

# Convolutional neural network

손경아

아주대학교

# Inductive bias

- Two components for solving a task
  - Priors (bias): things assumed beforehand
  - Learning (variance): things extracted from data

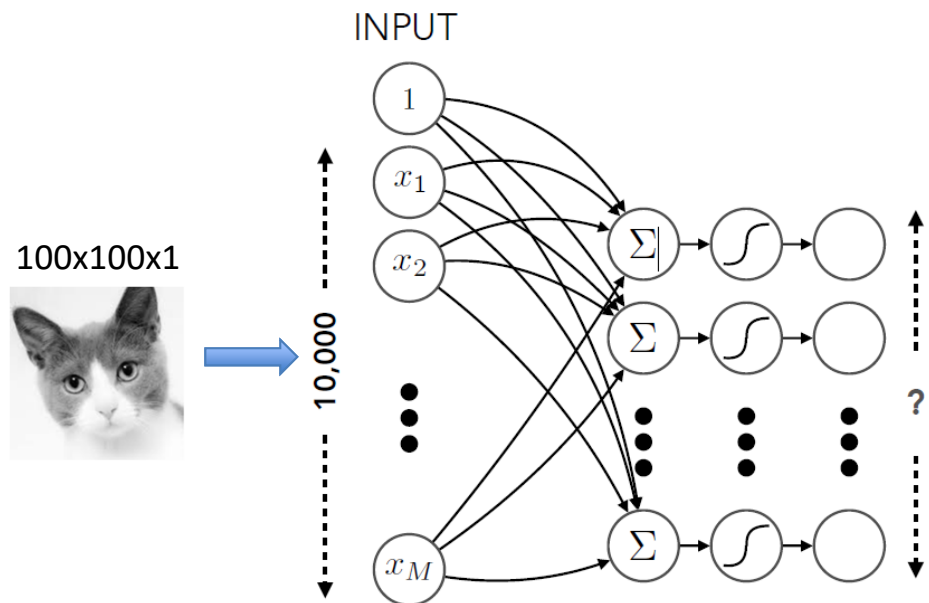
## **strong priors, minimal learning**

- fast/easy to learn and deploy
- may be too rigid, unadaptable

## **weak priors, much learning**

- slow/difficult to learn and deploy
- flexible, adaptable

# Neural network with images



**How many parameters?**

#units x 10,000 = # weights

1 -----→ 10,000

10 -----→ 100,000

100 -----→ 1,000,000

...

100,000 -----→ 1,000,000,000

# Convolutional networks

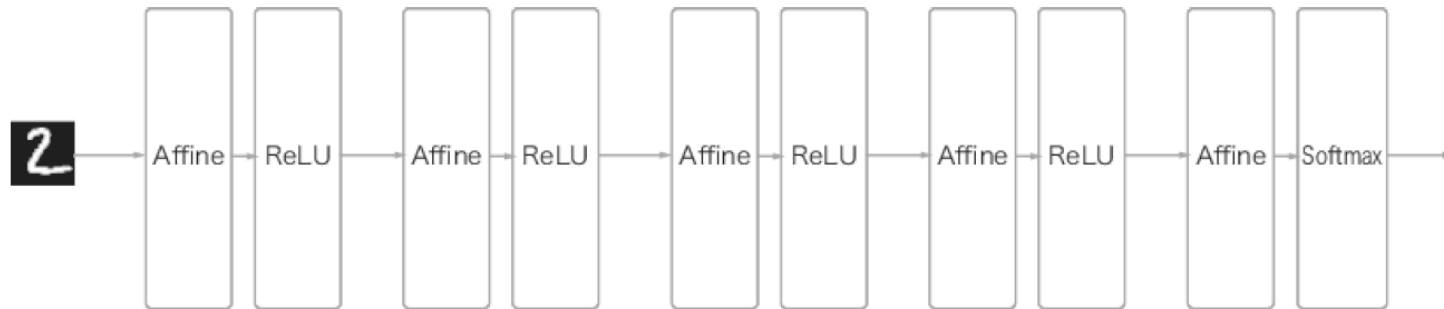
- Scale up neural networks to process very large images / video sequences
  - Sparse connections
  - Parameter sharing
- Automatically generalize across spatial translations of inputs
- Applicable to any input that is laid out on a grid (1-D, 2-D, 3-D, ...)

# Key idea

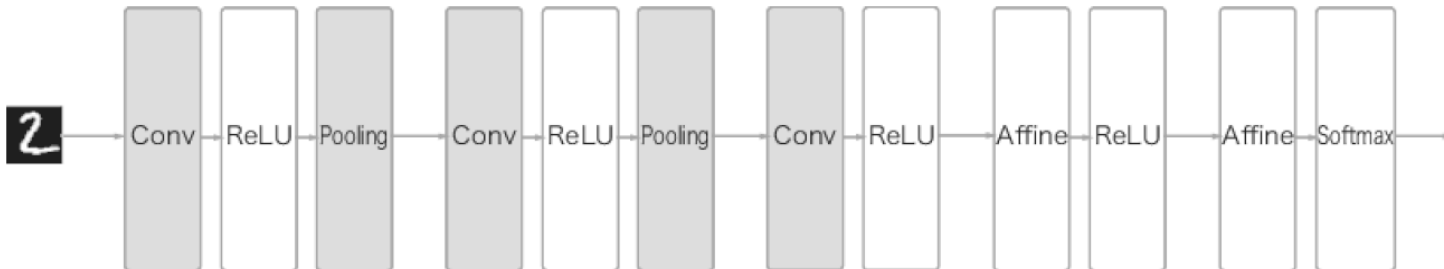
- Replace matrix multiplication in neural nets with [convolution](#)
- Everything else stays almost the same

# Convolutional networks

Network with Fully Connected (FC) layers



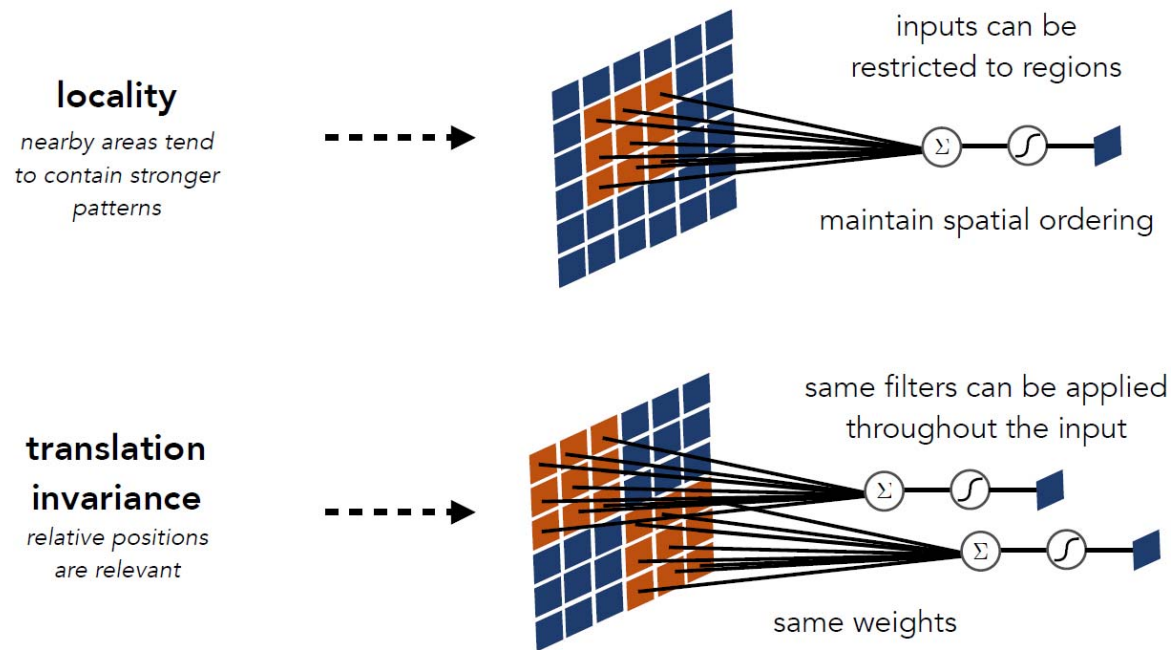
Network with convolutional layers



# Spatial structure in images

- Locality
  - Nearby areas tend to contain stronger patterns
- Translation invariance
  - Relative (rather than absolute) positions are relevant

# Exploit spatial structure in images

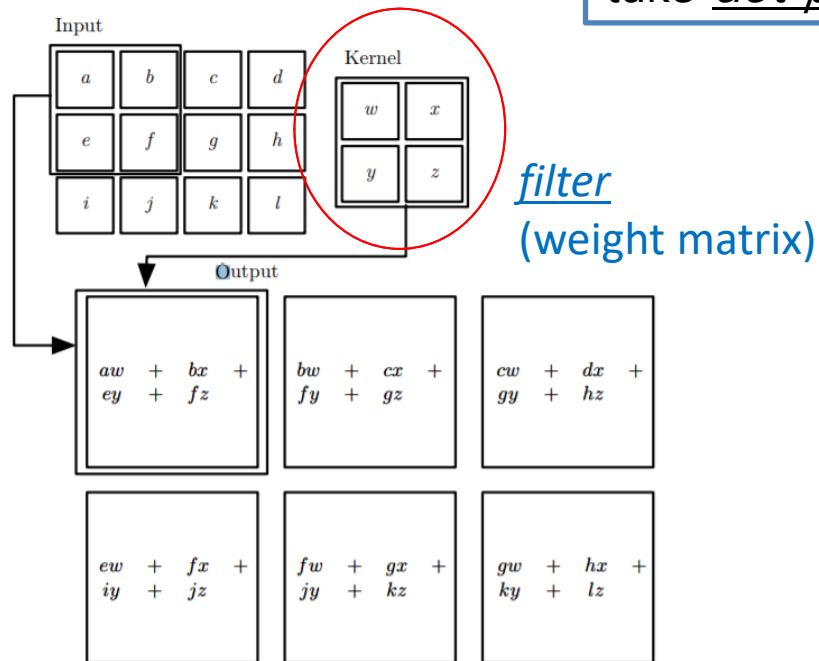




# Convolution

- 2D convolution

take dot product of filter and each input location



filter weights:  $\begin{pmatrix} 0 & 1 & 2 \\ 2 & 2 & 0 \\ 0 & 1 & 2 \end{pmatrix}$

3 <sub>0</sub>	3 <sub>1</sub>	2 <sub>2</sub>	1	0
0 <sub>2</sub>	0 <sub>2</sub>	1 <sub>0</sub>	3	1
3 <sub>0</sub>	1 <sub>1</sub>	2 <sub>2</sub>	2	3
2	0	0	2	2
2	0	0	0	1




