


# Large-Scale, Distributed Machine Learning

A decorative graphic featuring two overlapping circles, one blue on the left and one yellow on the right. A grey path with arrows starts from a grey oval on the right and moves towards the left, passing through the intersection of the circles. A dashed grey line also curves across the scene.

**CSE545 - Spring 2022**  
Stony Brook University

H. Andrew Schwartz

# Supervised Learning

(genes)

$X_1$

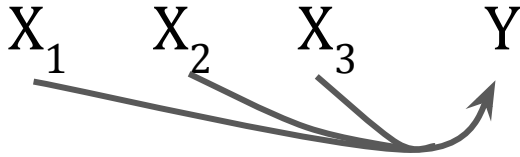
$X_2$

$X_3$

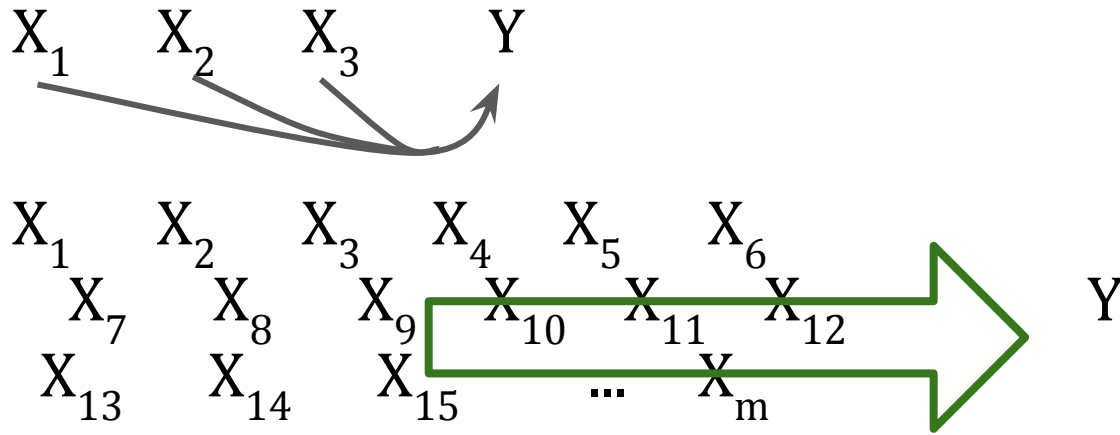
(health)

$Y$

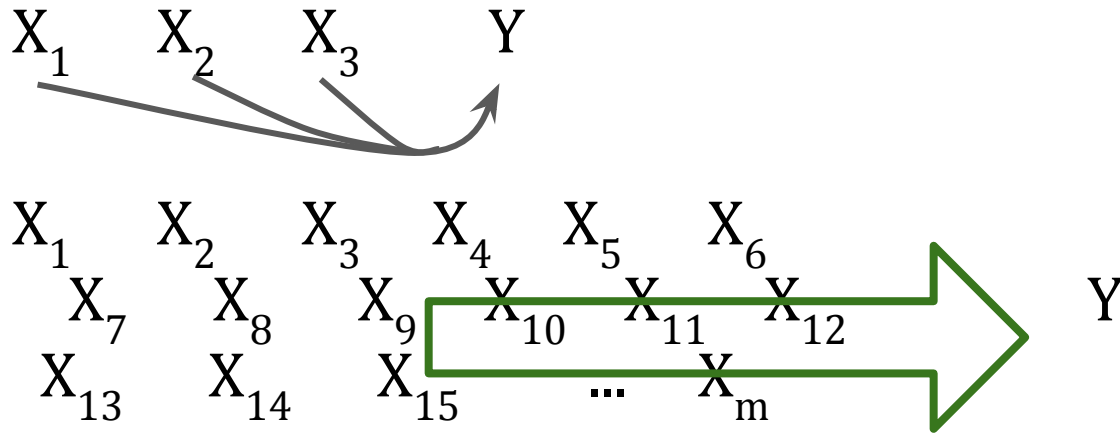
# Supervised Learning



# Supervised Learning



# Supervised Learning



Task: Determine a function,  $f$  (or parameters to a function) such that  $f(X) = Y$

# 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***

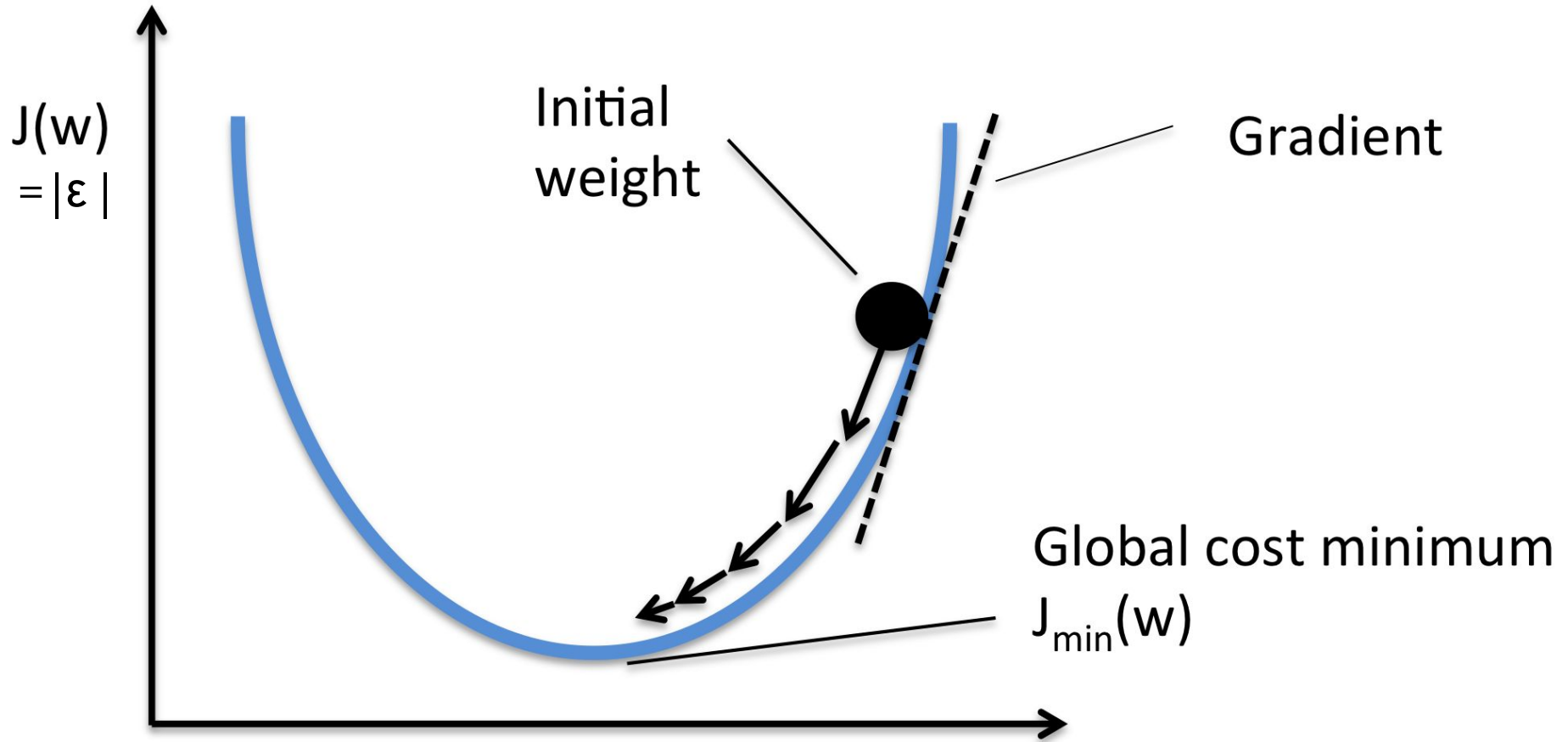
## ***session***

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

## ***devices***

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

# Review: Gradient Descent



$w$

(rasbt, [http://rasbt.github.io/mlxtend/user\\_guide/general\\_concepts/gradient-optimization/](http://rasbt.github.io/mlxtend/user_guide/general_concepts/gradient-optimization/))

# 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:

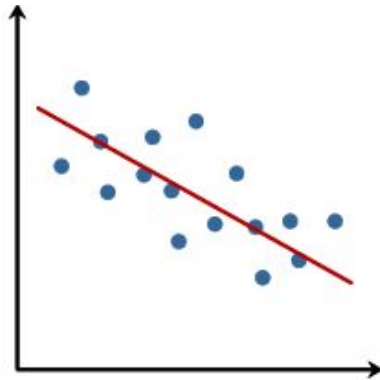
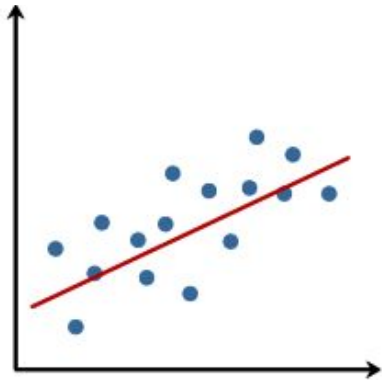
$$\hat{\beta} = \operatorname{argmin}_{\beta} \left\{ \sum_i^N (y_i - \hat{y}_i)^2 \right\}$$

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

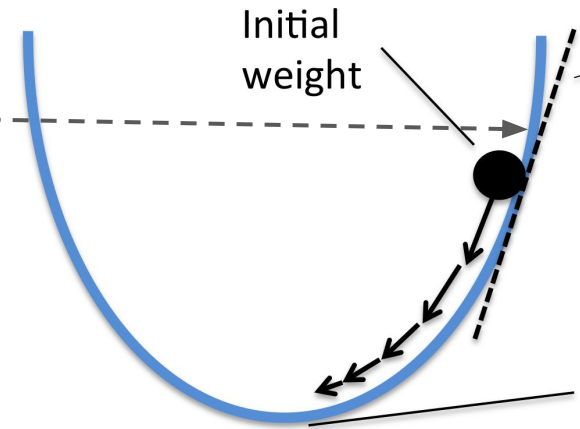
# Weights Derived from Gradients

**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?  $\beta_{\text{new}} = \beta_{\text{prev}} - \alpha * \text{grad}$   
(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,  $\lambda ||\beta||_2^2$ )

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

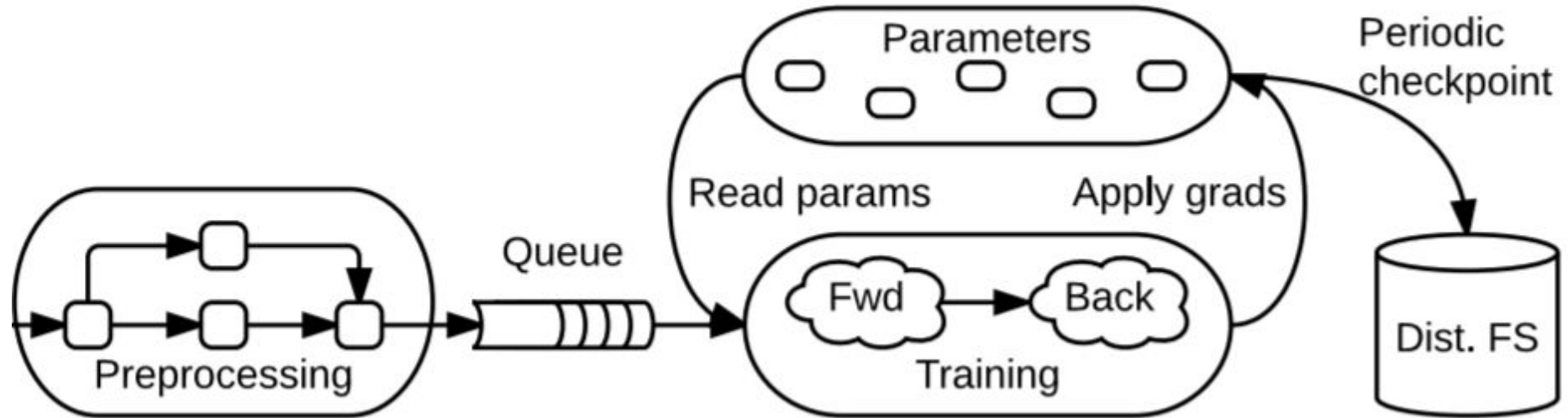
Gradient descent needs to solve.

(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 hyperparameters
- 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 hyperparameters
- 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 hyperparameters
- 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 hyperparameters
- 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 hyperparameters
- 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

- Train over all with different hyperparameters
- Train different folds per worker node.

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

## 2. Distribute data

- Each node finds parameters for subset of data
- Needs mechanism for updating parameters
  - Centralized parameter server
  - 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: Parallelism can be exploited

### Model Parallelism

Con: High communication for transferring Intermediar data.

# Model Parallelism

Multiple devices on multiple machines

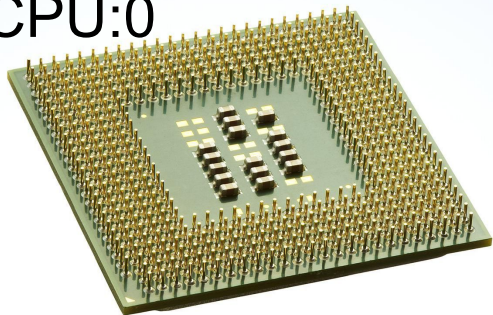
```
with tf.device("/cpu:1")  
    beta=tf.Variable(...)
```

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

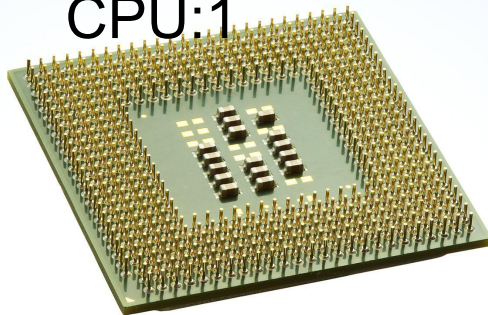
Transfer Tensors

Machine A

CPU:0

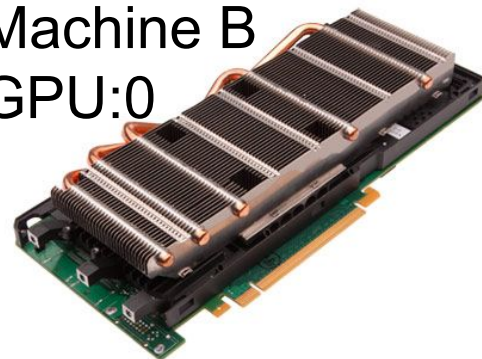


CPU:1



Machine B

GPU:0

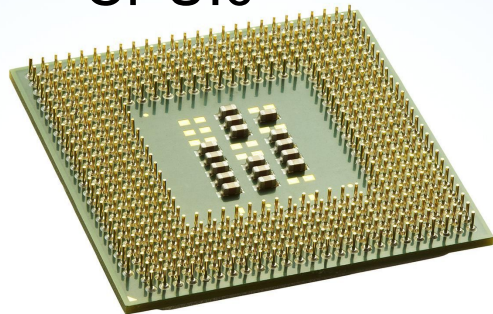


# Data Parallelism

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



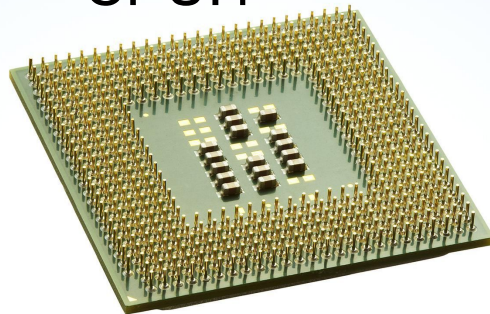
CPU:0



```
...  
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)
```

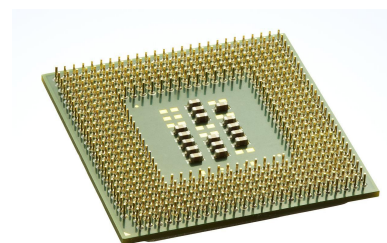
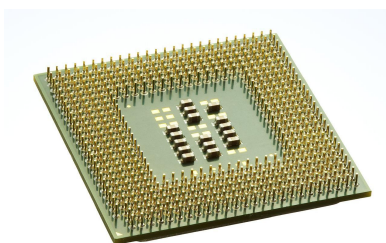
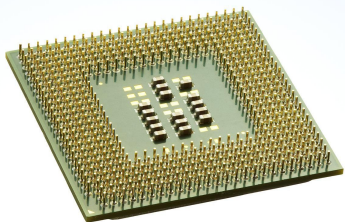
```
...  
beta=tf.Variable(...)  
pred=tf.matmul(beta,X)
```

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

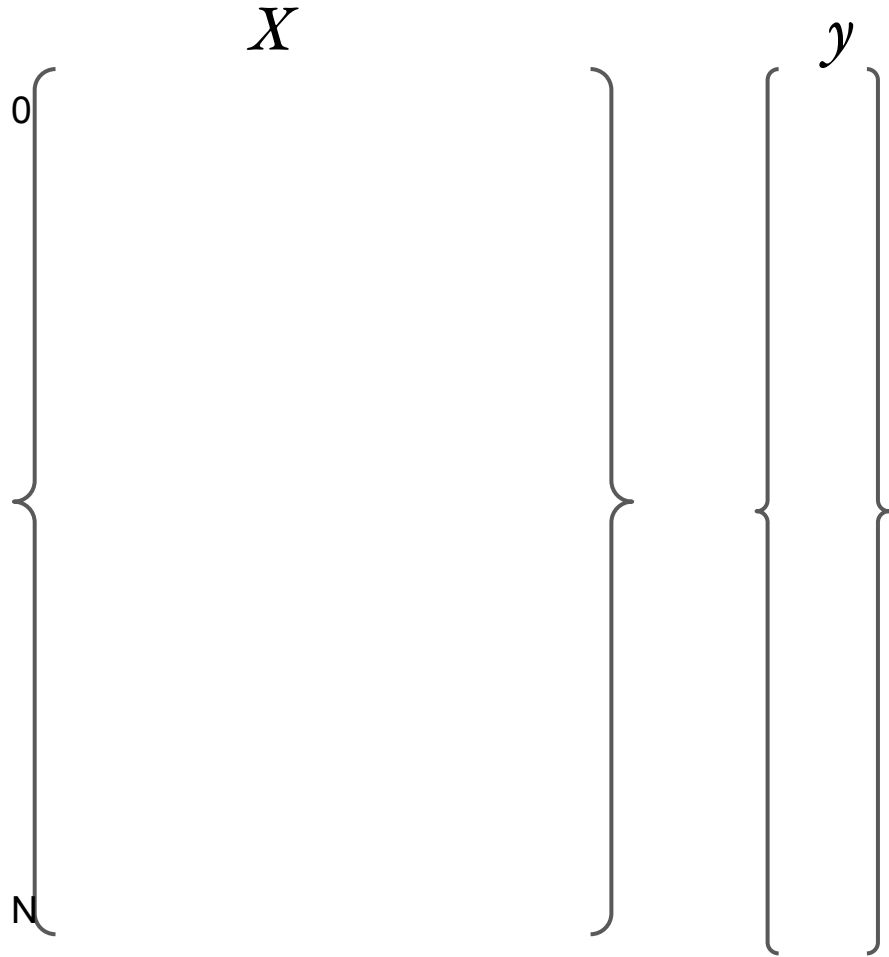
worker:0

worker:1

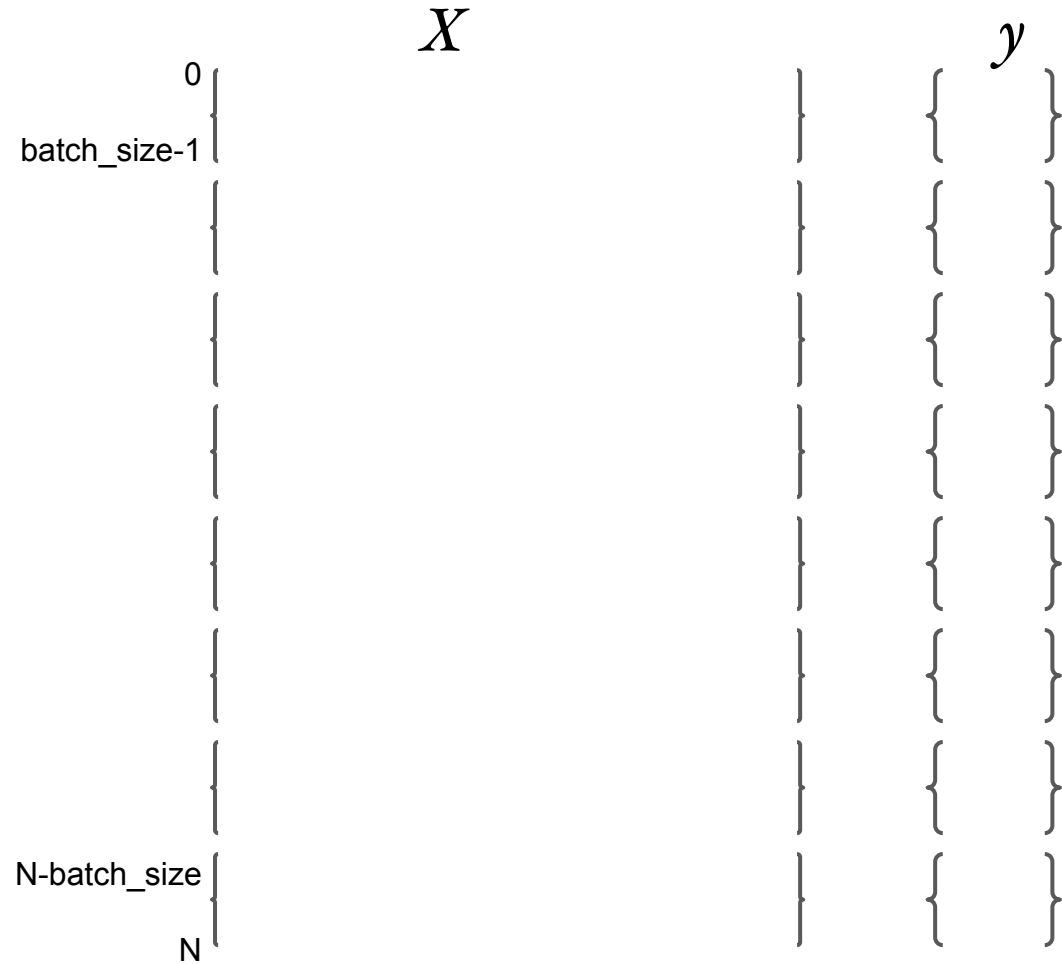
worker:2



# Distributing Data

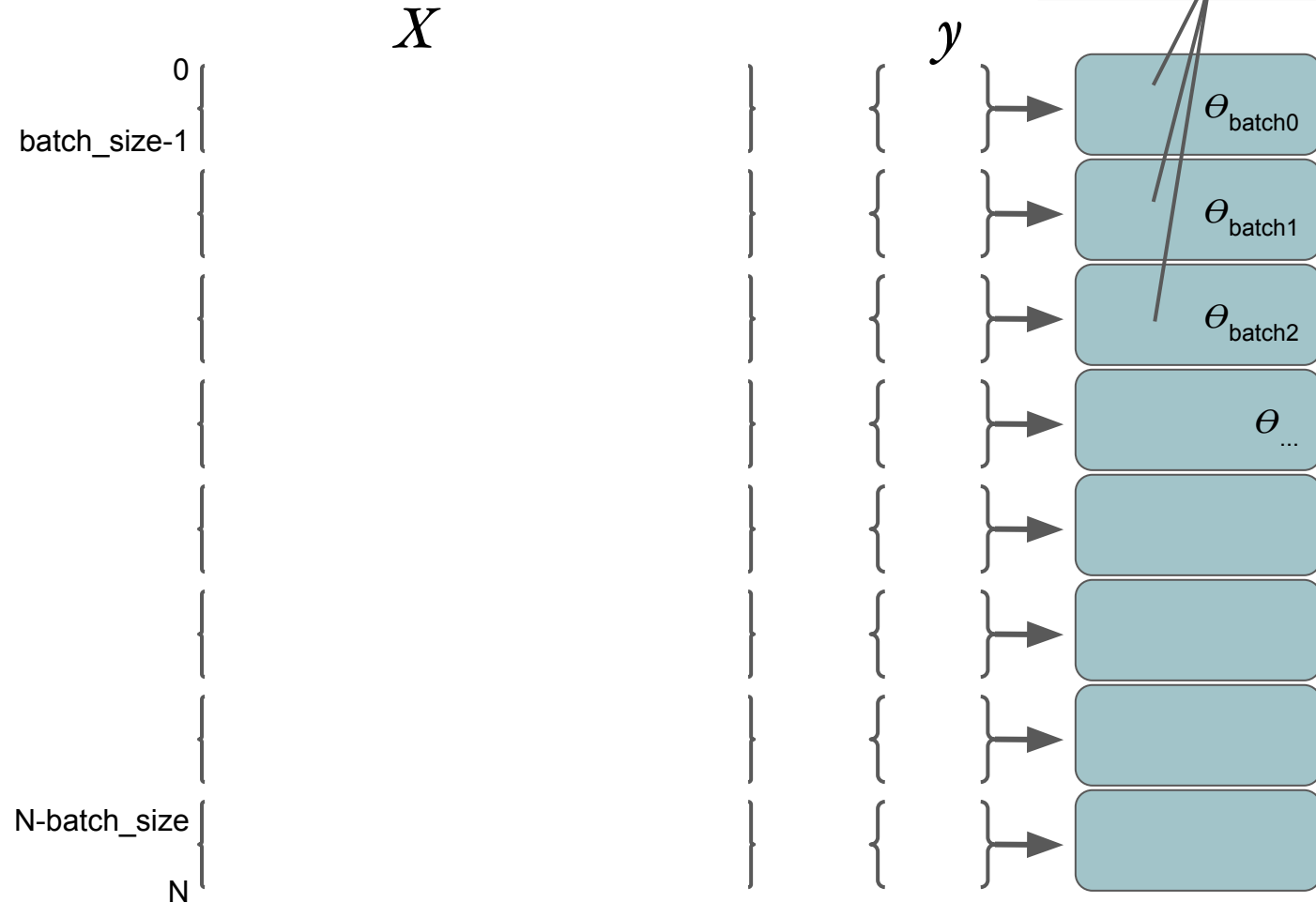


# Distributing Data

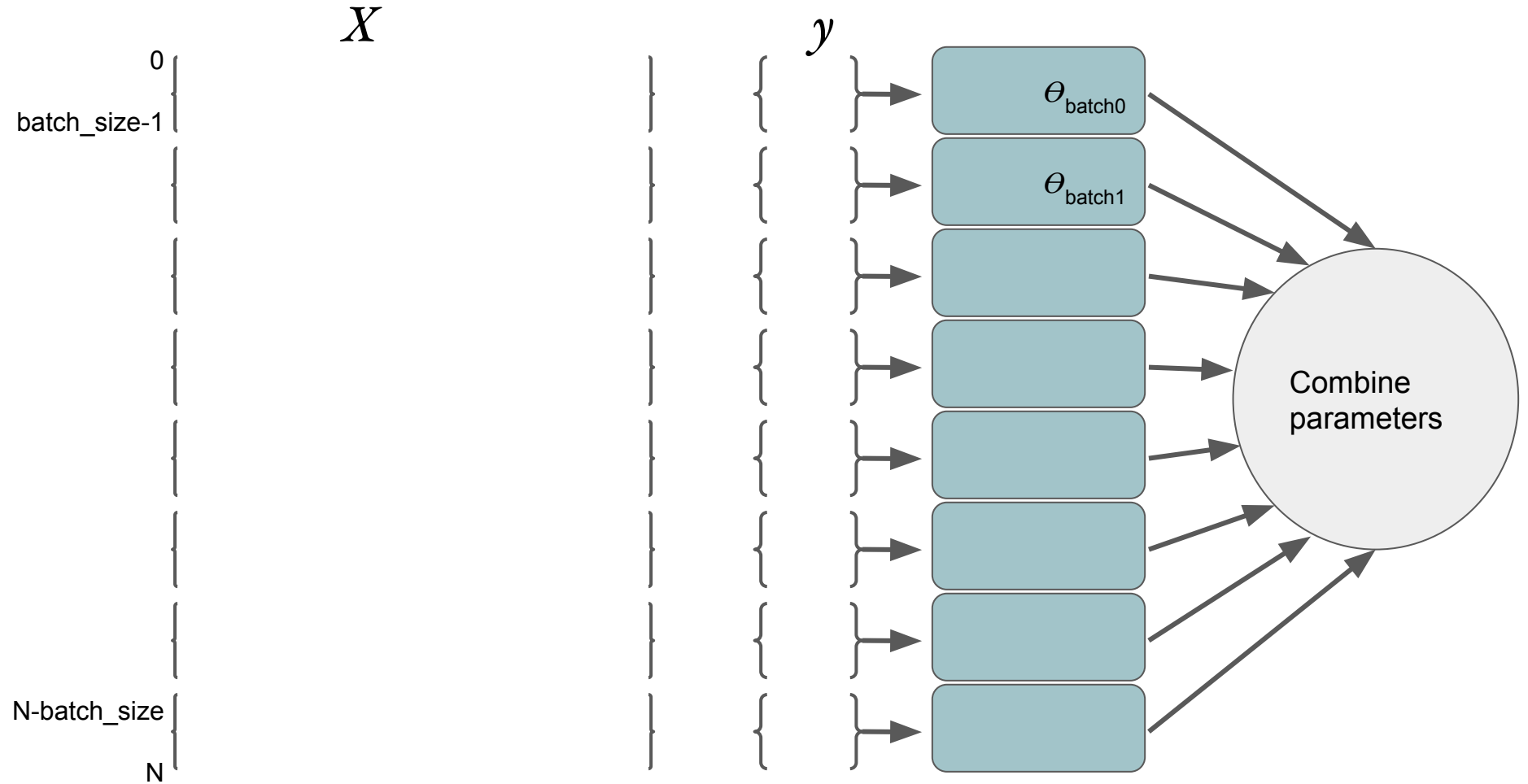


# Distributing Data

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

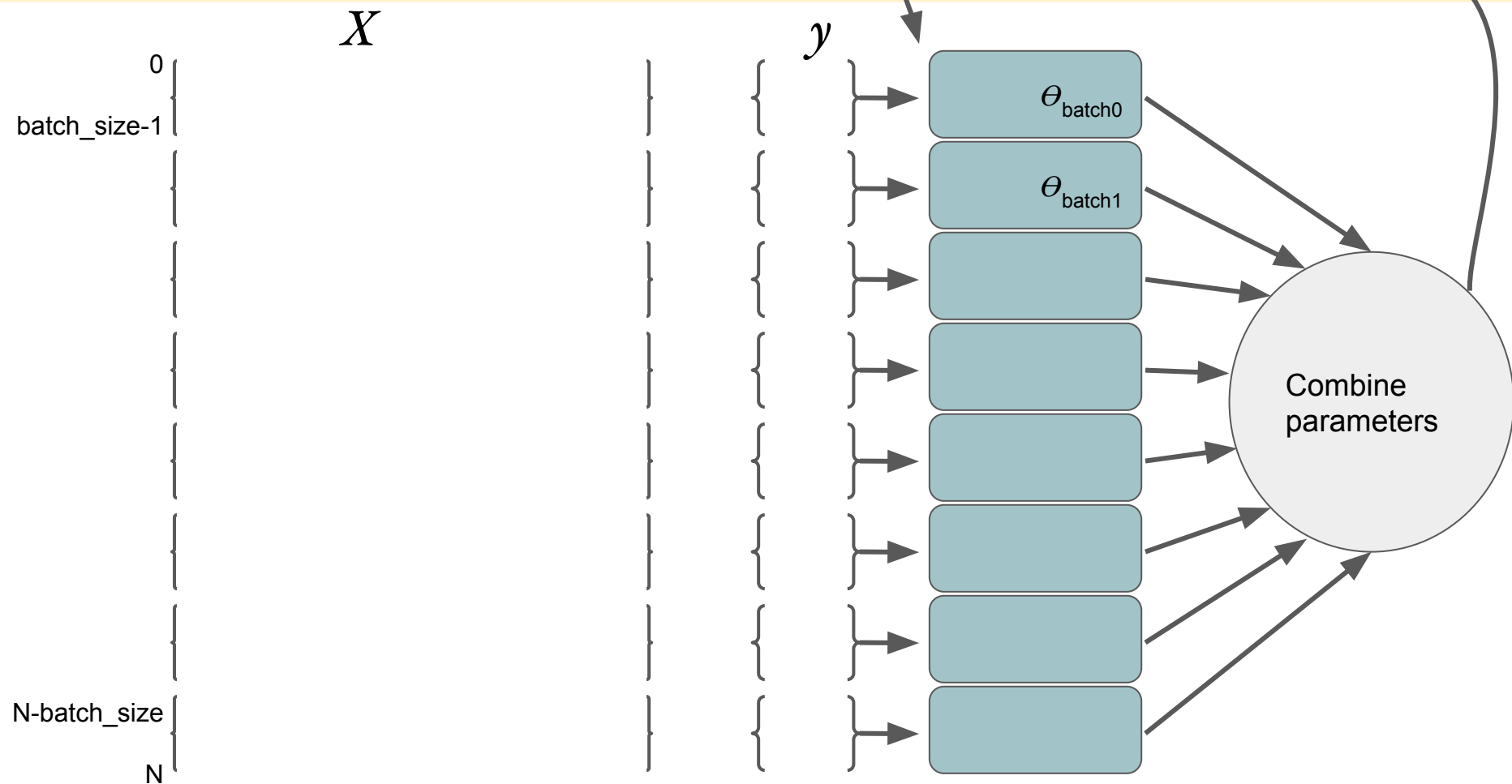


# Distributing Data

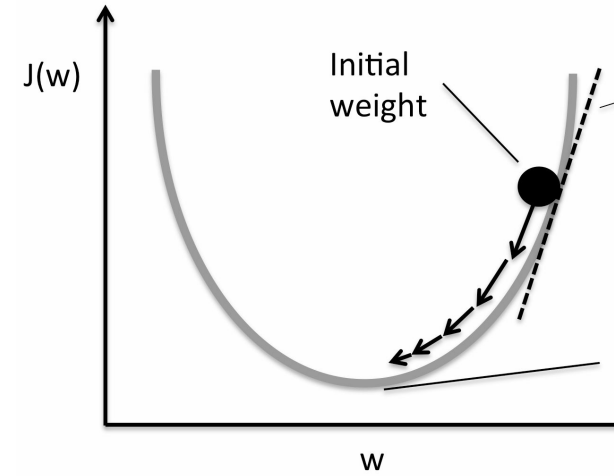




# Distributing Data



# Gradient Descent for Linear Regression



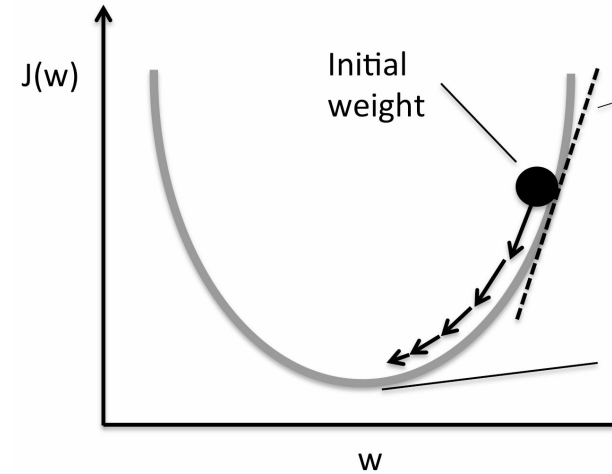
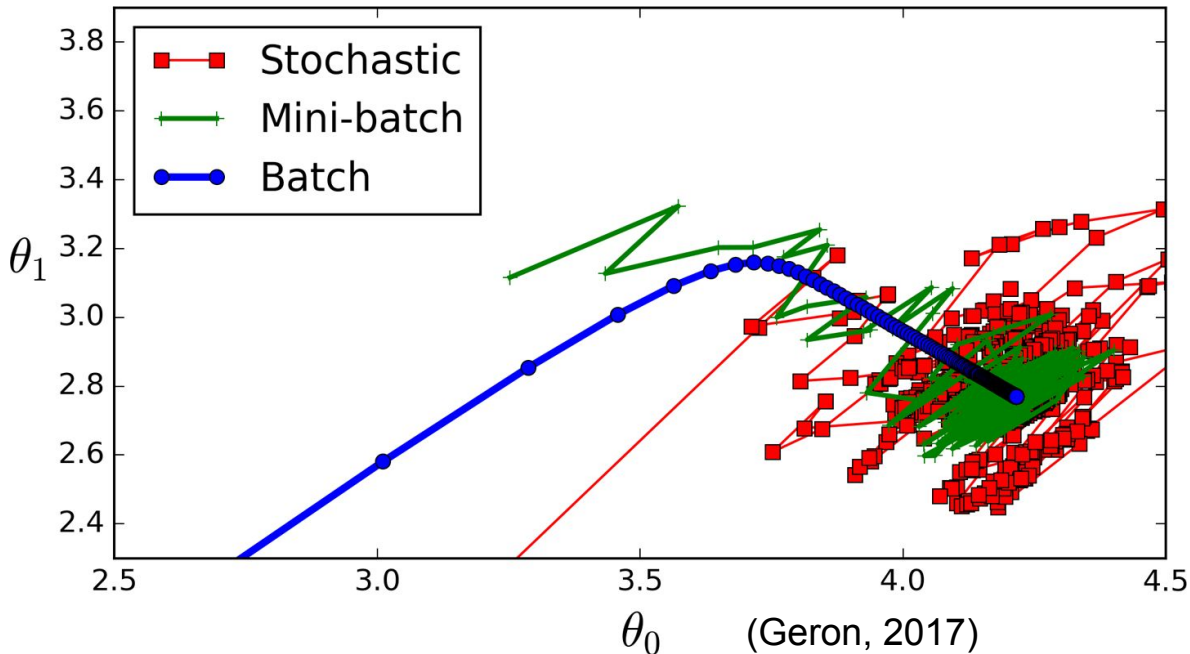
(Geron, 2017)

# 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.

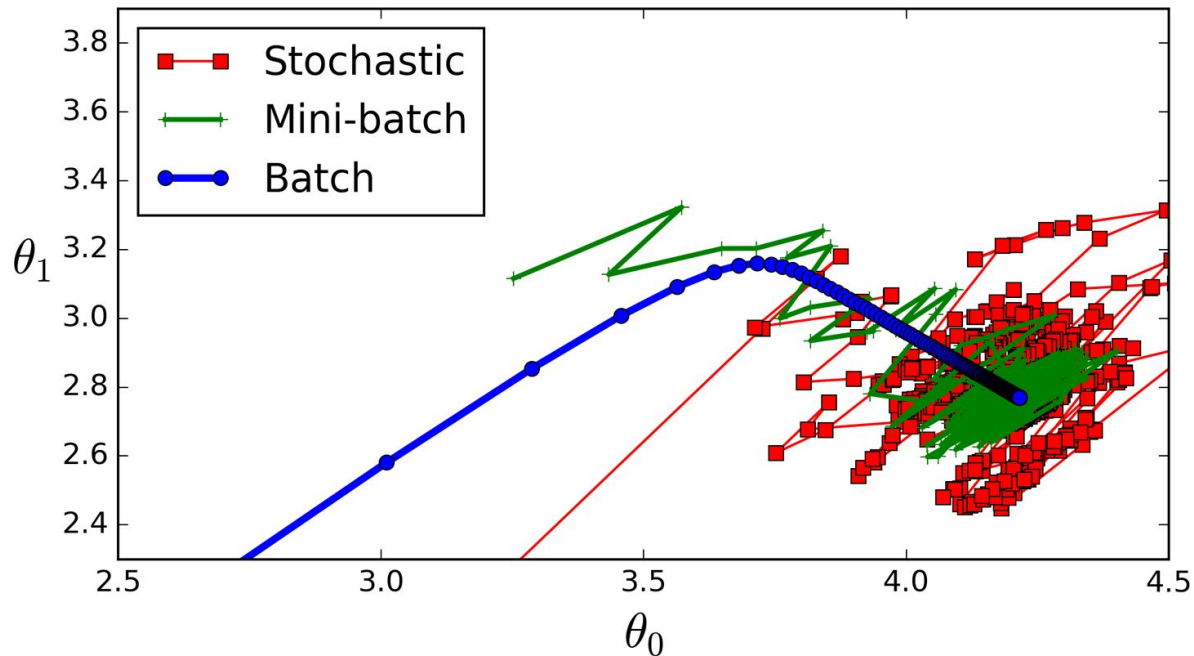


# Gradient Descent for Linear Regression

Batch Gradient Descent

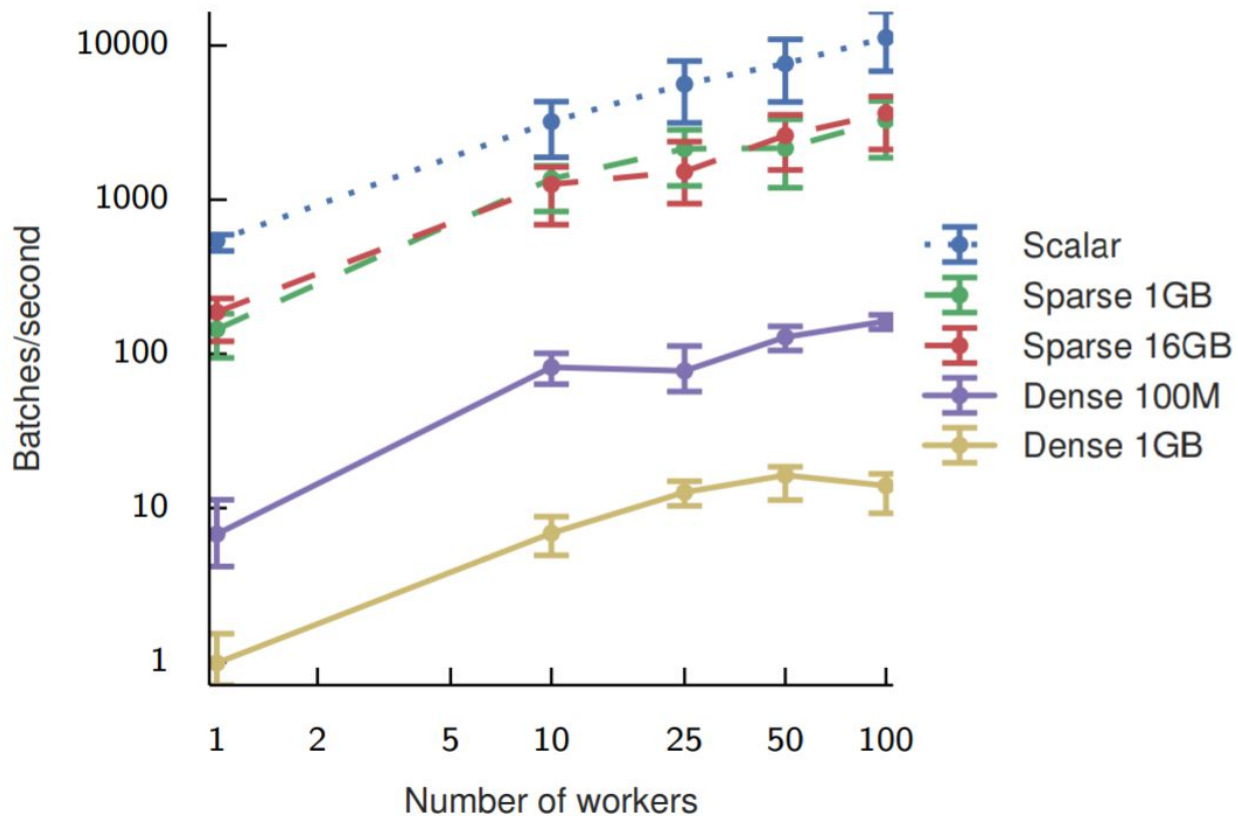
Stochastic Gradient Descent: One example at a time

