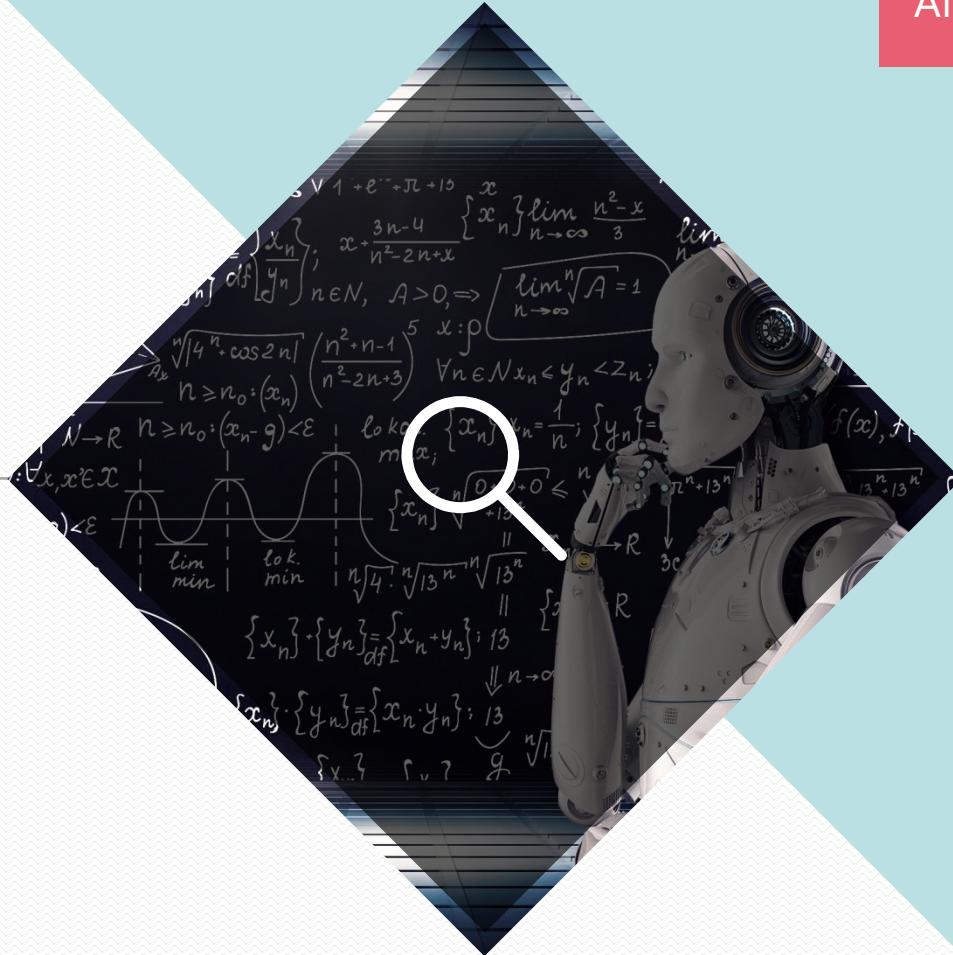


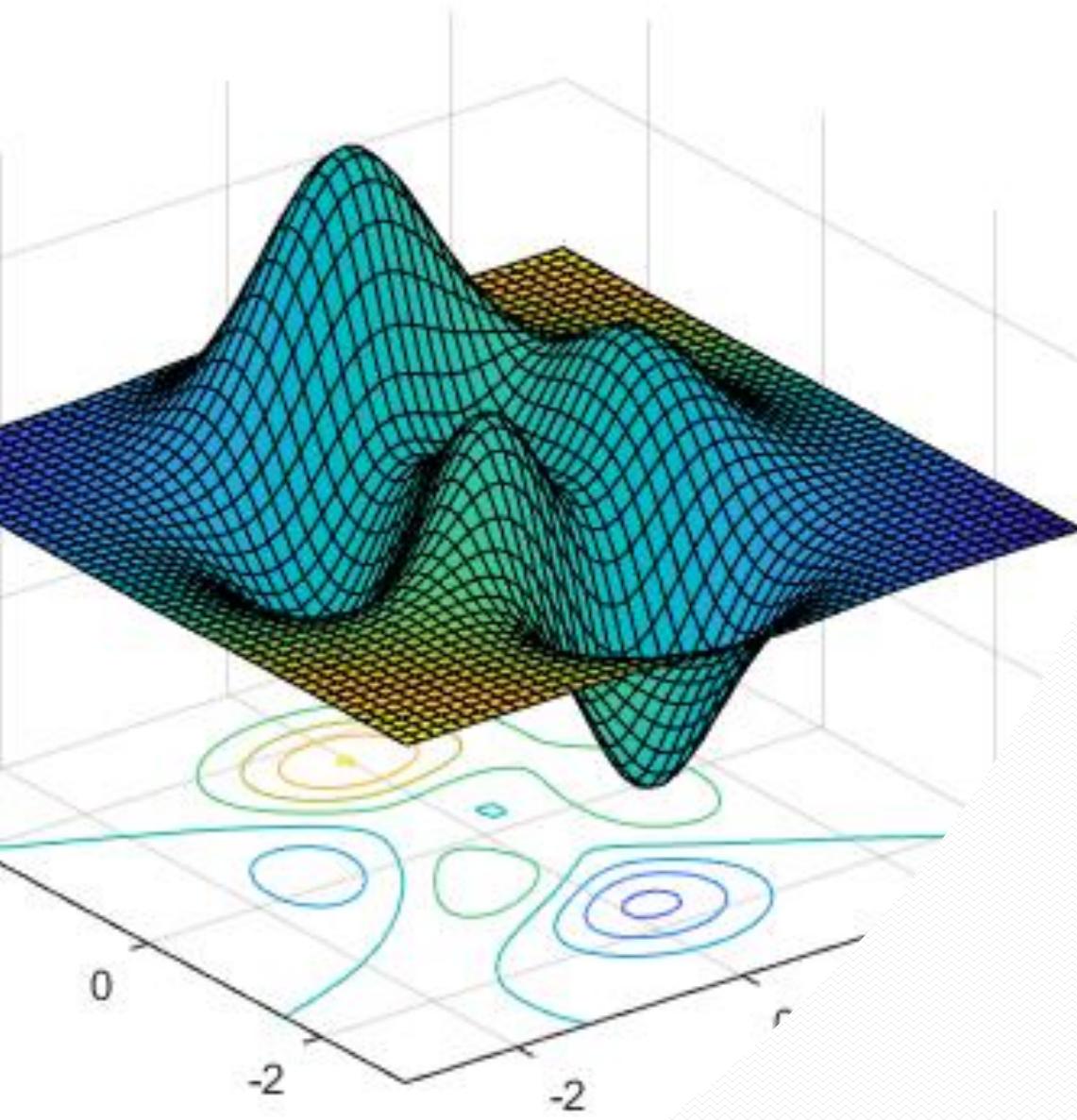
CNN

CNN & Convolution



Agenda

- *CNN*
 - *CNN*
 - *Use cases*
 - *Classification*
 - *Classification Activation Map*
 - *Segmentation*
 - *Object detection*



CNN

PROBLEMS ON ANN

REQUIRES MANY PARAMETERS

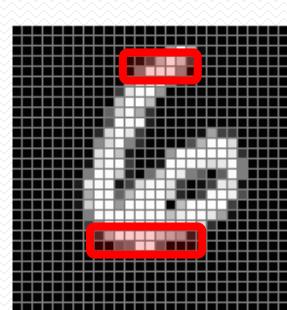
28x28 image → $(28 \times 28 + 1) \times 10 = 7,850 \approx 7.8\text{K}$

640x480 image → $(640 \times 480 + 1) \times 10 = 3,072,010 \approx 3\text{M}$

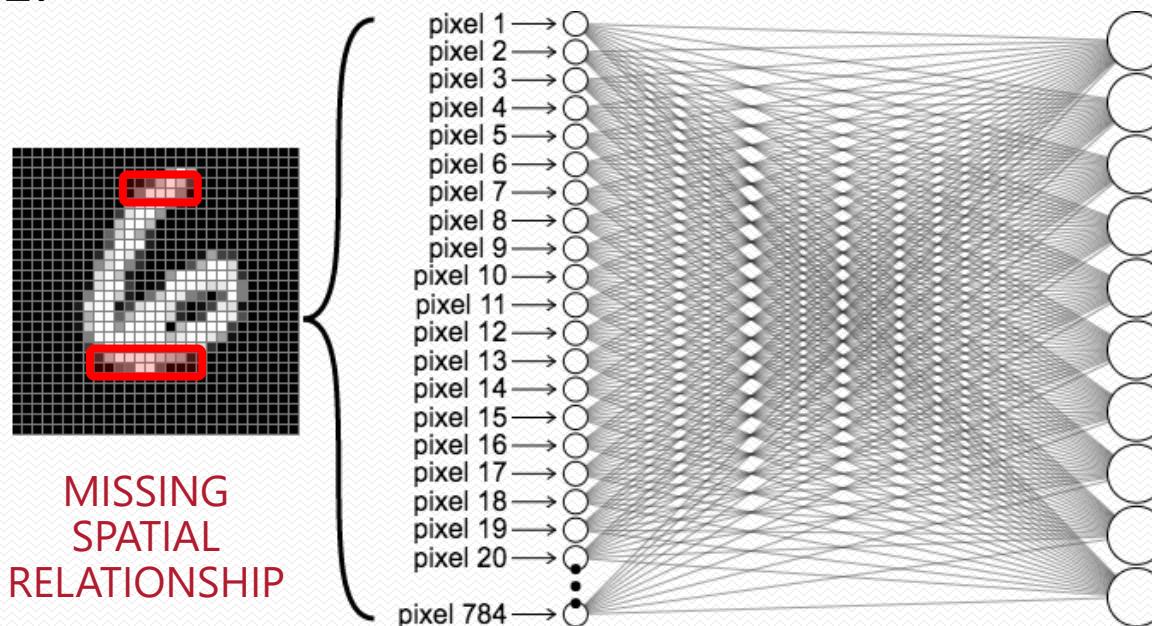
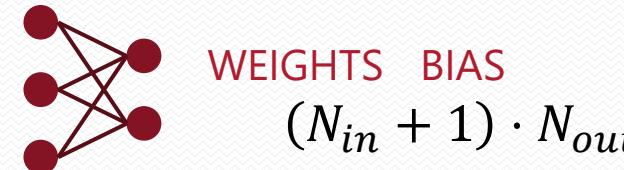
REQUIRES LOTS OF MEMORY AND COMPUTATION COST

TAKES 1D FORM OF INPUT ONLY

Lose spatial information!



MISSING SPATIAL RELATIONSHIP



CNN Classification 기본구조

- ✓ 특정 속성을 추출하는 과정후 **Flatten** 과정을 거쳐 **Fully Connected**됨
- ✓ 일반적인 CNN은 Convolution 연산, Activation 연산(대부분 ReLU), Pooling 연산의 반복으로 구성됨 (Feature Learning)
- ✓ 일정 횟수 이상의 Feature Learning 과정 이후에는 Flatten 과정을 통해 이미지가 1차원의 벡터로 변환됨

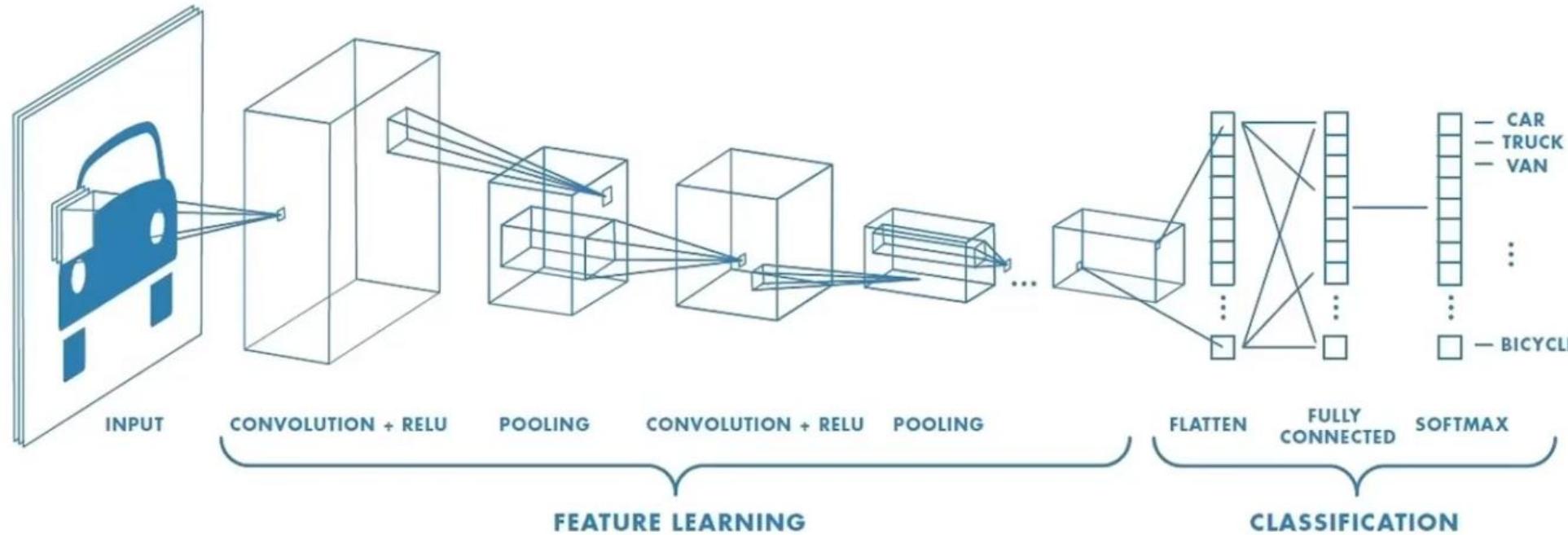
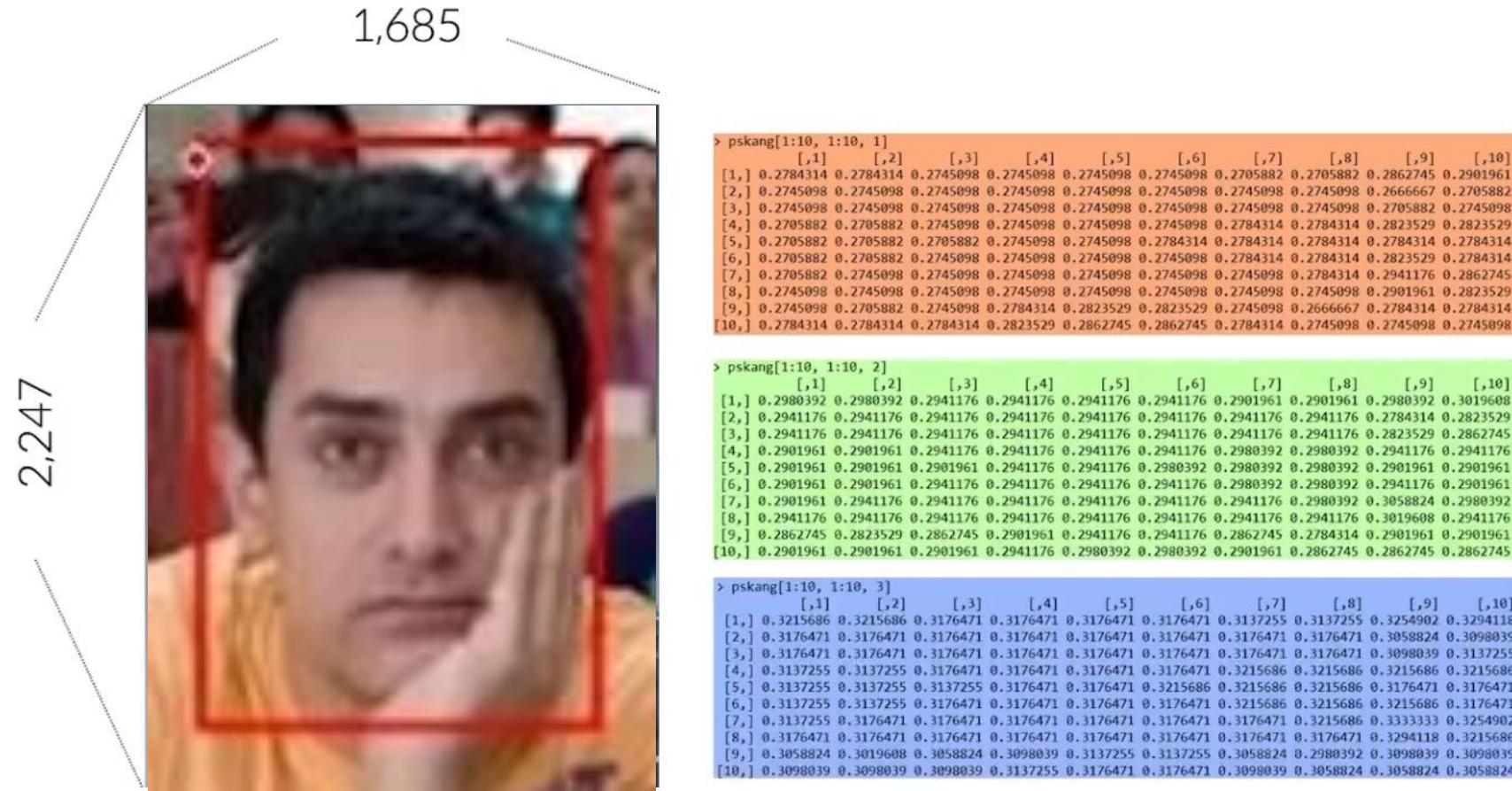


Image presentation

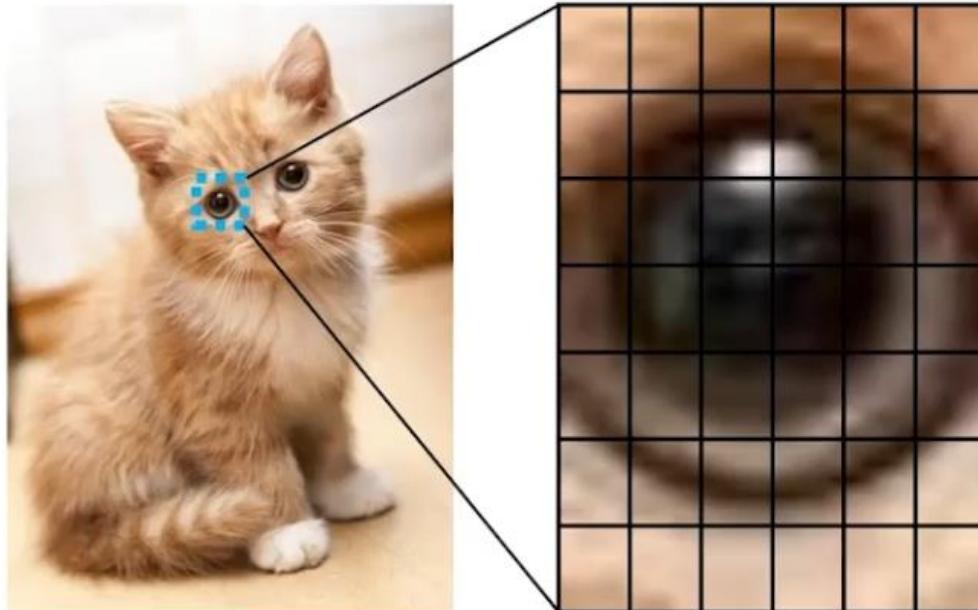
✓ 컬러 이미지는 3차원의 Tensor로 표현됨 : Width X Height X 3 (RGB)



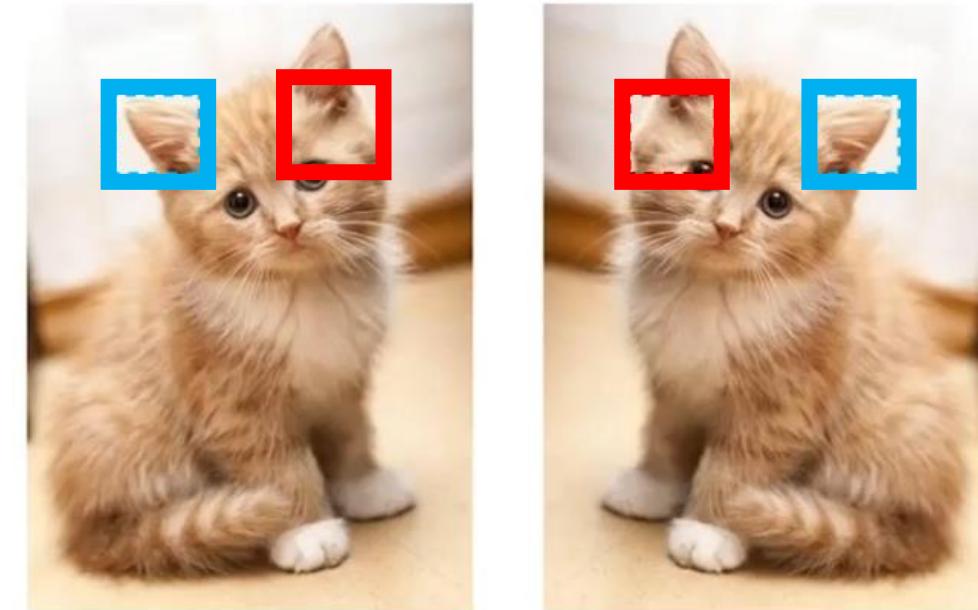
이미지 데이터의 특징

- ✓ 인접 픽셀간 높은 상관관계(spatially-local correlation)
- ✓ 이미지의 부분적 특성(e.g, 눈, 귀)은 고정된 위치에 등장하지 않음(feature invariance)

Spatially-Local Correlation

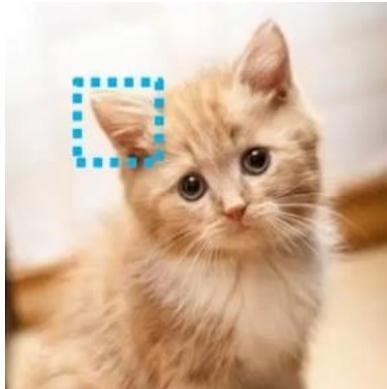


Invariant Feature



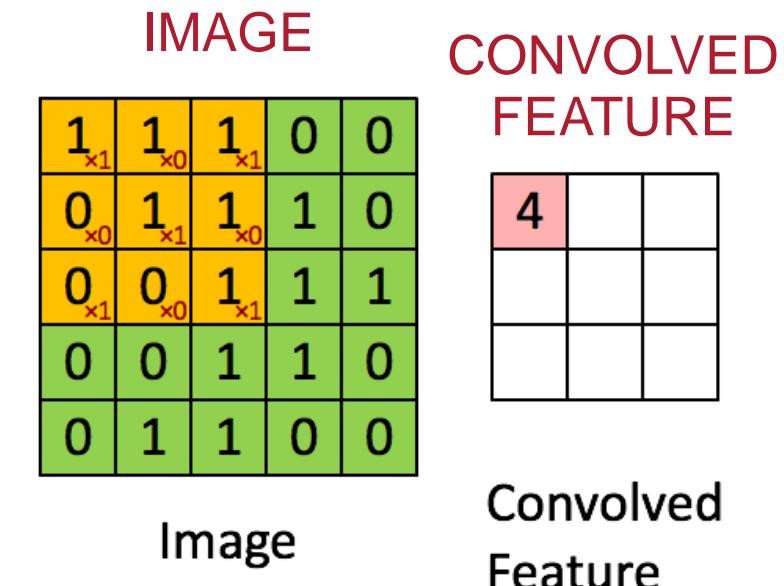
합성곱연산(Convolution)

- ✓ Spatially-local correlation을 고려하기 위한 Sparse connection 구성
- ✓ 인접한 변수만을 이용하여 새로운 Feature 생성
- ✓ Invariant feature를 추출하기 위해 Shared weight 개념 이용.
- ✓ 특정 특질을 추출하기 위한 도구로서 같은 대상 크기에는 위치가 다르더라도 동일한 Weight 적용



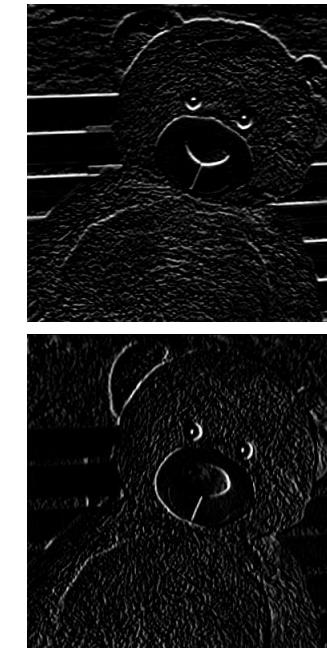
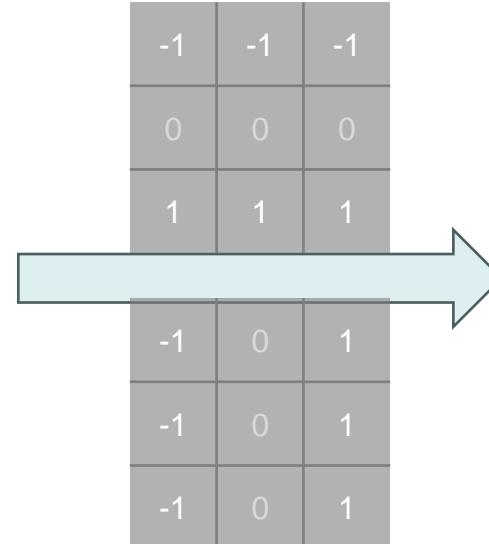
Sparse connection
전체 픽셀이 아닌 일부 픽셀만
가중치로 연결하여 연산수행

Shared weight
동일한 Filter는 대상
영역에 상관없이 같은
가중치를 사용하여 연산 수행



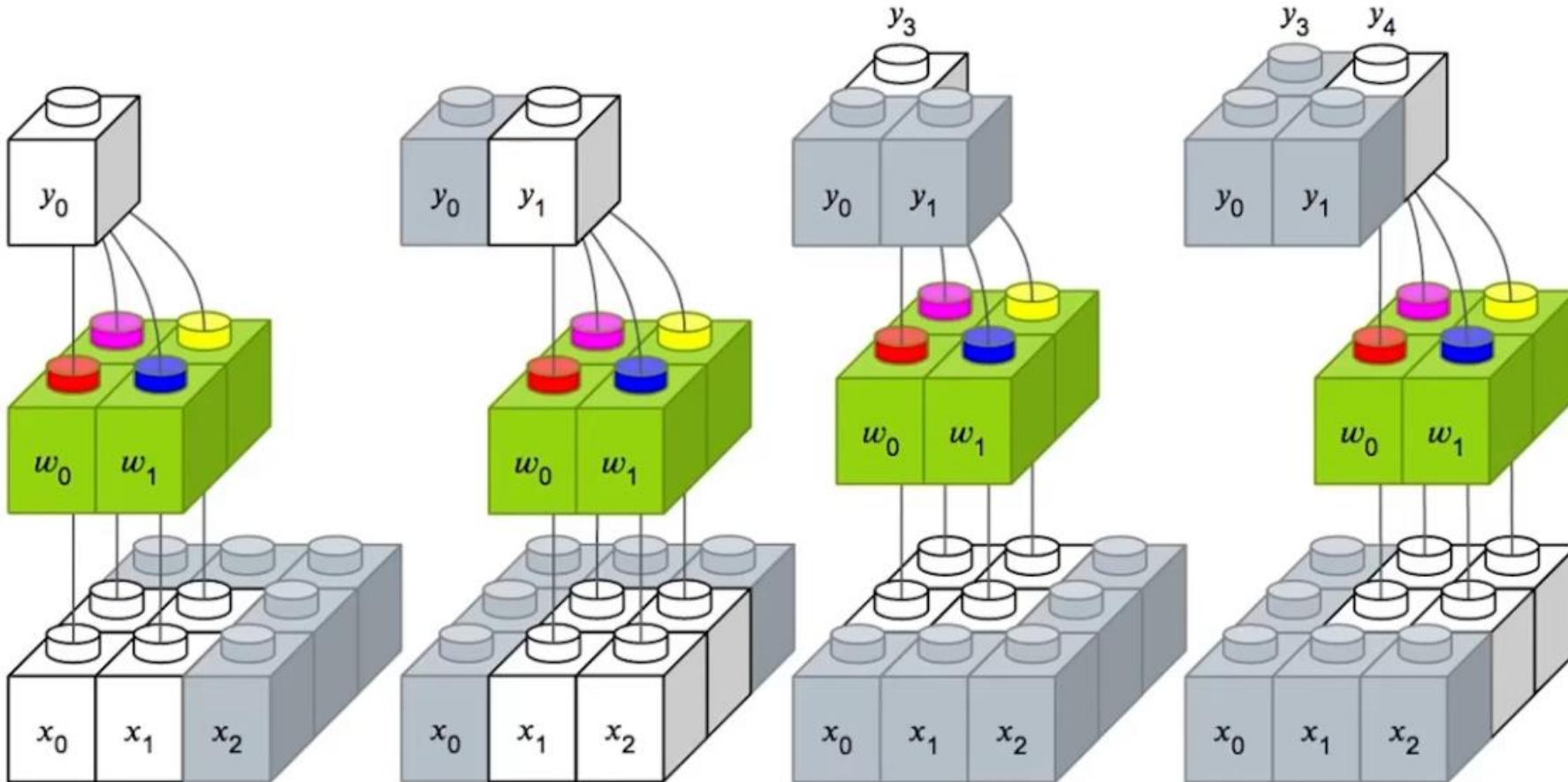
Feature Engineering

- ✓ Image Convolution(Filter, Kernel)
- ✓ 특정 속성을 탐지하는데 사용하는 matrix (예 : edge detection)
- ✓ CNN 등장 전 Computer Vision을 연구하는 연구자들의 핵심 주제는 "특정 특질을 잘 추출 Feature를 어떻게 설계해야 하는가?"였음



<https://www.software.ac.uk/blog/2016-10-03-how-photo-playboy-became-part-scientific-culture>

합성곱연산(Convolution)



https://tykimos.github.io/2017/01/27/CNN_Layer_Talk/

CONVOLUTION

MATHEMATICAL OPERATION

Slides one function over another
e.g., Photo filters



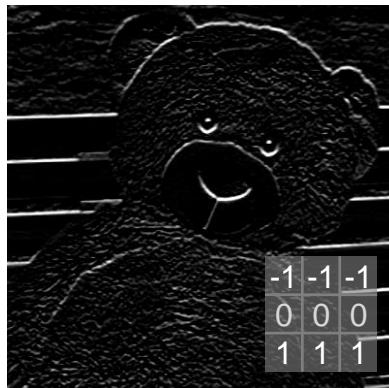
ORIGINAL



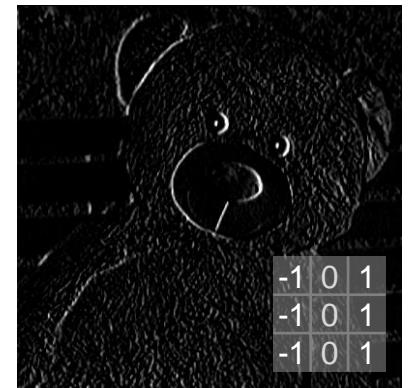
SHARPEN



EMBOSS



EXTRACT
HORIZONTAL LINE



EXTRACT
VERTICAL LINE

IMAGE

1 $\times 1$	1 $\times 0$	1 $\times 1$	0	0
0 $\times 0$	1 $\times 1$	1 $\times 0$	1	0
0 $\times 1$	0 $\times 0$	1 $\times 1$	1	1
0	0	1	1	0
0	1	1	0	0

CONVOLVED FEATURE

4		

CONVOLUTION : MULTI-CHANNEL

0	0	0	0	0	0	0	...
0	156	155	156	158	158	158	...
0	153	154	157	159	159	159	...
0	149	151	155	158	159	159	...
0	146	146	149	153	158	158	...
0	145	143	143	148	158	158	...
...

Input Channel #1 (Red)

0	0	0	0	0	0	0	...
0	167	166	167	169	169	169	...
0	164	165	168	170	170	170	...
0	160	162	166	169	170	170	...
0	156	156	159	163	168	168	...
0	155	153	153	158	168	168	...
0	154	152	152	157	167	167	...
...

Input Channel #2 (Green)

0	0	0	0	0	0	0	...
0	163	162	163	165	165	165	...
0	160	161	164	166	166	166	...
0	156	158	162	165	166	166	...
0	155	155	158	162	167	167	...
0	154	152	152	157	167	167	...
...

Input Channel #3 (Blue)

-1	-1	1
0	1	-1
0	1	1

Kernel Channel #1



308

1	0	0
1	-1	-1
1	0	-1

Kernel Channel #2



-498

0	1	1
0	1	0
1	-1	1

Kernel Channel #3



164

+

+ 1 = -25

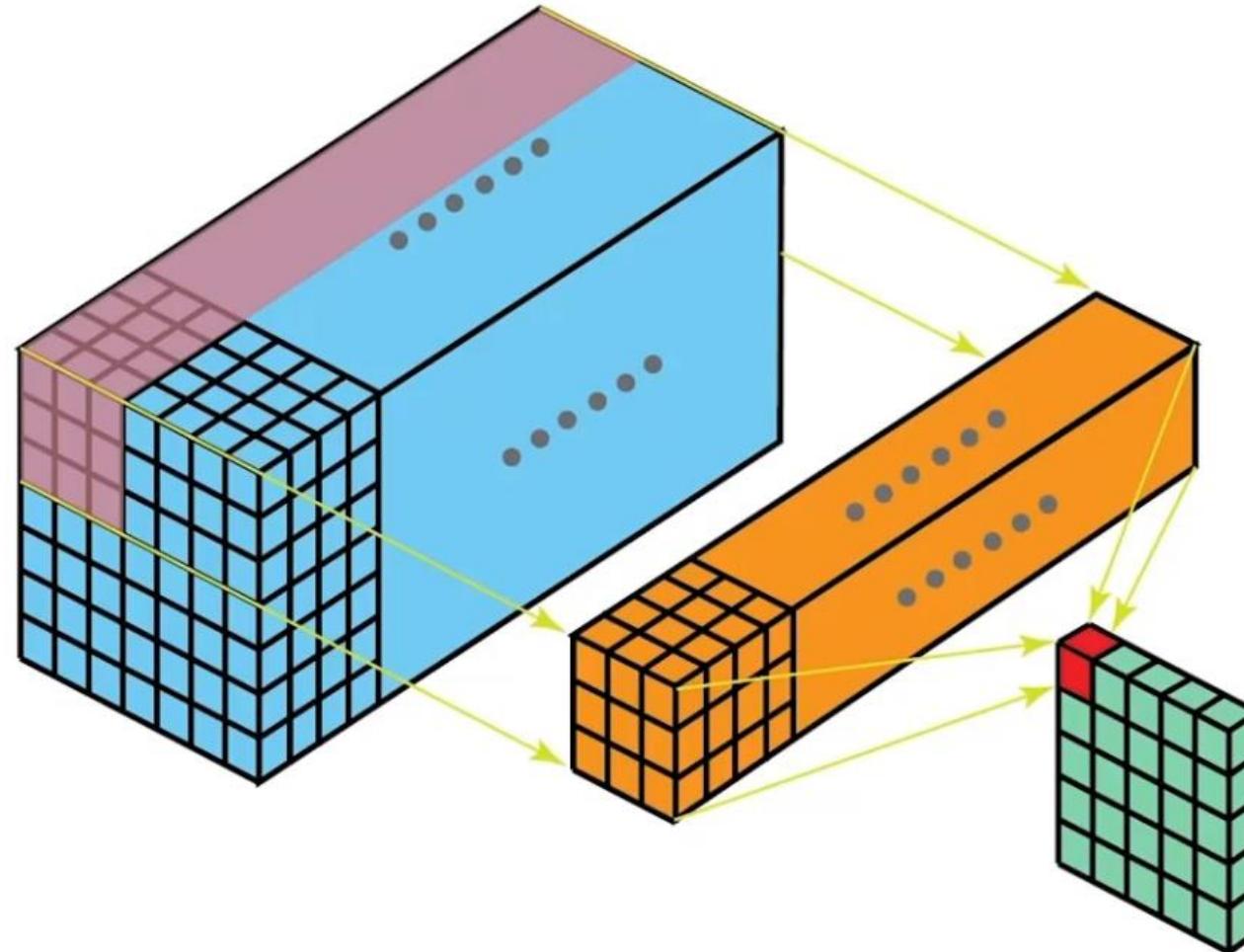


Bias = 1

-25				...
				...
				...
				...
...

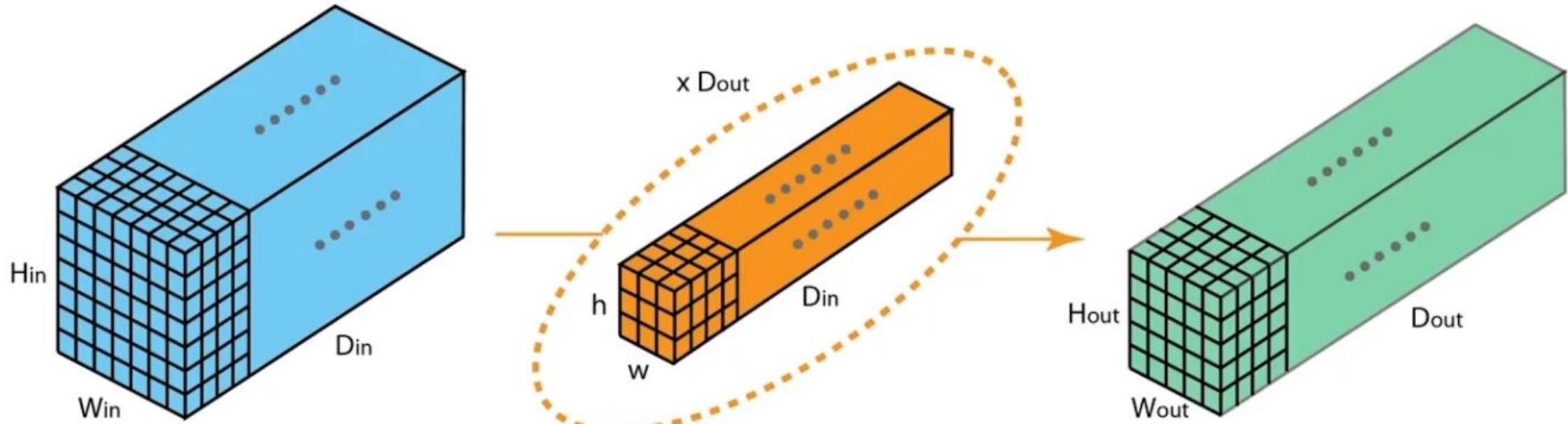
Input CH = Filter CH
Output CH = # of filters

합성곱연산(Convolution)



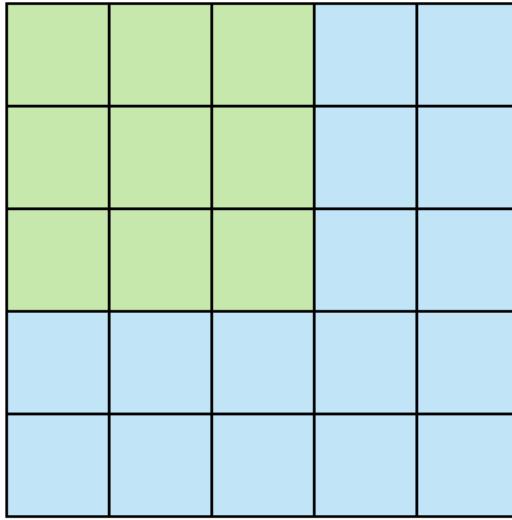
<https://towardsdatascience.com/a-comprehensive-introduction-to-different-types-of-convolutions-in-deep-learning-669281e58215>

합성곱연산(Convolution)

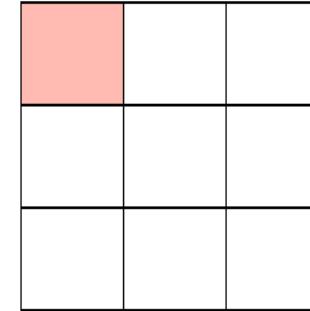


합성곱연산(Convolution: Stride)

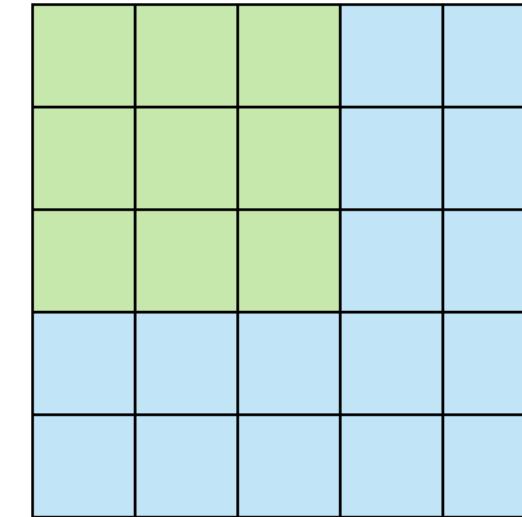
- ✓ 이슈 : 한번에 한 칸씩 이동하면 너무 오래 걸리지 않을까?
- ✓ 해결책 : Filter가 한번에 여러 칸 이동하도록 허용하자(Stride)



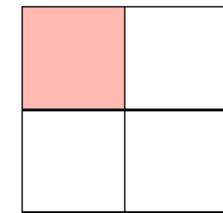
Stride 1



Feature Map



Stride 2



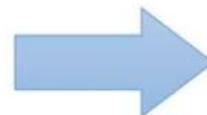
Feature Map

합성곱연산(Convolution: Padding)

- ✓ 문제점 : 가장자리에 있는 픽셀들은 중앙에 위치한 픽셀들에 비해 Convolution 연산이 적게 수행됨
- ✓ 해결책 : Padding (원래 이미지 테두리에 0의 값을 갖는 pad를 추가)

3	1	1	2	8	4
1	0	7	3	2	6
2	3	5	1	1	3
1	4	1	2	6	5
3	2	1	3	7	2
9	2	6	2	5	1

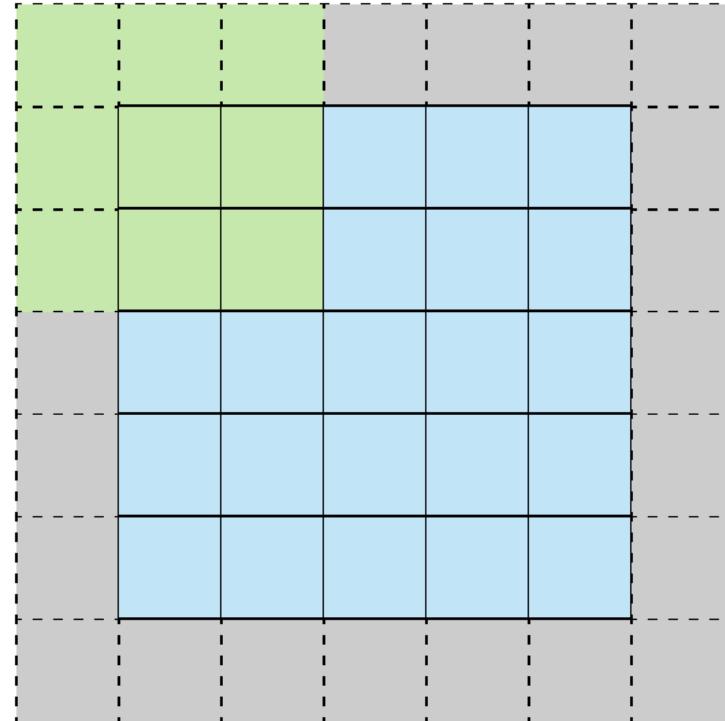
Padding
with size 1



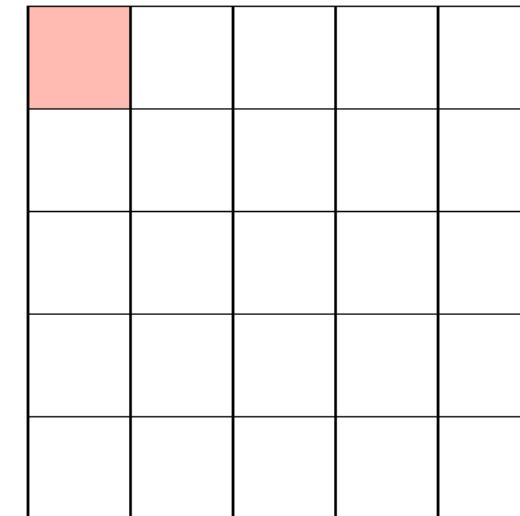
0	0	0	0	0	0	0	0
0	3	1	1	2	8	4	0
0	1	0	7	3	2	6	0
0	2	3	5	1	1	3	0
0	1	4	1	2	6	5	0
0	3	2	1	3	7	2	0
0	9	2	6	2	5	1	0
0	0	0	0	0	0	0	0

합성곱연산(Convolution)

Stride & padding

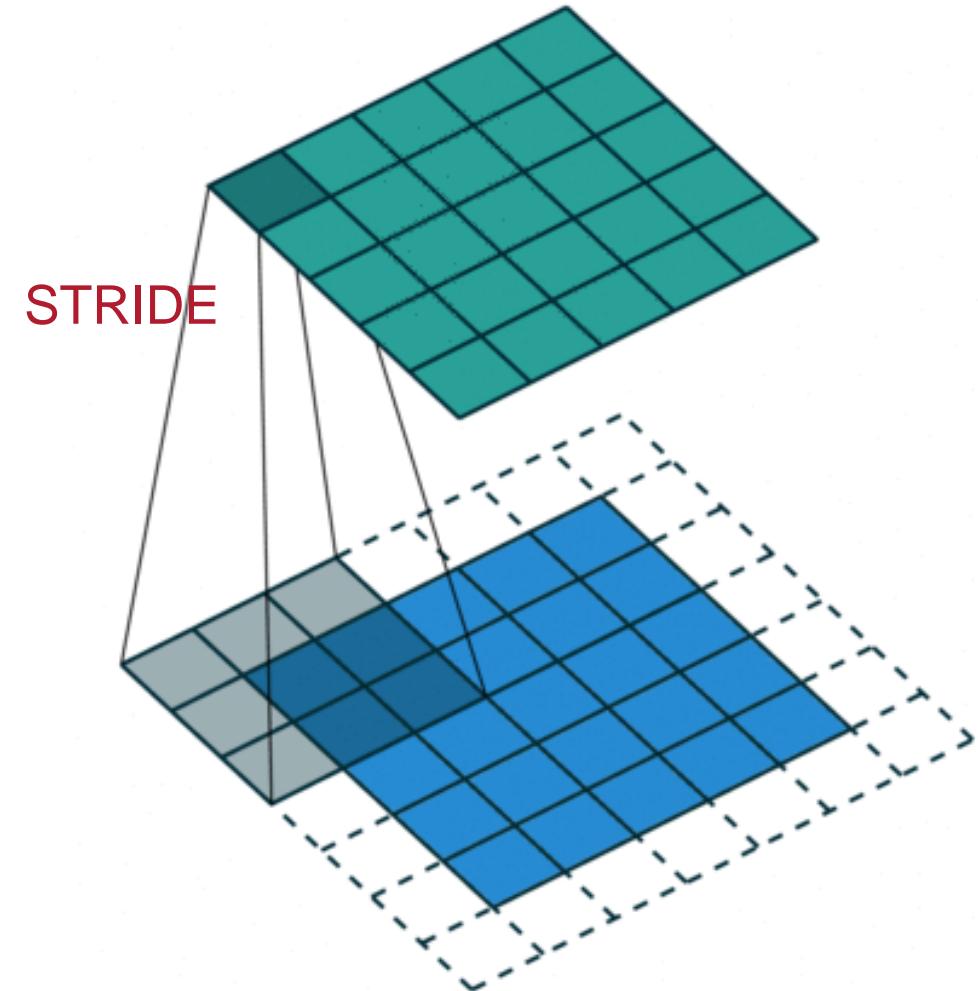
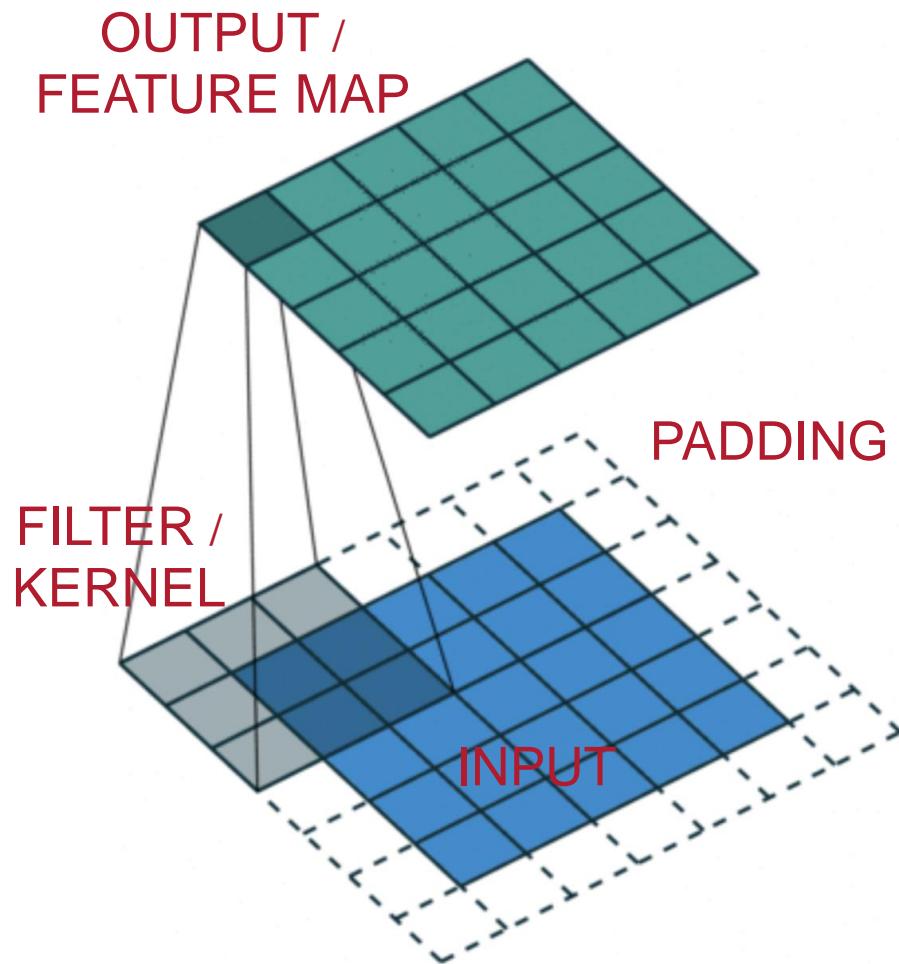


Stride 1 with Padding



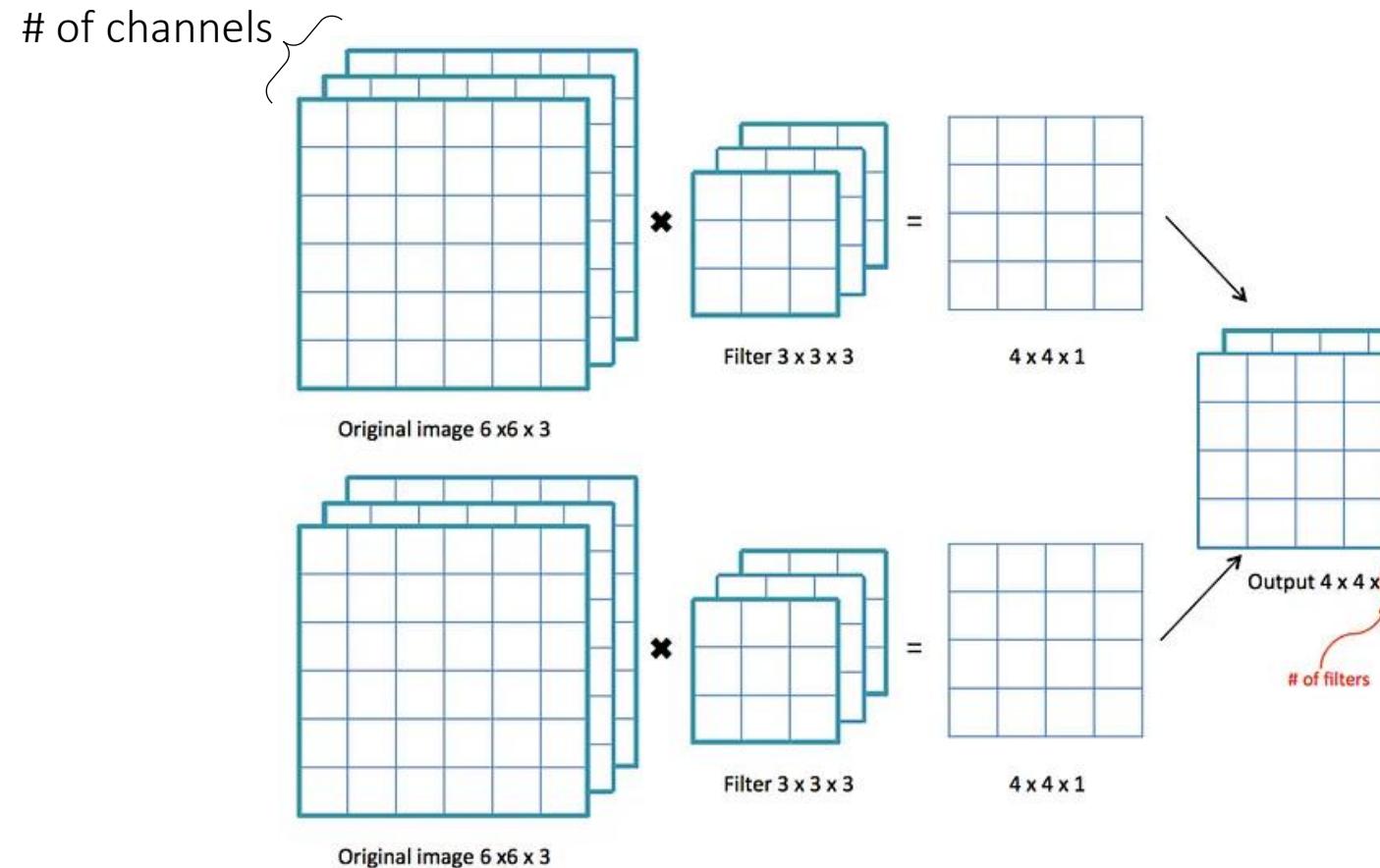
Feature Map

CONVOLUTIONAL LAYER / CONV



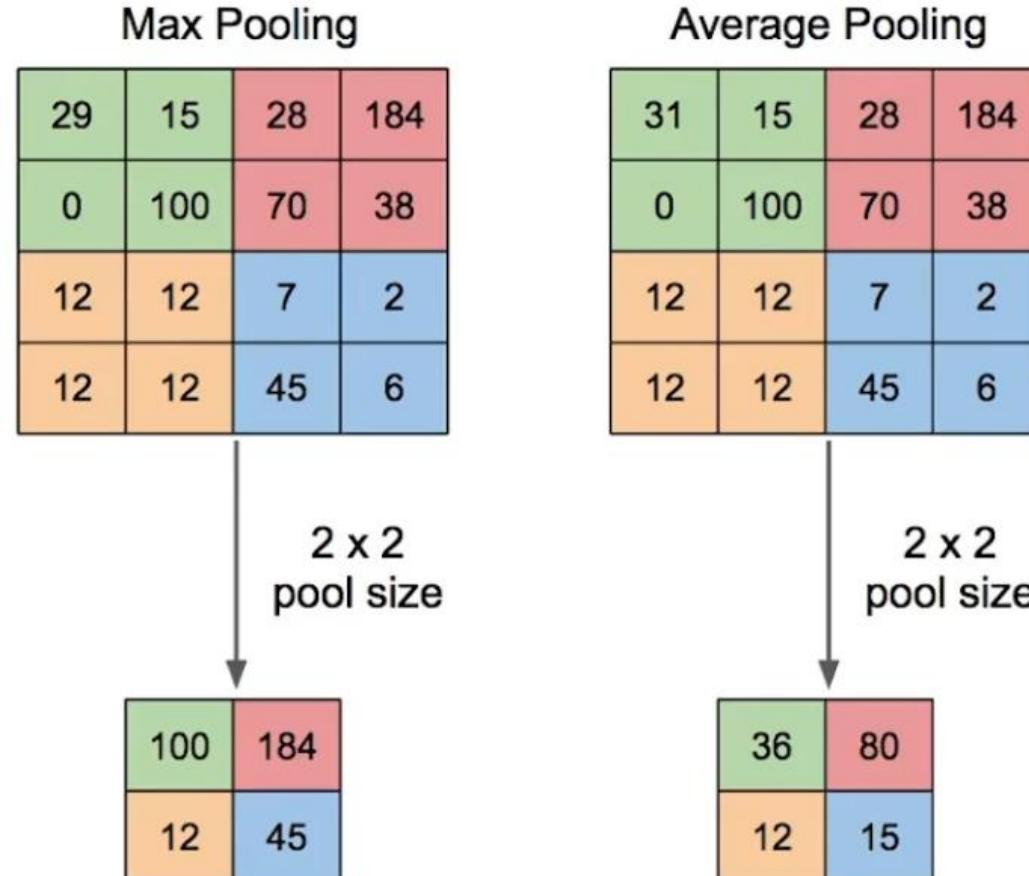
합성곱연산(Convolution)

$$\text{Output size} = \left(\frac{H+2p-f}{s} + 1 \right) \times \left(\frac{W+2p-f}{s} + 1 \right) \times (\# \text{ of filters})$$



