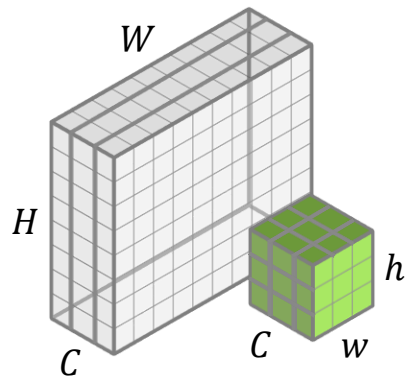


CNN Architectures

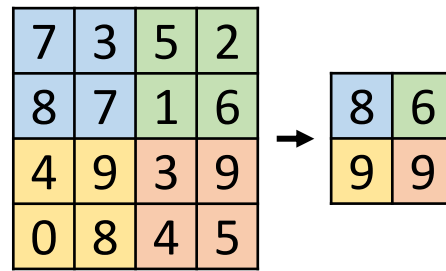
November 29th, 2022



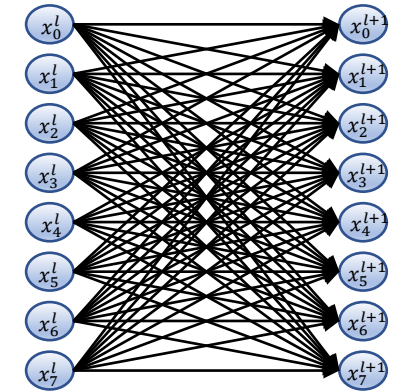
The ConvNet Building Blocks



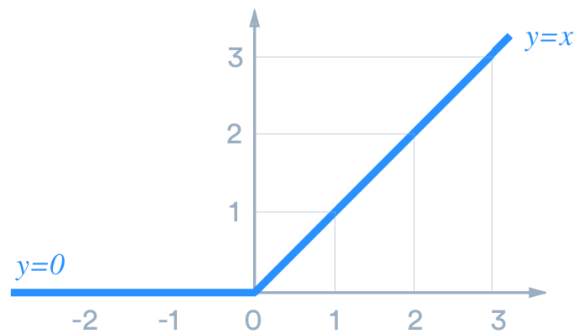
Conv Layer



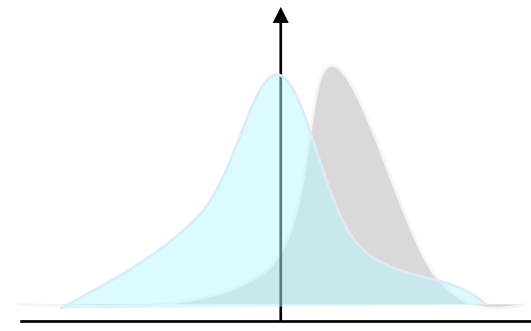
Max Pooling



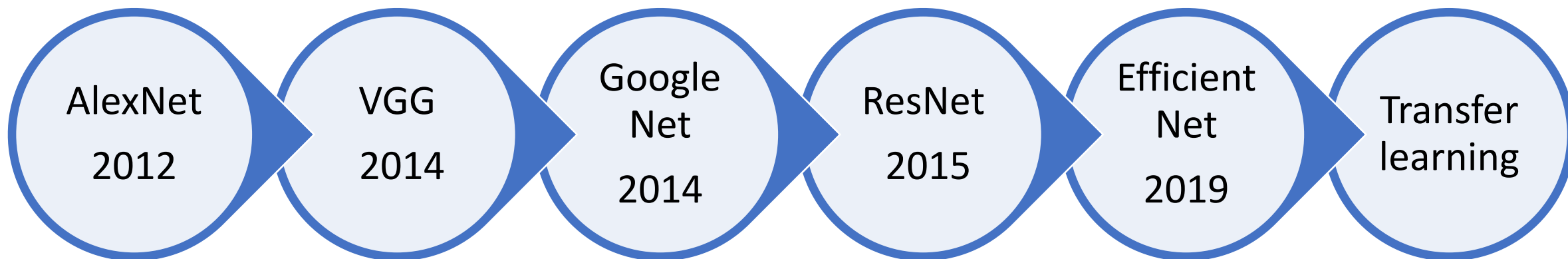
FC Layer



Activation Func.



Batch Norm

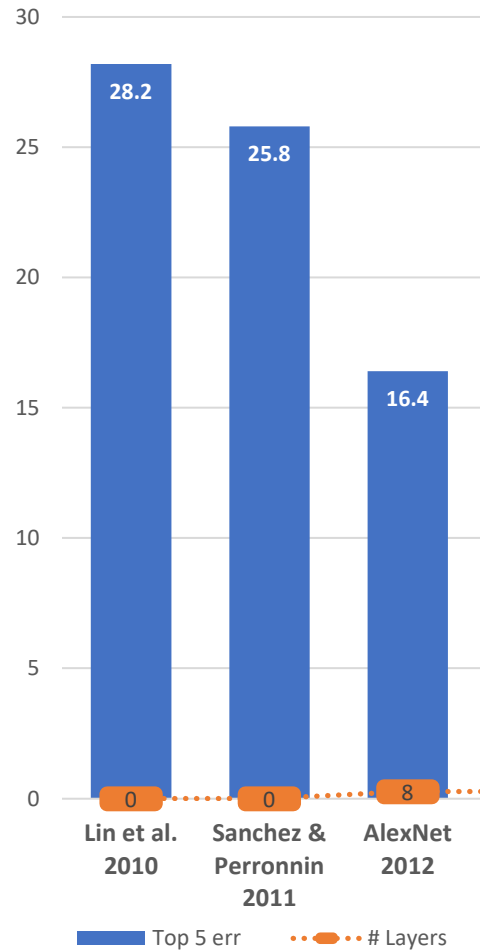
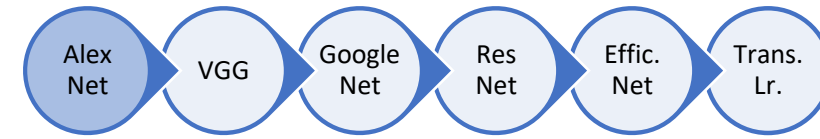




Visual Recognition Challenge (ILSVRC)

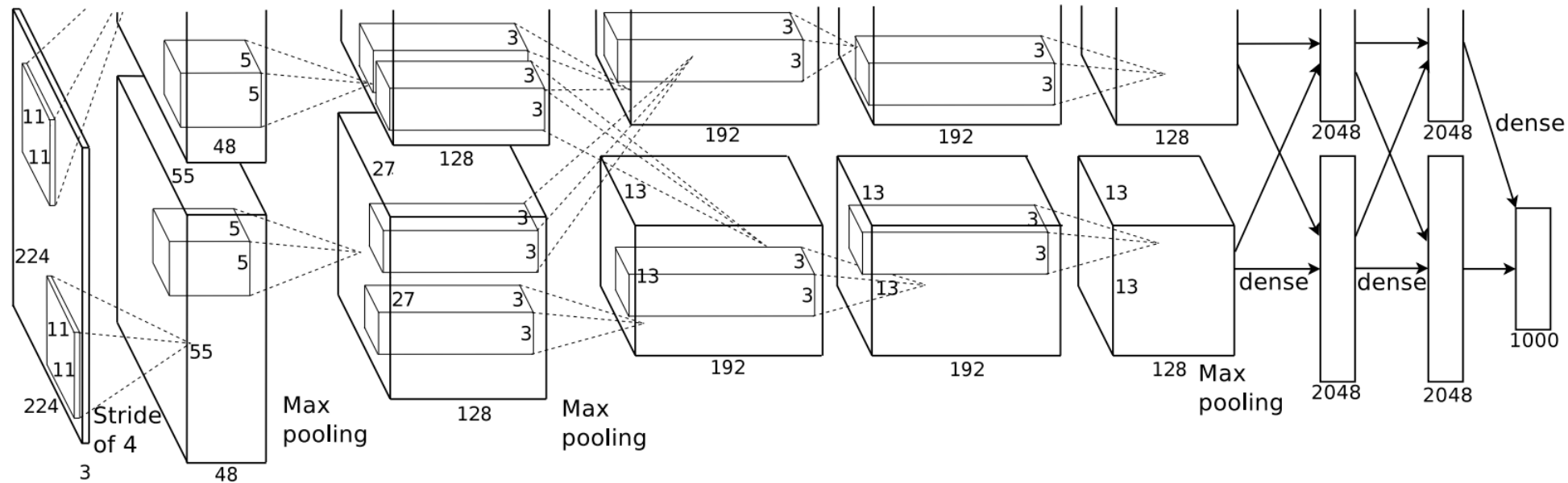
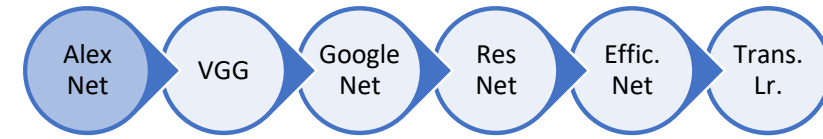
- 1.2M training images
- 100K testing images
- 1K categories

AlexNet



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

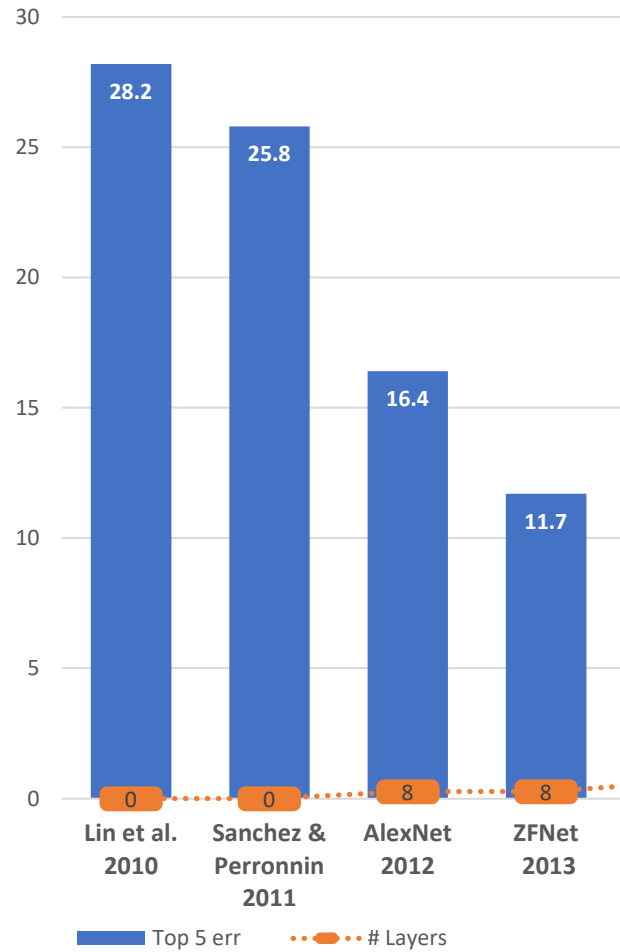
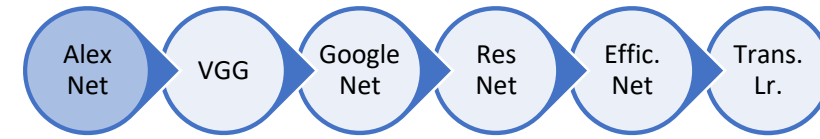
AlexNet



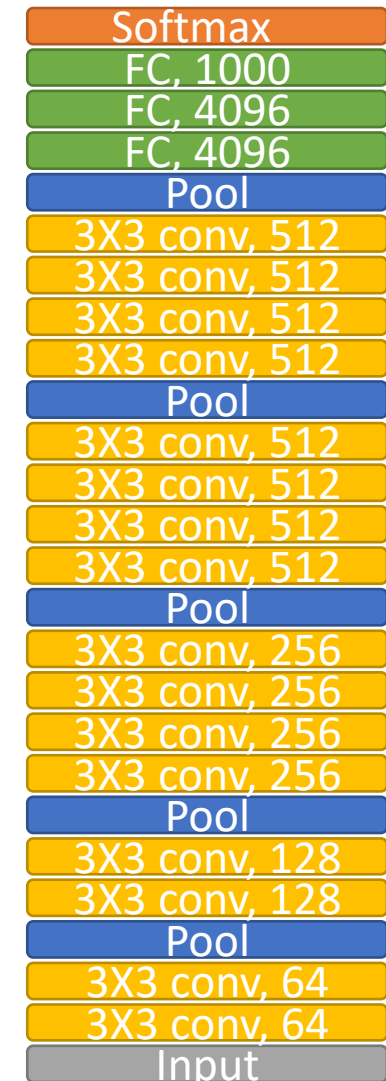
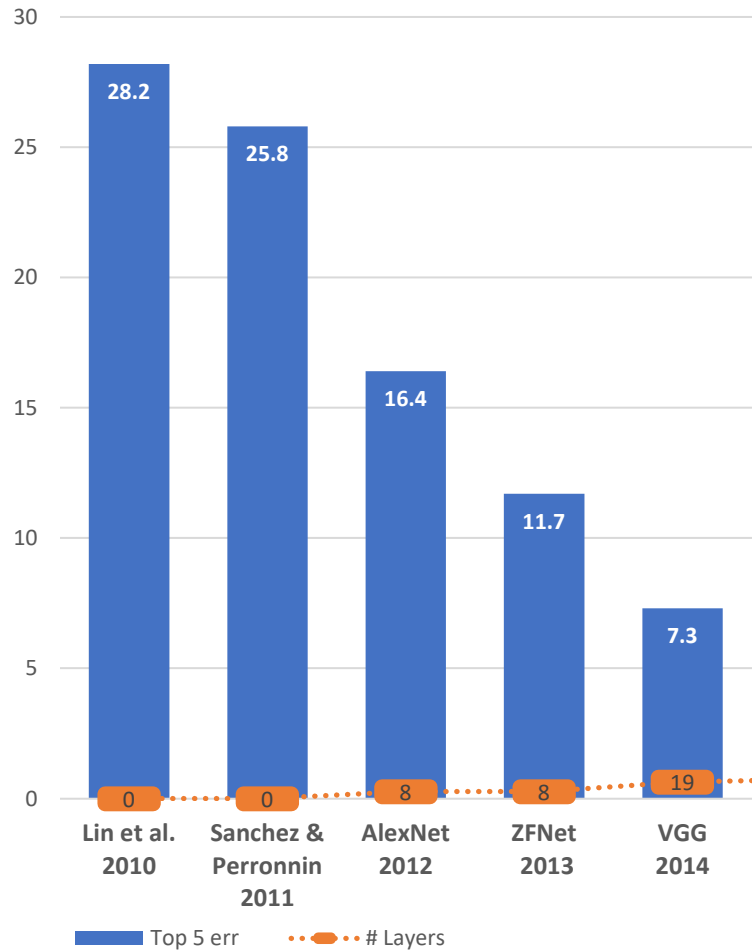
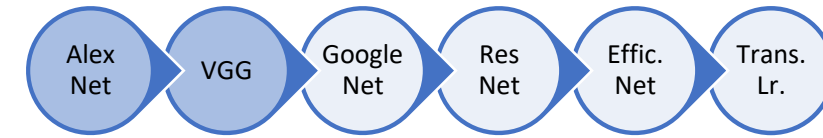
Details:

- ReLU activation
- Batch
- Max Pooling
- SGD Momentum
- Dropout

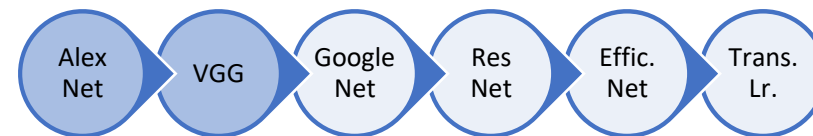
ZFNet



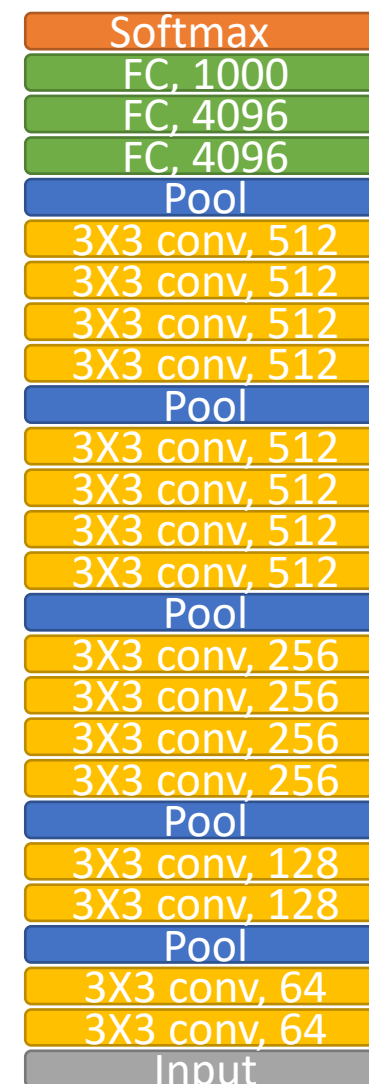
VGG



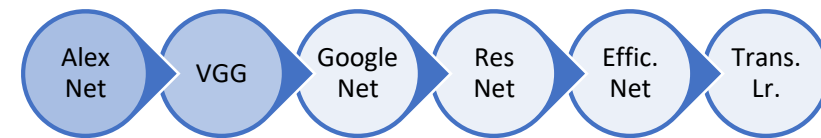
VGG



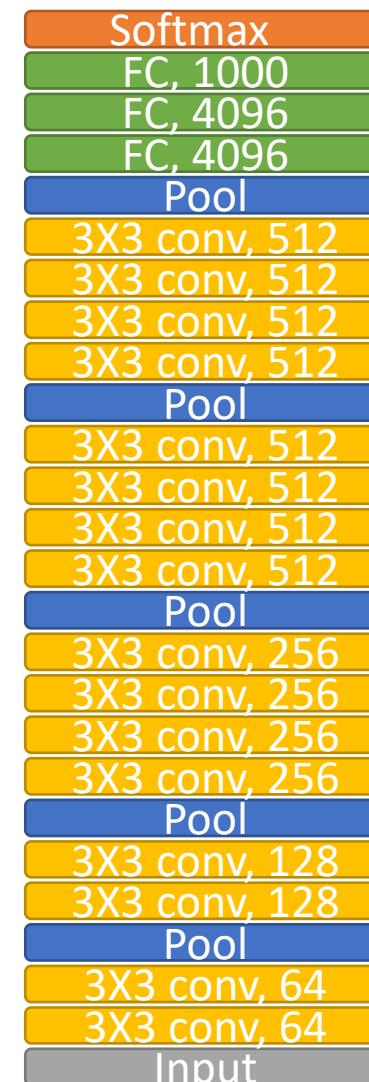
- All conv layers are 3x3, stride 1, pad 1
- All pooling layers are 2x2, stride 2
- After pooling, double the number of channels



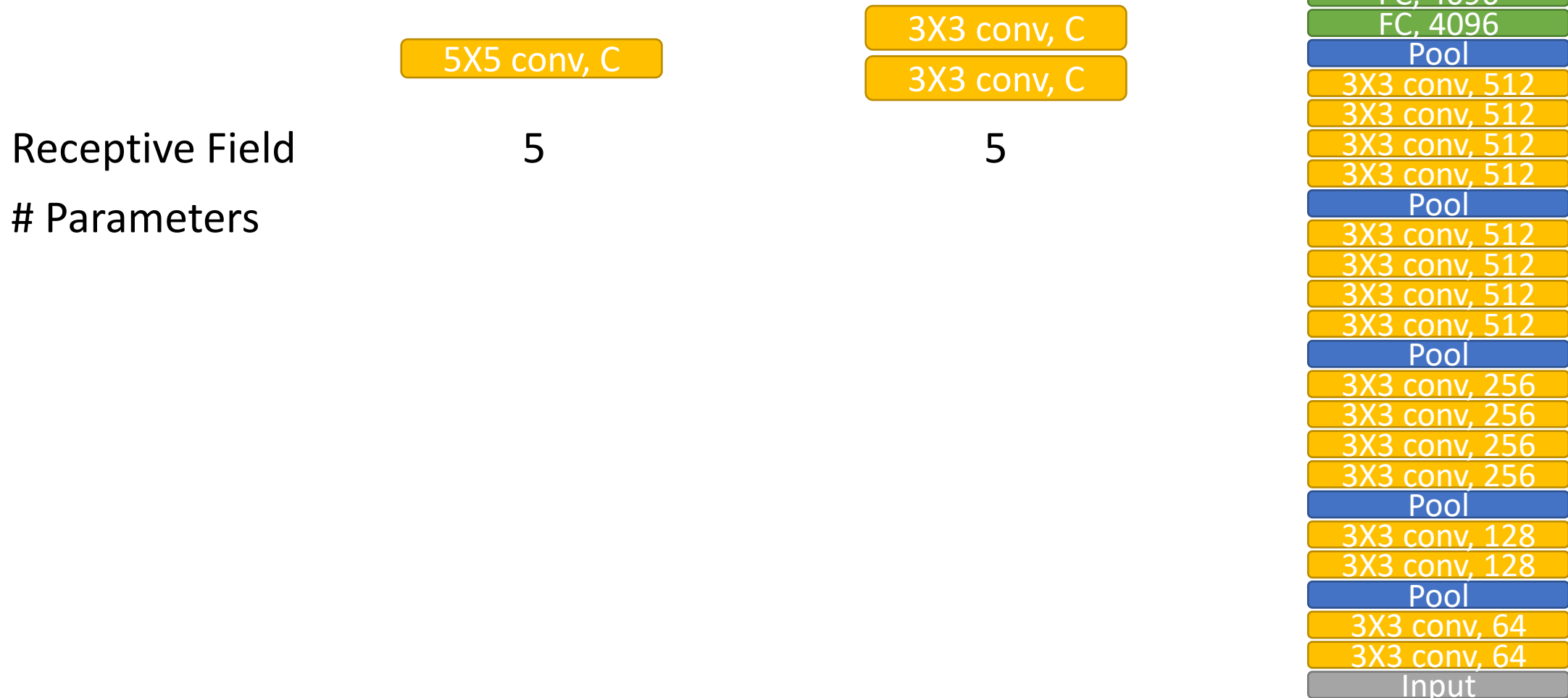
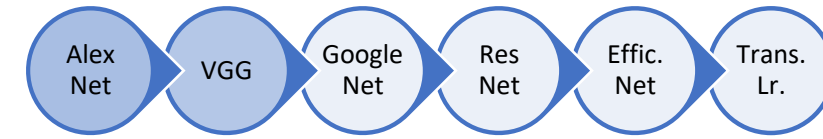
Why using 3x3 conv layers?



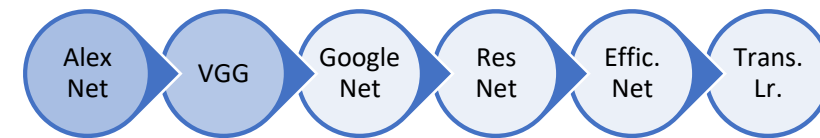
Receptive Field



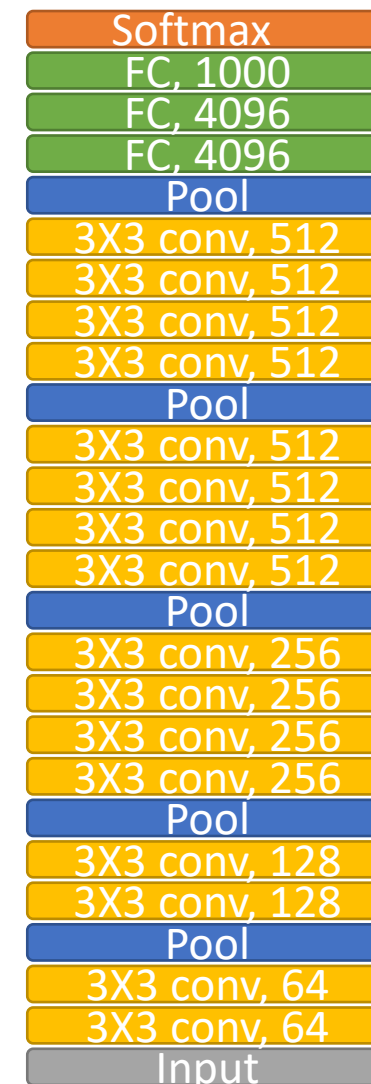
Why using 3x3 conv layers?



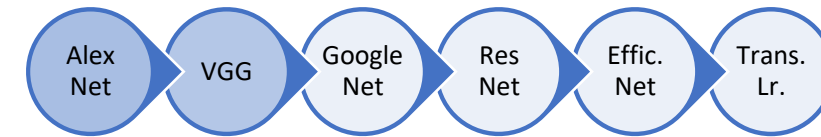
Why using 3x3 conv layers?



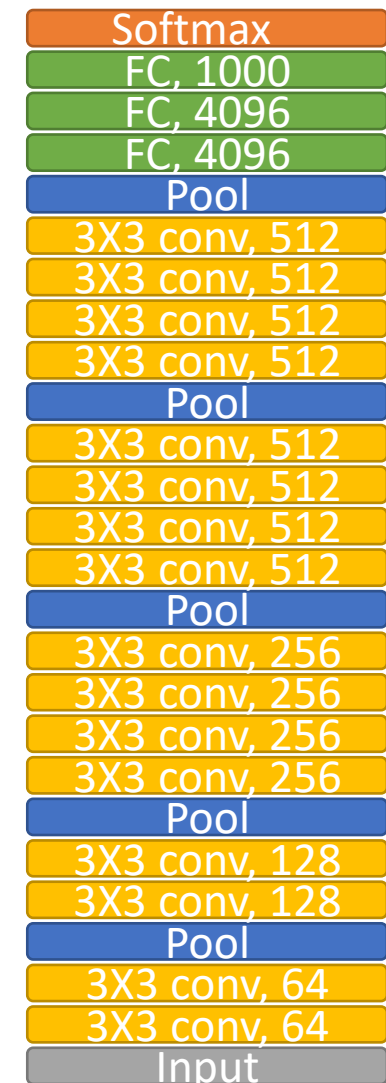
	5X5 conv, C	3X3 conv, C 3X3 conv, C
Receptive Field	5	5
# Parameters	$25C^2$	$9C^2 + 9C^2 = 18C^2$
# FLOPs		



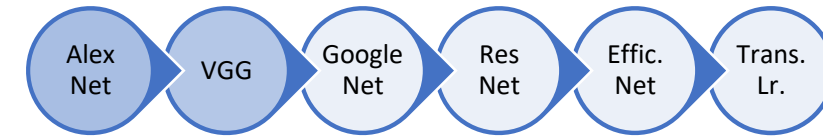
Why using 3x3 conv layers?



