

Deep Neural Network

(Convolutional Neural Network)

Convolutional Neural Network



- ❖ Hierarchical visual information processing of biological neural networks
 - David Hubble and Thorsten Niels Vissel published a study in 1959 on the structure and function of the visual cortex of animals

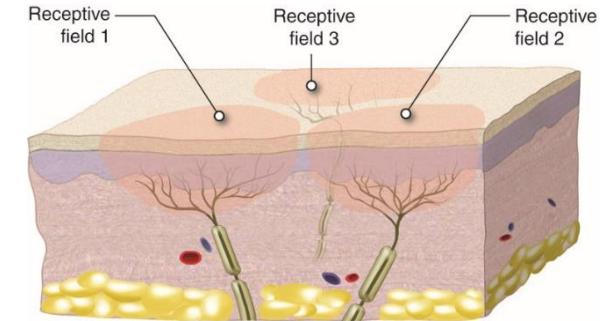
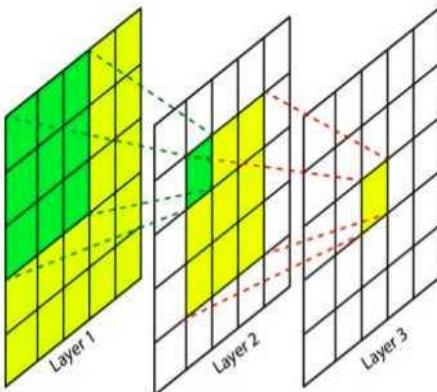
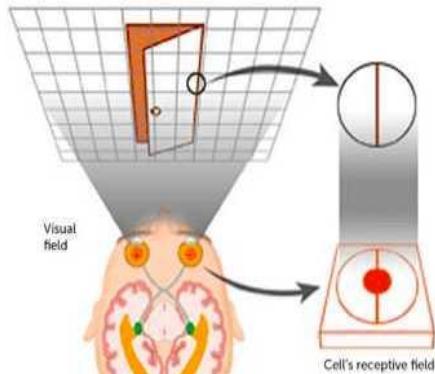
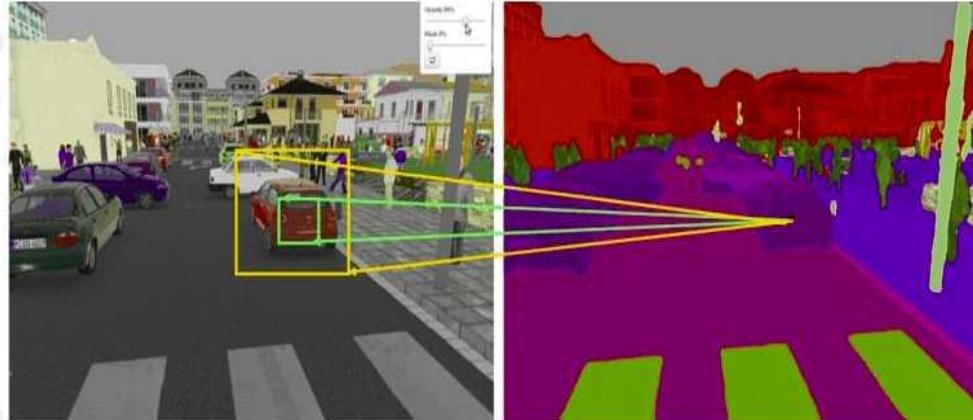


- Neurons respond to very narrow areas of stimulation
- The role of neurons in recognizing different shapes of characteristics
- Hierarchical processing of visual information

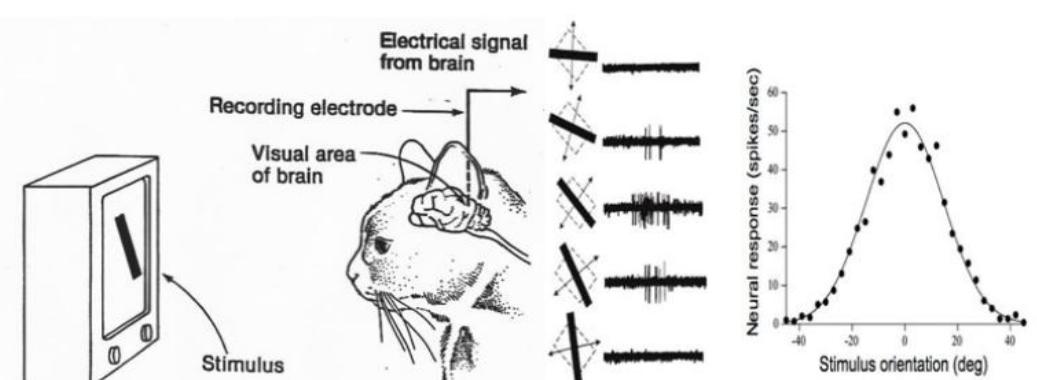
Convolutional Neural Network



❖ Principles of visual information processing



"고양이의 눈에서 답을 얻다."



출처: cs231

고 수준의 뉴런이 저 수준의 뉴런의 출력에 기반한다

출처: <https://brainconnection.brainhq.com/2004/03/06/overview-of-receptive-fields/>
<https://theaisummer.com/receptive-field/>

Convolution Layer



❖ Convolution

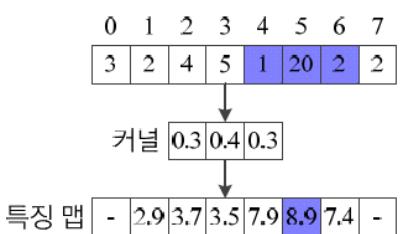
- Convolution is a linear operation that multiplies the corresponding elements and adds all the results
- u : kernel, z : input, s : output (feature map)

$$s(i) = z \circledast u = \sum_{x=-(h-1)/2}^{(h-1)/2} z(i+x)u(x)$$

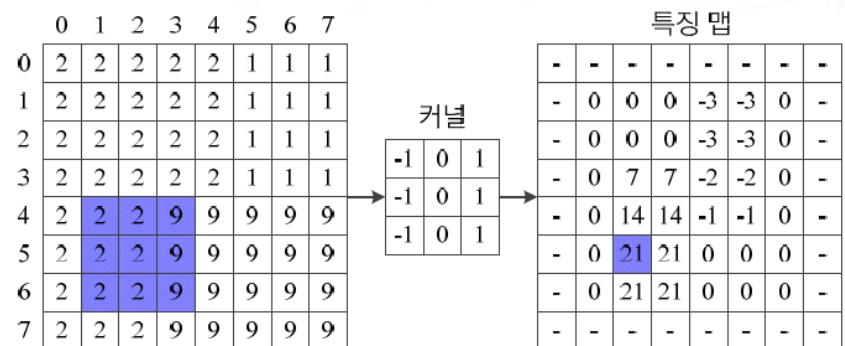
← 1D input

$$s(j, i) = z \circledast u = \sum_{y=-(h-1)/2}^{(h-1)/2} \sum_{x=-(h-1)/2}^{(h-1)/2} z(j+y, i+x)u(y, x)$$

← 2D input

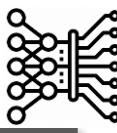


(a) 1차원 컨볼루션



(b) 2차원 컨볼루션

Convolution Layer



❖ Characteristics 1: Convolution Kernel

입력데이터 (4,4)

0	3	2	1
2	1	0	2
1	3	2	0
2	0	3	3

* 필터 (3,3)

1	0	1
0	1	0
1	0	1

= 6

(0x1)+(3x0)+(2x1)
+(2x0)+(1x1)+(0x0)
+(1x1)+(3x0)+(2x1)

입력데이터 (4,4)

0	3	2	1
2	1	0	2
1	3	2	0
2	0	3	3

* 필터 (3,3)

1	0	1
0	1	0
1	0	1

= 7

파처맵 (2,2)

6	7
10	8

입력데이터 (4,4)

0	3	2	1
2	1	0	2
1	3	2	0
2	0	3	3

* 필터 (3,3)

1	0	1
0	1	0
1	0	1

= 10

합성곱 예시

입력데이터 (4,4)

0	3	2	1
2	1	0	2
1	3	2	0
2	0	3	3

* 필터 (3,3)

1	0	1
0	1	0
1	0	1

= 8

x_{11}	x_{12}	x_{13}	x_{14}	x_{15}	x_{16}	x_{17}
x_{21}	x_{22}	x_{23}	x_{24}	x_{25}	x_{26}	x_{27}
x_{31}	x_{32}	x_{33}	x_{34}	x_{35}	x_{36}	x_{37}
x_{41}	x_{42}	x_{43}	x_{44}	x_{45}	x_{46}	x_{47}
x_{51}	x_{52}	x_{53}	x_{54}	x_{55}	x_{56}	x_{57}
x_{61}	x_{62}	x_{63}	x_{64}	x_{65}	x_{66}	x_{67}
x_{71}	x_{72}	x_{73}	x_{74}	x_{75}	x_{76}	x_{77}

필터의 슬라이딩 순서

Convolution Layer



❖ Padding

- Avoid reducing input at the edge

0	1	2	3	4	5	6	7
0	3	2	4	5	1	20	2

0.3	0.4	0.3
-----	-----	-----

1.8	2.9	3.7	3.5	7.9	8.9	7.4	1.4
-----	-----	-----	-----	-----	-----	-----	-----

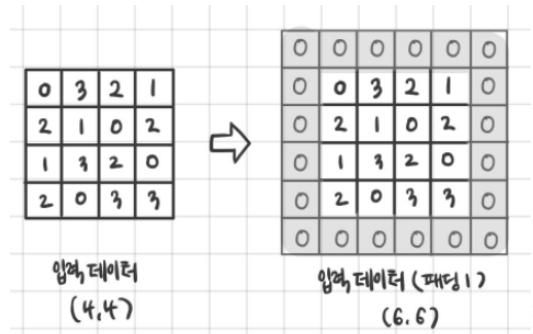
(a) 0 덧대기

0	1	2	3	4	5	6	7
3	3	2	4	5	1	20	2

0.3	0.4	0.3
-----	-----	-----

2.7	2.9	3.7	3.5	7.9	8.9	7.4	2.0
-----	-----	-----	-----	-----	-----	-----	-----

(b) 복사 덧대기



❖ Bias

0	1	2	3	4	5	6	7
3	2	4	5	1	20	2	2

1

0.2	0.3	0.4	0.3
-----	-----	-----	-----

(a) 1차원 컨볼루션

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

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

-	-	-	-	-	-	-	-
-	2	2	2	-1	-1	2	-
-	2	2	2	-1	-1	2	-
-	2	9	9	0	0	2	-
-	2	16	16	1	1	2	-
-	2	23	23	2	2	2	-
-	2	23	23	2	2	2	-
-	-	-	-	-	-	-	-

(b) 2차원 컨볼루션

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

커널

-	-	-	-	-	-	-	-
-	2	2	2	-1	-1	2	-
-	2	2	2	-1	-1	2	-
-	2	9	9	0	0	2	-
-	2	16	16	1	1	2	-
-	2	23	23	2	2	2	-
-	2	23	23	2	2	2	-
-	-	-	-	-	-	-	-

0	3	2	1
2	1	0	2
1	3	2	0
2	0	3	3

입력 테이터
(4,4)

1	0	1
0	1	0
1	0	1

필터
(3,3)

6	7
10	9

편향
(1,1)

8	9
12	10

퍼체인
(2,2)



