

# Deep Learning Programming

## Lecture 3.2: Brief history of CNN

Sangryul Jeon

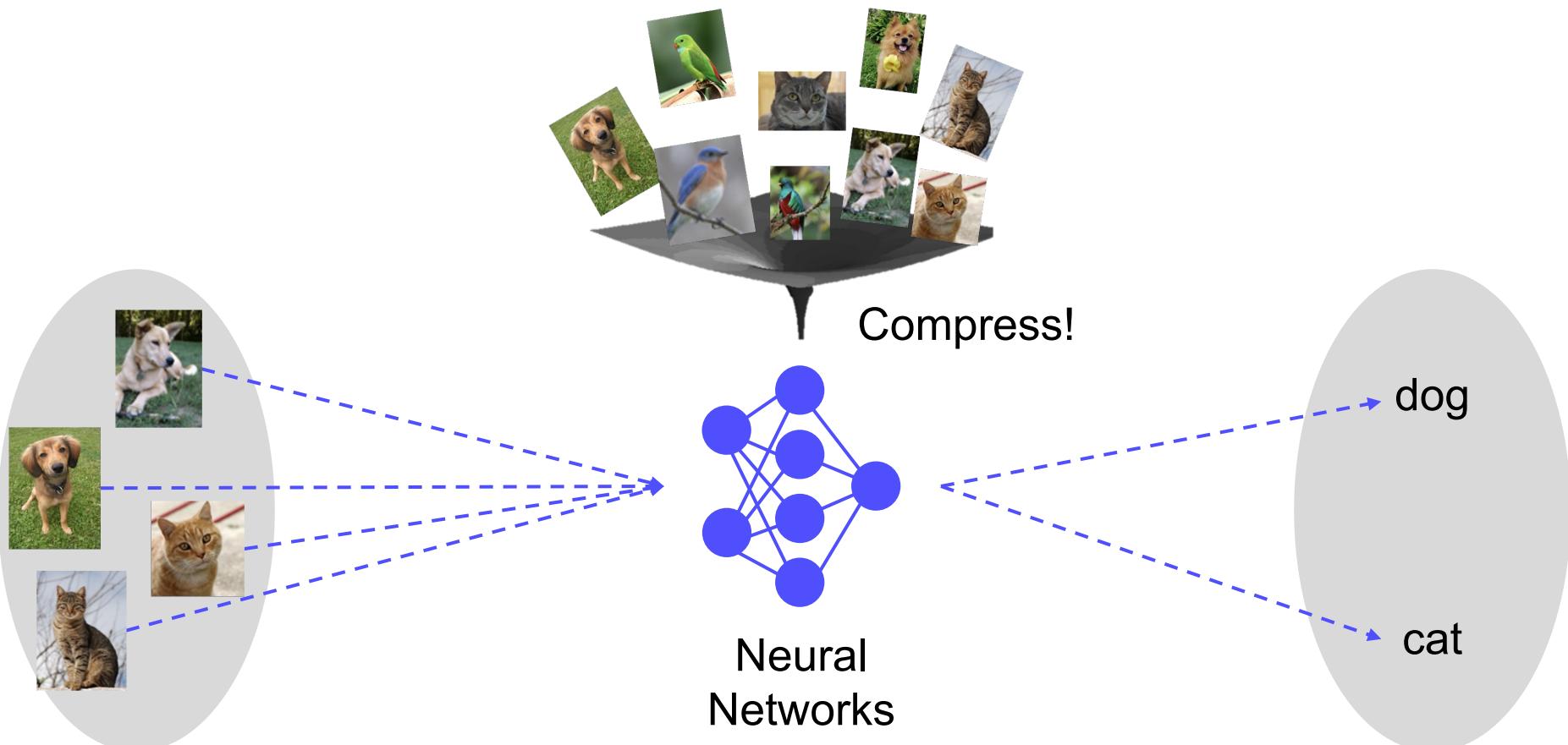
School of Computer Science and Engineering

[srjeonn@pusan.ac.kr](mailto:srjeonn@pusan.ac.kr)

# 1.1 Convolutional Neural Networks (CNN)

CNN architectures

Compress all the data we have into the neural network

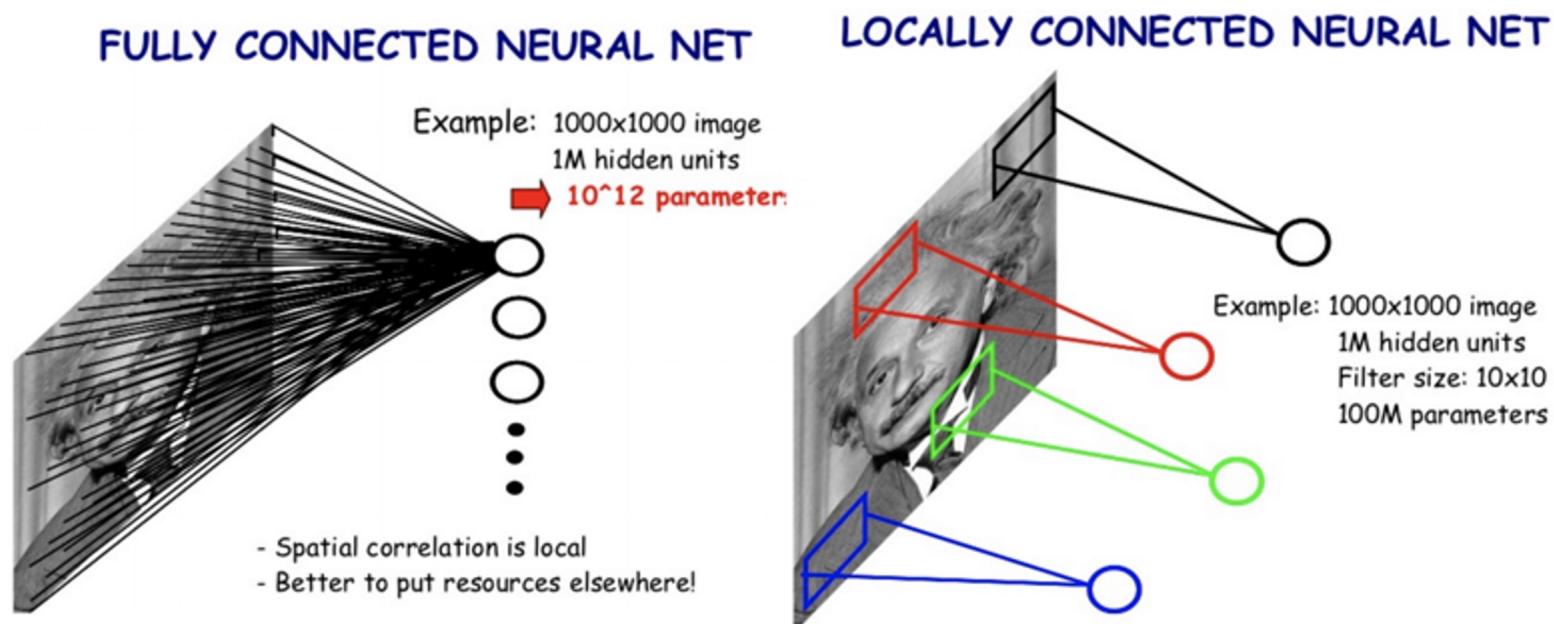


# 1.1 Convolutional Neural Networks (CNN)

CNN architectures

Convolution neural networks are fully locally connected neural networks

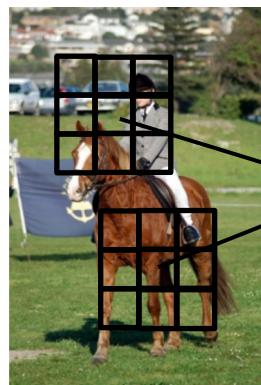
- Local feature learning
- Parameter sharing



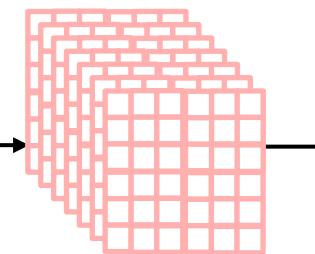
# 1.1 Convolutional Neural Networks (CNN)

CNN architectures

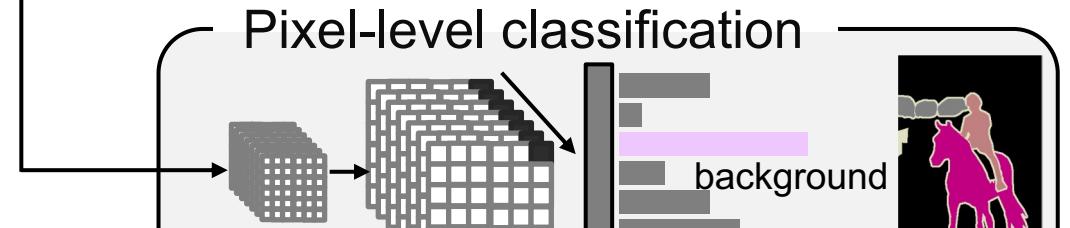
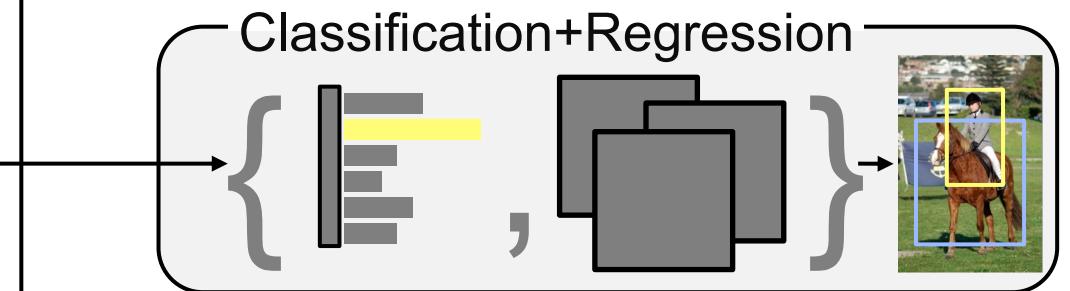
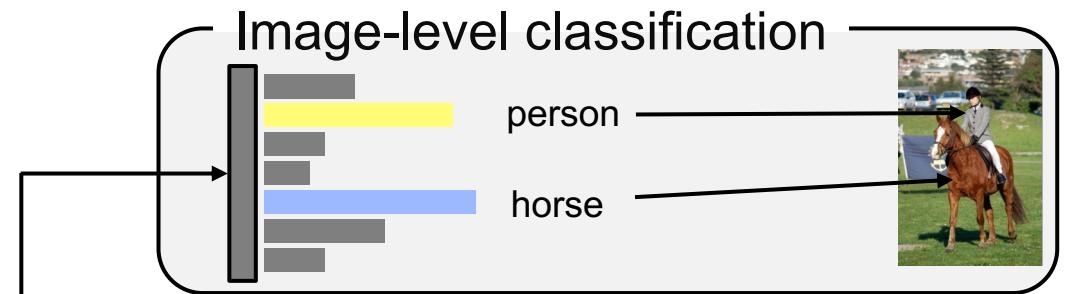
- CNN is used as a **backbone** of many CV tasks



CNN



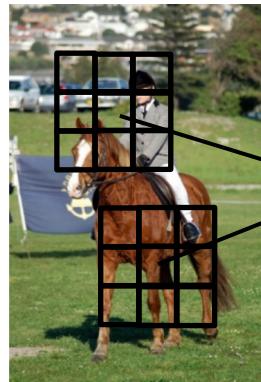
Feature map



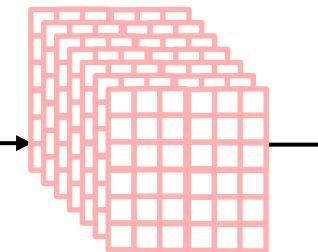
# 1.1 Convolutional Neural Networks (CNN)

CNN architectures

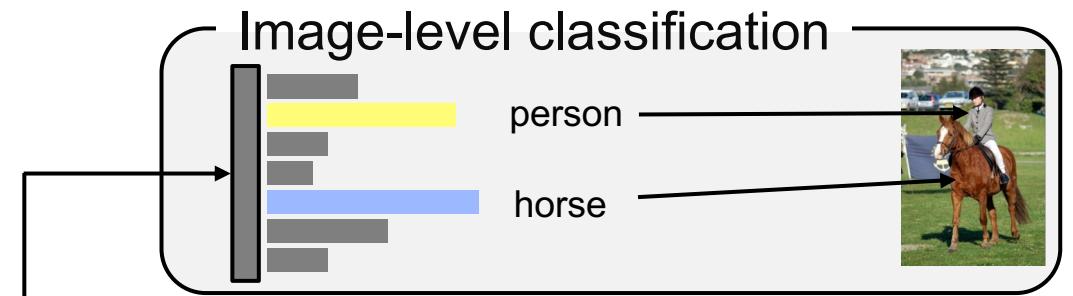
- CNN is used as a **backbone** of many CV tasks
- In this lecture, we will focus on **image-level classification**