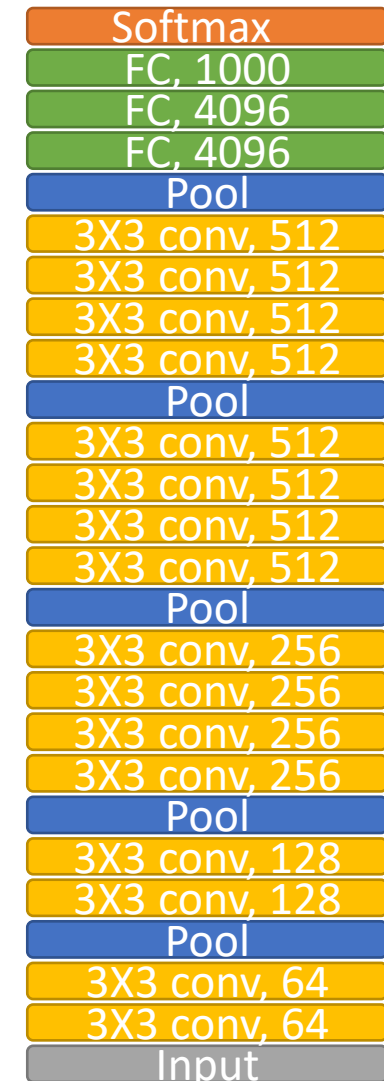
	5X5 conv, C	3X3 conv, C 3X3 conv, C
Receptive Field	5	5
# Parameters	$25C^2$	$9C^2 + 9C^2 = 18C^2$
# FLOPs	$25C^2HW$	$18C^2HW$



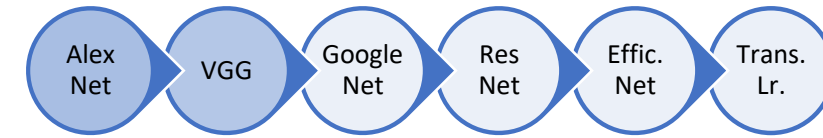
Why using 3x3 conv layers?



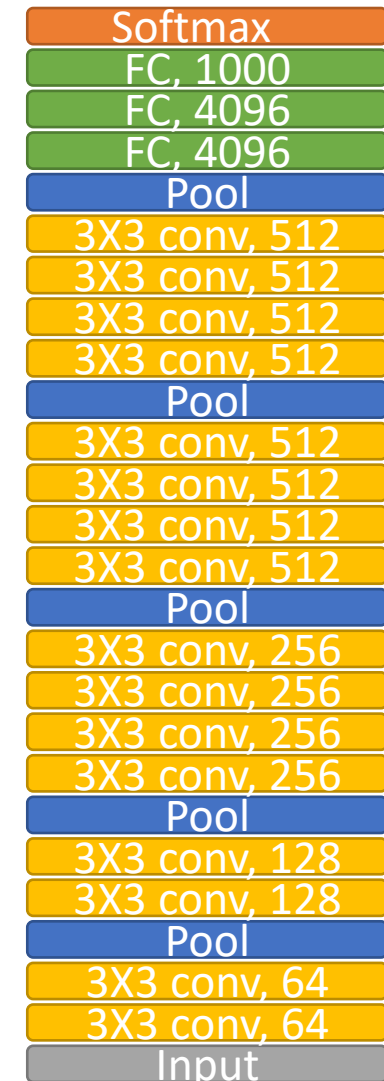
	5X5 conv, C	3X3 conv, C 3X3 conv, C
Receptive Field	5	5
# Parameters	$25C^2$	$9C^2 + 9C^2 = 18C^2$
# FLOPs	$25C^2HW$	$18C^2HW$
Memory Size (Output)		



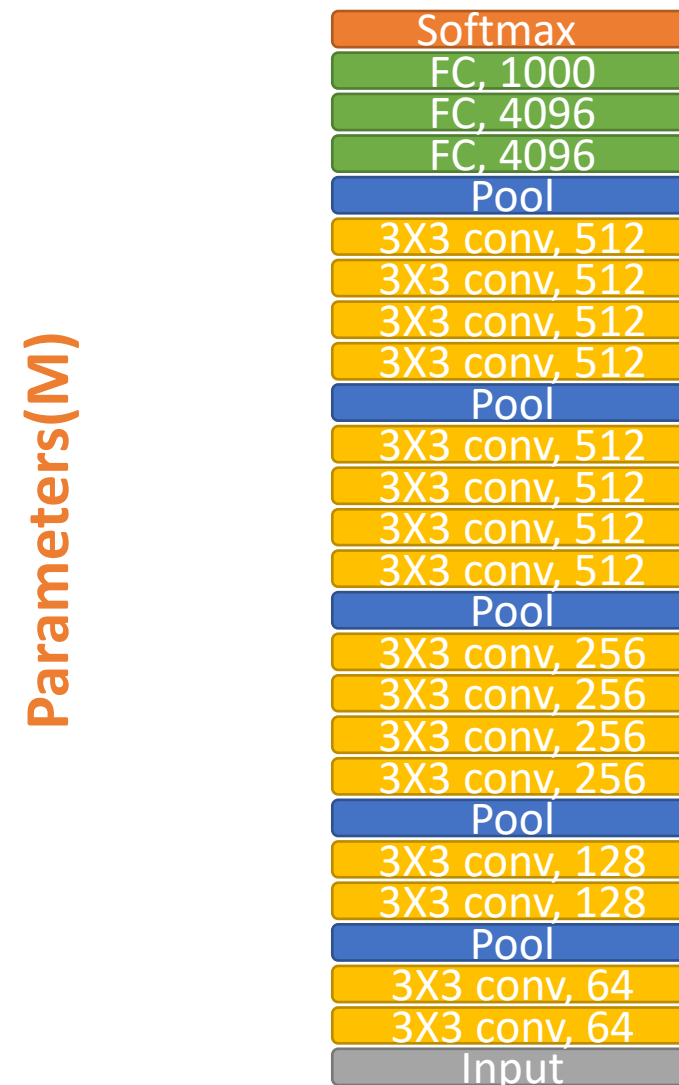
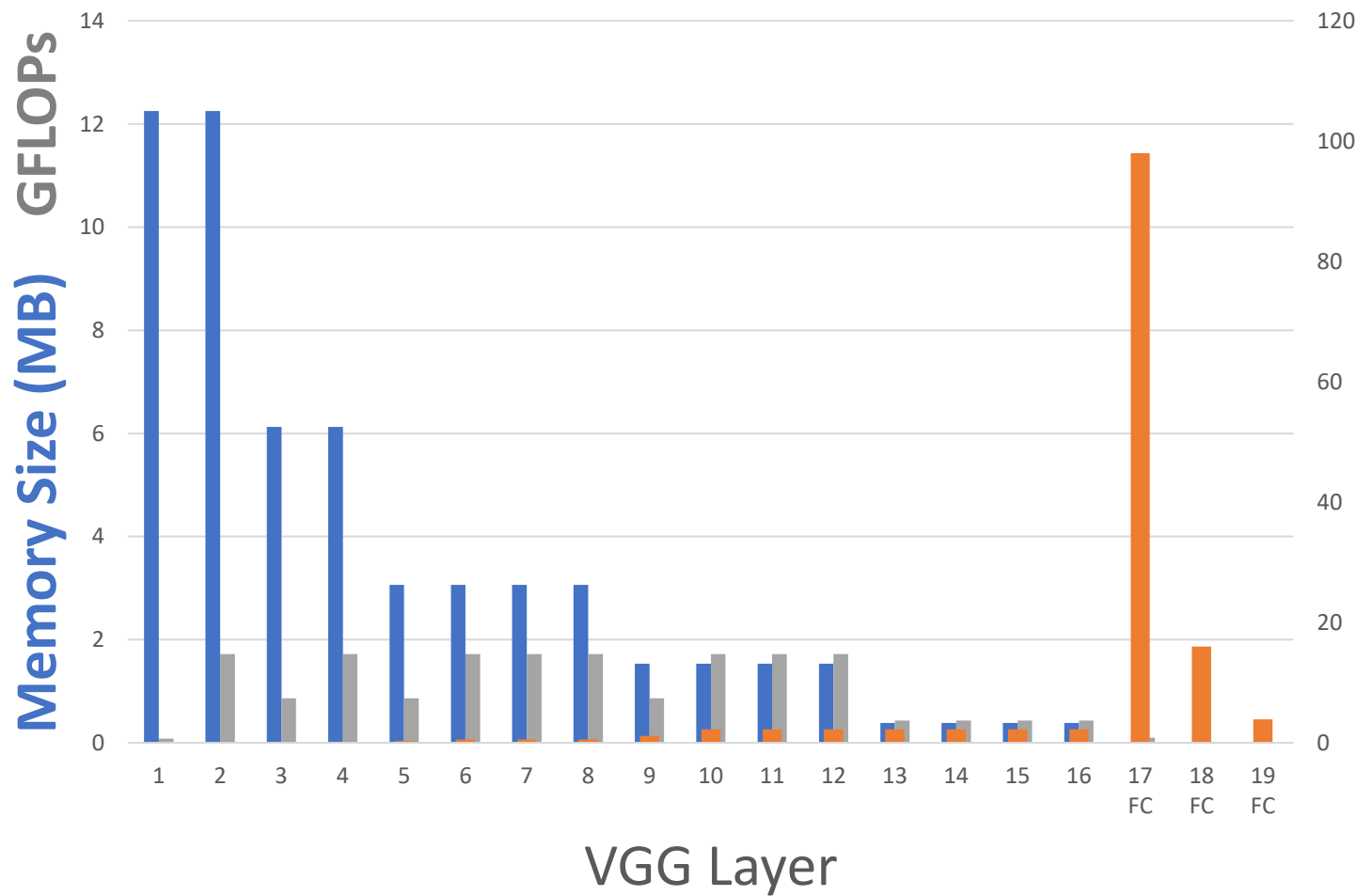
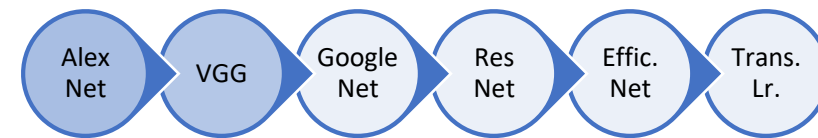
Why using 3x3 conv layers?



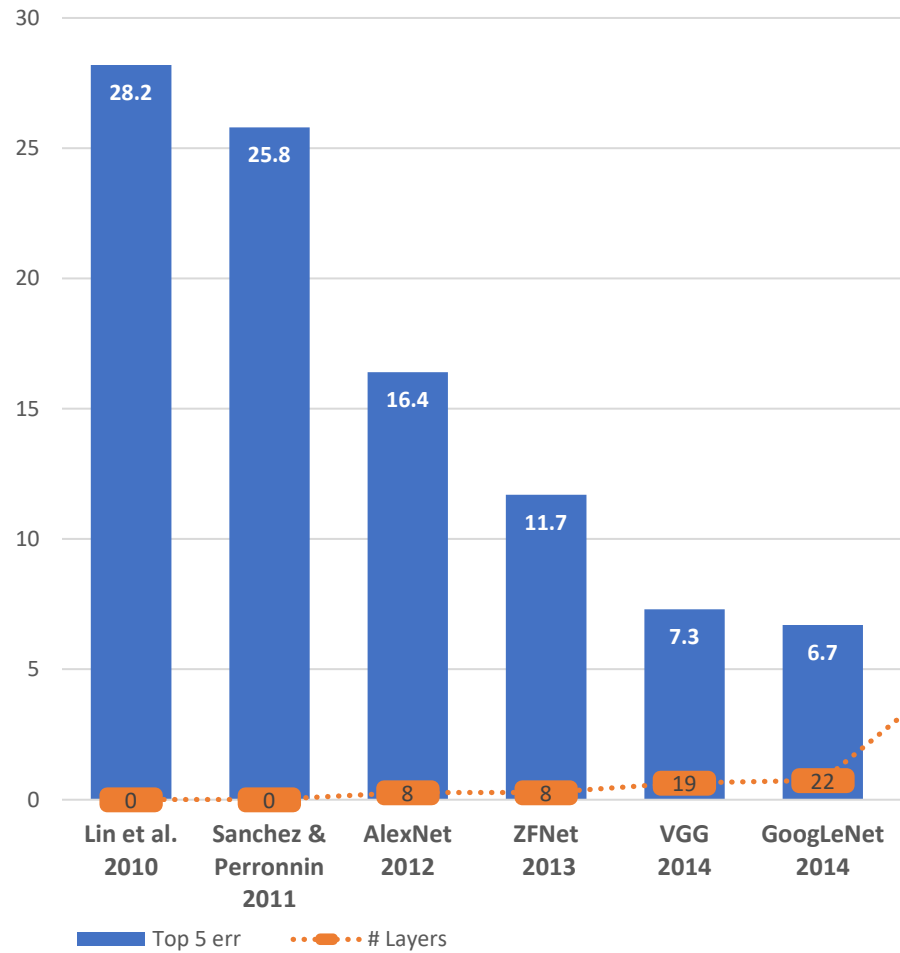
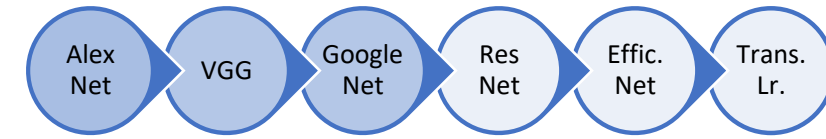
	5X5 conv, C	3X3 conv, C 3X3 conv, C
Receptive Field	5	5
# Parameters	$25C^2$	$9C^2 + 9C^2 = 18C^2$
# FLOPs	$25C^2HW$	$18C^2HW$
Memory Size (Output)	HWC	2HWC



