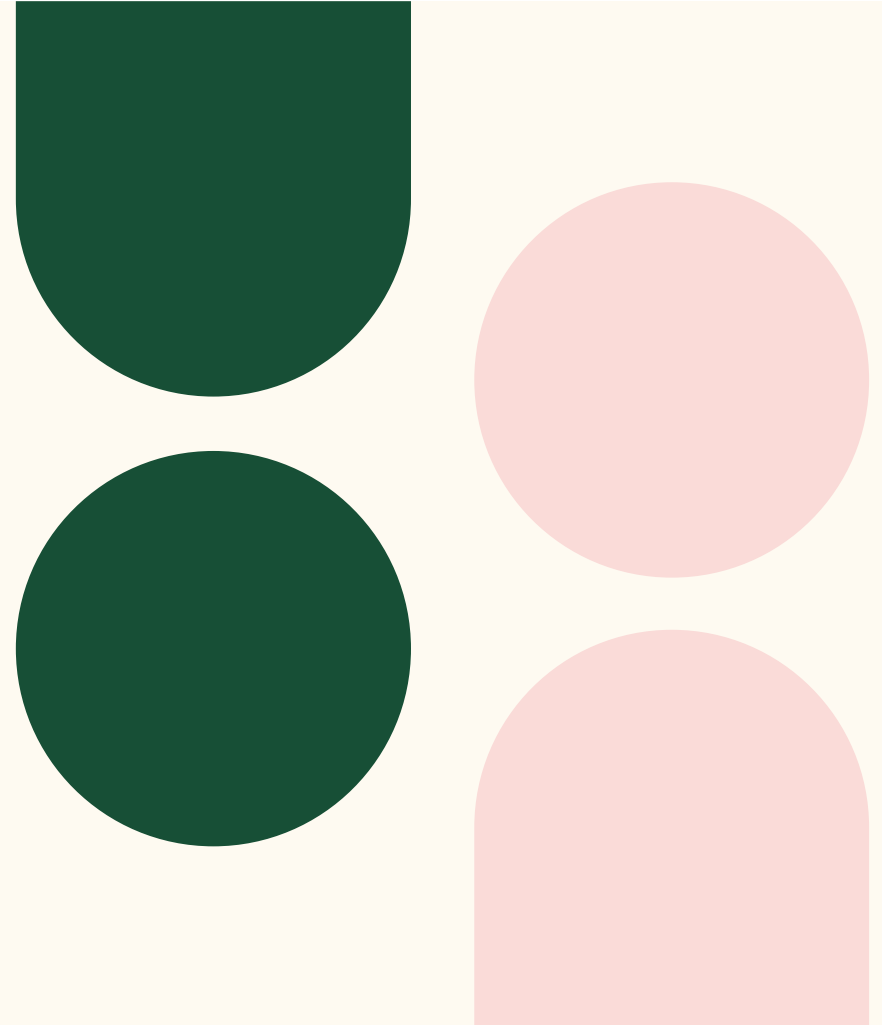


딥러닝 CNN (Convolutional Neural Net)



CONTENTS

01

CNN 구조

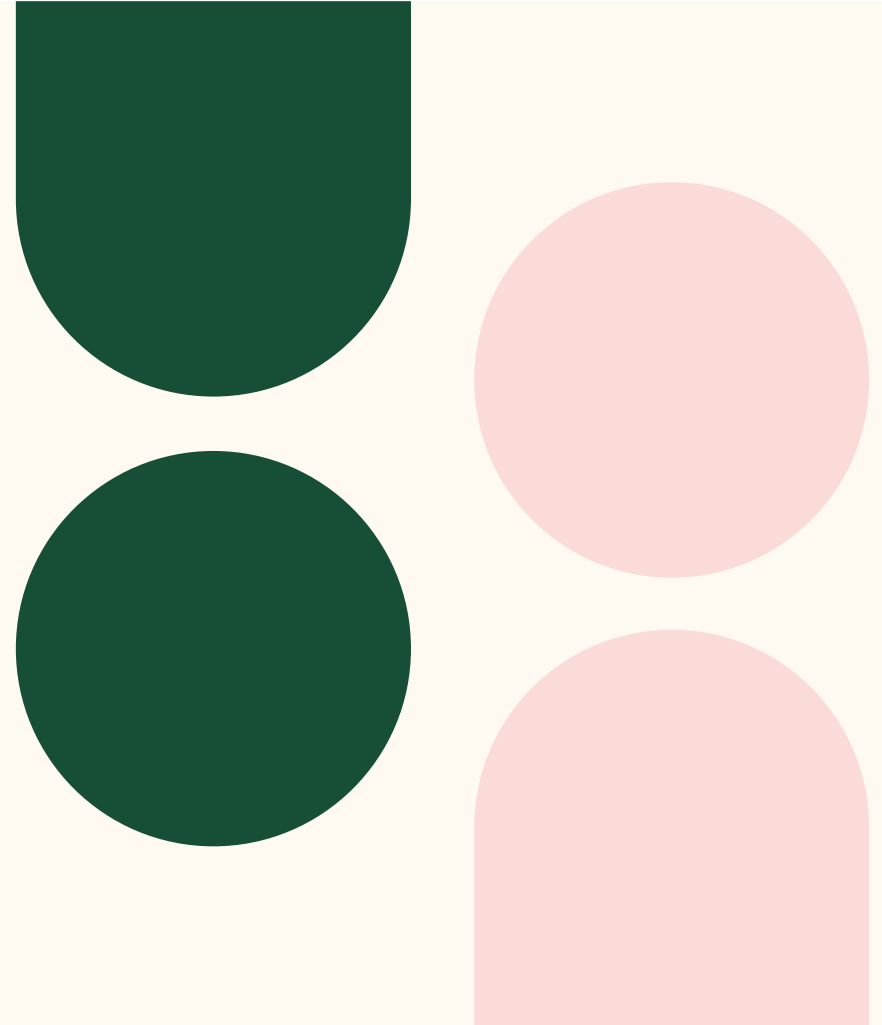
02

CNN 구현 설명 : 합성곱 연산

03

CNN 구현 설명 : Keras CNN Layer

01 CNN 구조



➤ CNN (Convolutional Neural Network)

