

Deep Learning with Apache MXNet

Julien Simon, Principal Technical Evangelist

@julsimon

Selected customers running AI on AWS



NETFLIX

Stanford



The Washington Post



Carnegie Mellon

Pinterest



C-SPAN



real networks



GoAnimate



图森 tu Simple

duolingo



zmags

Questions

- What's the business problem my IT has failed to solve?
 - That's probably where Deep Learning can help
- Should I design and train my own Deep Learning model?
 - Do I have the expertise?
 - Do I have enough time, data & compute to train it?
- Should I use a pre-trained model?
 - How well does it fit my use case?
 - On what data was it trained? How close is this to my own data?
- Should I use a SaaS solution?

Amazon AI for every developer

Services	Chat Amazon Lex		Speech Amazon Polly		Vision Amazon Rekognition	
Platforms	Amazon ML	Spark & EMR	Kinesis	Batch	ECS	
Engines	MXNet	TensorFlow	Caffe	Theano	Pytorch	CNTK
Infrastructure	GPU		CPU		IoT	Mobile

More information at aws.amazon.com/ai

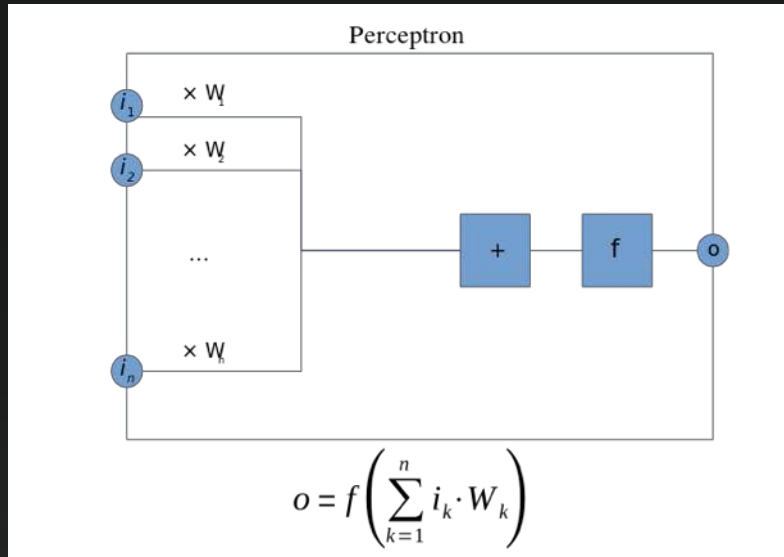


Neural Networks

1957 – The Perceptron

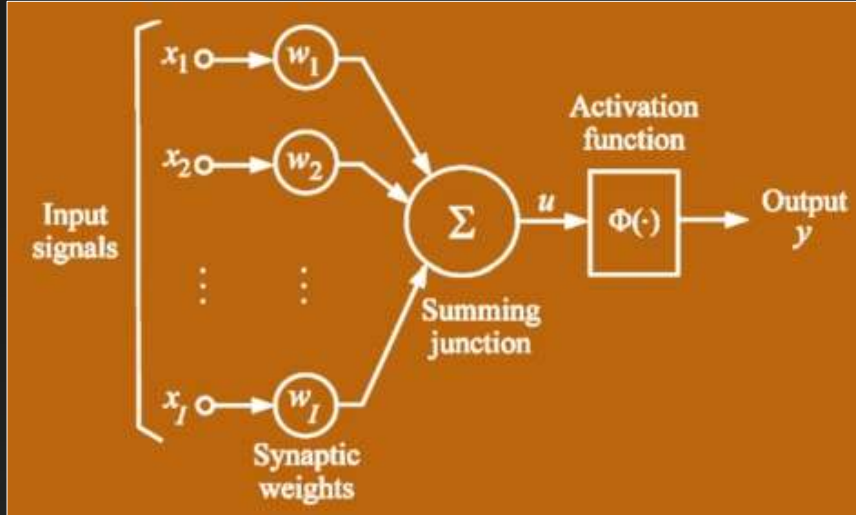


Frank Rosenblatt
Artificial Intelligence pioneer



Appropriate weights are applied to the inputs, and the resulting weighted sum is passed to a function that produces the output.