Pooling

- ✓ 문제점 : 고차원의 Tensor를 보다 Compact하게 축약해야 하지 않을까?
- ✓ Pooling : 일정 영역의 정보를 축약하는 역할



POOLING LAYER / POOL

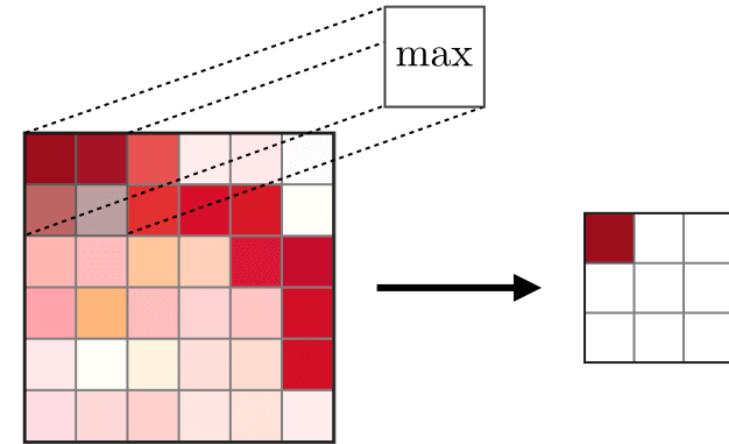
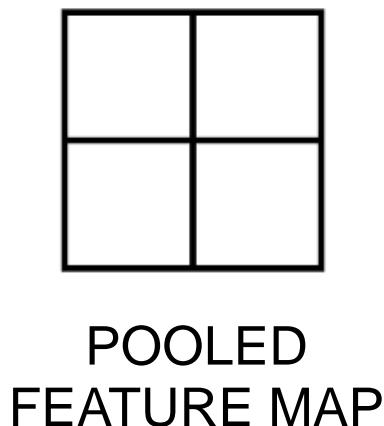
REDUCE THE SPATIAL SIZE

CONTROL OVERFITTING

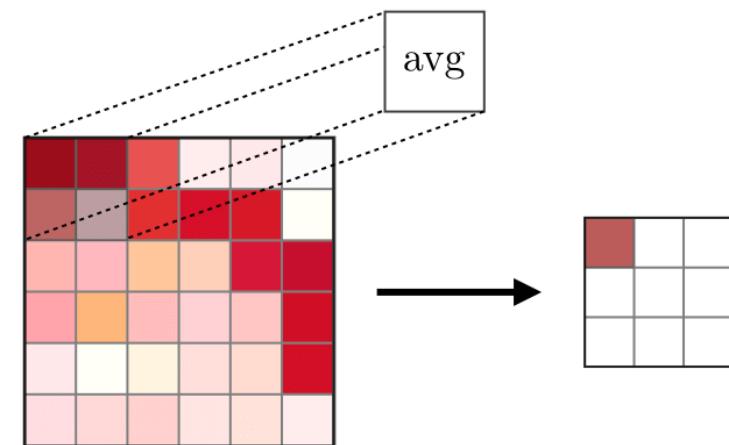
1	1	2	4
5	6	7	8
3	2	1	0
1	2	3	4

FEATURE MAP

TRAINING
NOT
HAPPEN

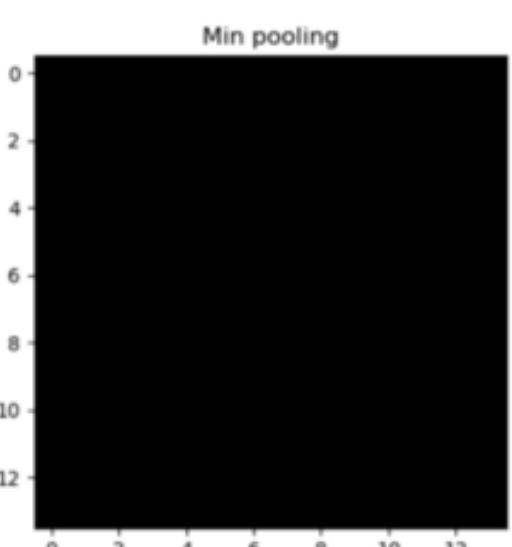
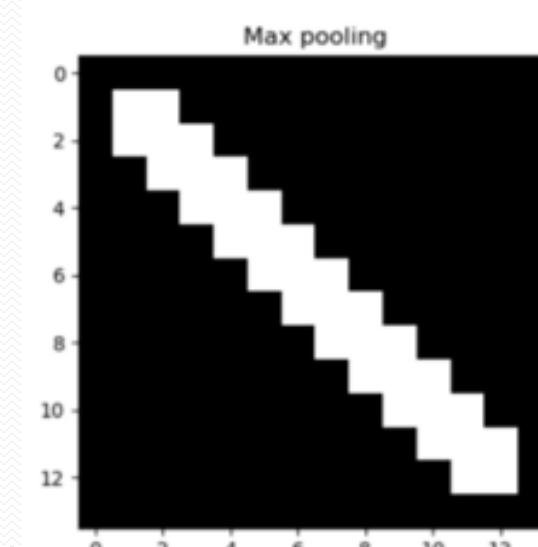
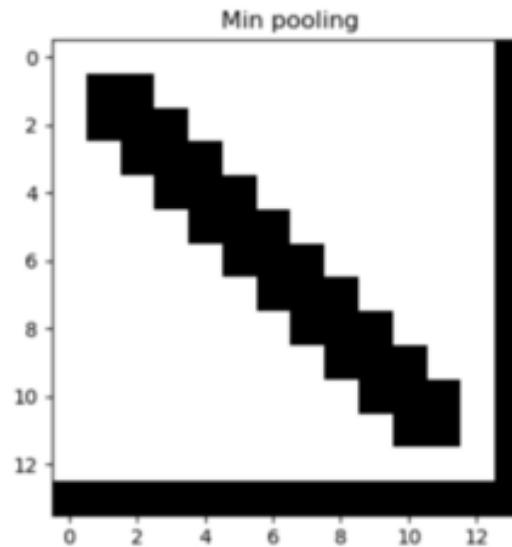
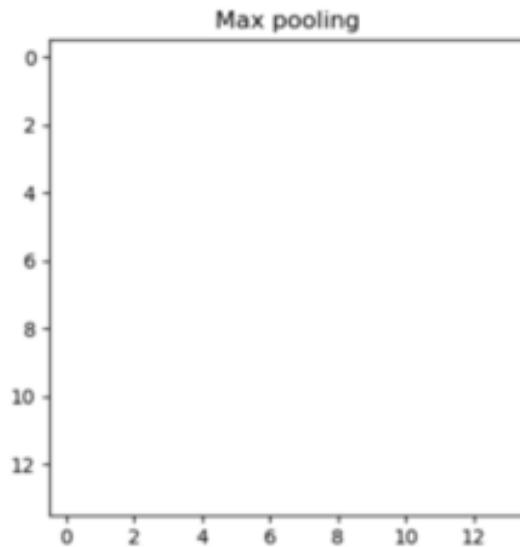
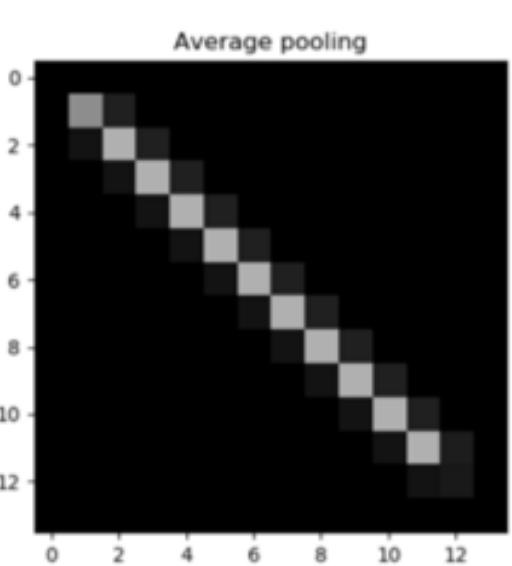
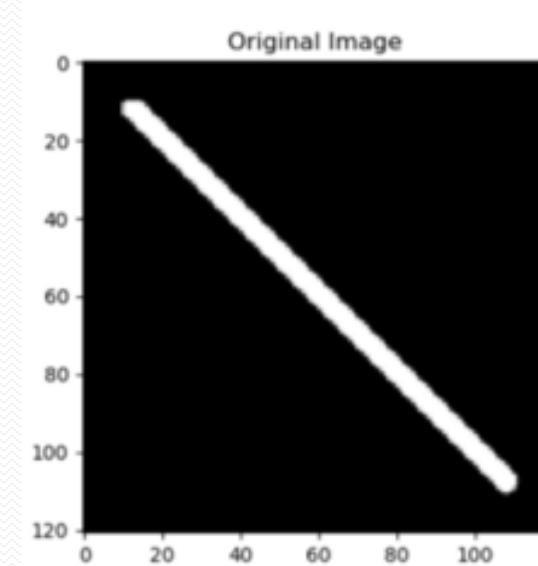
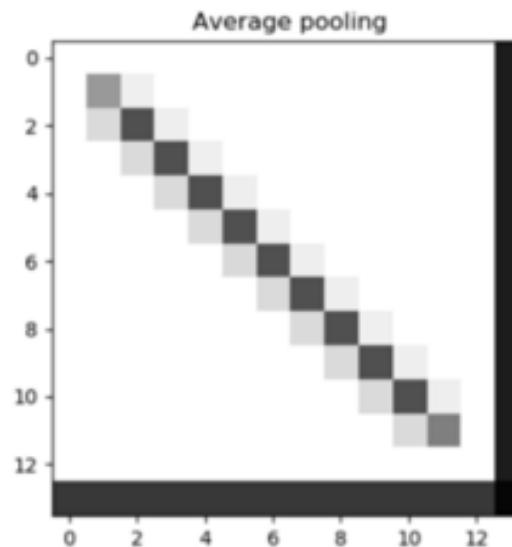
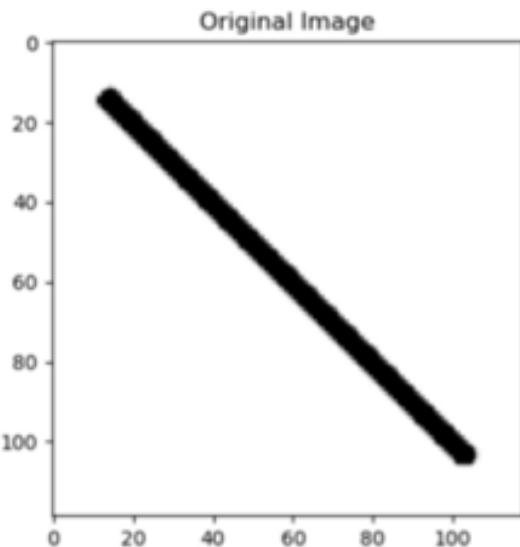


MAX POOLING

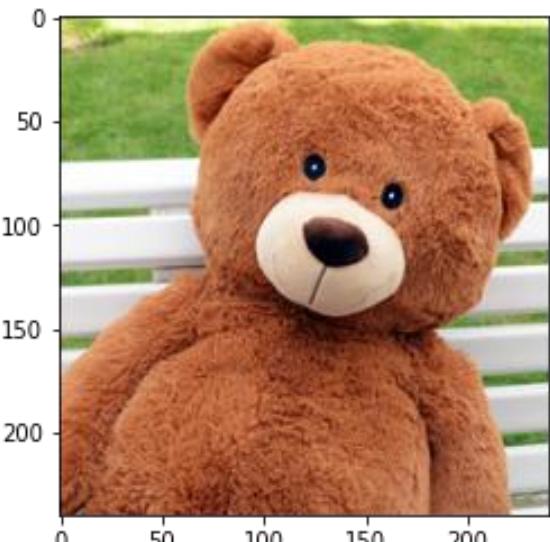


AVERAGE POOLING

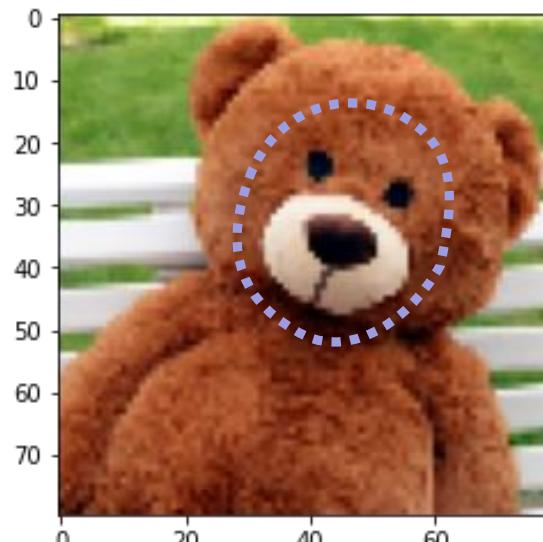
Compare Pooling algorithm (cont.)



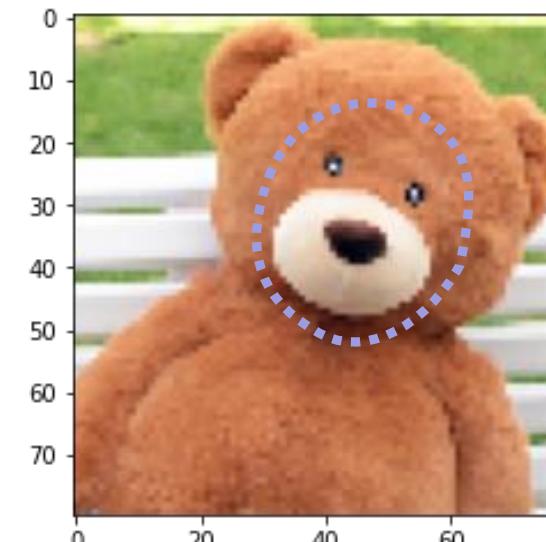
COMPARE POOLING ALGORITHM



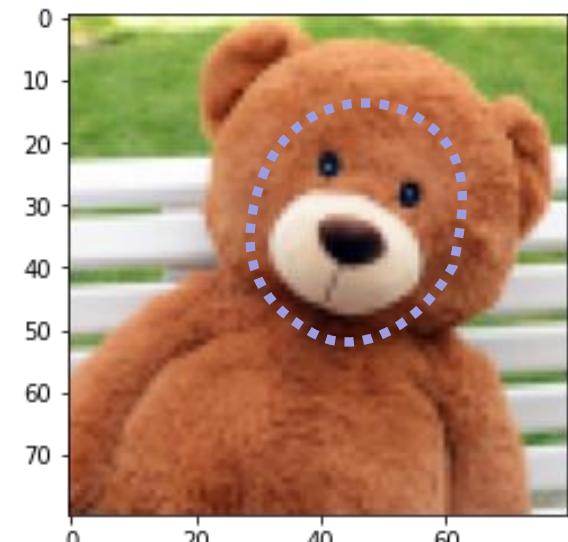
ORIGINAL



MIN

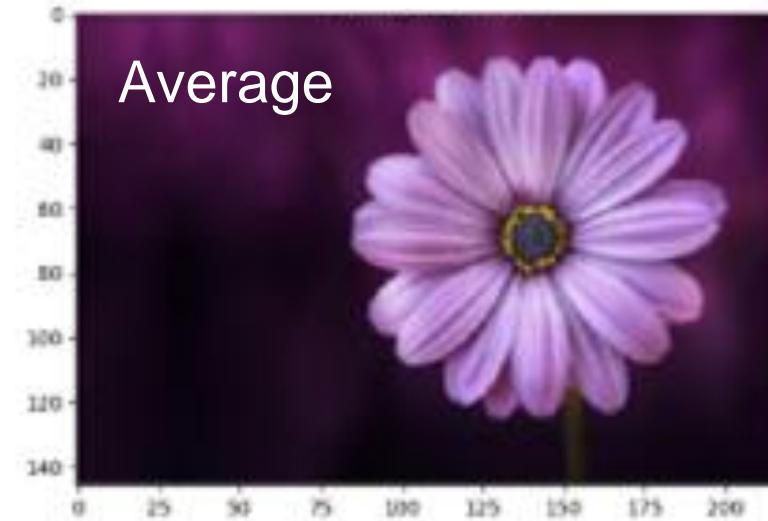


MAX



AVERAGE

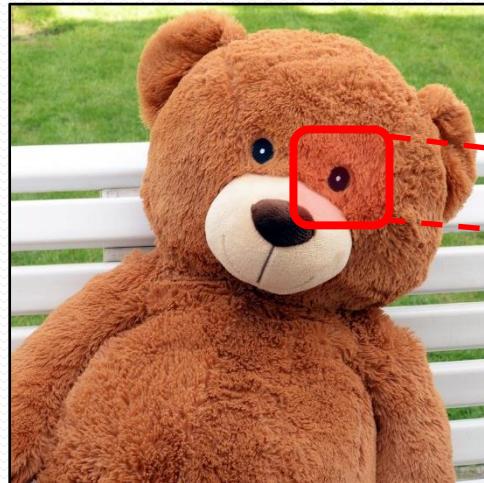
Compare Pooling algorithm



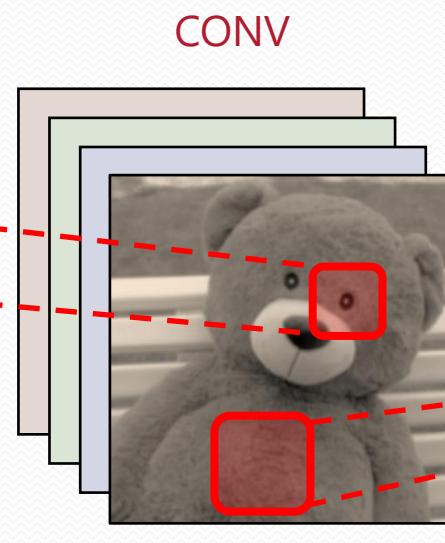
CONVOLUTIONAL NEURAL NETWORK

SPECIFIC TYPE OF NEURAL NETWORKS

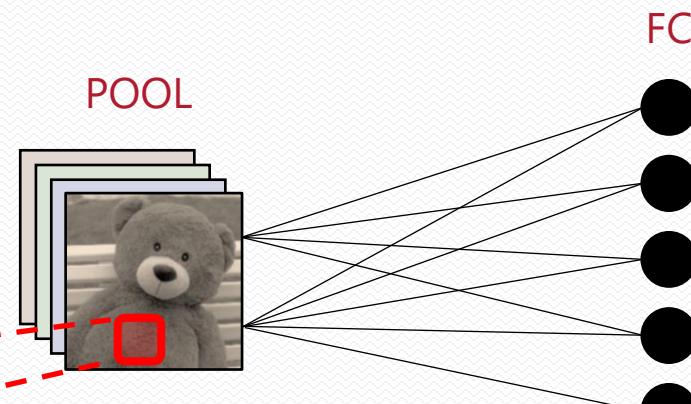
MAINLY USED FOR IMAGE CLASSIFICATION
GENERALLY COMPOSED OF THE 3 LAYERS



INPUT IMAGE



CONVOLUTIONS



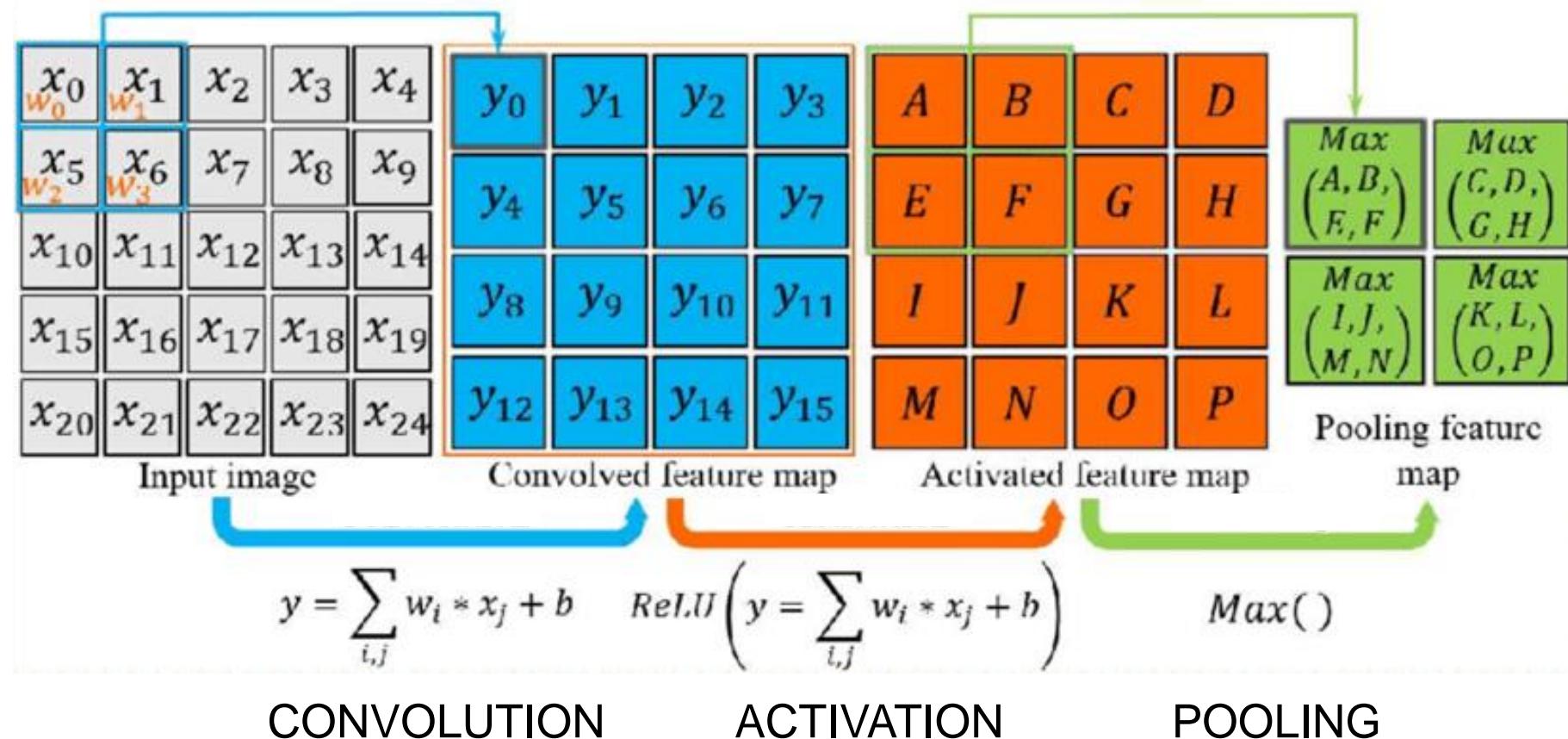
POOLING

FC



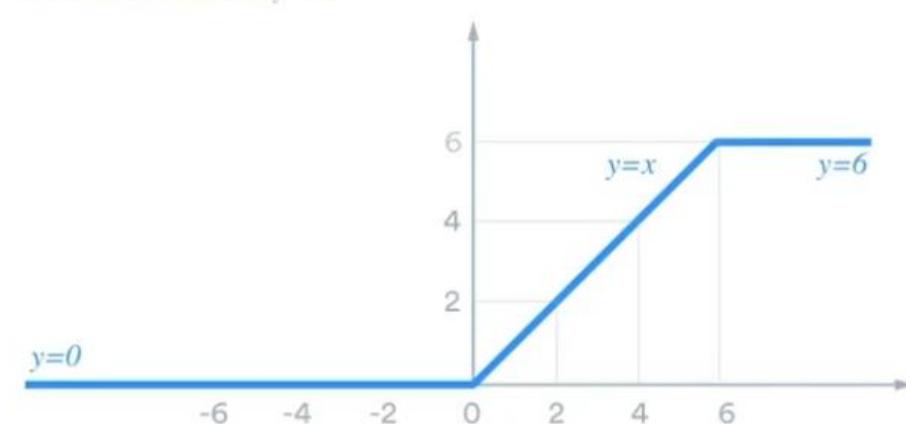
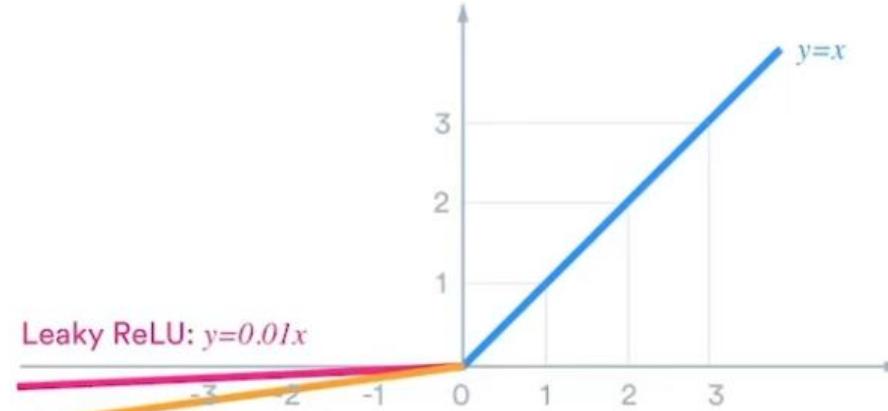
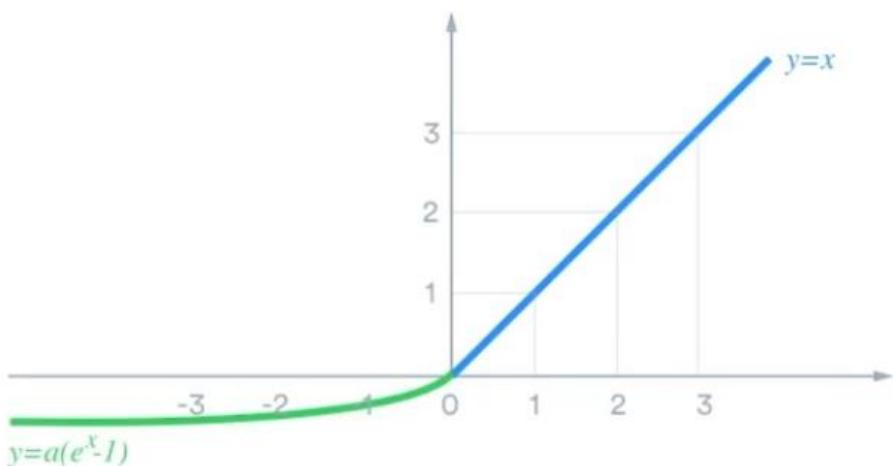
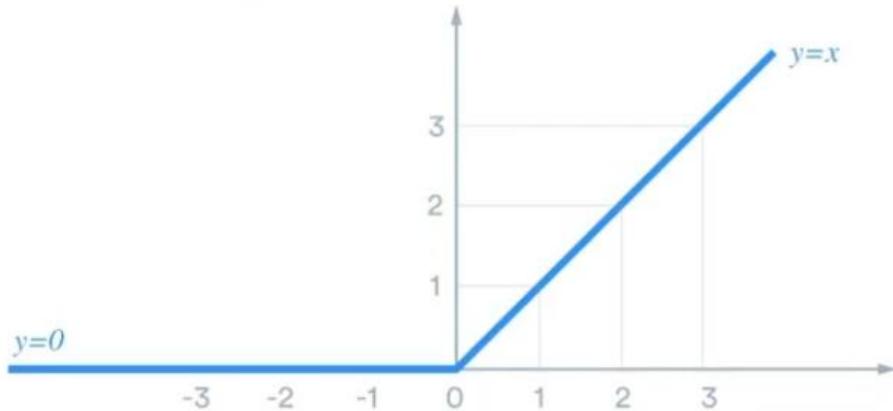
FULLY CONNECTED

PUT ALL TOGETHER



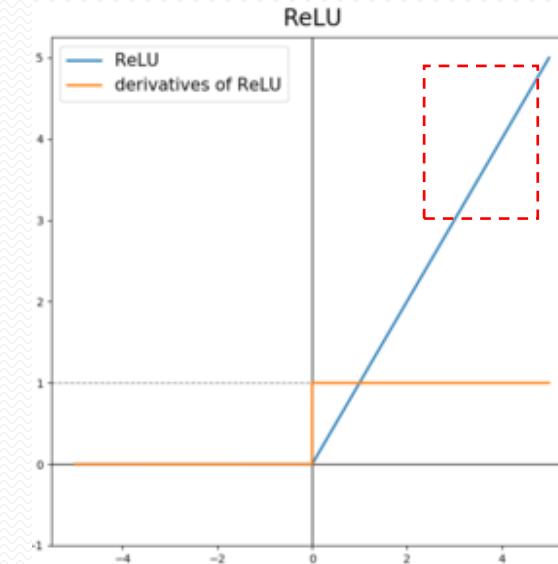
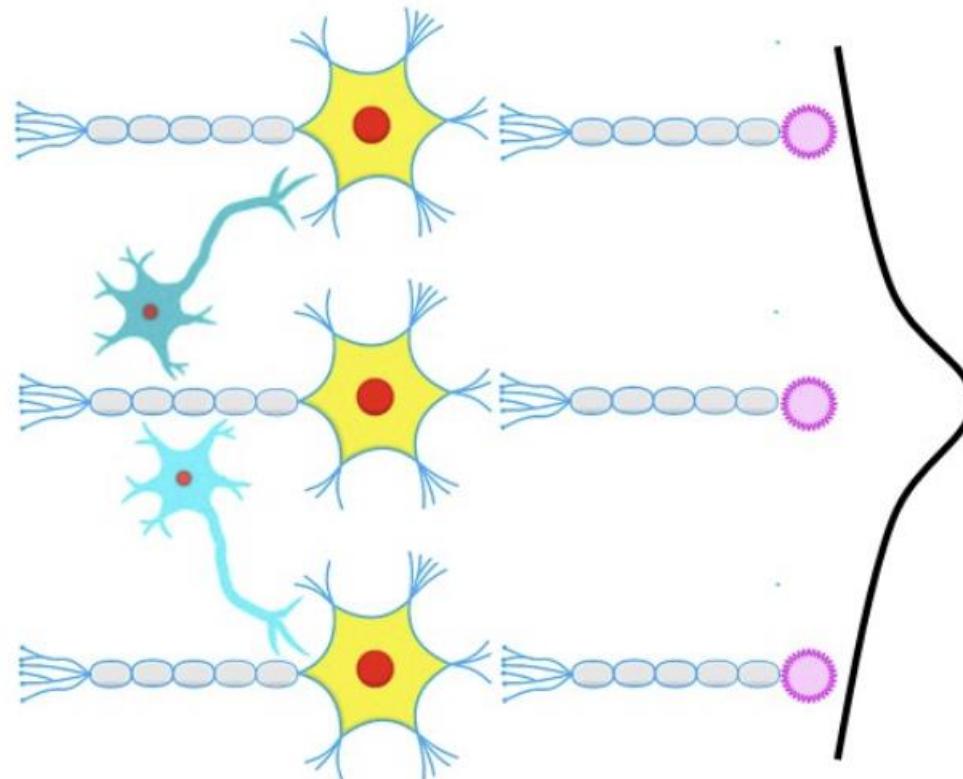
Activation

ReLU(Rectified Linear Unit) 사용 → 최근에는 변형 사용



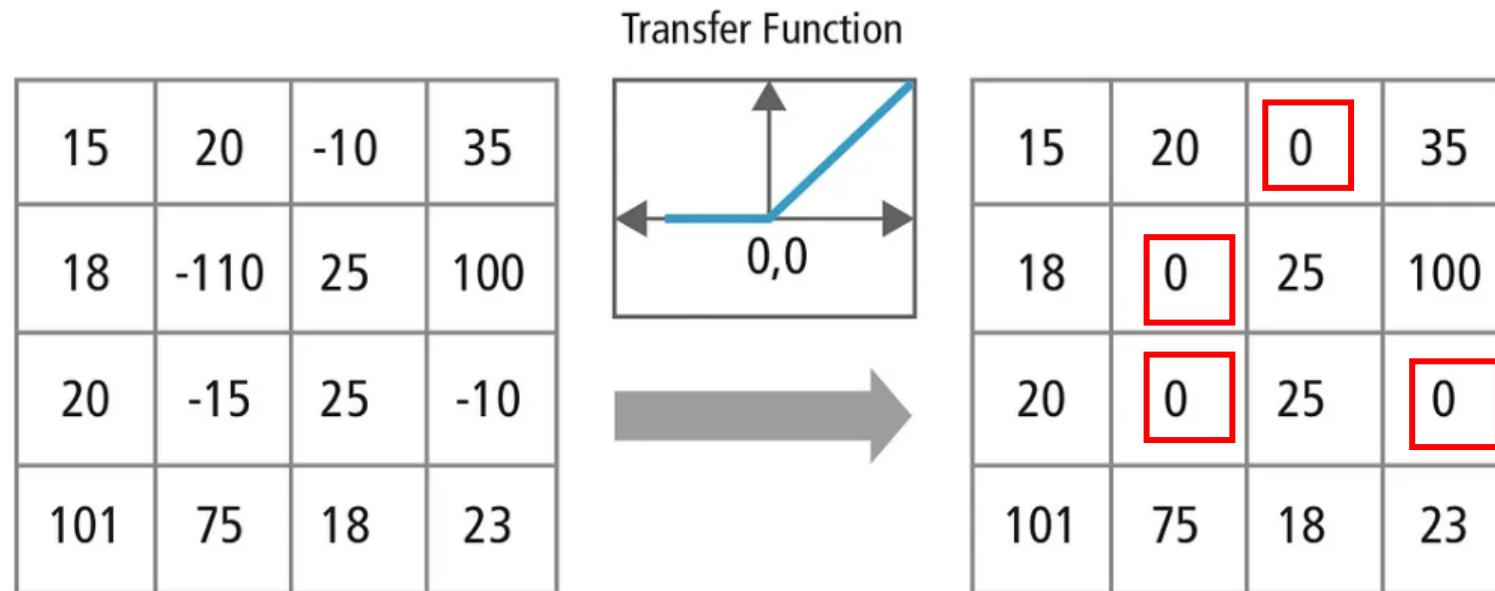
Activation

입출력단을 갖는 2개의 신경세포에서 각각의 신경세포의 입력 또는 출력이 서로 다른 신경세포의 입력 또는 출력에 의해 억제되는 현상.



Activation

- ✓ Convolution을 통해 학습된 값들의 비선형 변환을 수행
- ✓ 대부분 Rectified Linear Unit (ReLU)을 사용

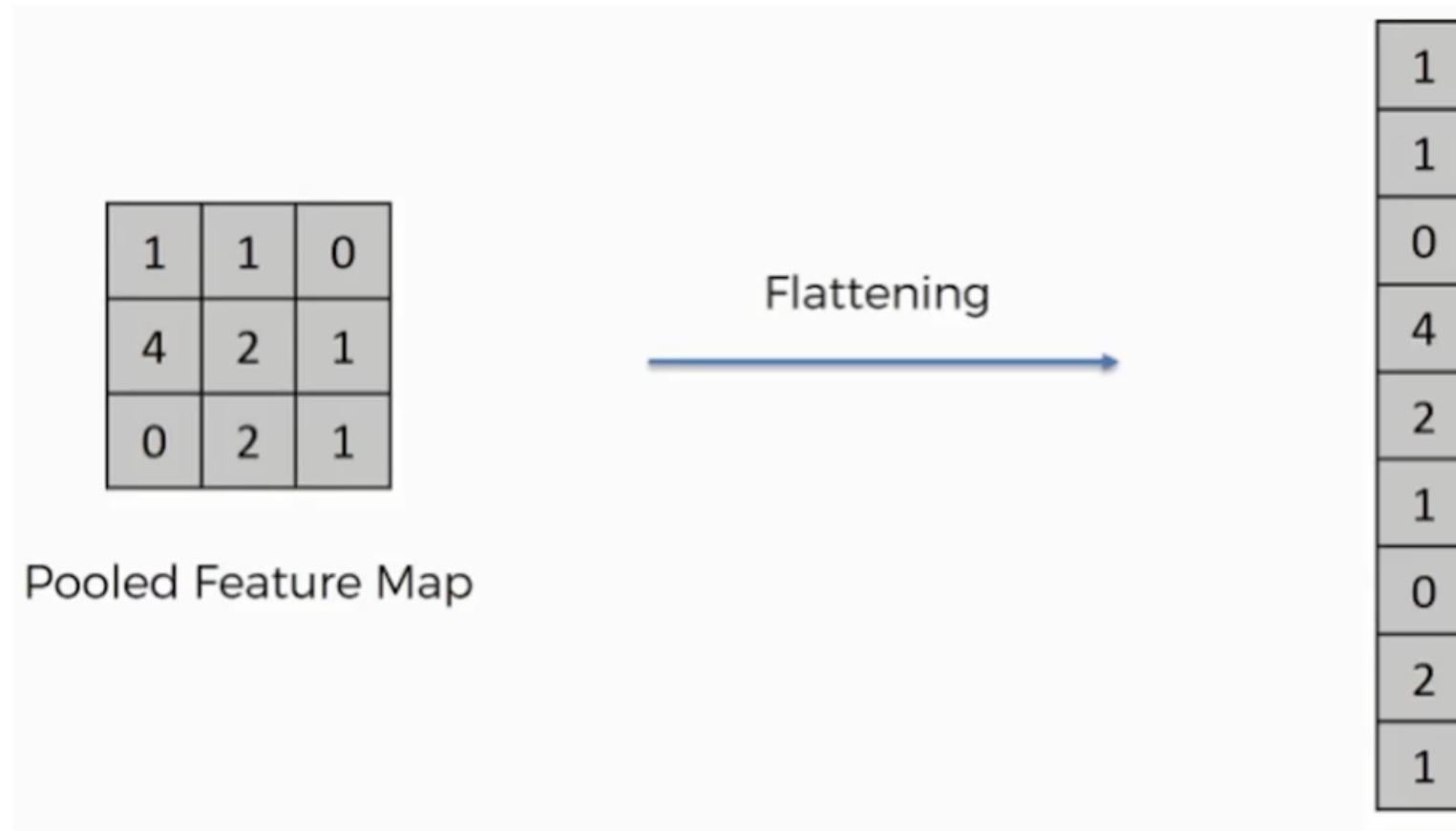


$$\text{Output} = \text{Max}(zero, \text{Input})$$

<https://medium.com/data-science-group-iitr/building-a-convolutional-neural-network-in-python-with-tensorflow-d251c3ca8117>

평탄화 : Flatten

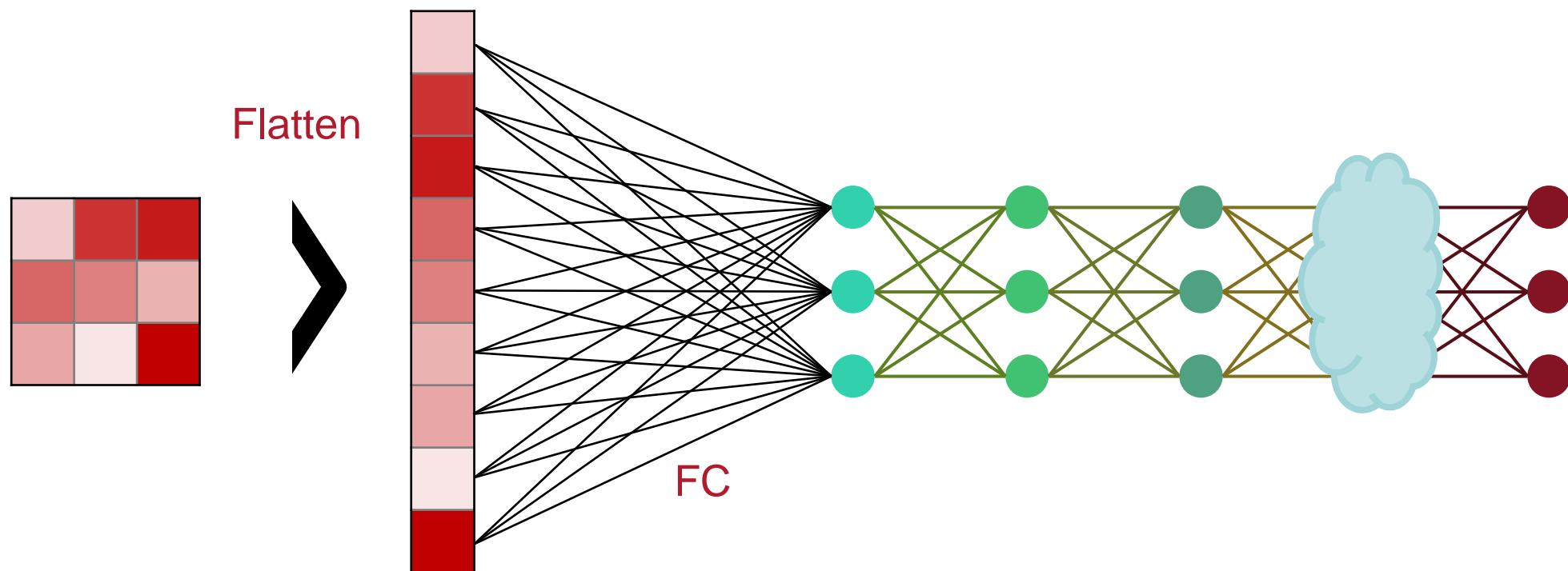
✓ 2차원/3차원의 Matrix/Tensor 구조를 1차원의 Vector로 변환하는 과정



FULLY CONNECTED LAYER / FC

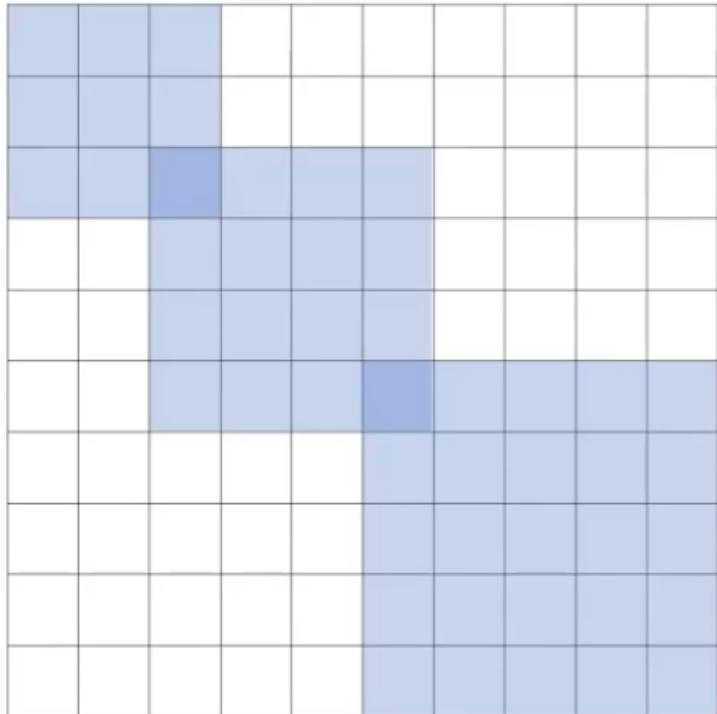
WORKS ON FLATTENED INPUTS WHERE EACH INPUT CONNECTS TO ALL NEURONS

FC LAYERS ARE USUALLY FOUND TOWARDS THE END OF CNN ARCHITECTURES



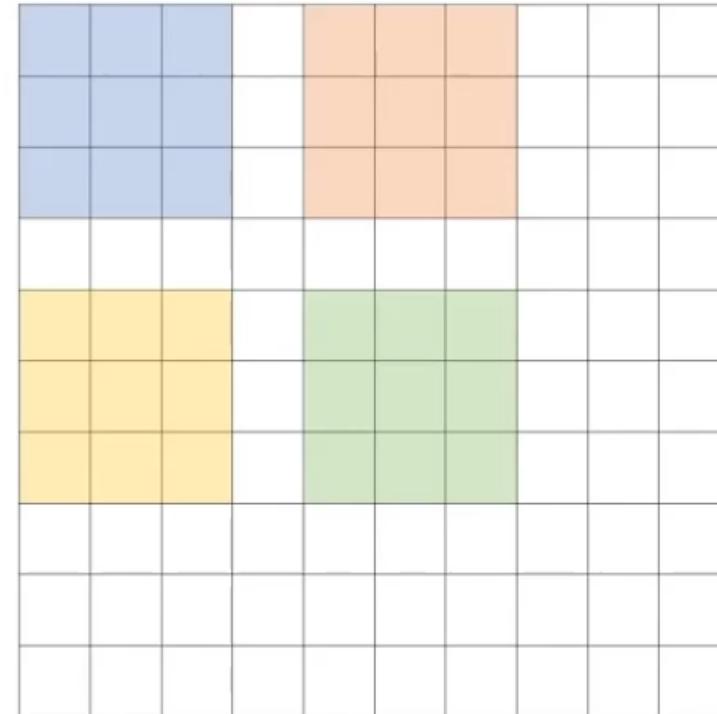
CNN의 하이퍼 파라미터

[Convolution Filter의 크기]



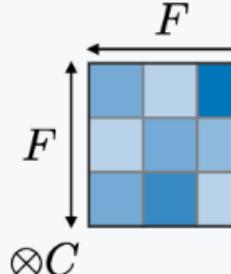
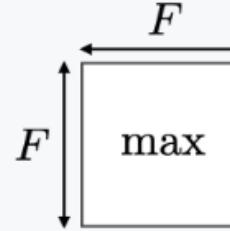
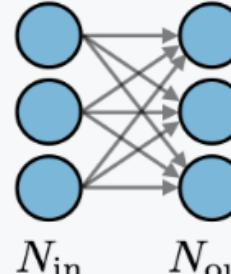
- Convolution Filter의 가로/세로 크기
- 이미지 데이터에는 주로 2-d Conv를 사용
(가로와 세로의 크기가 같은 Conv 2*2, 3*3 등)

[Convolution Filter의 수]



- Convolution Filter의 개수
- 이 수가 많을 수록 동일 영역에서 다양한 형태의 특질을 추출할 수 있음

Determine the N of parameters

	CONV	POOL	FC
Illustration	 $F \times F \times C \times K$	 $F \times F \times \text{max}$	 $N_{\text{in}} \times N_{\text{out}}$
Input size	$I \times I \times C$	$I \times I \times C$	N_{in}
Output size	$O \times O \times K$	$O \times O \times C$	N_{out}
Number of parameters	$(F \times F \times C + 1) \cdot K$	0	$(N_{\text{in}} + 1) \times N_{\text{out}}$
Remarks	<ul style="list-style-type: none">One bias parameter per filterIn most cases, $S < F$A common choice for K is $2C$	<ul style="list-style-type: none">Pooling operation done channel-wiseIn most cases, $S = F$	<ul style="list-style-type: none">Input is flattenedOne bias parameter per neuronThe number of FC neurons is free of structural constraints
S: Stride			

SIMPLE CNN EXAMPLE WITH KERAS

```
model = Sequential(name="Simple-CNN")
model.add(
    Conv2D(filters=8,
           kernel_size=(3, 3),
           padding="same",
           activation='relu',
           input_shape=(28, 28, 1)))
model.add(MaxPooling2D(pool_size=(2, 2)))

model.add(Conv2D(filters=16, kernel_size=(3, 3), activation='relu'))
model.add(MaxPooling2D(pool_size=(2, 2)))

model.add(Conv2D(filters=16, kernel_size=(3, 3), activation='relu'))

model.add(Flatten())
model.add(Dense(64, activation='sigmoid'))
model.add(Dense(10, activation='softmax'))
```

SIMPLE CNN EXAMPLE WITH KERAS / CONT.