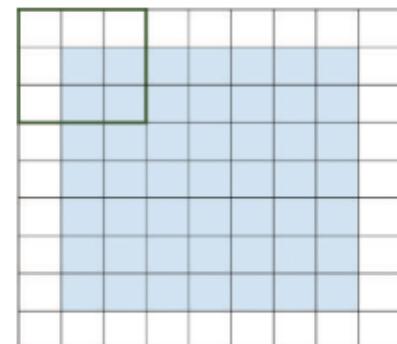
❖ Padding

- Output size of convolutional operations considering padding

$$O = \frac{(N+2 \times P) - F}{S} + 1$$

N: 입력 데이터 크기, P: 패딩, F: 콘볼루션 필터 크기, S: 스트라이드, O: 출력 데이터 크기

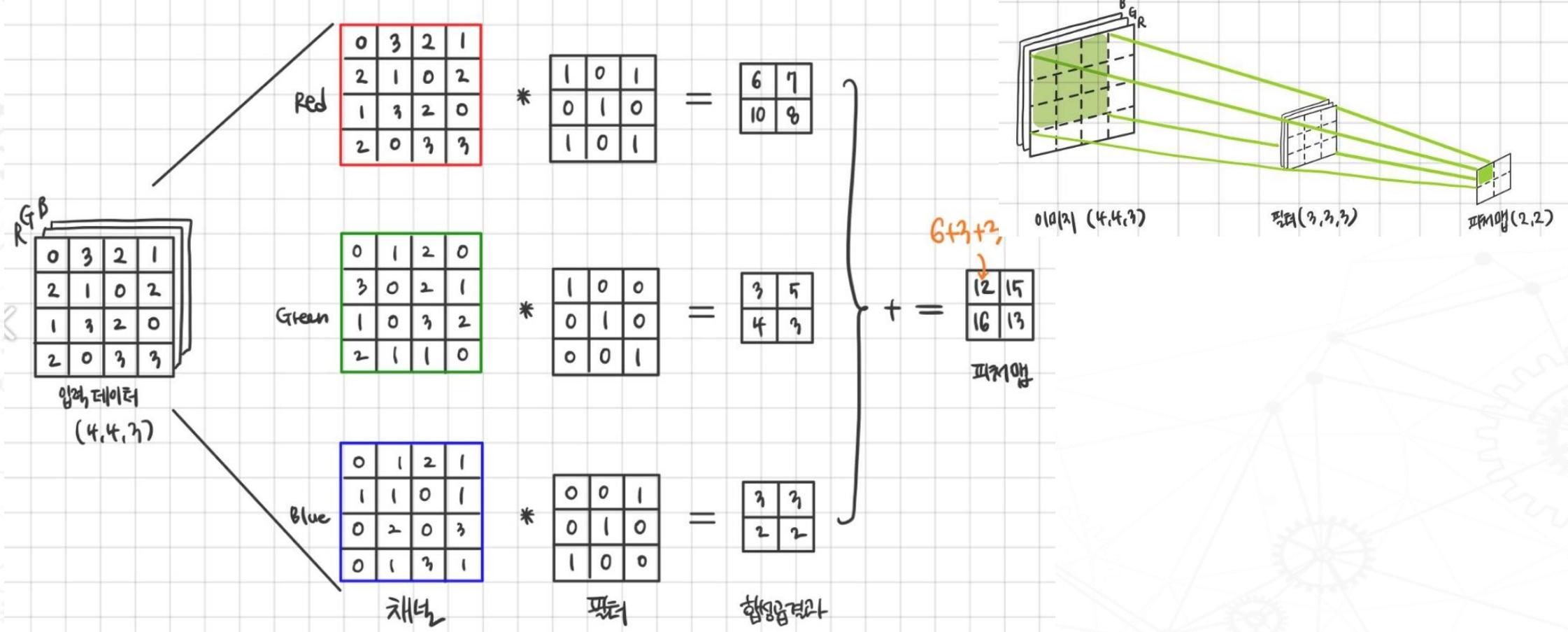
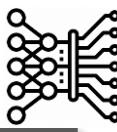
- **7×7** 이미지에 패딩을 해서 **9×9** 이미지로 만들어 보자



- 앞의 식에 N=7, P=1, F=3, S=1을 대입하면 컨볼루션 연산 후에도 이미지 크기가 7로 유지

$$O = \frac{(N+2 \times P) - F}{S} + 1 = \frac{(7+2 \times 1) - 3}{1} + 1 = 7$$

Convolution Layer

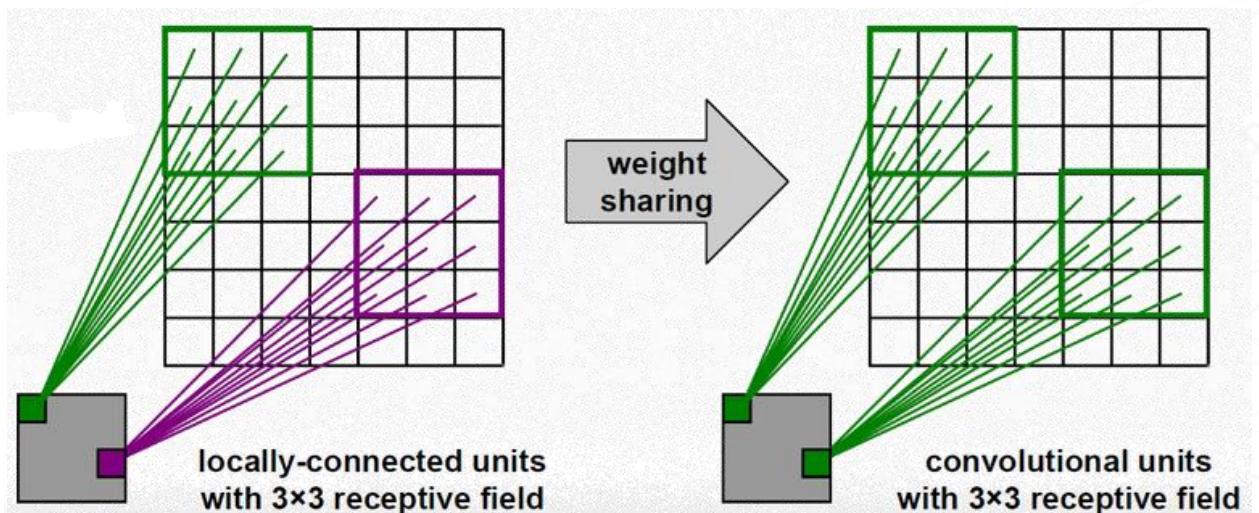
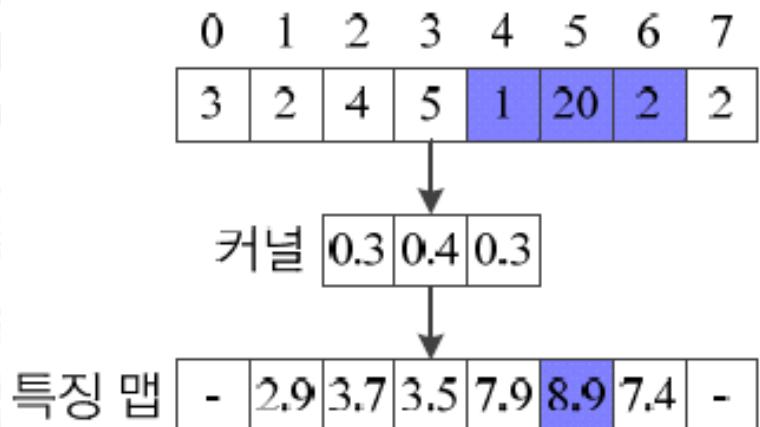


Convolution Layer



❖ Characteristics 2: Weight sharing

- All nodes use the same kernel
- Only three parameters (0.3, 0.4, 0.3) because they share weights
- Significantly reduced model complexity



Convolution Layer



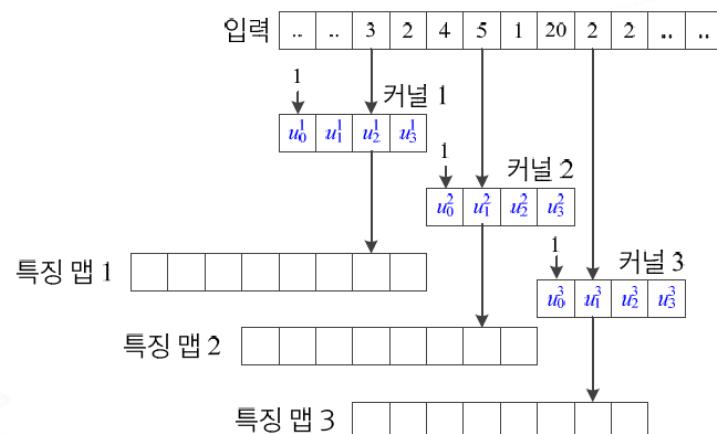
❖ Extract multiple feature maps

- Depending on the value of the kernel, the characteristics that the kernel extracts will vary
- Example) Vertical edge, horizontal edge

$$\begin{pmatrix} -1 & 0 & 1 \\ -1 & 0 & 1 \\ -1 & 0 & 1 \end{pmatrix}$$

$$\begin{pmatrix} -1 & -1 & -1 \\ 0 & 0 & 0 \\ 1 & 1 & 1 \end{pmatrix}$$

- Therefore, using only one kernel will extract too low-level features
- Figure shows a situation in which three feature maps are extracted using three kernels
- Hundreds of kernels in CNN



Convolution Layer

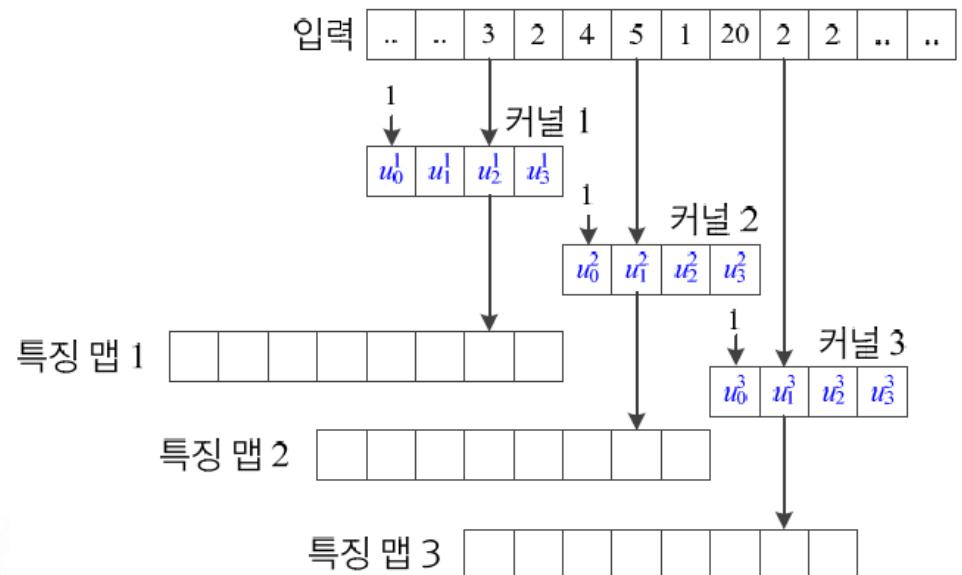


❖ Feature learning

- Kernel is not designed by humans, it is discovered by learning

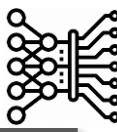
- u_i^k : i th parameter of k th kernel

- 예) 2차원 영상이 $7*7$ 커널을 64개 사용한다면, 학습 $(7*7+1)*64=3200$ 개의 매개변수를 찾아내야 함



- Training the kernel using backpropagation

Convolution Layer

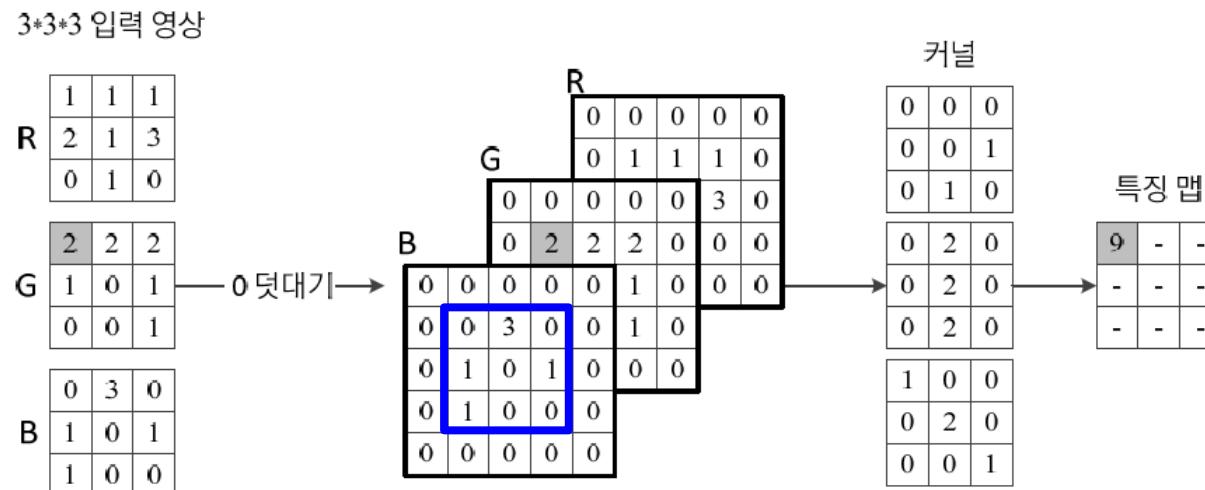


- Tensor: 벡터 계산을 단순화하기 위해 같은 성질의 여러 벡터를 한 행렬 안에 표기하고 그것을 단순화하여 표기한 것

❖ Apply to Tensor*

- Applicable to 3D or high-dimensional structures

- e.g., RGB 컬러 영상은 $3*m*n$ 의 3차원 텐서



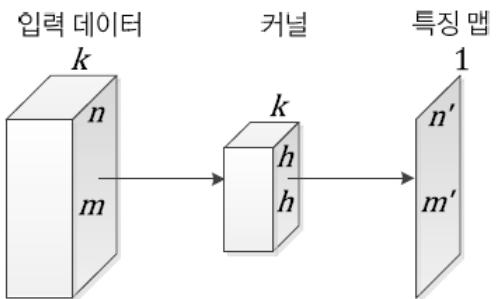
- Example of calculation of gray nodes in feature maps

