

Residual Network (ResNet)

By
Prof. Khaled Mostafa El Sayed
Faculty of Computers and Information
Cairo University

Many Slides are obtained from
Kaiming He, Xiangyu Zhang, Shaoqing, Jain Sun
Fei-Fei Li & Justin Johnson & Serena Yeung

ResNet Output Performance

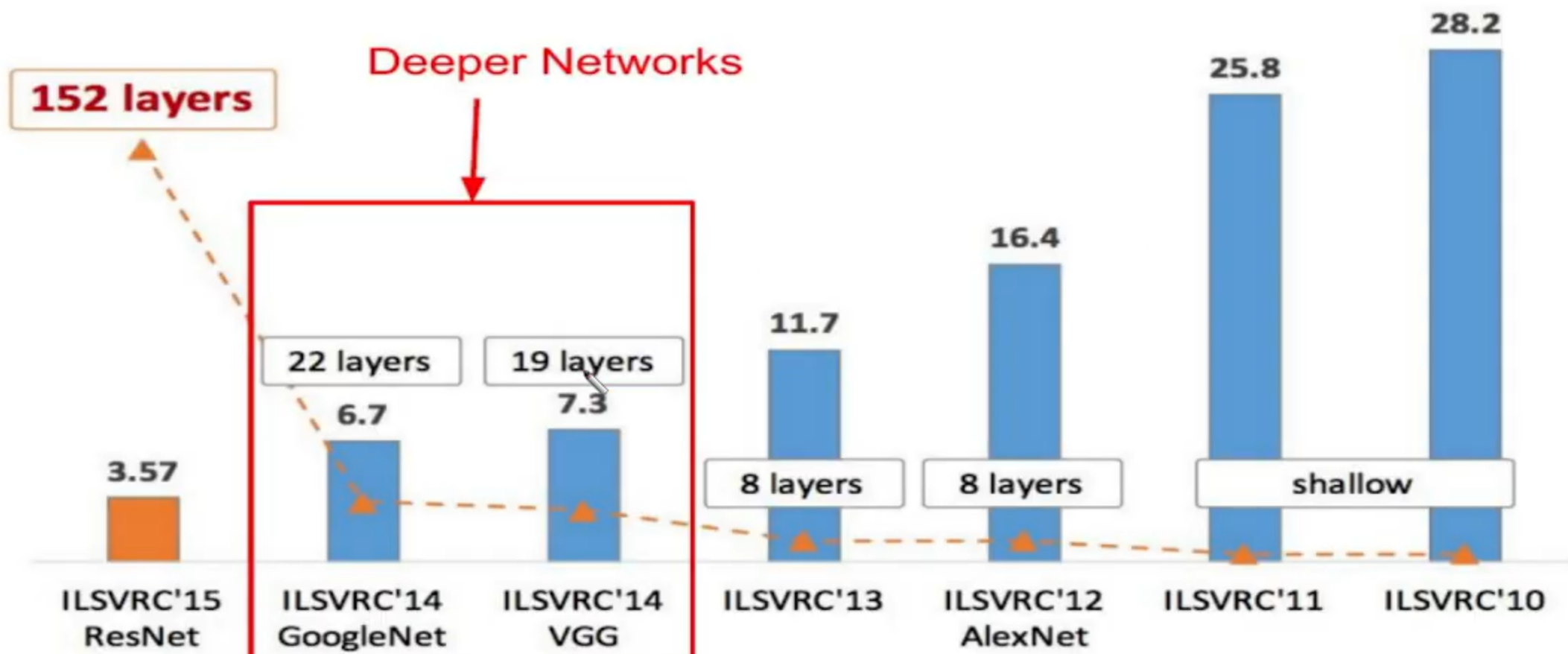
ResNet Performance

ResNet @ ILSVRC & COCO 2015 Competitions

1st places in all five main tracks

- ImageNet Classification: “Ultra-deep” 152-layer nets
- ImageNet Detection: 16% better than 2nd
- ImageNet Localization: 27% better than 2nd
- COCO Detection: 11% better than 2nd
- COCO Segmentation: 12% better than 2nd

