

What can('t) we do we 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



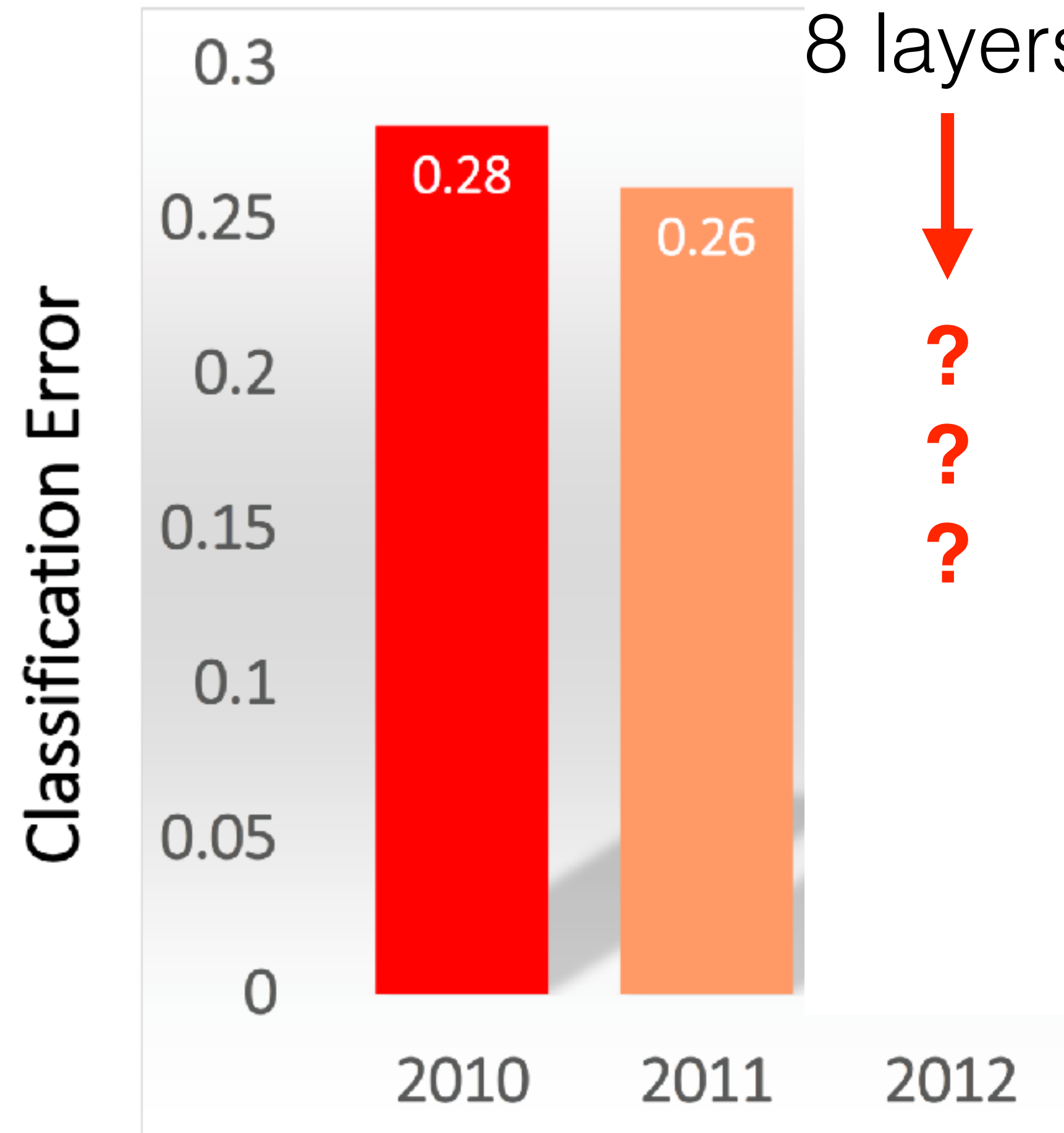
Output:
Scale
T-shirt
Giant panda
Drumstick
Mud turtle



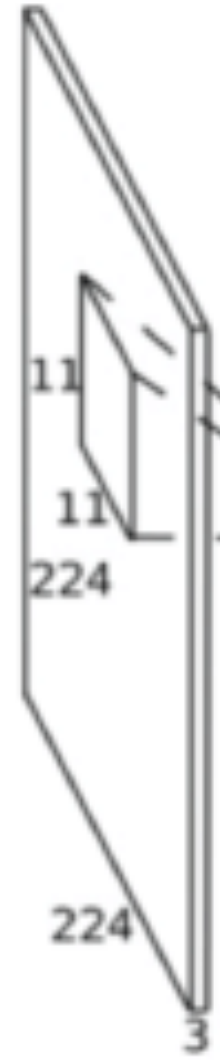
$$\text{Error} = \frac{1}{100,000} \sum_{100,000 \text{ images}} 1[\text{incorrect on image } i]$$

Classification results

AlexNet
8 layers



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

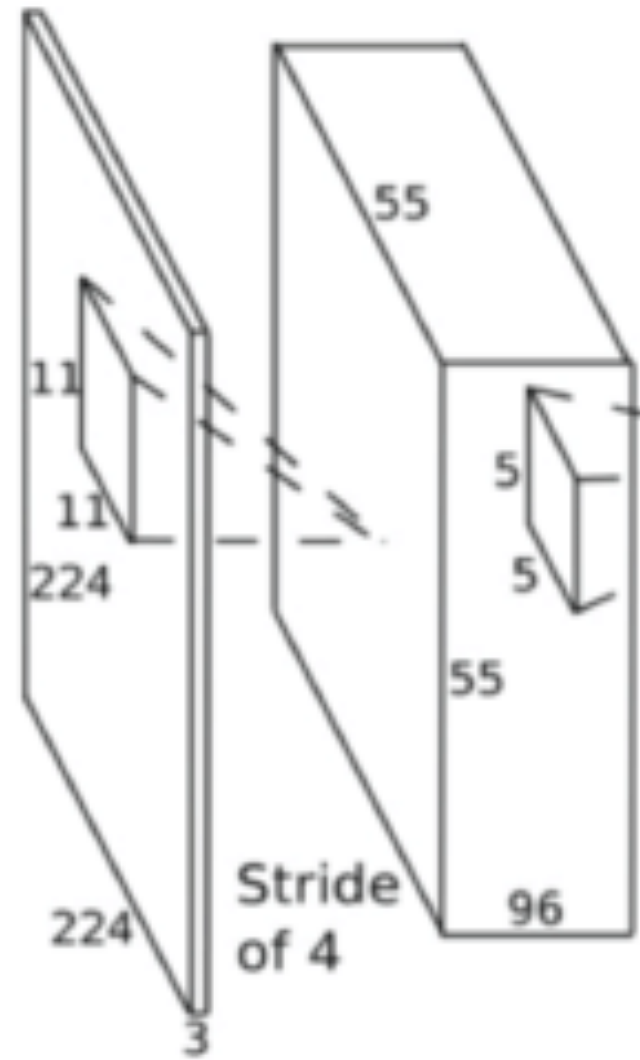


- Param in layer1 (conv, 96 11x11 filters, stride=4, pad=0)?

Alex Krizhevsky et al, Imagenet classification with deep convolutional neural networks, NIPS, 2012

<https://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks.pdf>

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

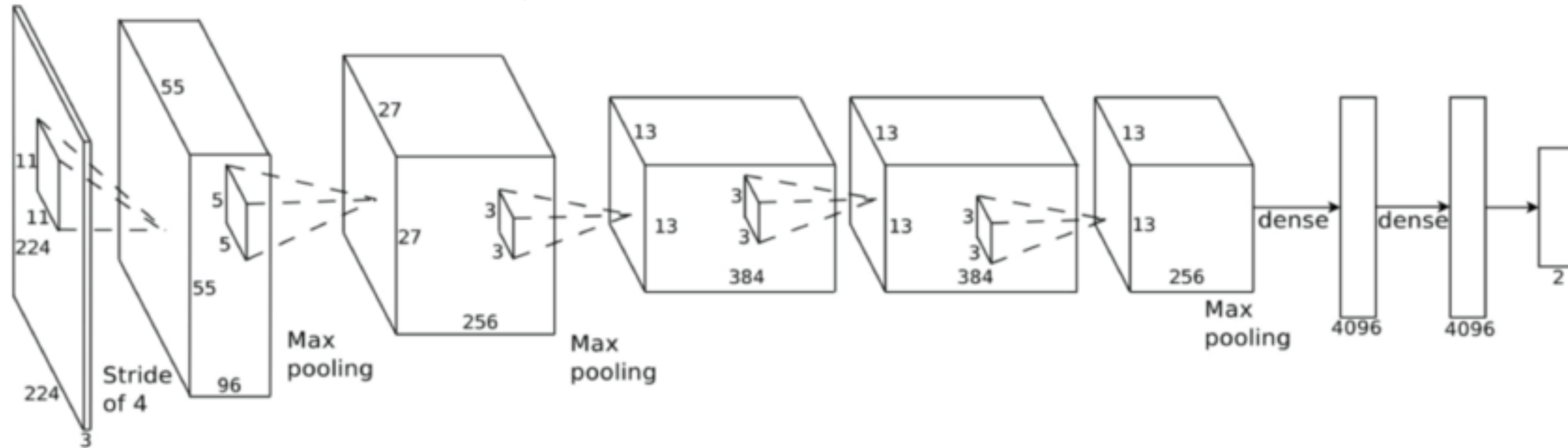


- Param in layer1 (conv, 96 11x11 filters, stride=4, pad=0)?
- Param in layer2 (maxp, 3x3 filters, stride=2, pad=0)?

Alex Krizhevsky et al, Imagenet classification with deep convolutional neural networks, NIPS, 2012

<https://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks.pdf>

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

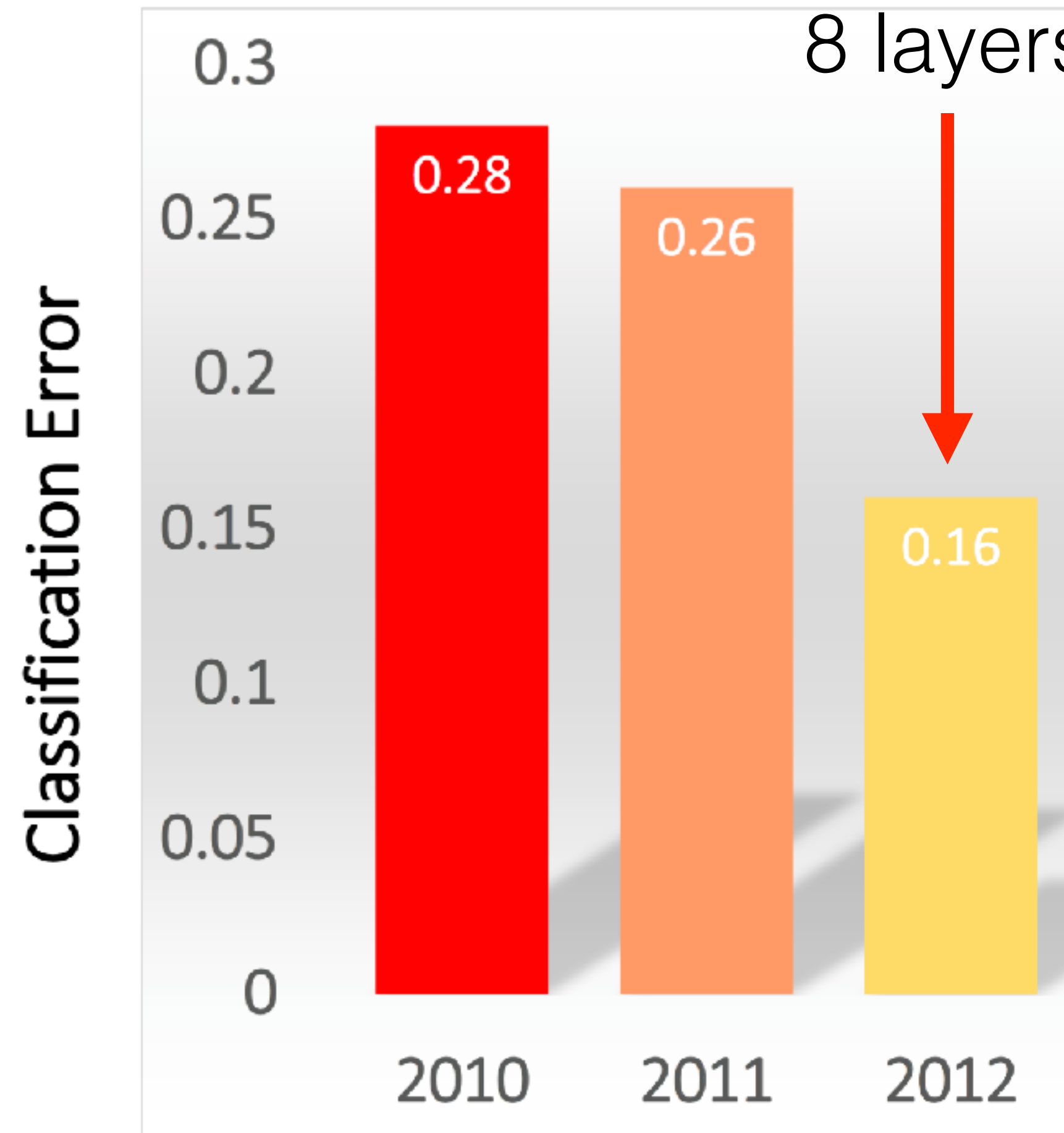
Alex Krizhevsky et al, Imagenet classification with deep convolutional neural networks, NIPS, 2012

<https://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks.pdf>

Classification results

AlexNet

8 layers

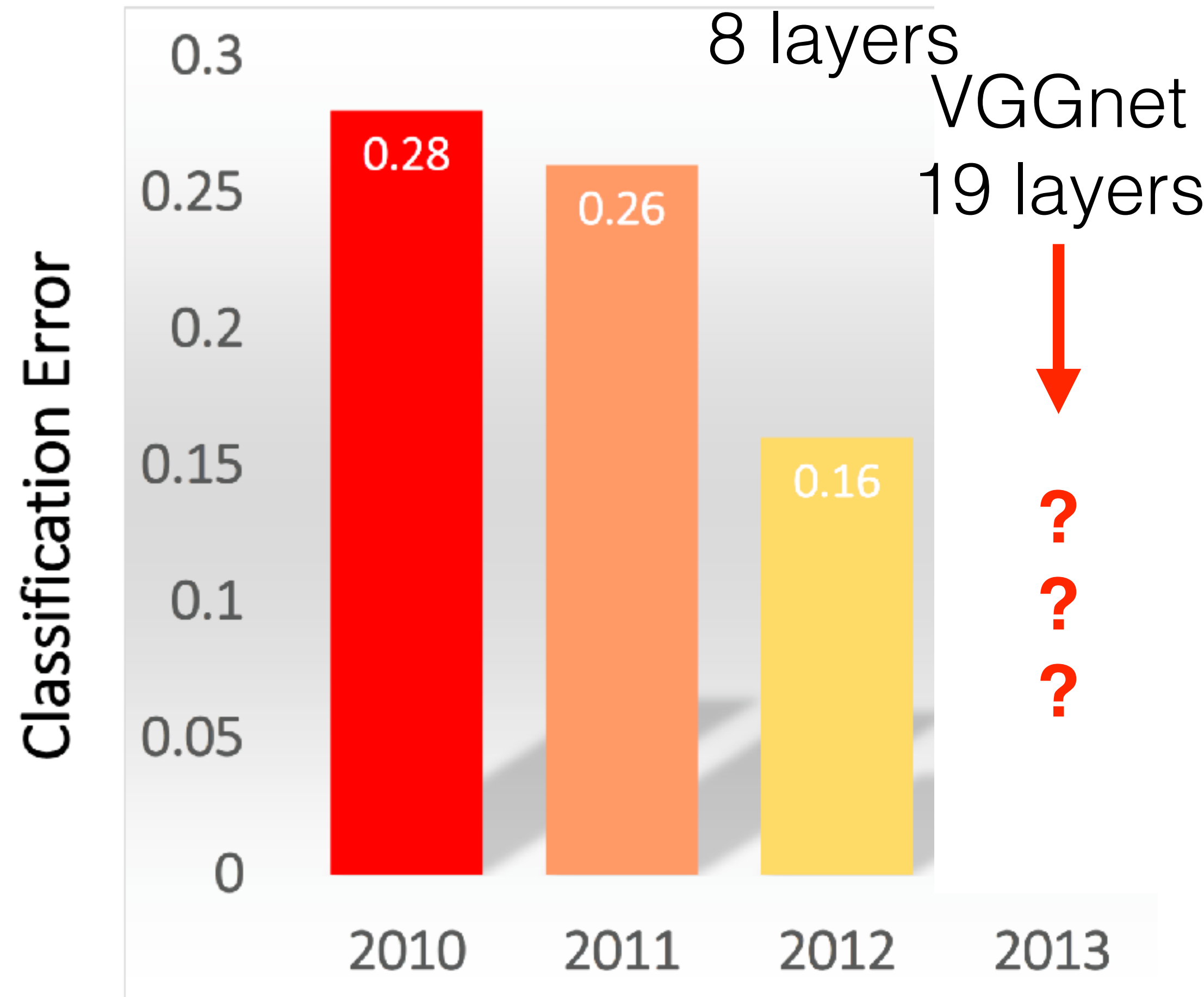


Classification results

AlexNet

8 layers

VGGnet
19 layers



VGGNet vs AlexNet



- large filters
- shallow (8 layers)



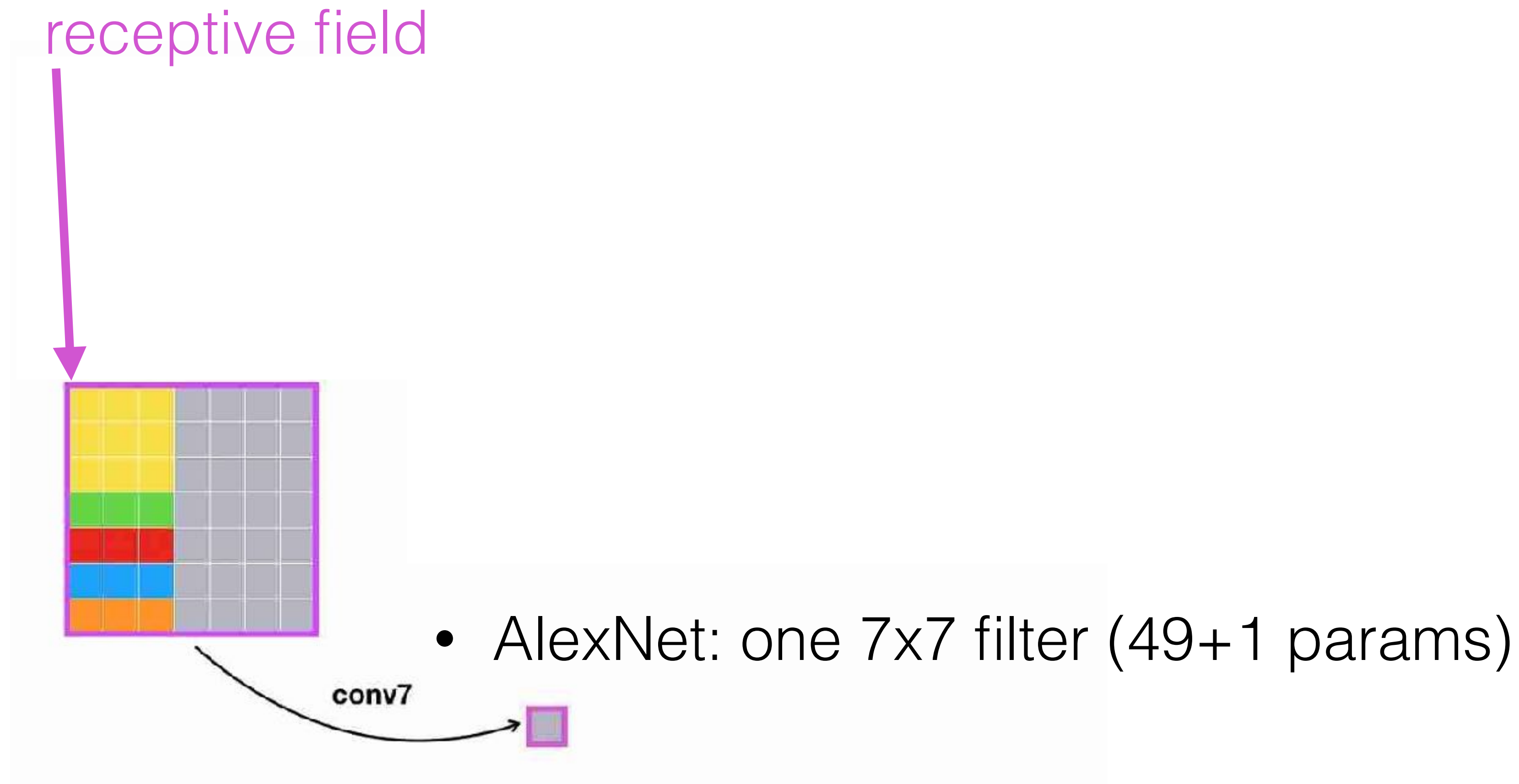
- small filters
- deeper (19 layers)

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

Simonyan and Zissermann, Very Deep Convolutional Networks for Large Scale Image Recognition, 2014

<https://arxiv.org/abs/1409.1556>

VGGNet vs AlexNet

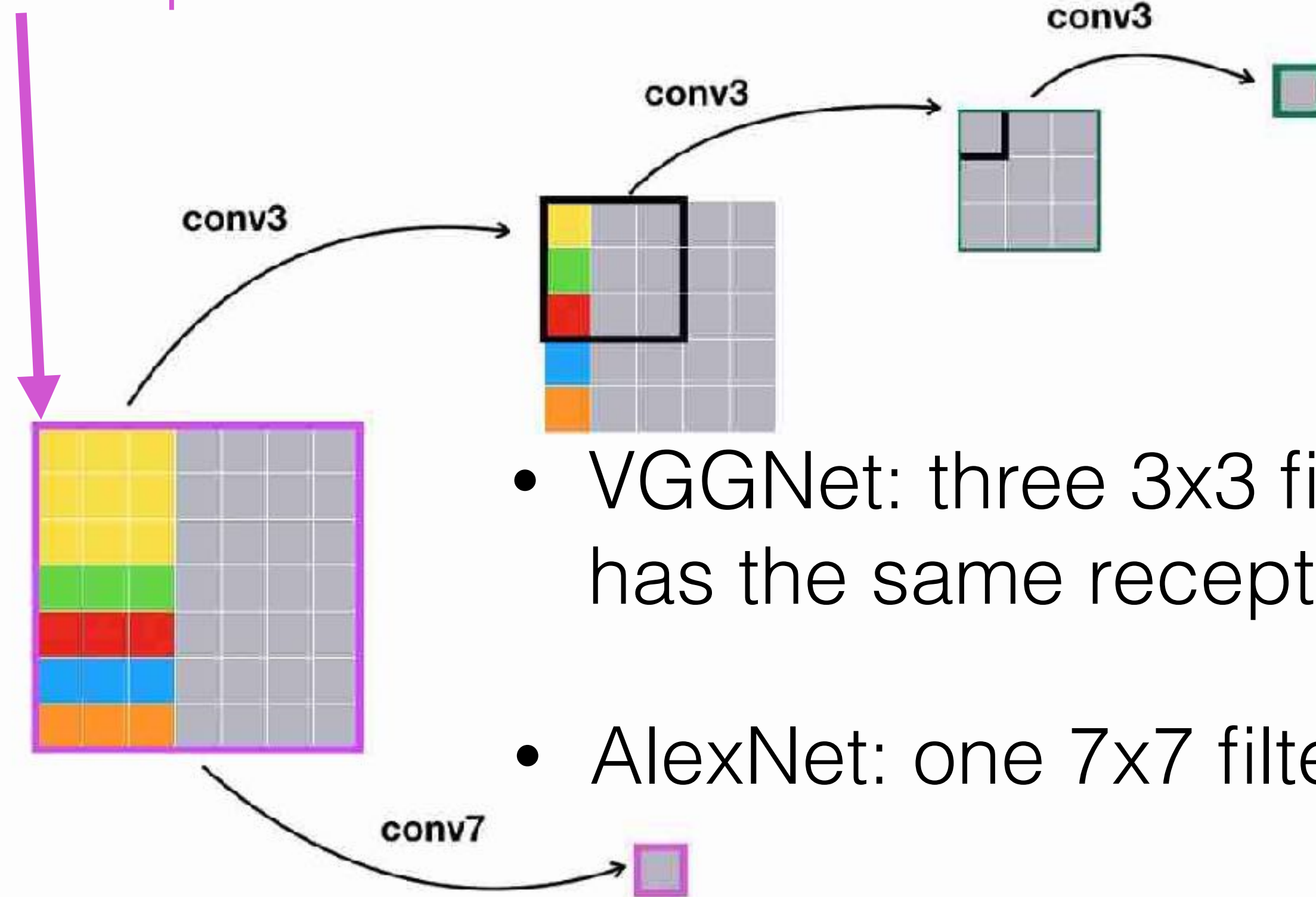


Receptive field of the neuron is a region of pixels which can influence its output

<https://arxiv.org/abs/1409.1556>

VGGNet vs AlexNet

receptive field



- VGGNet: three 3x3 filters ($3 \times 9 + 3$ params) has the same receptive field
- AlexNet: one 7x7 filter ($49 + 1$ params)

Receptive field of the neuron is a region of pixels which can influence its output