- 2D convolution

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

 $\otimes$ 

2	0	1
0	1	2
1	0	2

 $\rightarrow$ 

15	

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

 $\otimes$ 

2	0	1
0	1	2
1	0	2

 $\rightarrow$ 

15	16

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

 $\otimes$ 

2	0	1
0	1	2
1	0	2

 $\rightarrow$ 

15	16
6	

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

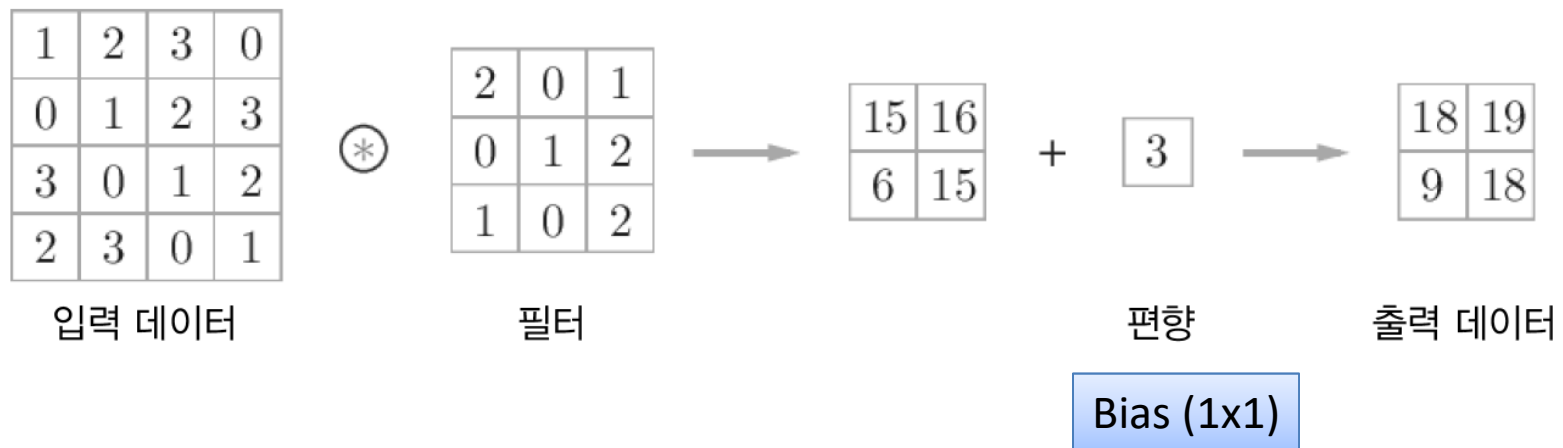
 $\otimes$ 

2	0	1
0	1	2
1	0	2

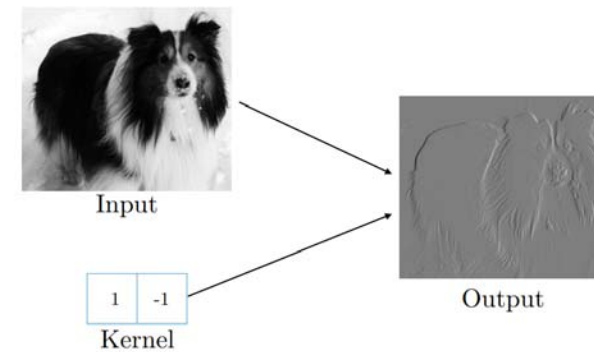
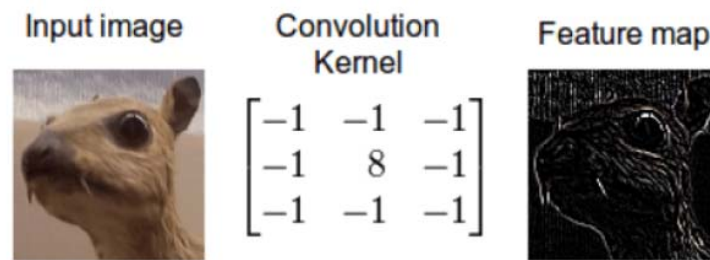
 $\rightarrow$ 

15	16
6	15

- Convolution with bias



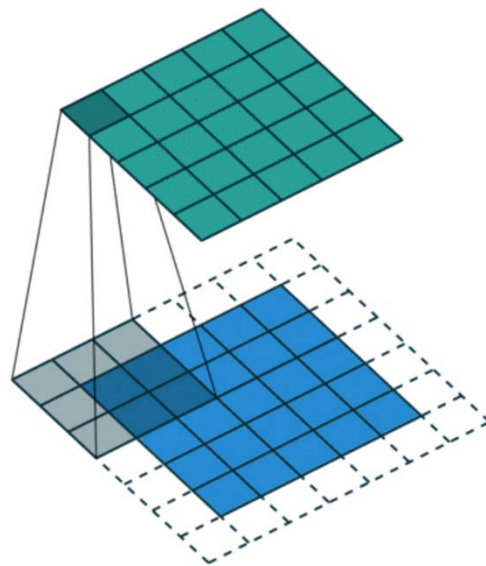
# Convolution example



- Acts like a high-pass filter: pixels near edges survive
- Can be performed very efficiently on modern libraries and GPUs

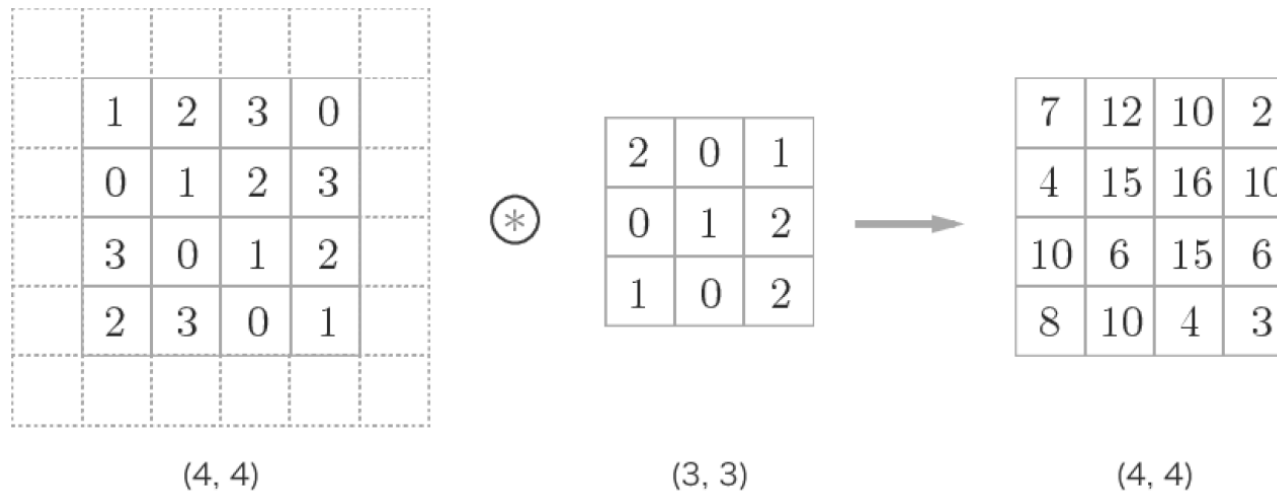
# Padding

- Use padding to preserve spatial size

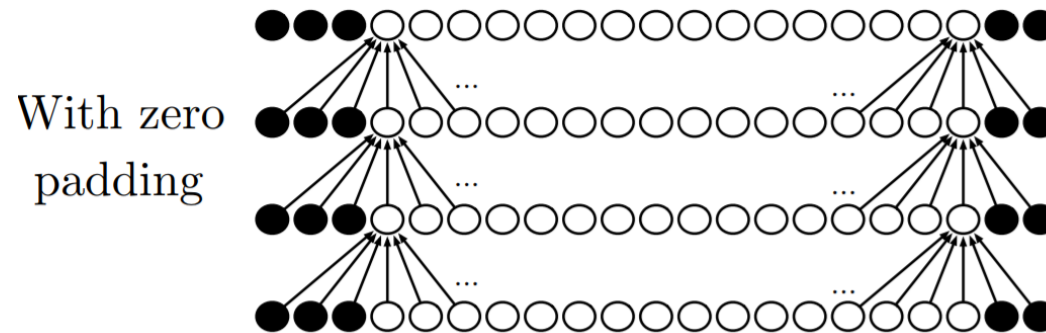
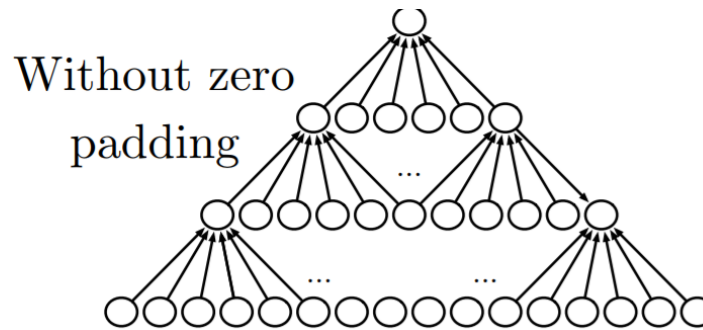