The neuron



$$u = \sum_{i=1}^n w_i x_i$$

$$\mathbf{x} = [x_1, x_2, \dots, x_I]$$

$$\mathbf{w} = [w_1, w_2, \dots, w_I]$$

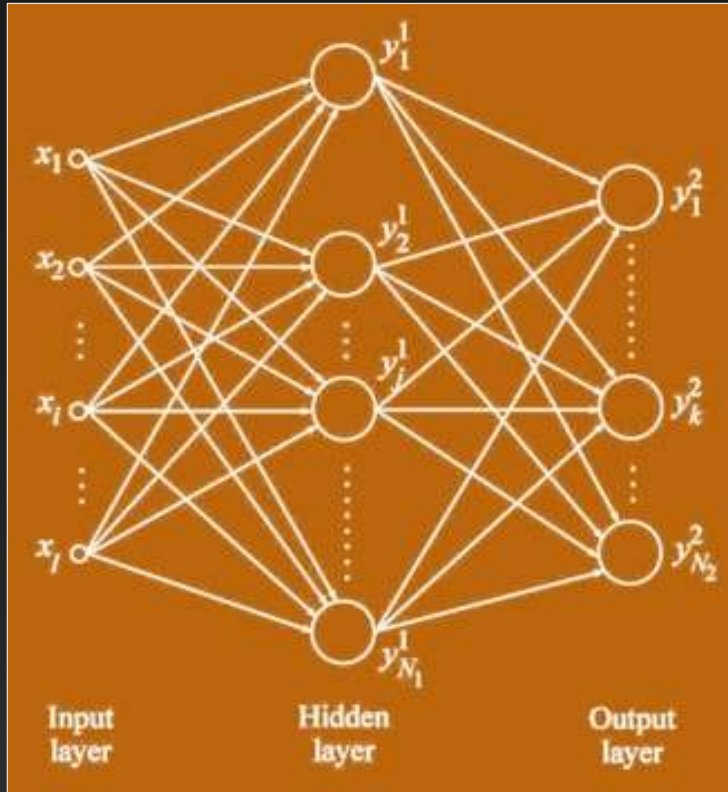
$$\mathbf{u} =$$

$$\mathbf{x} \cdot \mathbf{w}$$

Activation

Name	Plot	Equation
Identity		$f(x) = x$
Binary step		$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$
Logistic (a.k.a. Soft step)		$f(x) = \frac{1}{1 + e^{-x}}$
TanH		$f(x) = \tanh(x) = \frac{2}{1 + e^{-2x}} - 1$
ArcTan		$f(x) = \tan^{-1}(x)$
Softsign [7][8]		$f(x) = \frac{x}{1 + x }$
Rectified linear unit (ReLU) [9]		$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$

The neural network



$$X = \begin{bmatrix} x_{11}, x_{12}, \dots, x_{1l} \\ x_{21}, x_{22}, \dots, x_{2l} \\ \vdots \\ x_{m1}, x_{m2}, \dots, x_{ml} \end{bmatrix}$$

l features

m samples

$$y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix}$$

m labels

$$\text{Accuracy} = \frac{\text{Number of correct predictions}}{\text{Total number of predictions}}$$

The training process

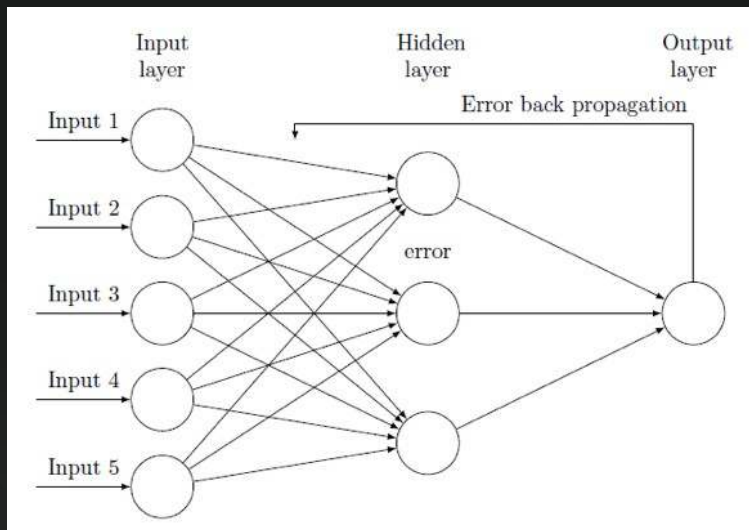
- The difference between the **predicted** output and the **actual** output (aka *ground truth*) is called the **prediction loss**.
- There are different ways to compute it (**loss functions**).
- The purpose of training is to iteratively **minimize loss** and **maximize accuracy** for a given data set.
- We need a way to **adjust weights** (aka parameters) in order to gradually minimize loss
 - Backpropagation + optimization algorithm