VGGNet has the same receptive field with less parameters

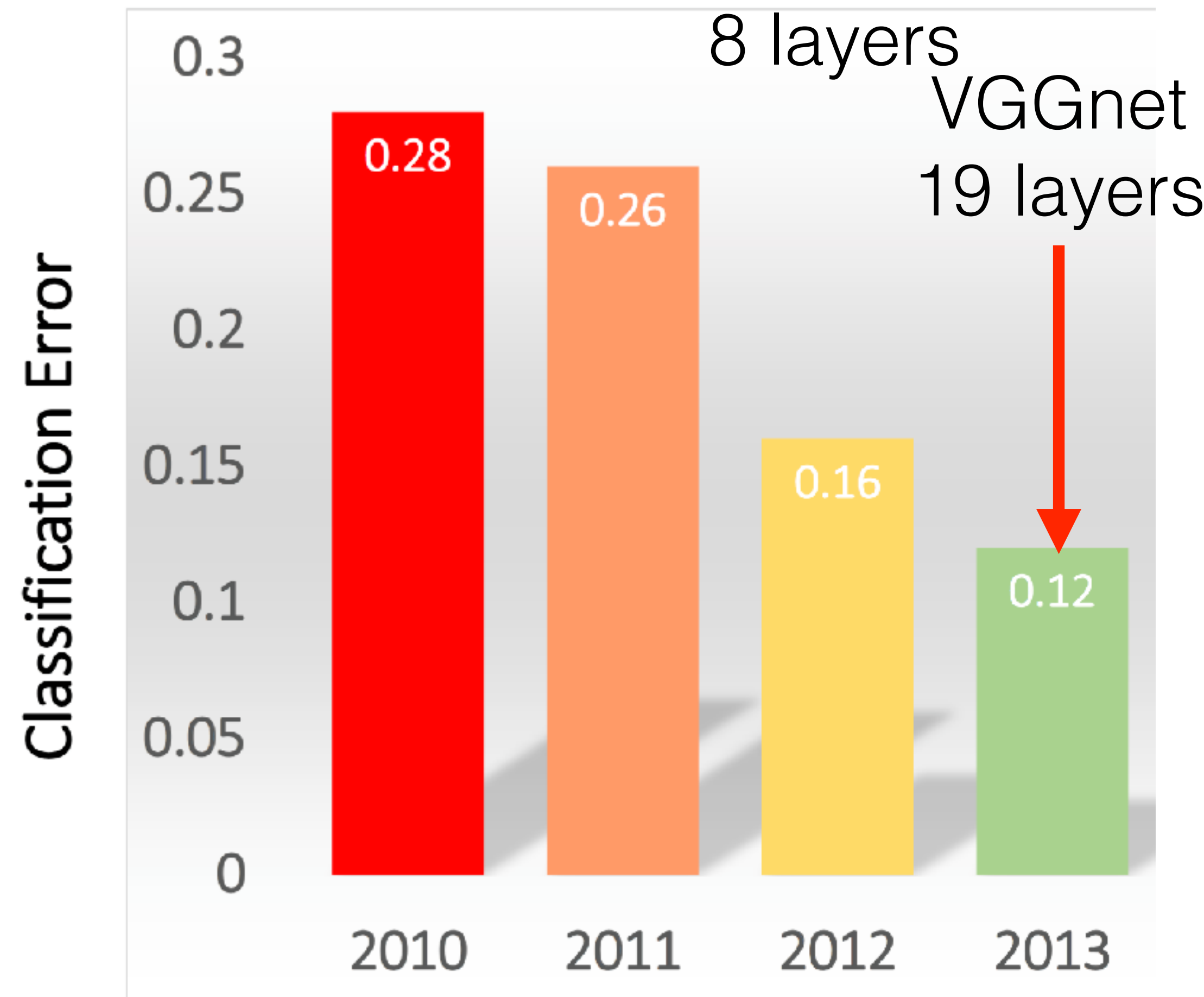
<https://arxiv.org/abs/1409.1556>

Classification results

AlexNet

8 layers

VGGnet
19 layers



Classification results

AlexNet

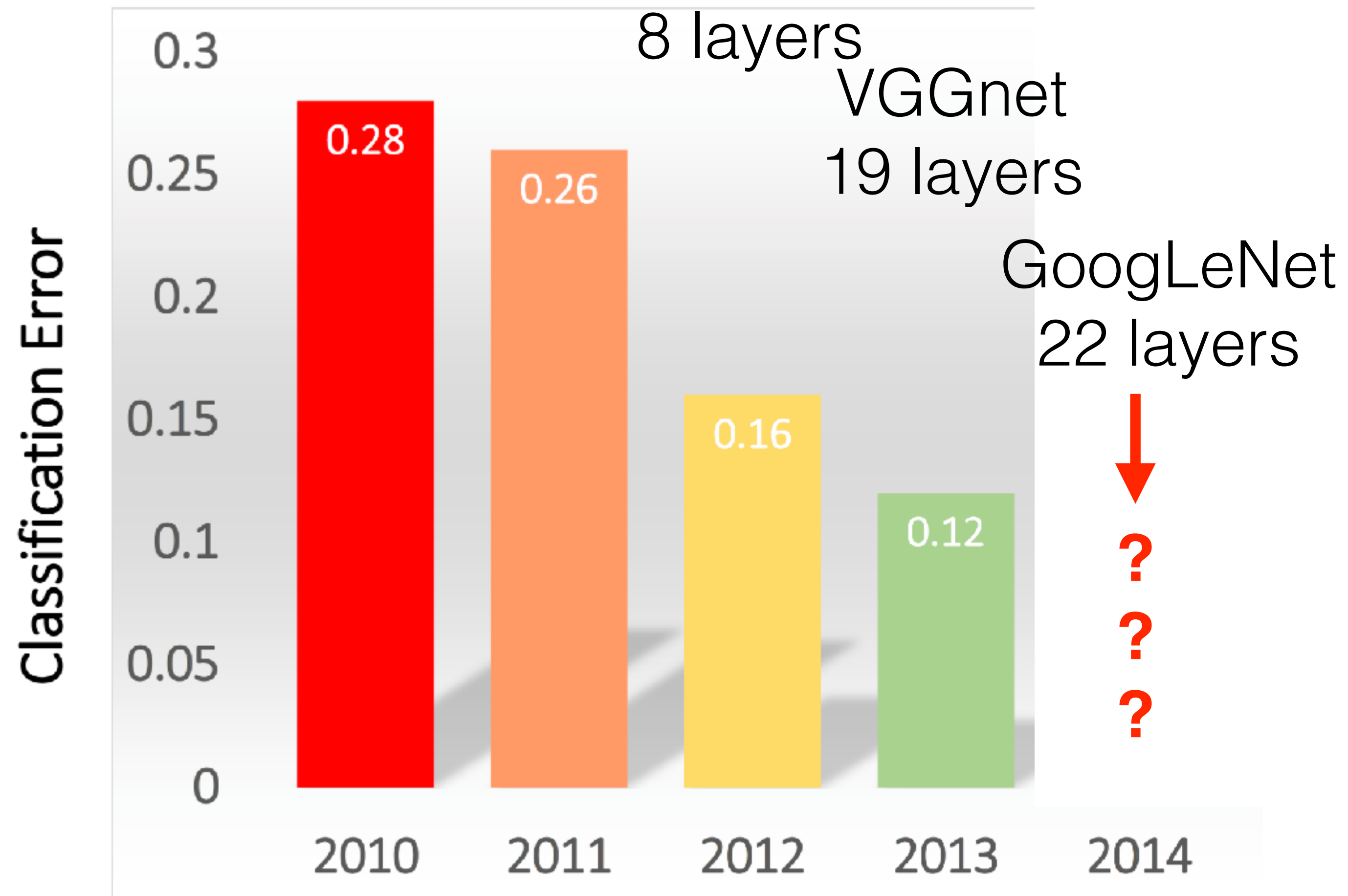
8 layers

VGGnet

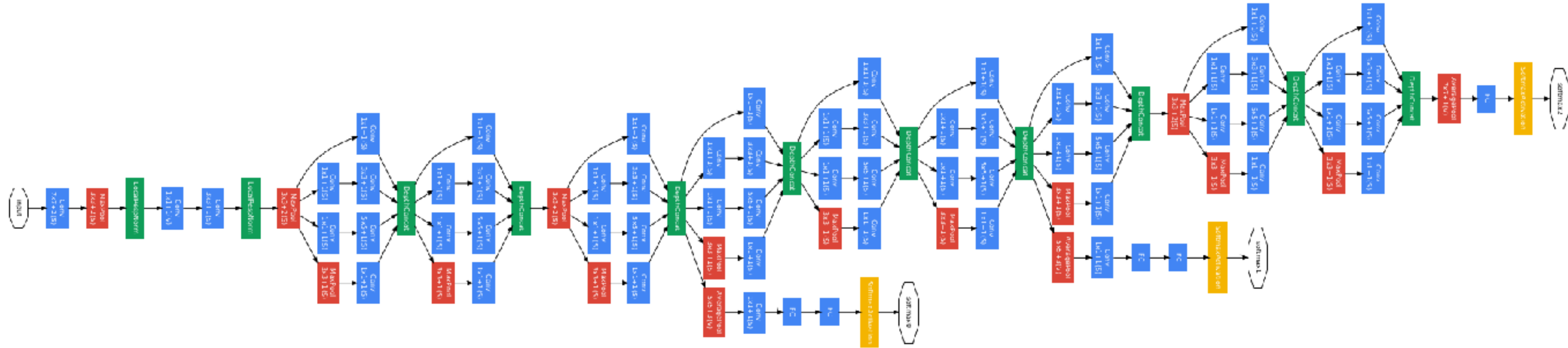
19 layers

GoogLeNet

22 layers

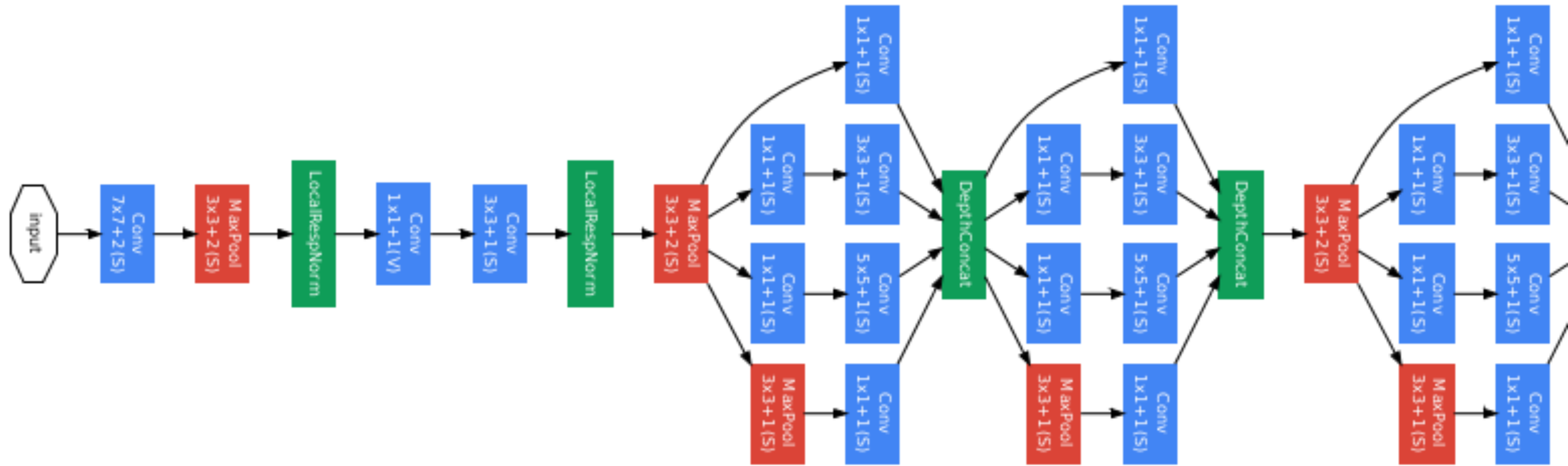


GoogLeNet



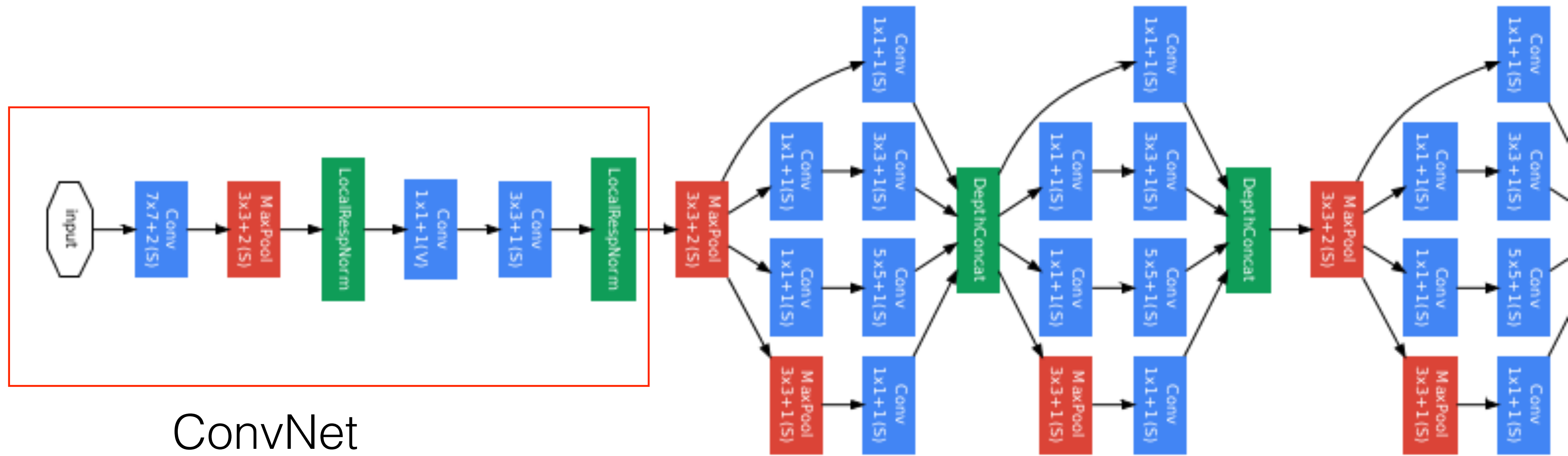
Szegedy et al. Going Deeper with Convolutions, CVPR, 2014
<https://arxiv.org/abs/1409.4842>

GoogLeNet



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

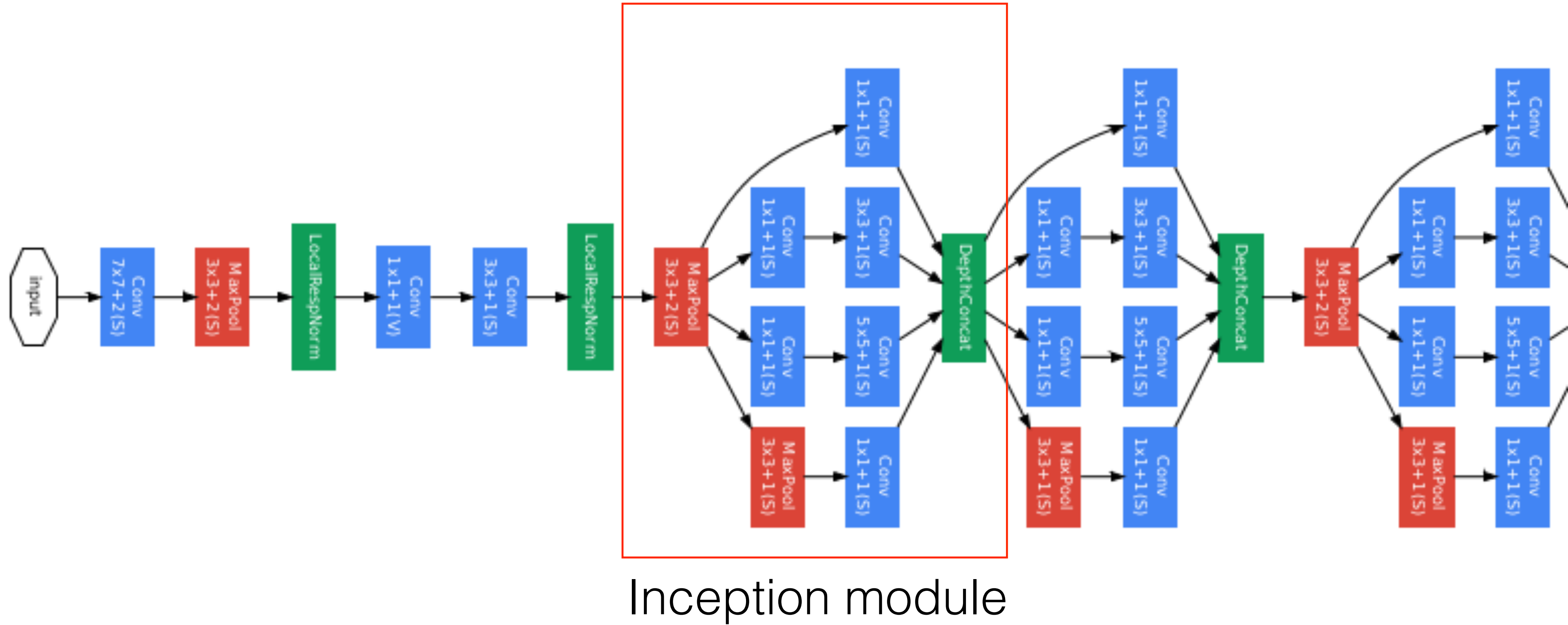
GoogLeNet



ConvNet

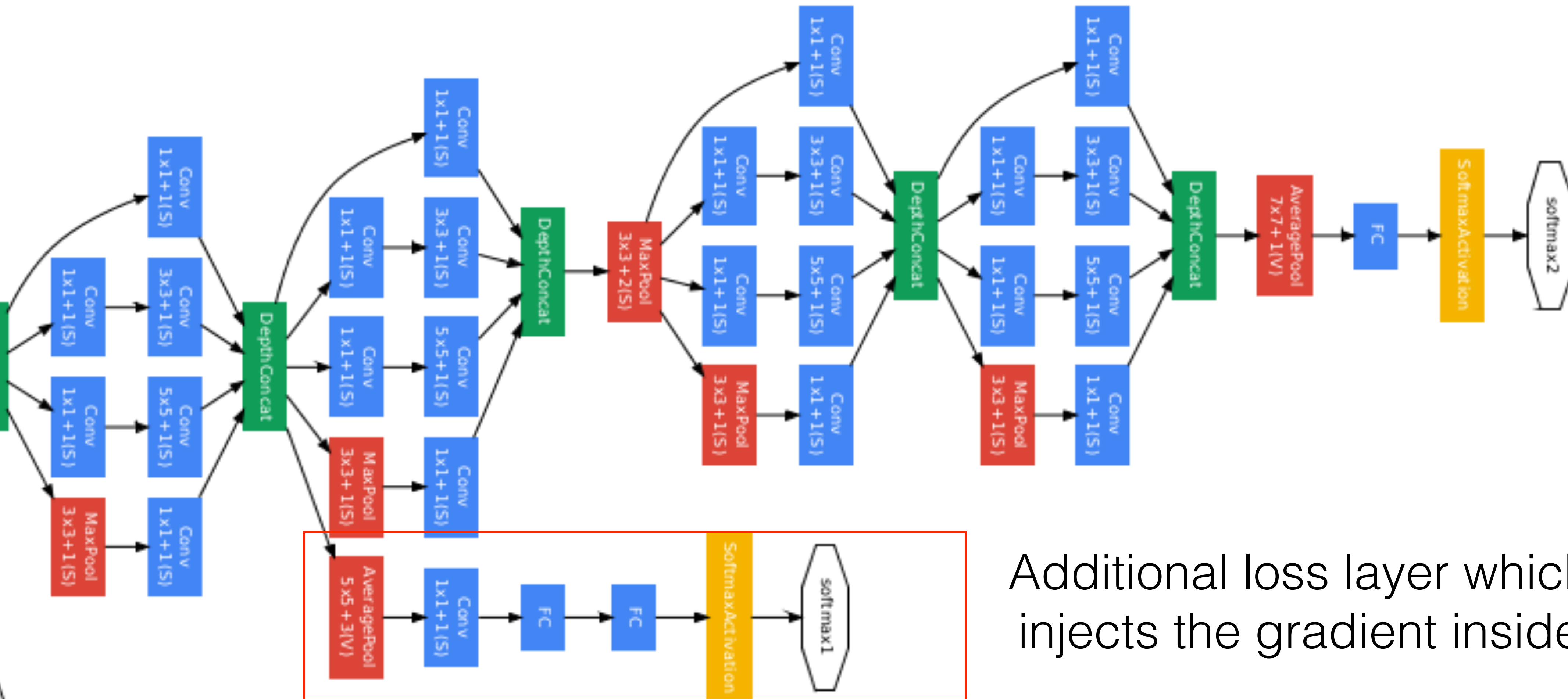
Szegedy et al. Going Deeper with Convolutions, CVPR, 2014
<https://arxiv.org/abs/1409.4842>

GoogLeNet



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

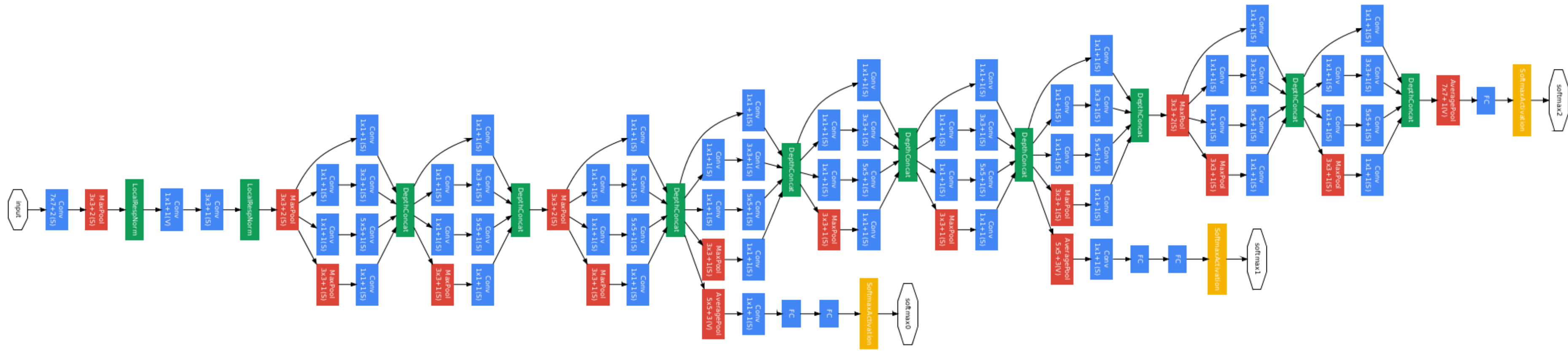
GoogLeNet



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

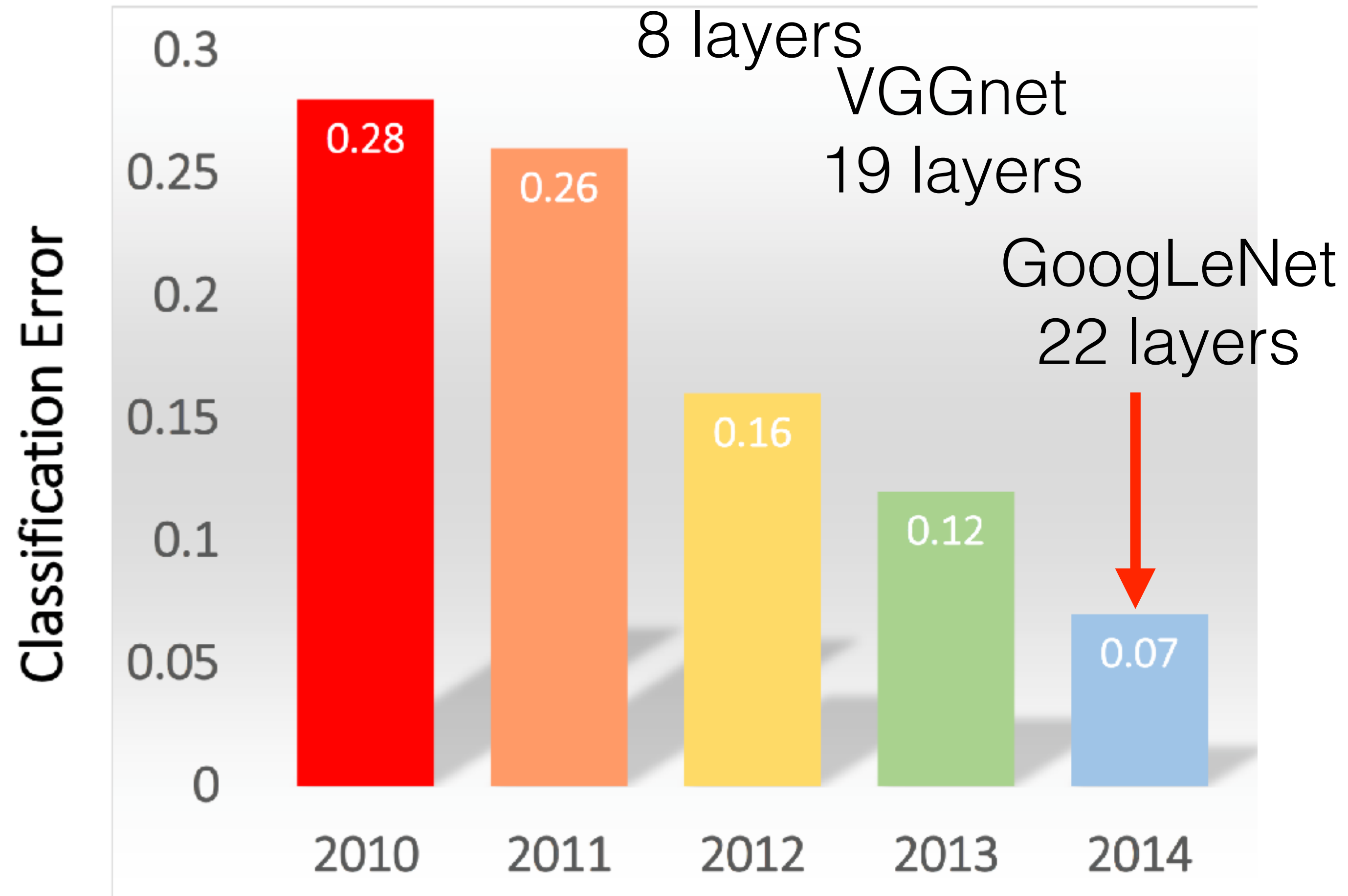
8 layers

VGGnet

19 layers

GoogLeNet

22 layers



Classification results

AlexNet

8 layers

VGGnet

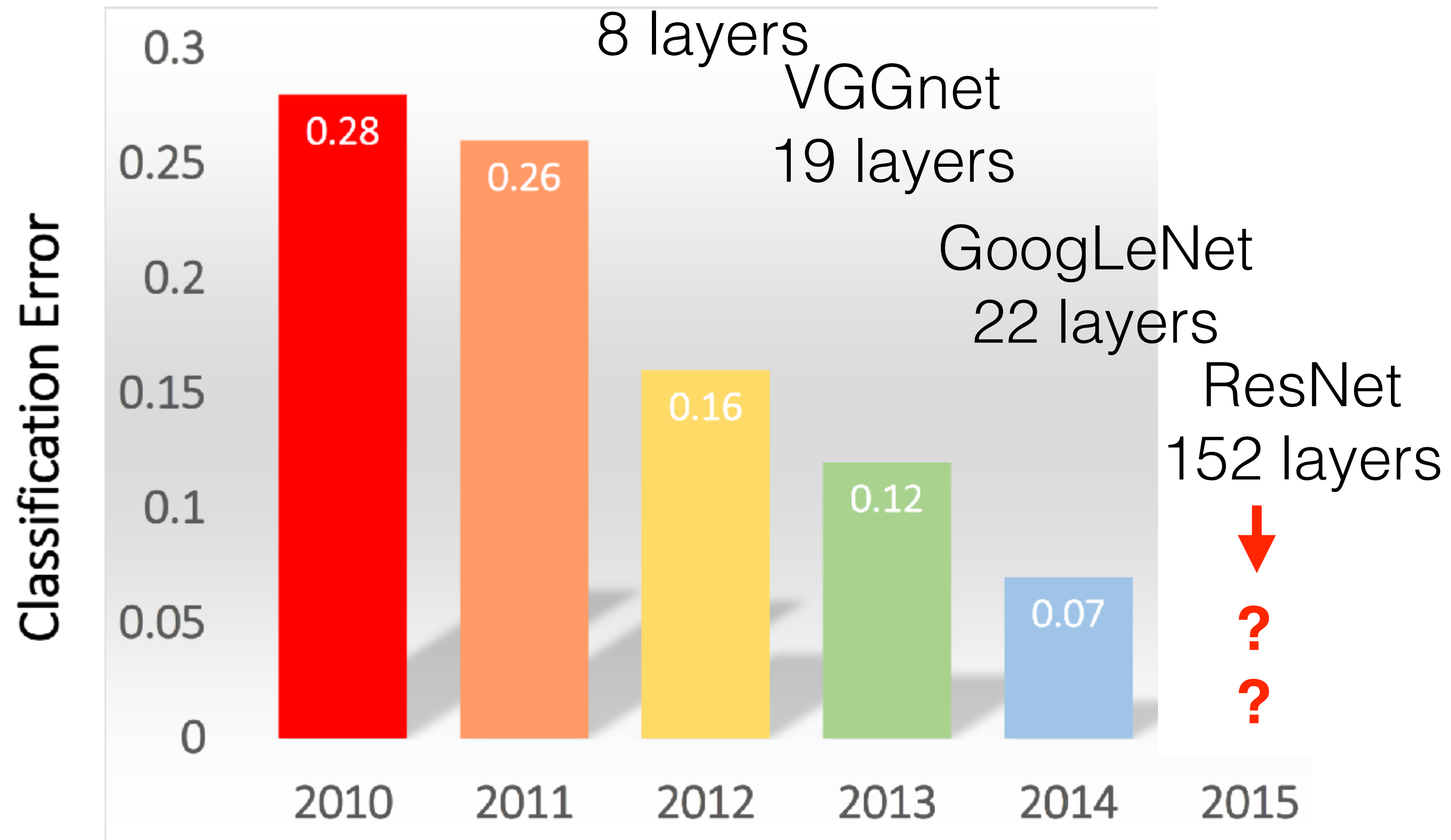
19 layers

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?**
- error was higher even in **training** => no overfitting, but vanishing gradient !