typically add *zeros* around  
the perimeter

# Zero-padding



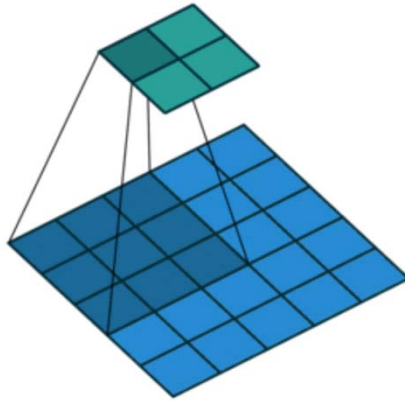
# Zero padding controls size



# Stride

- Use stride to downsample the input

stride = 2



Move 2 pixels at a time



# Stride

stride = 2

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

⊗

2	0	1
0	1	2
1	0	2



15		

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

⊗

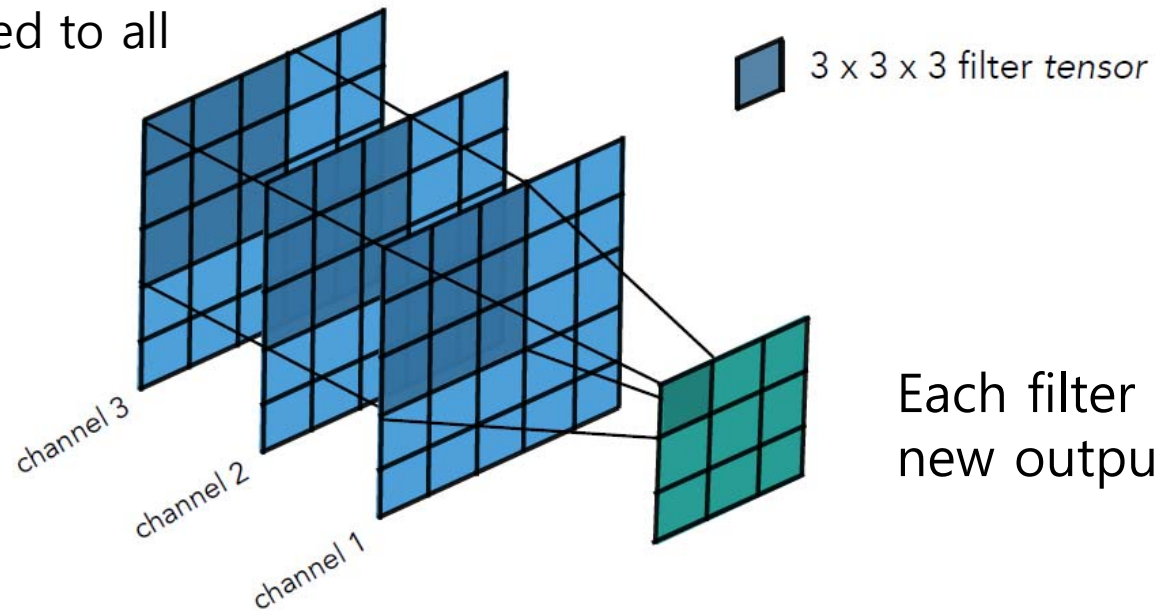
2	0	1
0	1	2
1	0	2



15	17	

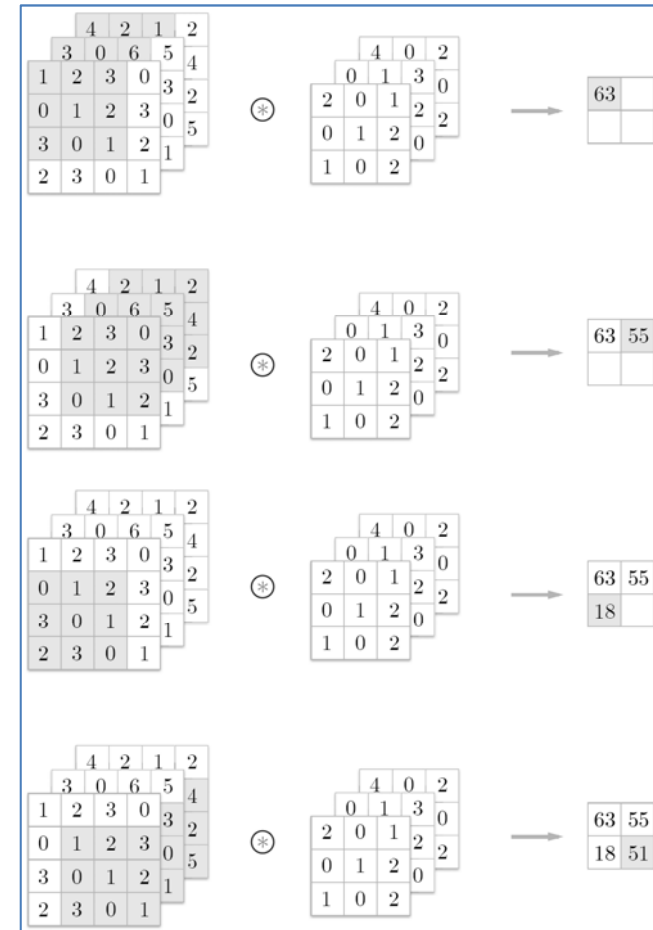
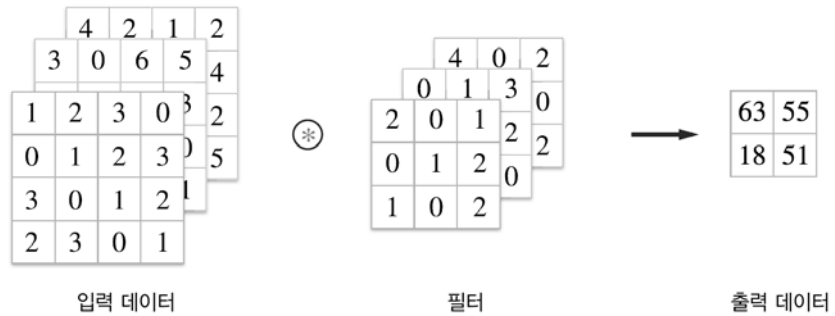
# Multiple channels

