

# Distributed TensorFlow



**CSE545 - Spring 2020**  
Stony Brook University

H. Andrew Schwartz

# Big Data Analytics, The Class

**Goal: Generalizations**  
*A model or summarization of the data.*

*Data Frameworks*

Hadoop File System ✓  
Streaming ✓  
MapReduce ✓  
Spark ✓  
Tensorflow

*Algorithms and Analyses*

Similarity Search  
Graph Analysis  
Recommendation Systems  
Deep Learning  
Hypothesis Testing

# Limitations of Spark

Spark is fast for being so flexible

- Fast: RDDs in memory + Lazy evaluation: optimized chain of operations.
- Flexible: Many transformations -- can contain any custom code.

# Limitations of Spark

Spark is fast for being so flexible

- Fast: RDDs in memory + Lazy evaluation: optimized chain of operations.
- Flexible: Many transformations -- can contain any custom code.

However:

- Hadoop MapReduce can still be better for extreme IO, data that will not fit in memory across cluster.

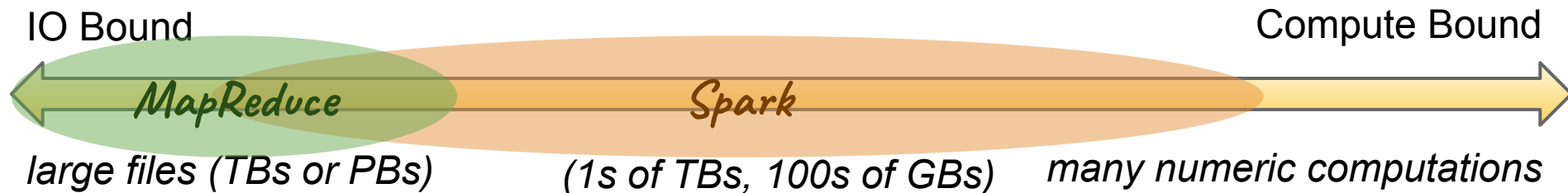
# Limitations of Spark

Spark is fast for being so flexible

- Fast: RDDs in memory + Lazy evaluation: optimized chain of operations.
- Flexible: Many transformations -- can contain any custom code.

However:

- Hadoop MapReduce can still be better for extreme IO, data that will not fit in memory across cluster.



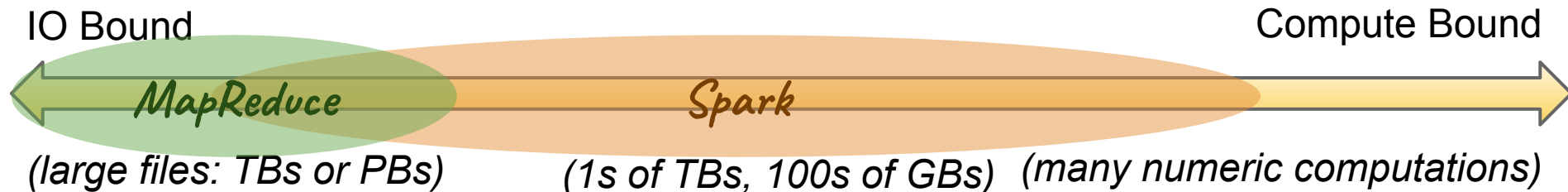
# Limitations of Spark

Spark is fast for being so flexible

- Fast: RDDs in memory + Lazy evaluation: optimized chain of operations.
- Flexible: Many transformations -- can contain any custom code.

However:

- Hadoop MapReduce can still be better for extreme IO, data that will not fit in memory across cluster.
- Modern machine learning (esp. Deep learning), a common big data task, requires heavy numeric computation.



\* this is the subjective approximation of the instructor as of February 2020. A lot of factors at play.

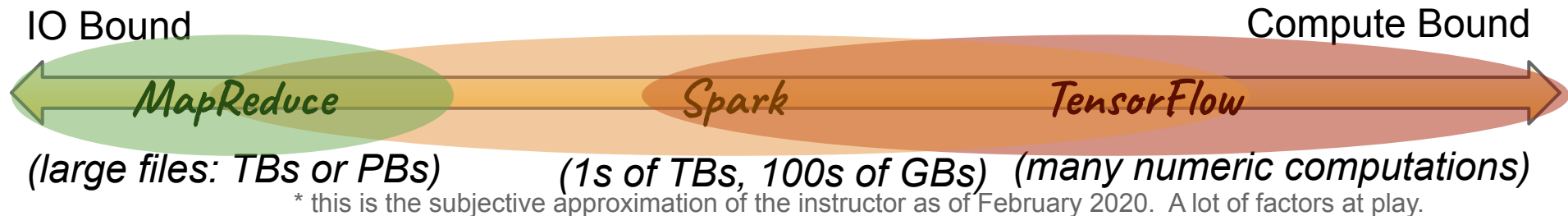
# Limitations of Spark

Spark is fast for being so flexible

- Fast: RDDs in memory + Lazy evaluation: optimized chain of operations.
- Flexible: Many transformations -- can contain any custom code.

However:

- Hadoop MapReduce can still be better for extreme IO, data that will not fit in memory across cluster.
- Modern machine learning (esp. Deep learning), a common big data task, requires heavy numeric computation.



# Learning Objectives

- Understand TensorFlow as a data workflow system.
  - Know the key components of TensorFlow.
  - Understand the key concepts of *distributed* TensorFlow.
- Execute basic distributed tensorflow program.
- Establish a foundation to distribute deep learning models:
  - Convolutional Neural Networks
  - Recurrent Neural Network (or LSTM, GRU)



# What is TensorFlow?

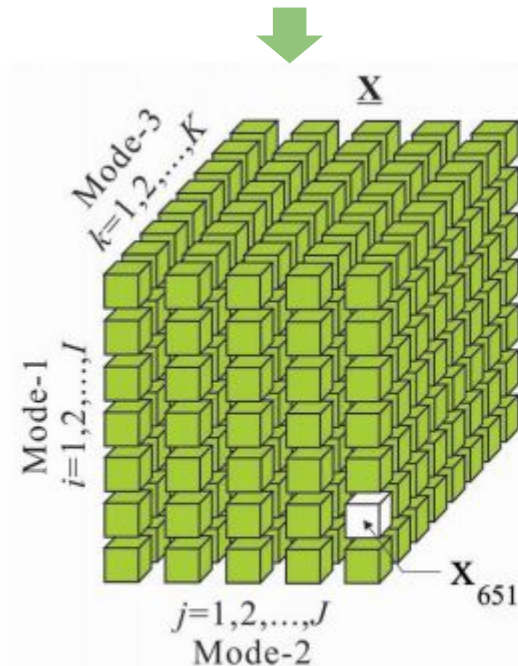
A workflow system catered to numerical computation.

One view: Like Spark, but uses *tensors* instead of *RDDs*.

# What is TensorFlow?

A workflow system catered to numerical computation.

One view: Like Spark, but uses *tensors* instead of *RDDs*.



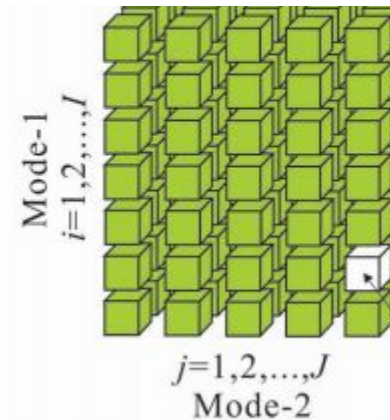
➔ A multi-dimensional matrix

(i.stack.imgur.com)

# What is TensorFlow?

A workflow system catered to numerical computation.

One view: Like Spark, but uses *tensors* instead of *RDDs*.



A 2-d tensor is just a matrix.

1-d: vector

0-d: a constant / scalar

Note: Linguistic ambiguity:  
Dimensions of a Tensor  $\neq$   
Dimensions of a Matrix

# What is TensorFlow?

A workflow system catered to numerical computation.

One view: Like Spark, but uses *tensors* instead of *RDDs*.



Examples > 2-d :

Image definitions in terms of RGB per pixel

Image[*row*][*column*][*rgb*]

Subject, Verb, Object representation of language:

Counts[*verb*][*subject*][*object*]

# What is TensorFlow?

A workflow system catered to numerical computation.

One view: Like Spark, but uses *tensors* instead of *RDDs*.



Technically, less abstract than *RDDs* which could hold tensors as well as many other data structures (dictionaries/HashMaps, Trees, ...etc...).

Then, why TensorFlow?

# What is TensorFlow?

Efficient, high-level built-in **linear algebra** and **machine learning optimization operations** (i.e. transformations).

enables complex models, like deep learning

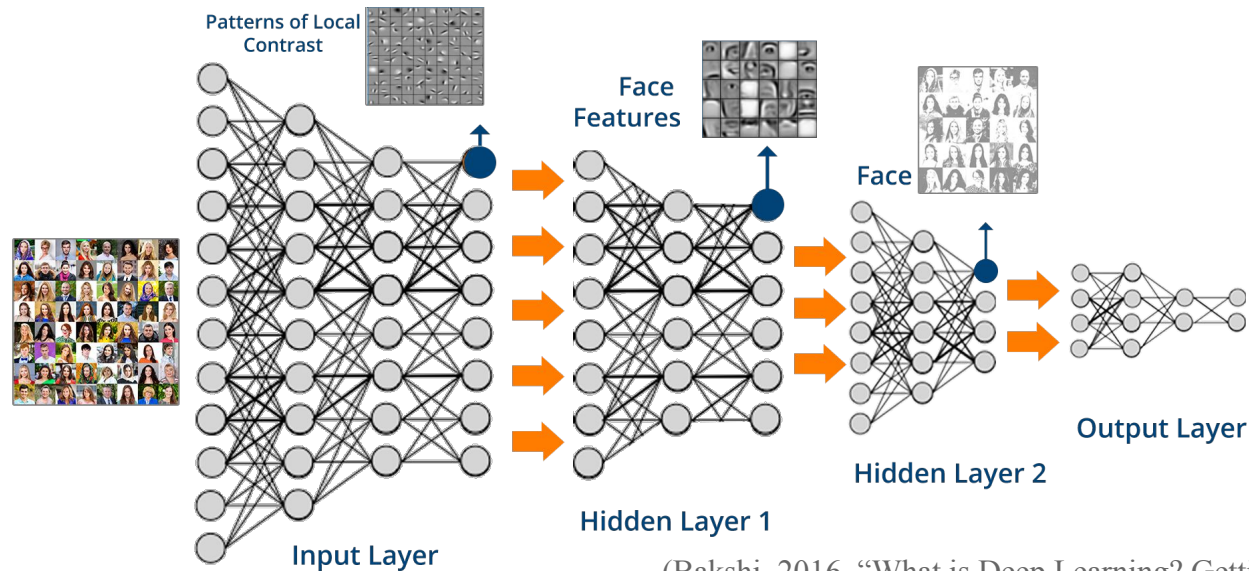


**Then, why TensorFlow?**

# What is TensorFlow?

Efficient, high-level built-in **linear algebra** and **machine learning optimization operations**.

enables complex models, like deep learning



(Bakshi, 2016, "What is Deep Learning? Getting Started With Deep Learning")

# What is TensorFlow?

Efficient, high-level built-in **linear algebra** and **machine learning operations**.

```
import tensorflow as tf

b = tf.Variable(tf.zeros([100])) # 100-d vector, init to zeroes
W = tf.Variable(tf.random_uniform([784,100],-1,1)) # 784x100 matrix w/rnd vals
x = tf.placeholder(name="x") # Placeholder for input
relu = tf.nn.relu(tf.matmul(W, x) + b) # Relu(Wx+b)
C = [...] # Cost computed as a function # of Relu

s = tf.Session()
for step in xrange(0, 10):
    input = ...construct 100-D input array ... # Create 100-d vector for input
    result = s.run(C, feed_dict={x: input}) # Fetch cost, feeding x=input
    print step, result
```



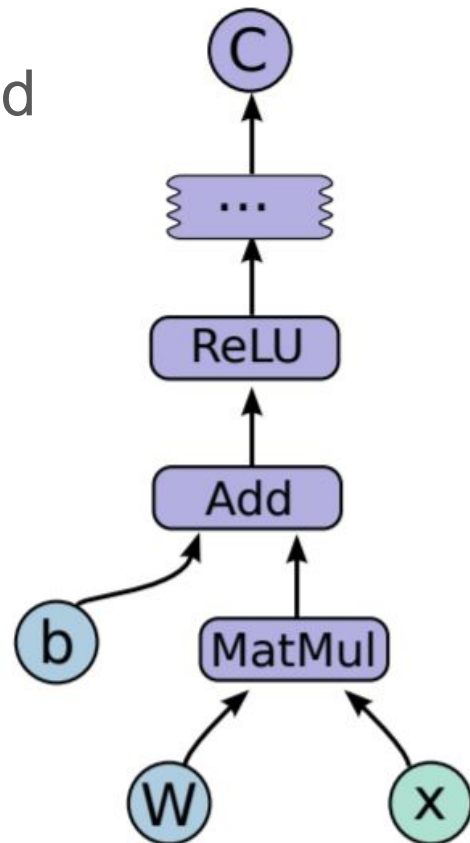
# TensorFlow

Operations on tensors are often conceptualized as **graphs**:

```
import tensorflow as tf

b = tf.Variable(tf.zeros([100]))
W = tf.Variable(tf.random_uniform([784,100],-1,1))
x = tf.placeholder(name="x")
relu = tf.nn.relu(tf.matmul(W, x) + b)
C = [...]

s = tf.Session()
for step in xrange(0, 10):
    input = ...construct 100-D input array ...
    result = s.run(C, feed_dict={x: input})
    print step, result
```