ResNet Performance

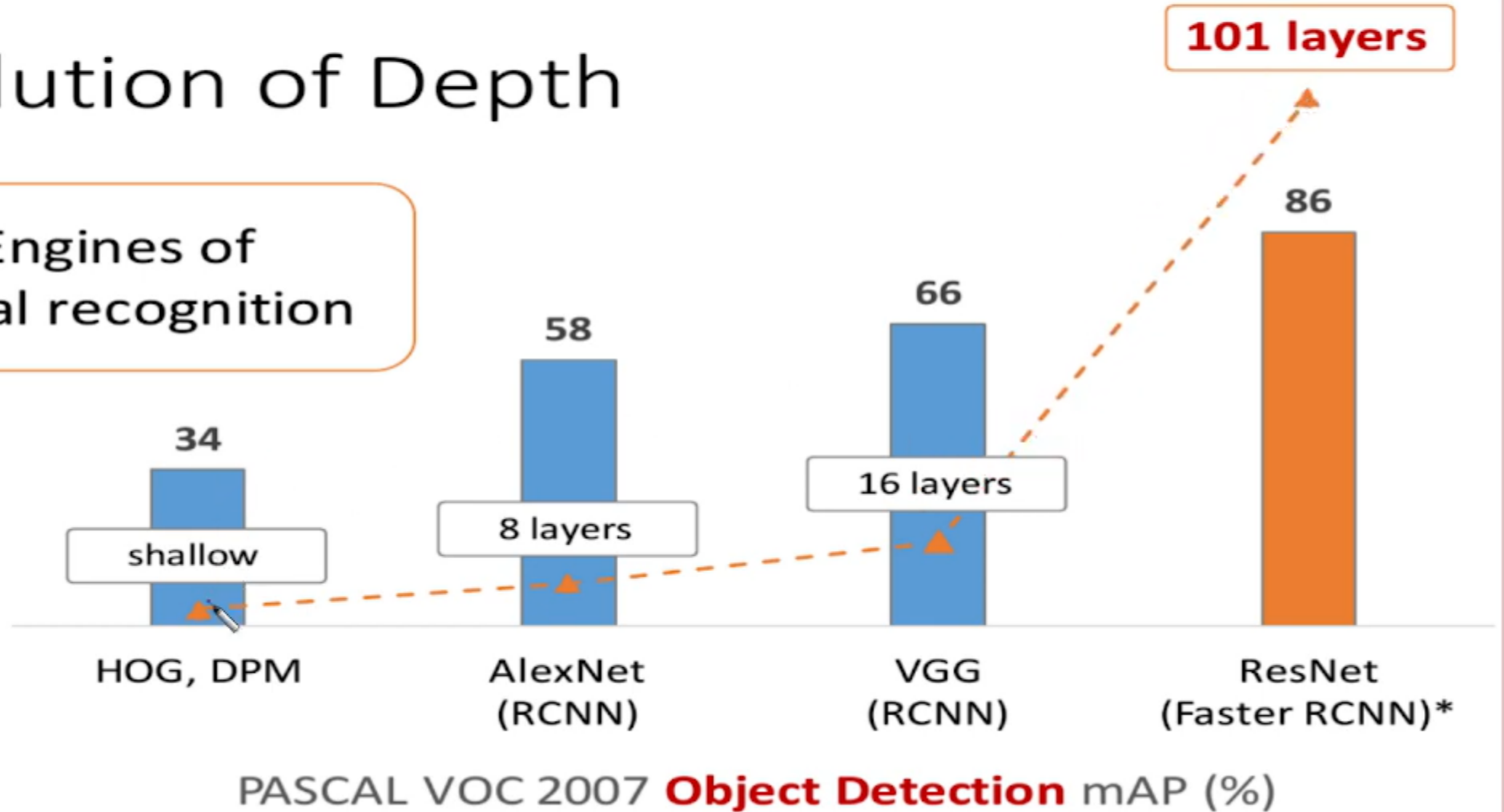


ImageNet Classification top-5 error (%)

Revolution of Depth

Revolution of Depth

Engines of
visual recognition



Revolution of Depth

Small filters, Deeper networks

8 layers (AlexNet)

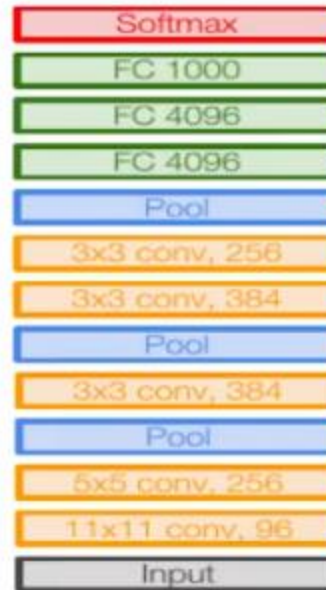
-> 16 - 19 layers (VGG16Net)

Only 3x3 CONV stride 1, pad 1
and 2x2 MAX POOL stride 2

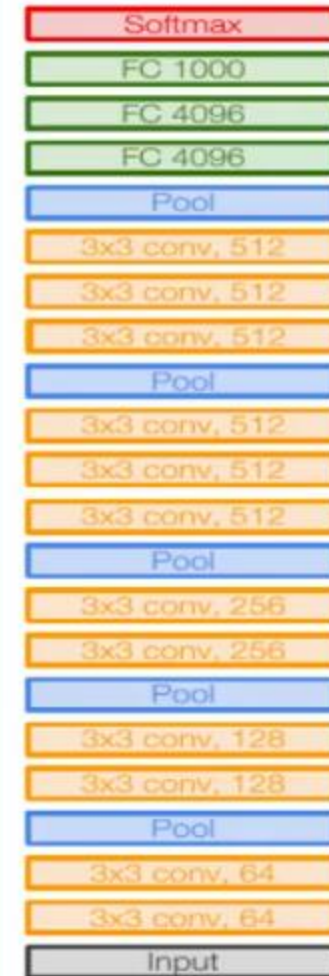
11.7% top 5 error in ILSVRC'13
(ZFNet)

-> 7.3% top 5 error in ILSVRC'14

Stack of three 3x3 conv (stride 1) layers
has same effective receptive field as
one 7x7 conv layer



AlexNet



VGG16



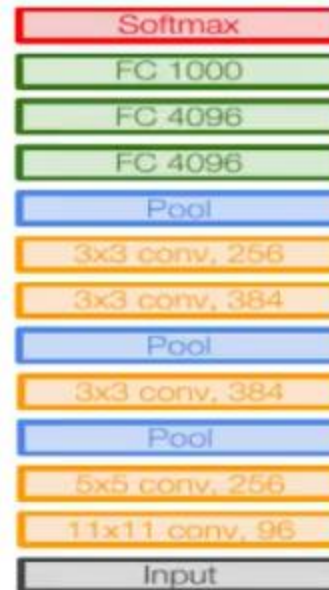
VGG19

Revolution of Depth

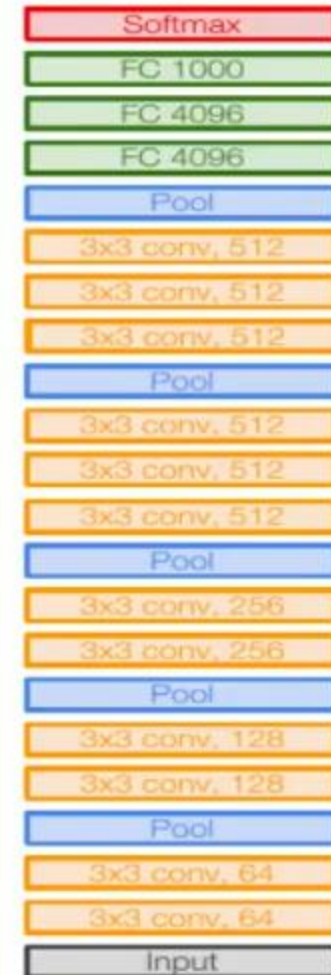
Stack of three 3x3 conv (stride 1) layers has same **effective receptive field** as one 7x7 conv layer

