

DEV DAY



DEV DAY

MLT 2

Deep Learning: concepts, algorithms & use cases

Julien Simon
Global Evangelist, AI & Machine Learning
Amazon Web Services

[@julsimon](#)



What to expect

An introduction to Deep Learning theory

- Neurons & Neural Networks

- The Training Process

- Backpropagation

- Optimizers

Common network architectures and use cases

- Convolutional Neural Networks

- Recurrent Neural Networks

- Long Short Term Memory Networks

- Generative Adversarial Networks

Getting started



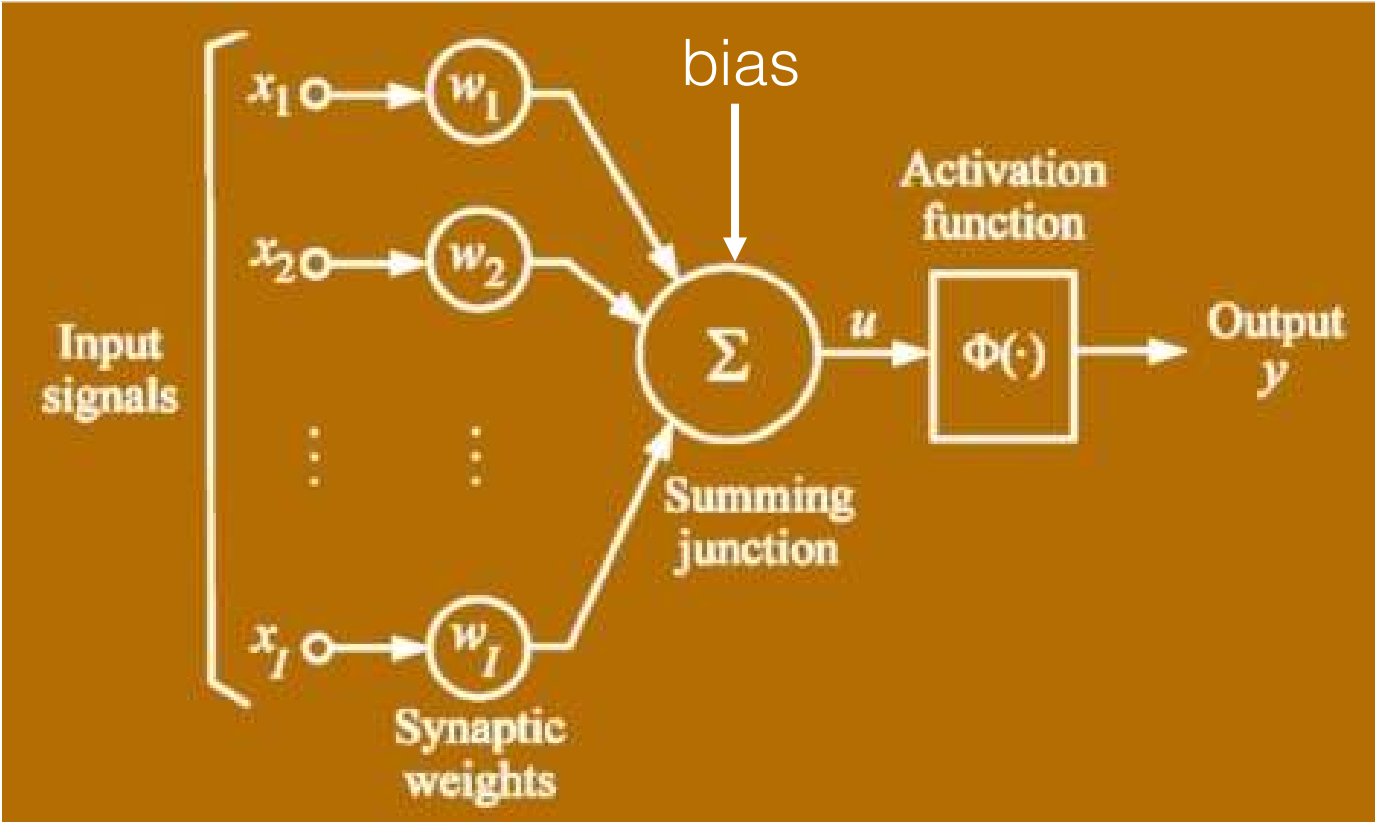
DEV DAY

An introduction to Deep Learning theory



© 2020, Amazon Web Services, Inc. or its affiliates. All rights reserved.

The neuron



$$\sum_{i=1}^I x_i * w_i + b = u$$

Activation functions

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]		$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$

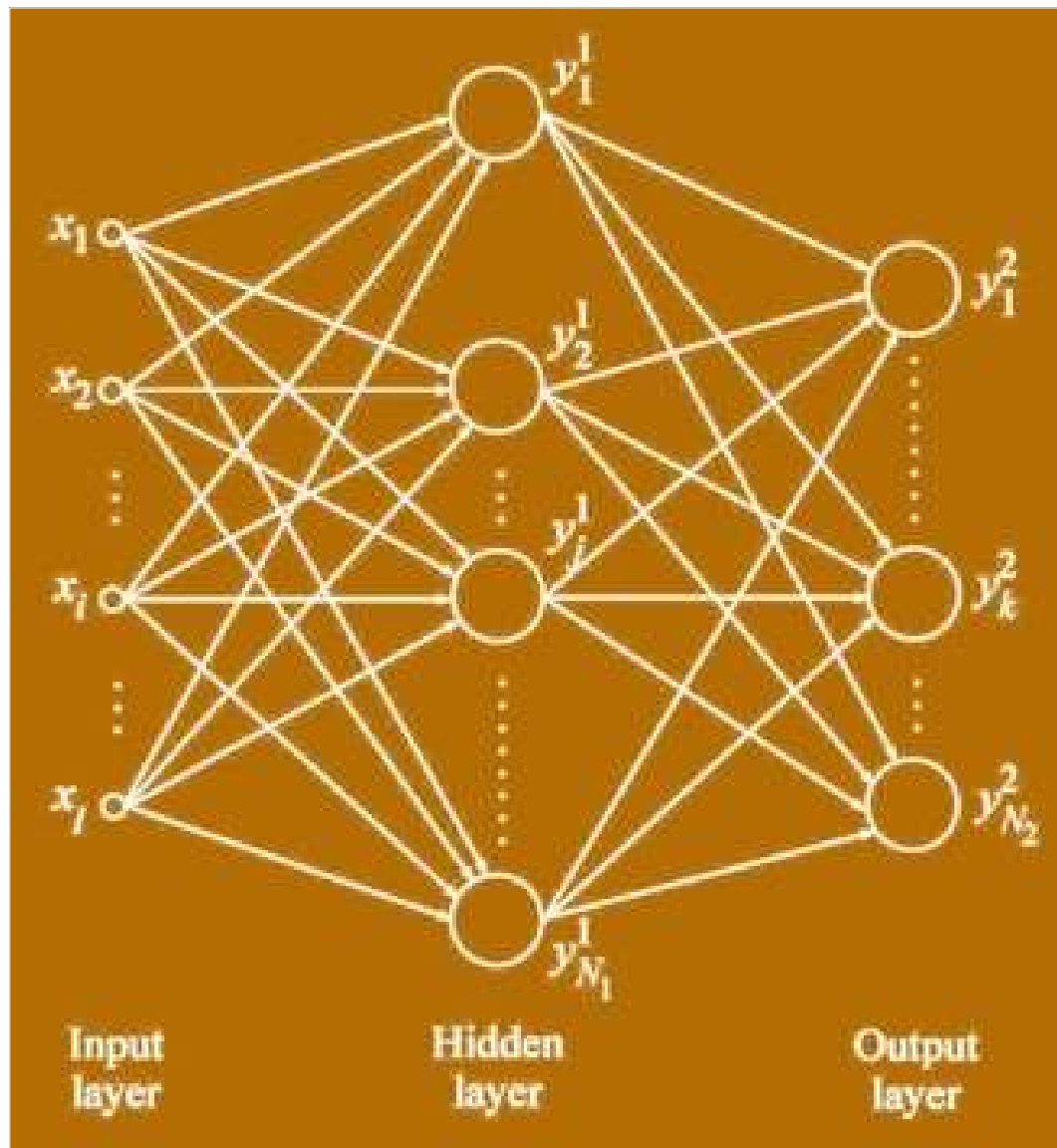
Source: Wikipedia

”Multiply and Accumulate”



Neural networks

Building a simple classifier



$$X = \begin{bmatrix} X_{11}, X_{12}, \dots, X_{1I} \\ X_{21}, X_{22}, \dots, X_{2I} \\ \dots \dots \dots \\ X_{m1}, X_{m2}, \dots, X_{mI} \end{bmatrix}$$

features

m samples

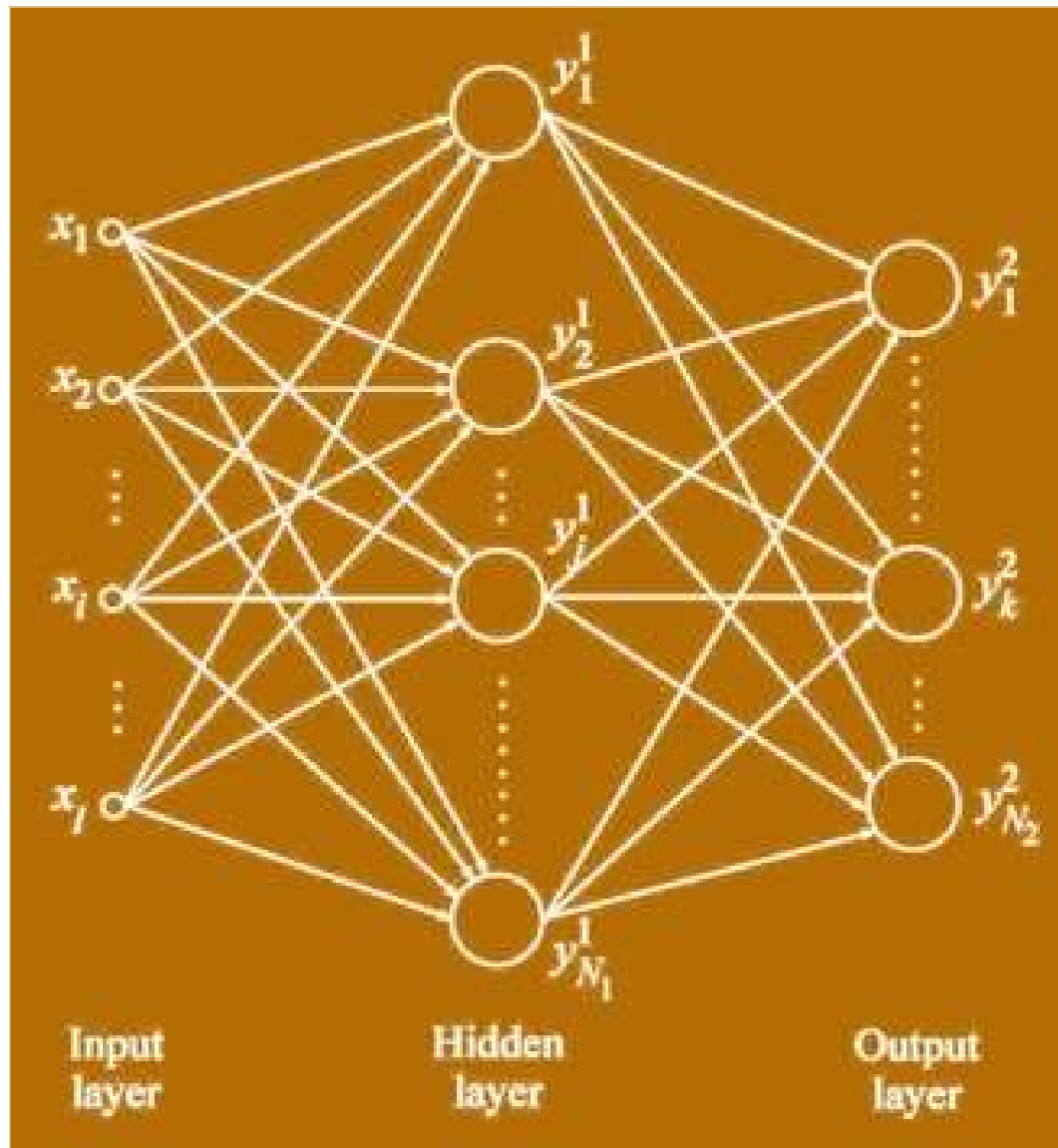
$$y = \begin{bmatrix} 2 \\ 0 \\ \dots \\ 4 \end{bmatrix} \quad \begin{bmatrix} 0,0,1,0,0,\dots,0 \\ 0,0,0,0,0,\dots,0 \\ \dots \\ 0,0,0,0,1,\dots,0 \\ 0 \end{bmatrix}$$