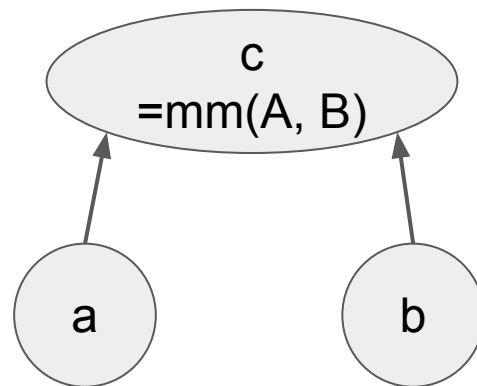


# TensorFlow

Operations on tensors are often conceptualized as graphs:

A simpler example:

```
c = tensorflow.matmul(a, b)
```



# TensorFlow

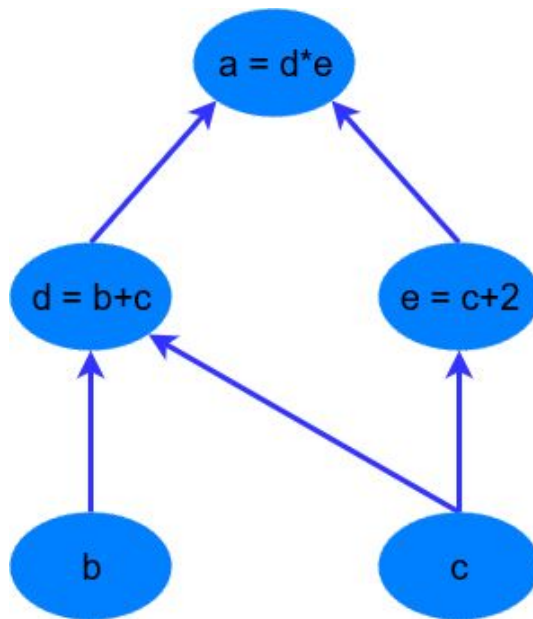
Operations on tensors are often conceptualized as graphs:

example:

$$d = b + c$$

$$e = c + 2$$

$$a = d * e$$



(Adventures in Machine Learning.  
*Python TensorFlow Tutorial*, 2017)

# Ingredients of a TensorFlow

## ***tensors\****

*variables* - persistent  
mutable tensors  
*constants* - constant  
*placeholders* - from data

## ***operations***

an abstract computation  
(e.g. matrix multiply, add)  
executed by device *kernels*

\* technically, still *operations*

***graph***

## ***session***

defines the environment in  
which operations *run*.  
(like a Spark context)

## ***devices***

the specific devices (cpus or  
gpus) on which to run the  
session.

# Ingredients of a TensorFlow

## ***tensors\****

*variables* - persistent  
mutable tensors  
*constants* - constant  
*placeholders* - from data

- `tf.Variable(initial_value, name)`
- `tf.constant(value, type, name)`
- `tf.placeholder(type, shape, name)`

*operations*  
an abstract computation  
(e.g. matrix multiply, add)  
executed by device *kernels*

\* technically, still *operations*

*graph*

## ***session***

defines the environment in  
which operations *run*.  
(like a Spark context)

## ***devices***

the specific devices (cpus or  
gpus) on which to run the  
session.

# Ingredients of a TensorFlow

## *tensors\**

*variables* - persistent  
mutable tensors  
*constants* - constant  
*placeholders* - from data

## *operations*

an abstract computation  
(e.g. matrix multiply, add)  
executed by device *kernels*

Category	Examples
Element-wise mathematical operations	Add, Sub, Mul, Div, Exp, Log, Greater, Less, Equal, ...
Array operations	Concat, Slice, Split, Constant, Rank, Shape, Shuffle, ...
Matrix operations	MatMul, MatrixInverse, MatrixDeterminant, ...
Stateful operations	Variable, Assign, AssignAdd, ...
Neural-net building blocks	SoftMax, Sigmoid, ReLU, Convolution2D, MaxPool, ...
Checkpointing operations	Save, Restore
Queue and synchronization operations	Enqueue, Dequeue, MutexAcquire, MutexRelease, ...
Control flow operations	Merge, Switch, Enter, Leave, NextIteration

# Ingredients of a TensorFlow

## *tensors\**

- Places operations on devices

*variables* - persistent

- Stores the values of variables (when not distributed)

*constants* - constant

- Carries out execution: `eval()` or `run()`

*placeholders* - from data

## *operations*

an abstract computation

(e.g. matrix multiply, add)

executed by device *kernels*

## *graph*

## *session*

defines the environment in which operations *run*.  
(like a Spark context)

## *devices*

the specific devices (cpus or gpus) on which to run the session.

# Ingredients of a TensorFlow

## ***tensors\****

*variables* - persistent  
mutable tensors  
*constants* - constant  
*placeholders* - from data

## ***operations***

an abstract computation  
(e.g. matrix multiply, add)  
executed by device *kernels*

***graph***

## ***session***

defines the environment in  
which operations *run*.  
(like a Spark context)

## ***devices***

the specific devices (cpus or  
gpus) on which to run the  
session.



# Distributed TensorFlow

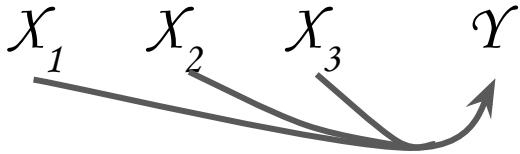
Typical use-case: (Supervised Machine Learning)

Determine weights,  $\mathcal{W}$ , of a function,  $f$ , such that  $\epsilon$  is minimized:  $f(\mathcal{X} | \mathcal{W}) = \mathcal{Y} + \epsilon$

# Distributed TensorFlow

Typical use-case:

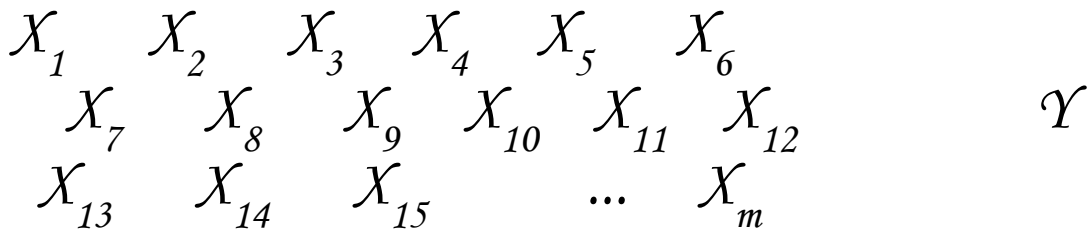
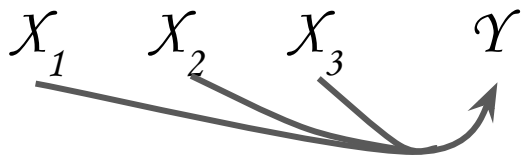
Determine weights,  $\mathcal{W}$ , of a function,  $f$ , such that  $\epsilon$  is minimized:  $f(\mathcal{X} | \mathcal{W}) = \mathcal{Y} + \epsilon$



# Distributed TensorFlow

Typical use-case:

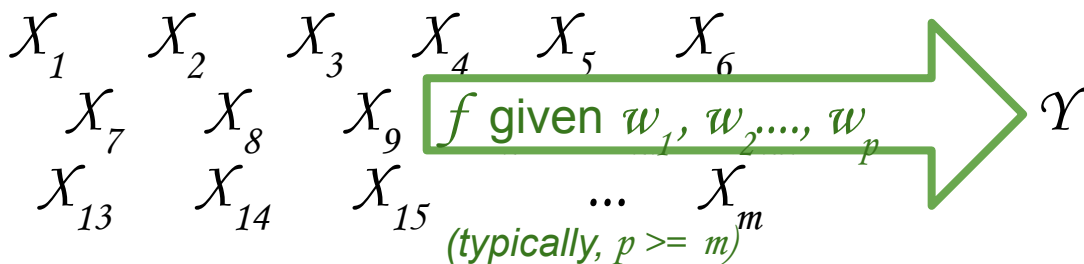
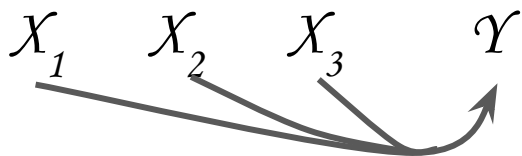
Determine weights,  $\mathcal{W}$ , of a function,  $f$ , such that  $\epsilon$  is minimized:  $f(\mathcal{X} | \mathcal{W}) = \mathcal{Y} + \epsilon$



# Distributed TensorFlow

Typical use-case:

Determine weights,  $\mathcal{W}$ , of a function,  $f$ , such that  $\epsilon$  is minimized:  $f(\mathcal{X} | \mathcal{W}) = \mathcal{Y} + \epsilon$

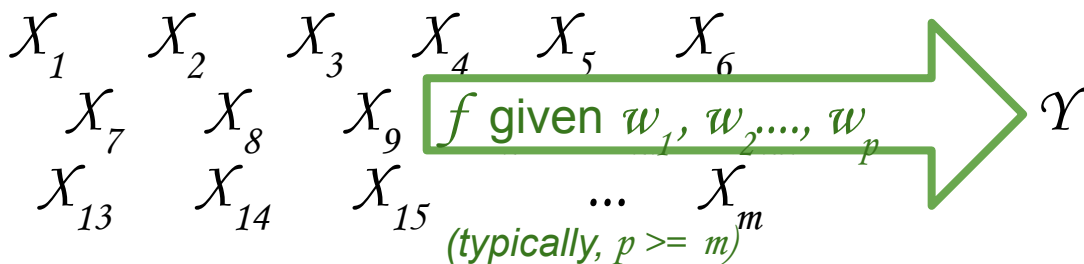
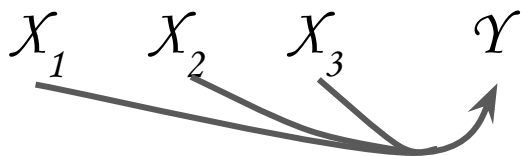


# Distributed TensorFlow

Typical use-case:

Determine weights,  $\mathcal{W}$ , of a function,  $f$ , such that  $|\epsilon|$  is minimized:

$$\begin{aligned}f(X/W) &= \hat{Y} \\ Y &= (X/W) + \epsilon \\ Y &= \hat{Y} + \epsilon \\ \epsilon &= \hat{Y} - Y\end{aligned}$$

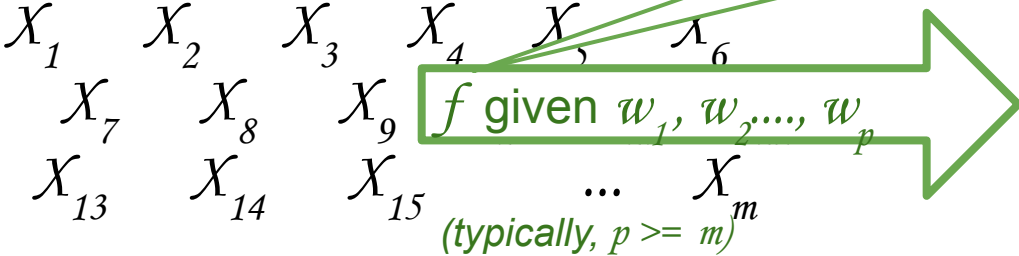
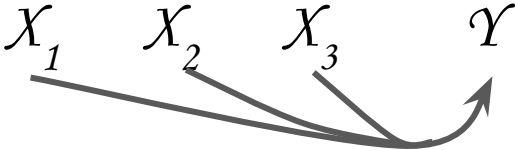


# Distributed TensorFlow

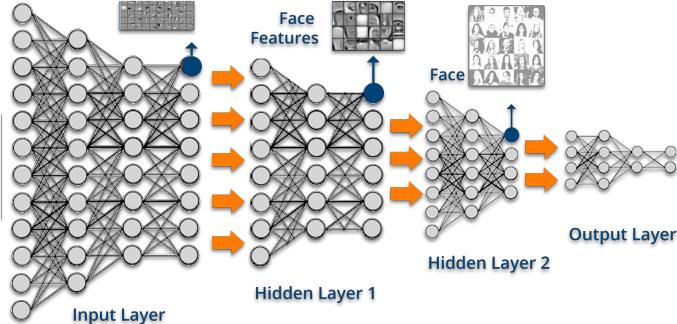
Typical use-case:

Determine weights,  $W$ , of a function,  $f$ , such that  $\epsilon$  is minimized:

$$\begin{aligned} f(X/W) &= \hat{Y} \\ Y &= (X/W) + \epsilon \\ Y &= \hat{Y} + \epsilon \\ \epsilon &= \hat{Y} - Y \end{aligned}$$



Typically, very complex!



