

What can('t) we do with ConvNets?

Classification architectures + Semantic segmentation

Karel Zimmermann

Czech Technical University in Prague

Faculty of Electrical Engineering, Department of Cybernetics



Outline

- Architectures of classification networks
- Architectures of segmentation networks
- Architectures of regression networks
- Architectures of detection networks
- Architectures of feature matching networks

Classification results

<http://image-net.org/challenges/LSVRC/2017/index>

Label: [Steel drum](#)



Classification results

<http://image-net.org/challenges/LSVRC/2017/index>

Label: Steel drum



Output:

- Scale
- T-shirt
- Steel drum
- Drumstick
- Mud turtle



Classification results

<http://image-net.org/challenges/LSVRC/2017/index>

Label: Steel drum



Output:
Scale
T-shirt
Steel drum
Drumstick
Mud turtle



Output:
Scale
T-shirt
Giant panda
Drumstick
Mud turtle



Classification results

<http://image-net.org/challenges/LSVRC/2017/index>

Label: Steel drum



Output:

- Scale
- T-shirt
- Steel drum
- Drumstick
- Mud turtle



Output:

- Scale
- T-shirt
- Giant panda
- Drumstick
- Mud turtle



$$\text{Error} = \frac{1}{100,000} \sum_{\substack{100,000 \\ \text{images}}} 1[\text{incorrect on image } i]$$

Classification results

AlexNet

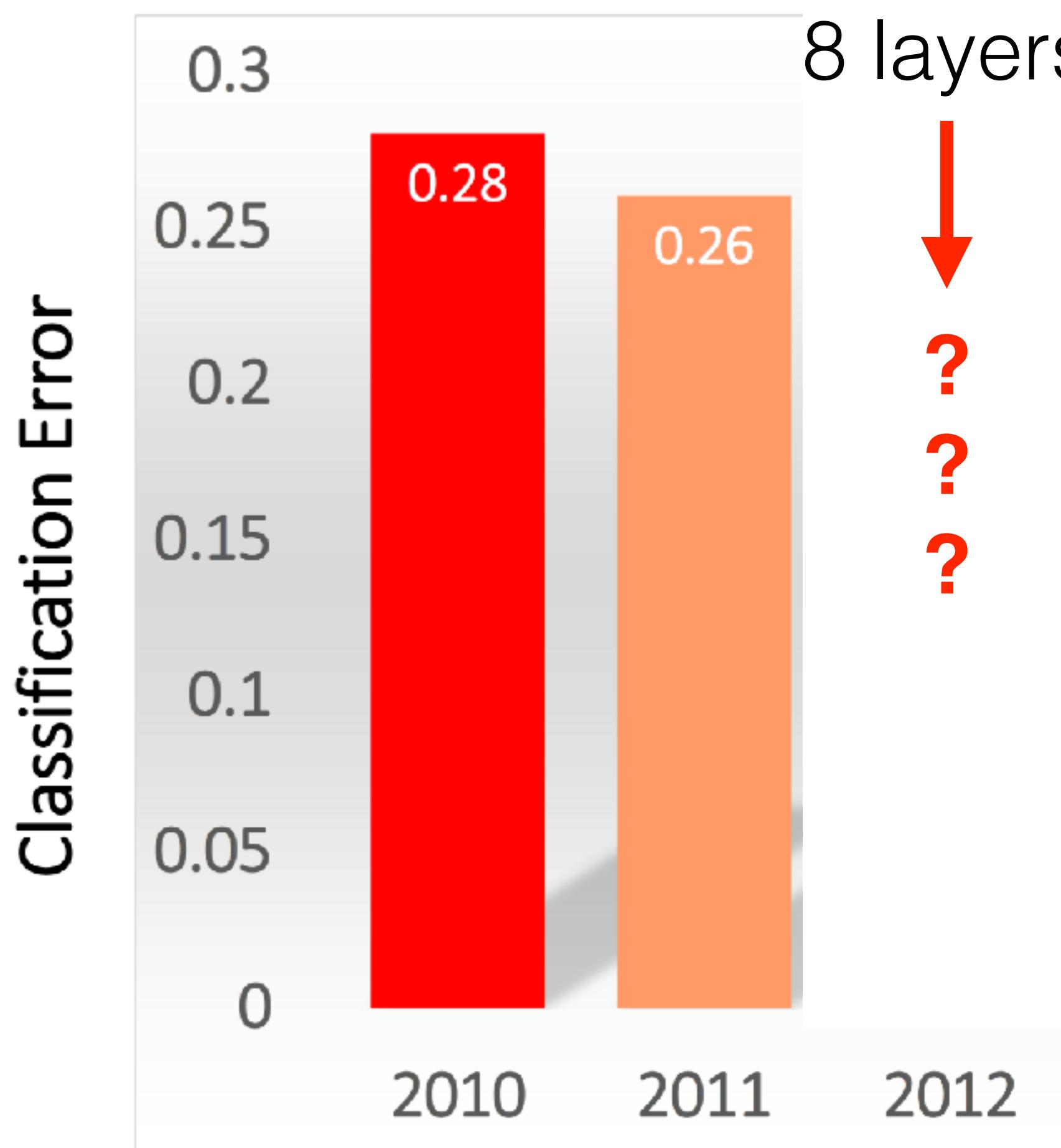
8 layers



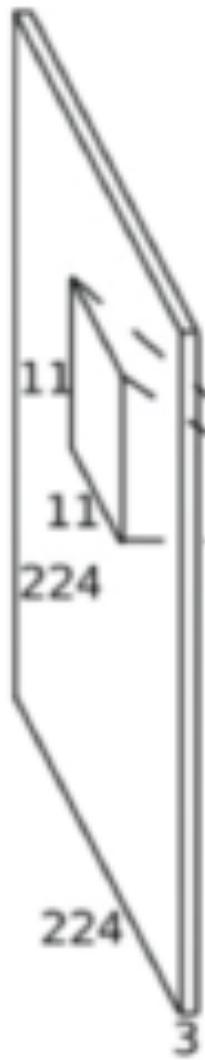
?

?

?

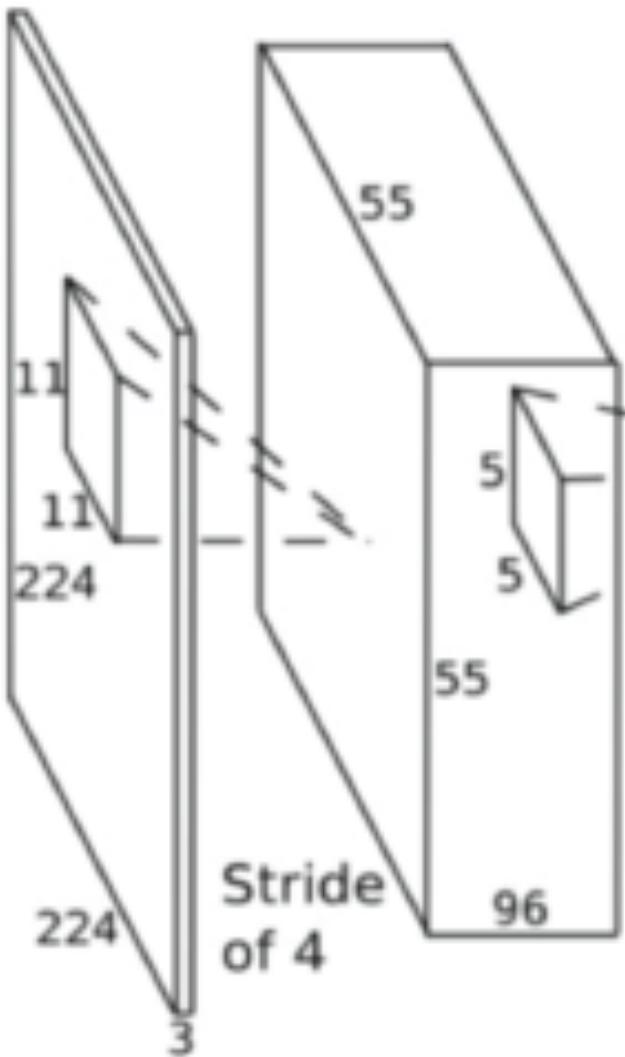


AlexNet on ImageNet 2012 (**over 27k citations !!!**)



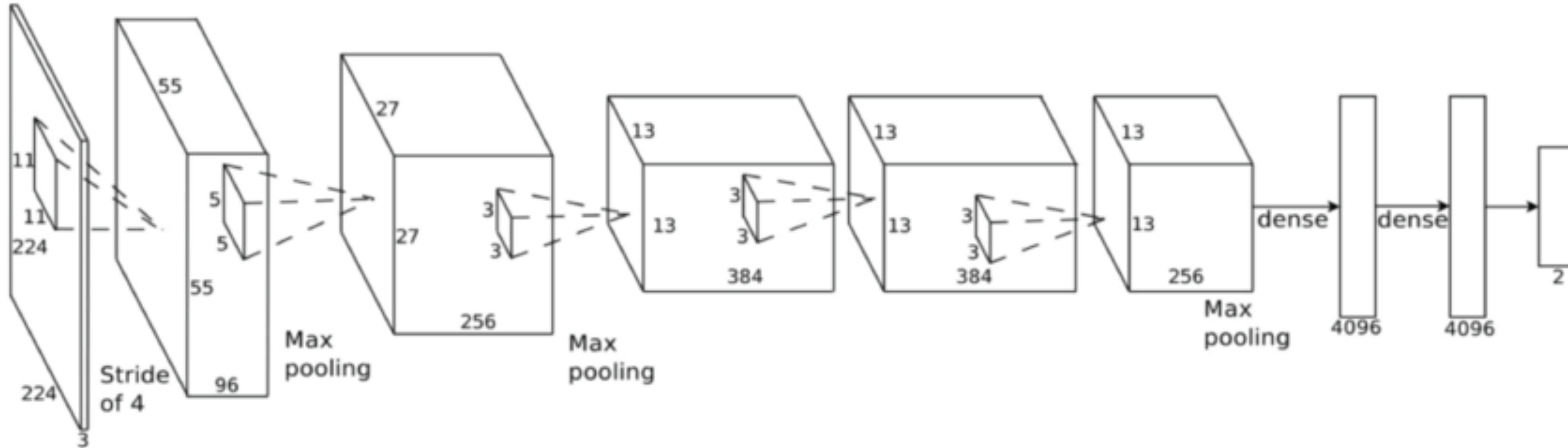
- Param in layer1 (conv, 96 11x11 filters, stride=4, pad=0)?

AlexNet on ImageNet 2012 (**over 27k citations !!!**)



- Param in layer1 (conv, 96 11x11 filters, stride=4, pad=0)?
- Param in layer2 (maxp,3x3 filters, stride=2, pad=0)?

AlexNet on ImageNet 2012 (**over 27k citations !!!**)

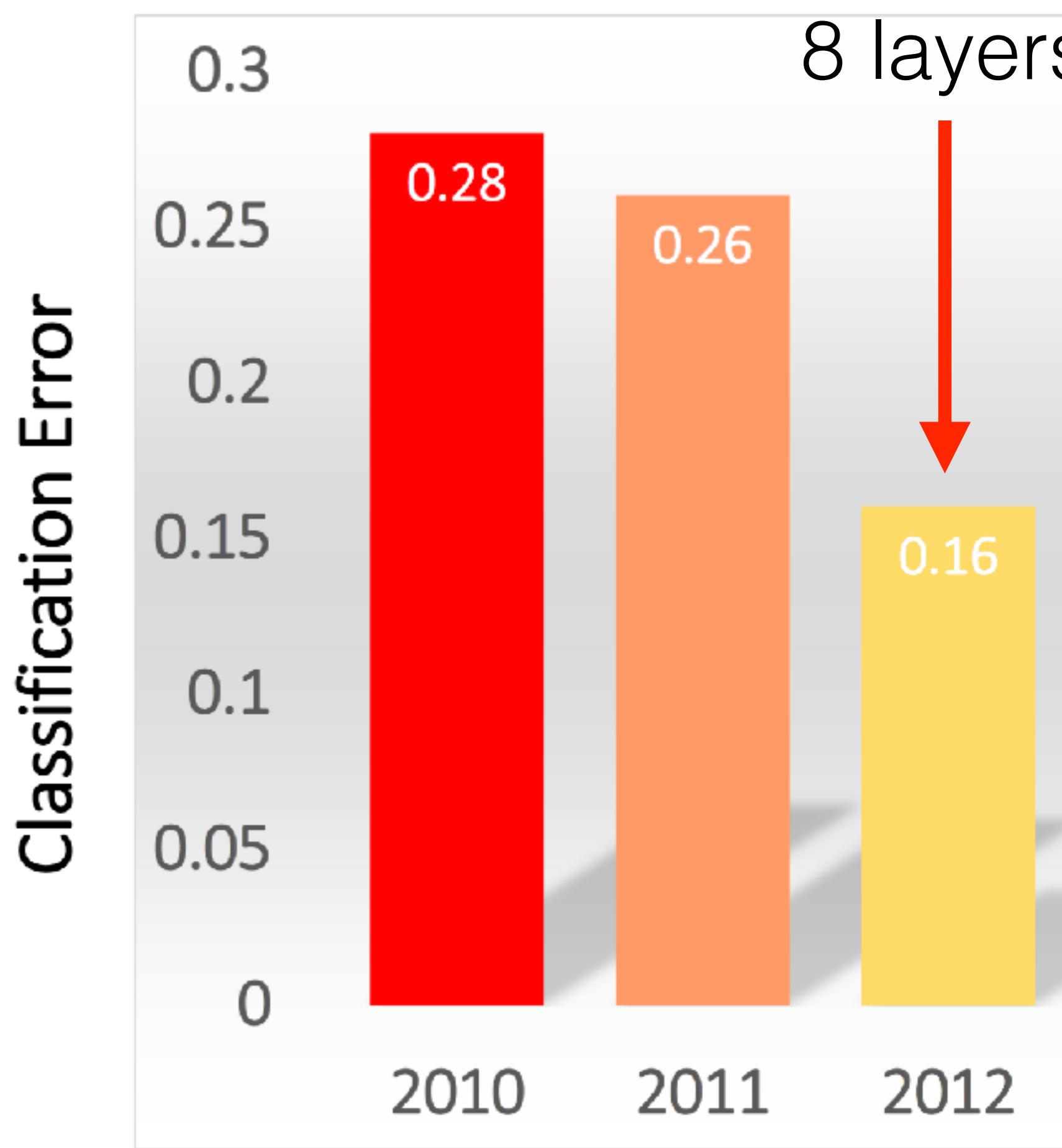


- Param in layer1 (conv, 96 11x11 filters, stride=4, pad=0)?
- Param in layer2 (maxp,3x3 filters, stride=2, pad=0)?
- Param in layer3 (conv, 256 5x5 filters, stride=1, pad=2?)
- Parameters in total: 60M, Depth: 8 layers

Classification results

AlexNet

8 layers



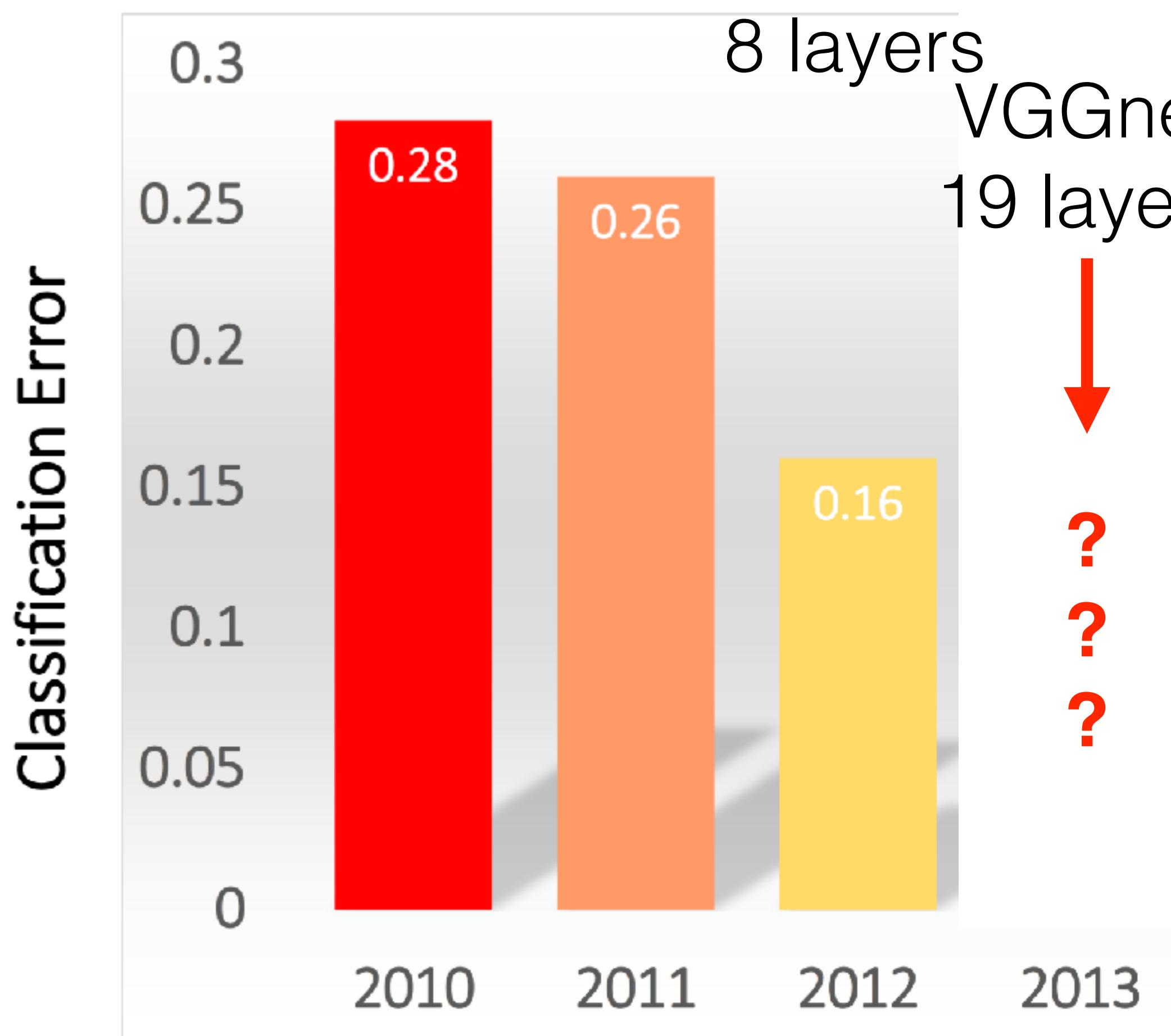
Classification results

AlexNet

8 layers

VGGnet

19 layers



VGGNet vs AlexNet

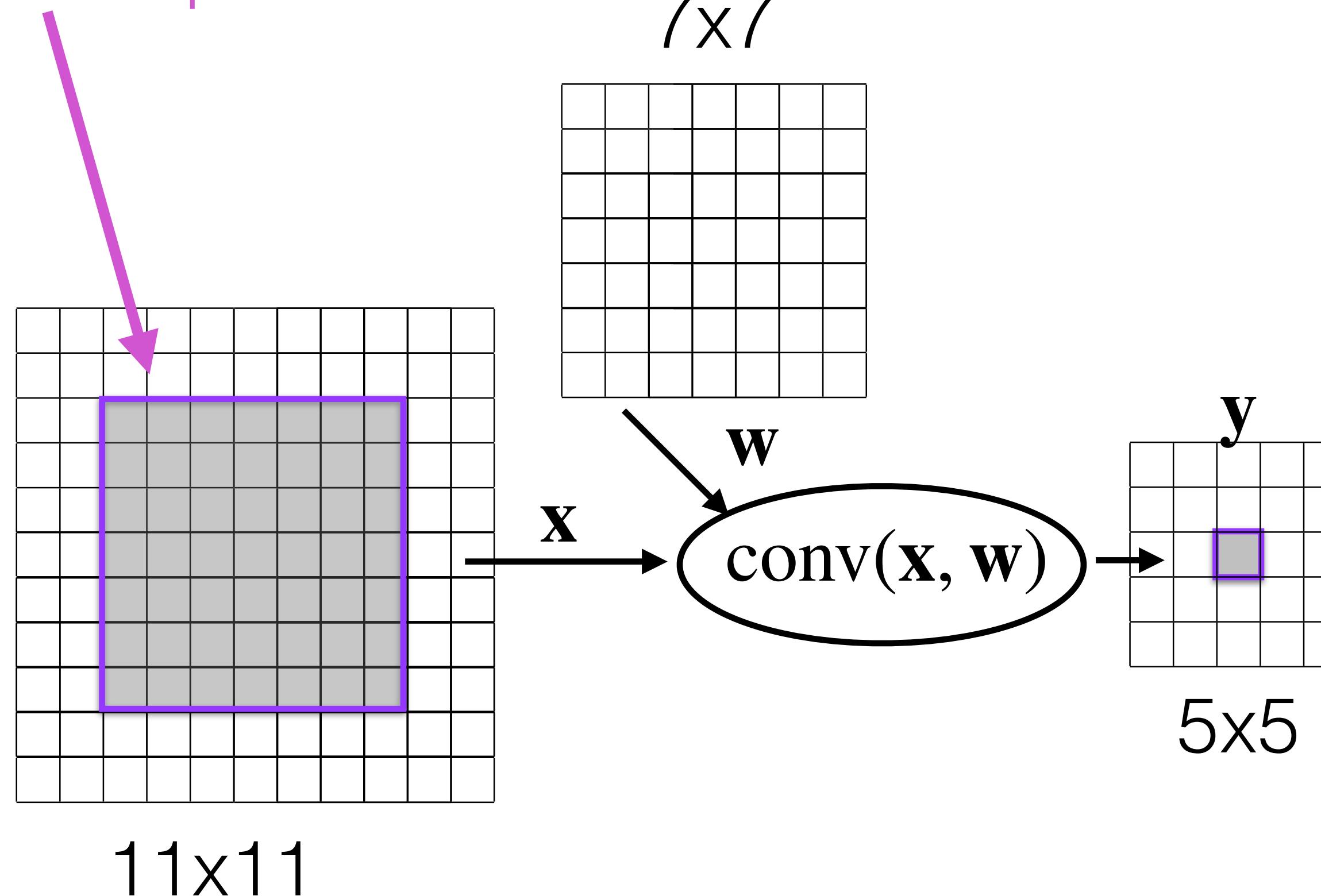


- large filters
 - shallow (8 layers)
- small filters
 - deeper (19 layers)

- Parameters in total: 138M, Depth: 19 layers

Receptive field

receptive field

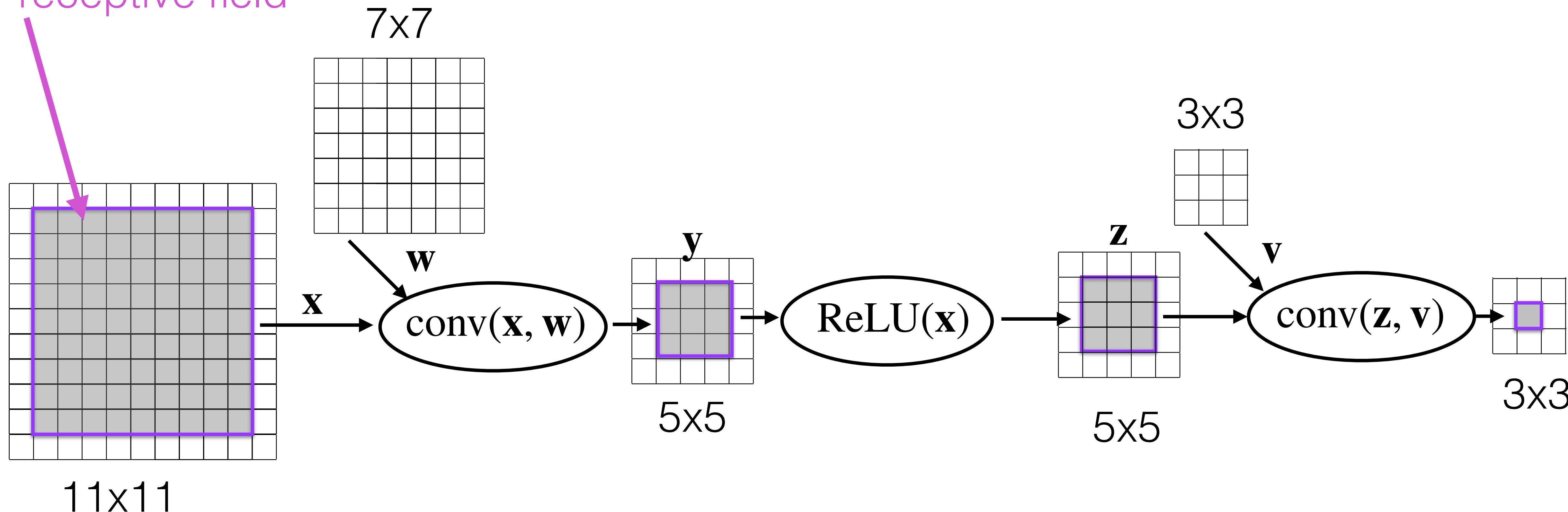


Receptive field: region in the input space that affect activation of a particular neuron.