$$\underbrace{\begin{pmatrix} 0 & 0 & 0 \\ 0 & 1 & 1 \\ 0 & 2 & 1 \end{pmatrix}}_{R} \underbrace{\begin{pmatrix} 0 & 0 & 0 \\ 0 & 2 & 2 \\ 0 & 1 & 0 \end{pmatrix}}_{G} \underbrace{\begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 3 \\ 0 & 1 & 0 \end{pmatrix}}_{B} \odot \underbrace{\begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{pmatrix}}_{c_1} \underbrace{\begin{pmatrix} 0 & 2 & 0 \\ 0 & 2 & 0 \\ 0 & 2 & 0 \end{pmatrix}}_{c_2} \underbrace{\begin{pmatrix} 1 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 1 \end{pmatrix}}_{c_3} = 9$$

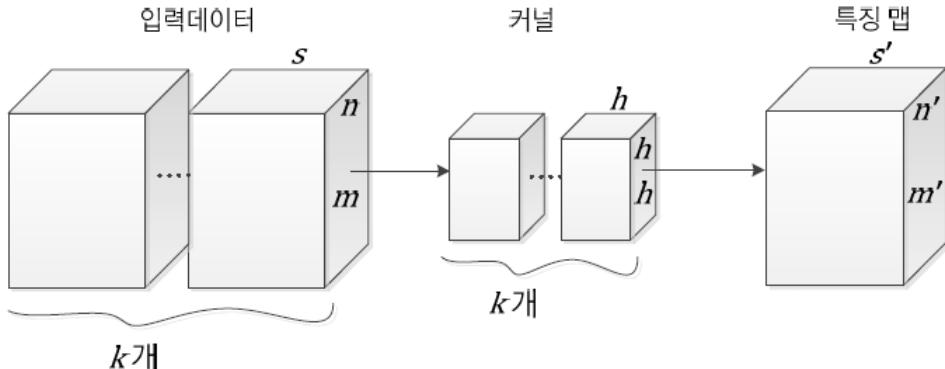
Convolution Layer



❖ Apply to data in 3D structure



(a) 다중채널 데이터(예: RGB 컬러 영상)



(b) 3차원 데이터(예: 동영상, MRI 뇌 영상)

- Calculate the size of the input data using the output size of layer

$$O = \frac{N-F}{S} + 1$$

$$N = (O-1) \times S + F$$

N : 입력 데이터 크기, F : 콘벌루션 필터 크기, S : 스트라이드, O : 출력 데이터 크기

Convolution Layer



❖ Example

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
conv3-64	conv3-64	conv3-64	conv3-64	conv3-64	conv3-64
	LRN	conv3-64	conv3-64	conv3-64	conv3-64
maxpool					
conv3-128	conv3-128	conv3-128	conv3-128	conv3-128	conv3-128
		conv3-128	conv3-128	conv3-128	conv3-128
maxpool					
conv3-256	conv3-256	conv3-256	conv3-256	conv3-256	conv3-256
conv3-256	conv3-256	conv3-256	conv3-256	conv1-256	conv3-256
					conv3-256
maxpool					
conv3-512	conv3-512	conv3-512	conv3-512	conv3-512	conv3-512
conv3-512	conv3-512	conv3-512	conv3-512	conv1-512	conv3-512
					conv3-512
maxpool					
conv3-512	conv3-512	conv3-512	conv3-512	conv3-512	conv3-512
conv3-512	conv3-512	conv3-512	conv3-512	conv1-512	conv3-512
					conv3-512
maxpool					
FC-4096					
FC-4096					
FC-1000					
soft-max					

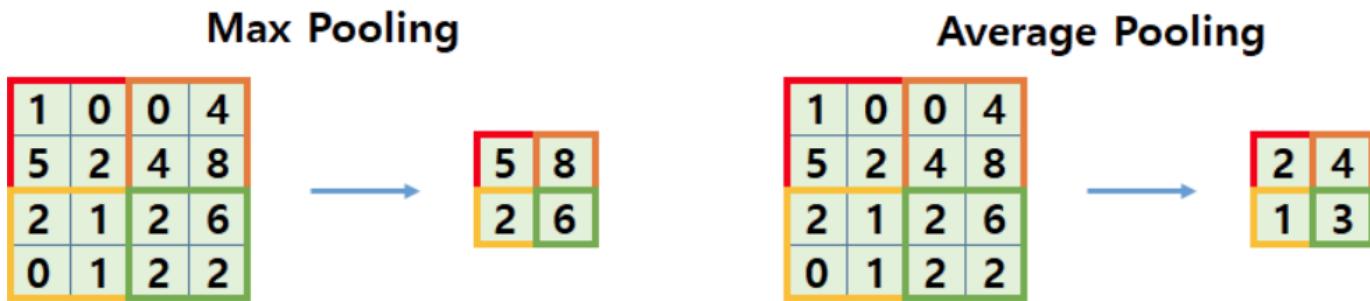
CNN-1	CNN-2	CNN-3
Output FC layer (44.29)	Output FC layer (91.46)	<i>Output FC layer (92.05))</i>
10 × 10 (88.67)	10 × 10 (91.14)	10 × 10 (91.03)
16 × 10 (88.72)	16 × 10 (91.58)	16 × 10 (91.77)
32 × 10 (88.93)	32 × 10 (91.99)	32 × 10 (92.02)
64 × 10 (89.72)	64 × 10 (91.82)	64 × 10 (91.8)
128 × 10 (89.2)	128 × 10 (91.86)	128 × 10 (89.2)
256 × 10 (89.23)	256 × 10 (92.02)	256 × 10 (89.23)
512 × 10 (88.95)	512 × 10 (90.98)	512 × 10 (91.78)
1024 × 10 (89.56)	1024 × 10 (91.54)	1024 × 10 (92.22)
2048 × 10 (87.4)	2048 × 10 (91.27)	2048 × 10 (91.59)
4096 × 10 (86.27)	4096 × 10 (87.51)	4096 × 10 (90.68)
64 × 64 × 10 (89.35)	256 × 256 × 10 (91.97)	1024 × 1024 × 10 (91.27)
128 × 64 × 10 (89.71)	512 × 256 × 10 (91.92)	2048 × 1024 × 10 (91.43)
256 × 64 × 10 (89.79)	1024 × 256 × 10 (91.53)	4096 × 1024 × 10 (91.94)
512 × 64 × 10 (89.88)	2048 × 256 × 10 (91.95)	-
1024 × 64 × 10 (90)	4096 × 256 × 10 (92.29)	-
2048 × 64 × 10 (90.28)	4096 × 4096 × 256 × 10 (91.64)	-
4096 × 64 × 10 (90.59)	-	-
4096 × 4096 × 64 × 10 (90.77)	-	-
4096 × 4096 × 4096 × 64 × 10 (90.74)	-	-

Pooling Layer



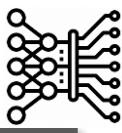
❖ Characteristics 3: Pooling

- All the data in the output feature map that went through the previous layers is not required
 - Pooling: Max pooling, Avg pooling, Weight avg pooling, and etc.



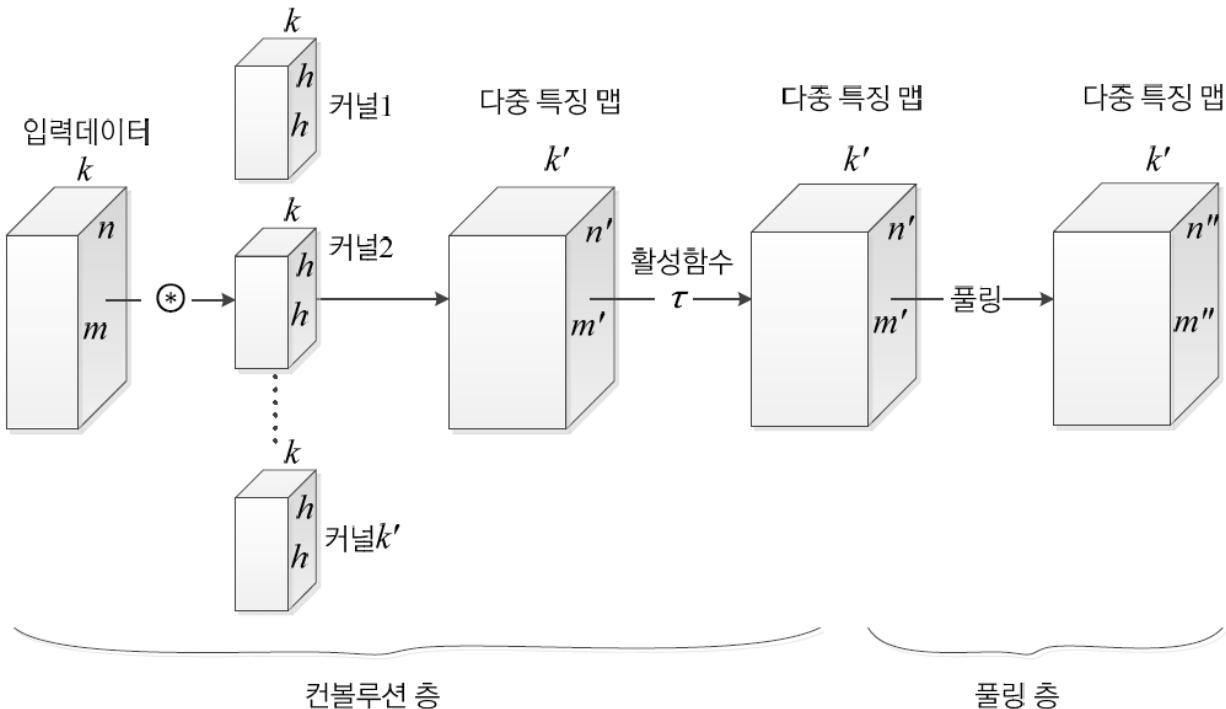
❖ Effects

- Reduce parameters, reduce network expression and suppress overfitting
- Reduced computation saves hardware resources and speeds up



❖ Building block

- CNN connects building blocks to create deep architecture
- Conv. → activation function → pooling
- Extracts multiple feature maps using multiple kernels