# Distributed TensorFlow

Typical use-case:

Determine weights,  $W$ , of a function,  $f$ , such that  $\epsilon$  is minimized:

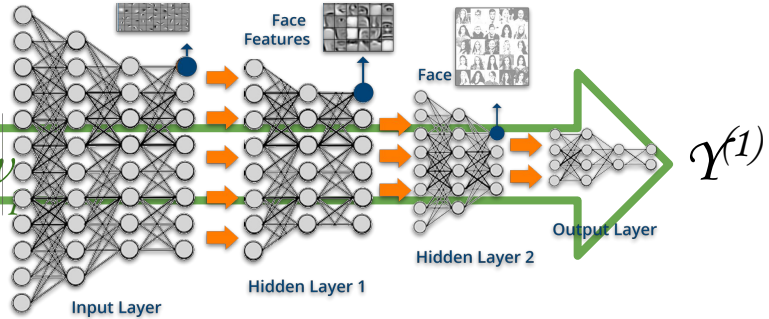
$W$  determined through *gradient descent*:

*back propagating* error across the network that defines  $f$ .

$$\begin{aligned} f(X/W) &= \hat{Y} \\ Y &= (X/W) + \epsilon \\ Y &= \hat{Y} + \epsilon \\ \epsilon &= \hat{Y} - Y \end{aligned}$$

$X_1^{(1)}$     $X_2$   
 $X_7$     $X_8$   
 $X_{13}$     $X_{14}$   
 $X_{15}$

$X_3$     $X_4$     $X_5$     $X_6$   
 $X_9$     $f$  given  $w_1, w_2, \dots, w_m$   
 $X_{15}$     $\dots$     $X_m$   
(typically,  $p \geq m$ )



# Distributed TensorFlow

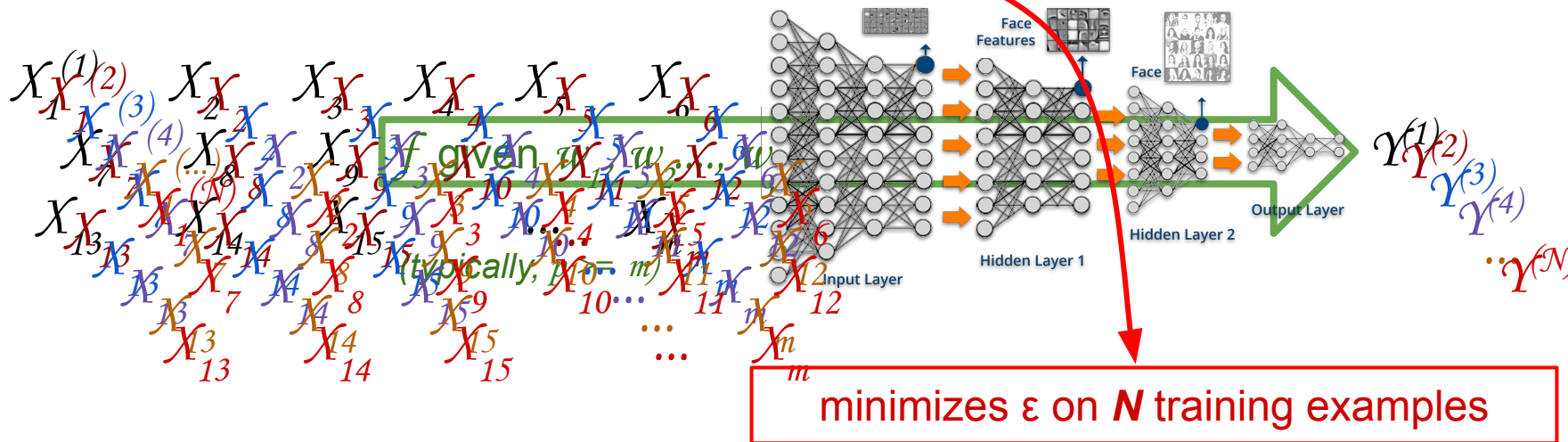
Typical use-case:

Determine weights,  $W$ , of a function,  $f$ , such that  $\epsilon$  is minimized:

$$\begin{aligned} f(X/W) &= \hat{Y} \\ Y &= (X/W) + \epsilon \\ Y &= \hat{Y} + \epsilon \\ \epsilon &= \hat{Y} - Y \end{aligned}$$

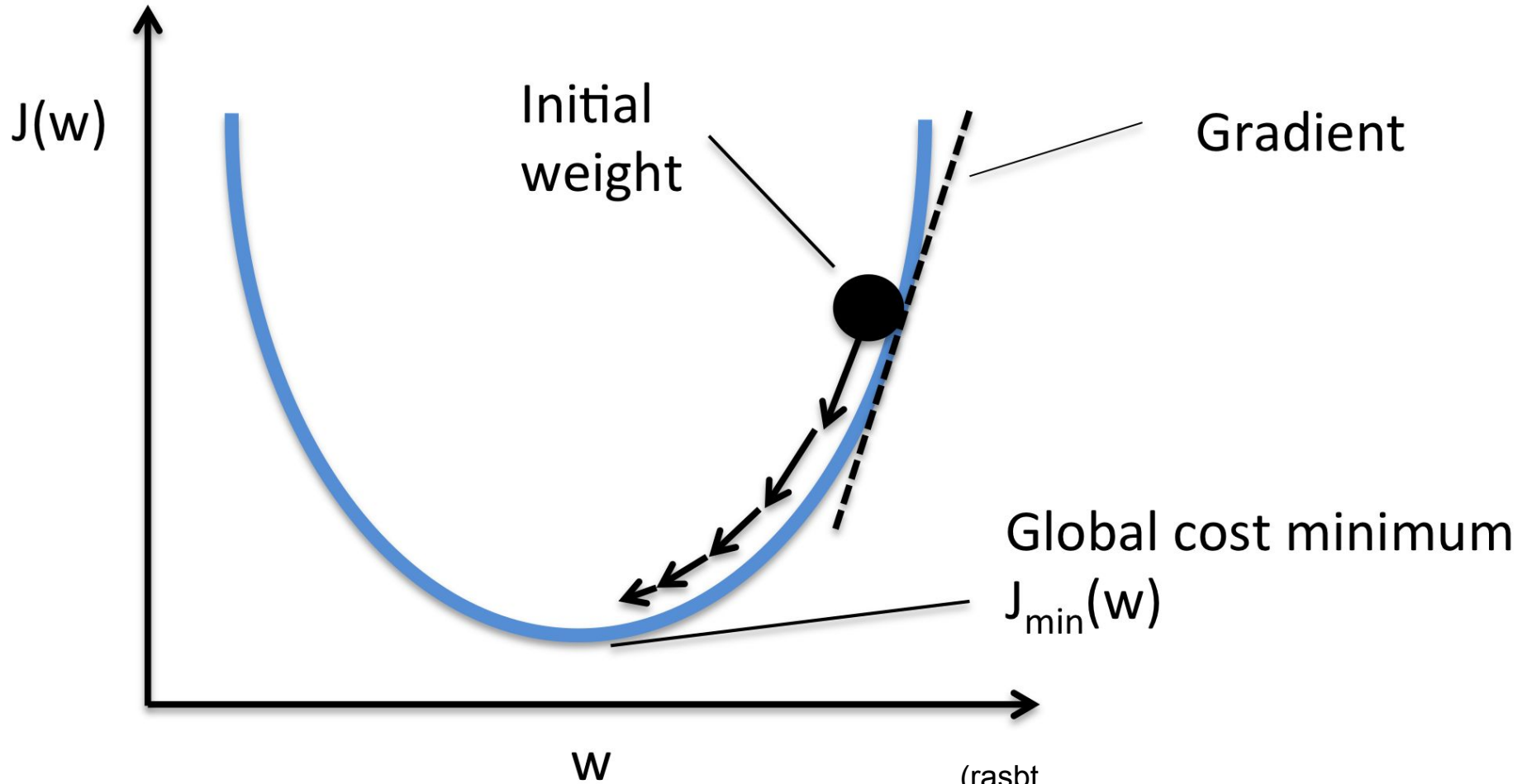
$W$  determined through *gradient descent*:

*back propagating* error across the network that defines  $f$ .





# Weights Derived from Gradients



# Weights Derived from Gradients

**Linear Regression:** Trying to find “betas” that minimize:

$$\hat{\beta} = \operatorname{argmin}_{\beta} \left\{ \sum_i^N (y_i - \hat{y}_i)^2 \right\}$$

# Weights Derived from Gradients

**Linear Regression:** Trying to find “betas” that minimize:

$$\hat{\beta} = \operatorname{argmin}_{\beta} \left\{ \sum_i^N (y_i - \hat{y}_i)^2 \right\}$$

*matrix multiply*

$$\hat{y}_i = X_i \beta$$

Thus:

$$\hat{\beta} = \operatorname{argmin}_{\beta} \left\{ \sum_{i=0}^N (y_i - X_i \beta)^2 \right\}$$

# Weights Derived from Gradients

**Linear Regression:** Trying to find “betas” that minimize:

$$\hat{\beta} = \operatorname{argmin}_{\beta} \left\{ \sum_i^N (y_i - \hat{y}_i)^2 \right\}$$

*matrix multiply*

$$\hat{y}_i = X_i \beta$$

Thus:

$$\hat{\beta} = \operatorname{argmin}_{\beta} \left\{ \sum_{i=0}^N (y_i - X_i \beta)^2 \right\}$$

In standard linear equation:

$$y = mx + b \quad \text{let } x' = x + [1, 1, \dots, 1]_N^T$$

then,  $y = mx'$

(if we add a column of 1s,  $mx + b$  is just  $\text{matmul}(m, x)$ )

# Weights Derived from Gradients

**Linear Regression:** Trying to find “betas” that minimize:

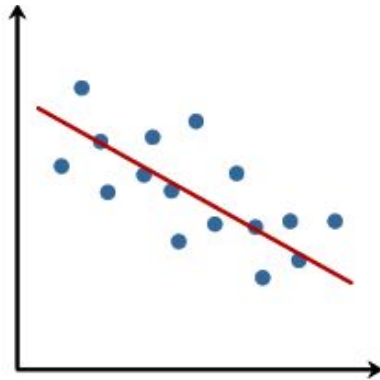
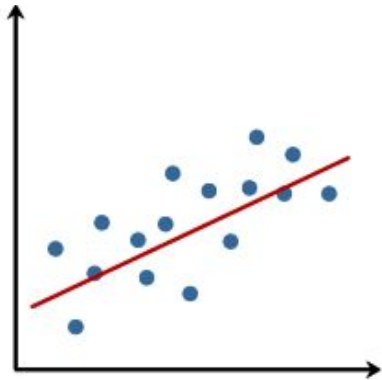
$$\hat{\beta} = \operatorname{argmin}_{\beta} \left\{ \sum_i^N (y_i - \hat{y}_i)^2 \right\}$$

*matrix multiply*

$$\hat{y}_i = X_i \beta$$

Thus:

$$\hat{\beta} = \operatorname{argmin}_{\beta} \left\{ \sum_{i=0}^N (y_i - X_i \beta)^2 \right\}$$



# Weights Derived from Gradients

**Linear Regression:** Trying to find “betas” that minimize:

$$\hat{\beta} = \operatorname{argmin}_{\beta} \left\{ \sum_i^N (y_i - \hat{y}_i)^2 \right\}$$

$$\hat{y}_i = X_i \beta \quad \text{Thus:} \quad \hat{\beta} = \operatorname{argmin}_{\beta} \left\{ \sum_{i=0}^N (y_i - X_i \beta)^2 \right\}$$

How to update?  $\beta_{\text{new}} = \beta_{\text{prev}} - \alpha * \text{grad}$

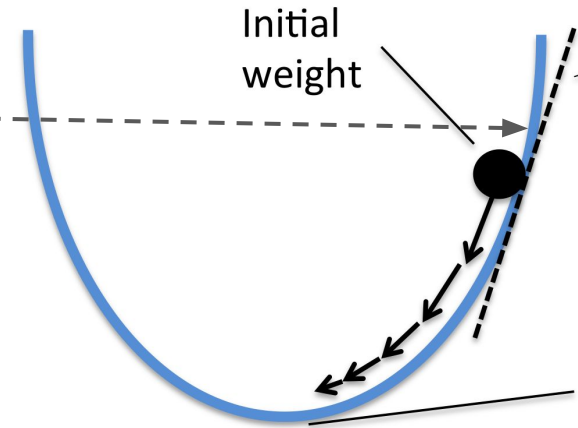
# Weights Derived from Gradients

**Linear Regression:** Trying to find “betas” that minimize:

$$\hat{\beta} = \operatorname{argmin}_{\beta} \left\{ \sum_i^N (y_i - \hat{y}_i)^2 \right\}$$

$$\hat{y}_i = X_i \beta \quad \text{Thus:} \quad \hat{\beta} = \operatorname{argmin}_{\beta} \left\{ \sum_{i=0}^N (y_i - X_i \beta)^2 \right\}$$

How to update?  $\beta_{new} = \beta_{prev} - \alpha * \operatorname{grad}$   
(for gradient descent)      “learning rate”



# Weights Derived from Gradients

Ridge Regression (L2 Penalized linear regression,  $\lambda \|\beta\|_2^2$ )

$$\hat{\beta}^{\text{ridge}} = \operatorname{argmin}_{\beta} \left\{ \sum_{i=1}^N (y_i - \sum_{j=1}^m x_{ij} \beta_j)^2 + \lambda \sum_{j=1}^m \beta_j^2 \right\}$$

1. Matrix Solution:

$$\hat{\beta}^{\text{ridge}} = (X^T X + \lambda I)^{-1} X^T y$$



# Weights Derived from Gradients

Ridge Regression (L2 Penalized linear regression,  $\lambda ||\beta||_2^2$ )

$$\hat{\beta}^{\text{ridge}} = \operatorname{argmin}_{\beta} \left\{ \sum_{i=1}^N (y_i - \sum_{j=1}^m x_{ij} \beta_j)^2 + \lambda \sum_{j=1}^m \beta_j^2 \right\}$$

## 2. Gradient descent solution

(Mirrors many parameter optimization problems.)