1974 – Backpropagation



Paul Werbos

Artificial Intelligence pioneer
IEEE Neural Network Pioneer Award

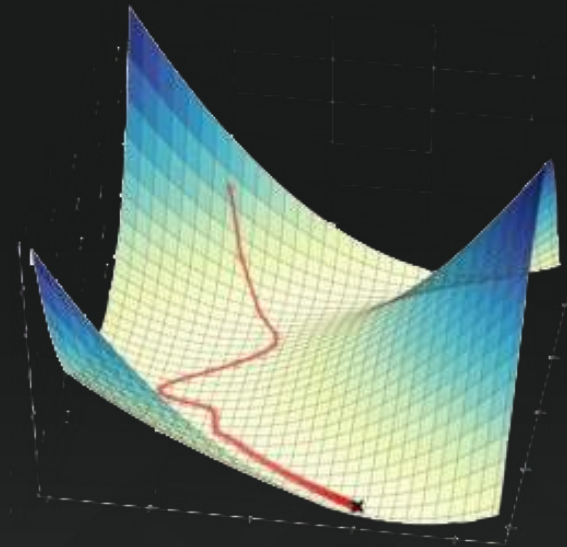


The back-propagation algorithm acts as an **error correcting** mechanism at each neuron level, thereby helping the network to learn effectively.

Stochastic Gradient Descent (SGD)

*Imagine you stand on top of a mountain with skis strapped to your feet. You want to **get down** to the valley **as quickly as possible**, but there is fog and you can only see your immediate surroundings. How can you get down the mountain as quickly as possible? You look around and **identify the steepest path down**, **go down** that path for **a bit**, again look around and find the new steepest path, go down that path, and **repeat**—this is exactly what gradient descent does.*

Tim Dettmers
University of
Lugano
2015

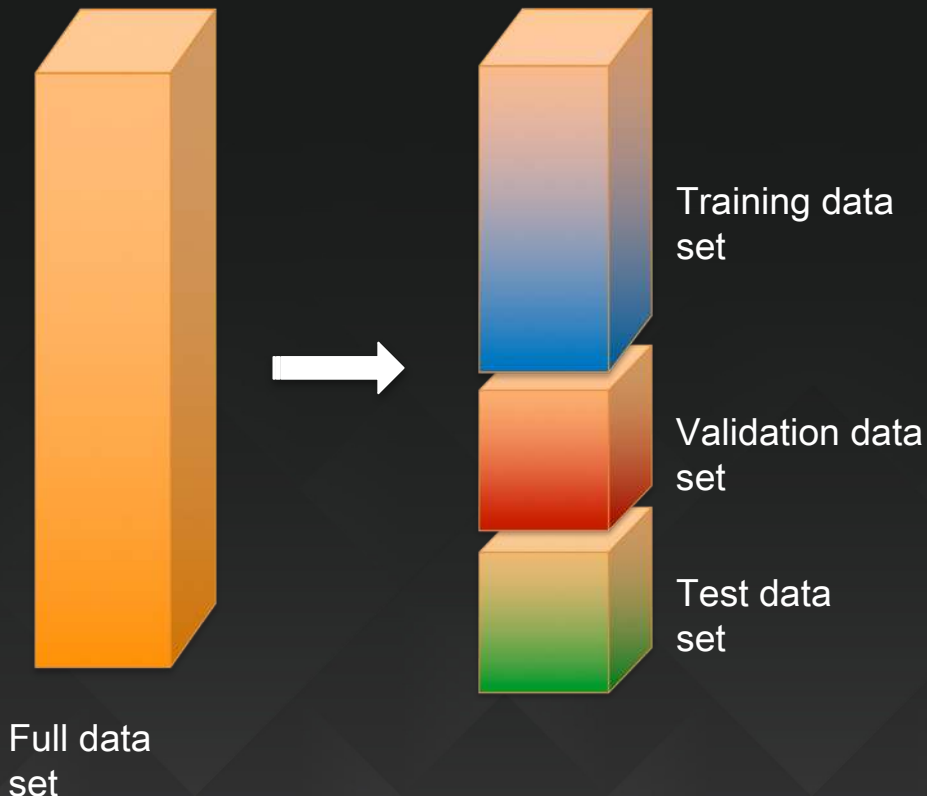


The « step size » is called
the **learning rate**

There is such a thing as « learning too well »

- If a network is **large enough** and given **enough time**, it will **perfectly** learn a data set (universal approximation theorem).
- But what about **new samples**? Can it also predict them correctly?
- Does the network **generalize** well or not?
- To prevent **overfitting**, we need to know when to **stop** training.
- The training data set is **not** enough.

