



Brad Miro - Google
@bradmiro
PyGotham 2019

Distributed Machine Learning with Python

Google Cloud

Agenda

Intro to Distributed Machine Learning

Hardware Considerations

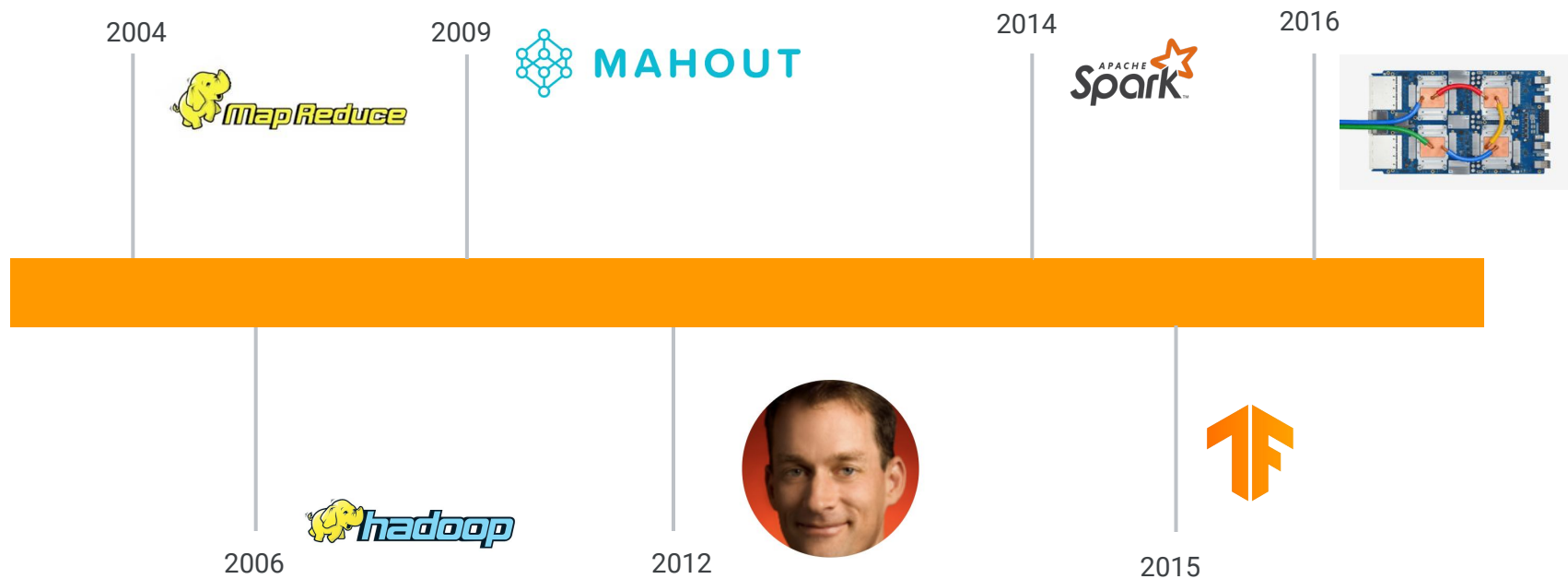
Paradigms

Software Considerations

Closing Thoughts

Intro to Distributed Machine Learning

Distributed Machine Learning



Machine Learning on One Machine

Machine Learning is **hard**

Optimizing linear algebra is computationally expensive

Large models may not fit into memory

Processing training data is time-consuming

What if machine crashes?



Machine Learning on Multiple Machines

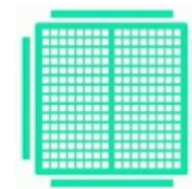
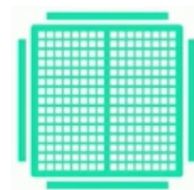
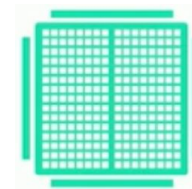
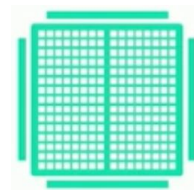
Linear algebra is parallelizable

Large models are split across multiple machines

Processing training data is parallelizable

Introduce fault tolerance

Provide near-linear additional performance per machine



Why?