He et al. Going Deeper with Convolutions, CVPR, 2015

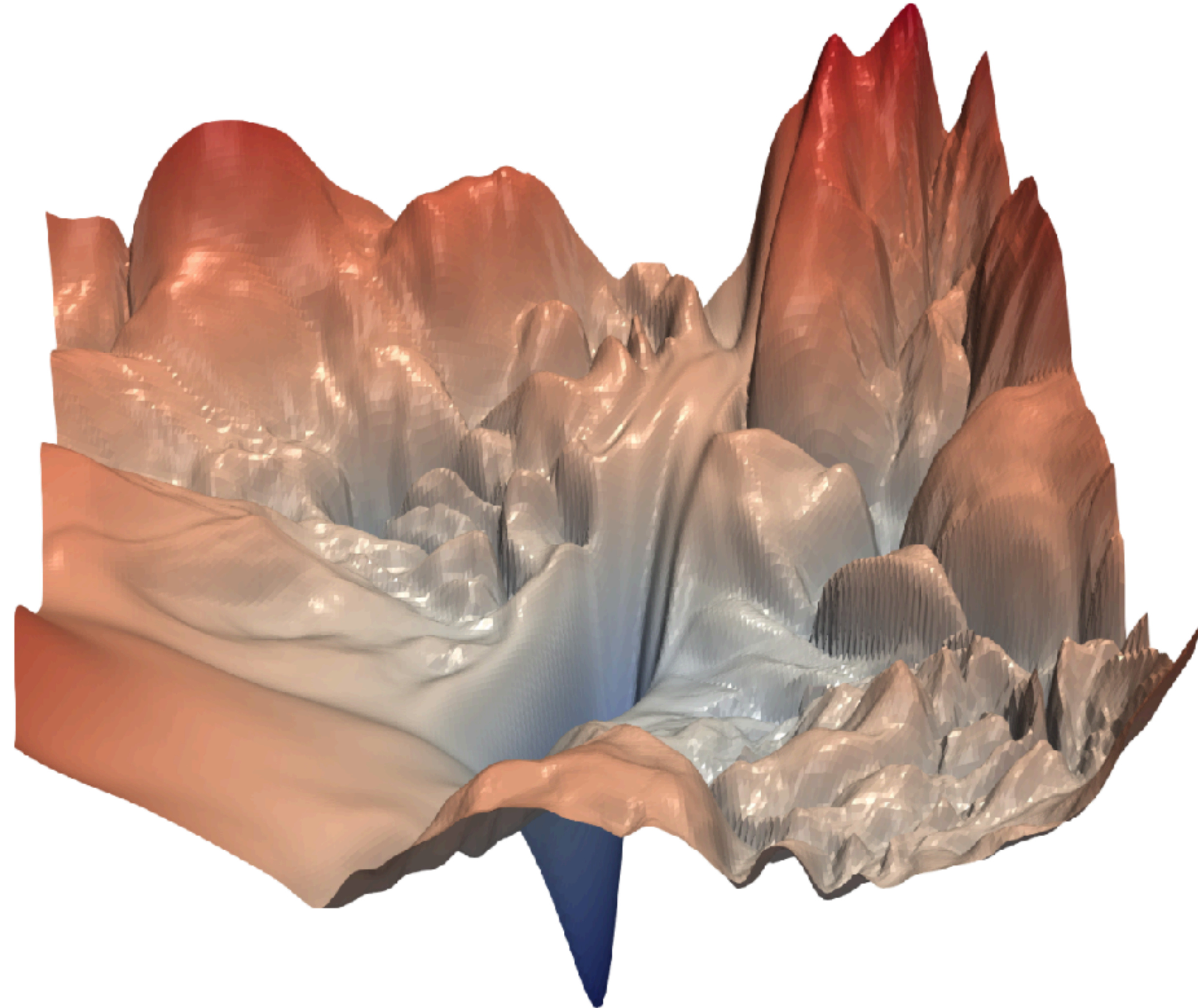
<https://arxiv.org/abs/1512.03385>

Visualizing Loss Landscape of Neural Nets

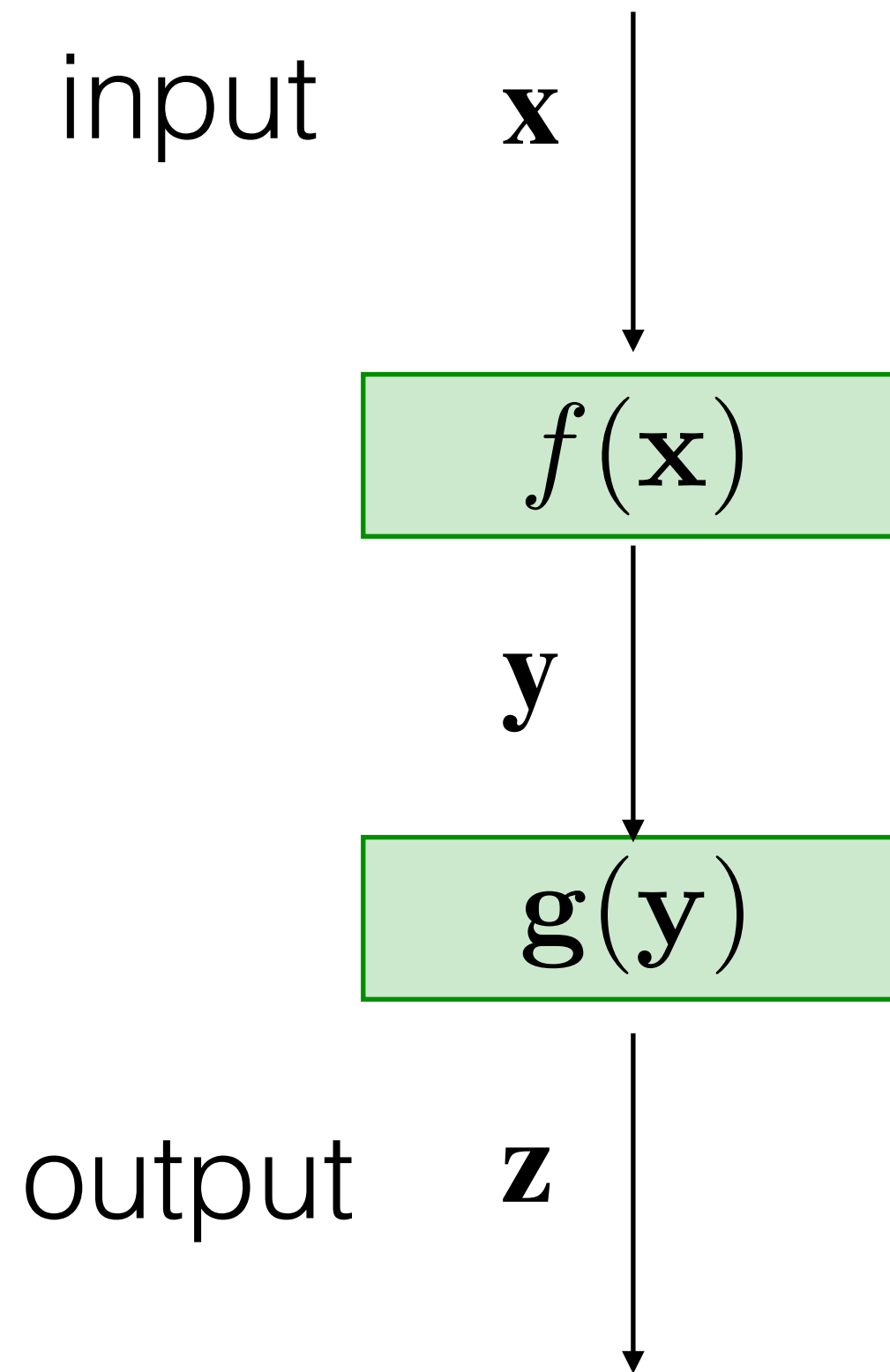
[Li et al, NIPS, 2018] <https://arxiv.org/pdf/1712.09913.pdf>

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

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



forward pass:



backward pass:

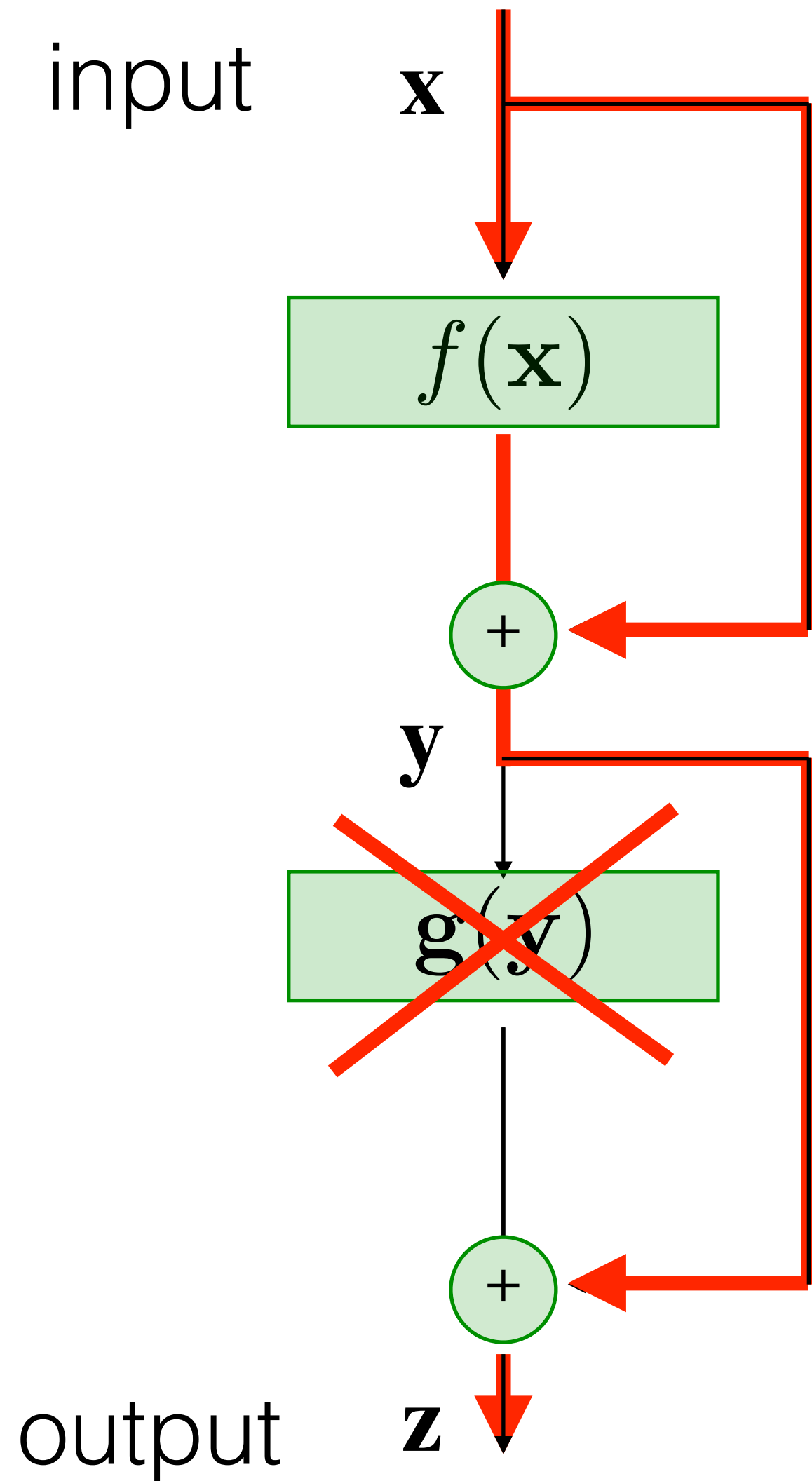
gradient $\frac{\partial \mathbf{z}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial \mathbf{z}}{\partial \mathbf{y}} & \frac{\partial \mathbf{z}}{\partial \mathbf{x}} \end{bmatrix} \approx \mathbf{0}$ then gradient is always zero

$\approx \mathbf{0}$ $\approx \mathbf{0}$

if any local gradient is zero

The diagram shows the backward pass gradient calculation. It features the equation $\frac{\partial \mathbf{z}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial \mathbf{z}}{\partial \mathbf{y}} & \frac{\partial \mathbf{z}}{\partial \mathbf{x}} \end{bmatrix} \approx \mathbf{0}$. The Jacobian matrix $\begin{bmatrix} \frac{\partial \mathbf{z}}{\partial \mathbf{y}} & \frac{\partial \mathbf{z}}{\partial \mathbf{x}} \end{bmatrix}$ is enclosed in a red rectangular box. Below the box, two red arrows point upwards towards the $\frac{\partial \mathbf{z}}{\partial \mathbf{y}}$ and $\frac{\partial \mathbf{z}}{\partial \mathbf{x}}$ terms, each accompanied by the text $\approx \mathbf{0}$. To the right of the main equation, the text "then gradient is always zero" is written. Below the entire equation block, the text "if any local gradient is zero" is written.

forward pass:



backward pass:

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

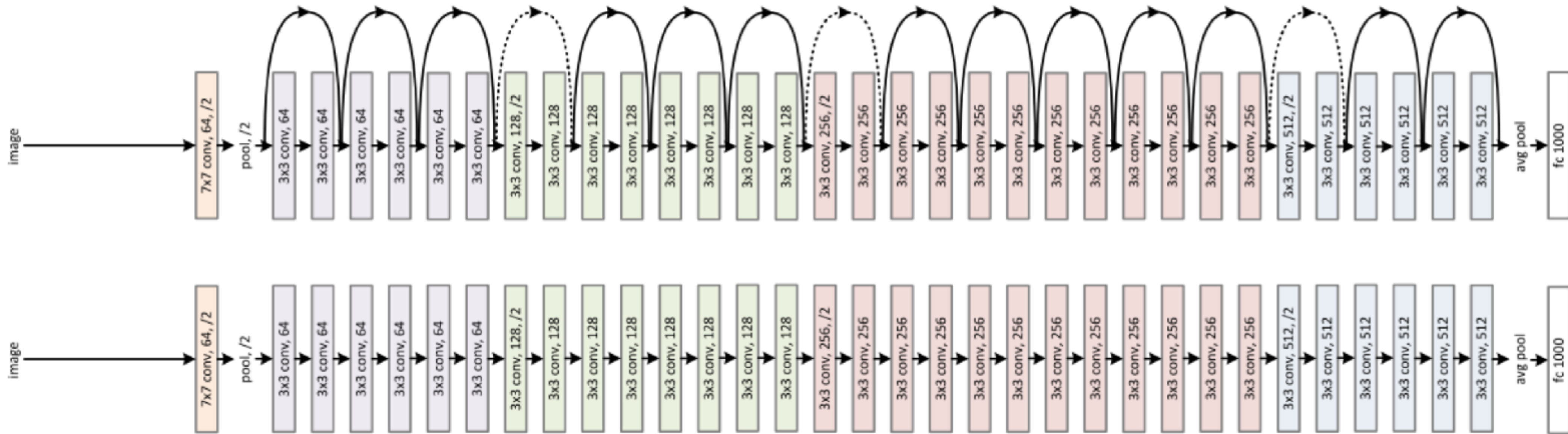
$\approx \mathbf{0}$ $\approx \mathbf{0}$

if any local gradient is zero

the gradient can still flow through another path

ResNet: deep ConvNet with skip connections

ResNet

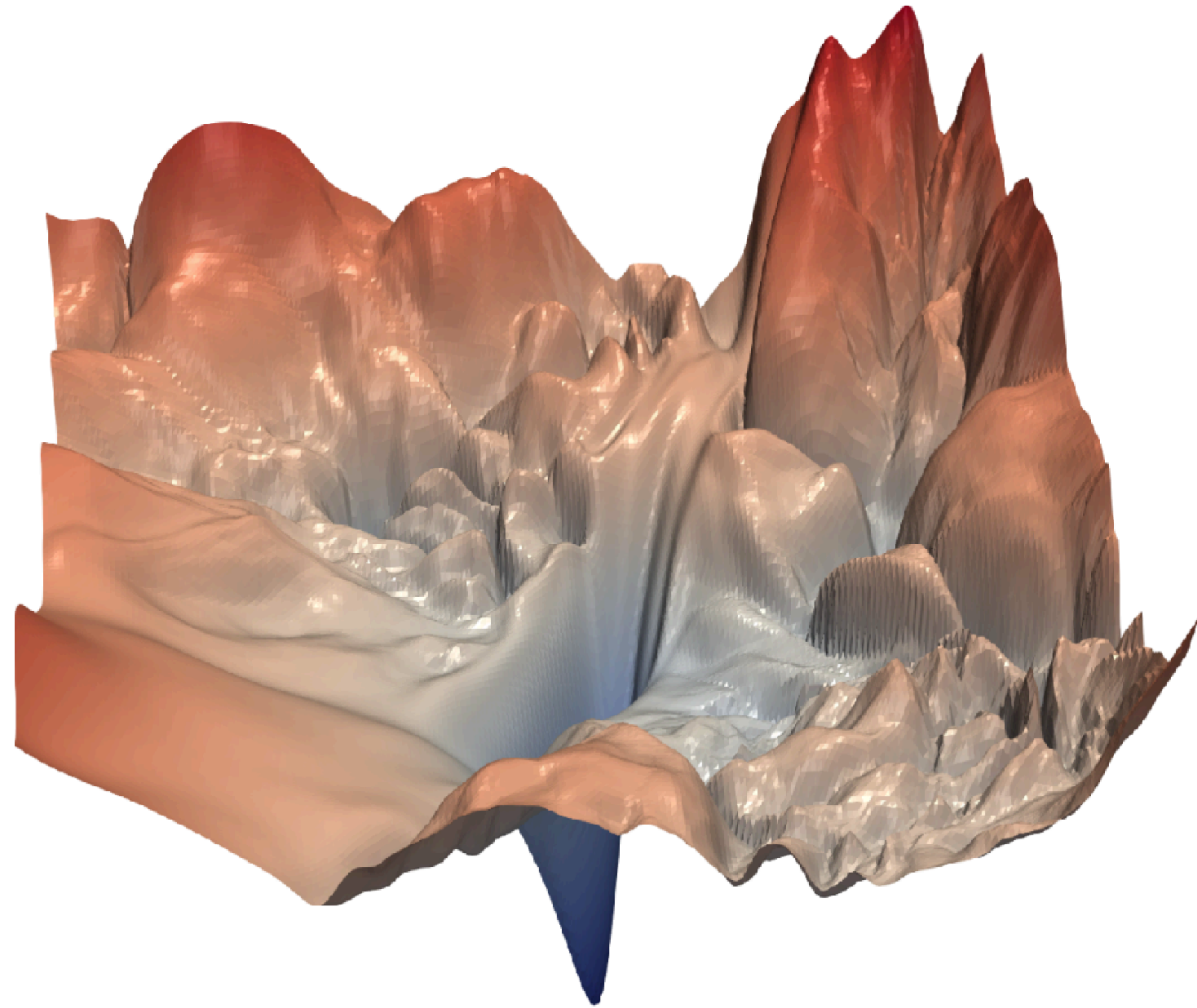


ResNet-NS



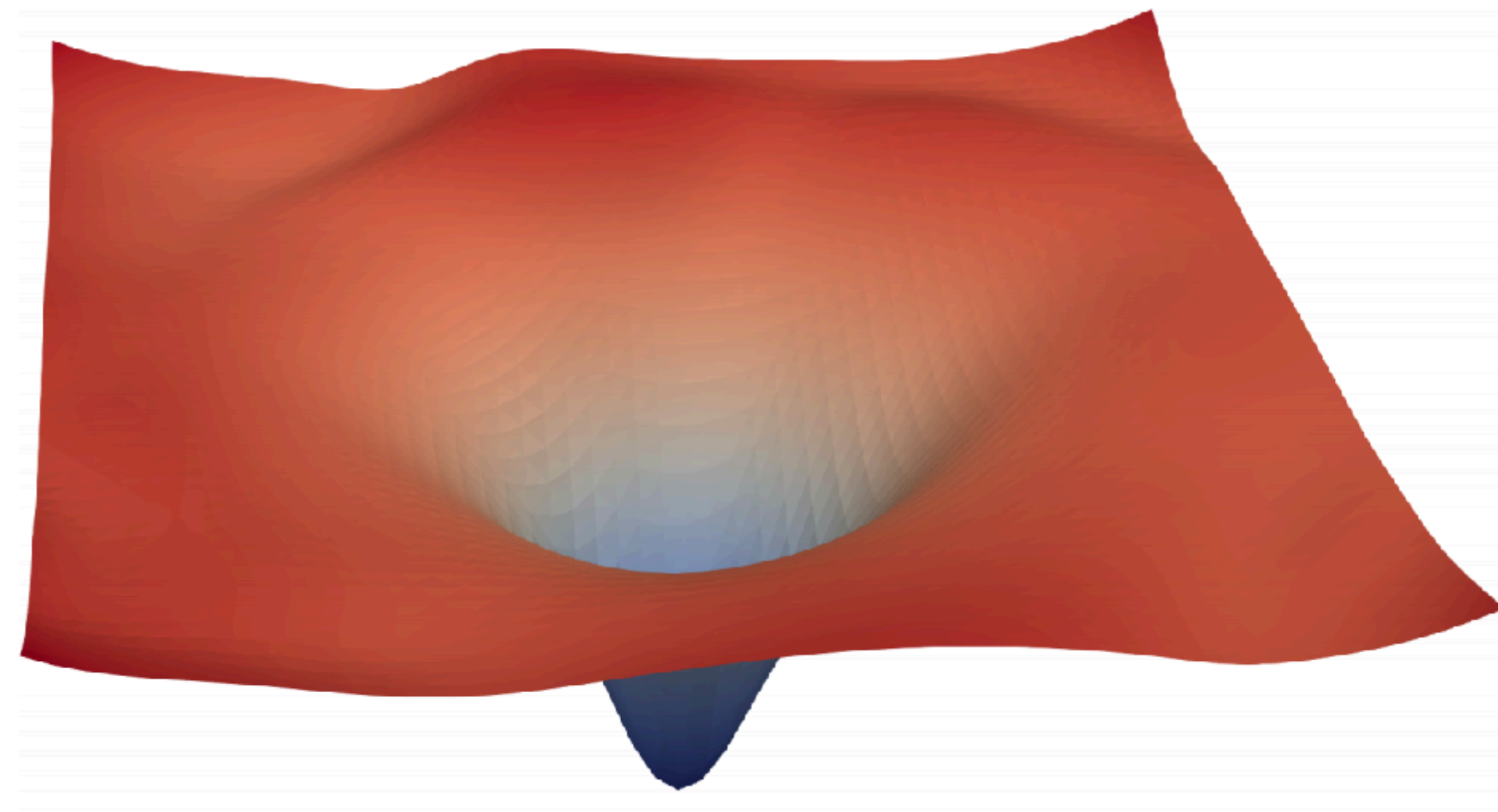
Visualizing Loss Landscape of Neural Nets

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



(a) without skip connections

ResNet-NS

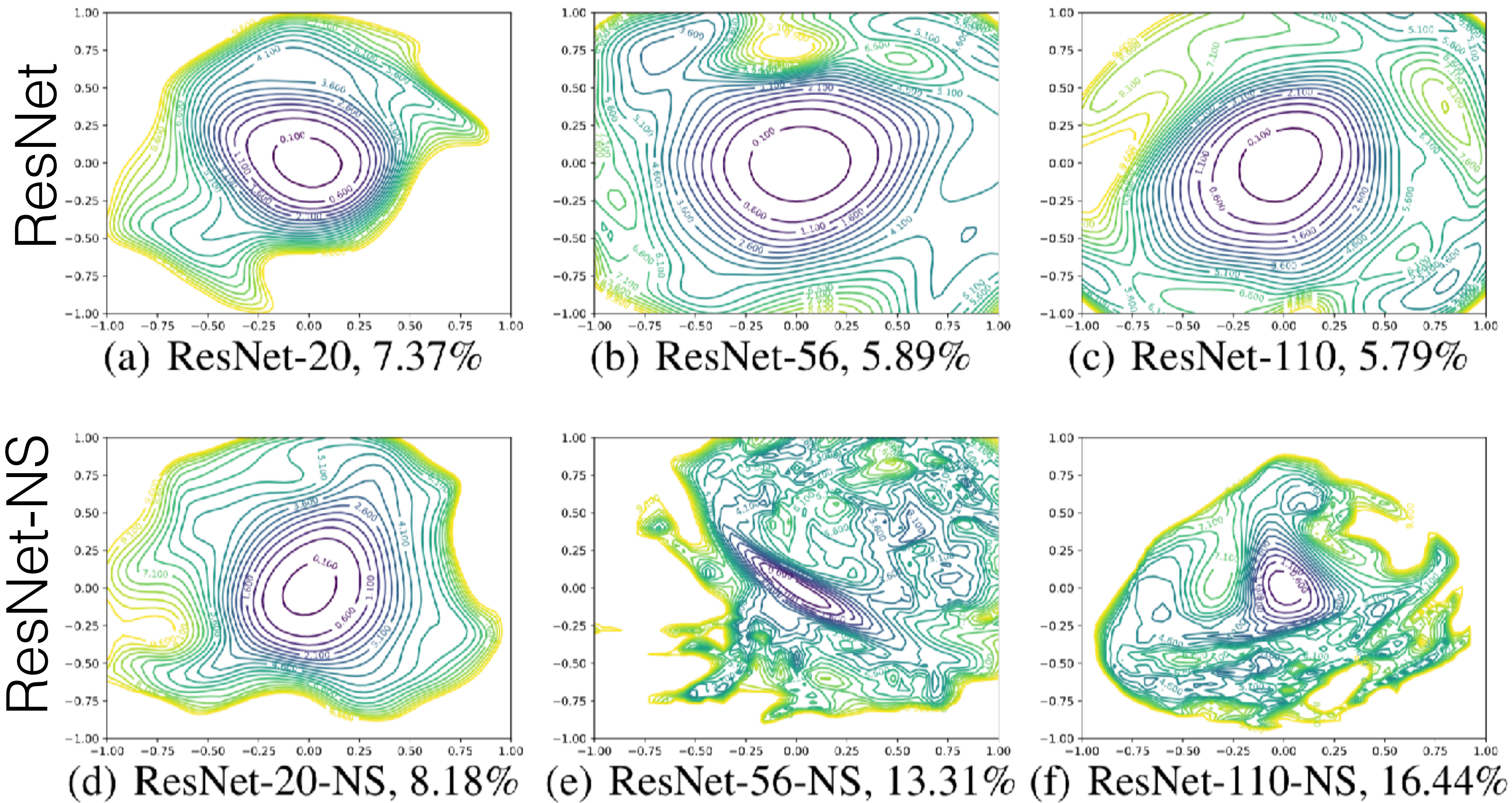


(b) with skip connections

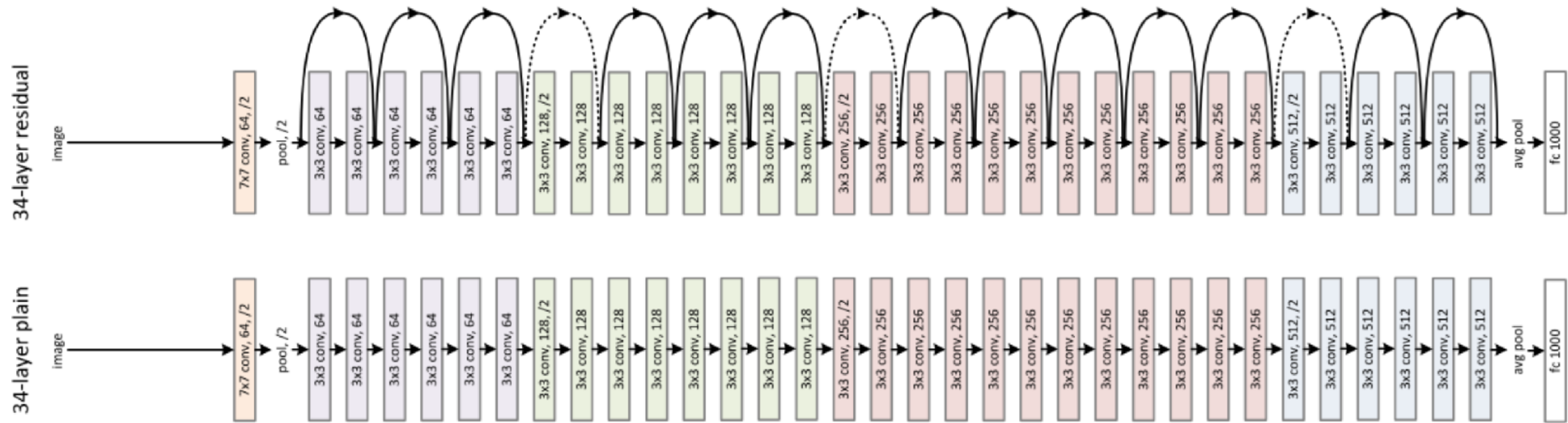
ResNet

ResNet: deep ConvNet with skip connections

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



ResNet: deep ConvNet with skip connections



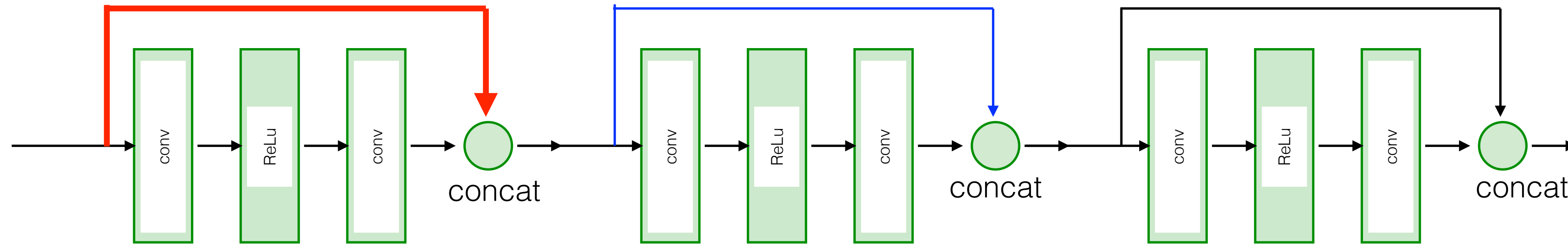
- 1k+ layers possible
- Initialization with zero weights is meaningful
- Better gradient flow (many independent paths, opt-friendly landscape)
- Robustness wrt noise and layer removal

