

# Deep Residual Learning

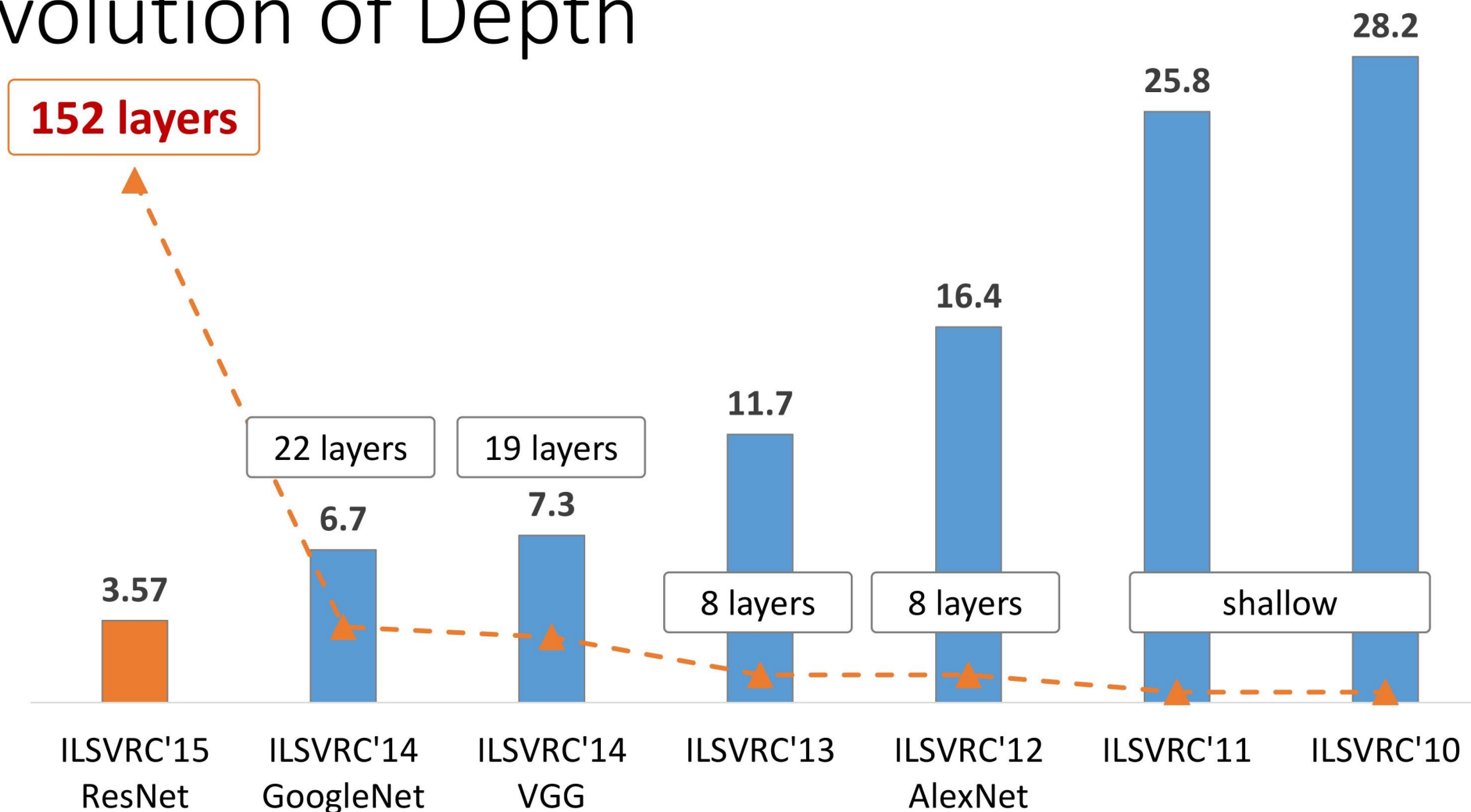
MSRA @ ILSVRC & COCO 2015 competitions

Kaiming He

with Xiangyu Zhang, Shaoqing Ren, Jifeng Dai, & Jian Sun

Microsoft Research Asia (MSRA)

# Revolution of Depth

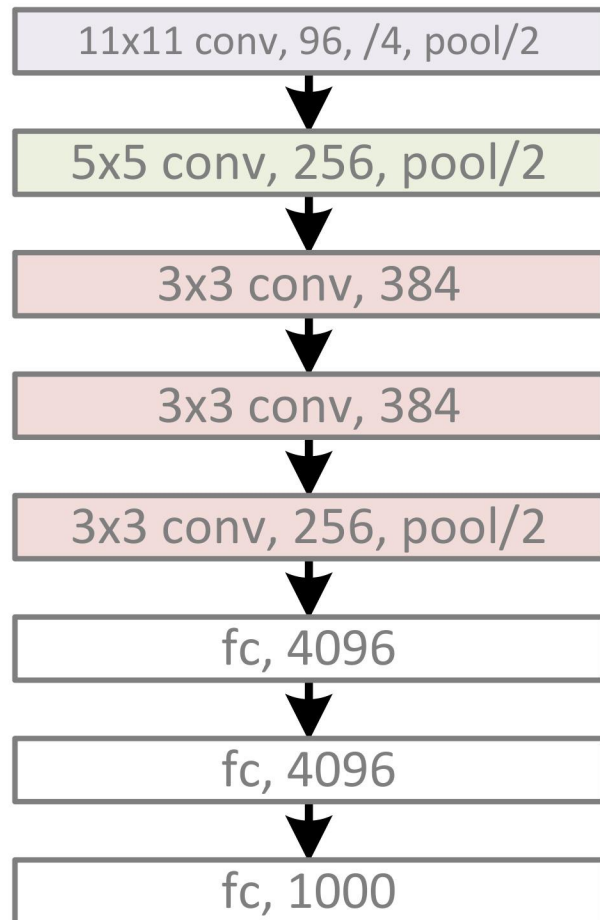


ImageNet Classification top-5 error (%)

Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". arXiv 2015.