CNN

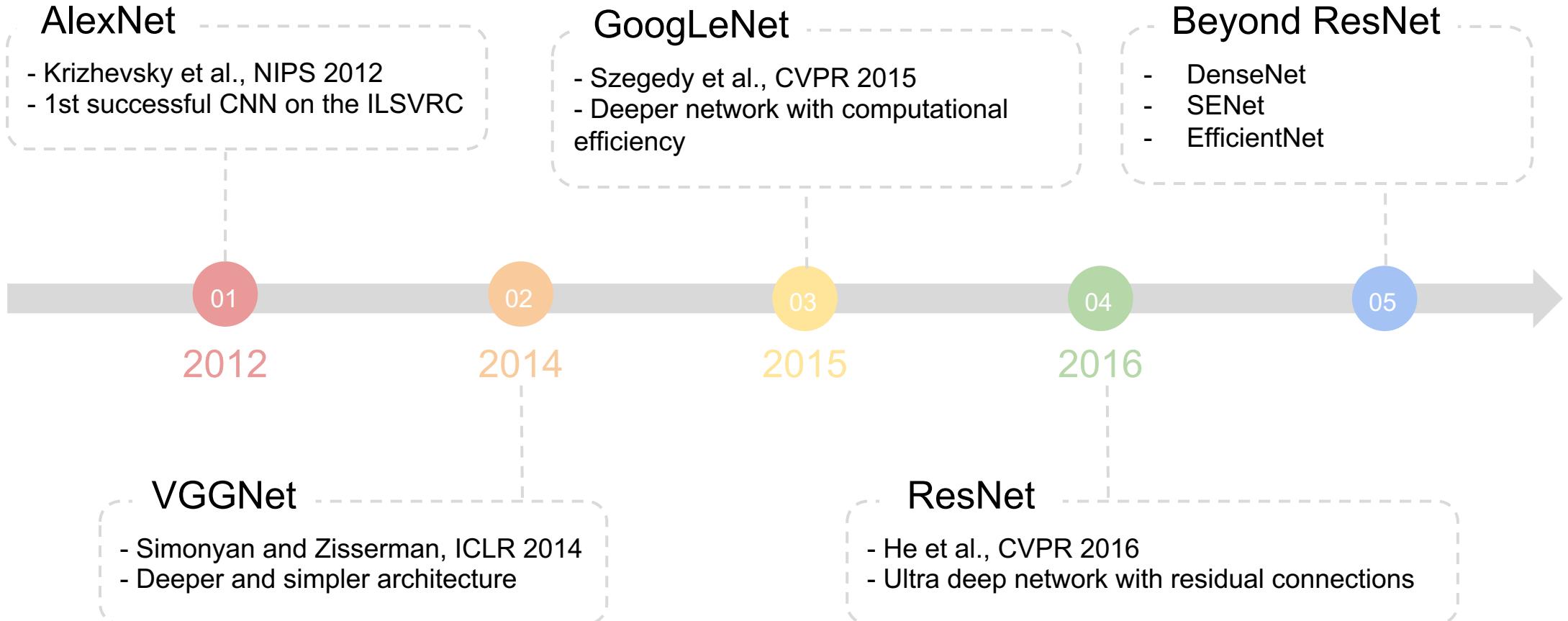


Feature map



## 1.2 Brief history

CNN architectures

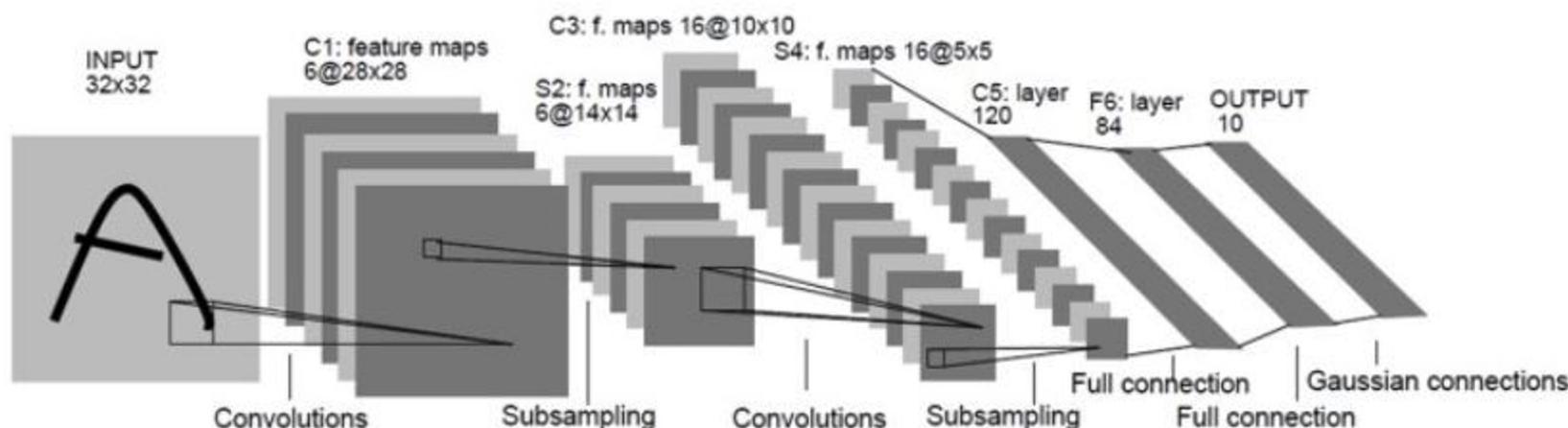


## 1.2 Brief history

### CNN architectures

#### LeNet-5

- A very simple CNN architecture introduced by Yann LeCun in 1998
  - Overall architecture: Conv - Pool - Conv - Pool - FC - FC
  - Convolution: 5x5 filters with stride 1
  - Pooling: 2x2 max pooling with stride 2



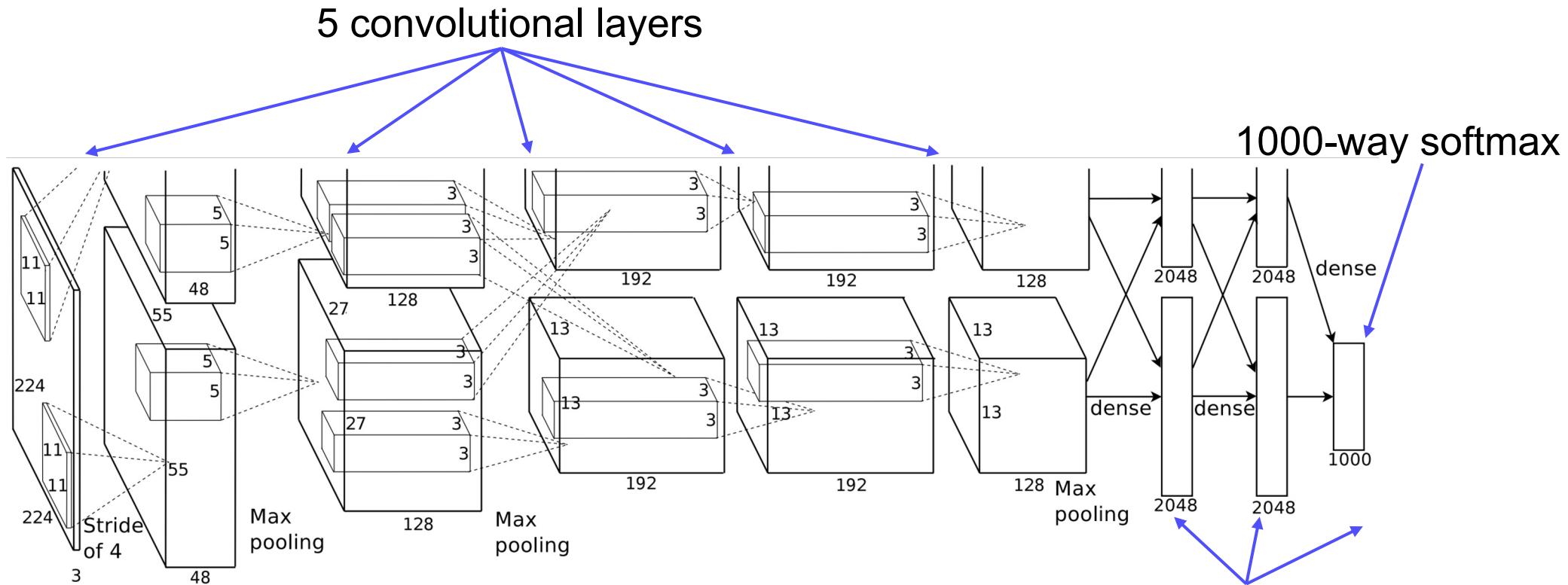
## 1.2 Brief history

CNN architectures

### AlexNet

- Conv - Pool - LRN - Conv - Pool - LRN - Conv - Conv - Pool - FC - FC - FC

\* LRN = Local Response Normalization

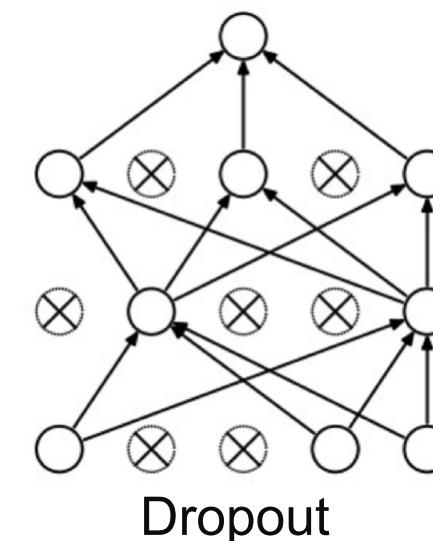
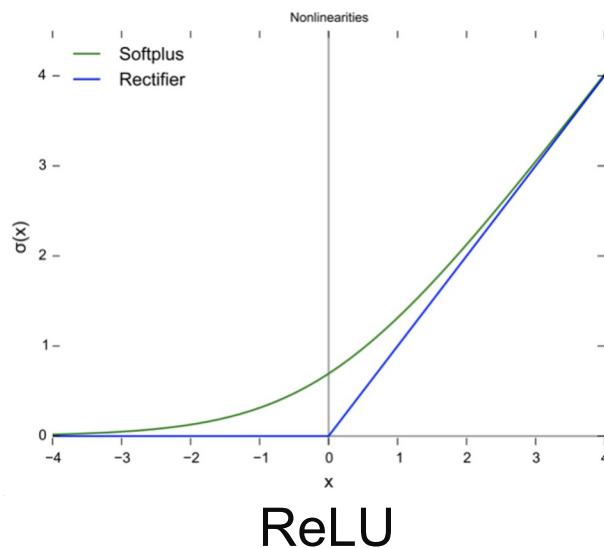


## 1.2 Brief history

### CNN architectures

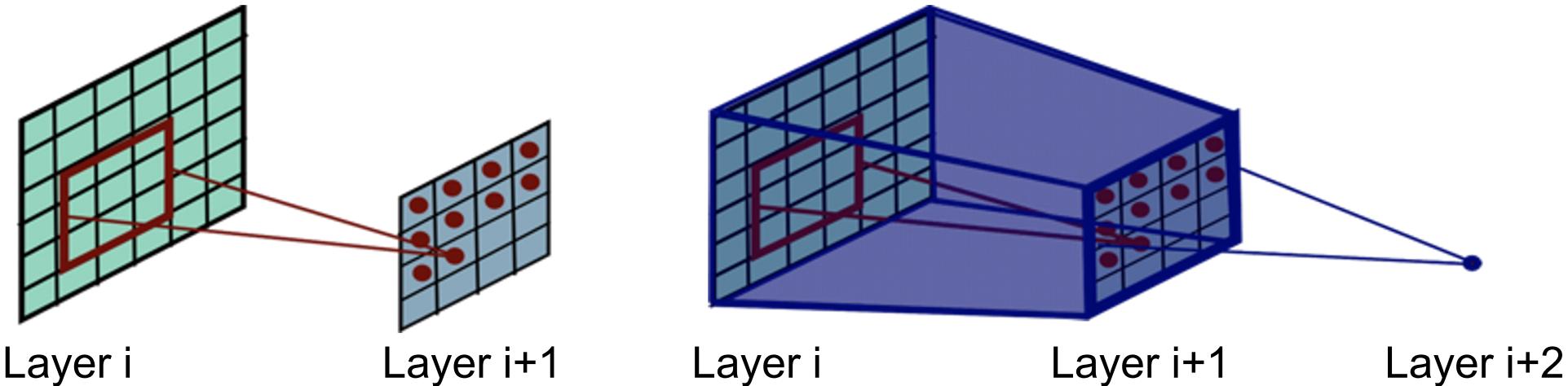
AlexNet - Similar with LeNet-5, but

- Bigger (7 hidden layers, 605k neurons, 60 million parameters)
- Trained with ImageNet (large amount of data, 1.2 millions)
- Using better activation function (**ReLU**) and regularization technique (**dropout**)