A BIT OF HISTORY

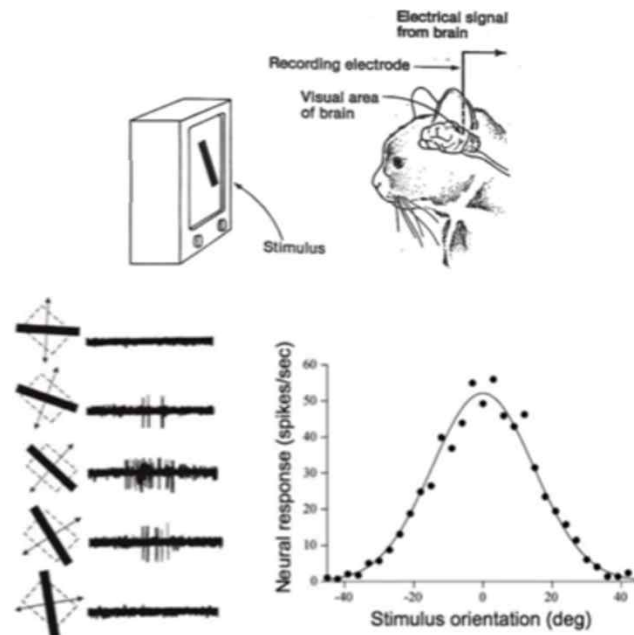
Hubel & Wiesel, 1959

RECEPTIVE FIELDS OF SINGLE
NEURONES IN
THE CAT'S STRIATE CORTEX

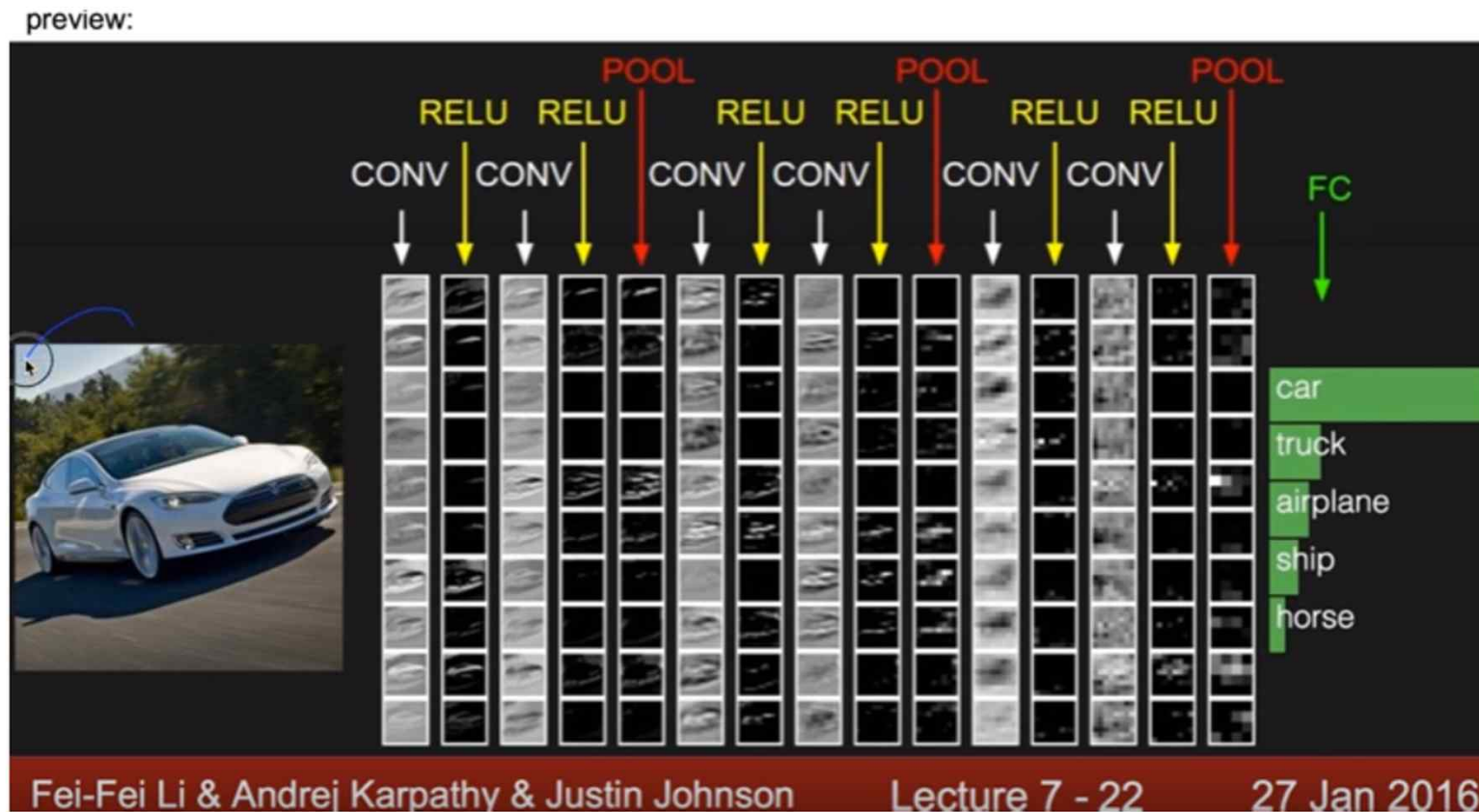
1962

RECEPTIVE FIELDS, BINOCULAR
INTERACTION
AND FUNCTIONAL ARCHITECTURE IN
THE CAT'S VISUAL CORTEX

1968...



➤ CNN (Convolutional Neural Network)



➤ CNN (Convolutional Neural Network)

CNN FILTER

Start with an image (width x height x depth)



32x32x3 image

➤ CNN (Convolutional Neural Network)

CNN FILTER

Let's focus on a small area only (5x5x3)



32x32x3 image

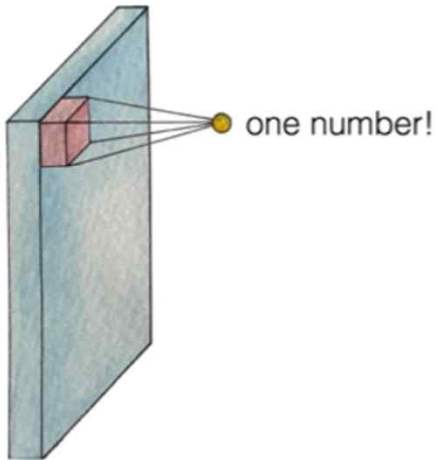


5x5x3 filter

➤ CNN (Convolutional Neural Network)

CNN FILTER

Let's look at other areas with the same filter (w)

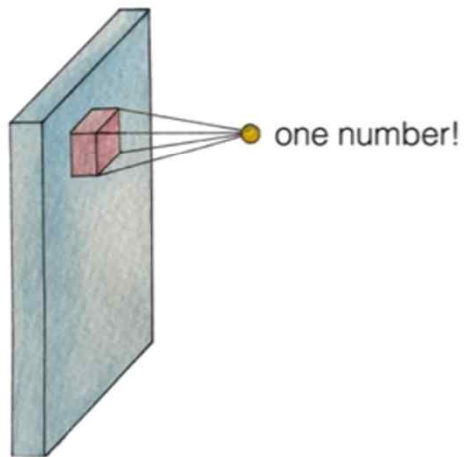


32x32x3 image

➤ CNN (Convolutional Neural Network)

CNN FILTER

Let's look at other areas with the same filter (w)



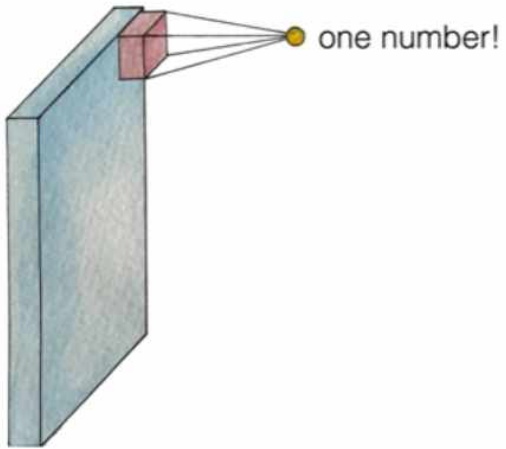
32x32x3 image

How many numbers
can we get?

➤ CNN (Convolutional Neural Network)

CNN FILTER

Let's look at other areas with the same filter (w)

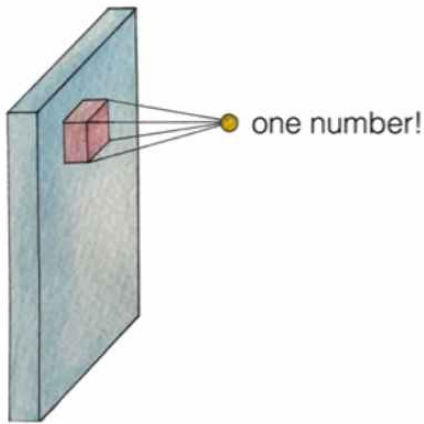


32x32x3 image

➤ CNN (Convolutional Neural Network)

CNN FILTER

Let's look at other areas with the same filter (w)



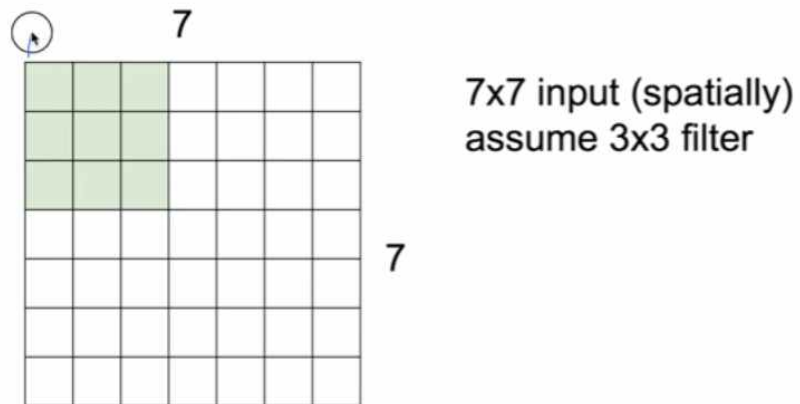
How many numbers
can we get?

32x32x3 image

➤ CNN (Convolutional Neural Network)

MOVING FILTER

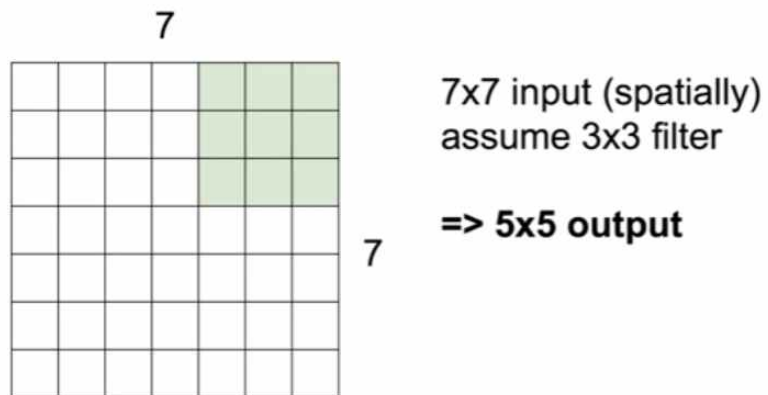
A closer look at spatial dimensions:



➤ CNN (Convolutional Neural Network)

MOVING FILTER

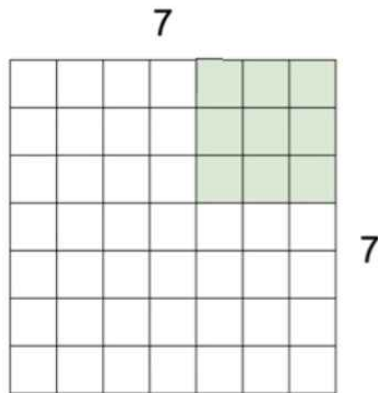
A closer look at spatial dimensions:



➤ CNN (Convolutional Neural Network)

MOVING FILTER

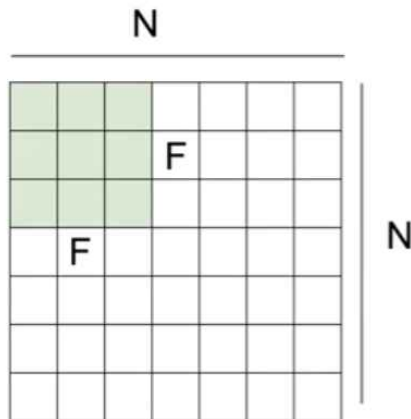
A closer look at spatial dimensions:



7x7 input (spatially)
assume 3x3 filter
applied **with stride 2**
=> 3x3 output!

➤ CNN (Convolutional Neural Network)

MOVING FILTER



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

e.g. $N = 7, F = 3$:

stride 1 $\Rightarrow (7 - 3) / 1 + 1 = 5$

stride 2 $\Rightarrow (7 - 3) / 2 + 1 = 3$

stride 3 $\Rightarrow (7 - 3) / 3 + 1 = 2.33 \therefore 2$

➤ CNN (Convolutional Neural Network)

ZERO PADDINGS

In practice: Common to zero pad the border

0	0	0	0	0	0			
0								
0								
0								
0								

e.g. input 7x7

3x3 filter, applied with **stride 1**

pad with 1 pixel border => what is the output?

