



CHULA ENGINEERING COMPUTER
Foundation toward Innovation



Introduction to Deep Learning

Asst. Prof. Peerapon Vateekul, Ph.D.

Department of Computer Engineering,
Faculty of Engineering, Chulalongkorn University
Peerapon.v@chula.ac.th

www.cp.eng.chula.ac.th/~peerapon/



Outline

- Introduction to Deep Learning
- Convolutional Neural Networks (CNN) [Imaging Task]
- Recurrent Neural Networks (RNN) [Time Series Forecasting]
- Language Technologies



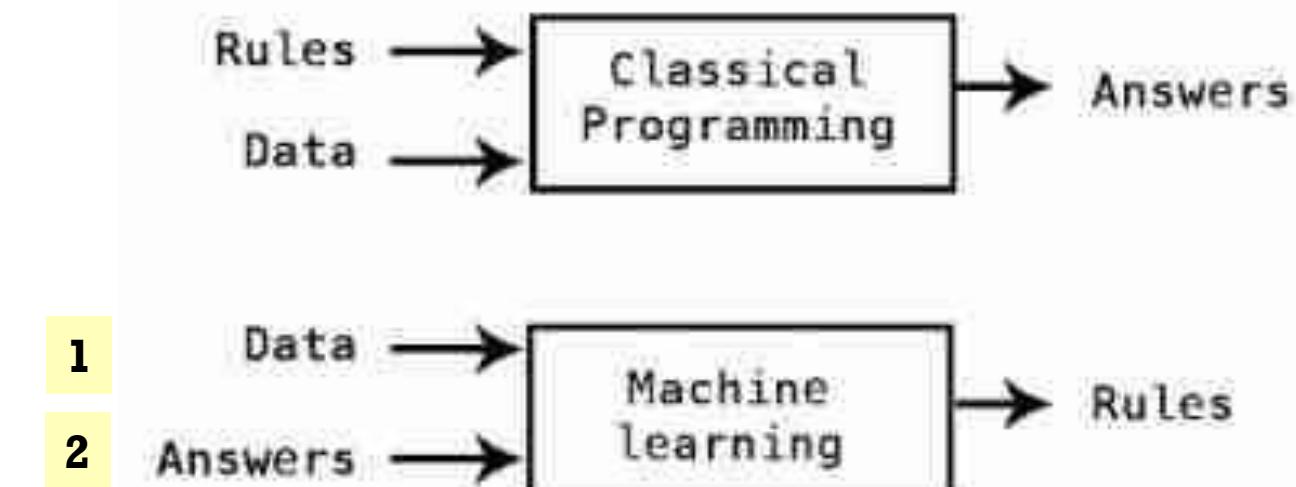
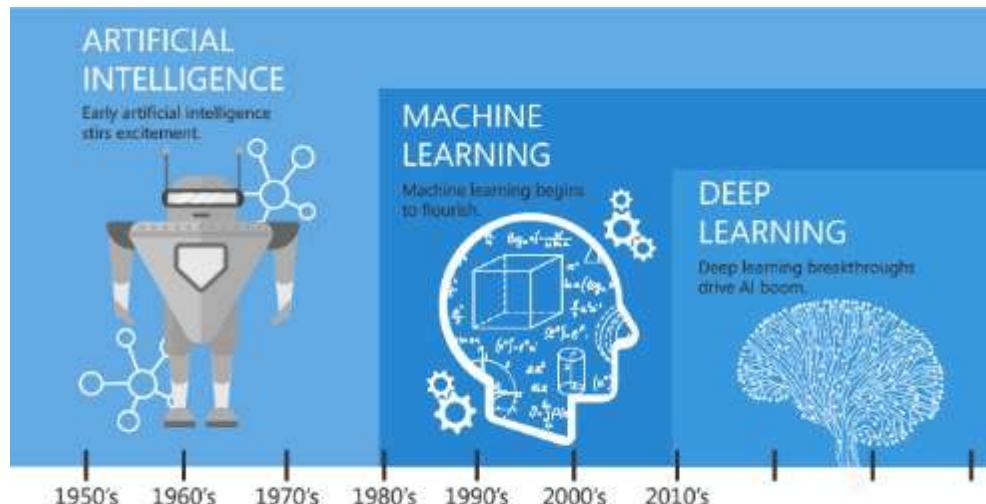


Introduction to Deep Learning



AI = Automation

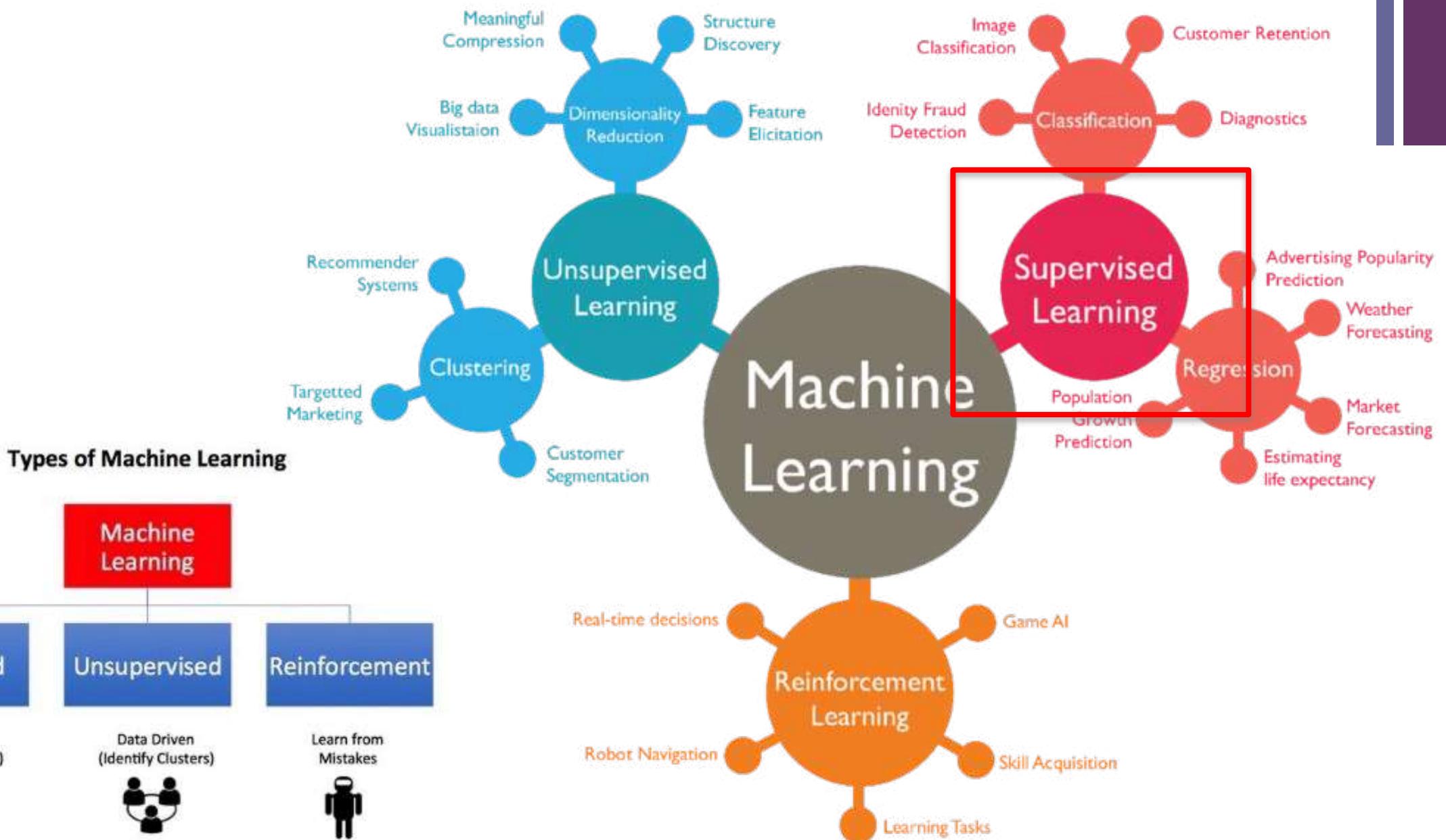
- 1) Rule-based AI (Symbolic AI)
- 2) Machine Learning



<https://mc.ai/machine-learning-basics-artificial-intelligence-machine-learning-and-deep-learning/>

+ Machine Learning

5





Task1: Supervised learning

Handcrafted features

Training Data



inputs					target
Age	Gender	BodyTemp	Cough	Corona	
12	Female	37	Yes	Yes	
35	Female	39	No	Yes	
32	Male	38	Yes	No	

Testing Data

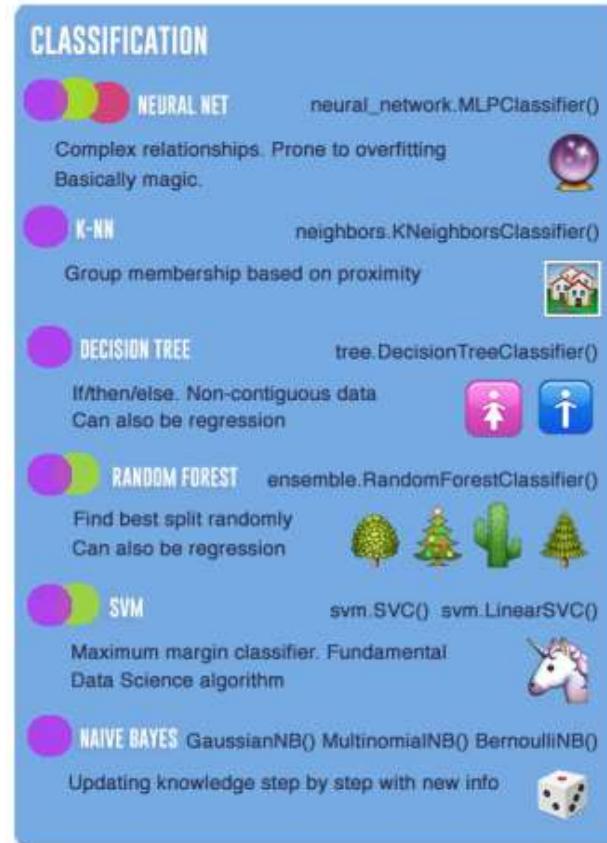
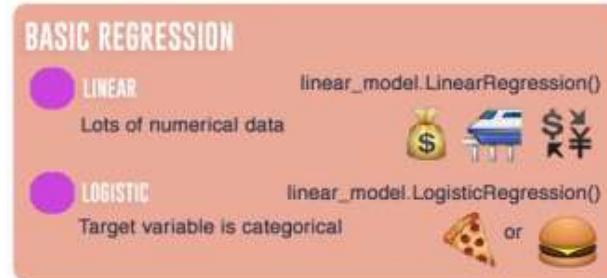


Age	Gender	BodyTemp	Cough	Corona
25	Male	40	No	?

Application: Corona Prediction

+ Prediction algorithms

- Decision Tree
- (Logistic) Regression
- kNN
- Support Vector Machine
- Neural Networks (NN)
- Deep Learning



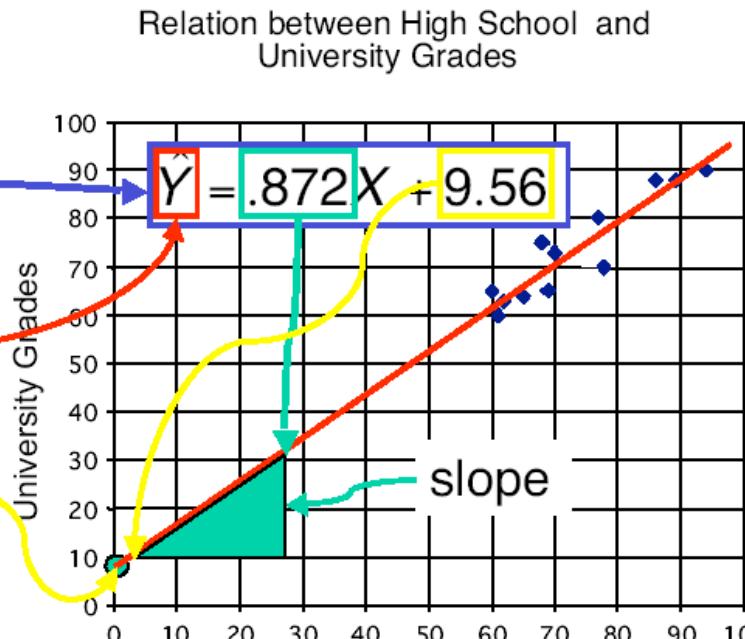
+

Regression – Linear Relationship

regression
equation

predicted
value of Y

y-intercept



weight, coefficient

$$\hat{y} = \hat{w}_0 + \hat{w}_1 x_1 + \hat{w}_2 x_2$$

target	intercept	input
--------	-----------	-------

- The least square method aims to minimize the following term

$$\sum_{\substack{training \\ data}} (y_i - \hat{y}_i)^2$$



Logistic Regression (cont.)

Linear Relationship

Training Data



inputs				target
Age	Gender	BodyTemp	Cough	Corona
12	Female (0)	37	Yes (1)	Yes
35	Female (0)	39	No (0)	Yes
32	Male (1)	38	Yes (1)	No

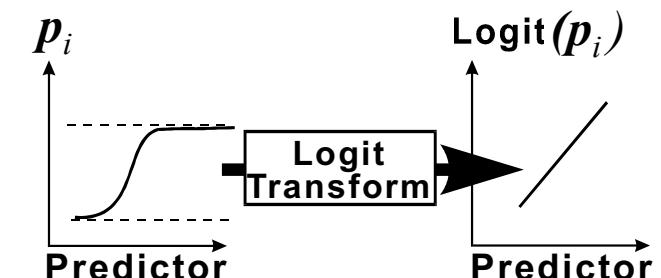
$$\text{Logit_score} = w_0 + w_1 * \text{Age} + w_2 * \text{Gender} + w_3 * \text{Temp} + w_4 * \text{Cough}$$

$$\text{Logit_score} = 0.01 - 0.3 * \text{Age} + 0.2 * \text{Gender} + 0.2 * \text{Temp} + 0.9 * \text{Cough}$$

Example

$$\text{Logit_score} = 0.01 - 0.3 * 12 + 0.2 * 0 + 0.2 * 37 + 0.9 * 1 = 4.71$$

prob = 0.9911 → "Yes"



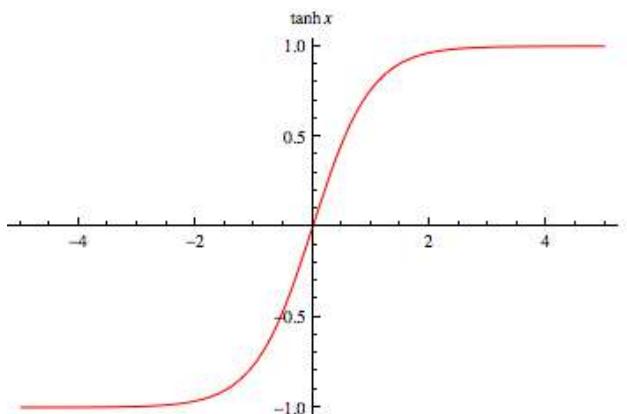
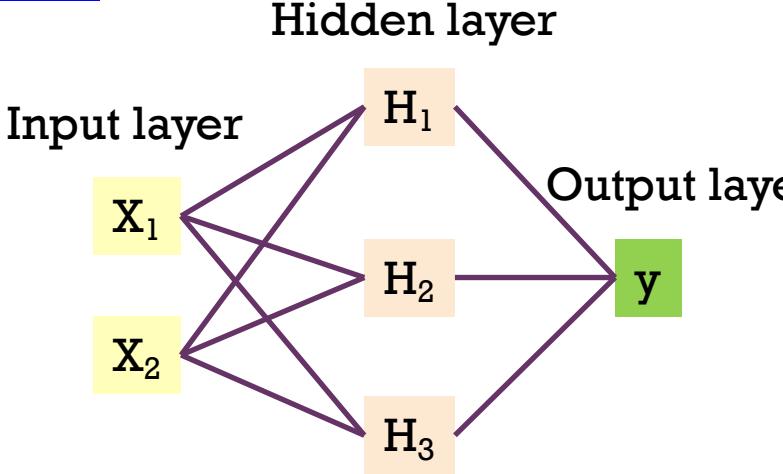
$$\hat{p} = \frac{1}{1 + e^{-\text{logit}(\hat{p})}}$$

Application: Corona Prediction



Neural Networks (universal approximator)

Non-linear relationship



$$\log\left(\frac{\hat{p}}{1-\hat{p}}\right) = \hat{w}_0 + \hat{w}_1 H_1 + \hat{w}_2 H_2 + \hat{w}_3 H_3$$

$$H_1 = \tanh(\hat{w}_{10} + \hat{w}_{11} x_1 + \hat{w}_{12} x_2)$$

$$H_2 = \tanh(\hat{w}_{20} + \hat{w}_{21} x_1 + \hat{w}_{22} x_2)$$

$$H_3 = \tanh(\hat{w}_{30} + \hat{w}_{31} x_1 + \hat{w}_{32} x_2)$$

Stop when?

- Converge (no change in loss)
- Max epochs

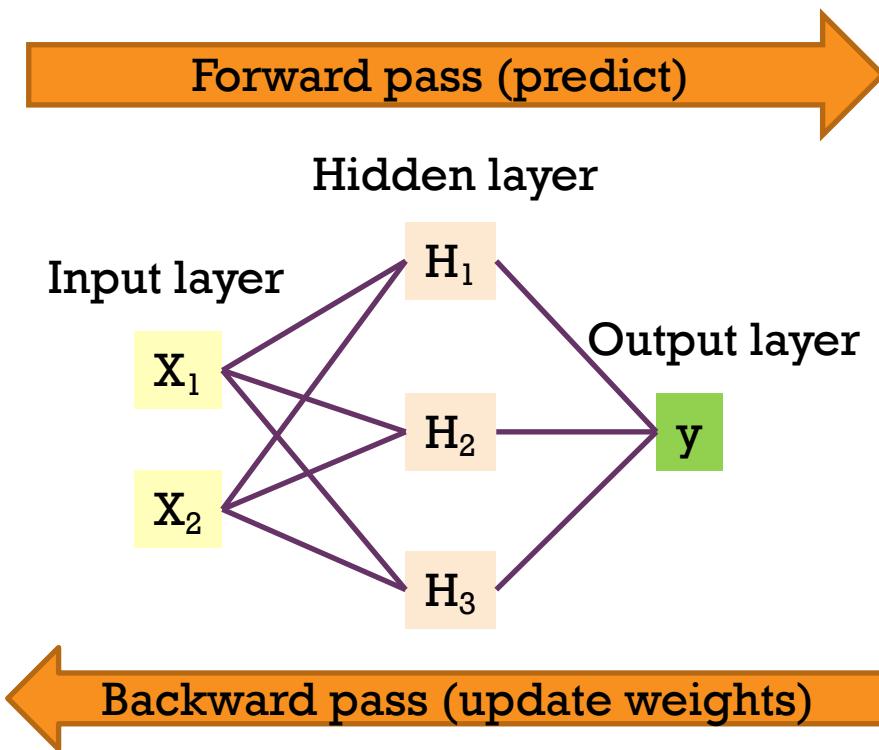
Important Params:

- #hidden units, #hidden layers
- Learning rate, Momentum, decay
- Seed number, etc.

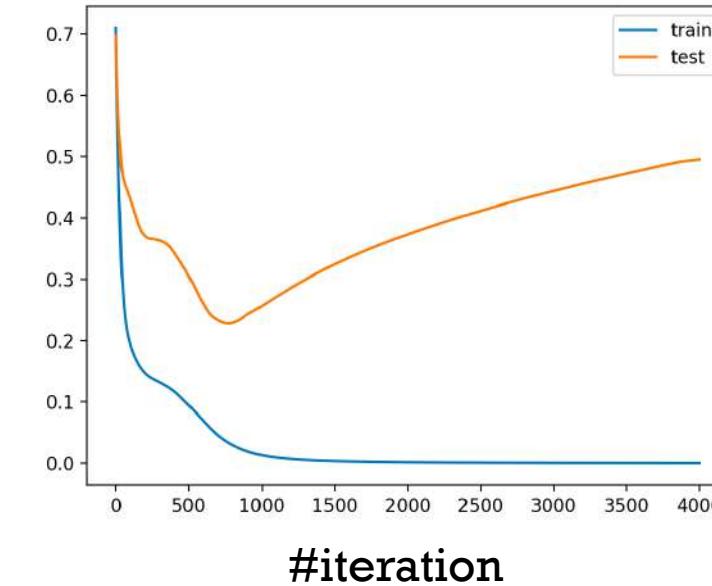


Neural Networks (cont.): Training

Non-linear relationship



Age	Income	Gender	Province	Corona
25	25,000	Female	Bangkok	Yes
35	50,000	Female	Nontaburi	Yes
32	35,000	Male	Bangkok	No



$$\log\left(\frac{\hat{p}}{1-\hat{p}}\right) = \hat{w}_0 + \hat{w}_1 H_1 + \hat{w}_2 H_2 + \hat{w}_3 H_3$$

$$H_1 = \tanh(\hat{w}_{10} + \hat{w}_{11} x_1 + \hat{w}_{12} x_2)$$

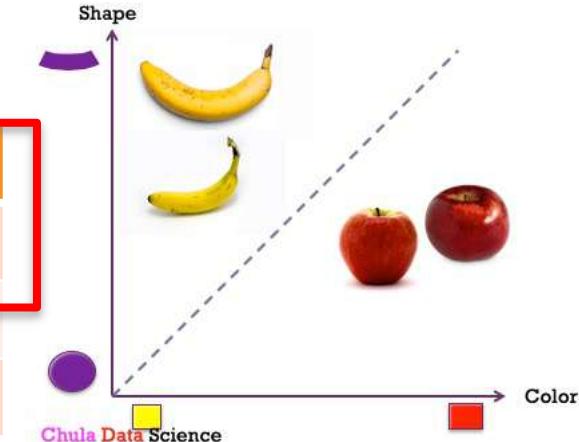
$$H_2 = \tanh(\hat{w}_{20} + \hat{w}_{21} x_1 + \hat{w}_{22} x_2)$$

$$H_3 = \tanh(\hat{w}_{30} + \hat{w}_{31} x_1 + \hat{w}_{32} x_2)$$



Handcrafted features

Age	Income	Gender	Province	Corona
25	25,000	Female	Bangkok	Yes
35	50,000	Female	Nonthaburi	Yes
32	35,000	Male	Bangkok	No



Can we still tell the features (columns)?

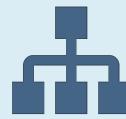




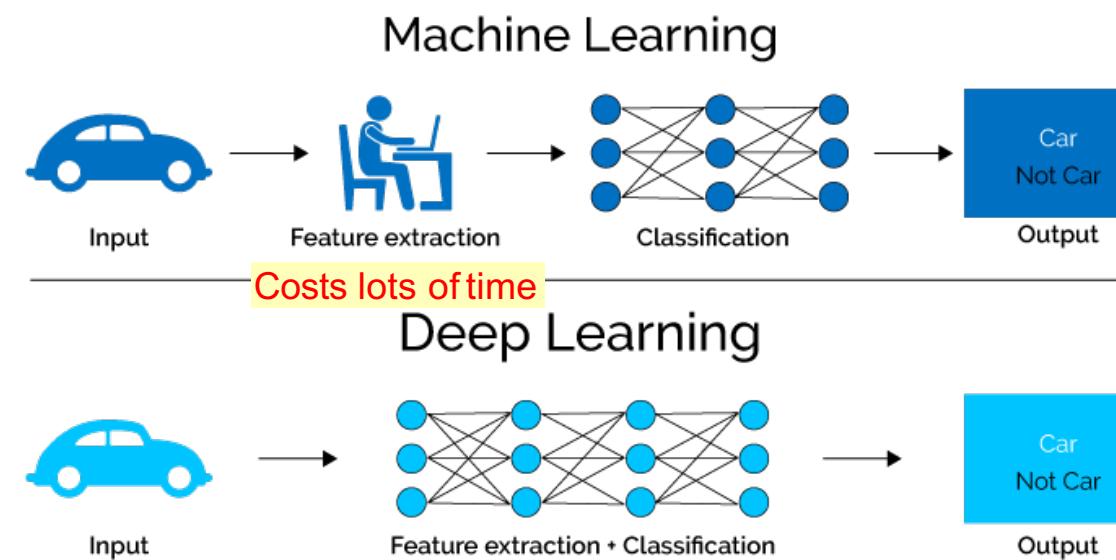
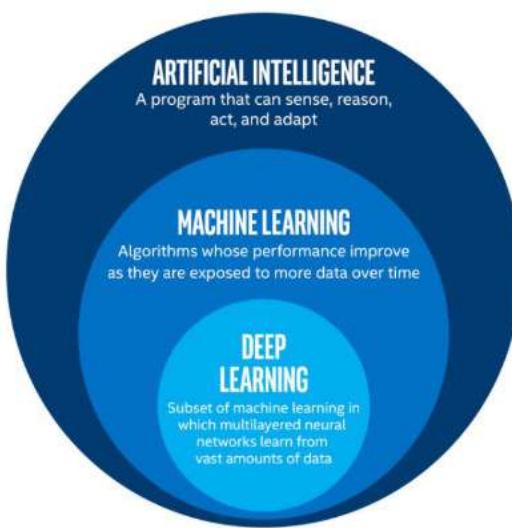
What is Deep Learning (DL)?



Part of the machine learning field of learning representations of data. Exceptional effective at learning patterns.



Utilizes learning algorithms that derive meaning out of data by using a **hierarchy** of multiple layers that **mimic the neural networks of our brain**.





Deep Learning – Basics (cont.)

What did it learn?

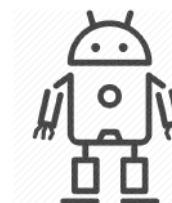
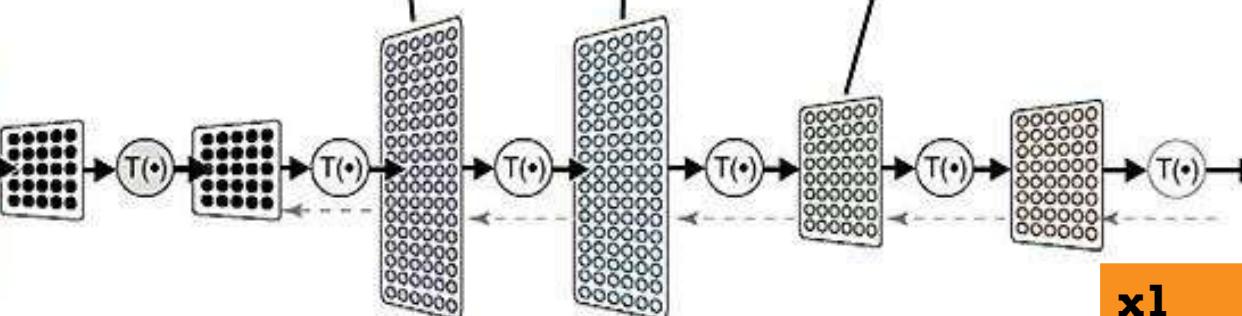
A deep neural network consists of a **hierarchy of layers**, whereby each layer **transforms the input data** into more abstract representations (e.g., edge -> nose -> face). The output layer combines those features to make predictions.



Age	Income	Gender	Province	Corona
25	25,000	Female	Bangkok	Yes



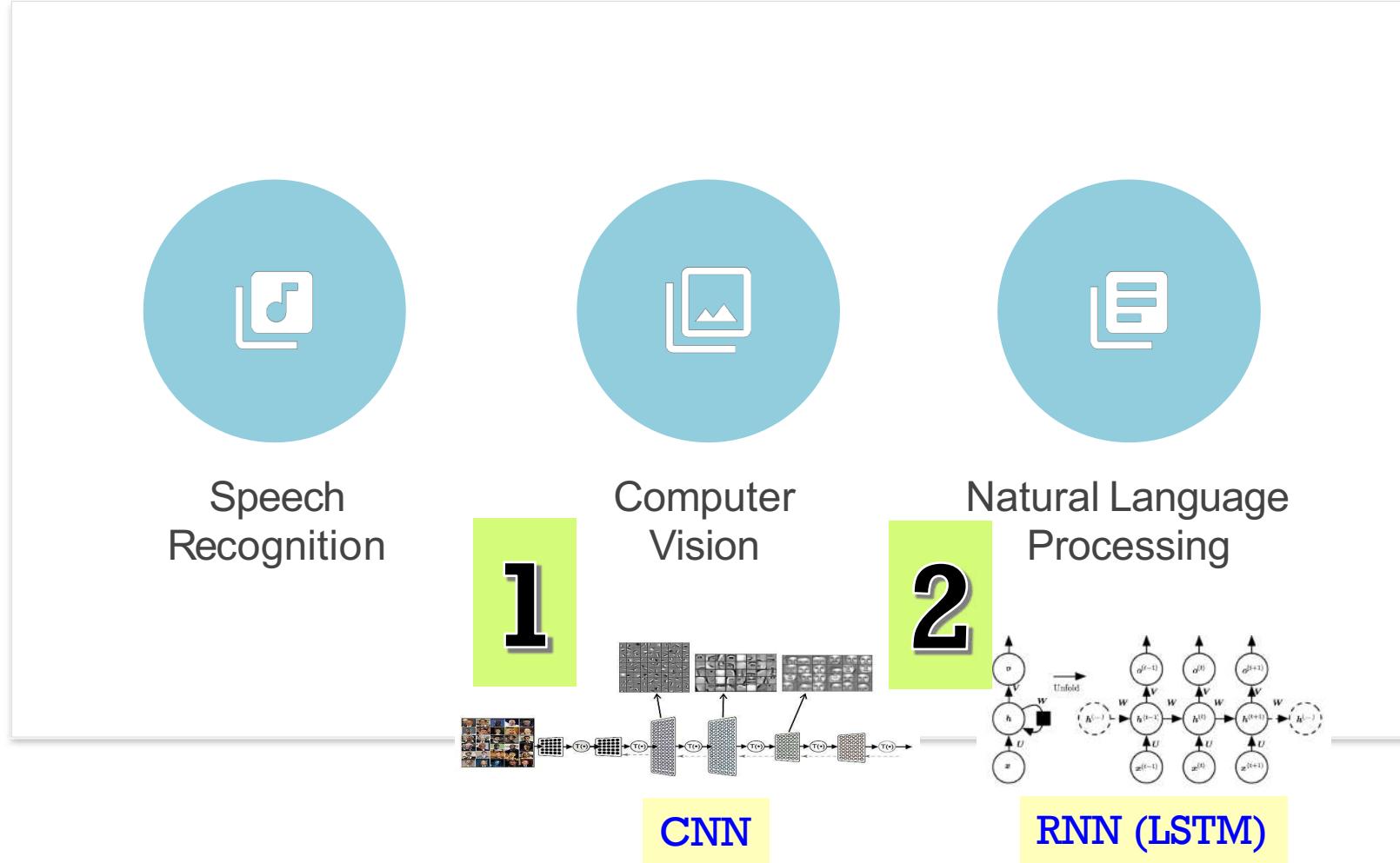
Edges Nose, Eye... Faces



x1	x2	x3	x4	Corona
0.7	0.2	-0.5	-0.1	Yes

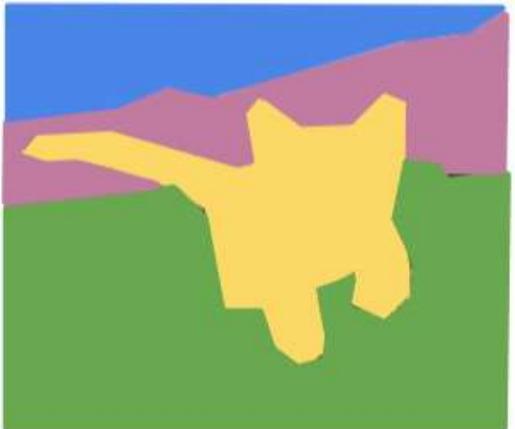


Deep Learning Application



Type of image tasks

Semantic Segmentation



GRASS, CAT,
TREE, SKY

No objects, just pixels

Classification + Localization



CAT

Single Object

Object Detection



DOG, DOG, CAT

Multiple Object

Instance Segmentation



DOG, DOG, CAT

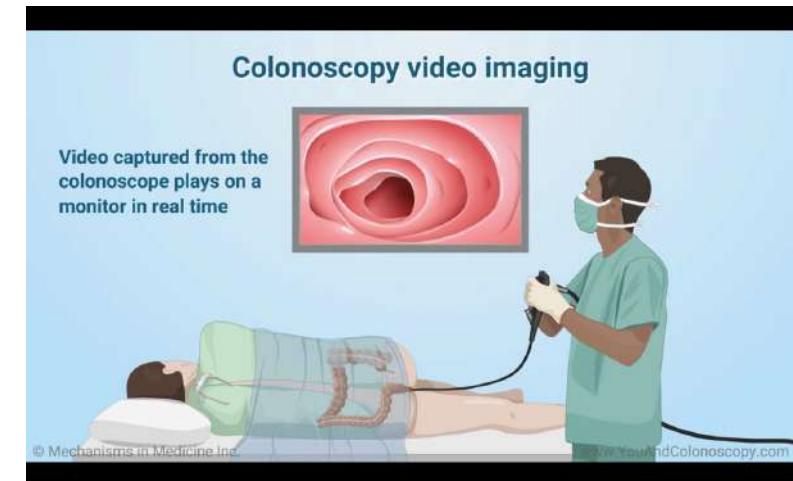
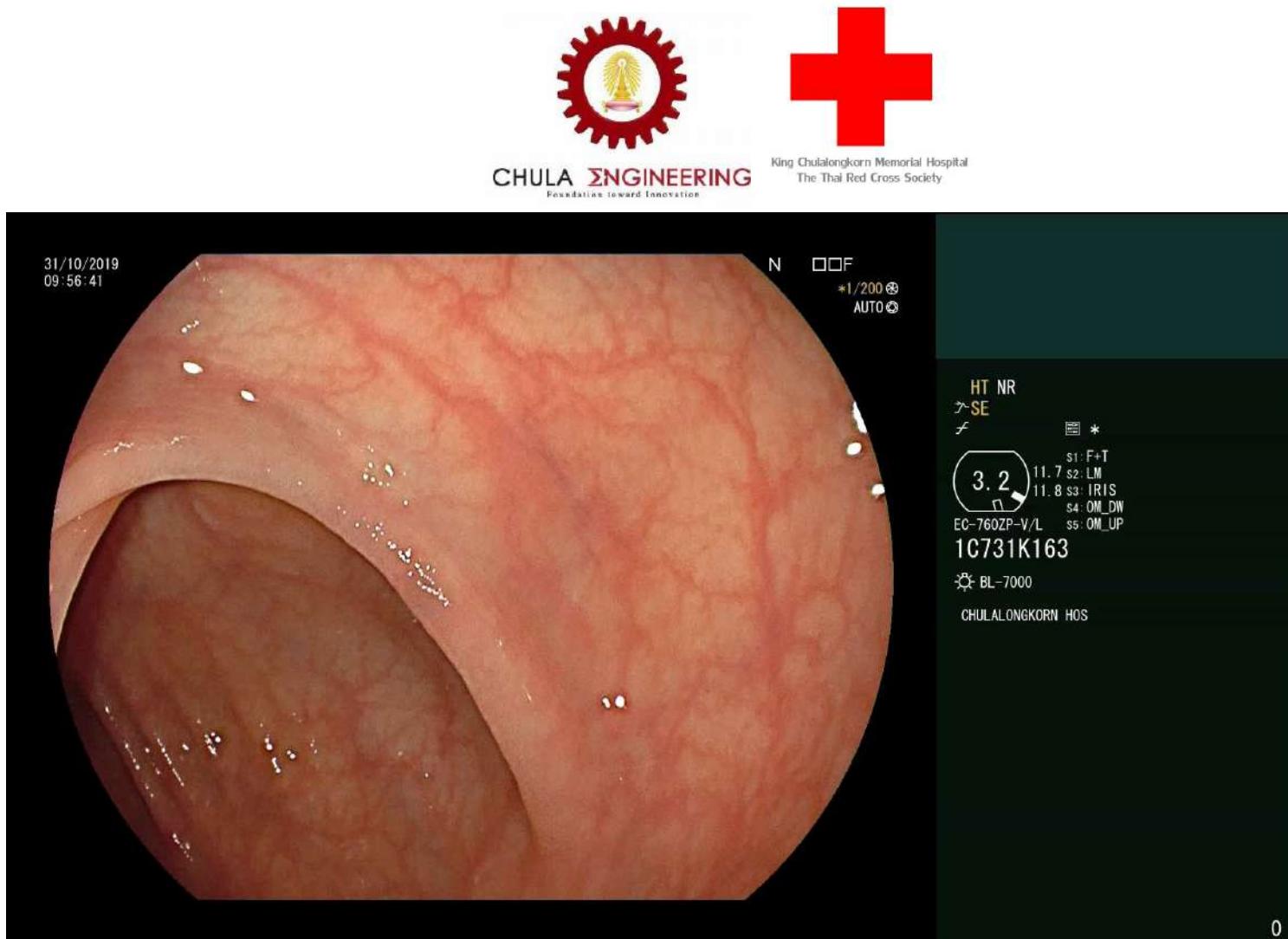
This image is CC0 public domain.

Face Recognition



www.youtube.com/workpointofficial

Smart Medical Diagnosis



+

CNN



Outline

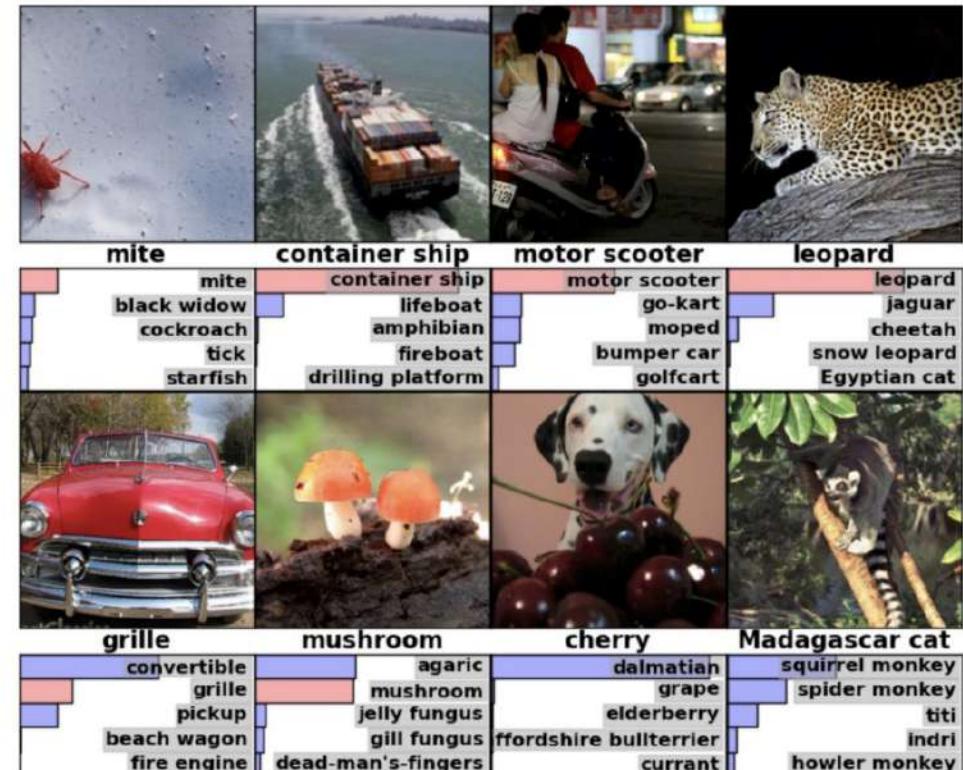
- Introduction
- Basic Image Processing
- Overview of CNN
- Image Processing Tasks with SOTA
- Case Study in Polyp Detection

The **ImageNet** project is a large visual database designed for use in visual object recognition software research. Over **14M** URLs of images have been hand-annotated by ImageNet to indicate what objects are pictured on **22K** categories.



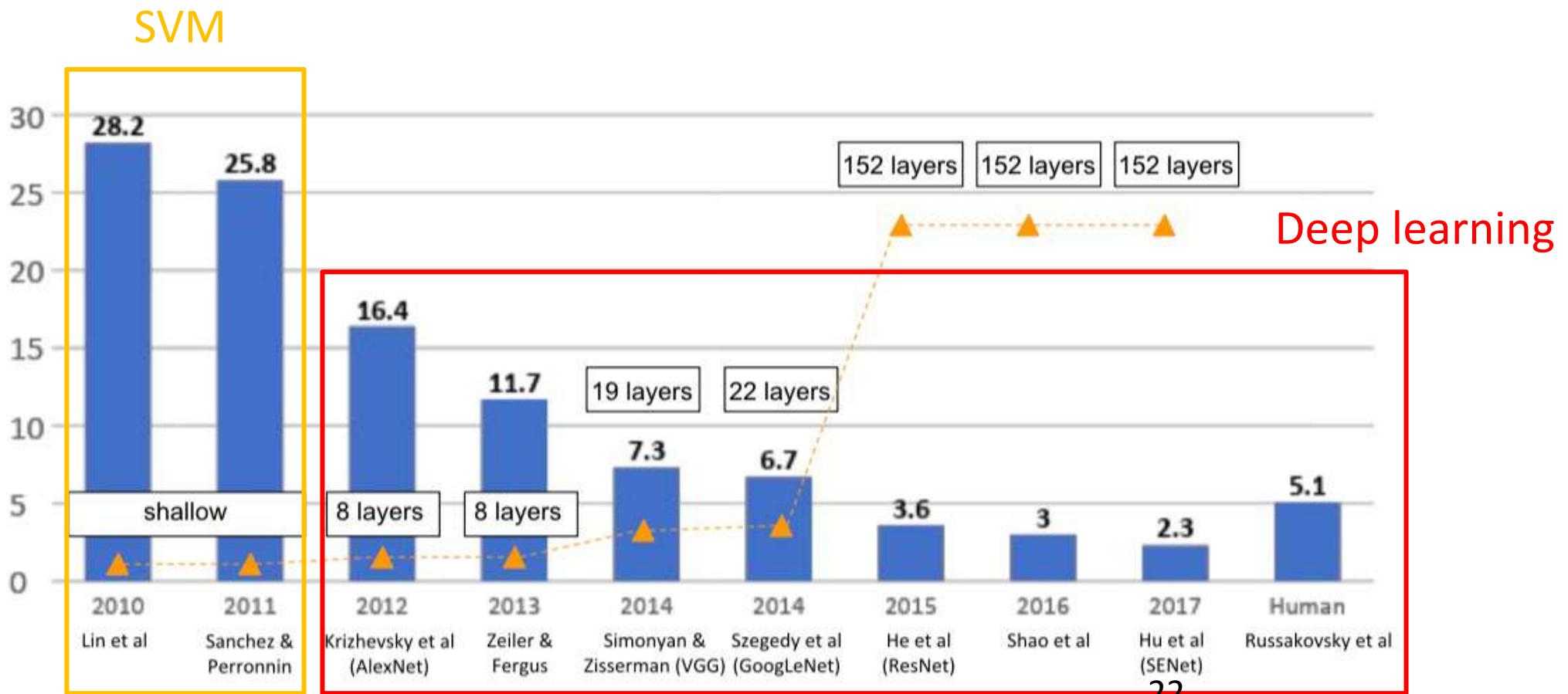
The Image Classification Challenge:
1,000 object classes
1,431,167 images
ImageNet 2017 is the last challenge.

ILSVRC



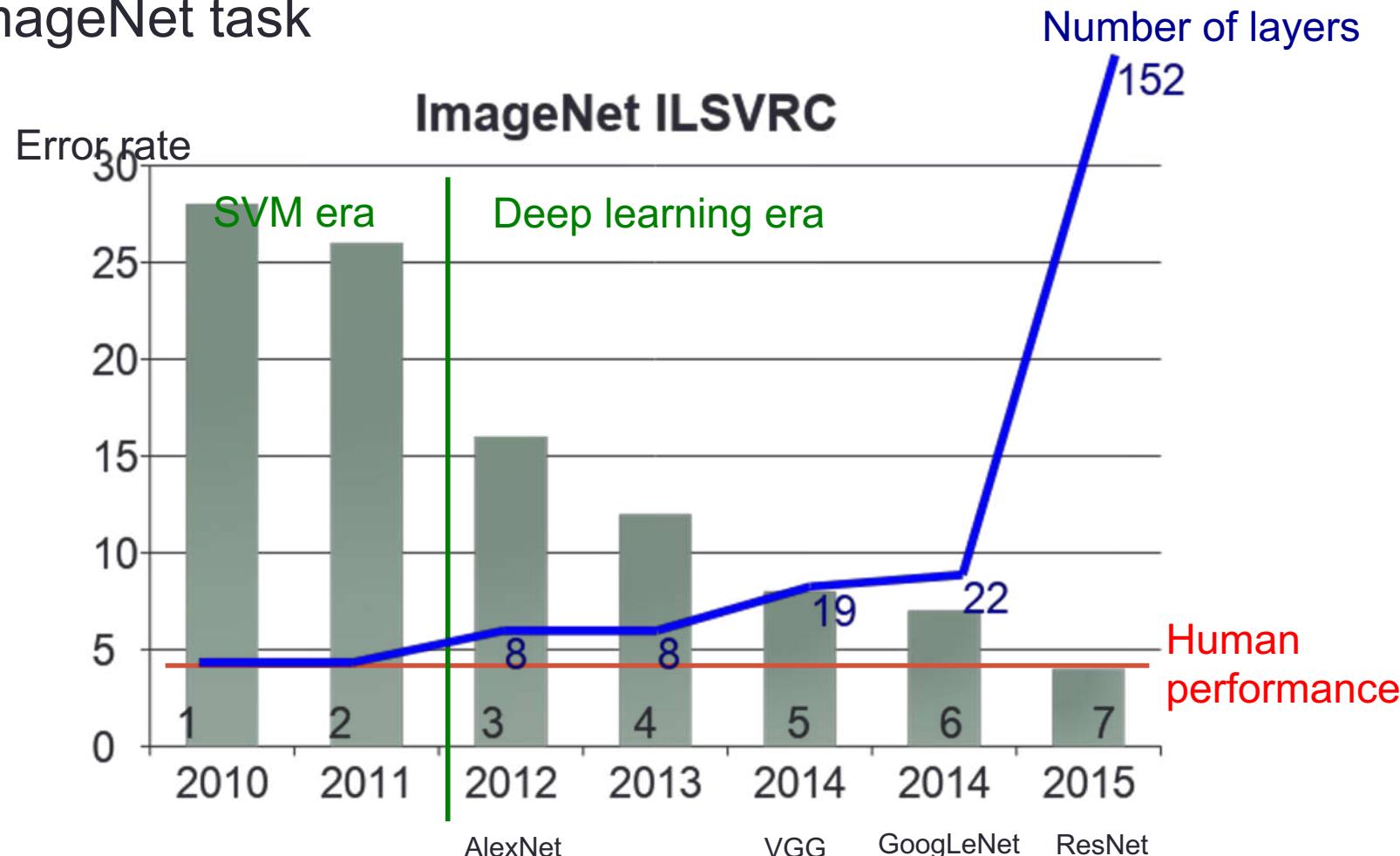
ImageNet Large Scale Visual Recognition Challenge (ILSVRC winners): From SVM to Deep Learning

CNN



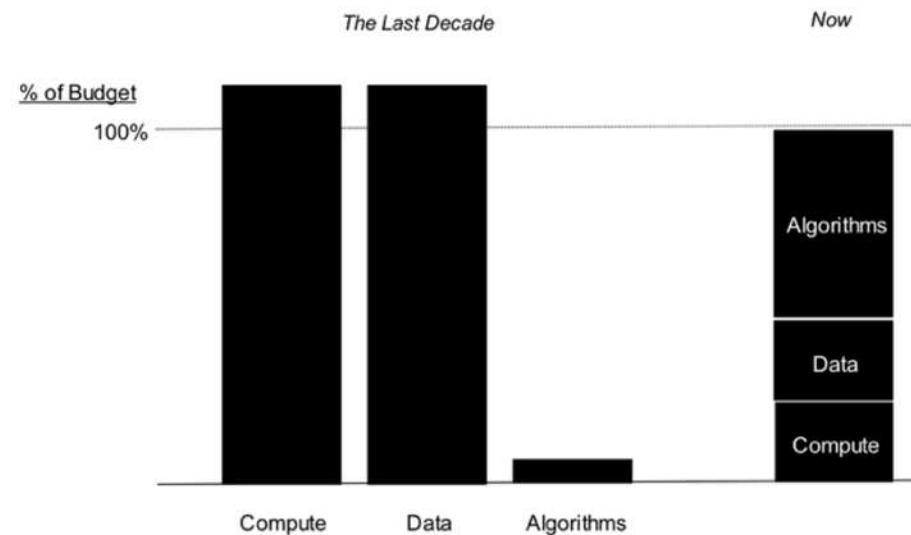
Wider and deeper networks (Beyond Human)

- ImageNet task



Why now

- Neural Networks has been around since 1990s
- **Big data** – DNN can take advantage of large amounts of data better than other models
- **GPU** – Enable training bigger models possible
- **Deep** – Easier to avoid bad local minima when the model is large



Application 1: Basic Image Processing

<https://softwaredevelopmentperestroika.wordpress.com/2014/02/11/image-processing-with-python-numpy-scipy-image-convolution/>

ภาควิชาวิศวกรรมคอมพิวเตอร์
จุฬาลงกรณ์มหาวิทยาลัย

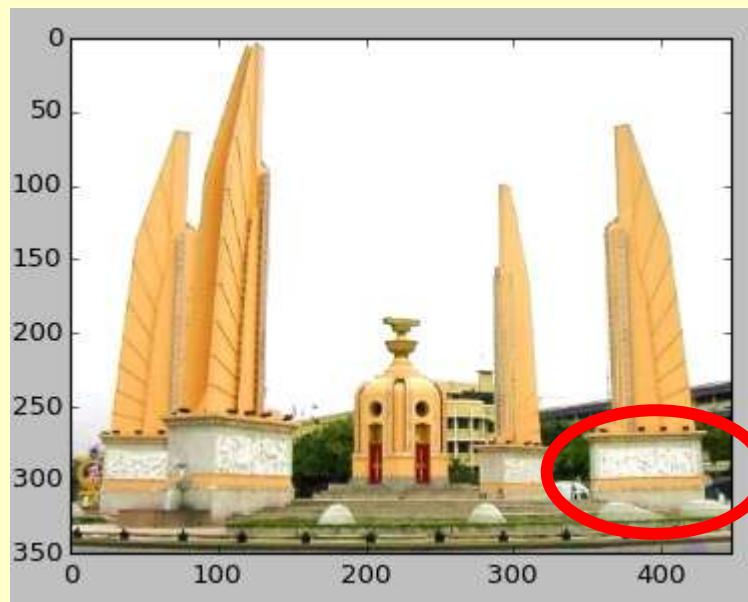
การแสดงรูปภาพใน python ด้วย matplotlib

```
import matplotlib.pyplot as plt  
import matplotlib.image as mpimg
```

เพื่อความง่ายโปรแกรมที่จะเขียน
จากนี้ไปใช้กับไฟล์ png เท่านั้น

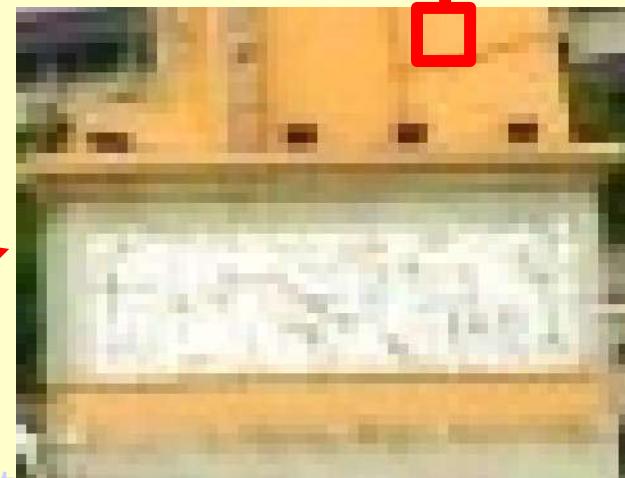
```
image = mpimg.imread('monument.png')  
plt.imshow(image)  
plt.show()
```

image เป็นอาร์เรย์ที่แต่ละช่อง
มีค่า 0 ถึง 1 แทนความเข้มของสี



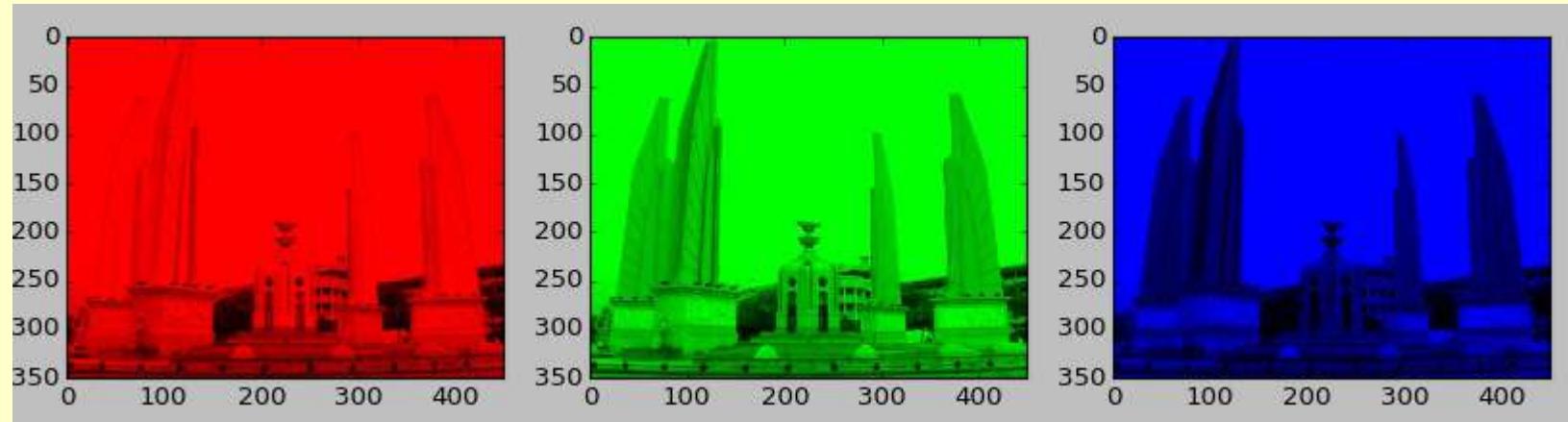
1 pixel = 1 จุดสี มี 3 สีอยู่

0.98762 0.72549 0.35294



1 ภาพ เป็นอาร์เรย์ 3 สามมิติเก็บค่าความเข้มของ R G B

```
>>> image.shape  
(350, 449, 3)
```



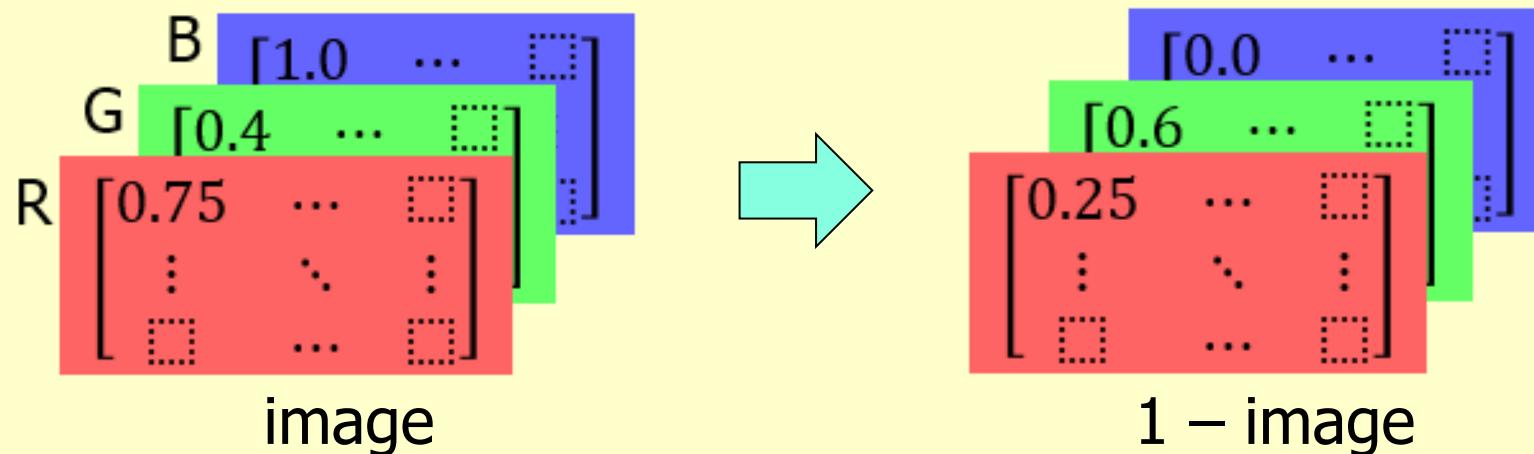
อาร์เรย์ 2 มิติสีแดง
 $\text{image}[:, :, 0]$

อาร์เรย์ 2 มิติสีเขียว
 $\text{image}[:, :, 1]$

อาร์เรย์ 2 มิติสีน้ำเงิน
 $\text{image}[:, :, 2]$

การประมวลผลภาพเบื้องต้น

- `image = mpimg.imread('monument.png')`
- `image.shape → (350, 449, 3)`
- `image = 1 - image` (broadcast & element-wise op.)