Larger models are often better for complex problems - audio, images, text

Models can often benefit from large amounts of data

You Shouldn't **ALWAYS** Distribute

Parallelization overhead

I/O limitation of older hardware

Smaller models can train faster on one machine

Questions to Ask

Do I REALLY need to distribute my training?

Will I expect to need to distribute my training?

Do I have a lot of data?

Will my data continue to grow?

Will my model continue to grow?



Hardware Considerations

Central Processing Unit (CPU)

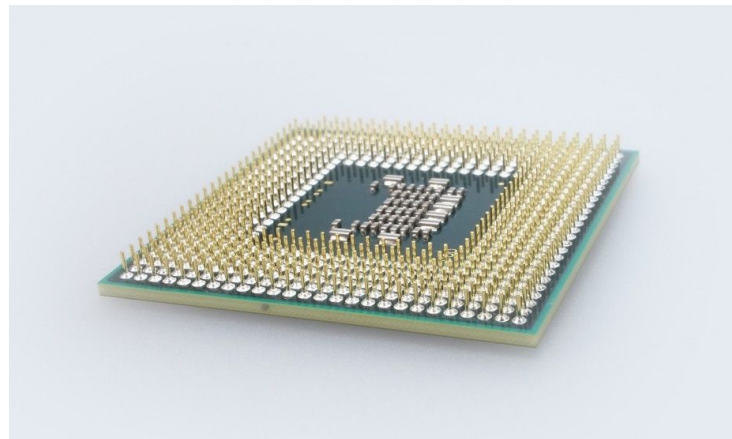
General-purpose programming

Designed for a wide-variety of use-cases

Designed for single-threaded operations

Operations happen synchronously

Can utilize multiple CPUs



Graphics Processing Unit (GPU)

Designed for parallel processing of information

Render graphics and mathematical computations

Most commonly used for distributed ML

Can utilize multiple GPUs



Tensor Processing Unit (TPU)

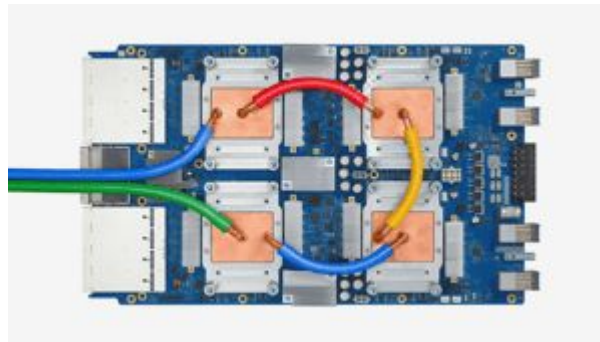
Designed specifically for ML use-cases

Built primarily for use with TensorFlow

Built by Google (who built TensorFlow)

Can access multiple via TPU pods

v1 2016

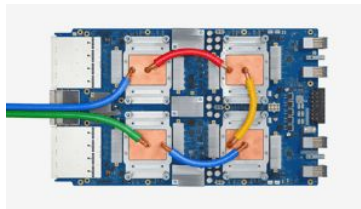
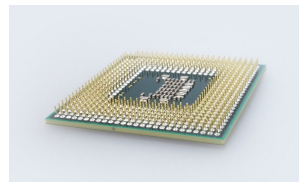


Which Should I Use?