7x7 output!

in general, common to see CONV layers with stride 1, filters of size $F \times F$, and zero-padding with $(F-1)/2$. (will preserve size spatially)

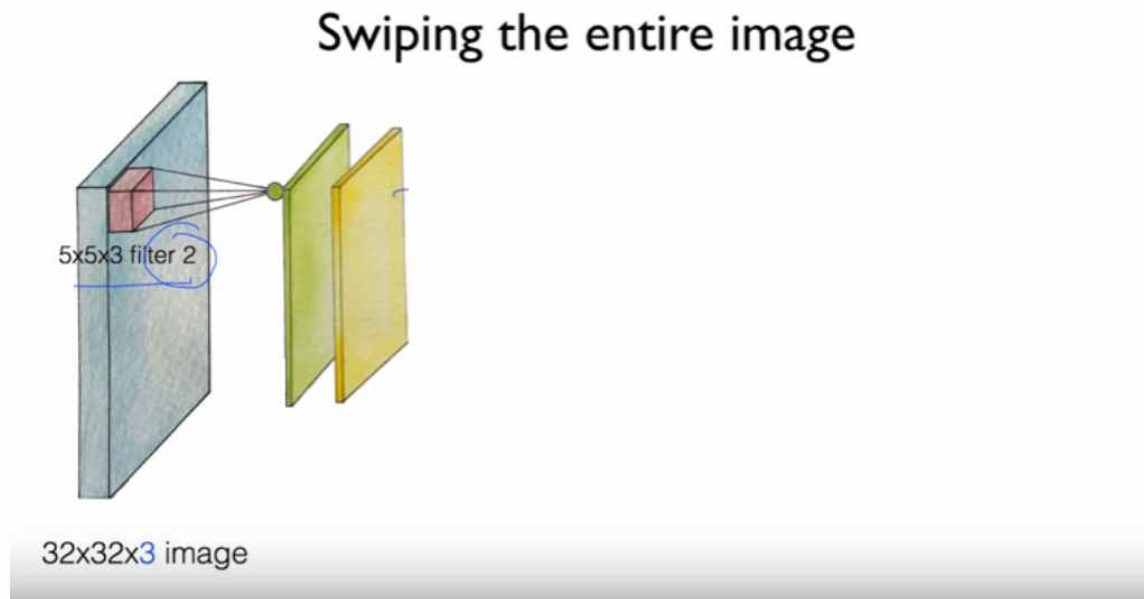
e.g. $F = 3 \Rightarrow$ zero pad with 1

$F = 5 \Rightarrow$ zero pad with 2

$F = 7 \Rightarrow$ zero pad with 3

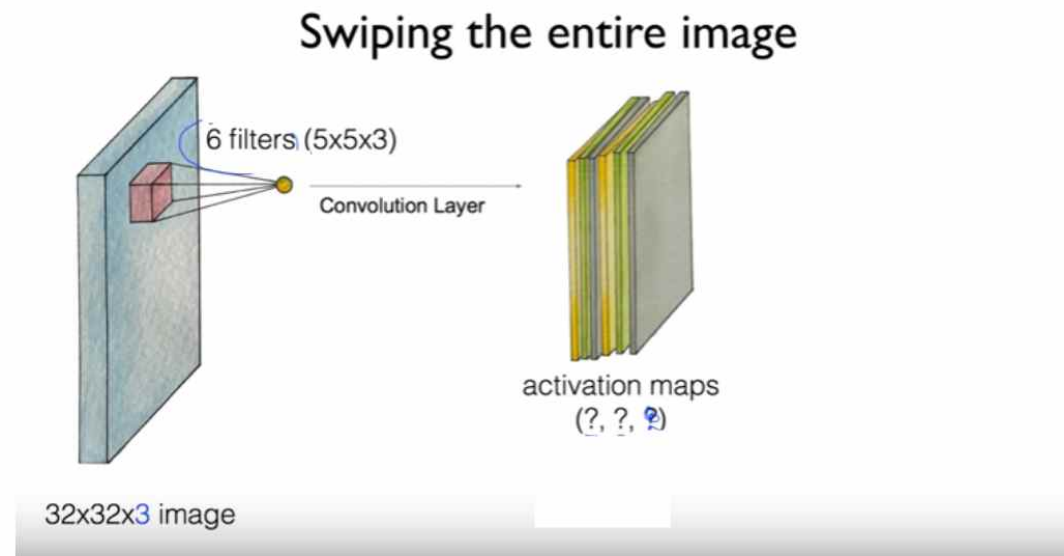
➤ CNN (Convolutional Neural Network)

CONVOLUTIONAL LAYERS



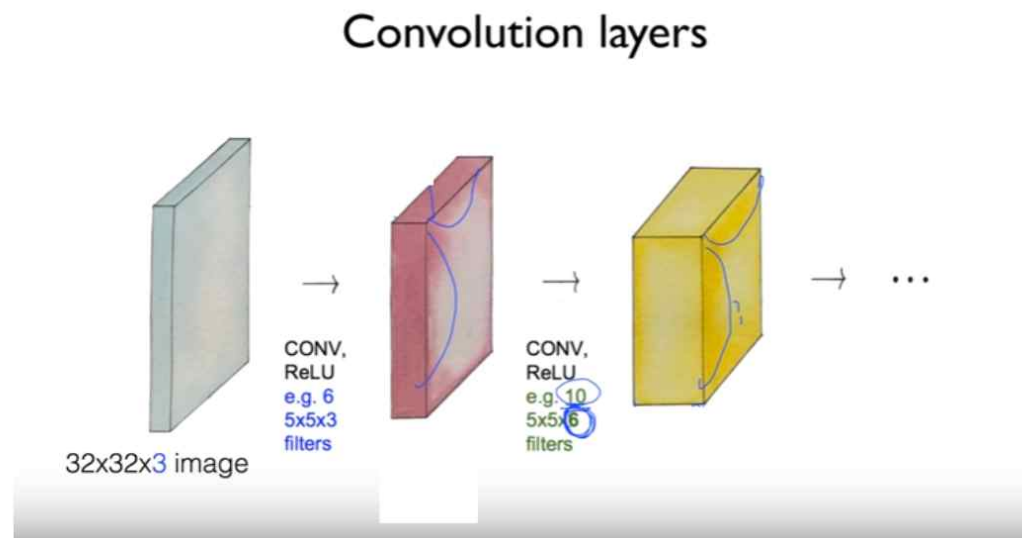
➤ CNN (Convolutional Neural Network)

CONVOLUTIONAL LAYERS



➤ CNN (Convolutional Neural Network)

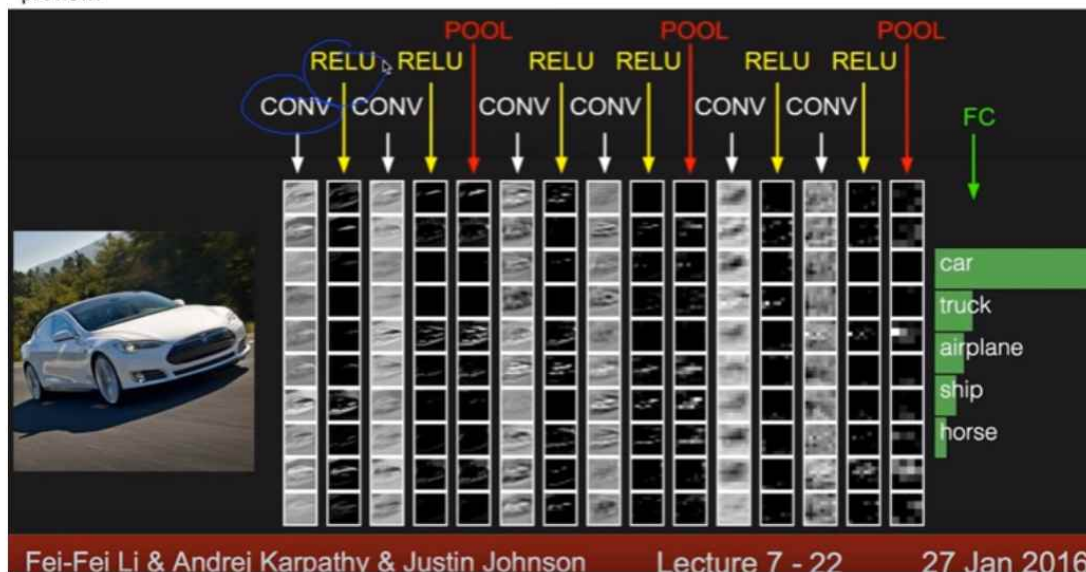
CONVOLUTIONAL LAYERS



➤ CNN (Convolutional Neural Network)

POOLING LAYER

preview:



Fei-Fei Li & Andrei Karpathy & Justin Johnson

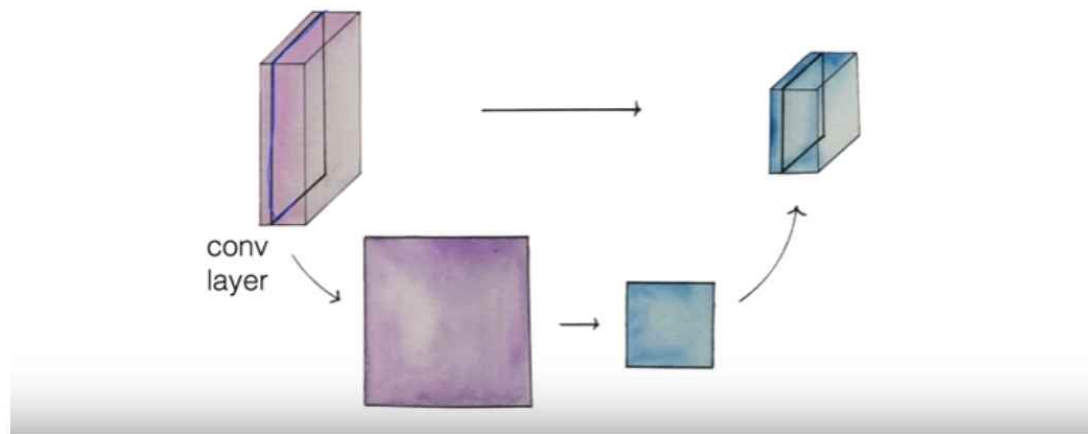
Lecture 7 - 22

27 Jan 2016

➤ CNN (Convolutional Neural Network)