ทำแบบนี้แล้ว ภาพเปลี่ยนแปลงไปอย่างไร ?

การประมวลผลภาพเบื้องต้น : ภาพ Negative

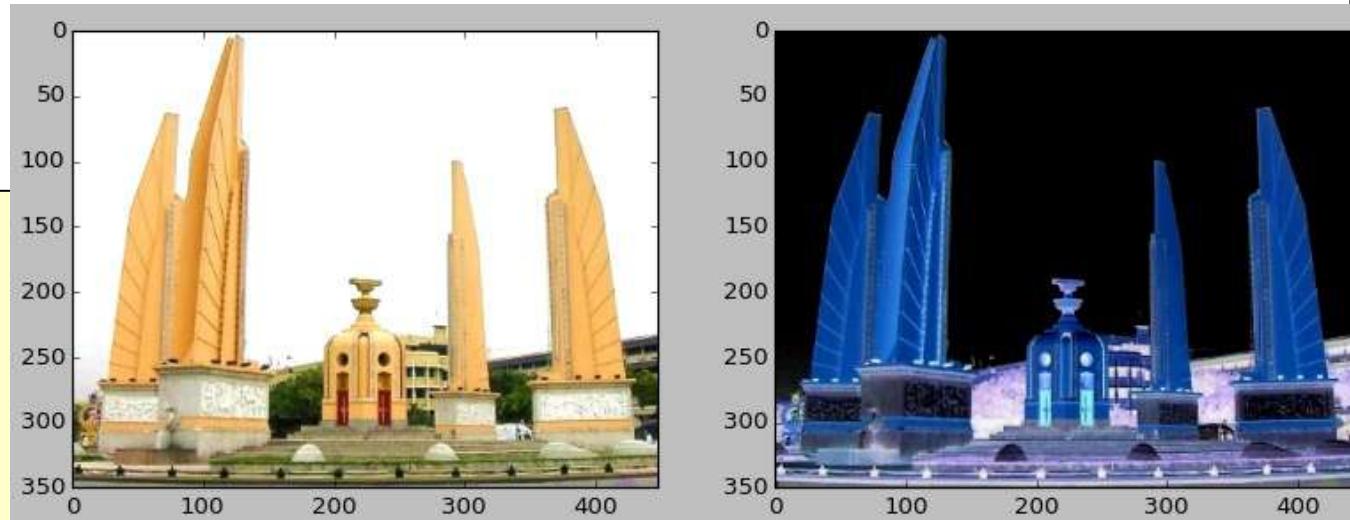
```
import matplotlib.pyplot as plt
import matplotlib.image as mpimg

image = mpimg.imread('monument.png')
plt.subplot(1, 2, 1)
plt.imshow(image)

negative = 1 - image
plt.subplot(1, 2, 2)
plt.imshow(negative)

plt.show()
```

1 – image คือนำค่าทุกช่องไปลบออกจาก 1
ความเข้มสีเปลี่ยนเป็นตรงข้าม
(ขาว → ดำ, เหลือง → น้ำเงิน, ...)
broadcast & element-wise subtract



การประมวลผลภาพเบื้องต้น : ภาพสีเทา

```
import numpy as np  
import matplotlib.pyplot as plt  
import matplotlib.image as mpimg
```

จุดสีที่ R G B มีค่าความเข้มเท่ากัน
ได้จุดสีเทา

```
image = mpimg.imread('monument.png')
```

```
plt.subplot(1, 2, 1)
```

```
plt.imshow(image)
```

```
gray = np.ndarray(image.shape)
```

```
gray[:, :, 0] = \
```

```
gray[:, :, 1] = \
```

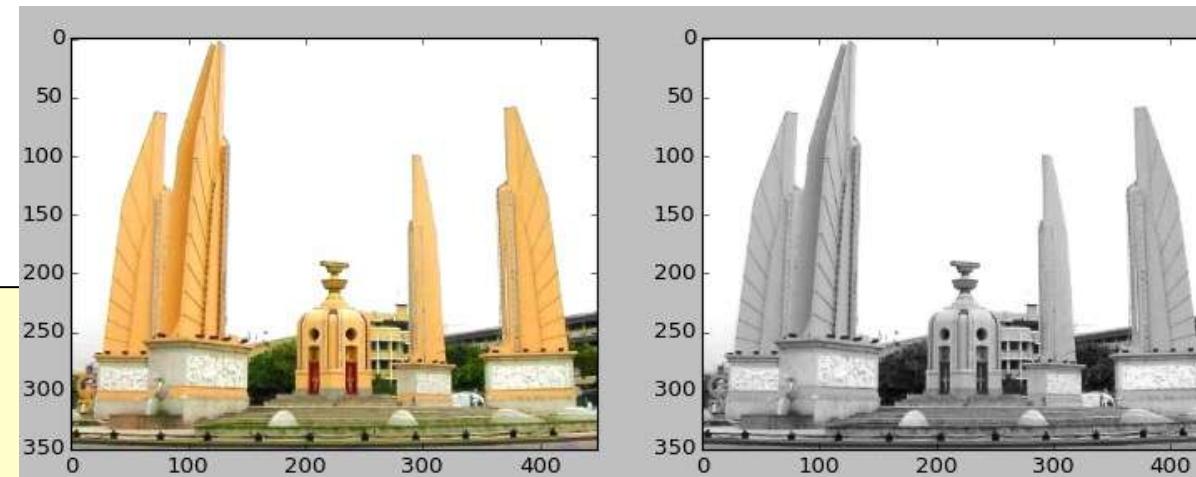
```
gray[:, :, 2] = (image[:, :, 0]+image[:, :, 1]+image[:, :, 2]) / 3
```

```
plt.subplot(1, 2, 2)
```

```
plt.imshow(gray)
```

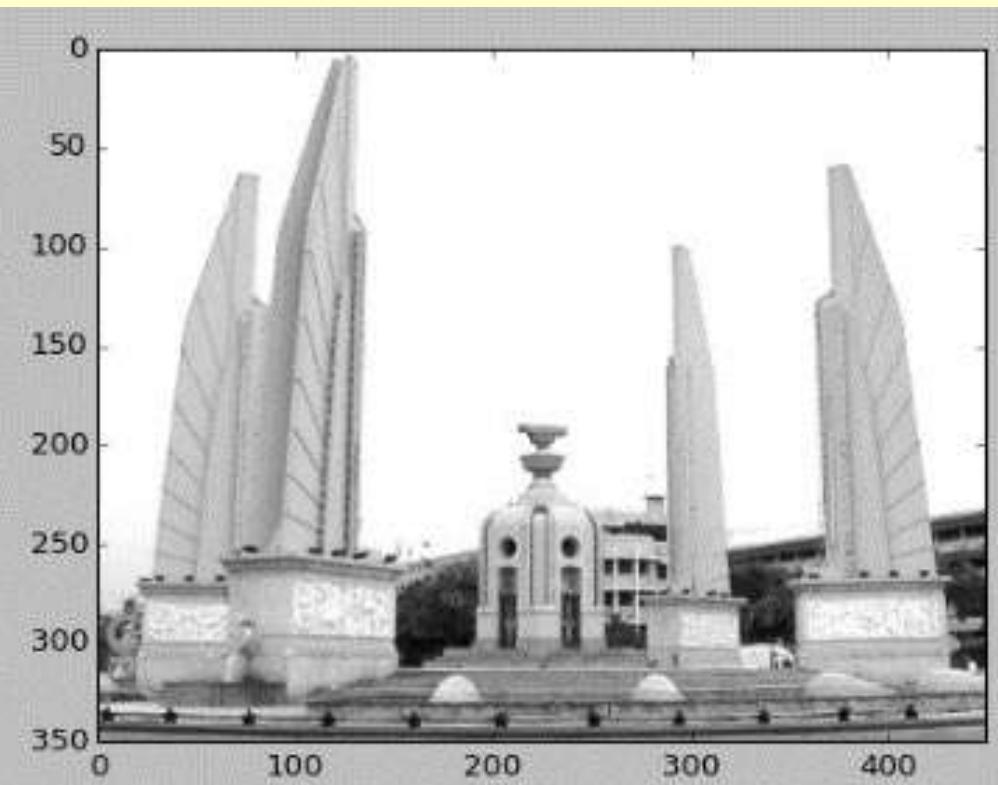
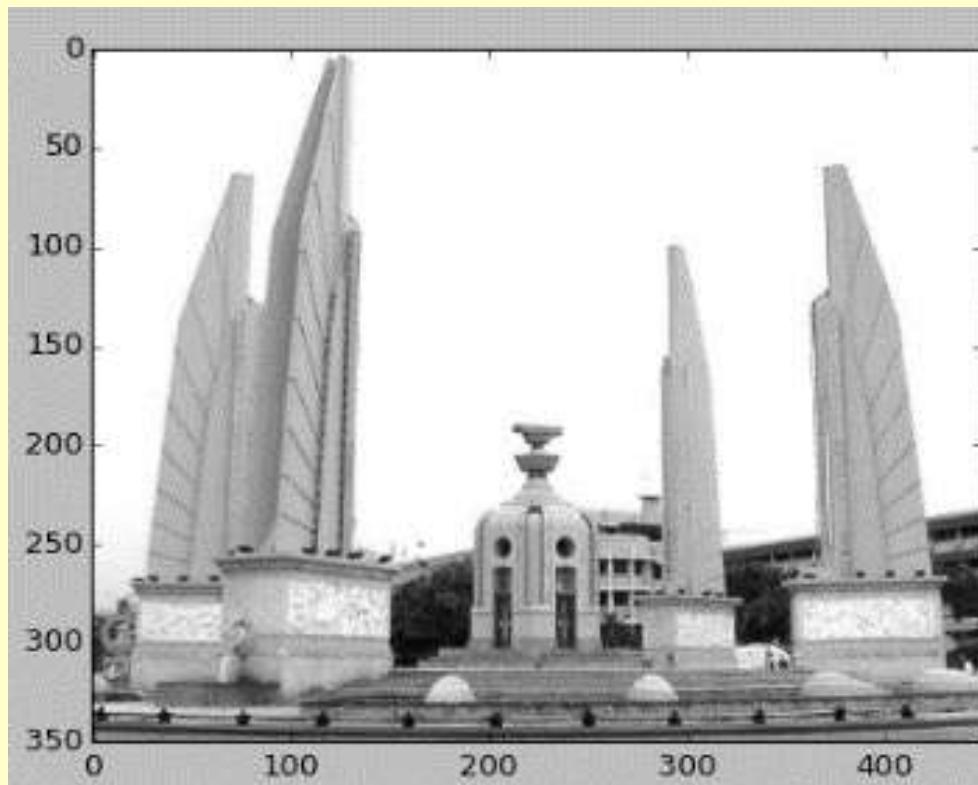
```
plt.show()
```

นำค่าความเข้มของ R G B มาหาค่าเฉลี่ย^{แล้วเปลี่ยน R G B ให้เป็นความเข้มระดับเดียวกัน}

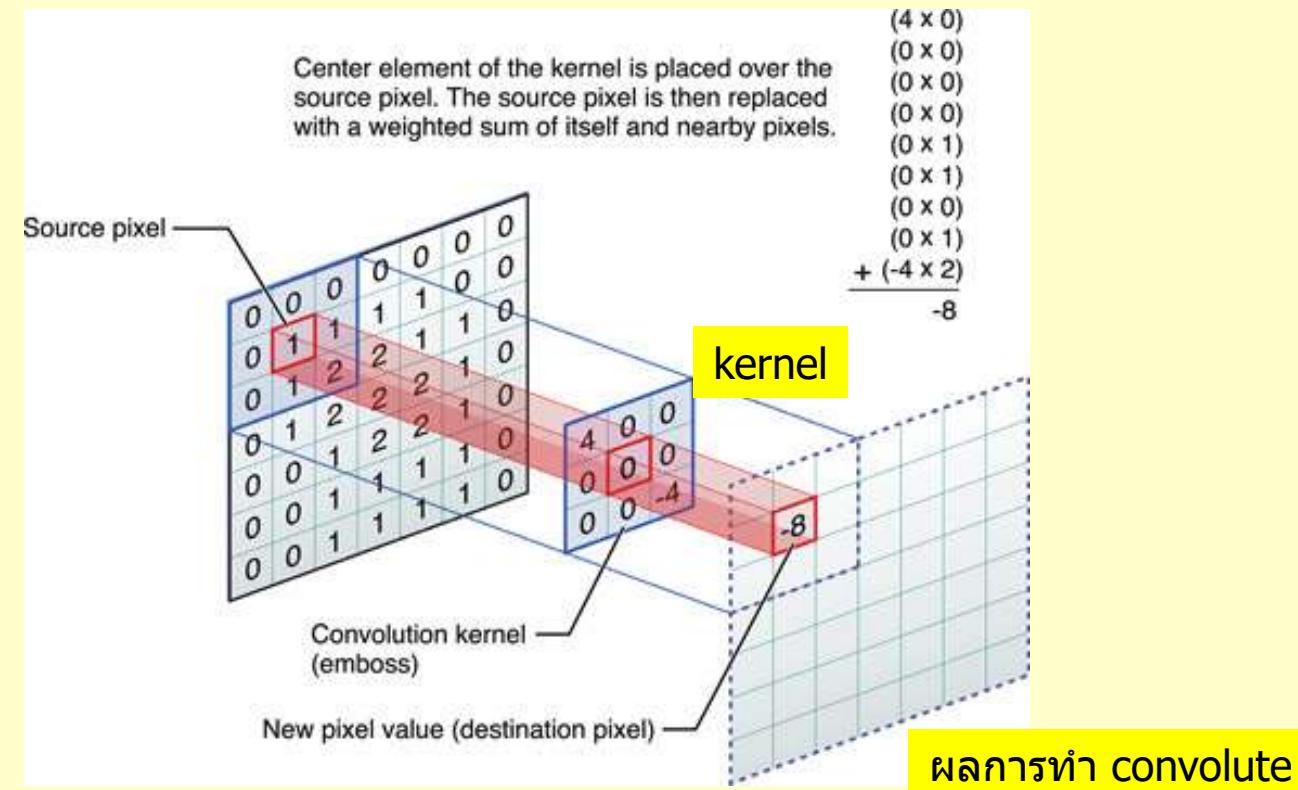


สีเทา เข้าว่าสูตรนี้ดีกว่า $0.299*R + 0.587*G + 0.114*B$

- CH08_6 จงเขียนโปรแกรมเพื่อแปลงรูปภาพจากรูปสีให้เป็นรูปสีเทา (gray scale) ด้วยสูตรข้างบนนี้ เทียบกับสูตร $(R + G + B)/3$



Convolution



เปลี่ยนค่าของเมทริกซ์ kernel
จะได้การประมวลผลภาพแบบอื่น ๆ

Image Convolution

- Image Convolution คือการนำรูปภาพมาผ่านตัวกรอง (kernel) เพื่อให้ได้ผลลัพธ์ตามที่ต้องการ เช่น blur, ขยายรูป, จับขอบรูป เป็นต้น

The diagram shows the convolution process between an original image and a mean filter kernel. The original image is a grayscale portrait of a person's face. The kernel is a 3x3 matrix with all elements equal to 1/9. The result is a blurred version of the original image, labeled "Blur (with a mean filter)".

Original

$*$ $\frac{1}{9} \begin{matrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{matrix}$ =

Blur (with a mean filter)

The diagram shows the convolution process between an original image and a shift-left kernel. The original image is a grayscale portrait of a person's face. The kernel is a 3x3 matrix where the central element is 1 and all other elements are 0. The result is a version of the original image shifted one pixel to the left, labeled "Shifted left By 1 pixel".

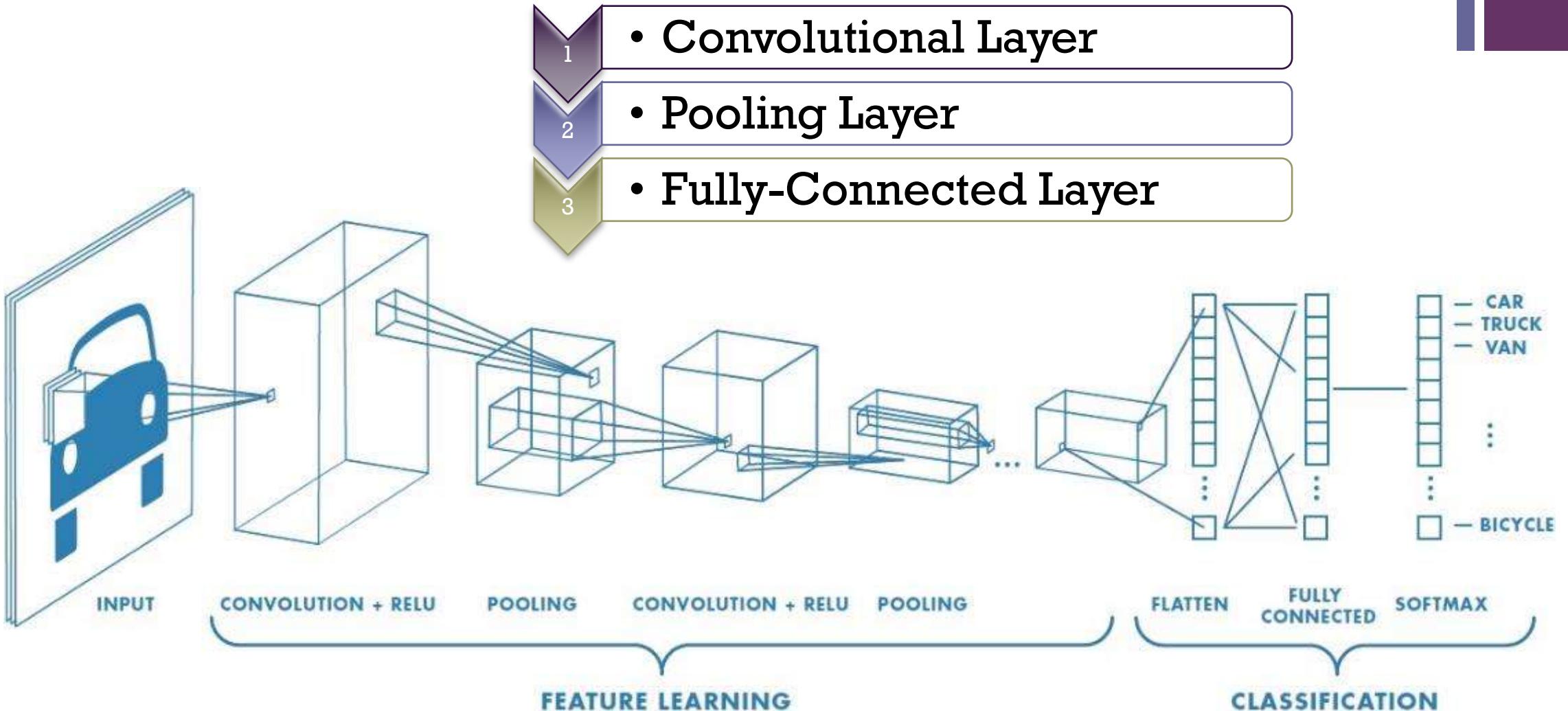
Original

$*$ $\begin{matrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 0 \end{matrix}$ =

Shifted left
By 1 pixel

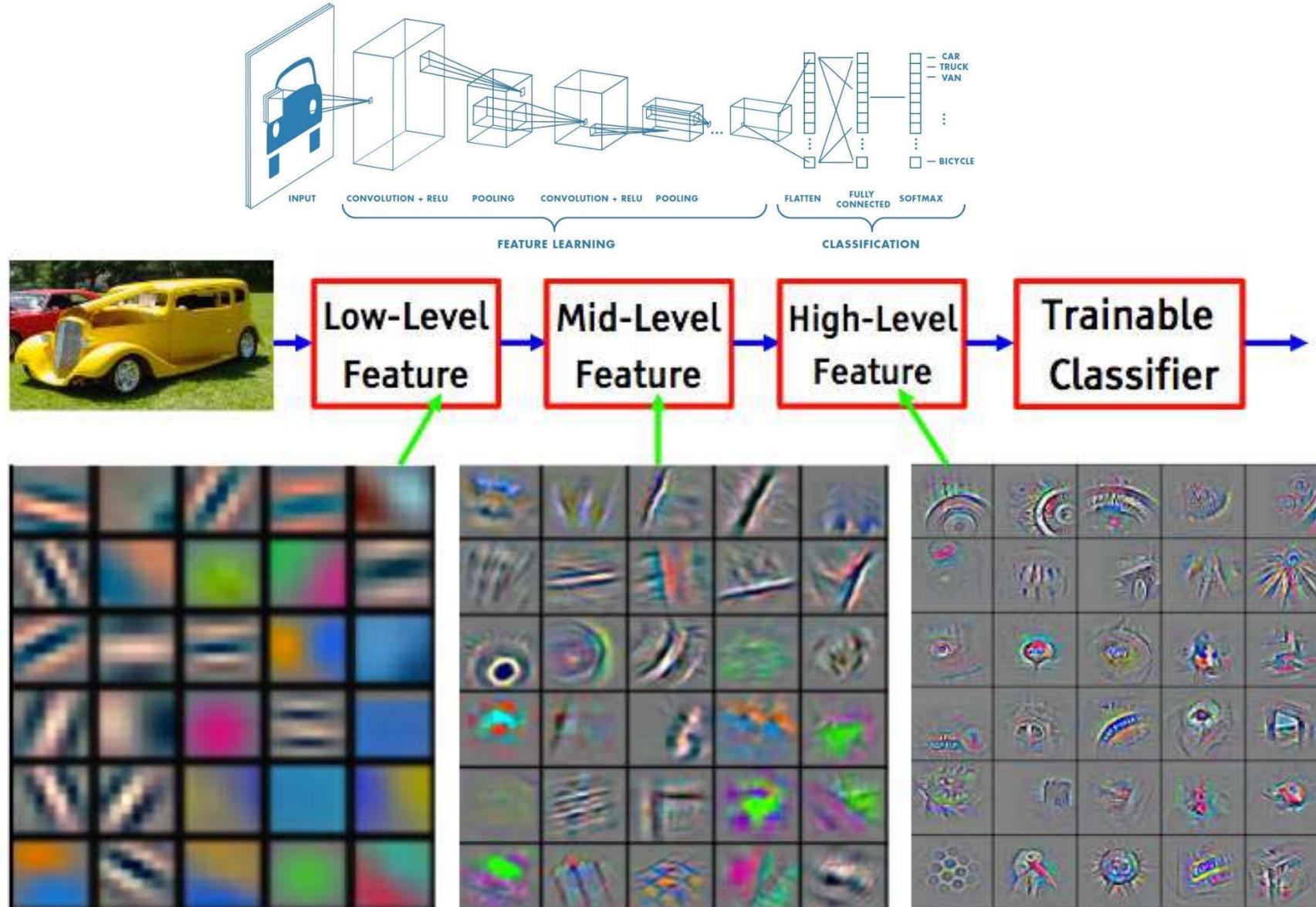


Convolutional Neural Networks (CNN)





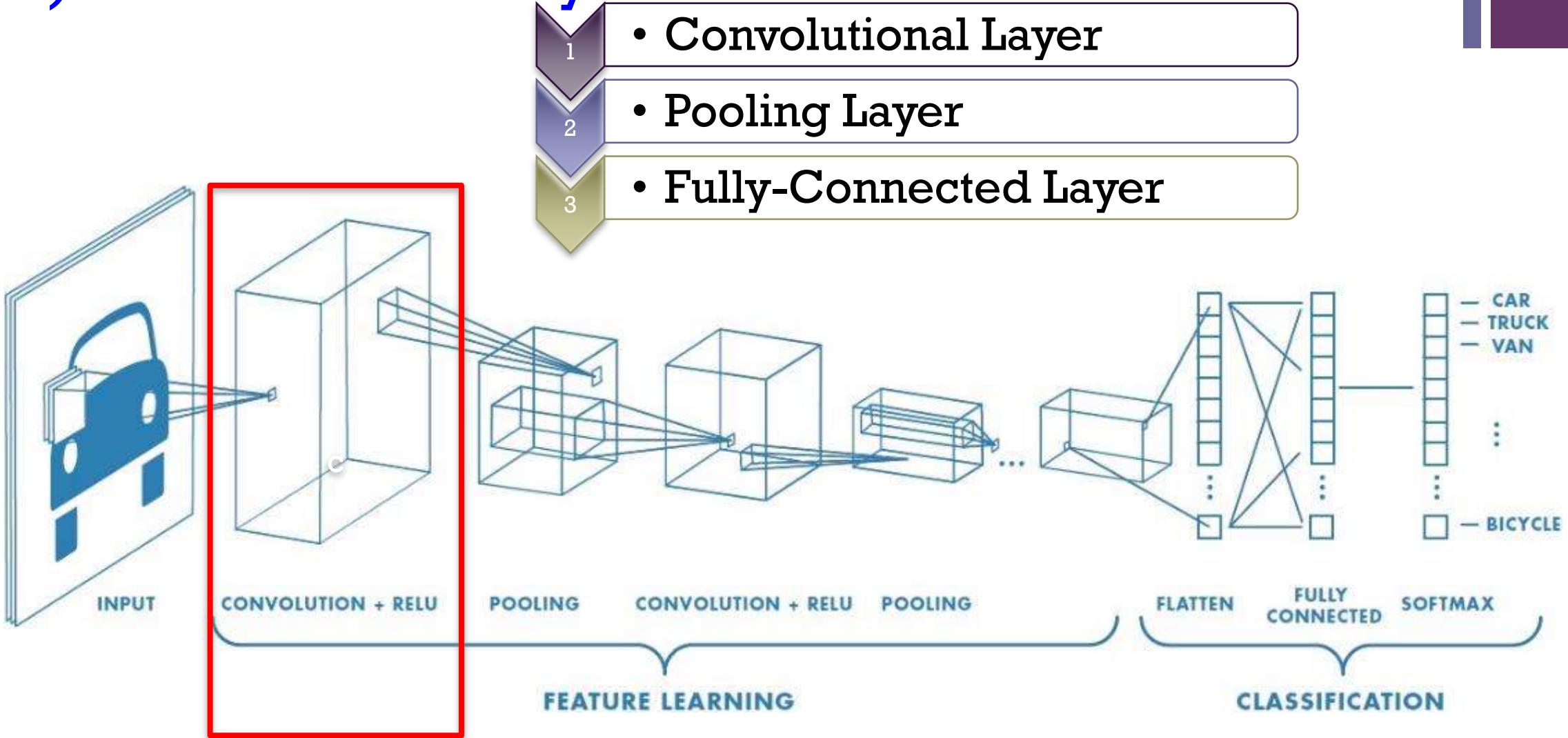
CNN (cont.)





Convolutional Neural Networks (CNN)

1) Convolutional Layer





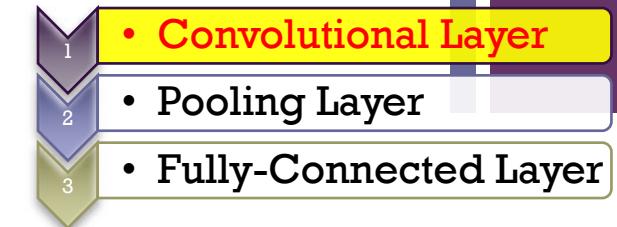
Layer1: Convolutional Layer

1	1	1	0	0
0	1	1	1	0
0	0	1	1	1
0	0	1	1	0
0	1	1	0	0

Image

(3*3 filter)

1	0	1
0	1	0
1	0	1



4		

Convolved
Feature



com Input Volume (+pad 1) (7x7x3)

x[:, :, 0]	0 0 0 0 0 0 0	0 1 0 2 2 2 0	0 0 0 0 0 0 0	0 2 0 2 2 2 0	0 1 0 0 0 0 0	0 1 0 0 2 1 0	0 0 0 0 0 0 0
x[:, :, 1]	0 0 0 0 0 0 0	0 1 0 1 0 1 0	0 0 0 0 0 1 0	0 0 0 2 0 0 0	0 2 0 0 0 0 0	0 0 0 1 0 0 0	0 0 0 0 0 0 0
x[:, :, 2]	0 0 0 0 0 0 0	0 0 2 2 0 0 0	0 0 0 0 0 1 0	0 1 1 2 0 2 0	0 2 0 0 0 0 0	0 0 1 1 0 1 0	0 0 0 0 0 0 0

Filter W0 (3x3x3)

w0[:, :, 0]	1 0 0	-1 0 0	0 -1 1	-1 -1 0	-1 -1 1	1 0 0	1 1 0
w0[:, :, 1]	0 0 0	0 1 -1	0 0 0	0 1 1	0 0 0	0 1 1	0 -1 0
w0[:, :, 2]	0 0 0	1 1 -1	0 1 1	1 1 1	1 1 1	1 1 -1	1 0 1
Bias b0 (1x1x1)	1						
b0[:, :, 0]							

Filter W1 (3x3x3)

w1[:, :, 0]	1 -1 0	-1 0 0	0 1 -1	0 1 -1	0 0 0	0 1 1	0 -1 0
w1[:, :, 1]	0 1 -1	-2 4 -1	3 3 0	-2 2 -1			
w1[:, :, 2]	1 1 -1	-1 0 1					
Bias b1 (1x1x1)							
b1[:, :, 0]	0						

Output Volume (3x3x2)

o[:, :, 0]	2 -2 -1	2 -3 -3	2 -1 -2	o[:, :, 1]	-2 4 -1	3 3 0	-2 2 -1

$$\begin{aligned}
 o(2,2,0) &= \sum x[:, :, 0] \times w[:, :, 0] + \sum x[:, :, 1] \times w[:, :, 1] + \sum x[:, :, 2] \times w[:, :, 2] + b_0 \\
 &= 0 \times 1 + 0 \times 0 + 0 \times 0 + 2 \times (-1) + 1 \times 0 + 0 \times 0 + 0 \times 0 + 0 \times (-1) + 0 \times 1 \\
 &\quad + 0 \times (-1) + 0 \times (-1) + 0 \times 0 + 0 \times (-1) + 0 \times (-1) + 0 \times 1 + 0 \times 1 + 0 \times 0 + 0 \times 1 \\
 &\quad + 0 \times 1 + 0 \times 1 + 0 \times 0 + 0 \times 0 + 1 \times (-1) + 0 \times 1 + 0 \times 1 + 0 \times 1 + 0 \times 1 \\
 &\quad + 1 \\
 &= -2
 \end{aligned}$$

CNNs: 1) Convolutional Layer

Feature Extraction

<http://cs231n.github.io/convolutional-networks/>

(3*3 filter)

$$w^T x + b$$

1	1	1	1	0	0
0	1	1	1	1	0
0	0	1	1	1	1
0	0	1	1	0	0
0	1	1	0	0	0

Image

Convolved
Feature
(feature map)

1	0	1
0	1	0
1	0	1

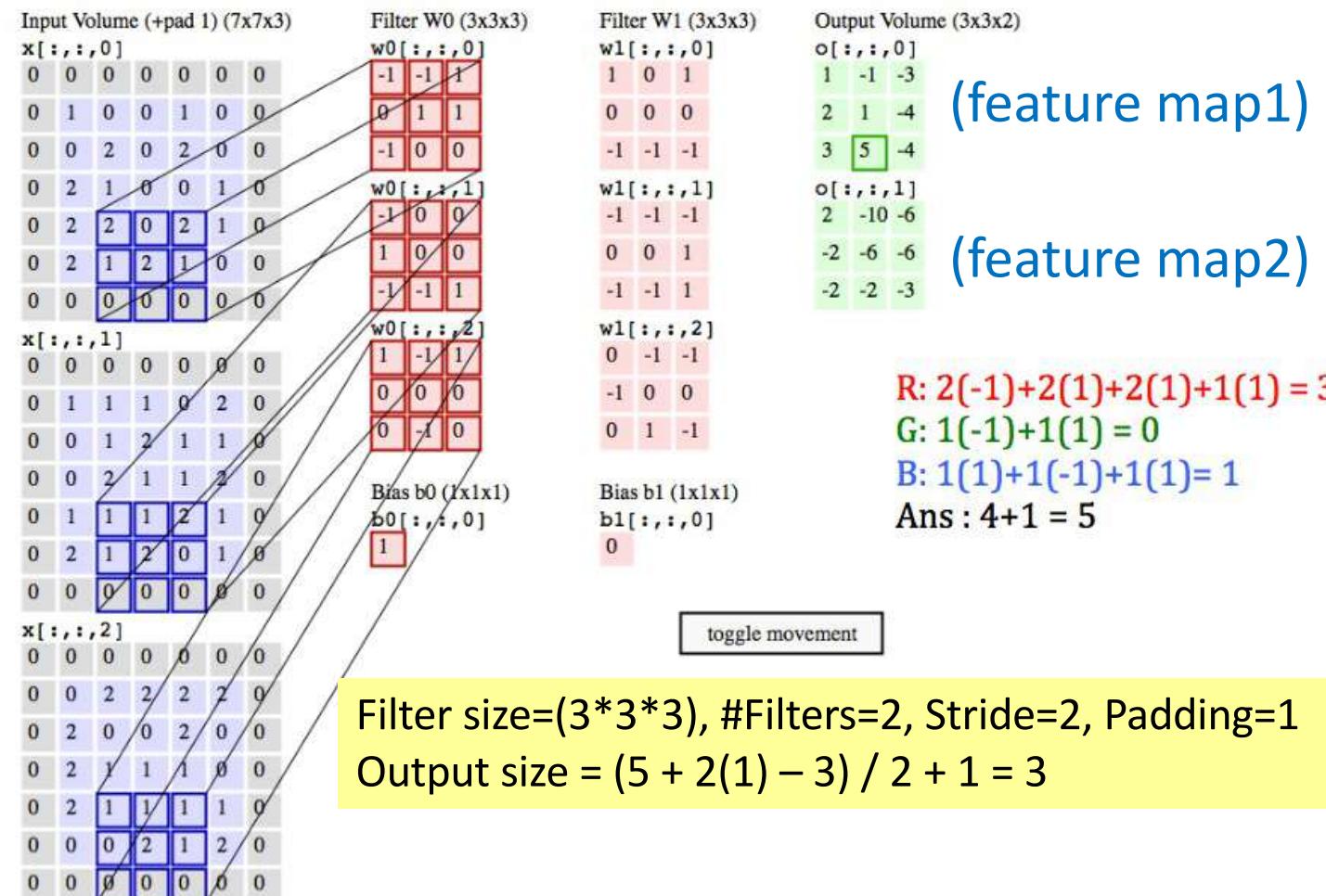
4		

Filter size=(3*3), #Filters=1, Stride=1, Padding=0
 Output size = $(5 + 2(0) - 3) / 1 + 1 = 3$

Summary: **4 parameters** for convolutional layer

- (1) Filter size, (2) #Filters, (3) Stride, (4) Padding

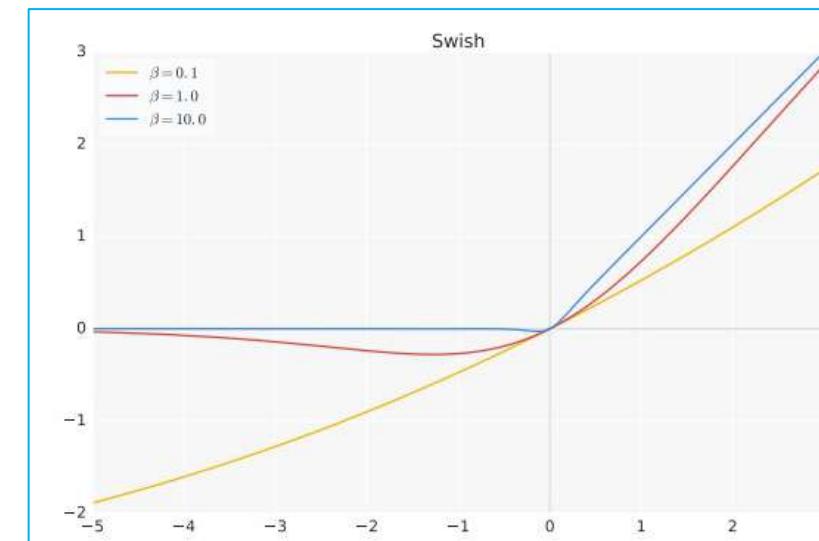
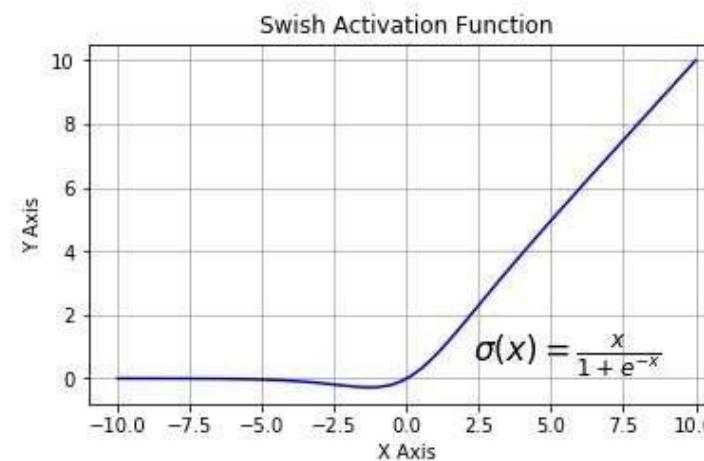
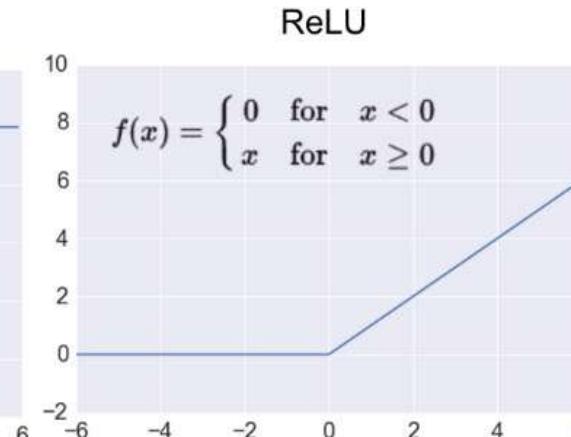
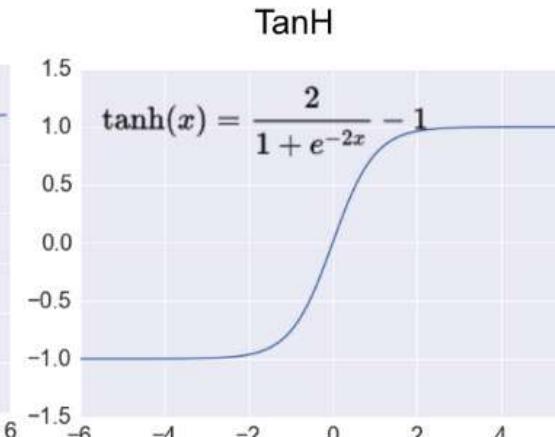
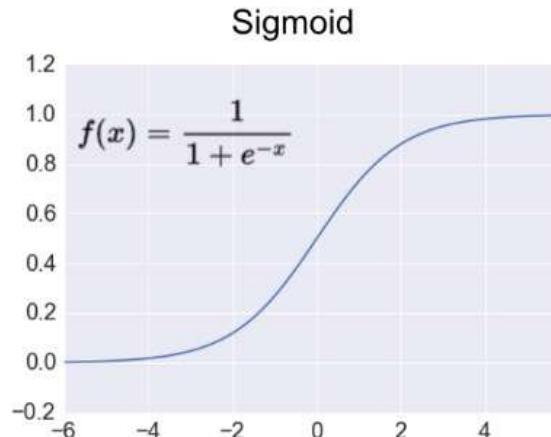
$$\text{Output size} = (N + 2*P - F) / S + 1$$





CNNs: 1) Convolutional Layer (cont.)

Non-Linear Activation Function

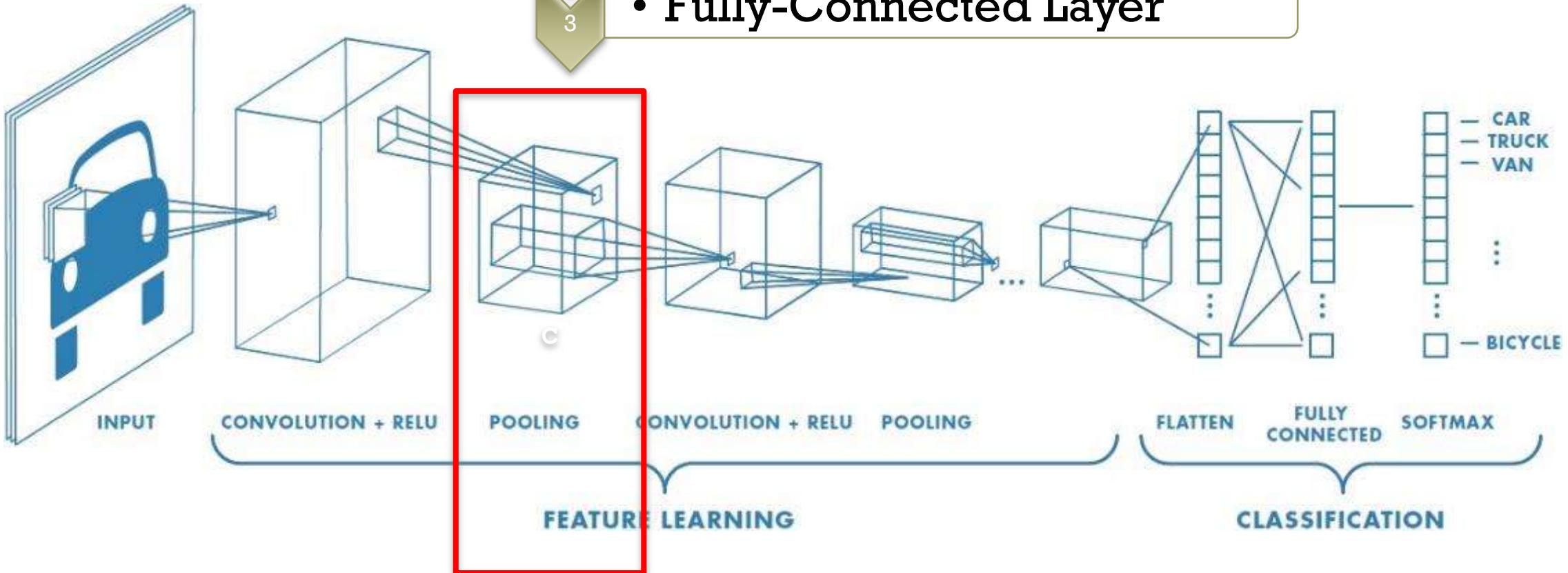




Convolutional Neural Networks (CNN)

2) Pooling Layer

- Convolutional Layer
- Pooling Layer
- Fully-Connected Layer



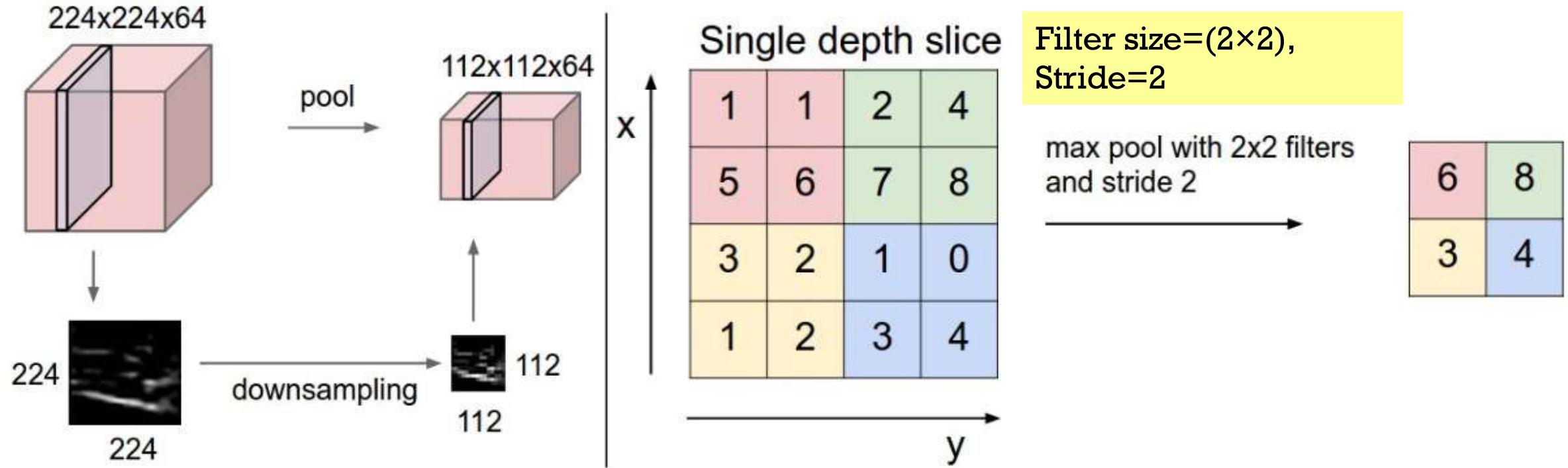
CNNs: 2) Pooling Layer

<http://cs231n.github.io/convolutional-networks/>

Summary: **2 parameters** for pooling layer

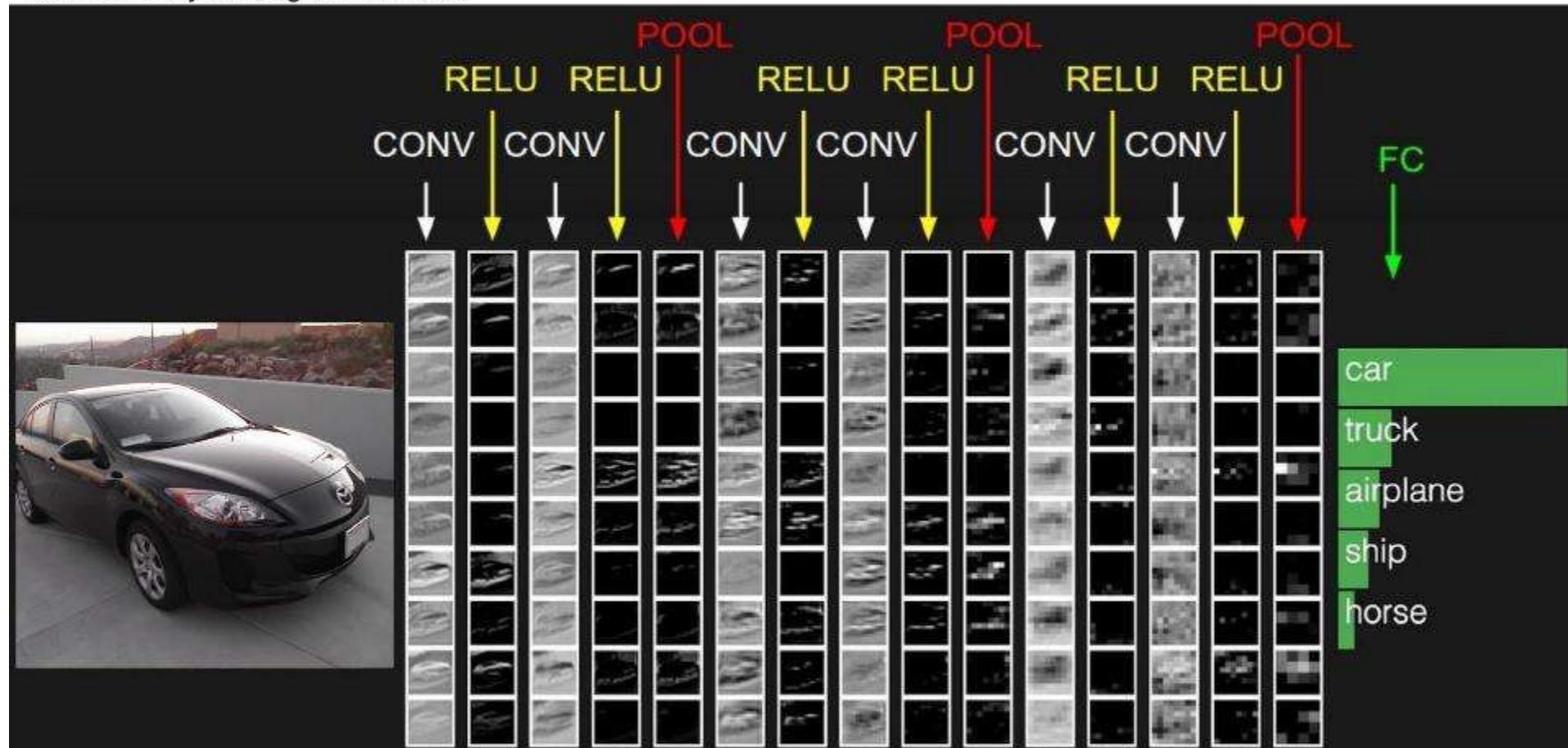
- Filter size, Stride

- Feature selection (reduction); downsampling





two more layers to go: POOL/FC

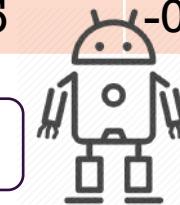




Convolutional Neural Networks

Embedding Vector

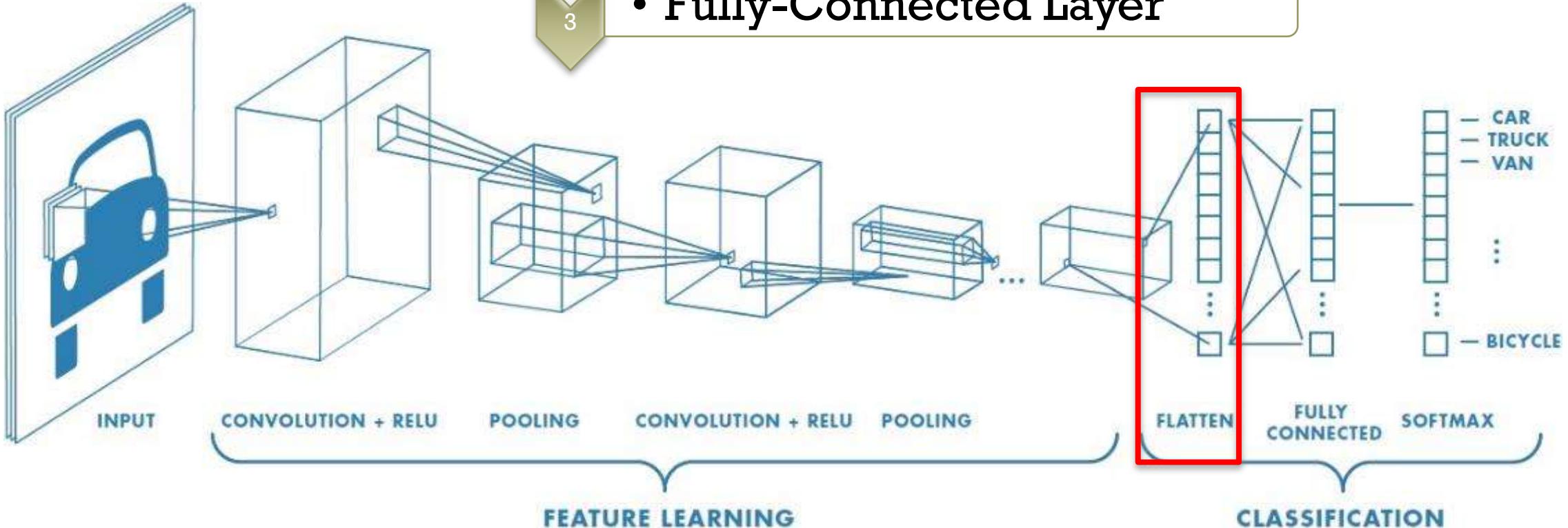
x1	x2	x3	x4	Corona
0.7	0.2	-0.5	-0.1	Yes