What is the size of the receptive field??? 7×7

Receptive field

receptive field



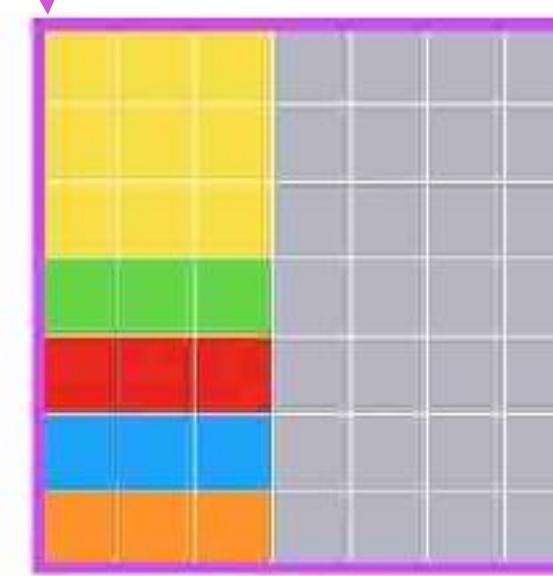
Receptive field: region in the input space that affect activation of a particular neuron.

What is the size of the receptive field??? **9x9**

VGGNet vs AlexNet

receptive field

7x7



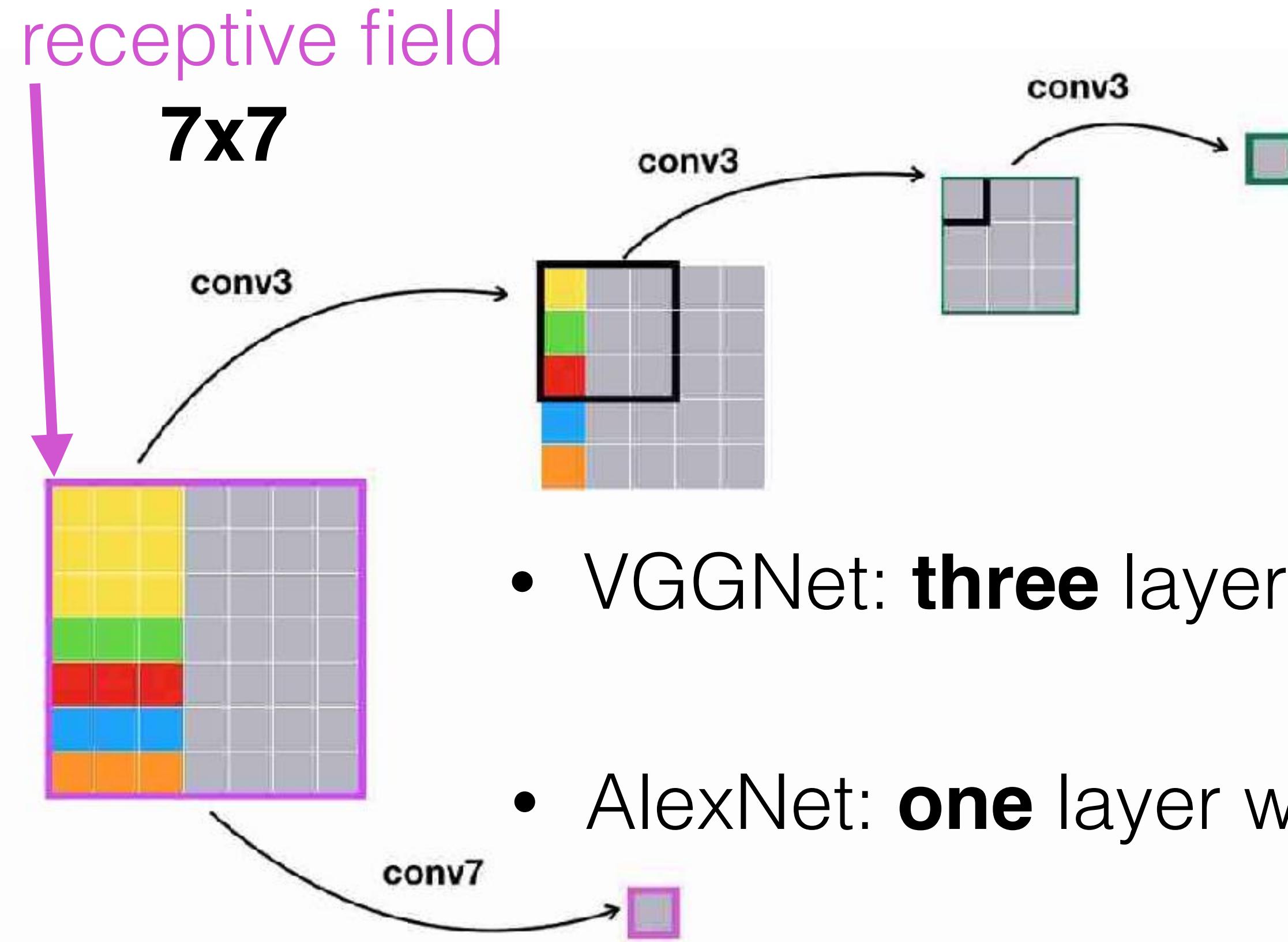
- AlexNet: **one** layer with **7x7** filter (49+1 params)

conv7



Receptive field: region in the input space that affect activation of a particular neuron.

VGGNet vs AlexNet



Receptive field: region in the input space that affect activation of a particular neuron.

VGGNet has **the same receptive field** with **less parameters**

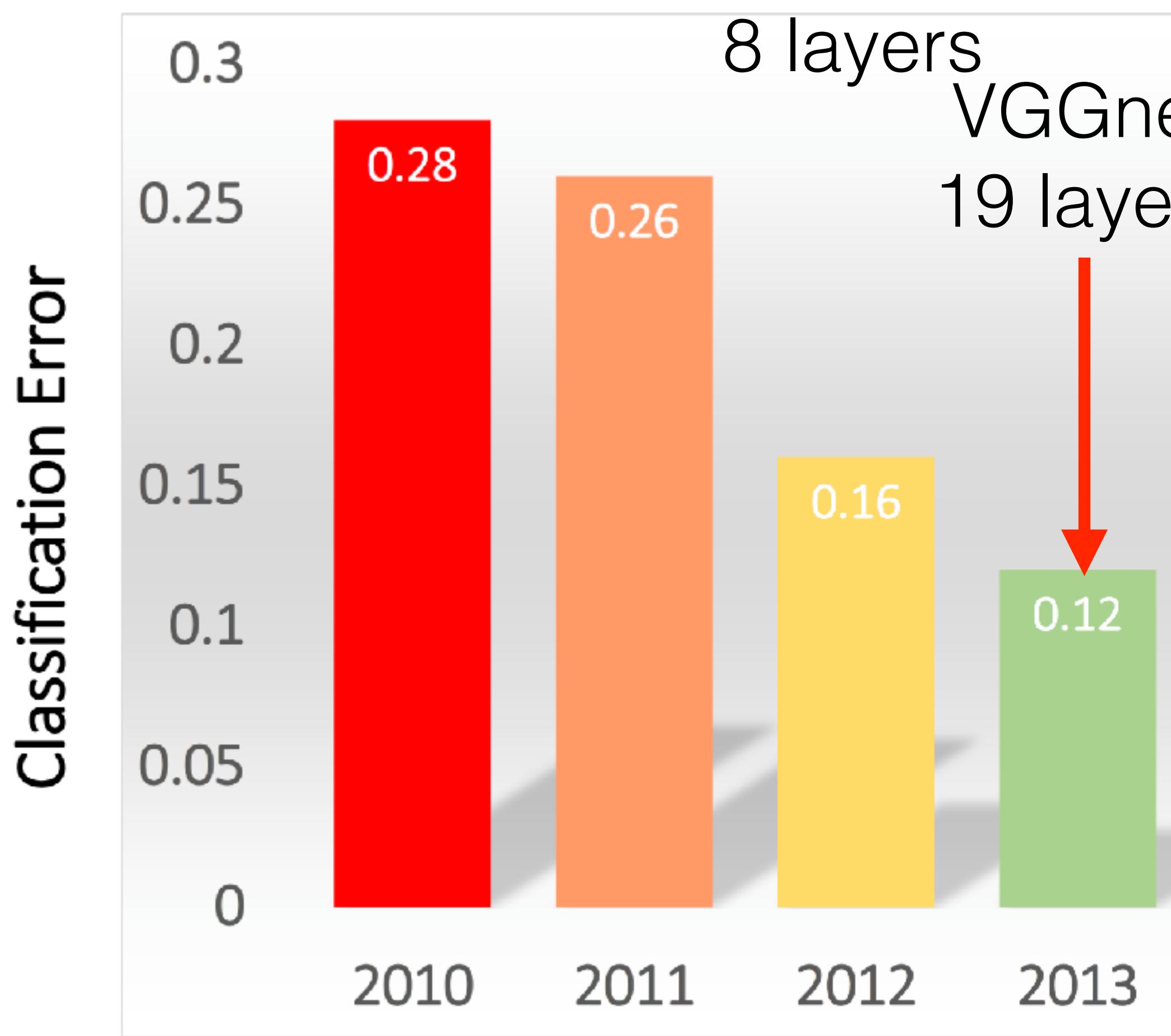
Classification results

AlexNet

8 layers

VGGnet

19 layers



Classification results

AlexNet

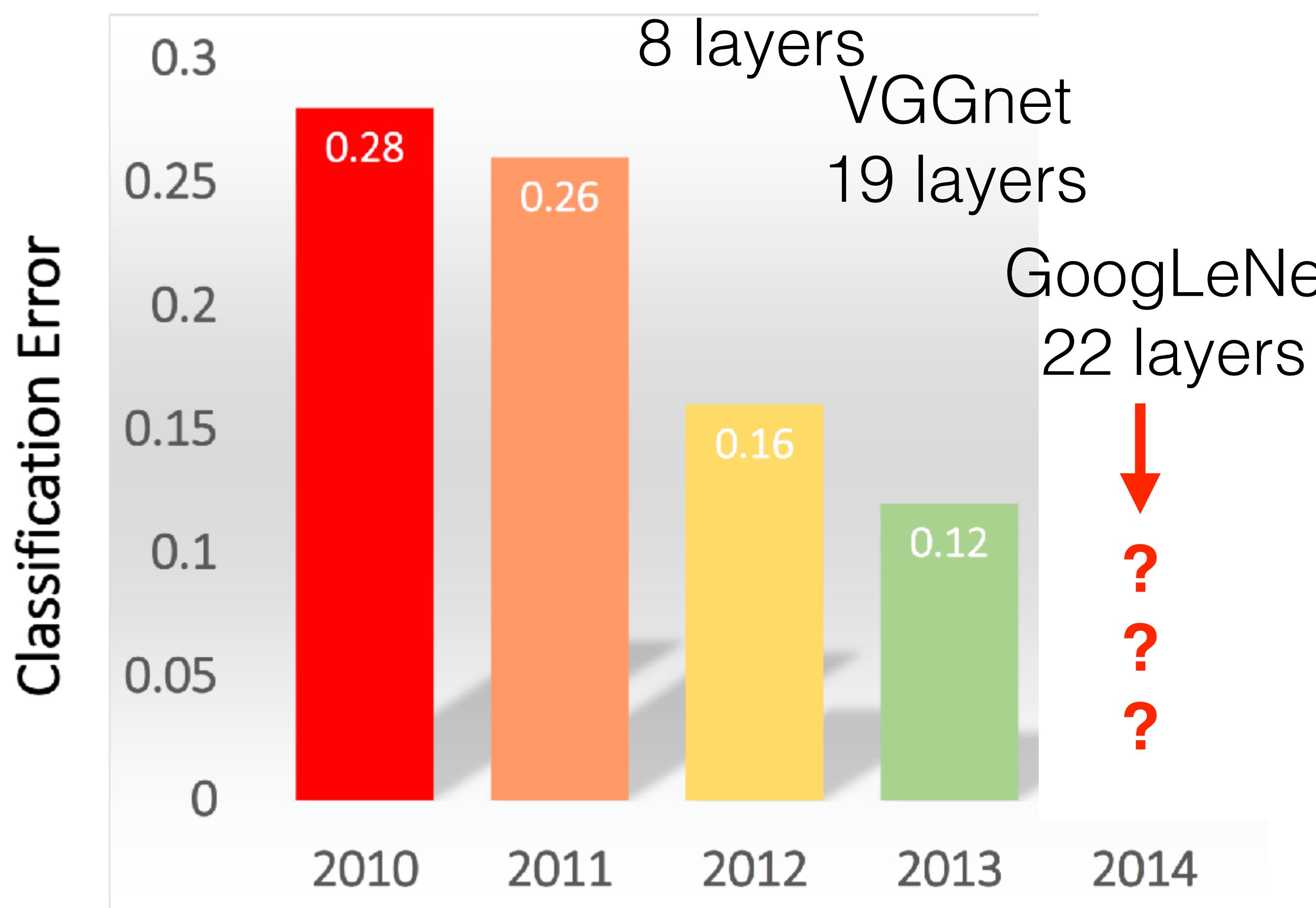
8 layers

VGGnet

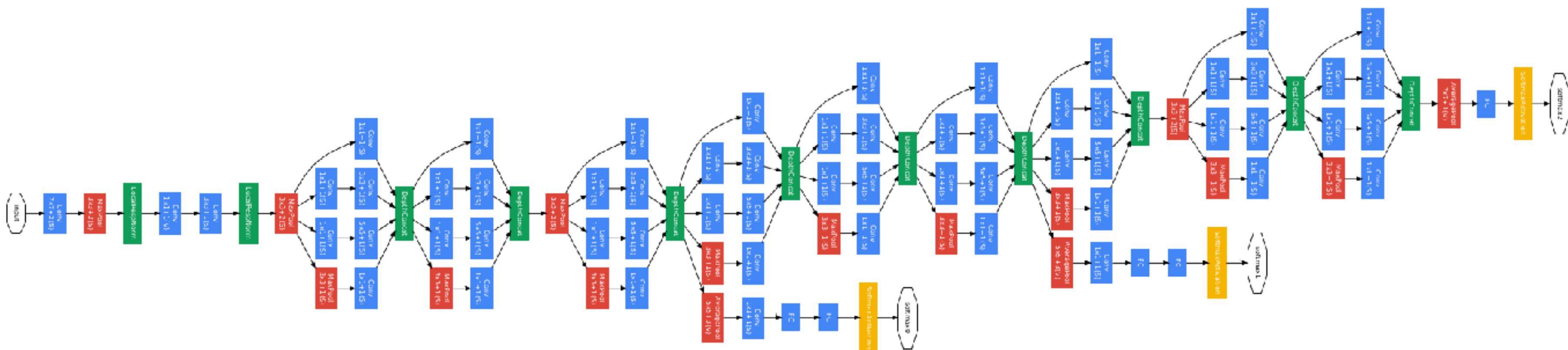
19 layers

GoogLeNet

22 layers

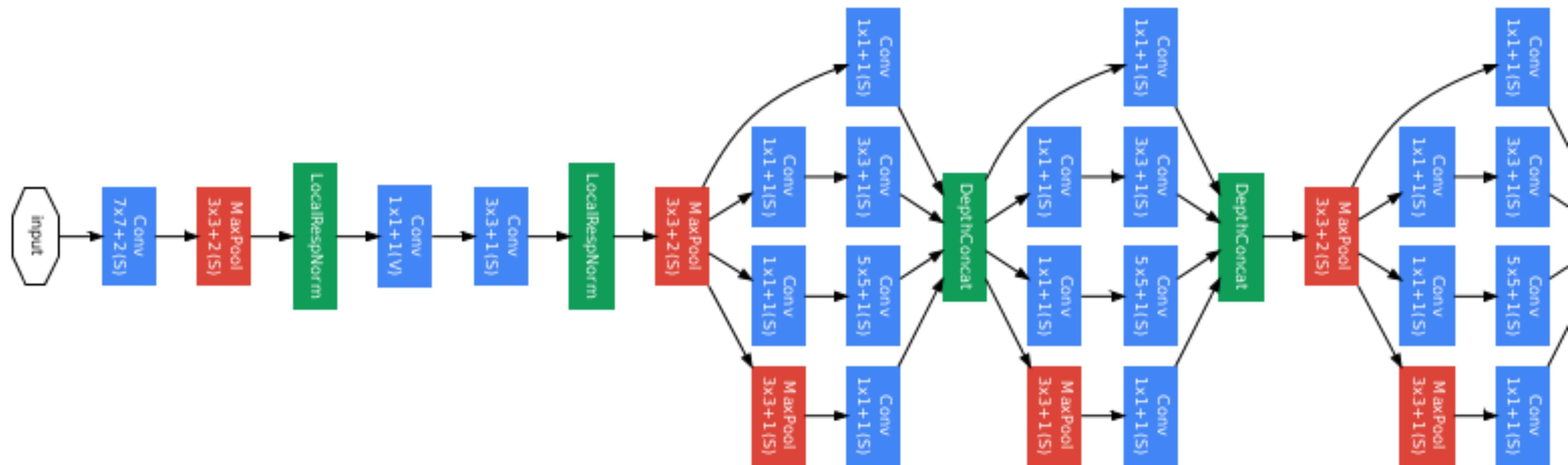


GoogLeNe



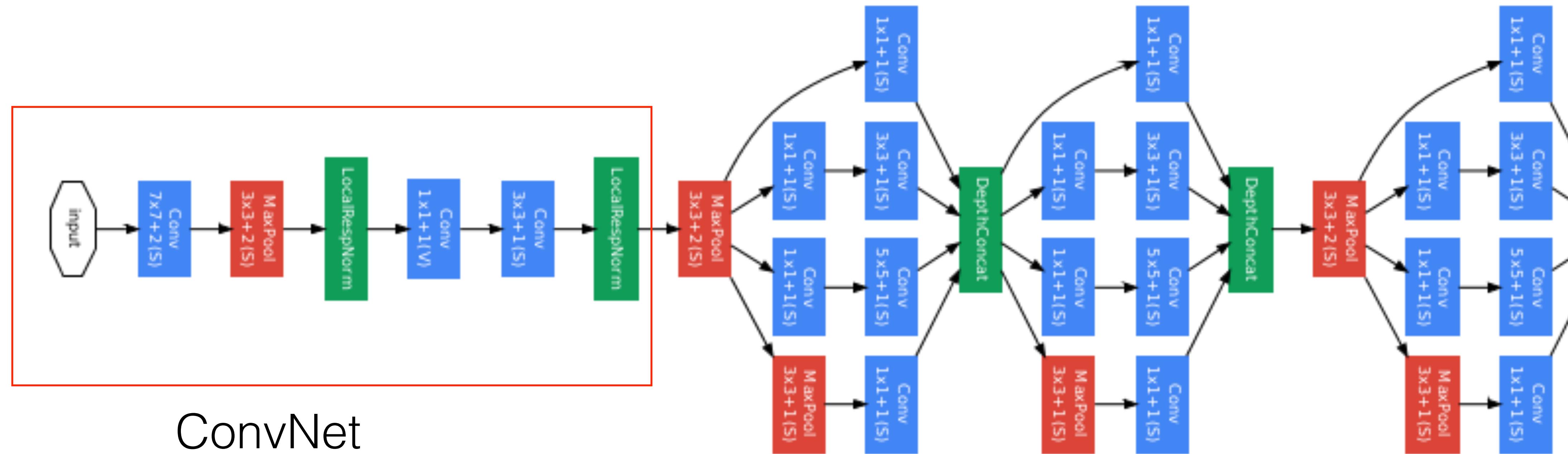
Szegedy et al. Going Deeper with Convolutions, CVPR, 2014
<https://arxiv.org/abs/1409.4842>

GoogLeNet



Szegedy et al. Going Deeper with Convolutions, CVPR, 2014
<https://arxiv.org/abs/1409.4842>

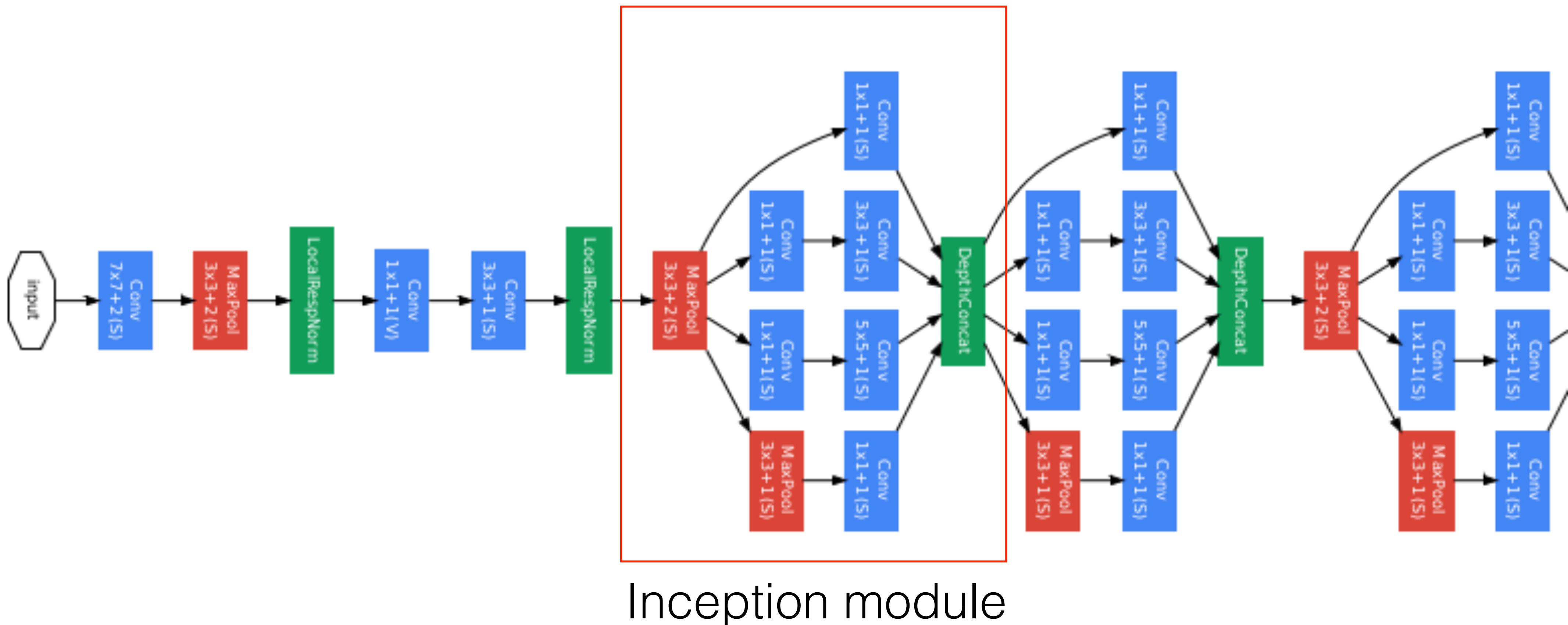
GoogLeNet



ConvNet

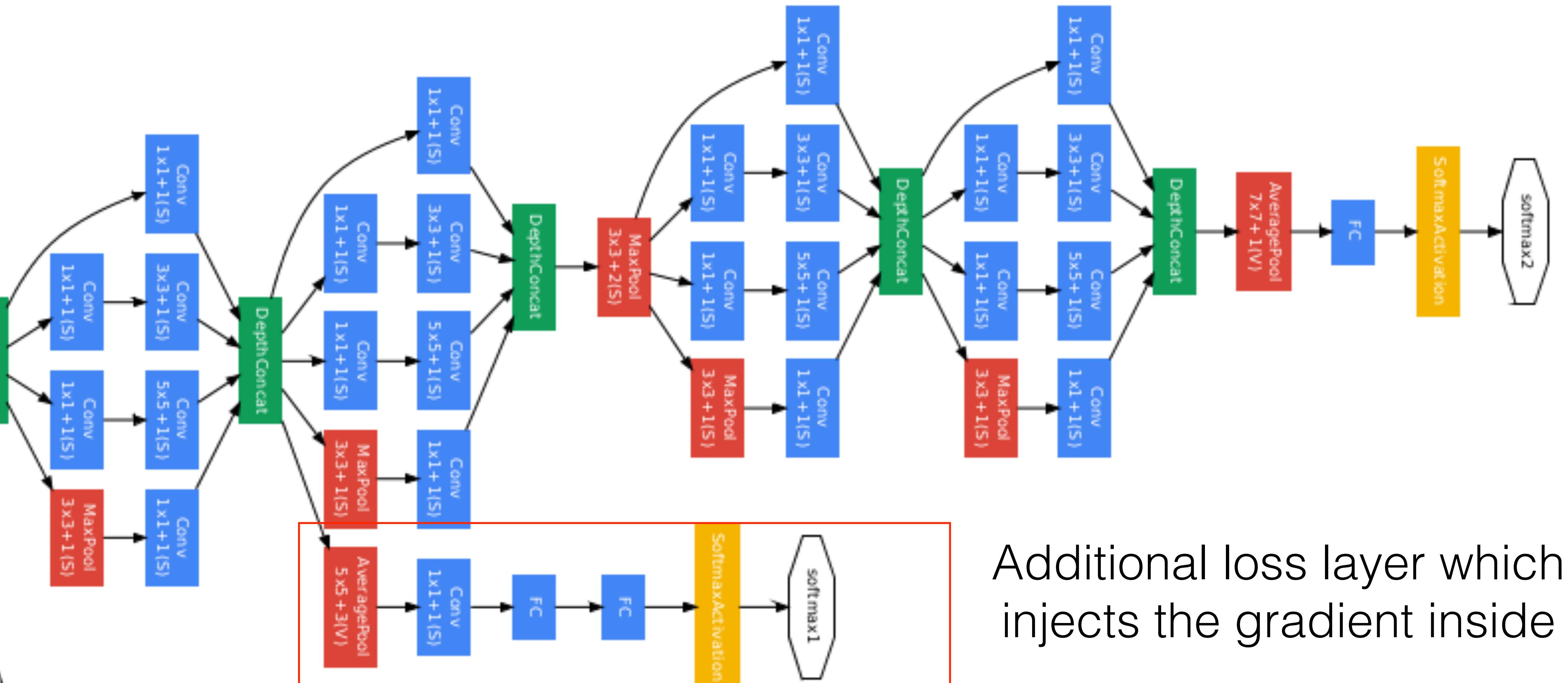
Szegedy et al. Going Deeper with Convolutions, CVPR, 2014
<https://arxiv.org/abs/1409.4842>

