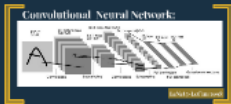
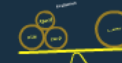
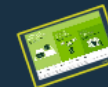
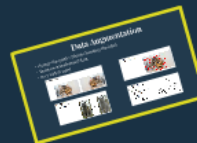


Introduction to Deep Learning (CNN)

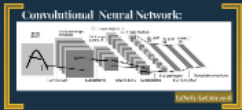
Oliver Schottmüller, Leibniz
UNi Hannover
Visual Neural Network (VNN) University



Deep Learning
What?
Why?
How?

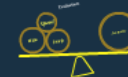
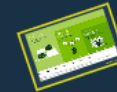


Maher, Mohammad; Ibrahim,
LCCN: April 2008
Mahid Rajeev; Koster; Richard; Chakraborty



Deep Learning

- What?
- Why?
- How?





Introduction to Deep Learning (CNNs)

Mahya Mohammadi Kashani

CSNN, April 2018

Shahid Rajaei Teacher Training University

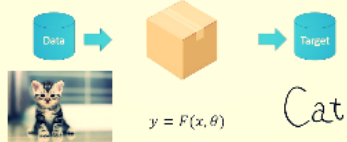
Out Line

- Review of Machine Learning (Classification)
- Why is Deep Learning?
- Applications
- Challenges
- Structures of CNNs
- Learning tricks
- CNN Architectures
- relation between DL and Brain!
- Implementations (TensorFlow , Keras)
- Deep models in AIA Task
- Conclusion
- References

Introduction

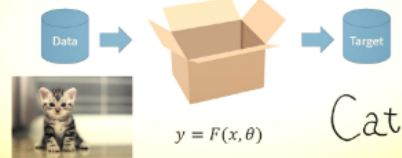
What is Deep Learning?

Deep Learning is a subfield of machine learning.



What is Deep Learning?

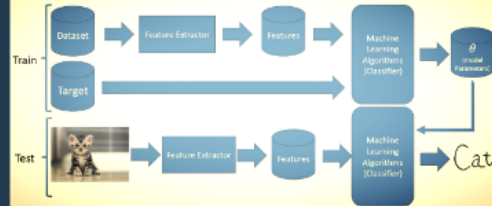
Deep Learning is a subfield of machine learning.



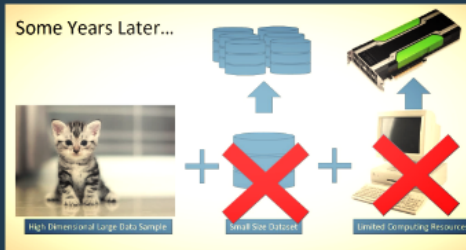
Machine Learning before being Deep...



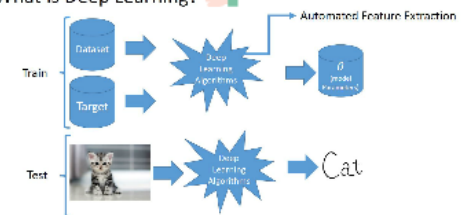
Machine Learning before being Deep...



Some Years Later...



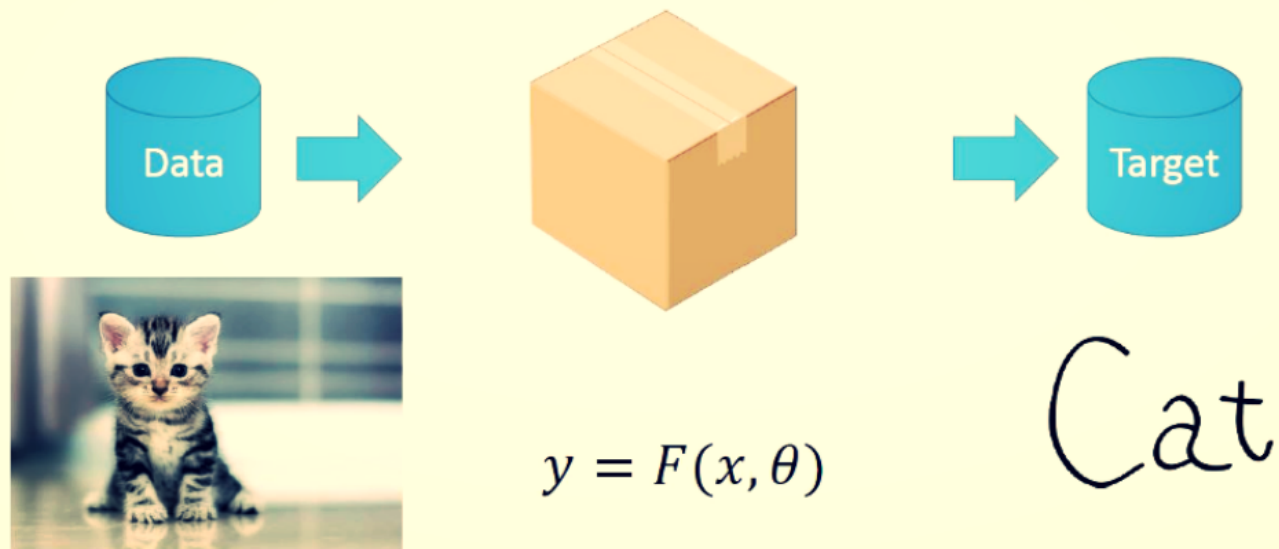
What is Deep Learning?



credit : UTDLSS2017

What is Deep Learning?

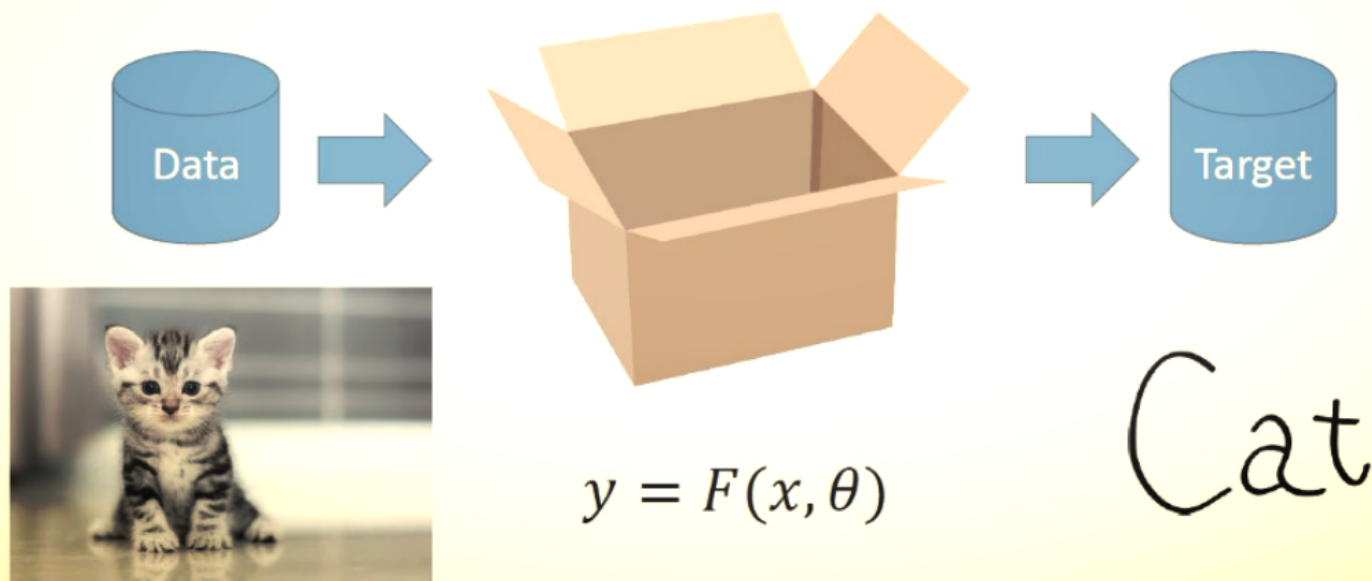
Deep Learning is a subfield of machine learning.



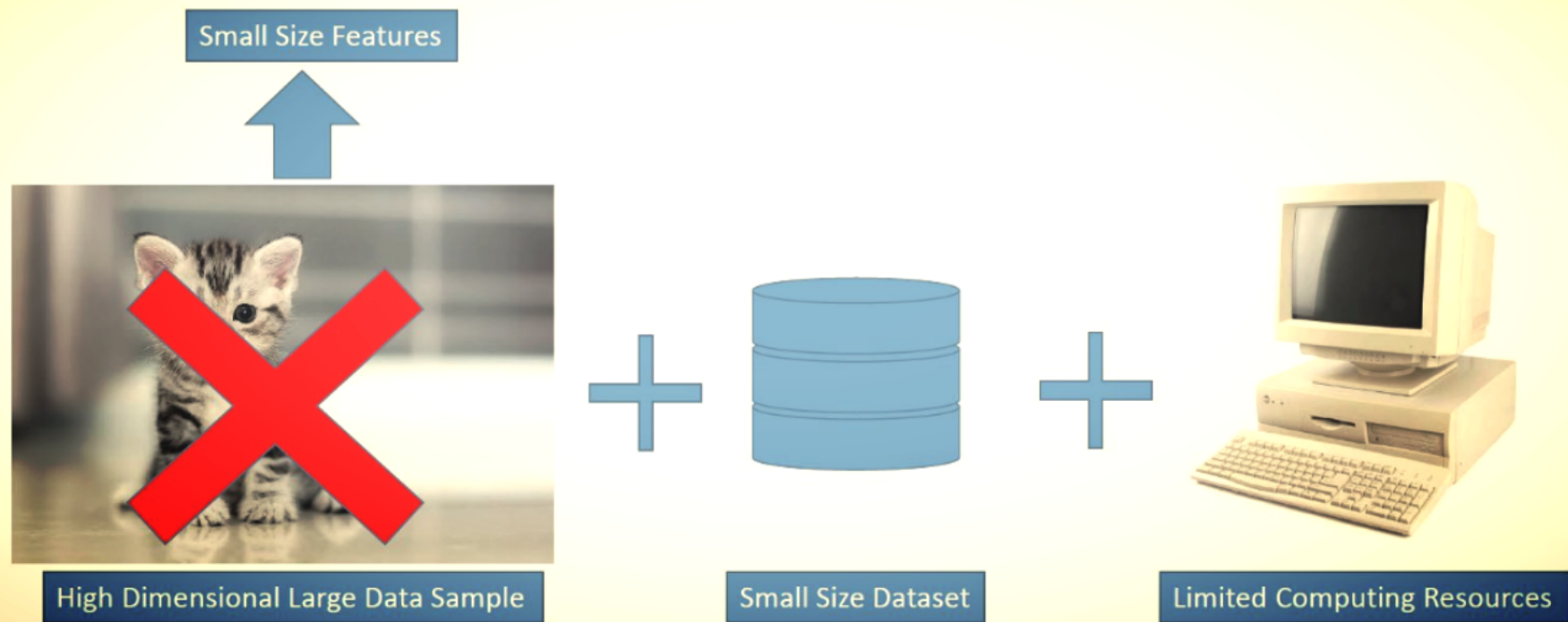
Introduction

What is Deep Learning?

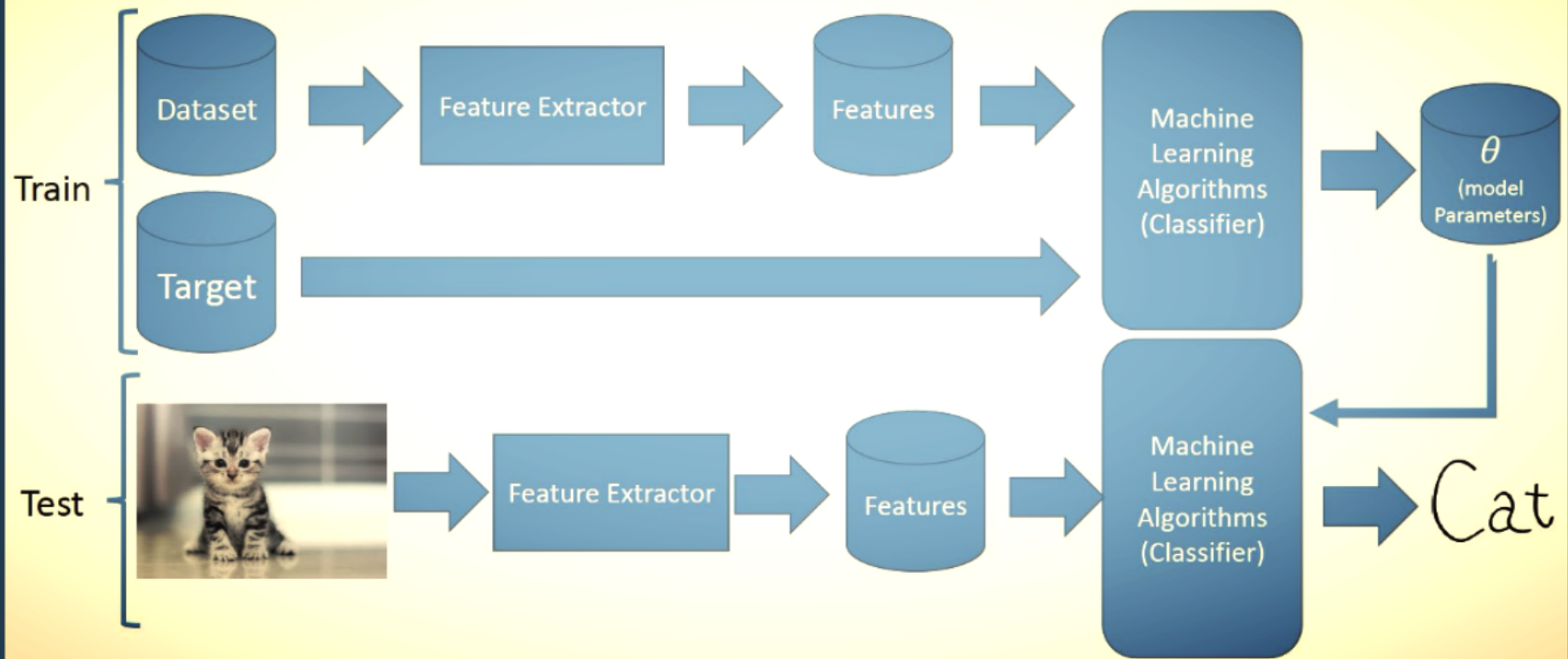
Deep Learning is a subfield of machine learning.



Machine Learning before being Deep...



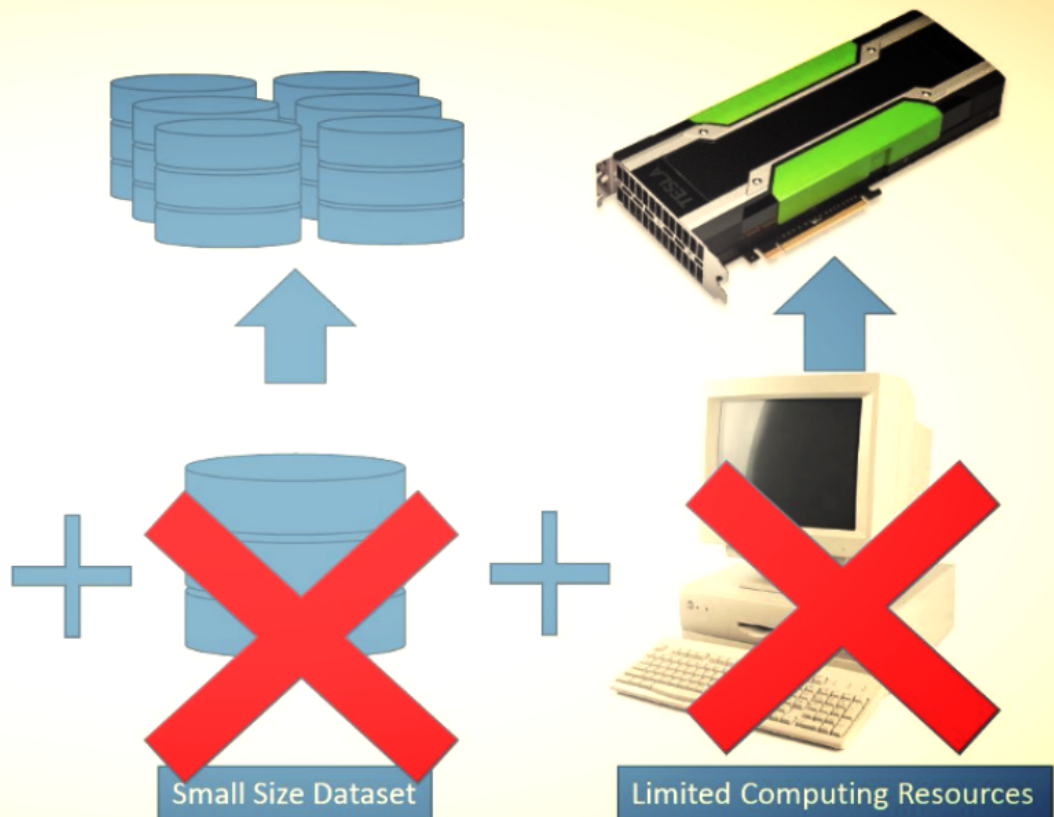
Machine Learning before being Deep...



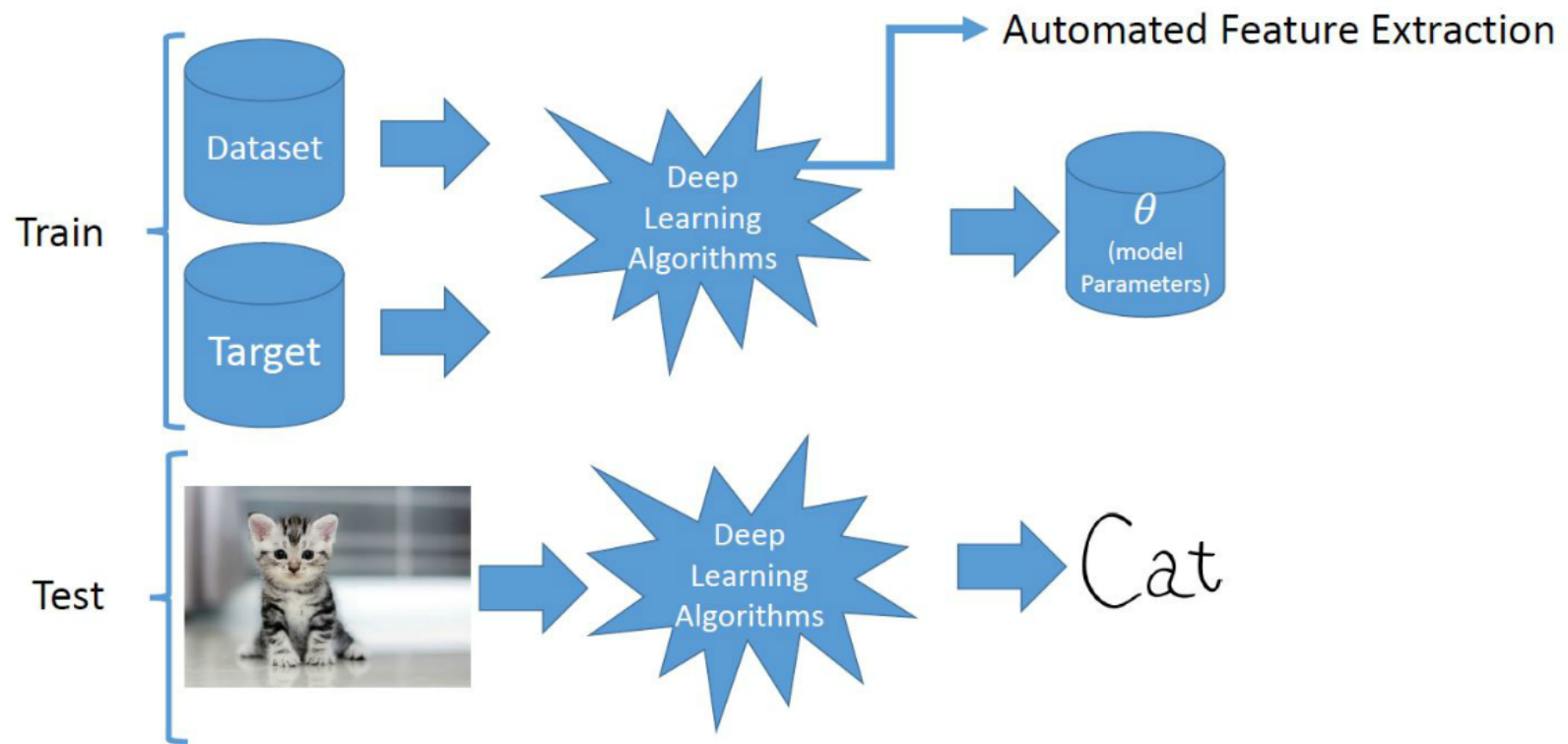
Some Years Later...



High Dimensional Large Data Sample



What is Deep Learning? 🍷



ARTIFICIAL INTELLIGENCE

Early artificial intelligence stirs excitement.



MACHINE LEARNING

Machine learning begins to flourish.



DEEP LEARNING

Deep learning breakthroughs drive AI boom.



1950's 1960's 1970's 1980's 1990's 2000's 2010's

Since an early flush of optimism in the 1950s, smaller subsets of artificial intelligence – first machine learning, then deep learning, a subset of machine learning – have created ever larger disruptions.

credit : Nvidia

Applications

Autonomous navigation of cars and drones



<https://www.cmu.edu/robotics/autonomous/>

Robots use deep learning

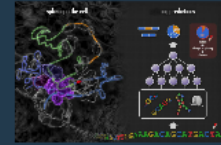


credit : Nvidia

Health Care



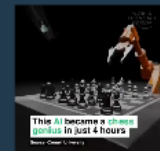
Deep Learning in Genomic



Deep Learning in Conversational Interface



Alpha Zero : Genius of Chess



This AI became a chess grandmaster in just 4 hours

Deep Learning In Speech Separation and Recognition



Autonomous navigation of cars and drones



Autonomous navigation of cars and drones

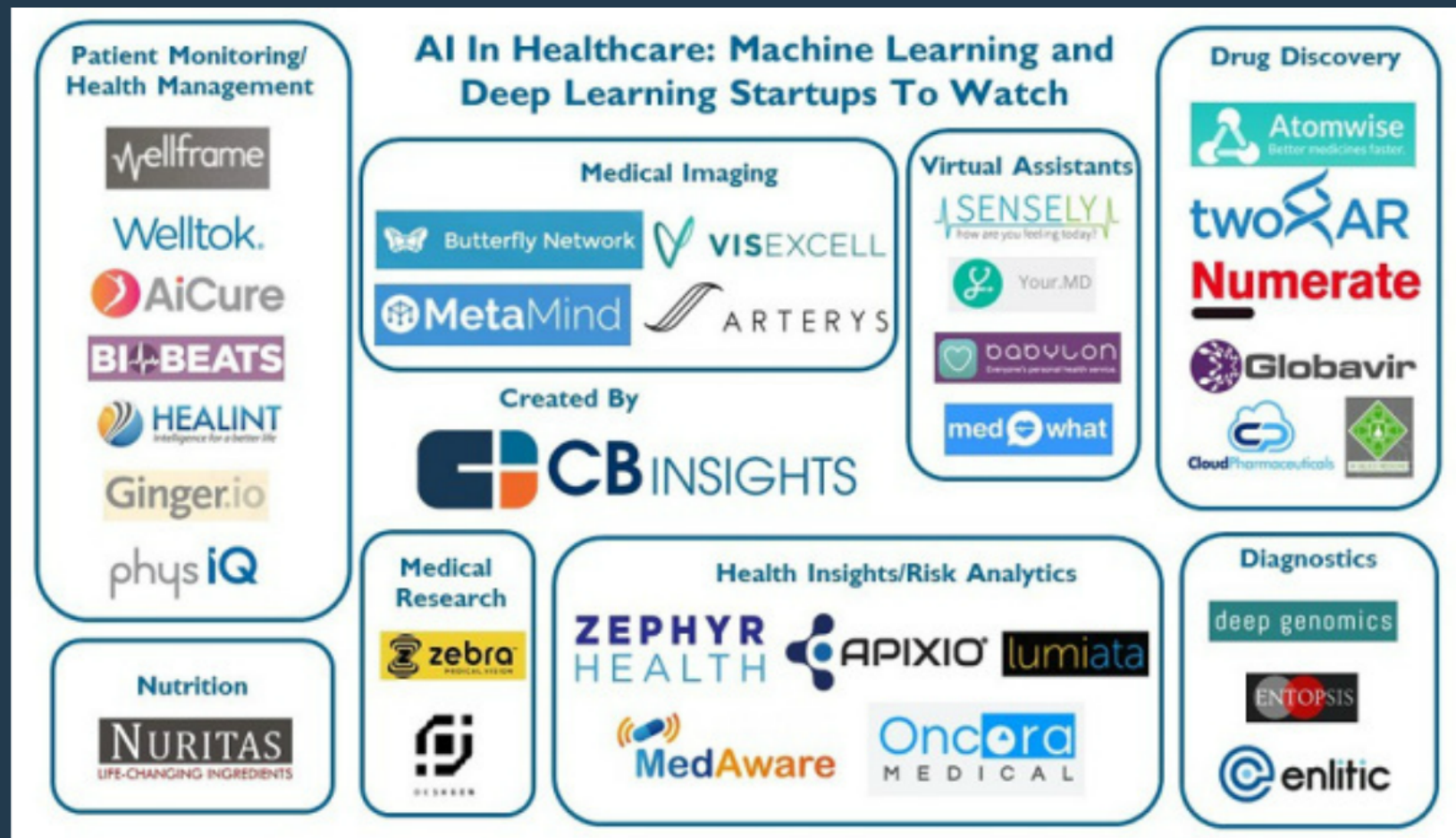


Robots use deep learning



credit : Nvidia

Health Care



Deep Learning in Conversational Interface



 alamy stock photo

JKF3XJ
www.alamy.com

Deep Learning in Speech Separation and Recognition

Input video + auto captions



<https://research.googleblog.com/2018/04/looking-to-listen-audio-visual-speech.html>

Deep Learning in Speech Separation and Recognition

Input video + auto captions



<https://research.googleblog.com/2018/04/looking-to-listen-audio-visual-speech.html>

Alpha Zero : Genius of Chess

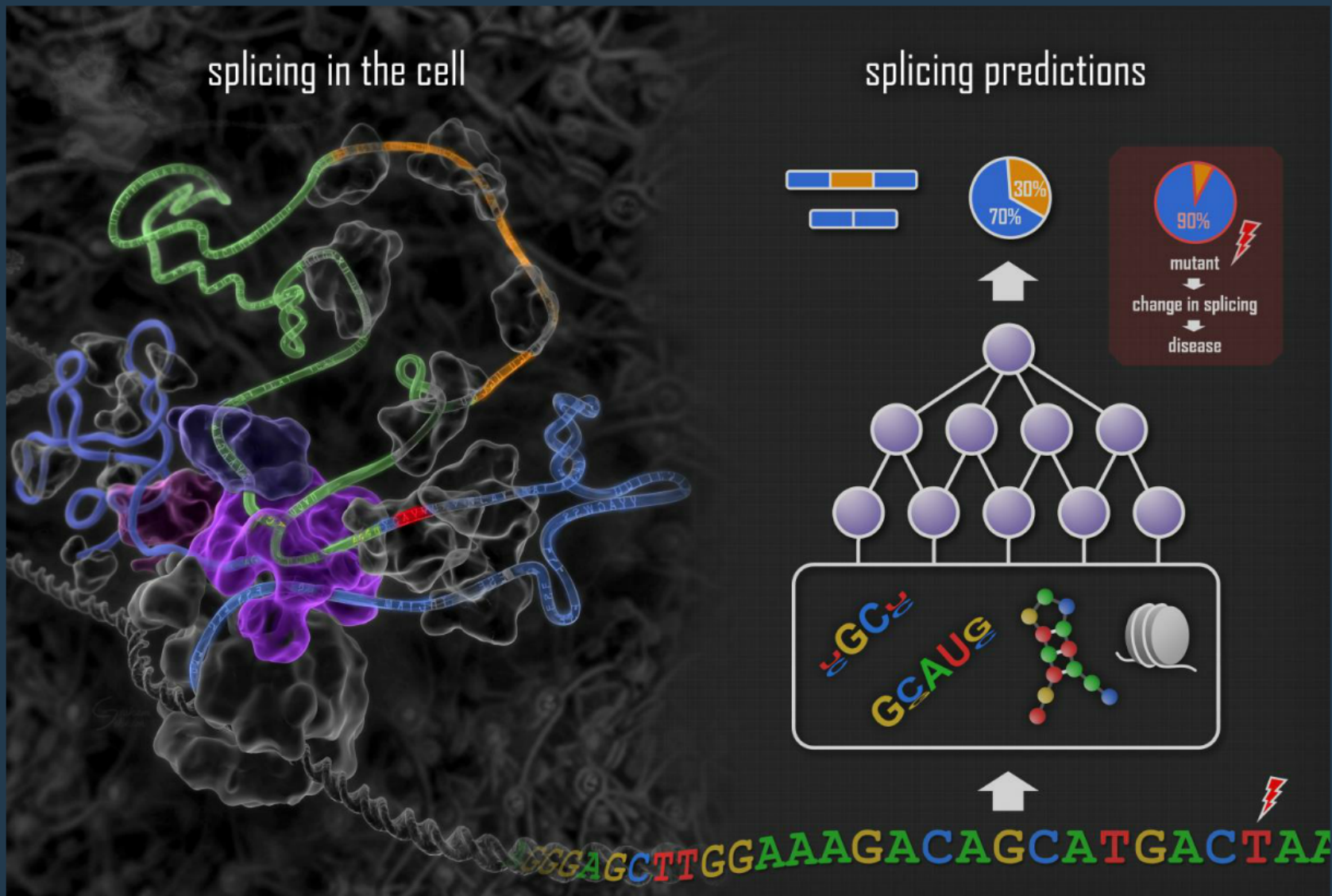




This AI became a chess genius in just 4 hours

Source: Cornell University

Deep Learning in Genomic



'Deep learning' reveals unexpected genetic roots of cancers, autism and other disorders

Image Classification



*Assume: you have given a set discrete labels :
{Cat, Dog, Bus, Tree, ... }*

Cat



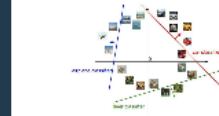
Classification

Machine Learning Problem



- Find model such that the model minimizes the loss on new data.