But deeper, more non-linearities

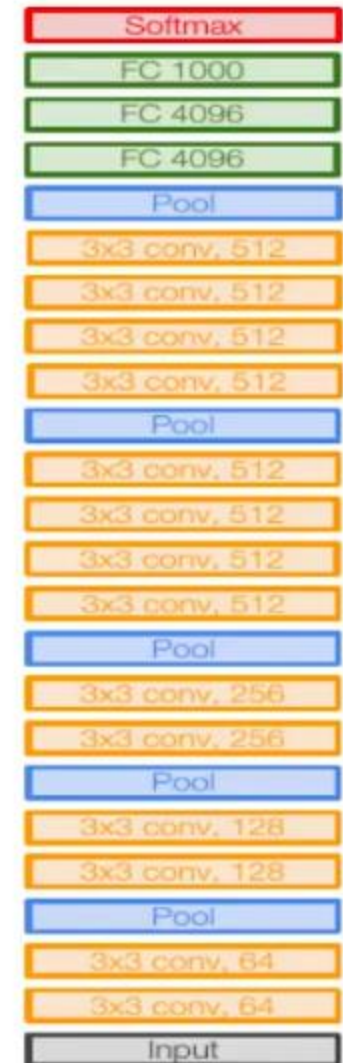
And fewer parameters: $3 * (3^2 C^2)$ vs. $7^2 C^2$ for C channels per layer



AlexNet



VGG16



VGG19

INPUT: [224x224x3] memory: $224 \times 224 \times 3 = 150\text{K}$ params: 0 (not counting biases)

CONV3-64: [224x224x64] memory: $224 \times 224 \times 64 = 3.2\text{M}$ params: $(3 \times 3 \times 3) \times 64 = 1,728$

CONV3-64: [224x224x64] memory: $224 \times 224 \times 64 = 3.2\text{M}$ params: $(3 \times 3 \times 64) \times 64 = 36,864$

POOL2: [112x112x64] memory: $112 \times 112 \times 64 = 800\text{K}$ params: 0

CONV3-128: [112x112x128] memory: $112 \times 112 \times 128 = 1.6\text{M}$ params: $(3 \times 3 \times 64) \times 128 = 73,728$

CONV3-128: [112x112x128] memory: $112 \times 112 \times 128 = 1.6\text{M}$ params: $(3 \times 3 \times 128) \times 128 = 147,456$

POOL2: [56x56x128] memory: $56 \times 56 \times 128 = 400\text{K}$ params: 0

CONV3-256: [56x56x256] memory: $56 \times 56 \times 256 = 800\text{K}$ params: $(3 \times 3 \times 128) \times 256 = 294,912$

CONV3-256: [56x56x256] memory: $56 \times 56 \times 256 = 800\text{K}$ params: $(3 \times 3 \times 256) \times 256 = 589,824$

CONV3-256: [56x56x256] memory: $56 \times 56 \times 256 = 800\text{K}$ params: $(3 \times 3 \times 256) \times 256 = 589,824$

POOL2: [28x28x256] memory: $28 \times 28 \times 256 = 200\text{K}$ params: 0

CONV3-512: [28x28x512] memory: $28 \times 28 \times 512 = 400\text{K}$ params: $(3 \times 3 \times 256) \times 512 = 1,179,648$

CONV3-512: [28x28x512] memory: $28 \times 28 \times 512 = 400\text{K}$ params: $(3 \times 3 \times 512) \times 512 = 2,359,296$

CONV3-512: [28x28x512] memory: $28 \times 28 \times 512 = 400\text{K}$ params: $(3 \times 3 \times 512) \times 512 = 2,359,296$

POOL2: [14x14x512] memory: $14 \times 14 \times 512 = 100\text{K}$ params: 0

CONV3-512: [14x14x512] memory: $14 \times 14 \times 512 = 100\text{K}$ params: $(3 \times 3 \times 512) \times 512 = 2,359,296$

CONV3-512: [14x14x512] memory: $14 \times 14 \times 512 = 100\text{K}$ params: $(3 \times 3 \times 512) \times 512 = 2,359,296$

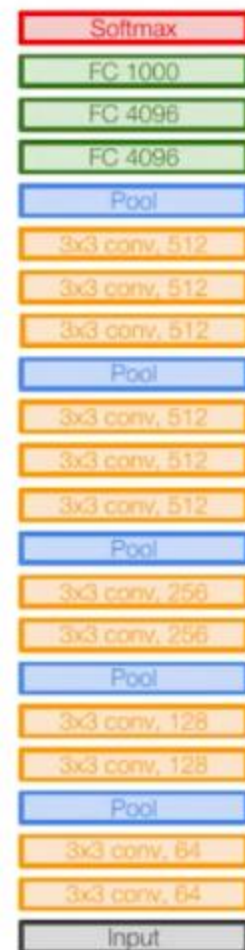
CONV3-512: [14x14x512] memory: $14 \times 14 \times 512 = 100\text{K}$ params: $(3 \times 3 \times 512) \times 512 = 2,359,296$

POOL2: [7x7x512] memory: $7 \times 7 \times 512 = 25\text{K}$ params: 0

FC: [1x1x4096] memory: 4096 params: $7 \times 7 \times 512 \times 4096 = 102,760,448$

FC: [1x1x4096] memory: 4096 params: $4096 \times 4096 = 16,777,216$

FC: [1x1x1000] memory: 1000 params: $4096 \times 1000 = 4,096,000$



VGG16

TOTAL memory: $24M * 4 \text{ bytes} \approx 96MB$ / image (only forward! $\sim *2$ for bwd)

TOTAL params: 138M parameters

INPUT: [224x224x3] memory: $224*224*3=150K$ params: 0 (not counting biases)

CONV3-64: [224x224x64] memory: $224*224*64=3.2M$ params: $(3*3*3)*64 = 1,728$

CONV3-64: [224x224x64] memory: $224*224*64=3.2M$ params: $(3*3*64)*64 = 36,864$

POOL2: [112x112x64] memory: $112*112*64=800K$ params: 0

CONV3-128: [112x112x128] memory: $112*112*128=1.6M$ params: $(3*3*64)*128 = 73,728$

CONV3-128: [112x112x128] memory: $112*112*128=1.6M$ params: $(3*3*128)*128 = 147,456$

POOL2: [56x56x128] memory: $56*56*128=400K$ params: 0

CONV3-256: [56x56x256] memory: $56*56*256=800K$ params: $(3*3*128)*256 = 294,912$

