

# Machine Learning

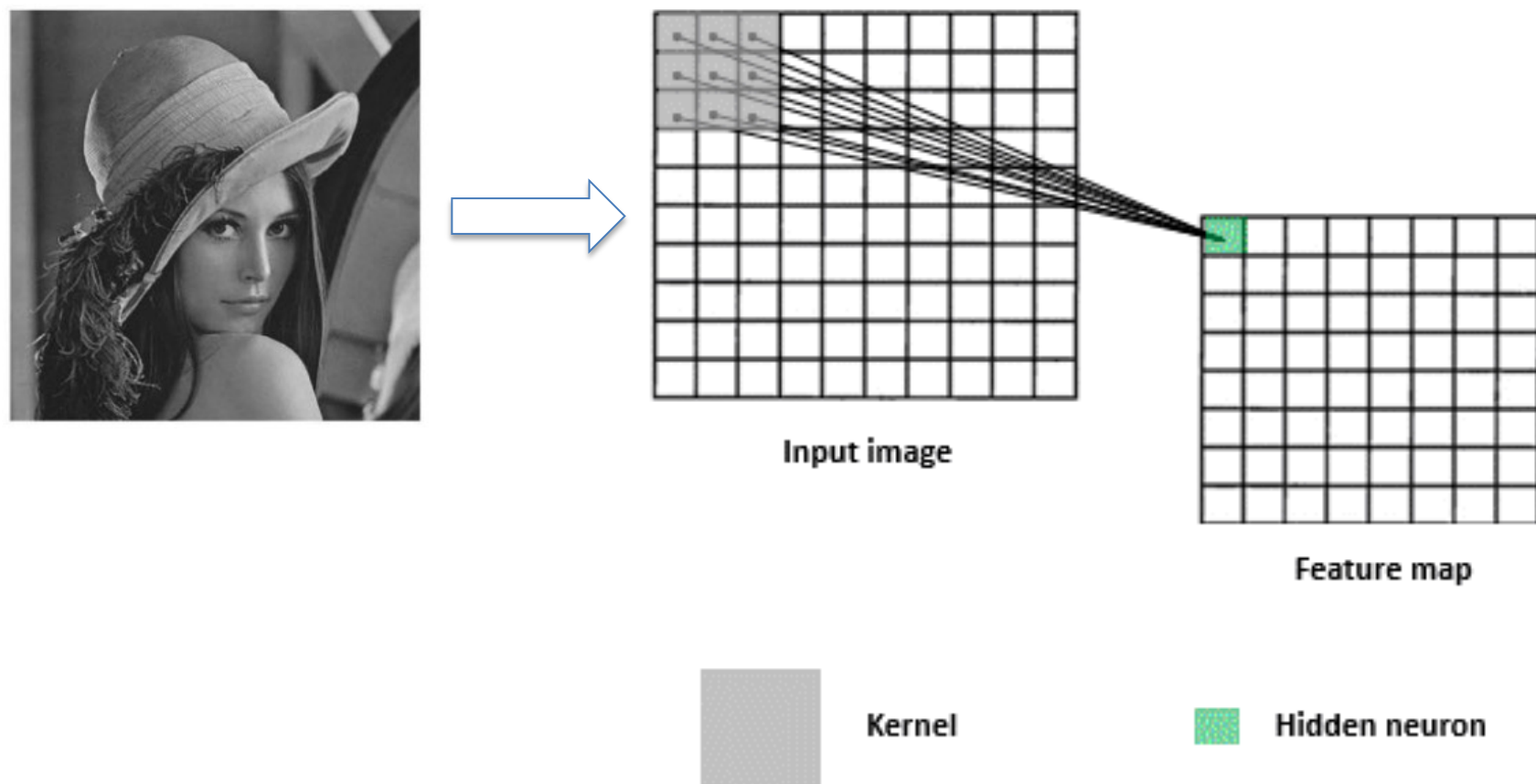
## Deep Learning

Instructor: Prof. Yi Fang  
yfang@nyu.edu

# CNN Architectures

- Basics of CNN
- Classic Networks

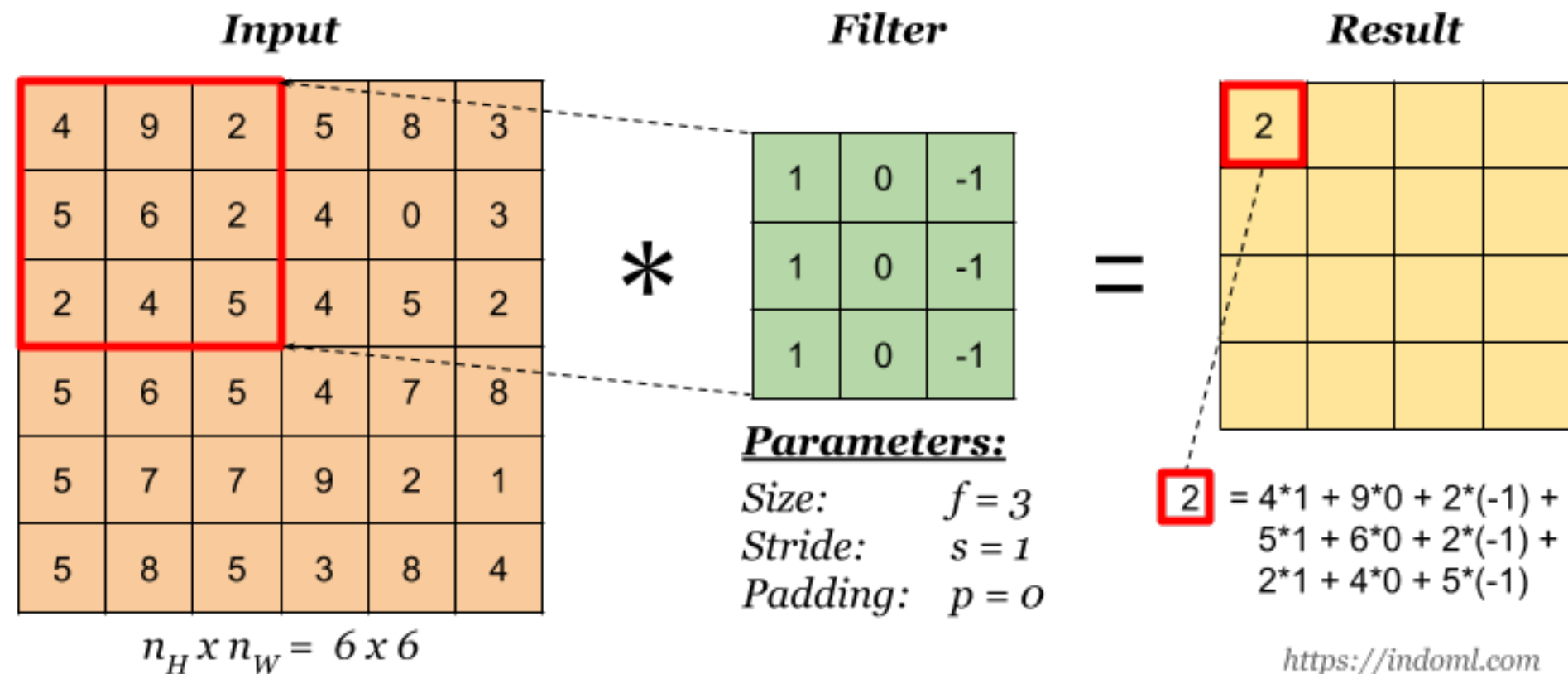
# Standard Convolution Operation

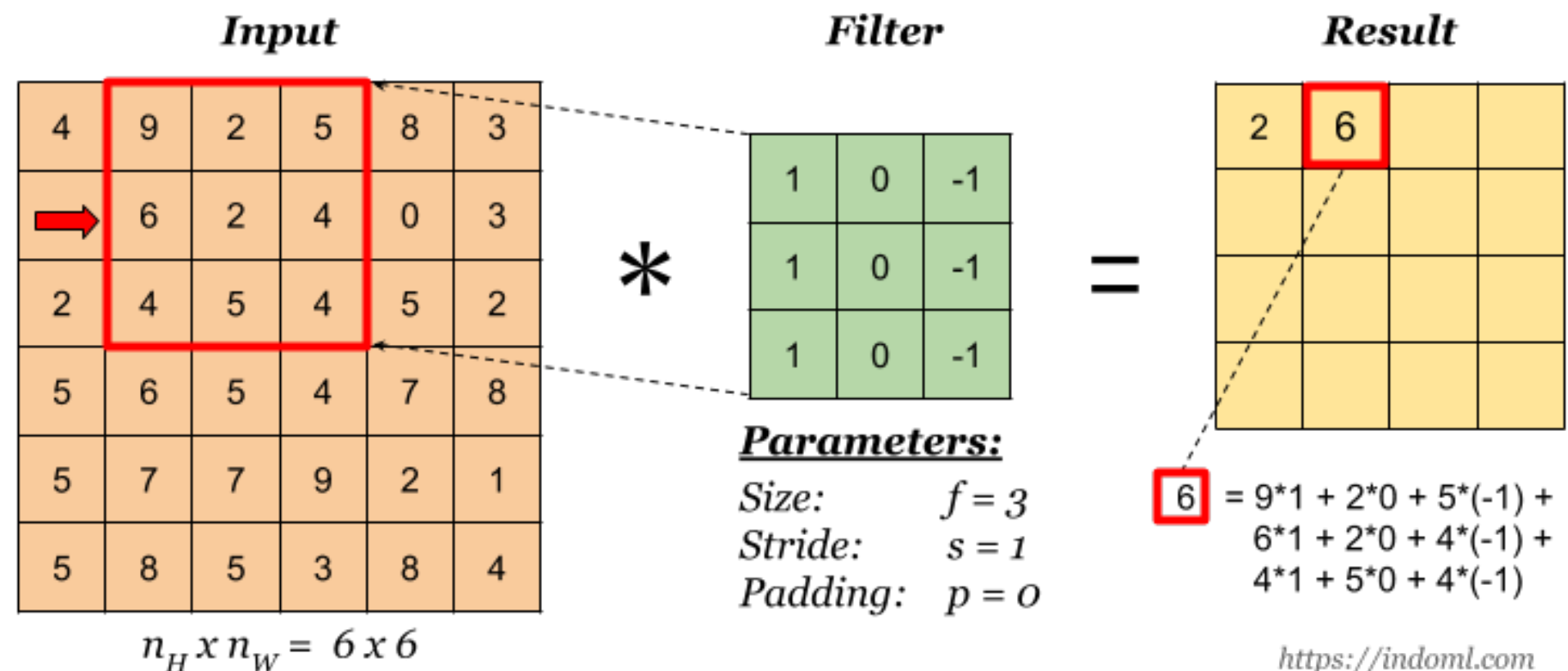


Source: <https://www.ahmedbesbes.com/blog/introduction-to-cnns>

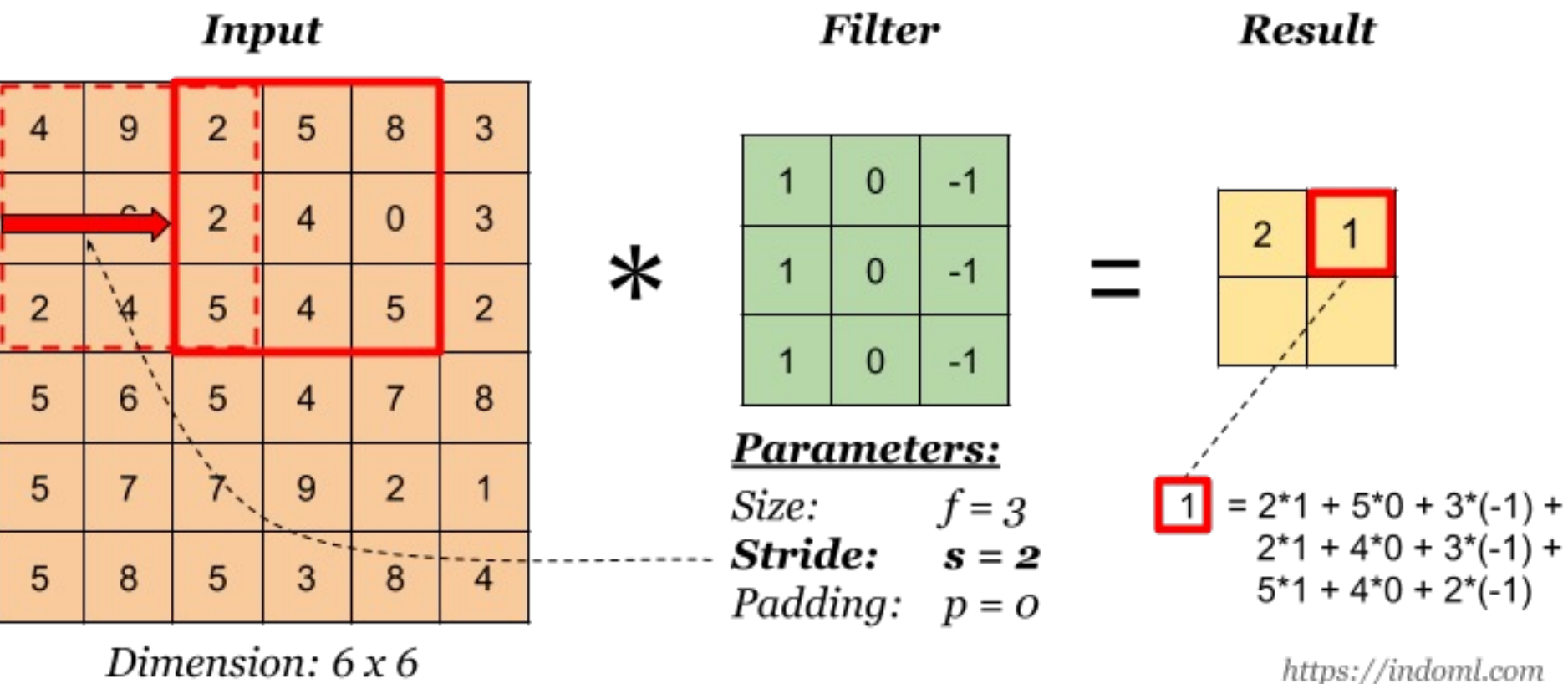