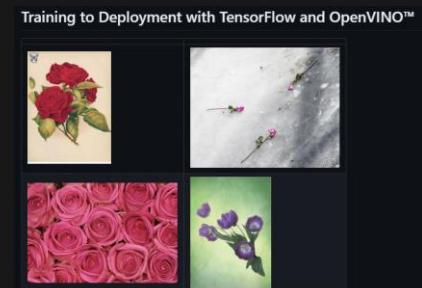
Model: "Simple-CNN"

Layer (type)	Output Shape	Param #
conv2d (Conv2D)	(None, 28, 28, 8)	80
max_pooling2d (MaxPooling2D)	(None, 14, 14, 8)	0
conv2d_1 (Conv2D)	(None, 12, 12, 16)	1168
max_pooling2d_1 (MaxPooling2D)	(None, 6, 6, 16)	0
conv2d_2 (Conv2D)	(None, 4, 4, 16)	2320
flatten (Flatten)	(None, 256)	0
dense (Dense)	(None, 64)	16448
dense_1 (Dense)	(None, 10)	650

SIMPLE CNN EXAMPLE WITH KERAS

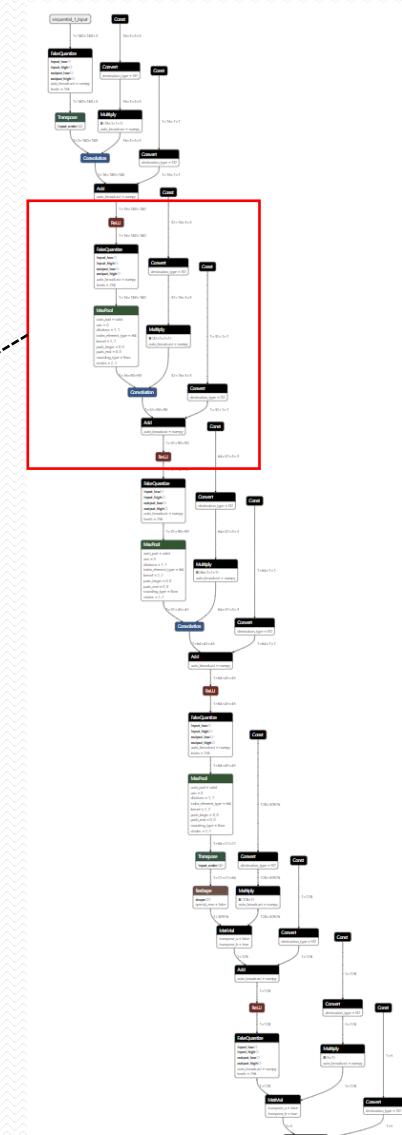
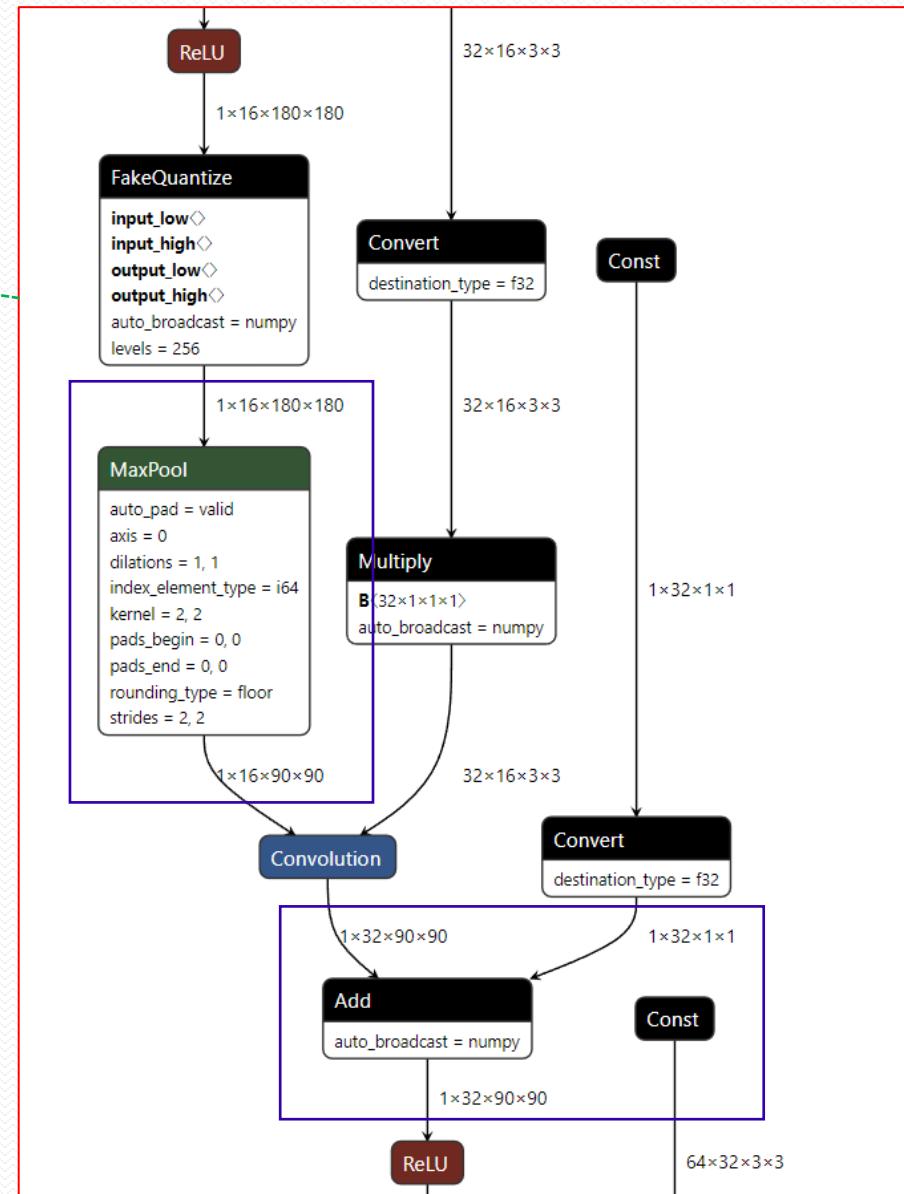
Layer (type)	Output Shape	Param #
<hr/>		
rescaling_1 (Rescaling)	(None, 180, 180, 3)	0
conv2d (Conv2D)	(None, 180, 180, 16)	448
max_pooling2d (MaxPooling2D)	(None, 90, 90, 16)	0
conv2d_1 (Conv2D)	(None, 90, 90, 32)	4640
max_pooling2d_1 (MaxPooling2D)	(None, 45, 45, 32)	0
conv2d_2 (Conv2D)	(None, 45, 45, 64)	18496
max_pooling2d_2 (MaxPooling2D)	(None, 22, 22, 64)	0
flatten (Flatten)	(None, 30976)	0
dense (Dense)	(None, 128)	3965056
dense_1 (Dense)	(None, 5)	645
<hr/>		
Total params: 3,989,285		
Trainable params: 3,989,285		
Non-trainable params: 0		

```
# 모델 만들기
model = tf.keras.Sequential()
model.add(
    tf.keras.layers.Rescaling(
        1./255, input_shape=(IMAGE_WIDTH, IMAGE_HEIGHT, 3)))
model.add(
    tf.keras.layers.Conv2D(16, 3, padding="same", activation="relu"))
model.add(
    tf.keras.layers.MaxPool2D())
model.add(
    tf.keras.layers.Conv2D(32, 3, padding="same", activation="relu"))
model.add(
    tf.keras.layers.MaxPool2D())
model.add(
    tf.keras.layers.Conv2D(64, 3, padding="same", activation="relu"))
model.add(
    tf.keras.layers.MaxPool2D())
model.add(
    tf.keras.layers.Flatten())
model.add(
    tf.keras.layers.Dropout(0.2))
)
model.add(
    tf.keras.layers.Dense(128, activation="relu"))
)
model.add(
    tf.keras.layers.Dense(5))
model.summary()
```



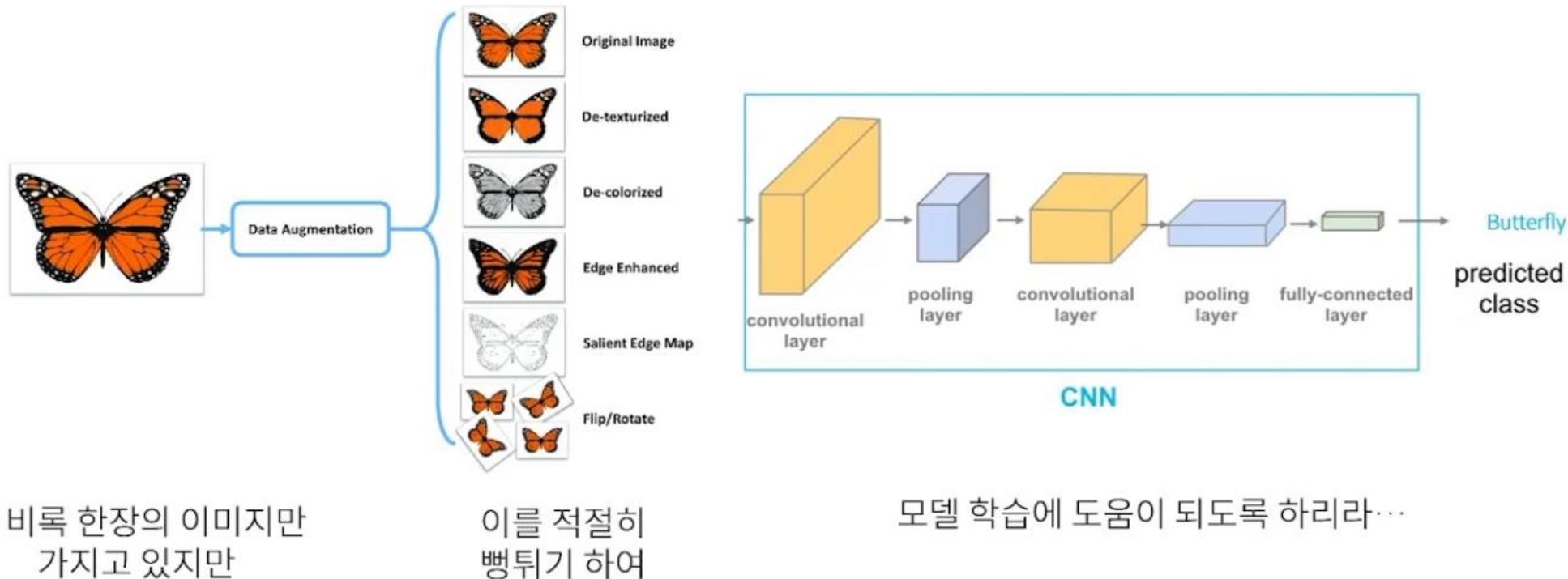
<https://netron.app/>

Layer (type)	Output Shape	Param #
<hr/>		
rescaling_1 (Rescaling)	(None, 180, 180, 3)	0
conv2d (Conv2D)	(None, 180, 180, 16)	448
max_pooling2d (MaxPooling2D)	(None, 90, 90, 16)	0
conv2d_1 (Conv2D)	(None, 90, 90, 32)	4640
max_pooling2d_1 (MaxPooling2D)	(None, 45, 45, 32)	0
conv2d_2 (Conv2D)	(None, 45, 45, 64)	18496
max_pooling2d_2 (MaxPooling2D)	(None, 22, 22, 64)	0
flatten (Flatten)	(None, 30976)	0
dense (Dense)	(None, 128)	3965056
dense_1 (Dense)	(None, 5)	645
<hr/>		
Total params:	3,989,285	
Trainable params:	3,989,285	
Non-trainable params:	0	

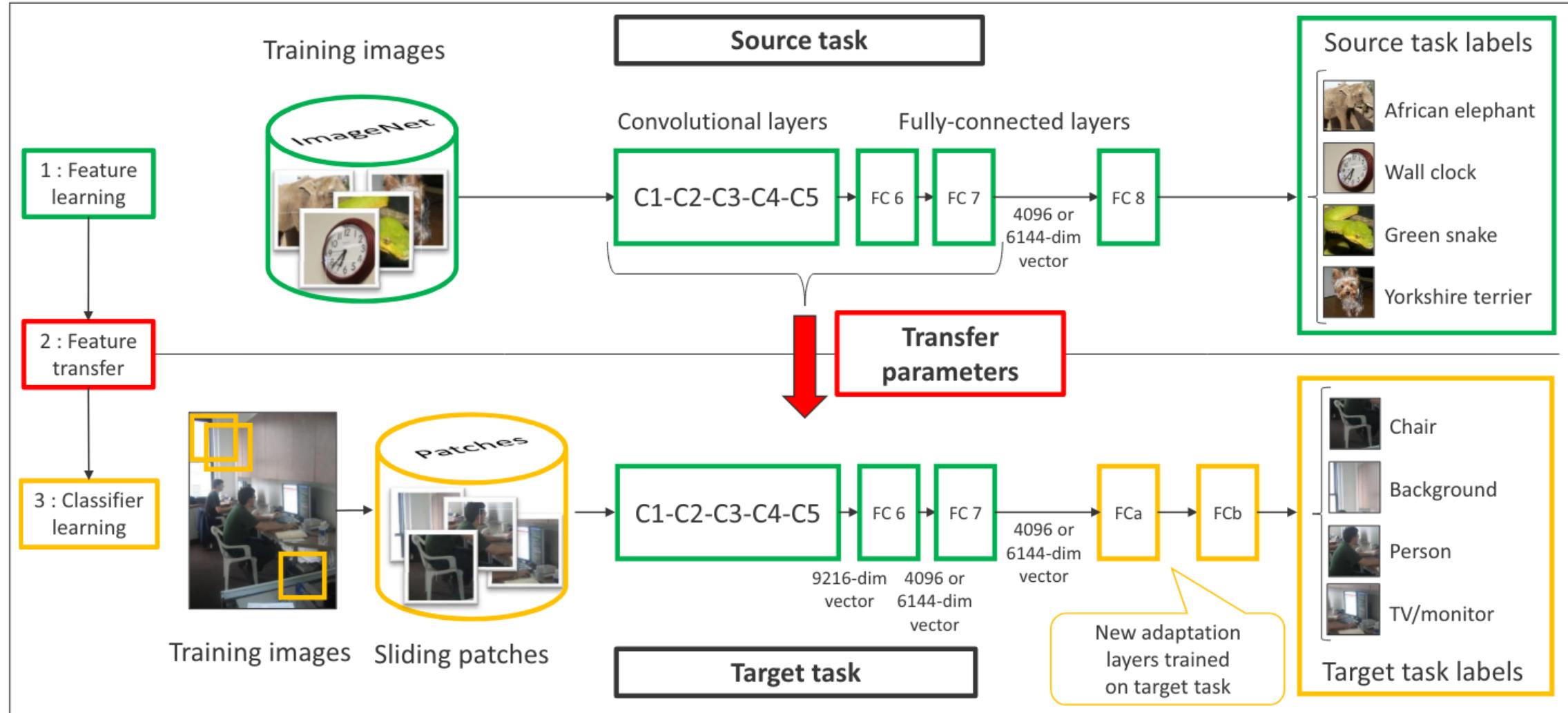


Data Augmentation

- ✓ 현실에서 충분히 확보하지 못한 데이터에 대해 적절한 변형을 가하여 인위적으로 데이터를 증강시키자!
- ✓ 그럴듯하게 증강시키면(이미지 회전, 상하 반전, 약간의 노이즈 추가 등) 오히려 보다 다양한 형태를 학습하게 됨으로써 더 높은 정확도를 얻을 수 있지 않을까?

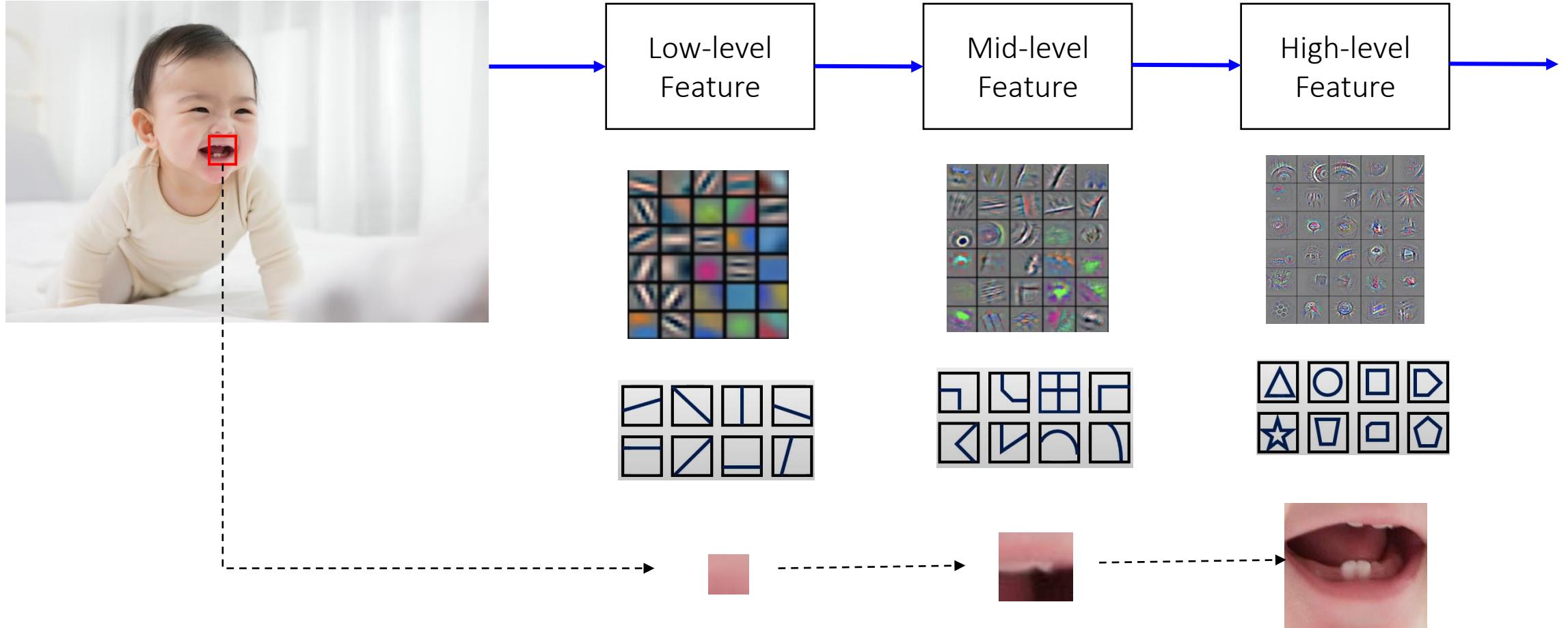


Transfer learning



Oquab, M., Bottou, L., Laptev, I., & Sivic, J. (2014). Learning and transferring mid-level image representations using convolutional neural networks. In Proceedings of the IEEE conference on computer vision and pattern recognition (pp. 1717-1724).

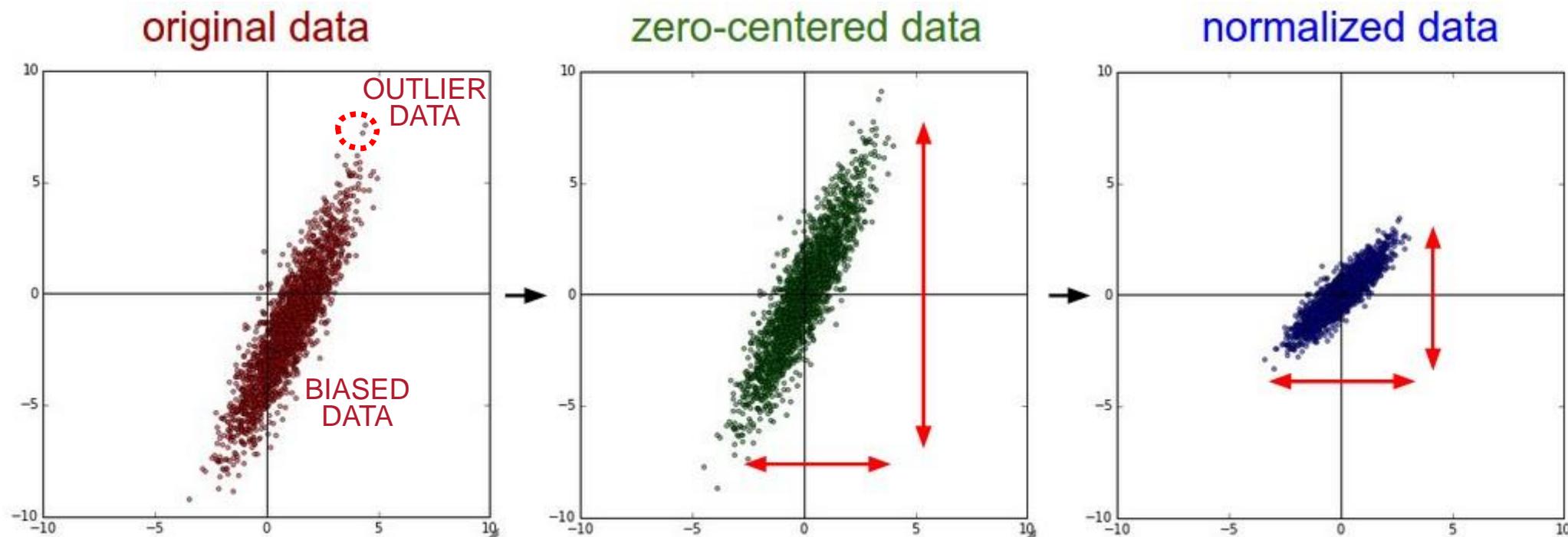
Transfer learning



BATCH NORMALIZATION

NORMALIZING THE DATA DIMENSIONS

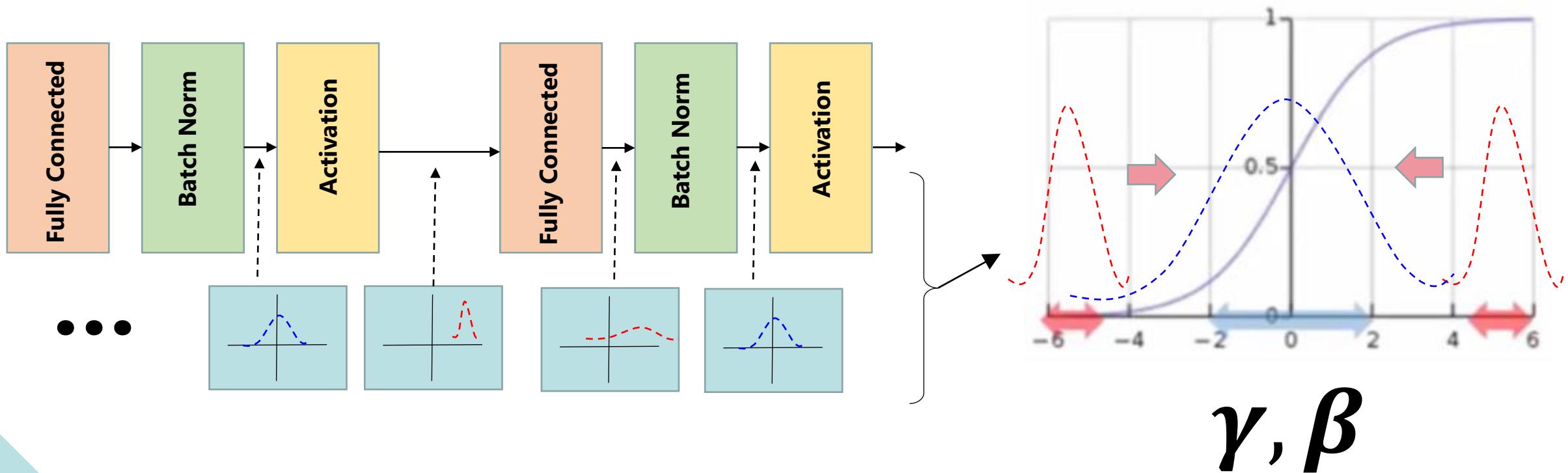
Method to help the network learn better



BATCH NORMALIZATION

NORMALIZING THE DATA DIMENSIONS

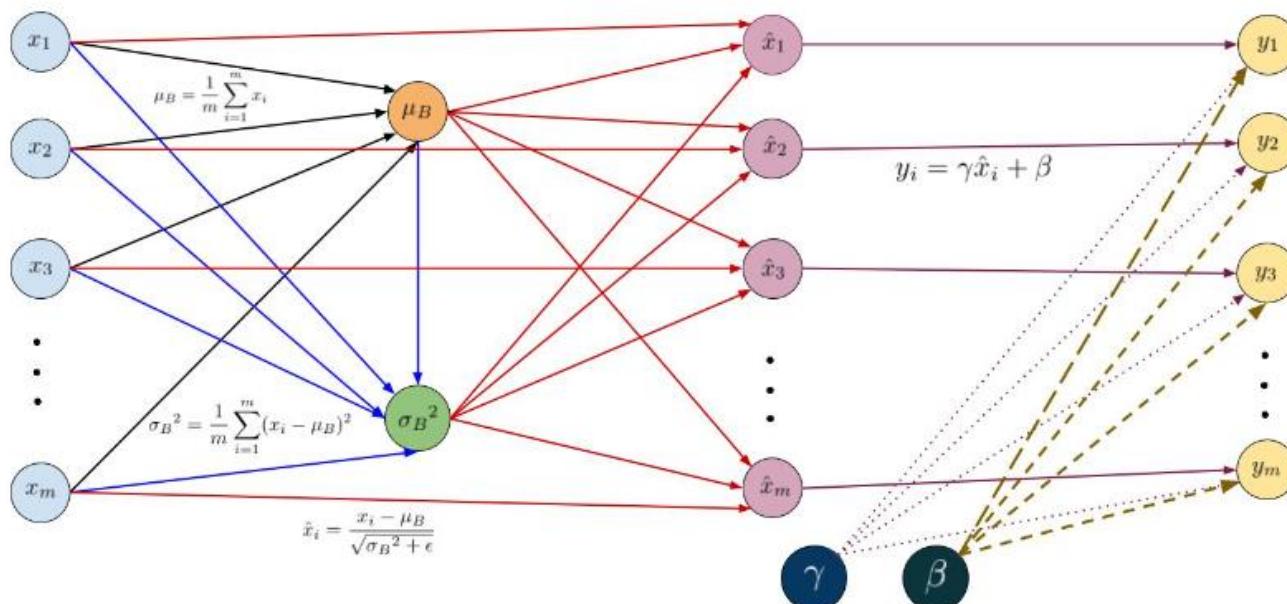
Method to help the network learn better



BATCH NORMALIZATION

NORMALIZING THE DATA DIMENSIONS

Method to help the network learn better



Input: Values of x over a mini-batch: $\mathcal{B} = \{x_1 \dots m\}$;
Parameters to be learned: γ, β

Output: $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{scale and shift}$$

BATCH NORMALIZATION

NORMALIZING THE DATA DIMENSIONS

Method to help the network learn better

Input: Values of x over a mini-batch: $\mathcal{B} = \{x_1 \dots m\}$;
Parameters to be learned: γ, β

Output: $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{scale and shift}$$

```
1 def batchnorm_forward(x, gamma, beta, eps=1e-5):
2     N, D = x.shape
3
4     sample_mean = x.mean(axis=0)
5     sample_var = x.var(axis=0)
6
7     std = np.sqrt(sample_var + eps)
8     x_centered = x - sample_mean
9     x_norm = x_centered / std
10    out = gamma * x_norm + beta
11
12    cache = (x_norm, x_centered, std, gamma)
13
14    return out, cache
```

BATCH NORMALIZATION

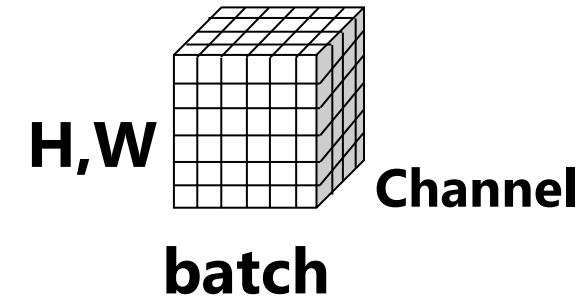
NORMALIZING THE DATA DIMENSIONS

Method to help the network learn better

$$\gamma \left(\frac{x_i - \mu}{\sqrt{\sigma^2 + \epsilon}} \right) + \beta \quad \rightarrow$$

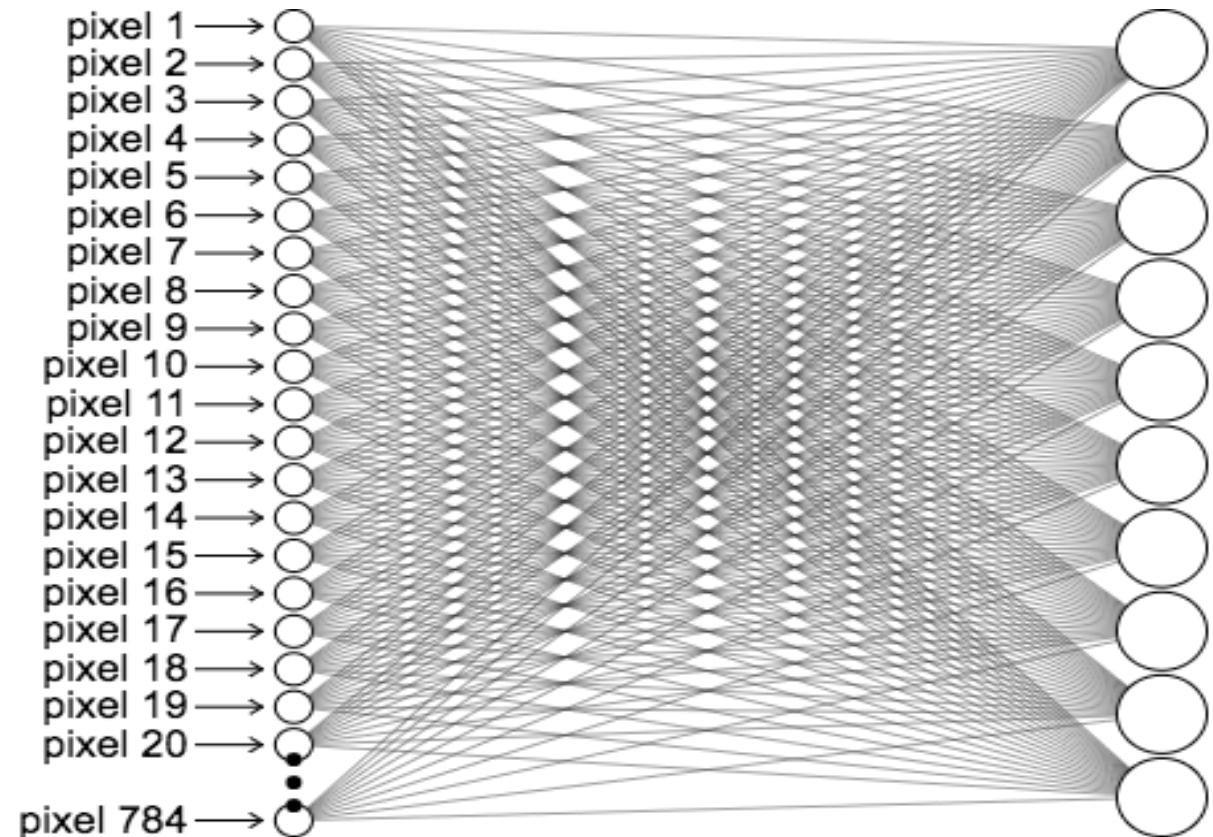
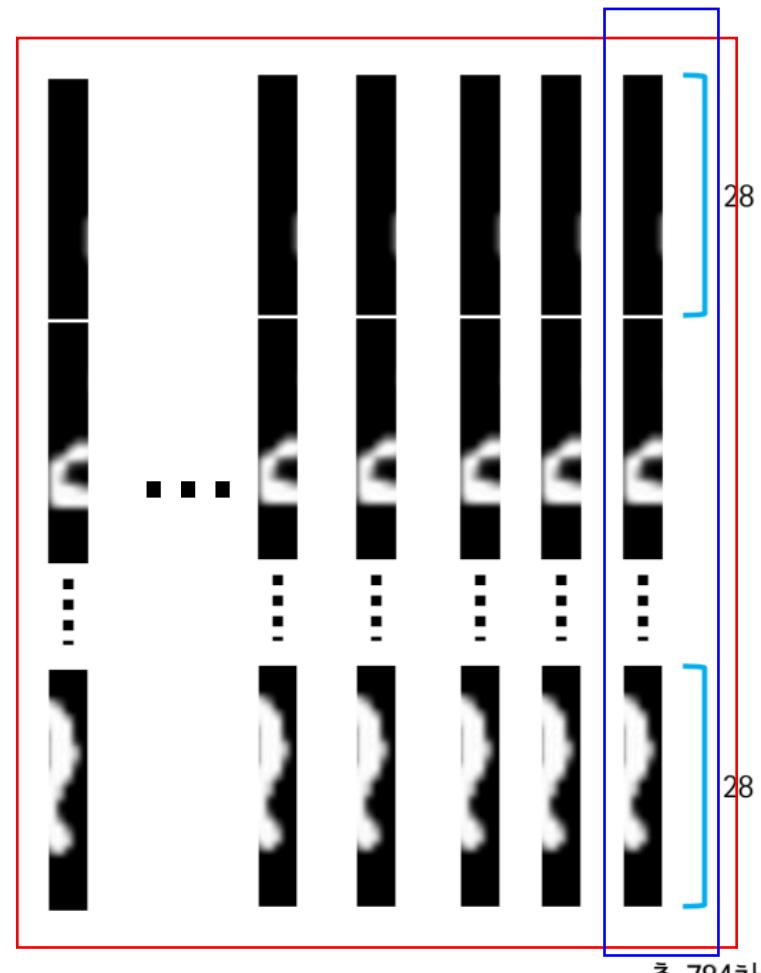
scale and shift
→ 분포재조정됨
- 평균: β
- 분산: γ^2

BATCH & LAYER NORMALIZATION



Batch Normalization

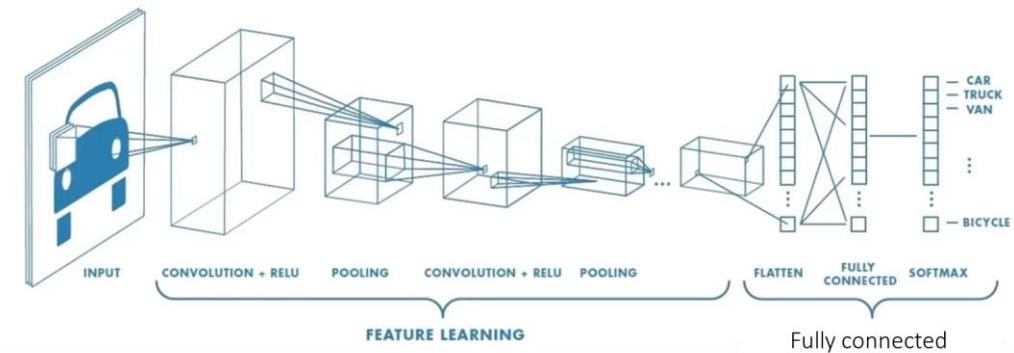
3 4 2 1 9 5 6 2 1 8
8 9 1 2 3 0 0 6 6 4
6 7 0 1 6 3 6 3 7 0
3 7 7 9 4 6 6 1 8 2
2 9 3 4 3 9 8 7 2 5
1 5 9 8 3 6 5 7 2 3
9 3 1 9 1 5 8 0 8 4
5 6 2 6 8 5 8 8 9 9
3 7 7 0 9 4 8 5 4 3
7 9 6 4 1 0 4 9 2 3



Layer Normalization

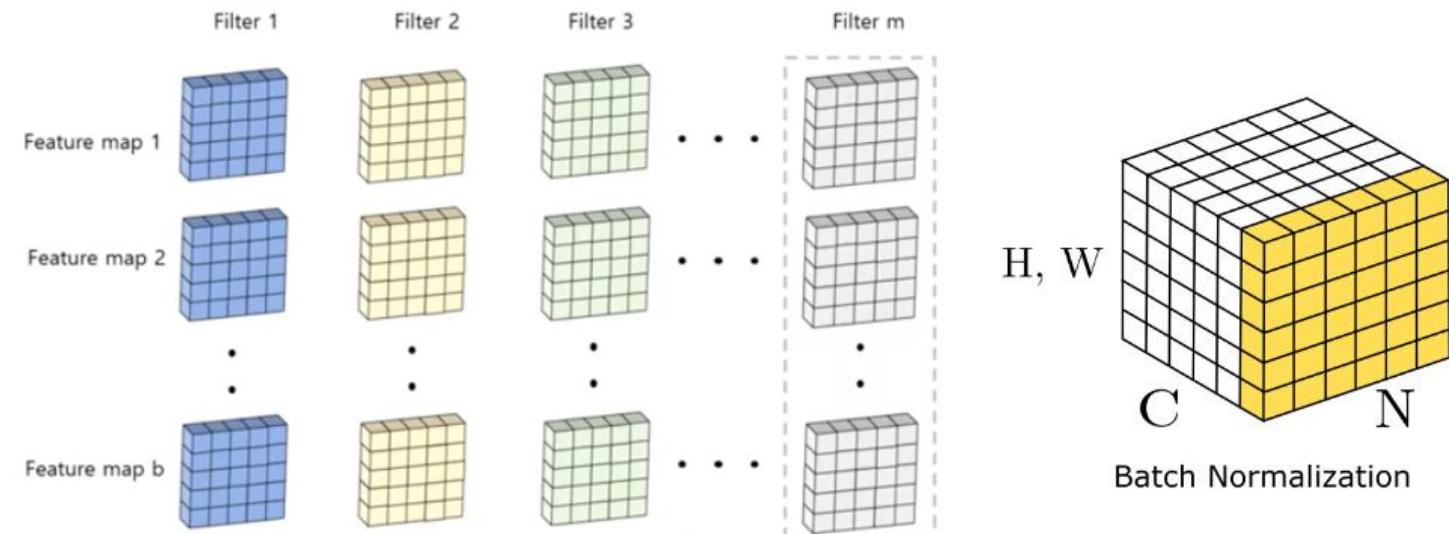
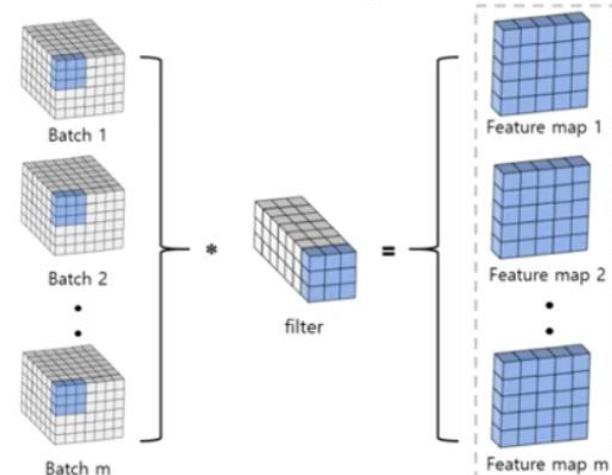
BATCH & LAYER NORMALIZATION

- 각 feature 별로, 한 mini-batch에 대한 summed input의 평균과 분산을 구한다.
- 일반화를 개선하고 overfitting을 감소시킨다.
- 학습의 수렴속도를 높이고 가중치 초기화에 대해 강건해진다.
- 학습의 안정성을 향상시키고 기울기 소실 문제를 완화한다.
- 피드포워드 네트워크 모델의 성능과 안정성 개선에 뛰어나다
- 미니배치 크기에 의존한다.
- RNN Model에 적용하기 힘들다



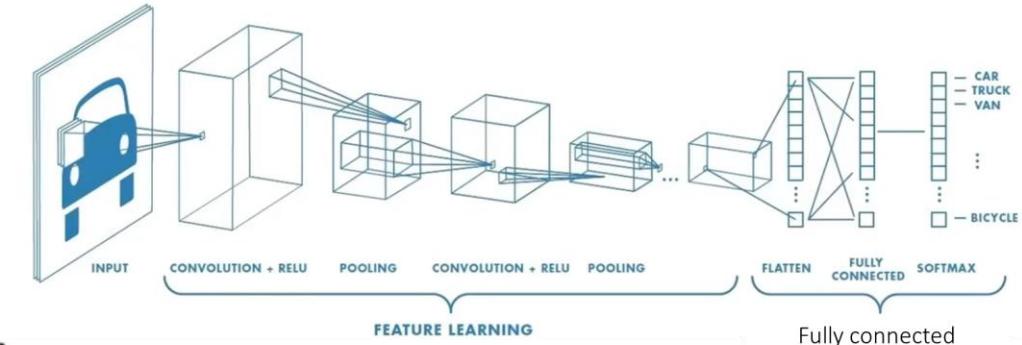
Batch Normalization

- m개의 입력 배치에 대해 m개의 feature map 도출

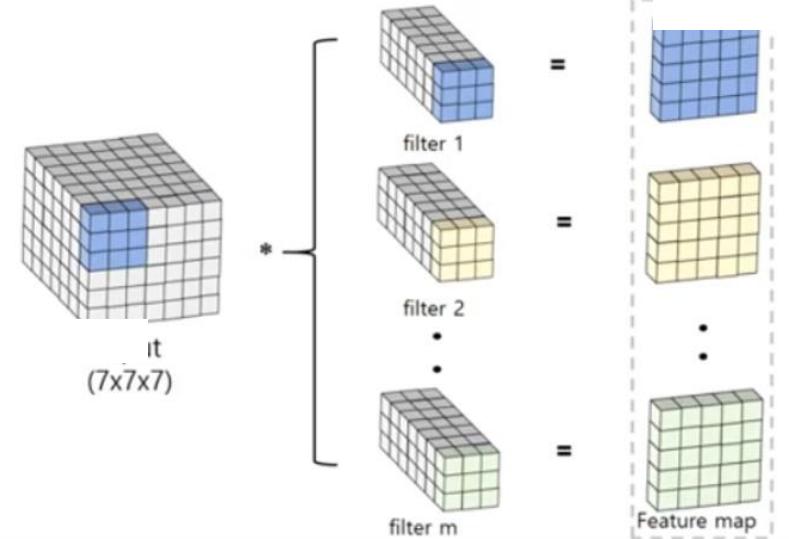


BATCH & LAYER NORMALIZATION

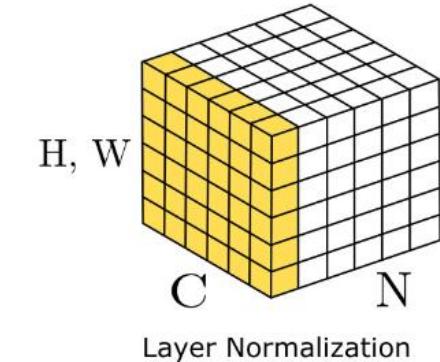
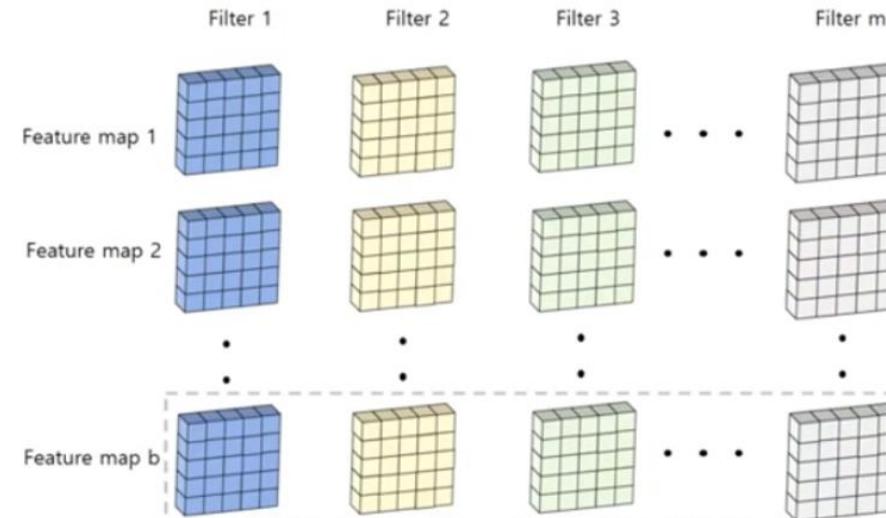
- Batch Normalization은 작은 배치 크기에서 극단적 결과를 내는데 반해, 작은 batch size에서도 효과적인 이용이 가능
- sequence에 따른 고정길이 정규화로 batch normalization에 비해 RNN 모델에 더 효과적인 방법이다.
- 일반화 성능 향상이 가능하다
- 피드포워드 네트워크 모델에서는 Batch Normalization 만큼 잘 동작하지 않을 수 있다.
- learning_rate, 가중치 초기화 같은 하이퍼파라미터 조정에 민감할 수 있다.
- 추가 계산 및 메모리 오버헤드가 발생할 수 있다



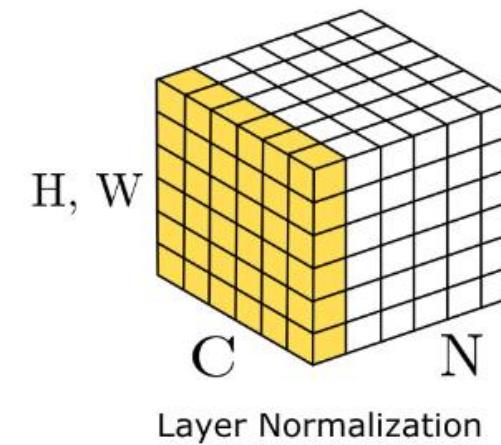
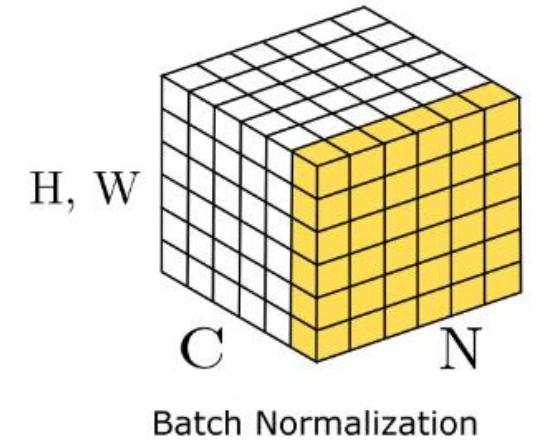
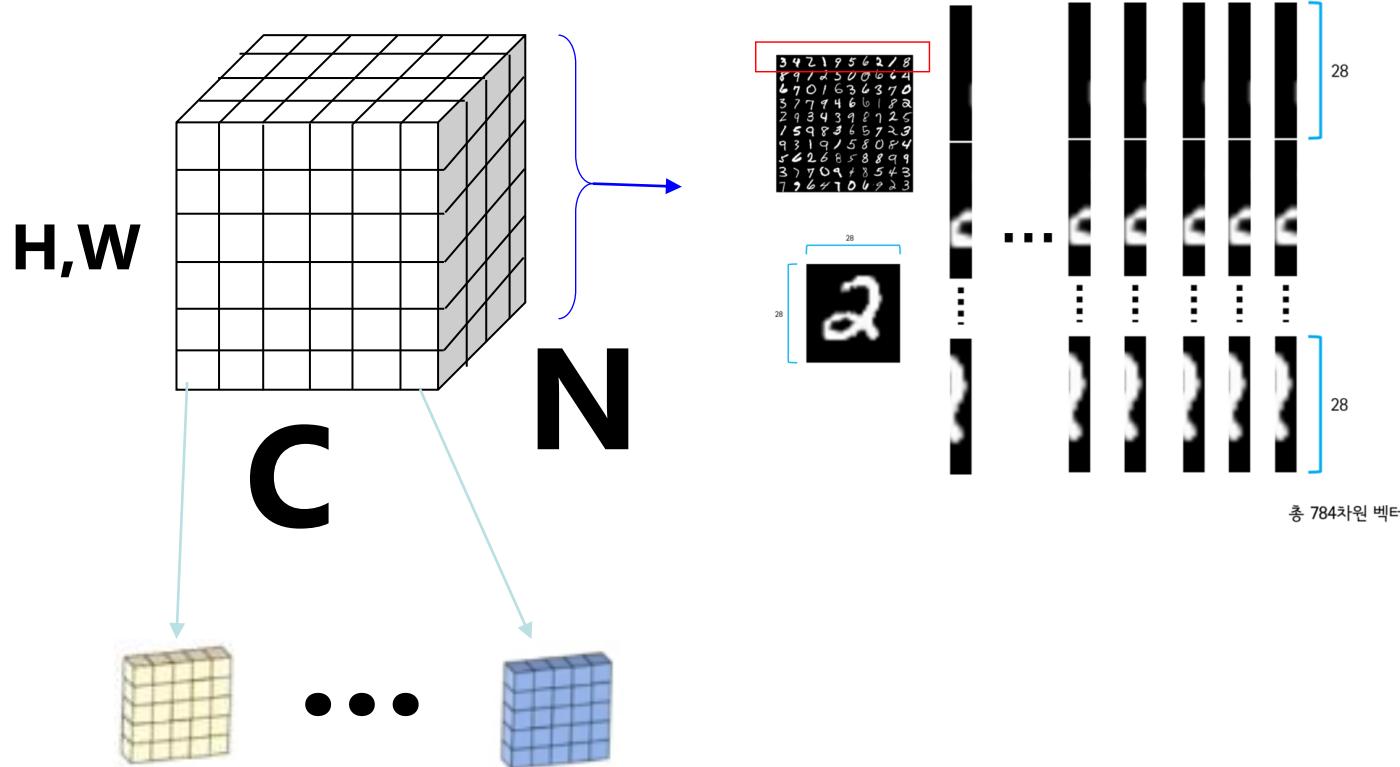
m개의 filter가 있다면 m개의 feature map이 결과로 도출



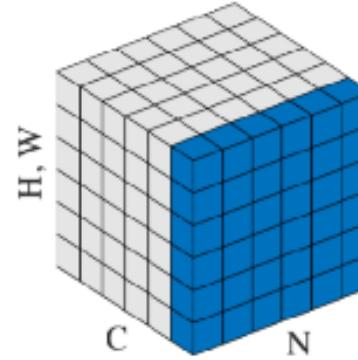
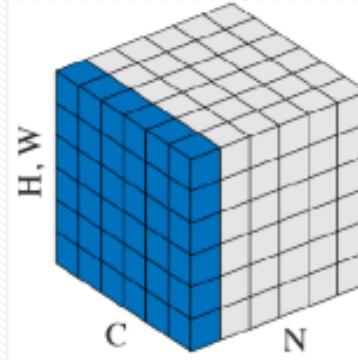
❖ Layer Normalization



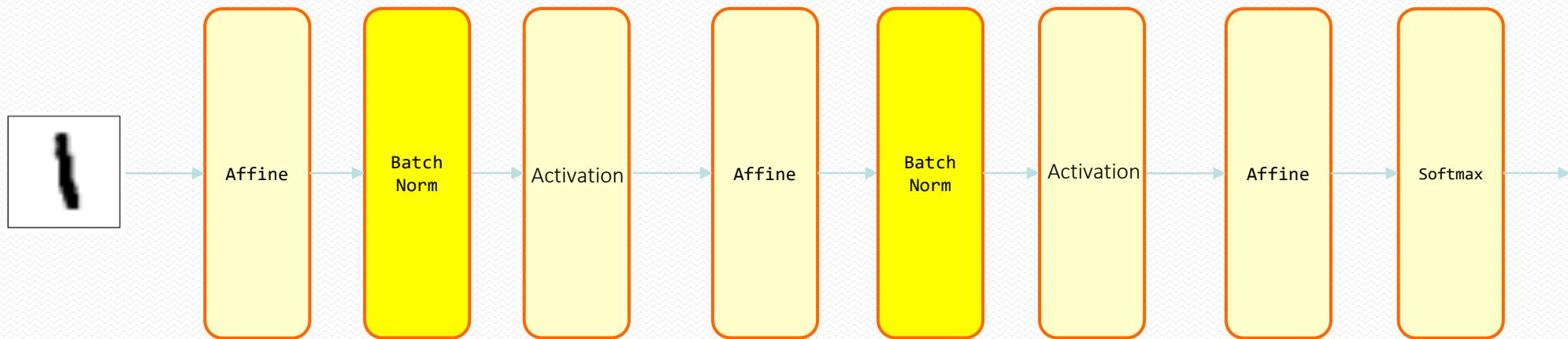
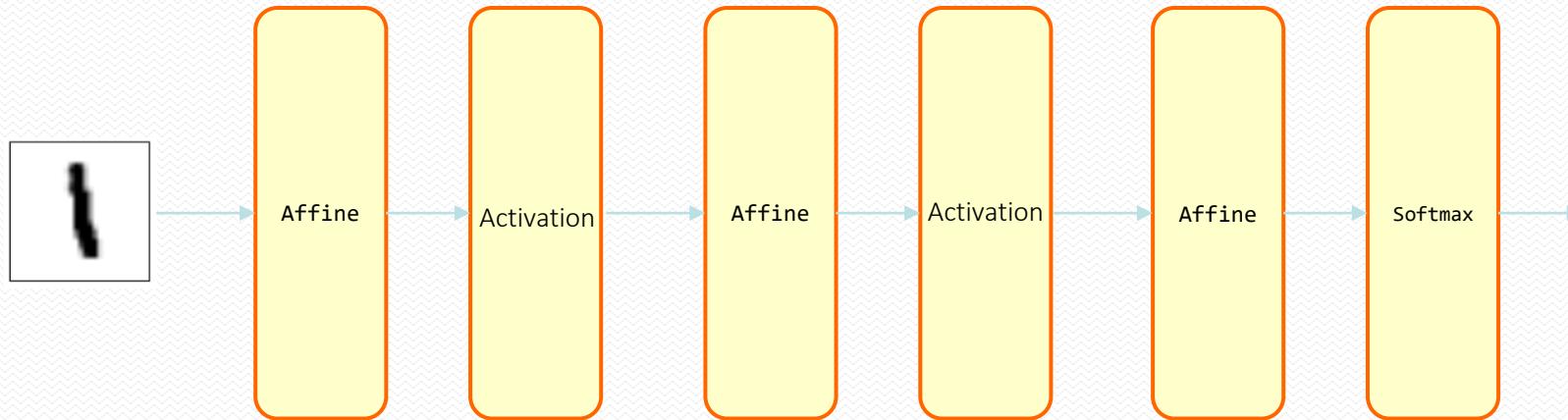
BATCH & LAYER NORMALIZATION



BATCH & LAYER NORMALIZATION

정규화	설명	그림	Shape 비교
Batch Norm	각 feature별로, 한 mini-batch에 대한 summed input의 평균과 분산을 구한다.		<ul style="list-style-type: none">fully-connected networks<ul style="list-style-type: none">$x : N \times D$$\mu, \sigma, \gamma, \beta : 1 \times D$convolutional networks<ul style="list-style-type: none">$x : N \times C \times H \times W$$\mu, \sigma, \gamma, \beta : 1 \times C \times 1 \times 1$
Layer Norm	각 training case별로, 한 레이어의 모든 뉴런에 대한 summed input의 평균과 분산을 구한다.		<ul style="list-style-type: none">fully-connected networks<ul style="list-style-type: none">$x : N \times D$$\mu, \sigma, \gamma, \beta : N \times 1$convolutional networks<ul style="list-style-type: none">$x : N \times C \times H \times W$$\mu, \sigma, \gamma, \beta : N \times 1 \times 1 \times 1$

BATCH & LAYER NORMALIZATION



$$\Sigma(\text{weight} \cdot \text{input}) + \text{bias}$$

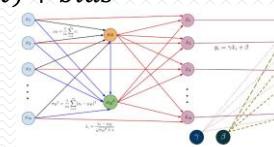
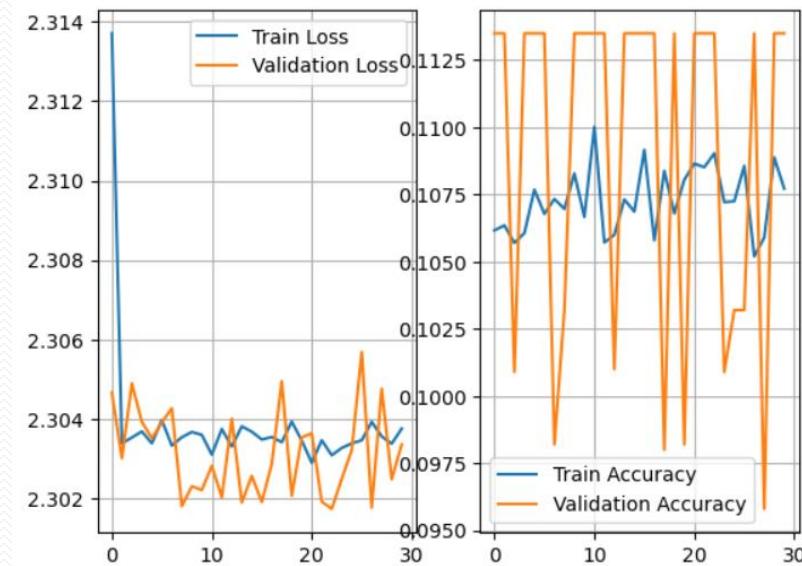
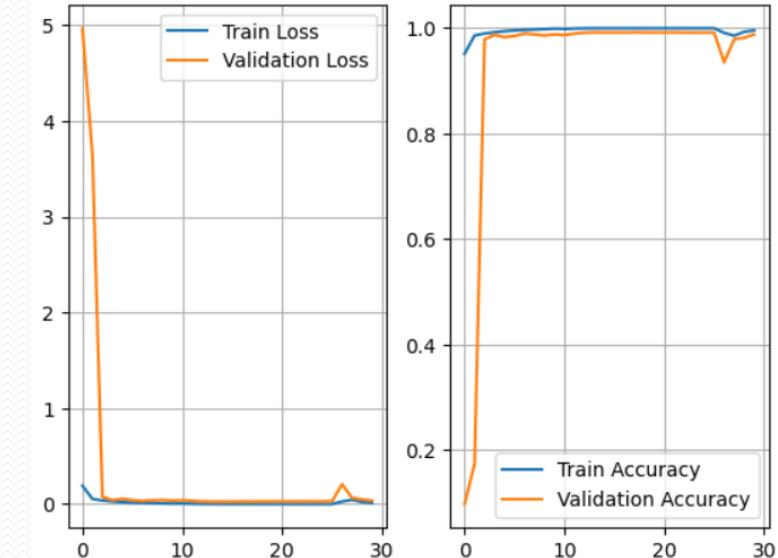


Fig 4. Flow of computation through Batch Normalization layer

배치정규화 알고리즘

```
1 model = Sequential()
2 model.add(Conv2D(32,(2,2),activation="sigmoid",input_shape=(28,28,1)))
3 model.add(BatchNormalization())
4 model.add(Conv2D(64,(2,2),activation="sigmoid"))
5 model.add(BatchNormalization())
6 model.add(Conv2D(128,(2,2),2,activation="sigmoid"))
7 model.add(BatchNormalization())
8 model.add(Conv2D(32,(2,2),activation="sigmoid"))
9 model.add(BatchNormalization())
10 model.add(Conv2D(64,(2,2),activation="sigmoid"))
11 model.add(BatchNormalization())
12 model.add(Conv2D(128,(2,2),2,activation="sigmoid"))
13 model.add(BatchNormalization())
14 model.add(Flatten())
15 model.add(Dense(128,activation="sigmoid"))
16 model.add(Dense(10,activation="softmax"))
```

```
1 model = Sequential()
2 model.add(Conv2D(32,(2,2),activation="sigmoid",input_shape=(28,28,1)))
3 # model.add(BatchNormalization())
4 model.add(Conv2D(64,(2,2),activation="sigmoid"))
5 # model.add(BatchNormalization())
6 model.add(Conv2D(128,(2,2),2,activation="sigmoid"))
7 # model.add(BatchNormalization())
8 model.add(Conv2D(32,(2,2),activation="sigmoid"))
9 # model.add(BatchNormalization())
10 model.add(Conv2D(64,(2,2),activation="sigmoid"))
11 # model.add(BatchNormalization())
12 model.add(Conv2D(128,(2,2),2,activation="sigmoid"))
13 # model.add(BatchNormalization())
14 model.add(Flatten())
15 model.add(Dense(128,activation="sigmoid"))
16 model.add(Dense(10,activation="softmax"))
```



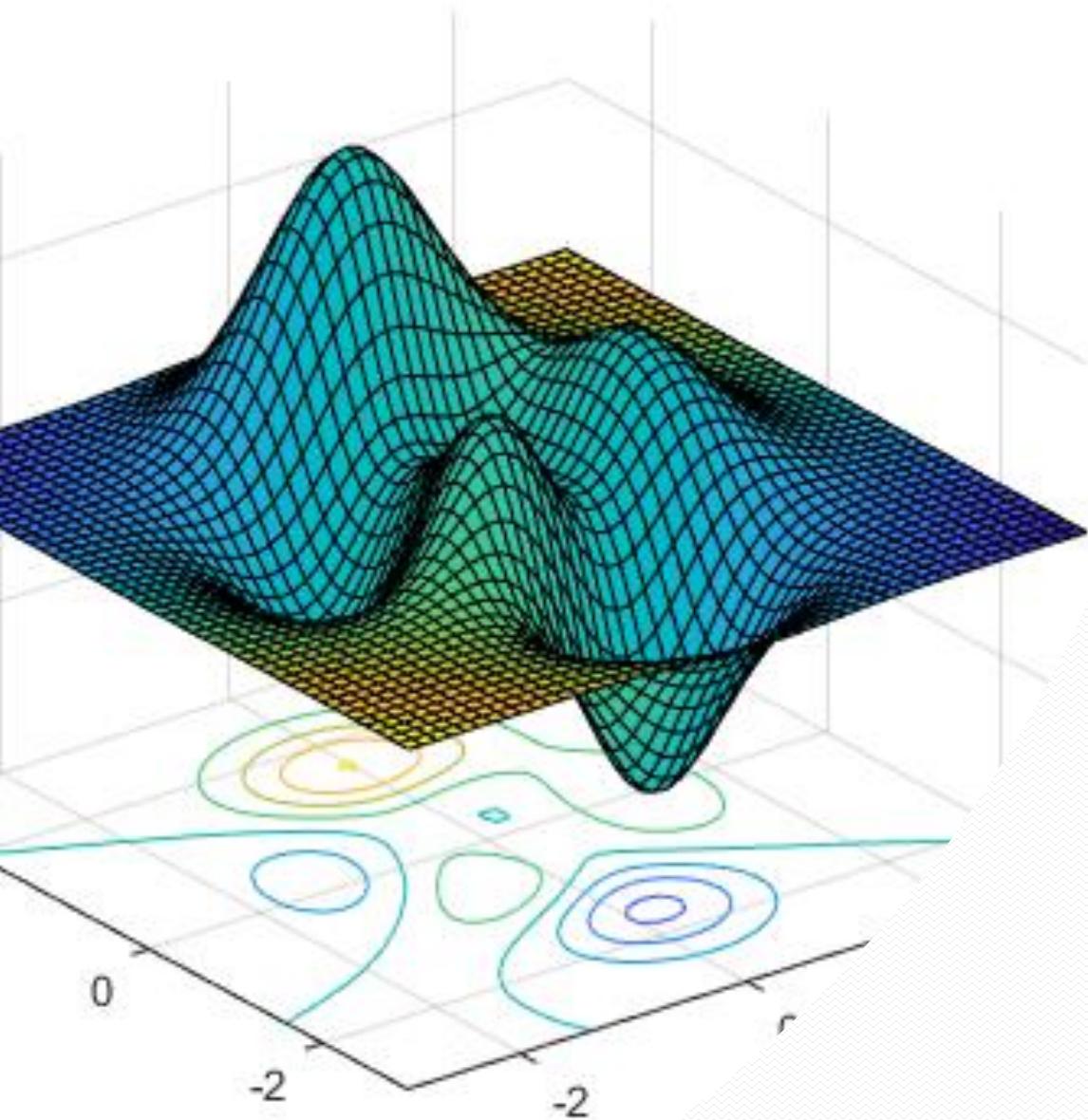
배치정규화 알고리즘

Layer (type)	Output Shape	Param #
<hr/>		
conv2d (Conv2D)	(None, 27, 27, 32)	160
batch_normalization (Batch Normalization)	(None, 27, 27, 32)	128
conv2d_1 (Conv2D)	(None, 26, 26, 64)	8256
batch_normalization_1 (BatchNormalization)	(None, 26, 26, 64)	256
conv2d_2 (Conv2D)	(None, 13, 13, 128)	32896
batch_normalization_2 (BatchNormalization)	(None, 13, 13, 128)	512
conv2d_3 (Conv2D)	(None, 12, 12, 32)	16416
batch_normalization_3 (BatchNormalization)	(None, 12, 12, 32)	128
conv2d_4 (Conv2D)	(None, 11, 11, 64)	8256
batch_normalization_4 (BatchNormalization)	(None, 11, 11, 64)	256
conv2d_5 (Conv2D)	(None, 5, 5, 128)	32896
batch_normalization_5 (BatchNormalization)	(None, 5, 5, 128)	512
flatten (Flatten)	(None, 3200)	0
dense (Dense)	(None, 128)	409728
dense_1 (Dense)	(None, 10)	1290
<hr/>		

Total params: 511690 (1.95 MB)
Trainable params: 510794 (1.95 MB)
Non-trainable params: 896 (3.50 KB)

```
1 model = Sequential()
2 model.add(Conv2D(32,(2,2),activation="sigmoid",input_shape=(28,28,1)))
3 model.add(BatchNormalization())
4 model.add(Conv2D(64,(2,2),activation="sigmoid"))
5 model.add(BatchNormalization())
6 model.add(Conv2D(128,(2,2),2,activation="sigmoid"))
7 model.add(BatchNormalization())
8 model.add(Conv2D(32,(2,2),activation="sigmoid"))
9 model.add(BatchNormalization())
10 model.add(Conv2D(64,(2,2),activation="sigmoid"))
11 model.add(BatchNormalization())
12 model.add(Conv2D(128,(2,2),2,activation="sigmoid"))
13 model.add(BatchNormalization())
14 model.add(Flatten())
15 model.add(Dense(128,activation="sigmoid"))
16 model.add(Dense(10,activation="softmax"))
```

None

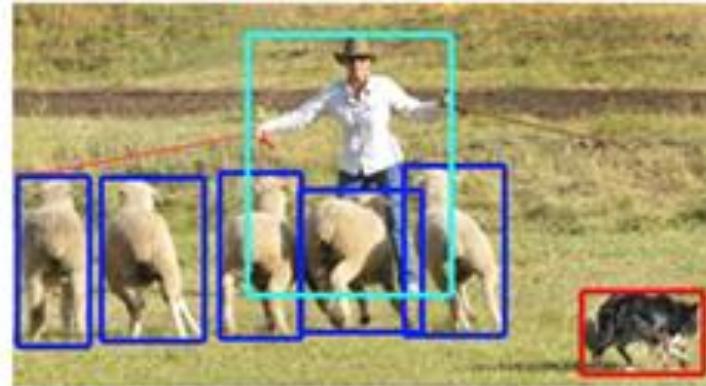


USE CASES

CNN USE CASES



CLASSIFICATION



DETECTION



SEGMENTATION



NEURAL STYLE TRANSFER



INPAINTING



COLORIZATION

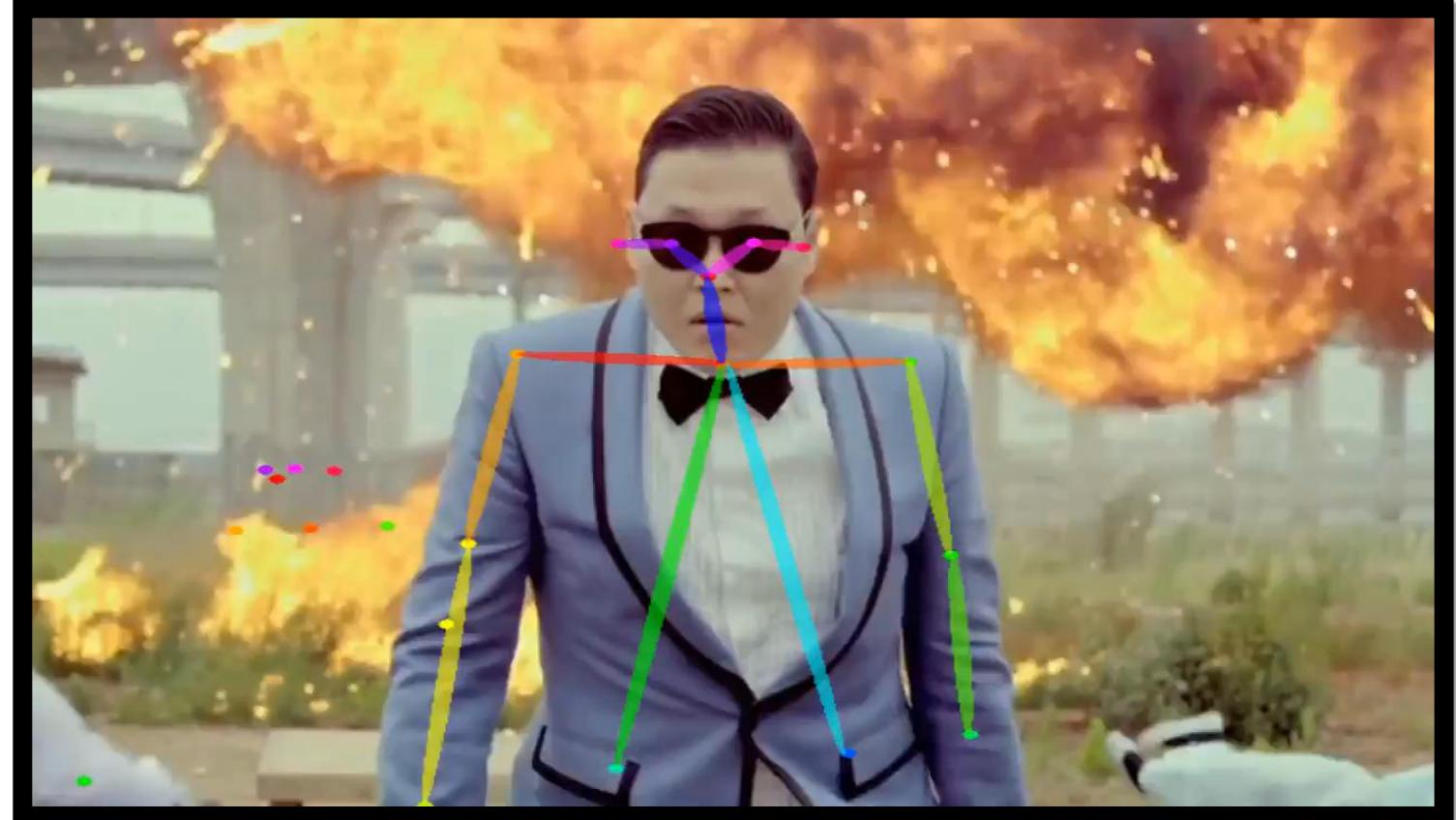
CNN USE CASES / CONT.



