

# 2. Convolutional Neural Networks

Donghwan Chi\*

Artificial Intelligence in Korea University(AIKU)

Department of Computer Science and Engineering, Korea University

# Recap : Dive into Deep Learning



Yoshua Bengio



Geoffrey Hinton



Yann LeCun

출처 : <https://awards.acm.org/about/2018-turing>

# Why deep learning?

---

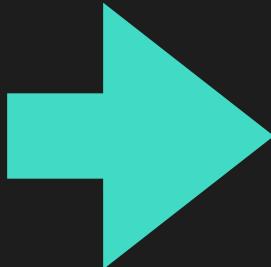
- 기존 방식들은 이미 충분히 강력하다.
  - 현실 세계의 많은 문제들이 통계학, ML로 해결되고, 충분하다.
  - 딥러닝 ‘향’ 첨가 프로젝트, 아이템들…
- 그러면 새롭게 등장한 DL이 왜 혁신적인 방법인가?
  - 그 중에서도 왜 CV, NLP가 제일 🔥한 분야인가?

# Image Classification

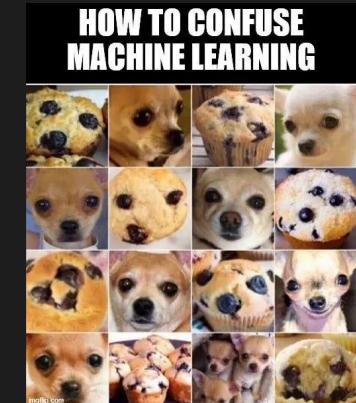


출처 : <https://medium.com/anubhav-shrima/dogs-vs-cats-image-classification-using-resnet-d2ed7e6db2bb>

# Why image classification is hard task?



```
>>> kitten.shape  
(484, 306, 4)  
>>> print(kitten[100:,59:,:,2])  
[[119 125 130 ... 255 255 255]  
 [116 122 129 ... 255 255 255]  
 [114 120 130 ... 255 255 255]  
 ...  
 [255 255 255 ... 255 255 255]  
 [255 255 255 ... 255 255 255]  
 [255 255 255 ... 255 255 255]]
```

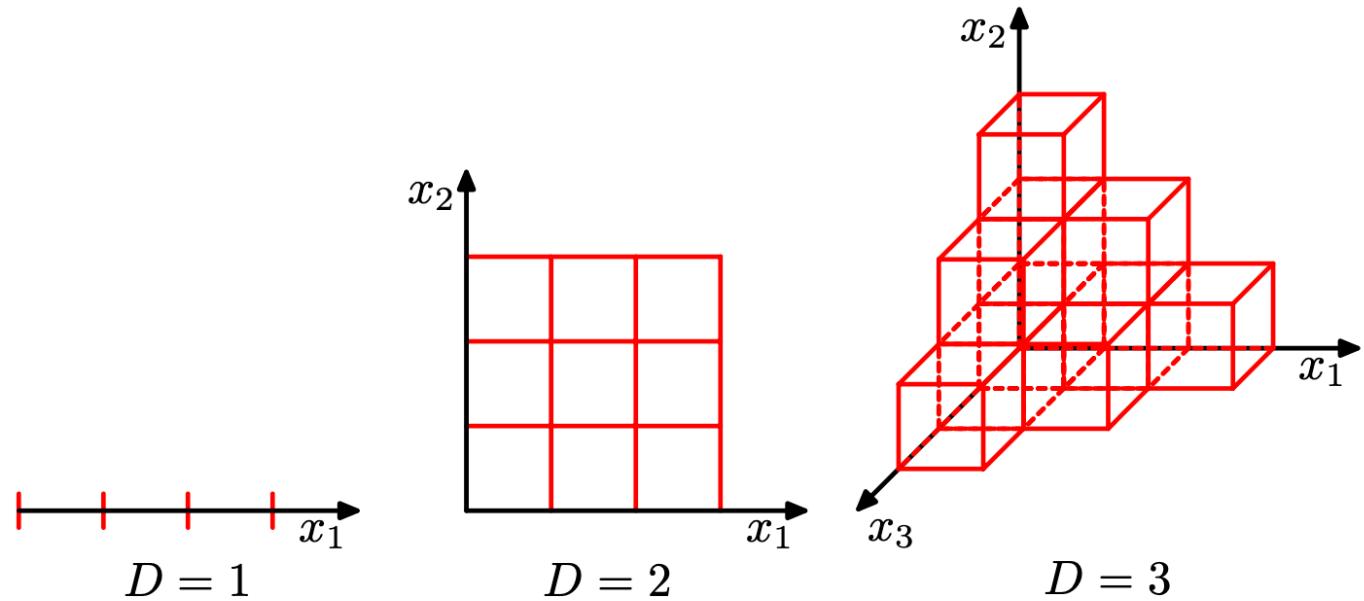


- 현실의 데이터는 고차원이다!

[출처](#)

# Curse of dimensionality

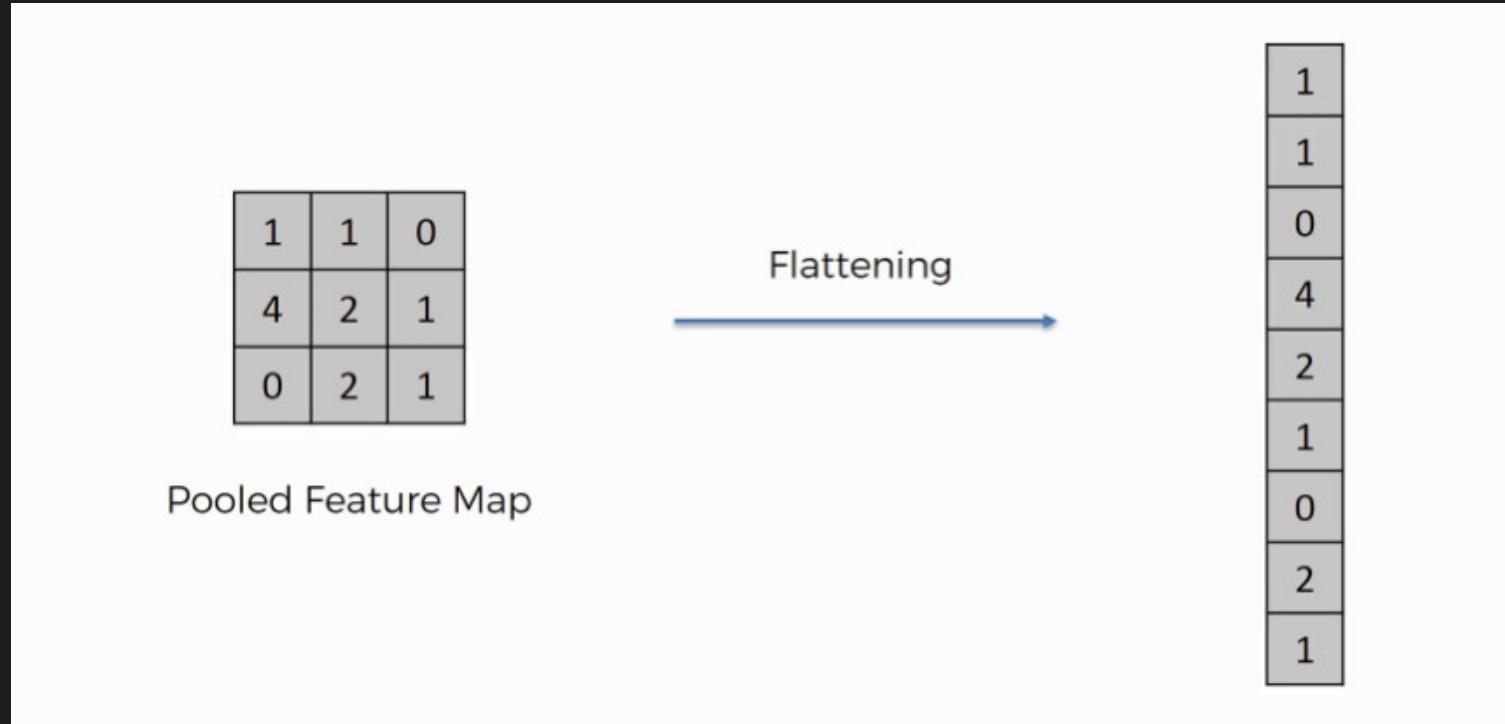
**Figure 1.21** Illustration of the curse of dimensionality, showing how the number of regions of a regular grid grows exponentially with the dimensionality  $D$  of the space. For clarity, only a subset of the cubical regions are shown for  $D = 3$ .



출처 : Pattern Recognition and Machine Learning

- 차원이 커질 수록, 문제의 계산법이 지수적으로 커진다.

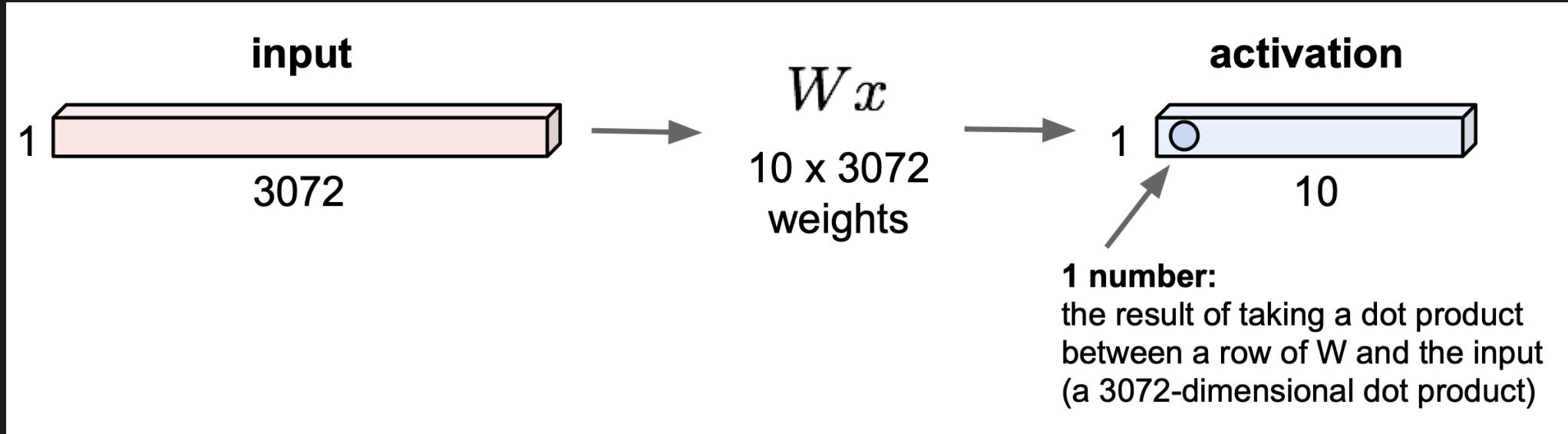
# Why not ‘MLP’?



출처 : <https://www.superdatascience.com/blogs/convolutional-neural-networks-cnn-step-3-flattening>

- MLP는 이미지를 ‘flatten’ 해야 한다.
  - 구조적인 정보가 파괴됨.

# Why not ‘MLP’?



출처 : CS231n (2017) Lecture 5

- MLP의 Parameters 수
  - CIFAR-10 :  $32 \times 32 \times 3 = 3072$  weights, FCN도 가능
  - 하지만 더 큰 image라면? → 과도한 Parameter → Overfitting

# Contents

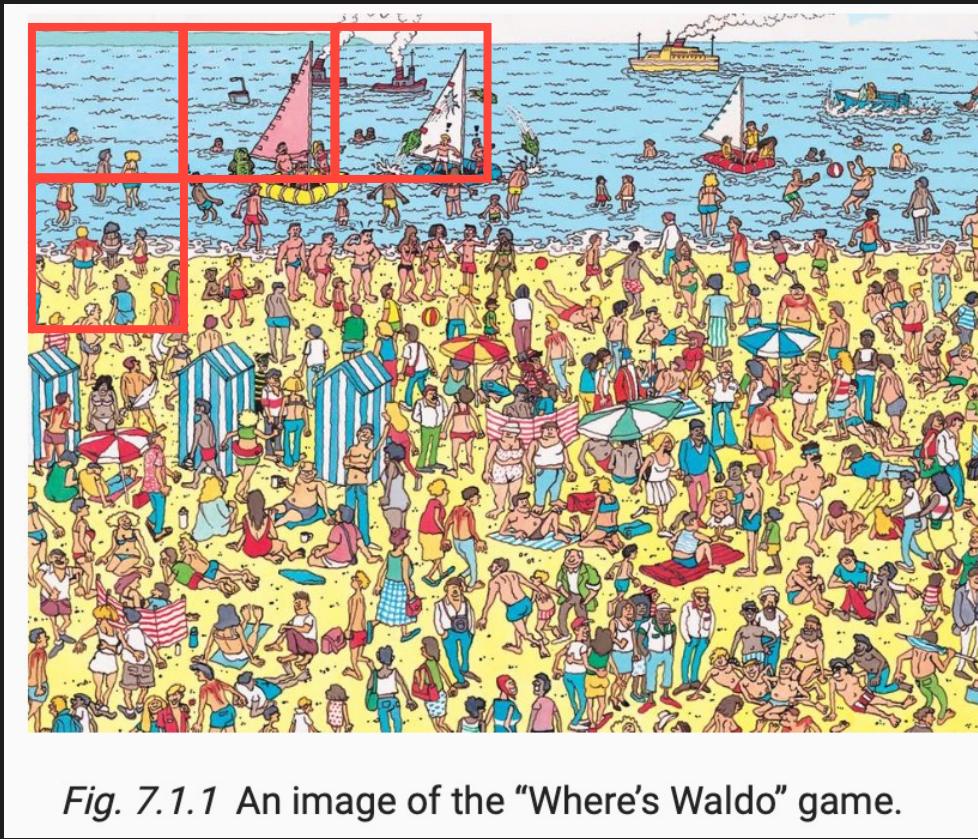
---

- Architecture Overview
- Convolutional Layer
  - Convolution
  - Spatial arrangement
  - Feature Map and Receptive Field
- Pooling Layer
- Case Study : Modern CNNs
- PyTorch hands-on

# Architecture Overview

# What is ‘Convolution’?

# Convolutional Neural Networks



출처 : [Dive Into Deep Learning](#)

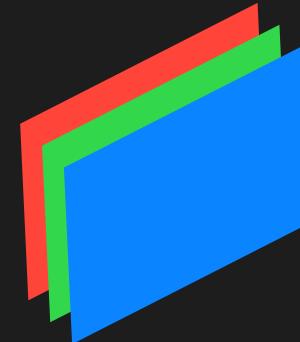
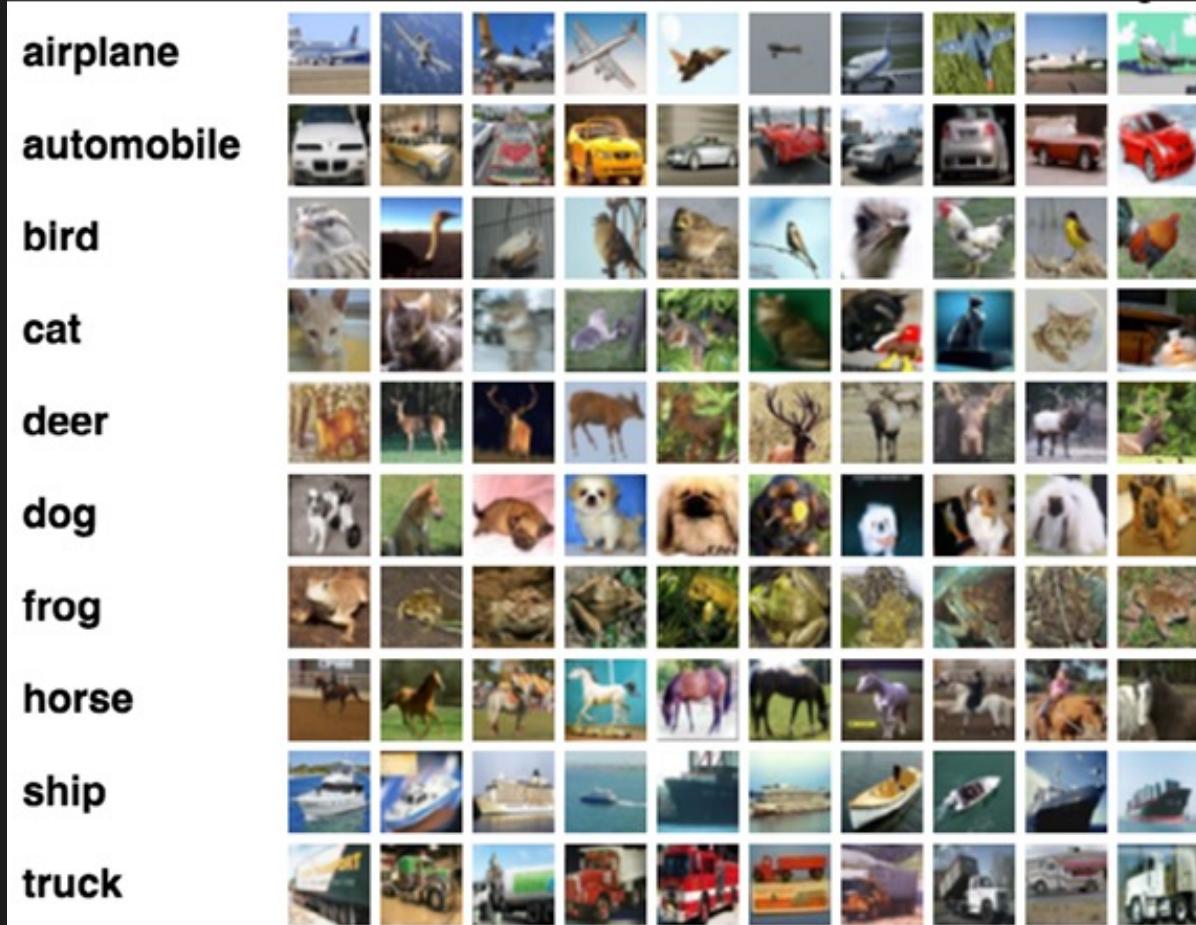
- 각 Patch에서 동일한 일을 한다.