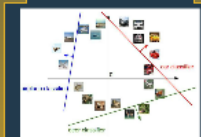
What is Classification?



Linear Classifier



Linear Classifier



Loss function

Example: Suppose 3 inputs, 2 classes

$f(x; W) = Wx$

Model output: $y = Wx$

Ground truth: y_{true}

Loss: $L = \sum_{i=1}^N \max(0, 1 - y_{true,i} \cdot y_i)$

Example: $y_{true} = [1, 0, 0]$, $y = [0.5, 0.5, 0.5]$

Loss: $L = 1.5$

Example: Suppose 3 inputs, 2 classes

$f(x; W) = Wx$

Model output: $y = Wx$

Ground truth: y_{true}

Loss: $L = \sum_{i=1}^N \max(0, 1 - y_{true,i} \cdot y_i)$

Example: $y_{true} = [1, 0, 0]$, $y = [0.5, 0.5, 0.5]$

Loss: $L = 1.5$

Example: Suppose 3 inputs, 2 classes

$f(x; W) = Wx$

Model output: $y = Wx$

Ground truth: y_{true}

Loss: $L = \sum_{i=1}^N \max(0, 1 - y_{true,i} \cdot y_i)$

Example: $y_{true} = [1, 0, 0]$, $y = [0.5, 0.5, 0.5]$

Loss: $L = 1.5$

Example: Suppose 3 inputs, 2 classes

$f(x; W) = Wx$

Model output: $y = Wx$

Ground truth: y_{true}

Loss: $L = \sum_{i=1}^N \max(0, 1 - y_{true,i} \cdot y_i)$

Example: $y_{true} = [1, 0, 0]$, $y = [0.5, 0.5, 0.5]$

Loss: $L = 1.5$

Loss Formula for all samples:

$$L(W) = \sum_{i=1}^N \max(0, 1 - y_{true,i} \cdot y_i)$$

There is a bug here

Weight Regularization

Example: Suppose 3 inputs, 2 classes

$f(x; W) = Wx$

Model output: $y = Wx$

Ground truth: y_{true}

Loss: $L = \sum_{i=1}^N \max(0, 1 - y_{true,i} \cdot y_i)$

Example: $y_{true} = [1, 0, 0]$, $y = [0.5, 0.5, 0.5]$

Loss: $L = 1.5$

Weight Regularization:

$$L(W) = \sum_{i=1}^N \max(0, 1 - y_{true,i} \cdot y_i) + \lambda \|W\|_2^2$$

Example: $y_{true} = [1, 0, 0]$, $y = [0.5, 0.5, 0.5]$

Loss: $L = 1.5 + \lambda \|W\|_2^2$

Weight Regularization Example:



Softmax Classifier:



Example: Suppose 3 inputs, 2 classes

$f(x; W) = Wx$

Model output: $y = Wx$

Ground truth: y_{true}

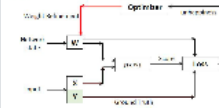
Loss: $L = \sum_{i=1}^N \max(0, 1 - y_{true,i} \cdot y_i)$

Example: $y_{true} = [1, 0, 0]$, $y = [0.5, 0.5, 0.5]$

Loss: $L = 1.5$

Optimization

Example: Suppose 3 inputs, 2 classes



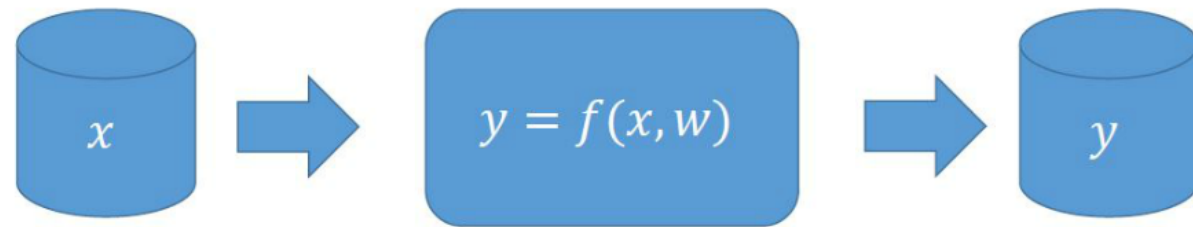
The loss is just a function of W

So just calculate $\frac{\partial L}{\partial W}$ using calculus! ☺



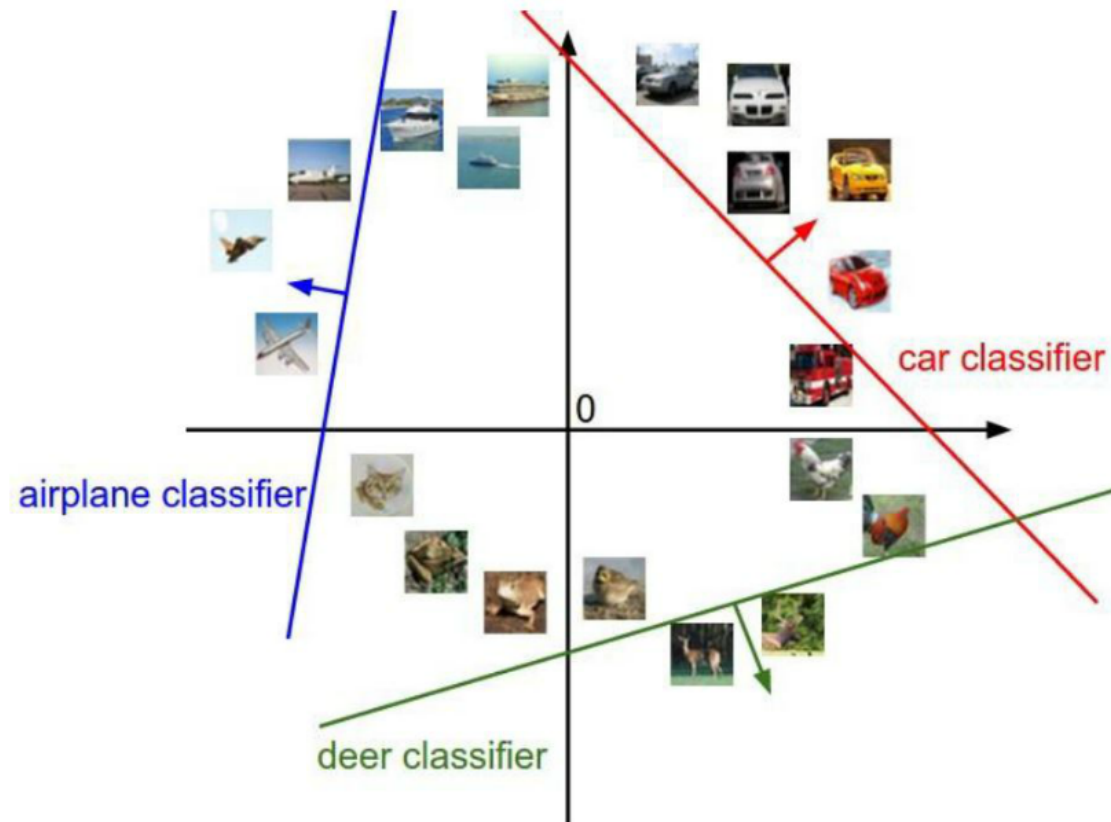
Credit: CS231N Stanford

Machine Learning Problem



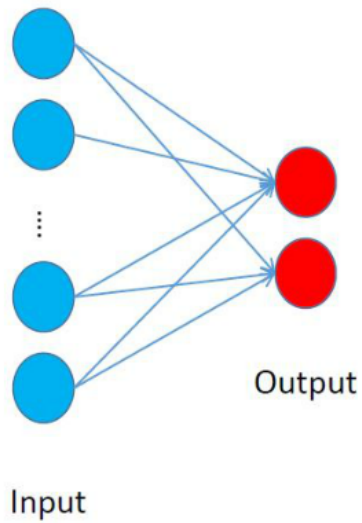
- $f()$ is pre-determined.
- w is the model parameters which need to be learned from data.

What is Classification:

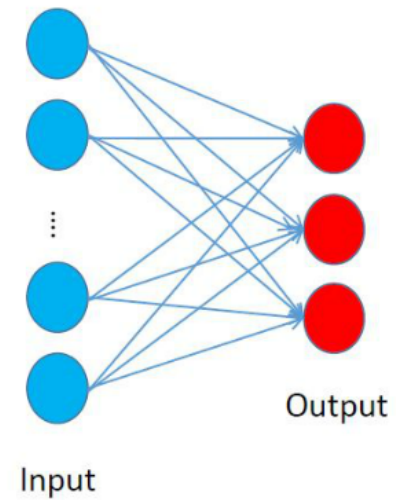


Linear Classifiers:

Two class Classification



Multiclass Classification: More than 2 outputs



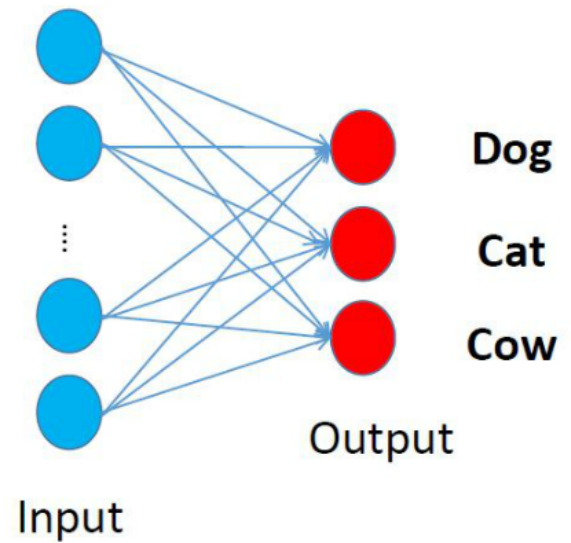
Linear Classifiers:



Input



Vectorization



Weight

Loss function

Example: Suppose 3 images, 3 classes

$$f(x; W) = Wx$$



Dog

3.2

Cat

5.1

Cow

-1.7

Multiclass SVM Loss:

Given an example (x_i, y_i) where x_i is the input sample and y_i is the (integer) label.

Also We define $s = f(x_i; W)$

The SVM formula is:

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

$$L_i = \max(0, 5.1 - 3.2 + 1) \\ + \max(0, -1.7 - 3.2 + 1)$$

$$L_i = 2.9 + 0 = 2.9$$

Example: Suppose 3 images, 3 classes

$$f(x; W) = Wx$$



Dog

1.3

Cat

4.9

Cow

2.0

Multiclass SVM Loss:

Given an example (x_i, y_i) where x_i is the input sample and y_i is the (integer) label.

Also We define $s = f(x_i; W)$

The SVM formula is:

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

$$L_i = \max(0, 1.3 - 4.9 + 1) \\ + \max(0, 2.0 - 4.9 + 1)$$

$$L_i = 0 + 0 = 0$$



Dog	1.3
Cat	4.9
Cow	2.0

Multiclass SVM Loss:

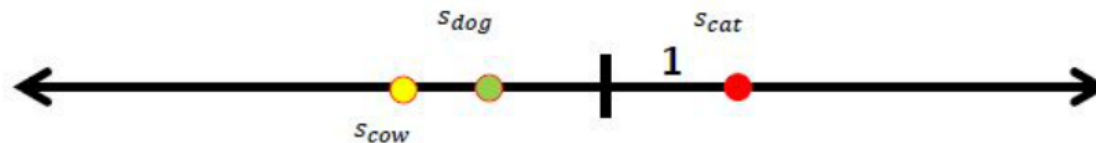
Given an example (x_i, y_i) where x_i is the input sample and y_i is the (integer) label.

Also We define $s = f(x_i; W)$

The SVM formula is:

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

Interpretation: when loss is zero



Credit: CS231N Stanford

Example: Suppose 3 images, 3 classes

$$f(x; W) = Wx$$



Dog

2.2

Cat

2.5

Cow

-3.1

Multiclass SVM Loss:

Given an example (x_i, y_i) where x_i is the input sample and y_i is the (integer) label.

Also We define $s = f(x_i; W)$

The SVM formula is:

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

$$L_i = \max(0, 2.2 + 3.1 + 1) \\ + \max(0, 2.5 + 3.1 + 1)$$

$$L_i = 6.3 + 6.6 = 12.9$$

Credit: CS231N Stanford

Loss Formula for all samples:

$$f(x; W) = Wx$$

$$L = \frac{1}{N} \sum_{i=1}^N \sum_{j \neq y_i} \max(0, f(x_i; W)_j - f(x_i; W)_{y_i} + 1)$$

There is a bug here

Suppose that we found a W such that $L=0$, is this W unique?

Weight Regularization

$$f(x; W) = Wx$$



Dog	1.3
Cat	4.9
Cow	2.0

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

Before:

$$L_i = \max(0, 1.3 - 4.9 + 1) \\ + \max(0, 2.0 - 4.9 + 1)$$

$$L_i = 0 + 0 = 0$$

With W twice as large:

$$L_i = \max(0, 2.6 - 9.8 + 1) \\ + \max(0, 4.0 - 9.8 + 1)$$

$$L_i = 0 + 0 = 0$$

Weight Regularization:

$$f(x; W) = Wx$$

$$L = \frac{1}{N} \sum_{i=1}^N \sum_{j \neq y_i} \max(0, f(x_i; W)_j - f(x_i; W)_{y_i} + 1) + \lambda R(W)$$

- **L2 norm:** $R(W) = \sum_i \sum_j W_{i,j}^2$
- **L1 norm:** $R(W) = \sum_i \sum_j |W_{i,j}|$

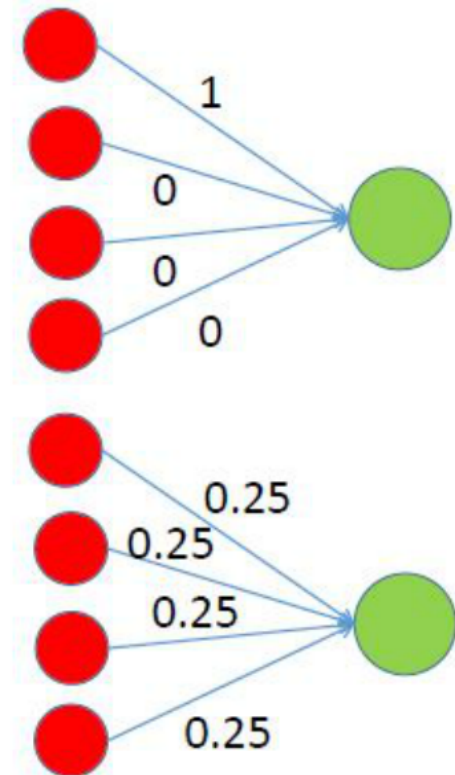
Weight Regularization Example:

$$x = [1, 1, 1, 1]$$

$$W_1 = [1, 0, 0, 0]$$

$$W_2 = [0.25, 0.25, 0.25, 0.25]$$

$$W_1^T \cdot x = W_2^T \cdot x = 1$$



Softmax Classifier:

$$f(x; W) = Wx$$

scores: un-normalized log probabilities of classes



Dog	3.2
Cat	5.1
Cow	-1.7

$$P(Y = k|X = x_i) = \frac{e^{s_i}}{\sum_j e^{s_j}} \quad \text{where} \quad s = f(x_i; W)$$

Want to maximize the log likelihood or (for loss function) minimize the negative log-likelihood of the correct class:

$$L_i = -\log P(Y = k|X = x_i)$$

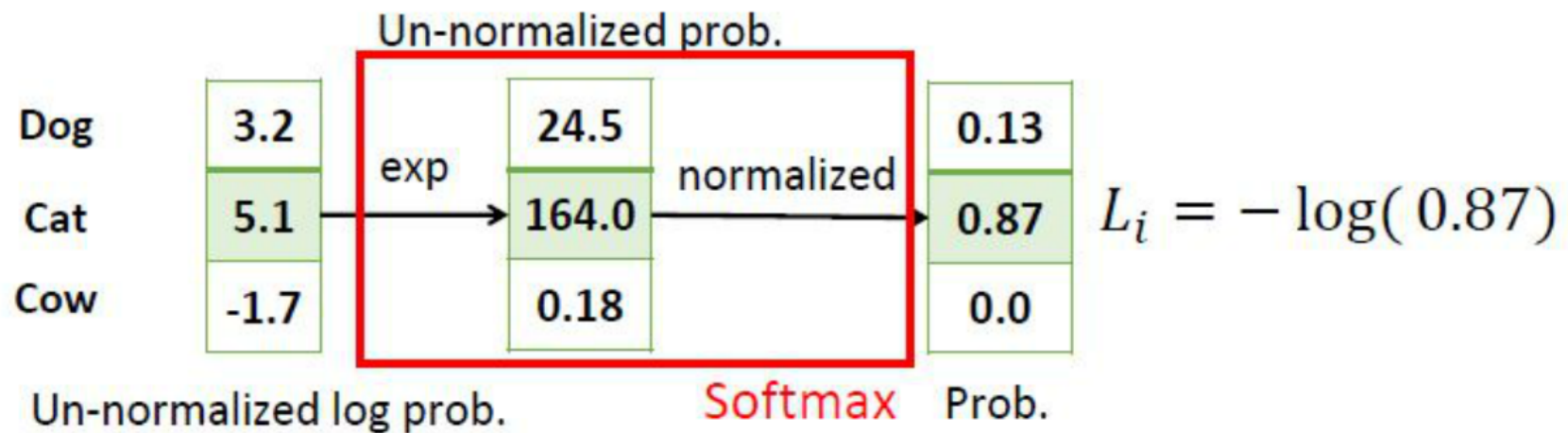
$$L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right)$$

Credit: CS231N Stanford

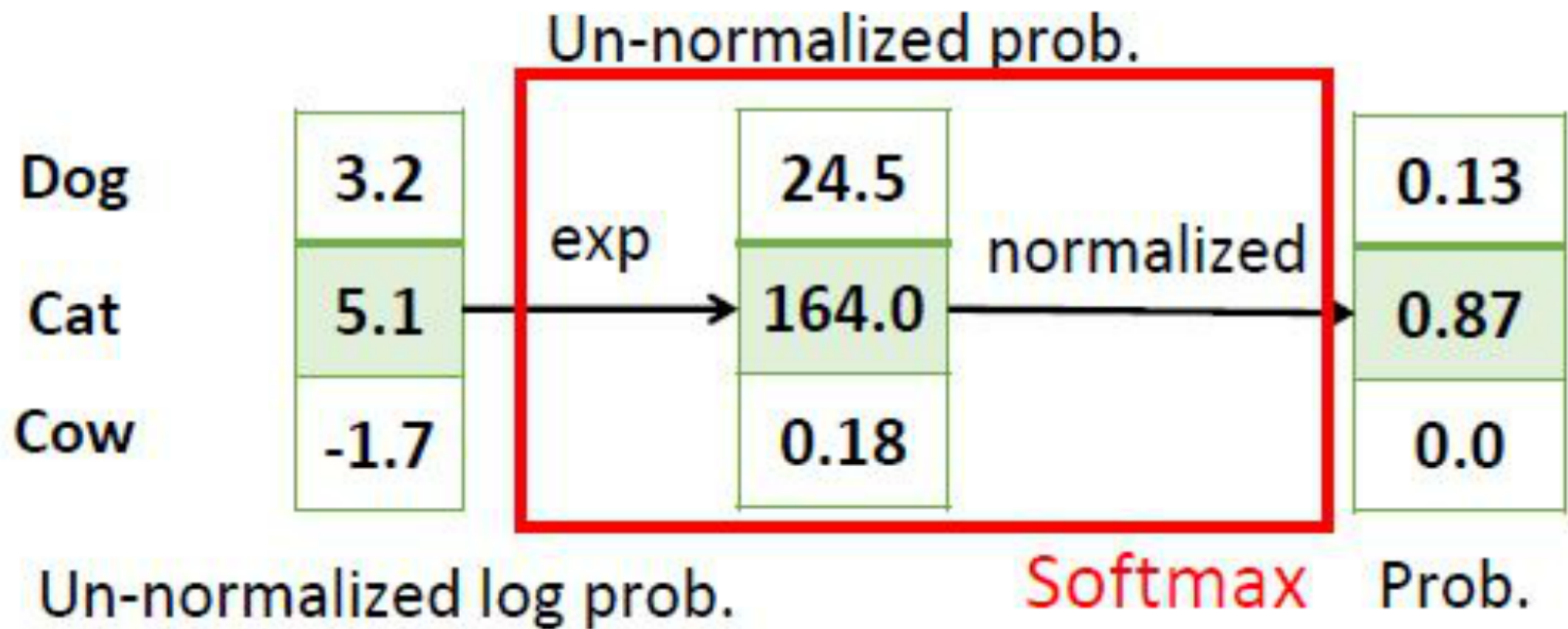


$$L_i = -\log\left(\frac{e^{s_{yi}}}{\sum_j e^{s_{yj}}}\right)$$

$$f(x; W) = Wx$$

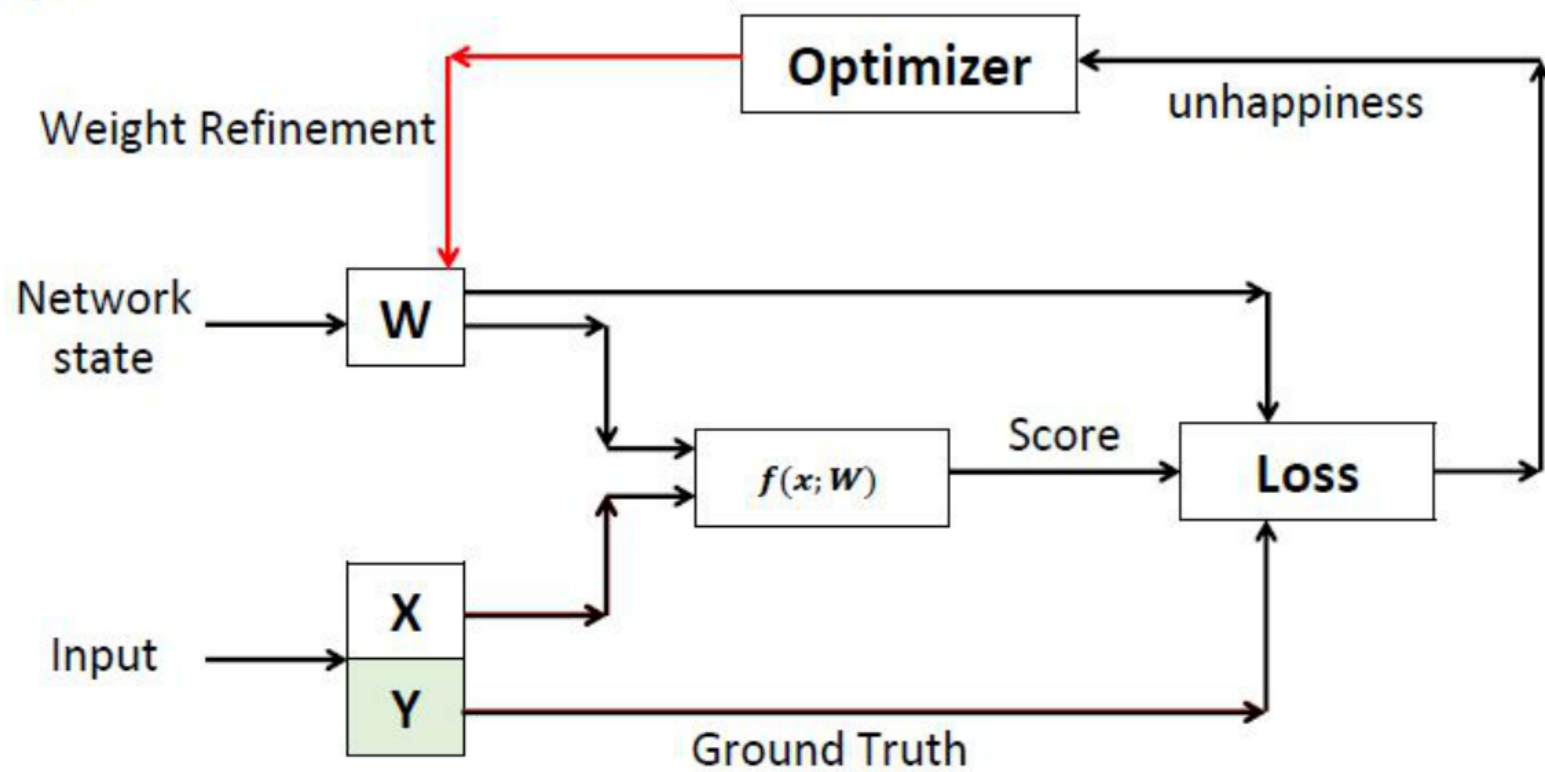


$$L_i = -\log\left(\frac{e^{s y_i}}{\sum_j e^{s_j y_j}}\right)$$



Optimization

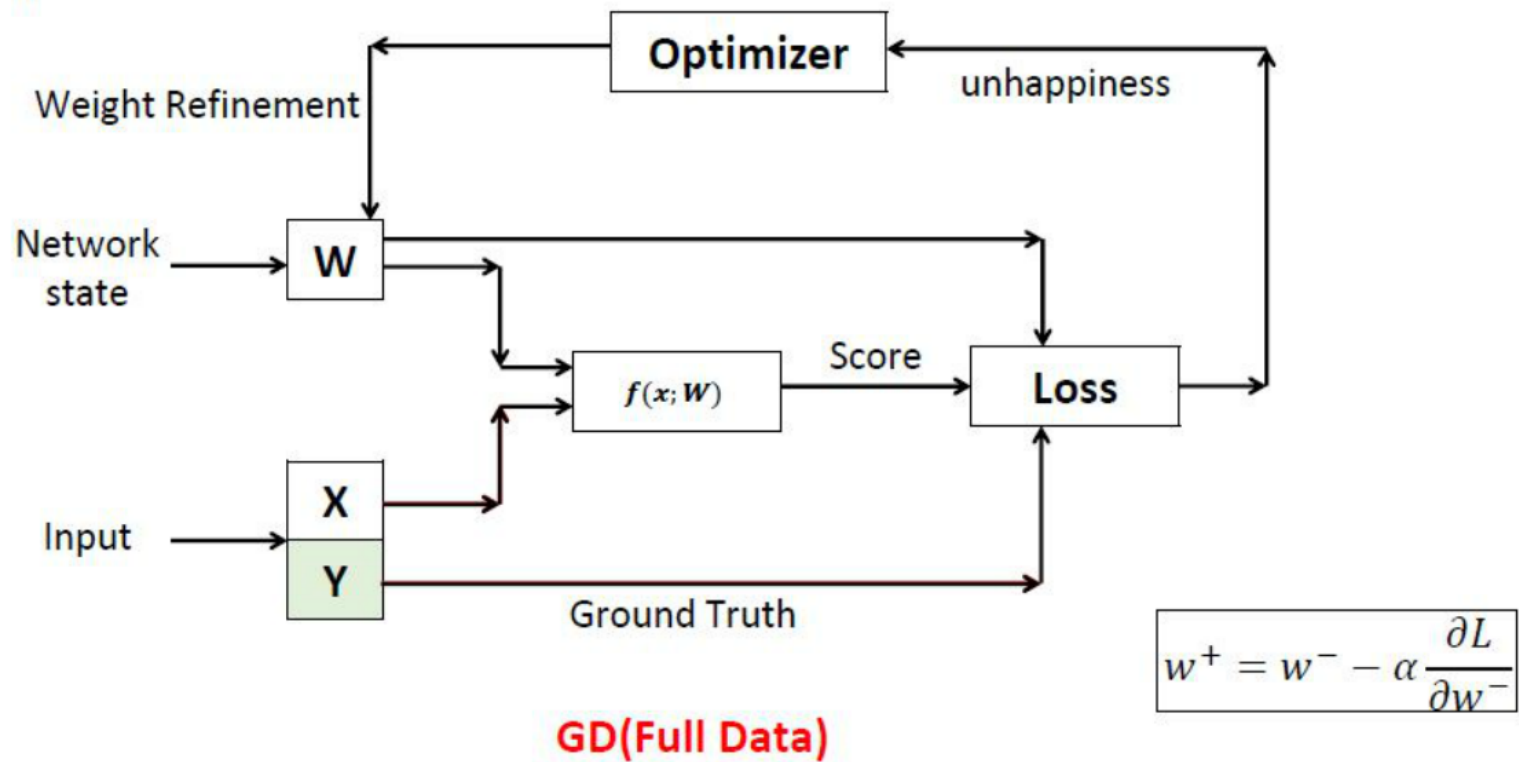
Recap:



