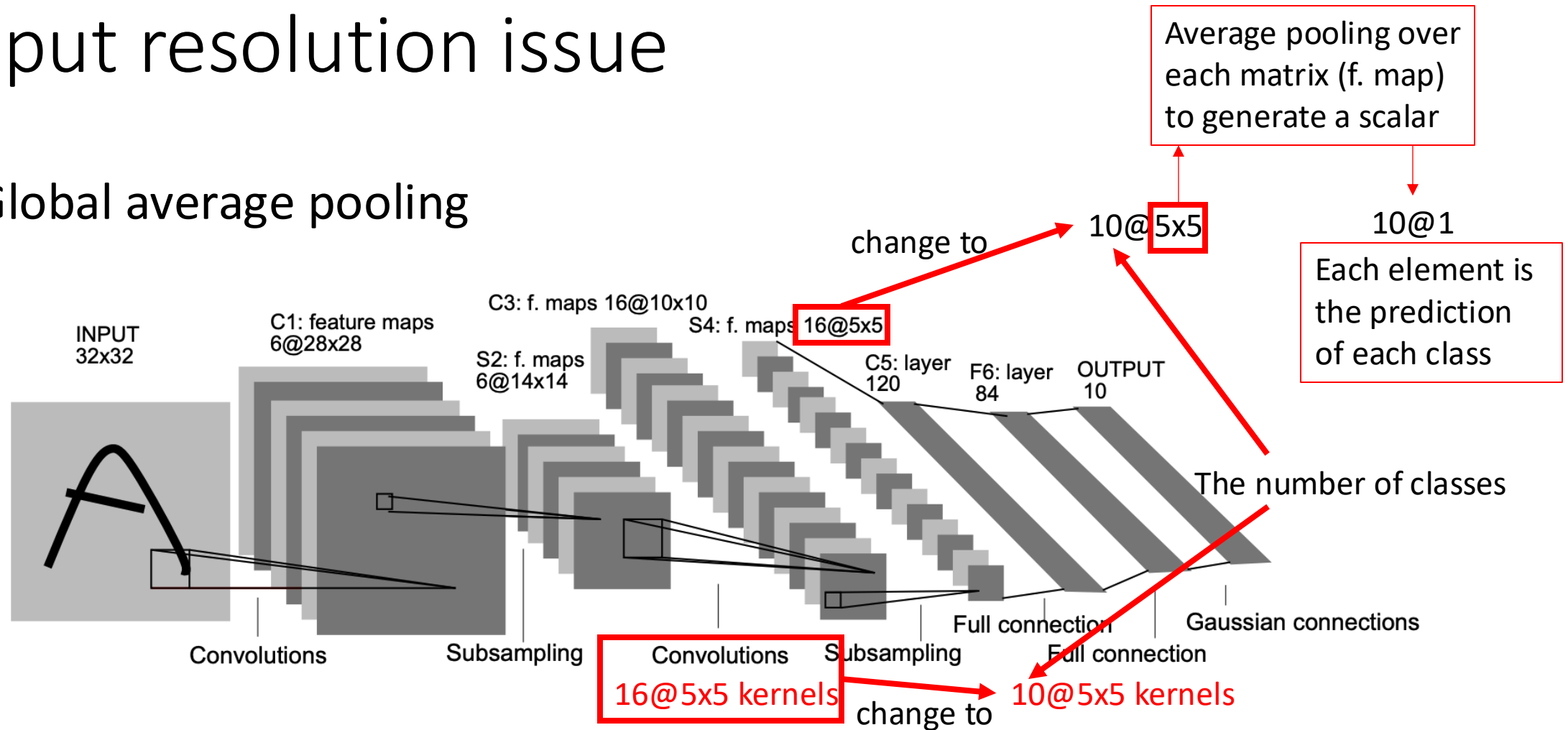


# CNN Architectures

Neural Networks Design And Application

# Input resolution issue

- Global average pooling

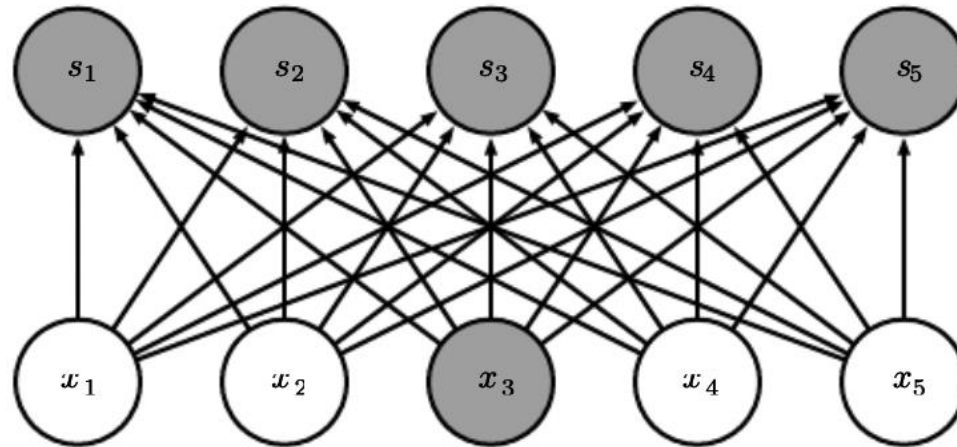


**Fig. 1.** Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

# Difference between ConvNet and MLP

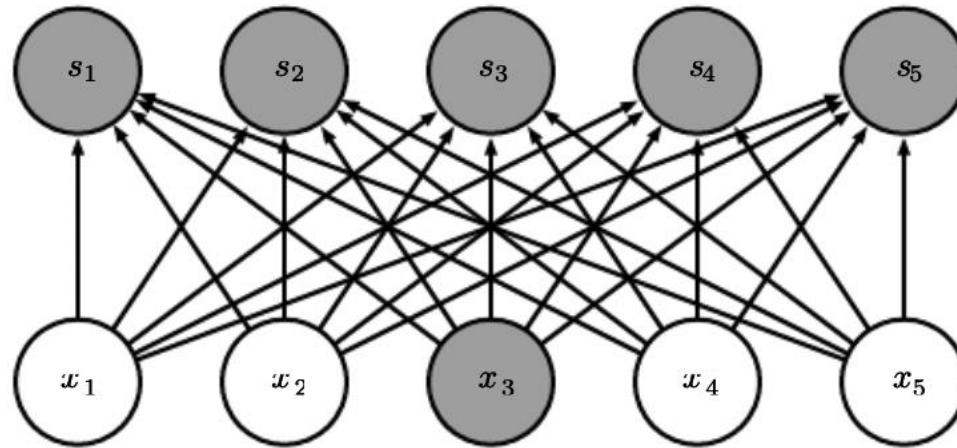
- Sparse connectivity
- Parameter sharing
- Equivariant representations

# Sparse connectivity of convolution



Feedforward network (fully connected layer)

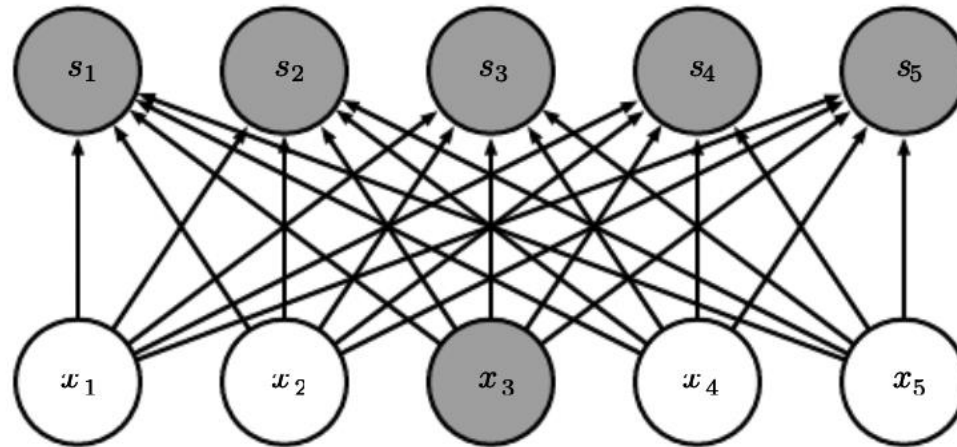
# Sparse connectivity of convolution



Feedforward network (fully connected layer)

Q: how many arrows we have?

# Sparse connectivity of convolution



Feedforward network (fully connected layer)

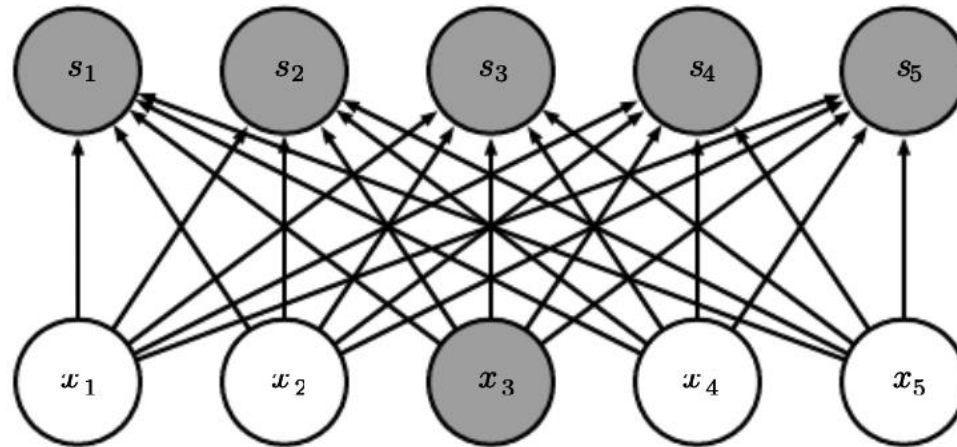
Q: how many arrows we have?

For each  $x_i$ : from  $s_1$  to  $s_5 \rightarrow 5$  arrows

$x_1, \dots, x_5 \rightarrow 25$  arrows

# Sparse connectivity of convolution

$$x'w \rightarrow s$$



Feedforward network (fully connected layer)

Q: how many arrows we have?

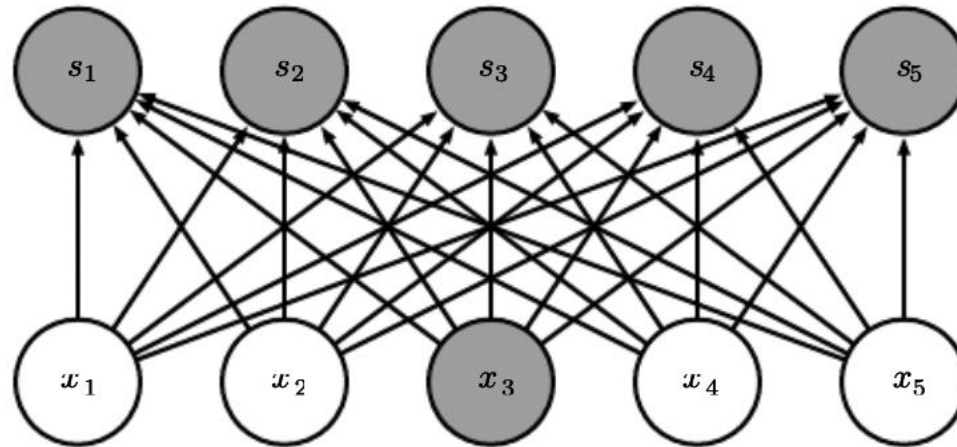
For each  $x_i$ : from  $s_1$  to  $s_5 \rightarrow 5$  arrows

$x_1, \dots, x_5 \rightarrow 25$  arrows

# Sparse connectivity of convolution

$x$  and  $w$  are vectors;  $s$  is a scalar number

$$x'w \rightarrow s$$



Feedforward network (fully connected layer)

Q: how many arrows we have?

For each  $x_i$ : from  $s_1$  to  $s_5 \rightarrow 5$  arrows

$x_1, \dots, x_5 \rightarrow 25$  arrows

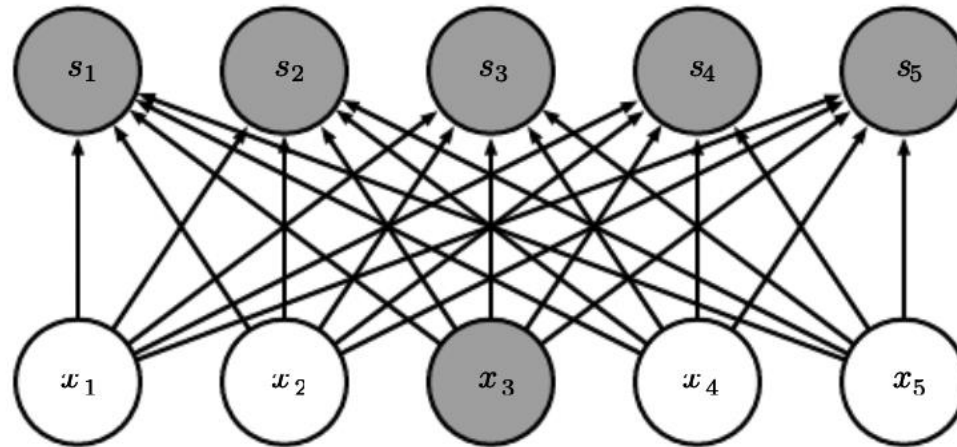


# Sparse connectivity of convolution

$x$  and  $w$  are vectors;  $s$  is a scalar number

$$x'w \rightarrow s$$

$$x'w_1 = \sum_{i=1}^5 x_i w_{1,i}$$



Feedforward network (fully connected layer)

Q: how many arrows we have?

For each  $x_i$ : from  $s_1$  to  $s_5 \rightarrow 5$  arrows

$x_1, \dots, x_5 \rightarrow 25$  arrows

# Sparse connectivity of convolution

$$x'w \rightarrow s$$

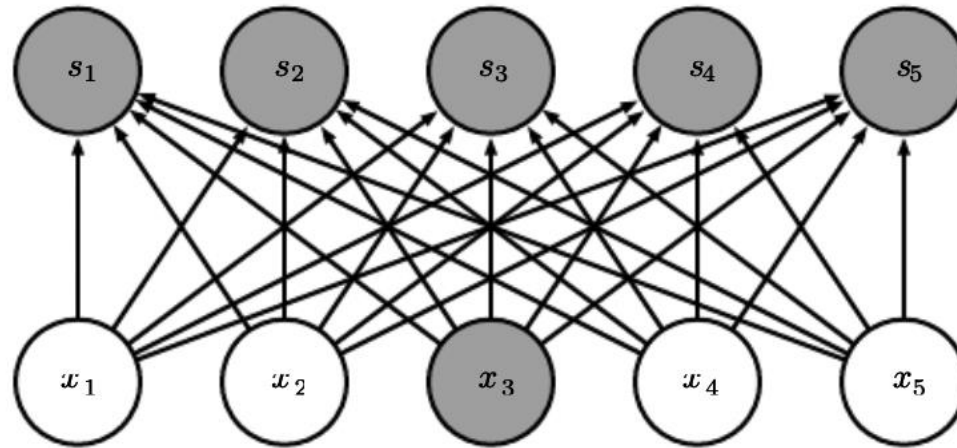
$$x'w_1 = \sum_{i=1}^5 x_i w_{1,i}$$

$$x'w_2 = \sum_{i=1}^5 x_i w_{2,i}$$

$$x'w_3 = \sum_{i=1}^5 x_i w_{3,i}$$

$$x'w_4 = \sum_{i=1}^5 x_i w_{4,i}$$

$$x'w_5 = \sum_{i=1}^5 x_i w_{5,i}$$



Feedforward network (fully connected layer)

Q: how many arrows we have?

For each  $x_i$ : from  $s_1$  to  $s_5 \rightarrow 5$  arrows

$x_1, \dots, x_5 \rightarrow 25$  arrows

# Sparse connectivity of convolution

$$x'w \rightarrow s$$

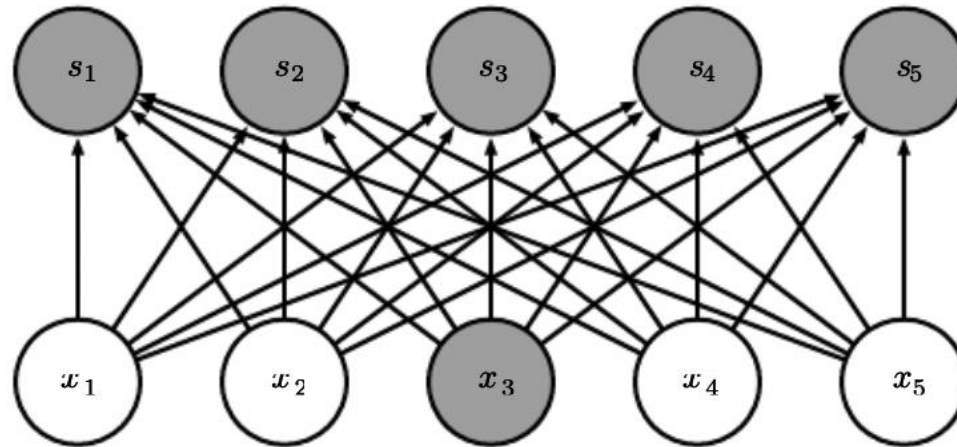
$$x'w_1 = \sum_{i=1}^5 x_i w_{1,i}$$

$$x'w_2 = \sum_{i=1}^5 x_i w_{2,i}$$

$$x'w_3 = \sum_{i=1}^5 x_i w_{3,i}$$

$$x'w_4 = \sum_{i=1}^5 x_i w_{4,i}$$

$$x'w_5 = \sum_{i=1}^5 x_i w_{5,i}$$



Feedforward network (fully connected layer)

Q: how many arrows we have?

For each  $x_i$ : from  $s_1$  to  $s_5 \rightarrow 5$  arrows

$x_1, \dots, x_5 \rightarrow 25$  arrows

# Sparse connectivity of convolution

$$x'w \rightarrow s$$

$$x'w_1 = \sum_{i=1}^5 x_i w_{1,i}$$

$$x'w_2 = \sum_{i=1}^5 x_i w_{2,i}$$

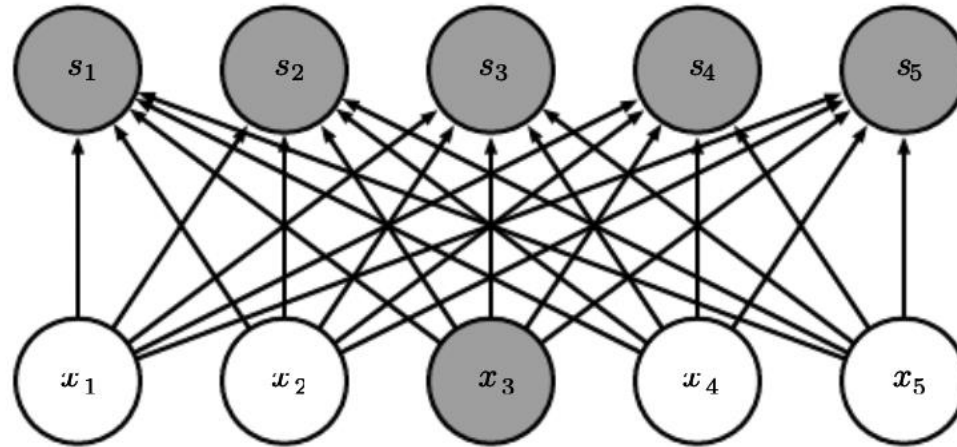
$$x'w_3 = \sum_{i=1}^5 x_i w_{3,i}$$

$$x'w_4 = \sum_{i=1}^5 x_i w_{4,i}$$

$$x'w_5 = \sum_{i=1}^5 x_i w_{5,i}$$

$W$   
||

$[w_1, w_2, w_3, w_4, w_5]$



Feedforward network (fully connected layer)

Q: how many arrows we have?

For each  $x_i$ : from  $s_1$  to  $s_5 \rightarrow 5$  arrows

$x_1, \dots, x_5 \rightarrow 25$  arrows

# Sparse connectivity of convolution

$$x'w \rightarrow s$$

$$x'w_1 = \sum_{i=1}^5 x_i w_{1,i}$$

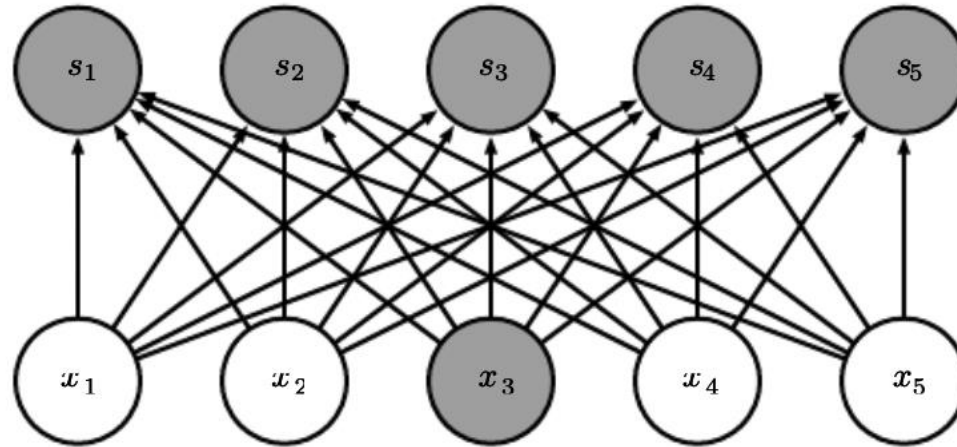
$$x'w_2 = \sum_{i=1}^5 x_i w_{2,i}$$

$$x'w_3 = \sum_{i=1}^5 x_i w_{3,i}$$

$$x'w_4 = \sum_{i=1}^5 x_i w_{4,i}$$

$$x'w_5 = \sum_{i=1}^5 x_i w_{5,i}$$

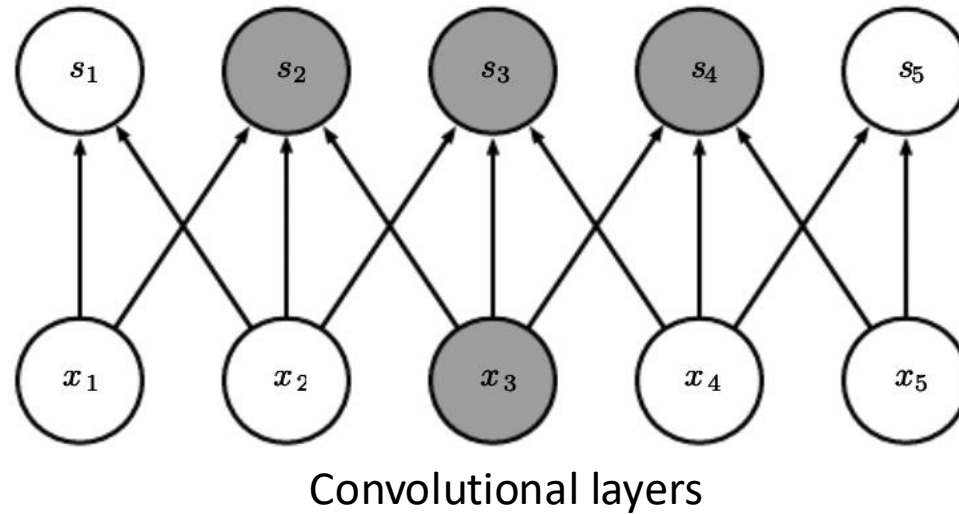
$W$   
 $\parallel$   
 $[w_1, w_2, w_3, w_4, w_5]$   
 Feedforward network (fully connected layer)  
 (5x5 weight matrix)



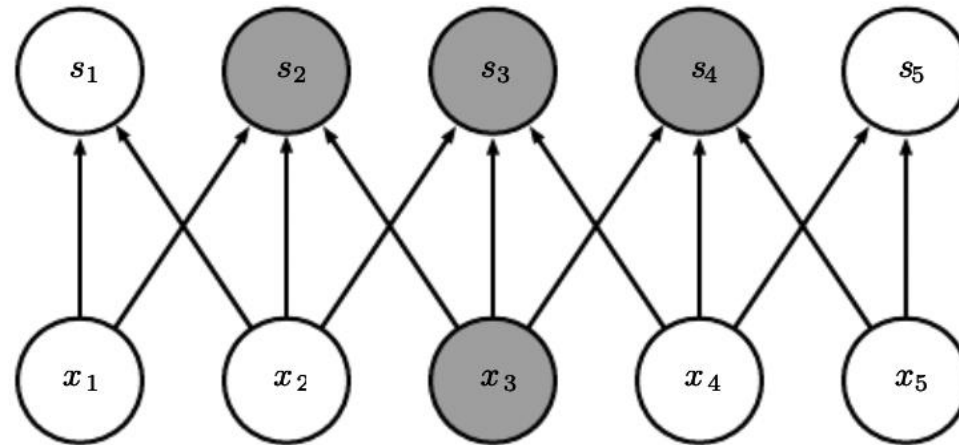
Q: how many arrows we have?

For each  $x_i$ : from  $s_1$  to  $s_5 \rightarrow 5$  arrows  
 $x_1, \dots, x_5 \rightarrow 25$  arrows

# Sparse connectivity of convolution



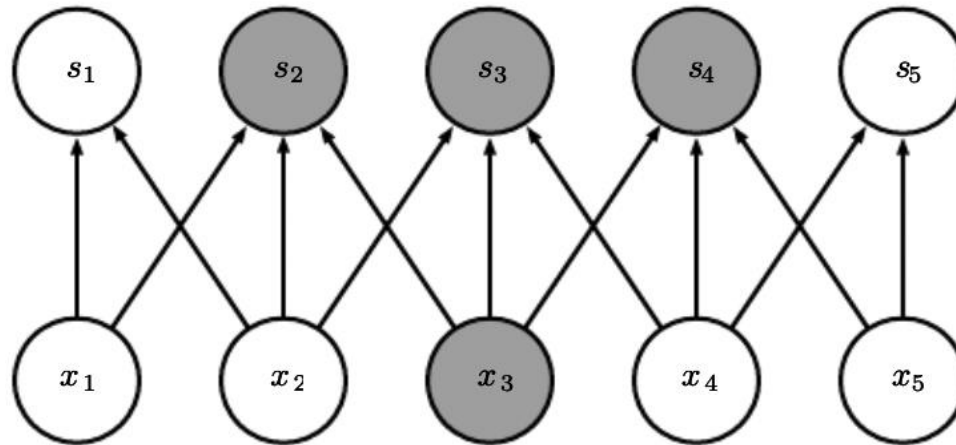
# Sparse connectivity of convolution



Convolutional layers

Q: how many arrows we have?

# Sparse connectivity of convolution



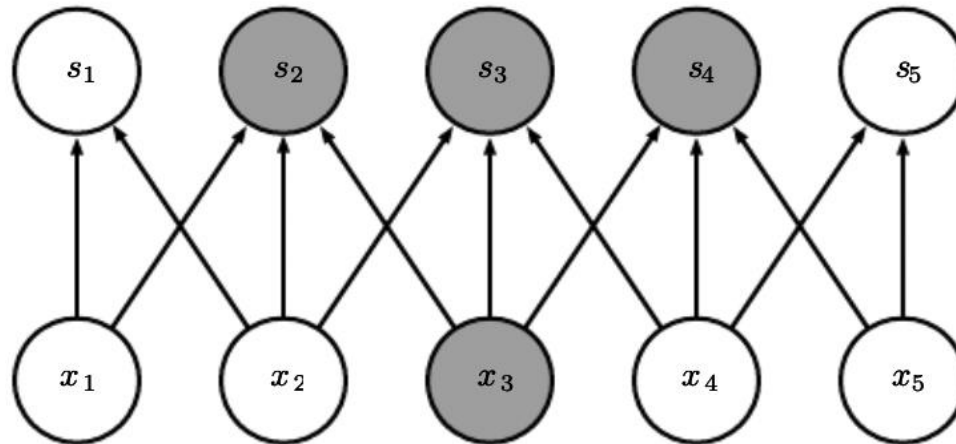
Convolutional layers

Q: how many arrows we have?

For each  $x_i$ : connect to 3  $s$  outputs



# Sparse connectivity of convolution



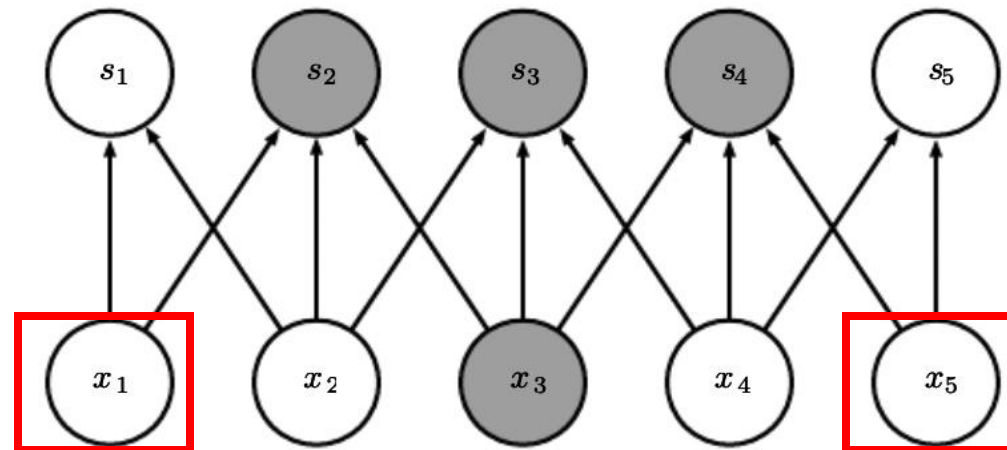
Convolutional layers

Q: how many arrows we have?

For each  $x_i$ : connect to 3 s outputs

$x_1, \dots, x_5 \rightarrow 3 \times 5 - 2 = 13$  arrows

# Sparse connectivity of convolution



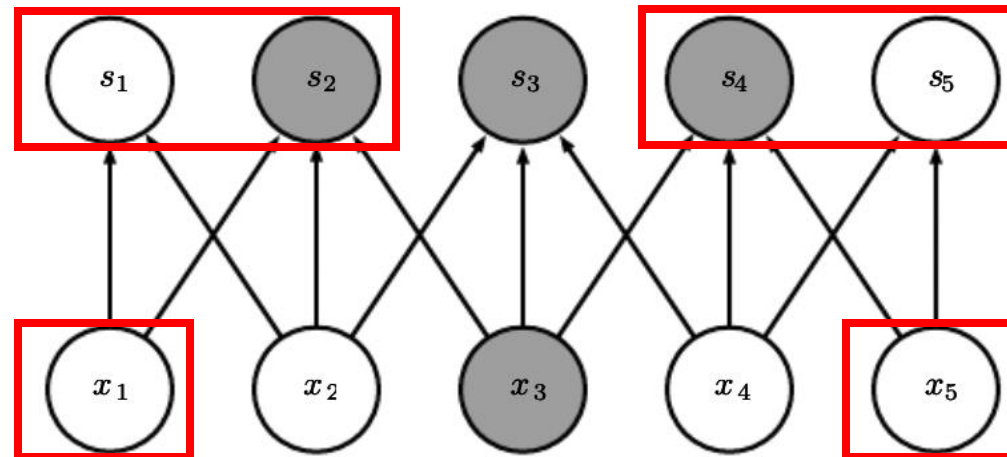
Convolutional layers

Q: how many arrows we have?

For each  $x_i$ : connect to 3  $s$  outputs

$x_1, \dots, x_5 \rightarrow 3 \times 5 - 2 = 13$  arrows

# Sparse connectivity of convolution



Convolutional layers

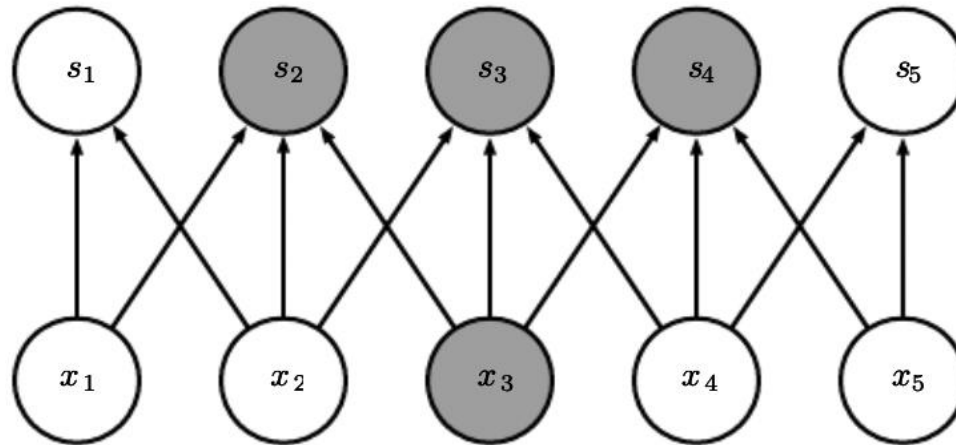
Q: how many arrows we have?

For each  $x_i$ : connect to 3 s outputs

$x_1, \dots, x_5 \rightarrow 3 \times 5 - 2 = 13$  arrows

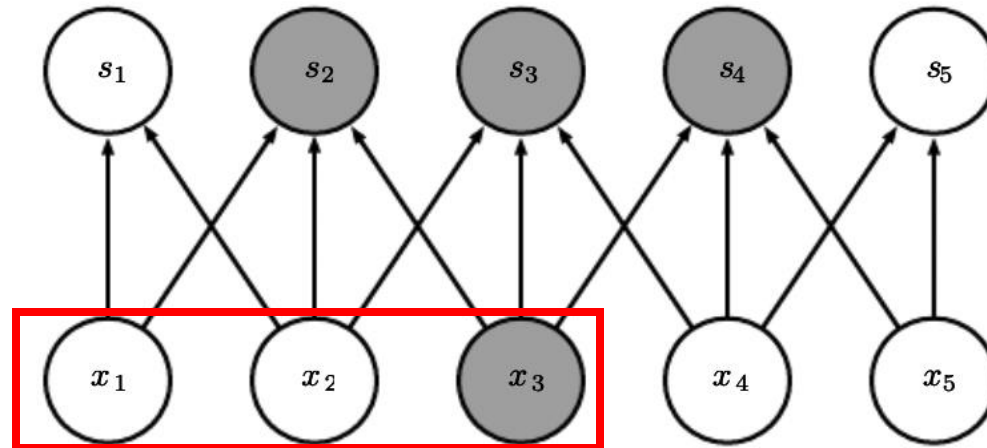
# Sparse connectivity of convolution

The view of convolutional kernel/filter



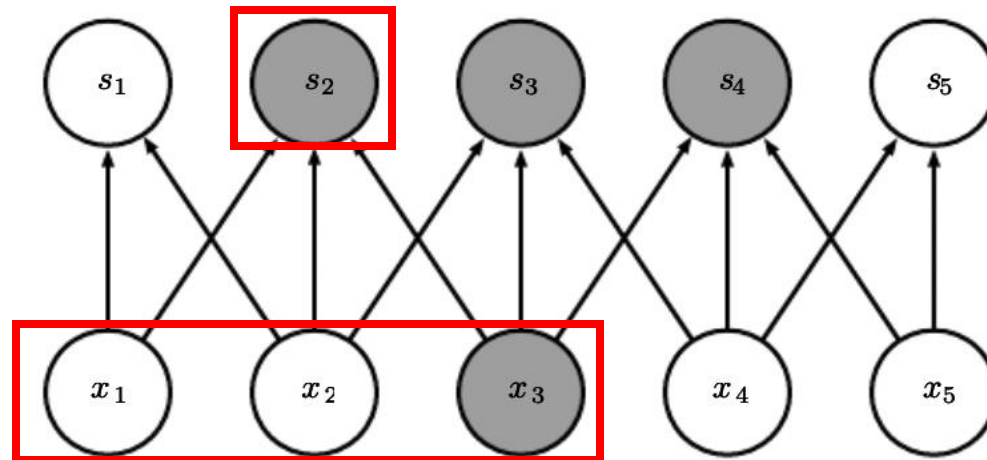
# Sparse connectivity of convolution

The view of convolutional kernel/filter



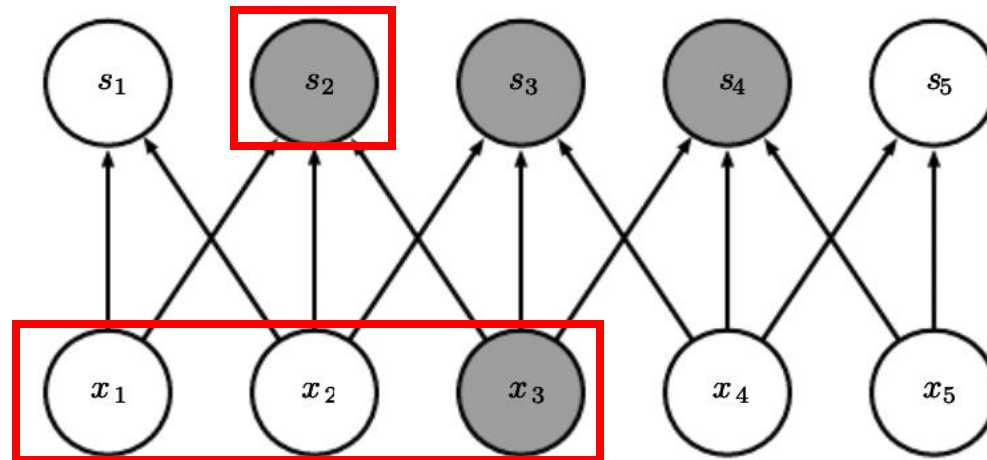
# Sparse connectivity of convolution

The view of convolutional kernel/filter



# Sparse connectivity of convolution

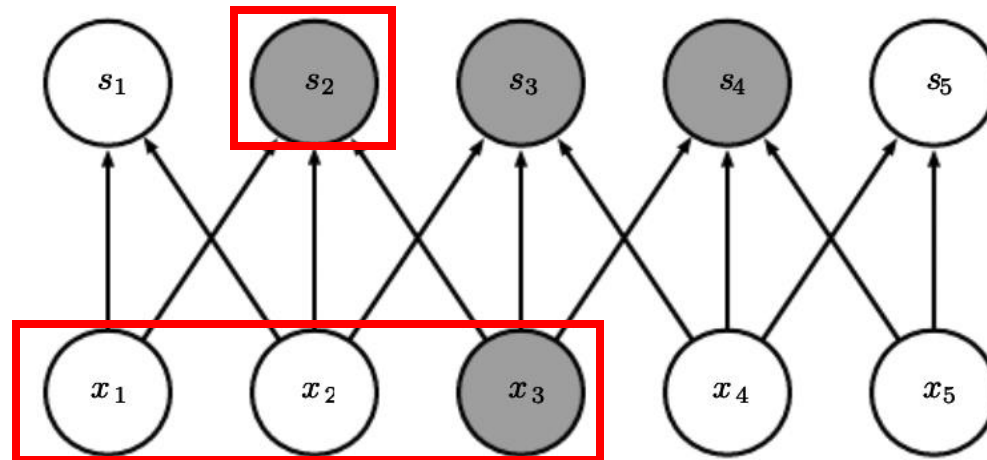
The view of convolutional kernel/filter



Q: filter size and stride?

# Sparse connectivity of convolution

The view of convolutional kernel/filter



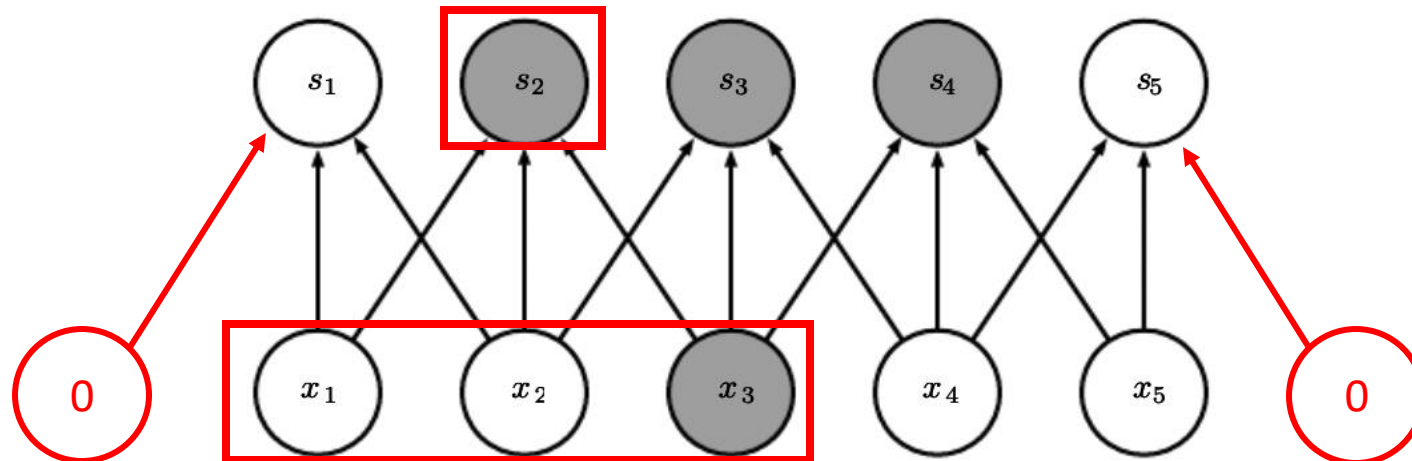
Q: filter size and stride?

Filter size = 3 + stride = 1 with 0-padding



# Sparse connectivity of convolution

The view of convolutional kernel/filter

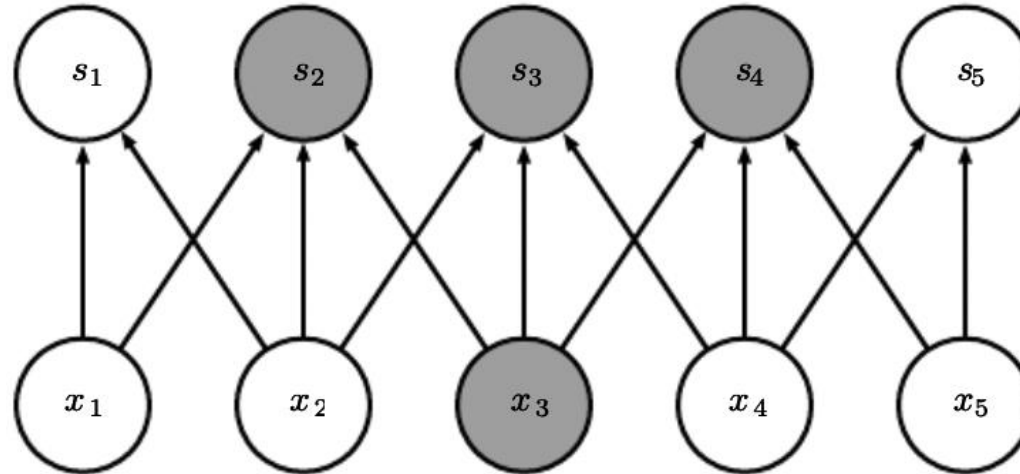


Q: filter size and stride?

