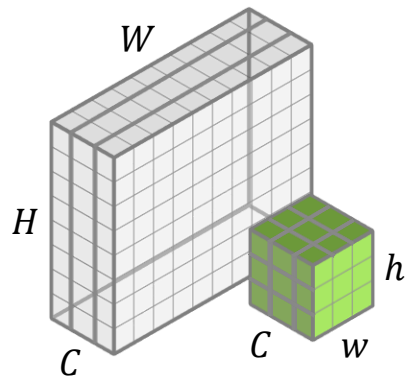


CNN Architectures

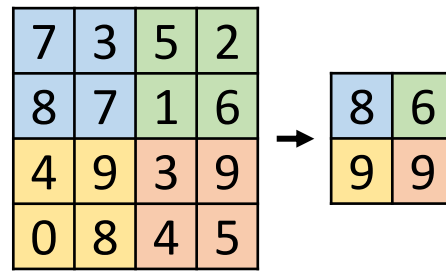
January 3rd, 2024



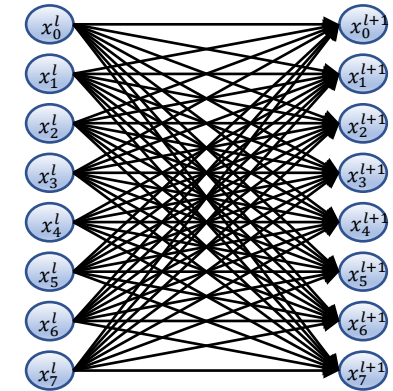
The ConvNet Building Blocks



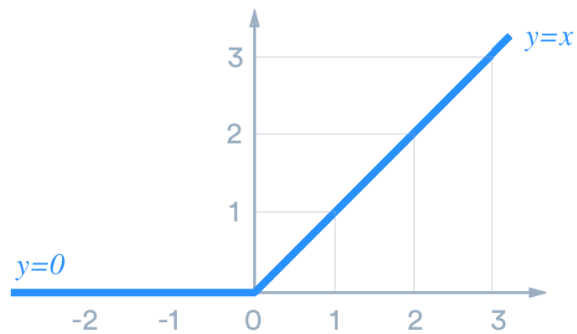
Conv Layer



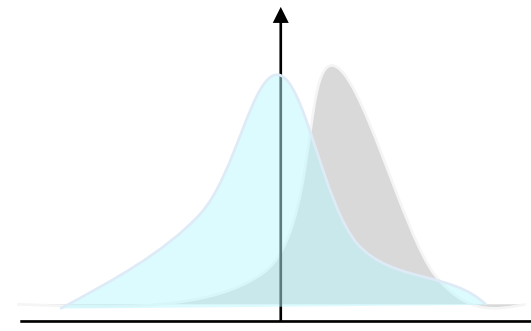
Max Pooling



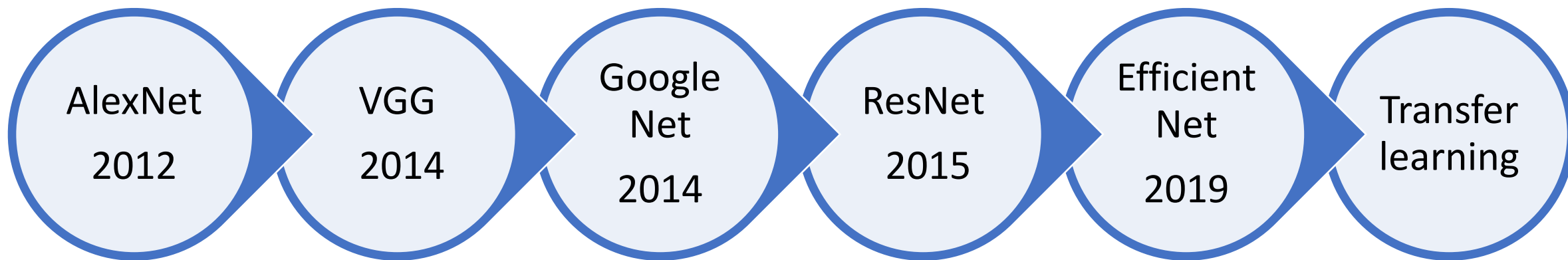
FC Layer



Activation Func.



Batch Norm

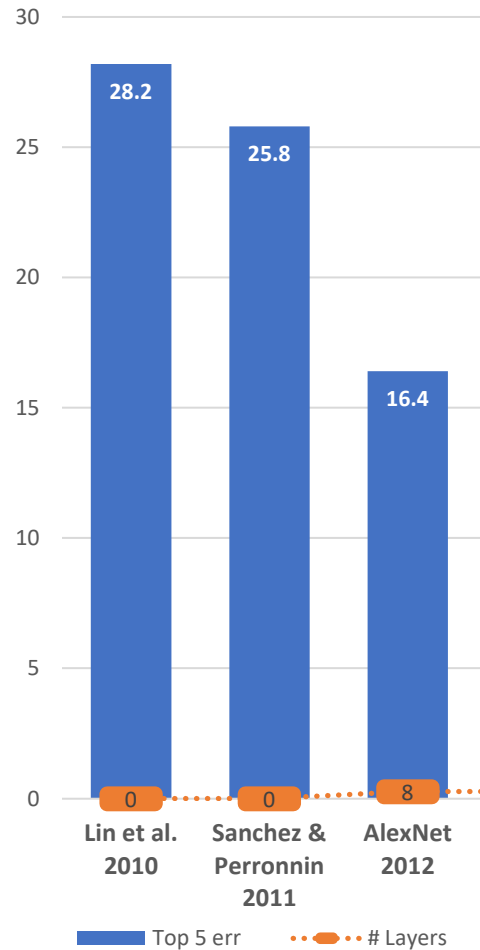
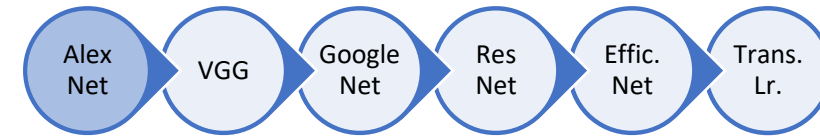




Visual Recognition Challenge (ILSVRC)

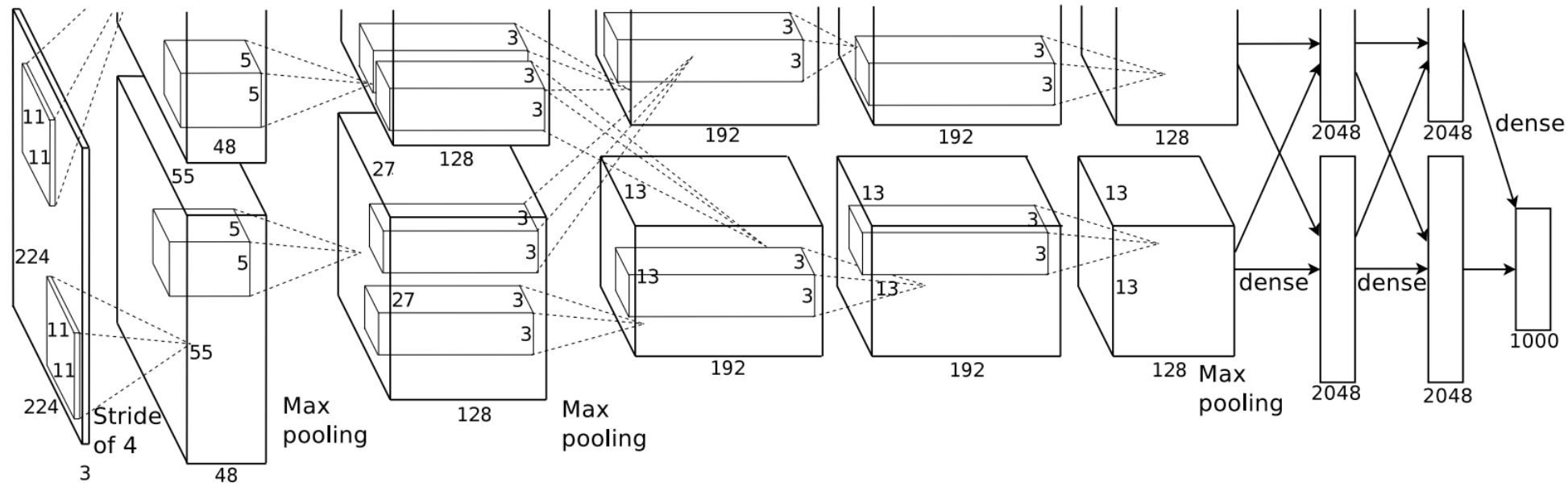
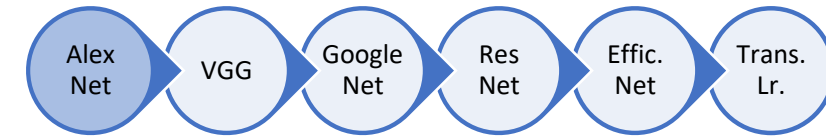
- 1.2M training images
- 100K testing images
- 1K categories

AlexNet



- Fukushima 1980
- LeCunn et al. 1989
- LeCunn et al. 1998

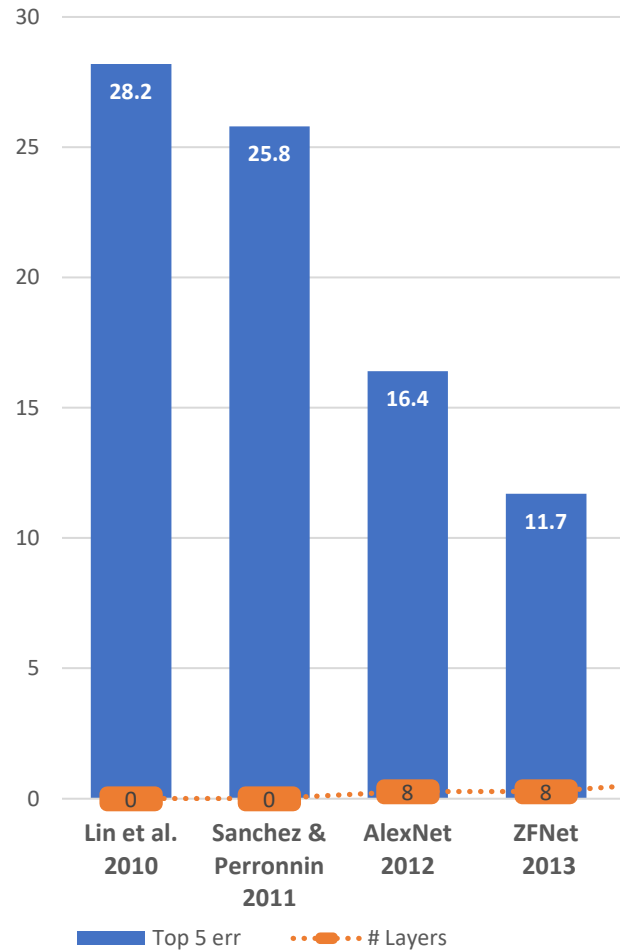
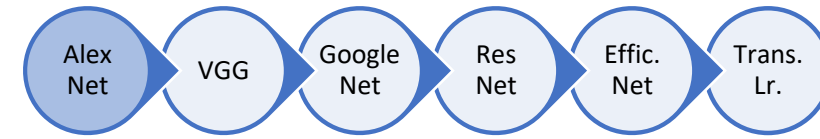
AlexNet



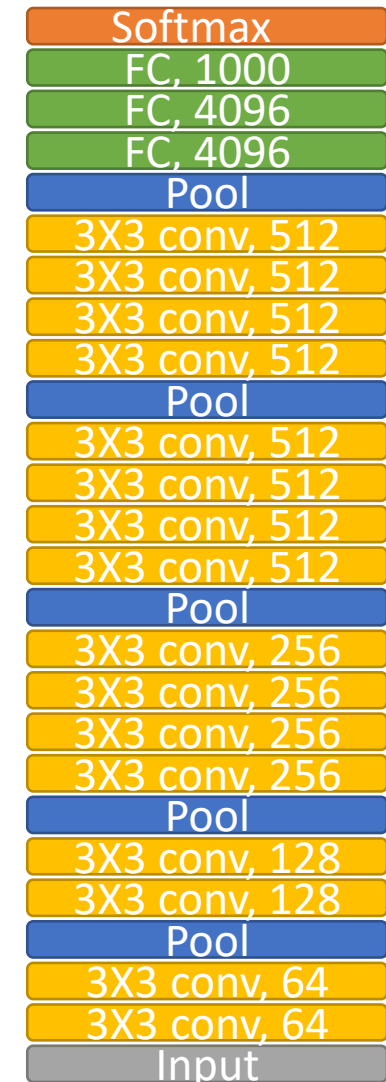
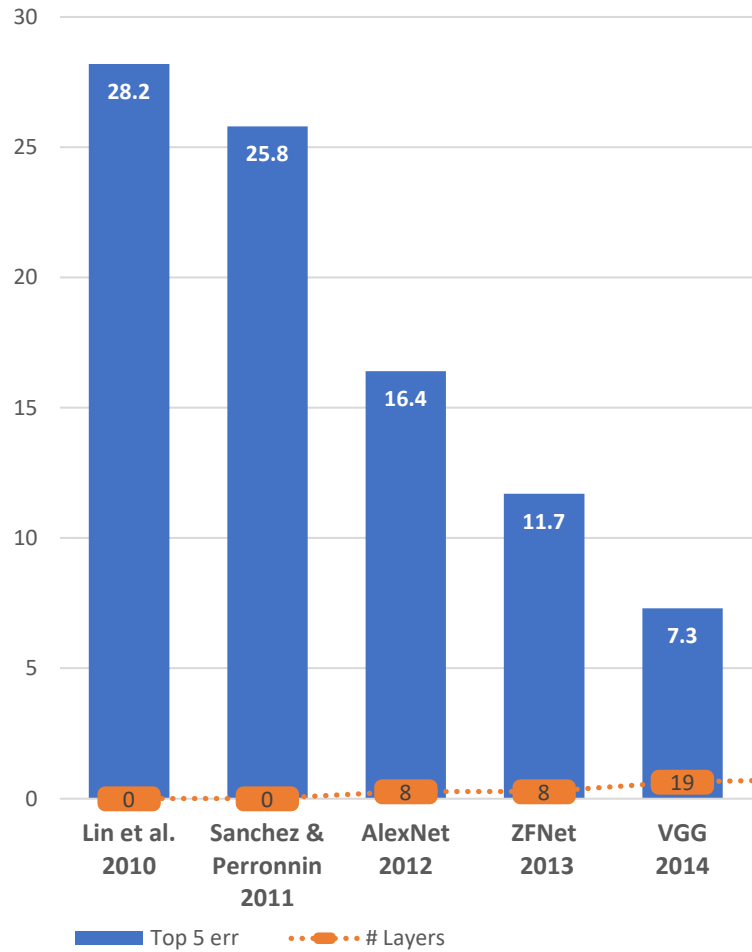
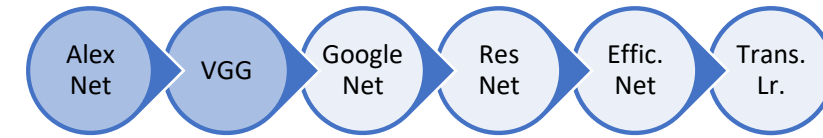
Architecture and training details:

- 5 conv + 3 fc layers
- ReLU activation
- Max Pooling
- Dropout
- SGD + Momentum

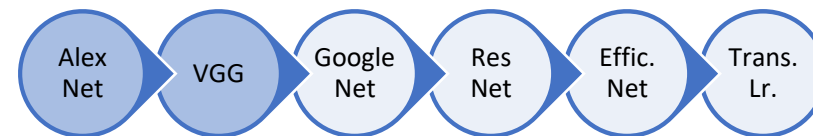
ZFNet



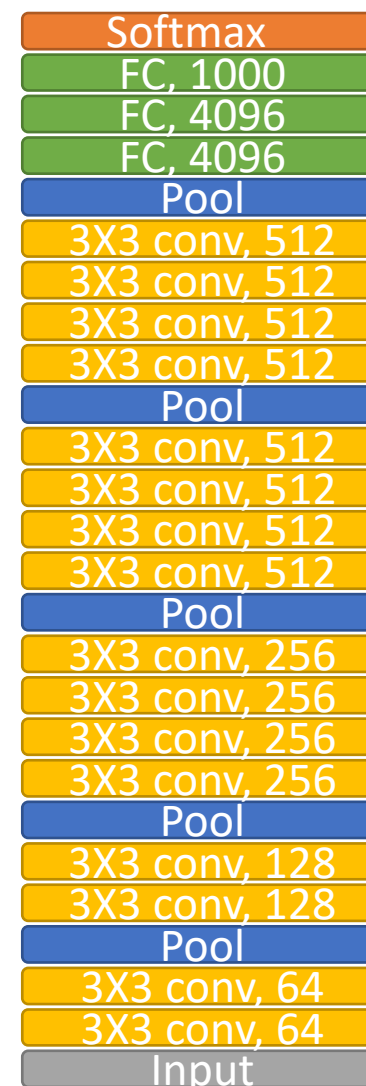
VGG



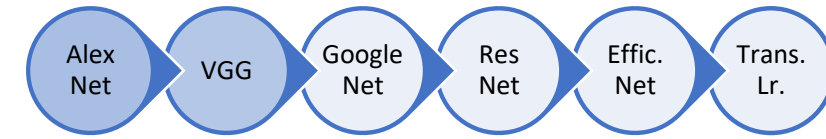
VGG



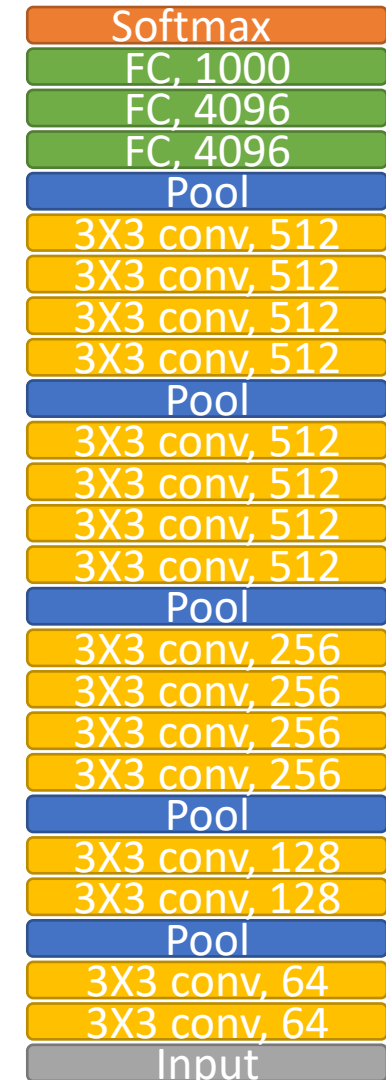
- All conv layers have 3x3 kernels, stride 1, pad 1
- All pooling layers are 2x2, stride 2
- After pooling, double the number of channels



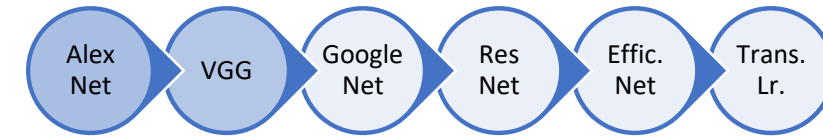
Why using 3x3 conv layers?



Receptive Field



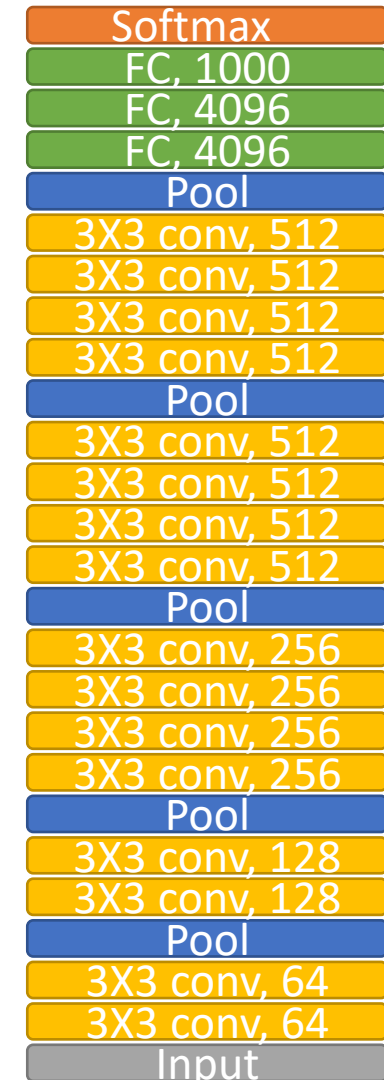
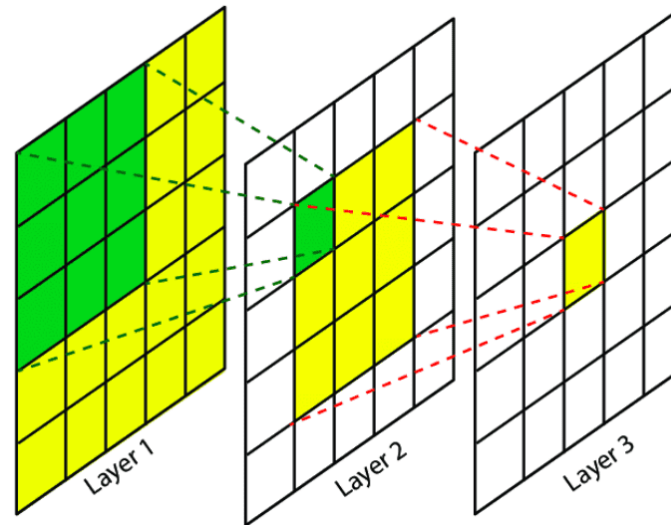
Why using 3x3 conv layers?



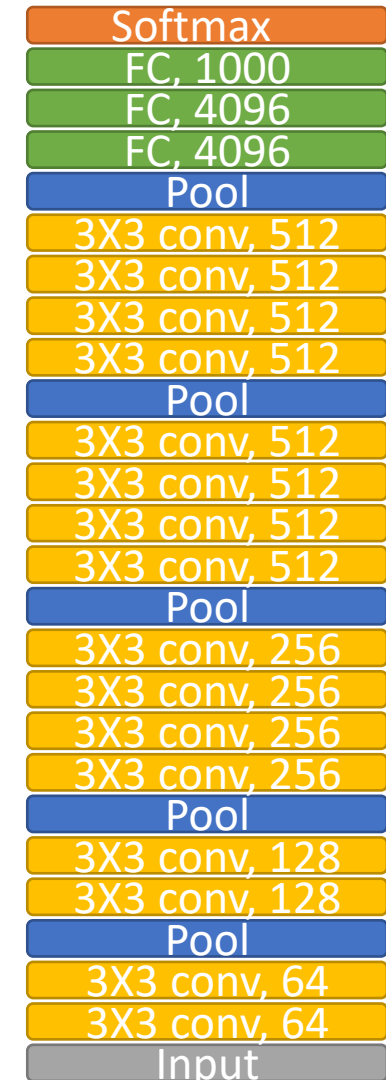
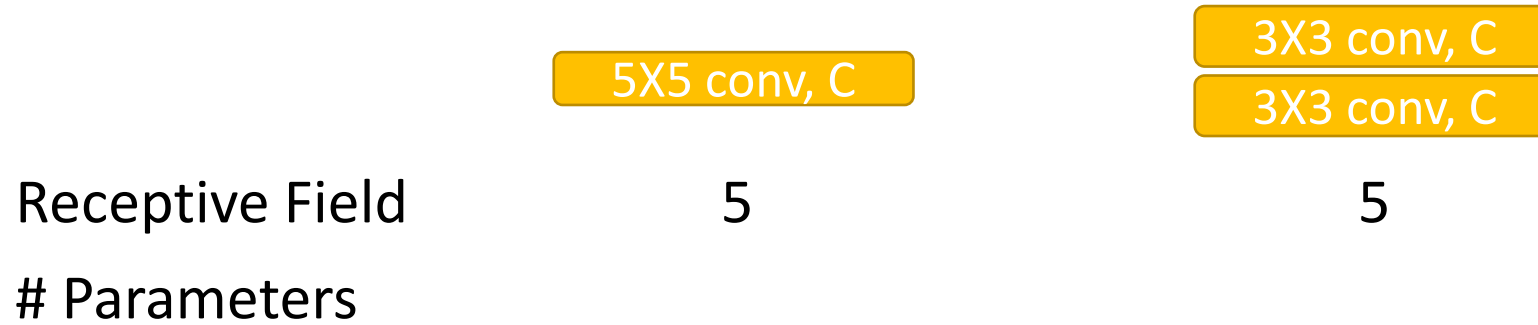
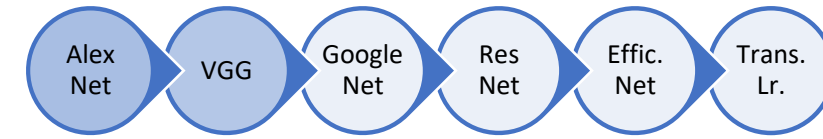
Receptive Field
Parameters



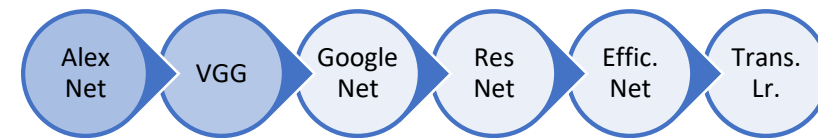
5



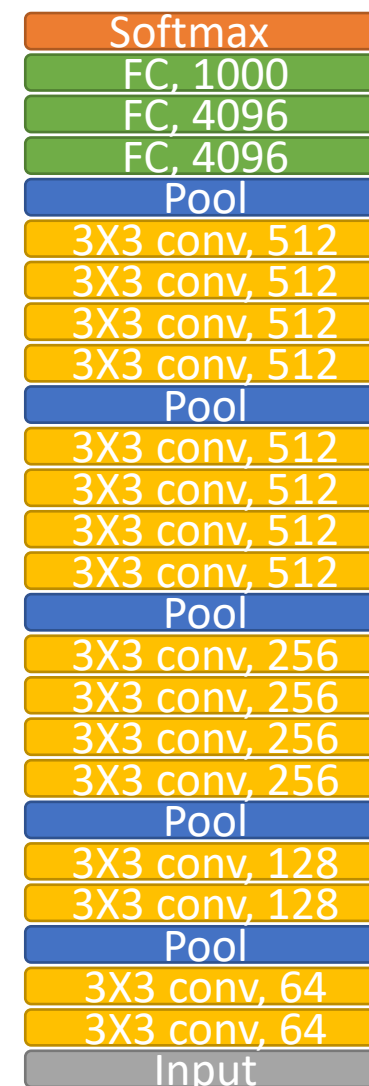
Why using 3x3 conv layers?



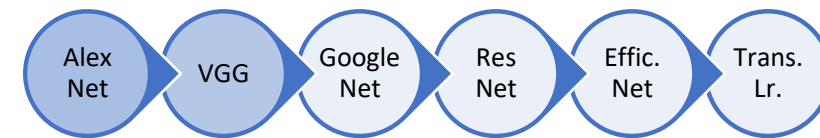
Why using 3x3 conv layers?



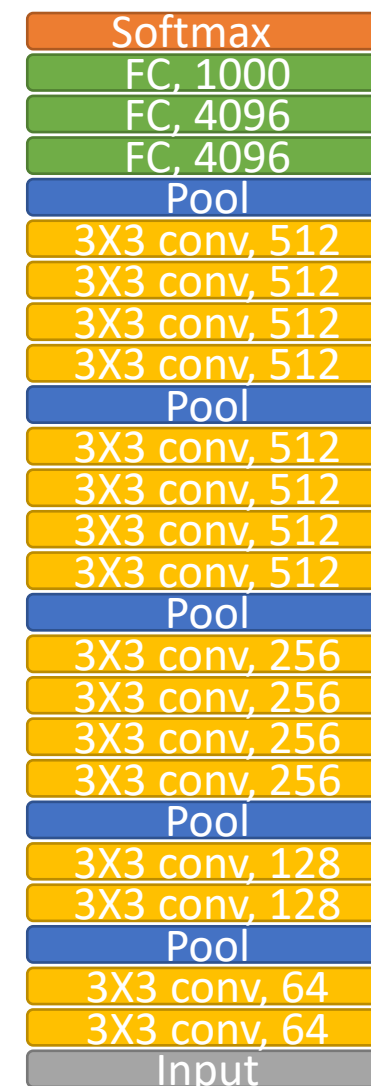
	5X5 conv, C	3X3 conv, C 3X3 conv, C
Receptive Field	5	5
# Parameters	$25C^2$	$9C^2 + 9C^2 = 18C^2$
# FLOPs		



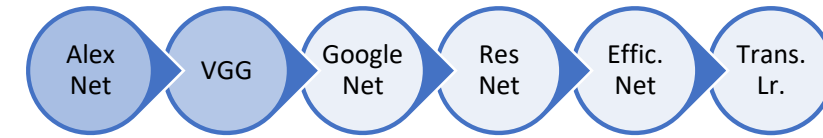
Why using 3x3 conv layers?