POOLING LAYER

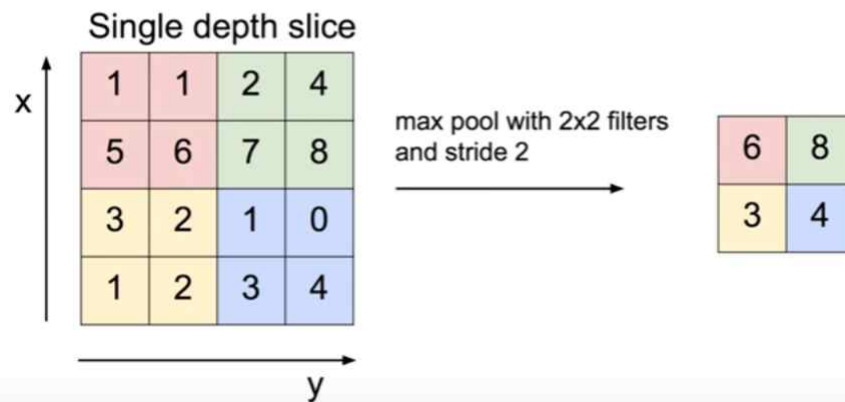
Pooling layer (sampling)



➤ CNN (Convolutional Neural Network)

POOLING LAYER

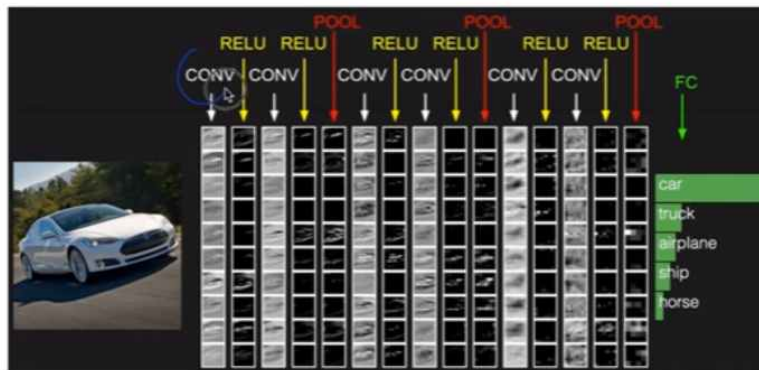
MAX POOLING



- CNN (Convolutional Neural Network)

FULLY CONNECTED LAYER

- Contains neurons that connect to the entire input volume, as in ordinary Neural Networks



➤ CNN (Convolutional Neural Network) History

ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners

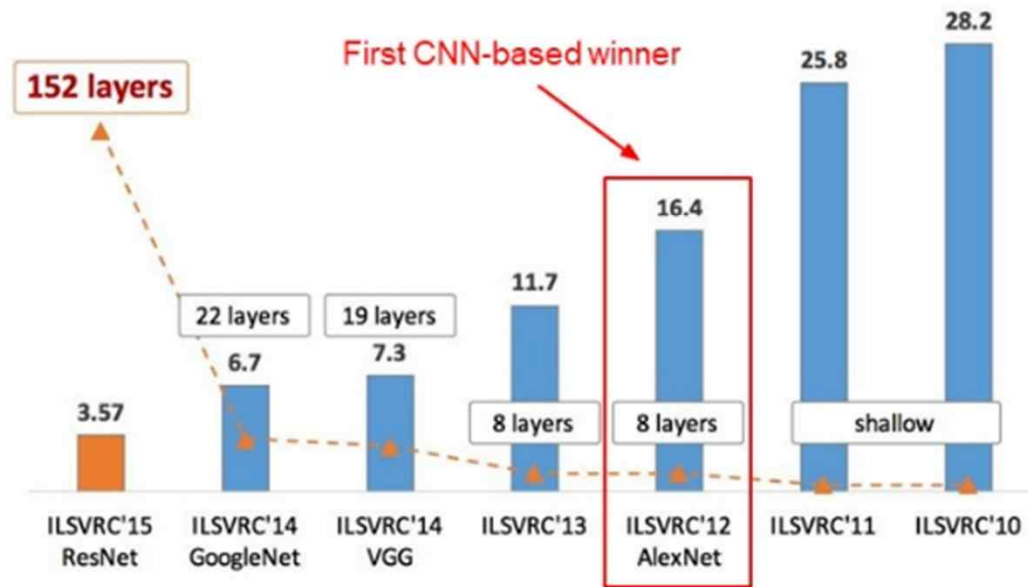


Figure copyright Kaiming He, 2016. Reproduced with permission.

<http://inha-kim.tistory.com/41>

➤ CNN (Convolutional Neural Network) History

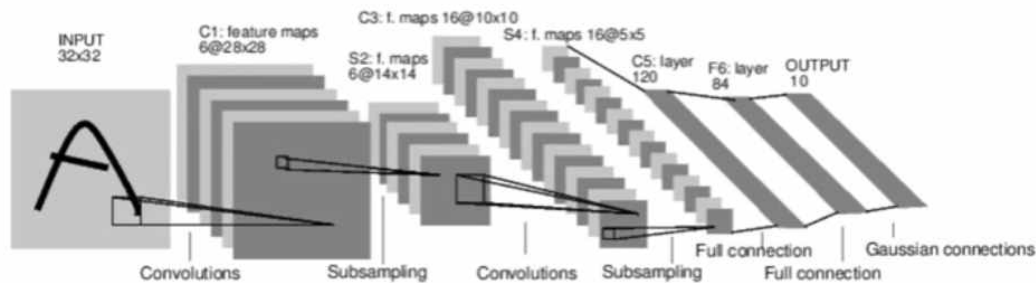
- ILSVRC 2010: SuperVision (에러율 28.2%)
- ILSVRC 2011: SuperVision (에러율 25.8%)
- ILSVRC 2012: AlexNet (에러율 15.3%)
- ILSVRC 2013: ZFNet (에러율 11.2%)
- ILSVRC 2014: GoogLeNet (에러율 6.7%)
- ILSVRC 2015: ResNet (에러율 3.6%)
- ILSVRC 2016: Ensemble of Multiple Models (에러율 2.99%)
<https://tensorflow.blog/2016/09/27/imagenet-ilsvrc-2016-results-out/>
- ILSVRC 2017: SENet (에러율 2.25%)
<https://deep-learning-study.tistory.com/539>
- ILSVRC 2018: SENet (에러율 2.13%)
- ILSVRC 2019: EfficientNet (에러율 1.5%)
<https://greeksharifa.github.io/computer%20vision/2022/03/01/EfficientNet/>

➤ CNN (Convolutional Neural Network)

CASE STUDY : LENET-5

Case Study: LeNet-5

[LeCun et al., 1998]



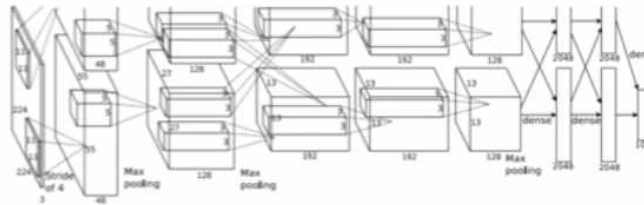
Conv filters were 5x5, applied at stride 1
 Subsampling (Pooling) layers were 2x2 applied at stride 2
 i.e. architecture is [CONV-POOL-CONV-POOL-CONV-FC]

➤ CNN (Convolutional Neural Network)

CASE STUDY : ALEXNET

Case Study: AlexNet

[Krizhevsky et al. 2012]



Input: 227x227x3 images

First layer (CONV1): 96 11x11 filters applied at stride 4

=>

Output volume **[55x55x96]**

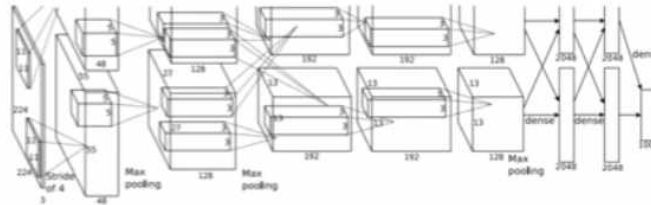
Parameters: $(11*11*3)*96 = 35K$

➤ CNN (Convolutional Neural Network)

CASE STUDY : ALEXNET

Case Study: AlexNet

[Krizhevsky et al. 2012]



Input: 227x227x3 images

After CONV1: 55x55x96

Second layer (POOL1): 3x3 filters applied at stride 2

Output volume: 27x27x96

Parameters: 0!