Filter size = 3 + stride = 1 with 0-padding

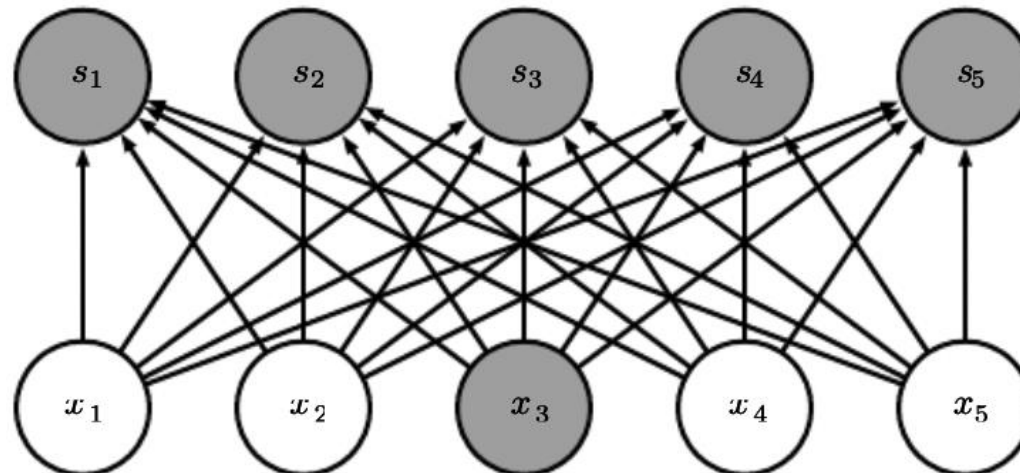
# Sparse connectivity of convolution

Sparse connection



Weights  $\rightarrow$  13 scalar numbers

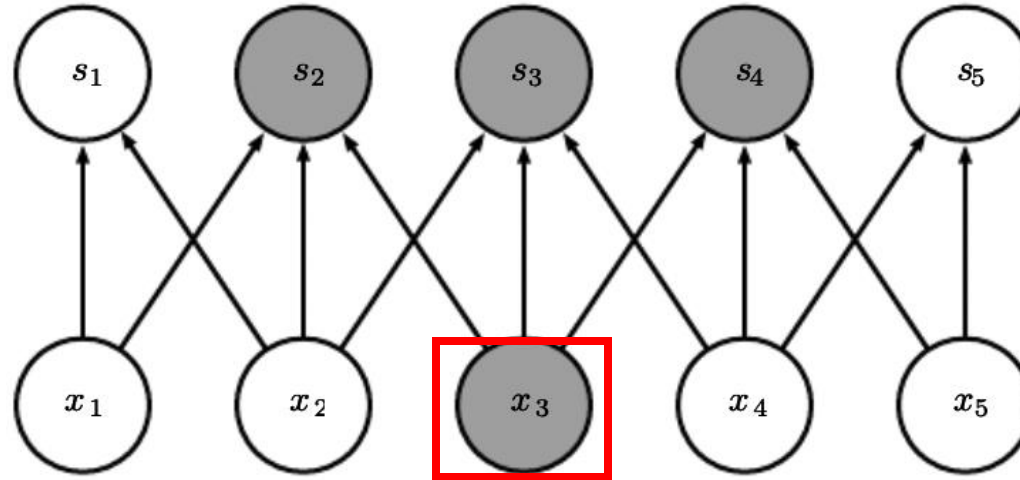
Dense connection



Weights  $\rightarrow$  25 scalar numbers

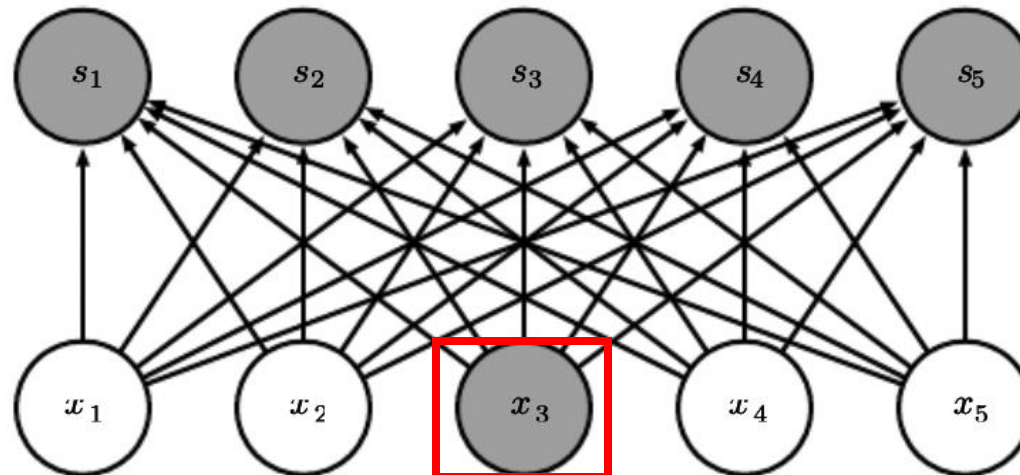
# Sparse connectivity of convolution

Sparse connection



Weights  $\rightarrow$  13 scalar numbers

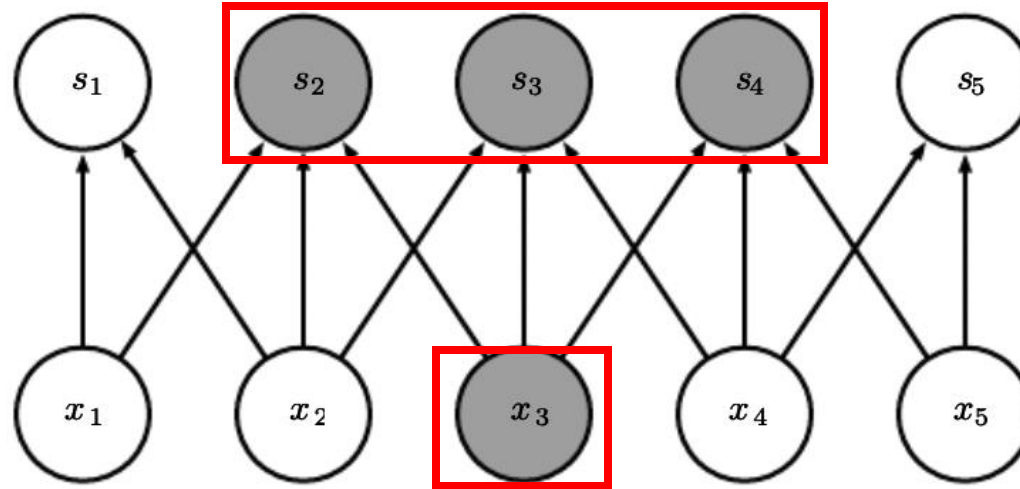
Dense connection



Weights  $\rightarrow$  25 scalar numbers

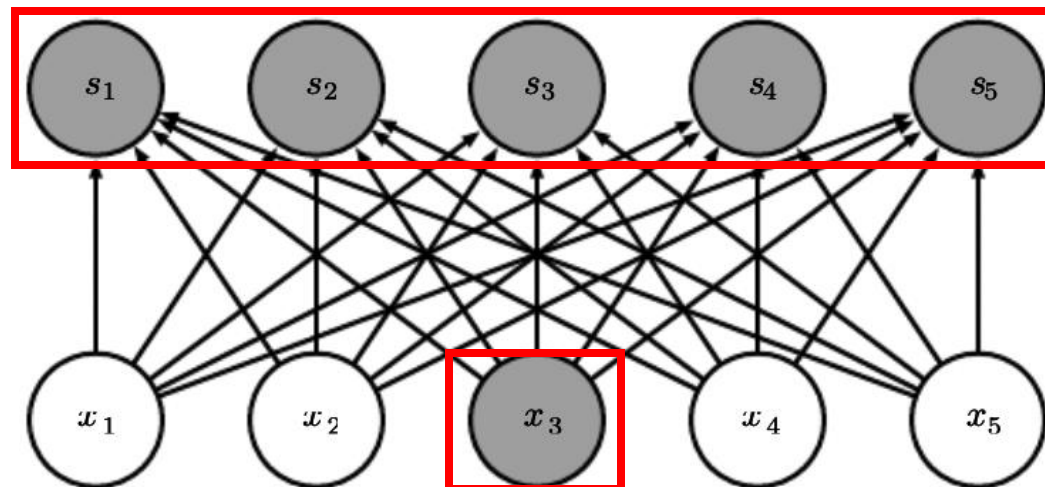
# Sparse connectivity of convolution

Sparse connection



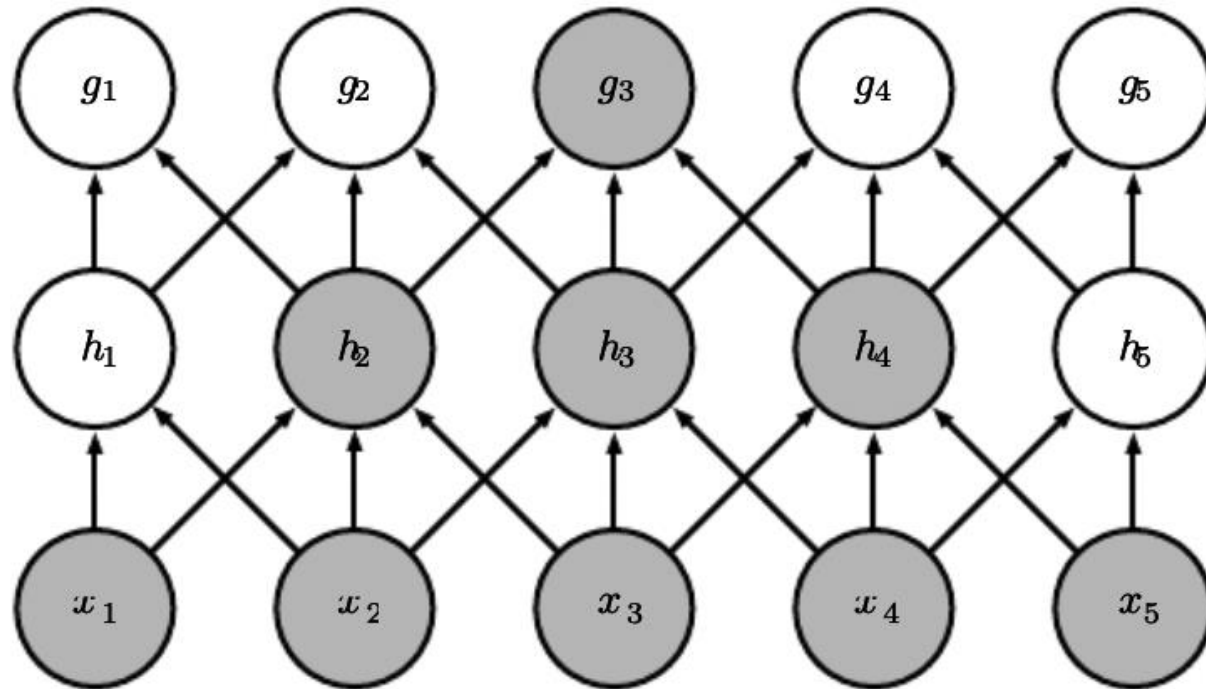
Weights  $\rightarrow$  13 scalar numbers

Dense connection

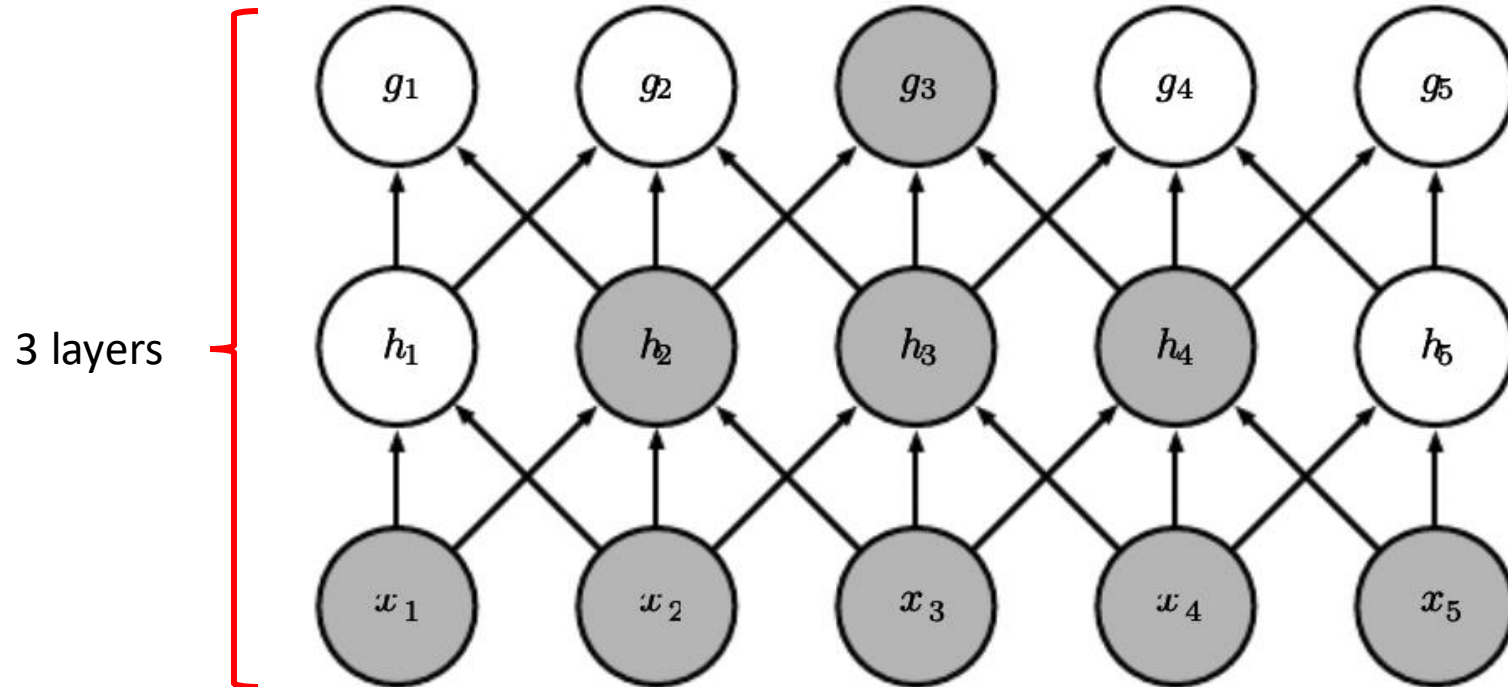


Weights  $\rightarrow$  25 scalar numbers

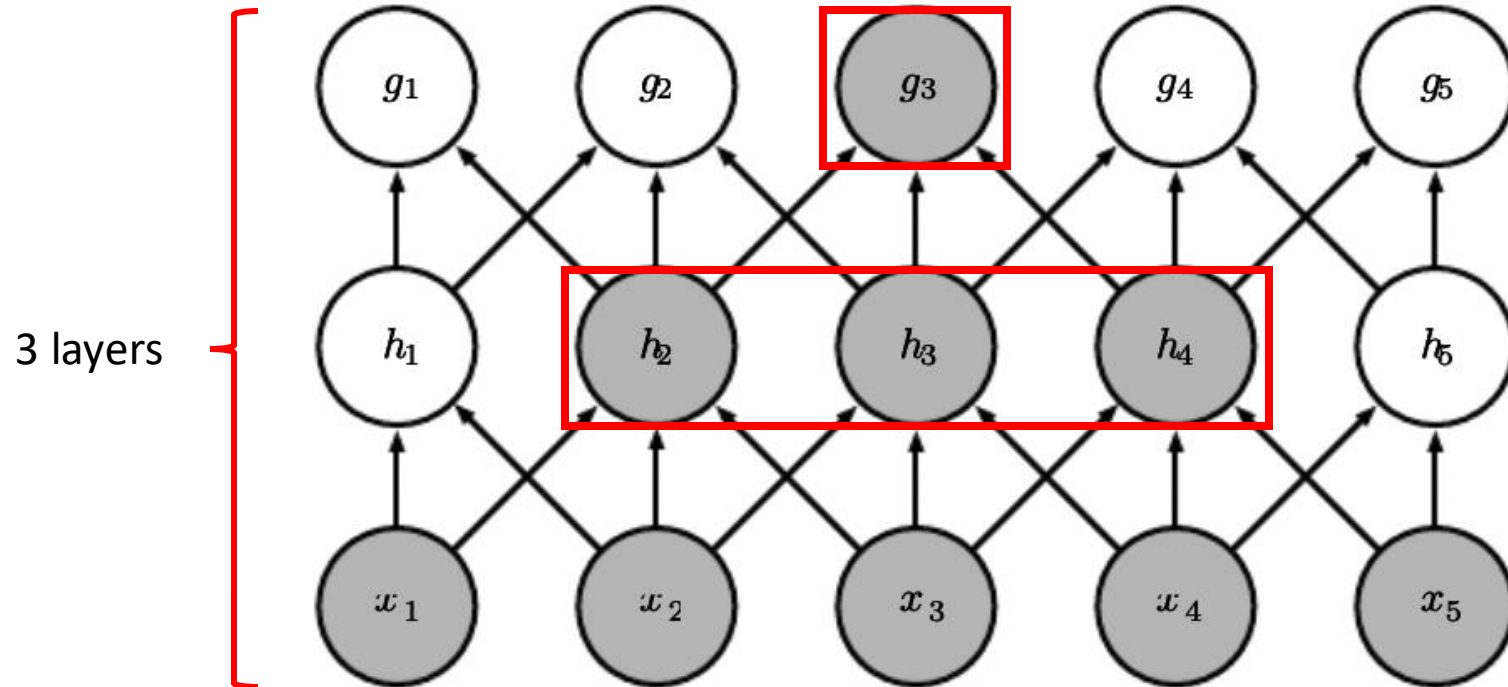
# Sparse connectivity of convolution



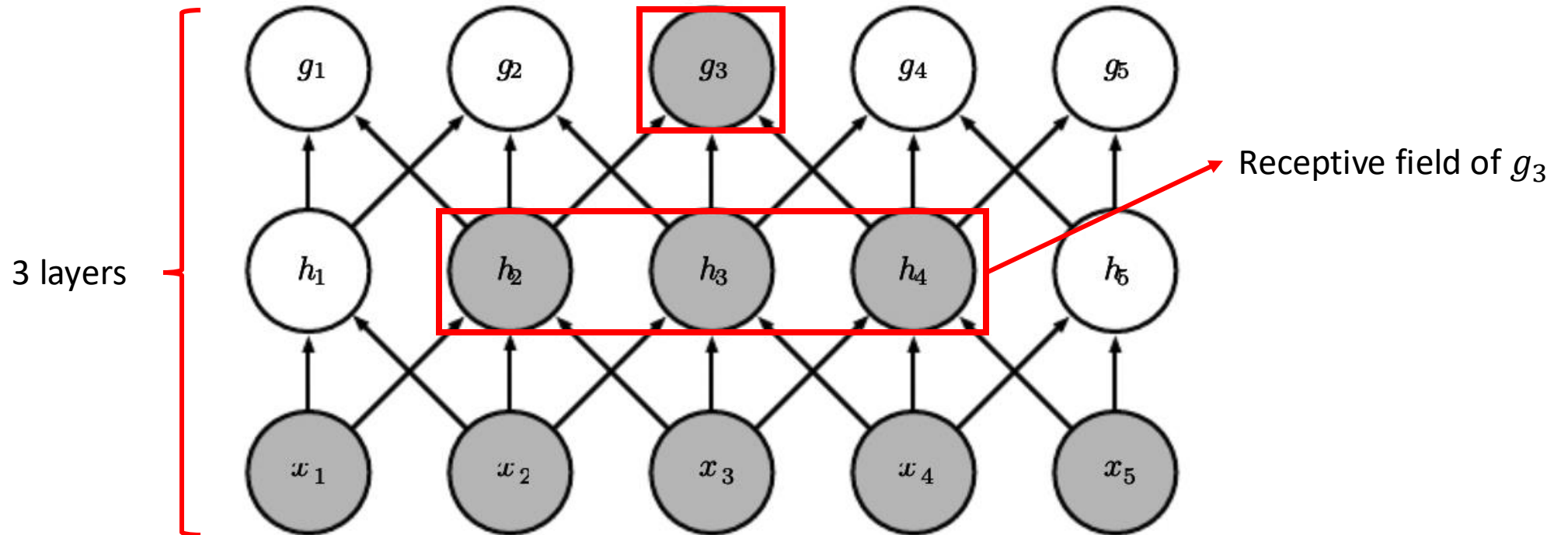
# Sparse connectivity of convolution



# Sparse connectivity of convolution

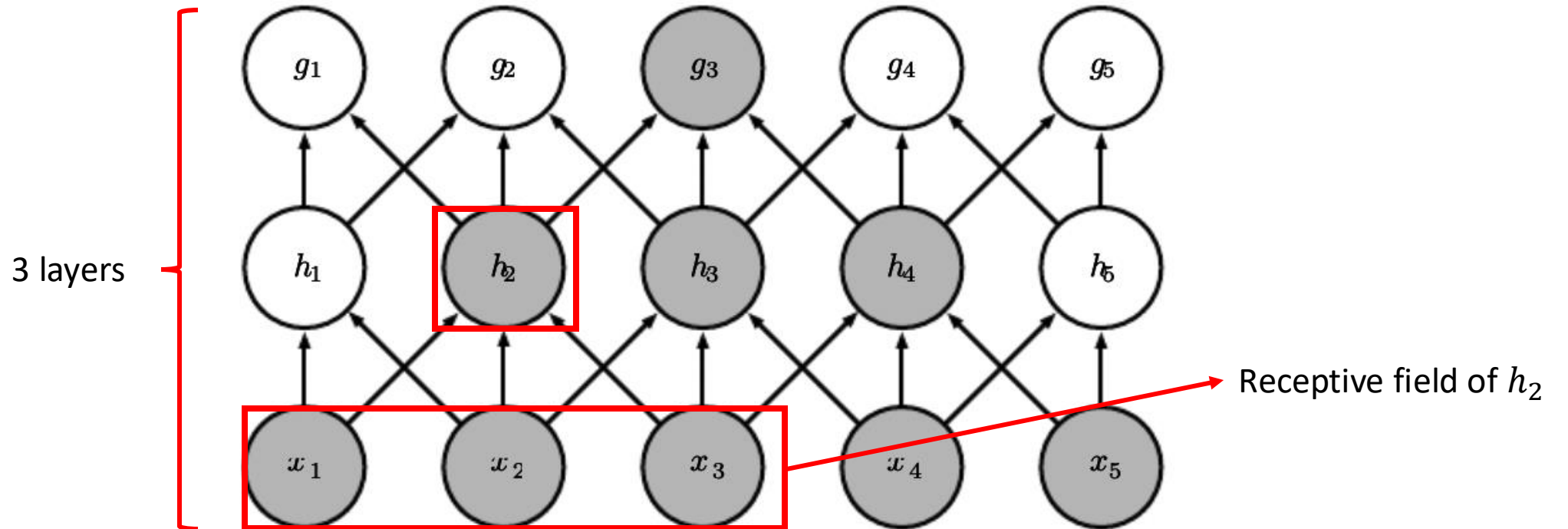


# Sparse connectivity of convolution

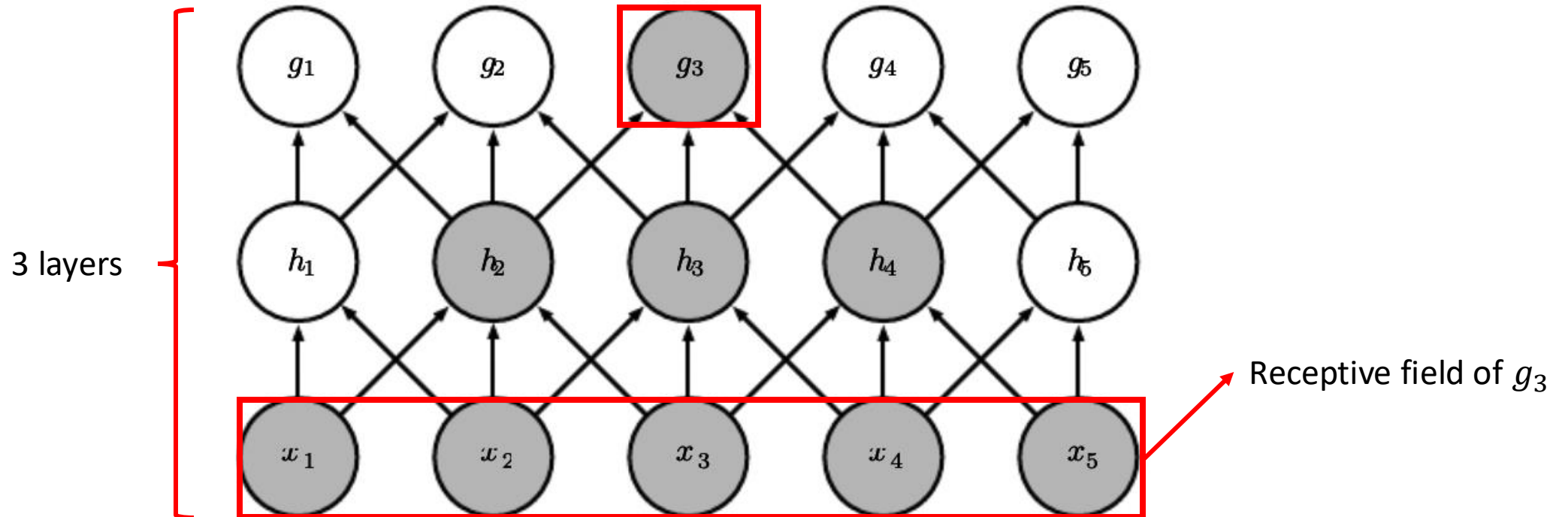




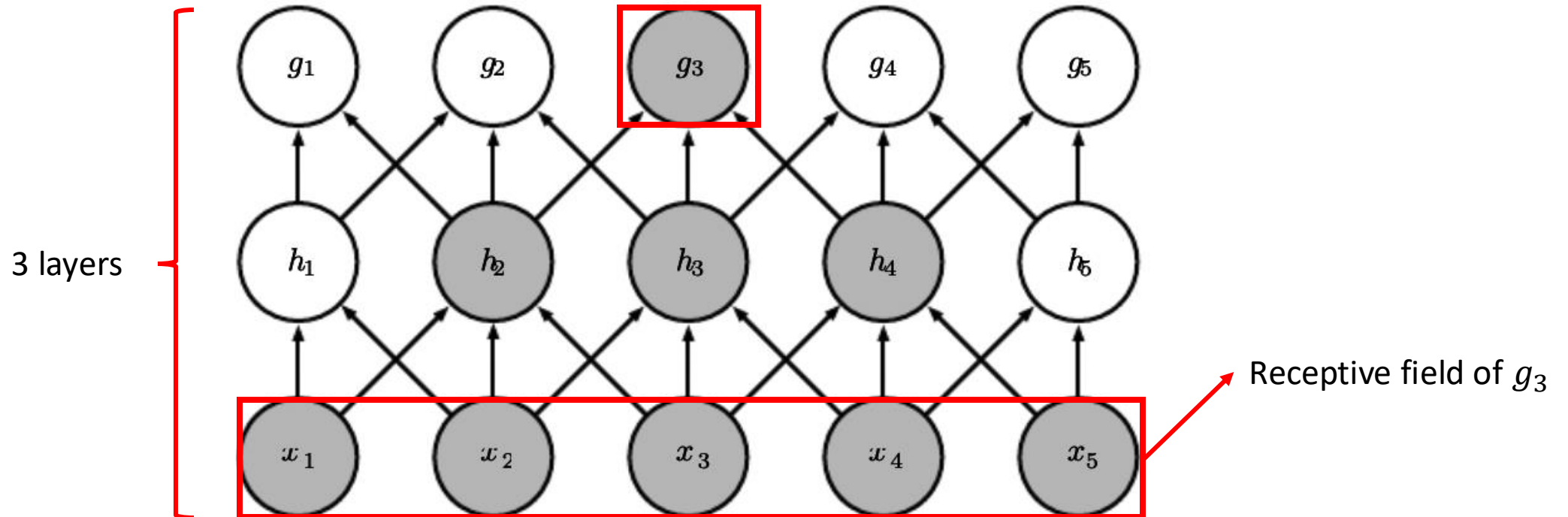
# Sparse connectivity of convolution



# Sparse connectivity of convolution

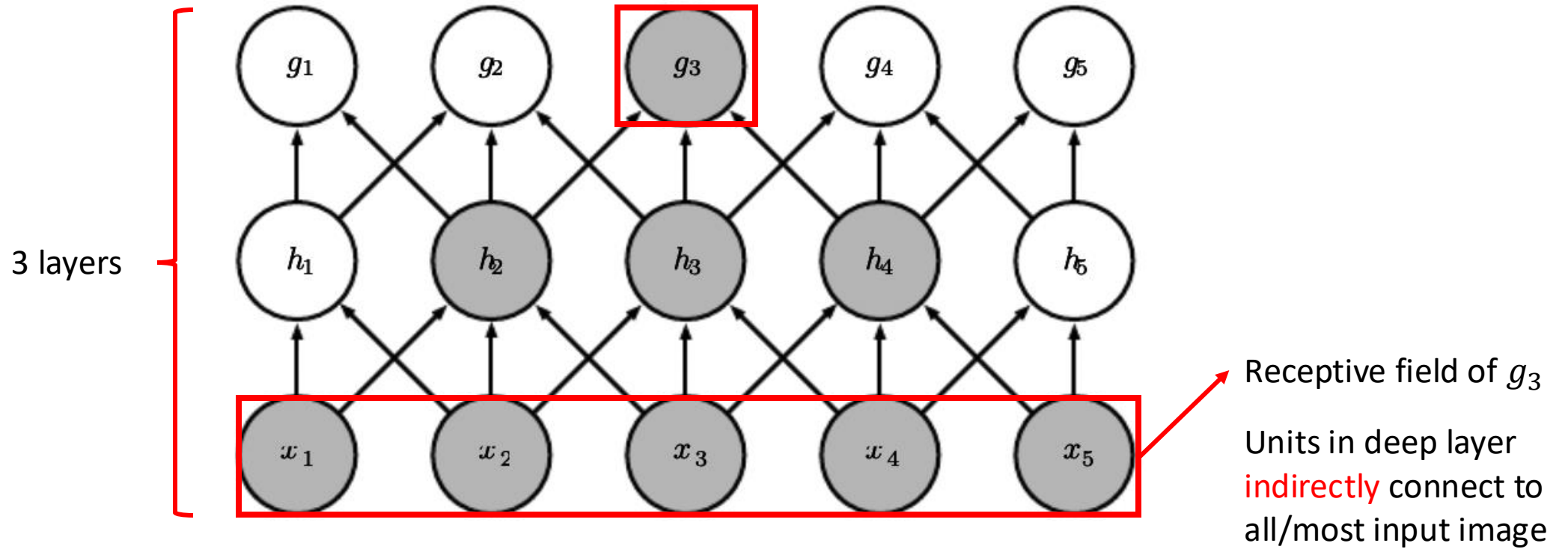


# Sparse connectivity of convolution



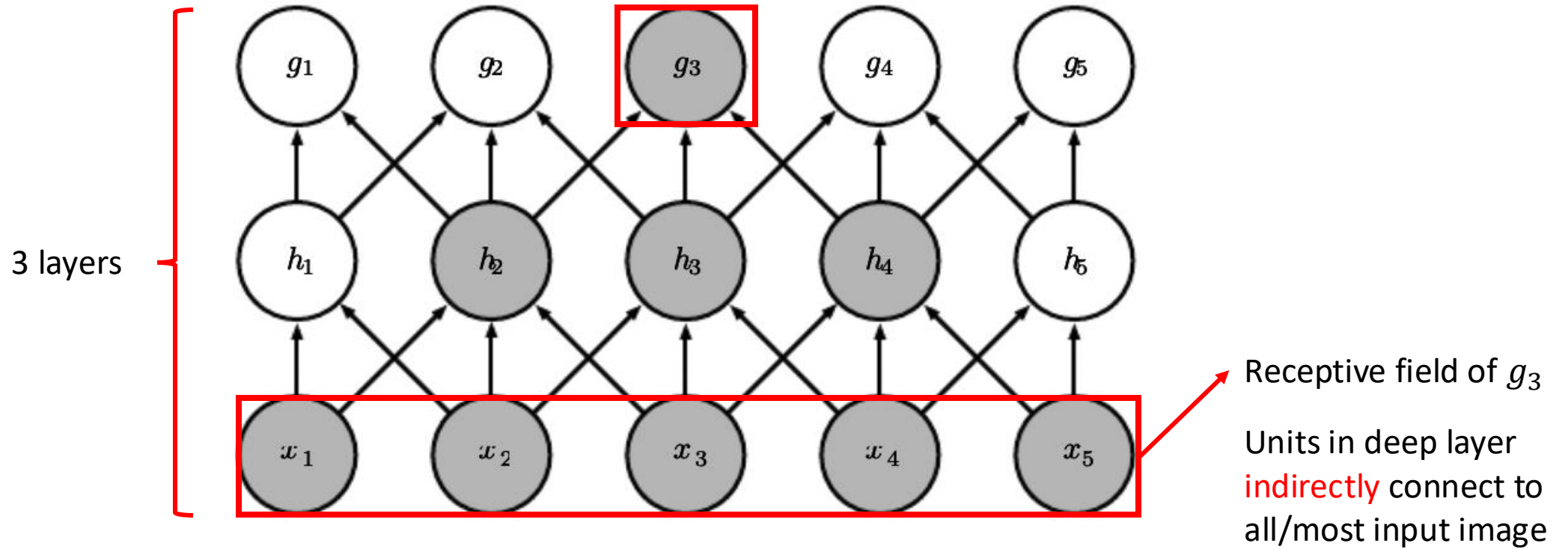
Deep layers has larger receptive field than shallow layers

# Sparse connectivity of convolution



Deep layers has larger receptive field than shallow layers

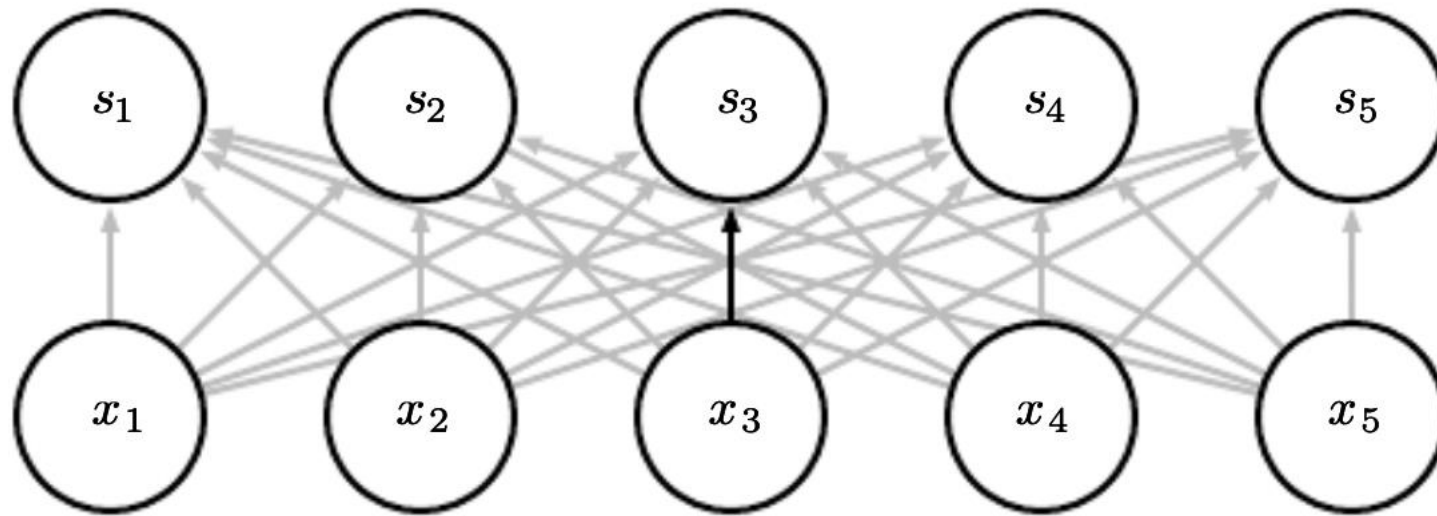
# Sparse connectivity of convolution



Deep layers has larger receptive field than shallow layers

Q: larger stride of convolution filter  $\rightarrow$  increase receptive field?

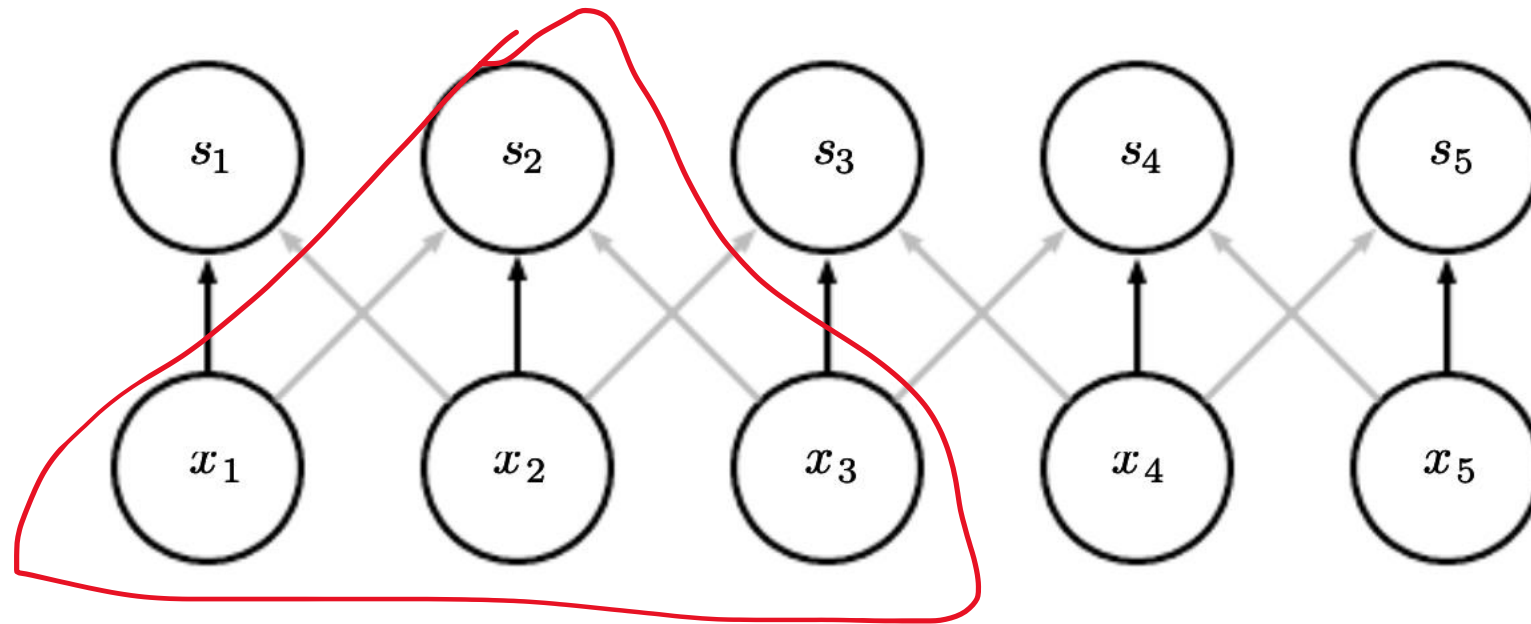
# Parameter sharing



In MLP (FC layer):  $w^T x$

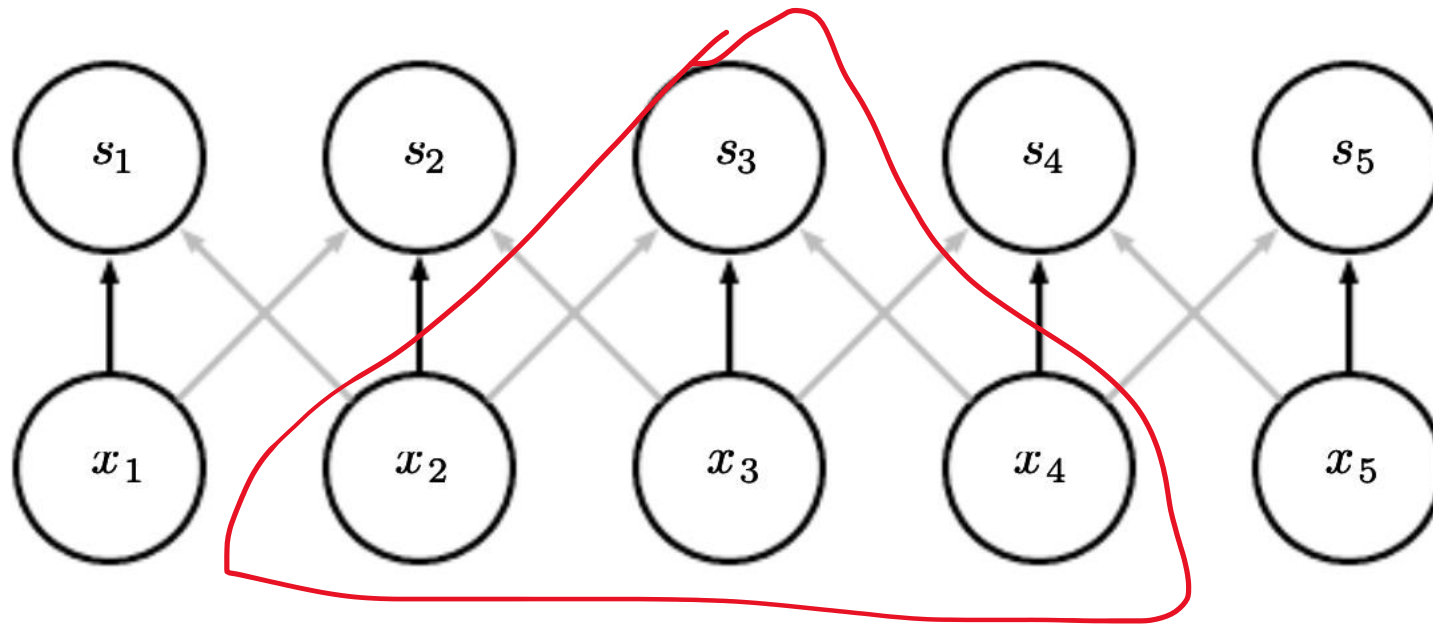
# Parameter sharing

Consider the same filter



# Parameter sharing

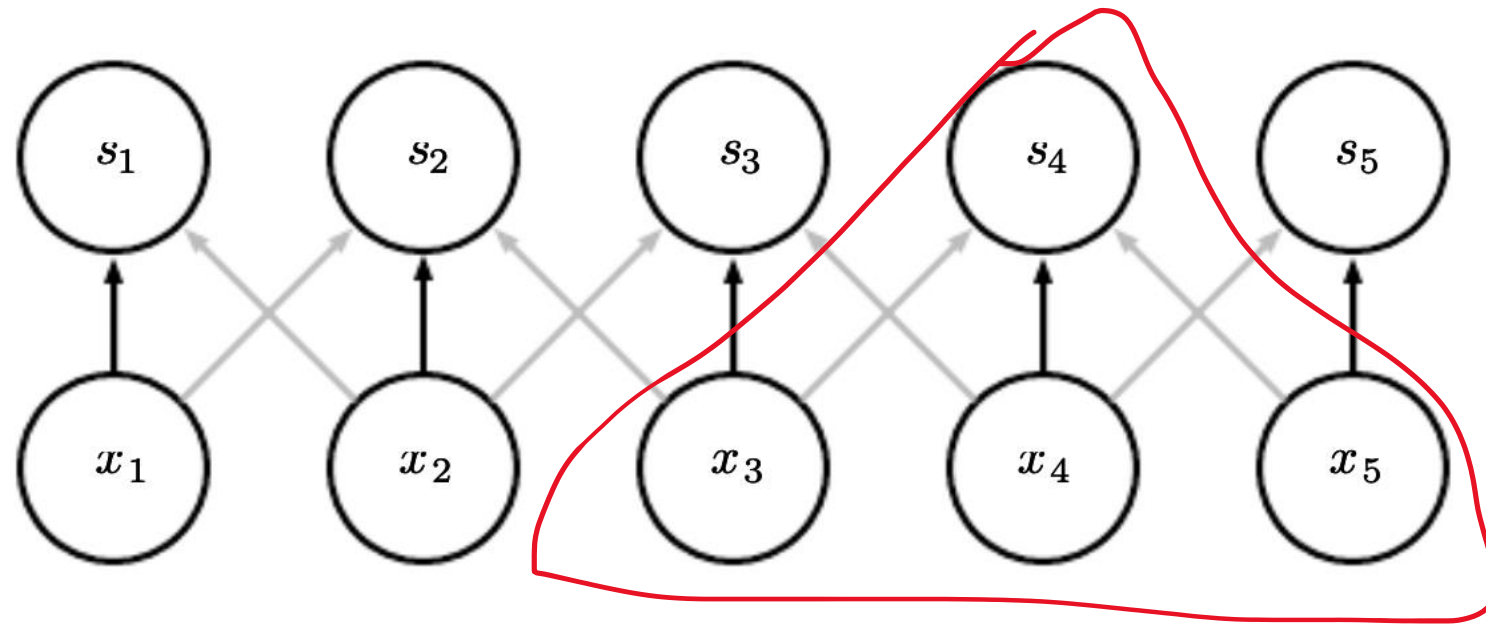
Consider the same filter





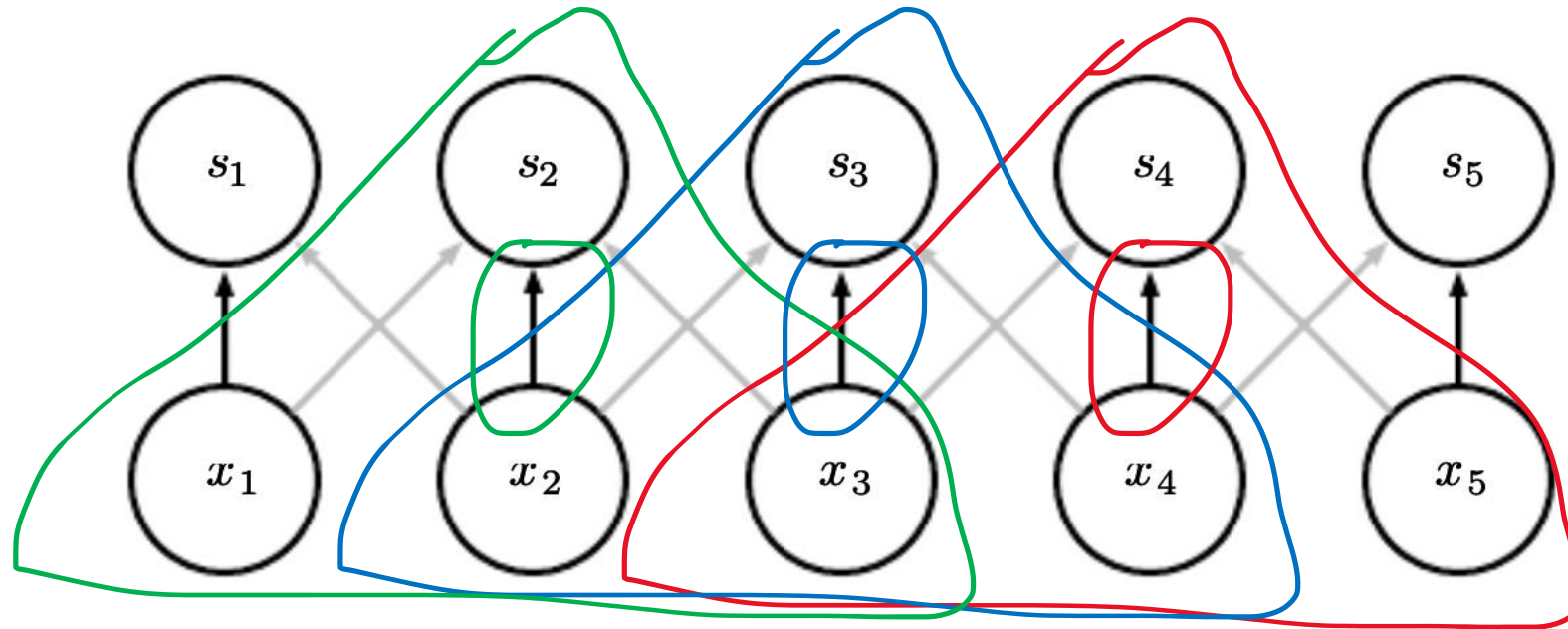
# Parameter sharing

Consider the same filter



# Parameter sharing

Consider the same filter (but different part of input feat. map)