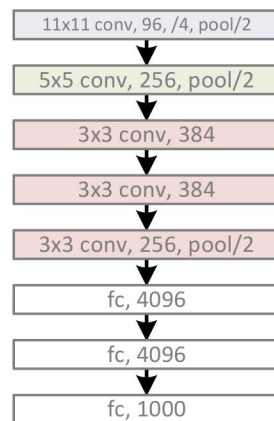
# Revolution of Depth

AlexNet, 8 layers  
(ILSVRC 2012)

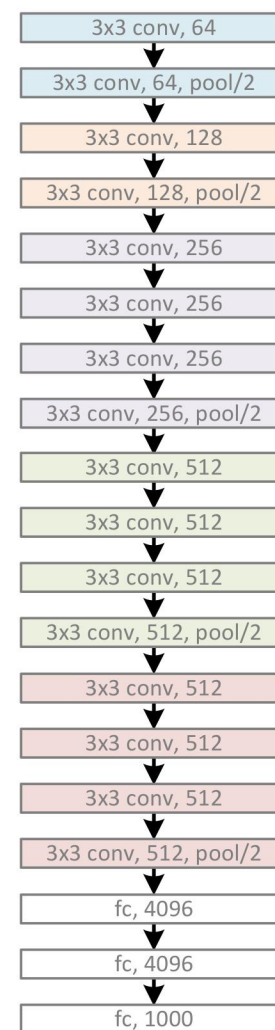


# Revolution of Depth

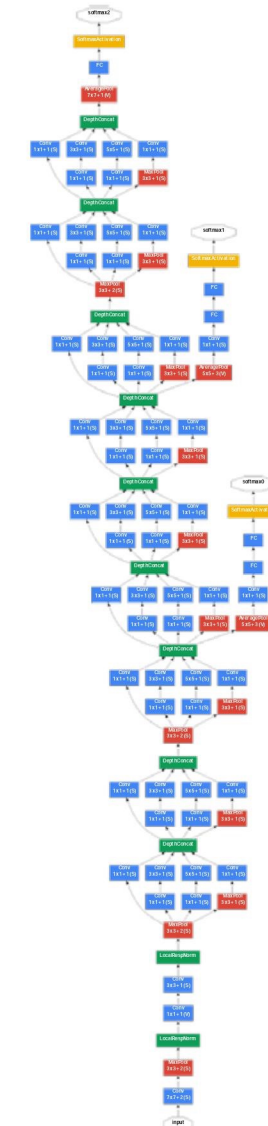
AlexNet, 8 layers (ILSVRC 2012)



VGG, 19 layers (ILSVRC 2014)



GoogleNet, 22 layers (ILSVRC 2014)



# Revolution of Depth

AlexNet, 8 layers  
(ILSVRC 2012)



VGG, 19 layers  
(ILSVRC 2014)

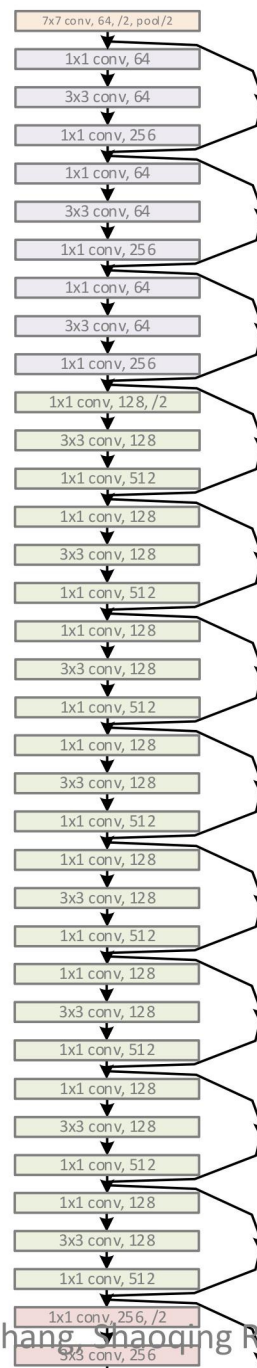


ResNet, **152 layers**  
(ILSVRC 2015)



# Revolution of Depth

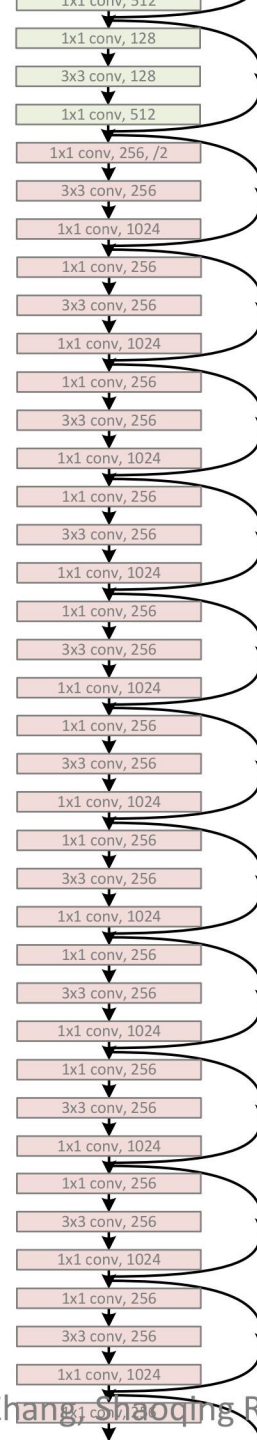
ResNet, 152 layers



(there was an animation here)

# Revolution of Depth

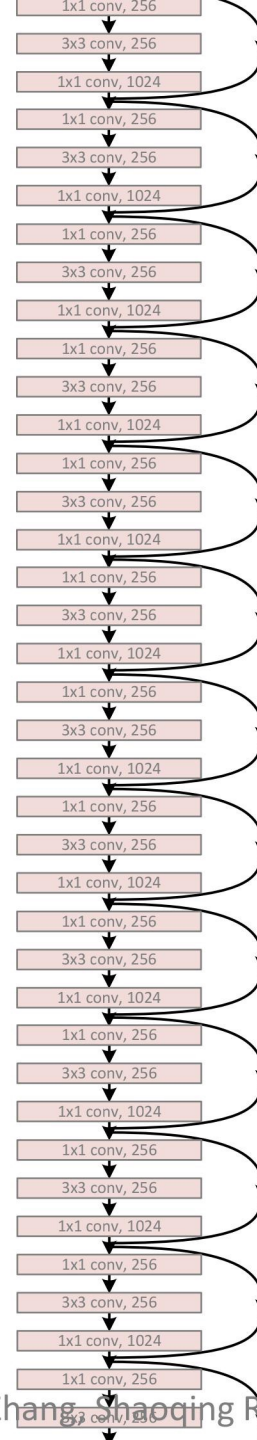
ResNet, 152 layers



(there was an animation here)

