

# Computer Vision

## Image classification II

---

Tae-Hyun Oh (오태현)

전자전기공학과

POSTECH

Slide by Sungbin Kim (김성빈)

TAs: {Dongmin Choi , Jongha Kim, Juyong Lee, Sungbin Kim} (in alphabetic order)

## 1. Problems with deeper layers

- 1.1 Going deeper with convolutions
- 1.2 Hard to optimize

## 2. CNN architectures for image classification 2

- 2.1 GoogLeNet
- 2.2 ResNet
- 2.3 Beyond ResNet

## 3. Summary of image classification

- 3.1 Summary of image classification
- 3.2 CNN backbones

1.

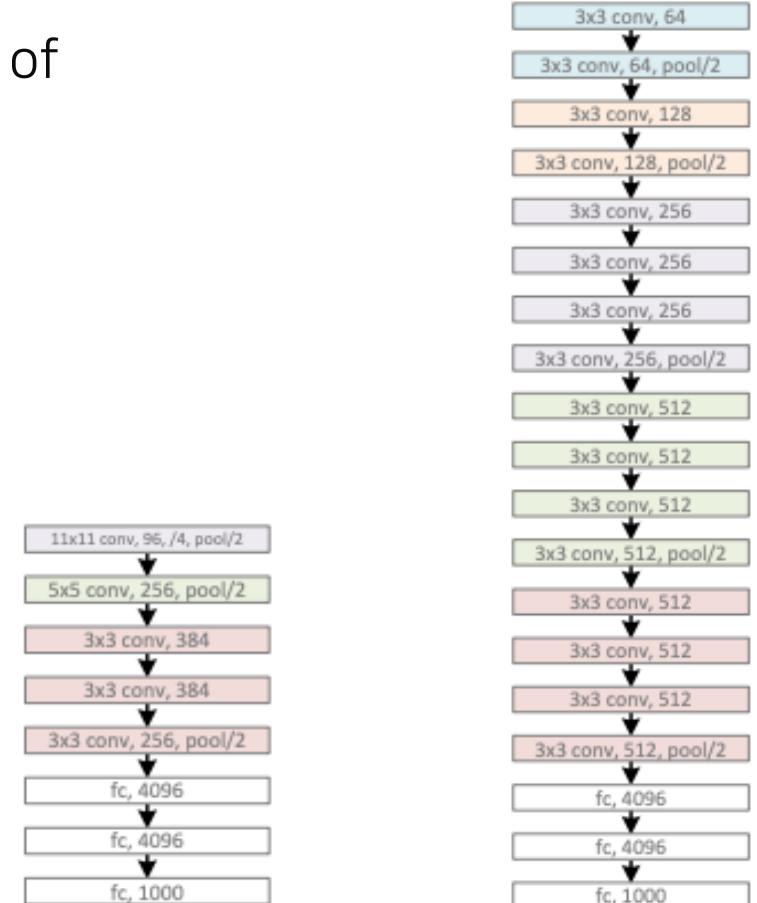
# Problems with deeper layers

# 1.1 Going deeper with convolutions

Problems with deeper layers

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)

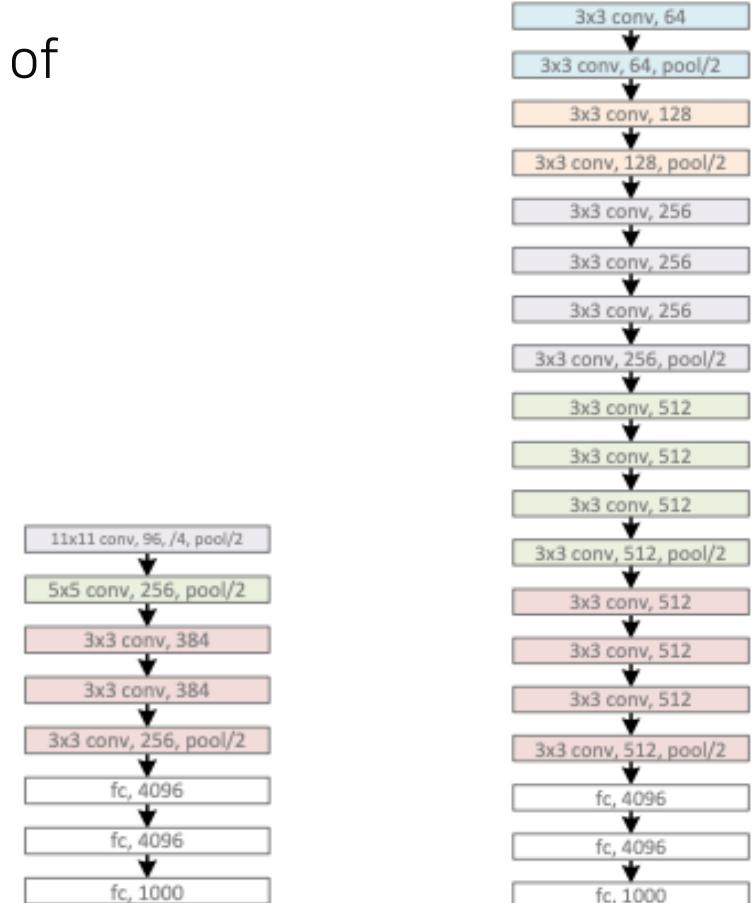


# 1.1 Going deeper with convolutions

Problems with deeper layers

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)

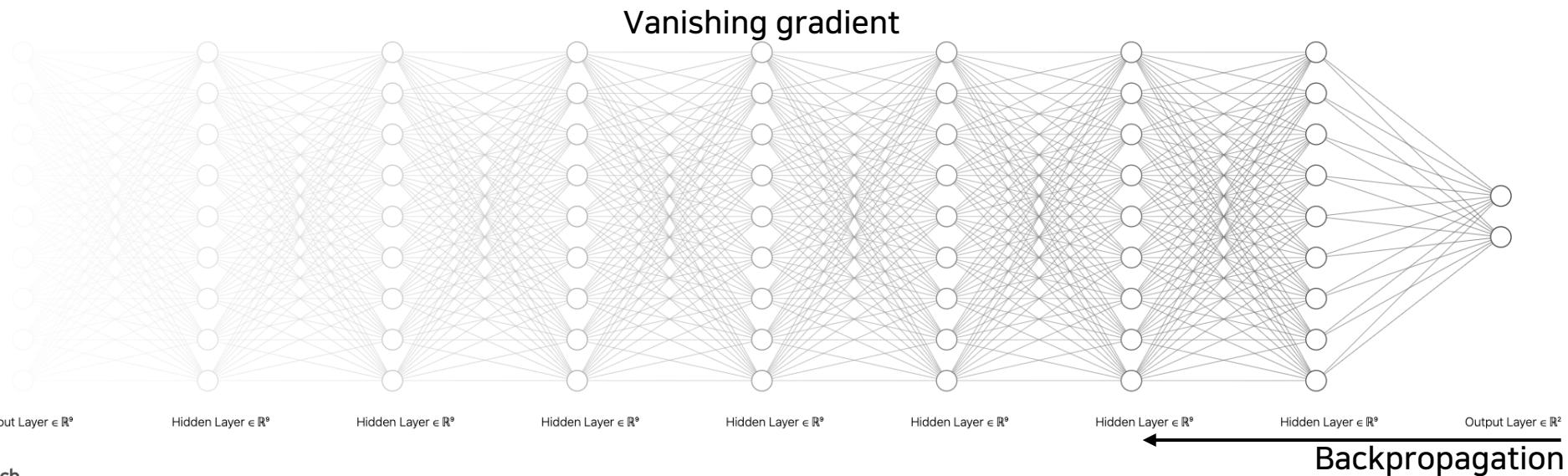
X

## 1.2 Hard to optimize

Problems with deeper layers

Deeper networks are harder to optimize

- Gradient vanishing / exploding

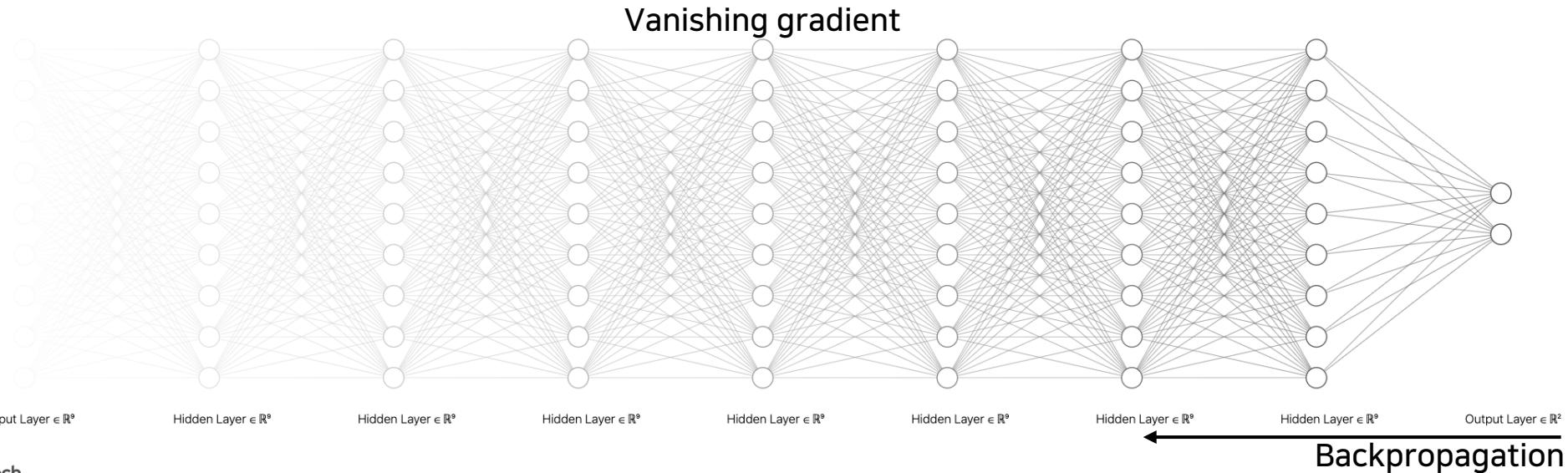


## 1.2 Hard to optimize

Problems with deeper layers

Deeper networks are harder to optimize

- Gradient vanishing / exploding
- Computationally complex

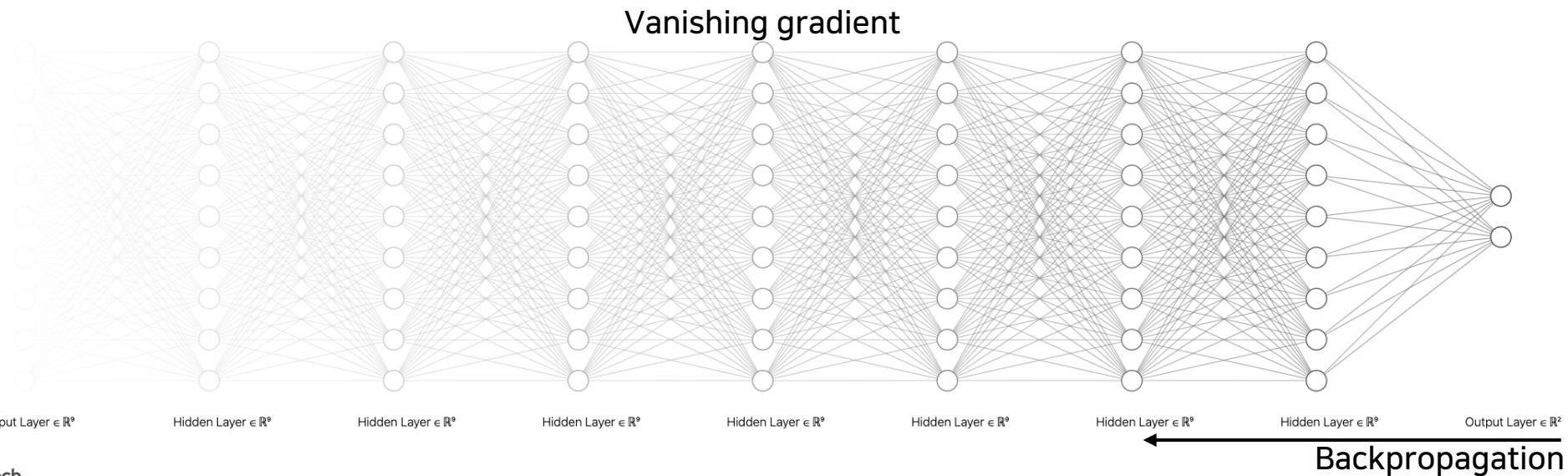


## 1.2 Hard to optimize

Problems with deeper layers

Deeper networks are harder to optimize

- Gradient vanishing / exploding
- Computationally complex
- Overfitting problem      Degradation problem



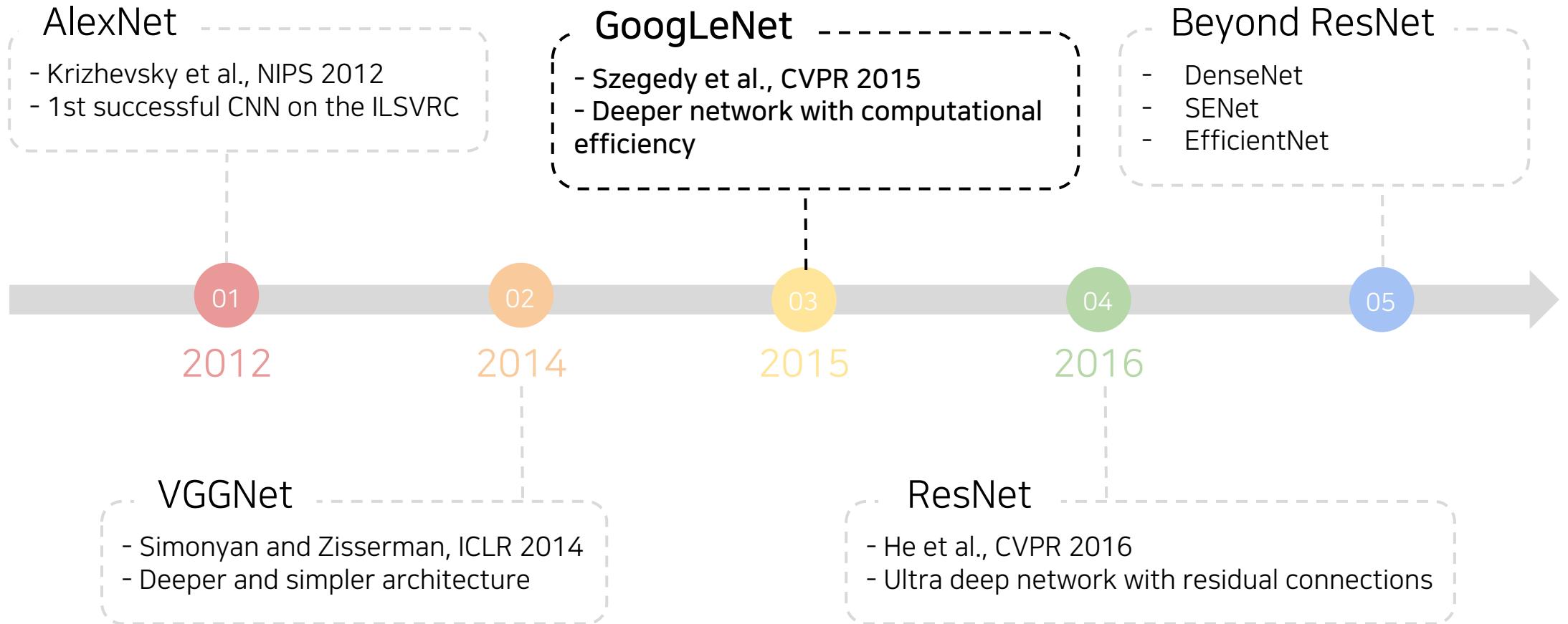
2.