# CNN Basics

## Basic Convolution Operation

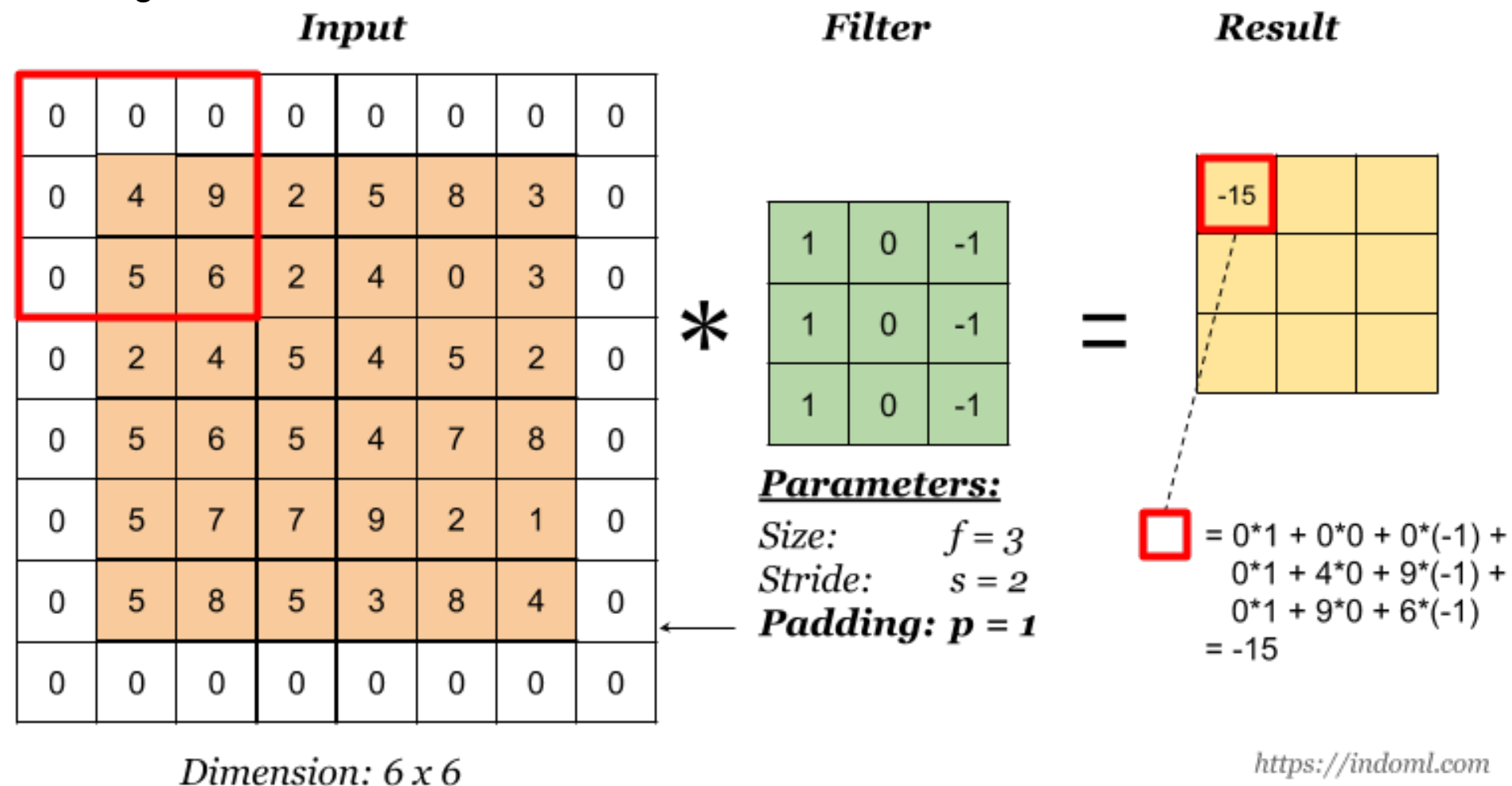




Stride:

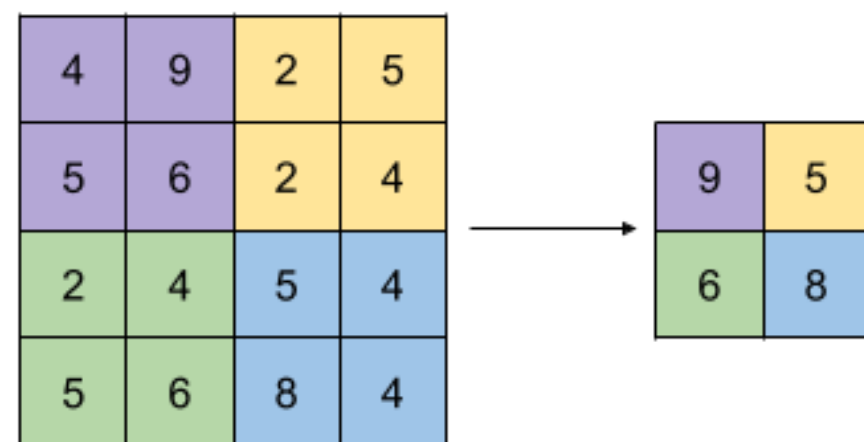


Padding:

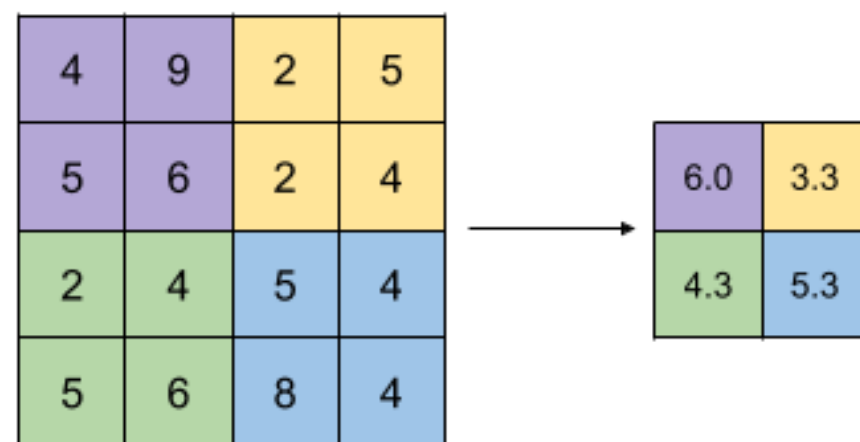


Pooling layer:

*Max Pooling*

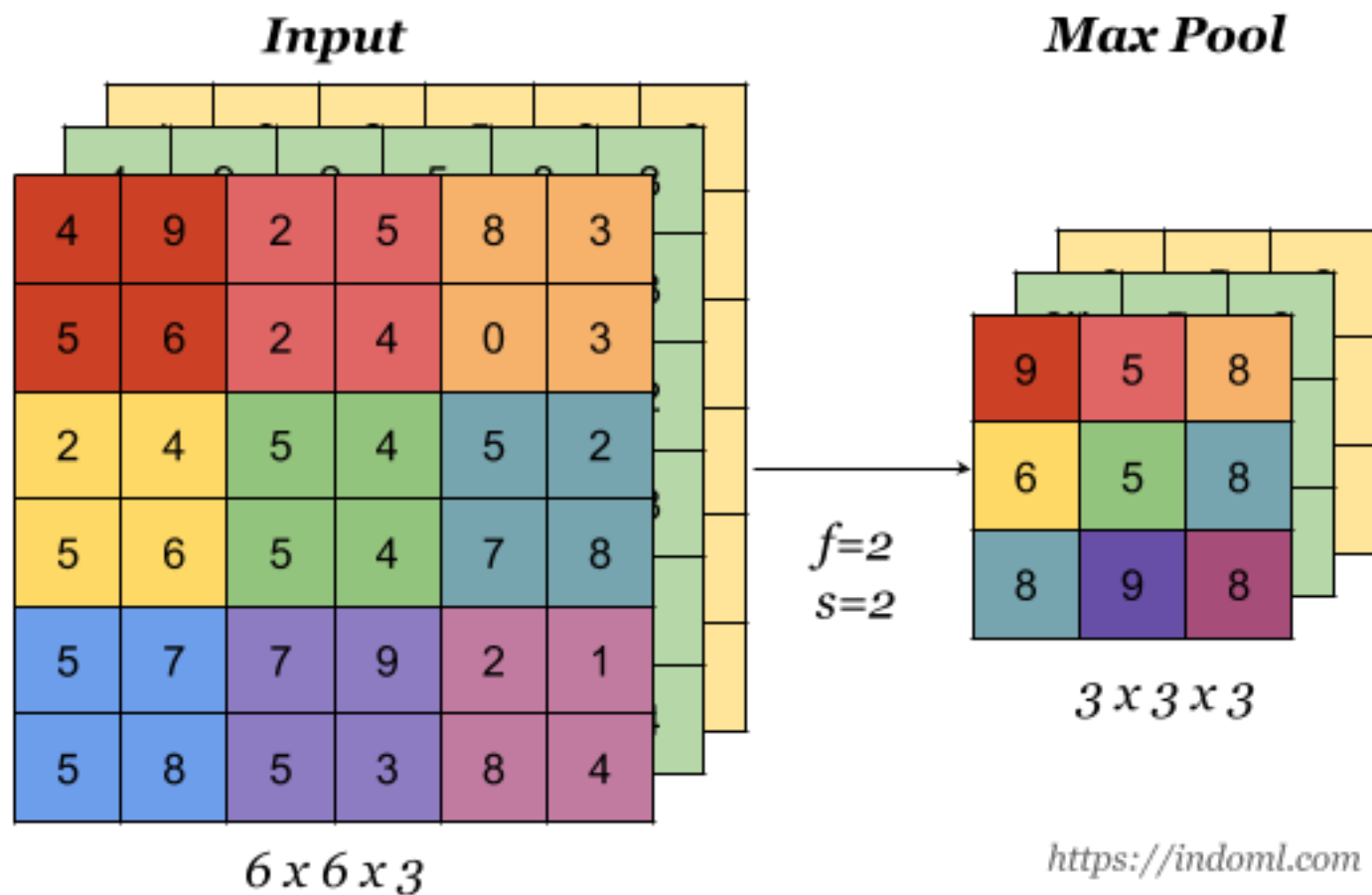


*Avg Pooling*



<https://indoml.com>





Softmax layer:

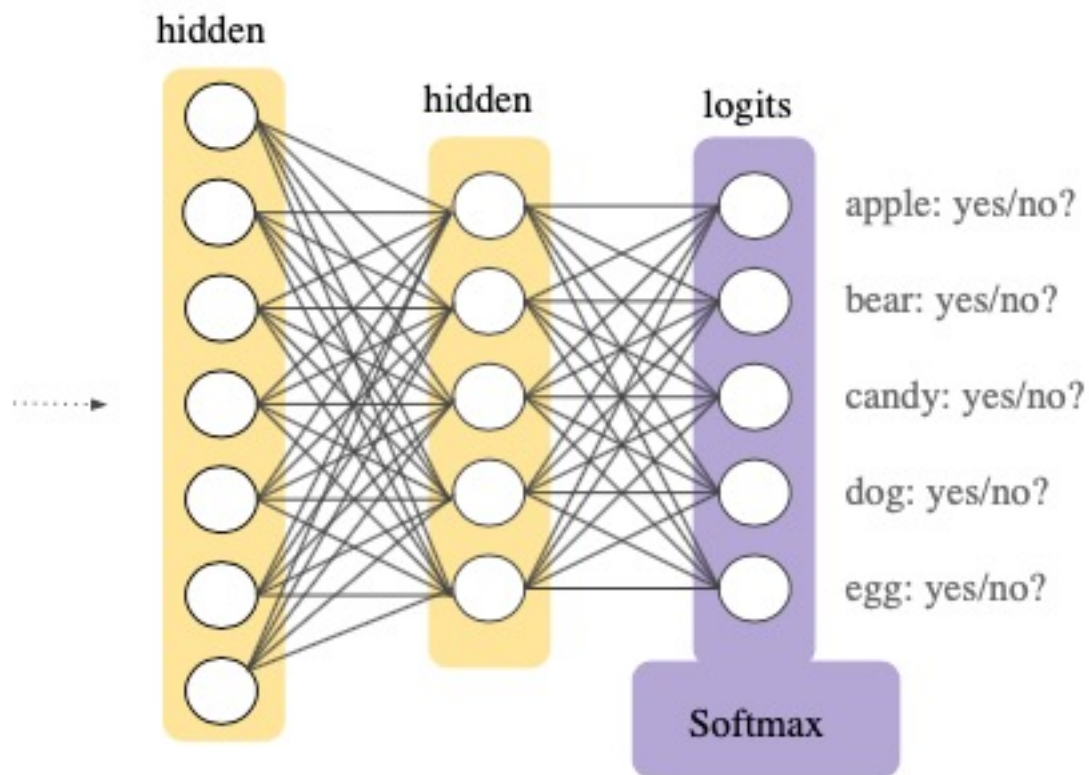
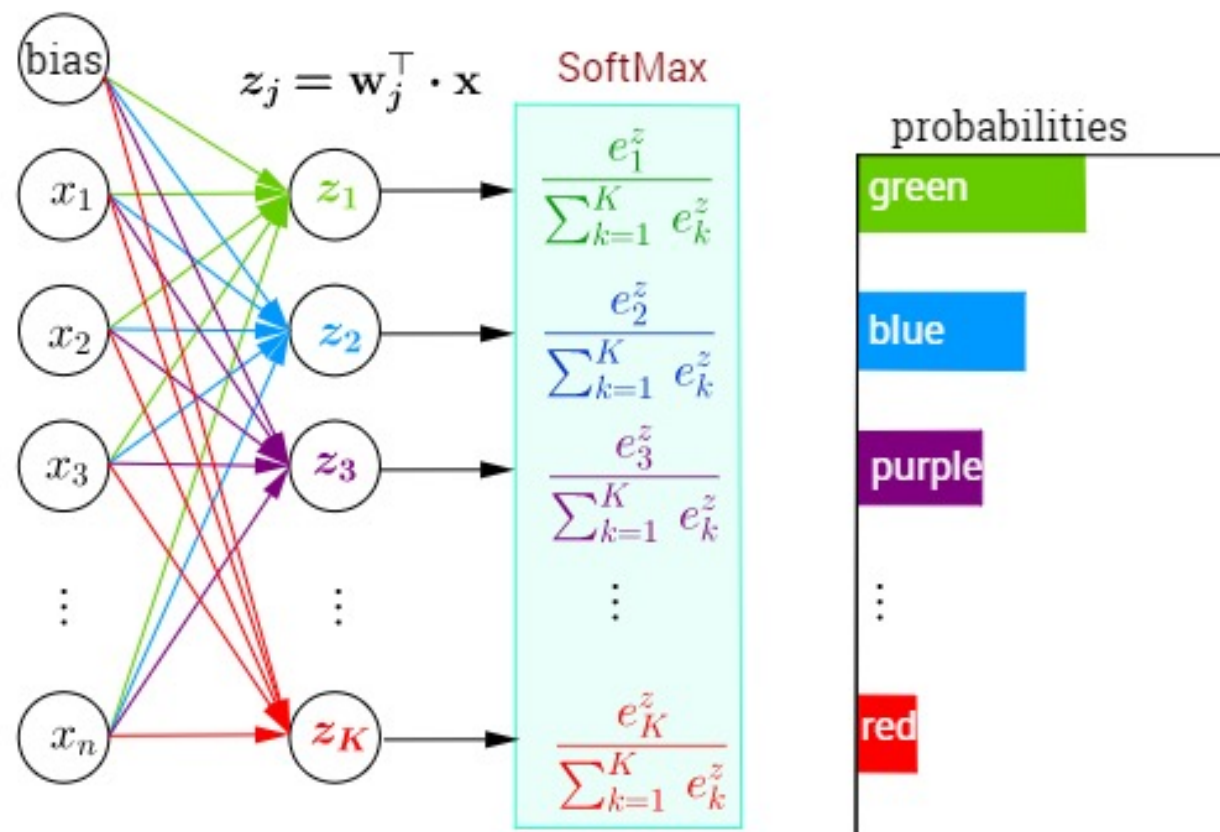