<https://www.kaggle.com/keras/resnet50/home>

He et al. Going Deeper with Convolutions, CVPR, 2015

<https://arxiv.org/abs/1512.03385>

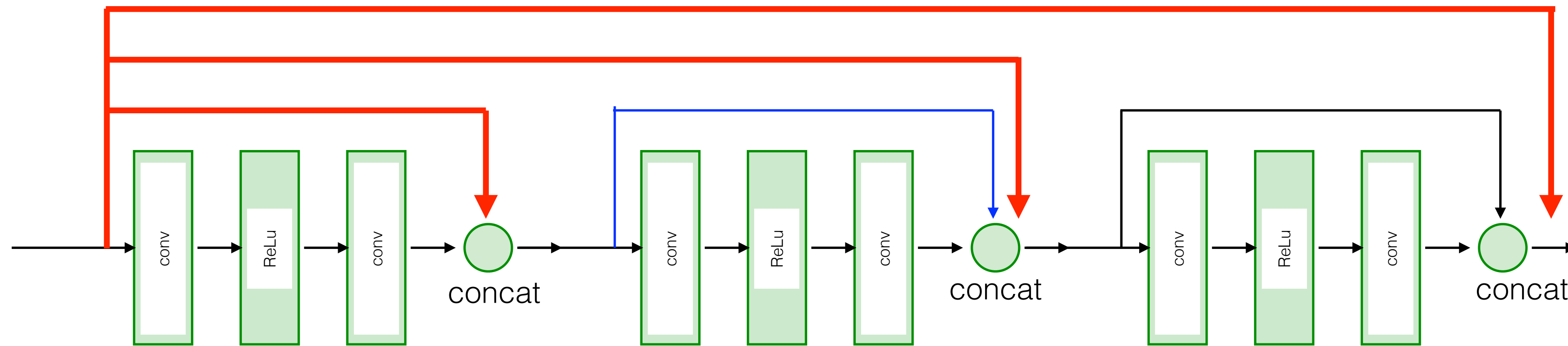
ResNet=>DenseNet



- 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

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

Classification results

AlexNet

8 layers

VGGnet

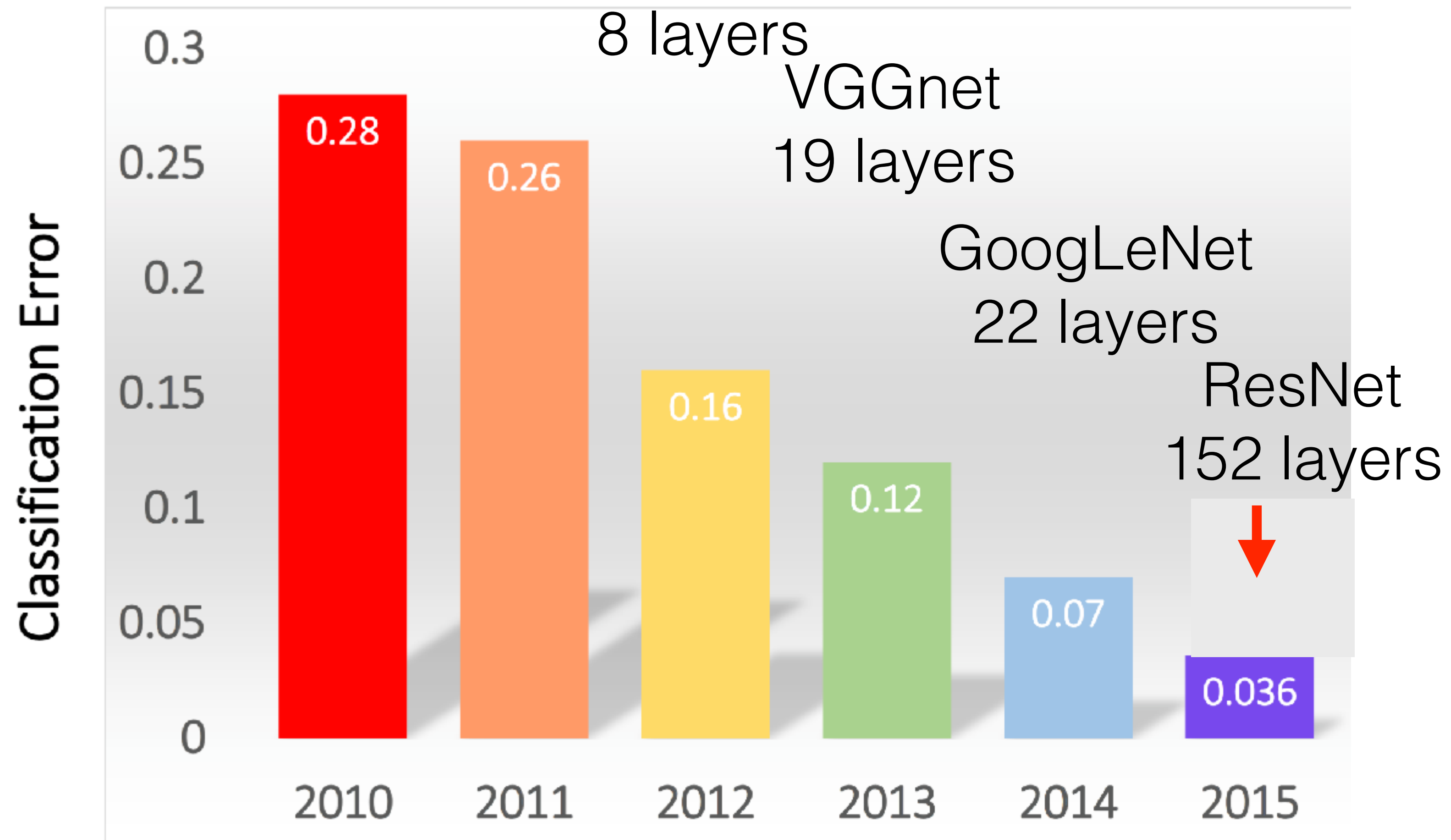
19 layers

GoogLeNet

22 layers

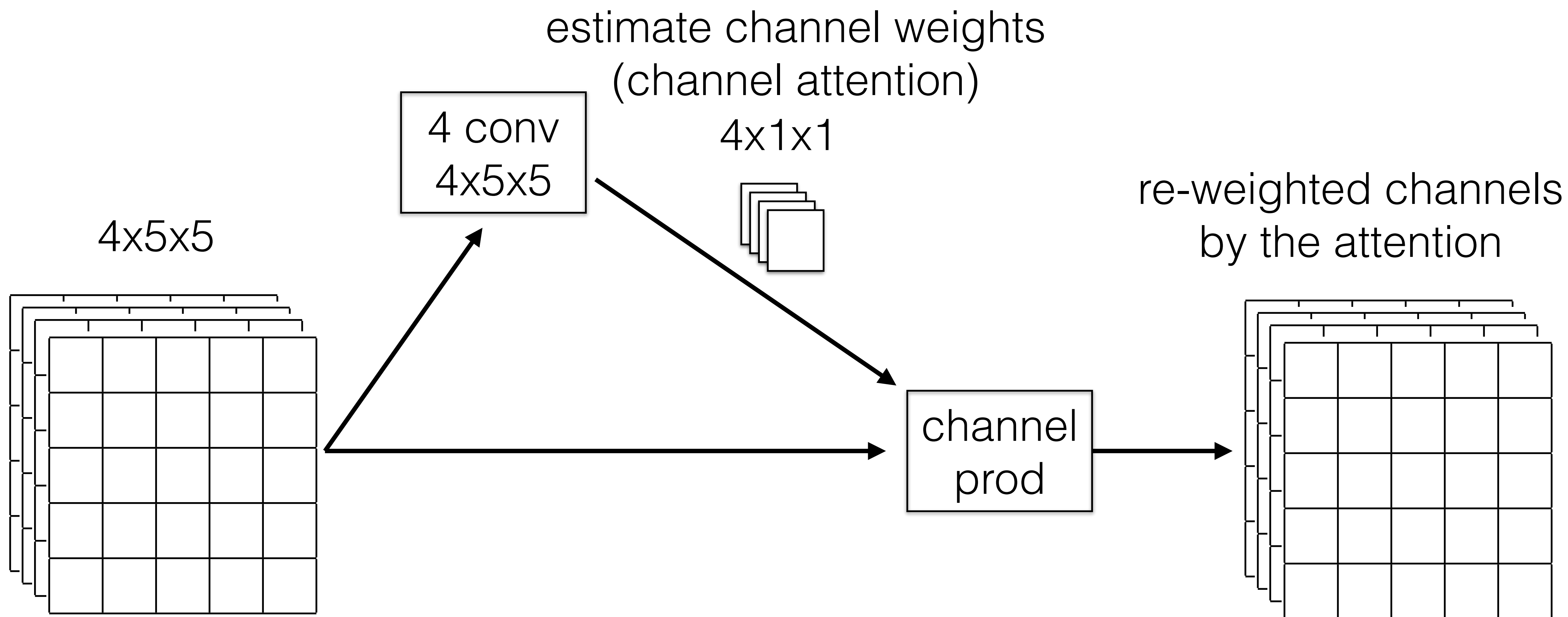
ResNet

152 layers



Human error around 5%

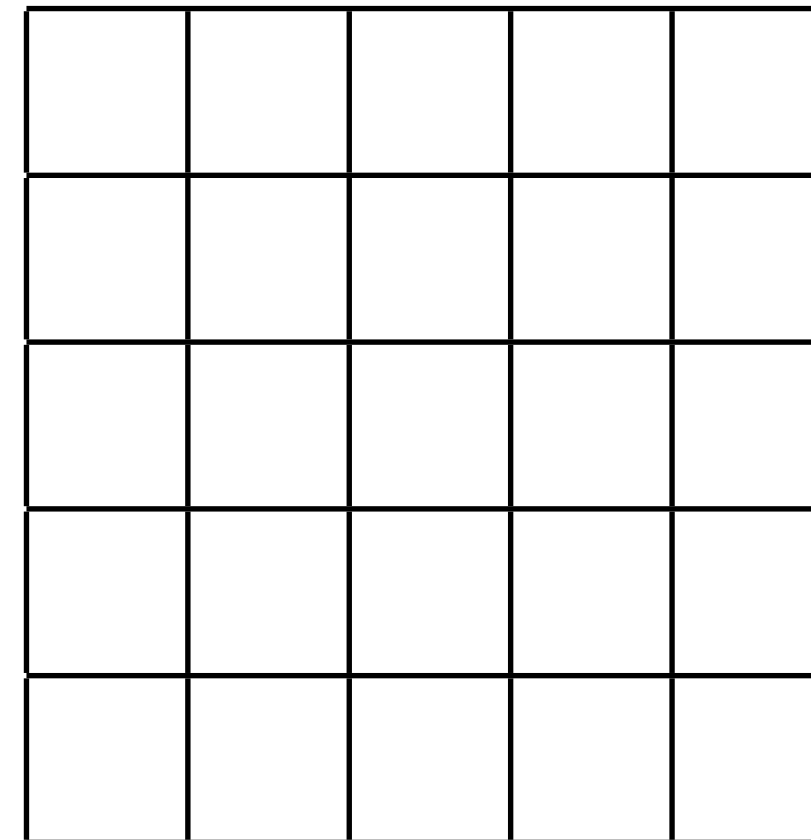
Channel attention



Spatial attention

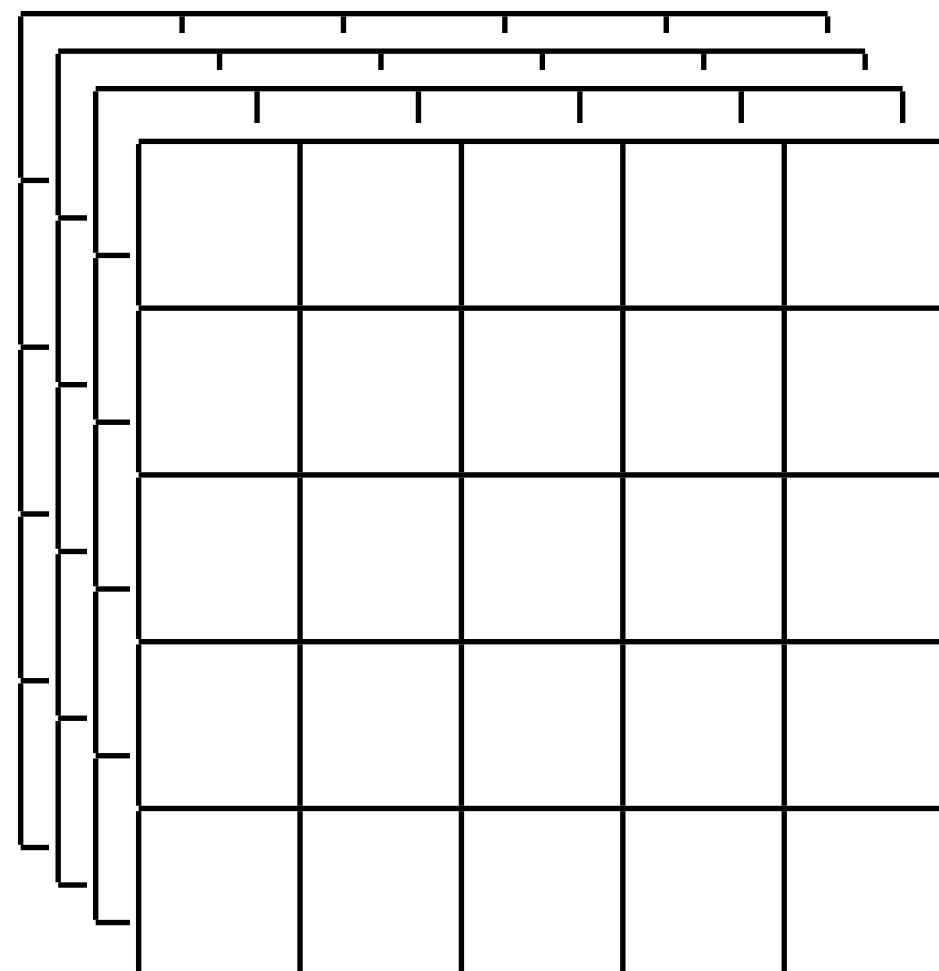
1x5x5

estimate pixel weights
(spatial attention)



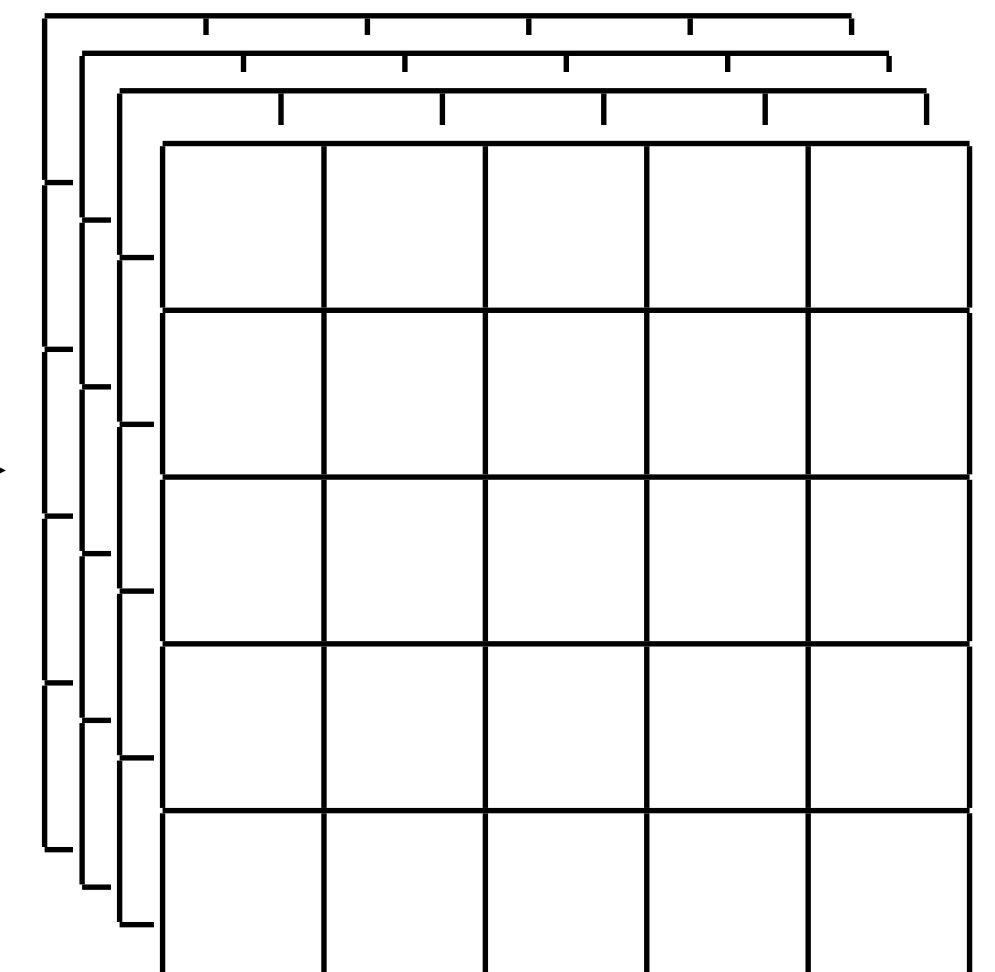
1 conv
4x1x1

4x5x5



re-weighted pixels
by spatial attention

pixel
prod



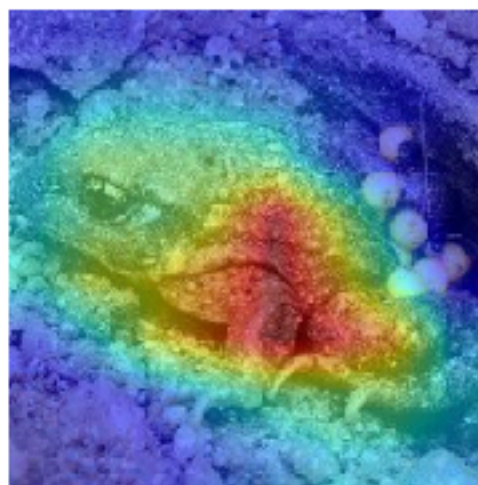
The attention map of classified object is better localized

Tailed frog

Input image

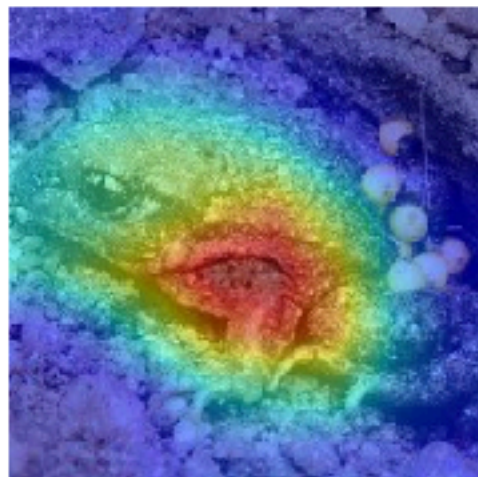


ResNet50



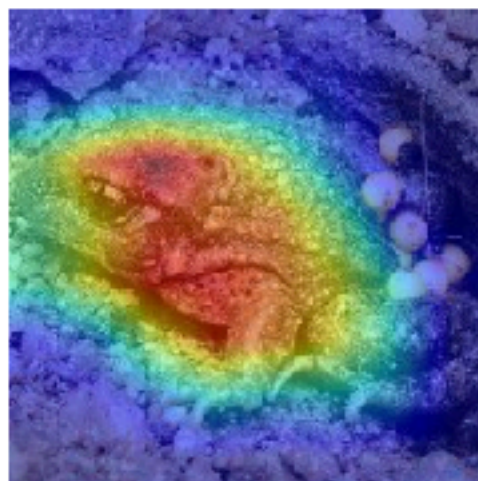
P=0.80736

ResNet50 + SE



P=0.87240

ResNet50 + CBAM



P = 0.96340

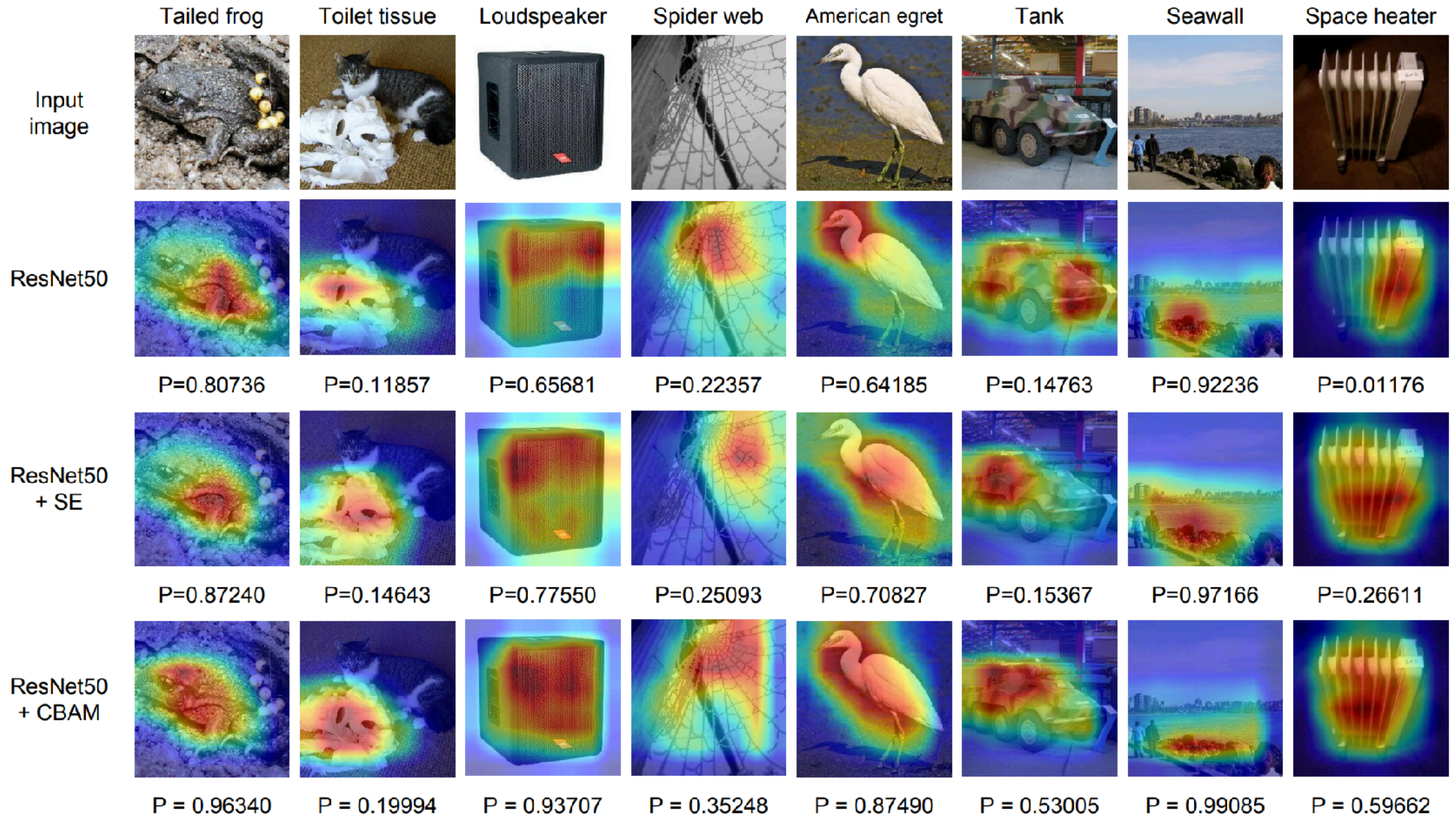
w/o attention

with channel attention

with spatial attention

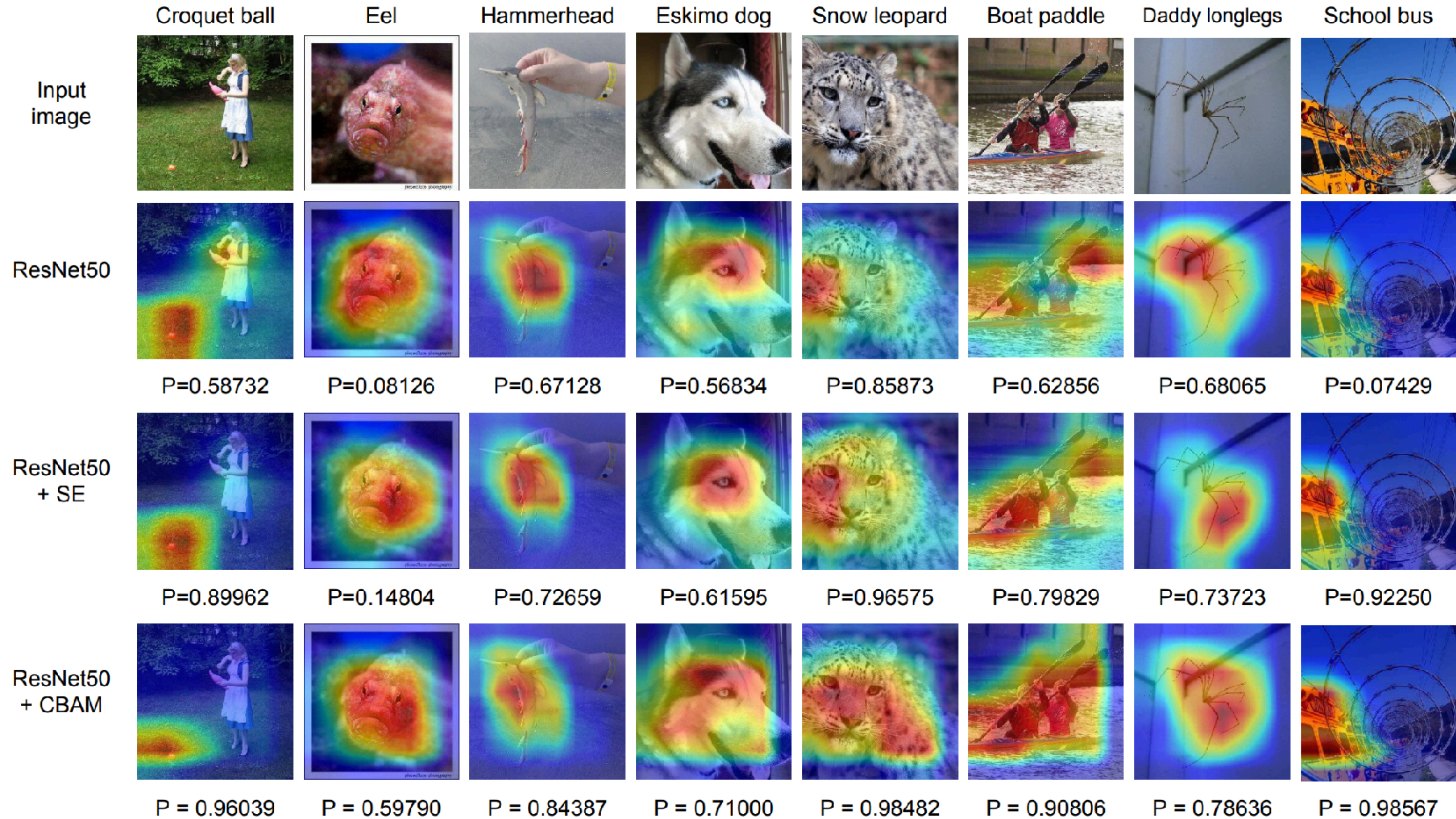
Attention modules [Woo et al, ECCV, 2018]