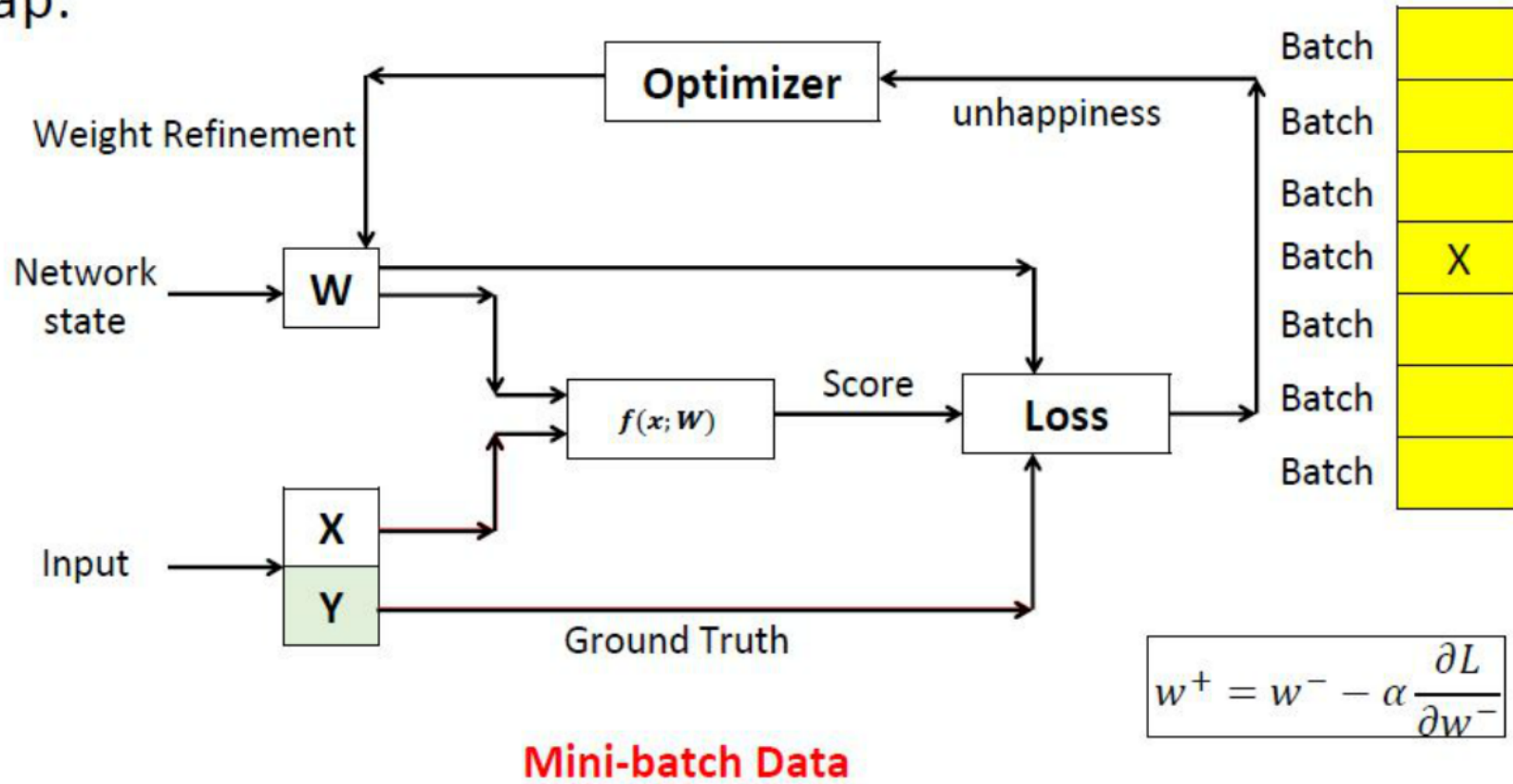
The loss is just a function of W

So just calculate $\frac{\partial L}{\partial w}$
using **calculus!** 😊

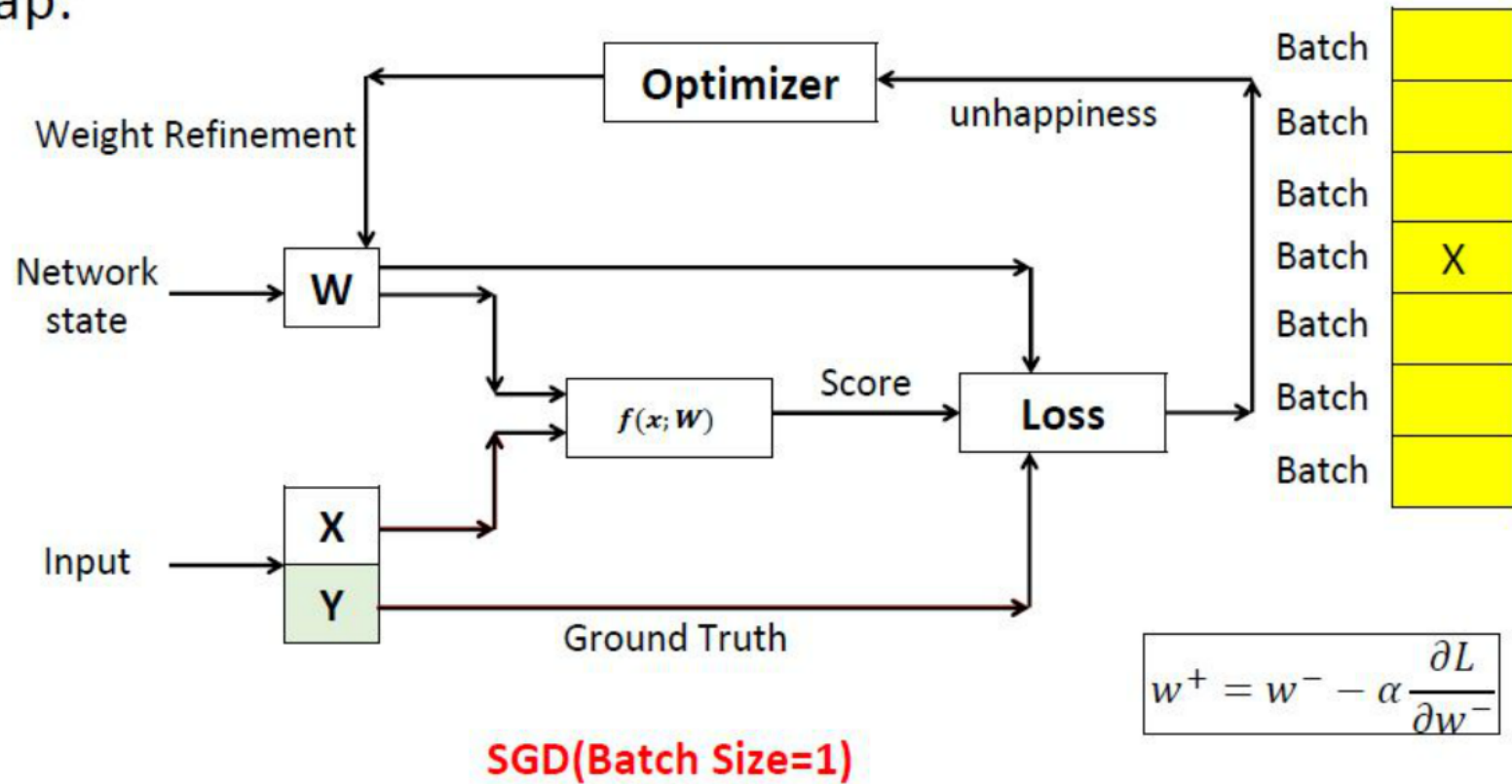
Recap:



Recap:



Recap:



Why Deep Learning?

Important Property of Neural Networks

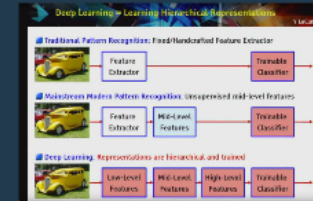
Results get better with

more data +
bigger models +
more computation

(Better algorithms, new insights and improved techniques always help, too!)

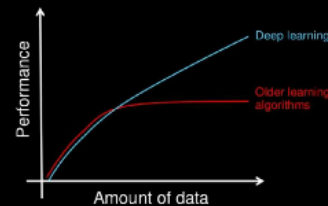


Results Get Better With More Data, Larger Model, More Computation, Slide by Jeff Dean



Deep Learning = Learning Hierarchical Representations Slide by Yann LeCun

Why deep learning



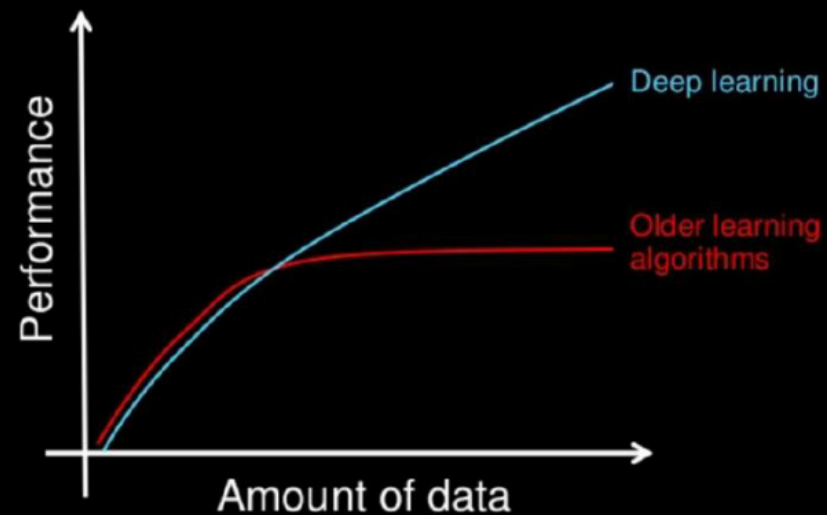
How do data science techniques scale with amount of data?

Why Deep learning, Slide by Andrew NG



Prezi

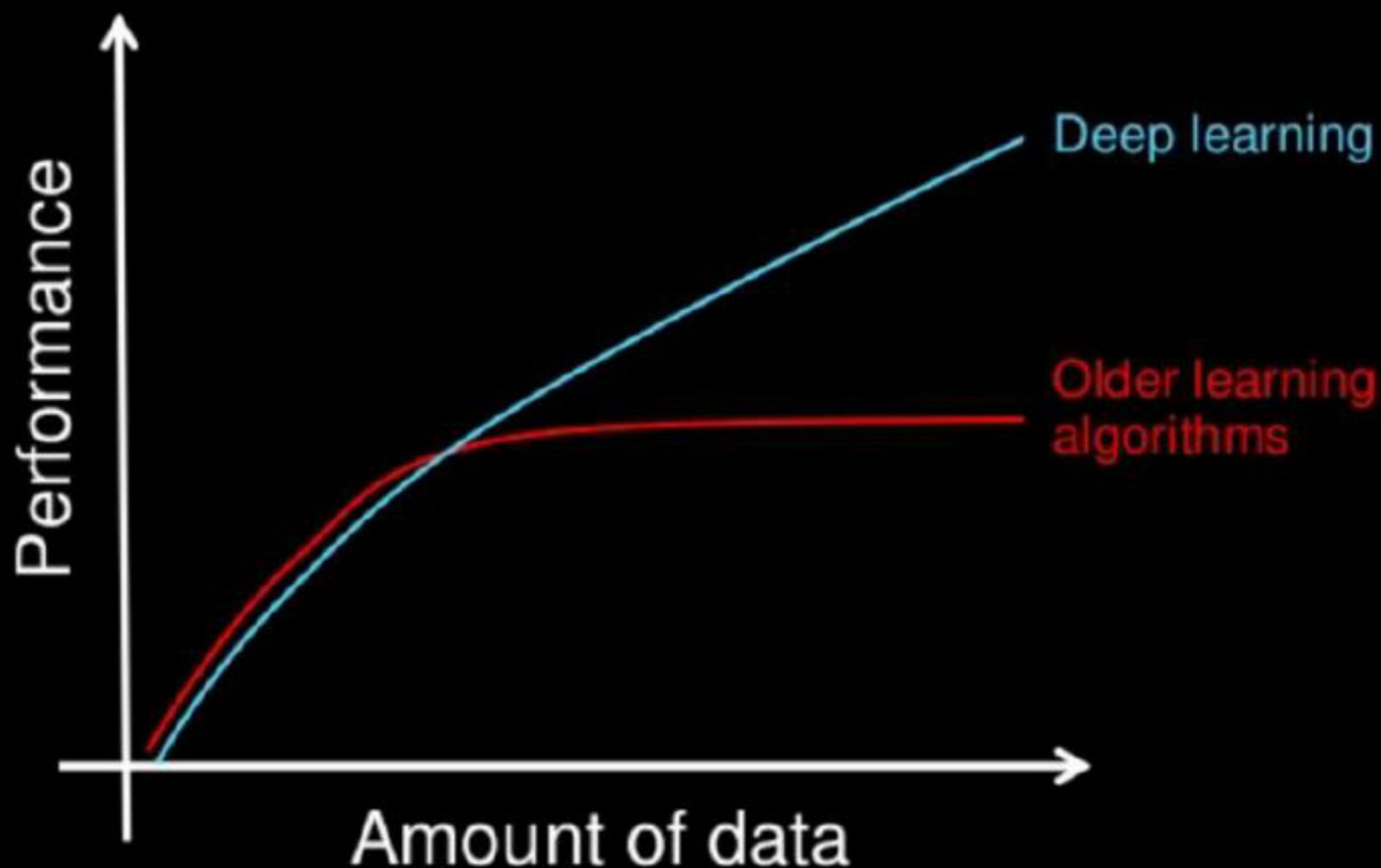
Why deep learning



How do data science techniques scale with amount of data?

Why Deep learning, Slide by Andrew NG

Why deep learning



How do data science techniques scale with amount of data?

Deep Learning = Learning Hierarchical Representations

Y LeCun

Traditional Pattern Recognition: Fixed/Handcrafted Feature Extractor



Mainstream Modern Pattern Recognition: Unsupervised mid-level features



Deep Learning: Representations are hierarchical and trained



Deep Learning = Learning Hierarchical Representations Slide by Yann LeCun

Important Property of Neural Networks

Results get better with

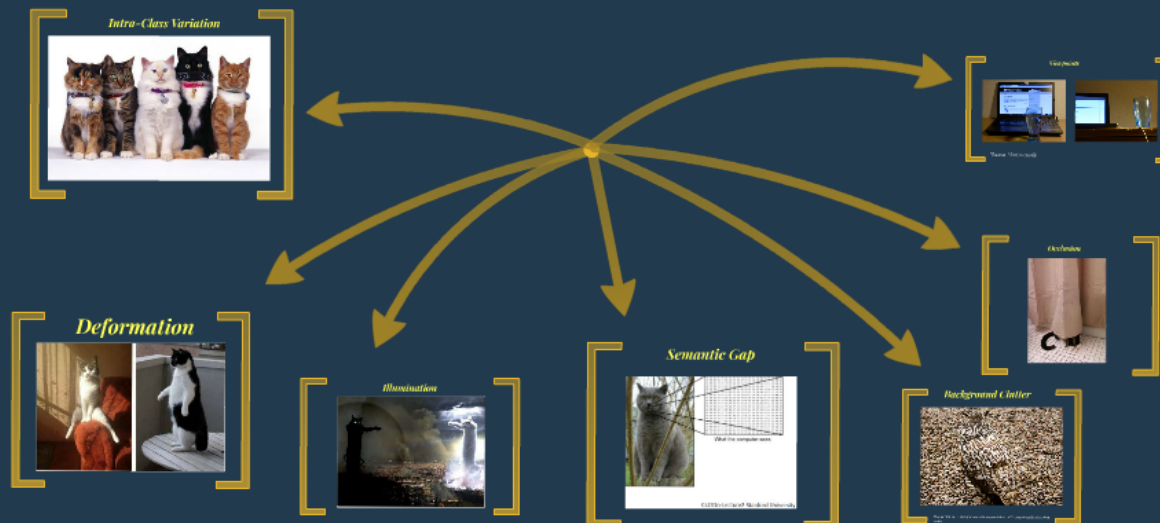
**more data +
bigger models +
more computation**

(Better algorithms, new insights and improved techniques always help, too!)



Result Get Better With More Data, Larger Model, More Computation, Slide by Jeff Dean

challenges



Semantic Gap

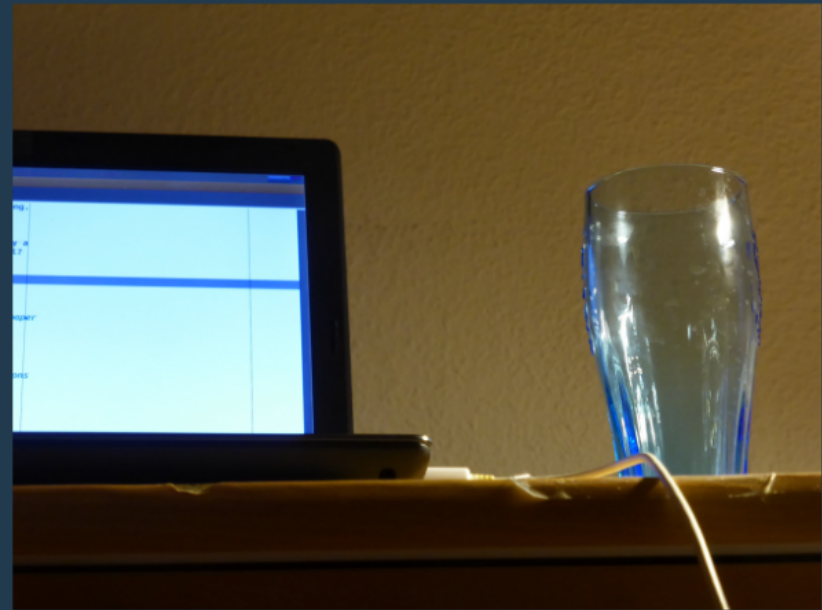


05	02	22	97	38	15	00	48	00	75	24	05	97	72	52	12	50	77	97	00
49	49	99	40	17	81	18	57	60	87	17	40	98	43	65	55	44	58	42	00
81	49	31	73	59	79	14	29	93	71	40	87	17	40	98	43	65	55	44	58
92	70	95	23	04	60	11	42	69	74	40	54	32	52	96	71	37	02	34	91
22	31	16	71	51	60	89	59	41	92	24	54	22	40	80	28	46	33	13	80
24	47	37	40	99	03	45	02	44	75	32	53	78	34	84	20	35	17	12	50
32	98	81	38	44	23	47	10	24	38	40	40	59	84	70	46	18	38	44	70
67	26	20	68	02	62	12	20	95	63	54	39	60	08	80	91	46	49	94	21
24	35	58	05	44	73	95	26	97	17	78	78	36	03	14	58	34	89	43	72
21	36	23	09	75	00	76	44	20	45	35	14	00	41	39	97	34	81	33	95
78	17	53	25	22	75	31	47	15	94	03	90	04	42	16	14	09	53	54	92
14	39	05	42	94	35	31	47	55	58	88	24	90	17	54	24	24	29	85	57
86	56	00	48	35	71	89	07	05	44	44	37	44	80	21	58	51	54	17	58
19	40	81	88	05	94	47	49	28	73	92	13	86	52	17	77	04	89	55	40
04	52	08	03	97	35	99	14	07	97	57	32	16	24	26	79	33	27	95	64
77	45	49	87	87	42	20	72	03	46	39	67	44	55	12	32	43	93	53	49
04	43	16	71	27	33	33	11	24	94	72	18	08	44	28	32	40	62	74	34
20	49	36	41	72	30	23	88	57	77	83	89	82	47	59	65	74	04	34	14
20	73	35	29	78	31	90	01	74	31	49	72	49	14	41	16	23	57	03	54
01	70	84	71	83	51	54	49	14	92	33	48	41	43	52	01	87	17	40	98

What the computer sees

Cs231n-Lecture2-Stanford University

Viewpoints

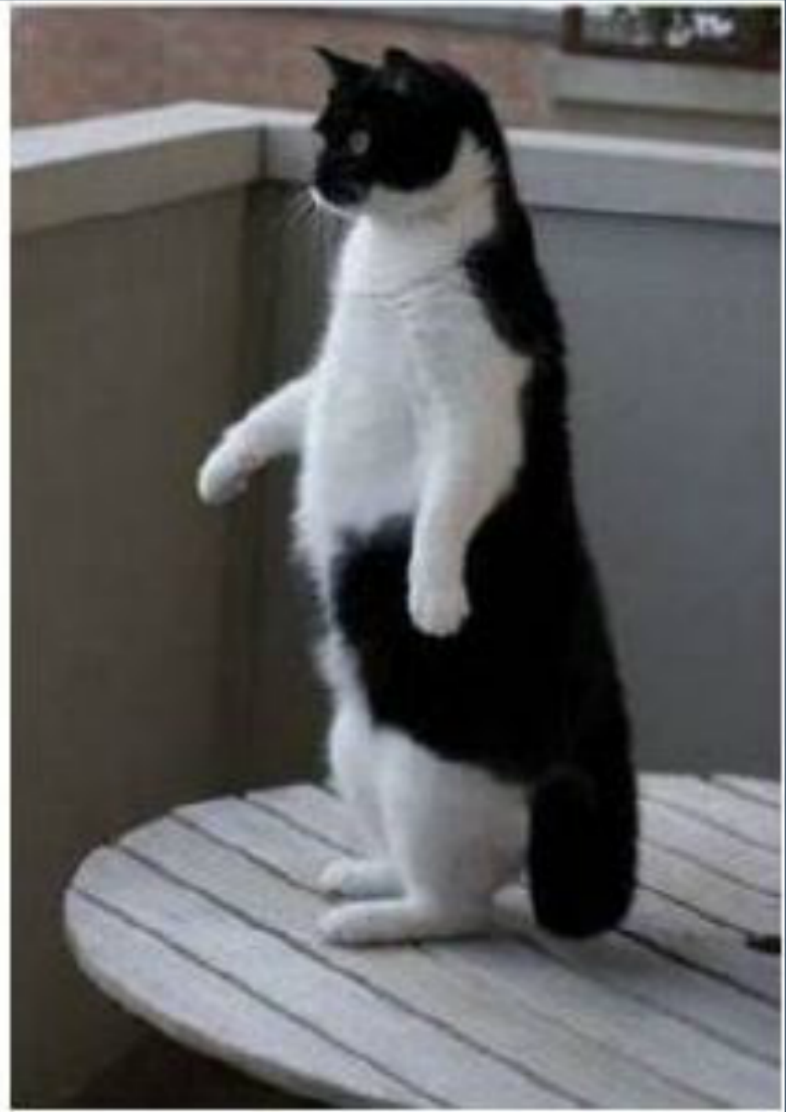


Thoma, Martin (2016).