# Revolution of Depth

ResNet, 152 layers

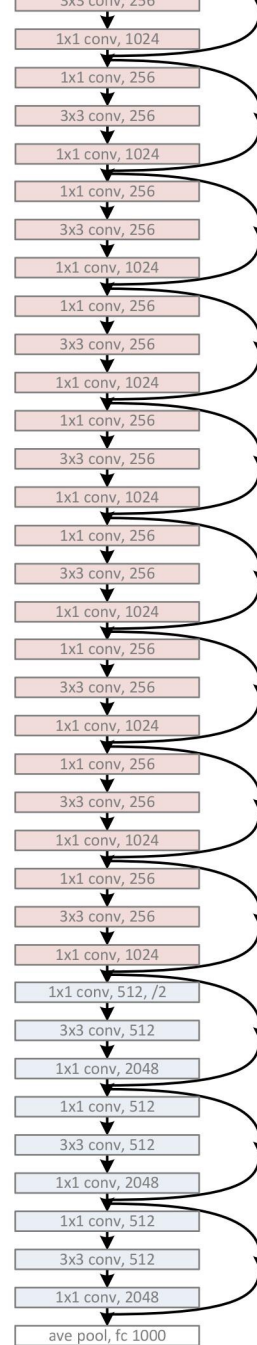


(there was an animation here)



# Revolution of Depth

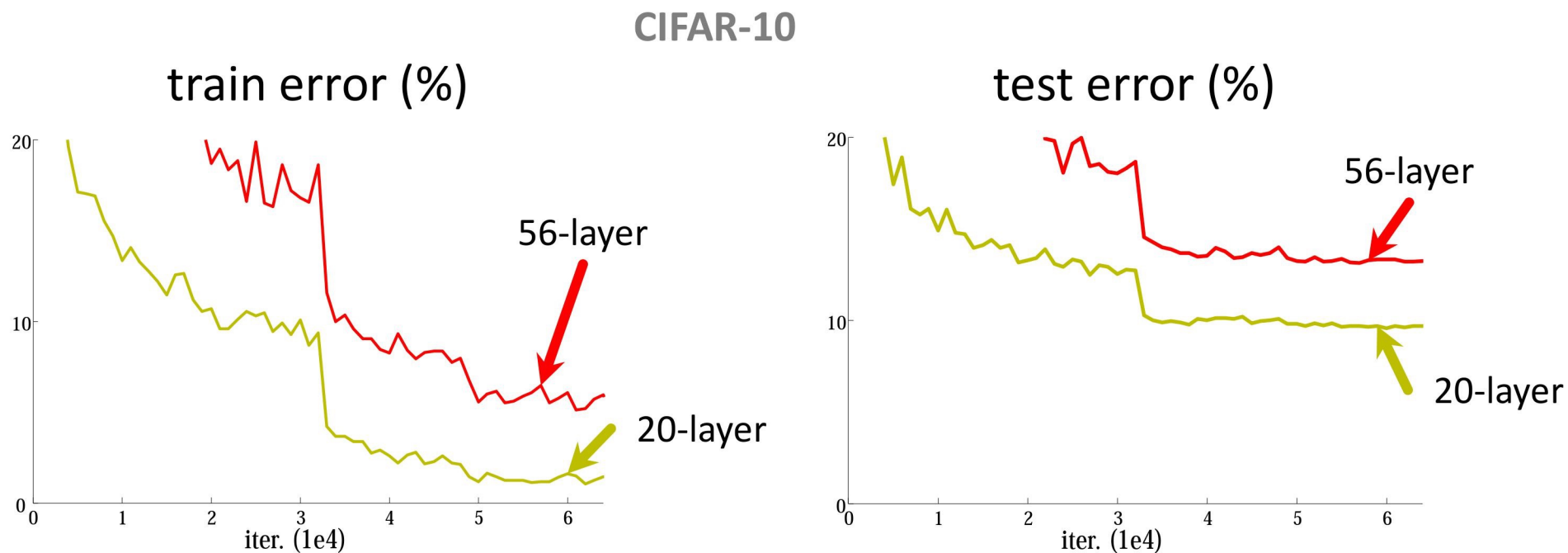
ResNet, 152 layers



(there was an animation here)

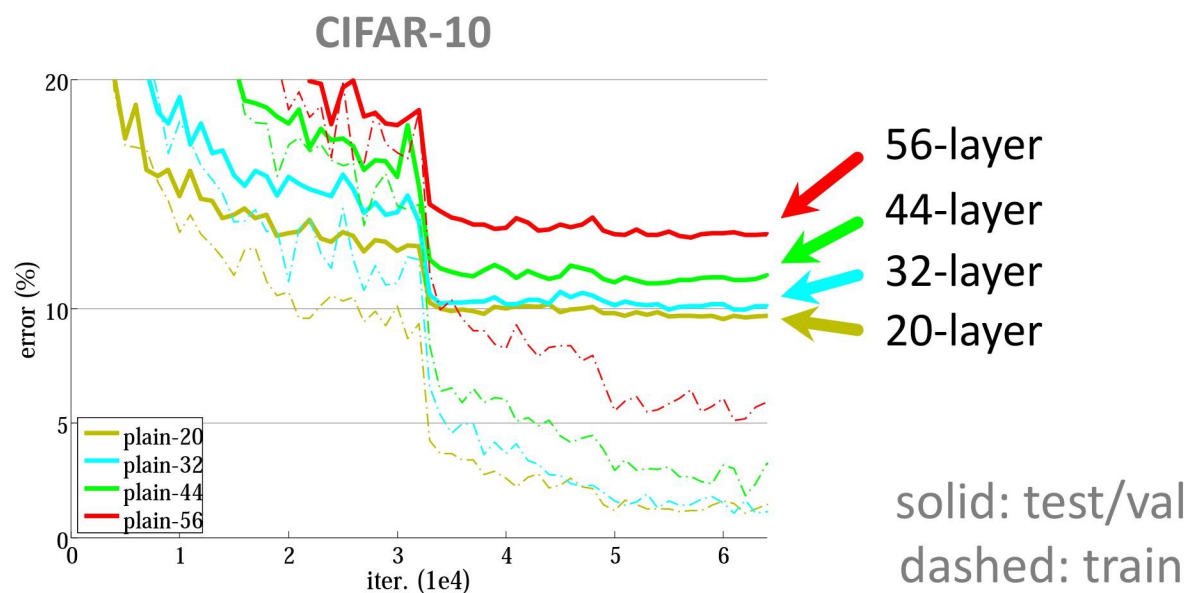
Is learning better networks  
as simple as stacking more layers?