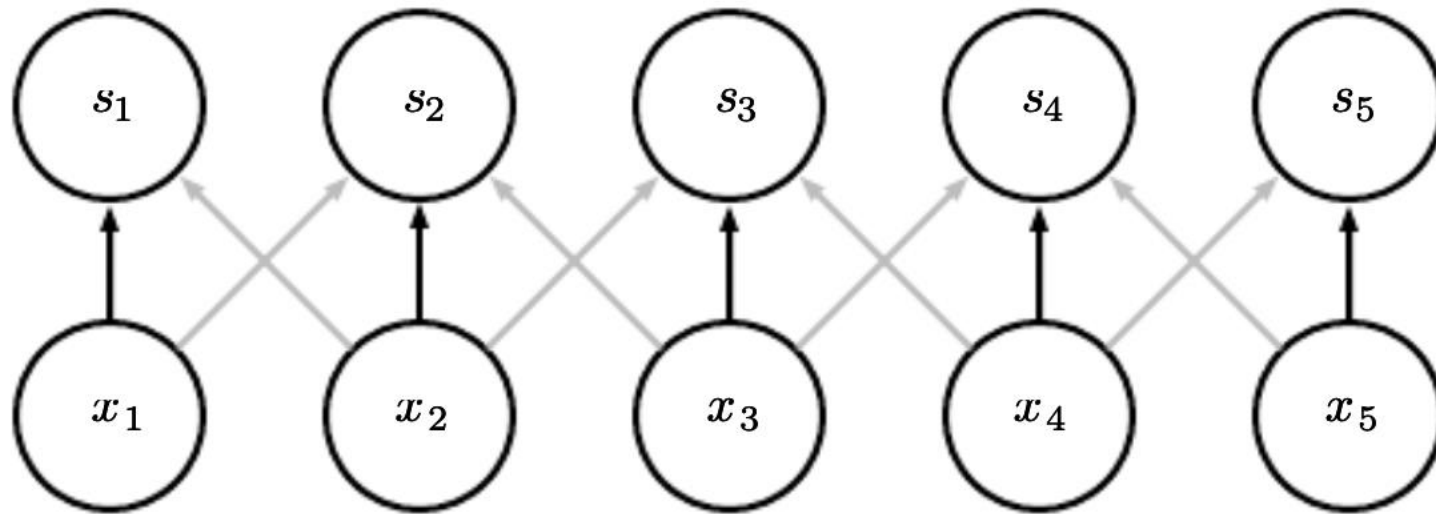
In convlayer:

$$w^T(x_1; x_2; x_3)$$

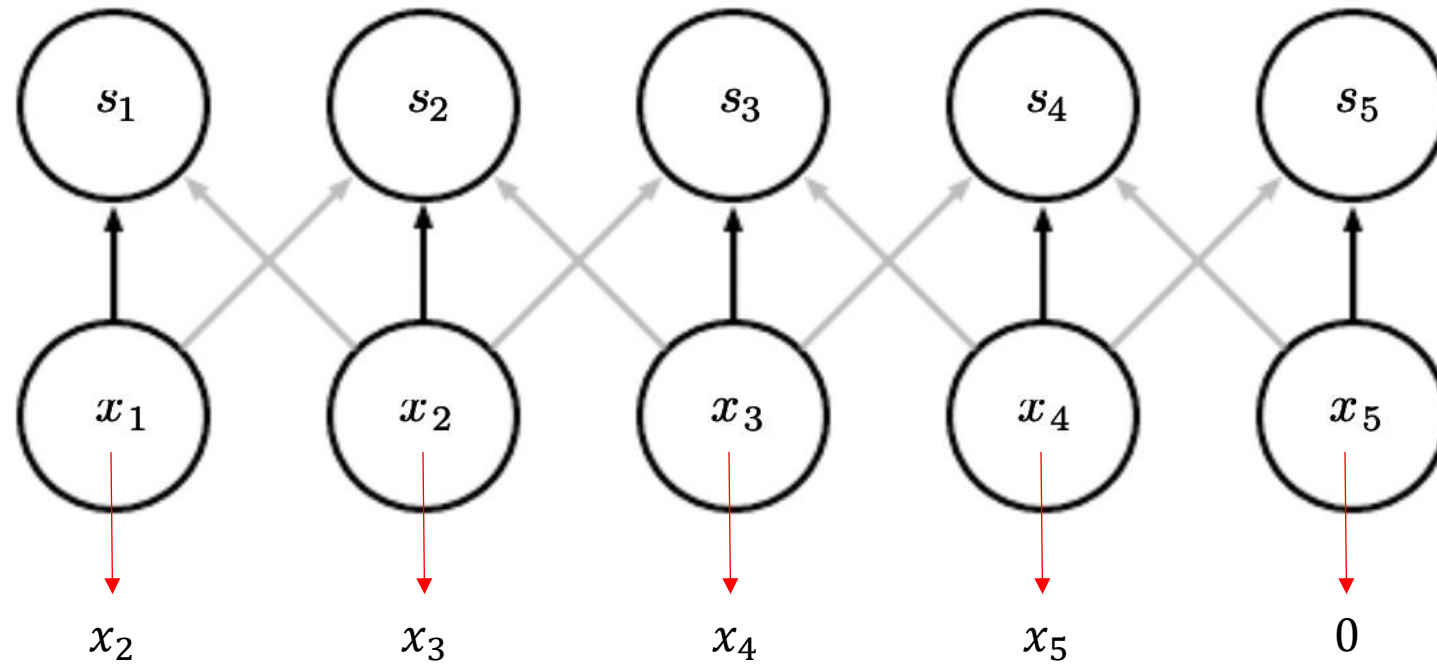
$$w^T(x_2; x_3; x_4)$$

$$w^T(x_3; x_4; x_5)$$

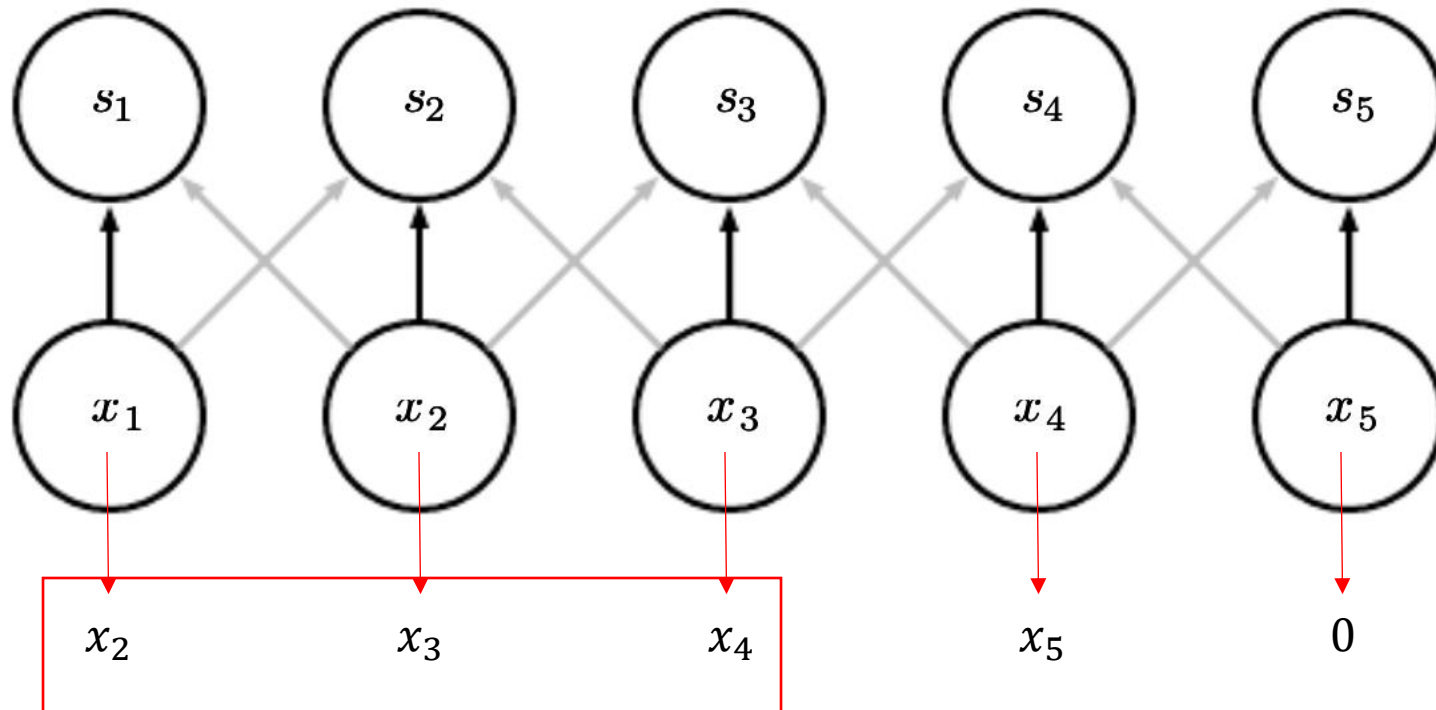
# Equivariance



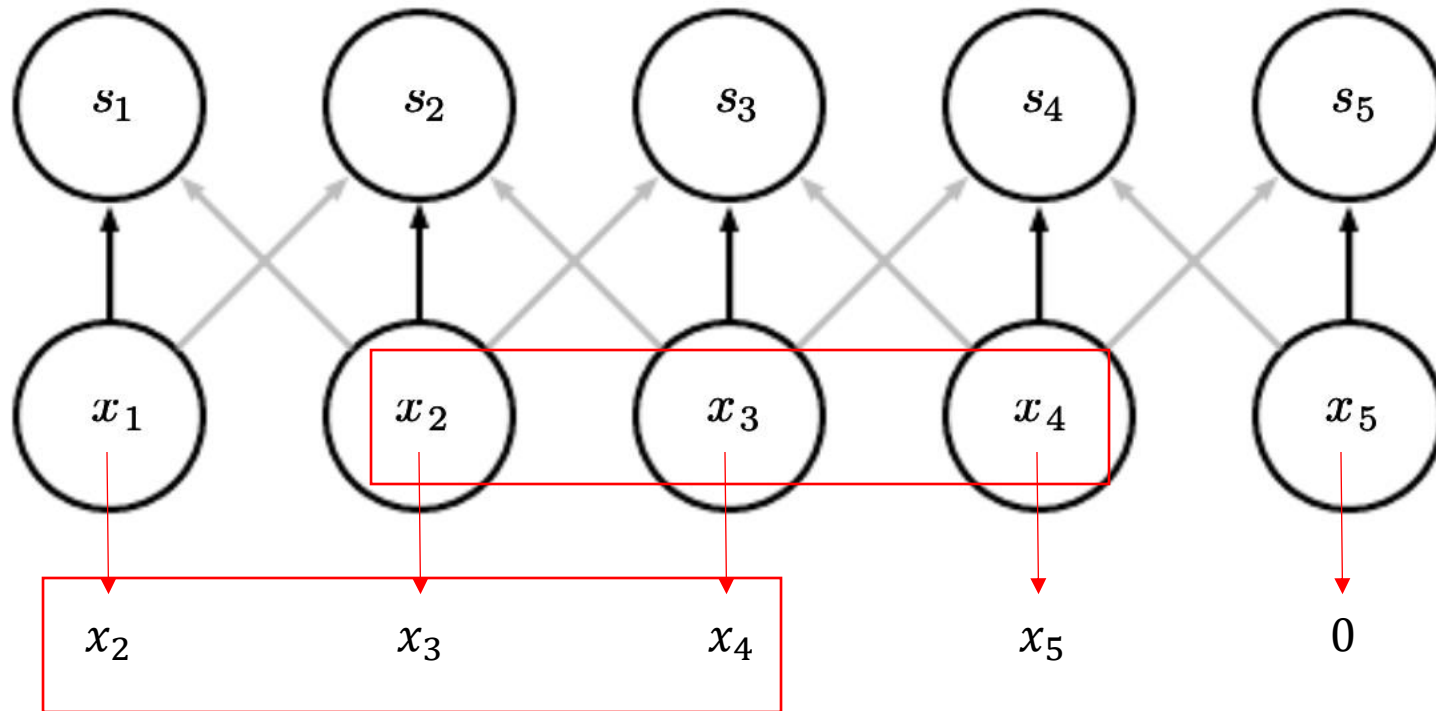
# Equivariance



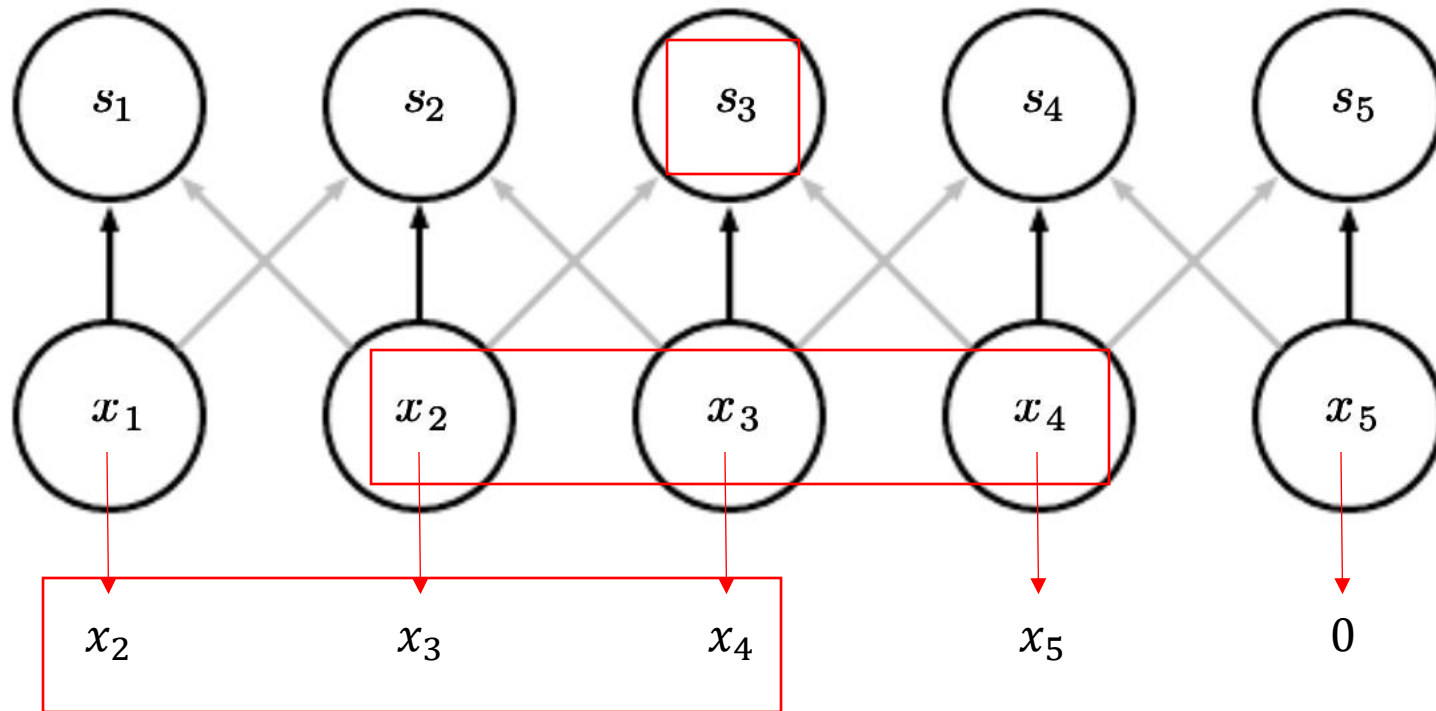
# Equivariance



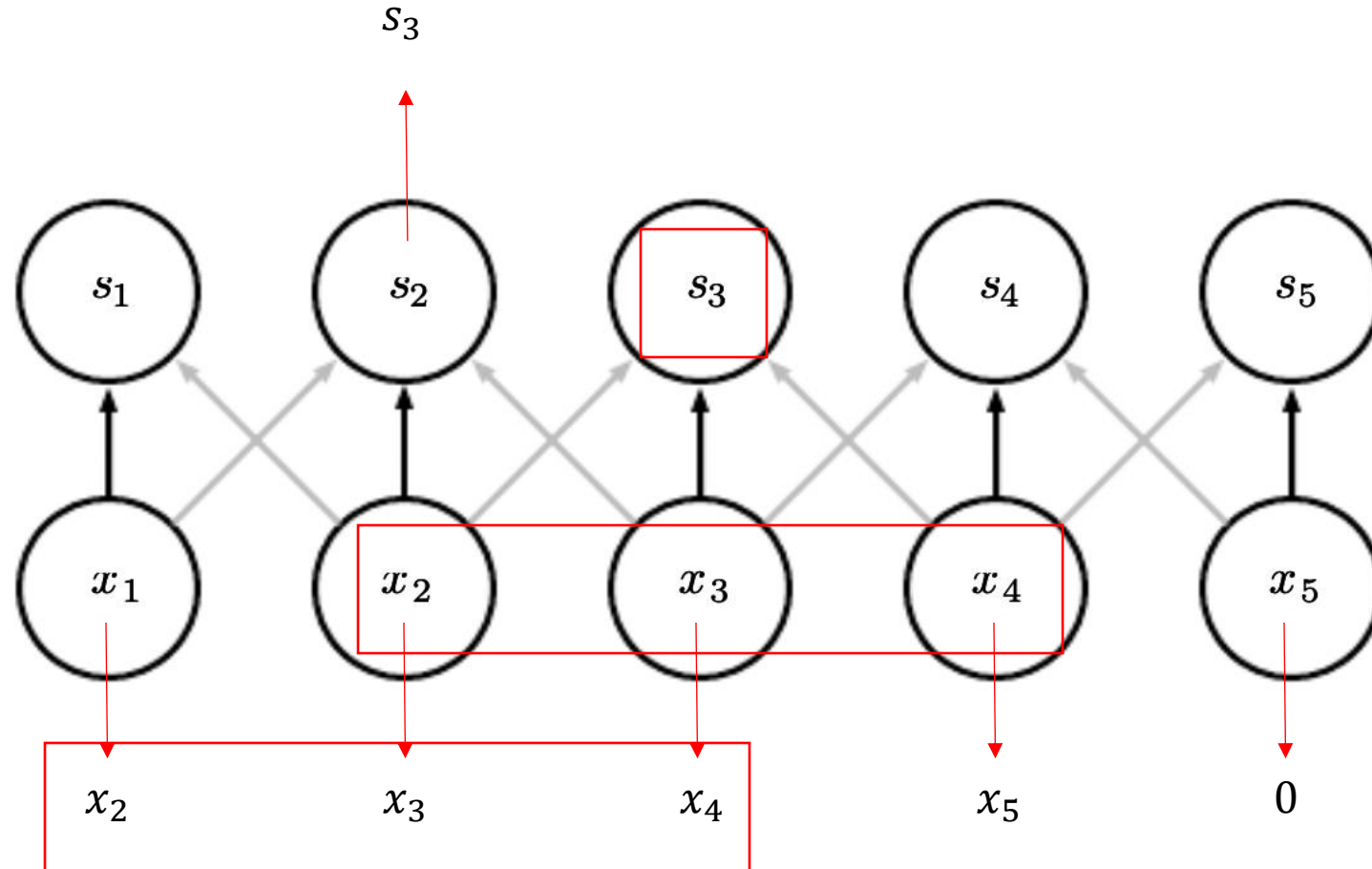
# Equivariance



# Equivariance

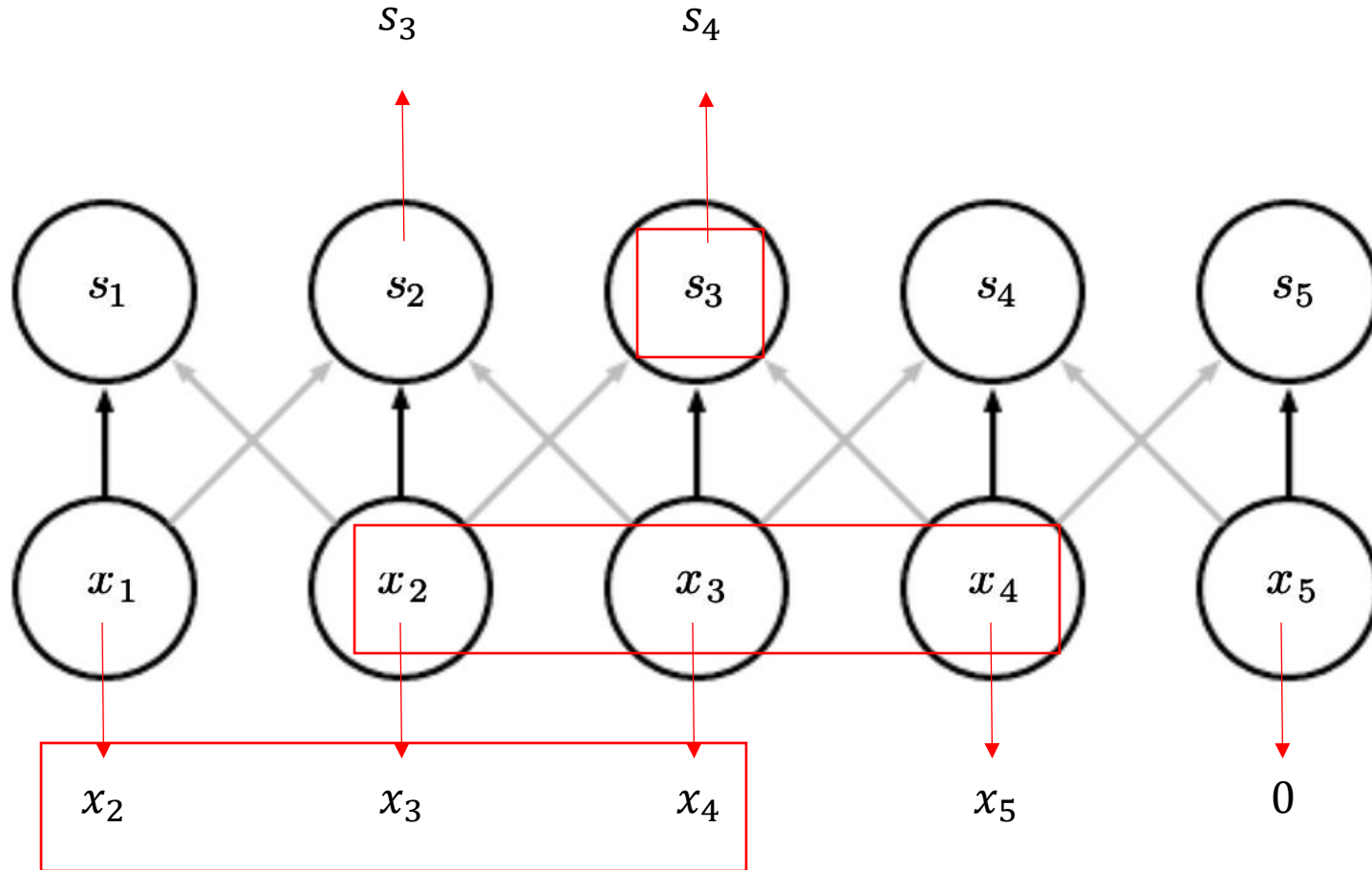


# Equivariance

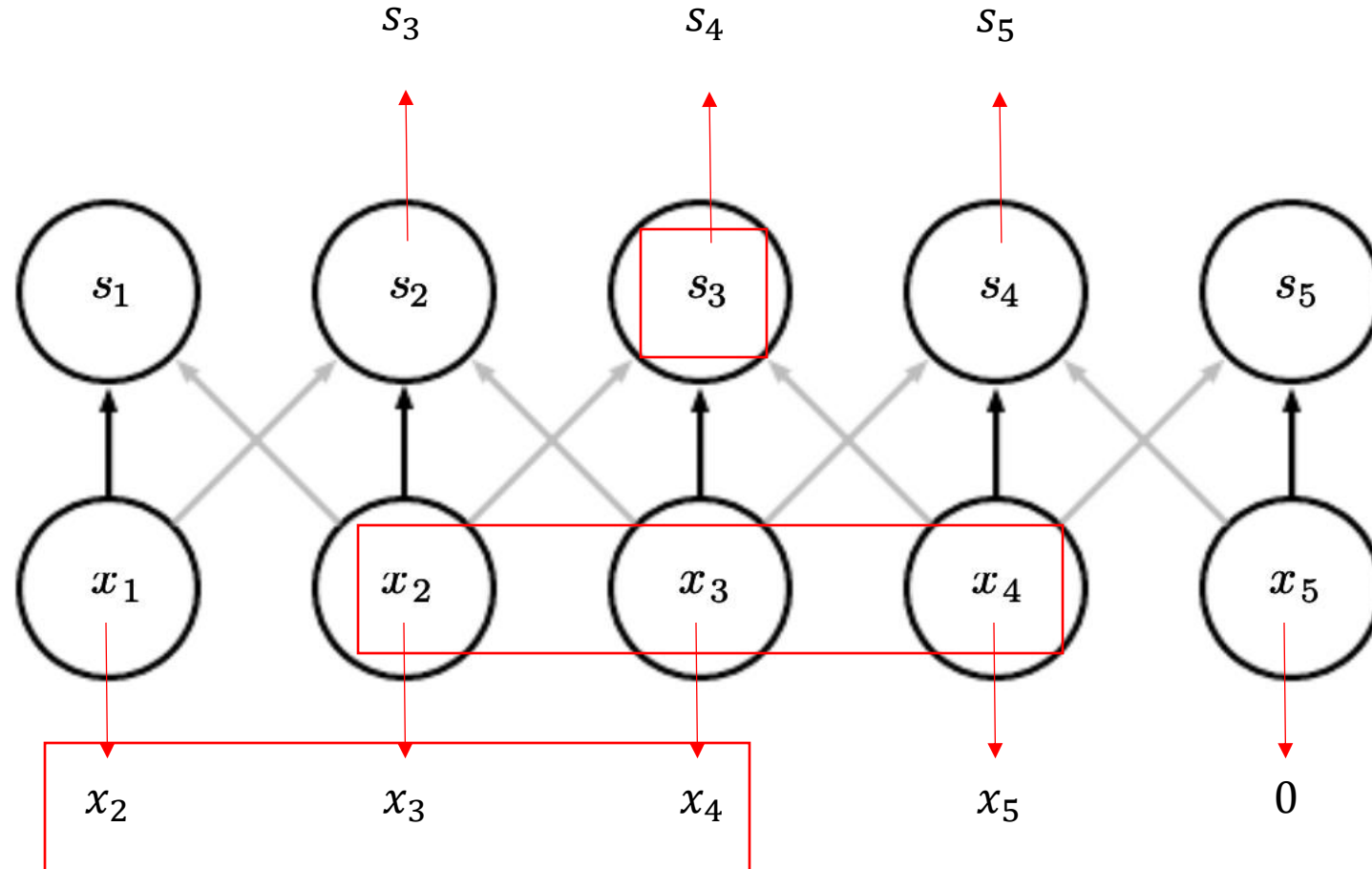




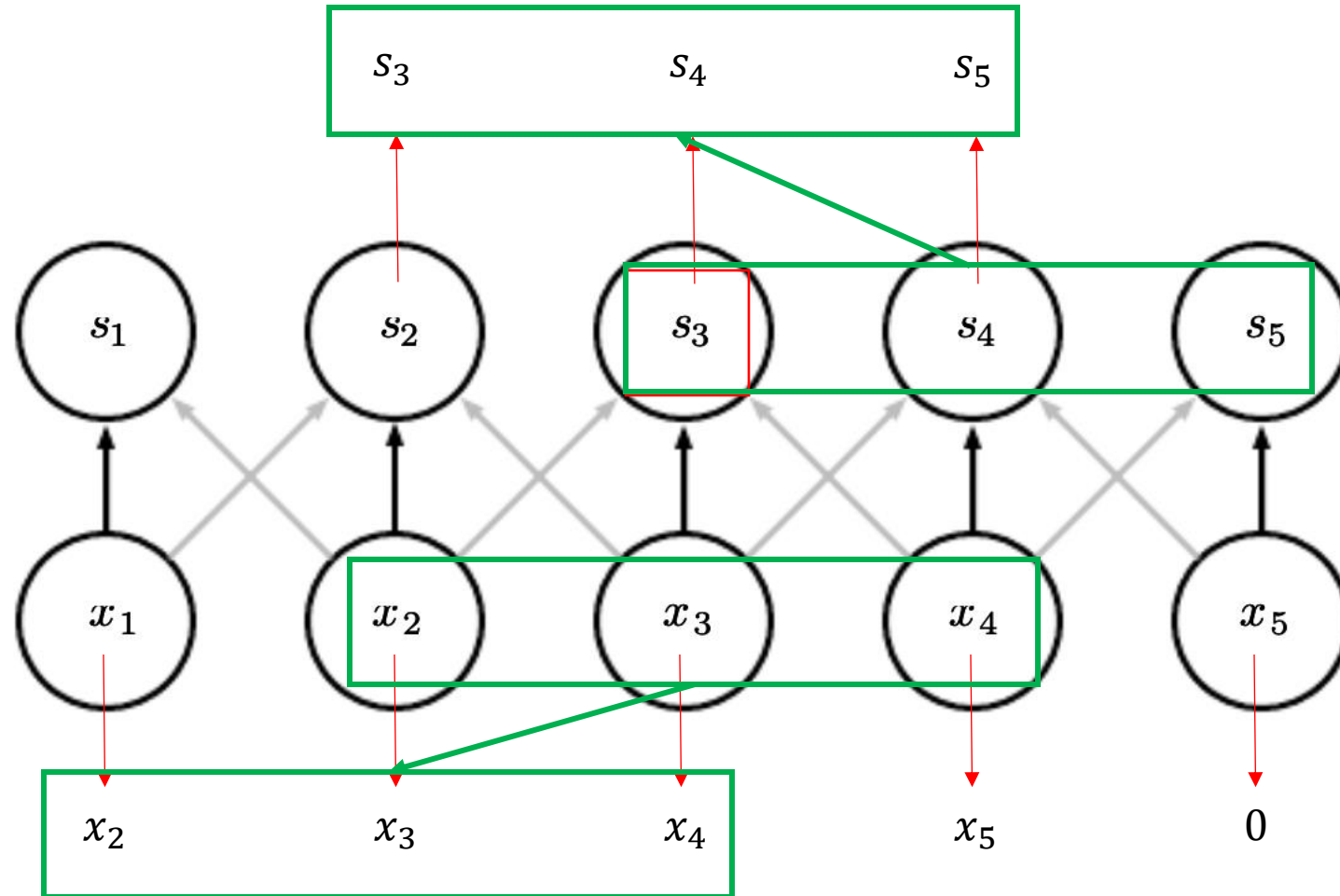
# Equivariance



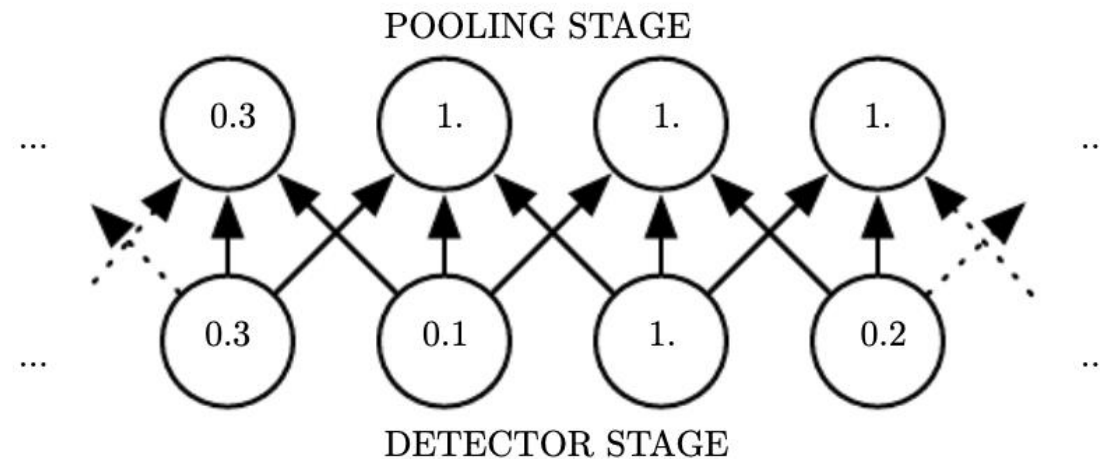
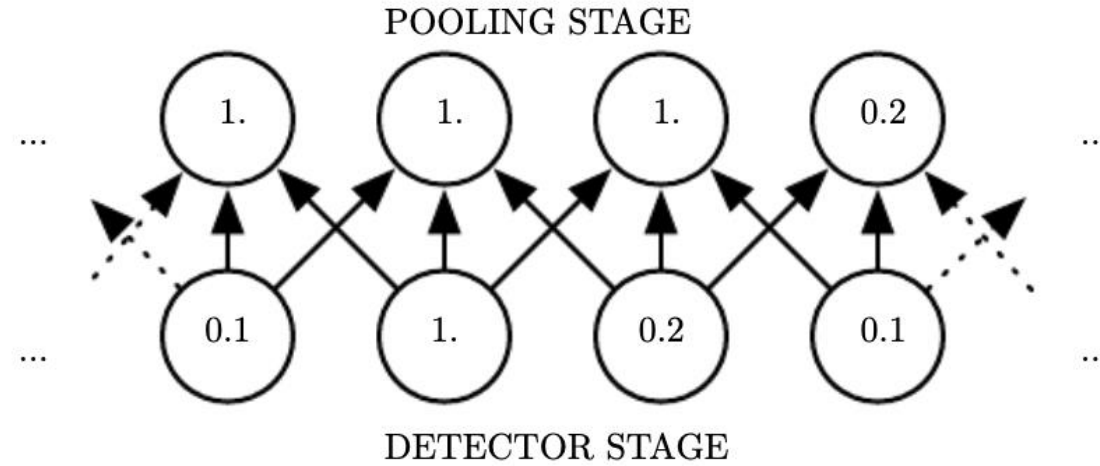
# Equivariance



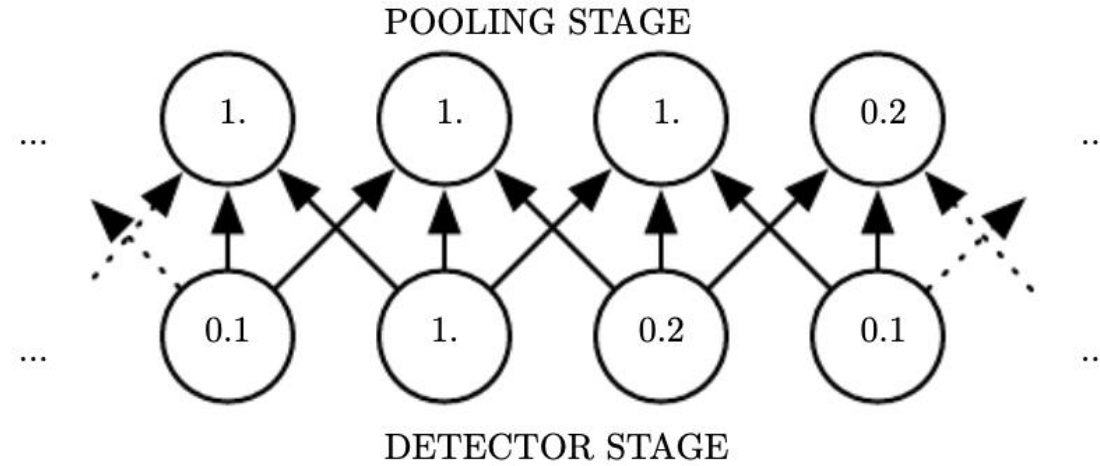
# Equivariance



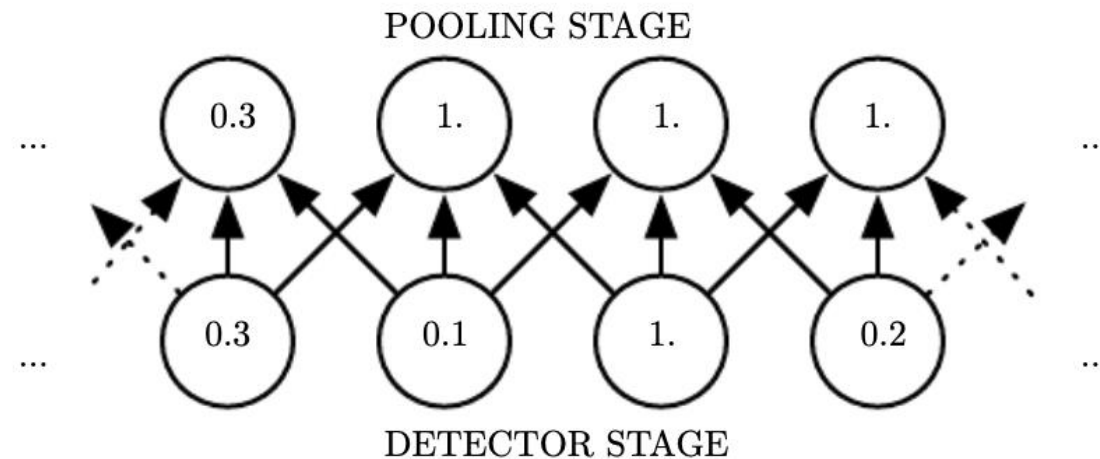
# Pooling: invariance to small translation



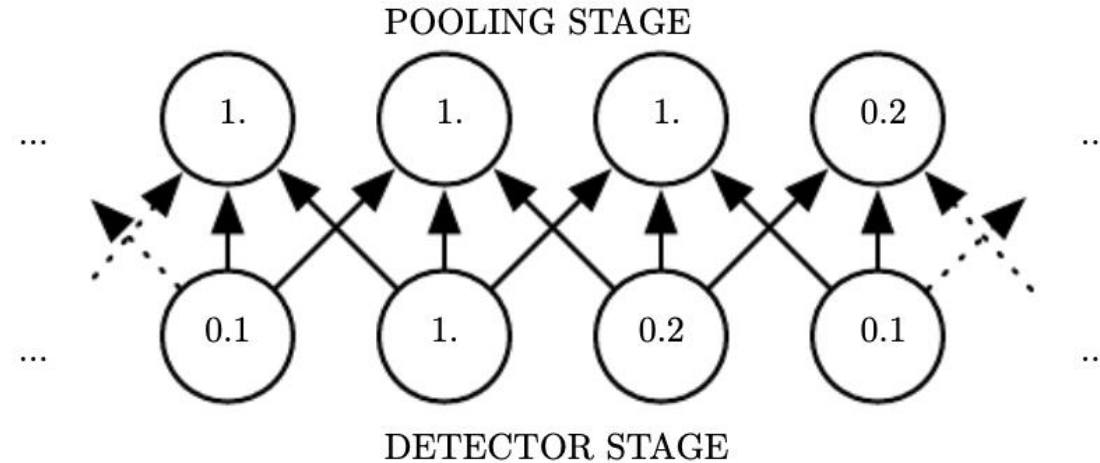
# Pooling: invariance to small translation



Q: what is type of pooling?  
Max or average pooling?