### Receptive field in CNN

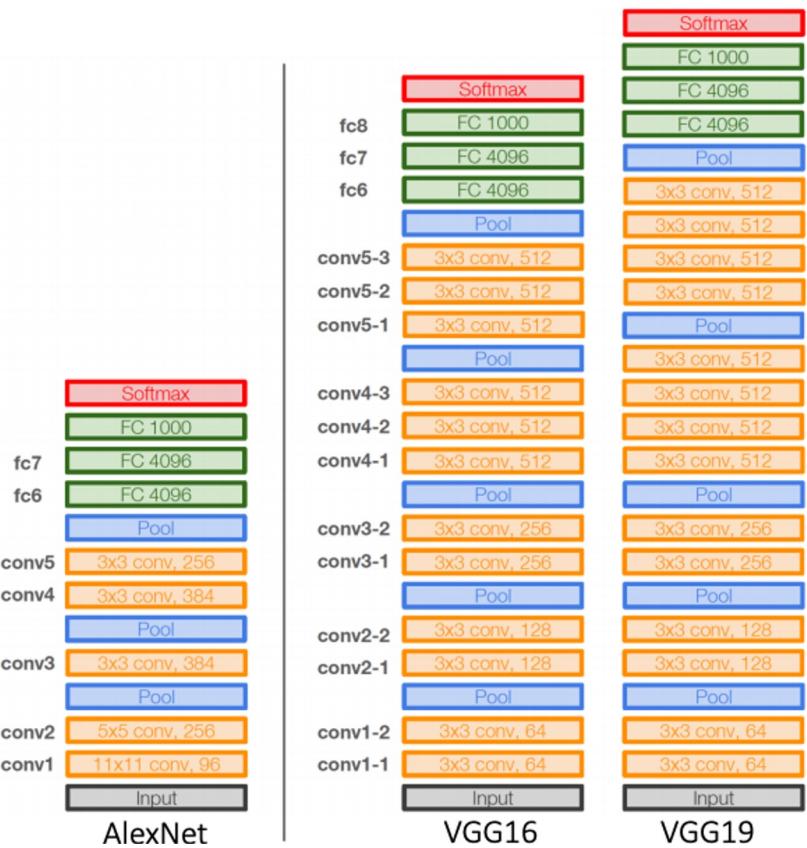
- The region in the input space that a particular CNN feature is looking at
- Suppose  $K \times K$  conv. filters with stride 1, and a pooling layer of size  $P \times P$ ,
  - then a value of each unit in the pooling layer depends on an input patch of size :  $(P+K-1) \times (P+K-1)$



## 1.2 Brief history

CNN architectures

### VGGNet

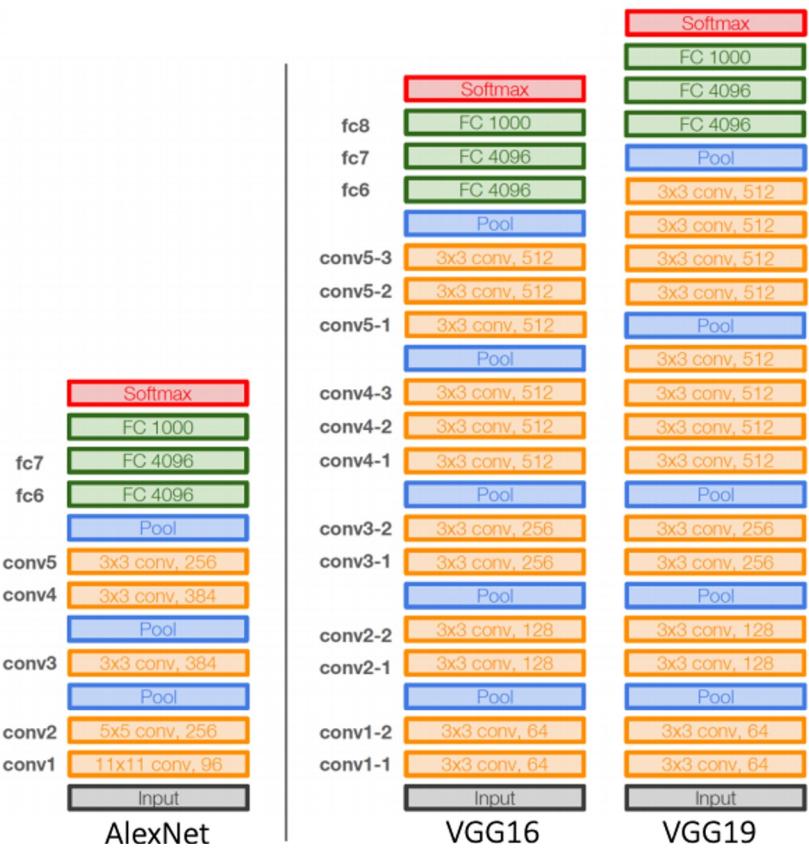


- Deeper architecture
  - 16 and 19 layers

## 1.2 Brief history

CNN architectures

### VGGNet

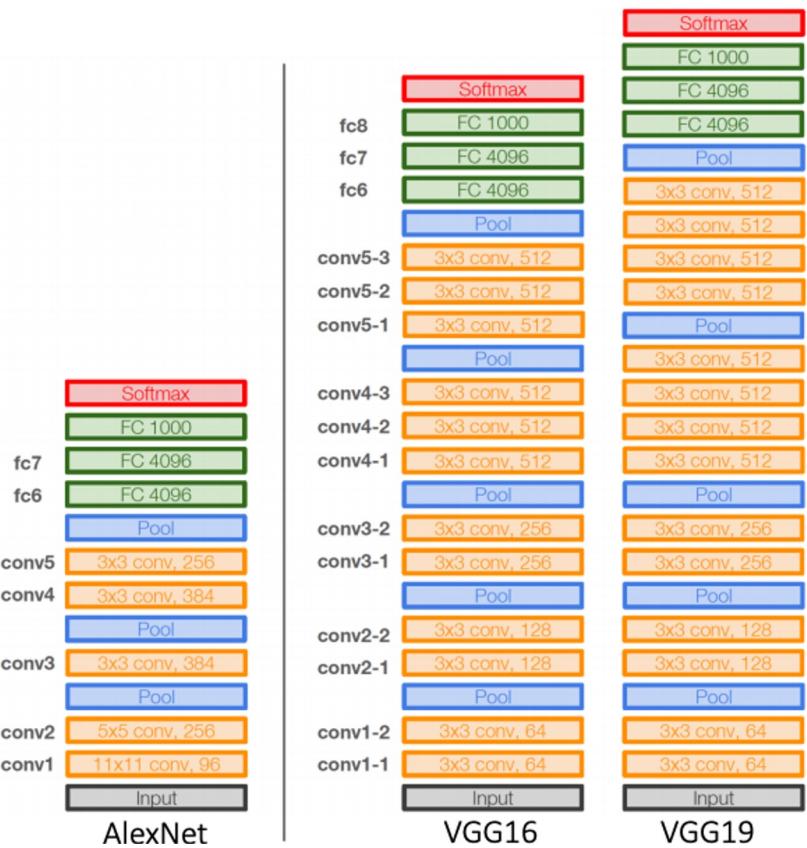


- Simpler architecture
  - No local response normalization
  - Only 3x3 conv filters blocks, 2x2 max pooling

## 1.2 Brief history

CNN architectures

### VGGNet

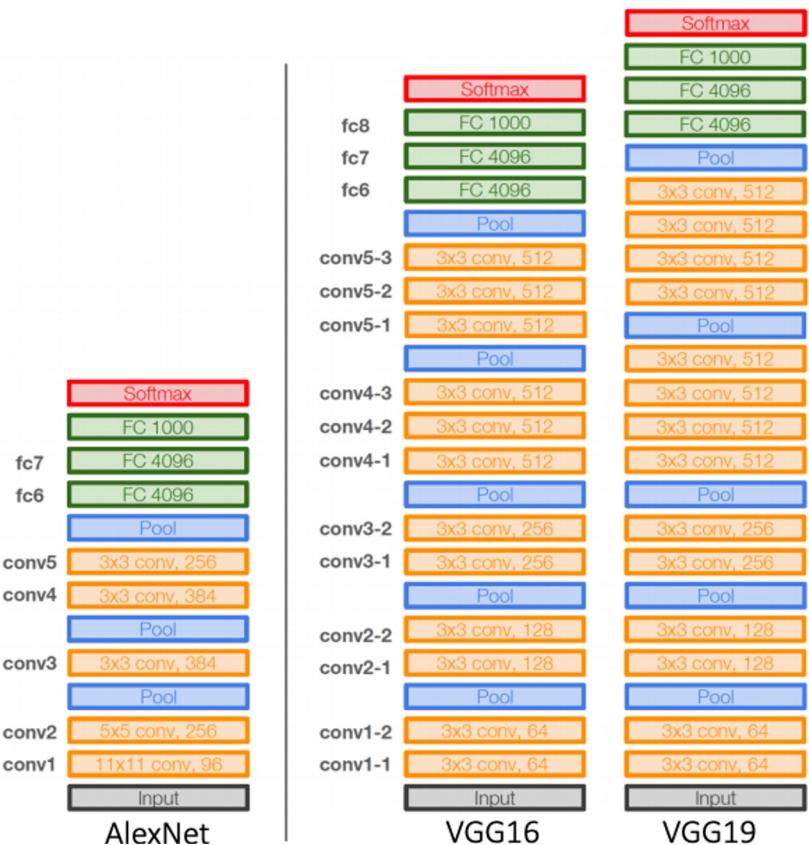


- Better performance
  - Significant performance improvement over AlexNet (2<sup>nd</sup> in ILSVRC14)

## 1.2 Brief history

CNN architectures

### VGGNet

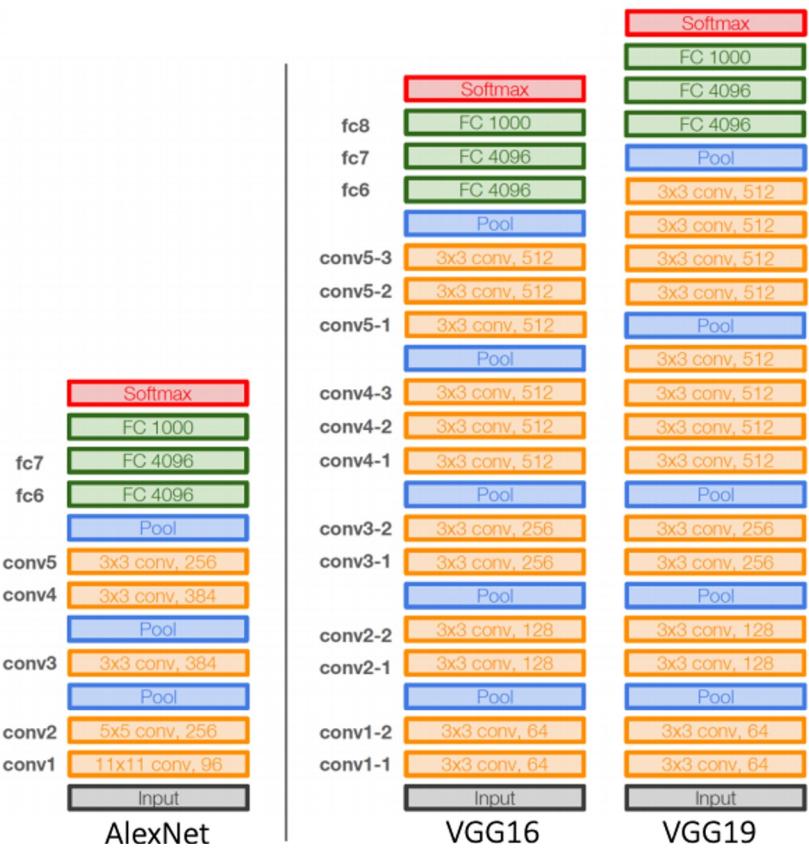


- Better generalization
  - Final features generalizing well to other tasks even without fine-tuning

## 1.2 Brief history

CNN architectures

### VGGNet



- Deeper architecture
- Simpler architecture
- Better performance
- Better generalization