CONV3-256: [56x56x256] memory: $56*56*256=800K$ params: $(3*3*256)*256 = 589,824$

CONV3-256: [56x56x256] memory: $56*56*256=800K$ params: $(3*3*256)*256 = 589,824$

POOL2: [28x28x256] memory: $28*28*256=200K$ params: 0

CONV3-512: [28x28x512] memory: $28*28*512=400K$ params: $(3*3*256)*512 = 1,179,648$

CONV3-512: [28x28x512] memory: $28*28*512=400K$ params: $(3*3*512)*512 = 2,359,296$

CONV3-512: [28x28x512] memory: $28*28*512=400K$ params: $(3*3*512)*512 = 2,359,296$

POOL2: [14x14x512] memory: $14*14*512=100K$ params: 0

CONV3-512: [14x14x512] memory: $14*14*512=100K$ params: $(3*3*512)*512 = 2,359,296$

CONV3-512: [14x14x512] memory: $14*14*512=100K$ params: $(3*3*512)*512 = 2,359,296$

CONV3-512: [14x14x512] memory: $14*14*512=100K$ params: $(3*3*512)*512 = 2,359,296$

POOL2: [7x7x512] memory: $7*7*512=25K$ params: 0

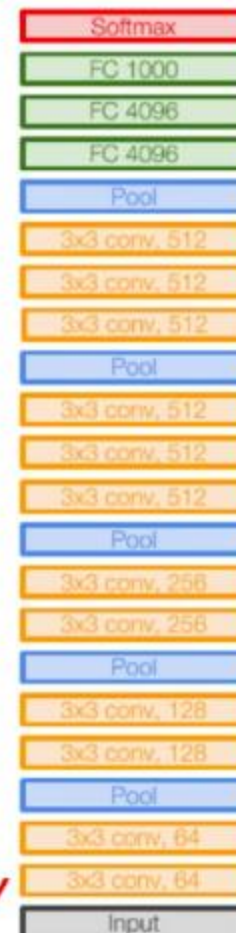
FC: [1x1x4096] memory: 4096 params: $7*7*512*4096 = 102,760,448$

FC: [1x1x4096] memory: 4096 params: $4096*4096 = 16,777,216$

FC: [1x1x1000] memory: 1000 params: $4096*1000 = 4,096,000$

TOTAL memory: $24M * 4 \text{ bytes} \approx 96MB$ / image (only forward! ~ 2 for bwd)

TOTAL params: 138M parameters



VGG16

Maximum Memory Usage

Maximum Number of Parameters

Go Deeper

AlexNet (2012)

VGG (2014)

GooLeNet (2014)

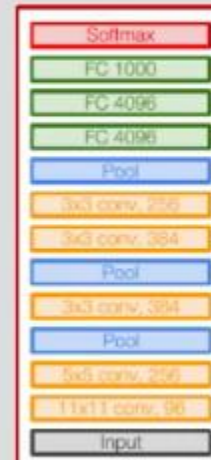
ResNet (2015)