# Architecture Overview

---

- 지금까지 배운 MLP
  - hidden layers!
- MLP의 Parameters 수
  - CIFAR-10 :  $32 \times 32 \times 3 = 3072$  weights, FCN도 가능
  - 더 큰 image라면?
- Convolutional Layer는 3D neurons을 가진다.
  - Figure에서 확인해 볼 것!

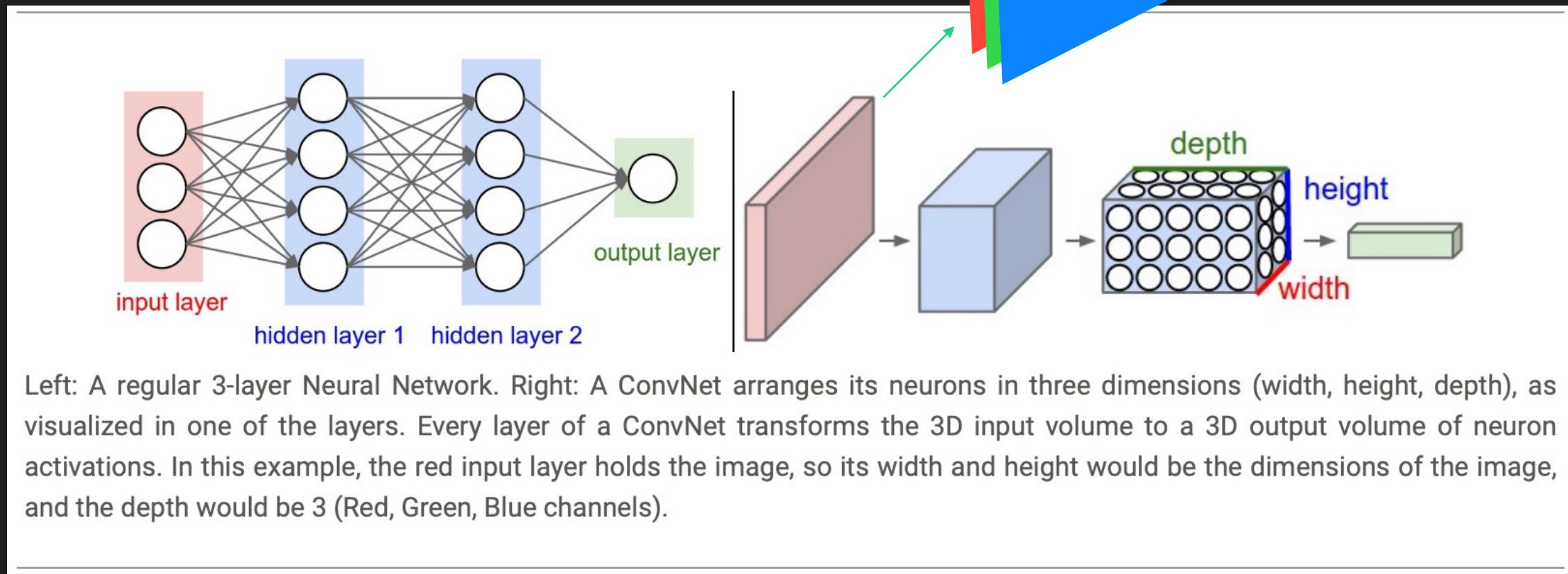
# CIFAR-10



- Canadian Institute For Advanced Research : ML, CV를 위해 일반적으로 사용되는 image의 모음
- 32x32x3. 10개의 class, 각 6000장.

출처 : <https://wikidocs.net/196917>

# Architecture Overview



Left: A regular 3-layer Neural Network. Right: A ConvNet arranges its neurons in three dimensions (width, height, depth), as visualized in one of the layers. Every layer of a ConvNet transforms the 3D input volume to a 3D output volume of neuron activations. In this example, the red input layer holds the image, so its width and height would be the dimensions of the image, and the depth would be 3 (Red, Green, Blue channels).

출처 : <https://cs231n.github.io/convolutional-networks/>

# Convolutional Layers

# Convolutional Layer

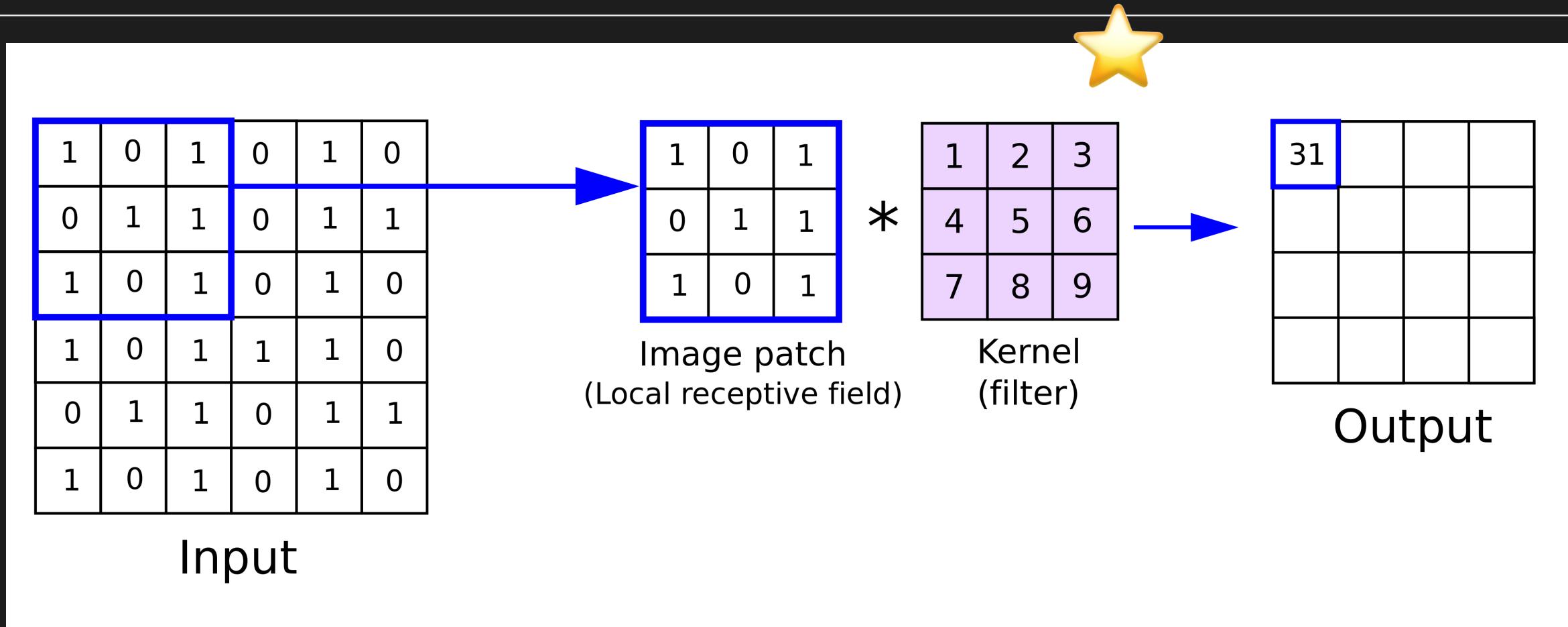
Input	Kernel	Output													
<table border="1" style="display: inline-table; vertical-align: middle;"><tr><td>0</td><td>1</td><td>2</td></tr><tr><td>3</td><td>4</td><td>5</td></tr><tr><td>6</td><td>7</td><td>8</td></tr></table>	0	1	2	3	4	5	6	7	8	$\ast$	<table border="1" style="display: inline-table; vertical-align: middle;"><tr><td>0</td><td>1</td></tr><tr><td>2</td><td>3</td></tr></table>	0	1	2	3
0	1	2													
3	4	5													
6	7	8													
0	1														
2	3														
	=	<table border="1" style="display: inline-table; vertical-align: middle;"><tr><td>19</td><td>25</td></tr><tr><td>37</td><td>43</td></tr></table>	19	25	37	43									
19	25														
37	43														

*Fig. 7.2.1 Two-dimensional cross-correlation operation. The shaded portions are the first output element as well as the input and kernel tensor elements used for the output computation:*

$$0 \times 0 + 1 \times 1 + 3 \times 2 + 4 \times 3 = 19.$$

출처 : [Dive Into Deep Learning](#)