Filters are applied to all input channels

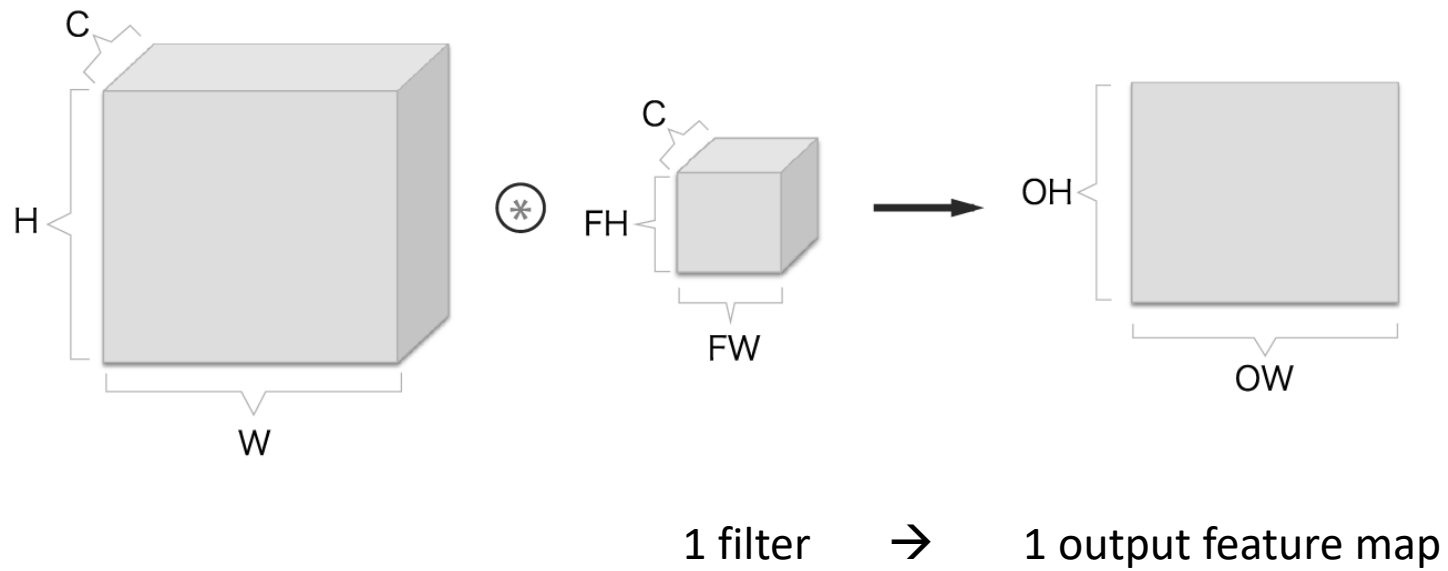


Each filter results in a new output channel

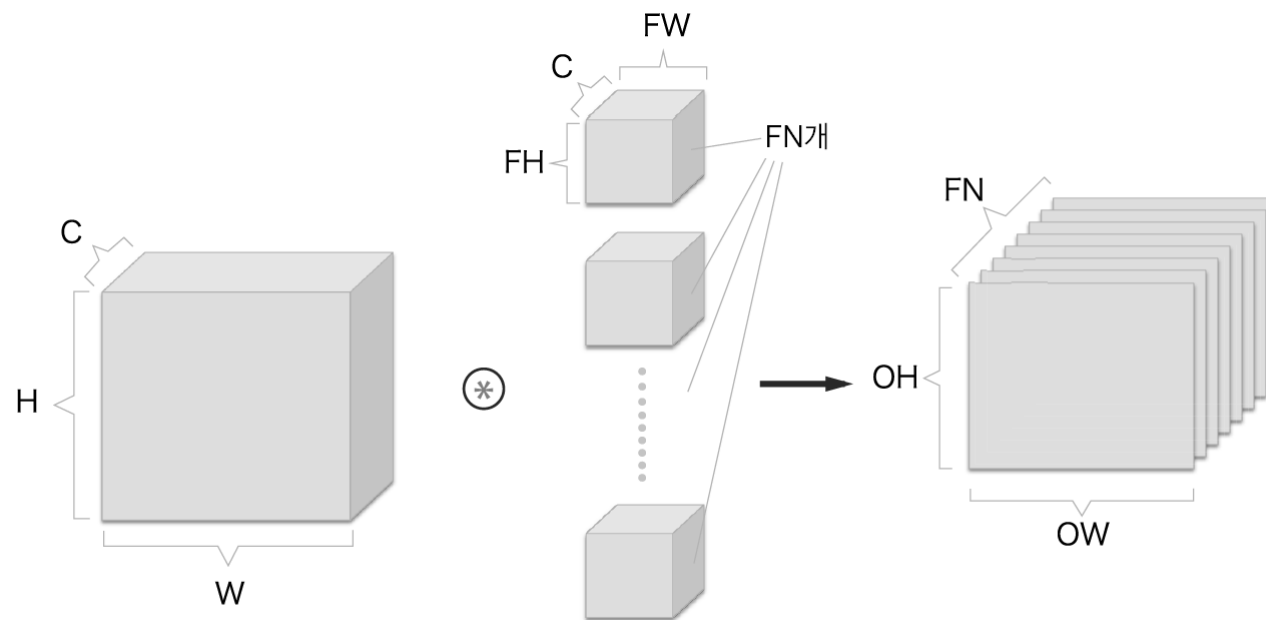
# 3D convolution

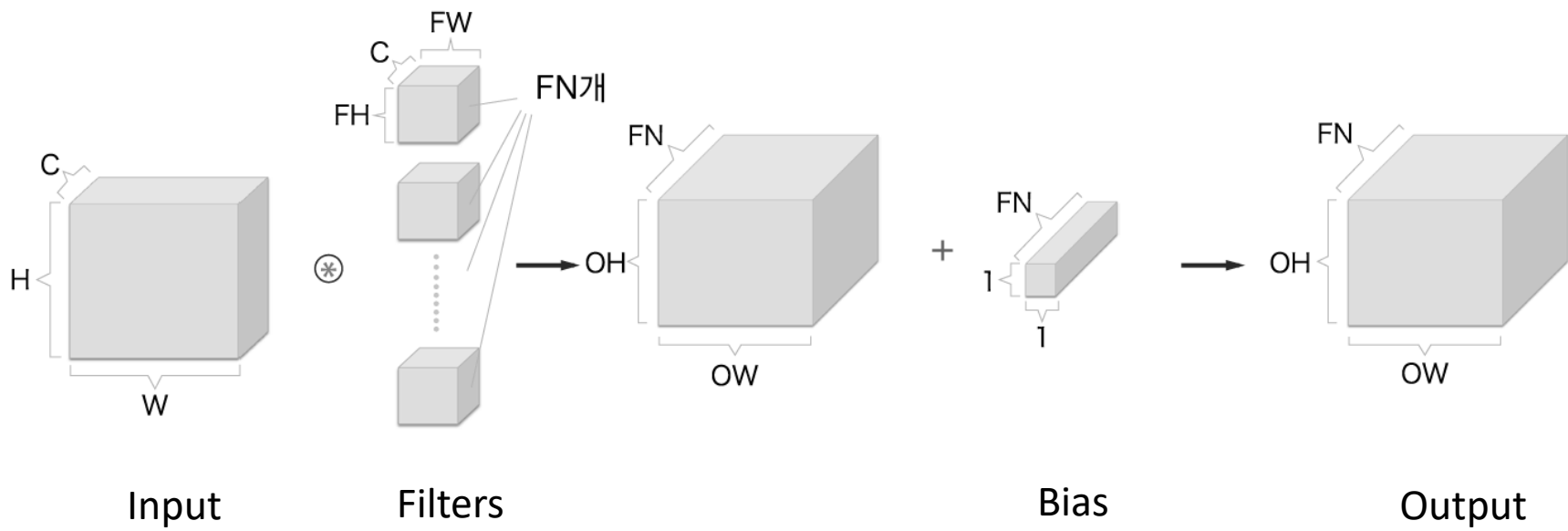


# Convolution



# Multiple filters

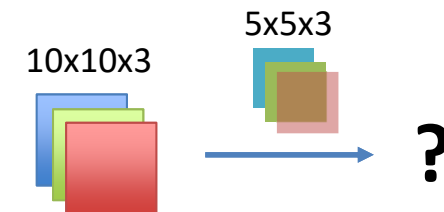




# Pop quiz 1

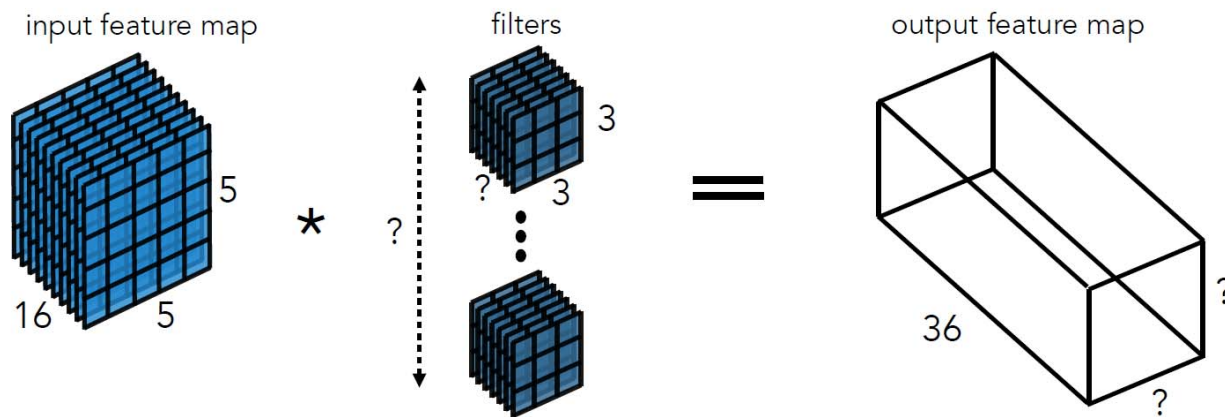
# input channels = 3, each 10 x 10

Kernel (filter) of size 5 x 5 x 3 is applied with stride 1



- What is the dimension of the output channels?
- What if the input maps are zero-padded?

# Pop quiz 2

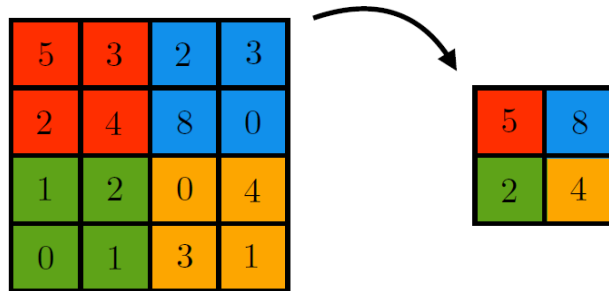


- If we use unit (1) stride and no padding,
  - What is each filter size?
  - How many filters are there?
  - What is the output filter map size?

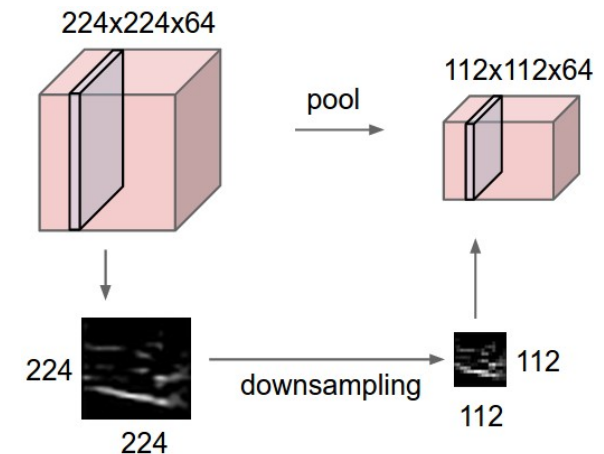


# Pooling

- Locally aggregates values in each feature map

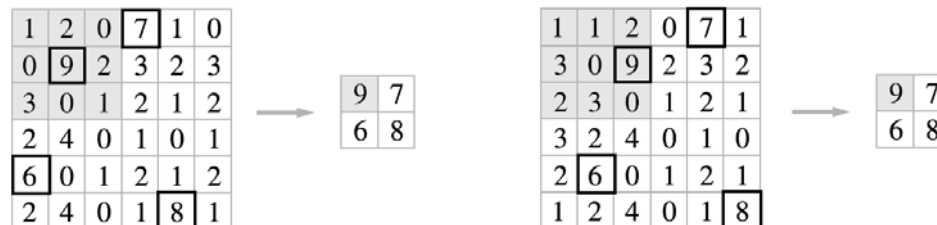
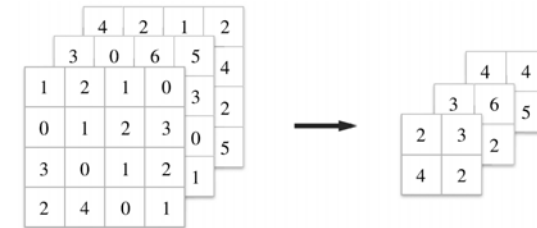


predefined operation: maximum, average, etc.

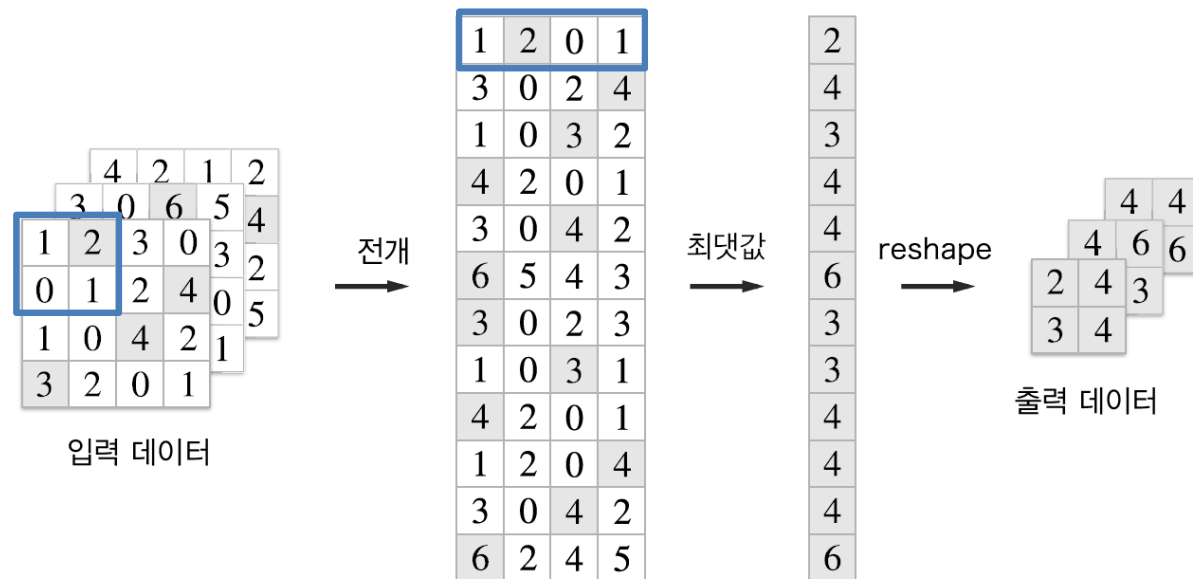


# Pooling layer

- No parameter to learn
- The number of channels does not change
- Robust to input variation