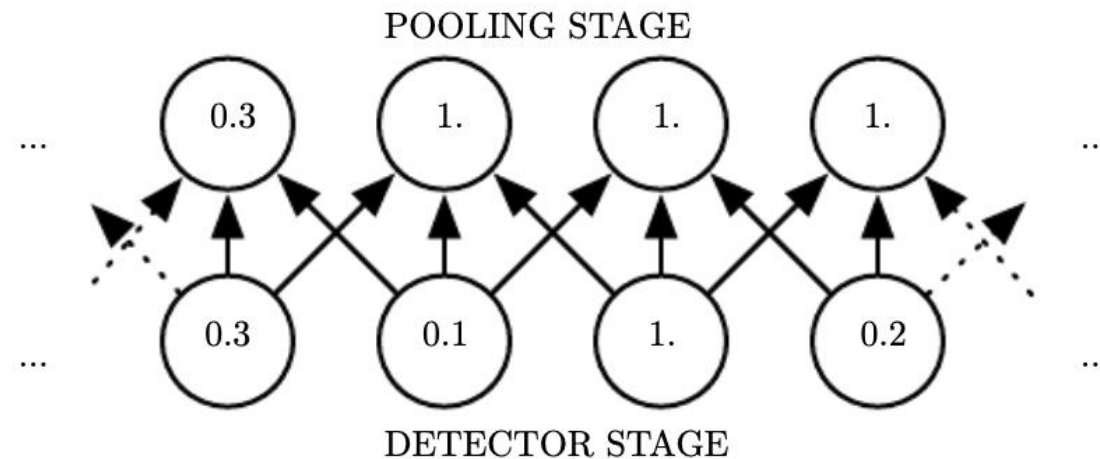


# Pooling: invariance to small translation

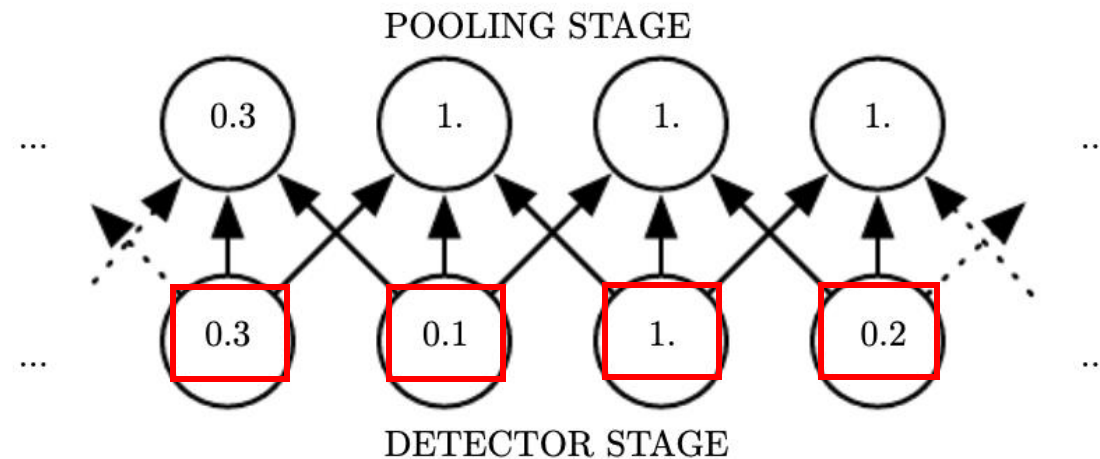
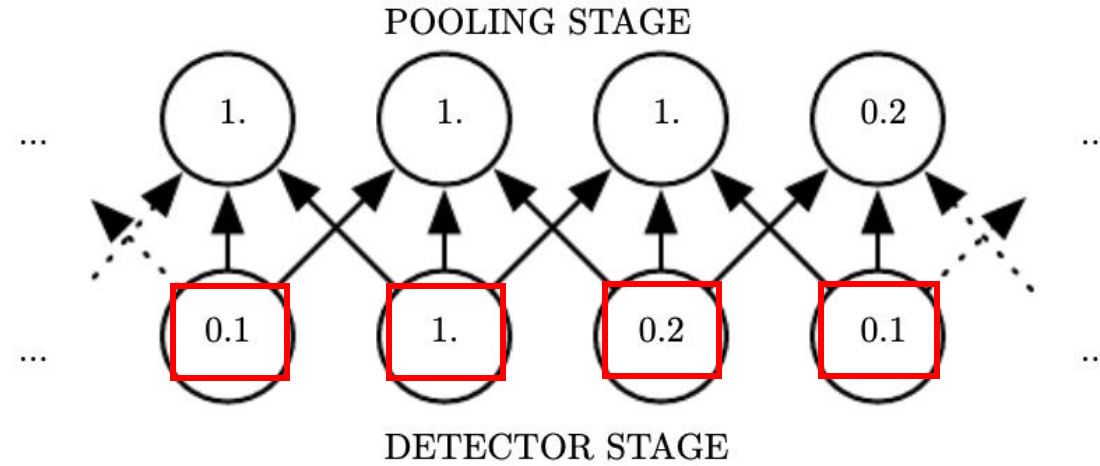


Q: what is type of pooling?  
Max or average pooling?

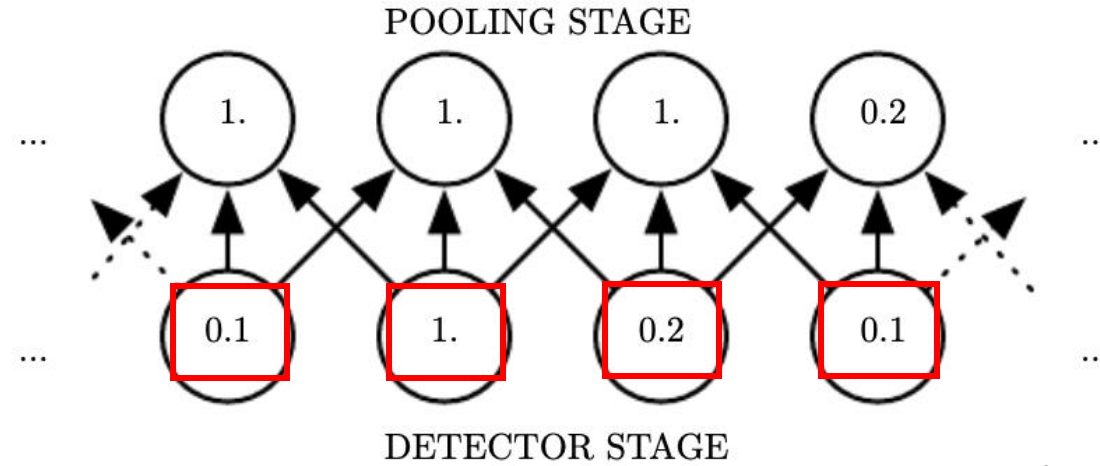
Max pooling



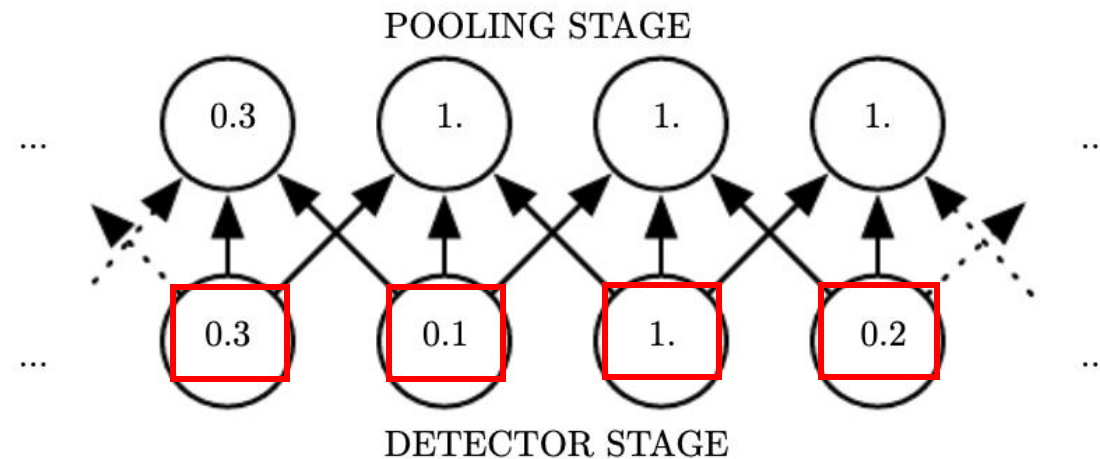
# Pooling: invariance to small translation



# Pooling: invariance to small translation

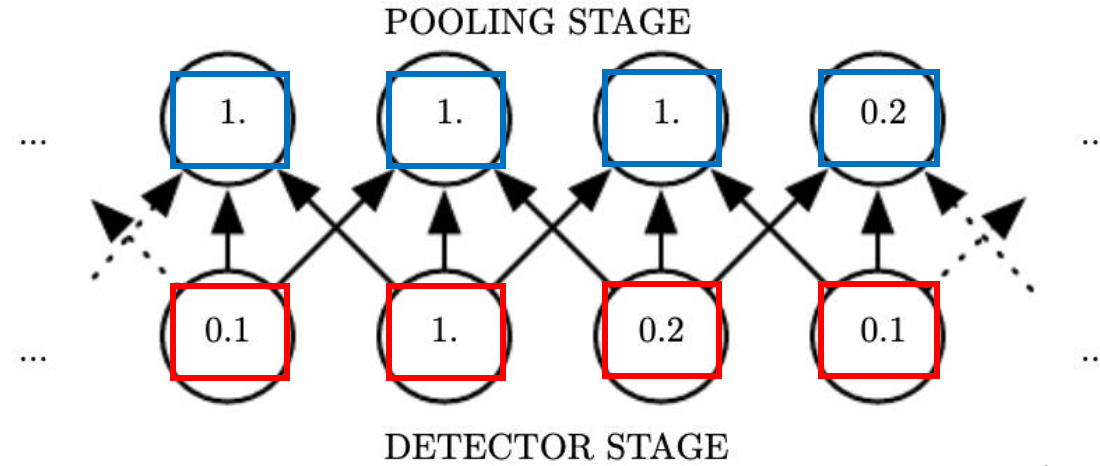


Translate: change the input value a little bit  
+ change their positions

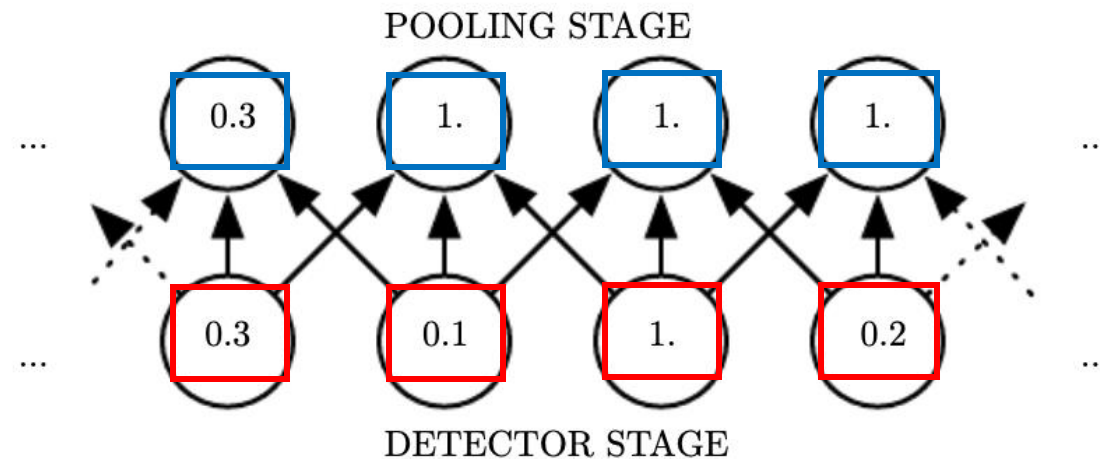




# Pooling: invariance to small translation



Translate: change the input value a little bit  
+ change their positions



Seagull = 1?



Seagull = 1?



Seagull is in the center

Seagull = 1?





Seagull = 1?



Seagulls are present, but not in the center

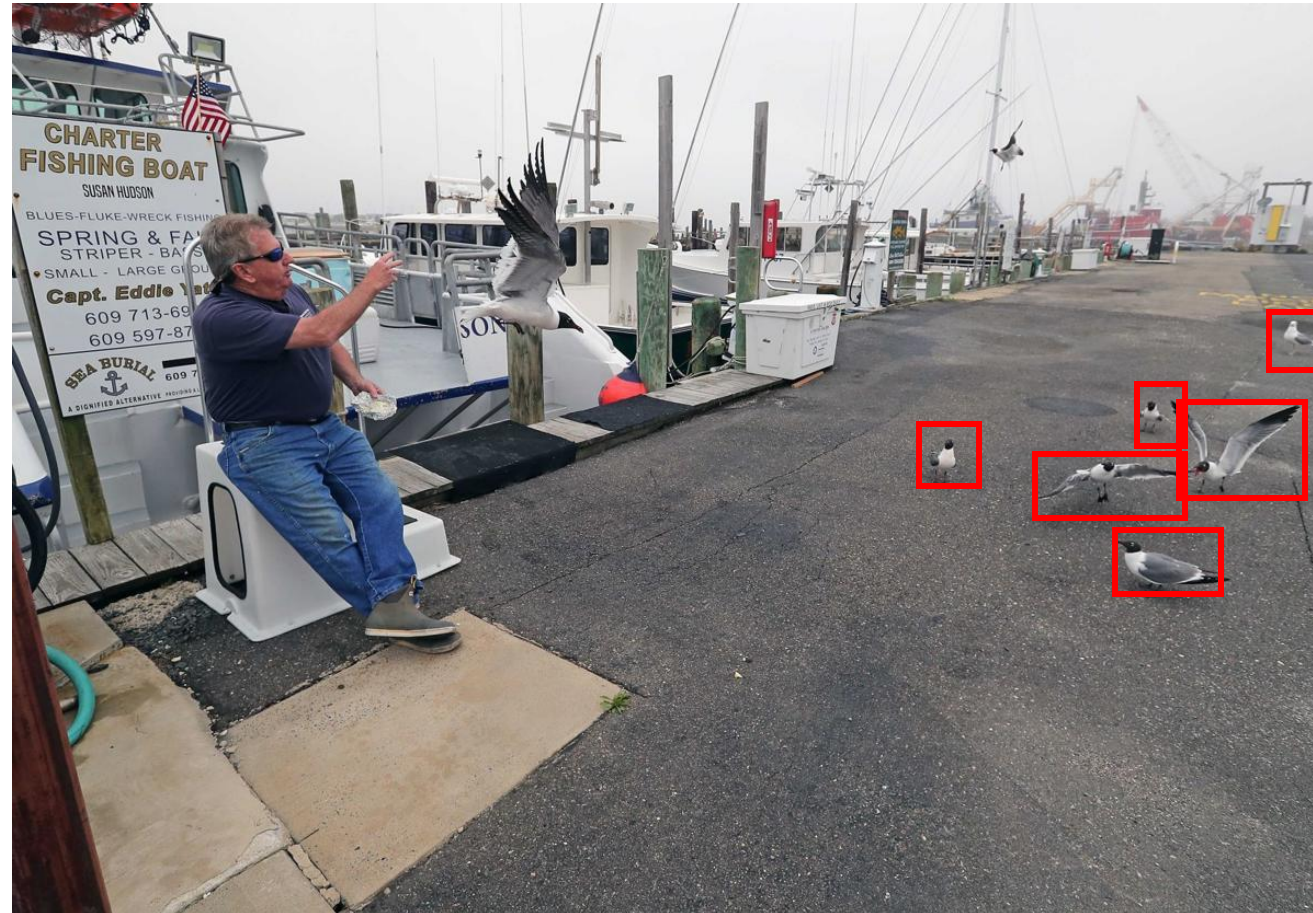
Seagull = 1?



Seagulls are present, but not in the center



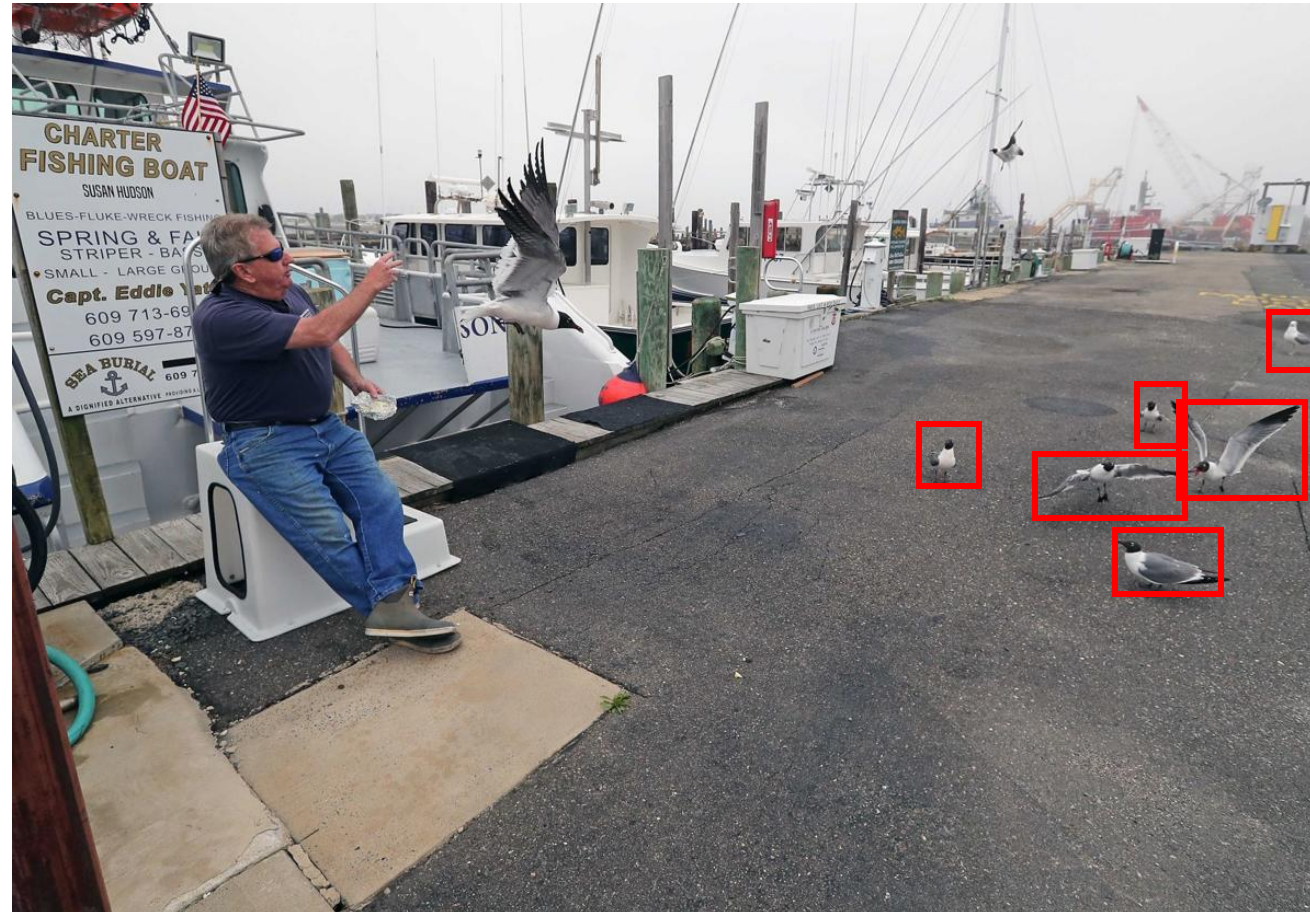
# Seagull = 1?



CNNs can tell:  
whether seagulls are present  
Not tell:  
their positions in the image

Seagulls are present, but not in the center

# Seagull = 1?

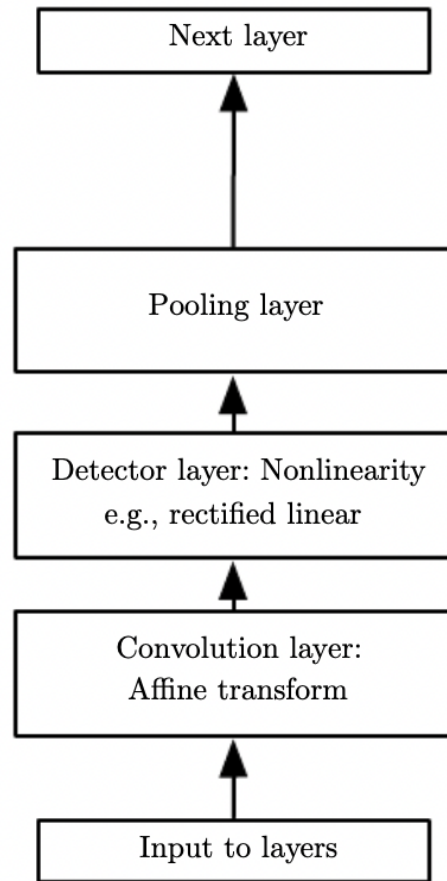


CNNs can tell:  
whether seagulls are present  
Not tell:  
their positions in the image  
Invariant to small translations

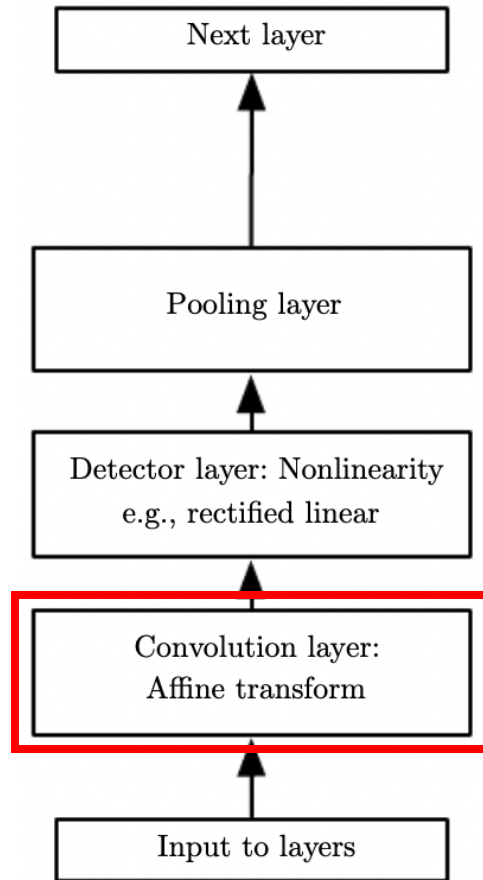
Seagulls are present, but not in the center