# Simply stacking layers?



- *Plain* nets: stacking 3x3 conv layers...
- 56-layer net has **higher training error** and test error than 20-layer net

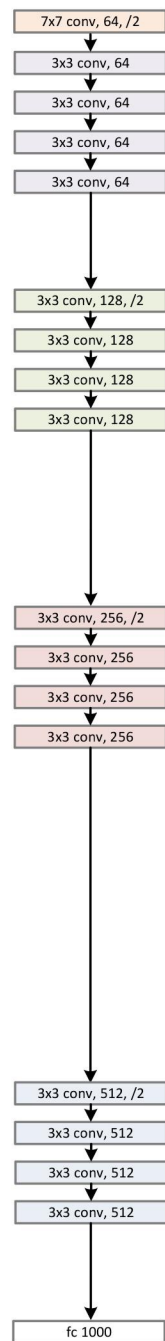
# Simply stacking layers?



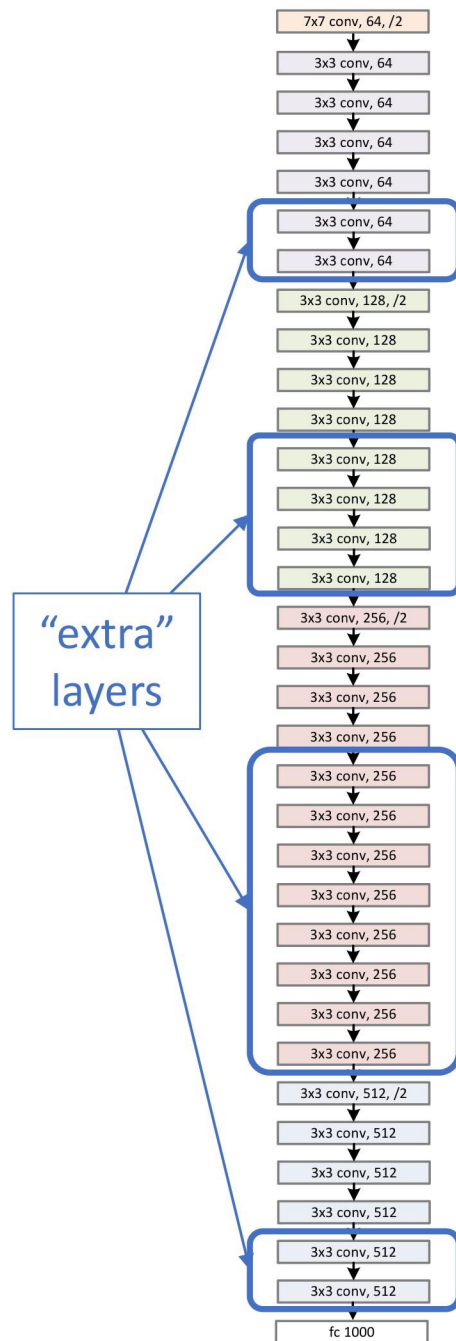
Vanishing/exploding gradient?

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

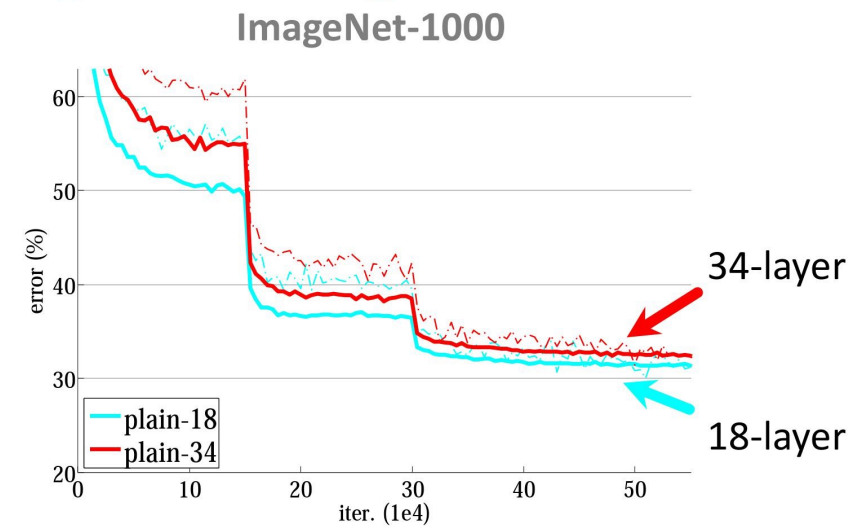
a shallower model  
(18 layers)



a deeper counterpart  
(34 layers)



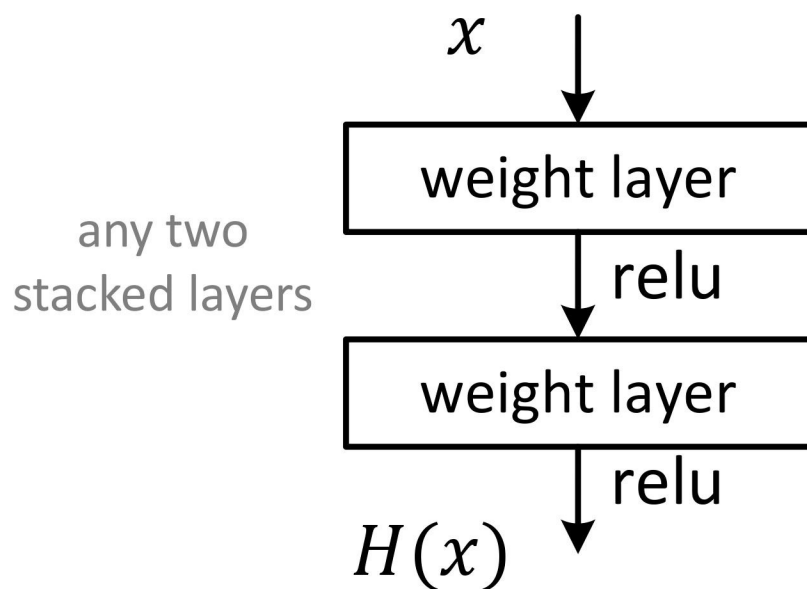
- A deeper model should not have **higher training error**



- Optimization difficulties:** solvers cannot find the solution when going deeper...