SUPER RESOLUTION



MONO DEPTH



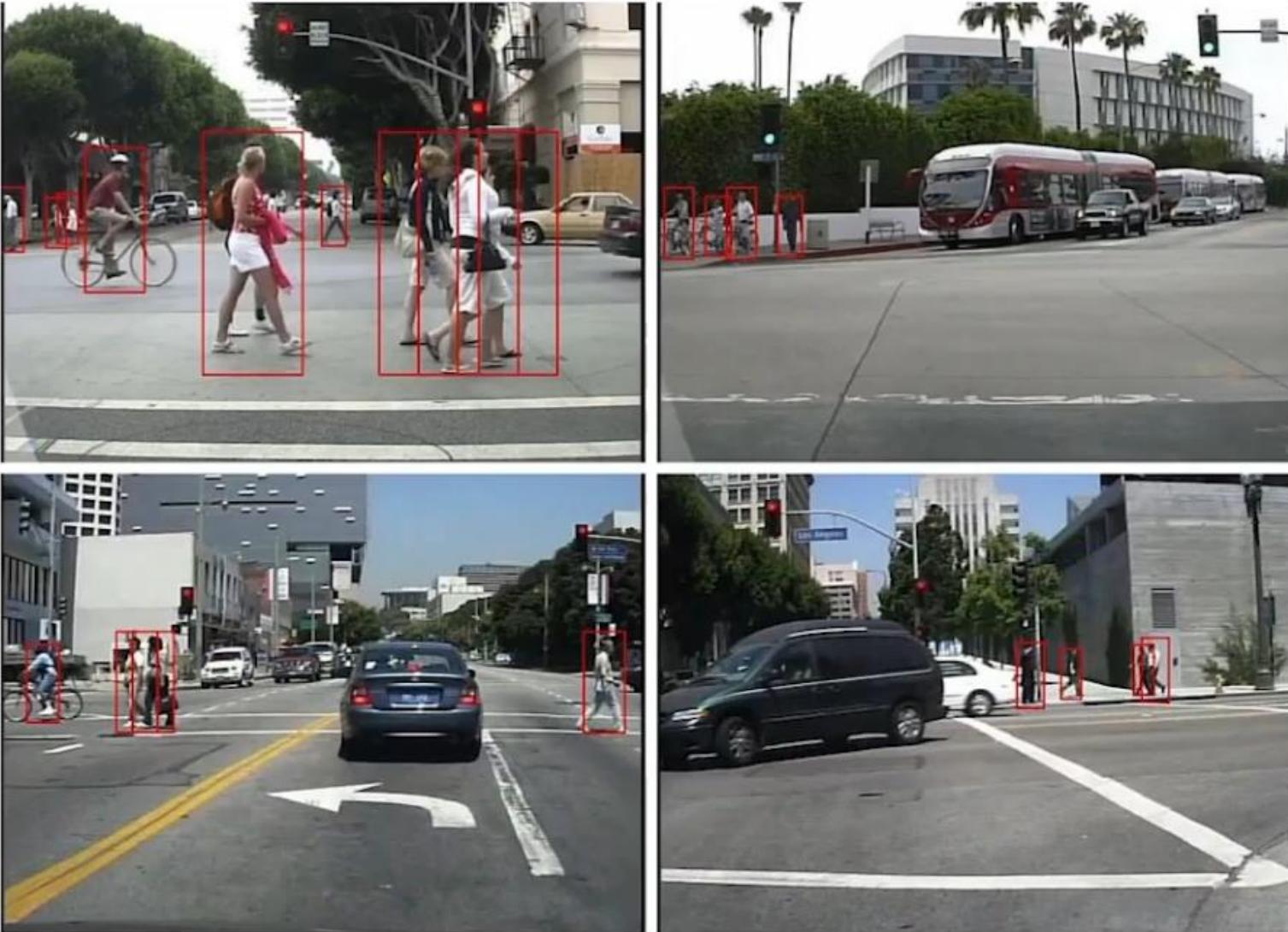
HUMAN POSE ESTIMATION

CNN USE CASES / CONT.

A young boy is playing basketball. 	Two dogs play in the grass. 	A dog swims in the water. 	A little girl in a pink shirt is swinging. 
A group of people walking down a street. 	A group of women dressed in formal attire. 	Two children play in the water. 	A dog jumps over a hurdle. 

IMAGE CAPTIONING

CNN USE CASES / CONT.



CNN USE CASES / CONT.

누가 지각/결석을 했지?

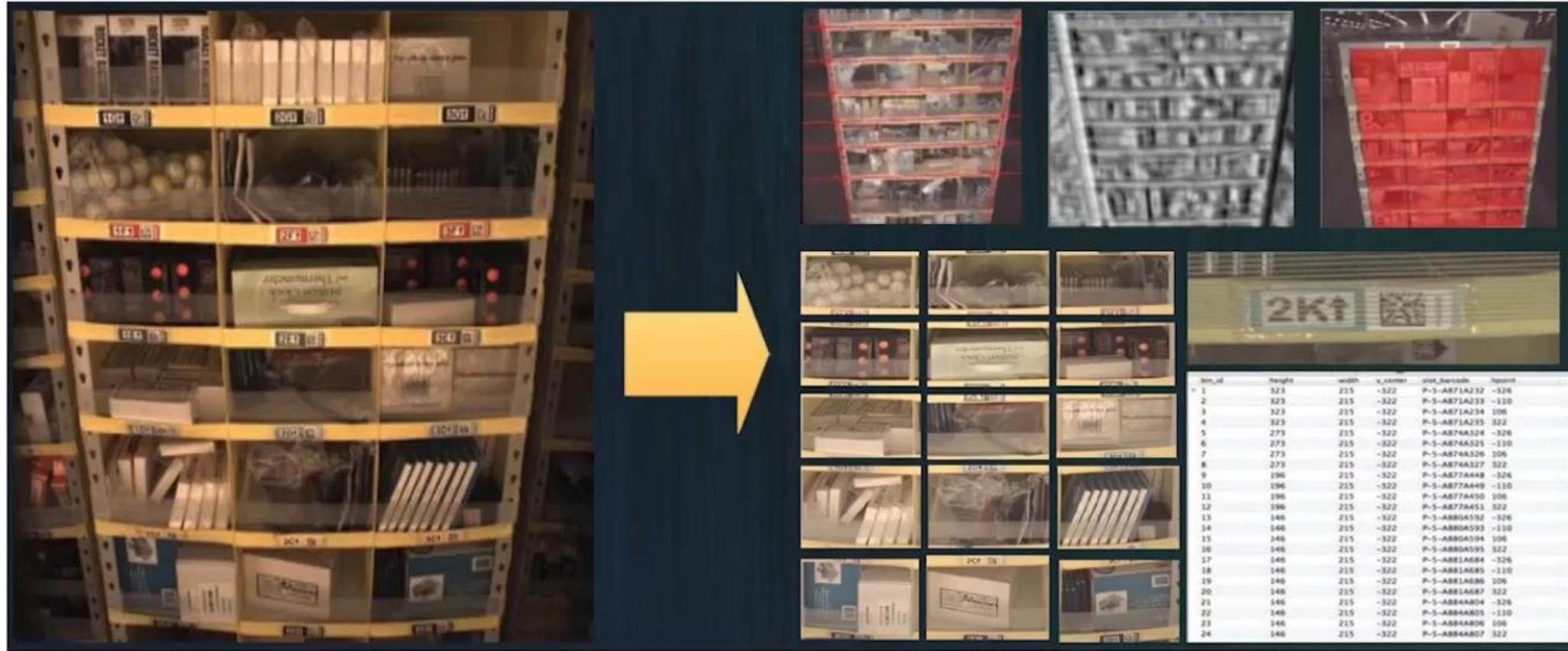


수업이 얼마나 따른한가?



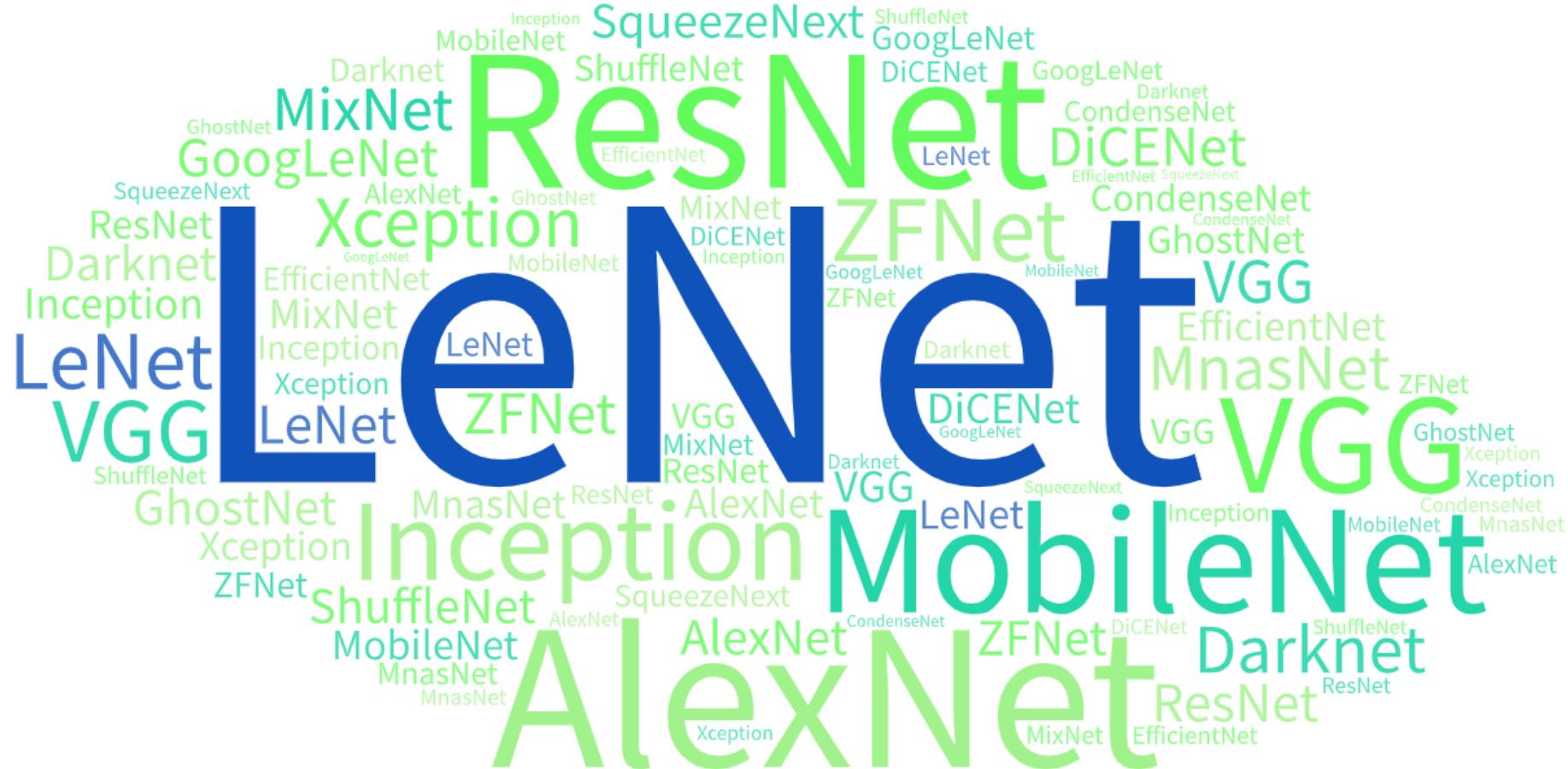
CNN USE CASES / CONT.

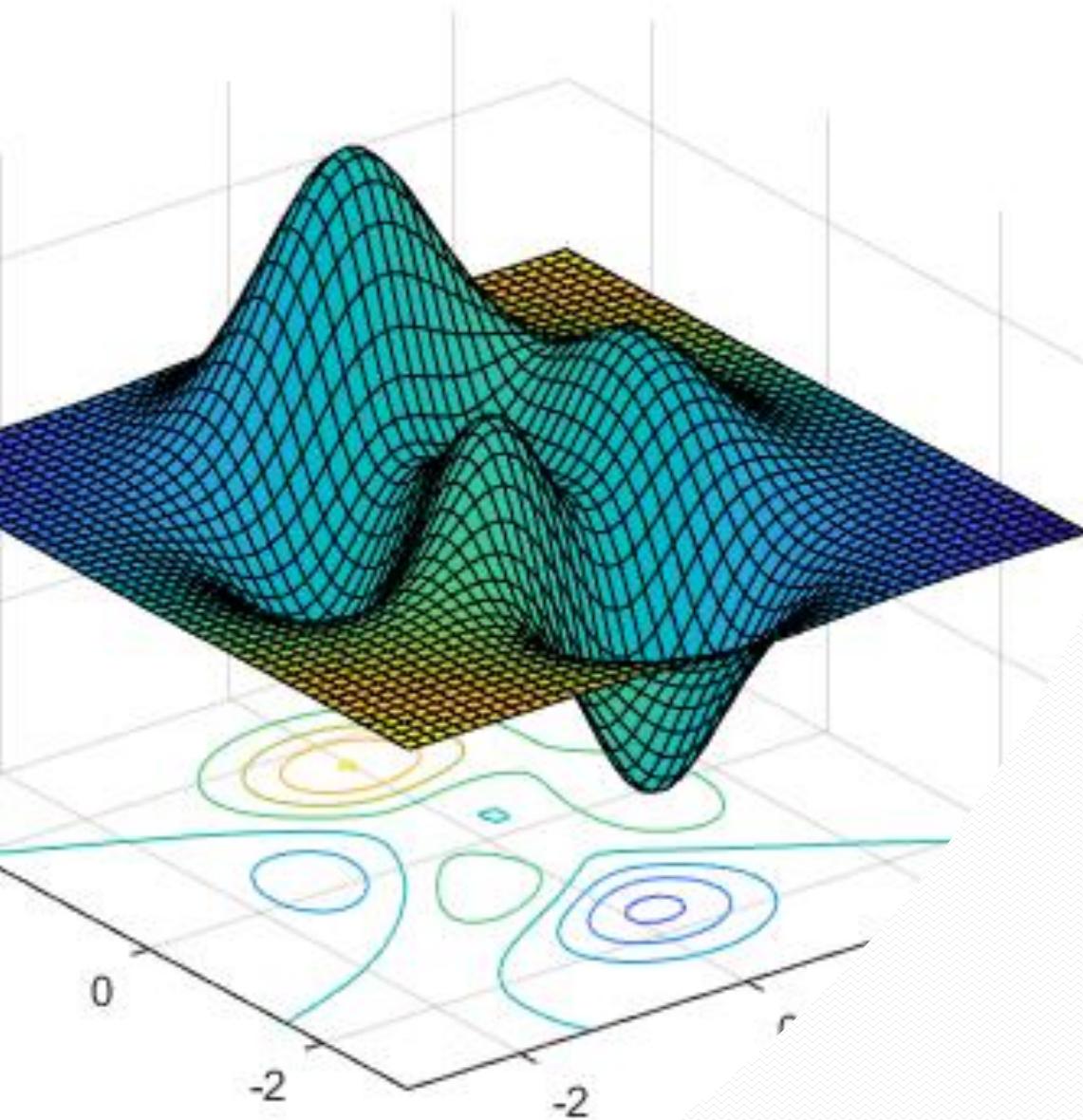
비즈니스 활용 사례 : 재고 파악



<https://www.slideshare.net/awskorea/amazon-deeplearningcasesandmlonaws>

NETWORKS WITH CNN





Classification

IMAGENET LARGE SCALE VISUAL RECOGNITION CHALLENGE ILSVRC

✓ Layer의 수가 많아질수록 이미지 분류 성능이 증가하는 추세를 보임

2010 - NEC-UIUC Lin et al.

2011 - XRCE Florent Perronnin, Jorge Sanchez

2012 - **AlexNet**

2013 - ZFNet

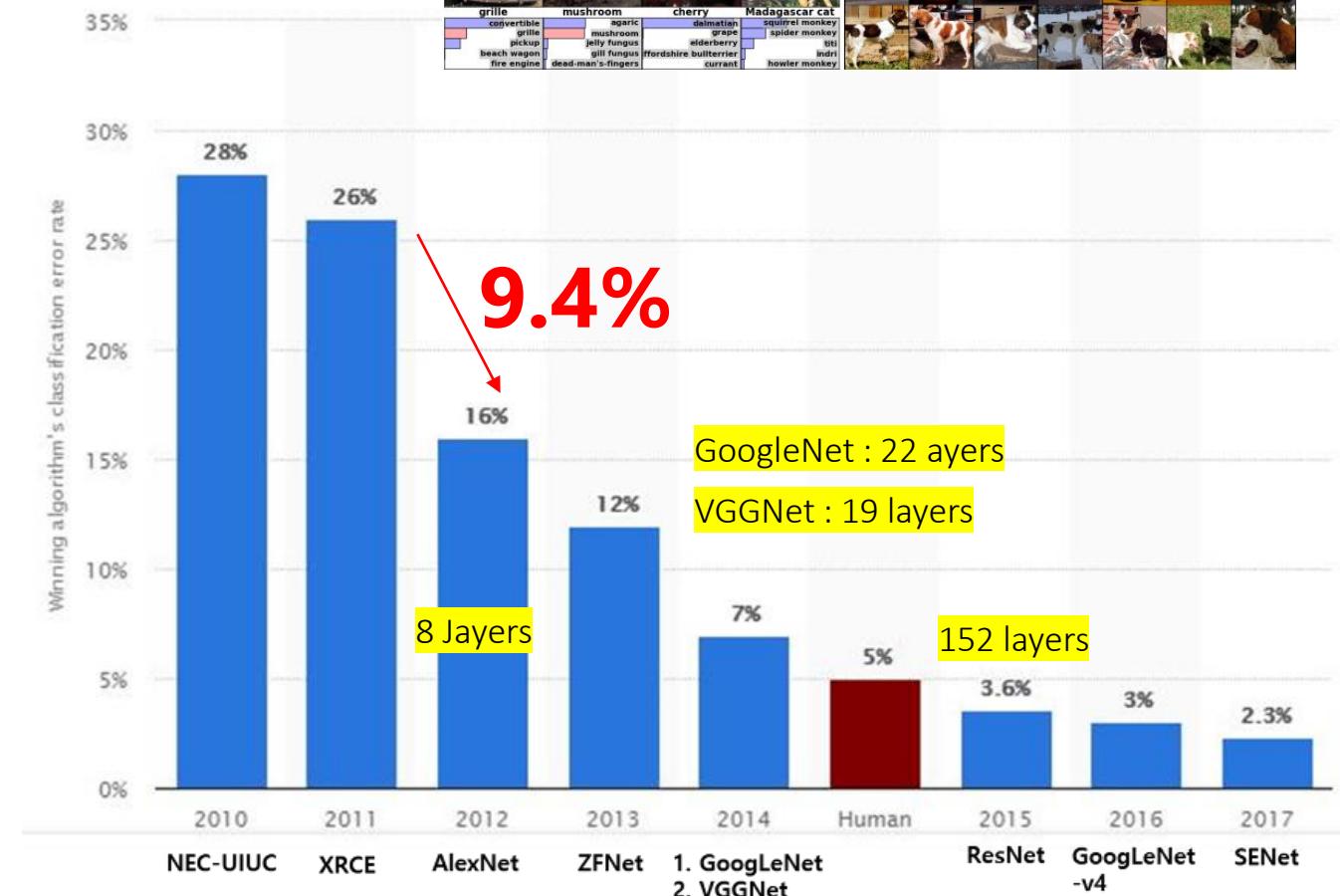
2014 – **GoogLeNet**

VGGNet Second Winner

2015 - **ResNet**

2016 - GoogLeNet-v4

2017 - SENet



LeNET-5 ARCHITECTURE

Yann LeCun은 CNN 개념을 최초로 개발했으며 LeNet은 Yann LeCun이 개발한 CNN architecture 이다.

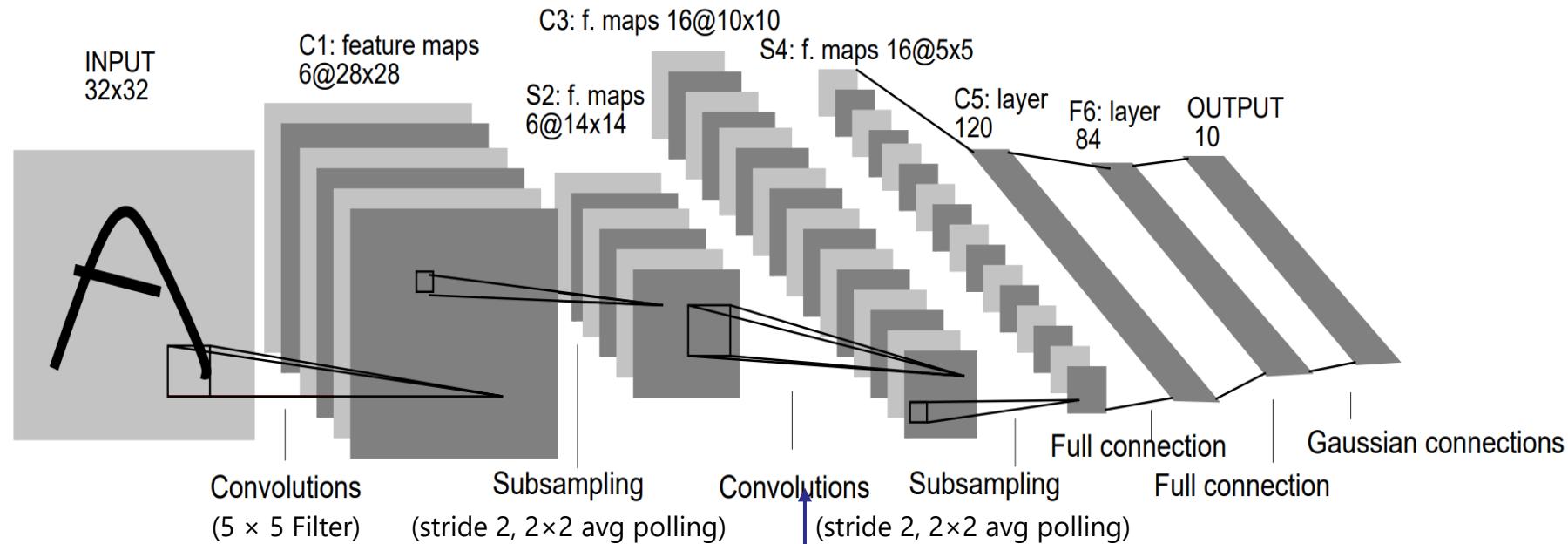


Fig. 2. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

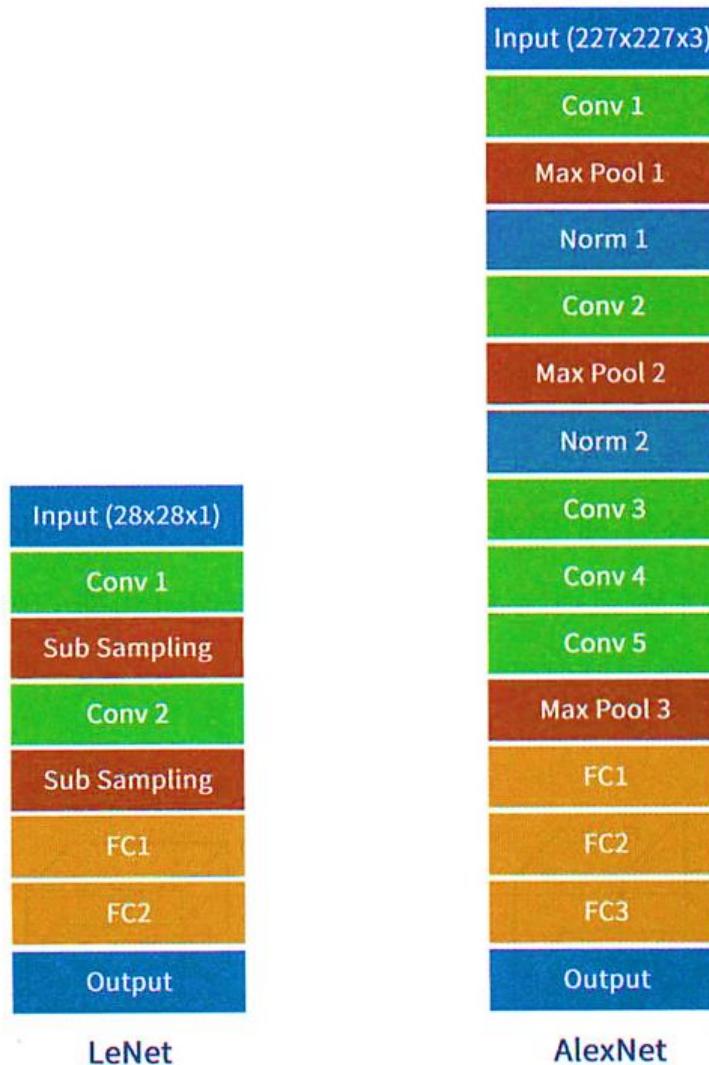
Input	Channel															
	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	X			X	X	X			X	X	X	X		X	X	X
1	X	X			X	X	X			X	X	X	X		X	
2	X	X	X			X	X	X			X		X	X	X	
3		X	X	X		X	X	X	X		X		X	X	X	
4		X	X	X			X	X	X	X		X	X	X		X
5			X	X	X			X	X	X	X		X	X	X	X

<-논문의 작성자가
임의로 선택한 값

EACH COLUMN INDICATES WHICH FEATURE MAP IN S2 ARE COMBINED
BY THE UNITS IN A PARTICULAR FEATURE MAP OF C3.

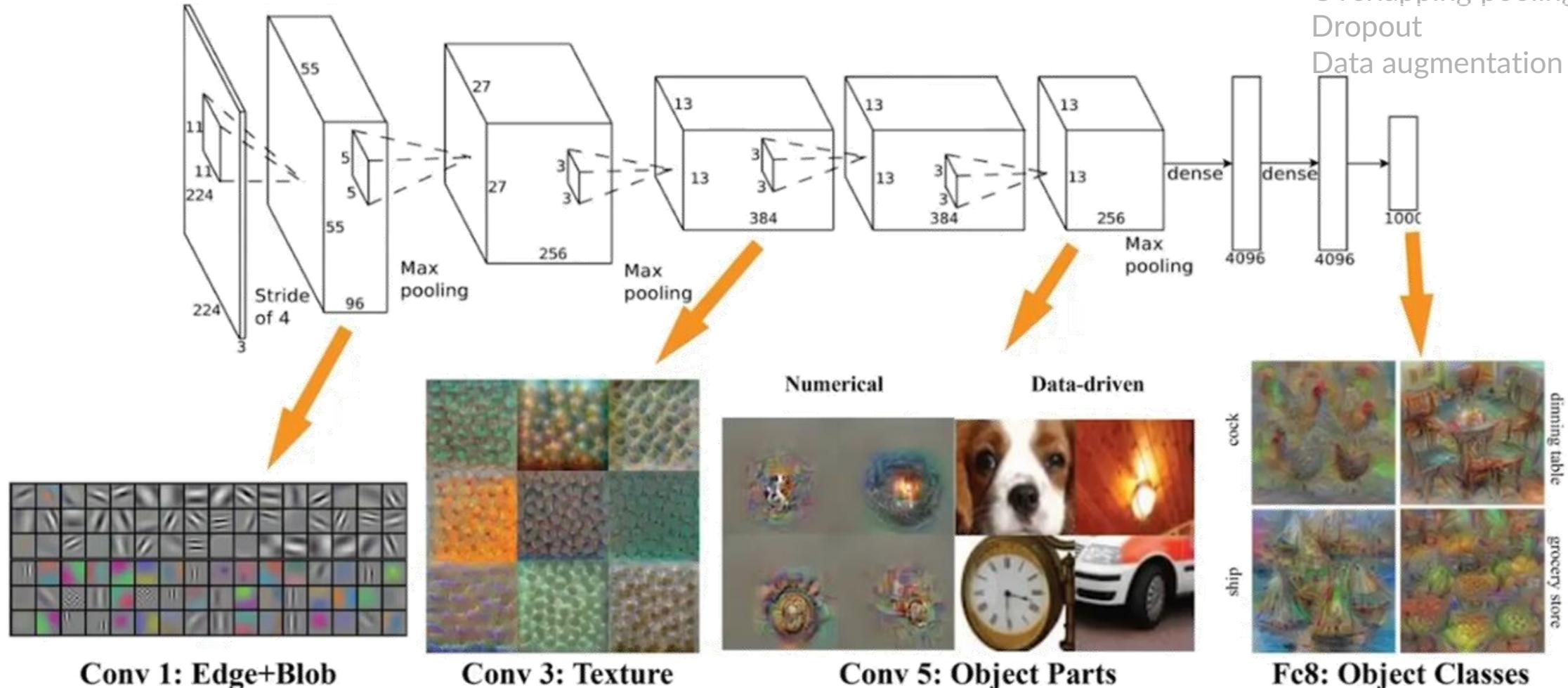
AlexNet ARCHITECTURE

ReLU + Local response normalization
Multiple GPU
Overlapping pooling
Dropout
Data augmentation



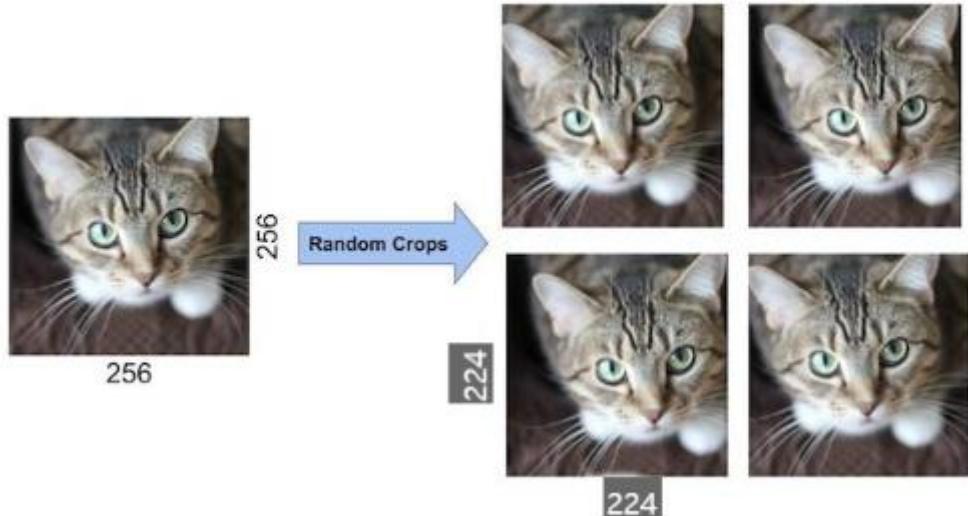
AlexNet ARCHITECTURE

AlexNet은 분류문제의 오차율을 25%에서 16%로 급격한 성능 향상을 보인 모델이다.

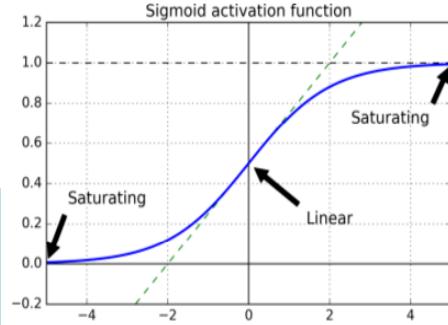


AlexNet ARCHITECTURE

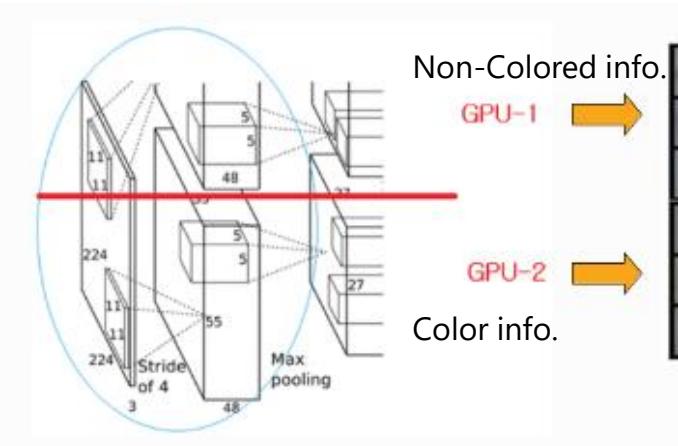
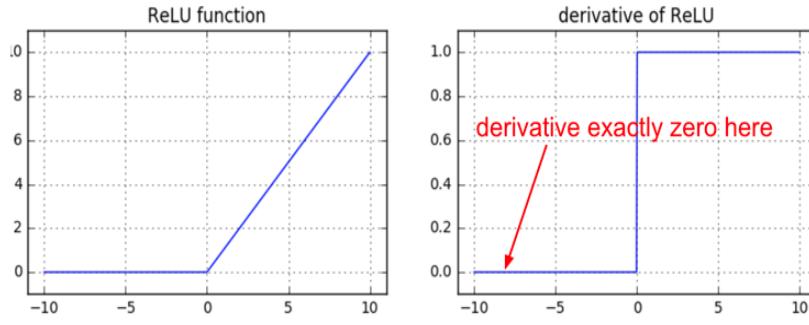
ReLU + Local response normalization
Multiple GPU
Overlapping pooling
Dropout
Data augmentation



Training speed slow down in the saturating point due to a small gradient.

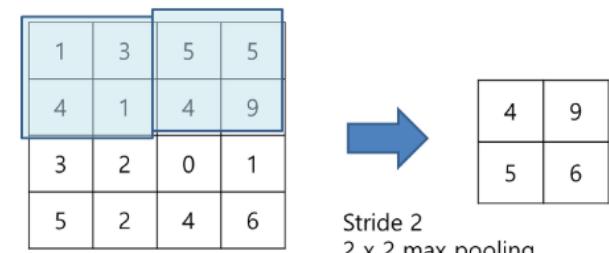


Specific values can increase infinitely.
Therefore, neighboring value can be ignored.
(this issue was solved by normalization (LRN)).



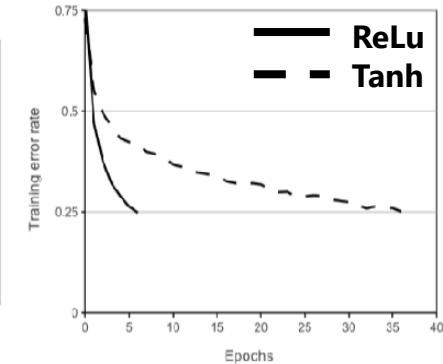
Filters in a first convolution layer

Non-overlapping pooling

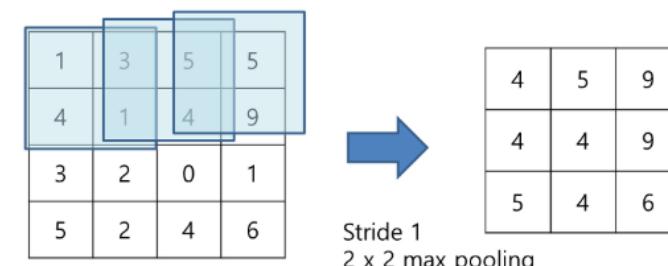


Stride 2
2 x 2 max pooling

ReLU shows faster training speed compared to Tanh.



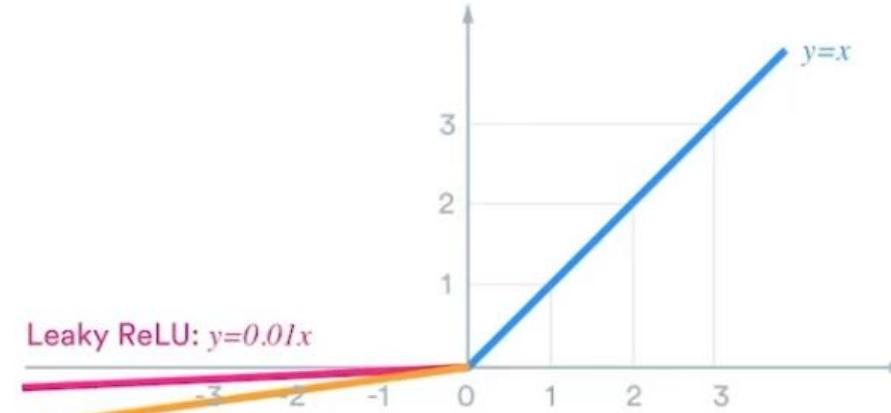
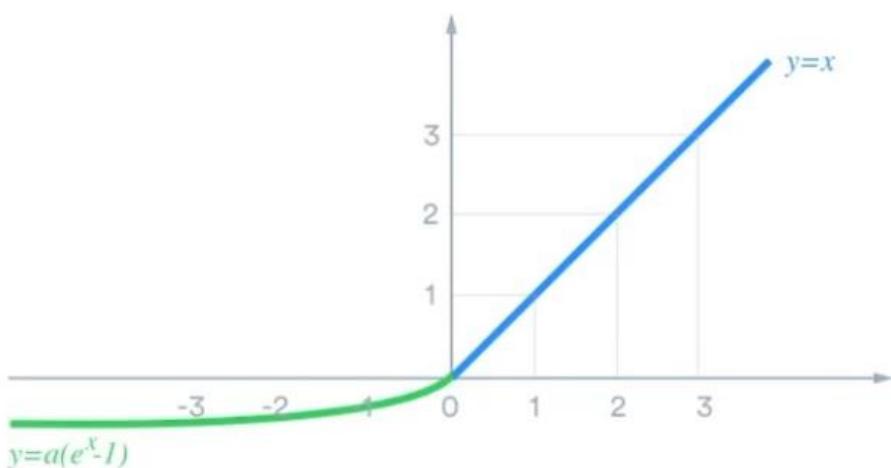
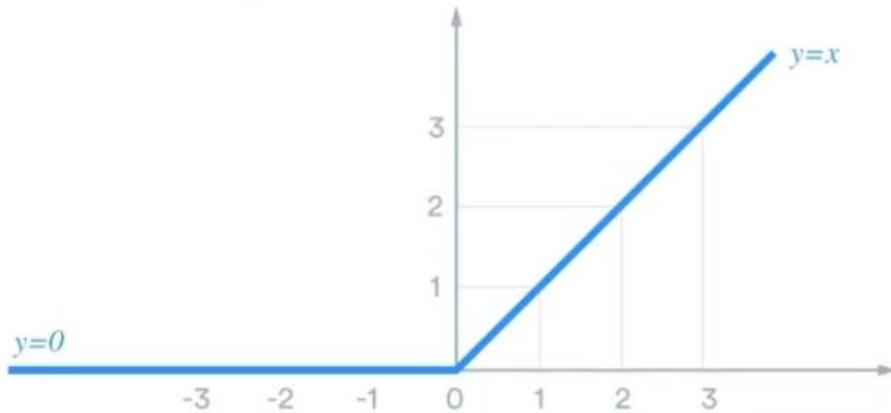
Overlapping pooling



Stride 1
2 x 2 max pooling

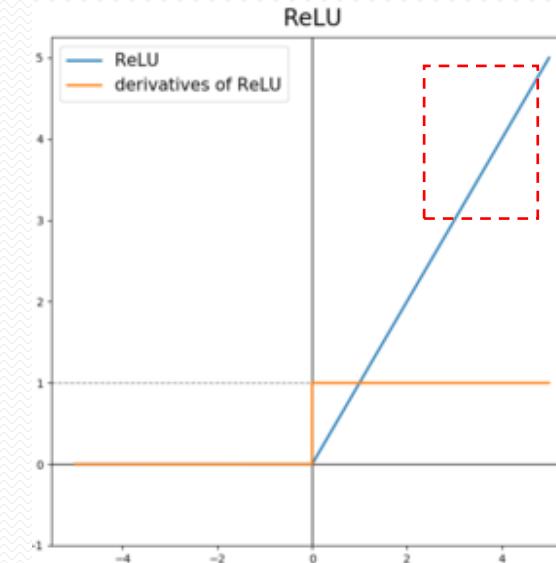
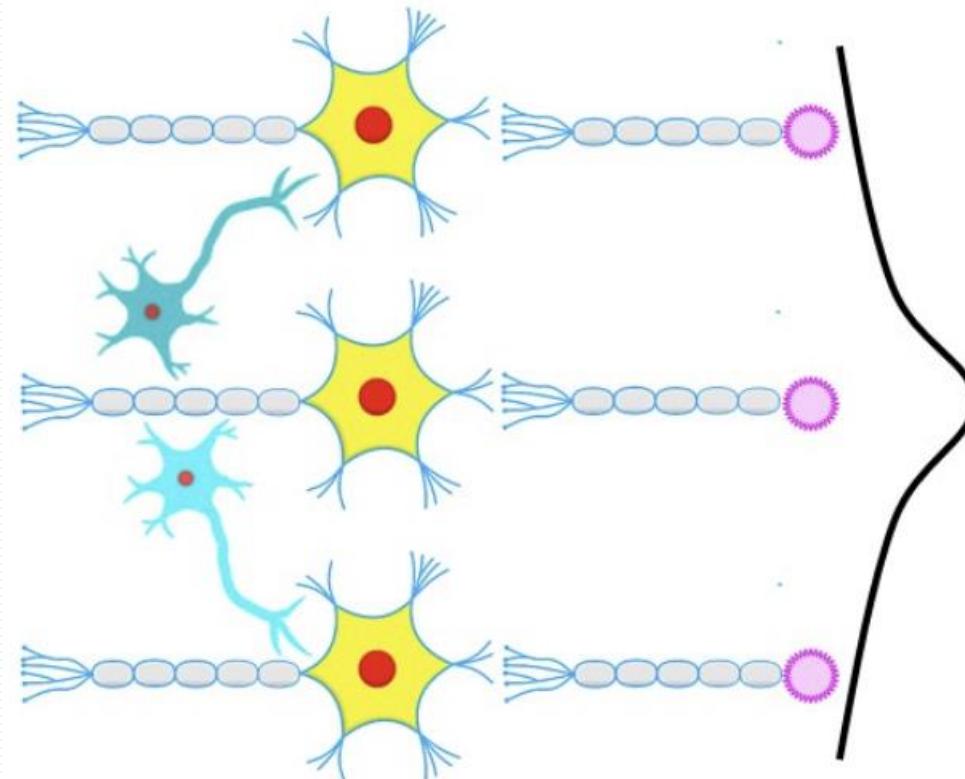
AlexNet ARCHITECTURE: ReLu

ReLU(Rectified Linear Unit) 사용 → 최근에는 변형 사용



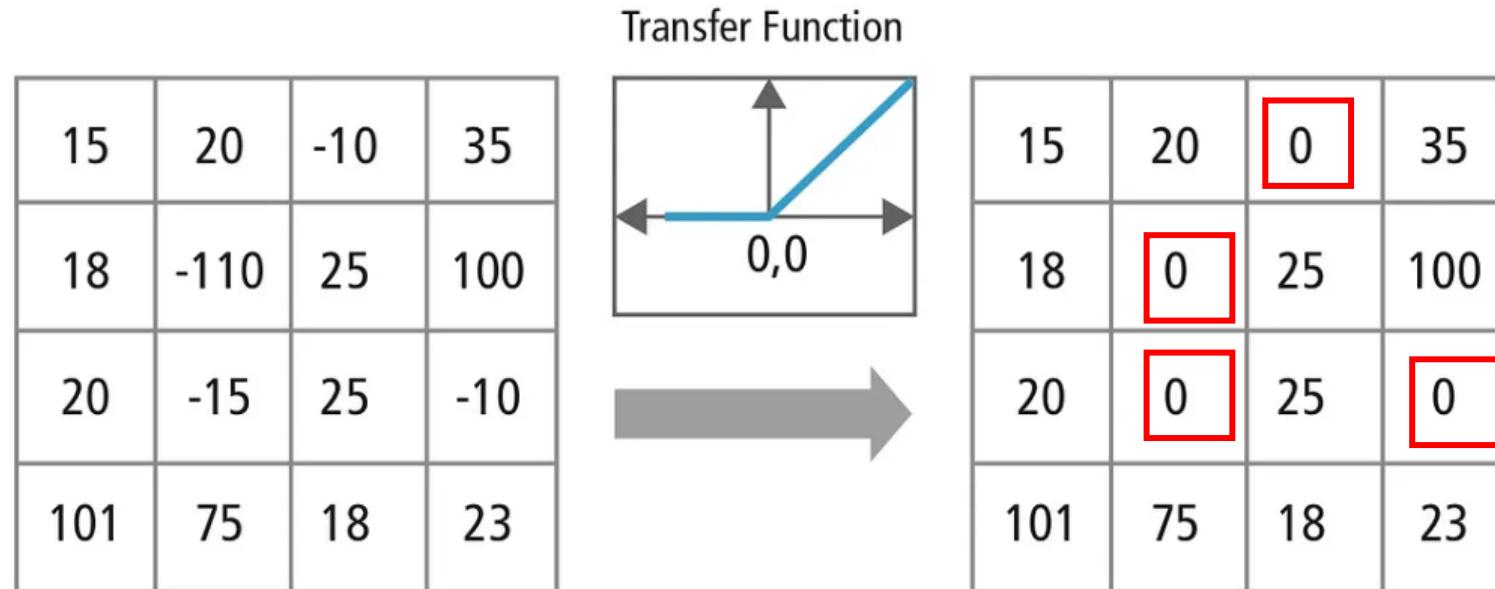
AlexNet ARCHITECTURE: ReLu vs 측면억제

입출력단을 갖는 2개의 신경세포에서 각각의 신경세포의 입력 또는 출력이 서로 다른 신경세포의 입력 또는 출력에 의해 억제되는 현상.