### VGGNet - Overall architecture



- Input
  - 224x224 RGB images (same with AlexNet)
  - Subtracting mean RGB values of training images

### VGGNet - Overall architecture

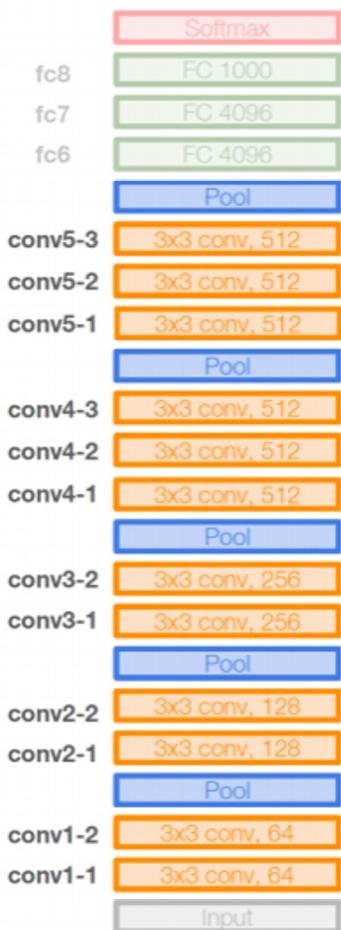


- Key design choices
  - 3x3 convolution filters with stride 1
  - 2x2 max pooling operations

## 1.2 Brief history

### CNN architectures

#### VGGNet - Overall architecture



- Key design choices
  - 3x3 convolution filters with stride 1
  - 2x2 max pooling operations



- Using many 3x3 conv layers instead of a small number of larger conv filters
- Keeping receptive field sizes large enough
  - Deeper with more non-linearities
  - Fewer parameters

### VGGNet - Overall architecture



- 3 fully-connected (FC) layers
- Other details
  - ReLU for non-linearity
  - No local response normalization

# 1.3 Going deeper with convolutions

CNN architectures

The neural network is getting deeper and wider

- Deeper networks learn more powerful features, because of
  - Larger receptive fields
  - More capacity and non-linearity



AlexNet (8 layers)

VGGNet (19 layers)

# 1.3 Going deeper with convolutions

CNN architectures

The neural network is getting deeper and wider

- Deeper networks learn more powerful features, because of
  - Larger receptive fields
  - More capacity and non-linearity
- But, getting deeper and deeper always works better?



AlexNet (8 layers)

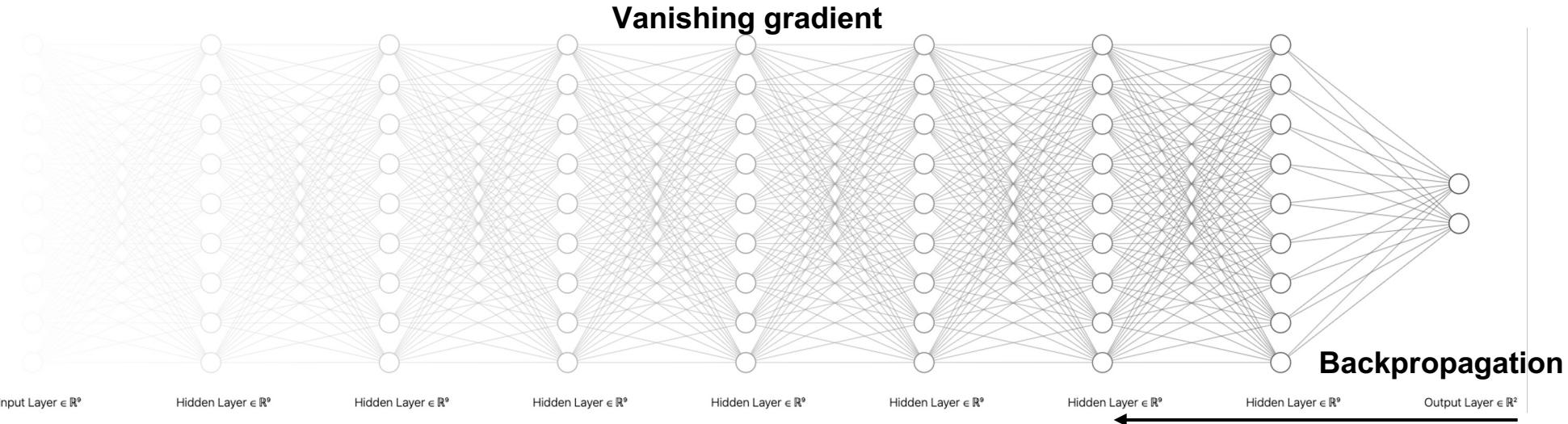
VGGNet (19 layers)

# 1.3 Going deeper with convolutions

CNN architectures

Deeper networks are harder to optimize

- Gradient vanishing / exploding

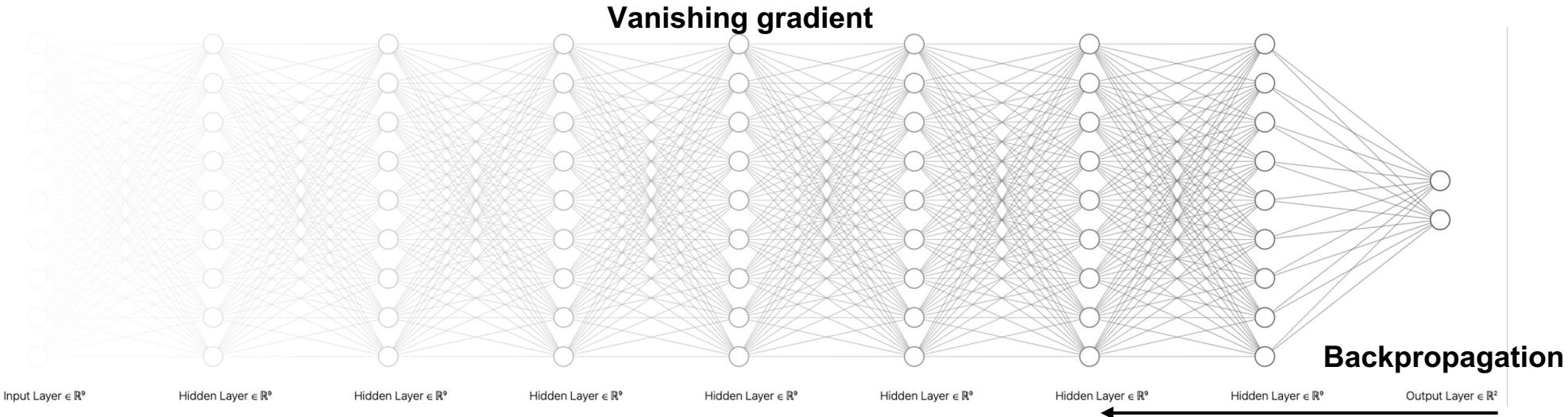


# 1.3 Going deeper with convolutions

CNN architectures

Deeper networks are harder to optimize

- Gradient vanishing / exploding
- ~~Overfitting problem~~      Degradation problem

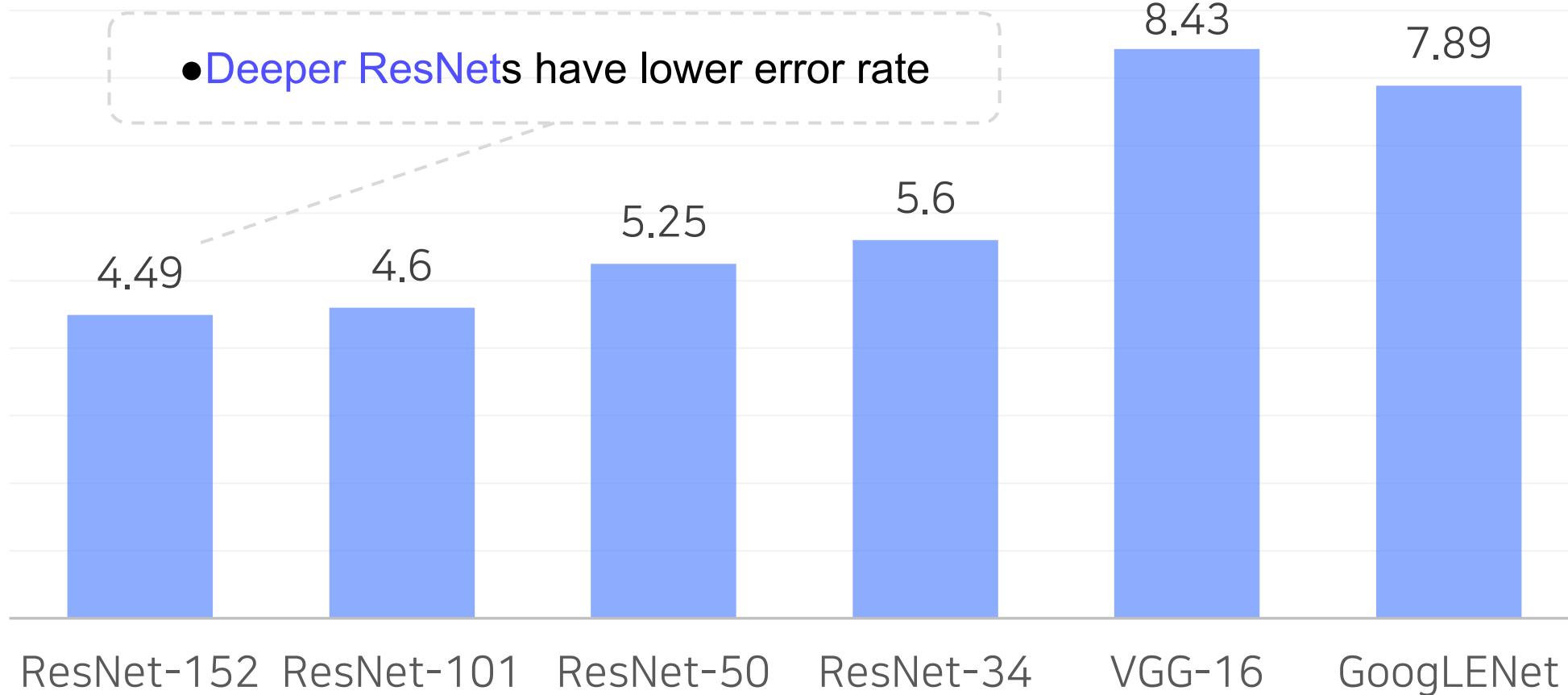


## 1.3 Going deeper with convolutions

CNN architectures

### ResNet

- Top-5 validation error on the ImageNet

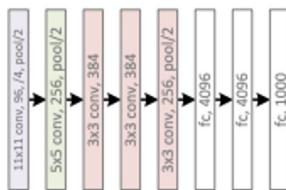


## 1.3 Going deeper with convolutions

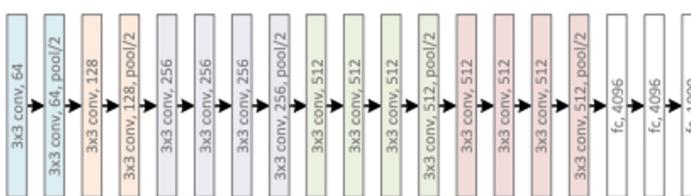
## CNN architectures

# Revolutions of depth

- Building ultra-deeper than any other networks
  - What makes it hard to build a very deep architecture?