## CNN architectures for image classification 2

## 2.1 GoogLeNet

CNN architectures for image classification 2

[Szegedy et al., CVPR 2015]



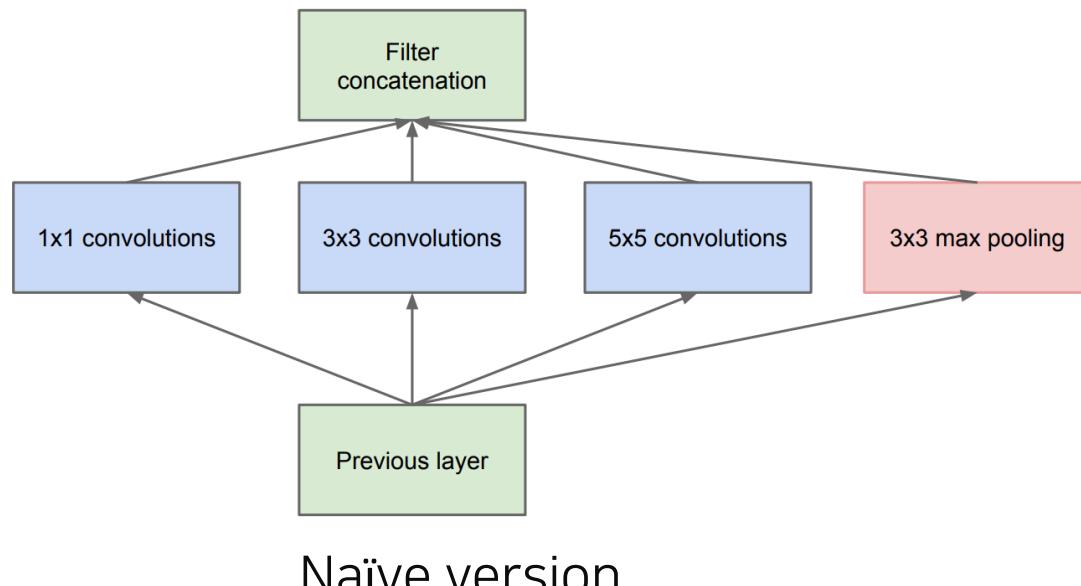
## 2.1 GoogLeNet

CNN architectures for image classification 2

### Inception module

[Szegedy et al., CVPR 2015]

- Apply multiple filter operations on input activation from the previous layer:
  - 1x1, 3x3, 5x5 convolution filters
  - 3x3 pooling operation
- **Concatenate** all filter outputs together along the channel axis



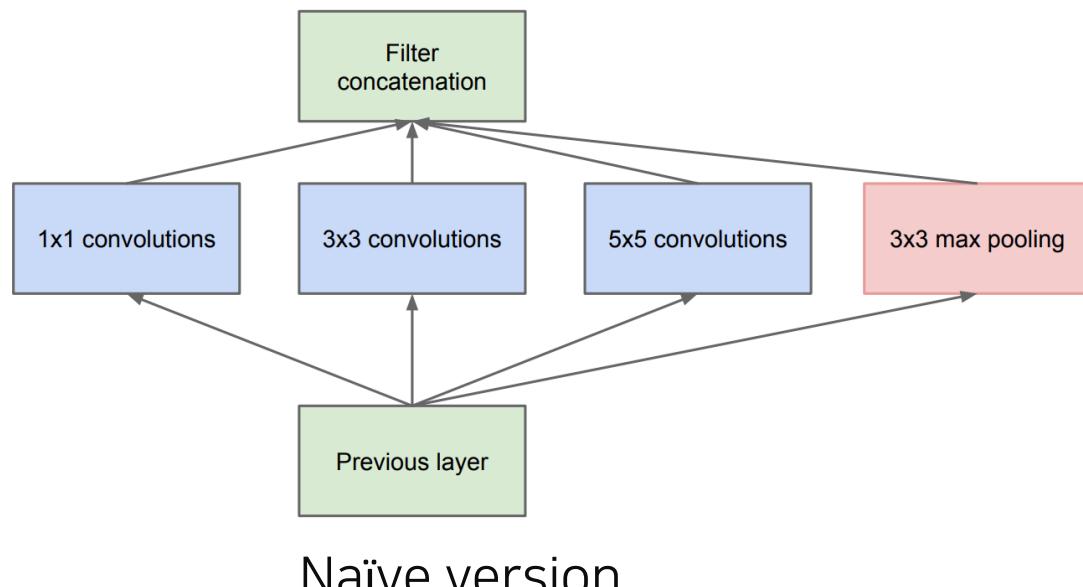
## 2.1 GoogLeNet

CNN architectures for image classification 2

### Inception module

[Szegedy et al., CVPR 2015]

- The increased network size increases the use of computational resources
- Use 1x1 convolutions!



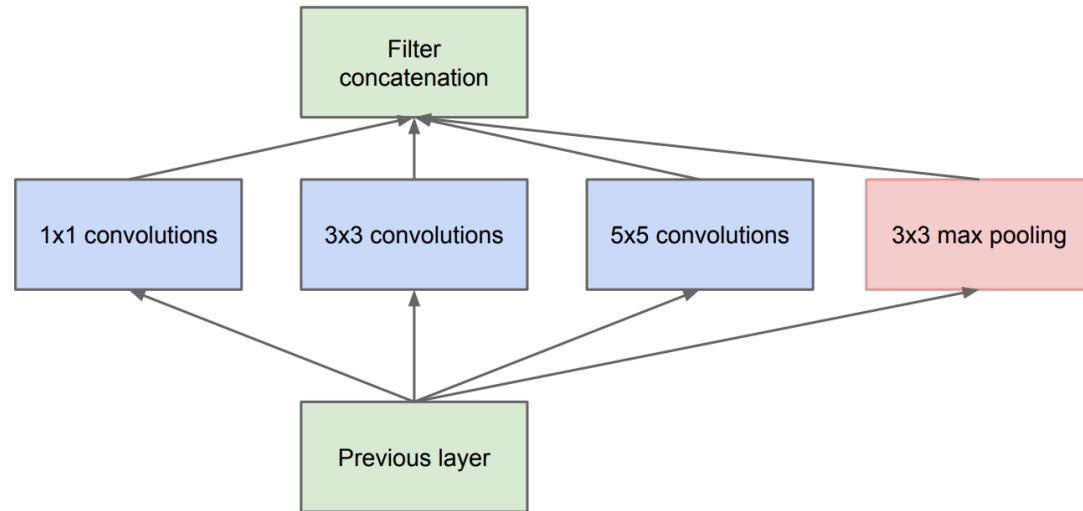
## 2.1 GoogLeNet

CNN architectures for image classification 2

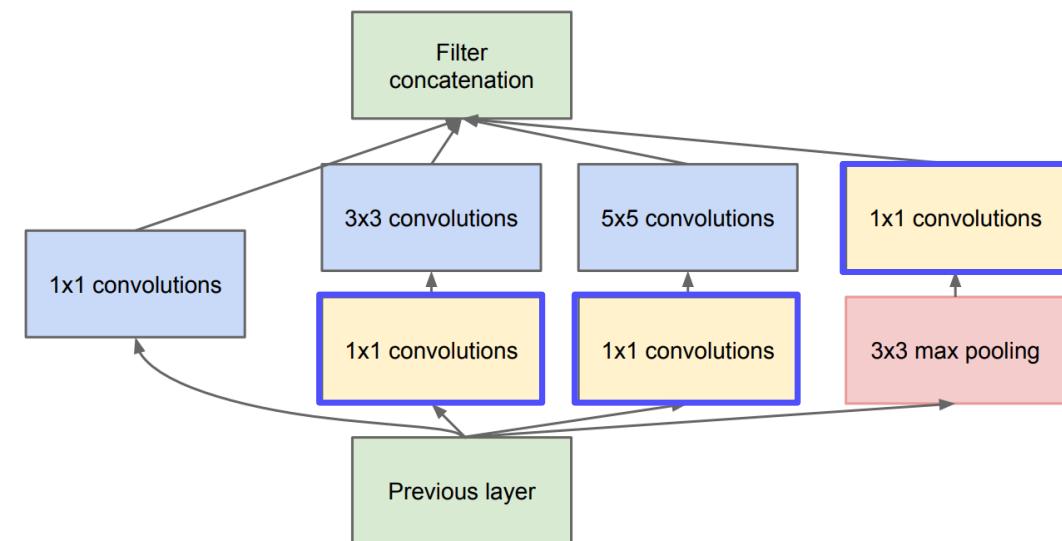
### Inception module

[Szegedy et al., CVPR 2015]

- Use 1x1 convolutions as "bottleneck" layers that reduce the number of channels



Naïve version



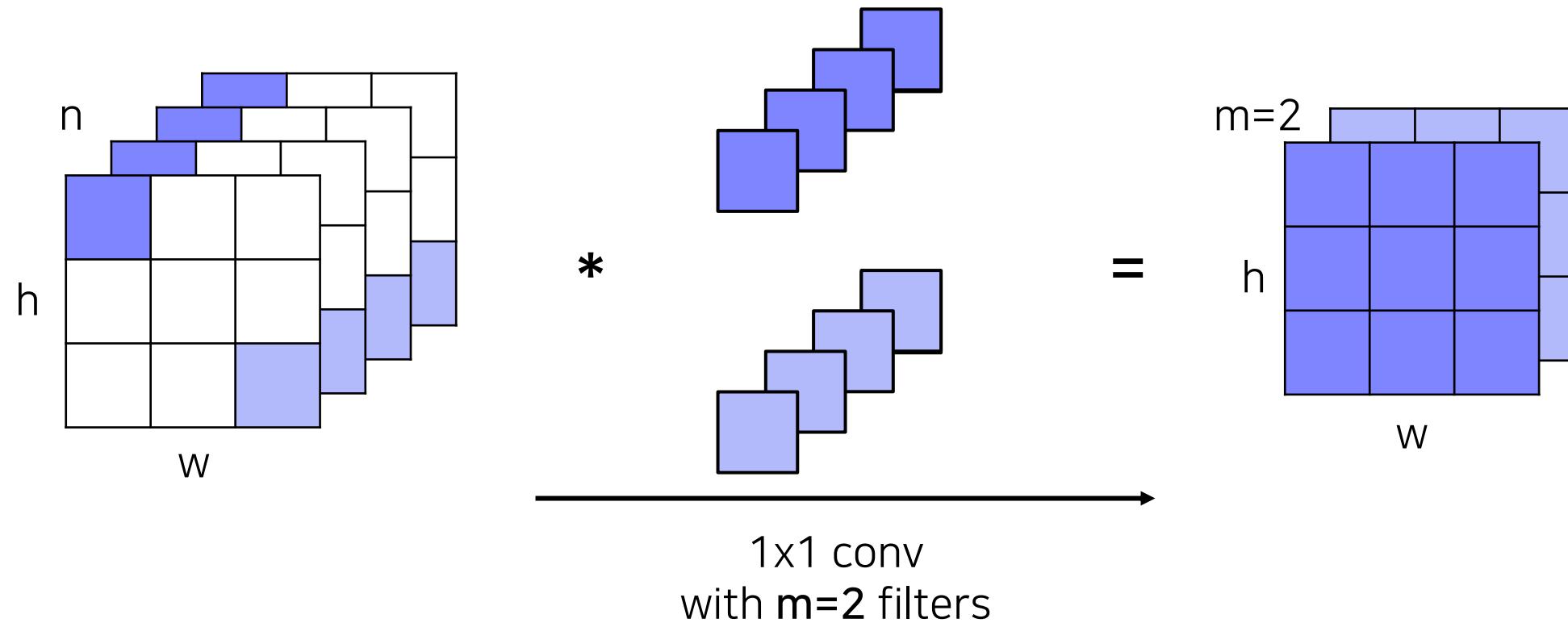
Dimension reduced version

## 2.1 GoogLeNet

CNN architectures for image classification 2

1x1 convolutions

[Szegedy et al., CVPR 2015]



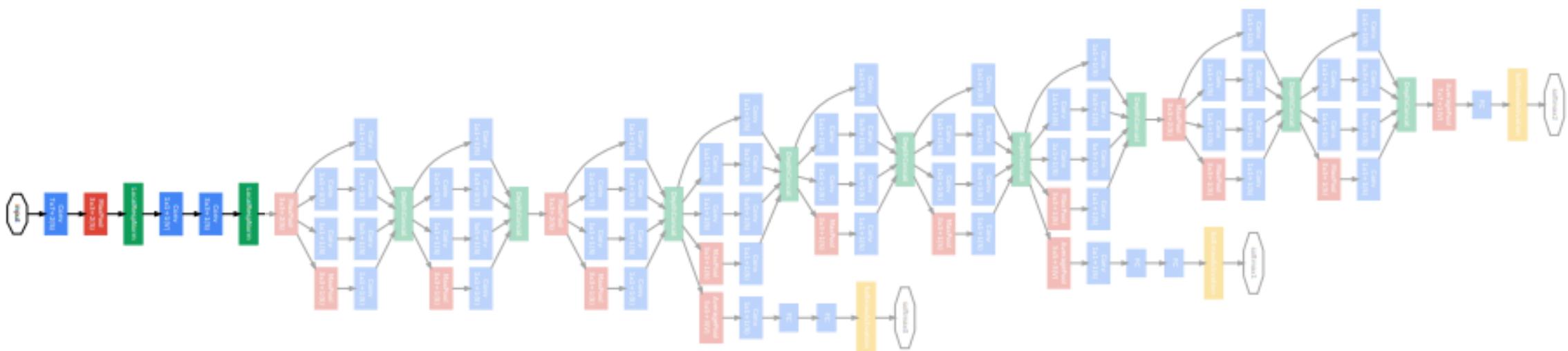
## 2.1 GoogLeNet

CNN architectures for image classification 2

Overall architecture

[Szegedy et al., CVPR 2015]

- Stem network: vanilla convolution networks



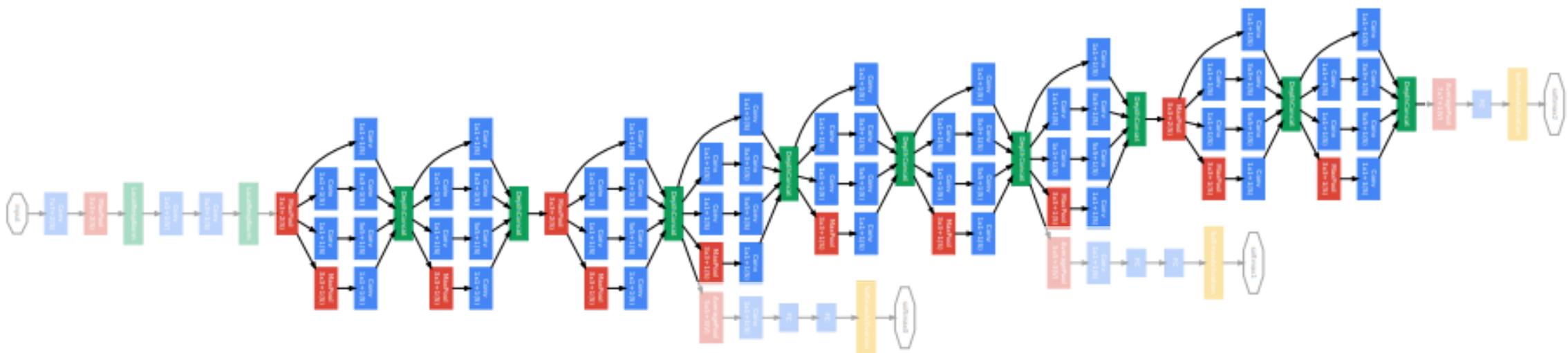
## 2.1 GoogLeNet

CNN architectures for image classification 2

Overall architecture

[Szegedy et al., CVPR 2015]

- Stem network: vanilla convolution networks
- Stacked inception modules



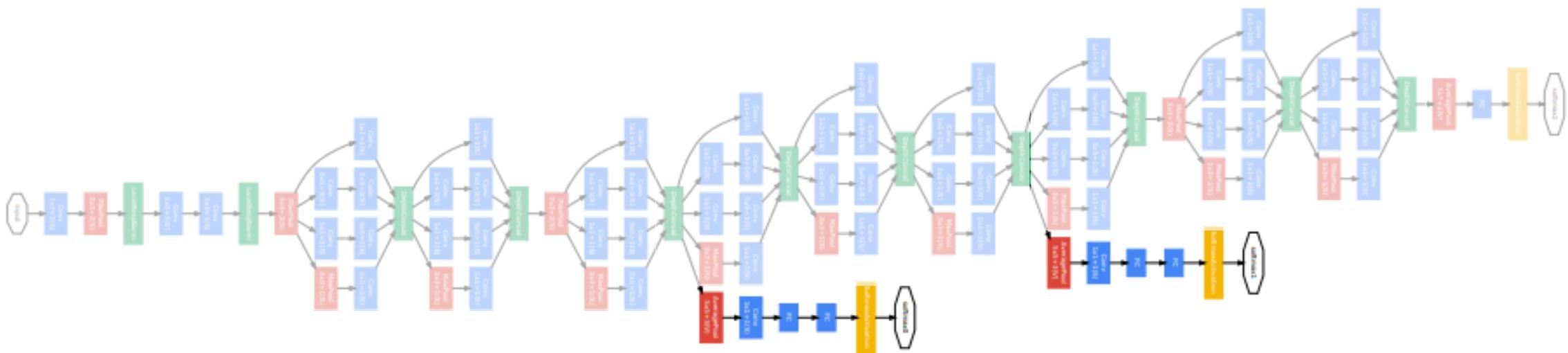
## 2.1 GoogLeNet

CNN architectures for image classification 2

Overall architecture

[Szegedy et al., CVPR 2015]

- Stem network: vanilla convolution networks
- Stacked inception modules
- Auxiliary classifiers