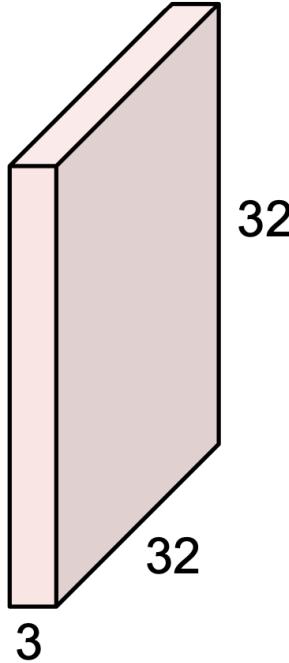
# Convolutional Layer



# Convolutional Layer

## Convolution Layer

32x32x3 image



Filters always extend the full depth of the input volume

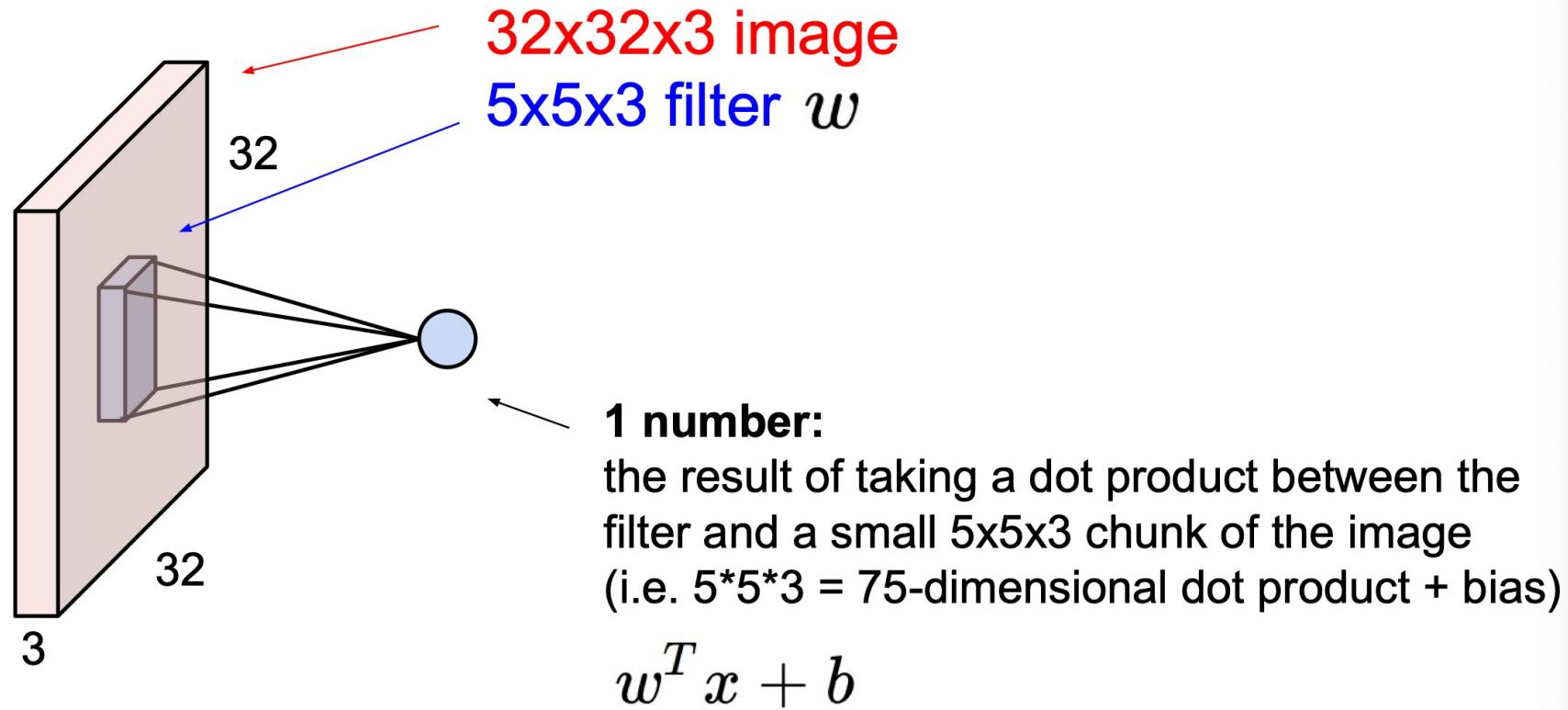
5x5x3 filter



**Convolve** the filter with the image  
i.e. “slide over the image spatially,  
computing dot products”

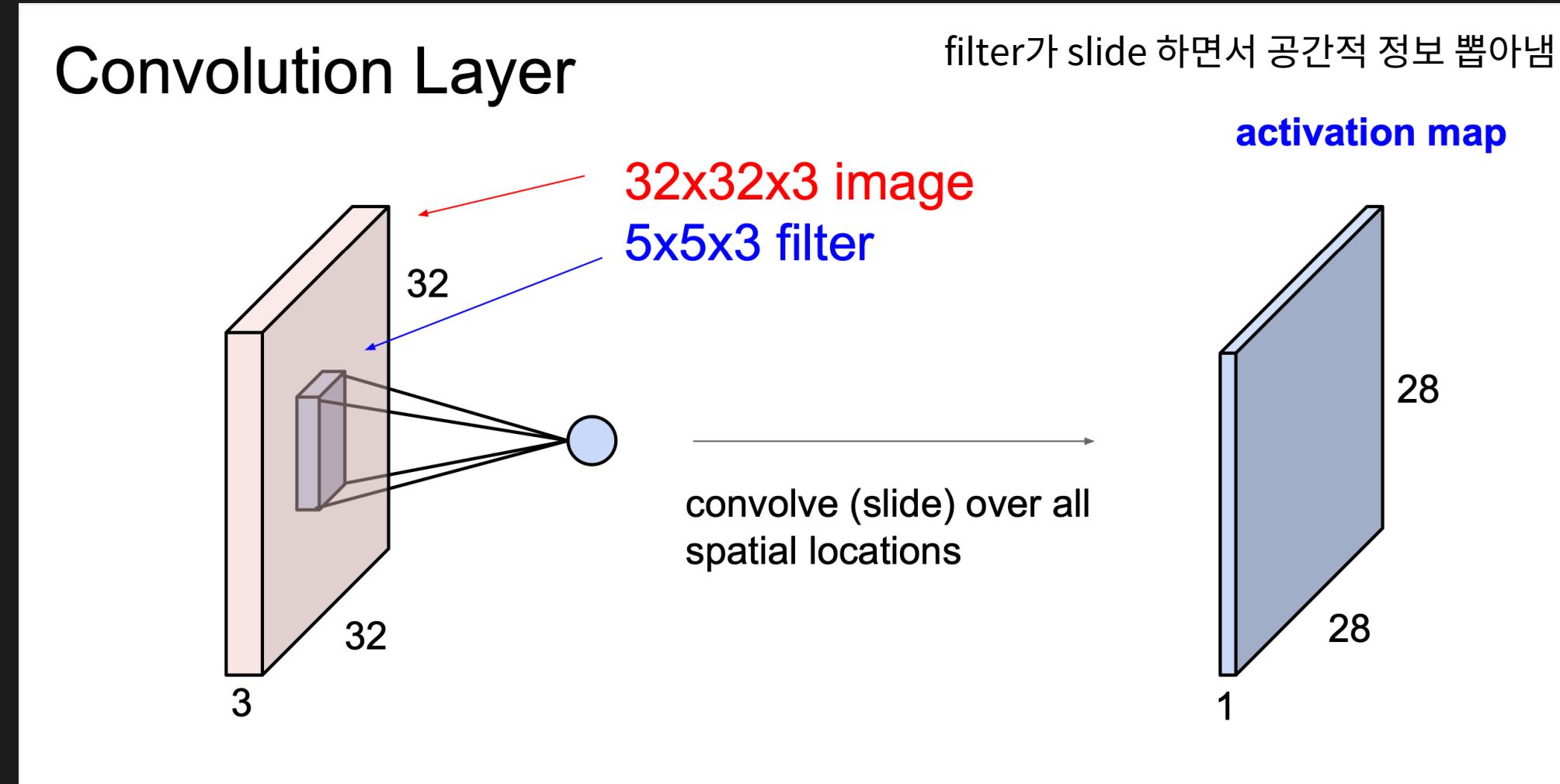
# Convolutional Layer

## Convolution Layer



출처 : CS231n (2017) Lecture 5

# Convolutional Layer

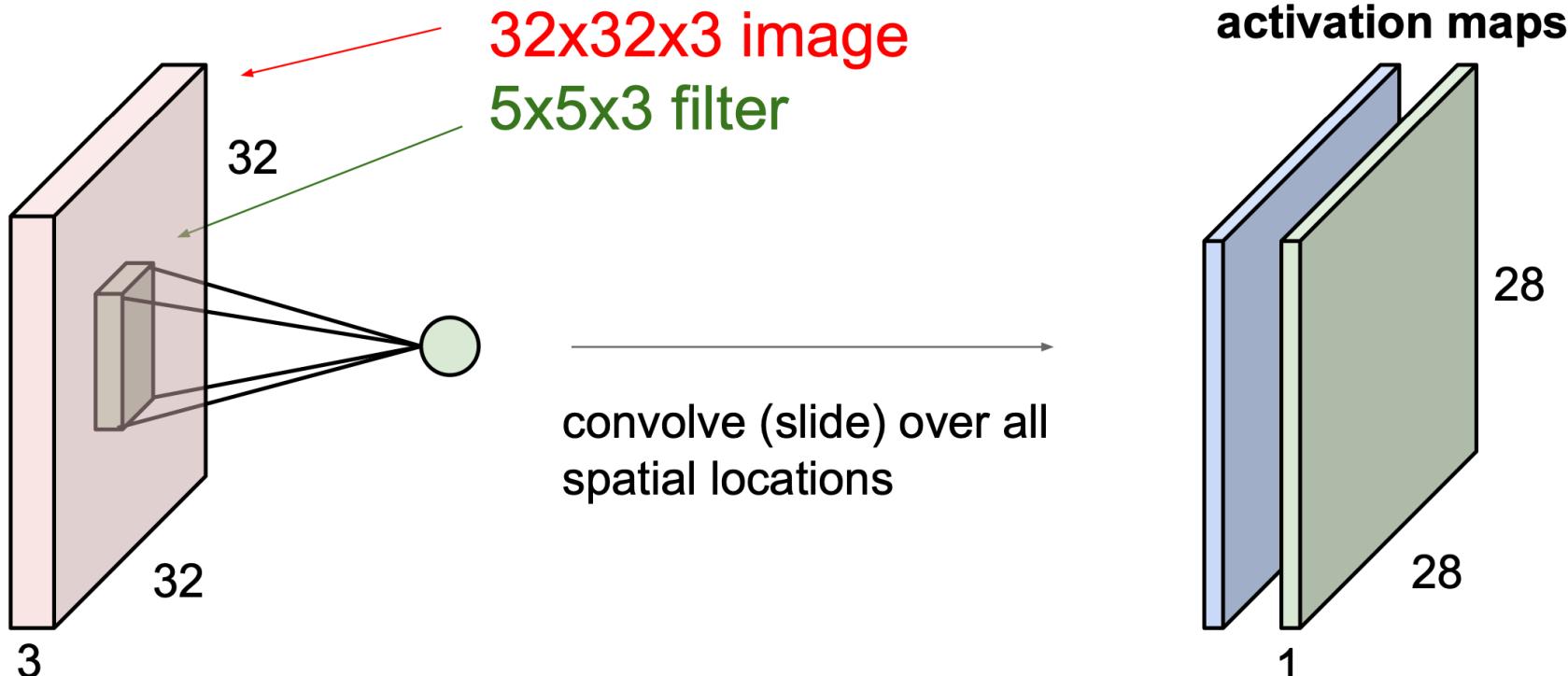


출처 : CS231n (2017) Lecture 5

# Convolutional Layer

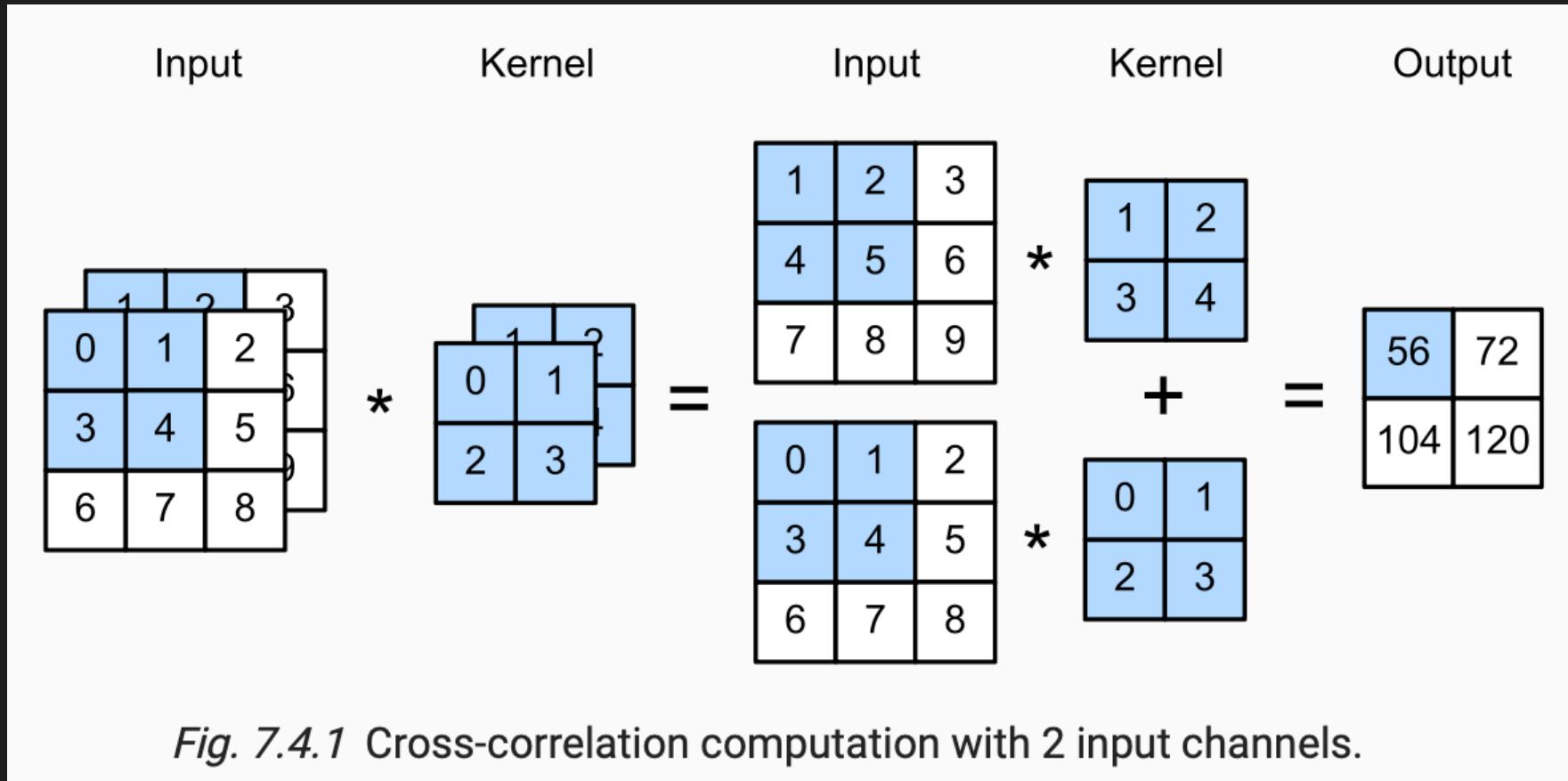
## Convolution Layer

consider a second, green filter



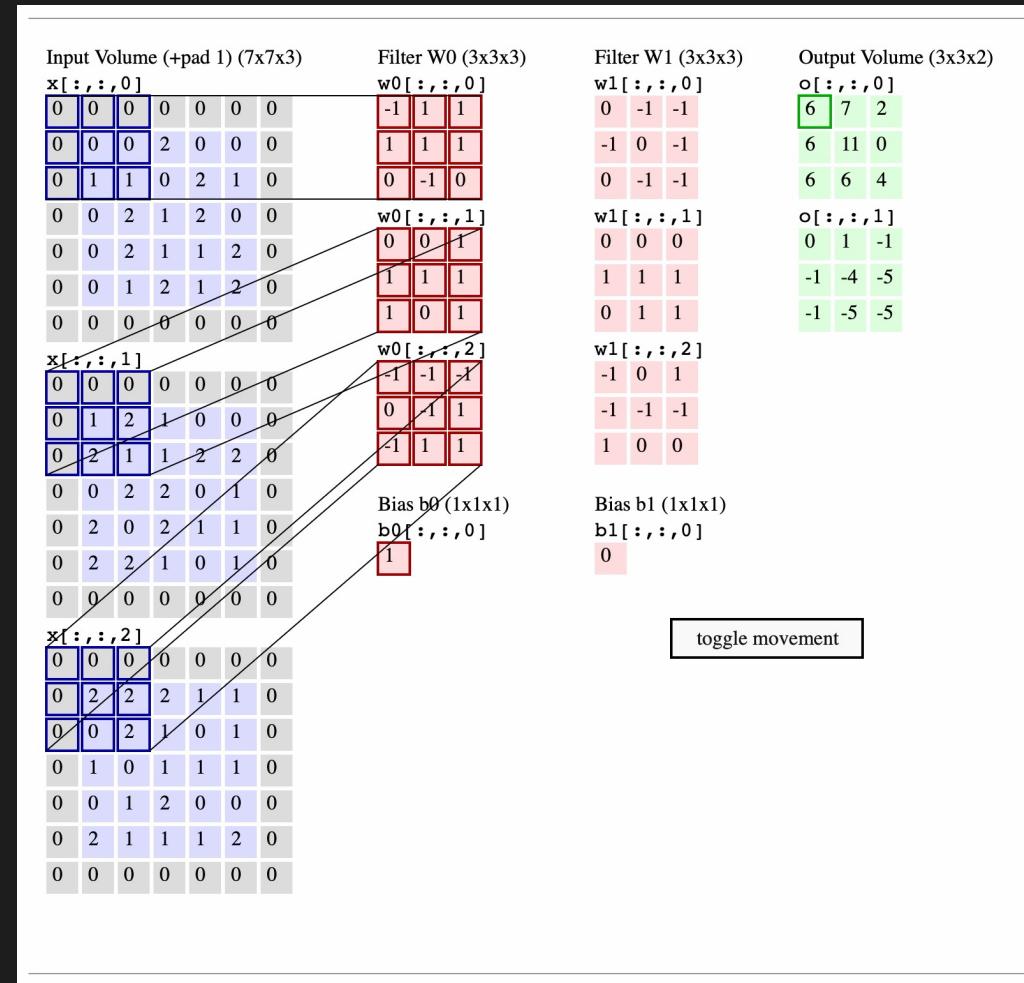
출처 : CS231n (2017) Lecture 5

# Convolutional Layer



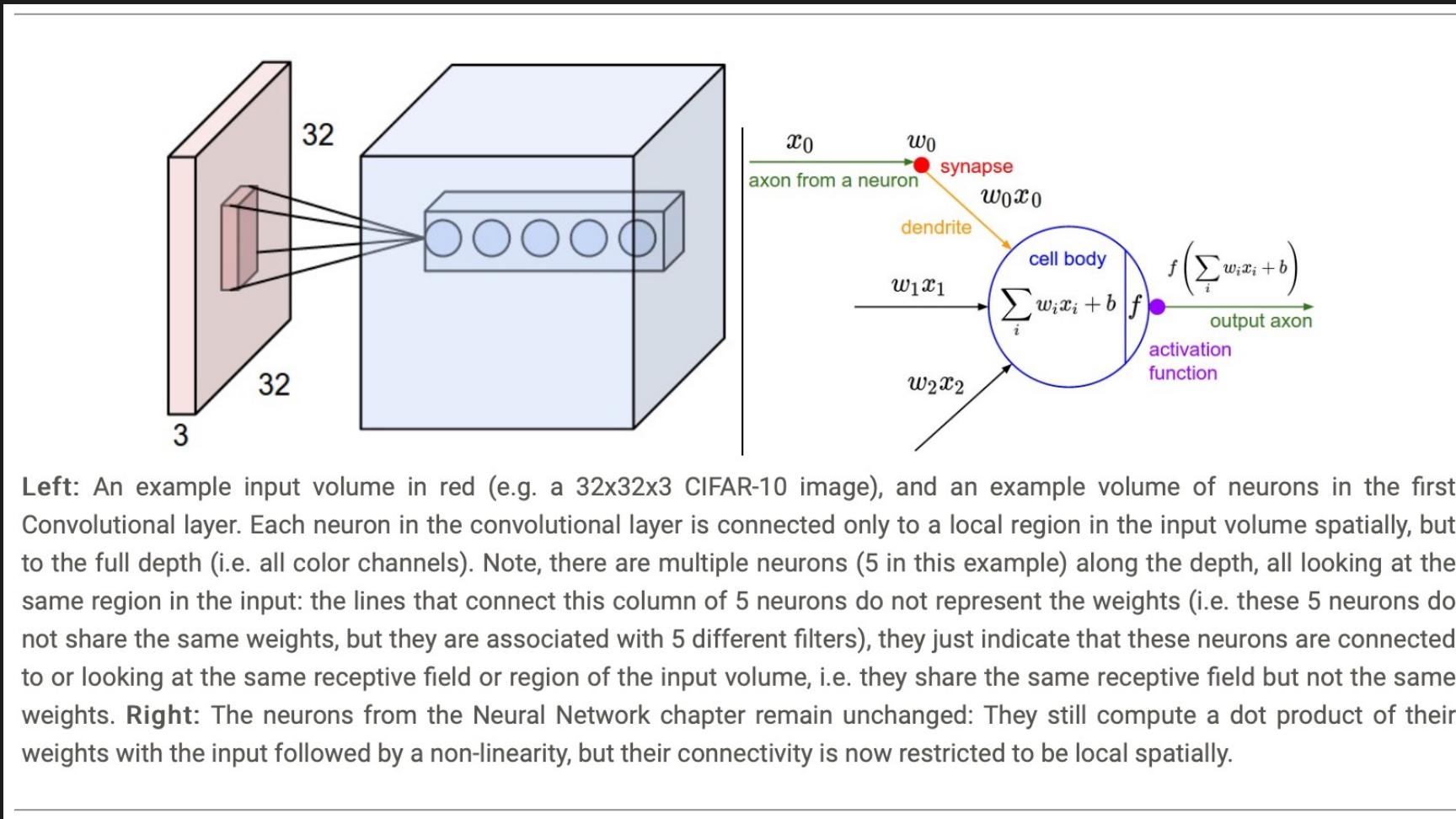
출처 : [Dive Into Deep Learning](#)

# Convolutional Layer



출처 : <https://cs231n.github.io/convolutional-networks/#conv>

# Convolutional Layer



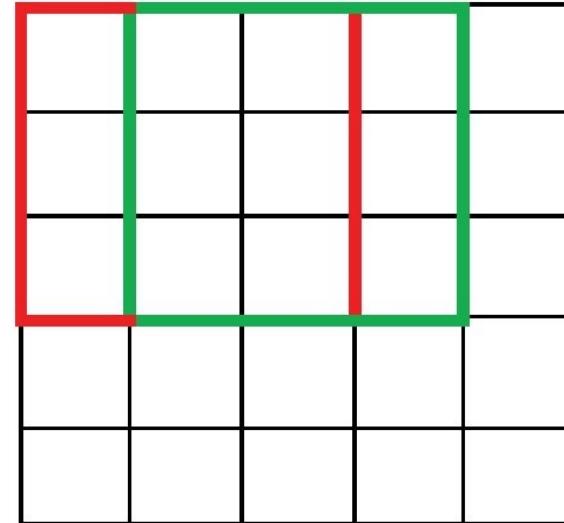
출처 : <https://cs231n.github.io/convolutional-networks/#conv>

# Convolutional Layer - Spatial arrangement

- Convolutional Layer는 input volume과 output volume에서 맞춰줘야 하는 값이 있다.
  - **depth of the output volume** : hyperparameter, 얼마나 많은 filter 쓰는지.
  - **stride** : filter를 얼마 만큼 slide 할 것인가? (1 또는 2 많이 씀)
  - **zero-padding** : 테두리에 0으로 채워 넣어 줌. output volume 조절에 유용.
- **input volume size (W), receptive field size (F), stride (S), zero padding (P)**
  - **Output volume** :  $O = \frac{W-F+2P}{S} + 1$
  - 직사각형이면?  $(OH, OW) = (\frac{W-FH+2P}{S} + 1, \frac{W-FW+2P}{S} + 1)$

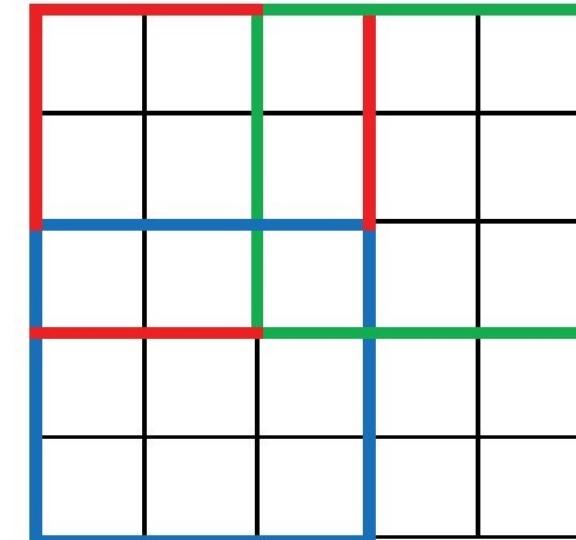
# Convolutional Layer - stride

Convolution  
with Stride=1



Output

Convolution  
with Stride=2



Output

출처 : <https://www.analyticsvidhya.com/blog/2022/03/basics-of-cnn-in-deep-learning/>

# Convolutional Layer - stride

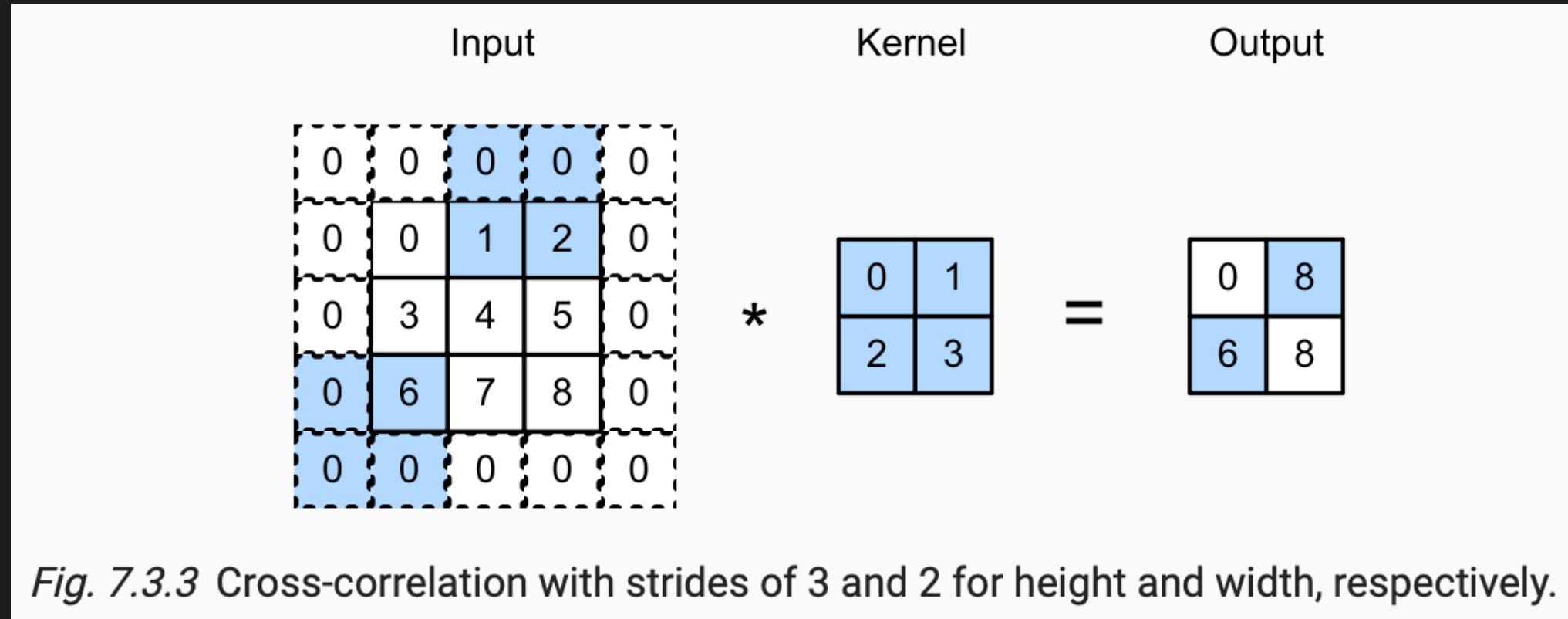


Fig. 7.3.3 Cross-correlation with strides of 3 and 2 for height and width, respectively.