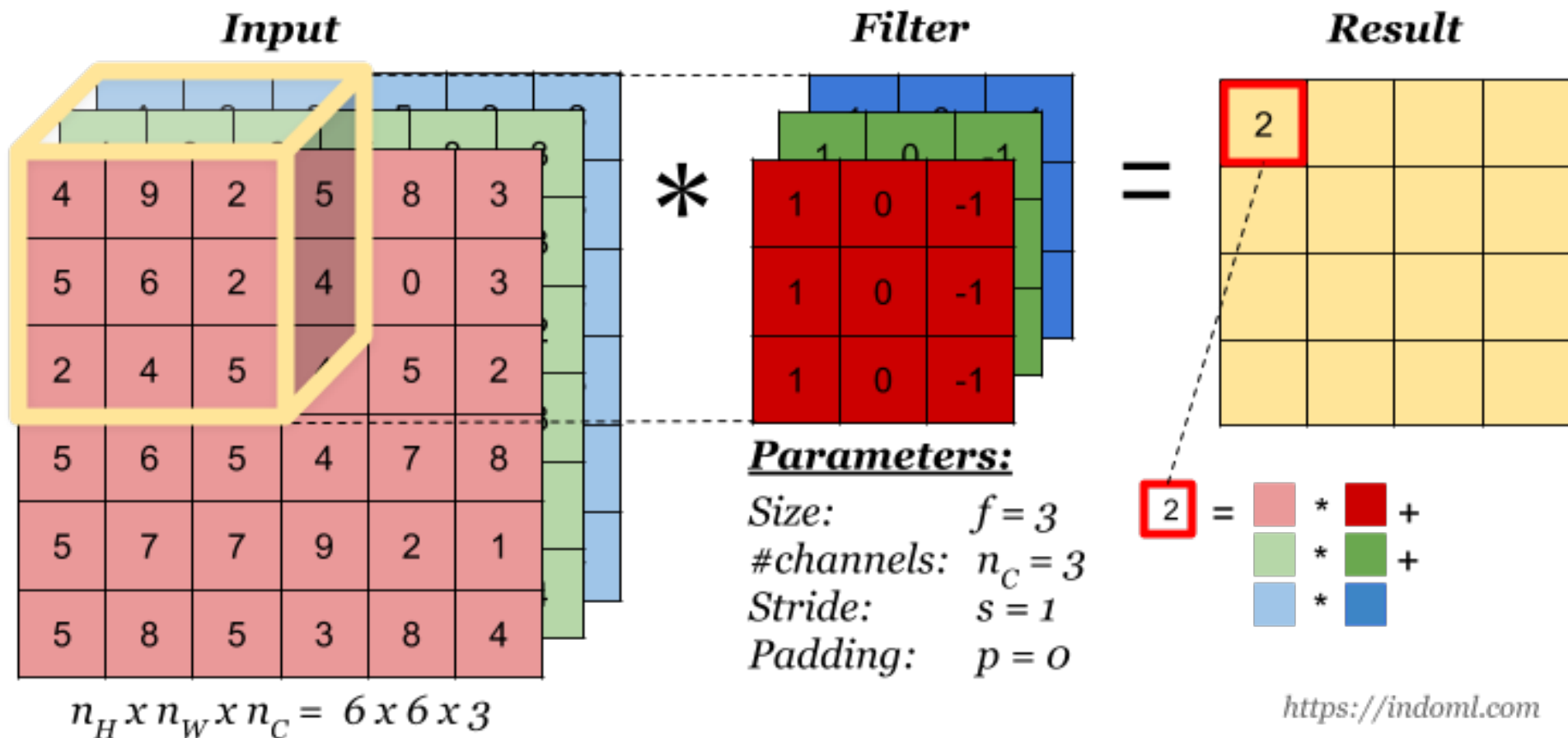
Image source: <https://developers.google.com/machine-learning/crash-course/multi-class-neural-networks/softmax>

$$\mathbf{z} = \begin{bmatrix} z_1 \\ z_2 \\ z_3 \\ \vdots \\ z_K \end{bmatrix} = \begin{bmatrix} \mathbf{w}_1^\top \\ \mathbf{w}_2^\top \\ \mathbf{w}_3^\top \\ \vdots \\ \mathbf{w}_K^\top \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{bmatrix}$$

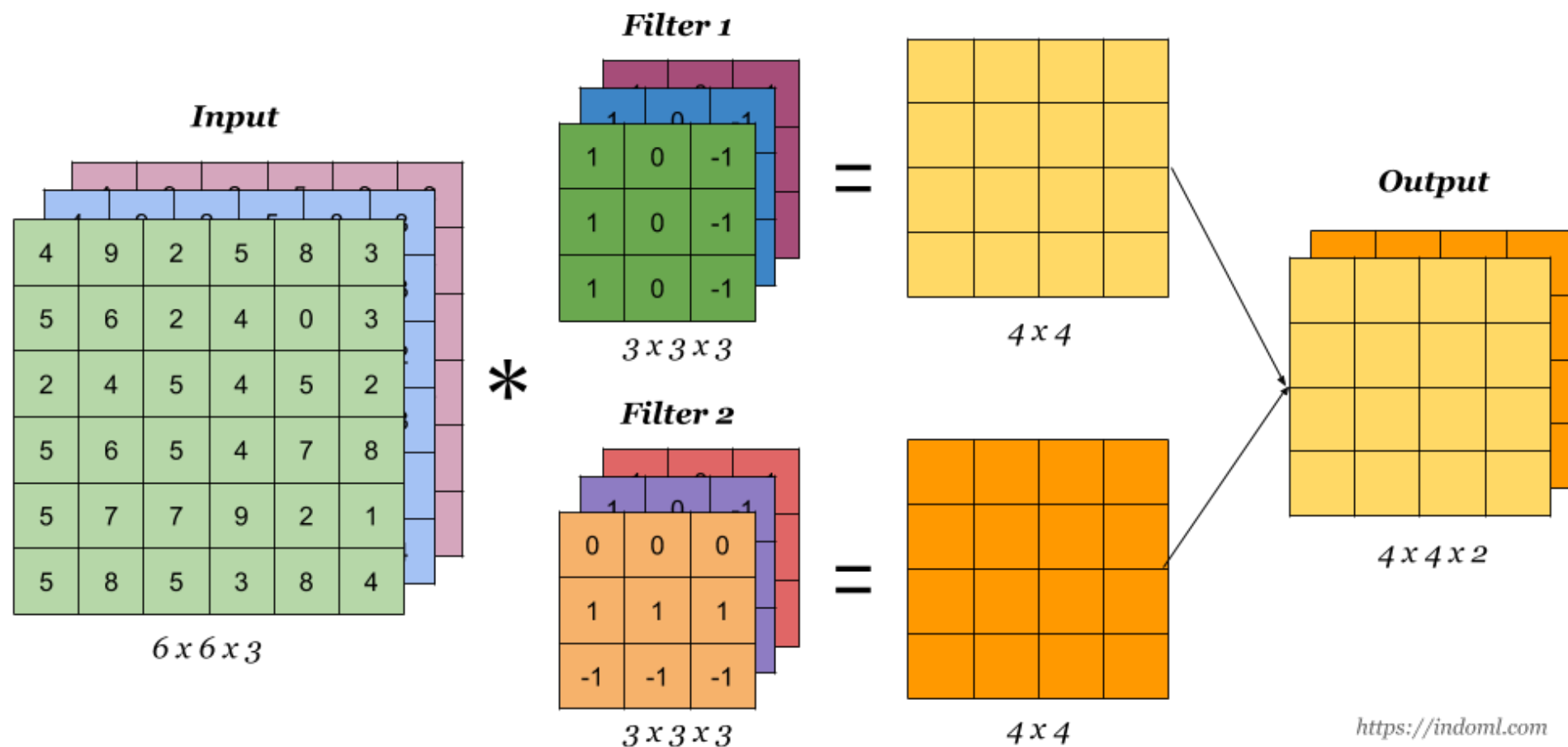


Source: <https://stats.stackexchange.com/questions/273465/neural-network-softmax-activation>

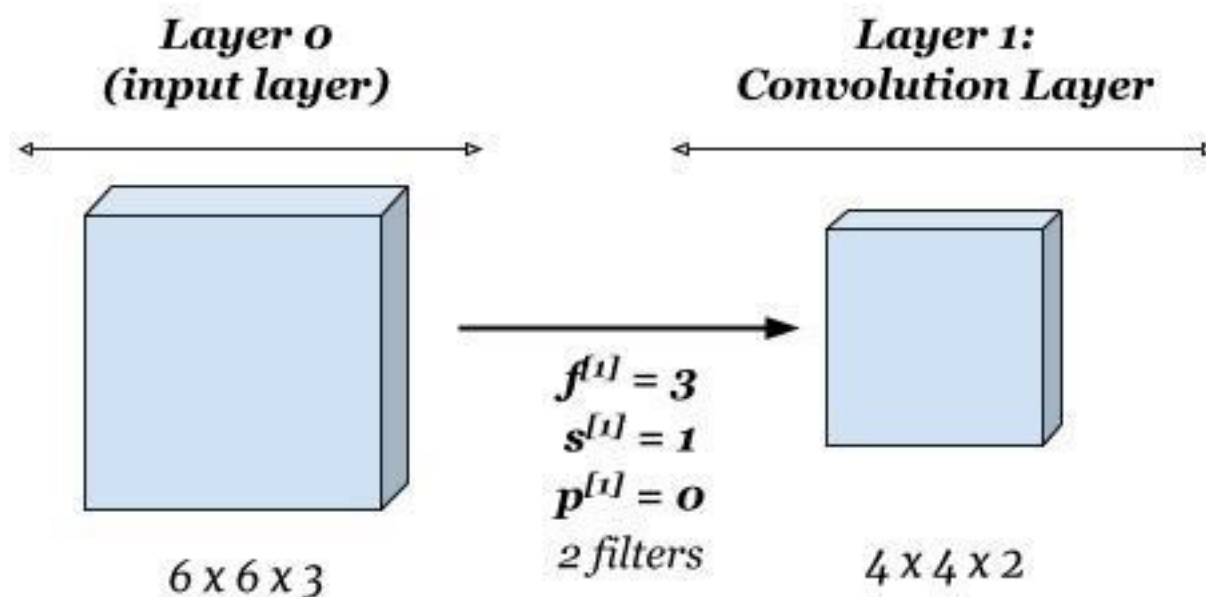
# CNN on Volume (Tensors)