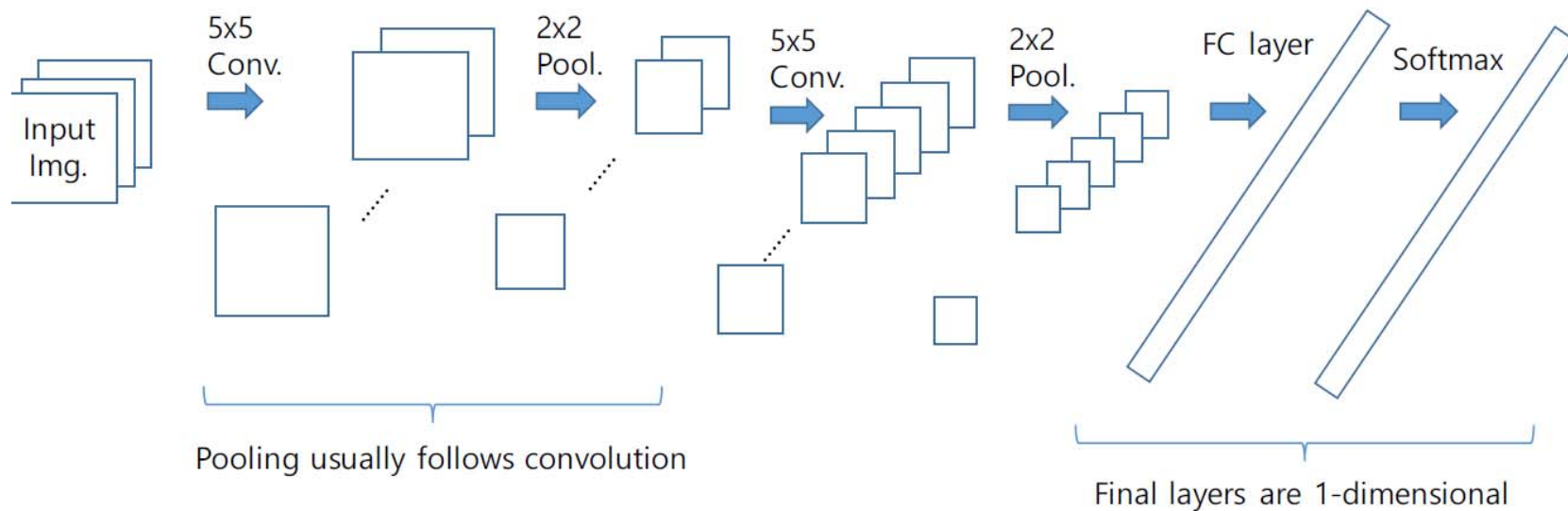


# Implementation: max pooling



# Convolutional Neural network

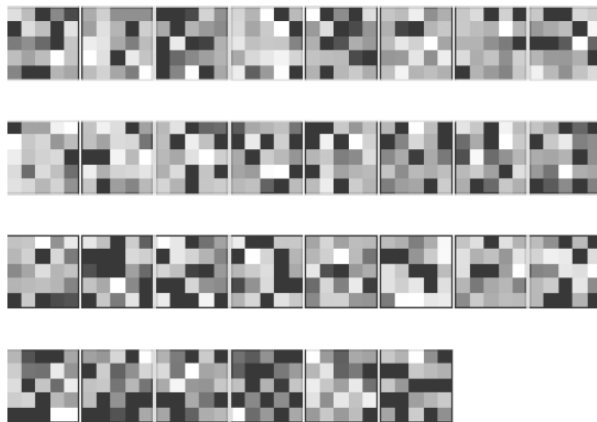
- A deep neural network composed of convolution-pooling layers



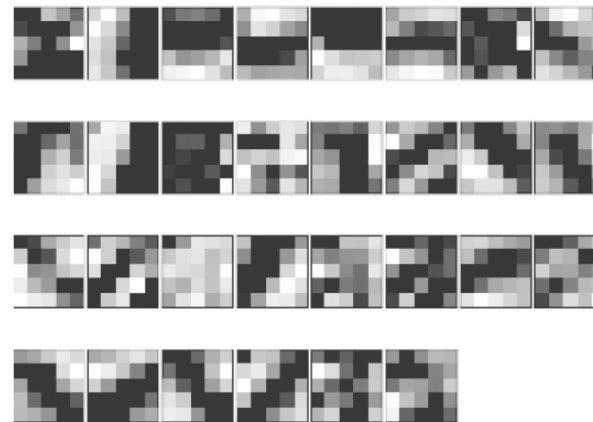
# Filter visualization

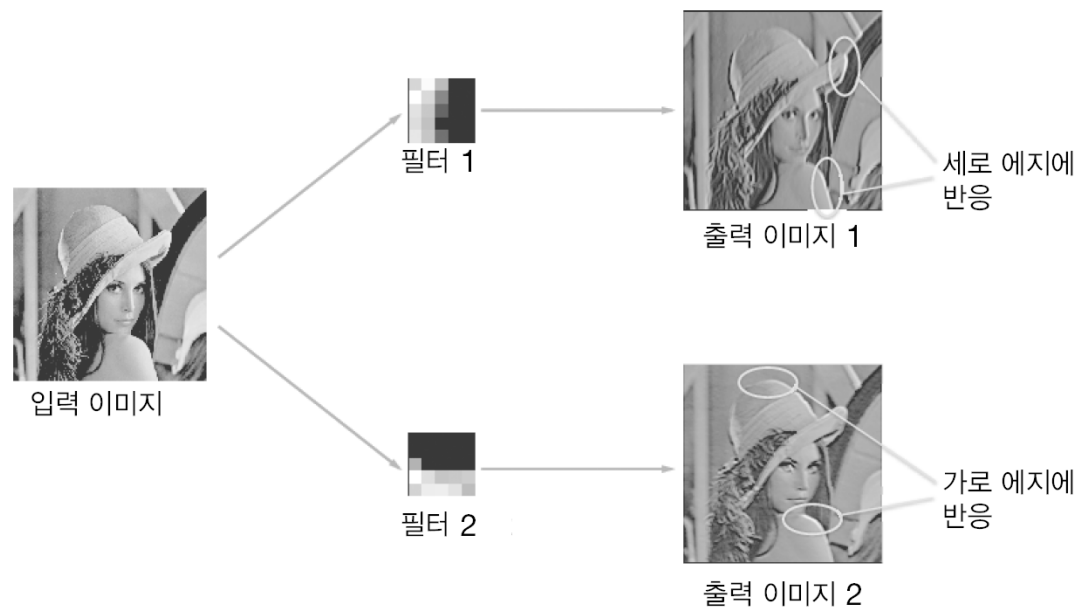
- Filters at 1<sup>st</sup> layer before and after training

학습 전



학습 후



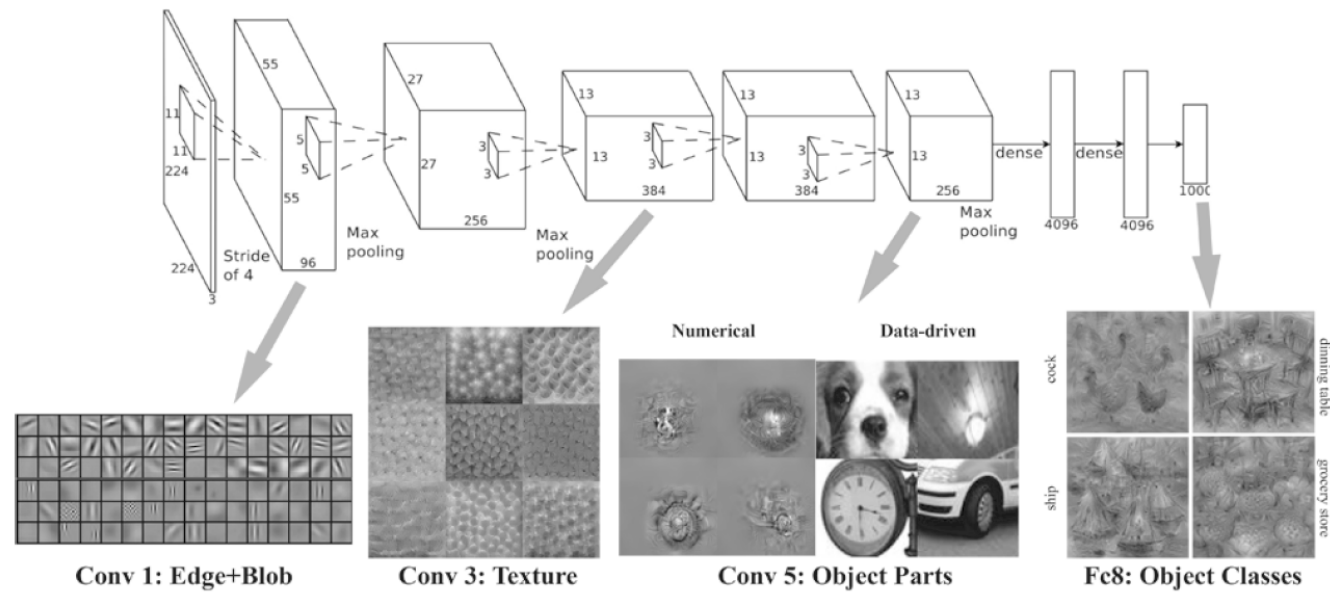


# Filter visualization



Each of the 96 filters is of size  $[11 \times 11 \times 3]$  and shared by the neurons in one depth slice

# Information extracted from each layer

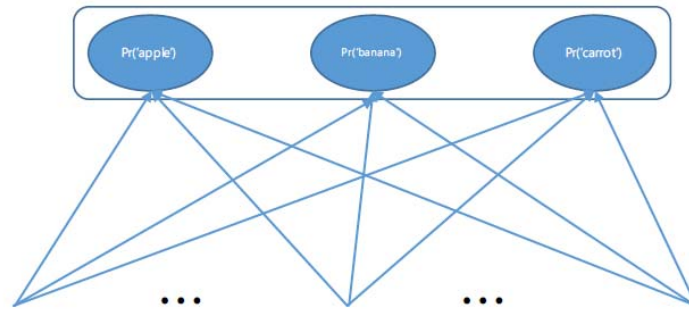




# Classification layer

- Many CNNs are used to perform classification
- The final layer of such CNN is used to represent the class probabilities

Prediction classes: {'apple', 'banana', 'carrot'}



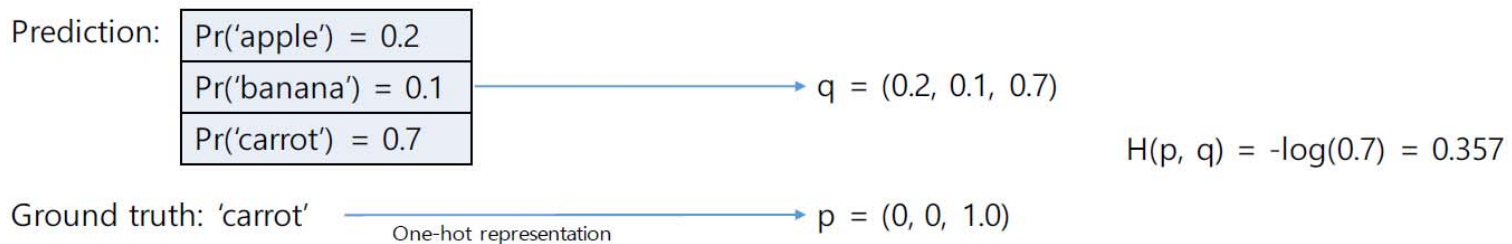
"Softmax" layer

$$\text{Pr}(o) = \frac{e^{a_o}}{\sum_{a'} e^{a'}}$$

# Classification layer

- Since the softmax layer specifies the probability, the cost should be computed differently
- Cross entropy: information-theoretic measure of difference between two distributions