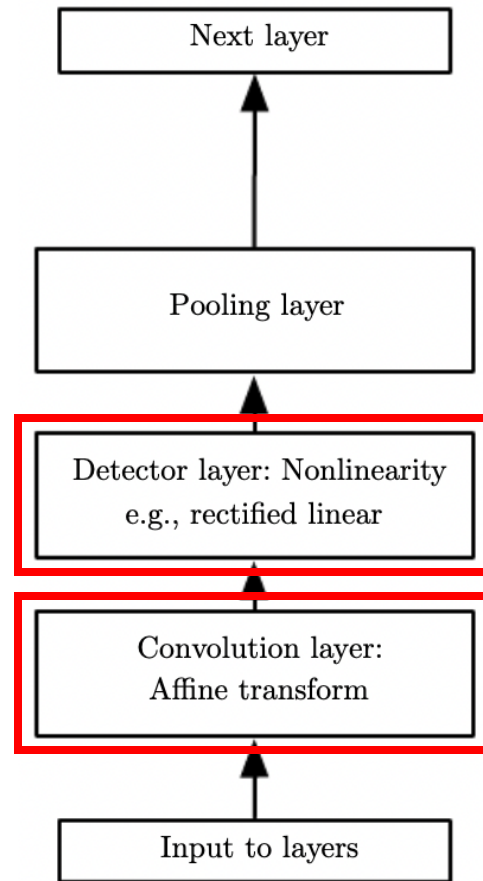
# A typical convolutional layer



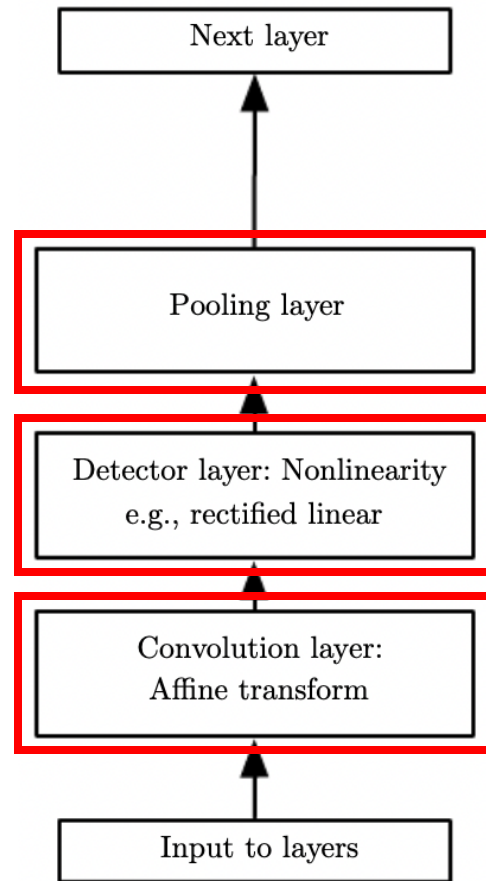
# A typical convolutional layer



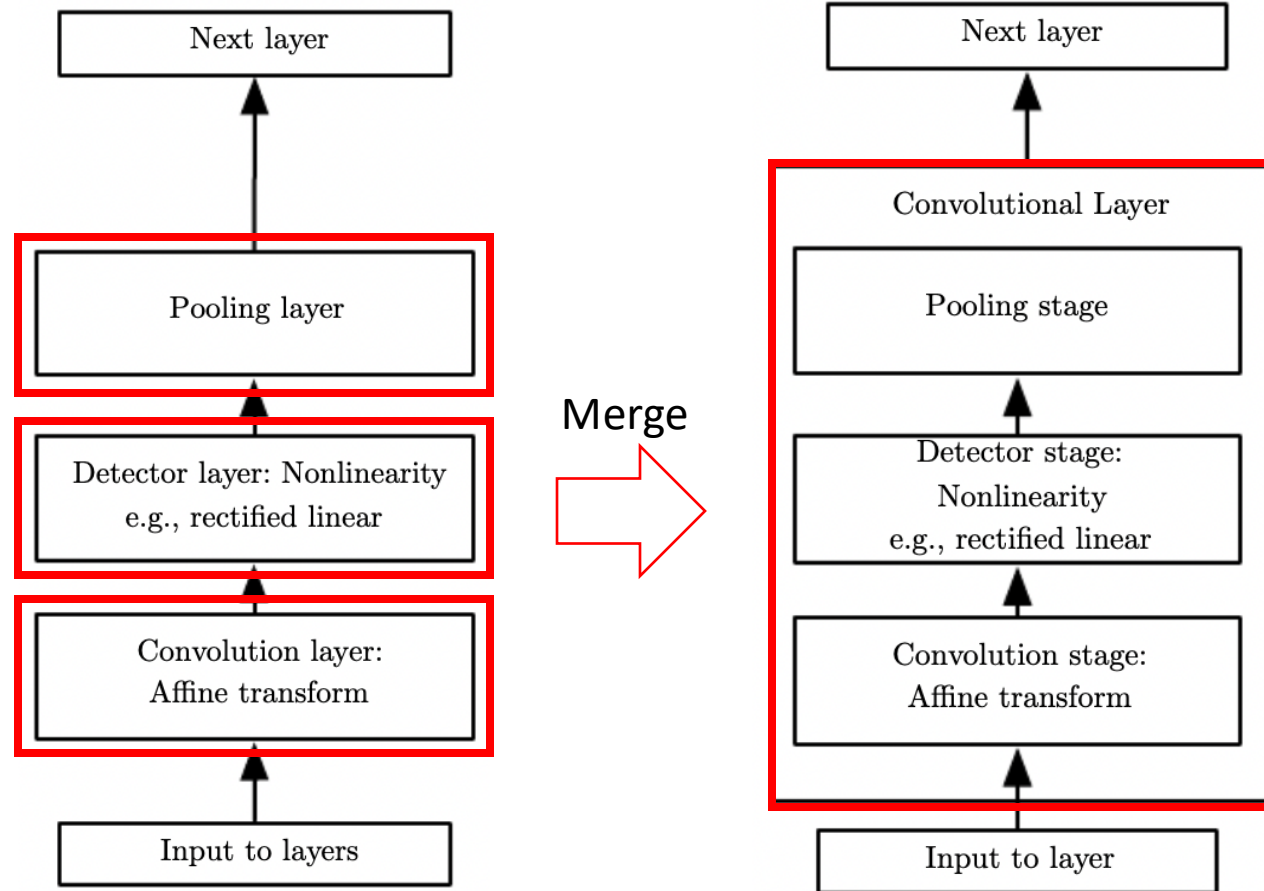
# A typical convolutional layer

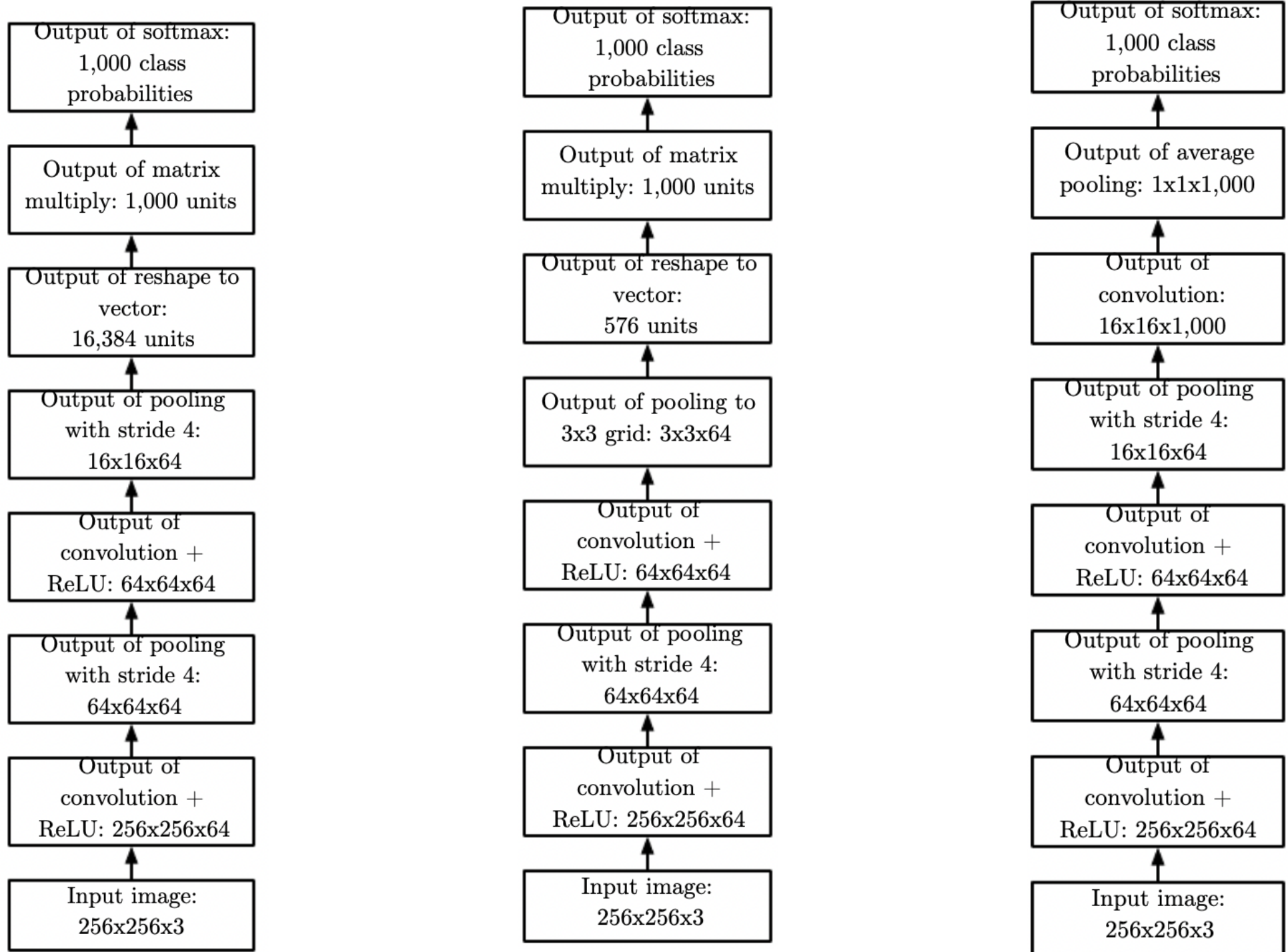


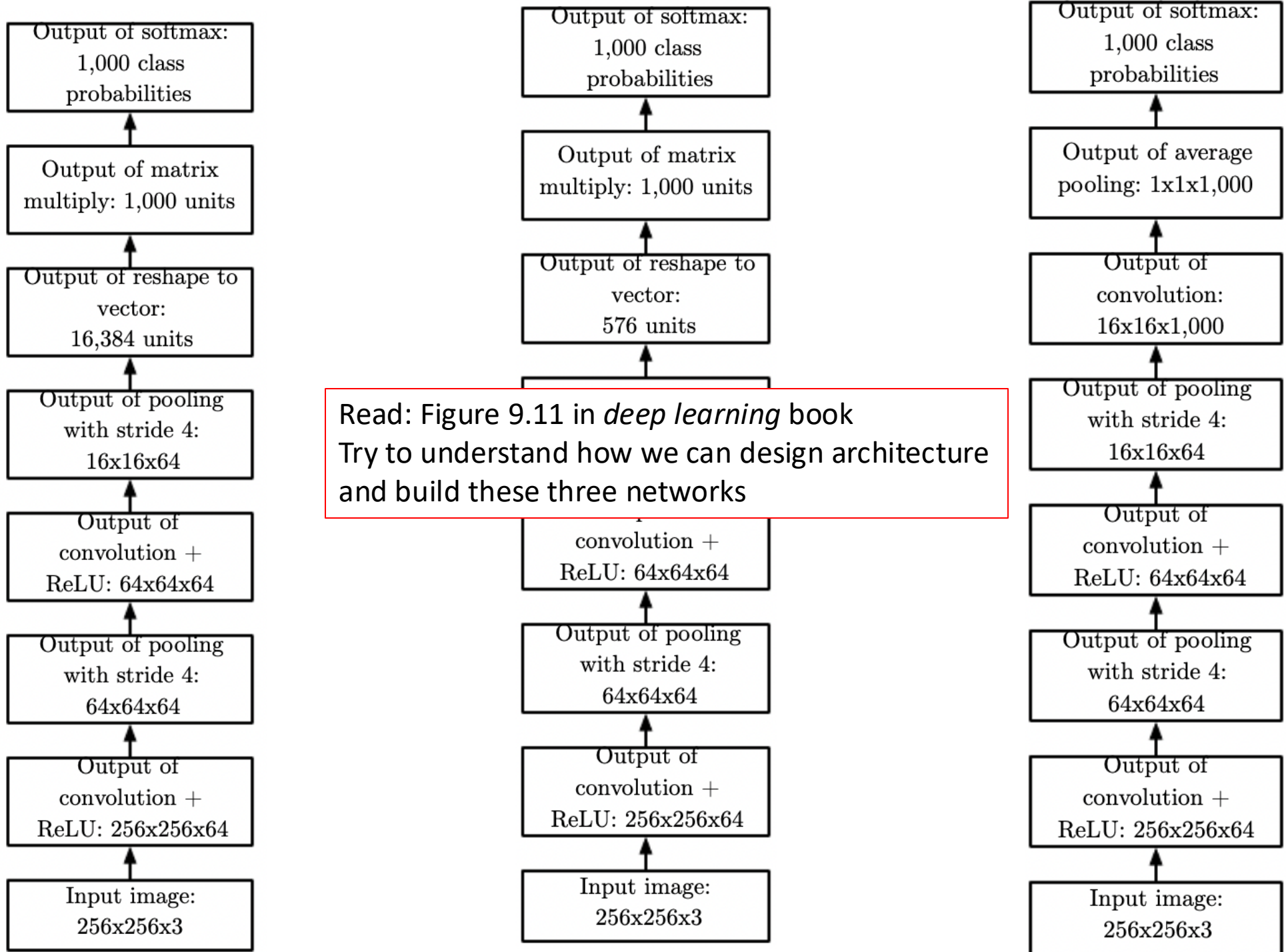
# A typical convolutional layer



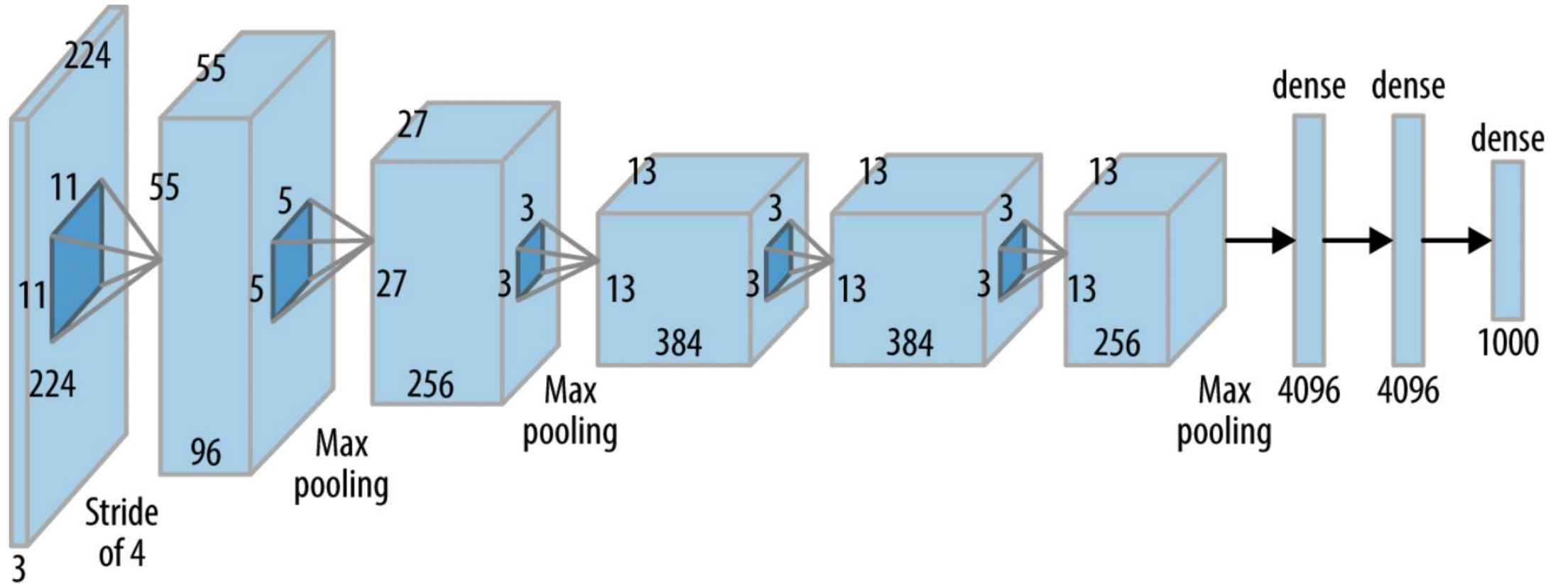
# A typical convolutional layer





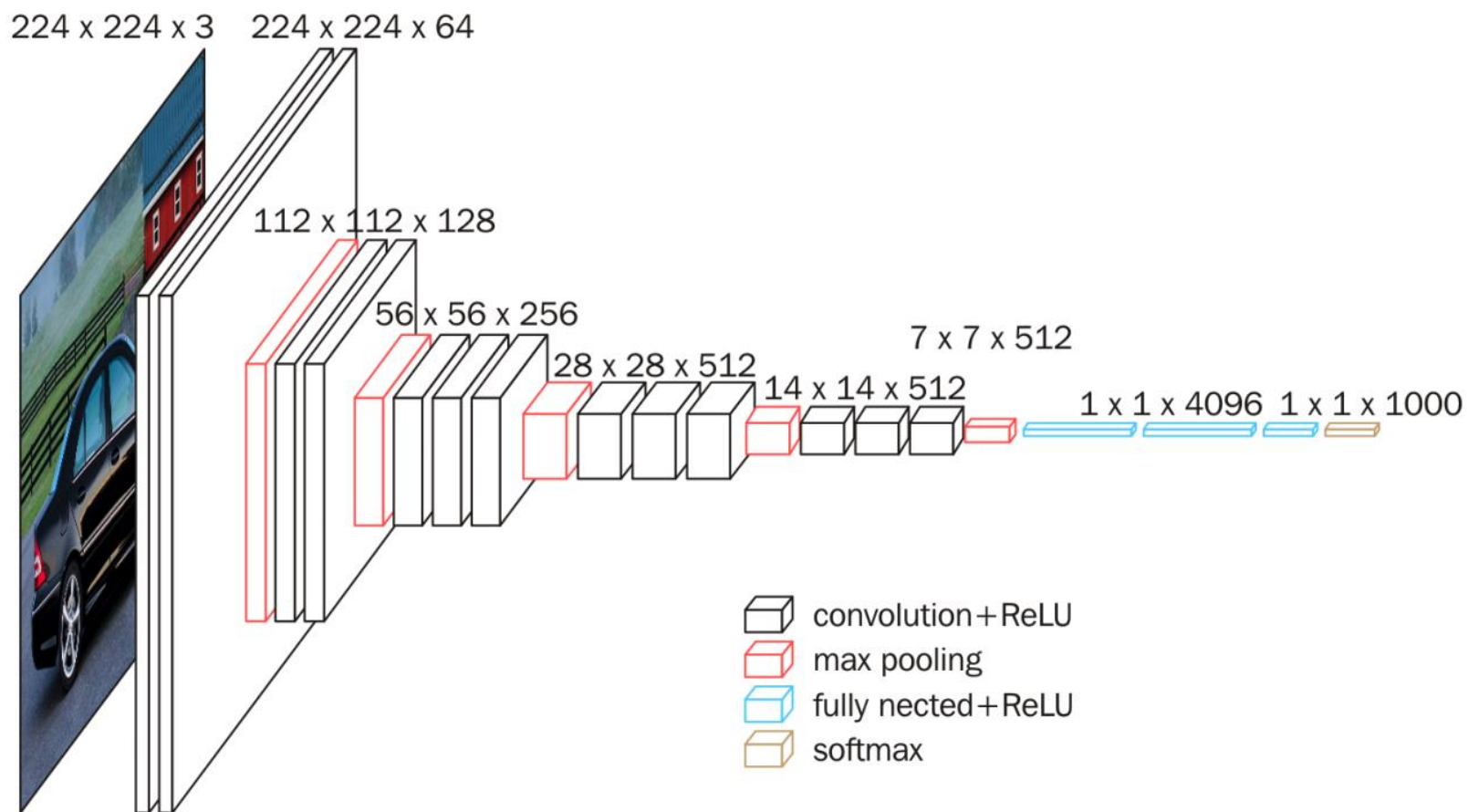


# AlexNet

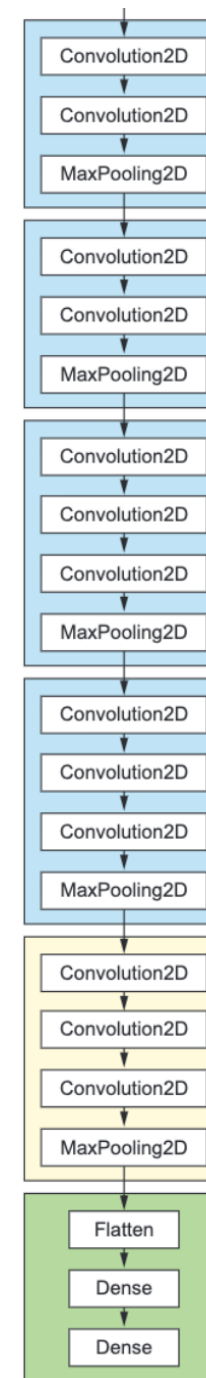




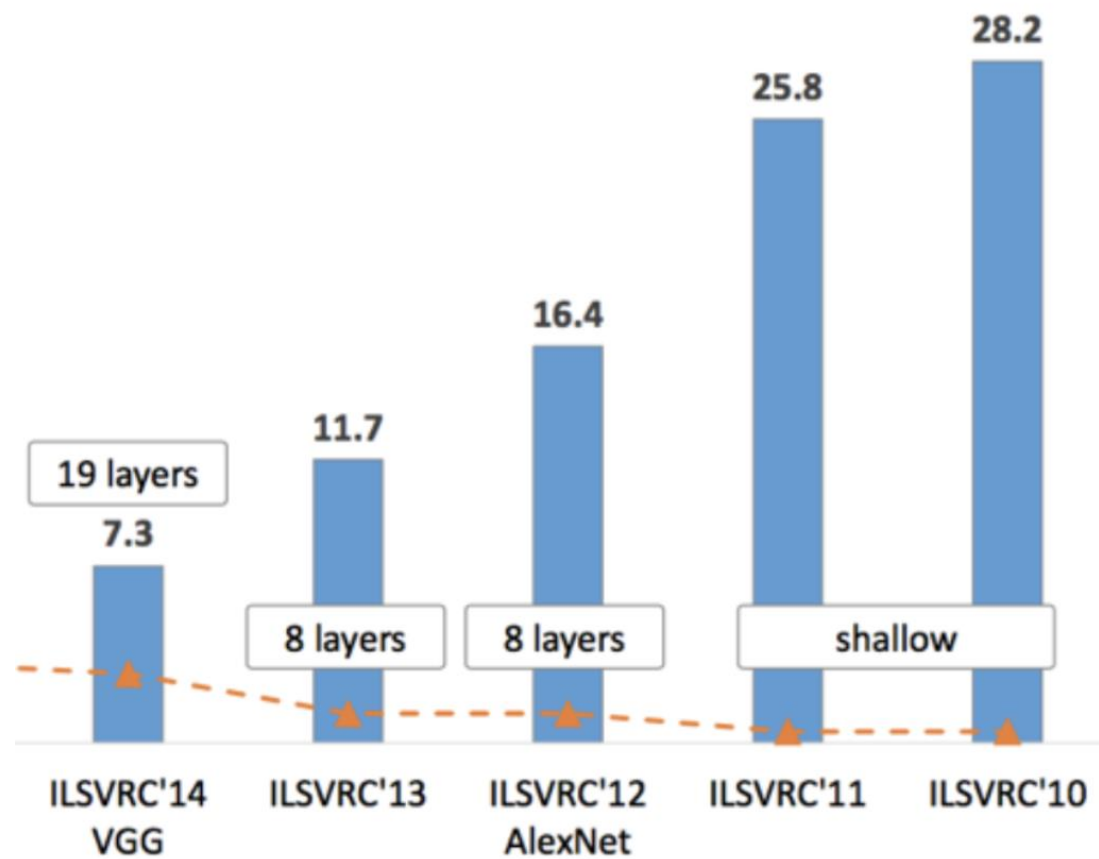
# VGG-16



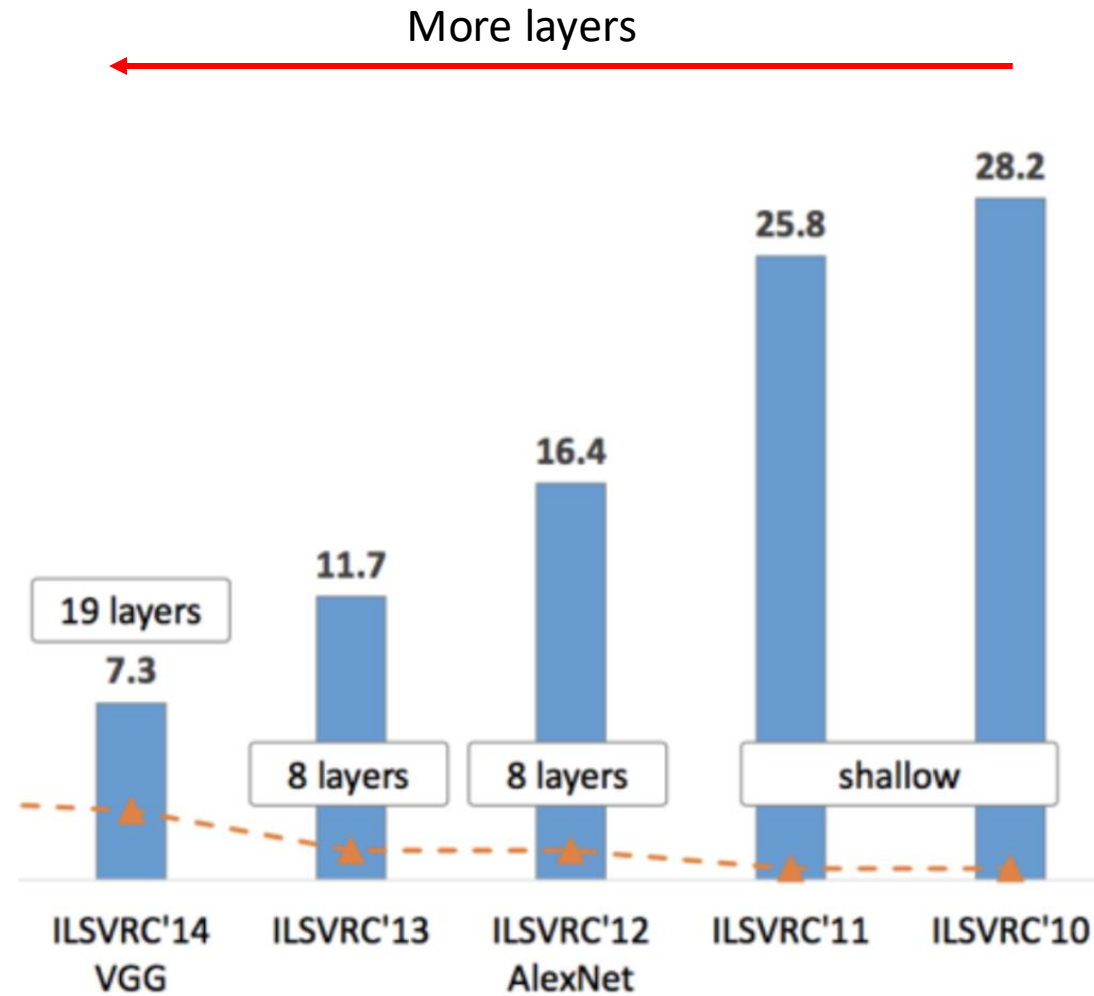
[VGG]



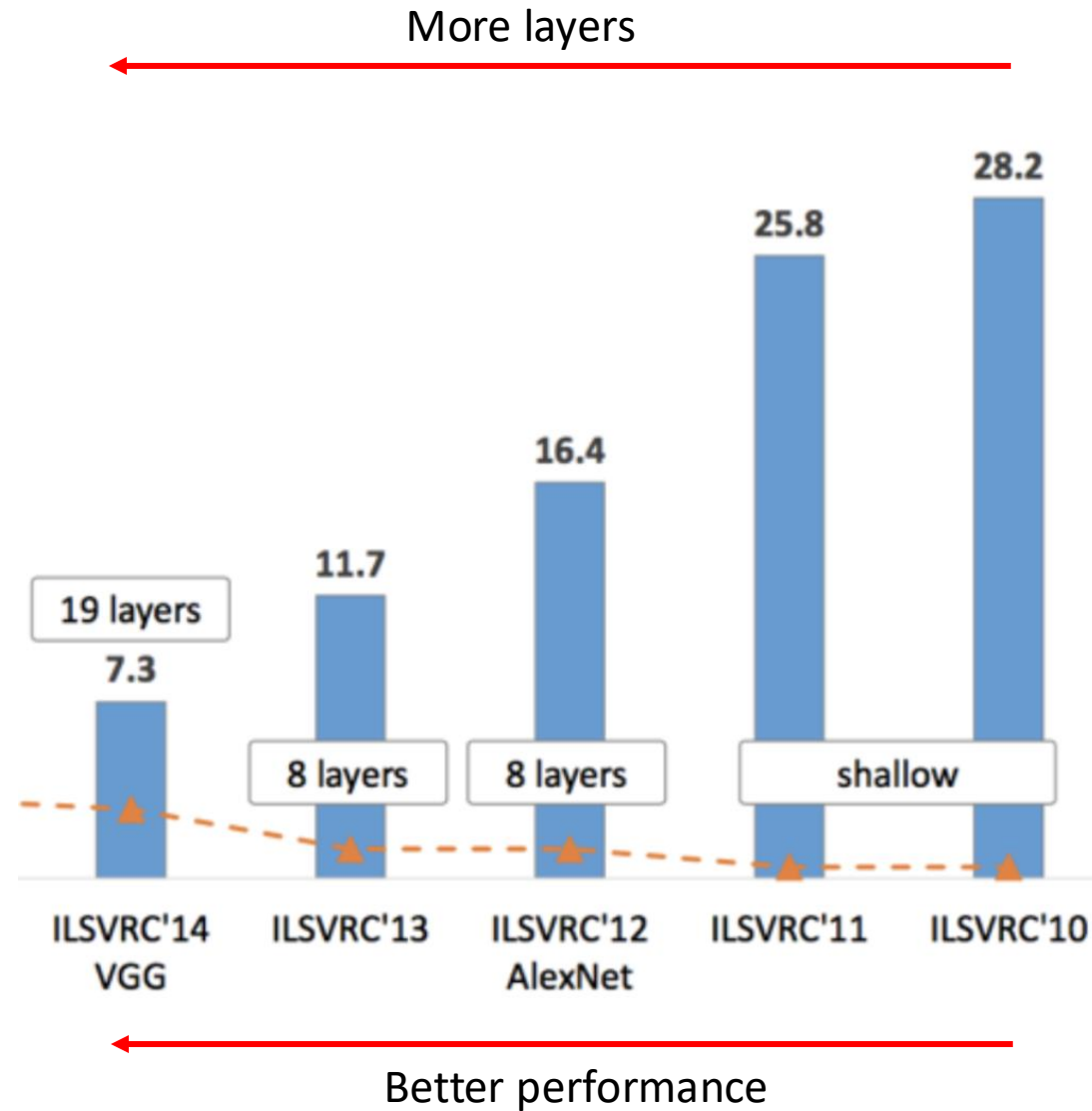
# ImageNet competition



# ImageNet competition



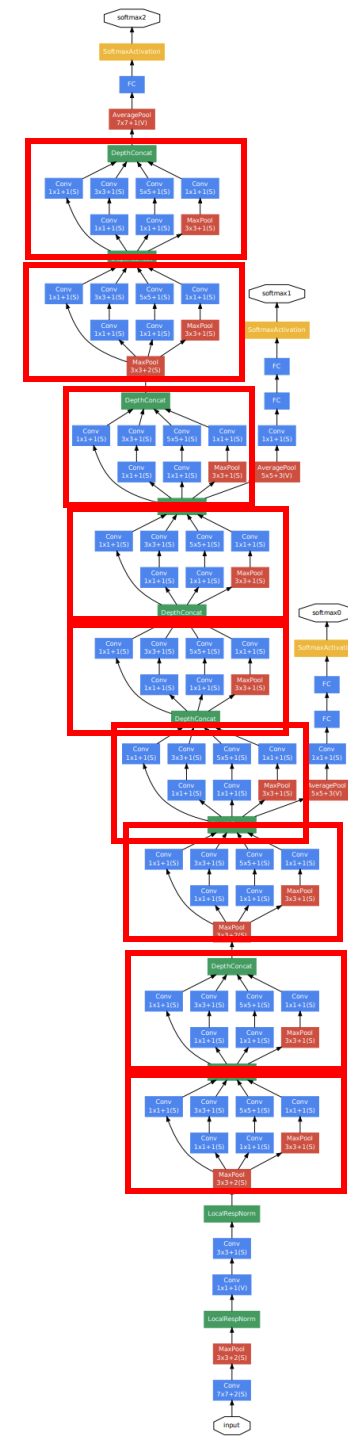
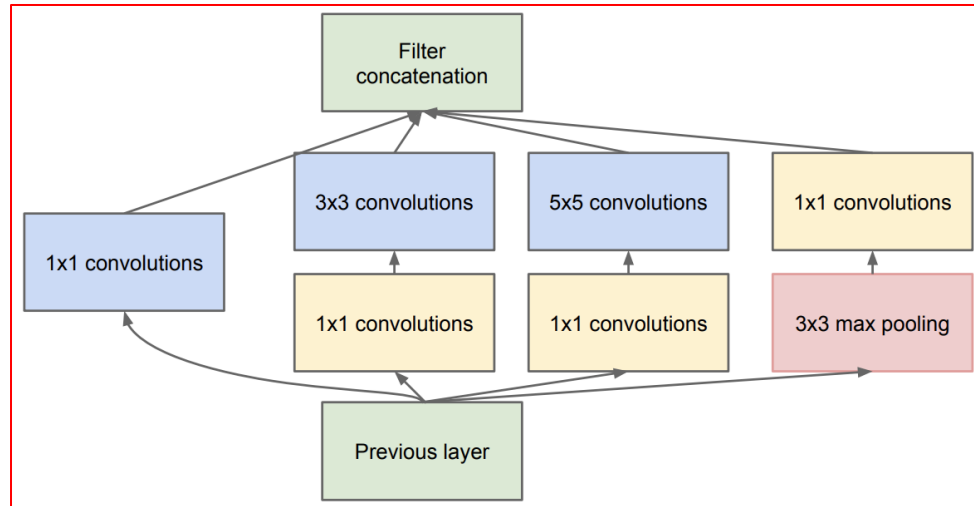
# ImageNet competition



# Inception (GoogLeNet)

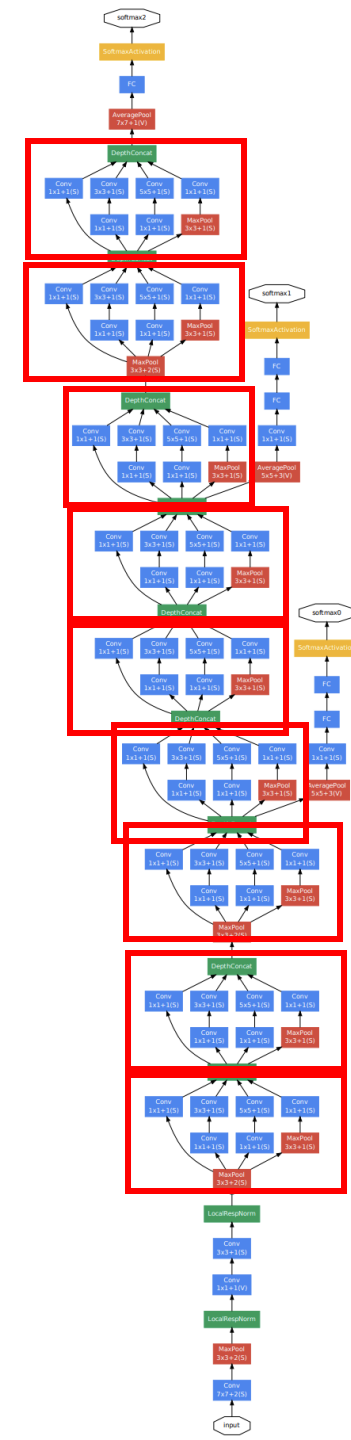
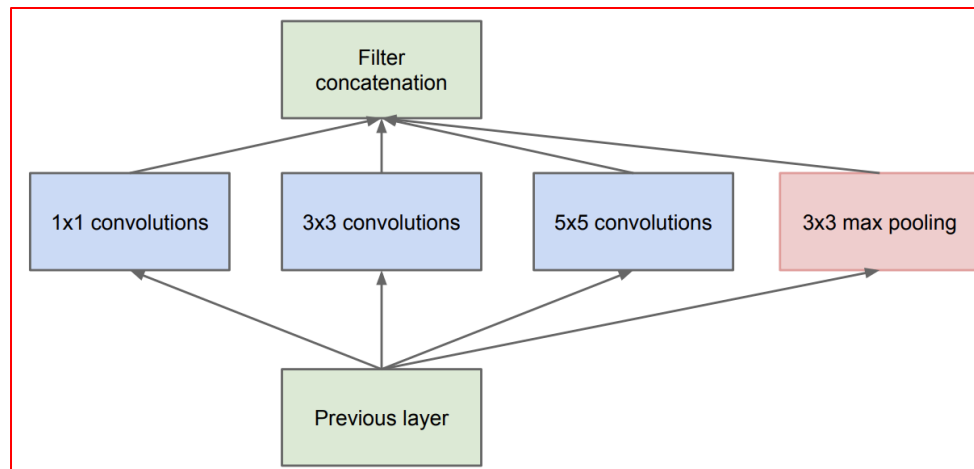
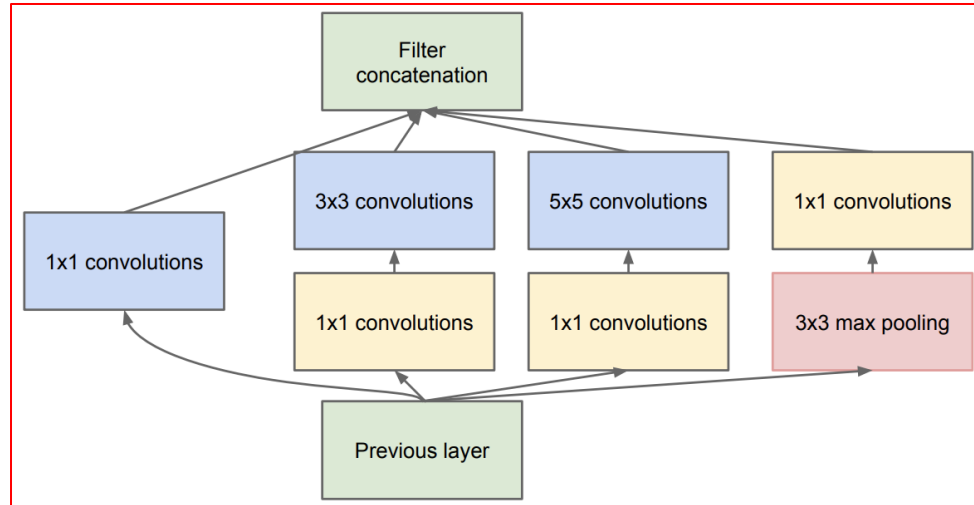


# Inception (GoogLeNet)



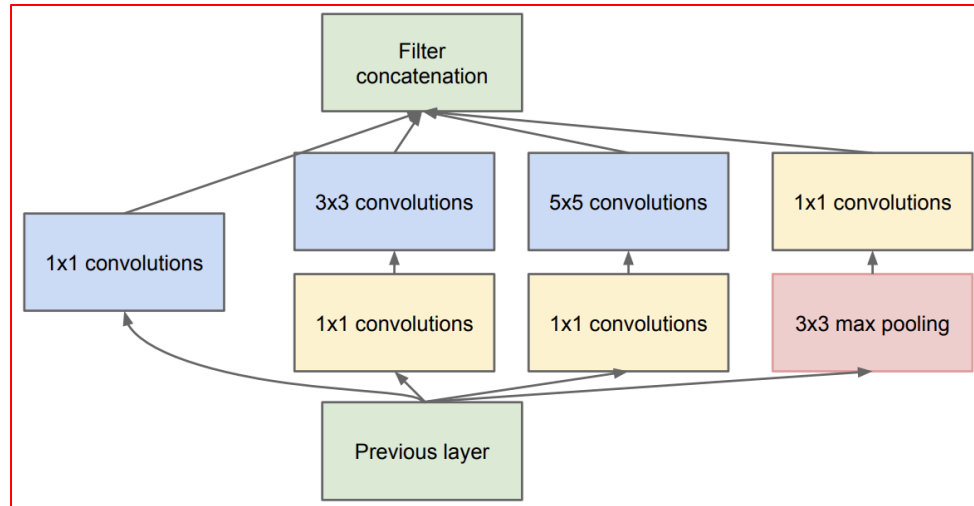


# Inception (GoogLeNet)

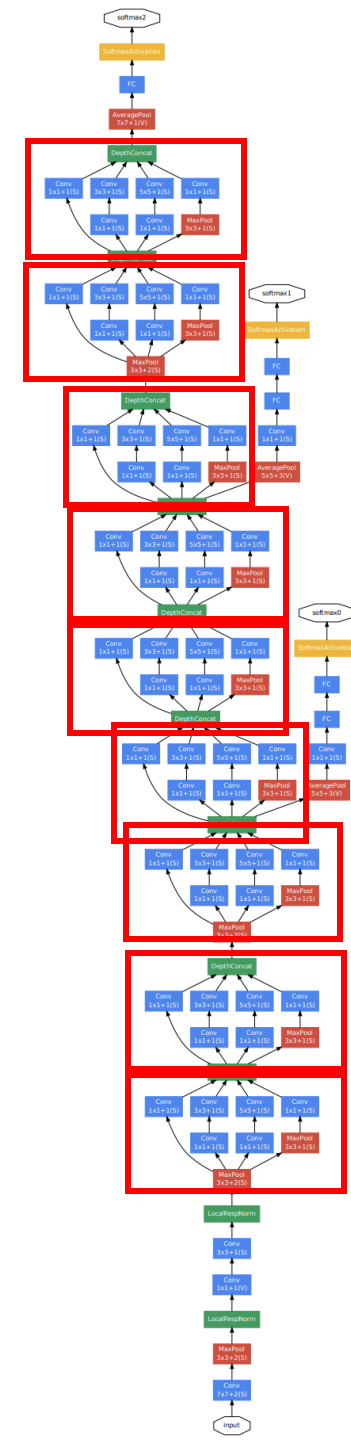
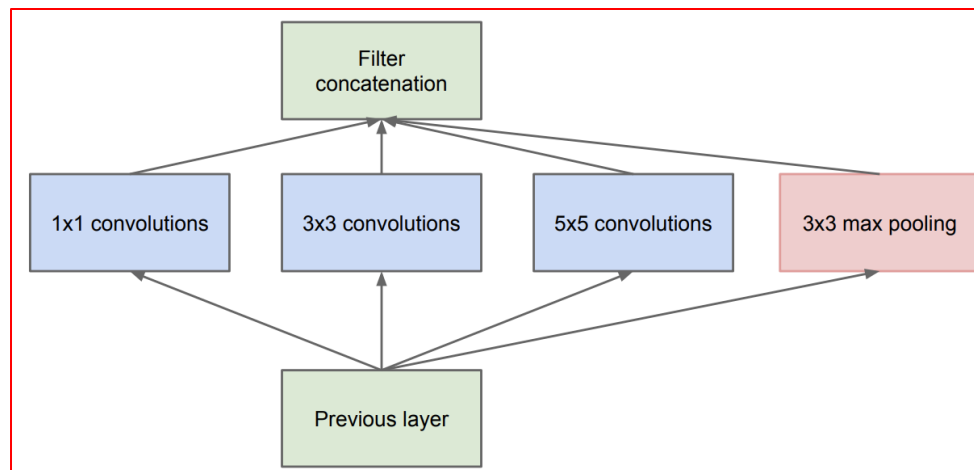




# Inception (GoogLeNet)

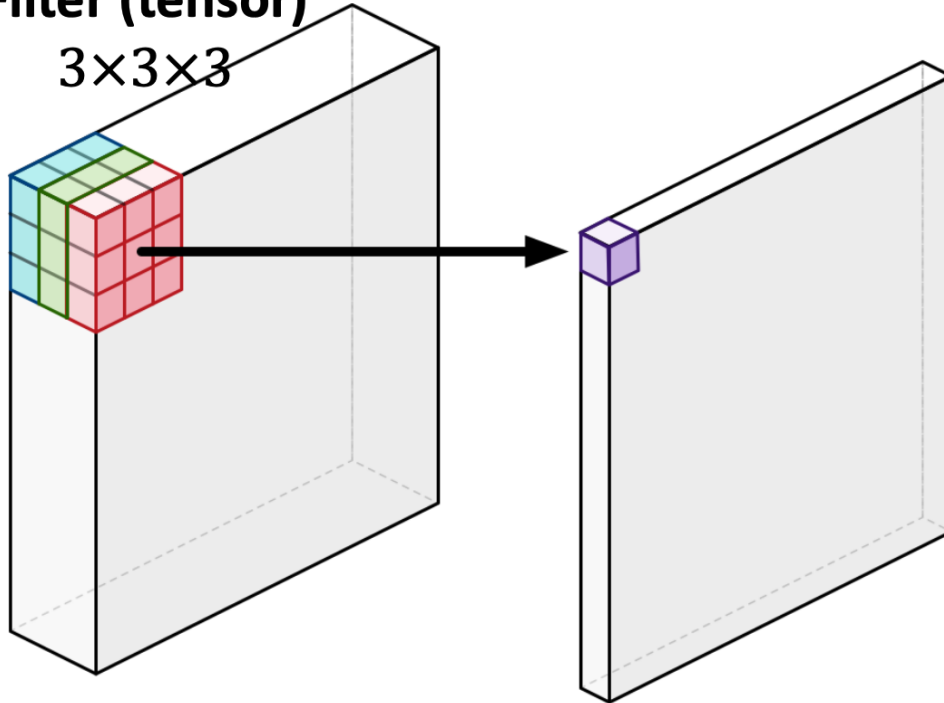


Q: difference between those two variants?



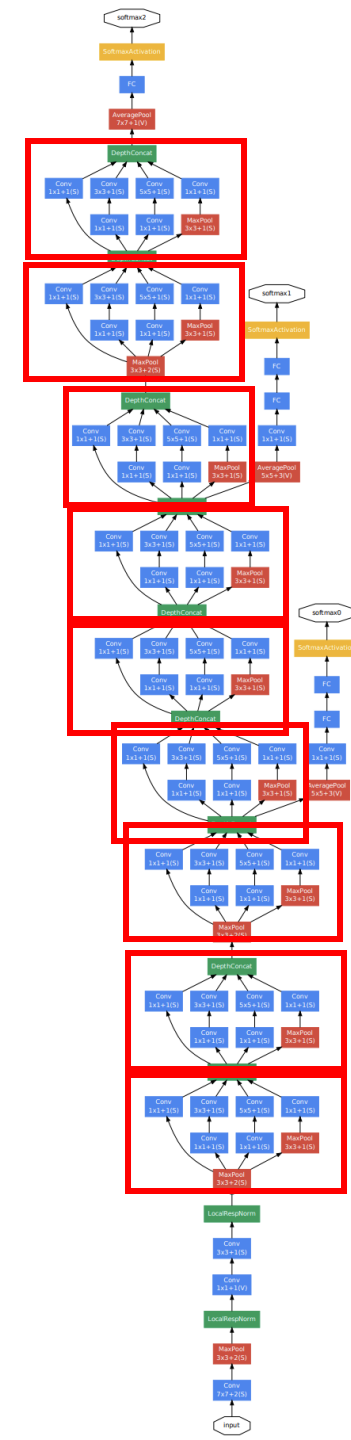
# Inception (GoogLeNet)

Filter (tensor)  
 $3 \times 3 \times 3$



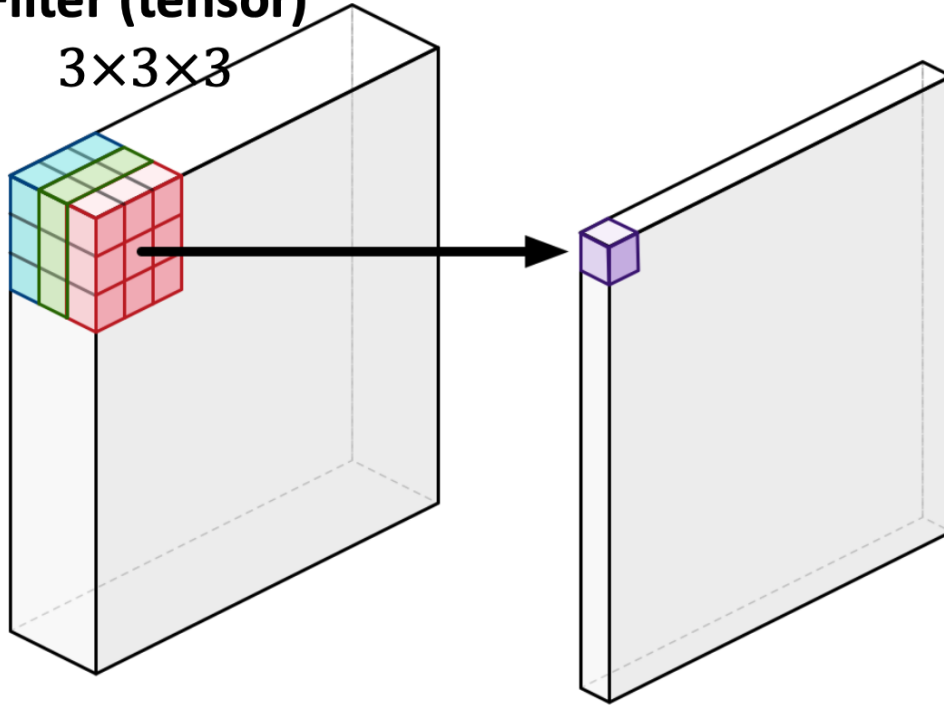
Input tensor  
 $d_1 \times d_2 \times 3$

Output matrix  
 $(d_1 - 2) \times (d_2 - 2)$



# Inception (GoogLeNet)

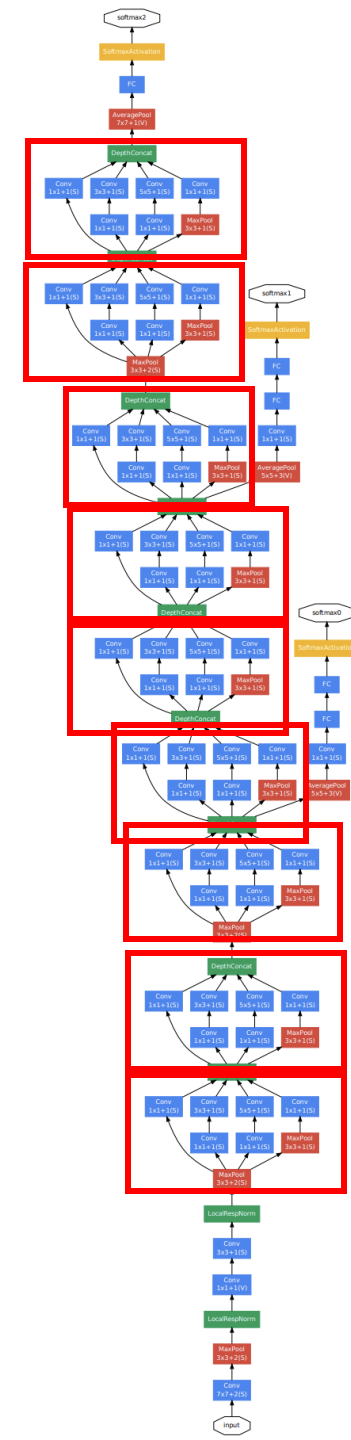
Filter (tensor)  
3×3×3



Input tensor  
 $d_1 \times d_2 \times 3$

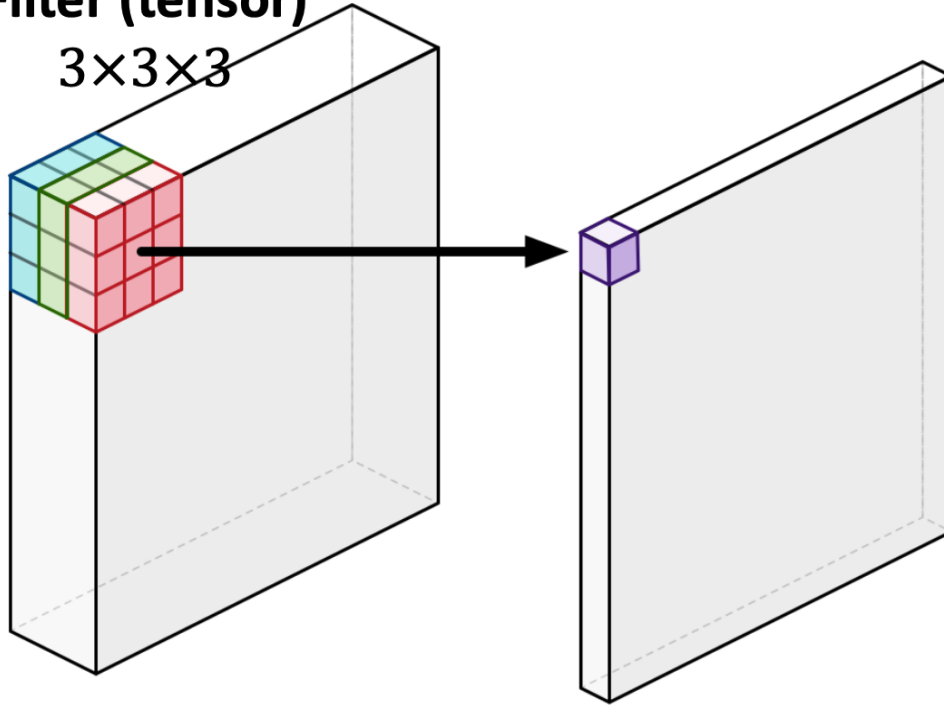
Output matrix  
 $(d_1 - 2) \times (d_2 - 2)$

a tensor  $\rightarrow$  a matrix (channel)  
 $\uparrow$   
 a filter/kernel



# Inception (GoogLeNet)

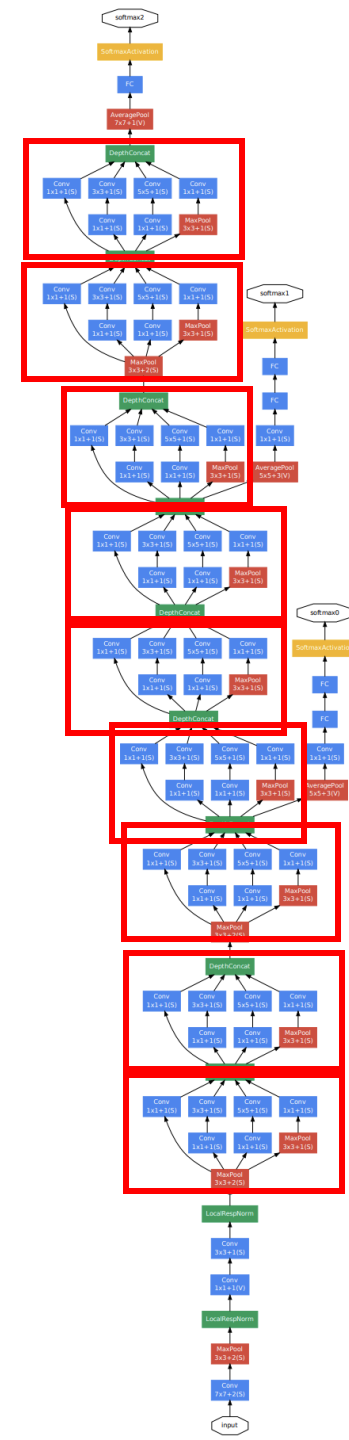
Filter (tensor)  
3×3×3



Input tensor  
 $d_1 \times d_2 \times 3$

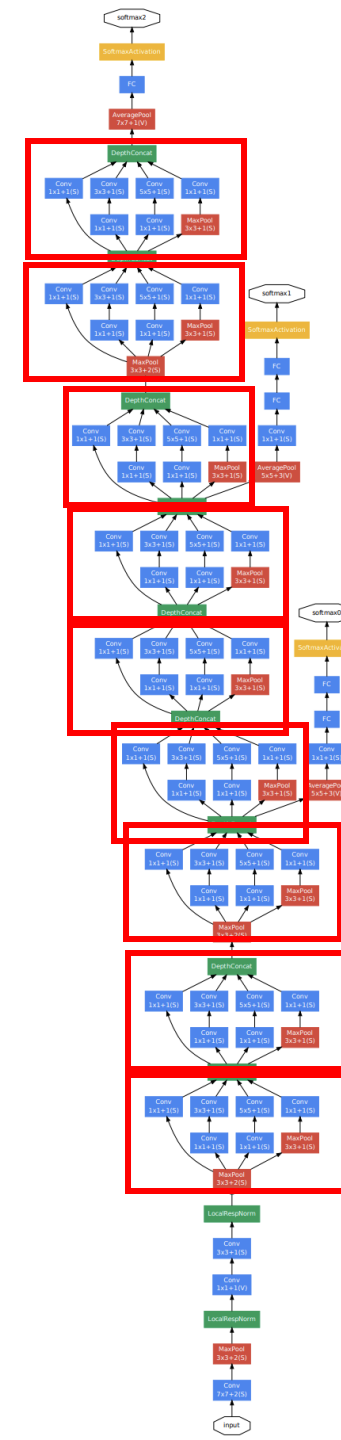
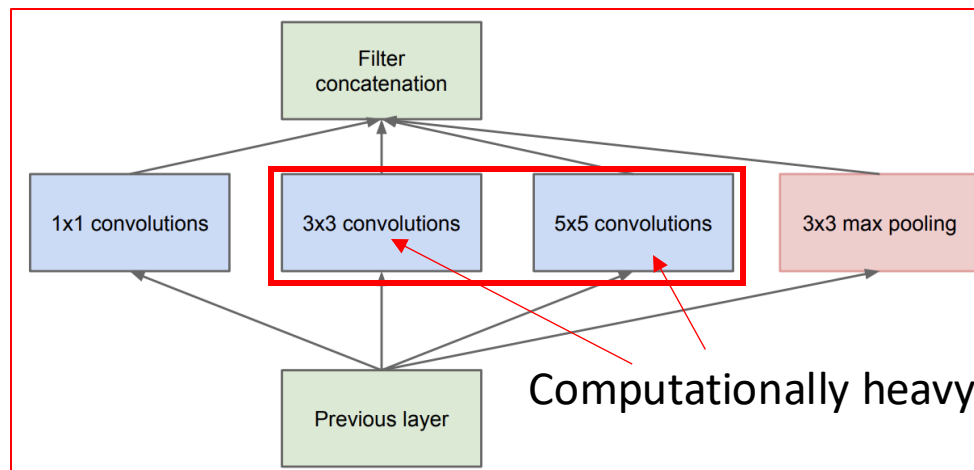
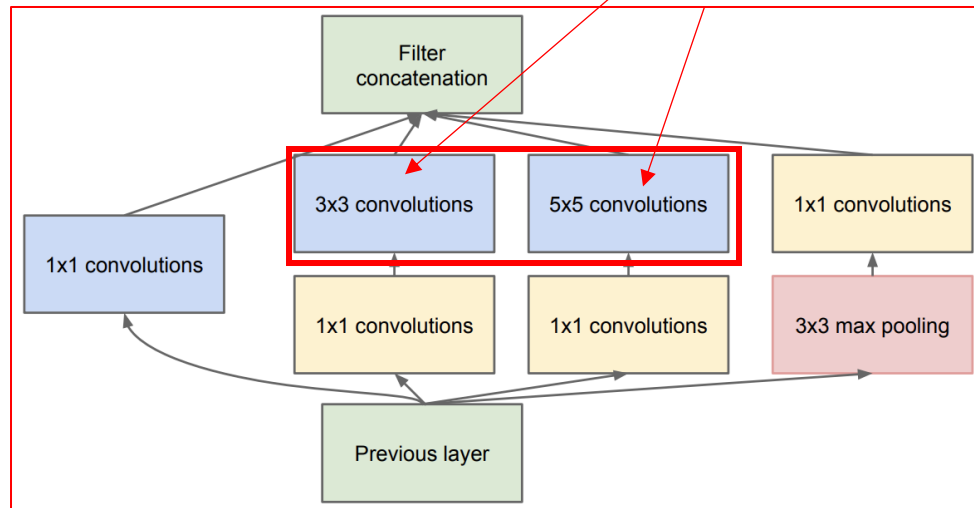
Output matrix  
 $(d_1 - 2) \times (d_2 - 2)$

a tensor  $\rightarrow$   $m$  matrices (channels)  
 $\uparrow$   
 $m$  filter/kernel



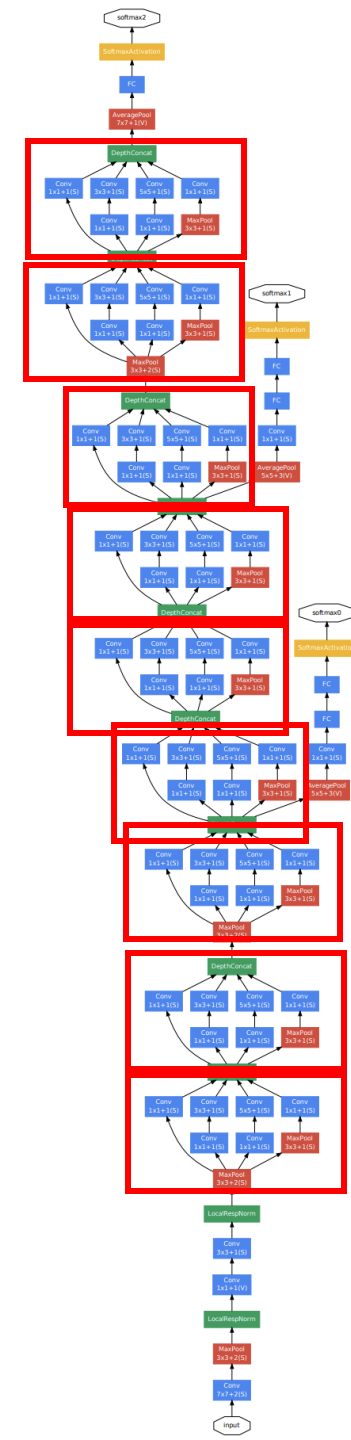
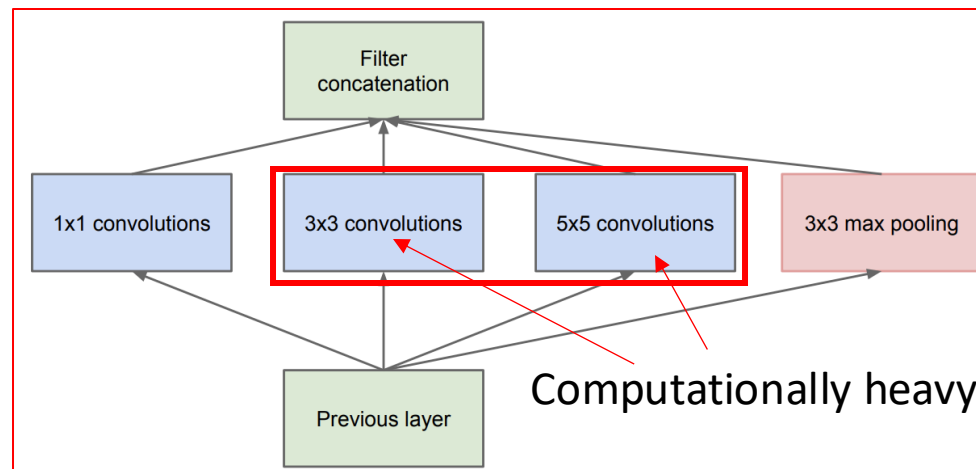
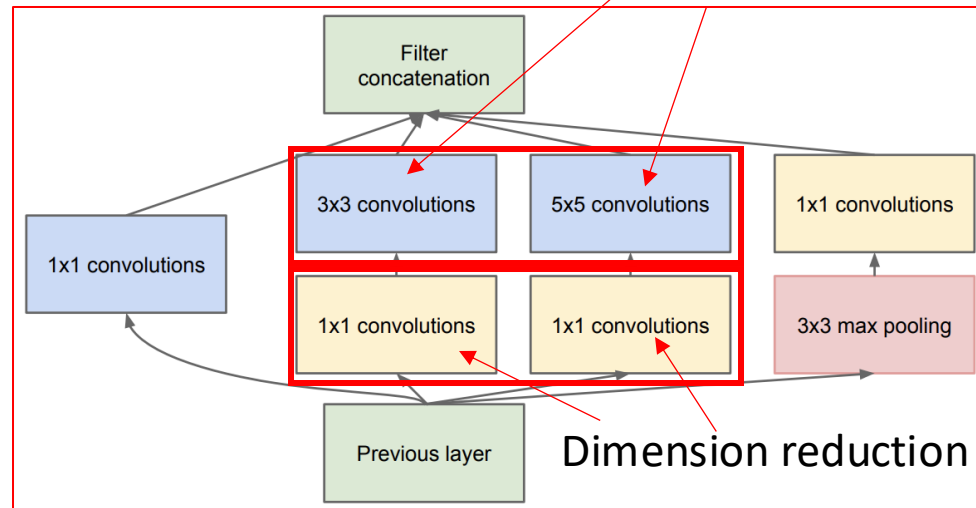
# Inception (GoogLeNet)