m labels, N_2 categories

One-hot encoding

Neural networks

Building a simple classifier



$$X = \begin{bmatrix} X_{11} & X_{12} & \dots & X_{1I} \\ X_{21} & X_{22} & \dots & X_{2I} \\ \dots & \dots & \dots & \dots \\ X_{m1} & X_{m2} & \dots & X_{mI} \end{bmatrix}$$

features

m samples

$$y = \begin{bmatrix} 2 \\ 0 \\ \dots \\ 4 \end{bmatrix}$$

$$= \begin{bmatrix} 0, 0, 1, \dots, \\ 0, 0, 0, \dots, \\ \dots \\ 0, 0, 0, \dots, \\ 0 \end{bmatrix}$$

m labels, N_2 categories

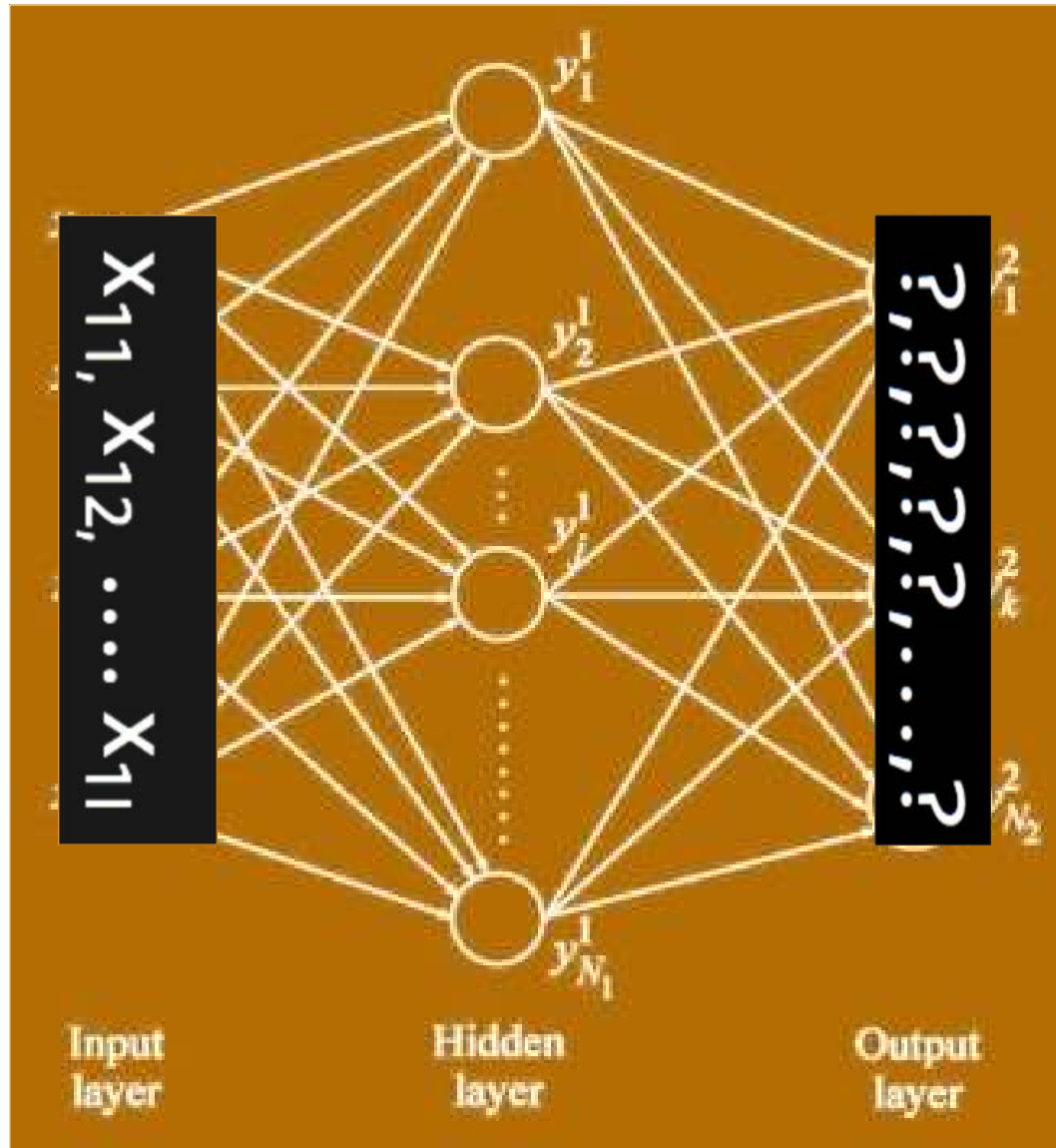
One-hot encoding

$$\text{Accuracy} =$$

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

Neural networks

Building a simple classifier



Initially, the network will **not** predict correctly
 $f(X_1) = Y'_1$

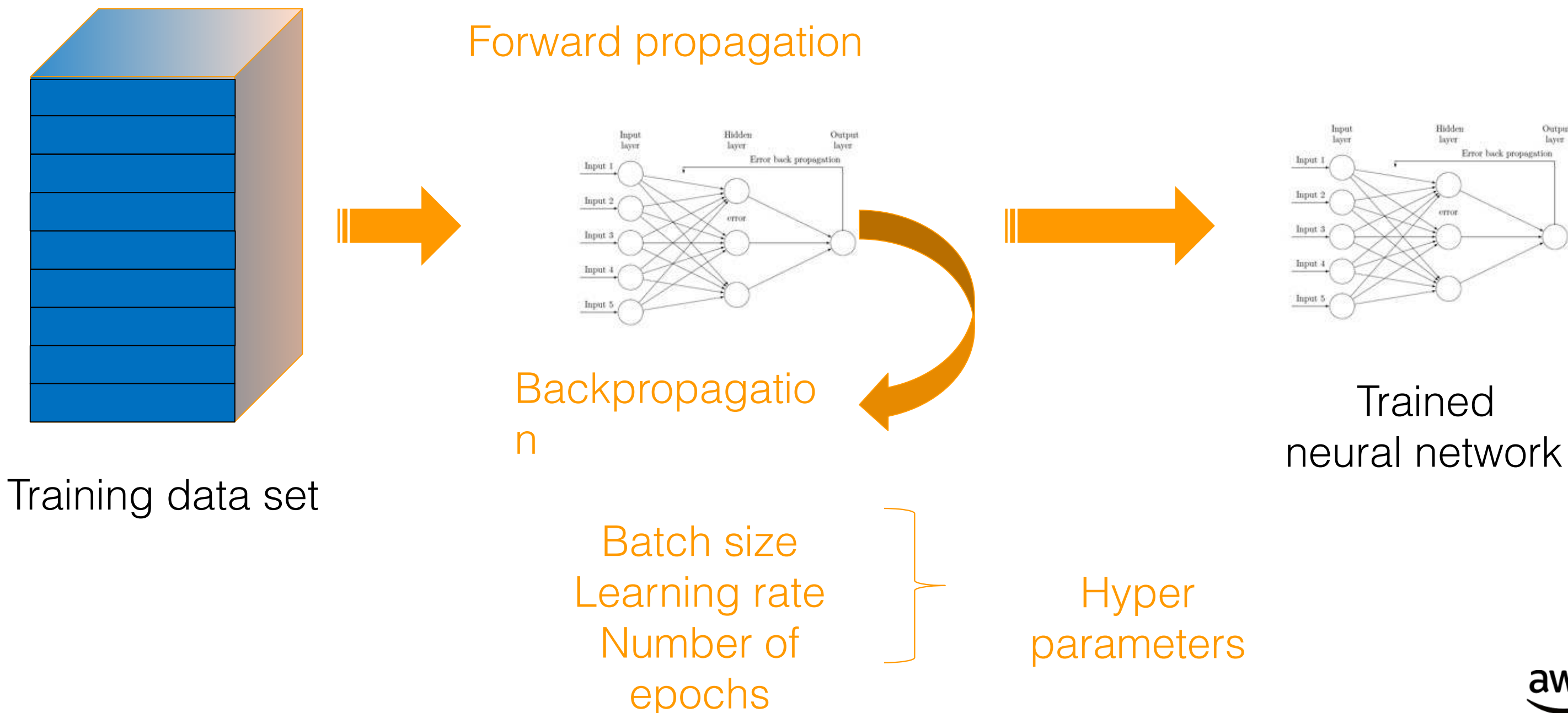
A **loss function** measures the difference between the **real label** Y_1 and the **predicted label** Y'_1
 $\text{error} = \text{loss}(Y_1, Y'_1)$

For a **batch** of samples:

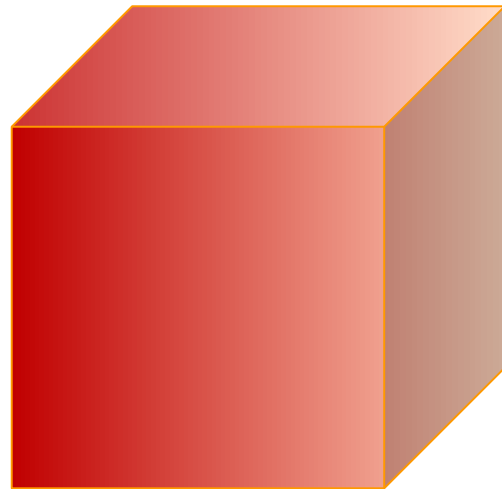
$$\sum_{i=1}^{\text{batch size}} \text{loss}(Y_i, Y'_i) = \text{batch error}$$

The purpose of the training process is to **minimize error** by gradually **adjusting weights**.

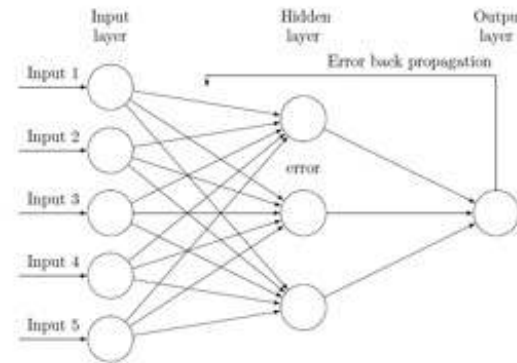
Mini-batch Training



Validation



Validation data set
(also called dev
set)



Neural network
in training

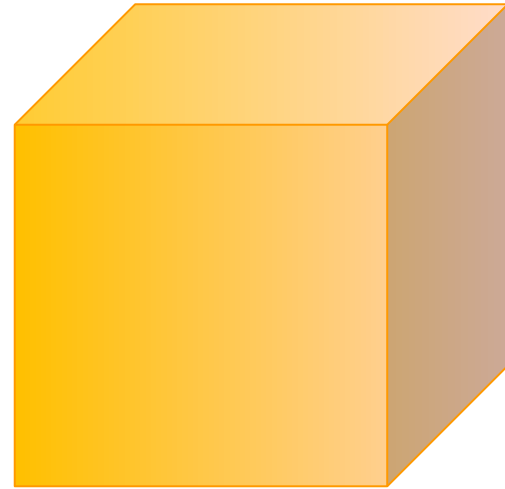


Validation
accuracy

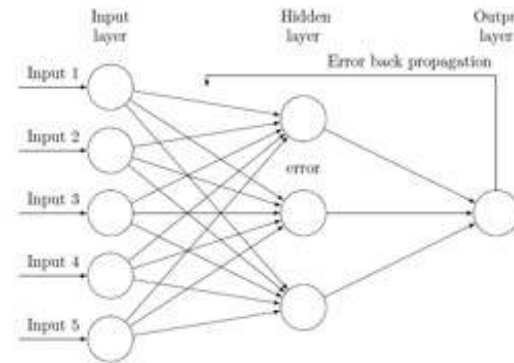
Prediction at
the end of
each epoch

*This data set must have the same distribution as real-life samples,
or else validation accuracy won't reflect real-life accuracy.*