CPU: Locally, parameter server, master

GPU: Mathematical computation

TPU: Optimized specifically for TensorFlow models



On-Prem vs The Cloud

Certain industries have stricter regulations around cloud-computing

Much easier to add more machines in the cloud vs on-prem

Maintaining physical hardware can be cumbersome

Paradigms

Single Node
Multi-core CPU



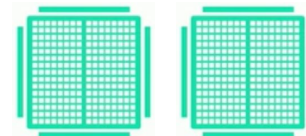
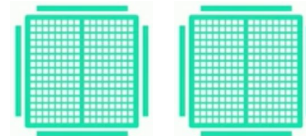
Single Node
Single GPU



Single Node
Multiple GPU



Multi-node
Multiple GPU



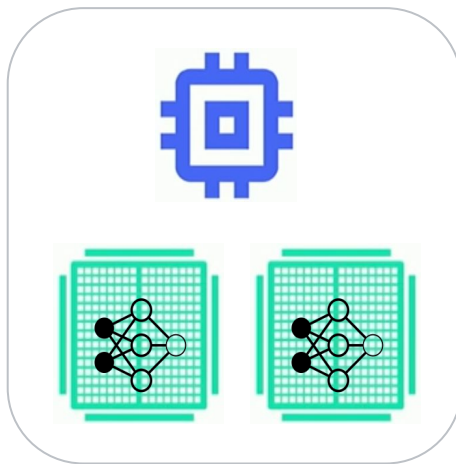
Data Parallelism

Data is partitioned across multiple machines

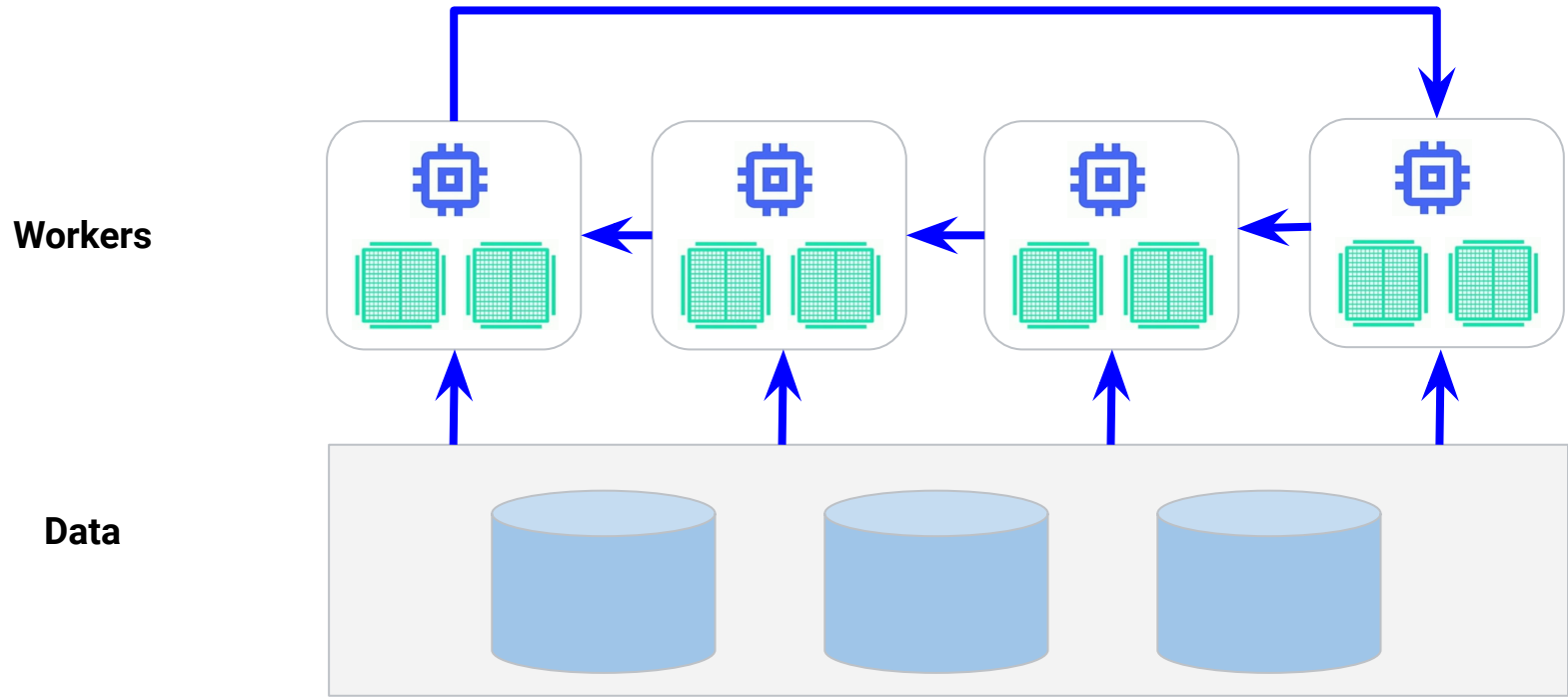
Processed in parallel

Each machine computes gradients / weights