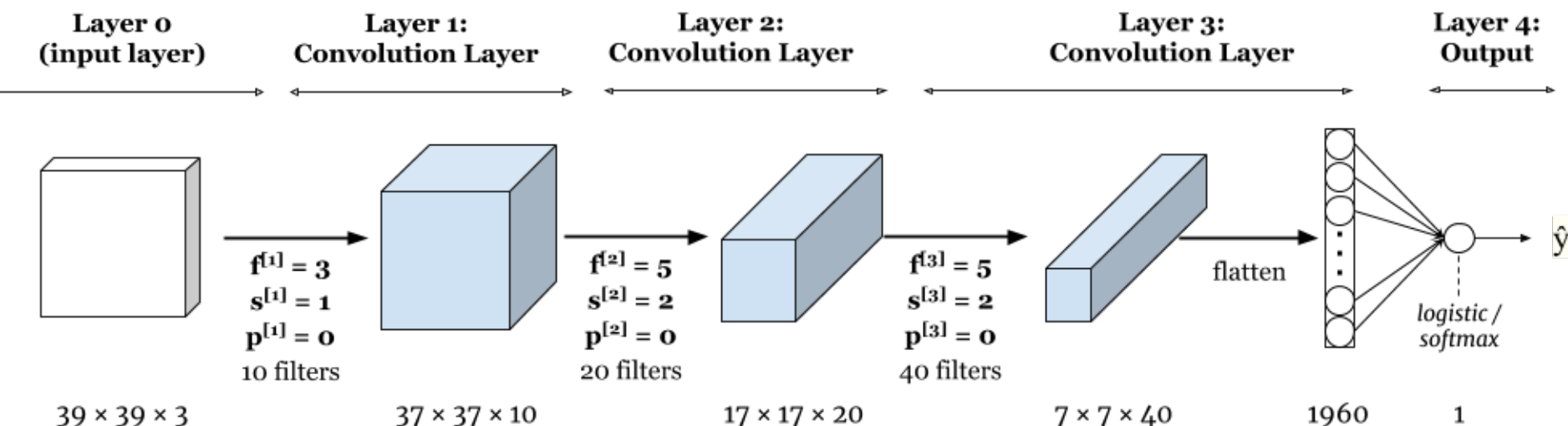
# Multiple Filters



# Shorthand Representation

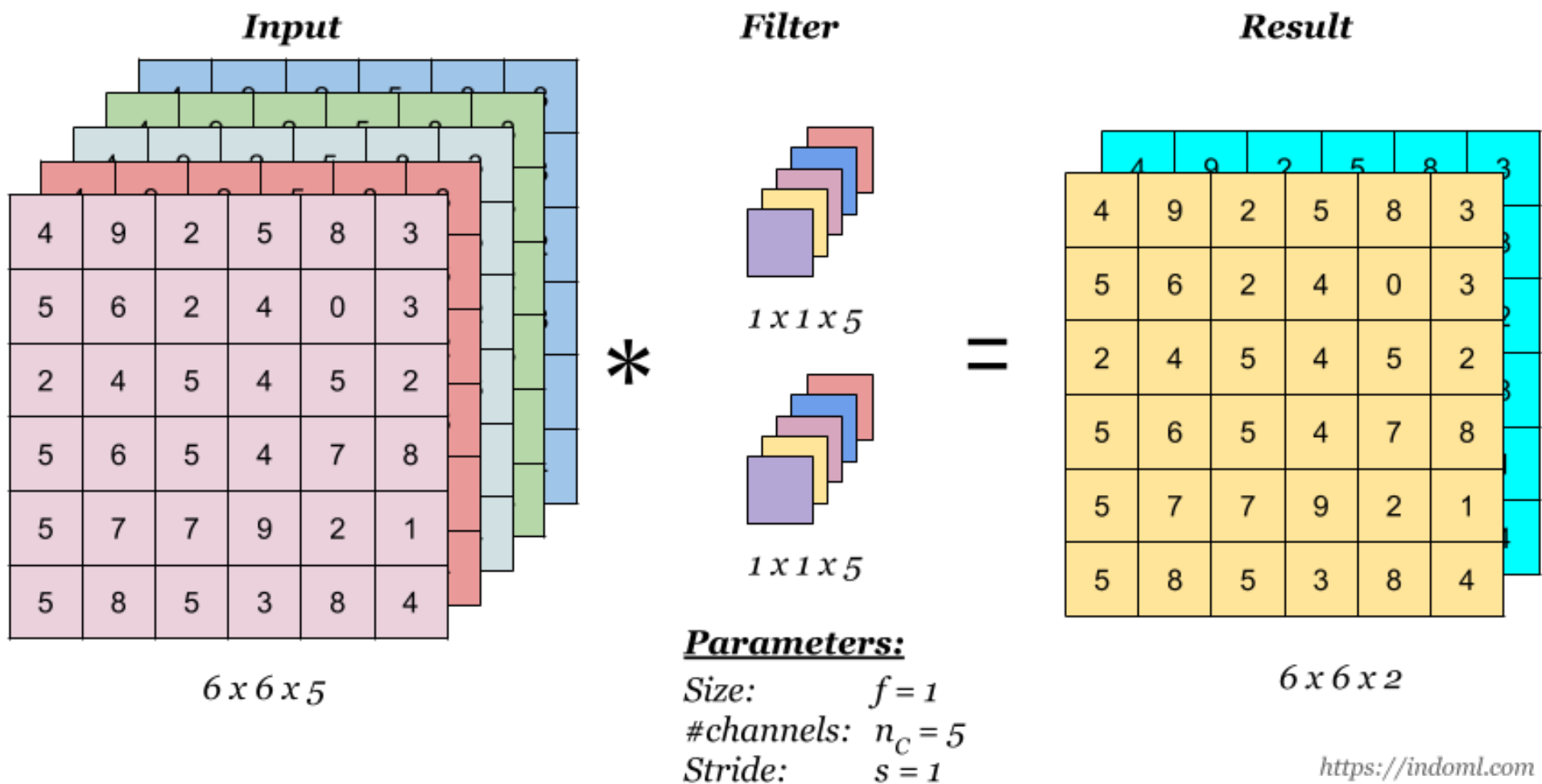


# Sample Network Structure



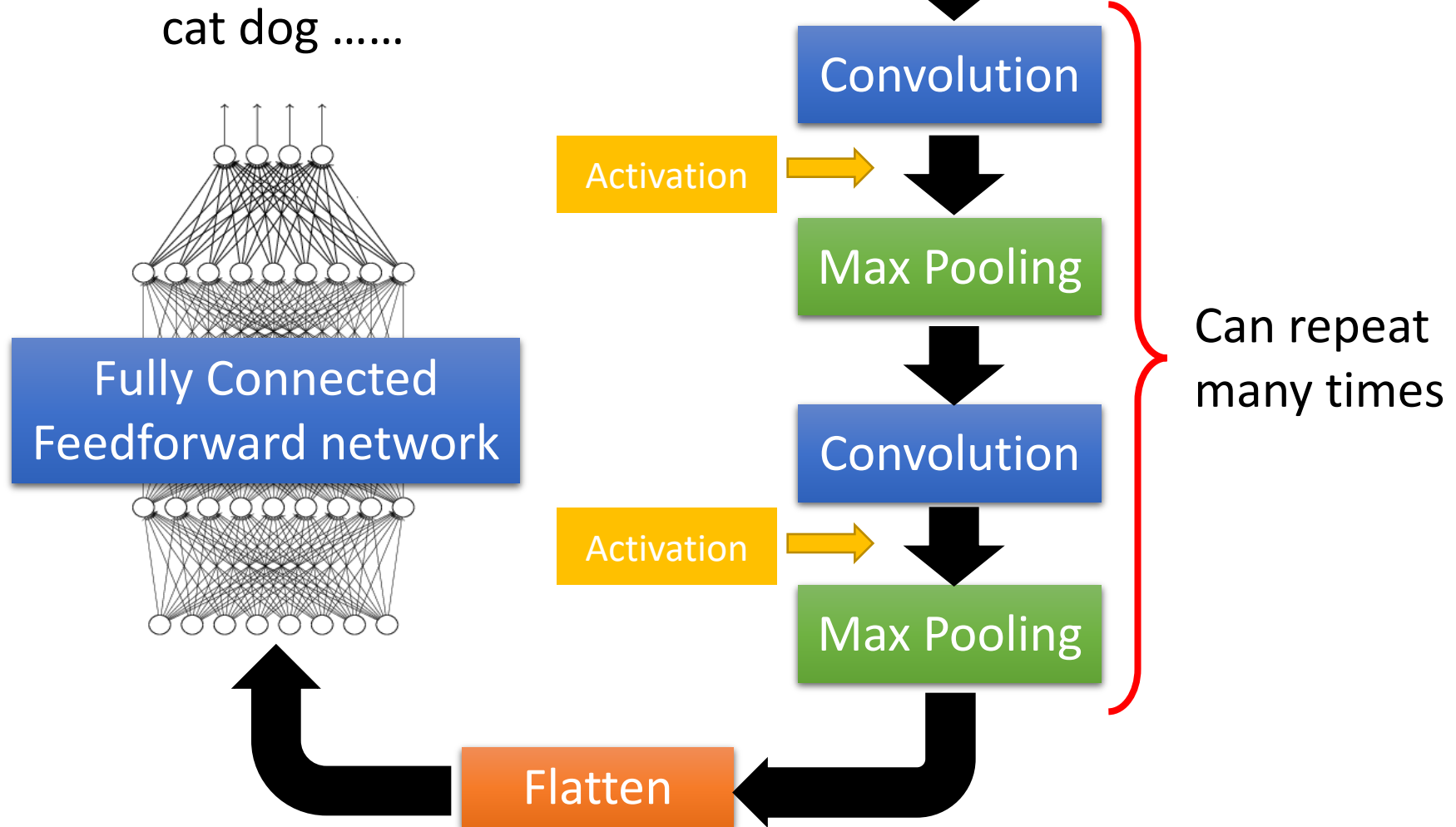
<https://indoml.com>

# 1\*1 Convolutions

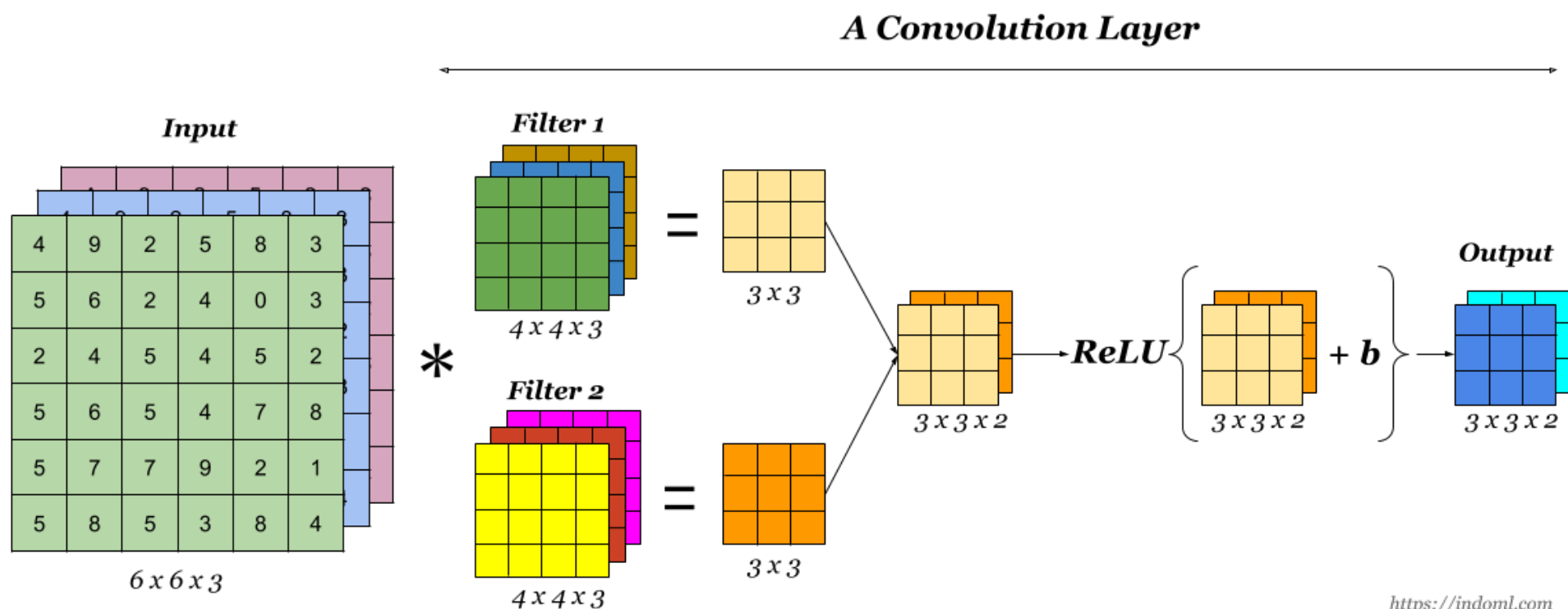




# The whole CNN



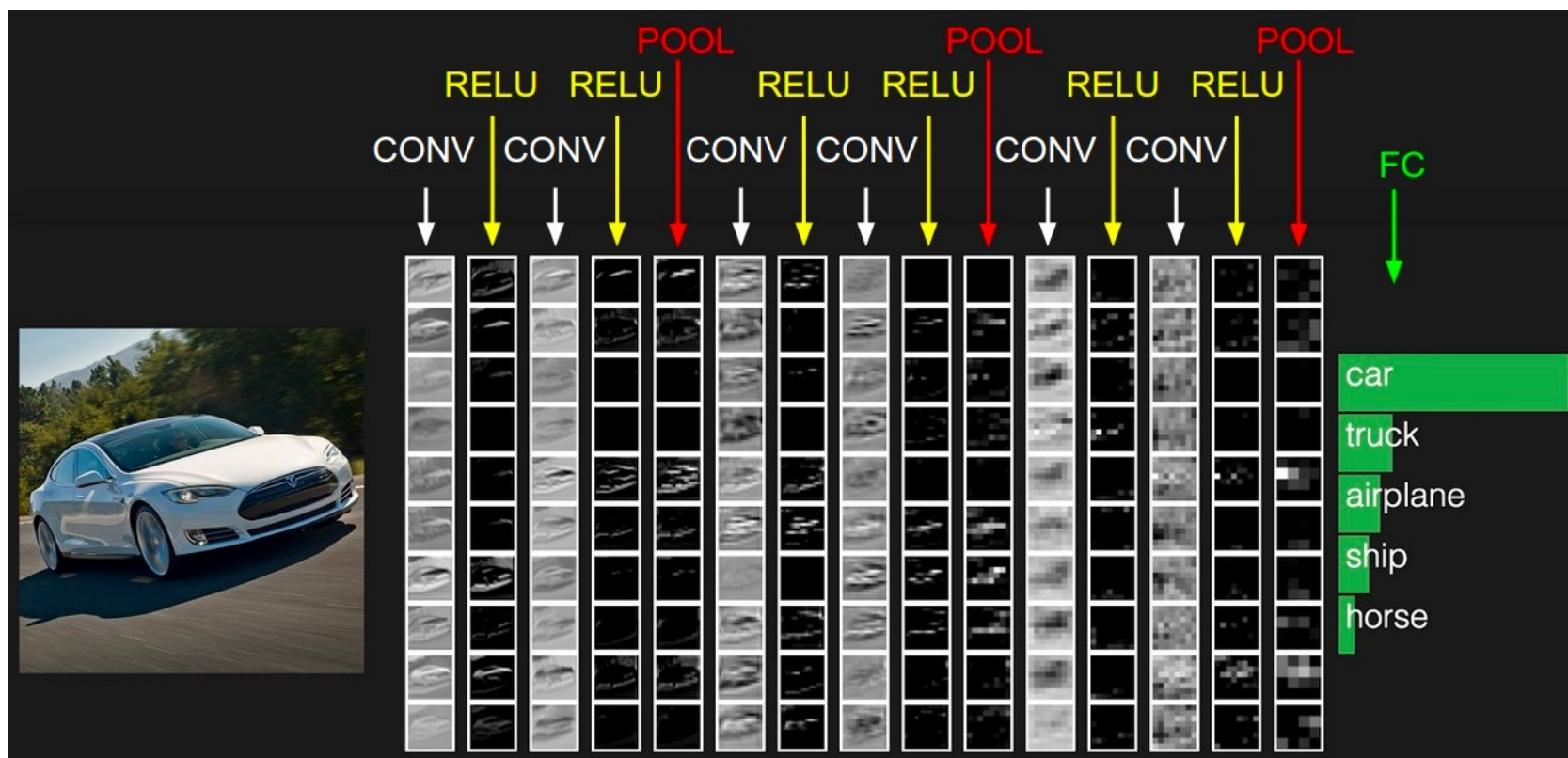
# One Layer Representation



<https://indoml.com>

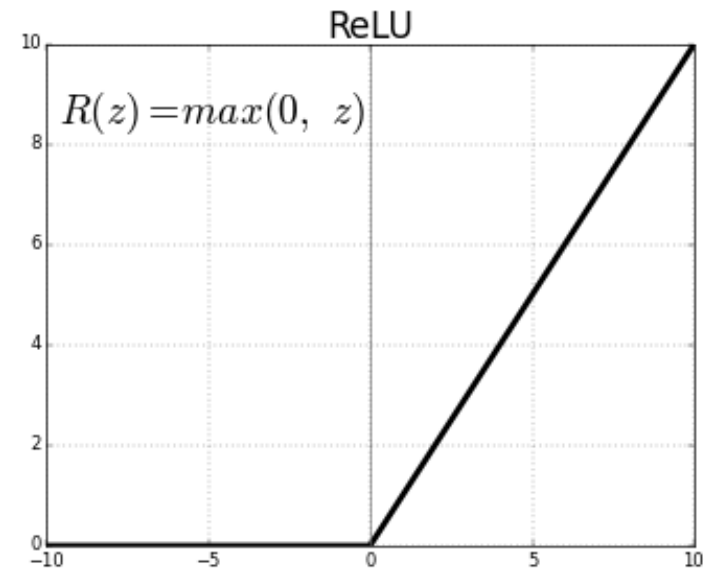
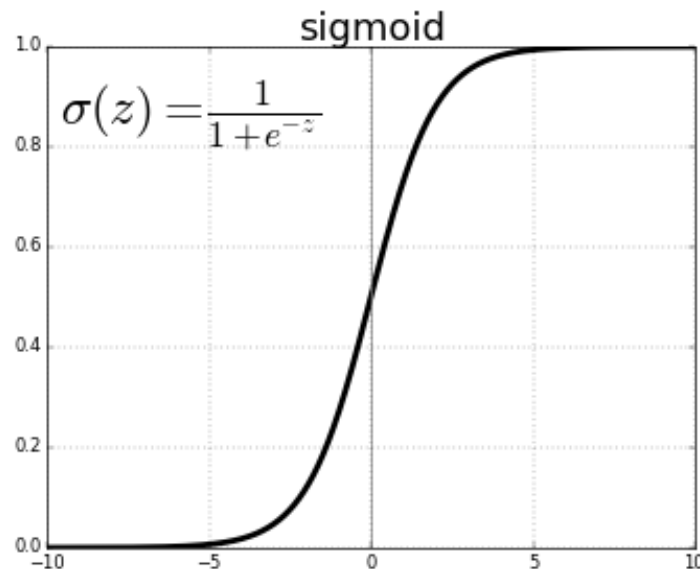
# Convolution Neural Network










- An example of ConvNet architecture.