# Deep Residual Learning

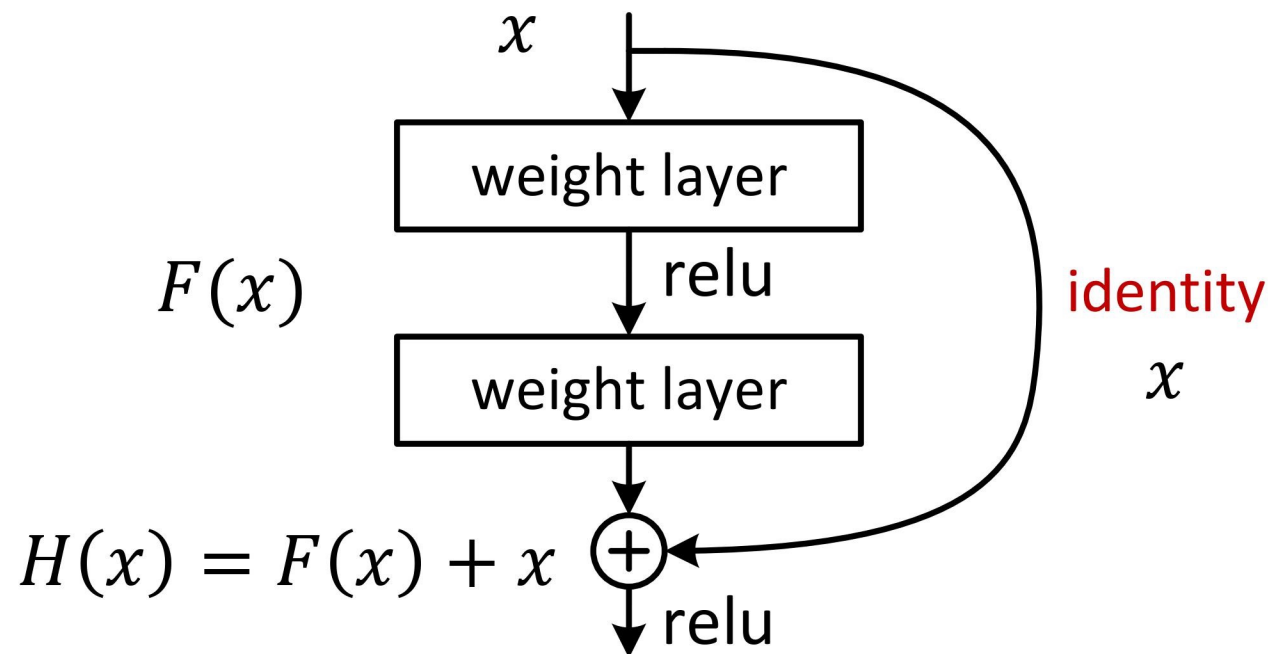
- Plain net



$H(x)$  is any desired mapping,  
hope the 2 weight layers fit  $H(x)$

# Deep Residual Learning

- Residual net



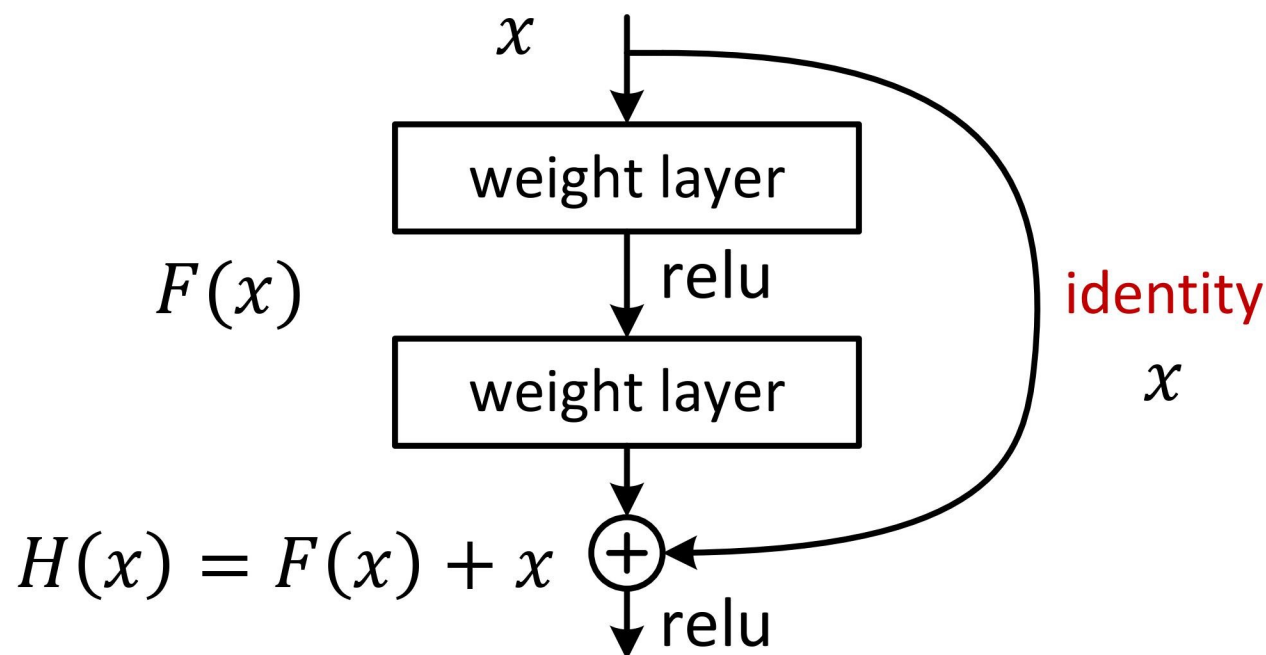
$H(x)$  is any desired mapping,  
~~hope the 2 weight layers fit  $H(x)$~~