Illumination



Deformation



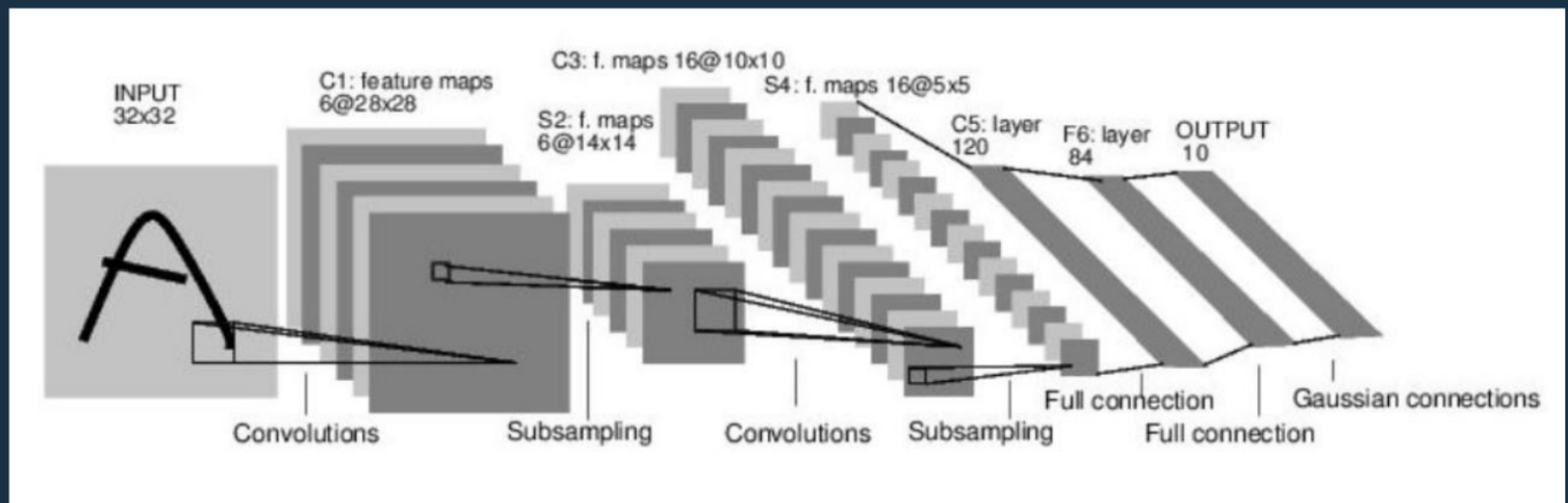
Background Clutter



Intra-Class Variation

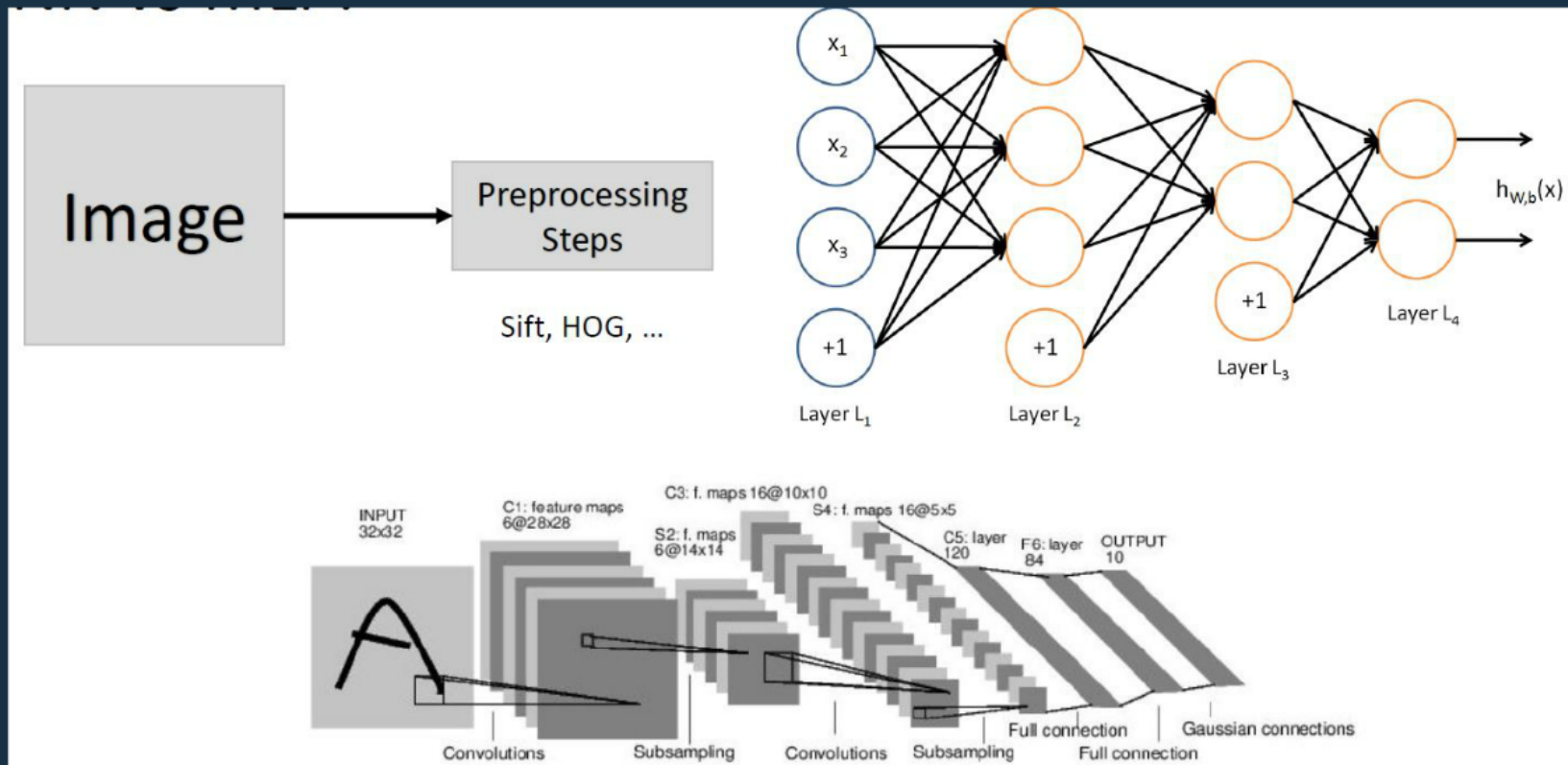


Convolutional Neural Network:



LeNet5-LeCun 1998

CNN vs MLP:



CNN Structure

Convolution Layers

No Layers

- Convolution Layer
- Subsampling Layer (Max Pooling Layer, Average Pooling Layer, ...)
- Very Basic Layer (Sigmoid, Tanh, ReLU, ...)

Convolutional Layer

Subsampling Layer

Activation Function Layer

Example of CNN Architecture

Convolutional Layer

Subsampling Layer

Activation Function Layer

Example of CNN Architecture

Convolutional Layer

Subsampling Layer

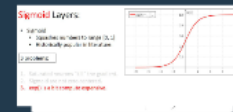
Activation Function Layer

Example of CNN Architecture

UCLAN 2017

Activation function Layer

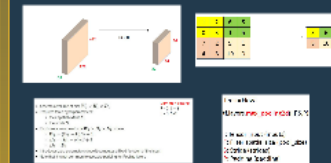
- Sigmoid
- Tanh
- ReLU
- Leaky ReLU



TensorFlow:
tf.nn.sigmoid()
tf.nn.tanh()
tf.nn.relu()

Cs231m-Lecture 5.7 Stanford

Pooling Layer



Cs231m-Lecture 5.7 Stanford

CNN on MNIST in TF:

Diagram: 28x28 input image (255 channels) is processed by a CNN to produce a 10-class output.

```

import tensorflow as tf
import numpy as np

# Define the input data
mnist = tf.keras.datasets.mnist
(train_images, train_labels), (test_images, test_labels) = mnist.load_data()

# Define the model architecture
model = tf.keras.Sequential([
    tf.keras.layers.Conv2D(10, [5, 5], activation='relu', input_shape=(28, 28, 1)),
    tf.keras.layers.MaxPooling2D([2, 2]),
    tf.keras.layers.Conv2D(10, [5, 5], activation='relu'),
    tf.keras.layers.MaxPooling2D([2, 2]),
    tf.keras.layers.Flatten(),
    tf.keras.layers.Dense(100, activation='relu'),
    tf.keras.layers.Dense(10)
])

# Compile the model
model.compile(optimizer='adam', loss='sparse_categorical_crossentropy', metrics=['accuracy'])

# Train the model
model.fit(train_images, train_labels, epochs=10)

# Evaluate the model
test_loss, test_acc = model.evaluate(test_images, test_labels)
print('Test accuracy:', test_acc)
    
```

deep visualization ToolBox

Deep Visualization Toolbox

yosinski.com/deepvis

#deepvis



Con

CNN Layers:

- Convolution Layer
- Sub sampling Layer (Max Pooling Layers, Average Pooling Layers,...)
- Non-linear layer (Sigmoid, Tanh, ReLU,...)

Example of convolution:



Prezi

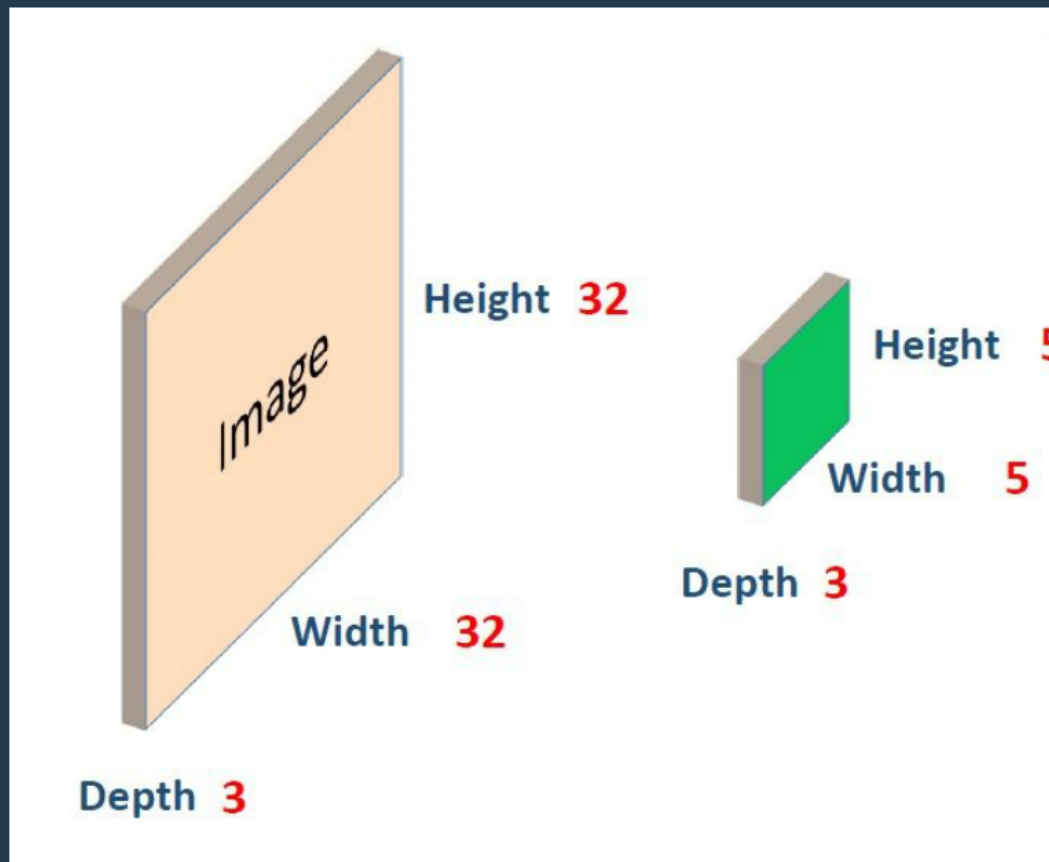


-1	-2	-1	-1	0	1
0	0	0	0	0	0

Covolutional Layer:



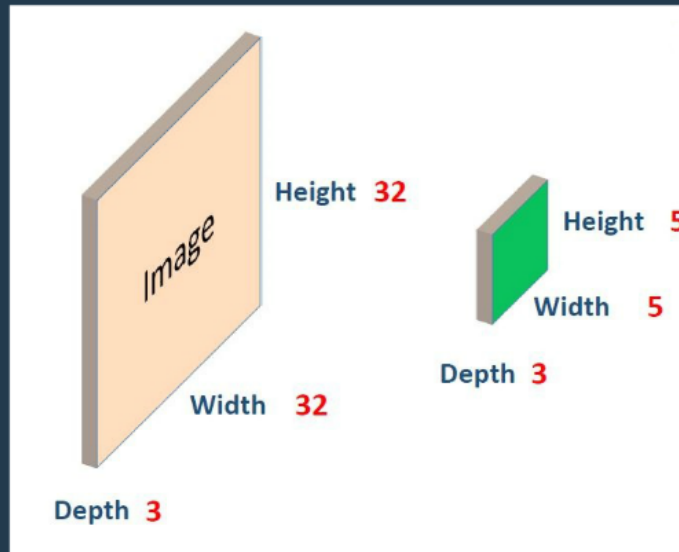
Covolutional Layer:



Credit: UTDLSS2017



Covolutional Layer:



Credit: UTDLSS2017

Example of convolution:

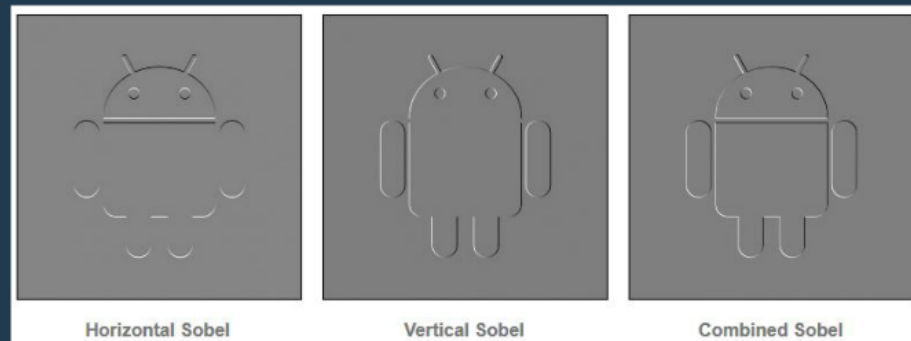


-1	-2	-1
0	0	0
1	2	1

Finds horizontals

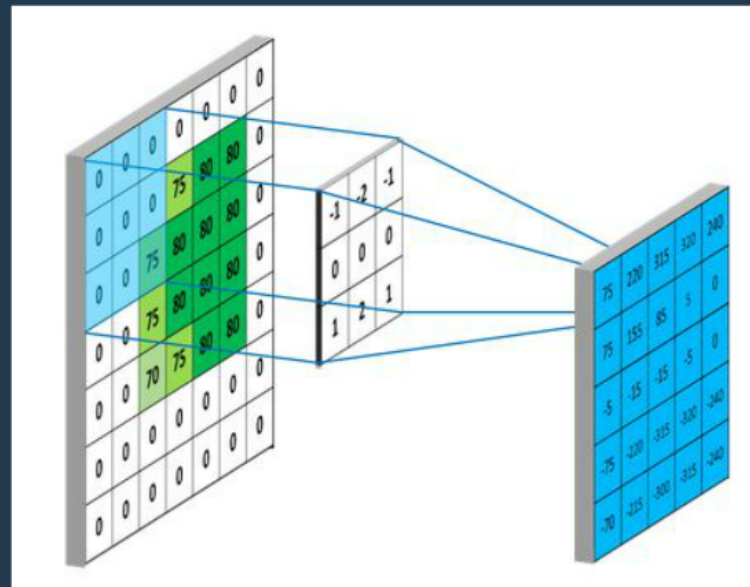
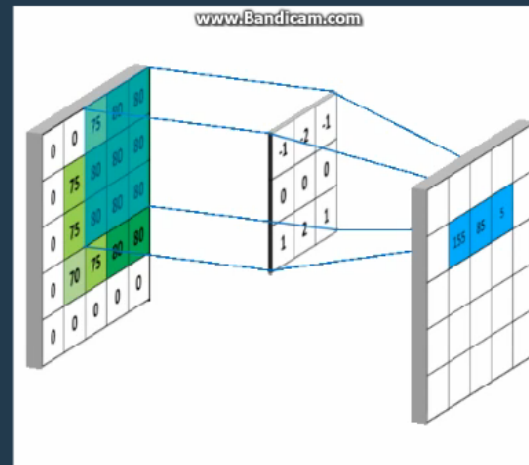
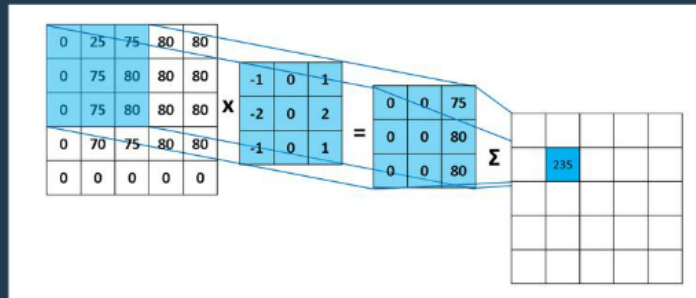
-1	0	1
-2	0	2
-1	0	1

Finds verticals



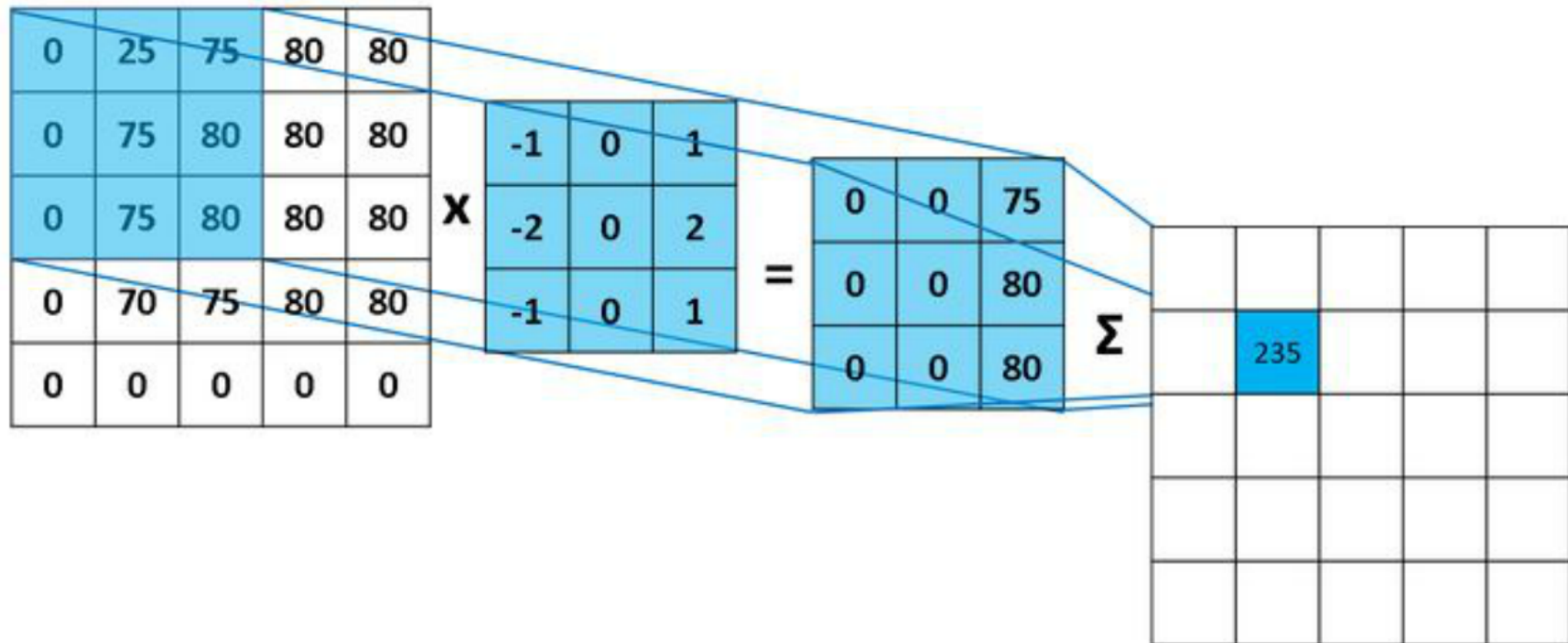
<https://mlnotebook.github.io/post/CNN1/>

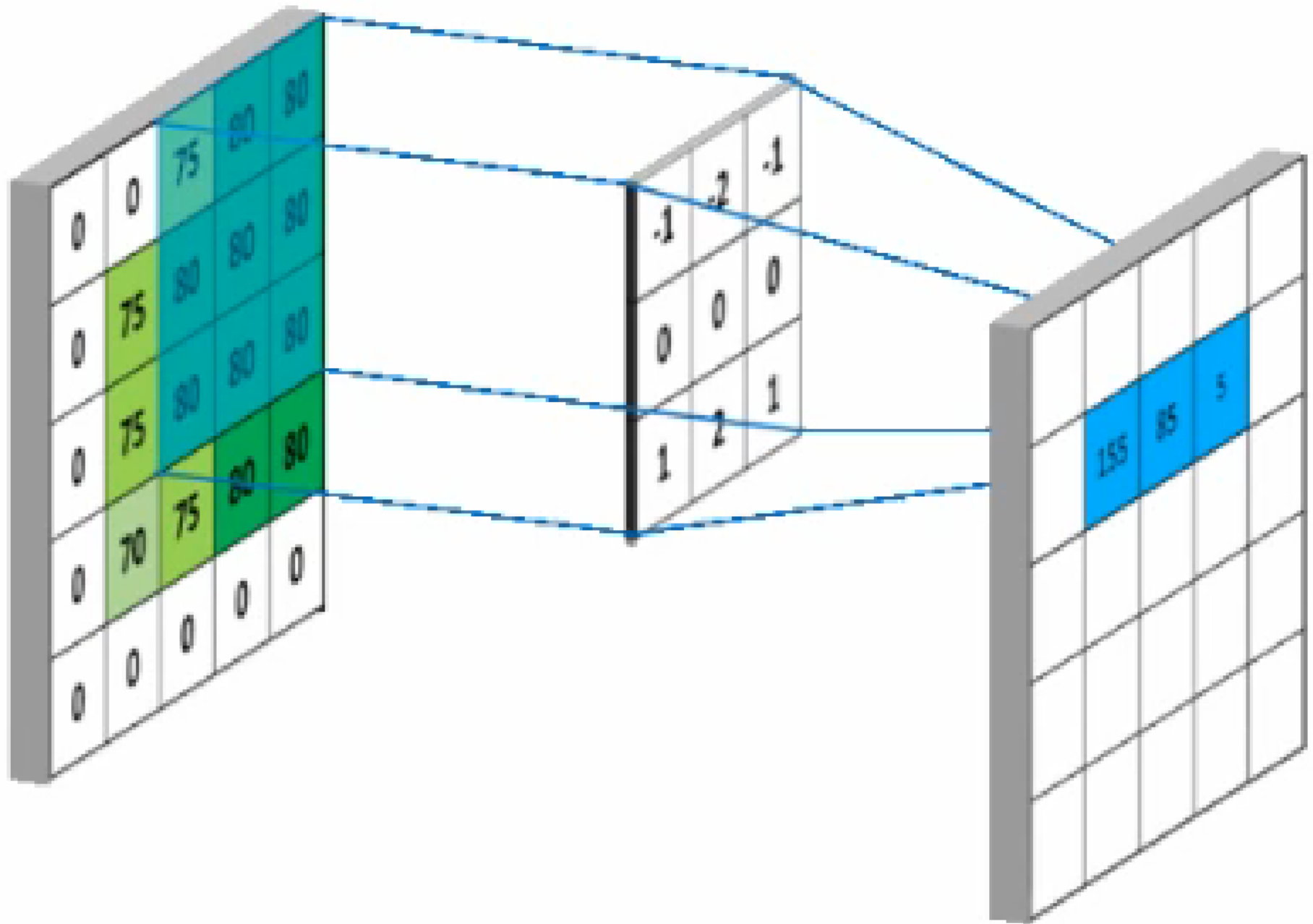
Covolutional Layer:

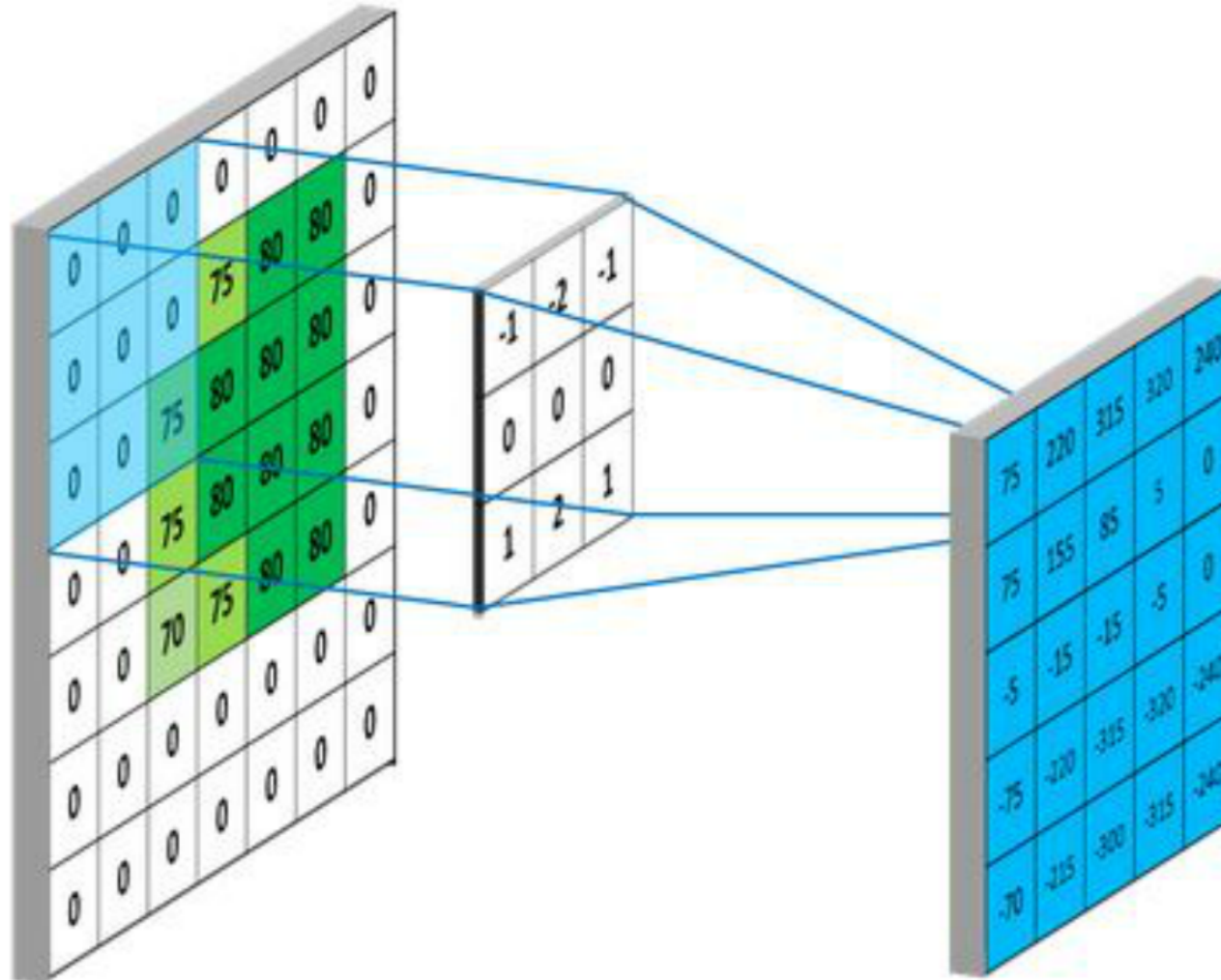


<https://mlnotebook.github.io/post/CNN1/>

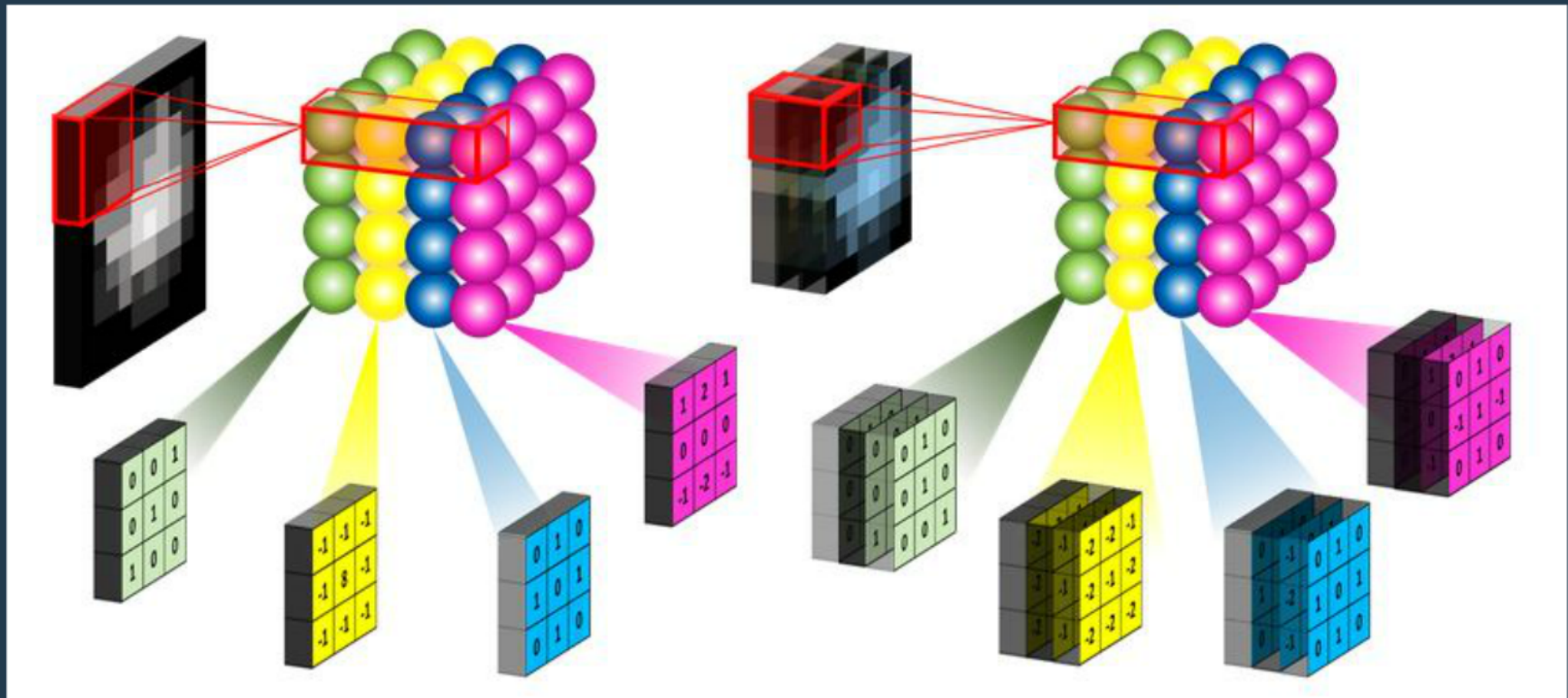
Convolutional Layer



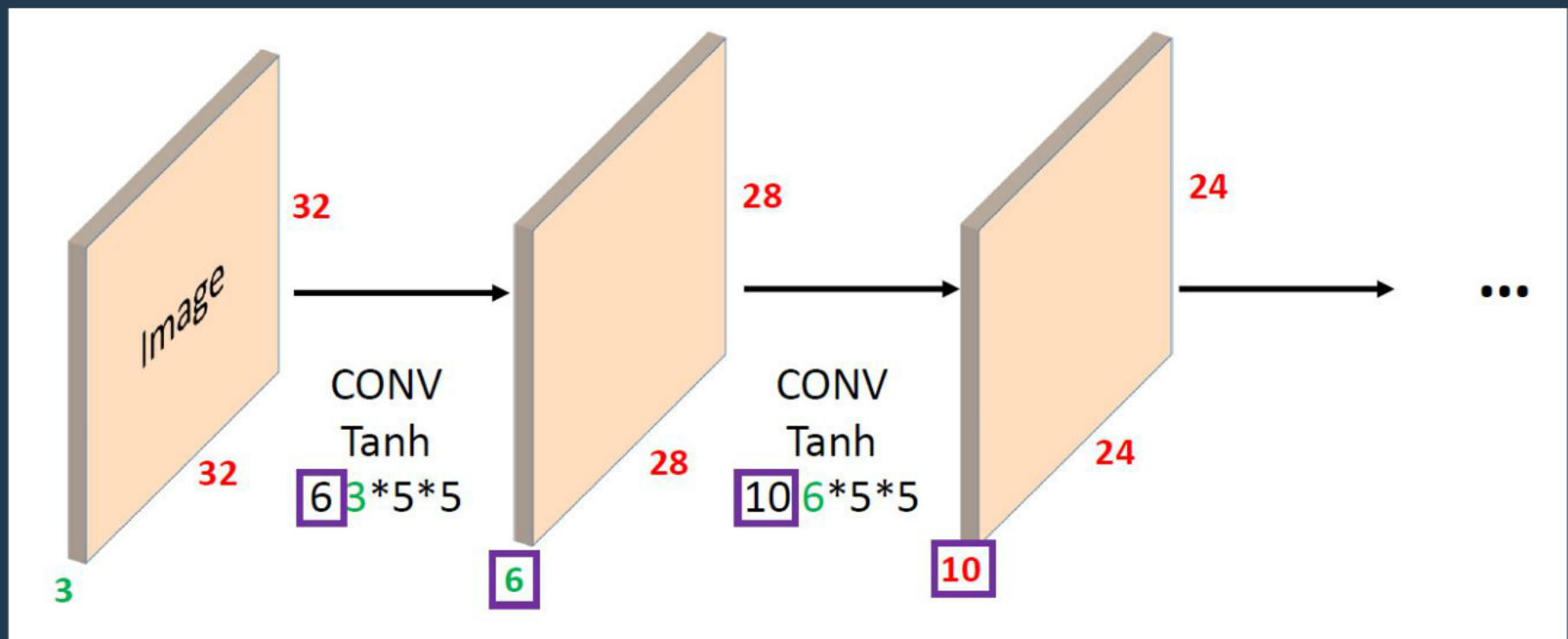




more than one Kernel:



ConvNet:





Features

Business

Explore

Marketplace

Pricing

This repository

Search

Sign in

or Sign up

vdumoulin / conv_arithmetic

Watch 197

Star 3,554

Fork 718

Code

Issues 3

Pull requests 1

Projects 0

Insights

Branch: master

conv_arithmetic / gif / no_padding_strides.gif

Find file

Copy path



John Iapsett Add dilation support

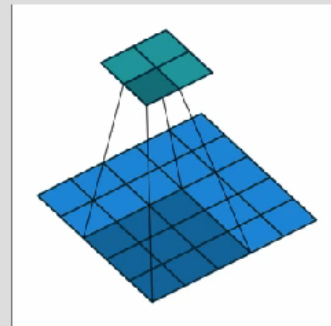
c595a98 on Oct 7, 2016

1 contributor

43.6 KB

Download

History



Closer Look at Convolution:

Practical Note: In practice, It is common to zero pad.

0	0	0	0	0	0	0	0	0
0								0
0								0
0								0
0								0
0								0
0								0
0								0
0	0	0	0	0	0	0	0	0

N

Output Size formula: $(N - F)/\text{stride} + 1$

$N = 9, S=1, F=3$

Output Size: $(9-3)/1+1=7$

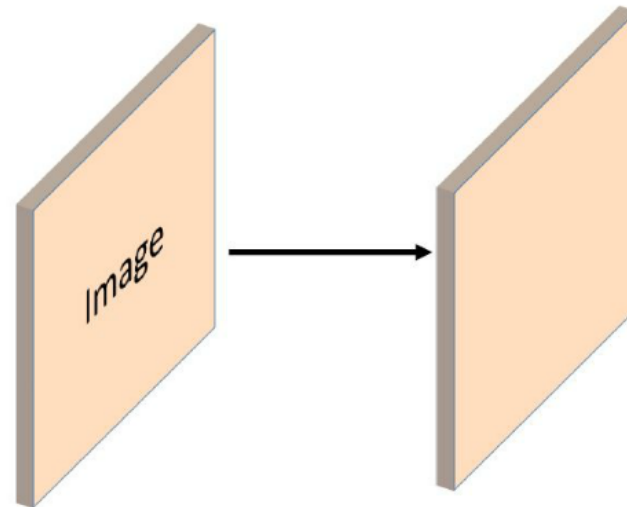
Convolution Example:

Input size: $3 \times 32 \times 32$

Filter size: $20 \times 3 \times 5 \times 5$, Padding=2 and Stride 1

Output size: $(32+4-5)/1+1=32 \rightarrow 20 \times 32 \times 32$

Number of Learnable Parameters: $20 \times 3 \times 5 \times 5 = 1500$

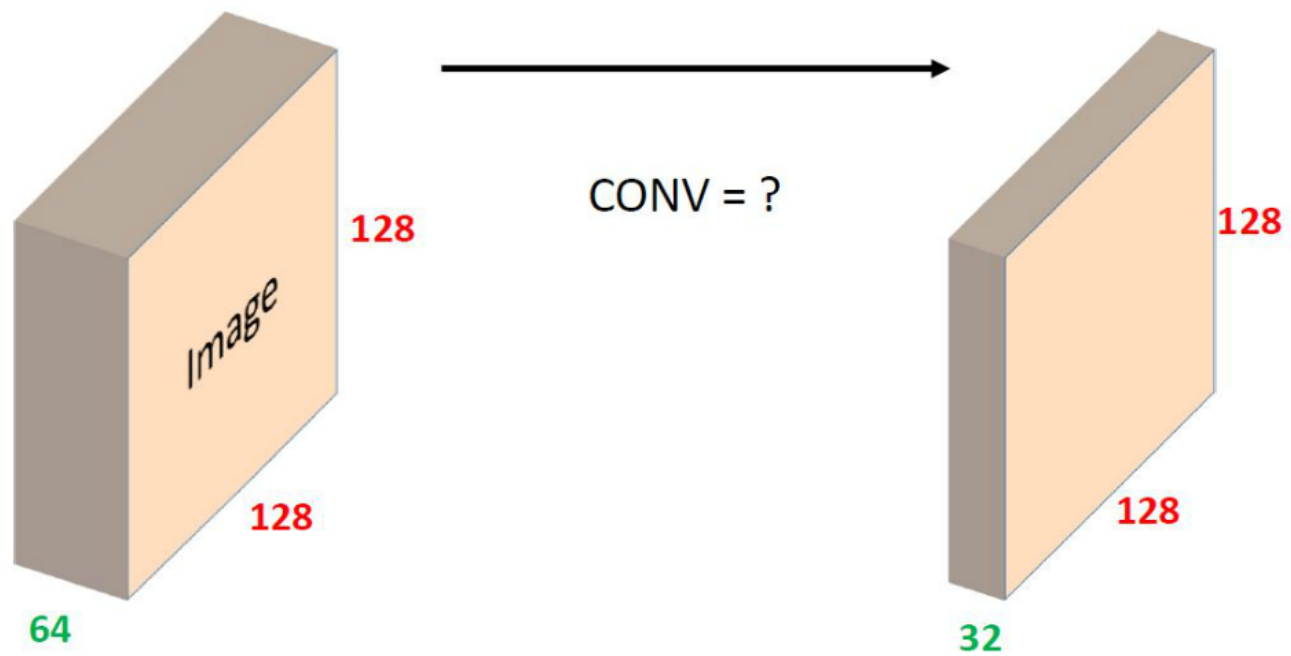


Convolution Summary:

Summary. To summarize, the Conv Layer:

- Accepts a volume of size $W_1 \times H_1 \times D_1$
- Requires four hyperparameters:
 - Number of filters K ,
 - their spatial extent F ,
 - the stride S ,
 - the amount of zero padding P .
- Produces a volume of size $W_2 \times H_2 \times D_2$ where:
 - $W_2 = (W_1 - F + 2P)/S + 1$
 - $H_2 = (H_1 - F + 2P)/S + 1$ (i.e. width and height are computed equally by symmetry)
 - $D_2 = K$
- With parameter sharing, it introduces $F \cdot F \cdot D_1$ weights per filter, for a total of $(F \cdot F \cdot D_1) \cdot K$ weights and K biases.
- In the output volume, the d -th depth slice (of size $W_2 \times H_2$) is the result of performing a valid convolution of the d -th filter over the input volume with a stride of S , and then offset by d -th bias.

Convolution Example:



TensorFlow:

```
tf.layers.conv2d(I, O, F, S, P, use_bias=True)
```

I: Tensor Input (inputs)

O: Output cube depth (filters)

F: Filter spatial size (Kernel_size)

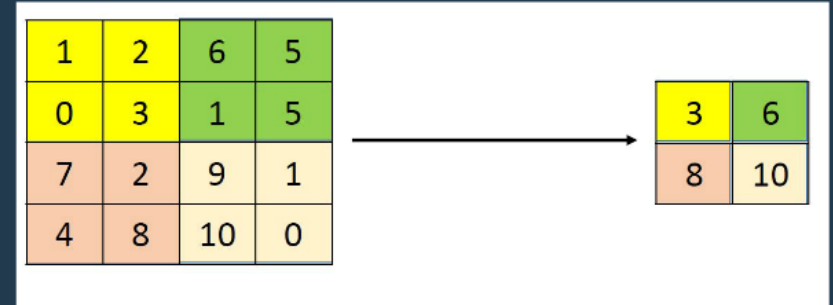
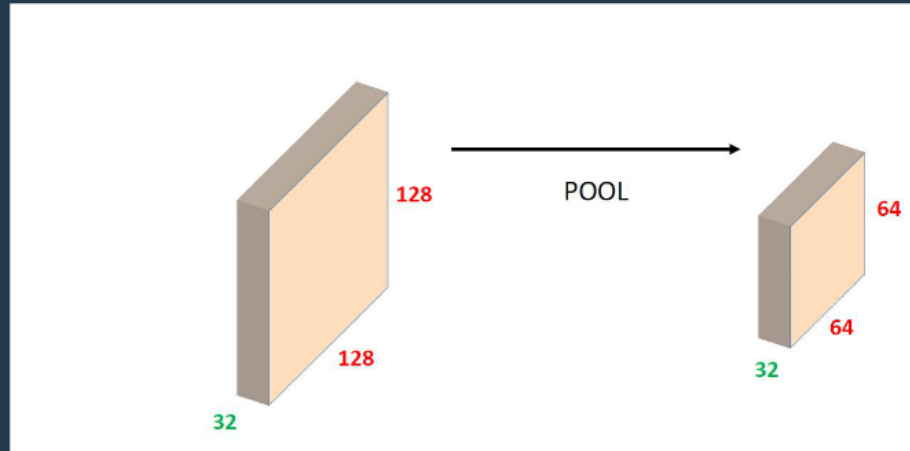
S: Stride (strides)

P: Padding (padding)

Use_bias: If is true, then you define a bias term for each filter

...

Pooling Layer



- Accepts a volume of size $W_1 \times H_1 \times D_1$
- Requires three hyperparameters:
 - their spatial extent F ,
 - the stride S ,
- Produces a volume of size $W_2 \times H_2 \times D_2$ where:
 - $W_2 = (W_1 - F) / S + 1$
 - $H_2 = (H_1 - F) / S + 1$
 - $D_2 = D_1$
- Introduces zero parameters since it computes a fixed function of the input
- Note that it is not common to use zero-padding for Pooling layers

Common Setting:

$F = 2, S = 2$

$F = 3, S = 2$

TensorFlow:

`tf.layers.max_pooling2d(I, F, S, P)`

I: Tensor Input (inputs)

F: Filter spatial size (pool_size)

S: Stride (strides)

P: Padding (padding)

Cs231n-Lecture7 Stanford

Activation function Layer

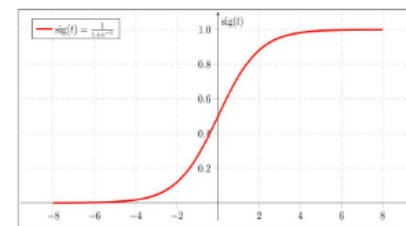
- Sigmoid
- Tanh
- ReLU
- Leaky ReLU

Sigmoid Layers:

- Sigmoid
 - Squashes numbers to range [0, 1]
 - Historically popular in literature.

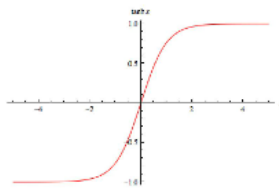
3 problems:

1. Saturated neurons "kill" the gradient.
2. Sigmoid are not zero centered.
3. $\exp()$ is a bit compute expensive.



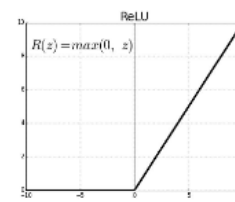
Tanh Layers:

- Tanh
 - Squashes numbers to range [-1, 1]
 - Zero centered (nice)
 - Still kill gradient when saturated



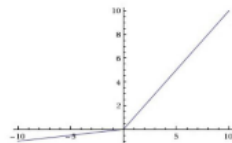
ReLU Layers:

- ReLU
 - Does not saturate (in +region)
 - Very computationally efficient
 - Converges much faster than sigmoid/tanh in practice (e.g. 6x)
 - Not zero centered output
 - An annoyance



Leaky ReLU Layers:

- ReLU
 - Does not saturate (in +region)
 - Very computationally efficient
 - Converges much faster than sigmoid/tanh in practice (e.g. 6x)
 - Will not die



TensorFlow:

`tf.nn.sigmoid()`

`tf.nn.tanh()`

`tf.nn.ReLU()`

CS231n-Lecture 5, 7 Stanford

- Sigmoid
- Tanh
- ReLU
- Leaky ReLU

Sigmoid Layers:

Sigmoid

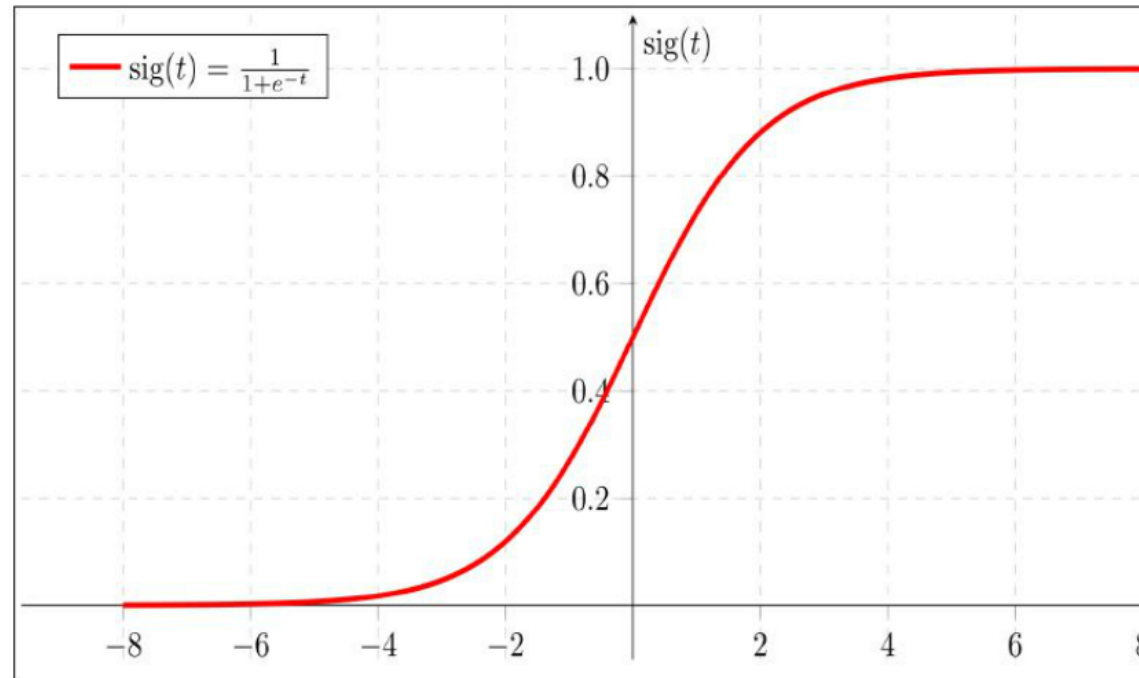
- Squashes numbers to range $[0, 1]$
- Historically popular in literature.

problems:

Saturated neurons “kill” the gradient.

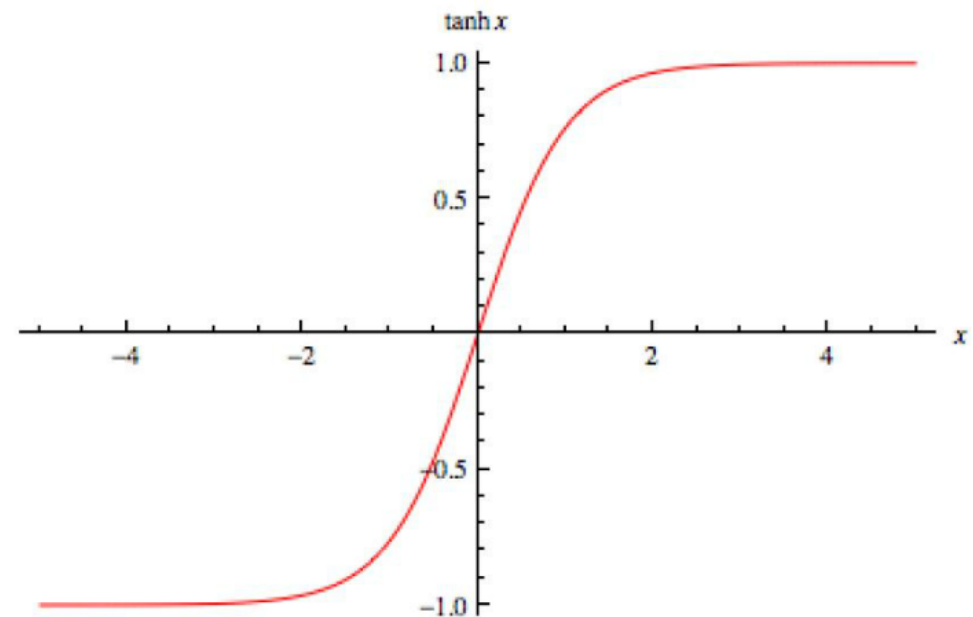
Sigmoid are not zero centered.

$\exp()$ is a bit compute expensive.



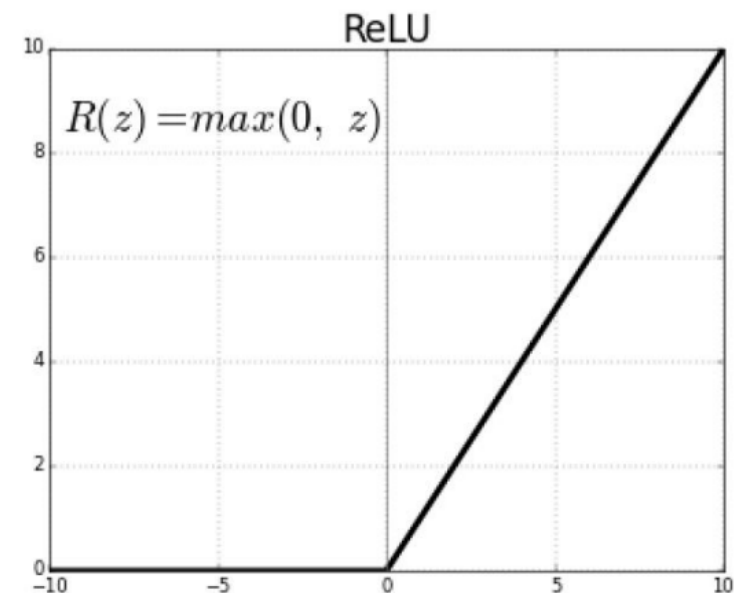
Tanh Layers:

- Tanh
 - Squashes numbers to range $[-1, 1]$
 - Zeros centered (nice)
 - Still kill gradient when saturated 😞

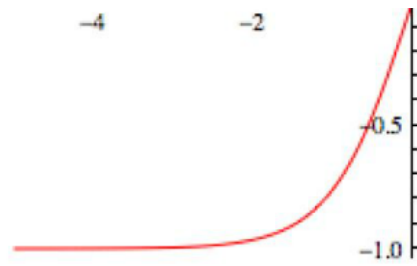


ReLU Layers:

- ReLU
 - Does not saturate (in +region)
 - Very computationally efficient
 - Converges much faster than sigmoid/tanh in practice (e.g. 6x)
 - Not zero centered output
 - An annoyance

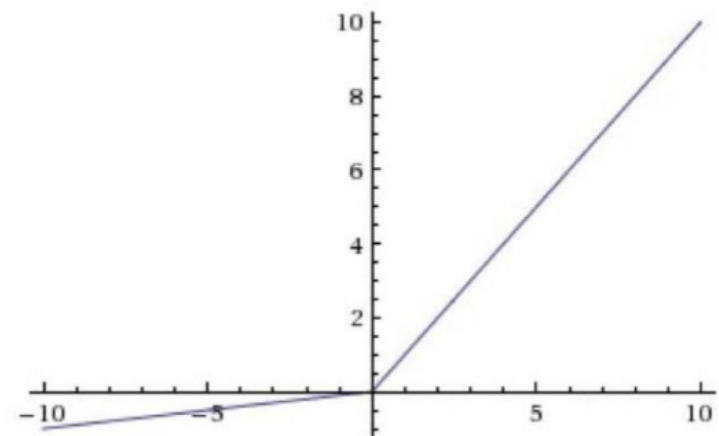


TensorFlow:



Leaky ReLU Layers:

- ReLU
 - Does not saturate (in +region)
 - Very computationally efficient
 - Converges much faster than sigmoid/tanh in practice (e.g. 6x)
 - Will not die



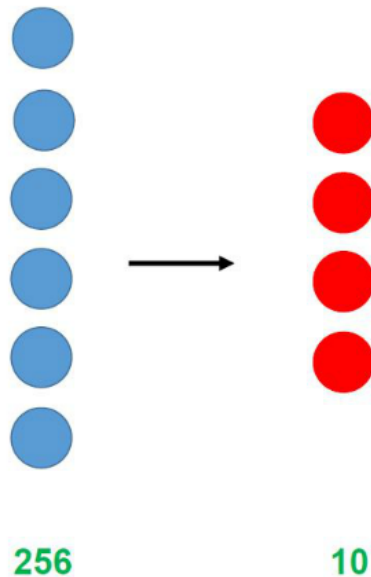
TensorFlow:

`tf.nn.sigmoid()`

`tf.nn.tanh()`

`tf.nn.ReLU()`

CNN on MNIST in TF:



```

input_layer = tf.reshape(x, [-1, 28, 28, 1])
# Define Convolution(& Pooling) Layer1
conv1 = tf.layers.conv2d(inputs=input_layer,
                        filters = 16,
                        kernel_size=[5, 5],
                        strides=[1,1],
                        activation=tf.nn.relu)
pool1 = tf.layers.max_pooling2d(inputs=conv1,
                                pool_size=[2,2],
                                strides=2)
# Define Convolution(& Pooling) Layer2
conv2 = tf.layers.conv2d(inputs=pool1,
                        filters = 64,
                        kernel_size=[5, 5],
                        strides=[1, 1],
                        activation=tf.nn.relu)
pool2 = tf.layers.max_pooling2d(inputs=conv2,
                                pool_size=[2, 2],
                                strides=2)
# Multiply The Shape of Previous Layer
dim = np.prod(pool2.get_shape().as_list()[1:])
# Define Fully Connected Layers
relu2_flat = tf.reshape(pool2, [-1, dim])
fc1 = tf.layers.dense(inputs=relu2_flat,
                    units=256,
                    activation=tf.nn.relu)
fc2 = tf.layers.dense(inputs=fc1,
                    units=n_classes)

```

deep visualization ToolBox

Deep Visualization Toolbox

yosinski.com/deepvis

#deepvis



Jason Yosinski



Jeff Clune



Anh Nguyen



Thomas Fuchs

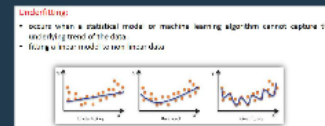
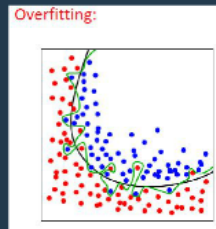


Hod Lipson



Advance info about Learning

- Learning Tricks
- Dropout
- Weight Initialization
- Batch Normalization



Learning Tricks:

Some Rules for train a classifier

Split your data into 3 parts:

1. Training data: for learning the weight of the network
2. Validation data: for detecting overfitting
3. Test data: for measuring the accuracy of your model

Validation data is not only used for checking overfitting, but it is useful for defining hyper-parameters (structural parameters):

1. # of hidden neurons
2. Best network architecture
3. Learning Rate
4. ...