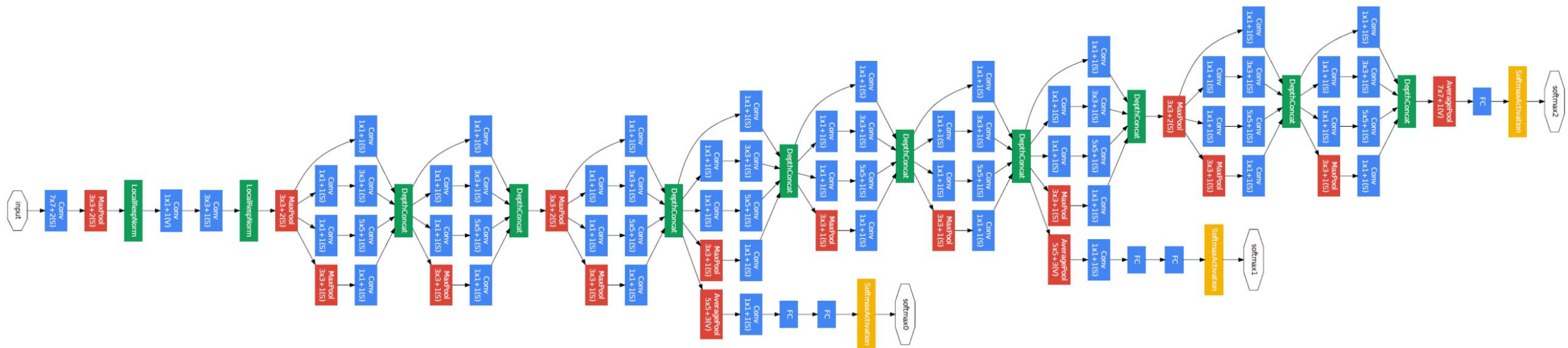
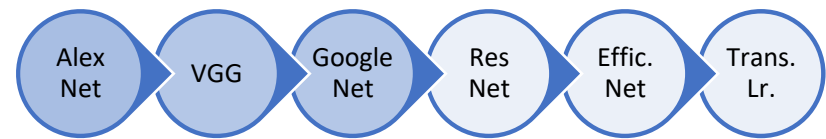
VGG resources



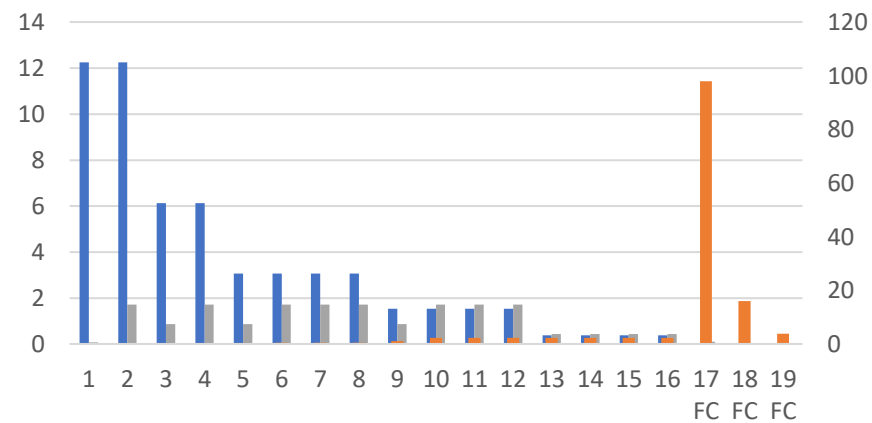
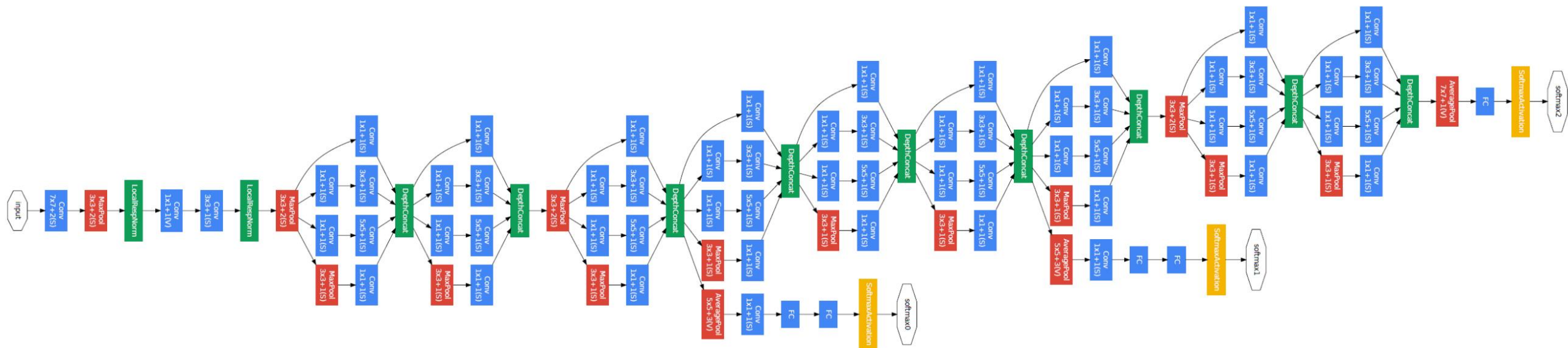
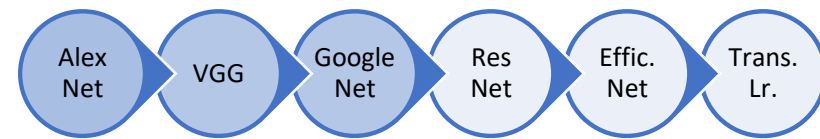
GoogLeNet



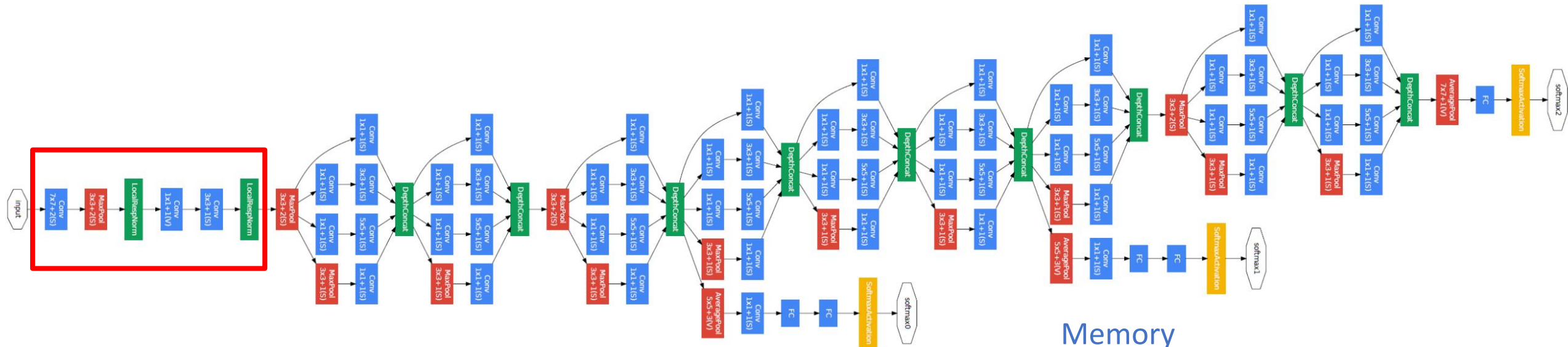
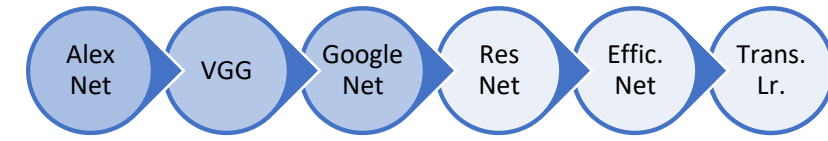
GoogLeNet



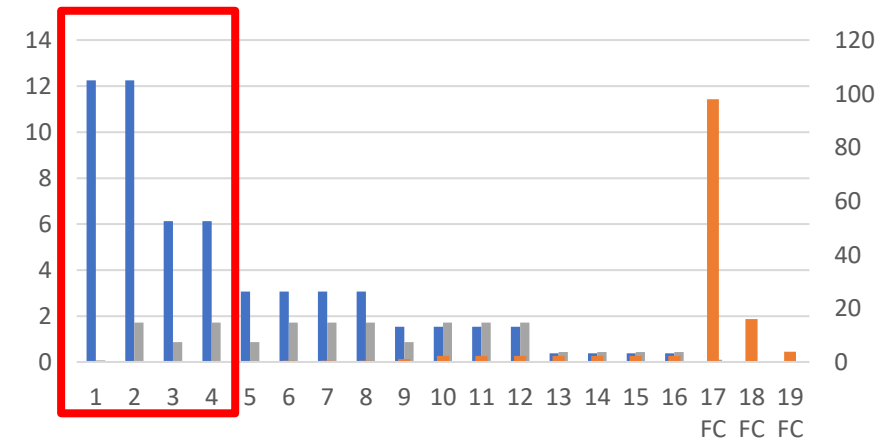
GoogLeNet



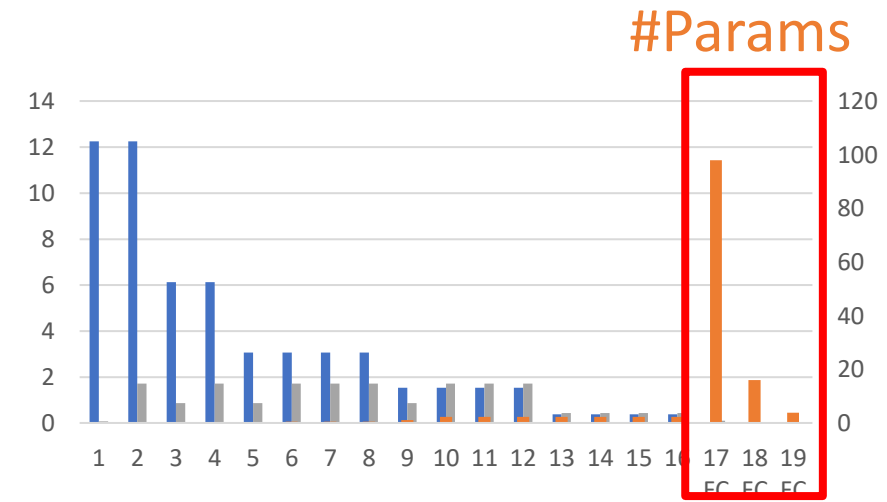
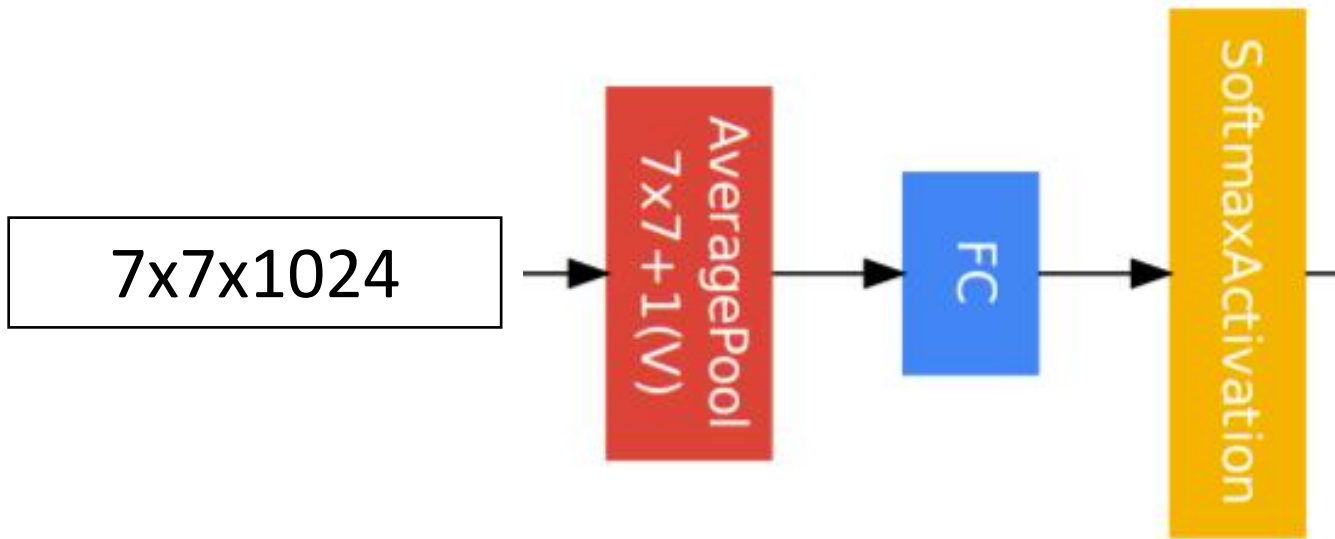
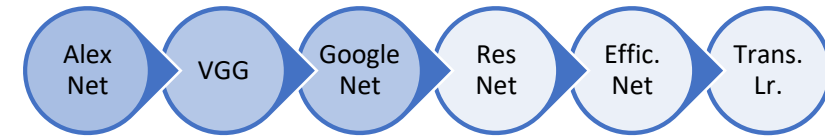
Aggressive downsampling