Computationally heavy



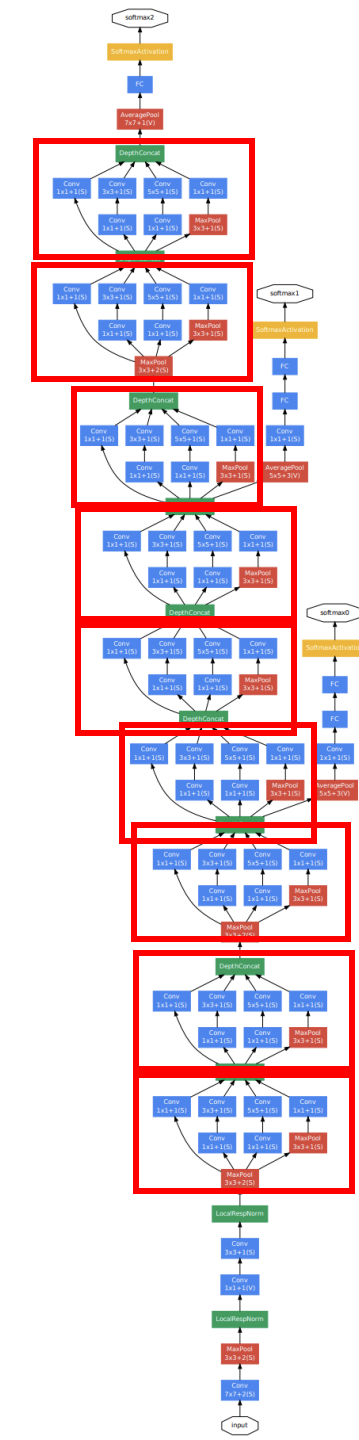
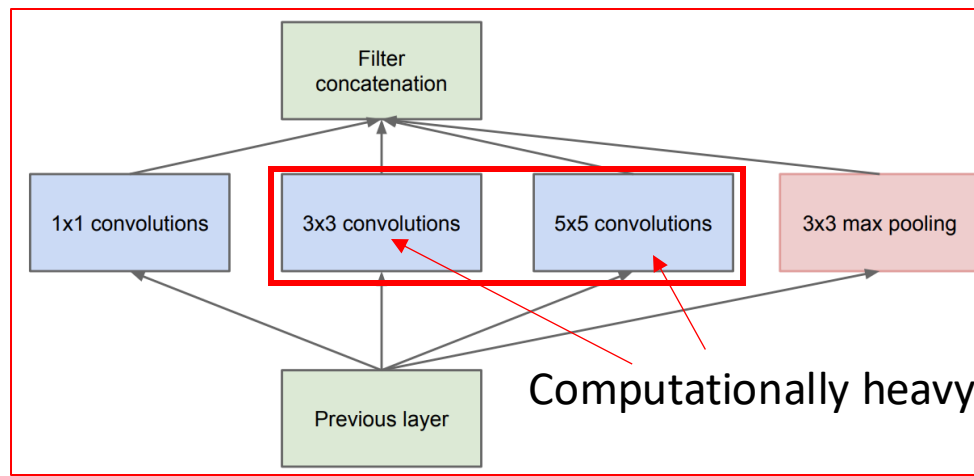
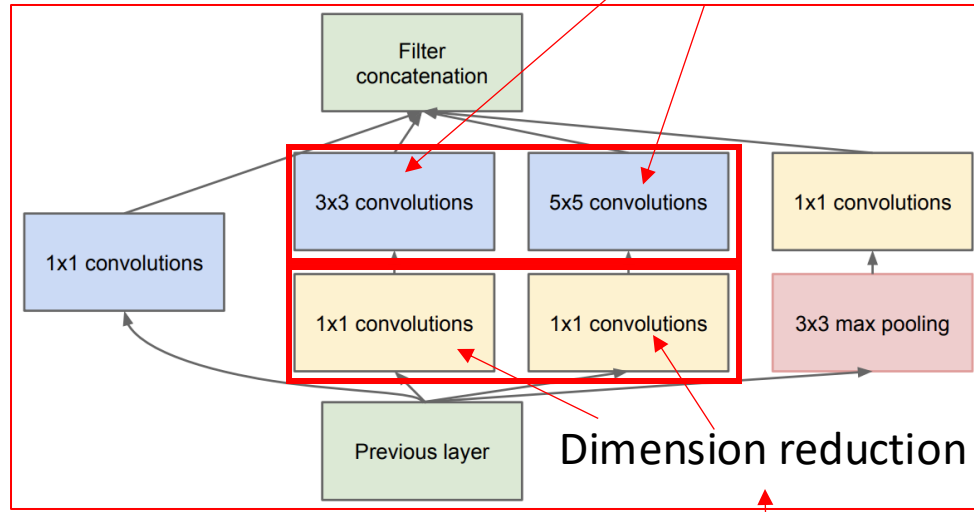
# Inception (GoogLeNet)

Computationally heavy

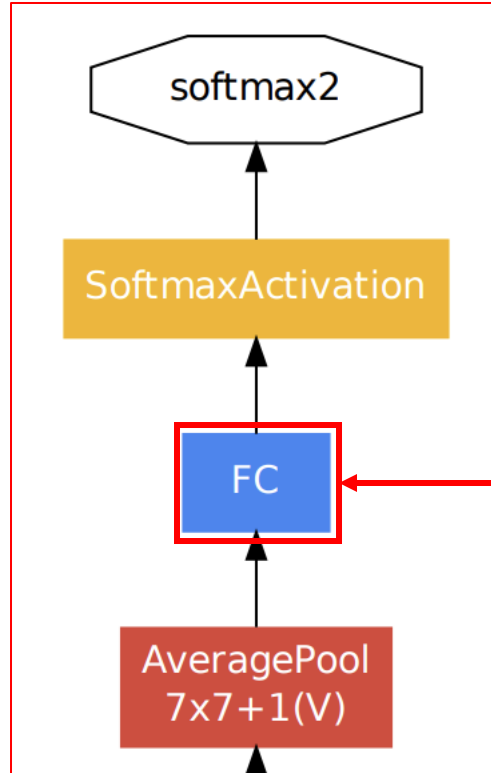


# Inception (GoogLeNet)

Computationally heavy



# Inception (GoogLeNet)

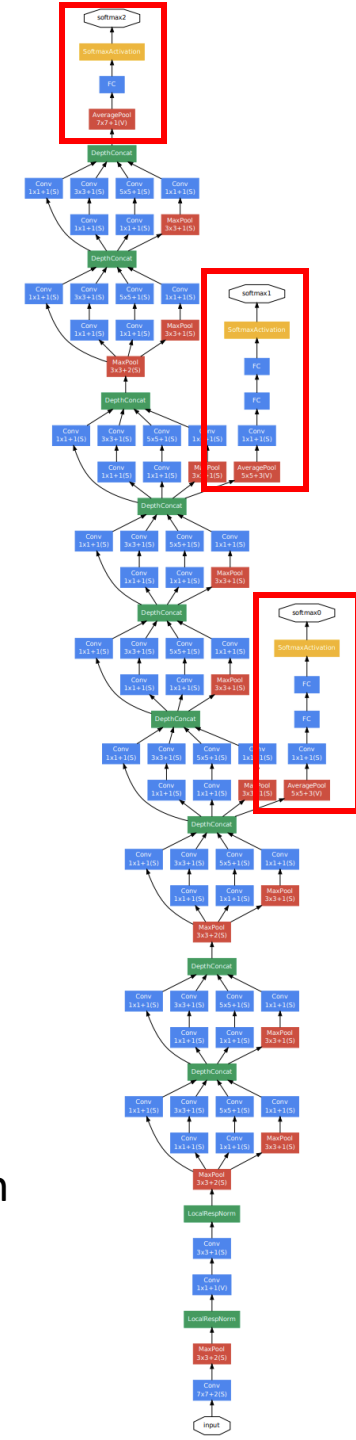
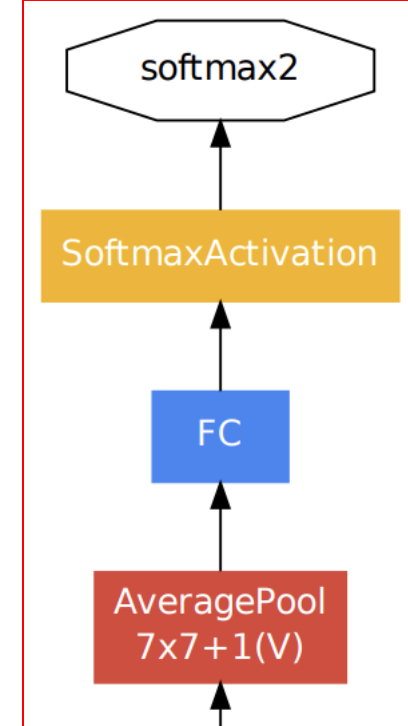
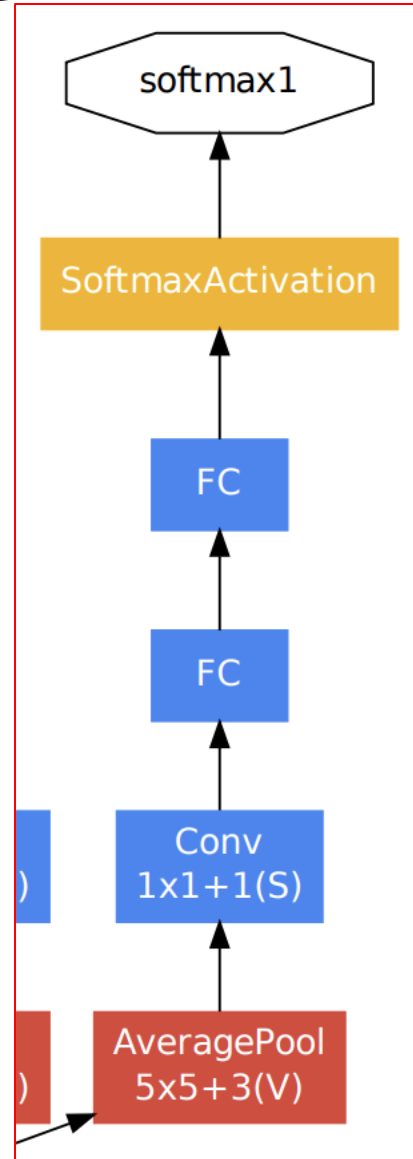
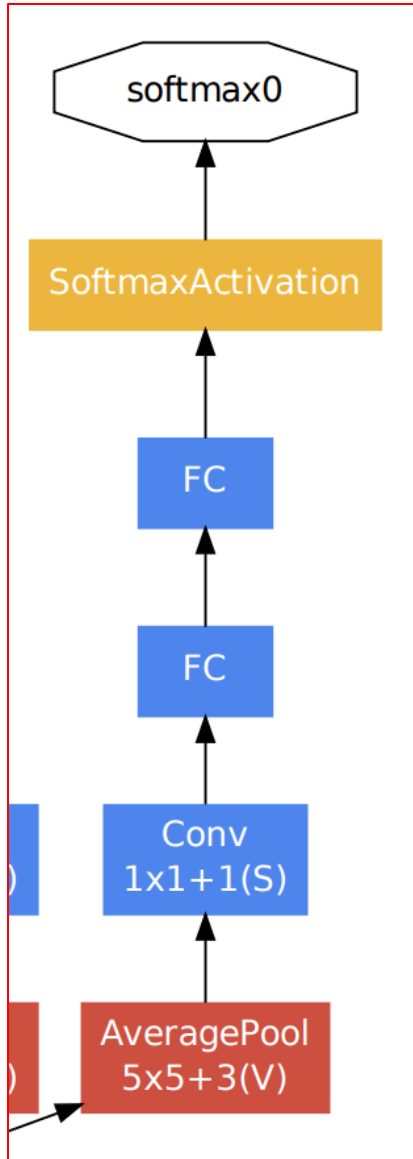






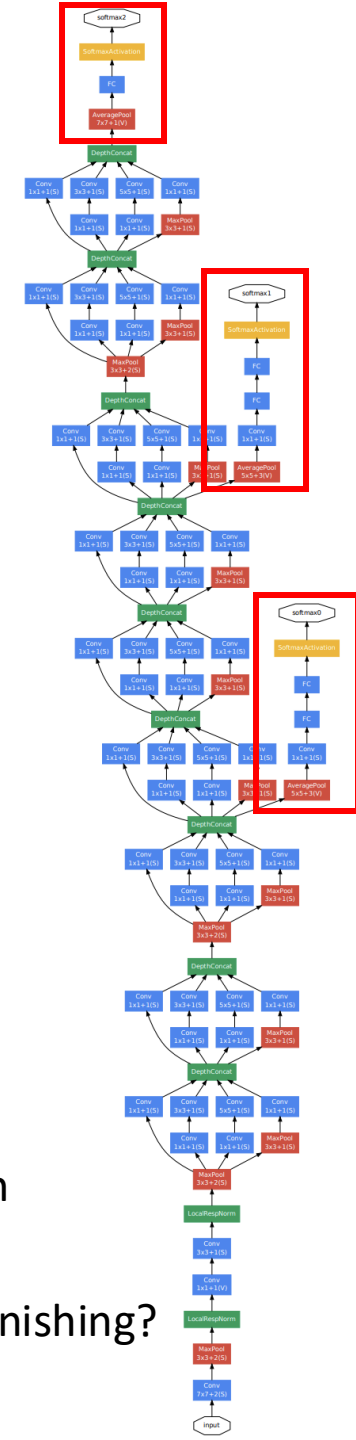
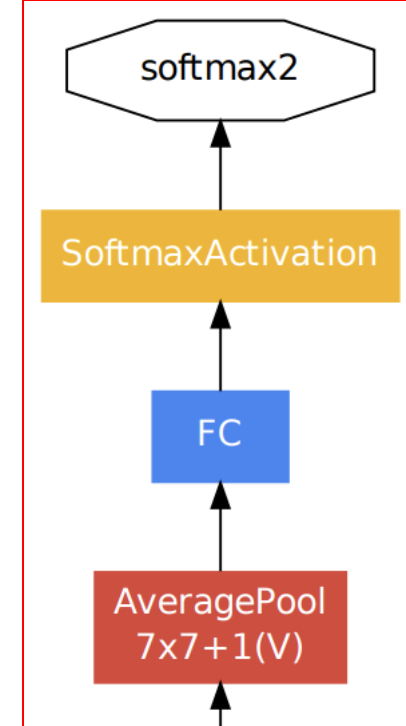
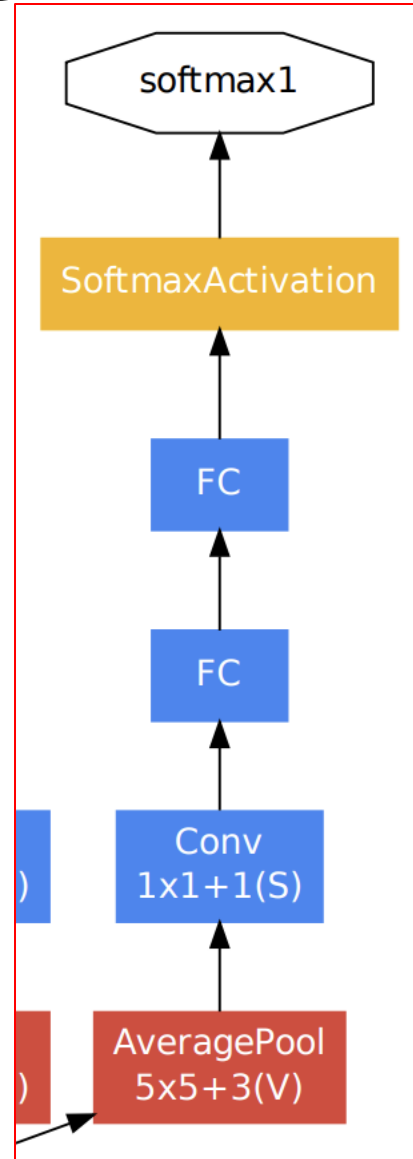
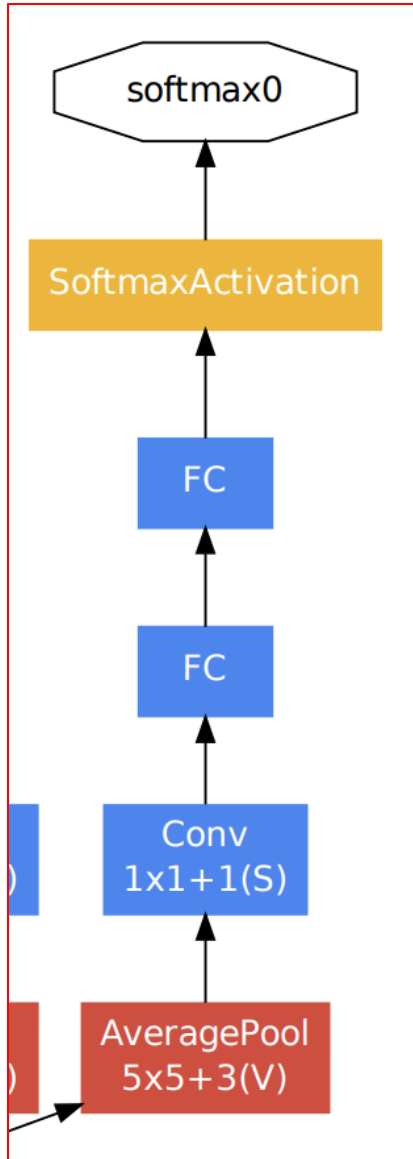


# Inception (GoogLeNet)



Q: why output prediction from lower layers?

# Inception (GoogLeNet)



**Q:** why output prediction from lower layers?

**Hint:** remember gradient vanishing?


# # of layers

- LeNet-5: 3 conv + 2 fc
- AlexNet: 5 conv + 2 fc
- VGG-16: 13 conv + 2 fc




More conv layers

# # of layers

- LeNet-5: 3 conv + 2 fc
  - AlexNet: 5 conv + 2 fc
  - VGG-16: 13 conv + 2 fc
- 
- More conv layers

Q: why they did not develop some architecture with more layers?

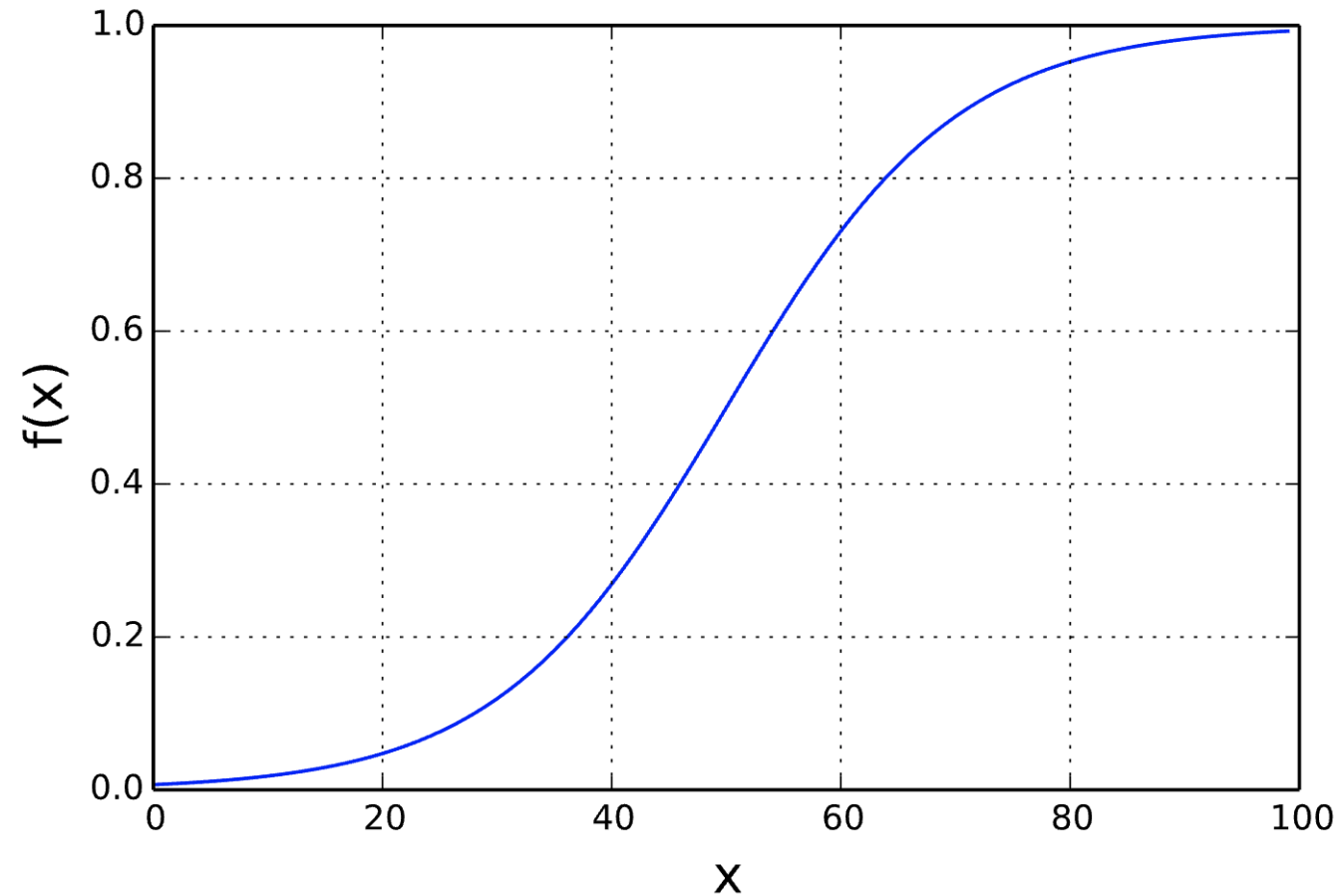
# # of layers

- LeNet-5: 3 conv + 2 fc
  - AlexNet: 5 conv + 2 fc
  - VGG-16: 13 conv + 2 fc
- 
- More conv layers

**Q:** why they did not develop some architecture with more layers?

**Hint (again):** remember gradient vanishing?

# Gradient vanish



Sigmoid function



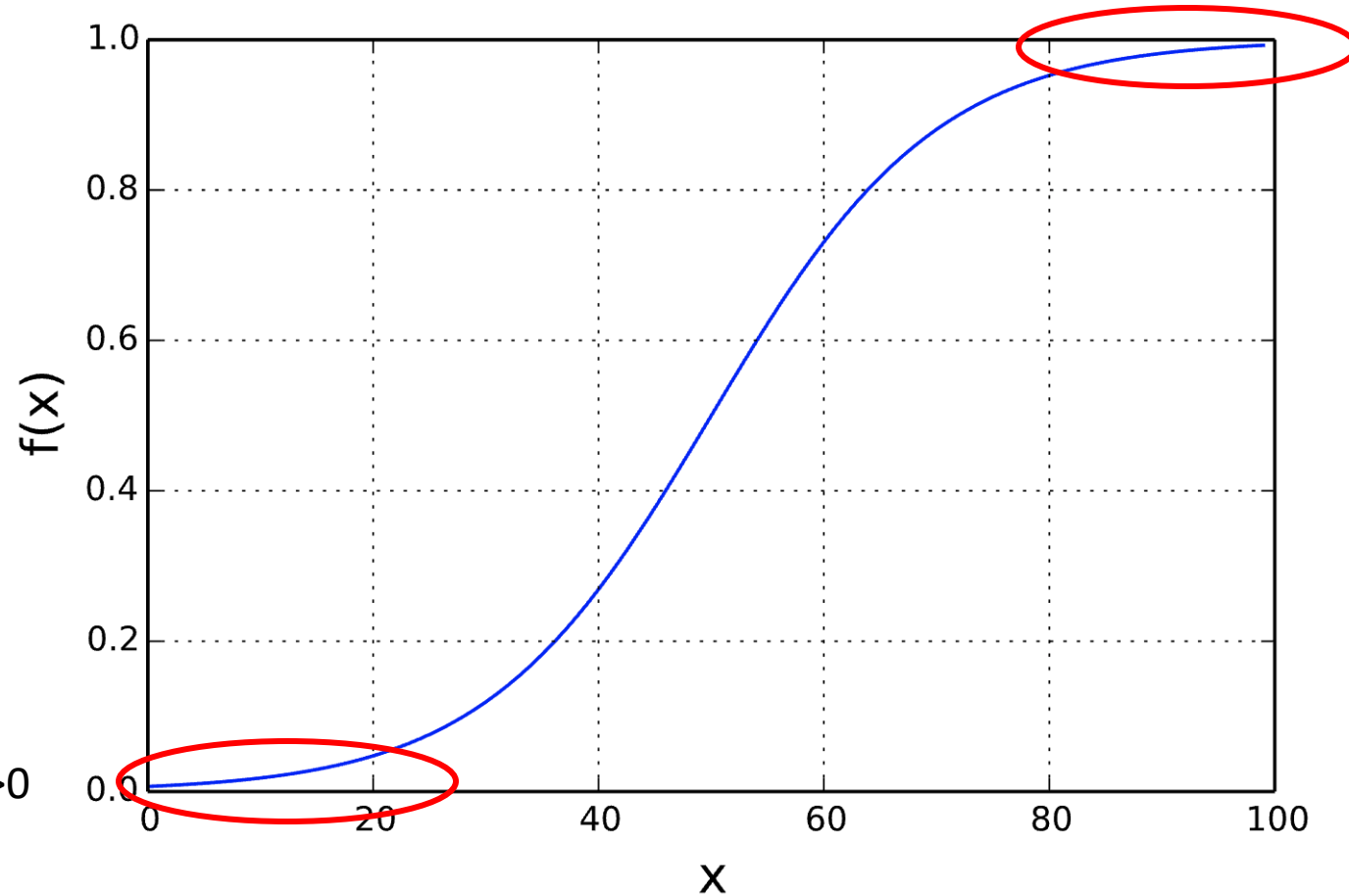
# Gradient vanish

$$f_n \left( \dots \left( f_2 \left( f_1(x) \right) \right) \right) \rightarrow ?$$

$$f_i \rightarrow x_i$$

$$\frac{dx_n}{dx_1} = \frac{dx_n}{dx_{n-1}} \cdot \dots \cdot \frac{dx_2}{dx_1} \cdot \frac{dx_1}{dx}$$

gradients->0



gradients->0

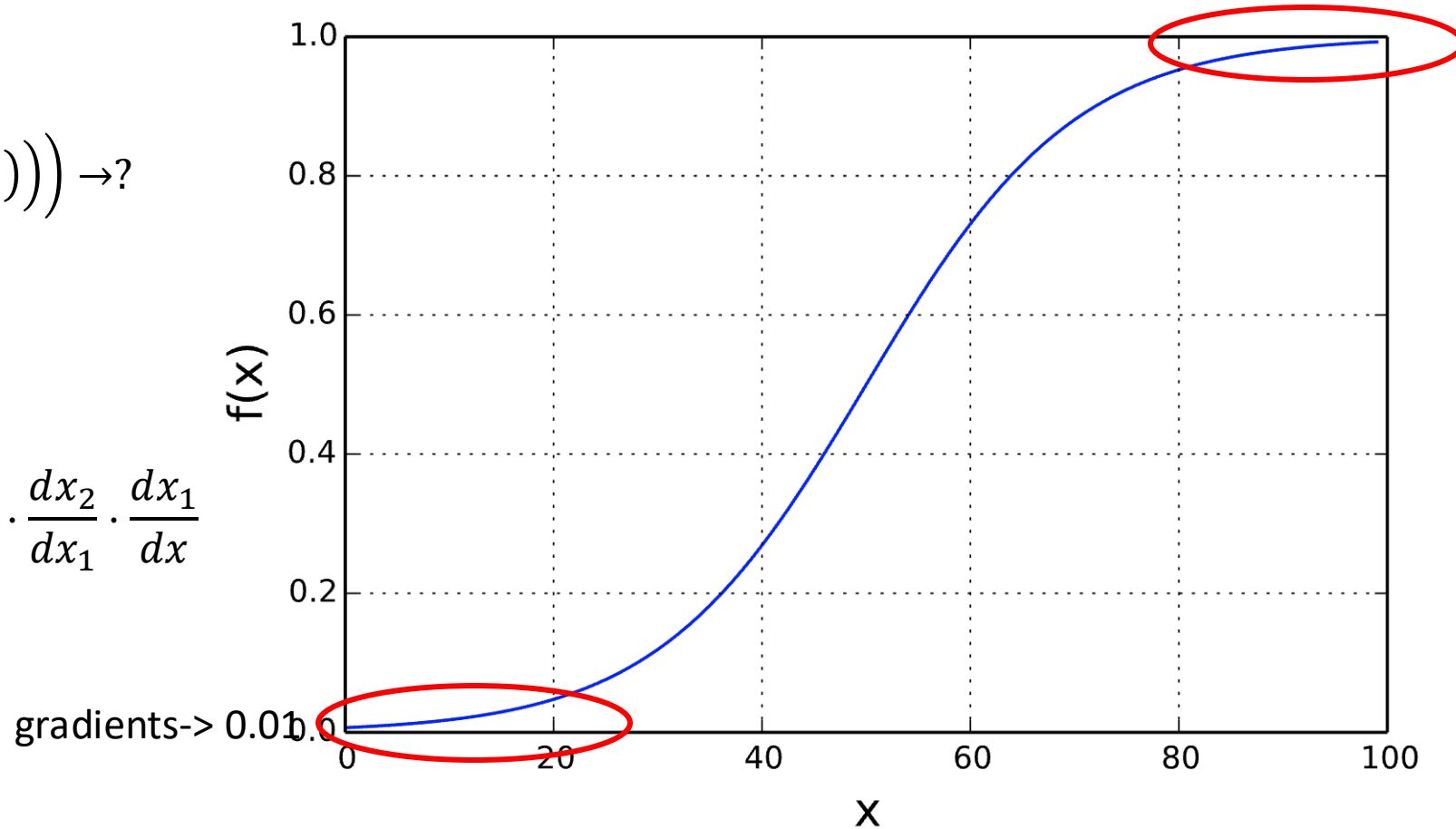
Sigmoid function

# Gradient vanish

$$f_n \left( \dots \left( f_2 \left( f_1(x) \right) \right) \right) \rightarrow ?$$

$$f_i \rightarrow x_i$$

$$\frac{dx_n}{dx_1} = \frac{dx_n}{dx_{n-1}} \cdot \dots \cdot \frac{dx_2}{dx_1} \cdot \frac{dx_1}{dx}$$



gradients- $\rightarrow$  0.01

gradients- $\rightarrow$  0.01

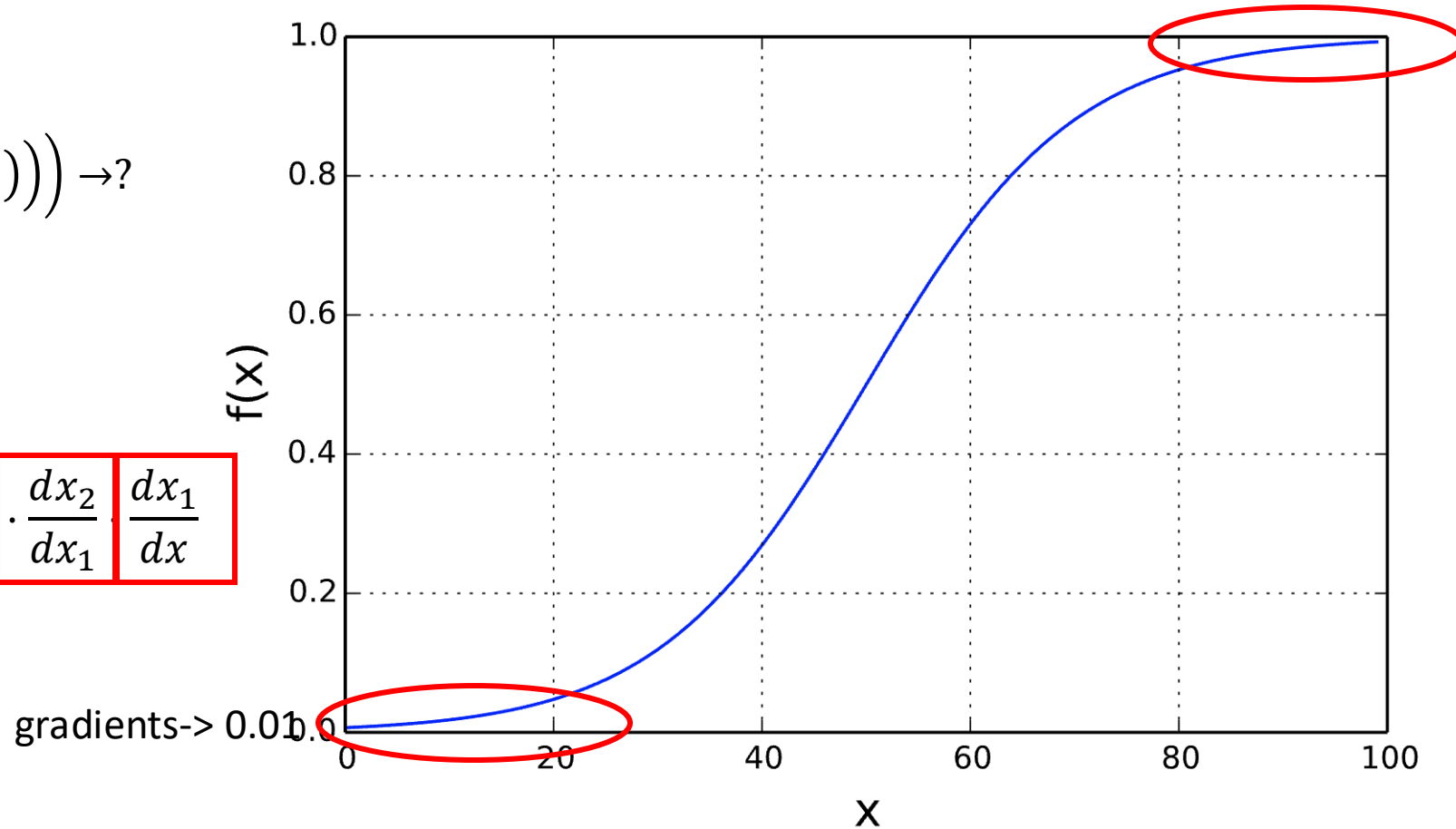
Sigmoid function

# Gradient vanish

$$f_n \left( \dots \left( f_2 \left( f_1(x) \right) \right) \right) \rightarrow ?$$

$$f_i \rightarrow x_i$$

$$\frac{dx_n}{dx_1} = \boxed{\frac{dx_n}{dx_{n-1}}} \cdots \boxed{\frac{dx_2}{dx_1} \frac{dx_1}{dx}}$$



Sigmoid function

# Gradient vanish

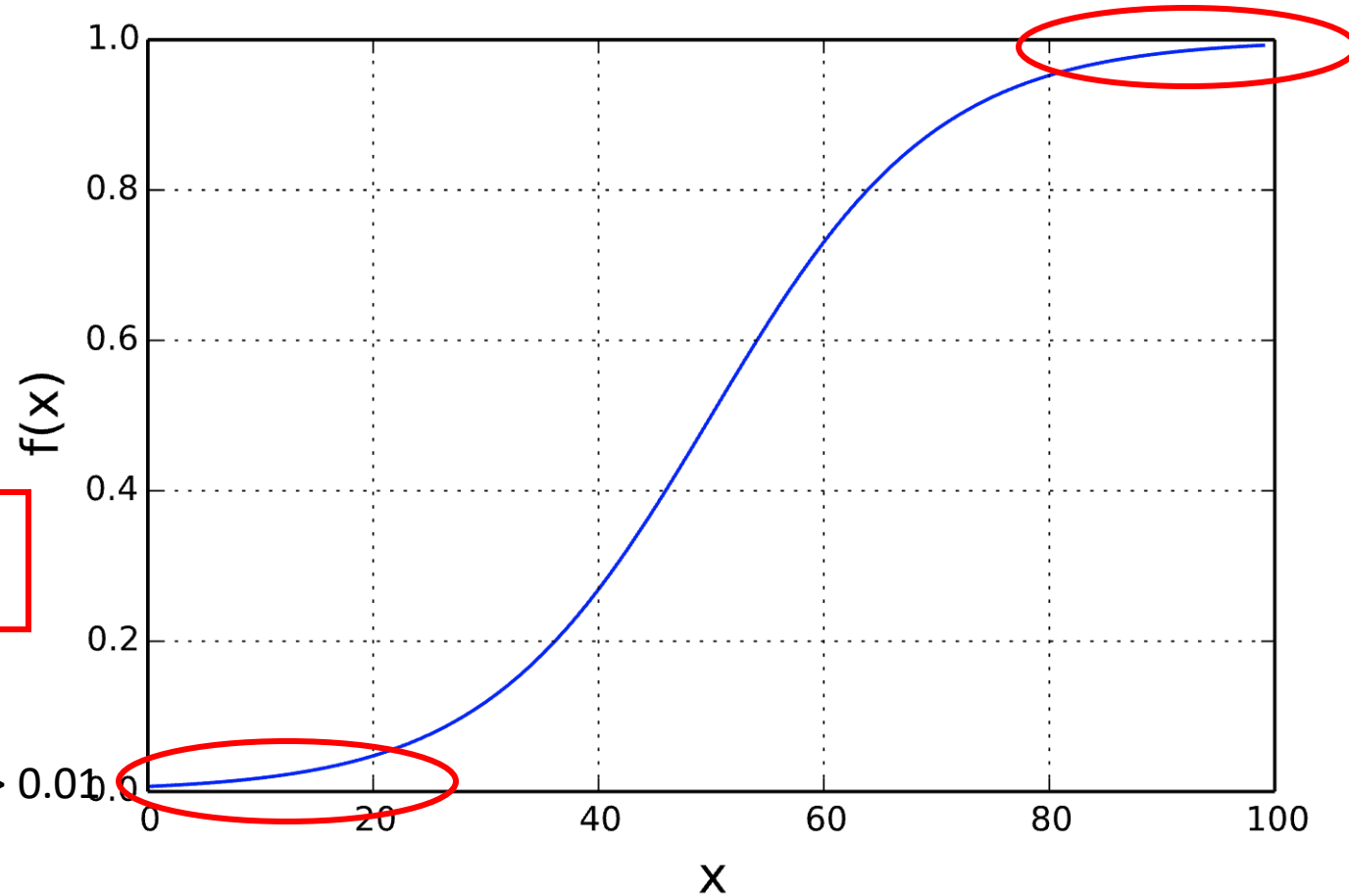
$$f_n \left( \dots \left( f_2 \left( f_1(x) \right) \right) \right) \rightarrow ?$$

$$f_i \rightarrow x_i$$

$$\frac{dx_n}{dx_1} = \boxed{\frac{dx_n}{dx_{n-1}}} \cdots \boxed{\frac{dx_2}{dx_1} \frac{dx_1}{dx}}$$

$\rightarrow 0.01^n$

gradients  $\rightarrow 0.01^n$



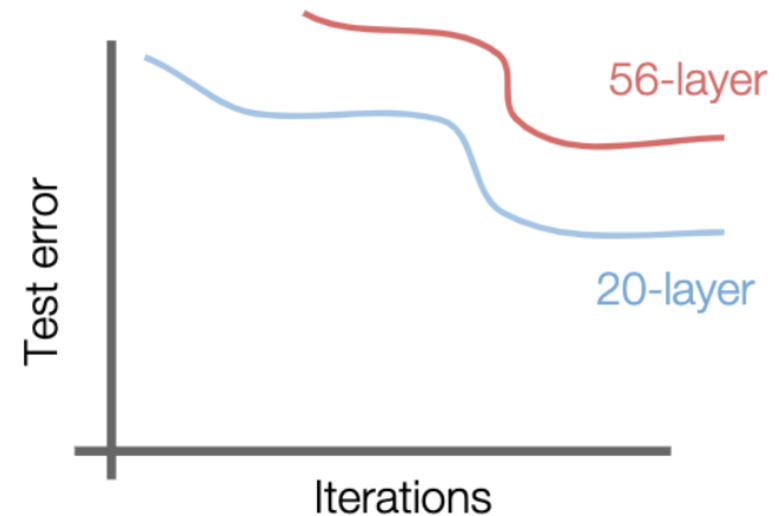
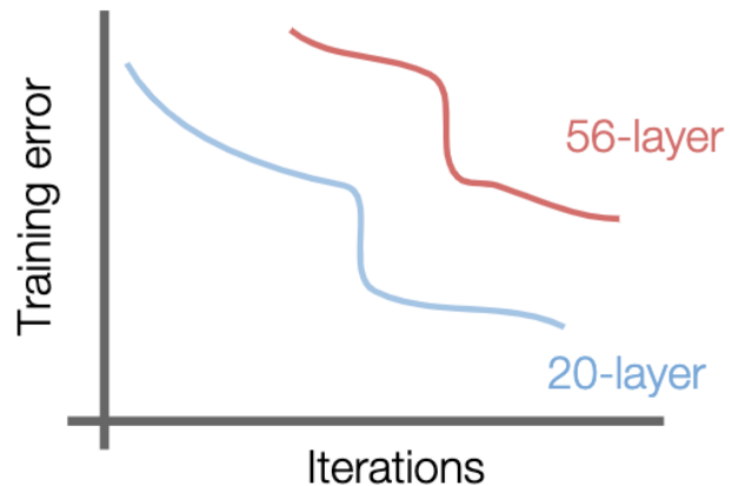
gradients  $\rightarrow 0.01$

Sigmoid function

# # of layers

- LeNet-5: 3 conv + 2 fc
- AlexNet: 5 conv + 2 fc
- VGG-16: 13 conv + 2 fc

Q: why they did not develop some architecture with more layers?  
Optimization may be difficult.