$$H(p, q) = - \sum_x p(x) \log q(x)$$



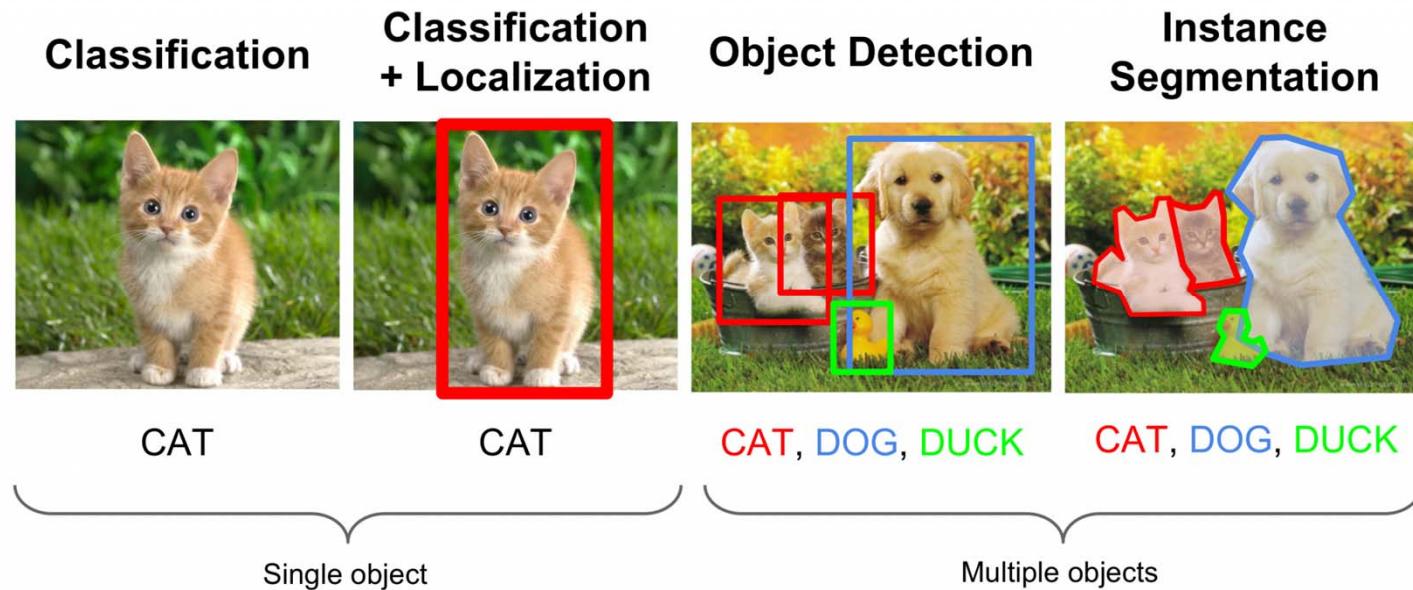
# Training CNNs

- Back-propagation
  - SGD, Adam, AdaDelta, AdaGrad, etc.
- (Trainable) parameters:
  - Kernel patch for each layer
  - bias vectors for each kernel
- (non-trainable) hyperparameters
  - Network depth
  - Number of kernels per layer, kernel size, stride, zero-padding
  - Which activation functions

# Natural image datasets

					...
<b>Caltech-101</b>	<b>Caltech-256</b>	<b>CIFAR-10</b>	<b>CIFAR-100</b>	<b>ImageNet</b>	
				<b>Competition</b>	<b>Full</b>
101 classes, 9,146 images	256 classes, 30,607 images	10 classes, 60,000 images	100 classes, 60,000 images	1,000 classes, 1.2 million images	21,841 classes, 14 million images

# Computer vision tasks



# Image classification



[This image is CC0 public domain](#)

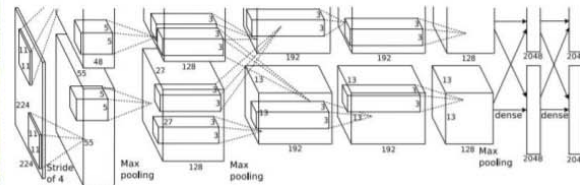


Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

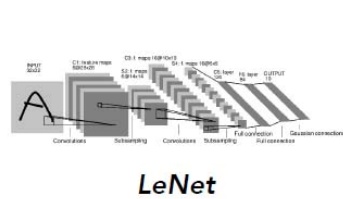
**Vector:**  
4096

→  
**Fully-Connected:**  
4096 to 1000

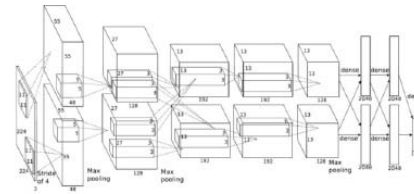
## Class Scores

Cat: 0.9  
Dog: 0.05  
Car: 0.01  
...

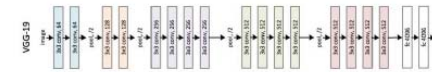
# CNNs for classification



LeNet



AlexNet



VGG



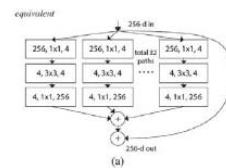
GoogLeNet



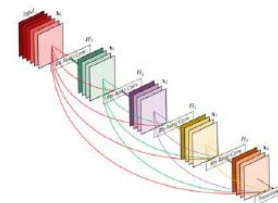
ResNet



Inception v4



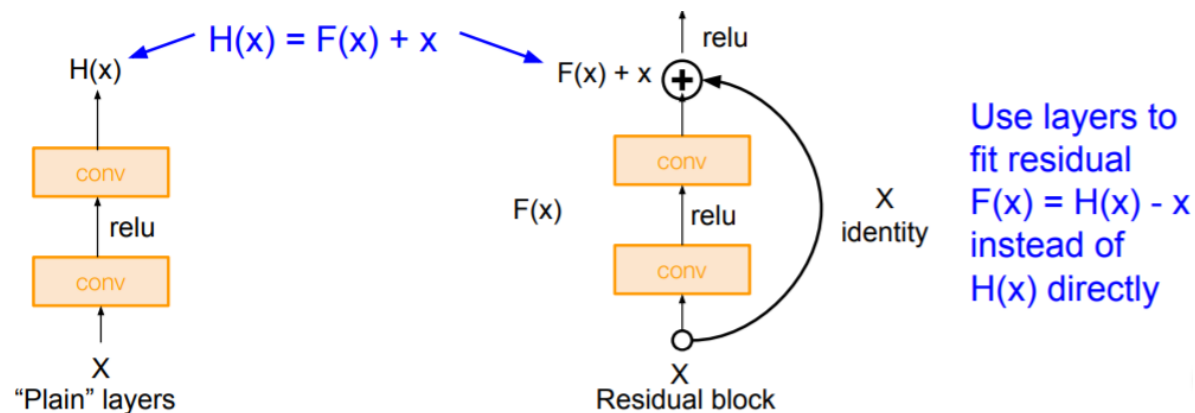
ResNeXt



DenseNet

# ResNet [He et al., 2015]

- Very deep networks using residual connections

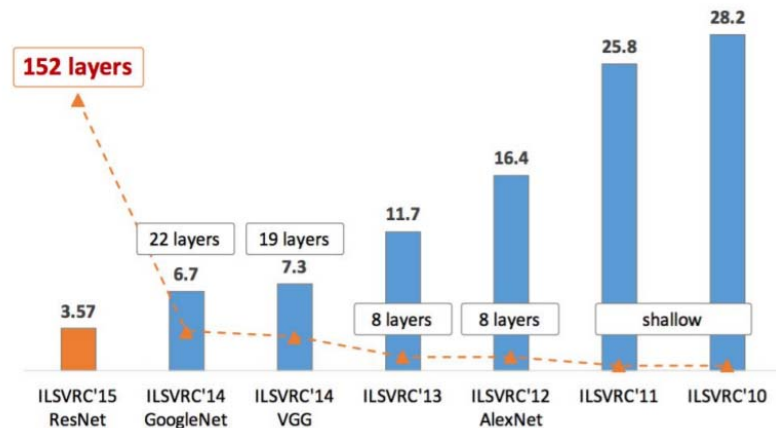




# ResNet [He et al., 2015]

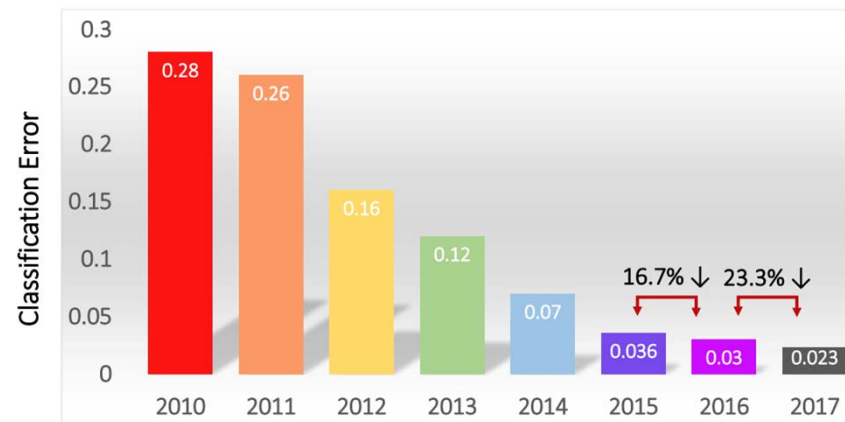
- Experiment
  - Able to train very deep networks without degrading (152 layers on ImageNet, 1202 on Cifar)
  - Deeper networks now achieve lower training error as expected
  - Swept 1st place in all ILSVRC and COCO 2015 competitions al results

ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



# ILSVRC (2010-2017)

## Classification Results

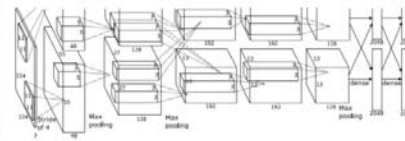


## Localization Results

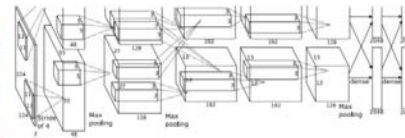


# Object detection

- As regression?

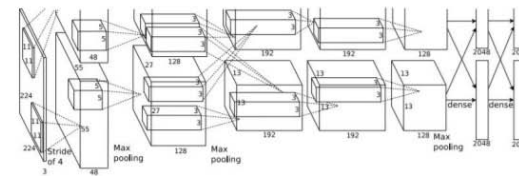


CAT: (x, y, w, h) 4 numbers



DOG: (x, y, w, h)  
DOG: (x, y, w, h) 16 numbers  
CAT: (x, y, w, h)

- As classification:  
Sliding window?

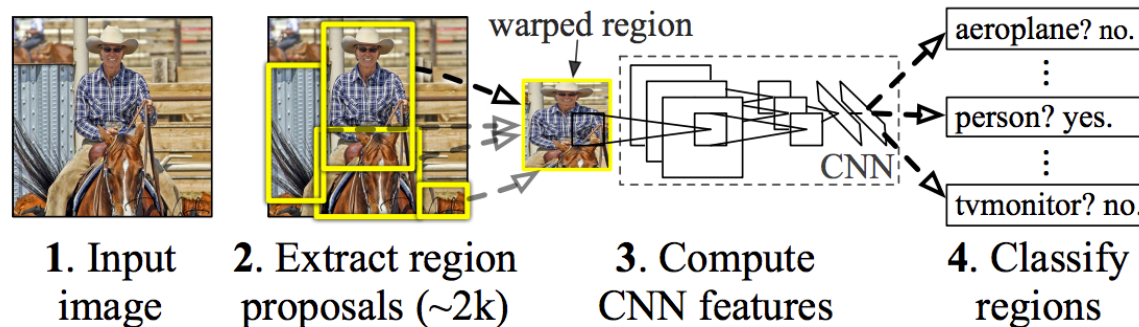


Dog? NO  
Cat? YES  
Background? NO

Problem: Need to apply CNN to huge number of locations and scales, very computationally expensive!