- Convolutional Layer

- Pooling Layer

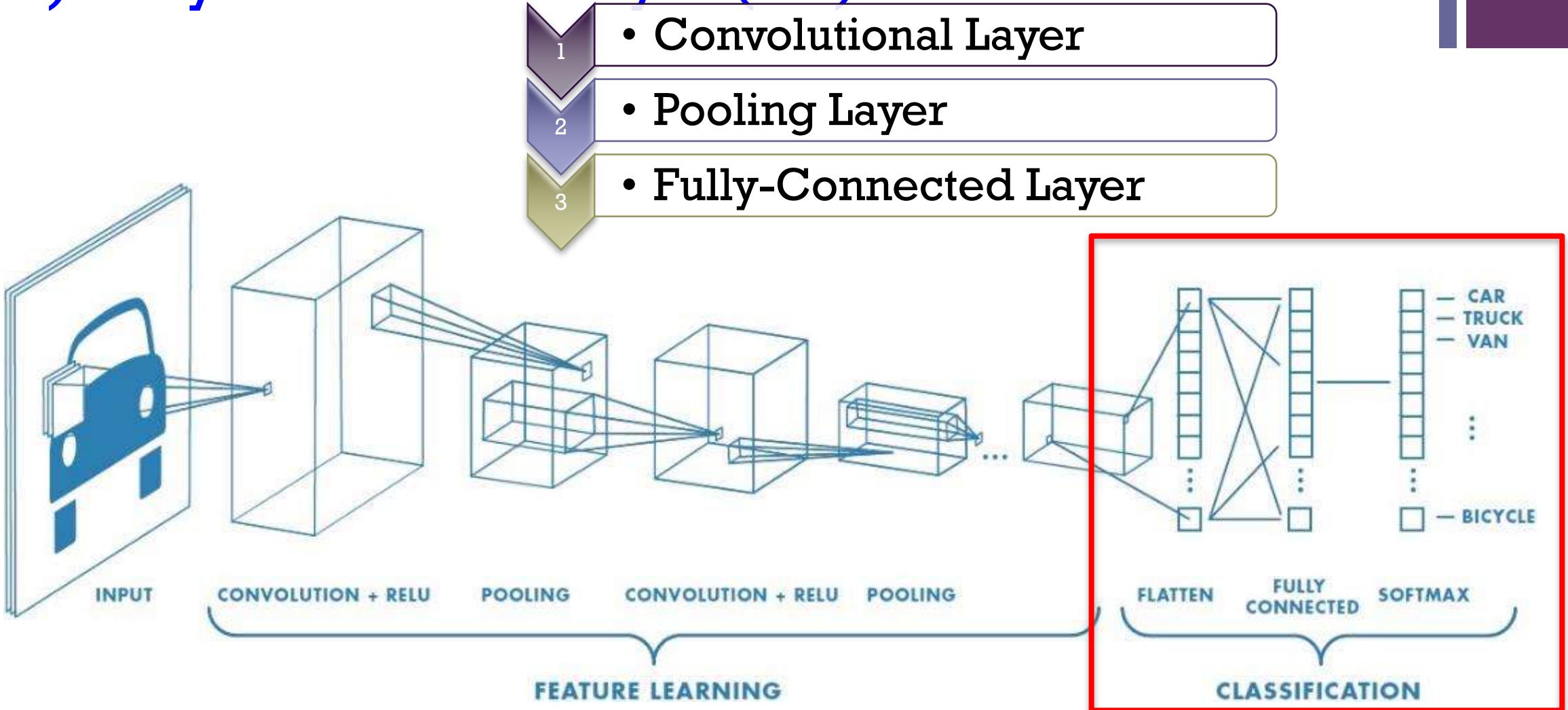
- Fully-Connected Layer



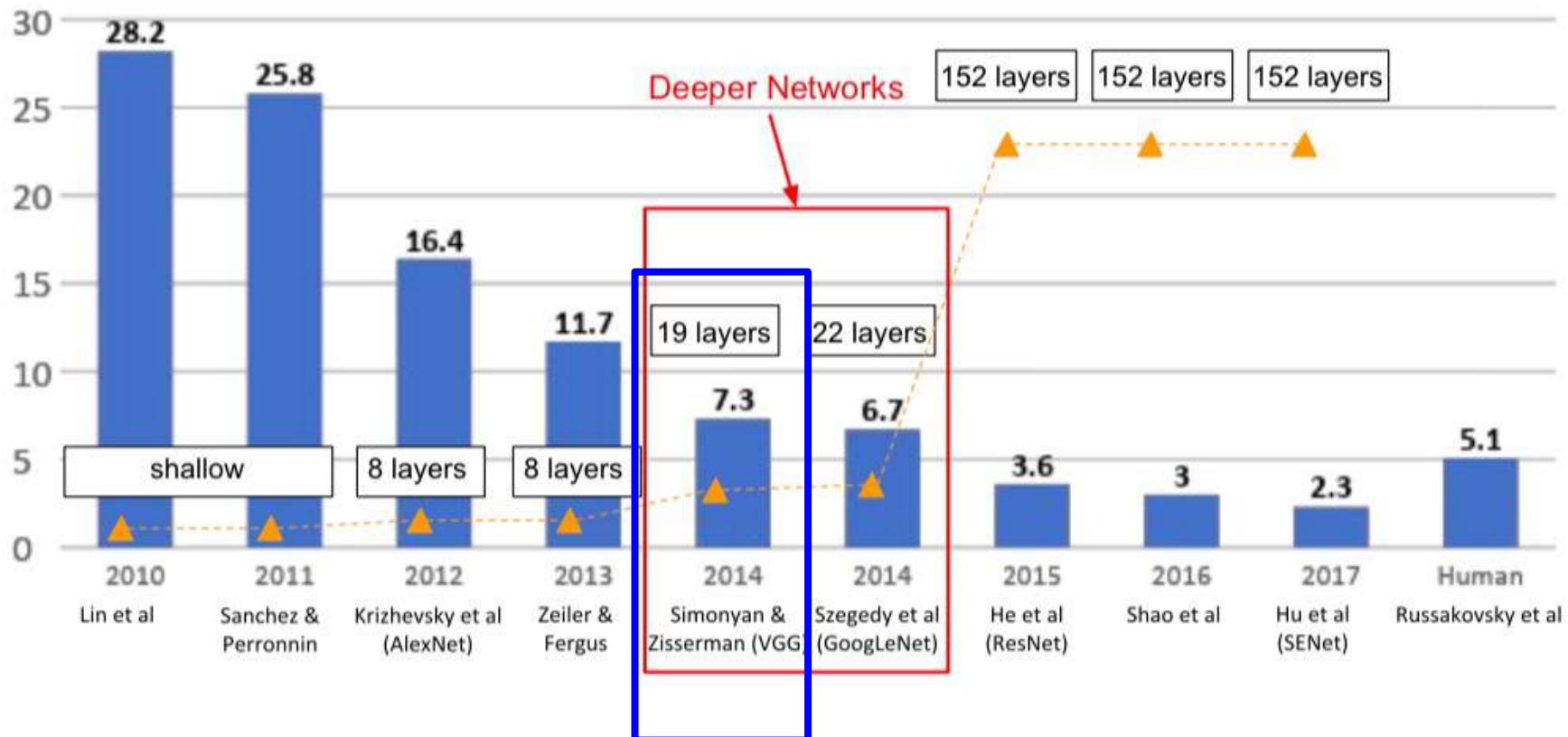


Convolutional Neural Networks (CNN)

3) Fully-Connected Layer (FC) = Neural Networks



ImageNet Large Scale Visual Recognition Challenge (ILSVRC 2nd runner-up in 2014): VGG19



Case Study: VGGNet (2014)

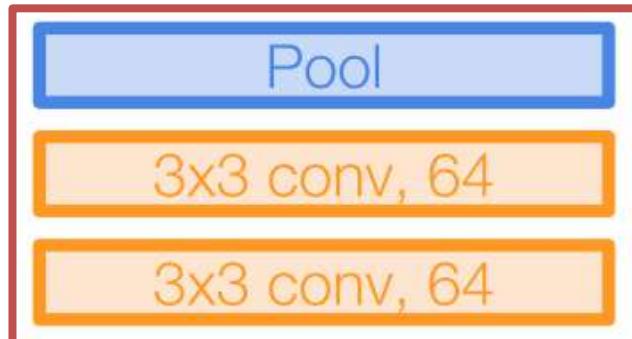
[Simonyan and Zisserman, 2014]

Small filters, Deeper networks

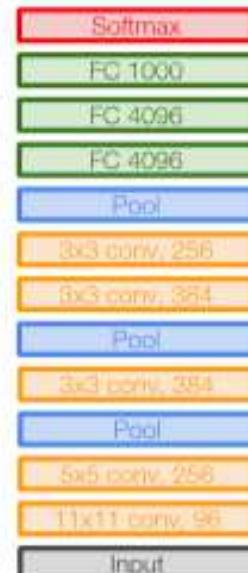
8 layers (AlexNet) → 16-19 layers (VGG16Net)

Only 3x3 CONV Stride 1, pad 1

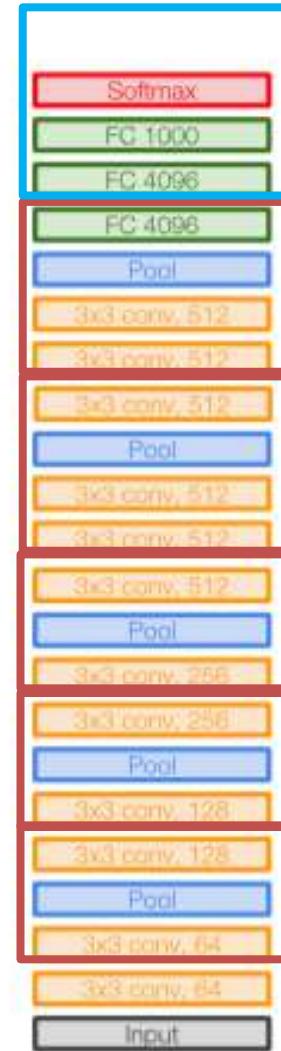
And 2x2 MAX POOL stride 2



1) Feature extraction block 2) Classification block

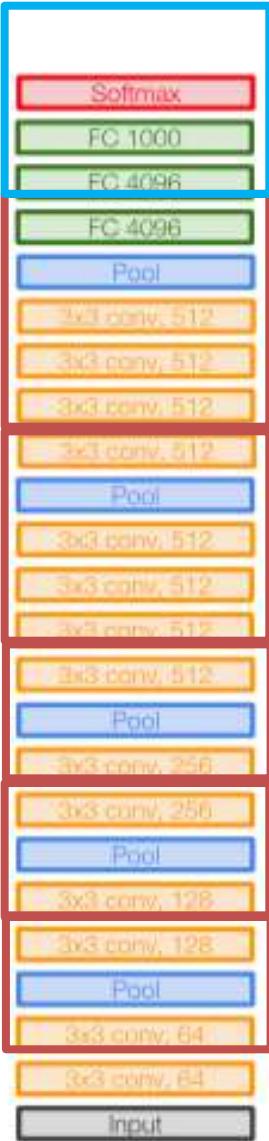


AlexNet



47

VGG16



VGG19

Case Study: VGGNet (2014)

[Simonyan and Zisserman, 2014]

Q: Why use smaller filters? (3x3 conv)

Stack of three 3x3 conv (stride 1) layers has same **effective receptive field** as one 7x7 conv layer in ZFNet

But deeper, more non-linearities

And fewer parameters: $3 \times (3^2 C^2)$ vs. $7^2 C^2$ for C channels per layer

K Keras

```
#VGG Block
conv1 = Conv2D(64, (3,3), strides=(1,1), padding='same', activation='ReLU')(x)
conv2 = Conv2D(64, (3,3), strides=(1,1), padding='same', activation='ReLU')(conv1)
maxpool1 = MaxPooling2D((2,2))(conv2)
```



AlexNet

VGG16

VGG19

Padding='valid'
No padding
48

Variants of VGG

- ConvNet configurations
 - The depth of the configurations increases from the left (A) to the right (E), as more layers are added (the added layers are shown in bold)
- The convolutional layer parameters are denoted as “conv<receptive field size>-<number of channels>”
- ILSVRC’14 2nd in classification, 1st in localization
- Use VGG16 or VGG 19 (VGG19 only slightly better, more memory)

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-4096	FC-1000	soft-max	
VGG16					
VGG19					

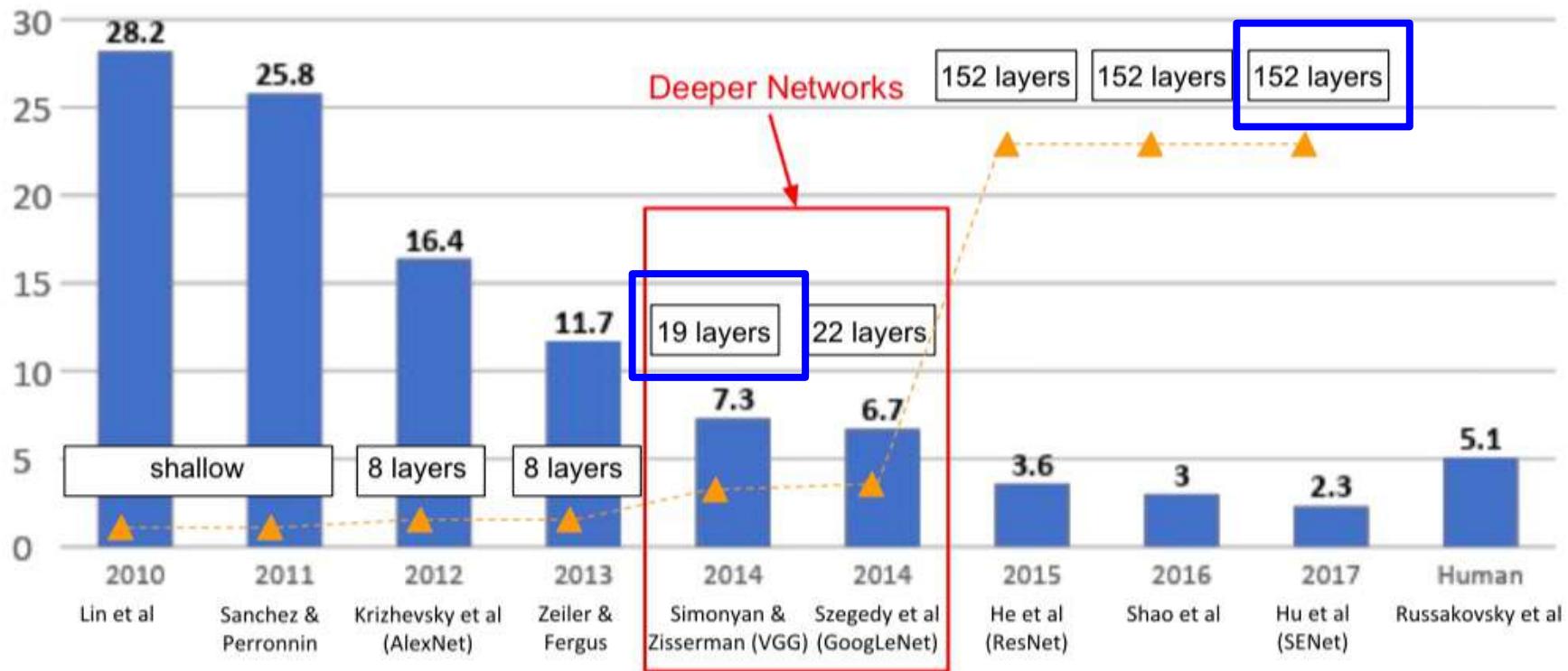
Calculate the number of parameters on VGG Net

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
maxpool					

Network	A	B	C	D	E
Number of parameters (Millions)	133	133	134	138	144
maxpool					
FC-4096					
FC-4096					
FC-1000					
soft-max					

<https://hoanglehaithanh.com/calculate-the-number-of-parameters-on-vgg-net/>

ImageNet Large Scale Visual Recognition Challenge (ILSVRC 2nd runner-up in 2014): Deeper than VGG19



[Overview](#)[Input](#)[Model](#)[Sequential](#)[activations](#)[applications](#)[Overview](#)[DenseNet121](#)[DenseNet169](#)[DenseNet201](#)[EfficientNetB0](#)[EfficientNetB1](#)[EfficientNetB2](#)[EfficientNetB3](#)[EfficientNetB4](#)[EfficientNetB5](#)[EfficientNetB6](#)[EfficientNetB7](#)[InceptionResNetV2](#)[InceptionV3](#)[MobileNet](#)[MobileNetV2](#)[MobileNetV3Large](#)[MobileNetV3Small](#)[NASNetLarge](#)[NASNetMobile](#)

Modules

[densenet](#) module: DenseNet models for Keras.

[efficientnet](#) module: EfficientNet models for Keras.

[imagenet_utils](#) module: Utilities for ImageNet data preprocessing & prediction decoding.

[inception_resnet_v2](#) module: Inception-ResNet V2 model for Keras.

[inception_v3](#) module: Inception V3 model for Keras.

[mobilenet](#) module: MobileNet v1 models for Keras.

[mobilenet_v2](#) module: MobileNet v2 models for Keras.

[nasnet](#) module: NASNet-A models for Keras.

[resnet](#) module: ResNet models for Keras.

[resnet50](#) module: Public API for tf.keras.applications.resnet50 namespace.

[resnet_v2](#) module: ResNet v2 models for Keras.

[vgg16](#) module: VGG16 model for Keras.

[vgg19](#) module: VGG19 model for Keras.

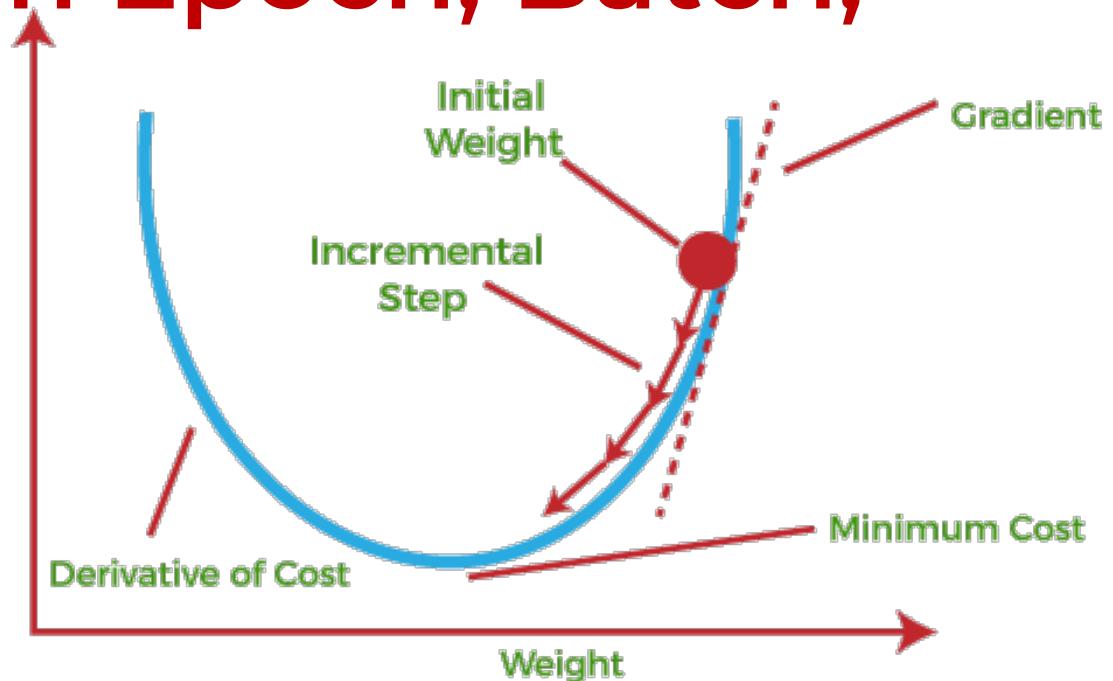
Model 1

- VGG19 (random initialized weights) + 2 Dense layers + Output layer

```
1 base_model = VGG19(weights=None, include_top=False, input_shape=(224, 224, 3))
2
3 for layer in base_model.layers:
4     layer.trainable = True
5
6 x = base_model.output
7 x = Flatten()(x)
8 x = Dense(1024)(x)
9 x = Dropout(0.5)(x)
10 x = Dense(512)(x)
11 x = Dropout(0.5)(x)
12 output = Dense(num_class, activation='softmax')(x)
13
```

Differences Between Epoch, Batch, and Mini-batch

- The gradient descent algorithm works in two steps that are performed continuously for a specific number of iterations:
 - First, we compute the **gradient**, which is the first-order derivative of the objective function with respect to the variables.
 - Then, we **update the variables** in the opposite direction of the gradient
 - Terms:
 - Epoch = use **all** training examples
 - Iteration = each round of **updating weight**

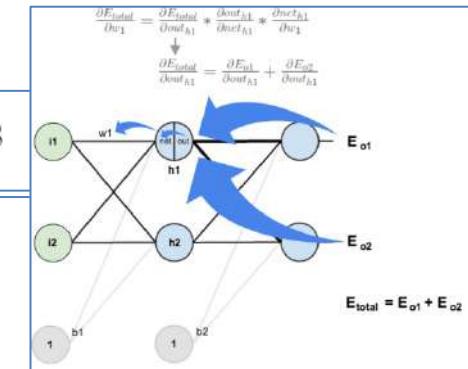


$$w_5^+ = w_5 - \eta * \frac{\partial E_{total}}{\partial w_5} = 0.4 - 0.5 * 0.082167041 = 0.35891648$$

Putting it all together:

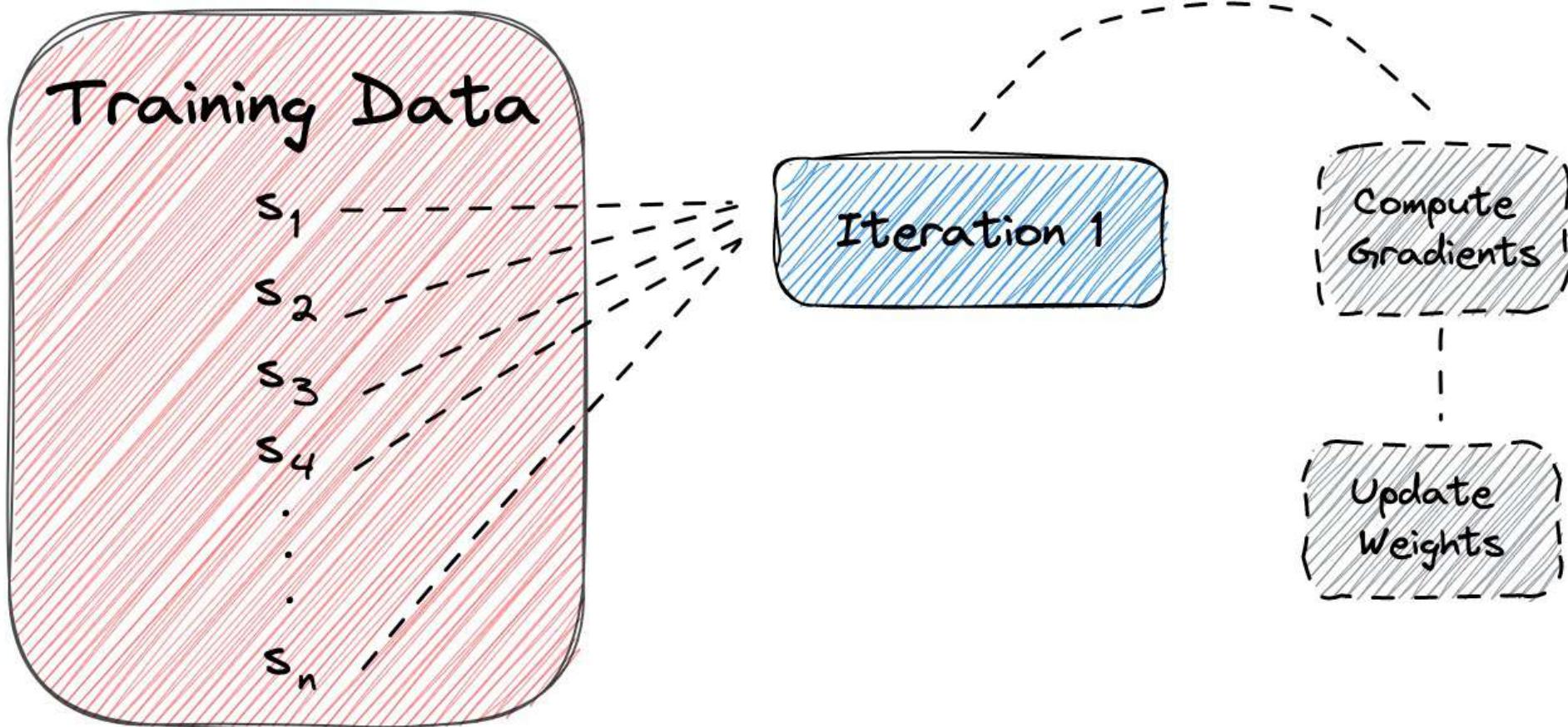
$$\frac{\partial E_{total}}{\partial w_5} = \frac{\partial E_{total}}{\partial out_{o1}} * \frac{\partial out_{o1}}{\partial net_{o1}} * \frac{\partial net_{o1}}{\partial w_5}$$

$$\frac{\partial E_{total}}{\partial w_5} = 0.74136507 * 0.186815602 * 0.593269992 = 0.082167041$$



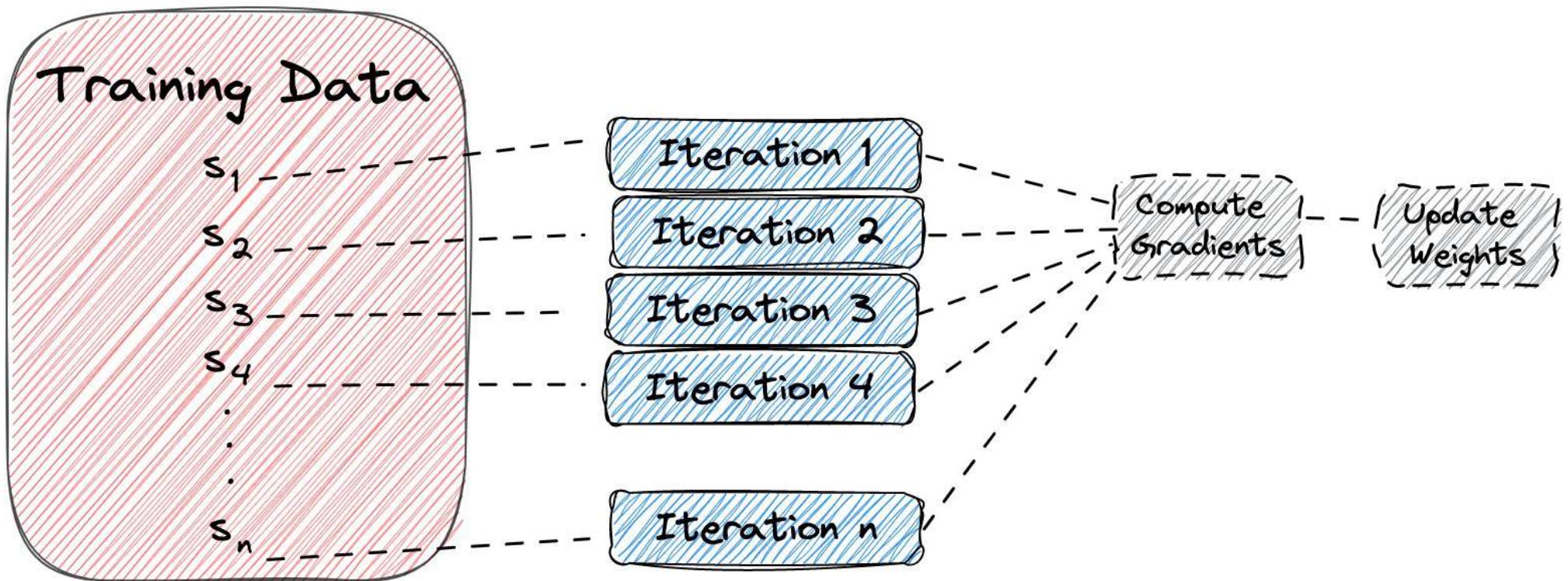
1) Batch Gradient Descent

- In batch gradient descent, we use all our training data in a single iteration of the algorithm.



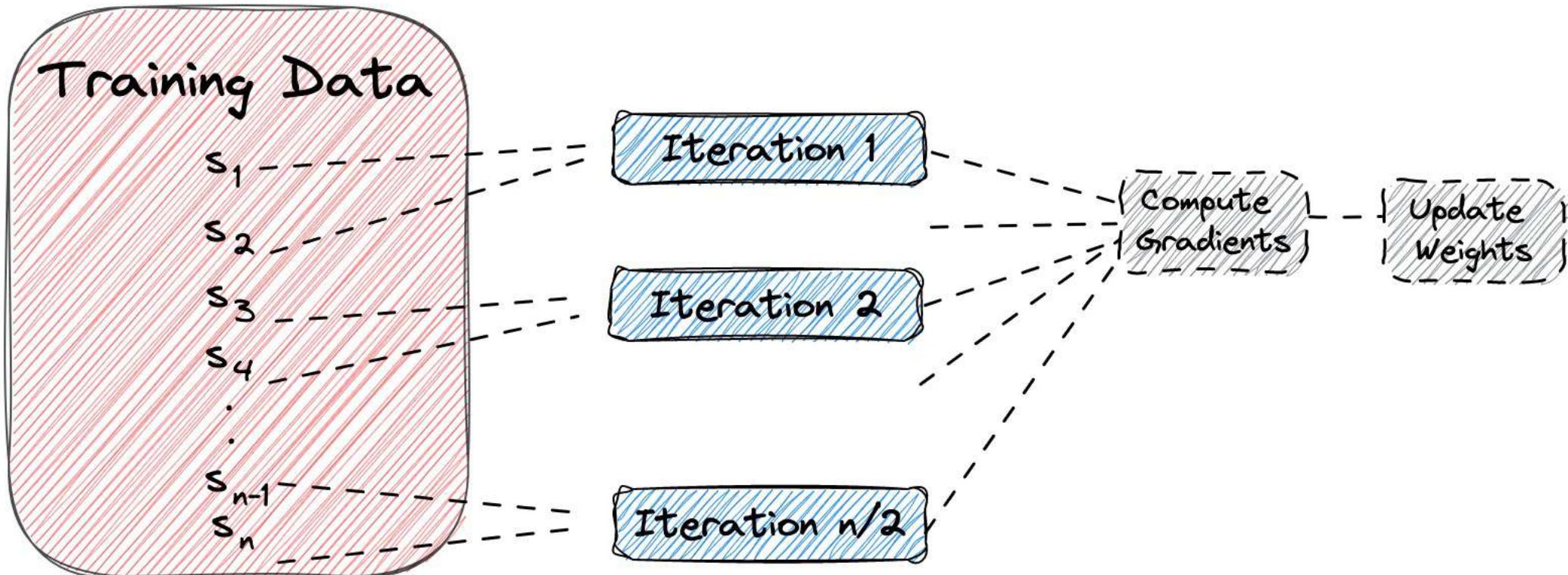
2) Stochastic Gradient Descent

- The previous method can be very time-consuming and inefficient in case the size of the training dataset is large. To deal with this, we use stochastic gradient descent, where we use just one sample in a single iteration of the algorithm.



3) Mini-Batch Gradient Descent

- Mini-batch gradient descent is a **combination of the previous methods** where we use a group of samples called mini-batch in a single iteration of the training algorithm.
- The **mini-batch** is a fixed number of training examples that is **less than the actual dataset**. So, in each iteration, we train the network on a different group of samples until all samples of the dataset are used.



Conclusion about batch

An epoch means that we have passed each sample of the training set one time through the network to update the parameters. Generally, the number of epochs is a hyperparameter that defines the number of times that gradient descent will pass the entire dataset.

If we look at the previous methods, we can see that:

- In batch gradient descent, one epoch corresponds to a single iteration.
- In stochastic gradient descent, one epoch corresponds to n iterations where n is the number of training samples.
- In mini-batch gradient descent, one epoch corresponds to $\frac{n}{b}$ iterations where b is the size of the mini-batch.

Finally, let's present a simple example to better understand the three terms.

Let's assume that we have a dataset with $n = 2000$ samples, and we want to train a deep learning model using gradient descent for 10 epochs and mini-batch size $b = 4$:

- In batch gradient descent, we'll update the network's parameters (using all the data) 10 times which corresponds to 1 time for each epoch.
- In stochastic gradient descent, we'll update the network's parameters (using one sample each time) $2000 * 10 = 20000$ times which corresponds to 2000 times for each epoch.
- In mini-batch gradient descent, we'll update the network's parameters (using $b = 4$ samples each time) $\frac{2000}{4} * 10 = 5000$ times that corresponds to $\frac{2000}{4} = 500$ times for each epoch.



Image Classification Tasks

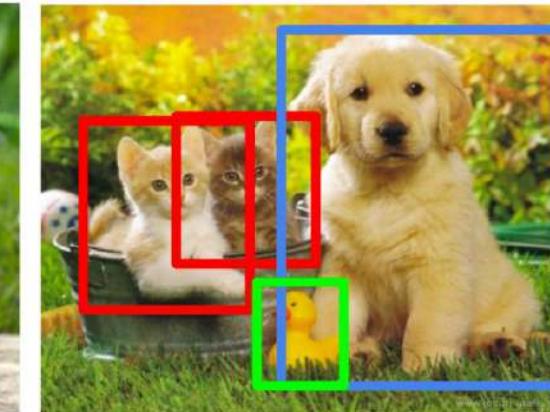
Classification



Classification + Localization



Object Detection



Instance Segmentation



CAT

CAT

CAT, DOG, DUCK

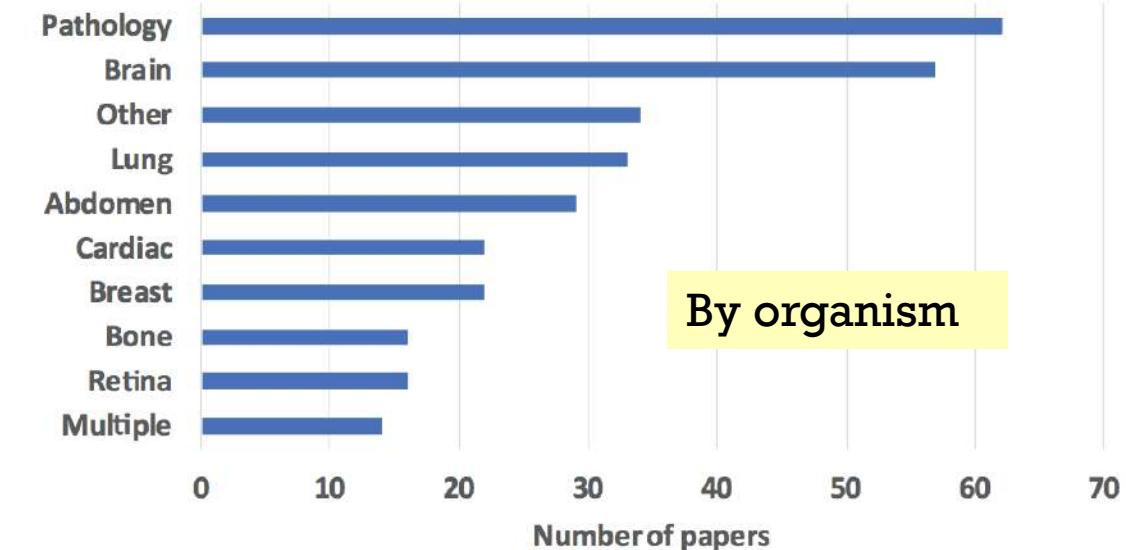
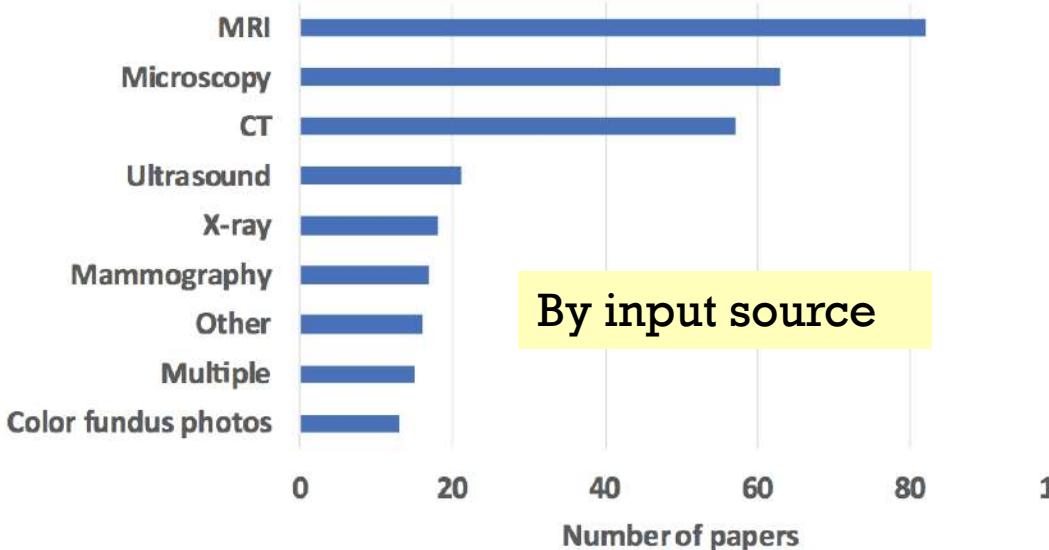
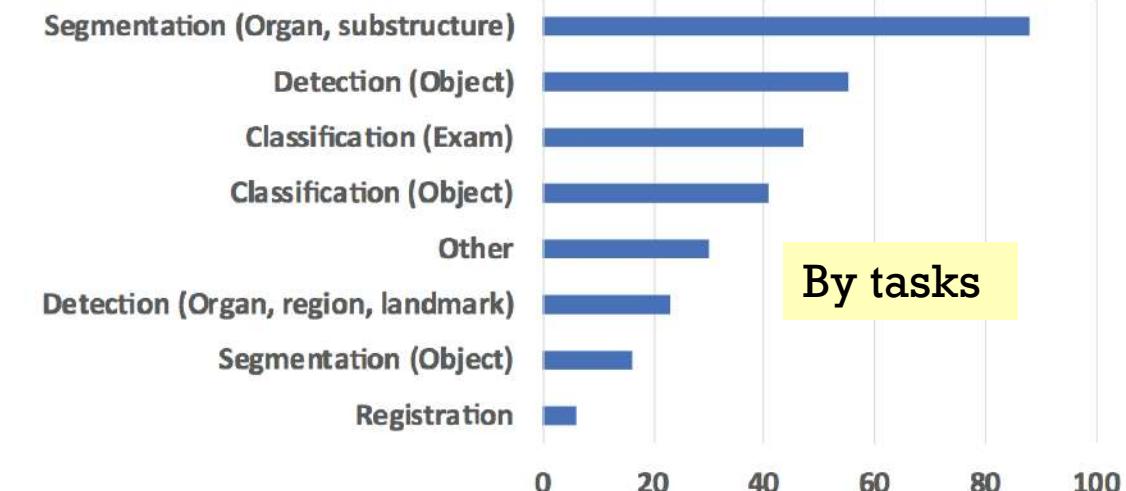
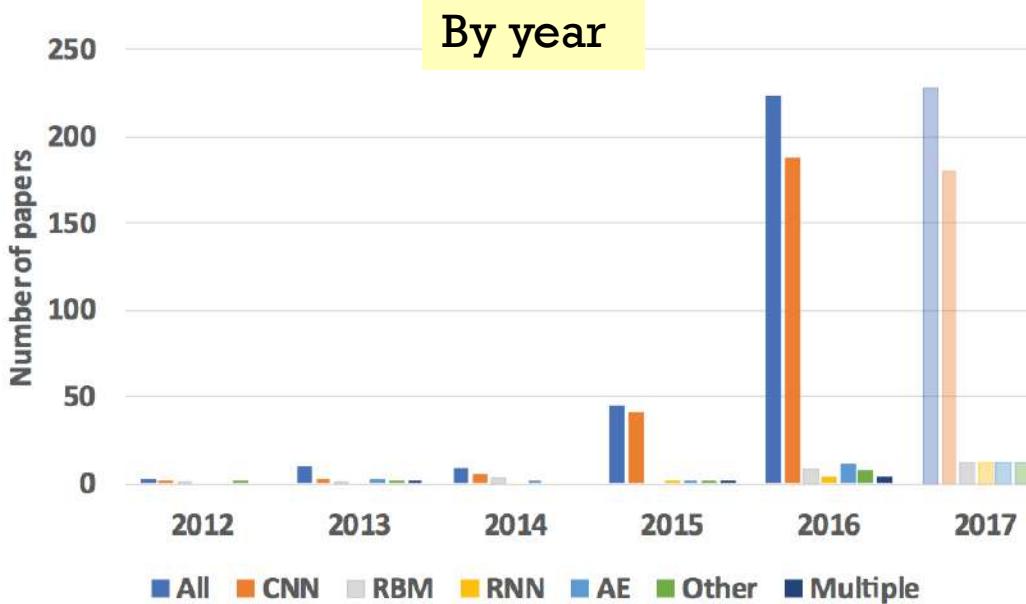
CAT, DOG, DUCK

Single object

Multiple objects

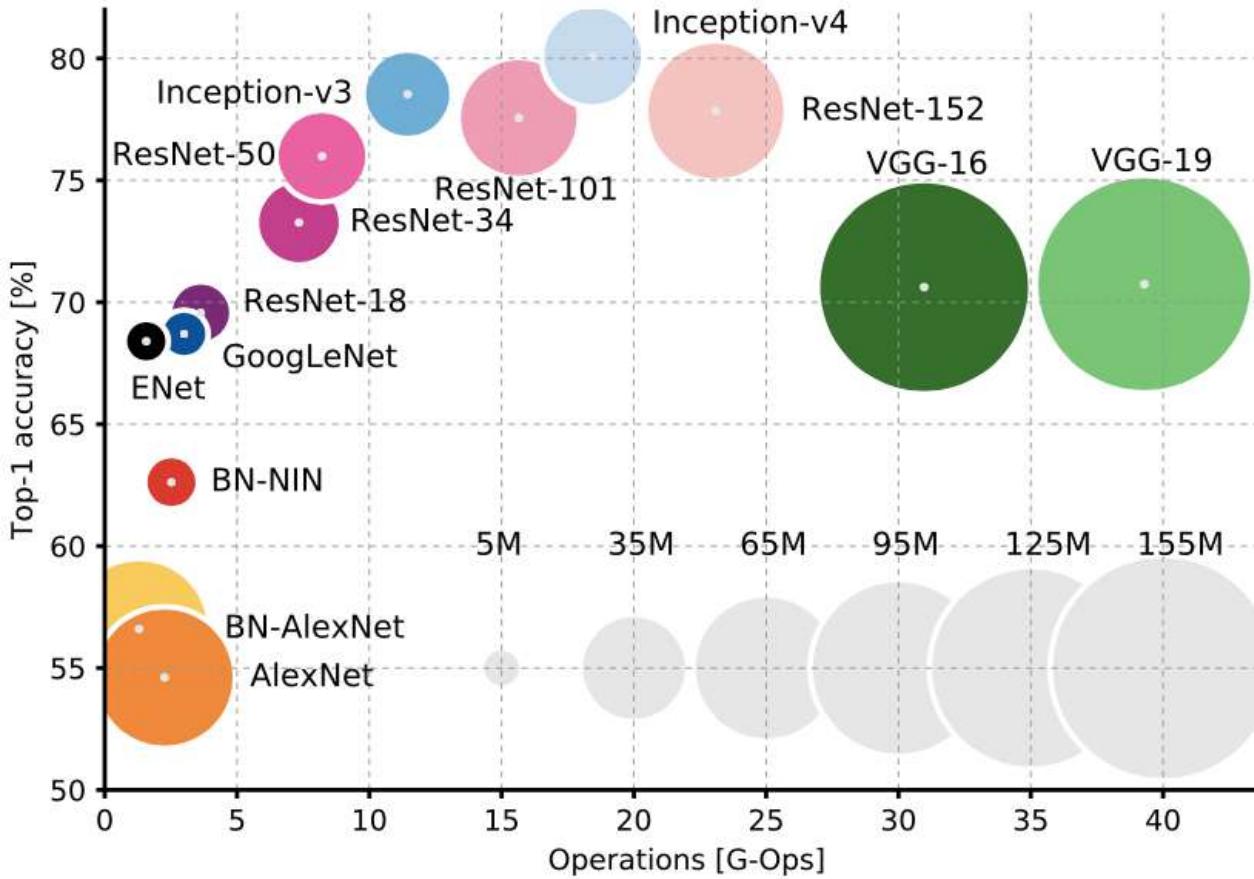


Deep Learning in Medical Image Analysis [arXiv 2017]

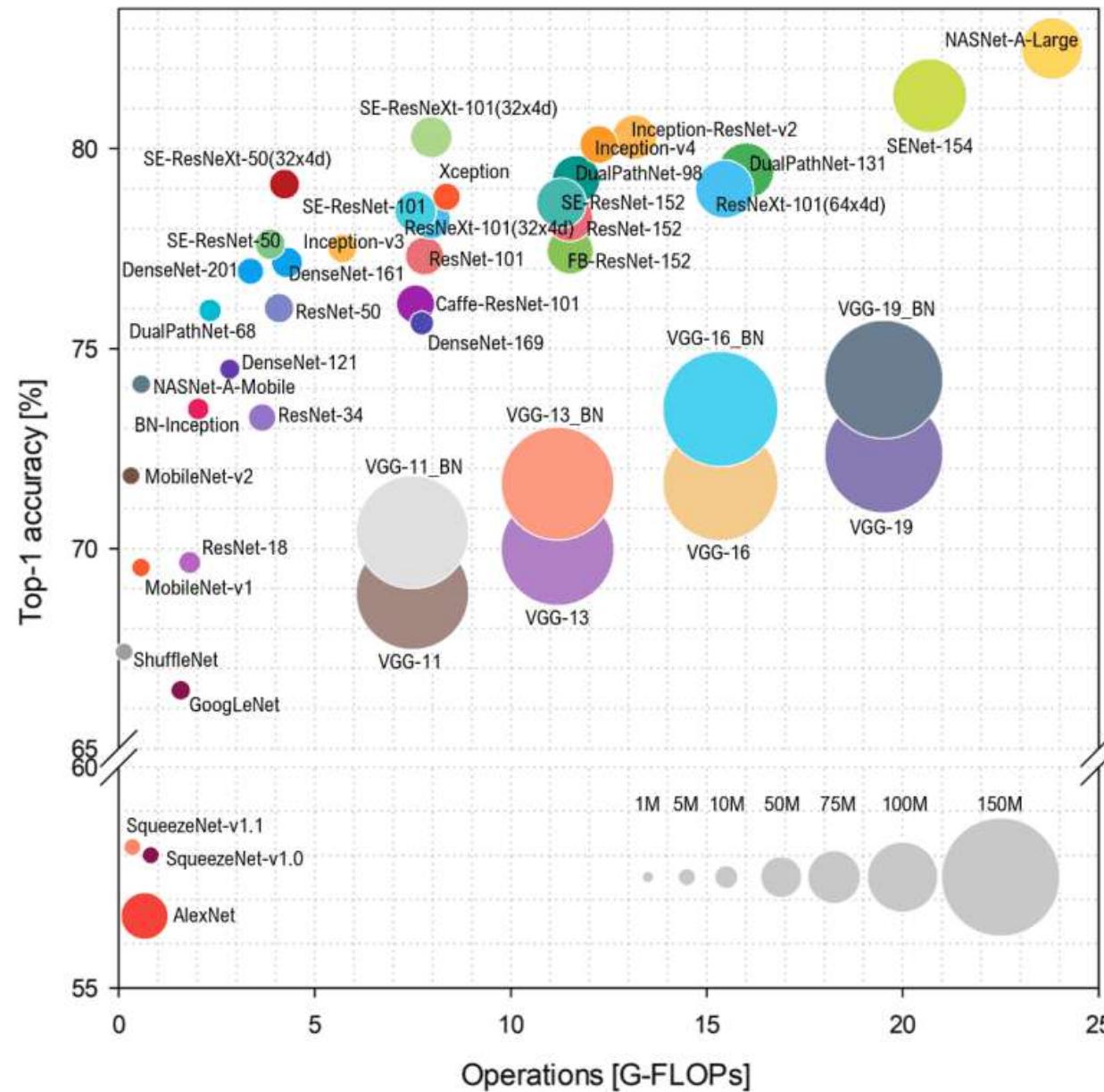




SOTA of Image Classification



https://blog.csdn.net/qq_34216467/article/details/83061692



<https://theaisummer.com/cnn-architectures/>



EfficientNet (May, 2019) → EffNetV2 (2021)

EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks

Mingxing Tan¹ Quoc V. Le¹

Abstract

Convolutional Neural Networks (ConvNets) are commonly developed at a fixed resource budget, and then scaled up for better accuracy if more resources are available. In this paper, we systematically study model scaling and identify that carefully balancing network depth, width, and resolution can lead to better performance. Based on this observation, we propose a new scaling method that uniformly scales all dimensions of depth/width/resolution using a simple yet highly effective *compound coefficient*. We demonstrate the effectiveness of this method on scaling up MobileNets and ResNet.

To go even further, we use neural architecture

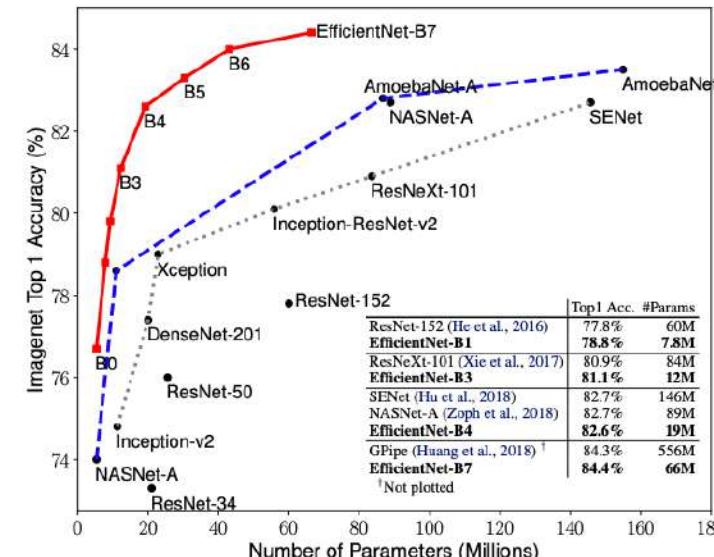
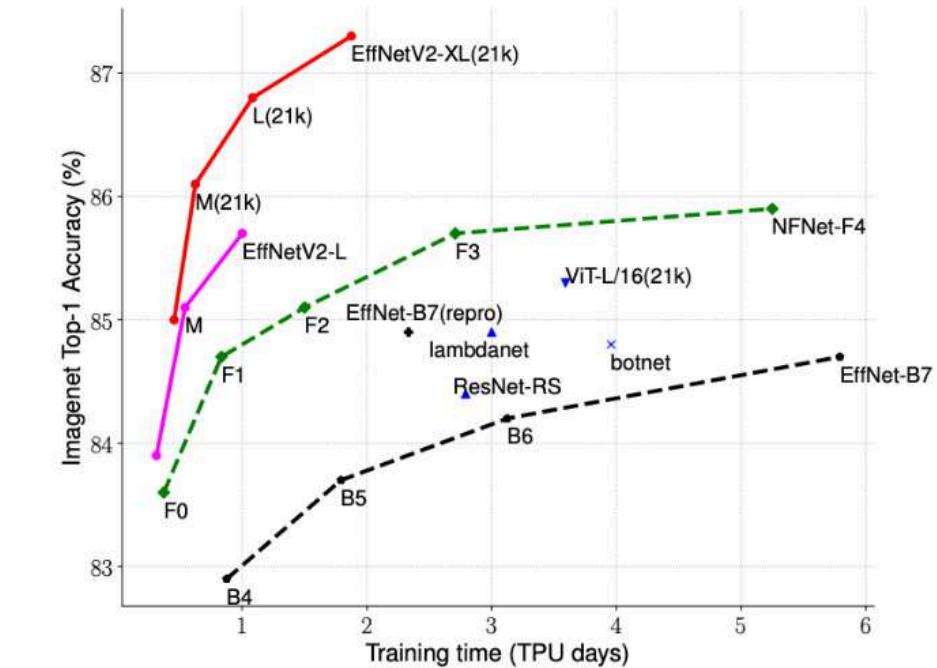


Figure 1. Model Size vs. ImageNet Accuracy. All numbers are



(a) Training efficiency.