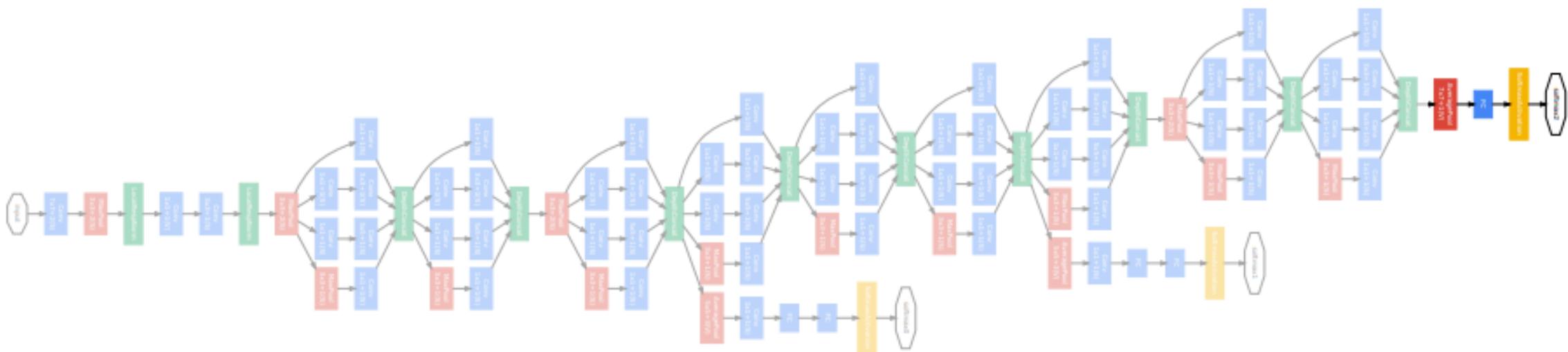
## 2.1 GoogLeNet

CNN architectures for image classification 2

Overall architecture

[Szegedy et al., CVPR 2015]

- Stem network: vanilla convolution networks
- Stacked inception modules
- Auxiliary classifiers
- Classifier output (a single FC layer)



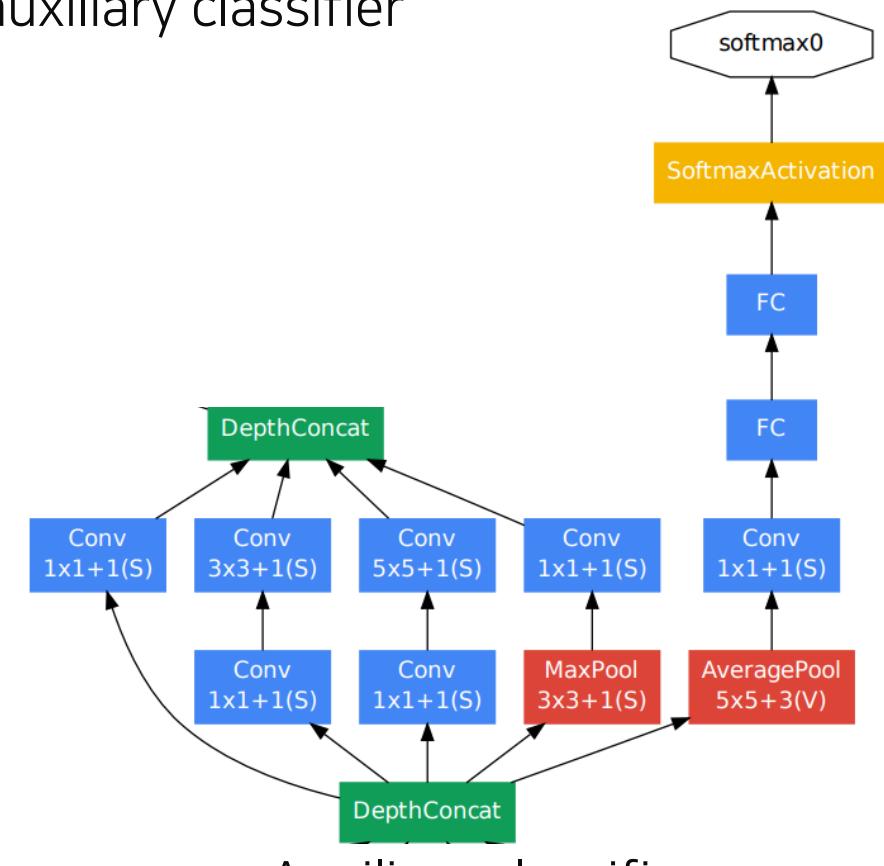
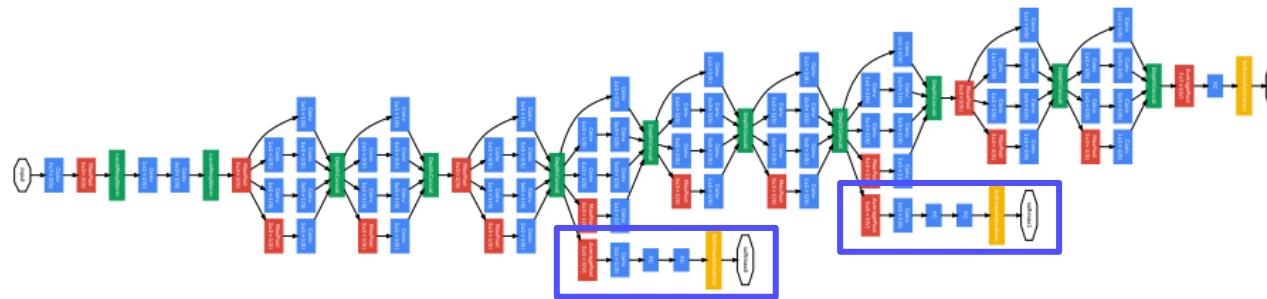
## 2.1 GoogLeNet

CNN architectures for image classification 2

### Auxiliary classifier

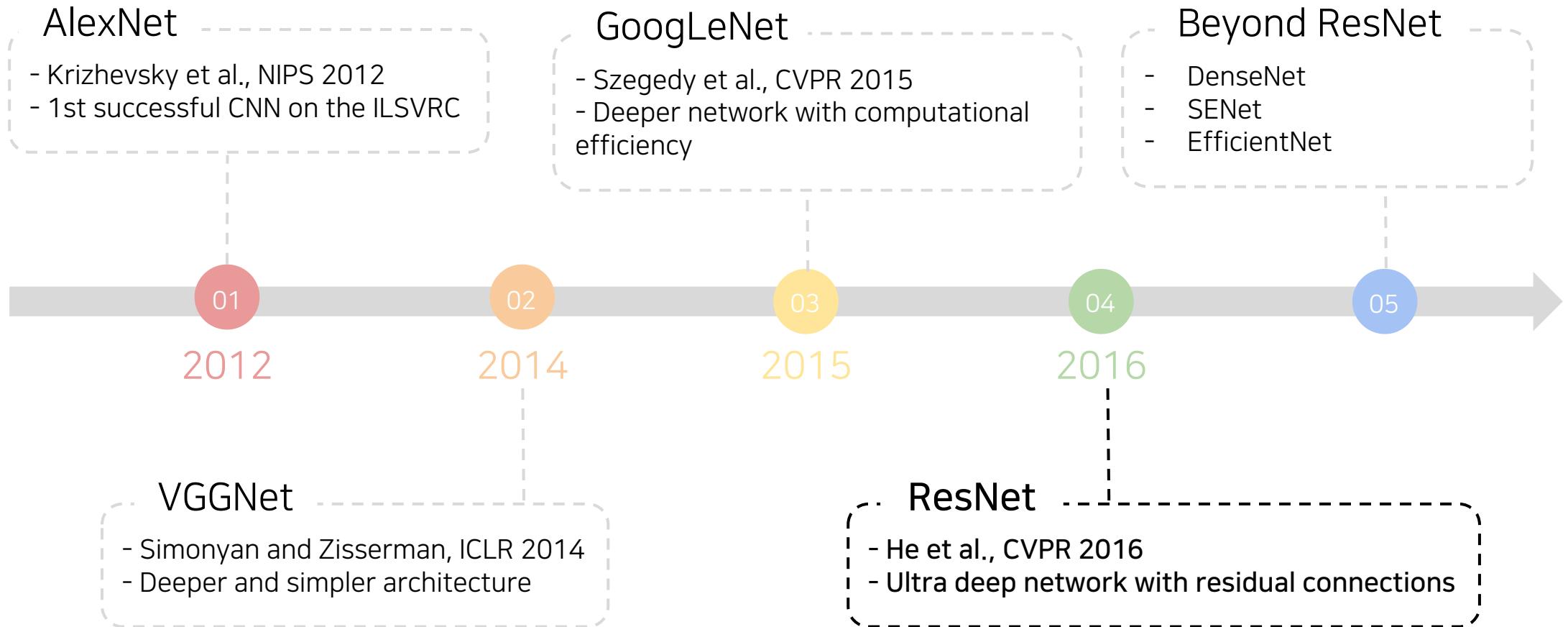
[Szegedy et al., CVPR 2015]

- The vanishing gradient problem is dealt with by the auxiliary classifier
- Injecting additional gradients into lower layers
- Used only during training, removed at testing time



## 2.2 ResNet

[He et al., CVPR 2015]

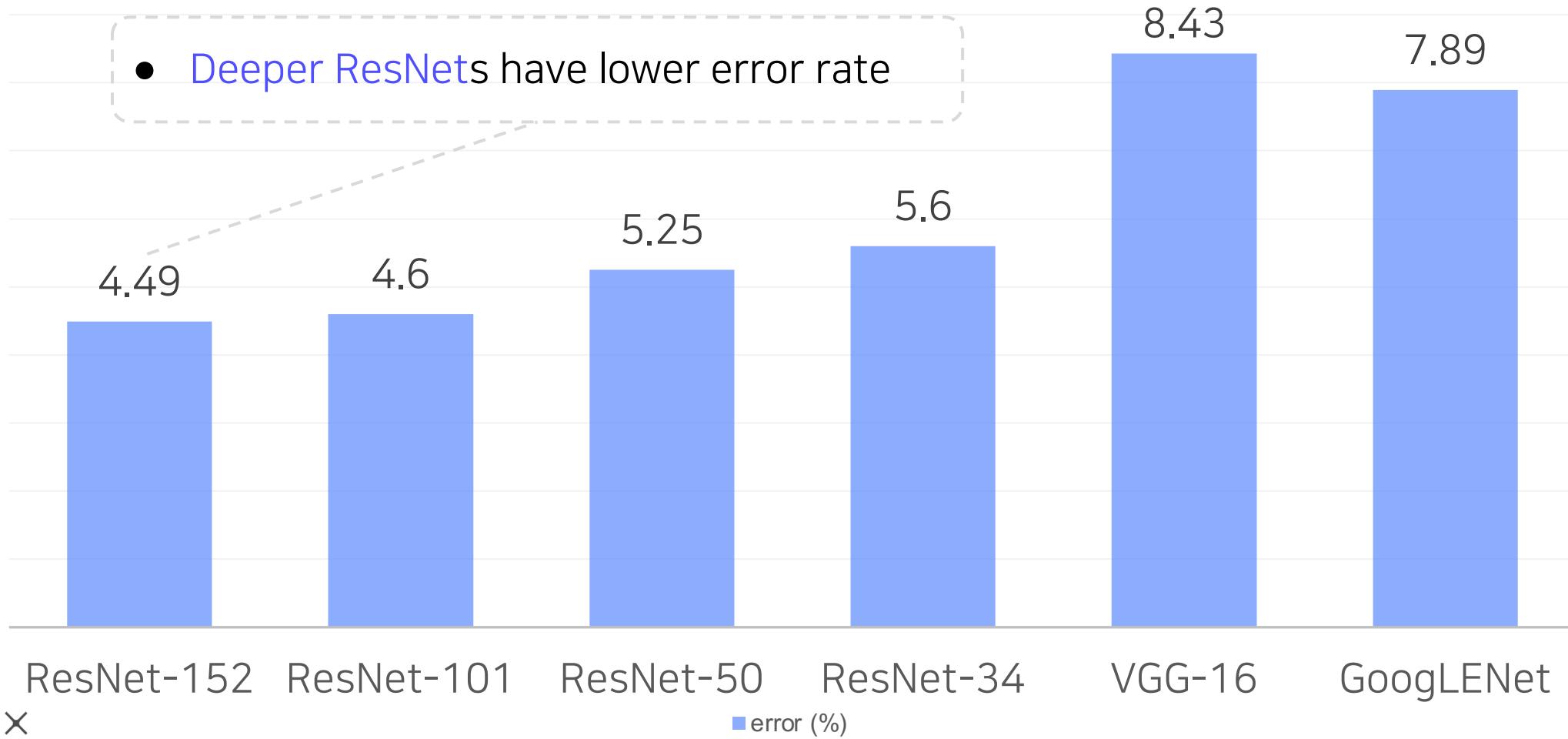


## 2.2 ResNet

CNN architectures for image classification 2

Top-5 validation error on the ImageNet

[He et al., CVPR 2016]

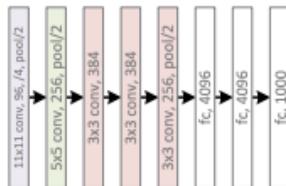


## 2.2 ResNet

# Revolutions of depth

[He et al., CVPR 2016]

- 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)



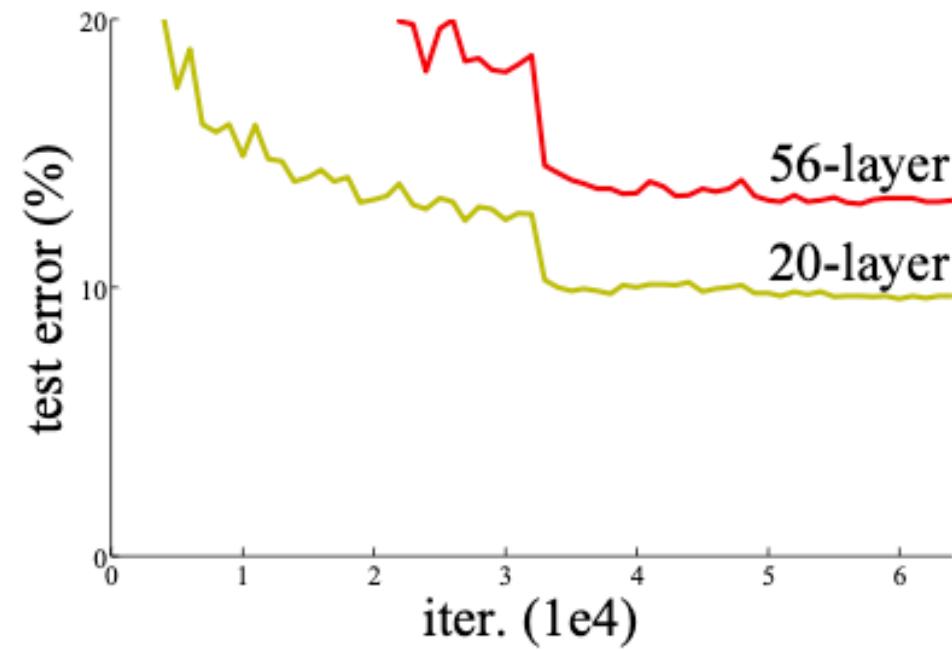
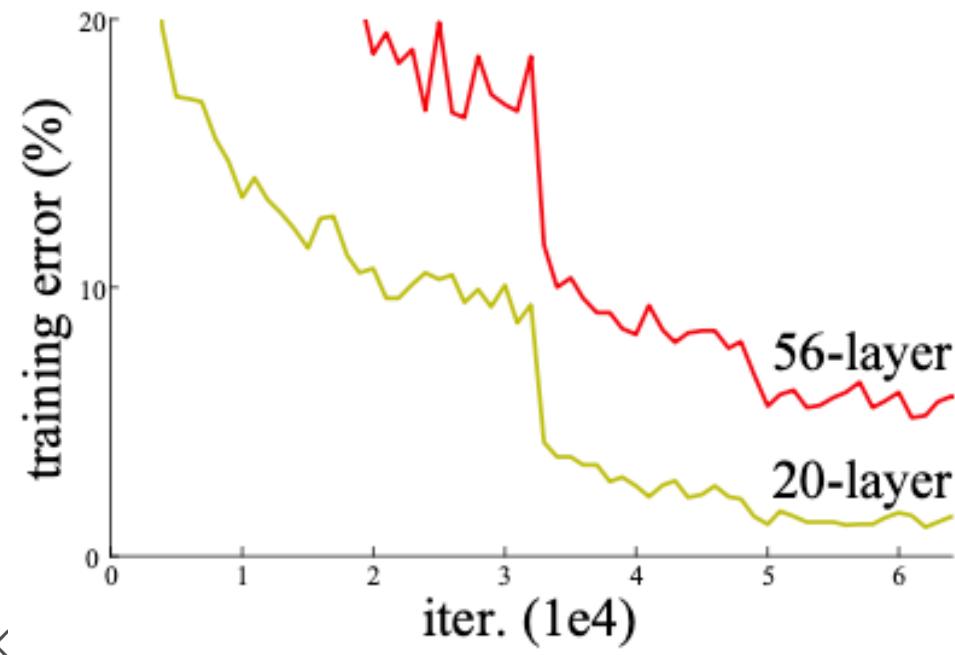
## 2.2 ResNet