8 Layers

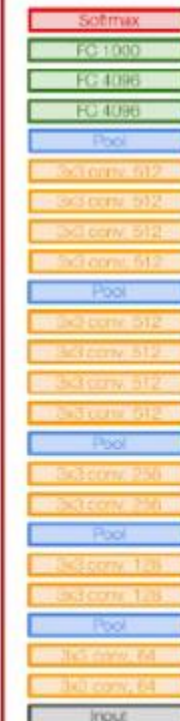
19 Layers

22 Layers

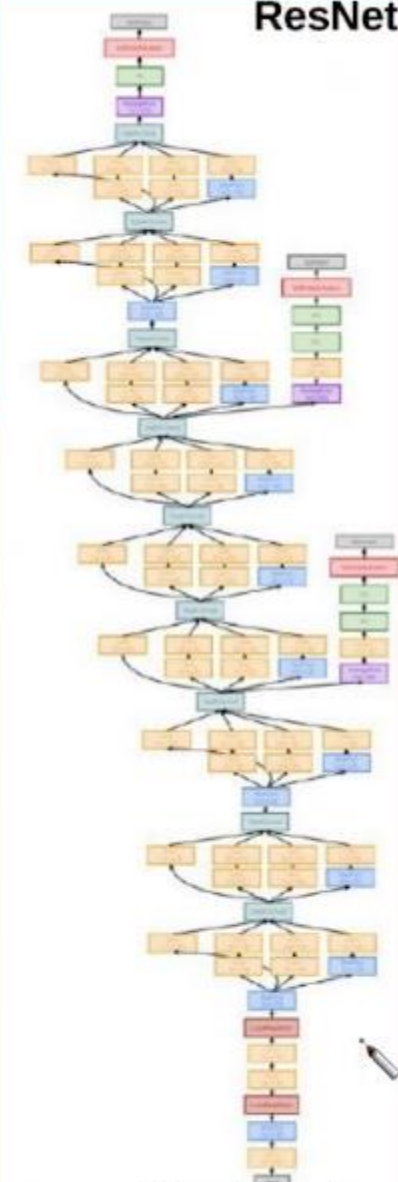
152 Layers



AlexNet



VGG19



GoogLeNet

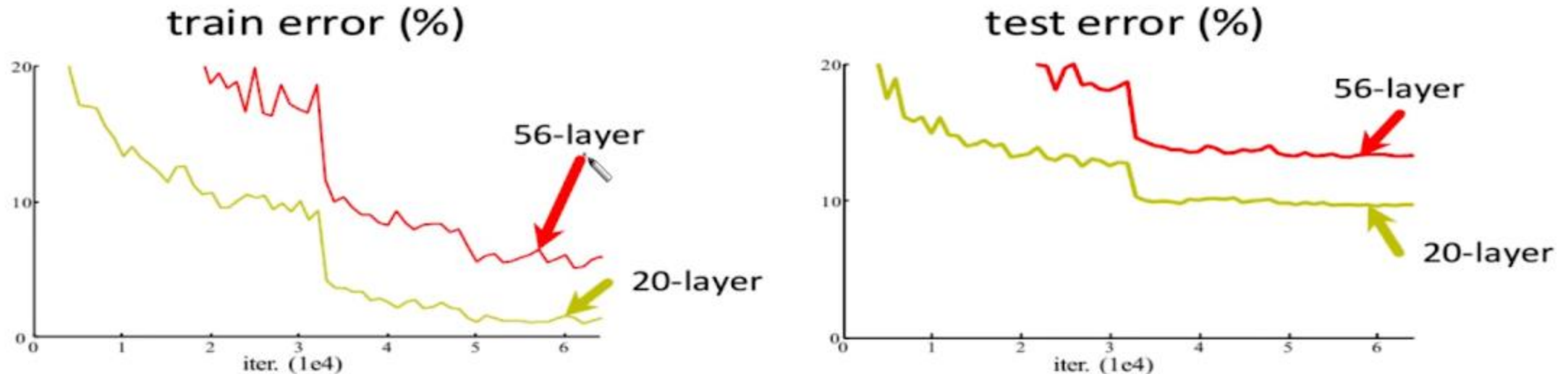


Training problem of “Stacking Layers” Deep Architectures



Error Behavior when Increasing Number of Layers (Train and Test Errors)

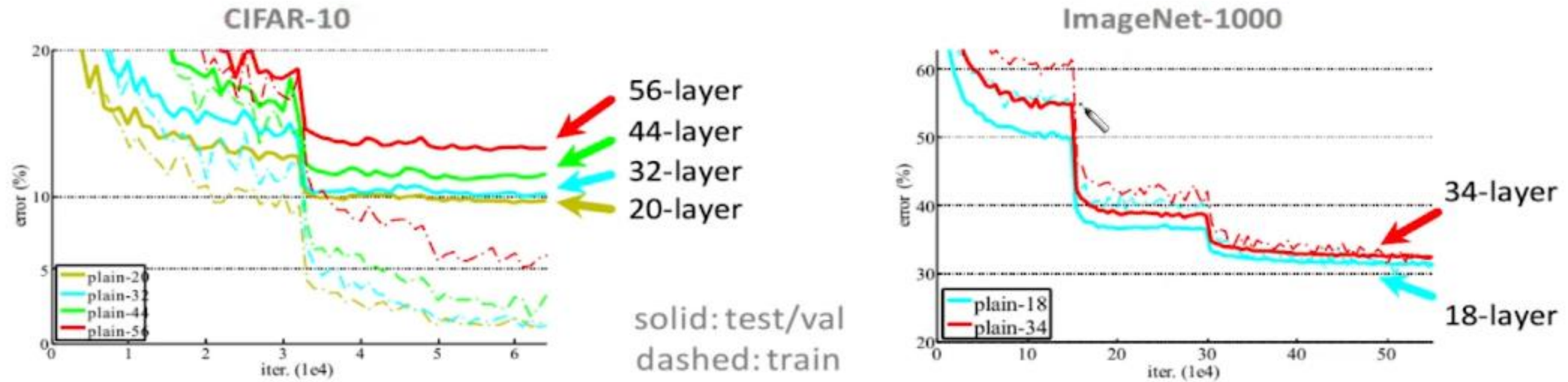
CIFAR-10



Plain nets: stacking 3x3 convlayers...

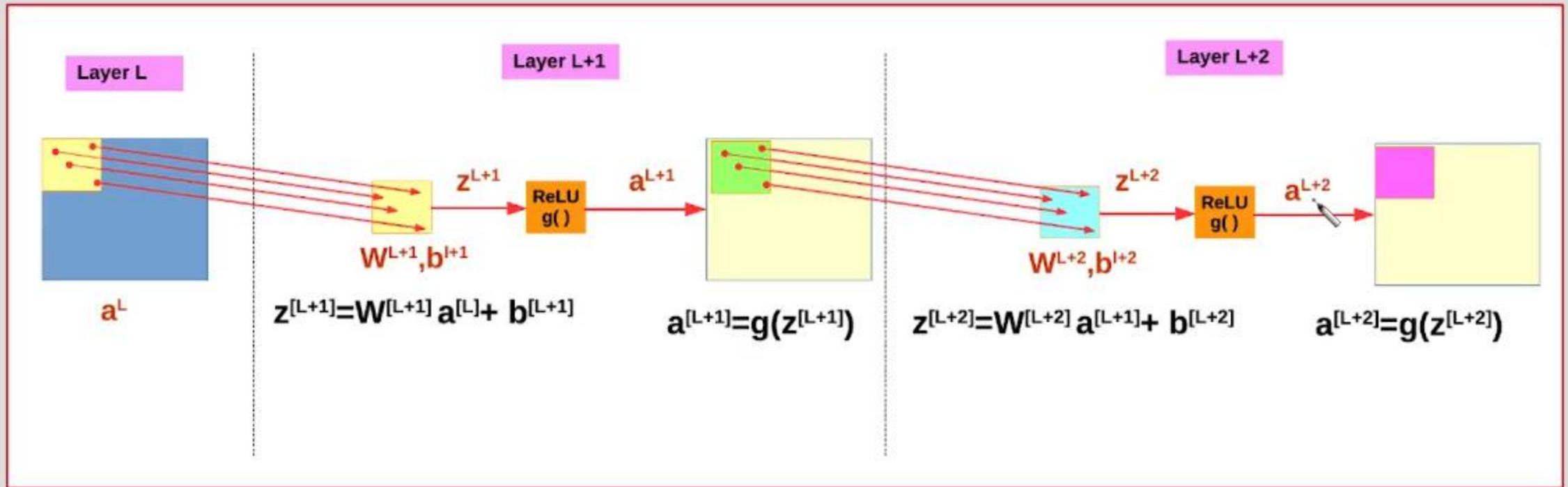
56-layer net has higher training error and test error than 20-layer net

Error Behavior when Increasing Number of Layers (Different Datasets)



“Overly deep” plain nets have higher training error
A general phenomenon, observed in many datasets

Apply Single Filter on One Block



$W^{[L+1]}$

Weights of Filter at Layer "L+1"

$b^{[L+1]}$

Bias of Filter at Layer "L+1"