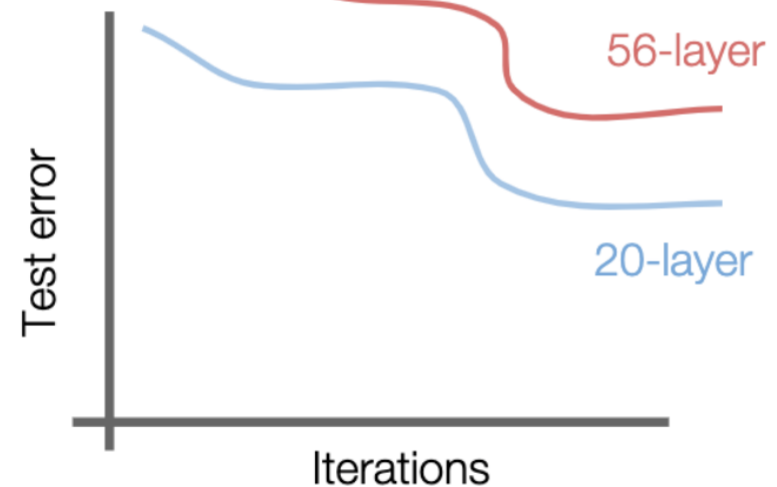
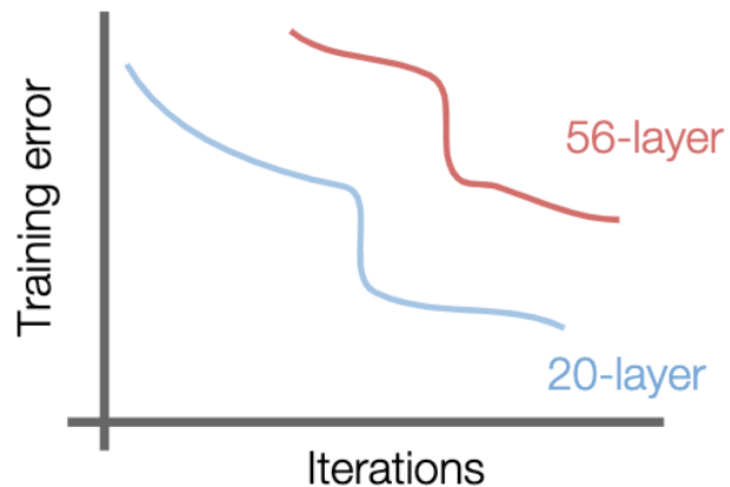


# # of layers

- LeNet-5: 3 conv + 2 fc
- AlexNet: 5 conv + 2 fc
- VGG-16: 13 conv + 2 fc

Q: why they did not develop some architecture with more layers?

Optimization may be difficult. We do not have a good solution as our model.

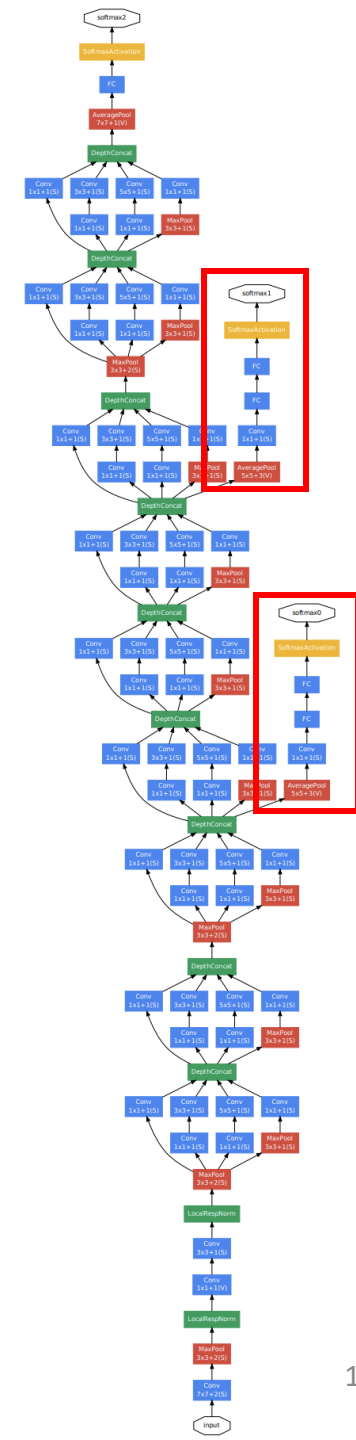


# Inception

$$f_n \left( \dots \left( f_2 \left( f_1(x) \right) \right) \right) \rightarrow ?$$

$$\frac{dx_n}{dx_1} = \frac{dx_n}{dx_{n-1}} \cdot \dots \cdot \frac{dx_2}{dx_1} \cdot \frac{dx_1}{dx}$$

→ 0.01<sup>n</sup>

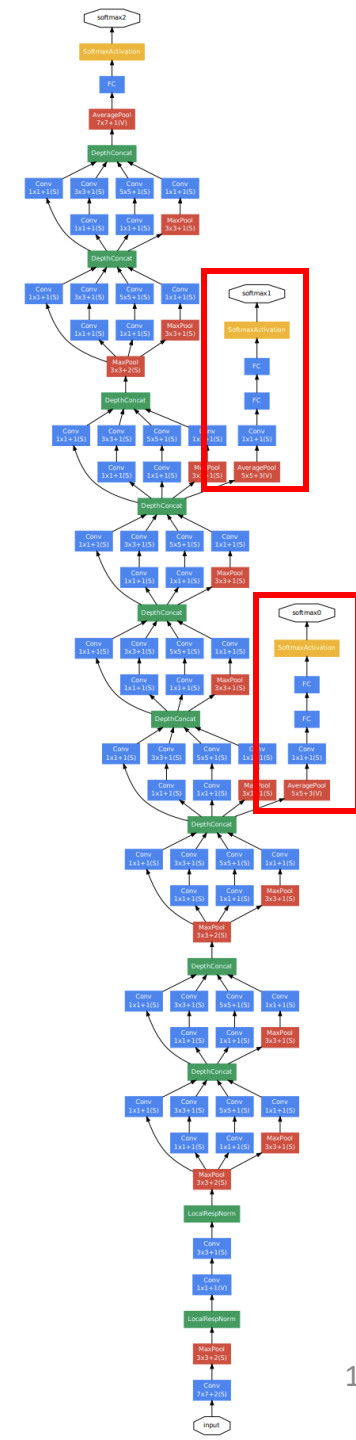


# Inception

$$f_n \left( \dots \left( f_2 \left( f_1(x) \right) \right) \right) + f_m \left( \dots \left( f_{n+2} \left( f_{n+1}(x) \right) \right) \right)$$

$$\frac{dx_n}{dx_1} = \frac{dx_n}{dx_{n-1}} \cdot \dots \cdot \frac{dx_2}{dx_1} \cdot \frac{dx_1}{dx}$$

→ 0.01<sup>n</sup>





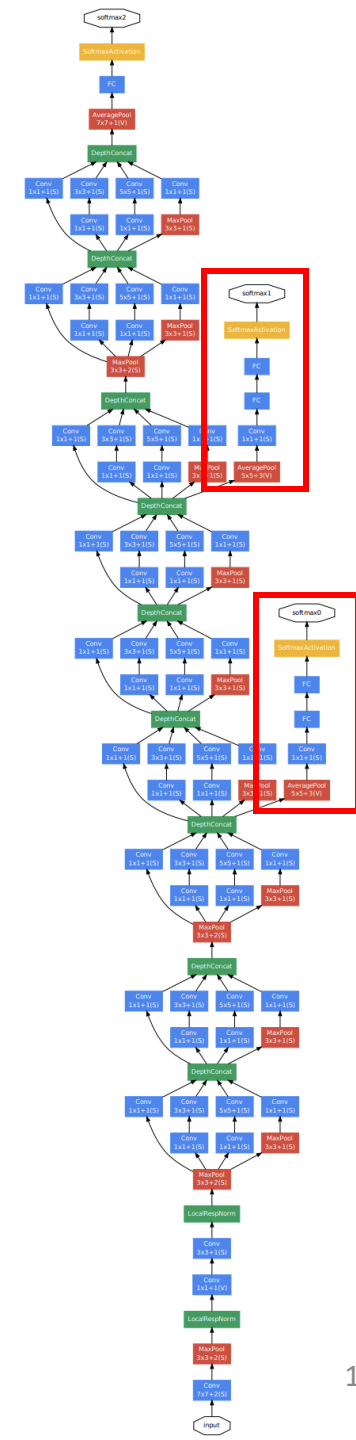
# Inception

$$f_n \left( \dots \left( f_2 \left( f_1(x) \right) \right) \right) + f_m \left( \dots \left( f_{n+2} \left( f_{n+1}(x) \right) \right) \right)$$

$$\frac{dx_n}{dx_1} = \frac{dx_n}{dx_{n-1}} \dots \frac{dx_2}{dx_1} \cdot \frac{dx_1}{dx} + \frac{dx_m}{dx_{m-1}} \dots \frac{dx_{n+2}}{dx_{n+1}} \cdot \frac{dx_{n+1}}{dx}$$

$\rightarrow 0.01^n$ 
 $\gg 0$

Will not be very small



# Residual neural networks (ResNet)

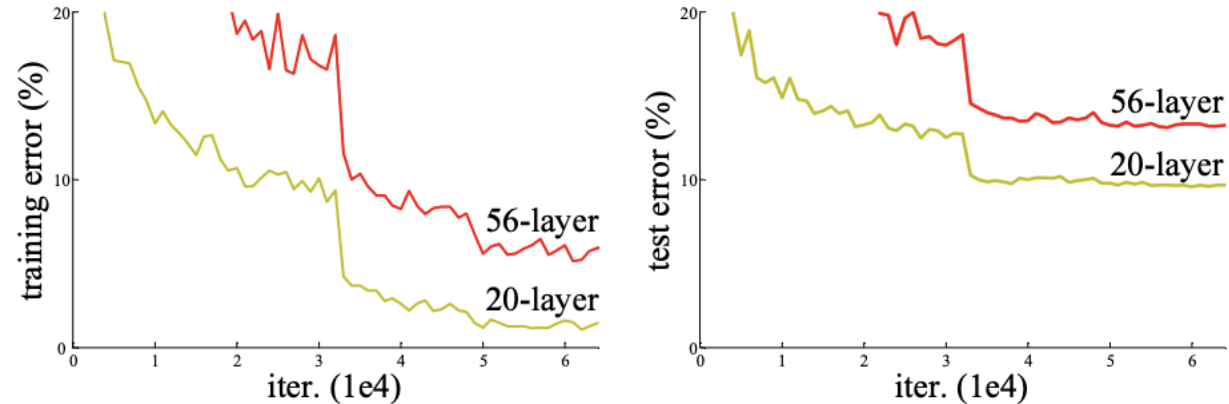


Figure 1. Training error (left) and test error (right) on CIFAR-10 with 20-layer and 56-layer “plain” networks. The deeper network has higher training error, and thus test error. Similar phenomena on ImageNet is presented in Fig. 4.

# Residual neural networks (ResNet)

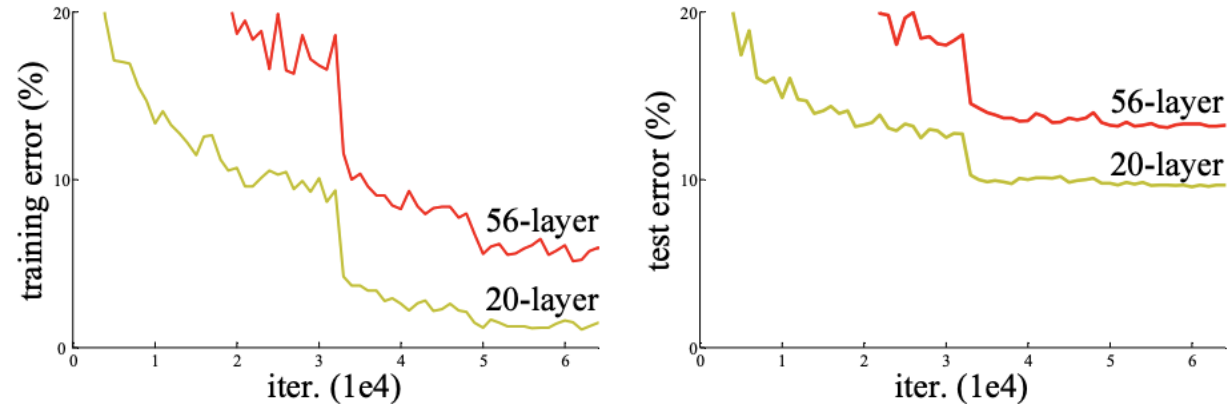


Figure 1. Training error (left) and test error (right) on CIFAR-10 with 20-layer and 56-layer “plain” networks. The deeper network has higher training error, and thus test error. Similar phenomena on ImageNet is presented in Fig. 4.

Without special structure other than conv/fc layers

# ResNet: shortcut connection

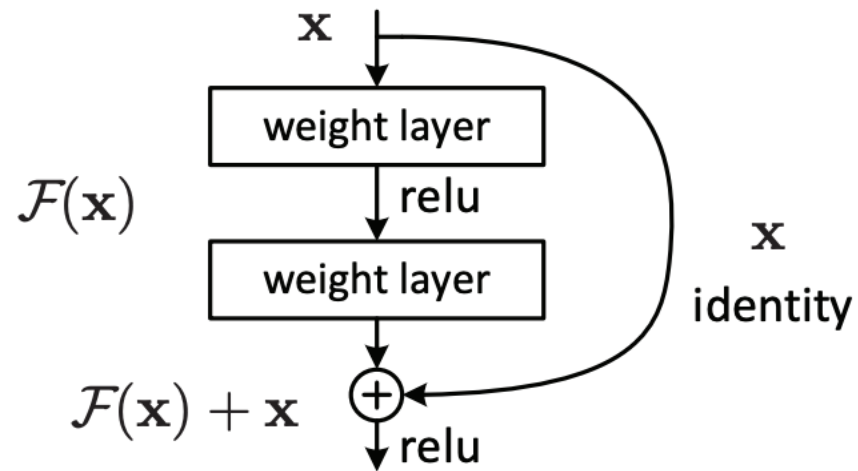


Figure 2. Residual learning: a building block.

# ResNet: shortcut connection

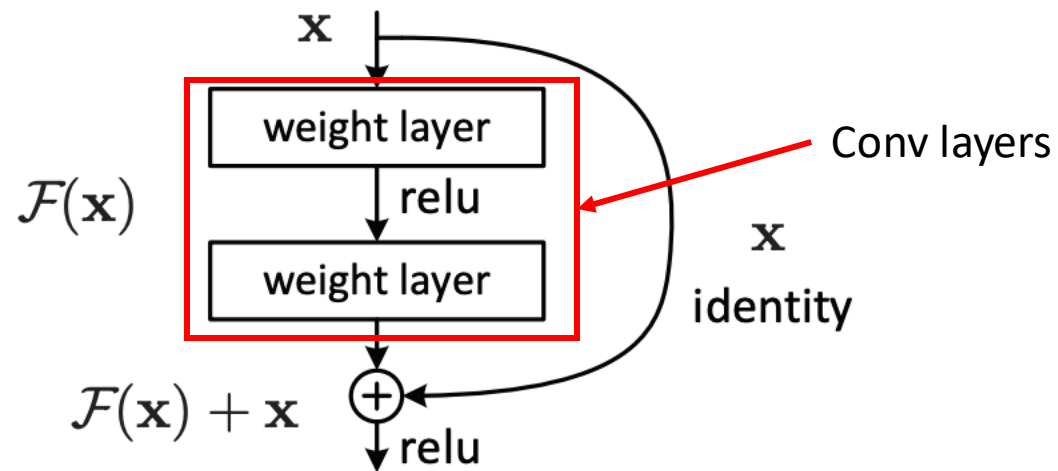


Figure 2. Residual learning: a building block.

# ResNet: shortcut connection

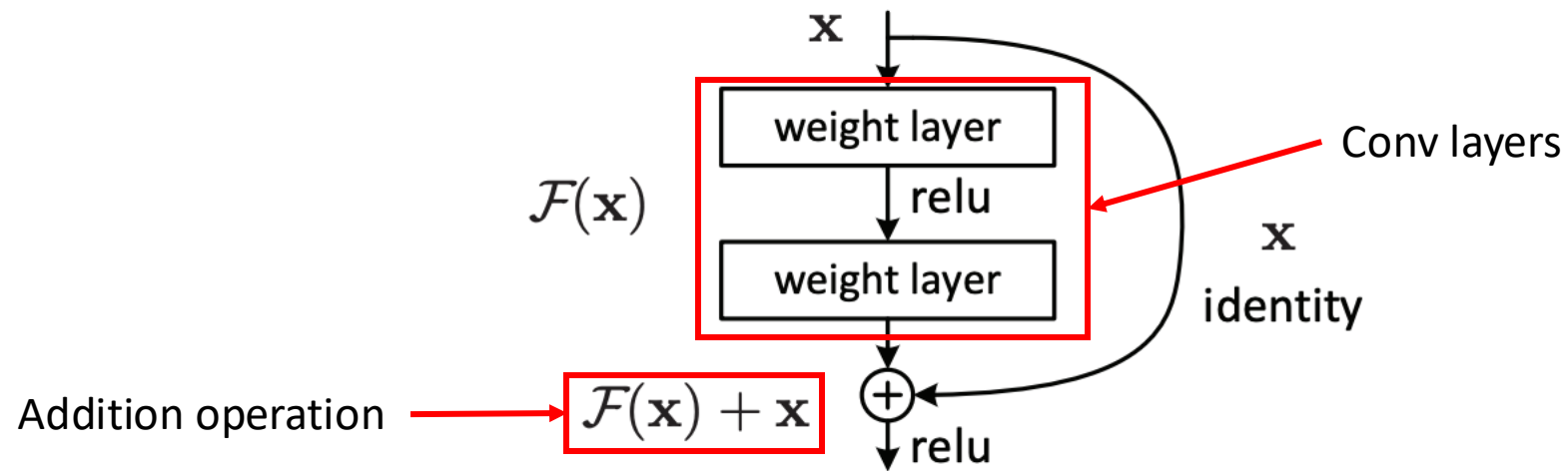


Figure 2. Residual learning: a building block.

# ResNet: shortcut connection

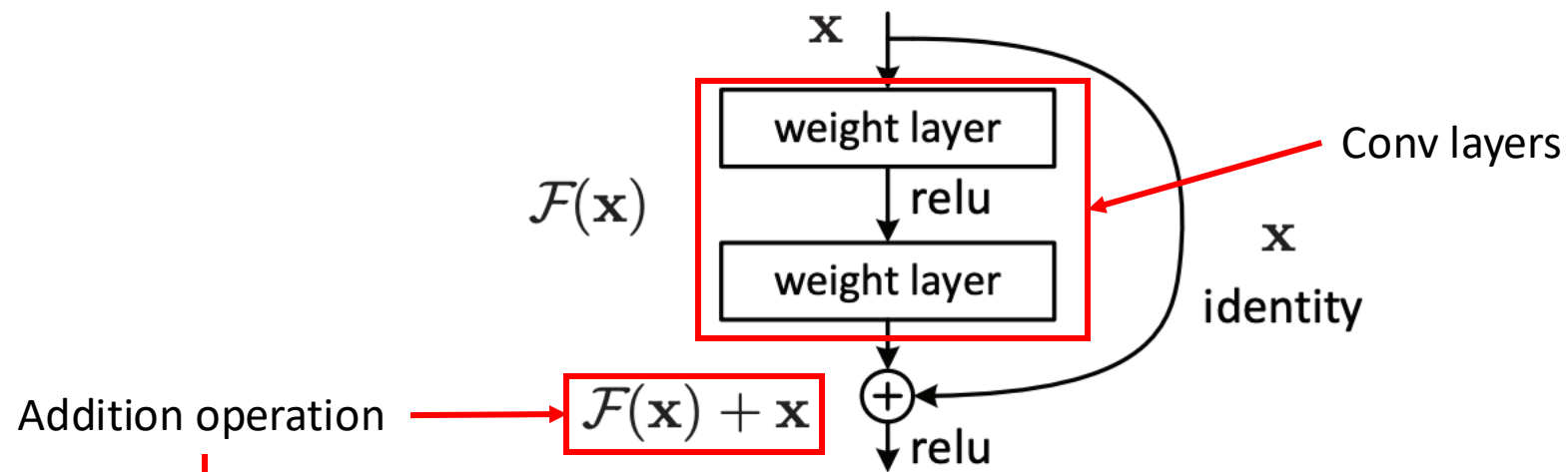


Figure 2. Residual learning: a building block.

implication: same dimension





# ResNet

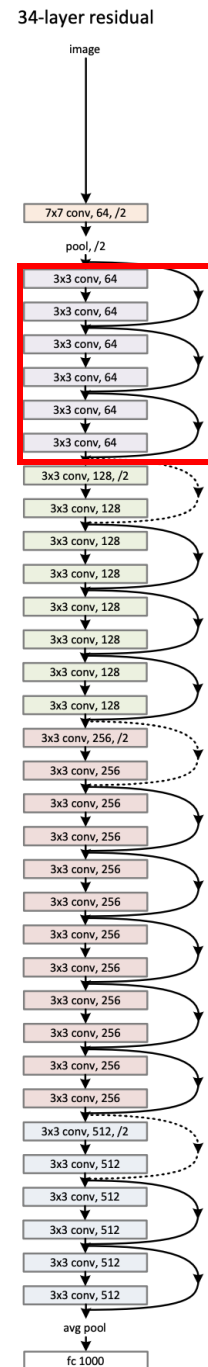
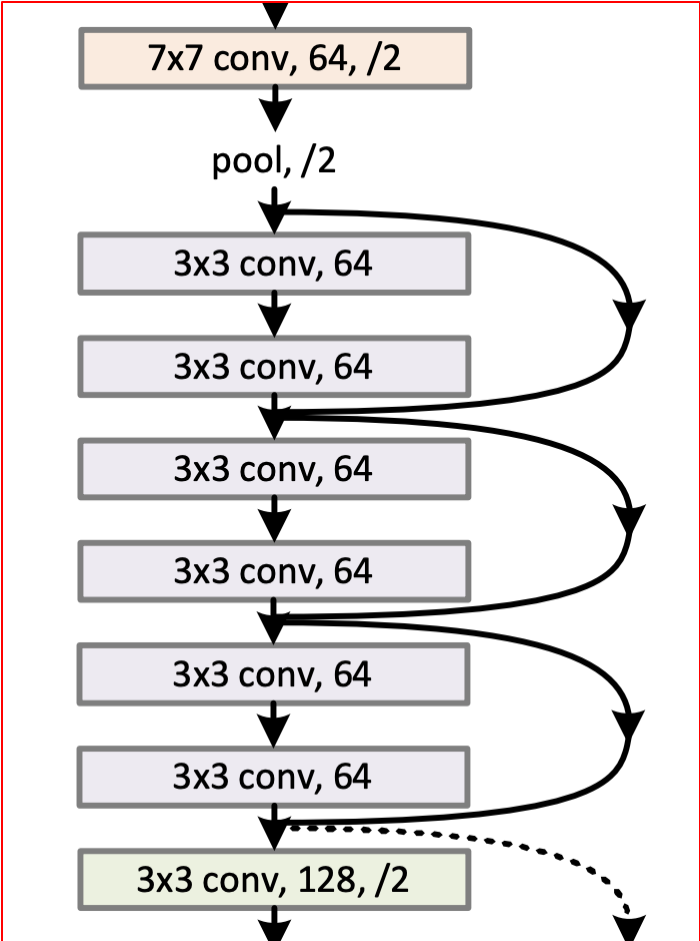


Figure 3. Example network architectures for ImageNet. **Left:** the VGG-19 model [41] (19.6 billion FLOPs) as a reference. **Middle:** a plain network with 34 parameter layers (3.6 billion FLOPs). **Right:** a residual network with 34 parameter layers (3.6 billion FLOPs). The dotted shortcuts increase dimensions. **Table 1** shows more details and other variants.

# ResNet



34-layer residual

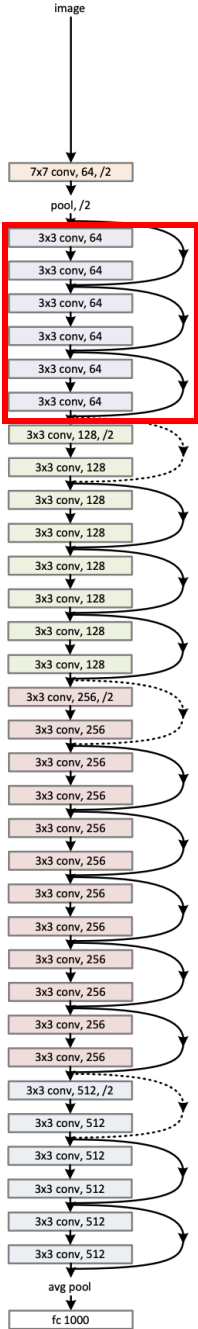
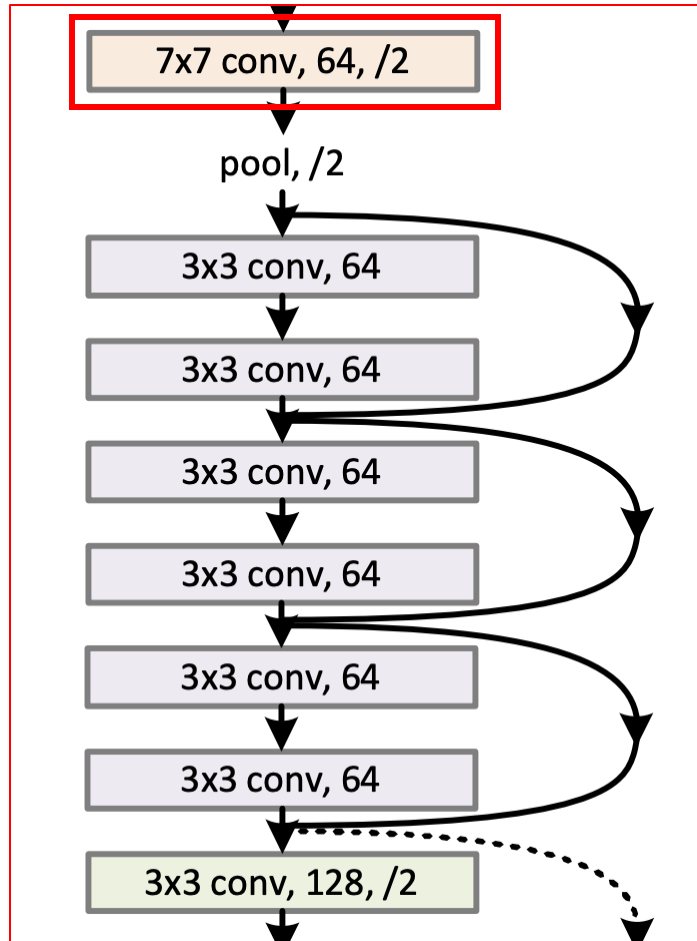


Figure 3. Example network architectures for ImageNet. **Left:** the VGG-19 model [41] (19.6 billion FLOPs) as a reference. **Middle:** a plain network with 34 parameter layers (3.6 billion FLOPs). **Right:** a residual network with 34 parameter layers (3.6 billion FLOPs). The dotted shortcuts increase dimensions. **Table 1** shows more details and other variants.

# ResNet



34-layer residual

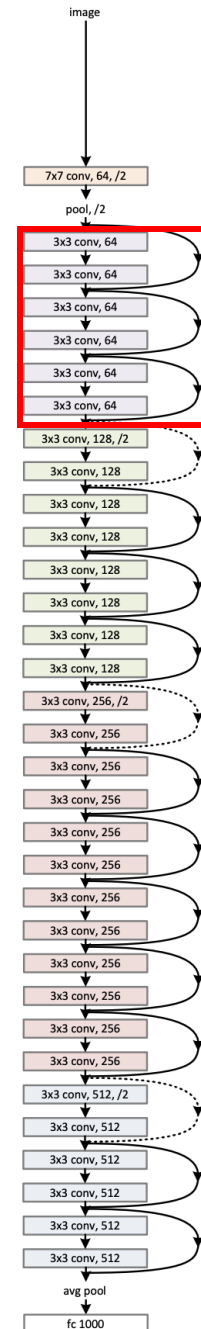
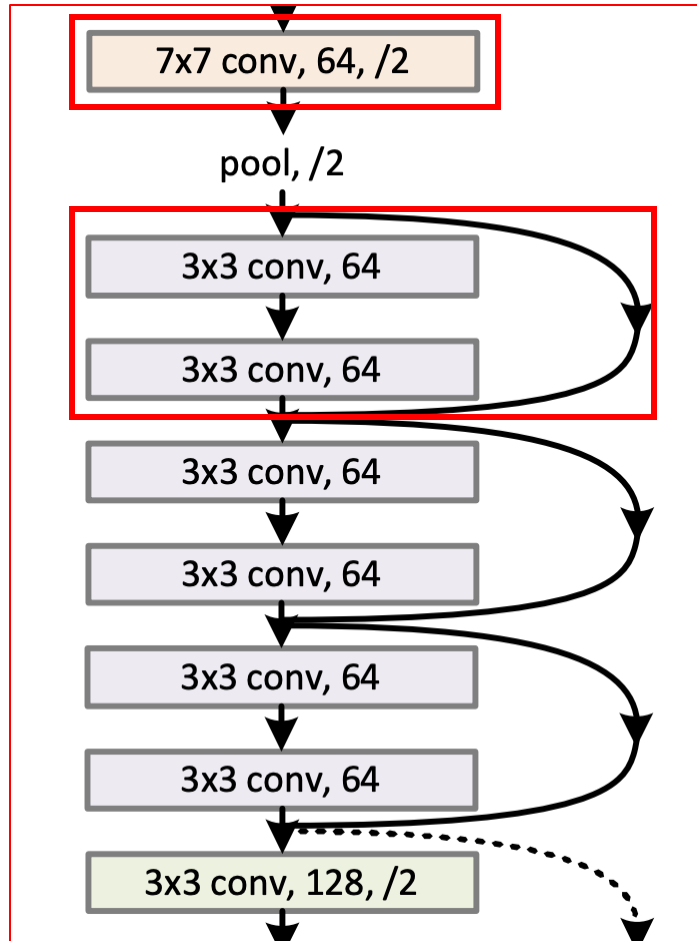


Figure 3. Example network architectures for ImageNet. **Left:** the VGG-19 model [41] (19.6 billion FLOPs) as a reference. **Middle:** a plain network with 34 parameter layers (3.6 billion FLOPs). **Right:** a residual network with 34 parameter layers (3.6 billion FLOPs). The dotted shortcuts increase dimensions. **Table 1** shows more details and other variants.

# ResNet



34-layer residual

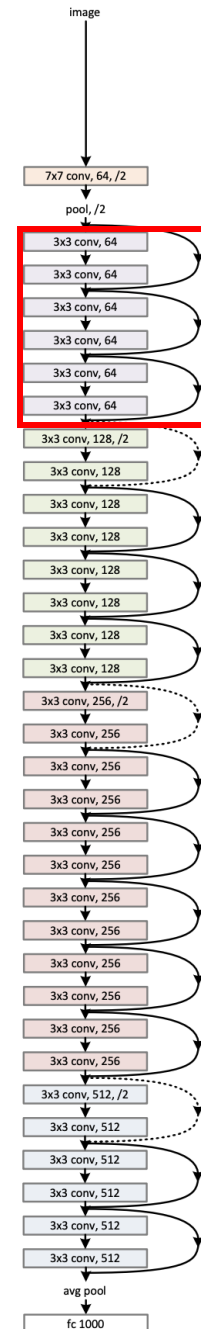
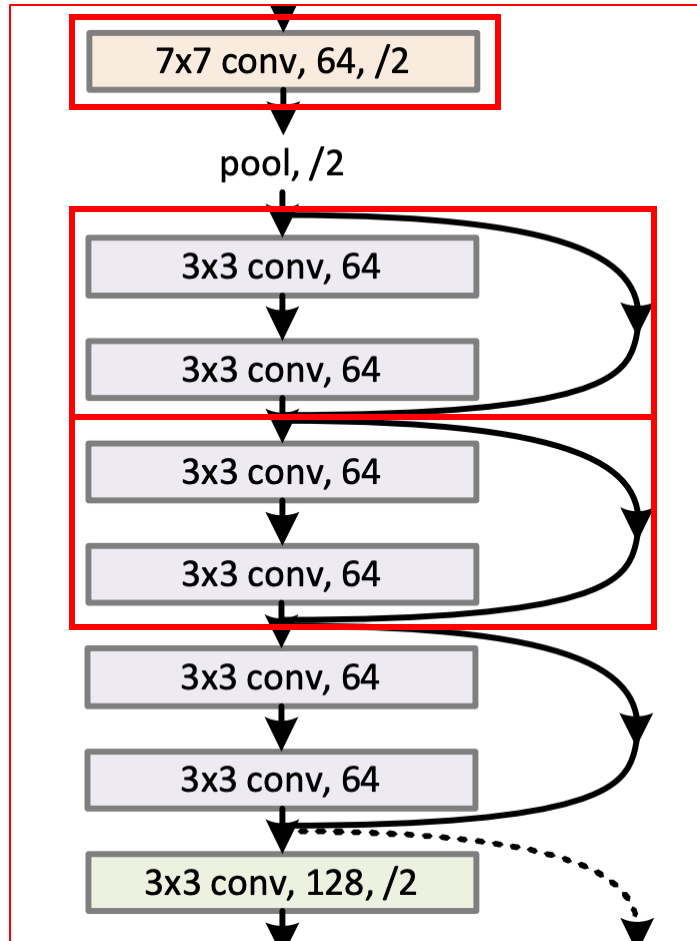


Figure 3. Example network architectures for ImageNet. **Left:** the VGG-19 model [41] (19.6 billion FLOPs) as a reference. **Middle:** a plain network with 34 parameter layers (3.6 billion FLOPs). **Right:** a residual network with 34 parameter layers (3.6 billion FLOPs). The dotted shortcuts increase dimensions. **Table 1** shows more details and other variants.

# ResNet



34-layer residual

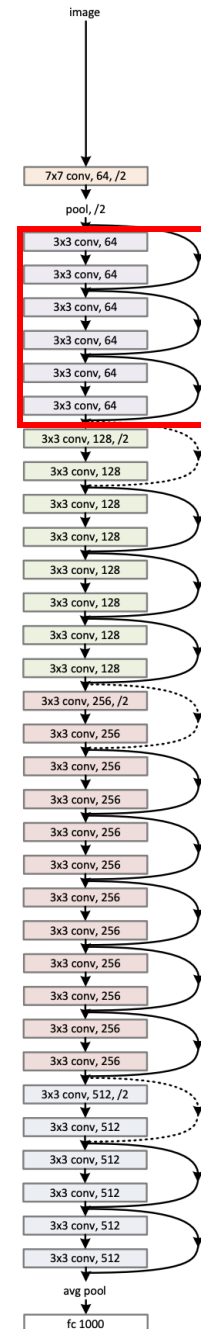
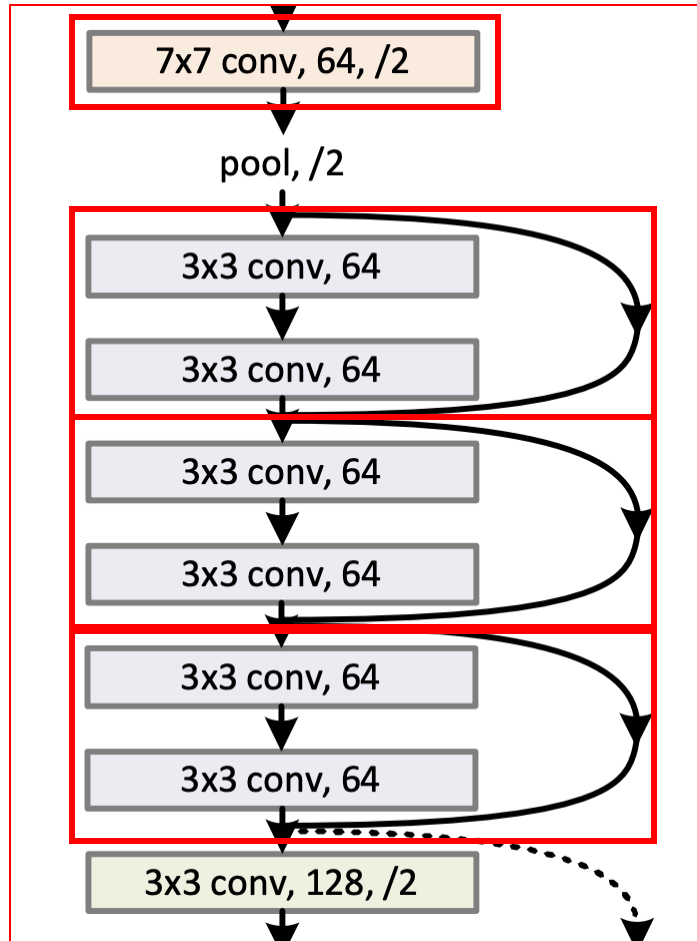


Figure 3. Example network architectures for ImageNet. **Left:** the VGG-19 model [41] (19.6 billion FLOPs) as a reference. **Middle:** a plain network with 34 parameter layers (3.6 billion FLOPs). **Right:** a residual network with 34 parameter layers (3.6 billion FLOPs). The dotted shortcuts increase dimensions. **Table 1** shows more details and other variants.

# ResNet



34-layer residual

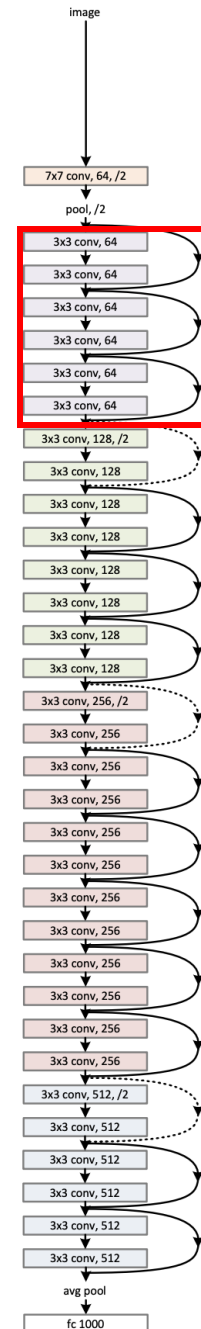


Figure 3. Example network architectures for ImageNet. **Left:** the VGG-19 model [41] (19.6 billion FLOPs) as a reference. **Middle:** a plain network with 34 parameter layers (3.6 billion FLOPs). **Right:** a residual network with 34 parameter layers (3.6 billion FLOPs). The dotted shortcuts increase dimensions. **Table 1** shows more details and other variants.







# ResNet

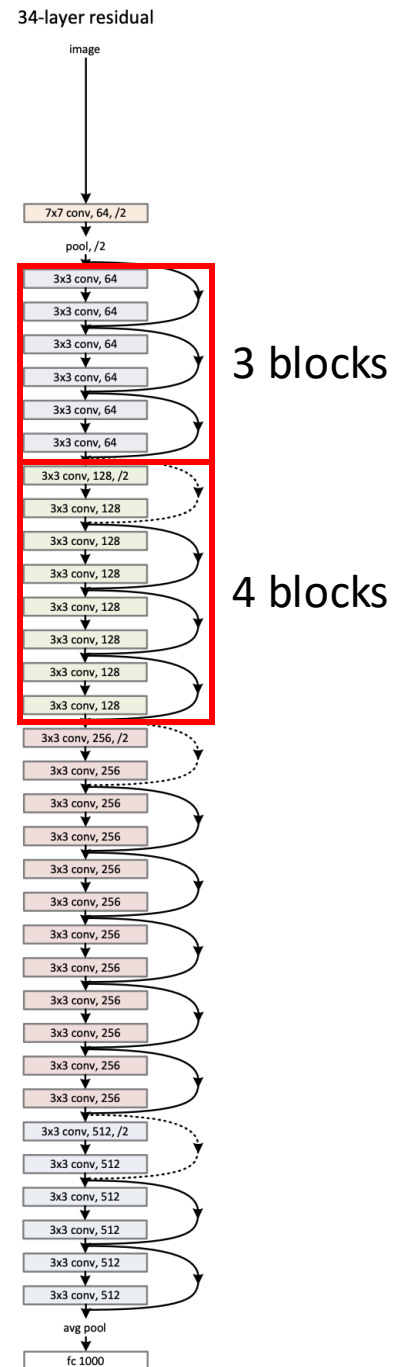


Figure 3. Example network architectures for ImageNet. **Left:** the VGG-19 model [41] (19.6 billion FLOPs) as a reference. **Middle:** a plain network with 34 parameter layers (3.6 billion FLOPs). **Right:** a residual network with 34 parameter layers (3.6 billion FLOPs). The dotted shortcuts increase dimensions. **Table 1** shows more details and other variants.

# ResNet

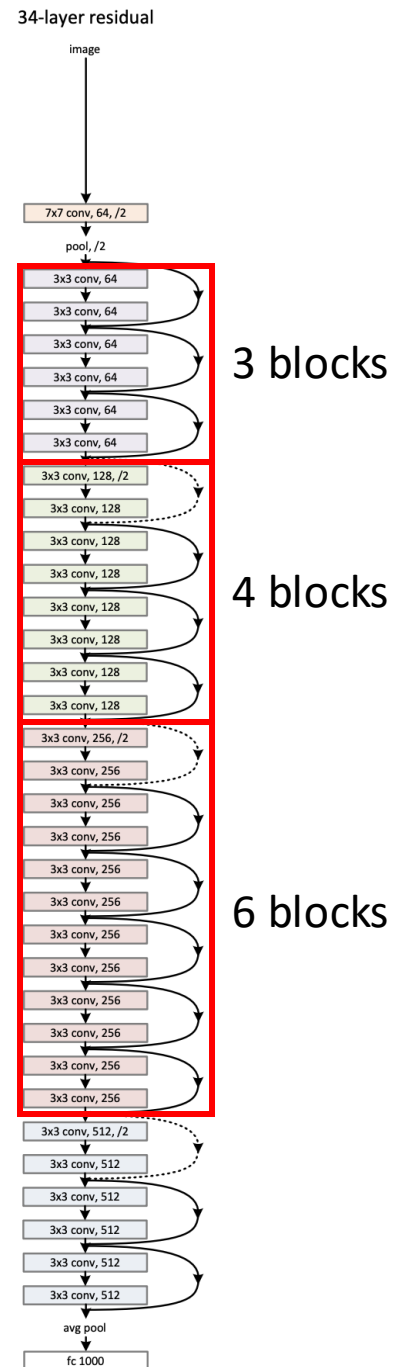


Figure 3. Example network architectures for ImageNet. **Left:** the VGG-19 model [41] (19.6 billion FLOPs) as a reference. **Middle:** a plain network with 34 parameter layers (3.6 billion FLOPs). **Right:** a residual network with 34 parameter layers (3.6 billion FLOPs). The dotted shortcuts increase dimensions. **Table 1** shows more details and other variants.

# ResNet

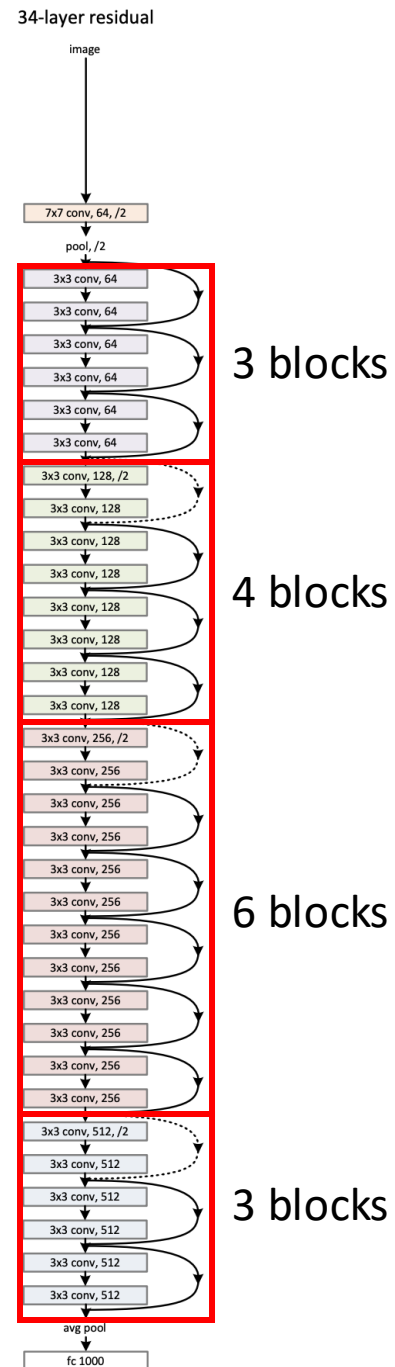


Figure 3. Example network architectures for ImageNet. **Left:** the VGG-19 model [41] (19.6 billion FLOPs) as a reference. **Middle:** a plain network with 34 parameter layers (3.6 billion FLOPs). **Right:** a residual network with 34 parameter layers (3.6 billion FLOPs). The dotted shortcuts increase dimensions. **Table 1** shows more details and other variants.