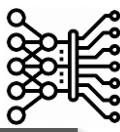
Overall Structure



❖ Configurable parameters

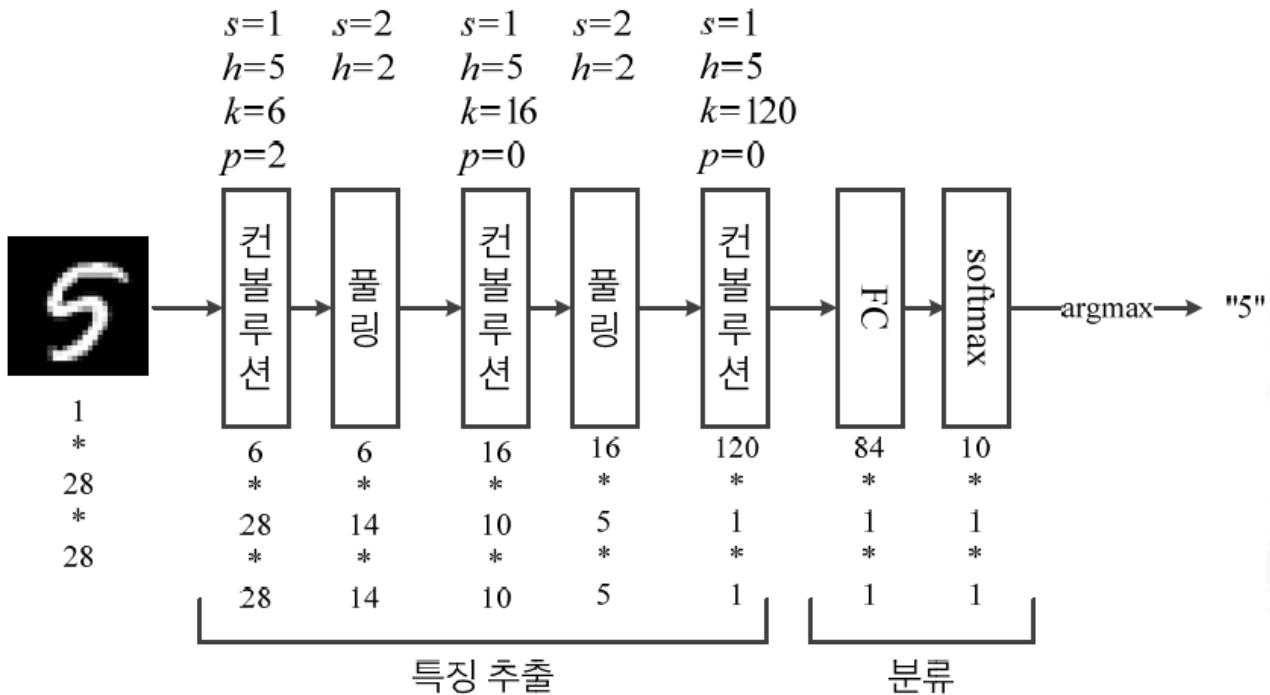
- Convolution filter
- Filter size
- Padding
- Stride
- Pooling layer

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

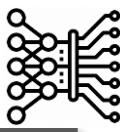


❖ LeNet-5

- Feature extraction: Converting a 28*28 contrast map into a 120-dimensional feature vector through five layers of C-P-C-P-C
- Classification: MLP with one hidden layer

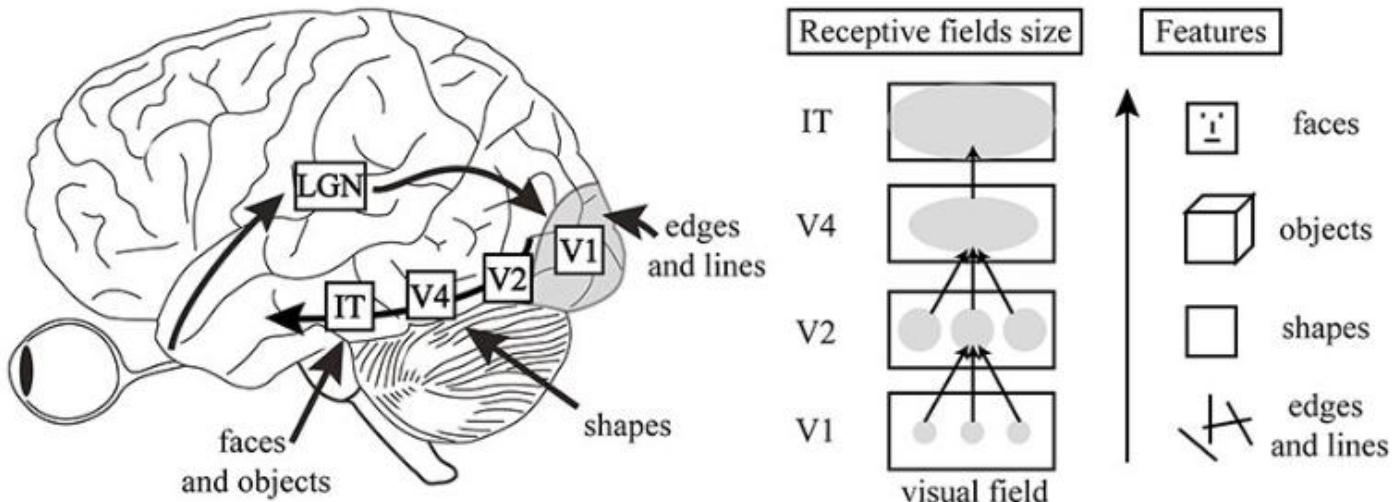


Overall Structure



❖ DMLP vs. CNN

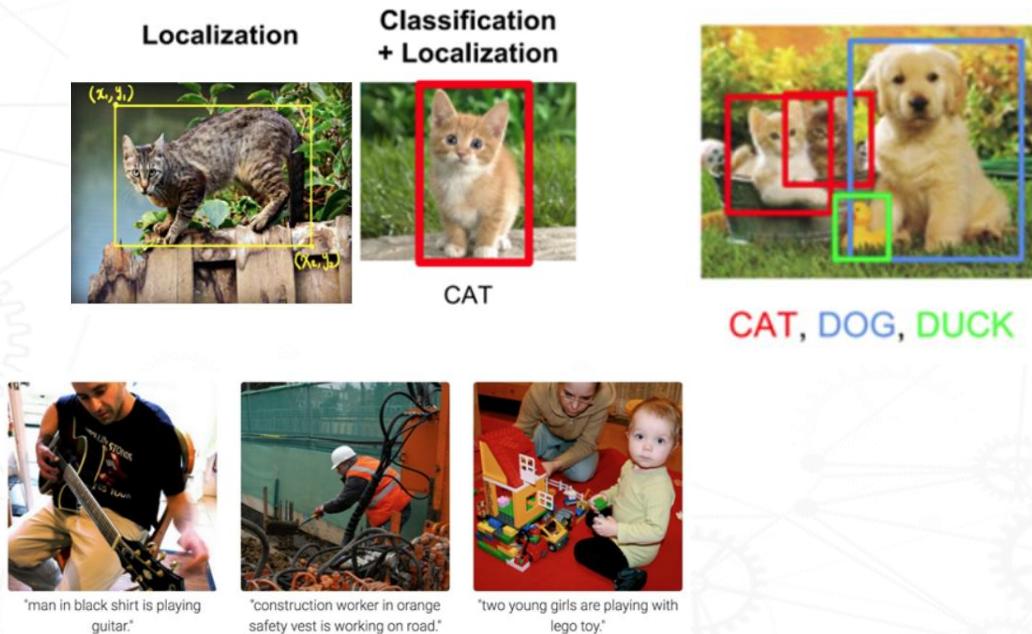
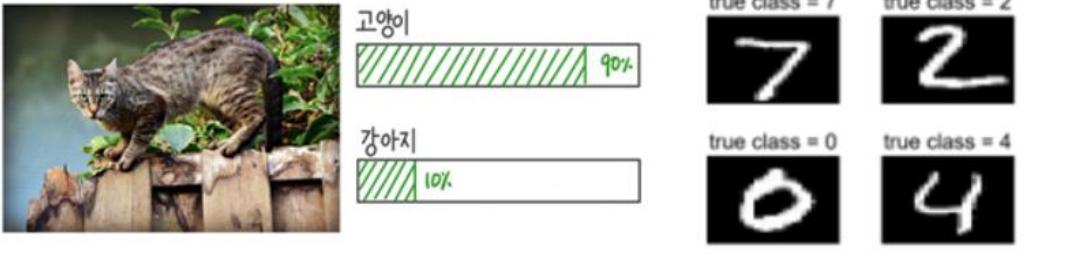
- DMLP cannot be calculated if the feature vector is different in size
- CNN has the advantage of handling variable sizes
 - Adjust the size of the feature map by adjusting the stride in the convolution layer or by adjusting the kernel or stride in the pooling layer



Application



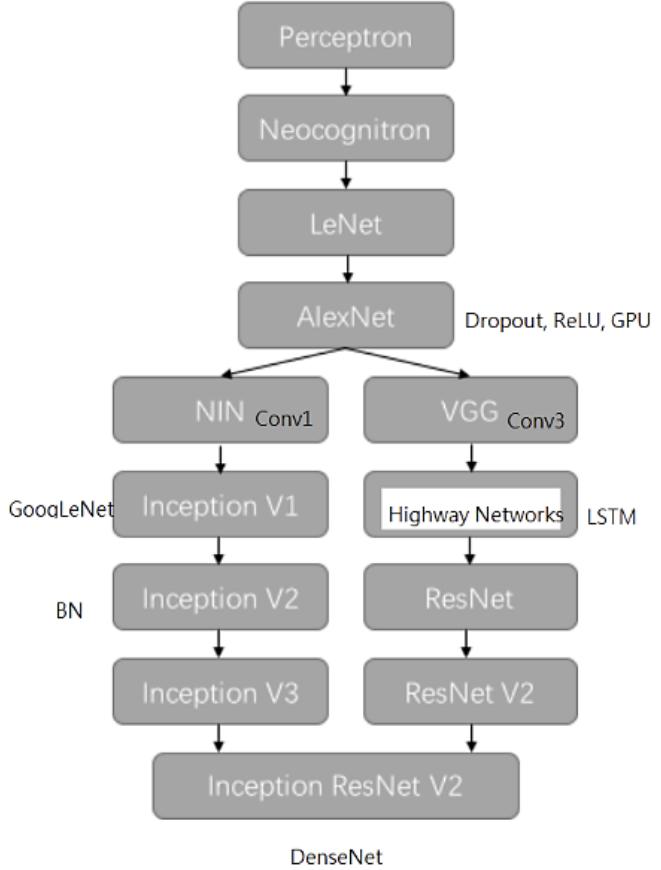
- ❖ Image classification
- ❖ Object localization
- ❖ Object detection
- ❖ Image captioning

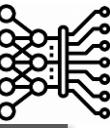


Representative CNN



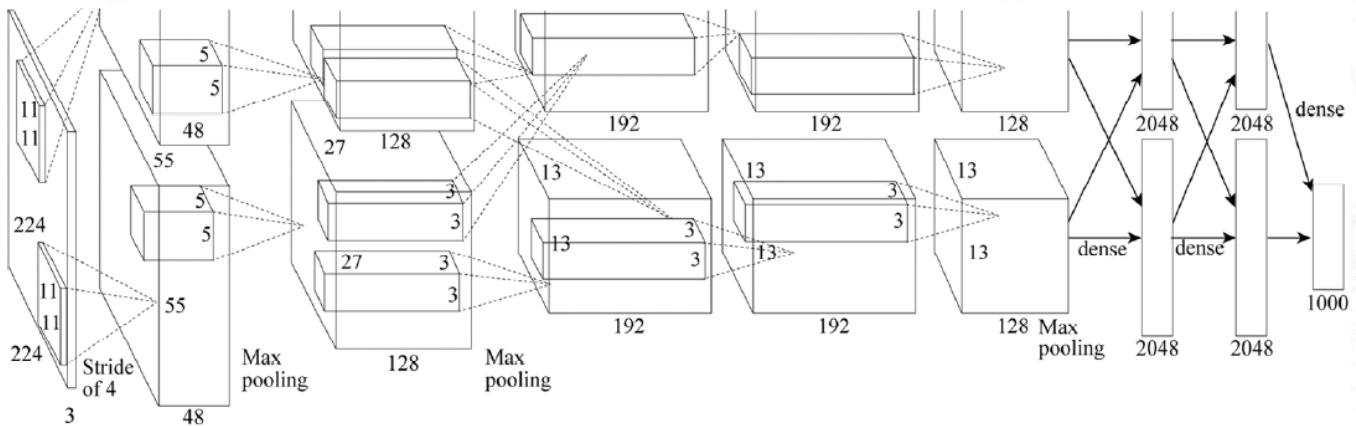
- ❖ AlexNet
- ❖ VGGNet
- ❖ GoogLeNet
- ❖ ResNet





❖ Structure

- 5 convolutional layers and 3 fully connected (FC) layers
 - 8 layer, 290400-186624-64896-43264-4096-4096-1000 nodes
- 2 million parameters for convolution layer and 65 million parameters for FC layer
- 30x more parameters for FC layer → In the future, CNN will advanced toward reducing the parameters of FC layer

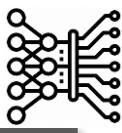




❖ What made AlexNet successful in learning?