# AlexNet (8 layers)



## VGGNet (19 layers)



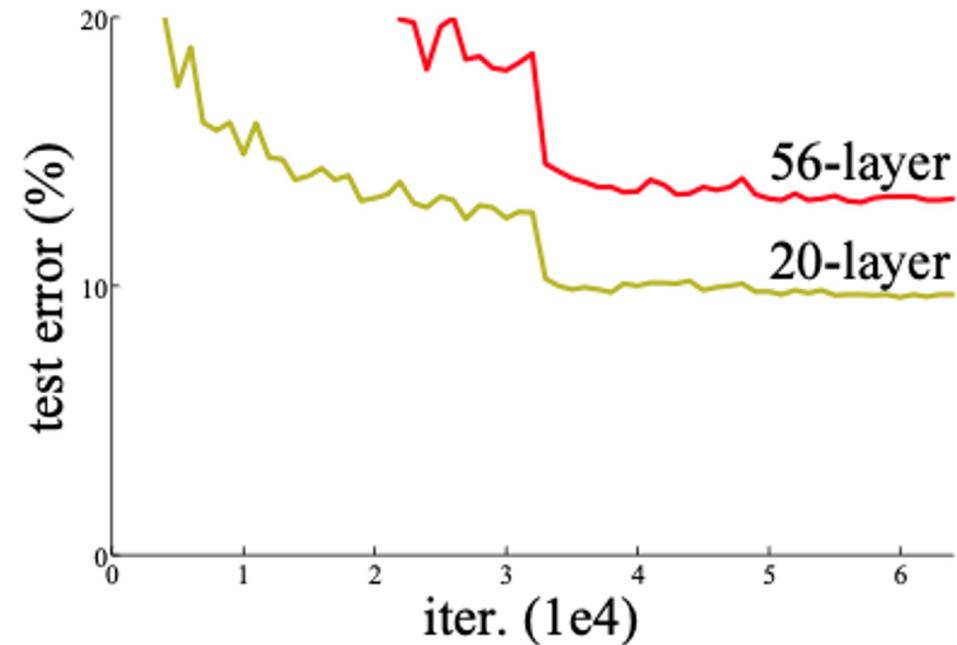
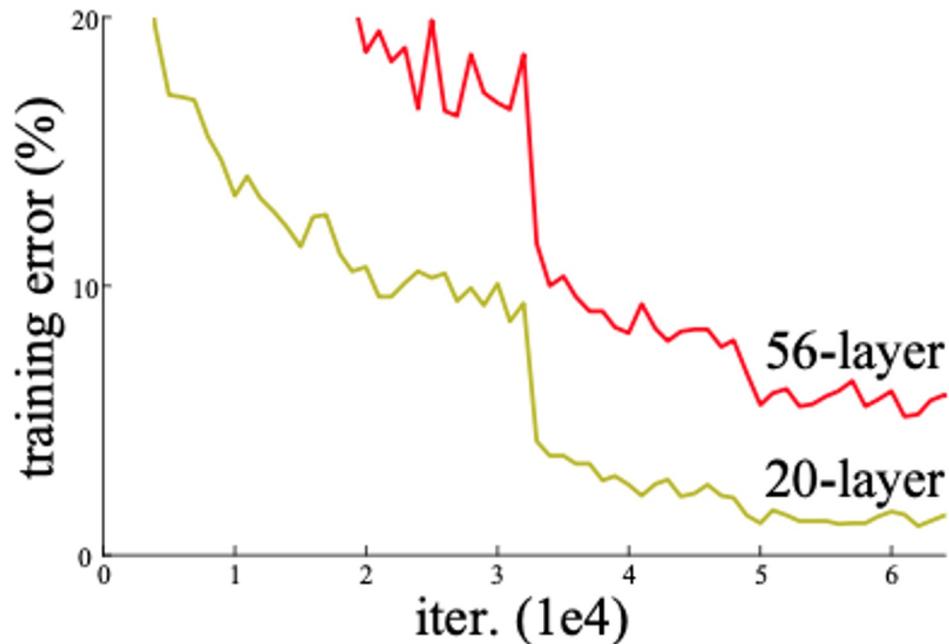
# ResNet (152 layers)

# 1.3 Going deeper with convolutions

CNN architectures

## Degradation problem

- As the network depth increases, accuracy gets saturated  $\Rightarrow$  **degrade rapidly**

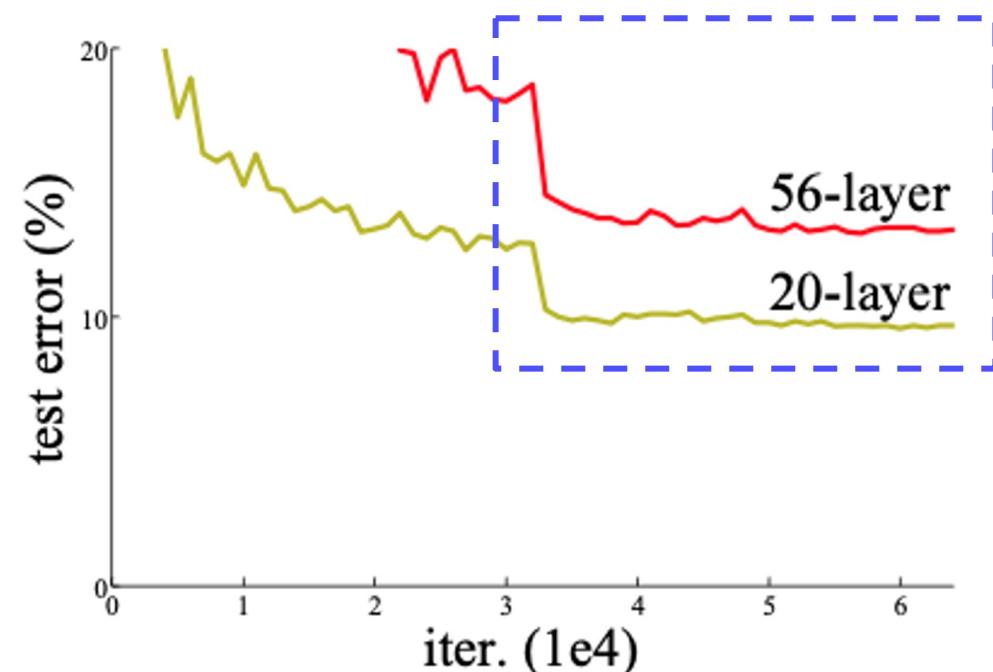
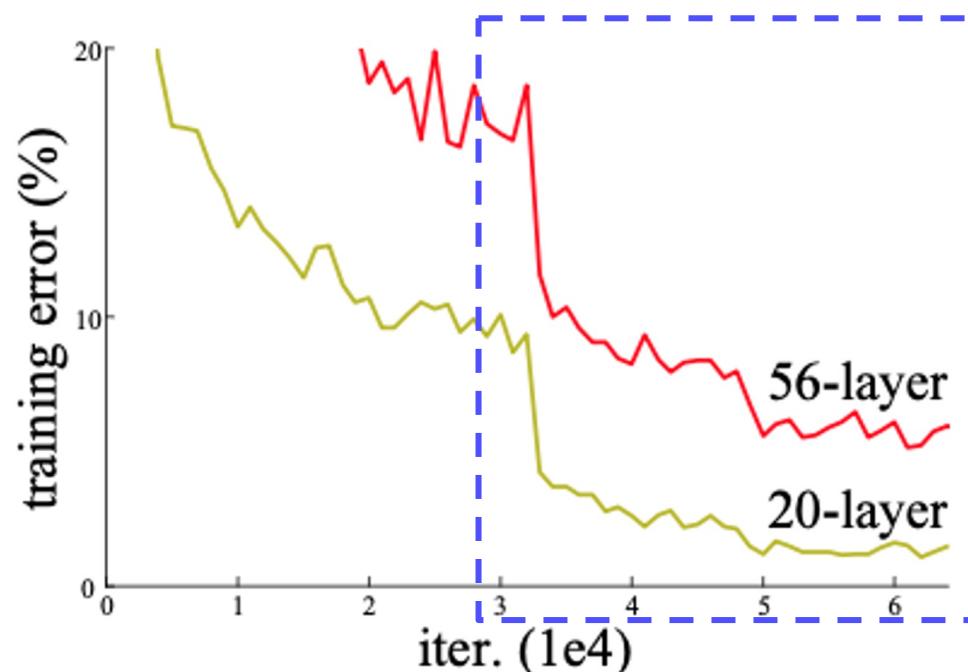


# 1.3 Going deeper with convolutions

CNN architectures

## Degradation problem

- As the network depth increases, accuracy gets saturated ⇒ **degrade rapidly**
- This is not caused by overfitting. The problem is **optimization!**



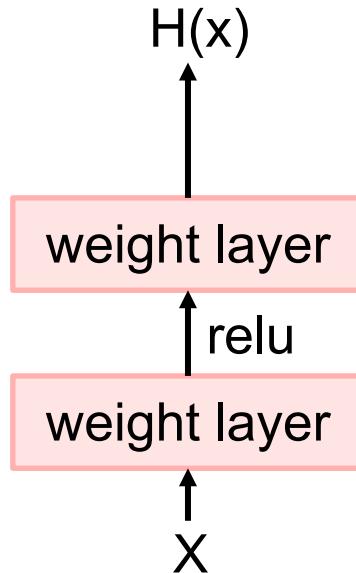
# 1.3 Going deeper with convolutions

CNN architectures

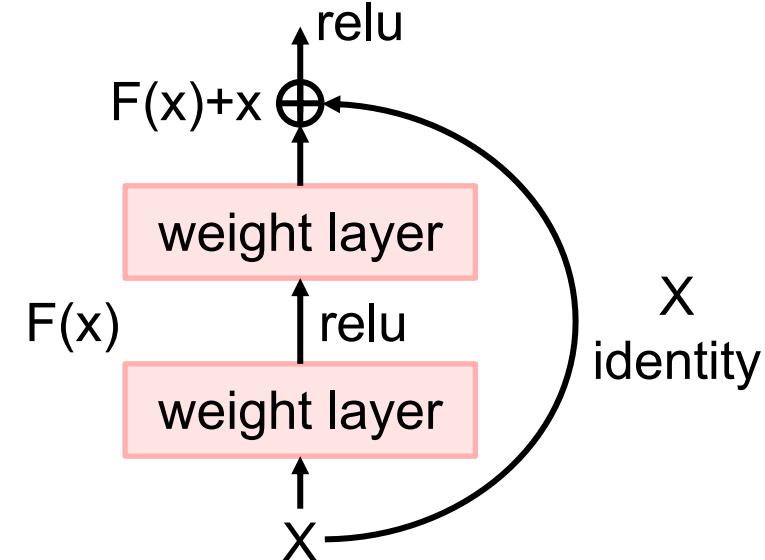
## ResNet - Hypothesis

- Plain layer: As the layers get deeper, it is hard to learn good  $H(x)$  directly

“Plain” layers



“Residual” block



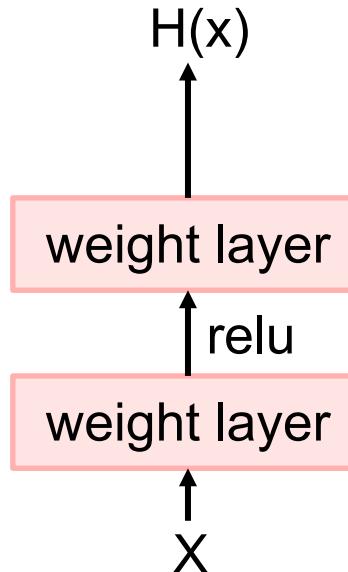
# 1.3 Going deeper with convolutions

CNN architectures

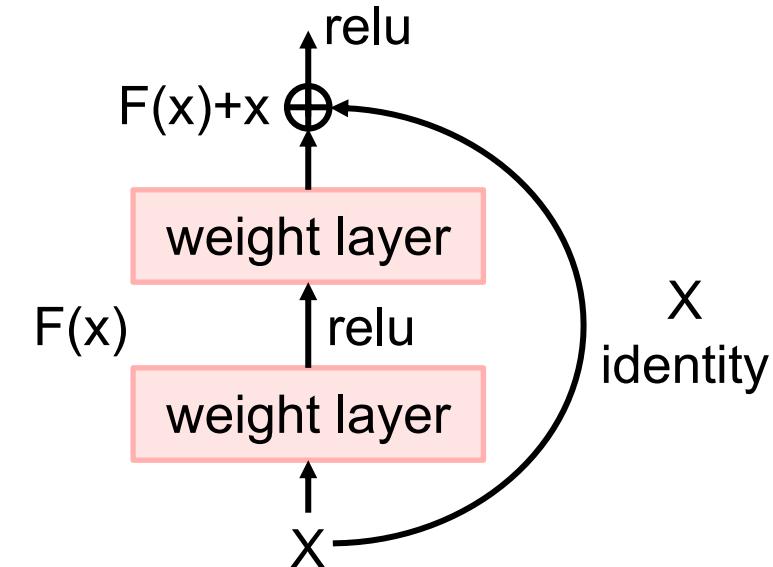
## ResNet - Hypothesis

- Plain layer: As the layers get deeper, it is hard to learn good  $H(x)$  directly
- Residual block: Instead, we learn residual
  - Target function :  $H(x) = F(x) + x$
  - Residual function :  $F(x) = H(x) - x$

“Plain” layers



“Residual” block

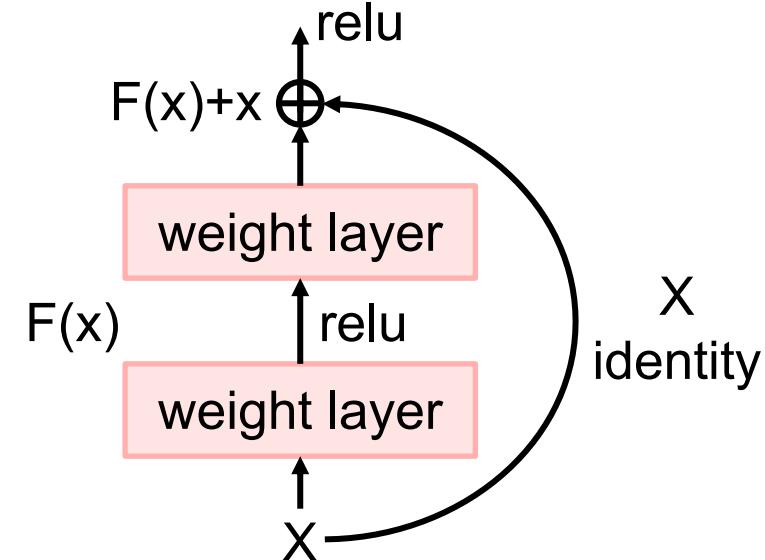
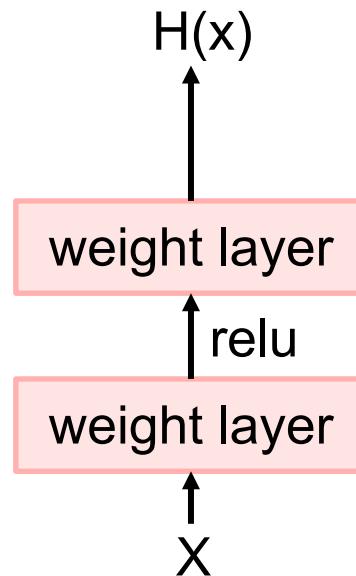


# 1.3 Going deeper with convolutions

CNN architectures

## ResNet - A solution: Shortcut connection

- Use layers to fit a residual mapping instead of directly fitting a desired underlying mapping
- The vanishing gradient problem is solved by shortcut connection
- Don't just stack layers up, instead use shortcut connection!