Test



Test data set



Fully trained
neural network



Test accuracy

Prediction at
the end of
experimentation

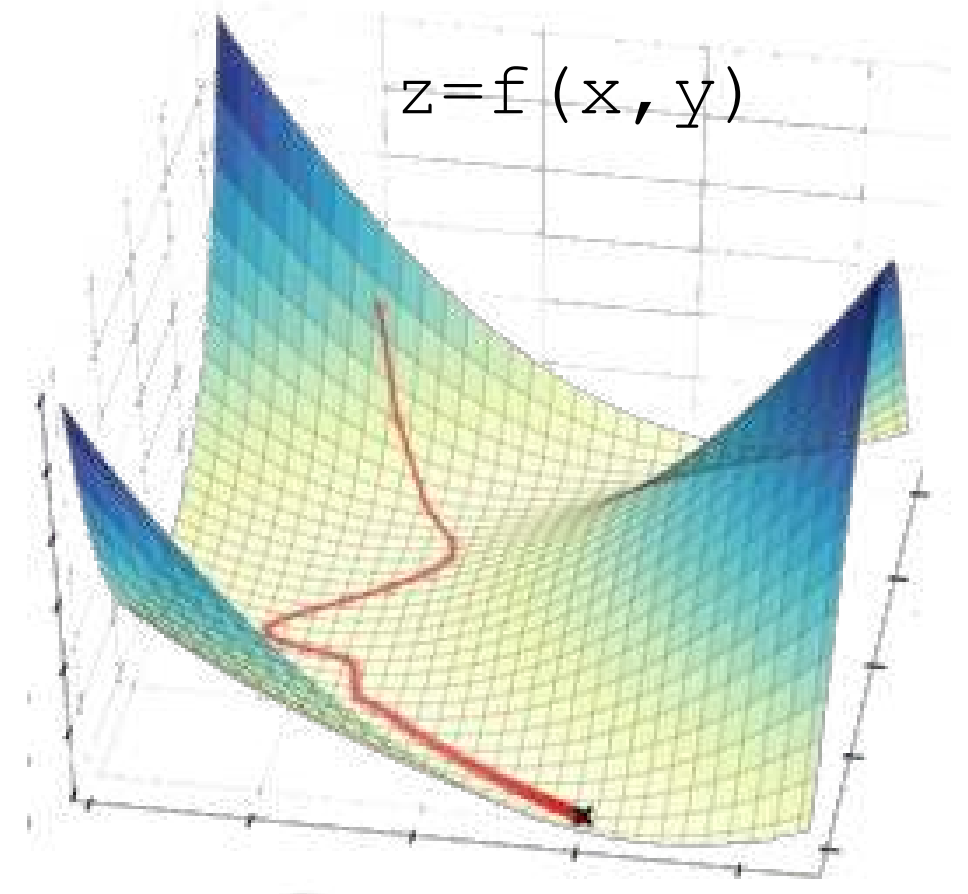
*This data set must have the same distribution as real-life samples,
or else test accuracy won't reflect real-life accuracy.*

Stochastic Gradient Descent (1951)

Imagine you stand on top of a mountain (...). 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.

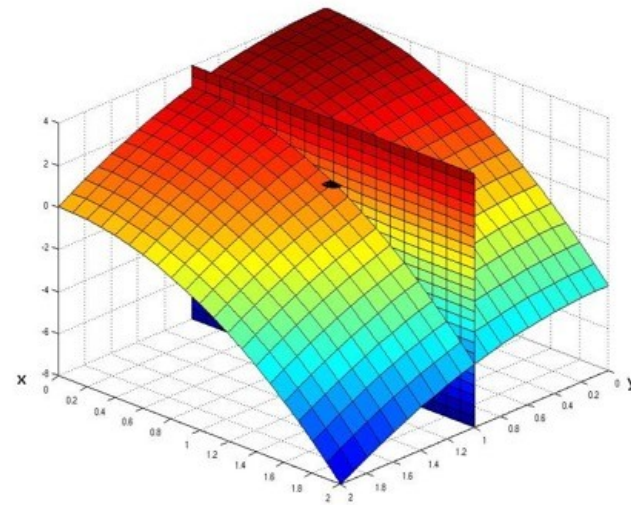
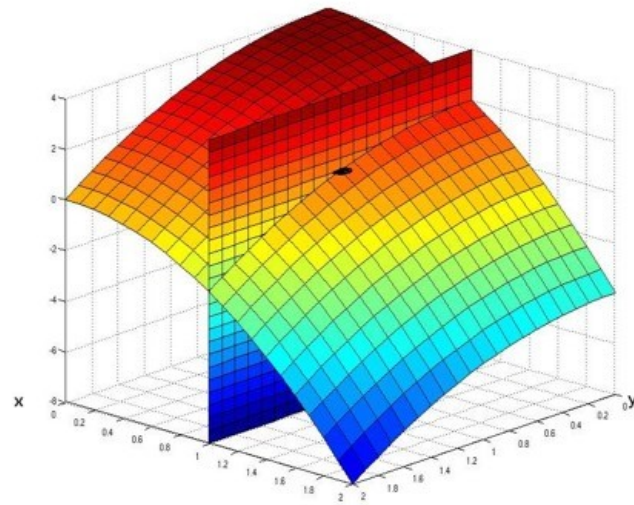
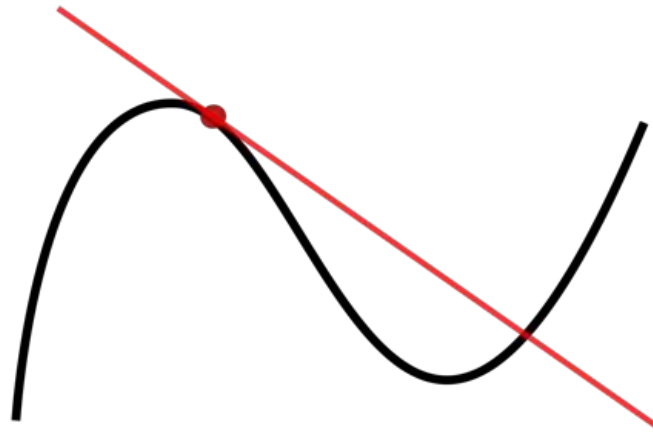
Tim Dettmers, University of Lugano, 2015



The « step size » depends on the **learning rate**

<https://devblogs.nvidia.com/parallelforall/deep-learning-nutshell-history-training/>

Finding the slope with Derivatives



Scalar-valued multivariable function

$$\nabla f(x_0, y_0, \dots) = \begin{bmatrix} \frac{\partial f}{\partial x}(x_0, y_0, \dots) \\ \frac{\partial f}{\partial y}(x_0, y_0, \dots) \\ \vdots \end{bmatrix}$$

Notation for gradient, called "nabla".

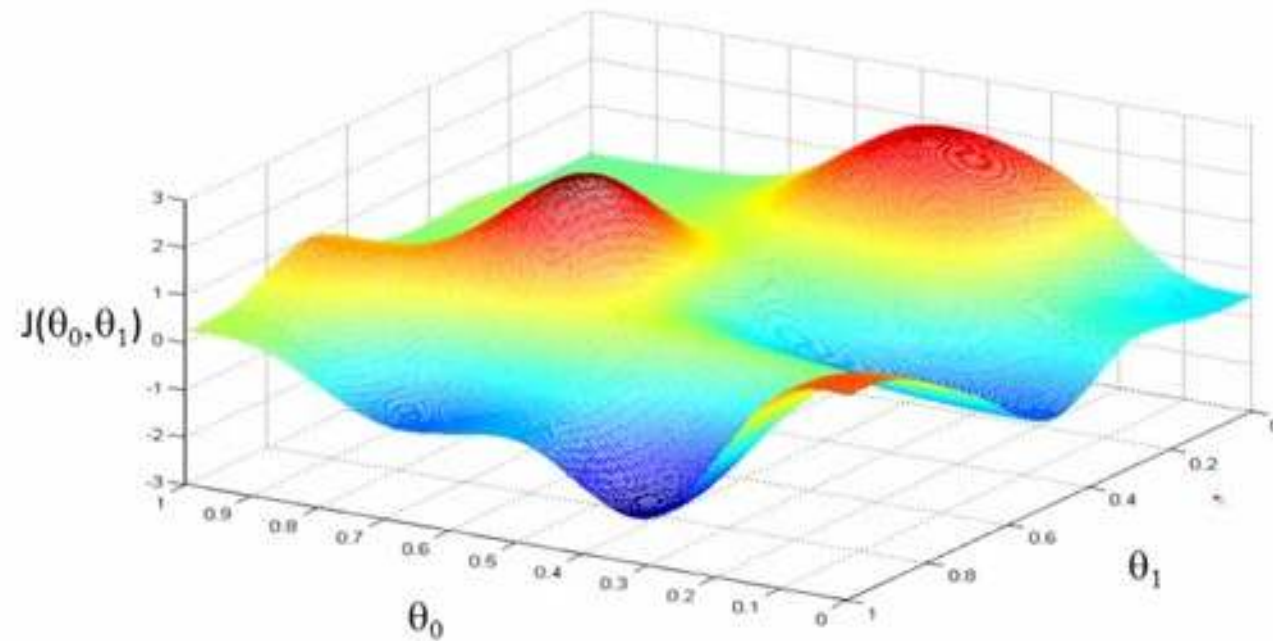
∇f takes the same type of inputs as f

∇f outputs a vector with all possible partial derivatives of f .

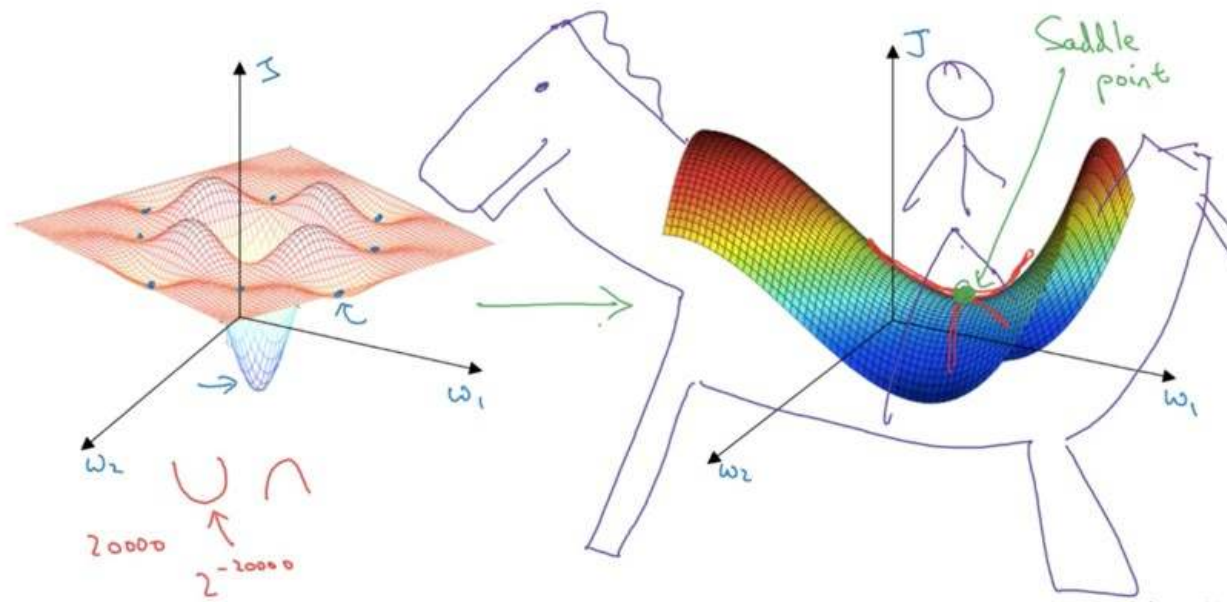
End-to-end example of computing backpropagation with partial derivatives:
<https://mattmazur.com/2015/03/17/a-step-by-step-backpropagation-example>



Local minima and saddle points



Local optima in neural networks



« Do neural networks enter and escape a series of local minima? Do they move at varying speed as they approach and then pass a variety of saddle points? Answering these questions definitively is difficult, but we present evidence strongly suggesting that the answer to all of these questions is no. »