	EfficientNet (2019)	ResNet-RS (2021)	DeiT/ViT (2021)	EfficientNetV2 (ours)
Top-1 Acc.	84.3%	84.0%	83.1%	83.9%
Parameters	43M	164M	86M	24M

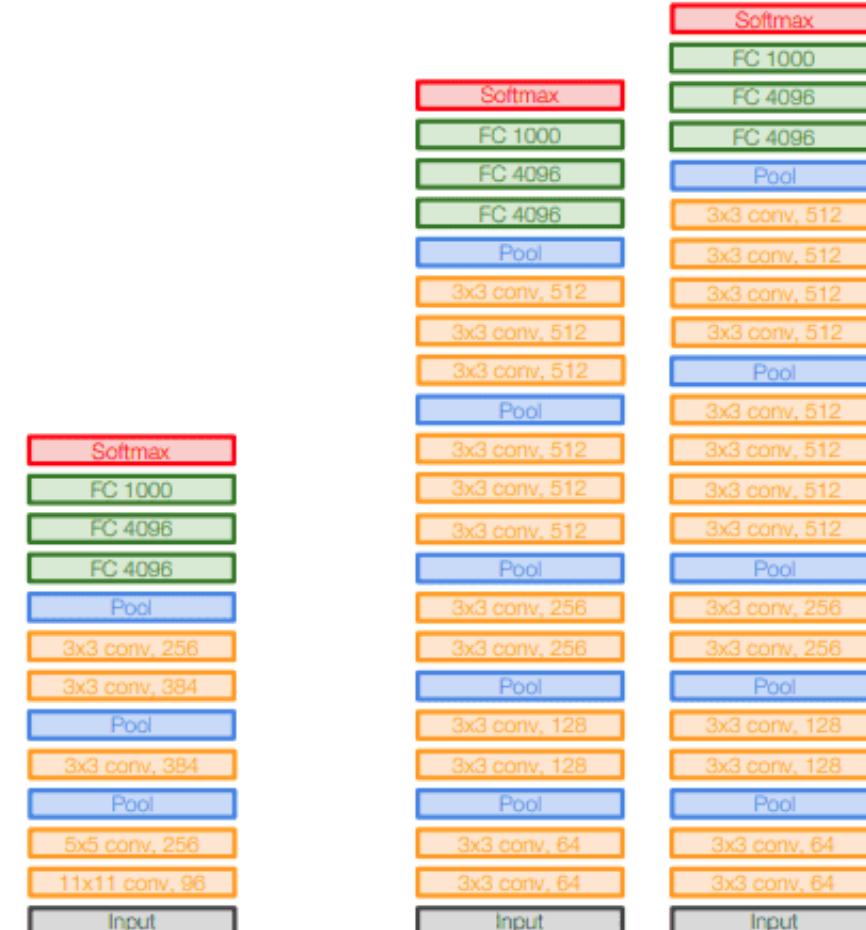
(b) Parameter efficiency.



Network series

1. VGG (2014)
2. Inception/GoogleNet (2014)
3. Inception V2, V3
4. ResNet
5. DenseNet
6. Inception-V4 & Inception-ResNet
7. EfficientNet V1, V2

1) VGG (2014)



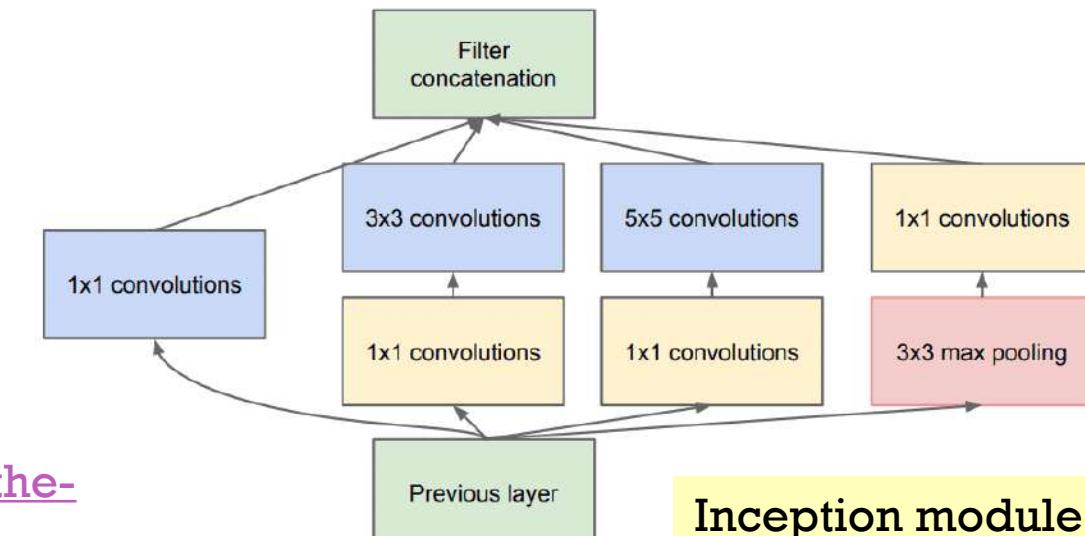
AlexNet

VGG16

VGG19

+ 2) Inception V1/GoogLeNet (2014)

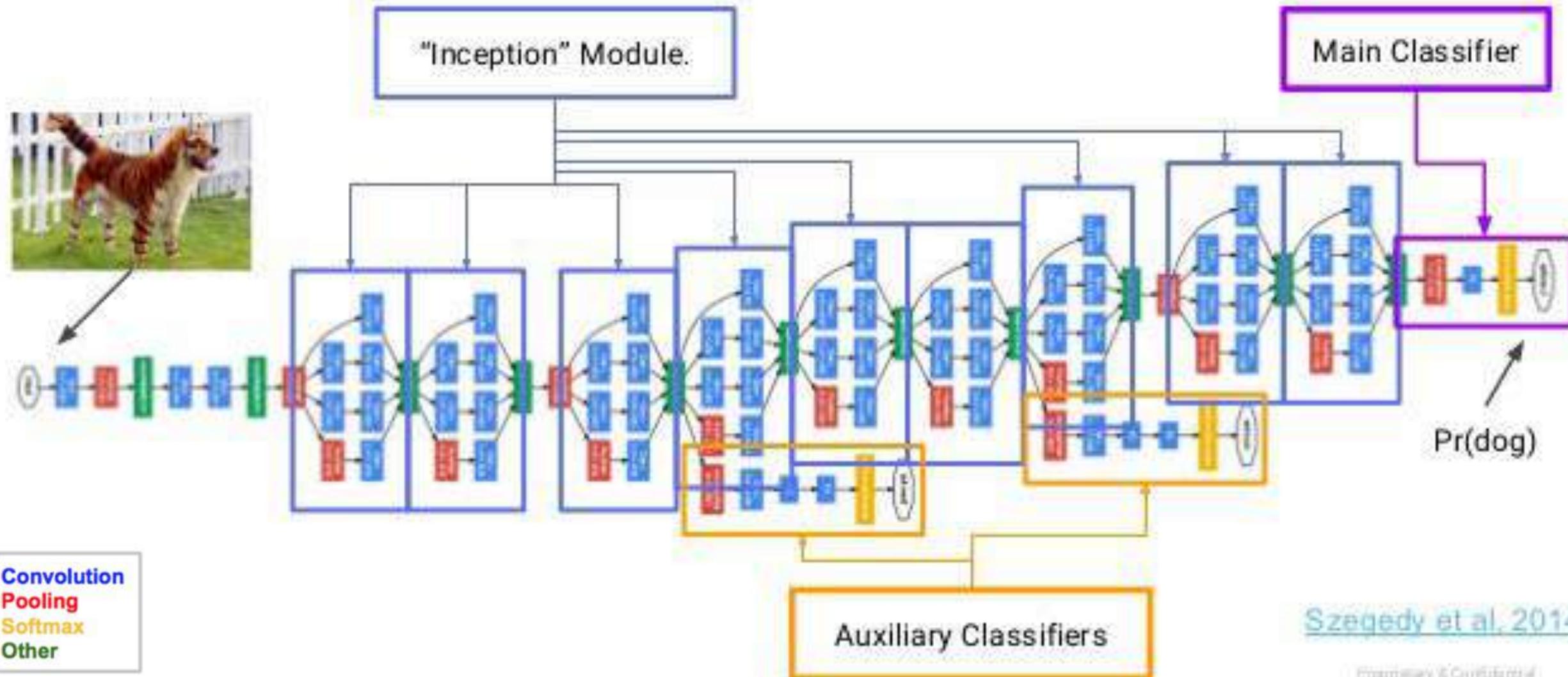
- After VGG, the paper “Going Deeper with Convolutions” [3] by Christian Szegedy et al. was a huge breakthrough. What about increasing both the depth (more layers) and width (more feature maps) of the network while keeping computations to a constant level?
- The answer is with **1×1 convolutions**! The main purpose is dimension reduction, by reducing the output channels of each convolution block.
 - 1×1 convolutions are used to compute reductions before the computationally expensive convolutions (3×3 and 5×5).
- Moreover, it uses convolutions of different kernel sizes (5×5, 3×3, 1×1) to capture details at multiple scales.
- The InceptionNet/GoogLeNet architecture consists of **9 inception modules stacked** together, with **max-pooling layers between** (to halve the spatial dimensions). It consists of 22 layers (**27 with the pooling layers**). It uses global average pooling after the last inception module.



GoogLeNet (aka “Inception”) Architecture



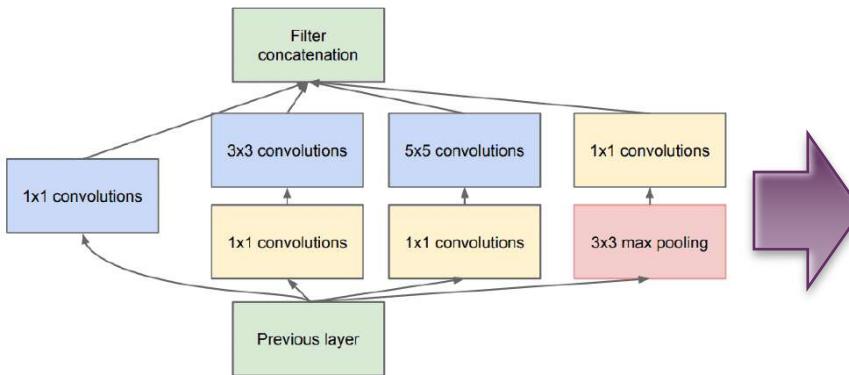
Inception is a deep network, to prevent the **middle part** of the network from “**dying out**”(vanishing gradient problem), the authors introduced **two auxiliary classifiers**. Softmax is applied in each of them and then **Auxiliary loss** is calculated on the same labels of the output classifier.



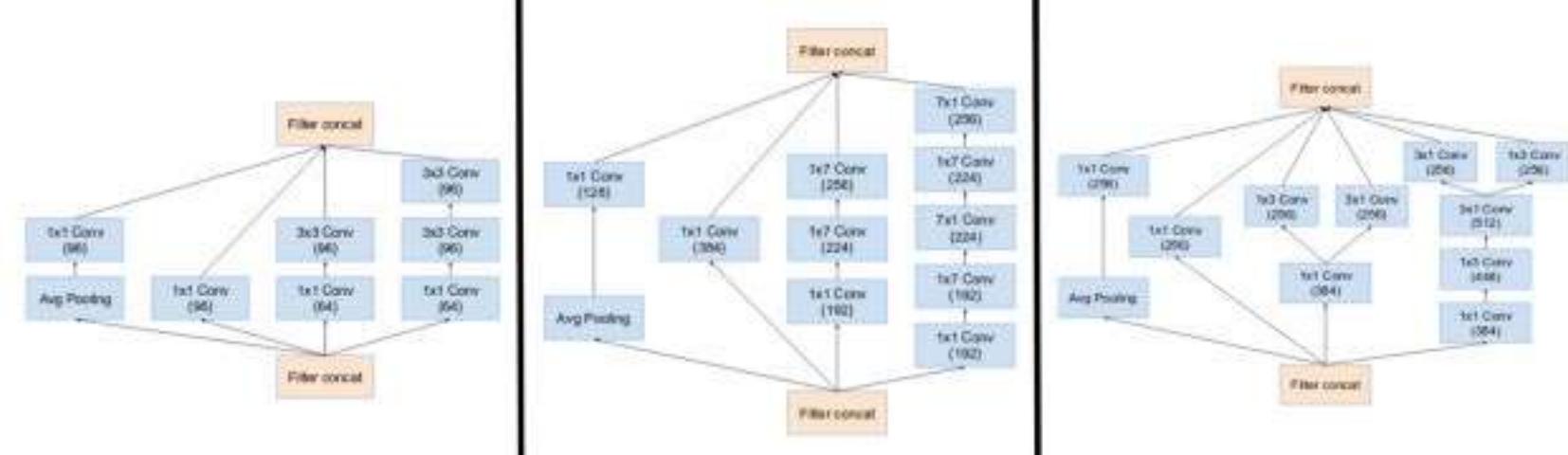


3) Inception V2, V3 (2015) – improve speed

- Factorize 5x5 and 7x7 (in InceptionV3) convolutions to two and three 3x3 sequential convolutions respectively. This improves computational speed. This is the same principle as VGG.
(analogy: $5 \times 5 = 25$ parameters vs. $2 \times (3 \times 3) = 18$ parameters)
- Simply, a 3x3 kernel is **decomposed** into two smaller ones: a 1x3 and a 3x1 kernel, which are applied sequentially. **(analogy: $3 \times 3 = 9$ parameters vs. $(1 \times 3) + (3 \times 1) = 6$ parameters)**
- The inception modules became **wider** (more feature maps).
- They added **batch normalization**.



Inception-V1 module



Inception Module (A, B, C)
in Inception V2, V3



Batch normalization

$$\hat{x} = \frac{x - \text{mean}(\bar{x})}{\text{std}(\bar{x})}$$

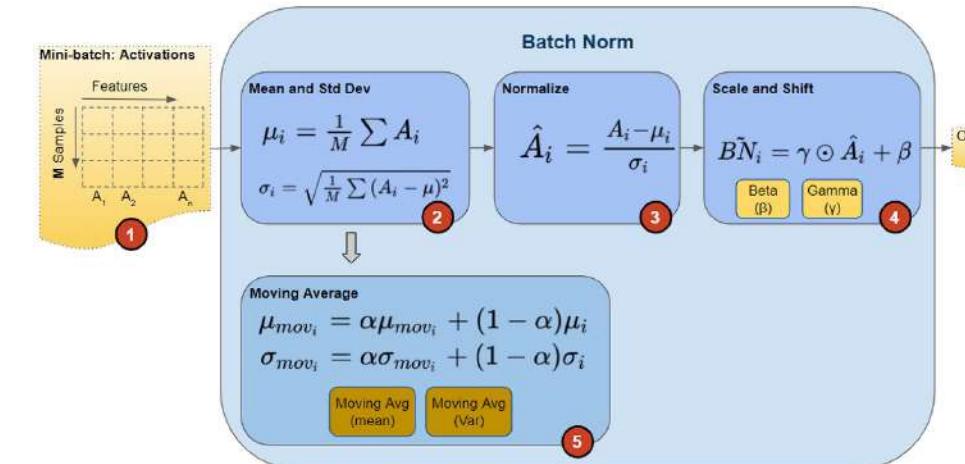
- Batch Normalization focuses on standardizing the inputs to **any particular layer** (i.e., activations from previous layers).

- Advantages of Batch Normalization Layer**

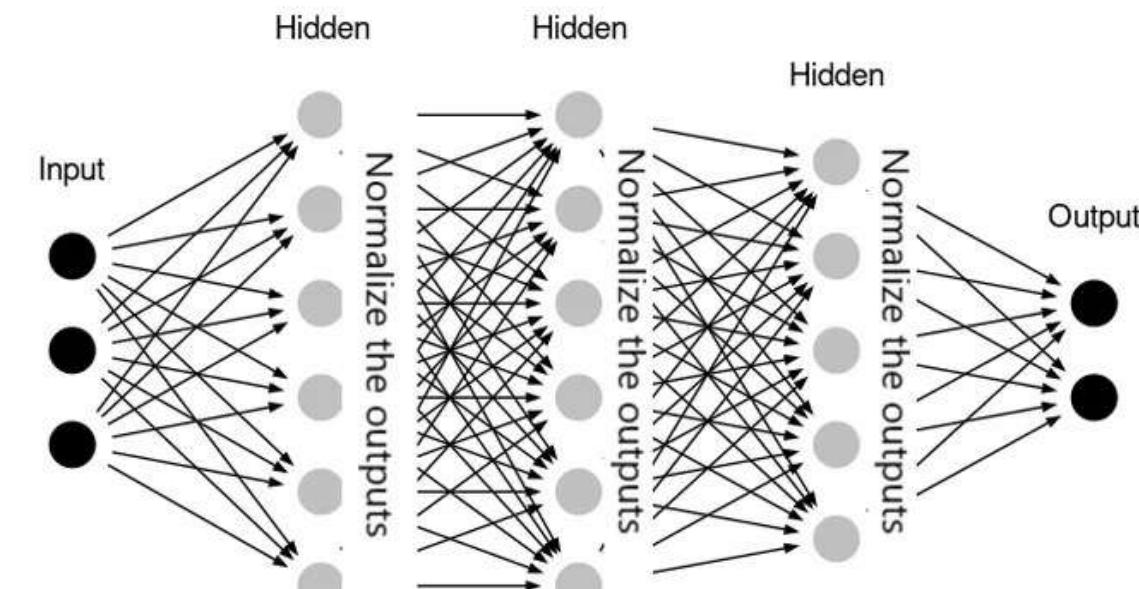
- Batch normalization improves the training time and accuracy of the neural network.
- It decreases the effect of weight initialization.
- It also adds a regularization effect on the network.
- It works better with the fully Connected Neural Network (FCN) and Convolutional Neural Network.

- Disadvantages of Batch Normalization Layer**

- Batch normalization is **dependent on mini-batch size** which means if the mini-batch size is small, it will have little to no effect
- Batch normalization does not work well with Recurrent Neural Networks (RNN)



<https://towardsdatascience.com/batch-norm-explained-visually-how-it-works-and-why-neural-networks-need-it-b18919692739>



Layer normalization

- Layer Normalization which addresses the drawbacks of batch normalization. This technique is not dependent on batches and the normalization is applied on the neuron for a single instance across all features.

■ Advantages of Layer Normalization

- It is not dependent on any batch sizes during training.

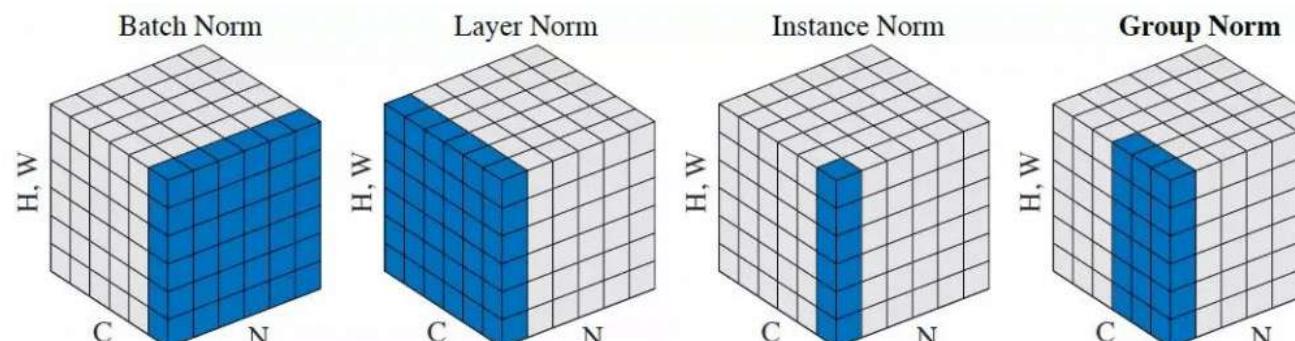
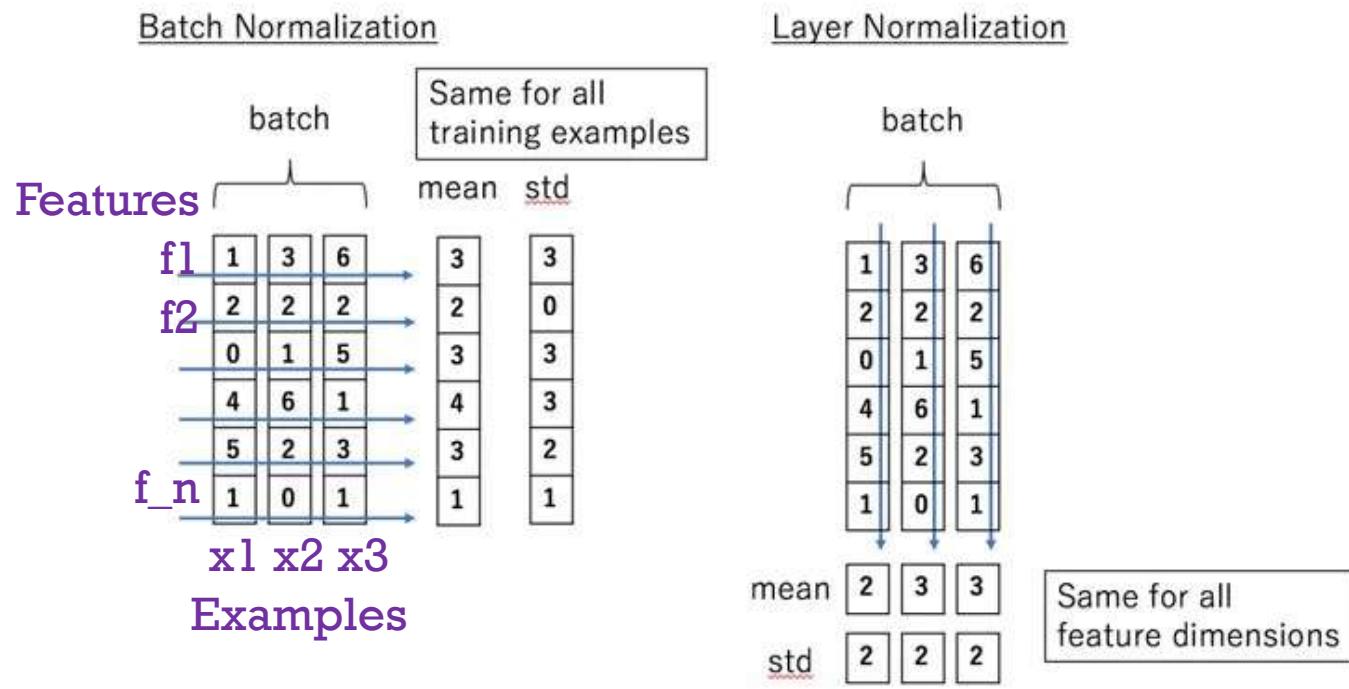
- It works better with Recurrent Neural Network.

■ Disadvantages of Layer Normalization

- It may not produce good results with Convolutional Neural Networks (CNN)

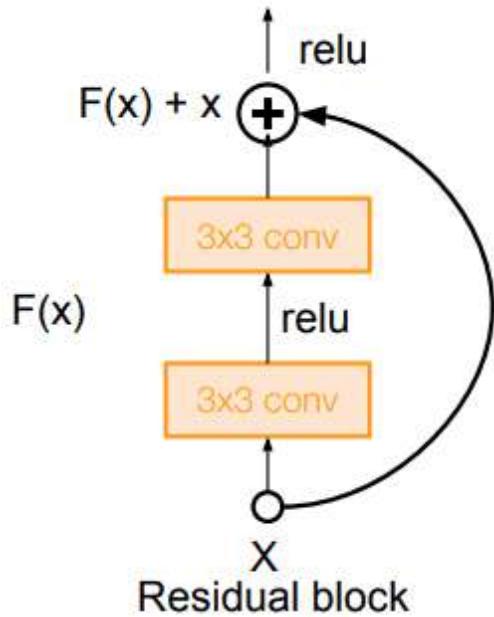
C = features (channels)

N = each mini-batch



+ 4) ResNet: Deep Residual Learning for Image Recognition (2015)

- All the pre-described issues such as vanishing gradients were addressed with two tricks:
 - **batch normalization** and
 - **short skip connections** - Instead of $H(x) = F(x)$, we ask them model to learn the difference (**residual**) $H'(x) = F(x) + x$, which means $H(x) - x = F(x)$ will be the residual part [4].
- Designed deeper architectures ranging from 18 (Resnet-18) to 150 (Resnet-150) layers.
- For the deepest models they adopted **1x1 convs**, as illustrated on the right:



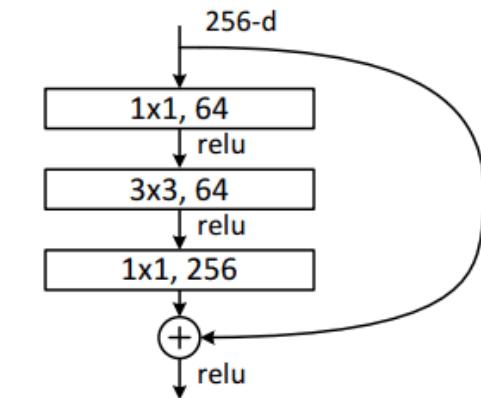
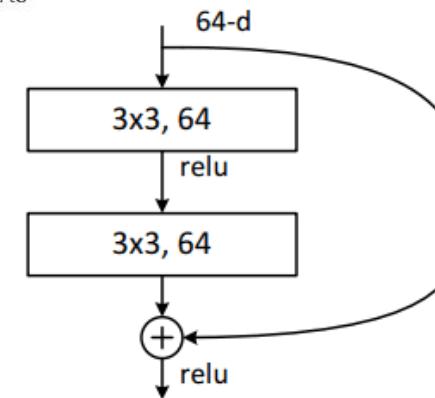
the layer followed by adding a bias term.

Then comes the activation function, $f()$ and we get the output as $H(x)$.

$H(x)=f(wx+b)$ or $H(x)=f(x)$

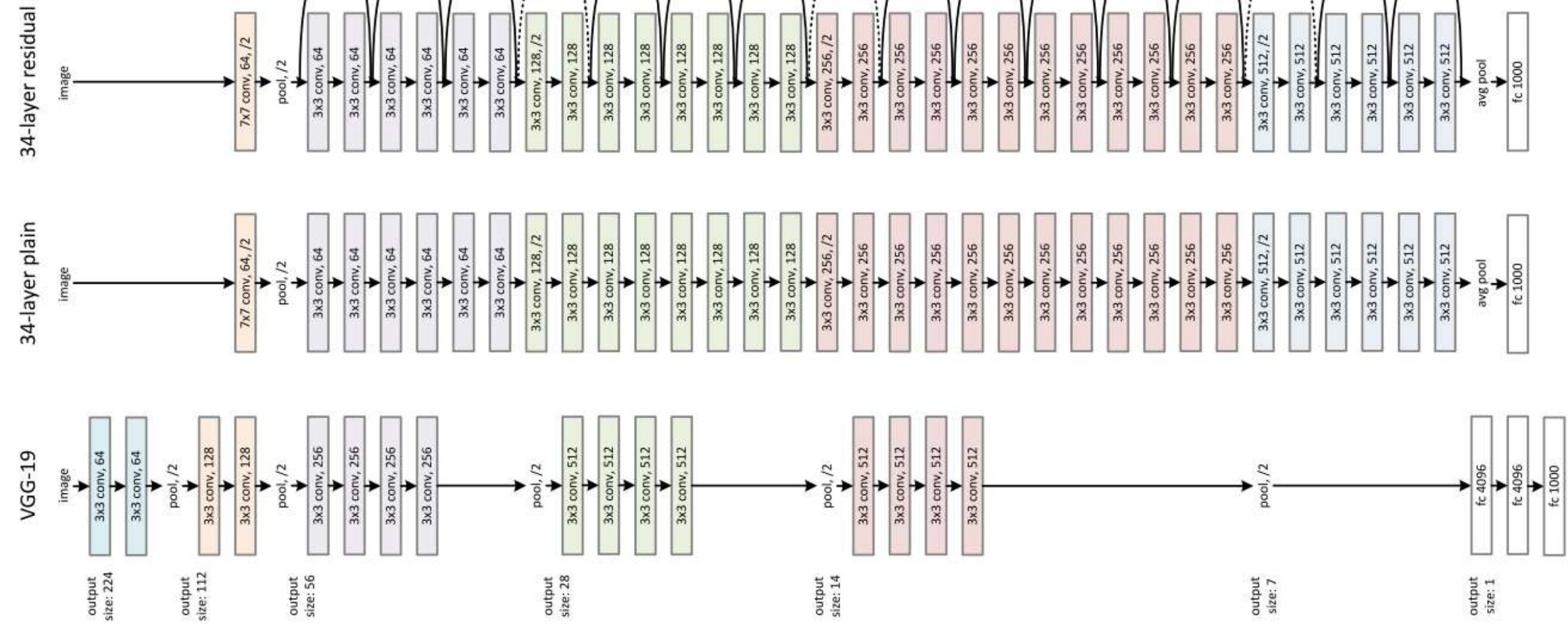
Now with the introduction of a new skip connection technique, the output is $H(x)$ is changed to

$H(x)=f(x)+x$





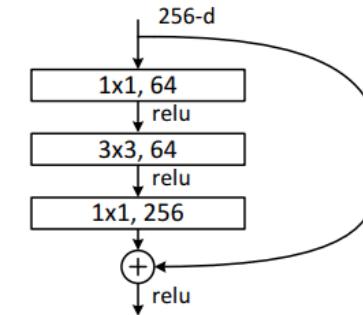
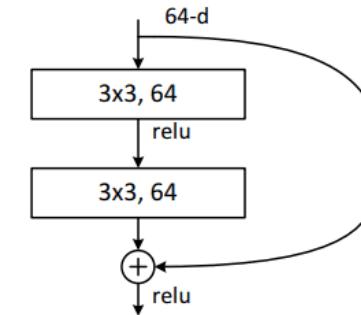
ResNet-34 vs. Plain Network (no skip connection)





ResNet (cont.)

ResNet-18, 34, 50, 101, 152



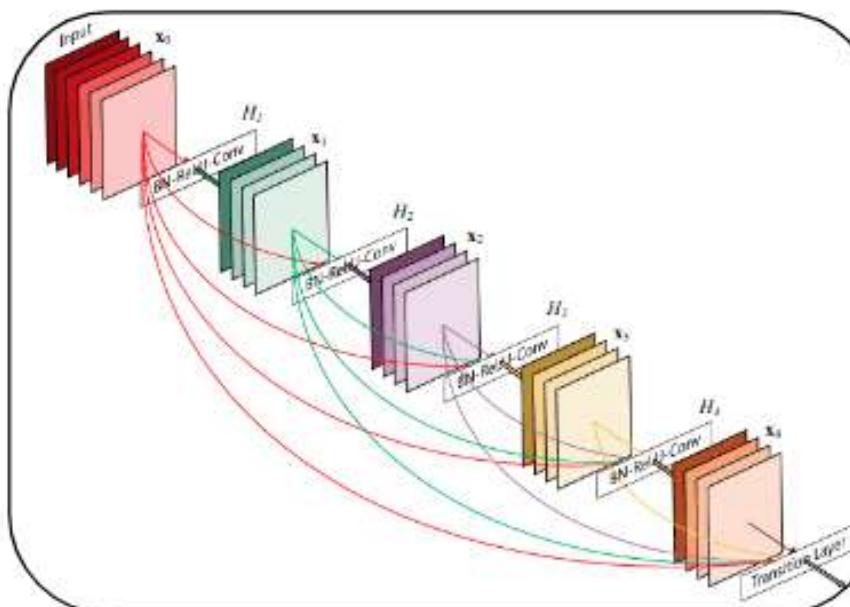
layer name	output size	18-layer	34-layer	50-layer	101-layer	152-layer
conv1	112×112			7×7, 64, stride 2		
conv2_x	56×56	$\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$
conv3_x	28×28	$\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 4$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 4$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 4$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 8$
conv4_x	14×14	$\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 6$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 6$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 23$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 36$
conv5_x	7×7	$\begin{bmatrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$
	1×1			average pool, 1000-d fc, softmax		
FLOPs		1.8×10^9	3.6×10^9	3.8×10^9	7.6×10^9	11.3×10^9



5) DenseNet: Densely Connected Convolutional Networks (2017)

- Skip connections are a pretty cool idea. Why don't we just skip-connect everything?

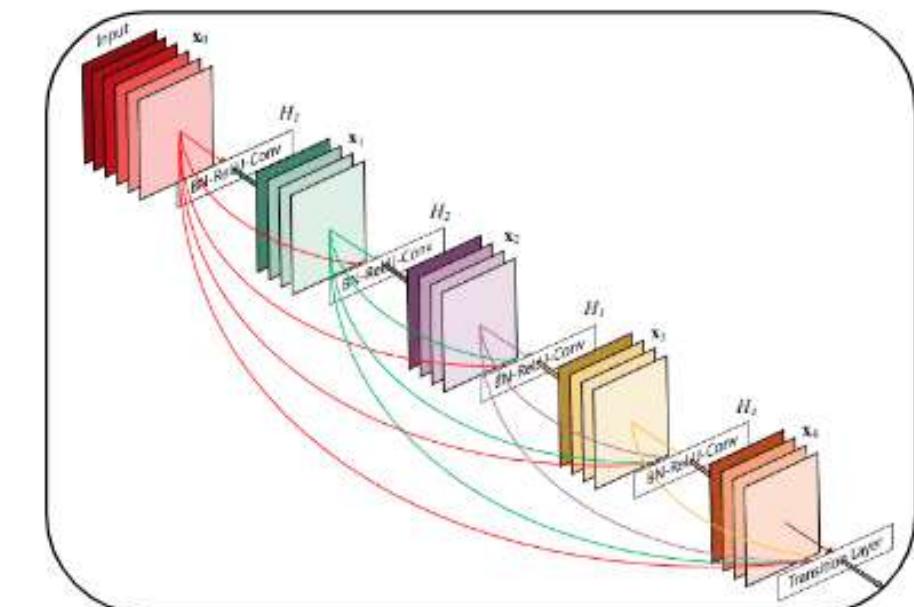
Dense Block 1



Transition Layer

$1 \times 1 \times K$ convs
+
 2×2 average
pooling

Dense Block 2



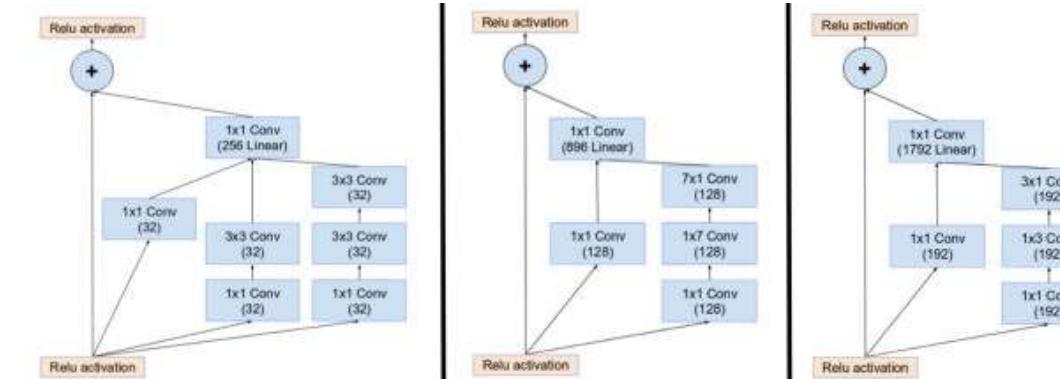
Spatial dims: 256x256
512 feature maps

Spatial dims: 128x128
K feature maps < 512



6) Inception-V4 & Inception-ResNet

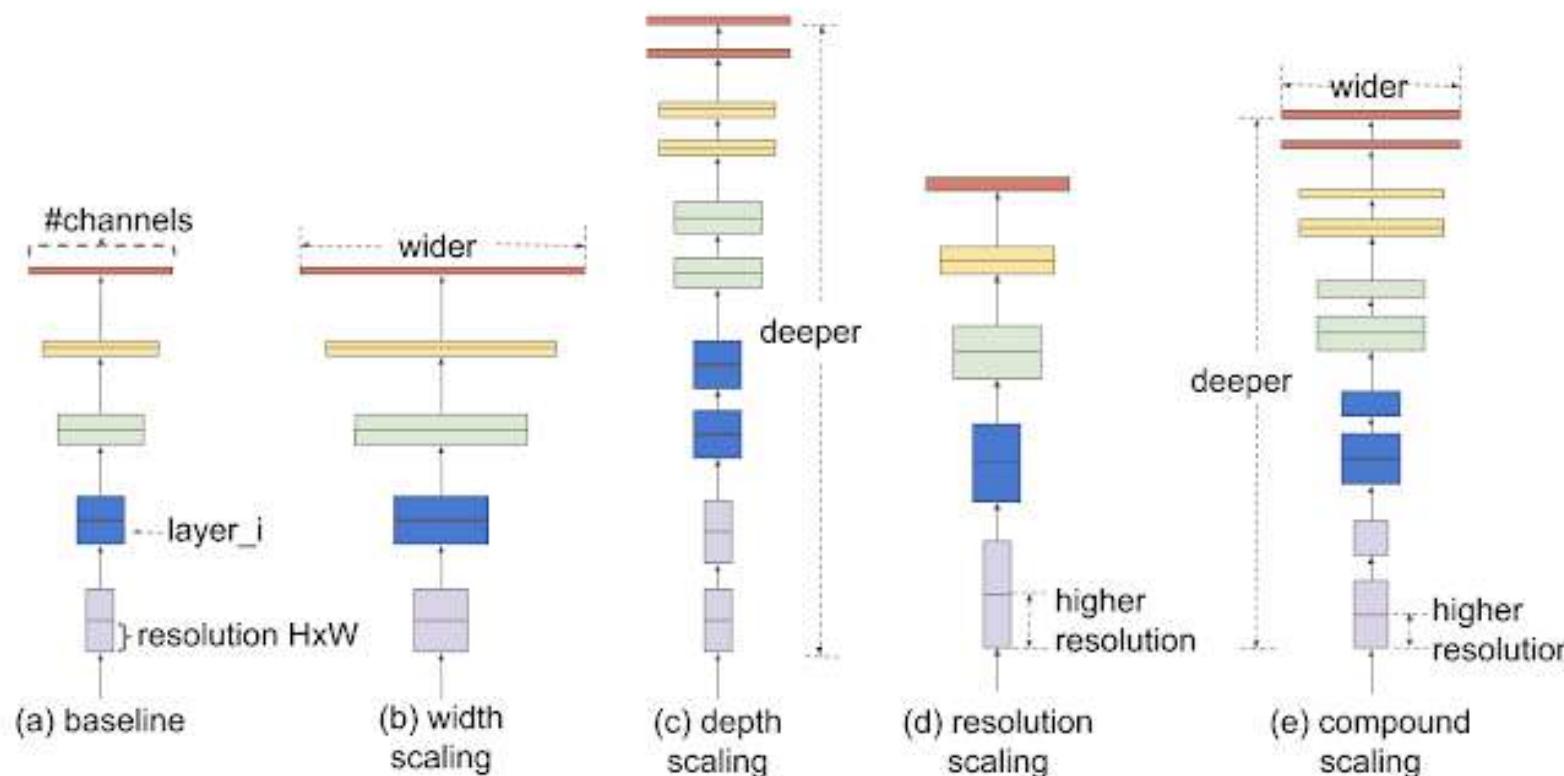
- Inception v4 and Inception-ResNet were introduced in the same paper.
- Inception-V4 (simplified previous version)
 - Make the modules more uniform. The authors also noticed that some of the modules were more complicated than necessary. **Thus; the “stem” of Inception v4 was modified (simplified).**
 - Inception v4 introduced specialized “**Reduction Blocks**” (**simplified**) which are used to change the width and height of the grid.
- Inception-ResNet V1, V2
 - For **residual addition** to work, the input and output after convolution must have the same dimensions. Hence, we use 1x1 convolutions after the original convolutions, to match the depth sizes (Depth is increased after convolution).
 - Inception-ResNet v1 has a computational cost that is similar to that of Inception v3.
 - Inception-ResNet v2 has a computational cost that is similar to that of Inception v4.





7) EfficientNet (May, 2019)

- So, let's instead scale up network **depth** (more layers), **width** (more feature maps), **resolution** (various sizes of filters) simultaneously. This is known as compound scaling.



For B0 to B7 base models, the input shapes are different. Here is a list of input shape expected for each model:

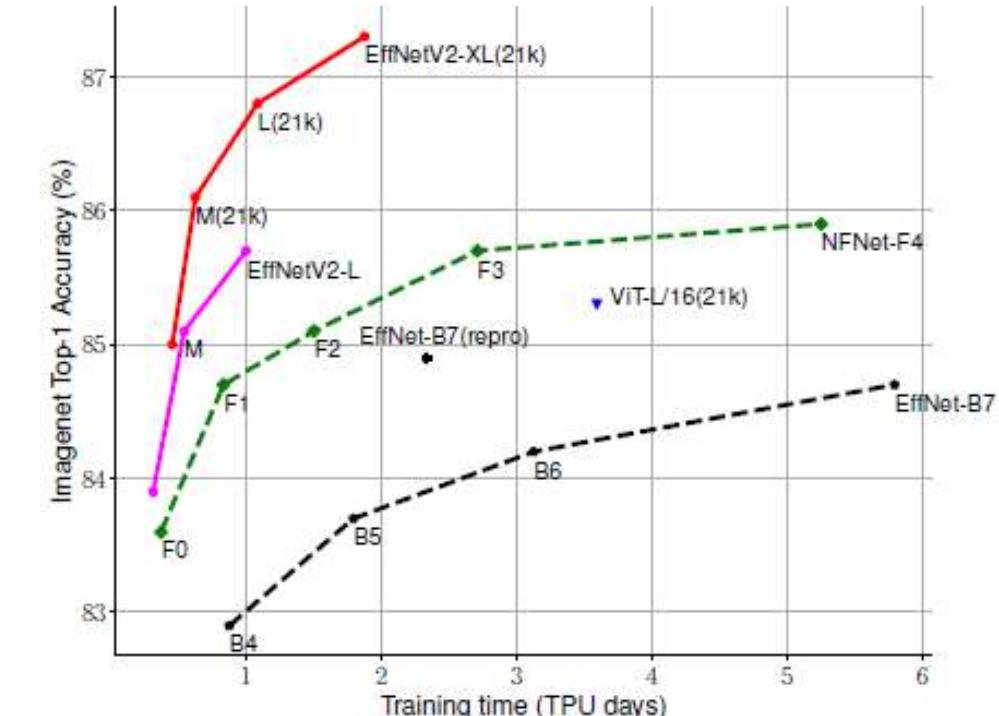
Base model	resolution
EfficientNetB0	224
EfficientNetB1	240
EfficientNetB2	260
EfficientNetB3	300
EfficientNetB4	380
EfficientNetB5	456
EfficientNetB6	528
EfficientNetB7	600



7) EfficientNetV2 (2021)

Smaller models & faster training

- Our experiments show that EfficientNetV2 models train much faster than state-of-the-art models while being up to **6.8x smaller**.
- There are many versions: small, medium, large, and extra-large (XL).

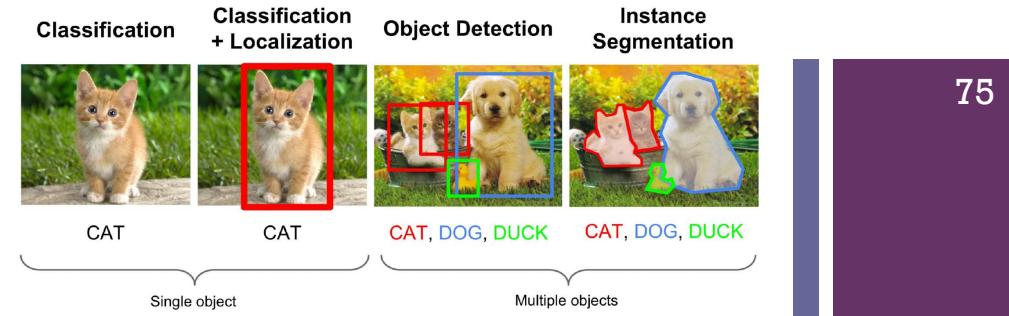


	EfficientNet (2019)	ResNet-RS (2021)	DeiT/ViT (2021)	EfficientNetV2 (ours)
Top-1 Acc.	84.3%	84.0%	83.1%	83.9%
Parameters	43M	164M	86M	24M

(b) Parameter efficiency.

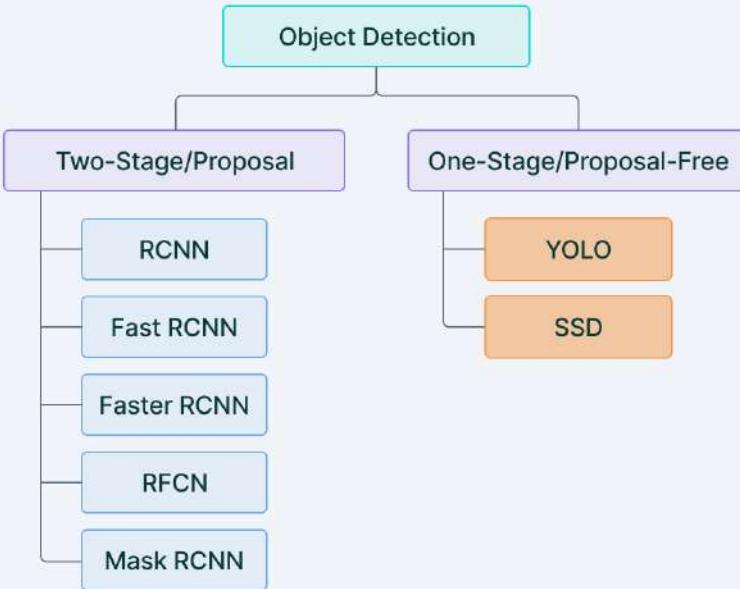


SOTA of Object Detection

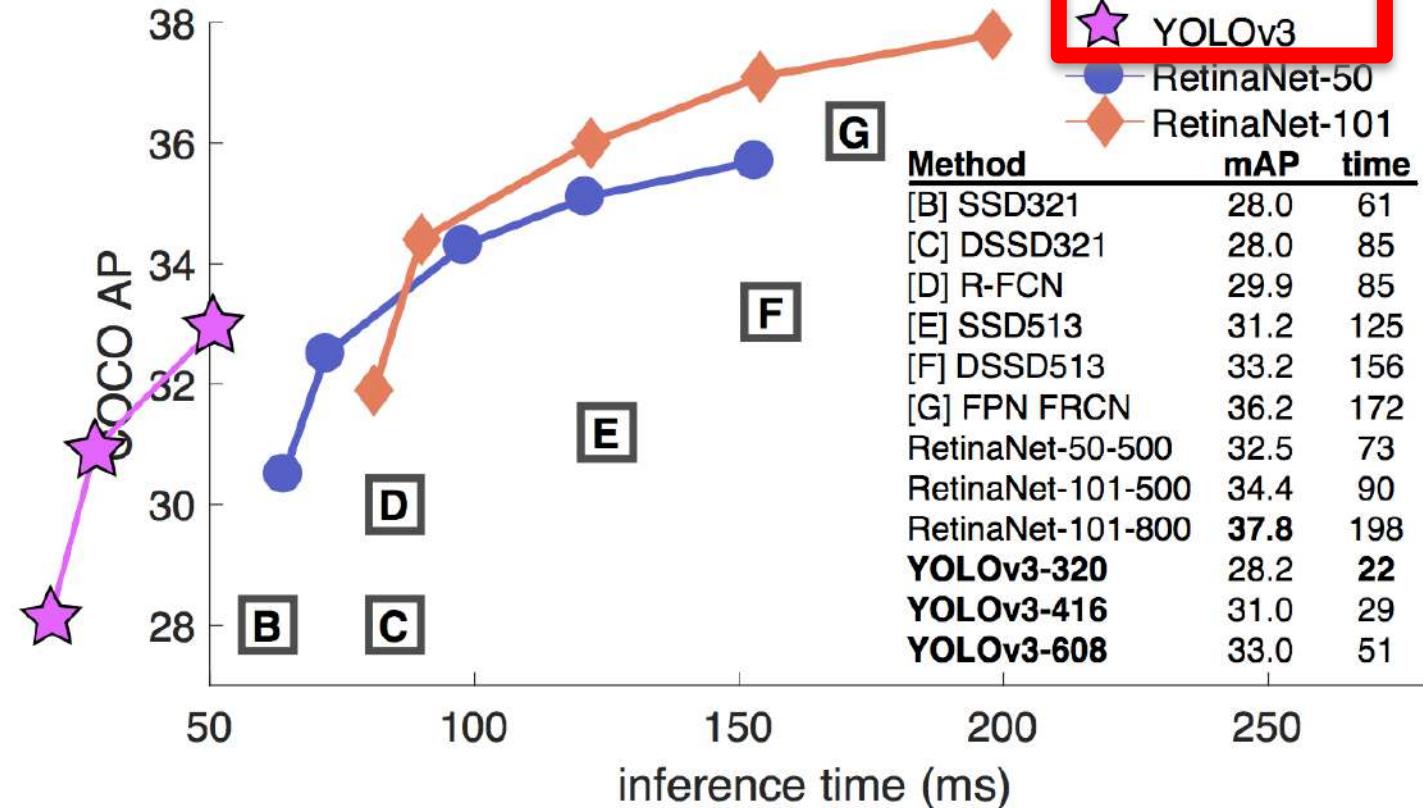


75

One and two stage detectors



V7 Labs



https://medium.com/@jonathan_hui/real-time-object-detection-with-yolo-yolov2-28b1b93e2088



YOLO Open Images in New York





YOLO's output

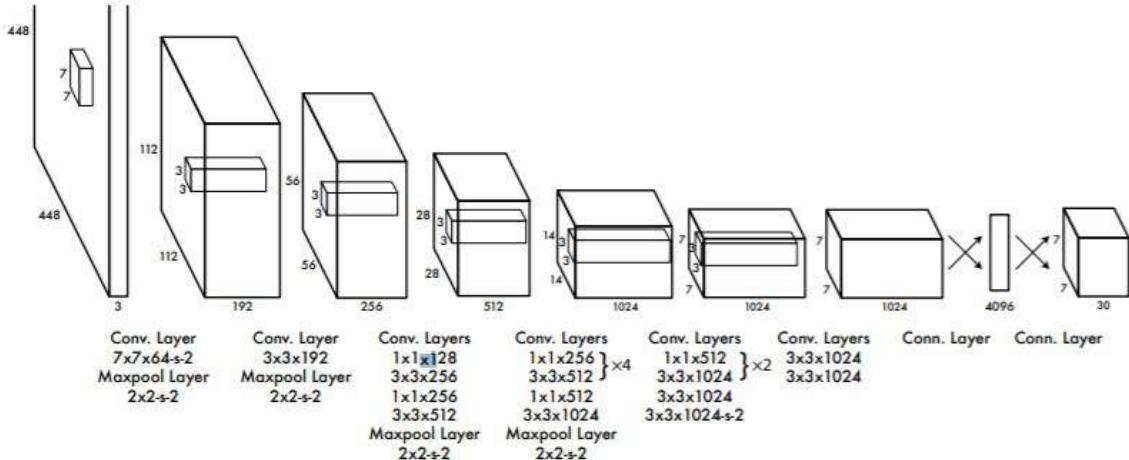
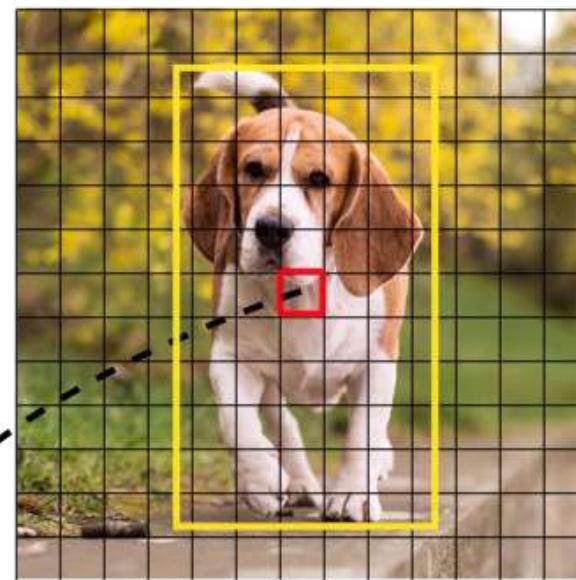
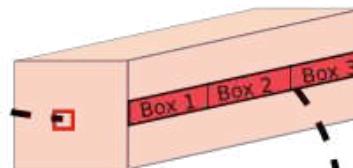


Figure 3: The Architecture. Our detection network has 24 convolutional layers followed by 2 fully connected layers. Alternating 1×1 convolutional layers reduce the features space from preceding layers. We pretrain the convolutional layers on the ImageNet classification task at half the resolution (224×224 input image) and then double the resolution for detection.