Activation

- ✓ Convolution을 통해 학습된 값들의 비선형 변환을 수행
- ✓ 대부분 Rectified Linear Unit (ReLU)을 사용



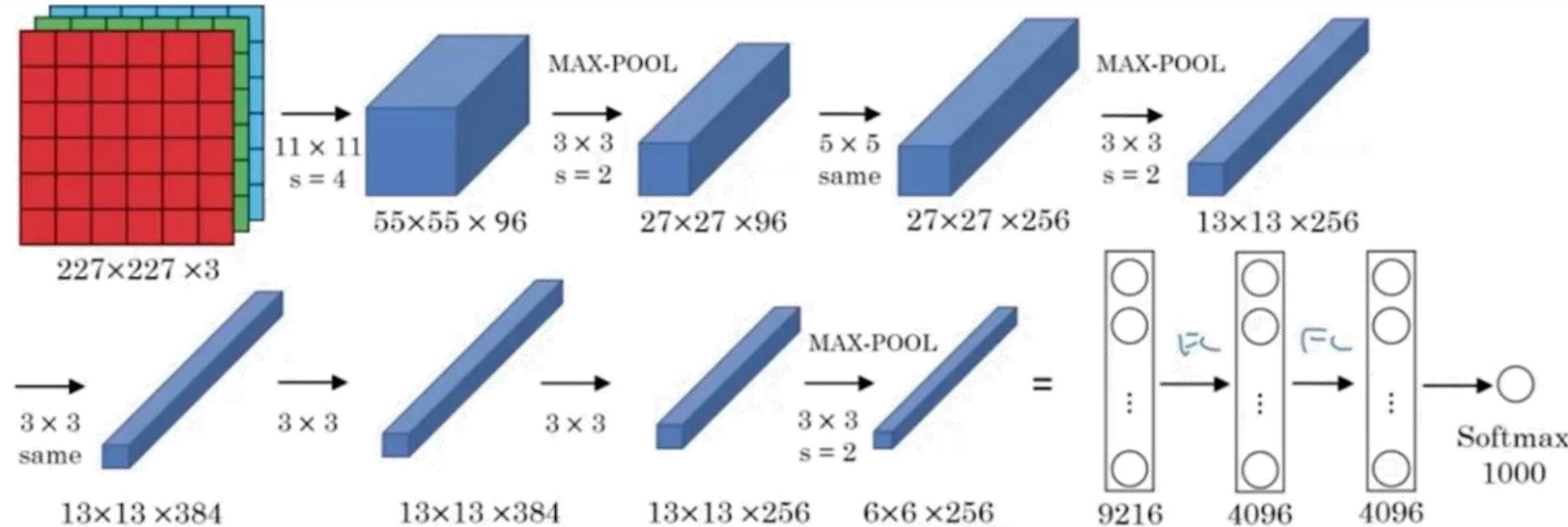
$$\text{Output} = \text{Max}(zero, \text{Input})$$

<https://medium.com/data-science-group-iitr/building-a-convolutional-neural-network-in-python-with-tensorflow-d251c3ca8117>

AlexNet ARCHITECTURE

- ✓ 224*224 크기의 이미지를 Input으로 사용
- ✓ 초기 단계에서는 큰 필터 사이즈와 Stride를 사용
- ✓ 상위 Layer로 갈수록 작은 필터 사이즈와 Stride를 사용
- ✓ 2개의 Fully connected layer 존재
- ✓ 파라미터 총 수 : 60 million

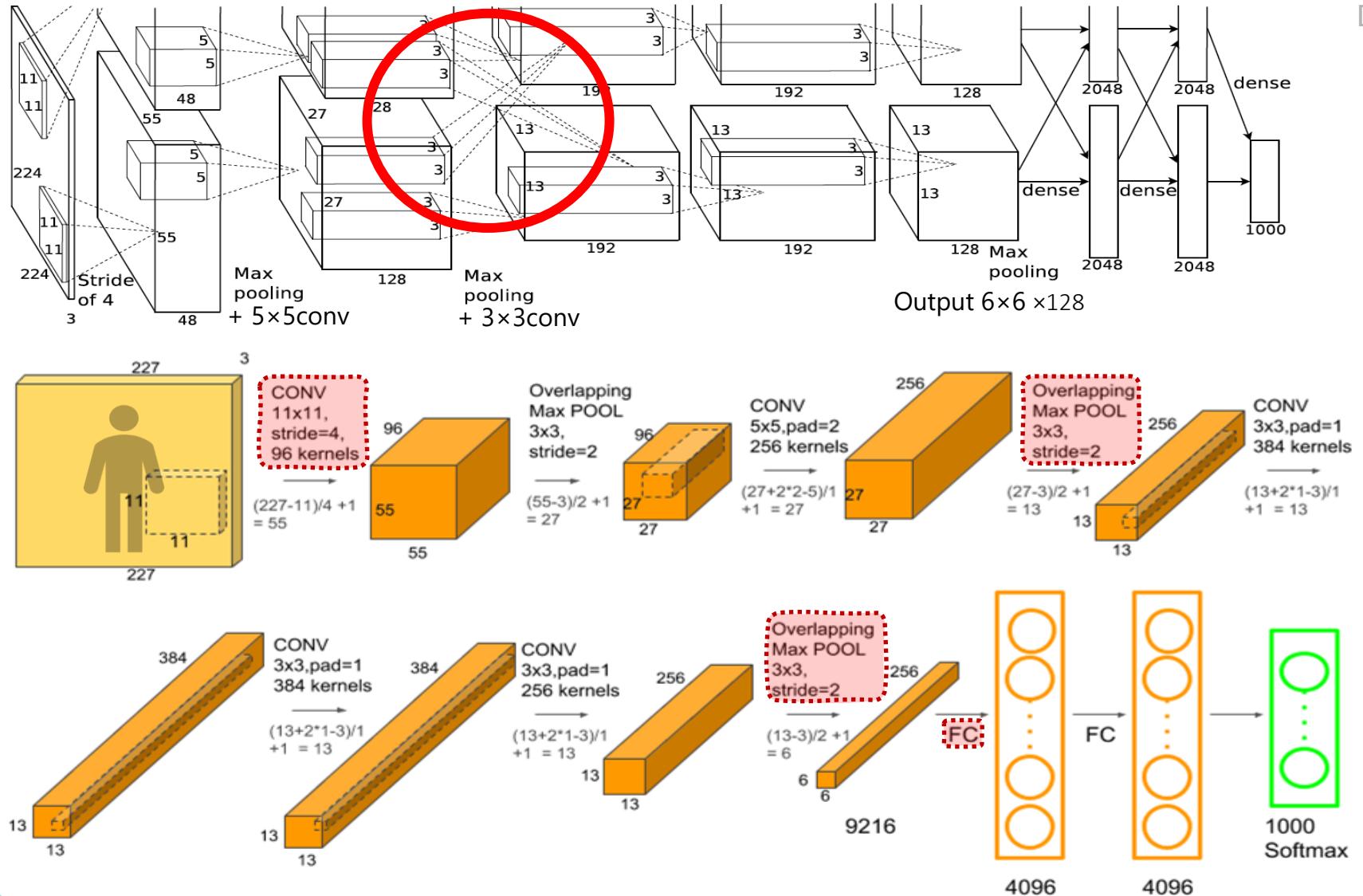
ReLU + Local response normalization
Multiple GPU
Overlapping pooling
Dropout
Data augmentation



AlexNet ARCHITECTURE

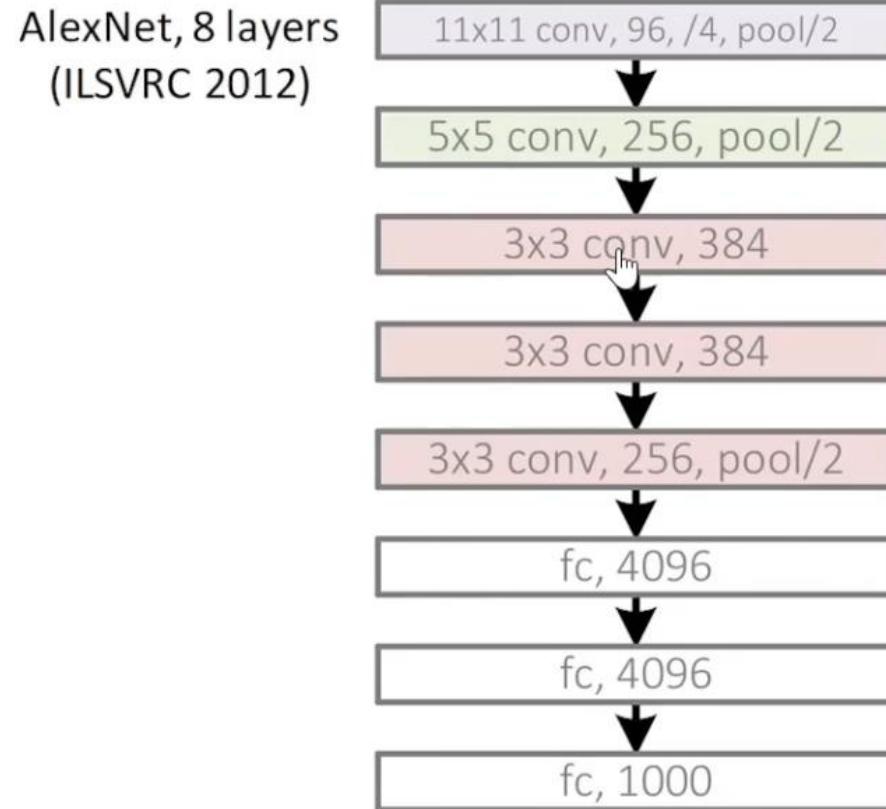
AlexNet은 분류문제의 오차율을 25%에서 16%로 급격한 성능 향상을 보인 모델이다.

ReLU + Local response normalization
Multiple GPU
Overlapping pooling
Dropout
Data augmentation



AlexNet ARCHITECTURE

- ✓ 최초의 ReLU 함수 제시
- ✓ 최초의 GPU 사용



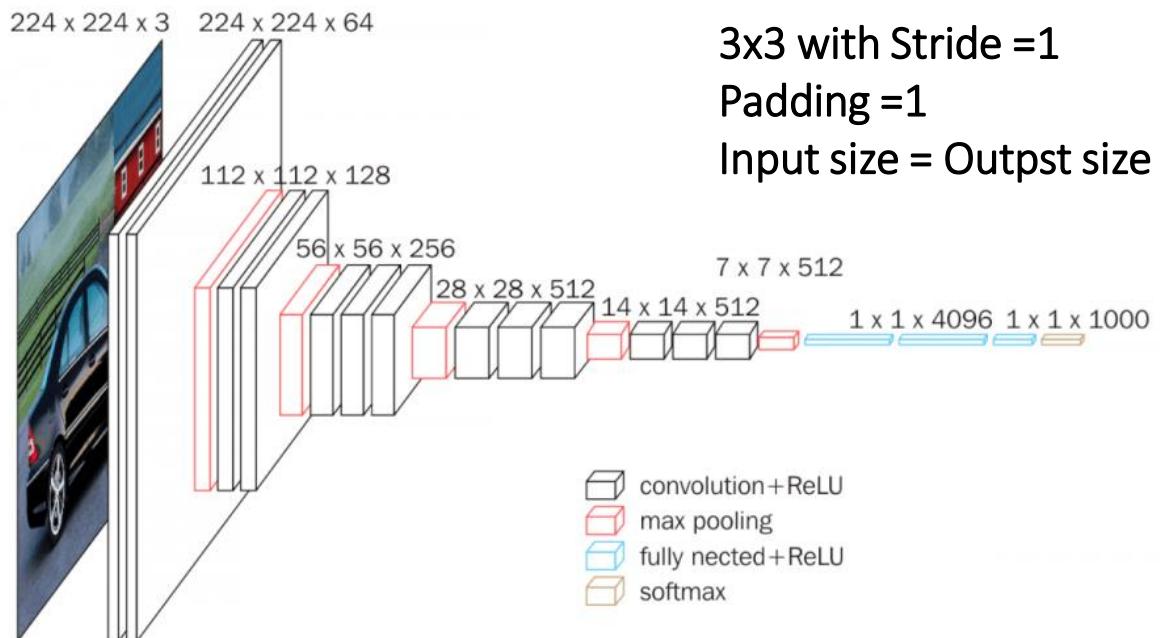
ReLU + Local response normalization
Multiple GPU
Overlapping pooling
`opout`
Data augmentation

VGG, 2014

Factorizing Conv.
Depth increased

- ✓ The key is to see how deeper the depth of the network affects performance
- ✓ AlexNet에 비해 단순하지만 깊은 구조로 단순하면서 성능이 좋은 네트워크
- ✓ 3 by 3 convolution with stride 1을 기본 연산으로 하며 중간 중간에 2 by 2 max polling을 수행
- ✓ 파라미터 총 수 : 138 million

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 conv3-128
maxpool					
conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256 conv1-256	conv3-256 conv3-256 conv3-256	conv3-256 conv3-256 conv3-256 conv3-256
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 conv1-512	conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512 conv3-512
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 conv1-512	conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512 conv3-512
maxpool					
FC-4096					
FC-4096					
FC-1000					
soft-max					

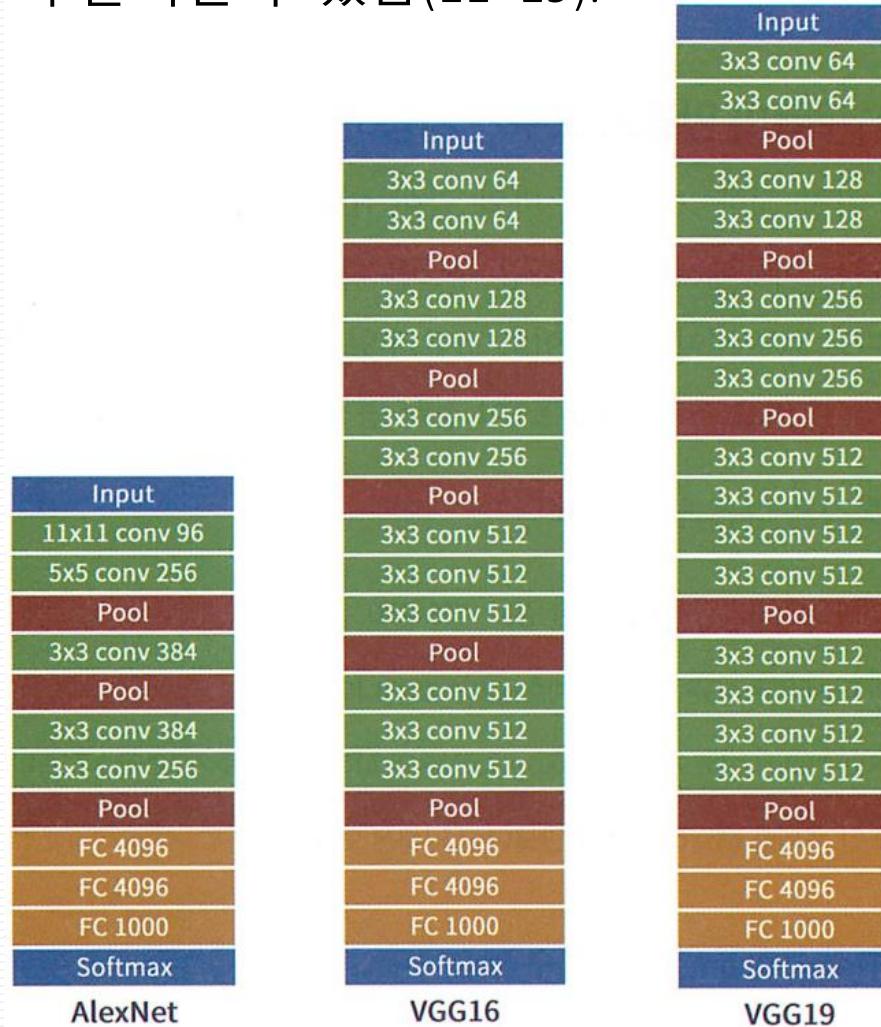


- conv[필터 사이즈]-[필터 개수]의 방식으로 표기.
- 이 중 16레이어와 19레이어로 구성된 VGG16과 VGG19가 오류율이 낮아서 많이 사용.

VGG, 2014

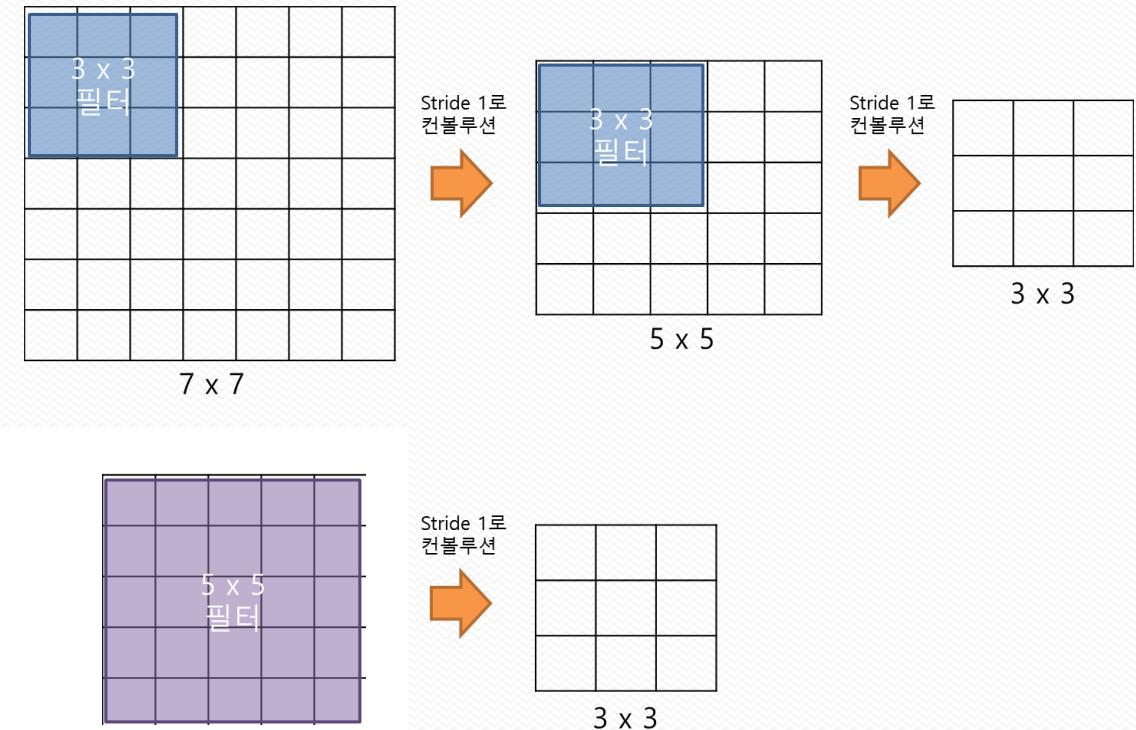
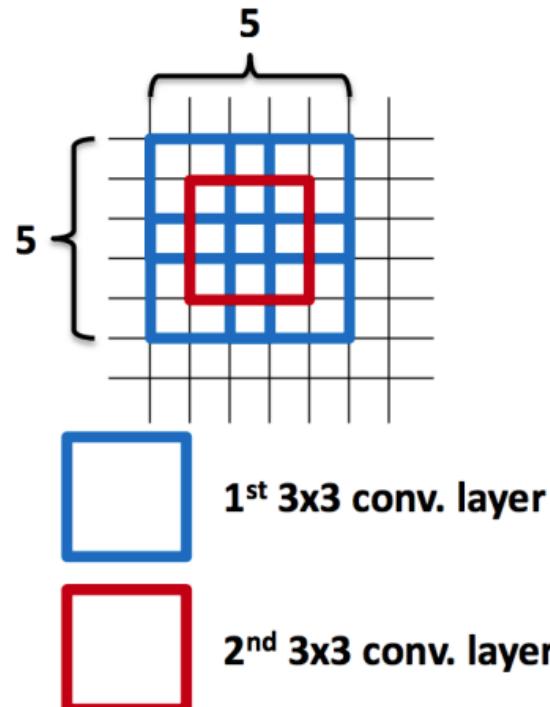
Factorizing Conv.
Depth increased

- ✓ VGGNet에서는 AlexNet에서 사용한 큰 사이즈의 필터가 없고 3×3 필터와 1×1 필터만 사용.
- ✓ AlexNet보다 레이어의 수가 늘어날 수 있음(11~19).
- ✓ 1×1 콘볼루션을 사용.



VGGNet Architecture

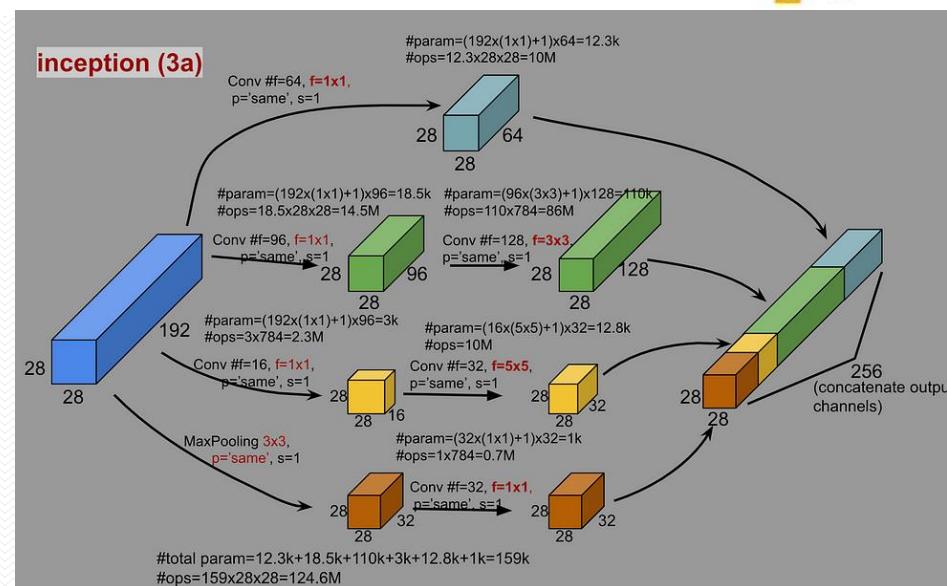
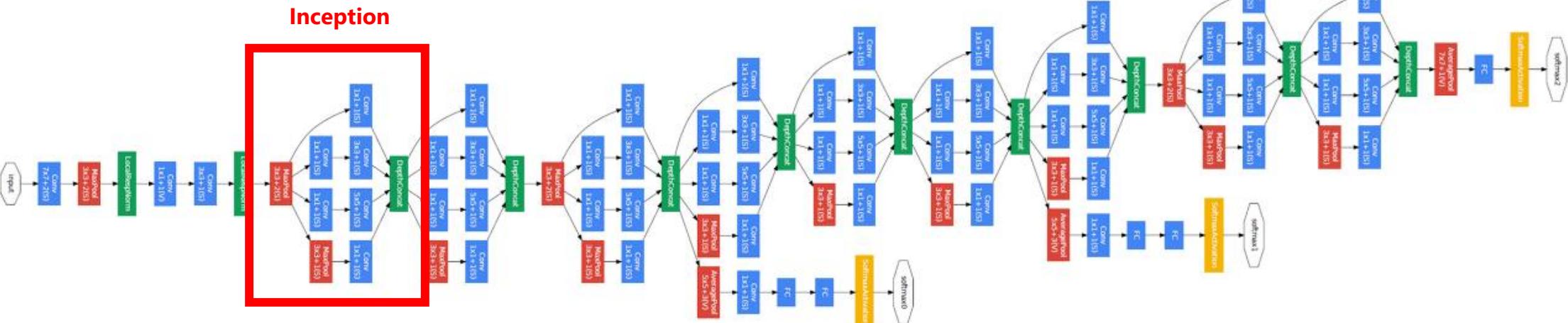
- 3x3 conv. Layer
 - Stacked conv. layers have a large receptive field
 - 3x3 (x2) – 5x5 receptive field
 - 3x3 (x3) – 7x7 receptive field
 - Less parameters to learn



GoogLeNet (Inception), 2014

- ✓ 구글이 만들고 모두가 사용하는 22계층으로 이뤄진 네트워크
- ✓ 핵심은 그래디언트 소실 문제를 해결하기 위해 만들어진 Inception module

Inception Module
Depth increased



GoogLeNet Architecture

1. Inception Module

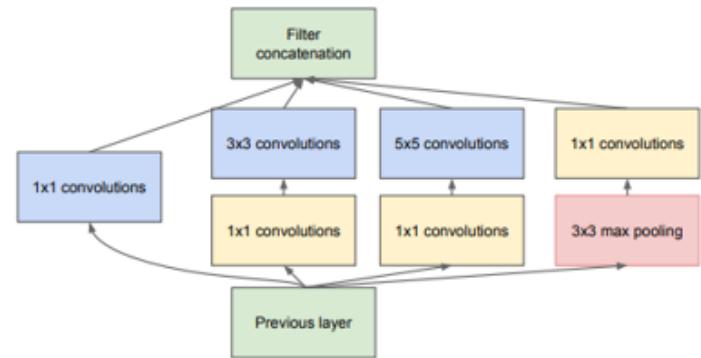
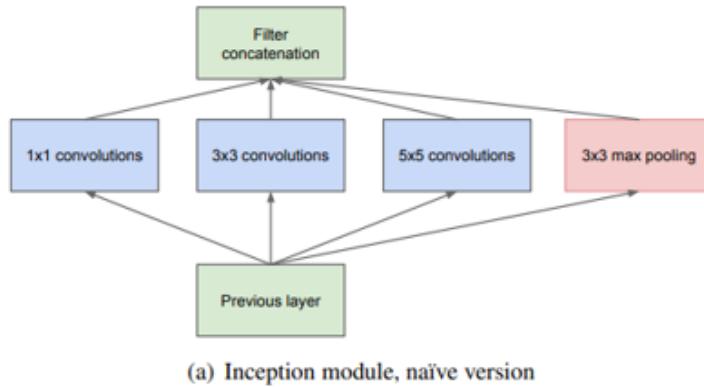
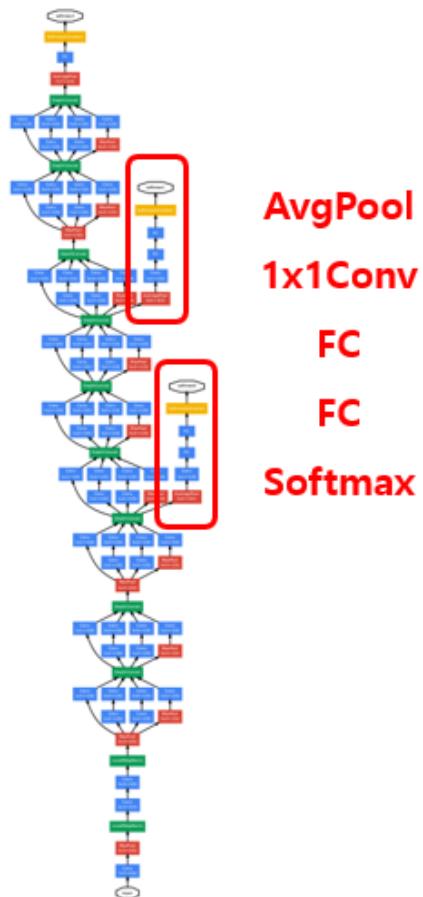
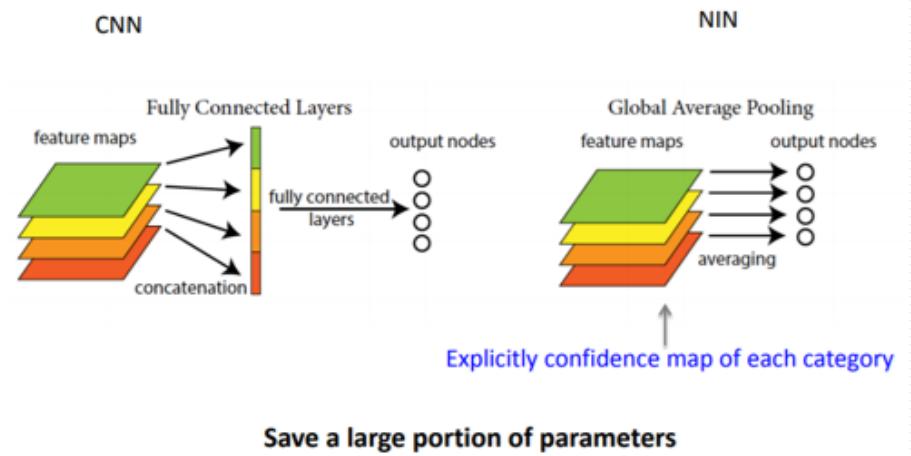


Figure 2: Inception module

2. Auxiliary Classifier



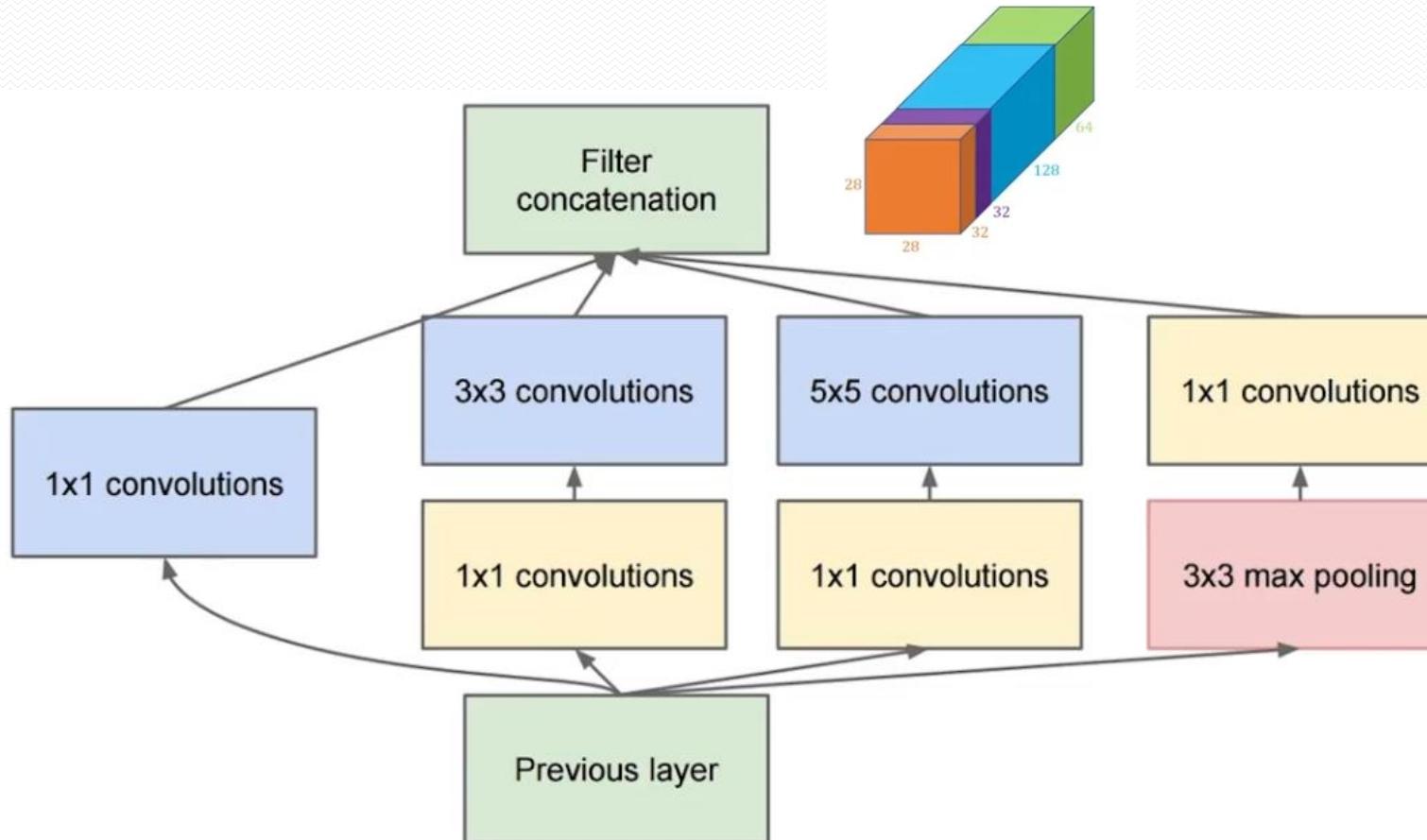
3. FC layer → Global Average Pooling



inception (5b)		$7 \times 7 \times 1024$
avg pool	$7 \times 7 / 1$	$1 \times 1 \times 1024$
dropout (40%)		$1 \times 1 \times 1024$
linear		$1 \times 1 \times 1000$
softmax		$1 \times 1 \times 1000$

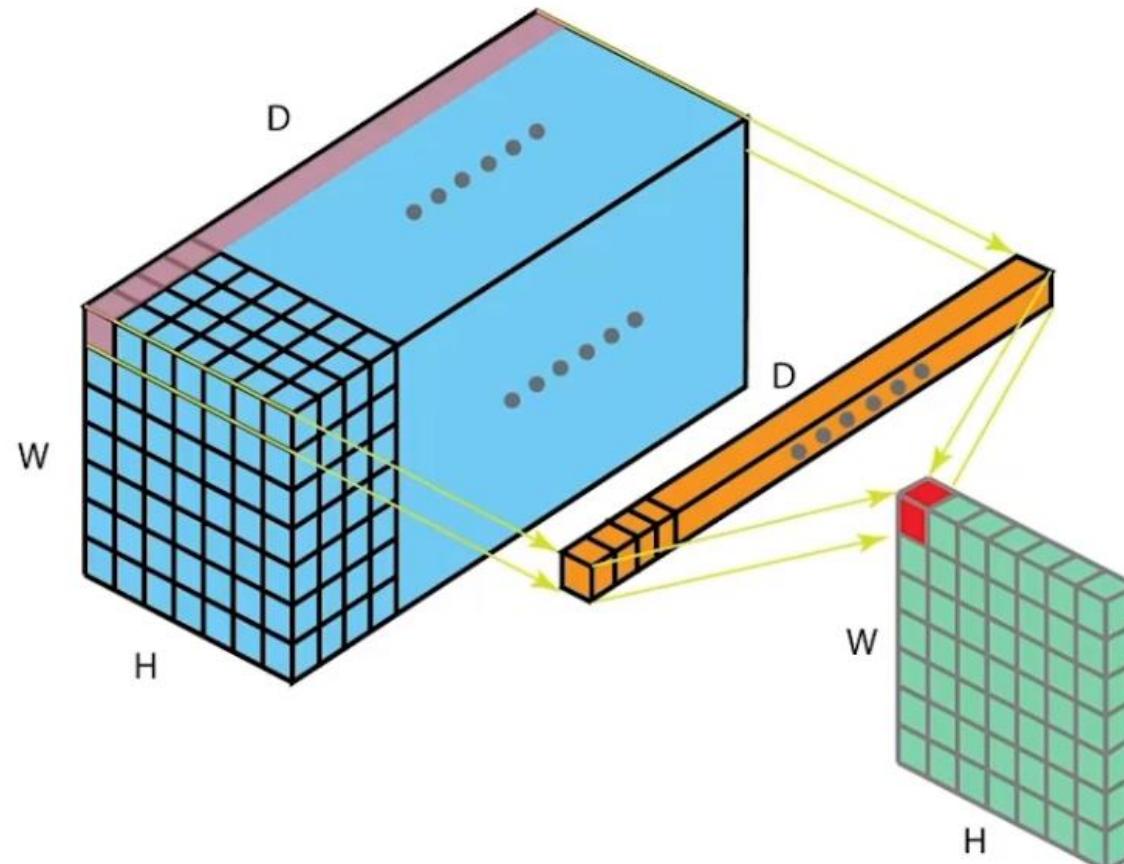
Inception module

- ✓ 왜 한 입력 image에 한 종류의 filter size만 적용해야 하는가?
- ✓ 뭐가 잘 될지 모르니 1 by 1, 3 by 3, 5 by 5, 3 by 3 max pooling을 다 해보자!



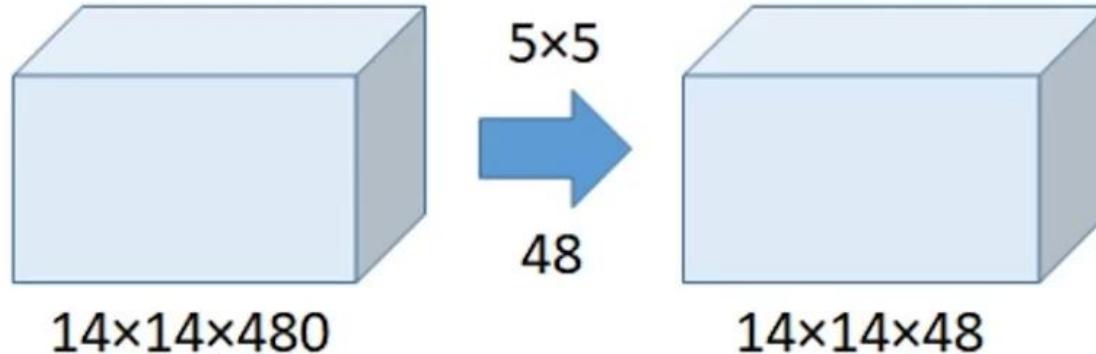
Inception module : 1x1 conv.

- ✓ 1 by 1 convolution의 역할 : feature map의 depth를 감소시켜 연산량을 감소시킴



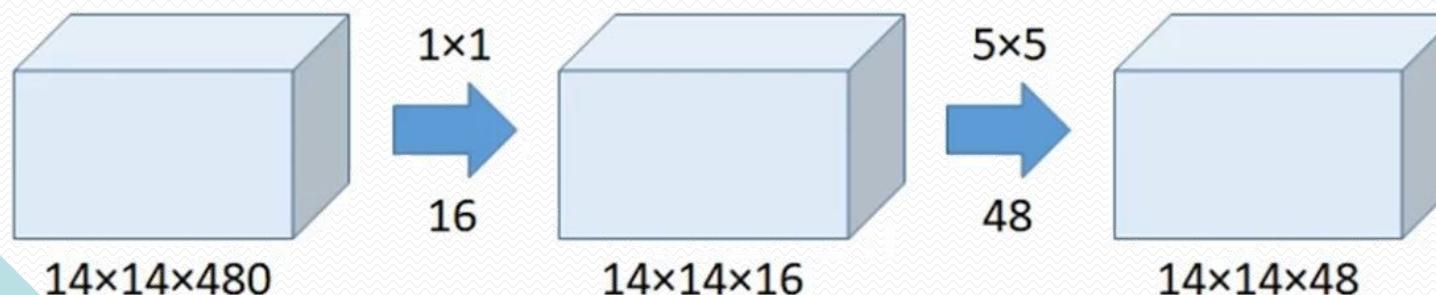
Inception module : 1x1 conv.

[1 by 1 convolution을 수행하지 않고 5 by 5 convolution을 수행]



✓ 연산량 : $(14 \times 14 \times 48) \times (5 \times 5 \times 480) = 112.9M$

[1 by 1 convolution을 수행후 5 by 5 convolution을 수행]



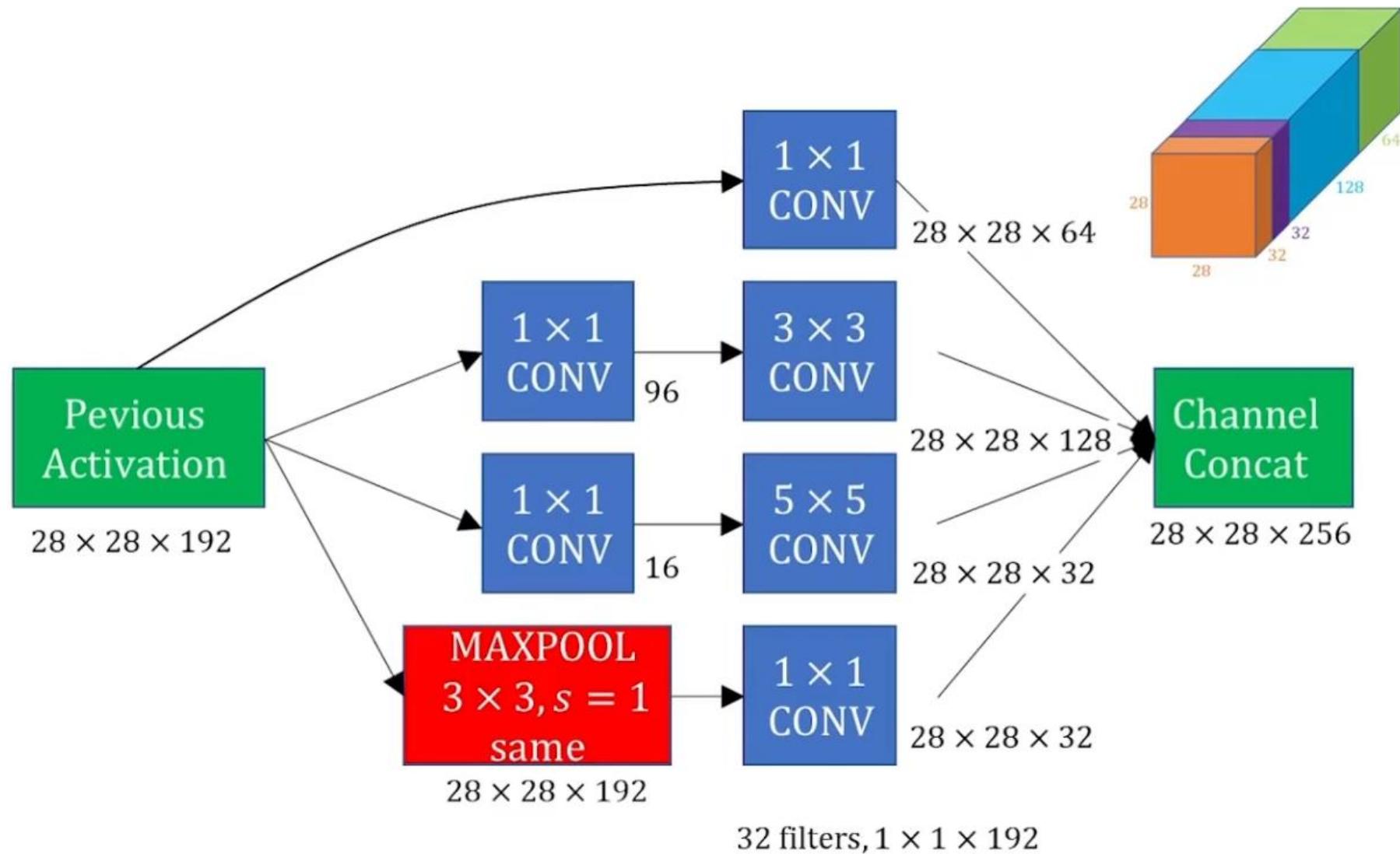
*연산량

- ✓ 1 by 1 Conv. : $(14 \times 14 \times 16) \times (1 \times 1 \times 480) = 1.5M$
- ✓ 5 by 5 Conv. : $(14 \times 14 \times 48) \times (5 \times 5 \times 16) = 3.8M$
- ✓ 총 연산량 : $1.5M + 3.8M = 5.3M << 112.9M$

✓ 1 by 1 convolution의 역할 : feature map의 depth를 감소시켜 연산량을 감소시킴

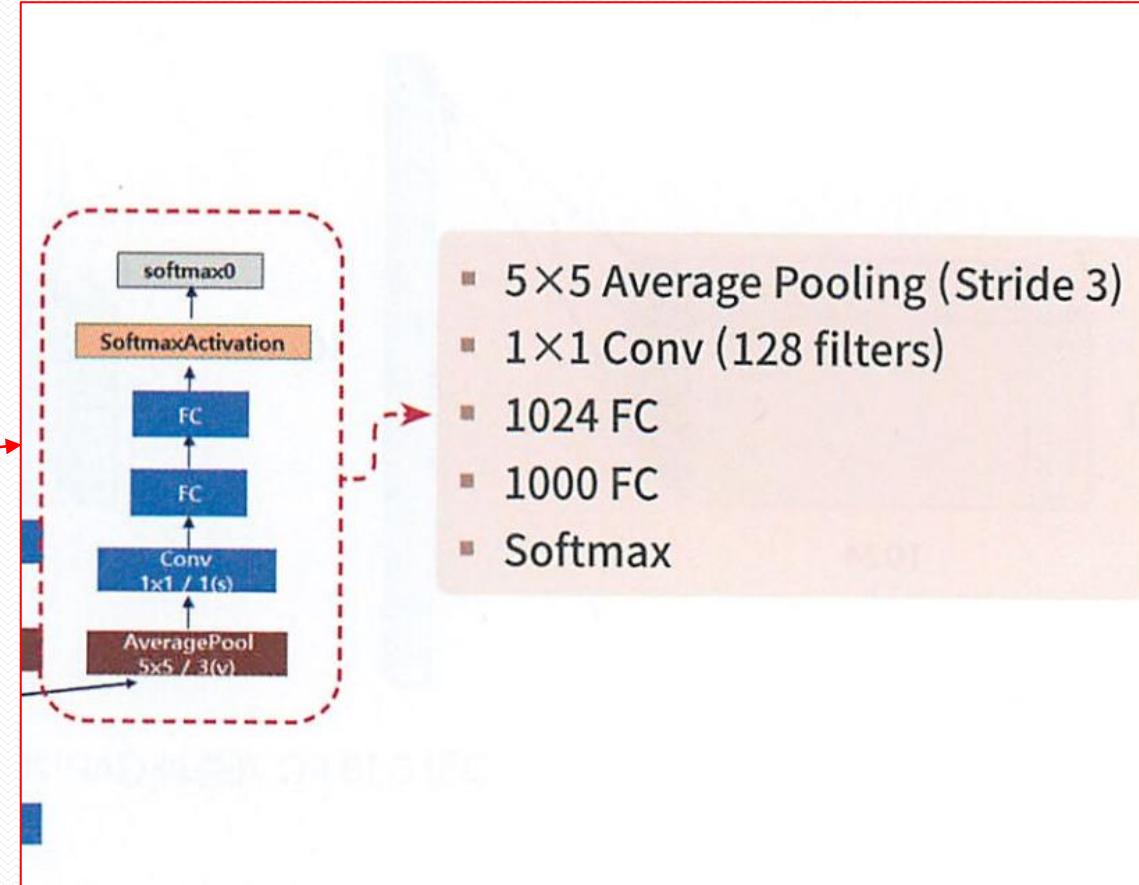
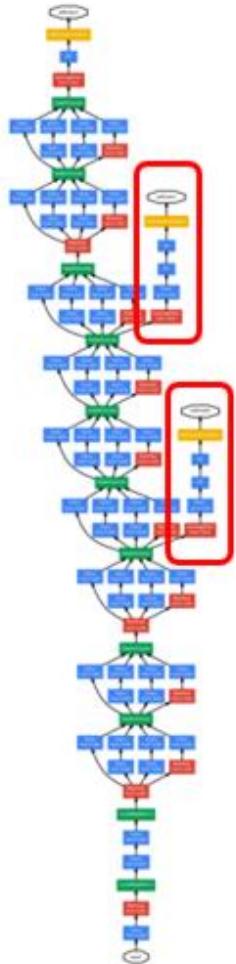
Inception module : concatnation

- ✓ Channel-wise concatnation



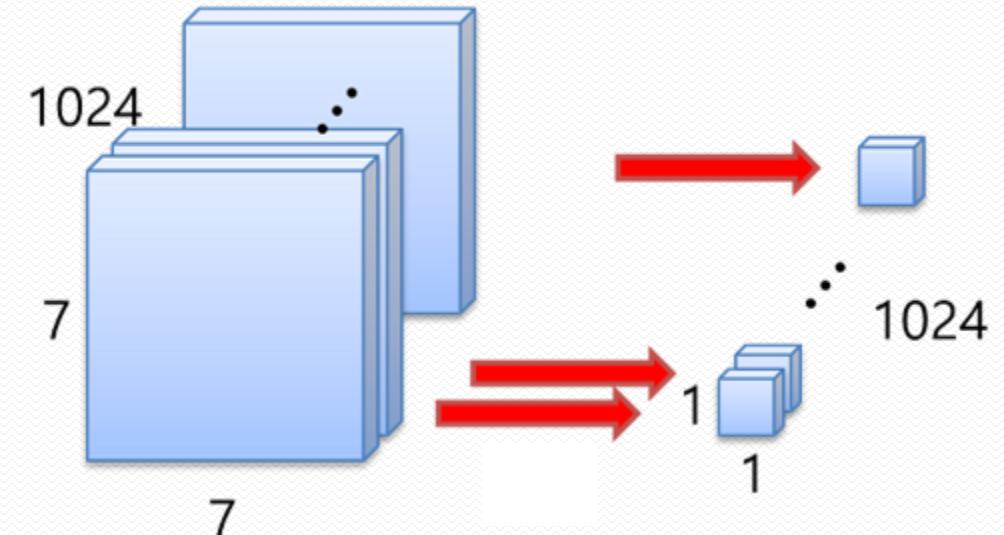
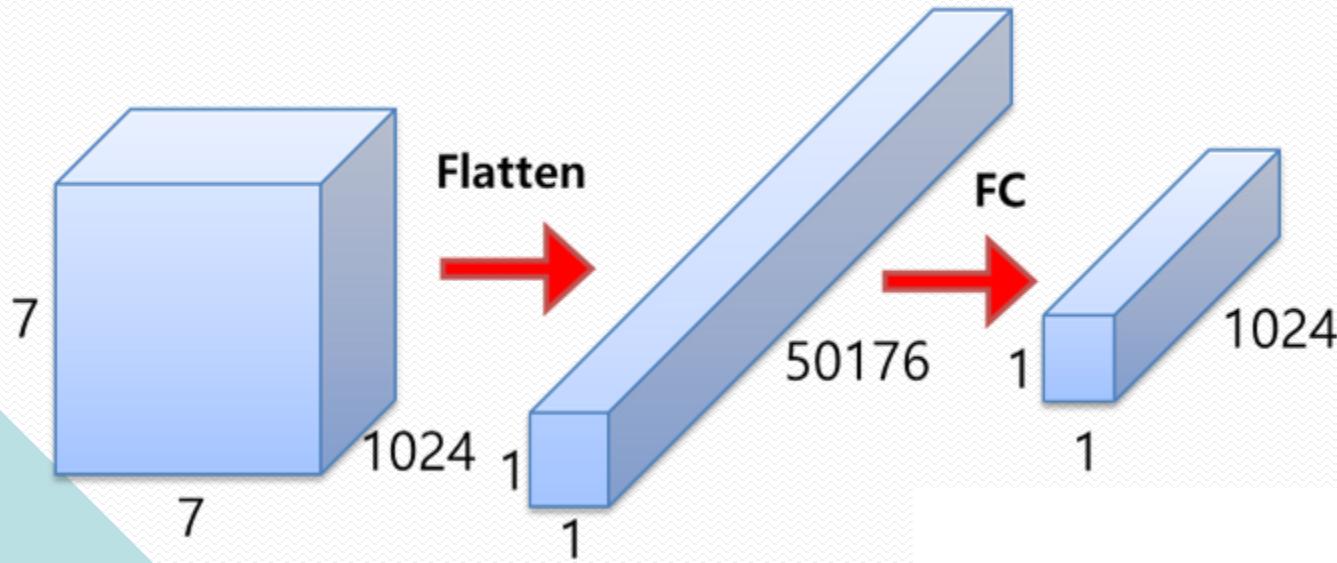
Auxiliary Classifiers

- 보조 분류기는 그래디언트 소실 문제를 해결하고 정규화를 제공

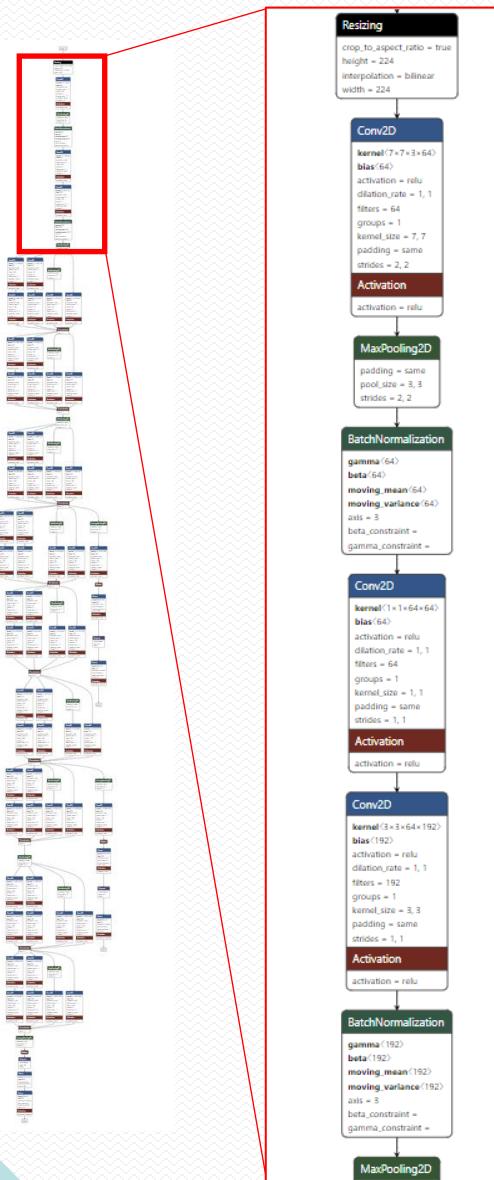


Global Average Pooling

- FC (Fully Connected)
 - Weights: $7 \times 7 \times 1024 \times 1024 = 51.3M$
- GAP (Global Average Pooling)
 - Weights: 0



GoogLeNet Architecture: <https://netron.app/>



```
# stage 1 (conv, polling, normalize, conv, conv, normalize, pooling)
model_in = Input(shape=f_image_train.shape[1:])
model = tf.keras.layers.Resizing(\n    224, 224, interpolation="bilinear", crop_to_aspect_ratio=True,\n    input_shape=f_image_train.shape[1:], name="Resizing_Layer")(model_in)

#model = tf.keras.layers.Rescaling(1./255)(model)
model = Conv2D(filters=64, kernel_size=(7,7), strides=2, padding = 'same', activation = 'relu')(model)
model = MaxPooling2D(pool_size=(3,3), strides=2, padding = 'same')(model)
model = BatchNormalization()(model)

# stage 2
model = Conv2D(filters=64, kernel_size=(1,1), strides=1, padding='same', activation='relu')(model)
model = Conv2D(filters=192, kernel_size=(3,3), strides=1, padding='same', activation='relu')(model)
model = BatchNormalization()(model)
model = MaxPooling2D(pool_size=(3,3), strides=2, padding='same')(model)

# stage 3
model = Inception(model, [64, (96, 128), (16, 32), 32])
model = Inception(model, [128, (128, 192), (32, 96), 64])
model = MaxPooling2D(pool_size=(3,3), strides=2, padding = 'same')(model)

# stage 4
model = Inception(model, [192, (96, 208), (16, 48), 64]) #a
aux1 = Auxiliary_classifier(model, name = 'aux1')
model = Inception(model, [160, (112, 224), (24, 64), 64])#b
model = Inception(model, [128, (128, 256), (24, 64), 64])#c
model = Inception(model, [112, (144, 288), (32, 64), 64])#d
aux2 = Auxiliary_classifier(model, name = 'aux2')
model = Inception(model, [256, (160, 320), (32, 128), 128])#e
model = MaxPooling2D(pool_size=(3, 3), strides = 2, padding = 'same')(model)

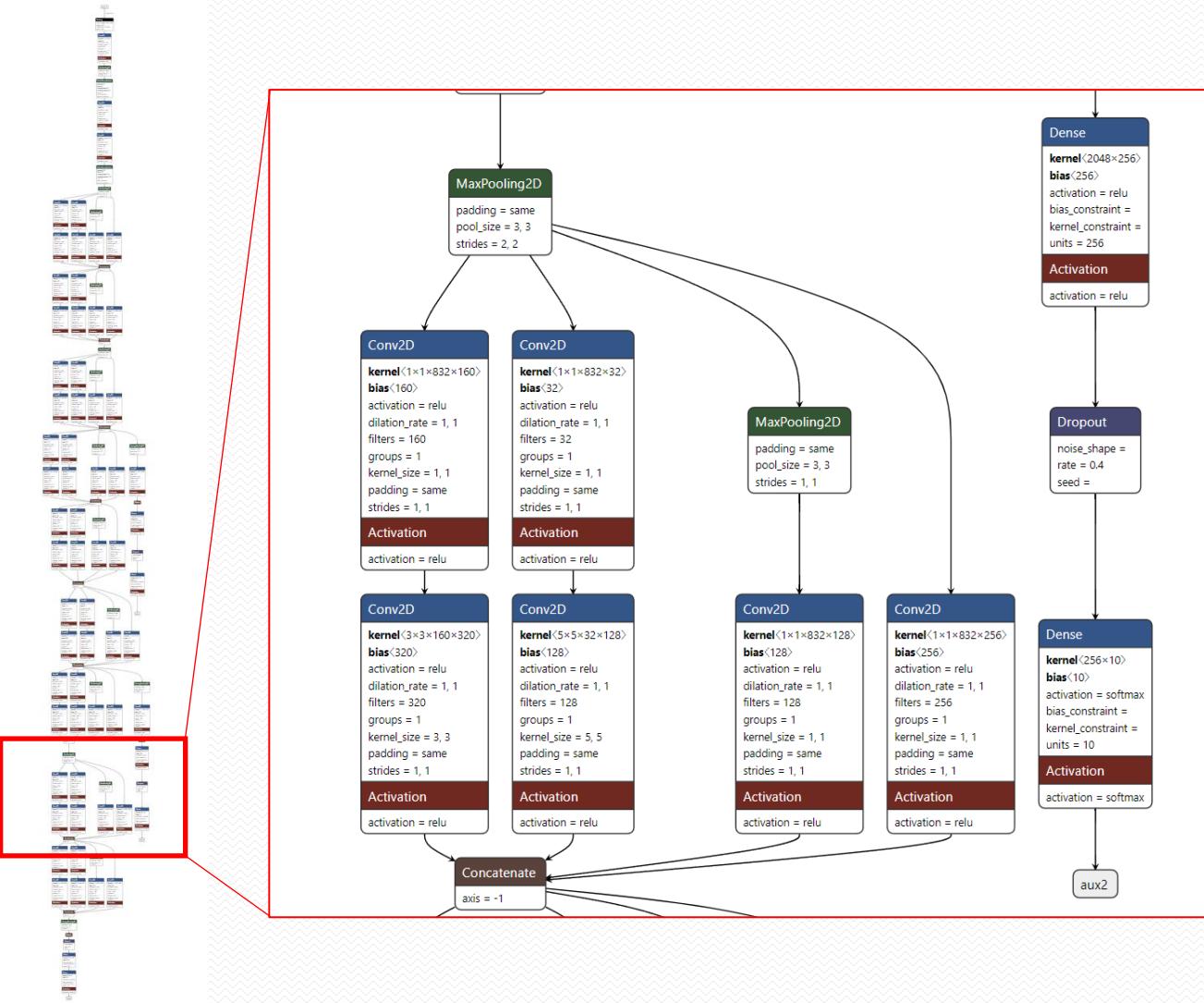
# stage 5
model = Inception(model, [256, (160, 320), (32, 128), 128])
model = Inception(model, [384, (192, 384), (48, 128), 128])
model = AveragePooling2D(pool_size=(7,7), strides = 1, padding = 'valid')(model)

#stage 6
model = Flatten()(model)
model = Dropout(0.4)(model)
#model = Dense(480, activation='linear')
model = Dense(units=256, activation='linear')(model)
main_branch = Dense(units=10, activation='softmax', name='main')(model)
model_fin = Model(inputs=model_in, outputs=[main_branch, aux1, aux2])

model_fin.compile(optimizer='adam',
                  loss='sparse_categorical_crossentropy',
                  metrics=['accuracy'])

model_fin.fit(f_image_train, f_label_train, epochs=100, batch_size = 5)
model_fin.save('GoogleNet_mnist.h5')
```