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

next slides

discussed  
previously

## Parallelisms:

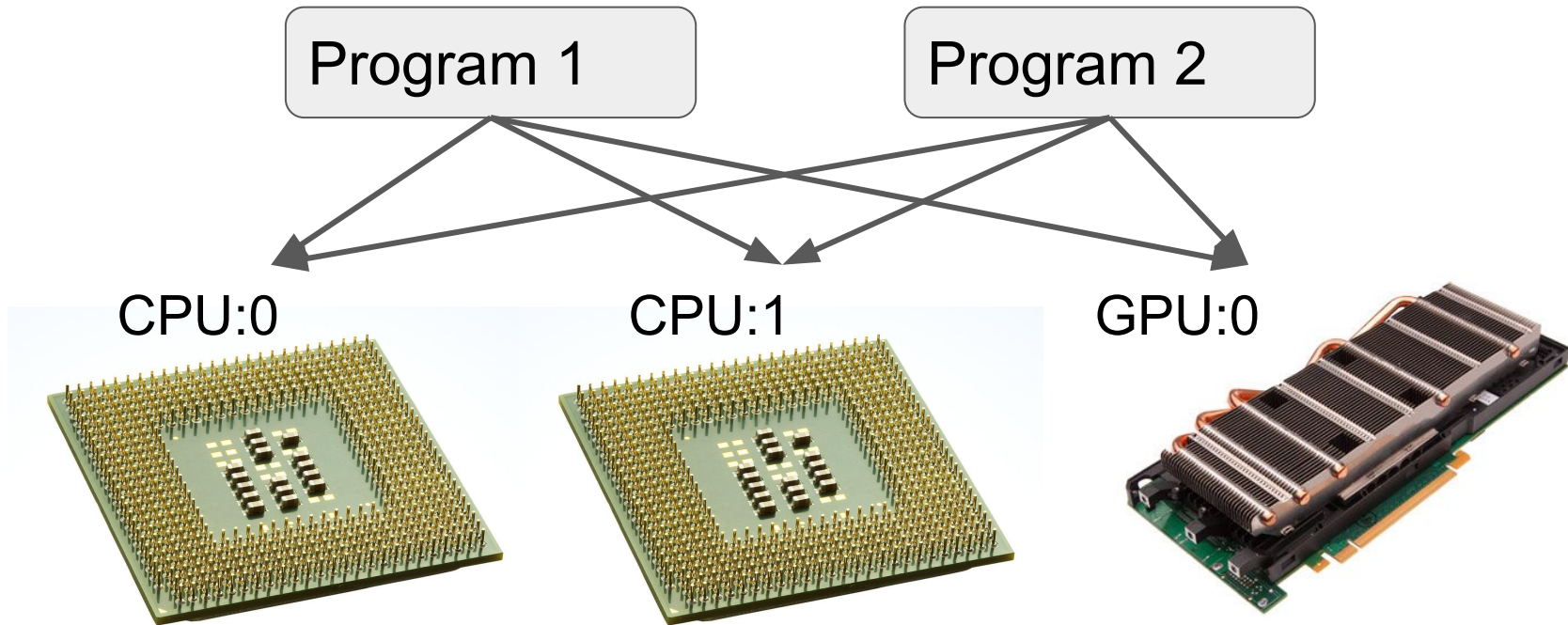
- Data Parallelism: All nodes doing same thing on different subsets of data
- Graph/Model Parallelism: Different portions of model on different devices

## Model Updates:

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

# Local Distribution

Multiple devices on single machine



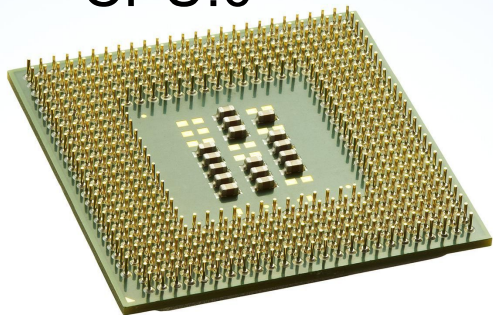
# Local Distribution

Multiple devices on single machine

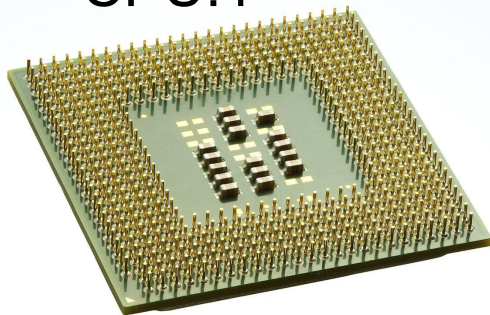
```
with tf.device("/cpu:1")  
    beta=tf.Variable(...)
```

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

CPU:0



CPU:1



GPU:0



# Cluster Distribution

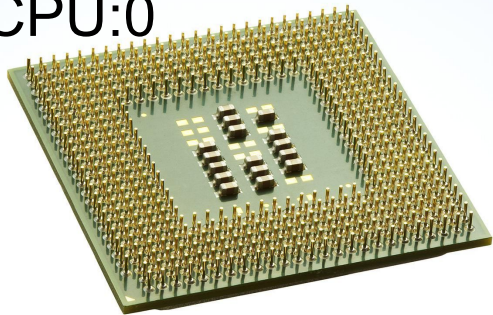
Multiple devices on multiple machines

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

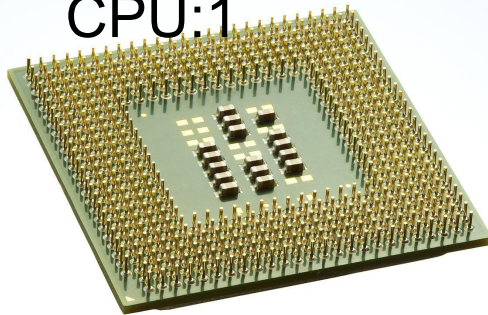
```
with tf.device("/gpu:0")  
    y_pred=tf.matmul(beta,X)
```

Machine A

CPU:0

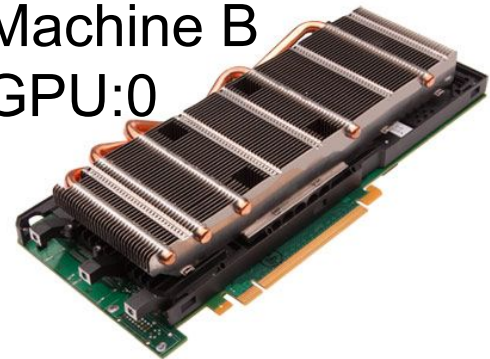


CPU:1



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

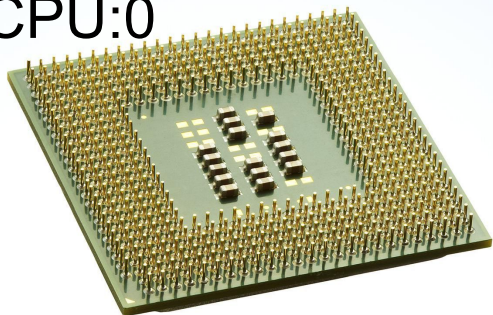
```
with tf.device("/cpu:1")  
    beta=tf.Variable(...)
```

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

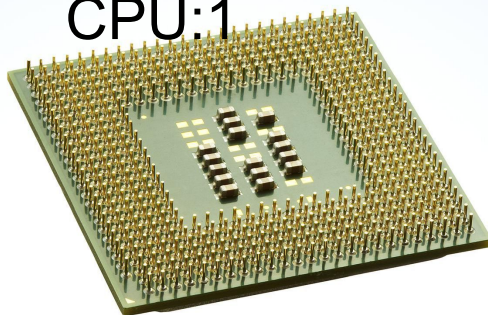
Transfer Tensors

Machine A

CPU:0

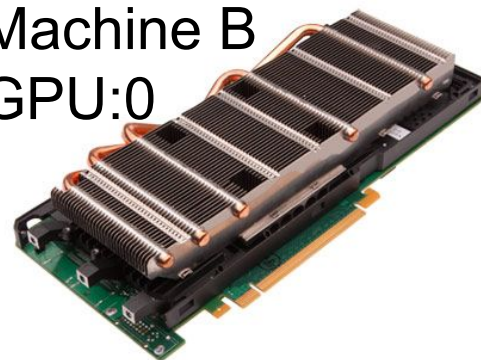


CPU:1



Machine B

GPU:0



# Parallelisms

Data Parallelism

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

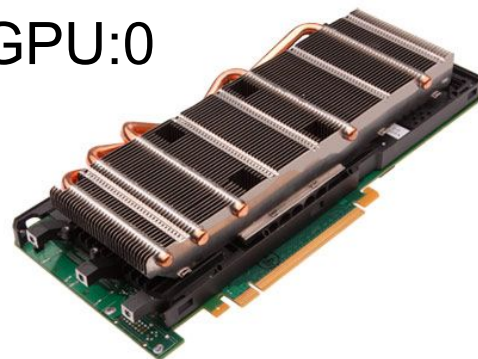
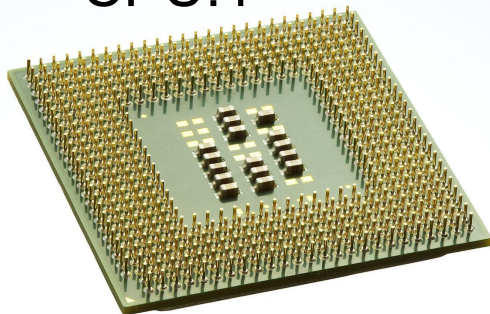
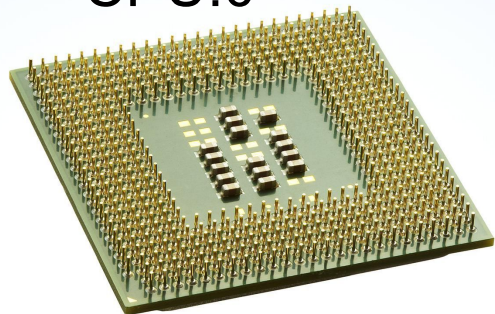
```
...  
beta=tf.Variable(...)  
pred=tf.matmul(beta,X)
```

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

CPU:0

CPU:1

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

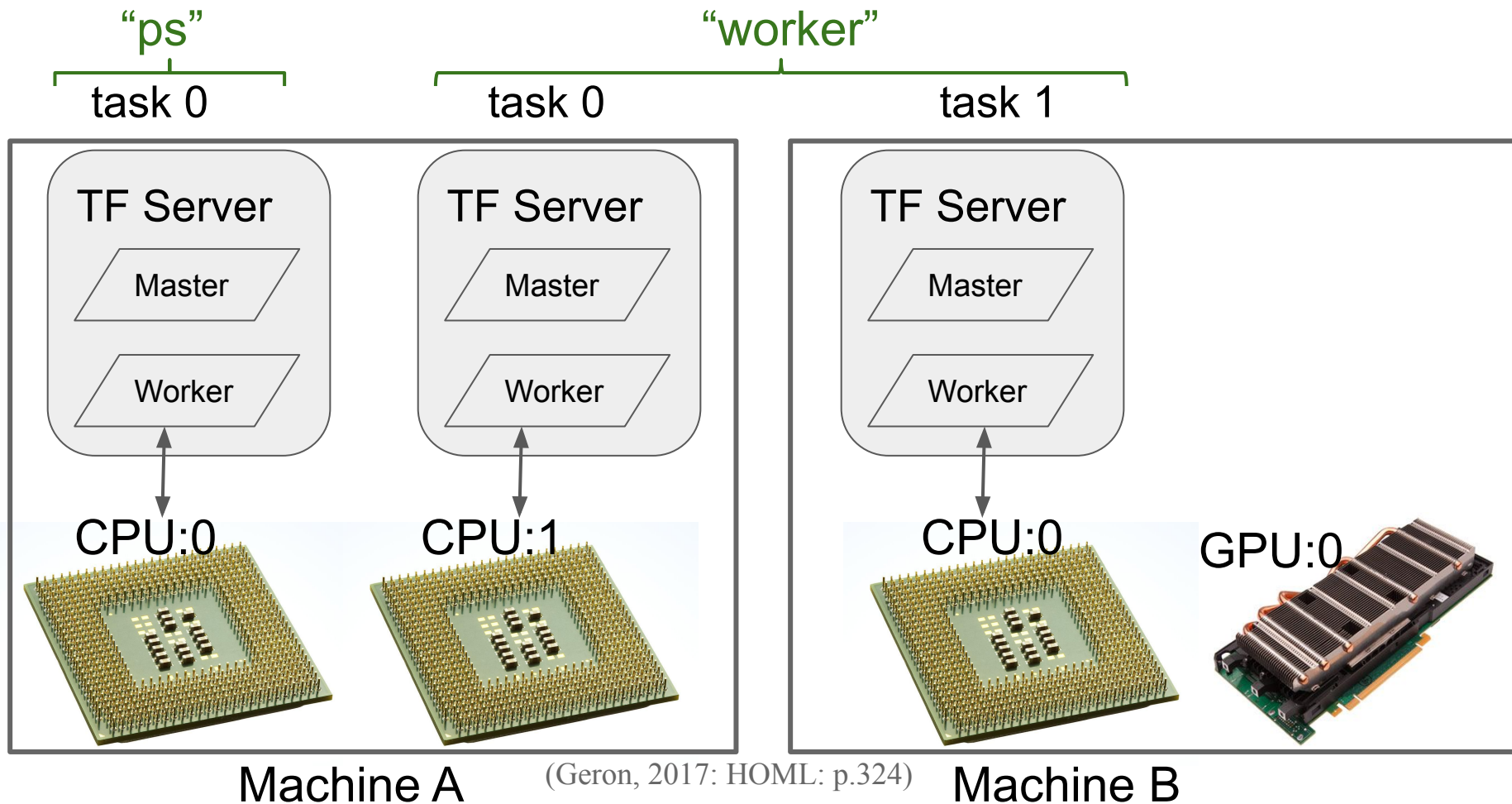
## Parallelisms:

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

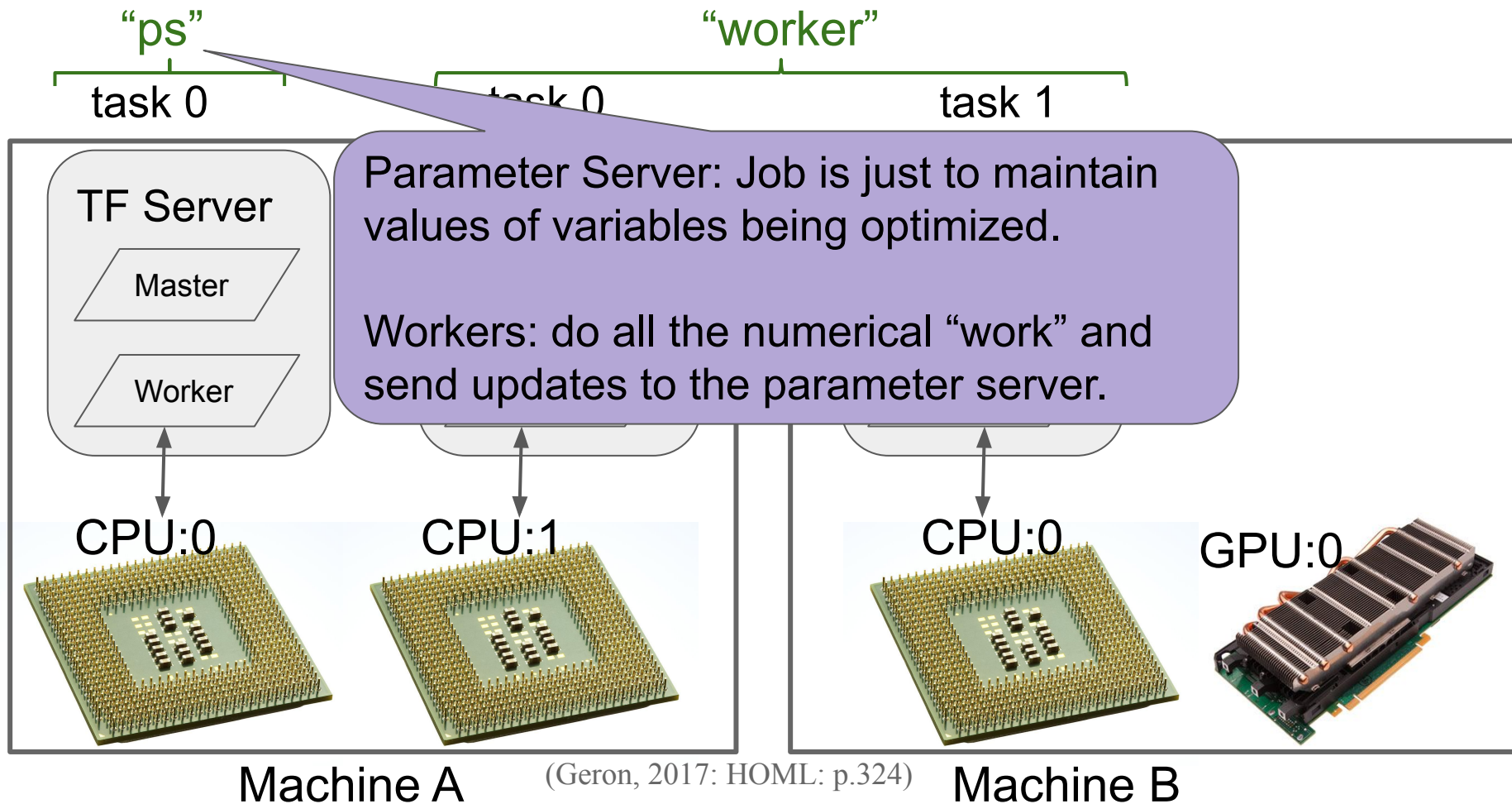
## Model Updates:

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

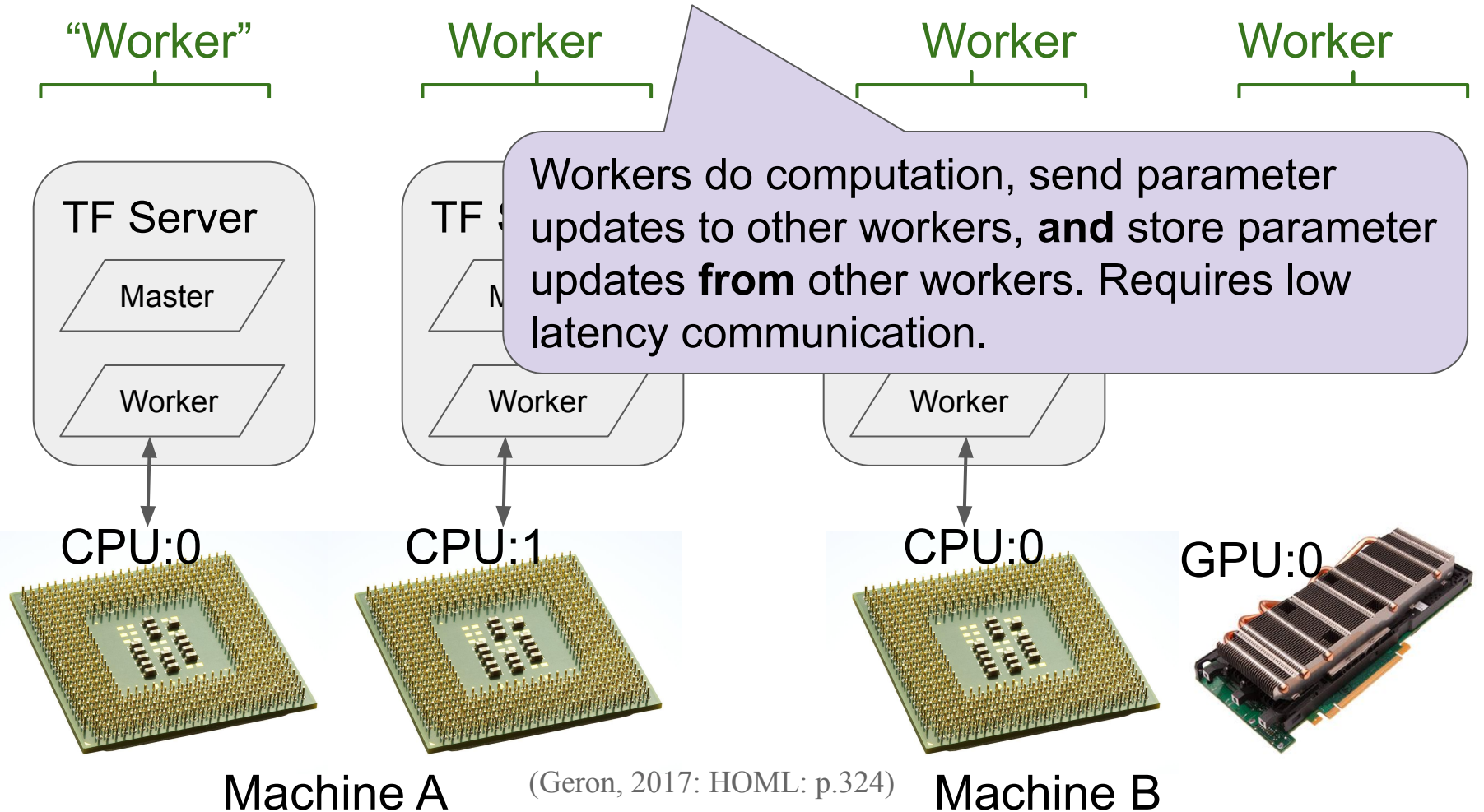
# Asynchronous Parameter Server



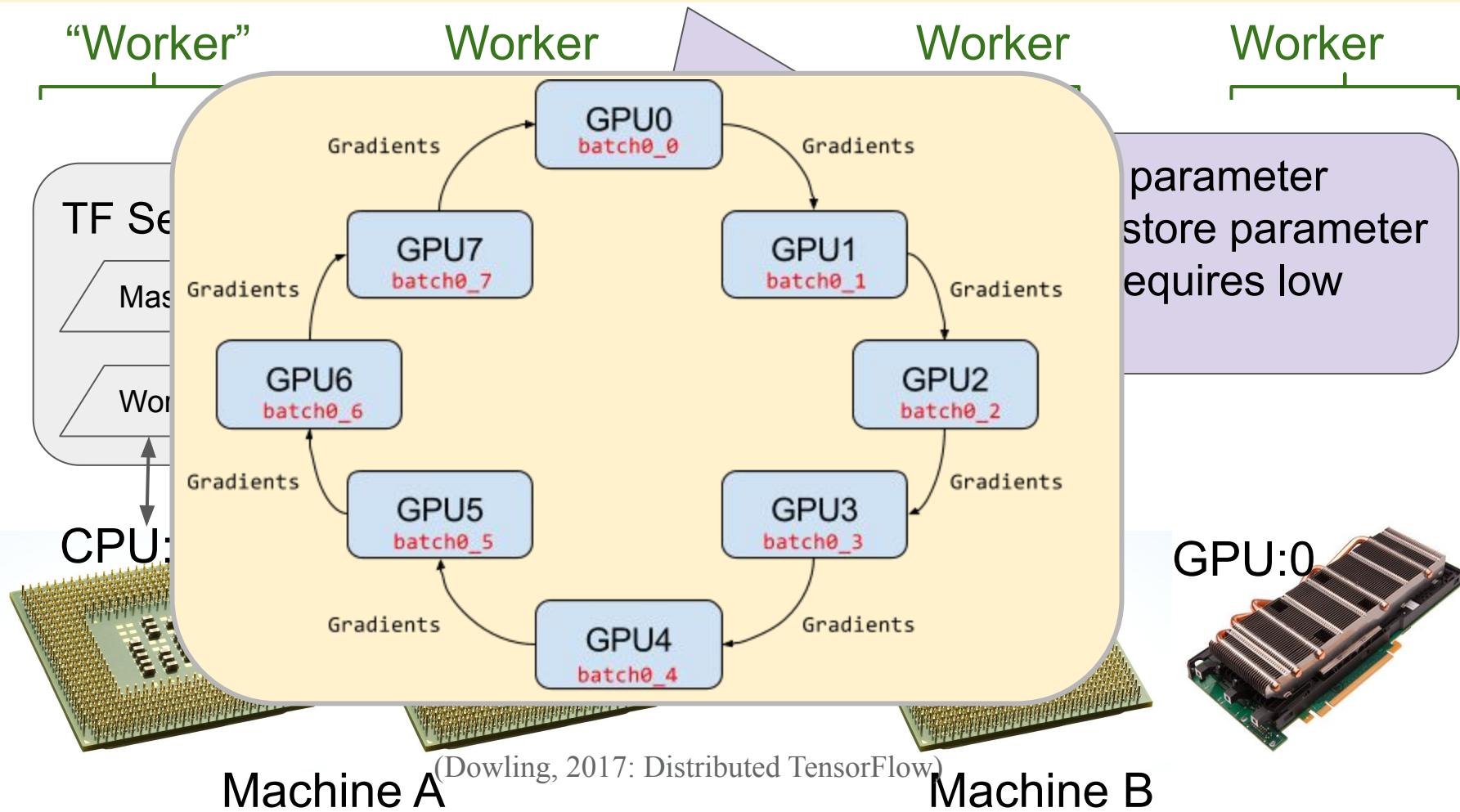
# Asynchronous Parameter Server



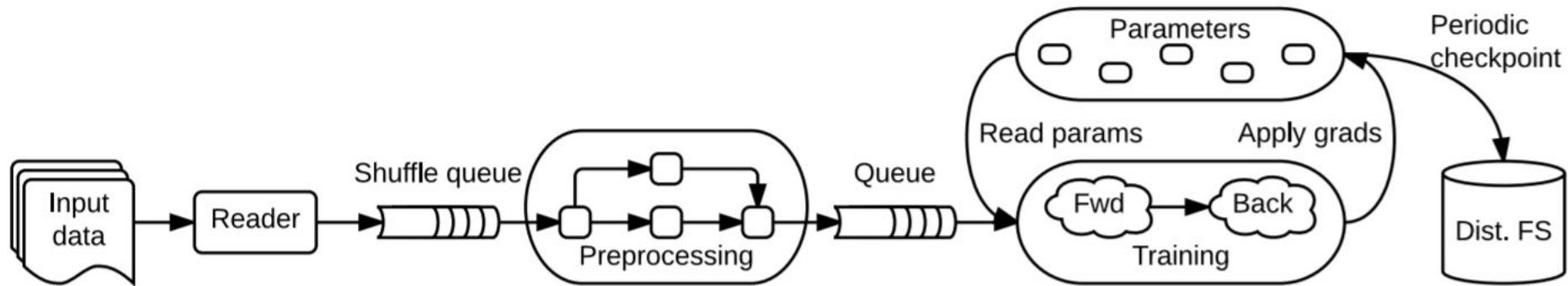
# Synchronous All Reduce



# Synchronous All Reduce



# Distributed TF: Full Pipeline



# Review: Distributed ML

Done very 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: Parallelism can be leveraged  
Con: High communication for transferring Intermediar data.

**Model Parallelism**



# Post-Exam2 Topics:

1. Research Ethics
2. Useful Plots
3. Machine Learning Cross Validation
4. Convolutional Neural Networks
5. Recurrent Neural Networks
6. Transformer Networks

# Post-Exam2 Topics:

1. **Research Ethics**
2. Useful Plots
3. Machine Learning Cross Validation
4. Convolutional Neural Networks
5. Recurrent Neural Networks
6. Transformer Networks

# Ethics in Big Data

Bias

Privacy

Ethical Research Practice

# Ethics in Big Data

## Types of bias:

- **Outcome Disparity:** Predicted distribution given  $A$ ,  
are dissimilar from ideal distribution given  $A$ 
  - Selection bias
  - Label bias
  - Over-amplification
- **Error Disparity:** Predicts less accurate for authors of given demographics.
- **Semantic Bias:** Representations of meaning store demographic associations.

# Two Examples

## The WSJ Effect

model  
accuracy

Jørgensen/Hovy/Sogaard, 2015  
Hovy & Sogaard, 2015

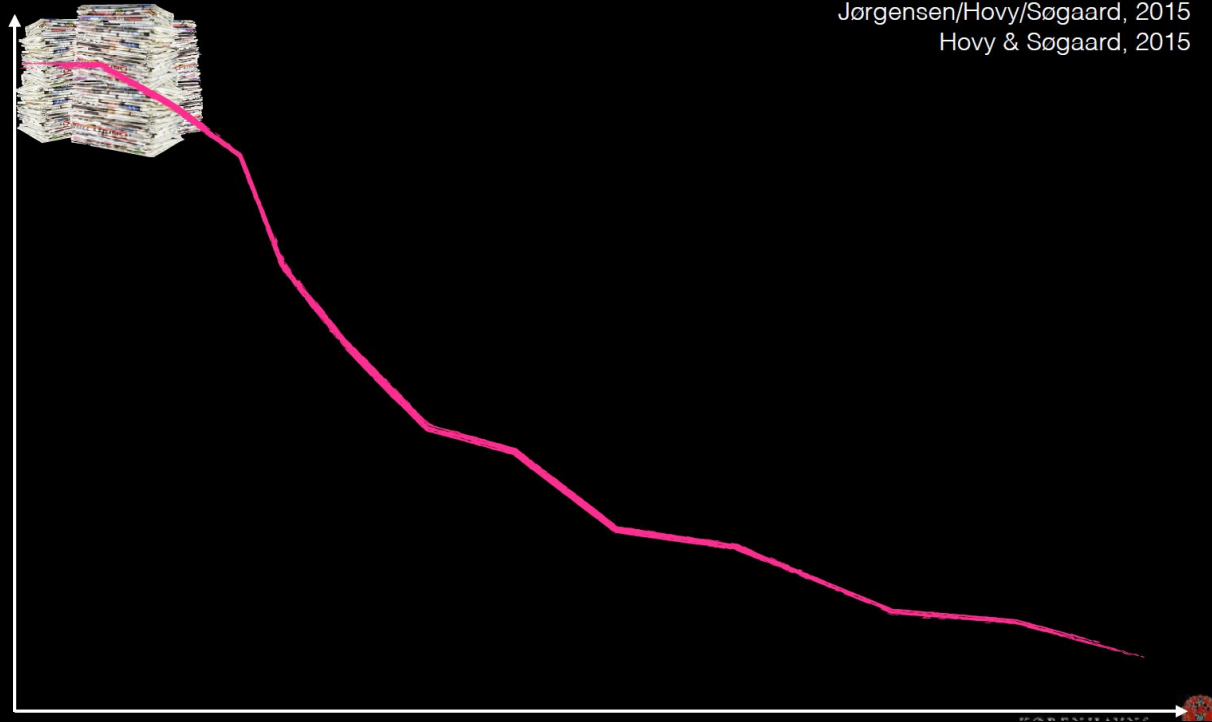


distance from "standard" WSJ author demographics

# Two Examples

## The WSJ Effect

model  
accuracy



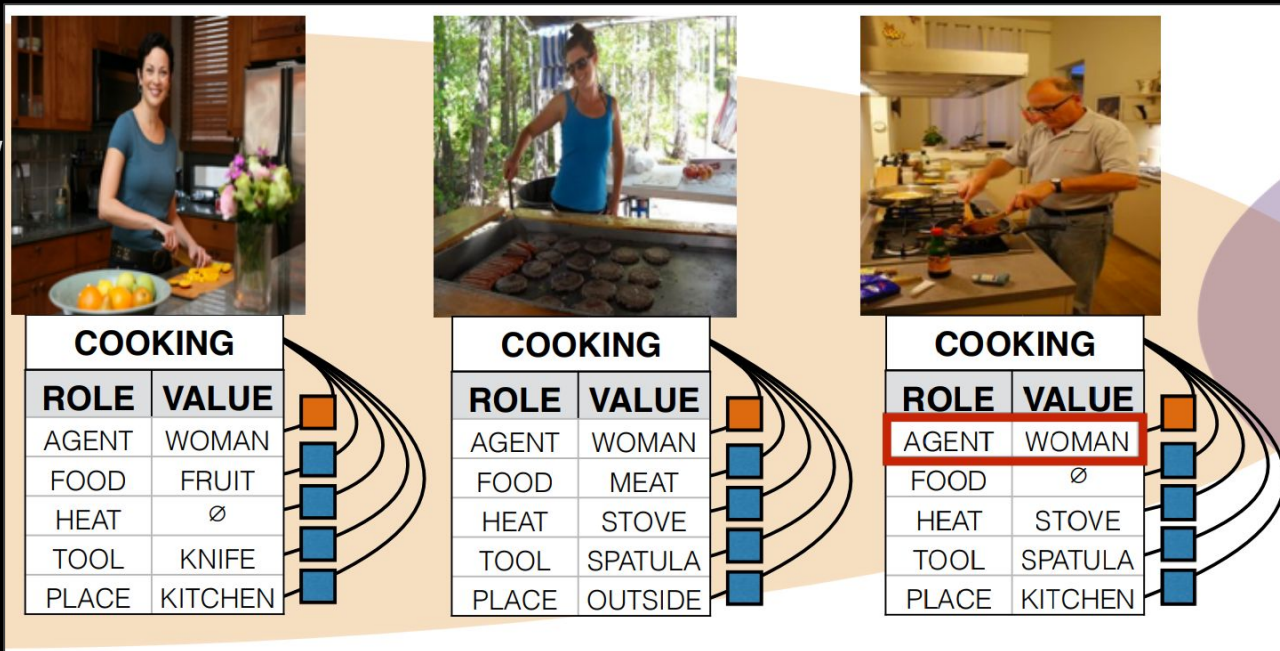
Jørgensen/Hovy/Sogaard, 2015  
Hovy & Sogaard, 2015

distance from "standard" WSJ author demographics

# Two Examples

## The W

model accuracy



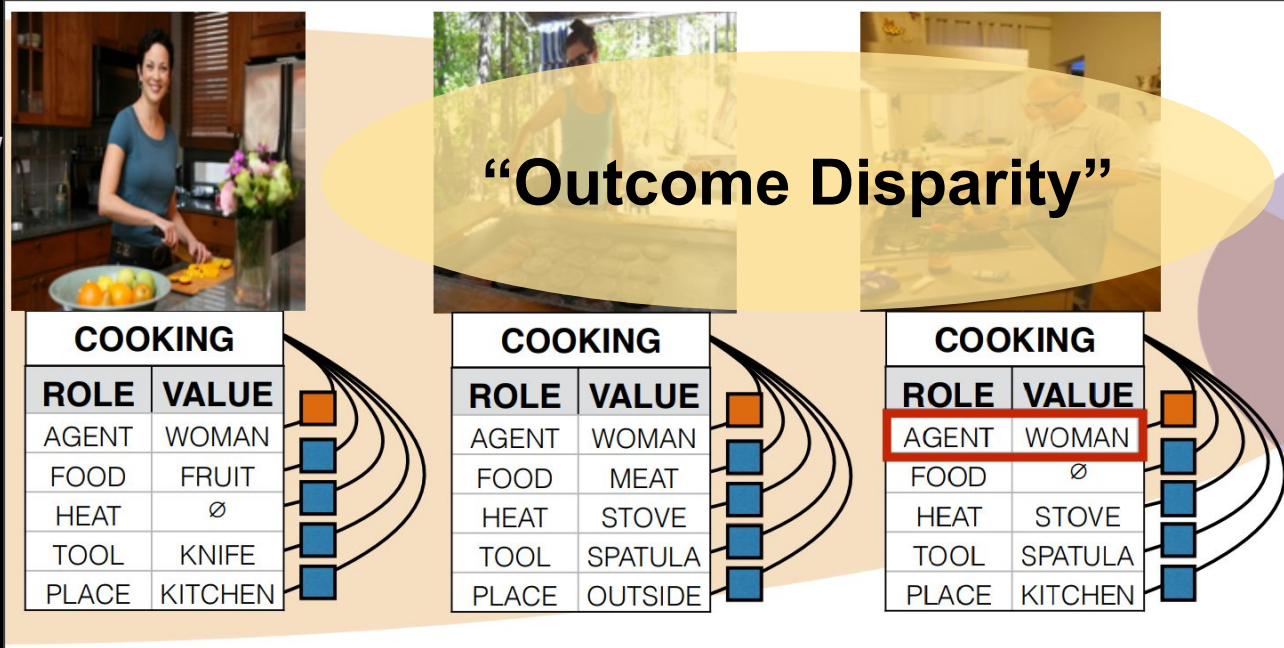
Zhao, Jieyu, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. "Men Also Like Shopping: Reducing Gender Bias Amplification using Corpus-level Constraints." In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*. 2017.

distance from "standard" WSJ author demographics

# Two Examples

## The W

model accuracy



**“Error Disparity”**

Zhao, Jieyu, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. "Men Also Like Shopping: Reducing Gender Bias Amplification using Corpus-level Constraints." In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*. 2017.

distance from "standard" WSJ author demographics



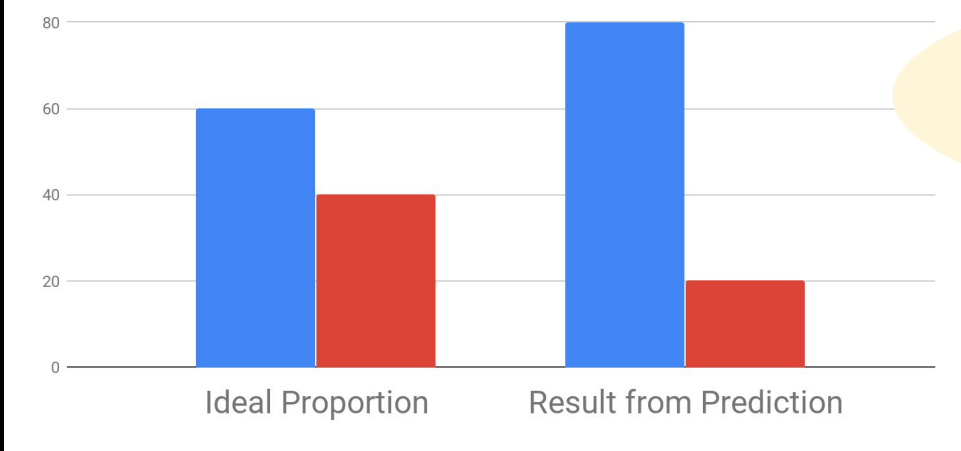
*Our data and models are (human) biased.*