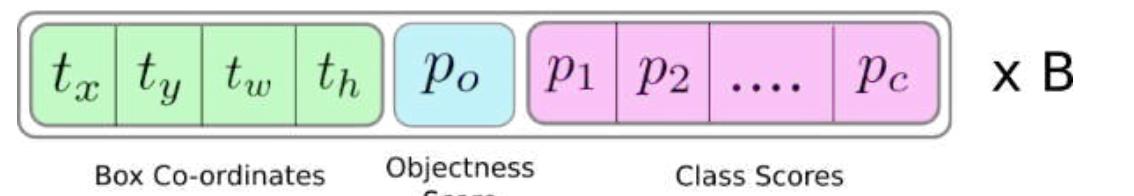
Image Grid. The Red Grid is responsible for detecting the dog



Prediction Feature Map



Attributes of a bounding box



<https://blog.paperspace.com/how-to-implement-a-yolo-object-detector-in-pytorch/>



SOTA of Object Detection (cont.)



Top YOLO Variants Of 2021



You can read more about YOLOP [here](#).

Rank	Model	IoU@AP	AP50	AP75	APS	APM	APL	Extra Training Data	Paper	Code	Result	Year	Tags
1	YOLOR-D6 (1280, single-scale, 30 fps)	57.3	75.0	62.7	40.4	61.2	69.2	✗	You Only Learn One Representation: Unified Network for Multiple Tasks			2021	
2	YOLOR-E6 (1280, single-scale, 37 fps)	56.4	74.1	61.6	39.1	60.1	68.2	✗	You Only Learn One Representation: Unified Network for Multiple Tasks			2021	
3	YOLOv4-P7 with TTA	55.8	73.2	61.2				✗	Scaled-YOLOv4: Scaling Cross Stage Partial Network			2020	
4	YOLOR-W6 (1280, single-scale, 47 fps)	55.5	73.2	60.6	37.6	59.5	67.7	✗	You Only Learn One Representation: Unified Network for Multiple Tasks			2021	
5	YOLOv4-P7 CSP-P7 (single-scale, 16 fps)	55.4	73.3	60.7	38.1	59.5	67.4	✗	Scaled-YOLOv4: Scaling Cross Stage Partial Network			2020	

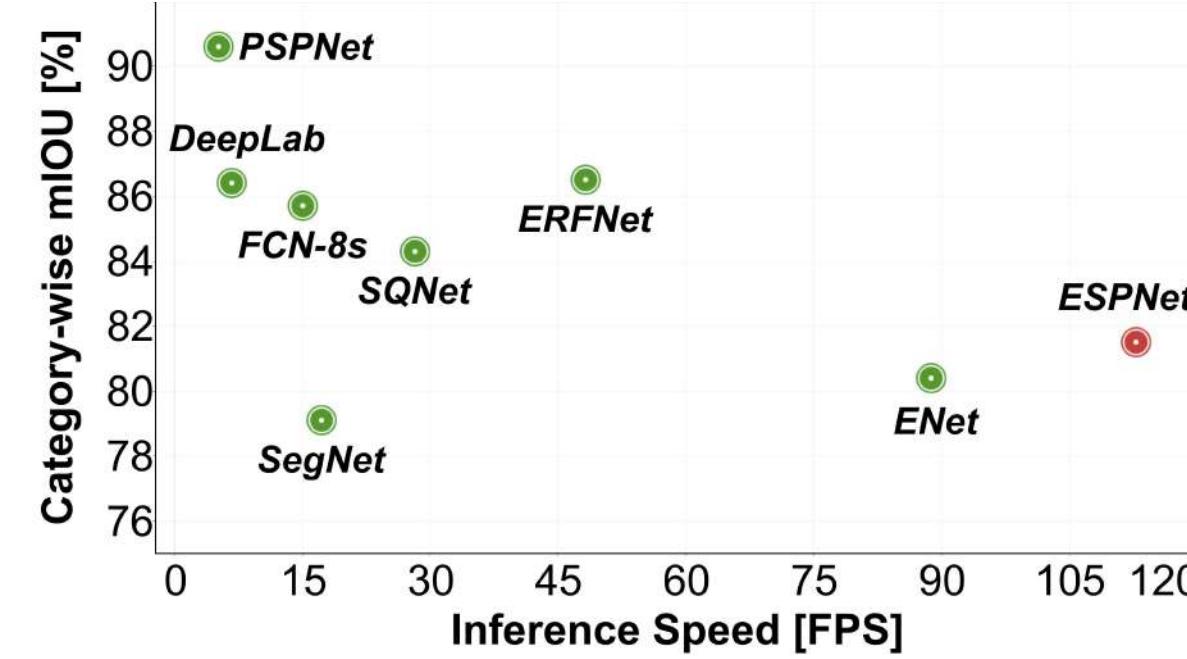
<https://medium.com/augmented-startups/top-yolo-variants-of-2021-19dddc23043c>



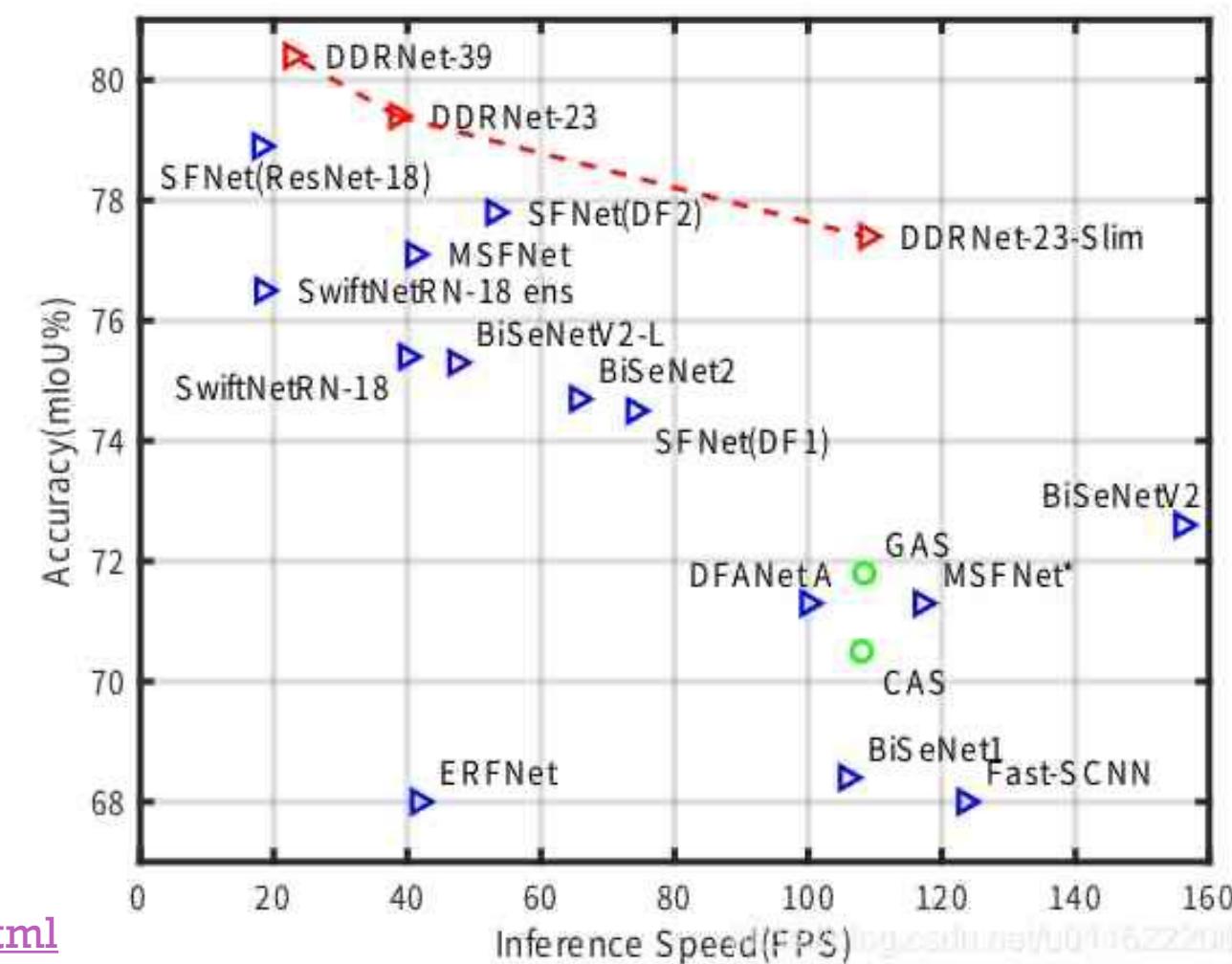
SOTA of Semantic Segmentation

[Personal open source] --

Real - time Semantic Segmentation ddrnet



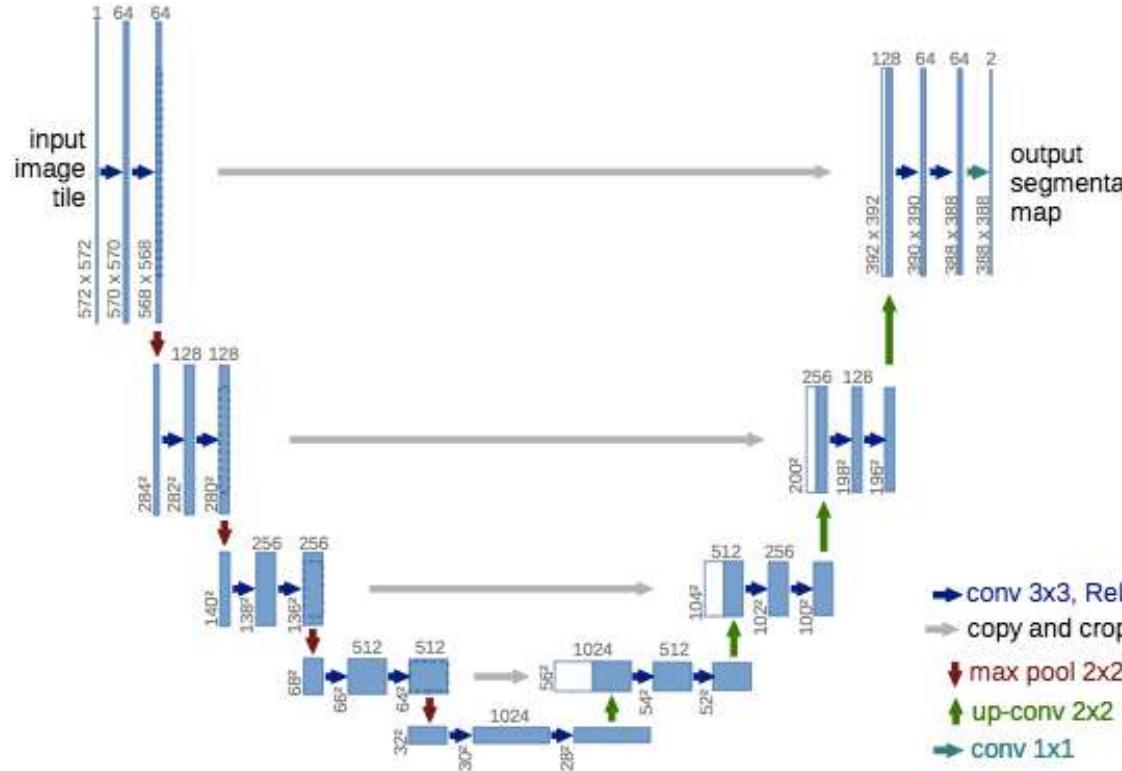
- **Real-time semantic segmentation**
- BiSeNet V1, V2
- DDRNet



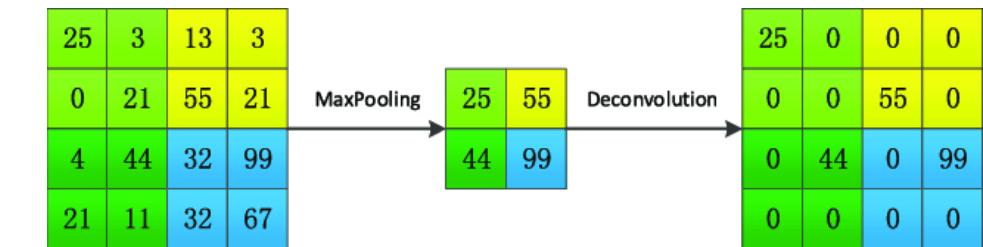


Review of Deep Learning Algorithms for Semantic Segmentation (cont.)

■ U-Net



Deconvolutional layer



Architecture of the U-net for a given input image. The blue boxes correspond to feature maps blocks with their denoted shapes. The white boxes correspond to the copied and cropped feature maps.

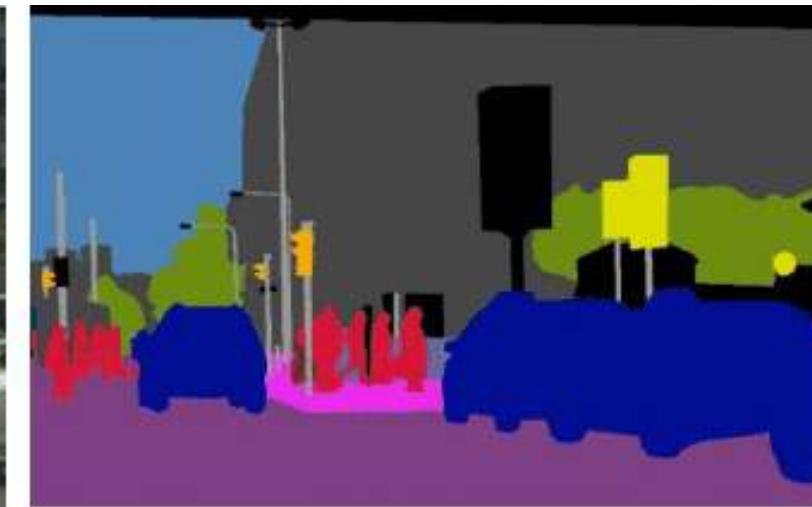
Source: [O. Ronneberger et al. \(2015\)](#)



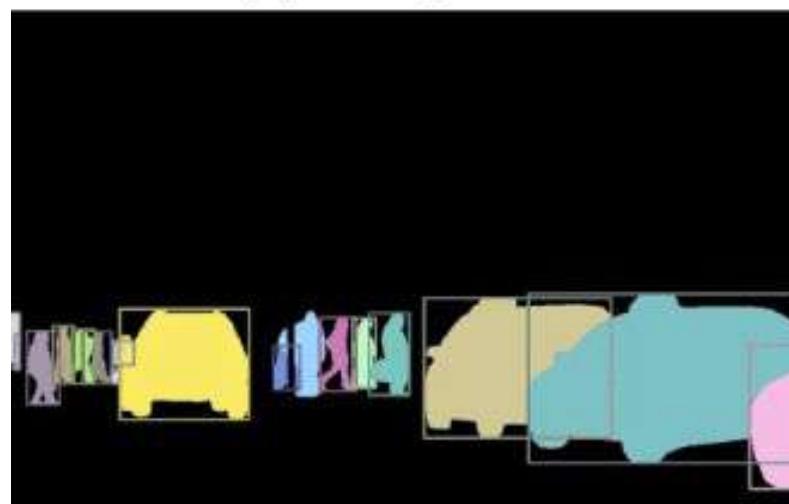
Panoptic Segmentation



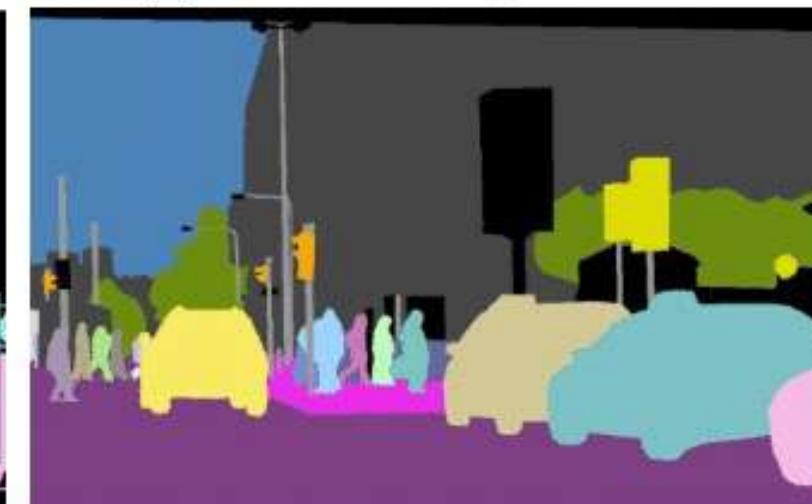
(a) Image



(b) Semantic Segmentation



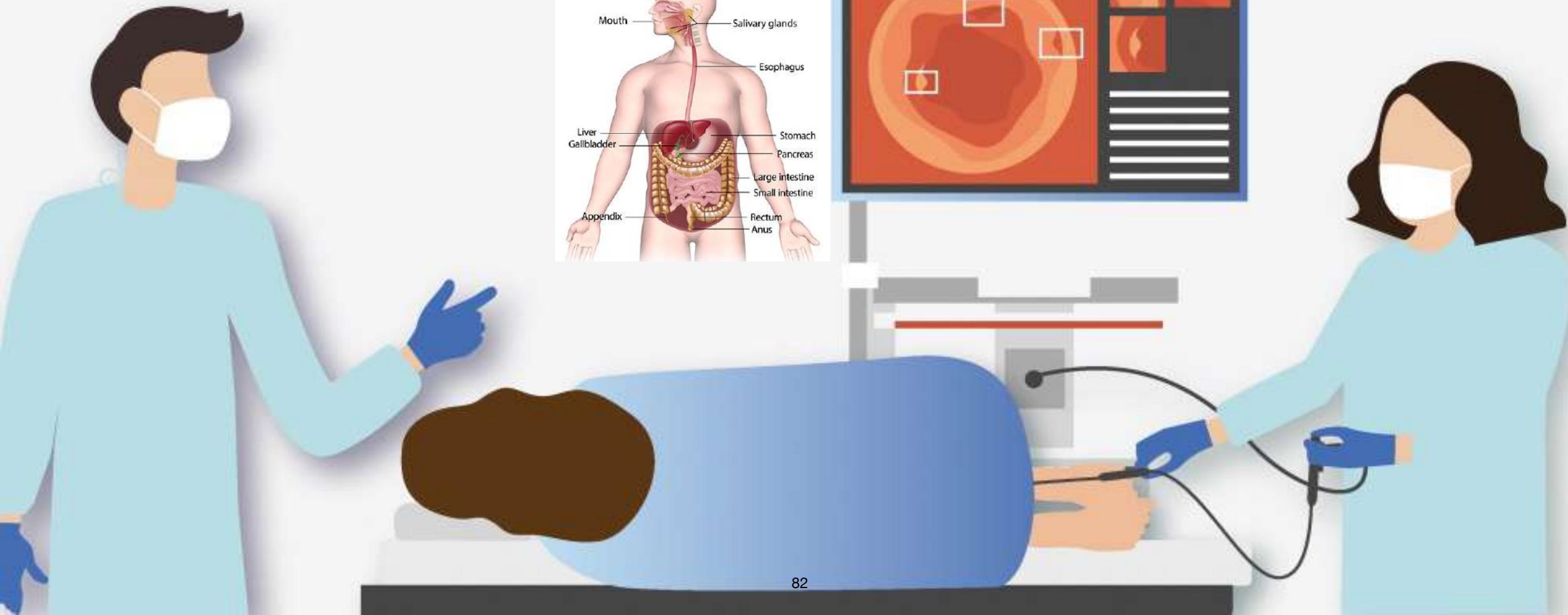
(c) Instance Segmentation



(d) Panoptic Segmentation

Real-Time Colonic Polyp Detection

An AI-assisted solution that aims to improve colonic polyp detection in real-time. It is compatible with all standard colonoscope systems.

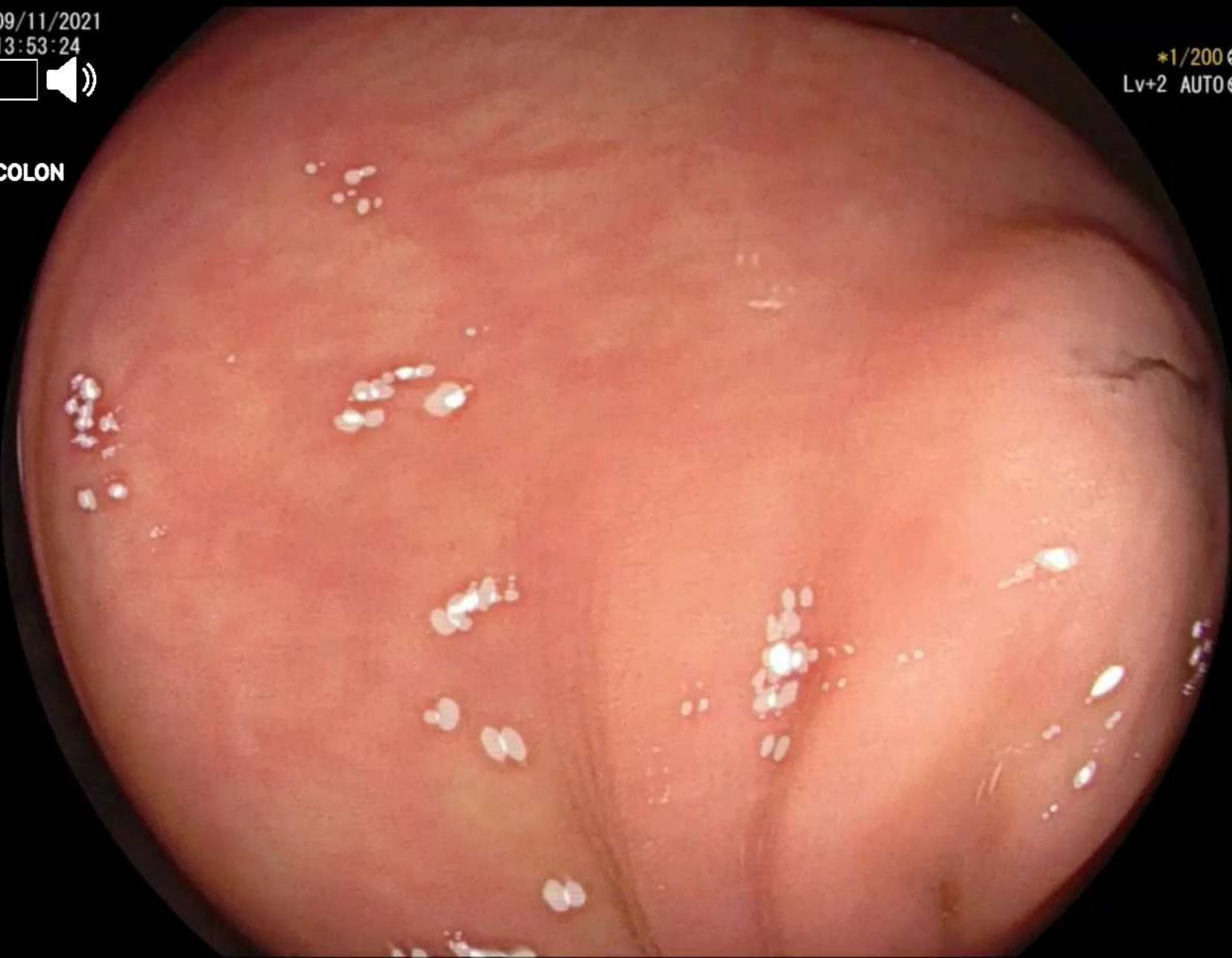


09/11/2021

13:53:24



COLON



*1/200
Lv+2 AUTO

HT NR
SE
f *
3.8 12.0 S1: F+T
12.0 S2: LM
S3: CAD
S4:
EC-760R-V/L
3C729K035
BL-7000

CHULALONGKORN HOS



NBI

threshold : 0.4
class : NBI

Name:

Sex: Age:

D.O.B.:

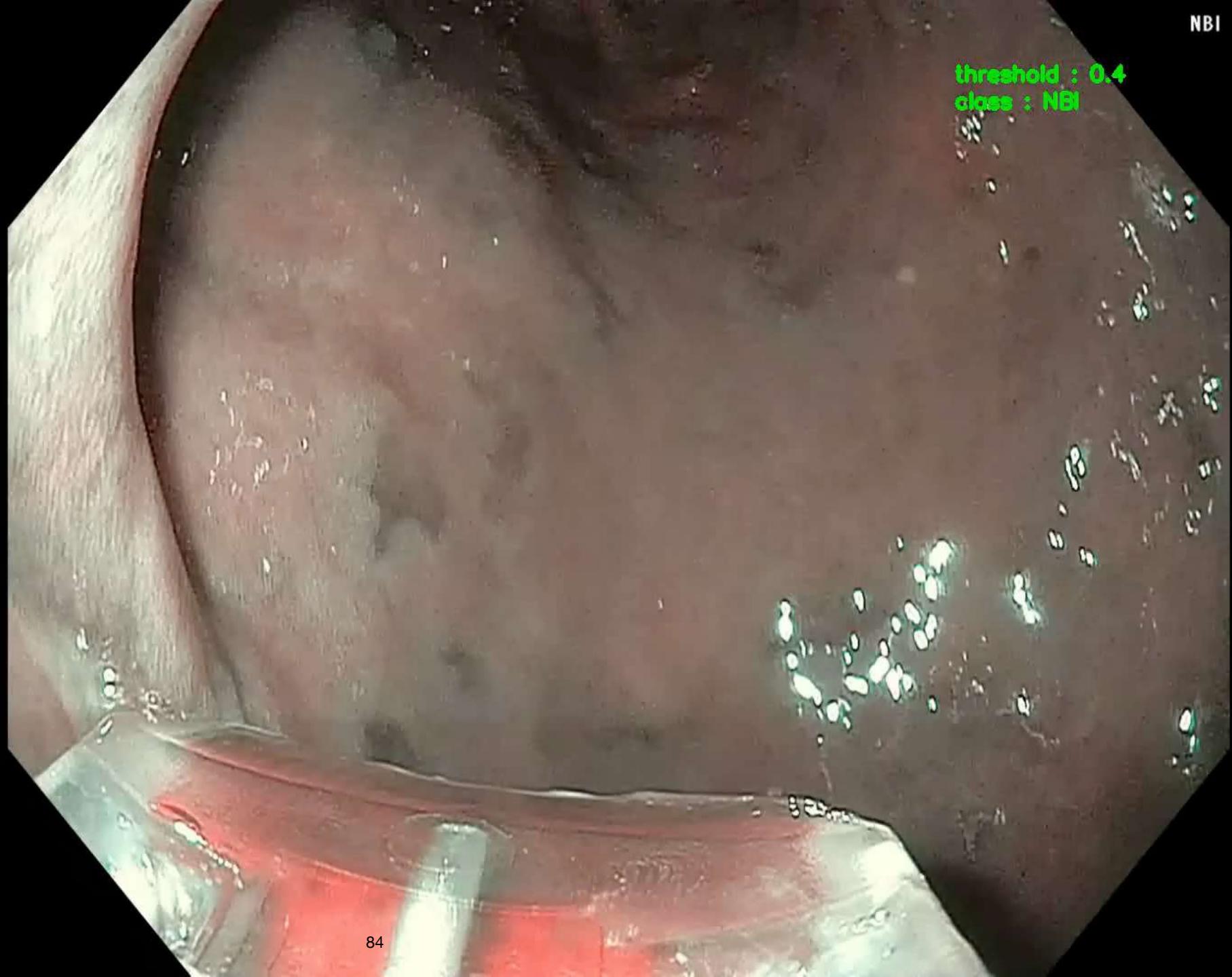
18/05/2021

08:29:36

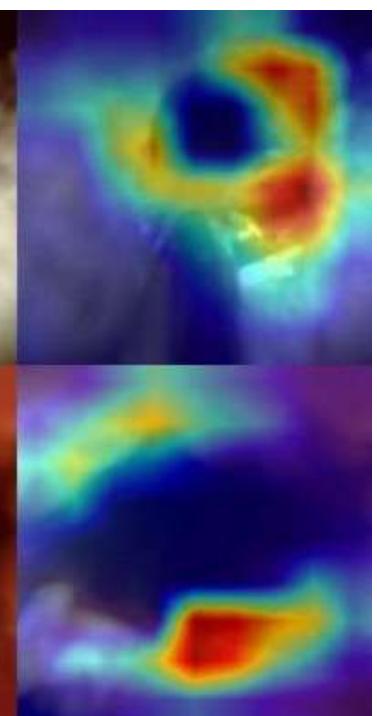
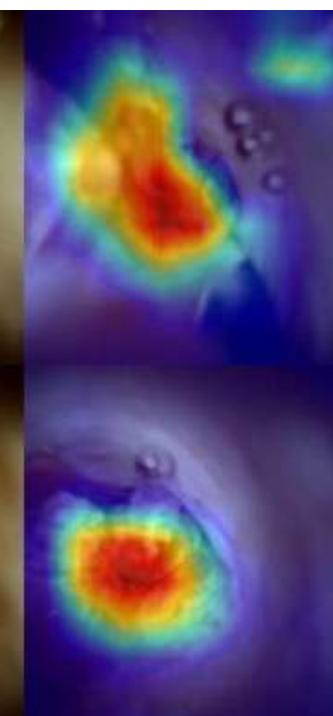
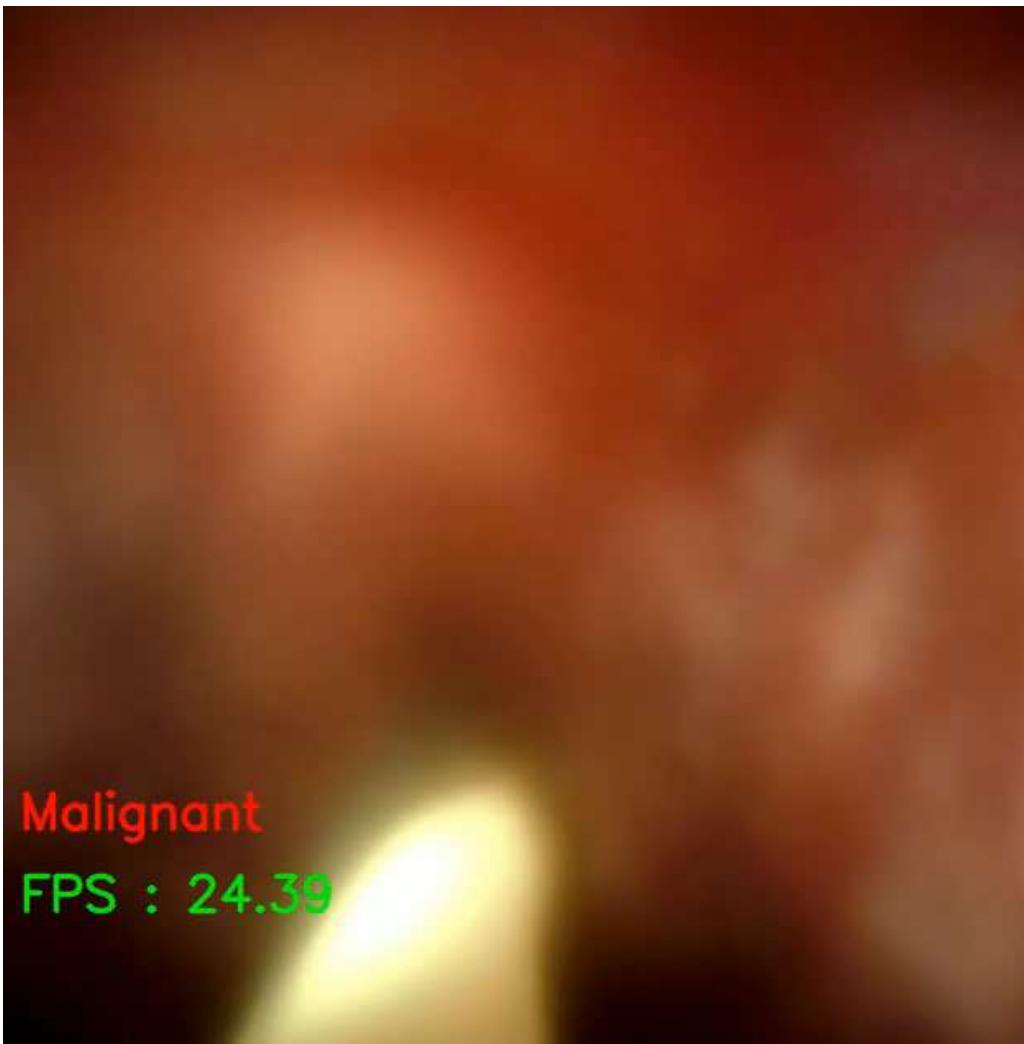
■■□/---(0/1)

Eh:B8 Cm:1

Comment:



Cholangioscopy





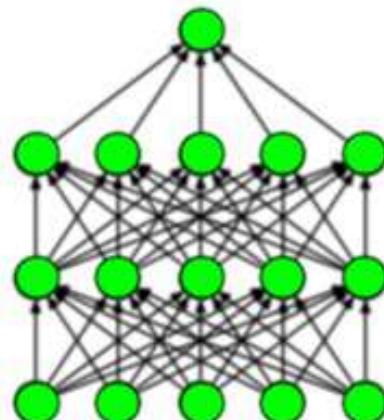
More techniques to prevent overfitting

- Dropout
- Augmentation
- Best validation loss

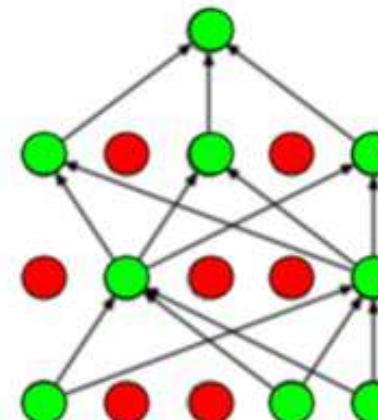


Dropout

- The Dropout layer prevents **overfitting** the model during training. Notably, Dropout **randomly deactivates** some neurons of a layer, thus nullifying their contribution to the output.
- Dropout is only used in training, so we don't want these weights to be fixed at this high a number during testing.



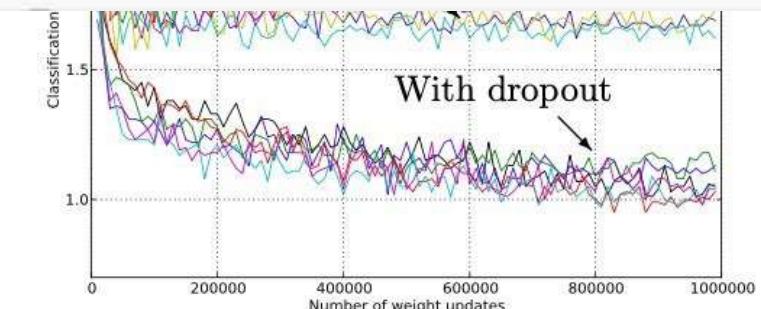
(a) Standard Neural Net



(b) After applying dropout.

```
#create model
model = Sequential()
#add model layers
model.add(Conv2D(64, kernel_size=3, activation='relu', input_shape=(28,28,1)))
model.add(Conv2D(32, kernel_size=3, activation='relu'))
model.add(Dropout(.5, noise_shape=None, seed=None)) #dropout layer
model.add(Conv2D(32, kernel_size=3, activation='relu'))

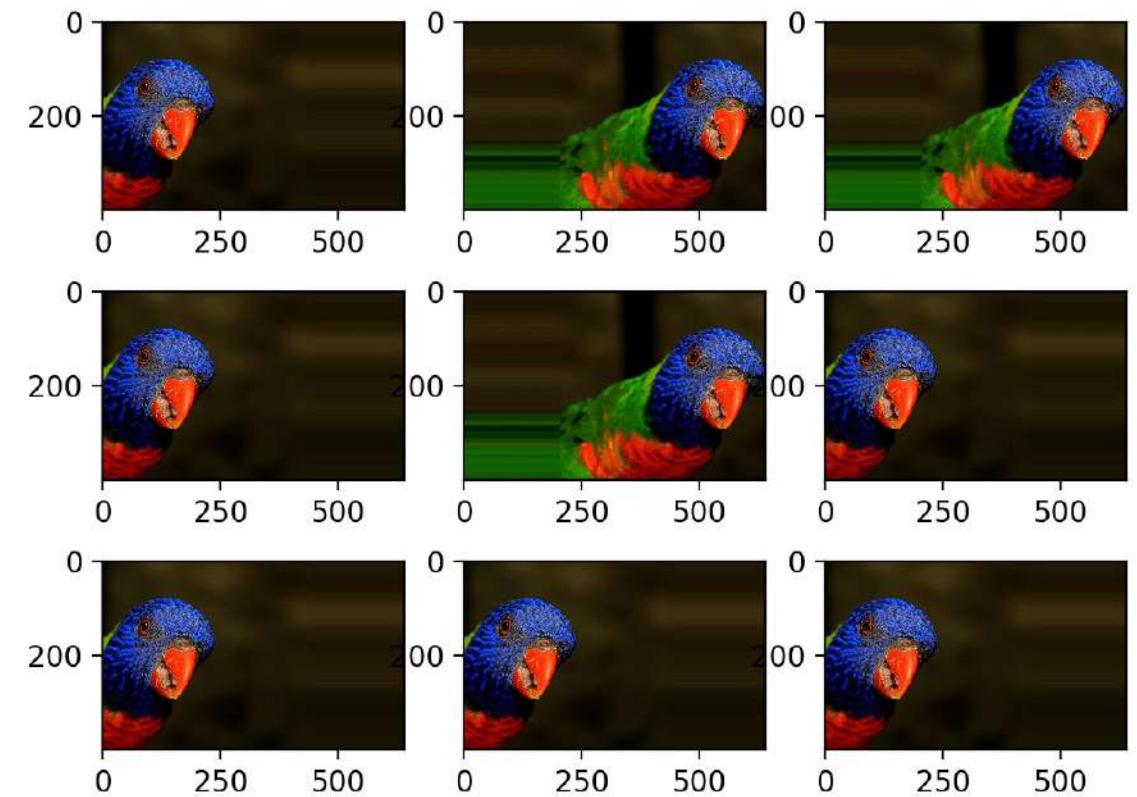
model.add(Flatten())
model.add(Dense(10, activation='softmax'))
```



<https://deepakbattini.medium.com/implementing-drop-out-regularization-in-neural-networks-ab2fc8d985e8>

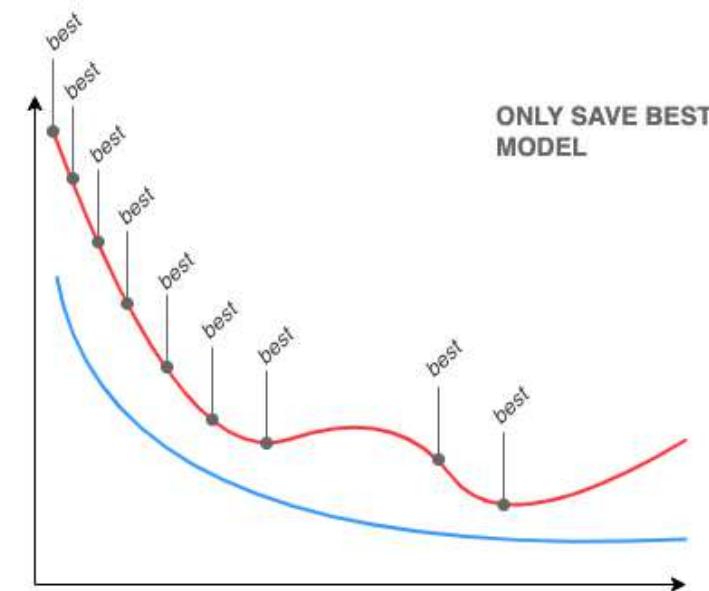
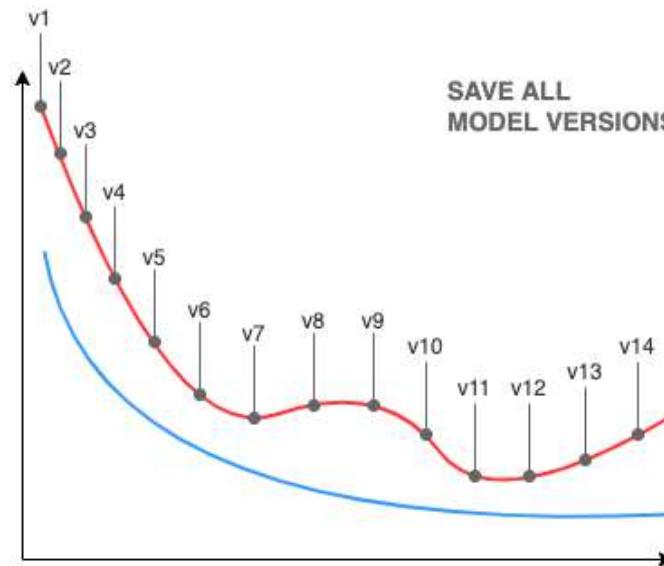
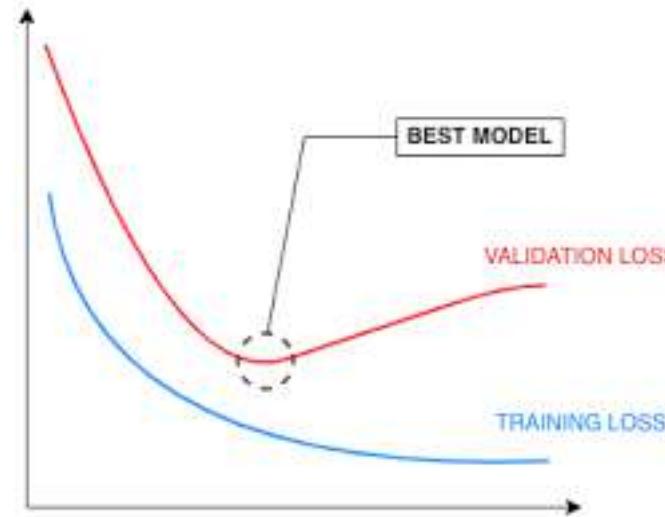


Augmentation





Best validation loss

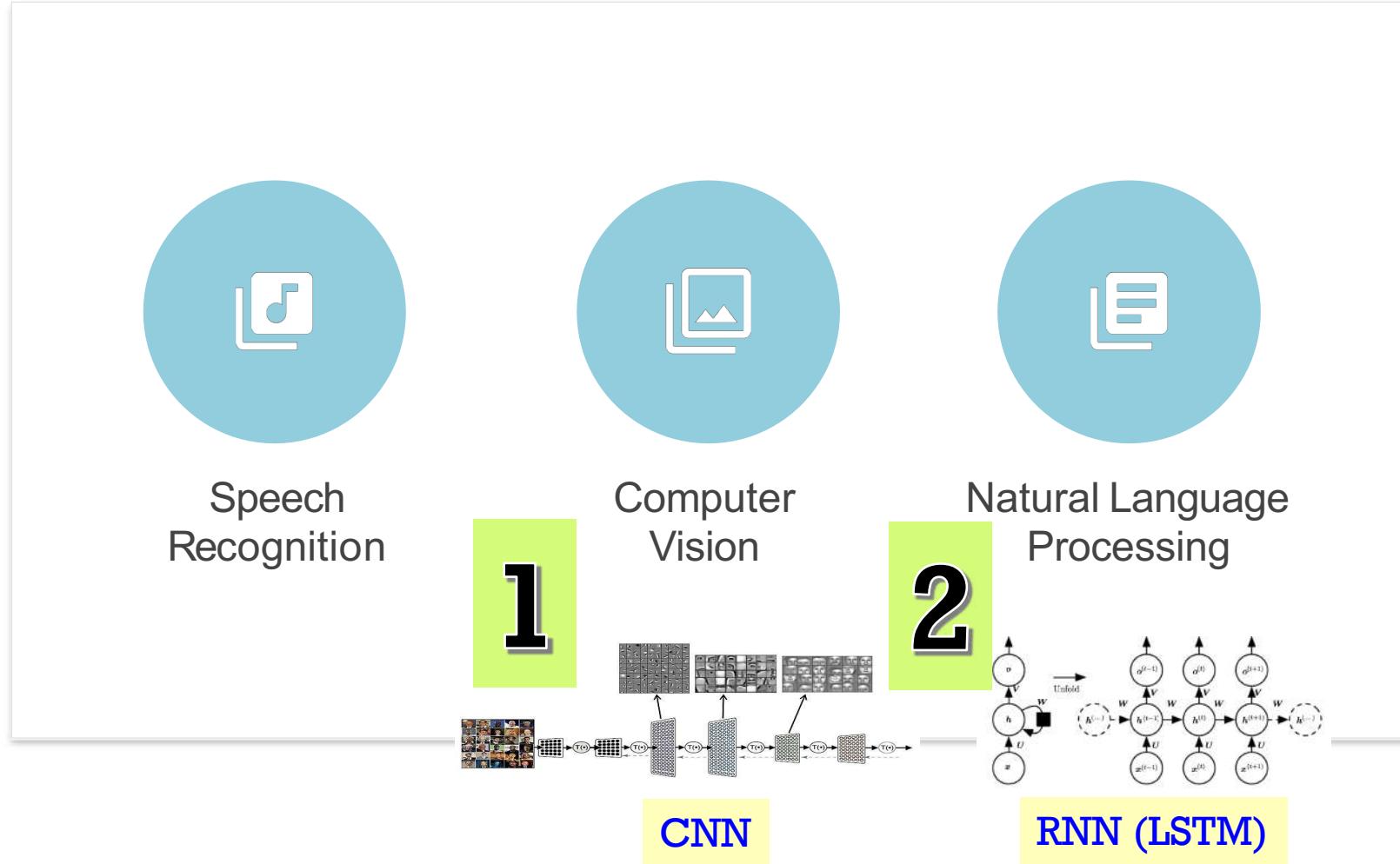


+

RNN



Deep Learning Application



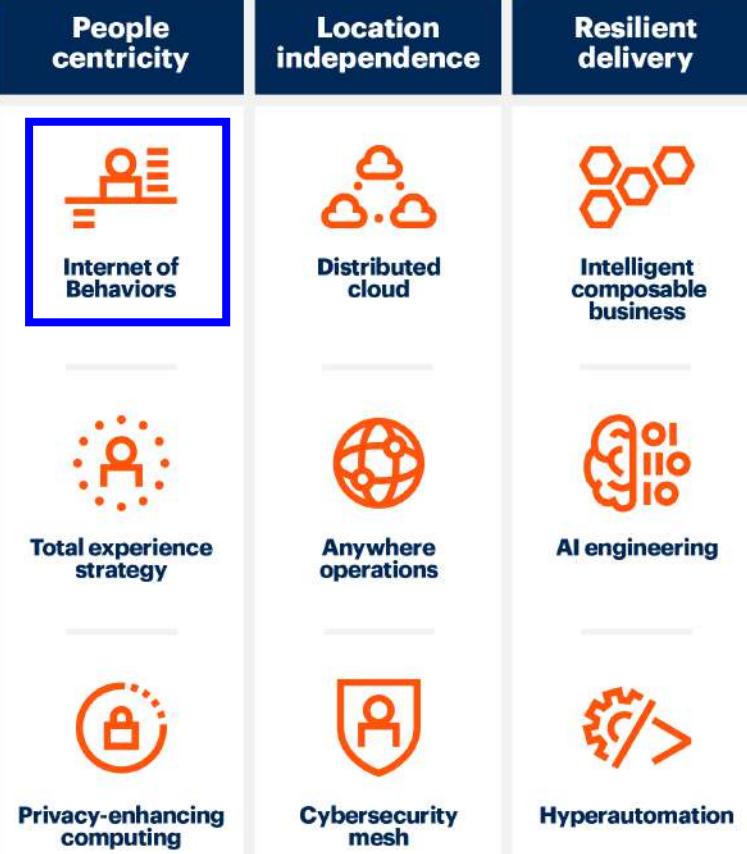
Top 10 Strategic Technology Trends for 2020



gartner.com/SmarterWithGartner

Gartner

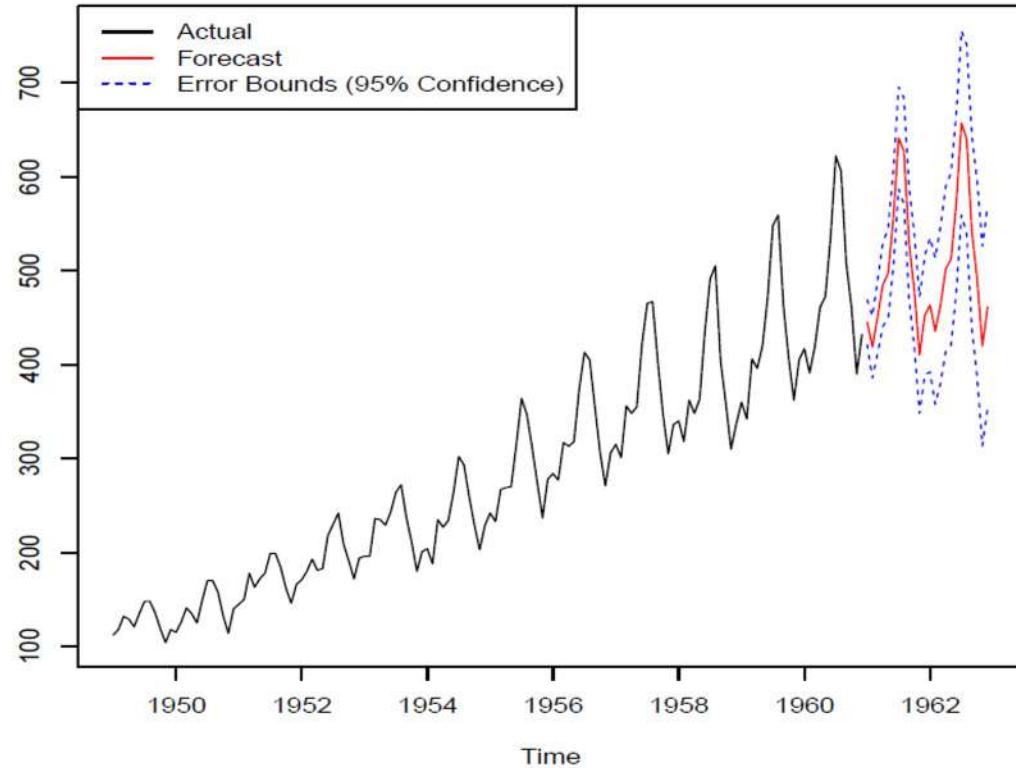
Gartner Top Strategic Technology Trends for 2021



Combinatorial innovation

gartner.com/SmarterWithGartner

Gartner



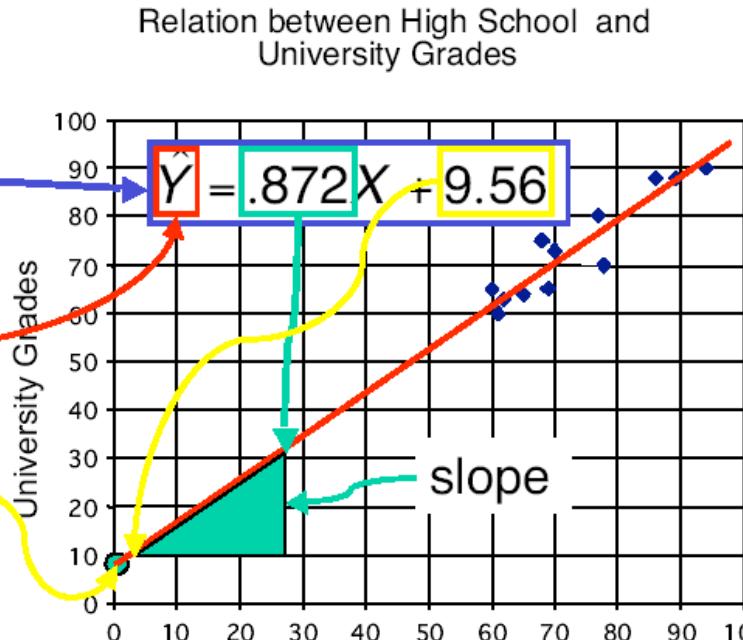


Regression – Linear Relationship

regression equation

predicted value of Y

y -intercept



weight, coefficient

$$\hat{y} = \hat{w}_0 + \hat{w}_1 x_1 + \hat{w}_2 x_2$$

target intercept input

- The least square method aims to minimize the following term

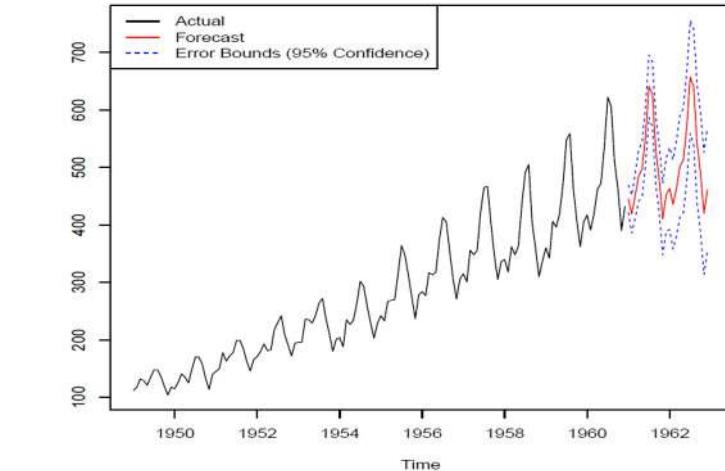
$$\sum_{\text{training data}} (y_i - \hat{y}_i)^2$$