Data sets

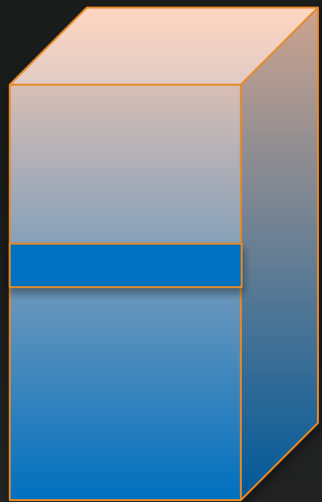


Training set: this data set is used to adjust the weights on the neural network.

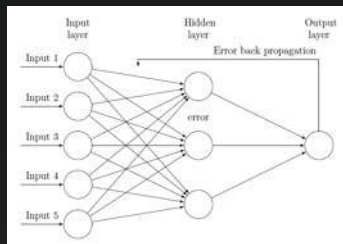
Validation set: this data set is used to minimize overfitting. You're not adjusting the weights of the network with this data set, you're just verifying that any increase in accuracy over the training data set actually yields an increase in accuracy over a data set that has not been shown to the network before.

Testing set: this data set is used only for testing the final weights in order to benchmark the actual predictive power of the network.

Training

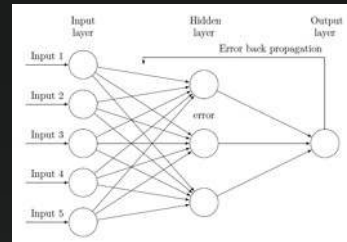
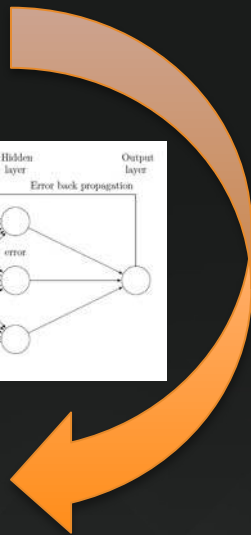


Training data
set



Trainin
g

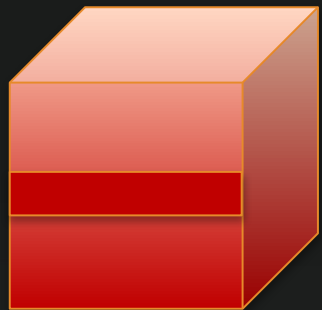
**Number of
epochs
Batch size
Learning rate**