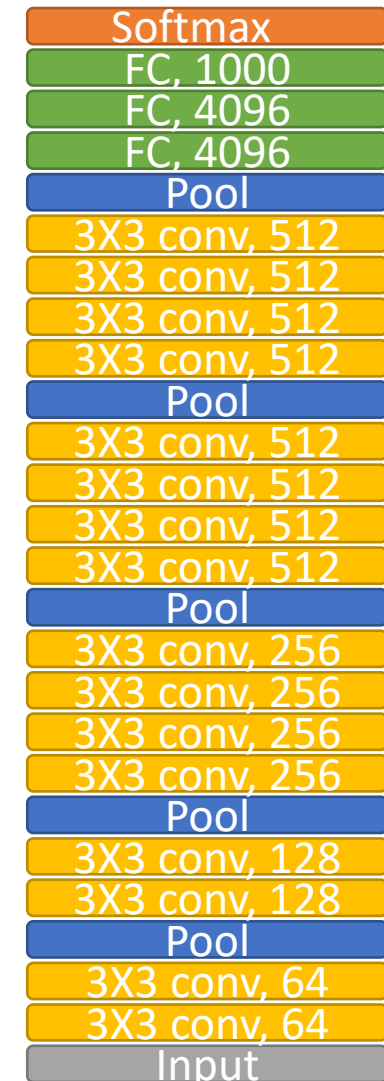
	5X5 conv, C	3X3 conv, C 3X3 conv, C
Receptive Field	5	5
# Parameters	$25C^2$	$9C^2 + 9C^2 = 18C^2$
# FLOPs	$\sim 25C^2HW$	$\sim 18C^2HW$



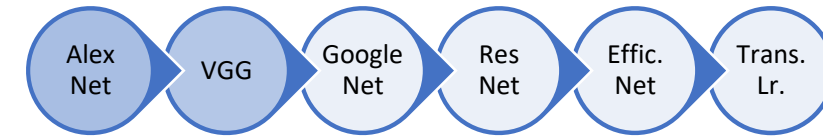
Why using 3x3 conv layers?



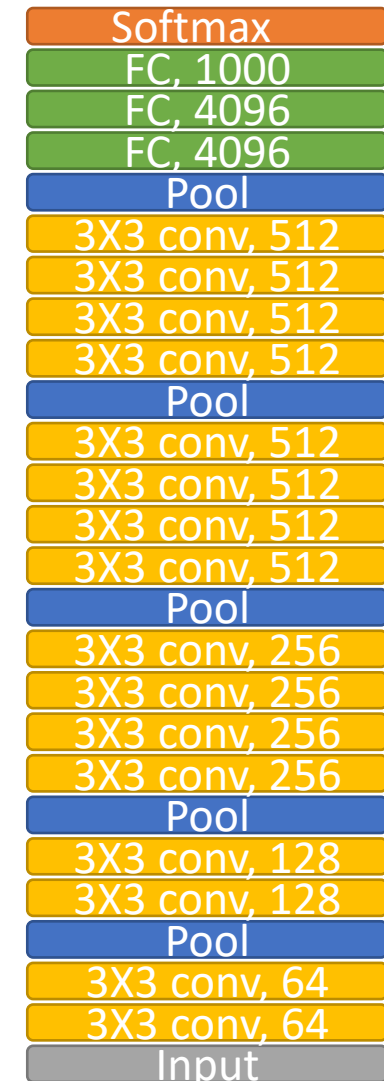
	5X5 conv, C	3X3 conv, C 3X3 conv, C
Receptive Field	5	5
# Parameters	$25C^2$	$9C^2 + 9C^2 = 18C^2$
# FLOPs	$\sim 25C^2HW$	$\sim 18C^2HW$
Memory Size (Output)		



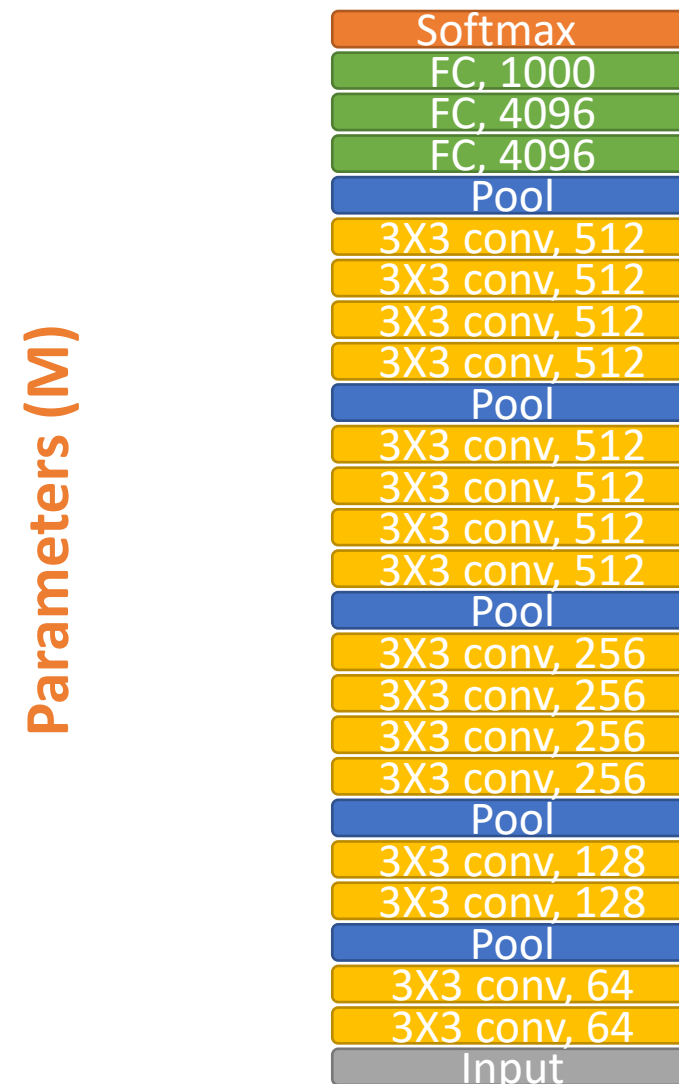
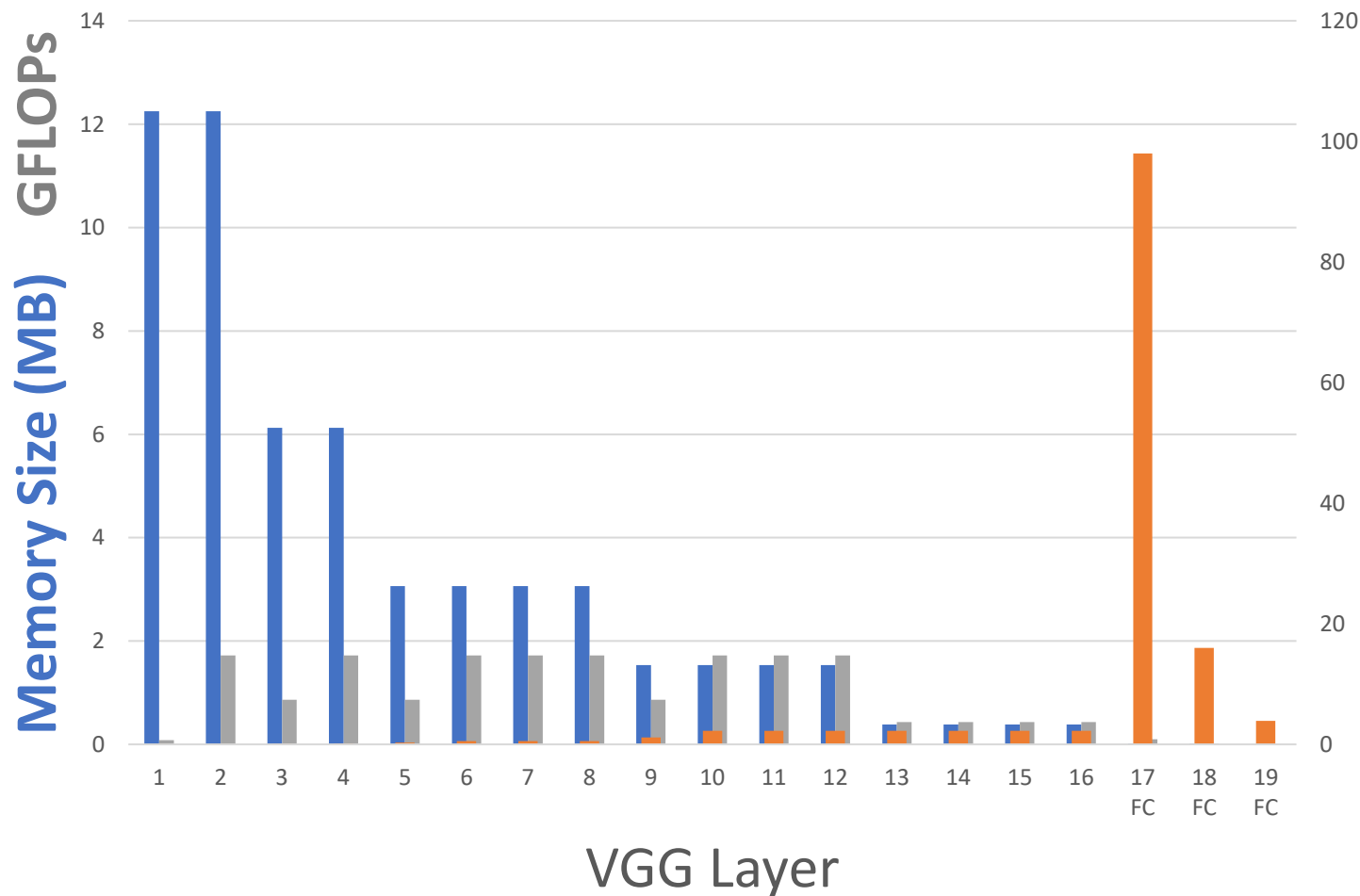
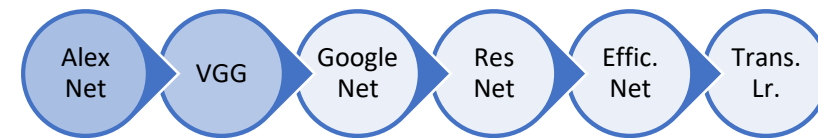
Why using 3x3 conv layers?



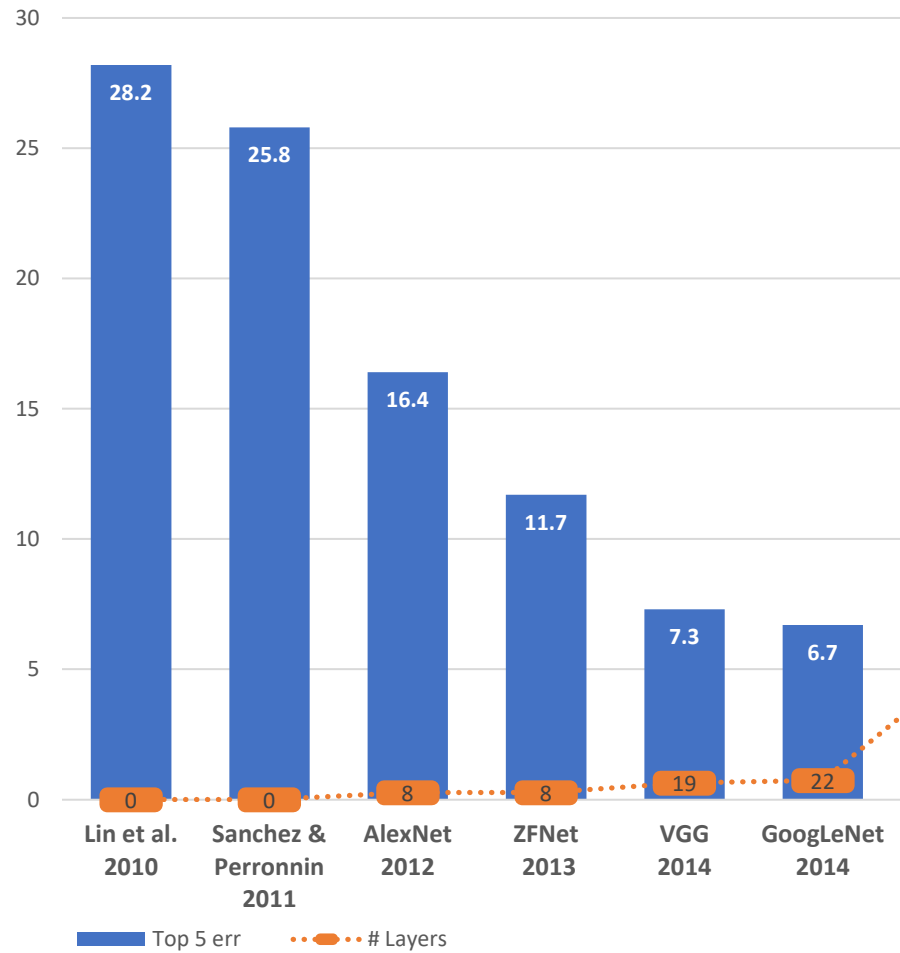
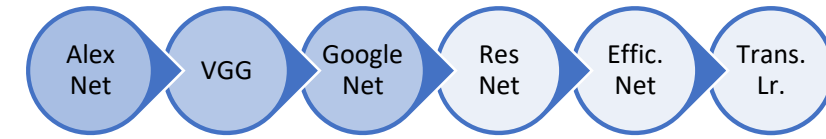
	5X5 conv, C	3X3 conv, C 3X3 conv, C
Receptive Field	5	5
# Parameters	$25C^2$	$9C^2 + 9C^2 = 18C^2$
# FLOPs	$\sim 25C^2HW$	$\sim 18C^2HW$
Memory Size (Output)	$\sim HWC$	$\sim 2HWC$



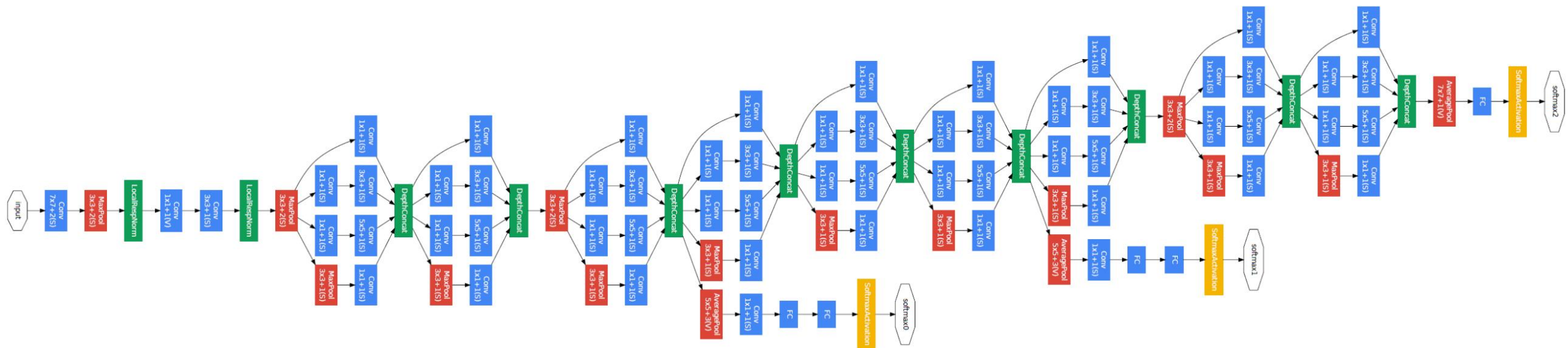
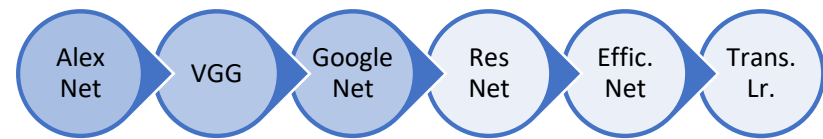
VGG resources