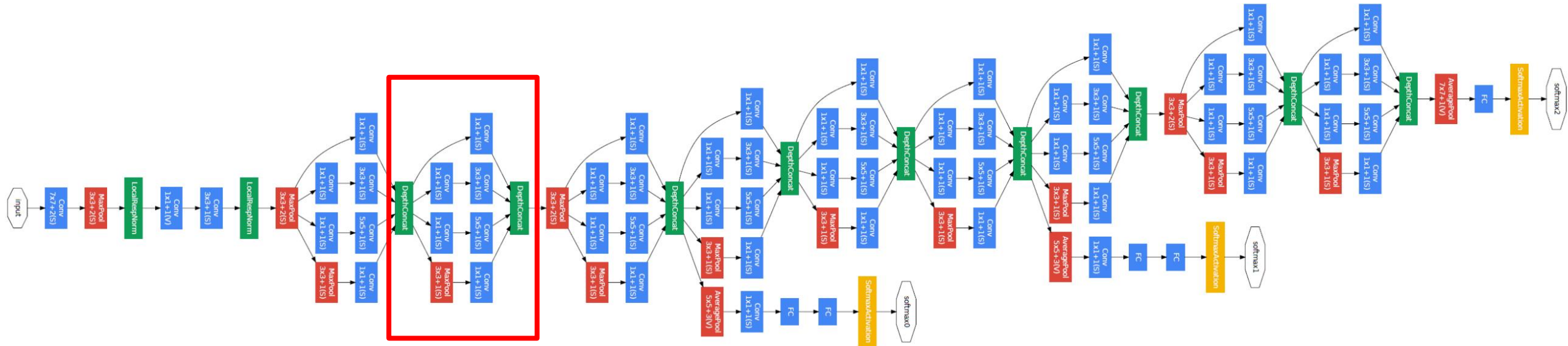
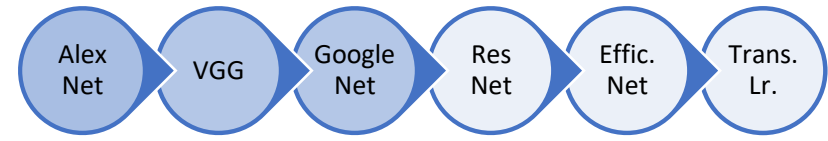
Memory



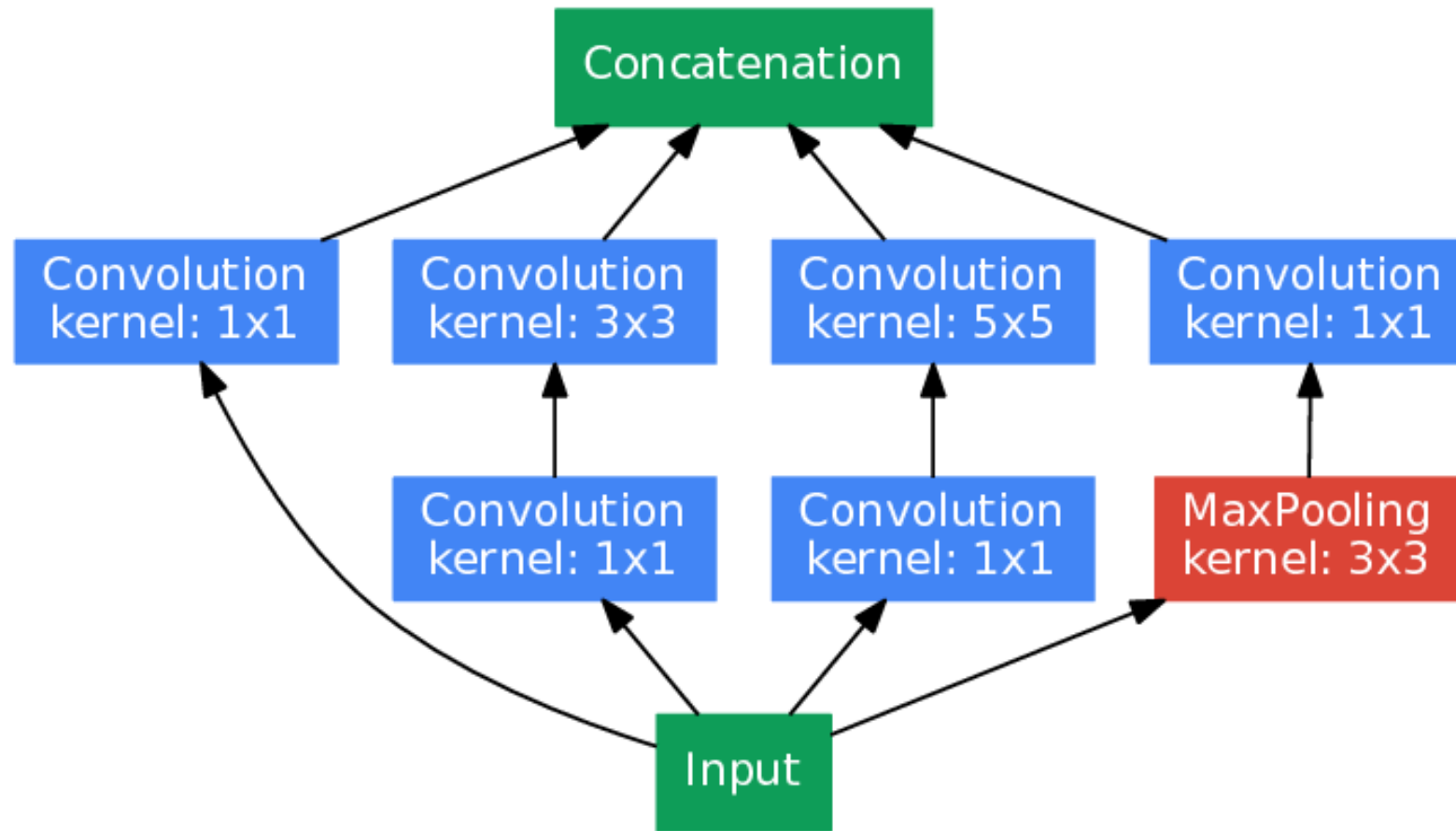
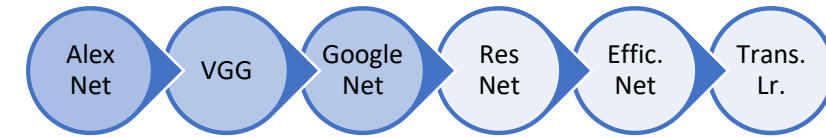
Global average pooling



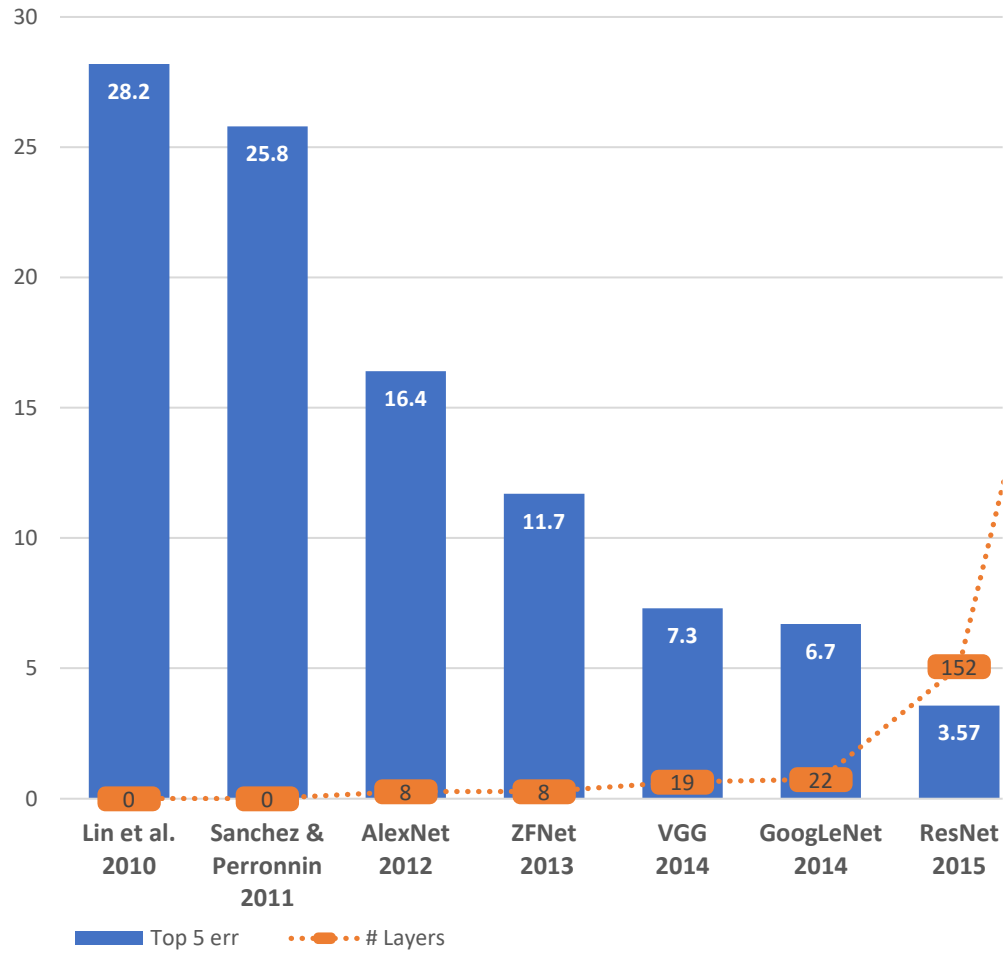
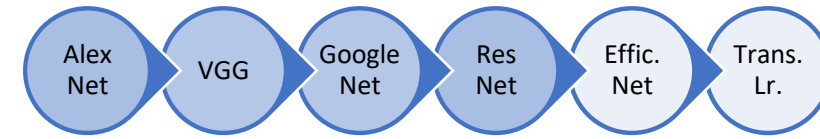
GoogLeNet



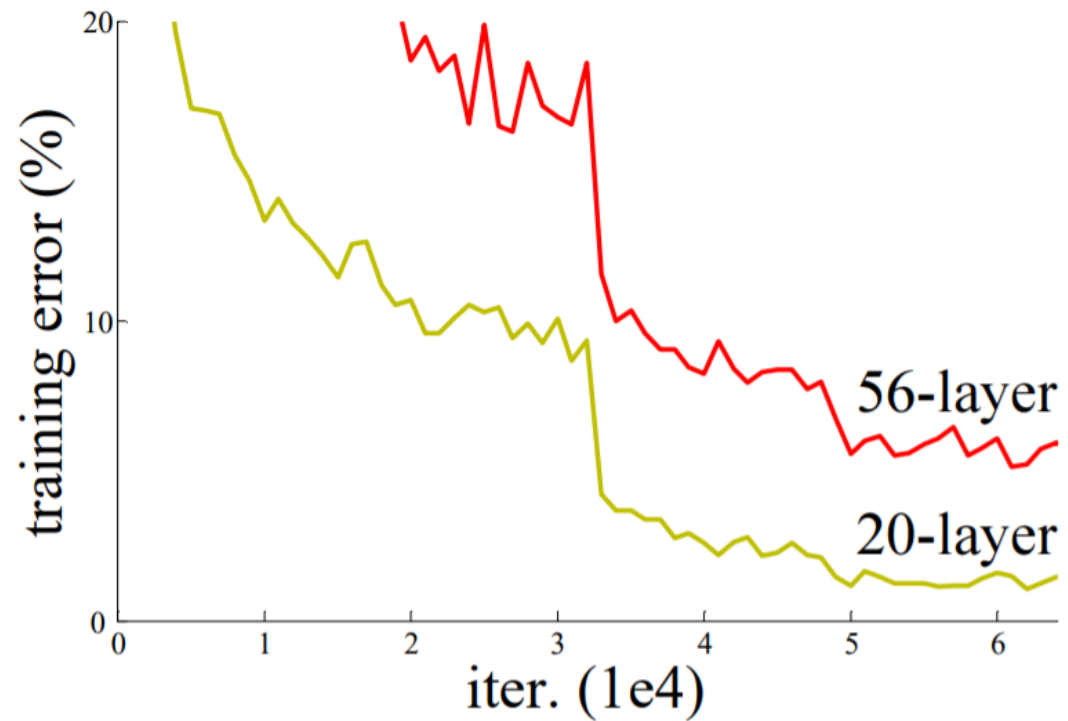
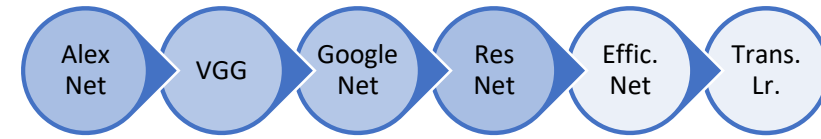
GoogLeNet: Inception blocks