CNN architectures for image classification 2

Degradation problem

[He et al., CVPR 2016]

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



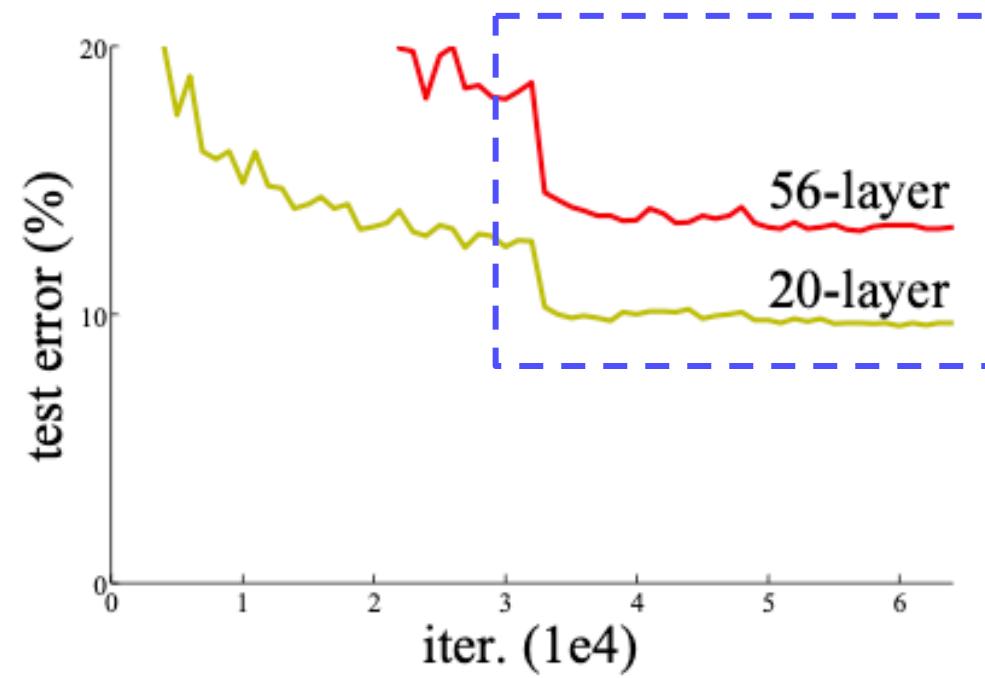
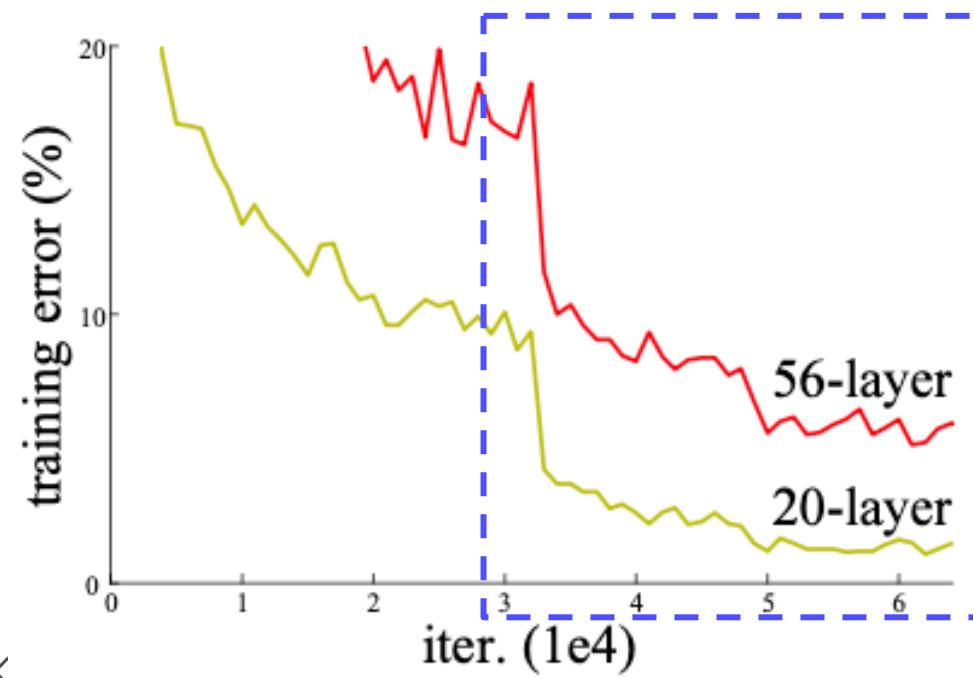
## 2.2 ResNet

CNN architectures for image classification 2

Degradation problem

[He et al., CVPR 2016]

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



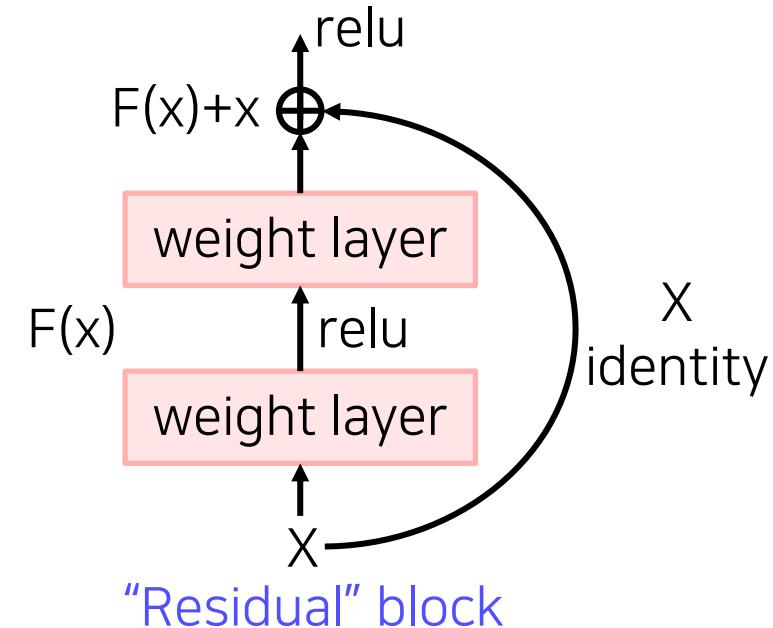
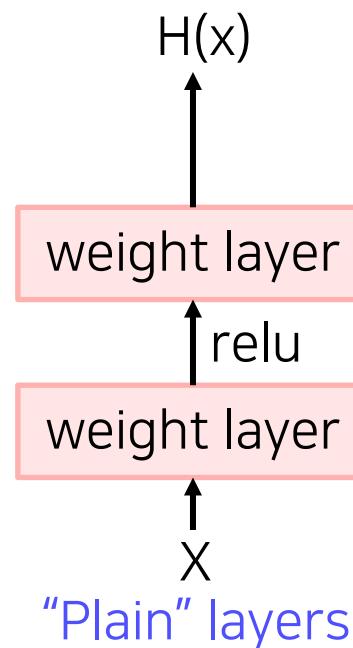
## 2.2 ResNet

CNN architectures for image classification 2

### Hypothesis

[He et al., CVPR 2016]

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



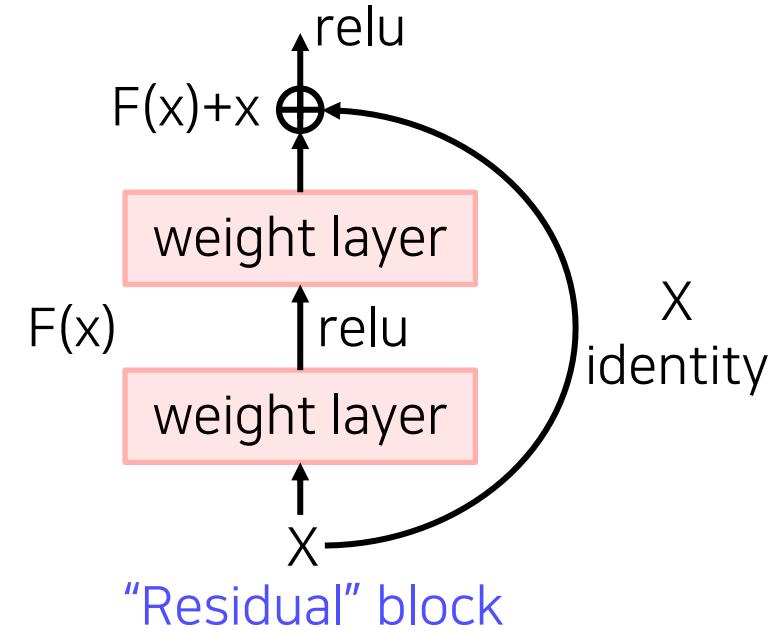
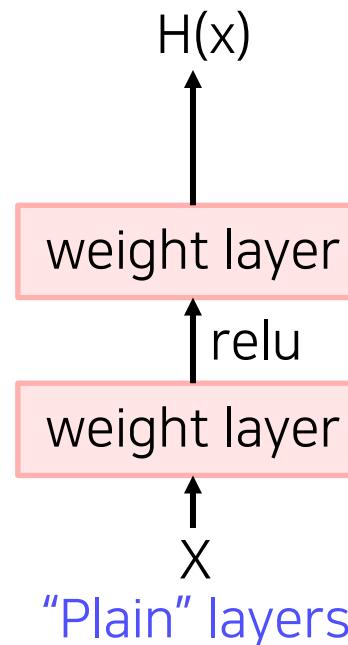
## 2.2 ResNet

CNN architectures for image classification 2

### Hypothesis

[He et al., CVPR 2016]

- 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 :  $F(x) = H(x) - x$



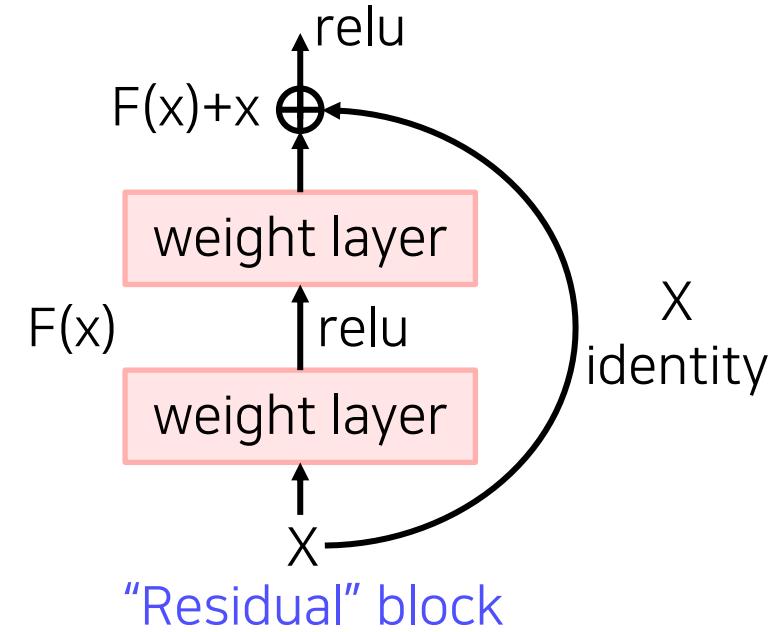
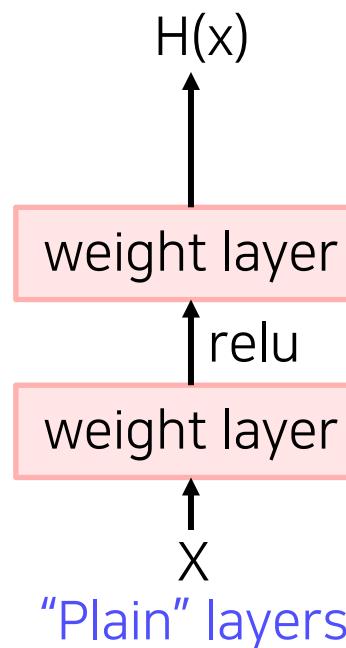
## 2.2 ResNet

CNN architectures for image classification 2

A solution: Shortcut connection

[He et al., CVPR 2016]

- 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!**



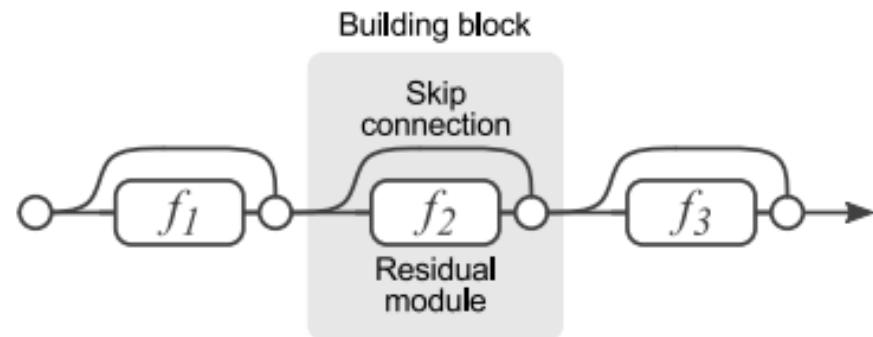
## 2.2 ResNet

CNN architectures for image classification 2

### Analysis of residual connection

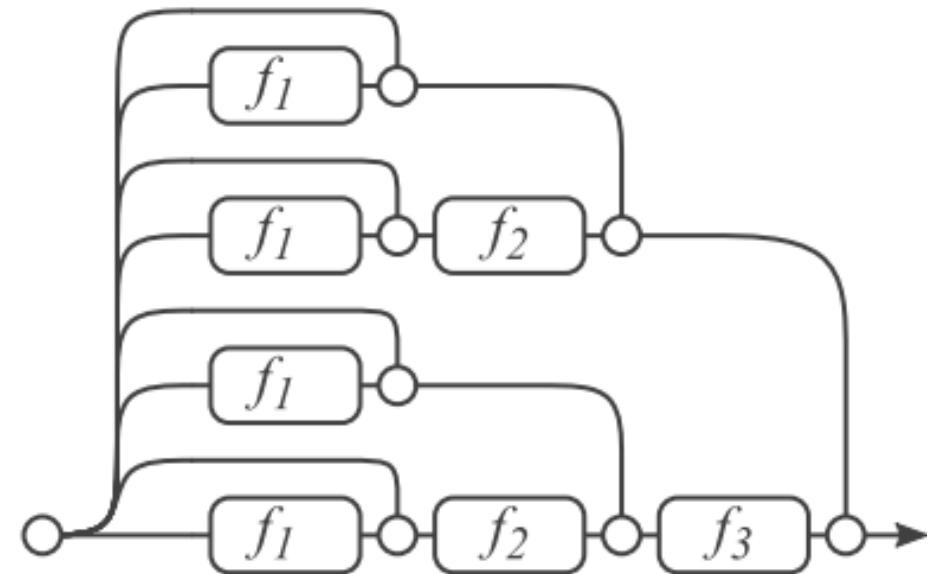
[Veit et al., NIPS 2016]

- 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)

## 2.2 ResNet

CNN architectures for image classification 2

### Overall architecture

[He et al., CVPR 2016]

- He initialization
- Additional conv layer at the beginning



## 2.2 ResNet

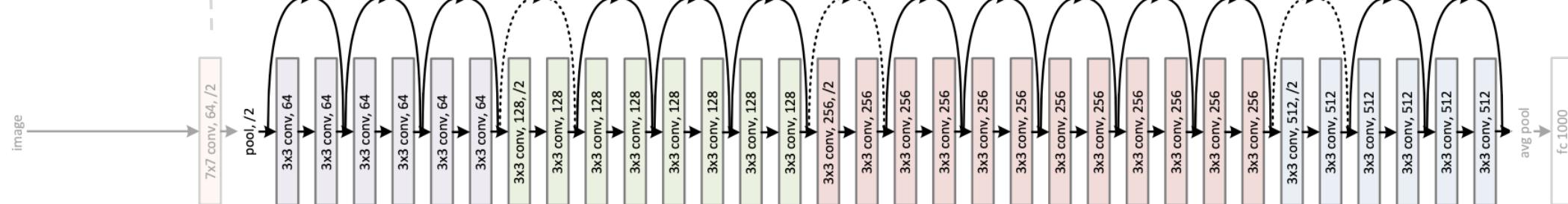
CNN architectures for image classification 2

### Overall architecture

[He et al., CVPR 2016]