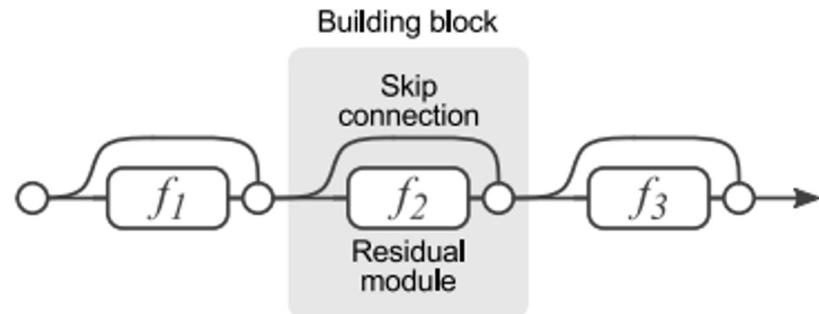


# 1.3 Going deeper with convolutions

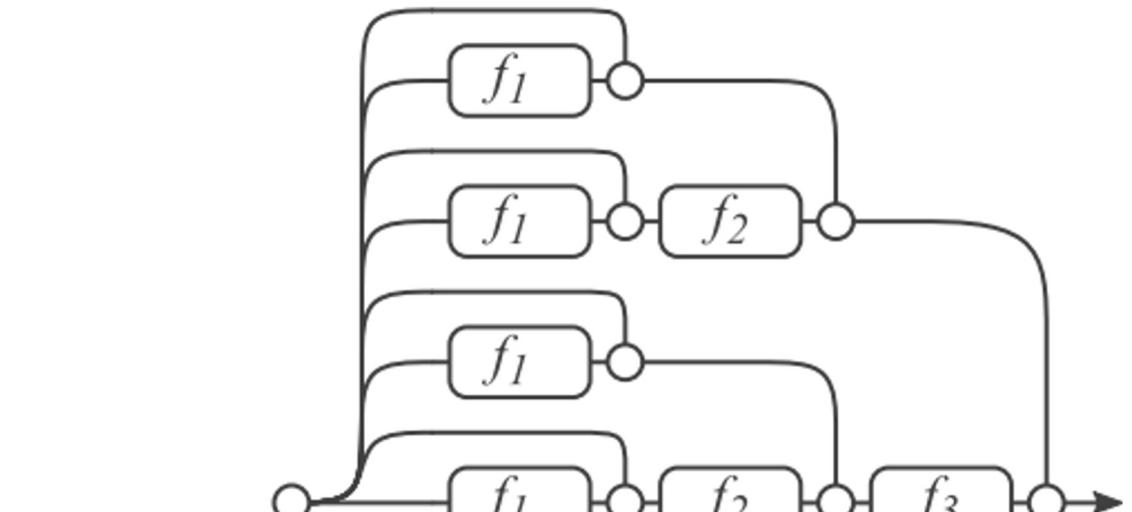
CNN architectures

## ResNet - Analysis of residual connection

- During training, gradients are mainly from relatively shorter paths
- Residual networks have  $O(2^n)$  implicit paths connecting input and output, and adding a block doubles the number of paths



(a) Conventional 3-block residual network



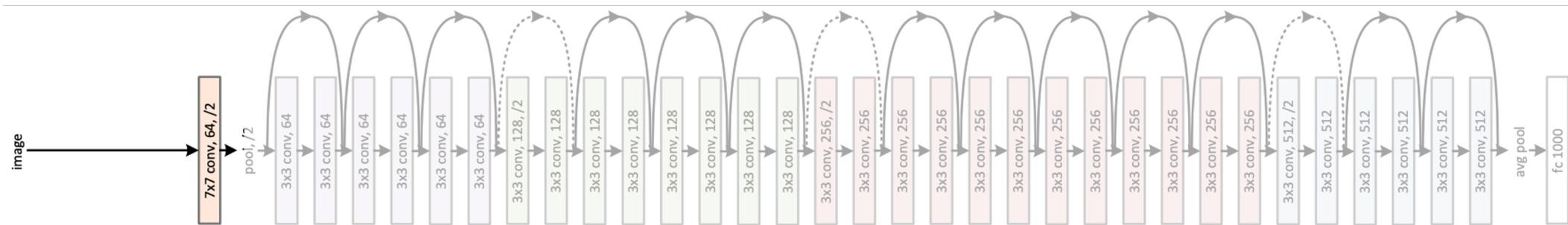
(b) Unraveled view of (a)

# 1.3 Going deeper with convolutions

CNN architectures

## ResNet - Overall architecture

- He initialization
- Additional conv layer at the beginning



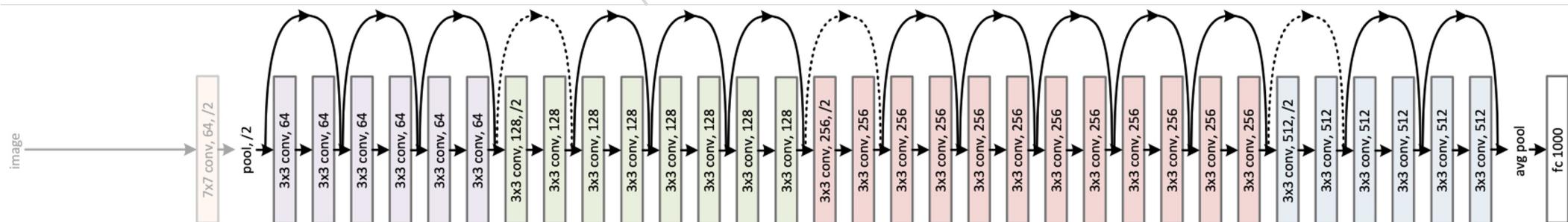
# 1.3 Going deeper with convolutions

CNN architectures

## ResNet - Overall architecture

- Stack residual blocks
- Every residual block has two 3x3 conv layers
- Batch norm after every conv layer

- He initialization
- Additional conv layer at the beginning



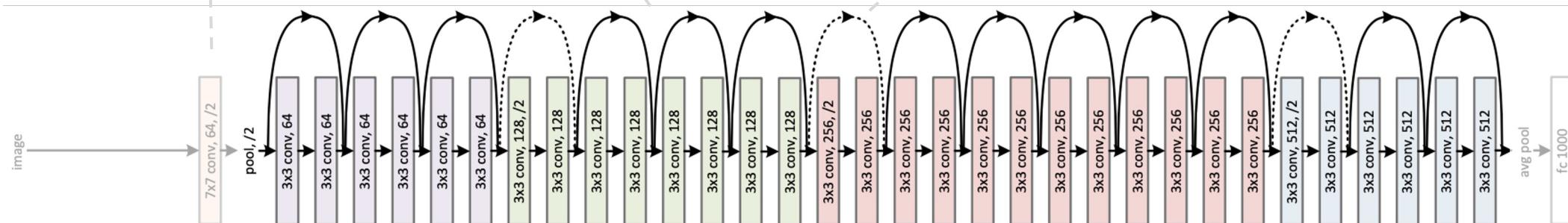
# 1.3 Going deeper with convolutions

CNN architectures

## ResNet - Overall architecture

- Stack residual blocks
- Every residual block has two 3x3 conv layers
- Batch norm after every conv layer
- Doubling the number of filters and spatially down-sampling by stride 2 instead of spatial pooling

- He initialization
- Additional conv layer at the beginning

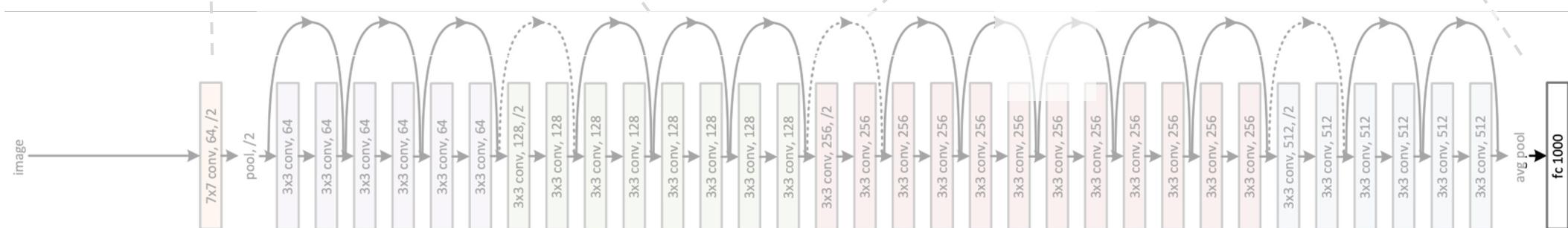


# 1.3 Going deeper with convolutions

CNN architectures

## ResNet - Overall architecture

- Stack residual blocks
- Every residual block has two 3x3 conv layers
- Batch norm after every conv layer
- He initialization
- Additional conv layer at the beginning
- Doubling the number of filters and spatially down-sampling by stride 2 instead of spatial pooling
- Only a single FC layer for output classes



# 1.3 Going deeper with convolutions

CNN architectures

## ResNet - Overall architecture

layer name	output size	18-layer	34-layer	50-layer	101-layer	152-layer
conv1	112×112			7×7, 64, stride 2		
conv2_x	56×56			3×3 max pool, stride 2		
conv3_x	28×28	$\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, \\ 3 \times 3, \\ 1 \times 1, \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, \\ 3 \times 3, \\ 1 \times 1, \end{bmatrix} \times 3$
conv4_x	14×14	$\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 4$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 4$	$\begin{bmatrix} 1 \times 1, \\ 3 \times 3, \\ 1 \times 1, \end{bmatrix} \times 4$	$\begin{bmatrix} 1 \times 1, \\ 3 \times 3, \\ 1 \times 1, \end{bmatrix} \times 4$
conv5_x	7×7	$\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 6$	$\begin{bmatrix} 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 6$	$\begin{bmatrix} 3 \times 3, 2 \\ 1 \times 1, 1024 \end{bmatrix} \times 6$	$\begin{bmatrix} 3 \times 3, 2 \\ 1 \times 1, 1024 \end{bmatrix} \times 6$
FLOPs	1×1			average pool, 1000-d fc, softmax		

```
def _forward_impl(self, x: Tensor):
    # See note [TorchScript super()]
    x = self.conv1(x)
    x = self.bn1(x)
    x = self.relu(x)
    x = self.maxpool(x)

    x = self.layer1(x)
    x = self.layer2(x)
    x = self.layer3(x)
    x = self.layer4(x)

    x = self.avgpool(x)
    x = torch.flatten(x, 1)
    x = self.fc(x)

    return x
```