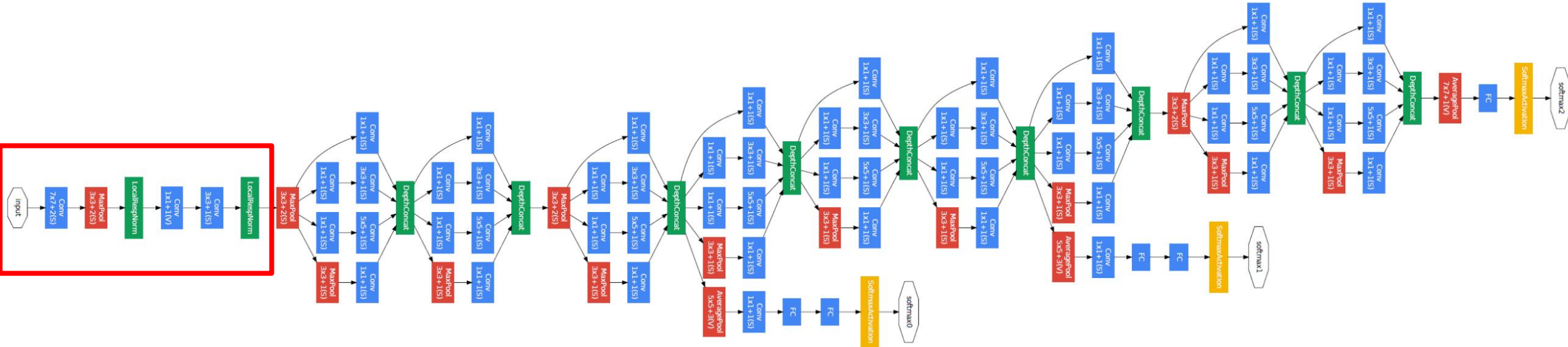
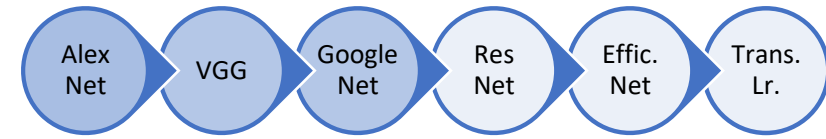
GoogLeNet



GoogLeNet

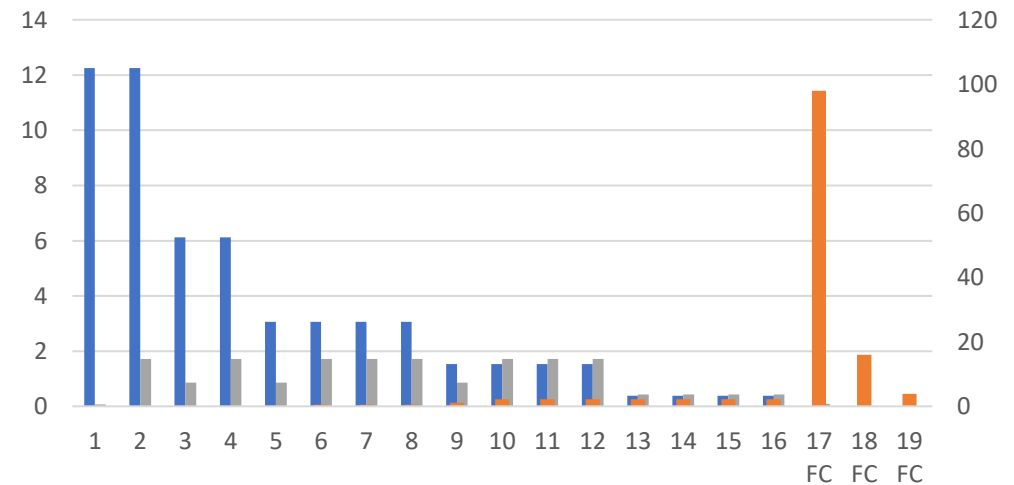


GoogLeNet

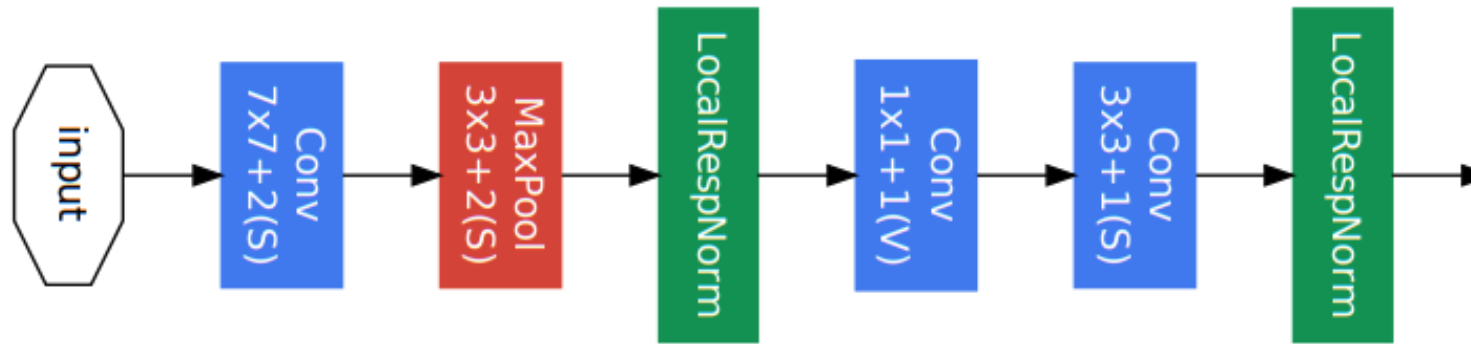
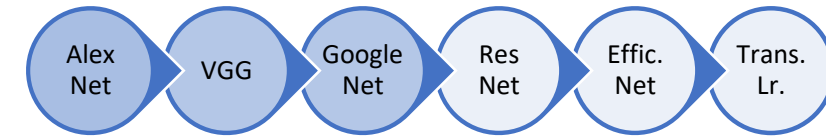


Want to avoid:

- High memory usage in the beginning
- High parameter count in the end

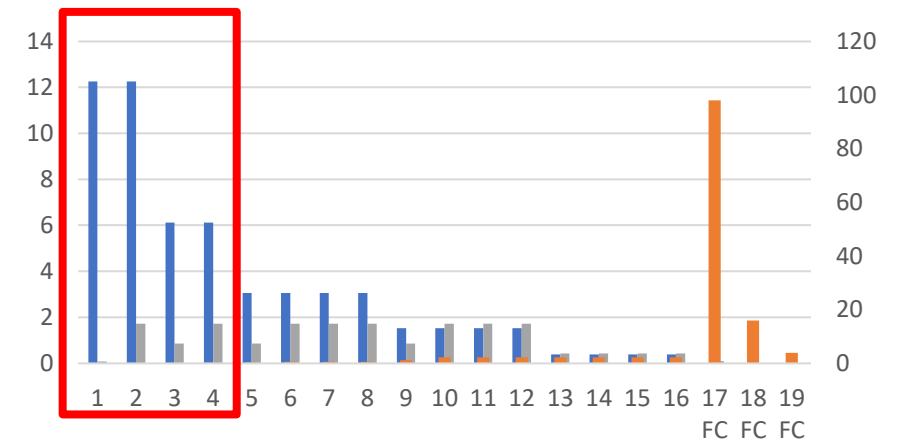


Reducing memory usage

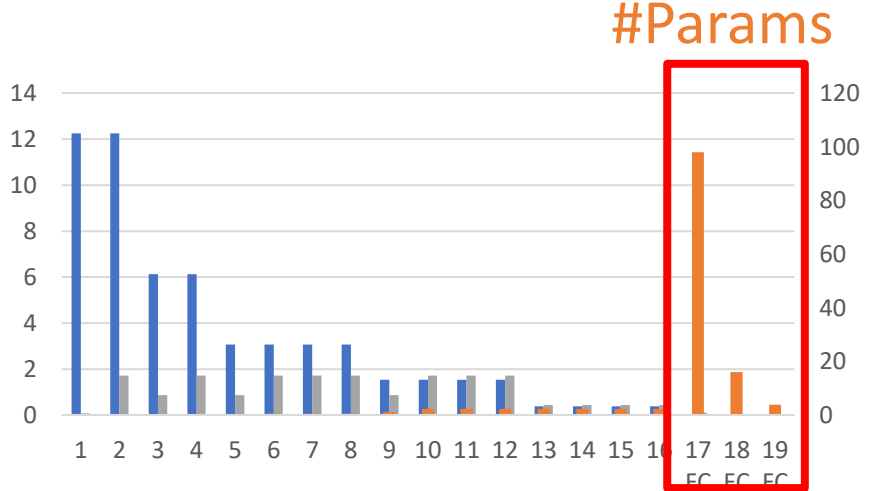
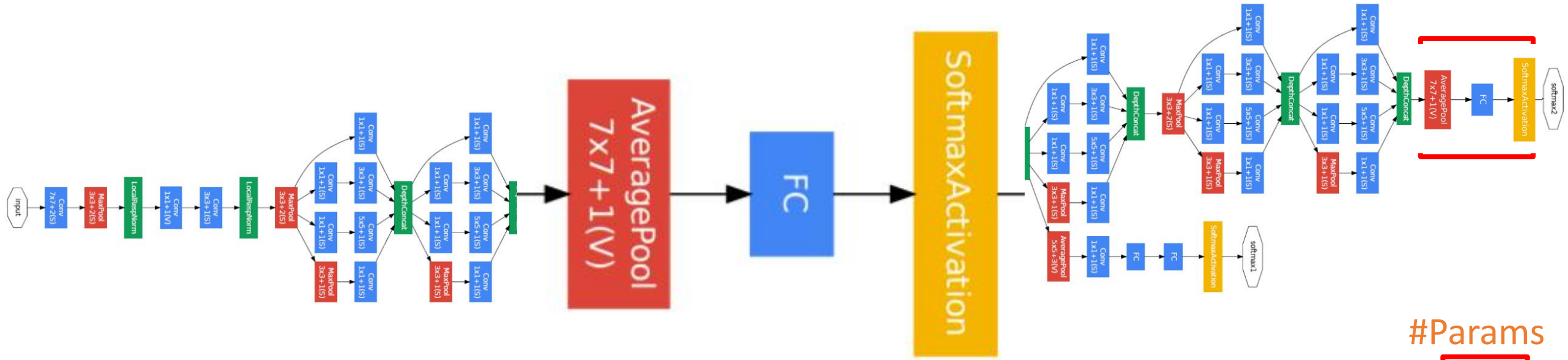
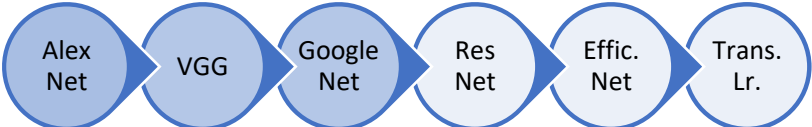


Downsample the input with 7*7 convolutions and MaxPool

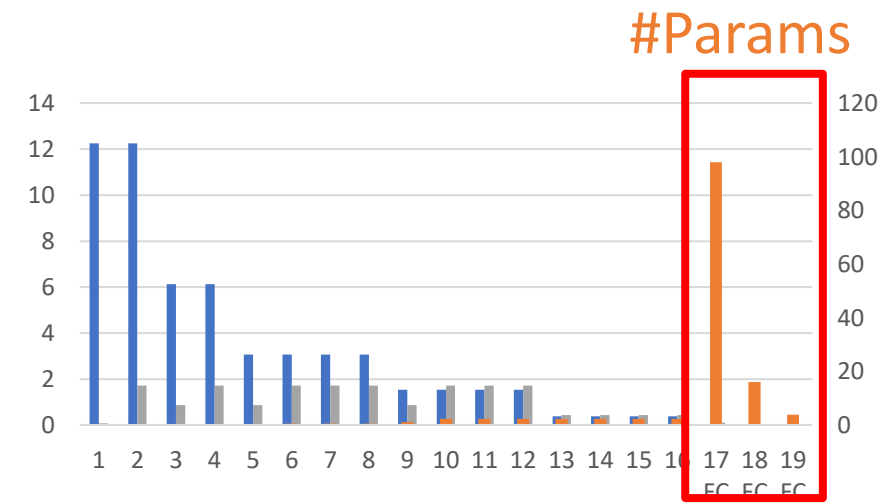
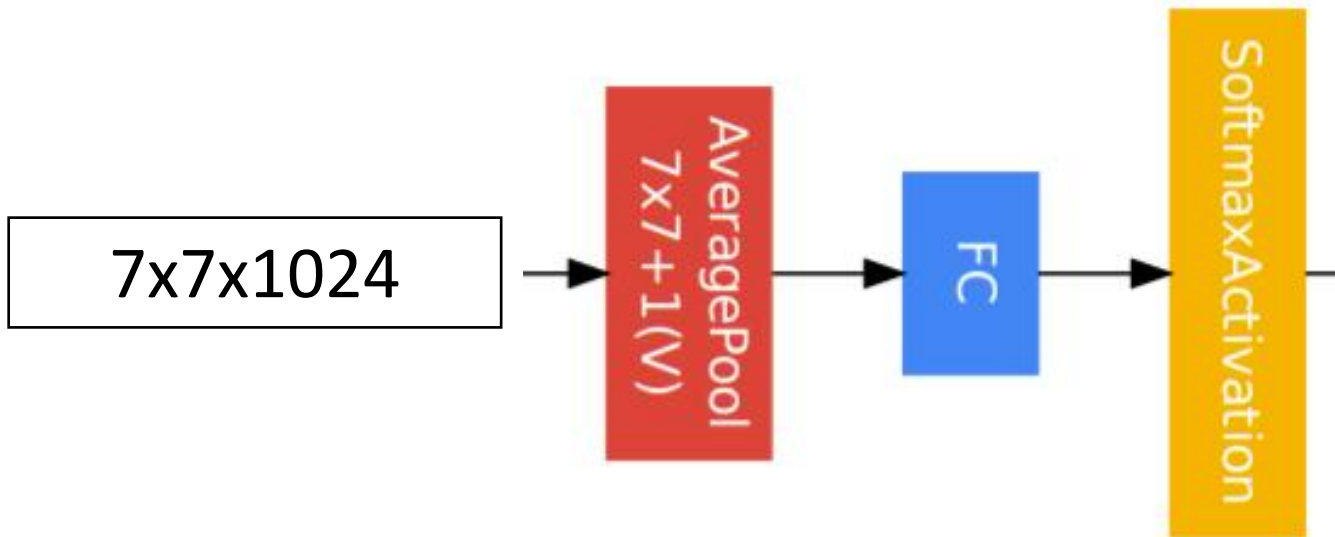
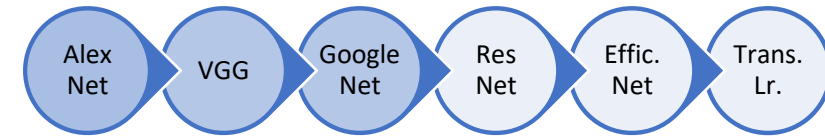
Memory