GoogLeNet Architecture



```

# stage 1 (conv, polling, normalize, conv, conv, normalize, pooling)
model_in = Input(shape=f_image_train.shape[1:])
model = tf.keras.layers.Resizing(\n    224, 224, interpolation="bilinear", crop_to_aspect_ratio=True,\n    input_shape=f_image_train.shape[1:], name="Resizing_Layer")(model_in)

#model = tf.keras.layers.Rescaling(1./255)(model)
model = Conv2D(filters=64, kernel_size=(7,7), strides=2, padding = 'same', activation = 'relu')(model)
model = MaxPooling2D(pool_size=(3,3), strides=2, padding = 'same')(model)
model = BatchNormalization()(model)

# stage 2
model = Conv2D(filters=64, kernel_size=(1,1), strides=1, padding='same', activation='relu')(model)
model = Conv2D(filters=192, kernel_size=(3,3), strides=1, padding='same', activation='relu')(model)
model = BatchNormalization()(model)
model = MaxPooling2D(pool_size=(3,3), strides=2, padding='same')(model)

# stage 3
model = Inception(model, [64, (96, 128), (16, 32), 32])
model = Inception(model, [128, (128, 192), (32, 96), 64])
model = MaxPooling2D(pool_size=(3,3), strides=2, padding = 'same')(model)

# stage 4
model = Inception(model, [192, (96, 208), (16, 48), 64]) # a
aux1 = Auxiliary_classifier(model, name = 'aux1')
model = Inception(model, [160, (112, 224), (24, 64), 64]) # b
model = Inception(model, [128, (128, 256), (24, 64), 64]) # c
model = Inception(model, [112, (144, 288), (32, 64), 64]) # d
aux2 = Auxiliary_classifier(model, name = 'aux2')
model = Inception(model, [256, (160, 320), (32, 128), 128]) # e
model = MaxPooling2D(pool_size=(3, 3), strides = 2, padding = 'same')(model)

# stage 5
model = Inception(model, [256, (160, 320), (32, 128), 128])
model = Inception(model, [384, (192, 384), (48, 128), 128])
model = AveragePooling2D(pool_size=(7,7), strides = 1, padding = 'valid')(model)

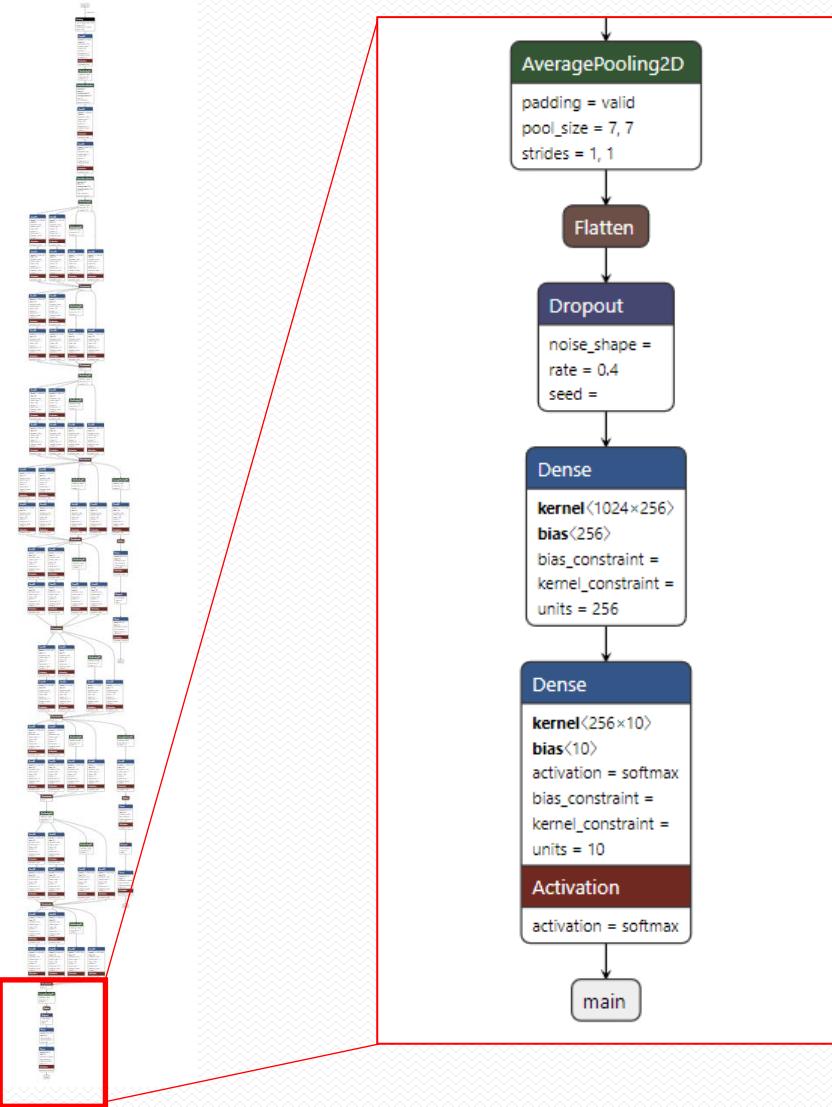
#stage 6
model = Flatten()(model)
model = Dropout(0.4)(model)
#model = Dense(480, activation='linear')
model = Dense(units=256, activation='linear')(model)
main_branch = Dense(units=10, activation='softmax', name='main')(model)
model_fin = Model(inputs=model_in, outputs=[main_branch, aux1, aux2])

model_fin.compile(optimizer='adam',
                  loss='sparse_categorical_crossentropy',
                  metrics=['accuracy'])

model_fin.fit(f_image_train, f_label_train, epochs=100, batch_size = 5)
model_fin.save('GoogleNet_mnist.h5')

```

GoogLeNet Architecture



```
# stage 1 (conv, polling, normalize, conv, conv, normalize, pooling)
model_in = Input(shape=f_image_train.shape[1:])
model = tf.keras.layers.Resizing(\n    224, 224, interpolation="bilinear", crop_to_aspect_ratio=True,\n    input_shape=f_image_train.shape[1:], name="Resizing_Layer")(model_in)

#model = tf.keras.layers.Rescaling(1./255)(model)
model = Conv2D(filters=64, kernel_size=(7,7), strides=2, padding = 'same', activation = 'relu')(model)
model = MaxPooling2D(pool_size=(3,3), strides=2, padding = 'same')(model)
model = BatchNormalization()(model)

# stage 2
model = Conv2D(filters=64, kernel_size=(1,1), strides=1, padding='same', activation='relu')(model)
model = Conv2D(filters=192, kernel_size=(3,3), strides=1, padding='same', activation='relu')(model)
model = BatchNormalization()(model)
model = MaxPooling2D(pool_size=(3,3), strides=2, padding='same')(model)

# stage 3
model = Inception(model, [64, (96, 128), (16, 32), 32])
model = Inception(model, [128, (128, 192), (32, 96), 64])
model = MaxPooling2D(pool_size=(3,3), strides=2, padding = 'same')(model)

# stage 4
model = Inception(model, [192, (96, 208), (16, 48), 64]) #a
aux1 = Auxiliary_classifier(model, name = 'aux1')
model = Inception(model, [160, (112, 224), (24, 64), 64])#b
model = Inception(model, [128, (128, 256), (24, 64), 64])#c
model = Inception(model, [112, (144, 288), (32, 64), 64])#d
aux2 = Auxiliary_classifier(model, name = 'aux2')
model = Inception(model, [256, (160, 320), (32, 128), 128])#e
model = MaxPooling2D(pool_size=(3, 3), strides = 2, padding = 'same')(model)

# stage 5
model = Inception(model, [256, (160, 320), (32, 128), 128])
model = Inception(model, [384, (192, 384), (48, 128), 128])
model = AveragePooling2D(pool_size=(7,7), strides = 1, padding = 'valid')(model)

#stage 6
model = Flatten()(model)
model = Dropout(0.4)(model)
#model = Dense(480, activation='linear')
model = Dense(units=256, activation='linear')(model)
main_branch = Dense(units=10, activation='softmax', name='main')(model)
model_fin = Model(inputs=model_in, outputs=[main_branch, aux1, aux2])

model_fin.compile(optimizer='adam',
                  loss='sparse_categorical_crossentropy',
                  metrics=['accuracy'])

model_fin.fit(f_image_train, f_label_train, epochs=100, batch_size = 5)
model_fin.save('GoogleNet_mnist.h5')
```

GoogLeNet Architecture

type	patch size/ stride	output size	depth	#1×1	#3×3 reduce	#3×3	#5×5 reduce	#5×5	pool proj	proj size
convolution	$7 \times 7/2$	$112 \times 112 \times 64$	1							
max pool	$3 \times 3/2$	$56 \times 56 \times 64$	0							
convolution	$3 \times 3/1$	$56 \times 56 \times 192$	2		64	192				
max pool	$3 \times 3/2$	$28 \times 28 \times 192$	0							
inception (3a)		$28 \times 28 \times 256$	2	64	96	128	16	32	32	
inception (3b)		$28 \times 28 \times 480$	2	128	128	192	32	96	64	
max pool	$3 \times 3/2$	$14 \times 14 \times 480$	0							
inception (4a)		$14 \times 14 \times 512$	2	192	96	208	16	48	64	
inception (4b)		$14 \times 14 \times 512$	2	160	112	224	24	64	64	
inception (4c)		$14 \times 14 \times 512$	2	128	128	256	24	64	64	
inception (4d)		$14 \times 14 \times 528$	2	112	144	288	32	64	64	
inception (4e)		$14 \times 14 \times 832$	2	256	160	320	32	128	128	
max pool	$3 \times 3/2$	$7 \times 7 \times 832$	0							
inception (5a)		$7 \times 7 \times 832$	2	256	160	320	32	128	128	
inception (5b)		$7 \times 7 \times 1024$	2	384	192	384	48	128	128	
avg pool	$7 \times 7/1$	$1 \times 1 \times 1024$	0							
dropout (40%)		$1 \times 1 \times 1024$	0							
linear		$1 \times 1 \times 1000$	1							
softmax		$1 \times 1 \times 1000$	0							

```

# stage 1 (conv, polling, normalize, conv, conv, normalize, pooling)
model_in = Input(shape=f_image_train.shape[1:])
model = tf.keras.layers.Resizing(\n    224, 224, interpolation="bilinear", crop_to_aspect_ratio=True,\n    input_shape=f_image_train.shape[1:], name="Resizing_Layer")(model_in)

#model = tf.keras.layers.Rescaling(1./255)(model)
model = Conv2D(filters=64, kernel_size=(7,7), strides=2, padding = 'same', activation = 'relu')(model)
model = MaxPooling2D(pool_size=(3,3), strides=2, padding = 'same')(model)
model = BatchNormalization()(model)
model = MaxPooling2D(pool_size=(3,3), strides=2, padding='same')(model)

# stage 2
model = Conv2D(filters=64, kernel_size=(1,1), strides=1, padding='same', activation='relu')(model)
model = Conv2D(filters=192, kernel_size=(3,3), strides=1, padding='same', activation='relu')(model)
model = BatchNormalization()(model)
model = MaxPooling2D(pool_size=(3,3), strides=2, padding='same')(model)

# stage 3
model = Inception(model, [64, (96, 128), (16, 32), 32])
model = Inception(model, [128, (128, 192), (32, 96), 64])
model = MaxPooling2D(pool_size=(3,3), strides=2, padding = 'same')(model)

# stage 4
model = Inception(model, [192, (96, 208), (16, 48), 64]) #a
aux1 = Auxiliary_classifier(model, name = 'aux1')
model = Inception(model, [160, (112, 224), (24, 64), 64])#b
model = Inception(model, [128, (128, 256), (24, 64), 64])#c
model = Inception(model, [112, (144, 288), (32, 64), 64])#d
aux2 = Auxiliary_classifier(model, name = 'aux2')
model = Inception(model, [256, (160, 320), (32, 128), 128])#e
model = MaxPooling2D(pool_size=(3, 3), strides = 2, padding = 'same')(model)

# stage 5
model = Inception(model, [256, (160, 320), (32, 128), 128])
model = Inception(model, [384, (192, 384), (48, 128), 128])
model = AveragePooling2D(pool_size=(7,7), strides = 1, padding = 'valid')(model)

#stage 6
model = Flatten()(model)
model = Dropout(0.4)(model)
#model = Dense(480, activation='linear')
model = Dense(units=256, activation='linear')(model)
main_branch = Dense(units=10, activation='softmax', name='main')(model)
model_fin = Model(inputs=model_in, outputs=[main_branch, aux1, aux2])

model_fin.compile(optimizer='adam',
                  loss='sparse_categorical_crossentropy',
                  metrics=['accuracy'])

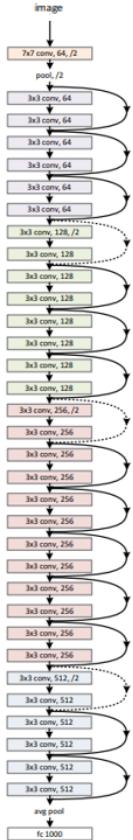
model_fin.fit(f_image_train, f_label_train, epochs=100, batch_size = 5)
model_fin.save('GoogleNet_mnist.h5')

```

ResNet, 2015

- Revolution of depth
- 가장 깊고 성능이 좋은 네트워크

34-layer residual



		18-layer	34-layer	50-layer	101-layer	152-layer
layer name	output size			7×7, 64, stride 2		
conv1	112×112					
conv2_x	56×56	$\left[\begin{matrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{matrix} \right] \times 2$	$\left[\begin{matrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{matrix} \right] \times 3$	$\left[\begin{matrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{matrix} \right] \times 3$	$\left[\begin{matrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{matrix} \right] \times 3$	$\left[\begin{matrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{matrix} \right] \times 3$
conv3_x	28×28	$\left[\begin{matrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{matrix} \right] \times 2$	$\left[\begin{matrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{matrix} \right] \times 4$	$\left[\begin{matrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{matrix} \right] \times 4$	$\left[\begin{matrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{matrix} \right] \times 4$	$\left[\begin{matrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{matrix} \right] \times 8$
conv4_x	14×14	$\left[\begin{matrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{matrix} \right] \times 2$	$\left[\begin{matrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{matrix} \right] \times 6$	$\left[\begin{matrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{matrix} \right] \times 6$	$\left[\begin{matrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{matrix} \right] \times 23$	$\left[\begin{matrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{matrix} \right] \times 36$
conv5_x	7×7	$\left[\begin{matrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{matrix} \right] \times 2$	$\left[\begin{matrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{matrix} \right] \times 3$	$\left[\begin{matrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{matrix} \right] \times 3$	$\left[\begin{matrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{matrix} \right] \times 3$	$\left[\begin{matrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{matrix} \right] \times 3$
	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

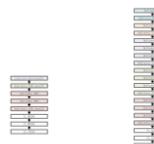
Residual Network

Great Depth

Better recognition than human

ResNet, 152 layers
(ILSVRC 2015)

VGG, 19 layers
(ILSVRC 2014)

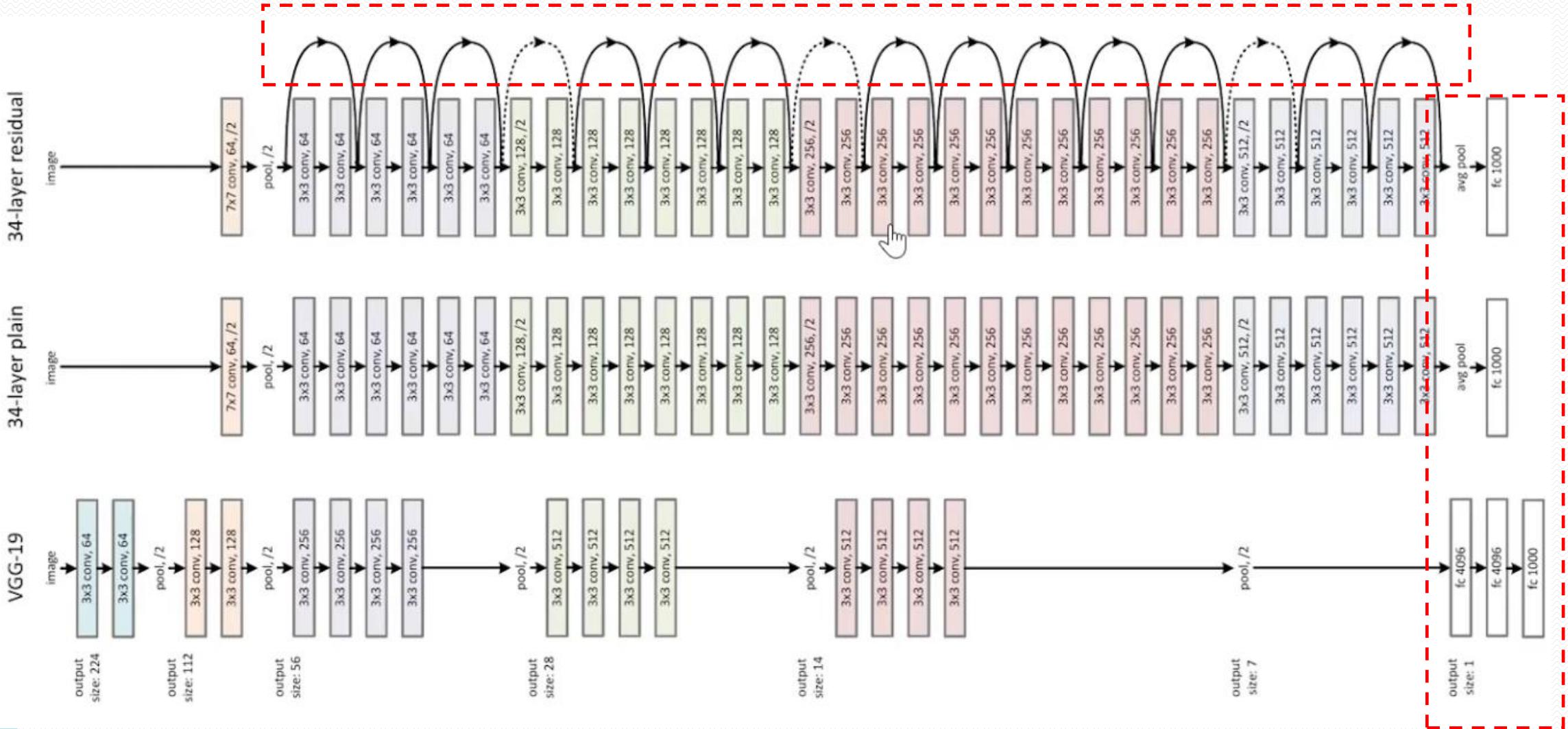


AlexNet, 8 layers
(ILSVRC 2012)



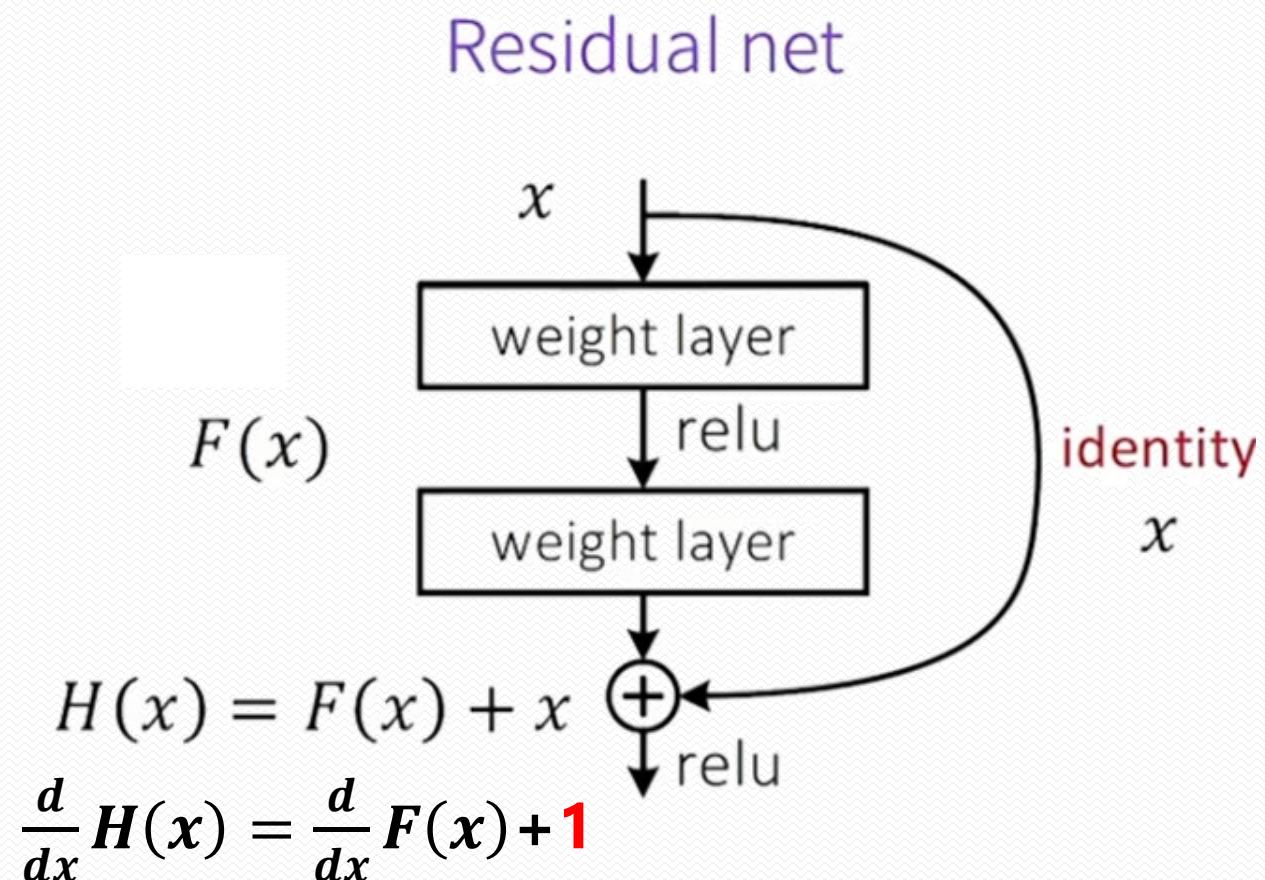
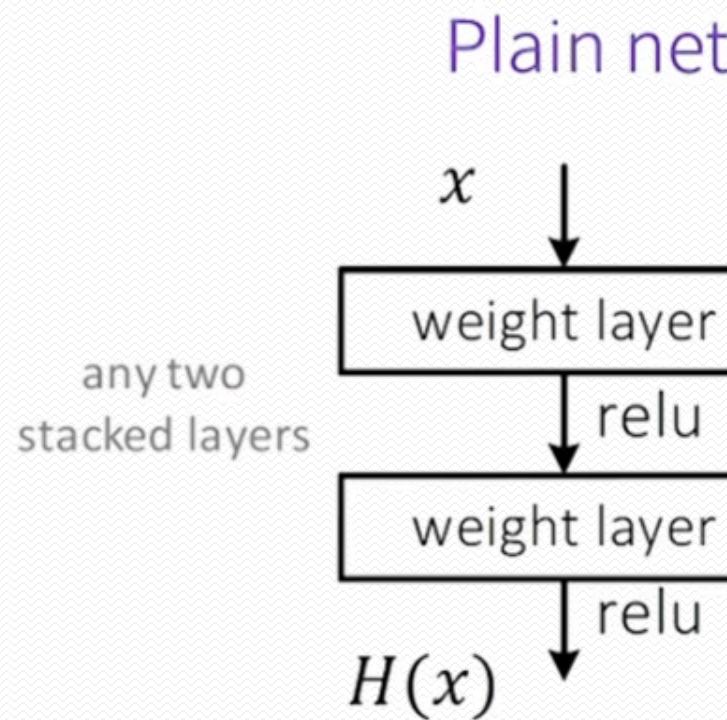
ResNet, 2015

- 핵심은 residual learning

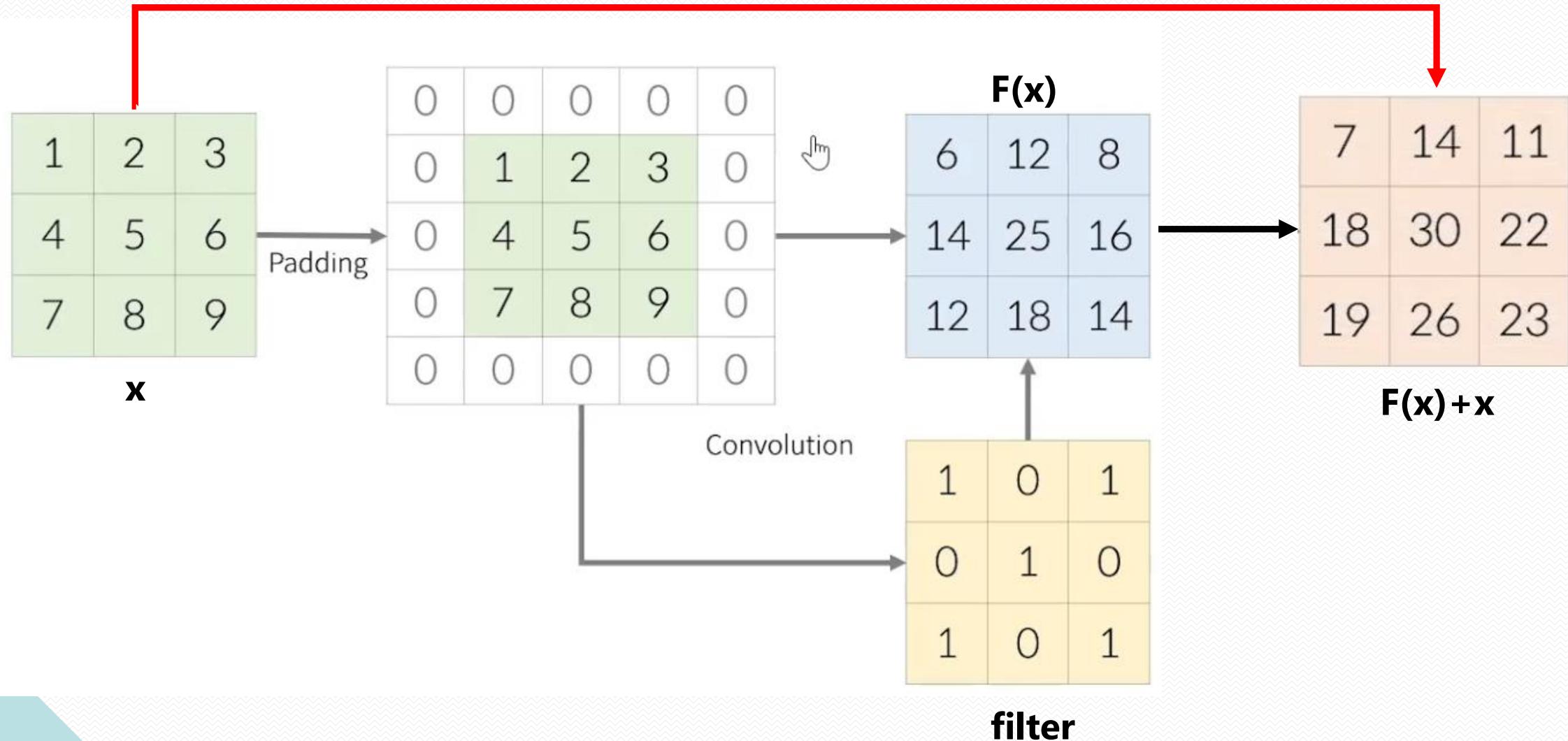


ResNet, 2015 : residual block (skip connection)

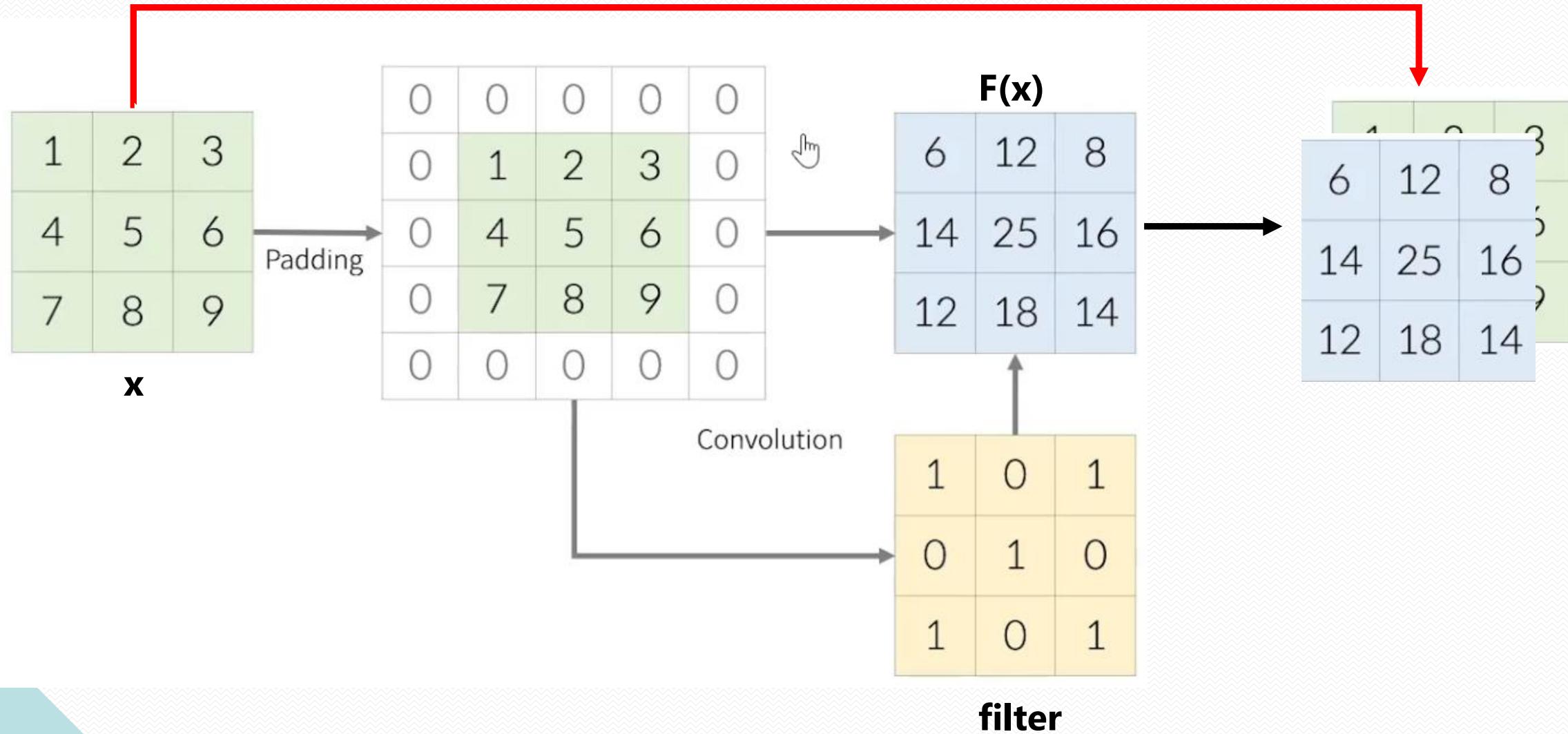
- ✓ 일정 시점마다 input x 자체를 skip connection을 통해 연결
- ✓ Gradient flow가 원활하게 이루어짐 (모델을 깊게 쌓는 것에 대한 부담이 줄어듦)
- ✓ Fully connected layer 제거 (모델의 파라미터 수를 줄임)



ResNet, 2015 : skip connection



DenseNet : channel-wise concatnation



Network Design

- Simple design
 - All 3x3 conv (almost)
 - No Hidden FC
 - Spatial size /2 → # filters x2

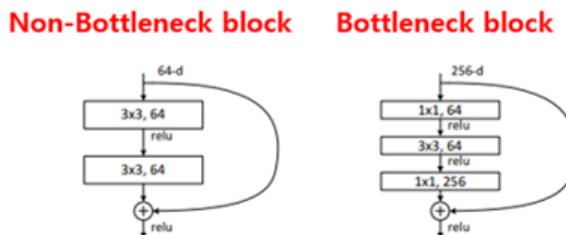


Figure 5. A deeper residual function \mathcal{F} for ImageNet. Left: a building block (on 56×56 feature maps) as in Fig. 3 for ResNet-34. Right: a “bottleneck” building block for ResNet-50/101/152.

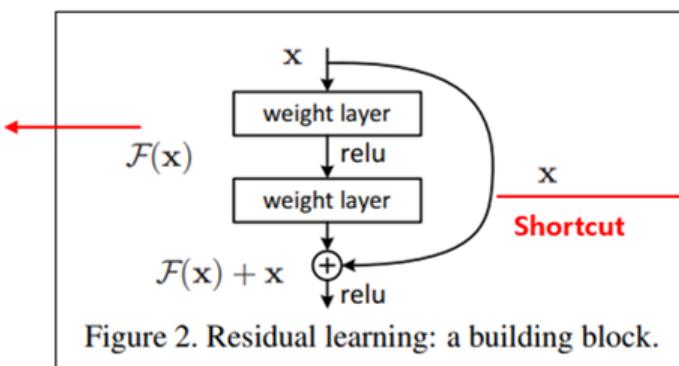


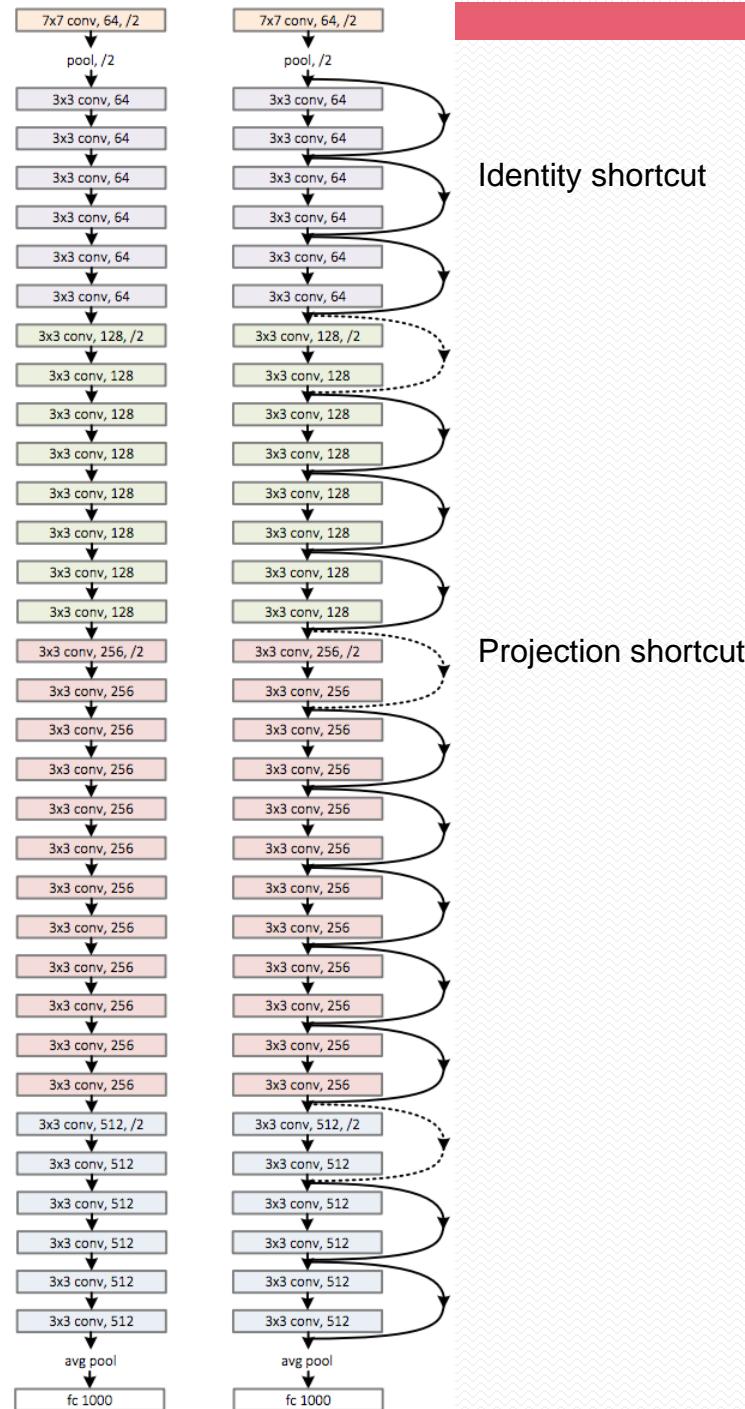
Figure 2. Residual learning: a building block.

$$\mathbf{y} = \mathcal{F}(\mathbf{x}, \{W_i\}) + \mathbf{x}$$

Identity Shortcut

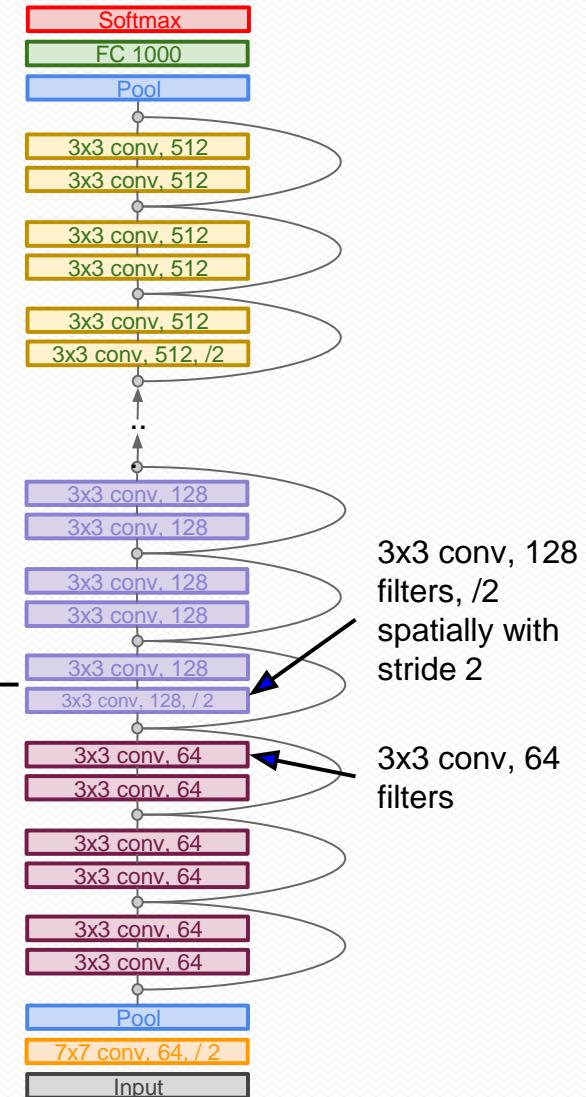
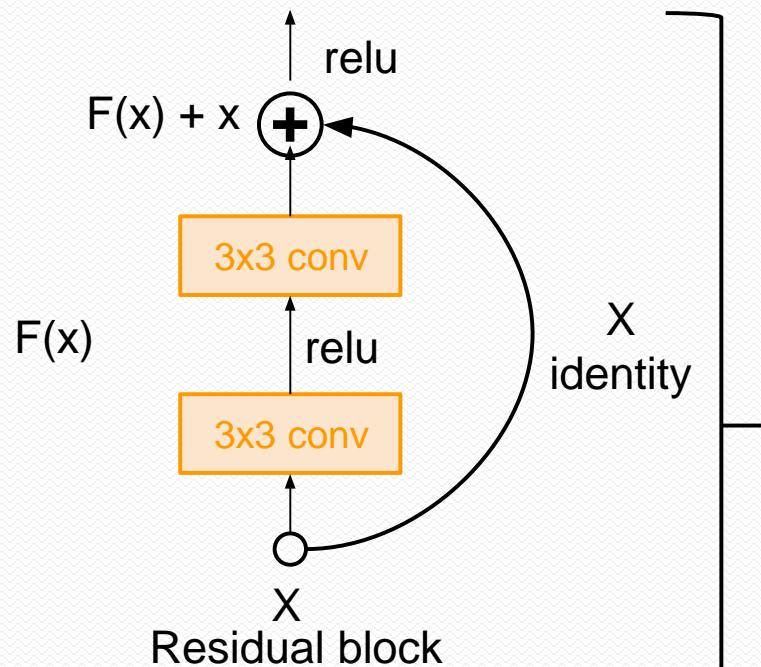
$$\mathbf{y} = \mathcal{F}(\mathbf{x}, \{W_i\}) + W_s \mathbf{x}$$

Projection Shortcut

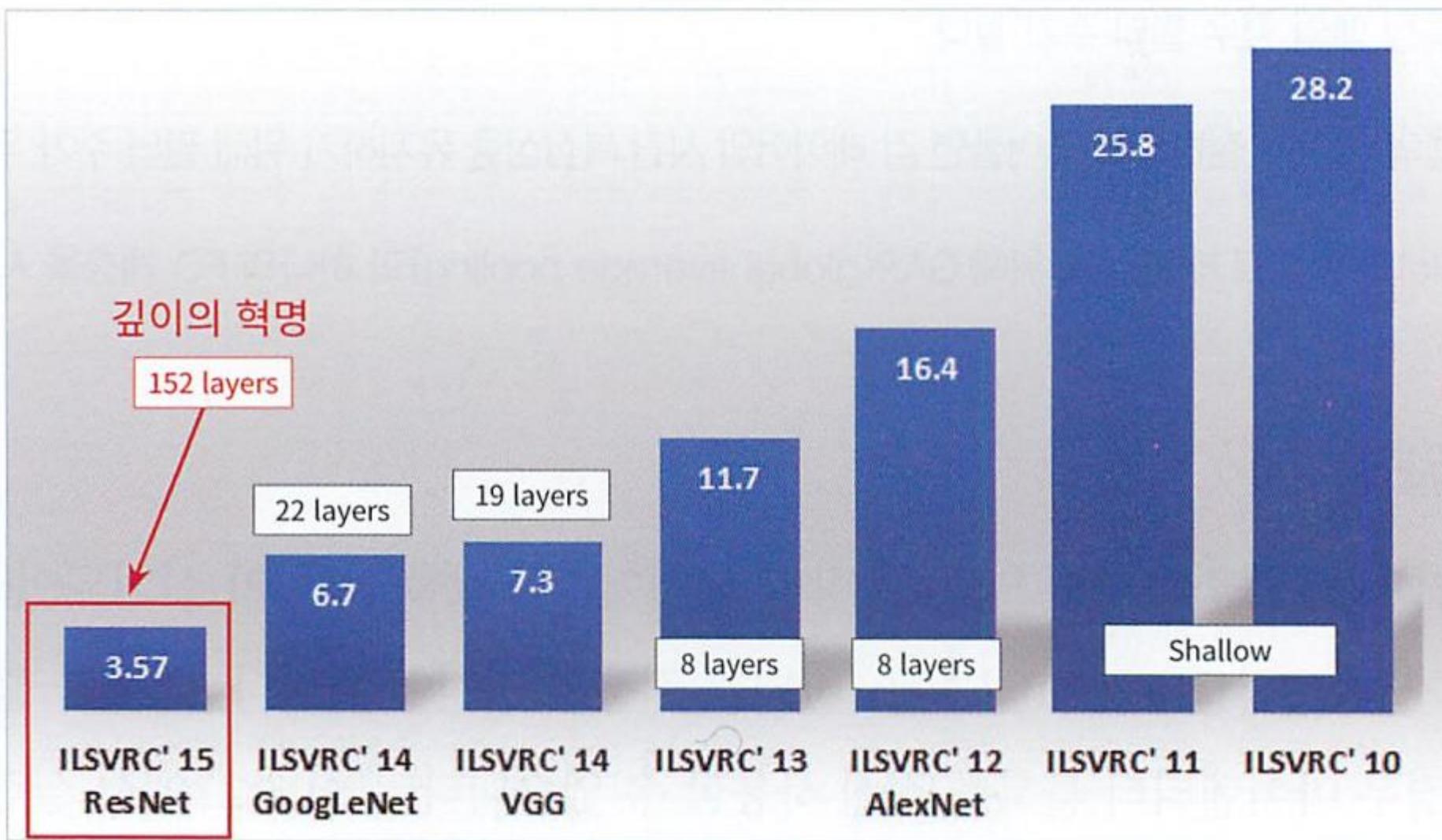


Full ResNet architecture

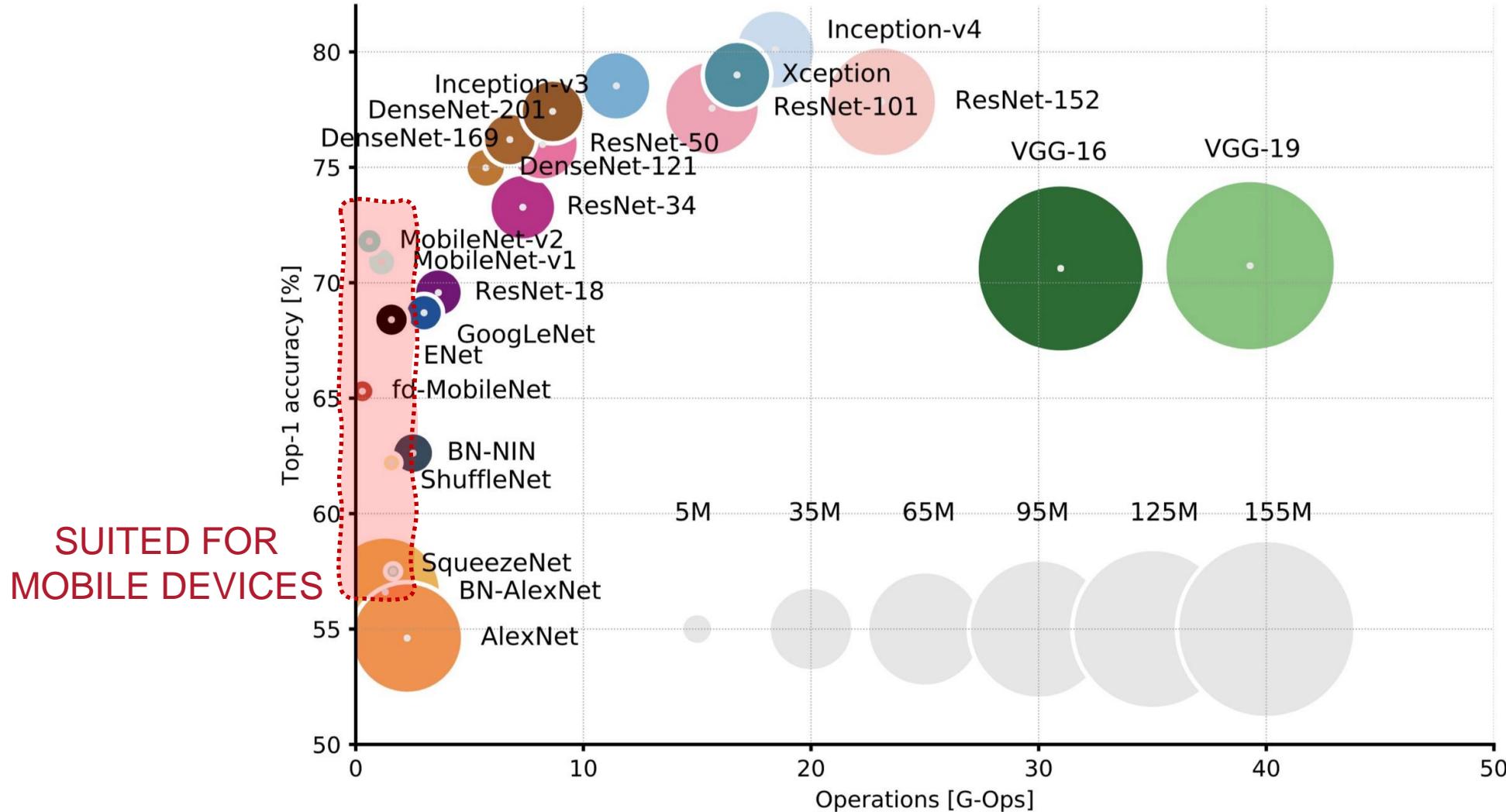
- Stack residual blocks
 - Consist of two 3x3 conv layers
- Periodically,
 - Double # of filters
 - Downsample spatially
By using stride of 2
 - (2 in each dimension)

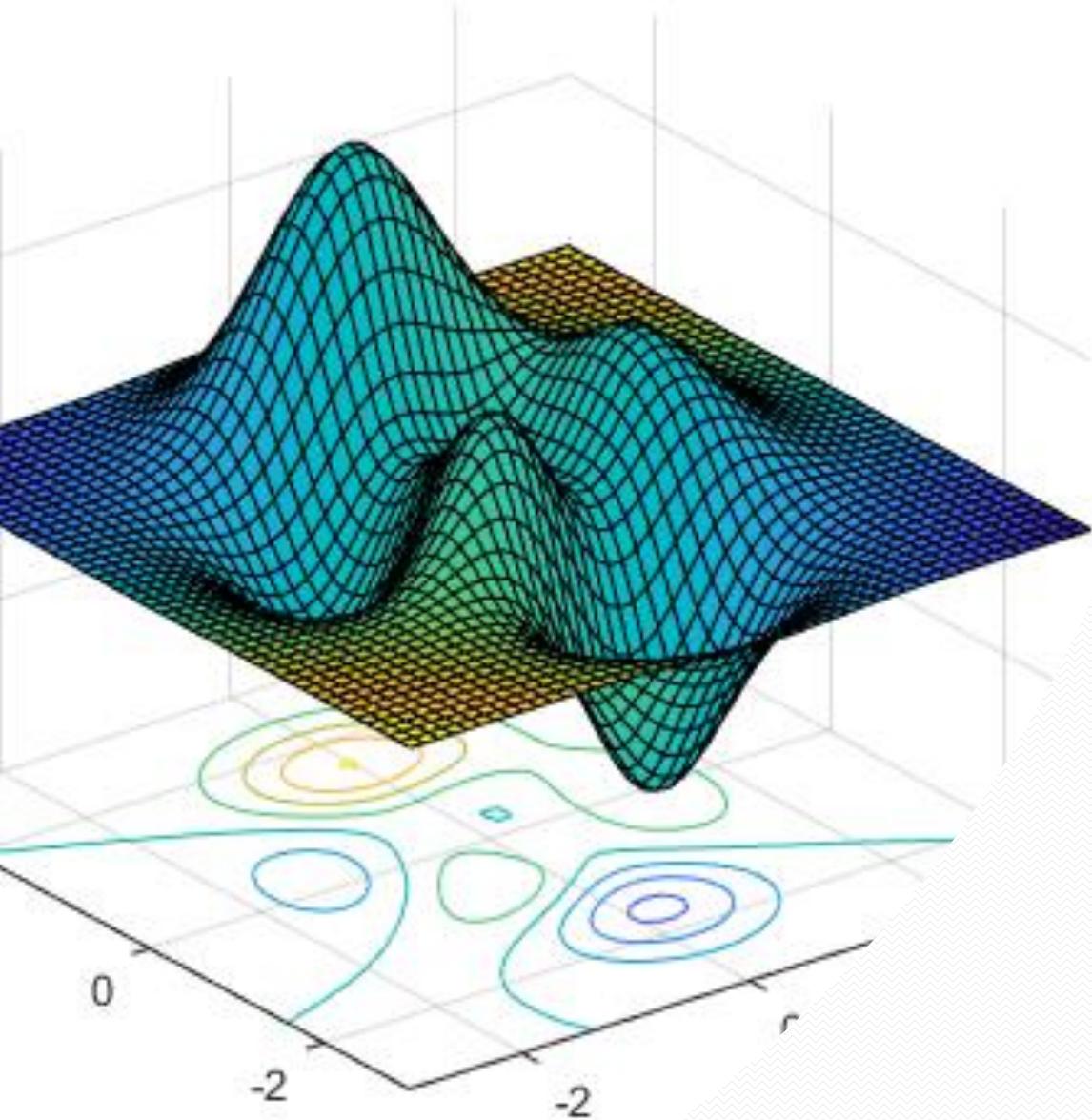


ResNet, 2015



TOP 1 VS. OPERATIONS, PARAMETERS

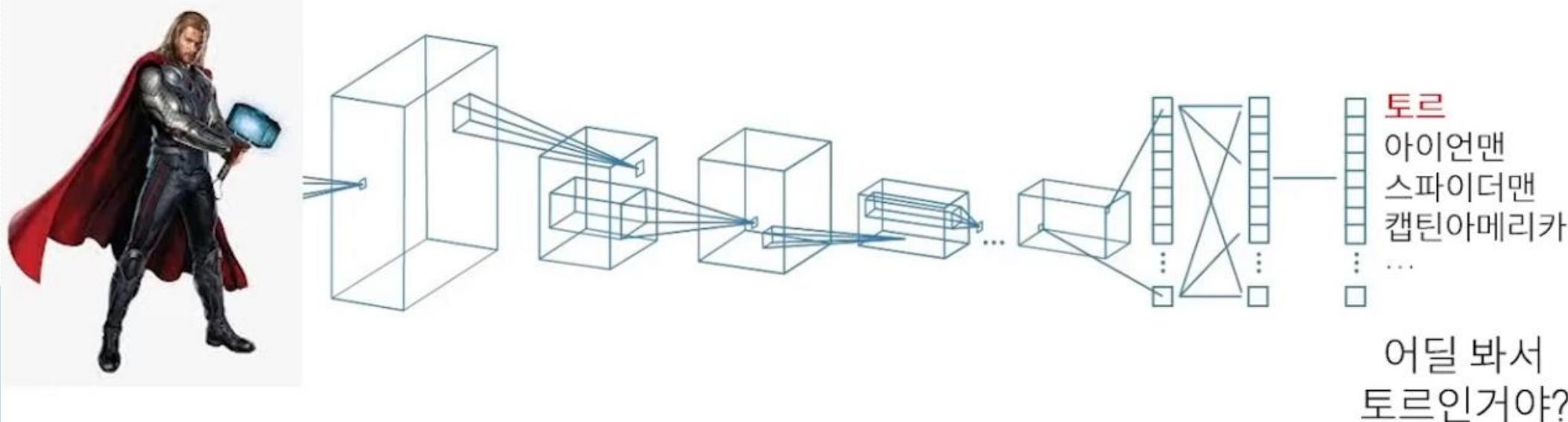
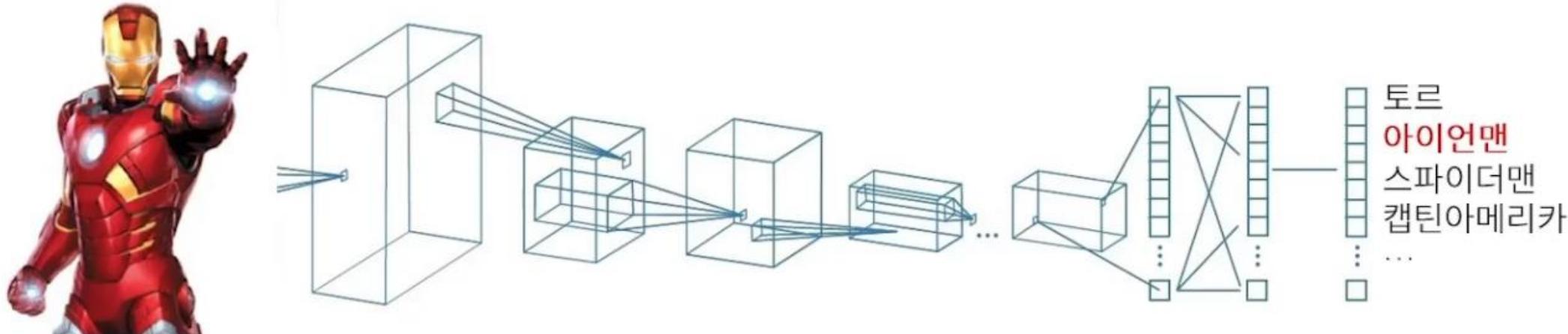