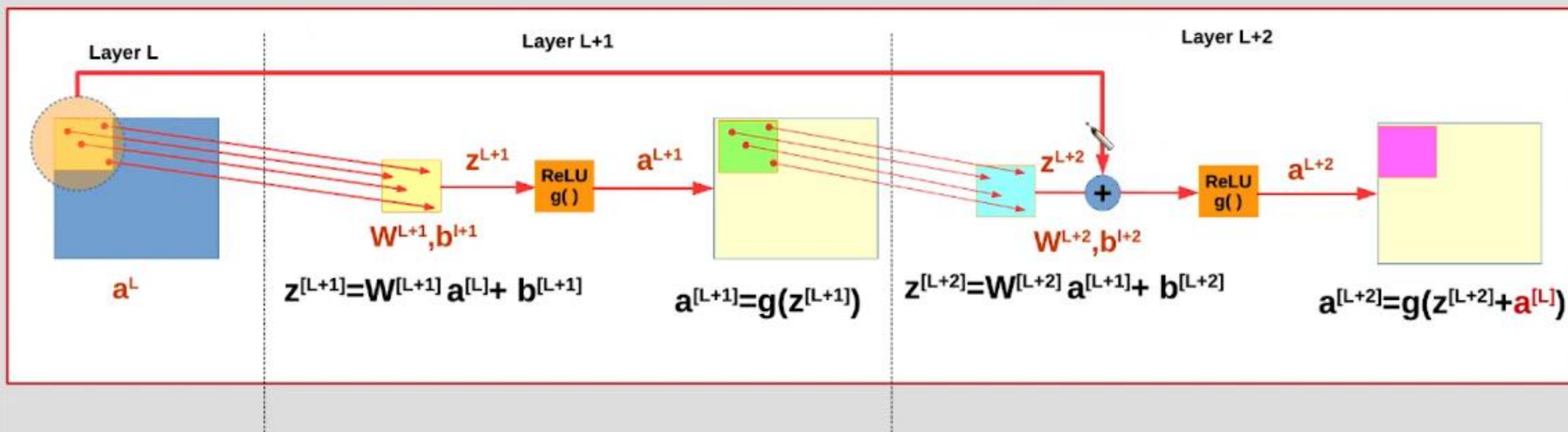
$z^{[L+1]}$

Output After Applying Filters on Layer "L"

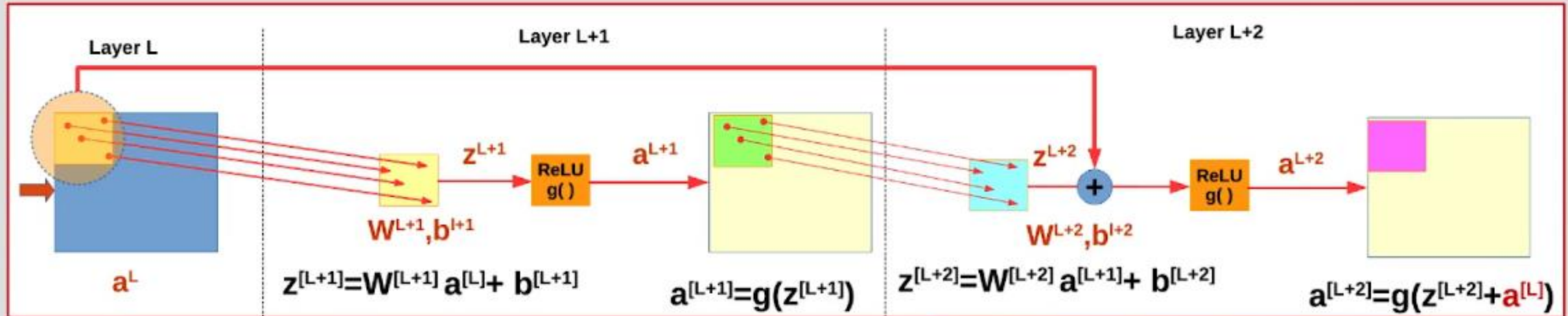
$a^{[L+1]}$

Output of ReLU $\{g(\cdot)\}$ at Layer "L+1"

Apply Single Filter on One Block and Feedback from Prev Input



Learn Identity Function Using Residual Block



Using **Relu** Activation function
 $g(a^{[L]}) = a^{[L]} \implies$ if **a** is positive value

Remember:

Output of ReLU is Always Positive

$a^{[L]}$ is Relu output of Previous Layer $\implies a^{[L]}$ Always Positive

If:

$$W^{L+2} = 0$$

$$b^{L+2} = 0$$

Then

$$z^{L+2} = 0$$

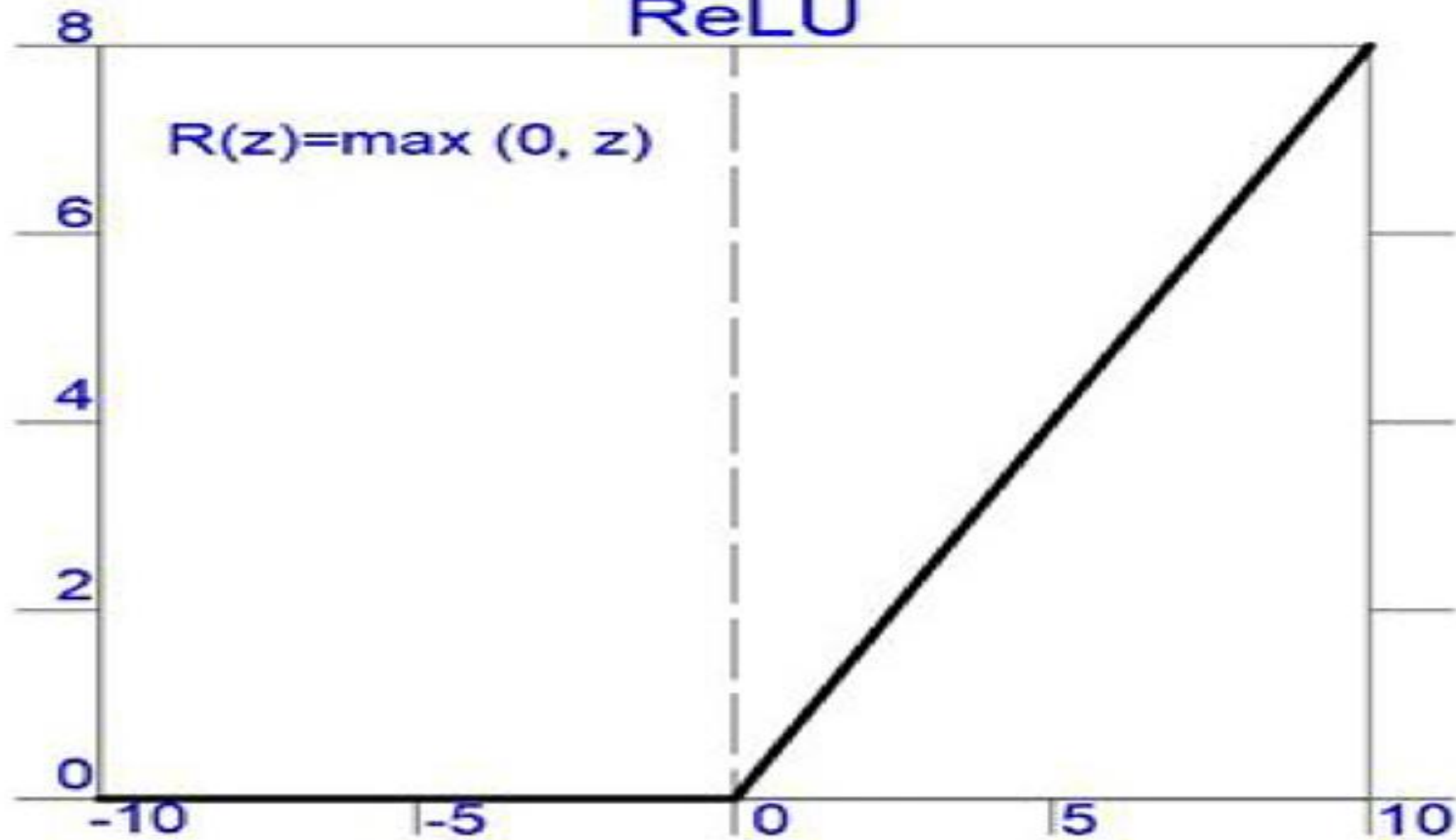
And

$$a^{L+2} = g(z^{L+2} + a^{[L]}) = a^{[L]}$$

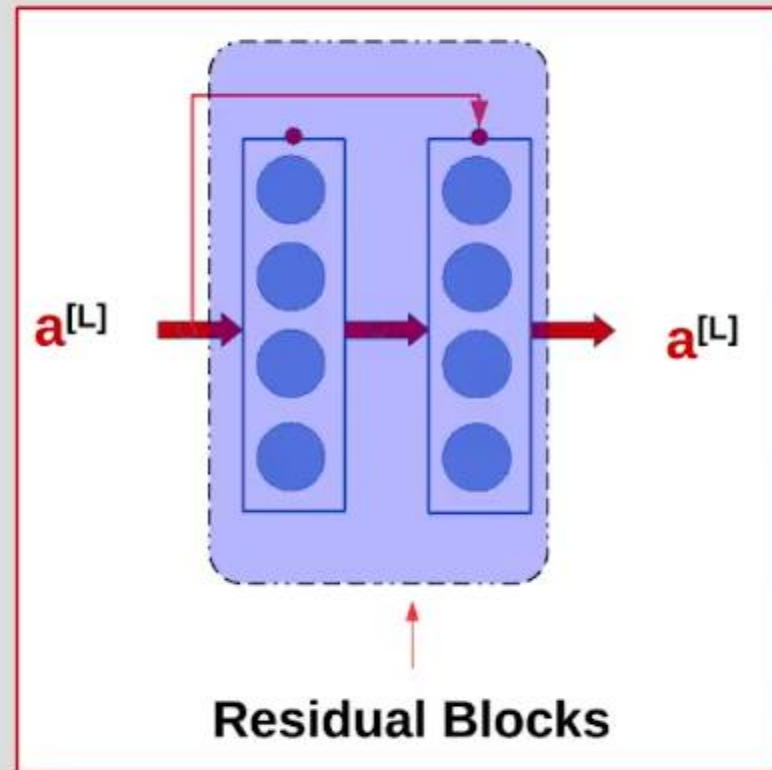
System Can Learn Identity Function

ReLU

$$R(z) = \max(0, z)$$

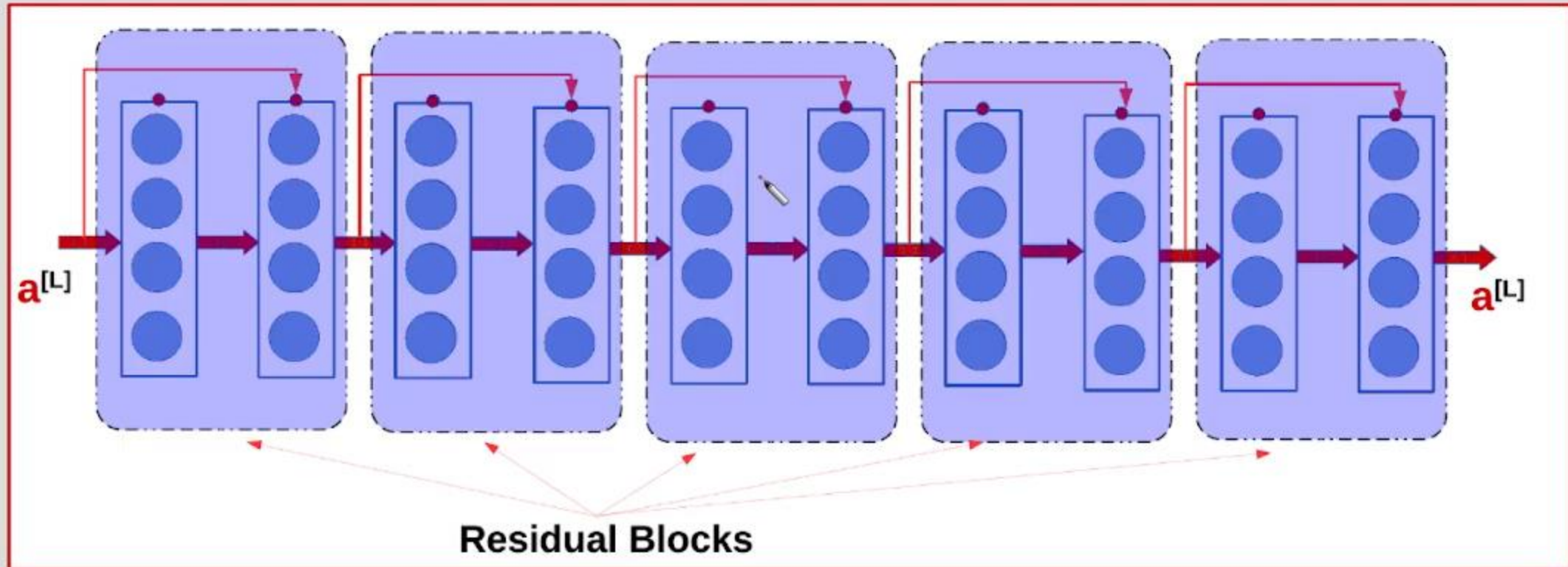


Residual Blocks Can Learn Identity Function