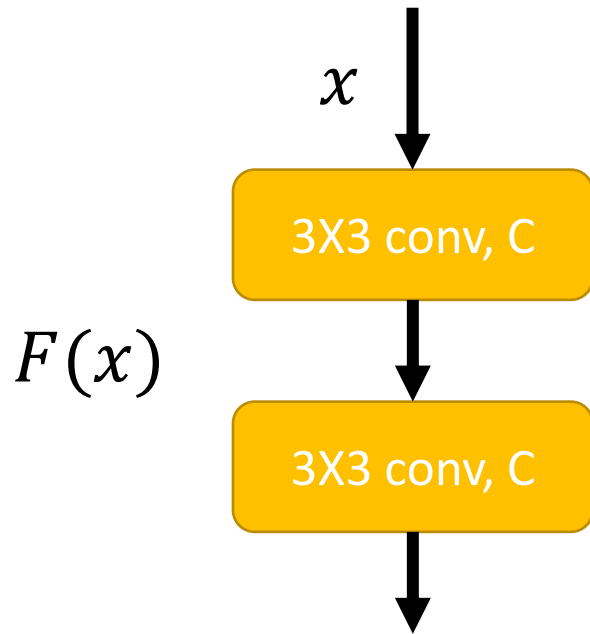
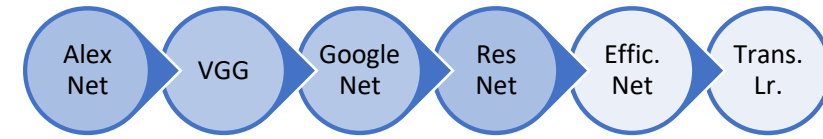
ResNet



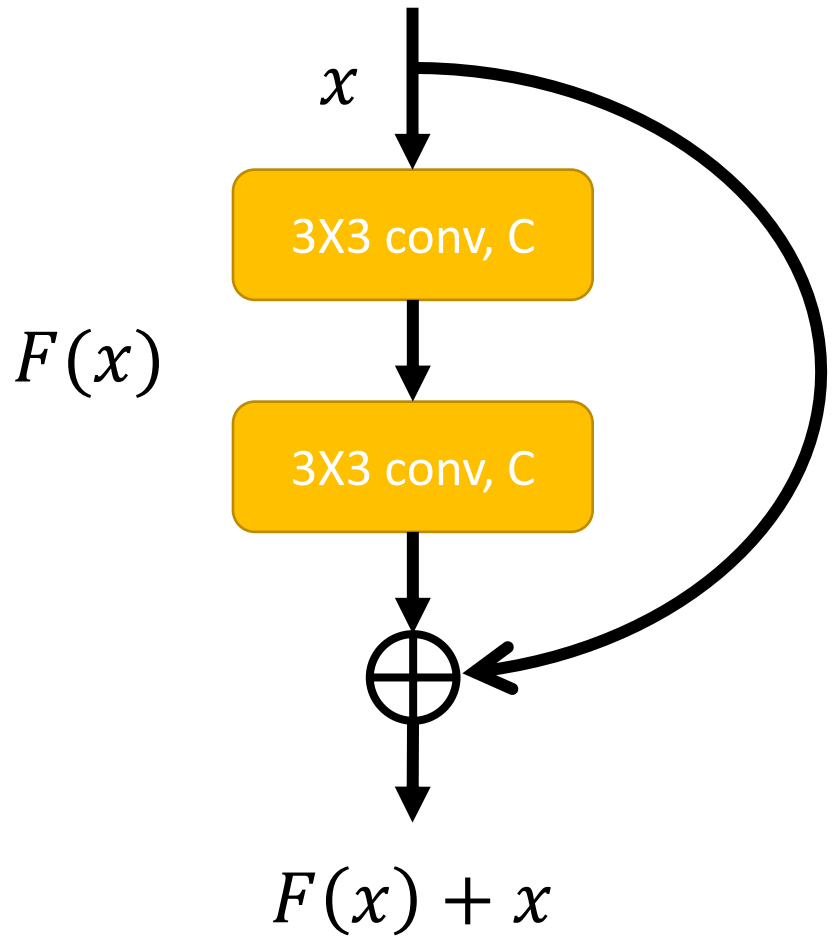
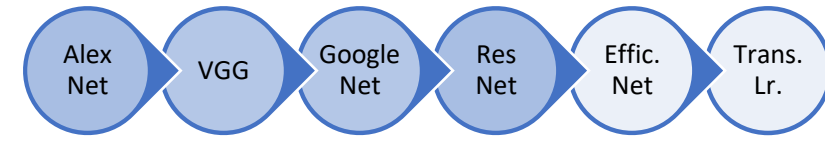
ResNet



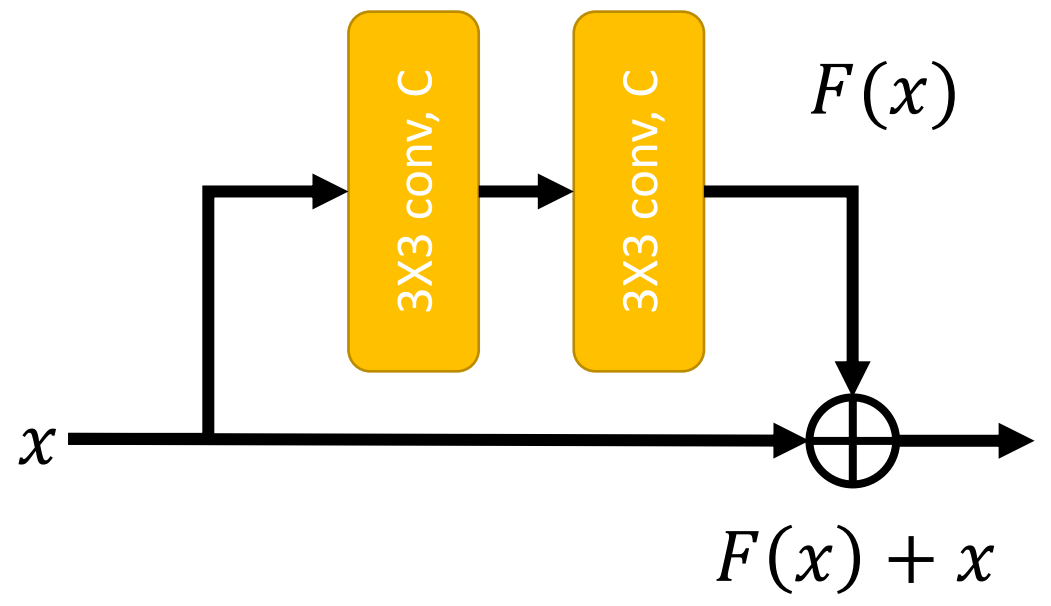
Residual Building Block



Residual Building Block

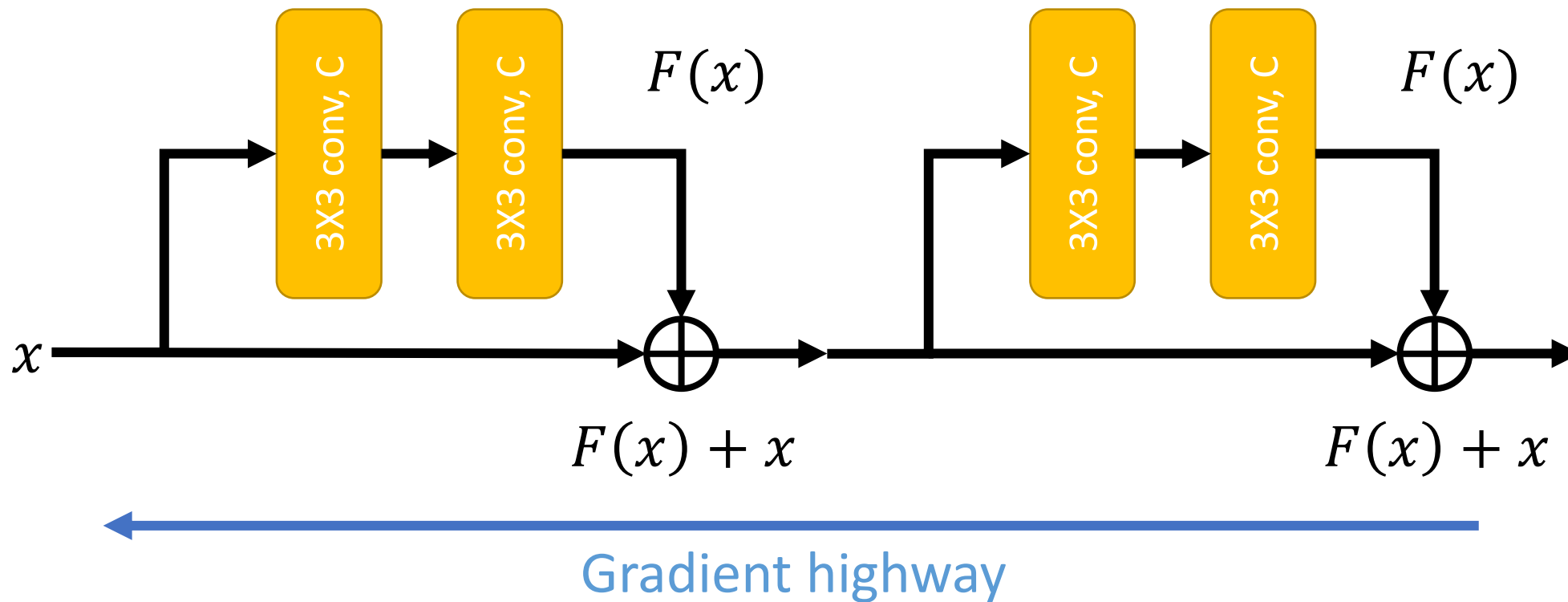
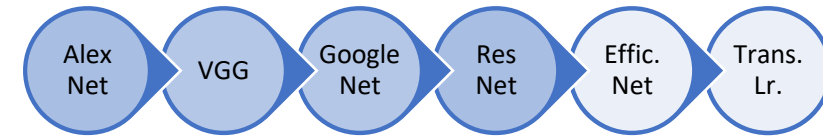


Same blocks
=

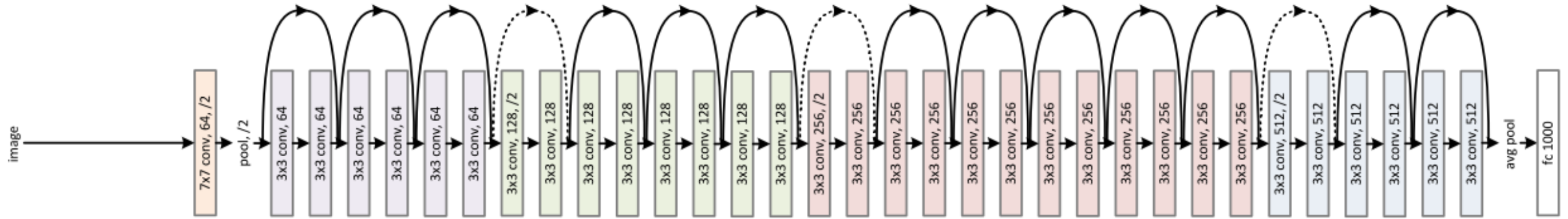
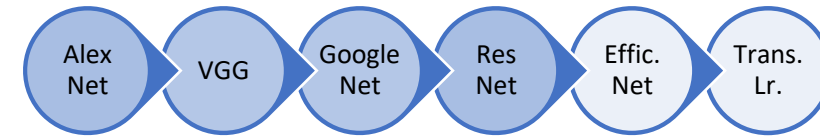