Class Localization

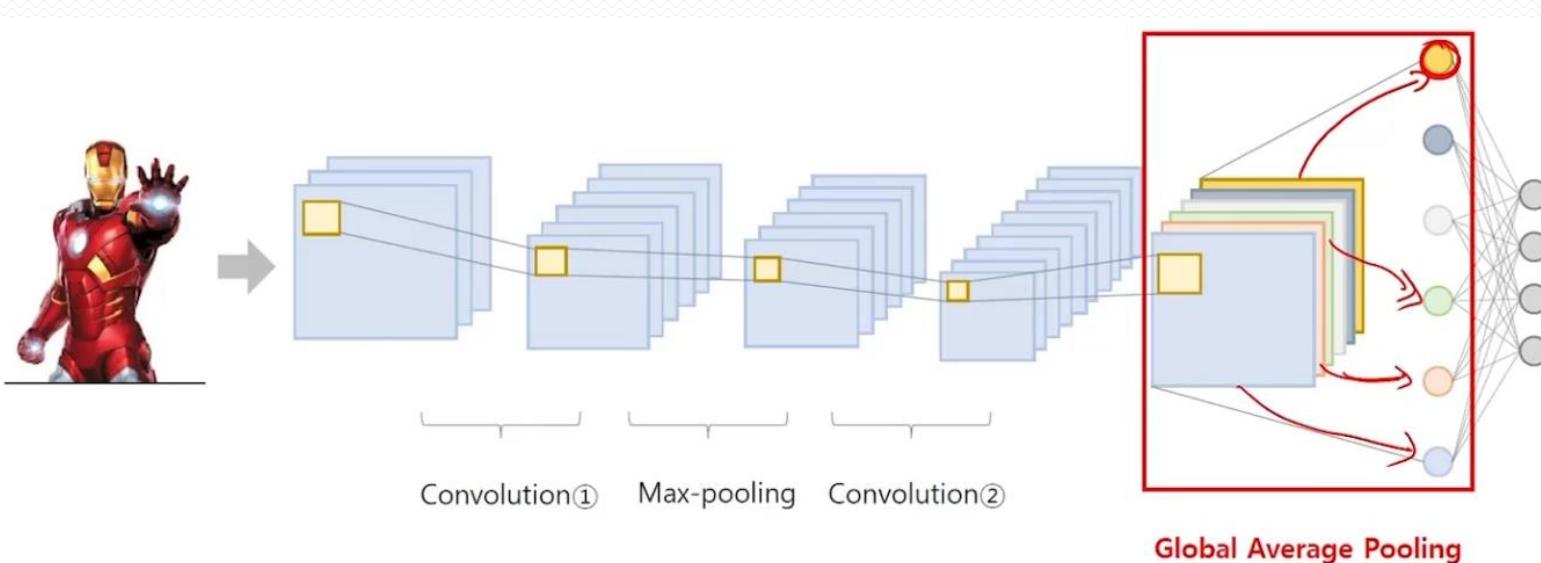
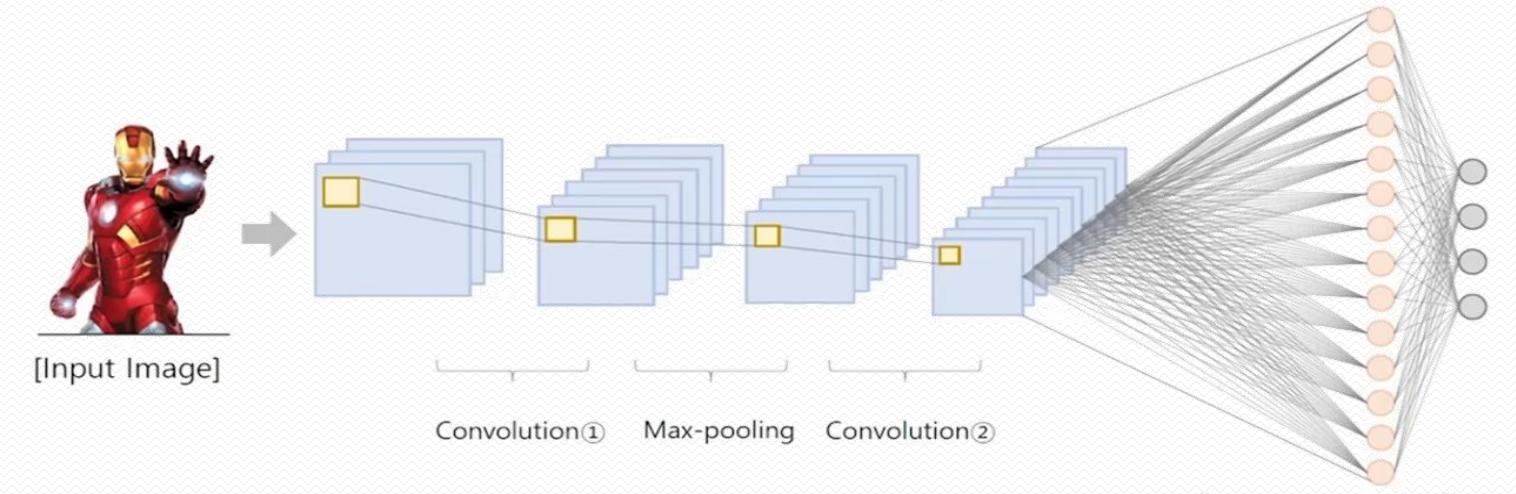
Class Activation Map

CNN Class Activation Map



CNN Class Activation Map

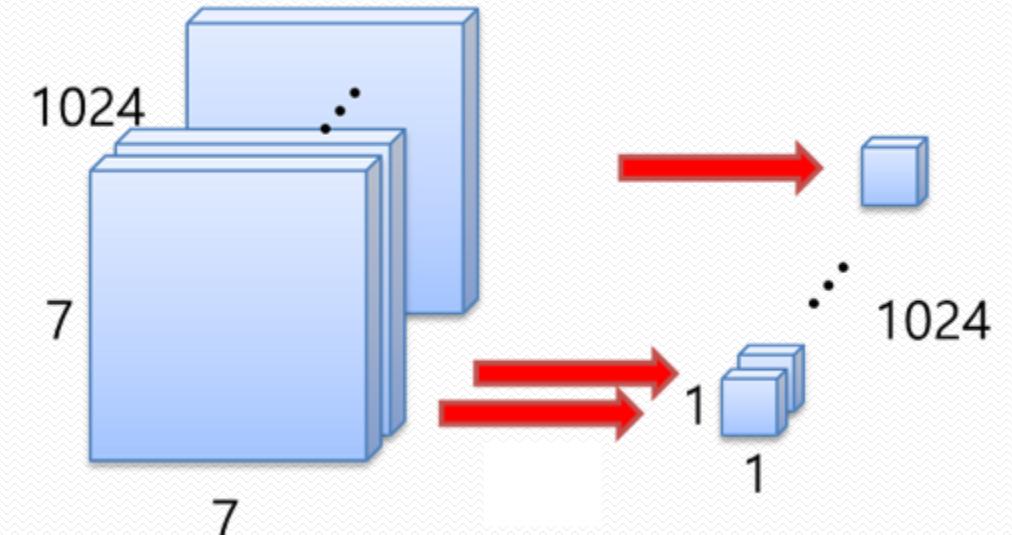
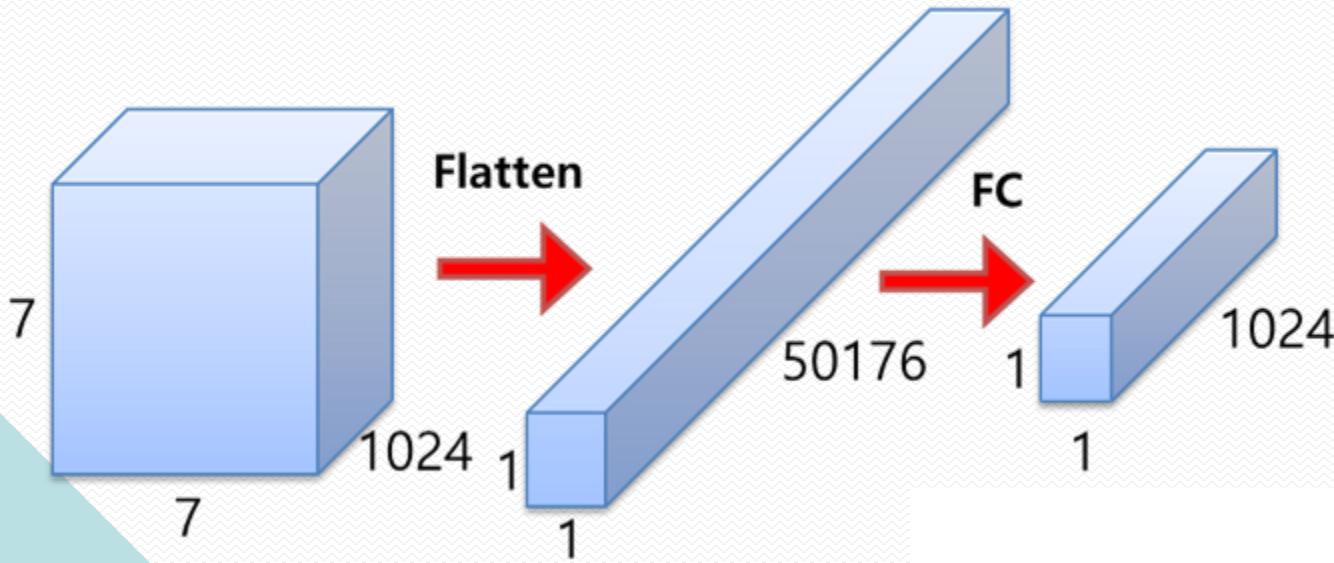
✓ Fully Connected Layer 대신 Global Average Pooling (GAP) 이라는 개념을 도입



Global Average Pooling

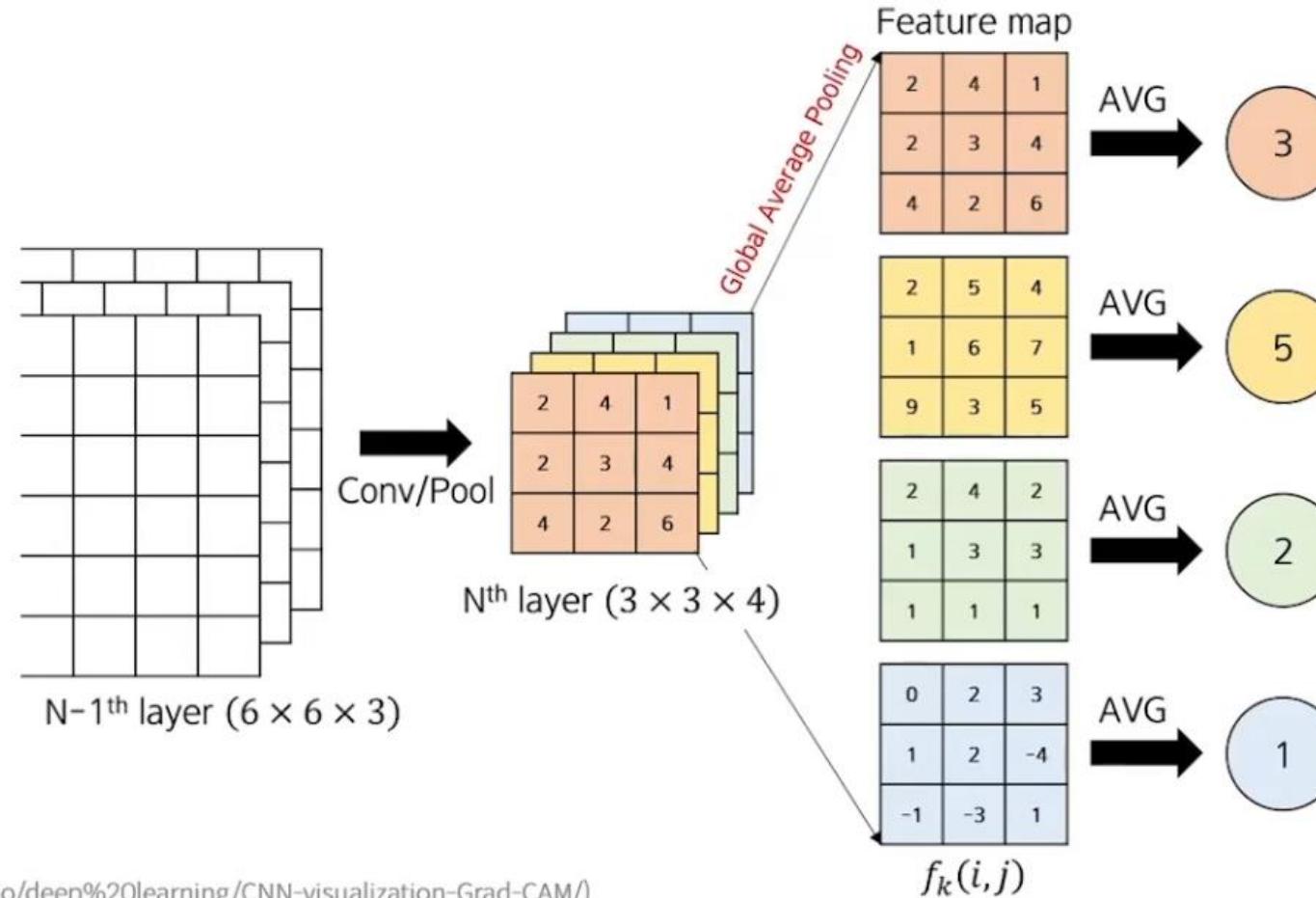
Global Average Pooling

- FC (Fully Connected)
 - Weights: $7 \times 7 \times 1024 \times 1024 = 51.3M$
- GAP (Global Average Pooling)
 - Weights: 0



CNN Class Activation Map

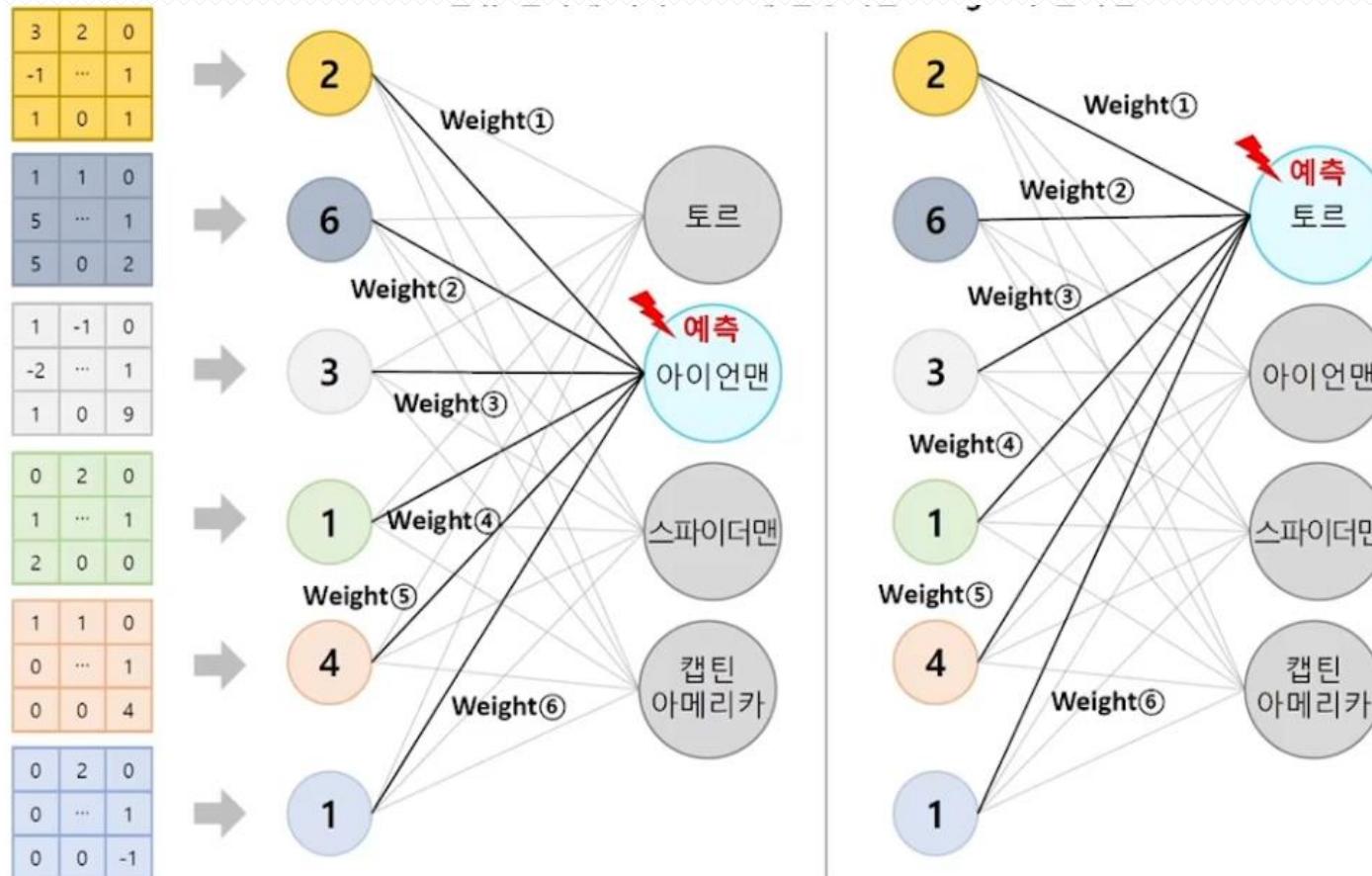
✓ Global Average Pooling : 각 Feature Map (채널)의 가중치 값들의 평균



(출처: <https://tyami.github.io/deep%20learning/CNN-visualization-Grad-CAM/>)

CNN Class Activation Map

- ✓ 각 가중치의 값이 클수록 해당 Feature map의 기여도가 크다.
- ✓ 분류 결과에 따라서 CAM에 활용되는 가중치가 달라짐



CNN Class Activation Map

3	2	0
-1	...	1
1	0	1

1	1	0
5	...	1
5	0	2

1	-1	0
-2	...	1
1	0	9

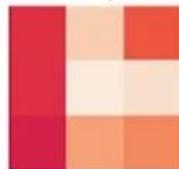
0	2	0
1	...	1
2	0	0

1	1	0
0	...	1
0	0	4

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

각 Feature map 별 heatmap을 그림

$x \text{ Weight}① =$



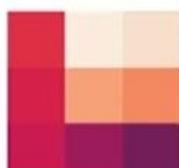
$x \text{ Weight}② =$



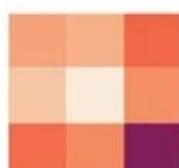
$x \text{ Weight}③ =$



$x \text{ Weight}④ =$



$x \text{ Weight}⑤ =$



$x \text{ Weight}⑥ =$



→ 동일 위치 Pixel별 합



→



아이언맨 얼굴을
예측 원인으로 판단

CNN Class Activation Map

- ✓ 마지막 Convolution Feature Map에 가장 많은 정보가 축약되어 있음
- ✓ 마지막 Feature Map 내의 1개 값은 원본 이미지에서 많은 부분이 요약된 결과물임

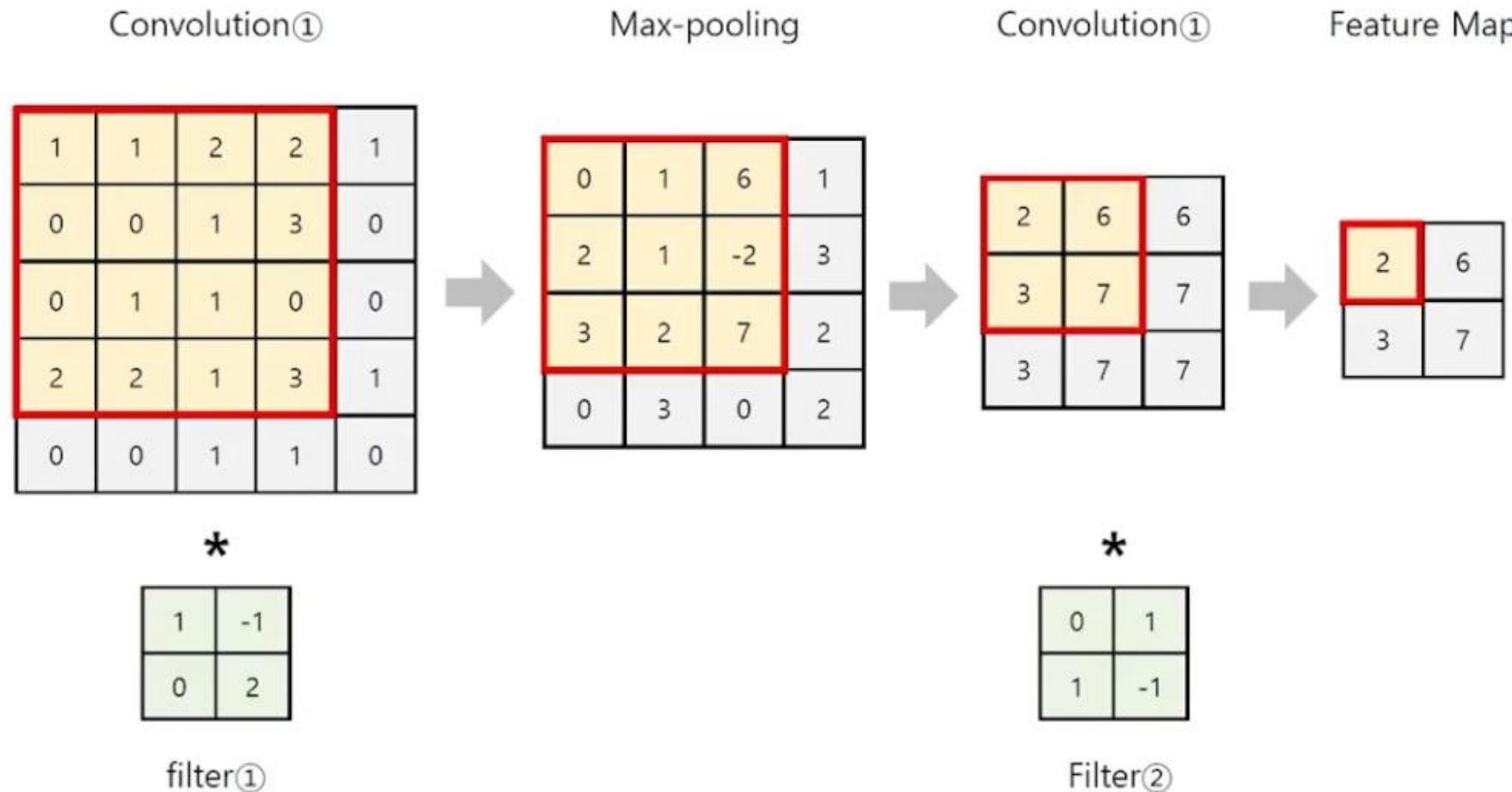


Image Localization: CNN Class Activation Map

✓ 예측된 범주 정보를 바탕으로 상대적으로 강하게 활성화가 되는 입력 이미지의 영역을 추정하는 것

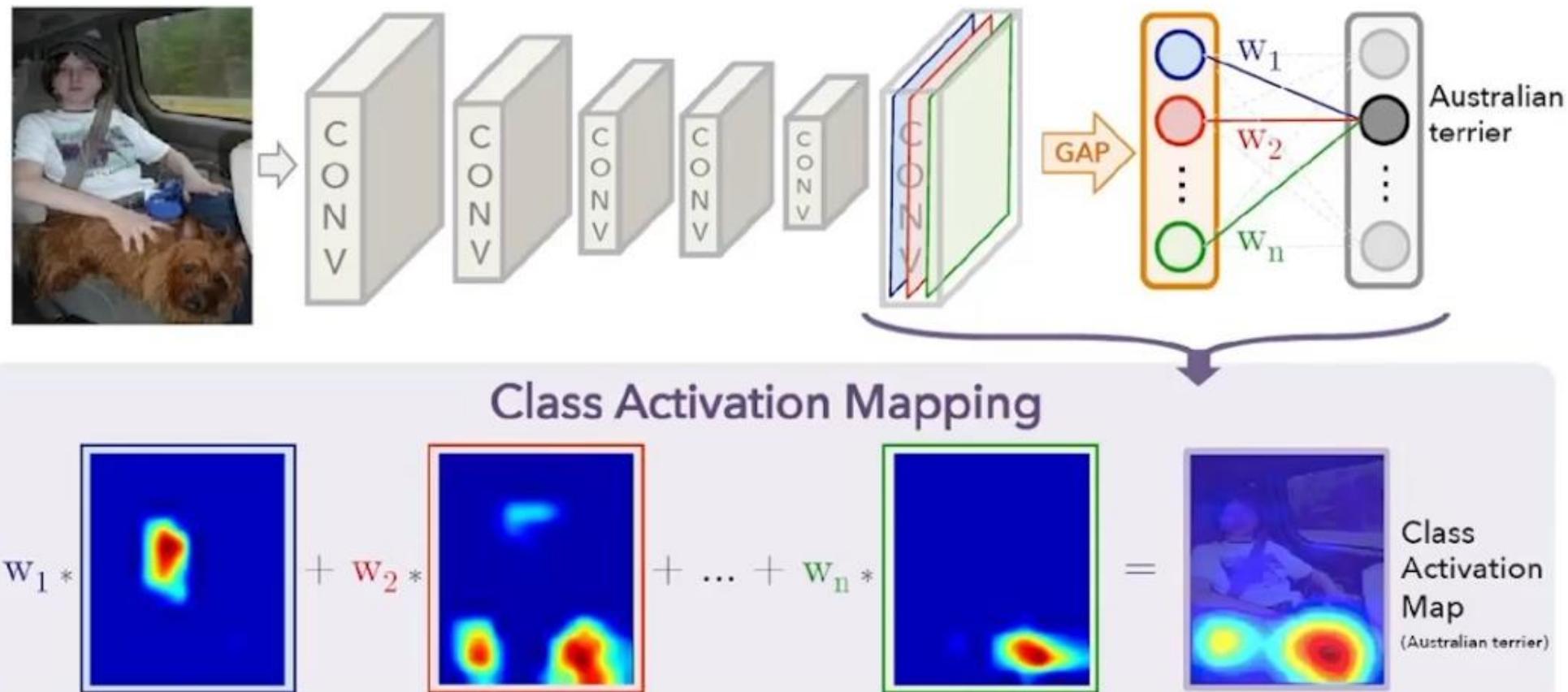
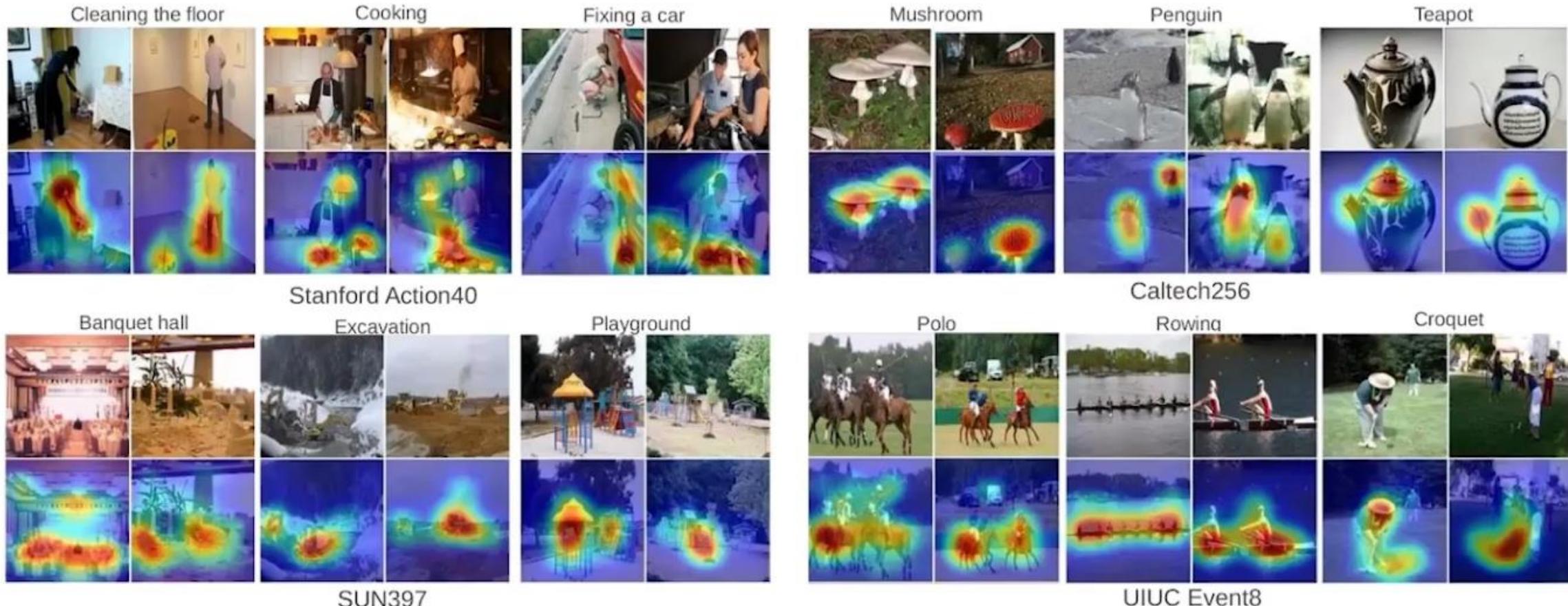
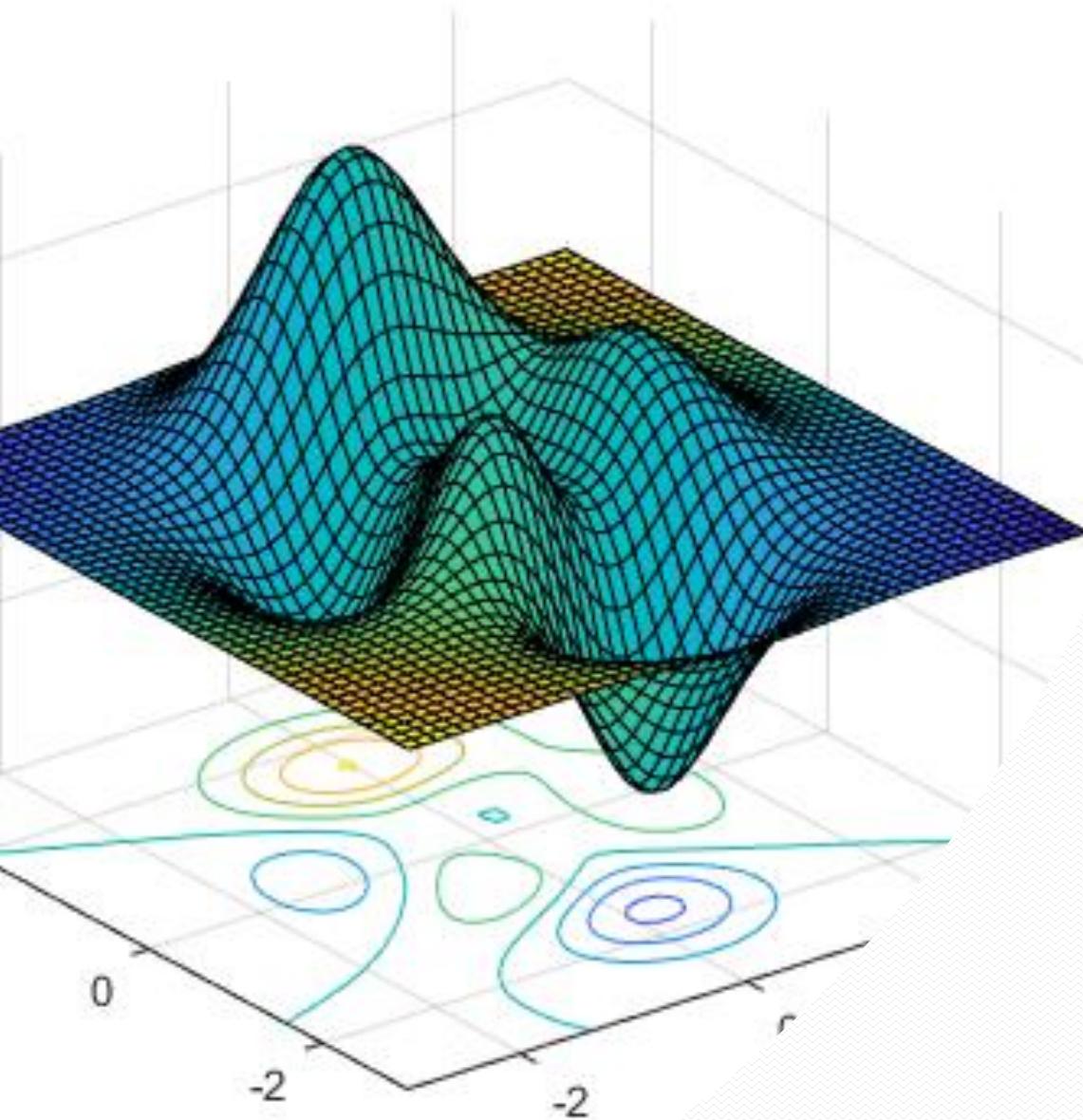


Image Localization: CNN Class Activation Map

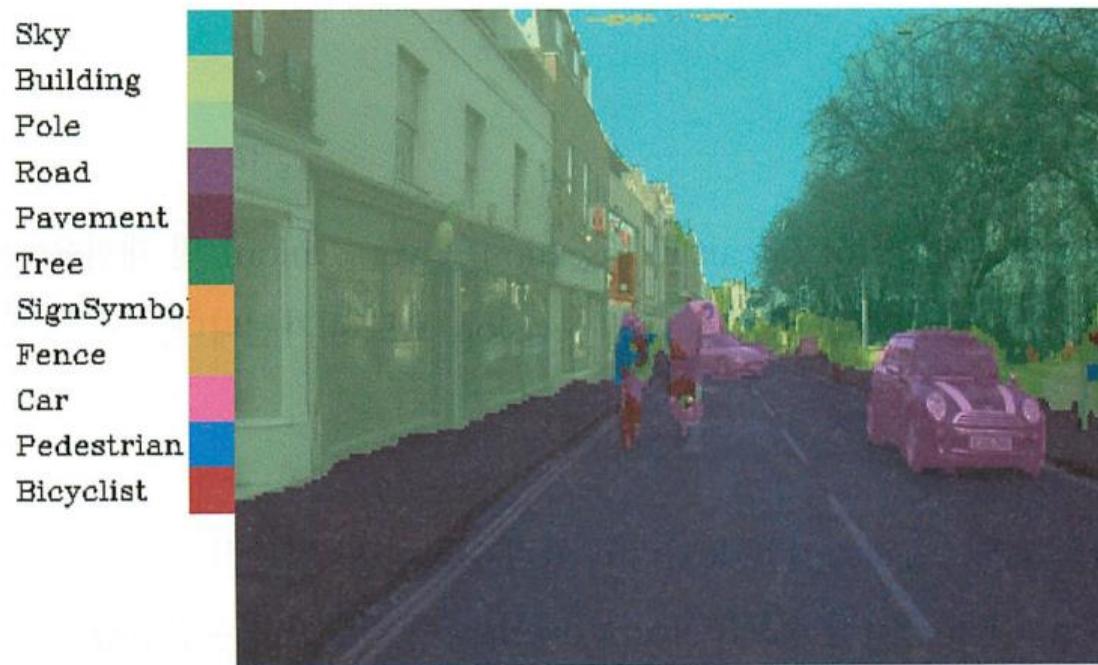




Segmentation

Segmentation

- ✓ 이미지의 각 픽셀은 범주에 따라 분류되므로 픽셀 수준 예측.
 - ✓ 덴스 예측(dense prediction)이라고 부르기도 함.



입력

segmented

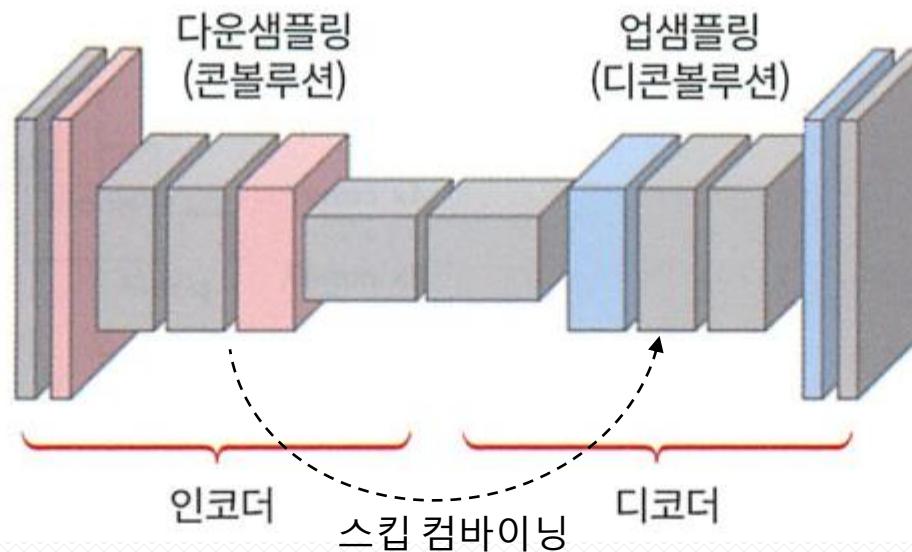
1: perso
2: grass
3: sky

출력

Segmentationmap

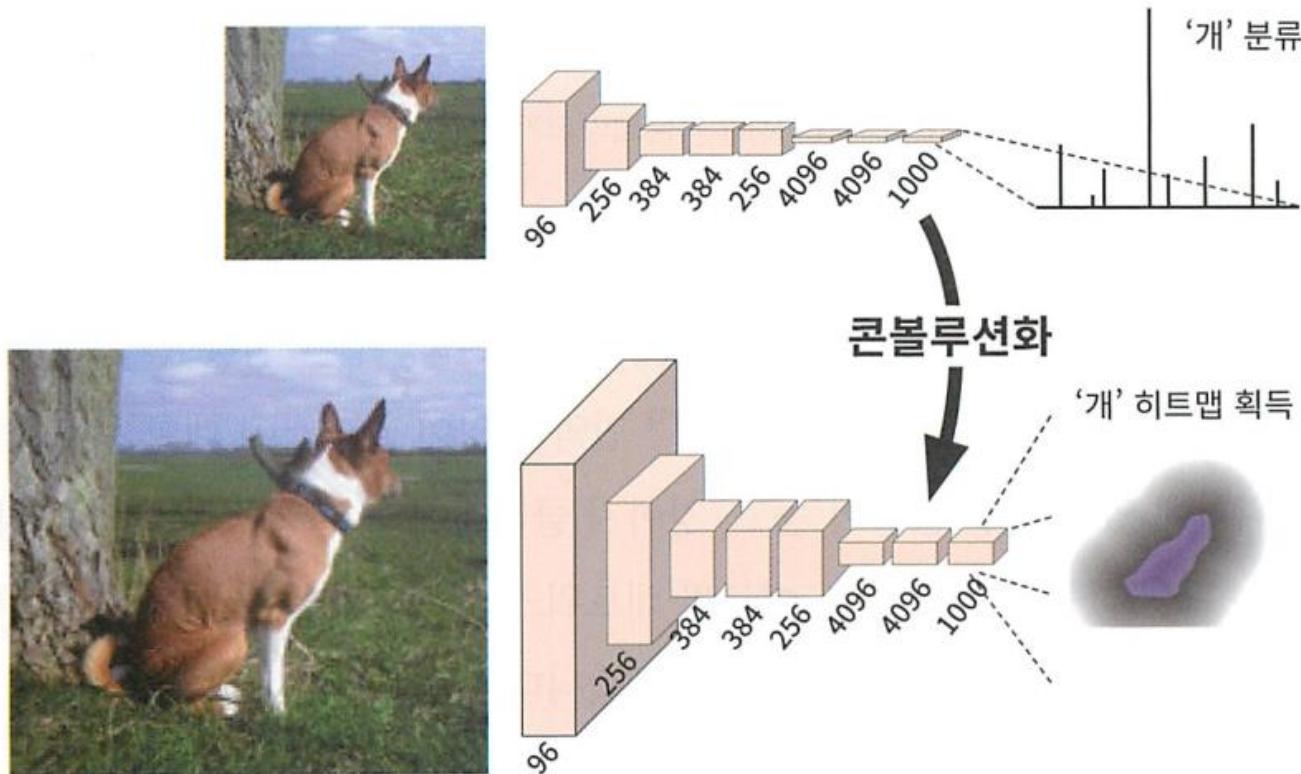
Segmentation

- ✓ 완전한 콘볼루션화(fully convolutionalization)
- ✓ 디콘볼루션(deconvolution)
- ✓ 스킵 아키텍처(skip architecture)



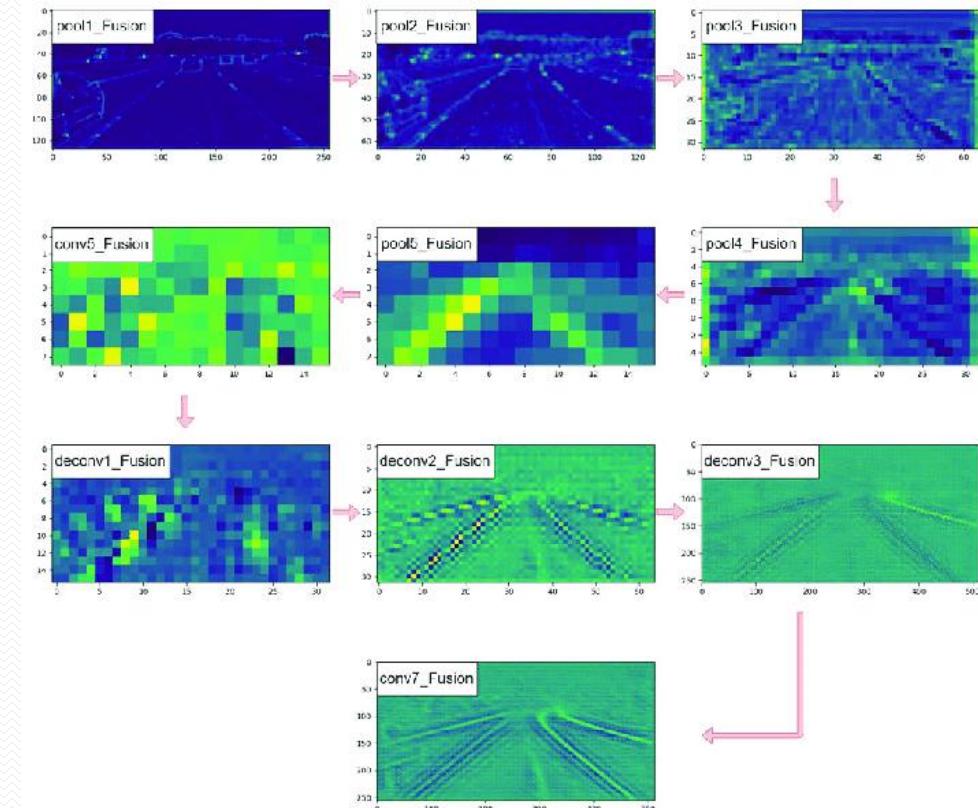
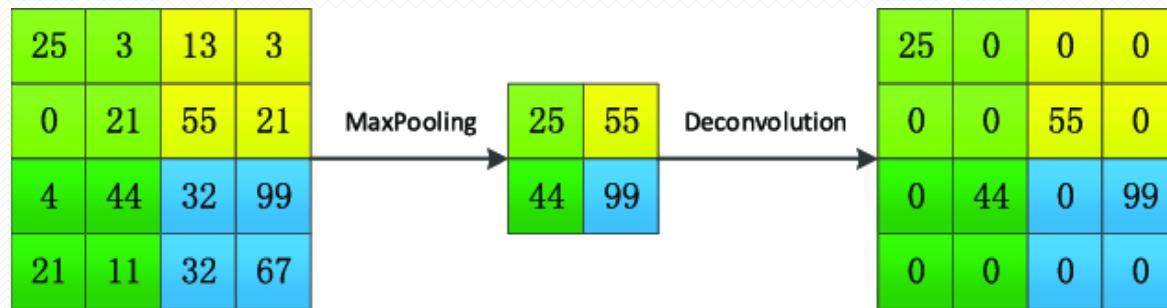
Fully convolutionalization

- ✓ FC 계층을 모두 Conv 계층으로 대체해 히트맵을 획득
- ✓ 마지막 특징 맵의 개수는 훈련된 클래스의 개수와 동일



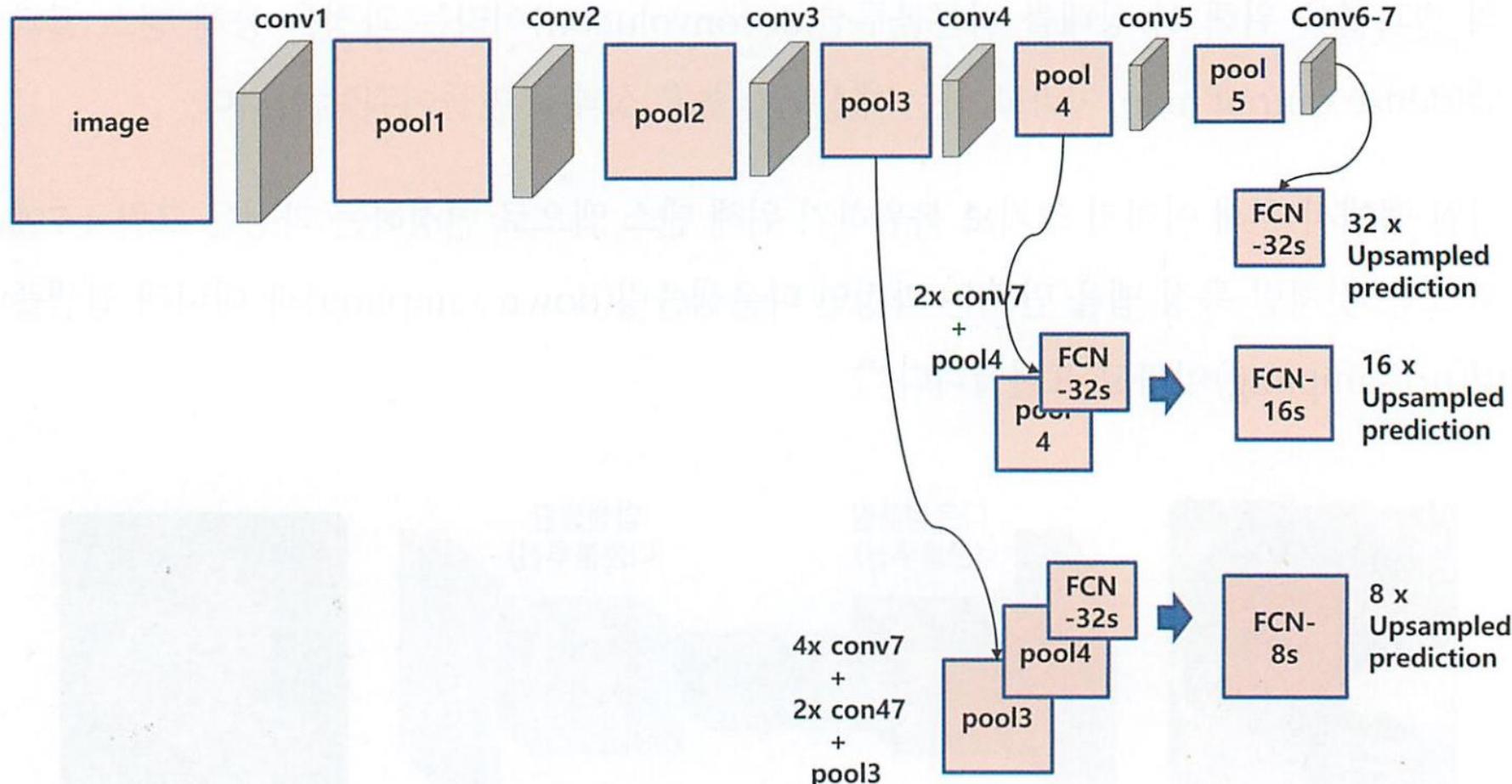
Deconvolution

- ✓ 특징 맵을 이용한 이미지 분할의 결과물은 이미지 원본 크기에서 픽셀 단위 예측에 사용하기에는 해상도가 너무나 부족
- ✓ 디콘볼루션(deconvolution)이라는 과정을 통해 덴스 결과 맵(dense output map)이라고 하는 해상도가 높은 상태로 만듬



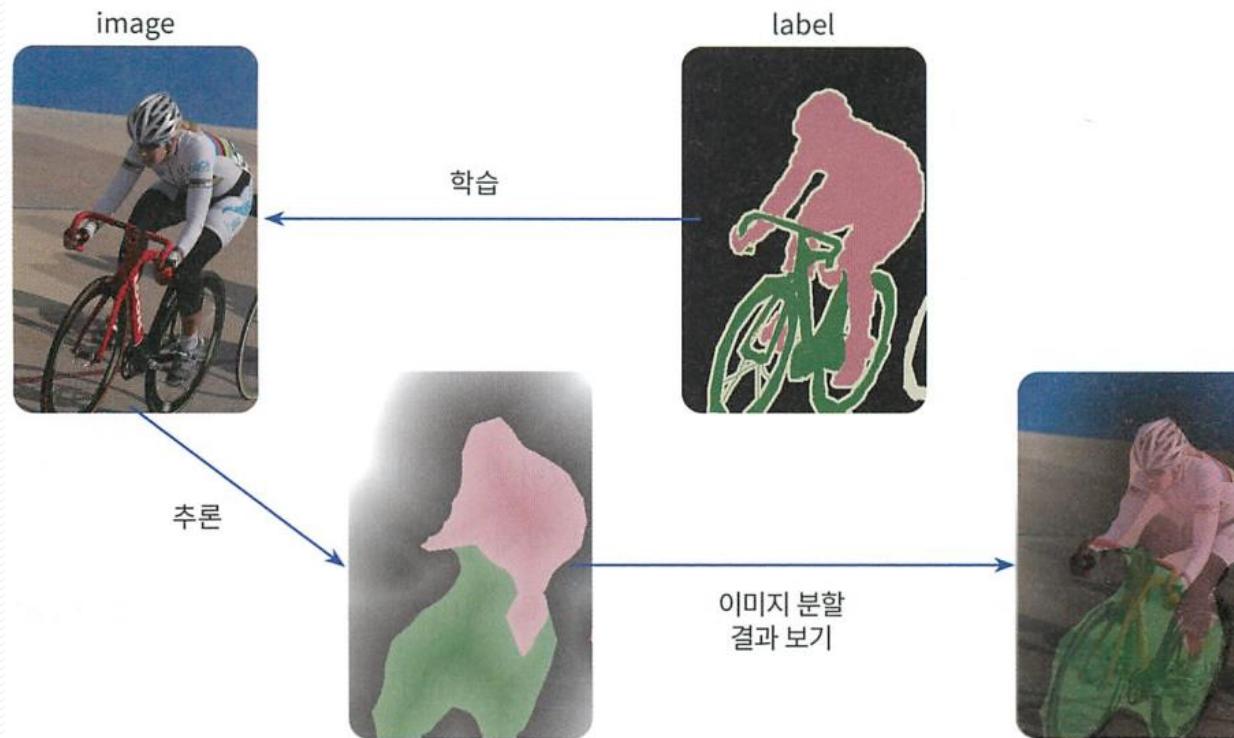
스킵 구조

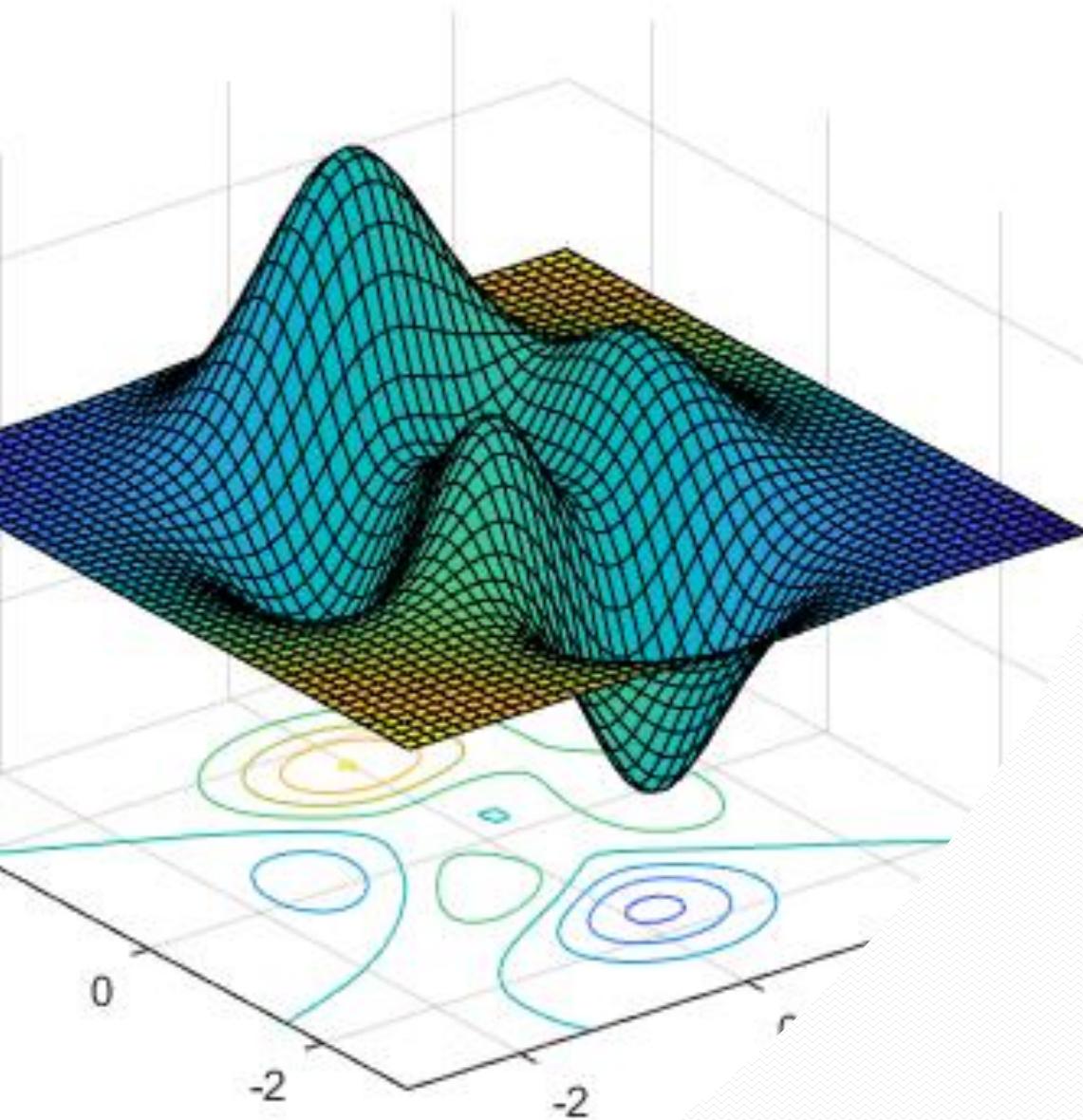
- ✓ 이전 콘볼루션 층의 특징 맵을 참고해 업샘플링해 줄 때 정확도를 높이는 방법
- ✓ 정확도를 높일 수 있는 이유는 이전 콘볼루션 층의 특징 맵들이 해상도 면에서는 더 좋기 때문



Segmentation

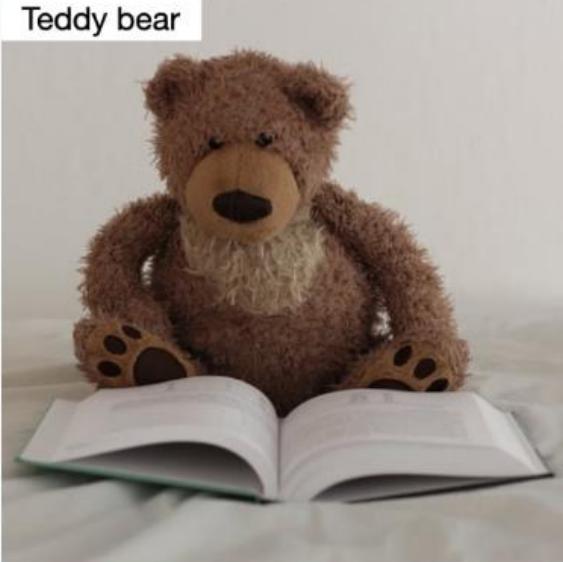
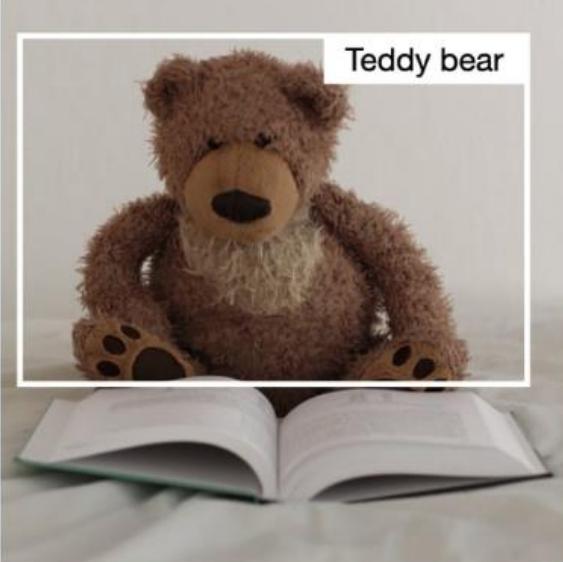
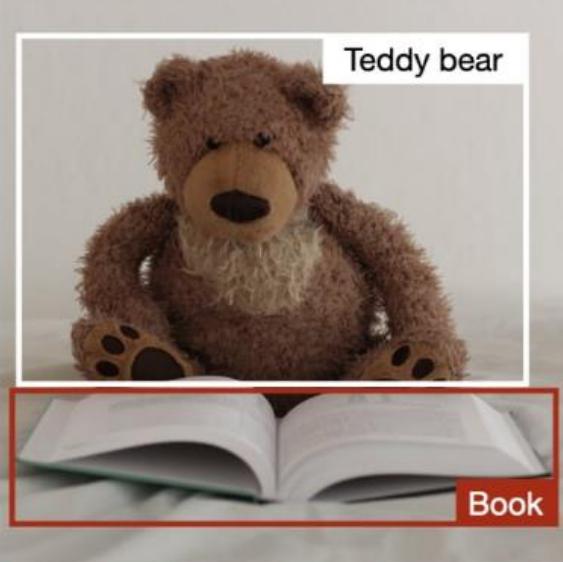
- ✓ 콘볼루션 과정의 입력은 사이즈 제약이 필요 없음.
- ✓ 콘볼루션의 결과인 특징 맵에는 사물의 위치 정보가 들어 있음.
- ✓ 최종 특징 맵뿐 아니라 중간 레이어의 특징 맵도 유용하게 이용할 수 있음.