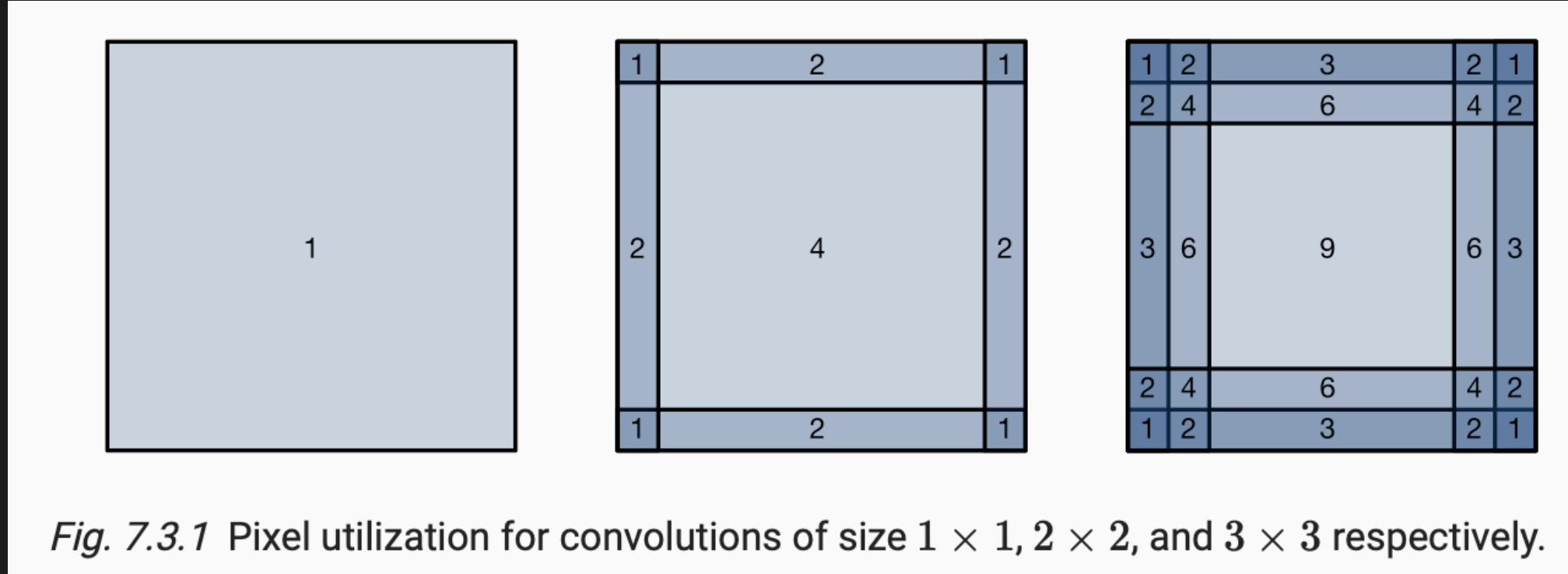
출처 : [Dive Into Deep Learning](#)

# Convolutional Layer - zero padding

0	0	0	0	0	0
0	35	19	25	6	0
0	13	22	16	53	0
0	4	3	7	10	0
0	9	8	1	3	0
0	0	0	0	0	0

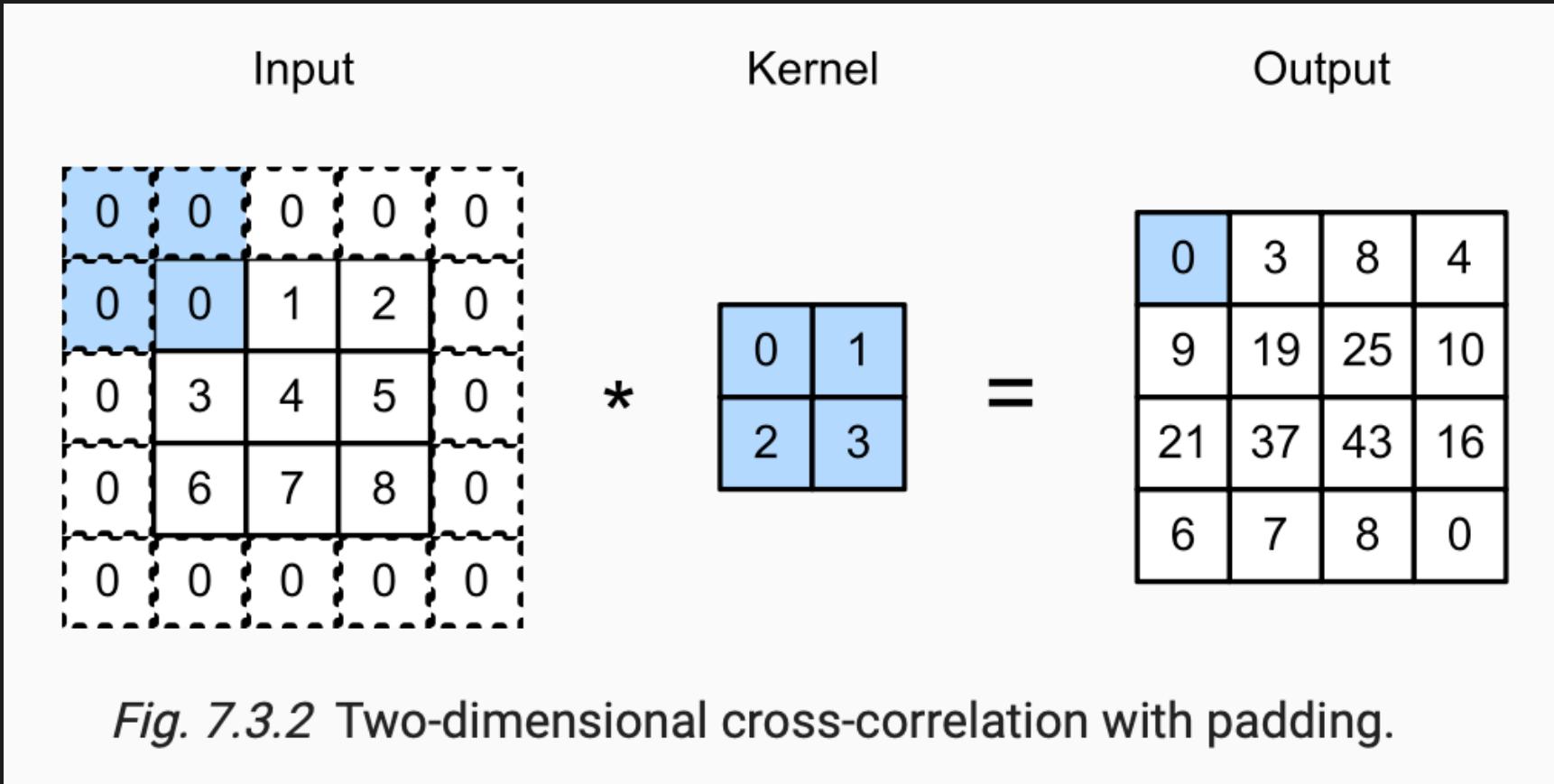
출처 : <https://blog.xrds.acm.org/2016/06/convolutional-neural-networks-cnns-illustrated-explanation/>

# Convolutional Layer - zero padding



출처 : [Dive Into Deep Learning](#)

# Convolutional Layer - zero padding



출처 : [Dive Into Deep Learning](#)

# Convolutional Layer - Spatial arrangement

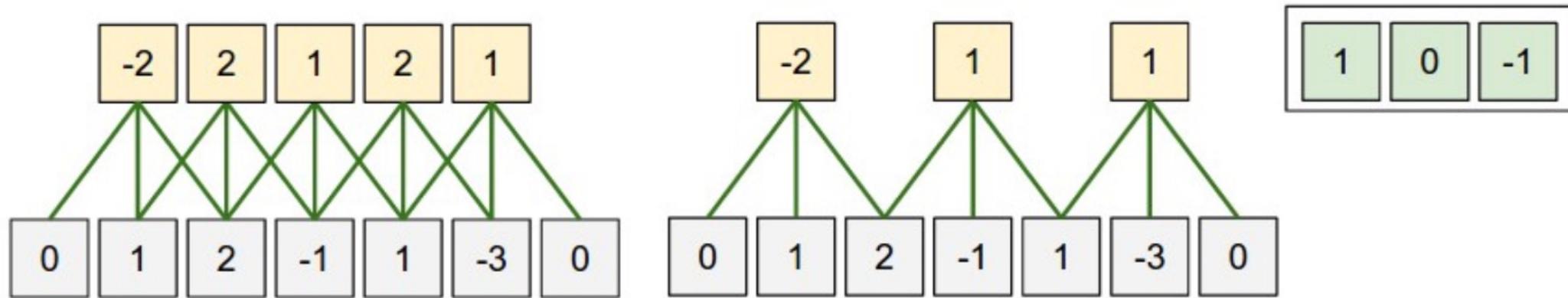


Illustration of spatial arrangement. In this example there is only one spatial dimension (x-axis), one neuron with a receptive field size of  $F = 3$ , the input size is  $W = 5$ , and there is zero padding of  $P = 1$ . **Left:** The neuron strided across the input in stride of  $S = 1$ , giving output of size  $(5 - 3 + 2)/1+1 = 5$ . **Right:** The neuron uses stride of  $S = 2$ , giving output of size  $(5 - 3 + 2)/2+1 = 3$ . Notice that stride  $S = 3$  could not be used since it wouldn't fit neatly across the volume. In terms of the equation, this can be determined since  $(5 - 3 + 2) = 4$  is not divisible by 3.

The neuron weights are in this example  $[1,0,-1]$  (shown on very right), and its bias is zero. These weights are shared across all yellow neurons (see parameter sharing below).

출처 : <https://cs231n.github.io/convolutional-networks/#conv>

# Convolutional Layer - Let's Check!

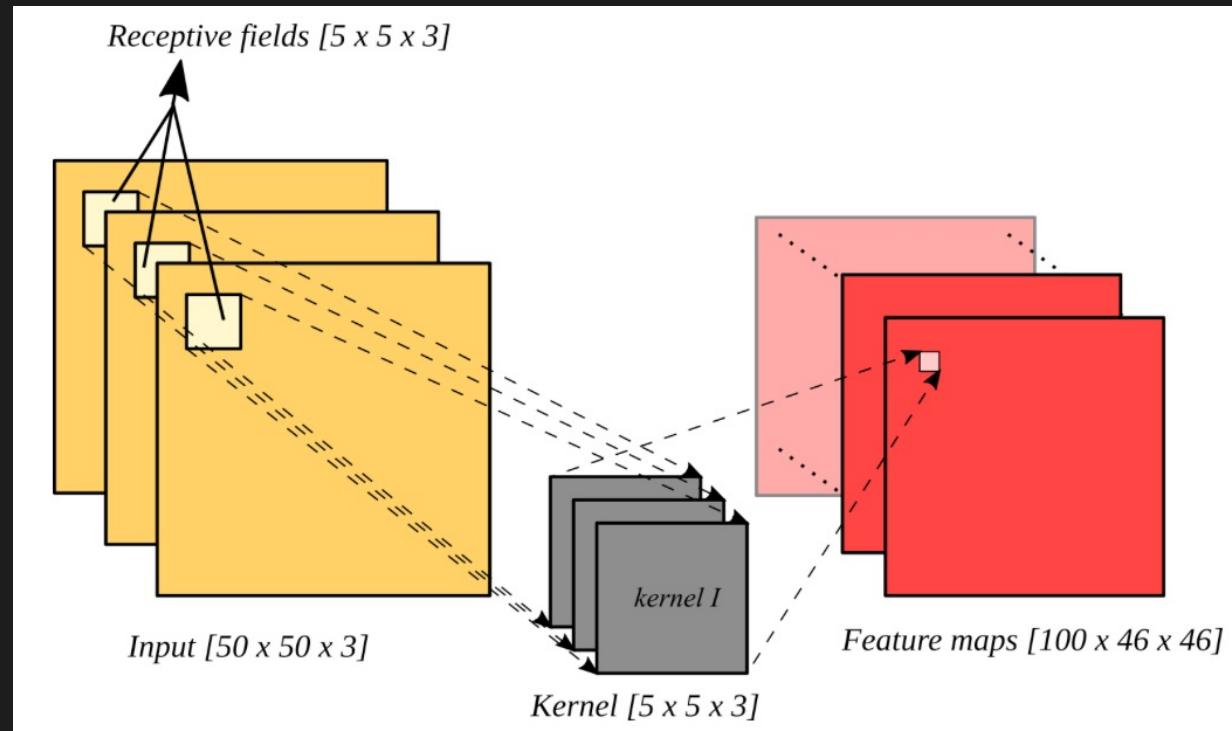
---

- Input volume :  $32 \times 32 \times 3$
- 10개의  $5 \times 5$  filters
- Stride 1
- zero padding 2
- 이 때, Output Volume size는?

# Convolutional Layer - Let's Check!

- Input volume :  $32 (W) \times 32 \times 3$
- 10개의  $5 (F)$  x 5 filters
- Stride 1 (S)
- zero padding 2 (P)
- 이 때, Output Volume size는?
  - Output volume =  $(W - F + 2P)/S + 1 = (32 - 5 + 2 * 2) / 1 + 1 = 32$
  - $32 \times 32 \times 10$

# Convolutional Layer - Feature Map and Receptive Field



출처 : <https://ai.stackexchange.com/questions/8701/what-is-the-difference-between-a-receptive-field-and-a-feature-map>

- feature map : convolution operation의 output
- receptive field : neuron의 activation 계산에 사용되는 부분

# Convolutional Layer - Feature Map and Receptive Field

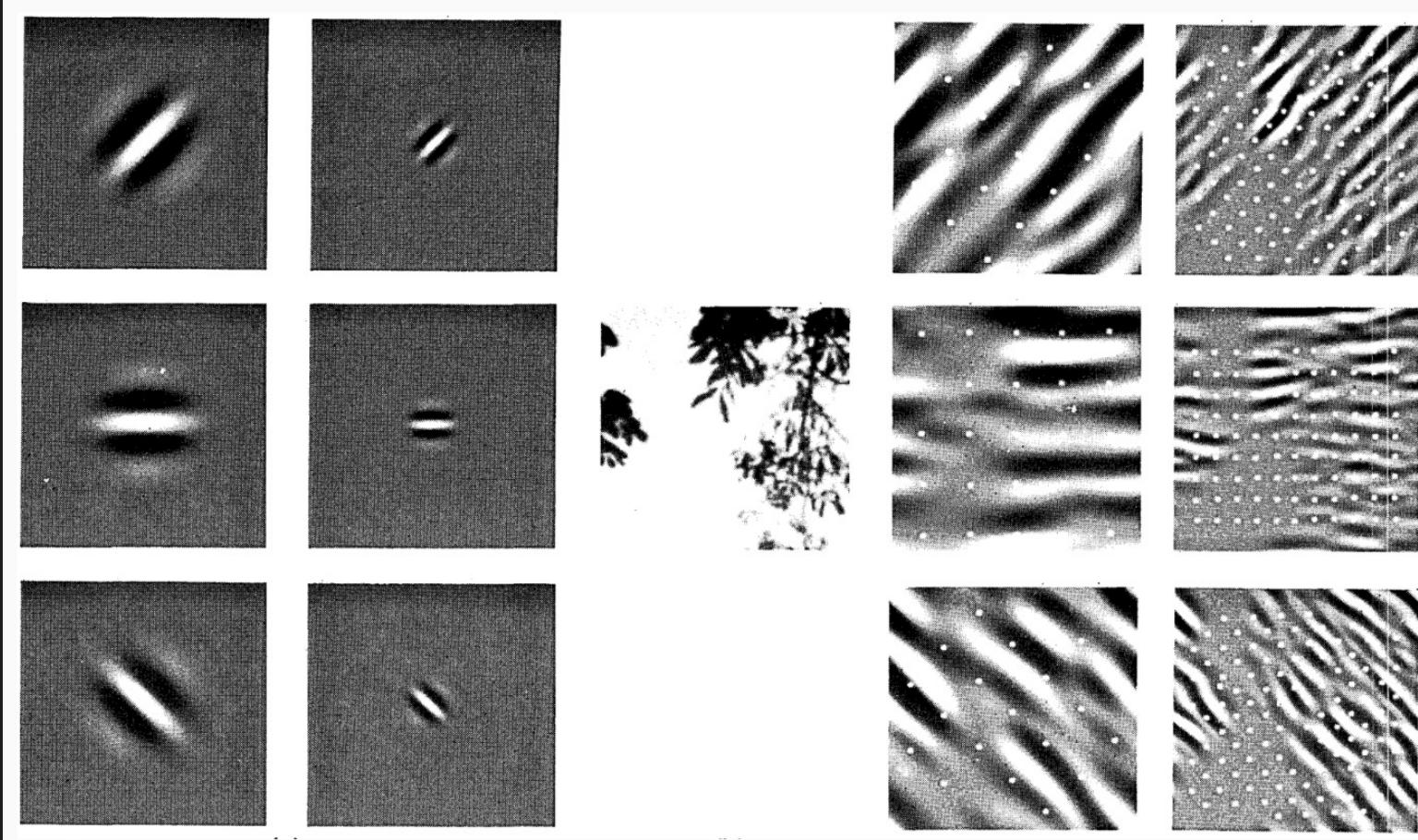
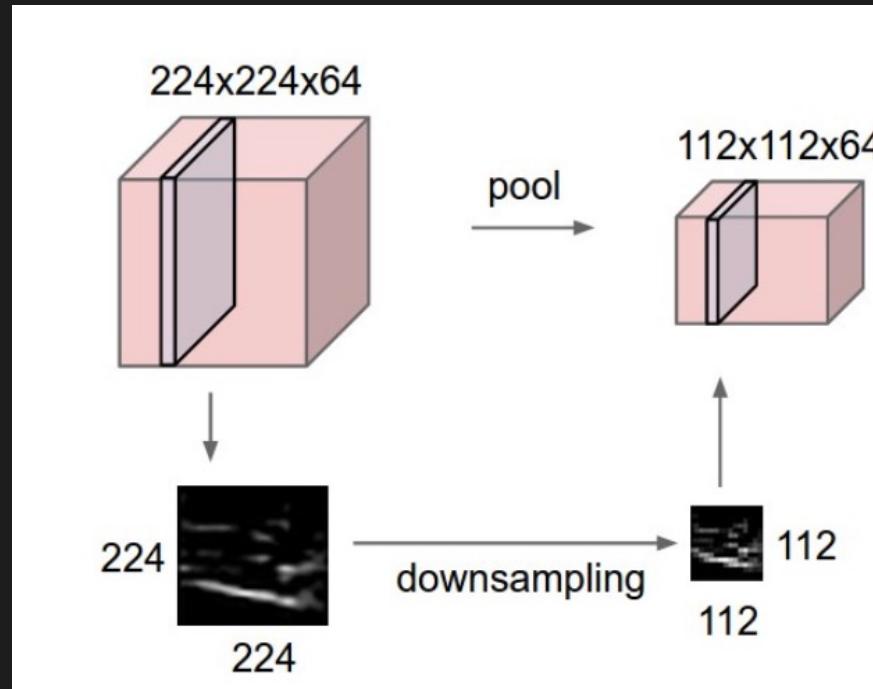


Fig. 7.2.2 Figure and caption taken from Field (1987): An example of coding with six different channels. (Left) Examples of the six types of sensor associated with each channel. (Right) Convolution of the image in (Middle) with the six sensors shown in (Left). The response of the individual sensors is determined by sampling these filtered images at a distance proportional to the size of the sensor (shown with dots). This diagram shows the response of only the even symmetric sensors.