Encourage identity mapping

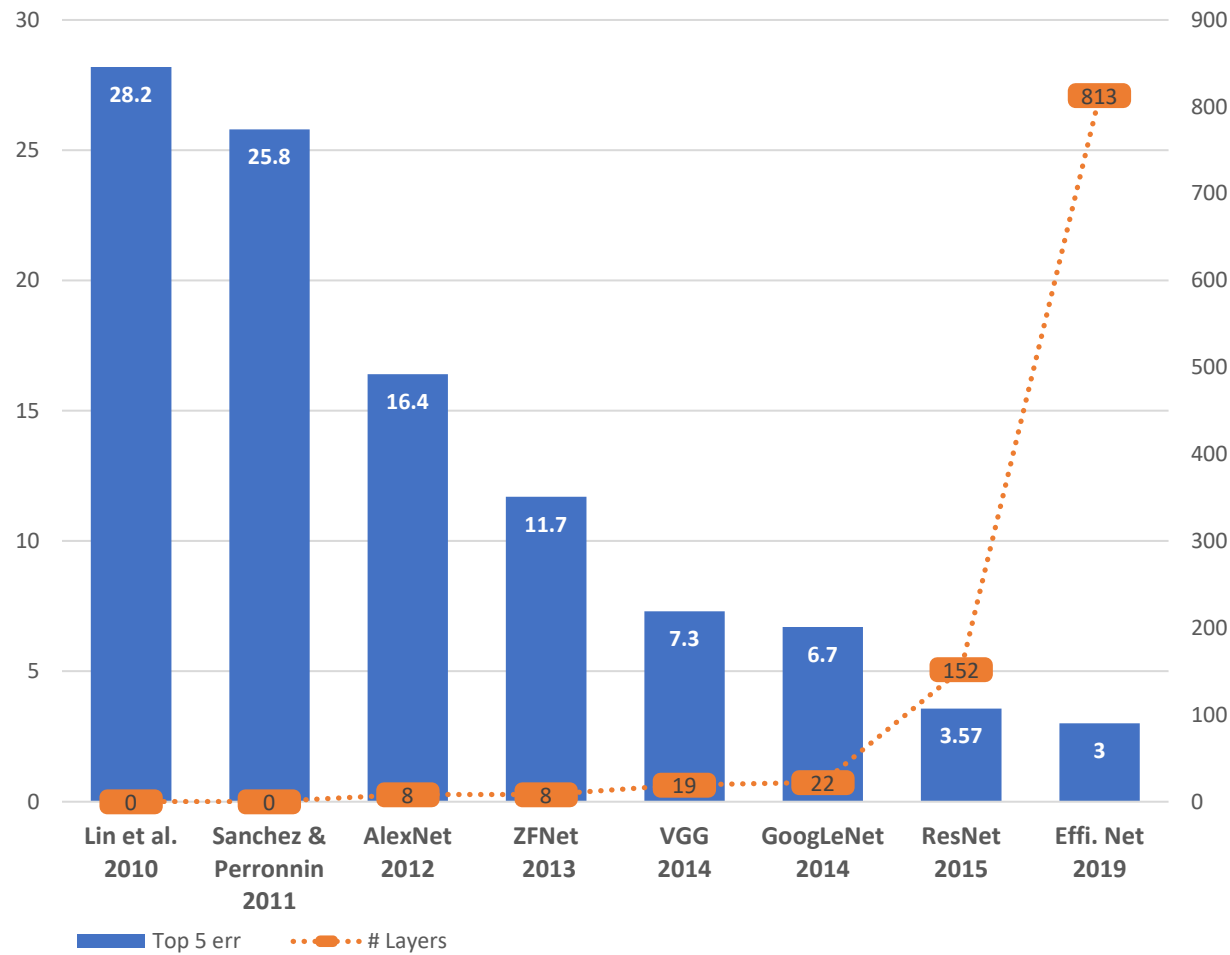
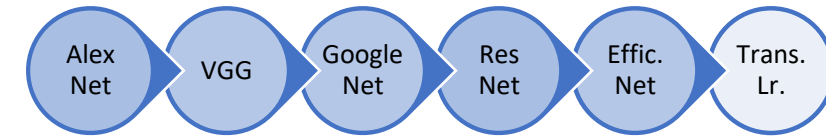
ResNet Gradient propagation



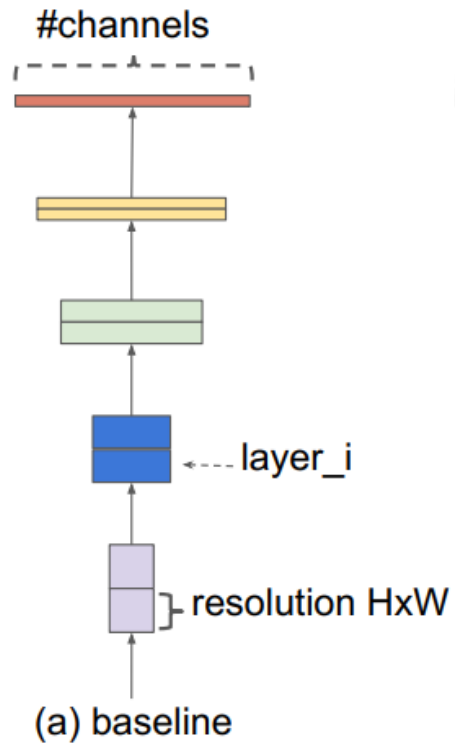
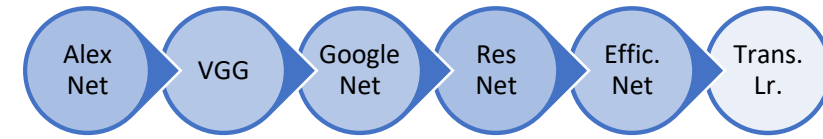
ResNet Architecture



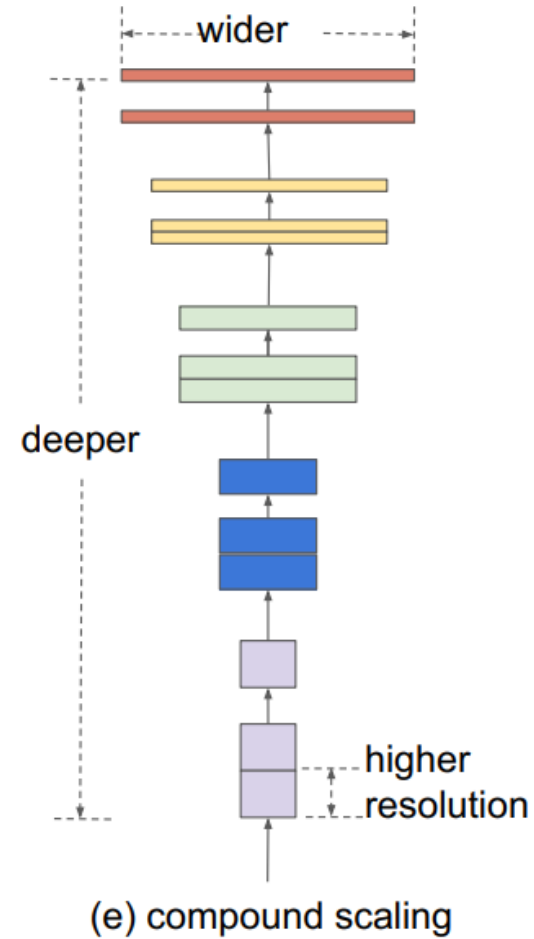
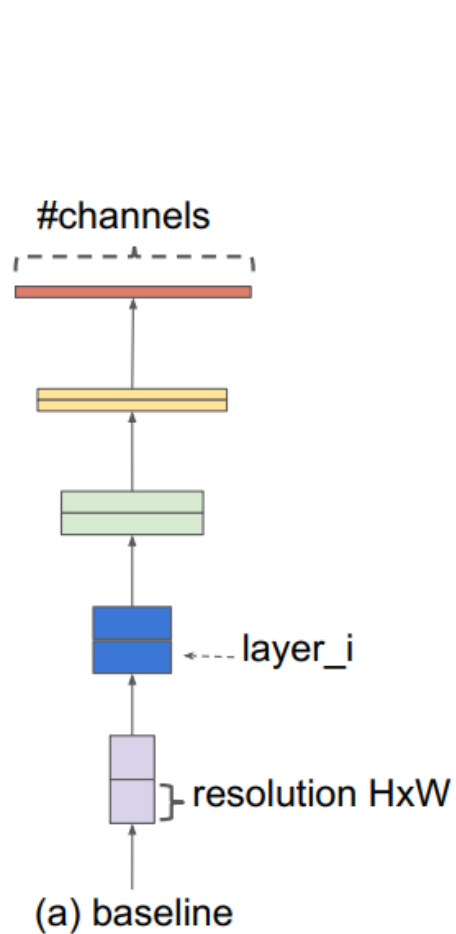
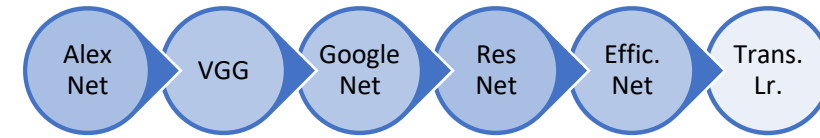
EfficientNet



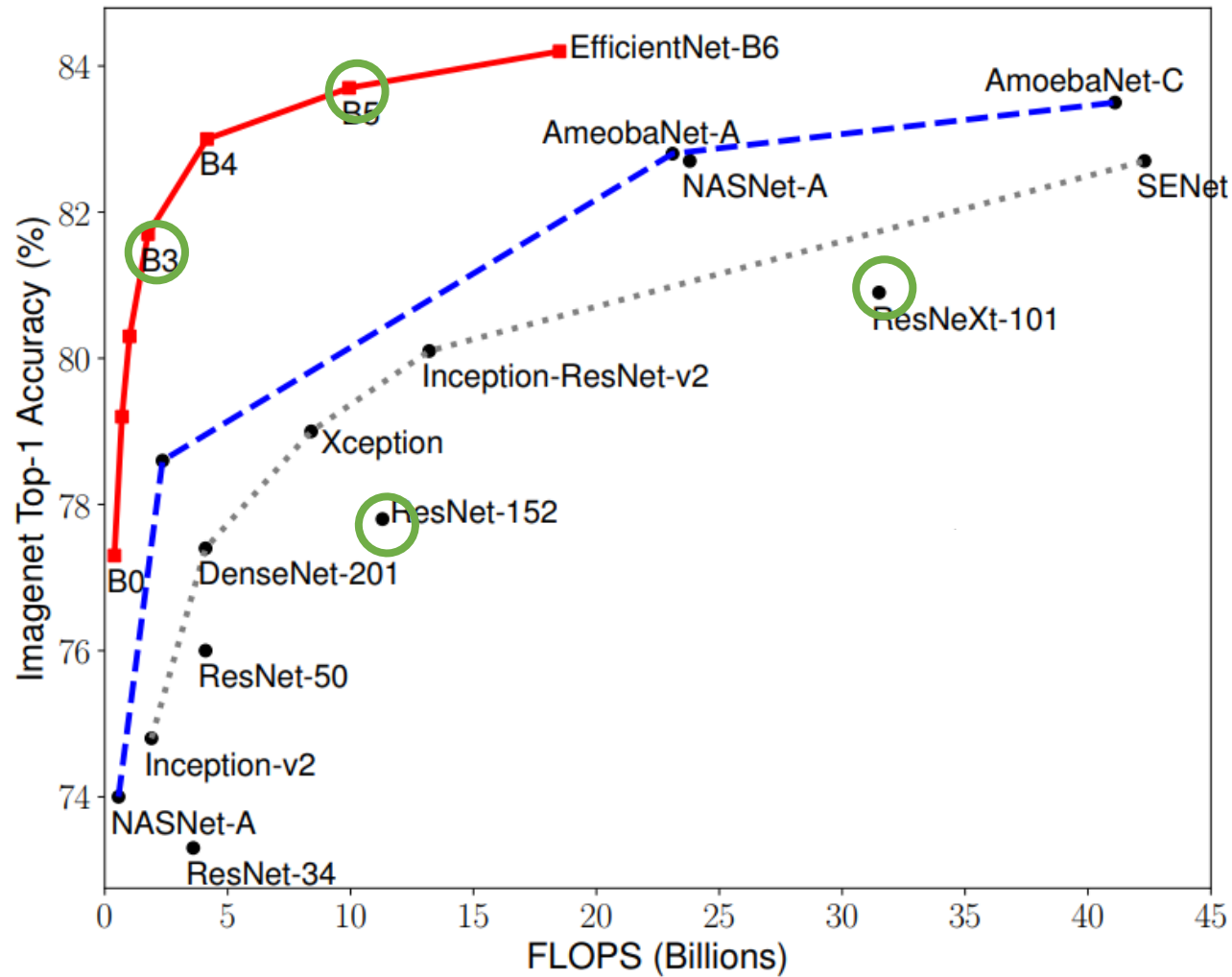
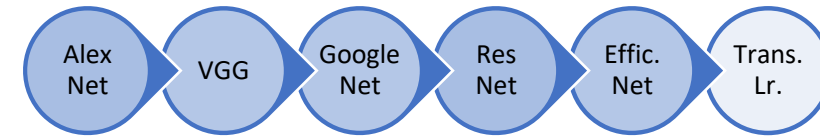
EfficientNet