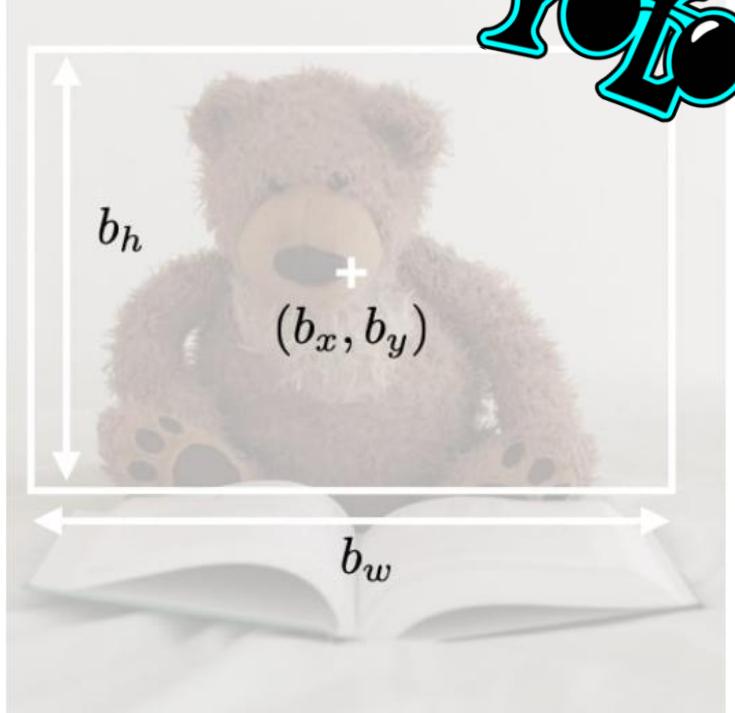
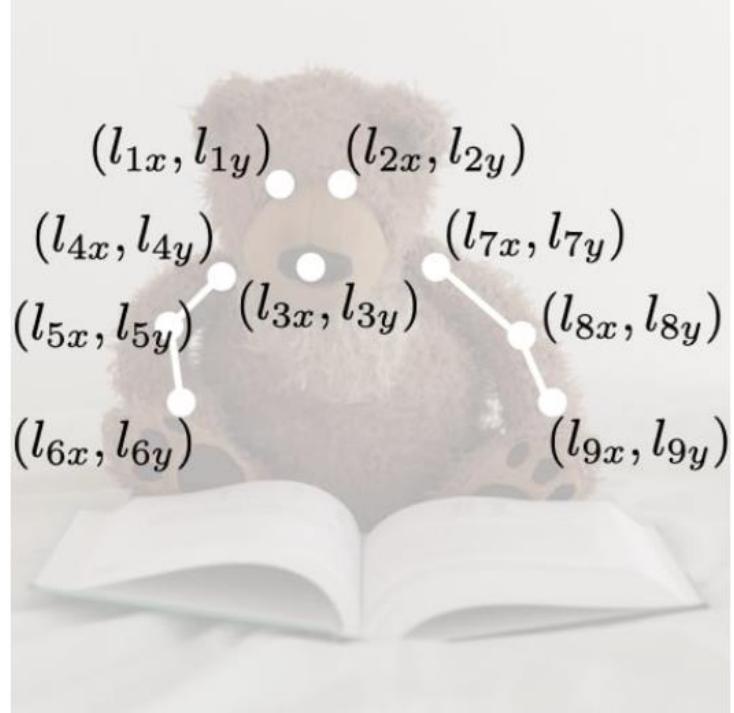


Object Detection

Object Detection

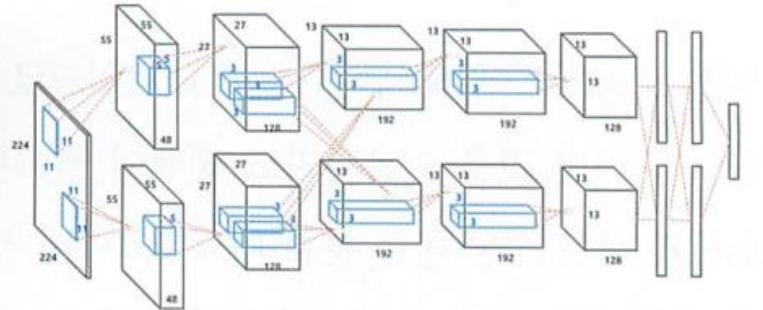
Image classification	Classification w. localization	Detection
<p>Teddy bear</p> 	<p>Teddy bear</p> 	<p>Teddy bear</p> 
<ul style="list-style-type: none">• Classifies a picture• Predicts probability of object	<ul style="list-style-type: none">• Detects an object in a picture• Predicts probability of object and where it is located	<ul style="list-style-type: none">• Detects up to several objects in a picture• Predicts probabilities of objects and where they are located
Traditional CNN	Simplified YOLO, R-CNN	YOLO R-CNN

Object Detection (cont.)

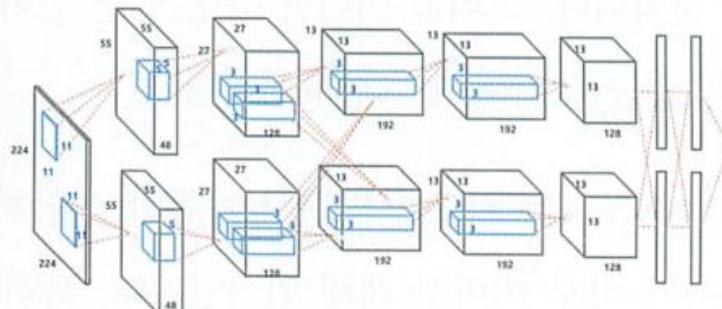
Bounding box detection	Landmark detection
<ul style="list-style-type: none">Detects the part of the image where the object is located 	<ul style="list-style-type: none">Detects a shape or characteristics of an object (e.g. eyes)More granular 
Box of center (b_x, b_y) , height b_h and width b_w	Reference points $(l_{1x}, l_{1y}), \dots, (l_{nx}, l_{ny})$

Object Detection

- ✓ 이미지에서 특정 대상을 검출하고자 할 때 멀티 레이블 이미지 분류와 바운딩 박스 회귀를 위한 기본적인 영역을 정하기 어려운 문제가 있음



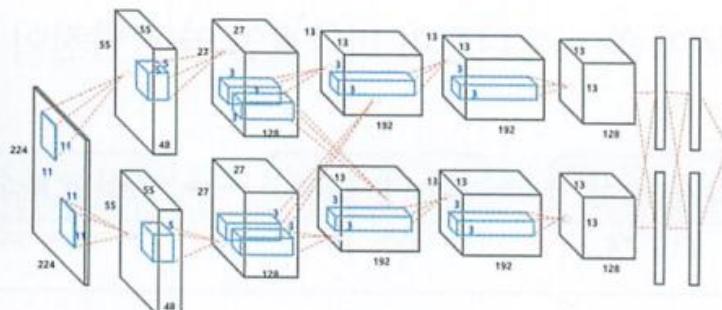
CAT



DOG

DOG

DOG (x, y, w, h)



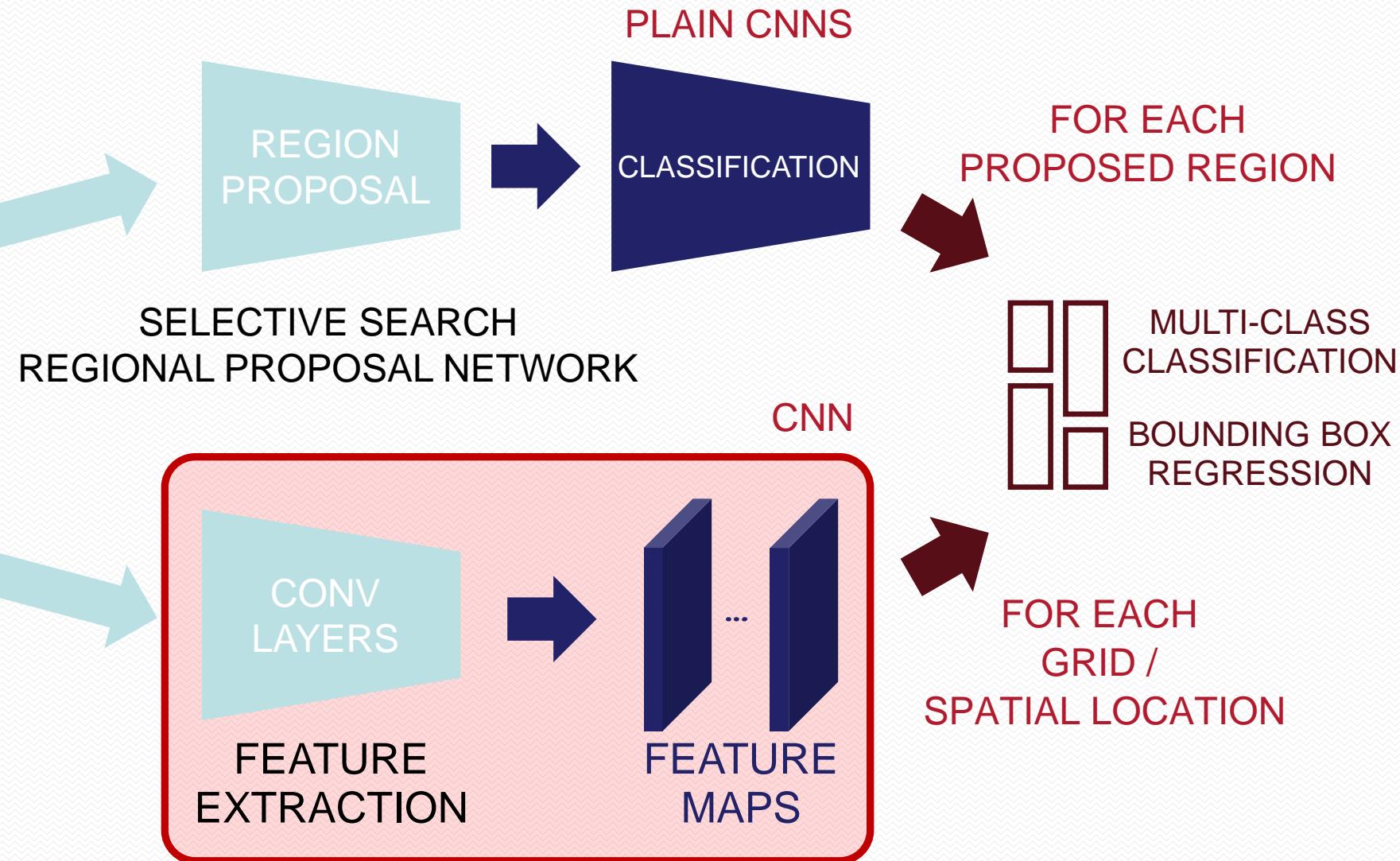
BUS (x, y, w, h)

BUS (x, y, w, h)

OBJECT DETECTION



2-STAGE



1-Stage & 2-Stage Detector

2-Stage



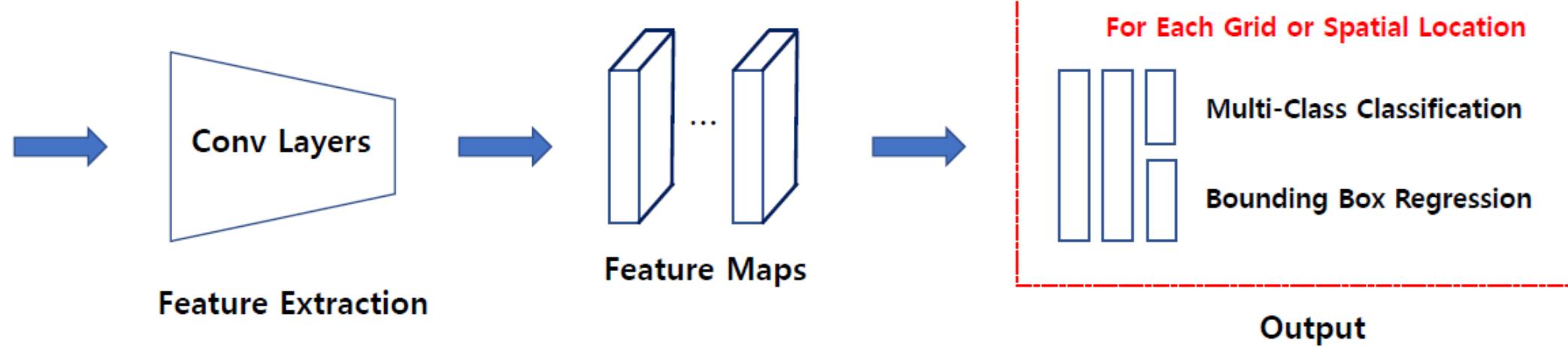
Input Image



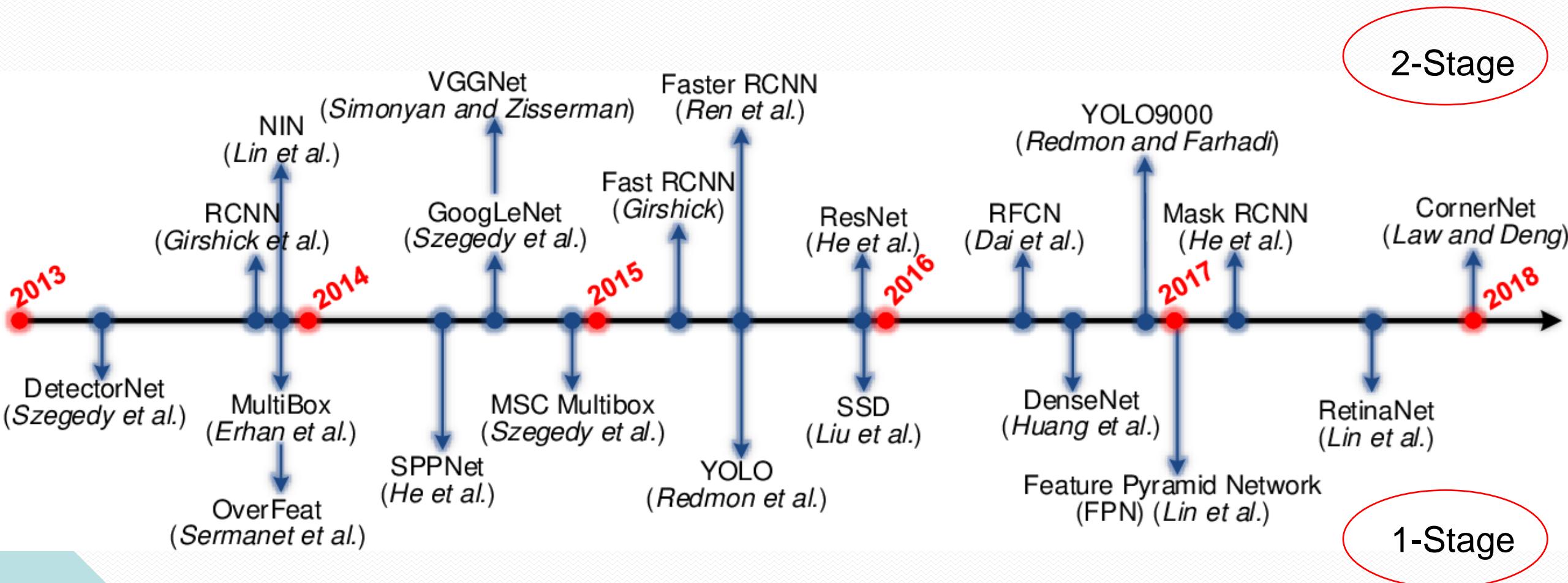
1-Stage

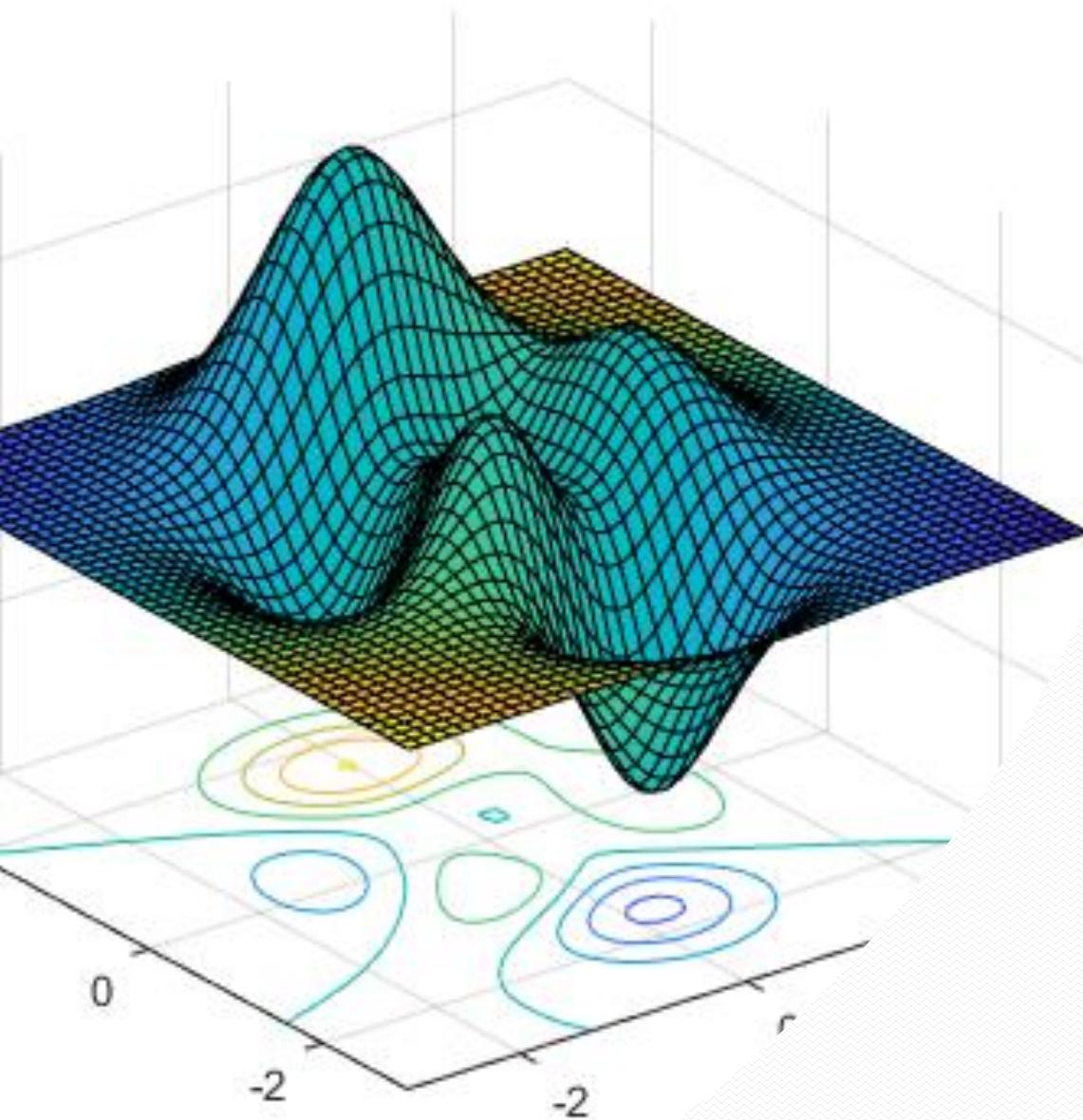


Input Image



1-Stage & 2-stage Detector (cont.)



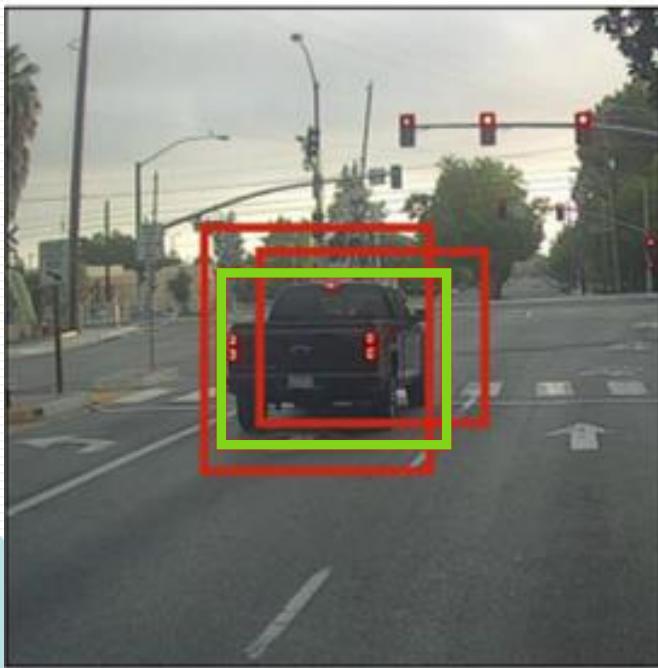


2 stage

NMS(Non-Maximum Suppression)

- If there are multiple boxes on the same object, only the box with the highest score is left

Before non-max suppression



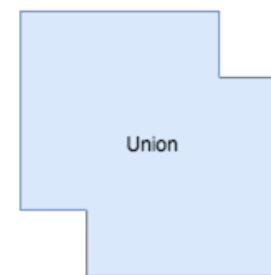
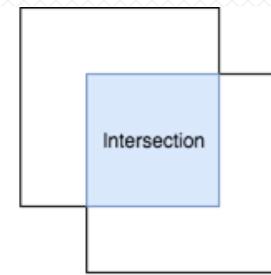
Non-Max
Suppression



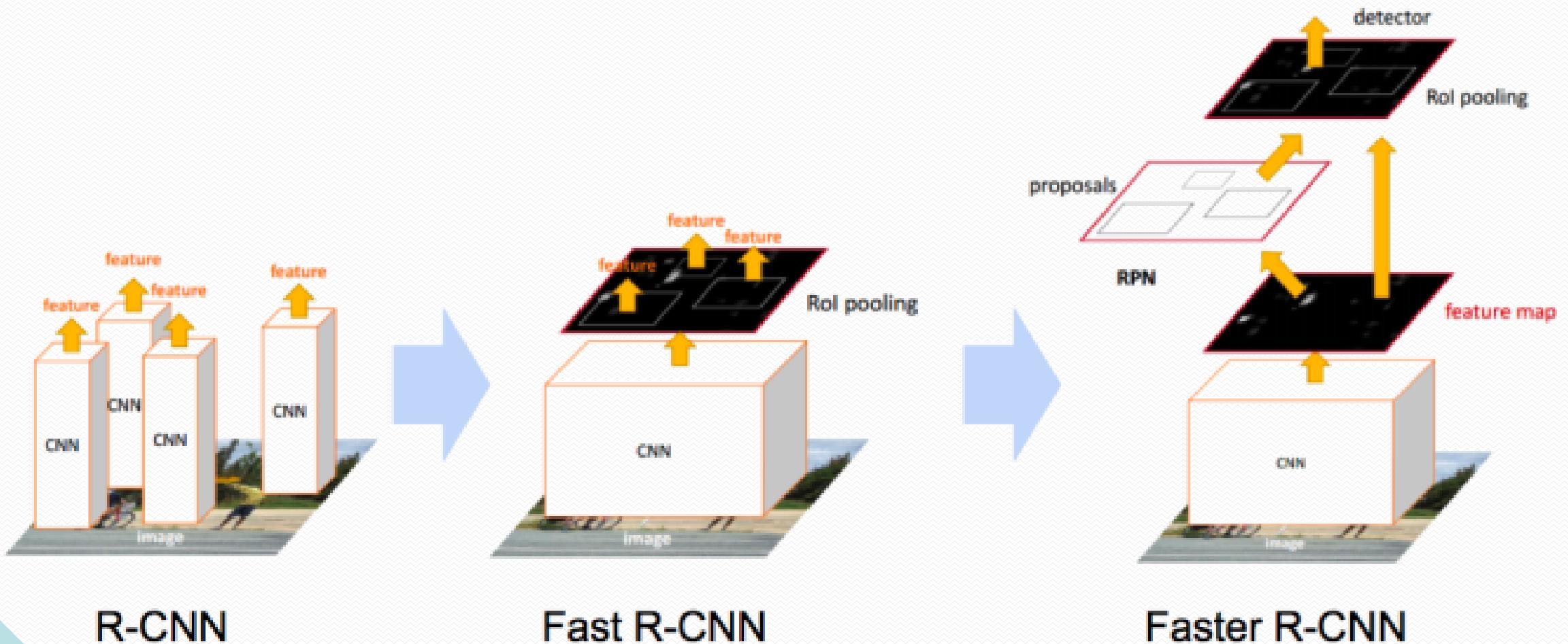
After non-max suppression



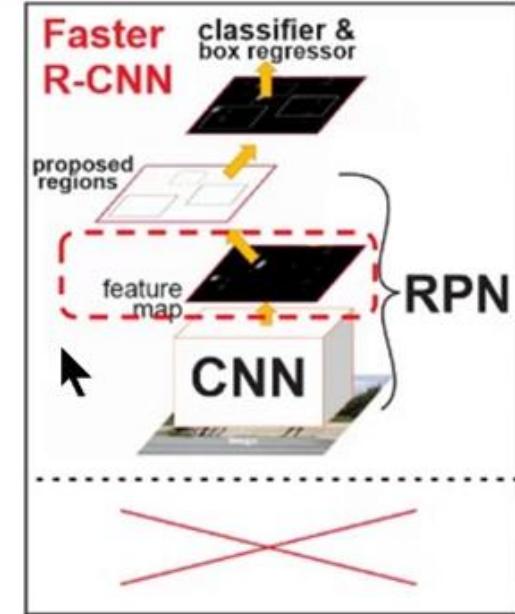
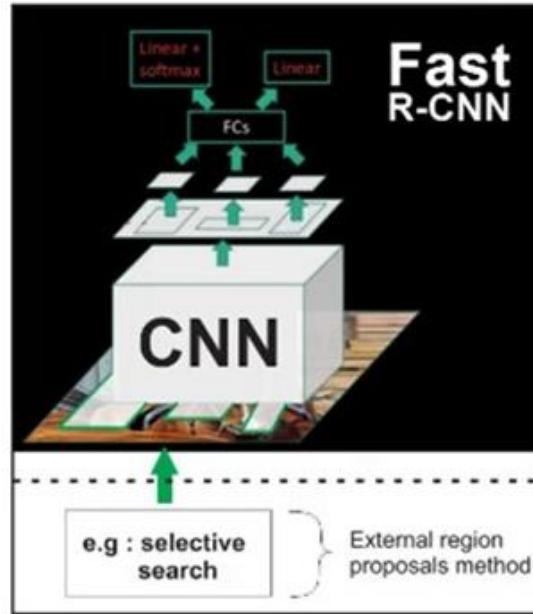
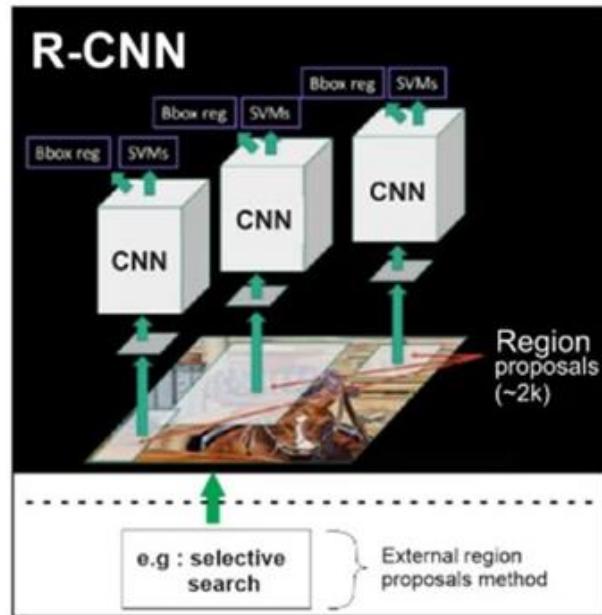
$$\text{IoU} = \frac{\text{Intersection}}{\text{Union}}$$



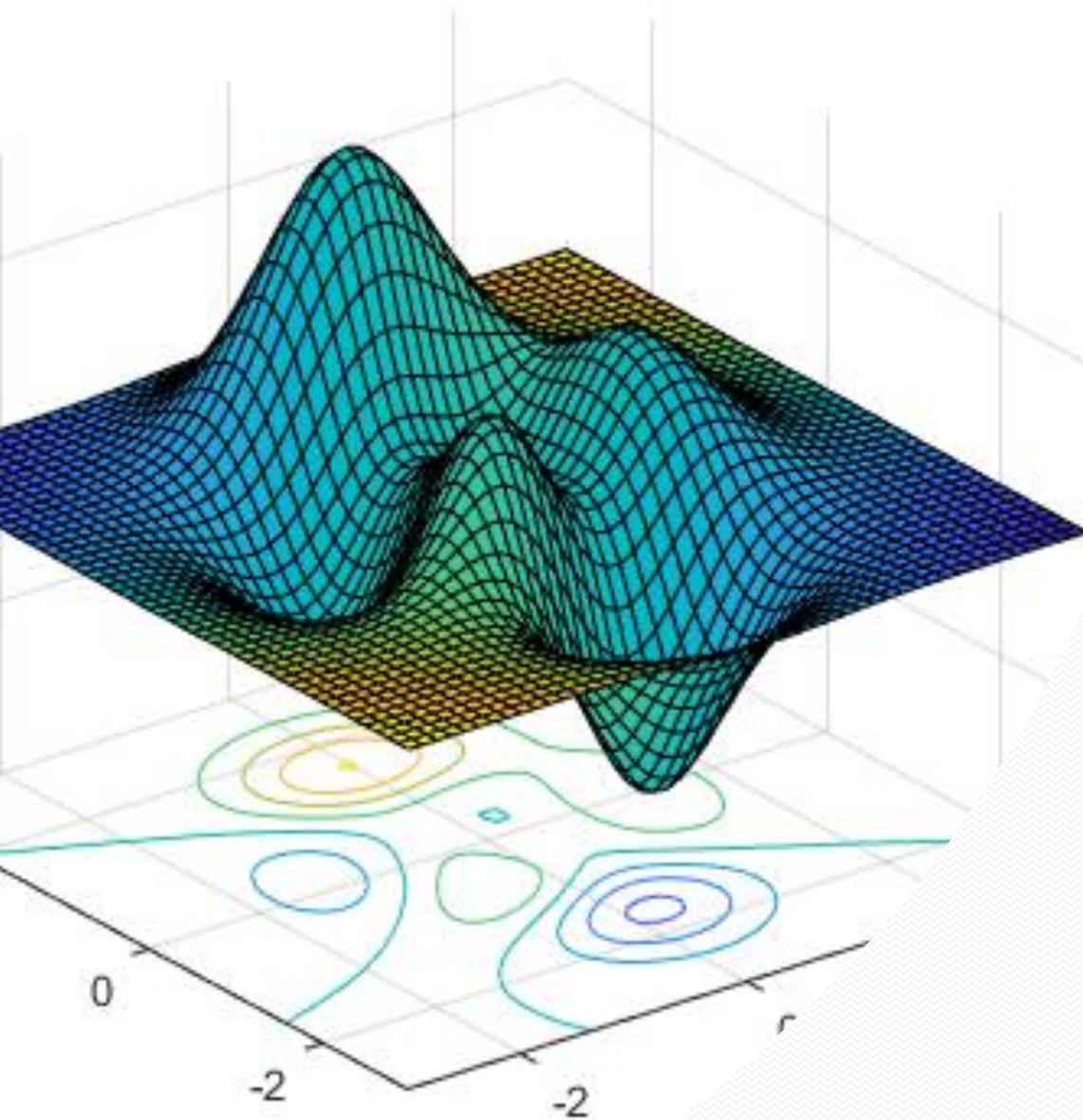
Evolution from R-CNN to Faster R-CNN



Compare performance

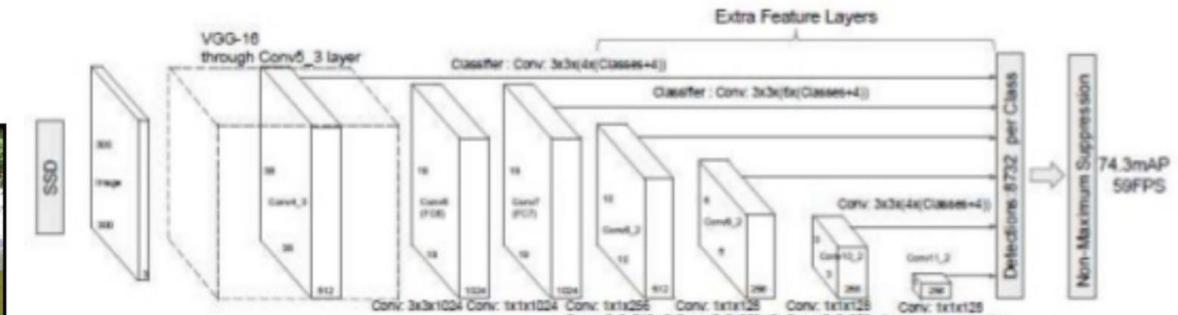
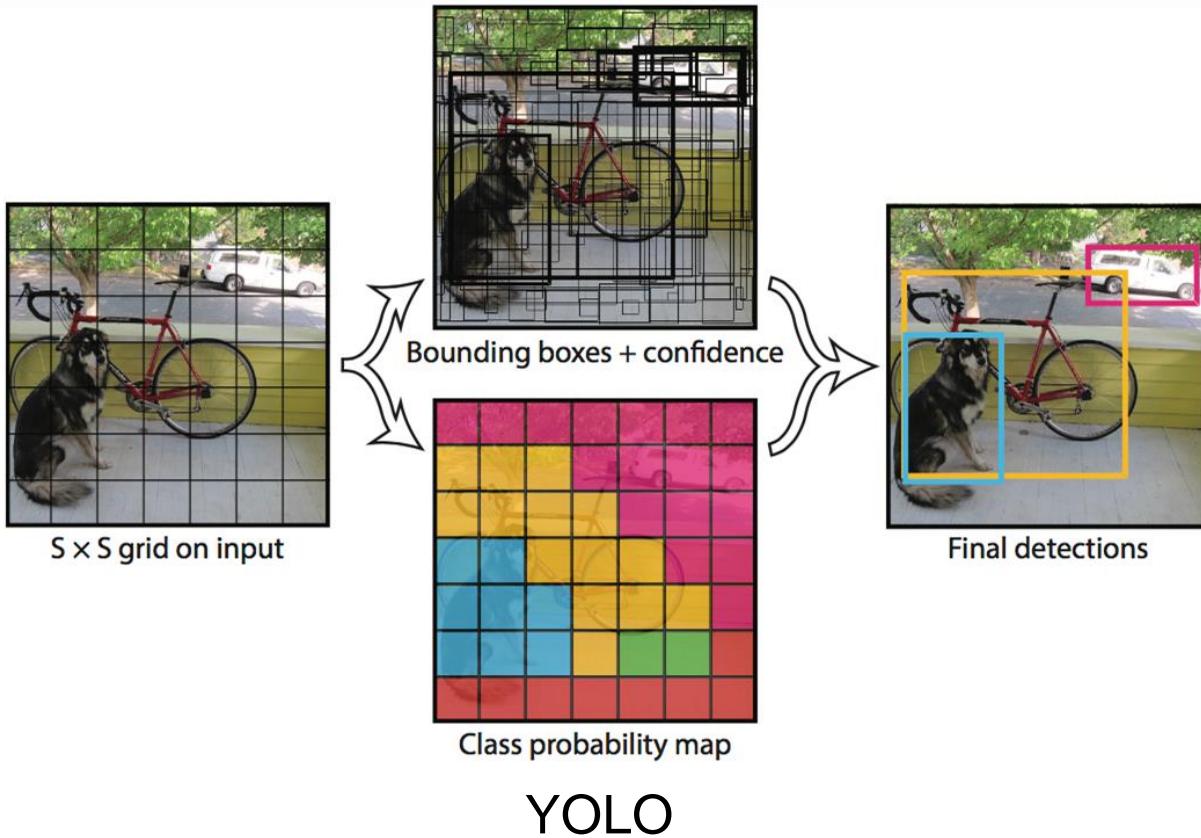


	R-CNN	Fast R-CNN	Faster R-CNN
Test time per image	50 seconds	2 seconds	0.2 seconds
Speed-up	1x	25x	250x
mAP (VOC 2007)	66.0%	66.9%	66.9%



1 stage

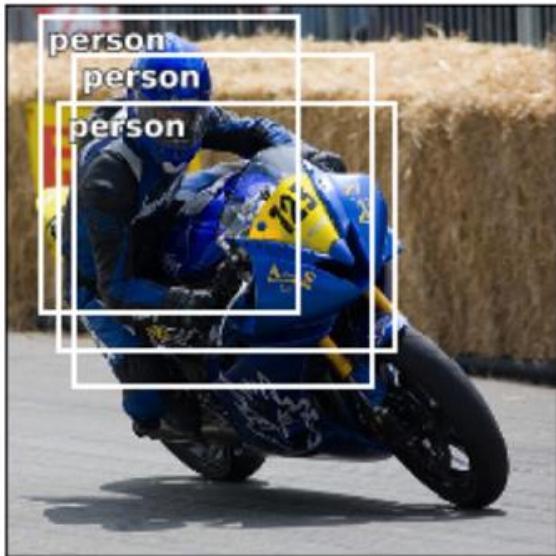
1 STAGE OBJECT DETECTION



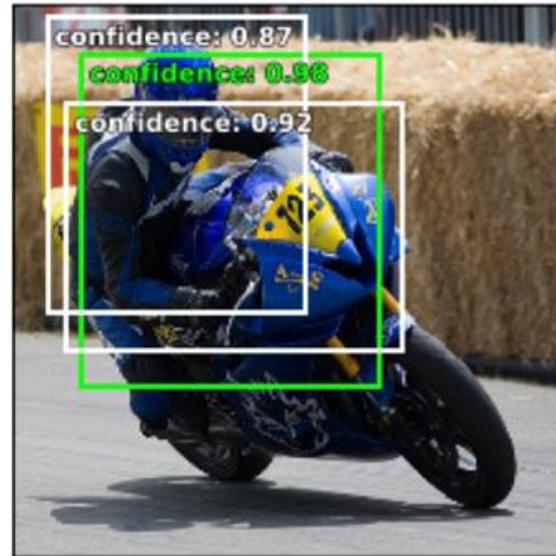
Predictions on a grid (cont.)

Repeat with next highest confidence prediction until no more boxes are being suppressed

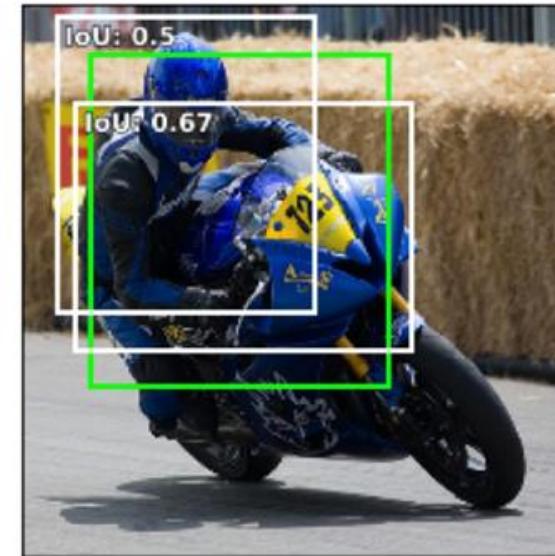
For each class...



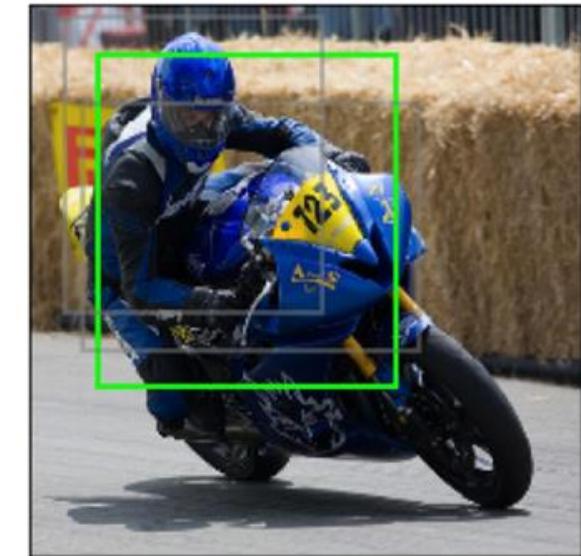
After filtering out low confidence predictions, we may still be left with **redundant detections**



Select the bounding box prediction with the **highest confidence**



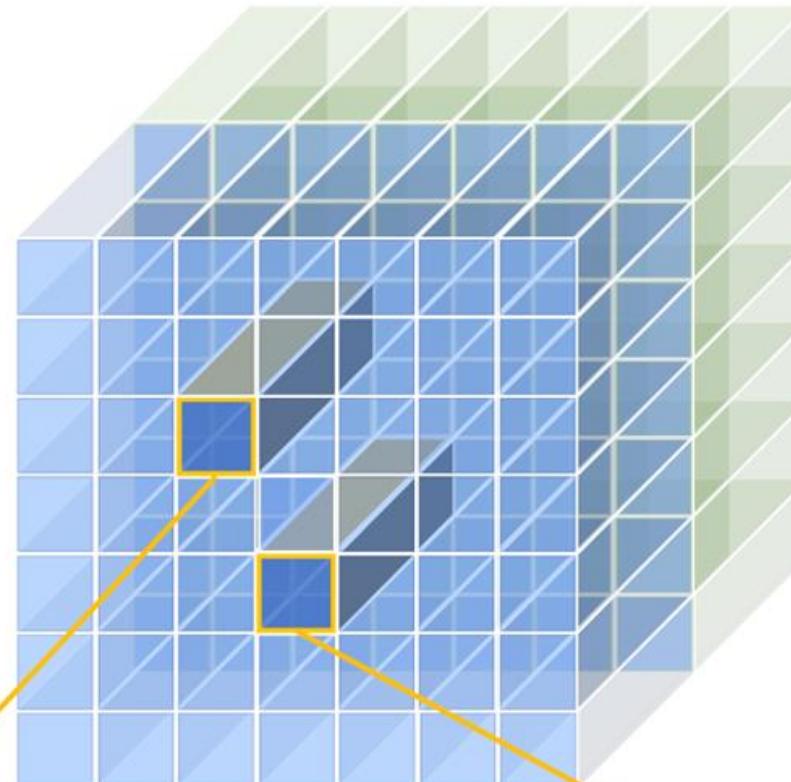
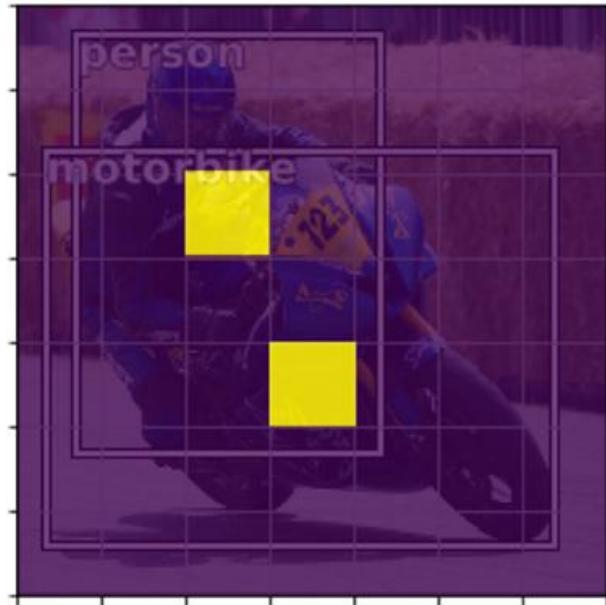
Calculate the IoU between the **selected box** and all remaining predictions



Remove any boxes which have an IoU score above some defined threshold

Predictions on a grid (cont.)

multiple objects can be detected in parallel



t_x	t_y	t_w	t_h	c_1	c_2	c_3	c_4	p_{obj}
-------	-------	-------	-------	-------	-------	-------	-------	-----------

Person bounding box descriptor

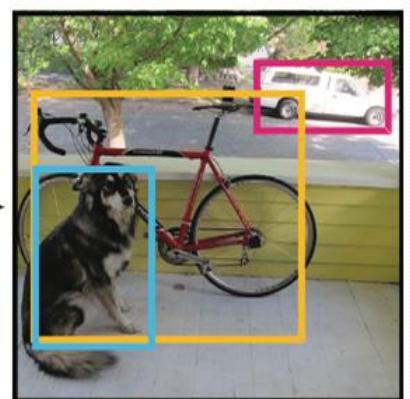
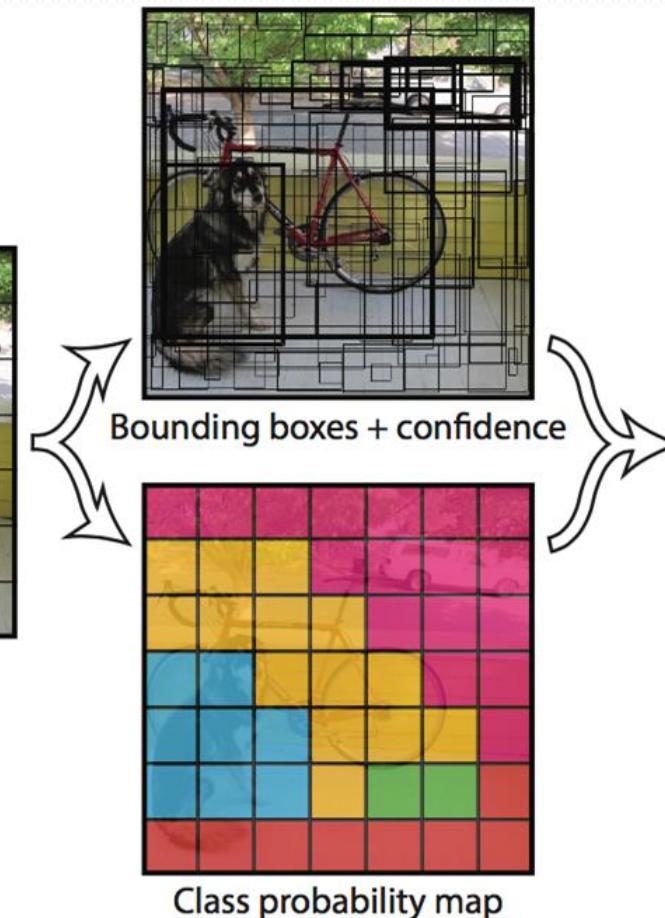
t_x	t_y	t_w	t_h	c_1	c_2	c_3	c_4	p_{obj}
-------	-------	-------	-------	-------	-------	-------	-------	-----------

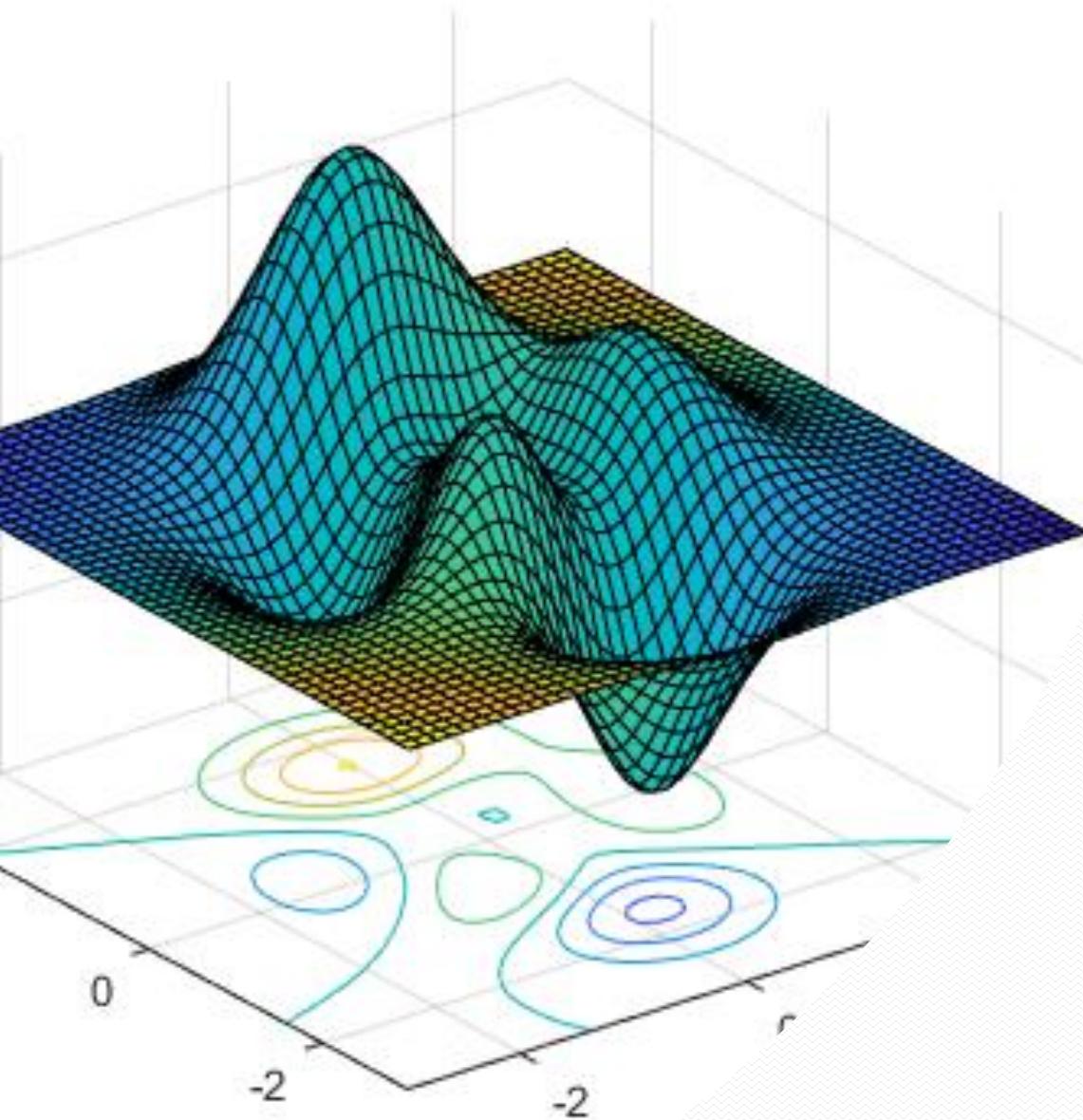
Motorbike bounding box descriptor

Yolo

- Predict one set of class probabilities per grid cell, regardless of the number of boxes B
- Image
 - **S** x **S** grids
- Grid cell
 - **B**: Bboxes & Confidence score
 - **C**: class probabilities with regard to #classes

x, y, w, h, confidence





CONFUSION MATRIX

CONFUSION MATRIX

QUANTIFIABLE MEASURE THAT IS USED TO TRACK AND ASSESS THE STATUS OF A SPECIFIC PROCESS

Model classification & answer → TRUE / FALSE

True Positive / TP

False Positive / FP

False Negative / FN

True Negative / TN

		ACTUAL	
		POSITIVE	NEGATIVE
PREDICTED	POSITIVE	TRUE POSITIVE	FALSE POSITIVE
	NEGATIVE	FALSE NEGATIVE	TRUE NEGATIVE

PRECISION / RECALL

Precision

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}} = \frac{\text{TP}}{\# \text{ predictions}}$$

Proportion of what the actual true among the model classifies as true

Accuracy rate

Recall

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}} = \frac{\text{TP}}{\# \text{ ground truths}}$$

Proportion of what the model predicted to be true among the actual true

Detection rate

Precision
= 100%

90	0
10	

		ACTUAL	
		강아지	고양이
PREDICTED	강아지	90	0
	고양이	10	90

Recall =
90%

90	0
10	

IOU INTERSECTION OVER UNION

EVALUATION METRIC USED TO MEASURE THE ACCURACY OF AN OBJECT DETECTOR

MOSTLY USED IN OBJECT DETECTION

PASCAL VOC challenge

$$IoU = \frac{TP}{(TP + FP + FN)}$$

PREREQUISITES

The ground-truth bounding boxes = hand-labeled bounding boxes

The predicted bounding boxes from the model.

$$IoU = \frac{\text{Area of Overlap}}{\text{Area of Union}}$$





THANK YOU