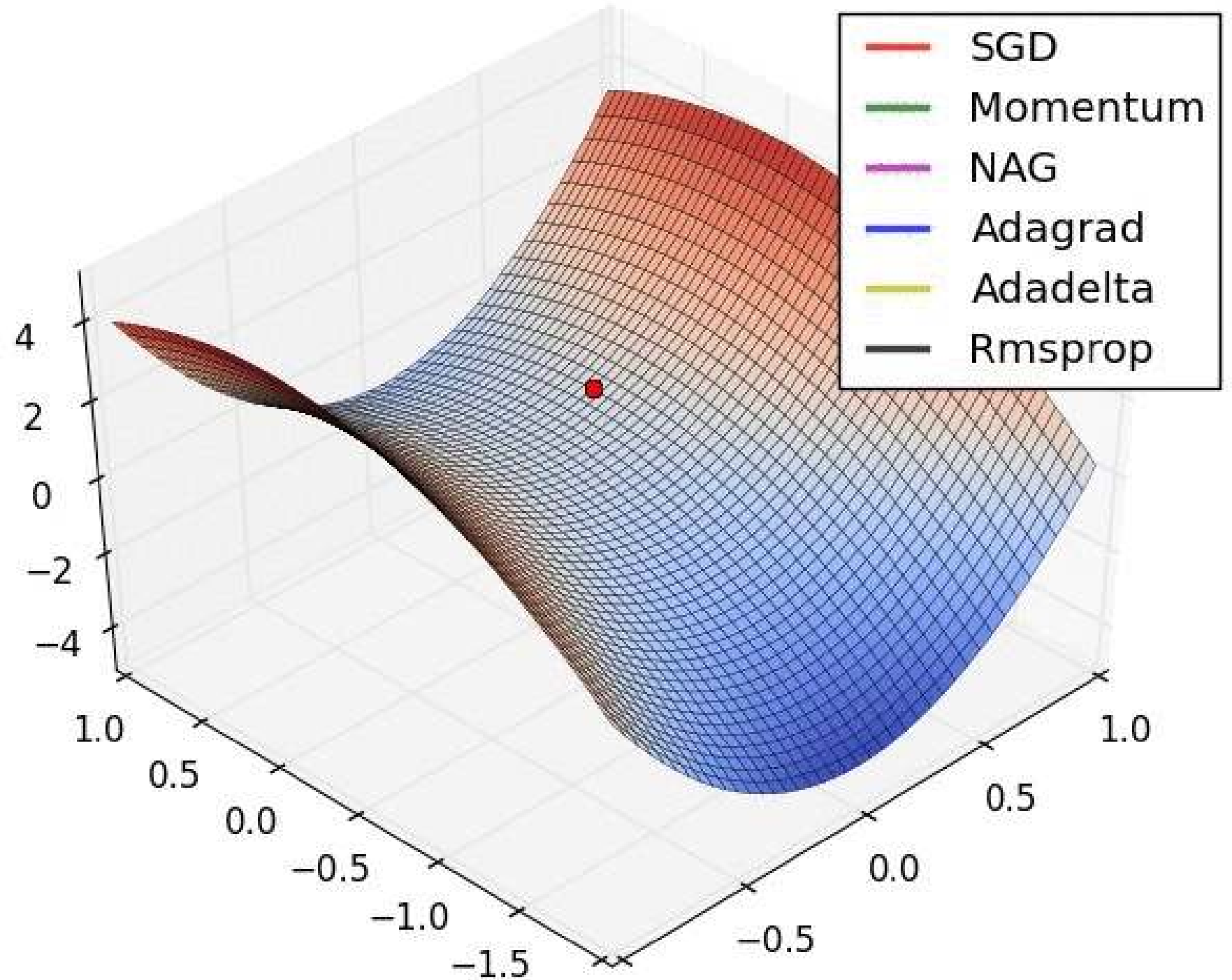
« Qualitatively characterizing neural network optimization problems », Goodfellow et al, 2015 <https://arxiv.org/abs/1412.6544>

Optimizers

SGD works remarkably well and is still widely used.

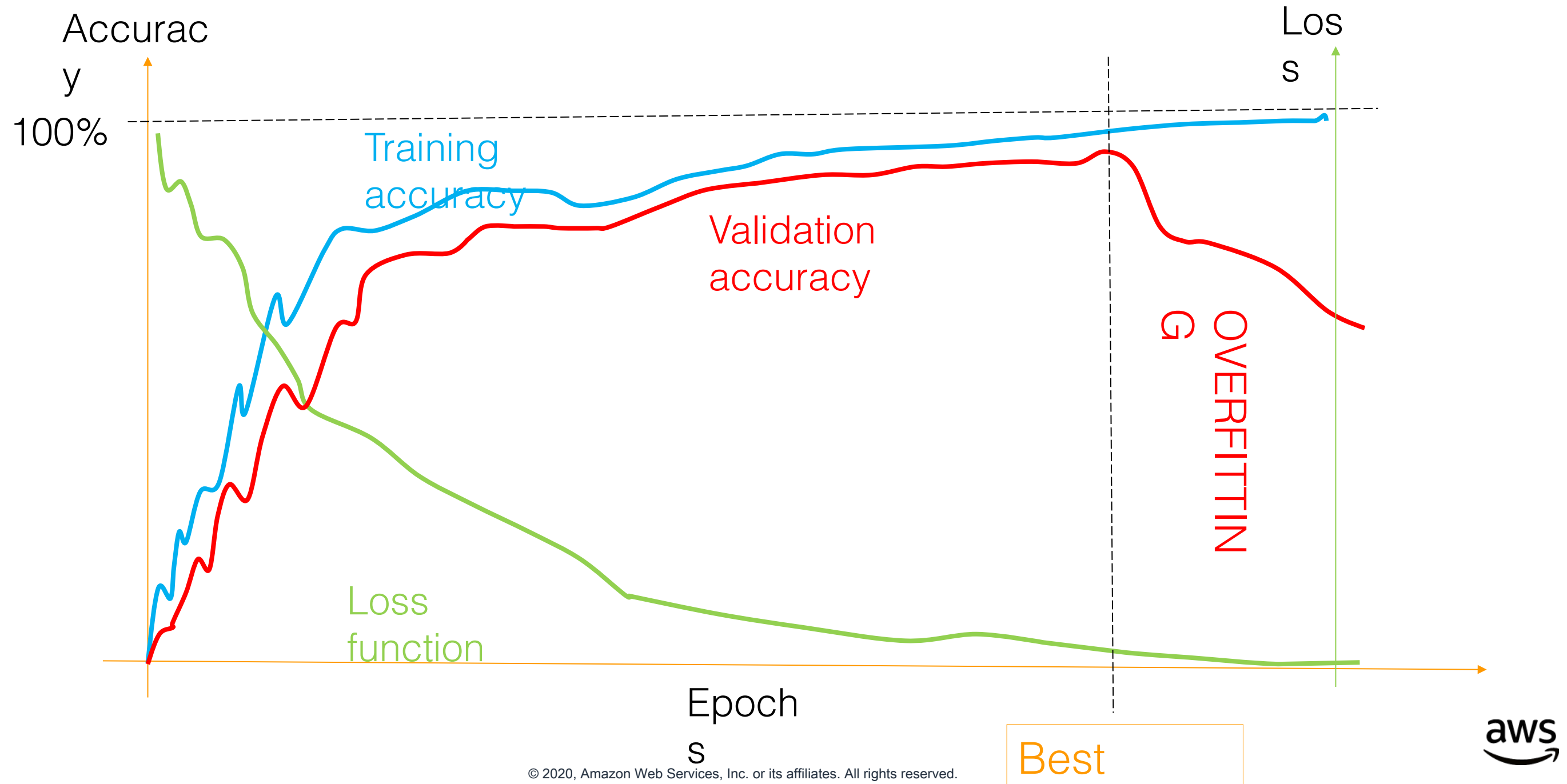
Adaptative optimizers use a variable learning rate.

Some even use a learning rate per dimension (Adam).



Summing things up

« Deep Learning ultimately is about finding a minimum that generalizes well, with bonus points for finding one fast and reliably », Sebastian Ruder



DEV DAY

Demo



Common network architectures and use cases

Fully Connected Networks are nice, but...

What if we need lots of layers in order to **extract complex features**?

The **number of parameters** increases very quickly with the number of layers

Overfitting is a constant problem

What about **large** data?

256x256 images = 65,535 input neurons ?

What about 2D/3D data ? Won't we **lose** lots of info by **flattening** it?

Images, videos, etc.

What about **sequential data**, where the order of samples is important?

Translating text

Predicting time series

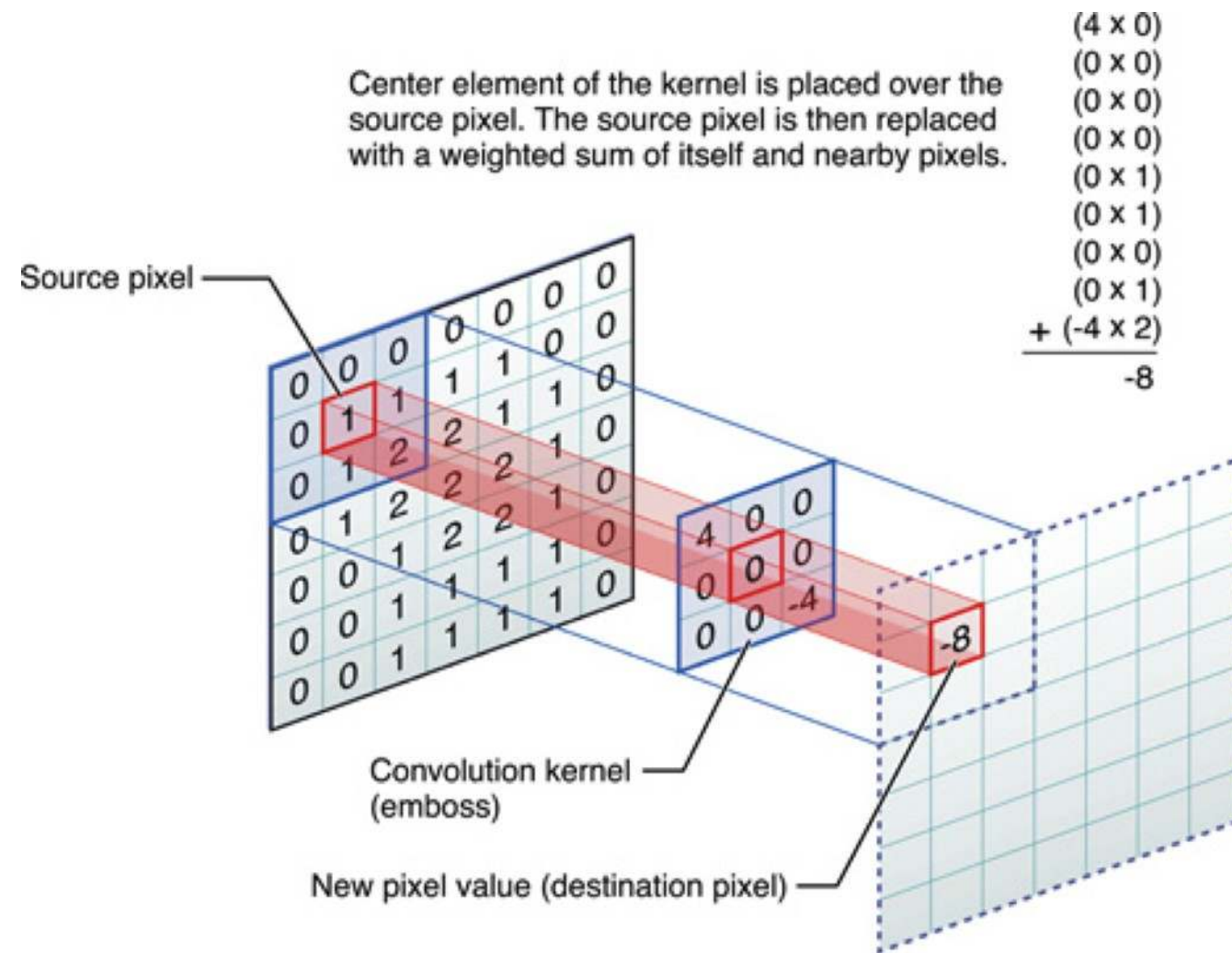
DEV DAY

Convolutional Neural Networks

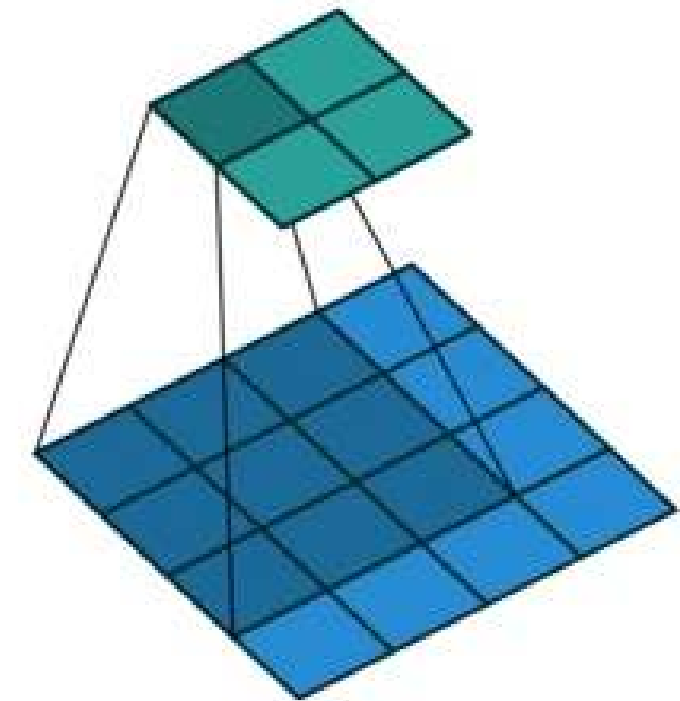


© 2020, Amazon Web Services, Inc. or its affiliates. All rights reserved.