GoogLeNet



Szegedy et al. Going Deeper with Convolutions, CVPR, 2014
<https://arxiv.org/abs/1409.4842>

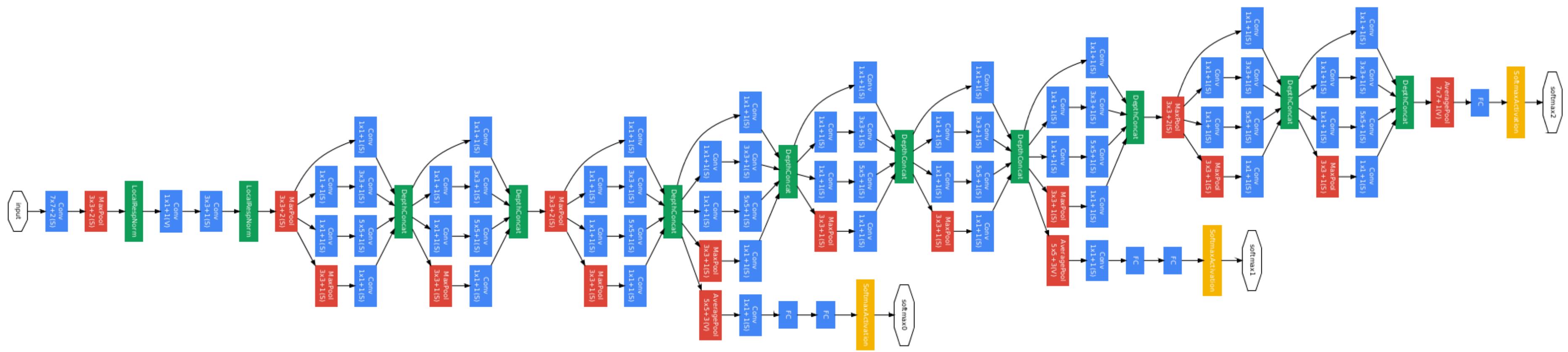
GoogLeNet



Additional loss layer which injects the gradient inside

Szegedy et al. Going Deeper with Convolutions, CVPR, 2014
<https://arxiv.org/abs/1409.4842>

GoogLeNet



- 12x fewer parameters than AlexNet
- depth 22 layers
- training: few high-end GPU about a week

Szegedy et al. Going Deeper with Convolutions, CVPR, 2014
<https://arxiv.org/abs/1409.4842>

Classification results

AlexNet

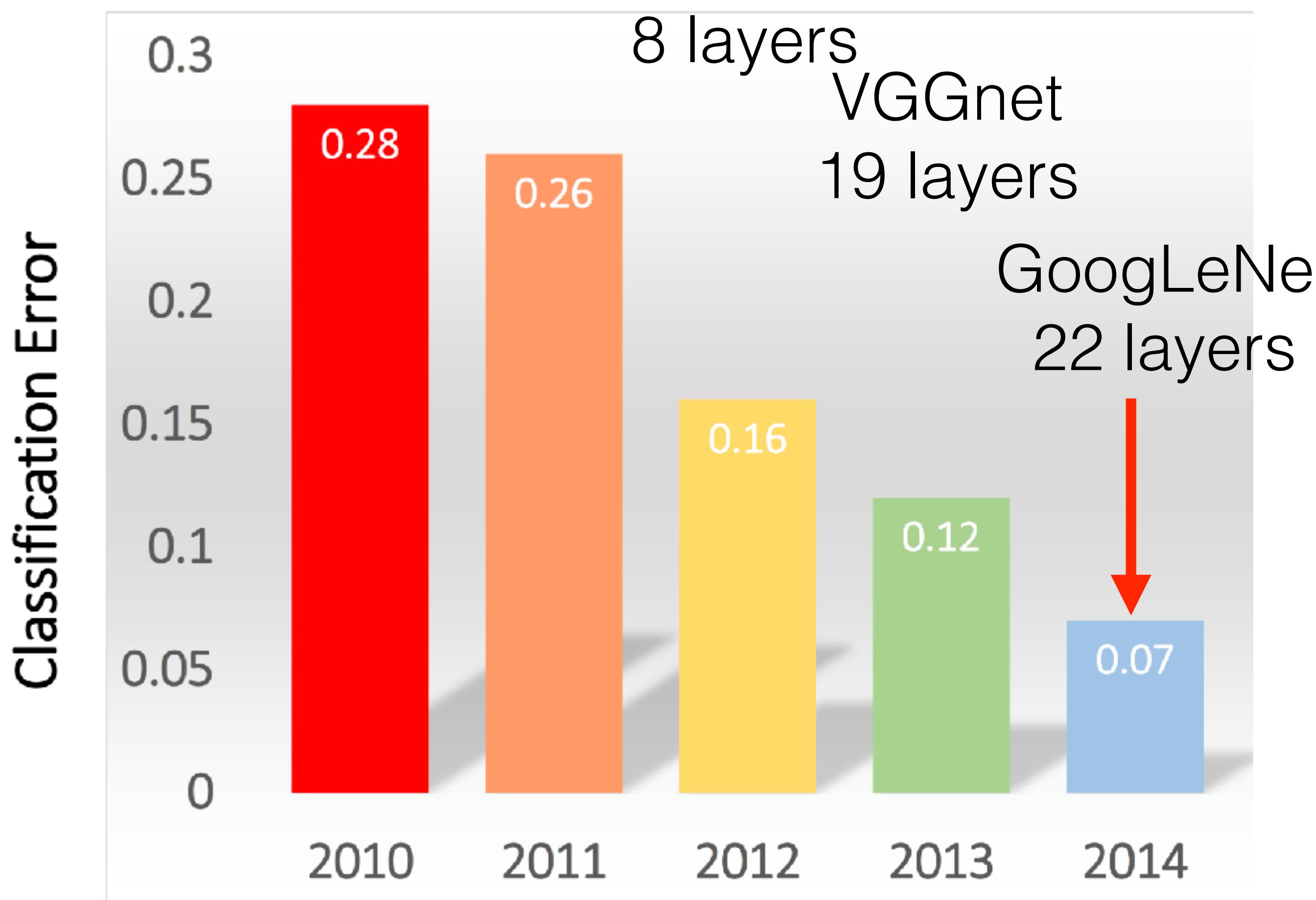
8 layers

VGGnet

19 layers

GoogLeNet

22 layers



Classification results

AlexNet

8 layers

VGGnet

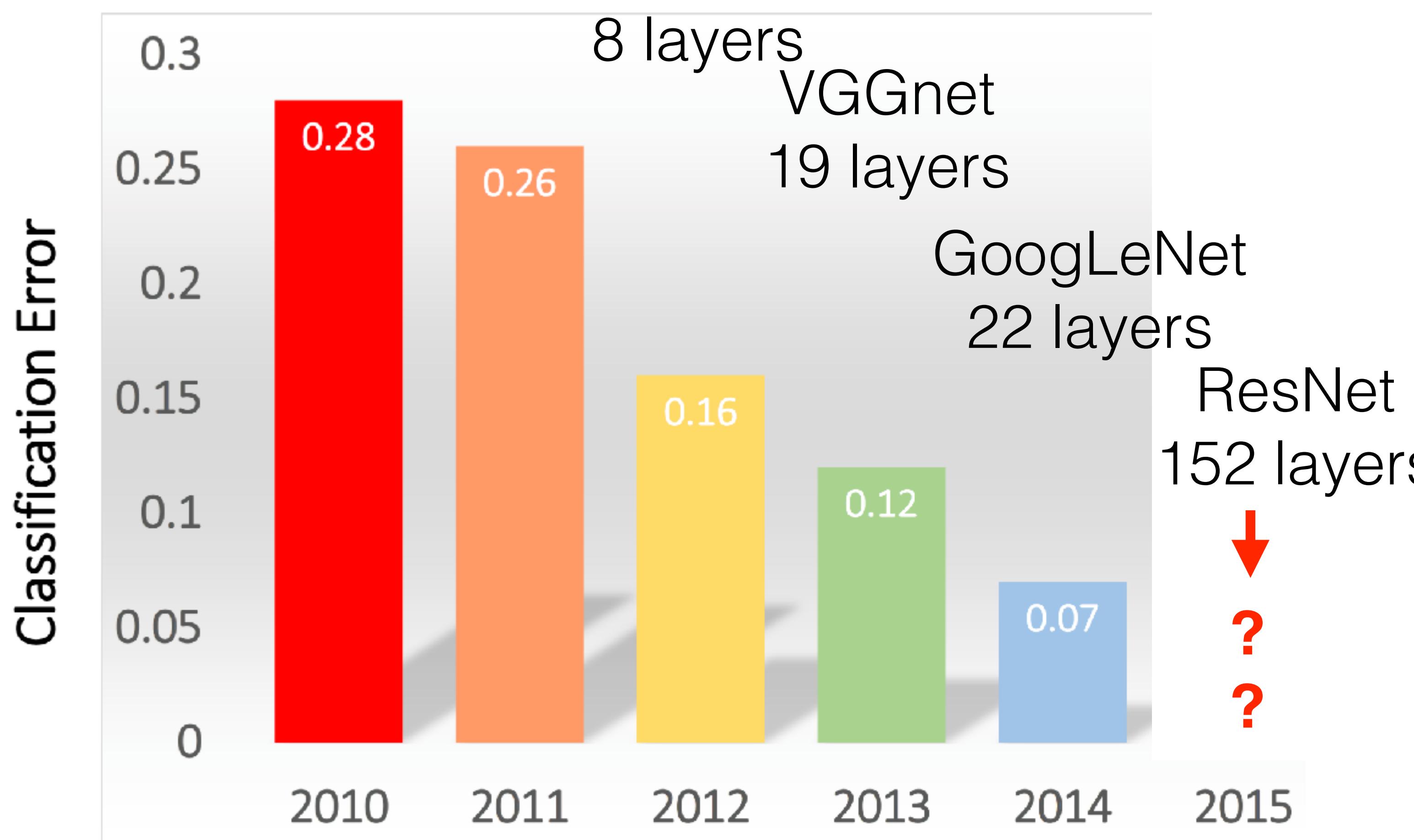
19 layers

GoogLeNet

22 layers

ResNet

152 layers



ResNet

Better results with smaller kernels + deeper architectures



Well said Leo, well said

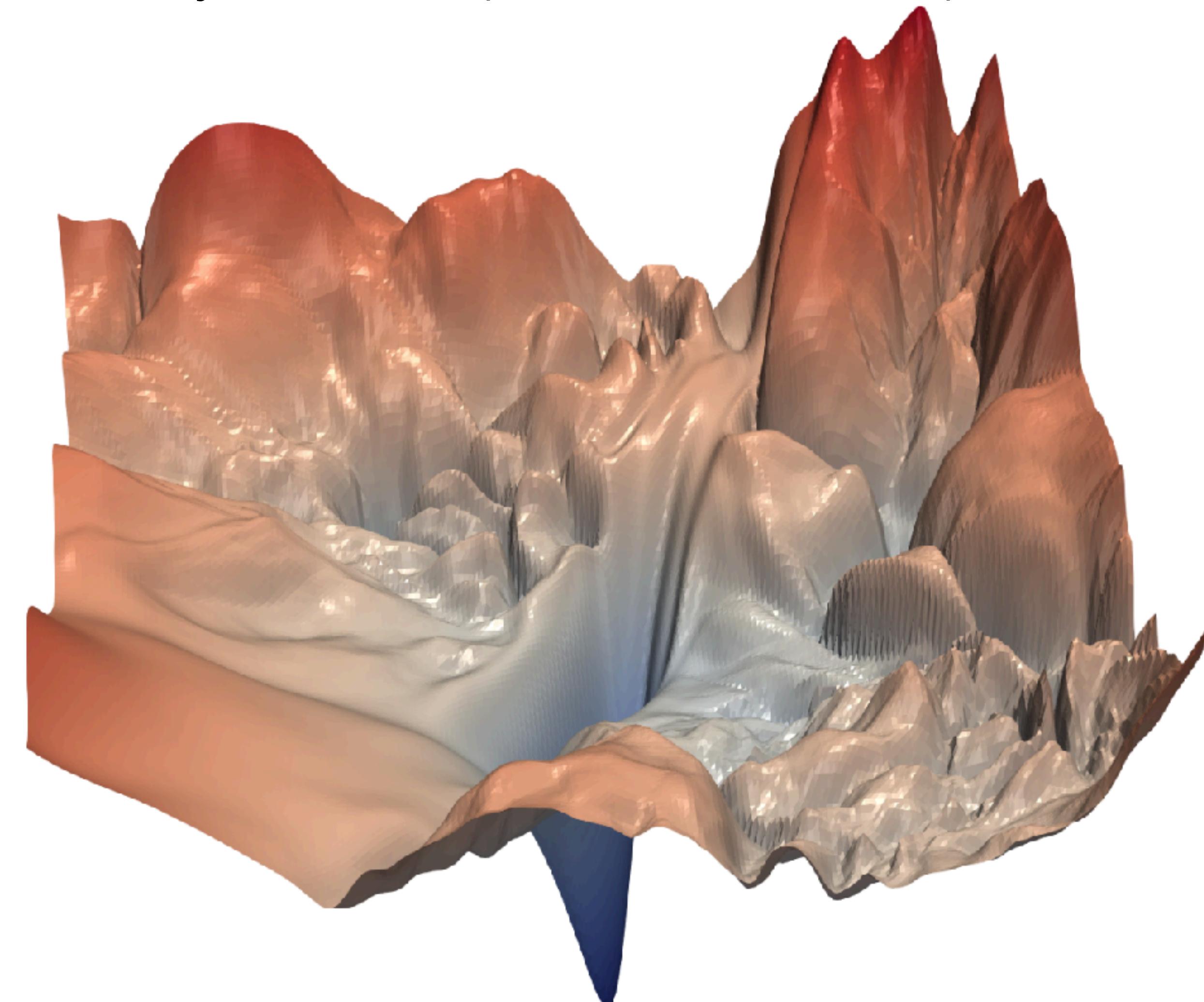
- deeper ConvNet architectures yielded **higher training errors**.=> **is it overfitting?**
- => no overfitting, but vanishing gradient !

He et al. Going Deeper with Convolutions, CVPR, 2015
<https://arxiv.org/abs/1512.03385>

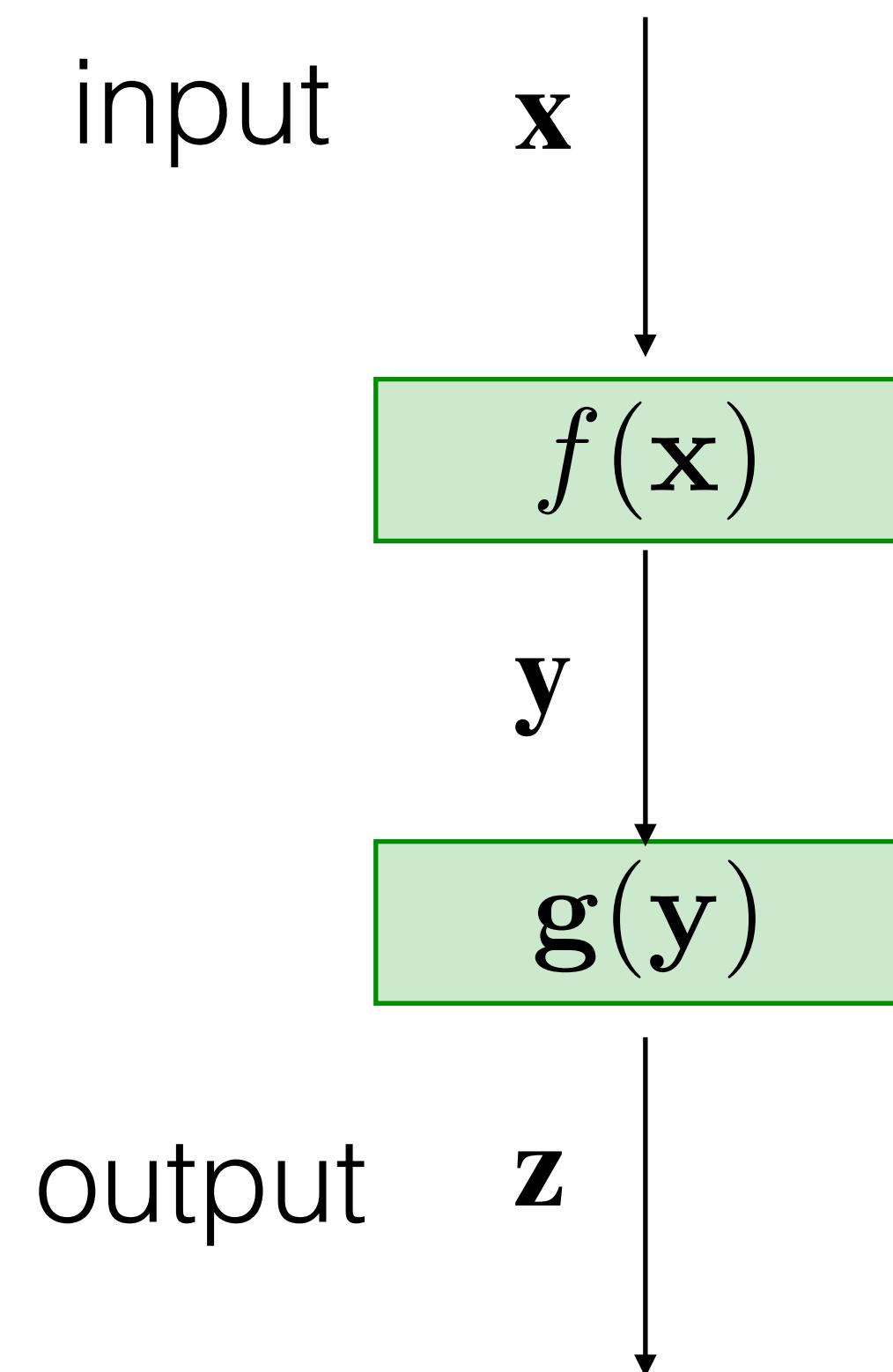
Visualizing Loss Landscape of Neural Nets

$$f(\alpha, \beta) = \mathcal{L}(\mathbf{w}^* + \alpha\mathbf{u} + \beta\mathbf{v})$$

for randomly chosen (and normalized) directions \mathbf{u}, \mathbf{v}



forward pass:

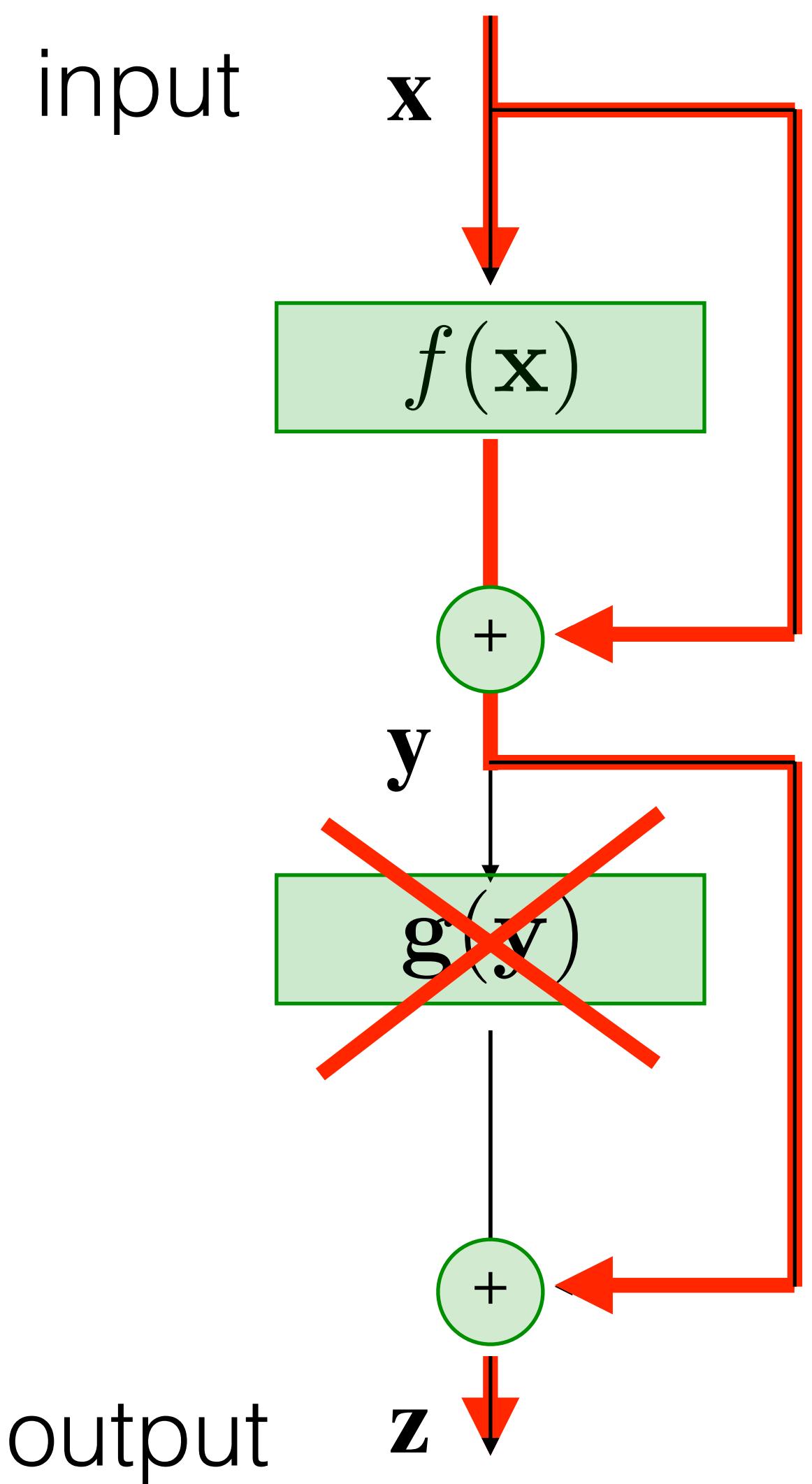


backward pass:

gradient $\frac{\partial \mathbf{z}}{\partial \mathbf{x}} = \begin{pmatrix} \frac{\partial \mathbf{z}}{\partial \mathbf{y}} & \frac{\partial \mathbf{y}}{\partial \mathbf{x}} \end{pmatrix} \approx \mathbf{0}$ then gradient is always zero
if any local gradient is zero

$$\approx \mathbf{0}$$

forward pass:



backward pass:

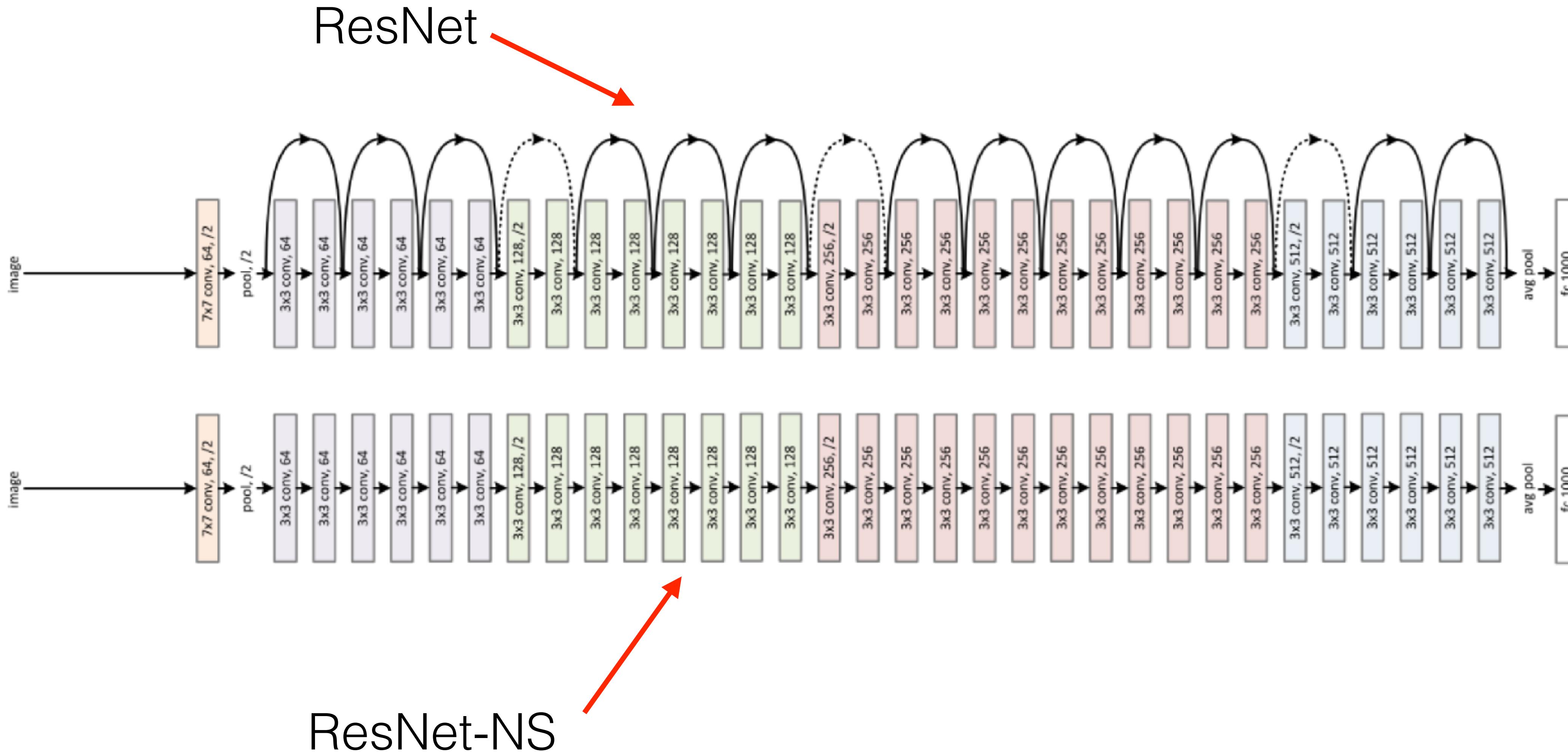
gradient $\frac{\partial z}{\partial x} = \left(\frac{\partial z}{\partial y} + 1 \right) \left(\frac{\partial y}{\partial x} + 1 \right) \neq 0$