Autoregressive Model – Linear Relationship

- Based on linear regression, but using previous timestep data

$$X_t = c + \sum_{i=1}^p \varphi_i X_{t-i}$$

- Where X_t is data at timestep t , c is constant, and each timestep data



are parameters for
 $\varphi_1, \varphi_2, \varphi_3, \dots, \varphi_p$

Example: When $p = 2$ (looking back two steps), the equation will be

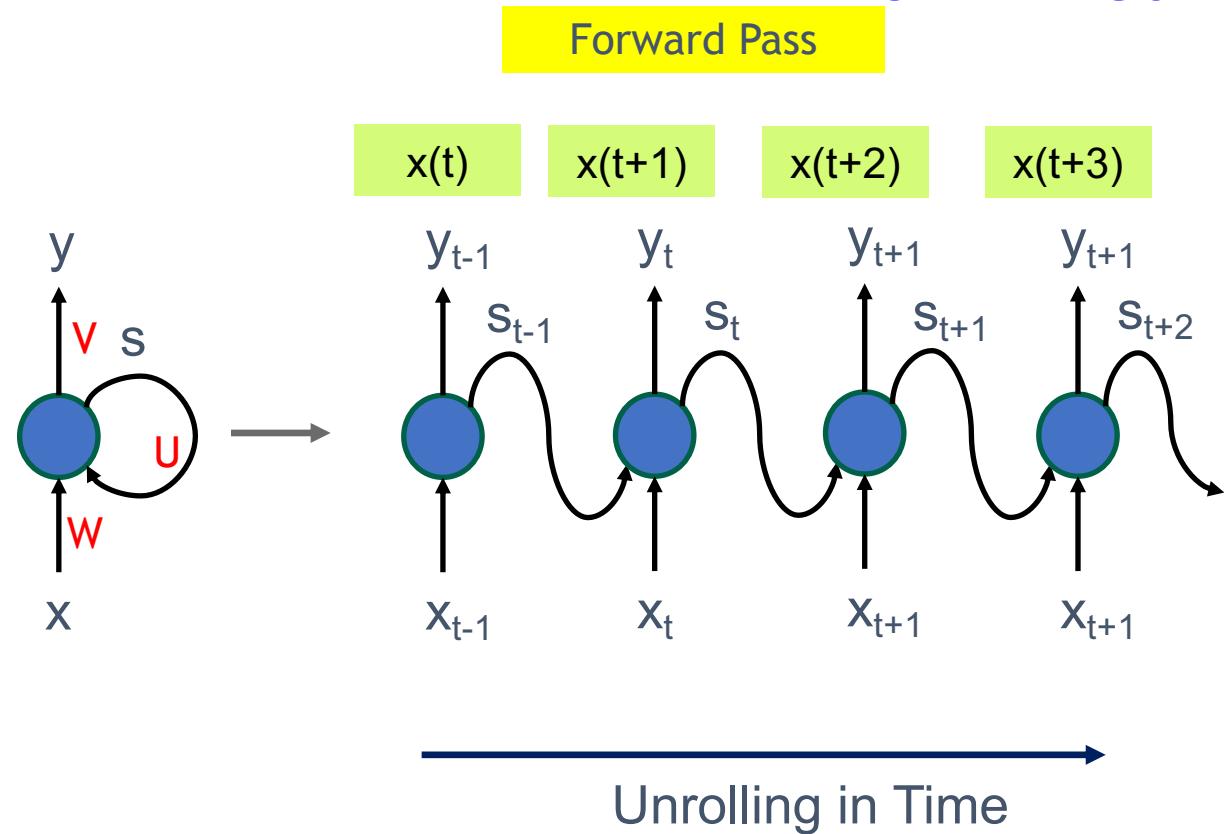
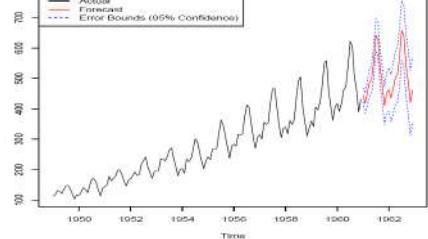
$$X_t = c + \varphi_1 X_{t-1} + \varphi_2 X_{t-2}$$

- Parameters are able to calculate using various method, such as ordinary least square procedure or Yule–Walker equations

$$x(t + 1) = \hat{w}_0 + \hat{w}_1 x(t) + \hat{w}_2 x(t - 1)$$

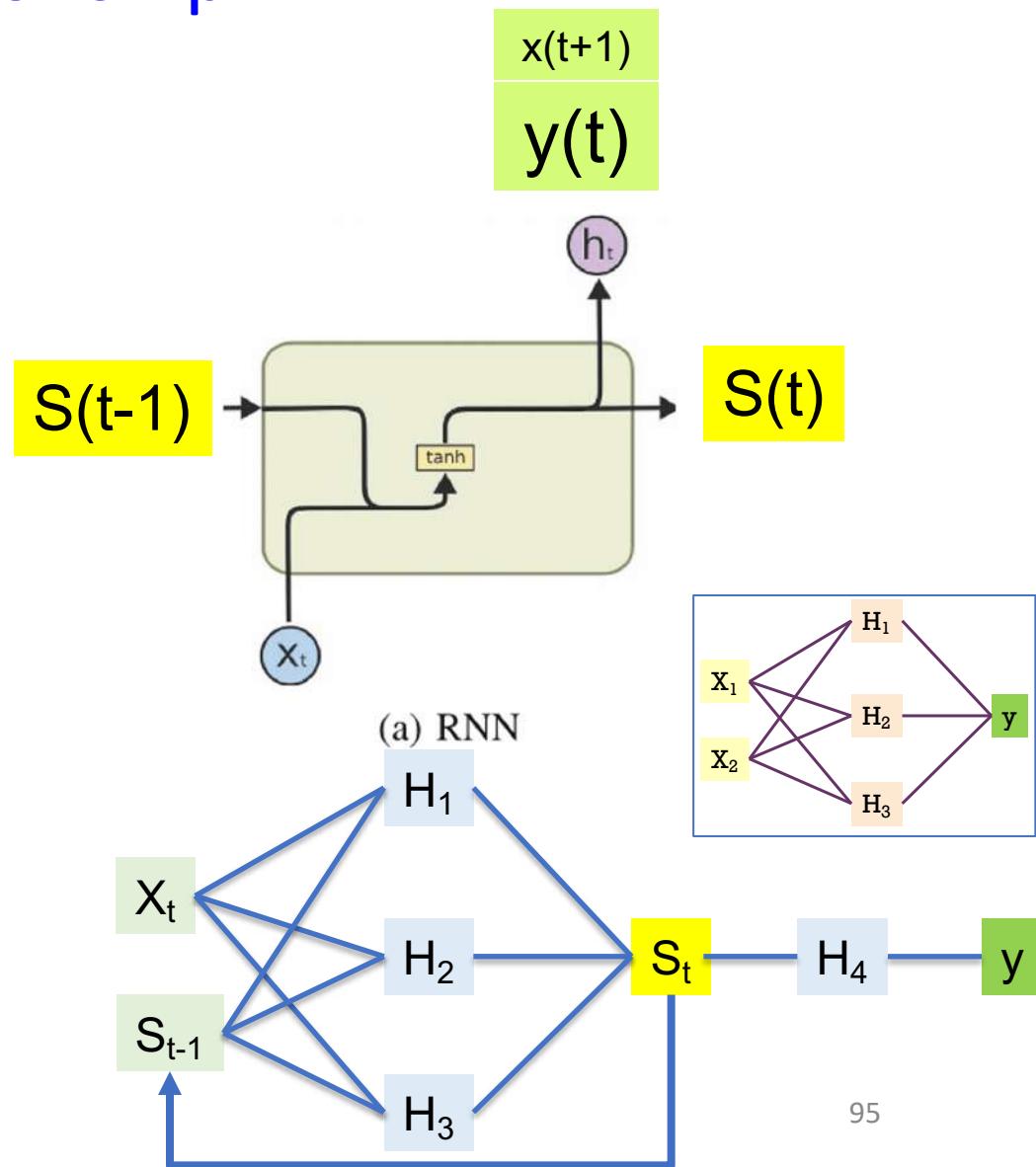
RECURRENT NEURONS (RNN)

Non-Linear Relationship

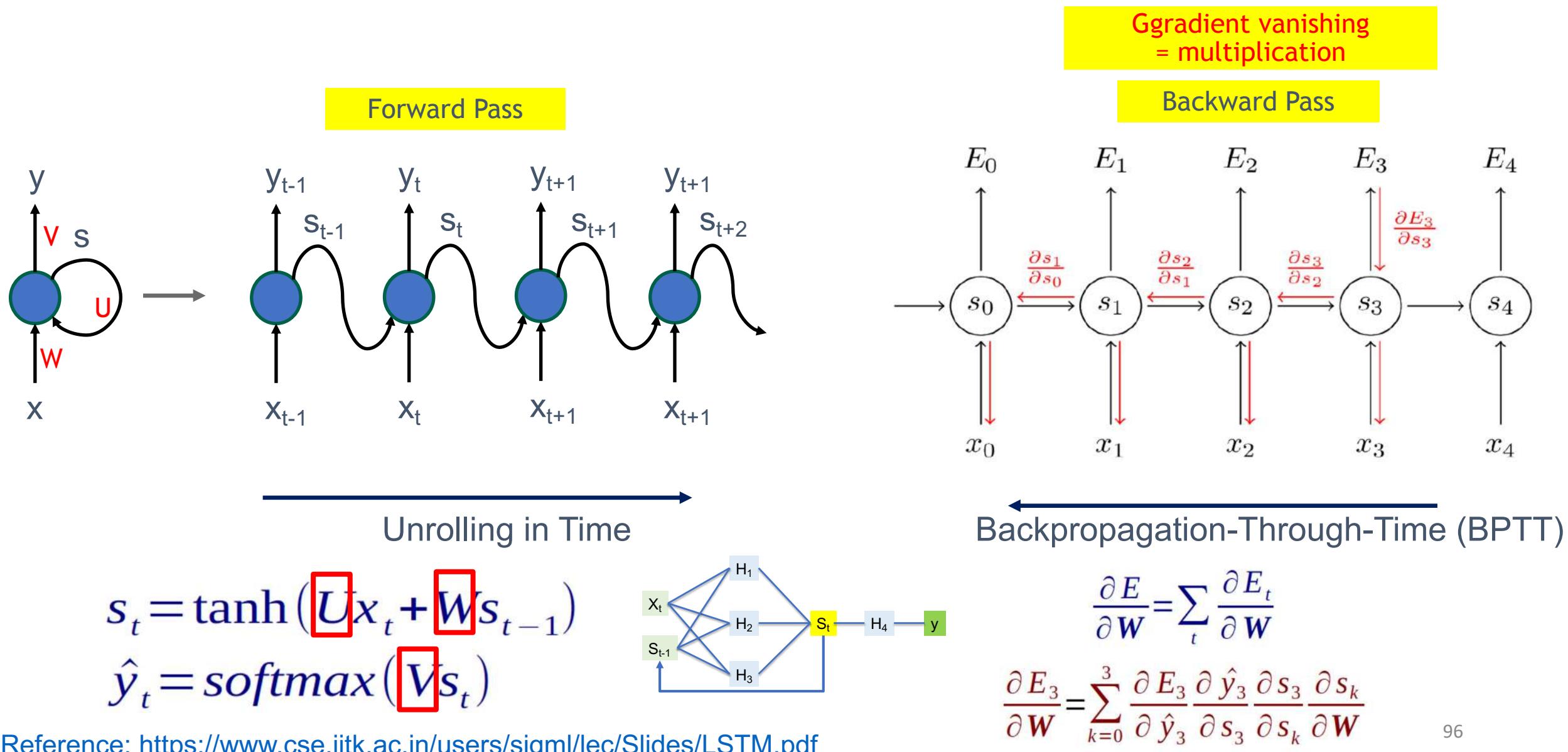


$$s_t = \tanh(Ux_t + Vs_{t-1})$$

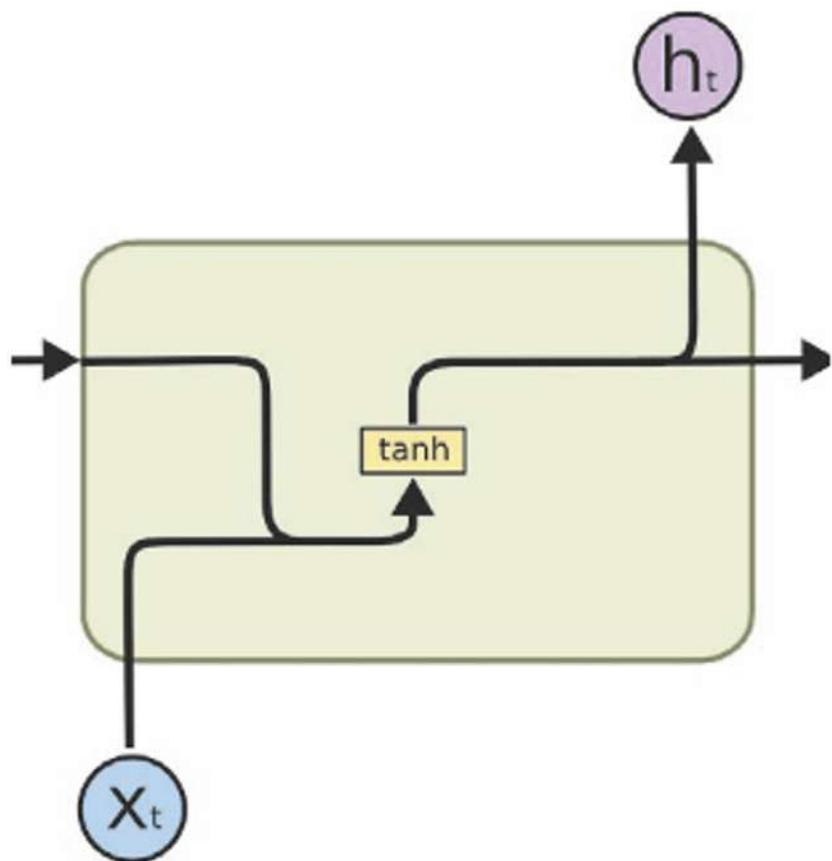
$$x(t+1) \hat{y}_t = \text{softmax}(Vs_t)$$



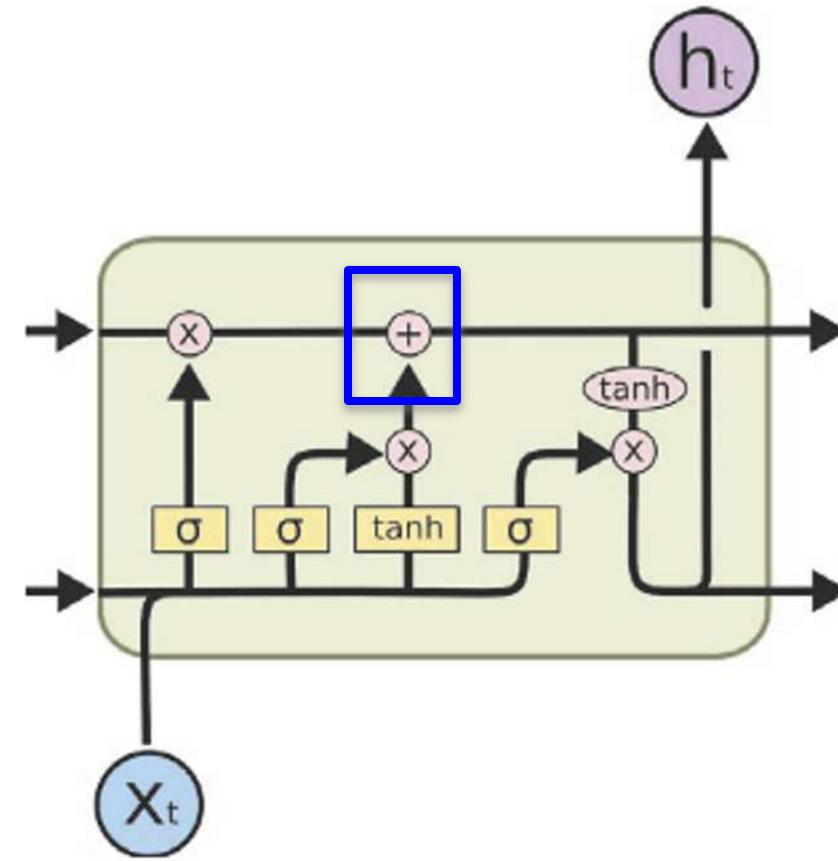
RECURRENT NEURONS (RNN)



RNN VS LSTM

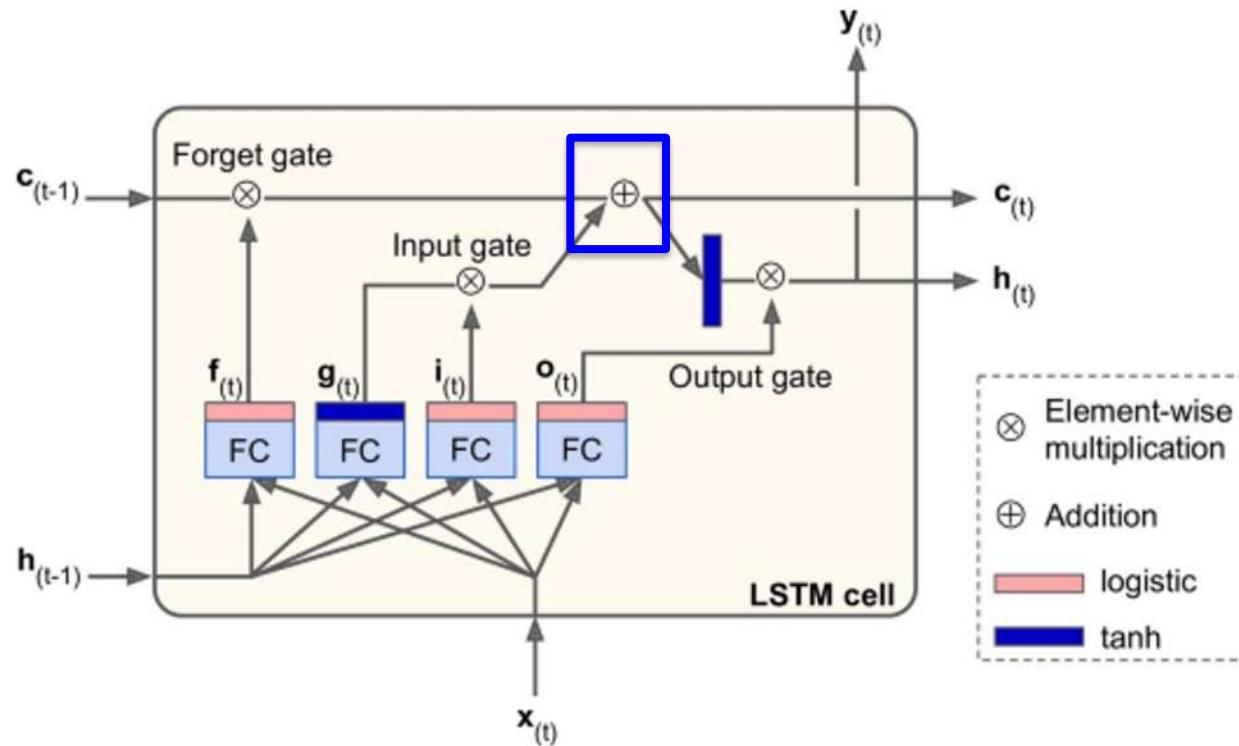


(a) RNN



(b) LSTM

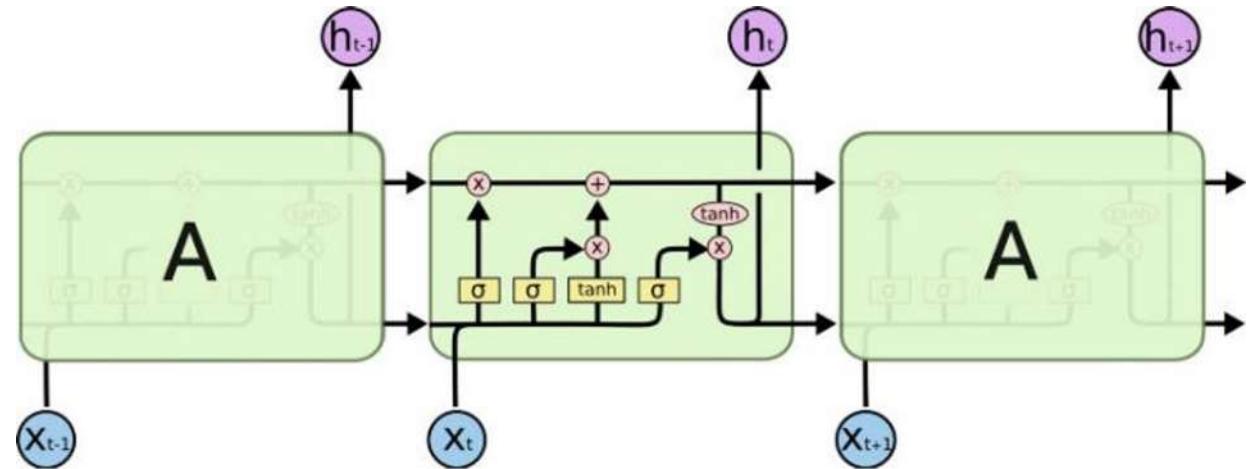
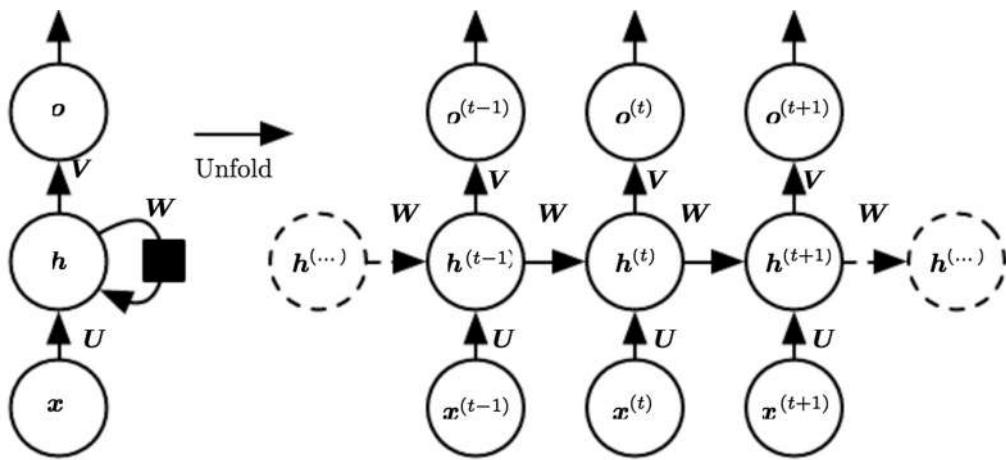
LONG SHORT TERM (LSTM) CELL



There are 3 gates with 2 outputs.

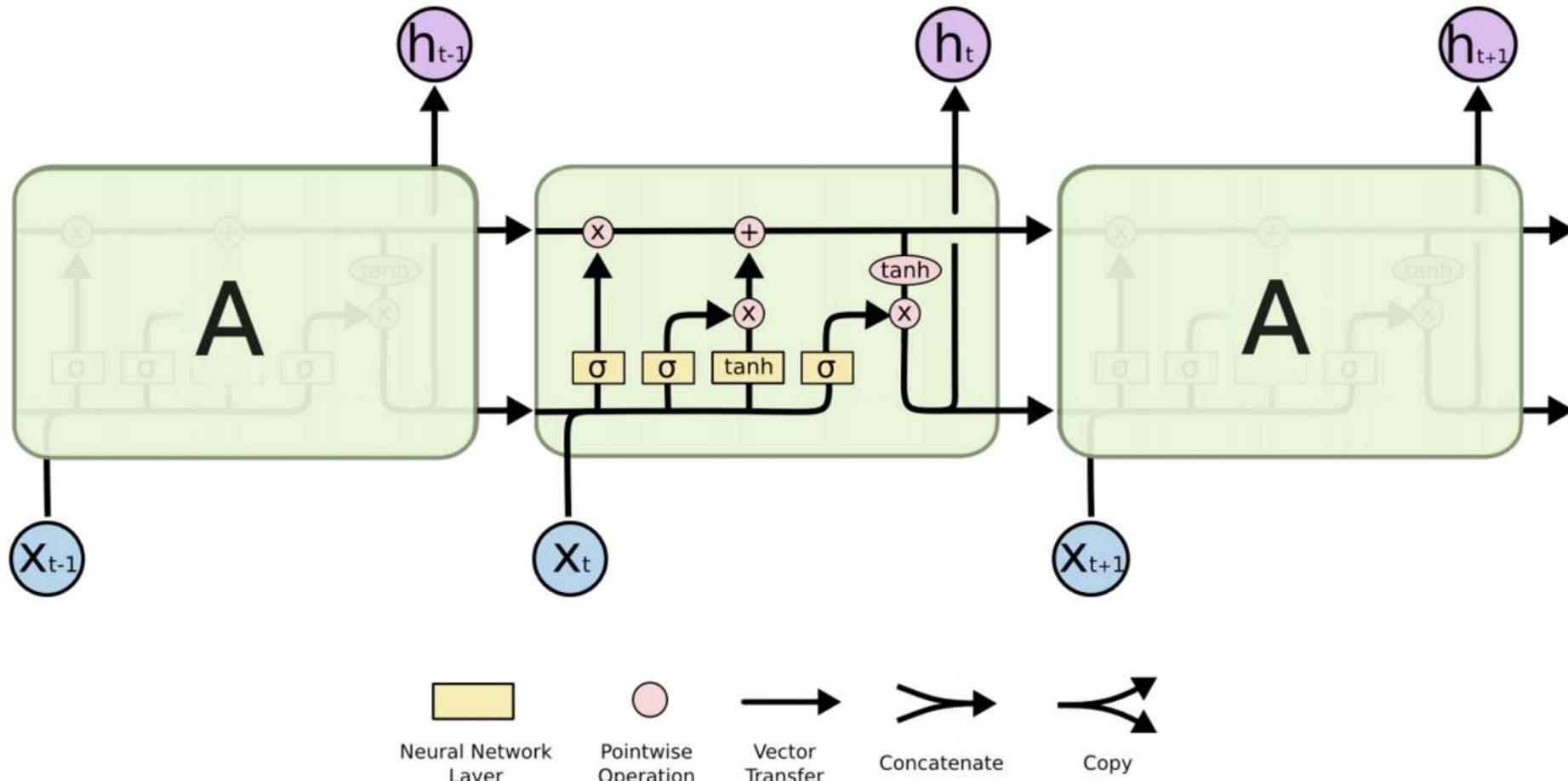
$$\boxed{\begin{aligned}\mathbf{i}_{(t)} &= \sigma(\mathbf{W}_{xi}^T \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hi}^T \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_i) \\ \mathbf{f}_{(t)} &= \sigma(\mathbf{W}_{xf}^T \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hf}^T \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_f) \\ \mathbf{o}_{(t)} &= \sigma(\mathbf{W}_{xo}^T \cdot \mathbf{x}_{(t)} + \mathbf{W}_{ho}^T \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_o) \\ \mathbf{g}_{(t)} &= \tanh(\mathbf{W}_{xg}^T \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hg}^T \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_g) \\ \mathbf{c}_{(t)} &= \mathbf{f}_{(t)} \otimes \mathbf{c}_{(t-1)} + \mathbf{i}_{(t)} \otimes \mathbf{g}_{(t)} \\ \mathbf{y}_{(t)} &= \mathbf{h}_{(t)} = \mathbf{o}_{(t)} \otimes \tanh(\mathbf{c}_{(t)})\end{aligned}}$$

Deep Learning Approach: LSTM

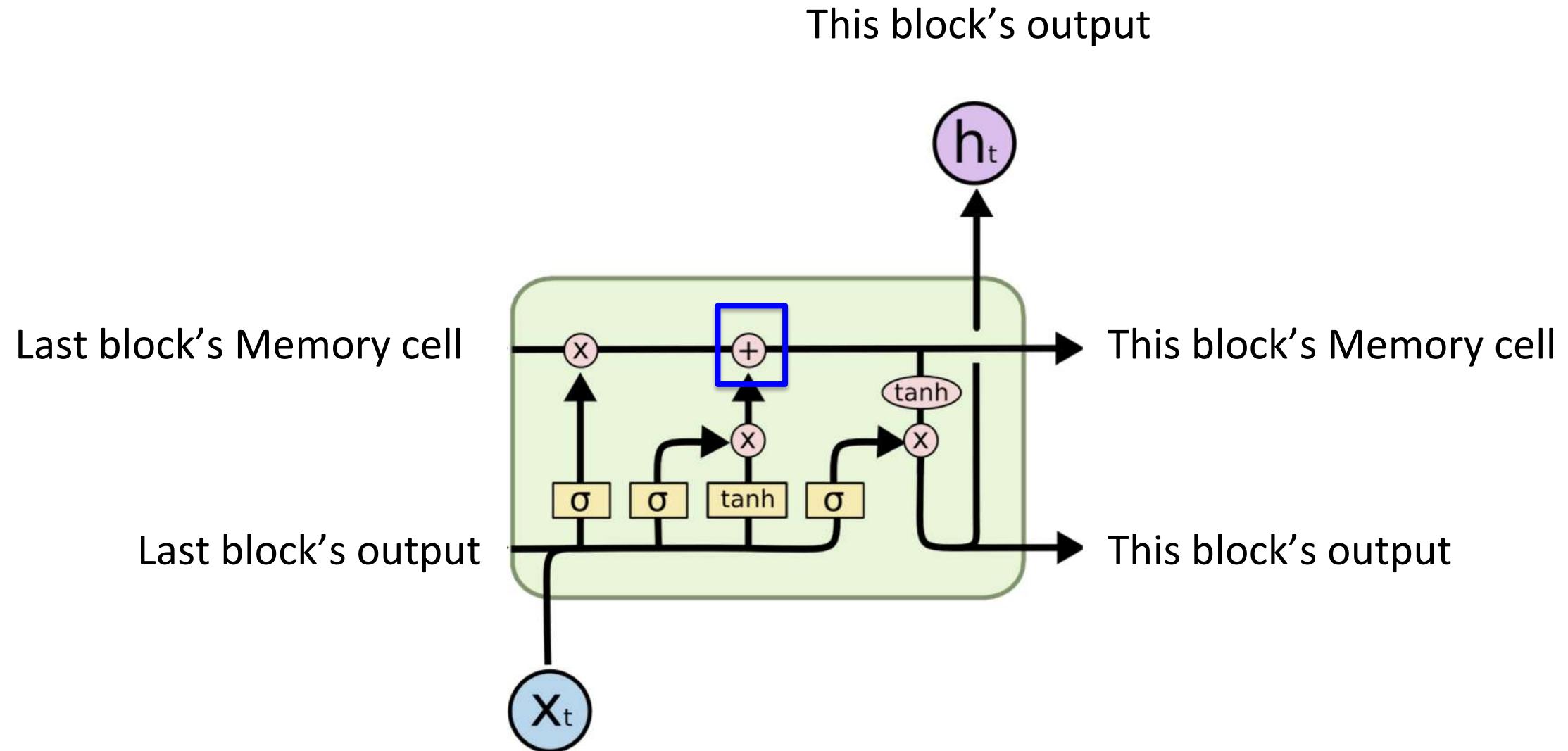


Reference: Ian Goodfellow, Yoshua Bengio, and Aaron Courville. 2016. Deep Learning. MIT Press.
<http://www.deeplearningbook.org>

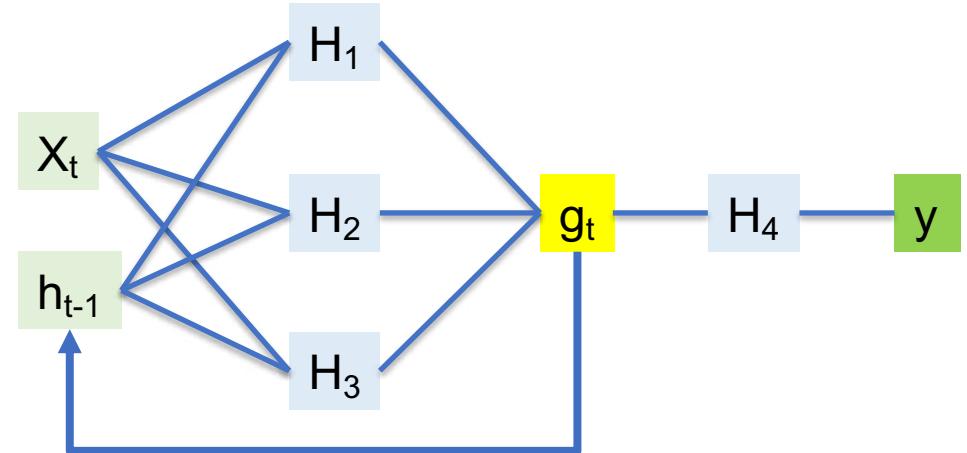
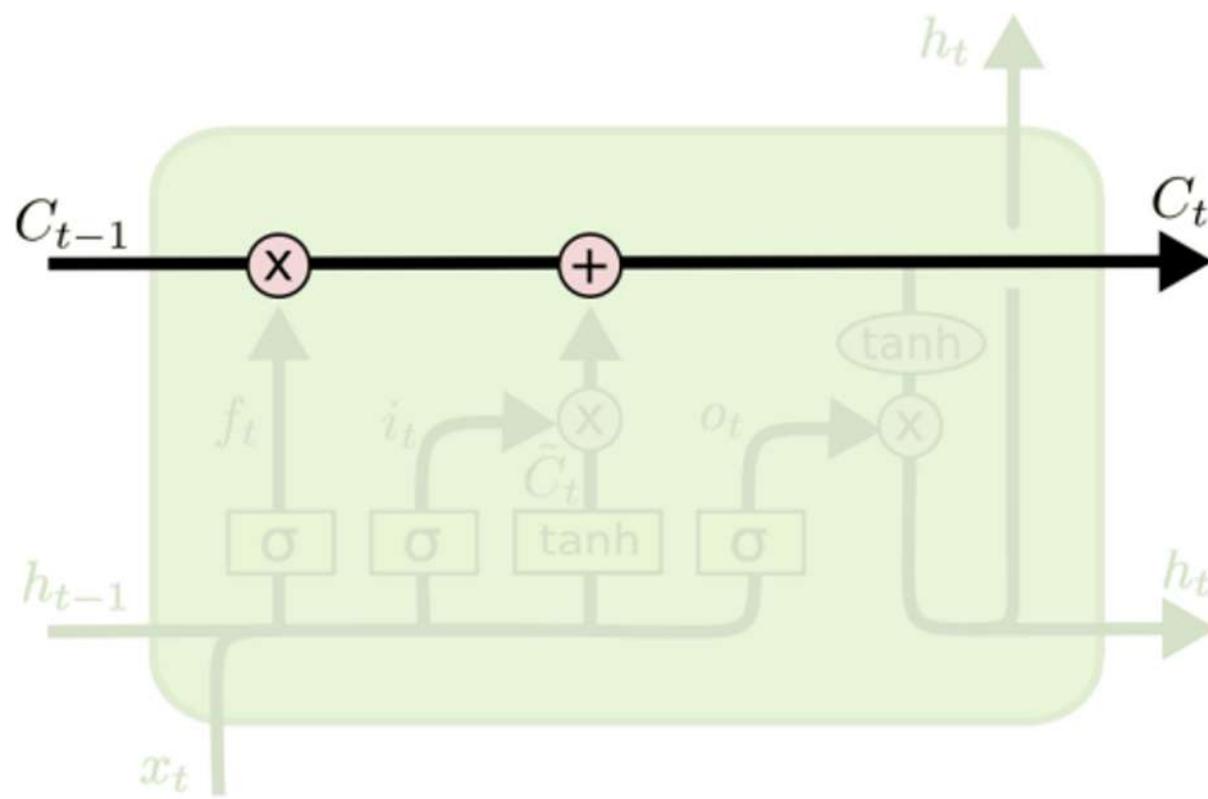
Inside LSTM



Inside LSTM



Inside LSTM: (1) Cell Memory



There are 3 gates with 2 outputs.

$$\mathbf{i}_{(t)} = \sigma(\mathbf{W}_{xi}^T \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hi}^T \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_i)$$

$$\mathbf{f}_{(t)} = \sigma(\mathbf{W}_{xf}^T \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hf}^T \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_f)$$

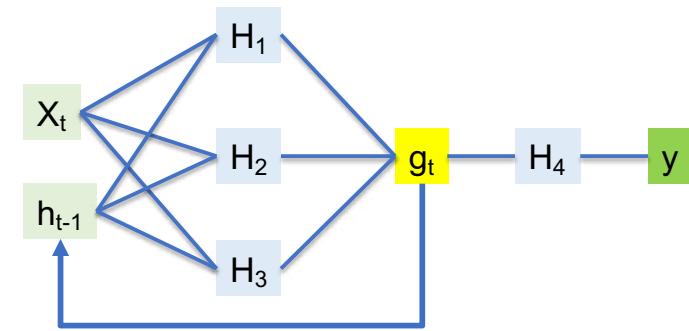
$$\mathbf{o}_{(t)} = \sigma(\mathbf{W}_{xo}^T \cdot \mathbf{x}_{(t)} + \mathbf{W}_{ho}^T \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_o)$$

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

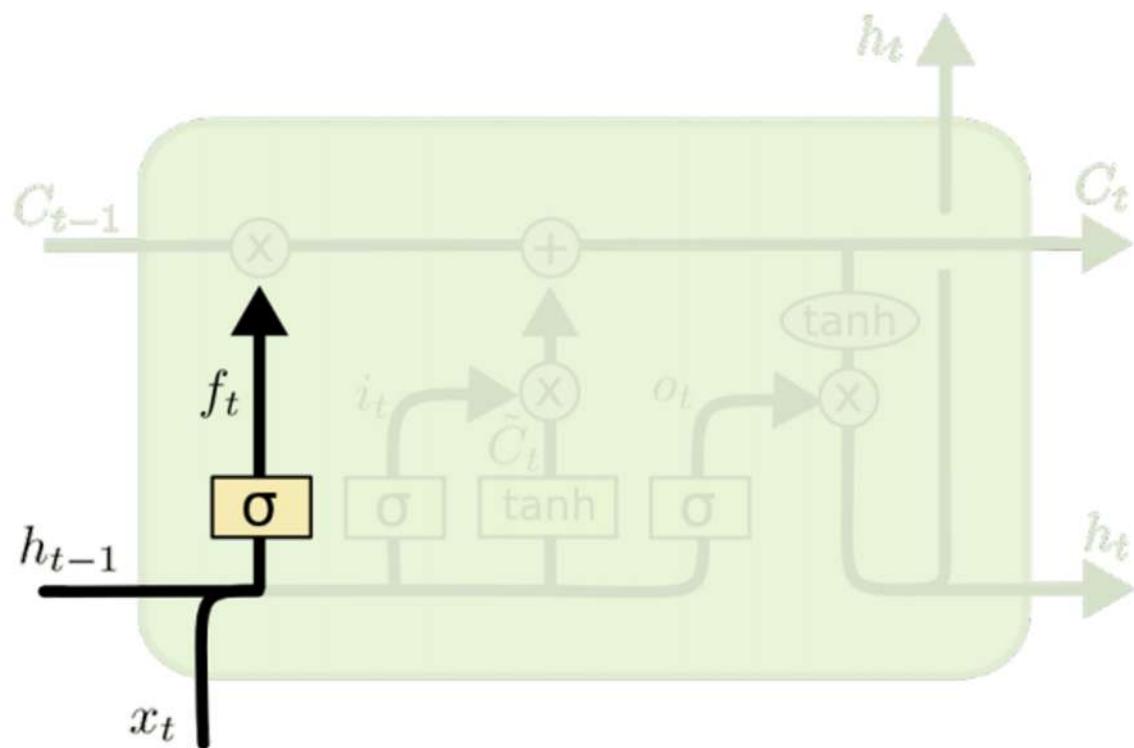
$$\mathbf{c}_{(t)} = \mathbf{f}_{(t)} \otimes \mathbf{c}_{(t-1)} + \mathbf{i}_{(t)} \otimes \mathbf{g}_{(t)}$$

$$\mathbf{y}_{(t)} = \mathbf{h}_{(t)} = \mathbf{o}_{(t)} \otimes \tanh(\mathbf{c}_{(t)})$$

Inside LSTM: (2) Forget Gate



Should we continue to remember this “bit” of information or not?



There are 3 gates with 2 outputs.

$$\mathbf{i}_{(t)} = \sigma(\mathbf{W}_{xi}^T \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hi}^T \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_i)$$

$$\boxed{\mathbf{f}_{(t)} = \sigma(\mathbf{W}_{xf}^T \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hf}^T \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_f)}$$

$$\mathbf{o}_{(t)} = \sigma(\mathbf{W}_{xo}^T \cdot \mathbf{x}_{(t)} + \mathbf{W}_{ho}^T \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_o)$$

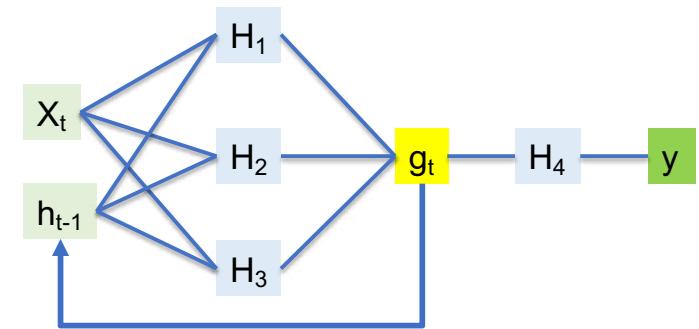
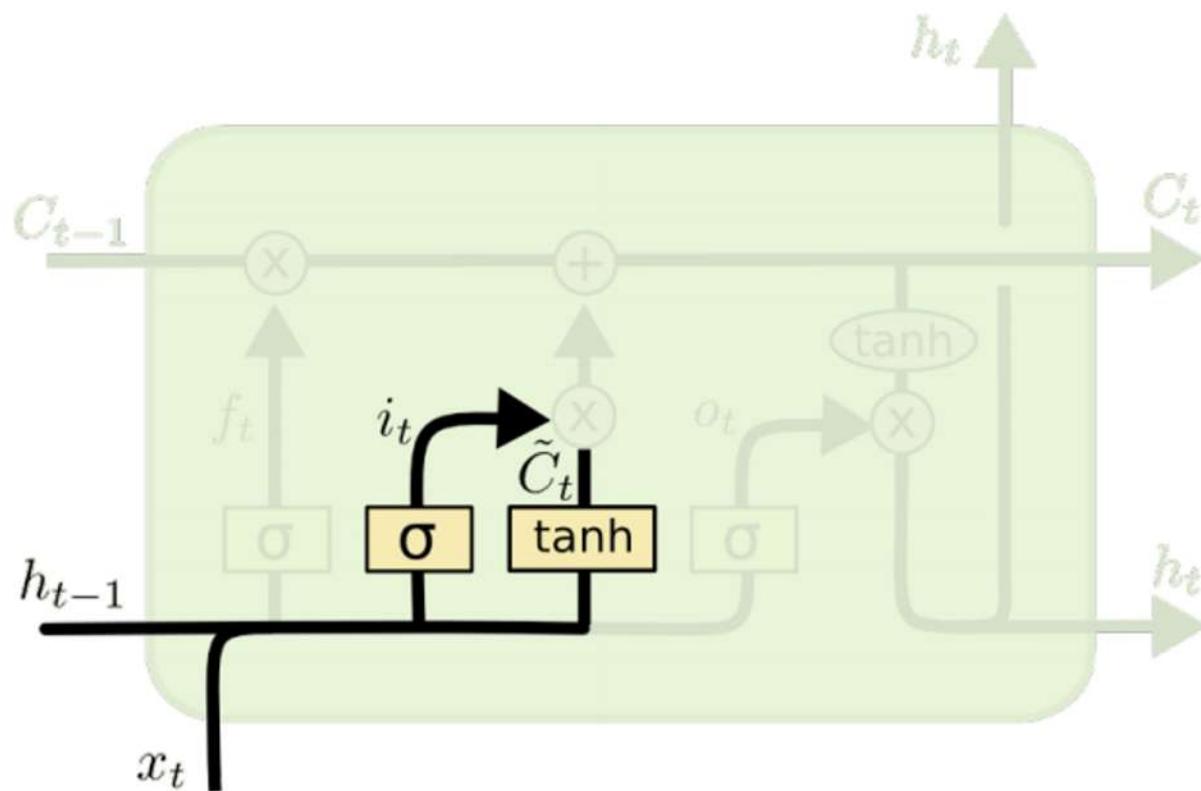
$$\mathbf{g}_{(t)} = \tanh(\mathbf{W}_{xg}^T \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hg}^T \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_g)$$

$$\mathbf{c}_{(t)} = \mathbf{f}_{(t)} \otimes \mathbf{c}_{(t-1)} + \mathbf{i}_{(t)} \otimes \mathbf{g}_{(t)}$$

$$\mathbf{y}_{(t)} = \mathbf{h}_{(t)} = \mathbf{o}_{(t)} \otimes \tanh(\mathbf{c}_{(t)})$$

Inside LSTM: (3) Input Gate

Should we update this “bit” of information or not?
If yes, then what should we remember?



There are 3 gates with 2 outputs.

$$\mathbf{i}_{(t)} = \sigma(\mathbf{W}_{xi}^T \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hi}^T \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_i)$$

$$\mathbf{f}_{(t)} = \sigma(\mathbf{W}_{xf}^T \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hf}^T \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_f)$$

$$\mathbf{o}_{(t)} = \sigma(\mathbf{W}_{xo}^T \cdot \mathbf{x}_{(t)} + \mathbf{W}_{ho}^T \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_o)$$

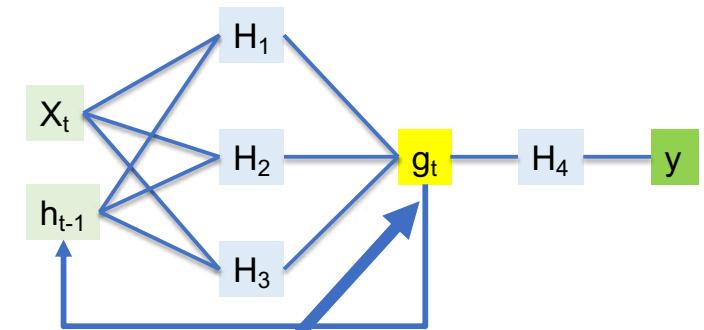
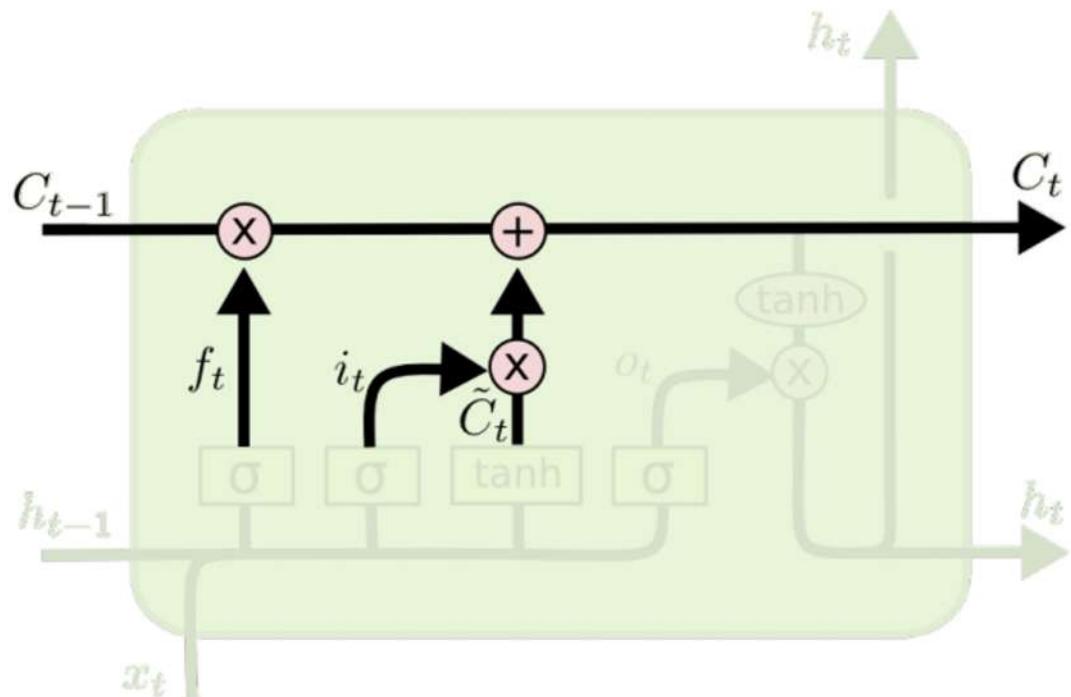
$$\mathbf{g}_{(t)} = \tanh(\mathbf{W}_{xg}^T \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hg}^T \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_g)$$

$$\mathbf{c}_{(t)} = \mathbf{f}_{(t)} \otimes \mathbf{c}_{(t-1)} + \mathbf{i}_{(t)} \otimes \mathbf{g}_{(t)}$$

$$\mathbf{y}_{(t)} = \mathbf{h}_{(t)} = \mathbf{o}_{(t)} \otimes \tanh(\mathbf{c}_{(t)})$$

Inside LSTM: (4) Memory Update

Forget what needs to be forgotten + memorize what needs to be remembered



There are 3 gates with 2 outputs.

$$\mathbf{i}_{(t)} = \sigma(\mathbf{W}_{xi}^T \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hi}^T \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_i)$$

$$\mathbf{f}_{(t)} = \sigma(\mathbf{W}_{xf}^T \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hf}^T \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_f)$$

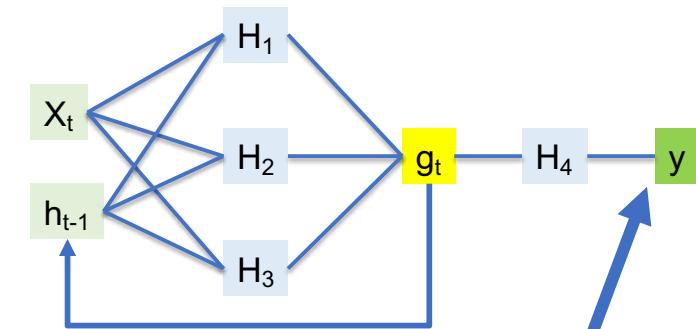
$$\mathbf{o}_{(t)} = \sigma(\mathbf{W}_{xo}^T \cdot \mathbf{x}_{(t)} + \mathbf{W}_{ho}^T \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_o)$$

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

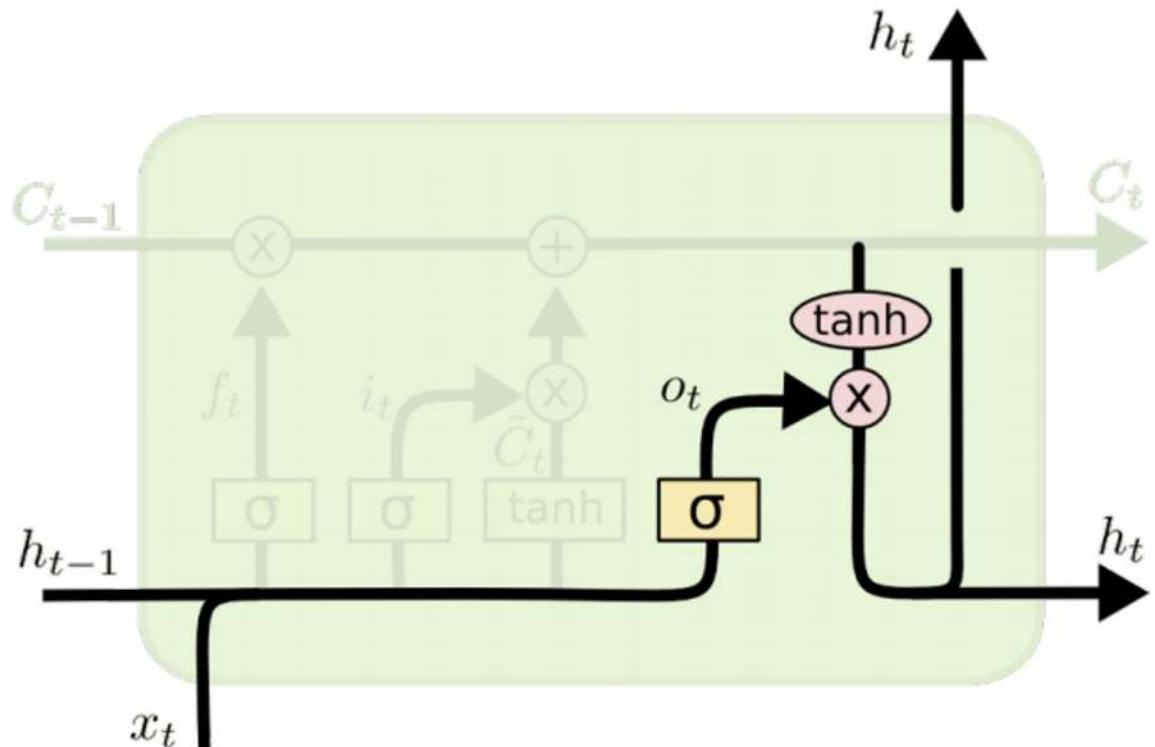
$$\mathbf{c}_{(t)} = \mathbf{f}_{(t)} \otimes \mathbf{c}_{(t-1)} + \mathbf{i}_{(t)} \otimes \mathbf{g}_{(t)}$$

$$\mathbf{y}_{(t)} = \mathbf{h}_{(t)} = \mathbf{o}_{(t)} \otimes \tanh(\mathbf{c}_{(t)})$$

Inside LSTM: Output Gate



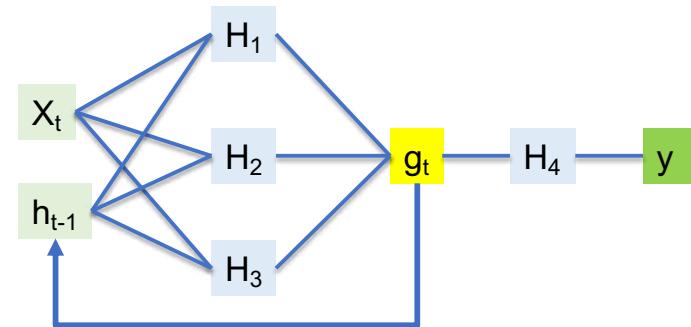
Should we output this bit of information (e.g., to “deeper” LSTM layers)?