<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

VGGnet

19 layers

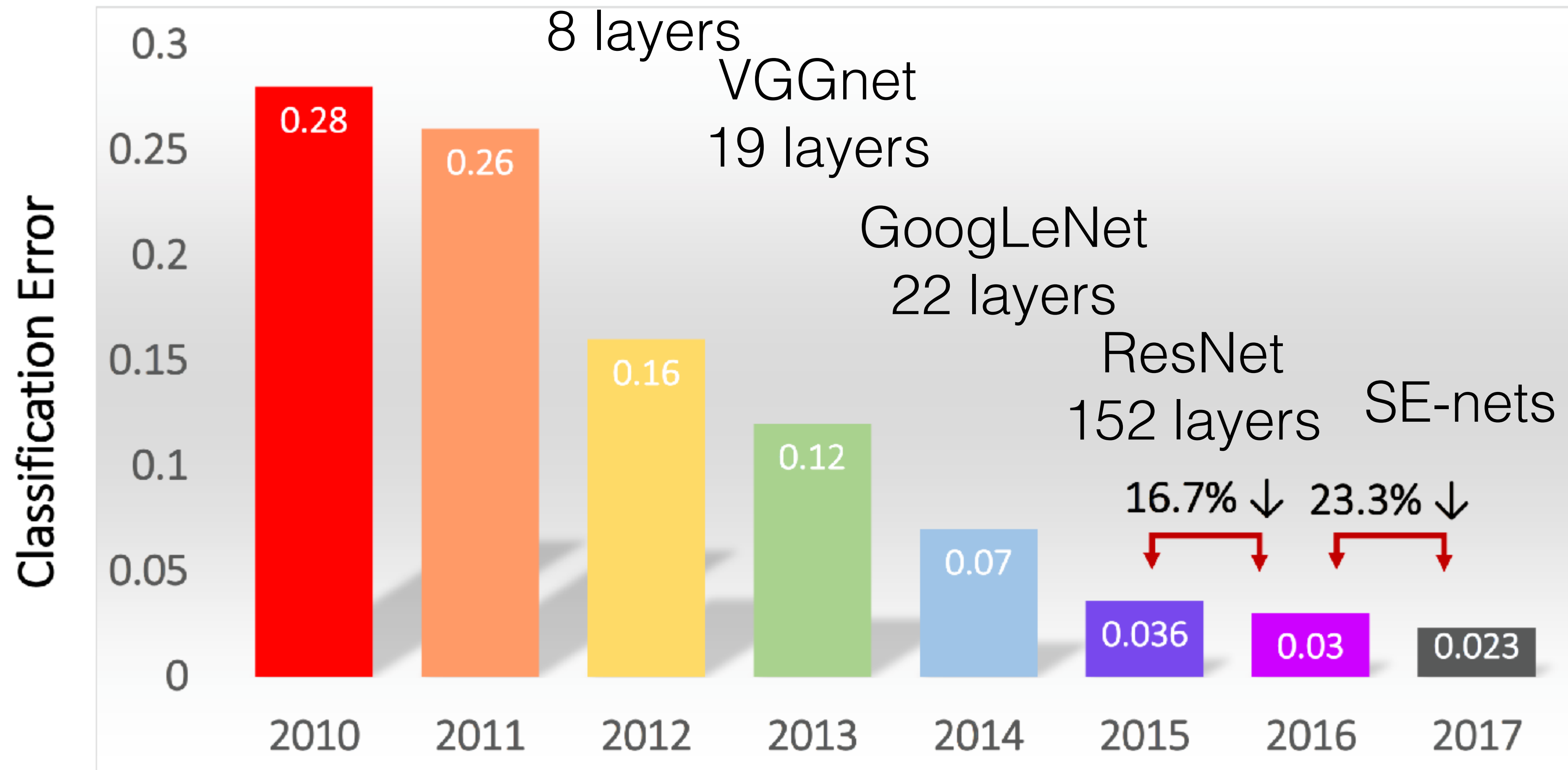
GoogLeNet

22 layers

ResNet

152 layers

SE-nets



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

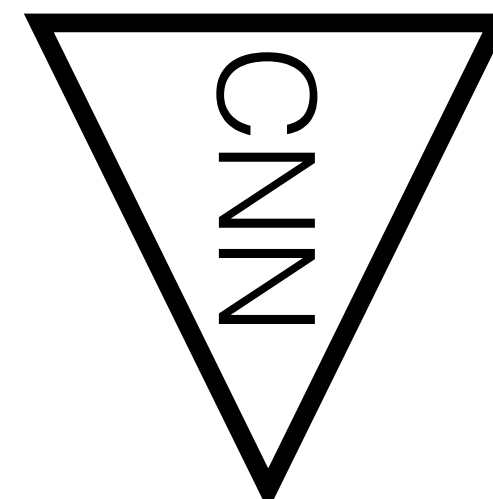
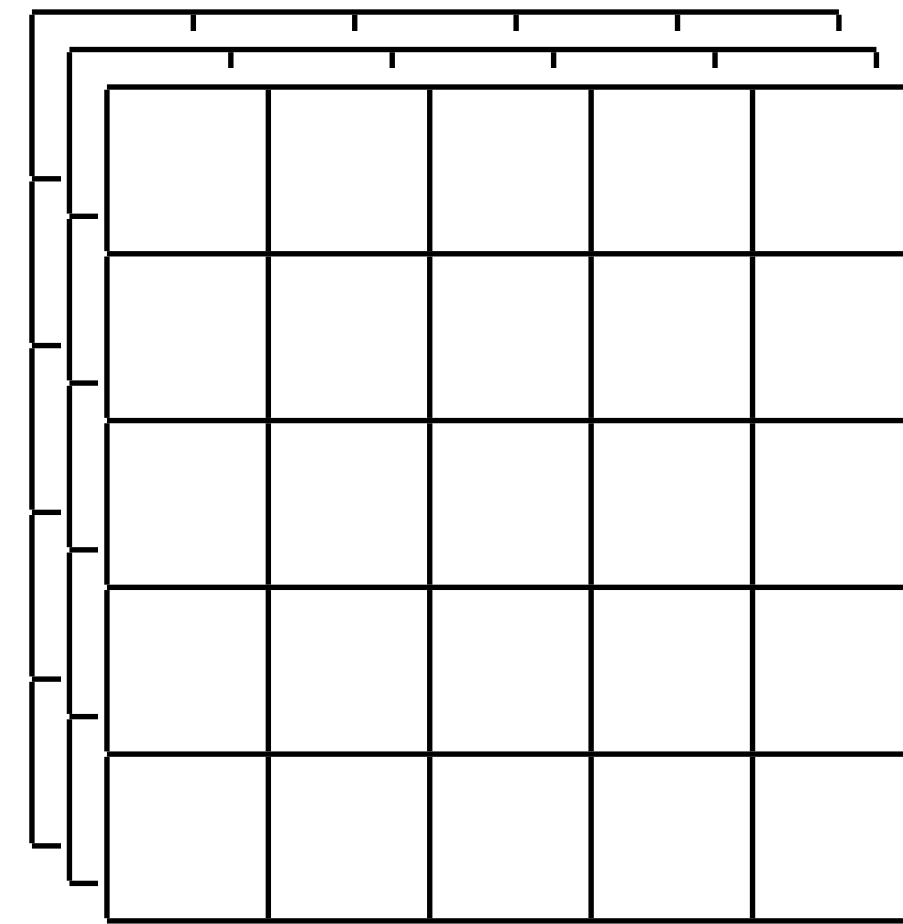
Semantic segmentation



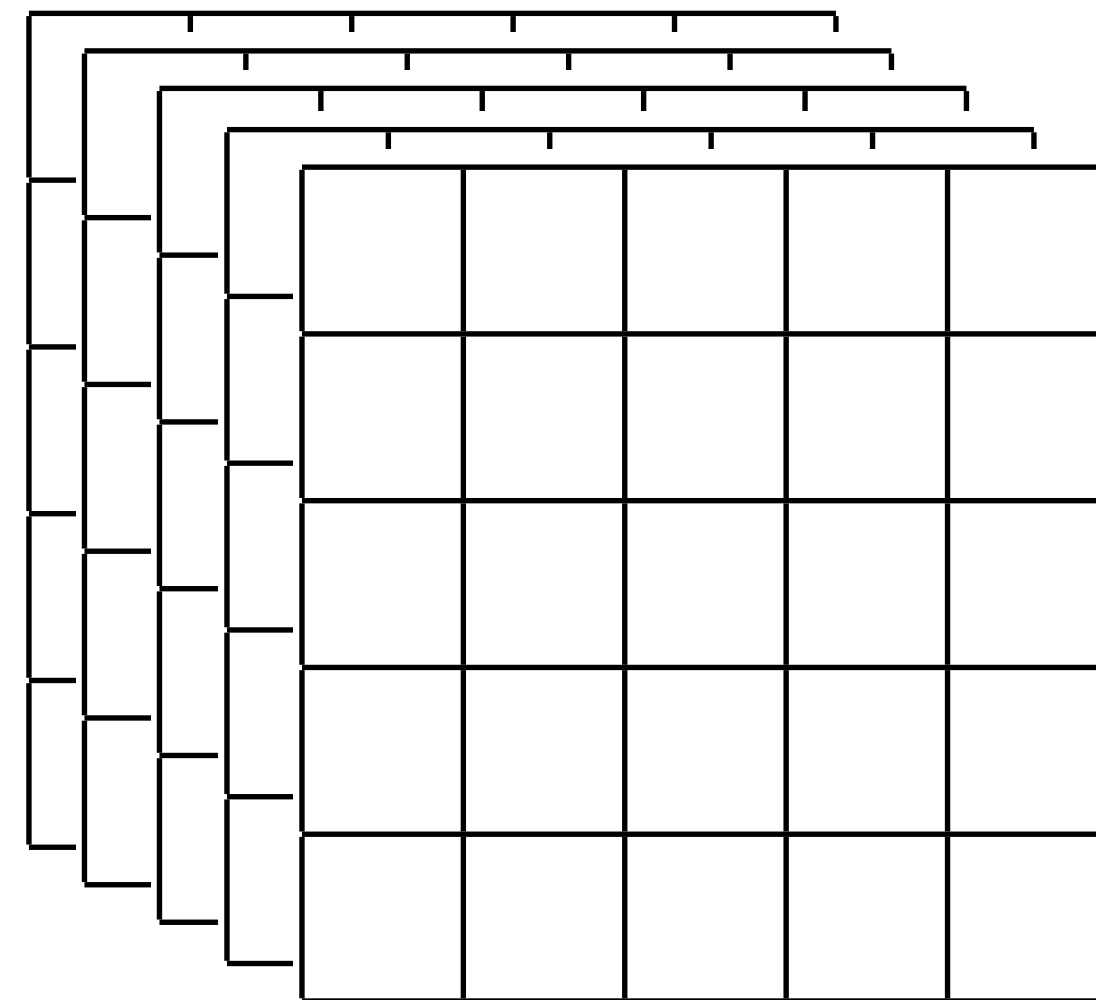
- road
- sideway
- pedestrian
- traffic sign
- trees
- sky

Semantic segmentation

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



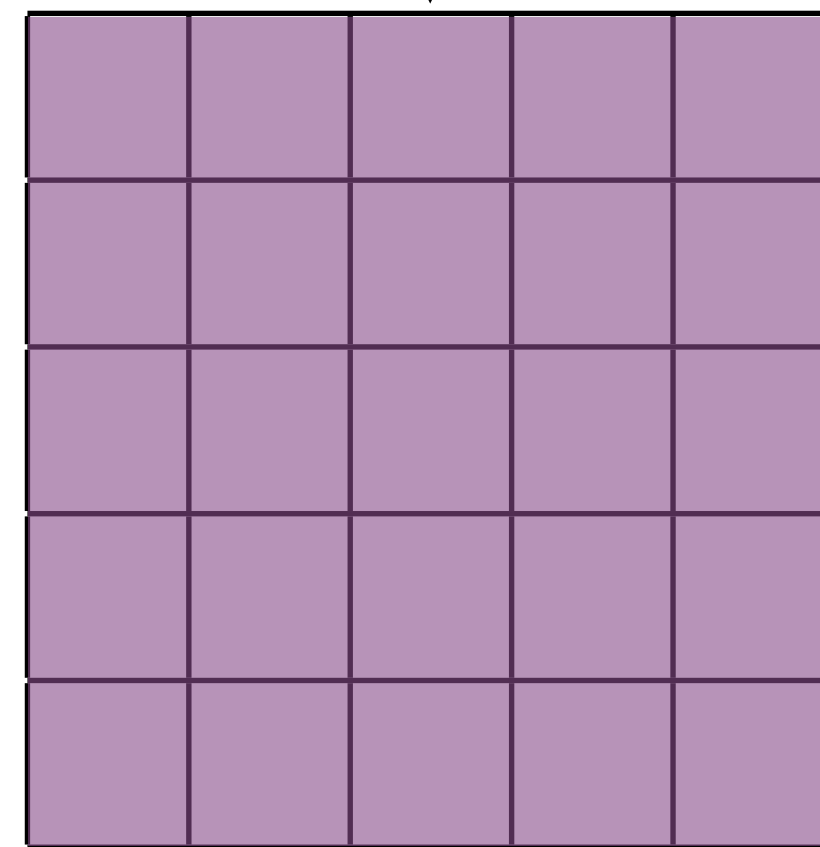
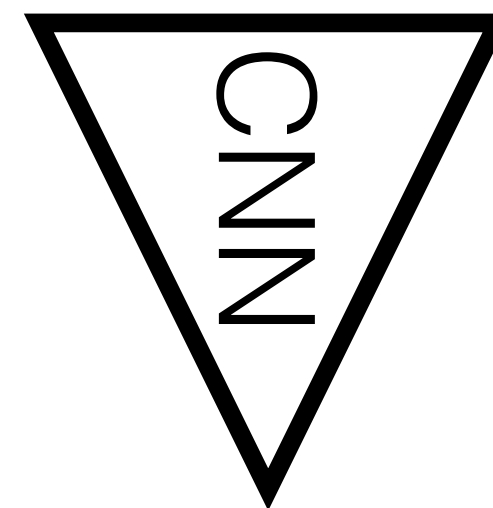
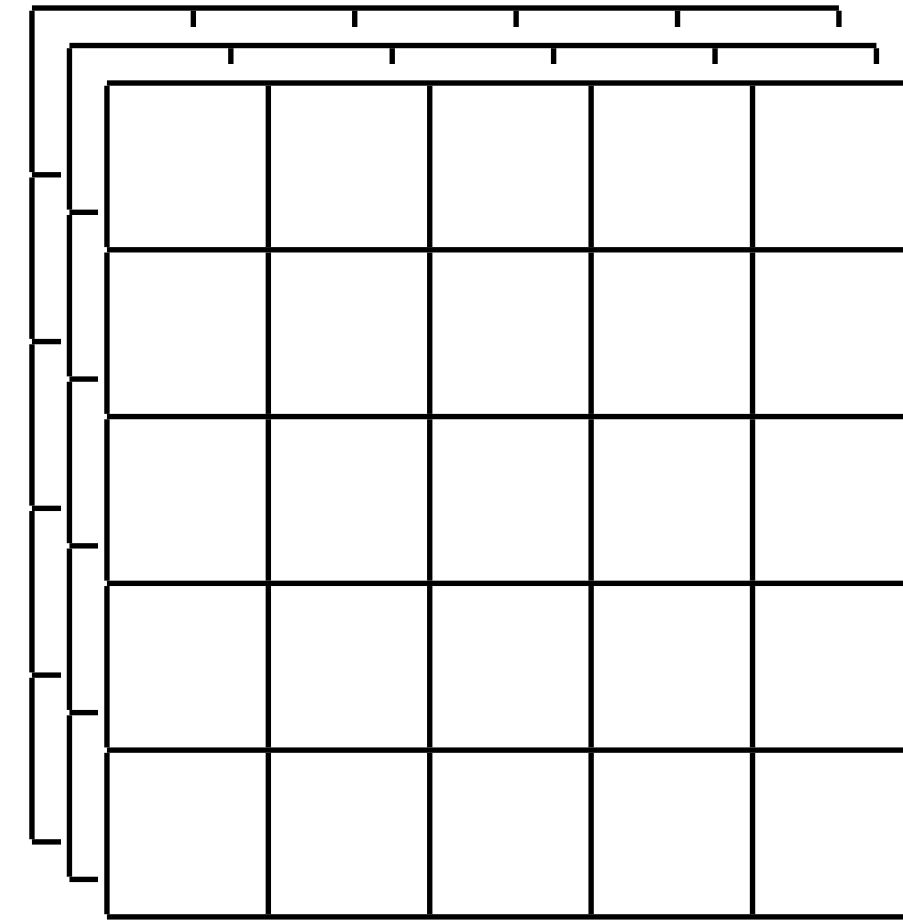
pixel-wise probs
($H \times W \times N$)



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

Semantic segmentation

RGB image
(HxWx3)



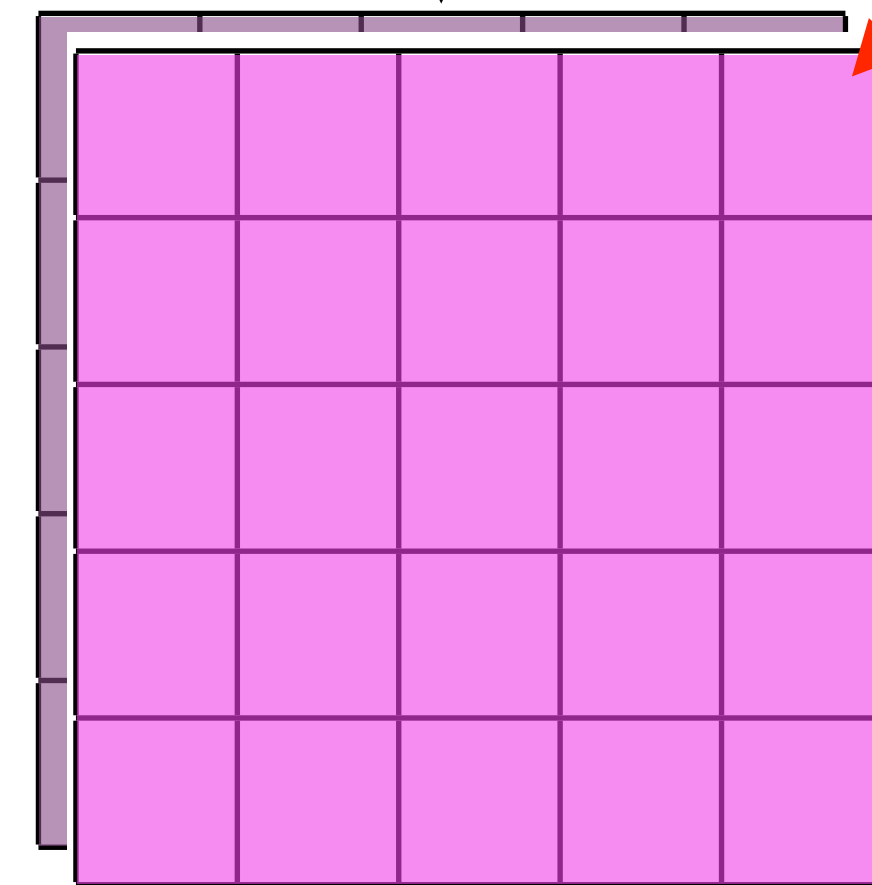
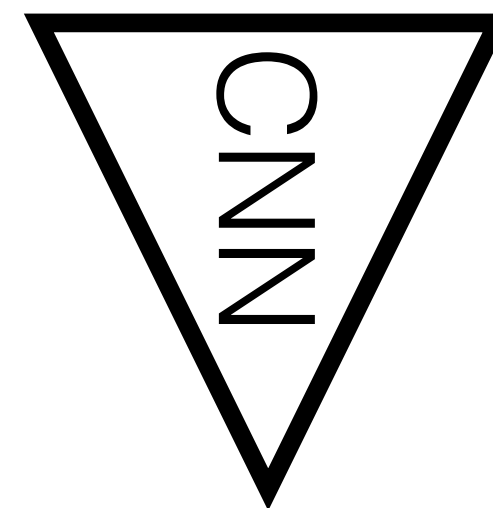
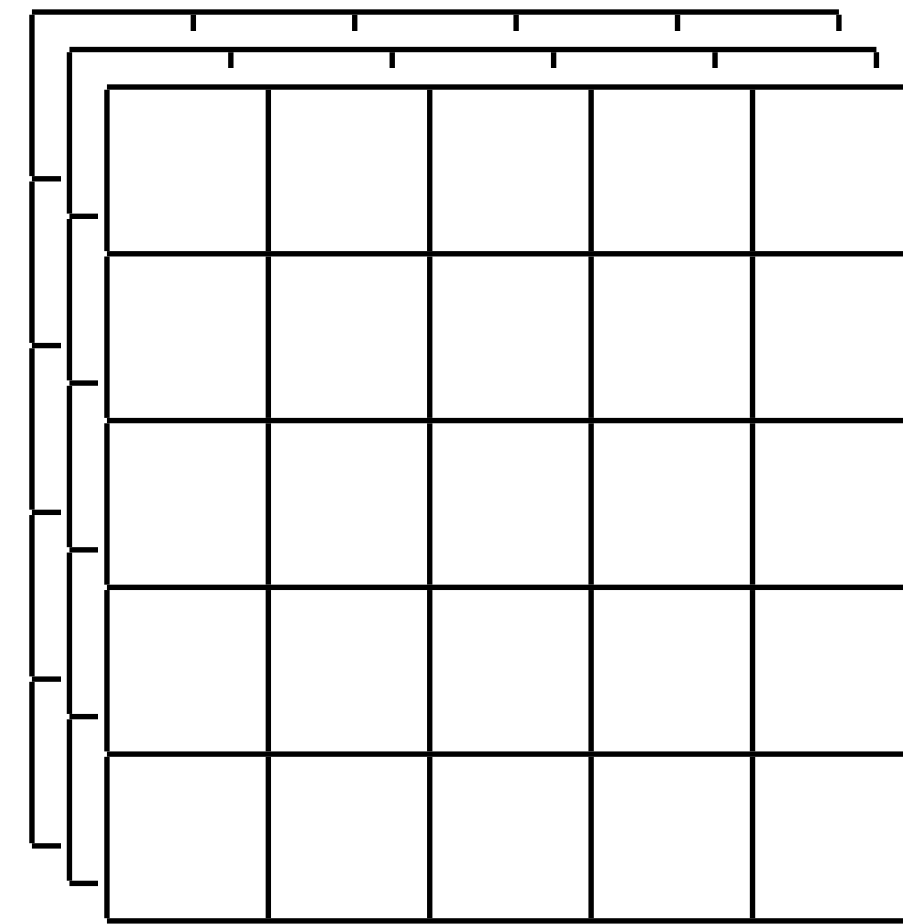
- road
- sideway
- pedestrian
- traffic sign
- trees
- sky

pixel-wise probability
of being **road**

channel 1

Semantic segmentation

RGB image
(HxWx3)



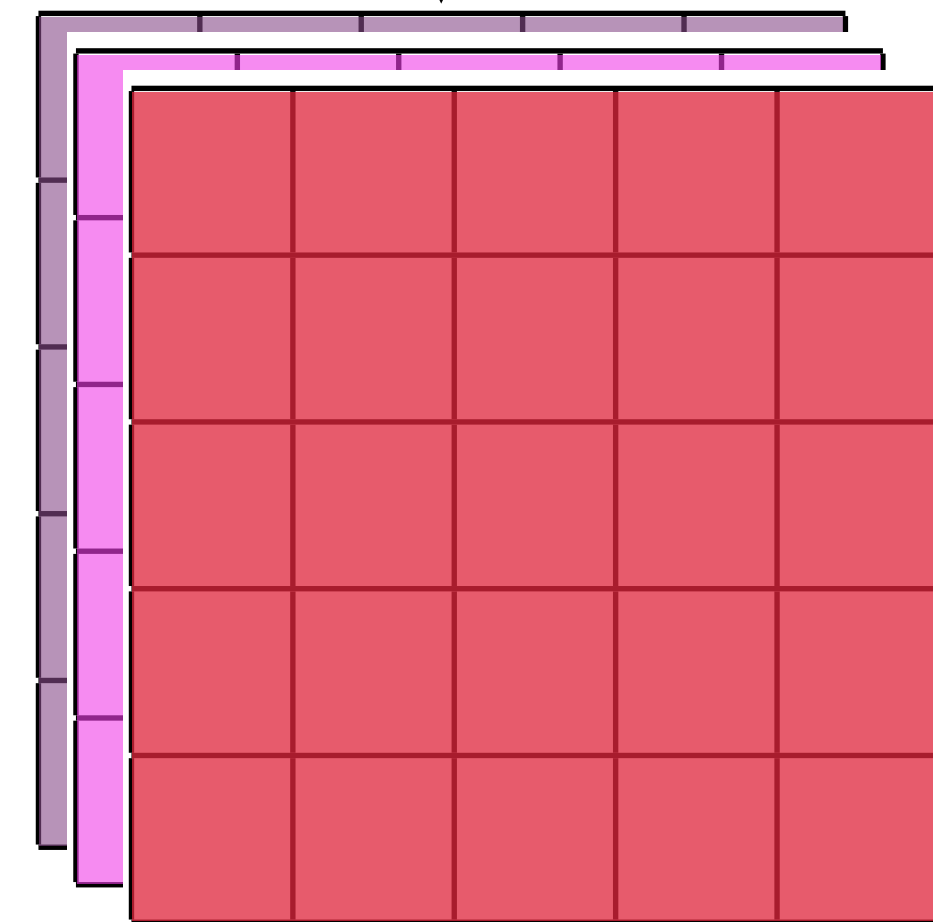
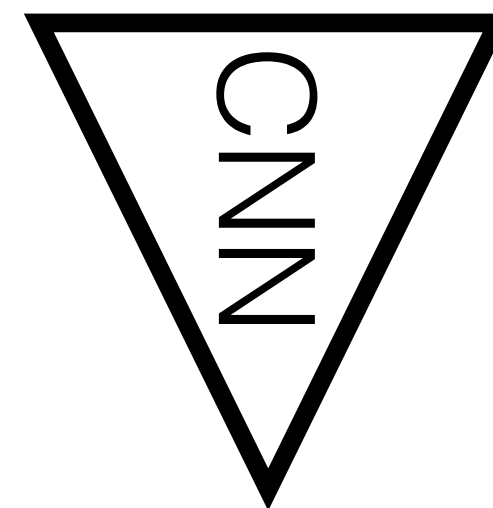
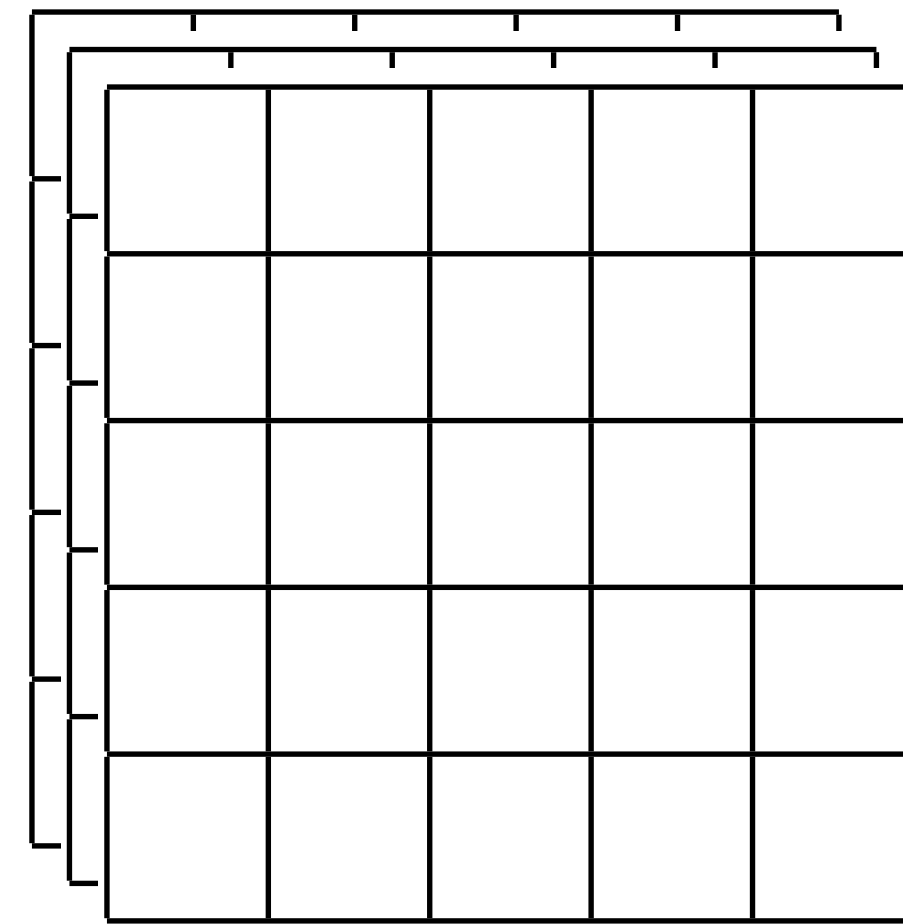
- road
- sideway
- pedestrian
- traffic sign
- trees
- sky

pixel-wise probability
of being **sideway**

channel 2

Semantic segmentation

RGB image
(HxWx3)