# R-CNN (2014): an early application of CNNs to Object Detection

- Inputs: Image
- Outputs: Bounding boxes + labels for each object in the image.

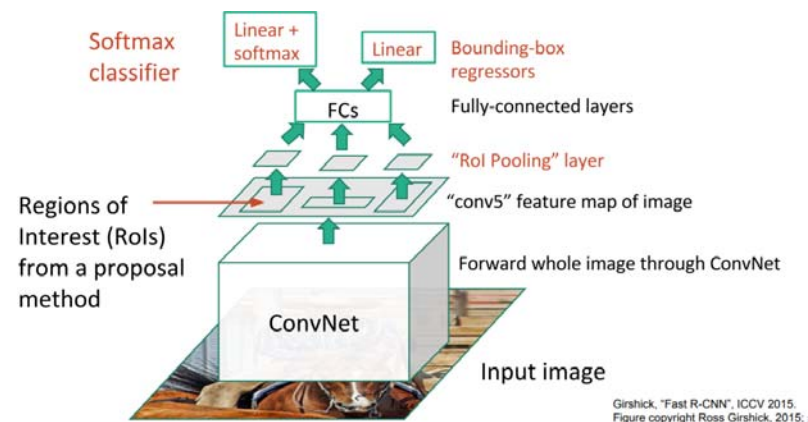


Source: <https://arxiv.org/abs/1311.2524>.

1. Generate a set of proposals for bounding boxes.
2. Run the images in the bounding boxes through a pre-trained AlexNet and finally an SVM to see what object the image in the box is.
3. Run the box through a linear regression model to output tighter coordinates for the box once the object has been classified.

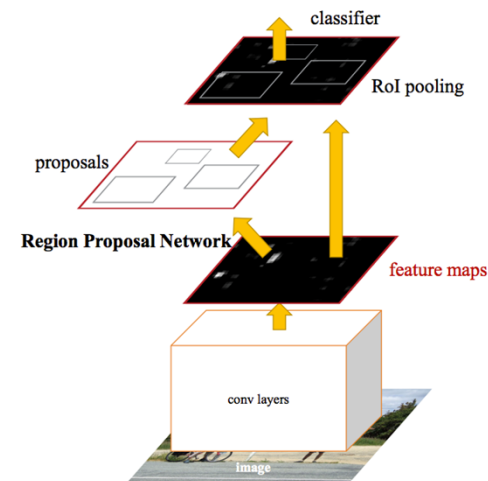
## Fast R-CNN (2015): Speeding up and Simplifying R-CNN

- Fast R-CNN combined the CNN, classifier, and bounding box regressor into one, single network



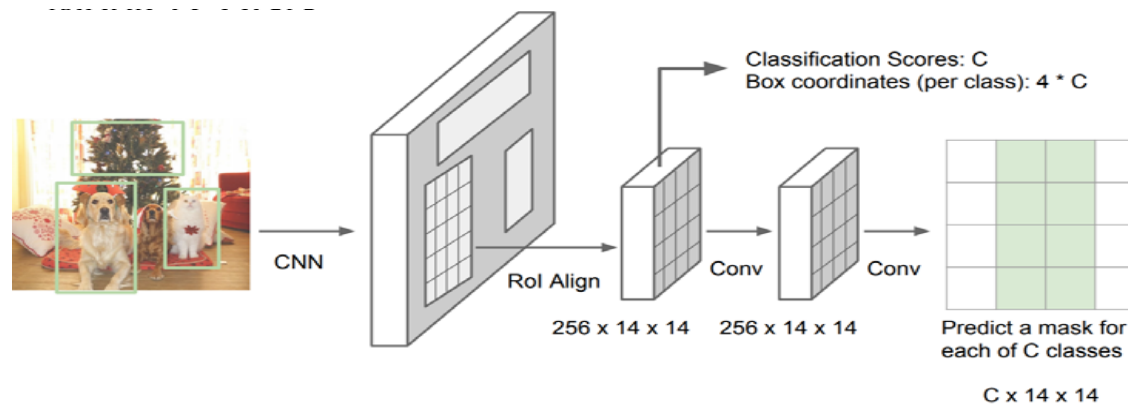
## Faster R-CNN (2016): Speeding up Region proposal

- In Faster R-CNN, a single CNN is used for region proposals, and classifications



# Mask R-CNN (2017): Extending Faster R-CNN for pixel level segmentation

- In Mask R-CNN, a Fully Convolutional Network (FCN) is added on top of the CNN features of Faster R-CNN to generate a **mask** (segmentation output)



# Mask R-CNN

Mask R-CNN is able to segment as well as classify the objects in an image.

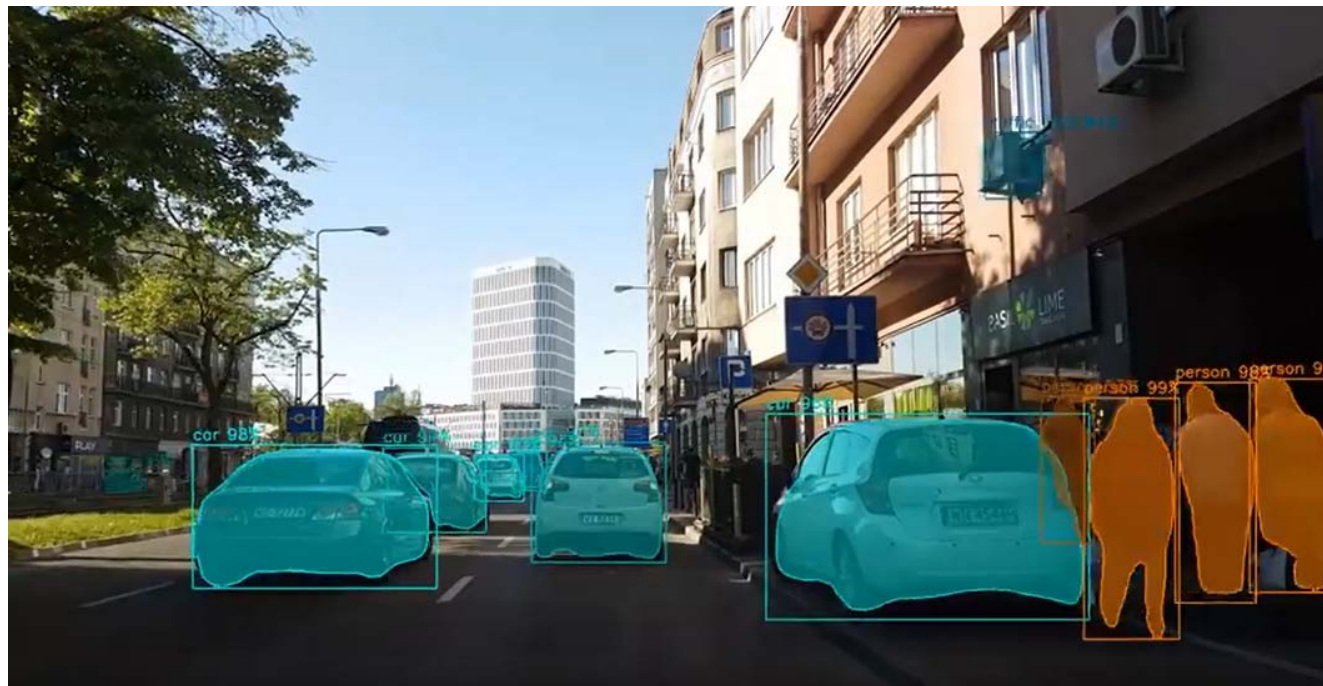


Source: <https://arxiv.org/abs/1703.06870>.



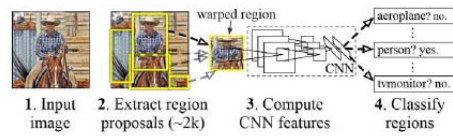
# Mask R-CNN COCO Object detection and segmentation

<https://www.youtube.com/watch?v=OOT3UIXZztE>

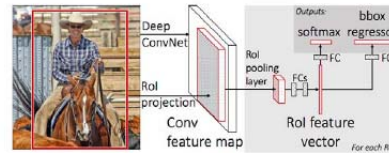




# CNNs for detection, segmentation, etc.



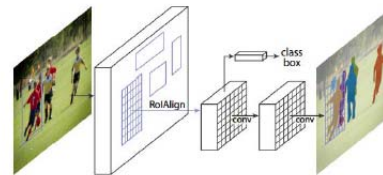
## R-CNN



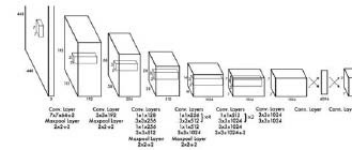
## Fast R-CNN



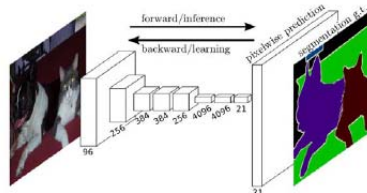
## Faster R-CNN



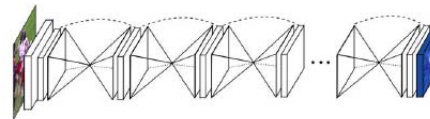
### Mask R-CNN



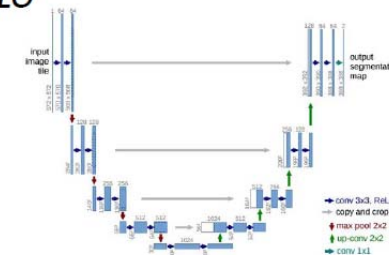
YOLO



## FCN



*Hourglass*



## U-Net

# Real-time 2D human pose estimation

<https://www.youtube.com/watch?v=pW6nZXeWlGM>



# Neural style transfer



Gatys, Ecker, and Bethge, "Image style transfer using convolutional neural networks", CVPR 2016

# Neural style transfer

- Loss function: content loss + style loss
  - Content loss = Euclidean distance between the feature map  $F$  from the content image and the one  $P$  from the generated output image
  - Style loss: take the so called “*Gram matrix*” of each feature map (from style image, and the output image), and compute their distances across layers

$$\mathcal{L}_{total} = \alpha \mathcal{L}_{content} + \beta \mathcal{L}_{style}$$