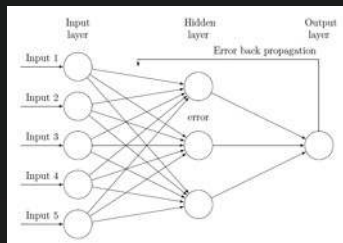
Trained
neural
network

**Hyper
parameters**

Validation



Validation Data
set



Trained
neural
network



**Validation
accuracy**

**Prediction
at the end
of each
epoch**

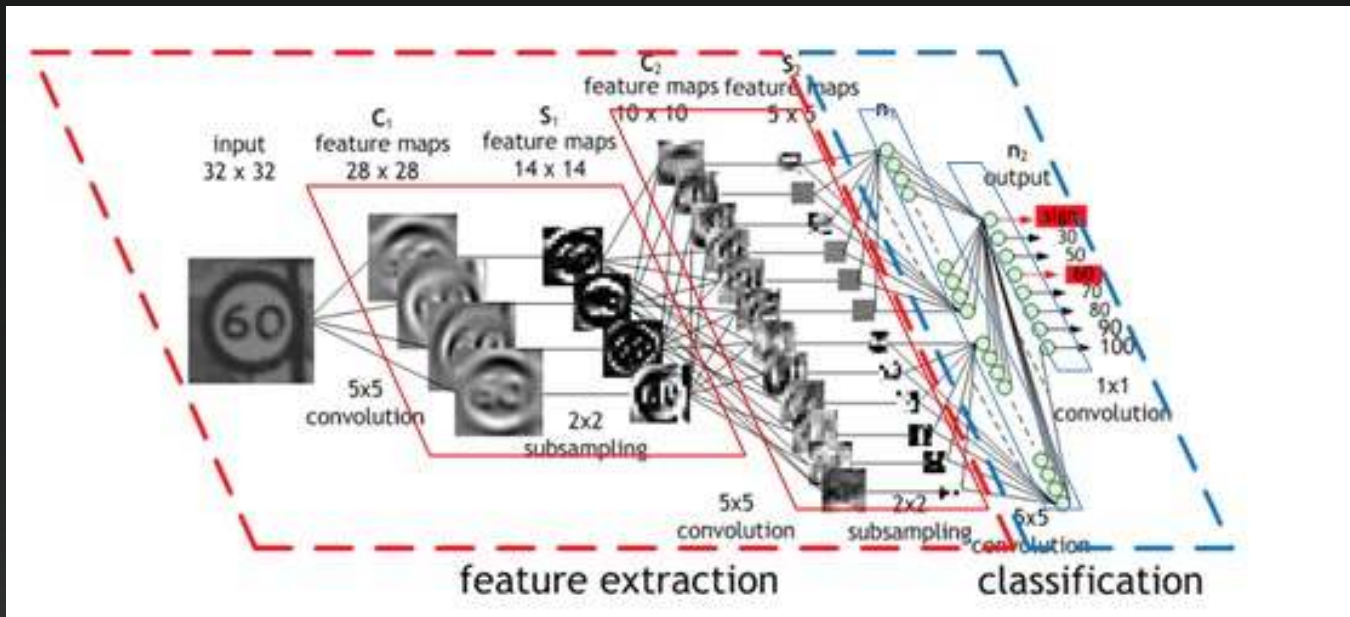
**Stop training when validation accuracy stops
increasing**

**Saving parameters at the end of each epoch is
a good idea :)**

Convolutional Neural Networks

Le Cun, 1998: handwritten digit recognition, 32x32 pixels

Feature extraction and **downsampling** allow smaller networks



<https://devblogs.nvidia.com/parallelforall/deep-learning-nutshell-core-concepts/>



Apache MXNet

Apache MXNet



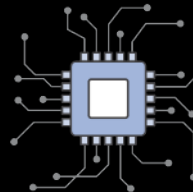
Programmable

Simple syntax,
multiple languages



Portable

Highly efficient
models for mobile
and IoT



High Performance

Near linear scaling
across hundreds of GPUs



Most Open

Accepted into the
Apache Incubator



Best On AWS

Optimized for
deep learning on
AWS

Imperative Programming

```
import numpy as np
a = np.ones(10)
b = np.ones(10) * 2
c = b * a
d = c + 1
```

Easy to tweak
in Python

PRO

S

- Straightforward and flexible.
- Take advantage of language native features (loop, condition, debugger).
- E.g. Numpy, Matlab, Torch, ...

CON

S

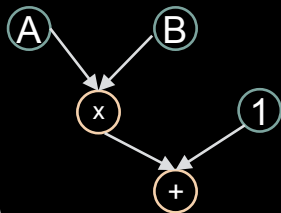
- Hard to optimize

Declarative Programming

```
A = Variable('A')  
B = Variable('B')  
C = B * A  
D = C + 1  
f = compile(D)  
d = f(A=np.ones(10),
```

```
B=np.ones(10)*2)
```

C can share memory with D
because C is deleted later



PRO

S

- More chances for optimization
- Cross different languages
- E.g. TensorFlow, Theano, Caffe