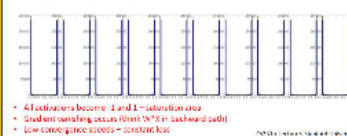
Regularization Technique:

- L2 norm
- L1 norm
- Dropout!



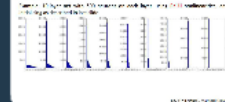
Covariate shift:

Example: 20-layer net with 500 neurons in each layer, using leaky ReLU, and initializing as described in last slide.



- All activations become 1 and 0 = saturation area
- Gradient vanishing occurs (then 50% in backward path)
- Low convergence speeds = stochastic loss

Covariate shift:

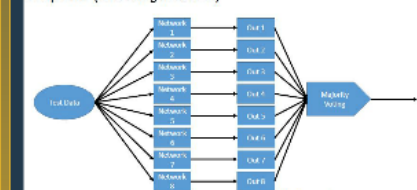


Batch Normalization:

Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift

Sergey Ioffe, Christian Szegedy - 2015

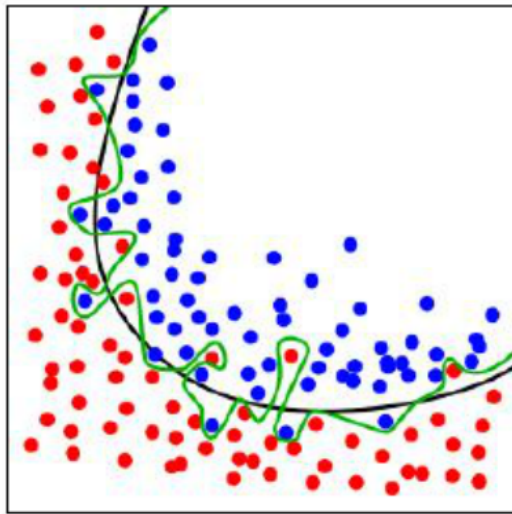
Dropout: (Overfitting Reduction)



a simple way to prevent neural networks from over fitting (nitish srivastava, 2014)

Advance info about Lea

Overfitting:



Underfitting:

- occurs when a statistical model or machine learning algorithm cannot capture the underlying trend of the data.
- fitting a linear model to non-linear data

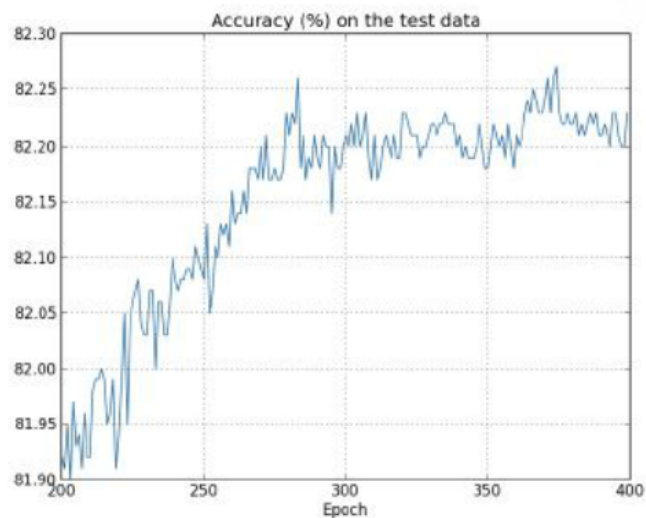
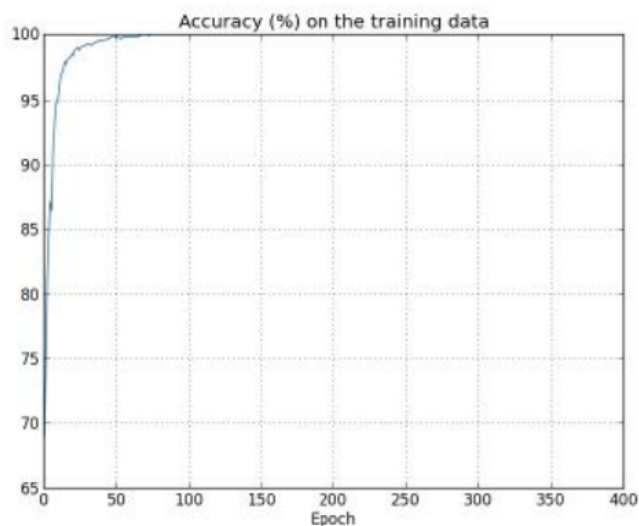


Regularization Technique:

- L2 norm
- L1 norm
- Dropout!

Learning Tricks:

Working as a detective:





Learning Tricks:

Some Rules for train a classifier

Split your data into 3 parts:

1. Training data: for learning the weight of the network
2. Validation data: For detecting overfitting
3. Test data: For measuring the accuracy of your model

Validation data is not only used for checking overfitting, but it is usable for defining **hyper-parameters** (structural parameters):

1. # epochs to train for
2. Best network architecture
3. Learning Rate
4. ...

Covariate shift:

Example: 10-layer net with 500 neurons on each layer, using tanh nonlinearities, and initializing as described in last slide

Regularization Technique:

- L2 norm
- L1 norm
- Dropout!

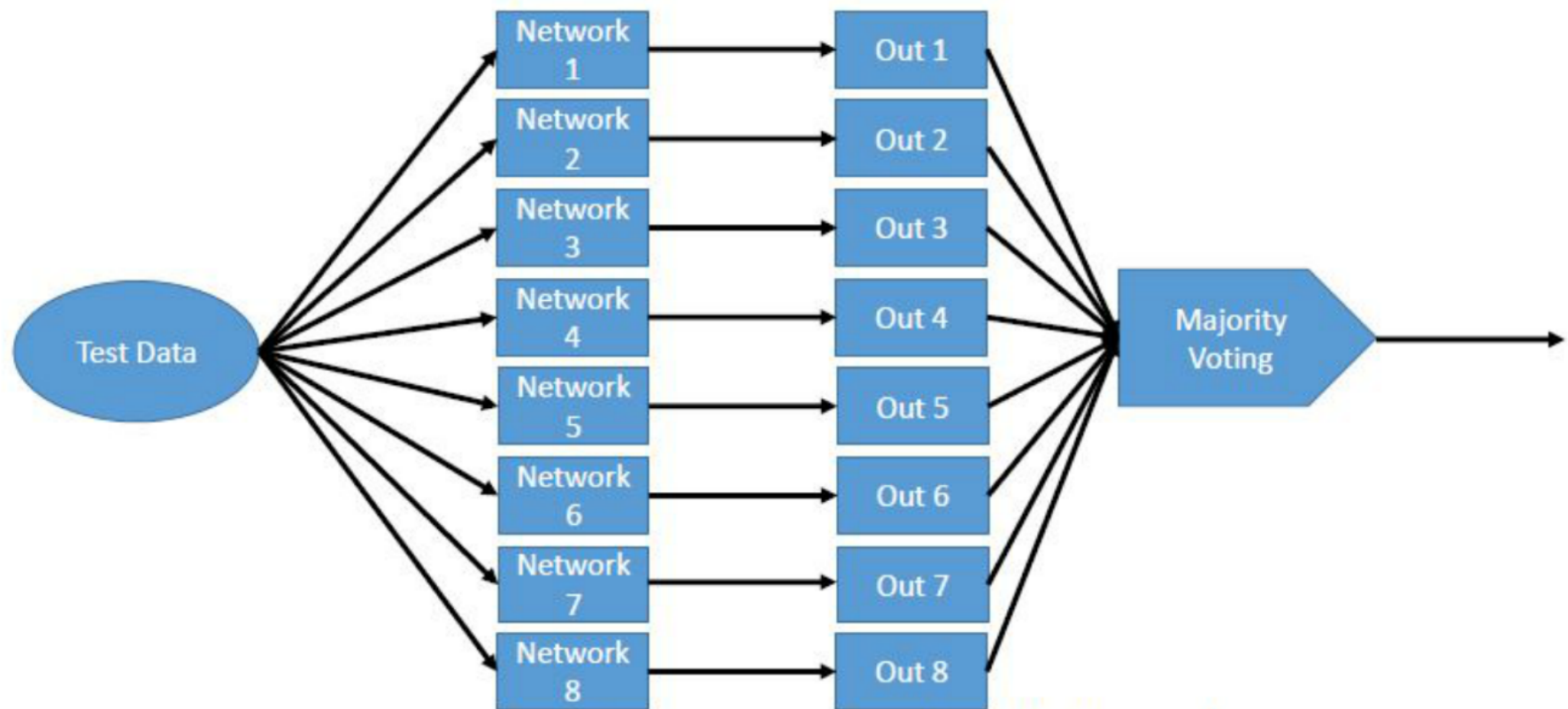




Prezi

Batch Normalization:

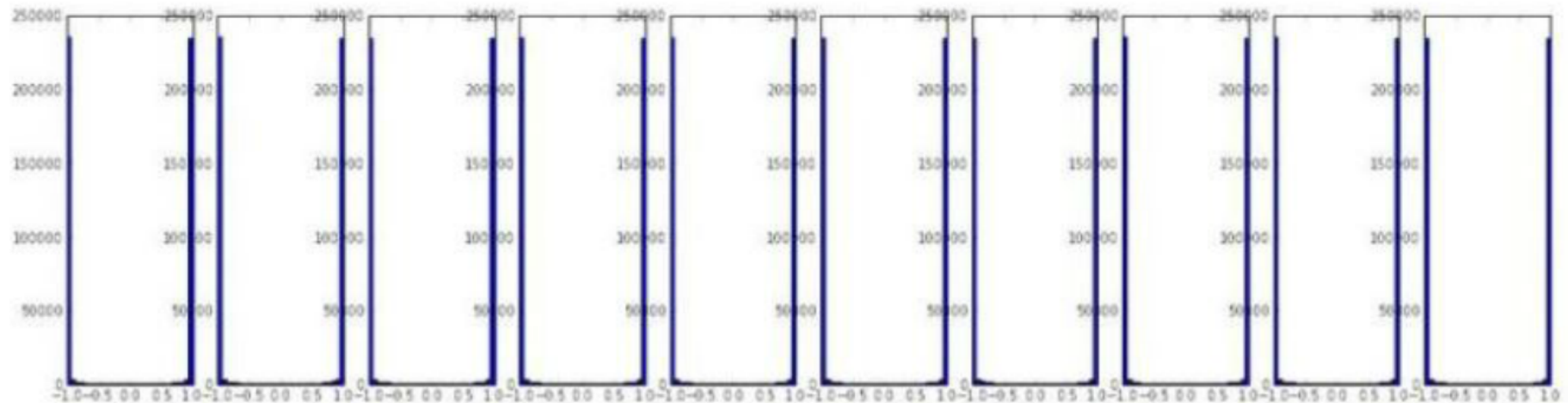
Dropout: (Overfitting Reduction)



a simple way to prevent neural networks from over fitting (nitish sirivasta, 2014)

Covariate shift:

Example: 10-layer net with 500 neurons on each layer, using tanh nonlinearities, and initializing as described in last slide

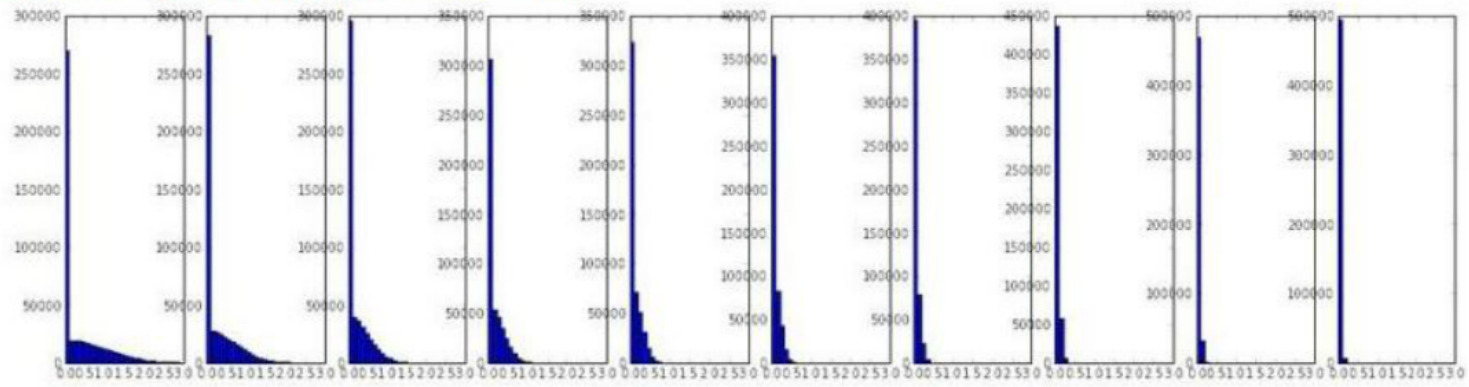


- All activations become -1 and 1 – saturation area
- Gradient vanishing occurs (think $W \cdot X$ in backward path)
- Low convergence speeds – constant loss

Cs231n-Lecture5-Stanford University

Covariate shift:

Example: 10-layer net with 500 neurons on each layer, using **ReLU** nonlinearities, and initializing as described in last slide



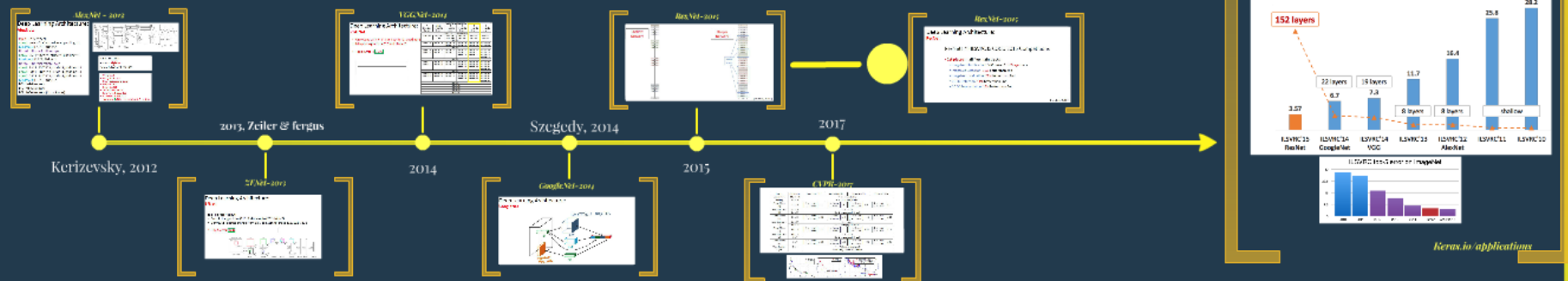
Cs231n-Lecture5-Stanford University

Batch Normalization:

Batch Normalization: **Accelerating** Deep Network Training by
Reducing Internal Covariate Shift

Sergey Ioffe, Christian Szegedy - 2015

Architecture of CNN



AlexNet - 2012

Deep Learning Architecture:

AlexNet:

Input: [3*227*227]

Conv1: 96 F=3*11*11, stride=4, padding=0

MaxPool1: F=3*3, Stride=2

Norm1 = Normalization Layer

Conv2: 256 F=96*5*5, stride=1, padding=2

MaxPool2: F=3*3, Stride=2

Norm2 = Normalization Layer

Conv3: 384 F=256*3*3, stride=1, padding=1

Conv4: 384 F=384*3*3, stride=1, padding=1

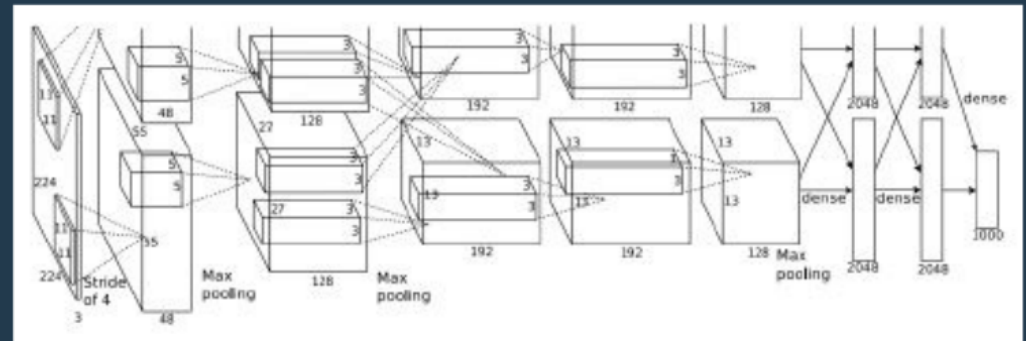
Conv5: 256 F=384*3*3, stride=1, padding=1

MaxPool3: F=3*3, Stride=2

FC6: 4096 Neurons

FC7: 4096 Neurons

FC8: 1000 Neurons (Class Scores)



Input: 256*13*13

Formula: $\frac{(N-F)}{S} + 1$

Output (MaxPool3): 256*6*6

- Using ReLU
- Using Norm Layer
- Heavy data augmentation
- Dropout = 0.5
- Batch size 128
- SGD momentum 0.9
- Lr=0.01, with decay 0.1
- L2 weight decay 0.0005
- Top 5 error: 18.2% -> ensemble 7 CNN (15.04%)

ZFNet-2013

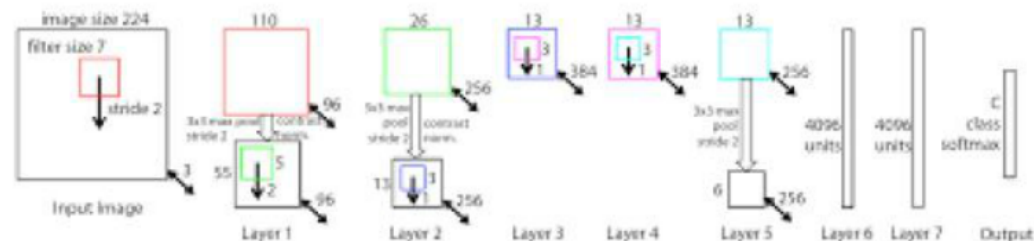
Deep Learning Architecture:

ZFNet:

It is like AlexNet but:

- Conv1 changes from (11*11 stride 4) to (7*7 stride 2)
- Conv3,4,5: depths change from 384, 384, 256 filters to 512, 1024, 512

- Top 5 error: 14.8



VGGNet-2014

Deep Learning Architecture:

VGGNet

- Only Conv with $F=3*3$ and stride=1, padding=1
- Maxpooling with $F=2*2$ and Stride 2
- Top 5 Error is: **7.3%**

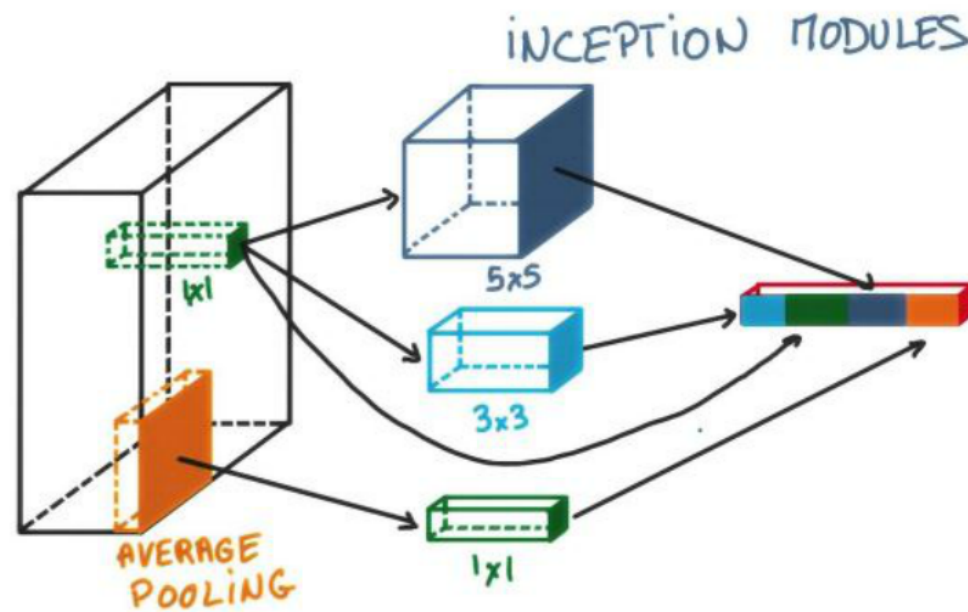
...

ConvNet Configuration					
A	A-LRN	B	C	D	E
11 weight layers	11 weight layers	13 weight layers	16 weight layers	16 weight layers	19 weight layers
input (224×224 RGB image)					
conv3-64	conv3-64 LRN	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64
maxpool					
conv3-128	conv3-128	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128
maxpool					
conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256 conv1-256	conv3-256 conv3-256 conv3-256	conv3-256 conv3-256 conv3-256 conv3-256
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 conv1-512	conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512 conv3-512
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 conv1-512	conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512 conv3-512
maxpool					
FC-4096					
FC-4096					
FC-1000					
soft-max					

GoogleNet-2014

Deep Learning Architecture:

GoogleNet

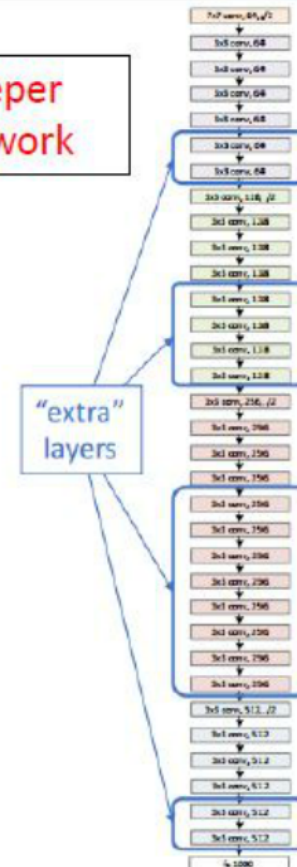


ResNet-2015

Shallow Network



Deeper Network



[He et al., 2015]

ResNet-2015

Deep Learning Architecture:

ResNet:

ResNets @ ILSVRC & COCO 2015 Competitions

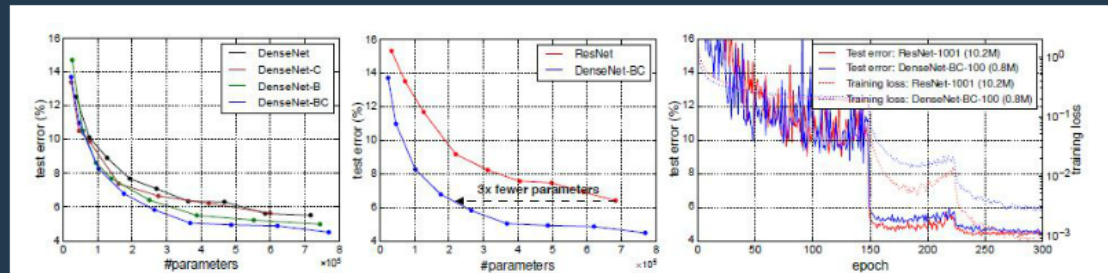
- **1st places in all five main tracks**

- ImageNet Classification: “Ultra-deep” 152-layer nets
- ImageNet Detection: 16% better than 2nd
- ImageNet Localization: 27% better than 2nd
- COCO Detection: 11% better than 2nd
- COCO Segmentation: 12% better than 2nd

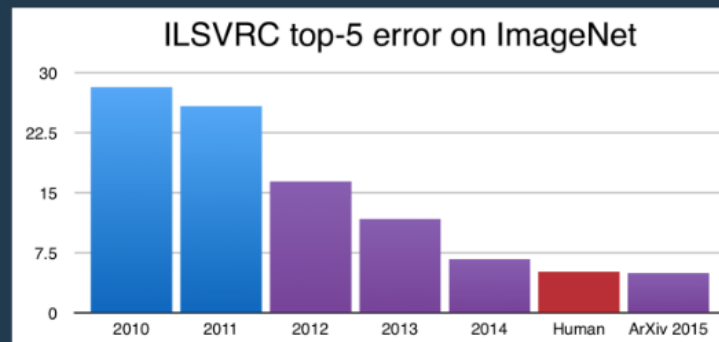
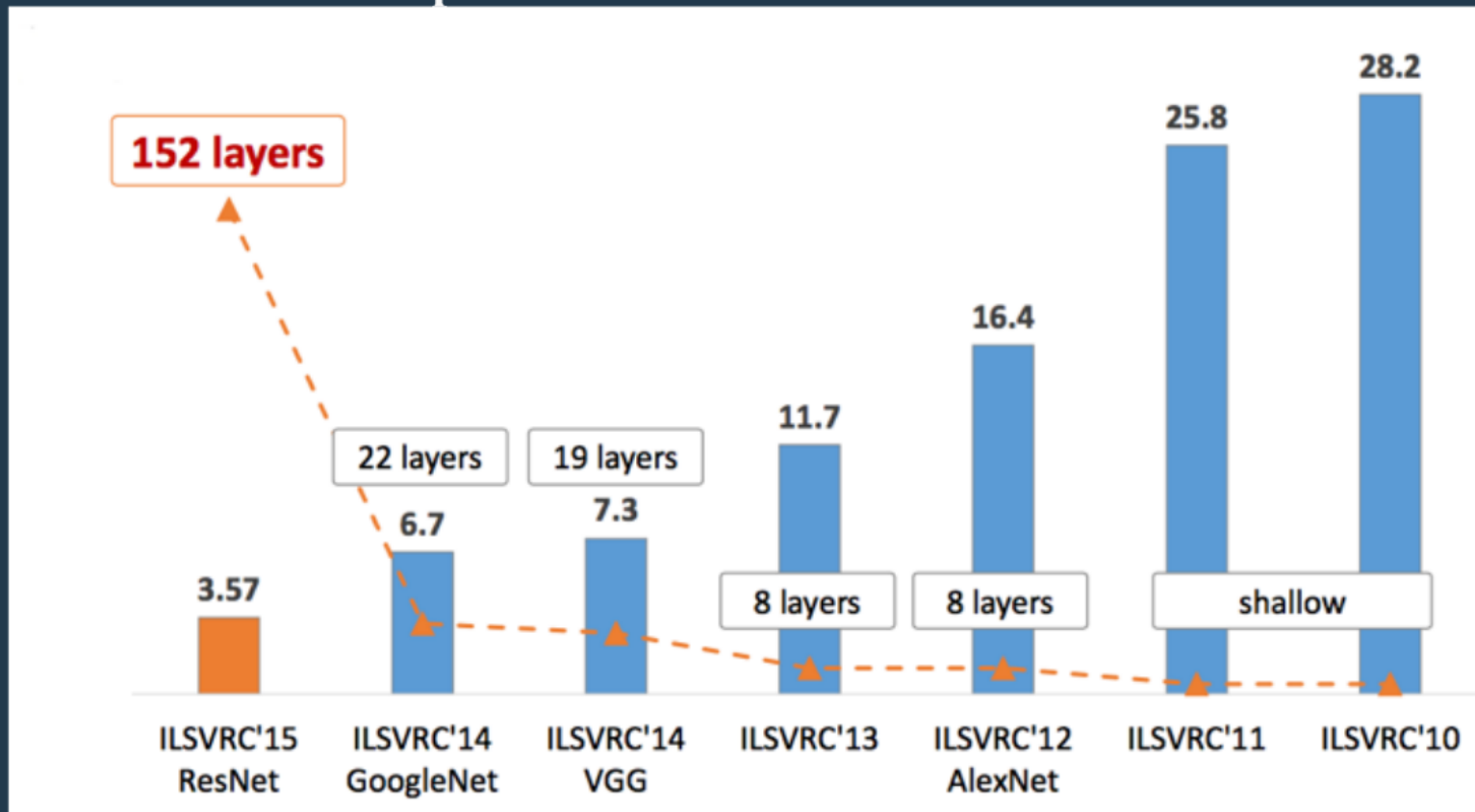
[He et al., 2015]

