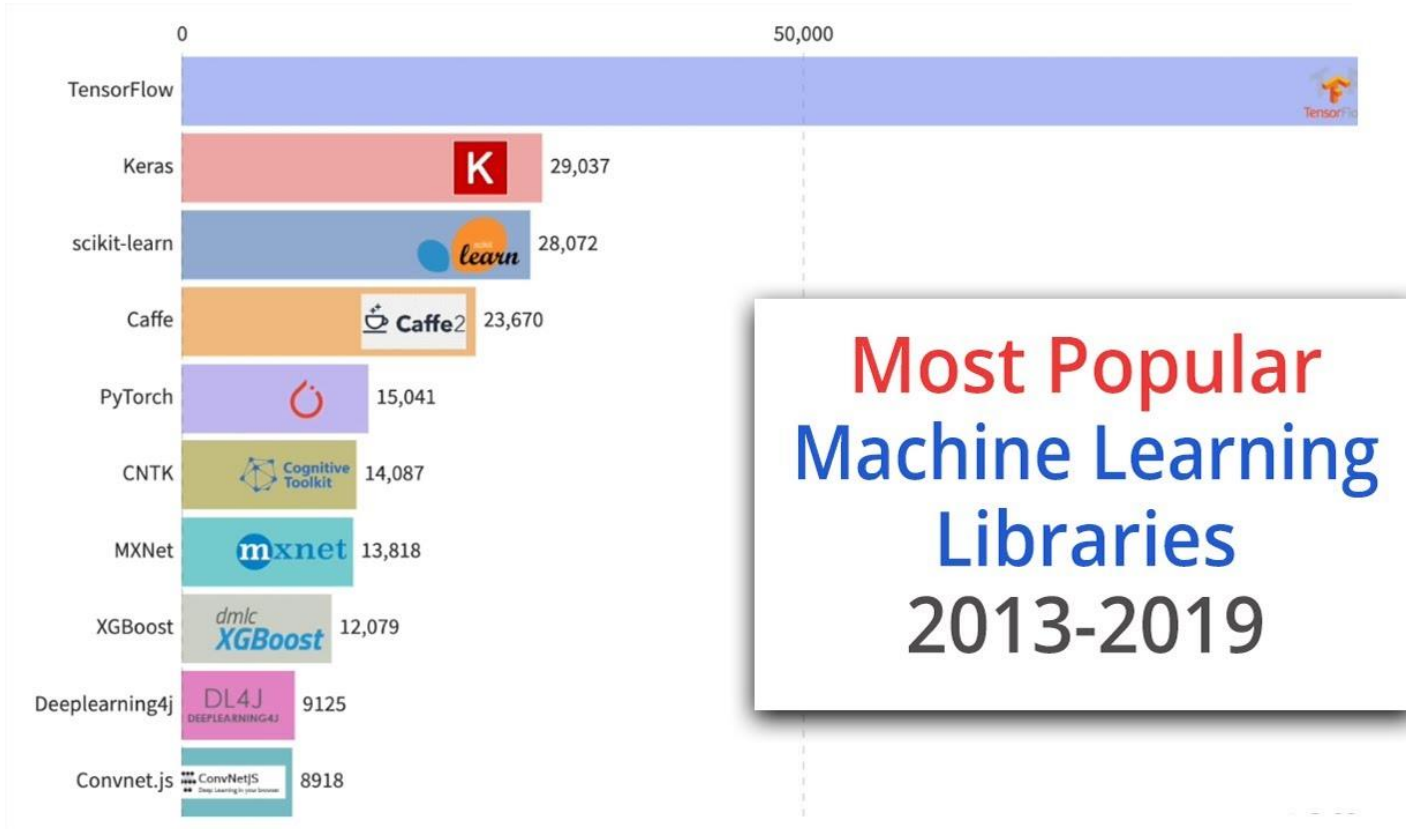


Hardware Accelerators



Machine Learning libraries and frameworks





TensorFlow

- Developed by Google, released in 2015
- Offers level of abstraction while building DL models
- Support for multiple programming languages
- Offers cool visualization metrics



TensorFlow

TensorBoard

TensorBoard.dev

SCALARS

My latest experiment

Simple comparison of several hyperparameters

- ☐ Show data download links
☒ Ignore outliers in chart scaling

Tooltip sorting method: **default**

Smoothing

0.6

Horizontal Axis

STEP

RELATIVE

WALL

Runs

Write a regex to filter runs

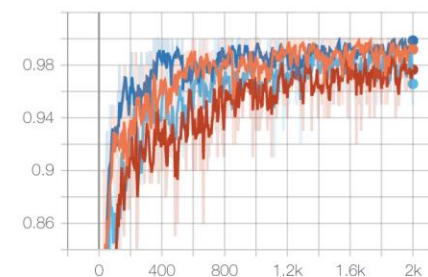
- ☒ ☐ lr_1E-03,conv=1,fc=2
- ☒ ☐ lr_1E-03,conv=2,fc=2
- ☒ ☐ lr_1E-04,conv=1,fc=2
- ☒ ☐ lr_1E-04,conv=2,fc=2