A diagram showing the backward pass gradient flow. It consists of two nodes connected by a red arrow pointing from right to left. The left node contains the expression $\frac{\partial z}{\partial y} + 1$ and is labeled ≈ 0 below it. The right node contains the expression $\frac{\partial y}{\partial x} + 1$ and is also labeled ≈ 0 below it. Red arrows point from each node towards the center, indicating the flow of the gradient.

if any local gradient is zero

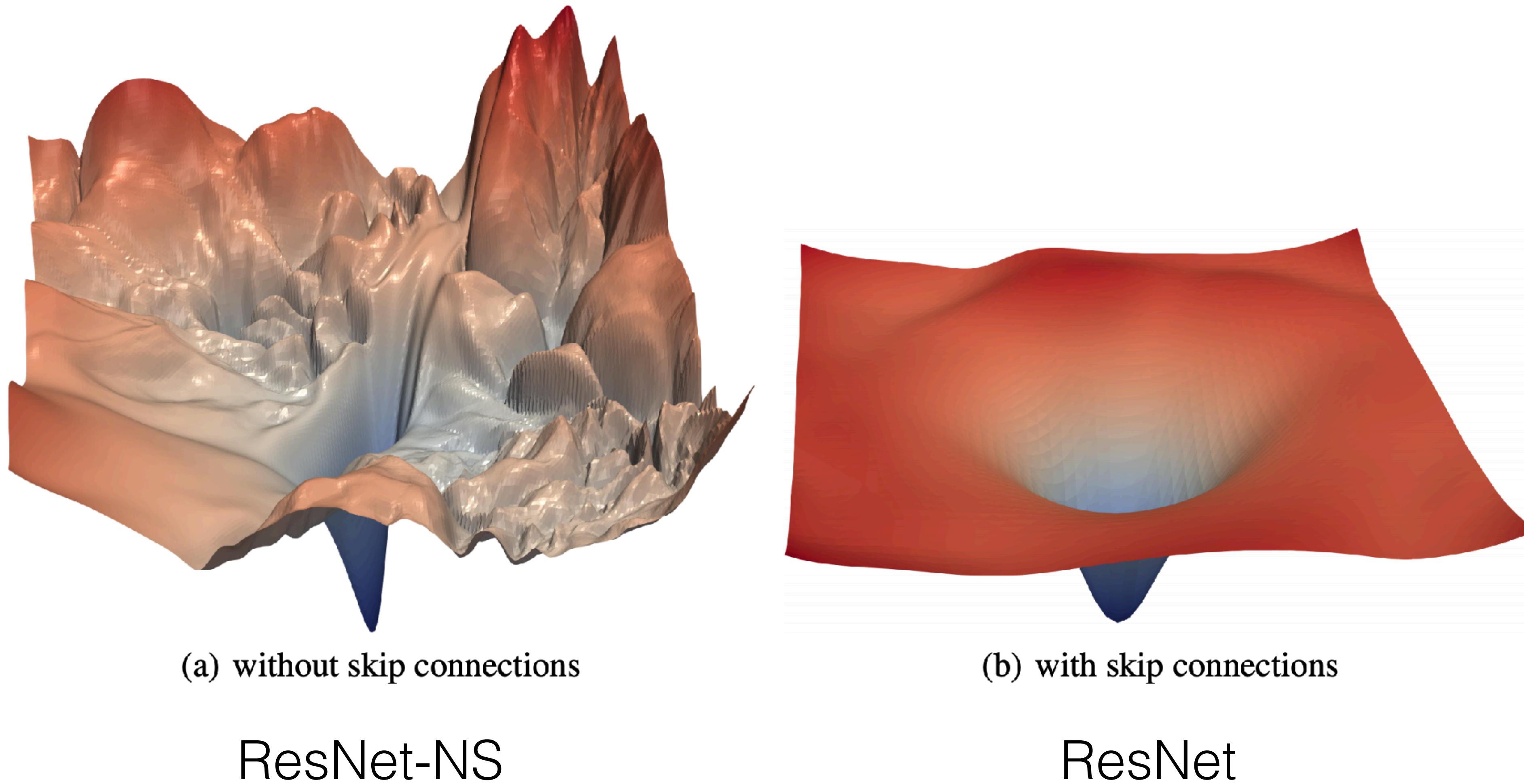
the gradient can still flow through another path

ResNet: deep ConvNet with skip connections



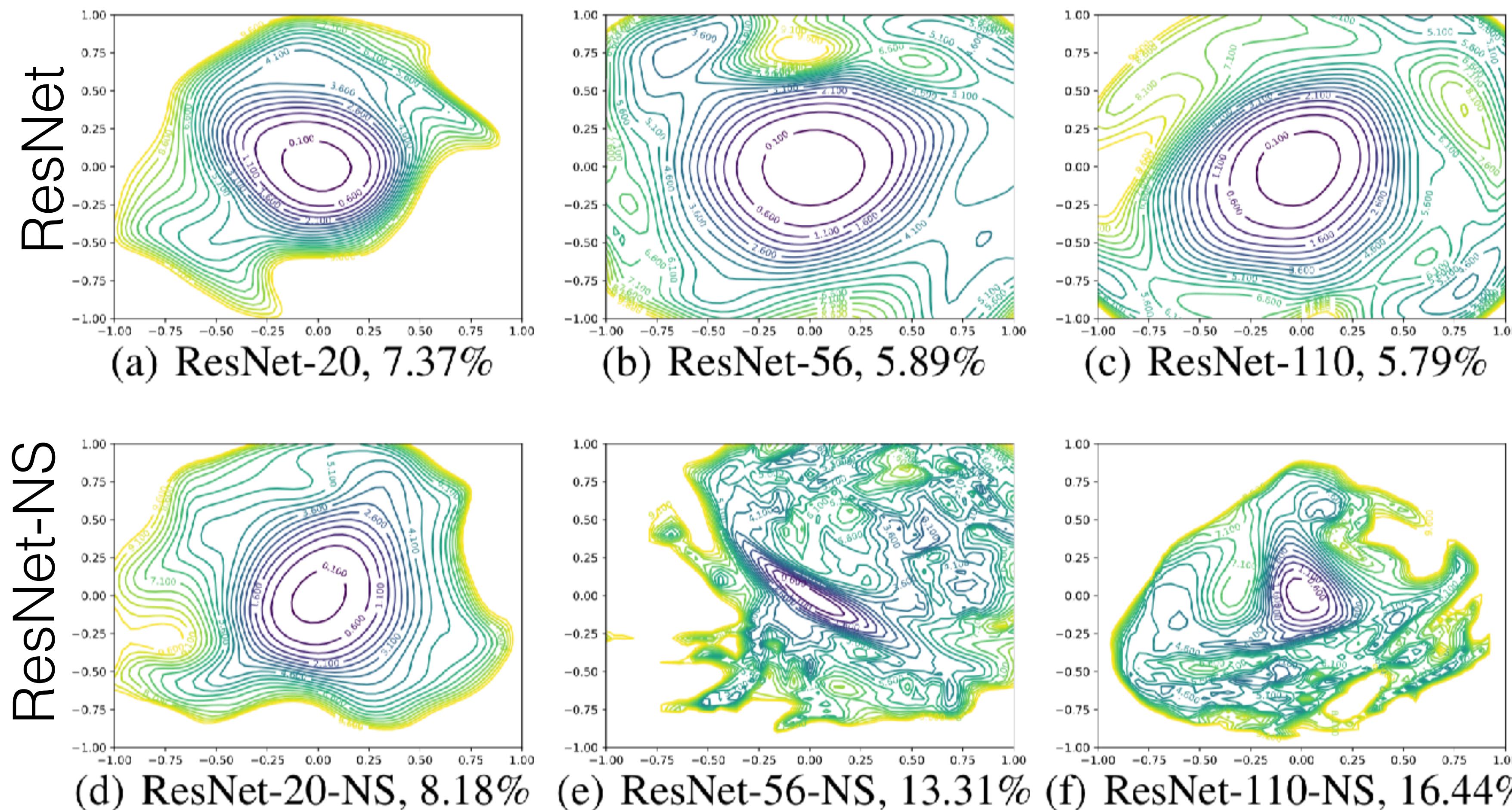
Visualizing Loss Landscape of Neural Nets

[Li et al, NIPS, 2018] <https://arxiv.org/pdf/1712.09913.pdf>

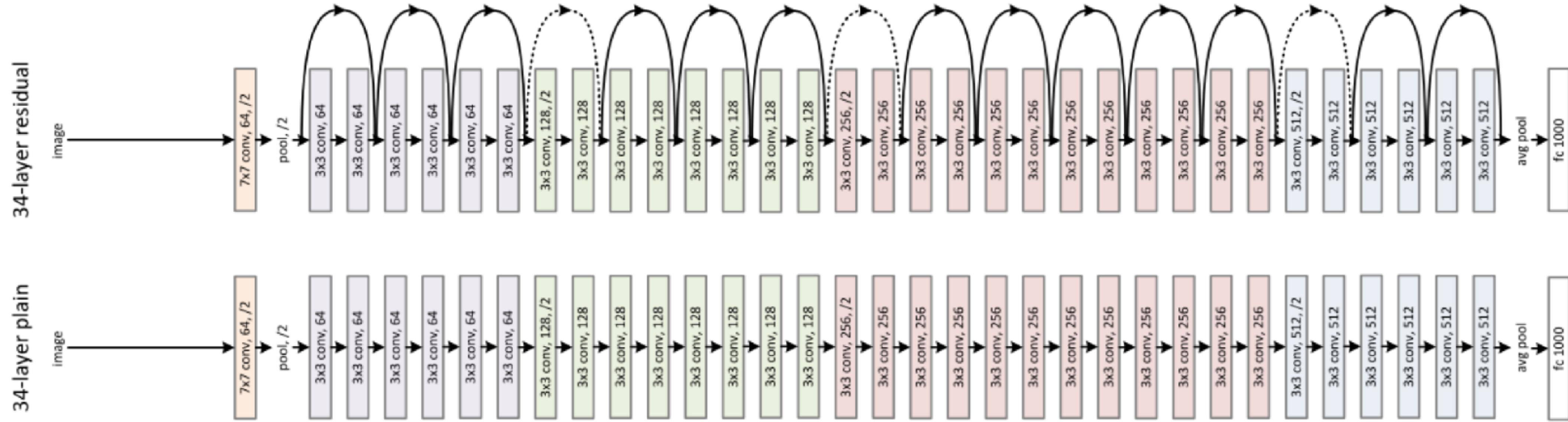


ResNet: deep ConvNet with skip connections

[Li et al, NIPS, 2018] <https://arxiv.org/pdf/1712.09913.pdf>

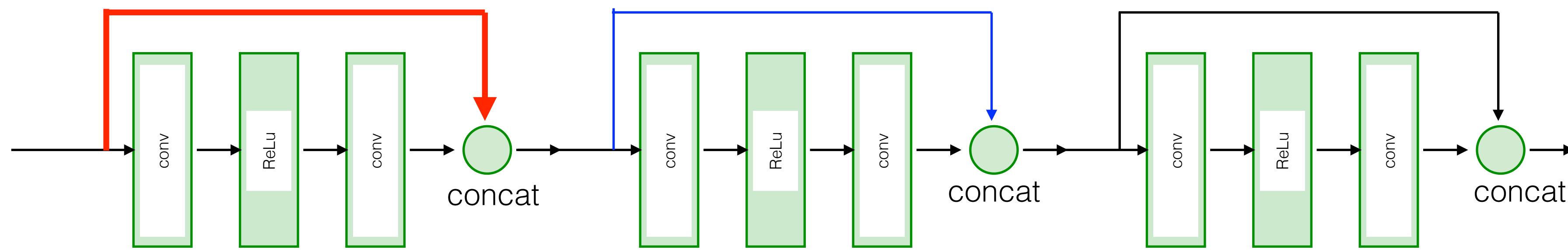


ResNet: deep ConvNet with skip connections



- **1k+ layers** possible
- **Initialization** with **zero** weights is meaningful
- **Better gradient flow:** many independent paths, opt-friendly landscape
- **Robustness** wrt **noise** and layer removal (if some important feature is not detected in a particular layer it can be detected in following layers)

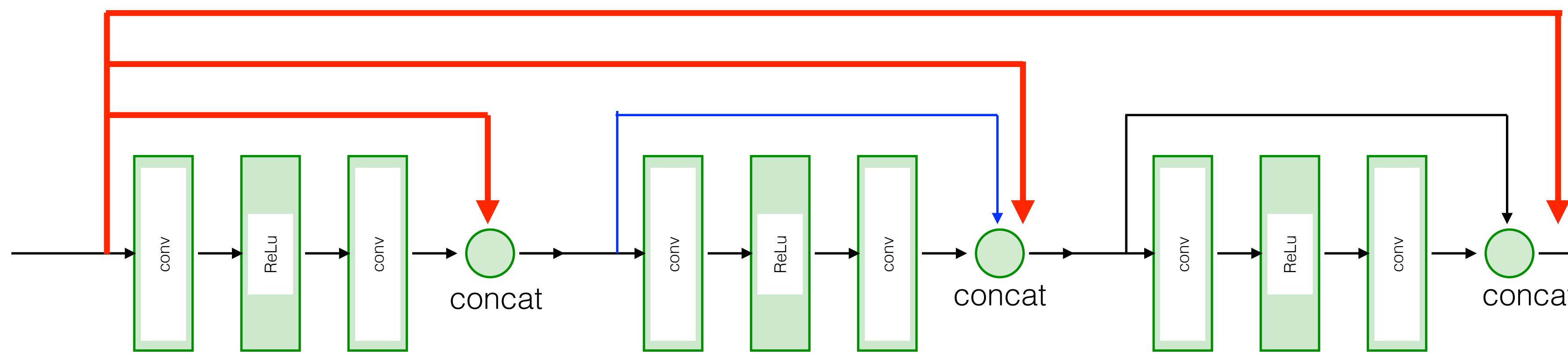
ResNet=>DenseNet



- Directly propagate each feature map to all following layers

Huang, Densely Connected Convolutional Networks, CVPR 2017. <https://arxiv.org/abs/1608.06993>

DenseNet



- Directly propagate each feature map to all following layers

Classification results

AlexNet

8 layers

VGGnet

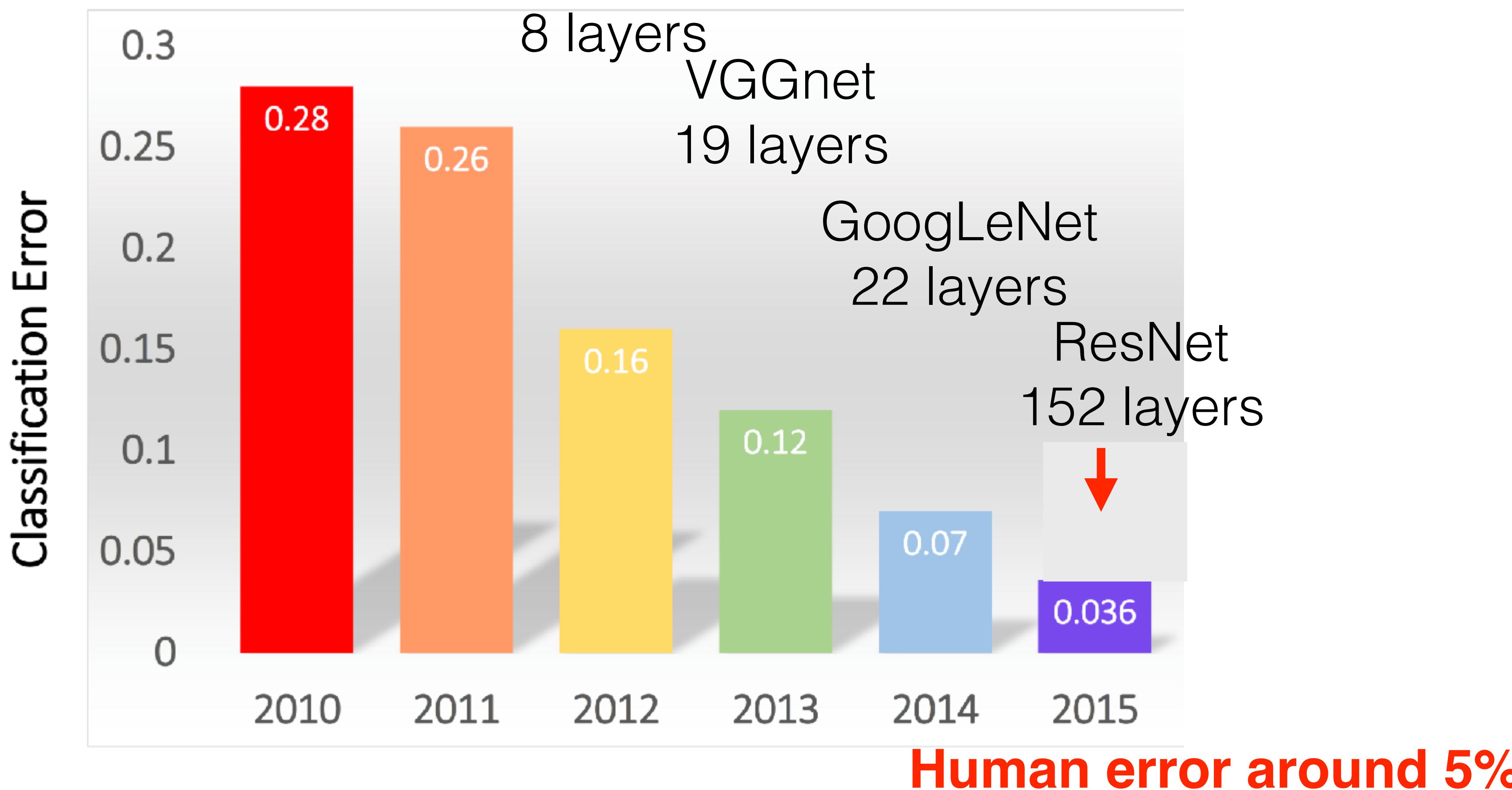
19 layers

GoogLeNet

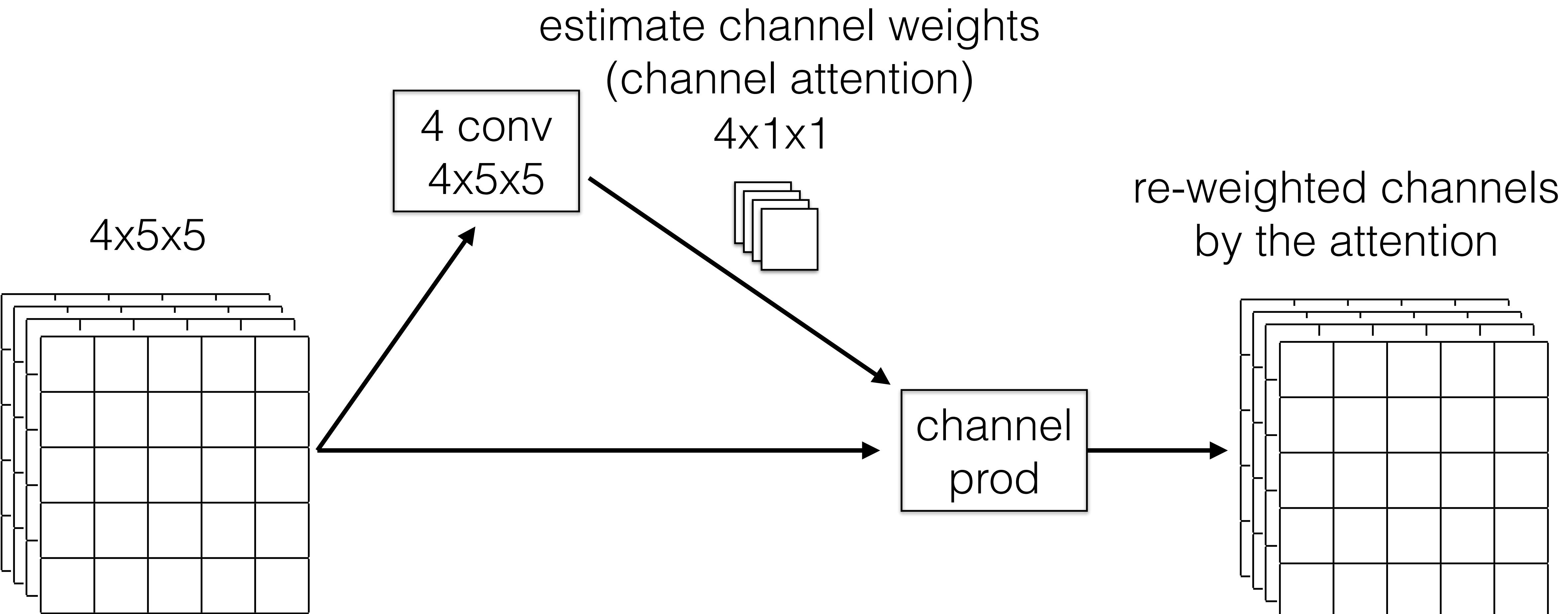
22 layers

ResNet

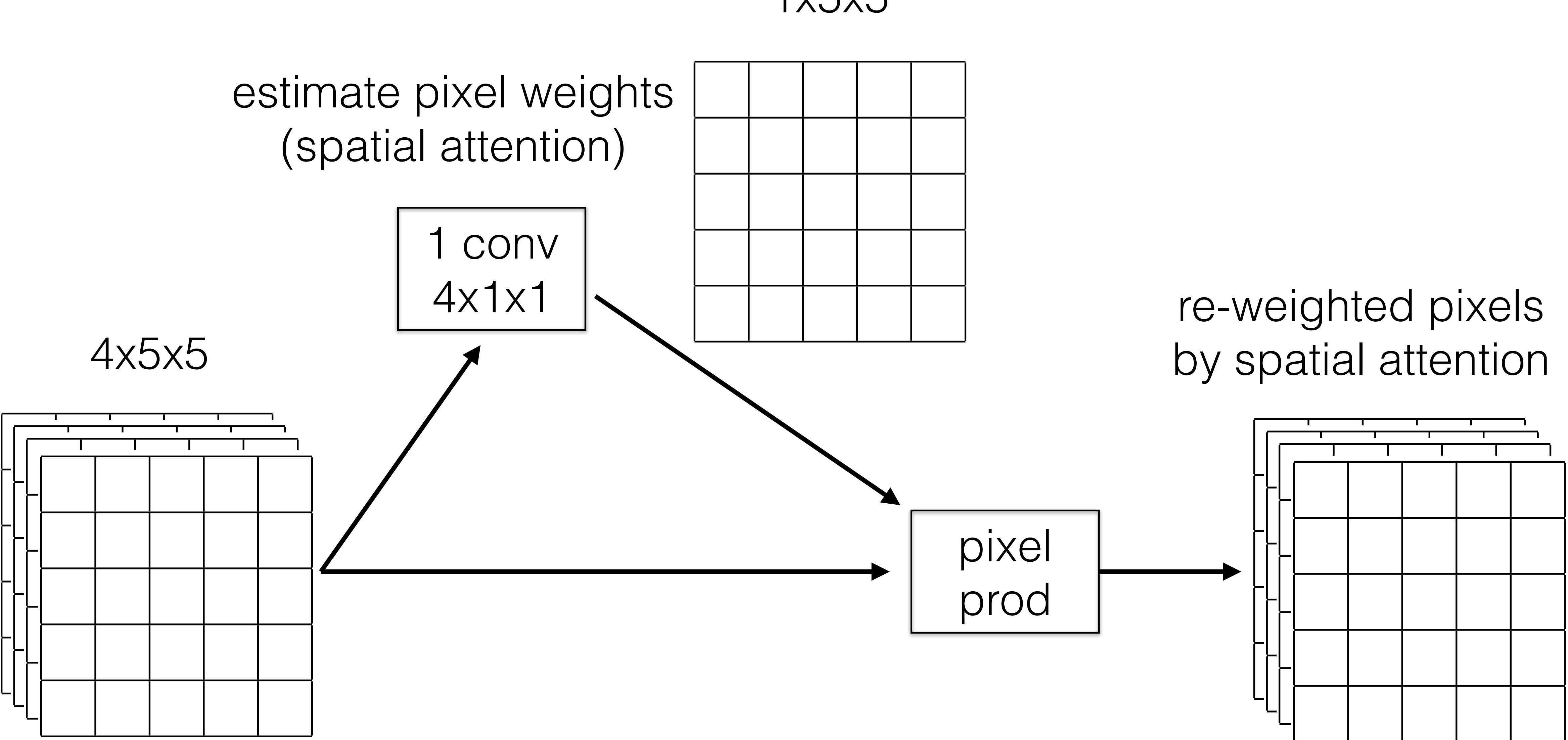
152 layers



Channel attention

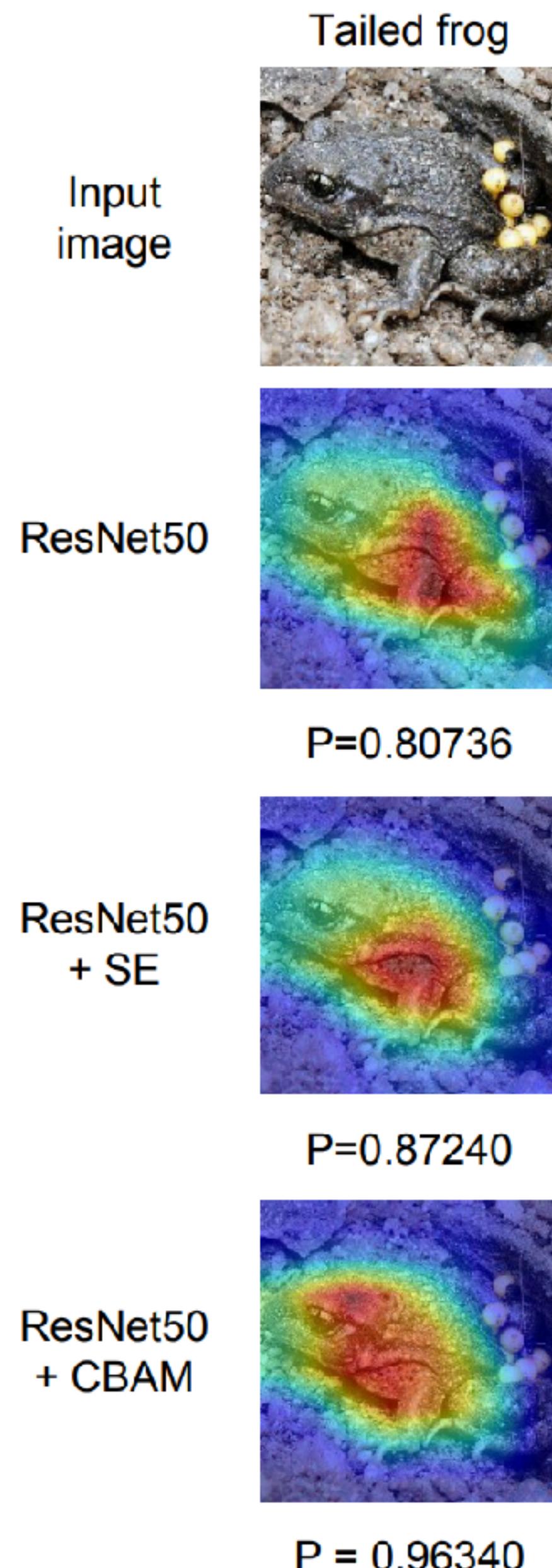


Spatial attention



Attention modules [Woo et al, ECCV, 2018]