Residual block Can LEARN to Output same input

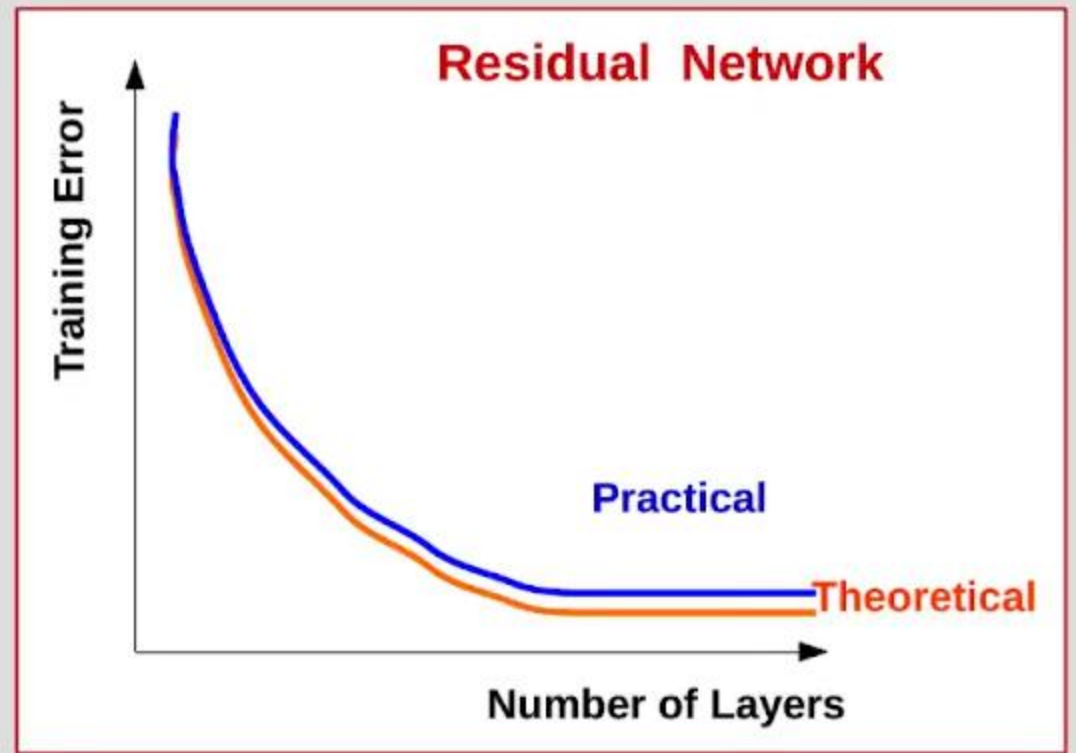
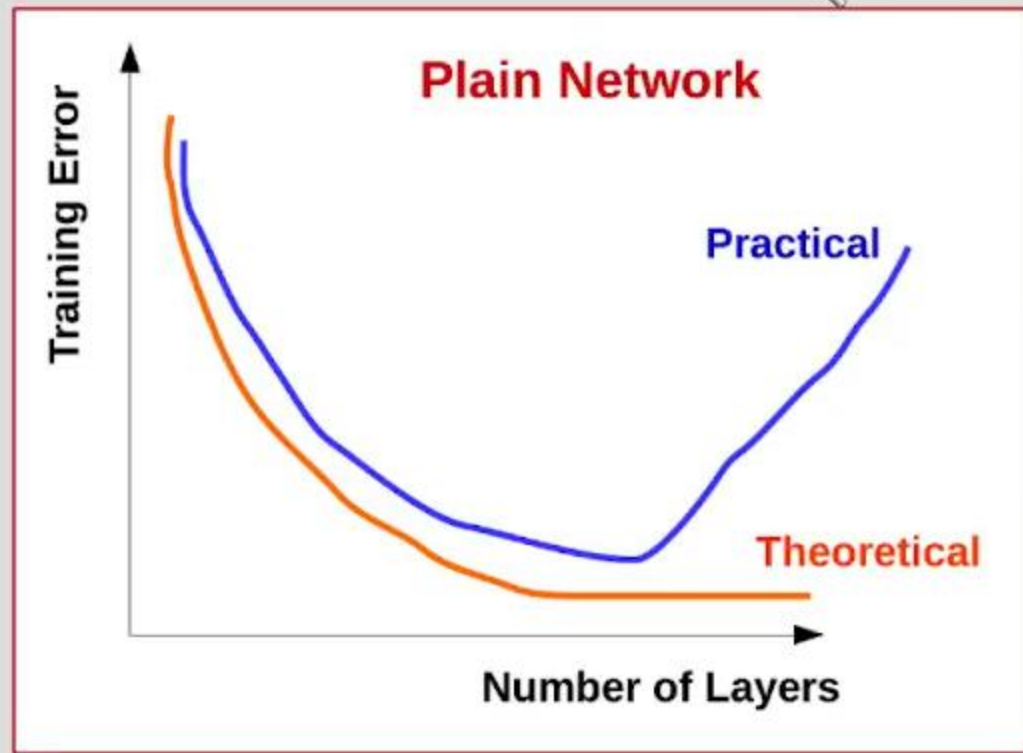
Cascaded Residual Blocks Can Learn Identity Function



Go Deeper will NOT Harm output of the Network

Plain vs Residual Network

effect of increasing Number of Layers



ResNet Architecture