➤ CNN (Convolutional Neural Network)

CASE STUDY : ALEXNET

Case Study: AlexNet

[Krizhevsky et al. 2012]

Full (simplified) AlexNet architecture:

[227x227x3] INPUT

[55x55x96] **CONV1**: 96 11x11 filters at stride 4, pad 0

[27x27x96] **MAX POOL1**: 3x3 filters at stride 2

[27x27x96] **NORM1**: Normalization layer

[27x27x256] **CONV2**: 256 5x5 filters at stride 1, pad 2

[13x13x256] **MAX POOL2**: 3x3 filters at stride 2

[13x13x256] **NORM2**: Normalization layer

[13x13x384] **CONV3**: 384 3x3 filters at stride 1, pad 1

[13x13x384] **CONV4**: 384 3x3 filters at stride 1, pad 1

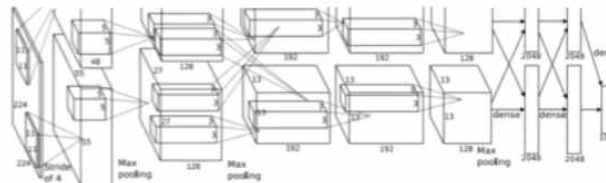
[13x13x256] **CONV5**: 256 3x3 filters at stride 1, pad 1

[6x6x256] **MAX POOL3**: 3x3 filters at stride 2

[4096] **FC6**: 4096 neurons

[4096] **FC7**: 4096 neurons

[1000] **FC8**: 1000 neurons (class scores)



Details/Retrospectives:

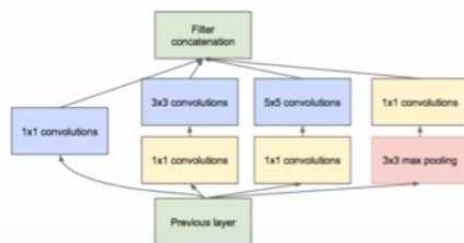
- first use of ReLU
- used Norm layers (not common anymore)
- heavy data augmentation
- dropout 0.5
- batch size 128
- SGD Momentum 0.9
- Learning rate 1e-2, reduced by 10 manually when val accuracy plateaus
- L2 weight decay 5e-4
- 7 CNN ensemble: 18.2% -> 15.4%

➤ CNN (Convolutional Neural Network)

CASE STUDY : GOOGLNET

Case Study: GoogLeNet

[Szegedy et al., 2014]



Inception module

ILSVRC 2014 winner (6.7% top 5 error)

➤ CNN (Convolutional Neural Network)

CASE STUDY : RESNET

Case Study: ResNet [He et al., 2015]
ILSVRC 2015 winner (3.6% top 5 error)

MSRA @ ILSVRC & COCO 2015 Competitions

- **1st places in all five main tracks**
 - ImageNet Classification: *"Ultra-deep"* (quote Yann) **152-layer** nets
 - ImageNet Detection: **16%** better than 2nd
 - ImageNet Localization: **27%** better than 2nd
 - COCO Detection: **11%** better than 2nd
 - COCO Segmentation: **12%** better than 2nd

*improvements are relative numbers

Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". arXiv 2015.

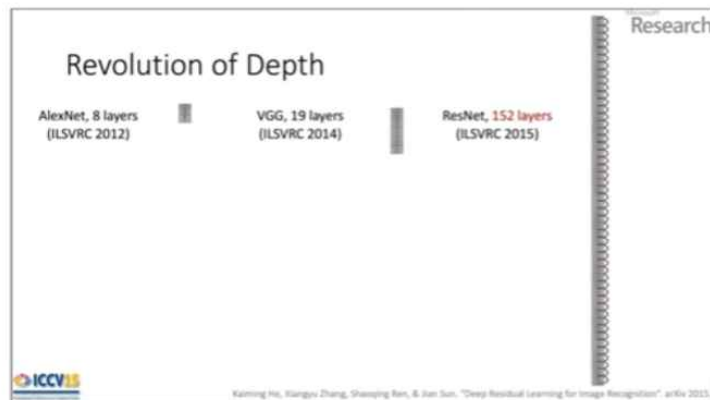
Slide from Kaiming He's recent presentation <https://www.youtube.com/watch?v=1PGLj-uKT1w>

➤ CNN (Convolutional Neural Network)

CASE STUDY : RESNET

Case Study: ResNet [He et al., 2015]

ILSVRC 2015 winner (3.6% top 5 error)



2-3 weeks of training
on 8 GPU machine

at runtime: faster
than a VGGNet!
(even though it has
8x more layers)

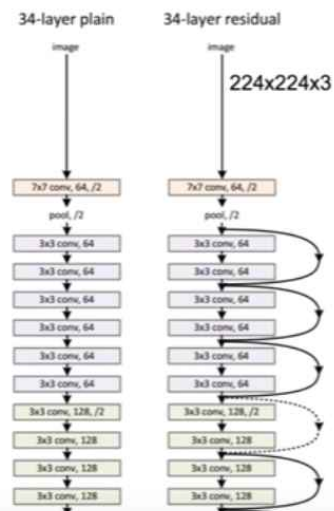
(slide from Kaiming He's recent presentation)

➤ CNN (Convolutional Neural Network)

CASE STUDY : RESNET

Case Study: ResNet

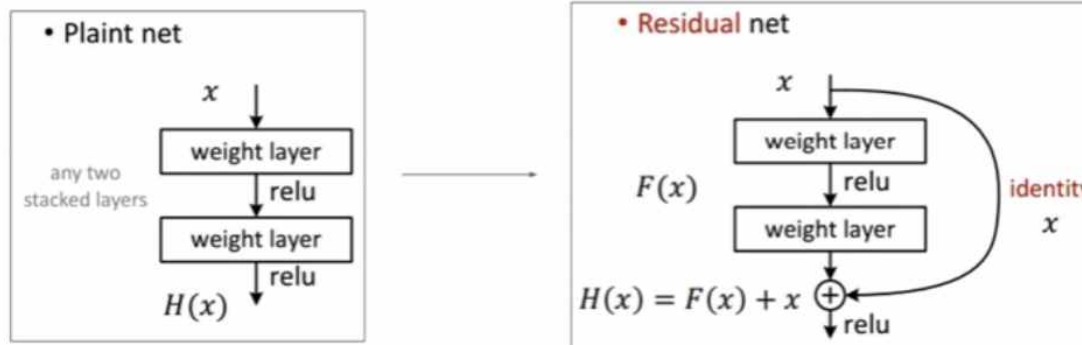
[He et al., 2015]



➤ CNN (Convolutional Neural Network)

CASE STUDY : RESNET

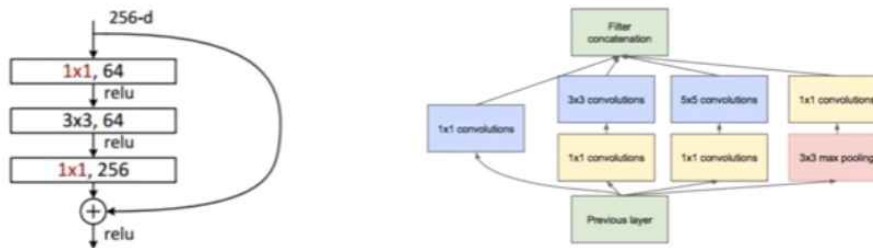
Case Study: ResNet [He et al., 2015]



► CNN (Convolutional Neural Network)

CASE STUDY : RESNET

Case Study: ResNet [He et al., 2015]



➤ CNN (Convolutional Neural Network)

CASE STUDY : SENTENCE CLASSIFICATION

Convolutional Neural Networks for Sentence Classification [Yoon Kim, 2014]

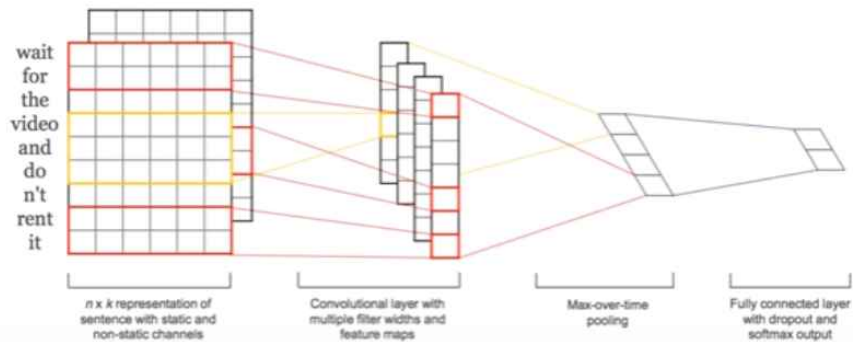
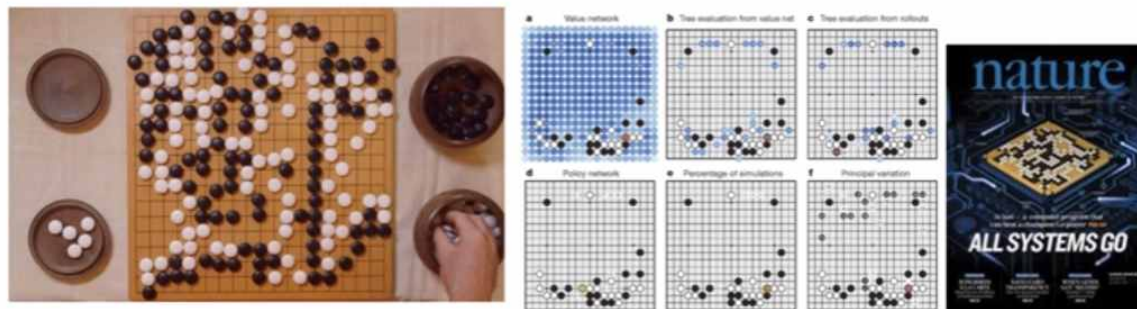


Figure 1: Model architecture with two channels for an example sentence.