- 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



## 2.2 ResNet

CNN architectures for image classification 2

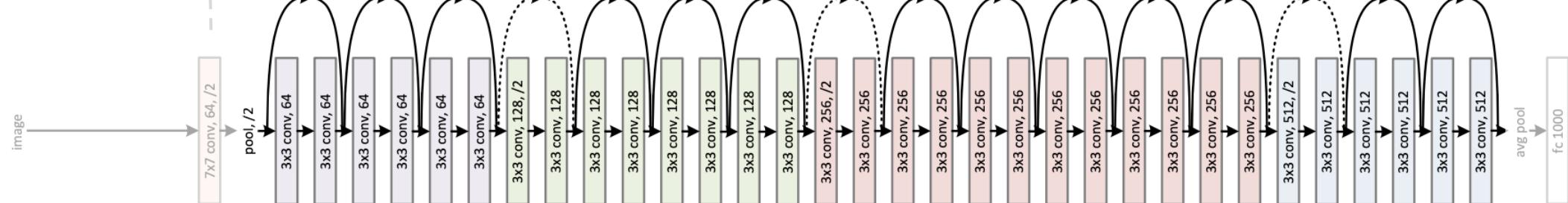
### Overall architecture

[He et al., CVPR 2016]

- 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



## 2.2 ResNet

CNN architectures for image classification 2

### Overall architecture

[He et al., CVPR 2016]

- 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



## 2.2 ResNet

CNN architectures for image classification 2

PyTorch code for ResNet

[He et al., CVPR 2016]

| 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       | $\begin{bmatrix} 3\times3, 64 \\ 3\times3, 64 \end{bmatrix} \times 2$   | $\begin{bmatrix} 3\times3, 64 \\ 3\times3, 64 \end{bmatrix} \times 3$                     | $\begin{bmatrix} 1\times1, 64 \\ 3\times3, 64 \\ 1\times1, 256 \end{bmatrix} \times 3$    | $\begin{bmatrix} 1\times1, 64 \\ 3\times3, 64 \\ 1\times1, 256 \end{bmatrix} \times 3$     | $\begin{bmatrix} 1\times1, 64 \\ 3\times3, 64 \\ 1\times1, 256 \end{bmatrix} \times 3$     |
| conv3_x    | 28×28       | $\begin{bmatrix} 3\times3, 128 \\ 3\times3, 128 \end{bmatrix} \times 2$ | $\begin{bmatrix} 3\times3, 128 \\ 3\times3, 128 \\ 1\times1, 512 \end{bmatrix} \times 4$  | $\begin{bmatrix} 1\times1, 128 \\ 3\times3, 128 \\ 1\times1, 512 \end{bmatrix} \times 4$  | $\begin{bmatrix} 1\times1, 128 \\ 3\times3, 128 \\ 1\times1, 512 \end{bmatrix} \times 4$   | $\begin{bmatrix} 1\times1, 128 \\ 3\times3, 128 \\ 1\times1, 512 \end{bmatrix} \times 8$   |
| conv4_x    | 14×14       | $\begin{bmatrix} 3\times3, 256 \\ 3\times3, 256 \end{bmatrix} \times 2$ | $\begin{bmatrix} 3\times3, 256 \\ 3\times3, 256 \\ 1\times1, 1024 \end{bmatrix} \times 6$ | $\begin{bmatrix} 1\times1, 256 \\ 3\times3, 256 \\ 1\times1, 1024 \end{bmatrix} \times 6$ | $\begin{bmatrix} 1\times1, 256 \\ 3\times3, 256 \\ 1\times1, 1024 \end{bmatrix} \times 23$ | $\begin{bmatrix} 1\times1, 256 \\ 3\times3, 256 \\ 1\times1, 1024 \end{bmatrix} \times 36$ |
| conv5_x    | 7×7         | $\begin{bmatrix} 3\times3, 512 \\ 3\times3, 512 \end{bmatrix} \times 2$ | $\begin{bmatrix} 3\times3, 512 \\ 3\times3, 512 \\ 1\times1, 2048 \end{bmatrix} \times 3$ | $\begin{bmatrix} 1\times1, 512 \\ 3\times3, 512 \\ 1\times1, 2048 \end{bmatrix} \times 3$ | $\begin{bmatrix} 1\times1, 512 \\ 3\times3, 512 \\ 1\times1, 2048 \end{bmatrix} \times 3$  | $\begin{bmatrix} 1\times1, 512 \\ 3\times3, 512 \\ 1\times1, 2048 \end{bmatrix} \times 3$  |
|            | 1×1         |   |   | average pool, 1000-d fc, softmax  |  |  |
| FLOPs      |             | $1.8 \times 10^9$   | $3.6 \times 10^9$   | $3.8 \times 10^9$   | $7.6 \times 10^9$  | $11.3 \times 10^9$   |

×

## 2.2 ResNet

CNN architectures for image classification 2

PyTorch code for ResNet

[He et al., CVPR 2016]

| 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       | $\begin{bmatrix} 3\times3, 64 \\ 3\times3, 64 \end{bmatrix} \times 2$   | $\begin{bmatrix} 3\times3, 64 \\ 3\times3, 64 \end{bmatrix} \times 3$                     | $\begin{bmatrix} 1\times1, 64 \\ 3\times3, 64 \\ 1\times1, 256 \end{bmatrix} \times 3$    | $\begin{bmatrix} 1\times1, 64 \\ 3\times3, 64 \\ 1\times1, 256 \end{bmatrix} \times 3$    | $\begin{bmatrix} 1\times1, 64 \\ 3\times3, 64 \\ 1\times1, 256 \end{bmatrix} \times 3$     |
| conv3_x    | 28×28       | $\begin{bmatrix} 3\times3, 128 \\ 3\times3, 128 \end{bmatrix} \times 2$ | $\begin{bmatrix} 3\times3, 128 \\ 3\times3, 128 \\ 1\times1, 512 \end{bmatrix} \times 4$  | $\begin{bmatrix} 1\times1, 128 \\ 3\times3, 128 \\ 1\times1, 512 \end{bmatrix} \times 4$  | $\begin{bmatrix} 1\times1, 128 \\ 3\times3, 128 \\ 1\times1, 512 \end{bmatrix} \times 4$  | $\begin{bmatrix} 1\times1, 128 \\ 3\times3, 128 \\ 1\times1, 512 \end{bmatrix} \times 8$   |
| conv4_x    | 14×14       | $\begin{bmatrix} 3\times3, 256 \\ 3\times3, 256 \end{bmatrix} \times 2$ | $\begin{bmatrix} 3\times3, 256 \\ 3\times3, 256 \\ 1\times1, 1024 \end{bmatrix} \times 6$ | $\begin{bmatrix} 1\times1, 256 \\ 3\times3, 256 \\ 1\times1, 1024 \end{bmatrix} \times 6$ | $\begin{bmatrix} 3\times3, 256 \\ 1\times1, 1024 \end{bmatrix} \times 23$                 | $\begin{bmatrix} 1\times1, 256 \\ 3\times3, 256 \\ 1\times1, 1024 \end{bmatrix} \times 36$ |
| conv5_x    | 7×7         | $\begin{bmatrix} 3\times3, 512 \\ 3\times3, 512 \end{bmatrix} \times 2$ | $\begin{bmatrix} 3\times3, 512 \\ 3\times3, 512 \\ 1\times1, 2048 \end{bmatrix} \times 3$ | $\begin{bmatrix} 1\times1, 512 \\ 3\times3, 512 \\ 1\times1, 2048 \end{bmatrix} \times 3$ | $\begin{bmatrix} 1\times1, 512 \\ 3\times3, 512 \\ 1\times1, 2048 \end{bmatrix} \times 3$ | $\begin{bmatrix} 1\times1, 512 \\ 3\times3, 512 \\ 1\times1, 2048 \end{bmatrix} \times 3$  |
|            | 1×1         | average pool, 1000-d fc, softmax  |   |   |   |  |
| FLOPs      |             | $1.8 \times 10^9$   | $3.6 \times 10^9$   | $3.8 \times 10^9$   | $7.6 \times 10^9$   | $11.3 \times 10^9$   |

×

## 2.2 ResNet

CNN architectures for image classification 2

PyTorch code for ResNet

[He et al., CVPR 2016]

| 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       | $\begin{bmatrix} 3\times3, 64 \\ 3\times3, 64 \end{bmatrix} \times 2$   | <pre>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)</pre> | $\begin{bmatrix} 1\times1, 64 \\ 3\times3, 64 \\ 1\times1, 256 \end{bmatrix} \times 3$     | $\begin{bmatrix} 1\times1, 64 \\ 3\times3, 64 \\ 1\times1, 256 \end{bmatrix} \times 3$     | $\begin{bmatrix} 1\times1, 64 \\ 3\times3, 64 \\ 1\times1, 256 \end{bmatrix} \times 3$     |
| conv3_x    | 28×28       | $\begin{bmatrix} 3\times3, 128 \\ 3\times3, 128 \end{bmatrix} \times 2$ |  | $\begin{bmatrix} 1\times1, 128 \\ 3\times3, 128 \\ 1\times1, 512 \end{bmatrix} \times 4$   | $\begin{bmatrix} 1\times1, 128 \\ 3\times3, 128 \\ 1\times1, 512 \end{bmatrix} \times 8$   | $\begin{bmatrix} 1\times1, 128 \\ 3\times3, 128 \\ 1\times1, 512 \end{bmatrix} \times 8$   |
| conv4_x    | 14×14       | $\begin{bmatrix} 3\times3, 256 \\ 3\times3, 256 \end{bmatrix} \times 2$ | <pre>x = self.layer1(x) x = self.layer2(x) x = self.layer3(x) x = self.layer4(x)</pre>   | $\begin{bmatrix} 1\times1, 256 \\ 3\times3, 256 \\ 1\times1, 1024 \end{bmatrix} \times 23$ | $\begin{bmatrix} 1\times1, 256 \\ 3\times3, 256 \\ 1\times1, 1024 \end{bmatrix} \times 36$ | $\begin{bmatrix} 1\times1, 256 \\ 3\times3, 256 \\ 1\times1, 1024 \end{bmatrix} \times 36$ |
| conv5_x    | 7×7         | $\begin{bmatrix} 3\times3, 512 \\ 3\times3, 512 \end{bmatrix} \times 2$ | <pre>x = self.avgpool(x) x = torch.flatten(x, 1) x = self.fc(x)</pre>  | $\begin{bmatrix} 1\times1, 512 \\ 3\times3, 512 \\ 1\times1, 2048 \end{bmatrix} \times 3$  | $\begin{bmatrix} 1\times1, 512 \\ 3\times3, 512 \\ 1\times1, 2048 \end{bmatrix} \times 3$  | $\begin{bmatrix} 1\times1, 512 \\ 3\times3, 512 \\ 1\times1, 2048 \end{bmatrix} \times 3$  |
|            | 1×1         |   |  | softmax  |  |  |
| FLOPs      |             | $1.8 \times 10^9$   |  |  | $7.6 \times 10^9$  | $11.3 \times 10^9$   |

X

## 2.2 ResNet

CNN architectures for image classification 2

PyTorch code for ResNet

[He et al., CVPR 2016]

| 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       | $\begin{bmatrix} 3\times3, 64 \\ 3\times3, 64 \end{bmatrix} \times 2$   |          | 3×3 max pool, stride 2 |  |  |
| conv3_x    | 28×28       | $\begin{bmatrix} 3\times3, 128 \\ 3\times3, 128 \end{bmatrix} \times 2$ |          |                        | $\begin{bmatrix} 1\times1, 64 \\ 3\times3, 64 \\ 1\times1, 256 \end{bmatrix} \times 3$ | $\begin{bmatrix} 1\times1, 64 \\ 3\times3, 64 \\ 1\times1, 256 \end{bmatrix} \times 3$ |
| conv4_x    | 14×14       | $\begin{bmatrix} 3\times3, 256 \\ 3\times3, 256 \end{bmatrix} \times 2$ |          |                        | $\begin{bmatrix} 1\times1, 128 \end{bmatrix}$  | $\begin{bmatrix} 1\times1, 128 \end{bmatrix}$  |
| conv5_x    | 7×7         | $\begin{bmatrix} 3\times3, 512 \\ 3\times3, 512 \end{bmatrix} \times 2$ |          |                        |  |  |
|            | 1×1         |   |          |                        |  |  |
| FLOPs      |             | $1.8 \times 10^9$   |          |                        | $7.6 \times 10^9$  | $11.3 \times 10^9$   |

```
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
```

## 2.2 ResNet

CNN architectures for image classification 2

PyTorch code for ResNet

[He et al., CVPR 2016]

| 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       | $\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 2$   |          | 3×3 max pool, stride 2 |           |           |
| conv3_x    | 28×28       | $\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 2$ |          |                        |           |           |
| conv4_x    | 14×14       | $\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 2$ |          |                        |           |           |
| conv5_x    | 7×7         | $\begin{bmatrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{bmatrix} \times 2$ |          |                        |           |           |
|            | 1×1         |   |          |                        |           |           |
| FLOPs      |             | $1.8 \times 10^9$   |          |                        |           |           |

```
def _make_layer(self, block: Type[Union[BasicBlock, Bottleneck]], planes: int, blocks: int,
                stride: int = 1, dilate: bool = False) -> nn.Sequential:
    norm_layer = self._norm_layer
    downsample = None
    previous_dilation = self.dilation
    if dilate:
        self.dilation *= stride
        stride = 1
    if stride != 1 or self.inplanes != planes * block.expansion:
        downsample = nn.Sequential(
            conv1x1(self.inplanes, planes * block.expansion, stride),
            norm_layer(planes * block.expansion),
        )
    layers = []
    layers.append(block(self.inplanes, planes, stride, downsample, self.groups,
                        self.base_width, previous_dilation, norm_layer))
    self.inplanes = planes * block.expansion
    for _ in range(1, blocks):
        layers.append(block(self.inplanes, planes, groups=self.groups,
                            base_width=self.base_width, dilation=self.dilation,
                            norm_layer=norm_layer))

    return nn.Sequential(*layers)
```

X

## 2.2 ResNet

CNN architectures for image classification 2

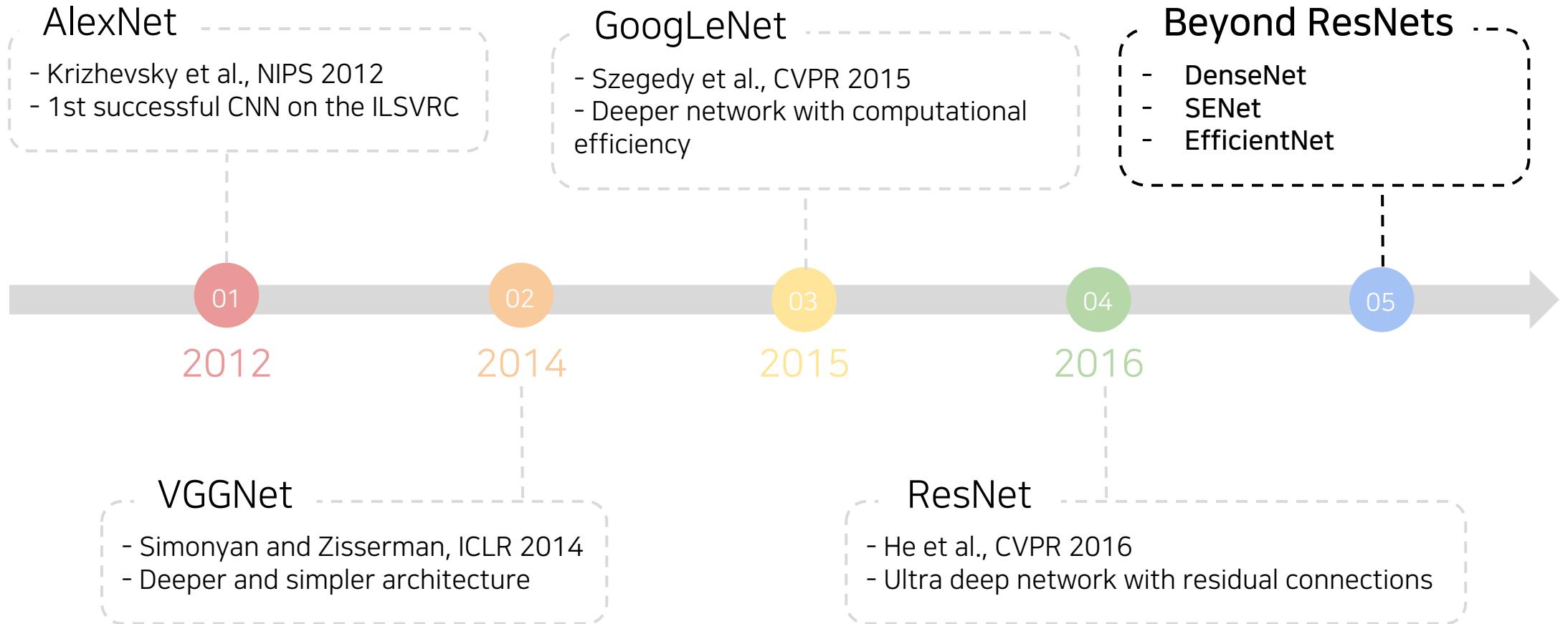
PyTorch code for ResNet

[He et al., CVPR 2016]

| 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       | $\left[ \begin{array}{l} 3 \times 3, 64 \\ 3 \times 3, 64 \end{array} \right] \times 2$   | $\left[ \begin{array}{l} 3 \times 3, 64 \\ 3 \times 3, 64 \end{array} \right] \times 3$   | $\left[ \begin{array}{l} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{array} \right] \times 3$    | $\left[ \begin{array}{l} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{array} \right] \times 4$  | $\left[ \begin{array}{l} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{array} \right] \times 4$  |
| conv4_x    | 14×14       | $\left[ \begin{array}{l} 3 \times 3, 256 \\ 3 \times 3, 256 \end{array} \right] \times 2$ | $\left[ \begin{array}{l} 3 \times 3, 256 \\ 3 \times 3, 256 \end{array} \right] \times 6$ | $\left[ \begin{array}{l} 3 \times 3, 256 \\ 1 \times 1, 1024 \end{array} \right] \times 6$                    | $\left[ \begin{array}{l} 3 \times 3, 256 \\ 1 \times 1, 1024 \end{array} \right] \times 6$                    | $\left[ \begin{array}{l} 3 \times 3, 256 \\ 1 \times 1, 1024 \end{array} \right] \times 6$                    |
| conv5_x    | 7×7         | $\left[ \begin{array}{l} 3 \times 3, 512 \\ 3 \times 3, 512 \end{array} \right] \times 2$ | $\left[ \begin{array}{l} 3 \times 3, 512 \\ 3 \times 3, 512 \end{array} \right] \times 3$ | $\left[ \begin{array}{l} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{array} \right] \times 3$ | $\left[ \begin{array}{l} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{array} \right] \times 3$ | $\left[ \begin{array}{l} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{array} \right] \times 3$ |
|            | 1×1         |   |   | average pool, 1000-d fc, softmax  |   |   |
| FLOPs      |             | $1.8 \times 10^9$   | $3.6 \times 10^9$   | $3.8 \times 10^9$   | $7.6 \times 10^9$   |   |