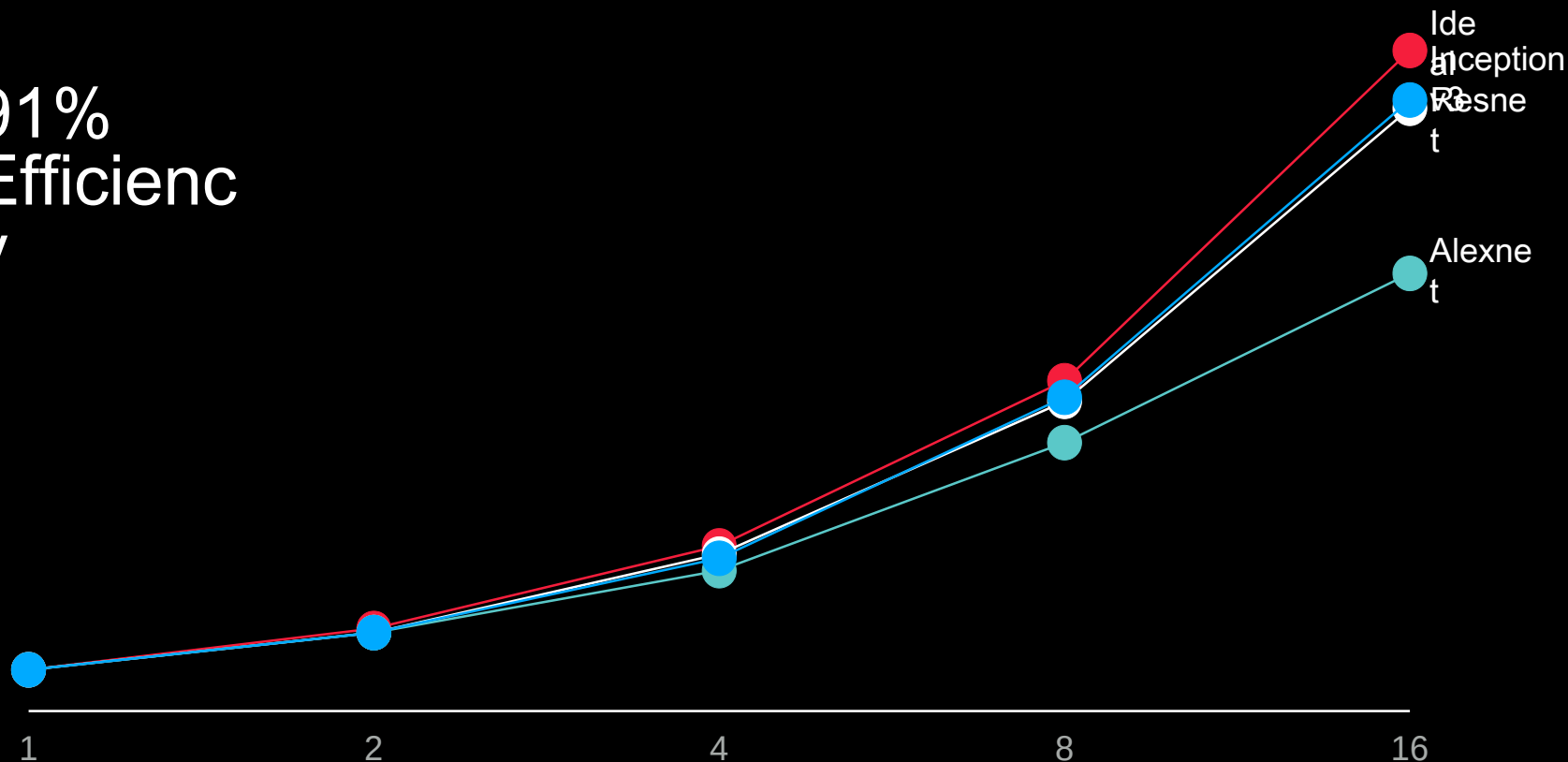
CON

S

- Less flexible

Multi-GPU Scaling With MXNet

91%
Efficiency



Input

Output

`mx.model.FeedForward`

`model.fit`

`mx.sym.SoftmaxOutput`



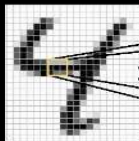
Speech



$\times =$

`mx.sym.Activation(data, act_type="xxxx")`

`mx.sym.FullyConnected(data, num_hidden=128)`



$\times =$

`mx.sym.Convolution(data, kernel=(5,5), num_filter=20)`



`mx.sym.Pooling(data, pool_type="max", kernel=(2,2),`

`stride=(2,2)`



\oplus

`lstm.lstm_unroll(num_lstm_layer, seq_len, len, num_hidden, num_embed)`

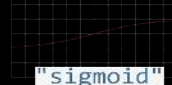


$\cos(w, queen) = \cos(w, king) - \cos(w, man) + \cos(w, woman)$

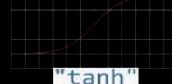
`mx.symbol.Embedding(data, input_dim, output_dim = k)`



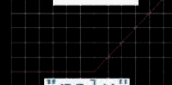
"sigmoid"



"tanh"



"relu"



"softrelu"



Let's code

1. Our first network
2. Using pre-trained models
3. Learning from scratch (MNIST, MLP, CNN)
4. Learning and fine-tuning (CIFAR-10, Resnext-101)
5. Using Apache MXNet as a Keras backend
6. Fine-tuning with Keras/MXNet (CIFAR-10, Resnet-50)
7. A quick word about Sockeye

<http://mxnet.io>

<http://medium.com/@julsimon>

<https://github.com/juliensimon/aws/tree/master/mxnet>



Thank you!

<http://aws.amazon.com/evangelists/julien-simon>
@julsimon