➤ CNN (Convolutional Neural Network)

CASE STUDY : DEEPMIND'S ALPHAGO

Case Study Bonus: DeepMind's AlphaGo



➤ CNN (Convolutional Neural Network)

CASE STUDY : DEEPMIND'S ALPHAGO

The input to the policy network is a $19 \times 19 \times 48$ image stack consisting of 48 feature planes. The first hidden layer zero pads the input into a 23×23 image, then convolves k filters of kernel size 5×5 with stride 1 with the input image and applies a rectifier nonlinearity. Each of the subsequent hidden layers 2 to 12 zero pads the respective previous hidden layer into a 21×21 image, then convolves k filters of kernel size 3×3 with stride 1, again followed by a rectifier nonlinearity. The final layer convolves 1 filter of kernel size 1×1 with stride 1, with a different bias for each position, and applies a softmax function. The match version of AlphaGo used $k = 192$ filters; [Fig. 2b](#) and [Extended Data Table 3](#) additionally show the results of training with $k = 128, 256$ and 384 filters.

policy network:

[19x19x48] Input

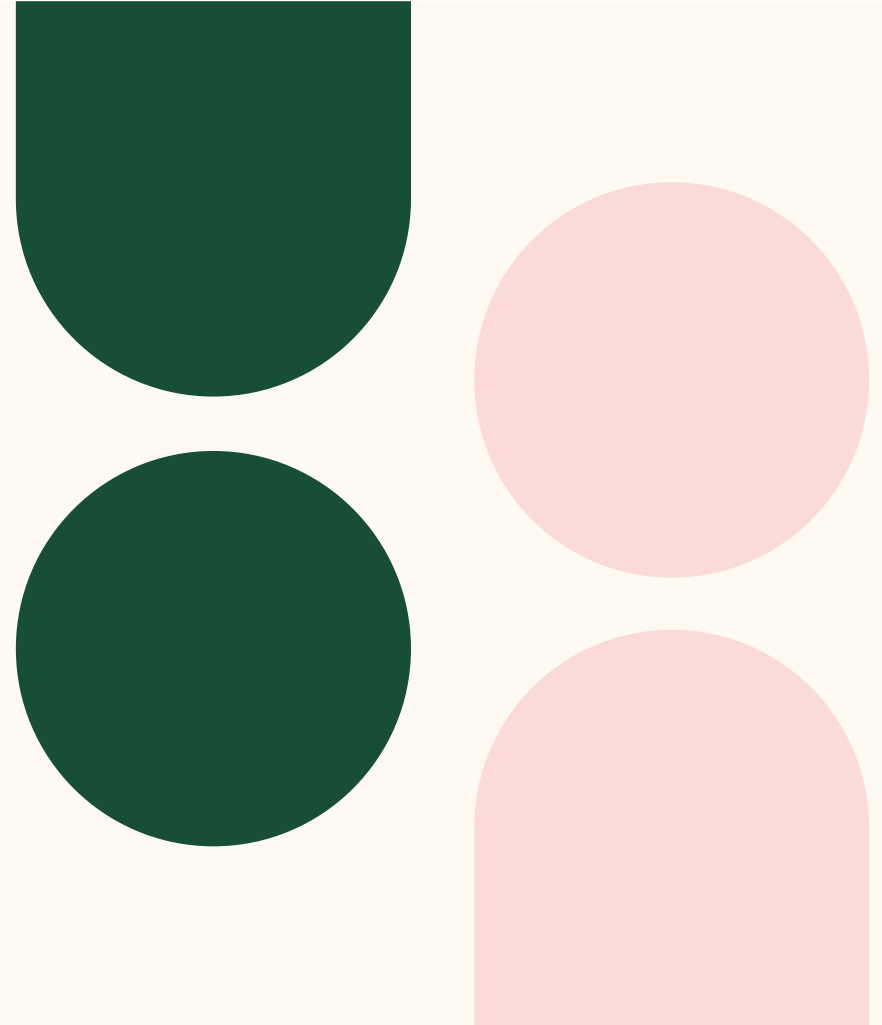
CONV1: 192 5x5 filters, stride 1, pad 2 => [19x19x192]

CONV2..12: 192 3x3 filters, stride 1, pad 1 => [19x19x192]

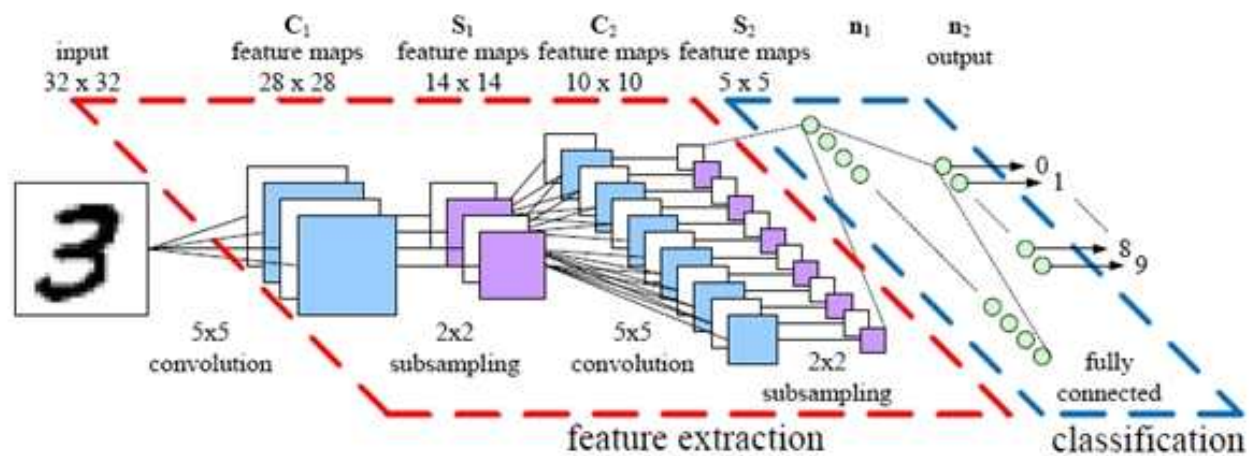
CONV: 1 1x1 filter, stride 1, pad 0 => [19x19] (*probability map of promising moves*)