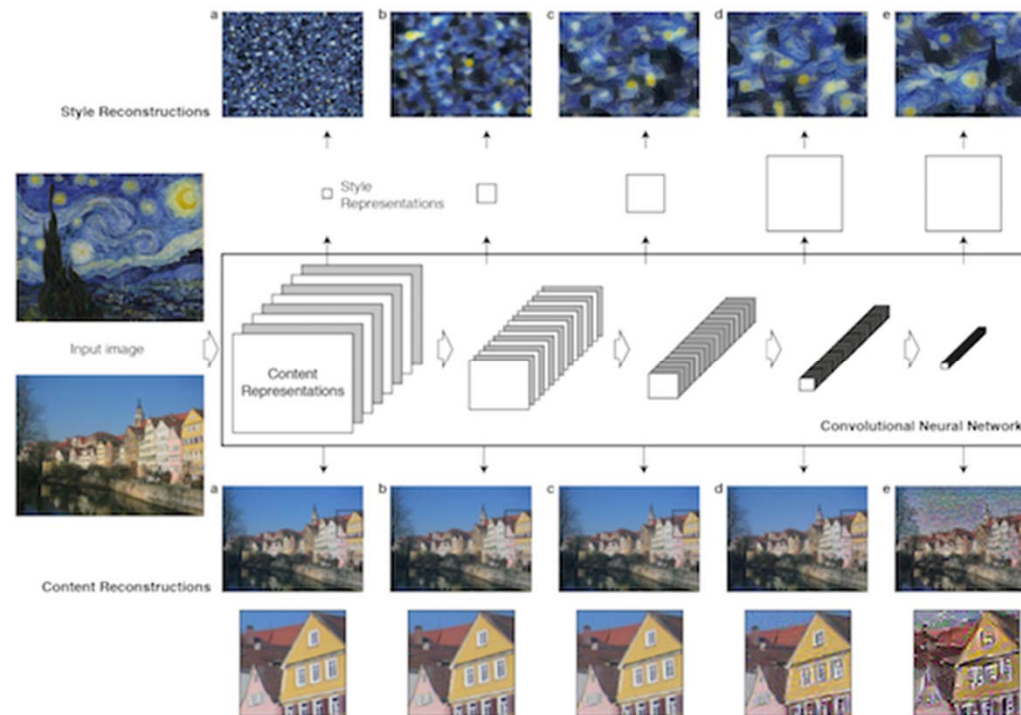


$$\mathcal{L}_{content} = \frac{1}{2} \sum_{ij} (F_{ij}^l - P_{ij}^l)^2$$

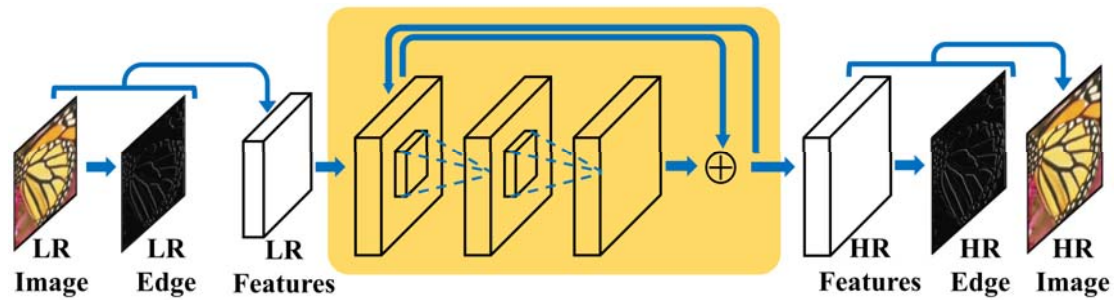
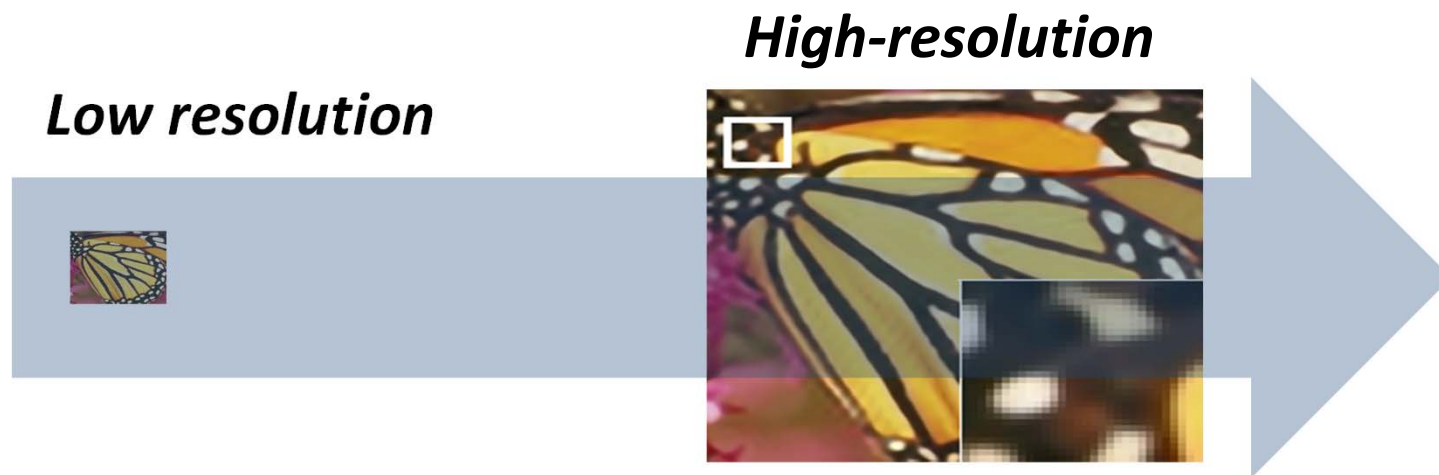


$$\mathcal{L}_{style} = \frac{1}{2} \sum_{l=0}^L (G_{ij}^l - A_{ij}^l)^2$$

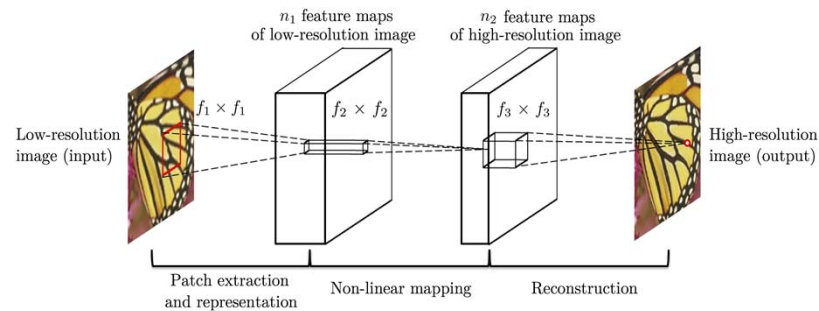
# Neural style transfer



# Image super-resolution

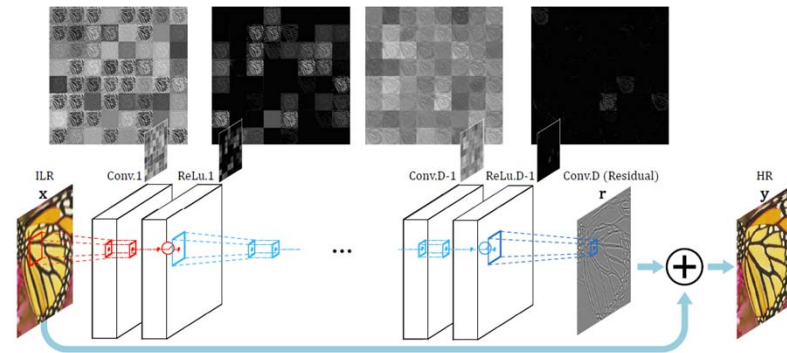


## SRCNN (ECCV 2014): Image Super-Resolution Using Deep Convolutional Networks



- Simple 3-layer CNN
  - End-to-end learning (LR - HR)
  - Convolution and ReLU

## VDSR (CVPR 2016): Accurate Image Super-Resolution Using Very Deep Convolutional Networks

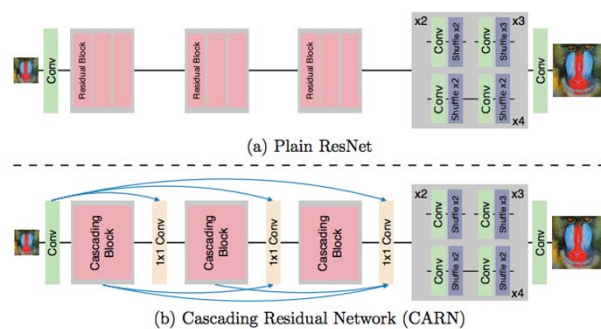


- Very deep CNN (20 layers)
  - Skip connection
  - No dimension reduction such as pooling



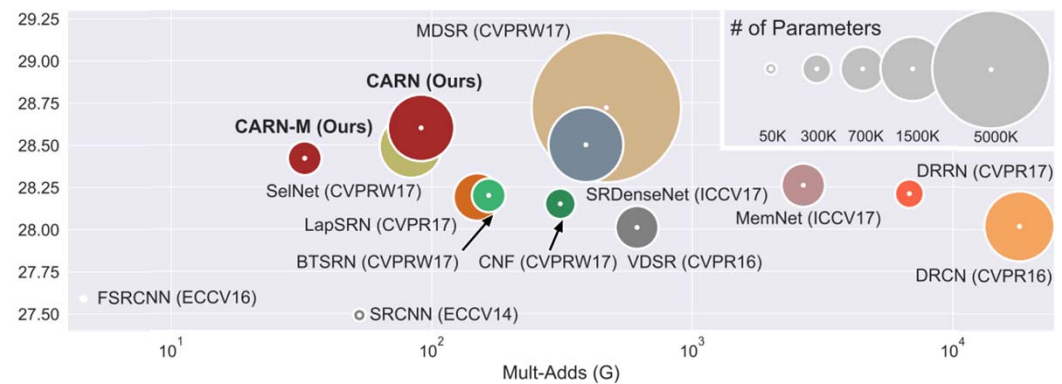
# CARN (Ahn et al. ECCV 2018): Fast, Accurate, and Lightweight Super-Resolution with Cascading Residual Network

- Deep learning-based super-resolution
  - (+) Accurate
  - (-) Slow, (-) Heavy
- To effectively perform super-resolution, we add many short-cut connections



See how the PSNR (y-axis) varies depending on:

- computational cost (x-axis) and
- model size (circle area)





# CNN baseline

- Download a pre-trained network
- Or copy-paste an architecture from a related task
- Or:
  - Deep residual network
  - Batch normalization
  - Adam

# Reference

- Books
  - Deep Learning, Ian Goodfellow, Yoshua Bengio, Aaron Courville
  - 밑바닥부터 시작하는 딥러닝, 사이토 고키
- Lecture notes
  - Stanford cs231n
  - Caltech CS/CNS/EE155