The convolution operation



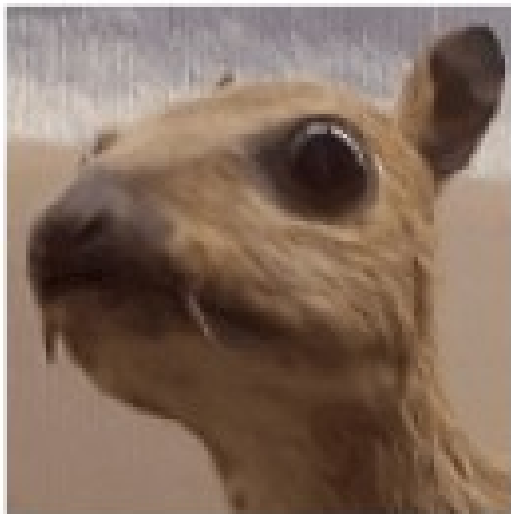
Source: <https://ikhlestov.github.io/pages/machine-learning/convolutions-types/>



Source: Theano documentation

Extracting features with convolution

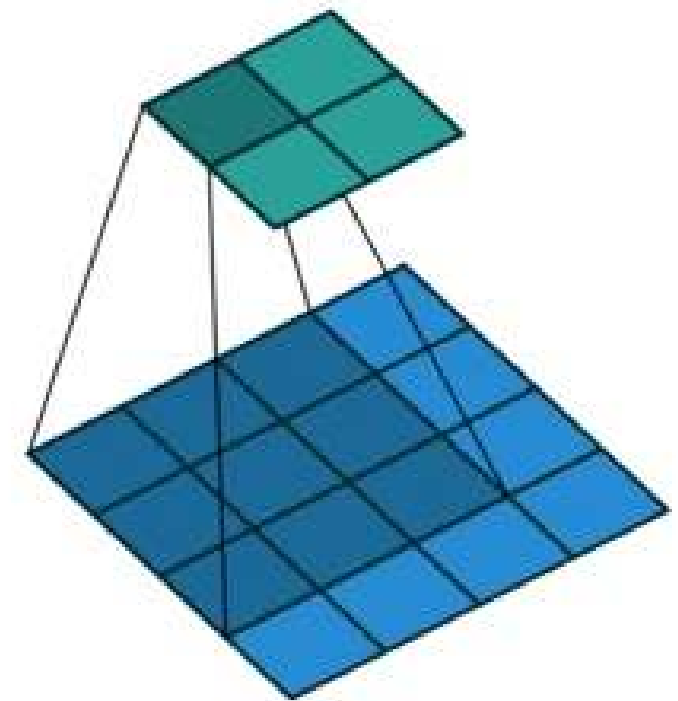
Input image



Convolution
Kernel

$$\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$$

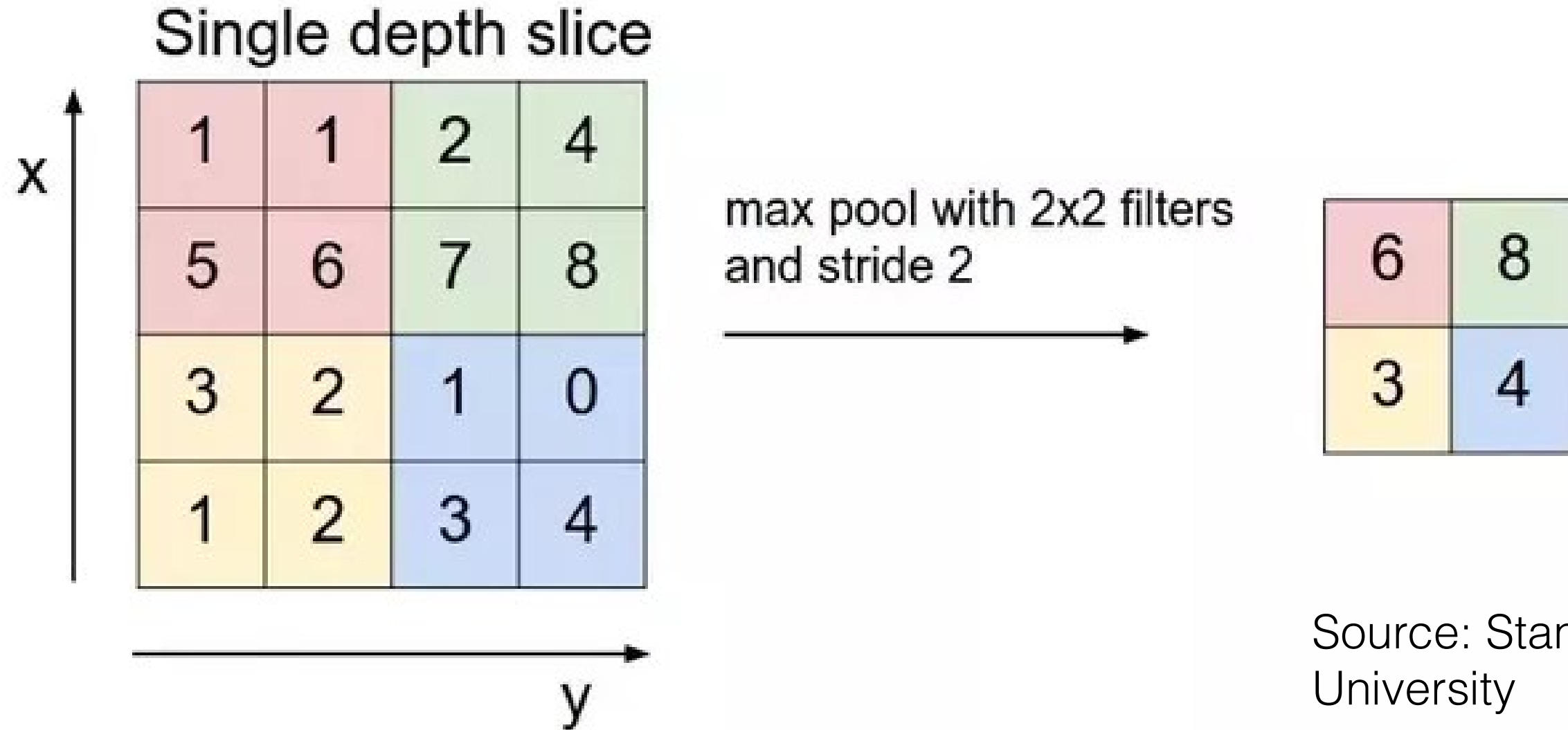
Feature map



Source: <http://timdettmers.com>

Convolution **extracts features** automatically.
Kernel parameters are **learned** during the training
process.

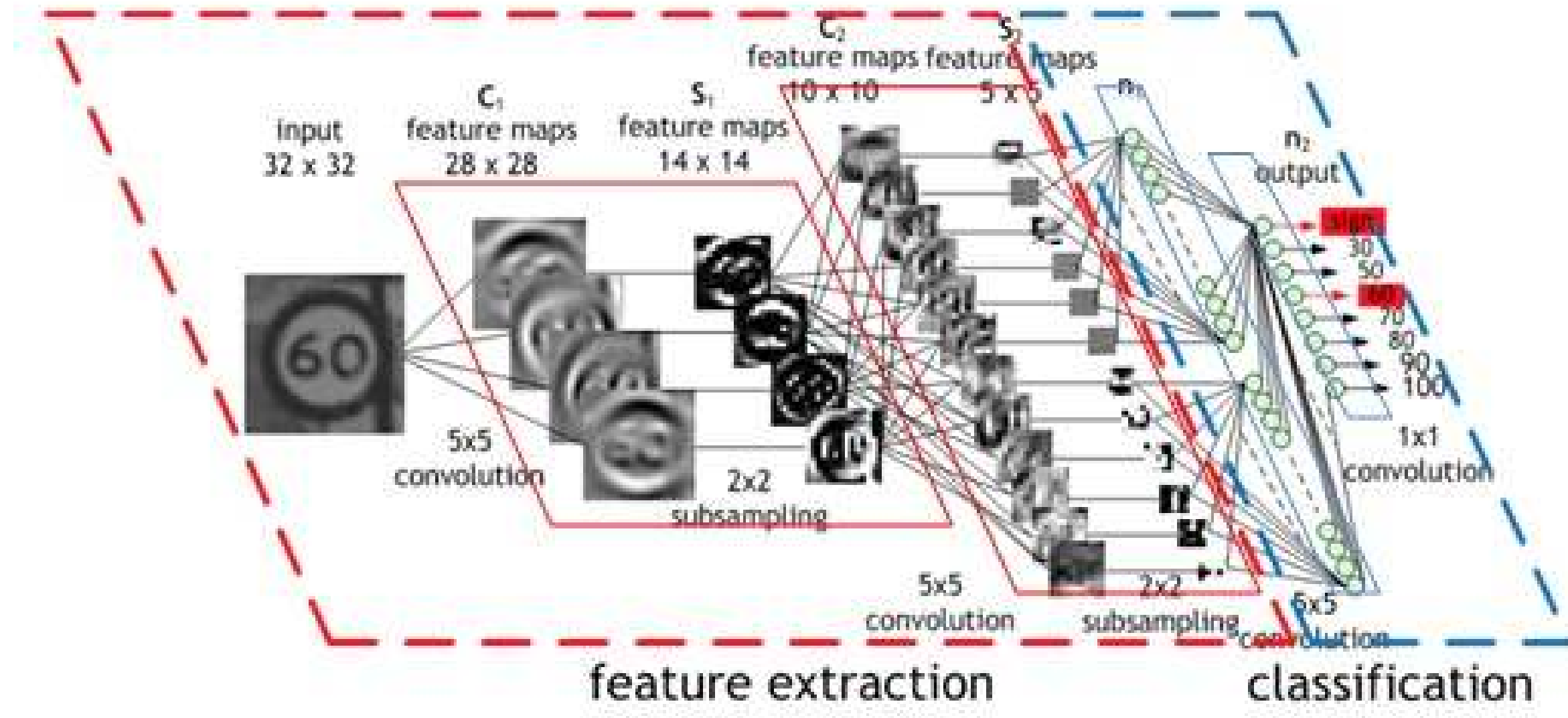
Downsampling images with pooling



Pooling shrinks images while preserving **significant** information.