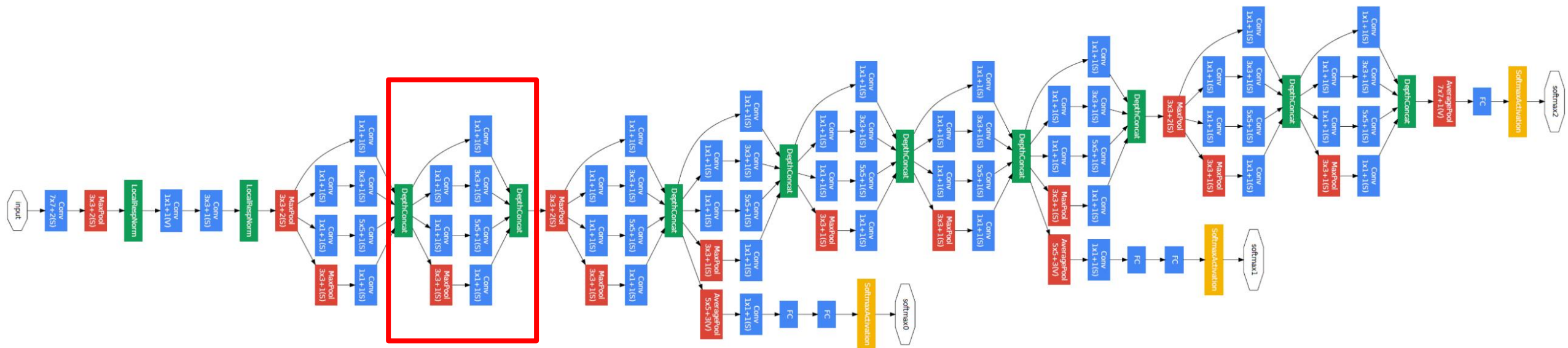
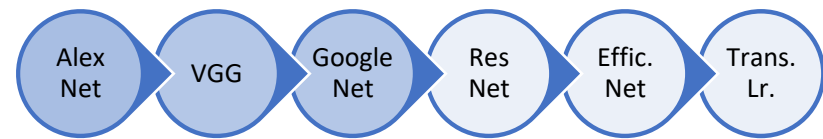
Reducing parameter count



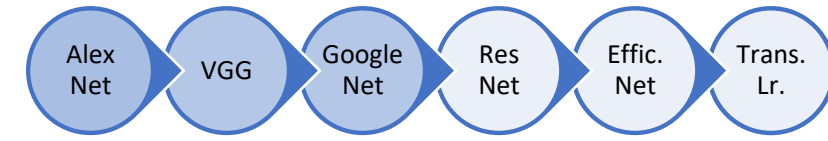
Global average pooling



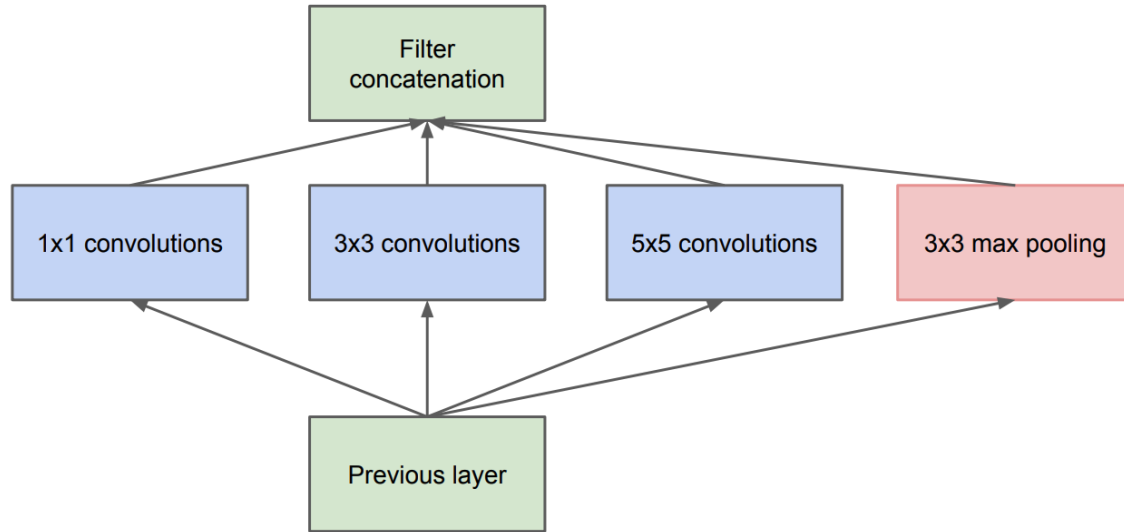
GoogLeNet



GoogLeNet: Inception blocks

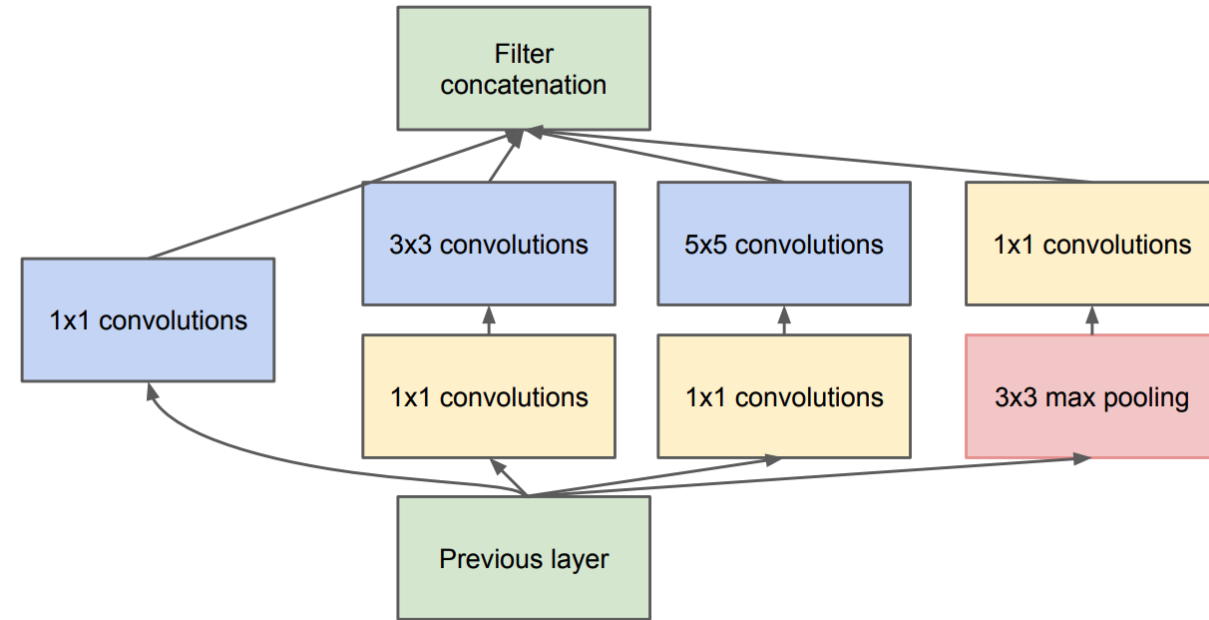


Inception block



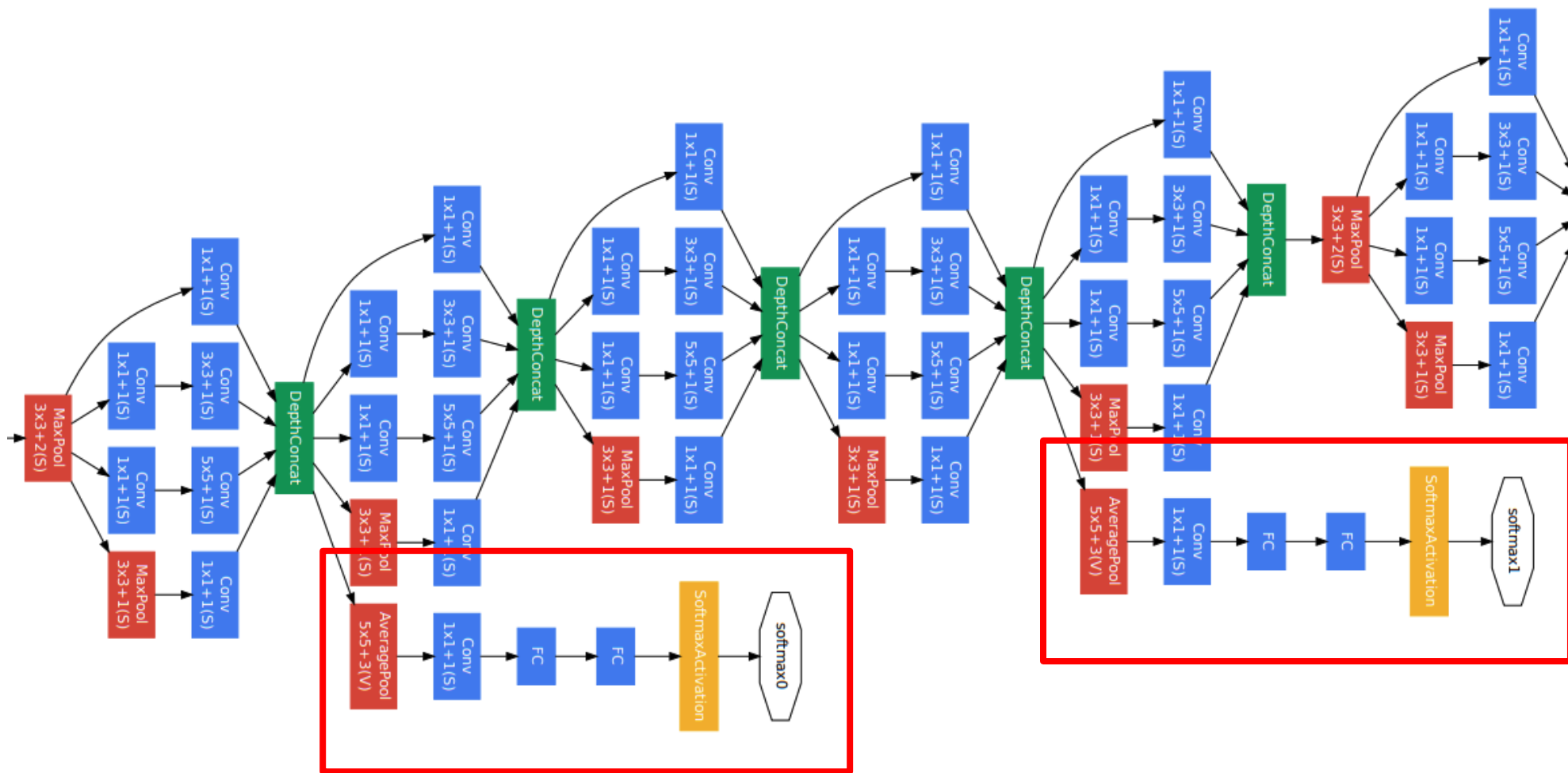
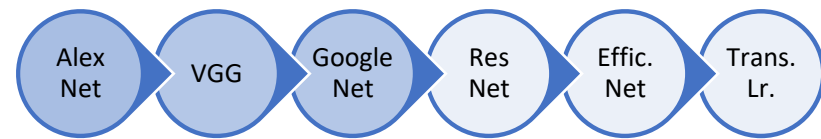
Problem: too computationally expensive

Inception block with dimension reductions

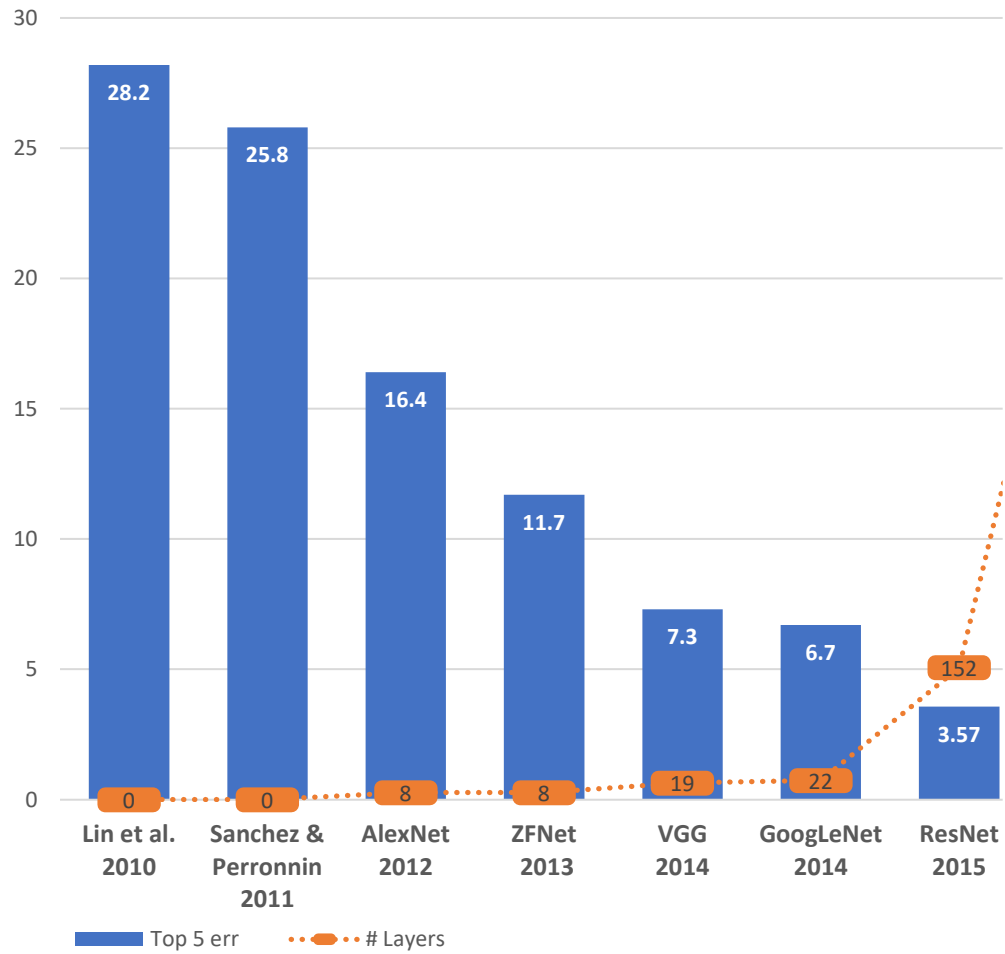
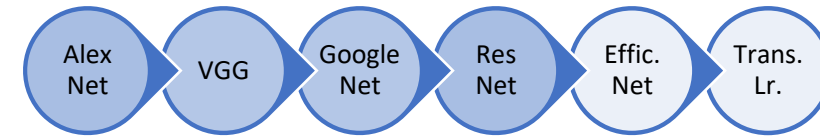