## 2.3 Beyond ResNets

CNN architectures for image classification 2



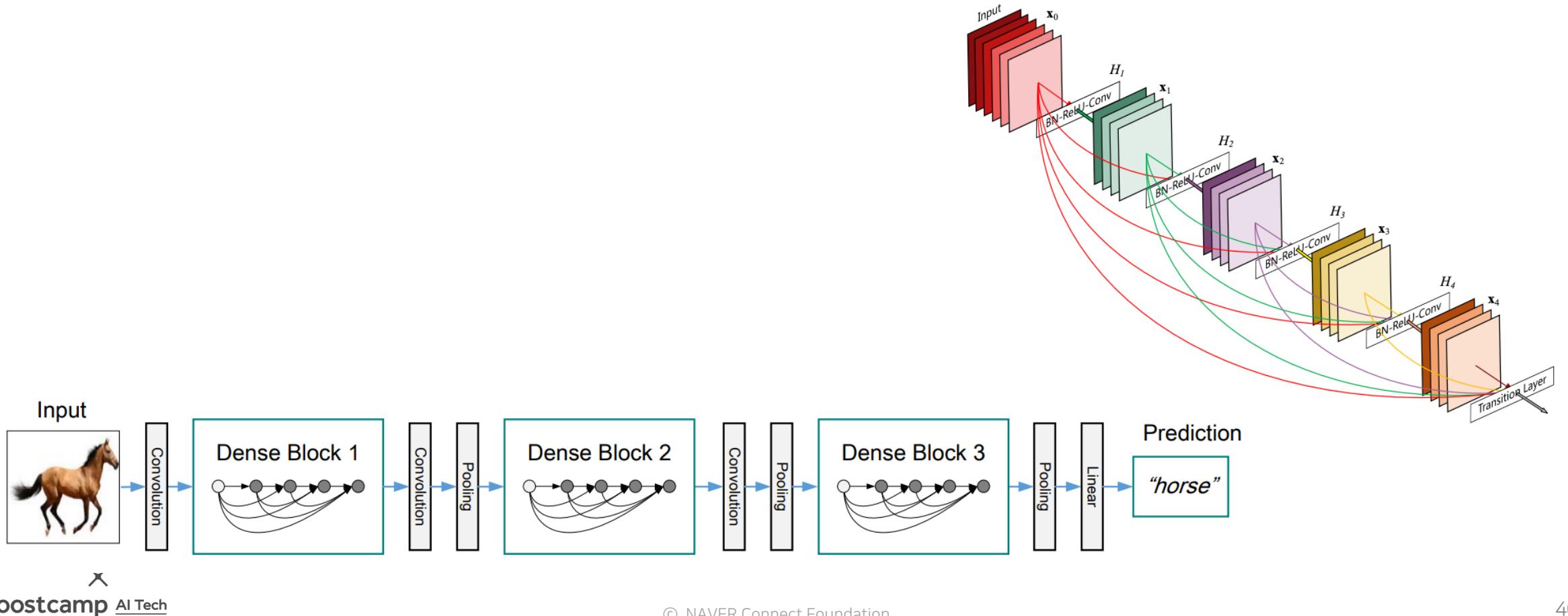
## 2.3 Beyond ResNets

CNN architectures for image classification 2

### DenseNet

[Huang et al., CVPR 2017]

- In ResNet, we added the input and the output of the layer element-wisely



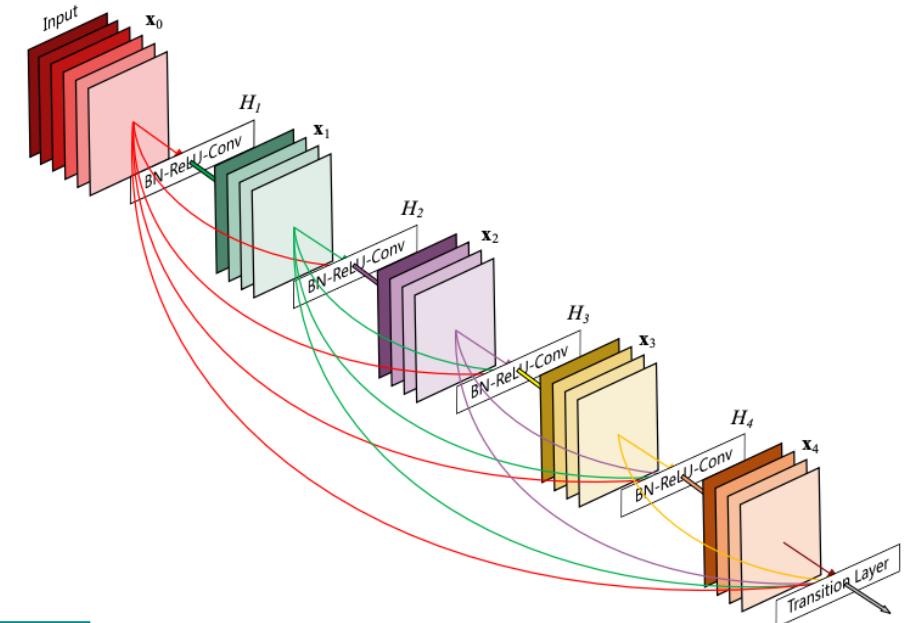
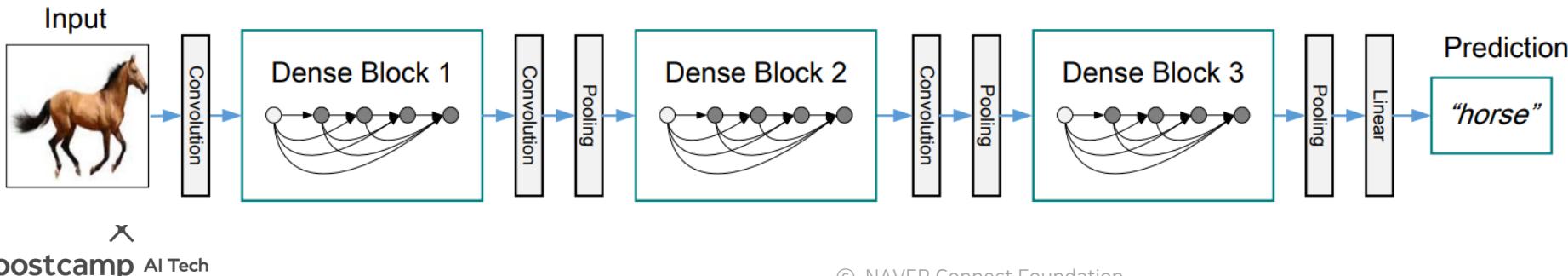
## 2.3 Beyond ResNets

CNN architectures for image classification 2

### DenseNet

[Huang et al., CVPR 2017]

- In ResNet, we added the input and the output of the layer element-wisely
- In the **Dense blocks**, every output of each layer is concatenated along the **channel axis**.
  - Alleviate vanishing gradient problem
  - Strengthen feature propagation
  - Encourage the reuse of features



## 2.3 Beyond ResNets

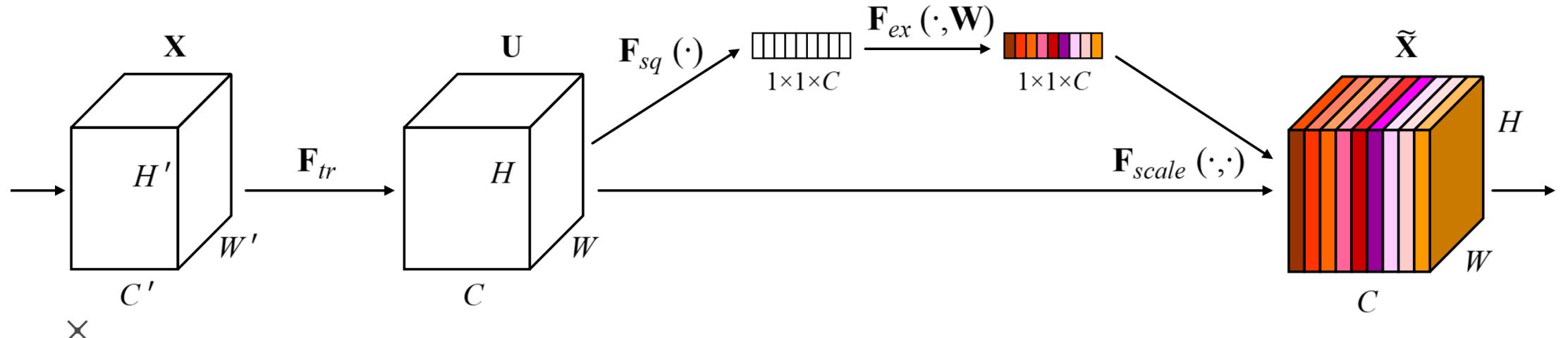
CNN architectures for image classification 2

SENet

[Hu et al., CVPR 2018]

\* SE : Squeeze and Excitation

- Attention across channels
- Recalibrates channel-wise responses by modeling interdependencies between channels



## 2.3 Beyond ResNets

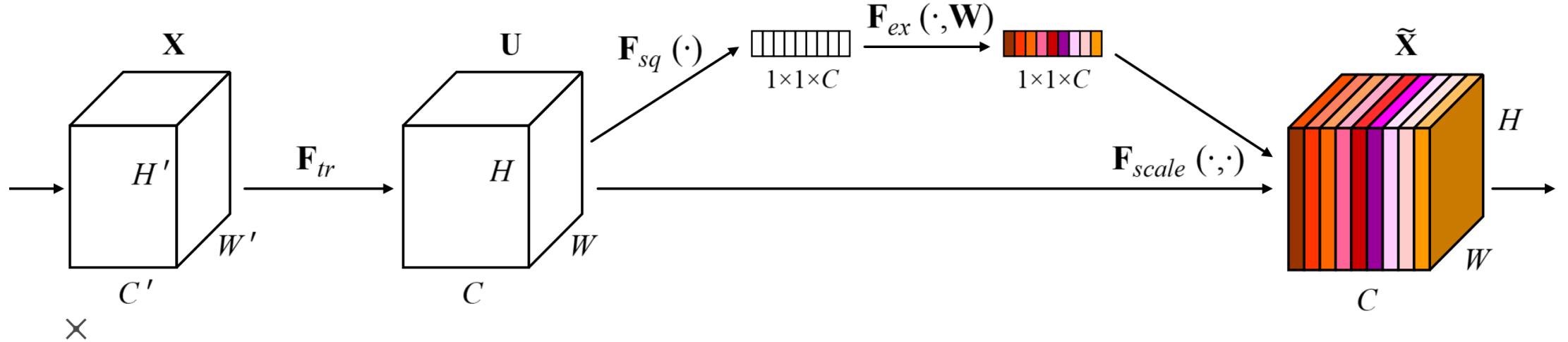
CNN architectures for image classification 2

SENet

[Hu et al., CVPR 2018]

\* SE : Squeeze and Excitation

- Attention across channels
- Recalibrates channel-wise responses by modeling interdependencies between channels
- Squeeze and excitation operations
  - **Squeeze**: capturing distributions of channel-wise responses by global average pooling
  - **Excitation**: gating channels by channel-wise attention weights obtained by a FC layer



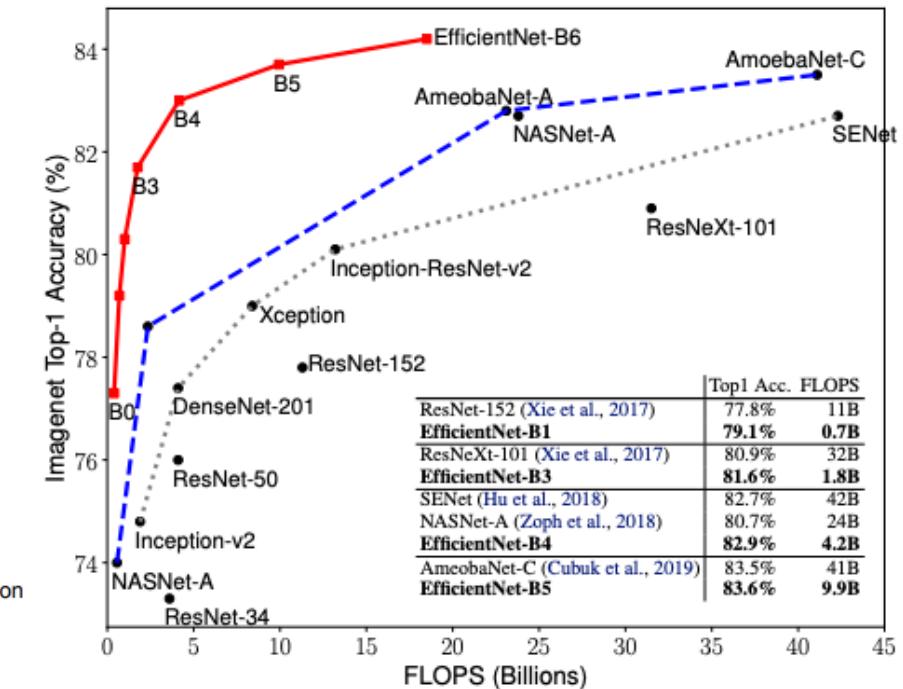
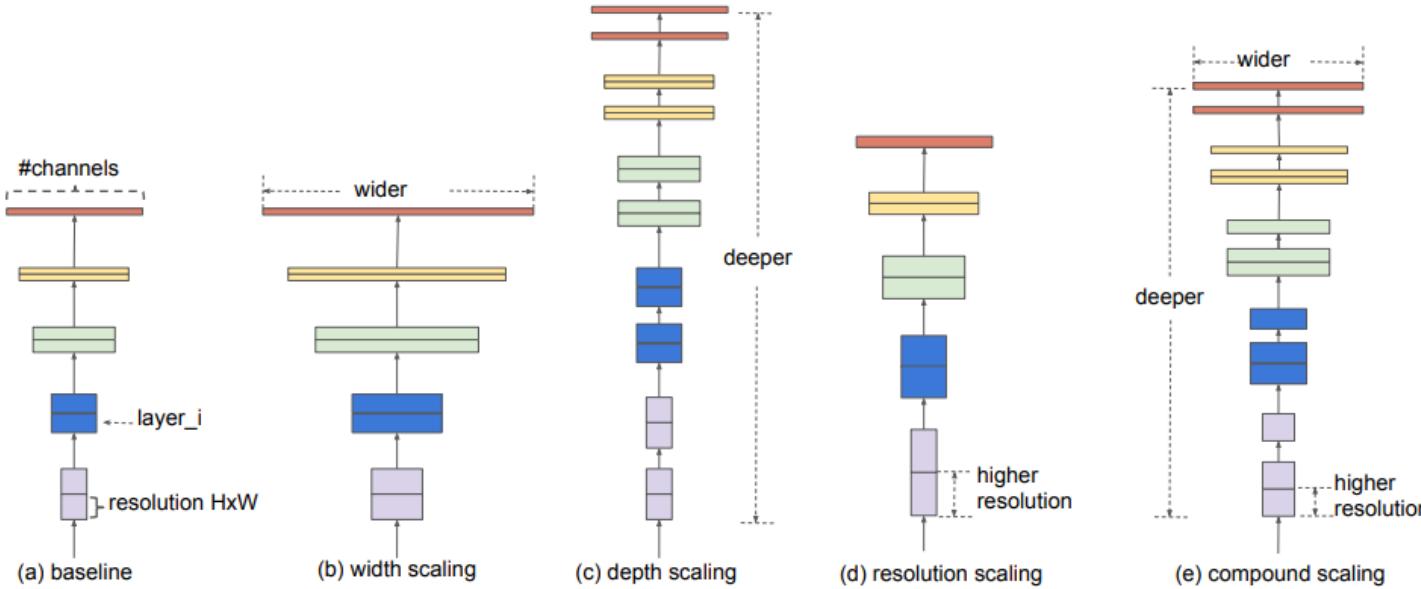
## 2.3 Beyond ResNets