## 1. Matrix Solution:

$$\hat{\beta}^{\text{ridge}} = (X^T X + \lambda I)^{-1} X^T y$$

# Weights Derived from Gradients

Ridge Regression (L2 Penalized linear regression,  $\lambda ||\beta||_2^2$ )

$$\hat{\beta}^{\text{ridge}} = \operatorname{argmin}_{\beta} \left\{ \sum_{i=1}^N (y_i - \sum_{j=1}^m x_{ij} \beta_j)^2 + \lambda \sum_{j=1}^m \beta_j^2 \right\}$$

**Gradient descent** needs to solve.

(Mirrors many parameter optimization problems.)

TensorFlow has built-in ability to derive gradients given a **cost function**.

# Weights Derived from Gradients

Ridge Regression (L2 Penalized linear regression,  $\lambda ||\beta||_2^2$ )

$$\hat{\beta}^{\text{ridge}} = \operatorname{argmin}_{\beta} \left\{ \sum_{i=1}^N (y_i - \sum_{j=1}^m x_{ij} \beta_j)^2 + \lambda \sum_{j=1}^m \beta_j^2 \right\}$$

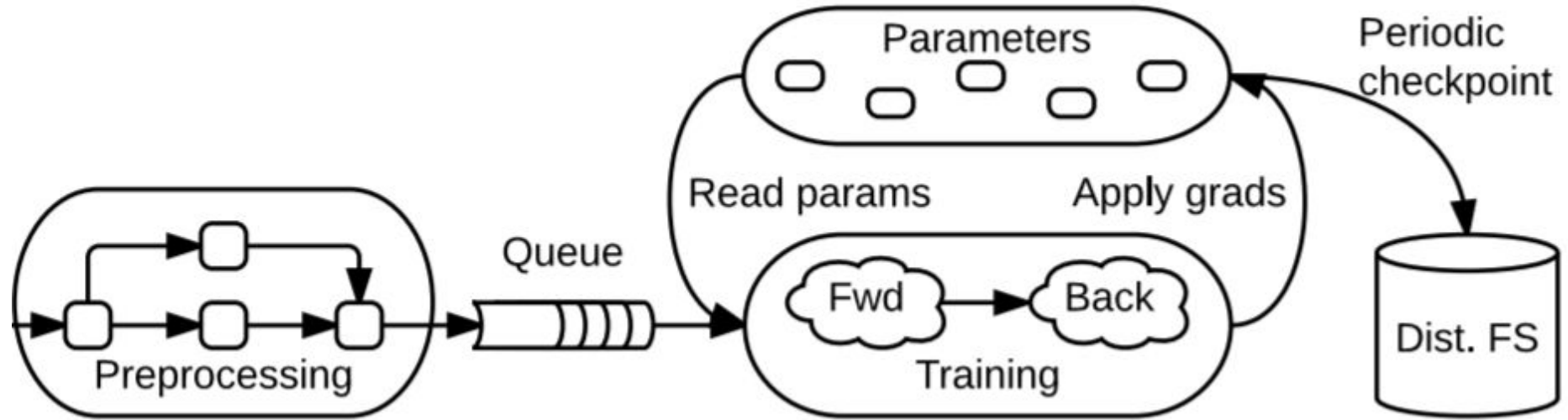
Gradient descent needs to solve.

(Mirrors many parameter optimization problems.)

TensorFlow has built-in ability to derive gradients given a cost function.

```
tf.gradients(cost, [params])
```

# Weights Derived from Gradients



TensorFlow has built-in ability to derive gradients given a cost function.

```
tf.gradients(cost, [params])
```

# Options for distribution

1. **Distribute copies of entire dataset**
  - a. Train over all with different hyperparameters
  - b. Train different folds per worker node.

# Options for distribution

1. **Distribute copies of entire dataset**
  - a. Train over all with different parameters
  - b. Train different folds per worker node.
  
2. **Distribute data**
  - a. Each node finds parameters for subset of data
  - b. Needs mechanism for updating parameters
    - i. Centralized parameter server
    - ii. Distributed All-Reduce

# Options for distribution

## 1. Distribute copies of entire dataset

- a. Train over all with different parameters
- b. Train different folds per worker node.

## 2. Distribute data

- a. Each node finds parameters for subset of data
- b. Needs mechanism for updating parameters
  - i. Centralized parameter server
  - ii. Distributed All-Reduce

## 3. Distribute model or individual operations (e.g. matrix multiply)

# Options for distribution

## 1. Distribute copies of entire dataset

- a. Train over all with different parameters
- b. Train different folds per worker node.

Pro: Easy; Good for compute-bound; Con: Requires data fit in worker memories

## 2. Distribute data

- a. Each node finds parameters for subset of data
- b. Needs mechanism for updating parameters
  - i. Centralized parameter server
  - ii. Distributed All-Reduce

Pro: Flexible to all situations; Con: Optimizing for subset is suboptimal

## 3. Distribute model or individual operations (e.g. matrix multiply)

Pro: Parameters can be localized Con: High communication for transferring Intermediar data.



# Options for distribution

Done often in practice. Not talked about much because it's mostly as easy as it sounds.

## 1. Distribute copies of entire dataset

- a. Train over all with different parameters
- b. Train different folds per worker node.

Pro: Easy; Good for compute-bound; Con: Requires data fit in worker memories

## 2. Distribute data

- a. Each node finds parameters for subset of data
- b. Needs mechanism for updating parameters
  - i. Centralized parameter server
  - ii. Distributed All-Reduce

Pro: Flexible to all situations; Con: Optimizing for subset is suboptimal

## 3. Distribute model or individual operations (e.g. matrix multiply)

Pro: Parameters can be localized Con: High communication for transferring Intermediar data.

# Options for distribution

Done often in practice. Not talked about much because it's mostly as easy as it sounds.

## 1. Distribute copies of entire dataset

- a. Train over all with different parameters
- b. Train different folds per worker node.

Pro: Easy; Good for compute-bound; Con: Requires data fit in worker memories

## 2. Distribute data

- a. Each node finds parameters for subset of data
- b. Needs mechanism for updating parameters
  - i. Centralized parameter server
  - ii. Distributed All-Reduce

Preferred method for big data or very complex models (i.e. models with many internal parameters).

Pro: Flexible to all situations; Con: Optimizing for subset is suboptimal

## 3. Distribute model or individual operations (e.g. matrix multiply)

Pro: Parameters can be localized Con: High communication for transferring Intermediar data.

# Options for distribution

Done often in practice. Not talked about much because it's mostly as easy as it sounds.

1. **Distribute copies of entire dataset**
  - a. Train over all with different parameters
  - b. Train different folds per worker node.

Pro: Easy; Good for compute-bound; Con: Requires data fit in worker memories

2. **Distribute data**
  - a. Each node finds parameters for subset of data
  - b. Needs mechanism for updating parameters
    - i. Centralized parameter server
    - ii. Distributed All-Reduce

## Data Parallelism

Preferred method for big data or very complex models (i.e. models with many internal parameters).

Pro: Flexible to all situations; Con: Optimizing for subset is suboptimal

3. **Distribute model or individual operations** (e.g. matrix multiply)

Pro: Parameter servers can be shared

## Model Parallelism

Con: High communication for transferring Intermediar data.

# Model Parallelism

Multiple devices on multiple machines

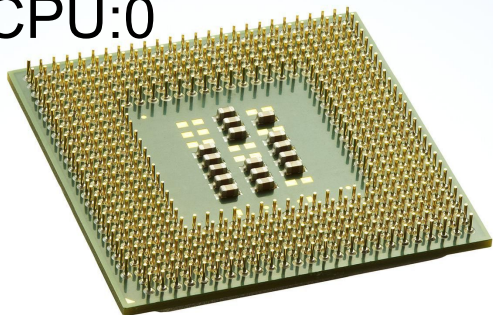
```
with tf.device("/cpu:1")  
    beta=tf.Variable(...)
```

```
with tf.device("/gpu:0")  
    y_pred=tf.matmul(beta,X)
```

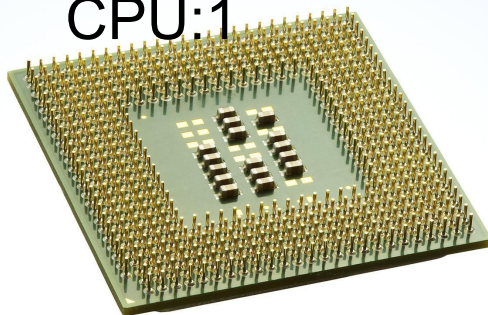
Transfer Tensors

Machine A

CPU:0

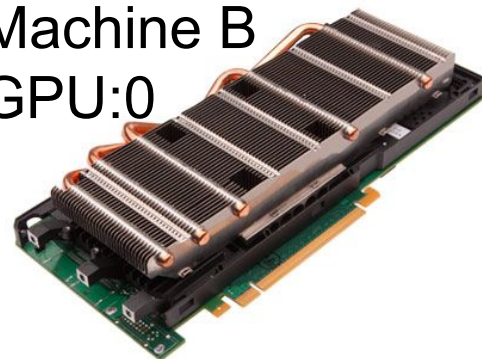


CPU:1



Machine B

GPU:0

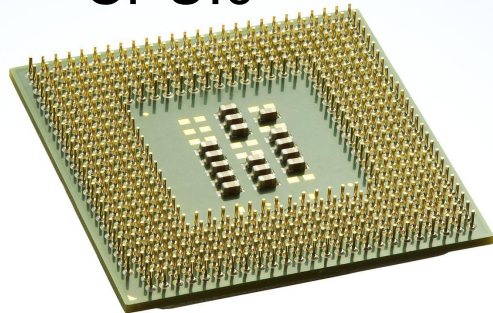


# Data Parallelism

```
...  
beta=tf.Variable(...)  
pred=tf.matmul(beta,X)
```



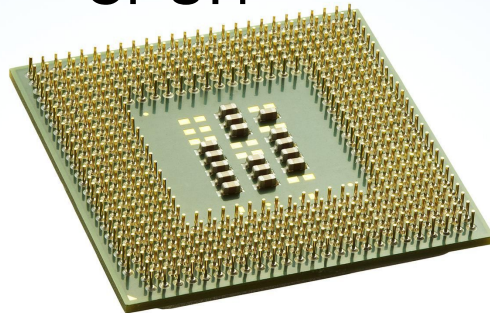
CPU:0



```
...  
beta=tf.Variable(...)  
pred=tf.matmul(beta,X)
```

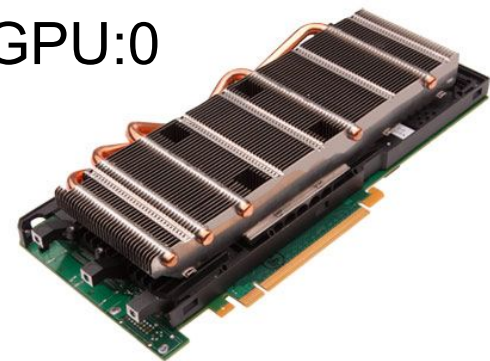


CPU:1



```
...  
beta=tf.Variable(...)  
pred=tf.matmul(beta,X)
```

GPU:0



# Data Parallelism

```
...  
beta=tf.Variable(...)  
pred=tf.matmul(beta,X)
```

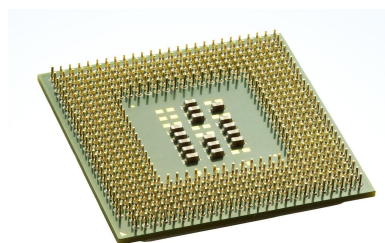
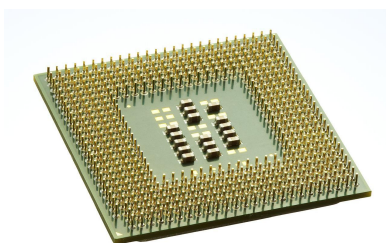
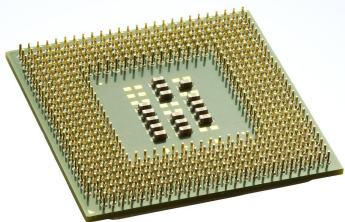
```
...  
beta=tf.Variable(...)  
pred=tf.matmul(beta,X)
```

```
...  
beta=tf.Variable(...)  
pred=tf.matmul(beta,X)
```