# Activation Functions

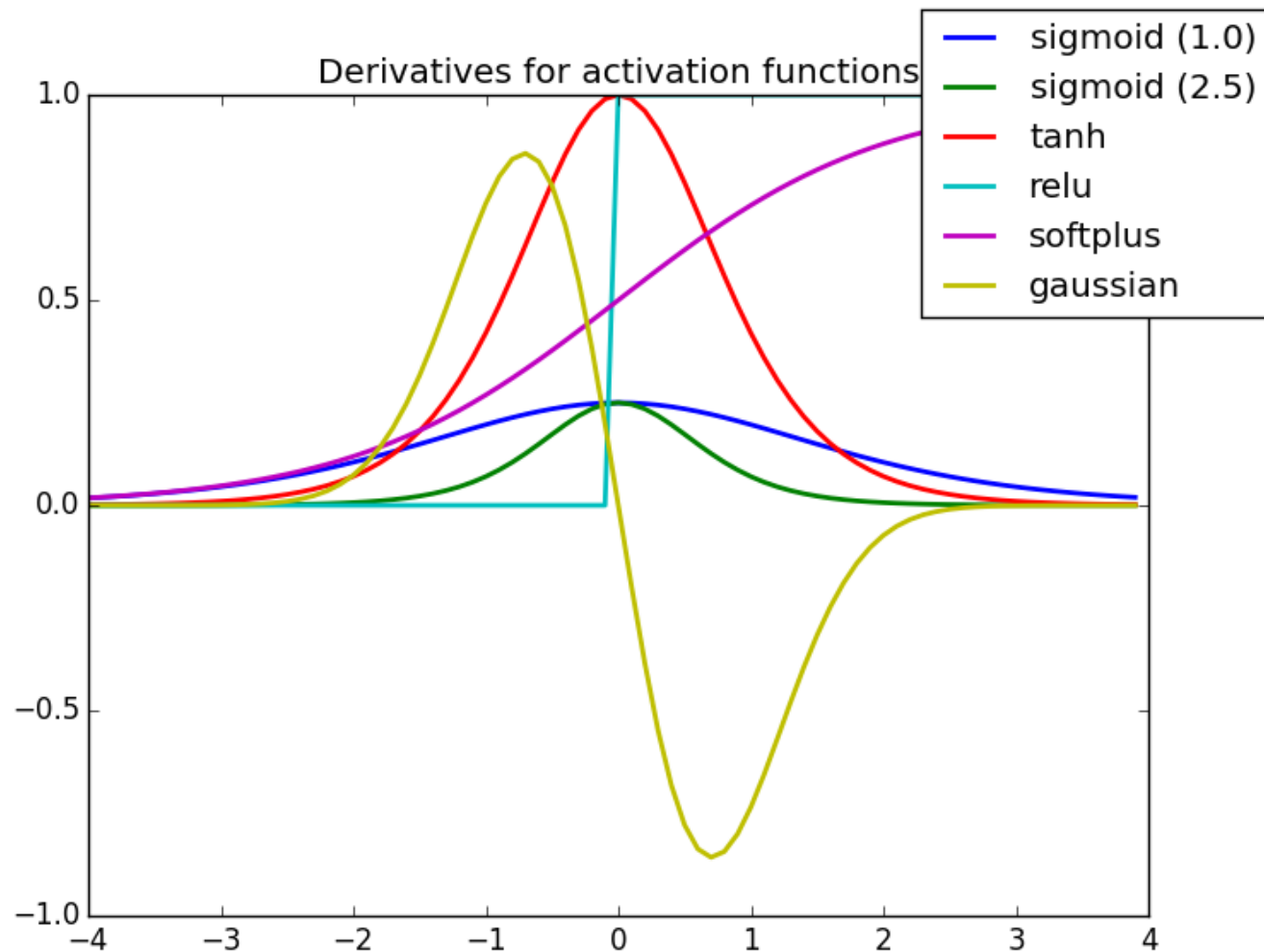
- Sigmoid Function
- ReLU Function
- More



Name	Plot	Equation	Derivative
Identity		$f(x) = x$	$f'(x) = 1$
Binary step		$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} 0 & \text{for } x \neq 0 \\ ? & \text{for } x = 0 \end{cases}$
Logistic (a.k.a Soft step)		$f(x) = \frac{1}{1 + e^{-x}}$	$f'(x) = f(x)(1 - f(x))$
Tanh		$f(x) = \tanh(x) = \frac{2}{1 + e^{-2x}} - 1$	$f'(x) = 1 - f(x)^2$
ArcTan		$f(x) = \tan^{-1}(x)$	$f'(x) = \frac{1}{x^2 + 1}$
Rectified Linear Unit (ReLU)		$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$
Parameteric Rectified Linear Unit (PReLU) [2]		$f(x) = \begin{cases} \alpha x & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} \alpha & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$
Exponential Linear Unit (ELU) [3]		$f(x) = \begin{cases} \alpha(e^x - 1) & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} f(x) + \alpha & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$
SoftPlus		$f(x) = \log_e(1 + e^x)$	$f'(x) = \frac{1}{1 + e^{-x}}$

Source: <https://towardsdatascience.com/activation-functions-neural-networks-1cbd9f8d91d6>

# Derivatives of Activation Functions

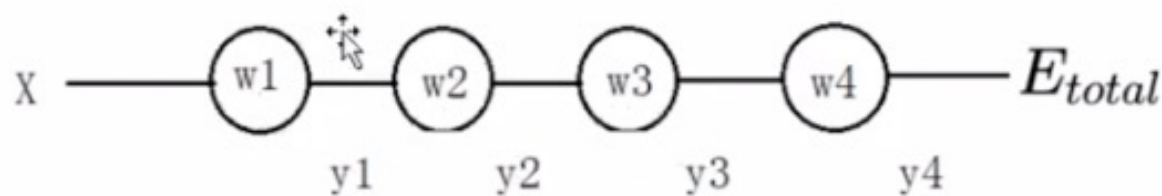


Source: <https://www.picswe.com/pics/relu-activation-7d.html>

# Gradient Vanishing and Exploding

$$w_1^+ = w_1 - \eta \frac{\partial E_{total}}{\partial w_1}$$

$$y_i = \sigma(z_i) = \sigma(w_i x_i + b_i)$$



$$\begin{aligned} \frac{\partial E_{total}}{\partial w_1} &= \frac{\partial E_{total}}{\partial y_4} \frac{\partial y_4}{\partial z_4} \frac{\partial z_4}{\partial x_4} \frac{\partial x_4}{\partial z_3} \frac{\partial z_3}{\partial x_3} \frac{\partial x_3}{\partial z_2} \frac{\partial z_2}{\partial x_2} \frac{\partial x_2}{\partial z_1} \frac{\partial z_1}{\partial w_1} \\ &= \frac{\partial E_{total}}{\partial y_4} \sigma'(z_4) w_4 \sigma'(z_3) w_3 \sigma'(z_2) w_2 \sigma'(z_1) x_1 \end{aligned}$$

# Data Normalization for Deep Learning training



# Batch Normalization

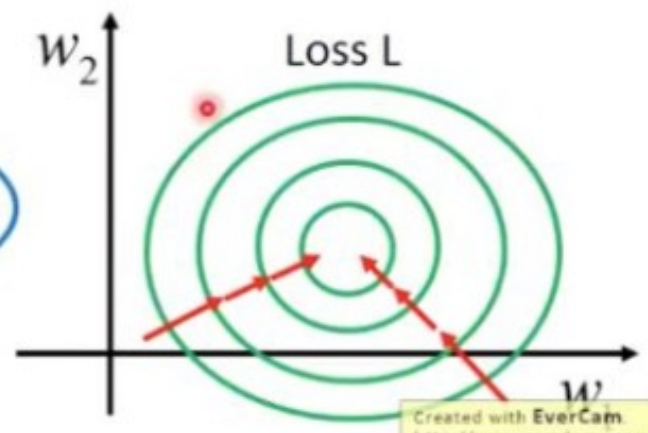
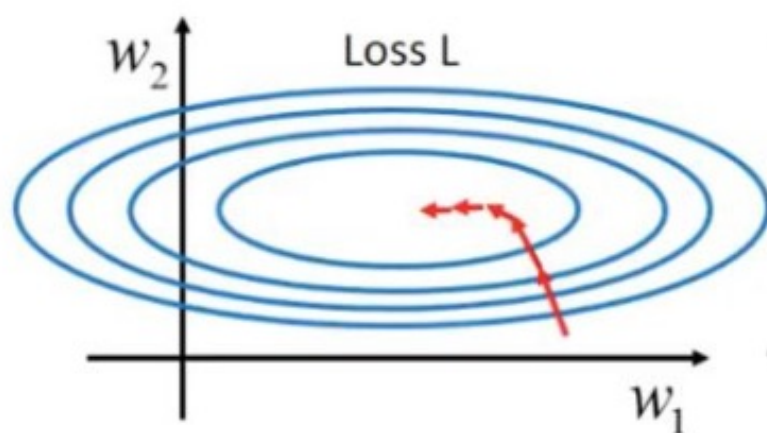
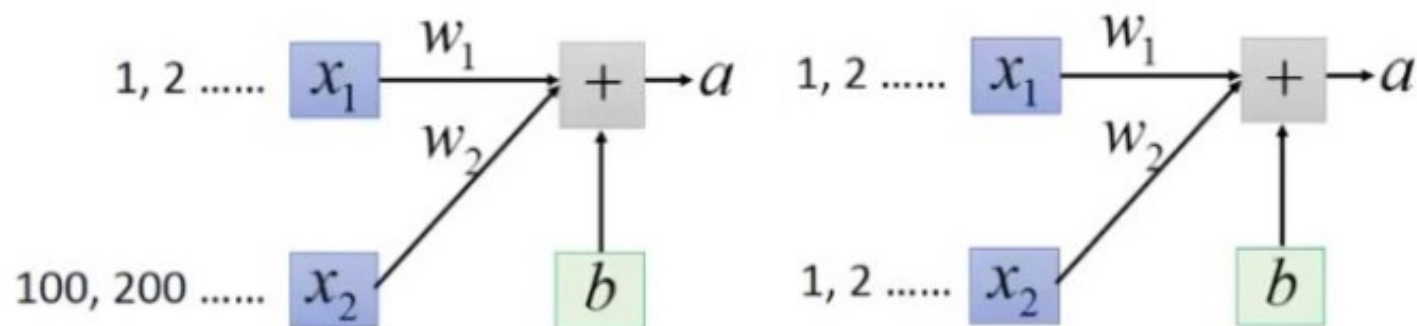
- Batch normalization speeds up the training by setting high learning rate
- Batch normalization prevents the gradients from getting too small or too large by normalizing data across each batch, as the name suggests. It also acts as a regularization method, similar to dropout.

Cited: <https://zhuanlan.zhihu.com/p/34480619> and Prof. Lee from TW Univ.

# Feature Scaling

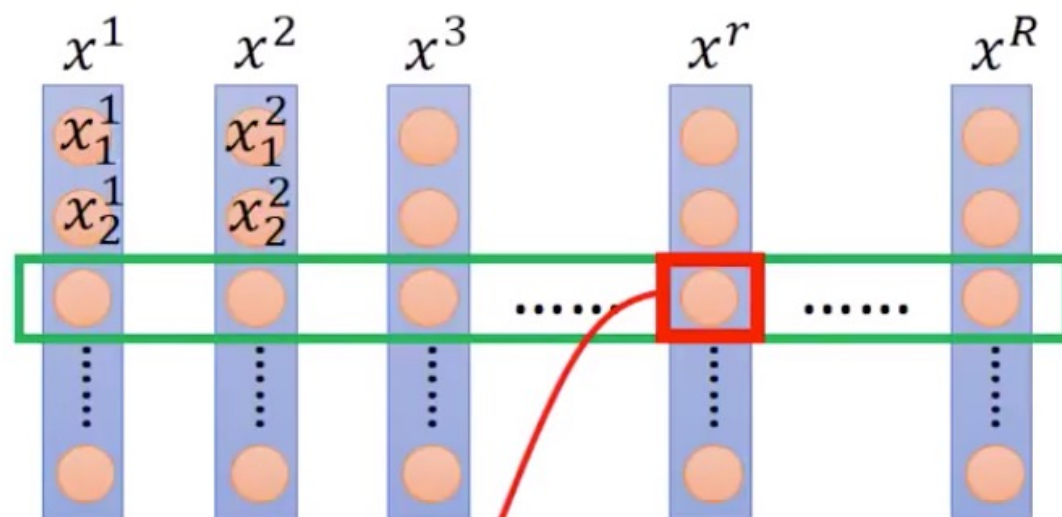
## Feature Scaling

Make different features have the same scaling



Created with EverCam  
<http://www.camdemy.com>

# Feature Scaling



For each dimension  $i$ :

mean:  $m_i$

standard deviation:  $\sigma_i$

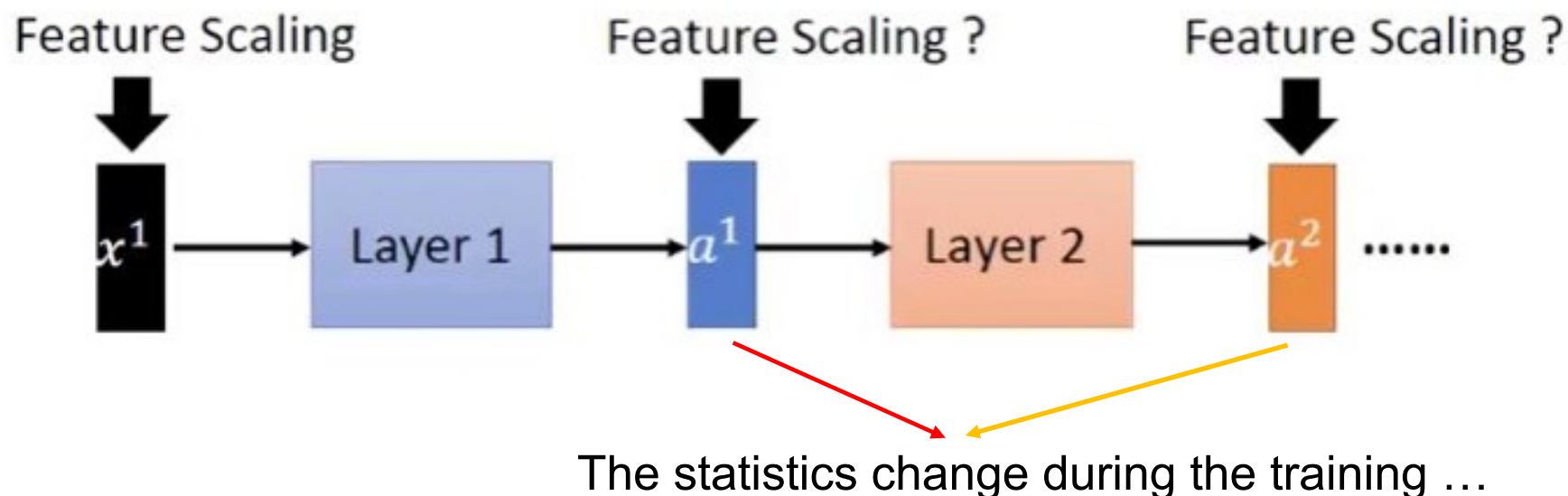
$$x_i^r \leftarrow \frac{x_i^r - m_i}{\sigma_i}$$