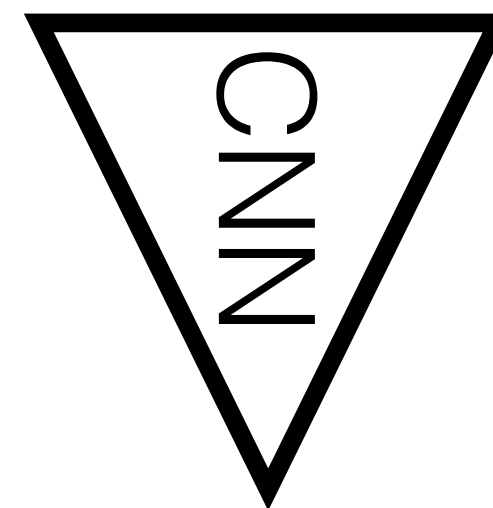
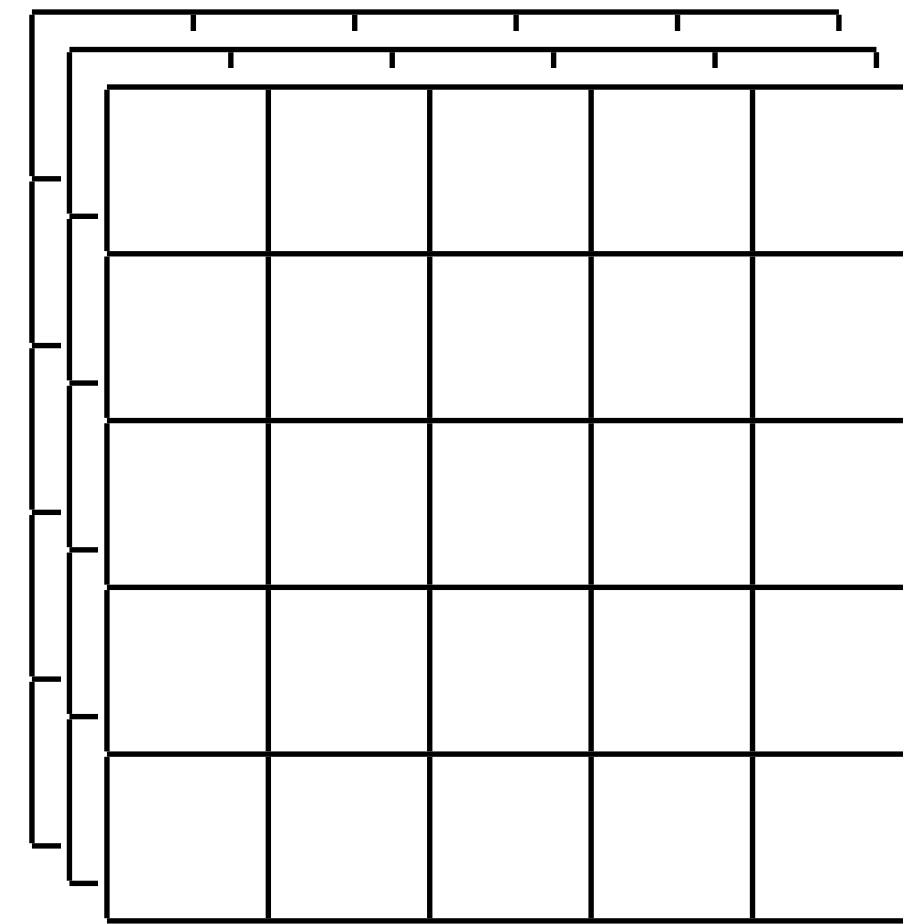


- road
- sideway
- pedestrian
- traffic sign
- trees
- sky

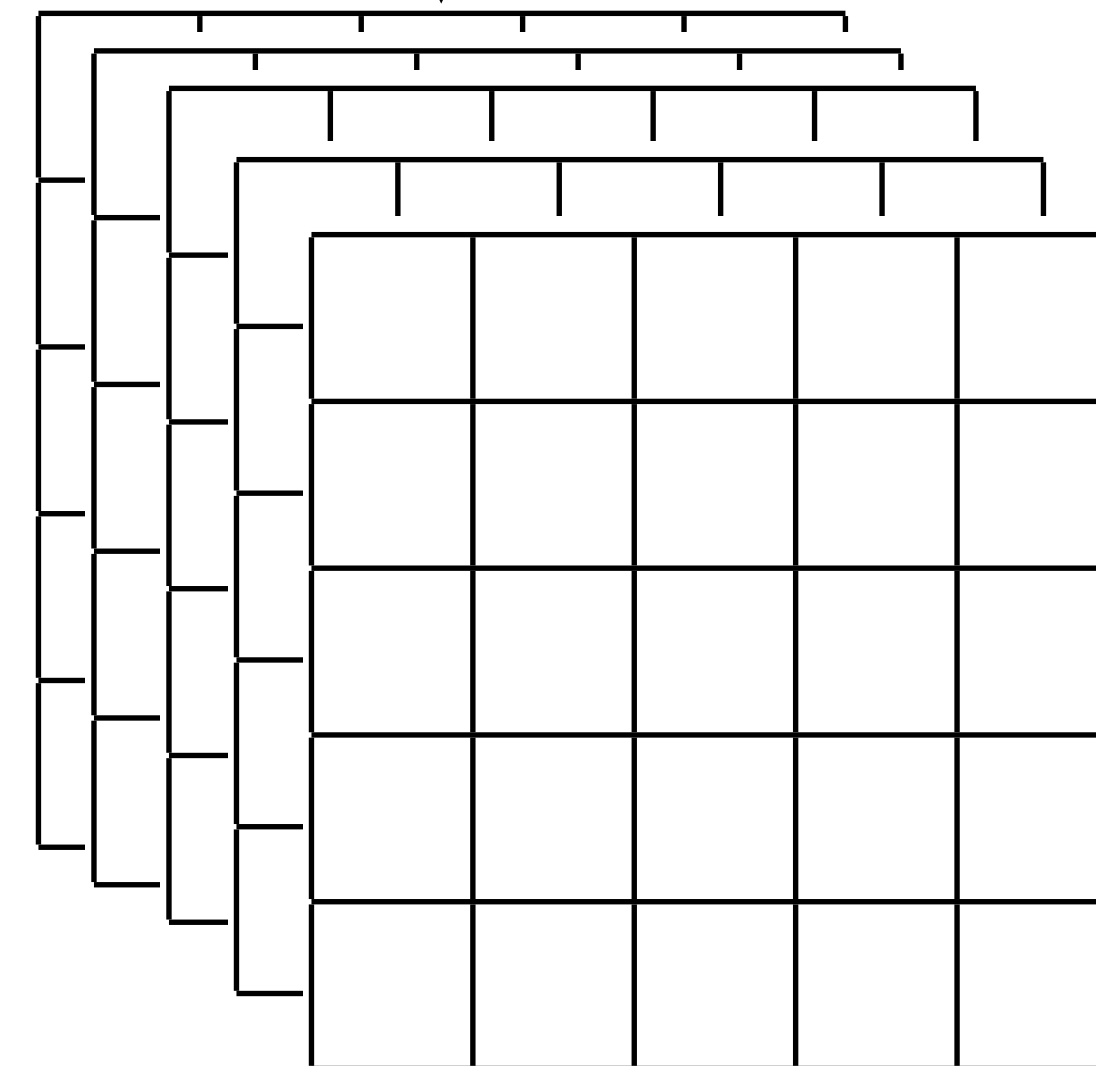
pixel-wise probability
of being **pedestrian**
channel 3

Semantic segmentation

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



pixel-wise probs
($H \times W \times N$)

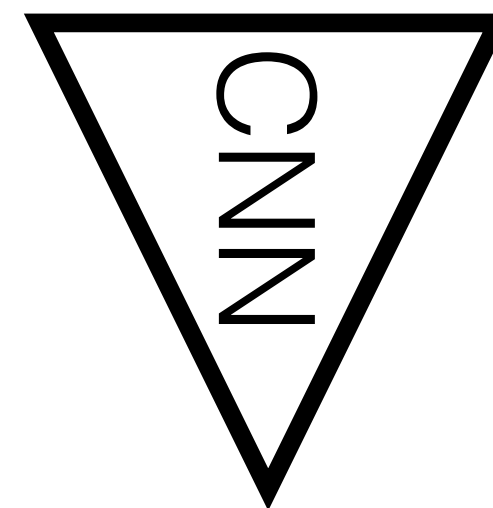
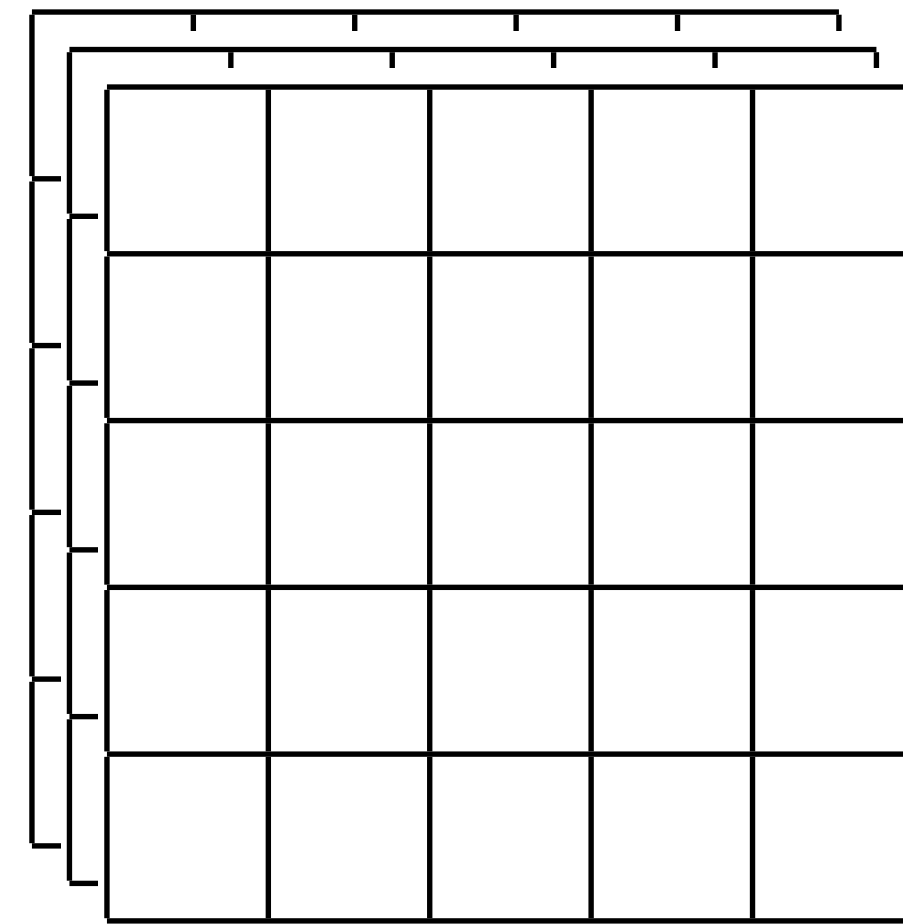


- road
- sideway
- pedestrian
- traffic sign
- trees
- sky

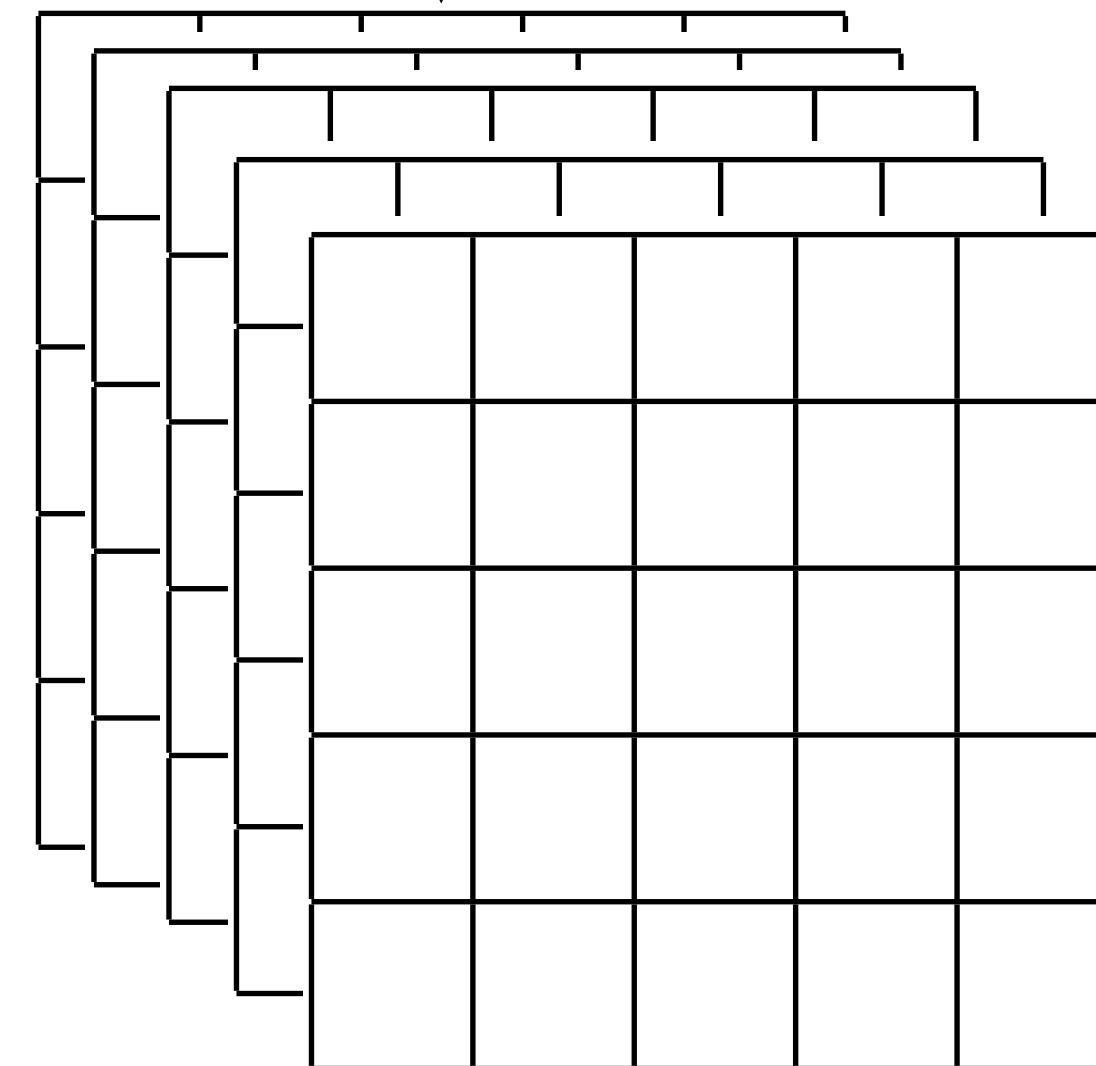
as many output channels
as semantic labels

Semantic segmentation

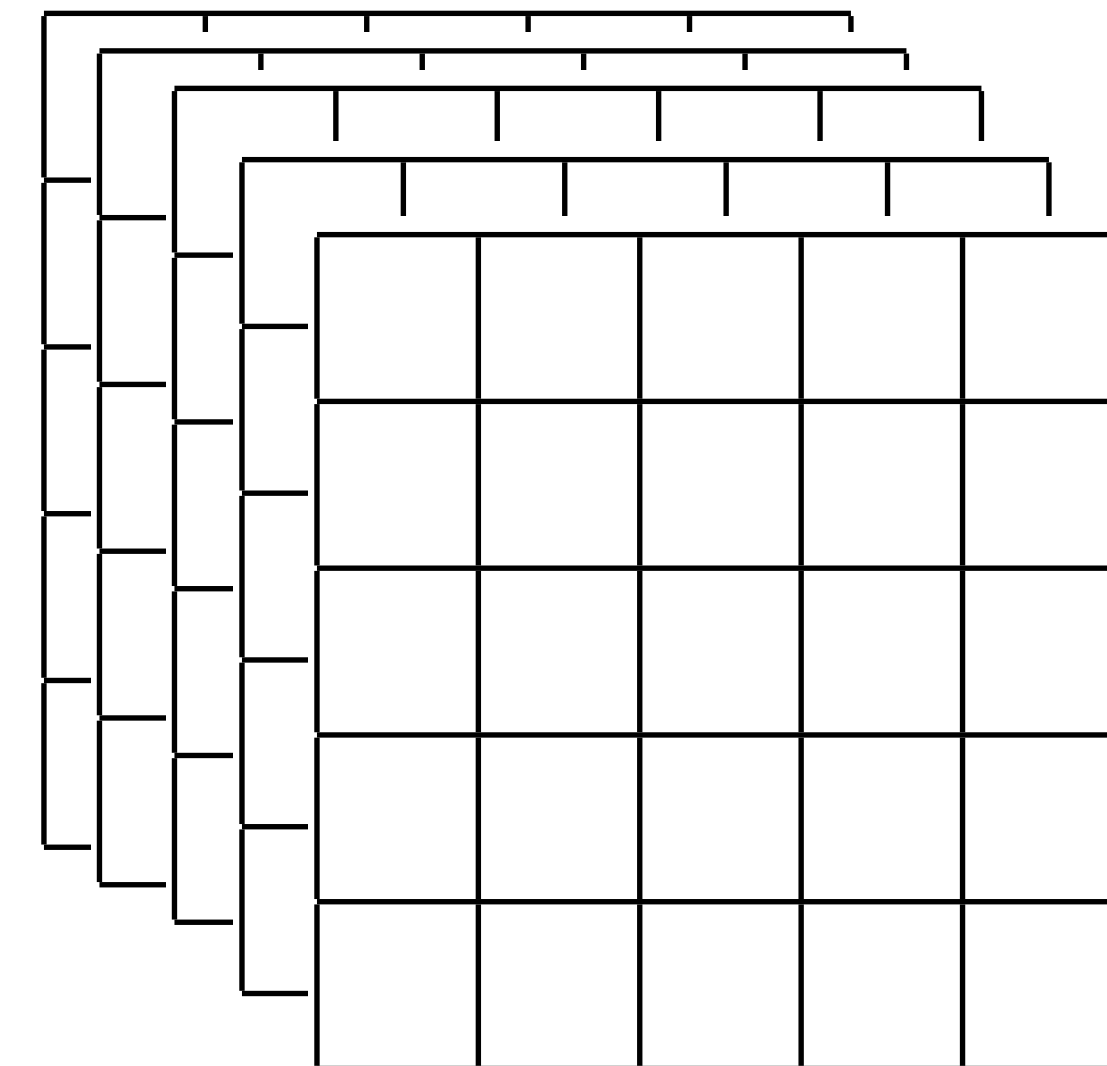
RGB image
($H \times W \times 3$)



pixel-wise probs
($H \times W \times N$)



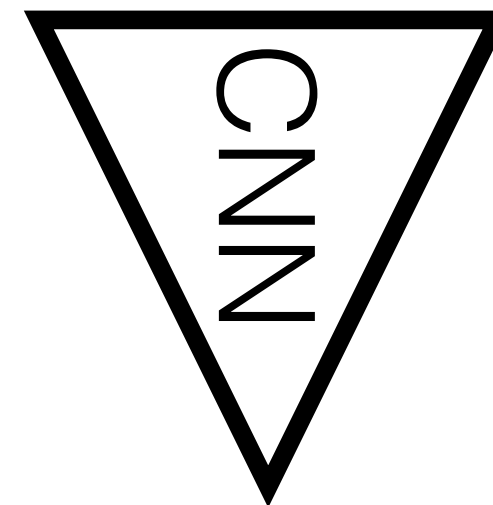
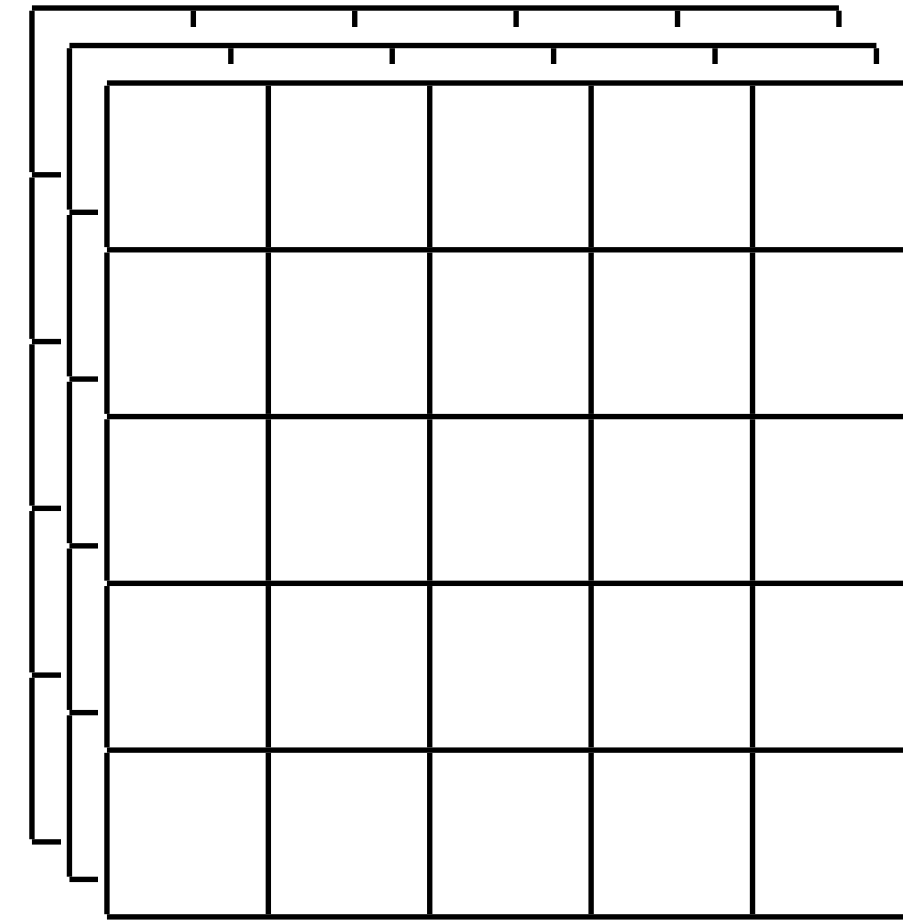
ground truth (1-hot encoding)



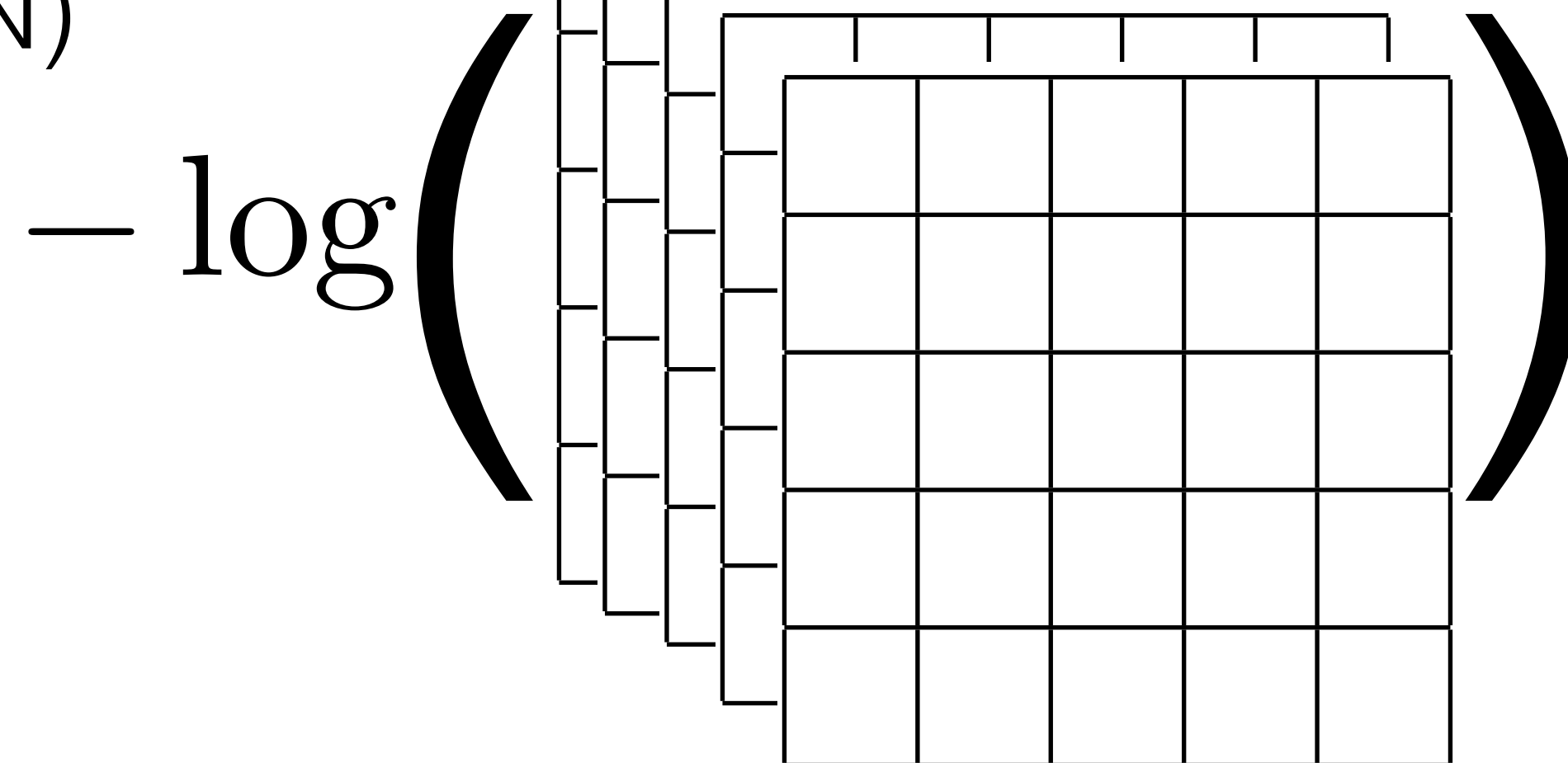
annotations
($H \times W \times N$)

Semantic segmentation

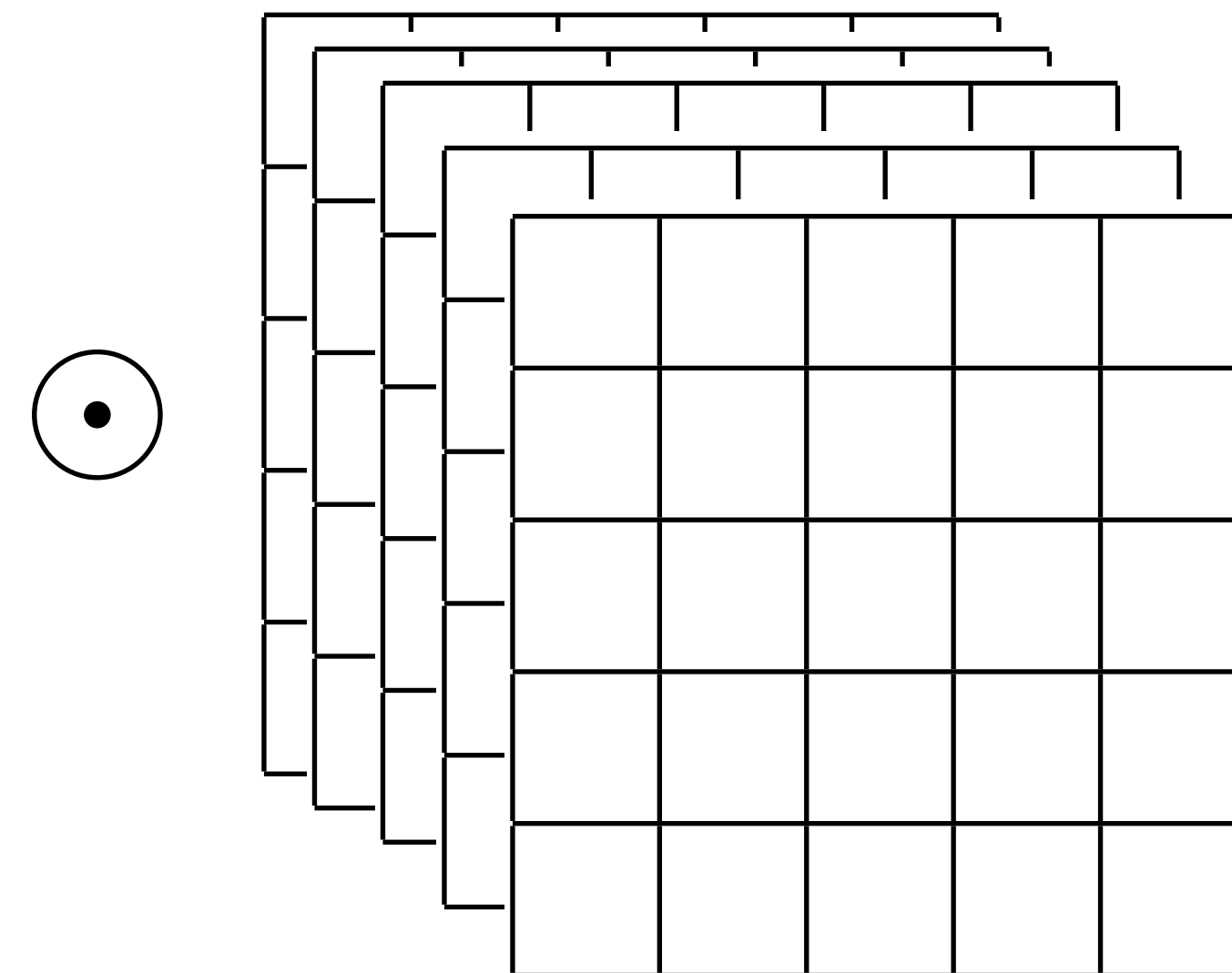
RGB image
(HxWx3)



pixel-wise CE-loss
(HxWxN)