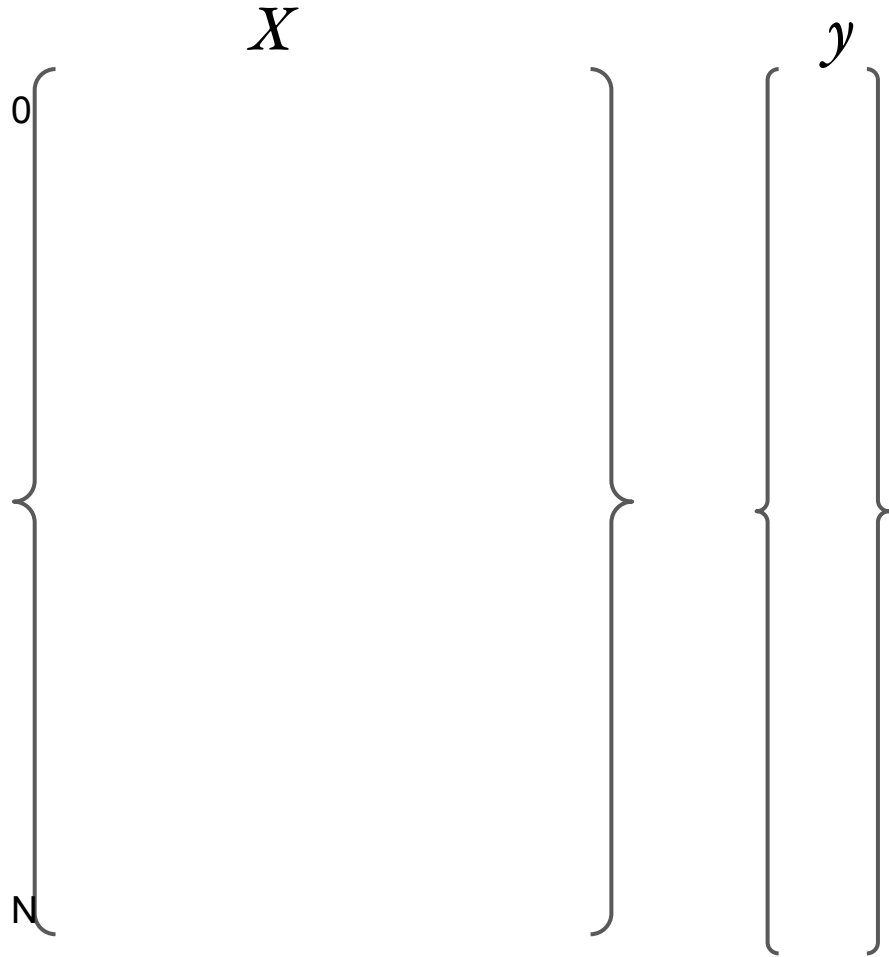
worker:0

worker:1

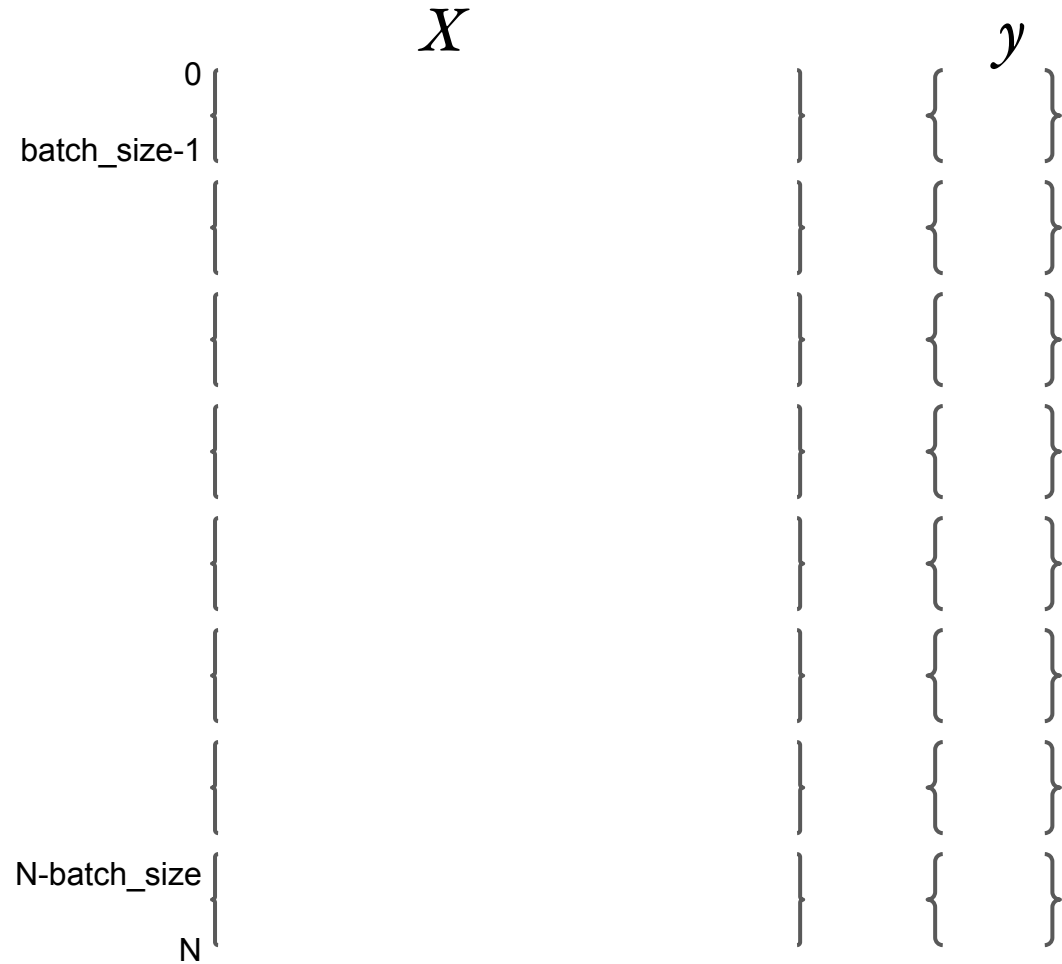
worker:2



# Distributing Data



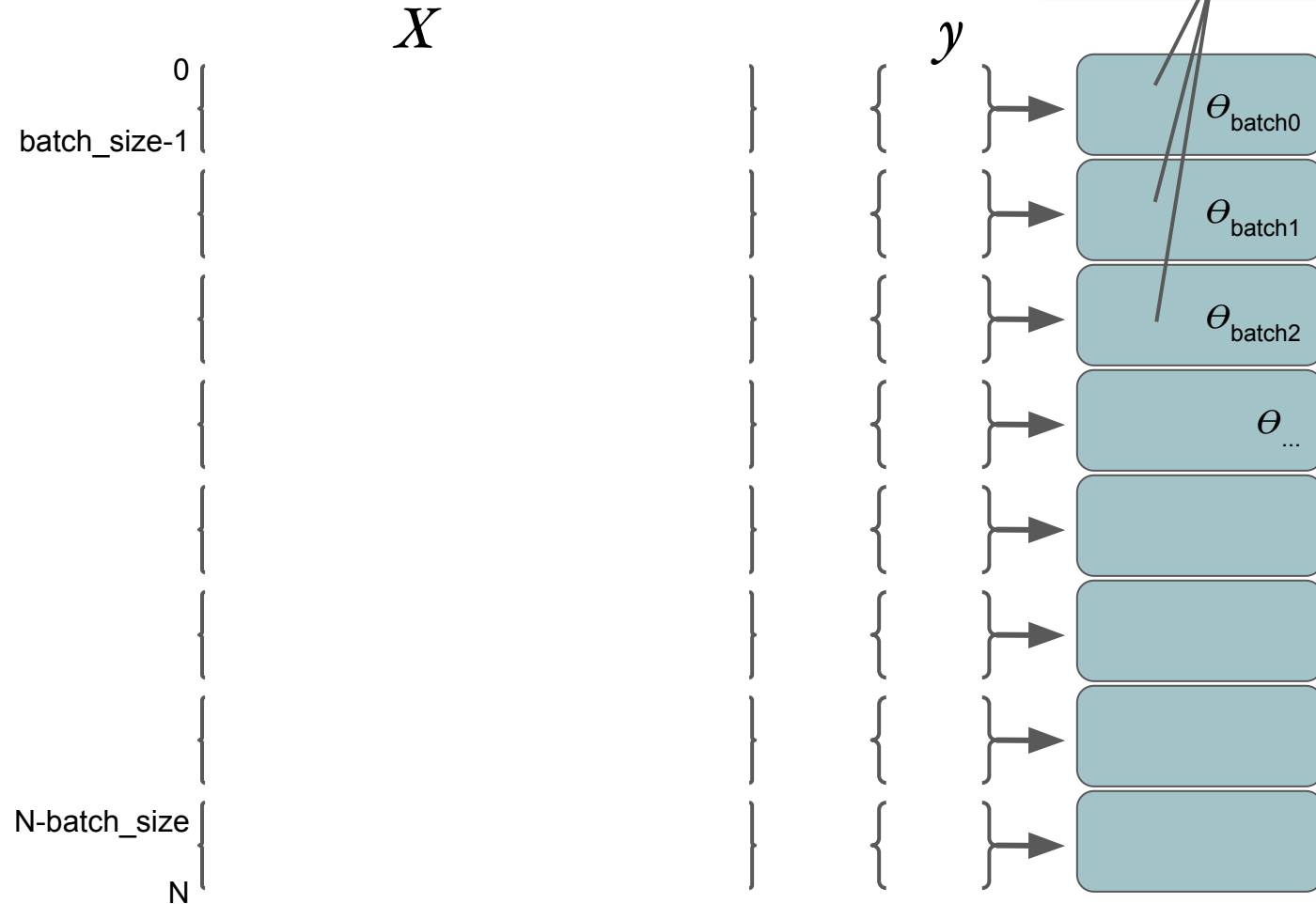
# Distributing Data



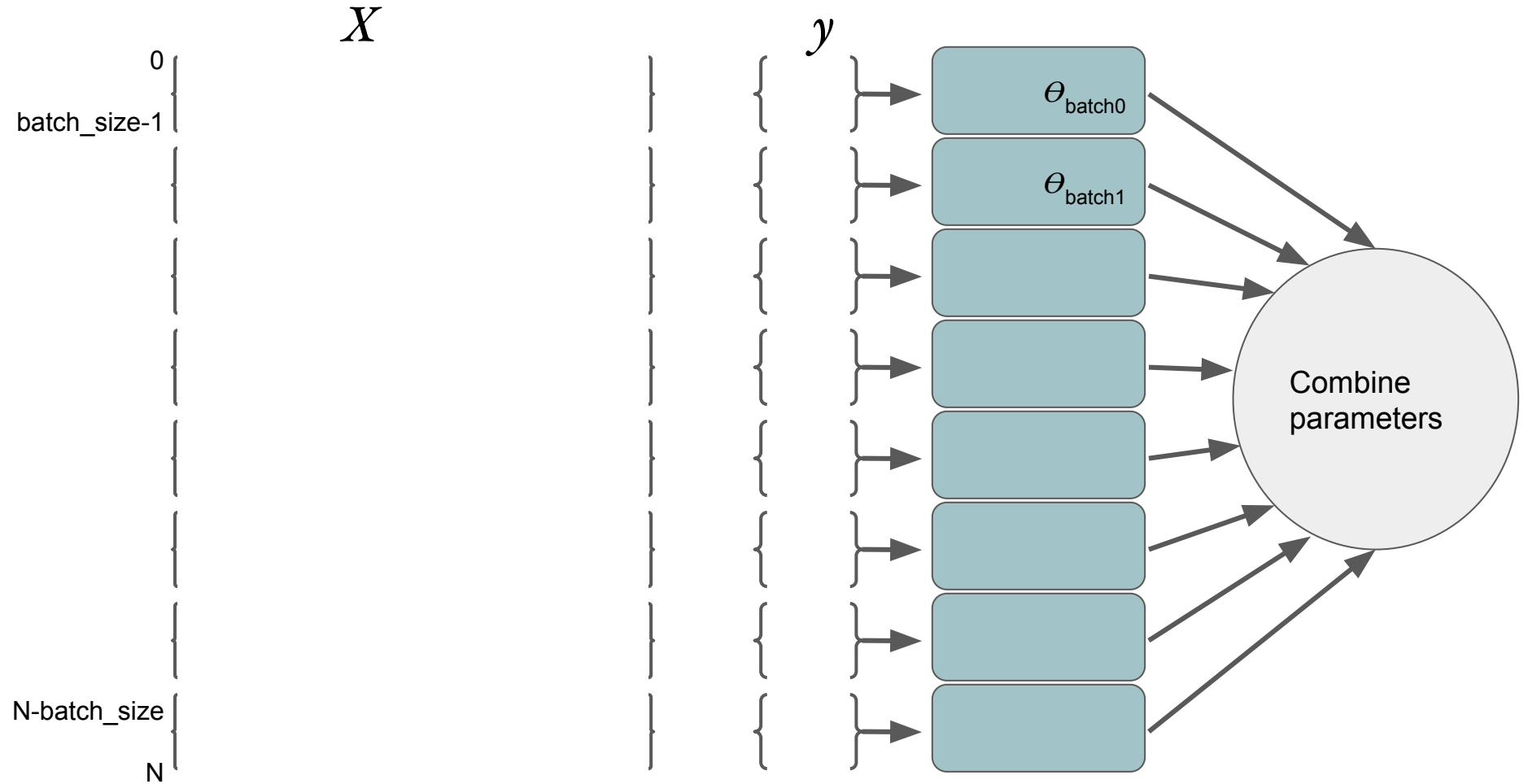


# Distributing Data

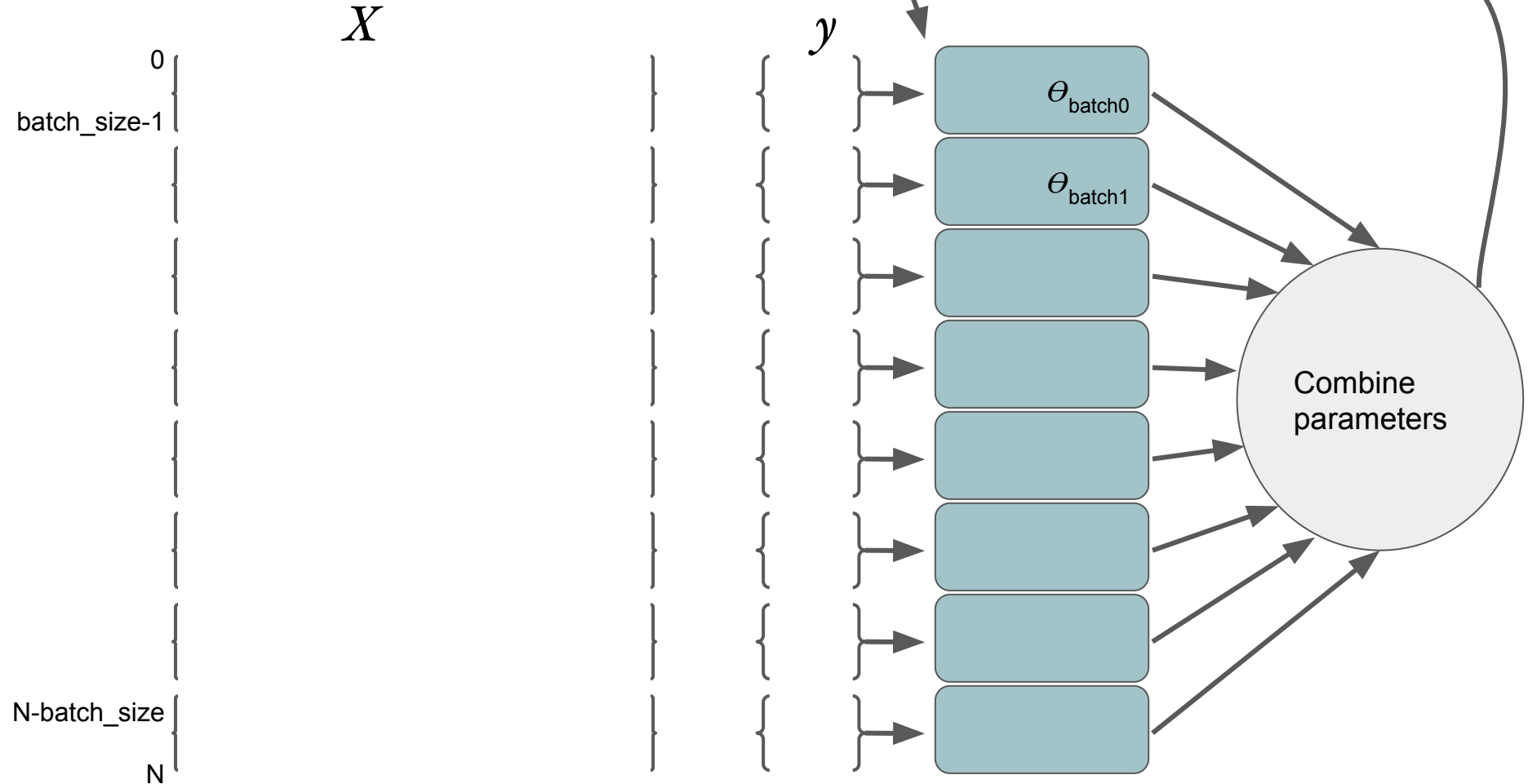
learn parameters (i.e. weights),  
given graph with cost function  
and *optimizer*



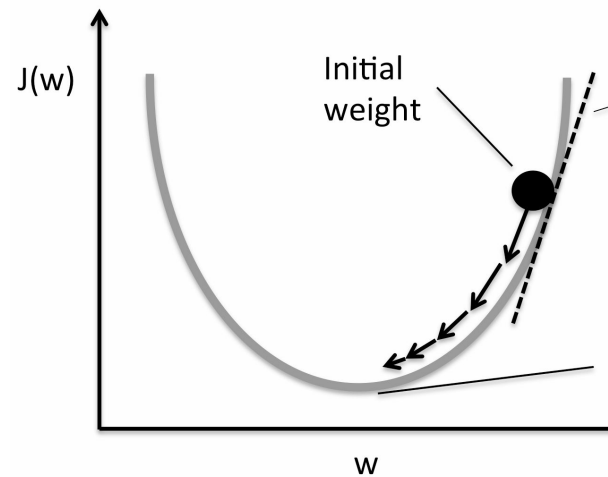
# Distributing Data



# Distributing Data



# Gradient Descent for Linear Regression



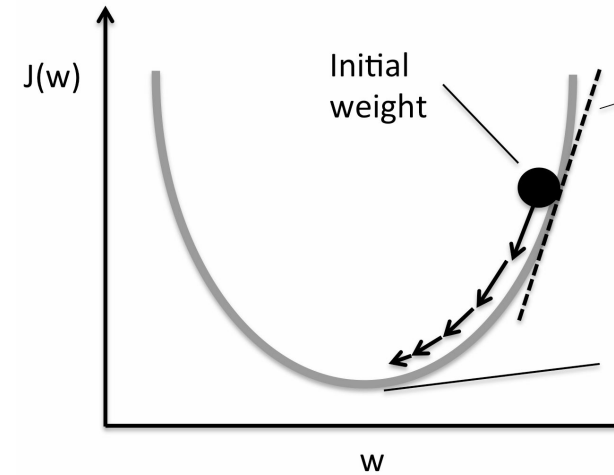