Synchronous vs Asynchronous



SYNCHRONOUS

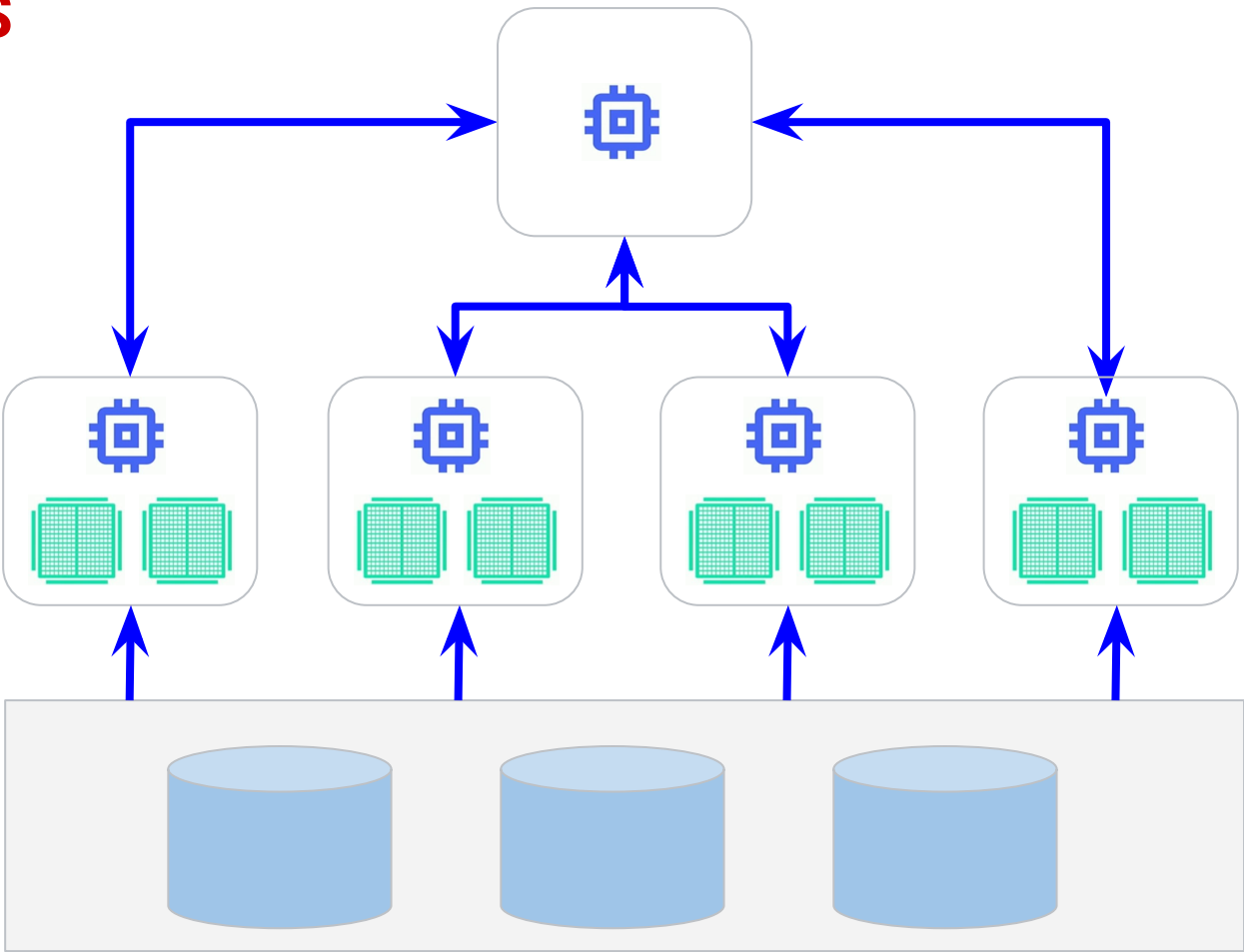


ASYNCHRONOUS

Parameter Server

Workers

Data

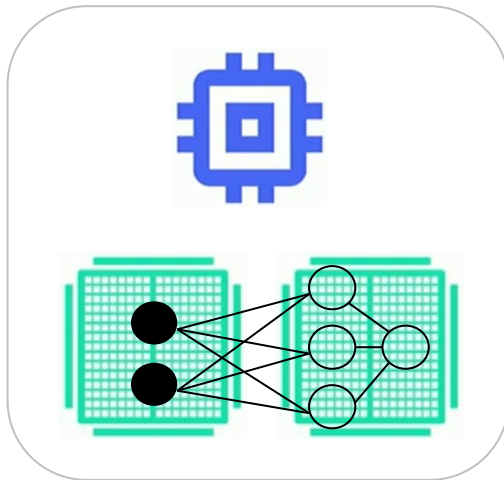


Model Parallelism

Model is split across multiple machines

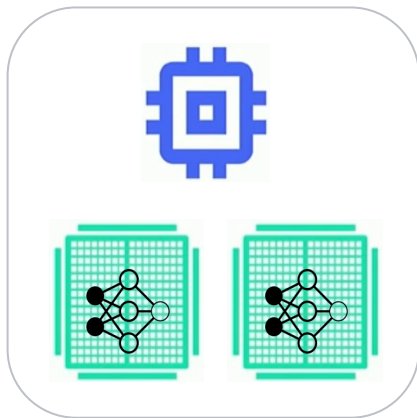
Process multiple layers in parallel

Commonly used for Deep Learning

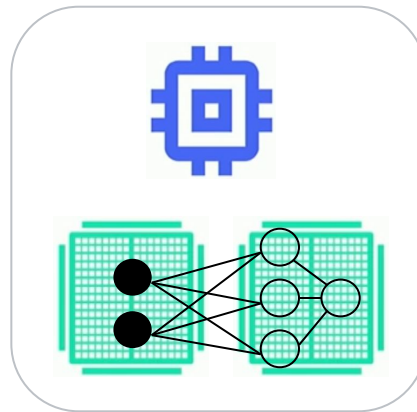


When to Use Which

Data Parallelism:
model **can** fit onto
a single machine



Model Parallelism:
model **cannot** fit onto
a single machine



Software Considerations

Why Python?

Rich ecosystem for scientific computing

Ease of use

Interacts with engines written in other languages (Java / C++)





Apache Spark

OSS “Unified analytics engine for large-scale data processing”

In-memory distributed data processing

Rich ecosystem - MLlib

Python, Java, Scala, and R

Abstracted parallelization



Ecosystem

**Spark
SQL**

**Spark
Streaming**

MLlib

GraphX

**SPARK
-NLP**

Apache Spark

Runtimes

**Spark
Standalone**

Yarn

Mesos

Kubernetes

```
spark = SparkSession.builder.appName("my_app").getOrCreate()

training = spark.read.format("csv").load("gs://my_bucket/my_data.csv")

lr = LinearRegression(maxIter=10, regParam=0.2)

model = lr.fit(training)

rmse = model.summary.rootMeanSquaredError

print(f"RMSE: {rmse}")
```



TensorFlow

OSS Machine Learning Framework

Rich ecosystem

Keras as high-level API

Easy to distribute

Python, Javascript, Swift, Java...