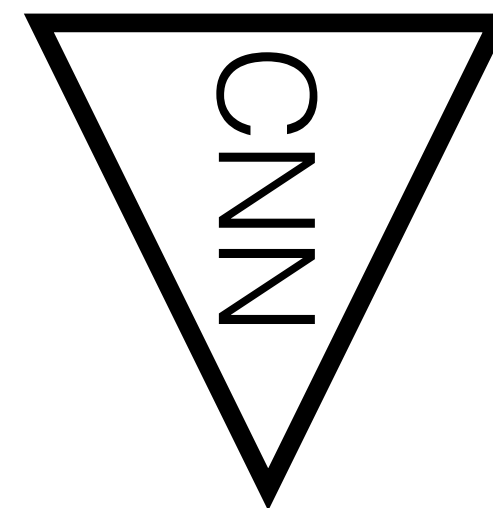
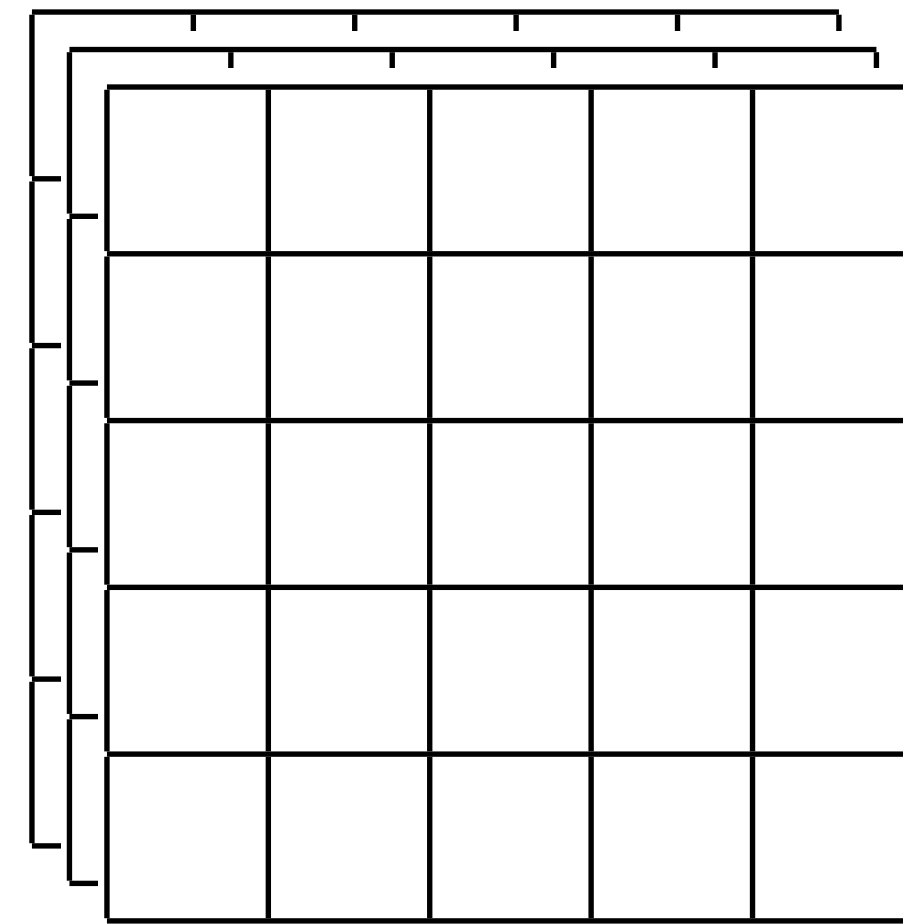
ground truth (1-hot encoding)



Semantic segmentation

Per-pixel classifier with the loss summed over all pixels

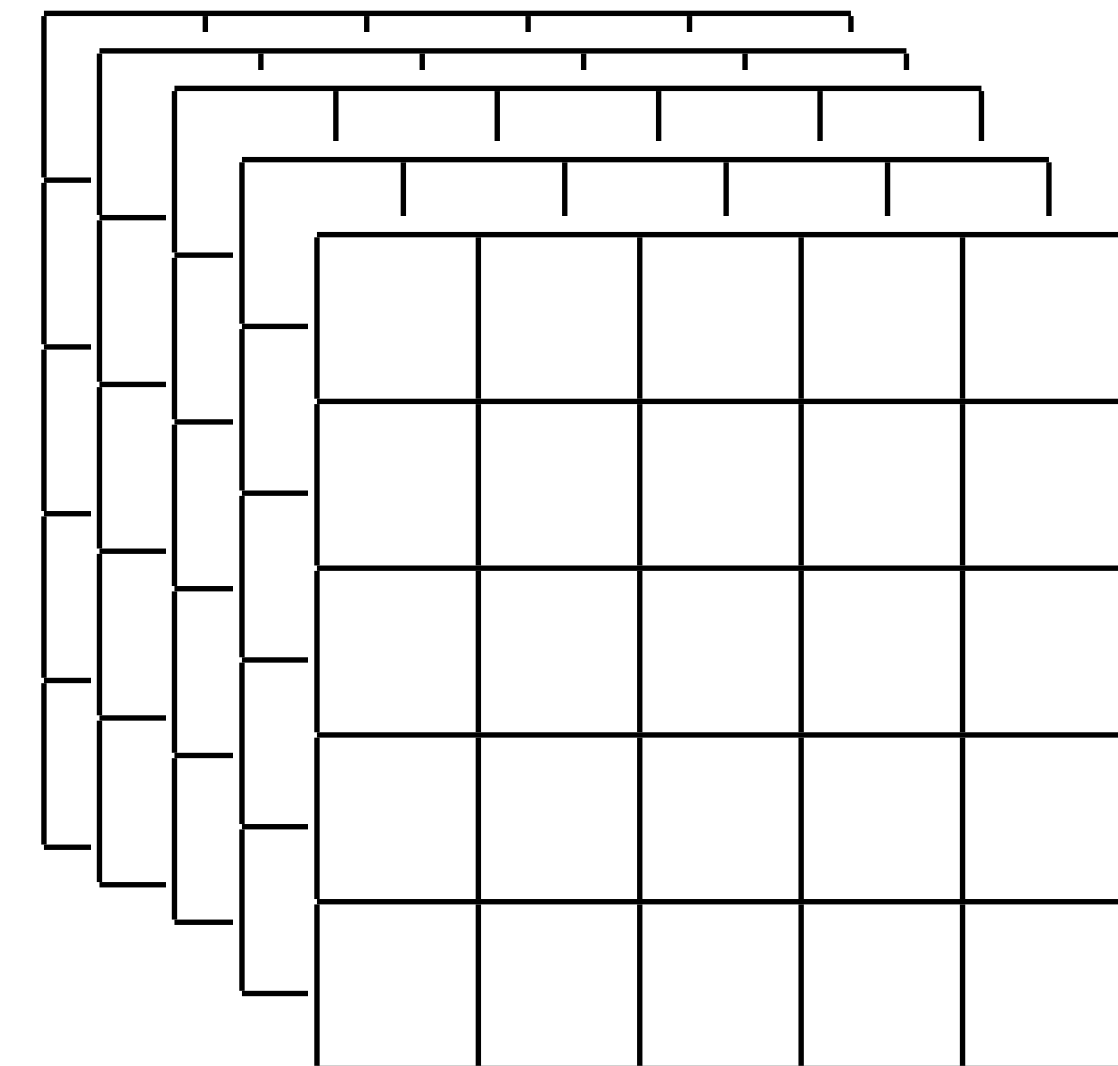
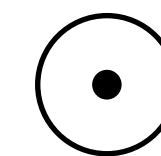
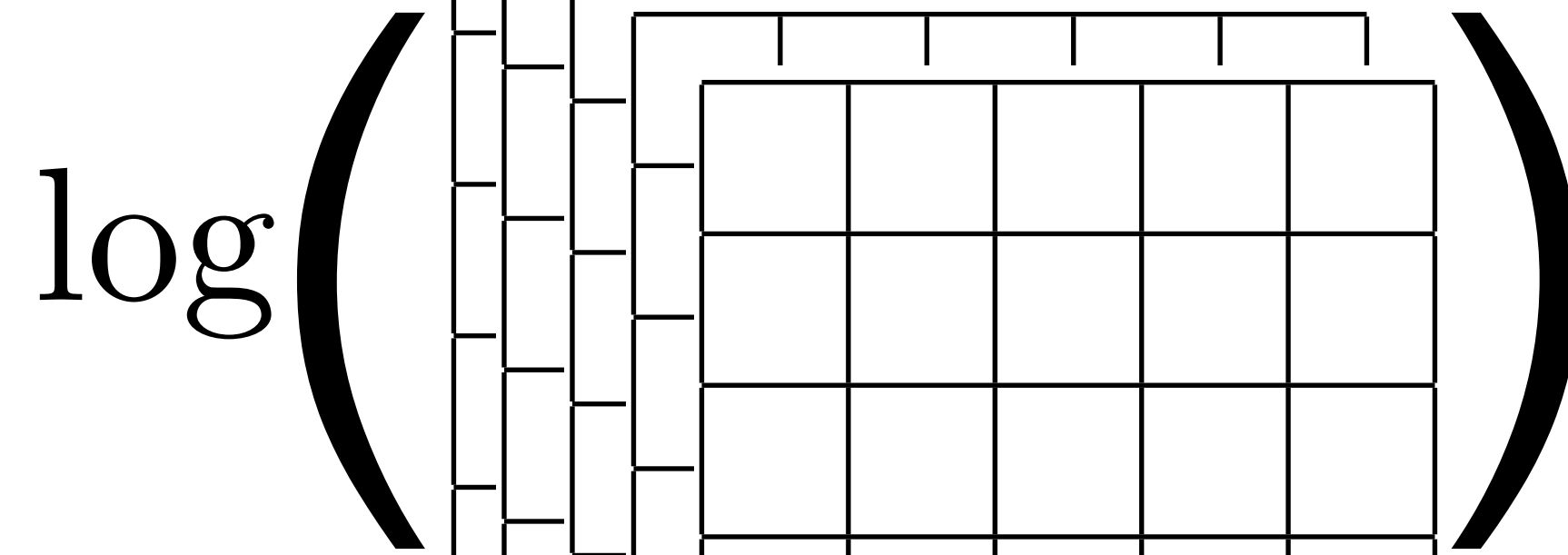
RGB image
(HxWx3)



cross-entropy loss

ground truth (1-hot encoding)

$$\sum_{\text{pixels}} -\log$$



U-net architecture



Mirror the usual classification network

[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

high spatial resolution

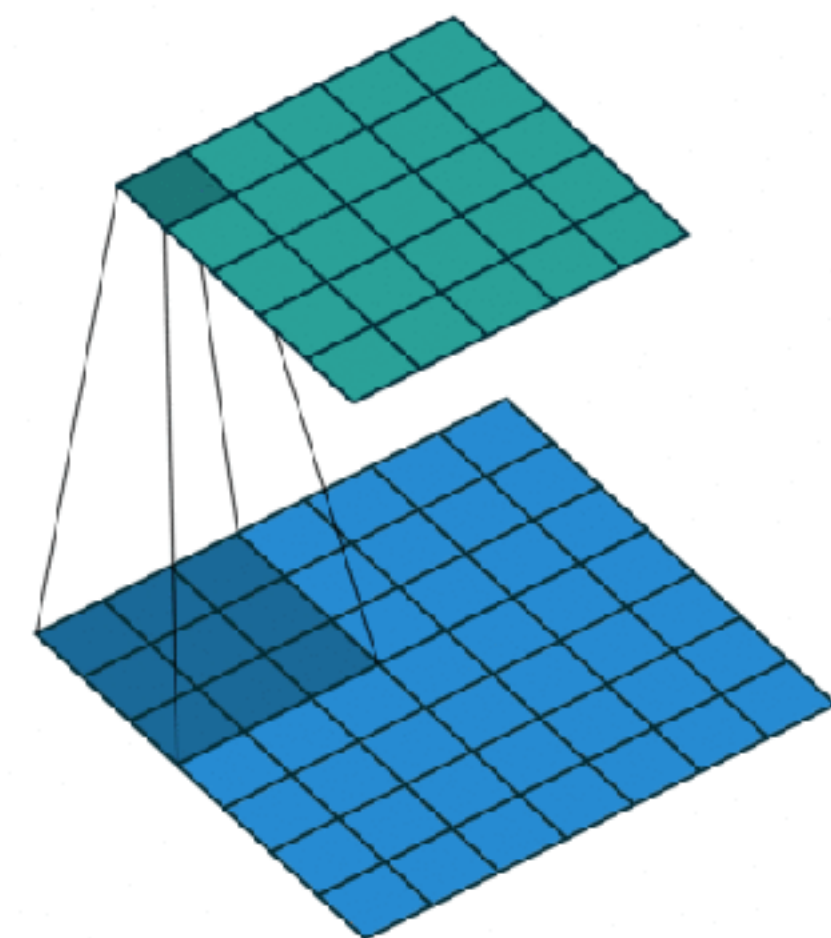
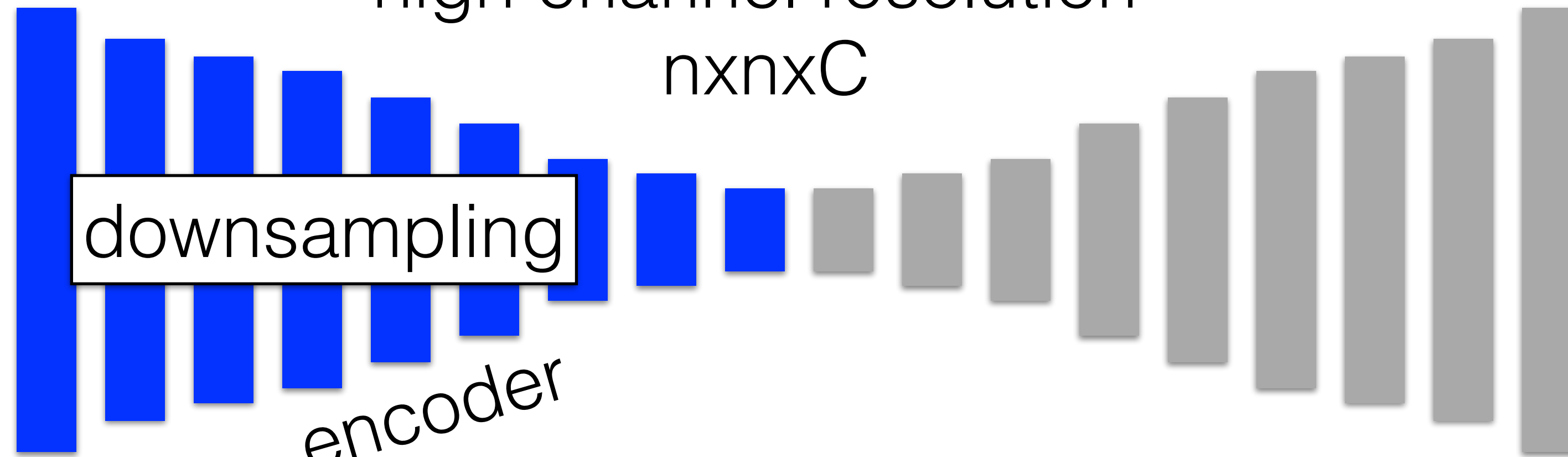
$N \times N \times c$

image embedding

high channel resolution

$n \times n \times C$

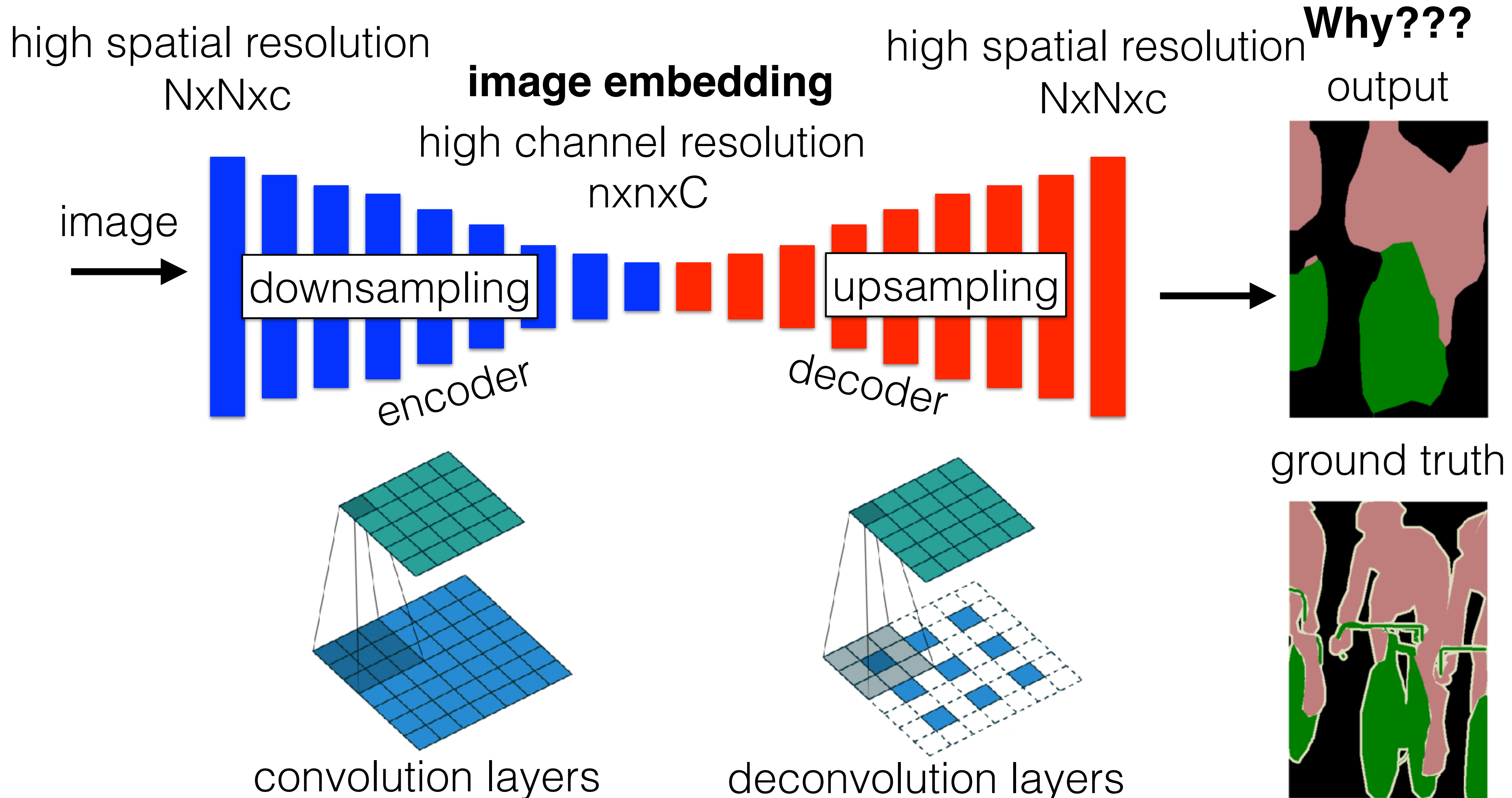
image



convolution layers

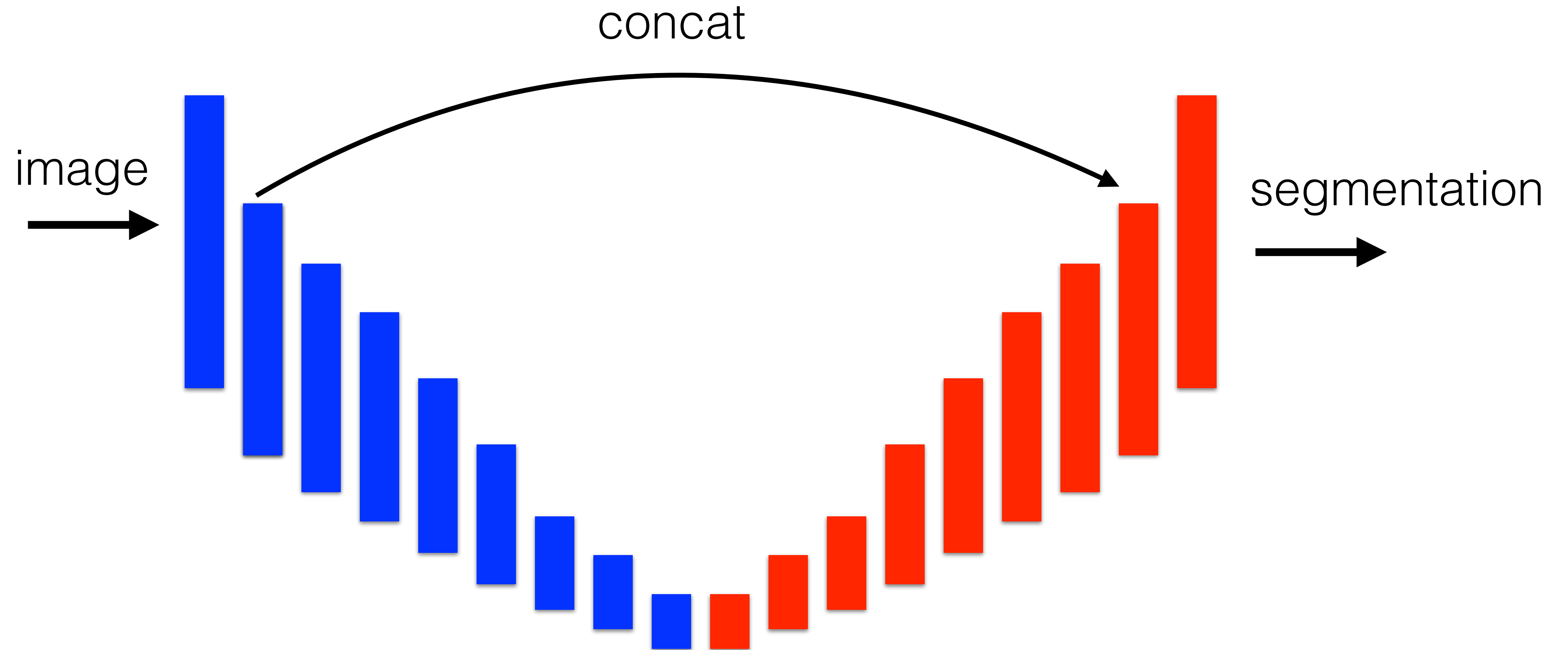
[Noh et al ICCV 2015] <https://arxiv.org/pdf/1505.04366.pdf>

U-net architecture



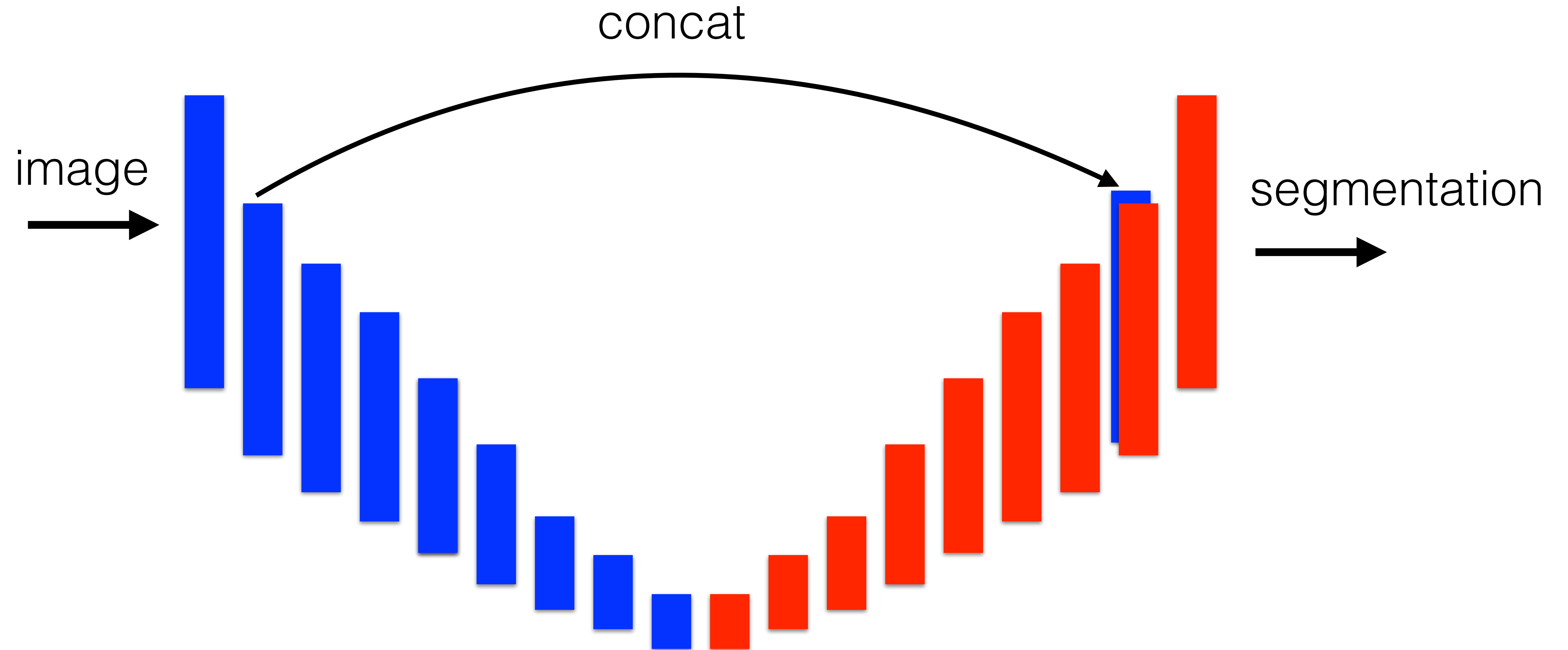
[Noh et al ICCV 2015] <https://arxiv.org/pdf/1505.04366.pdf>

U-net architecture



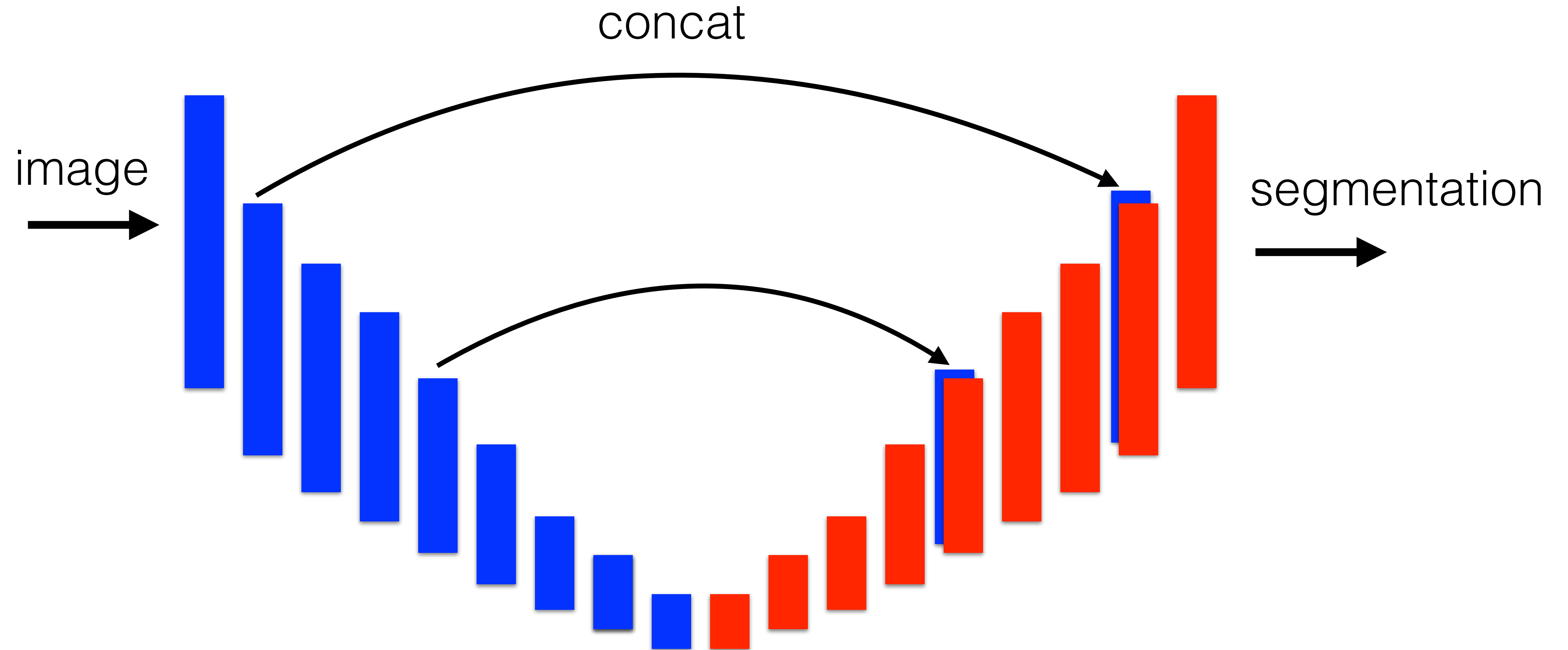
[Noh et al ICCV 2015] <https://arxiv.org/pdf/1505.04366.pdf>