The means of all dimensions are 0, and the variances are all 1

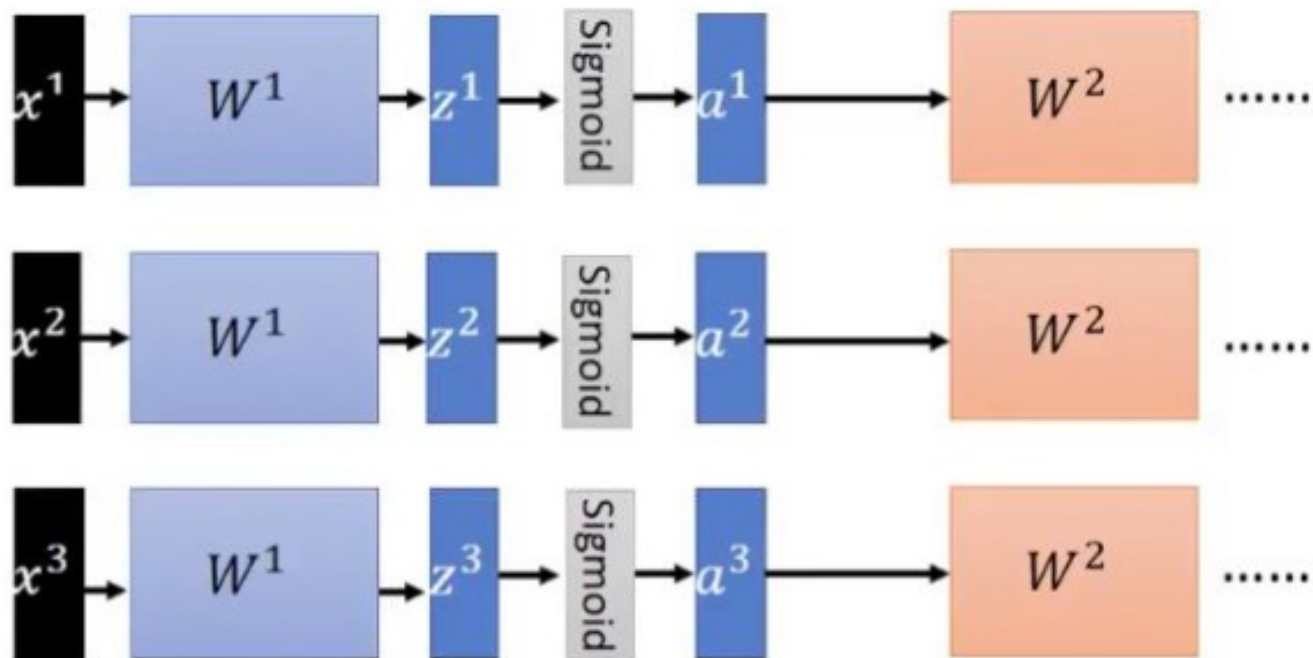
In general, gradient descent converges much faster with feature scaling than without it.

# Scaling at Hidden Layer

- Internal Covariate Shift



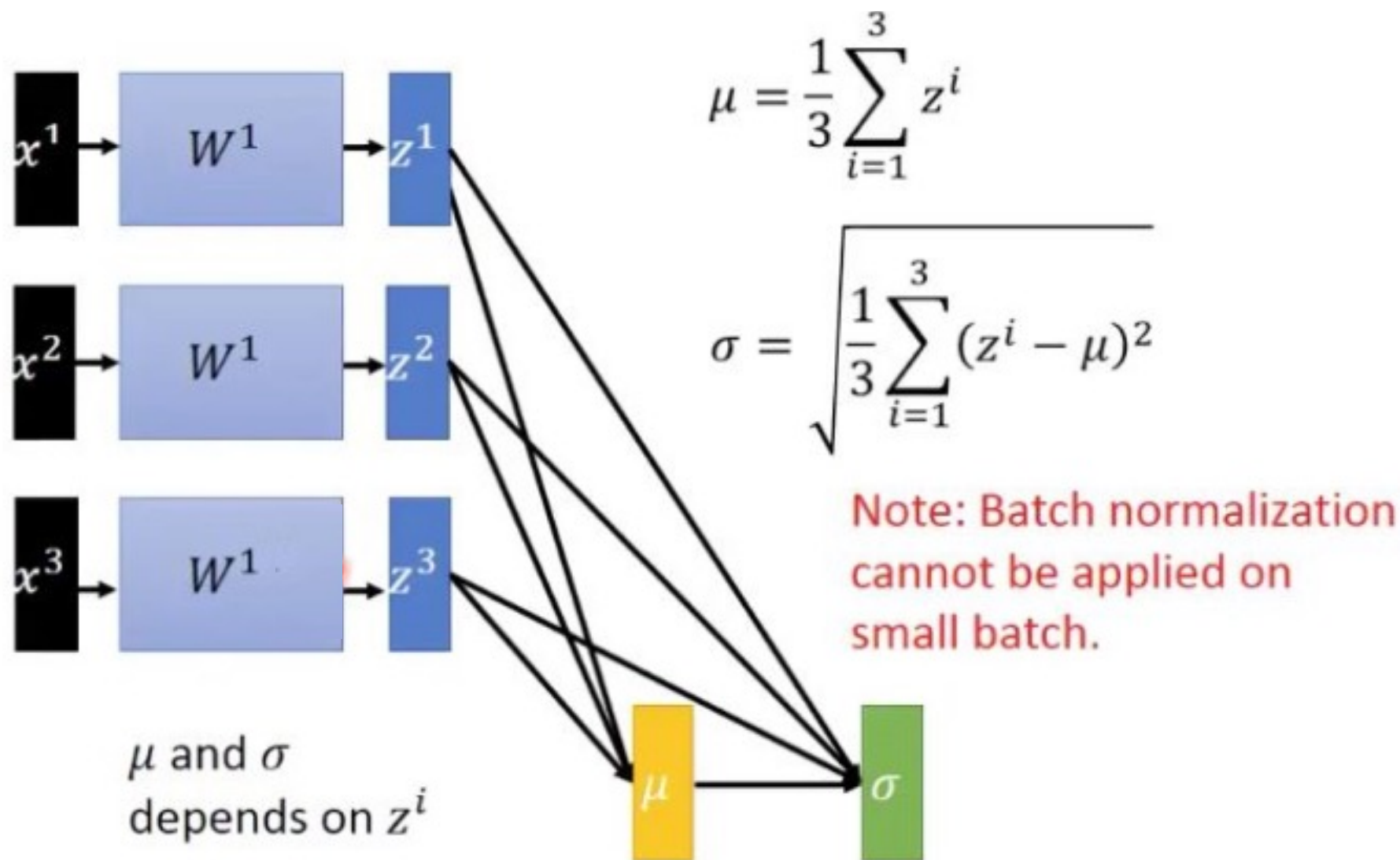
# Batch of Data



**Batch**

$$\begin{bmatrix} z^1 & z^2 & z^3 \end{bmatrix} = \begin{bmatrix} W^1 \end{bmatrix} \begin{bmatrix} x^1 & x^2 & x^3 \end{bmatrix}$$

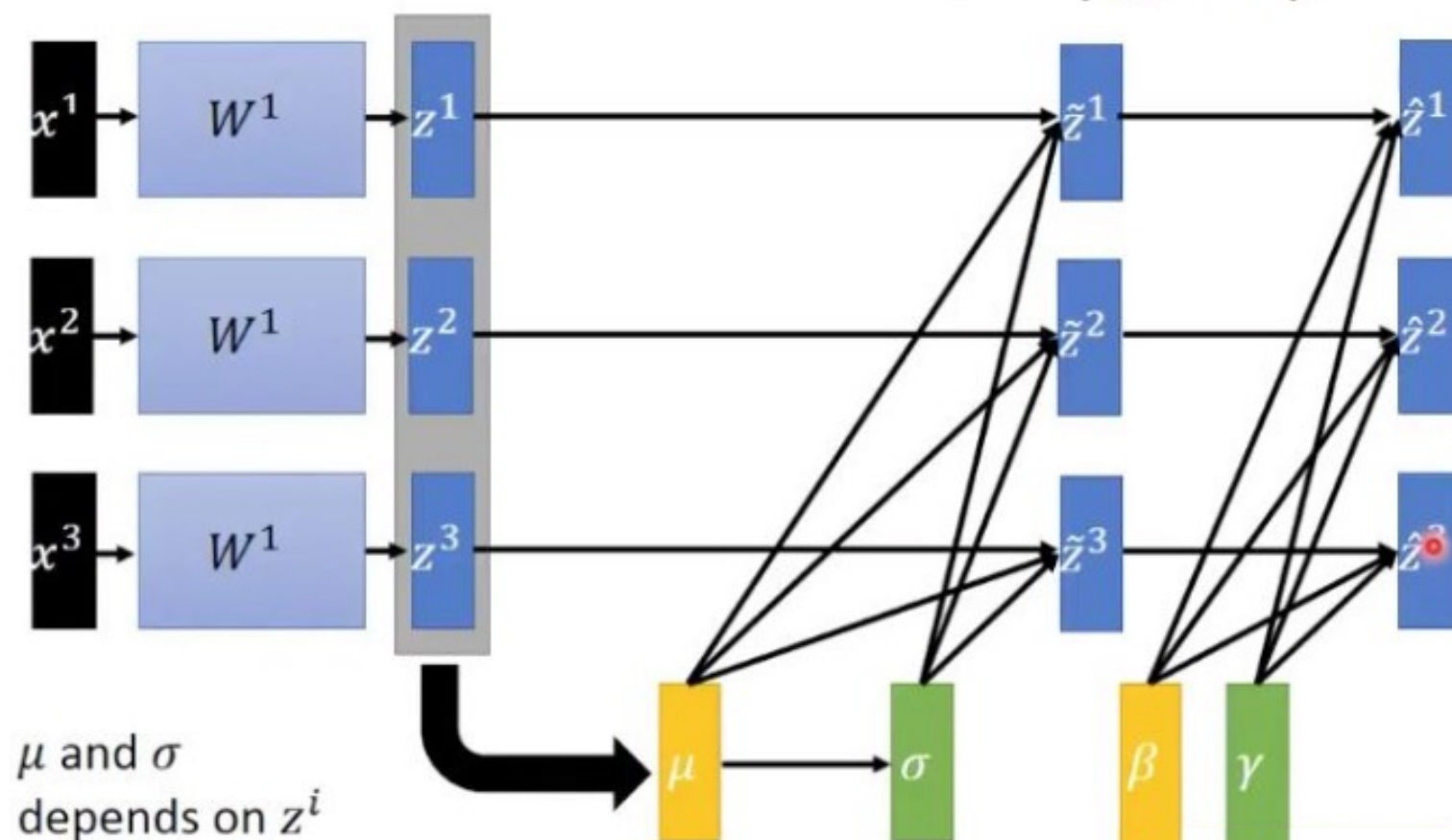
# Batch normalization



# Batch normalization

$$\tilde{z}^i = \frac{z^i - \mu}{\sigma}$$

$$\hat{z}^i = \gamma \odot \tilde{z}^i + \beta$$

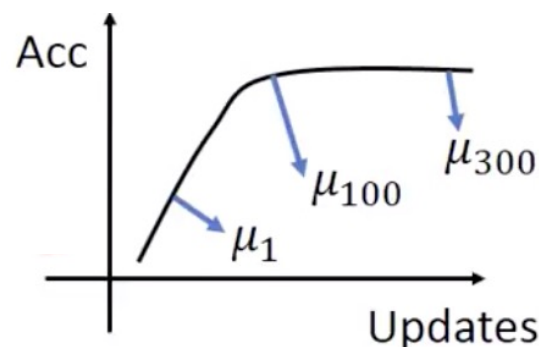
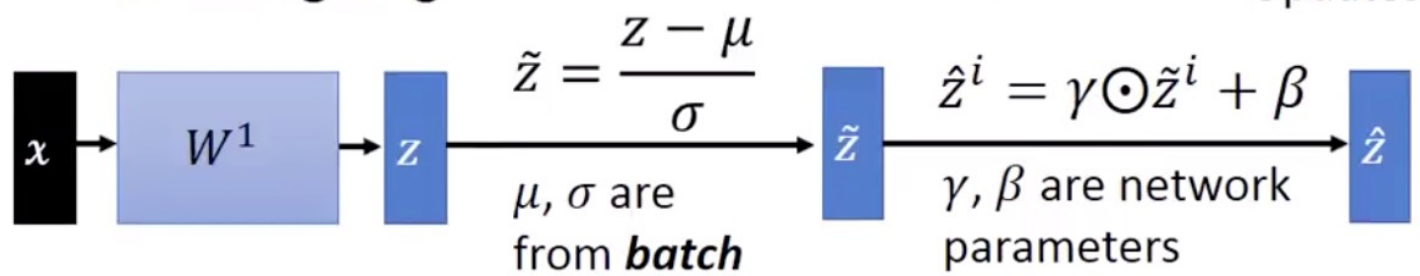




# BN at Testing

## Batch normalization

- At testing stage:



We do not have batch at testing stage.

Ideal solution:

Computing  $\mu$  and  $\sigma$  using the whole training dataset.

Practical solution:

Computing the moving average of  $\mu$  and  $\sigma$  of the batches during training.