**“Outcome Disparity”**

Person-level  
■ attribute = 1  
■ attribute = 2

**“Error Disparity”**

*Our data and models are (human) biased.*

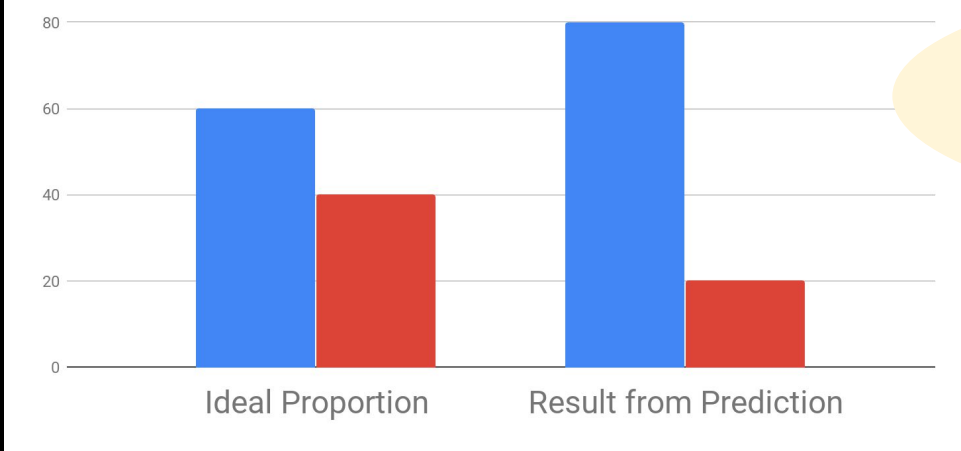


**“Outcome Disparity”**

Person-level  
■ attribute = 1  
■ attribute = 2

**“Error Disparity”**

*Our data and models are (human) biased.*



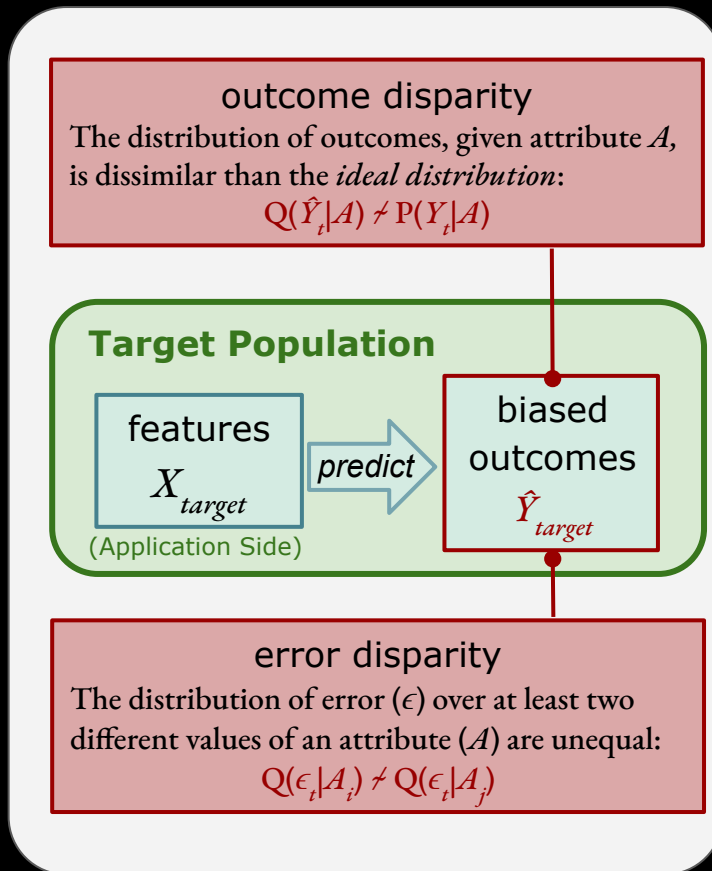
**“Outcome Disparity”**

Person-level  
■ attribute = 1  
■ attribute = 2

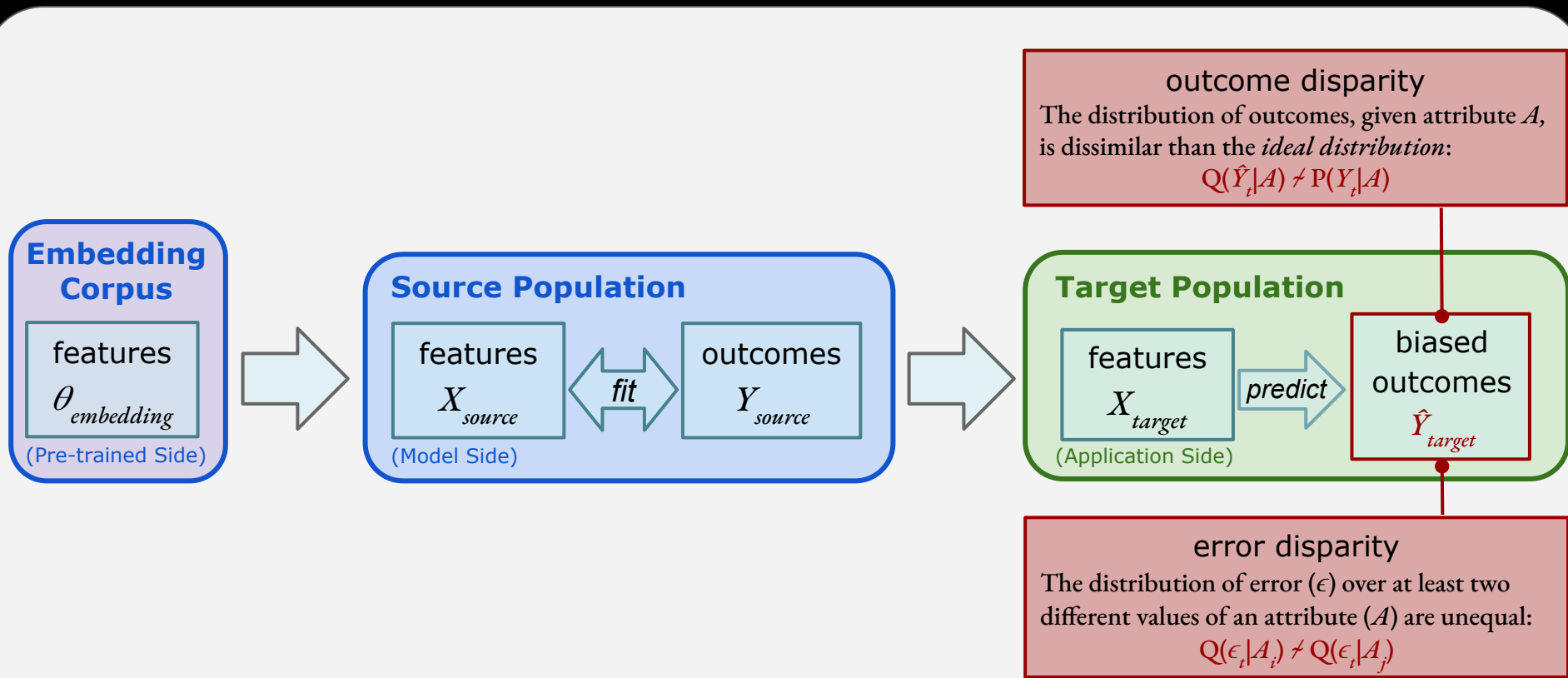
**“Error Disparity”**



# Predictive Bias Framework



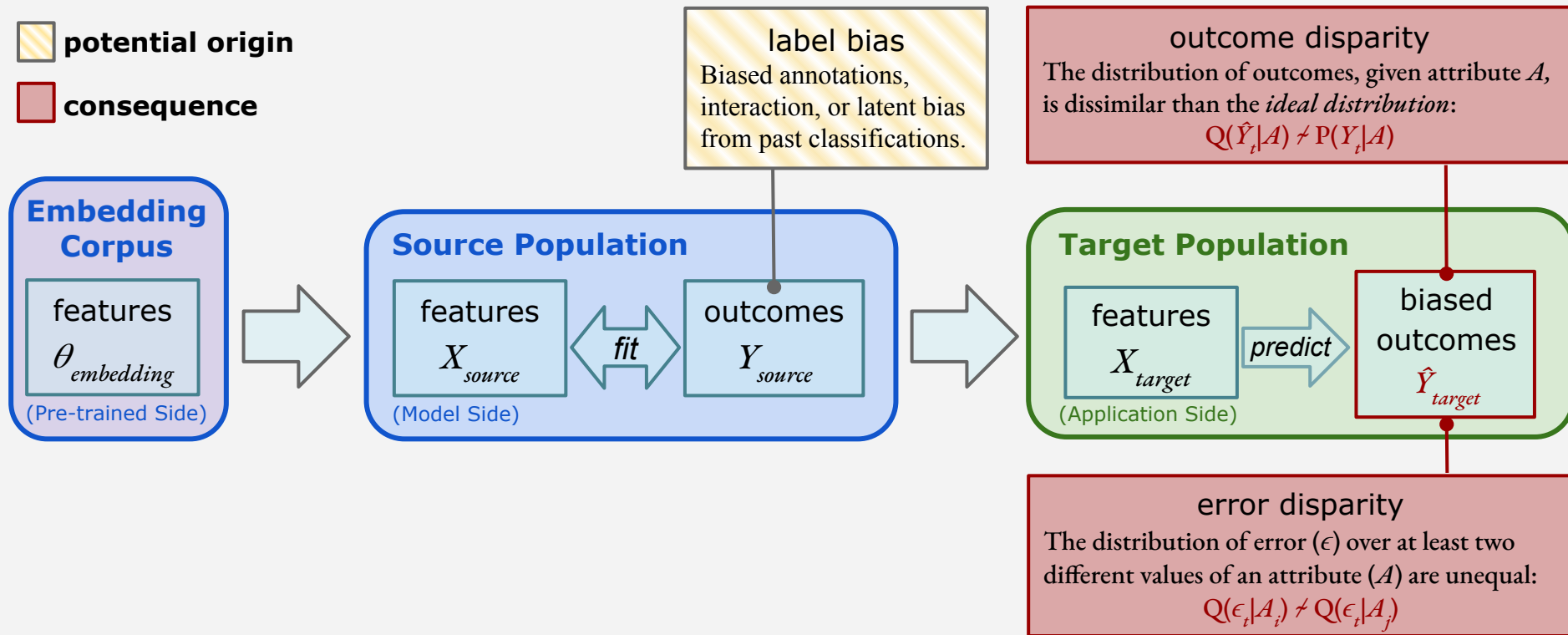
# Predictive Bias Framework



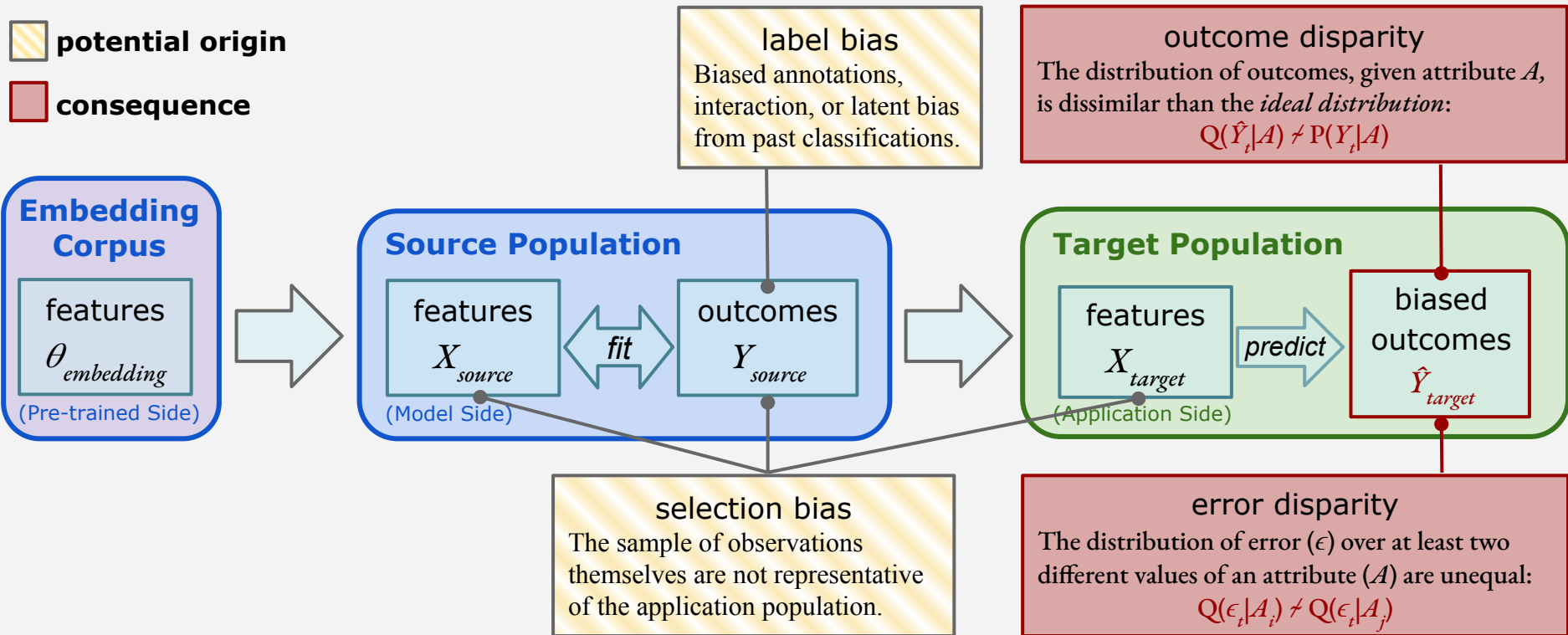
# Predictive Bias Framework

 **potential origin**

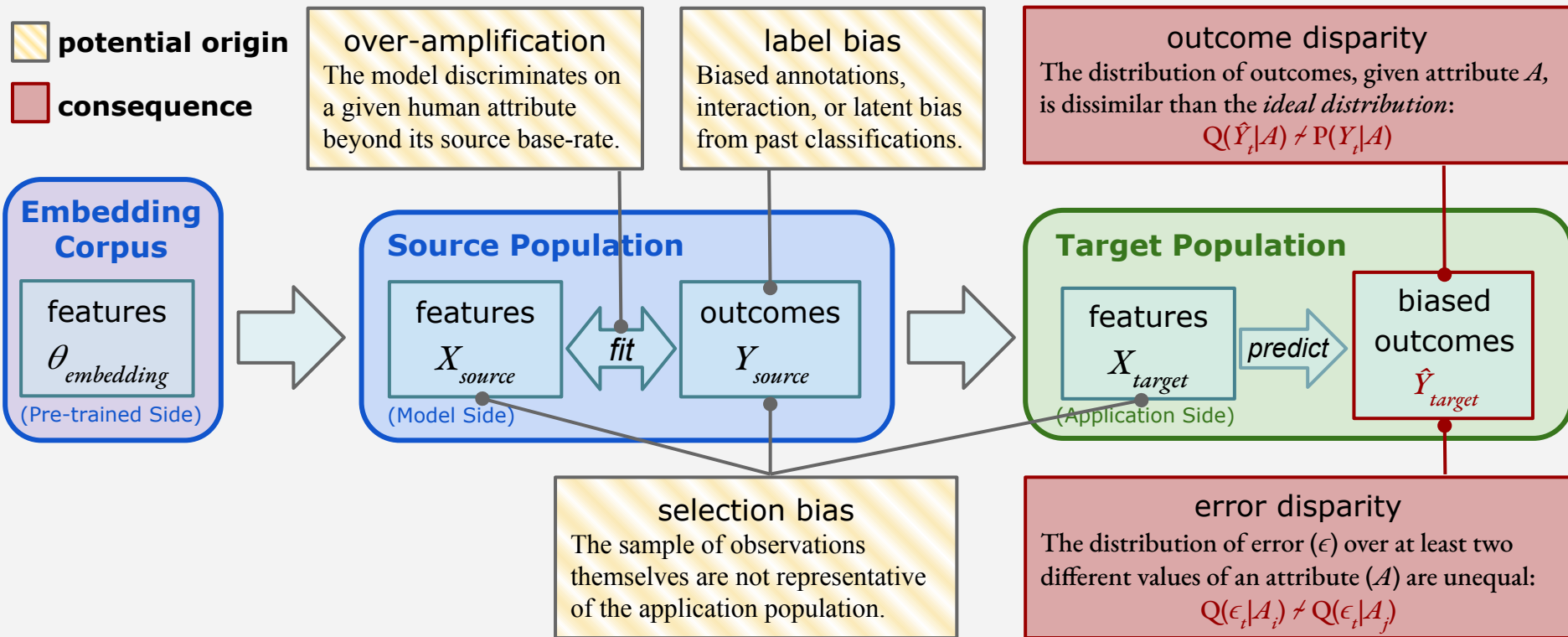
 **consequence**



# Predictive Bias Framework

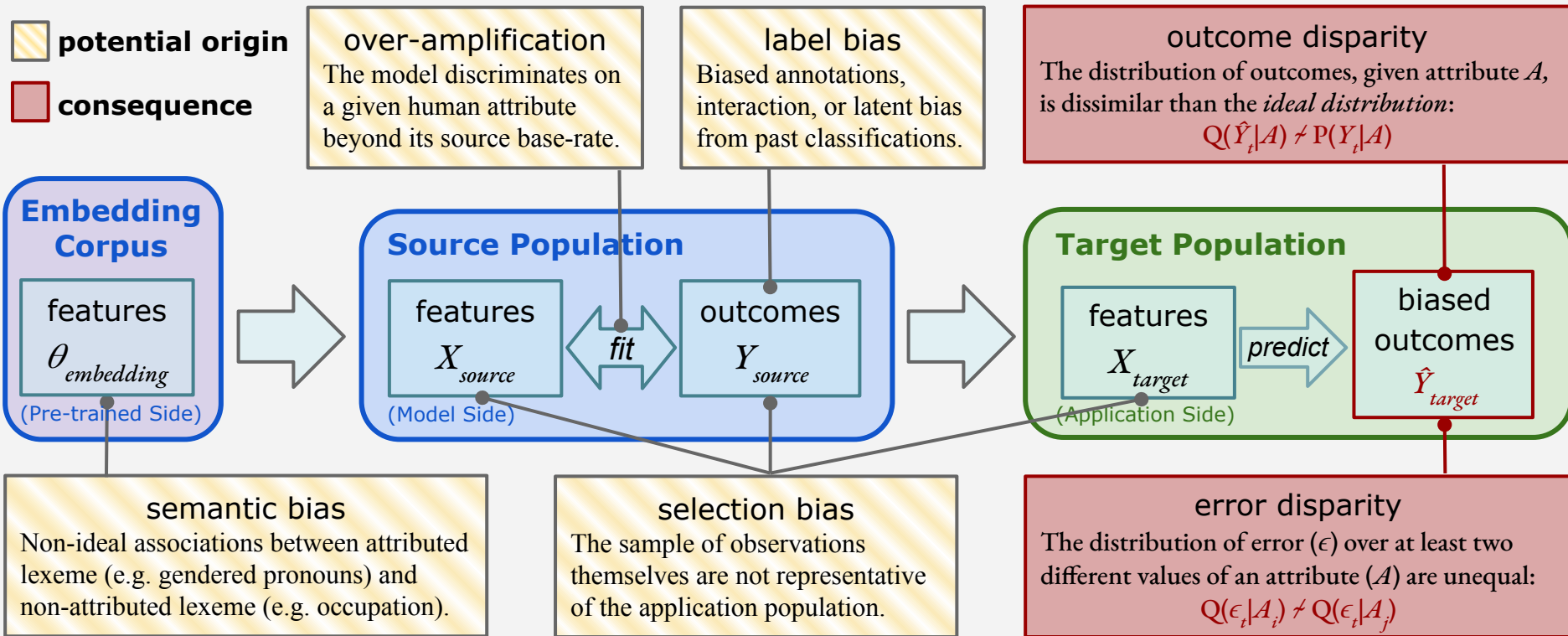


# Predictive Bias Framework





# Predictive Bias Framework



# Predictive Bias Framework

E.g. Coreference resolution:  
connecting entities to references (i.e. pronouns).

*“The doctor told Mary that she had run some blood tests.”*

## semantic bias

Non-ideal associations between attributed lexeme (e.g. gendered pronouns) and non-attributed lexeme (e.g. occupation).

## selection bias

The sample of observations themselves are not representative of the application population.

## error disparity

The distribution of error ( $\epsilon$ ) over at least two different values of an attribute ( $A$ ) are unequal:

$$Q(\epsilon_t | A_i) \neq Q(\epsilon_t | A_j)$$

# Ethics in Big Data

## Types of bias:

- **Outcome Disparity:** Predicted distribution given  $A$ ,  
are dissimilar from ideal distribution given  $A$ 
  - Selection bias
  - Label bias
  - Over-amplification
- **Error Disparity:** Predicts less accurate for authors of given demographics.
- **Semantic Bias:** Representations of meaning store demographic associations.

# Ethics in Big Data

## Privacy

- Risk Categories:
  - Revealing unintended private information
  - Targeted persuasion



# Ethics in Big Data

## Privacy

- Risk Categories:
  - Revealing unintended private information
  - Targeted persuasion
- Mitigation strategies:
  - Informed consent -- let participants know
  - Do not share / secure storage
  - *Federated learning* -- separate and obfuscate to the point of preserving privacy
  - Transparency in information targeting  
“You are being shown this ad because ...”



# Ethics in Big Data

Human Subjects Research

Observational versus Interventional

# Ethics in Big Data

## Human Subjects Research

### Observational versus Interventional

(The Belmont Report, 1979)

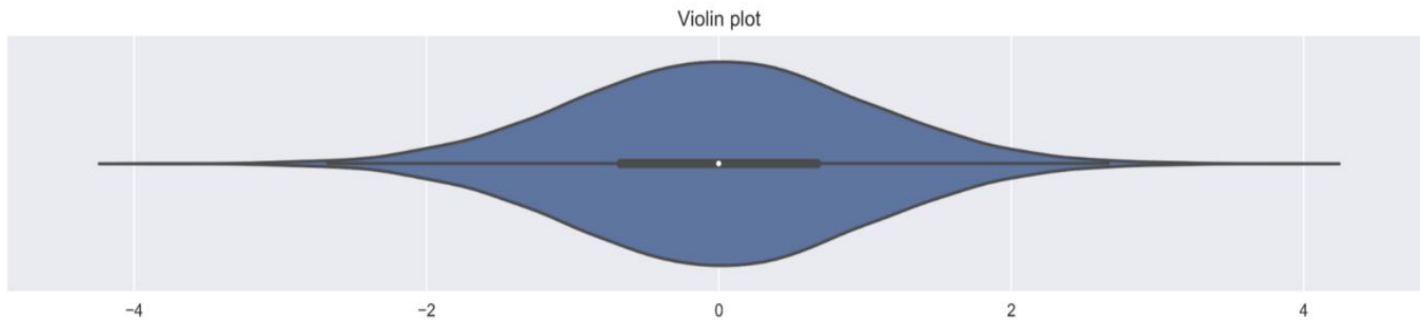
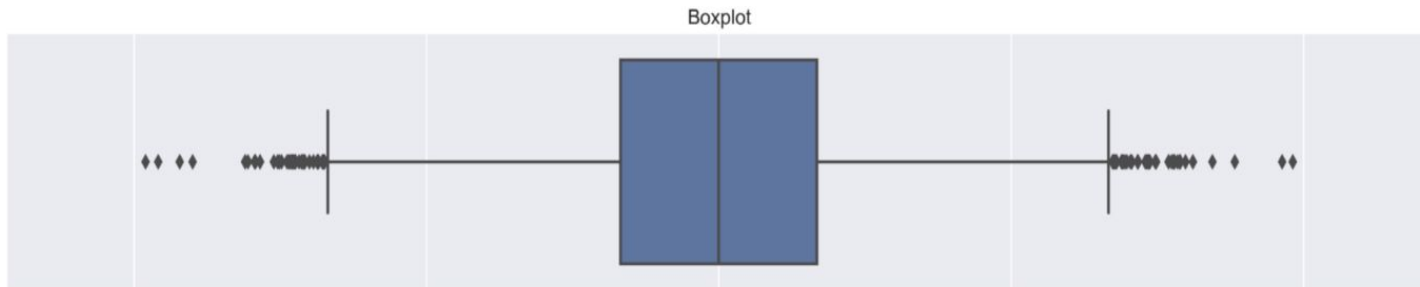
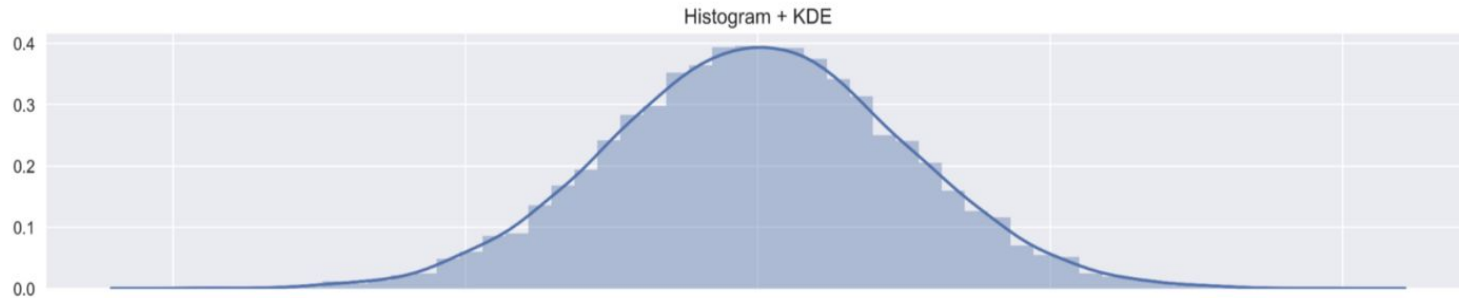
- (i) Distinction of research from practice.
- (ii) Risk-Benefit criteria
- (iii) Appropriate selection of human subjects for participation in research
- (iv) Informed consent in various research settings.

# Post-Exam2 Topics:

1. Research Ethics
2. **Useful Plots**
3. Machine Learning Cross Validation
4. Convolutional Neural Networks
5. Recurrent Neural Networks
6. Transformer Networks



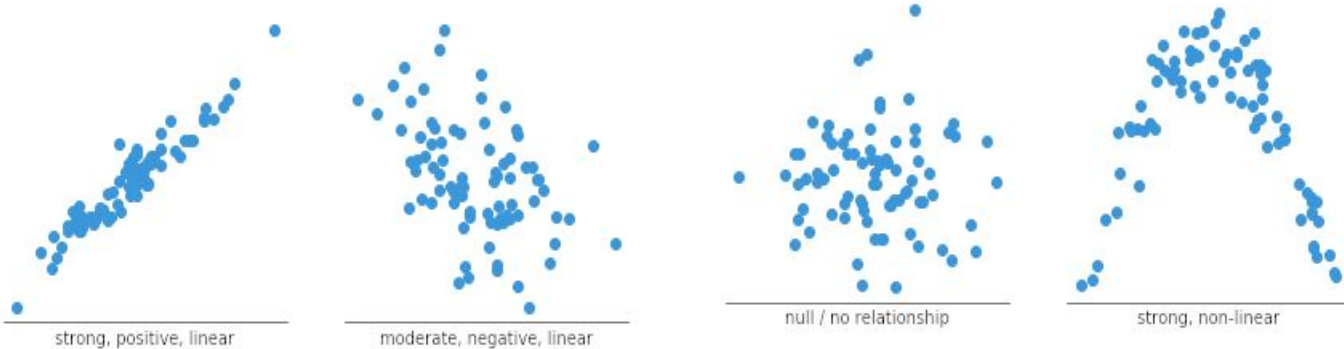
# Useful Plots: For distributions



[\(Lewinson, 2019\)](#)

# Useful Plots: Correlation

**Scatter Plot:** for two variables expected to be associated (with optional regression line)



(Chartio)

**Correlation Matrix:** for comparing associations between many variables (use Bonferroni correction if hyp testing)

	FriendSize	Intelligence Quotient	Income	Sat W/ Life	Depression
F1	0.03	0.04	0.12	0.02	-0.1
F2	0.04	-0.26	-0.19	-0.09	0.11
F3	-0.07	-0.13	0.02	-0.02	-0.02
F4	-0.03	0.27	-0.08	-0.12	0.11
F5	-0.01	0.23	0.29	0.07	-0.21

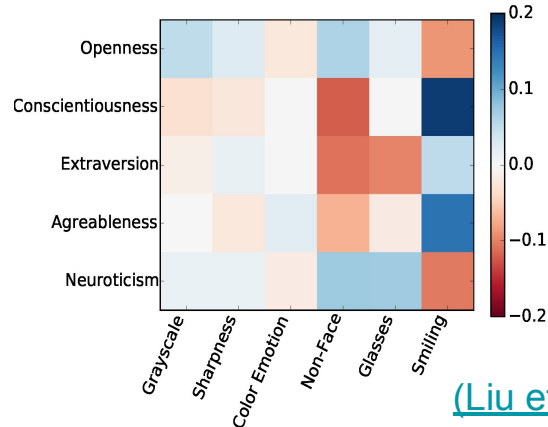


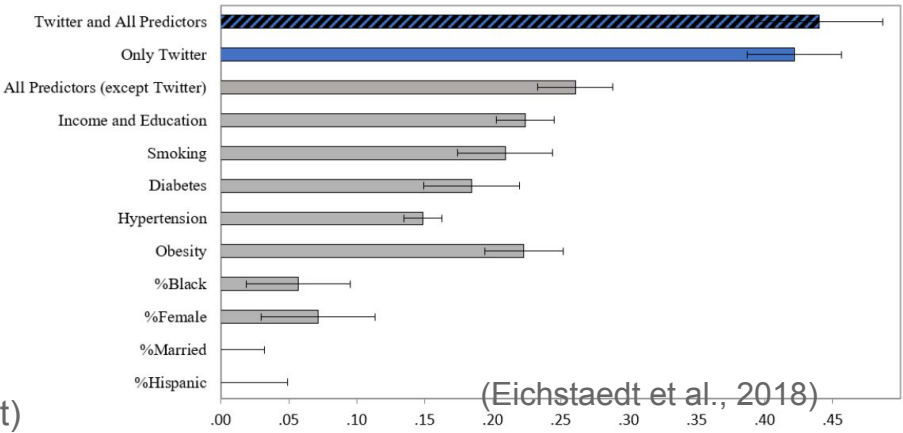
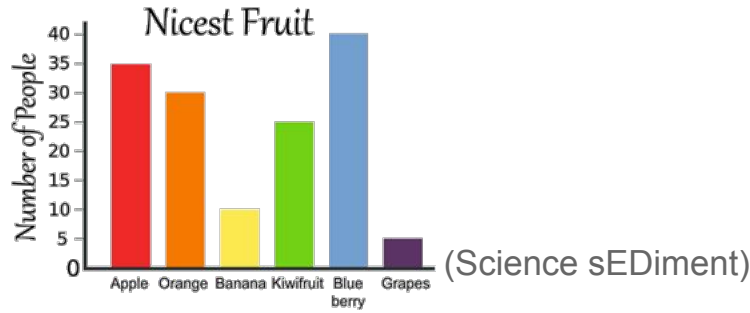
Fig 3. Individual factor correlations with outcomes. Note how F4 which captures the use of swear words negatively correlates with Satisfaction with Life (SWL).

<https://doi.org/10.1371/journal.pone.0201703.g003>

(Liu et al., 2016)

# Useful Plots: Any Values

**Bar Plot:** To visually compare values under different selections/conditions.



**Line Plot:** When one variable has a natural ordering (e.g. time)

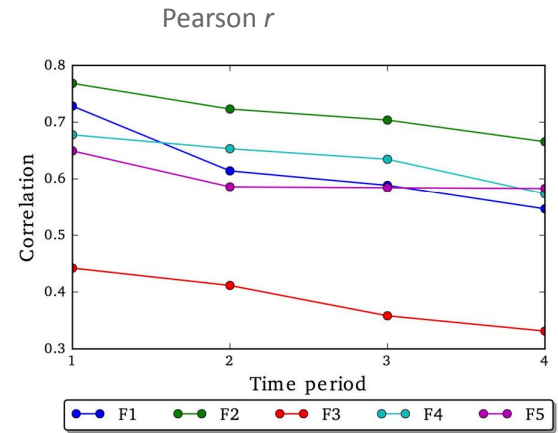
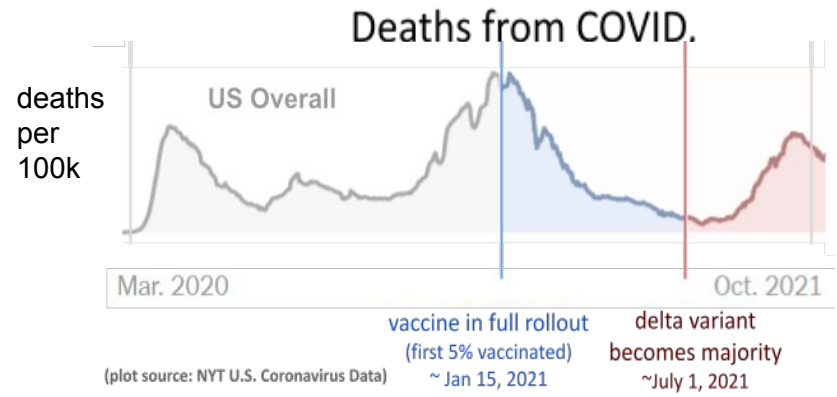
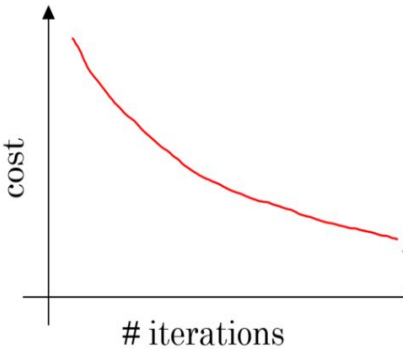


Fig 6. Test re-test validity of our learned factors.  
<https://doi.org/10.1371/journal.pone.0201703.g006>

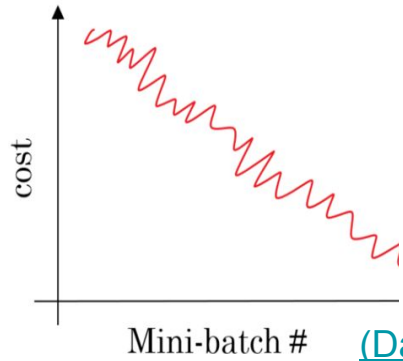
# Useful Plots: Prediction

**Learning Curve:** for plotting error from gradient descent.

for a model with convex optimization (i.e. linear regression)

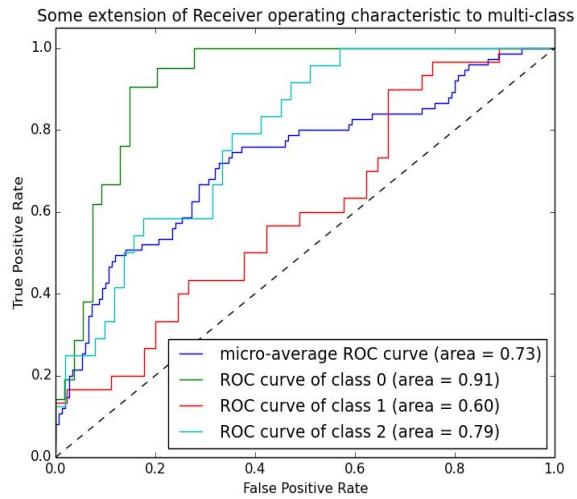


for a model with non-convex optimization (i.e. most deep learning)

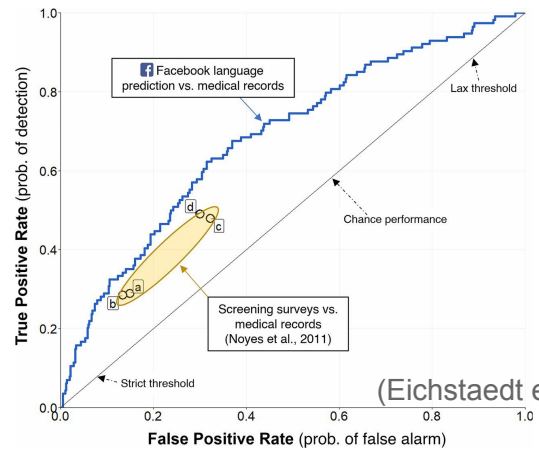


[\(Dabura, 2017\)](#)

**ROC Plot:** for visualizing true-positive to false-positive rates (used for AUC metric)



[\(PLOT ROC\)](#)

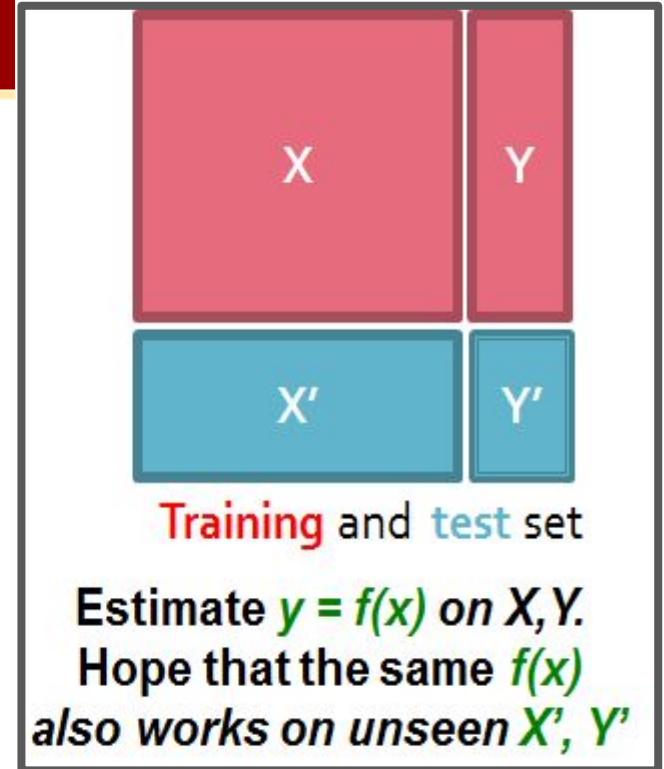
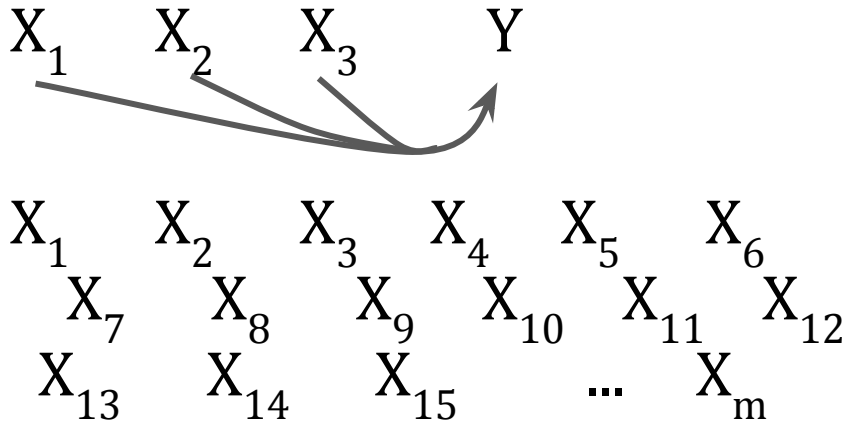


[\(Eichstaedt et al., 2018\)](#)

# Post-Exam2 Topics:

1. Research Ethics
2. Useful Plots
3. **Machine Learning Cross Validation**
4. Convolutional Neural Networks
5. Recurrent Neural Networks
6. Transformer Networks

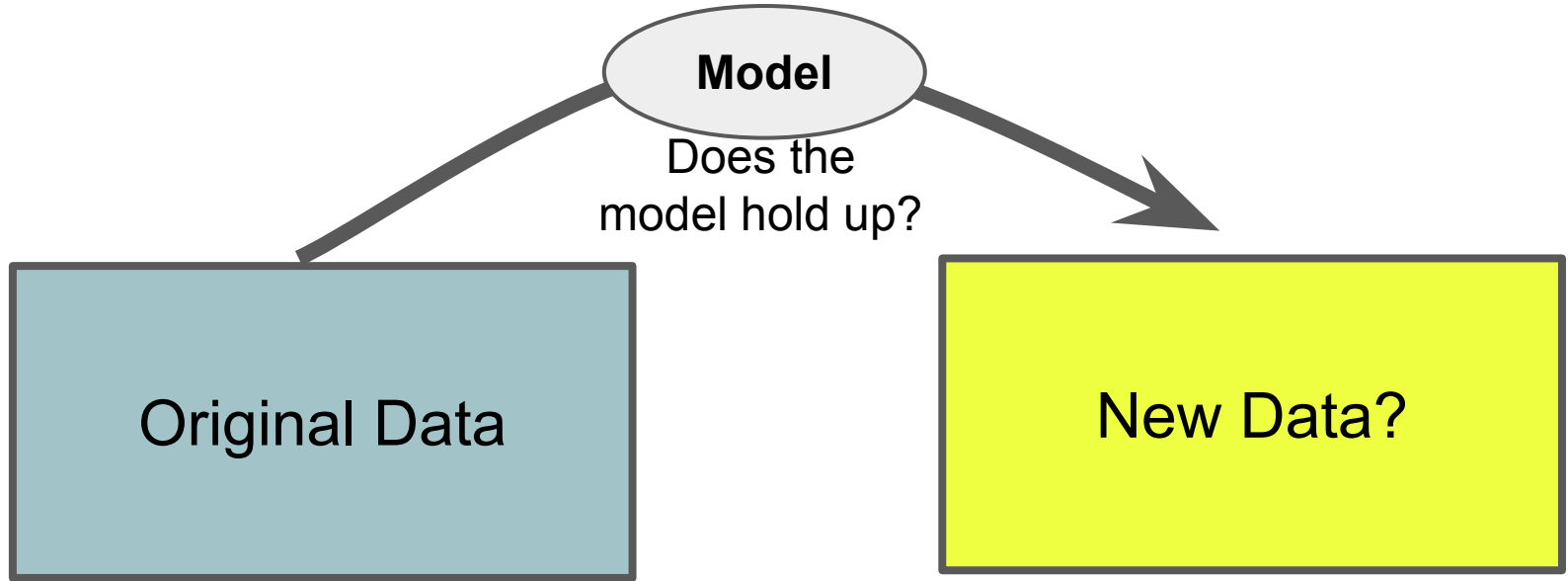
# Supervised Learning



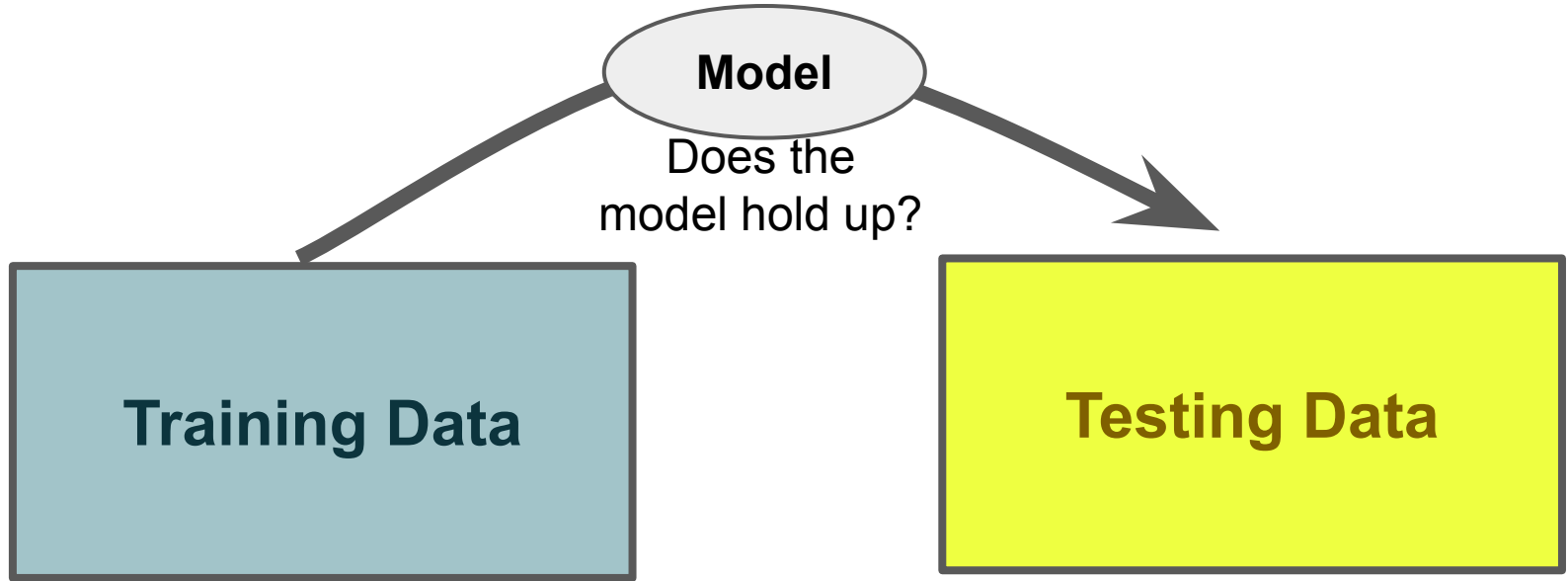
J. Leskovec, A. Rajaraman, J. Ullman: Mining of Massive Datasets, <http://www.mmms.org>

Task: Determine a function,  $f$  (or parameters to a function) such that  $f(X) = Y$

# Common Goal: Generalize to new data

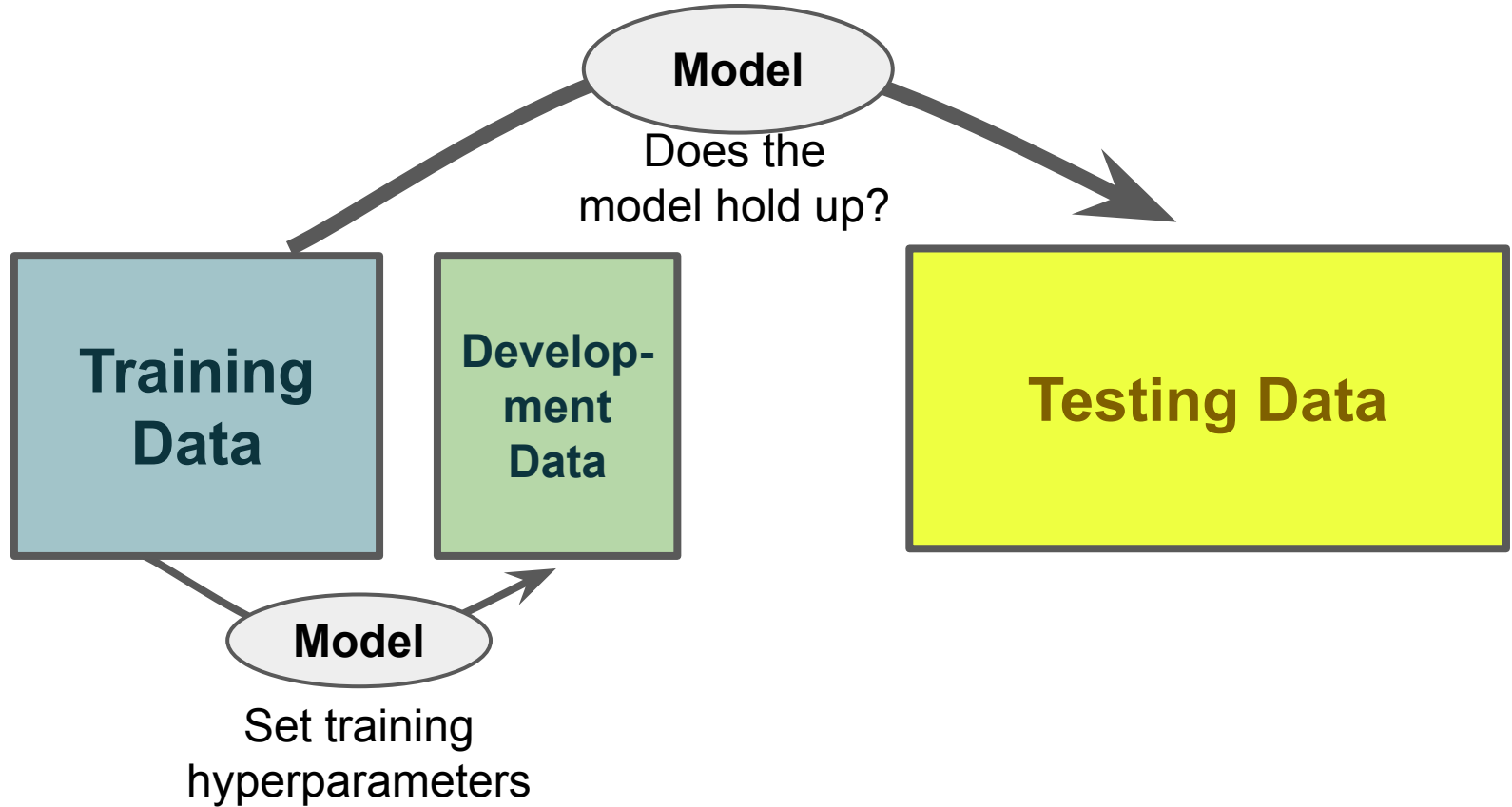


# Common Goal: Generalize to new data





# ML: GOAL



# N-Fold Cross Validation

Goal: Decent estimate of model accuracy



Iter 1



Iter 2



Iter 3



....



observed  
dep  
variable

<b>test</b>	<b>test</b>	<b>test</b>	<b>test</b>	<b>test</b>
-------------	-------------	-------------	-------------	-------------

estimated  
dep  
variable

<b>ptest</b>	<b>ptest</b>	<b>ptest</b>	<b>ptest</b>	<b>ptest</b>
--------------	--------------	--------------	--------------	--------------

# Review: Distributed ML

Done very 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

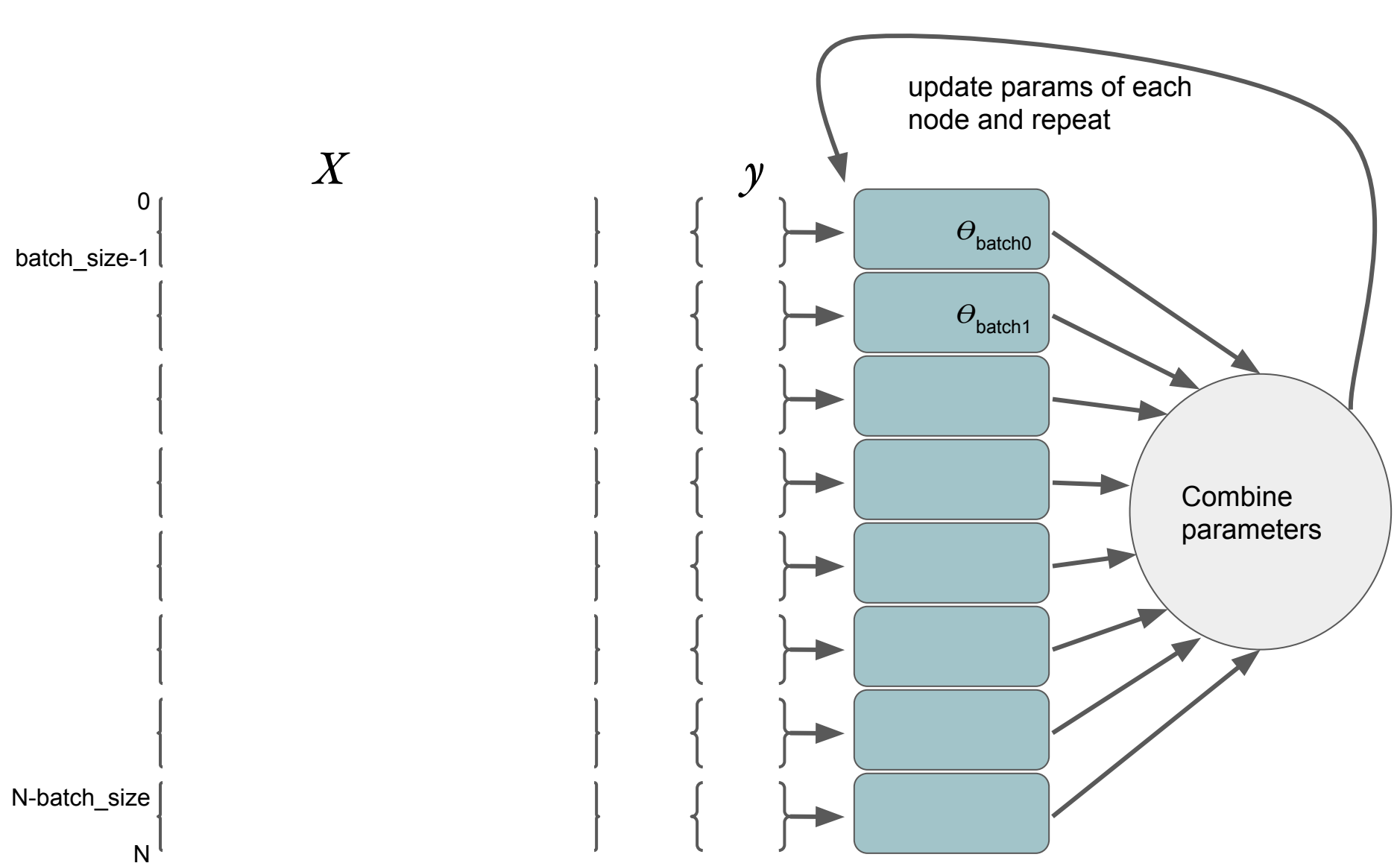
Preferred method for big data or very complex models (i.e. models with many internal parameters).