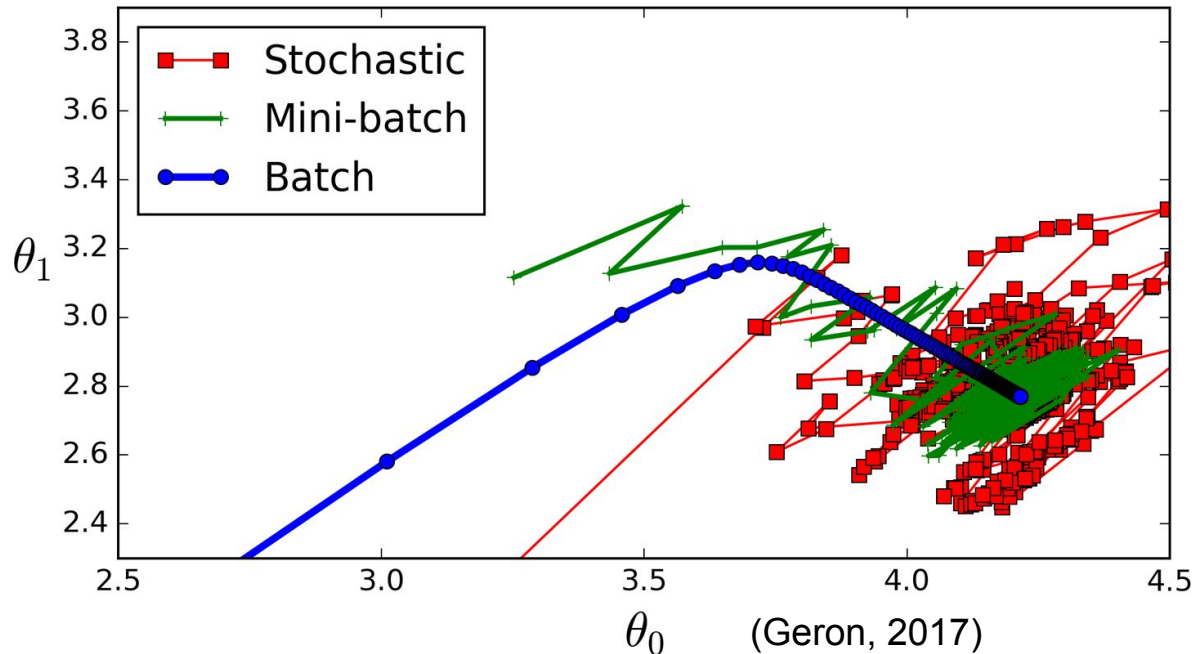
(Geron, 2017)

# Gradient Descent for Linear Regression

Batch Gradient Descent

Stochastic Gradient Descent: One example at a time

Mini-batch Gradient Descent:  $k$  examples at a time.

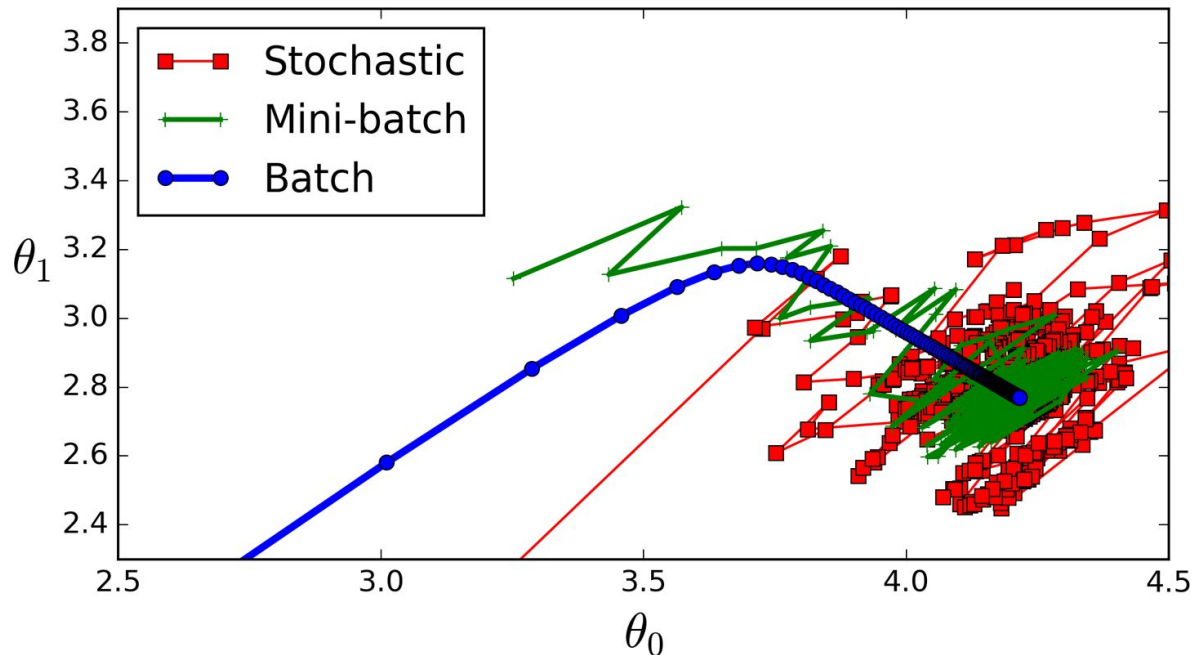


# Gradient Descent for Linear Regression

Batch Gradient Descent

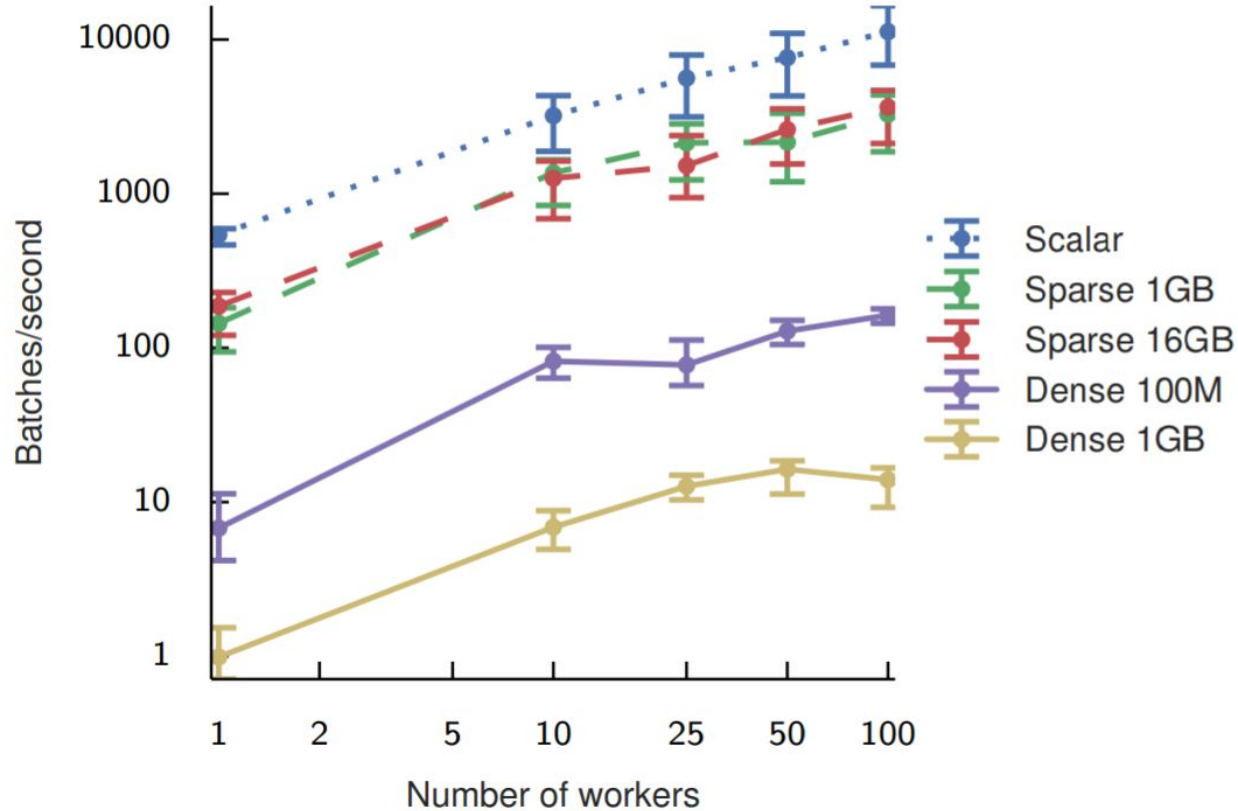
Stochastic Gradient Descent: One example at a time

Mini-batch Gradient Descent:  $k$  examples at a time.