CNN architectures for image classification 2

### EfficientNet

[Tan and Le, ICML 2019]

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



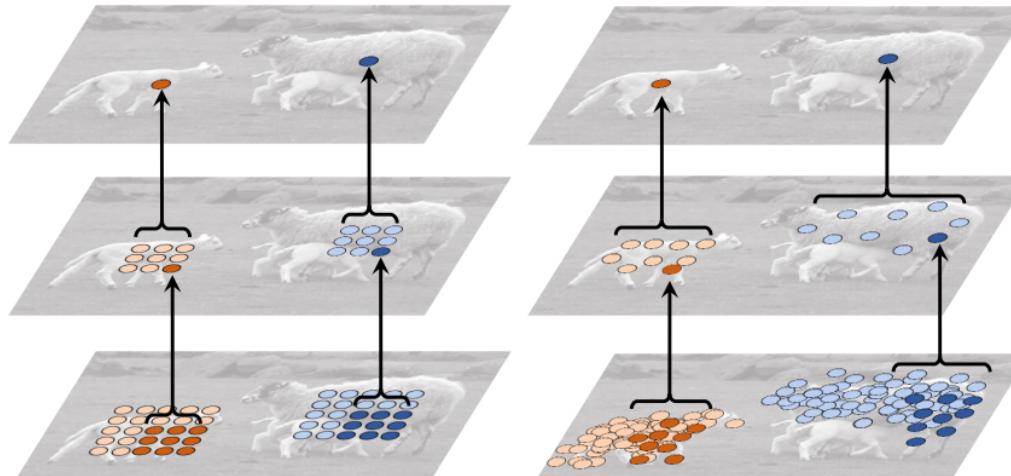
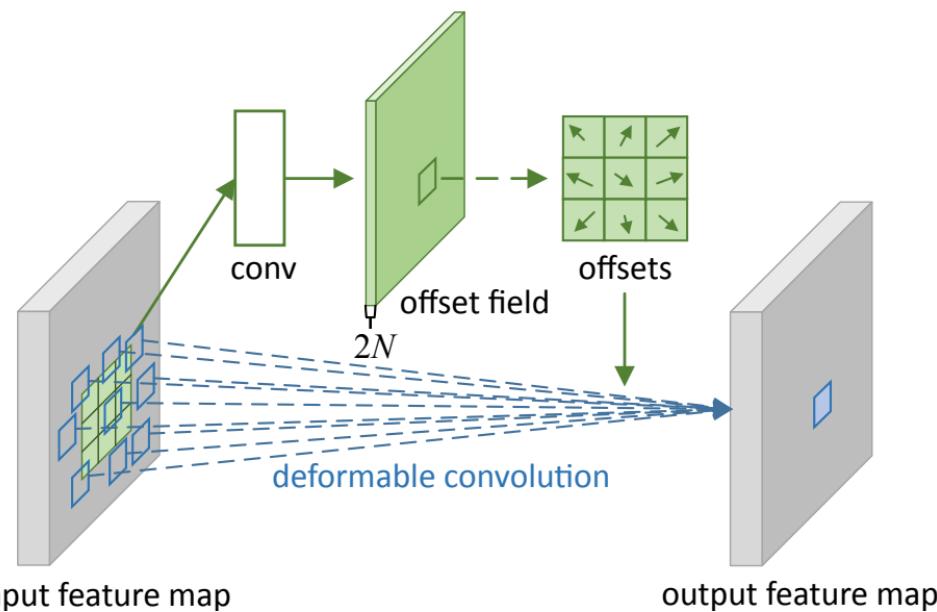
## 2.3 Beyond ResNets

CNN architectures for image classification 2

### Deformable convolution

[Dai et al., ICCV 2017]

- 2D spatial offset prediction for irregular convolution
- Irregular grid sampling with 2D spatial offsets
- Implemented by **standard CNN** and grid sampling with **2D offsets**



(a) standard convolution

(b) deformable convolution

3.

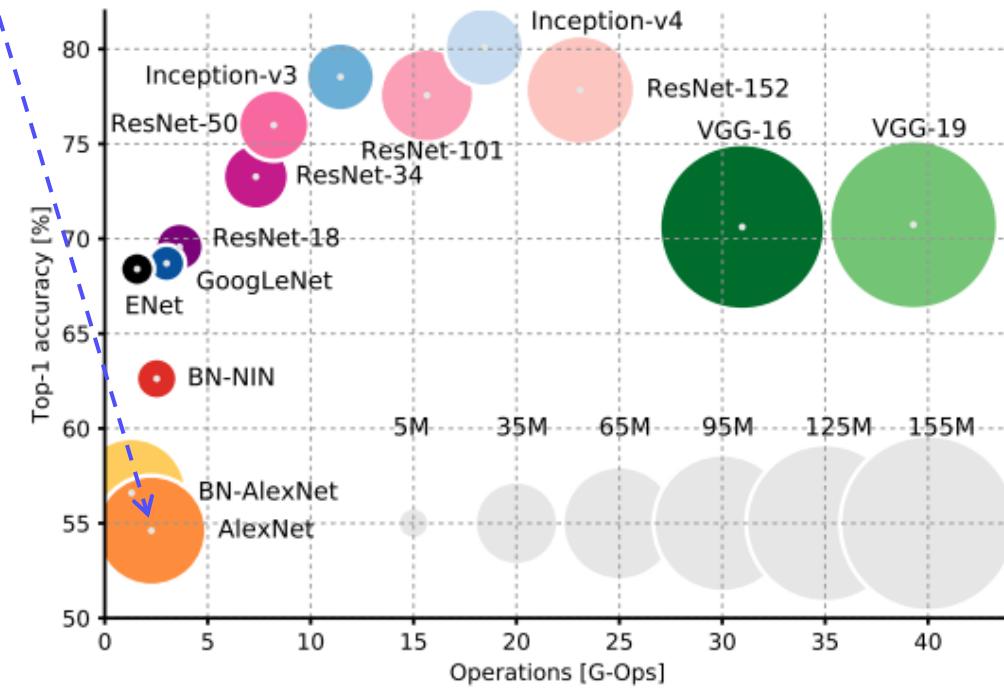
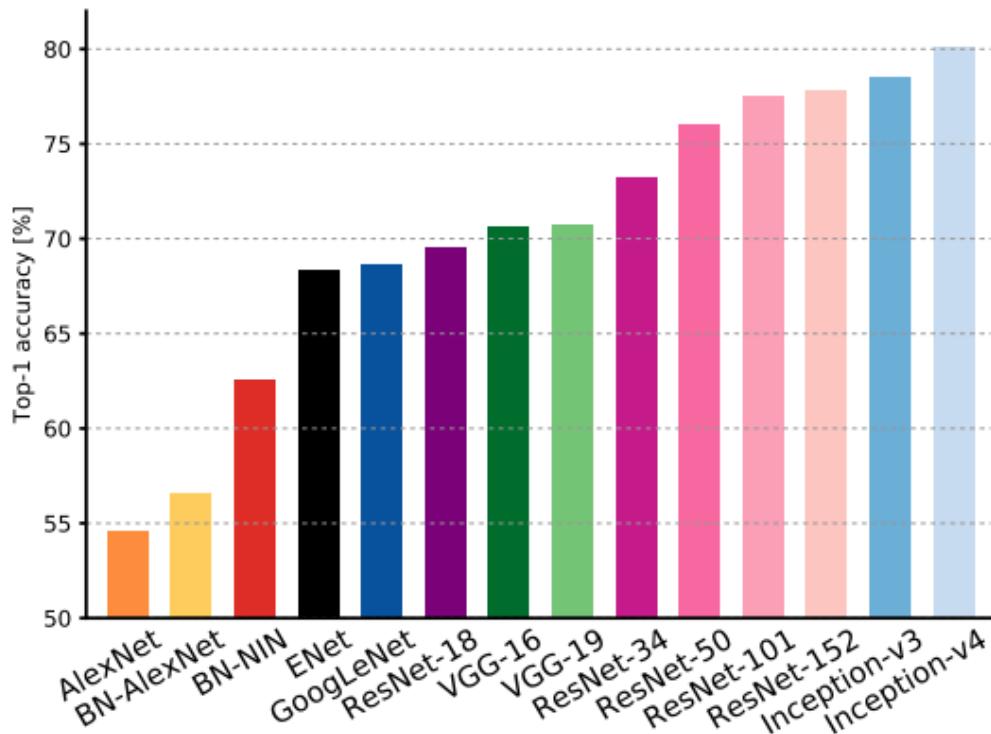
## Summary of image classification

## 3.1 Summary of image classification

Summary of image classification

[Canziani et al., CVPR 2016]

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

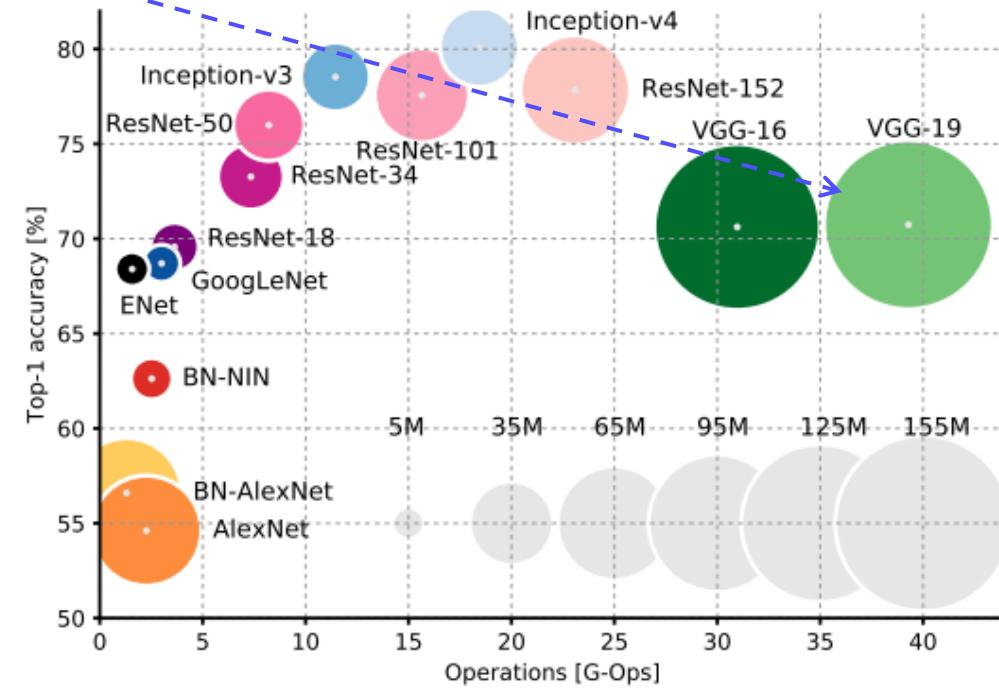
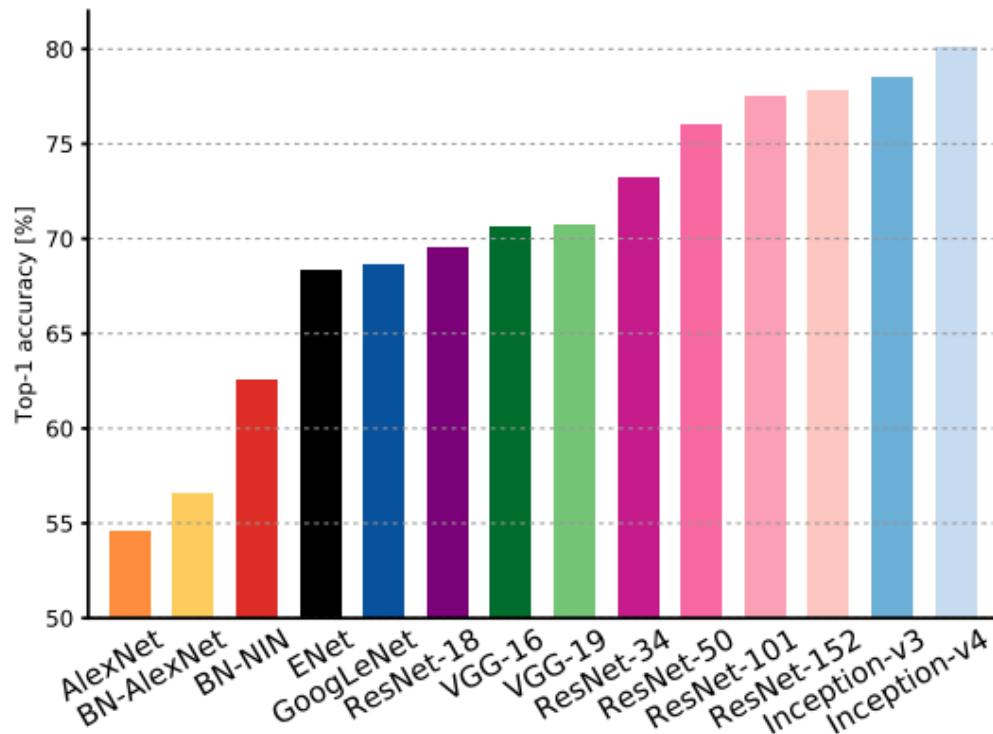


## 3.1 Summary of image classification

Summary of image classification

[Canziani et al., CVPR 2016]