02 CNN 구현 설명

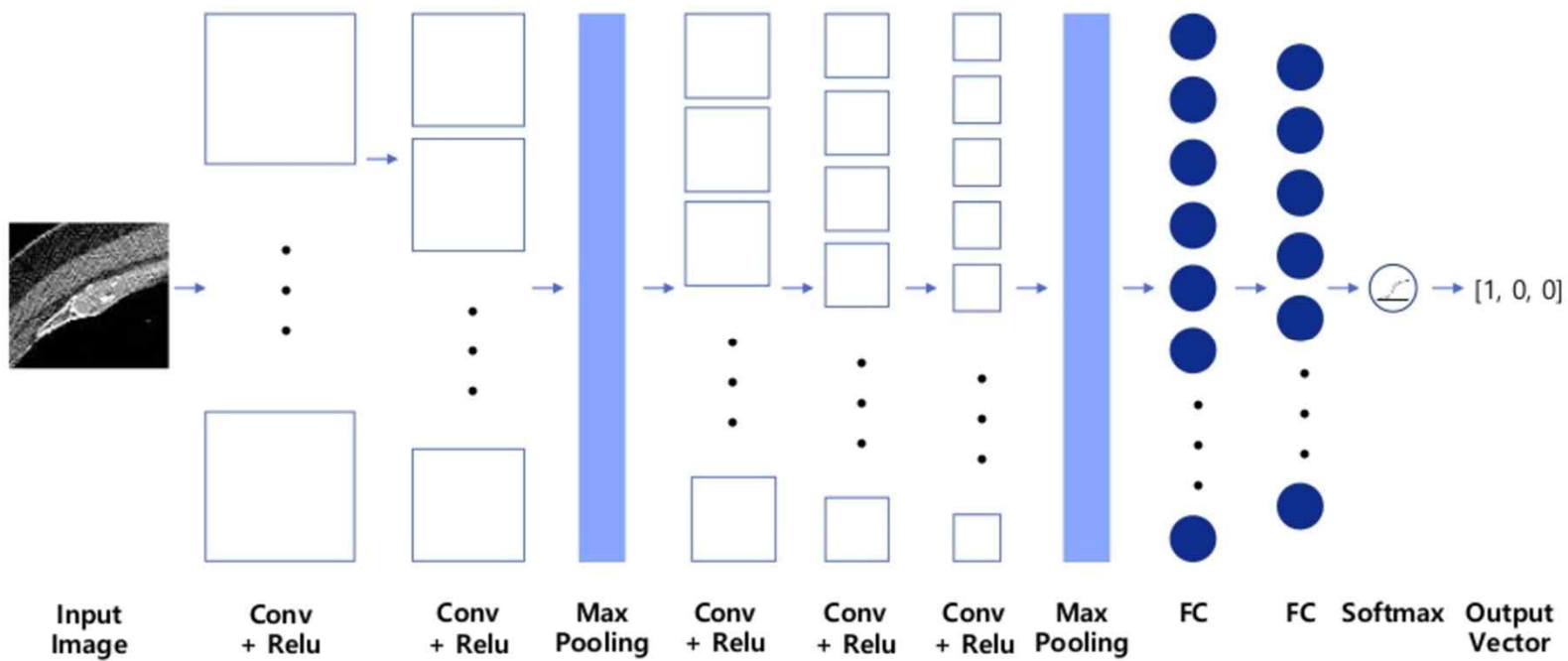
합성곱 연산



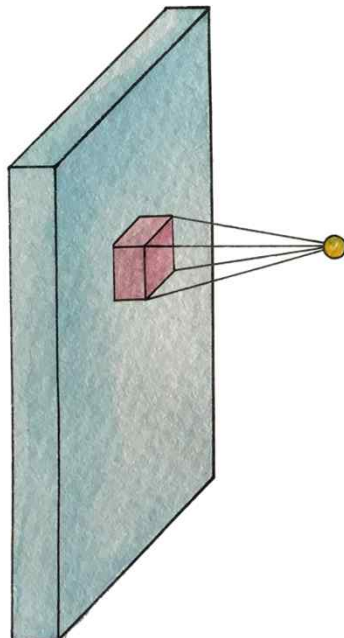
➤ CNN



➤ CNN for CT images



➤ Convolution layer and max pooling



Single depth slice

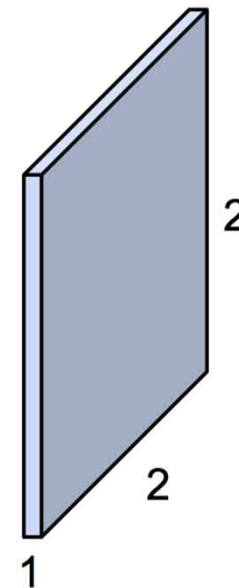
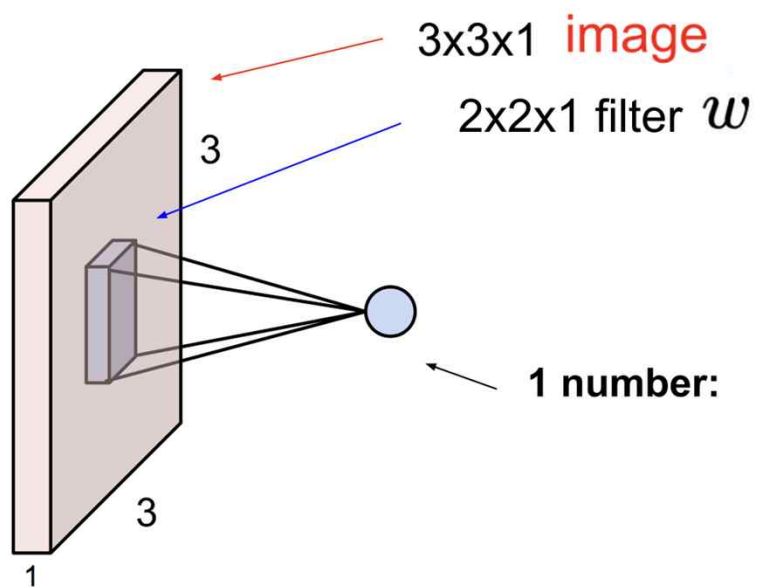
1	1	2	4
5	6	7	8
3	2	1	0
1	2	3	4

max pool with 2x2 filters
and stride 2

6	8
3	4

➤ Simple convolution layer

Stride: 1x1

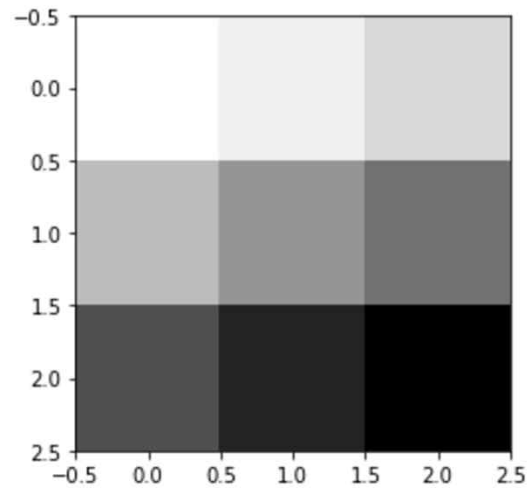


➤ Simple convolution layer

```
In [2]: sess = tf.InteractiveSession()
image = np.array([[[[1],[2],[3]],
                    [[4],[5],[6]],
                    [[7],[8],[9]]]], dtype=np.float32)
print(image.shape)
plt.imshow(image.reshape(3,3), cmap='Greys')

(1, 3, 3, 1)
```

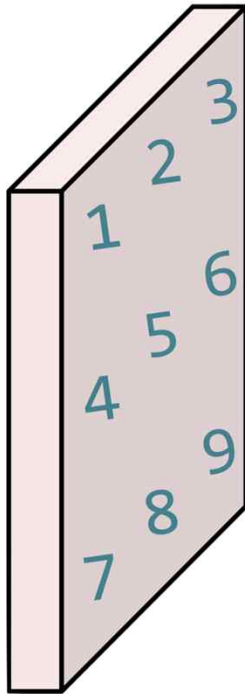
```
Out[2]: <matplotlib.image.AxesImage at 0x10db67dd8>
```



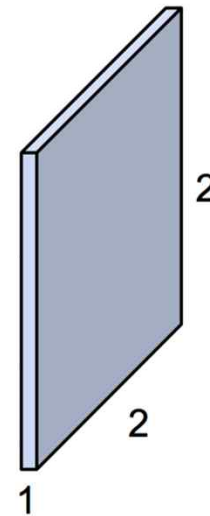
Toy image

➤ Simple convolution layer

Image: 1,3,3,1 image, Filter: 2,2,1,1, Stride: 1x1, Padding: VALID



```
[[[1.],[1.]],  
 [[1.],[1.]]]  
shape=(2,2,1,1)
```

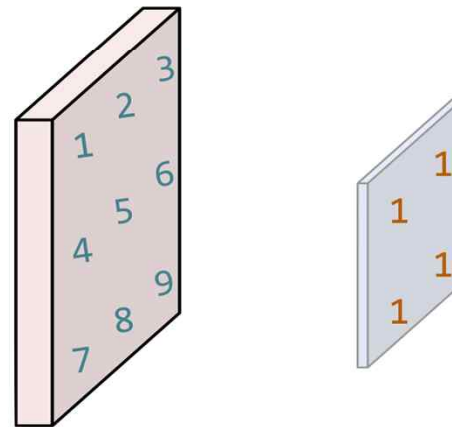
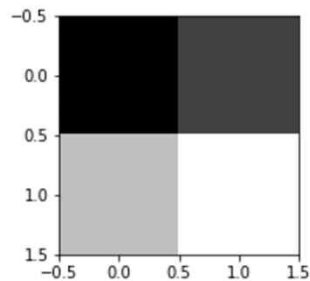


➤ Simple convolution layer

Image: 1,3,3,1 image, Filter: 2,2,1,1, Stride: 1x1, Padding: VALID

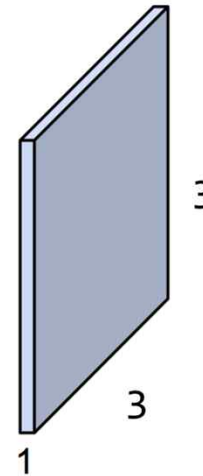
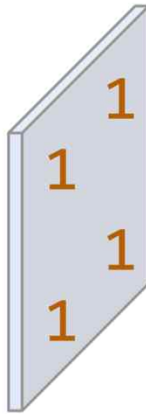
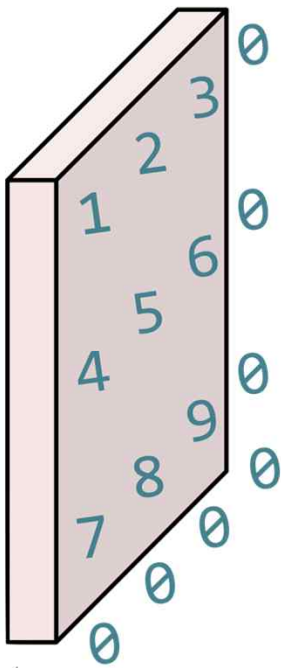
```
# print("imag:\n", image)
print("image.shape", image.shape)
weight = tf.constant([[[[1.]], [[1.]]],
                      [[[1.]], [[1.]]]])
print("weight.shape", weight.shape)
conv2d = tf.nn.conv2d(image, weight, strides=[1, 1, 1, 1], padding='VALID')
conv2d_img = conv2d.eval()
print("conv2d_img.shape", conv2d_img.shape)
conv2d_img = np.swapaxes(conv2d_img, 0, 3)
for i, one_img in enumerate(conv2d_img):
    print(one_img.reshape(2,2))
    plt.subplot(1,2,i+1), plt.imshow(one_img.reshape(2,2), cmap='gray')
```

```
image.shape (1, 3, 3, 1)
weight.shape (2, 2, 1, 1)
conv2d_img.shape (1, 2, 2, 1)
[[ 12.  16.]
 [ 24.  28.]]
```



➤ Simple convolution layer

Image: 1,3,3,1 image, Filter: 2,2,1,1, Stride: 1x1, Padding: **SAME**



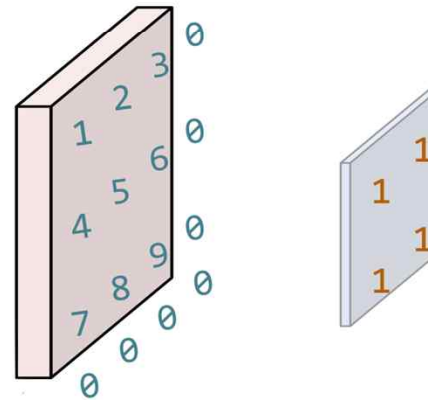
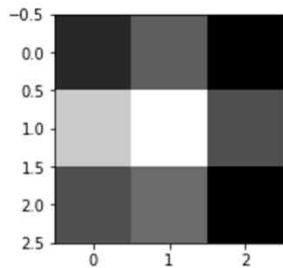
➤ Simple convolution layer

Image: 1,3,3,1 image, Filter: 2,2,1,1, Stride: 1x1, Padding: SAME

```
# print("imag:\n", image)
print("image.shape", image.shape)

weight = tf.constant([[[[1.]], [1.]],
                      [[1.]], [1.]]],
                      dtype=tf.float32)
print("weight.shape", weight.shape)
conv2d = tf.nn.conv2d(image, weight, strides=[1, 1, 1, 1], padding='SAME')
conv2d_img = conv2d.eval()
print("conv2d_img.shape", conv2d_img.shape)
conv2d_img = np.swapaxes(conv2d_img, 0, 3)
for i, one_img in enumerate(conv2d_img):
    print(one_img.reshape(3,3))
    plt.subplot(1,2,i+1), plt.imshow(one_img.reshape(3,3), cmap='gray')
```

```
image.shape (1, 3, 3, 1)
weight.shape (2, 2, 1, 1)
conv2d_img.shape (1, 3, 3, 1)
[[ 12.  16.   9.]
 [ 24.  28.  15.]
 [ 15.  17.   9.]]
```