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

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

Pro: Parallelism can be leveraged  
Con: High communication for transferring Intermediar data.

**Model Parallelism**



# Post-Exam2 Topics:

1. Research Ethics
2. Useful Plots
3. Machine Learning Cross Validation
4. **Recurrent Neural Networks**
5. Convolutional Neural Networks
6. Transformer Networks

# From Linear Models to Neural Nets

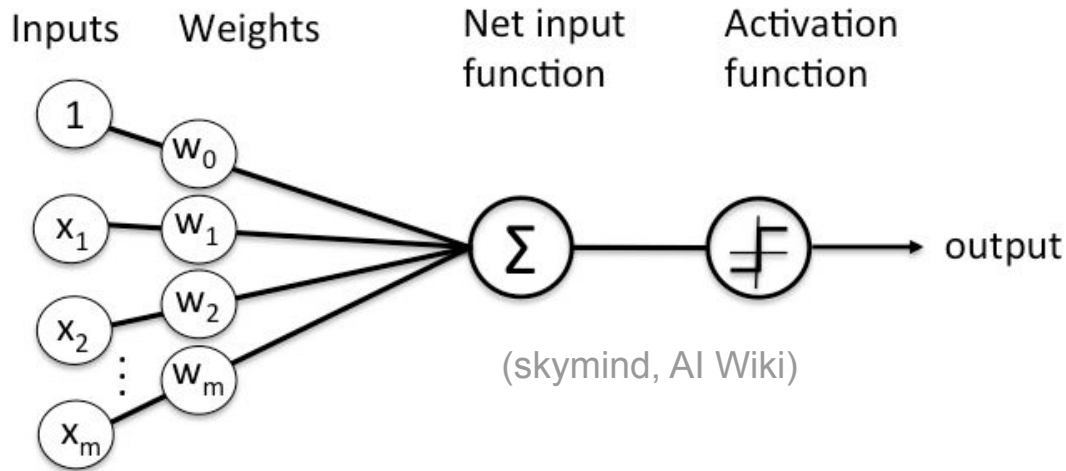
Linear Regression:  $y = wX$

Neural Network Nodes:  $output = f(wX)$

# From Linear Models to Neural Nets

Linear Regression:  $y = wX$

Neural Network Nodes:  $output = f(wX)$

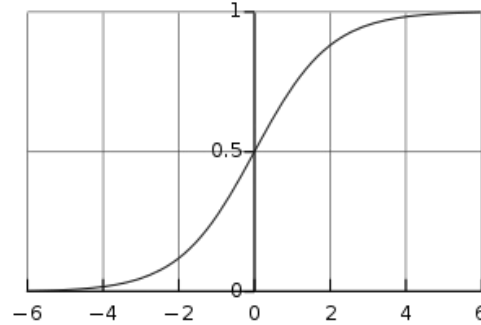




# Common Activation Functions

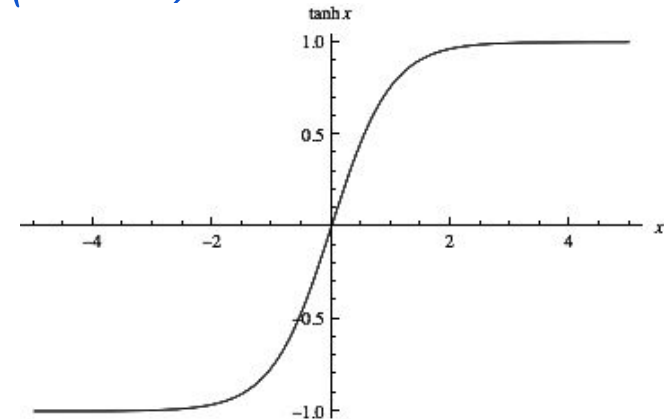
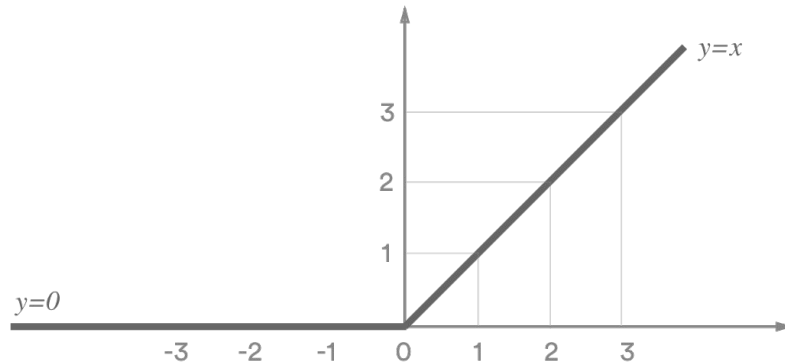
$$z = wX$$

Logistic:  $\sigma(z) = 1 / (1 + e^{-z})$



Hyperbolic tangent:  $\tanh(z) = 2\sigma(2z) - 1 = (e^{2z} - 1) / (e^{2z} + 1)$

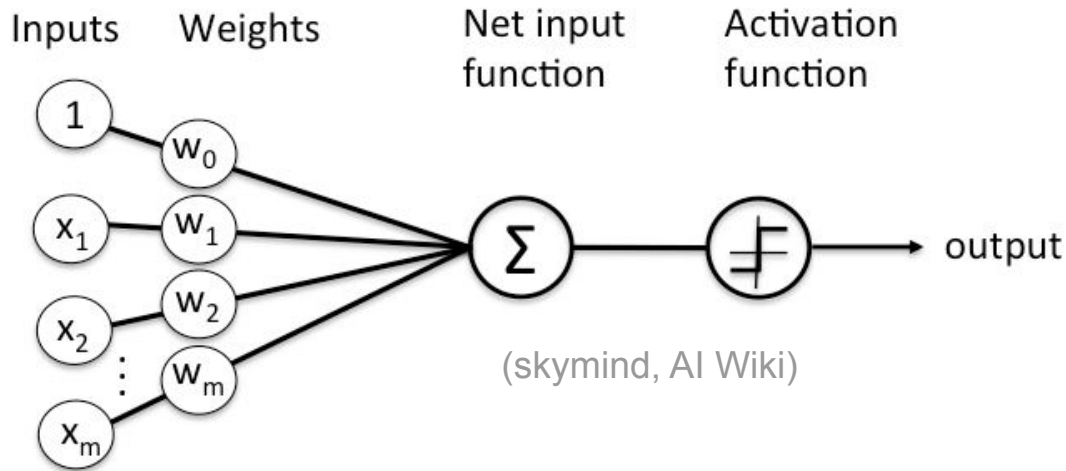
Rectified linear unit (ReLU):  $ReLU(z) = \max(0, z)$



# From Linear Models to Neural Nets

Linear Regression:  $y = wX$

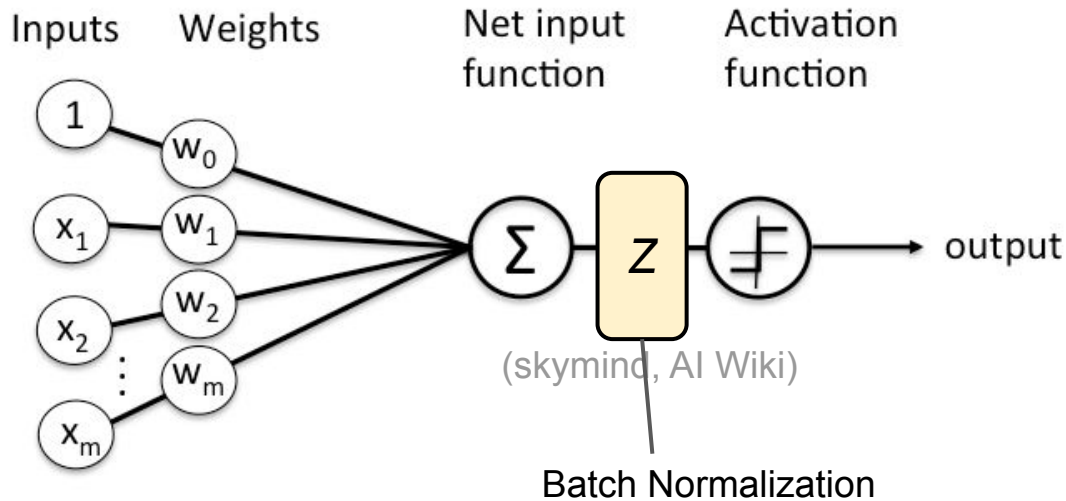
Neural Network Nodes:  $output = f(wX)$



# From Linear Models to Neural Nets

Linear Regression:  $y = wX$

Neural Network Nodes:  $output = f(wX)$



# Batch Normalization

**Input:** Values of  $x$  over a mini-batch:  $\mathcal{B} = \{x_{1\dots m}\}$ ;

Parameters to be learned:  $\gamma, \beta$

**Output:**  $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{ mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{ mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{ normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{ scale and shift}$$

(Ioffe and Szegedy, 2015)

# Batch Normalization

This is just standardizing!  
(but within the current batch of observations)

**Input:** Values of  $x$  over a mini-batch:  $\mathcal{B} = \{x_{1\dots m}\}$ ;

Parameters to be learned:  $\gamma, \beta$

**Output:**  $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{ mini-batch mean}$$

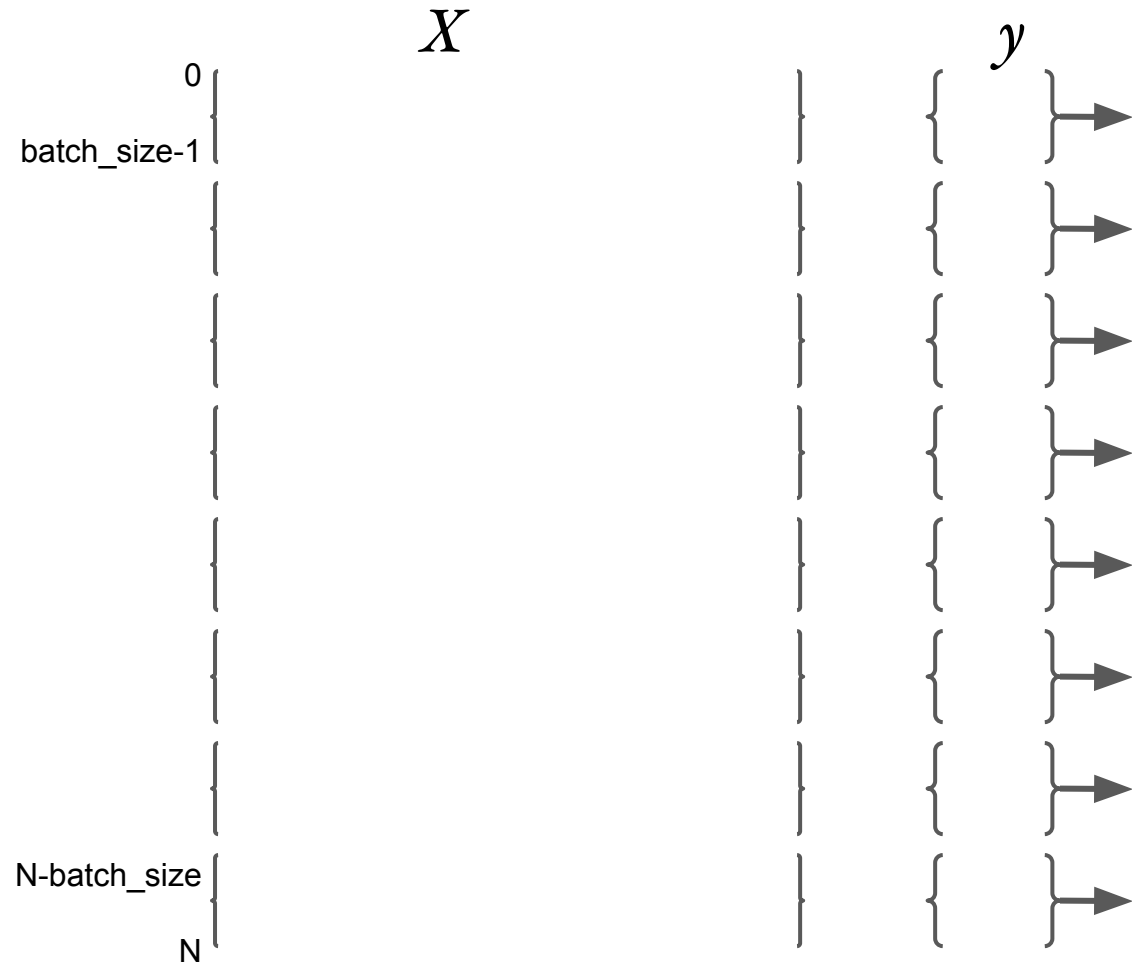
$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{ mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{ normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{ scale and shift}$$

(Ioffe and Szegedy, 2015)

# Batch Normalization



# Batch Normalization

This is just standardizing!  
(but within the current batch of observations)

**Input:** Values of  $x$  over a mini-batch:  $\mathcal{B} = \{x_{1\dots m}\}$ ;  
Parameters to be learned:  $\gamma, \beta$

**Output:**  $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{ mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{ mini-batch variance}$$

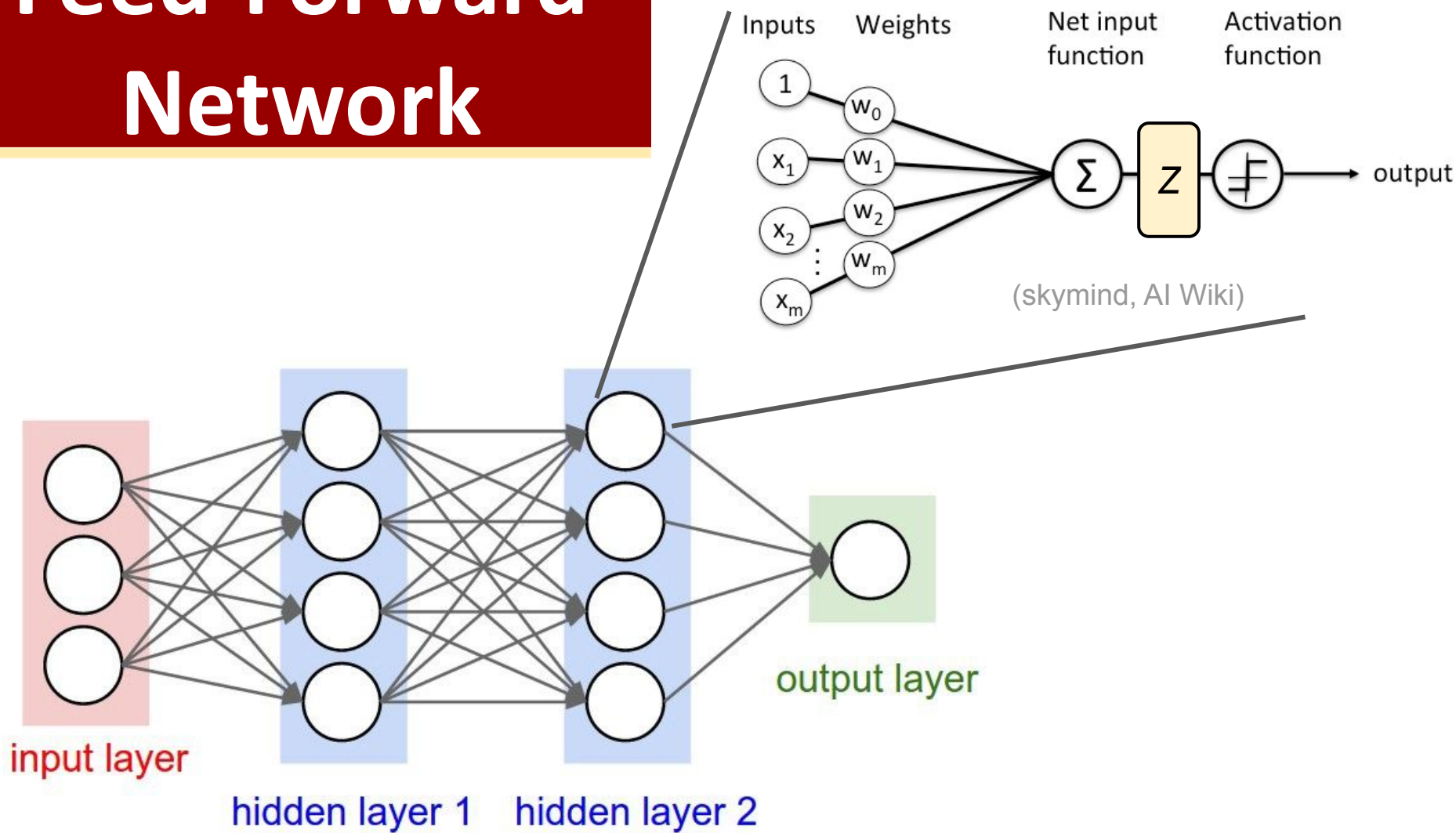
$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{ normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{ scale and shift}$$

## Why?

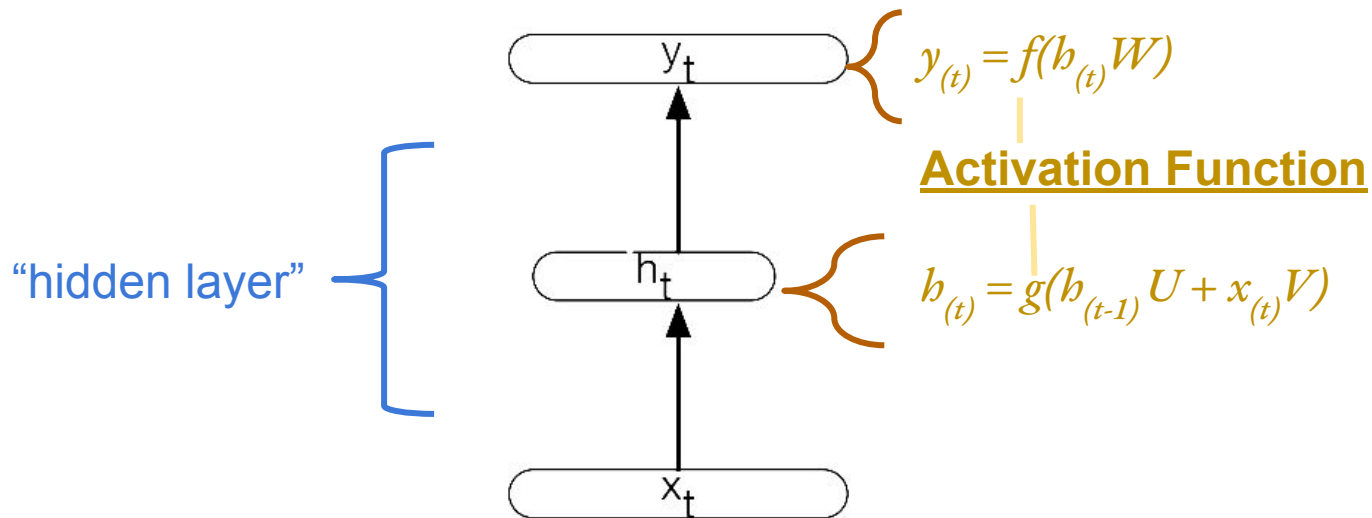
- Empirically, it works!
- Conceptually, generally good for weight optimization to keep data within a reasonable range (dividing by sigma) and such that positive weights move it up and negative down (centering).
- Small effect: When done over mini-batches, adds regularization due to differences between batches.

# Feed-Forward Network



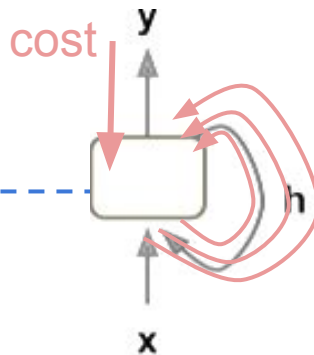


# Recurrent Neural Network



**Figure 9.2** Simple recurrent neural network after Elman (Elman, 1990). The hidden layer includes a recurrent connection as part of its input. That is, the activation value of the hidden layer depends on the current input as well as the activation value of the hidden layer from the previous timestep. (Jurafsky, 2019)

# RNN: Optimization



## Backward Propagation through Time

...

*#define forward pass graph:*

$h_{(0)} = \emptyset$

for  $i$  in range(1, len(x)):

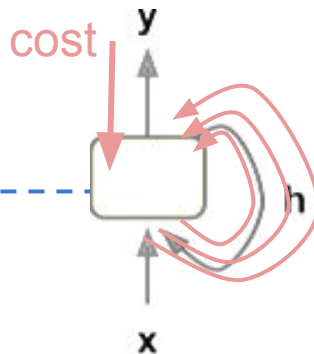
$h_{(i)} = \text{tf.tanh}(\text{tf.matmul}(U, h_{(i-1)}) + \text{tf.matmul}(W, x_{(i)}))$  *#update hidden state*

$y_{(i)} = \text{tf.softmax}(\text{tf.matmul}(V, h_{(i)}))$  *#update output*

...

$\text{cost} = \text{tf.reduce\_mean}(-\text{tf.reduce\_sum}(y * \text{tf.log}(y\_pred)))$

# RNN: Optimization



## Backward Propagation through Time

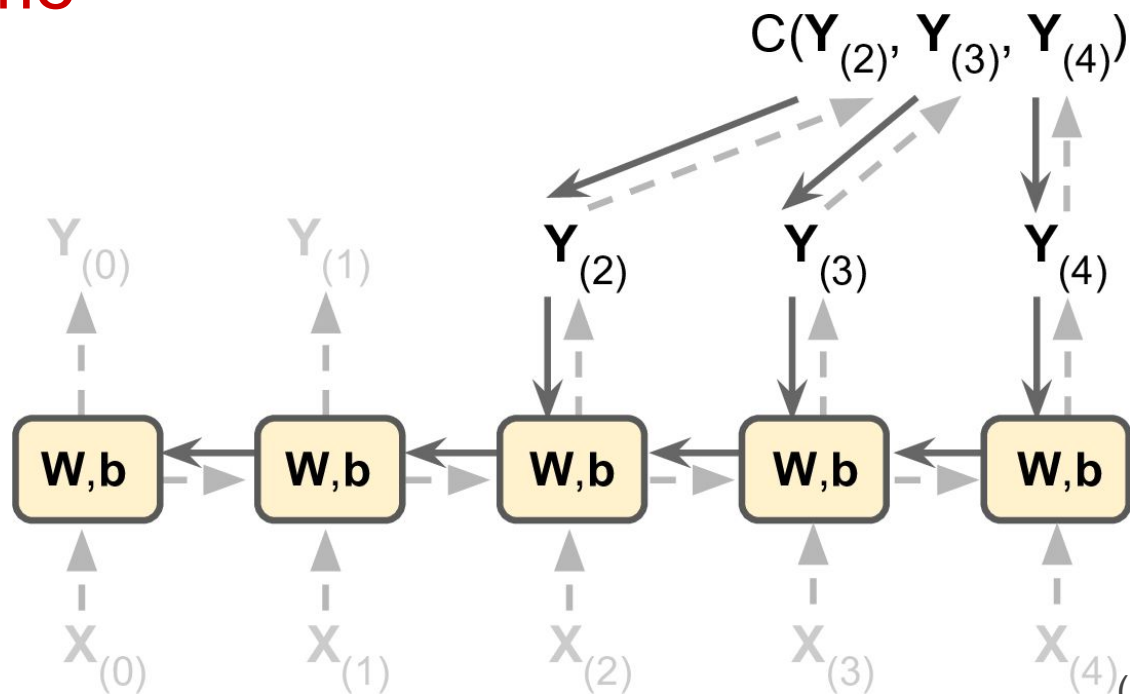
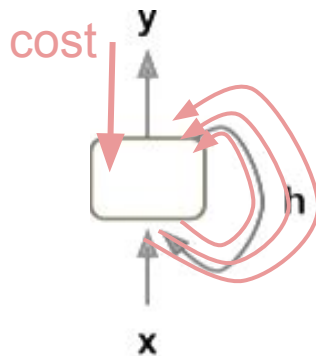
```
...  
#define forward pass graph:  
h(0) = 0  
for i in range(1, len(x)):  
    h(i) = tf.tanh(tf.matmul(U,  
state  
    y(i) = tf.softmax(tf.matmul  
...  
cost = tf.reduce_mean(-tf.reduce
```

To find the gradient for the overall graph, we use **back propogation**, which *essentially* chains together the gradients for each node (function) in the graph.

With many recursions, the gradients can vanish or explode (become too large or small for floating point operations).

# RNN: Optimization

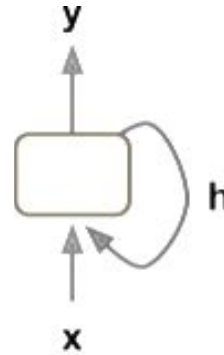
## Backward Propagation through Time



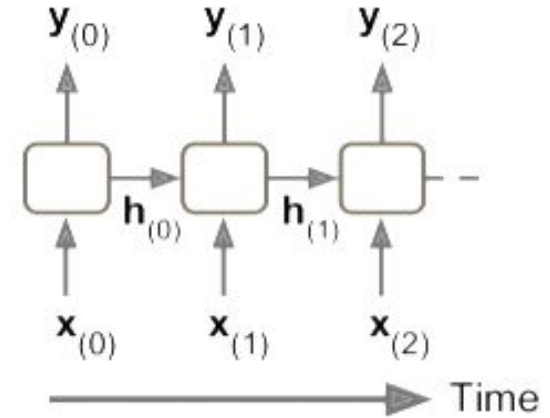
(Geron, 2017)

# How to Addressing Vanishing Gradient?

Dominant approach: Use Long Short Term Memory Networks (LSTM)



RNN model

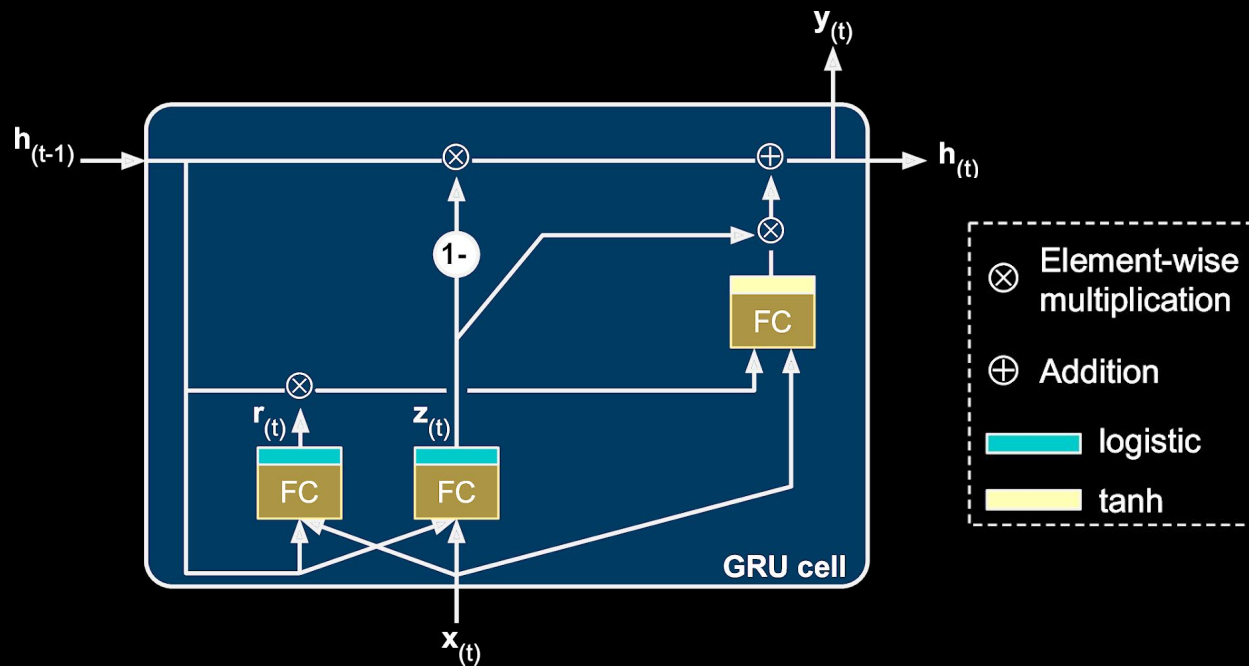


“unrolled” depiction

(Geron, 2017)

# RNN: The GRU

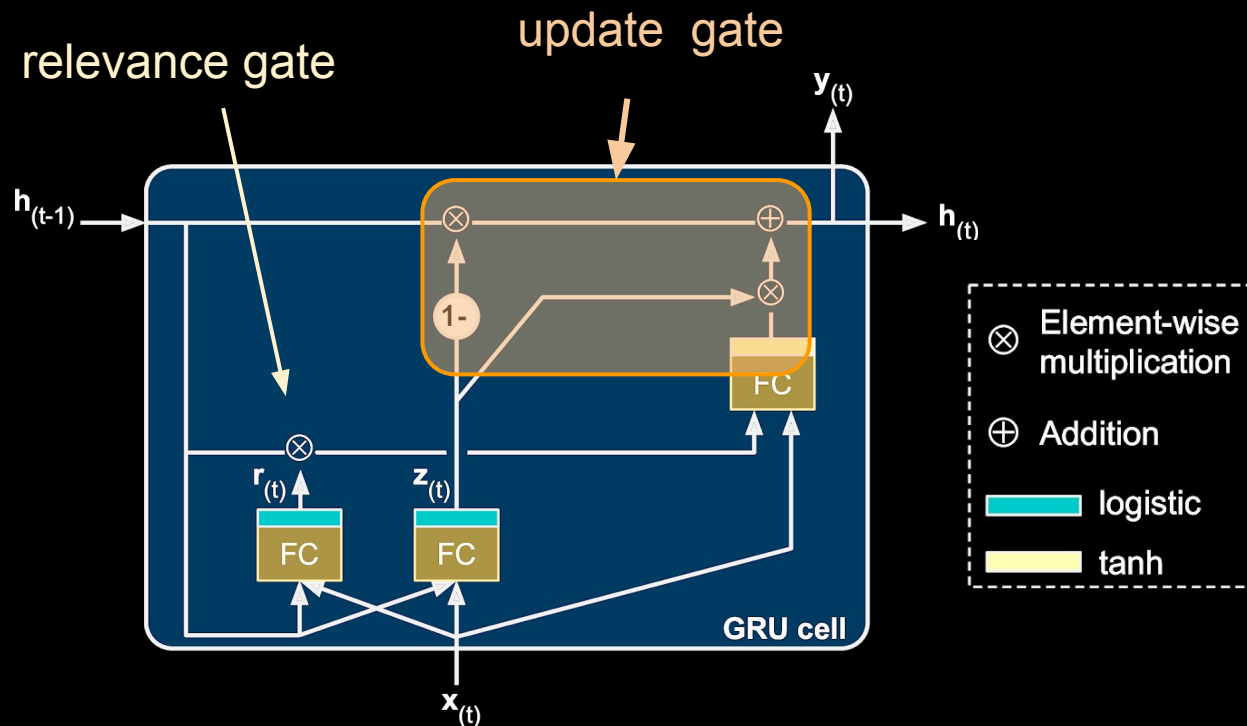
## Gated Recurrent Unit



(Geron, 2017)

# RNN: The GRU

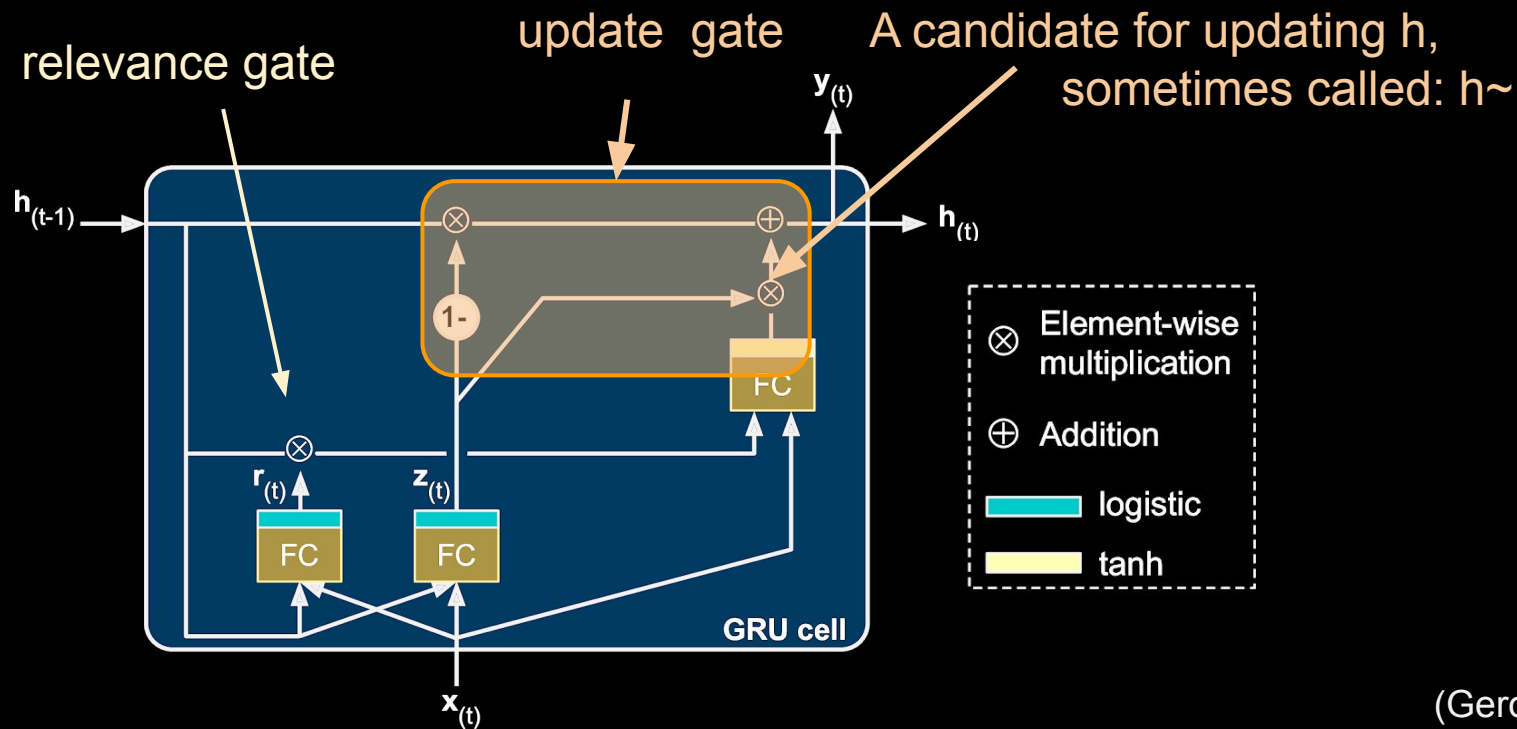
## Gated Recurrent Unit



(Geron, 2017)

# RNN: The GRU

## Gated Recurrent Unit



(Geron, 2017)



# RNN: The GRU

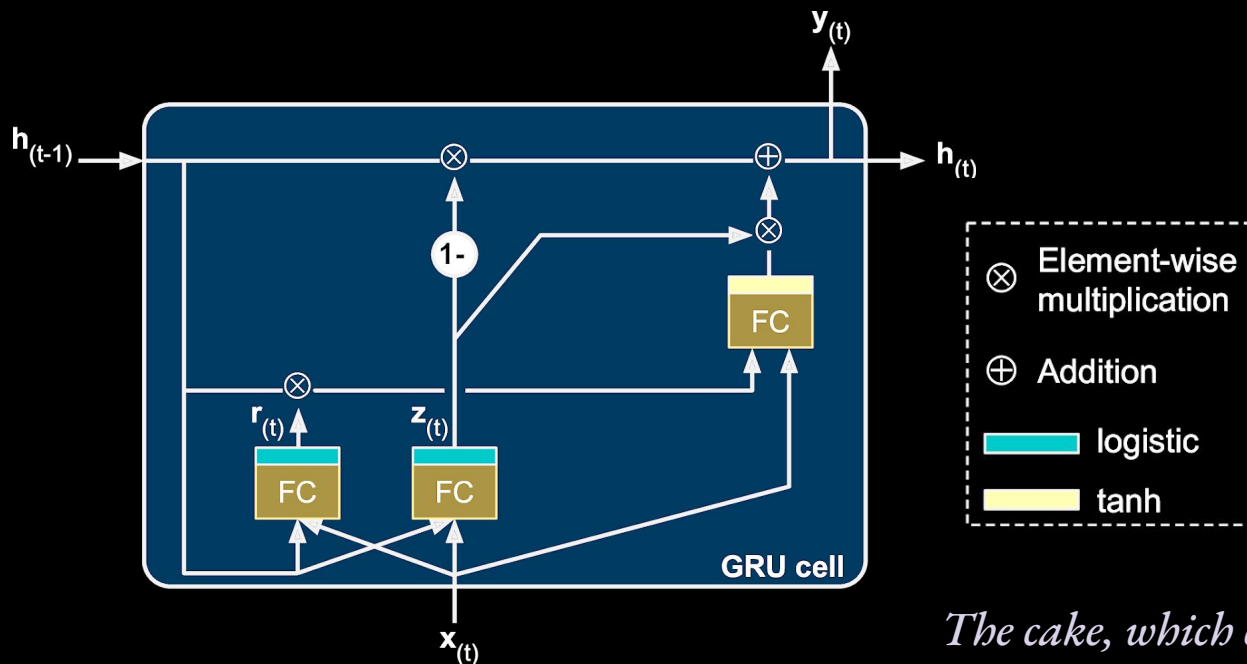
## Gated Recurrent Unit

$$\mathbf{z}_{(t)} = \sigma(\mathbf{W}_{xz}^T \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hz}^T \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_z)$$

$$\mathbf{r}_{(t)} = \sigma(\mathbf{W}_{xr}^T \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hr}^T \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_r)$$

$$\mathbf{g}_{(t)} = \tanh(\mathbf{W}_{xg}^T \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hg}^T \cdot (\mathbf{r}_{(t)} \otimes \mathbf{h}_{(t-1)}) + \mathbf{b}_g)$$

$$\mathbf{h}_{(t)} = \mathbf{z}_{(t)} \otimes \mathbf{h}_{(t-1)} + (1 - \mathbf{z}_{(t)}) \otimes \mathbf{g}_{(t)}$$



*The cake, which contained candles, was eaten.*

# What about the gradient?

$$\mathbf{z}_{(t)} = \sigma(\mathbf{W}_{xz}^T \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hz}^T \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_z)$$

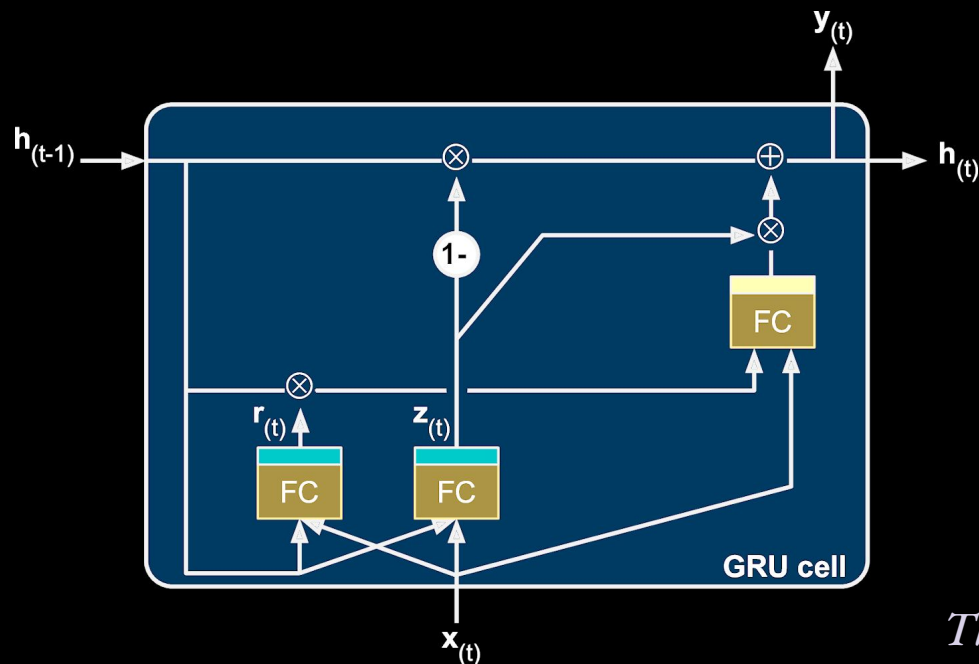
$$\mathbf{r}_{(t)} = \sigma(\mathbf{W}_{xr}^T \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hr}^T \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_r)$$

$$\mathbf{g}_{(t)} = \tanh(\mathbf{W}_{xg}^T \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hg}^T \cdot (\mathbf{r}_{(t)} \otimes \mathbf{h}_{(t-1)}) + \mathbf{b}_g)$$

$$\mathbf{h}_{(t)} = \mathbf{z}_{(t)} \otimes \mathbf{h}_{(t-1)} + (1 - \mathbf{z}_{(t)}) \otimes \mathbf{g}_{(t)}$$

The gates (i.e. multiplications based on a logistic) often end up keeping the hidden state exactly (or nearly exactly) as it was. Thus, for most dimensions of  $\mathbf{h}$ ,

$$\mathbf{h}_{(t)} \approx \mathbf{h}_{(t-1)}$$



*The cake, which contained candles, was eaten.*

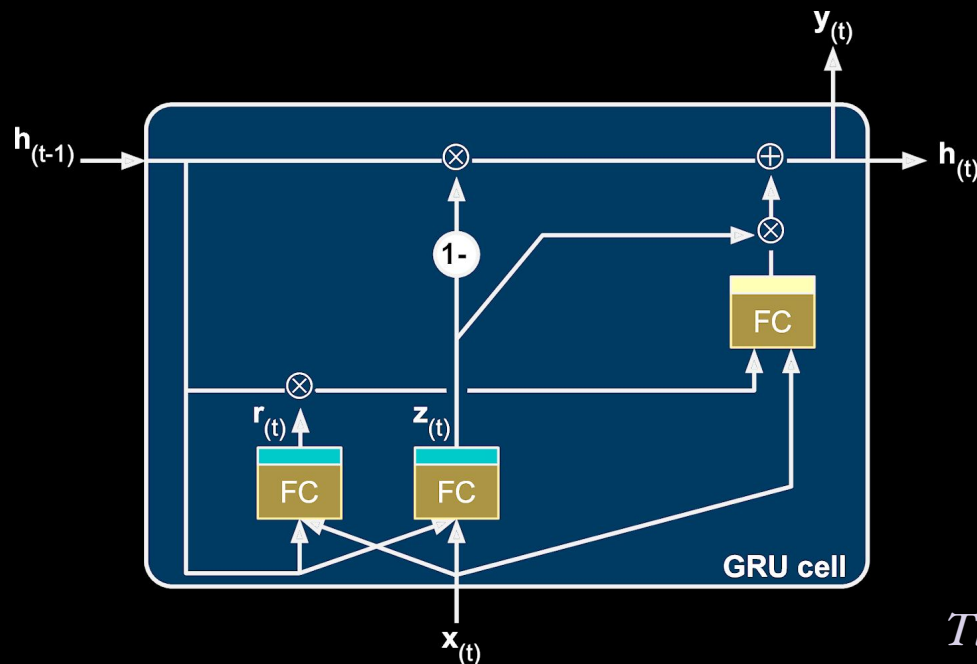
# What about the gradient?

$$\mathbf{z}_{(t)} = \sigma(\mathbf{W}_{xz}^T \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hz}^T \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_z)$$

$$\mathbf{r}_{(t)} = \sigma(\mathbf{W}_{xr}^T \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hr}^T \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_r)$$

$$\mathbf{g}_{(t)} = \tanh(\mathbf{W}_{xg}^T \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hg}^T \cdot (\mathbf{r}_{(t)} \otimes \mathbf{h}_{(t-1)}) + \mathbf{b}_g)$$

$$\mathbf{h}_{(t)} = \mathbf{z}_{(t)} \otimes \mathbf{h}_{(t-1)} + (1 - \mathbf{z}_{(t)}) \otimes \mathbf{g}_{(t)}$$



The gates (i.e. multiplications based on a logistic) often end up keeping the hidden state exactly (or nearly exactly) as it was. Thus, for most dimensions of  $\mathbf{h}$ ,

$$\mathbf{h}_{(t)} \approx \mathbf{h}_{(t-1)}$$

This tends to keep the gradient from vanishing since the same values will be present through multiple times in backpropagation through time. (The same idea applies to LSTMs but is easier to see here).