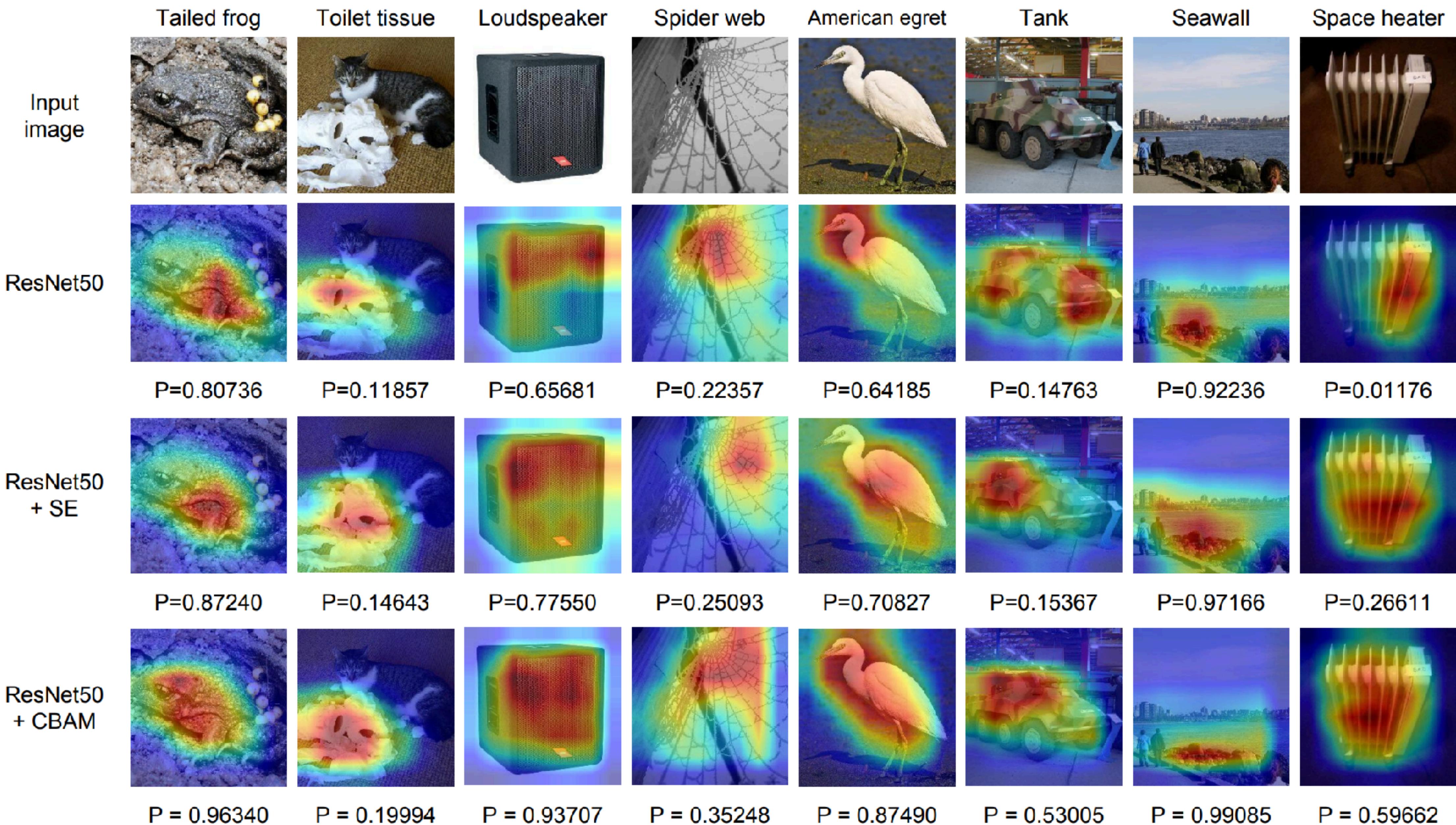
<https://arxiv.org/pdf/1807.06521v2.pdf>



The attention map of classified object is better localized

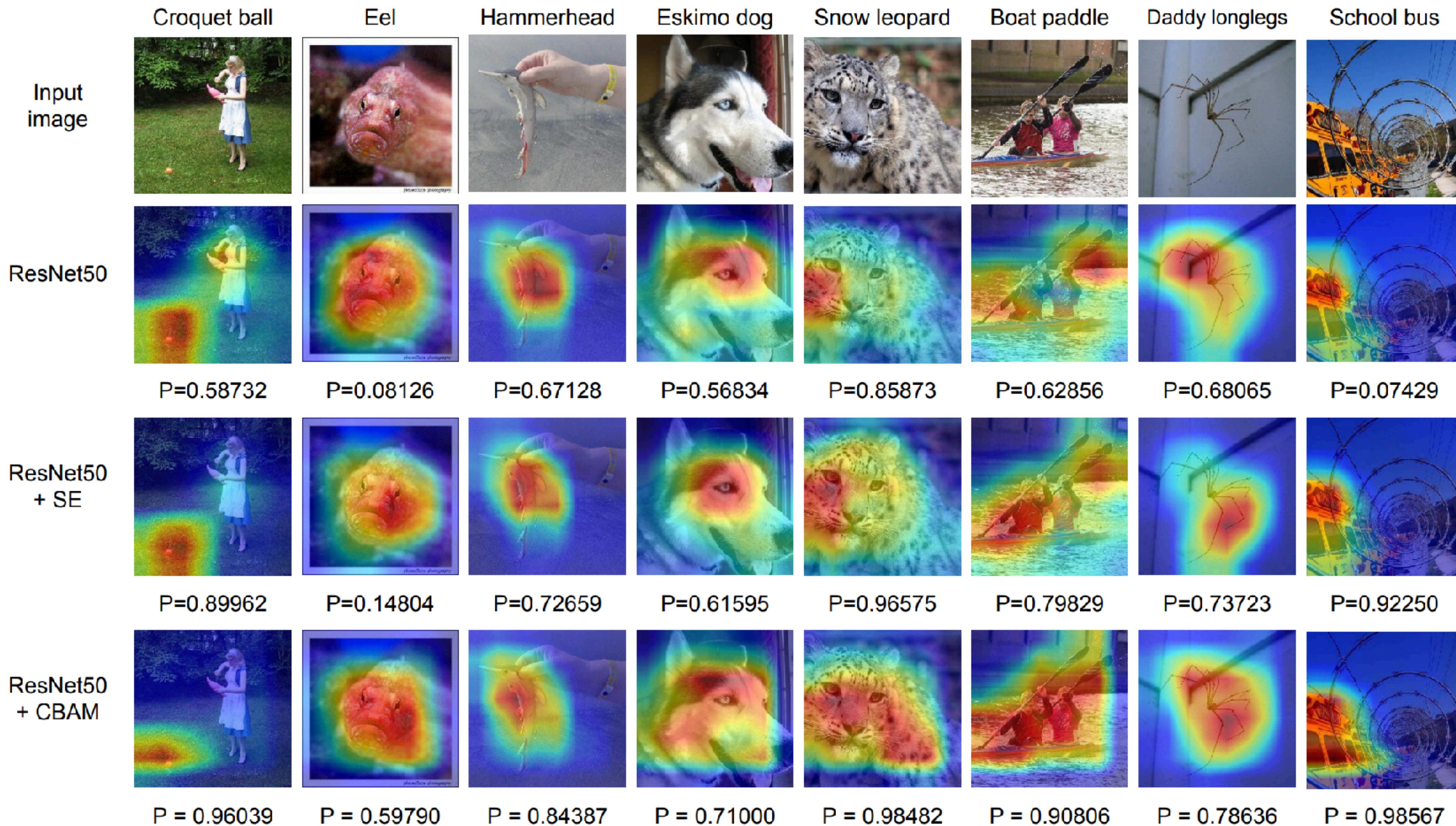
Attention modules [Woo et al, ECCV, 2018]

<https://arxiv.org/pdf/1807.06521v2.pdf>



Attention modules [Woo et al, ECCV, 2018]

<https://arxiv.org/pdf/1807.06521v2.pdf>



GradCAM [Selvaraju et al,ICCV, 2018]



Classification results

AlexNet

8 layers

VGNet

19 layers

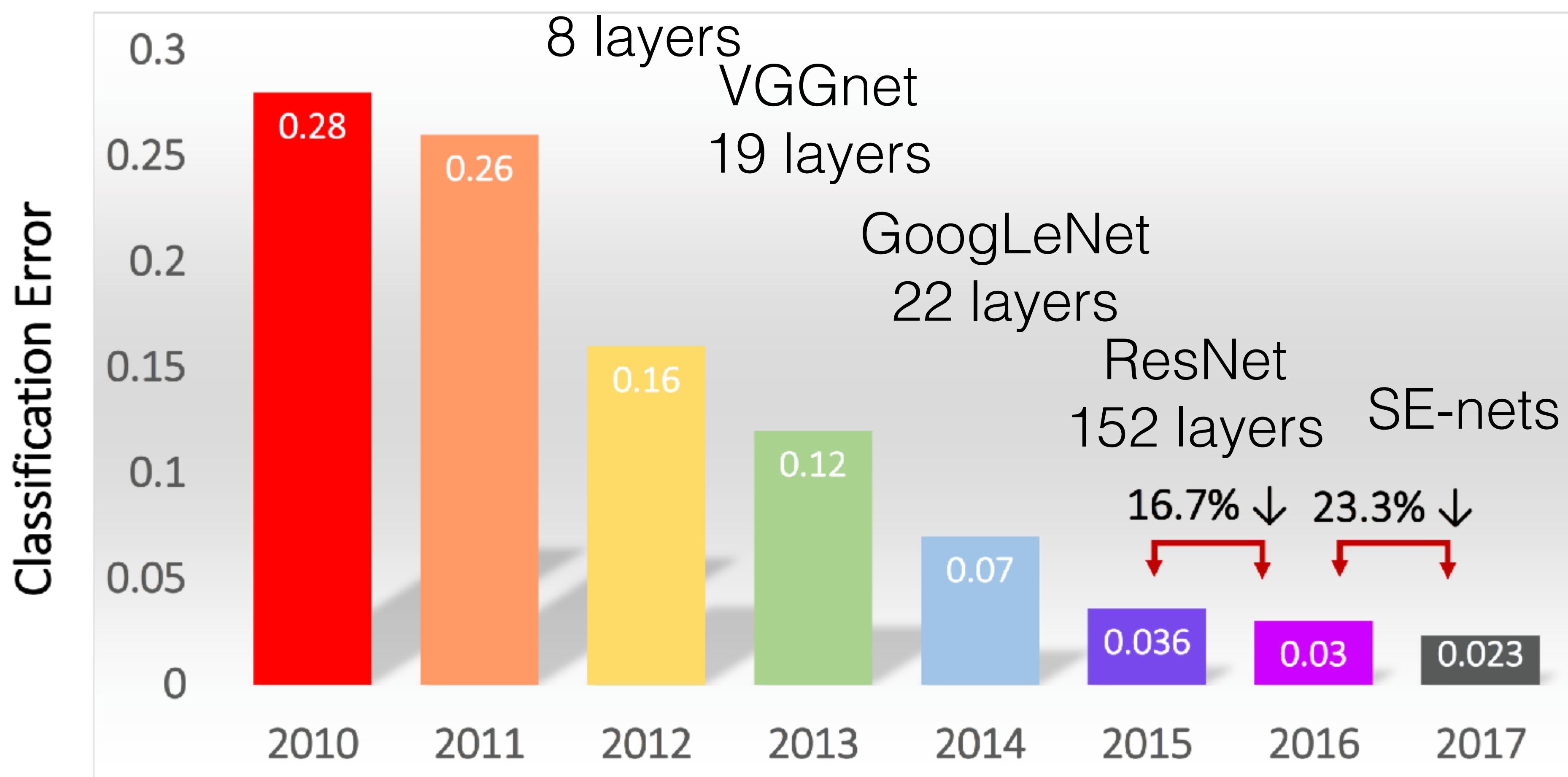
GoogLeNet

22 layers

ResNet

152 layers SE-nets

16.7% ↓ 23.3% ↓



Summary classification

- Injecting gradient
- Skip connections
- Receptive field (“one 7x7” vs “three 3x3 convs”)
- Spatial and Channel Attention

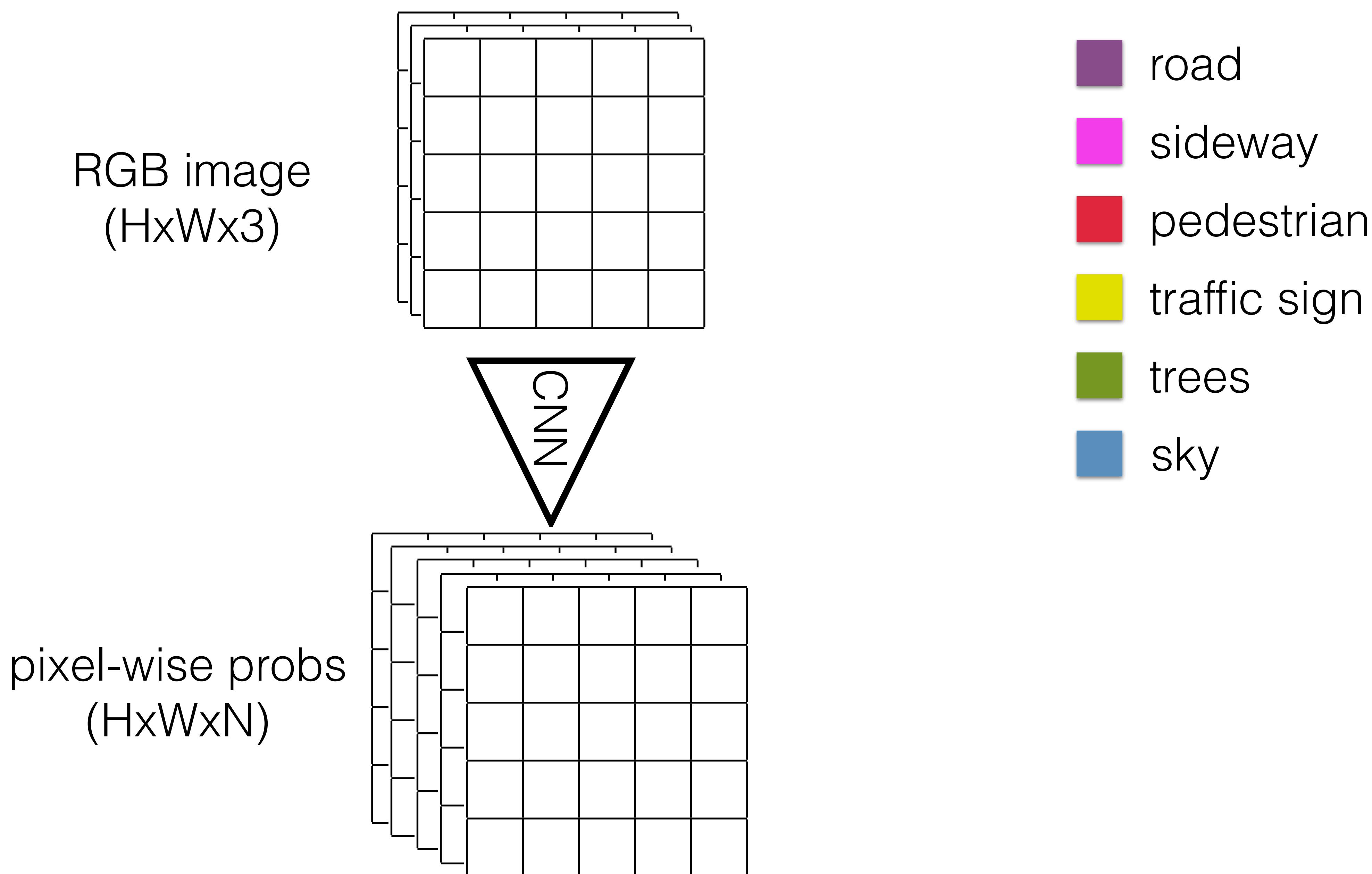
Outline

- Architectures of classification networks
- Architectures of segmentation networks
- Architectures of regression networks
- Architectures of detection networks
- Architectures of regression networks
- Architectures of feature matching networks

Semantic segmentation

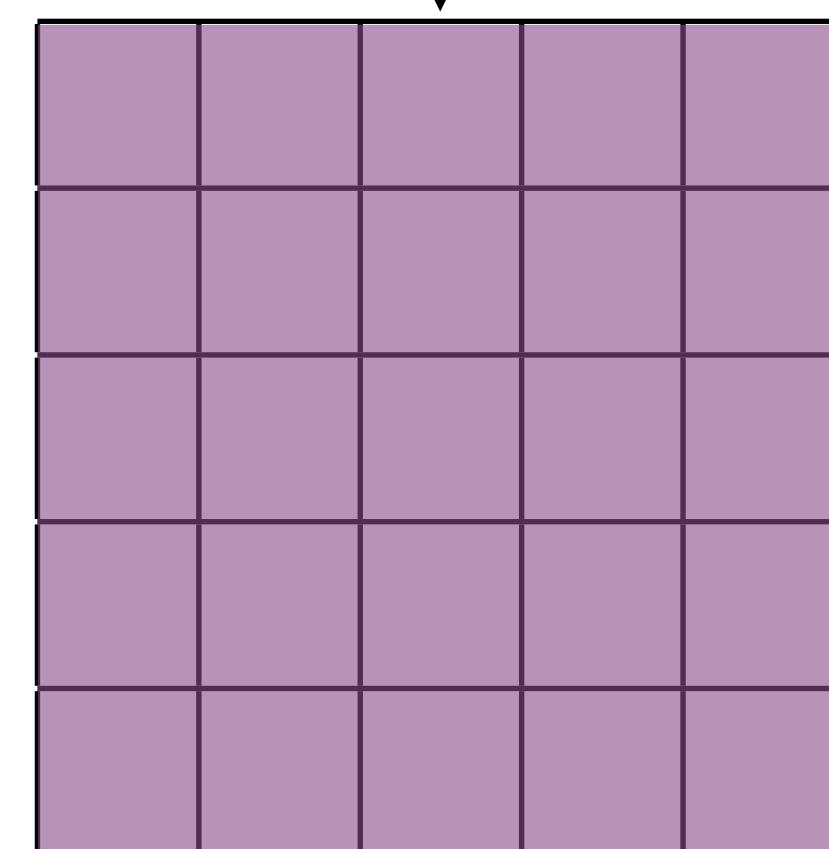
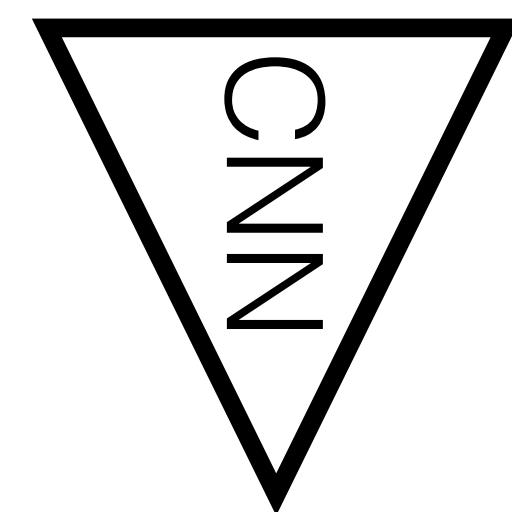
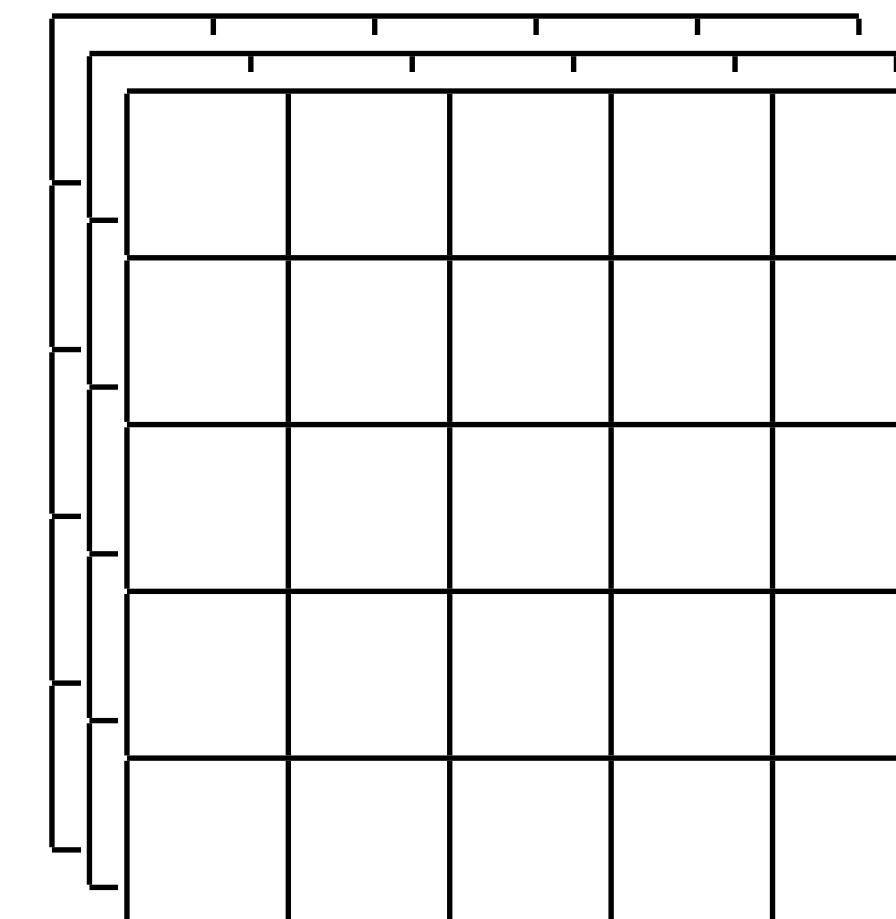


Semantic segmentation



Semantic segmentation

RGB image
(HxWx3)