출처 : [Dive Into Deep Learning](#)

# Pooling Layer

# Pooling?



출처 : <https://cs231n.github.io/convolutional-networks/#conv>

- 계속해서 Convolutional Layer의 정보를 더해 간다.
  - Layer가 깊어질 수록, Receptive field가 넓어진다.
- Convolutional Layer의 민감성을 완화하고, 공간적으로 downsampling 하기 위해!

# Max-pooling

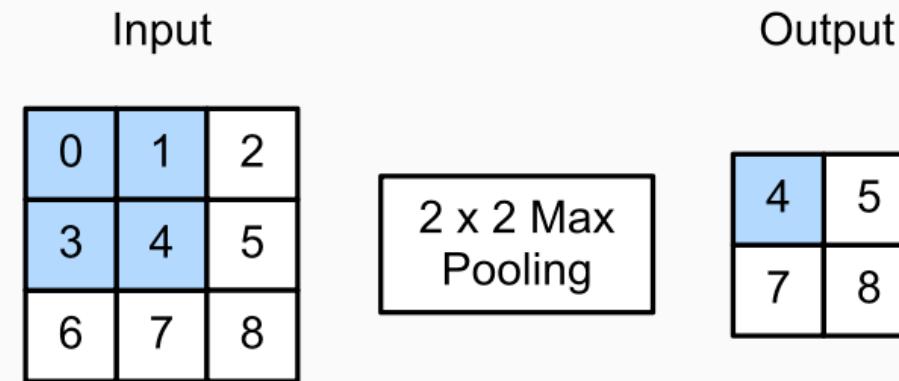
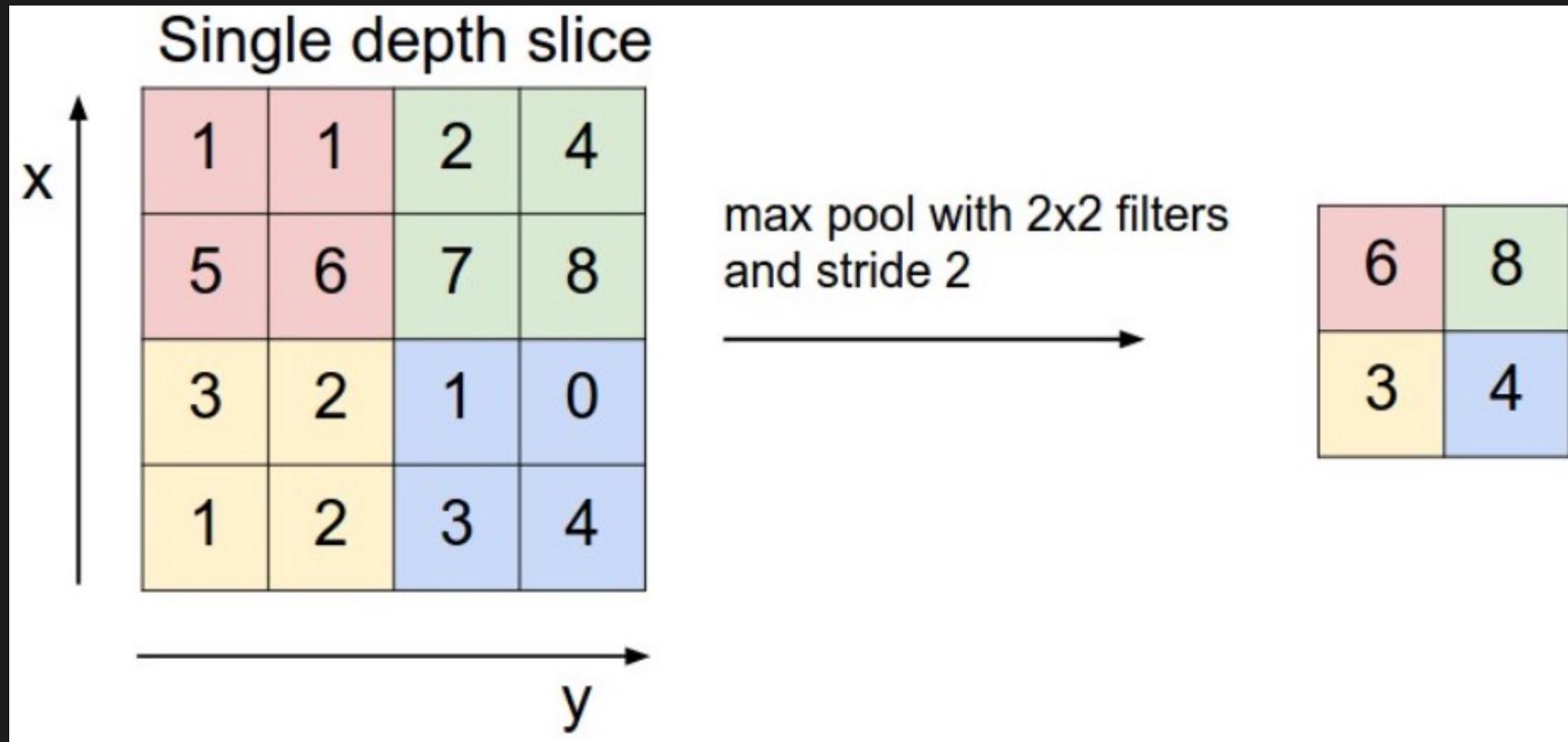


Fig. 7.5.1 Max-pooling with a pooling window shape of  $2 \times 2$ . The shaded portions are the first output element as well as the input tensor elements used for the output computation:

$$\max(0, 1, 3, 4) = 4.$$

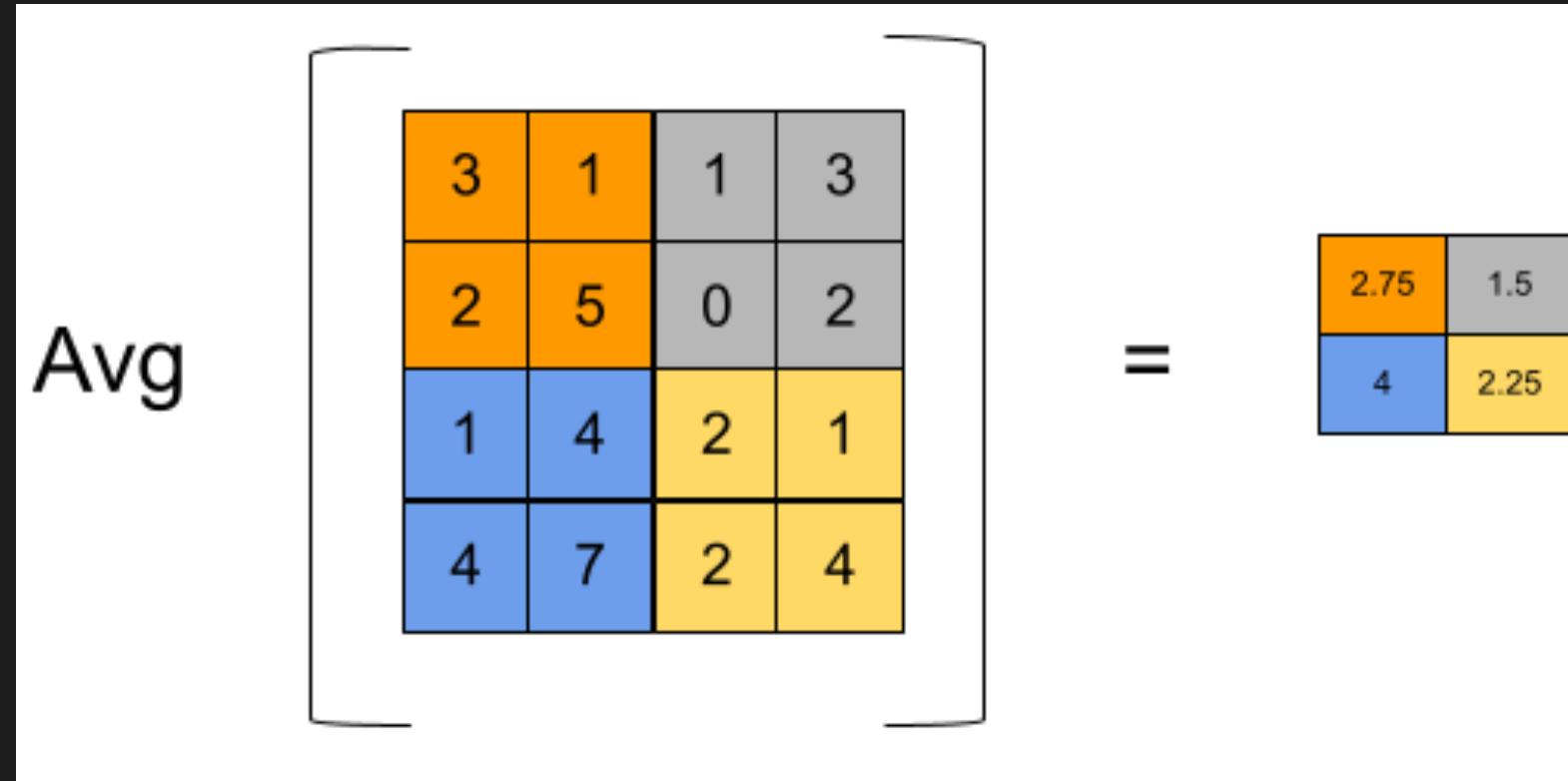
출처 : [Dive Into Deep Learning](#)

# Max-pooling



출처 : <https://cs231n.github.io/convolutional-networks/#conv>

# Average Pooling



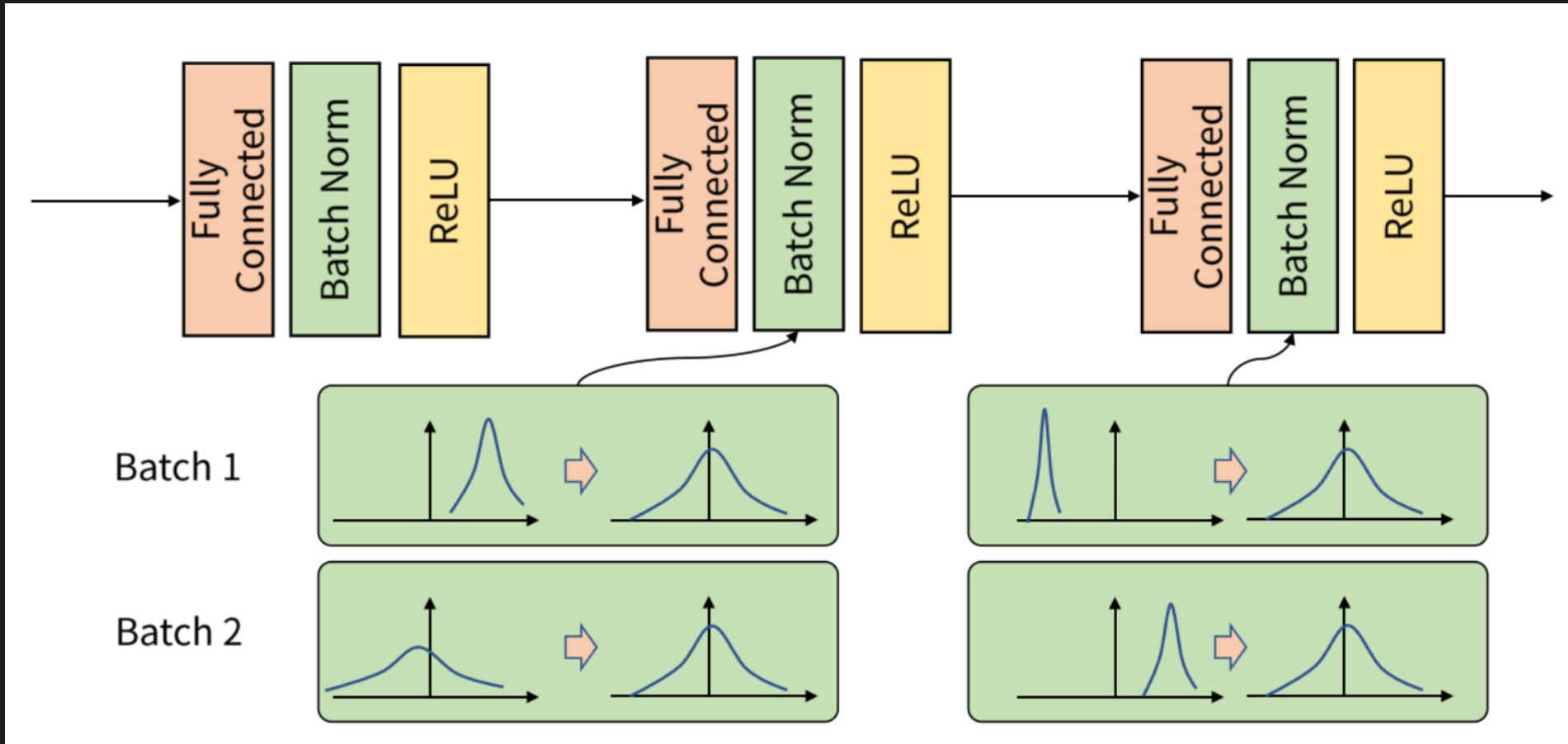
- Image 분야에선 Max Pooling이 일반적

# Other Techniques

---

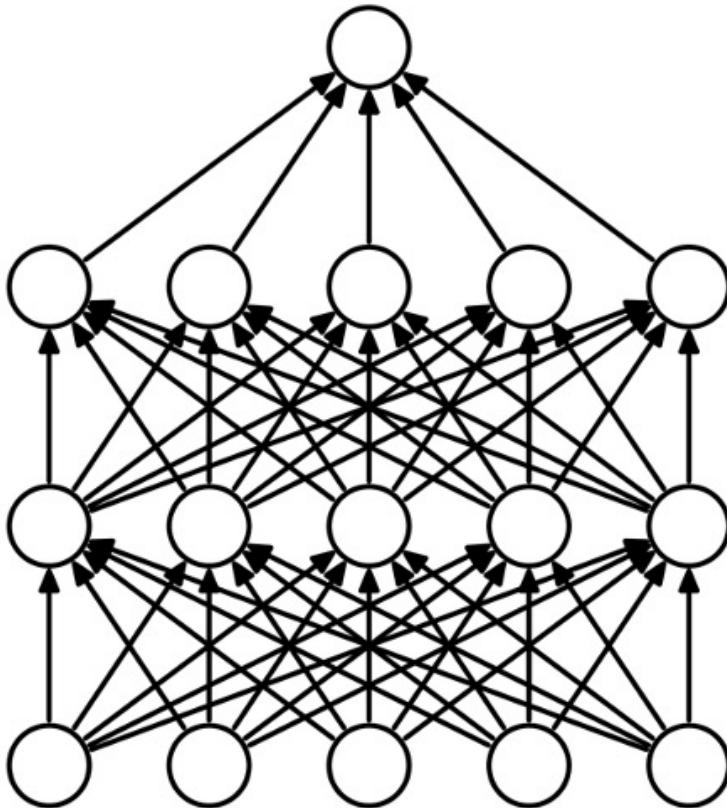
- Batch Normalization
  - Stochastic Gradient Descent 하는 단위 : Batch
  - Batch data들의 distribution이 다르다. 이를 맞춰주자!
- Dropout
  - Parameter를 Random하게 학습 시키지 않아 Overfitting 방지
  - Ensemble과 같은 효과

# Batch Normalization

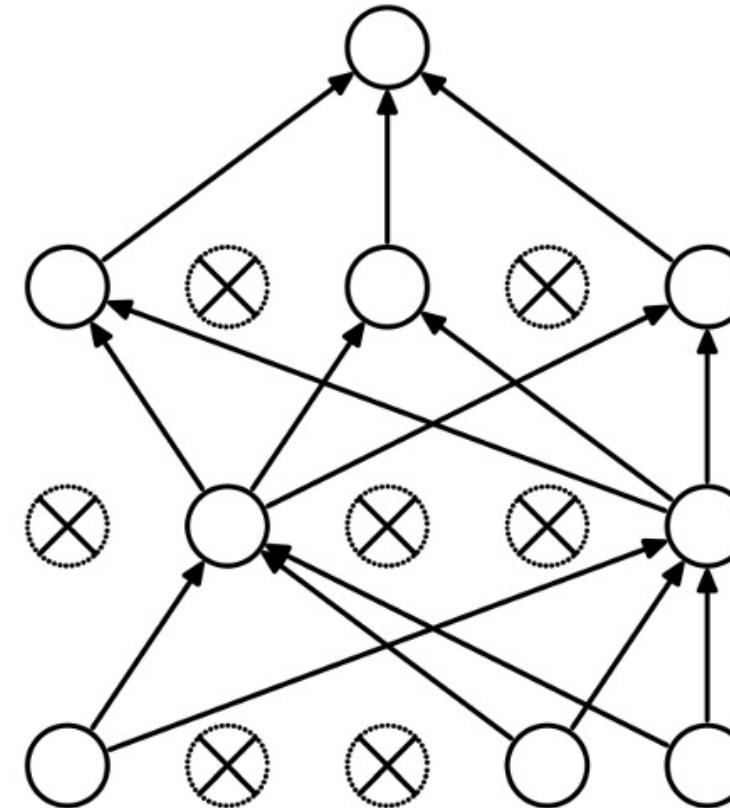


출처 : <https://gaussian37.github.io/dl-concept-batchnorm/>

# Dropout



(a) Standard Neural Net



(b) After applying dropout.