Convolutional Neural Networks (CNN)

LeCun, 1998: handwritten digit recognition, 32x32 pixels



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

DEV DAY

Demo



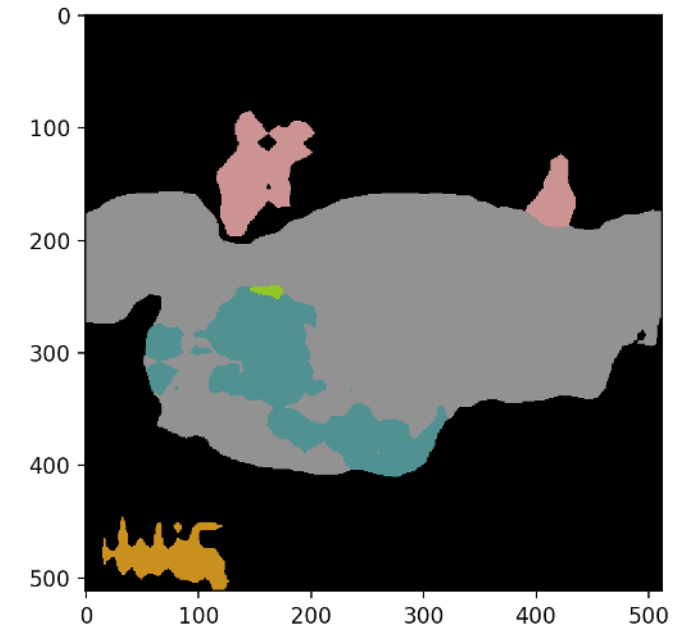
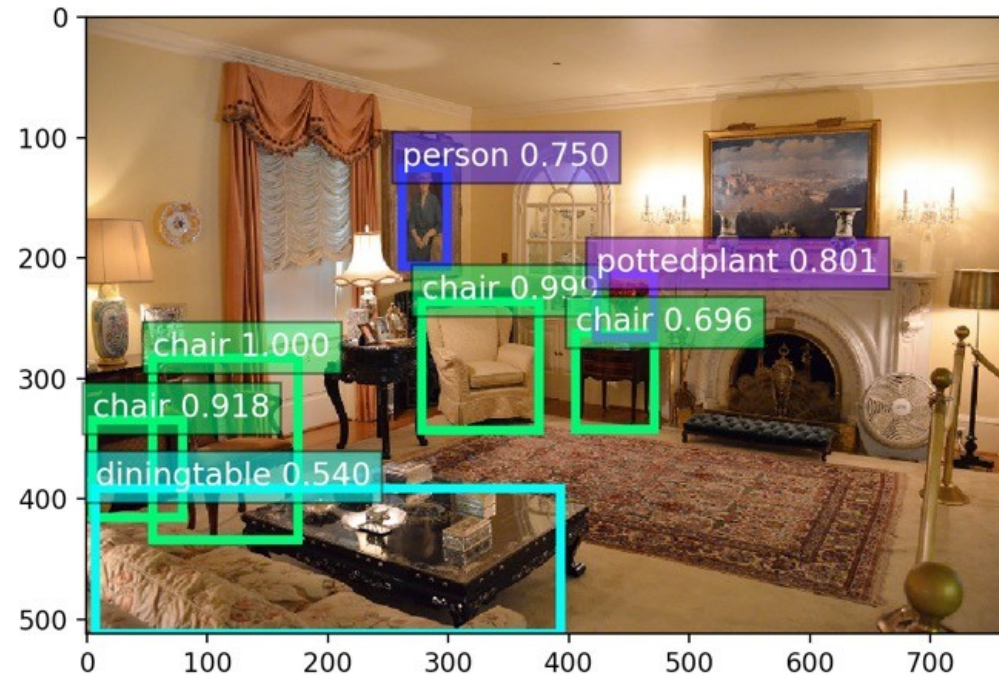
Classification, detection, segmentation



[electric_guitar],
with probability 0.671

<https://github.com/dmlc/gluon-cv>

Based on models published in 2015-2017



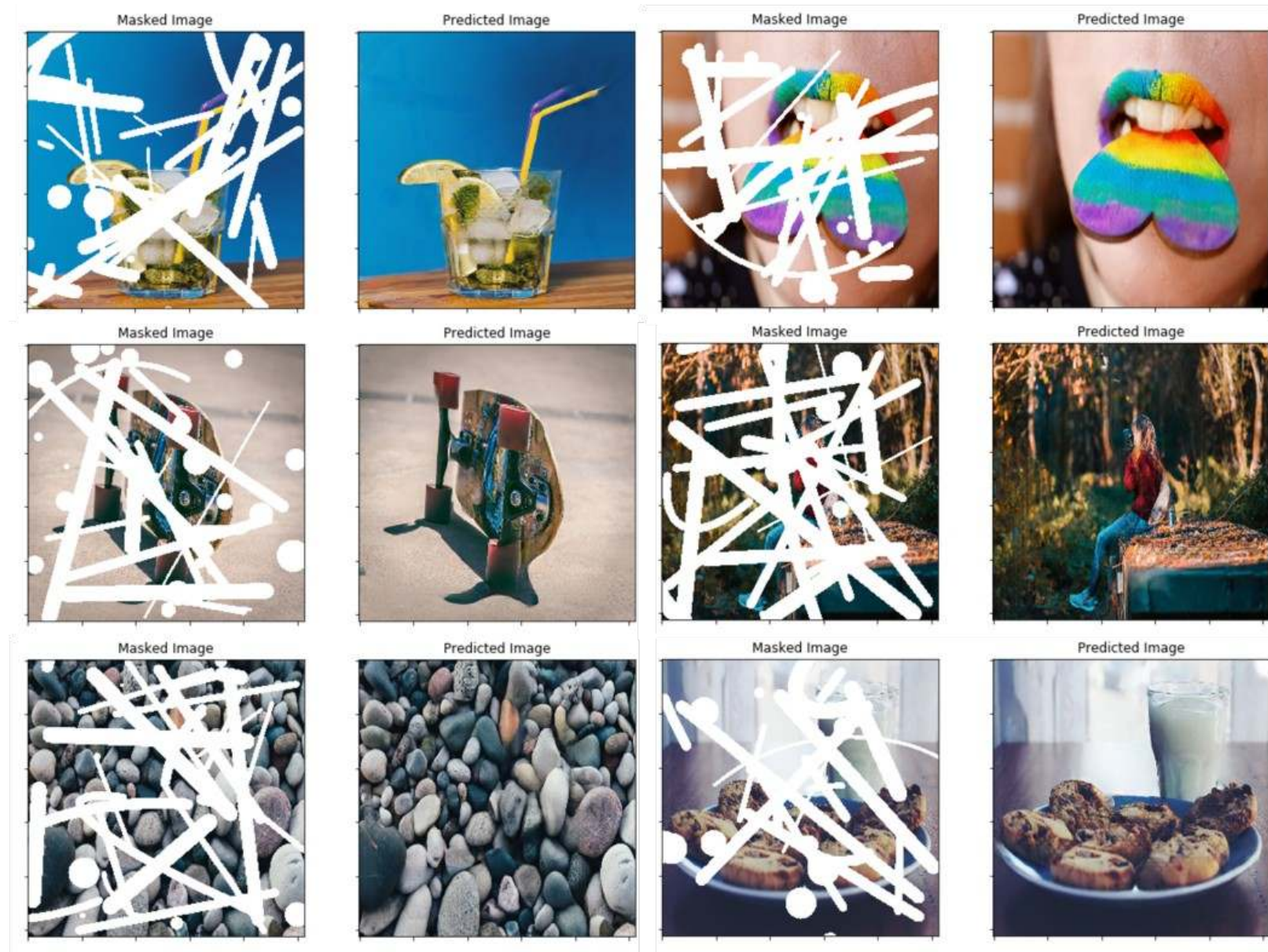
DEV DAY

Demo





Image Inpainting



<https://github.com/MathiasGruber/PConv-Keras>
<https://arxiv.org/abs/1804.07723>

April 2018

© 2020, Amazon Web Services, Inc. or its affiliates. All rights reserved.



Detectron2



<https://github.com/facebookresearch/detectron2> November 2019

Demo on the AWS Deep Learning AMI <https://www.youtube.com/watch?v=7R8-VAk0ruk>

DEV DAY

Demo



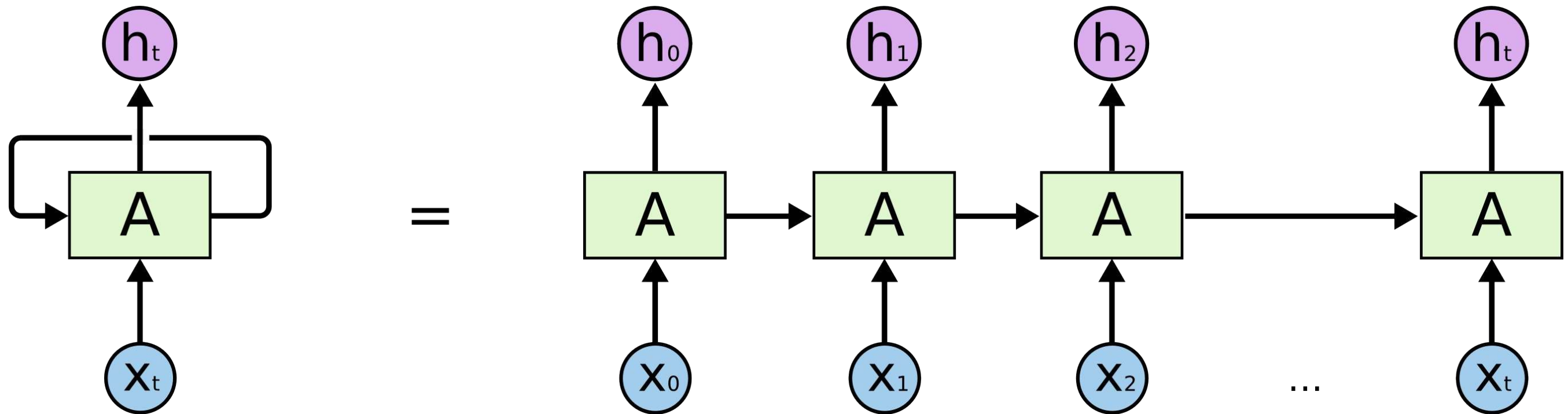
DEV DAY

Recurrent Neural Networks



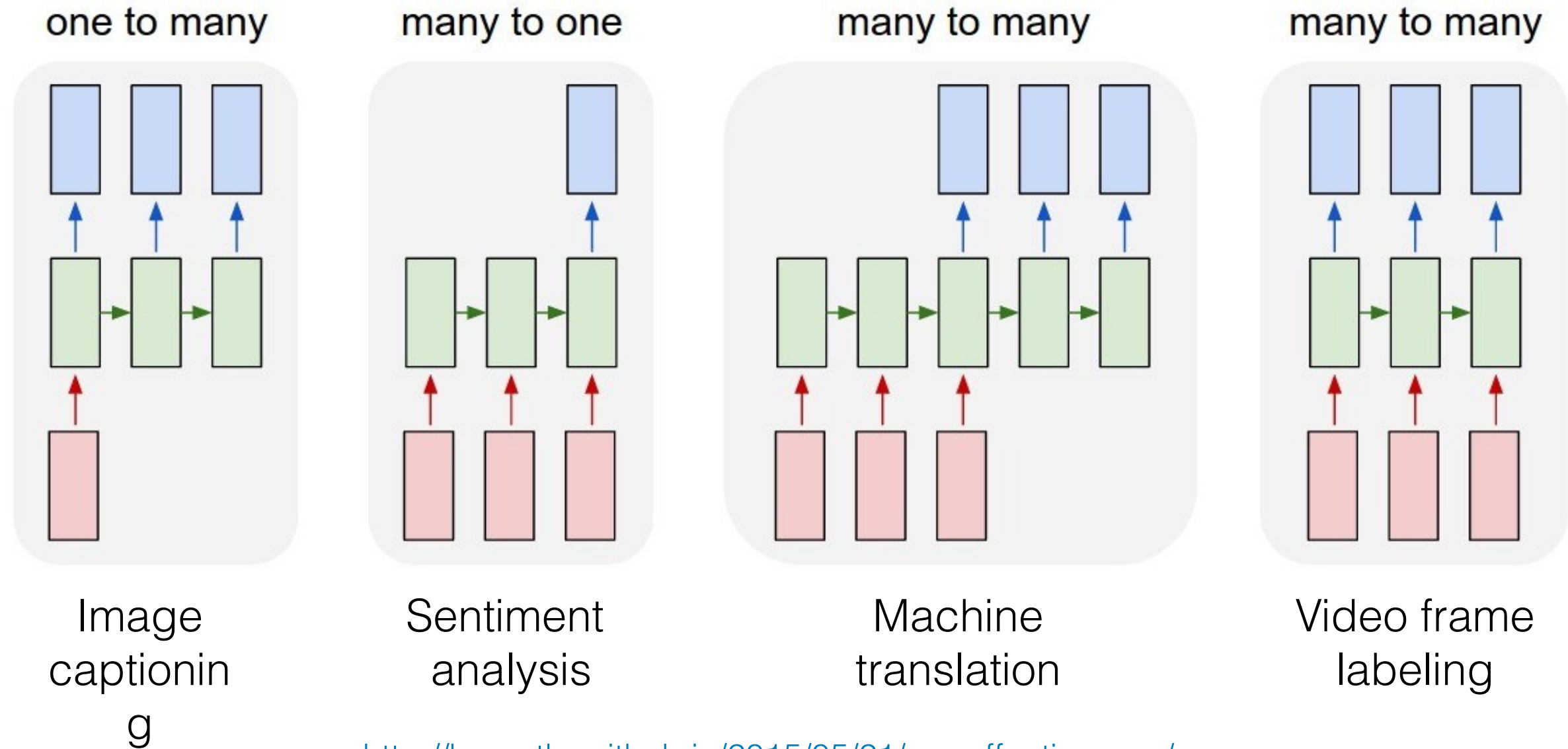
© 2020, Amazon Web Services, Inc. or its affiliates. All rights reserved.