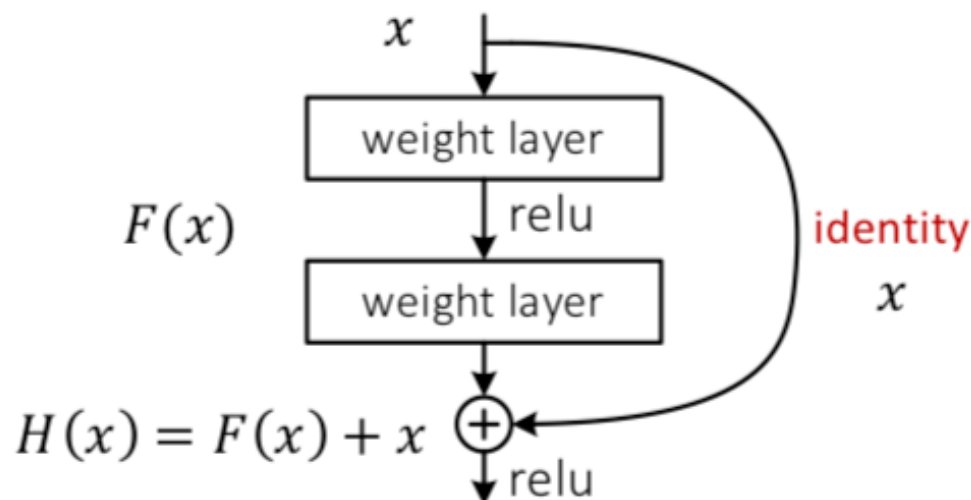
# BN: Benefits

- Reduce the covariate shift, speed up the training
- Reduce the vanishing/exploding gradients
- Less affected by weights initialization

# Residual Training

## Deep Residual Learning

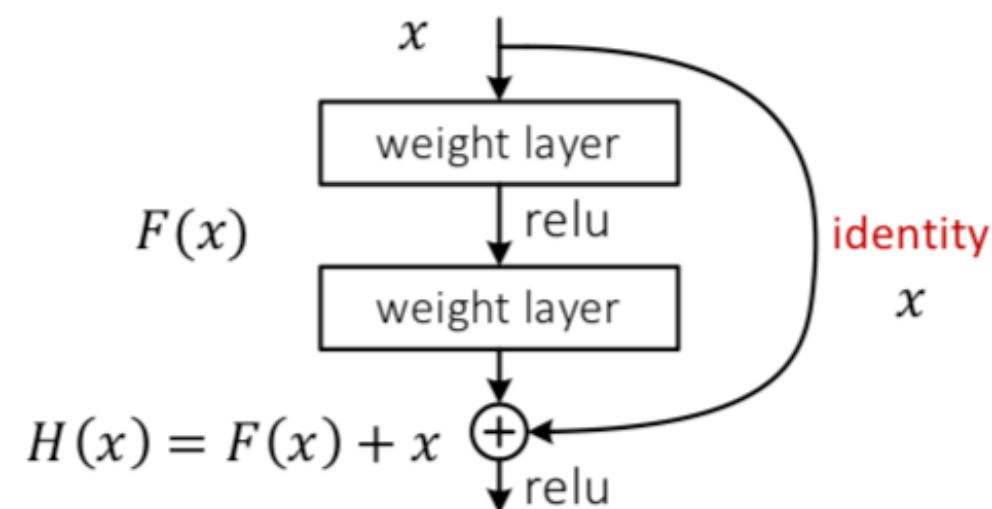
- **Residual** net



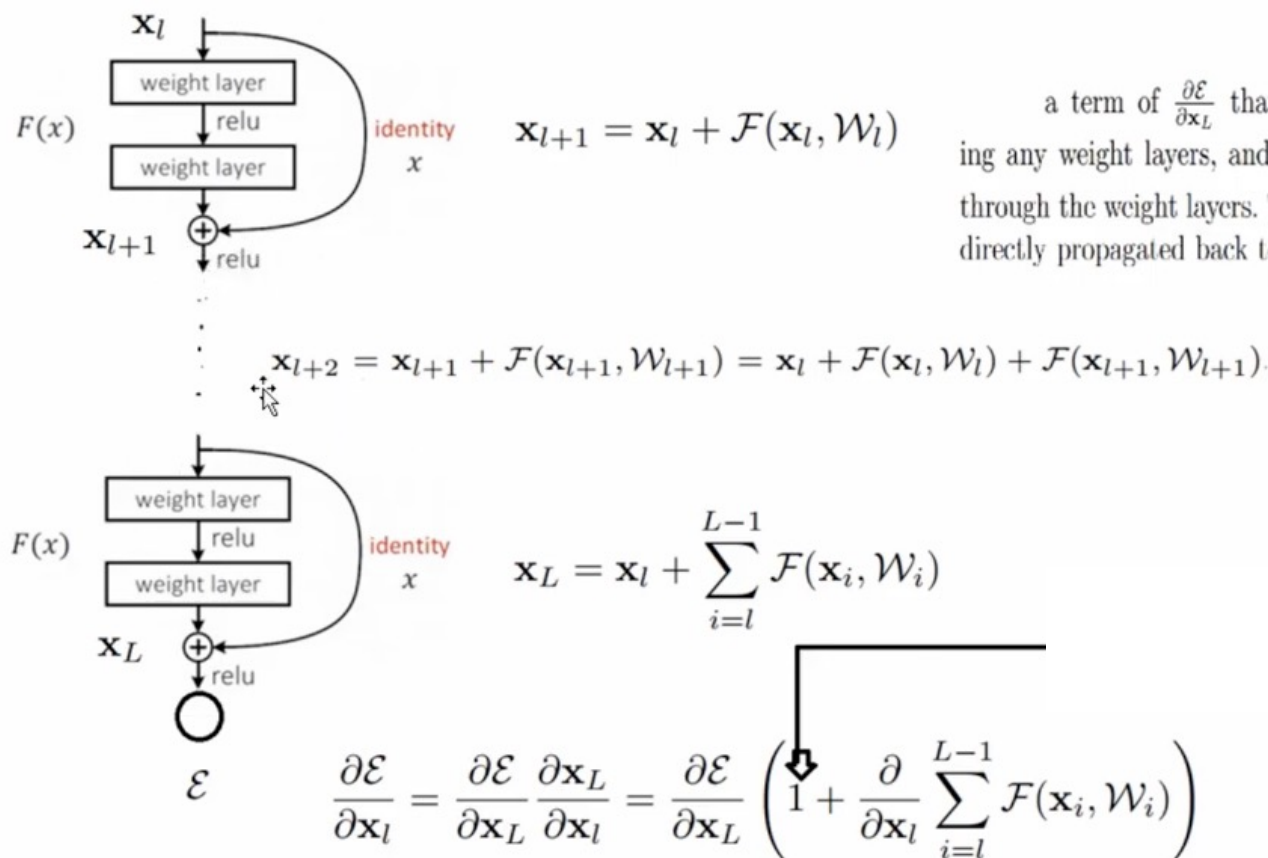
$H(x)$  is any desired mapping,  
~~hope the 2 weight layers fit  $H(x)$~~   
 hope the 2 weight layers fit  $F(x)$   
 let  $H(x) = F(x) + x$

# Deep Residual Learning

- $F(x)$  is a **residual** mapping w.r.t. **identity**

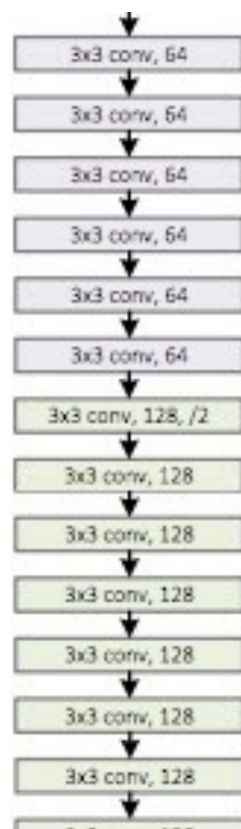


- If identity were optimal, easy to set weights as 0
- If optimal mapping is closer to identity, easier to find small fluctuations

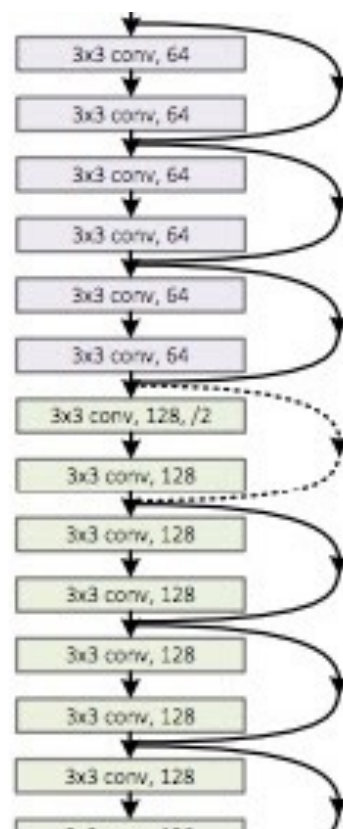


a term of  $\frac{\partial \mathcal{E}}{\partial \mathbf{x}_L}$  that propagates information directly without concerning any weight layers, and another term of  $\frac{\partial \mathcal{E}}{\partial \mathbf{x}_L} \left( \frac{\partial}{\partial \mathbf{x}_l} \sum_{i=l}^{L-1} \mathcal{F} \right)$  that propagates through the weight layers. The additive term of  $\frac{\partial \mathcal{E}}{\partial \mathbf{x}_L}$  ensures that information is directly propagated back to *any shallower unit*  $l$ .

plain net



ResNet



Network Design

# Benefits

- Speed up the training
- Deeper network with reduced gradient vanishing

# References

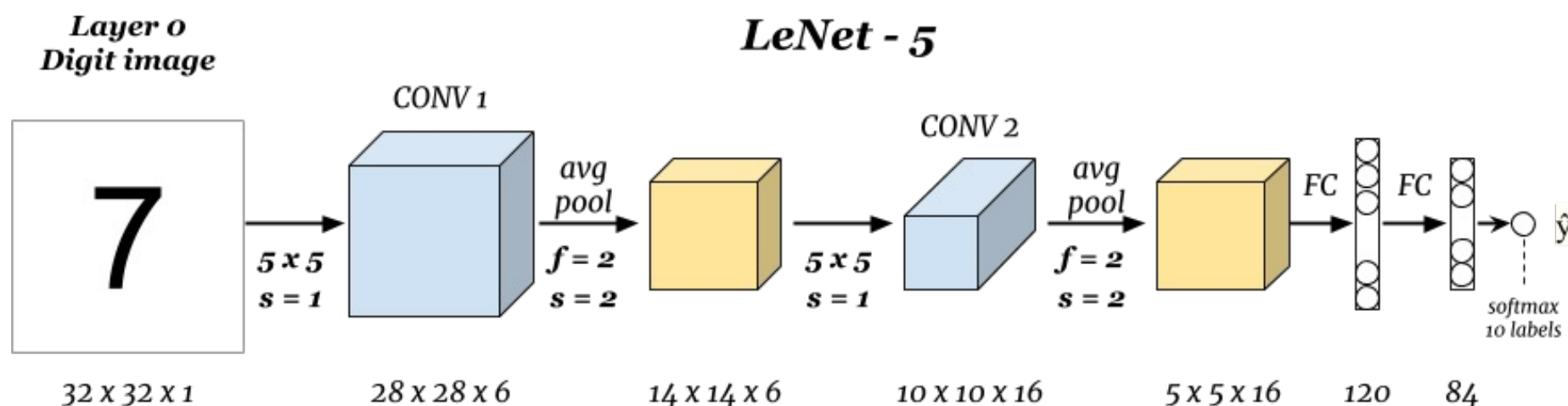
- CNN: ML Lecture 10: Convolutional Neural Network by Hongyi Li (youtube)
- Chapter 20 from Berkeley Artificial Intelligence: A Modern Approach (link: <http://aima.eecs.berkeley.edu/>)
- Module 2: Convolutional Neural Networks from course notes of Stanford CS231n (link: <http://cs231n.github.io/convolutional-networks/>)
- Full connection blog (link: <https://www.superdatascience.com/blogs/convolutional-neural-networks-cnn-step-4-full-connection>)
- <http://aima.eecs.berkeley.edu/>
- <http://cs231n.github.io/convolutional-networks/>

# Classic Networks

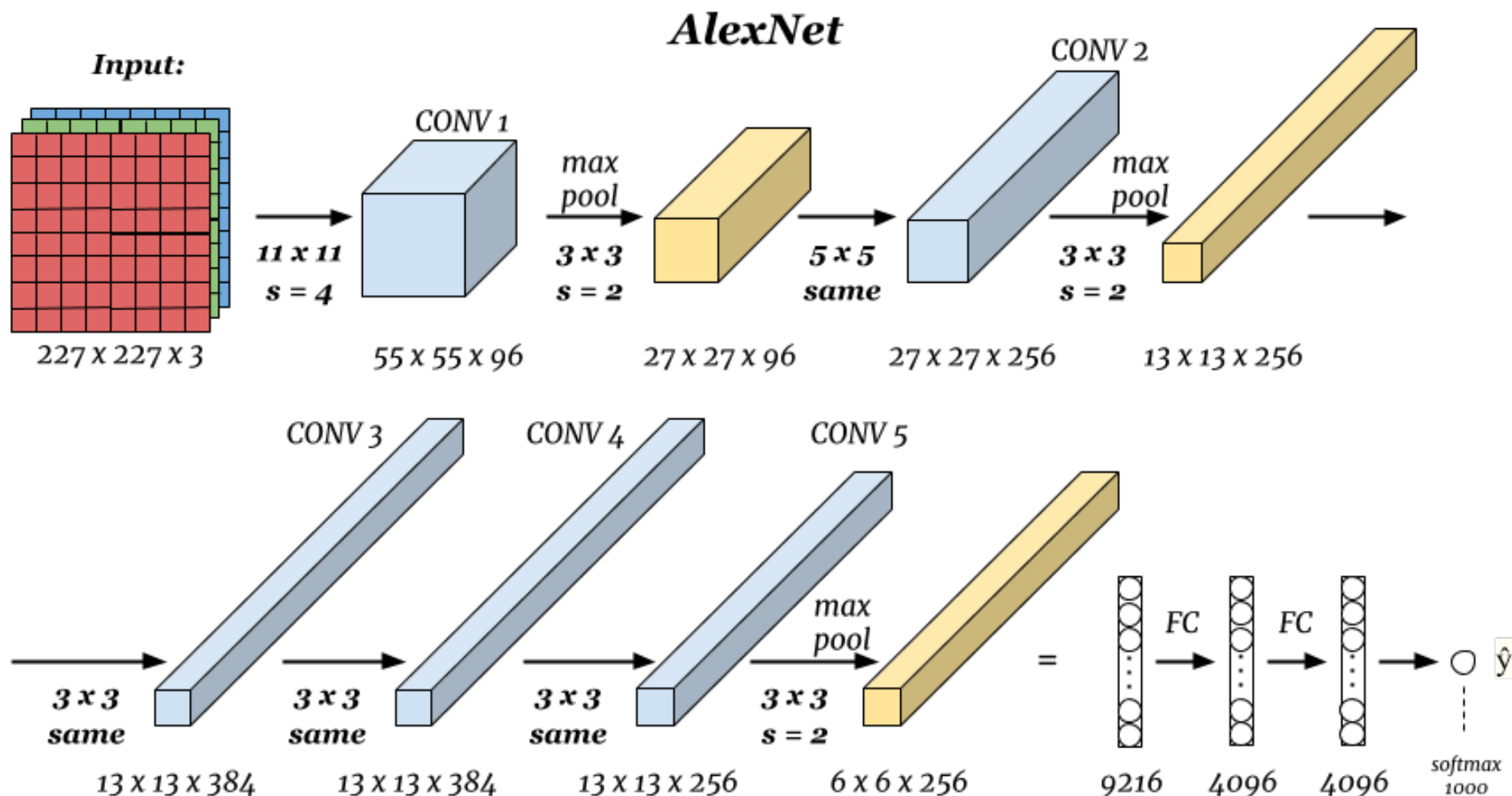
- LeNet
- AlexNet
- VGG Net
- ResNet
- GooLeNet