(Geron, 2017)

# Distributed TensorFlow



# Distributed TensorFlow

## Distributed:

- Locally: Across processors (cpus, gpus, tpus)
- Across a Cluster: Multiple machine with multiple processors



# Distributed TensorFlow

## Distributed:

- Locally: Across processors (cpus, gpus, tpus)
- Across a Cluster: Multiple machine with multiple processors

## Parallelisms:

- Data Parallelism: All nodes doing same thing on different subsets of data
- Graph/Model Parallelism: Different portions of model on different devices

# Distributed TensorFlow

## Distributed:

- Locally: Across processors (cpus, gpus, tpus)
- Across a Cluster: Multiple machine with multiple processors

## Parallelisms:

- Data Parallelism: All nodes doing same thing on different subsets of data
- Graph/Model Parallelism: Different portions of model on different devices

## Model Updates:

- Asynchronous Parameter Server
- Synchronous AllReduce (doesn't work with Model Parallelism)

# Distributed TensorFlow

## Distributed:

- Locally: Across processors (cpus, gpus, tpus)
- Across a Cluster: Multiple machine with multiple processors

## Parallelisms:

discussed  
previously

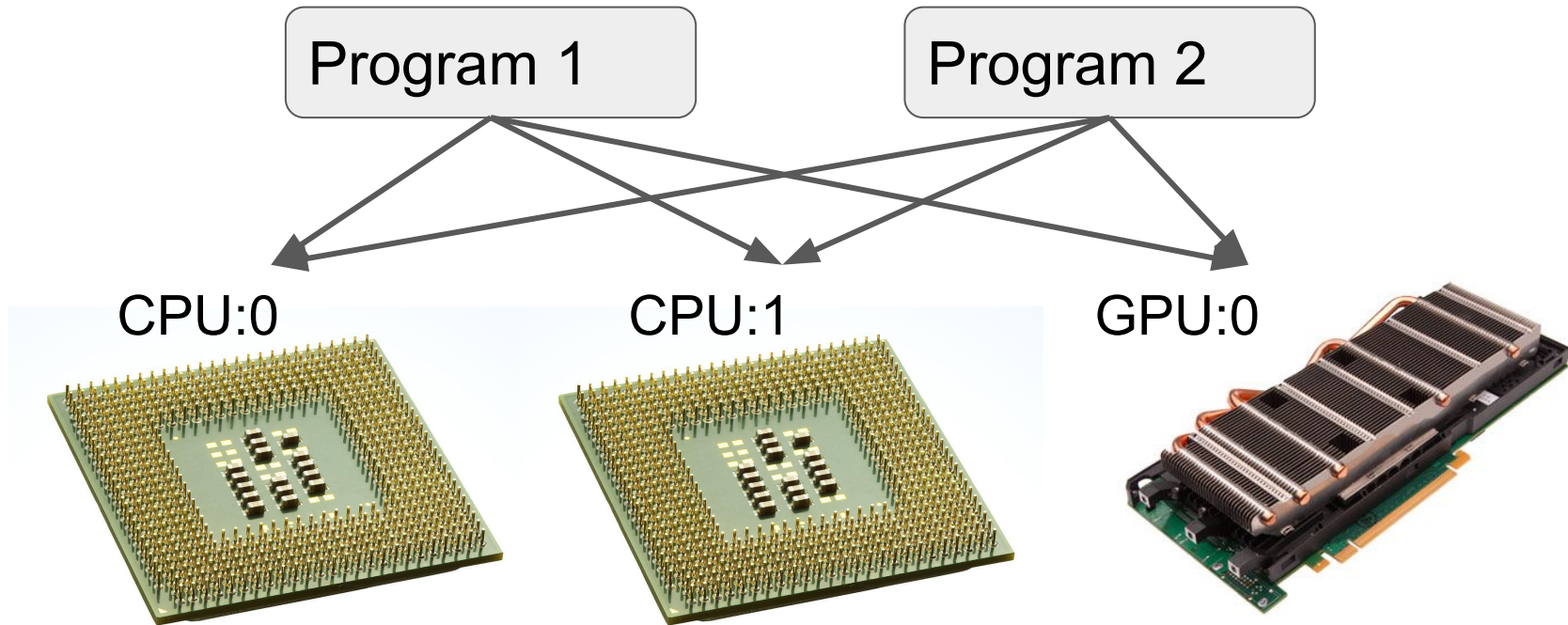
- Data Parallelism: All nodes doing same thing on different subsets of data
- Graph/Model Parallelism: Different portions of model on different devices

## Model Updates:

- Asynchronous Parameter Server
- Synchronous AllReduce (doesn't work with Model Parallelism)

# Local Distribution

Multiple devices on single machine



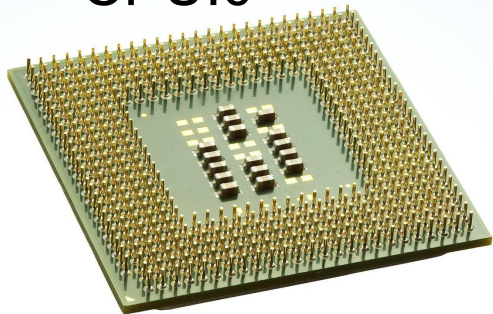
# Local Distribution

Multiple devices on single machine

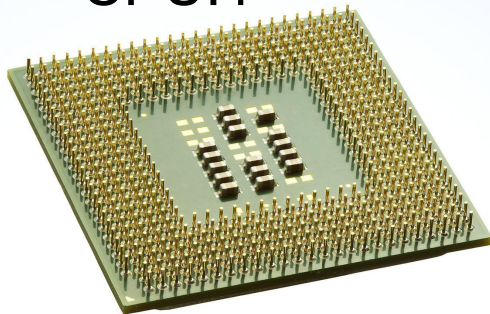
```
with tf.device("/cpu:1")  
    beta=tf.Variable(...)
```

```
with tf.device("/gpu:0")  
    y_pred=tf.matmul(beta,X)
```

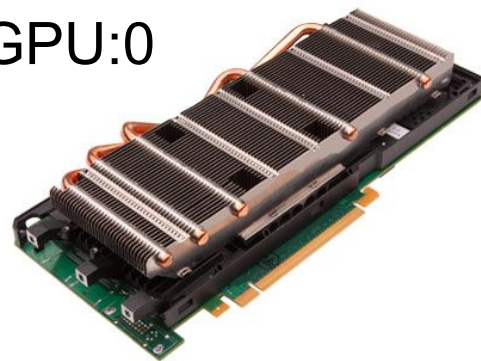
CPU:0



CPU:1



GPU:0

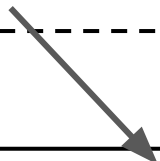


# Local Distribution

Multiple devices on multiple machines

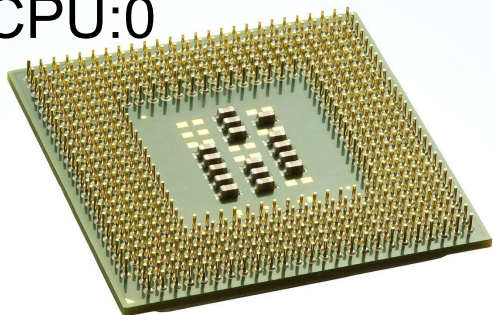
```
with tf.device("/cpu:1")  
    beta=tf.Variable(...)
```

```
with tf.device("/gpu:0")  
    y_pred=tf.matmul(beta,X)
```

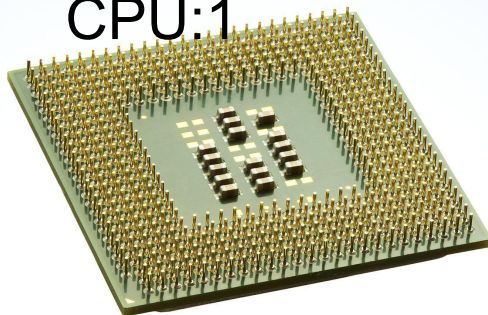


Machine A

CPU:0

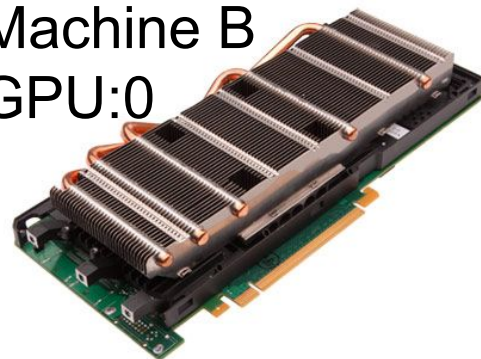


CPU:1



Machine B

GPU:0



# Distributed TensorFlow

## Distributed:

- Locally: Across processors (cpus, gpus, tpus)
- Across a Cluster: Multiple machine with multiple processors

## Parallelisms:

- Data Parallelism: All nodes doing same thing on different subsets of data
- Graph/Model Parallelism: Different portions of model on different devices

## Model Updates:

- Asynchronous Parameter Server
- Synchronous AllReduce (doesn't work with Model Parallelism)

# Distributed TensorFlow

## Distributed:

- Locally: Across processors (cpus, gpus, tpus)
- Across a Cluster: Multiple machine with multiple processors

## Parallelisms:

- Data Parallelism: All nodes doing same thing on different subsets of data
- Graph/Model Parallelism: Different portions of model on different devices

## Model Updates:

- Asynchronous Parameter Server
- Synchronous AllReduce (doesn't work with Model Parallelism)



# Parallelisms

Model Parallelism

Multiple devices on multiple machines

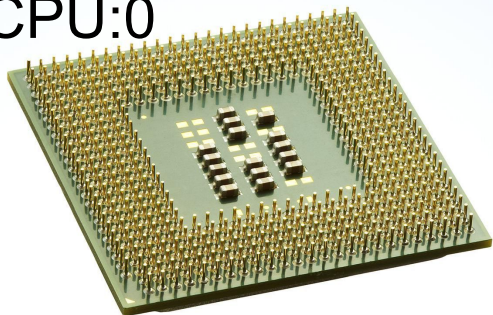
```
with tf.device("/cpu:1")  
    beta=tf.Variable(...)
```

```
with tf.device("/gpu:0")  
    y_pred=tf.matmul(beta,X)
```

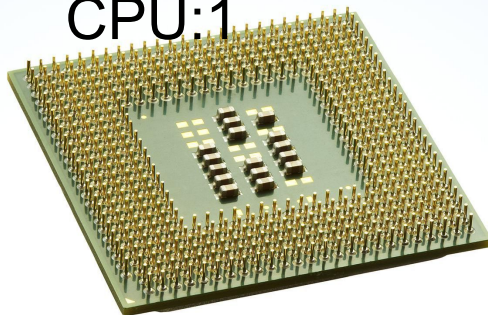
Transfer Tensors

Machine A

CPU:0

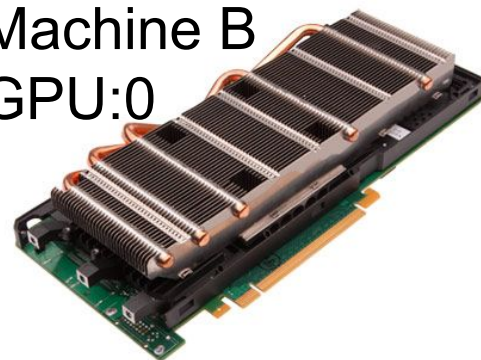


CPU:1



Machine B

GPU:0



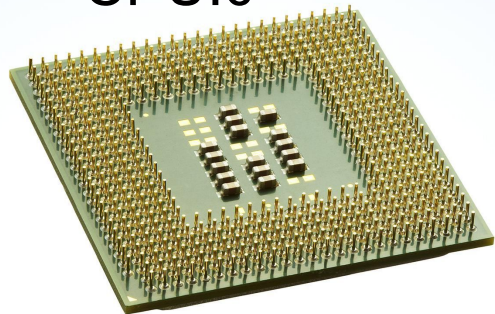
# Parallelisms

## Data Parallelism

```
...  
beta=tf.Variable(...)  
pred=tf.matmul(beta,X)
```



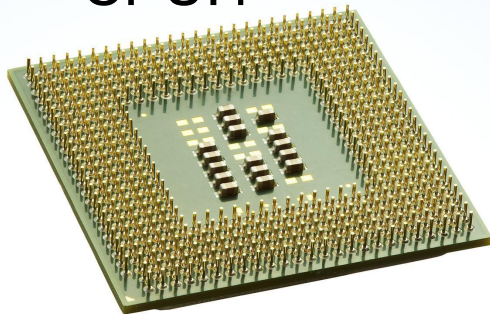
CPU:0



```
...  
beta=tf.Variable(...)  
pred=tf.matmul(beta,X)
```

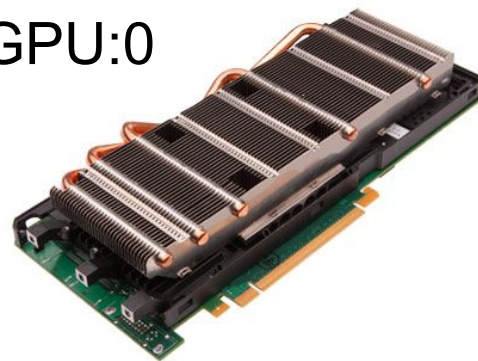


CPU:1



```
...  
beta=tf.Variable(...)  
pred=tf.matmul(beta,X)
```

GPU:0



# Distributed TensorFlow

## Distributed:

- Locally: Across processors (cpus, gpus, tpus)
- Across a Cluster: Multiple machine with multiple processors

## Parallelisms:

- Data Parallelism: All nodes doing same thing on different subsets of data
- Graph/Model Parallelism: Different portions of model on different devices

## Model Updates:

- Asynchronous Parameter Server
- Synchronous AllReduce (doesn't work with Model Parallelism)

# Distributed TensorFlow

## Distributed:

- Locally: Across processors (cpus, gpus, tpus)
- Across a Cluster: Multiple machine with multiple processors