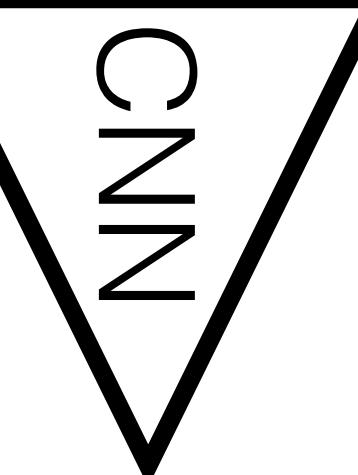
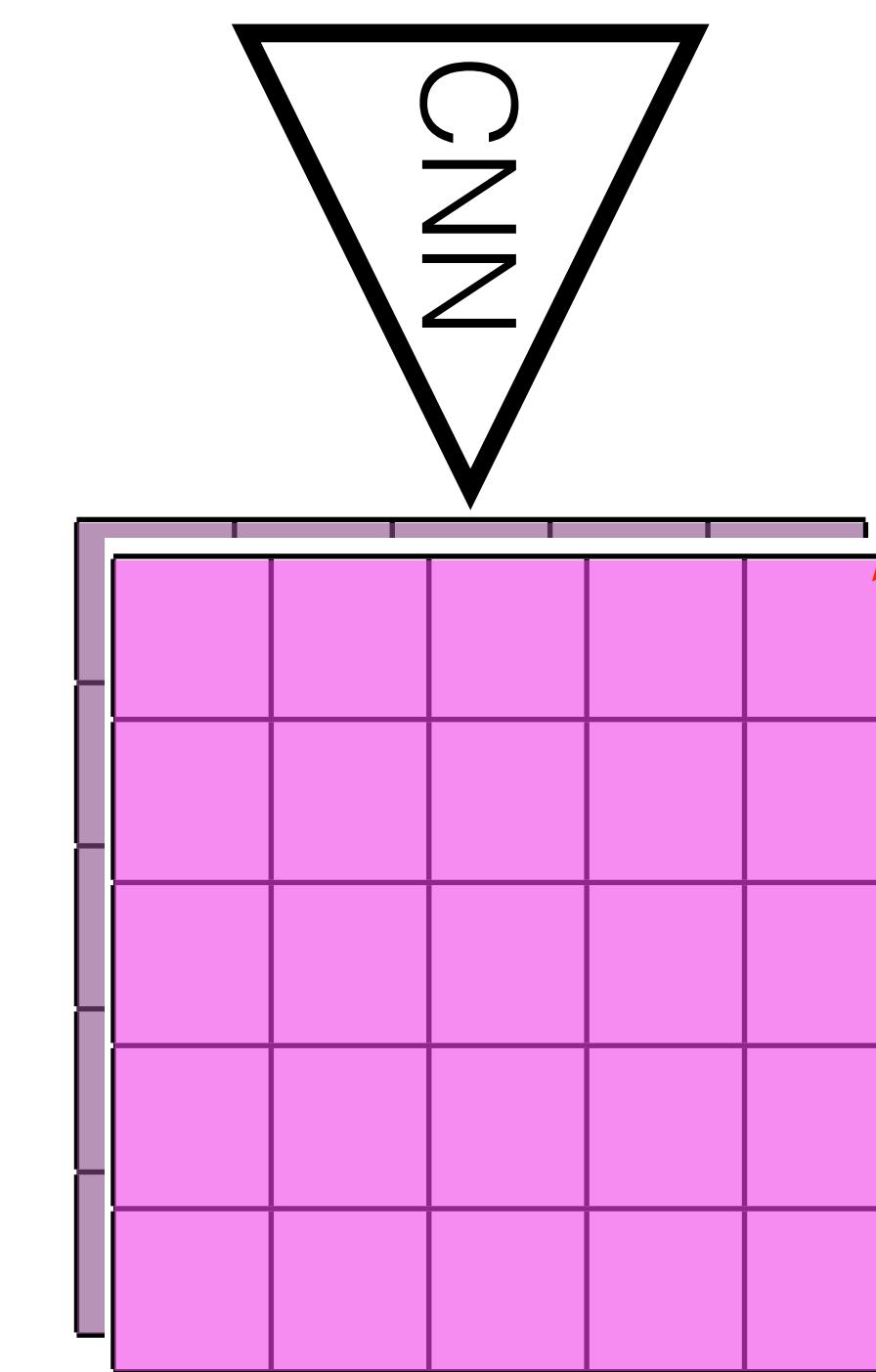
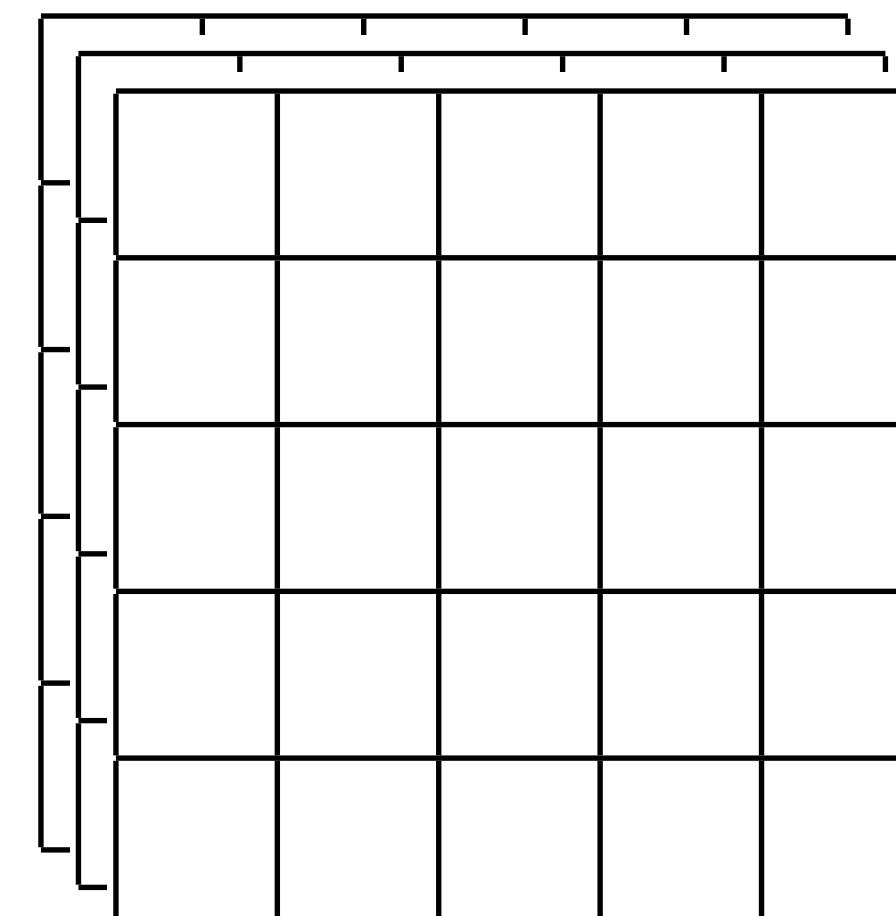
- █ road
- █ sideway
- █ pedestrian
- █ traffic sign
- █ trees
- █ sky

pixel-wise probability
of being **road**

channel 1

Semantic segmentation

RGB image
(HxWx3)



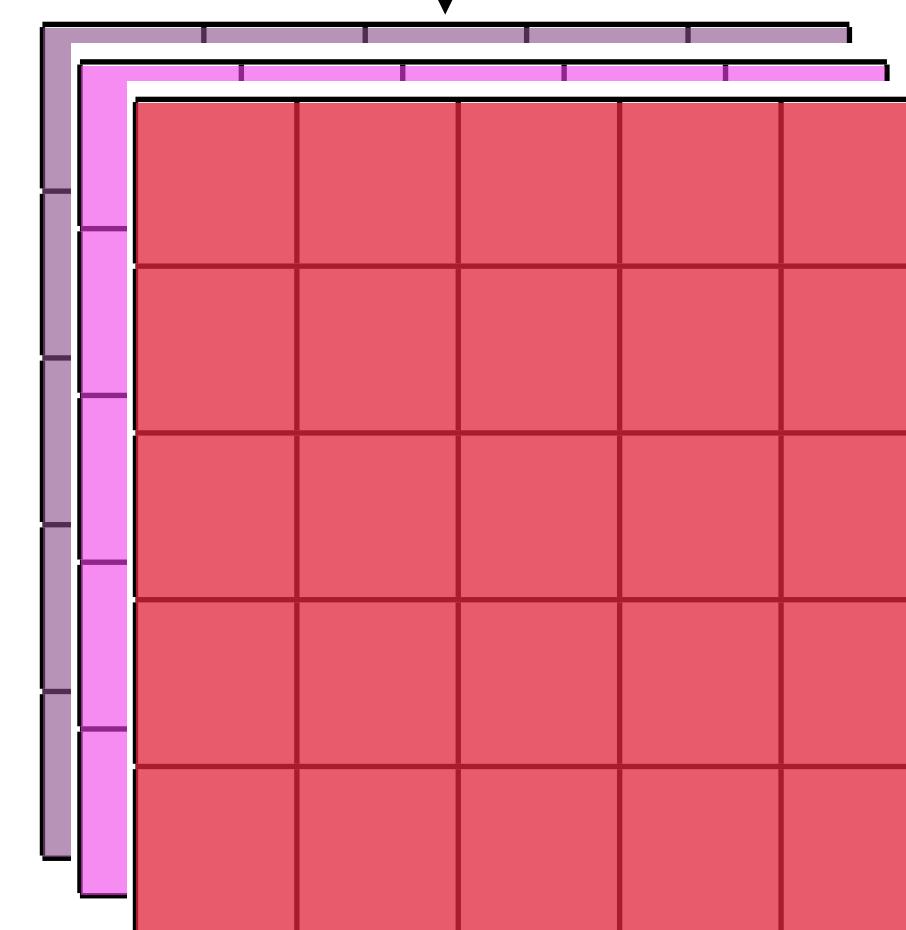
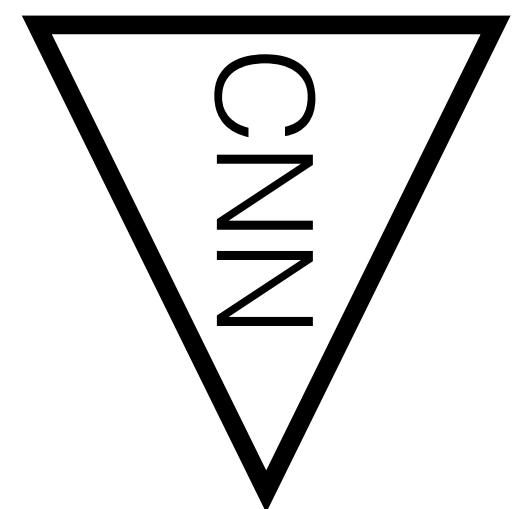
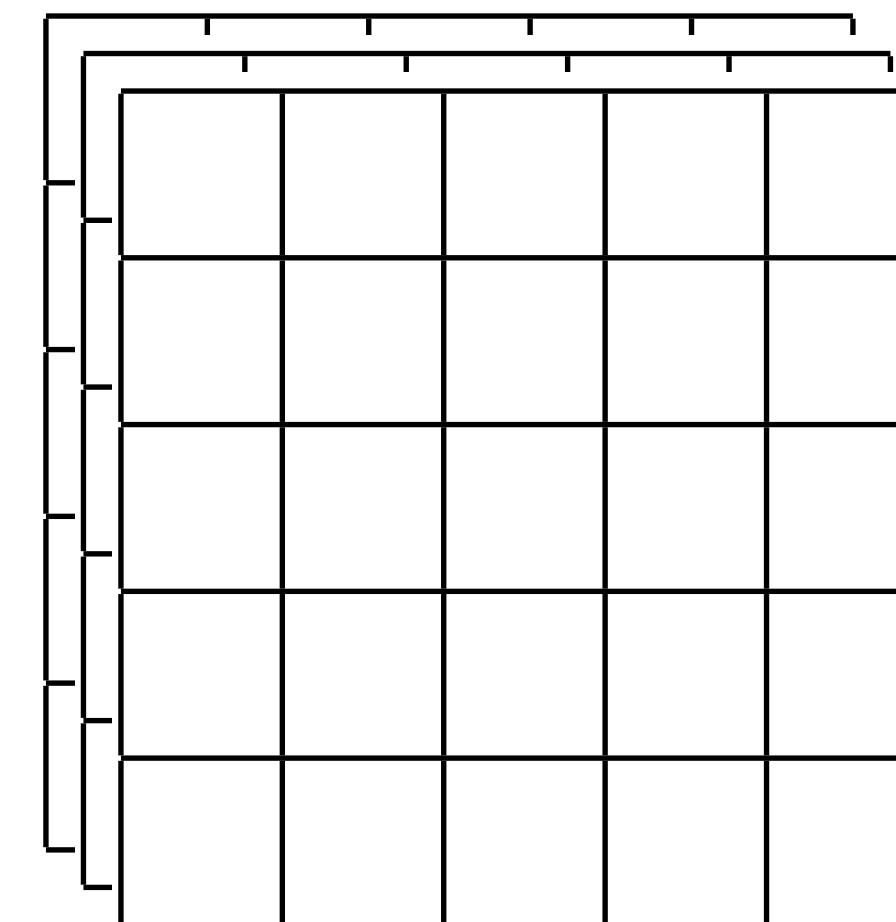
pixel-wise probability
of being **sideway**

channel 2

- █ road
- █ sideway
- █ pedestrian
- █ traffic sign
- █ trees
- █ sky

Semantic segmentation

RGB image
(HxWx3)