Recurrent Neural Networks (RNN)



<http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

Recurrent Neural Networks (RNN)

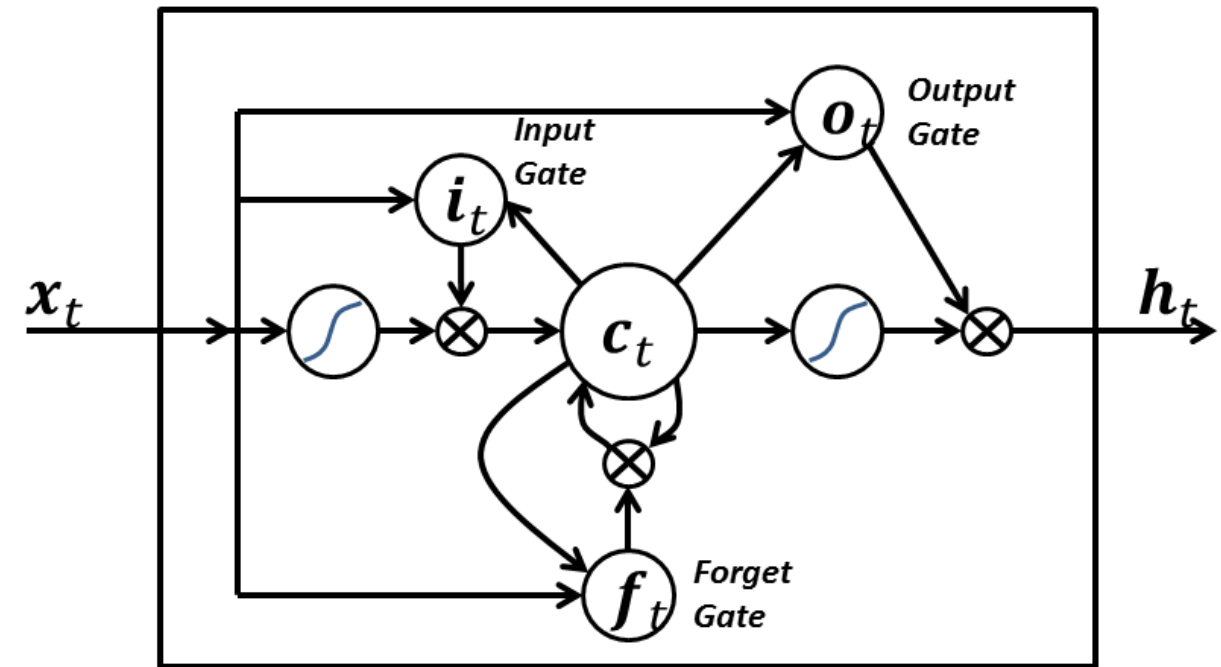


<http://karpathy.github.io/2015/05/21/rnn-effectiveness/>

Long Short Term Memory Networks (LSTM)

Hochreiter and Schmidhuber, 1997

- A LSTM neuron computes the output based on the input and a **previous state**
- LSTM neurons have « **short-term memory** »
- They do a better job than RNN at predicting longer **sequences** of data (i.e. hundreds of steps)
- Gated Recurrent Units aka **GRU** (Cho et al. 2014) are easier to train, and just as efficient.



DEV DAY

Generative Adversarial Networks

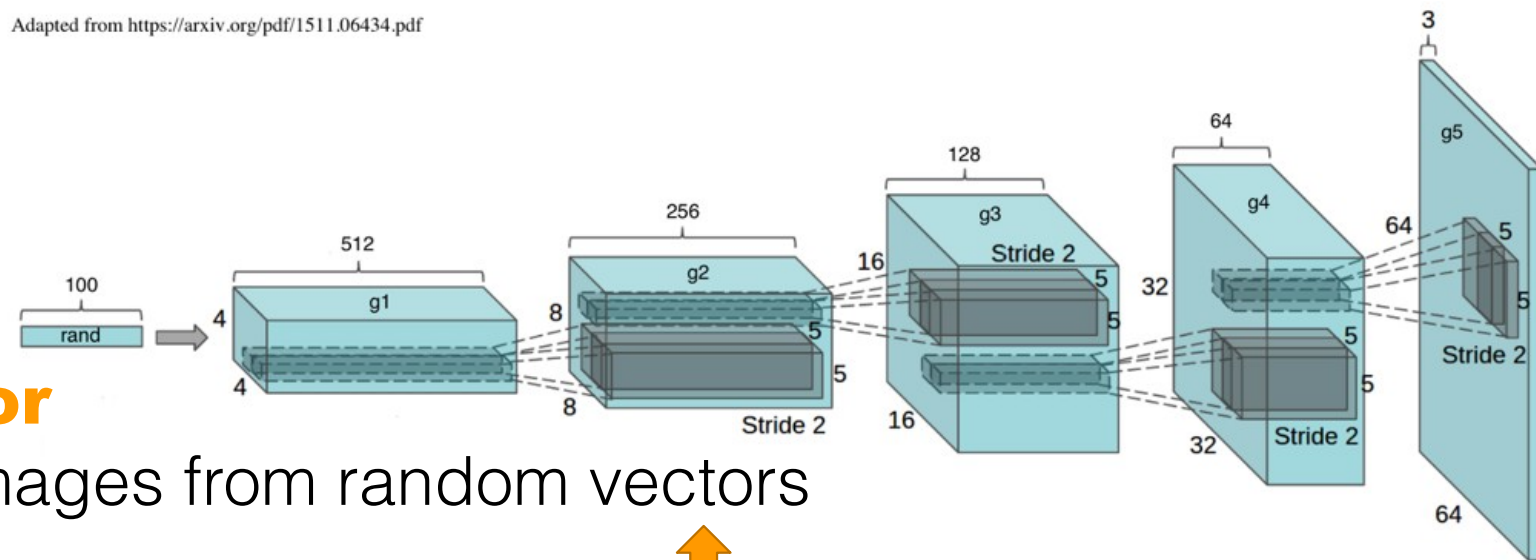


© 2020, Amazon Web Services, Inc. or its affiliates. All rights reserved.

Generative Adversarial Networks

Goodfellow, 2014 <https://arxiv.org/abs/1406.2661>

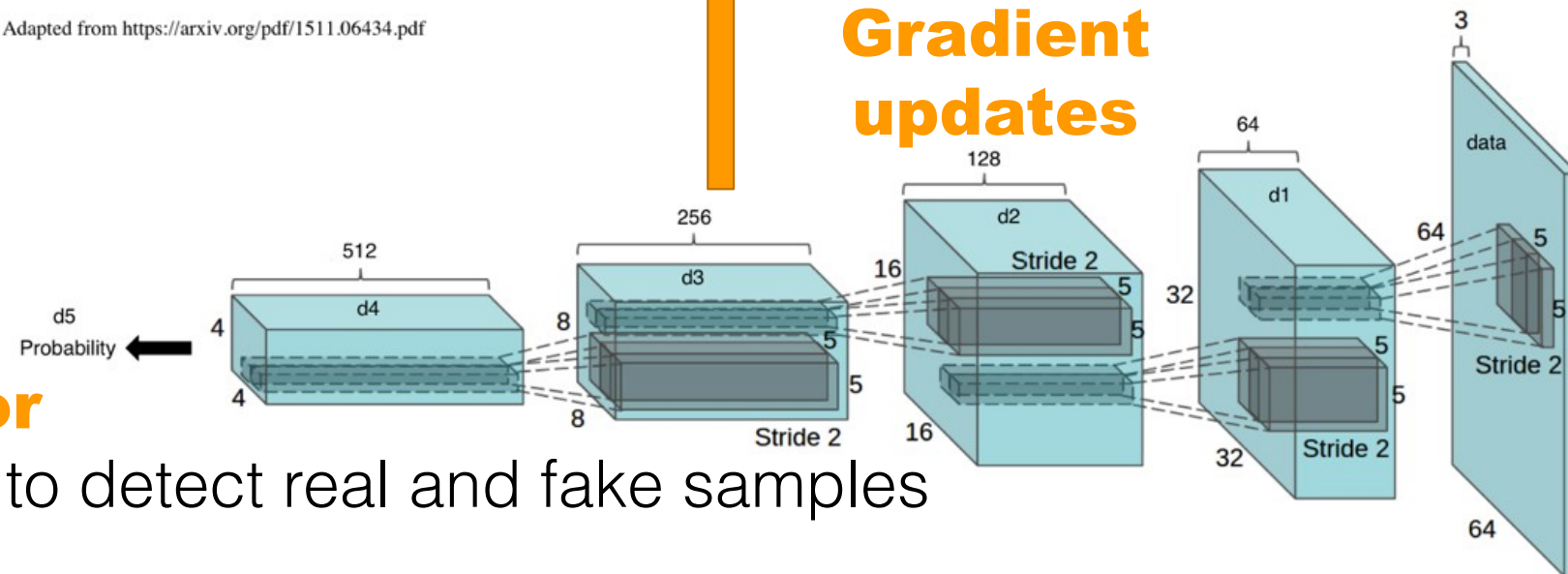
Adapted from <https://arxiv.org/pdf/1511.06434.pdf>



Generator

Building images from random vectors

Adapted from <https://arxiv.org/pdf/1511.06434.pdf>



Detector

Learning to detect real and fake samples

Gradient updates

Fake images

Real images



<https://medium.com/@julsimon/generative-adversarial-networks-on-apache-mxnet-part-1-b6d39e6b5df1>



GAN: Welcome to the (un)real world, Neo

TF



PyTorch



Generating new "celebrity" faces https://github.com/tkarras/progressive_growing_of_gans