TensorFlow

TF Probability

TF Agents

TF Ranking

TF Text

TF Federated



TensorFlow Distribution Strategies

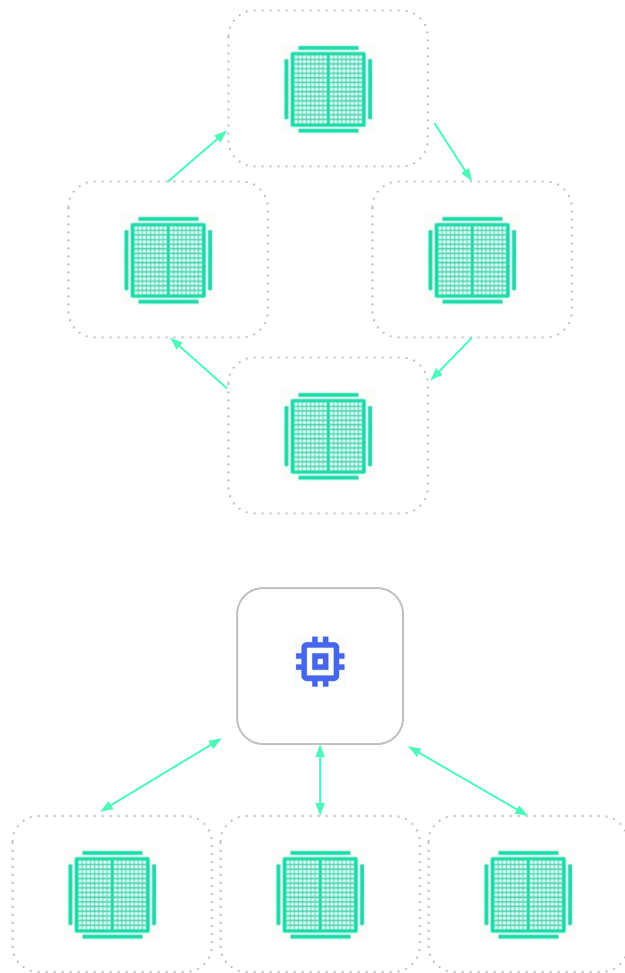
Designed to make ML distribution easy

Useful for researchers, practitioners, etc.

Provide good performance out of the box

Easy to switch between strategies

Little code changes




```
strategy = tf.distribute.MirroredStrategy()

with strategy.scope():
    model = tf.keras.models.Sequential([
        tf.keras.layers.Dense(64, input_shape=[10]),
        tf.keras.layers.Dense(64, activation='relu'),
        tf.keras.layers.Dense(10, activation='softmax')])

    model.compile(optimizer='adam',
                  loss='categorical_crossentropy',
                  metrics=['accuracy'])

    model.fit(training_data)
```

Closing Thoughts

Takeaways

Data Parallelism distributes
training data

CPUs are the brains

Apache Spark MLlib for
Machine Learning

Model Parallelism distributes
model parameters

GPUs and TPUs do the math

TensorFlow for Deep Learning

In the Cloud

Apache Spark - Cloud Dataproc



TensorFlow - AI Platform



Continued Resources

Large Scale Distributed Deep Networks - Dean et al. (2012)

Deep Learning with COTS HPC - Coates et al. (2013)

Spark: Cluster Computing with Working Sets - Zaharia et al. (2014)



Thank you!

Brad Miro - Google
@bradmiro
PyGotham 2019

Google Cloud