Solution: reduce # channels by 1*1 convolutions

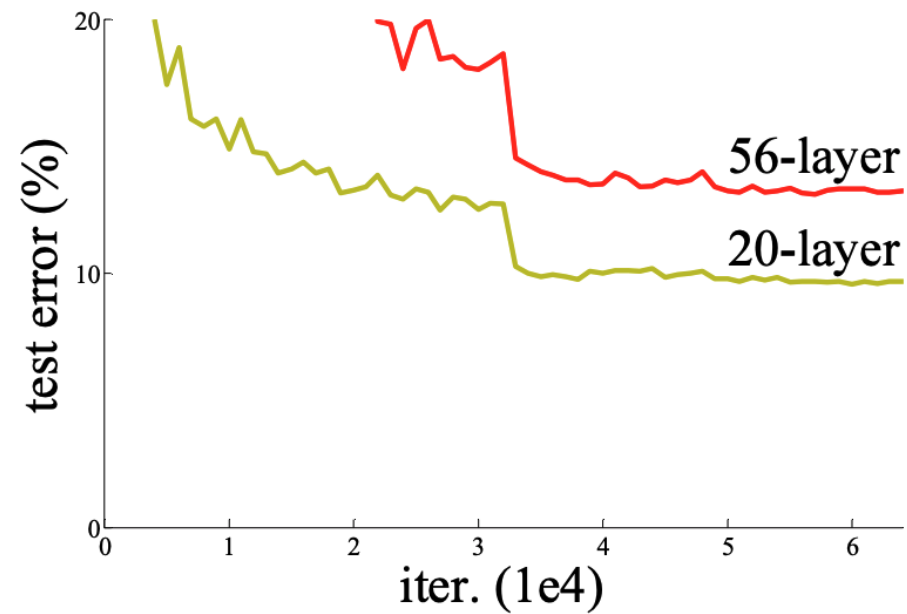
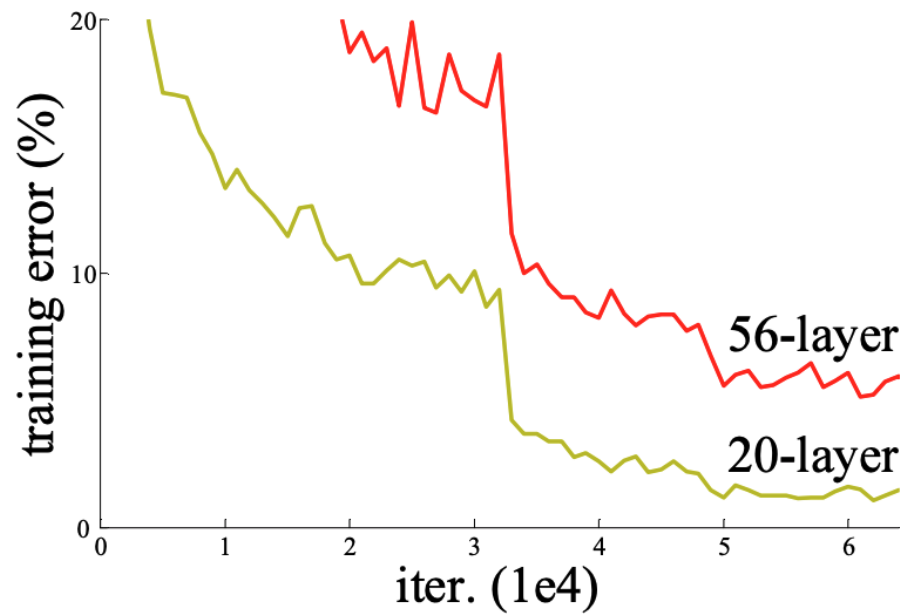
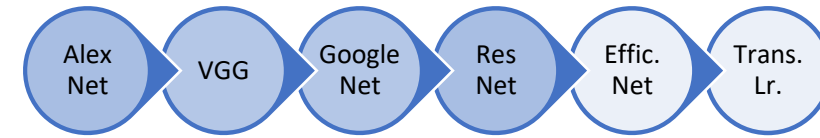
GoogLeNet



ResNet

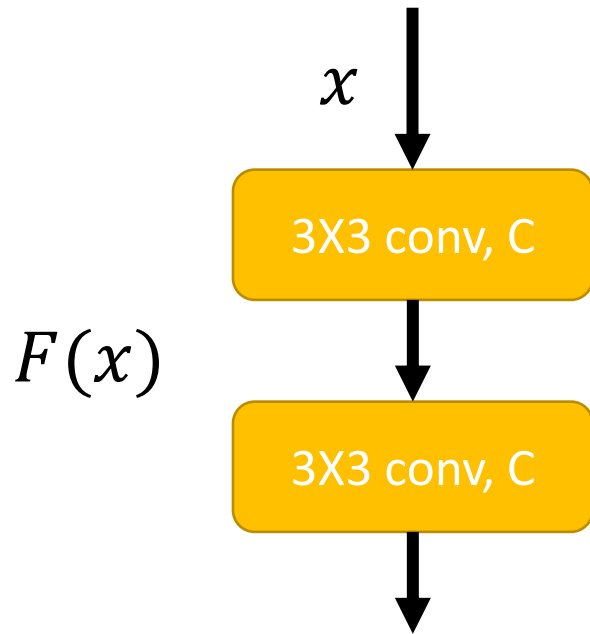
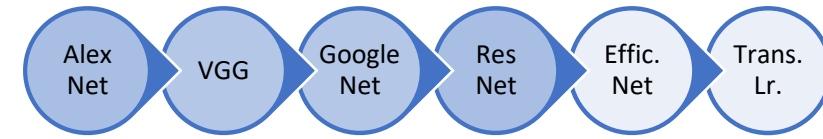


ResNet

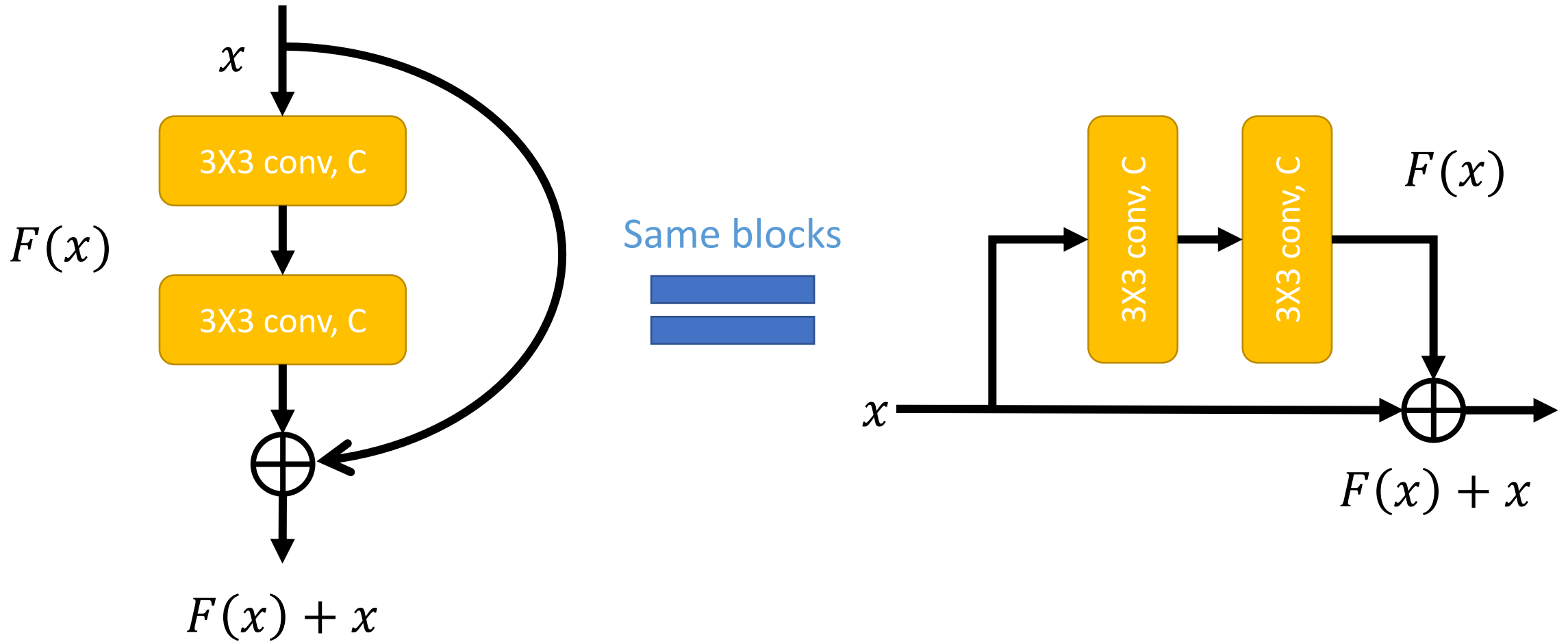
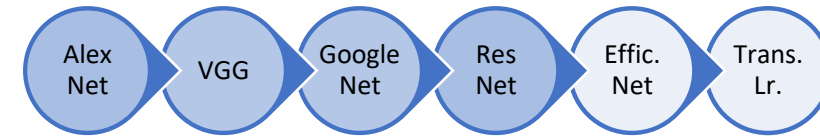


Deeper models are harder to optimize - vanishing gradients problem

Residual Building Block

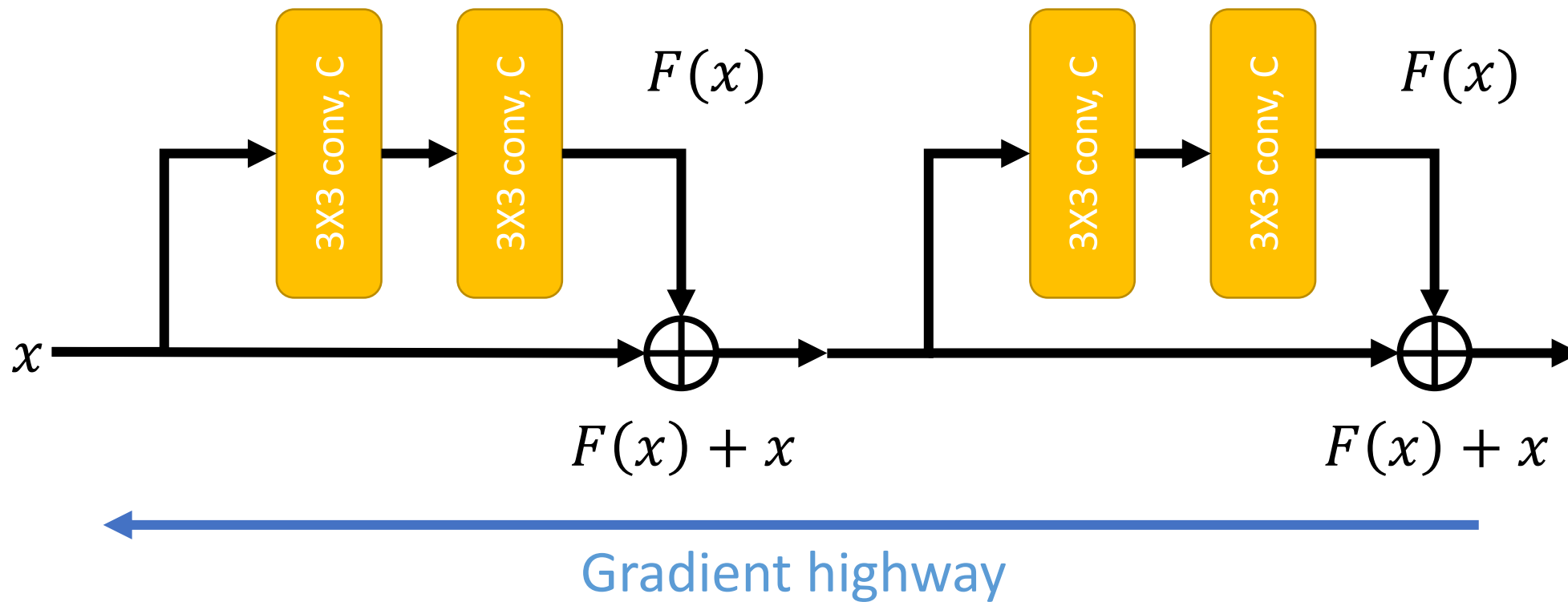
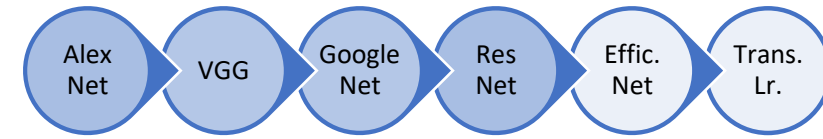


Residual Building Block

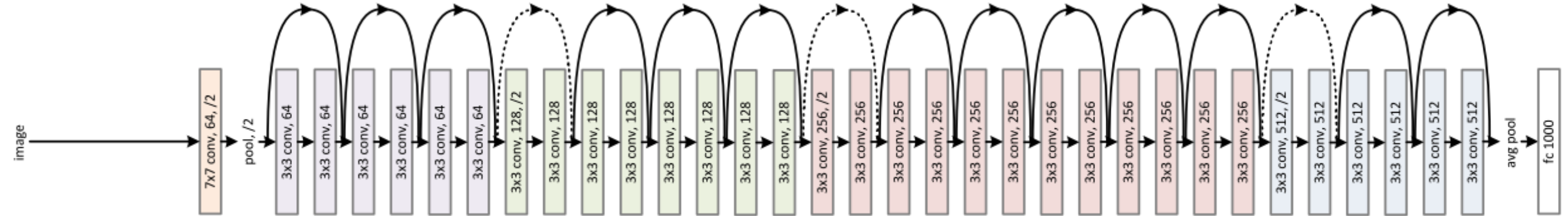
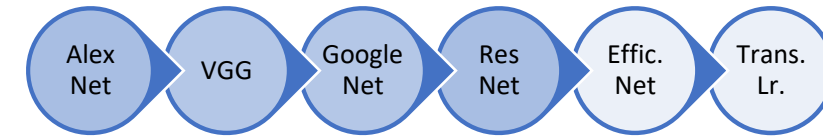