April 2018

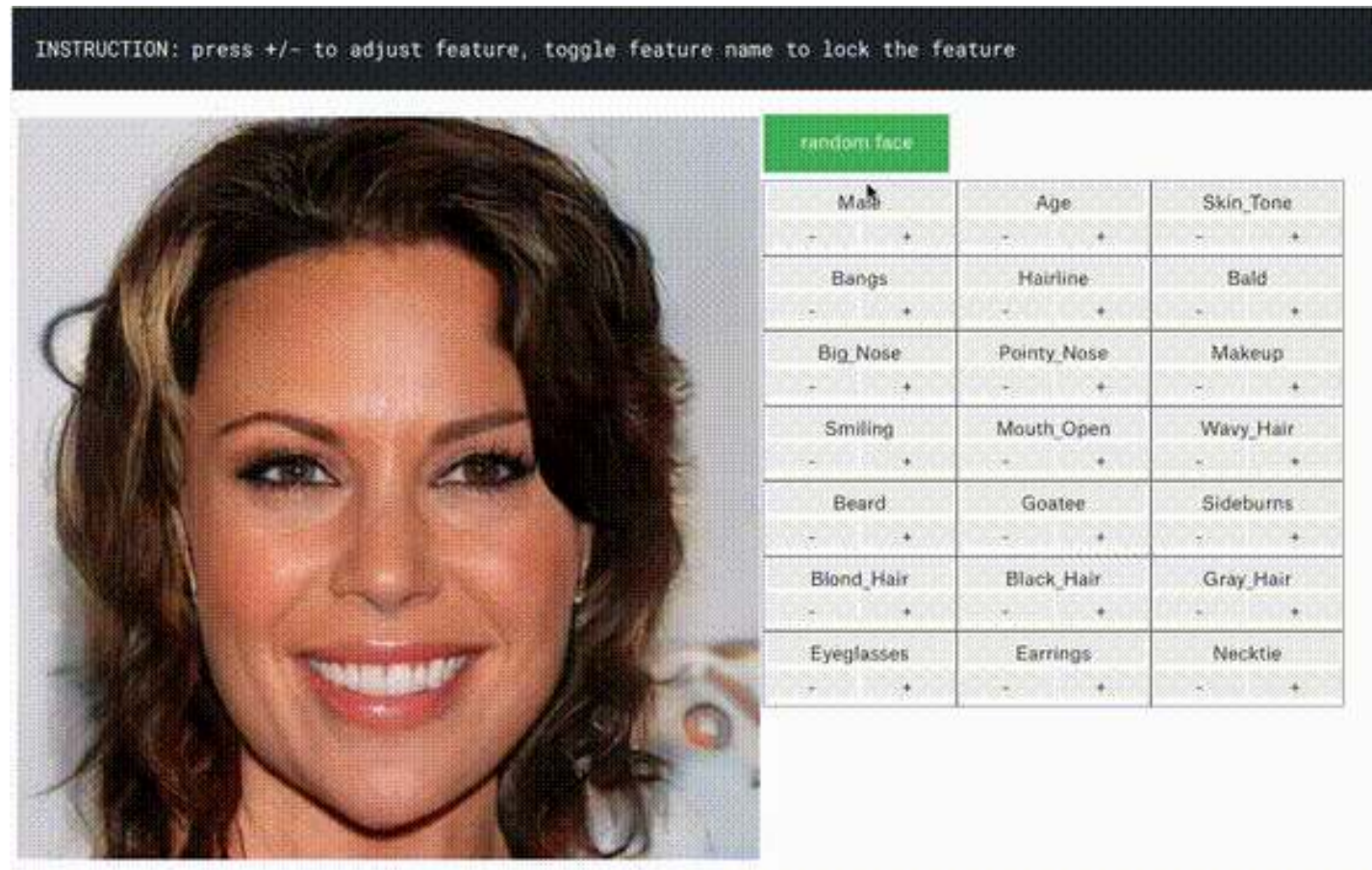
From semantic map to 2048x1024 picture

<https://tcwang0509.github.io/pix2pixHD/>

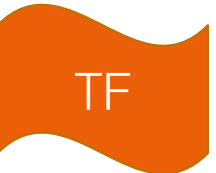
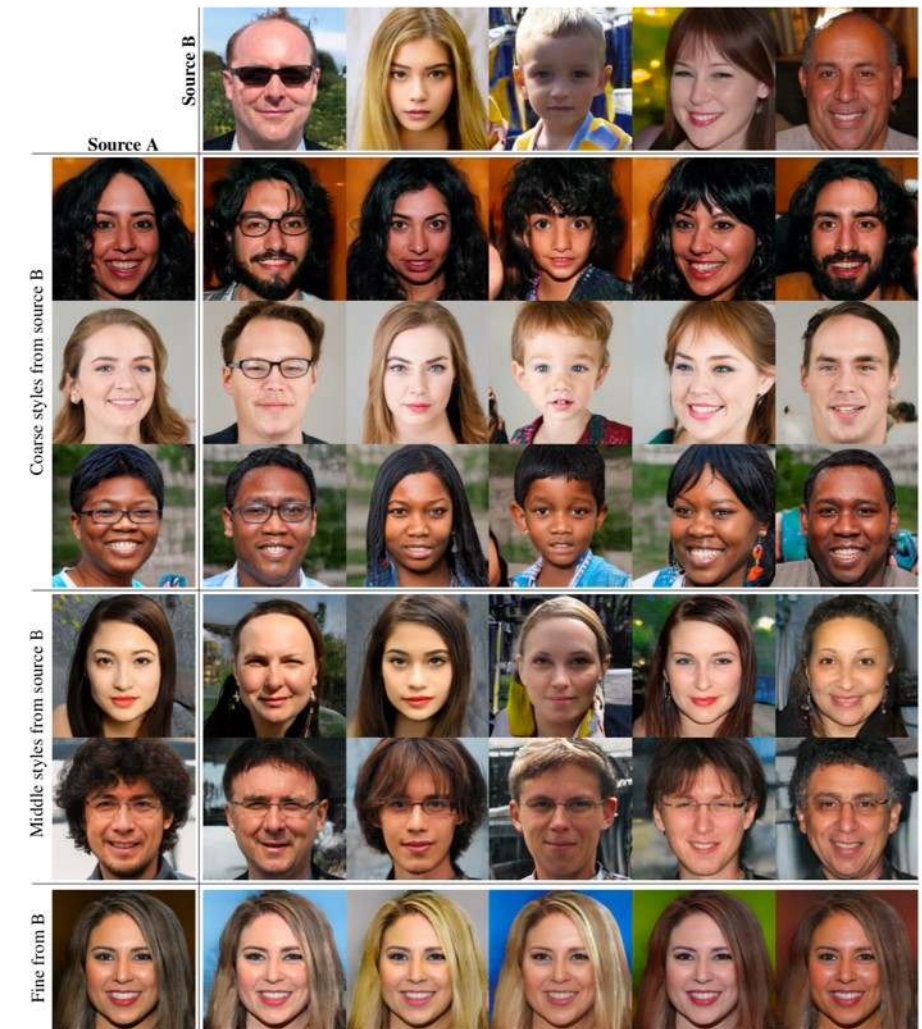
November 2017



More face generation with GANs



Controlled Image Generation with TL-GAN
https://github.com/SummitKwan/transparent_latent_gan
October 2018



Applying the style of a face to another face <https://github.com/NVlabs/stylegan>
<https://www.youtube.com/watch?v=kSLJriaOumA>
March 2019

GAN: Everybody dance now



<https://arxiv.org/abs/1808.07371>
<https://www.youtube.com/watch?v=PCBTZh41Ris>

August 2018

© 2020, Amazon Web Services, Inc. or its affiliates. All rights reserved.



DEV DAY

Getting started



Resources

<https://aws.training/machinelearning>

<https://deeplearning.ai>

<https://fast.ai>

<http://www.deeplearningbook.org/>

<https://d2l.ai/>

<https://gluon.mxnet.io>

<https://keras.io>

<https://gitlab.com/juliensimon/dlnotebooks>



DEV DAY

Thank you!



Appendix – Apache MXNet demos

Demo – Image classification: using a pre-trained model

*** VGG16

```
[(0.46811387, 'n04296562 stage'), (0.24333163, 'n03272010 electric guitar'), (0.045918692, 'n02231487 walking stick, walkingstick, stick insect'), (0.03316205, 'n04286575 spotlight, spot'), (0.021694135, 'n03691459 loudspeaker, speaker, speaker unit, loudspeaker system, speaker system')]
```

*** ResNet-152

```
[(0.8726753, 'n04296562 stage'), (0.046159592, 'n03272010 electric guitar'), (0.041658506, 'n03759954 microphone, mike'), (0.018624334, 'n04286575 spotlight, spot'), (0.0058045341, 'n02676566 acoustic guitar')]
```

*** Inception v3

```
[(0.44991142, 'n04296562 stage'), (0.43065304, 'n03272010 electric guitar'), (0.067580454, 'n04456115 torch'), (0.012423956, 'n02676566 acoustic guitar'), (0.0093934005, 'n03250847 drumstick')]
```



Demo – Image classification: fine-tuning a model

CIFAR-10 data set

60,000 images in 10 classes
32x32 color images

Initial training

Resnet-50 CNN

200 epochs

82.12% validation

Cars vs. horses

88.8% validation accuracy

airplane



automobile



bird



cat



deer



dog



frog



horse



ship



truck



Demo – Image classification: fine-tuning a model

Freezing all layers but the last one

Fine-tuning on « cars vs. horses » for 10 epochs

2 minutes on 1 GPU

98.8% validation accuracy

```
Epoch 10/10
```

```
10000/10000 [=====] - 12s
```

```
loss: 1.6989 - acc: 0.9994 - val_loss: 1.7490 - val_acc: 0.9880
```

```
2000/2000 [=====] - 2s
```

```
[1.7490020694732666, 0.9879999999999999]
```

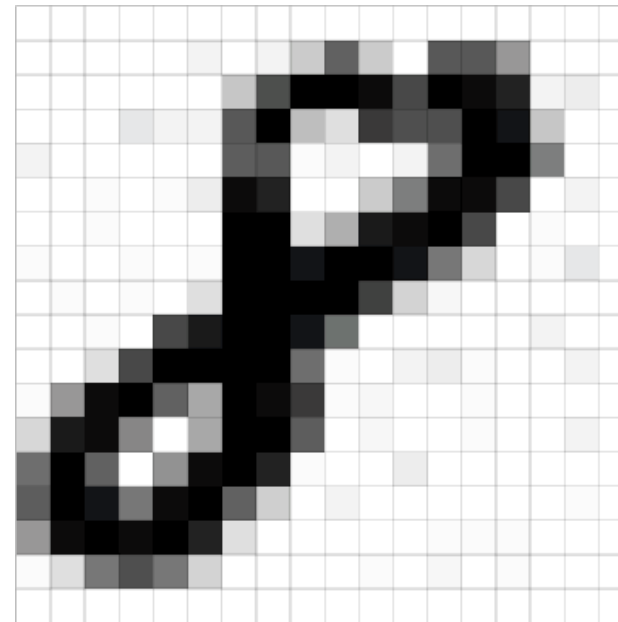


Demo – Image classification: learning from scratch

MNIST data set

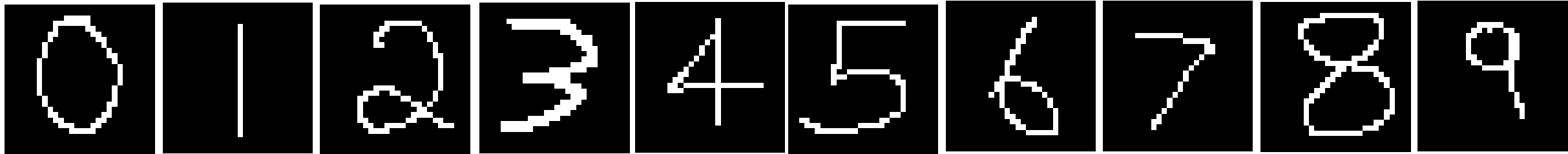
70,000 hand-written digits

28x28 grayscale images

[illegible]

Multi-Layer Perceptron vs. Handmade-Digits-From-Hell™

784/128/64/10, Relu, AdaGrad, 100 epochs → 97.51% validation accuracy

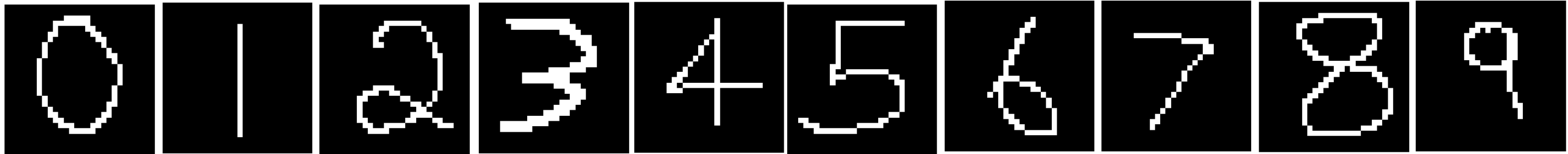


[[0.839	0.034	0.039	0.009	0.	0.008	0.066	0.002	0.	0.004]]
[[0.	0.988	0.001	0.003	0.001	0.001	0.002	0.003	0.001	0.002]]
[[0.006	0.01	0.95	0.029	0.	0.001	0.004	0.	0.	0.]]
[[0.	0.	0.	1.	0.	0.	0.	0.	0.	0.]]
[[0.	0.001	0.005	0.001	0.982	0.001	0.	0.007	0.	0.002]]
[[0.001	0.001	0.	0.078	0.	0.911	0.01	0.	0.	0.]]
[[0.003	0.	0.019	0.	0.005	0.004	0.863	0.	0.105	0.001]]
[[0.001	0.008	0.098	0.033	0.	0.	0.	0.852	0.004	0.004]]
[[0.001	0.	0.006	0.	0.	0.001	0.002	0.	0.991	0.]]
[[0.002	0.158	0.007	0.117	0.082	0.001	0.	0.239	0.17	0.224]]



LeNet vs. Handmade-Digits-From-Hell™

ReLu instead of tanh, 20 epochs, AdaGrad → 99.20% validation accuracy



[[1.	0.	0.	0.	0.	0.	0.	0.	0.	0.]]
[[0.	1.	0.	0.	0.	0.	0.	0.	0.	0.]]
[[0.	0.	1.	0.	0.	0.	0.	0.	0.	0.]]
[[0.	0.	0.	1.	0.	0.	0.	0.	0.	0.]]
[[0.	0.	0.001	0.	0.998	0.	0.	0.001	0.	0.]]
[[0.	0.	0.	0.	0.	1.	0.	0.	0.	0.]]
[[0.	0.	0.	0.	0.	0.	1.	0.	0.	0.]]
[[0.	0.	0.	0.001	0.	0.	0.	0.999	0.	0.]]
[[0.	0.	0.006	0.	0.	0.	0.	0.	0.994	0.]]
[[0.	0.	0.	0.001	0.001	0.	0.	0.001	0.001	0.996]]]