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

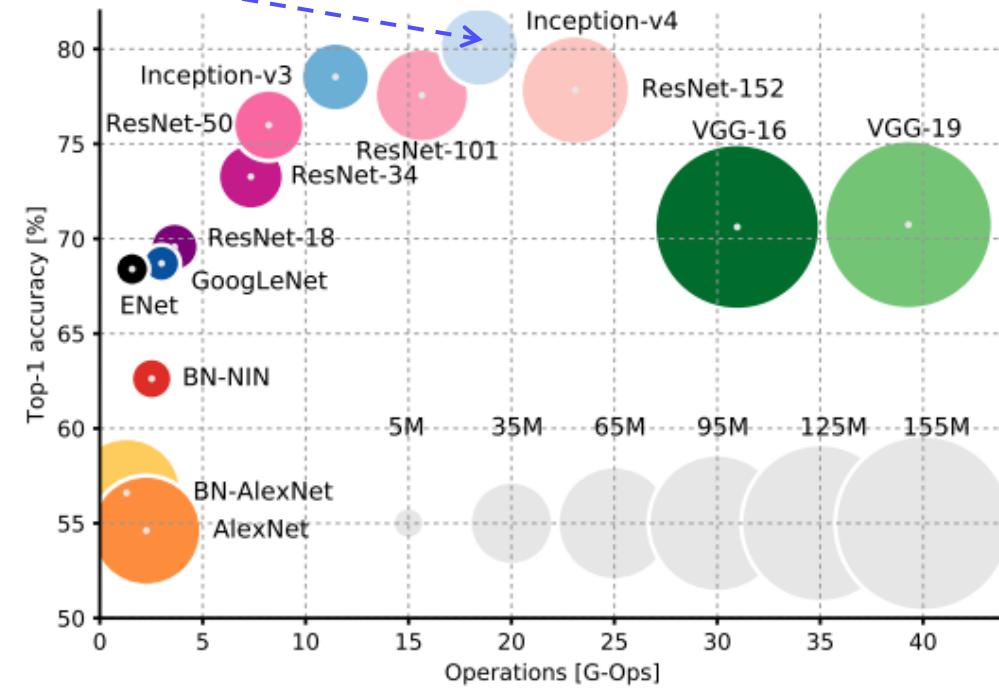
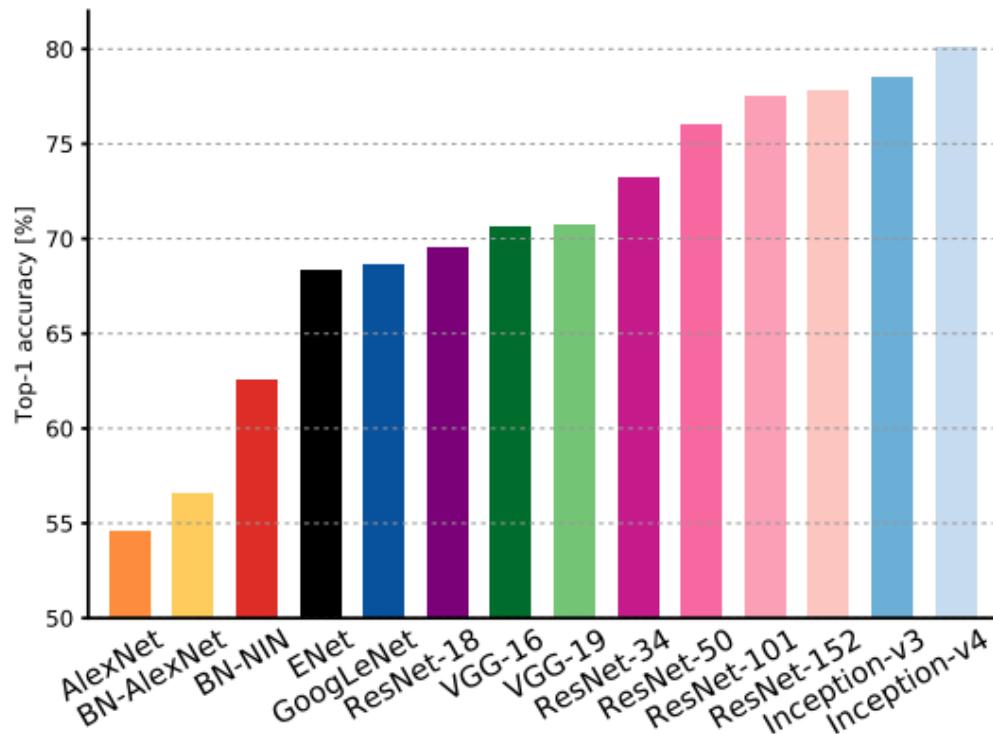


## 3.1 Summary of image classification

Summary of image classification

[Canziani et al., CVPR 2016]

- **GoogLeNet**: inception module and auxiliary classifier

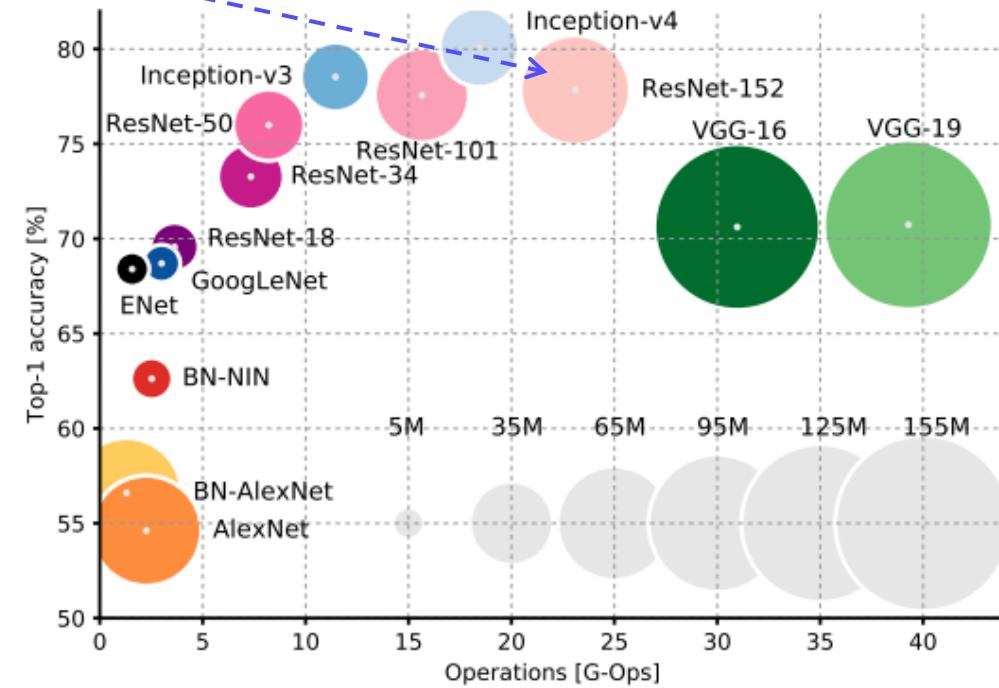
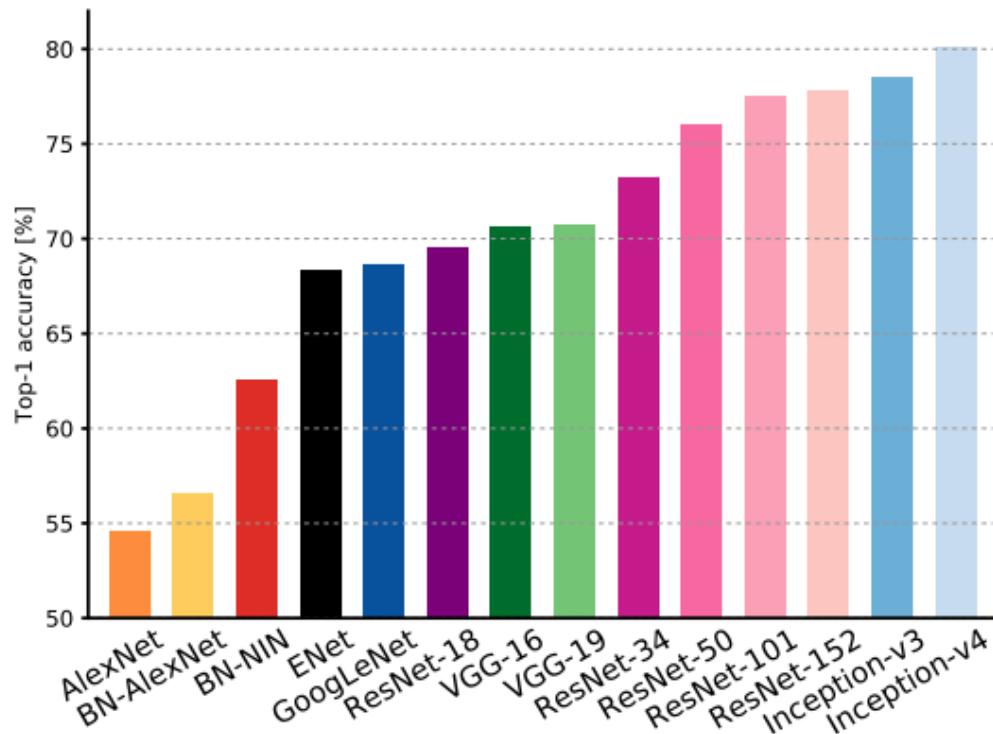


## 3.1 Summary of image classification

Summary of image classification

[Canziani et al., CVPR 2016]

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

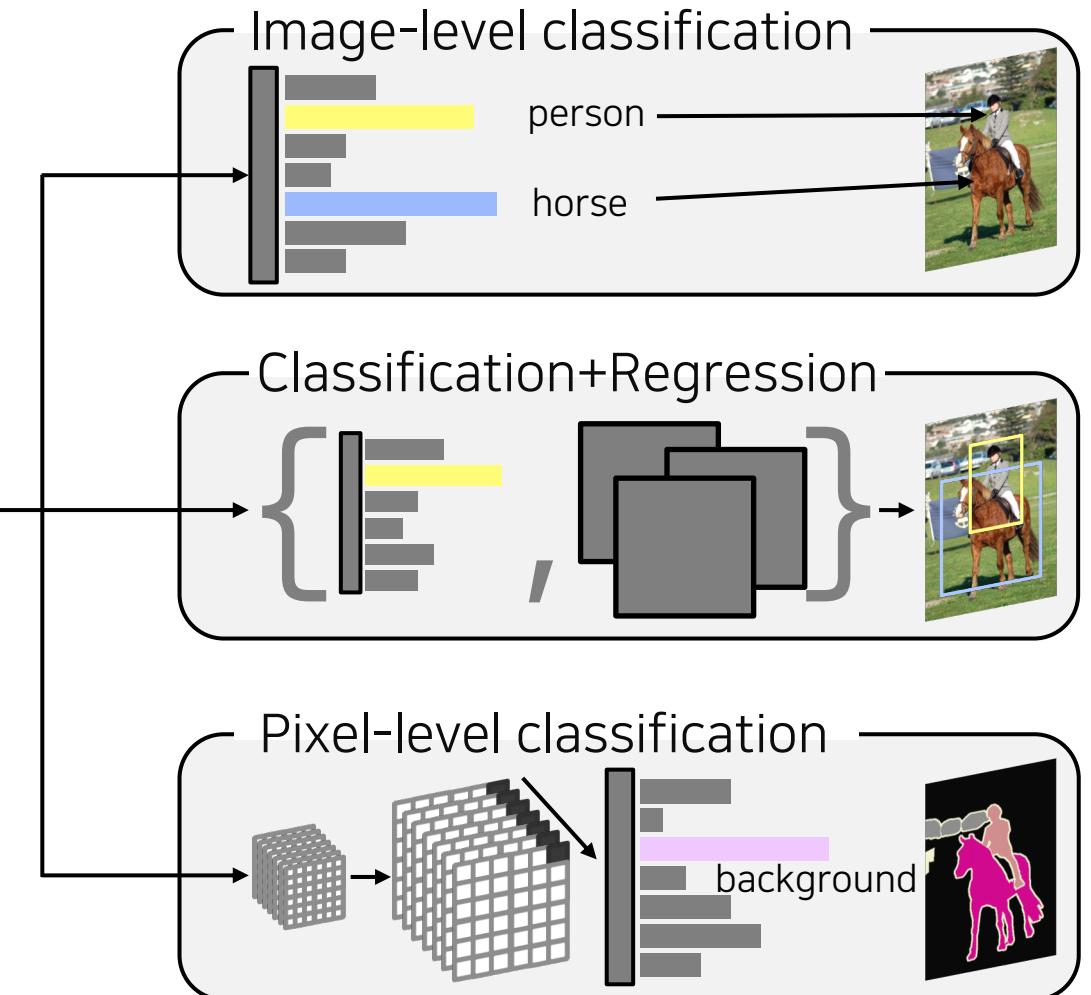
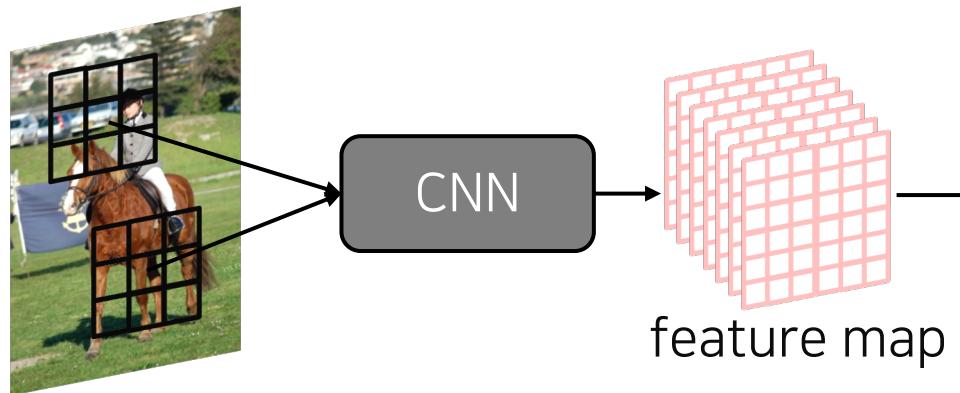


## 3.2 CNN backbones

Summary of image classification

Simple but powerful backbone model

- GoogLeNet is the most efficient CNN model out of {AlexNet, VGG, ResNet}
- But, it is complicated to use

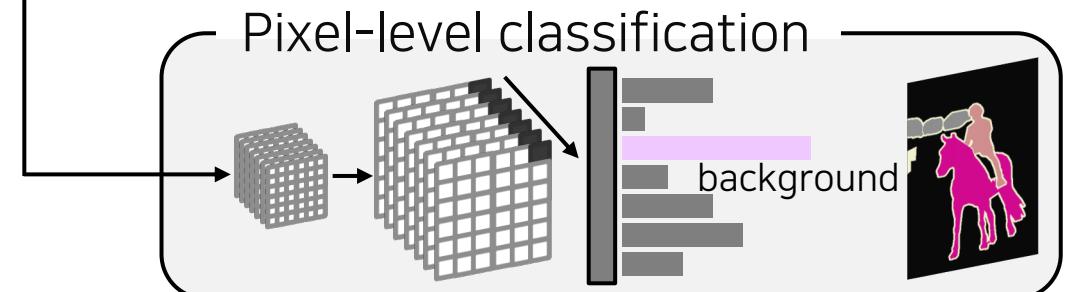
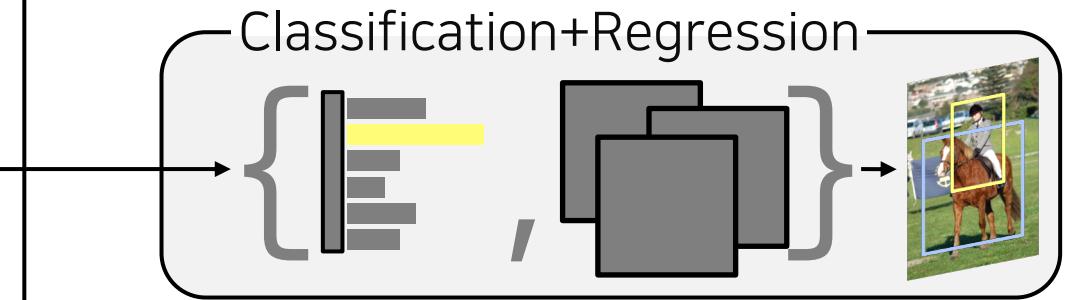
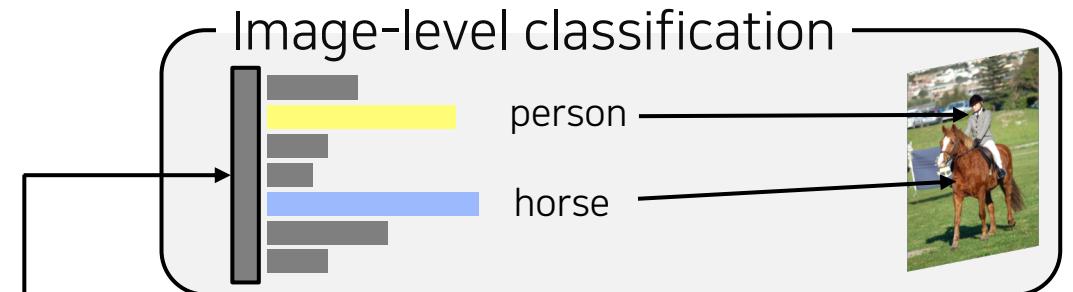
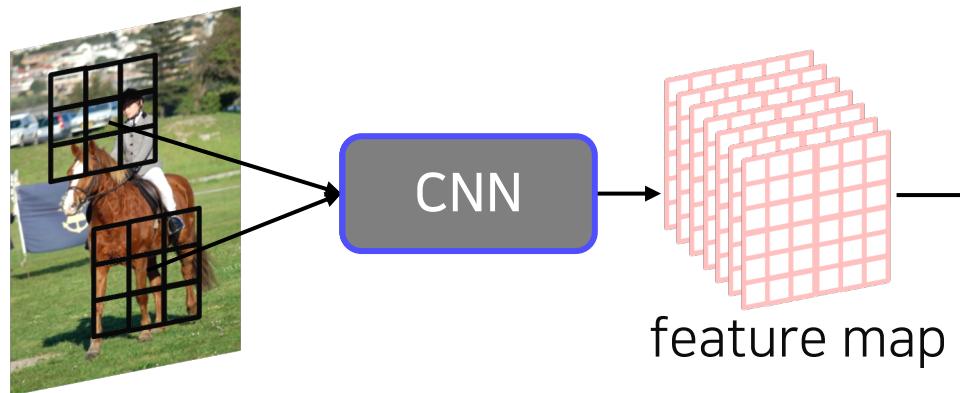


## 3.2 CNN backbones

Summary of image classification

Simple but powerful backbone model

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



# Reference

---

## 2. CNN architectures for image classification 2

- 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

## 3. Summary of image classification

- Canziani et al., An Analysis of Deep Neural Network Models for Practical Applications, CVPR 2016

# End of Document

## Thank You.