Use network layers to learn a residual mapping

ResNet Gradient propagation

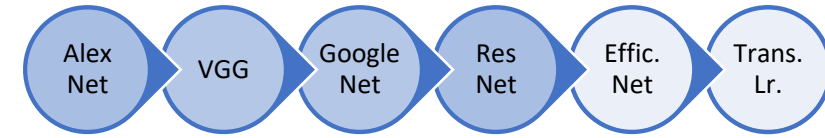


ResNet Architecture



- Residual connections help gradient propagation to initial layers
- No fully connected layers at the end
- Rapid downsampling of input

More architectures



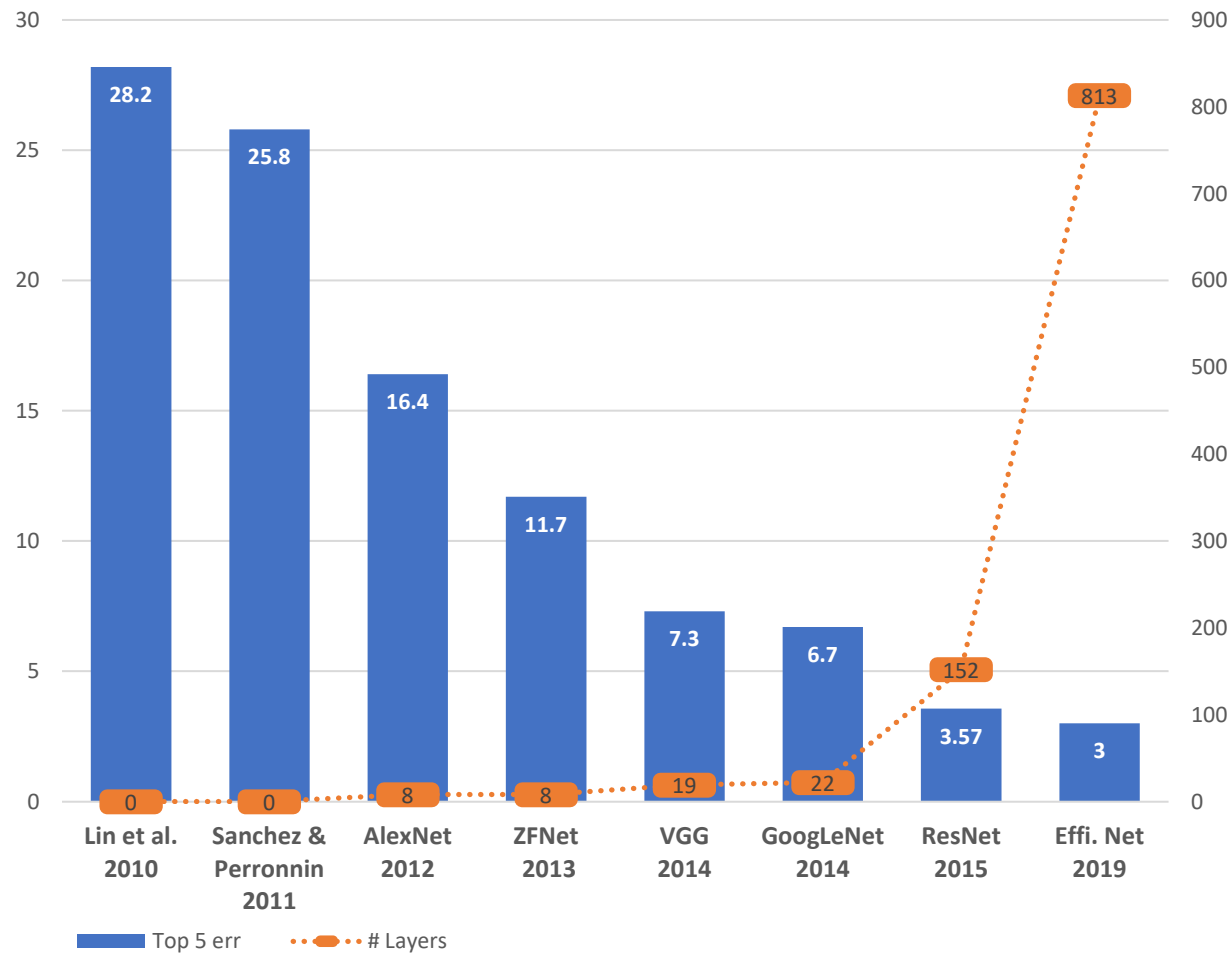
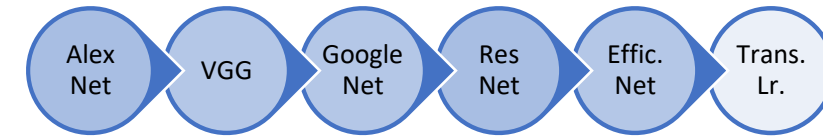
Accuracy

- DenseNet, Huang et al. 2017
- ResNext, Xie et al. 2017
- SENet, Hu et al. 2018

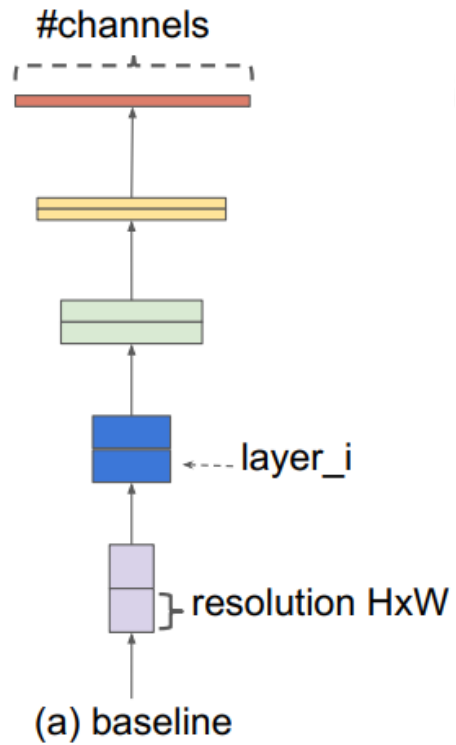
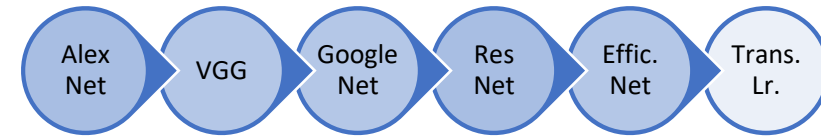
Efficiency

- MobileNet, Howard et al. 2017
- ShuffleNet, Zhang et al. 2018
- Neural Architecture Search, Zoph and Le. 2017

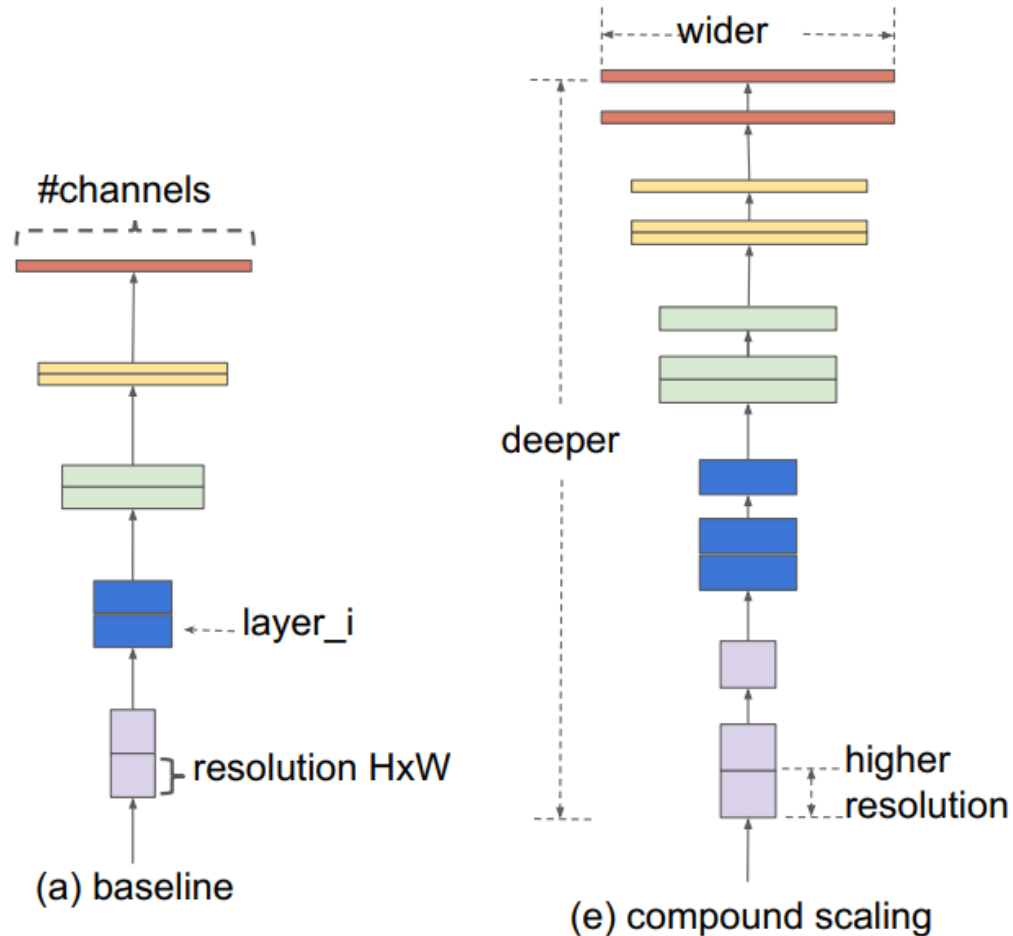
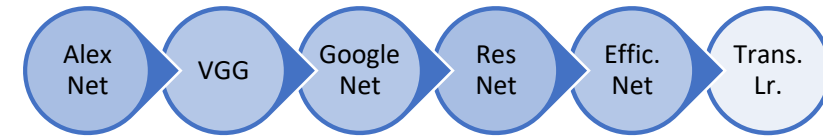
EfficientNet



EfficientNet



EfficientNet



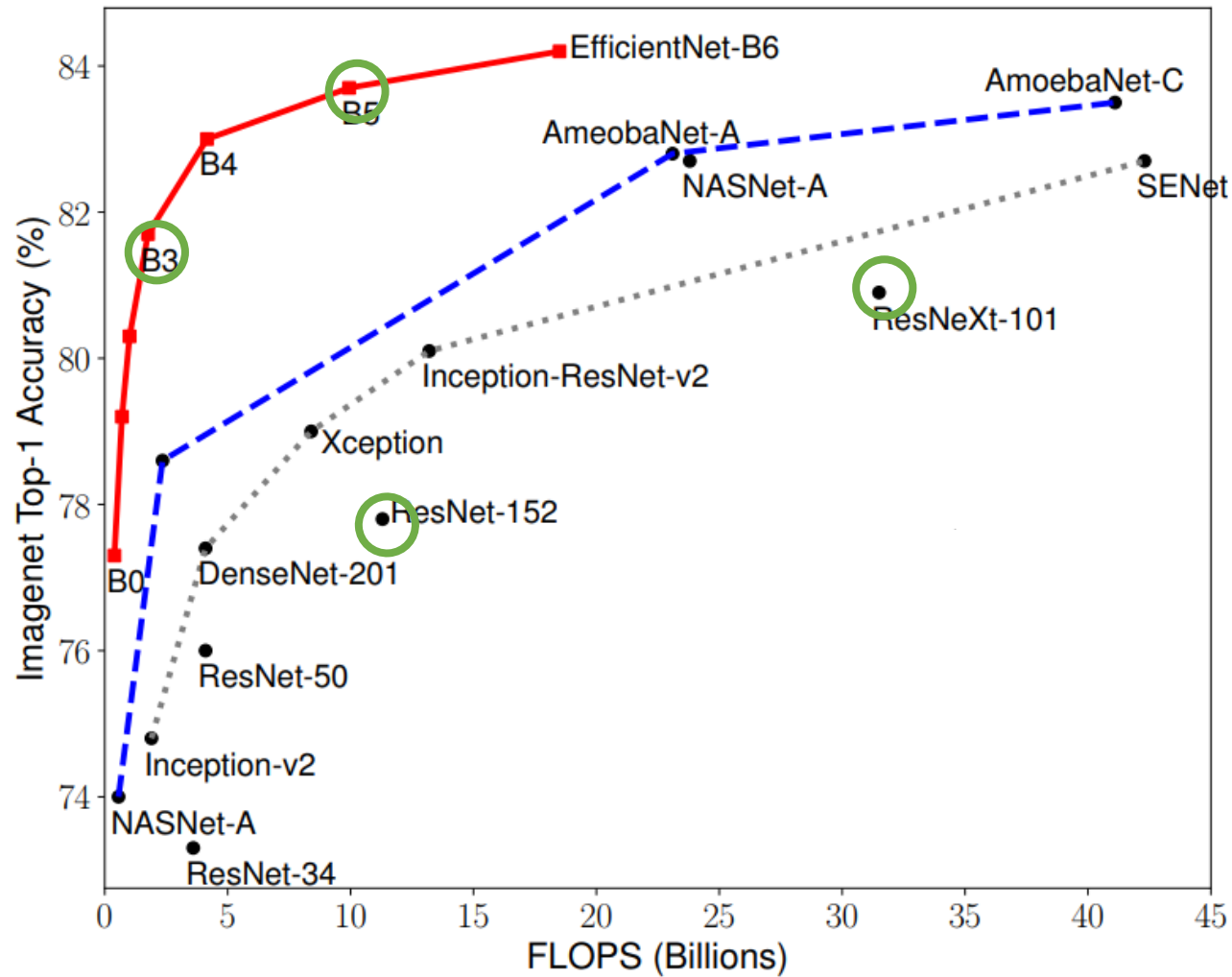
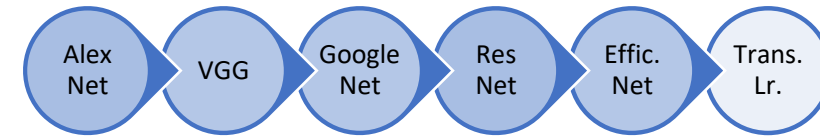
If we want to use 2^N more resources (FLOPs):

- Increase depth by α^N
- Increase width by β^N
- Increase input resolution by γ^N

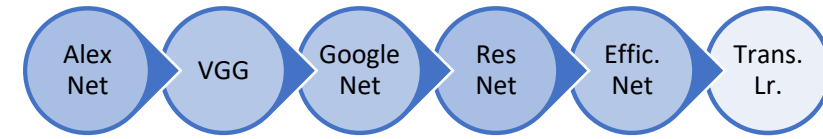
α, β, γ – constant scaling coefficients determined by a small grid search on the original small model

(they found $\alpha = 1.2, \beta = 1.1, \gamma = 1.15$)