*The cake, which contained candles, was eaten.*

# The GRU (LSTM): Zoomed out

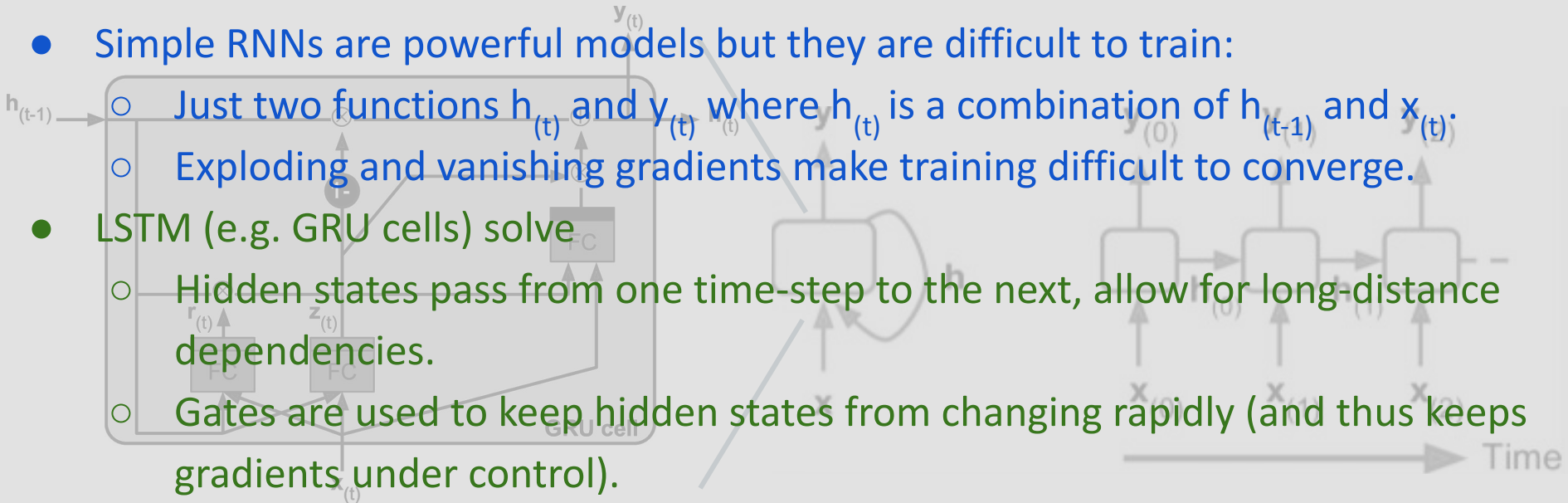
## Take-Aways

- Simple RNNs are powerful models but they are difficult to train:

- Just two functions  $h_{(t)}$  and  $y_{(t)}$  where  $h_{(t)}$  is a combination of  $h_{(t-1)}$  and  $x_{(t)}$ .
- Exploding and vanishing gradients make training difficult to converge.

- LSTM (e.g. GRU cells) solve

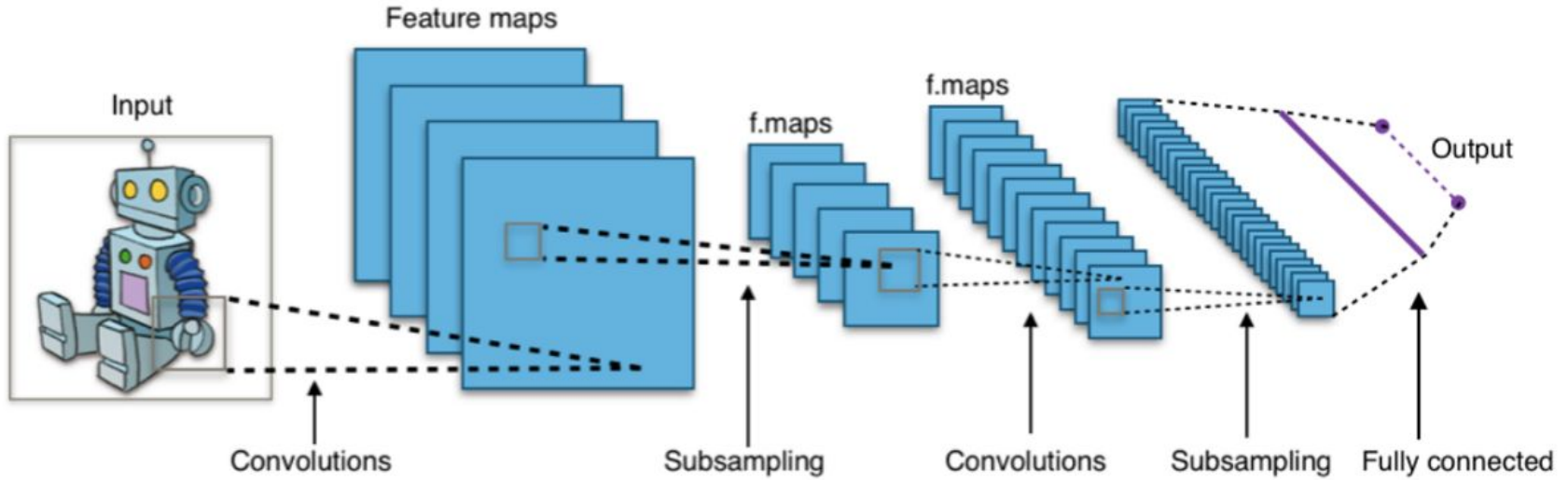
- Hidden states pass from one time-step to the next, allow for long-distance dependencies.
- Gates are used to keep hidden states from changing rapidly (and thus keeps gradients under control).
- To train: mini-batch stochastic gradient descent over cross-entropy cost



# Post-Exam2 Topics:

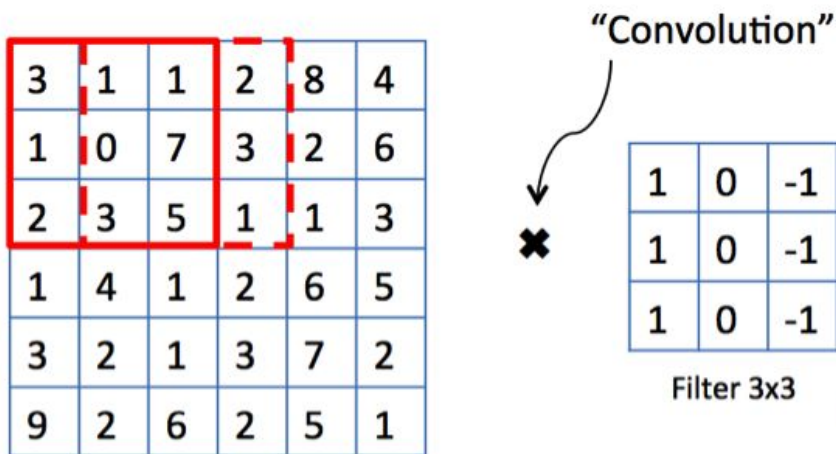
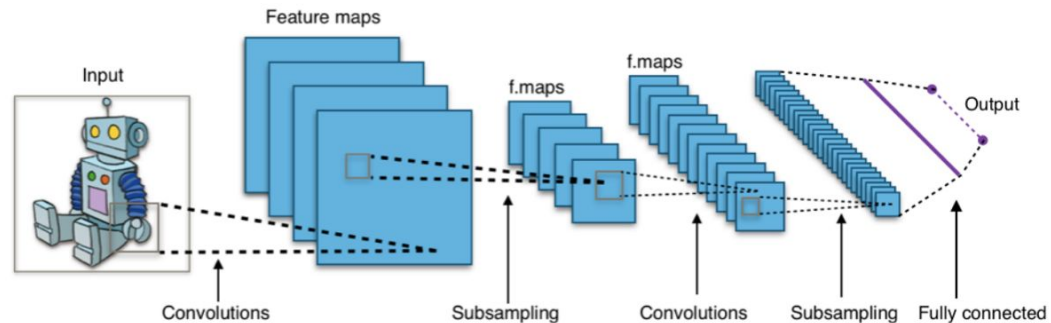
1. Research Ethics
2. Useful Plots
3. Machine Learning Cross Validation
4. Recurrent Neural Networks
- 5. Convolutional Neural Networks**
6. Transformer Networks

# Convolutional Neural Networks



(wikipedia)

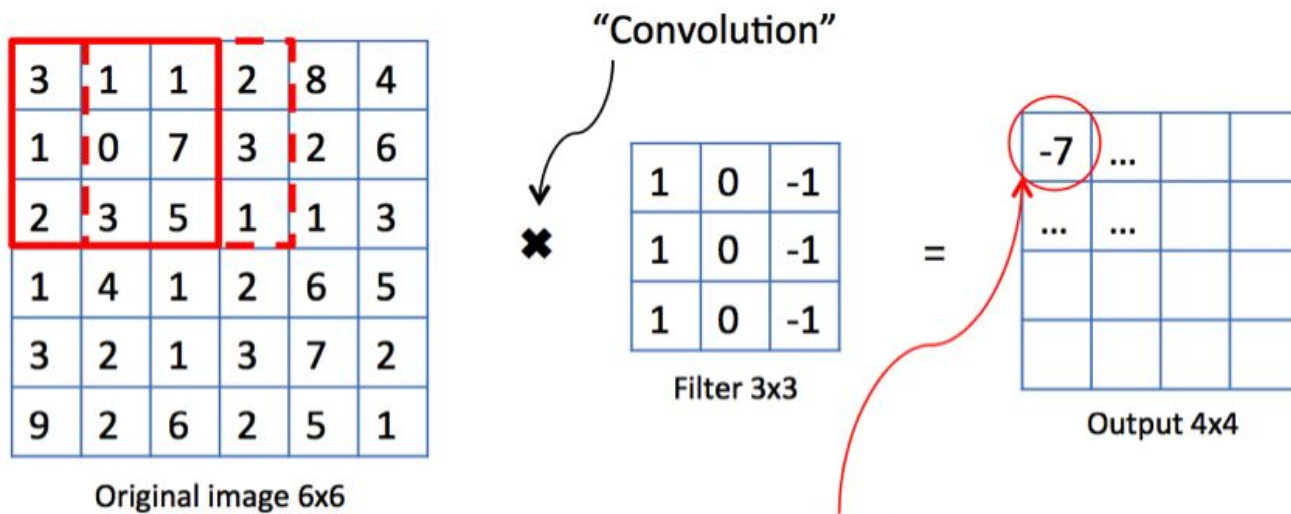
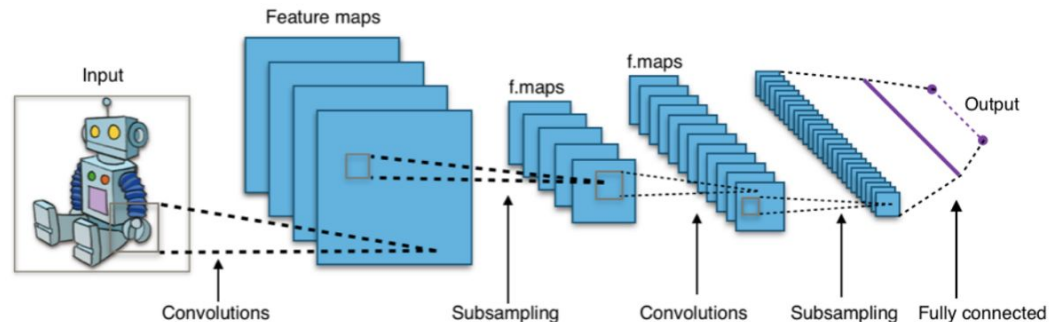
# Convolution Layer



Original image 6x6

(Barter, 2018)

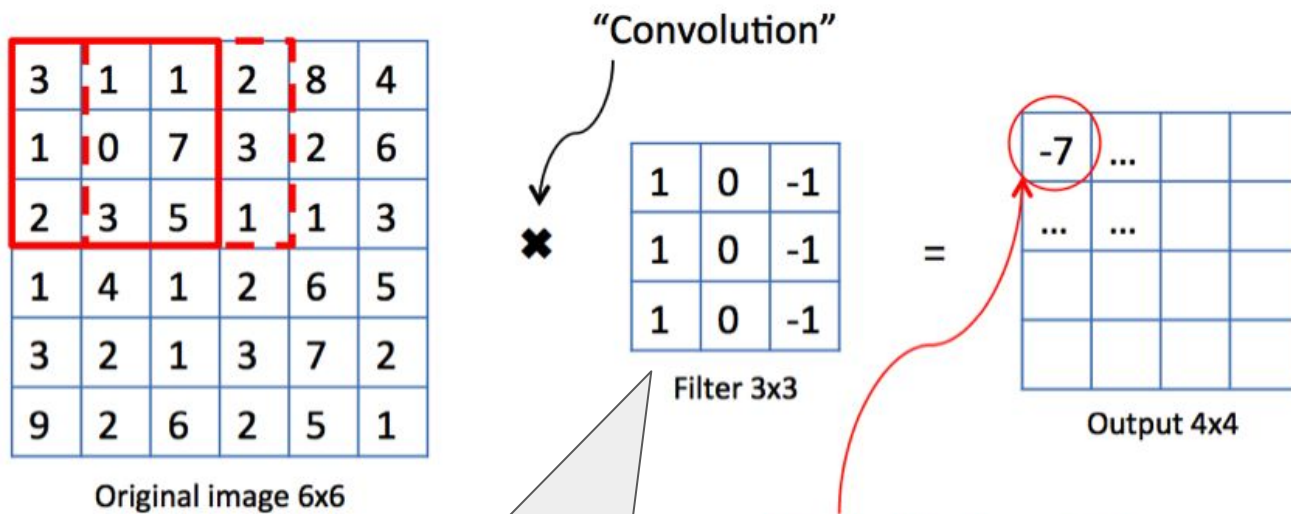
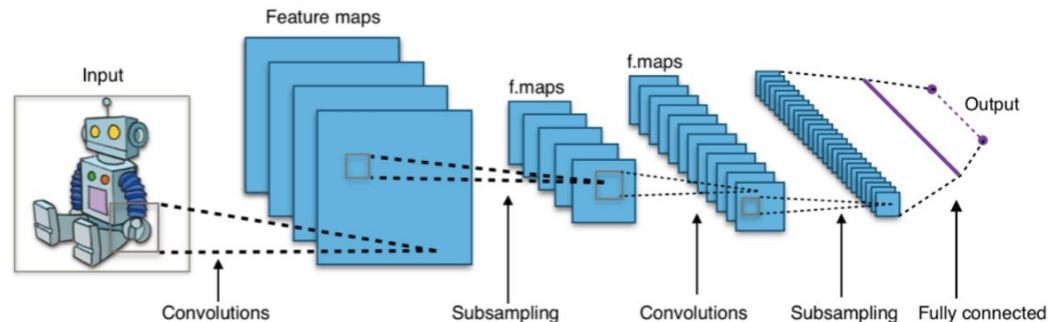
# Convolution Layer



(Barter, 2018)



# Convolution Layer

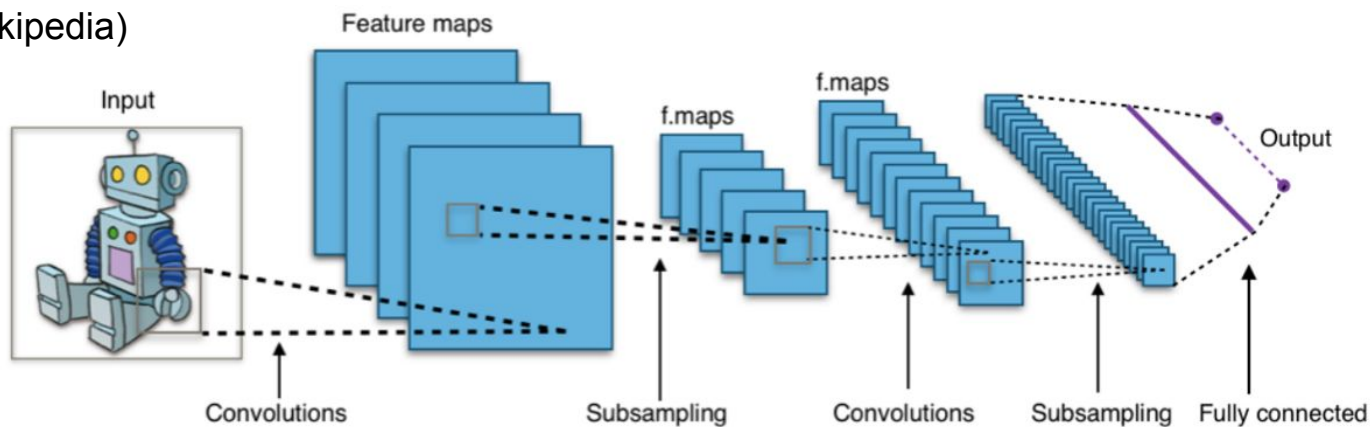


Breakthrough in image classification: Let the model automatically learn the filter weights!

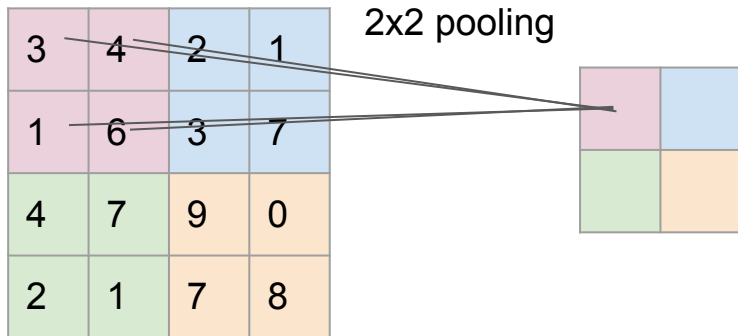
Result of the element-wise product and sum of the filter matrix and the original image

# Subsampling (Pooling)

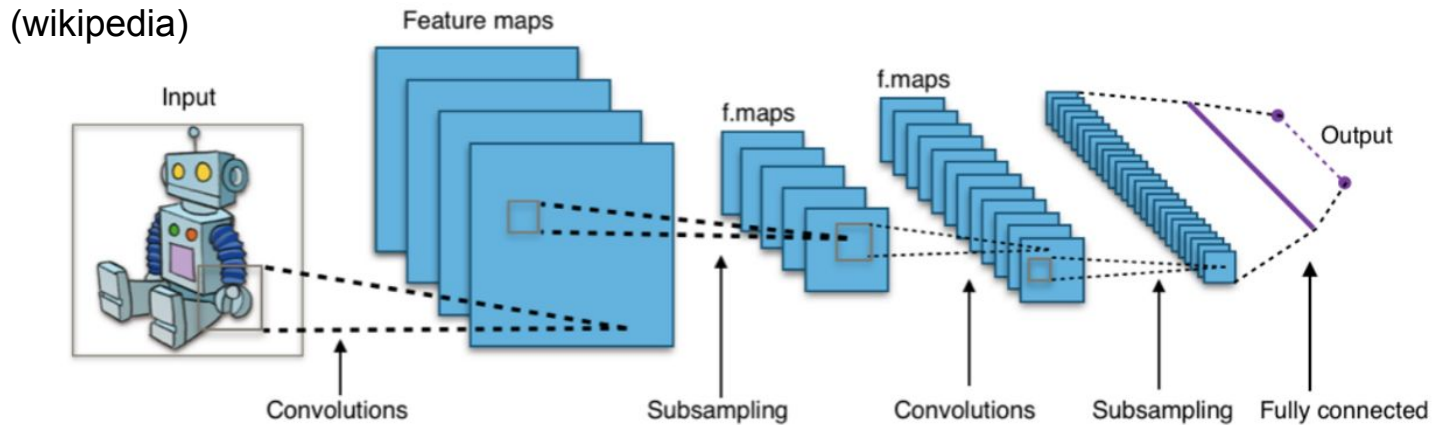
(wikipedia)



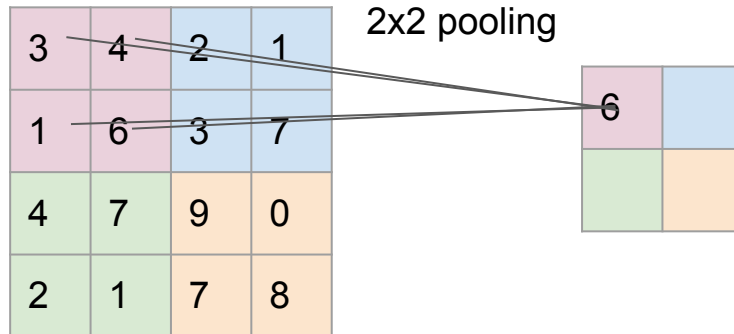
Subsampling -- reducing total grid size (i.e. reducing parameters for next layer)



# Subsampling (Pooling)



Subsampling -- reducing total grid size (i.e. reducing parameters for next layer)

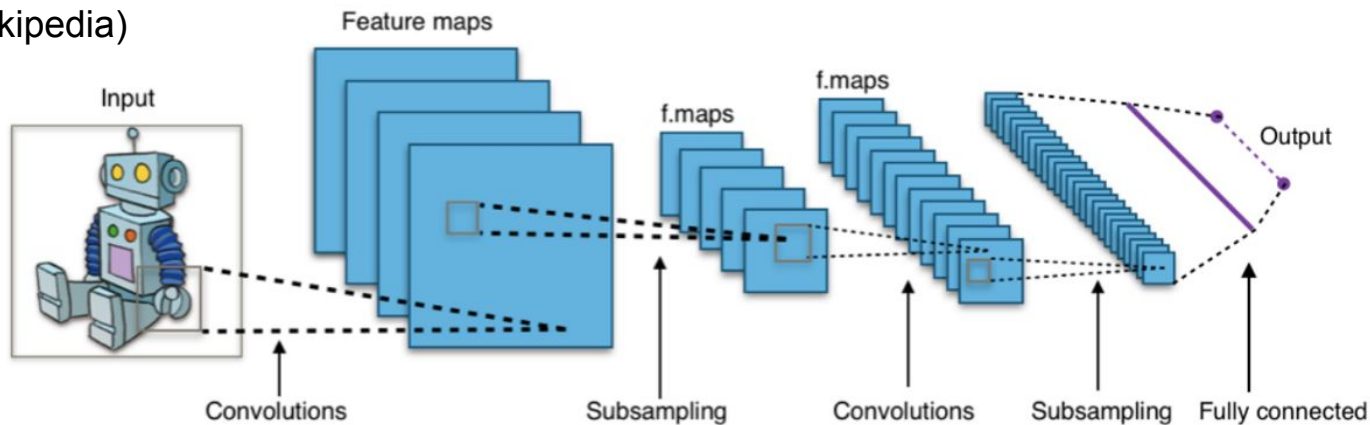


Types of pooling

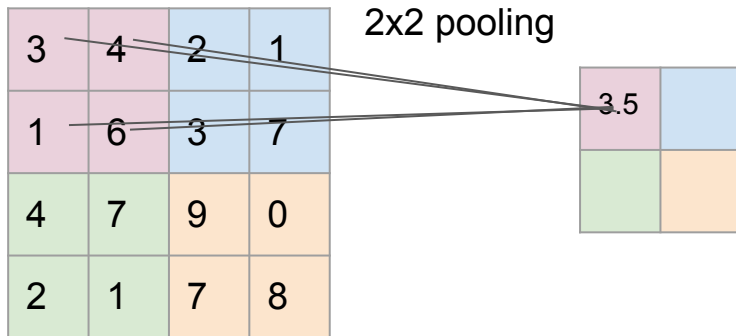
- max
- avg

# Subsampling (Pooling)

(wikipedia)



Subsampling -- reducing total grid size (i.e. reducing parameters for next layer)



Types of pooling

- max
- **avg**

# Standard Training Loss Function

```
RNN_cost = tf.reduce_mean(-tf.reduce_sum(y*tf.log(y_pred))
```

```
#where did this come from?
```

Logistic Regression Likelihood:  $L(\beta_0, \beta_1, \dots, \beta_k | X, Y) = \prod_{i=1}^n p(x_i)^{y_i} (1 - p(x_i))^{1-y_i}$

Final Cost Function:  $J^{(t)} = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^{|V|} y_{i,j}^{(t)} \log \hat{y}_{i,j}^{(t)}$  -- "cross entropy error"

# Standard Training Loss Function

```
RNN_cost = tf.reduce_mean(-tf.reduce_sum(y*tf.log(y_pred))
```

```
#where did this come from?
```

Logistic Regression Likelihood:  $L(\beta_0, \beta_1, \dots, \beta_k | X, Y) = \prod_{i=1}^n p(x_i)^{y_i} (1 - p(x_i))^{1-y_i}$

Log Likelihood:  $\ell(\beta) = \sum_{i=1}^N y_i \log p(x_i) + (1 - y_i) \log (1 - p(x_i))$

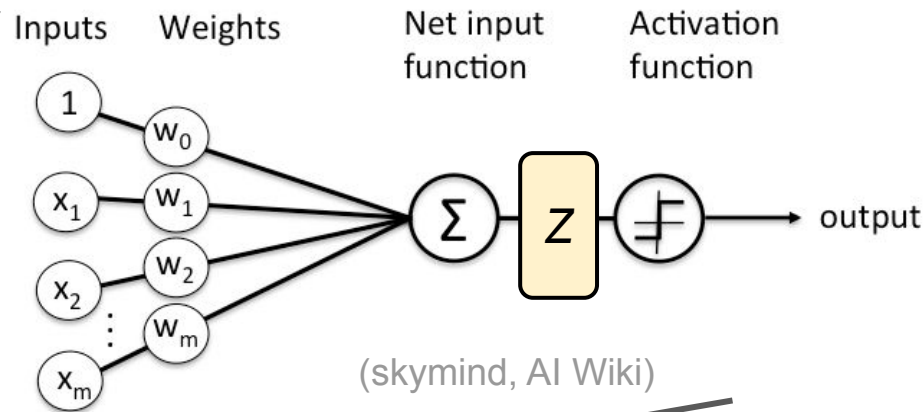
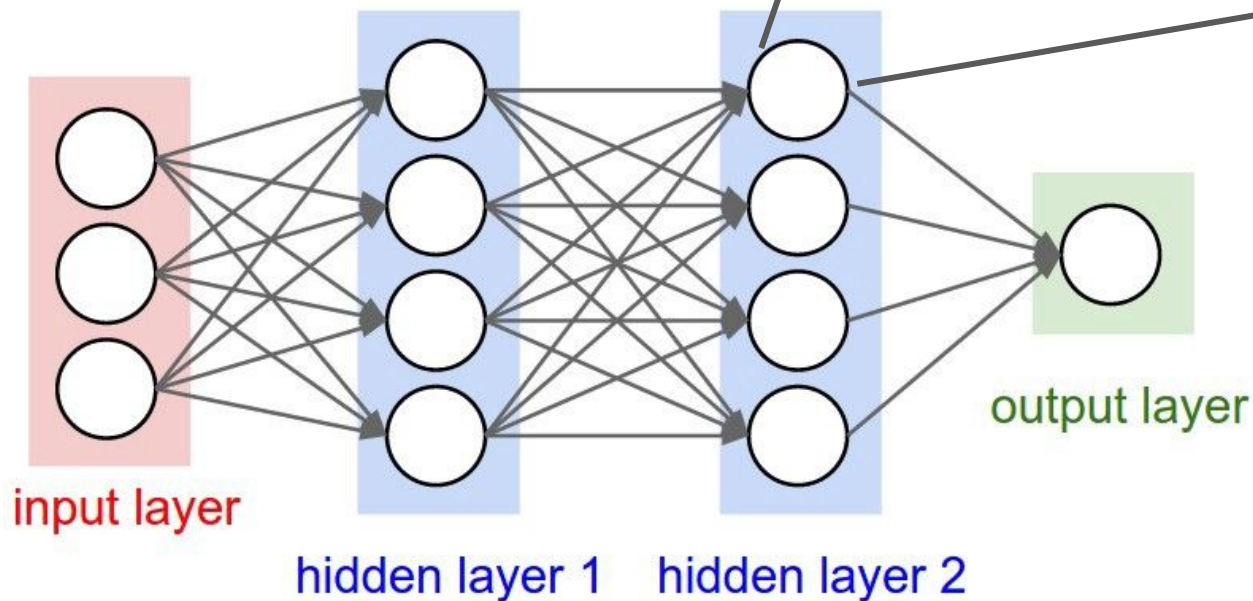
Log Loss:  $J(\beta) = -\frac{1}{N} \sum_{i=1}^N y_i \log p(x_i) + (1 - y_i) \log (1 - p)(x_i)$

Cross-Entropy Cost:  $J = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^{|V|} y_{i,j} \log p(x_{i,j})$  (a "multiclass" log loss)

Final Cost Function:  $J^{(t)} = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^{|V|} y_{i,j}^{(t)} \log \hat{y}_{i,j}^{(t)}$  -- "cross entropy error"

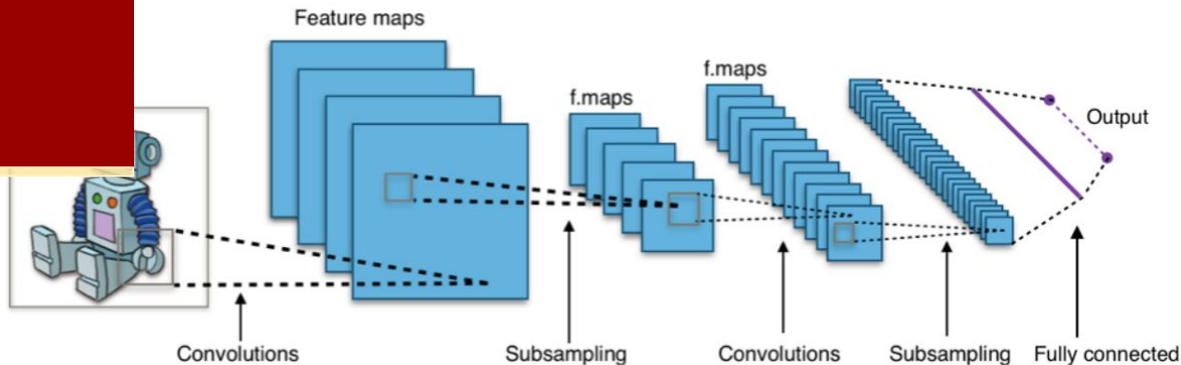
# Review

## Feed Forward Network (full-connected)



# Review

## Convolutional NN



3	1	1	2	8	4
1	0	7	3	2	6
2	3	5	1	1	3
1	4	1	2	6	5
3	2	1	3	7	2
9	2	6	2	5	1

Original image 6x6

“Convolution”

$\otimes$

1	0	-1
1	0	-1
1	0	-1

Filter 3x3

=

-7	...		
...	...		

Output 4x4

Result of the element-wise product and sum of the filter matrix and the original image

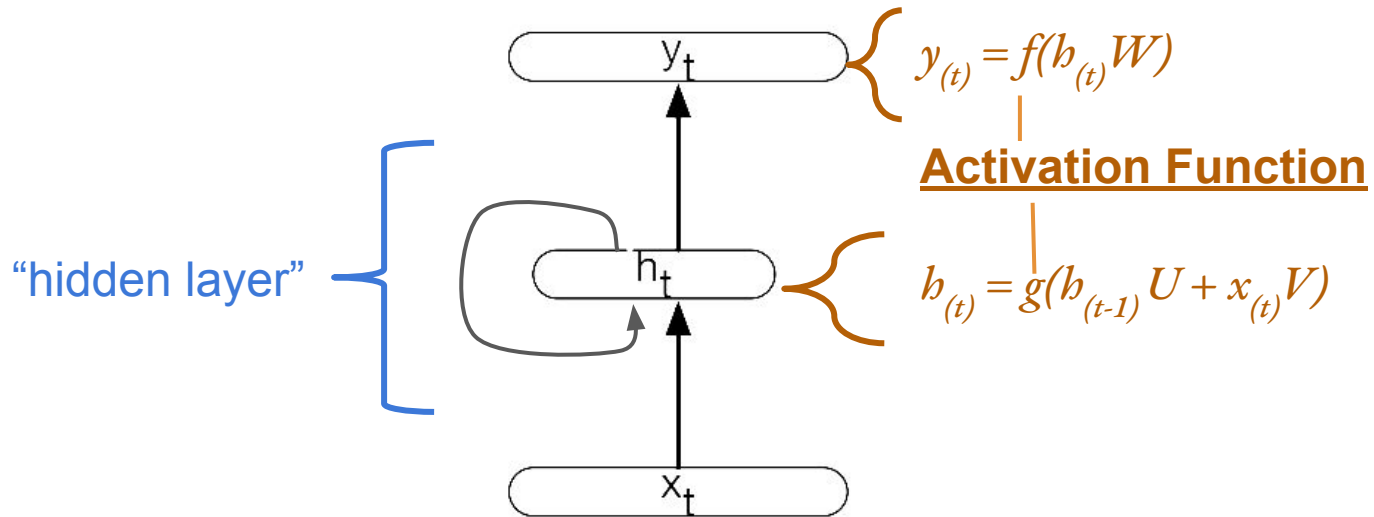


# Post-Exam2 Topics:

1. Research Ethics
2. Useful Plots
3. Machine Learning Cross Validation
4. Recurrent Neural Networks
5. Convolutional Neural Networks
6. **Transformer Networks**

# Review

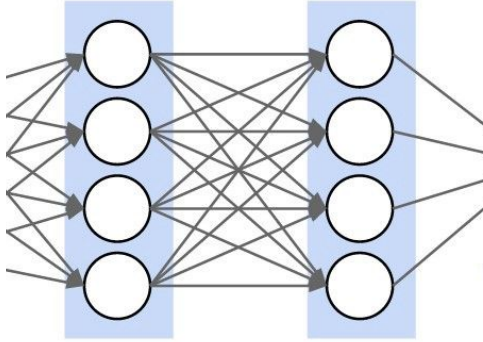
## Recurrent Neural Network



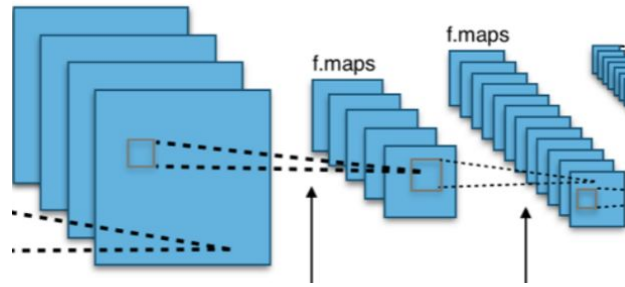
**Figure 9.2** Simple recurrent neural network after Elman (Elman, 1990). The hidden layer includes a recurrent connection as part of its input. That is, the activation value of the hidden layer depends on the current input as well as the activation value of the hidden layer from the previous timestep.

(Jurafsky, 2019)

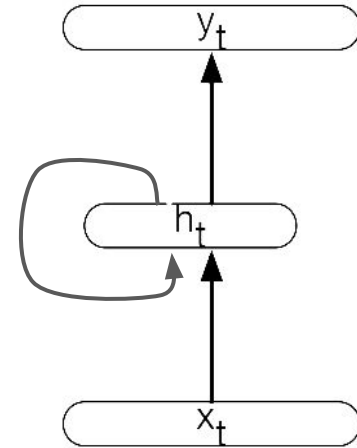
# FFN



# CNN



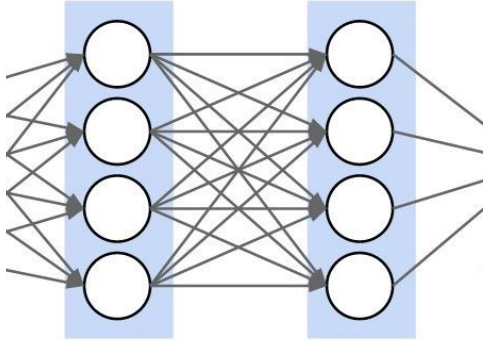
# RNN



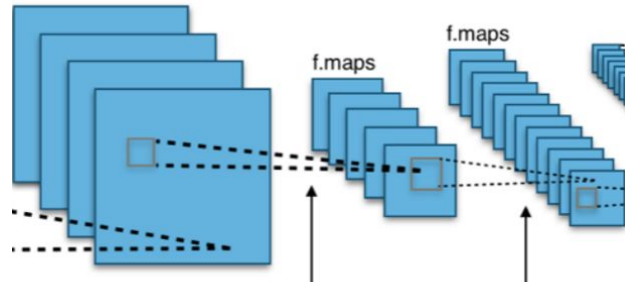
*Can model computation (e.g. matrix operations for a single input) be parallelized?*



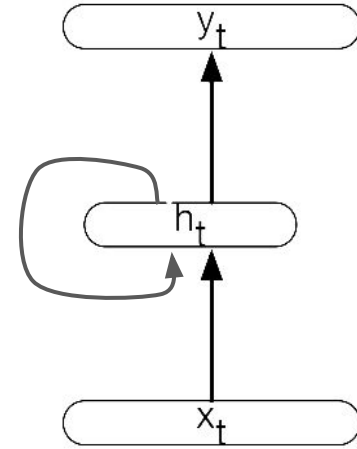
# FFN



# CNN



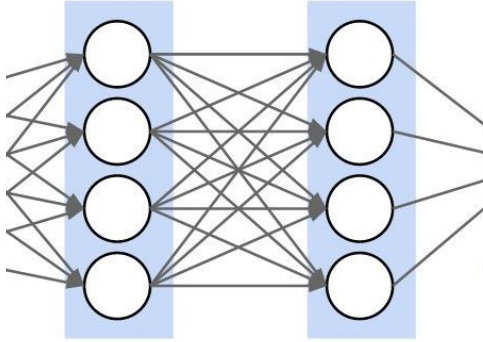
# RNN



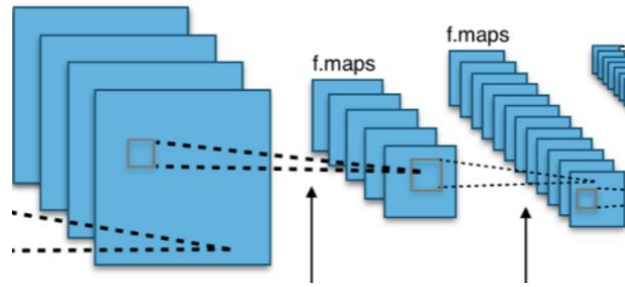
*Can model computation (e.g. matrix operations for a single input) be parallelized?*



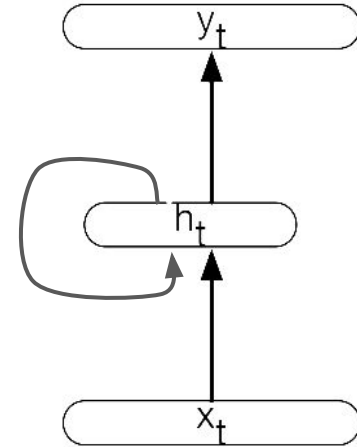
# FFN



# CNN



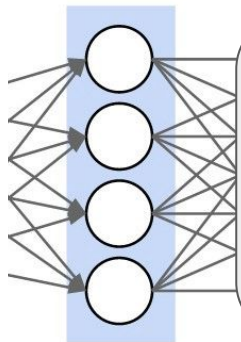
# RNN



*Can model computation (e.g. matrix operations for a single input) be parallelized?*



## FFN

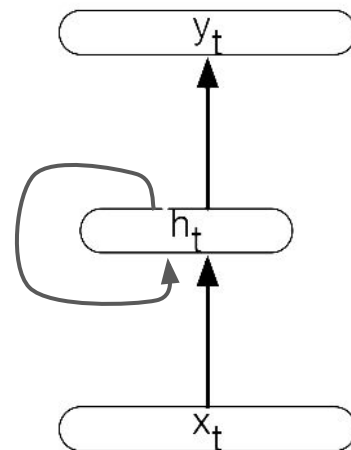


## CNN

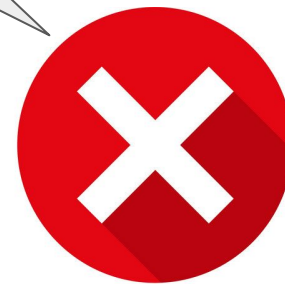


*Ultimately limits how complex the model can be (i.e. it's total number of paramers/weights) as compared to a CNN.*

## RNN



*Can model computation (e.g. matrix operations for a single input) be parallelized?*



# The Transformer: Attention-only Models

Can handle sequences and long-distance dependencies, but....