hope the 2 weight layers fit  $F(x)$

$$\text{let } H(x) = F(x) + x$$

# Deep Residual Learning

- $F(x)$  is a **residual** mapping w.r.t. **identity**

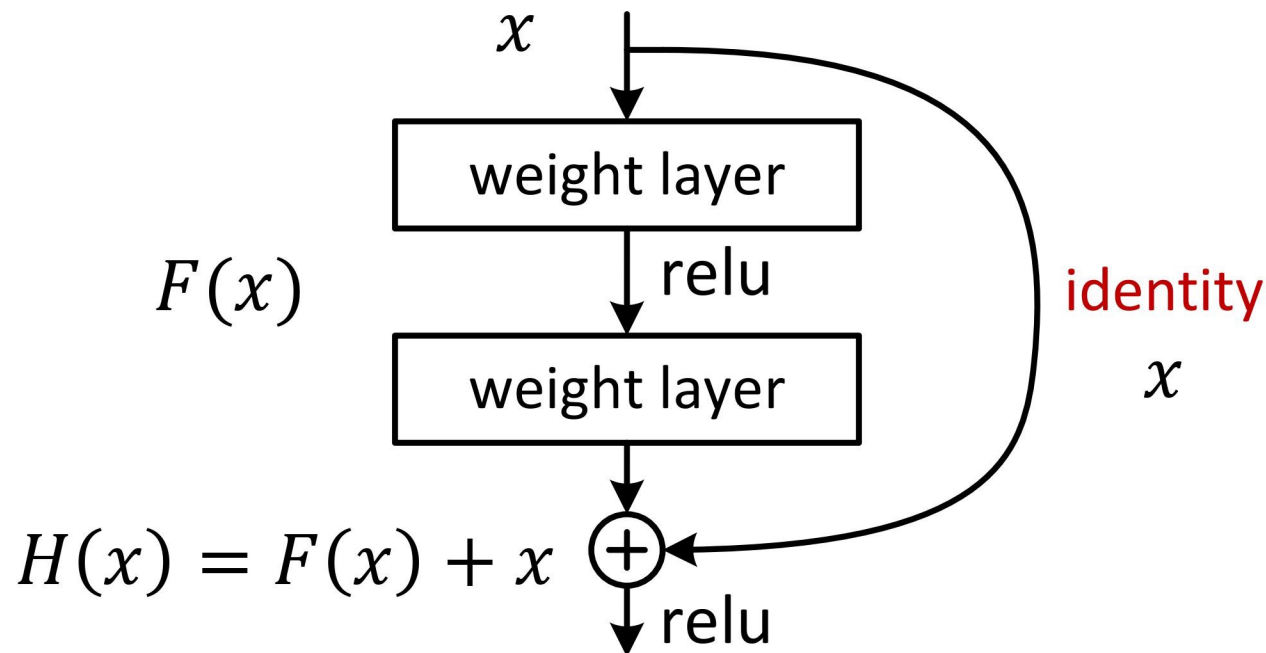


- It's an element wise addition, between the feature maps and channel by channel
- Resnet of 1 layer do not improve, minimum 2



# Deep Residual Learning

- $F(x)$  is a **residual** mapping w.r.t. **identity**



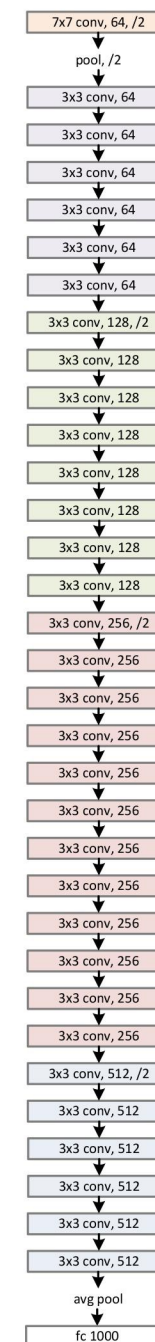
- If the dimension of  $x$  and  $F$  is different:

$$H(x) = \mathcal{F}(x, \{W_i\}) + W_s x$$

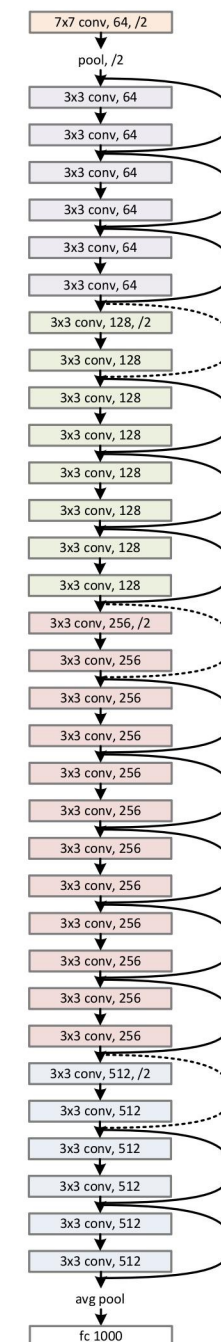
# Network “Design”

- Keep it simple
- Our basic design (VGG-style)
  - all 3x3 conv (almost)
  - spatial size /2 => # filters x2
  - **Simple design; just deep!**
- Other remarks:
  - no max pooling (almost)
  - no hidden fc
  - no dropout