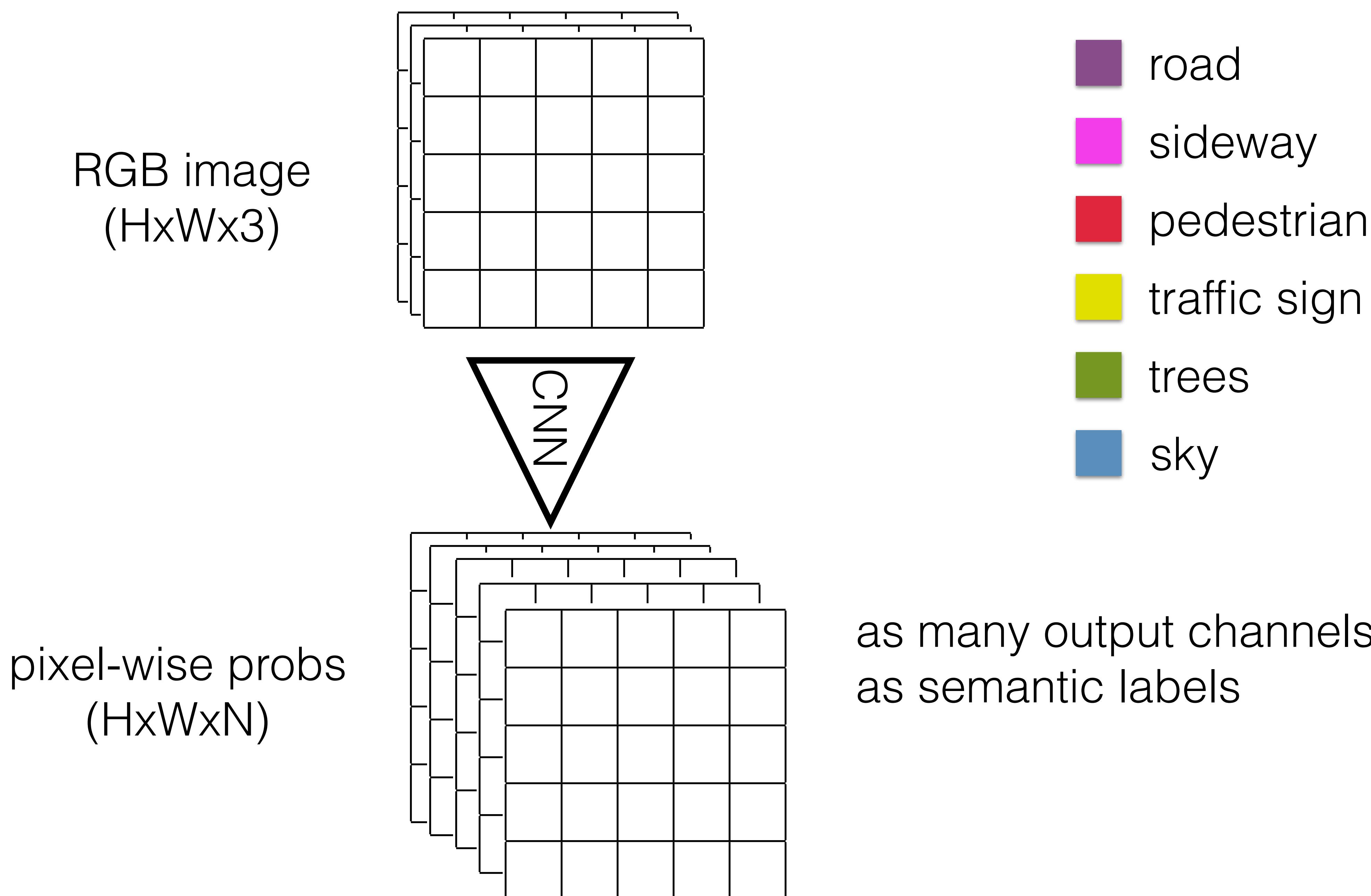


pixel-wise probability
of being **pedestrian**

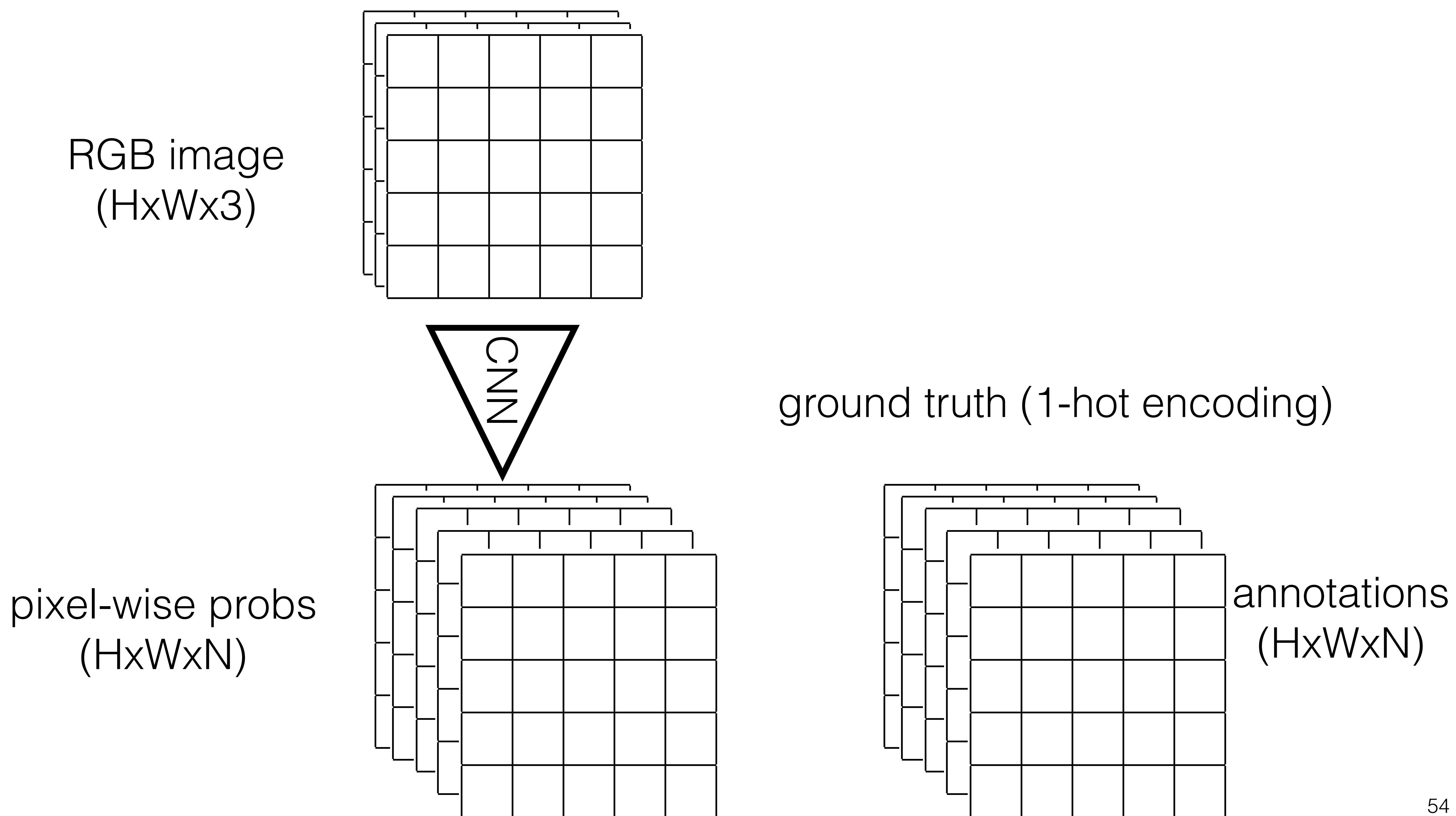
channel 3

- █ road
- █ sideway
- █ pedestrian
- █ traffic sign
- █ trees
- █ sky

Semantic segmentation

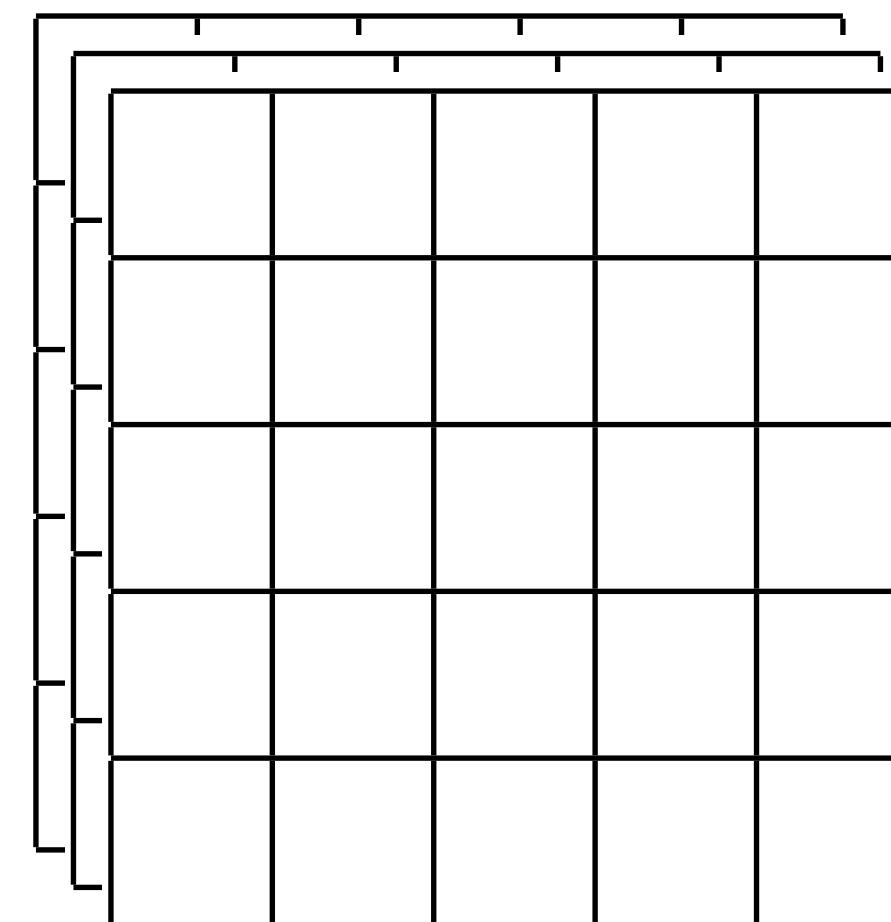


Semantic segmentation

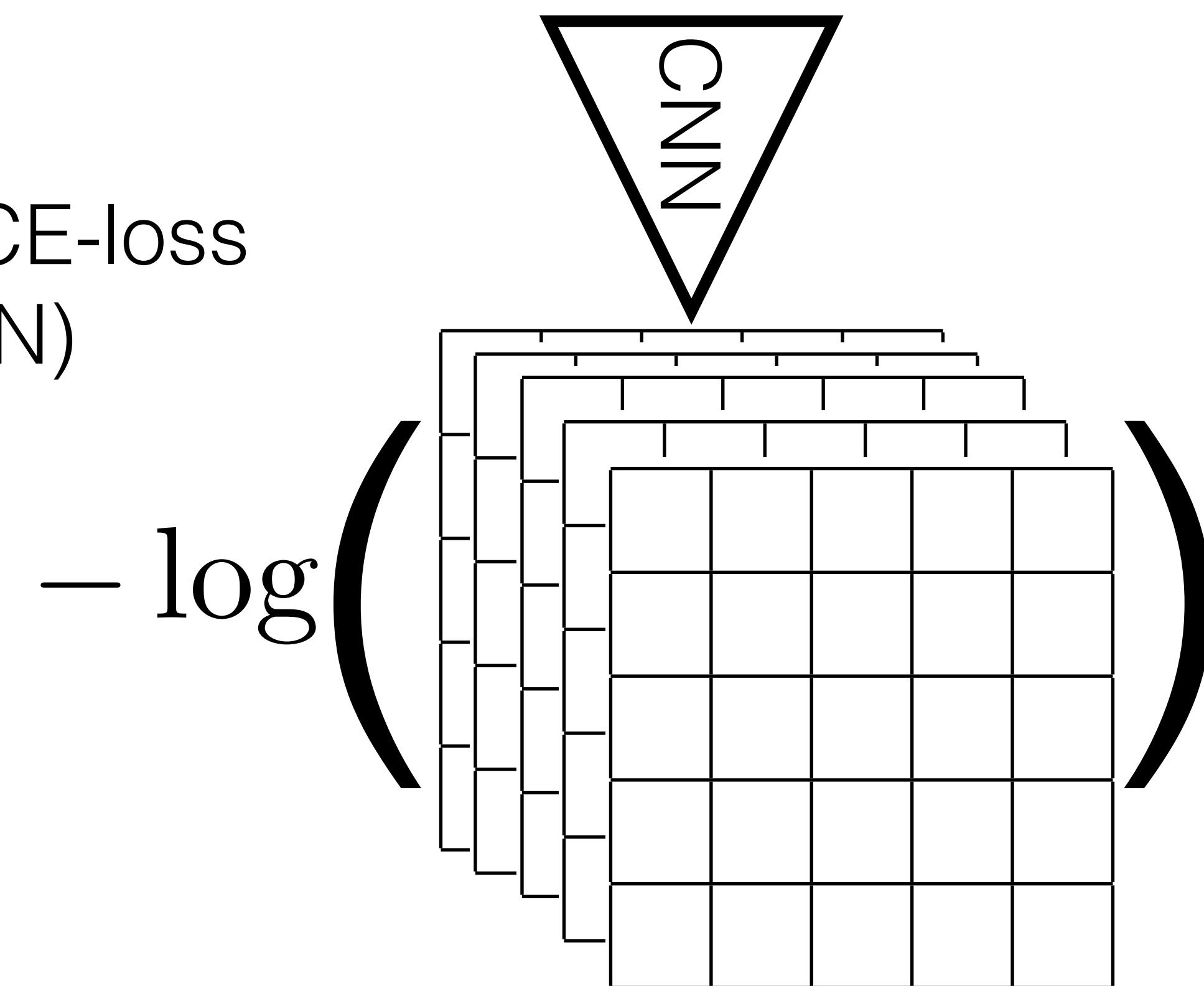


Semantic segmentation

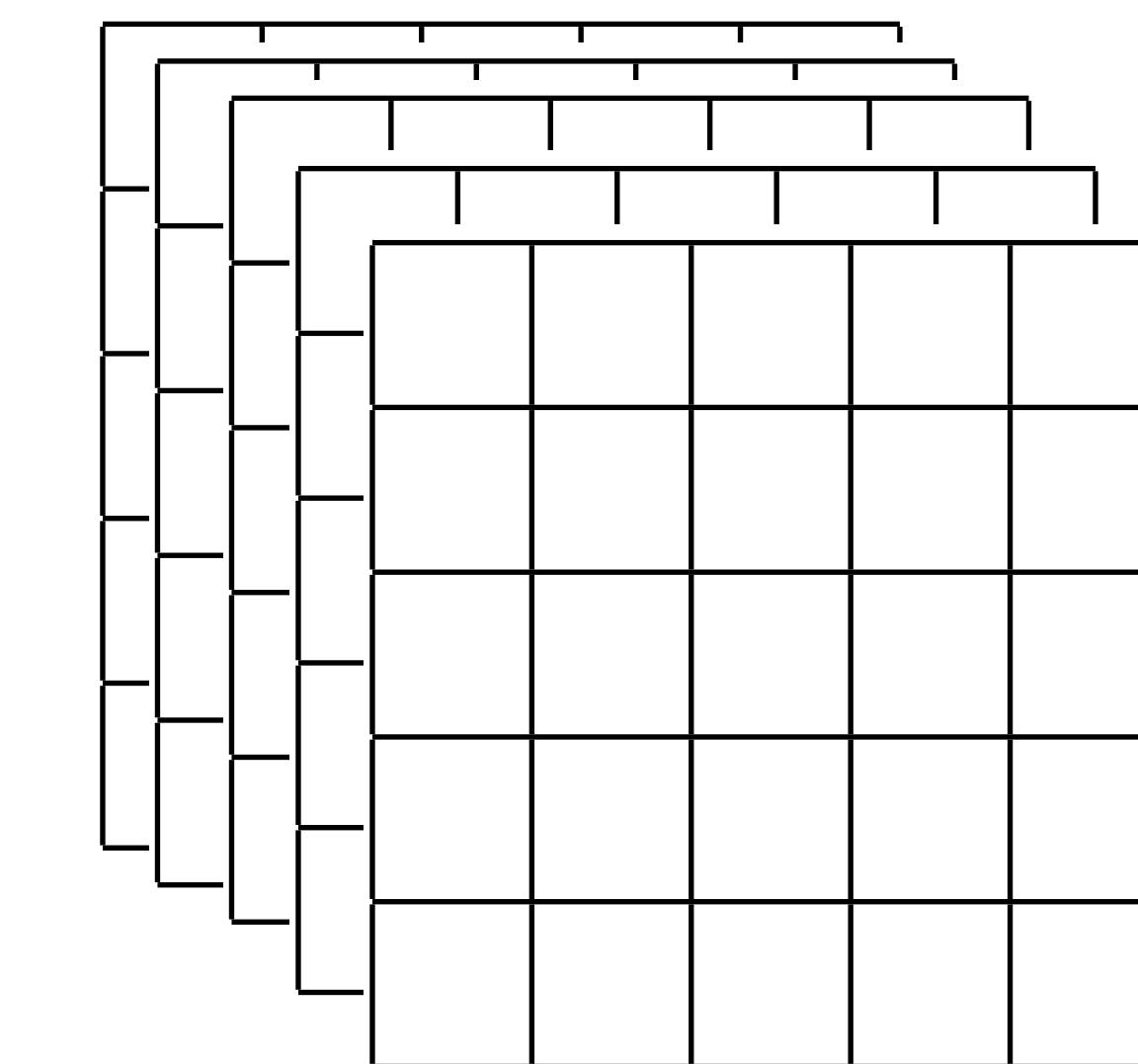
RGB image
($H \times W \times 3$)



pixel-wise CE-loss
($H \times W \times N$)



ground truth (1-hot encoding)

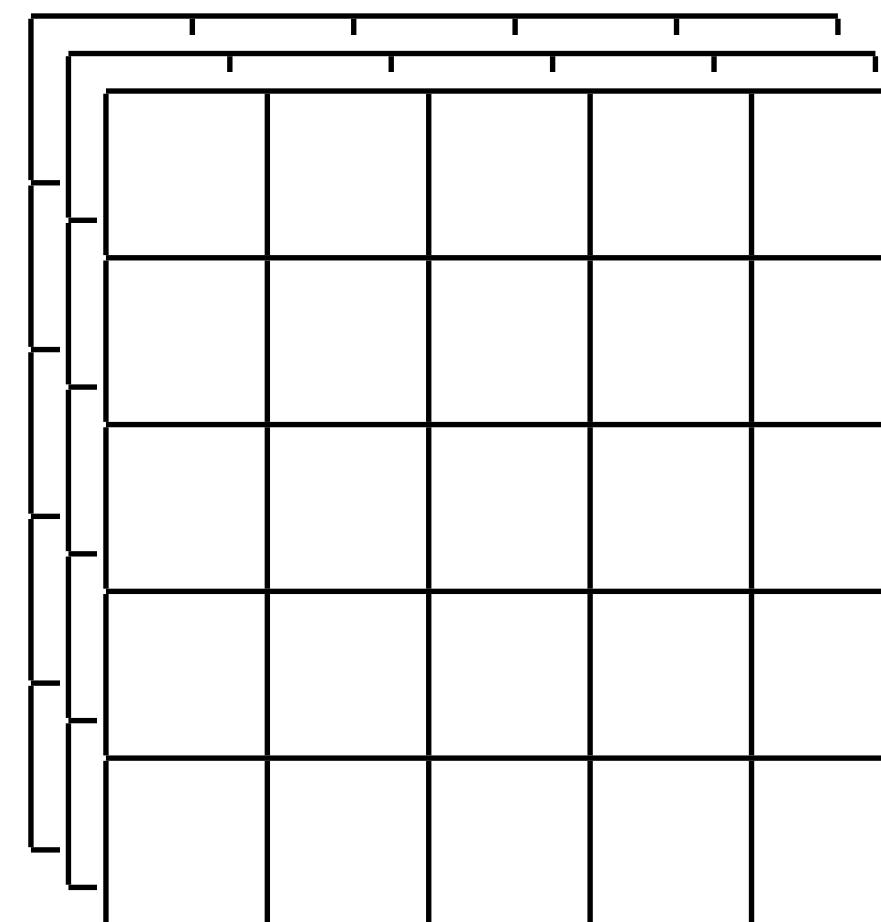


annotations
($H \times W \times N$)

Semantic segmentation

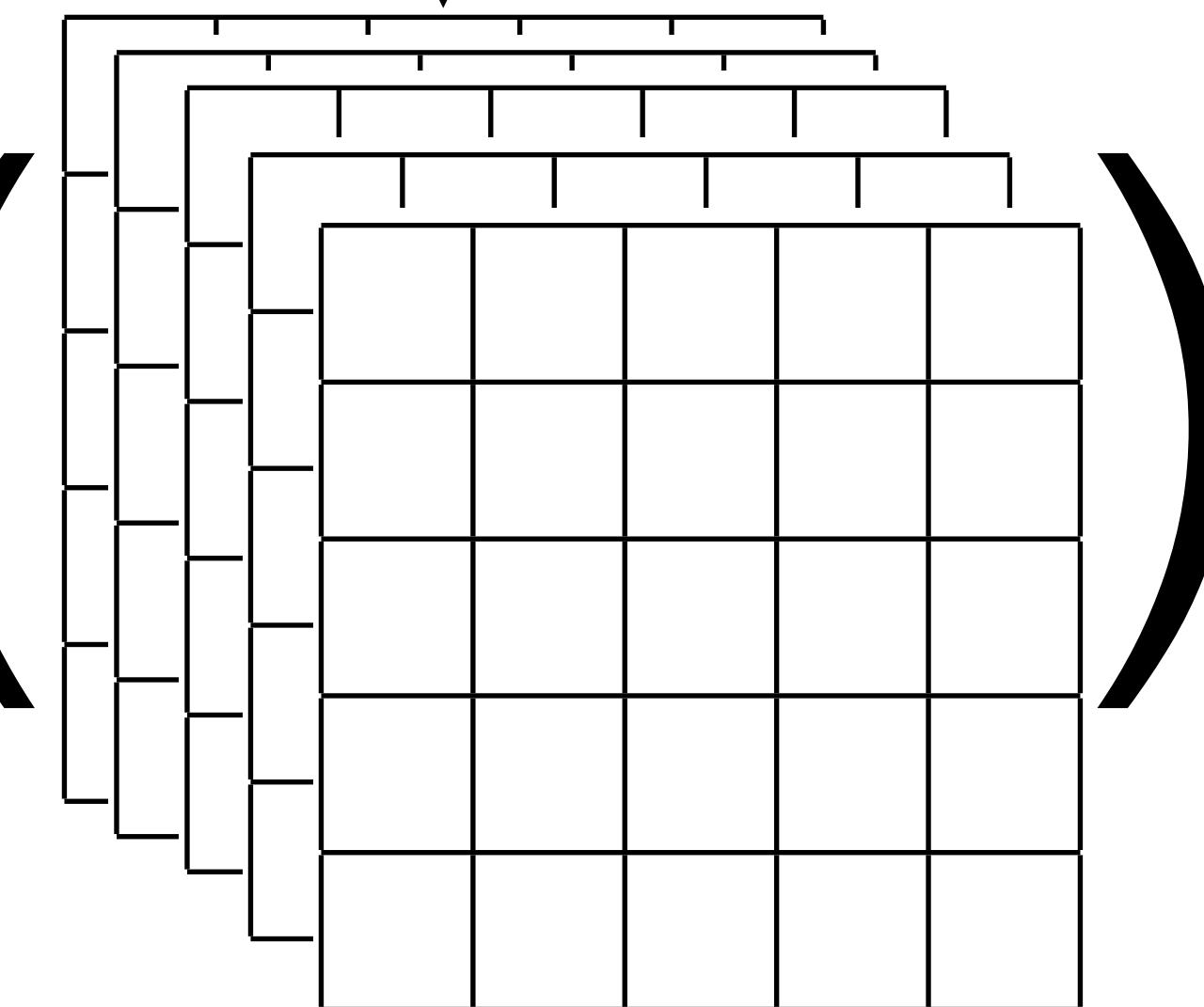
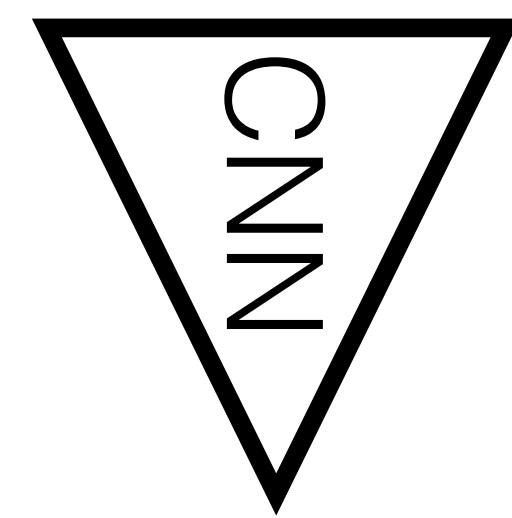
Pixel-wise classifier with the loss summed over all pixels

RGB image
($H \times W \times 3$)

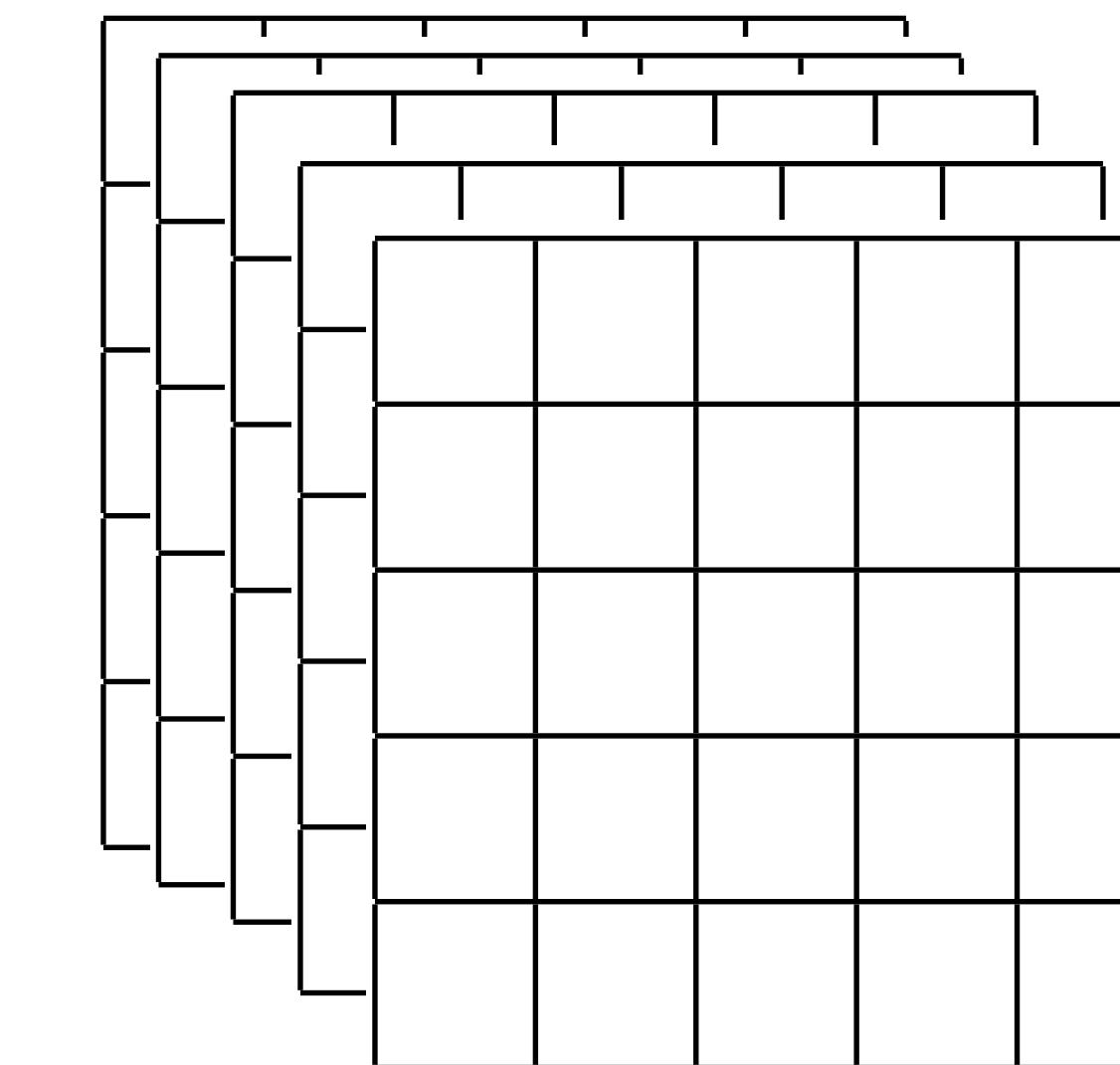
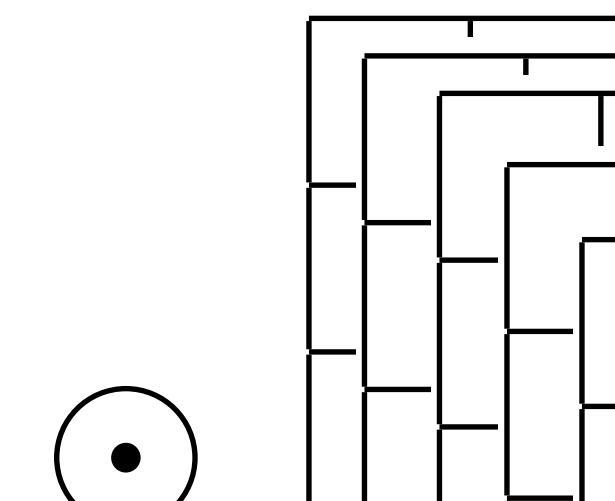


pixel-wise CE-loss
($H \times W \times N$)

$\sum_{\text{pixels}} - \log \left(\frac{\delta}{\text{channels}} \right)$



ground truth (1-hot encoding)



annotations
($H \times W \times N$)

U-net architecture



Mirror the usual classification network

U-net architecture



[Noh et al ICCV 2015] <https://arxiv.org/pdf/1505.04366.pdf>

U-net architecture

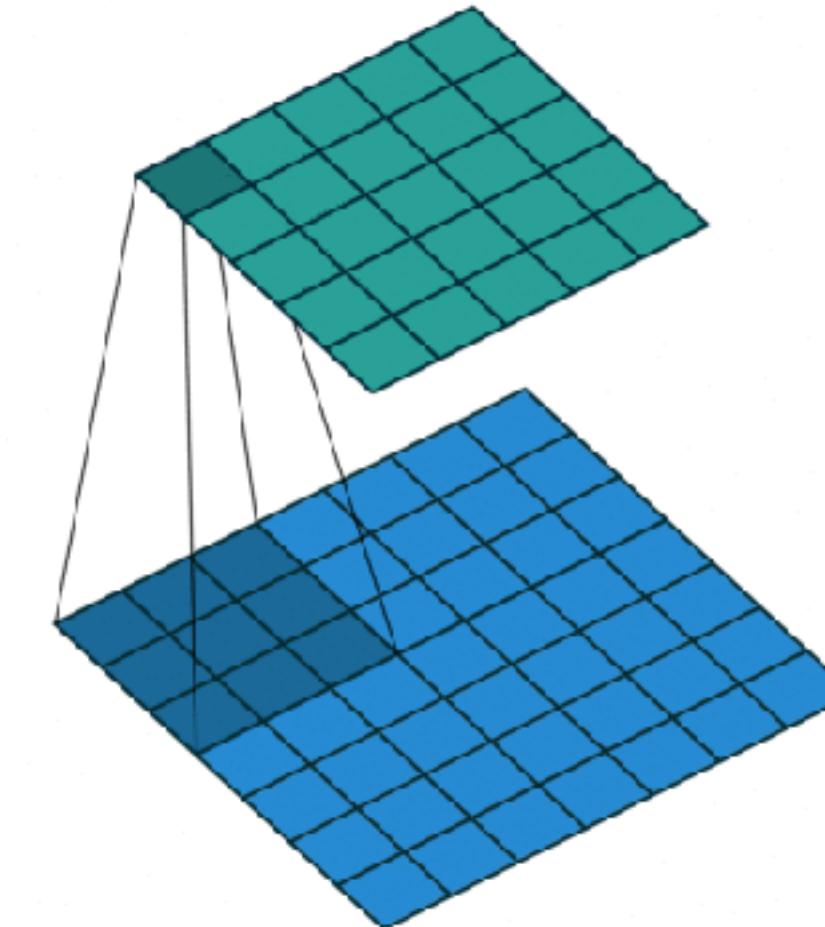
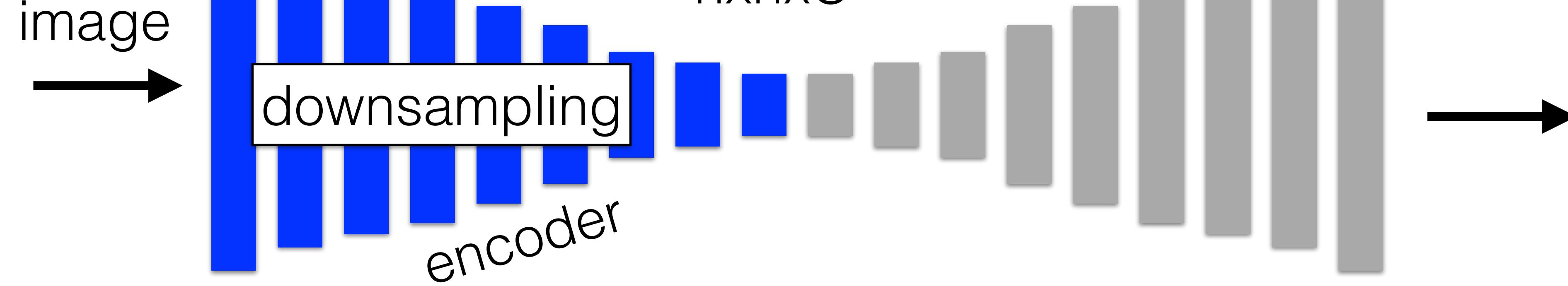
high spatial resolution

$N \times N \times C$

image embedding

high channel resolution

$n \times n \times C$

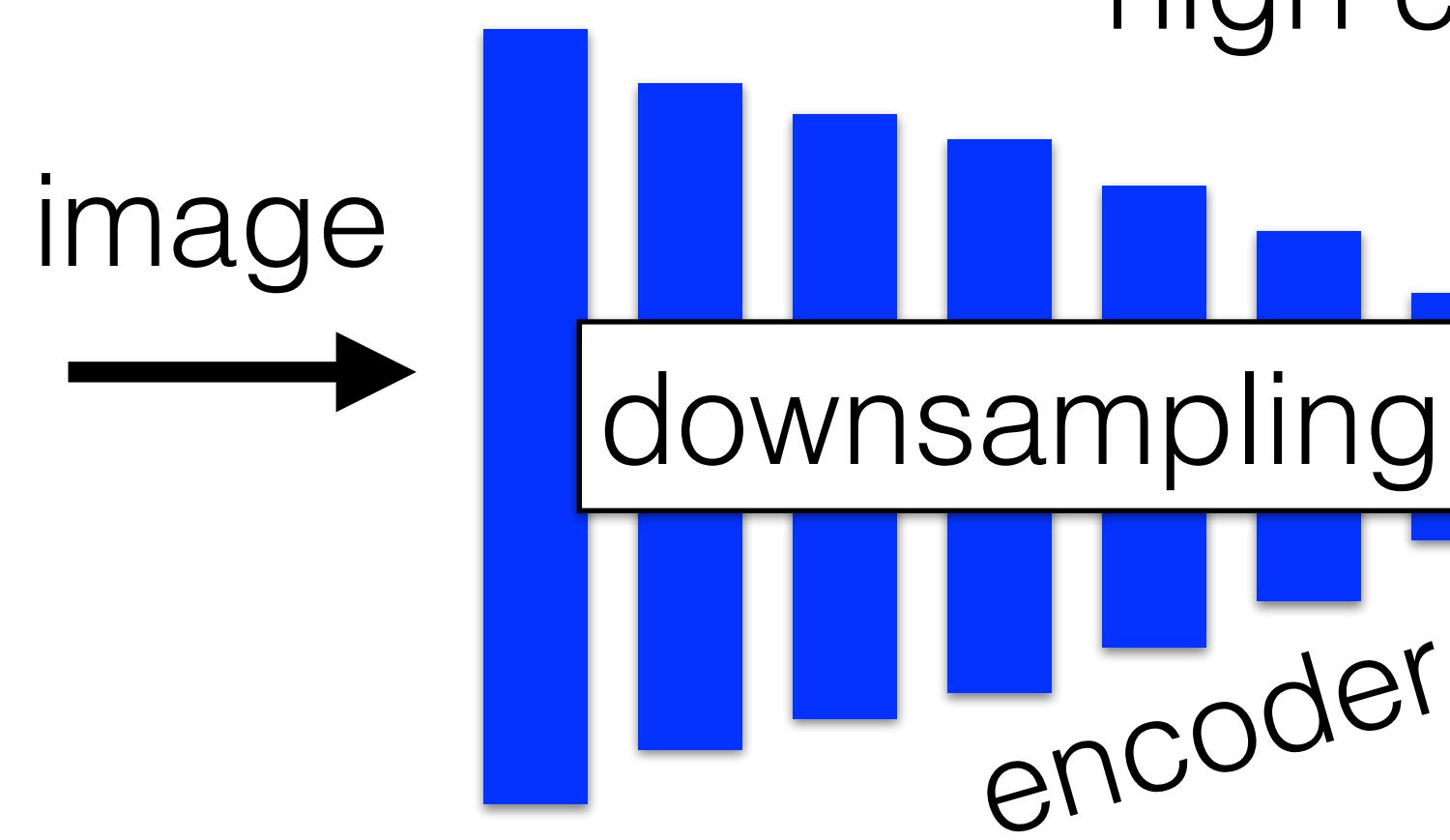


convolution layers

U-net architecture

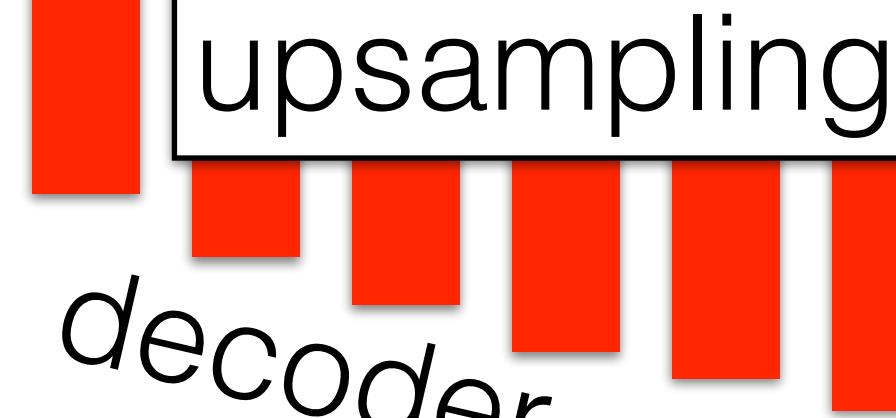
high spatial resolution
 $N \times N \times C$