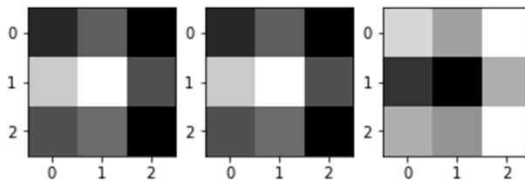
➤ Simple convolution layer

3 filters (2,2,1,3)

```
# print("imag:\n", image)
print("image.shape", image.shape)

weight = tf.constant([[[[1.,10.,-1.]],[[1.,10.,-1.]],[[1.,10.,-1.]],[[1.,10.,-1.]]]],
                      dtype=tf.float32)
print("weight.shape", weight.shape)
conv2d = tf.nn.conv2d(image, weight, strides=[1, 1, 1, 1], padding='SAME')
conv2d_img = conv2d.eval()
print("conv2d_img.shape", conv2d_img.shape)
conv2d_img = np.swapaxes(conv2d_img, 0, 3)
for i, one_img in enumerate(conv2d_img):
    print(one_img.reshape(3,3))
    plt.subplot(1,3,i+1), plt.imshow(one_img.reshape(3,3), cmap='gray')
```

```
image.shape (1, 3, 3, 1)
weight.shape (2, 2, 1, 3)
conv2d_img.shape (1, 3, 3, 3)
[[ 12.  16.   9.]
 [ 24.  28.  15.]
 [ 15.  17.   9.]]
[[ 120.  160.   90.]
 [ 240.  280.  150.]
 [ 150.  170.   90.]]
[[-12. -16.  -9.]
 [-24. -28. -15.]
 [-15. -17.  -9.]]
```



➤ Max Pooling

4	3
2	1

```
In [19]: image = np.array([[[[4],[3]],
                             [[2],[1]]]], dtype=np.float32)
pool = tf.nn.max_pool(image, ksize=[1, 2, 2, 1],
                      strides=[1, 1, 1, 1], padding='SAME')
print(pool.shape)
print(pool.eval())
```

(1, 2, 2, 1)
 [[[4.]
 [3.]

 [[2.]
 [1.]]]]

SAME: Zero paddings

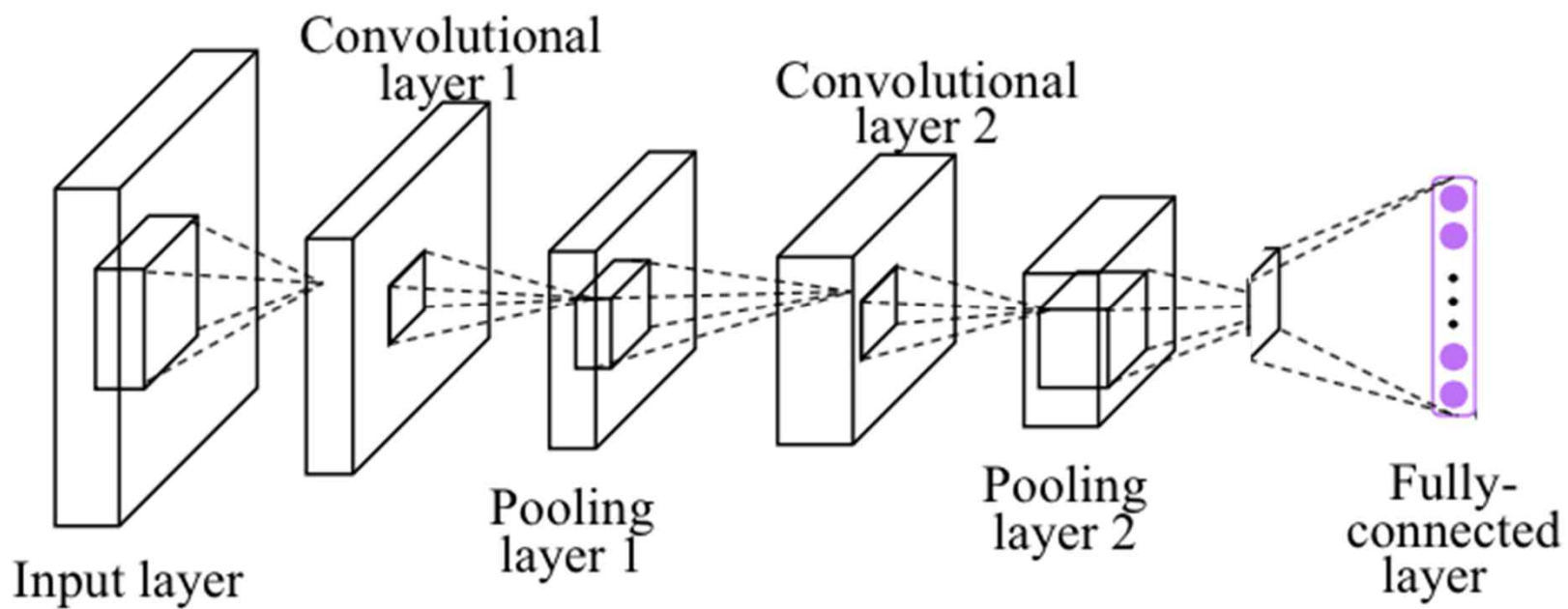
4	3	0
2	1	0
0	0	0

4	3	0
2	1	0
0	0	0

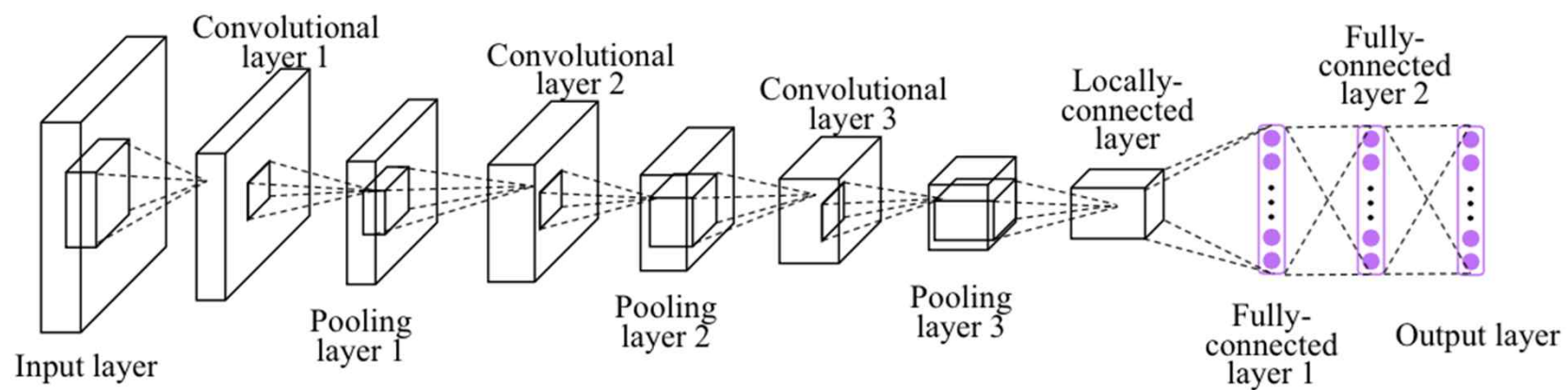
4	3	0
2	1	0
0	0	0

4	3	0
2	1	0
0	0	0

➤ Simple CNN



➤ Deep CNN



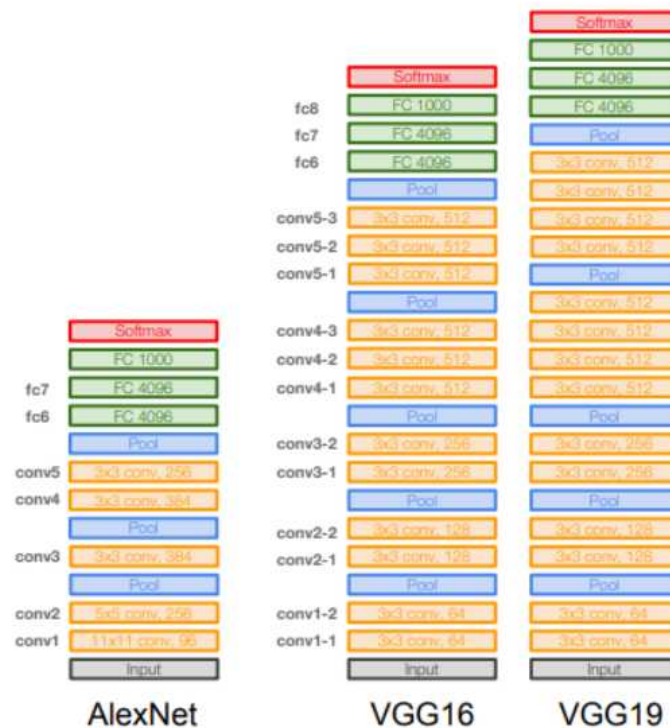
➤ VGGNet

Case Study: VGGNet

[Simonyan and Zisserman, 2014]

Details:

- ILSVRC'14 2nd in classification, 1st in localization
- Similar training procedure as Krizhevsky 2012
- No Local Response Normalisation (LRN)
- Use VGG16 or VGG19 (VGG19 only slightly better, more memory)
- Use ensembles for best results
- FC7 features generalize well to other tasks



VGG : 2014 2nd winner, Visual Geometry Group in Oxford
`tf.keras.applications.VGG19()`