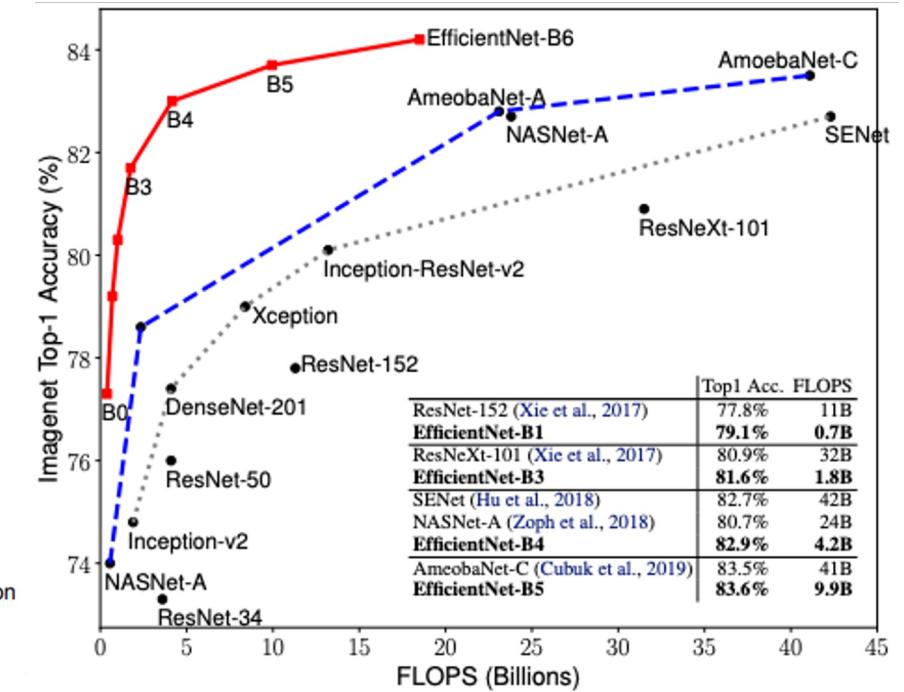
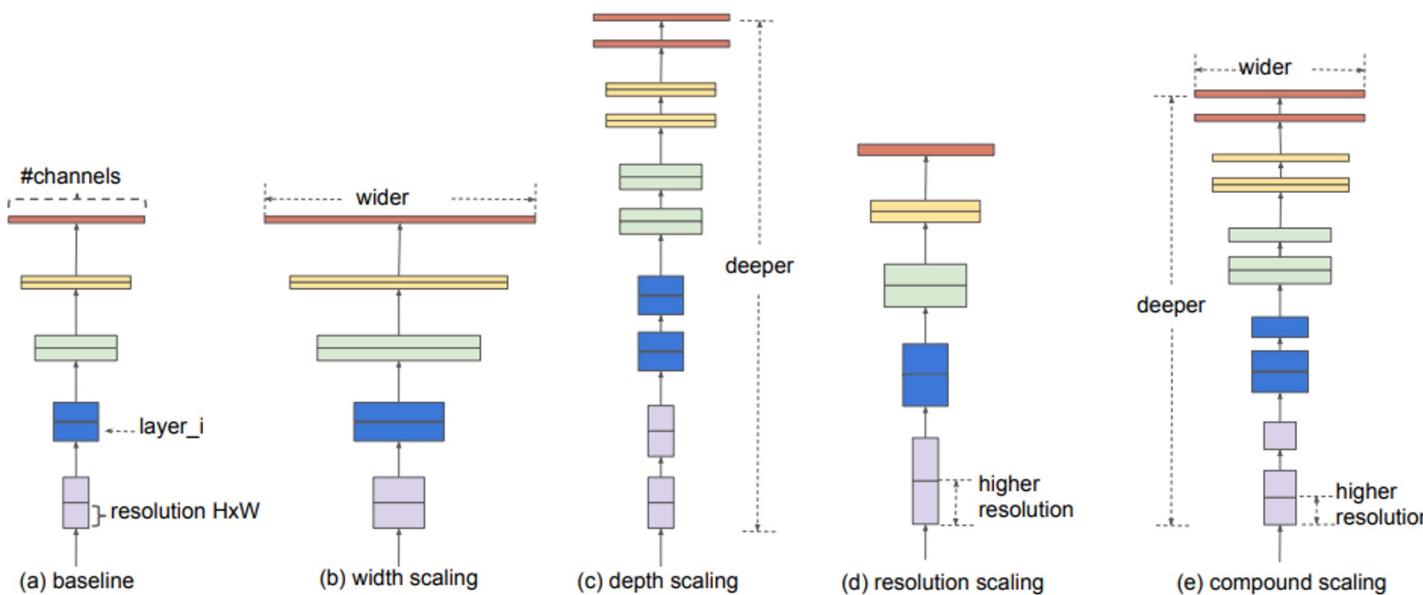
1

# 1.3 Going deeper with convolutions

CNN architectures

## Beyond ResNets - EfficientNet

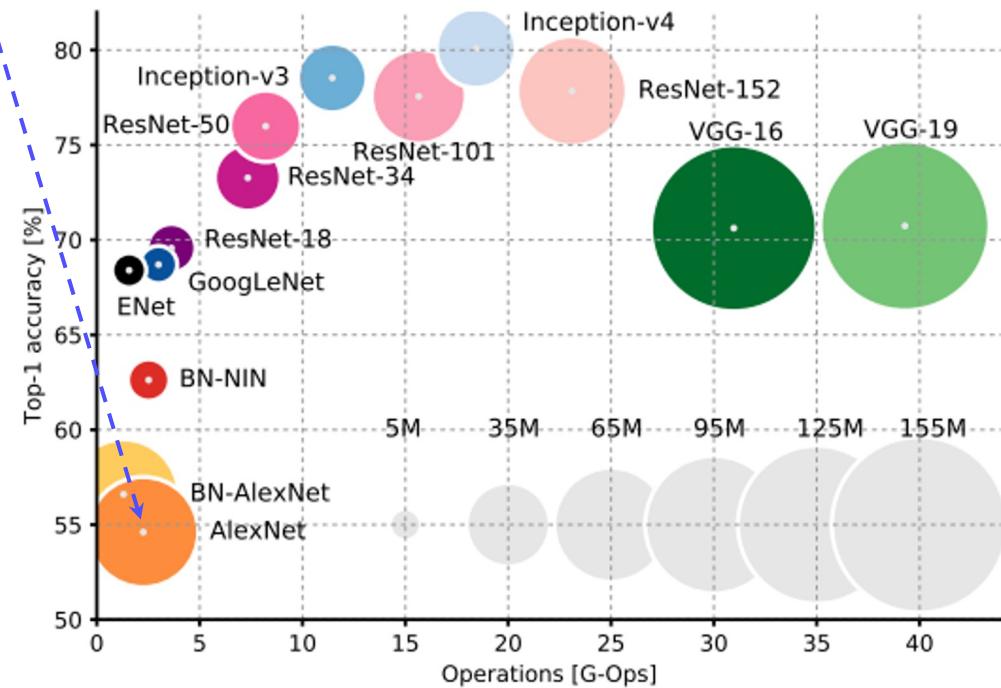
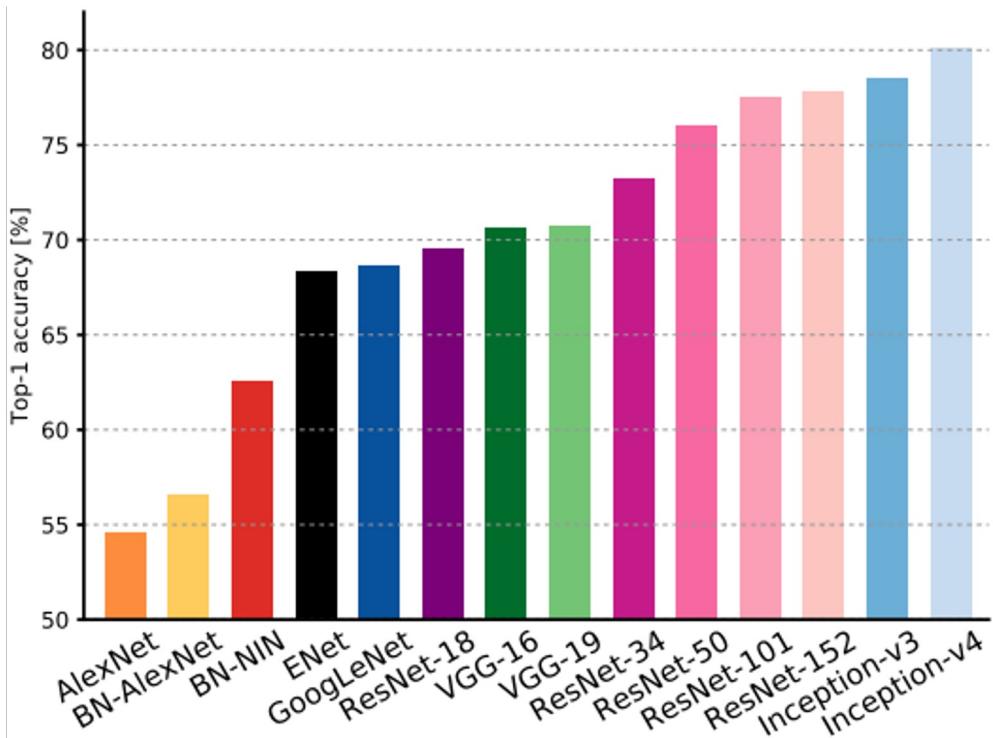
- Building **deep**, **wide**, and **high resolution** networks in an efficient way



## 1.4 CNN backbones

CNN architectures

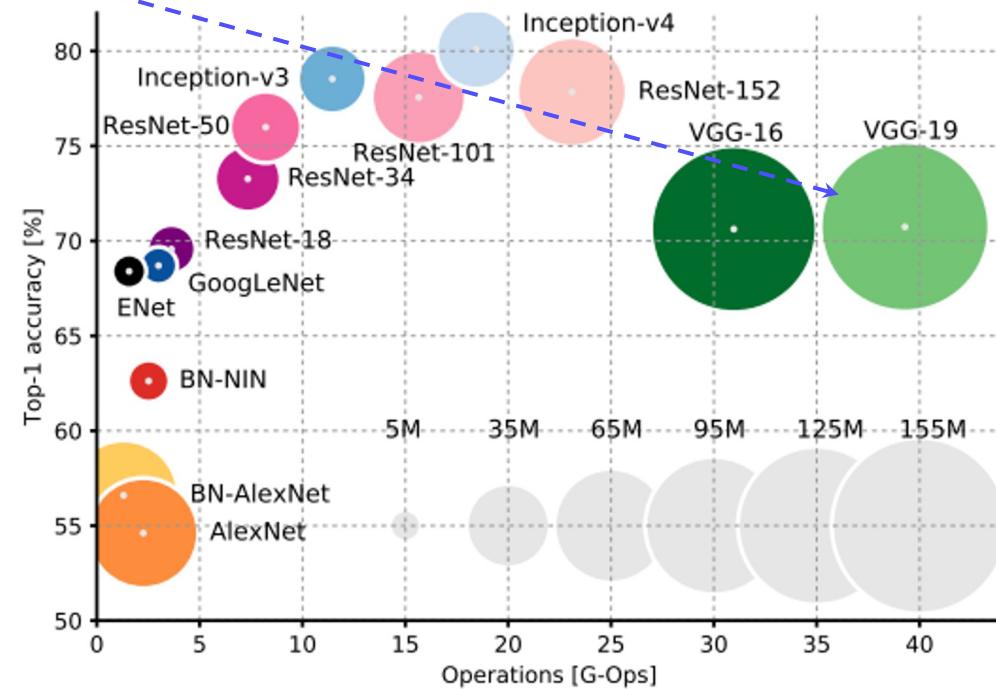
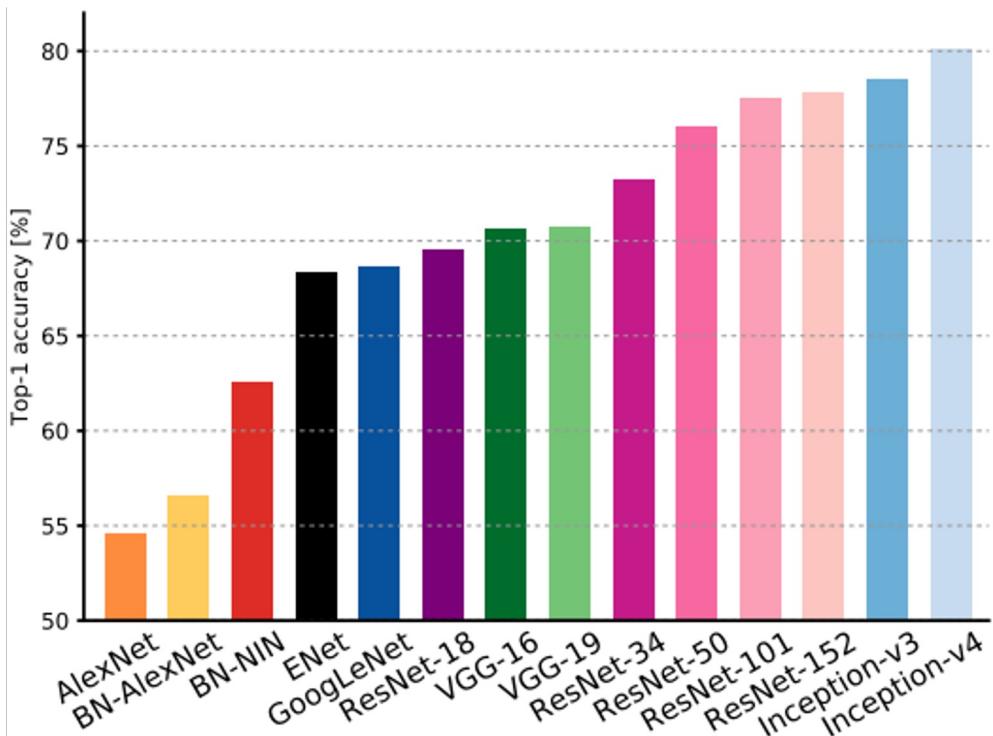
- **AlexNet**: simple CNN architecture
- Simple computation, but heavy memory size
- Low accuracy