출처 : Dropout: a simple way to prevent neural networks from overfitting, Srivastava et al., Journal of Machine Learning Research, 2014

# Case Study : Modern CNNs

# LeNet Overview

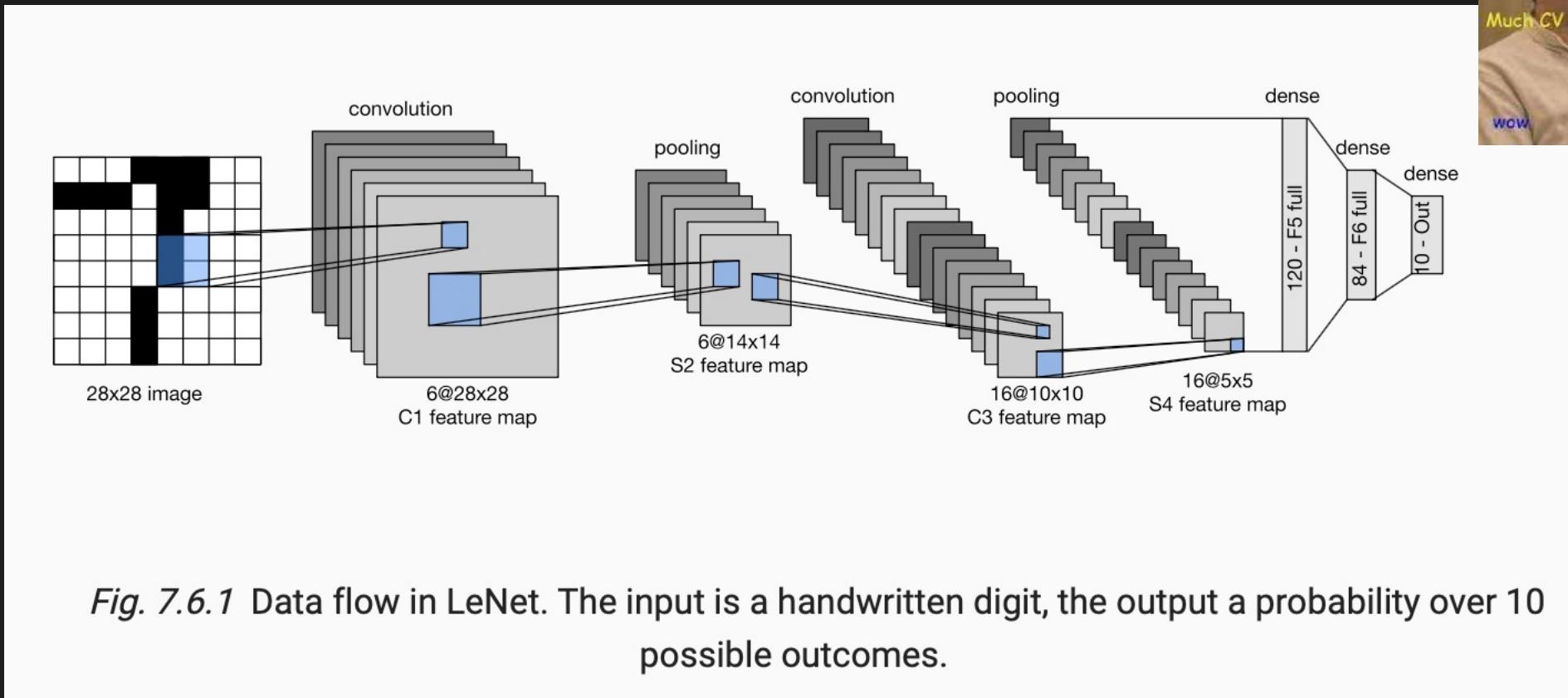
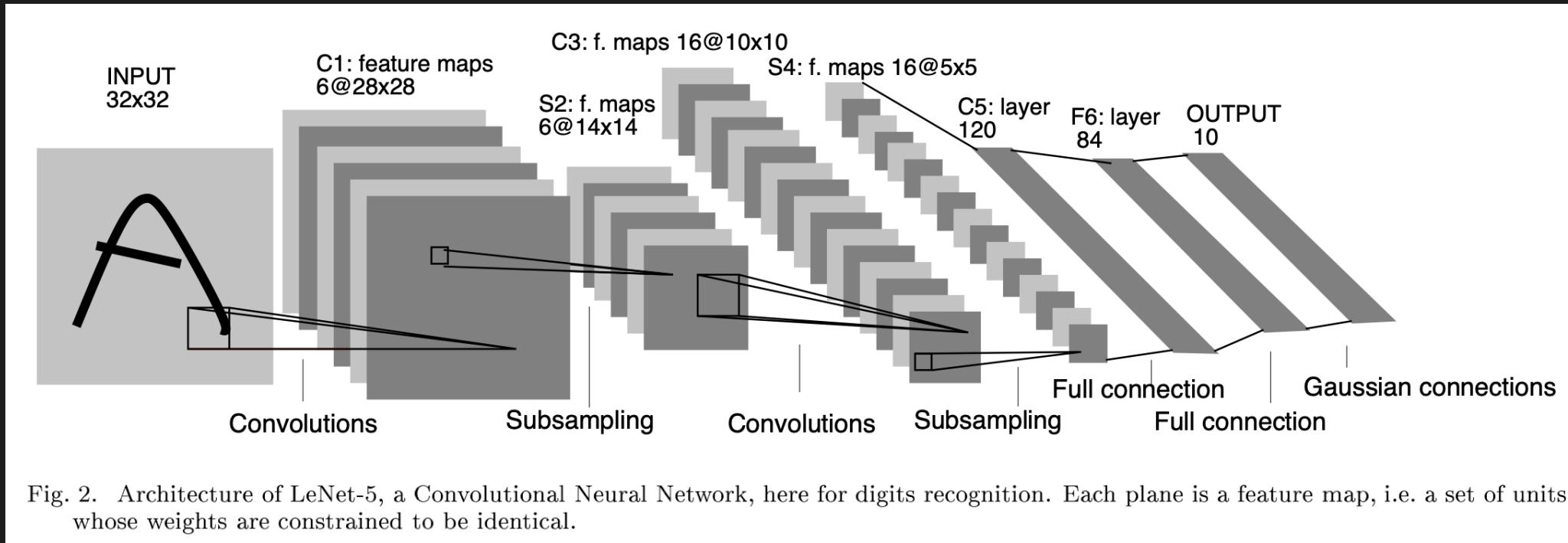


Fig. 7.6.1 Data flow in LeNet. The input is a handwritten digit, the output a probability over 10 possible outcomes.

출처 : [Dive Into Deep Learning](#)

- Yann LeCun!

# LeNet-5



출처 : [Gradient Based Learning Applied to Document Recognition \(LeCun et al., 1995\)](#)

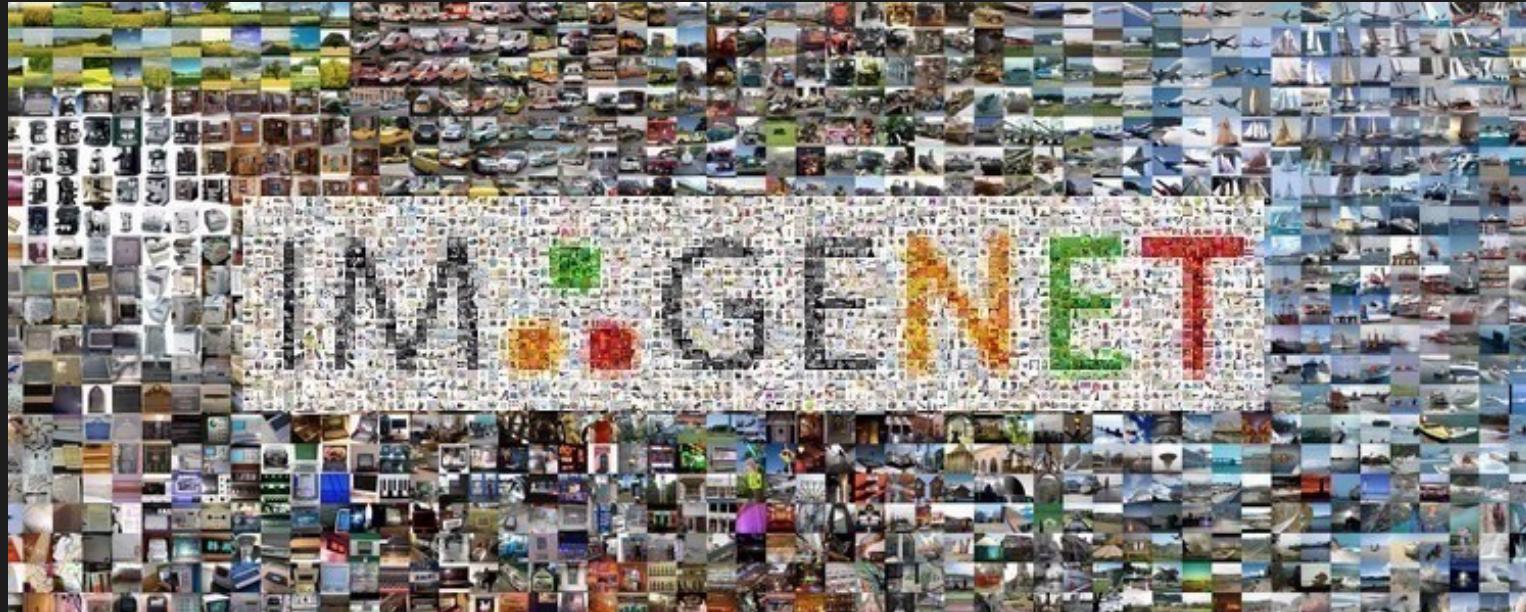
- Breakthrough for Neural Networks

# Limitation of LeNet

---

- 더 크고, realistic dataset에 적용 시키기 어려웠다.
  - 현실 세계의 image는 더 고차원
  - representation learning
- 왜 LeNet(1995년)에는 불가능했지만, AlexNet(2012년)에는 가능했는가?
  - Missing Ingredient
    - Data
    - GPU

# Missing Ingredient: Data



출처 : <https://www.kaggle.com/datasets/sautkin/imagenet1kvalid>

- ImageNet dataset (2009년 released)
  - 100만 장 (CIFAR-100은 6만장)
  - 224 x 224 pixels

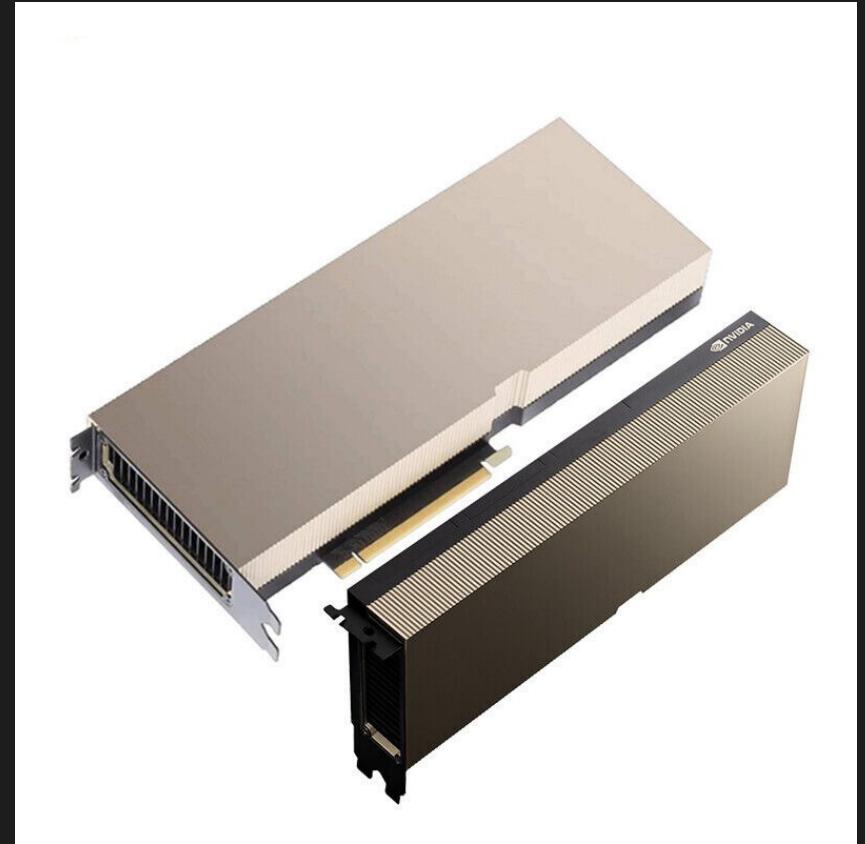
# Missing Ingredient: Hardware

NVIDIA GeForce GTX 580



출처 : <https://www.techpowerup.com/gpu-specs/geforce-gtx-580.c270>

NVIDIA Tesla A100 Ampere



출처 : <https://www.ebay.com/p/26042475154>

# Deep Convolutional Neural Networks (AlexNet)

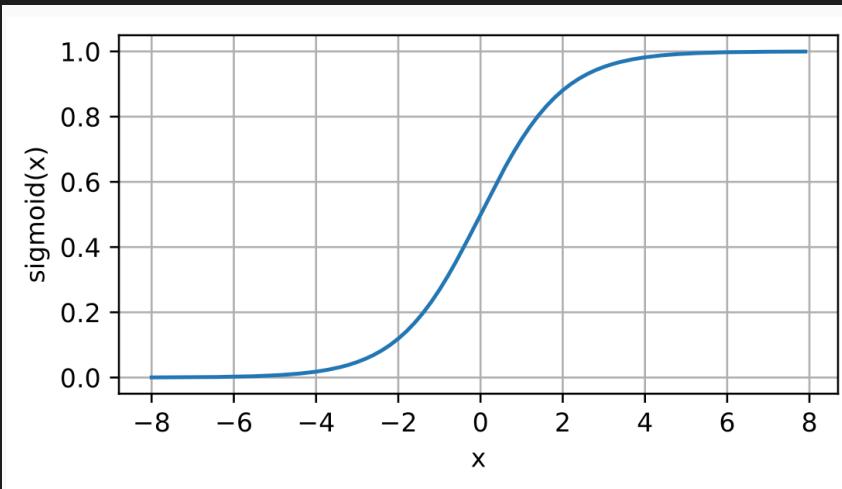


- ImageNet Challenge 2012년에 혜성 같이 등장한 ‘AlexNet’
- Geoffrey Hinton
- 두 개의 GPU에서 학습 되도록 설계
- LeNet과의 차이점
  - Deeper!
  - No sigmoid, Yes ReLU!

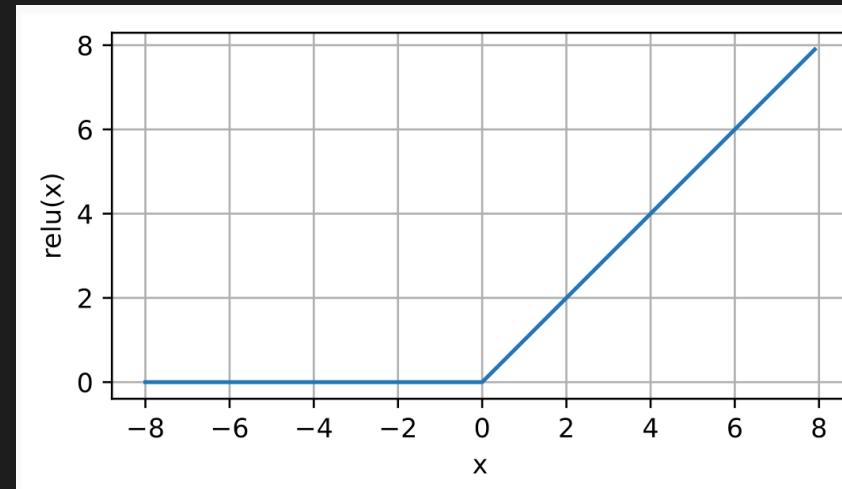
출처 : [Dive Into Deep Learning](#)

# Recap: Sigmoid vs ReLU

Sigmoid



ReLU

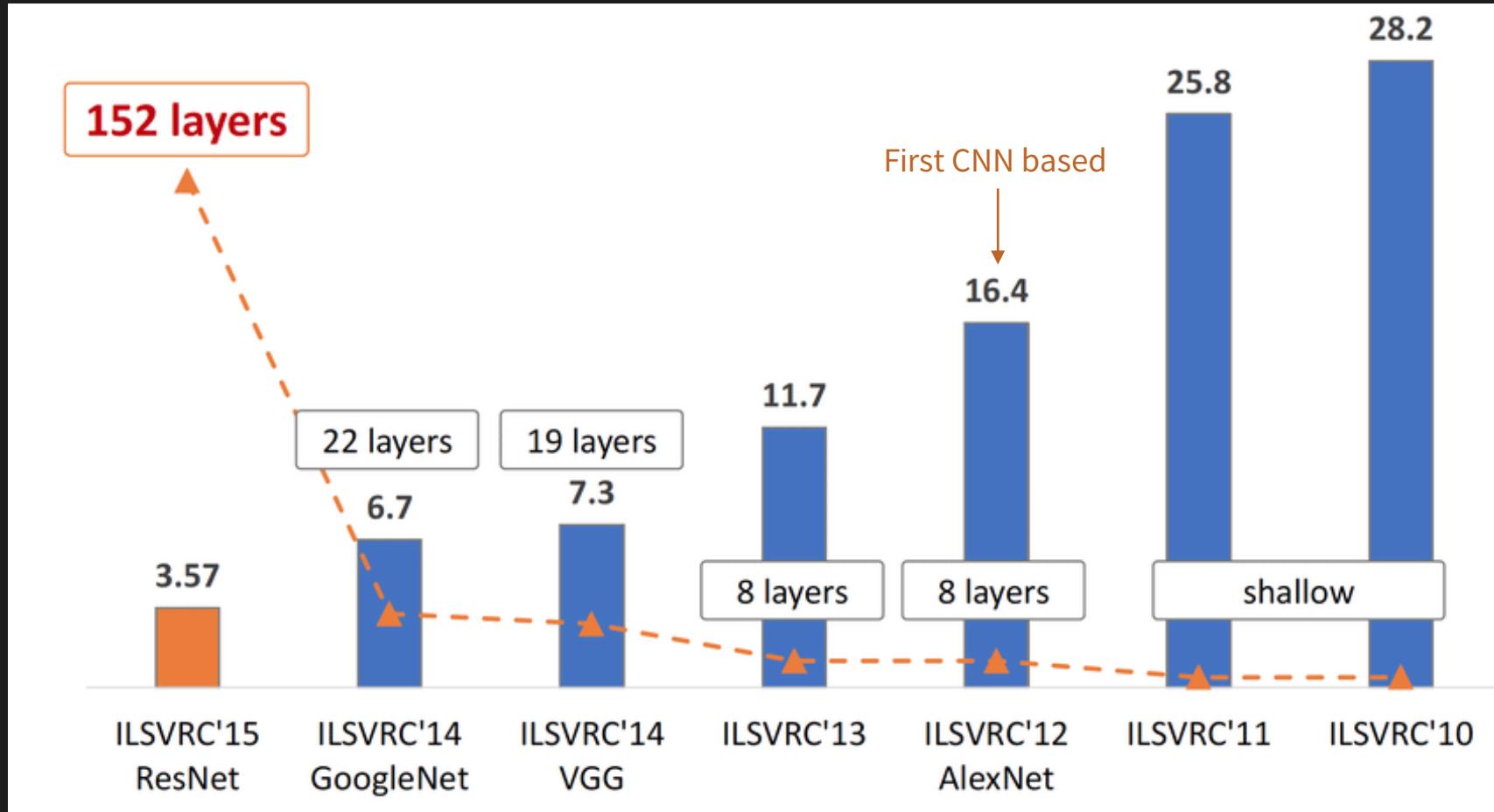


출처 : [Dive Into Deep Learning](#)

출처 : [Dive Into Deep Learning](#)

- 더 간단한 ReLU (exponentiation  $\times$ ) → 미분이 쉽다!
- Sigmoid( $x$ )가 0 또는 1에 가까우면 기울기가 0이 된다. initialized가 잘 되어야 함.
  - ReLU는 기울기 전달 가능. 학습이 쉽다.