- External factors
 - A large database called ImageNet
 - Parallel processing using GPUs
- Internal factors
 - Use ReLU as active function
 - Applying local response normalization techniques
 - Apply multiple regulatory techniques to prevent overfitting
 - Data augmentation (2048x growth with truncation and reversal)
 - Dropout, etc.
- Ensemble application during test phase
- 2 ~ 3% reduction in error rates

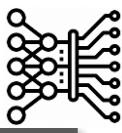




❖ The number of parameters

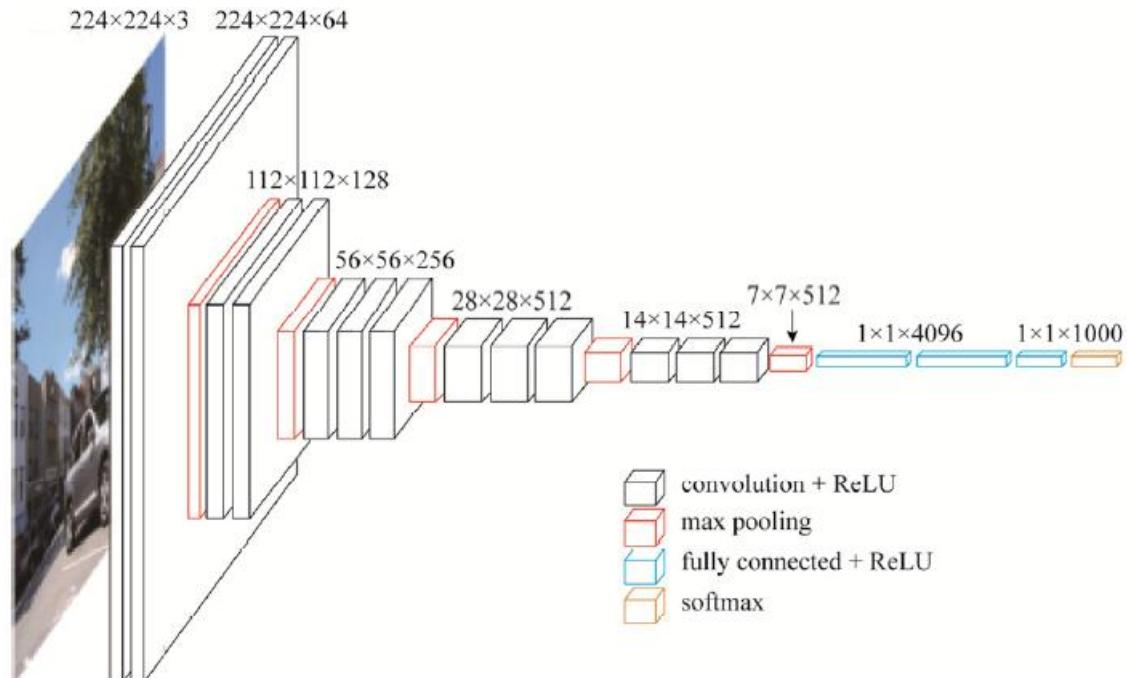
계층	텐서 크기	가중치	편향	파라미터
입력 이미지	$227 \times 227 \times 3$	0	0	0
CONV1	$55 \times 55 \times 96$	34,848	96	34,944
MaxPool-1	$27 \times 27 \times 96$	0	0	0
CONV2	$27 \times 27 \times 256$	614,400	256	614,656

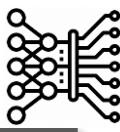
계층	텐서 크기	가중치	편향	파라미터
MaxPool-2	$13 \times 13 \times 256$	0	0	0
CONV3	$13 \times 13 \times 384$	884,736	384	885,120
CONV4	$13 \times 13 \times 384$	1,327,104	384	1,327,488
CONV5	$13 \times 13 \times 256$	884,736	256	884,992
MaxPool-3	$6 \times 6 \times 256$	0	0	0
FC1	4096×1	37,748,736	4,096	37,752,832
FC2	4096×1	16,777,216	4,096	16,781,312
FC3	1000×1	4,096,000	1,000	4,097,000
출력	1000×1	0	0	0
전체				62,378,344



❖ Structure

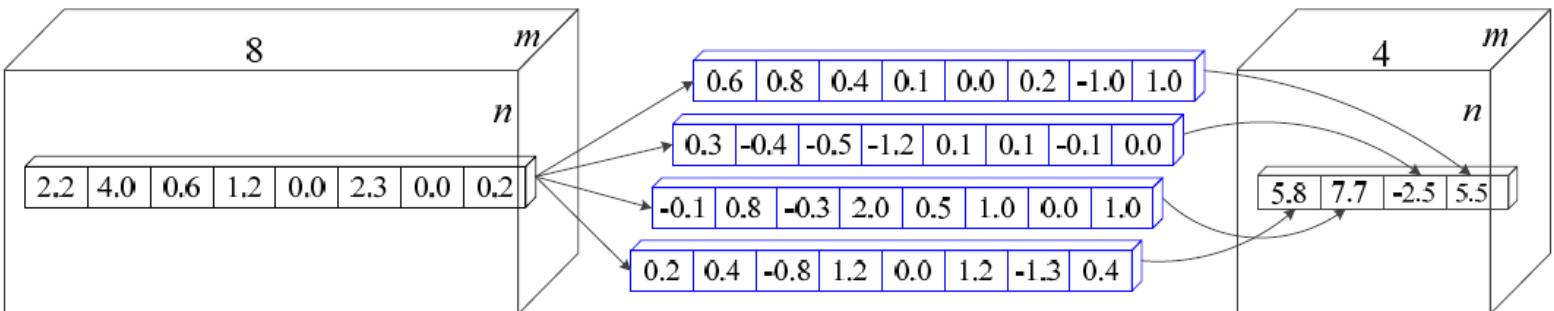
- Deepen neural networks using a small 3×3 kernel
- 8 to 16 convolution layers for 2 to 3 times deeper than 5 on AlexNet
- 16 blocks for VGG-16(13 Conv. layer + 3 FC layer)

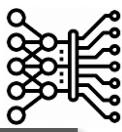




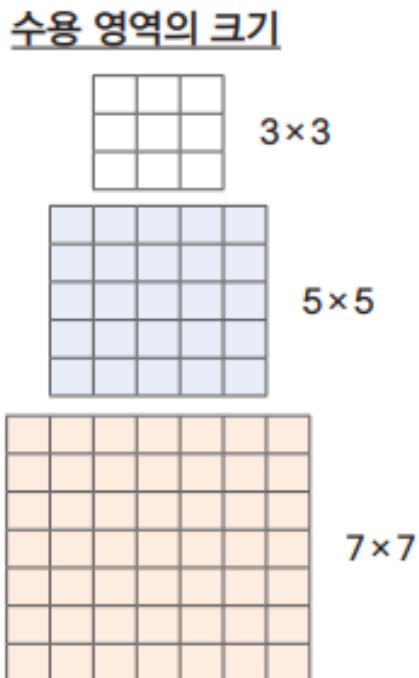
❖ 1*1 kernel

- Dimension reduction
 - Apply four 1*1 kernels to 8 feature maps of $m \times n \rightarrow 4$ feature maps of $m \times n$
 - In other words, four $m \times n$ tensors are output by applying four $8 \times 1 \times 1 \times 1$ kernels to $8 \times m \times n$ tensors
- Applying nonlinear active functions such as ReLU increases the discernment of feature maps

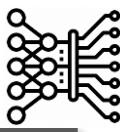




❖ Comparison with large and small convolution filters

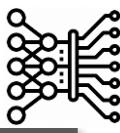


	필터 1개 사용	여러 계층에 3×3 필터 사용
	3×3 필터	3×3 필터 \times 1개 계층
파라미터 수:	9개	$9 \times 1 = 9$ 개
	5×5 필터	3×3 필터 \times 2개 계층
파라미터 수:	25개	$9 \times 2 = 18$ 개
	7×7 필터	3×3 필터 \times 3개 계층
파라미터 수:	49개	$9 \times 3 = 27$ 개



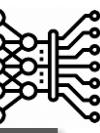
❖ Specification

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



❖ The number of parameters

구분 계층	계층 이름	출력 크기	액티베이션 맵	파라미터 수
	입력 이미지	224×224×3	224×224×3=150K	0
Conv1	Conv1-1	224×224×64	224×224×64=3.2M	$(3 \times 3 \times 3) \times 64 = 1,728$
	Conv1-2	224×224×64	224×224×64=3.2M	$(3 \times 3 \times 64) \times 64 = 36,864$
	MaxPool-1	112×112×64	112×112×64=800K	0
Conv2	Conv2-1	112×112×128	112×112×128=1.6M	$(3 \times 3 \times 64) \times 128 = 73,728$
	Conv2-2	112×112×128	112×112×128=1.6M	$(3 \times 3 \times 128) \times 128 = 147,456$
	MaxPool-2	56×56×128	56×56×128=400K	0
Conv3	Conv3-1	56×56×256	56×56×256=800K	$(3 \times 3 \times 128) \times 256 = 294,912$
	Conv3-2	56×56×256	56×56×256=800K	$(3 \times 3 \times 256) \times 256 = 589,824$
	Conv3-3	56×56×256	56×56×256=800K	$(3 \times 3 \times 256) \times 256 = 589,824$
	MaxPool-3	28×28×256	28×28×256=200K	0
Conv4	Conv4-1	28×28×512	28×28×512=400K	$(3 \times 3 \times 256) \times 512 = 1,179,648$
	Conv4-2	28×28×512	28×28×512=400K	$(3 \times 3 \times 512) \times 512 = 2,359,296$
	Conv4-3	28×28×512	28×28×512=400K	$(3 \times 3 \times 512) \times 512 = 2,359,296$
	MaxPool-4	14×14×512	14×14×512=100K	
Conv5	Conv5-1	14×14×512	4×14×512=100K	$(3 \times 3 \times 512) \times 512 = 2,359,296$
	Conv5-2	14×14×512	4×14×512=100K	$(3 \times 3 \times 512) \times 512 = 2,359,296$
	Conv5-3	14×14×512	4×14×512=100K	$(3 \times 3 \times 512) \times 512 = 2,359,296$
	MaxPool-5	7×7×512	7×7×512=25K	
FC	FC-1	4096	4096	$7 \times 7 \times 512 \times 4096 = 102,760,448$
	FC-2	4096	4096	$4096 \times 4096 = 16,777,216$
	FC-3	1000	1000	$4096 \times 1000 = 4,096,000$
	출력	1000	0	
	전체		24M × 4 bytes ≈ 96MB	138MB