plain net



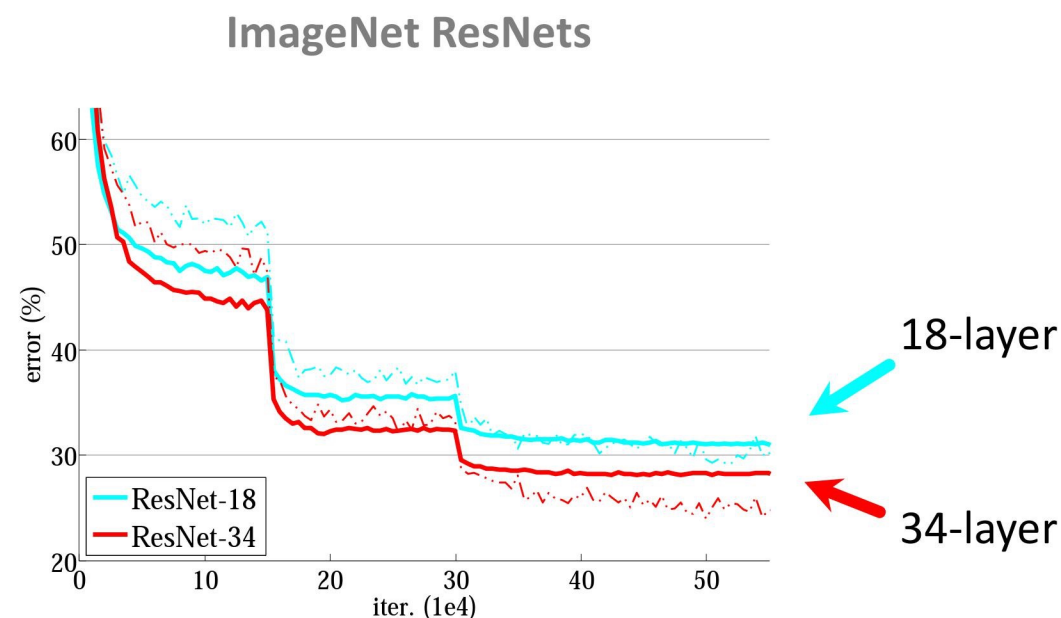
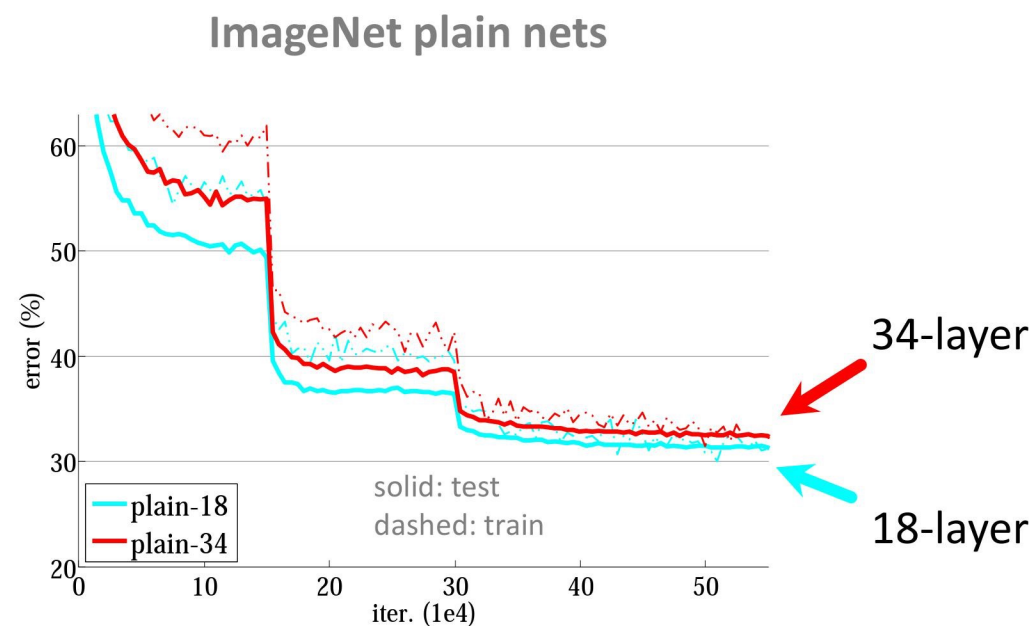
ResNet



# Training

- All plain/residual nets are trained **from scratch**
- All plain/residual nets use Batch Normalization
- Standard augmentation

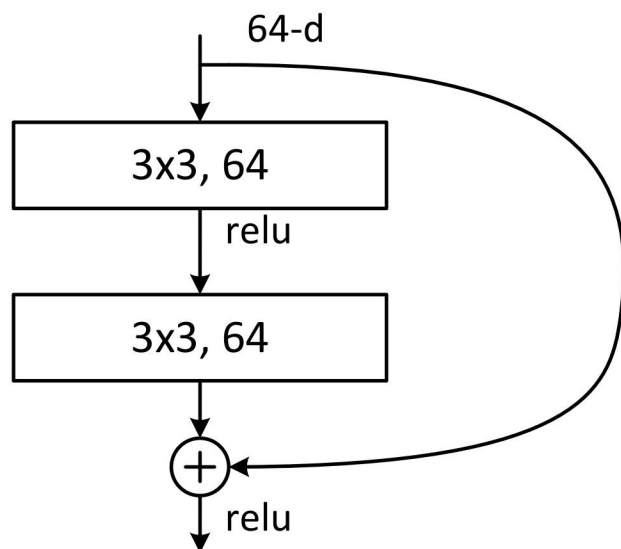
# ImageNet experiments



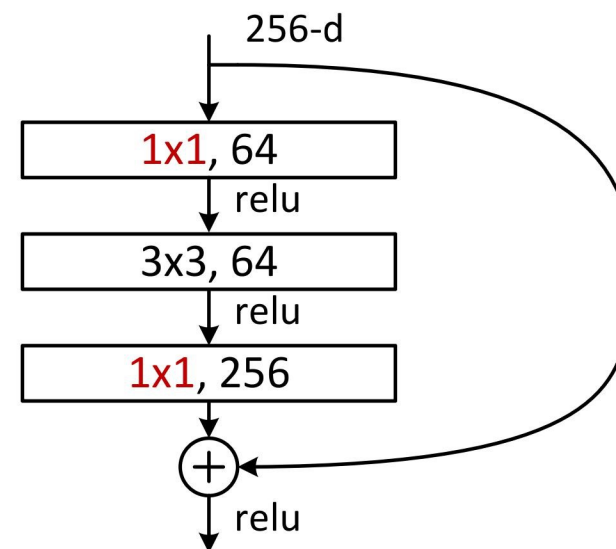
- Deep ResNets can be trained without difficulties
- Deeper ResNets have **lower training error**, and also lower test error