# ResNet

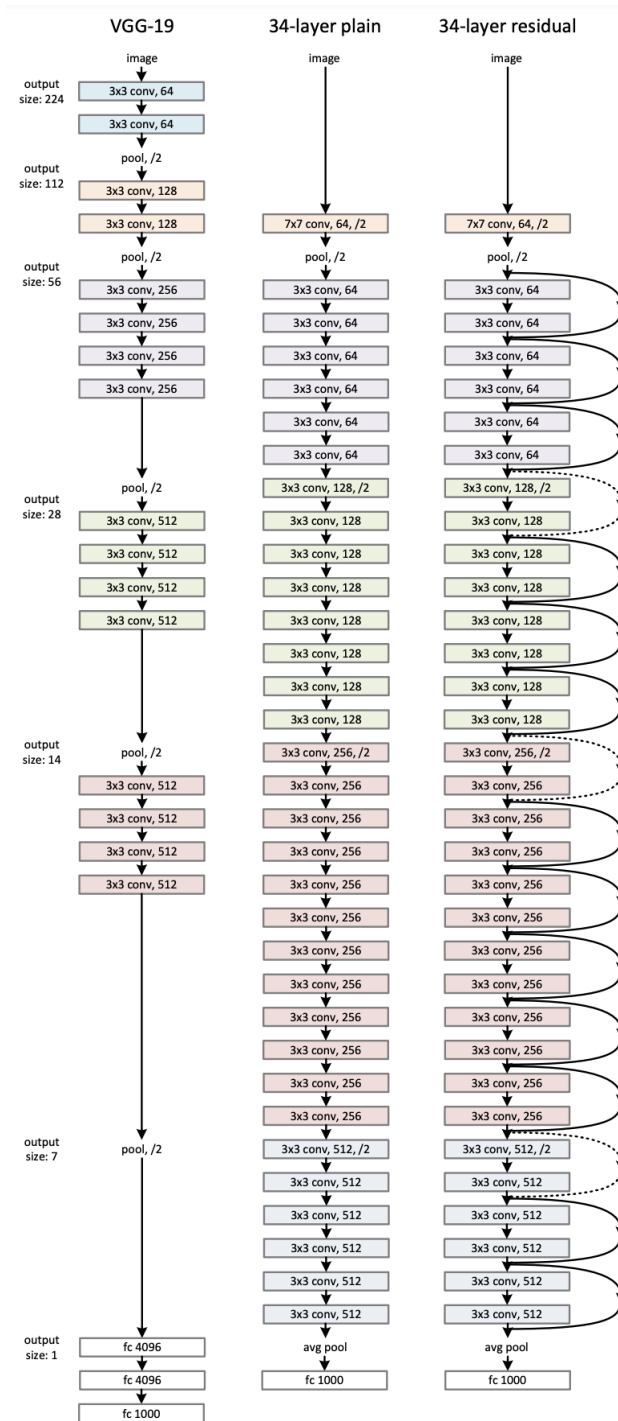


Figure 3. Example network architectures for ImageNet. **Left:** the VGG-19 model [41] (19.6 billion FLOPs) as a reference. **Middle:** a plain network with 34 parameter layers (3.6 billion FLOPs). **Right:** a residual network with 34 parameter layers (3.6 billion FLOPs). The dotted shortcuts increase dimensions. **Table 1** shows more details and other variants.

# ResNet

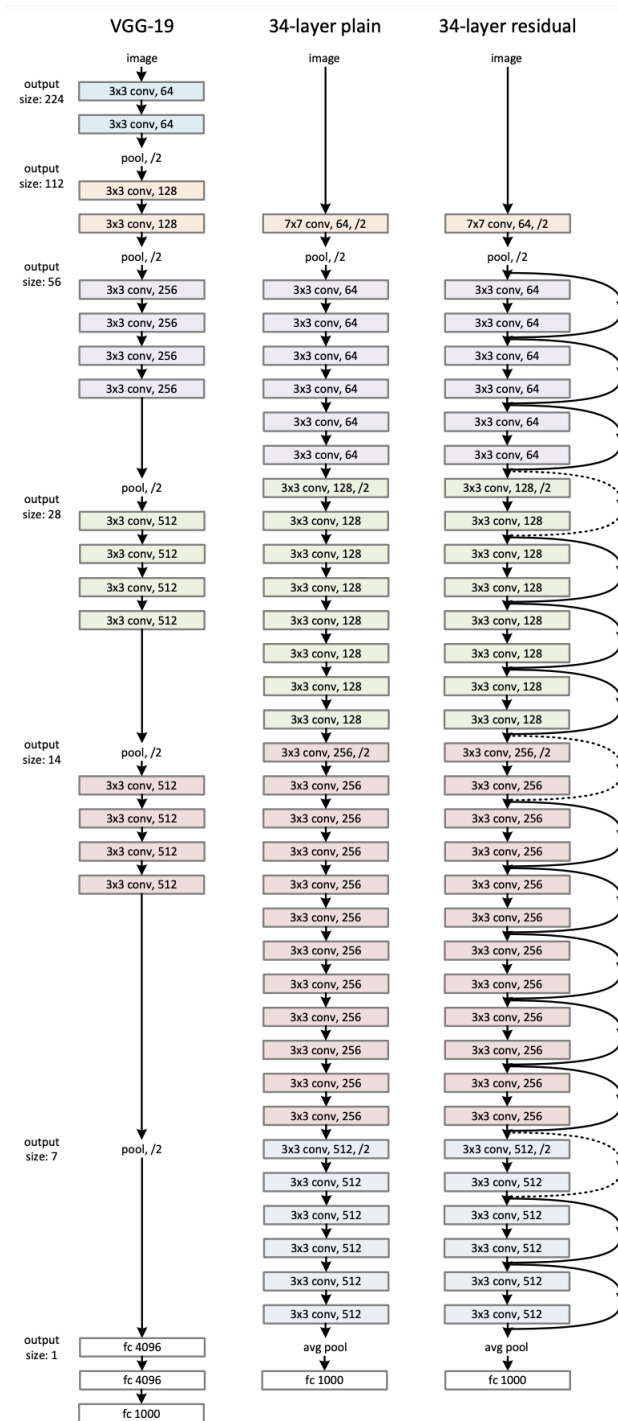


Figure 3. Example network architectures for ImageNet. **Left:** the VGG-19 model [41] (19.6 billion FLOPs) as a reference. **Middle:** a plain network with 34 parameter layers (3.6 billion FLOPs). **Right:** a residual network with 34 parameter layers (3.6 billion FLOPs). The dotted shortcuts increase dimensions. **Table 1** shows more details and other variants.



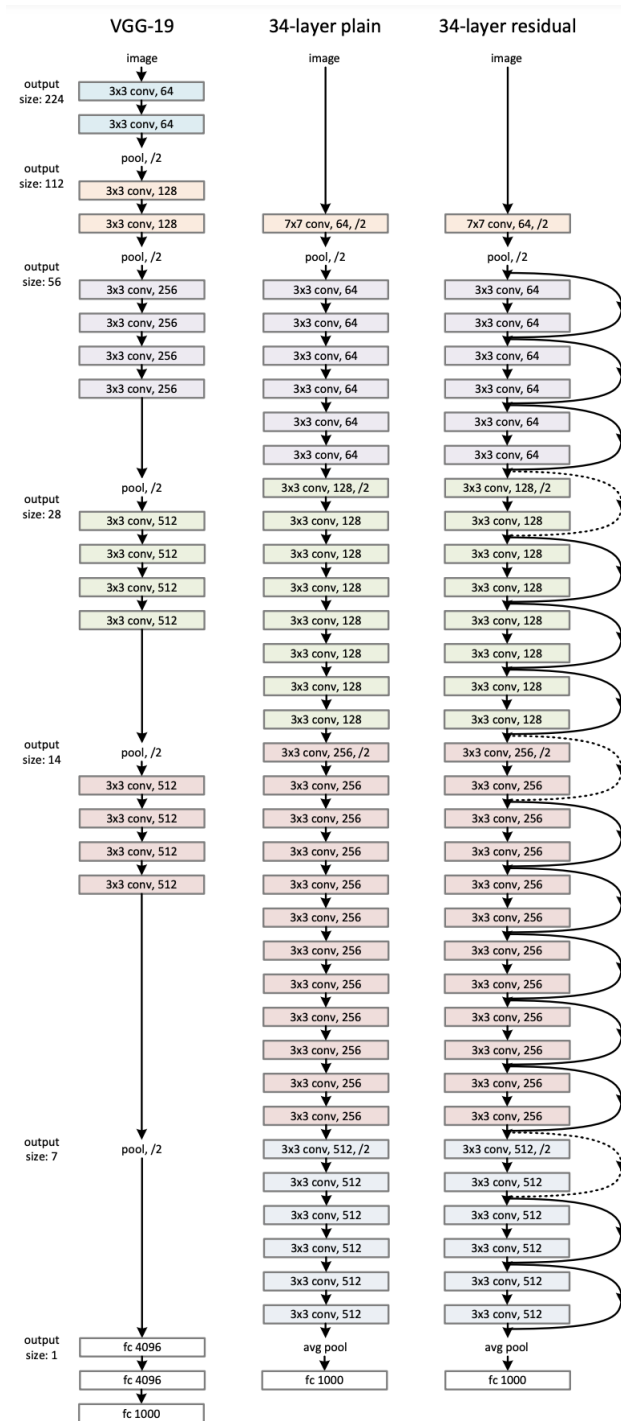








# ResNet



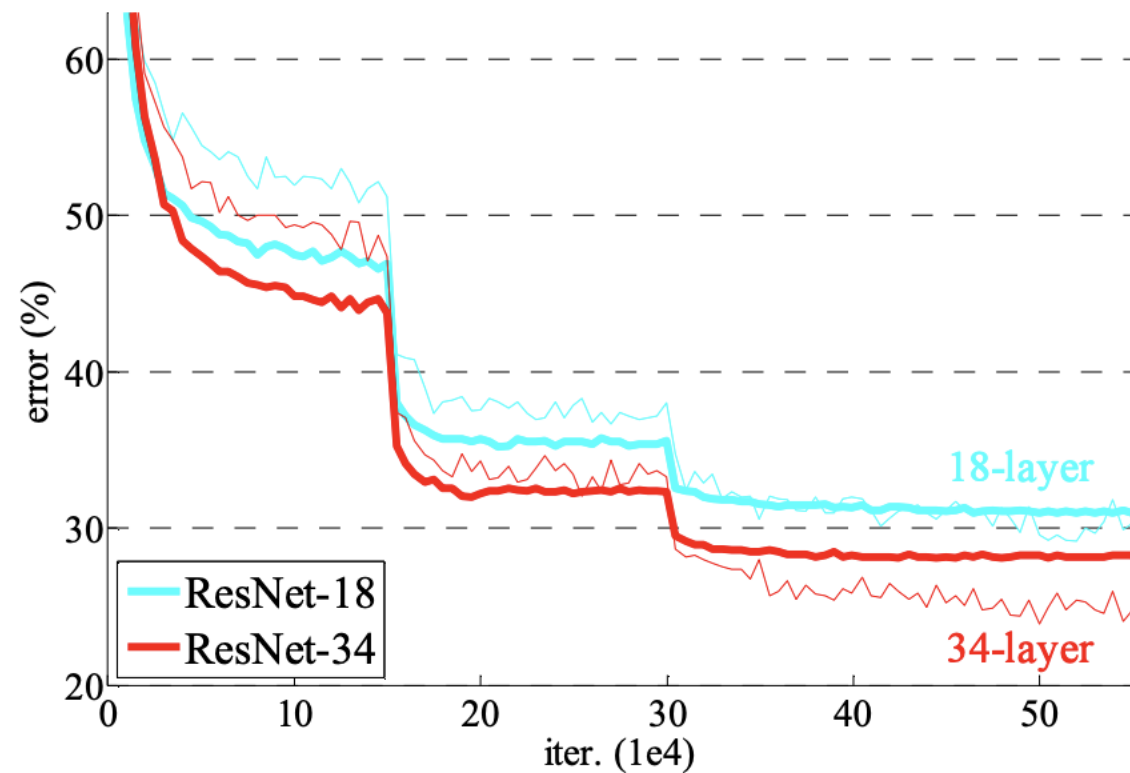
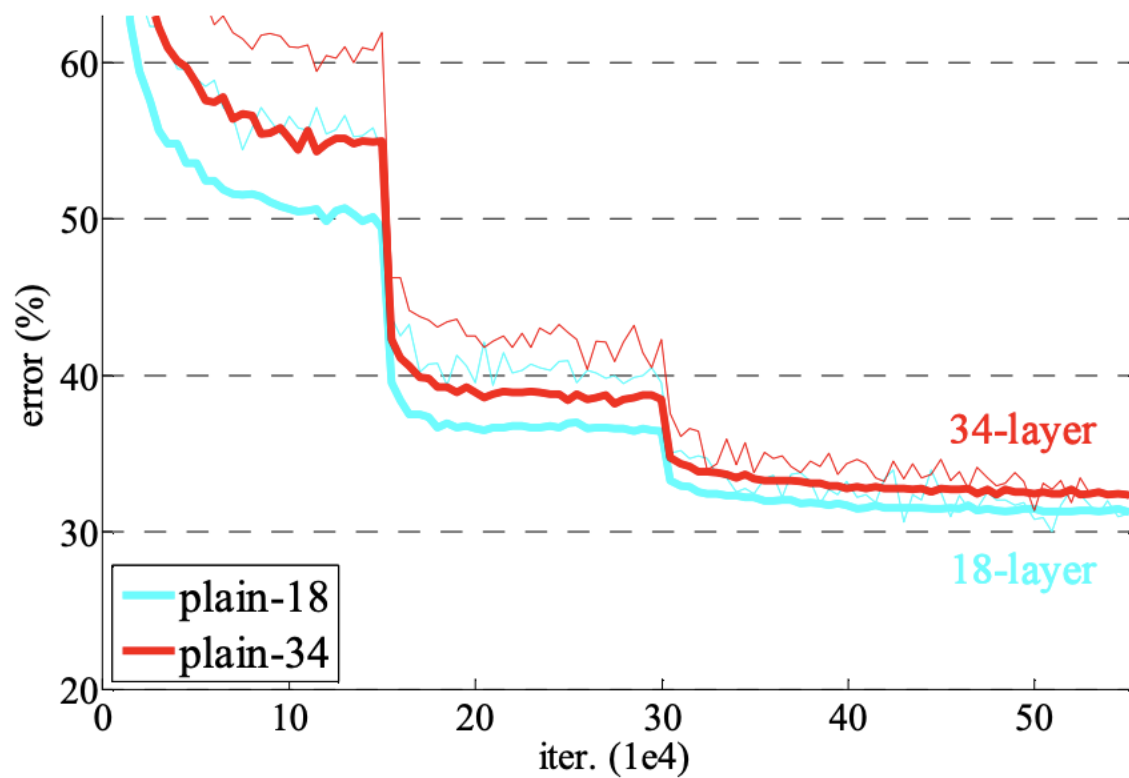
Measure: how complicated the model is

Figure 3. Example network architectures for ImageNet. **Left:** the VGG-19 model [41] (19.6 billion FLOPs) as a reference. **Middle:** a plain network with 34 parameter layers (3.6 billion FLOPs). **Right:** a residual network with 34 parameter layers (3.6 billion FLOPs). The dotted shortcuts increase dimensions. **Table 1** shows more details and other variants.

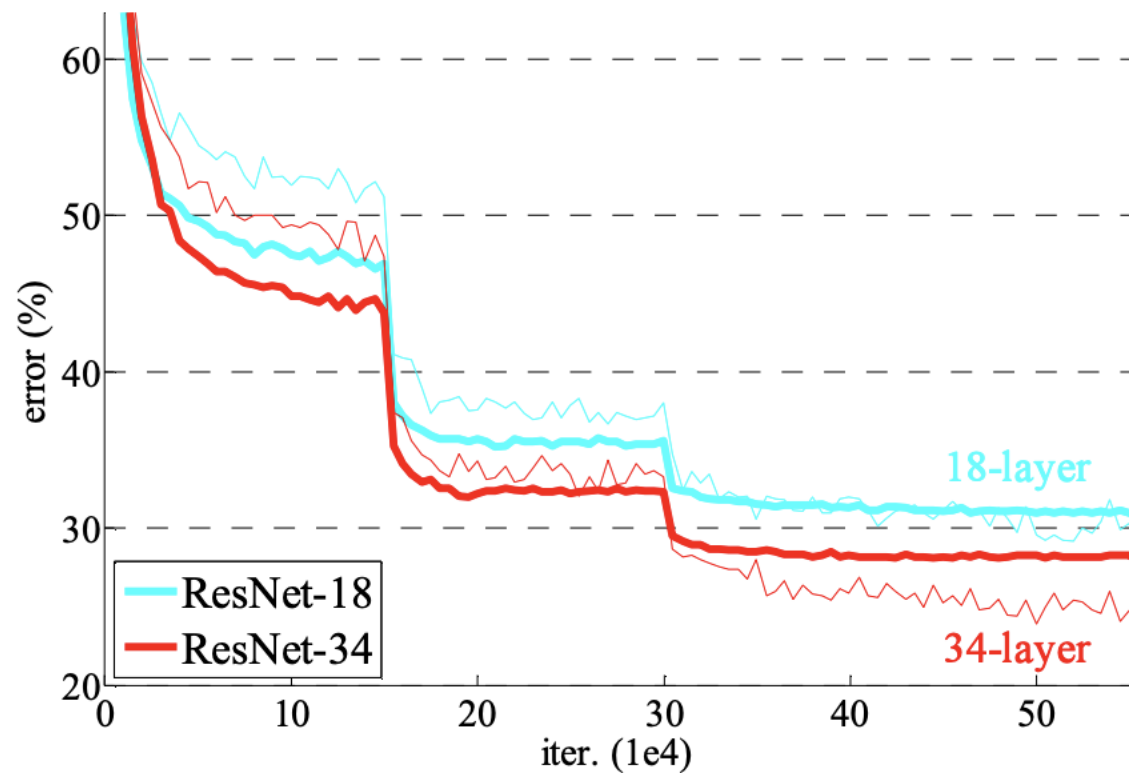
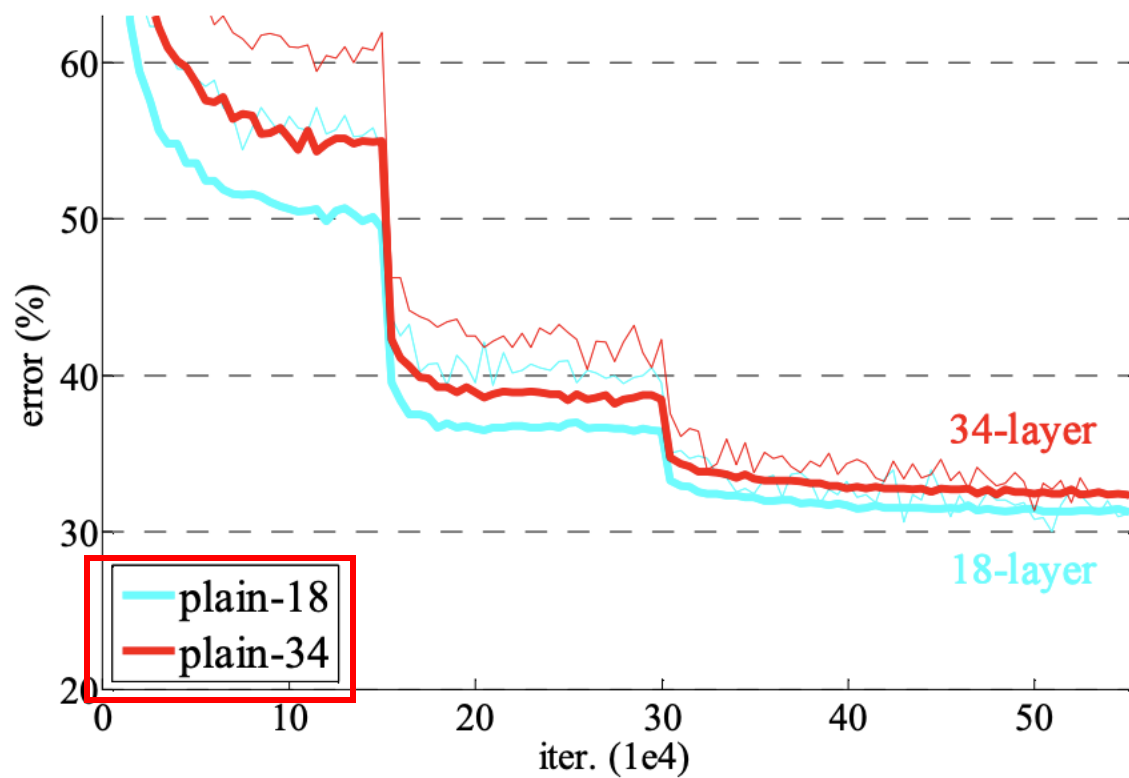
Q: VGG-19 has much more FLOPs than 34-layer plain network and 34-layer ResNet?



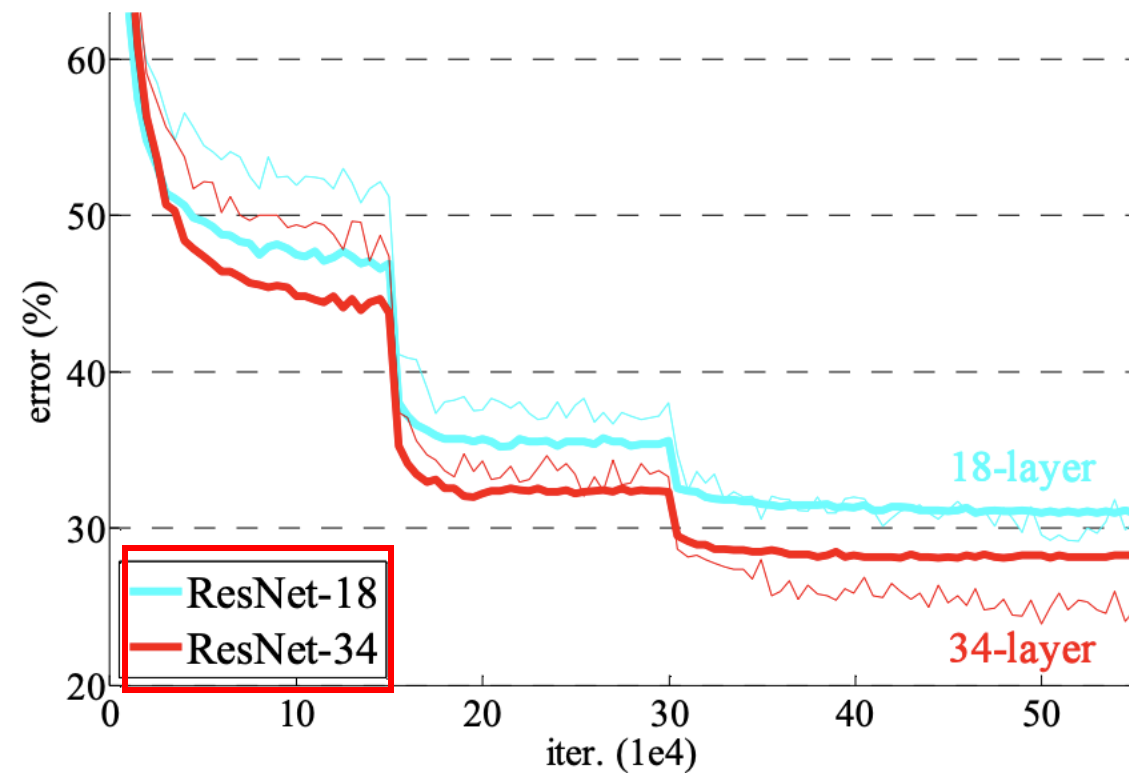
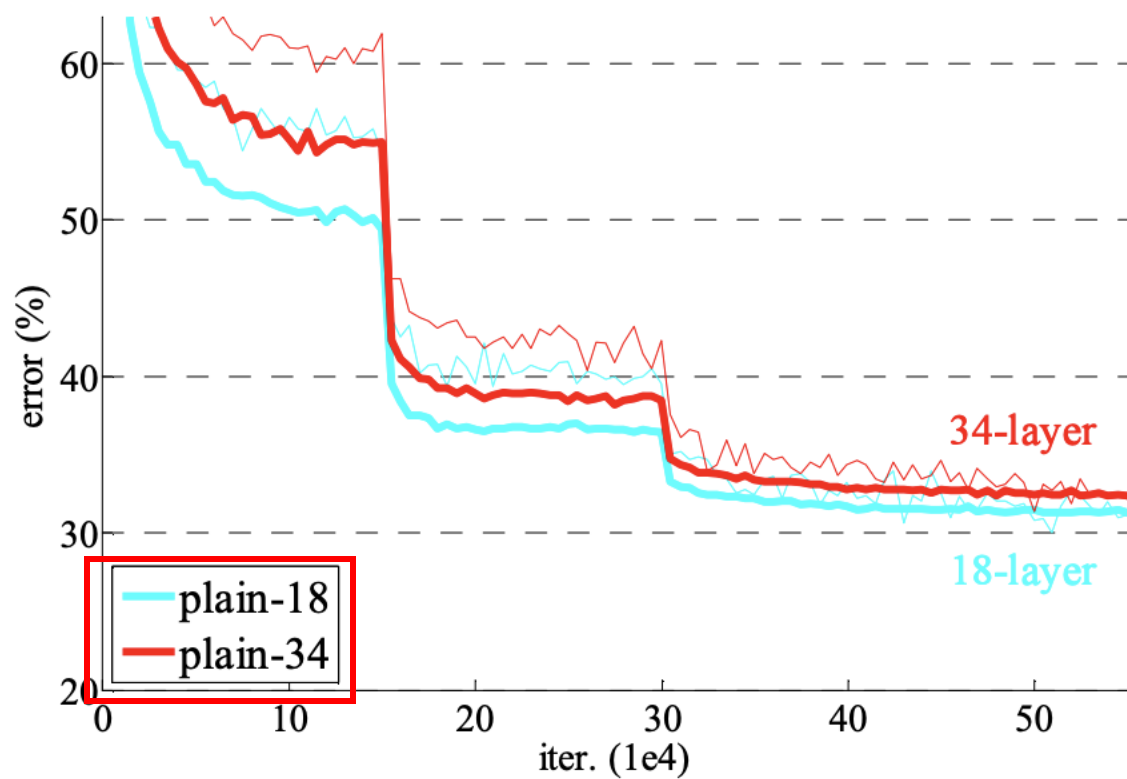
# ResNet



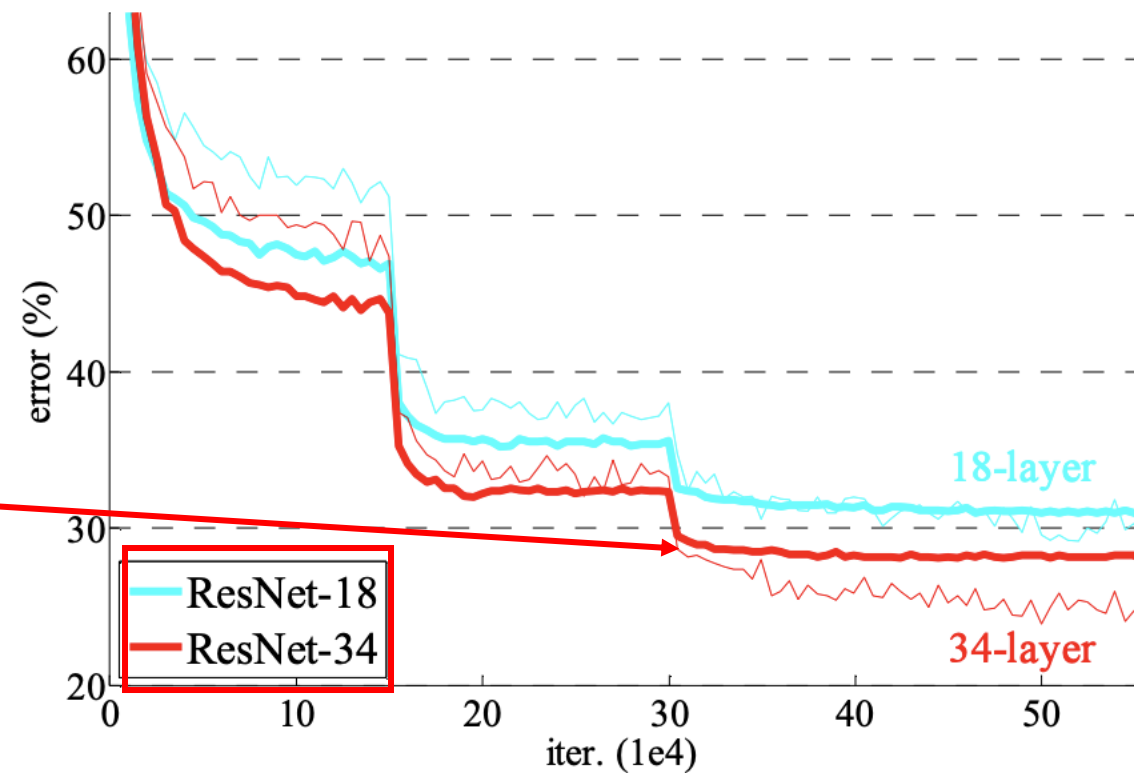
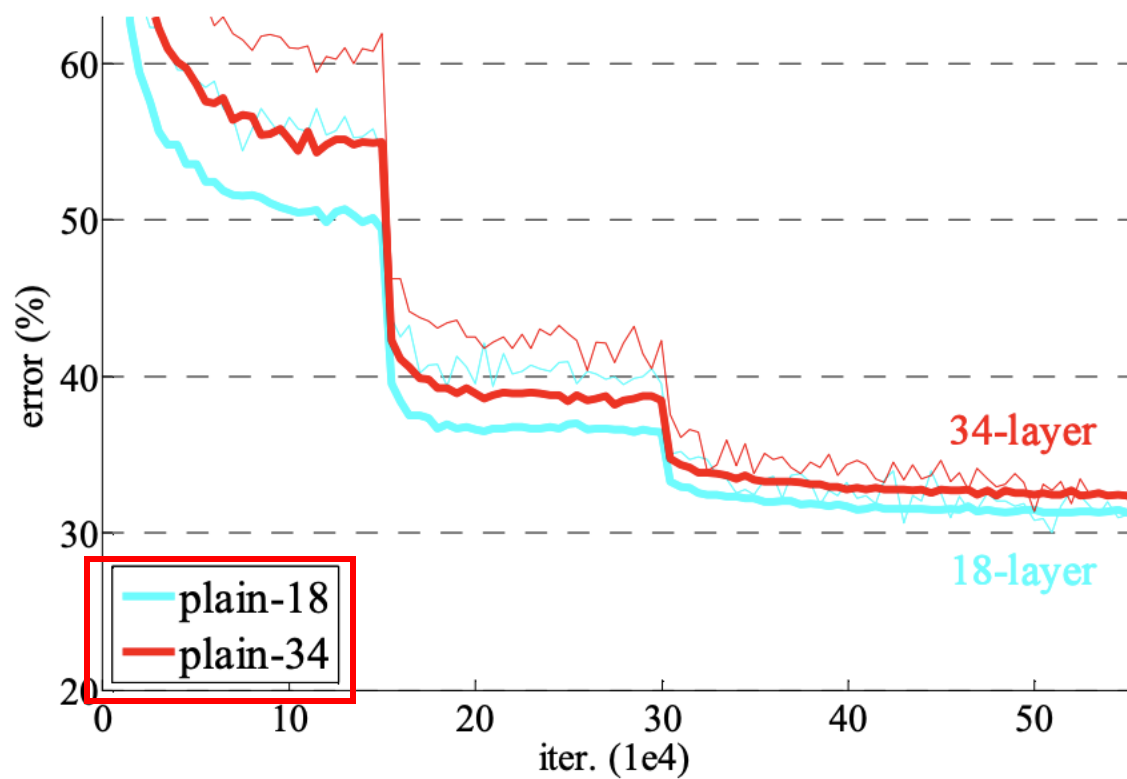
# ResNet



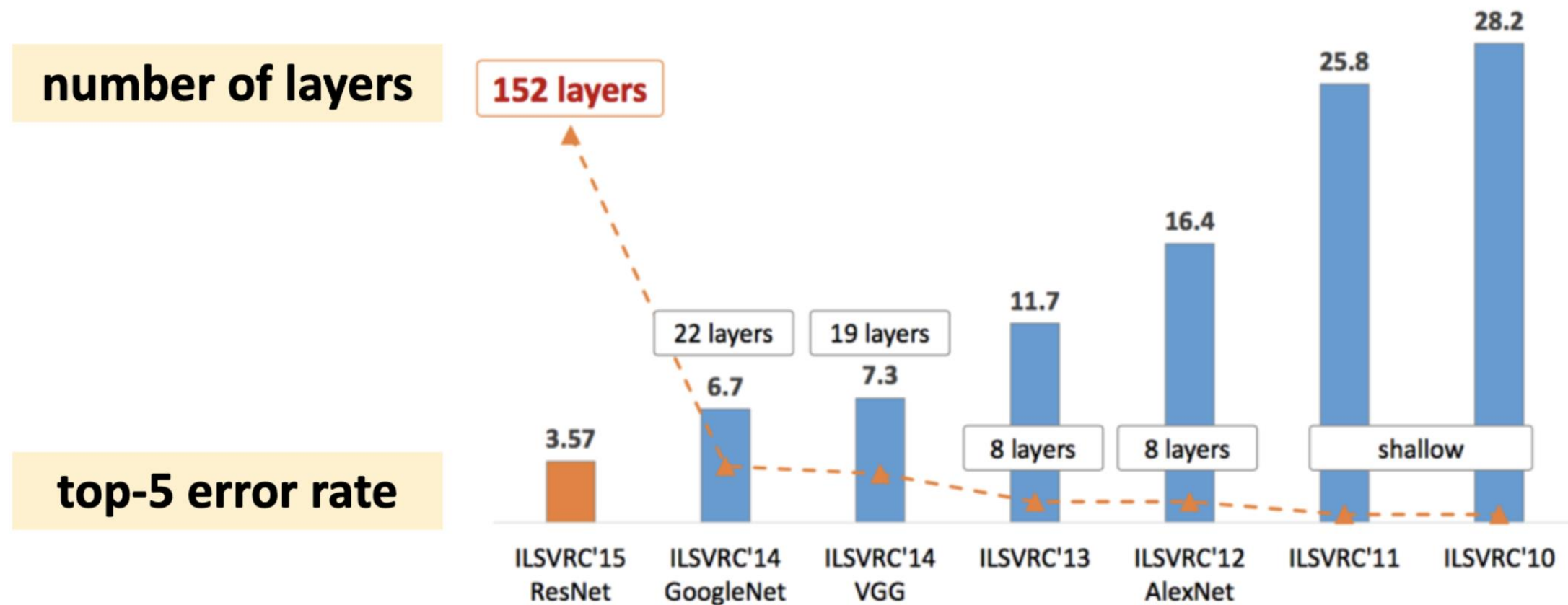
# ResNet



# ResNet

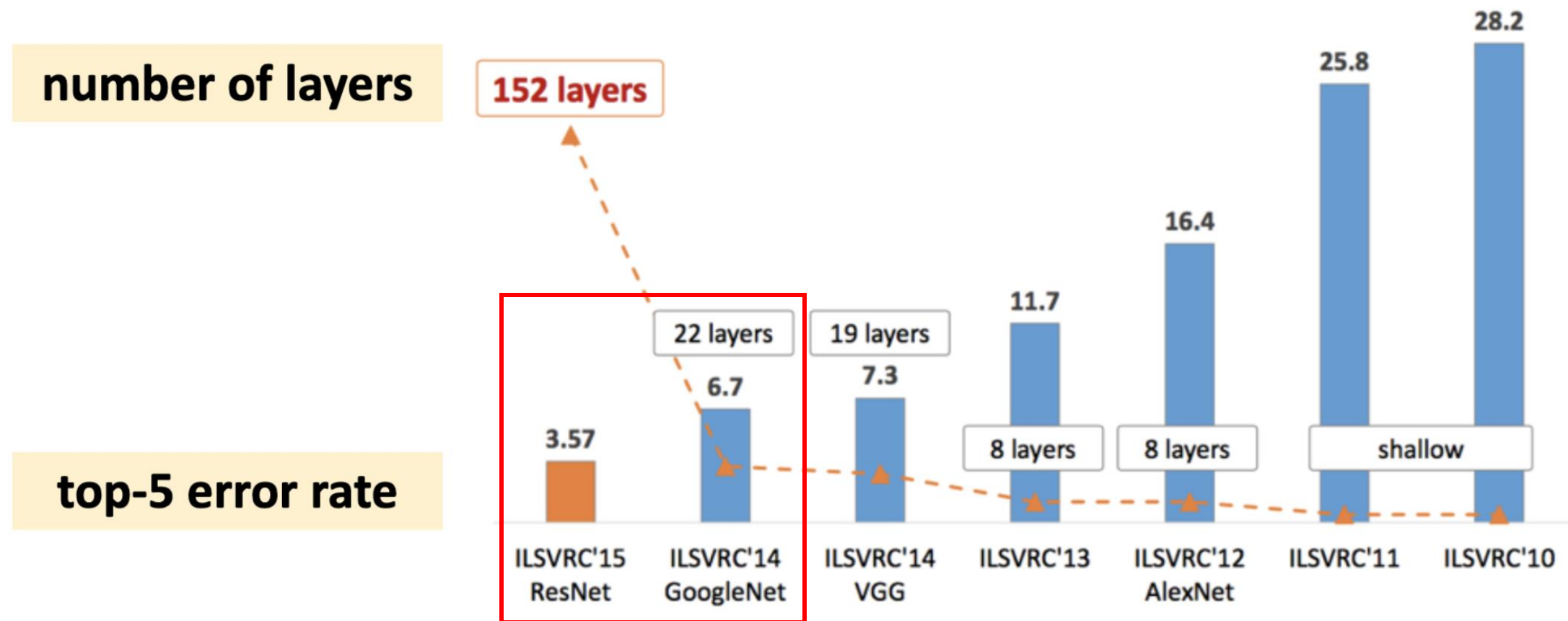


# ImageNet competition winners

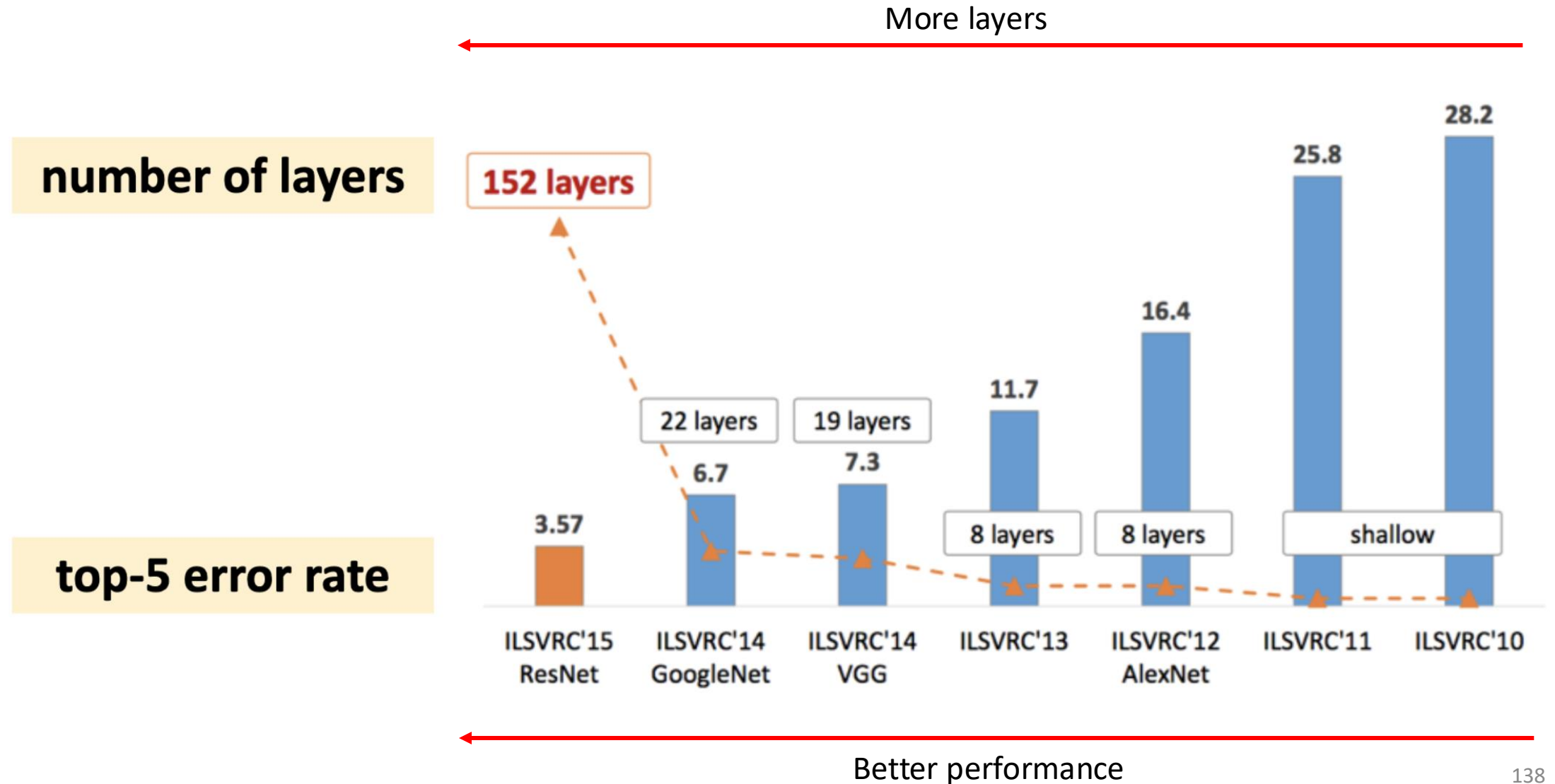




# ImageNet competition winners



# ImageNet competition winners



# References

- [LetNet-5] LeCun, Yann, Léon Bottou, Yoshua Bengio, and Patrick Haffner. "Gradient-based learning applied to document recognition." *Proceedings of the IEEE* 86, no. 11 (1998): 2278-2324.
- [AlexNet] Krizhevsky, Alex, Ilya Sutskever, and Geoffrey E. Hinton. "Imagenet classification with deep convolutional neural networks." *Advances in neural information processing systems* 25 (2012): 1097-1105.
- [VGG] Simonyan, Karen, and Andrew Zisserman. "Very deep convolutional networks for large-scale image recognition." *arXiv preprint arXiv:1409.1556* (2014). Arxiv at <https://arxiv.org/pdf/1409.1556.pdf>

# References

- [Inception] Szegedy, Christian, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. "Going deeper with convolutions." In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1-9. 2015. ArXiv at <https://arxiv.org/pdf/1409.4842.pdf> (Section 4 and 5)
- [ResNet] He, Kaiming, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. "Deep residual learning for image recognition." In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770-778. 2016. ArXiv at <https://arxiv.org/pdf/1512.03385.pdf> (Section 3.1, 3.2 and 3.3)