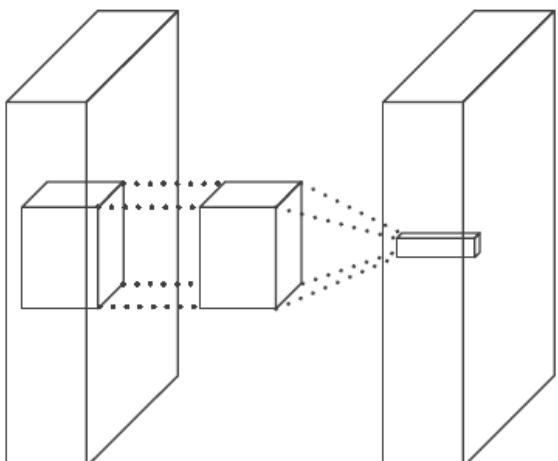


❖ Structure

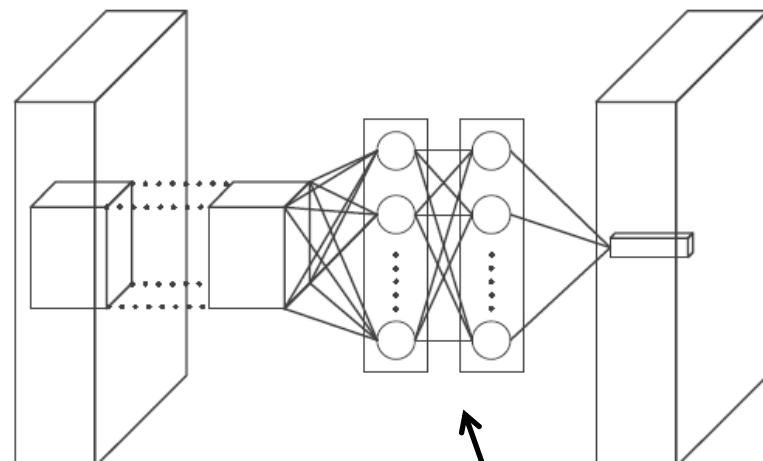
- Based on NIN(Network in Network)

❖ NIN structure

- MLPconv layer replaces convolution operation
- MLPconv performs forward computation of MLP by moving the kernel

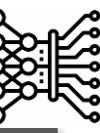


(a) 기존 컨볼루션층



(b) NIN의 MLPconv층

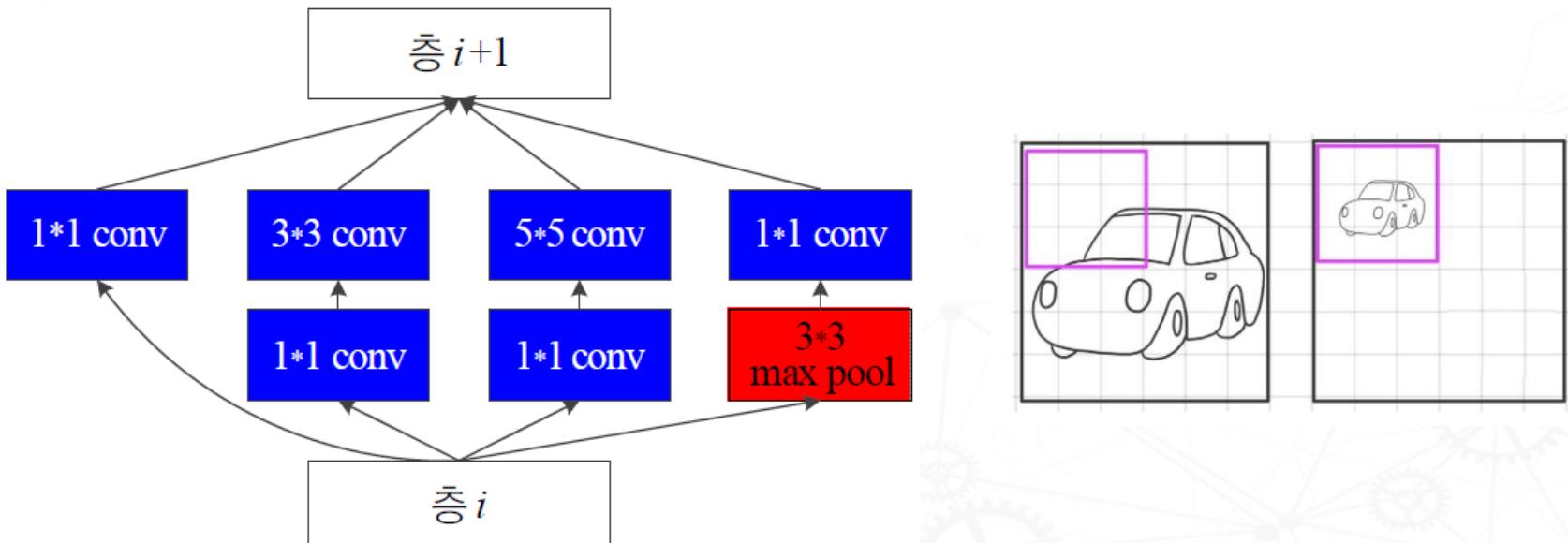
MLPconv layer (Micro network)

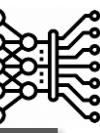


- ❖ GoogLeNet is a neural network that extends NIN ideas

- Inception module

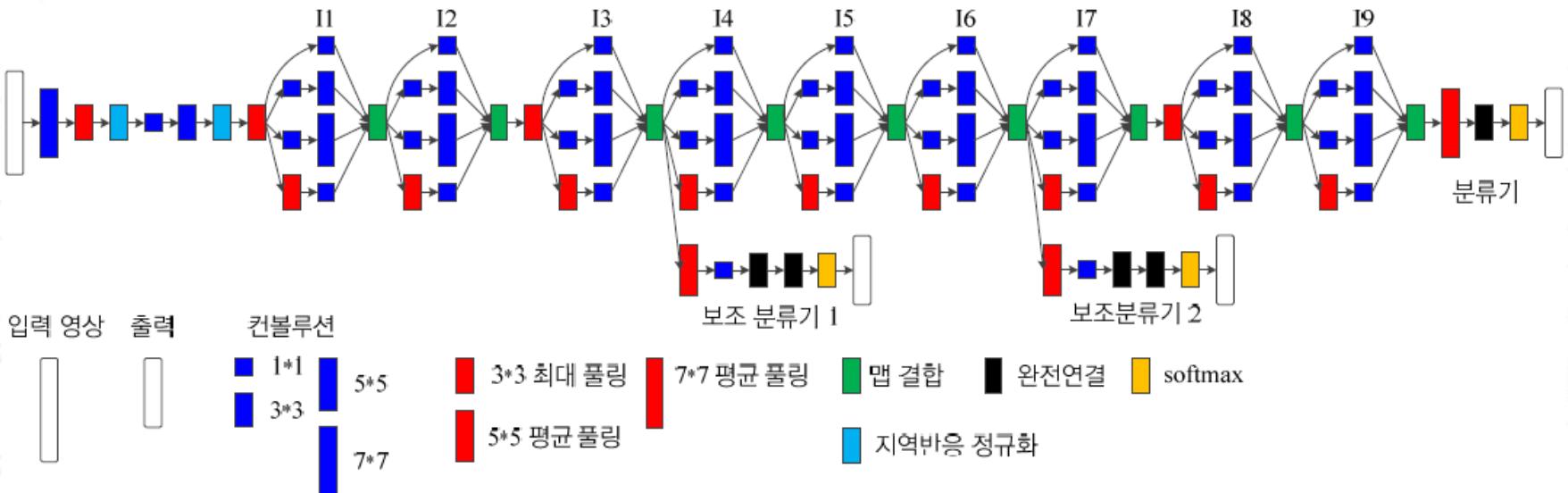
- Micronetwork uses four kinds of convolution operations instead of MLPconv
 - Use 1*1 convolution to reduce dimensions

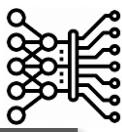




❖ GoogLeNet combines 9 inception modules

- 22 layers with parameters and 5 layers without parameters
- 1 FC layer





❖ Specification

type	필터 크기/ 스트라이드	출력 크기	깊이	1×1 필터수	3×3 병목계층 필터수	3×3 필터수	5×5 병목계층 필터수	5×5 필터수	맥스 풀링 병목계층 필터수	파라미터 수	연산 횟수
콘벌루션	7×7/2	112×112×64	1							2.7k	34M
맥스 풀링	3×3/2	56×56×64	0								
콘벌루션	3×3/1	56×56×192	2		64	192				112k	360M
맥스 풀링	3×3/2	28×28×192	0								
인셉션(3a)		28×28×256	2	64	96	128	16	32	32	159k	128M
인셉션(3b)		28×28×480	2	128	128	192	32	96	64	380k	304M
맥스 풀링	3×3/2	14×14×480	0								
인셉션(4a)		14×14×512	2	192	96	208	16	48	64	364k	73M
인셉션(4b)		14×14×512	2	160	112	224	24	64	64	437k	88M
인셉션(4c)		14×14×512	2	128	128	256	24	64	64	463k	100M
인셉션(4d)		14×14×528	2	112	144	288	32	64	64	580k	119M
인셉션(4e)		14×14×832	2	256	160	320	32	128	128	840k	170M
맥스 풀링	3×3/2	7×7×832	0								
인셉션(5a)		7×7×832	2	256	160	320	32	128	128	1072k	54M
인셉션(5b)		7×7×1024	2	384	192	384	48	128	128	1388k	71M
평균 풀링	7×7/1	1×1×1024	0								
드롭 아웃(40%)		1×1×1024	0								
완전 연결		1×1×1000	1							1000k	1M
소프트맥스		1×1×1000	0								



❖ ResNet

- Using the idea of residual learning to significantly increase the number of floors while avoiding performance degradation

- Basic CNN

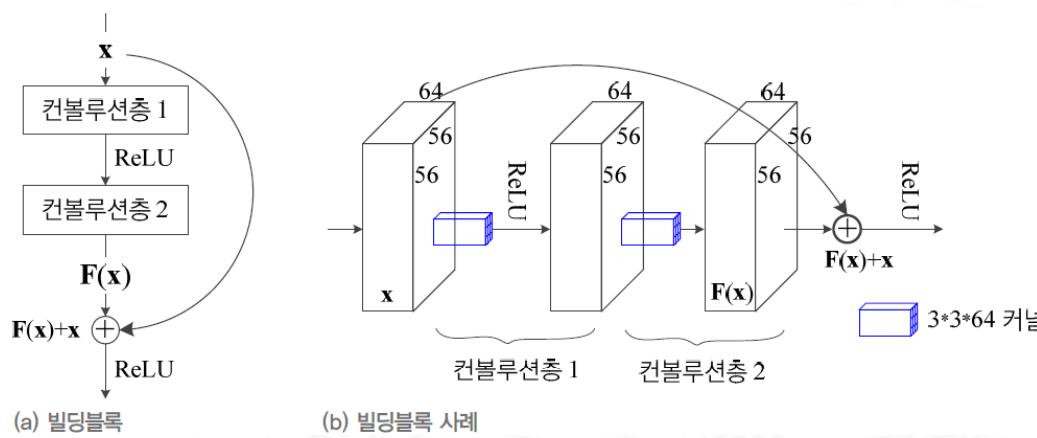
$$\mathbf{F}(\mathbf{x}) = \tau(\mathbf{x} \odot \mathbf{w}_1) \odot \mathbf{w}_2$$

$$\mathbf{y} = \tau(\mathbf{F}(\mathbf{x}))$$

- Residual learning applies τ to $\mathbf{x} + \mathbf{F}$ with shortcut connected \mathbf{x}