# ImageNet experiments

- A practical design of going deeper



all-3x3



**bottleneck**  
(for ResNet-50/101/152)

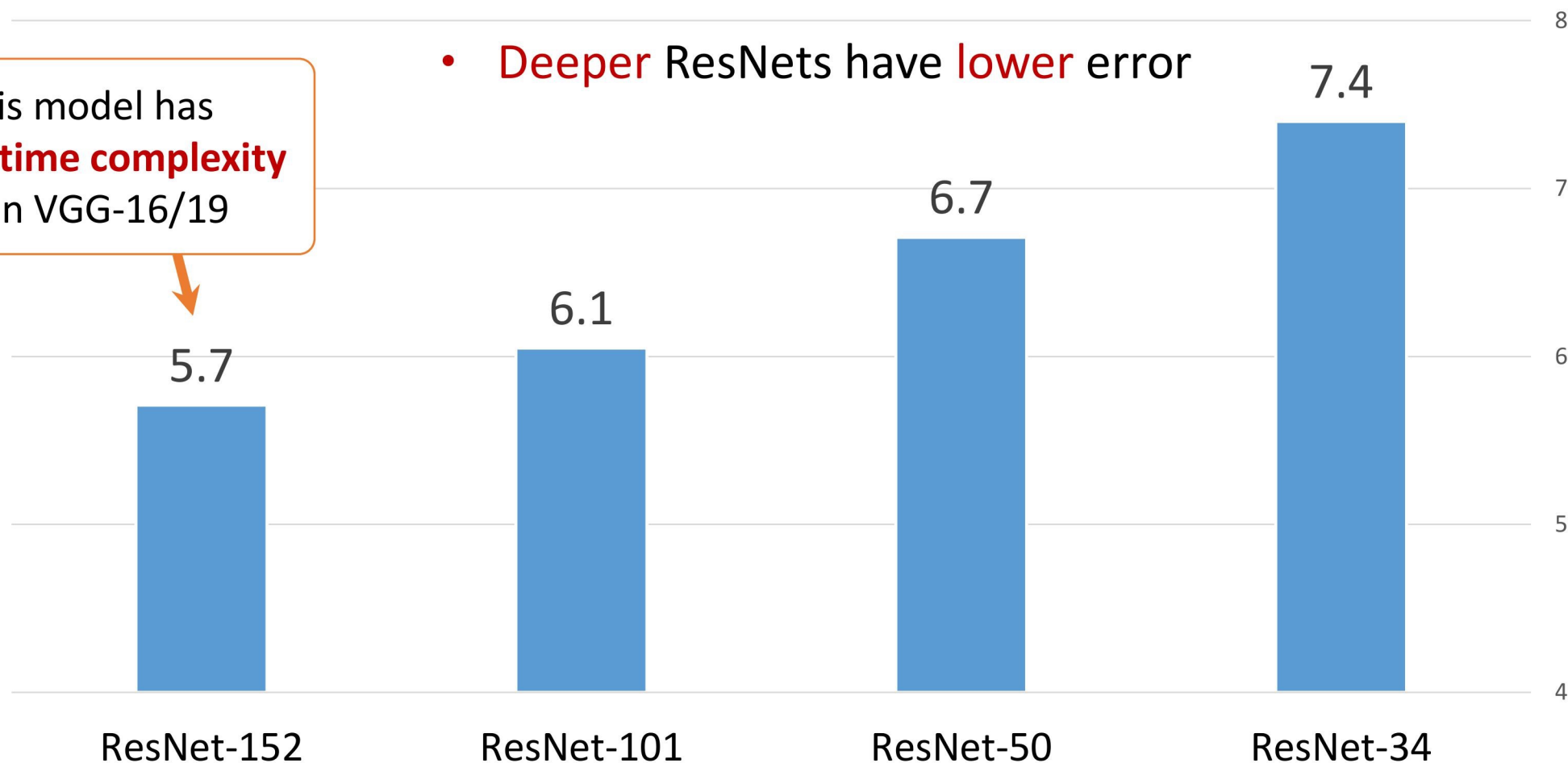
layer name	output size	18-layer	34-layer	50-layer	101-layer	152-layer
conv1	112×112	7×7, 64, stride 2				
conv2_x	56×56	3×3 max pool, stride 2				
		$\begin{bmatrix} 3\times 3, 64 \\ 3\times 3, 64 \end{bmatrix} \times 2$	$\begin{bmatrix} 3\times 3, 64 \\ 3\times 3, 64 \end{bmatrix} \times 3$	$\begin{bmatrix} 1\times 1, 64 \\ 3\times 3, 64 \\ 1\times 1, 256 \end{bmatrix} \times 3$	$\begin{bmatrix} 1\times 1, 64 \\ 3\times 3, 64 \\ 1\times 1, 256 \end{bmatrix} \times 3$	$\begin{bmatrix} 1\times 1, 64 \\ 3\times 3, 64 \\ 1\times 1, 256 \end{bmatrix} \times 3$
conv3_x	28×28	$\begin{bmatrix} 3\times 3, 128 \\ 3\times 3, 128 \end{bmatrix} \times 2$	$\begin{bmatrix} 3\times 3, 128 \\ 3\times 3, 128 \end{bmatrix} \times 4$	$\begin{bmatrix} 1\times 1, 128 \\ 3\times 3, 128 \\ 1\times 1, 512 \end{bmatrix} \times 4$	$\begin{bmatrix} 1\times 1, 128 \\ 3\times 3, 128 \\ 1\times 1, 512 \end{bmatrix} \times 4$	$\begin{bmatrix} 1\times 1, 128 \\ 3\times 3, 128 \\ 1\times 1, 512 \end{bmatrix} \times 8$
conv4_x	14×14	$\begin{bmatrix} 3\times 3, 256 \\ 3\times 3, 256 \end{bmatrix} \times 2$	$\begin{bmatrix} 3\times 3, 256 \\ 3\times 3, 256 \end{bmatrix} \times 6$	$\begin{bmatrix} 1\times 1, 256 \\ 3\times 3, 256 \\ 1\times 1, 1024 \end{bmatrix} \times 6$	$\begin{bmatrix} 1\times 1, 256 \\ 3\times 3, 256 \\ 1\times 1, 1024 \end{bmatrix} \times 23$	$\begin{bmatrix} 1\times 1, 256 \\ 3\times 3, 256 \\ 1\times 1, 1024 \end{bmatrix} \times 36$
conv5_x	7×7	$\begin{bmatrix} 3\times 3, 512 \\ 3\times 3, 512 \end{bmatrix} \times 2$	$\begin{bmatrix} 3\times 3, 512 \\ 3\times 3, 512 \end{bmatrix} \times 3$	$\begin{bmatrix} 1\times 1, 512 \\ 3\times 3, 512 \\ 1\times 1, 2048 \end{bmatrix} \times 3$	$\begin{bmatrix} 1\times 1, 512 \\ 3\times 3, 512 \\ 1\times 1, 2048 \end{bmatrix} \times 3$	$\begin{bmatrix} 1\times 1, 512 \\ 3\times 3, 512 \\ 1\times 1, 2048 \end{bmatrix} \times 3$
	1×1	average pool, 1000-d fc, softmax				
FLOPs		$1.8 \times 10^9$	$3.6 \times 10^9$	$3.8 \times 10^9$	$7.6 \times 10^9$	$11.3 \times 10^9$

The 152-layers ResNet still has lower complexity than VGG-16/19 (15.3/19.6 billion FLOPs) !!!

# ImageNet experiments

- Deeper ResNets have **lower** error

this model has  
**lower time complexity**  
than VGG-16/19



**10-crop** testing, top-5 val error (%)

“*Features matter.*” (quote [Girshick et al. 2014], the R-CNN paper)

task	2nd-place winner	MSRA	margin (relative)
ImageNet Localization (top-5 error)	12.0	9.0	<b>27%</b>
ImageNet Detection (mAP@.5)	53.6	62.1	<b>16%</b>
COCO Detection (mAP@.5:.95)	33.5	37.3	<b>11%</b>
COCO Segmentation (mAP@.5:.95)	25.1	28.2	<b>12%</b>

- Our results are all based on **ResNet-101**
- Our features are **well transferrable**