CVPR-2017

Layers	Output Size	DenseNet-121	DenseNet-169	DenseNet-201	DenseNet-264
Convolution	112×112	7×7 conv, stride 2			
Pooling	56×56	3×3 max pool, stride 2			
Dense Block (1)	56×56	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 6$	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 6$	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 6$	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 6$
Transition Layer (1)	56×56	1×1 conv			
	28×28	2×2 average pool, stride 2			
Dense Block (2)	28×28	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 12$	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 12$	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 12$	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 12$
Transition Layer (2)	28×28	1×1 conv			
	14×14	2×2 average pool, stride 2			
Dense Block (3)	14×14	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 24$	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 32$	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 48$	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 64$
Transition Layer (3)	14×14	1×1 conv			
	7×7	2×2 average pool, stride 2			
Dense Block (4)	7×7	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 16$	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 32$	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 32$	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 48$
Classification Layer	1×1	7×7 global average pool			
		1000D fully-connected, softmax			



Deep Architectures Performance



Keras.io/applications

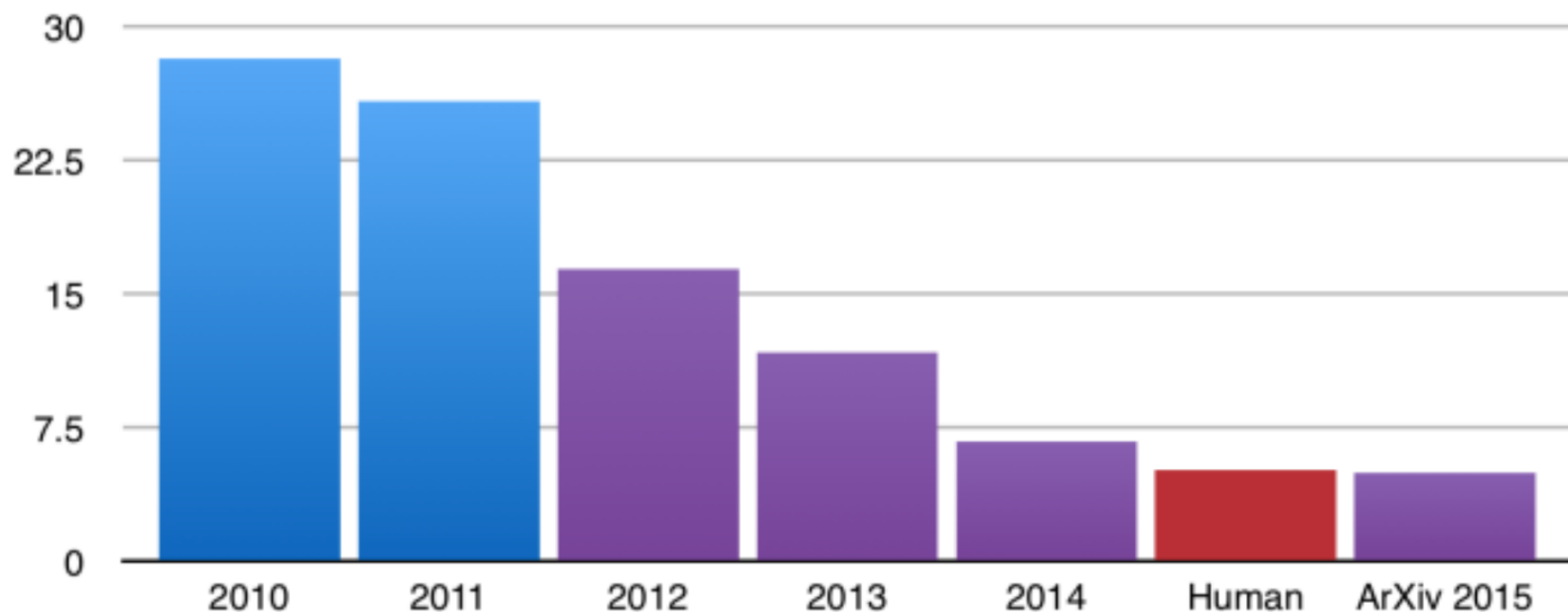
ILSVRC'14
GoogleNet

ILSVRC'14
VGG

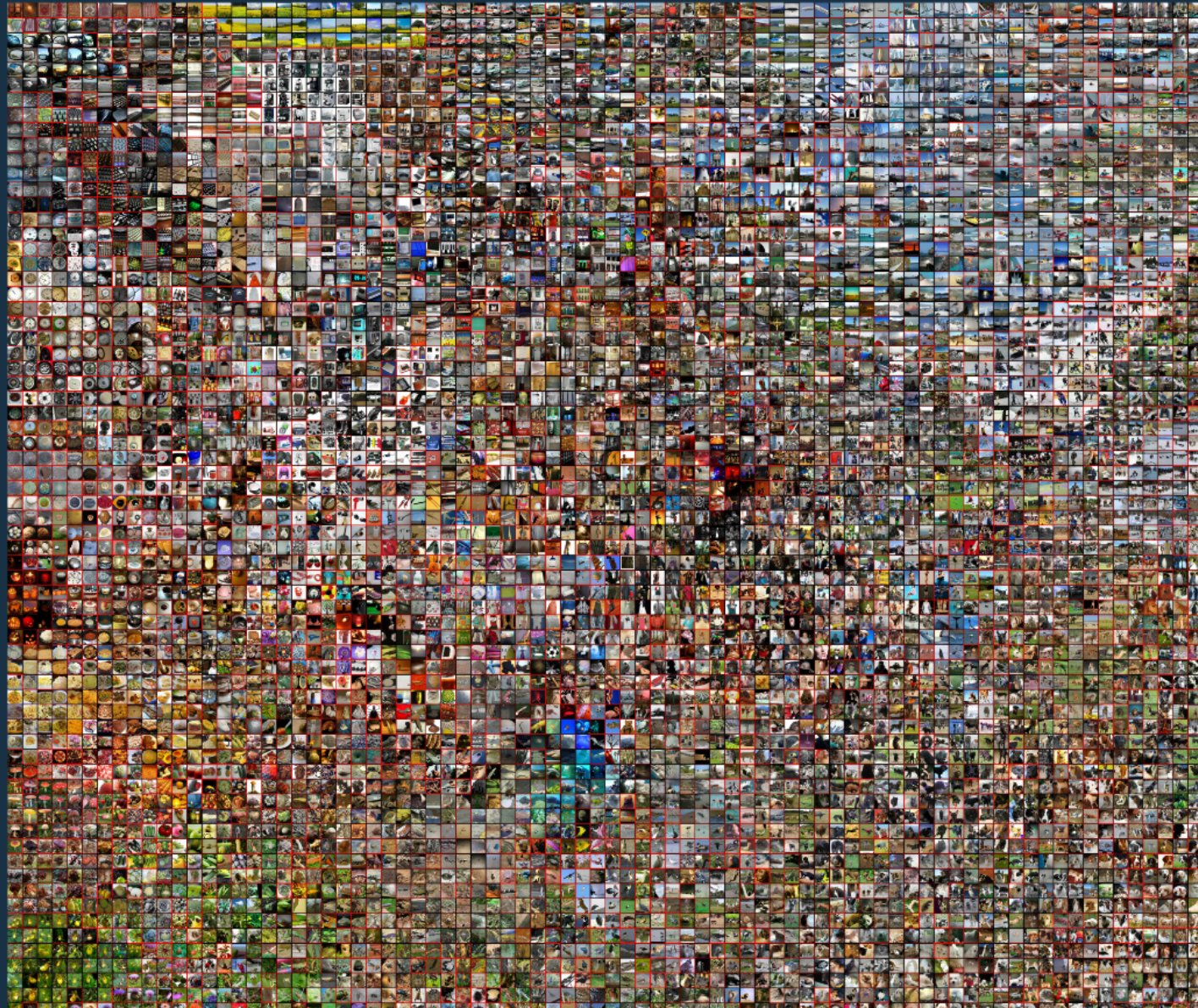
ILSVRC'13

ILSVRC'12
AlexNet

ILSVRC top-5 error on ImageNet



t-SNE visualization of CNN codes



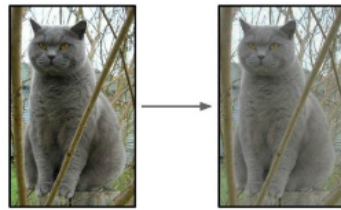
Data Augmentation

- change the pixel without changing the label
- Train on transformed data
- very widely used

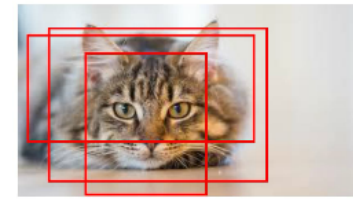
1. Horizontal flips



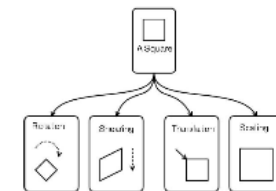
1. Horizontal flips
2. Random Crops / Scales
3. Random jitter



1. Horizontal flips
2. Random Crops / Scales



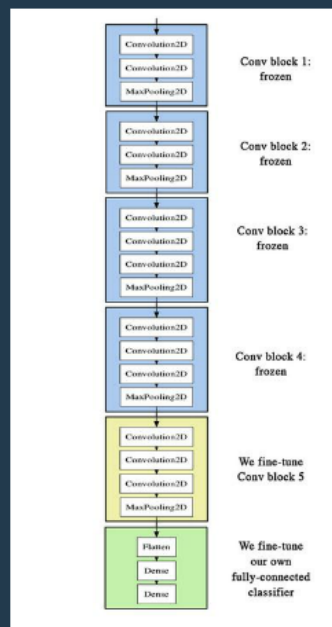
1. Horizontal flips
2. Random Crops / Scales
3. Random jitter
4. Be Creative:
 - Random Combination of:
 1. Translation
 2. Rotation
 3. Stretching
 4. Shearing
 5. Lens distortion



Transfer Learning

Transfer Learning or Domain Adaptation:

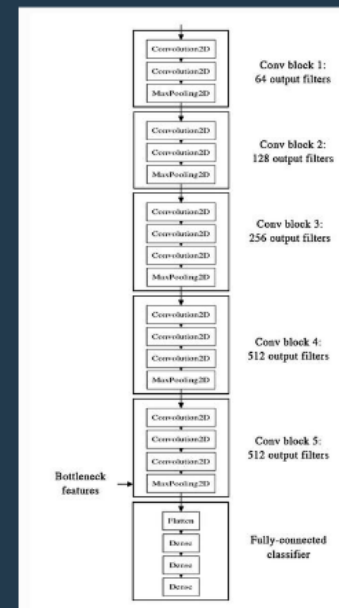
- Very few people train an entire CNN from scratch, because we do not have enough data
- It is common to use from pretrained model trained on very large dataset which contains 1.2 million large images with 1000 categories.
- Using ConvNet as an initialization or feature extractor!
- Two different scenario for using pretrained model:
 - ConvNet as fixed feature extractor: **CNN Codes**
 - Fine-tuning the Convnet:
 - Replace the last fully-connected layers with new one with random weight and train their weights again.
 - If you have more data, you can retrain more layers with “backprop”



Transfer Learning or Domain Adaptation:

When and how to fine tune:

1. New dataset is small and is similar to original dataset:
 - No need to retrain high level features in CNN – If you do, you overfit
 - Just need to retrain the classifier part of model
2. New dataset is large and is similar to original dataset:
 - Retrain all the weights in the network
3. New dataset is small and are totally different from original dataset:
 - Train the classifier not on top of the network but on top of earlier layers
4. New dataset is large and are totally different:
 - Just fine tune all the weights but using the pretrained model as an initialization



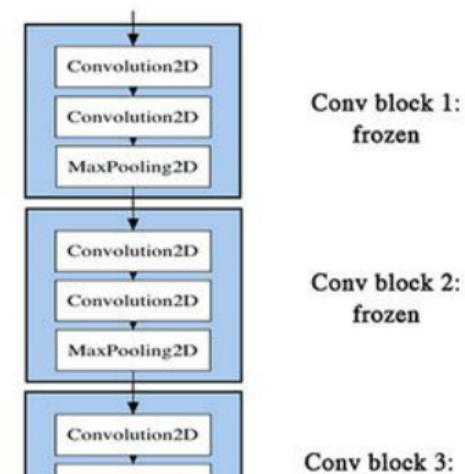
Transfer

Transfer Learning or Domain Adaptation:

- Very few people train an entire CNN from scratch, because we do not have enough data
- It is common to use from pretrained model trained on very large dataset which contains 1.2 million large images with 1000 categories.
- Using ConvNet as an initialization or feature extractor!
- Two different scenario for using pretrained model:
 - ConvNet as fixed feature extractor: **CNN Codes**
 - Fine-tuning the Convnet:
 - Replace the last fully-connected layers with new one with random weight and train their weights again.
 - If you have more data, you can retrain more layers with “backprop”

Transfer Learning or Domain Adaptation:

- Very few people train an entire CNN from scratch, because we do not have enough data
- It is common to use from pretrained model trained on very large dataset which contains 1.2 million large images with 1000 categories.
- Using ConvNet as an initialization or feature extractor!
- Two different scenario for using pretrained model:
 - ConvNet as fixed feature extractor: **CNN Codes**
 - Fine-tuning the Convnet:
 - Replace the last fully-connected layers with new one with random weight and train their weights again.
 - If you have more data, you can retrain more layers with “backprop”



er Learning

Transfer Learning or Domain Adaptation:

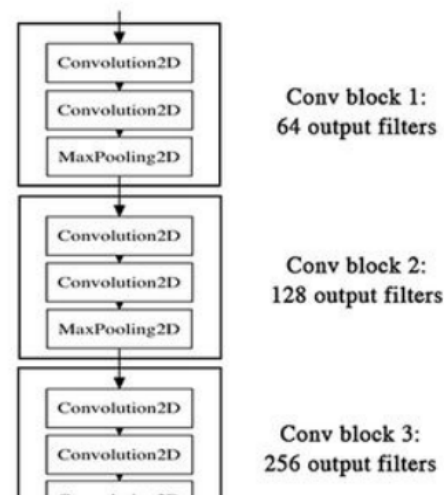
When and how to fine tune:

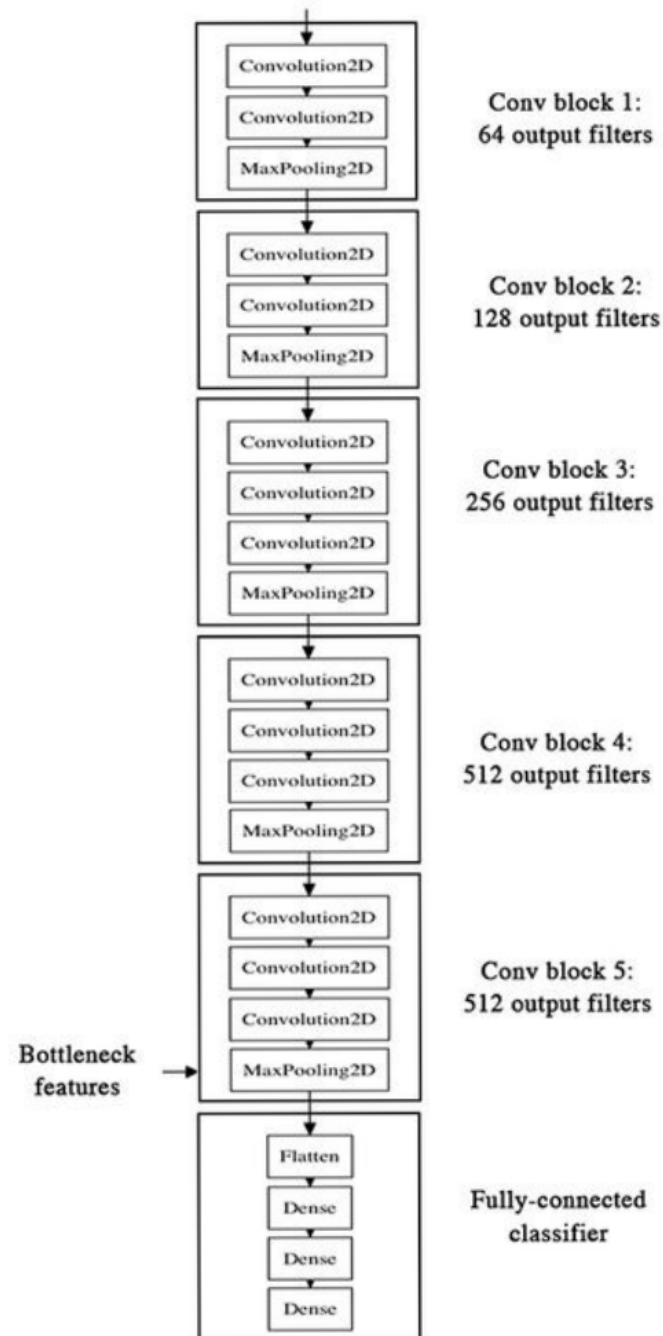
1. New dataset is small and is similar to original dataset:
 - No need to retrain high level features in CNN – If you do, you overfit
 - Just need to retrain the classifier part of model
2. New dataset is large and is similar to original dataset:
 - Retrain all the weights in the network
3. New dataset is small and are totally different from original dataset:
 - Train the classifier not on top of the network but on top of earlier layers
4. New dataset is large and are totally different
 - Just fine tune all the weights but using the pretrained model as an initialization

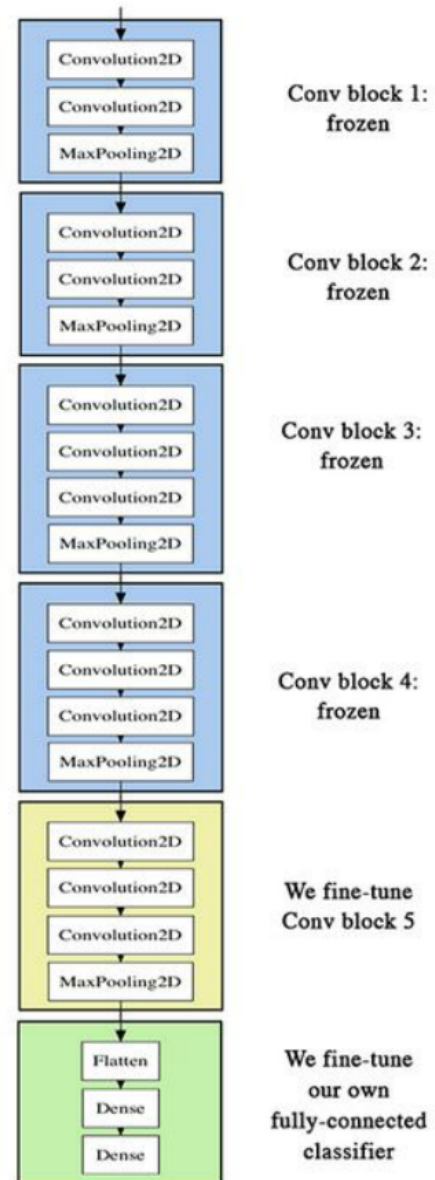
Transfer Learning or Domain Adaptation.

When and how to fine tune:

1. New dataset is small and is similar to original dataset:
 - No need to retrain high level features in CNN – If you do, you overfit
 - Just need to retrain the classifier part of model
2. New dataset is large and is similar to original dataset:
 - Retrain all the weights in the network
3. New dataset is small and are totally different from original dataset:
 - Train the classifier not on top of the network but on top of earlier layers
4. New dataset is large and are totally different:
 - Just fine tune all the weights but using the pretrained model as an initialization



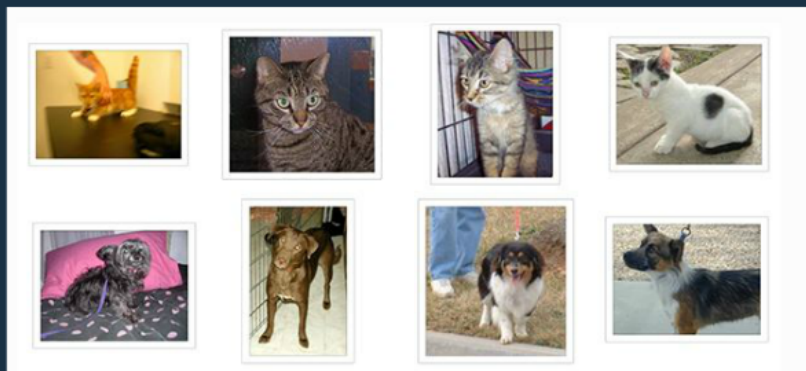




Bottom
feature

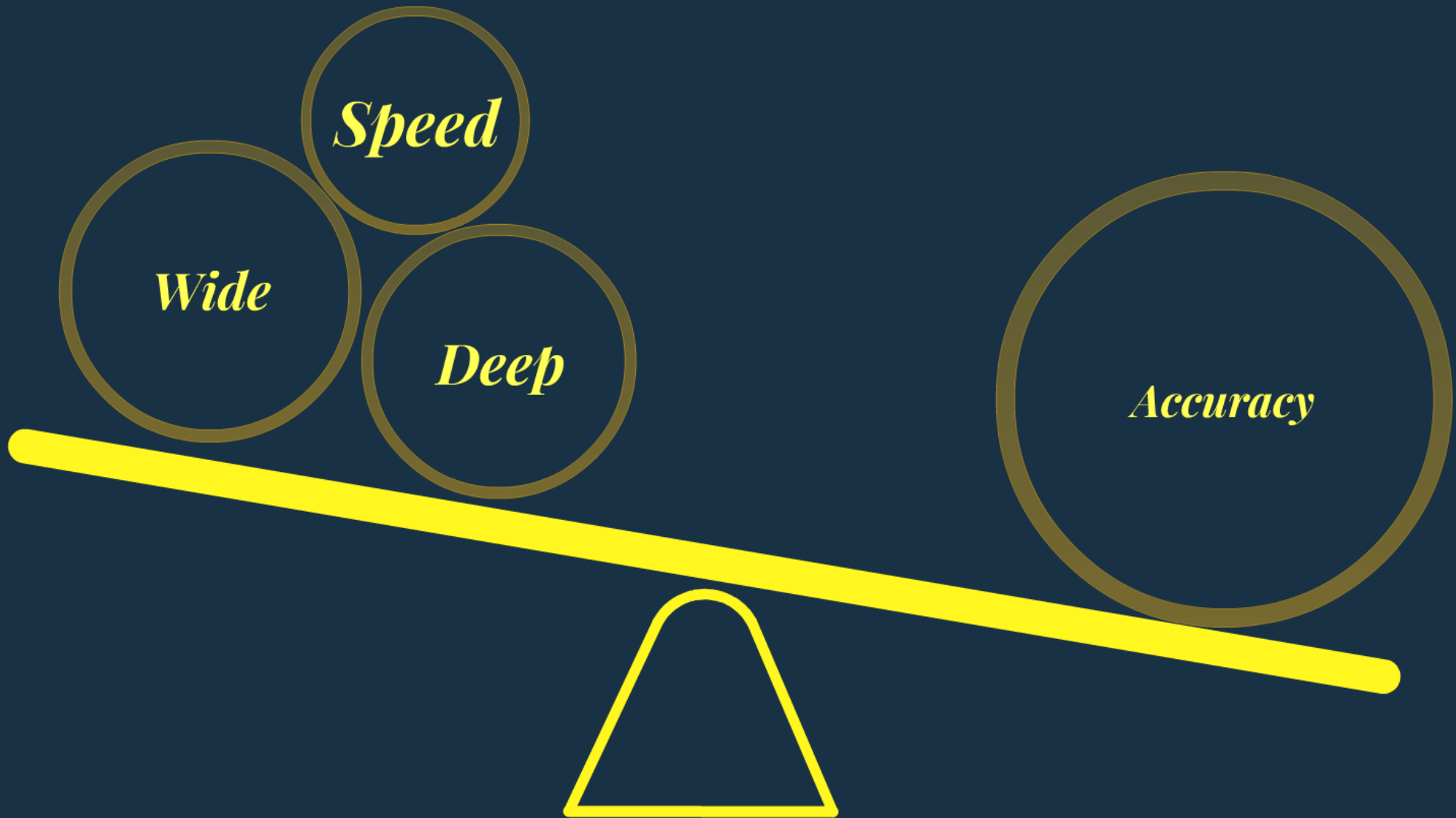
some useful links:

- <https://mitpress.mit.edu/books/deep-learning>
- <https://github.com/Alireza-Akhavan/TensorFlow-Examples>
- <http://bigdataworkgroup.ir/>
- <https://flyyufelix.github.io/2016/10/03/fine-tuning-in-keras-part1.html>
- <https://github.com/Alireza-Akhavan/class.vision>
-
-
- Telegram:
-
- @cvision
-
- @deeplearning.ir
-
- @http://qa.deeplearning.ir
-
- @irandeeplearning