# LeNet

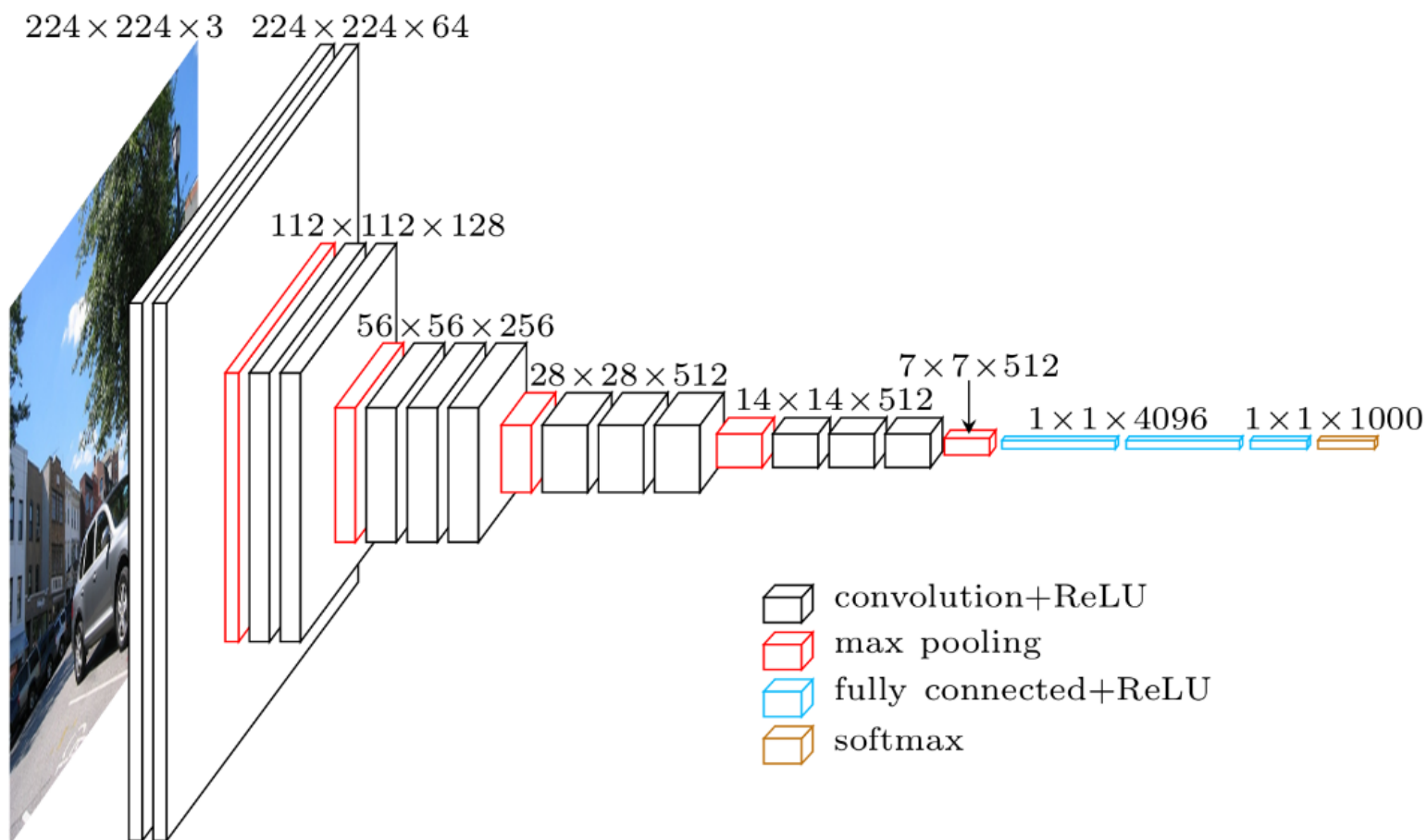


# AlexNet

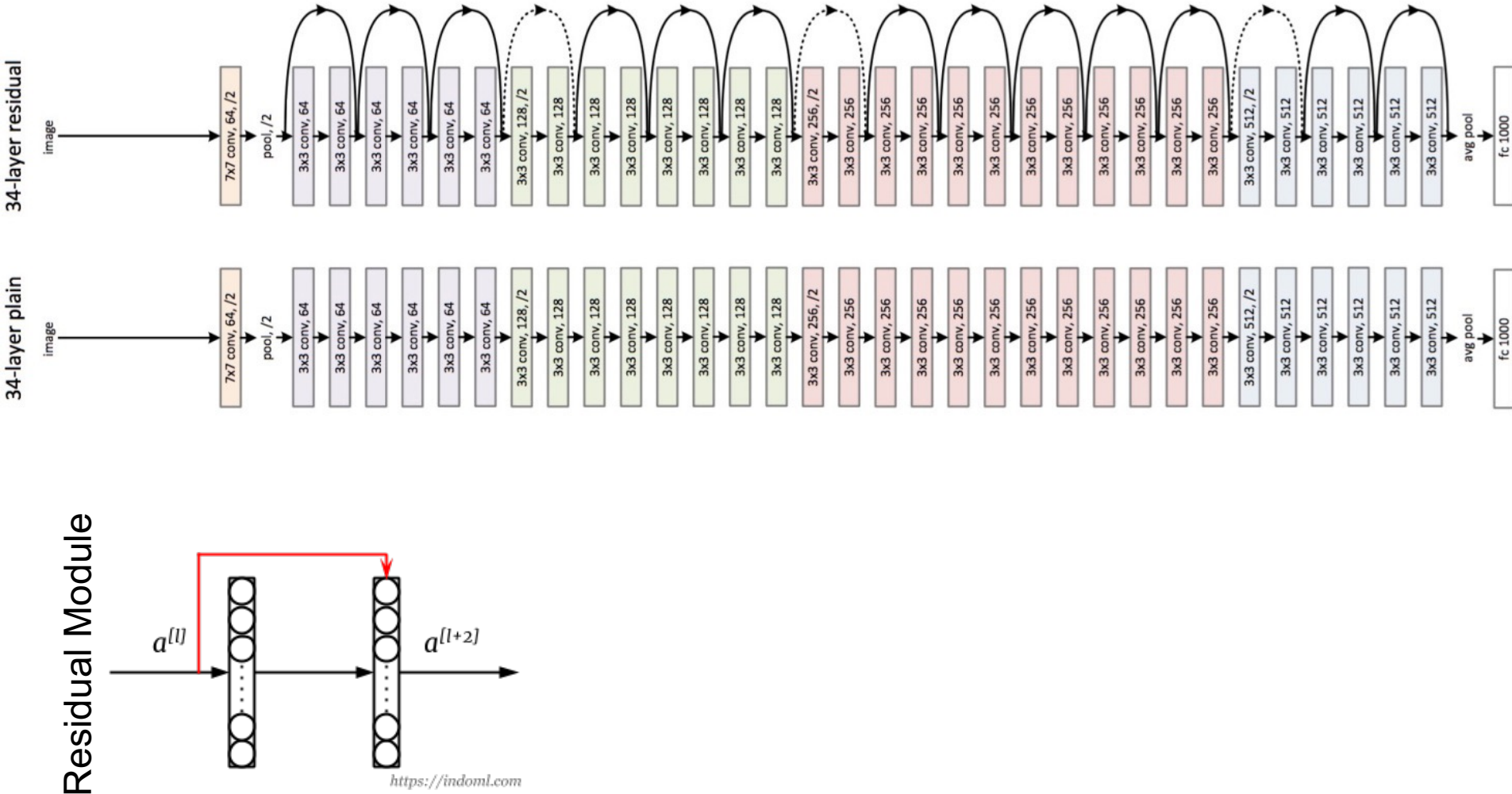


<https://indoml.com>

# VGG-16

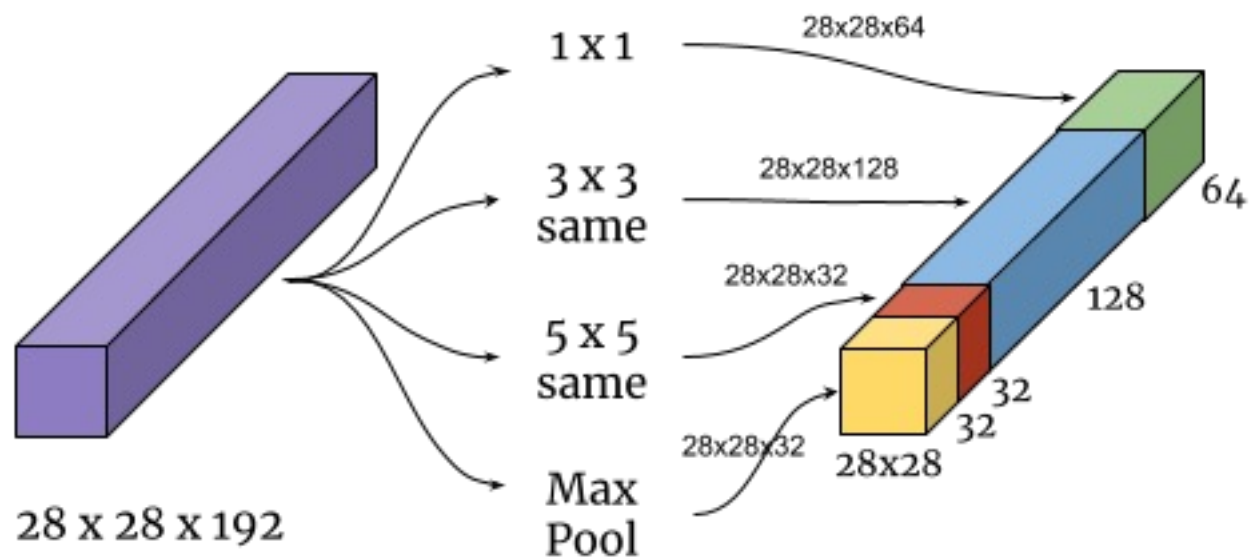


# ResNet

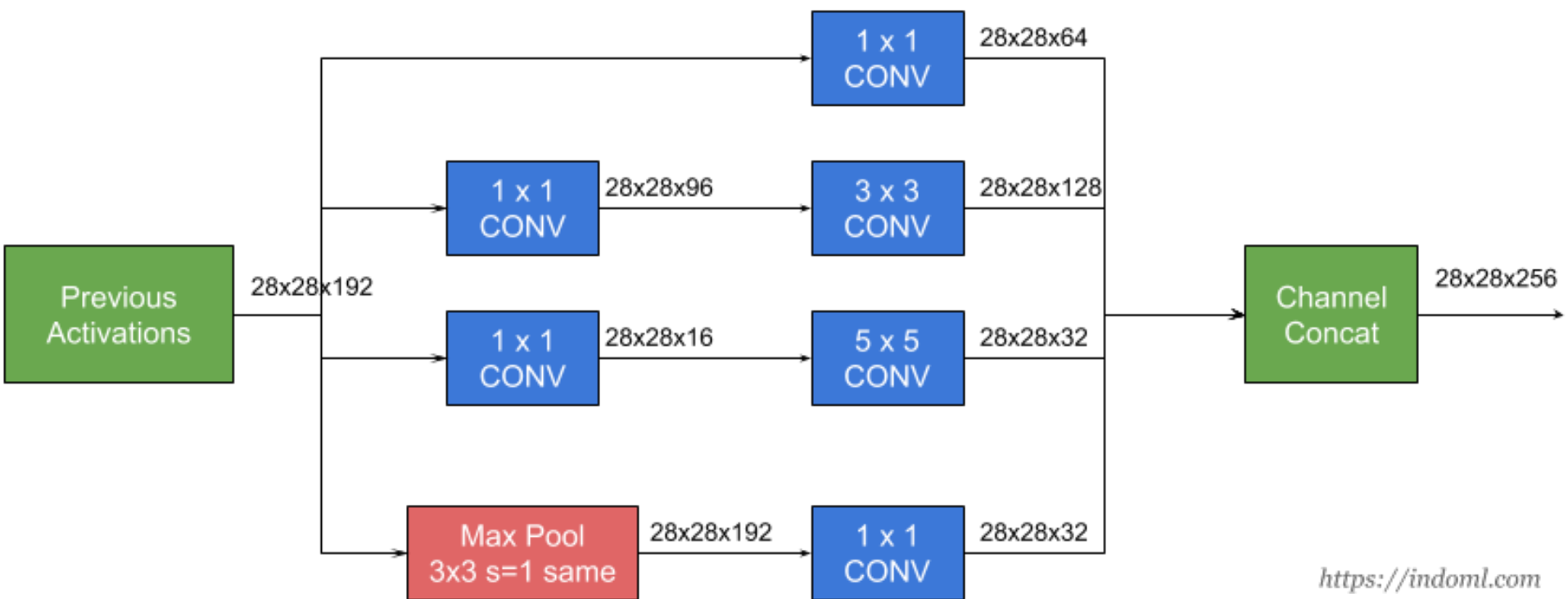


# GoogLeNet

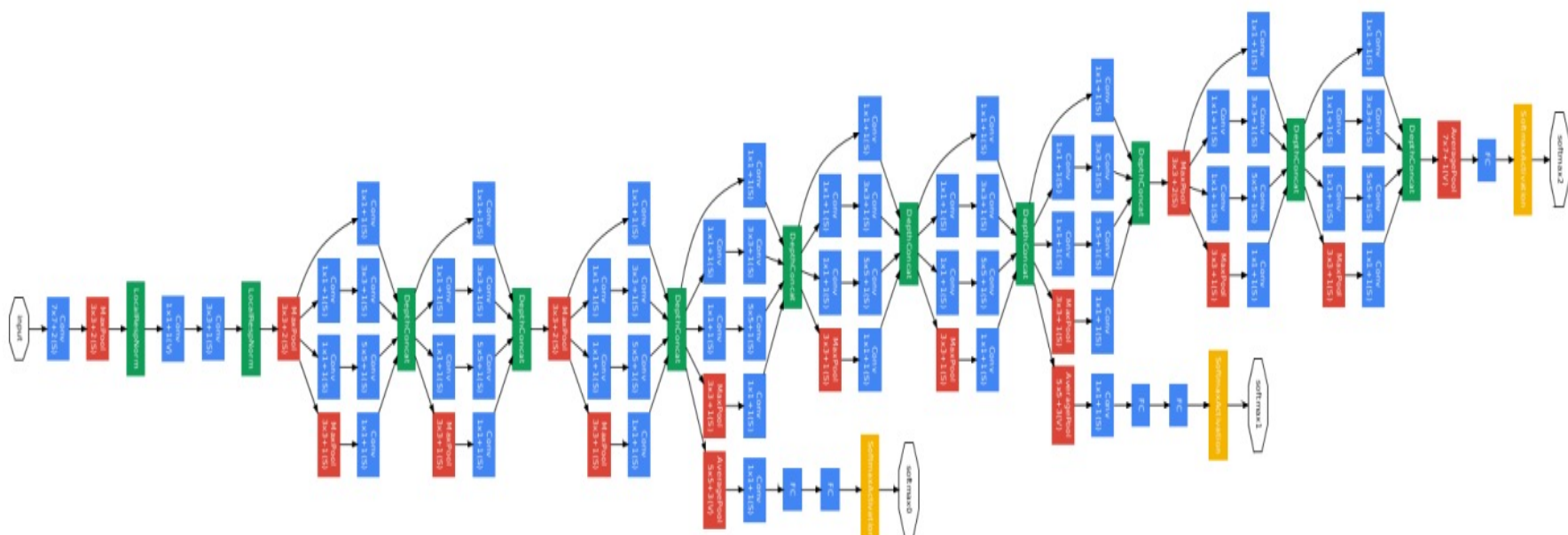
## Inception Module



<https://indoml.com>



# GoogLeNet:



# Summary of Networks

Year	CNN	Developed by	Place	Top-5 error rate	No. of parameters
1998	LeNet(8)	Yann LeCun et al			60 thousand
2012	AlexNet(7)	Alex Krizhevsky, Geoffrey Hinton, Ilya Sutskever	1st	15.3%	60 million
2014	GoogLeNet(19)	Google	1st	6.67%	4 million
2014	VGG Net(16)	Simonyan, Zisserman	2nd	7.3%	138 million
2015	<u>ResNet(152)</u>	Kaiming He	1st	3.6%	