U-net architecture



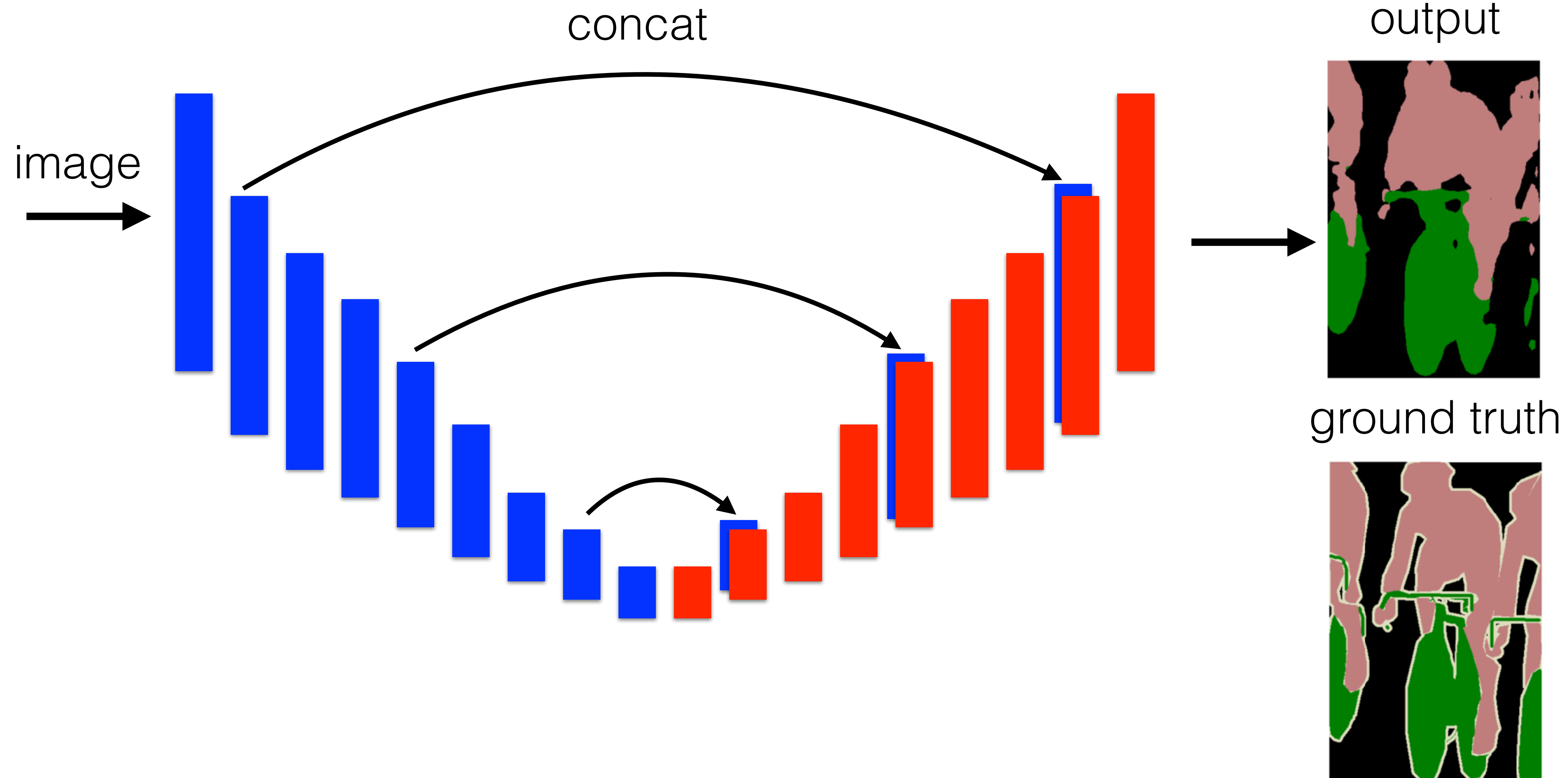
[Noh et al ICCV 2015] <https://arxiv.org/pdf/1505.04366.pdf>

U-net architecture



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

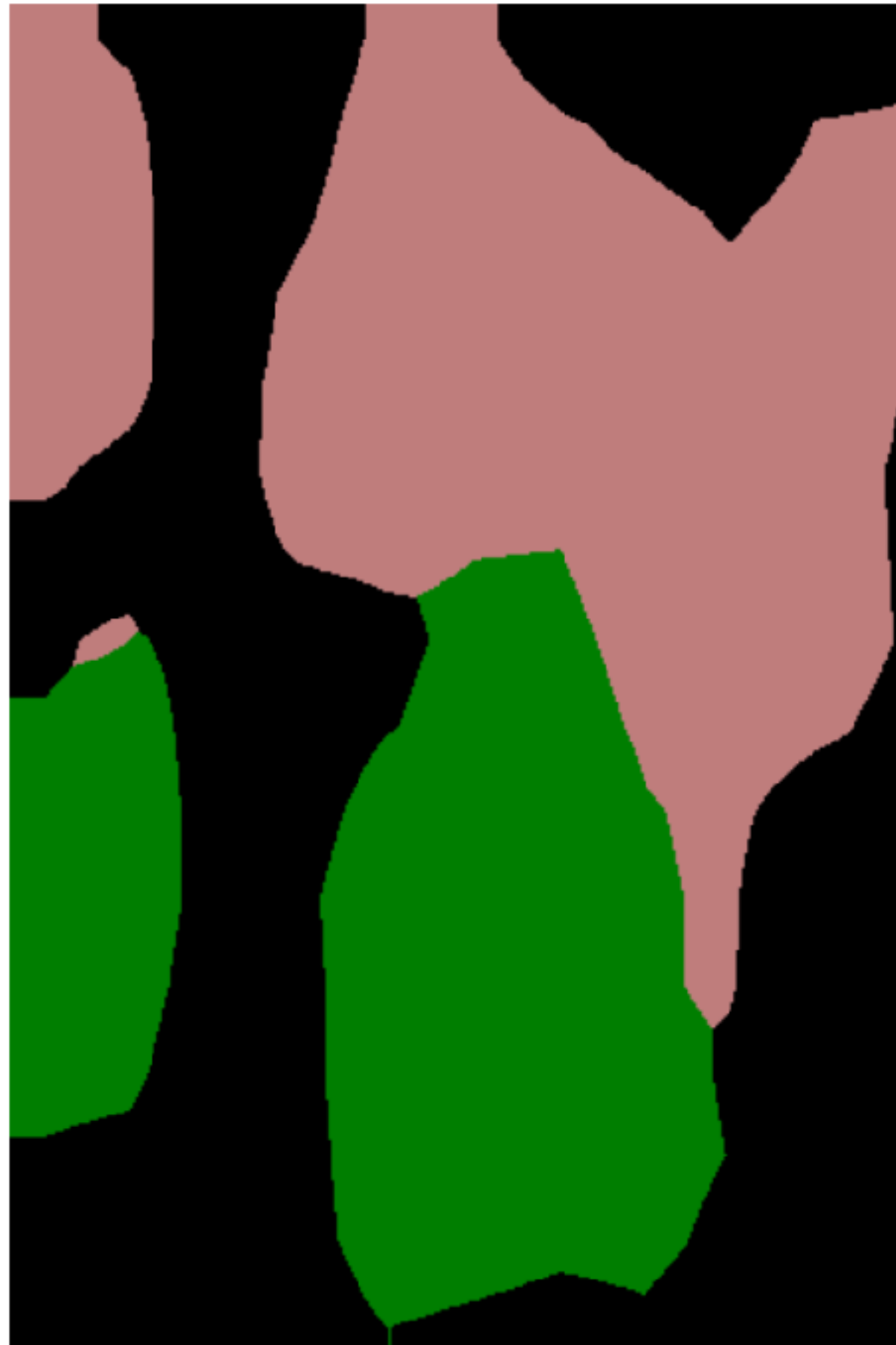
U-net architecture



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

U-net architecture

no skip connections



with skip connections



ground truth

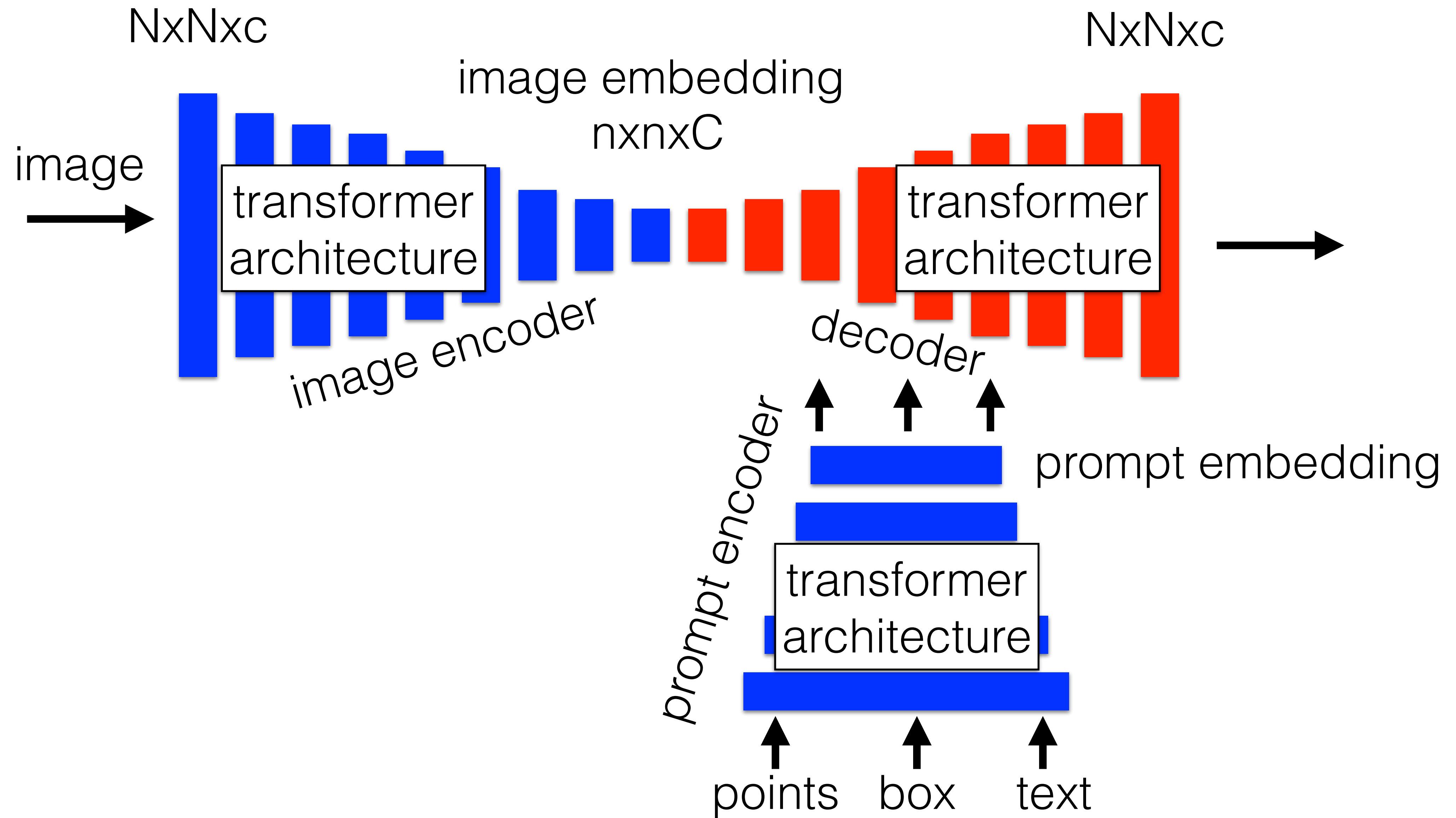


[Noh et al ICCV 2015] <https://arxiv.org/pdf/1505.04366.pdf>

Segmentation architectures

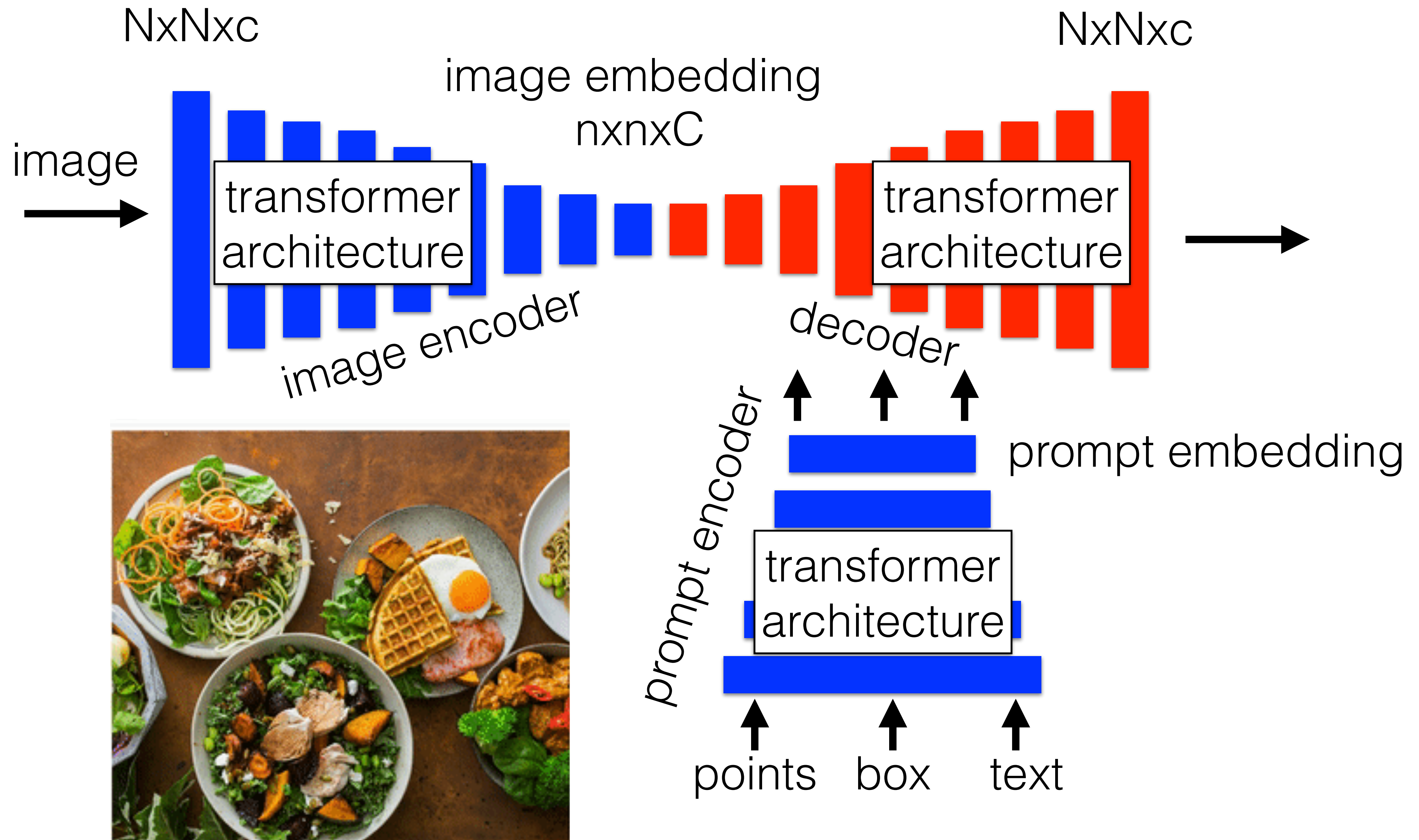
- **FCN** (Fully Convolutional Network): LongCVPR 2015, FCN was one of the pioneering architectures for end-to-end pixel-wise prediction.
- **U-Net**: symmetric architecture and skip connections, which helps in capturing both global and local features.
- **DeepLab**: DeepLab uses dilated convolutions 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 extended for segmentation task
- **SAM**: Segment anything architecture from Facebook 2023

SAM: Segment anything



[Facebook 2023] <https://arxiv.org/pdf/2304.02643.pdf>

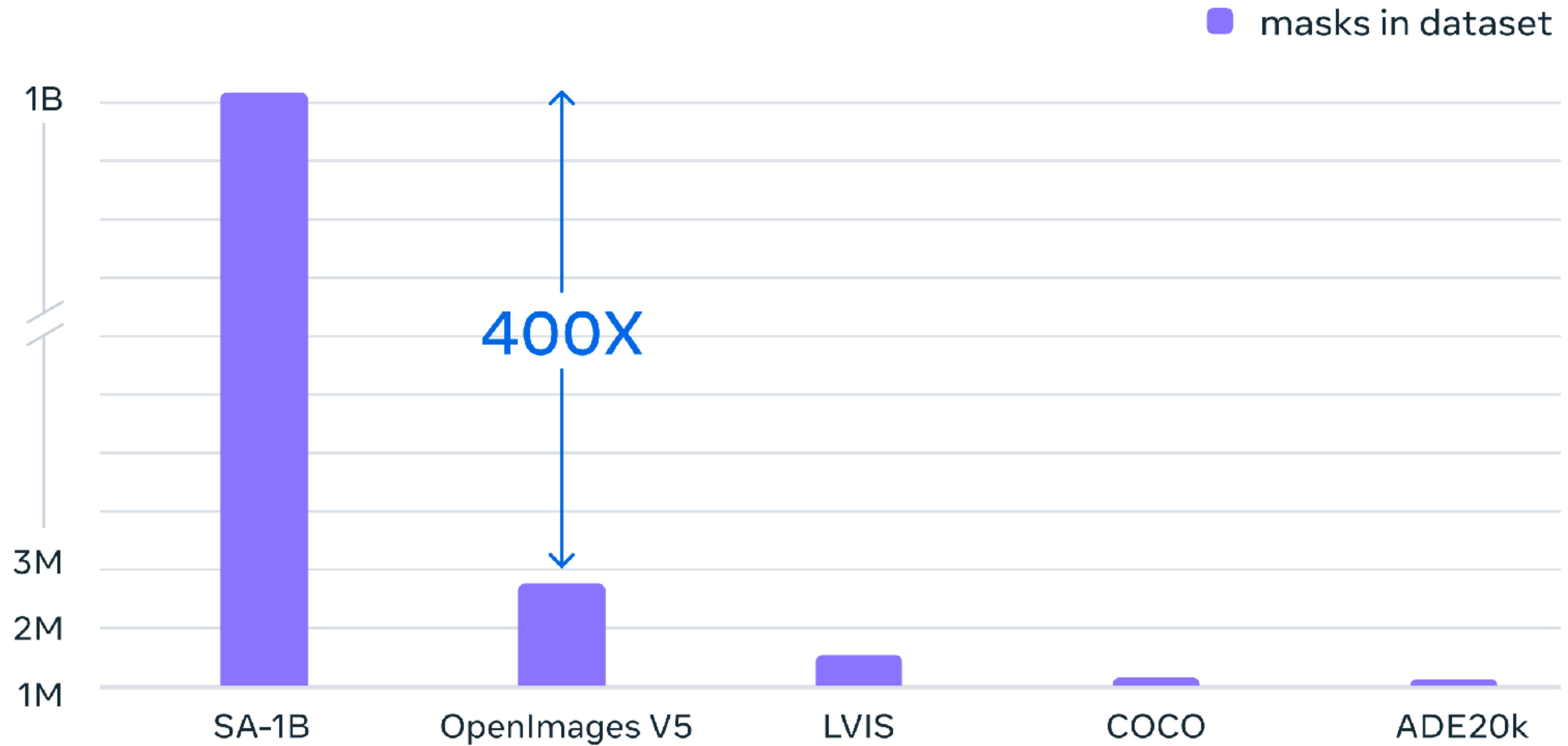
SAM: Segment anything



[Facebook 2023]

<https://arxiv.org/pdf/2304.02643.pdf>

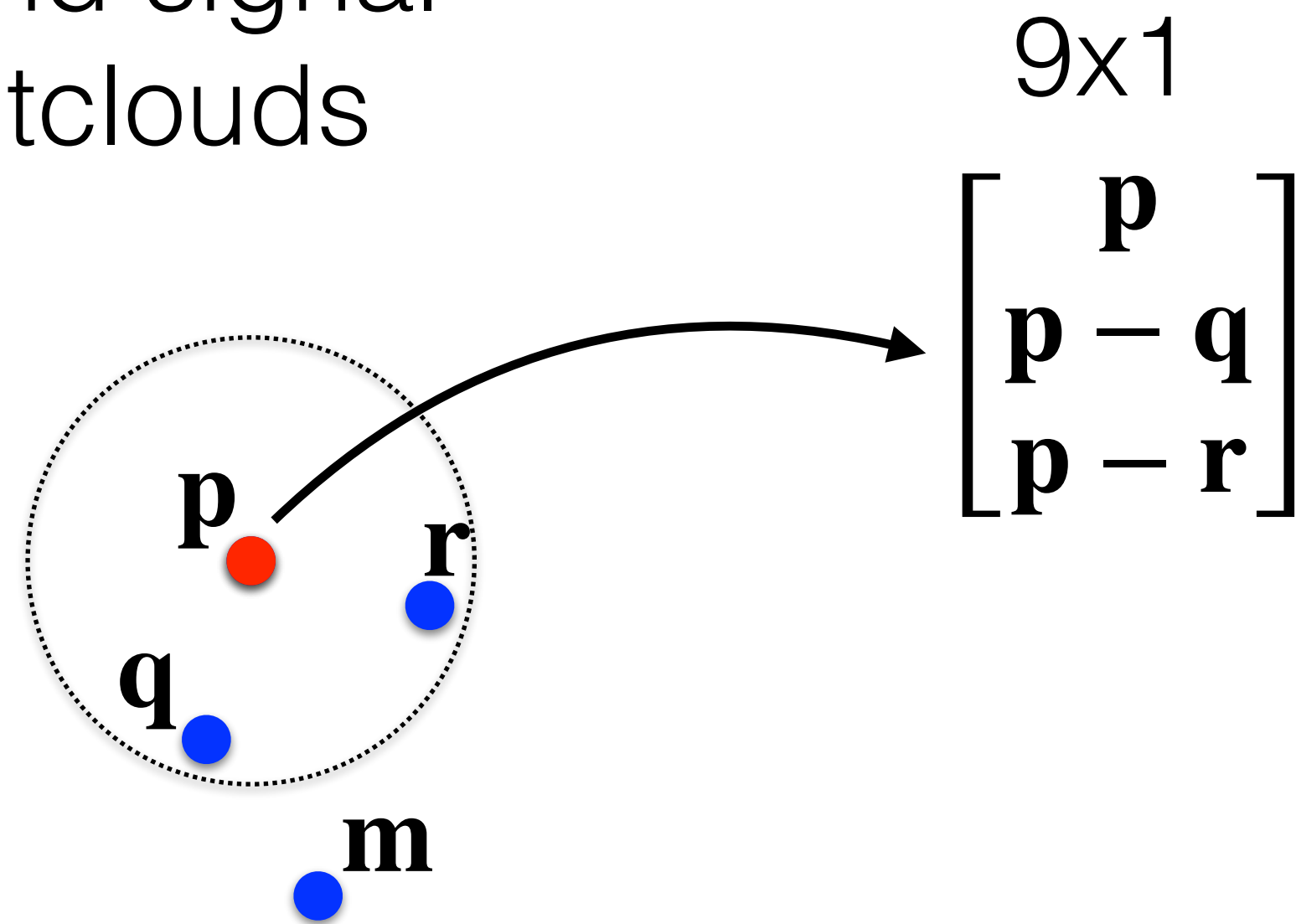
SAM: Segment anything



[Facebook 2023] <https://arxiv.org/pdf/2304.02643.pdf>

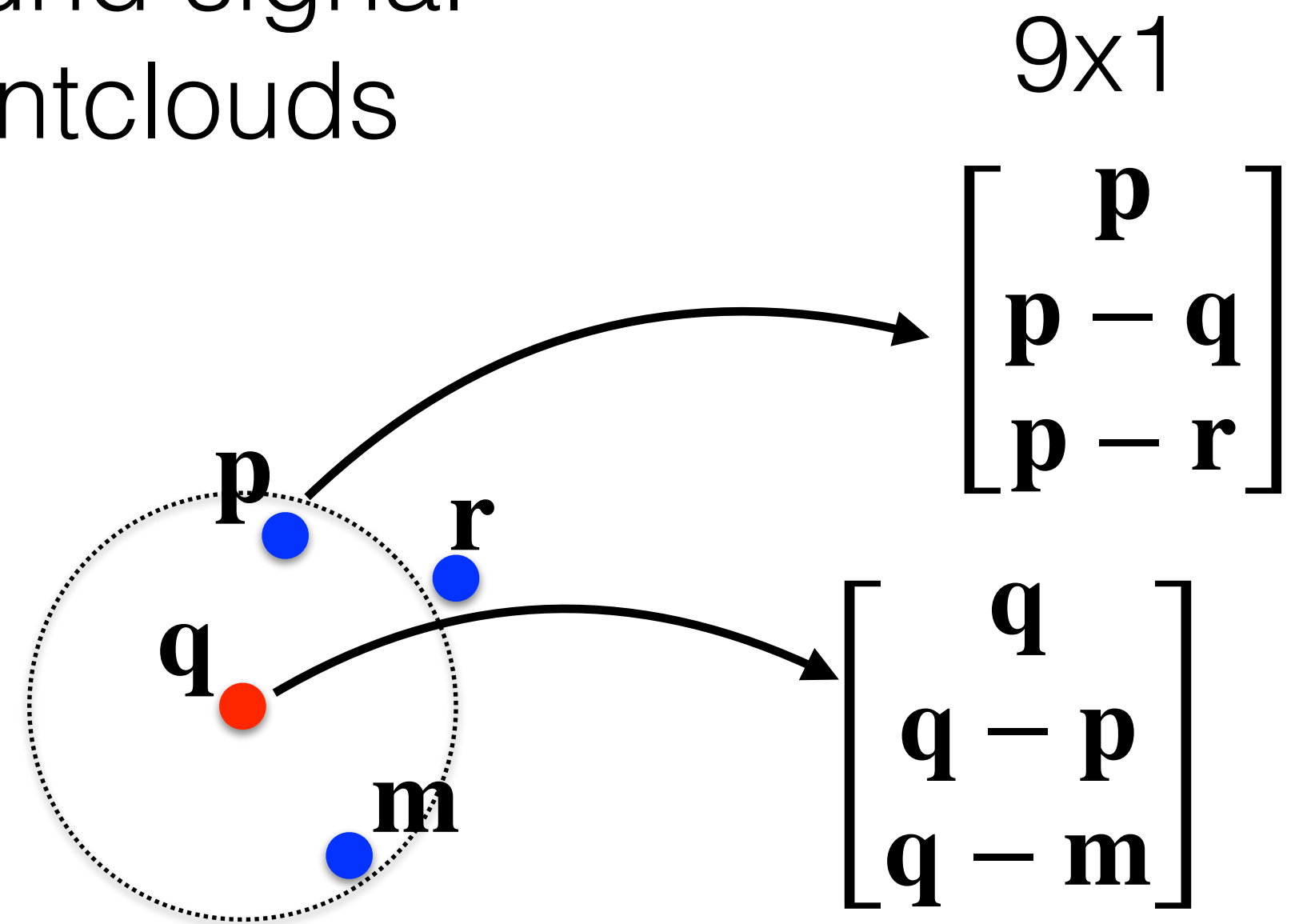
Different data modalities

- Images (RGB, thermal, RGBD, event-based cameras, 3D MRI scans)
- Videos (sequences of images)
- Text (phrases)
- Sound signal
- Pointclouds



Different data modalities

- Images (RGB, thermal, RGBD, event-based cameras, 3D MRI scans)
- Videos (sequences of images)
- Text (phrases)
- Sound signal
- Pointclouds



Different data modalities

- Images (RGB, thermal, RGBD, event-based cameras, 3D MRI scans)
- Videos (sequences of images)
- Text (phrases)
- Sound signal
- Pointclouds