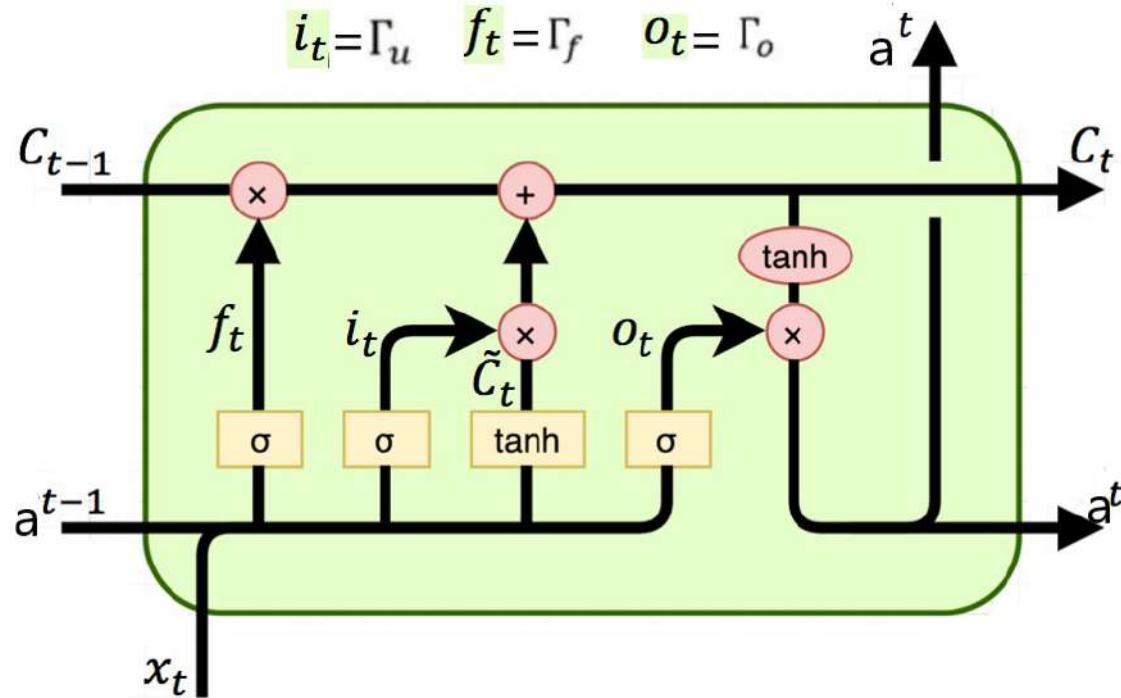
There are 3 gates with 2 outputs.

$$\begin{aligned}\mathbf{i}_{(t)} &= \sigma(\mathbf{W}_{xi}^T \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hi}^T \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_i) \\ \mathbf{f}_{(t)} &= \sigma(\mathbf{W}_{xf}^T \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hf}^T \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_f) \\ \mathbf{o}_{(t)} &= \sigma(\mathbf{W}_{xo}^T \cdot \mathbf{x}_{(t)} + \mathbf{W}_{ho}^T \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_o) \\ \mathbf{g}_{(t)} &= \tanh(\mathbf{W}_{xg}^T \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hg}^T \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_g) \\ \mathbf{c}_{(t)} &= \mathbf{f}_{(t)} \otimes \mathbf{c}_{(t-1)} + \mathbf{i}_{(t)} \otimes \mathbf{g}_{(t)} \\ \mathbf{y}_{(t)} &= \mathbf{h}_{(t)} = \mathbf{o}_{(t)} \otimes \tanh(\mathbf{c}_{(t)})\end{aligned}$$

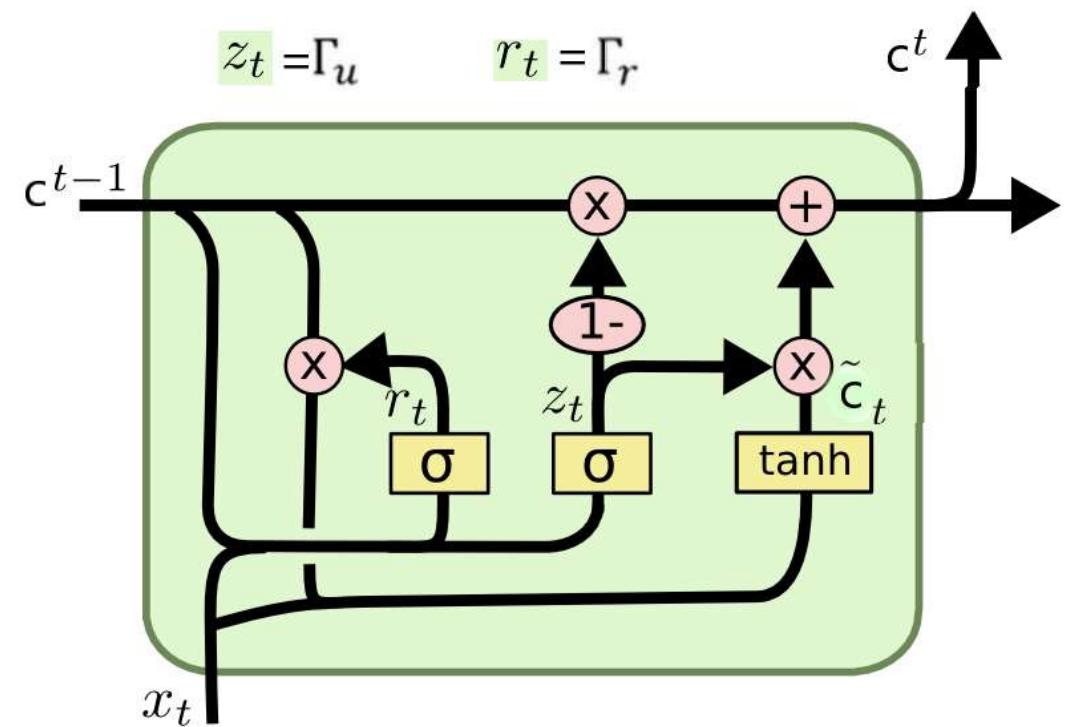
LSTM vs. Gated Recurrent Units (GRU)



There are 3 gates with 2 outputs.
3 inputs ($x(t)$, $C(t-1)$, $a(t-1)$)

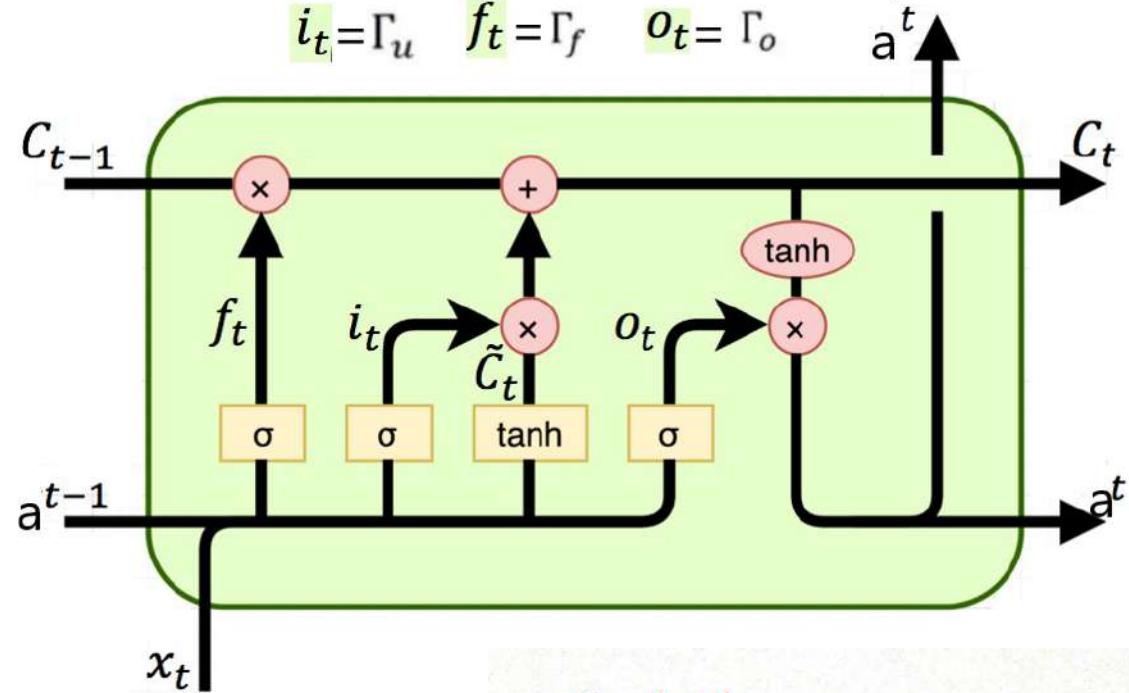


- There are 2 gates (update & reset gates)
- 1 output
- 2 inputs ($x(t)$, $C(t-1)$, $a(t-1)$)



There are 3 gates with 2 outputs.

$$i_t = \Gamma_u \quad f_t = \Gamma_f \quad o_t = \Gamma_o$$



$$\begin{aligned} i_{(t)} &= \sigma(W_{xi}^T \cdot x_{(t)} + W_{hi}^T \cdot h_{(t-1)} + b_i) \\ f_{(t)} &= \sigma(W_{xf}^T \cdot x_{(t)} + W_{hf}^T \cdot h_{(t-1)} + b_f) \\ o_{(t)} &= \sigma(W_{xo}^T \cdot x_{(t)} + W_{ho}^T \cdot h_{(t-1)} + b_o) \\ g_{(t)} &= \tanh(W_{xg}^T \cdot x_{(t)} + W_{hg}^T \cdot h_{(t-1)} + b_g) \\ c_{(t)} &= f_{(t)} \otimes c_{(t-1)} + i_{(t)} \otimes g_{(t)} \\ y_{(t)} &= h_{(t)} = o_{(t)} \otimes \tanh(c_{(t)}) \end{aligned}$$

Candidate cell $\rightarrow \tilde{c}^{(t)} = \tanh(W_c[a^{(t-1)} + x^{(t)}] + b_c)$

Update / Input gate $\rightarrow T_u = \sigma(W_u[a^{(t-1)} + x^{(t)}] + b_u)$

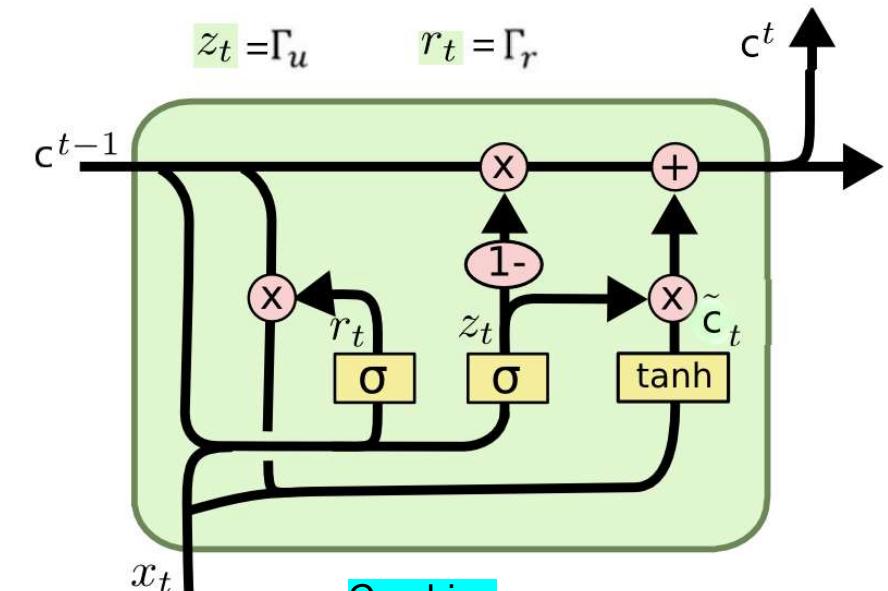
Forget gate $\rightarrow T_f = \sigma(W_f[a^{(t-1)} + x^{(t)}] + b_f)$

Output gate $\rightarrow T_o = \sigma(W_o[a^{(t-1)} + x^{(t)}] + b_o)$

Cell state $\rightarrow c^{(t)} = T_u * \tilde{c}^{(t)} + T_f * c^{(t-1)}$

Activation / output $\rightarrow a^{(t)} = T_o * \tanh(c^{(t)})$

- There are 2 gates (update & reset gates)
- 1 output
- 2 inputs ($x(t)$, $C(t-1)$, $a(t-1)$)



Combine

- Previous \rightarrow how much to reset $C(t-1)$
- Current [$x(t)$]

Candidate cell $\rightarrow \tilde{c}^{(t)} = \tanh(W_c[T_r * c^{(t-1)} + x^{(t)}] + b_c)$

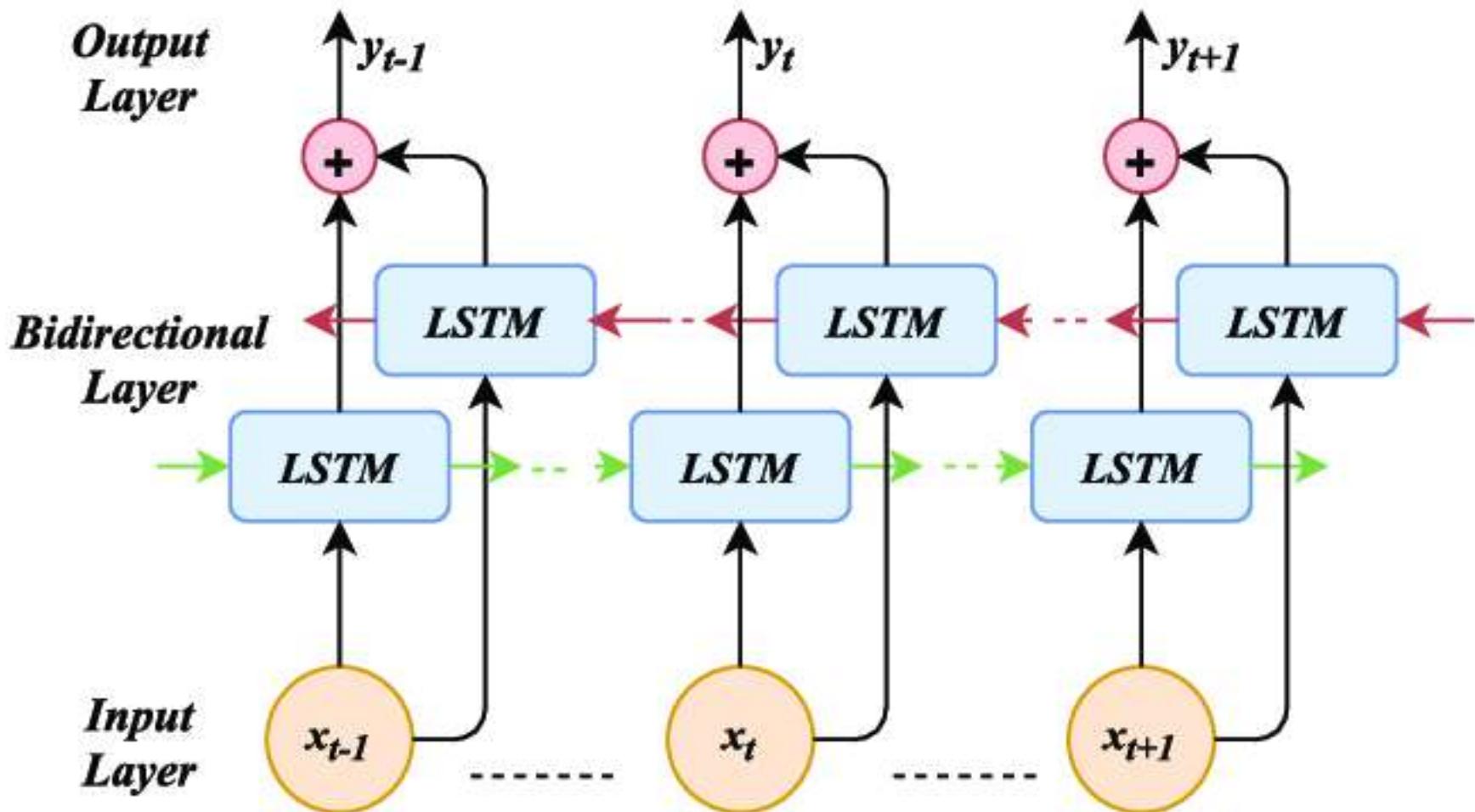
Update gate $\rightarrow T_u = \sigma(W_u[c^{(t-1)} + x^{(t)}] + b_u)$

Reset gate $\rightarrow T_r = \sigma(W_r[c^{(t-1)} + x^{(t)}] + b_r)$

Cell state $\rightarrow c^{(t)} = T_u * \tilde{c}^{(t)} + (1 - T_u) * c^{(t-1)}$

Activation / output $\rightarrow a^{(t)} = c^{(t)}$ *Update g(t) & C(t-1)*

Bidirectional LSTM

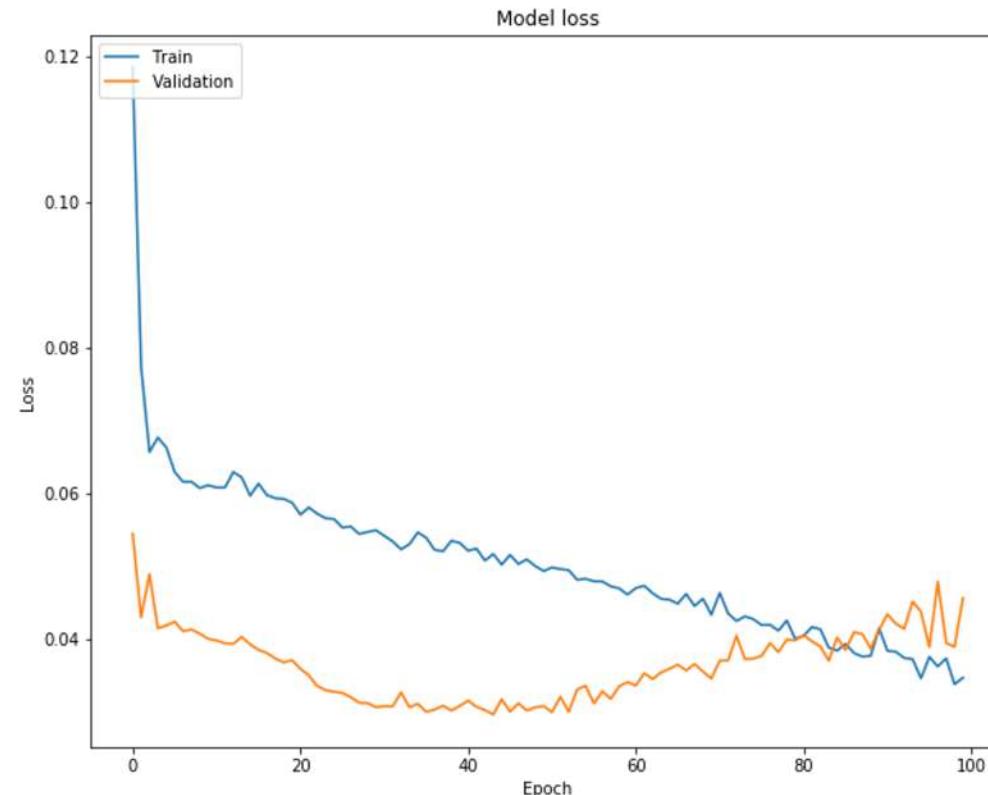


LSTM Parameters Tuning

- Epochs
- Learning Rate
- LSTM Parameters
 - Number of hidden Unit
 - Number of layers
- Loss function

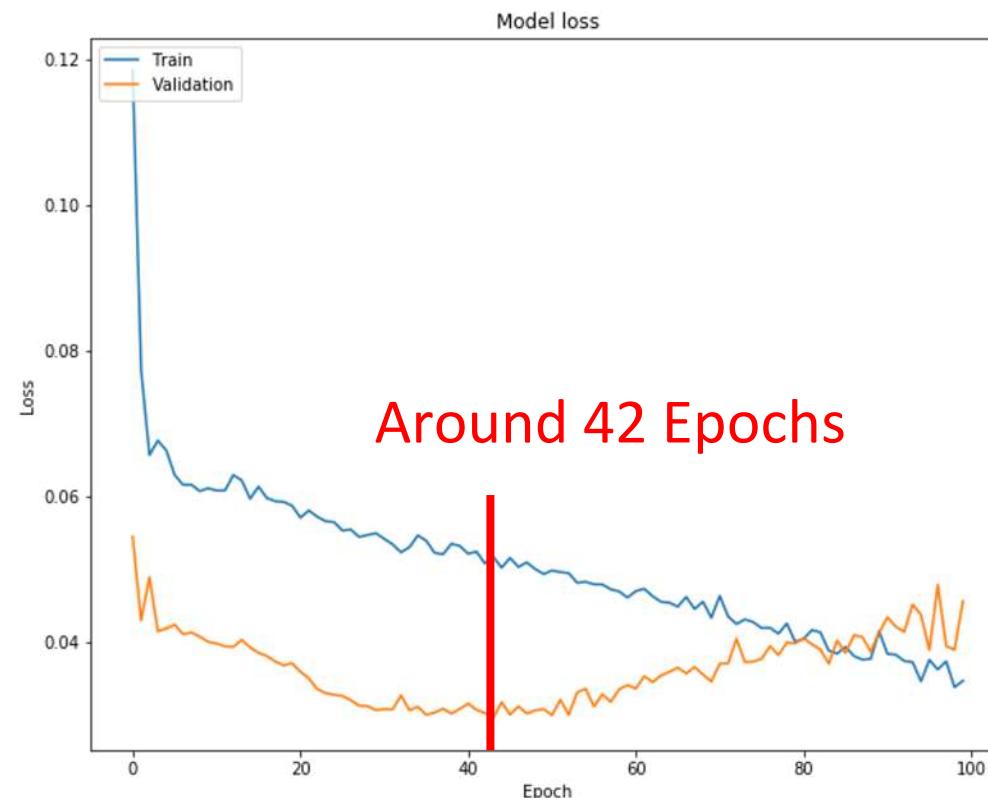
LSTM Parameters Tuning: Epoch

- The single time you see all examples in the dataset.
- Choosing by Plotting Train, Validation loss



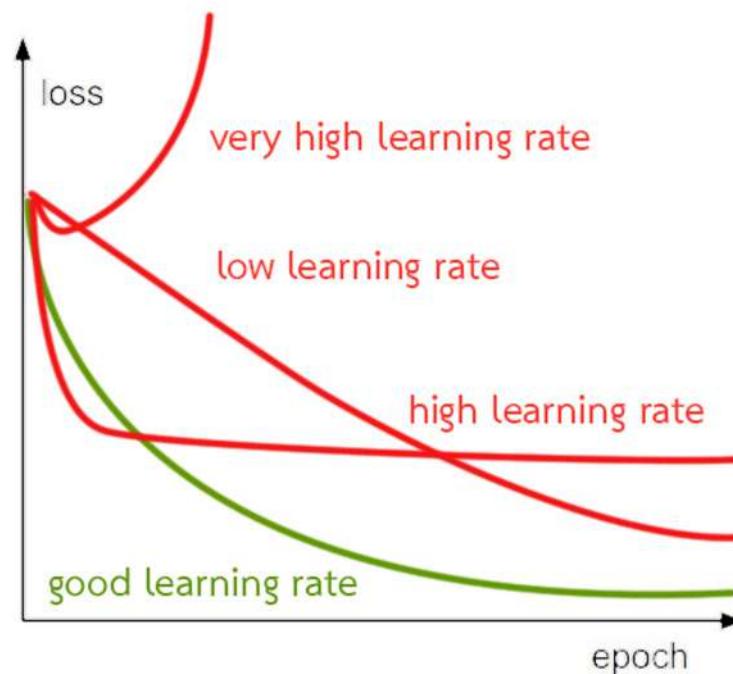
LSTM Parameters Tuning: Epoch (cont.)

- The single time you see all examples in the dataset.
- Choosing by Plotting Train, Validation loss



LSTM Parameters Tuning: Learning Rate

- The learning rate is a hyperparameter that controls how much to change the model in response to the estimated error each time the model weights are updated.



LSTM Parameters Tuning: Loss Function

- We are currently use Mean Square Error

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \tilde{y}_i)^2$$

NLP Lecture2 (2020s2) by Aj.Ekapol

<https://youtu.be/WxiO3wvKhCE>

The image shows a video conference interface with a presentation slide. The slide has a red header 'Dictionary-based drawbacks' with a small '2/1' in the top right corner. Below the header is a bulleted list:

- Cannot handle words outside of the dictionary (Out-of-Vocabulary, OOV words)
- Performs worse than machine-learning-based approach

Below the list is a graph on lined paper. The graph has a vertical axis labeled from 0 to 1000 and a horizontal axis. A blue line starts at approximately (0, 850) and slopes downward to about (1000, 150). A pink circle highlights the start of the line at (0, 850). To the right of the graph is a small table:

Parameter	Value
Optimal size	1000
Freshness	1000
Time frame	1000

At the bottom of the slide, there is a footer with the text 'Majumdar Kumar, Dipan Chaitanya, OPTIMAL SIZE, FRESHNESS AND TIME FRAME FOR VOICE SEARCH CAPABILITY'.



Language Technologies

NVIDIA AI TECHNOLOGY CENTER (NVAITC)

2019



NVAITC COLLABORATORS

Singapore



NATIONAL
RESEARCH
FOUNDATION
PRIME MINISTER'S OFFICE
SINGAPORE



Taiwan



National
Taiwan
University



南科AIROBOT 自造基地

Europe



Hartree Centre
Science & Technology Facilities Council

Thailand



Australia
MONASH
University

China



The University of
Nottingham

UNITED KINGDOM - CHINA - MALAYSIA

Hong Kong



India



IIT Hyderabad

#AI #GTC20

Intelligent End-to-End AI Chatbot with Audio-Driven Facial Animation

39,101 views • May 16, 2020



NVIDIA 689K subscribers



1

During the GTC 2020 virtual keynote, NVIDIA CEO and founder Jensen Huang interacted with Misty, a conversational AI weather chatbot, to demonstrate an end-to-end pipeline to create an AI driven 3D digital avatar. The NVIDIA Jarvis multimodal application framework includes pre-
[SHOW MORE](#)



ASR



NLP



TTS

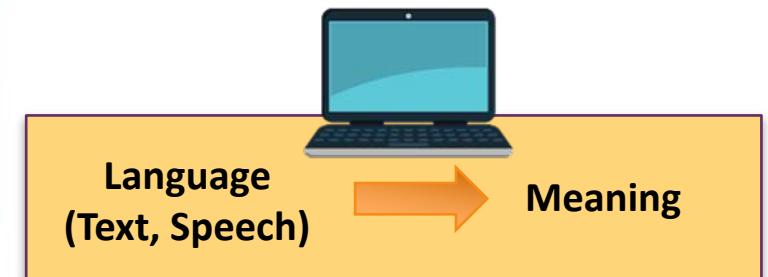
Outlines

- Part1: Overview of NLP & Its Tasks
- Part2: Word Segmentation
- Part3: Text Classification

Part 1: Overview of NLP & Its Tasks

Natural Language Processing (NLP)

NLP is a subfield in AI, where the goal is to bridge the gap between how people communicate and what machines **understand** in order to perform useful tasks, e.g. making appointments, buying things, question answering, etc.



* Not just string matching

Level of understanding in NLP

https://www.tutorialspoint.com/artificial_intelligence/artificial_intelligence_natural_language_processing.htm

Lexical Analysis:

Text→ Paragraphs, Sentences, and Words

Syntactic Analysis (Parsing):

Grammar/Relationship between words

Semantic Analysis:

Exact meaning of the sentence

Discourse Integration:

Meaning of the sentence based on the previous sentence (pronouns)

Pragmatic Analysis:

Actual Meaning based on **the context** and real-world knowledge

Discourse

Semantics

Syntax: Constituents

Syntax: Part of Speech

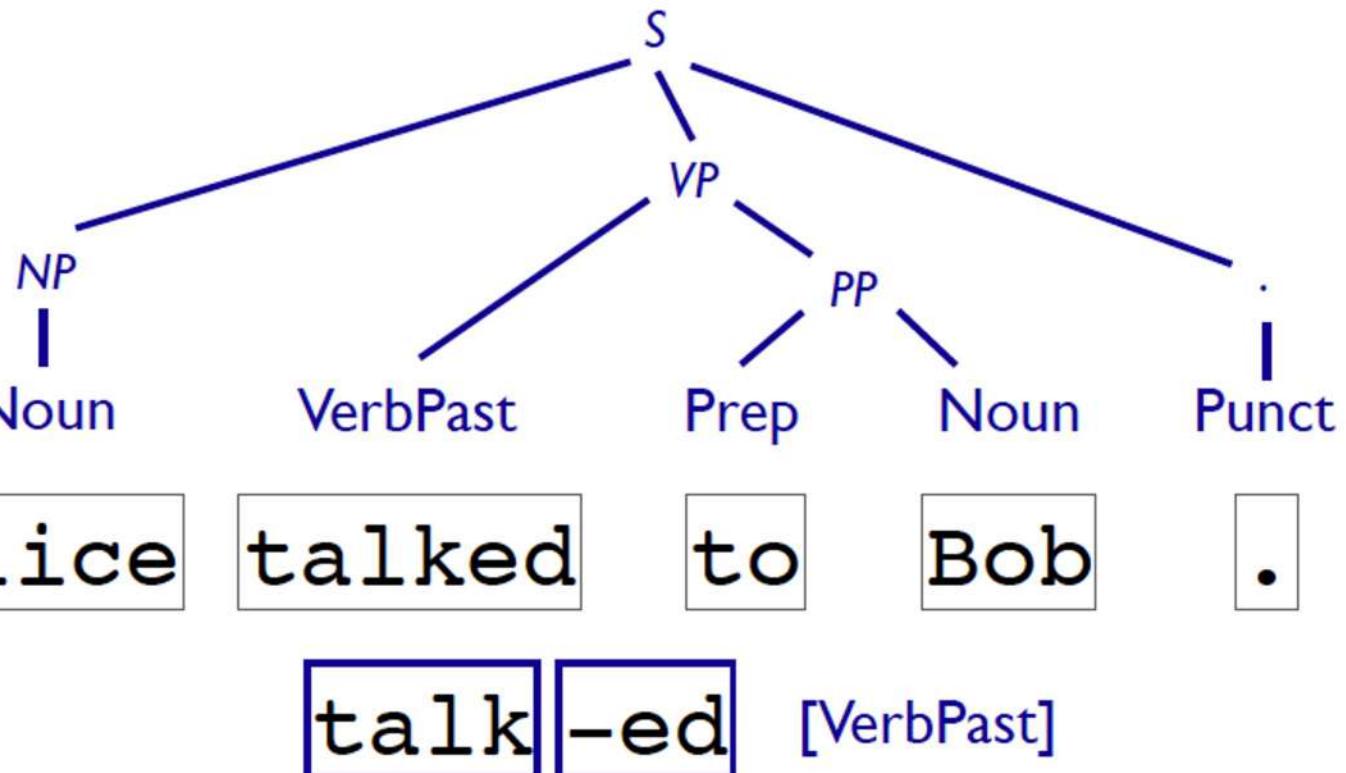
Words

Morphology

Characters

CommunicationEvent(e)
Agent(e, Alice)
Recipient(e, Bob)

SpeakerContext(s)
TemporalBefore(e, s)



Alice talked to Bob.

Tokenization

- Input: ขสมก. เลึง จัดหารรถ
- Output: ขสมก., เลึง, จัดหา, รถ

Part of Speech tagging

- Input: [ขสมก., เลึง, จัดหา, รถ]
- Output: [(ขสมก., NR), (เลึง, VV), (จัดหา, VV), (รถ, NN)]

NER

- Input: [(ขสมก., NR), (เลึง, VV), (จัดหา, VV), (รถ, NN)]
- Output: [(ขสมก., NR, ORG), (เลึง, VV, O), (จัดหา, VV, O), (รถ, NN, O)]

Application

- Word Cloud (Named Entity Only)
- Text Classification
- Etc.

Downstream task

PENN Part Of Speech Tags

- NR – proper noun
 - VV - Main verbs in clauses, verb-form
 - NN – Non-proper noun
- Ref: BEST2010 dataset

Named Entity Tags

- PER –Person
- LOC – Location
- ORG – Organization
- O – Other



Machine Translation (MT)

Google google translate

All Images Maps News Videos More Settings Tools

About 1,180,000,000 results (0.39 seconds)

English ▾ Thai ▾

As the new year gets underway, expert commentators give their view on what 2018 holds in store.

Here are three big themes to watch out for over the next 12 months.

Can the stock market rally go on? The new year has begun with stock markets in the UK and US hitting new record highs.

The Dow Jones Industrial Average rose above 25,000 points for the first time this week, while the broader S&P 500 is also at historic highs.

เป็นปีใหม่ที่กำลังได้รับการแสดงความคิดเห็นของผู้เชี่ยวชาญให้มุมมองของพวากษาเกี่ยวกับสิ่งที่ 2018 เก็บไว้ในร้าน

ต่อไปนี้เป็นหัวข้อใหญ่สามข้อที่ควรระวังในช่วง 12 เดือนข้างหน้า

การซัมมูมตลาดหุ้นสามารถดำเนินต่อไปได้หรือไม่? ปีใหม่เริ่มมีตลาดหุ้นในสหรัฐอาณาจักรและสหราชอาณาจักรที่สูงสุดเป็นประวัติการณ์

ดัชนีเฉลี่ยอุตสาหกรรมดาว โจนส์ปรับตัวสูงขึ้นกว่า 25,000 จุดเป็นครั้งแรกในสัปดาห์นี้ขณะที่ดัชนี S & P 500 ที่ใหญ่ขึ้นก็อยู่ในระดับสูงเป็นประวัติการณ์

Markets, Brexit and Bitcoin: 2018's themes

By Chris Johnston
Business reporter

5 January 2018

f t m e Share



As the new year gets underway, expert commentators give their view on what 2018 holds in store.

<http://www.bbc.com/news/business-42581934>

Question Answering (QA)



IBM Watson wowed the tech industry and a corner of U.S. pop culture with its 2011 win against two of Jeopardy's greatest champions. Here's how IBM pulled it off and a look at what Watson's real career is going to be.

<https://www.techrepublic.com/article/ibm-watson-the-inside-story-of-how-the-jeopardy-winning-supercomputer-was-born-and-what-it-wants-to-do-next/>

Ref: Prof. Regina Barzilay, NLP @MIT

Chatbots

Number of Websites with a Chatbot by Country

Country	Number of Chatbot Sites	Population	Chatbots/1000 people
USA	15,515	325,000,000	0.0477
Australia	1009	24,000,000	0.0420
UK	2667	65,000,000	0.0410
Netherlands	684	17,000,000	0.0402
Canada	1265	36,000,000	0.0351
France	569	66,000,000	0.0086
Germany	476	82,000,000	0.0058
Brazil	399	207,000,000	0.0019
India	592	1,300,000,000	0.0005

Sites Using Chatbots Gathered from NerdyData



<https://research.aimultiple.com/chatbot-stats/>

An advertisement for Botnoi Consulting's Enterprise Chatbot Solution. It features a dark background with a smartphone displaying a chat interface. The text reads: "Enterprise Chatbot Solution" and "White label chatbot solutions powered by Natural Language Processing that provide the best chat experience with a uniquely human touch. Allow your agents to work together seamlessly with a chatbot." At the bottom right is a blue button labeled "Chatbot Solution →".

<https://botnoigroup.com/>

Search/Summarization

Google aquaman

All Images Videos News Maps More Settings Tools

About 164,000,000 results (0.69 seconds)

Showtimes for Aquaman

Aquaman Movie Official Website - In theaters December 21, 2018
<https://www.aquamanmovie.com/> ▾
Aquaman - #AquamanMovie- In theaters December 21st, 2018.

Aquaman (2018) - Rotten Tomatoes
https://www.rottentomatoes.com/m/aquaman_2018/ ▾
★★★★★ Rating: 64% - 298 reviews
Dec 21, 2018 - Critic Consensus: Aquaman swims with its entertainingly ludicrous tide, offering up CGI superhero spectacle that delivers energetic action with ...

Aquaman (film) - Wikipedia
[https://en.wikipedia.org/wiki/Aquaman_\(film\)](https://en.wikipedia.org/wiki/Aquaman_(film)) ▾
Aquaman is a 2018 American superhero film based on the DC Comics character of the same name, and distributed by Warner Bros. Pictures. It is the sixth ...
Amber Heard · James Wan · Ocean Master · Yahya Abdul-Mateen II

Top stories



8:30pm

อควาแมน เจ้าสมุทร

พ.ศ. 2561 · ภาพยนตร์แนวแฟนตาซี/ภาพยนตร์รัตนวิทยา
วิทยาศาสตร์ · 2 ชม. 22 นาที

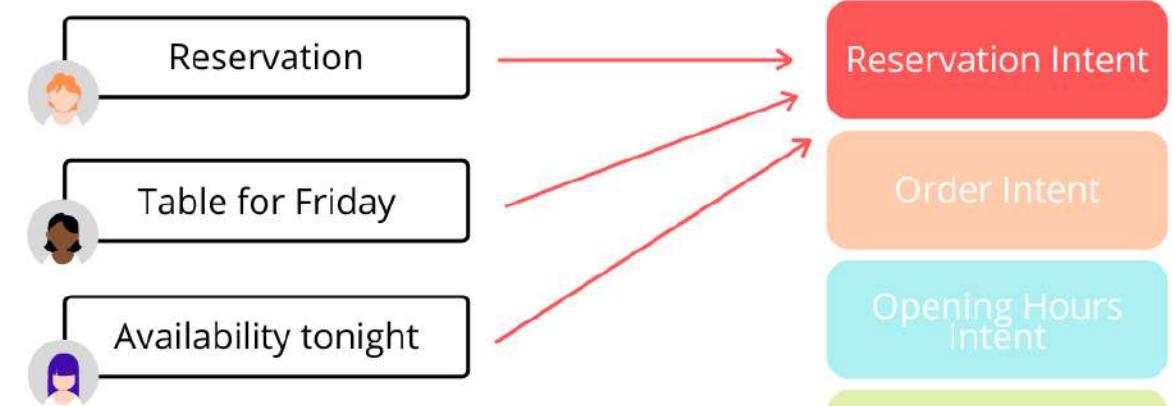
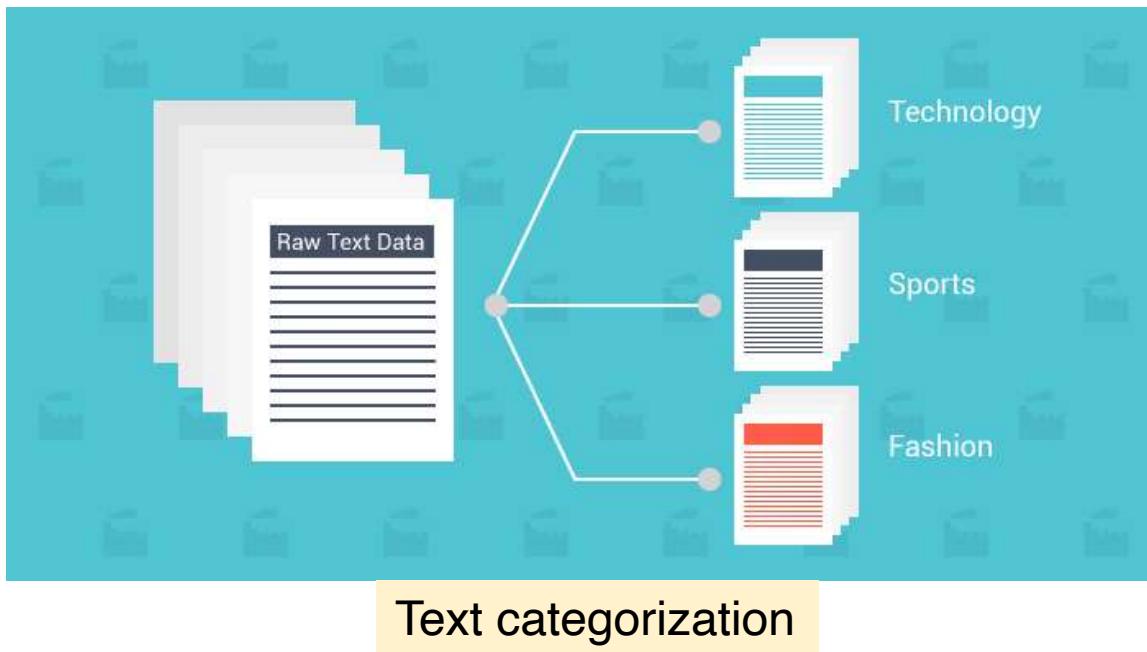
7.6/10 IMDb 64% Rotten Tomatoes 55% Metacritic

94% ชอบภาพยนตร์เรื่องนี้
ผู้ใช้ Google

อควาแมน เจ้าสมุทร เป็นภาพยนตร์ชุดเบอร์ชิริจากประเทศสหราชอาณาจักร ดัดแปลงจากตัวละคร "อควาแมน" ของตีชุดมิกส์ ถูกวางแผนให้เป็นภาพยนตร์ลำดับที่ 6 ในชุดภาพยนตร์จากจักรวาลขยายตัวเดียวกันโดย James Wan เป็นบทโดย David Leslie Johnson-McGoldrick และ Will Beall แต่งเรื่องโดย Wan, Beall, และ Geoff Johns นำแสดงโดย Jason Momoa ...



Text Classification



Intent classification

"I am happy with this water bottle."



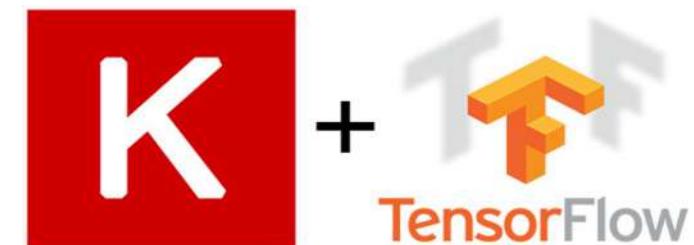
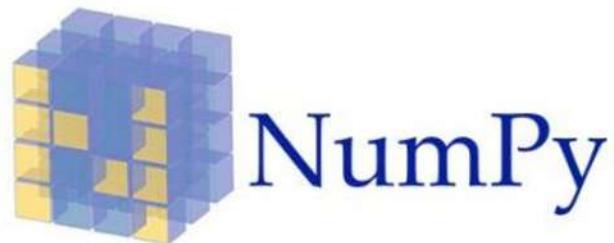
"This is a bad investment."



"I am going to walk today."

Sentiment classification

Demo





HUGGING FACE

On a mission to solve NLP,
one commit at a time.

Star 39,335

text-classification

t: Most downloads ▾

token-classification

question-answering ?

multiple-choice

masked-lm

causal-lm

summarization

translation

conversational

zero-shot-classification

ExBERT

NLP Libraries

NLTK 3.5 documentation

[NEXT](#) | [MODULES](#) | [INDEX](#)

Natural Language Toolkit

NLTK is a leading platform for building Python programs to work with human language data. It provides easy-to-use interfaces to [over 50 corpora and lexical resources](#) such as WordNet, along with a suite of text processing libraries for classification, tokenization, stemming, tagging, parsing, and semantic reasoning, wrappers for industrial-strength NLP libraries, and an active [discussion forum](#).

The screenshot shows the Gensim project page on the Hugging Face platform. At the top, there's a navigation bar with links for Help, Sponsors, Log in, and Register. Below the navigation is a search bar labeled "Search projects" with a magnifying glass icon. A prominent feature is a large green button for "pip install gensim" with a pip icon. To the right of this button is a "Latest version" button with a checkmark and a dropdown menu showing "Released: Apr 1, 2021". Below the main header, the text "Python framework for fast Vector Space Modelling" is displayed. The page is divided into sections: "Navigation" (with "Project description" highlighted), "Project description" (containing a "build failing" badge and a "wheel yes" badge), and "Release history" (with a link to "gensim-4.0.1"). A brief description at the bottom states: "Gensim is a Python library for topic modelling, document indexing and similarity retrieval with large corpora. Target audience is the natural language processing (NLP) and information retrieval (IR) community."

NLP Libraries for Thai

<https://thainlp.org/pythainlp/docs/2.0/index.html>

The screenshot shows the PyThaiNLP documentation homepage. The sidebar on the left contains links for 'NOTES' (Command Line, Getting Started, Installation, From PyThaiNLP 1.7 to PyThaiNLP 2.0) and 'PACKAGE REFERENCE' (pythainlp.corpus, pythainlp.soundex, pythainlp.spell, pythainlp.summarize, pythainlp.tag, pythainlp.tokenize, pythainlp.tools, pythainlp.transliterate, pythainlp.ulmfit, pythainlp.util, pythainlp.word_vector). The main content area includes a 'PyThaiNLP documentation' section, a 'Project description' section featuring the PyThaiNLP logo (a stylized blue and yellow bird-like character), and a 'Package reference:' section listing all the package components.

Docs » PyThaiNLP documentation

[View page source](#)

PyThaiNLP documentation

PyThaiNLP is a Python library for natural language processing (NLP) of Thai language.

Notes

- Command Line
- Getting Started
- Installation
- From PyThaiNLP 1.7 to PyThaiNLP 2.0

Package reference:

- pythainlp.corpus
- pythainlp.soundex
- pythainlp.spell
- pythainlp.summarize
- pythainlp.tag
- pythainlp.tokenize
- pythainlp.tools
- pythainlp.transliterate
- pythainlp.ulmfit
- pythainlp.util
- pythainlp.word_vector



Project description



PyThaiNLP is a Python library for Thai natural language processing. The library provides functions like word tokenization, part-of-speech tagging, transliteration, soundex generation, and spell checking.

WangchanBERTa โมเดลประมวลผลภาษาไทยที่ใหญ่และก้าวหน้าที่สุดในขณะนี้

VISTEC-depa AI Research Institute of Thailand [Follow](#)
Jan 24 · 5 min read

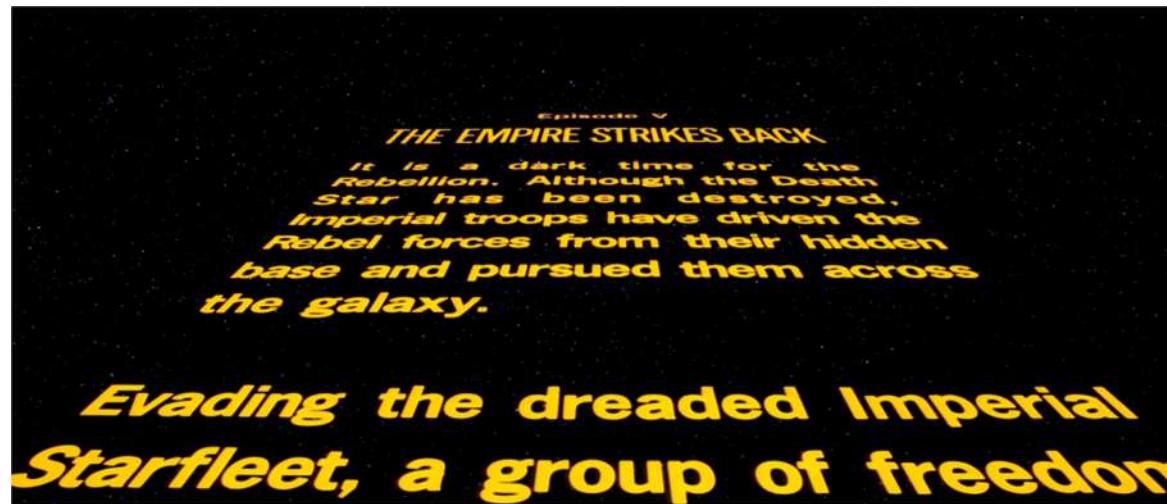


เปิดให้ทุกคนใช้ฟรีโดย AIResearch.in.th และ VISTEC ภายใต้ลิขสิทธิ์อนุญาต CC-BY-SA 4.0

Part2: Word Segmentation

The need for segmentation

- Text as a stream of characters



- We need a way to understand the meaning of text
 - Break into words (assign meaning to word) ← Tokenization
 - Break into sentences (put word meanings back to sentence meaning)

Tokenization - Thai

- Thai has **no** space between **words**
- Thai has **no** clear **sentence boundaries**

- เริ่มจากชั้นของสังคมแรกรสุด สามารถพัฒนา ทำการปิดล้อม-รุกรานดาวนานุ ส่งผลต่อมาขยายเป็นสังคมโคลนอันมีฝ่ายแบ่งแยกินแดนเป็นผู้ชั้นในสังคม หมายโน่นล้มฝ่ายสาธารณรัฐ นานาไปกับเรื่องราวพัฒนาการของของเด็กน้อยคนหนึ่งผู้มีพลังสกิดแรงมาก ชาวดาวหะเหลறายทากอินนาม "อนาคต สถาบันล็อกเกอร์" ผู้ถูกคาดการณ์ว่าคือผู้ถูกเลือกในตานานของเจ้า หลังจากสังคมยุทธการดาวนานุ อนาคต ก็ได้รับฝึกฝนในวิถีเจ้า โดย อ.เจ้า "โอบีวัน" แต่ เมื่ออนาคตได้เป็นหนุ่มก็ต้นแรกกูเจ้าโดยแอบมีความสัมพันธ์ซึ่งสาวลับๆกับ ราชินี "อมิตาลา" แห่งดาวนานุ จนเหอตั้งครรภ์ลูกแฟด ... และอนาคตก็อยู่ๆถูกกลืนเข้าสู่ด้านมืดของพลังจันกลายเป็นชีวลอร์ด 'ไดจายา "ดาร์ธ เวเดอร์" ภายใต้การโน้มน้าวชั้นนำของ ชีวลอร์ดลีกลับนาม "ดาร์ธ ชีเดียส" ชี้่งเผยแพร่ในตอนท้ายว่า ชีเดียส ไม่ใช่ใครที่ไหน แต่คือท่าน "พลพาร์ทีน" สมุหนายกผู้นำสูงสุดของฝ่ายสาธารณรัฐเสียเอง และแท้จริงก็ยัง เป็นผู้นำลับชักใฝ่ฝ่ายแบ่งแยกินแดนด้วยอีกต่างหาก ... สังคมจะลงที่ฝ่ายสาธารณรัฐพ่ายแพ้ล่มสลาย อมิตาลา ก็ตายหลังคลอดลูกแฟด ทั้งเหล่า อัศวินเจ้าได้ถูกฆ่าขาดล้างสิ้นแบบไม่ทันตั้งตัว เหลือแต่ อ.โอบีวัน กับ อ.โยดา ต้องลี้ภัยหลบหนีช่อนตัว โดยลูกแฟดของอนาคตได้ถูกส่งไปสู่ที่หลบ ช่อนลับเช่นกัน ... และแล้ว พลพาร์ทีน ก็กินรวบทั้งกระดาน เปลี่ยนการปกครองจาก ระบบสาธารณรัฐเติมไปเป็น ระบบเพดีจการจักรพรรดิชีวแทน ตั้ง ตนเป็นจักรพรรดิปักษ์กลางซึ่งทั้งปวง โดยมี ดาร์ธ เวเดอร์ เป็นขุนพลชีวลอร์ดเคียงข้างนับแต่นั้นมา

[Search docs](#)

NOTES

[Command Line](#)[Getting Started](#)[Installation](#)[From PyThaiNLP 1.7 to PyThaiNLP 2.0](#)

PACKAGE REFERENCE:

[pythainlp.corpus](#)[pythainlp.soundex](#)[pythainlp.spell](#)[pythainlp.summarize](#)[pythainlp.tag](#)

pythainlp.tokenize

[Modules](#)

Tokenization Engines

[pythainlp.tools](#)[pythainlp.transliterate](#)[pythainlp.ulmfit](#)[pythainlp.util](#)

`pythainlp.tokenize.word_tokenize(text: str, custom_dict: marisa_trie.Trie = None, engine: str = 'newmm', keep_whitespace: bool = True) → List[str]` [\[source\]](#)

This function tokenizes running text into words.