## 1.4 CNN backbones

CNN architectures

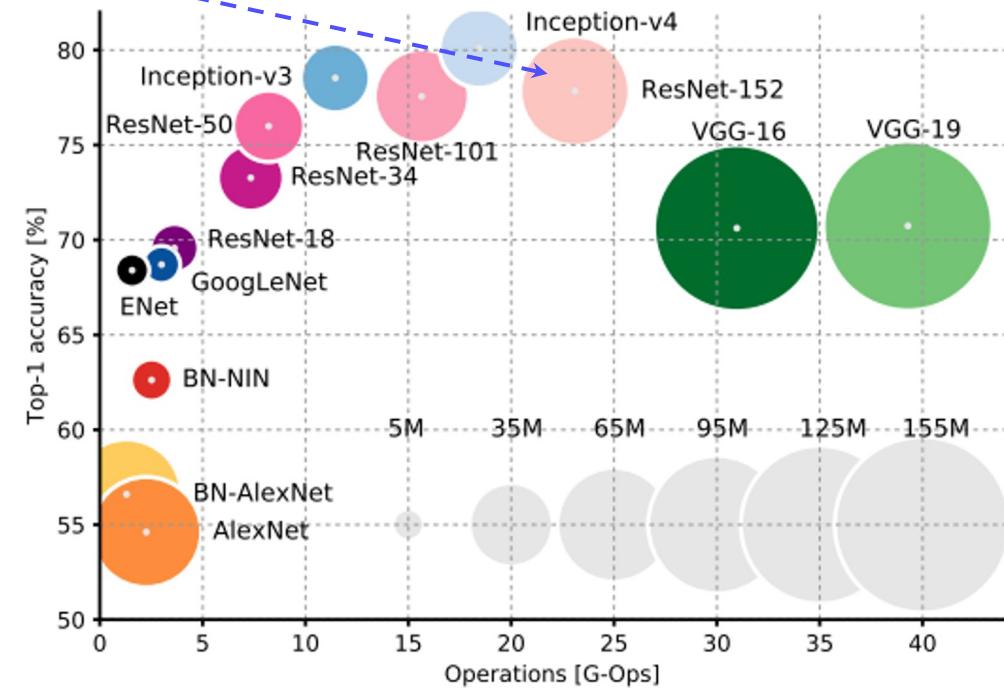
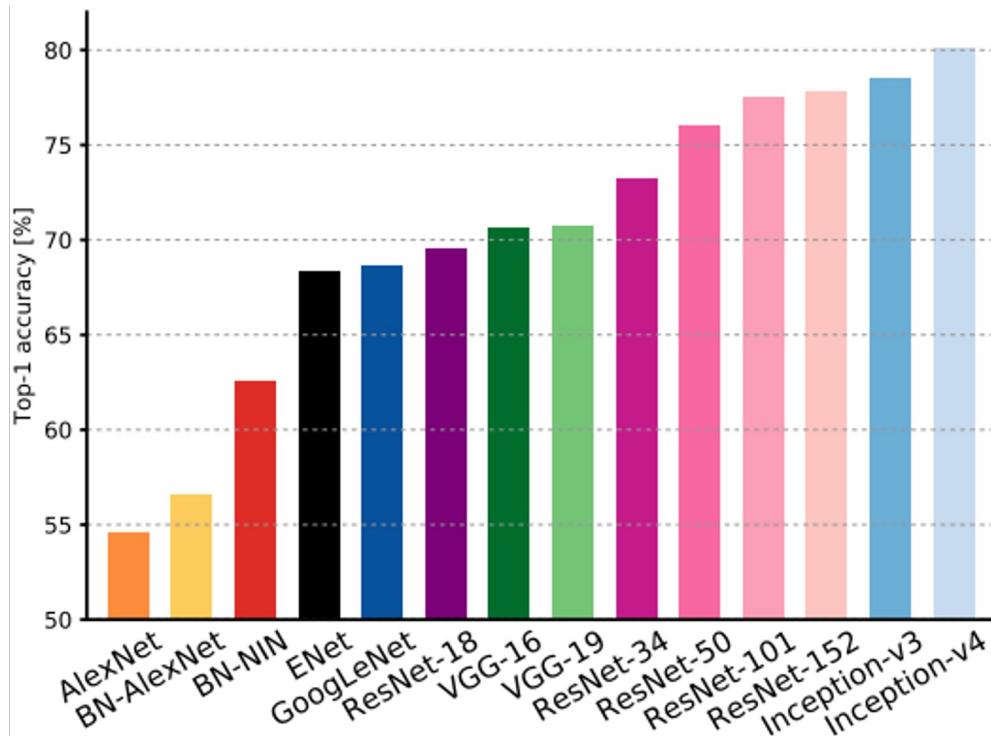
- **VGGNet**: simple with 3x3 convolutions
- Highest memory, the heaviest computation



## 1.4 CNN backbones

CNN architectures

- **ResNet**: deeper layers with residual blocks
- Moderate efficiency (depending on the model)

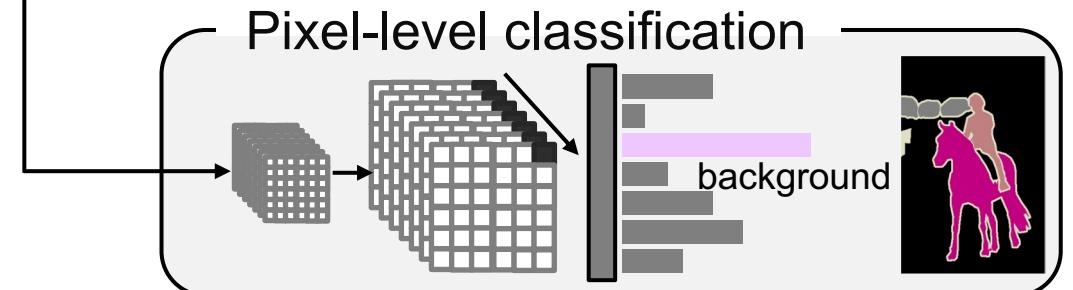
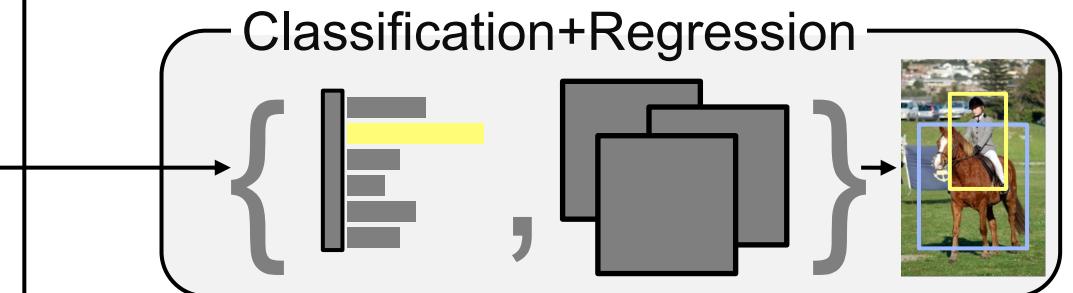
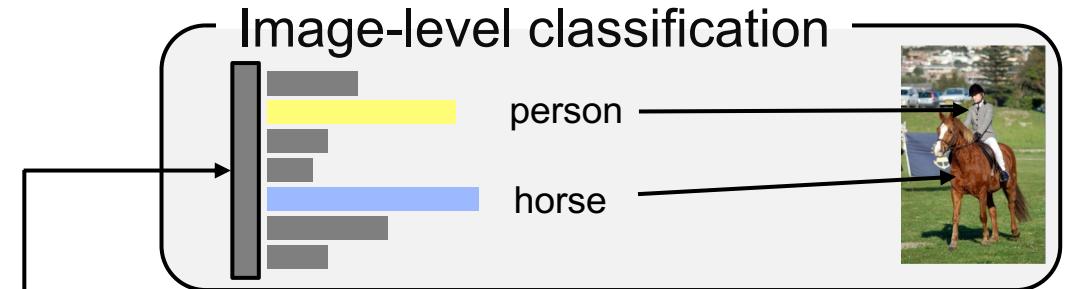
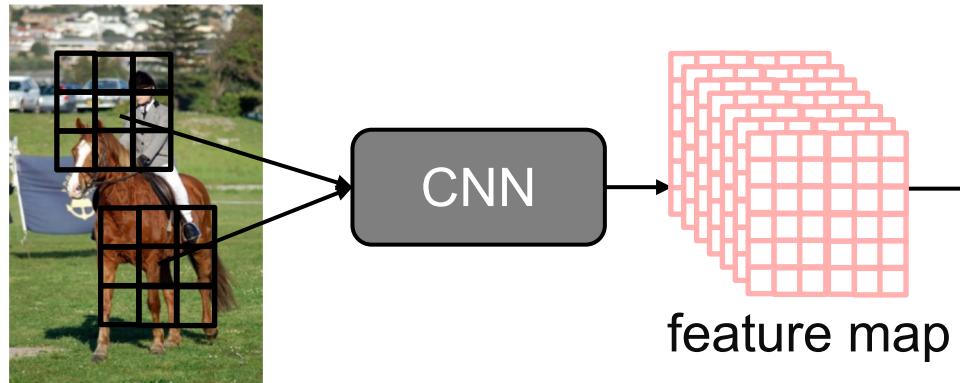


## 1.4 CNN backbones

CNN architectures

### Simple but powerful backbone model

- **VGGNet** and **ResNet** are typically used as a backbone model for many tasks
- Constructed with simple 3x3 conv layers



# Reference

---

## 1. CNN architectures

- Lecun et al., Gradient-based Learning Applied to Document Recognition, Proceedings of the IEEE 1998
- Krizhevsky et al., ImageNet Classification with Deep Convolutional Neural Networks, NeurIPS 2012
- Simonyan and Zisserman, Very Deep Convolutional Networks for Large-Scale Image Recognition, ICLR 2015
- Szegedy et al., Going Deeper with Convolution, CVPR 2015
- He et al., Deep Residual Learning for Image Recognition, CVPR 2015
- Veit et al., Residual Networks Behave Like Ensembles of Relatively Shallow Networks, NIPS 2016
- Huang et al., Densely Connected Convolutional Networks, CVPR 2017
- Hu et al., Squeeze-and-Excitation Networks, CVPR 2018
- Tan and Le, EfficientNet: Rethinking Model Scalining for Convolutional Neural Networks, ICML 2019
- Dai et al., Deformable Convolutional Networks, ICCV 2017
- Canziani et al., An Analysis of Deep Neural Network Models for Practical Applications, CVPR 2016

# Reference

---

## 2. Vision Transformers

- Jeffery L., Finding structure in time, Cognitive science, 1990
- Mike et al., Bidirectional Recurrent Neural Networks, Transactions on signal processing, 1997
- Vaswani et al., Attention is all you need, NeurIPS, 2017
- Dosovitskiy et al., An image is worth 16x16 words: Transformers for image recognition at scale, ICLR, 2021
- Devlin et al., BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding, NAACL, 2019
- Kaplan et al. "Scaling laws for neural language models." arXiv preprint arXiv 2020
- Zhai et al. "Scaling vision transformers." CVPR 2022
- Liu et al. "Swin transformer: Hierarchical vision transformer using shifted windows." ICCV 2021
- He et al. "Masked autoencoders are scalable vision learners." CVPR 2022
- Caron et al. "Emerging properties in self-supervised vision transformers." ICCV 2021

# Thank you

**Prof. Jeon, Sangryul**

Computer Vision Lab.

Pusan National University, Korea

Tel: +82-51-510-2423

Web: <http://sr-jeon.github.io/>

E-mail: [srjeonn@pusan.ac.kr](mailto:srjeonn@pusan.ac.kr)