image embedding

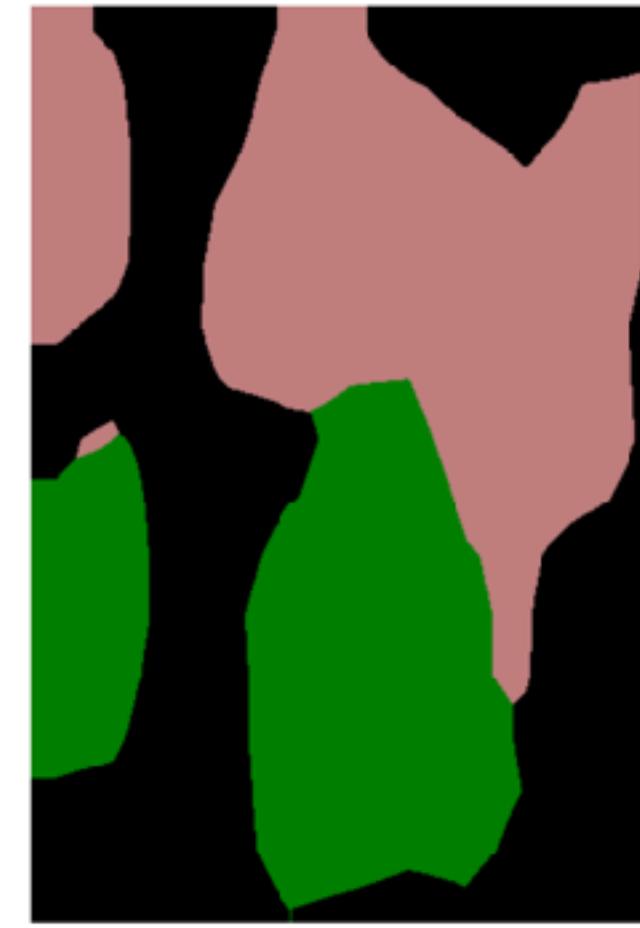


high channel resolution
 $n \times n \times C$

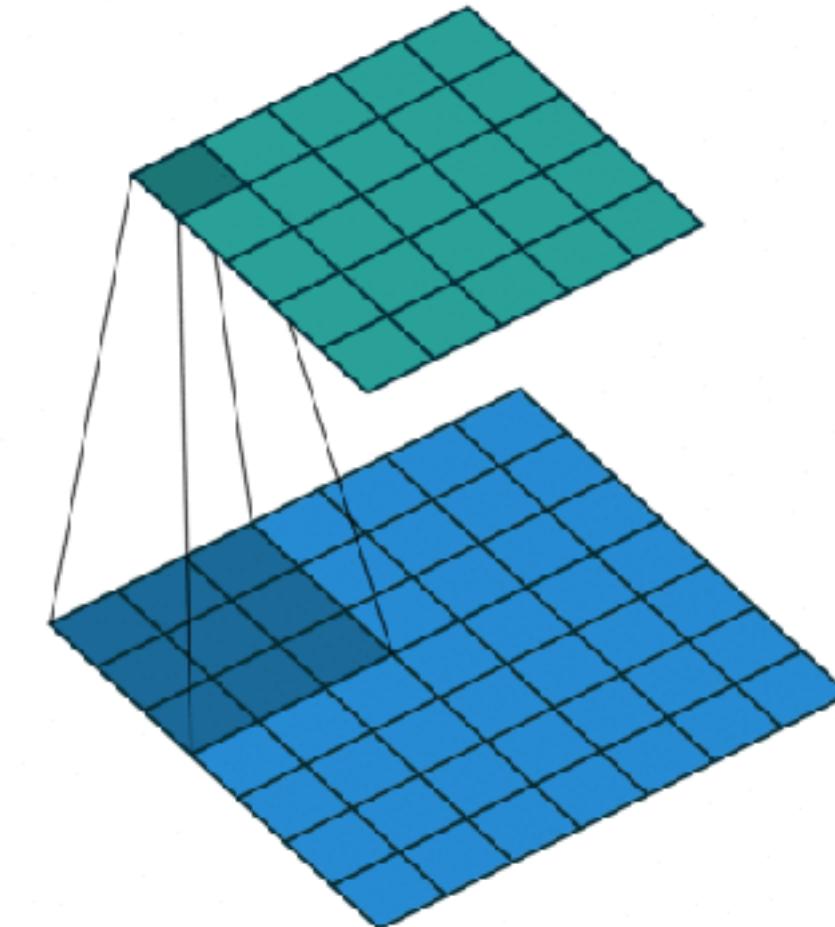
high spatial resolution
 $N \times N \times C$



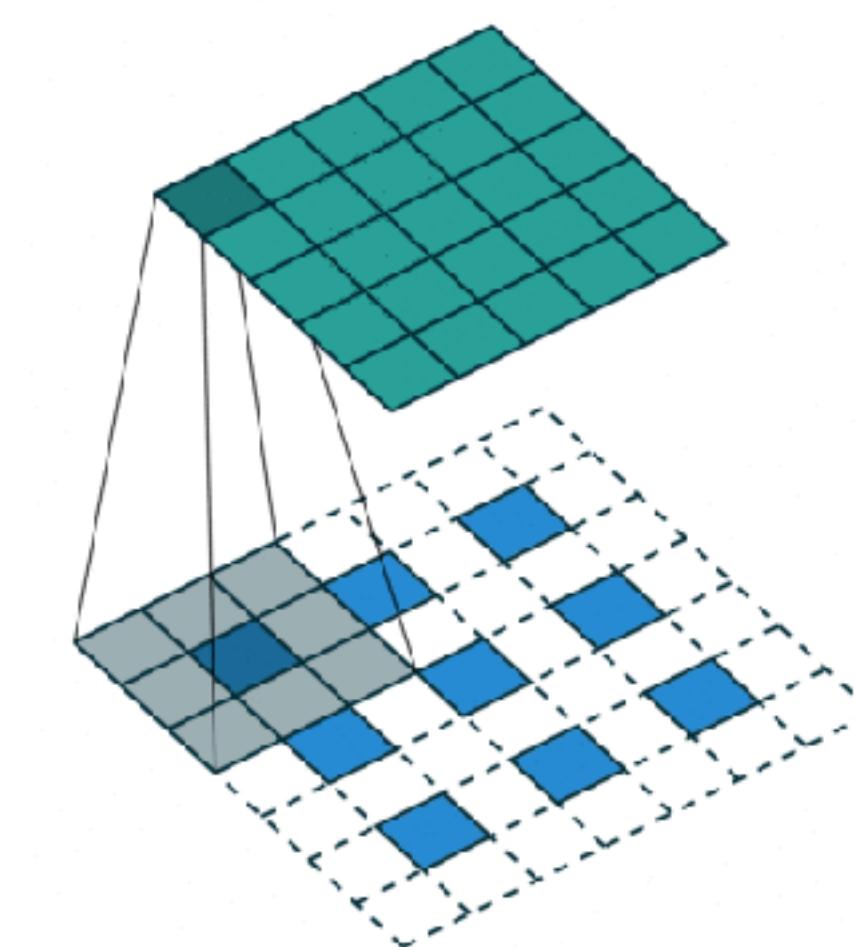
Why???
output



ground truth



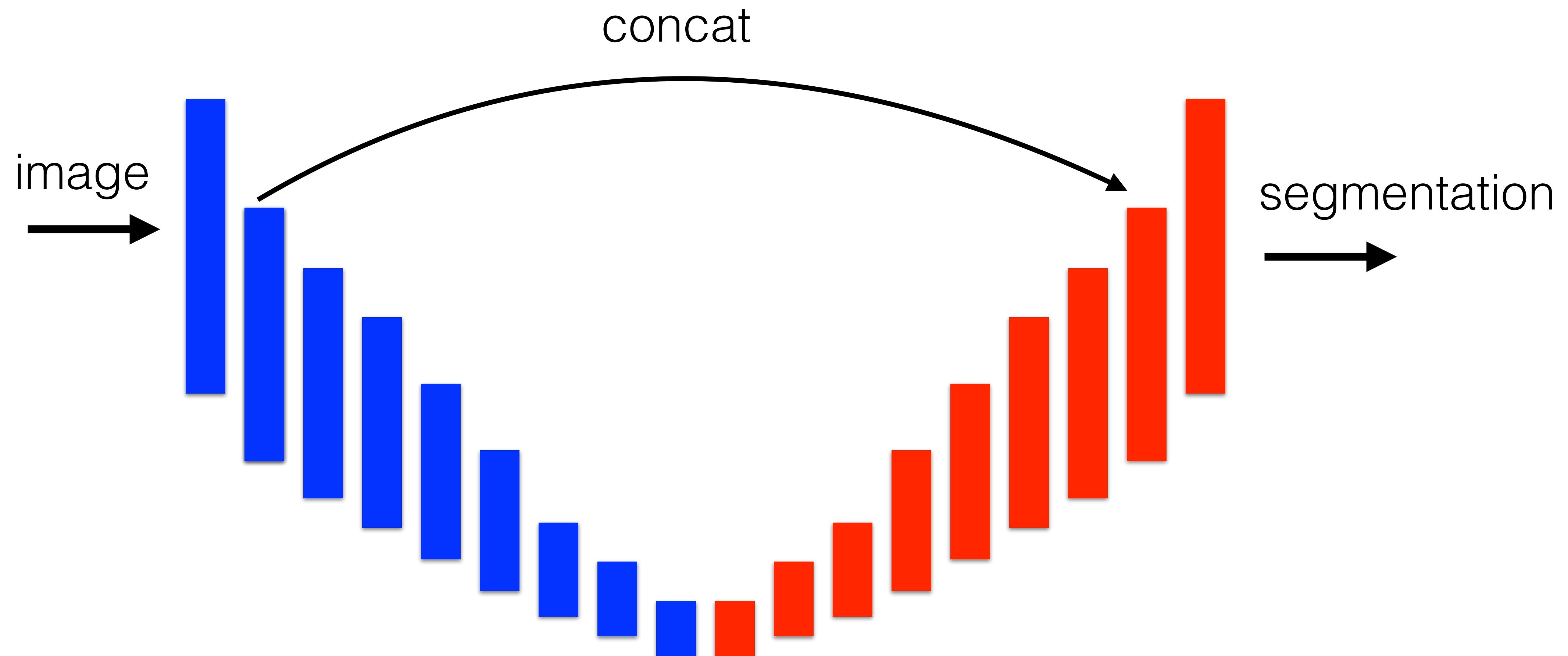
convolution layers



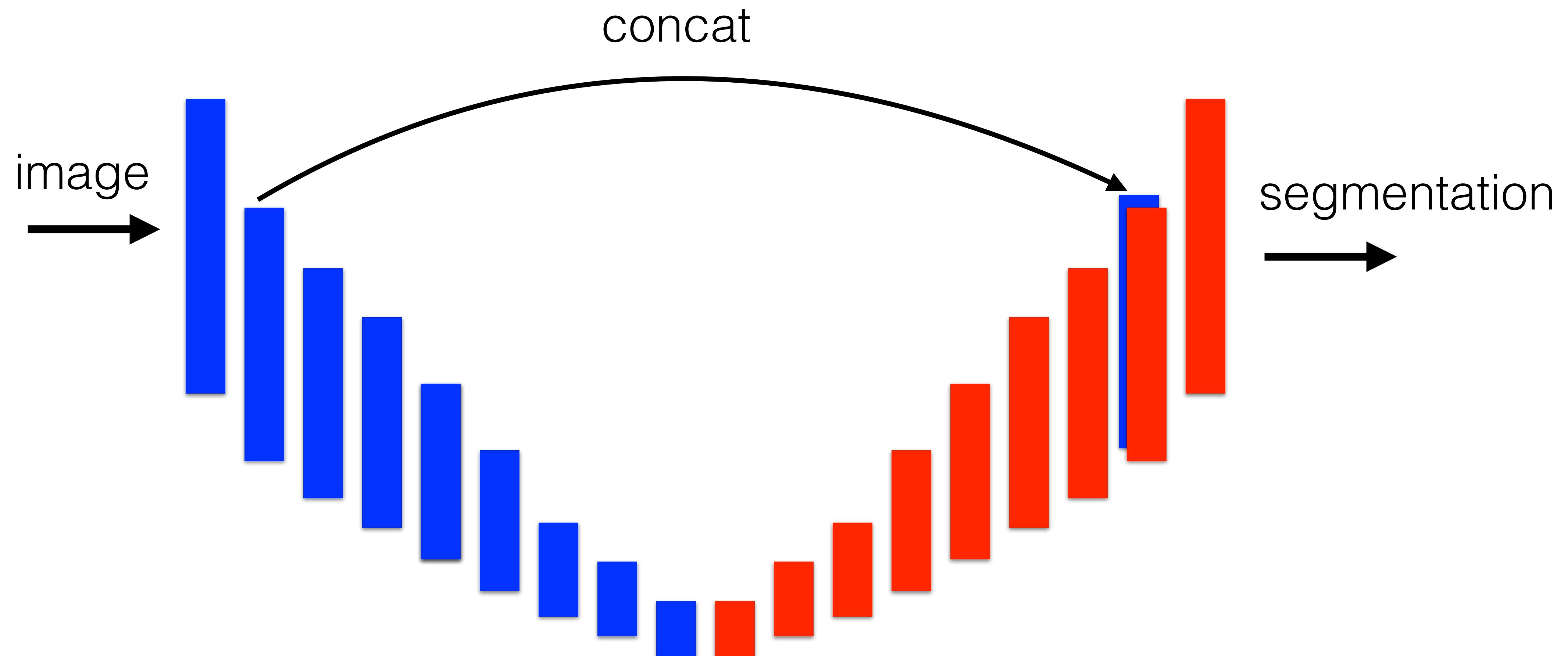
deconvolution layers



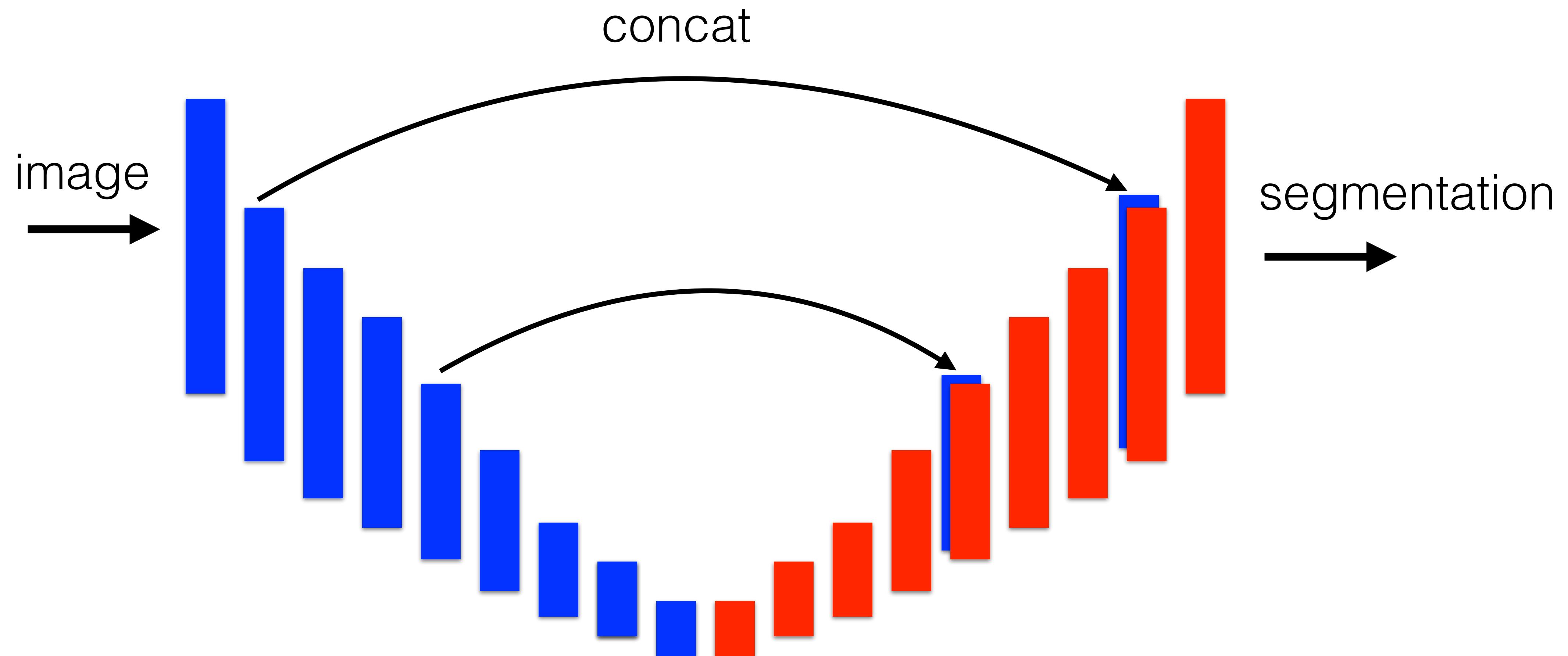
U-net architecture



U-net architecture

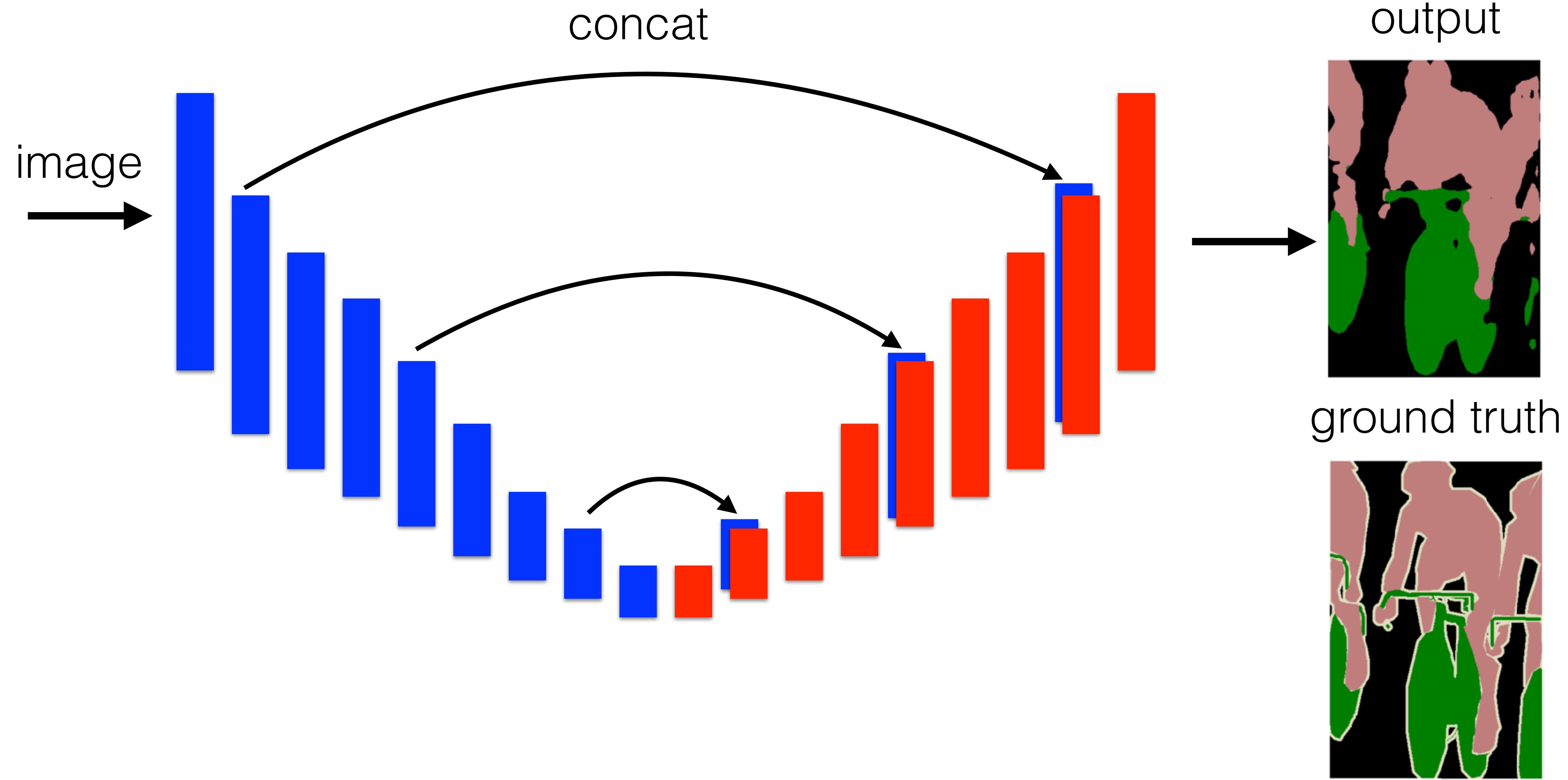


U-net architecture



[Noh et al ICCV 2015] <https://arxiv.org/pdf/1505.04366.pdf>

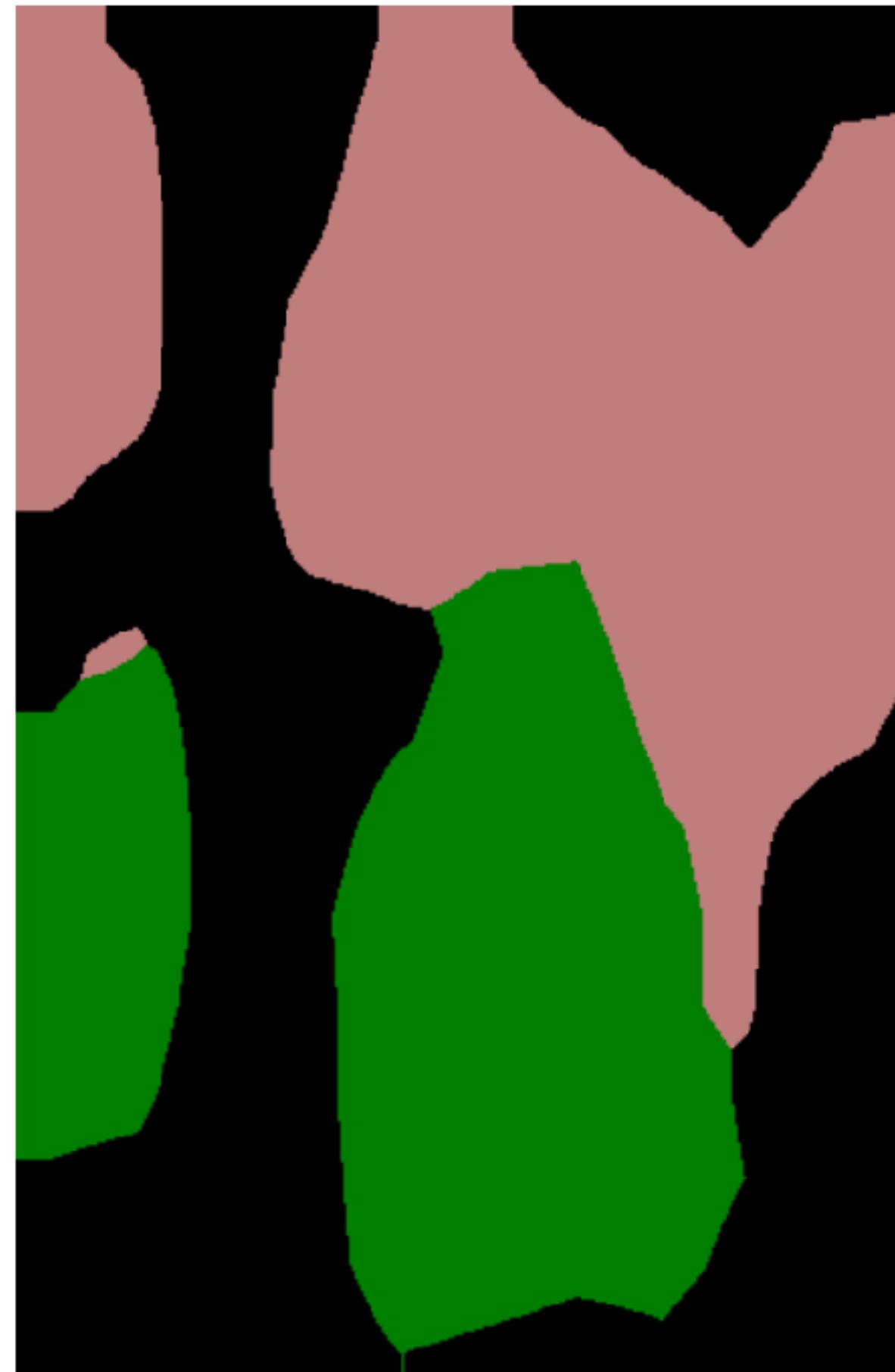
U-net architecture



[Noh et al ICCV 2015] <https://arxiv.org/pdf/1505.04366.pdf>

U-net architecture

no skip connections



with skip connections



ground truth



Segmentation architectures

- **FCN** (Fully Convolutional Network): LongCVPR 2015, FCN was one of the pioneering architectures for end-to-end pixel-wise prediction.
- **U-Net**: symmetric architecture and skip connections, which helps in capturing both global and local features.
- **DeepLab**: DeepLab uses dilated convolutions to enlarge the field of view and capture multi-scale information. DeepLabv3 is an improved version.
- **MobileNet**: lightweight architecture designed for real-time semantic segmentation. It is known for its efficiency and speed.
- **LinkNet**: LinkNet employs a skip connection structure with a series of encoder and decoder blocks for segmentation tasks.
- **HRNet** (High-Resolution Network): HRNet focuses on maintaining high-resolution representations throughout the network, which can be beneficial for capturing fine details.
- **ViTS**: Vision transformer exptended for segmentation task
- **SAM**: Segment anything architecture from Facebook 2023