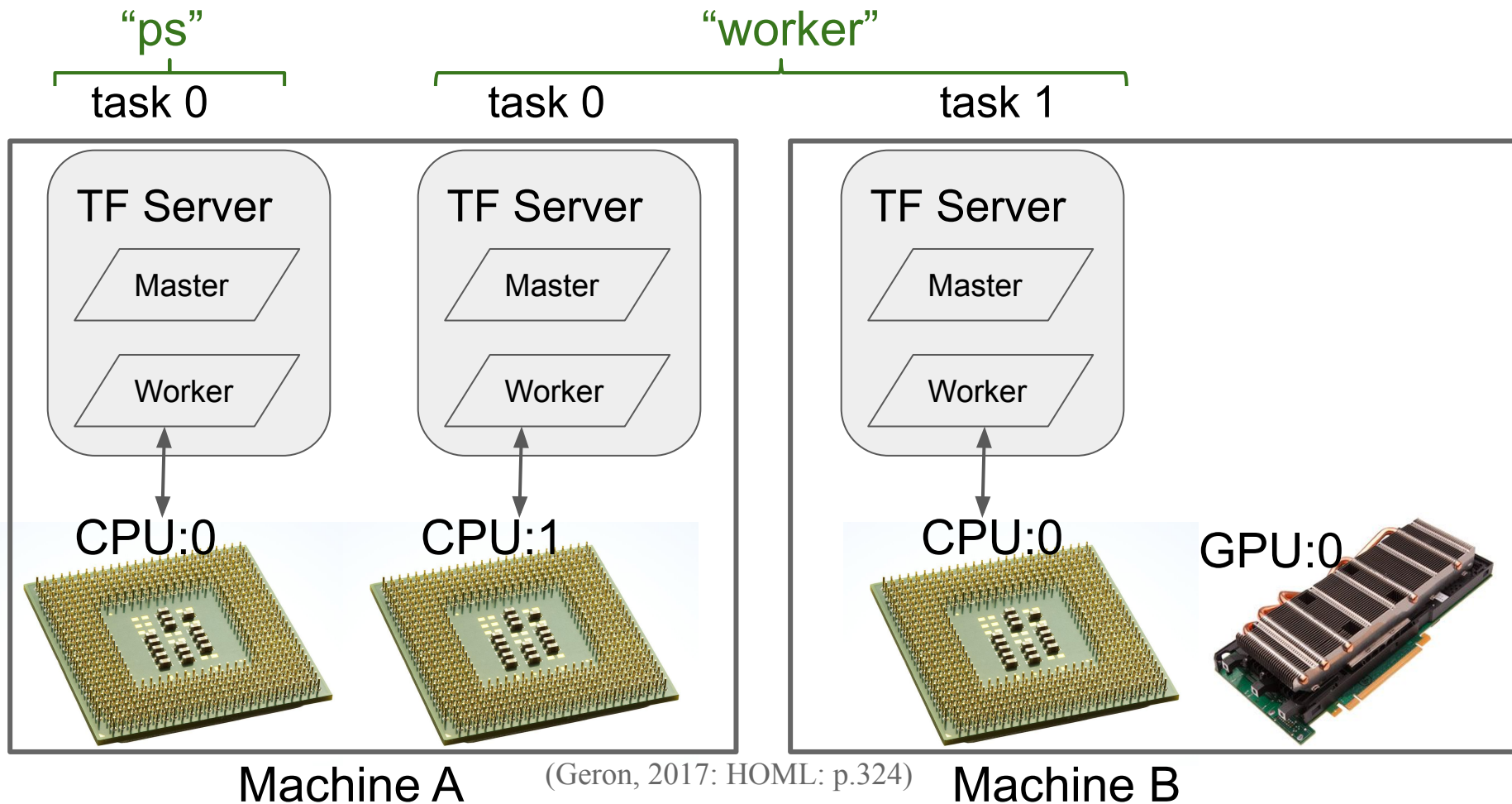
## Parallelisms:

- Data Parallelism: All nodes doing same thing on different subsets of data
- Graph/Model Parallelism: Different portions of model on different devices

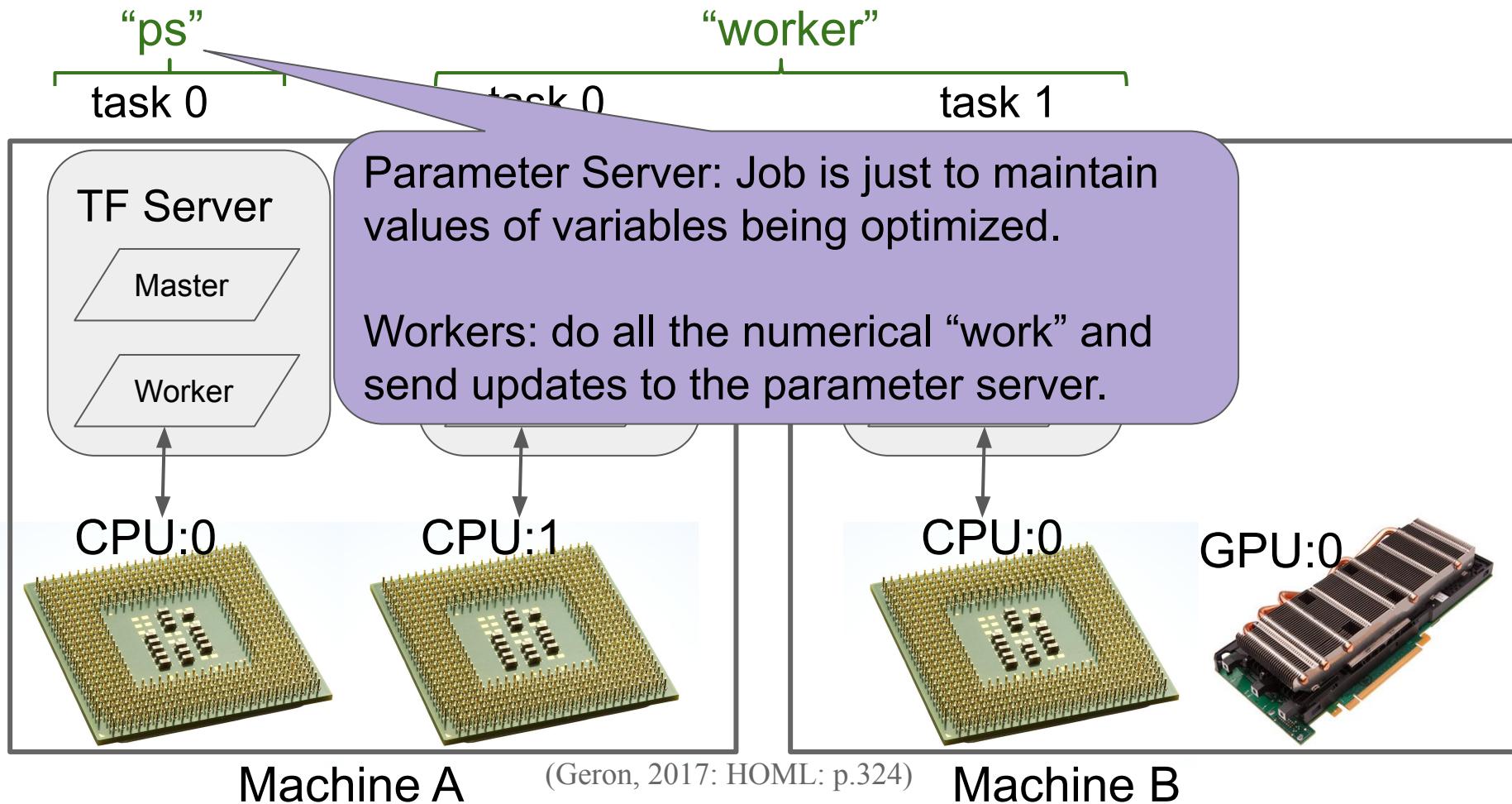
## Model Updates:

- Asynchronous Parameter Server
- Synchronous AllReduce (doesn't work with Model Parallelism)

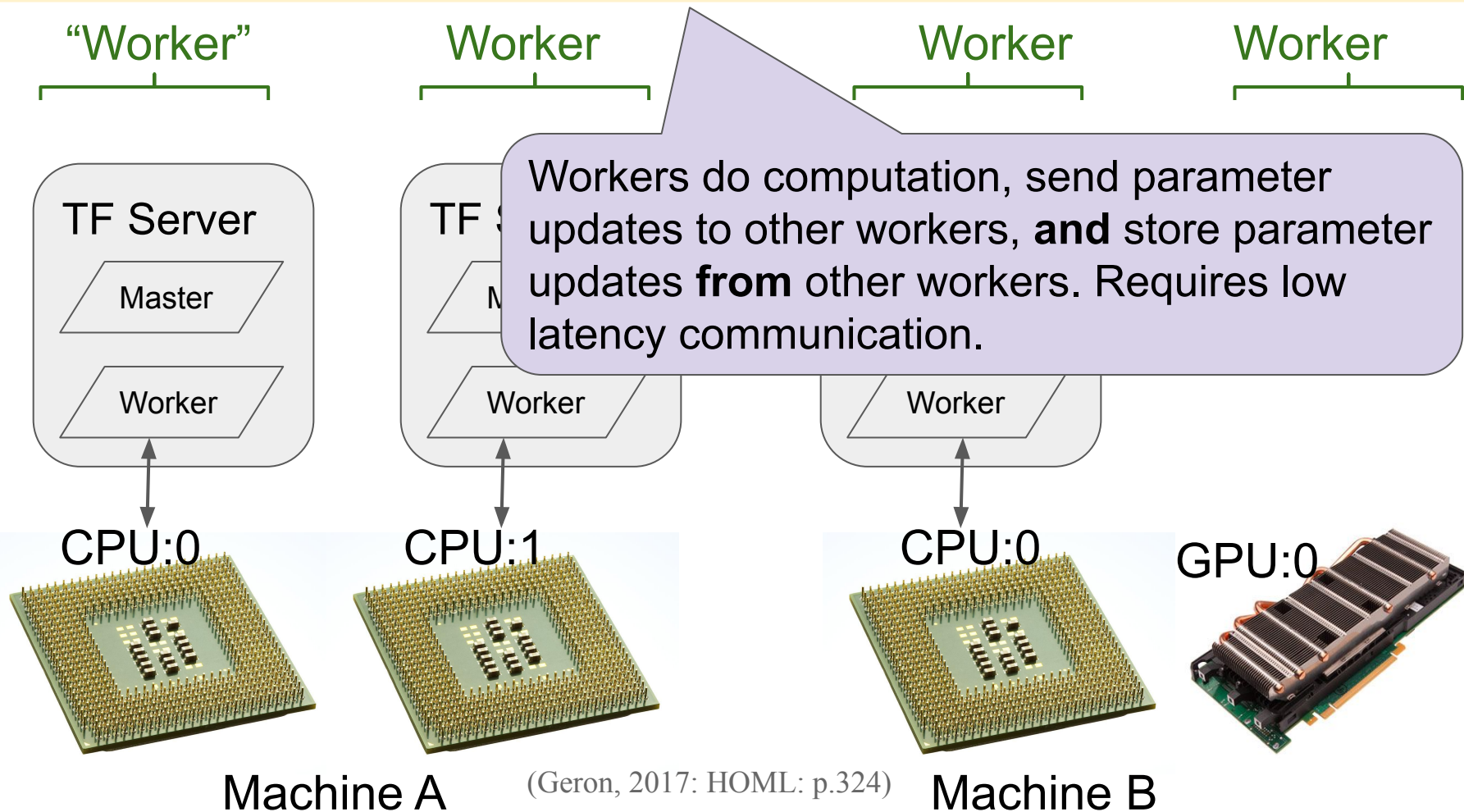
# Asynchronous Parameter Server



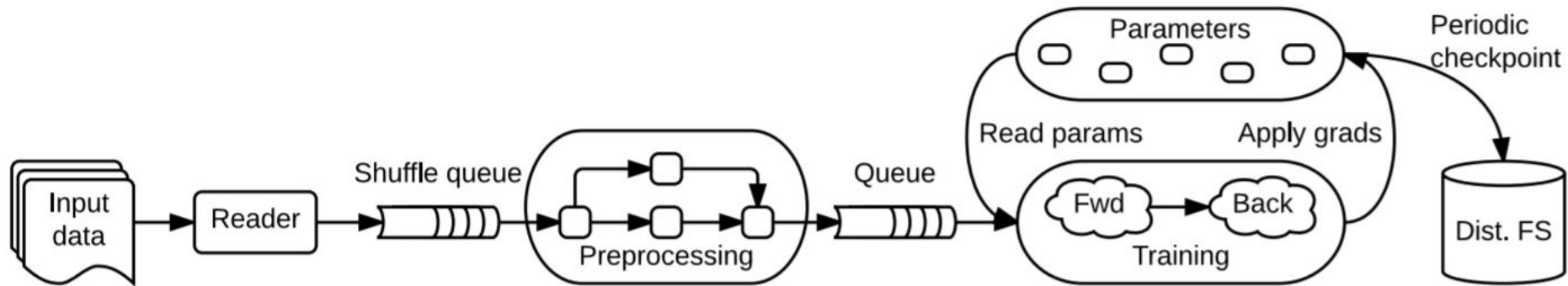
# Asynchronous Parameter Server



# Synchronous All Reduce



# Distributed TF: Full Pipeline





# Summary

- TF is a workflow system, where records are always tensors
  - *operations* applied to tensors (as either Variables, constants, or placeholder)
- Optimized for numerical / linear algebra
  - automatically finds gradients
  - custom kernels for given devices
- “Easily” distributes
  - Within a single machine (local: many devices))
  - Across a cluster (many machines and devices)
  - Jobs broken up as parameter servers / workers makes coordination of data efficient