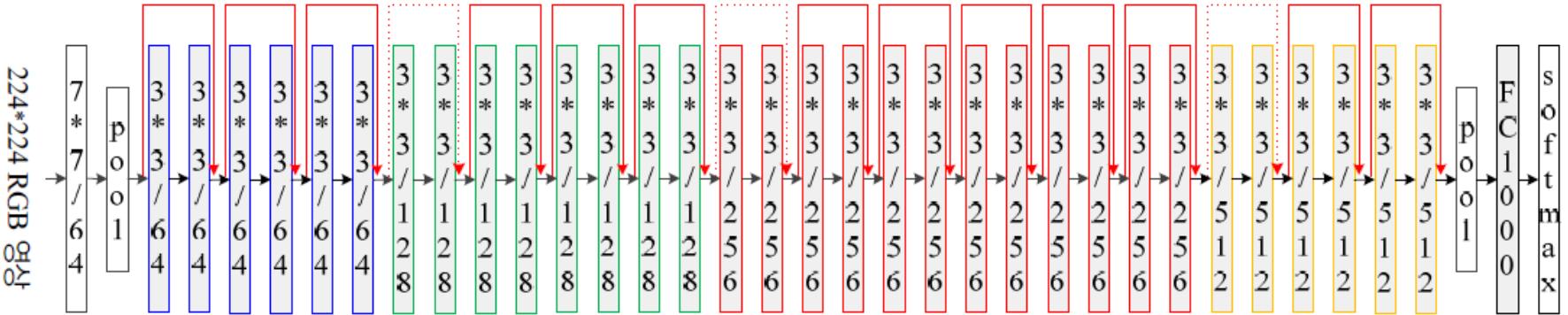
$\mathbf{F}(\mathbf{x})$ remains

$$\mathbf{y} = \tau(\mathbf{F}(\mathbf{x}) + \mathbf{x})$$





❖ Example architecture



❖ Same as VGGNet

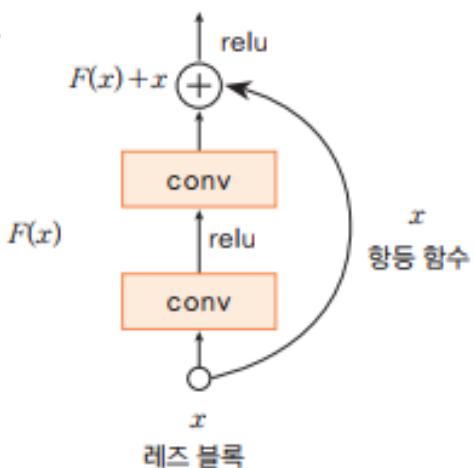
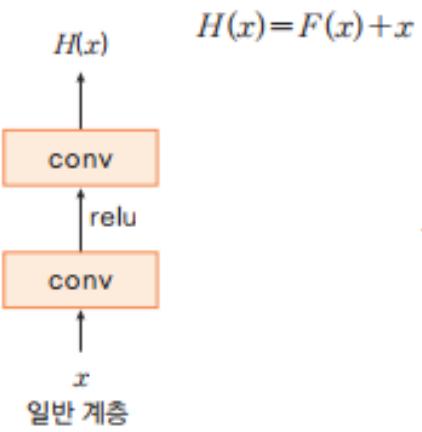
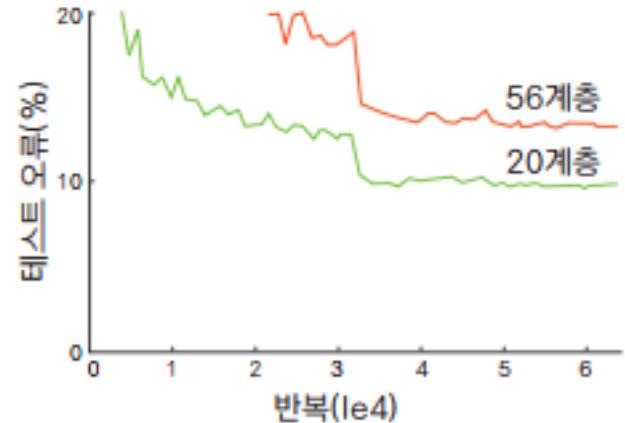
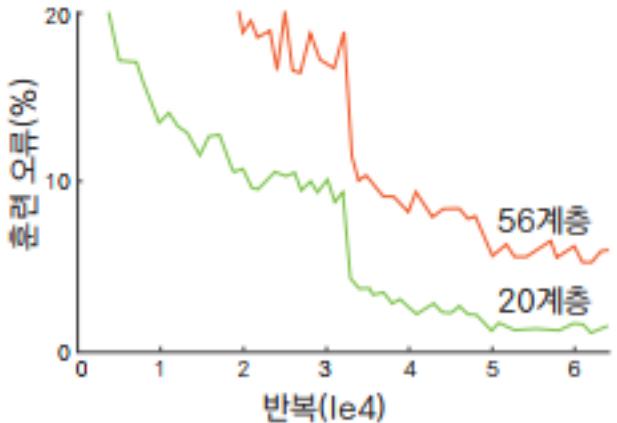
- Using 3*3 kernel

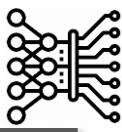
❖ Difference as VGGNet

- Residual learning
 - Using global average pooling (remove FC layer)
 - Apply to batch normalization



❖ Architecture Design





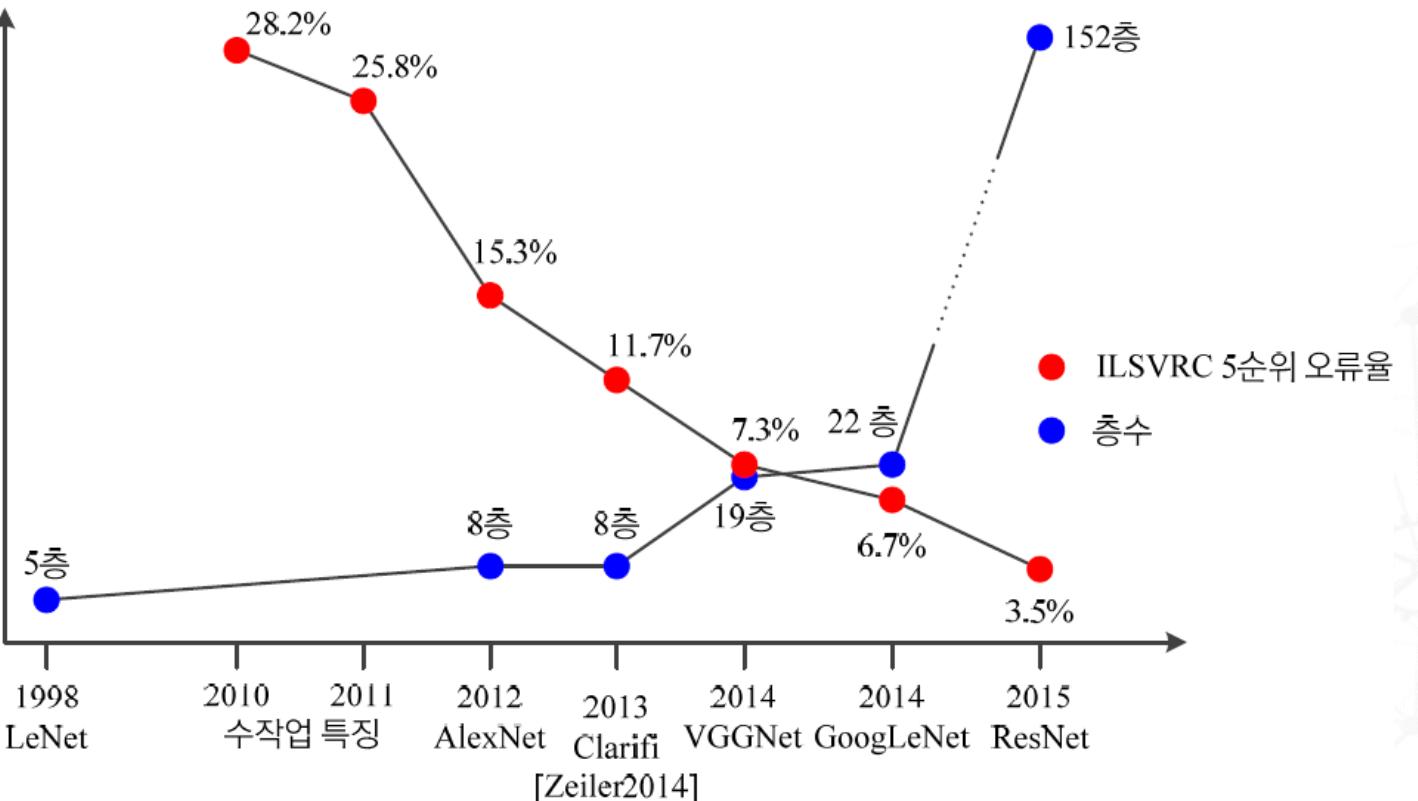
❖ Specification

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



❖ Grade in ILSVRC competition

- 2012: AlexNet (15.3% error rate)
- 2015: ResNet (3.5% error rate)





❖ Restoring colors to black-and-white pictures and images

Let there be Color!: Joint End-to-end Learning of Global and Local Image Priors for Automatic Image Colorization with Simultaneous Classification

Satoshi Iizuka* Edgar Simo-Serra* Hiroshi Ishikawa (*equal contribution)

SIGGRAPH 2016



Colorado National Park, 1941



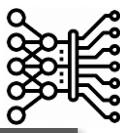
Textile Mill, June 1937



Berry Field, June 1909



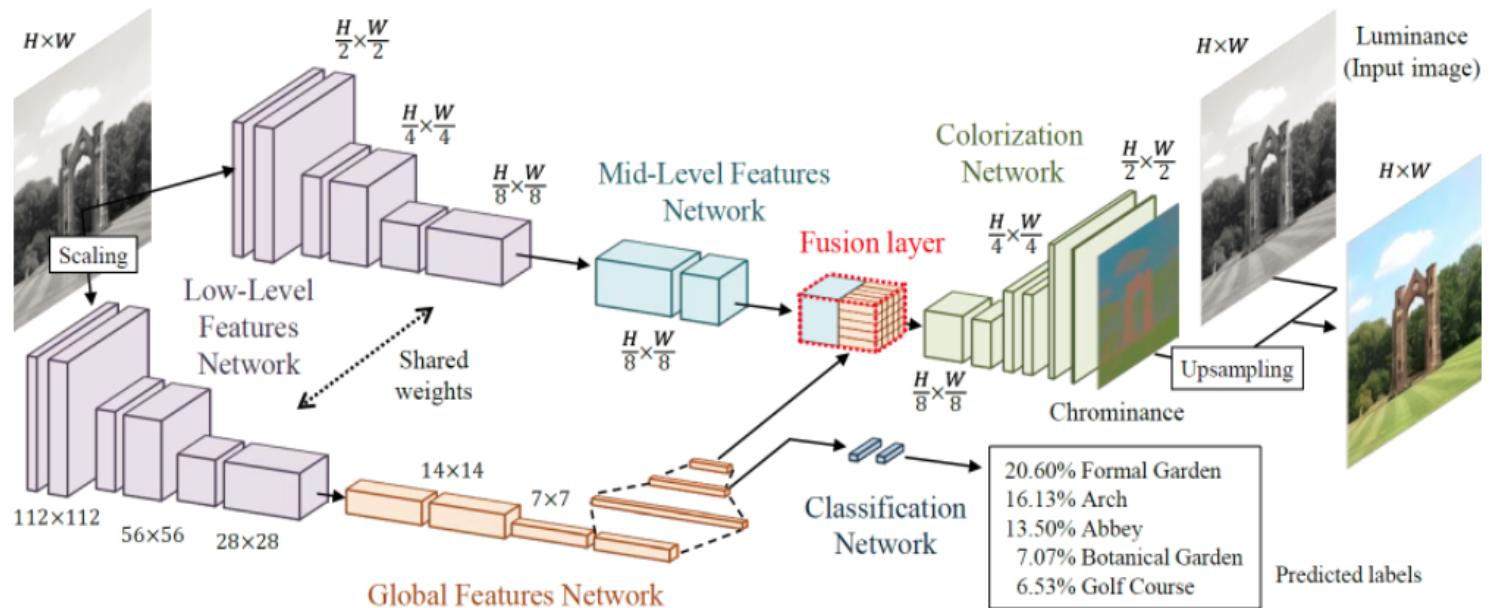
Hamilton, 1936



❖ Restoring colors to black-and-white pictures and images

Colorization Architecture:

Our model consists of four main components: a low-level features network, a mid-level features network, a global features network, and a colorization network. The components are all tightly coupled and trained in an end-to-end fashion. The output of our model is the chrominance of the image which is fused with the luminance to form the output image.