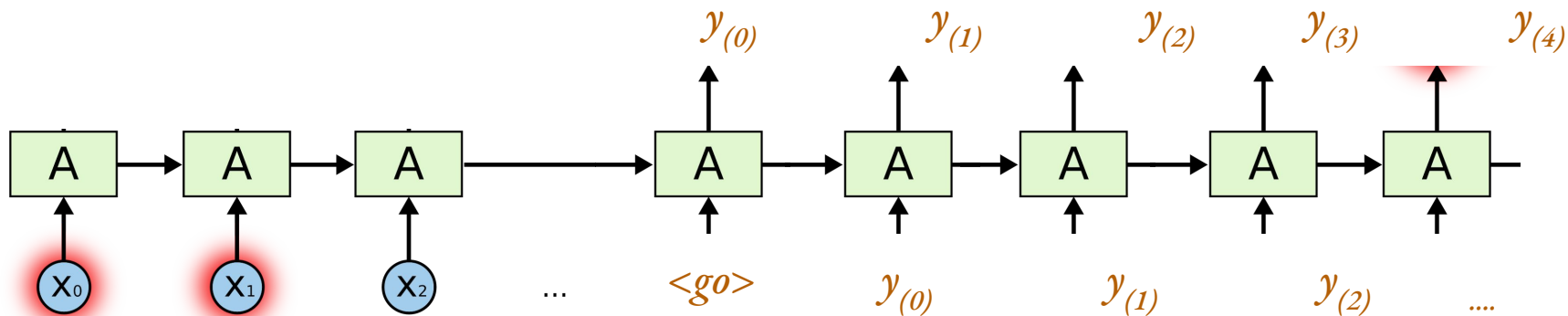
- Don't want complexity of LSTM/GRU cells
- Constant num edges between input steps
- Enables “interactions” (i.e. adaptations) between words
- **Easy to parallelize -- don't need sequential processing.**

# The Transformer: Attention-only Models

Challenge:

*The ball was kicked by kayla.*

- Long distance dependency when translating:



*Kayla kicked the ball.*

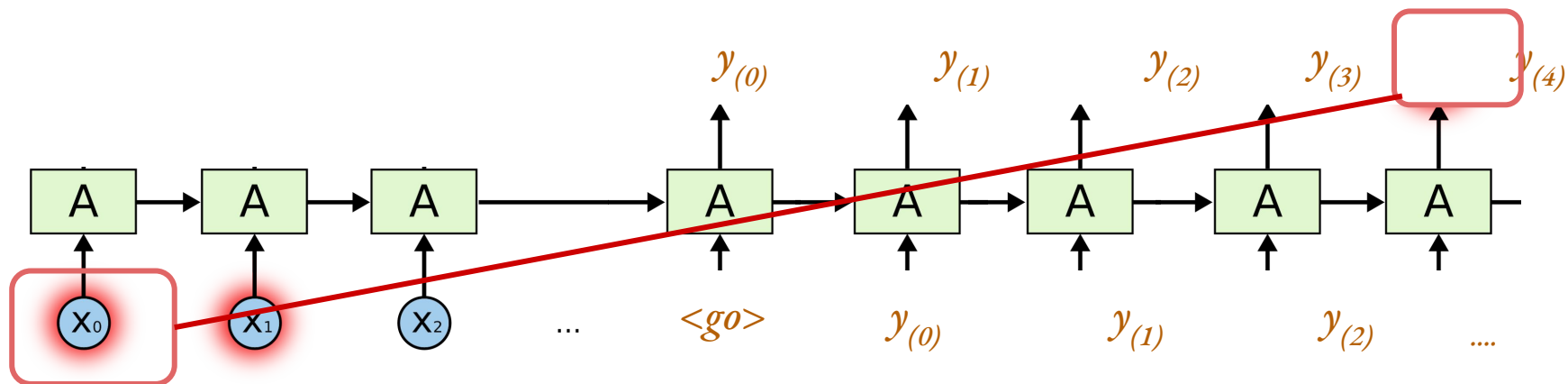


# The Transformer: Attention-only Models

Challenge:

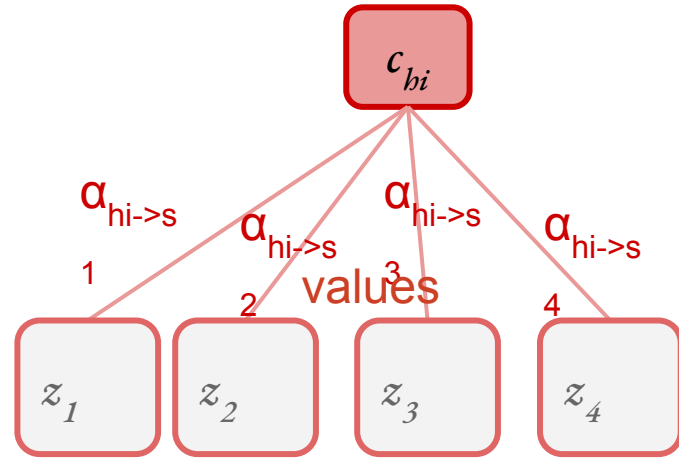
*The ball was kicked by kayla.*

- Long distance dependency when translating:

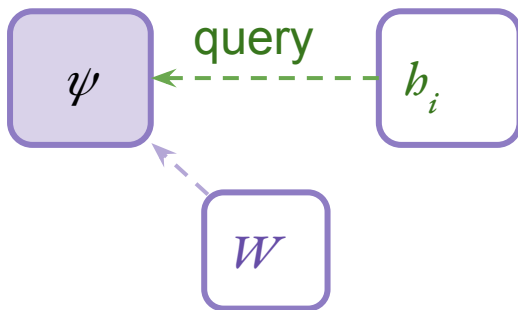
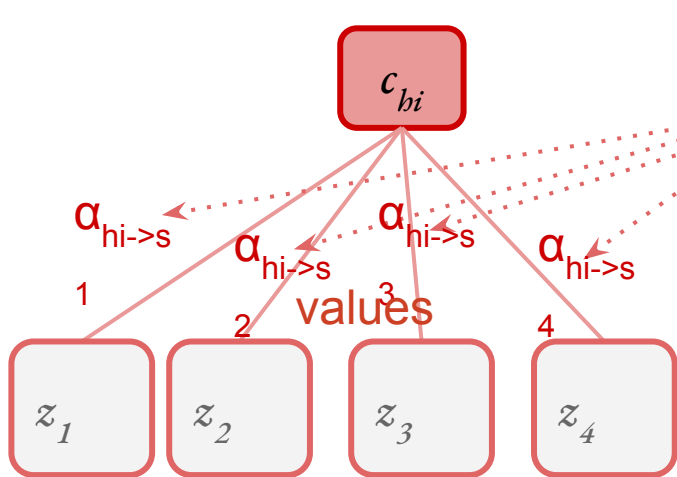


*Kayla kicked the ball.*

# Attention



# Attention

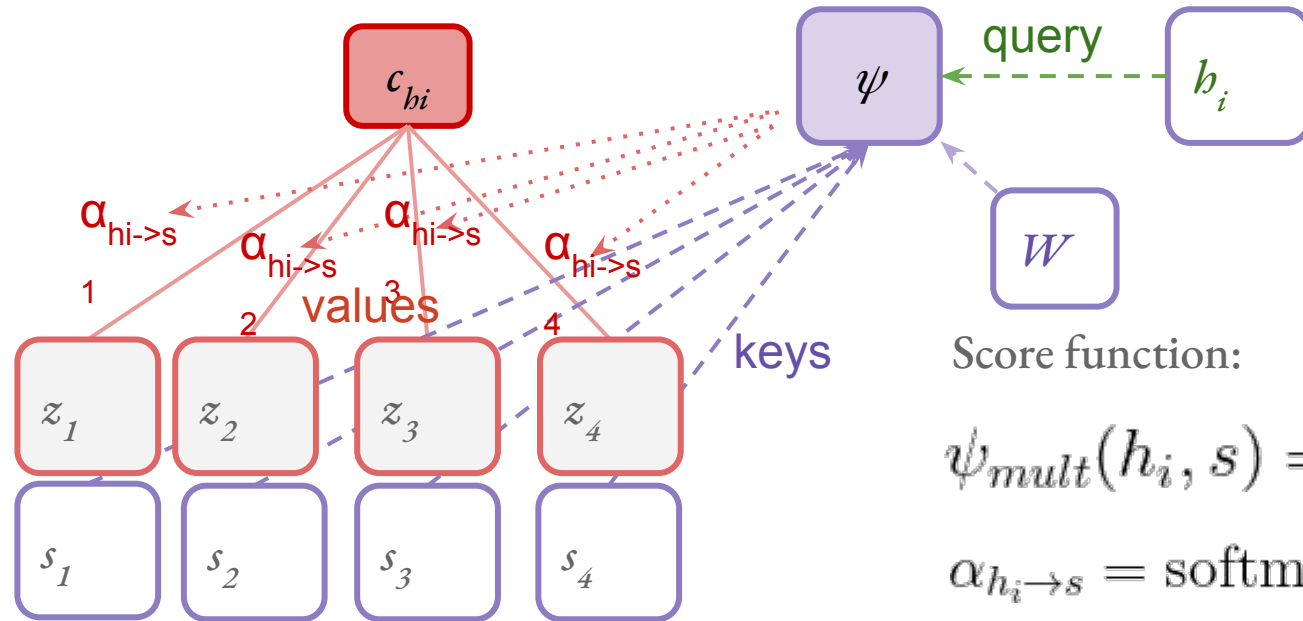


Score function:

$$\psi_{mult}(h_i, s) = s^T W h_i$$

$$\alpha_{h_i \rightarrow s} = \text{softmax}(\psi(h_i, s))$$

# Attention



Score function:

$$\psi_{mult}(h_i, s) = s^T W h_i$$

$$\alpha_{h_i \rightarrow s} = \text{softmax}(\psi(h_i, s))$$

$$c_{h_i} = \sum_{n=1}^{|s|} \alpha_{h_i \rightarrow s_n} z_n$$

# The Transformer: Attention-only Models

attention-only models

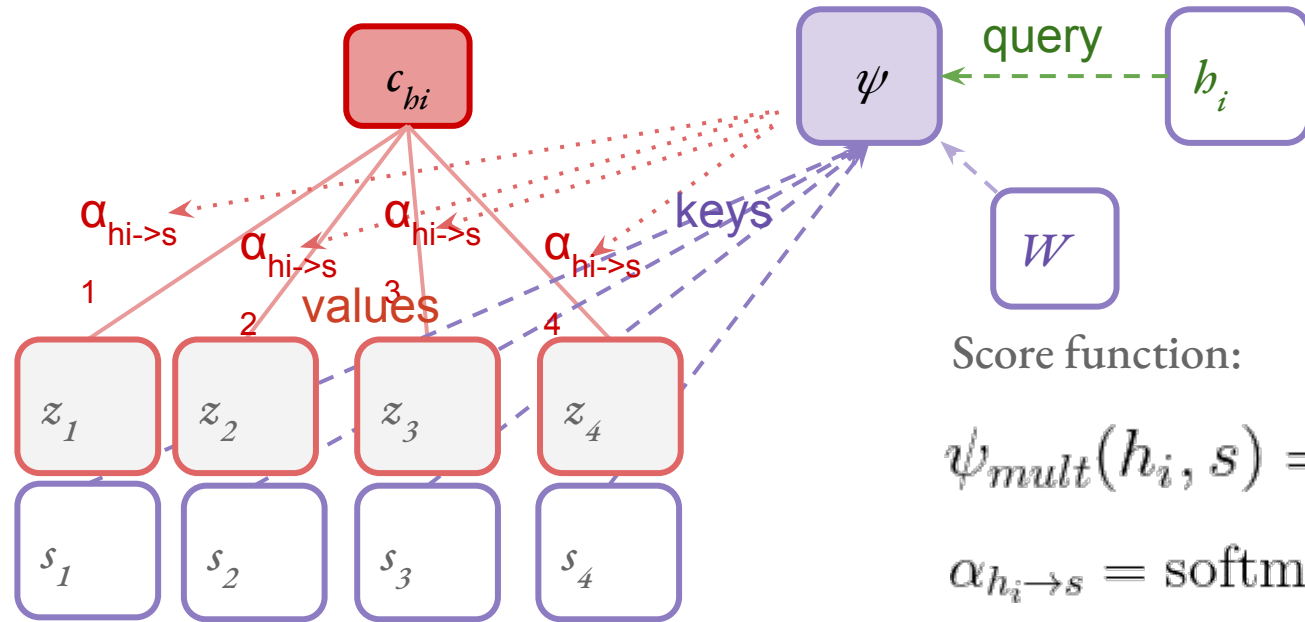
Challenge:

- Long distance dependency when translating:

Attention came about for encoder decoder models.

Then self-attention was introduced:

# Attention



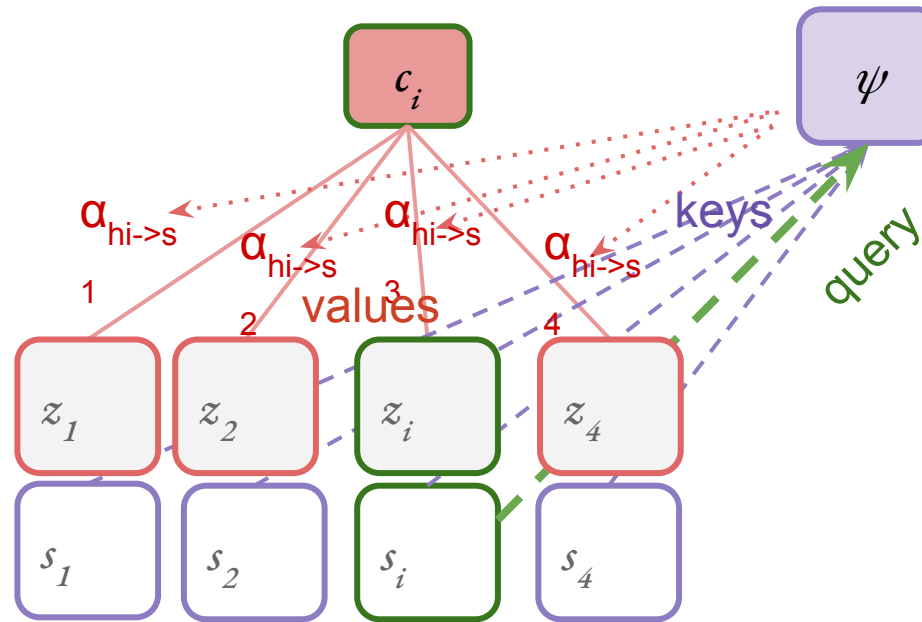
Score function:

$$\psi_{mult}(h_i, s) = s^T W h_i$$

$$\alpha_{h_i \rightarrow s} = \text{softmax}(\psi(h_i, s))$$

$$c_{h_i} = \sum_{n=1}^{|s|} \alpha_{h_i \rightarrow s_n} z_n$$

# Attention



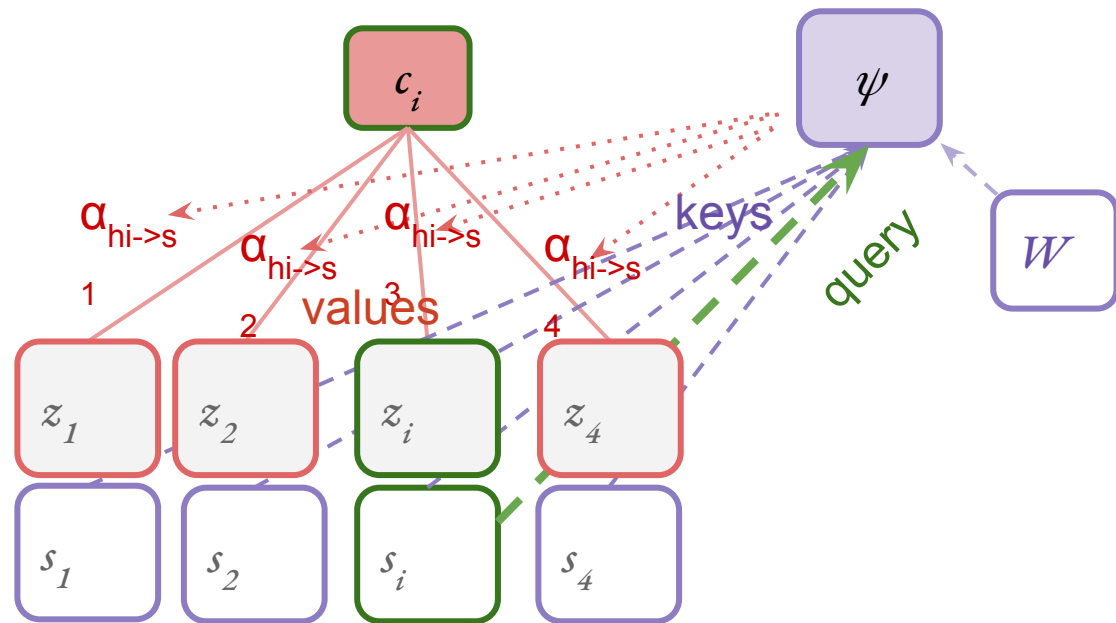
Score function:

$$\psi_{mult}(h_i, s) = s^T W h_i$$

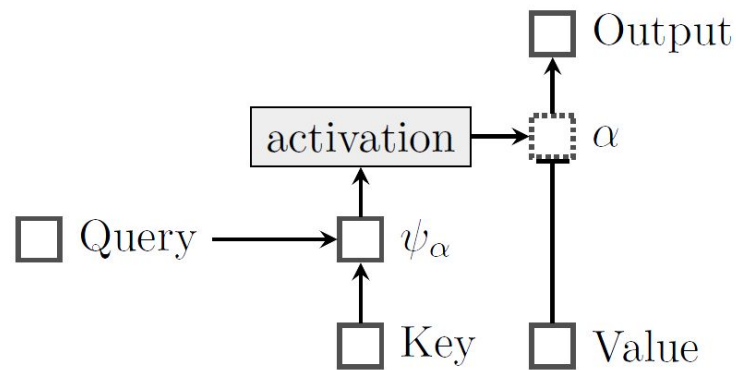
$$\alpha_{h_i \rightarrow s} = \text{softmax}(\psi(h_i, s))$$

$$c_{h_i} = \sum_{n=1}^{|s|} \alpha_{h_i \rightarrow s_n} z_n$$

# Attention



Attention as weighting a value based on a query and key:



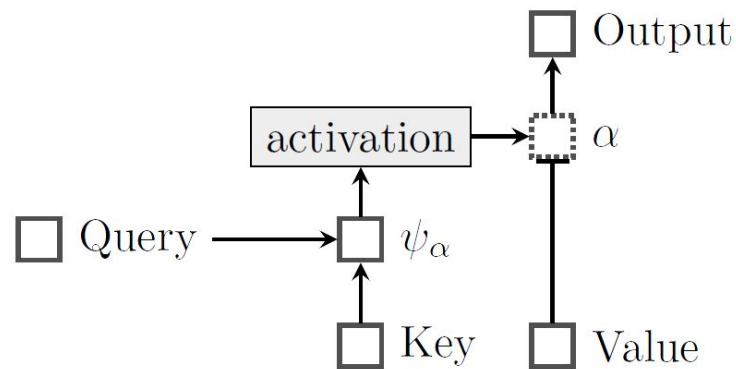
(Eisenstein, 2018)



# The Transformer: Attention-only Models

Attention-only Models

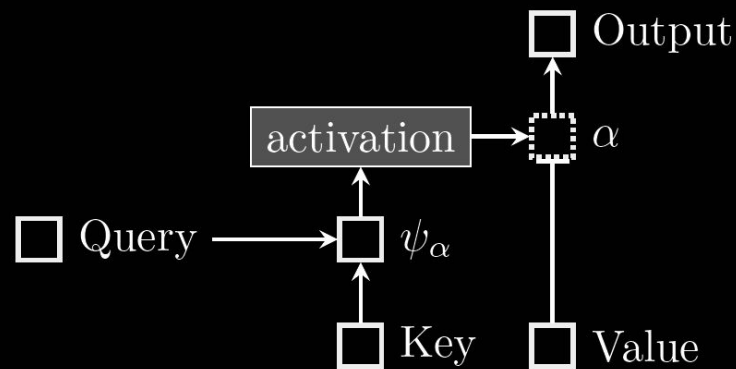
Attention as weighting a value based on a query and key:



(Eisenstein, 2018)

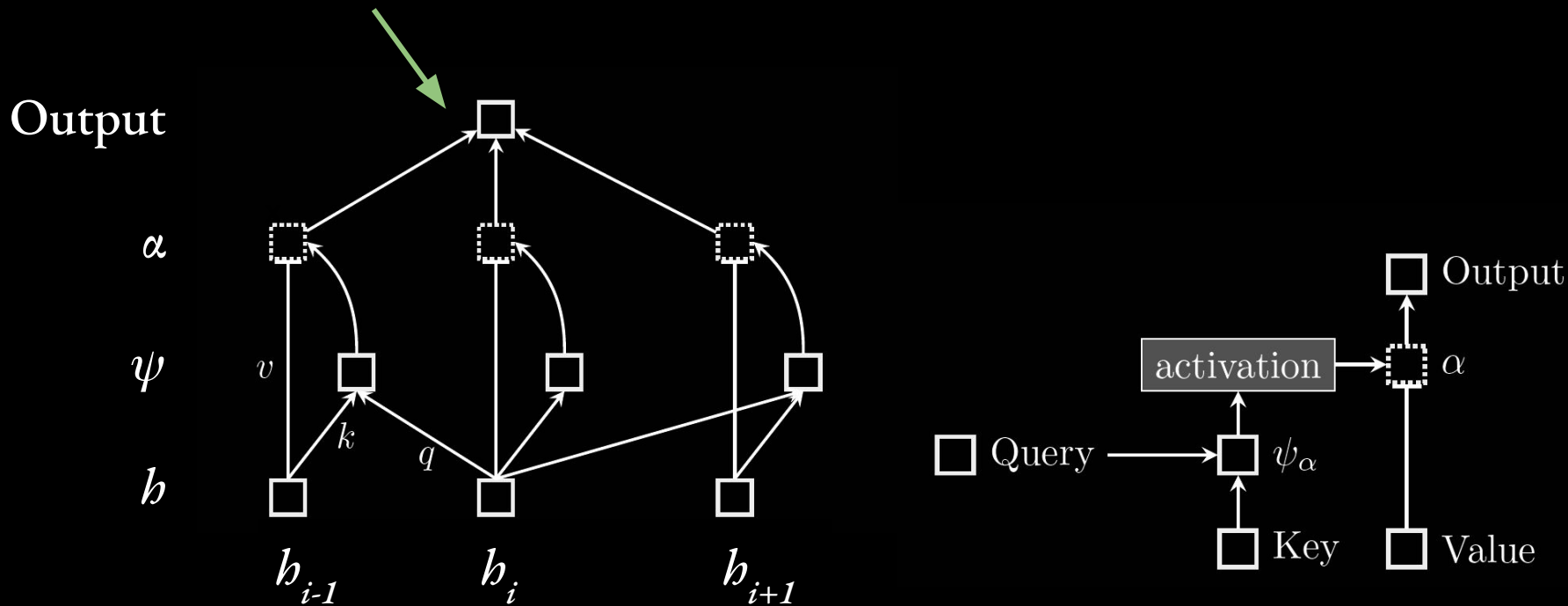
# The Transformer: Attention-only Models

CSE 545 Supplemental Lecture  
Will begin at 2:00pm



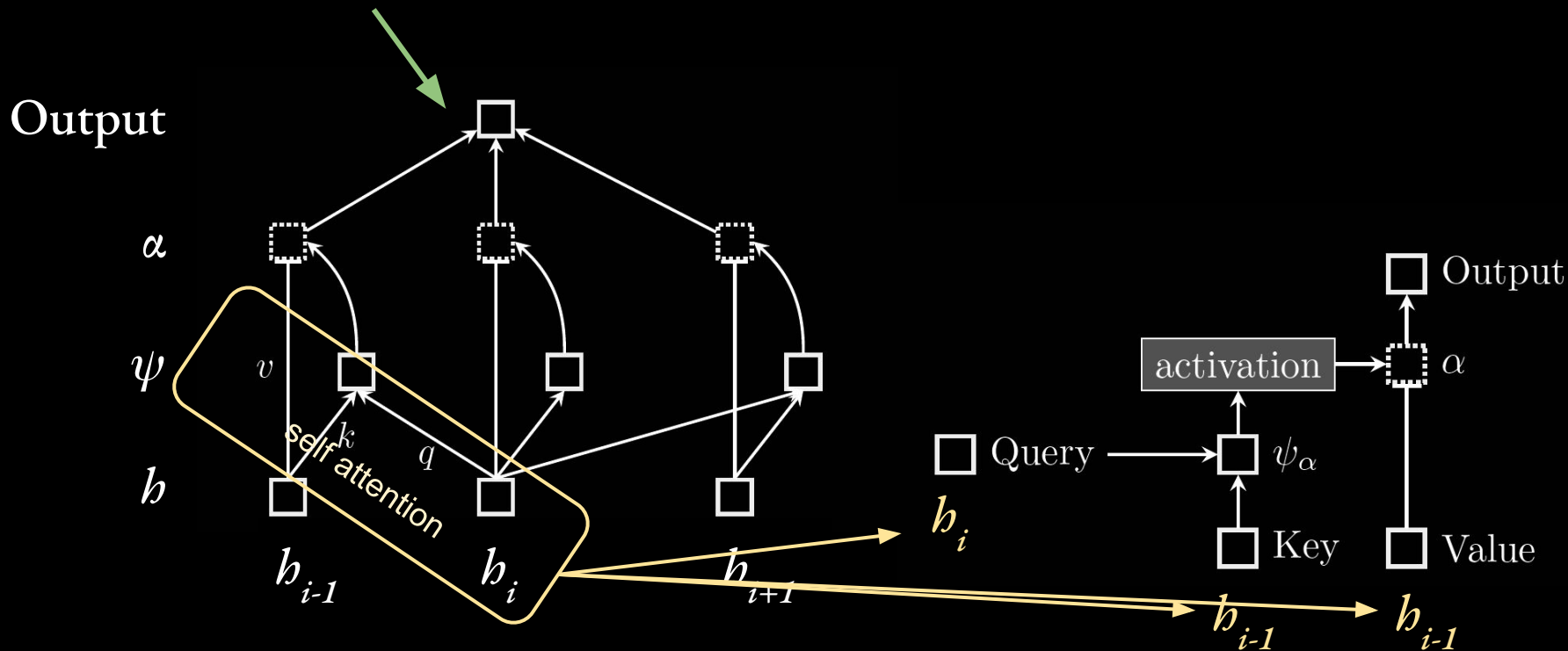
(Eisenstein, 2018)

# The Transformer: Attention-only Models



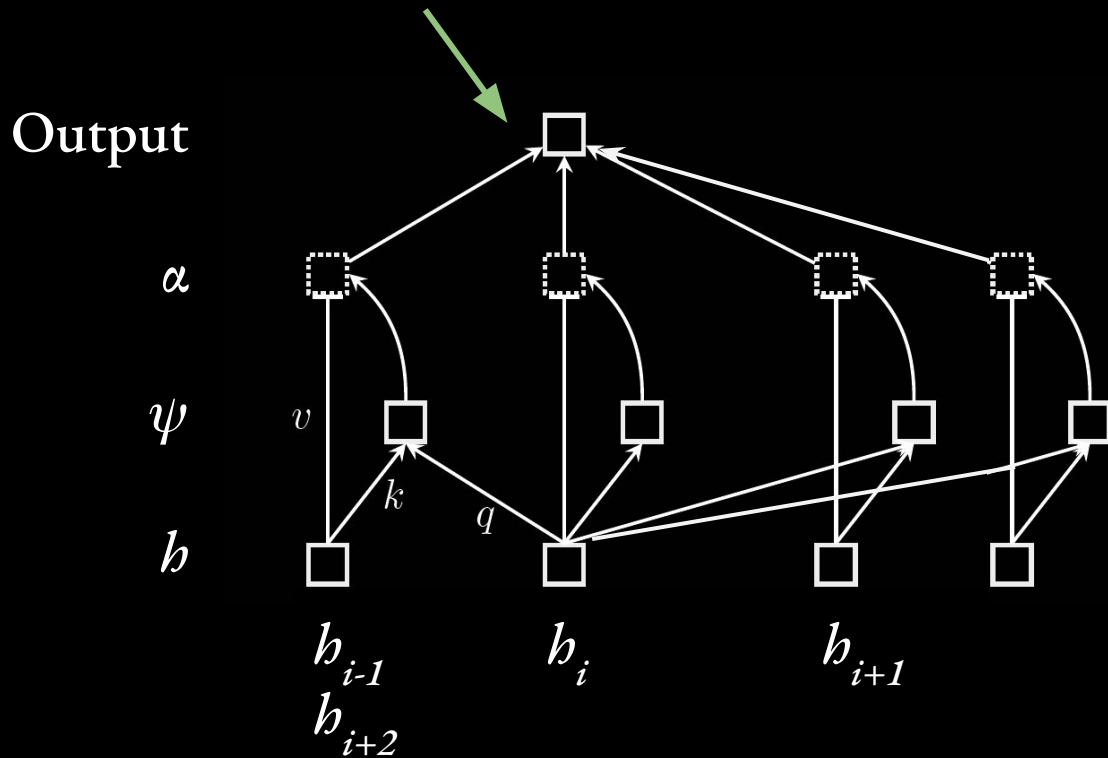
(Eisenstein, 2018)

# The Transformer: Attention-only Models

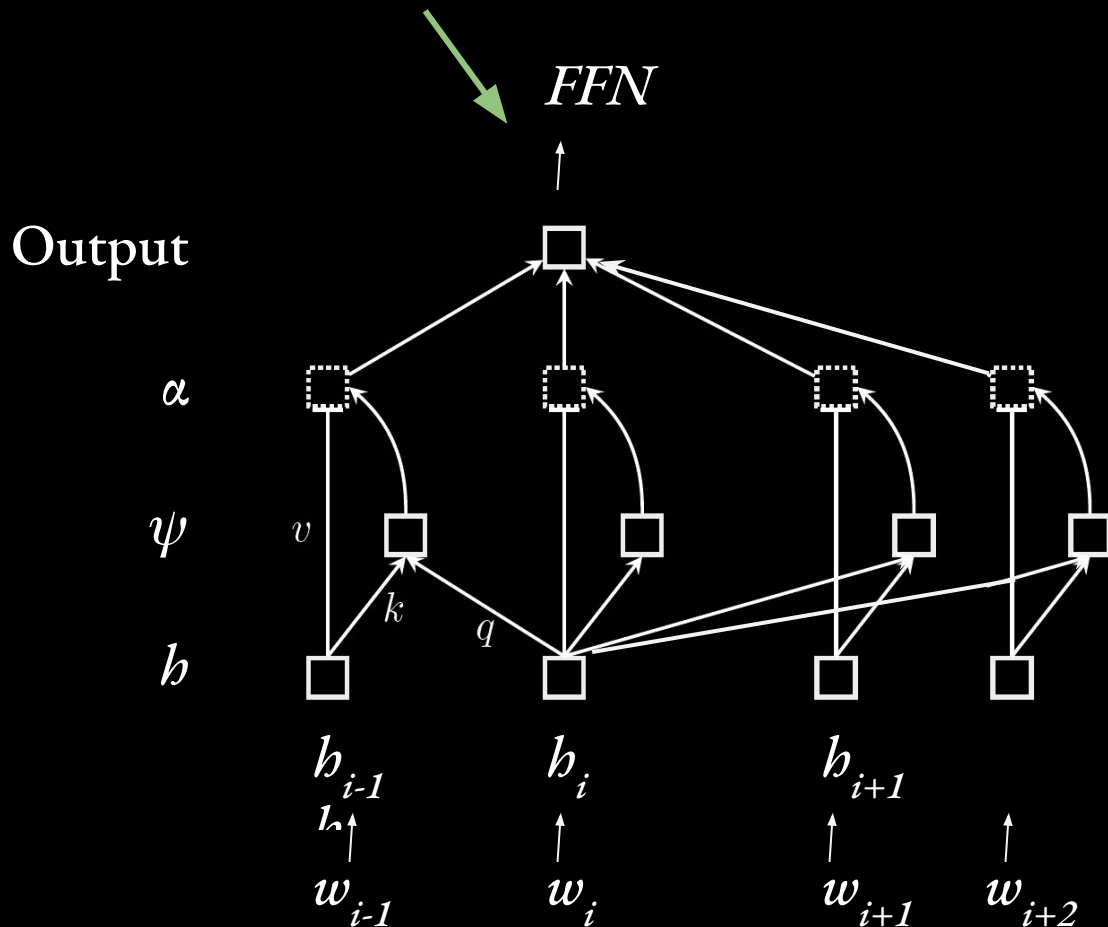


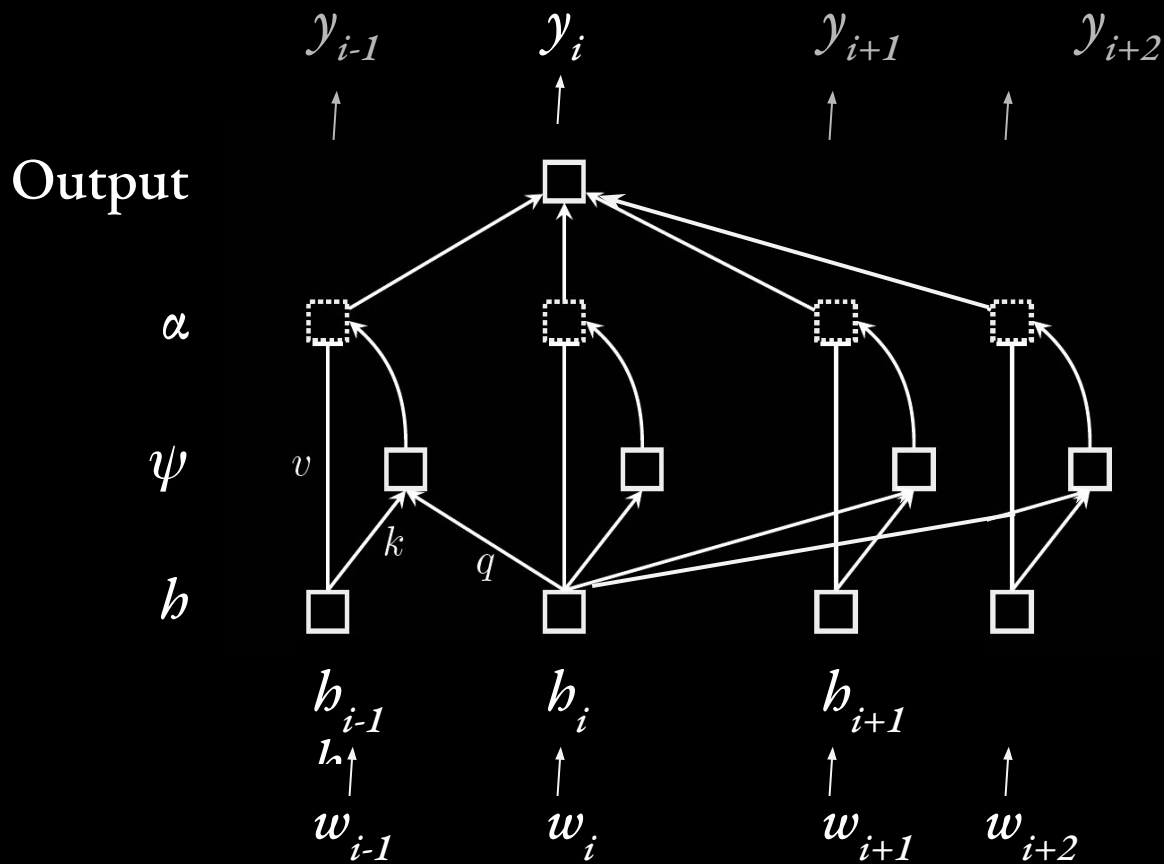
(Eisenstein, 2018)

# The Transformer: Attention-only Models

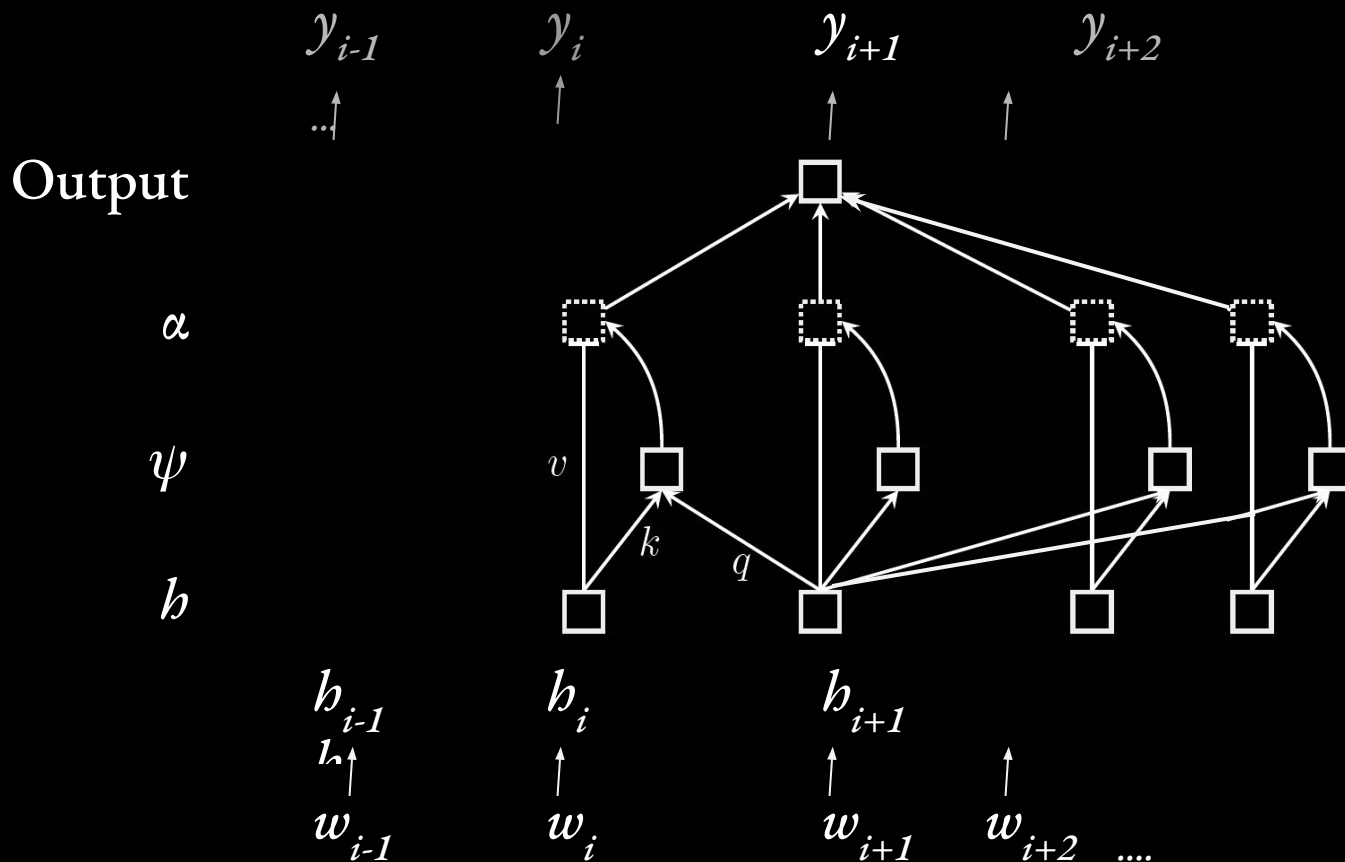


# The Transformer: Attention-only Models



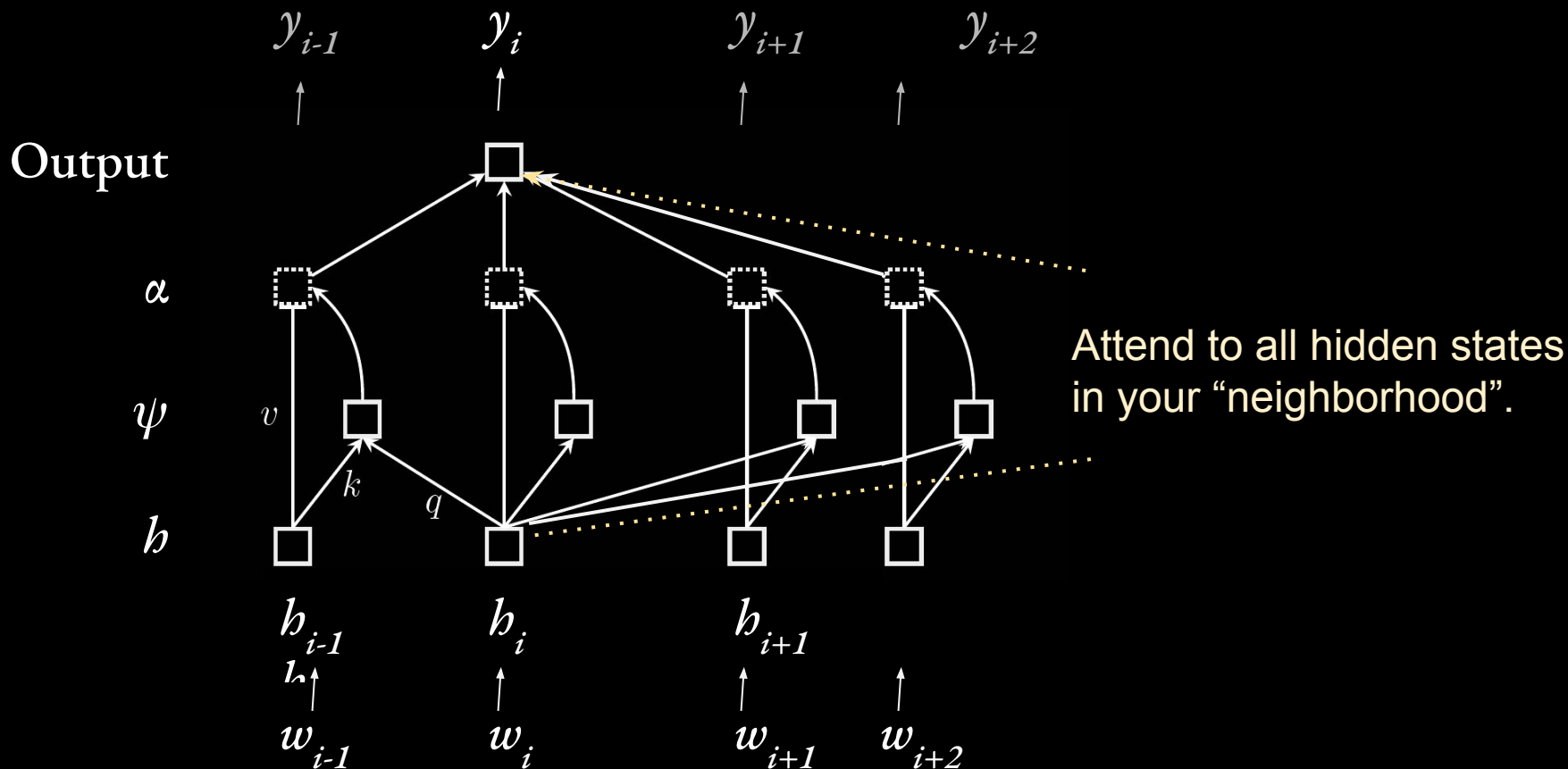


# The Transformer: “Attention-only” models





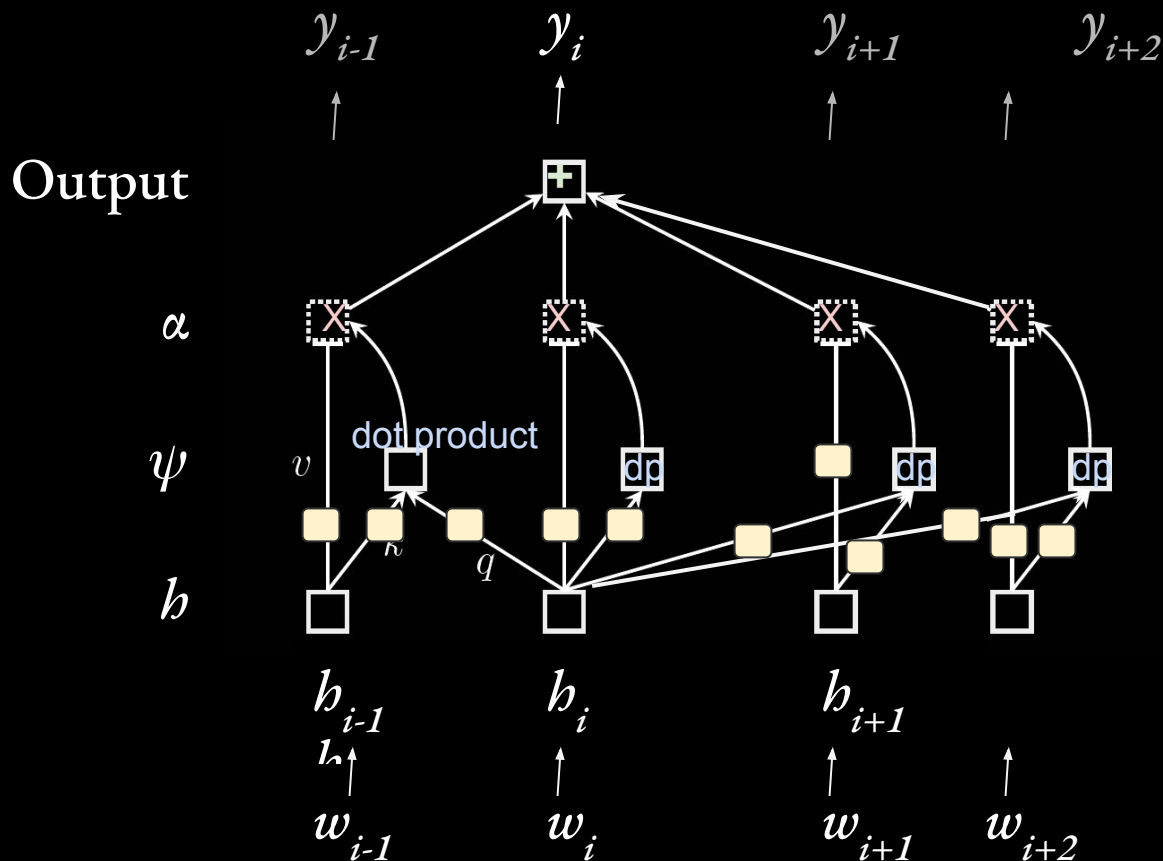
# The Transformer: “Attention-only” models