# ImageNet Challenge - Much Deeper!



출처

출처 : [https://www.researchgate.net/figure/The-evolution-of-the-winning-entries-on-the-ImageNet-Large-Scale-Visual-Recognition\\_fig1\\_321896881](https://www.researchgate.net/figure/The-evolution-of-the-winning-entries-on-the-ImageNet-Large-Scale-Visual-Recognition_fig1_321896881)

# VGG Network

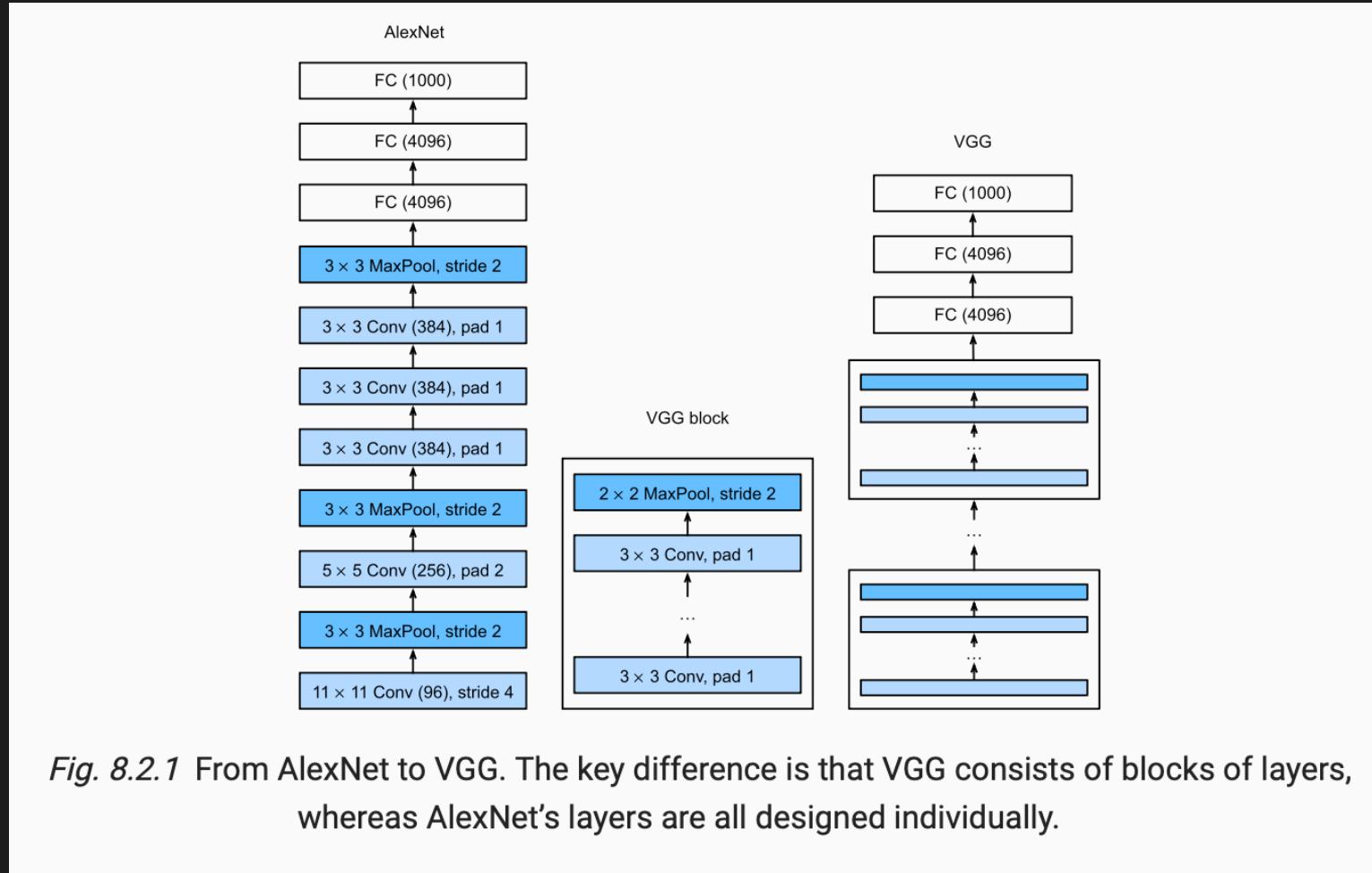


Fig. 8.2.1 From AlexNet to VGG. The key difference is that VGG consists of blocks of layers, whereas AlexNet's layers are all designed individually.

- 더 깊어진 Network
- Heuristic concepts

출처 : [Dive Into Deep Learning](#)

# VGG Network

## VGG: Deeper, Regular design

VGG design rules:

All conv are 3x3 stride 1

All max pool 2x2 stride 2

after pool, double # channels

- two 3x3 vs one 5x5
- covers same local area (called receptive field)
- lower resource
- more nonlinearity

Option 1:  
Conv (5x5, C->C)

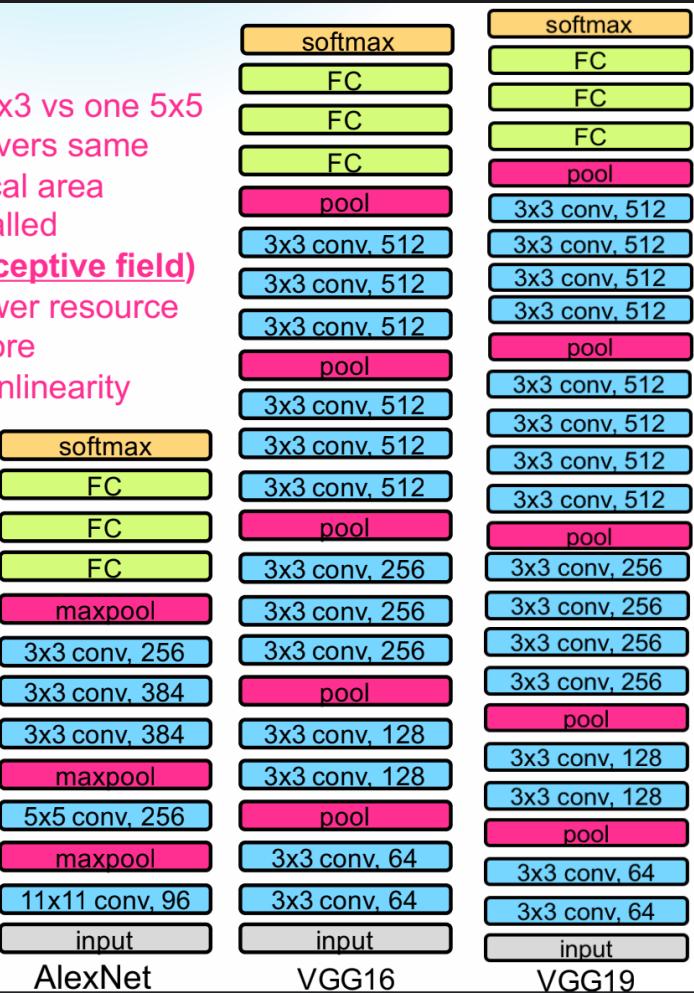
Params:  $25C^2$

FLOPs:  $25C^2HW$

Option 2:  
Conv (3x3, C->C)  
Conv (3x3, C->C)

Params:  $18C^2$

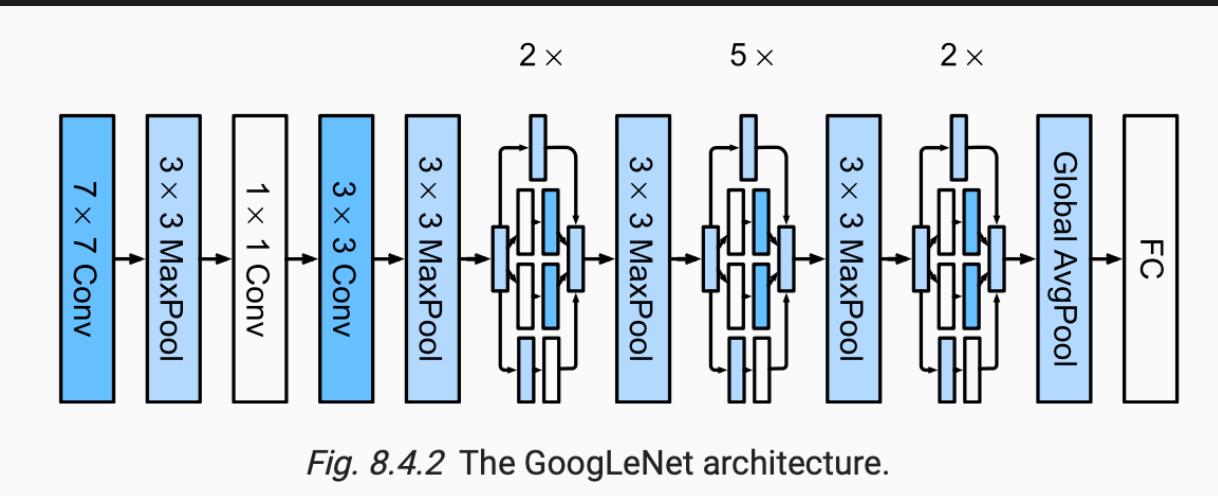
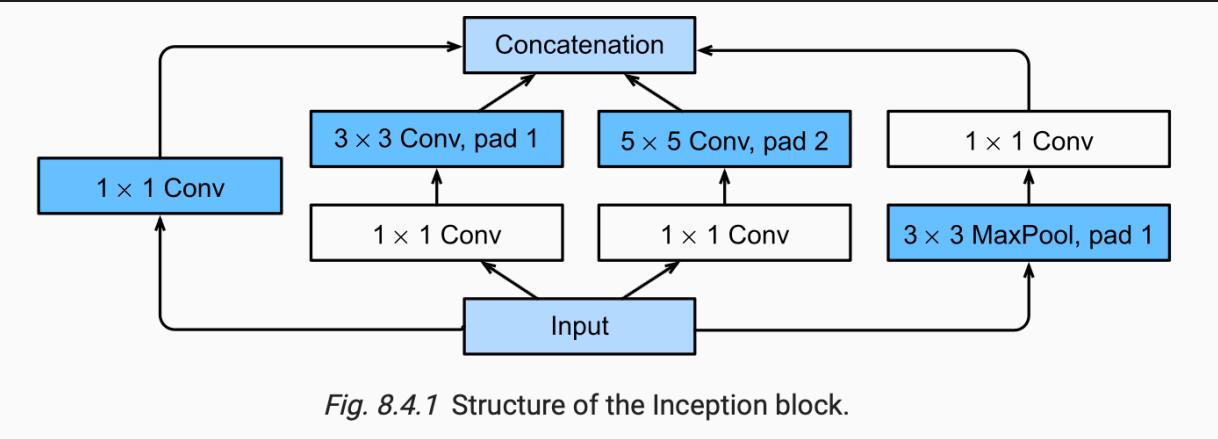
FLOPs:  $18C^2HW$



출처 미상

- $3 \times 3$  두 개 ==  $5 \times 5$  한 개
  - 같은 receptive field
- 하지만  $3 \times 3$ 이 더 적은 resource로 더 많은 유동성.
- $3 \times 3$  Conv로 통일!

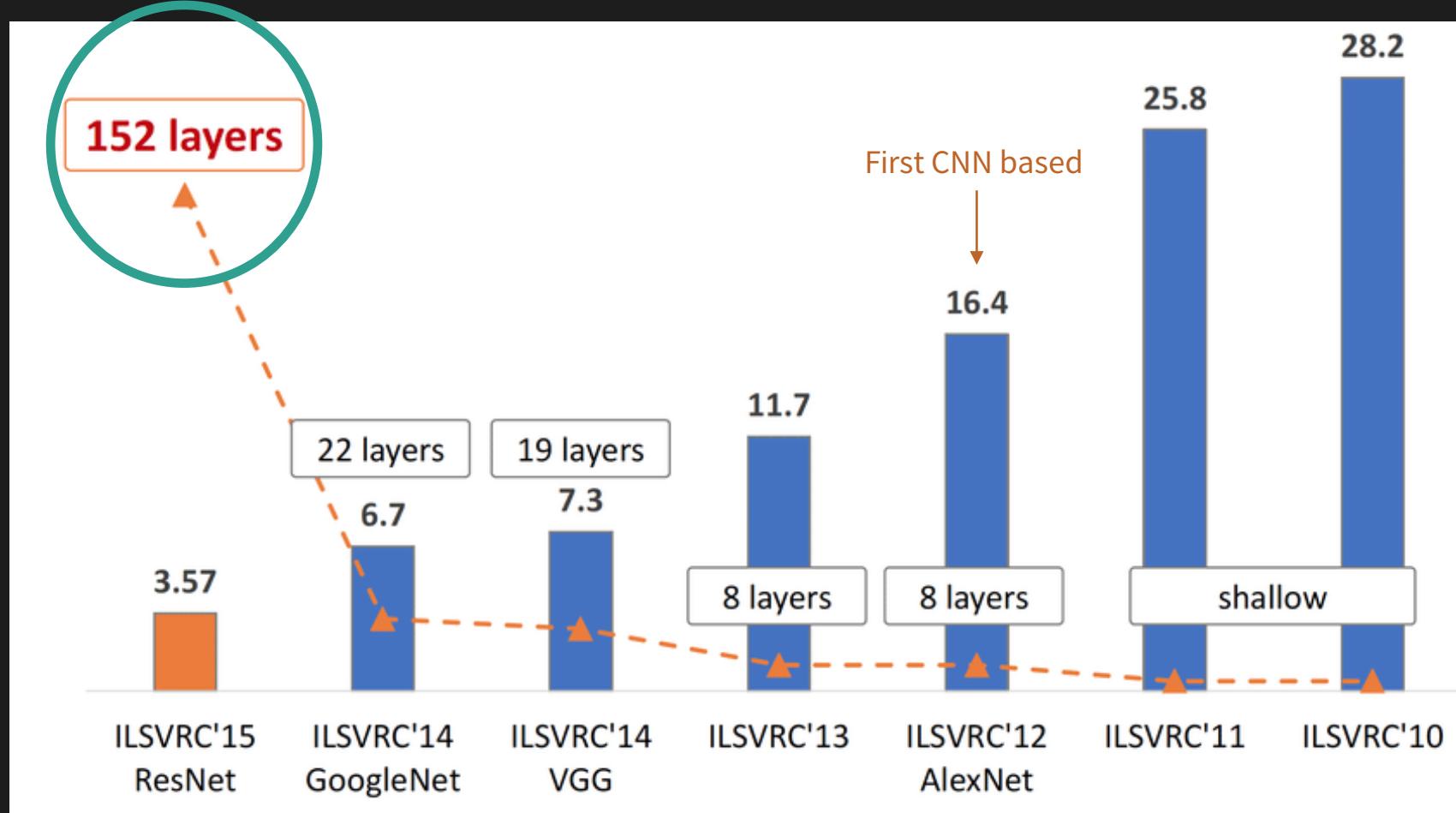
# GoogLeNet



- ‘Network in Network’ (NiN) 구조
  - FC가 많은 Parameter 차지
  - 중간에 FC Layer를 못 넣음
- 1 x 1 Conv로 nonlinearity
- global average pooling

출처 : [Dive Into Deep Learning](#)

# What the hack??



출처 : [https://www.researchgate.net/figure/The-evolution-of-the-winning-entries-on-the-ImageNet-Large-Scale-Visual-Recognition\\_fig1\\_321896881](https://www.researchgate.net/figure/The-evolution-of-the-winning-entries-on-the-ImageNet-Large-Scale-Visual-Recognition_fig1_321896881)

# Much, Much, Much Deeper - Residual Connection

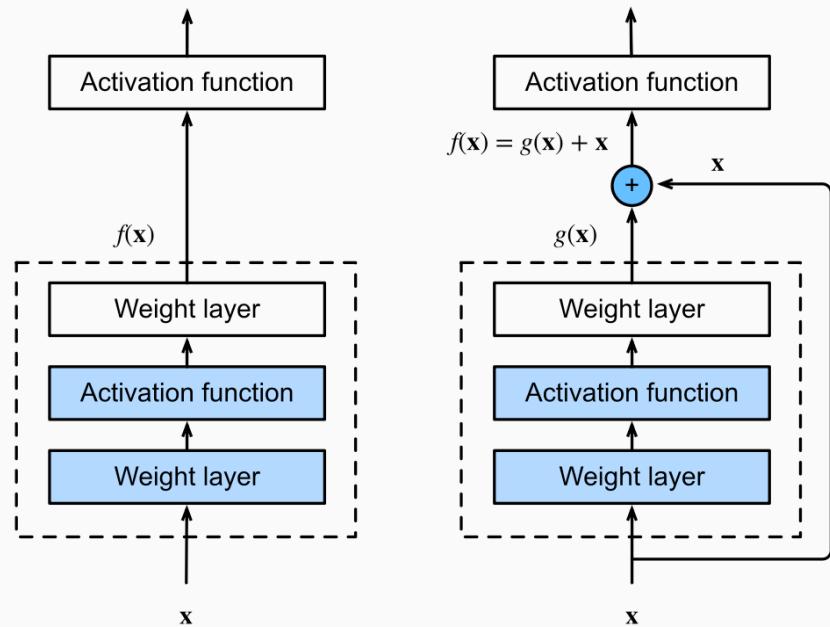


Fig. 8.6.2 In a regular block (left), the portion within the dotted-line box must directly learn the mapping  $f(x)$ . In a residual block (right), the portion within the dotted-line box needs to learn the residual mapping  $g(x) = f(x) - x$ , making the identity mapping  $f(x) = x$  easier to learn.