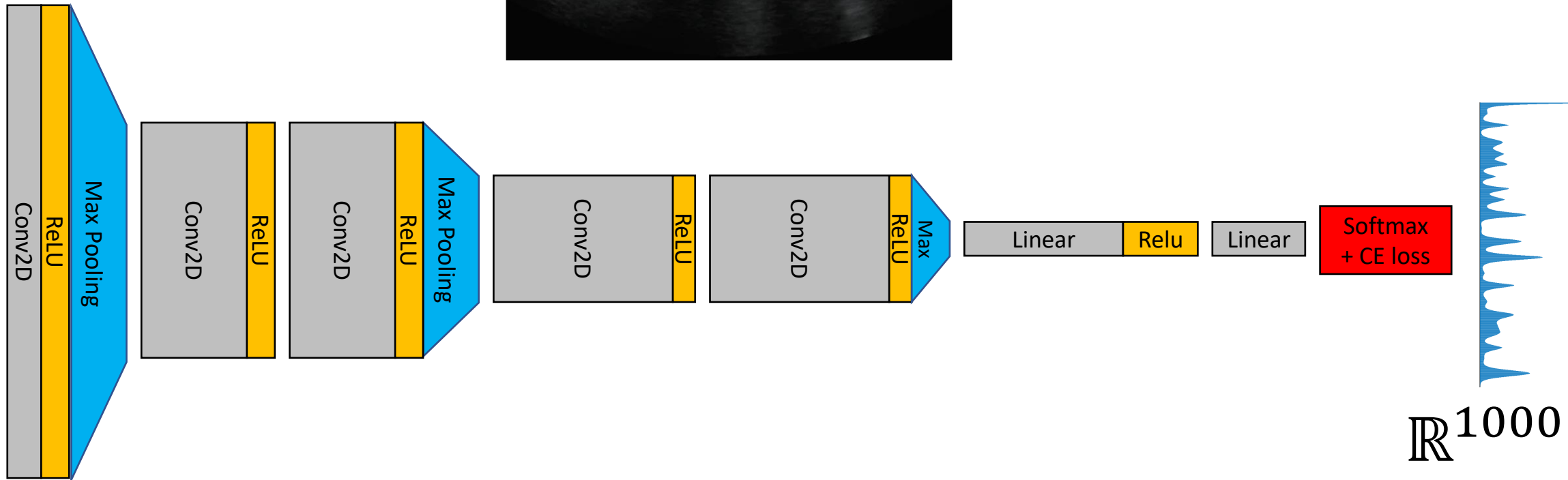
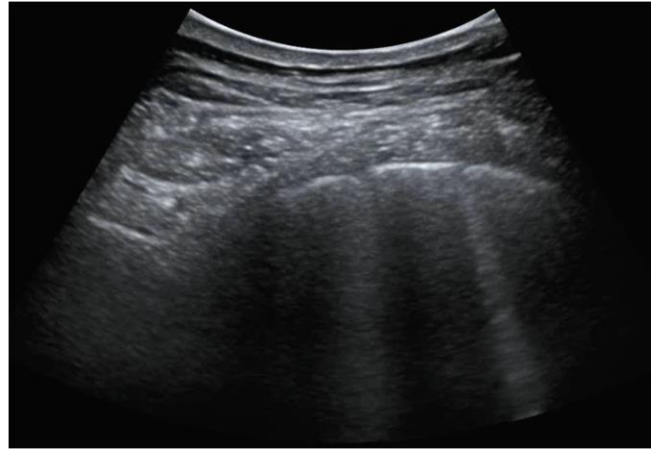
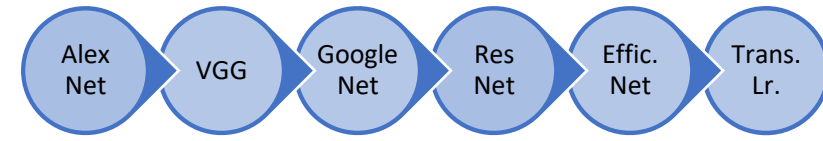
EfficientNet



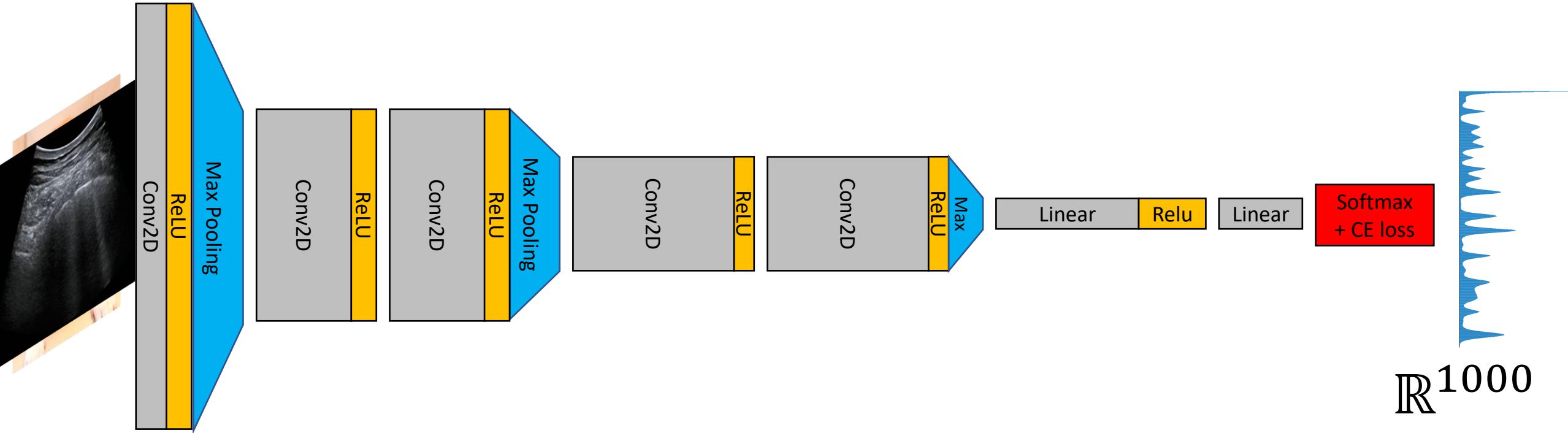
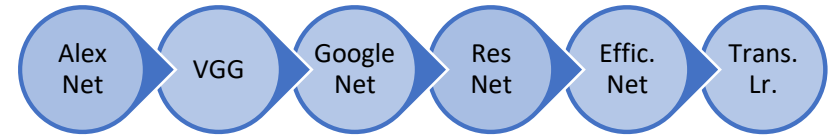
Results



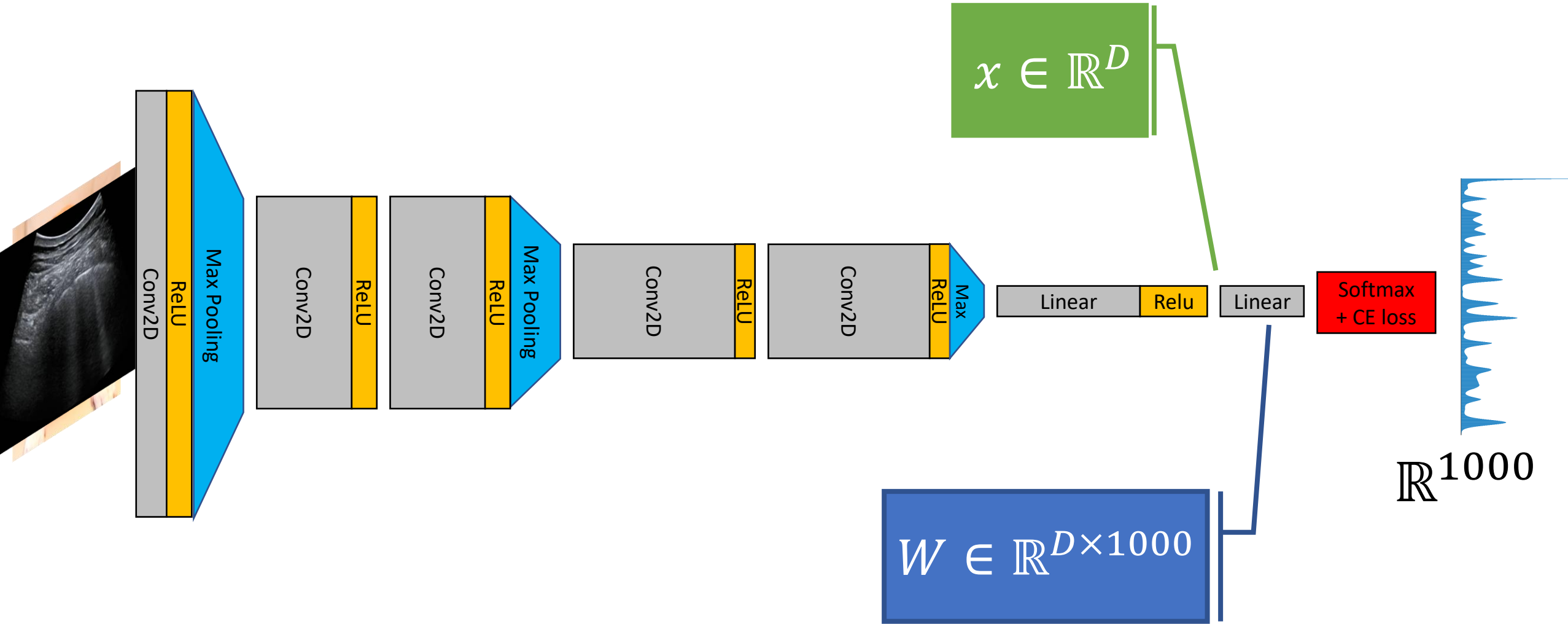
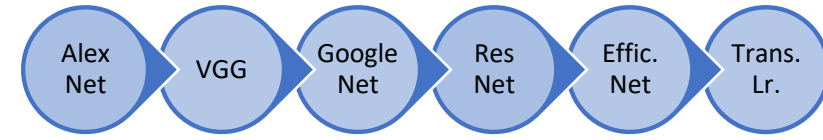
Transfer Learning



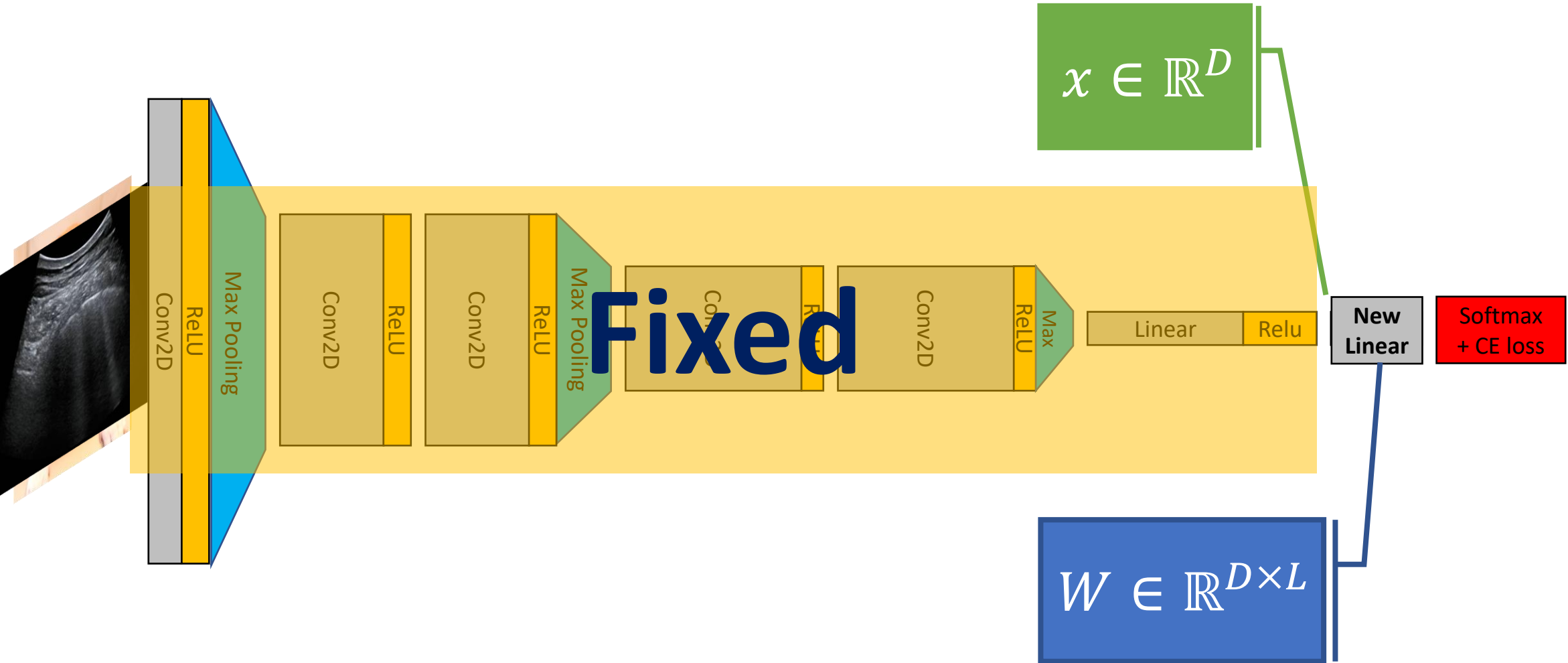
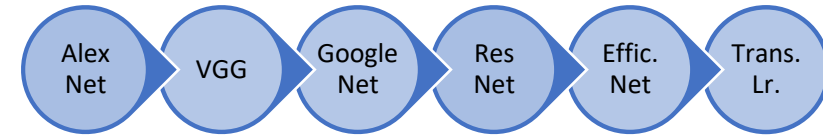
Transfer Learning



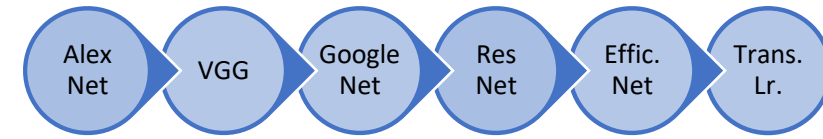
Transfer Learning



Transfer Learning



Transfer Learning



- In practice:

```
from torchvision.prototype.models import resnet50, ResNet50_Weights

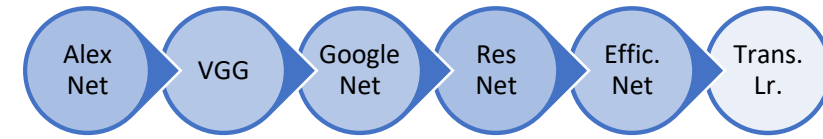
# Legacy weights with accuracy 76.130%
model = resnet50(weights=ResNet50_Weights.IMAGENET1K_V1)

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

# Best available weights (currently alias for IMAGENET1K_V2)
model = resnet50(weights=ResNet50_Weights.DEFAULT)

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

Architecture summary



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

Questions?