# The Transformer: "Attention-only" models



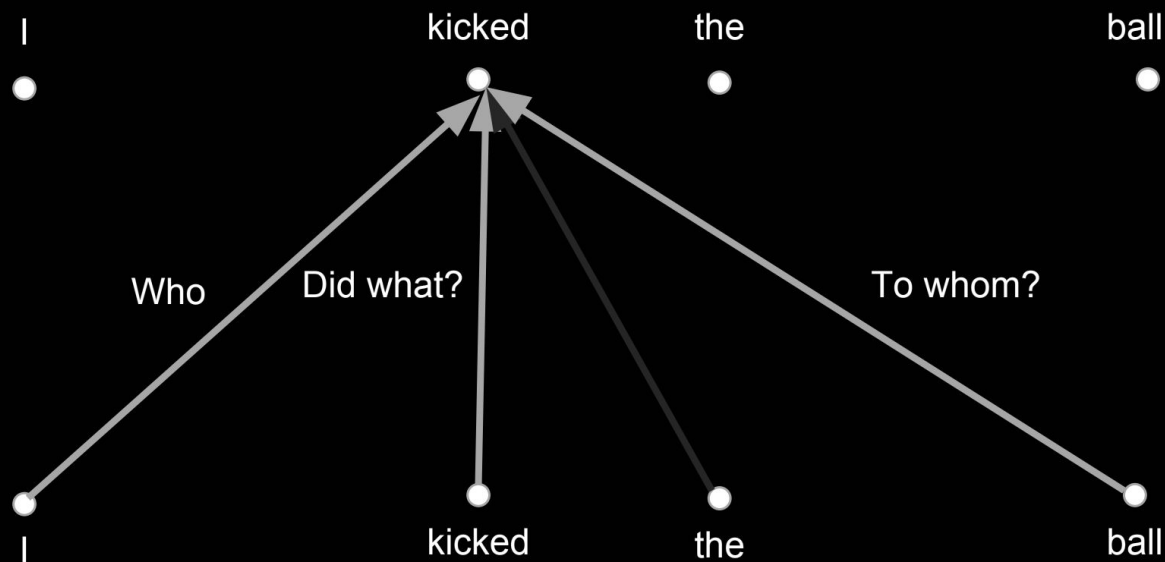
$$\psi_{dp}(k, q) = (k^t q) \sigma$$

Linear layer:  
 $W^T X$

One set of weights for  
each of for K, Q, and V

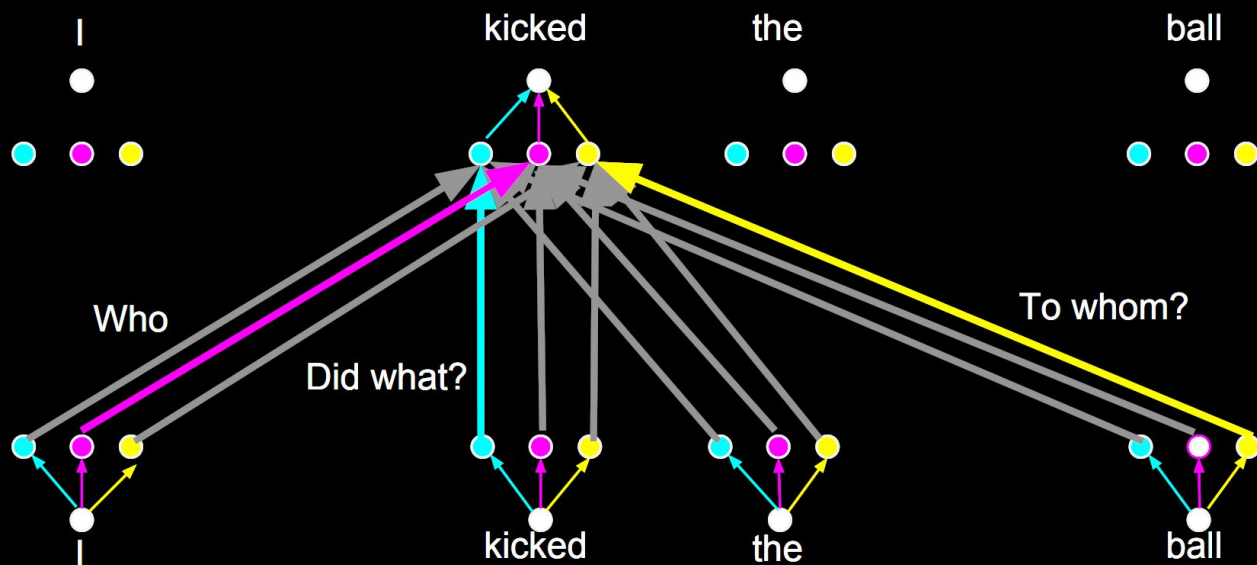
# The Transformer

Limitation (thus far): Can't capture multiple types of dependencies between words.

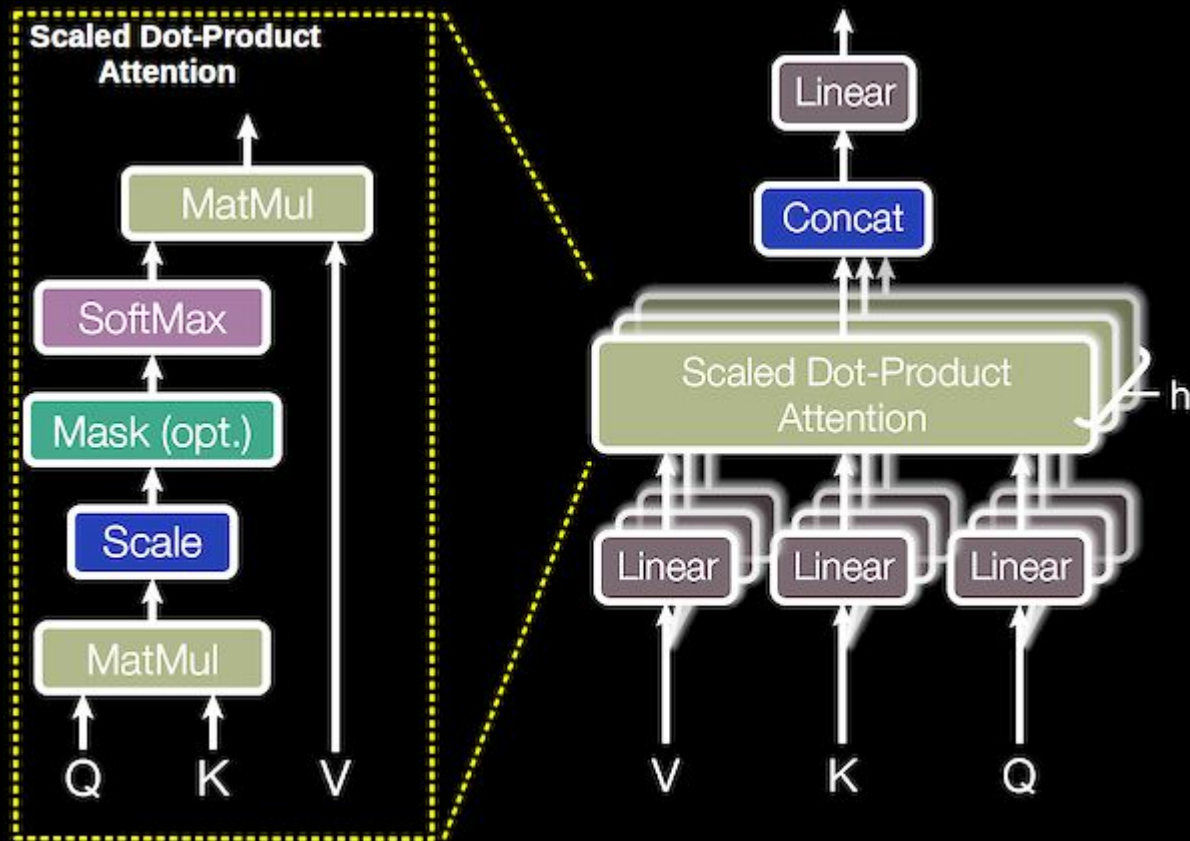


# The Transformer

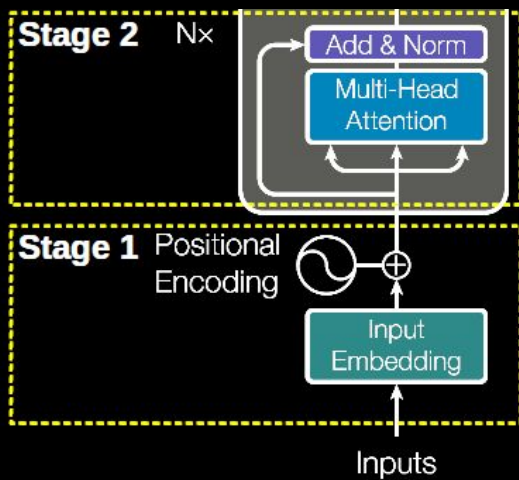
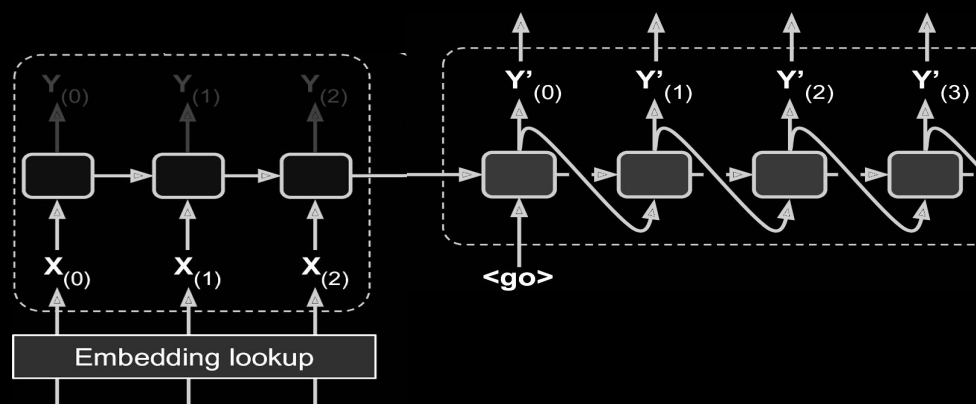
Solution: Multi-head attention



# Multi-head Attention

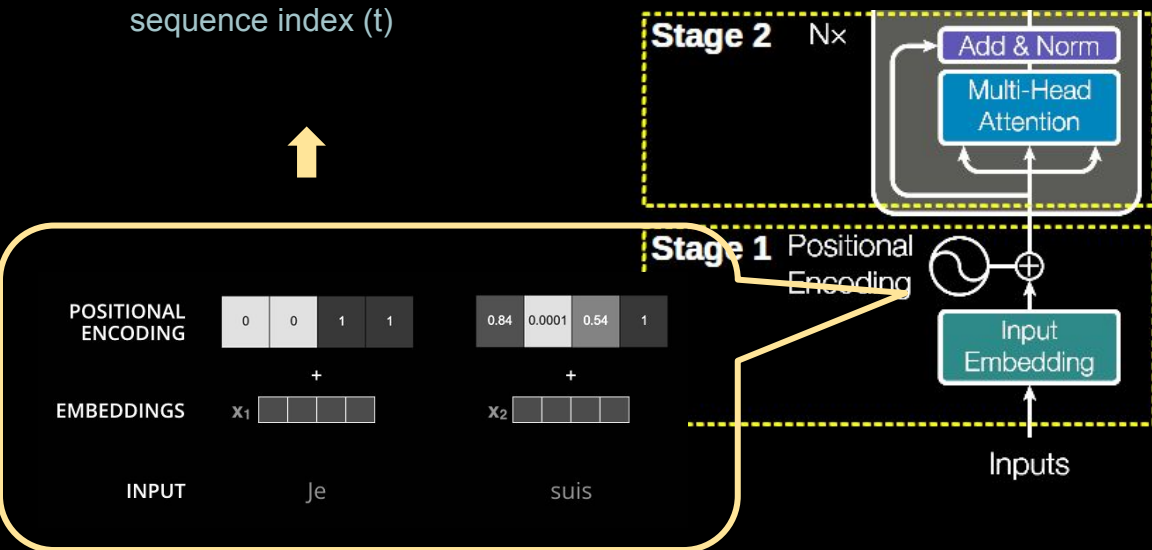
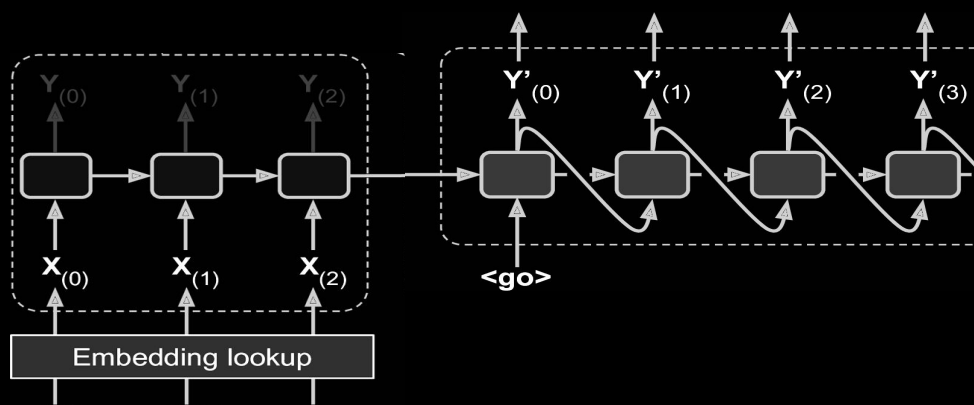
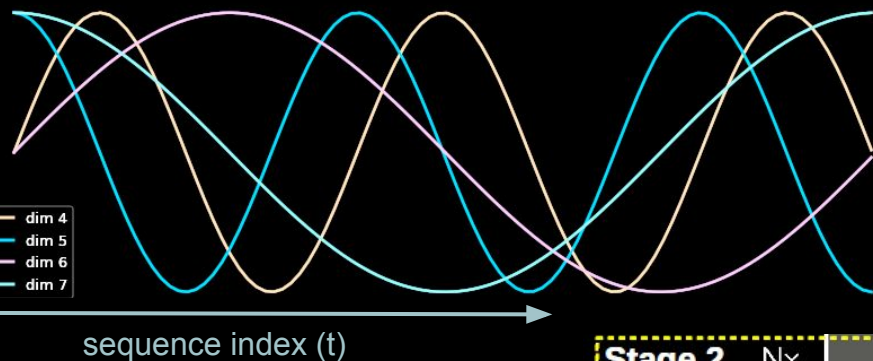


# Transformer for Encoder-Decoder

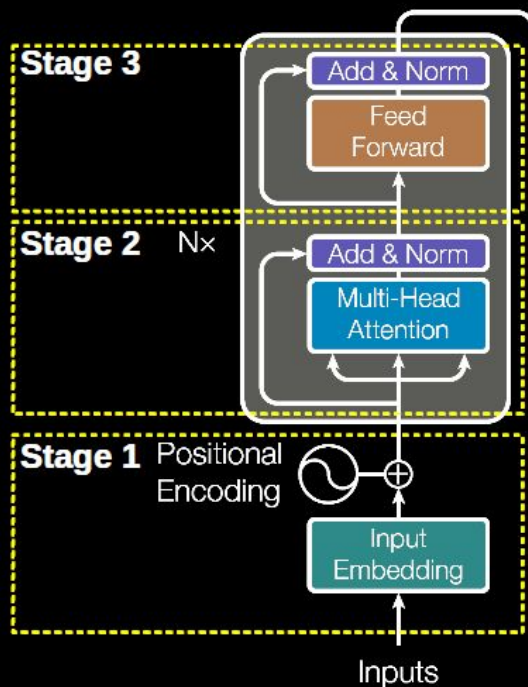
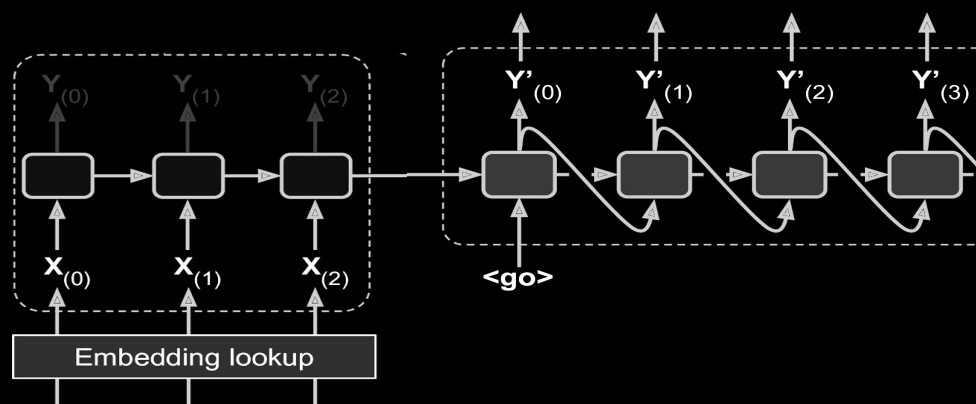




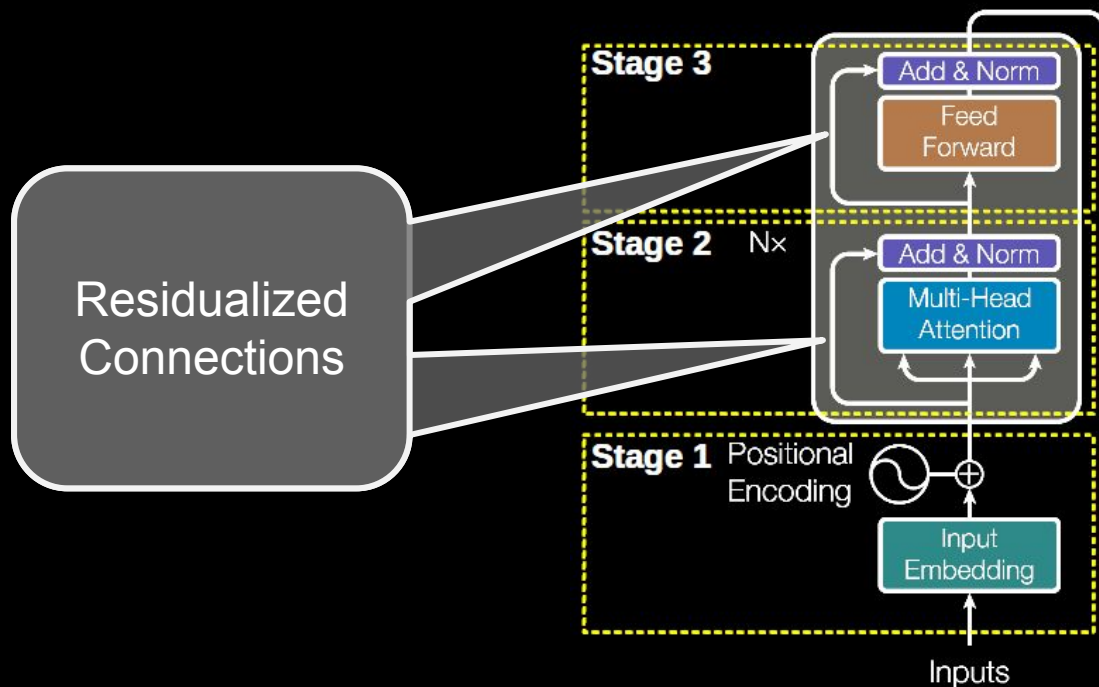
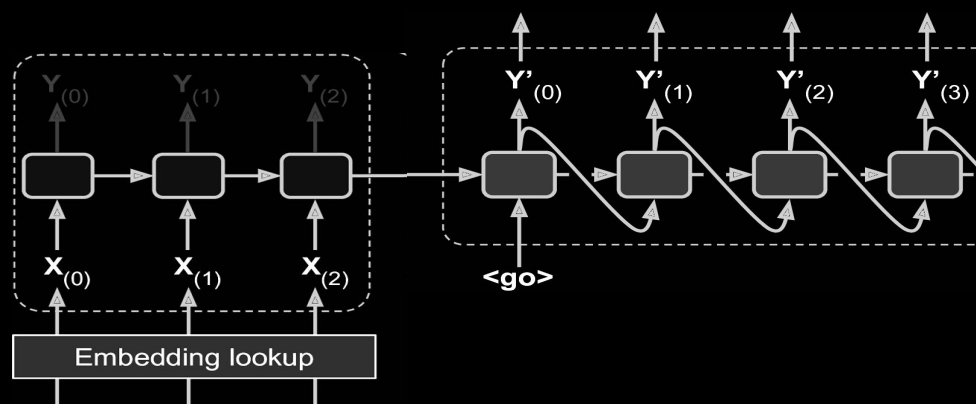
# Transformer for Encoder-Decoder



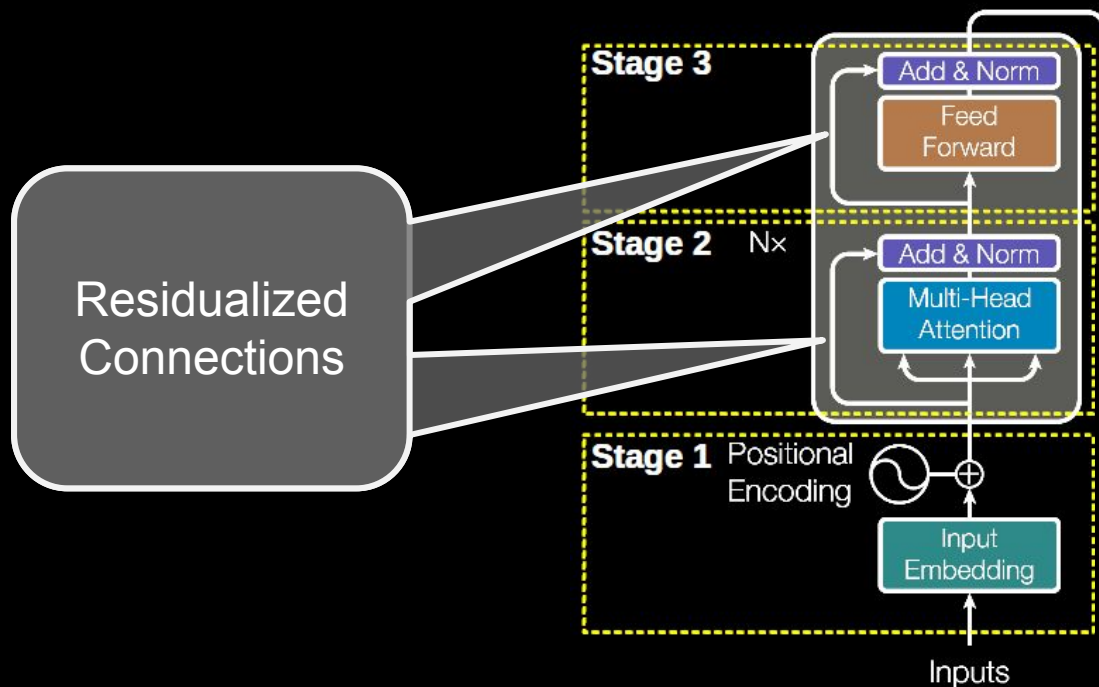
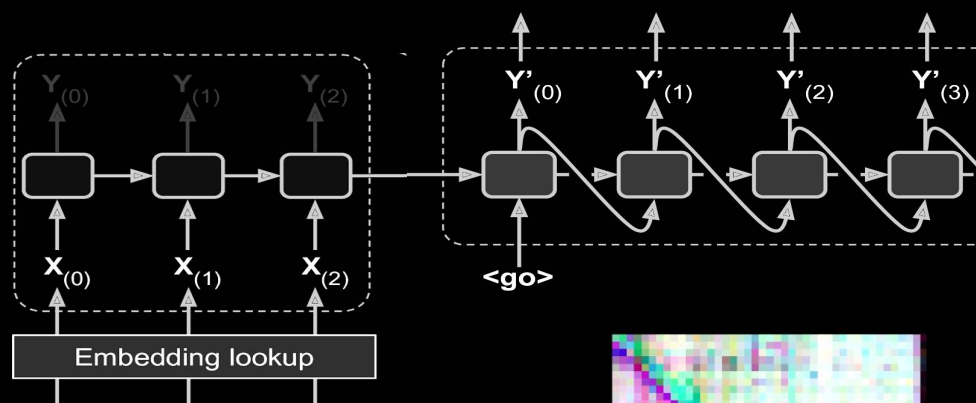
# Transformer for Encoder-Decoder



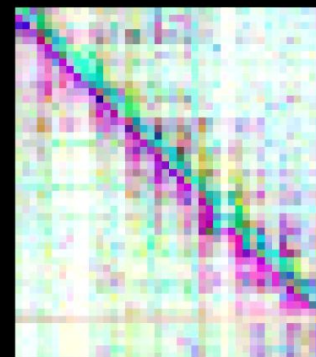
# Transformer for Encoder-Decoder



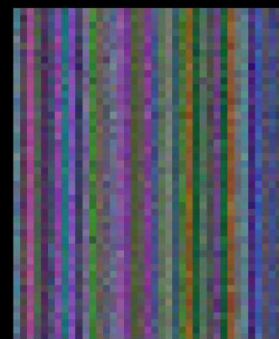
# Transformer for Encoder-Decoder



residuals enable positional information to be passed along

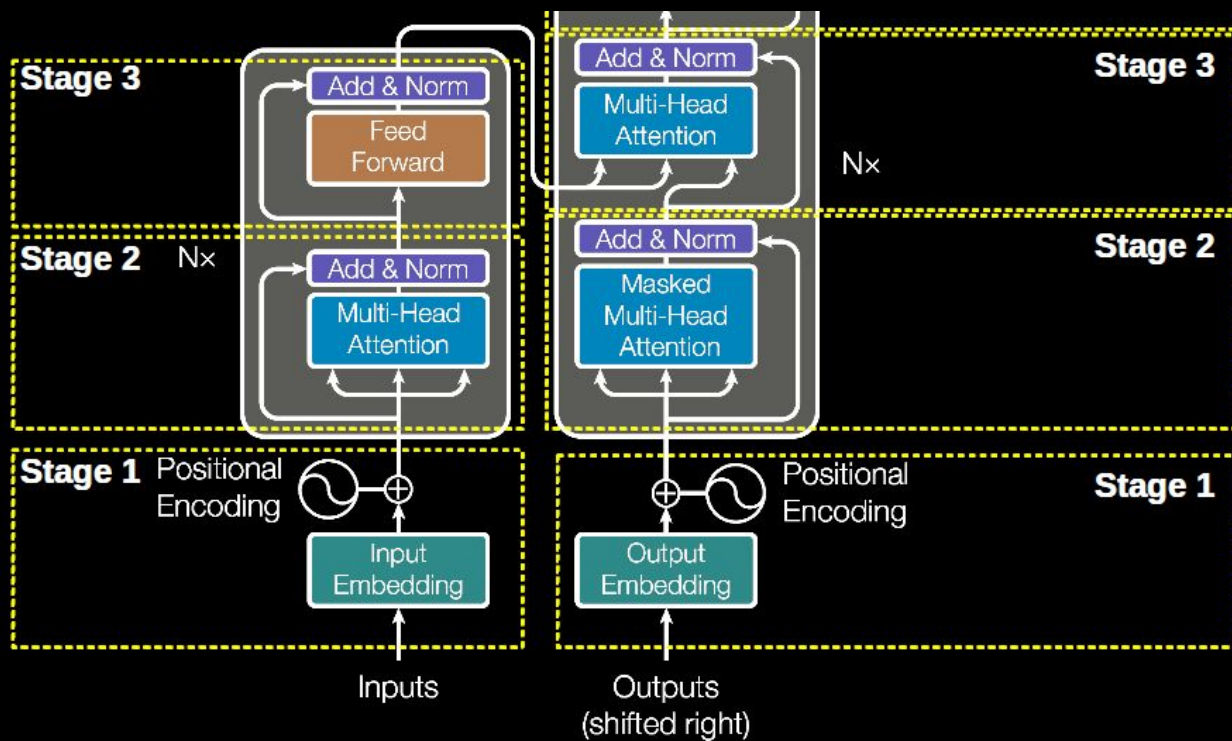
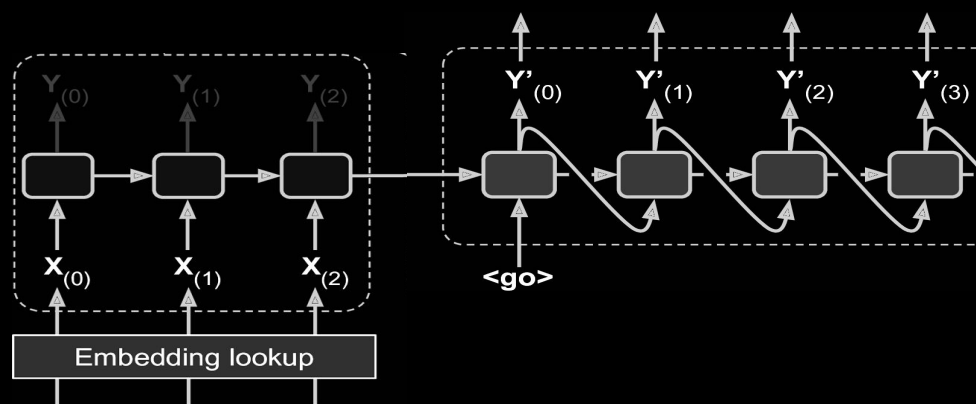


With residuals

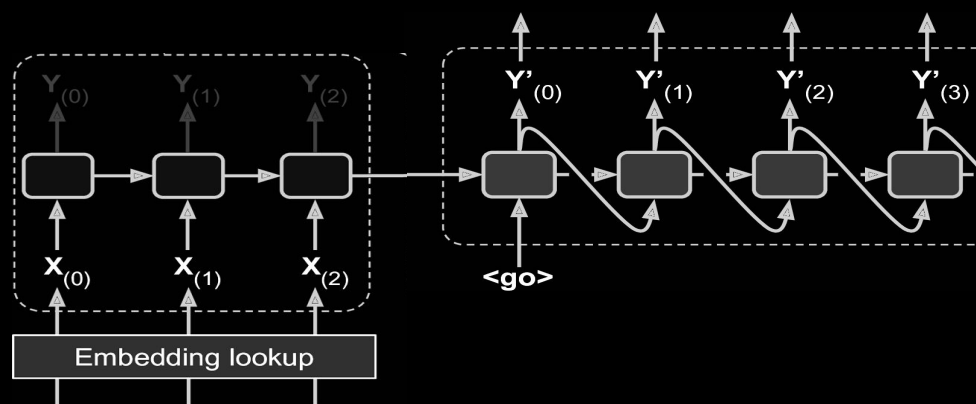


Without residuals

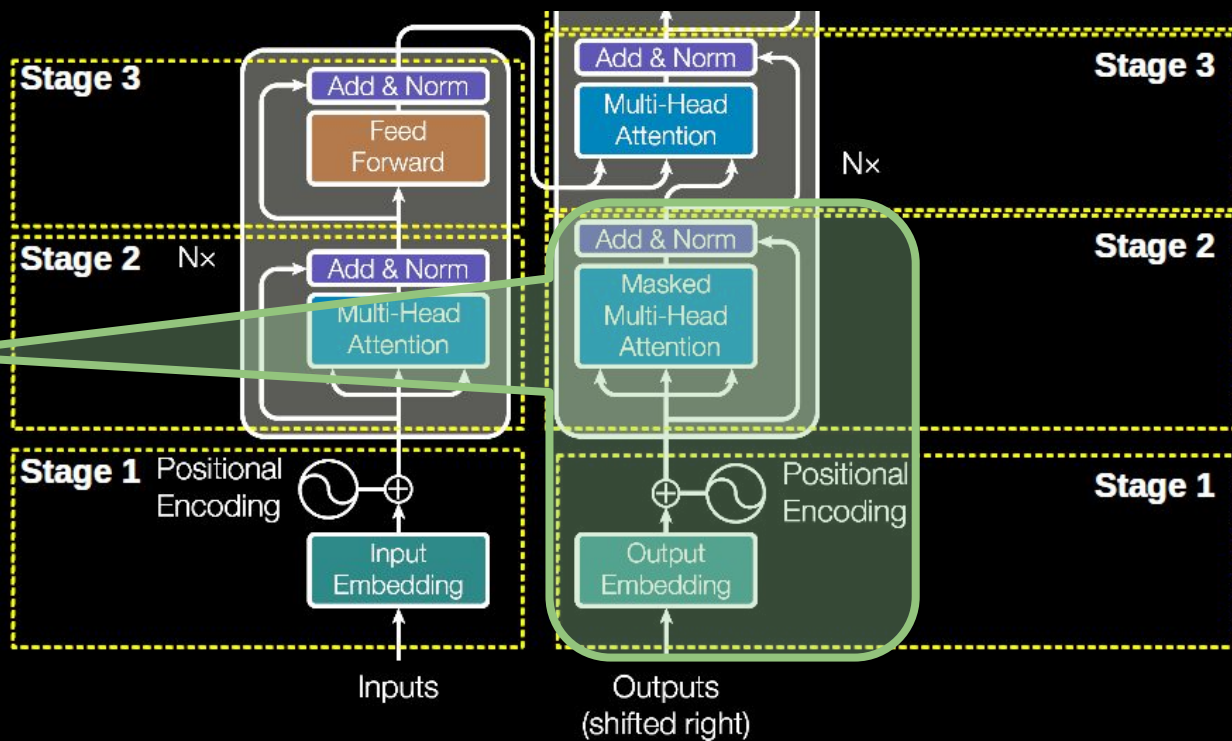
# Transformer for Encoder-Decoder



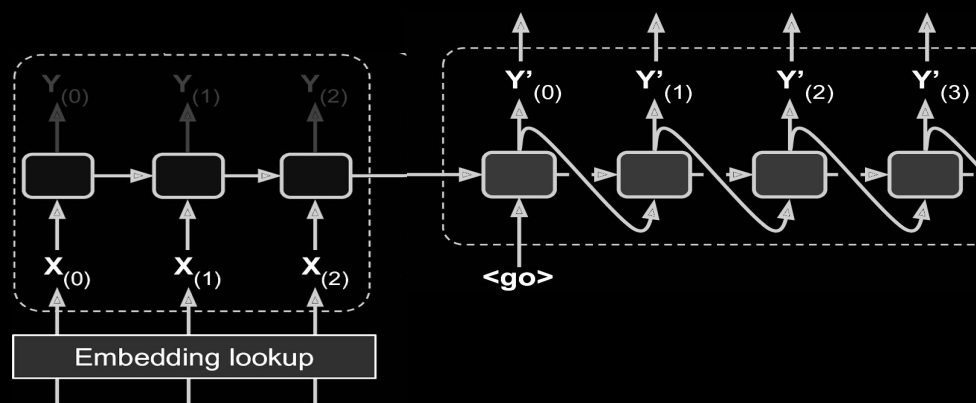
# Transformer for Encoder-Decoder



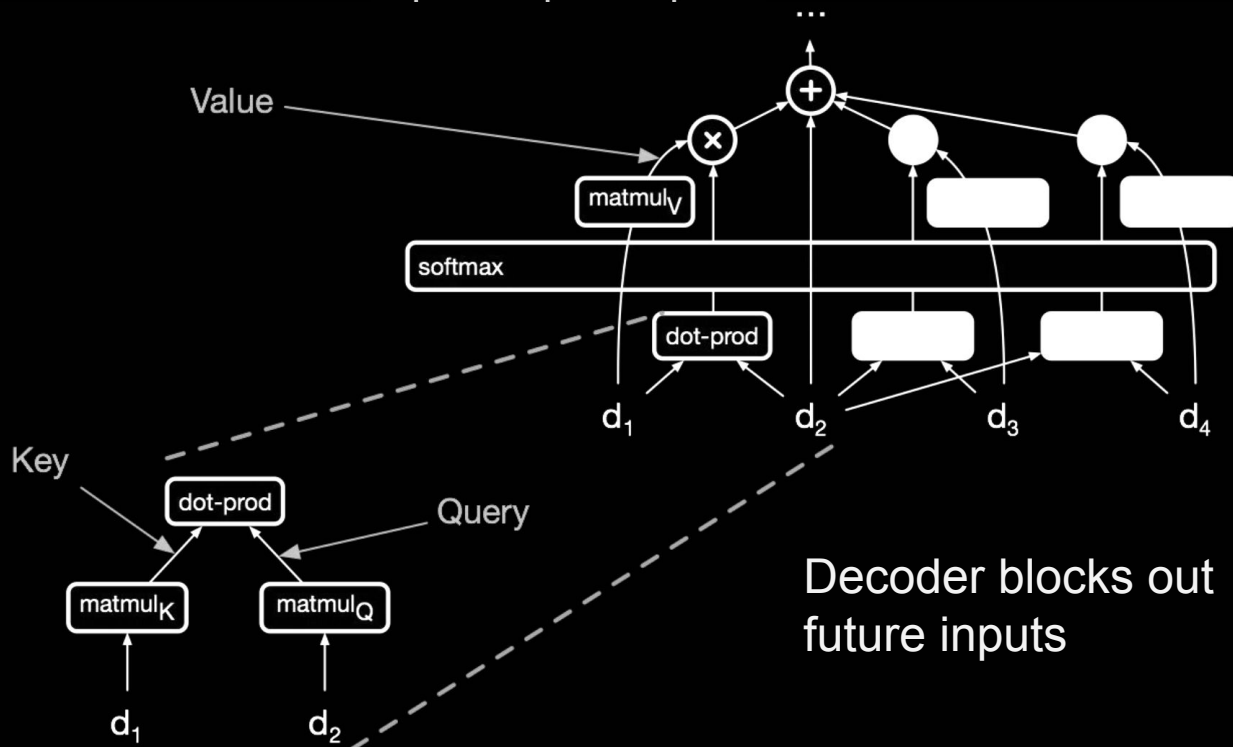
essentially, a language model



# Transformer for Encoder-Decoder

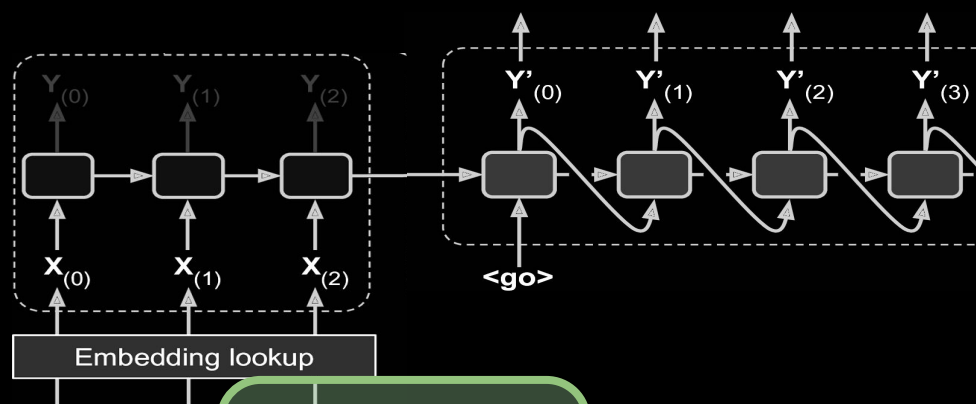


essentially, a language model



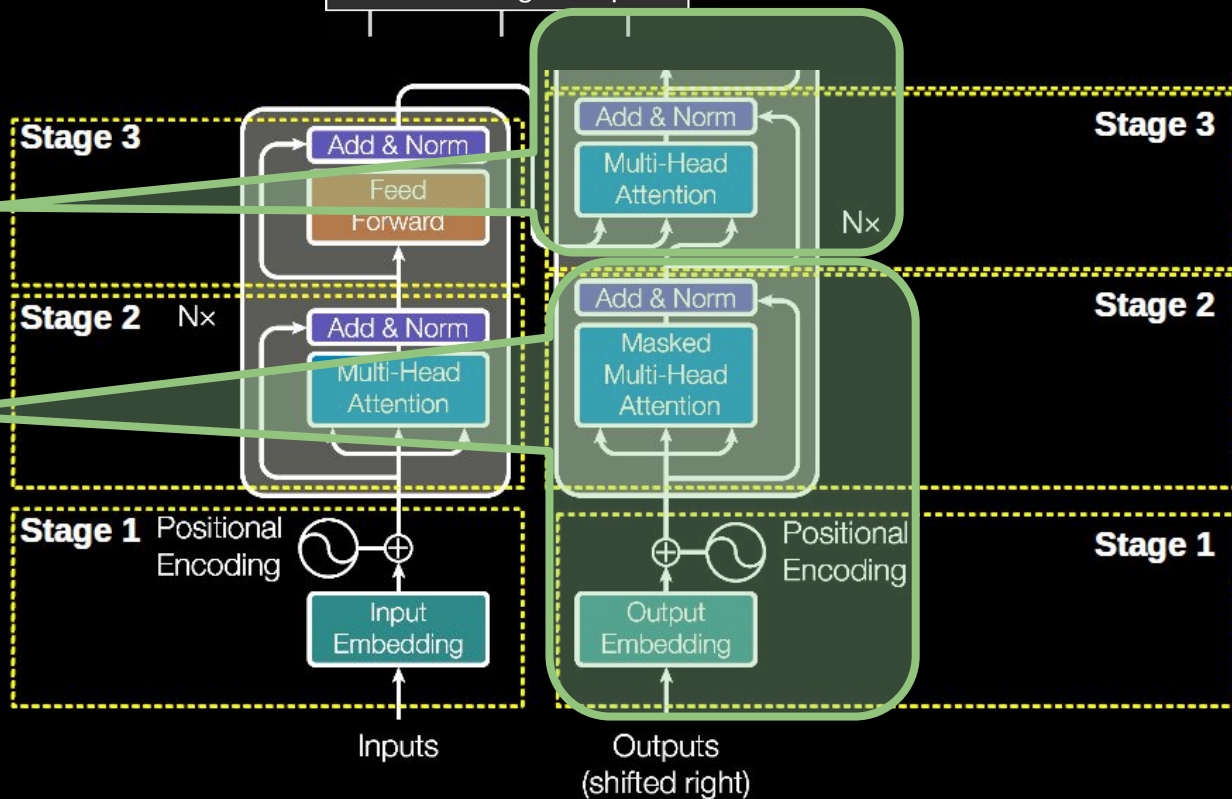
Decoder blocks out future inputs

# Transformer for Encoder-Decoder



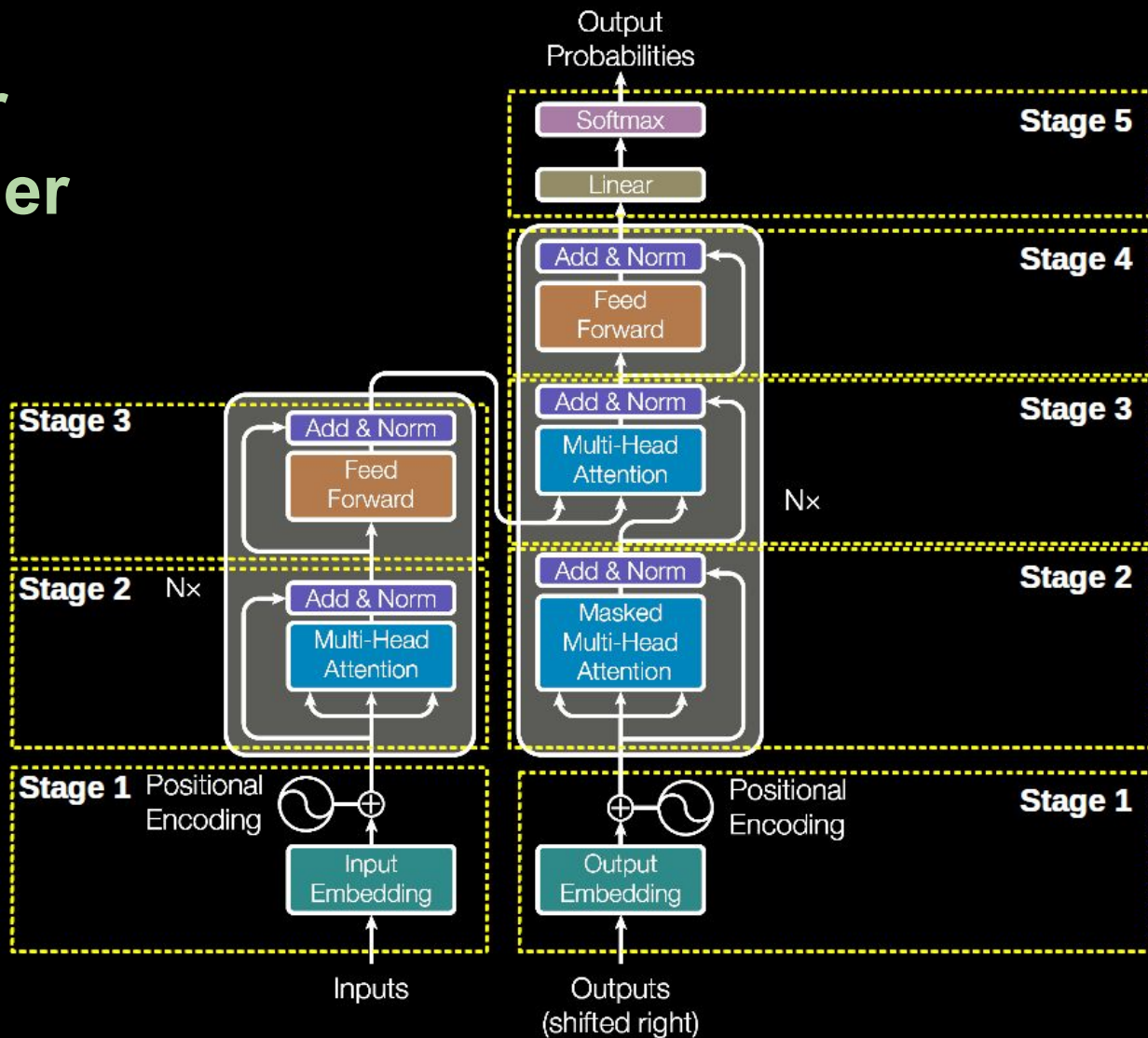
Add conditioning of the LM based on the encoder