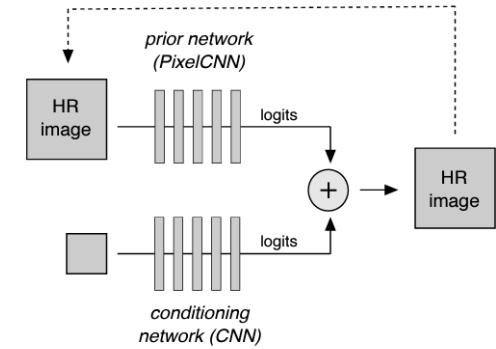
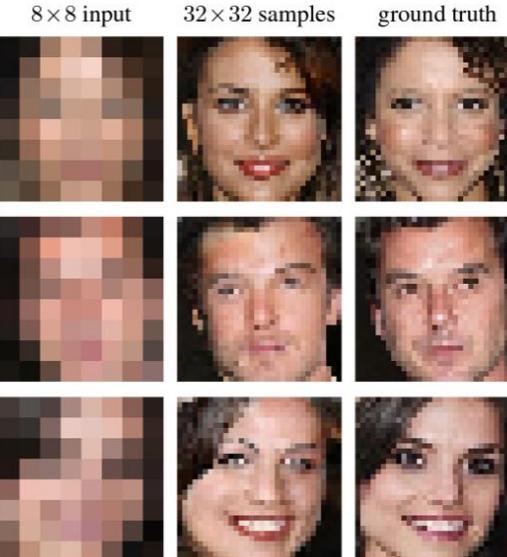
❖ Restore pixels

Pixel Recursive Super Resolution

Ryan Dahl * Mohammad Norouzi Jonathon Shlens
Google Brain
`{rld,mnorouzi,shlens}@google.com`

Abstract

Super resolution is the problem of artificially enlarging a low resolution photograph to recover a plausible high resolution version. In the regime of high magnification factors, the problem is dramatically underspecified and many plausible, high resolution images may match a given low resolution image. In particular, traditional super resolution techniques fail in this regime due to the multimodality of the problem and strong prior information that must be imposed on image synthesis to produce plausible high resolution images. In this work we propose a new probabilistic deep network architecture, a pixel recursive super resolution model, that is an extension of PixelCNNs to address this problem. We demonstrate that this model produces a diversity of plausible high resolution images at large magnification factors. Furthermore, in human evaluation studies we demonstrate how previous methods fail to fool human observers. However, high resolution images sampled from this probabilistic deep network do fool a naive human observer a significant fraction of the time.





❖ Estimating the movement of multiple people in real time

Realtime Multi-Person 2D Pose Estimation using Part Affinity Fields *

Zhe Cao Tomas Simon Shih-En Wei Yaser Sheikh
The Robotics Institute, Carnegie Mellon University
{zhecao, shihenw}@cmu.edu {tsimon, yaser}@cs.cmu.edu

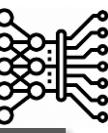
Abstract

We present an approach to efficiently detect the 2D pose of multiple people in an image. The approach uses a non-parametric representation, which we refer to as Part Affinity Fields (PAFs), to learn to associate body parts with individuals in the image. The architecture encodes global context, allowing a greedy bottom-up parsing step that maintains high accuracy while achieving realtime performance, irrespective of the number of people in the image. The architecture is designed to jointly learn part locations and their association via two branches of the same sequential prediction process. Our method placed first in the inaugural COCO 2016 keypoints challenge, and significantly exceeds the previous state-of-the-art result on the MPII Multi-Person benchmark, both in performance and efficiency.

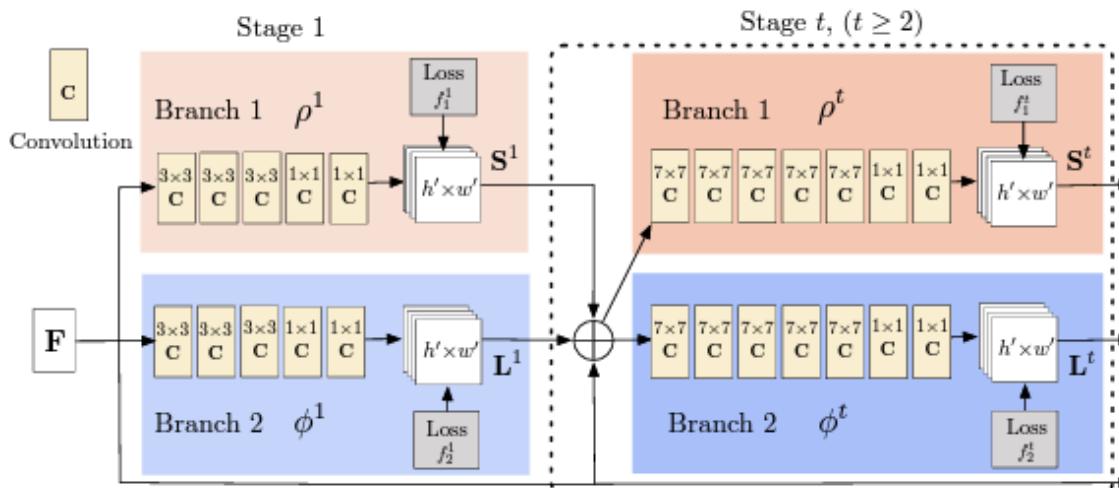
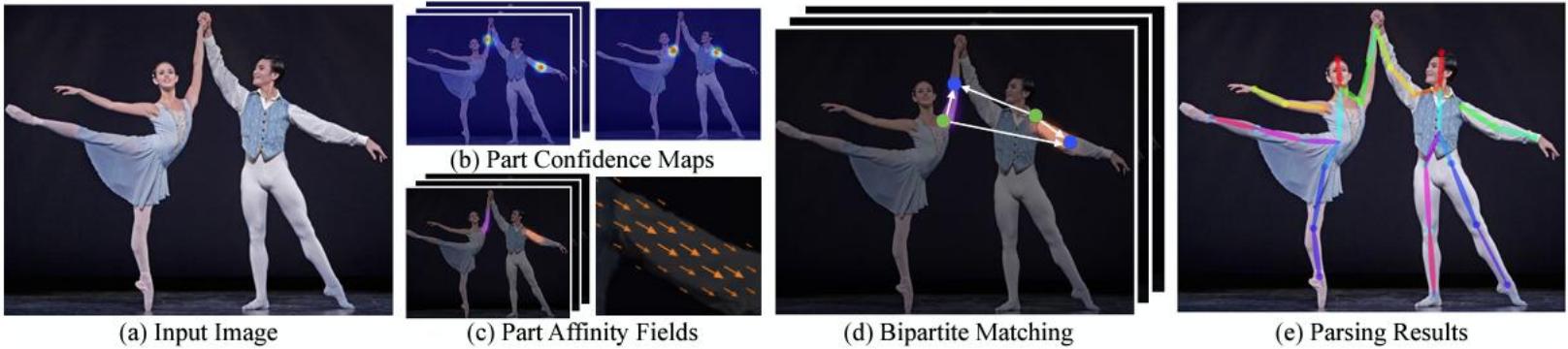


Figure 1. **Top:** Multi-person pose estimation. Body parts belonging to the same person are linked. **Bottom left:** Part Affinity Fields (PAFs) corresponding to the limb connecting right elbow and right wrist. The color encodes orientation. **Bottom right:** A zoomed in view of the predicted PAFs. At each pixel in the field, a 2D vector encodes the position and orientation of the limbs.

Applications



❖ Estimating the movement of multiple people in real time





❖ Describe a picture



This CVPR2015 paper is the Open Access version, provided by the Computer Vision Foundation.
The authoritative version of this paper is available in IEEE Xplore.

Deep Visual-Semantic Alignments for Generating Image Descriptions

Andrej Karpathy Li Fei-Fei
Department of Computer Science, Stanford University
{karpathy,feifeili}@cs.stanford.edu

Abstract

We present a model that generates natural language descriptions of images and their regions. Our approach leverages datasets of images and their sentence descriptions to learn about the inter-modal correspondences between language and visual data. Our alignment model is based on a novel combination of Convolutional Neural Networks over image regions, bidirectional Recurrent Neural Networks over sentences, and a structured objective that aligns the two modalities through a multimodal embedding. We then describe a Multimodal Recurrent Neural Network architecture that uses the inferred alignments to learn to generate novel descriptions of image regions. We demonstrate that our alignment model produces state of the art results in retrieval experiments on Flickr8K, Flickr30K and MSCOCO datasets. We then show that the generated descriptions significantly outperform retrieval baselines on both full images and on a new dataset of region-level annotations.

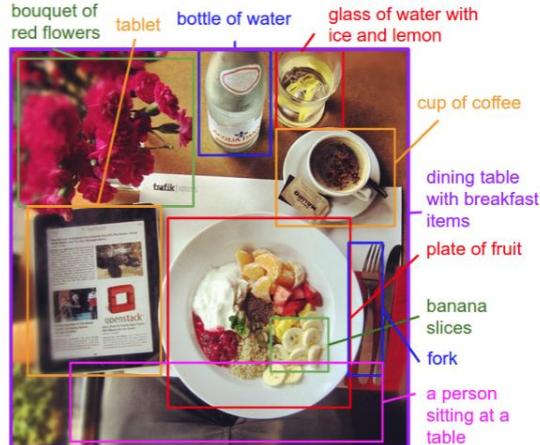


Figure 1. Motivation/Concept Figure: Our model treats language as a rich label space and generates descriptions of image regions.

Applications



❖ Describe a picture

NeuralTalk Sentence Generation Results
Showing results for coco on 1000 images

The image displays a grid of 10 image-sentence pairs, each consisting of a small thumbnail image followed by a caption and a log probability value. The images represent various scenes: food, birds, public transportation, animals, sports, urban environments, and indoor settings.

Image Description	Caption	logprob:
	a plate of food with a fork and a sandwich	-9.16
	a bird is standing on a rock in the water	-9.30
	a bus is parked on the side of the street	-7.76
	a black bear is standing in the grass	-8.37
	a young boy is holding a baseball bat	-7.11
	a bus is parked on the side of the road	-7.57
	a bathroom with a toilet sink and mirror	
	a woman is holding a cell phone to her ear	-8.22
	a clock on a pole in front of a building	
	a woman is standing in front of a refrigerator	



❖ Real-time behavioral analysis





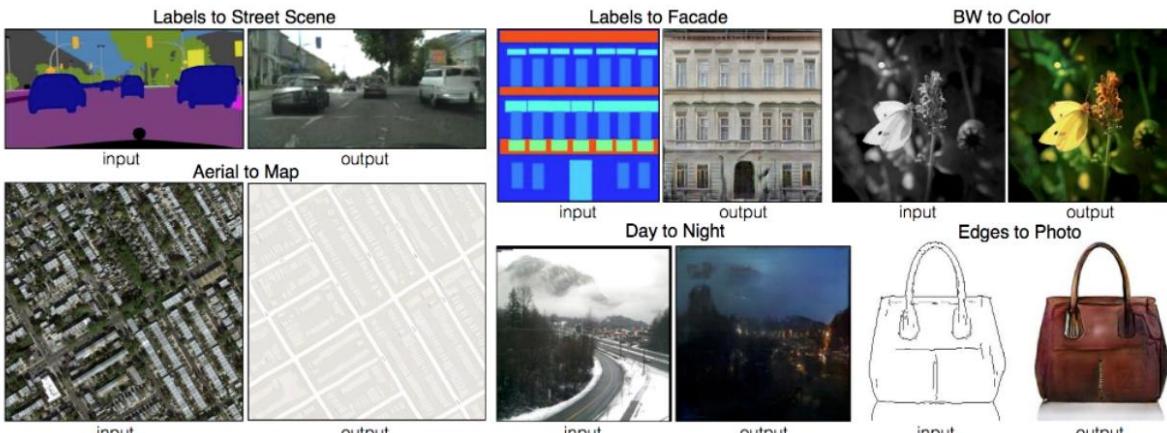
❖ Generating a new image

Image-to-Image Translation with Conditional Adversarial Nets

Philip Isola Jun-Yan Zhu Tinghui Zhou Alexei A. Efros

University of California, Berkeley
In CVPR 2017

[Paper] [GitHub]



Example results on several image-to-image translation problems. In each case we use the same architecture and objective, simply training on different data.

Abstract

We investigate conditional adversarial networks as a general-purpose solution to image-to-image translation problems. These networks not only learn the mapping from input image to output image, but also learn a loss function to train this mapping. This makes it possible to apply the same generic approach to problems that traditionally would require very different loss formulations. We demonstrate that this approach is effective at synthesizing photos from label maps, reconstructing objects from edge maps, and colorizing images, among other tasks. As a community, we no longer hand-engineer our mapping functions, and this work suggests we can achieve reasonable results without hand-engineering our loss functions either.



❖ Generating a new image



사진이나 신발의 아웃라인을 그리면 딥러닝을 이용한 모형이 알아서 색을 칠해줄 것입니다.