- Vanishing Gradient Problem
  - Function Classes
- Fit layer to ‘residual’ :  $g(x) = f(x) - x$

출처 : [Dive Into Deep Learning](#)

# ResNet

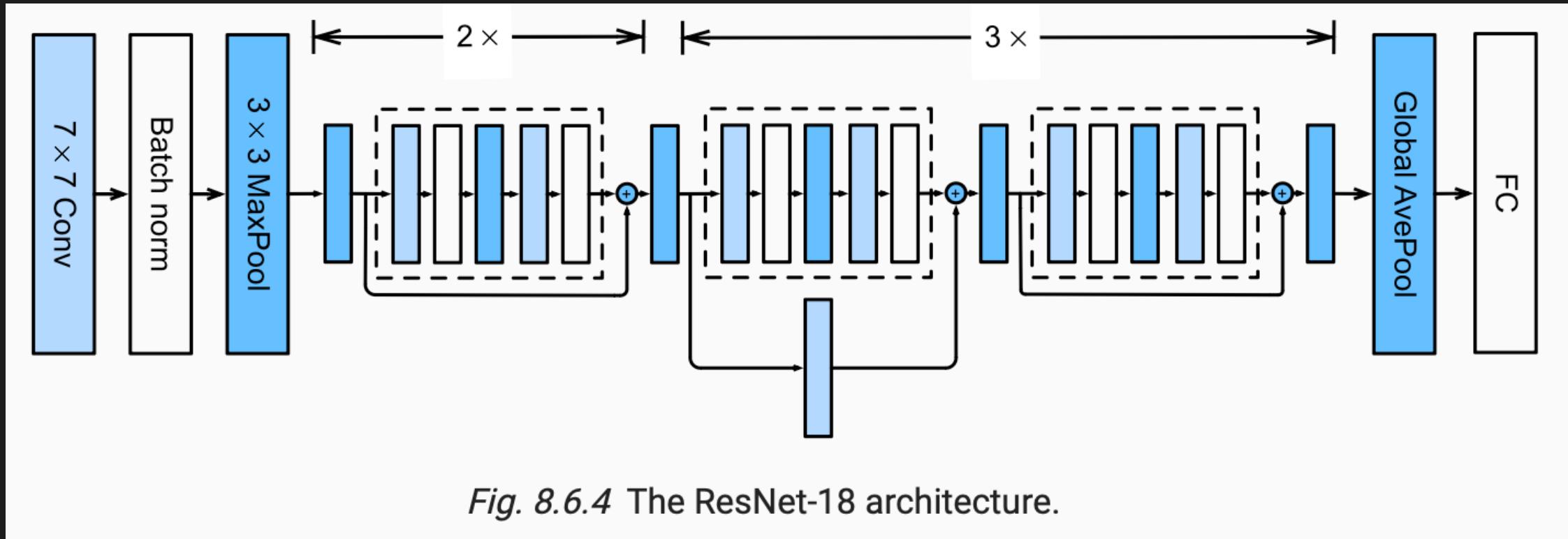


Fig. 8.6.4 The ResNet-18 architecture.

출처 : [Dive Into Deep Learning](#)

# Other Image Tasks

---

- Classic
  - Classification, Object Detection, Semantic Segmentation, Instance Segmentation
- New
  - ViT, Image-Text multimodality, Video, 3D Vision, Generation Model

# Now then...

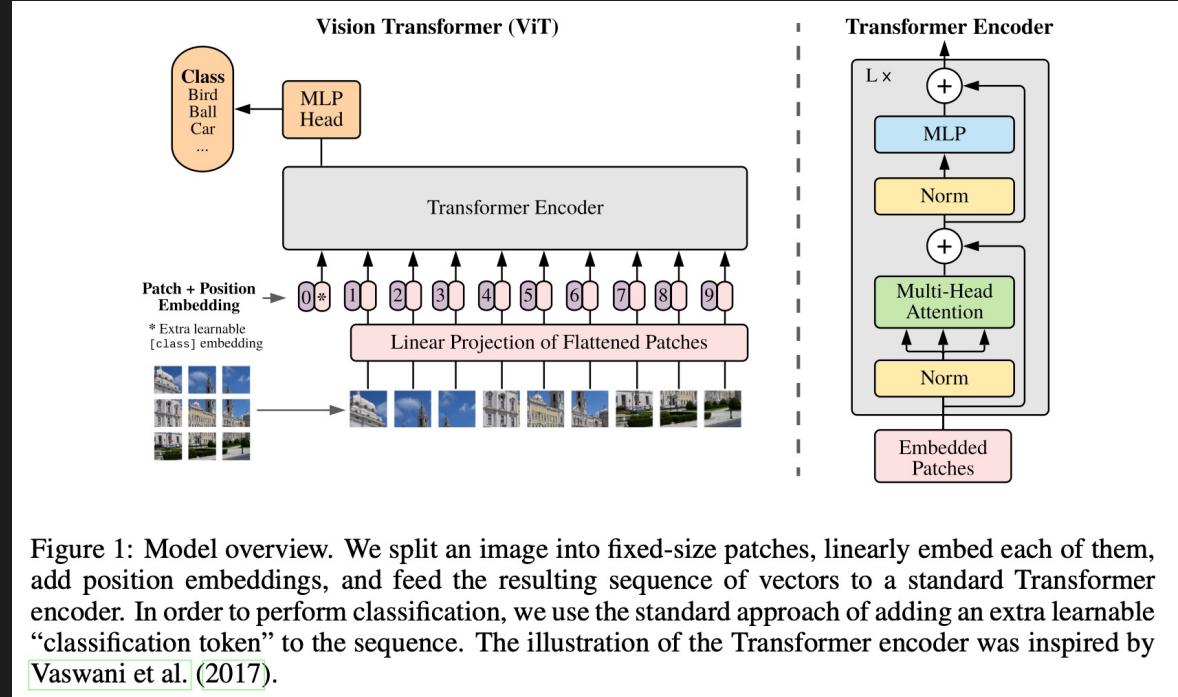
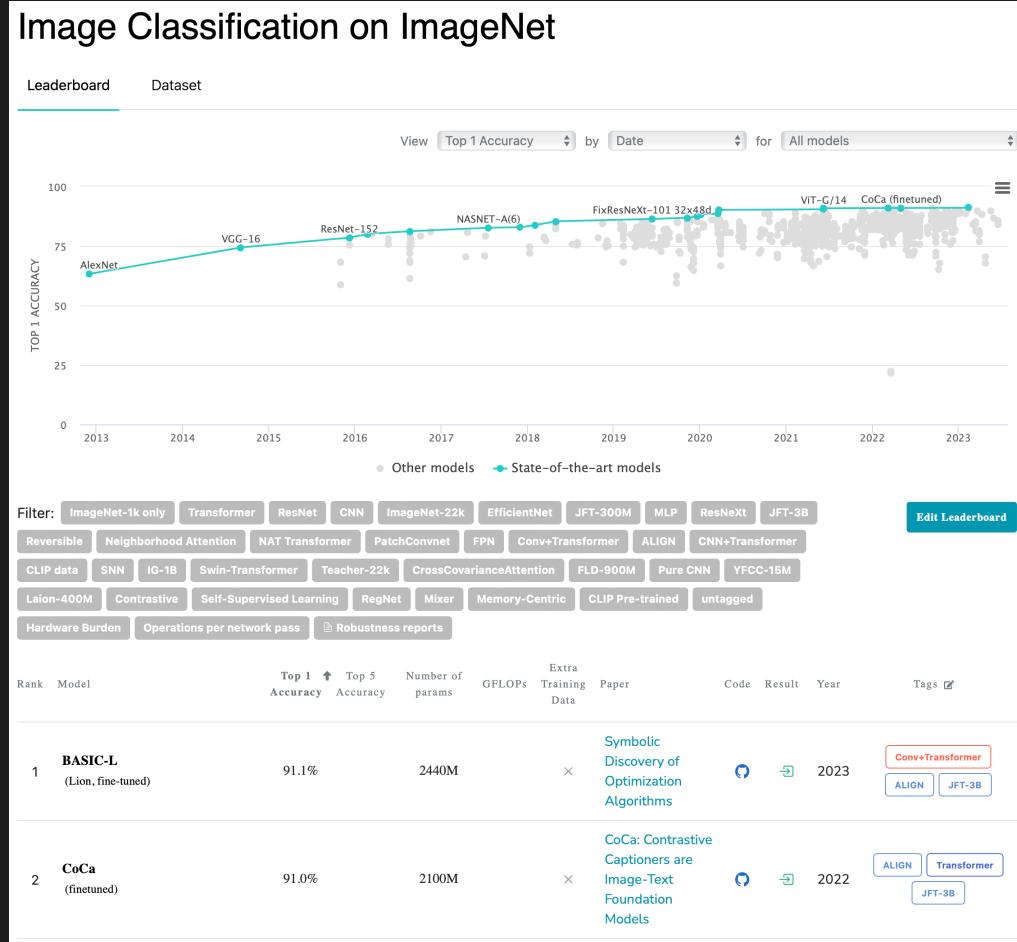


Figure 1: Model overview. We split an image into fixed-size patches, linearly embed each of them, add position embeddings, and feed the resulting sequence of vectors to a standard Transformer encoder. In order to perform classification, we use the standard approach of adding an extra learnable “classification token” to the sequence. The illustration of the Transformer encoder was inspired by Vaswani et al. (2017).

- Transformer 전성 시대!

출처 : [Papers with code](#)

# Discussion Assignment

- CNN을 이미지 이외의 자연어처리 등 분야에 사용할 수 없을까?
- CNN이 가정한 'Spatial invariance'라는 가정은 정말 좋을까?
- Due day ~다음 수업 시작 전까지
  - AIKU Notion → 일정 → DeepIntoDeep 해당 회차
  - 본인 칸에 자신의 생각 적기!

# PyTorch hands-on

# PyTorch?



출처 : <https://tutorials.pytorch.kr/>

- 오픈소스 딥러닝 라이브러리
- Tensor를 GPU 상에 올려 연산 가능.
- 추천하는 사이트 : [파이토치 한국어 튜토리얼](#)
- [Documentation](#)
- Rival : TensorFlow
  - PyTorch가 모델을 수정하기 더 쉬워 연구자들 사이에서는 더 일반적.

# 개발 환경



출처 : Google

- 서버에서 작업한다면 PyTorch + CUDA
  - 보통 Anaconda를 사용해 가상 환경 사용
- 하지만… PyTorch 버전 관리가 어려움
  - Colab!

# Tensor

- Tensor
  - array, matrix.
    - Numpy의 ndarray와 유사. GPU 상에서 돌아간다는 차이만! (.to(device))
    - Autograd에 최적화 됨.
  - 참고 : Tensor 차원 관리가 어렵다면? [einops](#)

```
import torch
import numpy as np

# Directly
data = [[1, 2], [3, 4]]
x_data = torch.tensor(data)

# From Numpy
np_array = np.array(data)
x_np = torch.from_numpy(np_array)
```

# Data, DataLoader

---

- Data를 불러오고, 학습에 사용하기 편하도록 묶어주는 class
  - Custom dataset : `__init__`, `__len__`, `__getitem__`
  - 과제에서는 주어짐
- 전반적인 학습의 과정
  - Dataloader 만들기
  - 모델 정의하기
  - train 함수 만들기
- 튜토리얼 꼭 읽어 보시길 추천 드립니다!

# FAQ

---

- PyTorch Document 어디를 보면 되나?
  - [`torch.nn`](#), [`torch.nn.functional`](#), [`torch.optim`](#)을 주로 봅니다.
  - [`Dataloader`](#), `train` 함수는 비슷하게 반복되어 사용됩니다.
- nn과 nn.functional은 뭐가 다른가?
  - nn은 Class로 정의되어 있고, functional은 함수로 정의되어 있습니다.
- `model.eval()`이 무슨 뜻인가요? `torch.no_grad()`는?
  - eval : model을 test 모드로 전환해 줍니다.
  - no\_grad : test 할 때에는 model의 Gradient를 추적하는 것을 꺼 버립니다. [참고](#)

# 과제 설명

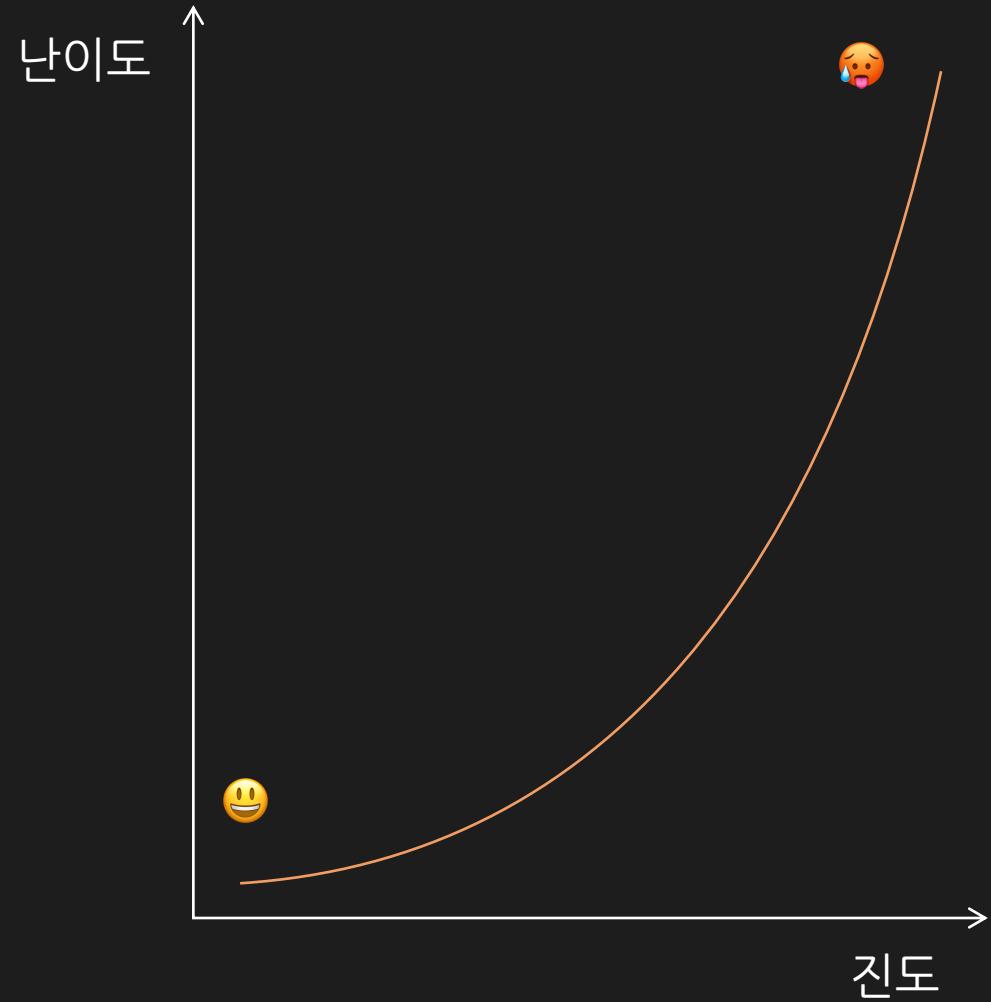
- Due Day ~2023.07.23. 제출은 [github classroom](#)
- 과제 다운로드, 본인 Google Drive에 과제 업로드
- Colab으로 열기.
  - 런타임 → 런타임 유형 변경 → 하드웨어 가속기 ‘GPU’
- 맨 위 코드 블럭에 ‘FOLDERNAME’ 변경하기
  - 자신이 업로드 한 경로 제대로 적기!! → 안 되면 문의
- 위에서부터 하나씩 실행해 가면서 코드 이해해 보기
- ‘TODO’로 감싸져 있는 부분만 코드 추가하기!
- 중간 중간에 있는 설명 작성 부분에 자기 생각 적기

# 주목해야 할 Class / Function

---

- Class
  - [Linear](#)
  - [Conv2D](#)
  - [SGD](#)
  - [Dropout2D](#)
  - [MaxPool2D](#)
  - [BatchNorm2D](#)
- Function
  - [Conv2D](#)
  - [ReLU](#)

# Welcome to Deep Learning



- 지금까지 배운 것
  - 딥러닝의 시작
  - CNN
- 앞으로 배울 것
  - 만약 데이터가 Sequence가 있다면?  
RNN
  - Attention is all you need.  
Transformer
  - 다양한 Vision, NLP Model들
  - **Welcome to Deep Learning!**

# 감사합니다!