essentially, a language model





# Transformer for Encoder-Decoder



# Transformer (as of 2017)

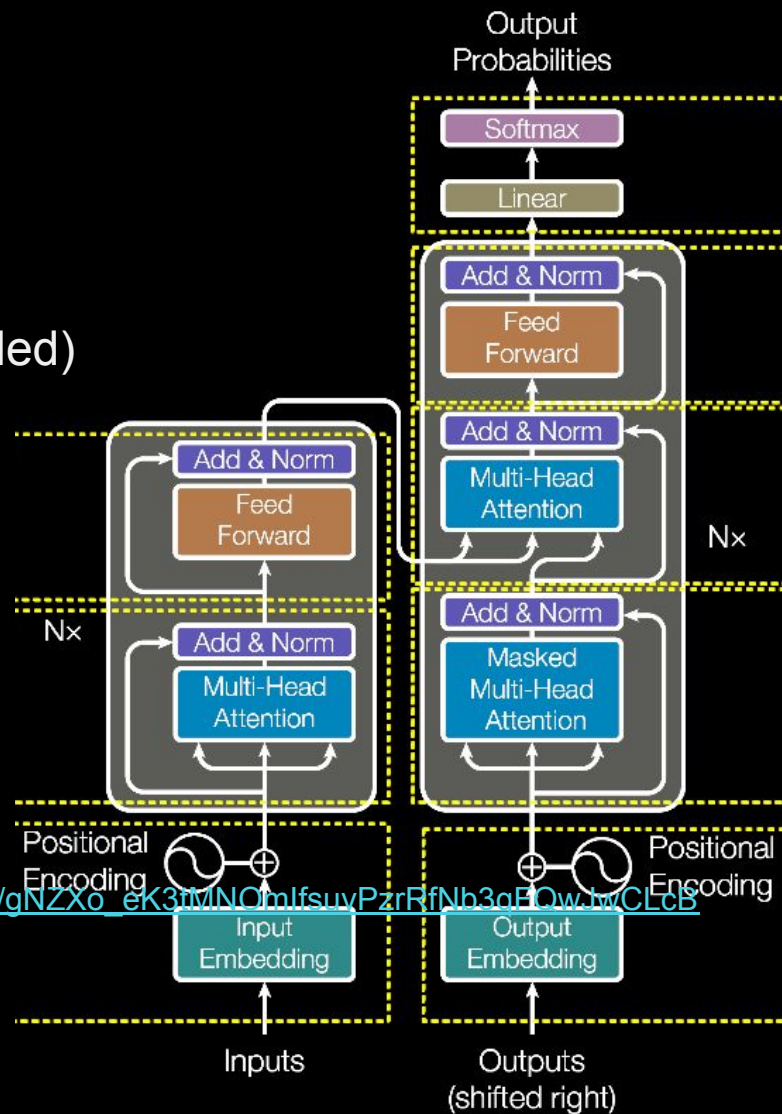
“WMT-2014” Data Set. BLEU scores:

	EN-DE	EN-FR
GNMT (orig)	24.6	39.9
ConvSeq2Seq	25.2	40.5
Transformer*	<b>28.4</b>	<b>41.8</b>

# Transformer

- Utilize Self-Attention
- Simple att scoring function (dot product, scaled)
- Added linear layers for Q, K, and V
- Multi-head attention
- Added positional encoding
- Added residual connection
- Simulate decoding by masking

[https://4.bp.blogspot.com/-OlrV-PAtEkQ/W3RkOJCBkaI/AAAAAAAAADOg/gNZXo\\_eK3tMNOmIfsuyPzrRfNb3qEQwIwCLcB/GAs/s640/image1.gif](https://4.bp.blogspot.com/-OlrV-PAtEkQ/W3RkOJCBkaI/AAAAAAAAADOg/gNZXo_eK3tMNOmIfsuyPzrRfNb3qEQwIwCLcB/GAs/s640/image1.gif)



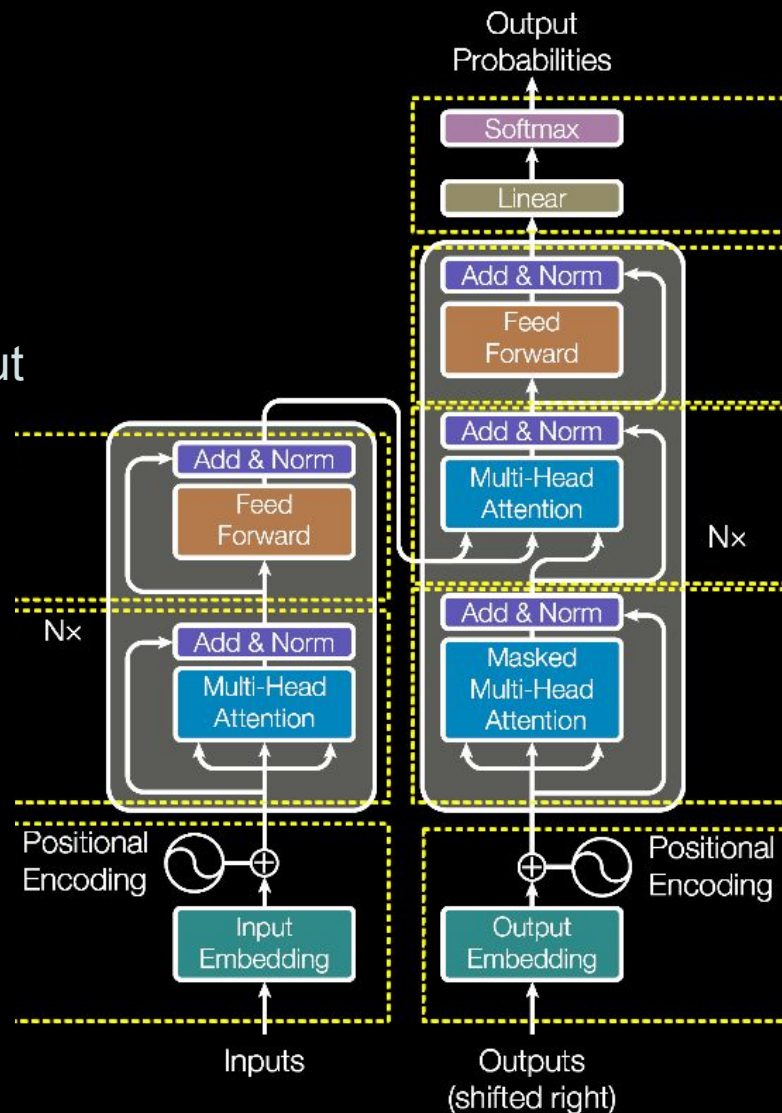
# Transformer

## Why?

- Don't need complexity of LSTM/GRU cells
- Constant num edges between words (or input steps)
- Enables “interactions” (i.e. adaptations) between words
- Easy to parallelize -- don't need sequential processing.

## Drawbacks:

- Only unidirectional by default
- Only a “single-hop” relationship per layer (multiple layers to capture multiple)



# BERT

## Bidirectional Encoder Representations from Transformers

Produces contextualized embeddings  
(or pre-trained contextualized encoder)

### **Drawbacks of Vanilla Transformers:**

- Only unidirectional by default
- Only a “single-hop” relationship per layer  
(multiple layers to capture multiple)

# BERT

## Bidirectional Encoder Representations from Transformers

Produces contextualized embeddings  
(or pre-trained contextualized encoder)

- Bidirectional context by “masking” in the middle
- A lot of layers, hidden states, attention heads.

### **Drawbacks of Vanilla Transformers:**

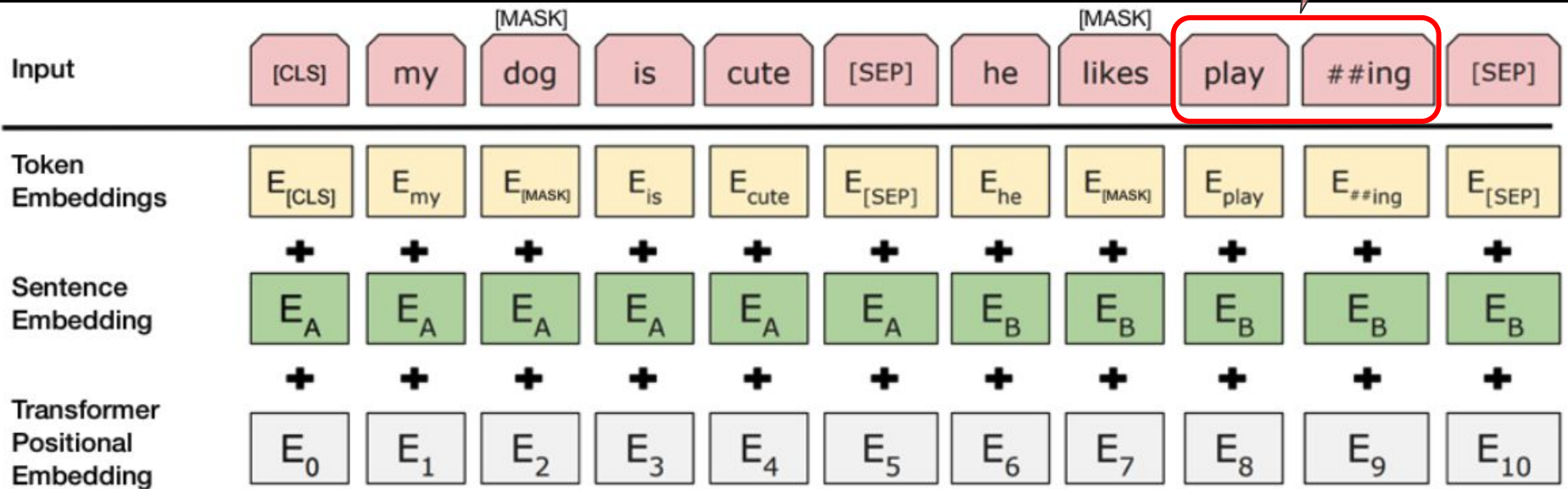
- Only unidirectional by default
- Only a “single-hop” relationship per layer  
(multiple layers to capture multiple)

# BERT

**Sentence A** = The man went to the store.  
**Sentence B** = He bought a gallon of milk.  
**Label** = IsNextSentence

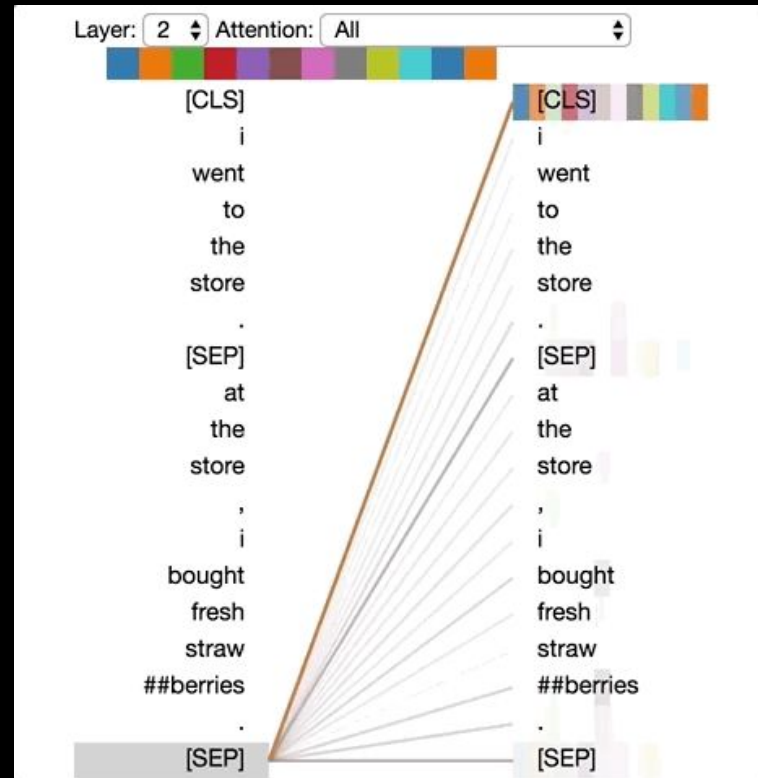
**Sentence A** = The man went to the store.  
**Sentence B** = Penguins are flightless.  
**Label** = NotNextSentence

tokenize into "word pieces"



# Bert: Attention by Layers

<https://colab.research.google.com/drive/1vIOJ1lhdujVjfH857hvYKIdKPTD9Kid8>

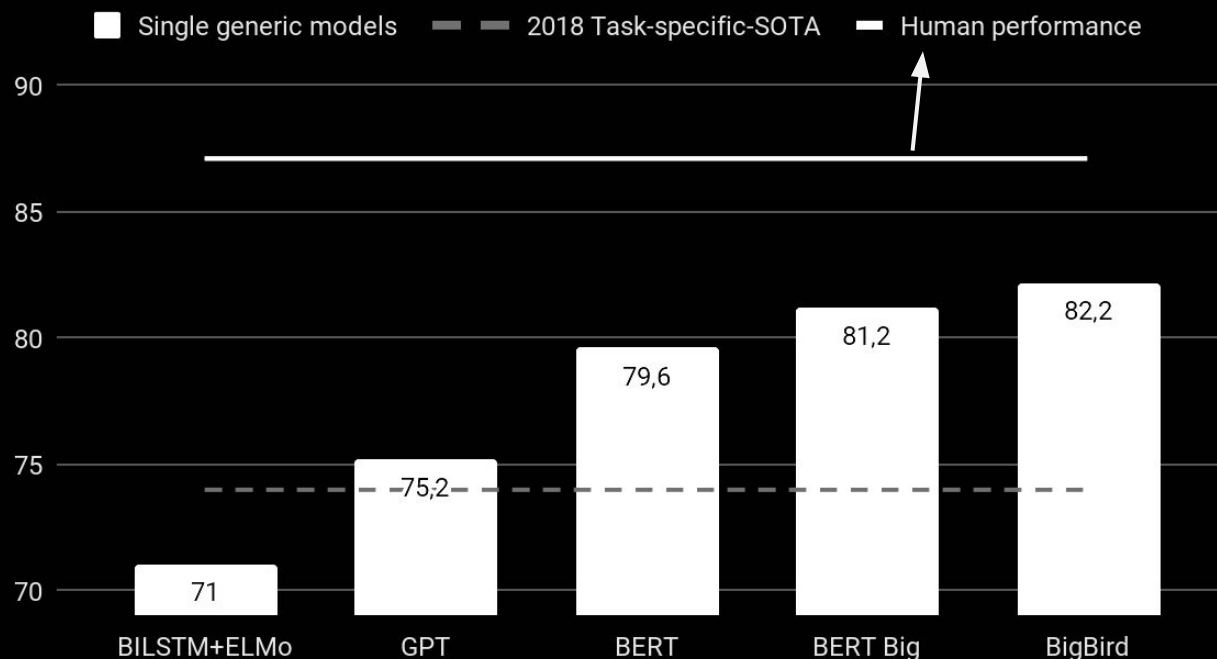


(Vig, 2019)



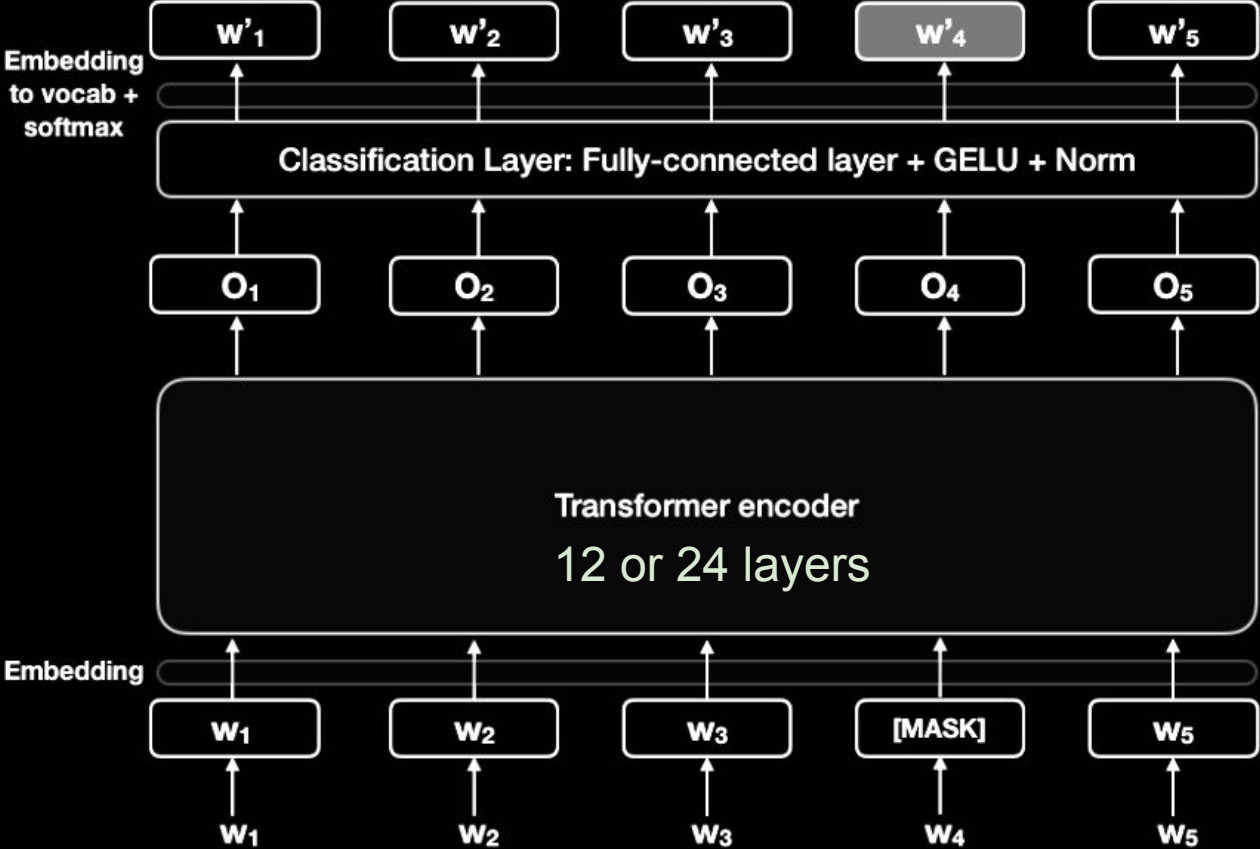
# BERT Performance: e.g. Question Answering

GLUE scores evolution over 2018-2019

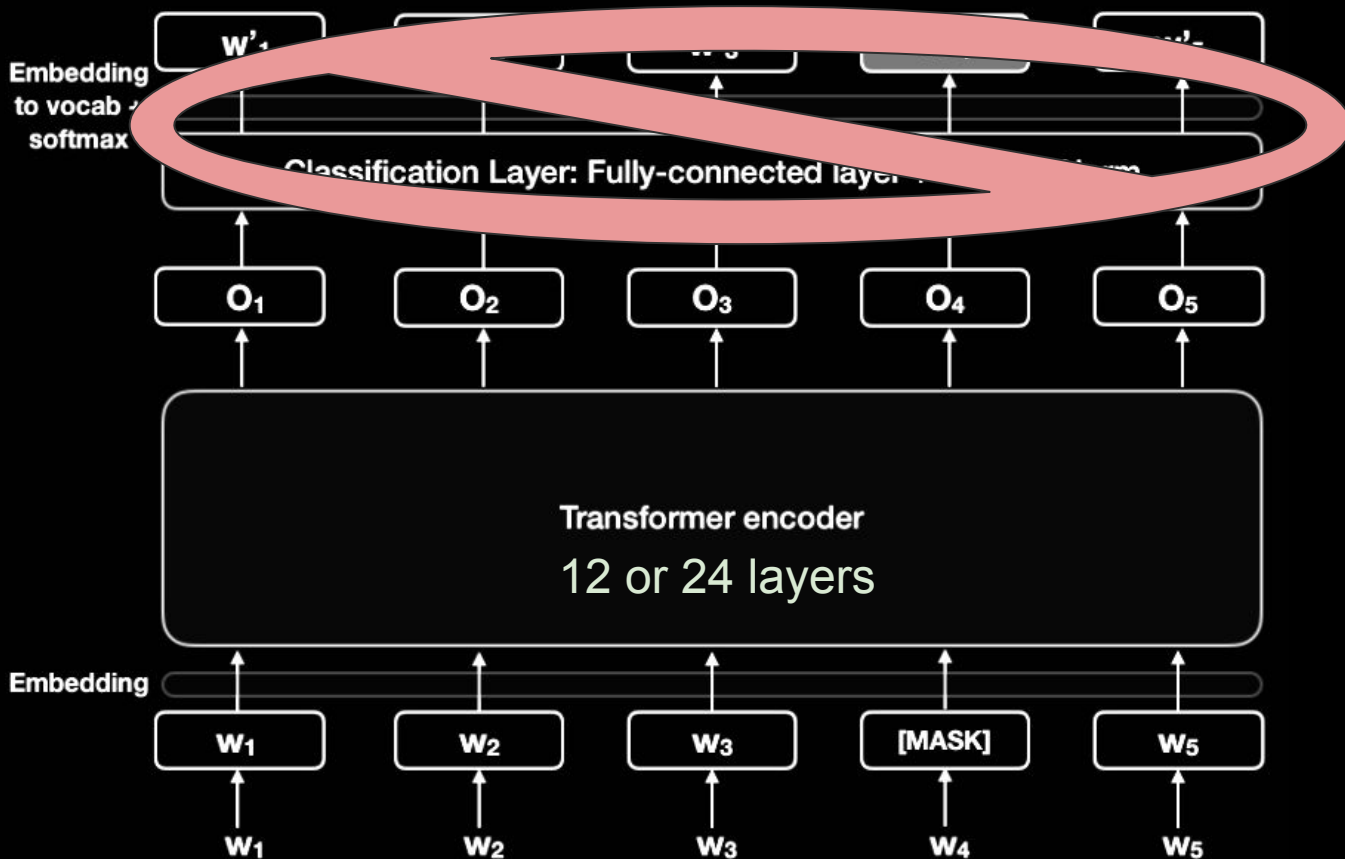


<https://rajpurkar.github.io/SQuAD-explorer/>

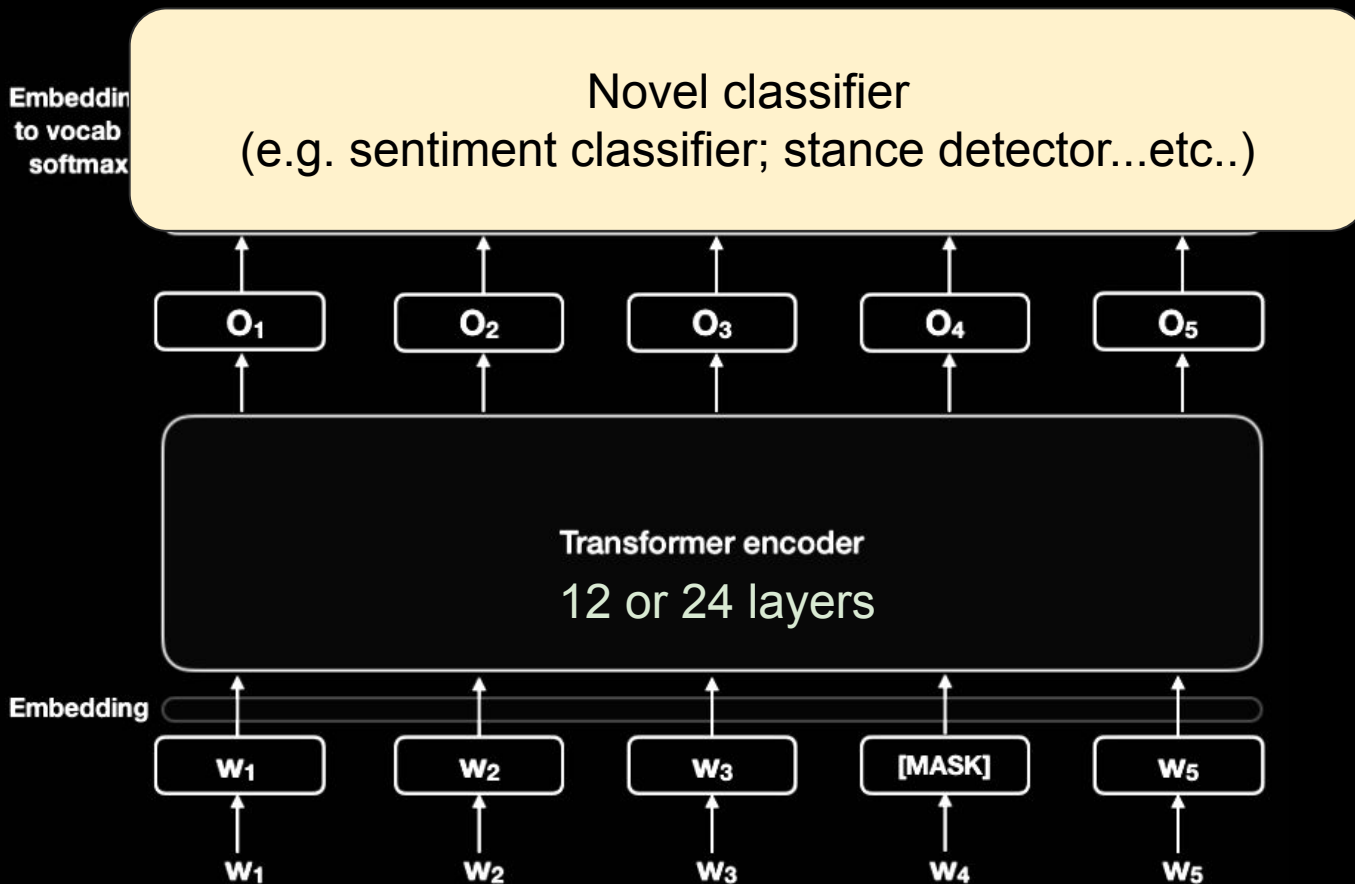
# BERT: Pre-training; Fine-tuning



# BERT: Pre-training; Fine-tuning



# BERT: Pre-training; Fine-tuning



# Neural Network Summary

- Goal is accurate prediction of  $y$  (outcome) given features ( $x$ )
- Use L1 or L2 penalization (as a regularization) to avoid overfit
- Reason for Train, Dev, Test split
- Components of a neural network
- Batch Normalization
- Distribution options: why is data parallelism preferred?
- Recurrent Neural Network
- Convolution Operation with Filters

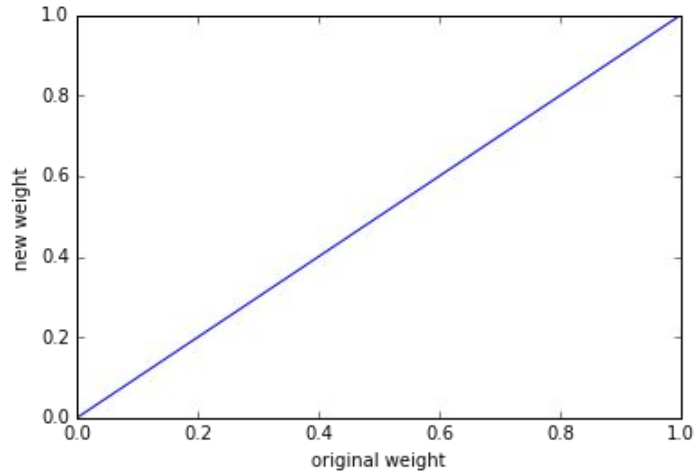
# Feature Selection / Subset Selection

## (bad) solution to overfit problem

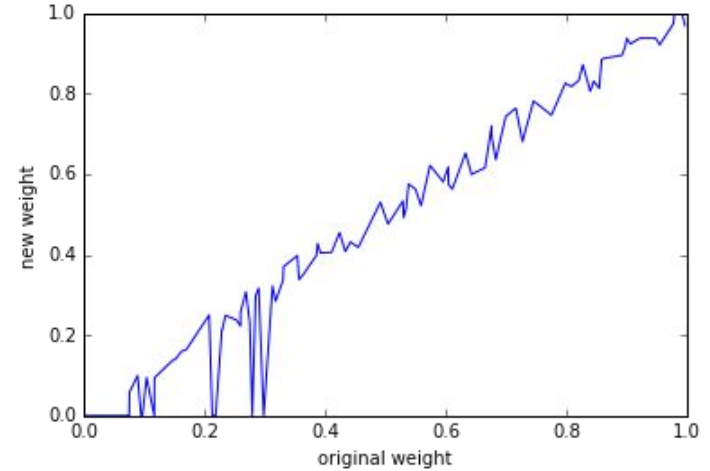
Use less features based on Forward Stepwise Selection:

- start with `current_model` just has the intercept (mean)  
`remaining_predictors = all_predictors`  
for `i` in `range(k)`:  
    #find best `p` to add to `current_model`:  
    for `p` in `remaining_predictors`  
        refit `current_model` with `p`  
        #add best `p`, based on  $RSS_p$  to `current_model`  
    #remove `p` from `remaining predictors`

# Regularization (Shrinkage)



No selection (weight= $\beta$ )



forward stepwise

Why just keep or discard features?

# Regularization (L2, Ridge Regression)

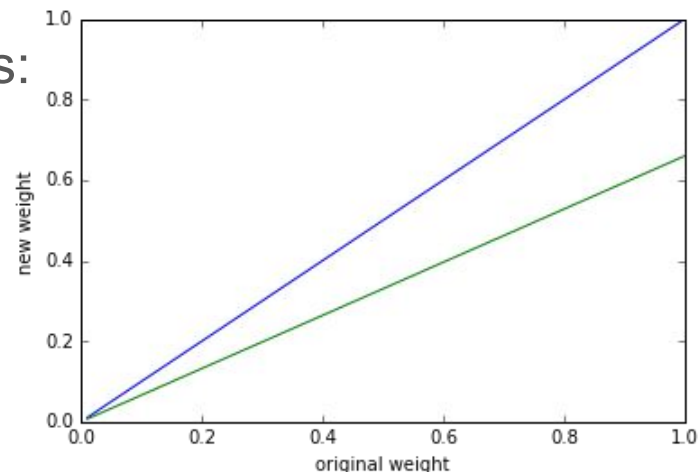
Idea: Impose a penalty on size of weights:

Ordinary least squares objective:

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

Ridge regression:

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





# Regularization (L2, Ridge Regression)

Idea: Impose a penalty on size of weights:

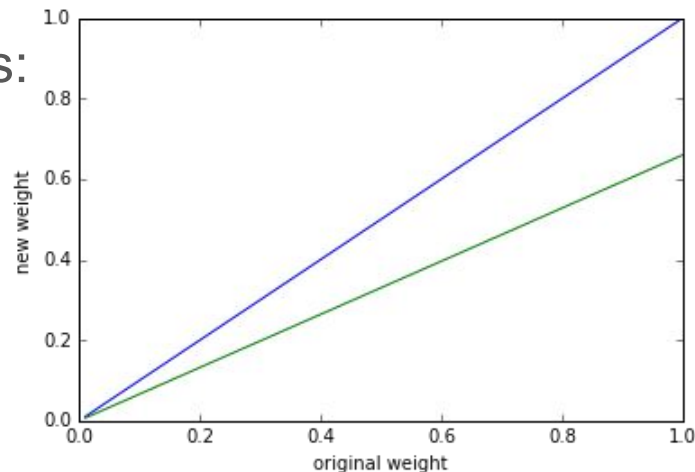
Ordinary least squares objective:

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

Ridge regression:

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

$$\lambda \sum_{j=1}^m \beta_j^2$$



→  $\lambda \|\beta\|_2^2$

# Regularization (L2, Ridge Regression)

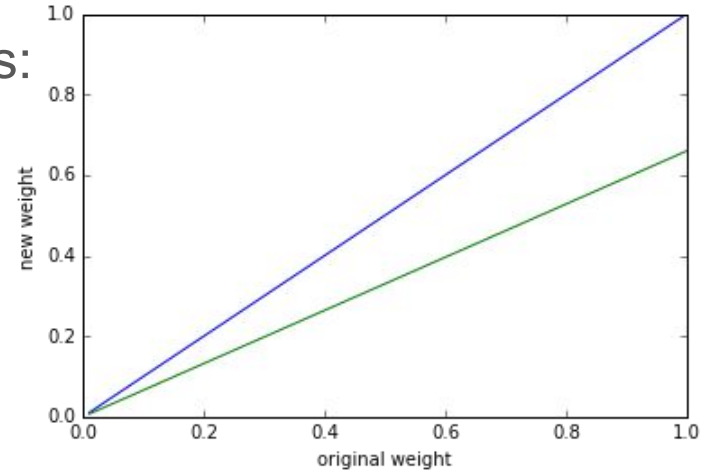
Idea: Impose a penalty on size of weights:

Ordinary least squares objective:

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

Ridge regression:

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



In Matrix Form:

$$\text{RSS}(\lambda) = (y - X\beta)^T (y - X\beta) + \lambda \beta^T \beta$$

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

$I$ :  $m \times m$  identity matrix

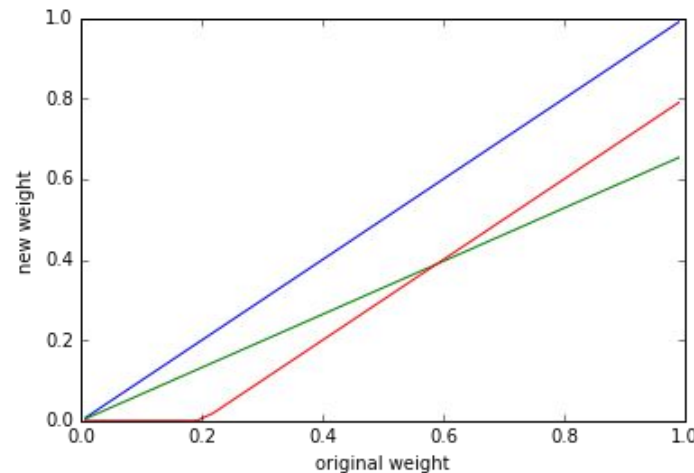
$$\lambda \|\beta\|_2^2$$

# Regularization (L1, The “Lasso”)

Idea: Impose a penalty and zero-out some weights

The Lasso Objective:

$$\hat{\beta}^{\text{lasso}} = \underset{\beta}{\operatorname{argmin}} \left\{ \frac{1}{2} \sum_{i=1}^N (Y_i - \sum_{j=1}^m x_{ij} \beta_j)^2 + \lambda \sum_{j=1}^m |\beta_j| \right\}$$



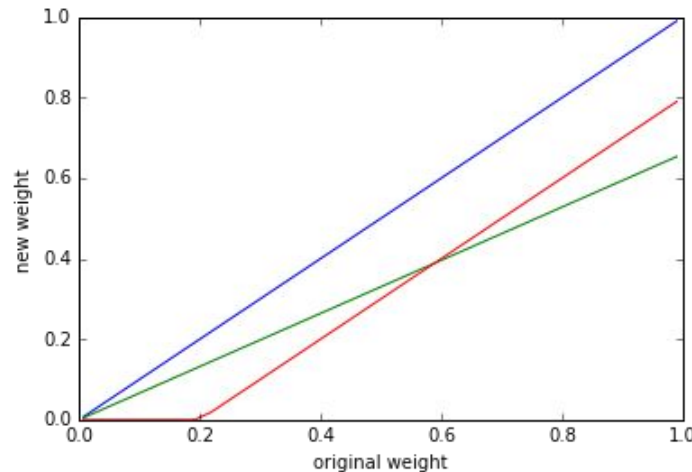
→  $\lambda \|\beta\|_1$

# Regularization (L1, The “Lasso”)

Idea: Impose a penalty and zero-out some weights

The Lasso Objective:

$$\hat{\beta}^{\text{lasso}} = \underset{\beta}{\operatorname{argmin}} \left\{ \frac{1}{2} \sum_{i=1}^N (Y_i - \sum_{j=1}^m x_{ij} \beta_j)^2 + \lambda \sum_{j=1}^m |\beta_j| \right\}$$



No closed form matrix solution, but often solved with coordinate descent.

$$\lambda \|\beta\|_1$$

Application:  $p \approx n$  or  $p \gg n$  (p: features; n: observations)

# Cluster Distribution

Model Parallelism

Multiple devices on multiple machines

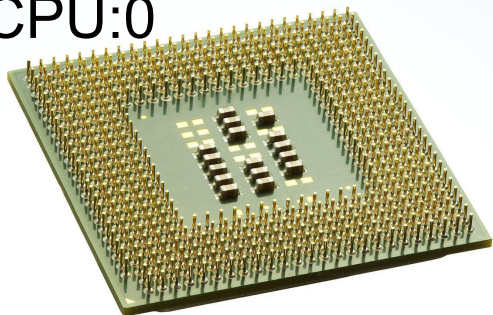
```
with tf.device("/cpu:1")  
    beta=tf.Variable(...)
```

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

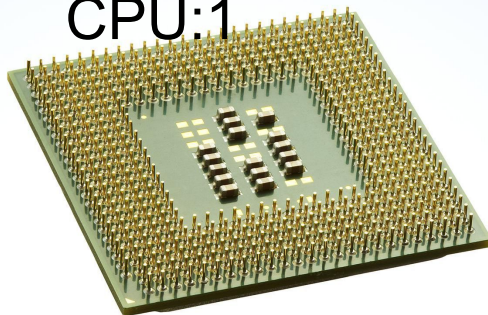
Transfer Tensors

Machine A

CPU:0



CPU:1



Machine B

GPU:0