TOGGLE ALL RUNS

experiment AdYd1TgeTlaLWXx6l8JUba

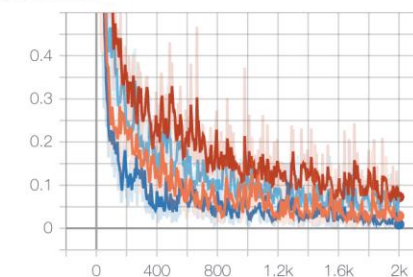
accuracy

accuracy
tag: accuracy/accuracy



xent

xent_1
tag: xent/xent_1



PyTorch

- Developed by Facebook, released in 2016
- Great support for distributed training and cloud computations
- Good documentation and online support community



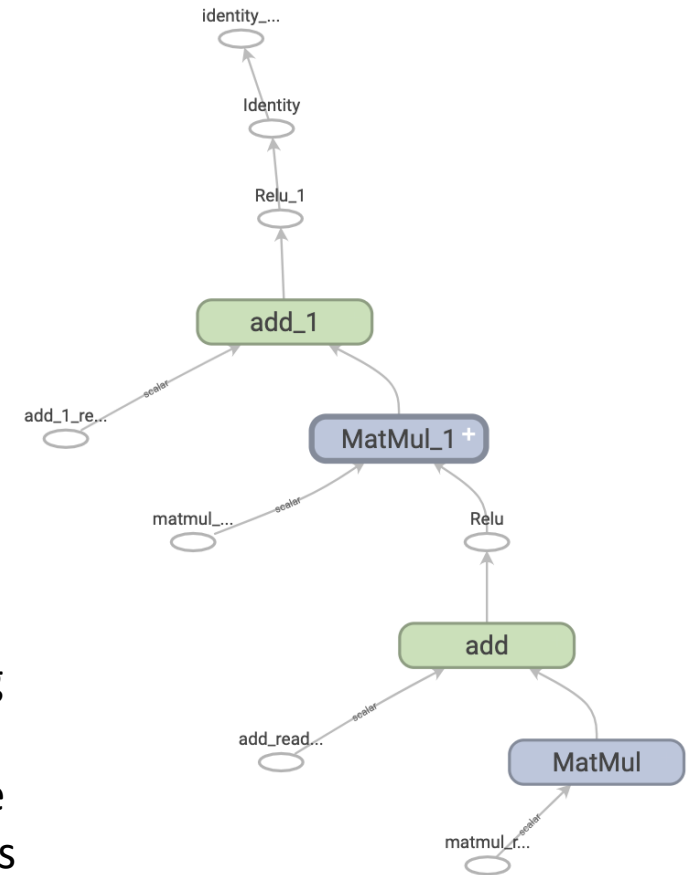
Graph execution vs. eager execution

TensorFlow offers both eager execution and **graph execution** modes.

Graphs are data structures that contain a set of `tf.Operation` objects, which represent units of computation; and `tf.Tensor` objects, which represent the units of data that flow between operations.

Graphs are also easily optimized:

- Statically infer the value of tensors by folding constant nodes
- Separate sub-parts of a computation that are independent and split them between threads or devices.
- Simplify arithmetic operations by eliminating common subexpressions.



Automatic differentiation

Task: compute gradient for function $y=x^2$ at $x=3$

TensorFlow 

```
x = tf.Variable(3.0)
with tf.GradientTape() as tape:
    y = x * x
dy_dx = tape.gradient(y, x)
```

PyTorch 

```
x = torch.tensor([3.], requires_grad=True)
Q = x**2
external_grad = torch.tensor([1.])
Q.backward(gradient=external_grad)
x.grad
```

Simple neural network

Task: create a two-layer neural network

TensorFlow 

```
model = tf.keras.Sequential([
    tf.keras.layers.Dense(128, activation='relu'),
    tf.keras.layers.Dense(10)
])
model.compile(
    optimizer="adam",
    loss=tf.keras.losses.SparseCategoricalCrossentropy(),
    metrics=["accuracy"],
)
```

PyTorch 

```
class Net(nn.Module):

    def __init__(self):
        super(Net, self).__init__()
        self.fc1 = nn.Linear(128, 64)
        self.fc2 = nn.Linear(64, 10)

    def forward(self, x):
        x = F.relu(self.fc1(x))
        x = self.fc2(x)
        return x

net = Net()
```

Questions?

Thank you for attention