Results

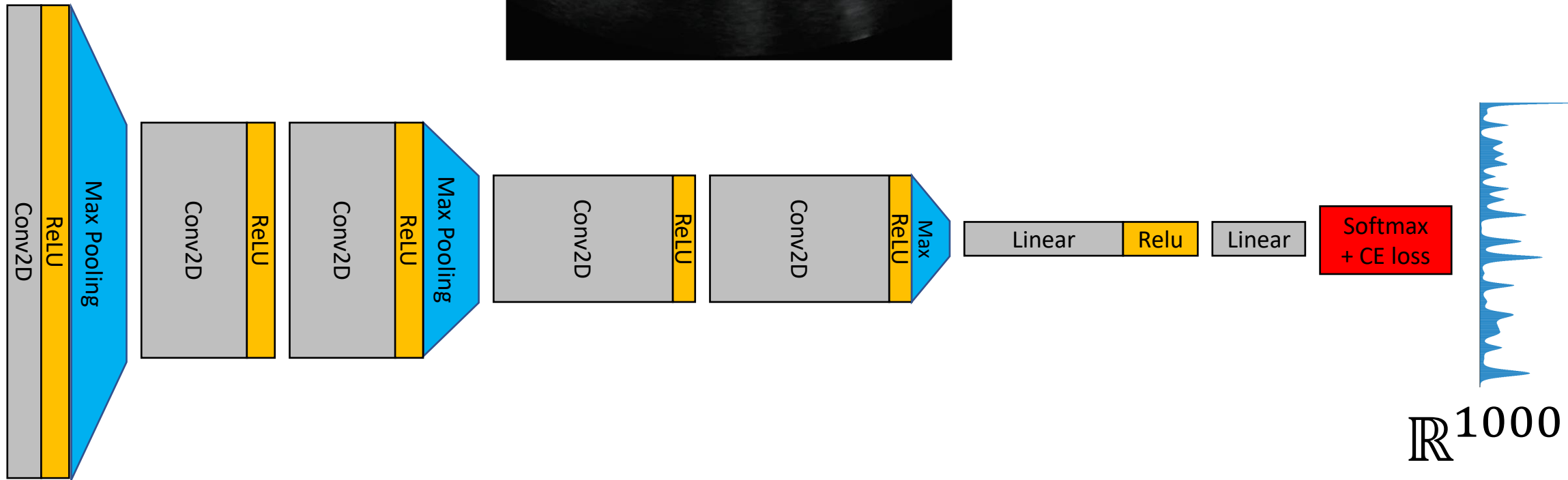
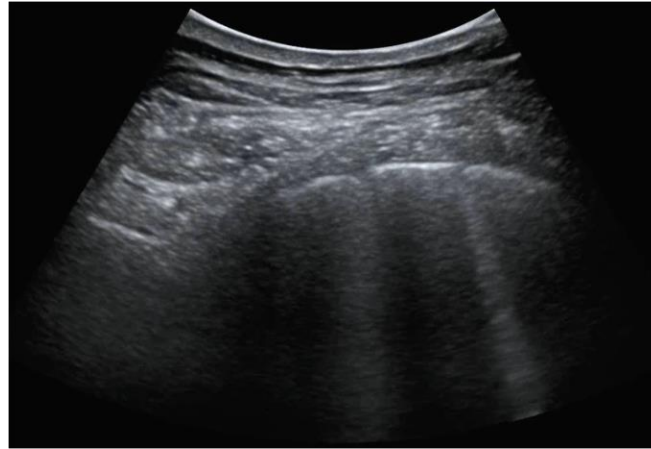
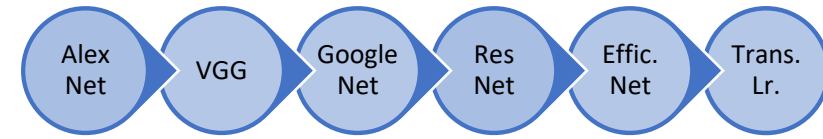


Architectures summary

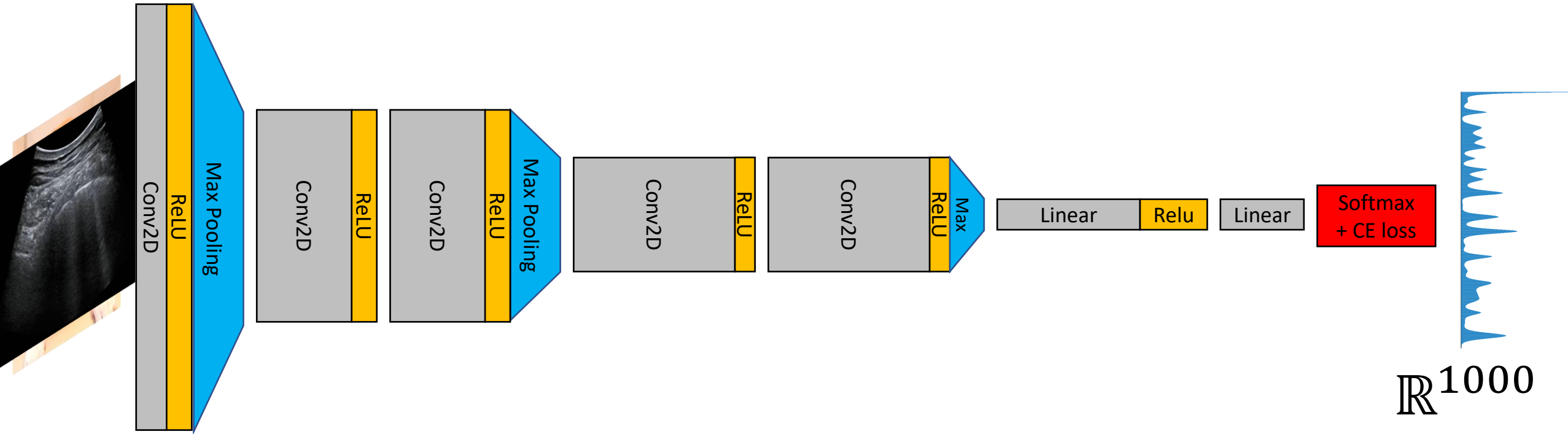
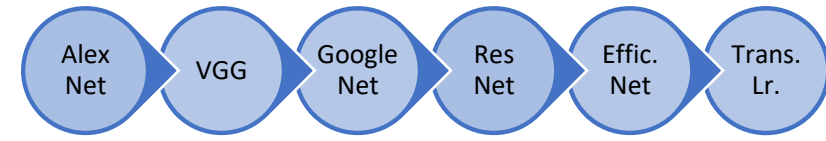


- Design your network according to your task and resources
- Take into consideration
 - # Parameters
 - # FLOPs
 - Memory Size
- Use skip connections :)
- Use **existing architectures** when possible
- Use **pre-trained models** when possible

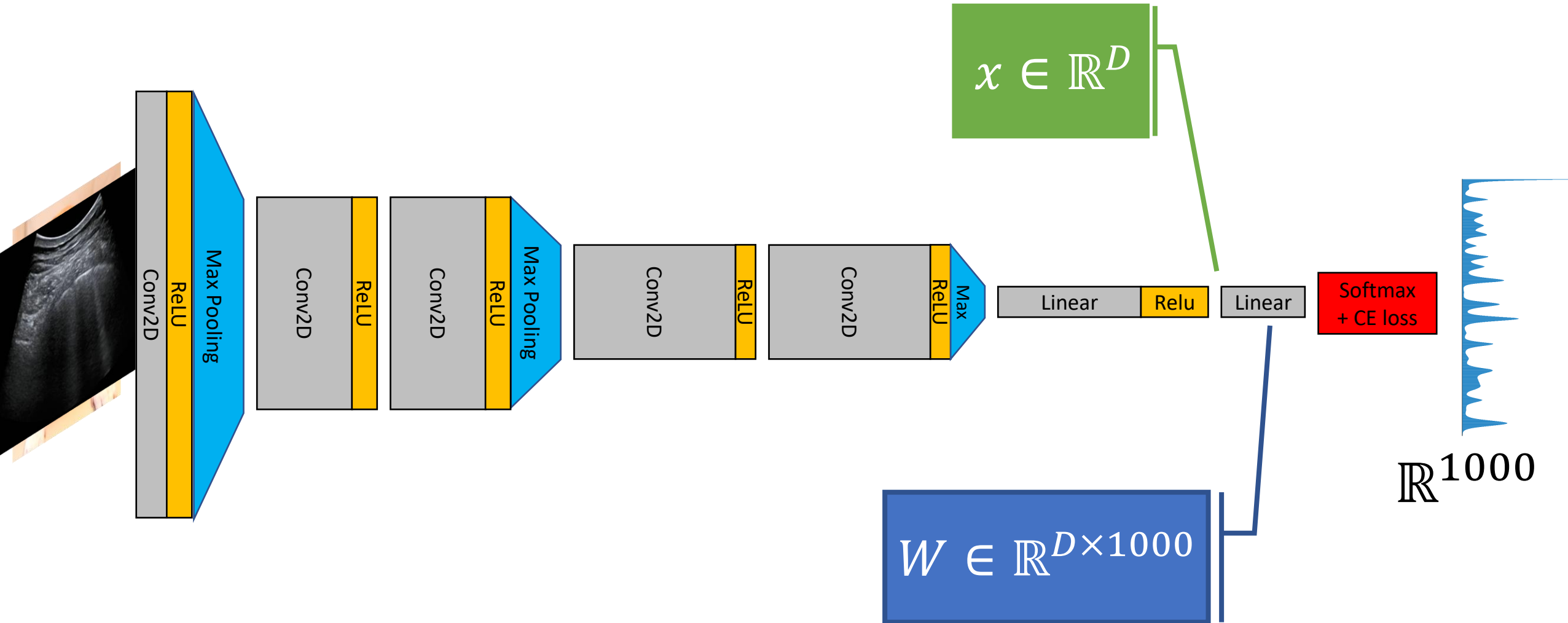
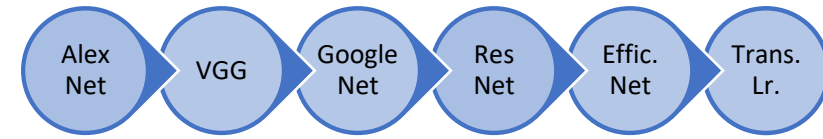
Transfer Learning



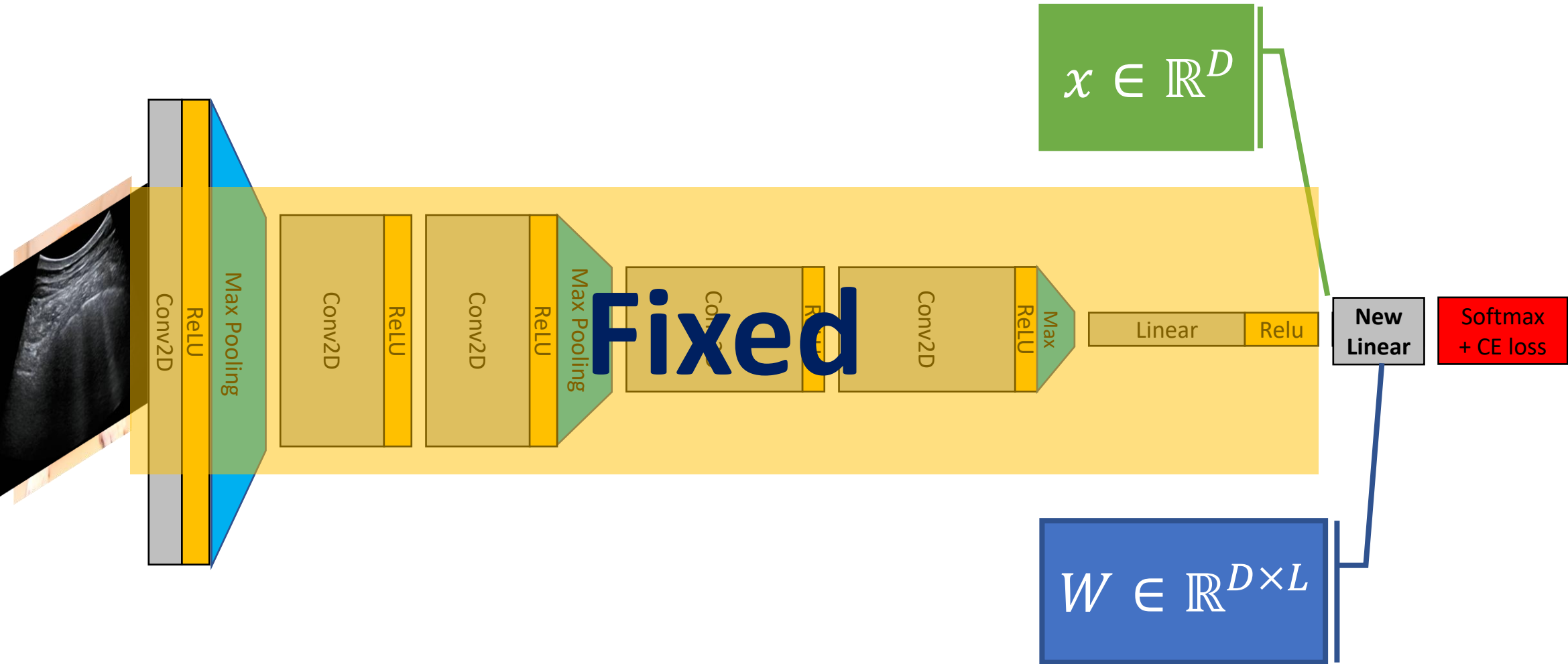
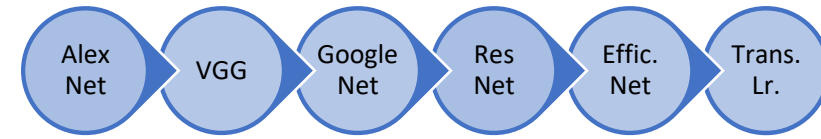
Transfer Learning



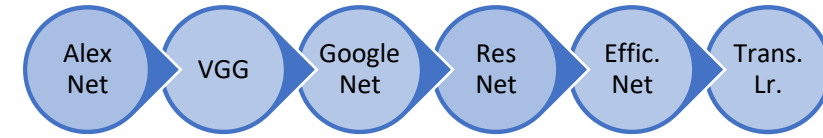
Transfer Learning



Transfer Learning



Transfer Learning



- In practice:

```
from torchvision.models import resnet50, ResNet50_Weights

# Old weights with accuracy 76.130%
resnet50(weights=ResNet50_Weights.IMAGENET1K_V1)

# New weights with accuracy 80.858%
resnet50(weights=ResNet50_Weights.IMAGENET1K_V2)

# Best available weights (currently alias for IMAGENET1K_V2)
# Note that these weights may change across versions
resnet50(weights=ResNet50_Weights.DEFAULT)

# Strings are also supported
resnet50(weights="IMAGENET1K_V2")

# No weights - random initialization
resnet50(weights=None)
```

Questions?