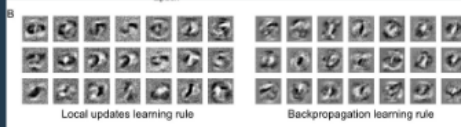
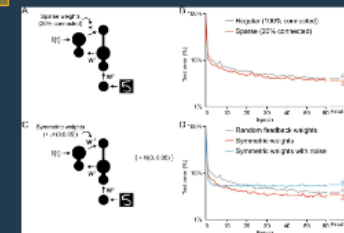
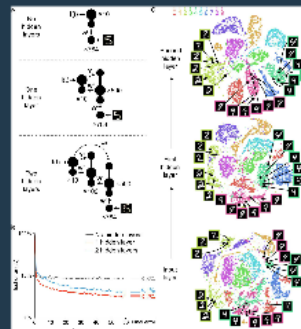
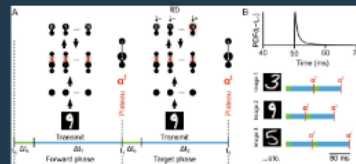
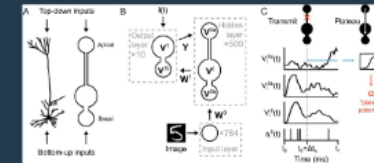
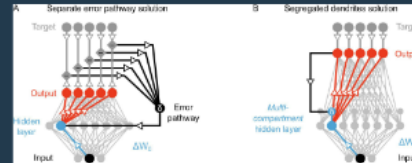
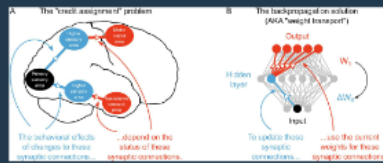


Kaggle.com

Evaluation

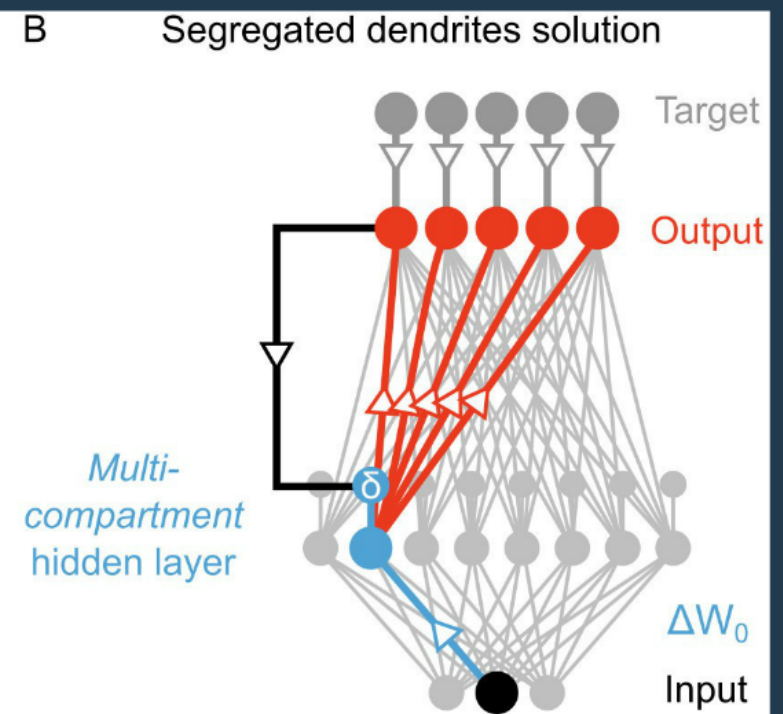
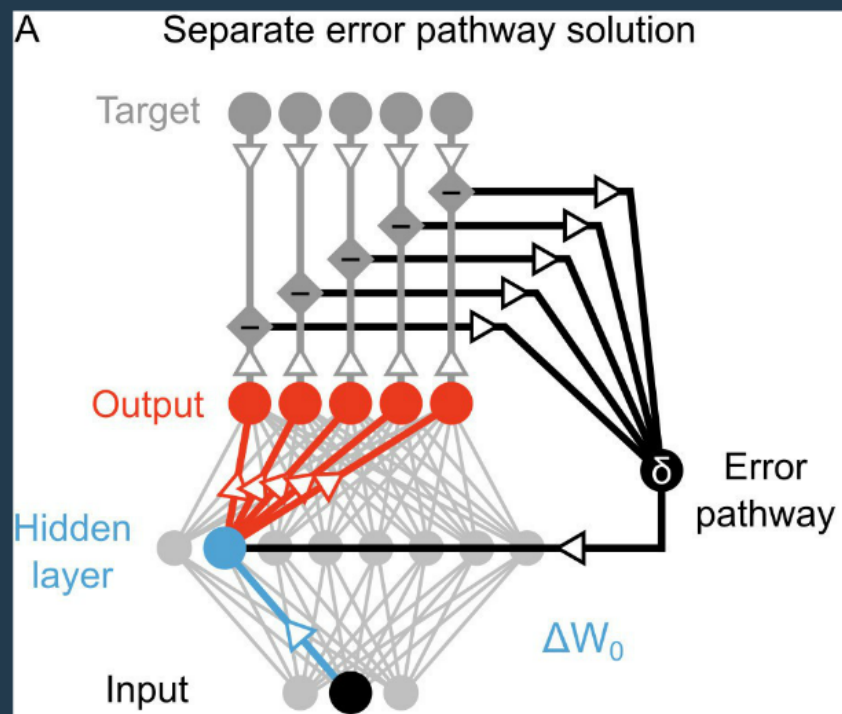


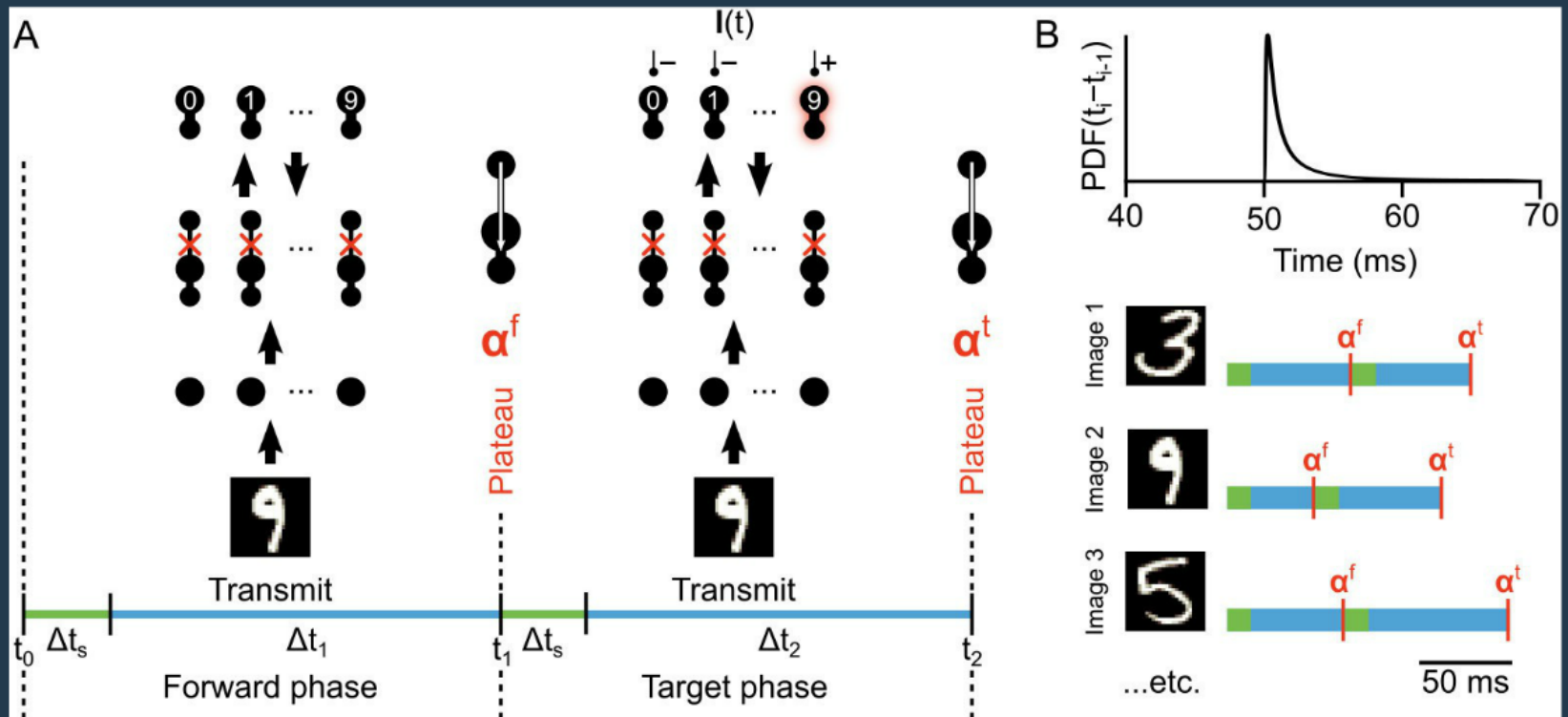
Do Our Brains Use Deep Learning to Make Sense of the World?

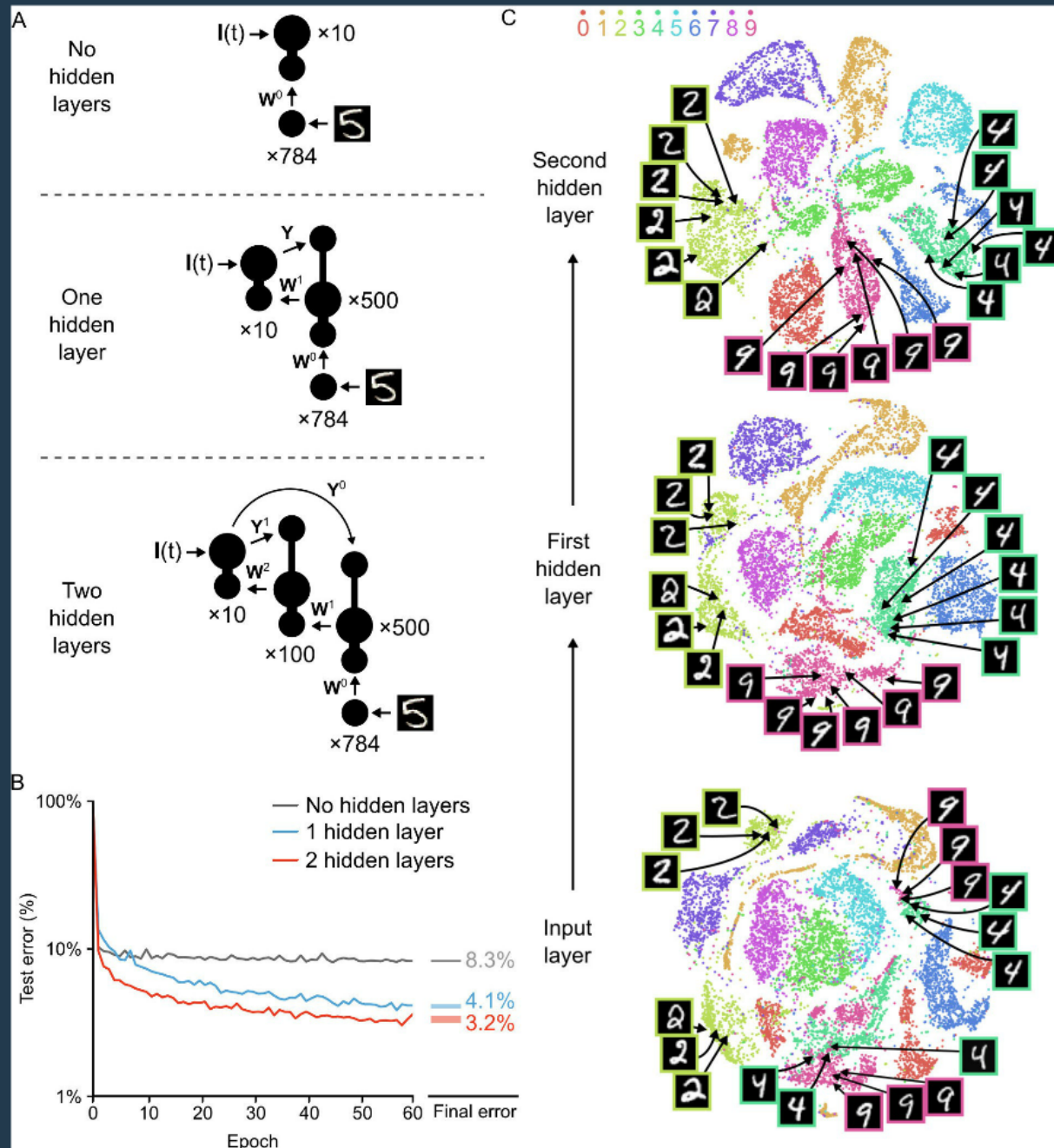


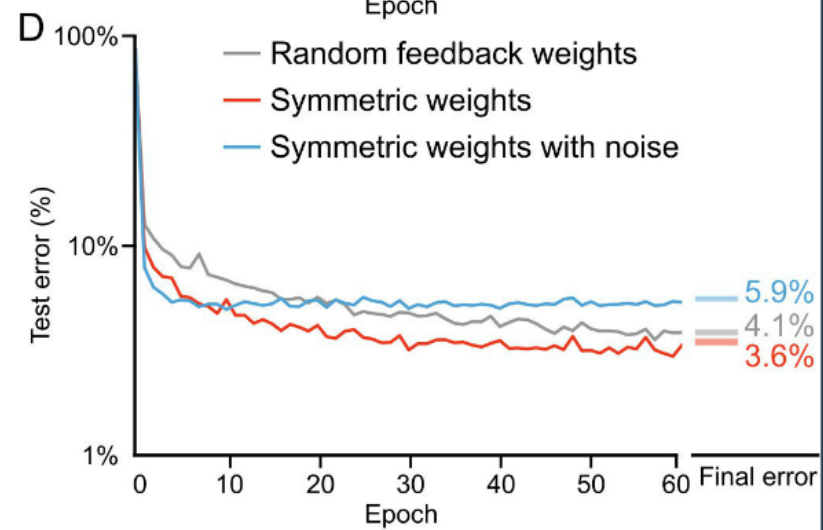
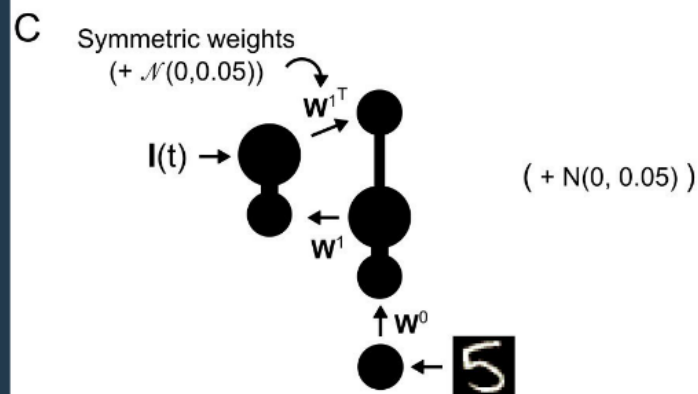
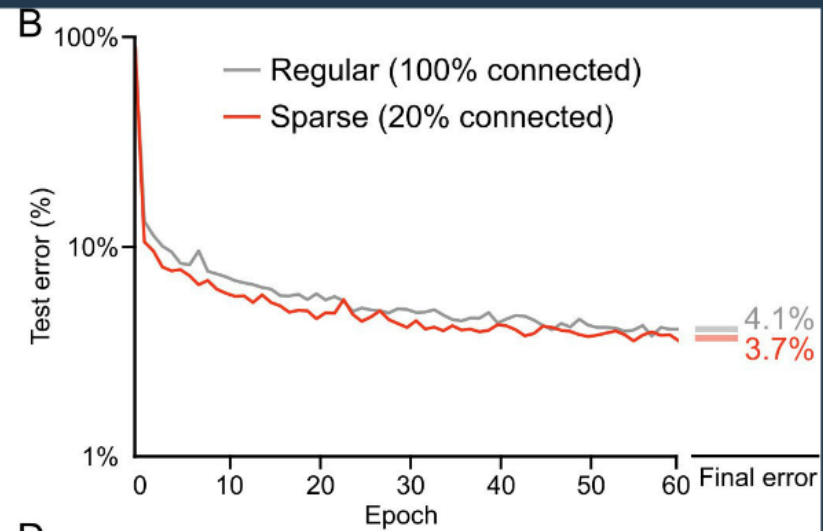
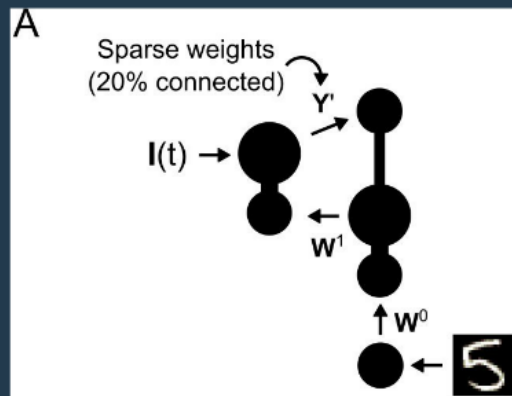
<https://elifesciences.org/articles/22901/figures>

Do Our Brains Use Deep Learning to Make Sense of the World?

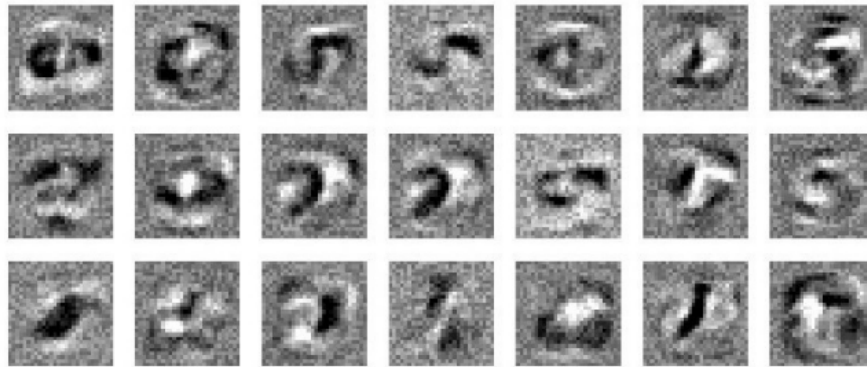




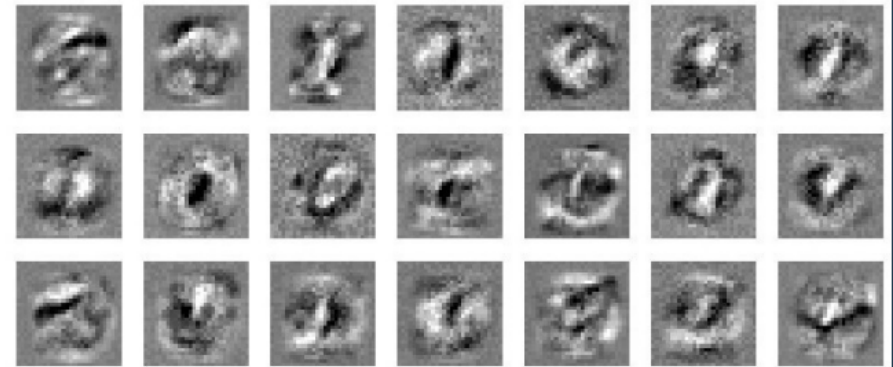




B



Local updates learning rule



Backpropagation learning rule

Deep Learning Frameworks

TensorFlow

Caffe

Keras

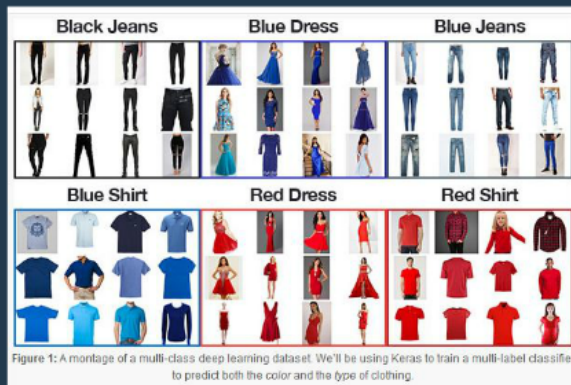
Theano

Lasagne

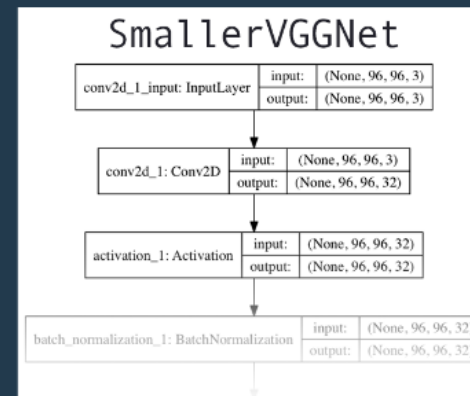


Practice!

- **Transfer Learning (Code)**
- **Classify new category (Code)**
- **multi label classification**
- **Speech Recognition**



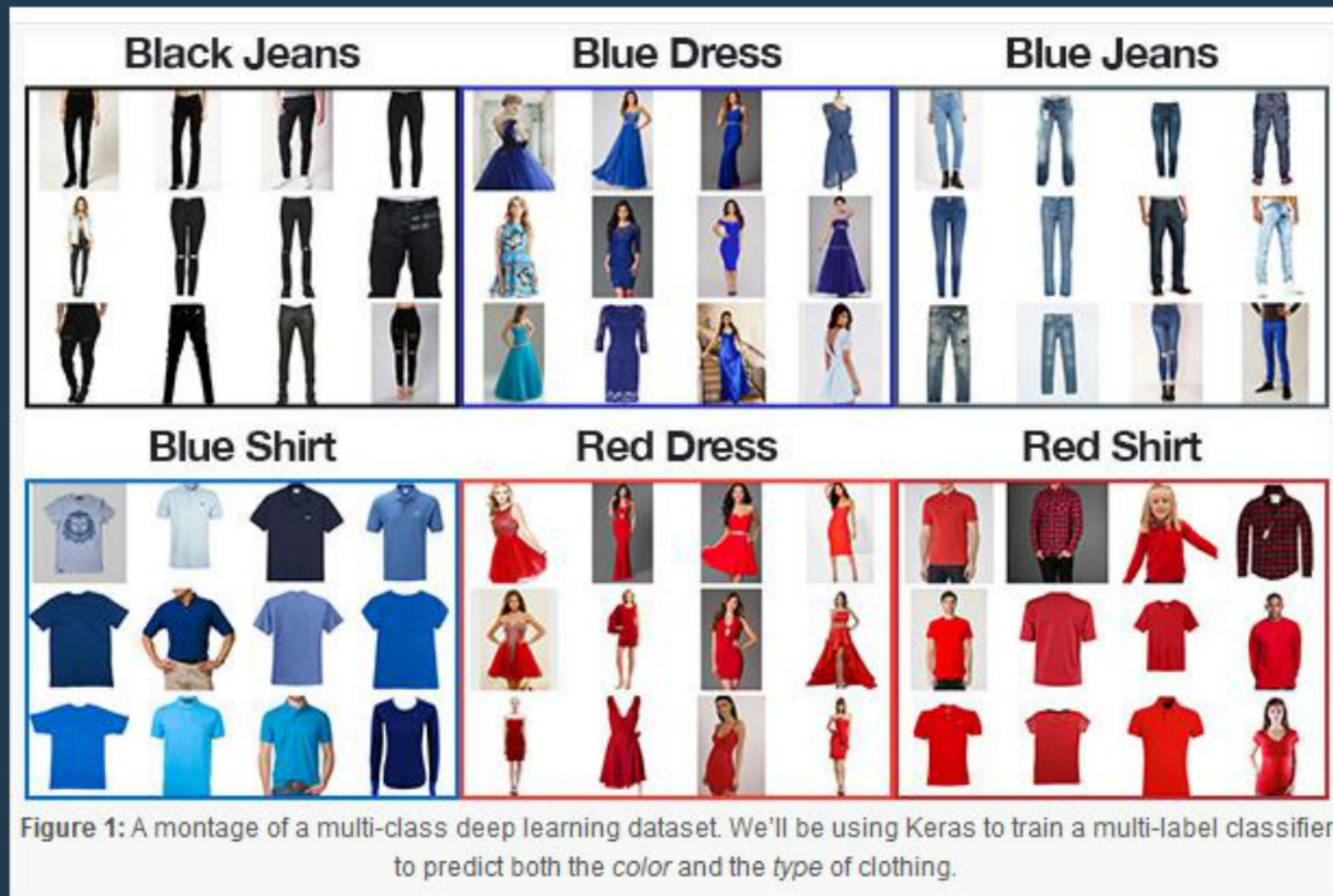
pyimagesearch.com



pyimagesearch.com



pyimagesearch.com



pyimagesearch.com

pyimagesearch.com



classifier

SmallerVGGNet

conv2d_1_input: InputLayer	input:	(None, 96, 96, 3)
	output:	(None, 96, 96, 3)



conv2d_1: Conv2D	input:	(None, 96, 96, 3)
	output:	(None, 96, 96, 32)



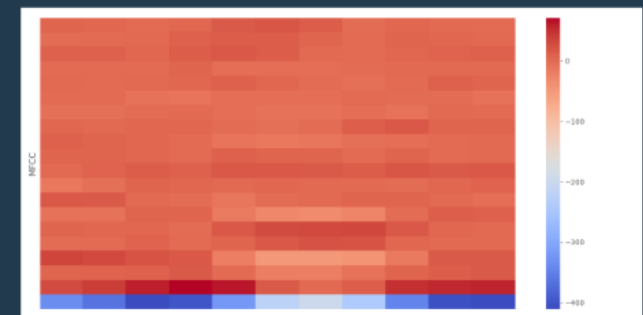
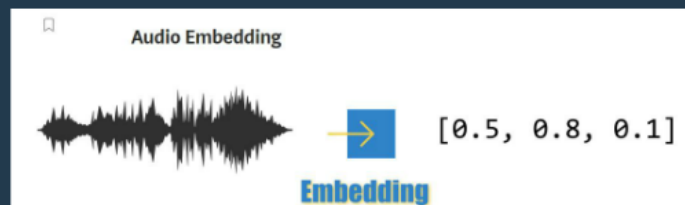
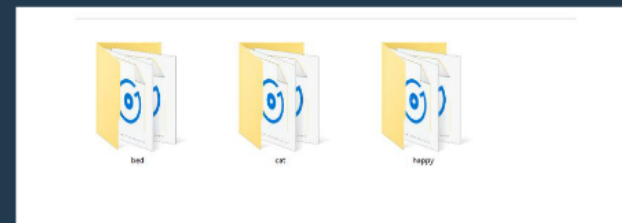
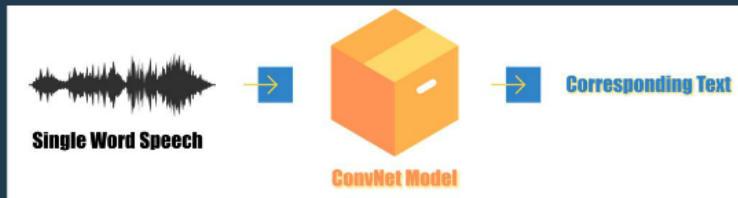
activation_1: Activation	input:	(None, 96, 96, 32)
	output:	(None, 96, 96, 32)



batch_normalization_1: BatchNormalization	input:	(None, 96, 96, 32)
	output:	(None, 96, 96, 32)



pyimagesearch.com



<https://blog.manash.me/building-a-dead-simple-word-recognition-engine-using-convnet-in-keras-25e72c19c12b>



Single Word Speech



ConvNet Model



Corresponding Text



Audio Embedding



bed



cat



happy

Single Word Speech

ConvNet Model

Audio Embedding



$[0.5, 0.8, 0.1]$

Embedding



cat

happy



Deep Learning

What?

Why?

How?

Thank You for attention :)

Any Question?

Thank You for attention :)

Any Question?