Parameters

- `text (str)` – text to be tokenized
- `engine (str)` – name of the tokenizer to be used
- `custom_dict (marisa_trie.Trie)` – marisa dictionary trie
- `keep_whitespace (bool)` – True to keep whitespaces, a common mark for end of phrase in Thai. Otherwise, whitespaces are omitted.

Returns

list of words

Return type

`list[str]`

Options for engine

- `newmm` (default) - dictionary-based, Maximum Matching + Thai Character Cluster
- `longest` - dictionary-based, Longest Matching
- `deepcut` - wrapper for `deepcut`, language-model-based
- `icu` - wrapper for ICU (International Components for Unicode, using PyICU), dictionary-based
- `ulmfit` - for `thai2fit`

1) Longest Matching

- Scan a sentence from left to right
 - Keep finding a word from the starting point, until no word matched (**longest**), then move to the next point
 - Backtrack if current segmentation leads to an un-segmentable chunk
-
- ป้ายกลับรถ Start scanning with “ป” as the starting point
 - ป้ายกลับรถ Keep scanning ...
 - ป้าย/**กลับรถ** No more words start with “ป้าย”, move to the next point
 - ...
 - ป้าย/**กลับ**/**รถ**

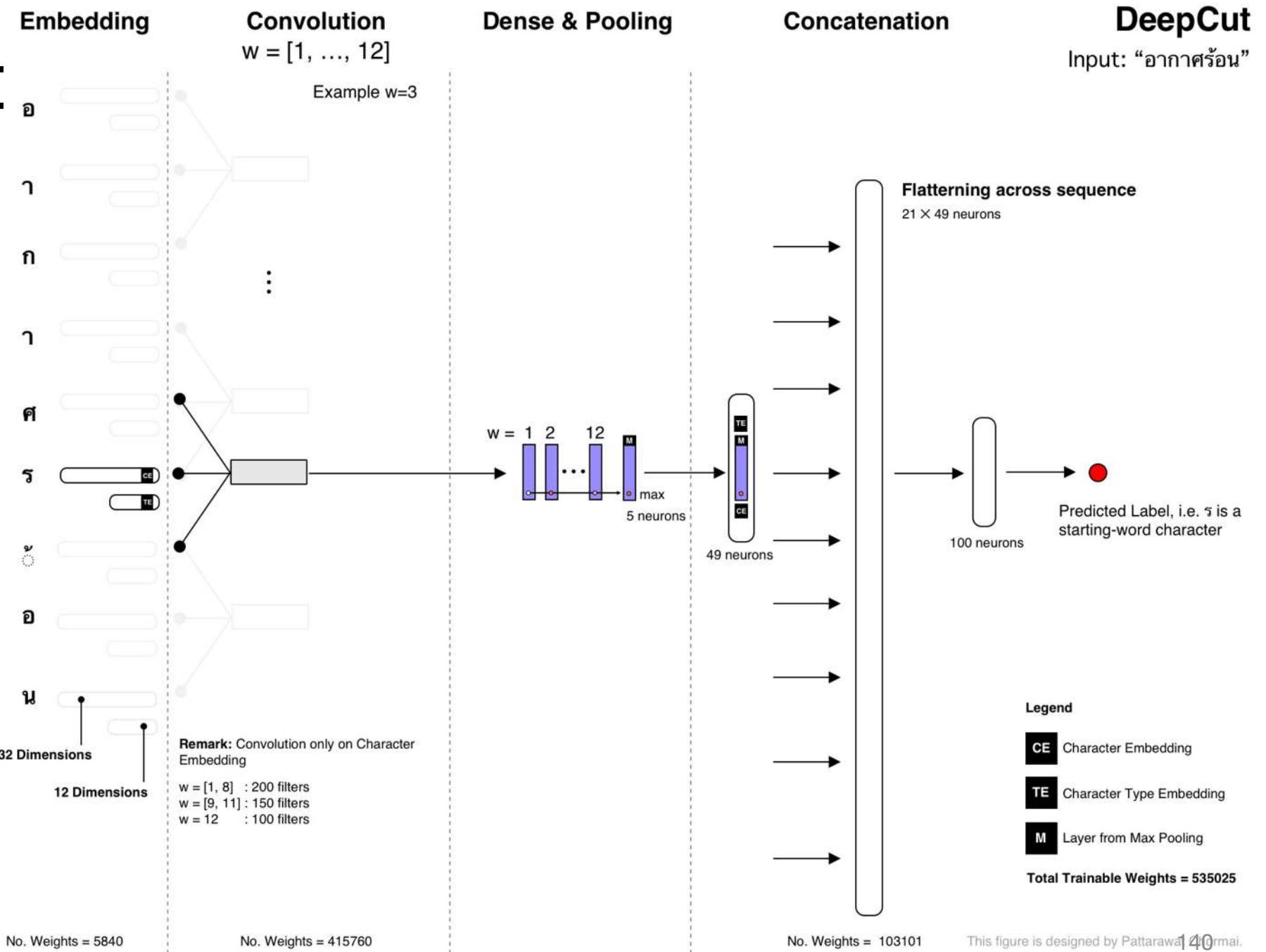
2) Maximal Matching

- Generate all possible segmentations
- Select the segmentations with the **fewest words**

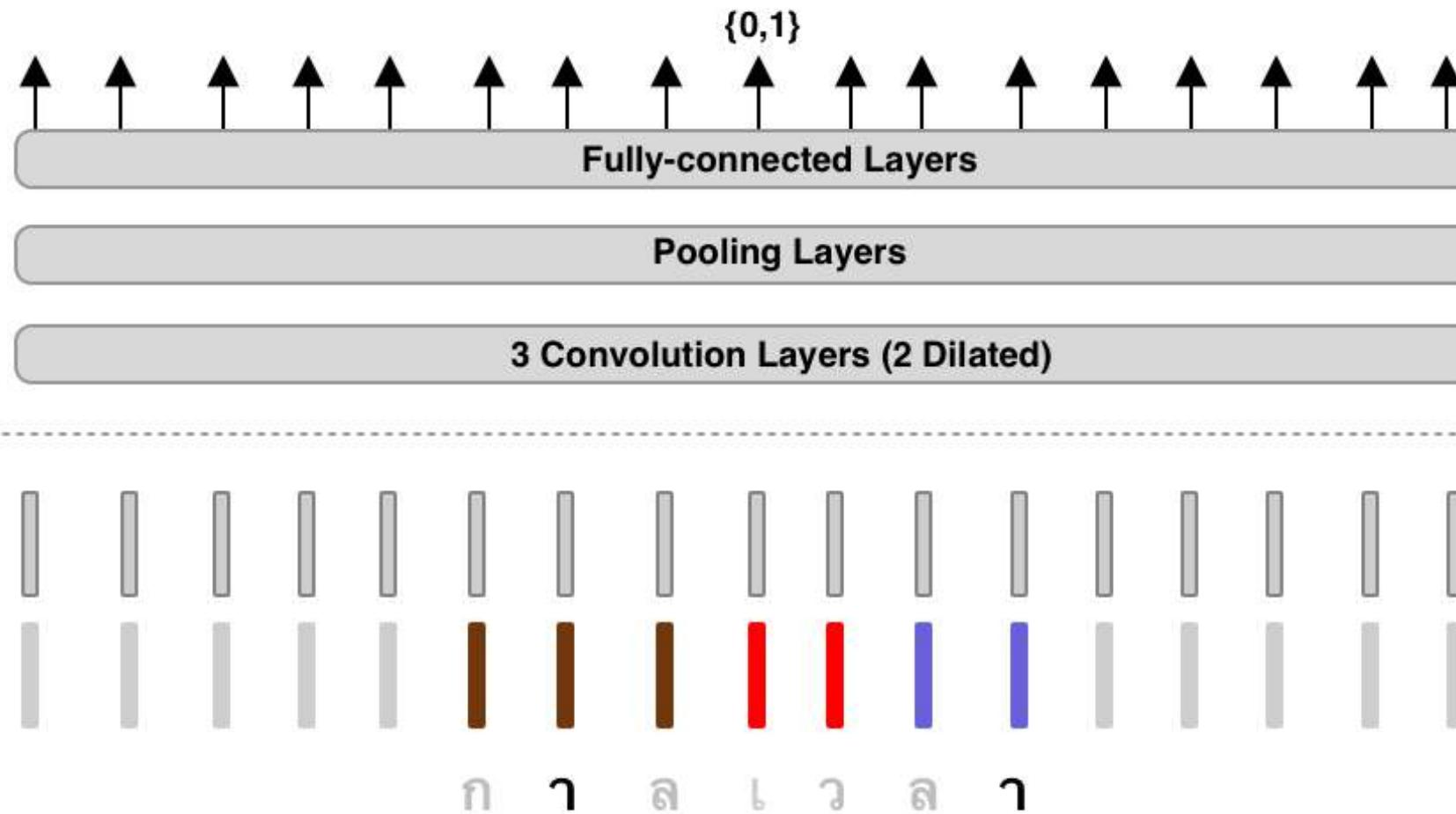


Haruechayasak, Choochart, Sarawoot Kongyoung, and Matthew Dailey. "A comparative study on thai word segmentation approaches." Electrical Engineering/Electronics, Computer, Telecommunications and Information Technology, 2008. ECTI-CON 2008. 5th International Conference on. Vol. 1. IEEE, 2008.

3) DeepCut



4) AttaCut

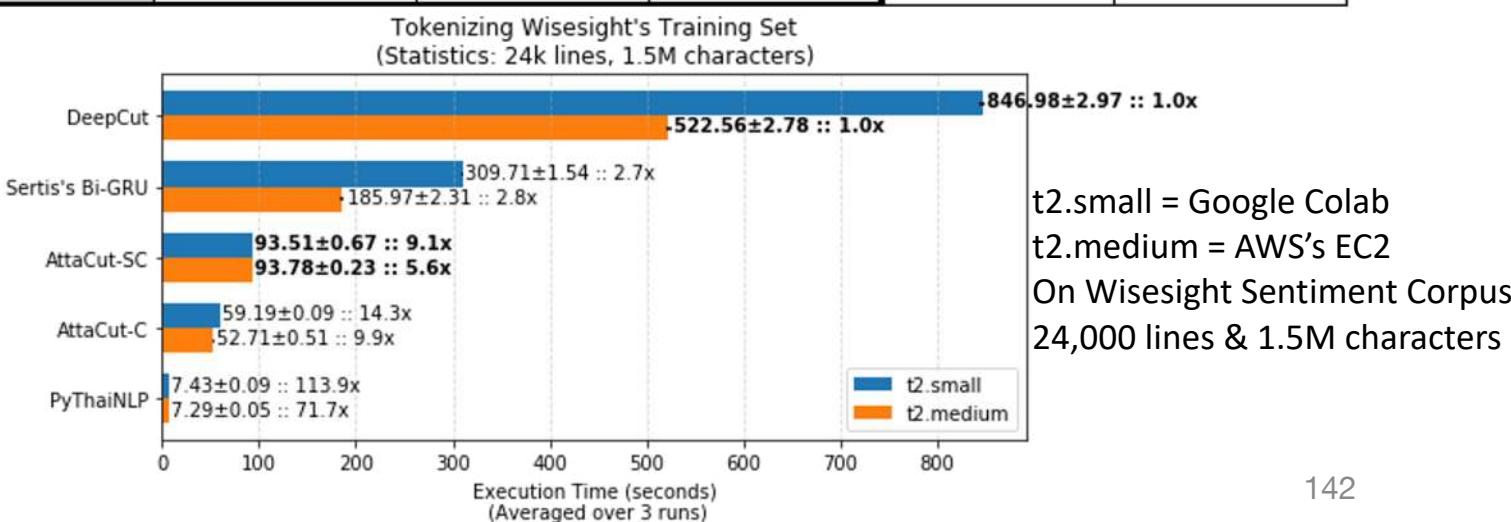
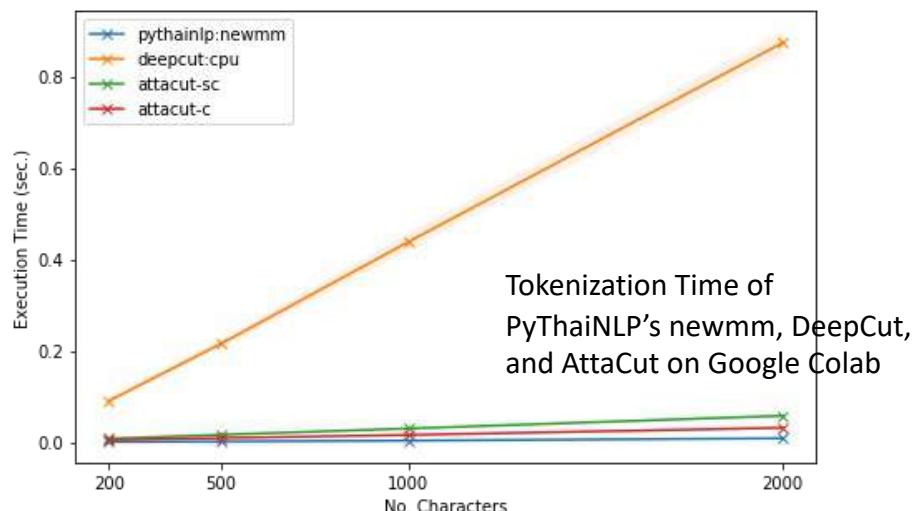


TL;DR: 3-Layer Dilated CNN on syllable and character features. It's **6x faster** than DeepCut (SOTA) while its WL-f1 on BEST is **91%**, only 2% lower.

<https://github.com/PyThaiNLP/attacut>

NECTEC. BEST: Benchmark for Enhancing the Standard of Thai language processing, 2010.

		Others		Ours	
Last Updated: 29/08/2019	PyThaiNLP newmm	Sertis Bi-GRU	DeepCut	AttaCut-C	AttaCut-SC
BEST Validation Set					
Character-Level	precision	0.94±0.11	0.95±0.10	0.99±0.05	0.97±0.07
	recall	0.83±0.09	0.99±0.02	0.99±0.03	0.98±0.04
	f1	0.88±0.08	0.97±0.07	0.99±0.04	0.98±0.05
Word-Level	precision	0.73±0.16	0.91±0.14	0.97±0.07	0.94±0.10
	recall	0.65±0.16	0.94±0.10	0.97±0.07	0.94±0.09
	f1	0.68±0.15	0.93±0.12	0.97±0.07	0.94±0.10
BEST Test Set					
Character-Level	precision	0.91±0.15	0.92±0.11	0.96±0.08	0.94±0.10
	recall	0.85±0.09	0.98±0.04	0.98±0.04	0.98±0.04
	f1	0.86±0.11	0.95±0.08	0.97±0.06	0.96±0.07
Word-Level	precision	0.70±0.19	0.85±0.18	0.92±0.14	0.88±0.17
	recall	0.64±0.18	0.90±0.14	0.93±0.12	0.91±0.14
	f1	0.67±0.19	0.87±0.16	0.93±0.13	0.89±0.16



Part3: Text Classification

(Common) Text classification pipeline

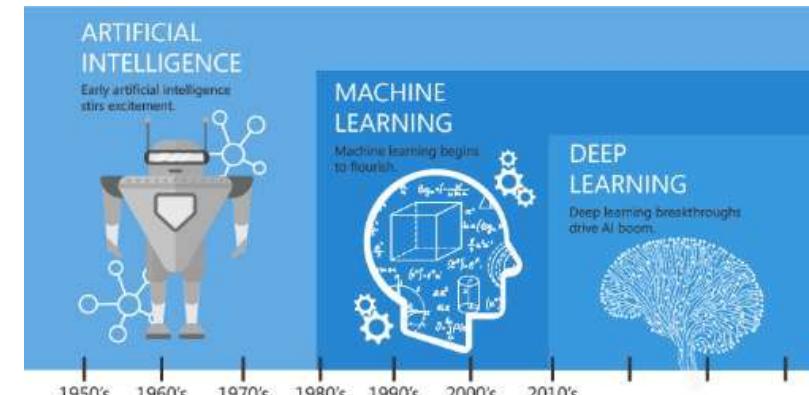
1) Data preparation



2) Document representation

Comments	Good	Like	Hate	Sentiment
Tweet1	7	8	0	😊
Tweet2	1	0	10	😢
Tweet3	2	9	1	😊

3) Supervised learning model



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

Truevoice-intent (action 7 classes)

truevoice-intent: destination

We benchmark **truevoice-intent** by using `destination` as target and construct a 7-class multi-class classification. The performance is measured by micro-averaged and macro-averaged accuracy and F1 score. Codes can be run to confirm performance at this [notebook](#). We also provide performance metrics by class in the notebook.

```
1 # Show first 10 lines of the data
2 !head clean-phone-data-for-students.csv
```

```
1 display(data_df['cleaned_labels'].unique())

array(['enquire', 'report', 'cancel', 'buy', 'activate', 'request',
       'garbage', 'change'], dtype=object)
```

Sentence Utterance,Action,Object

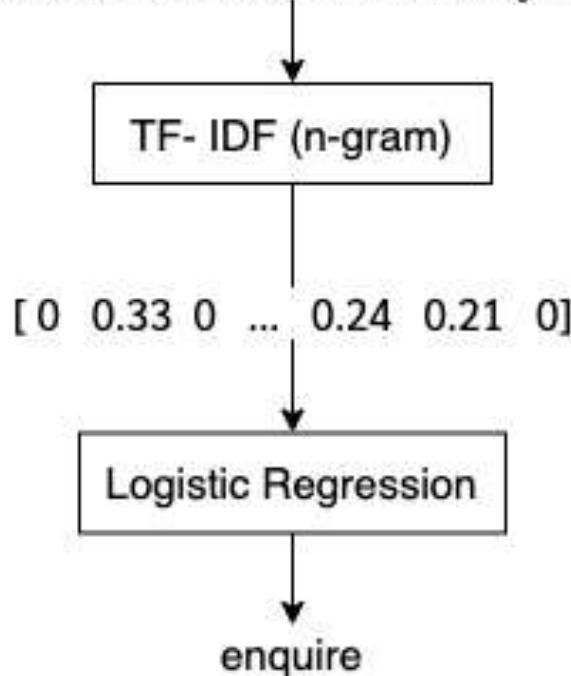
<PHONE_NUMBER_REMOVED> ผ่านไปจ่ายเงินที่ Counter Services เค้าเข็ต 3276.25 บาท เมื่อวานที่ผ่านมาเช็คที่ศูนย์บอกมียอด 3057.79 บาท,enquire,payment
internet ยังความเร็วอยู่เท่าไหร่ ครับ,enquire,package
ตะกี้ไปชำระค่าบริการไปแล้ว แต่ยังใช้งานไม่ได้ ครับ,report,suspend
พี่ค่ะยังใช้ internet ไม่ได้เลยค่ะ เป็นเครื่อง โน้ตบุ๊ก,enquire,internet
อาโทล คะ พอดีว่าเมื่อวานเปิดซิมทรูมูฟ แต่พักโทรศัพท์ไม่ได้ค่ะ แต่เล่นเน็ตได้ค่ะ,report,phone_issues
*2222 ใช้งานยังไง ขอรายละเอียดการสมัครหน่อย,enquire,service
<PHONE_NUMBER_REMOVED> เคยมีช่างมาซ่อมที่บ้าน แล้วโทรศัพท์ใช้งานไม่ได้ครับ,enquire,nonTrueMove
<PHONE_NUMBER_REMOVED> ต้องค่าบริการเท่าไหร่ครับ,enquire,balance
<PHONE_NUMBER_REMOVED> อินเตอร์เน็ตไฟ Adsl ไม่มีสัญญาณครับ,enquire,nonTrueMove

	Sentence	Utterance	Action	Object
0	<PHONE_NUMBER_REMOVED>	ผ่านไปจ่ายเงินที่ Counter...	enquire	payment
1		internet ยังความเร็วอยู่เท่าไหร่ คร...	enquire	package
2		ตะกี้ไปชำระค่าบริการไปแล้ว แต่ยัง...	report	suspend
3		พี่ค่ะยังใช้ internet ไม่ได้เลยค่ะ เป็น...	enquire	internet
4		อาโทล คะ พอดีว่าเมื่อวานเปิดซิมทรูมู...	report	phone_issues

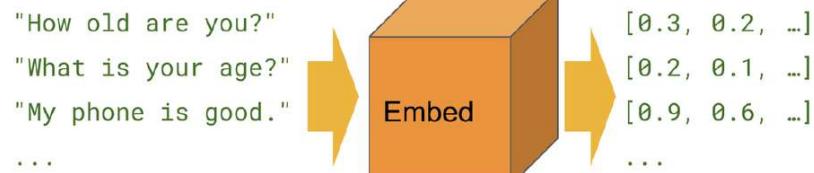
Text Classification Approach

1) TF-IDF + Classifier

ให้รหัสไวไฟมา ใส่เท่าไหรก็เด้งไม่ถูกต้อง



2) USE + Classifier



3) BERT



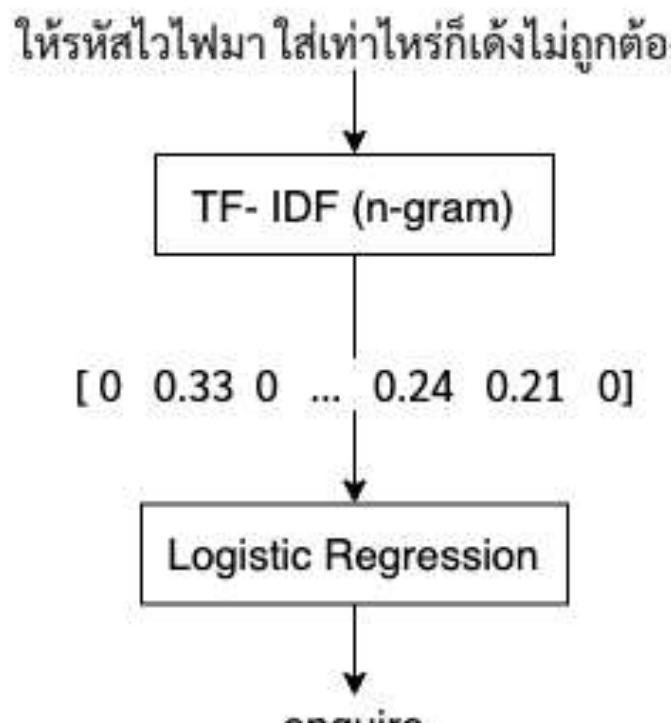
**The AI community
building the future.**

Build, train and deploy state of the art models powered by the reference open source in natural language processing.

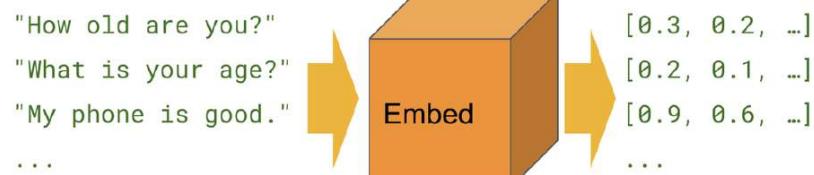
Star 47,792

Text Classification Approach

1) TF-IDF + Classifier



2) USE + Classifier



3) BERT



**The AI community
building the future.**

Build, train and deploy state of the art models powered by the reference open source in natural language processing.

Star 47,792

Sparse representation: Term Frequency (TF)

- Each row represents a word in the vocabulary and term-document matrix
- Each column represents a document.

vocabulary	As You Like It	Twelfth Night	Julius Caesar	Henry V
battle	1	1	8	15
soldier	2	2	12	36
fool	37	58	1	5
clown	5	117	0	0

Figure 15.1 The term-document matrix for four words in four Shakespeare plays. Each cell contains the number of times the (row) word occurs in the (column) document.

Sparse representation: TF-IDF

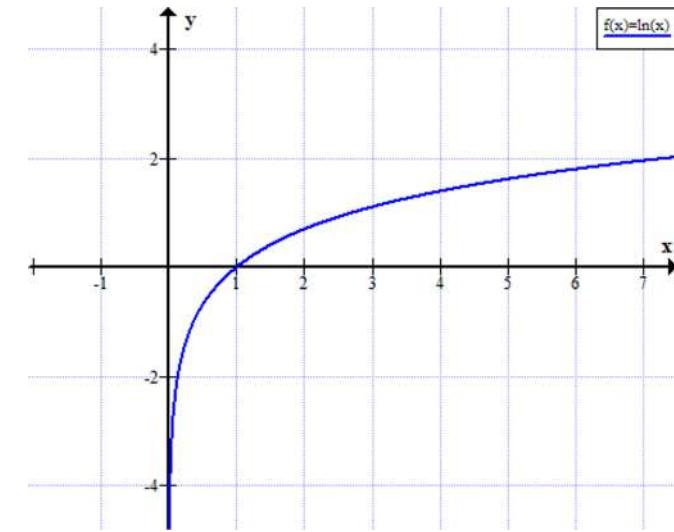
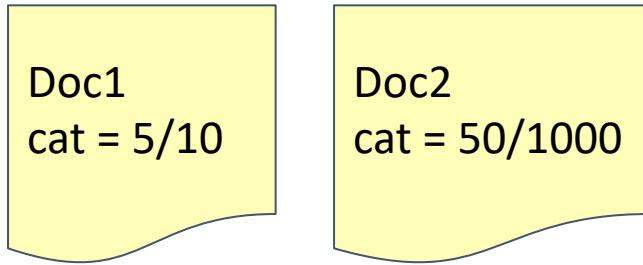
Need for normalization in TF

- Term Frequency (TF) – per each document

$$TF(w) = \frac{\text{Frequency of word } w \text{ in a document}}{\text{Total number of words in the document}}$$

- Inverse Document Frequency (IDF) – per corpus (all documents)

$$IDF(w) = \log_e \left(\frac{\text{Total number of documents}}{\text{Number of documents that contain word } w} \right)$$



penalty score
i.e., a, an, the

- TF-IDF

$$TFIDF(w) = TF(w) * IDF(w)$$

```
1 # TF-IDF
2 tfidf = TfidfVectorizer(
3     ngram_range=(1,2),           # Use unigram and bigram
4     tokenizer=word_tokenize,      # Use `word_tokenize` method from pythainlp for tokenizer
5     min_df=2,                   # The word found less than three times in dataset is ignored
6     max_df=0.9,                 # The word found more than 90% of entries is ignore
7     use_idf=True,
8     smooth_idf=True,
9     sublinear_tf=True
10 )
11 # Logistic regression
12 model = LogisticRegression(C=4, max_iter=300, random_state=42)
```

Sparse representation: TF-IDF (cont.)

TF				
	As You Like It	Twelfth Night	Julius Caesar	Henry V
battle	1	1	8	15
soldier	2	2	12	36
fool	37	58	1	5
clown	5	117	0	0

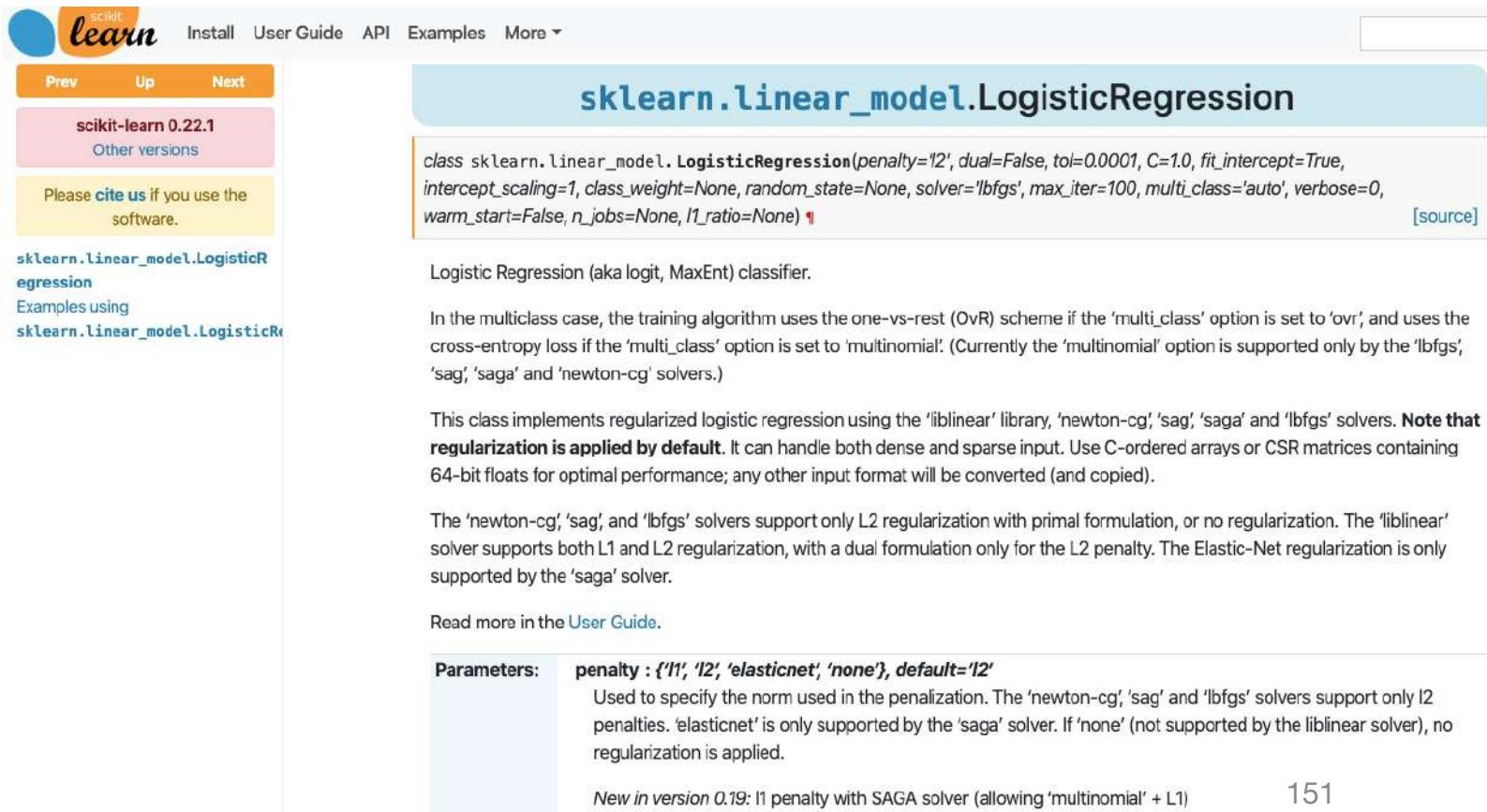
↓

TF-IDF				
	As You Like It	Twelfth Night	Julius Caesar	Henry V
battle	0.074	0	0.22	0.28
good	0	0	0	0
fool	0.019	0.021	0.0036	0.0083
wit	0.049	0.044	0.018	0.022

What classifier?

- Any classifier you like
- k-NN
- Naïve Bayes
- Logistic regression
- SVM
- Neural networks

```
1 # TF-IDF
2 tfidf = TfidfVectorizer(
3     ngram_range=(1,2),           # Use unigram and bigram
4     tokenizer=word_tokenize,    # Use `word_tokenize` method from pythainlp for tokenizer
5     min_df=2,                  # The word found less than three times in dataset is ignored
6     max_df=0.9,                # The word found more than 90% of entries is ignore
7     use_idf=True,
8     smooth_idf=True,
9     sublinear_tf=True
10 )
11 # Logistic regression
12 model = LogisticRegression(C=4, max_iter=300, random_state=42)
```



The screenshot shows the scikit-learn documentation page for the `sklearn.linear_model.LogisticRegression` class. The top navigation bar includes links for Install, User Guide, API, Examples, and More. A sidebar on the left provides links to Prev, Up, Next, scikit-learn 0.22.1, and Other versions. It also encourages users to cite the software. The main content area is titled "sklearn.linear_model.LogisticRegression" and contains the class definition:

```
class sklearn.linear_model.LogisticRegression(penalty='l2', dual=False, tol=0.0001, C=1.0, fit_intercept=True, intercept_scaling=1, class_weight=None, random_state=None, solver='lbfgs', max_iter=100, multi_class='auto', verbose=0, warm_start=False, n_jobs=None, l1_ratio=None) ¶
```

[source]

The text explains that Logistic Regression (aka logit, MaxEnt) classifier is used. It details the one-vs-rest (OvR) scheme for multiclass cases and the cross-entropy loss for multinomial cases. The class implements regularized logistic regression using various solvers like 'liblinear', 'newton-cg', 'sag', 'saga', and 'lbfgs'. It supports both L1 and L2 regularization.

The 'newton-cg', 'sag', and 'lbfgs' solvers support only L2 regularization with primal formulation, or no regularization. The 'liblinear' solver supports both L1 and L2 regularization, with a dual formulation only for the L2 penalty. The Elastic-Net regularization is only supported by the 'saga' solver.

Read more in the [User Guide](#).

Parameters: `penalty : {'l1', 'l2', 'elasticnet', 'none'}, default='l2'`

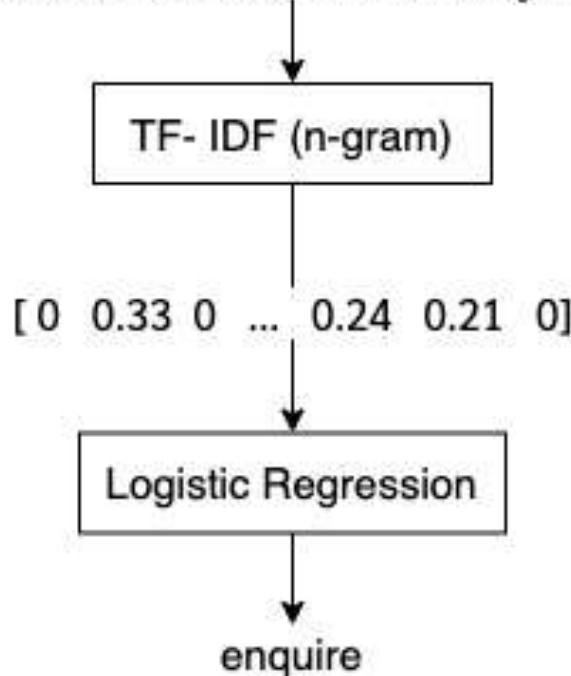
Used to specify the norm used in the penalization. The 'newton-cg', 'sag' and 'lbfgs' solvers support only L2 penalties. 'elasticnet' is only supported by the 'saga' solver. If 'none' (not supported by the liblinear solver), no regularization is applied.

New in version 0.19: l1 penalty with SAGA solver (allowing 'multinomial' + L1)

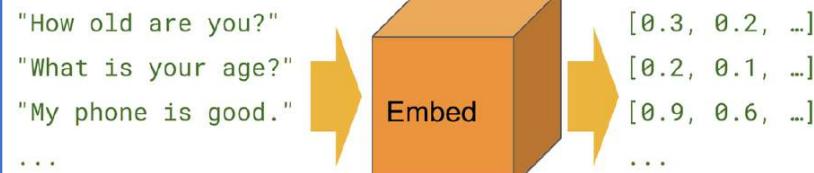
Text Classification Approach

1) TF-IDF + Classifier

ให้รหัสไวไฟมา ใส่เท่าไหรก็เด้งไม่ถูกต้อง



2) USE + Classifier



3) BERT



**The AI community
building the future.**

Build, train and deploy state of the art models powered by the reference open source in natural language processing.

Star 47,792

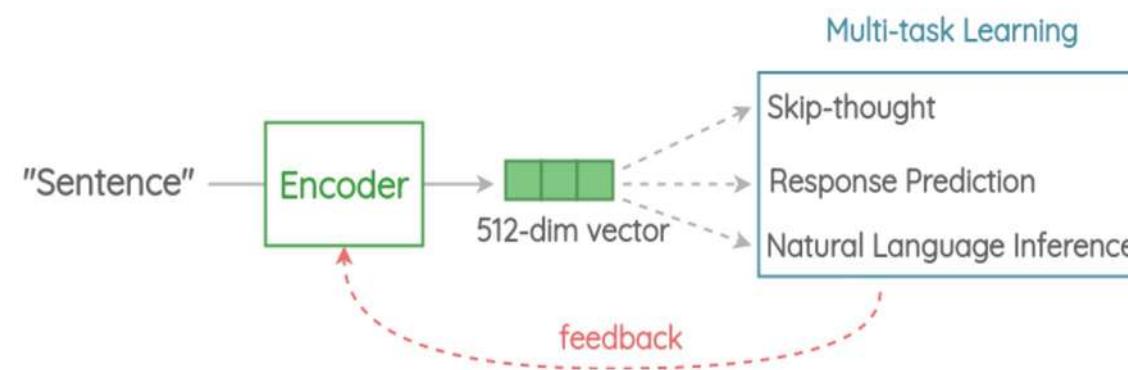
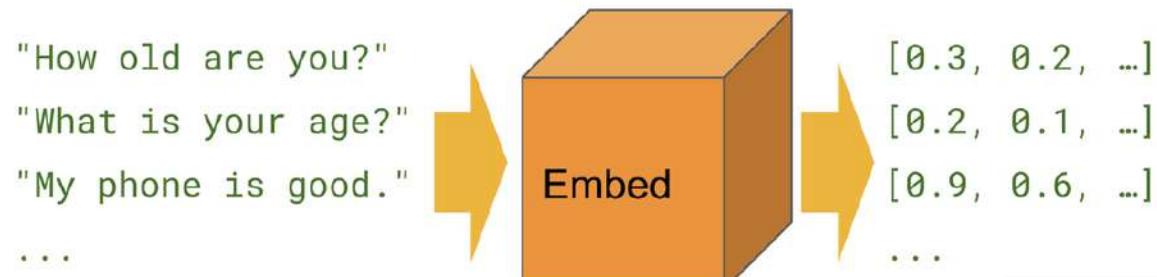
Universal Sentence Encoder (USE)

A model focusing on [sentence representation](#)

Use [sentencepiece tokenization](#)

Pre-trained then used anywhere

Based on (1) DAN (lite version) or (2) Transformer



A terminal window demonstrating the USE model. The code uses the TensorFlow Hub API to process a list of strings and prints the resulting 512-dimensional vector and its type. The output shows the vector size and type as expected.

```
1 test_str = "\u0E2D\u0E22\u0E32\u0E01\u0E40\u0E23\u0E1B"
2 input = [ test_str ] # List of strings
3 output = use_model(input)
4 print( vector size: , output.shape)
5 print( Object type: , type(output))
```

test_str: "อย่างเรียน NLP จังเลยจ้า"

Vector size: (1, 512)
Object type: <class 'tensorflow.python.framework.ops.EagerTensor'>

Official implementation with pretrained weights

<https://tfhub.dev/google/collections/universal-sentence-encoder/1>

<https://ai.googleblog.com/2018/05/advances-in-semantic-textual-similarity.html>

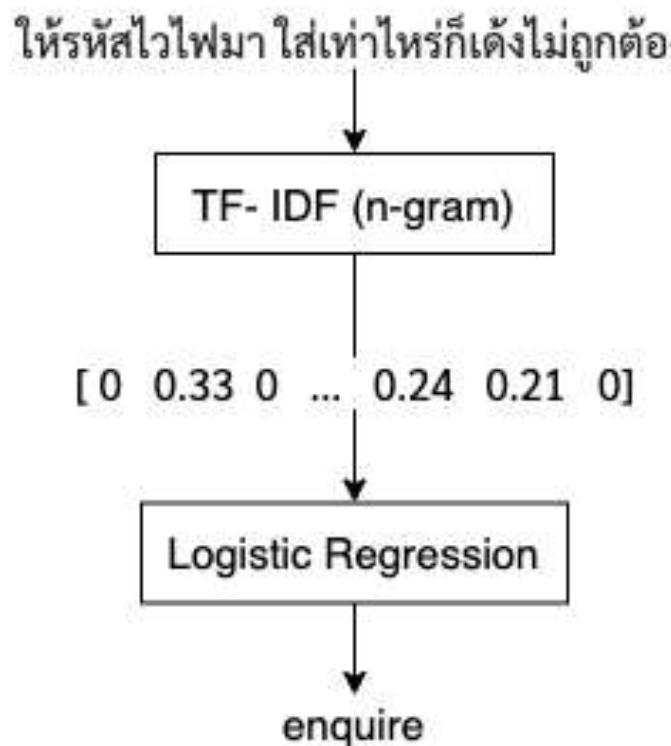
Download-ables

Model	Comments
universal-sentence-encoder	
universal-sentence-encoder-large	
universal-sentence-encoder-lite	
universal-sentence-encoder-qa	Question answering
universal-sentence-encoder-multilingual	16 languages
universal-sentence-encoder-multilingual-large	16 languages
universal-sentence-encoder-multilingual-qa	16 languages , Question answering

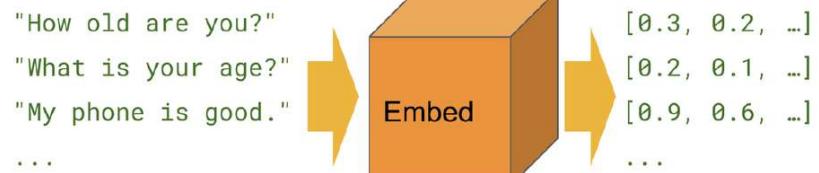
<https://tfhub.dev/google/collections/universal-sentence-encoder/1>

Text Classification Approach

1) TF-IDF + Classifier



2) USE + Classifier



3) BERT



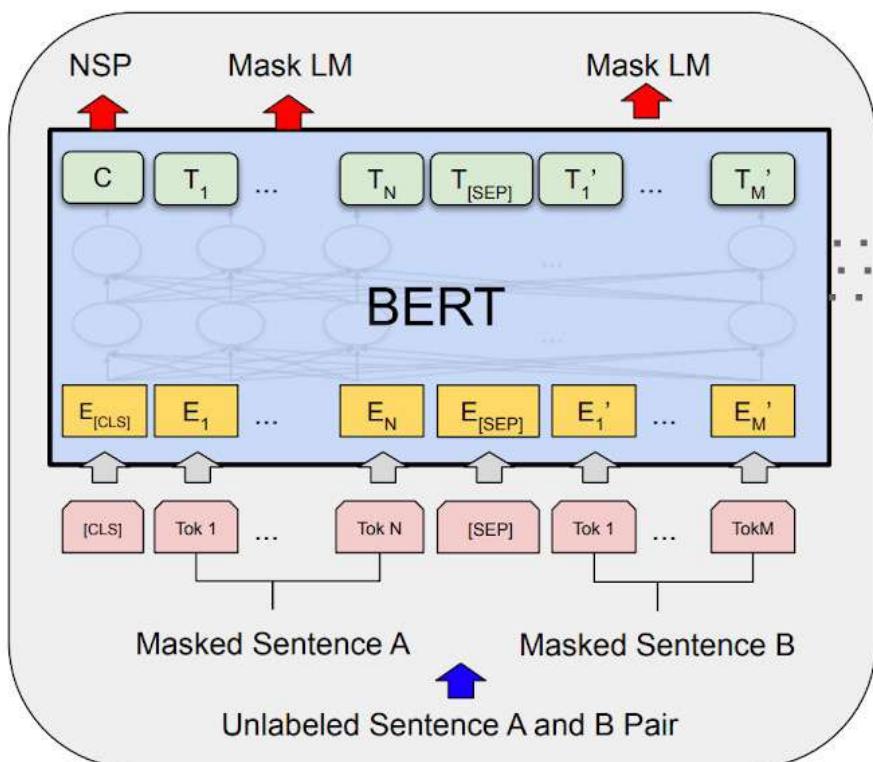
**The AI community
building the future.**

Build, train and deploy state of the art models powered by the reference open source in natural language processing.

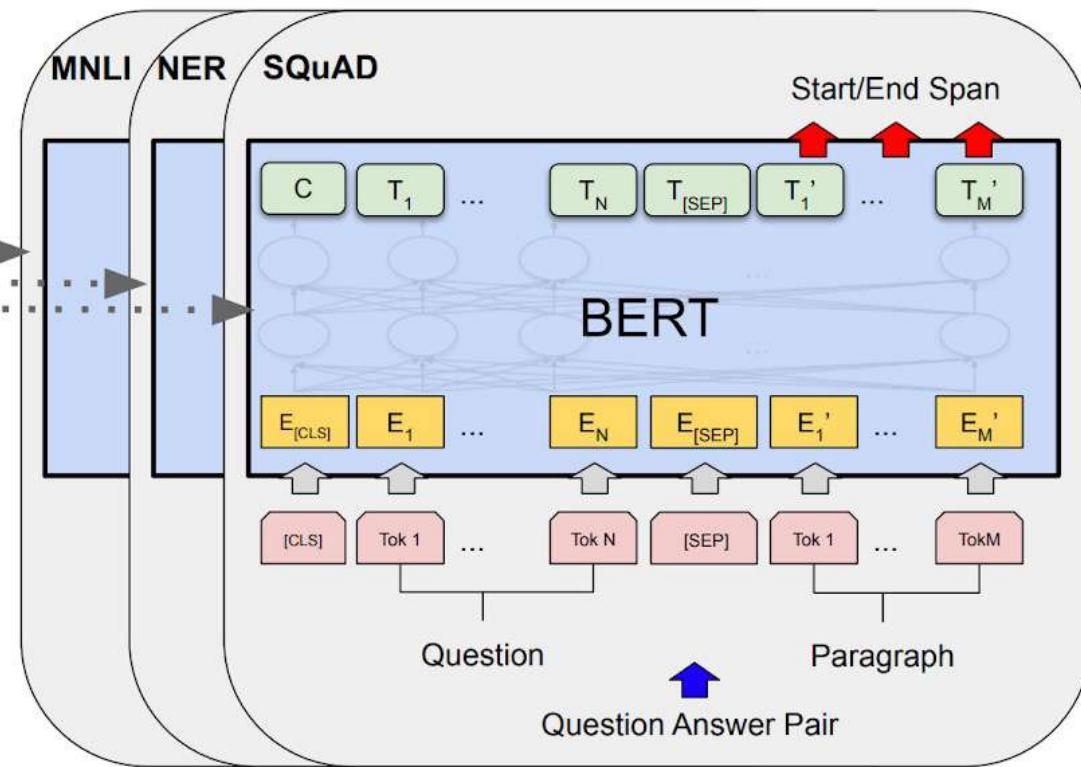
Star 47,792

BERT

- Pretrained language model based on **transformers**
 - Can be used in many NLP tasks



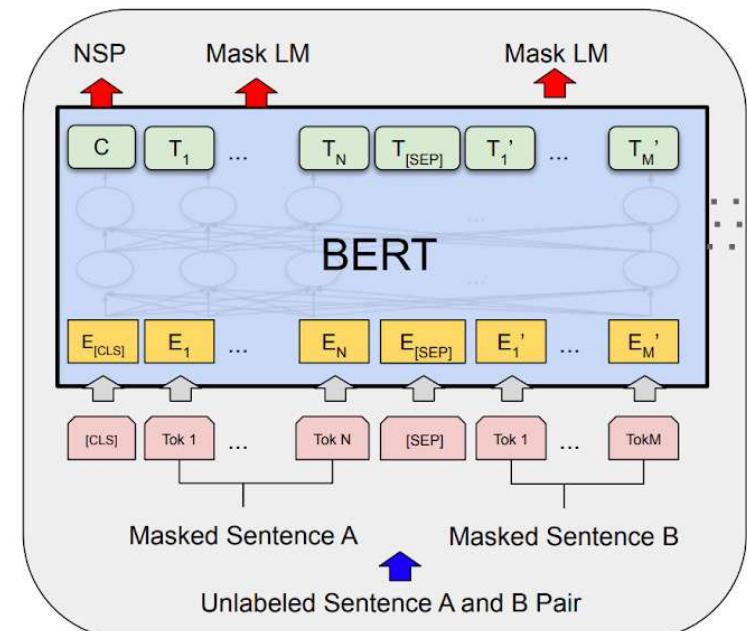
Pre-training



Fine-Tuning

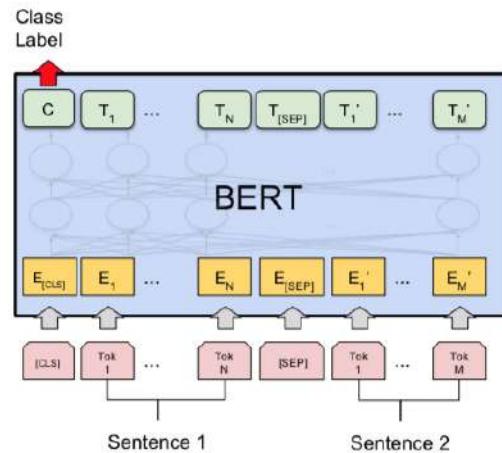
Pre-training BERT

- Predicting **masked words** in a sentence
 - The quick brown fox jumps over the [MASK]
 - Variants: predict correct word or not, predict swapped words, etc.
- **Next sentence prediction**
 - A: The cat is scared. B: It hides under the table.
 - A: The apple is on the table. B: It always rain.
 - Variants: sentence order prediction

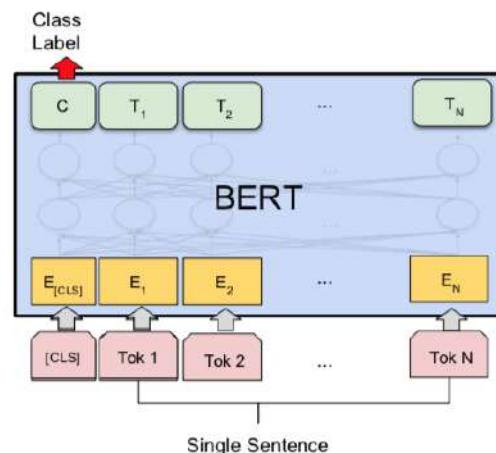


Pre-training

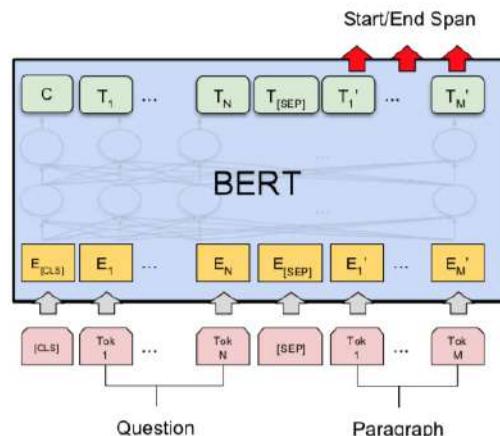
Downstream tasks with BERT



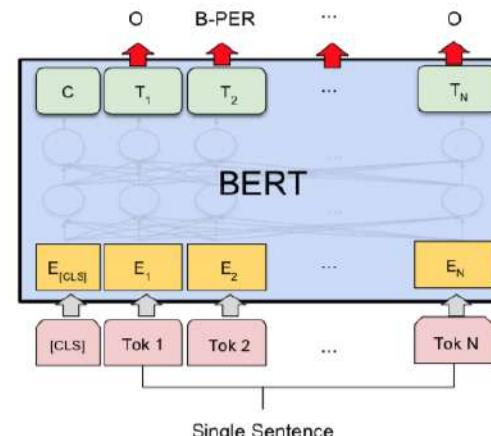
(a) Sentence Pair Classification Tasks:
MNLI, QQP, QNLI, STS-B, MRPC,
RTE, SWAG



(b) Single Sentence Classification Tasks:
SST-2, CoLA



(c) Question Answering Tasks:
SQuAD v1.1



(d) Single Sentence Tagging Tasks:
CoNLL-2003 NER

Huggingface

- An opensource library for transformer-related models
- Has datasets, models, scripts, deployment solutions
- New official online course
<https://huggingface.co/course/chapter1>



The AI community building the future.

Build, train and deploy state of the art models powered by the reference open source in natural language processing.

Star 47,792

Thai pre-trained BERTs

- WangchanBERTa
- XLM-RoBERTa
- MT5

<https://huggingface.co/airesearch/wangchanberta-base-att-spm-uncased>

https://huggingface.co/transformers/model_doc/xlmroberta.html

https://huggingface.co/transformers/model_doc/mt5.html

Results

- WangchanBERTa is the winner.
- Bidirectional LSTM is the second rank, but it is the slowest one.
- Fine-tuning can improve a performance.
- USE is pretty easy to use and can show a promising performance.
- NB-SVM is a strong traditional ML model.

#	Model	Accuracy	macro F1	speed (ms)
Model1	n-grams + NB	0.825	0.739	2.4
Model2	TF-IDF + Logistic	0.864	0.760	2.9
Model3	NB-SVM	0.858	0.764 4	17.2
Model3.1	LTW2V + CNN	0.854	0.725	58.3
Model3.2	LTW2V + Bidirectional LSTM	0.870	0.779 2	63.3
Model3.3	LTW2V + DAN	0.736	0.472	
Model4	LTW2V + USE	0.870	0.766 3	23.9
Model3.1 (f)	LTW2V + CNN (fine-tune)	0.849	0.741	
Model3.2 (f)	LTW2V + Bidirectional LSTM (fine-tune)	0.872	0.791 2	
Model5	WangchanBERTa	0.891	0.814 1	10.3