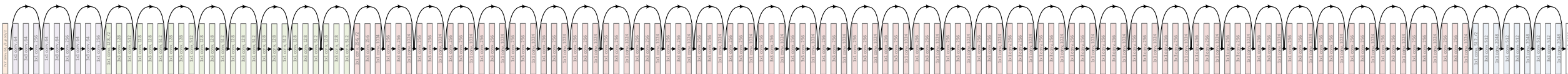# Deep Residual Learning for Image Recognition

**Kaiming He**, Xiangyu Zhang, Shaoqing Ren, Jian Sun

work done at
Microsoft Research Asia
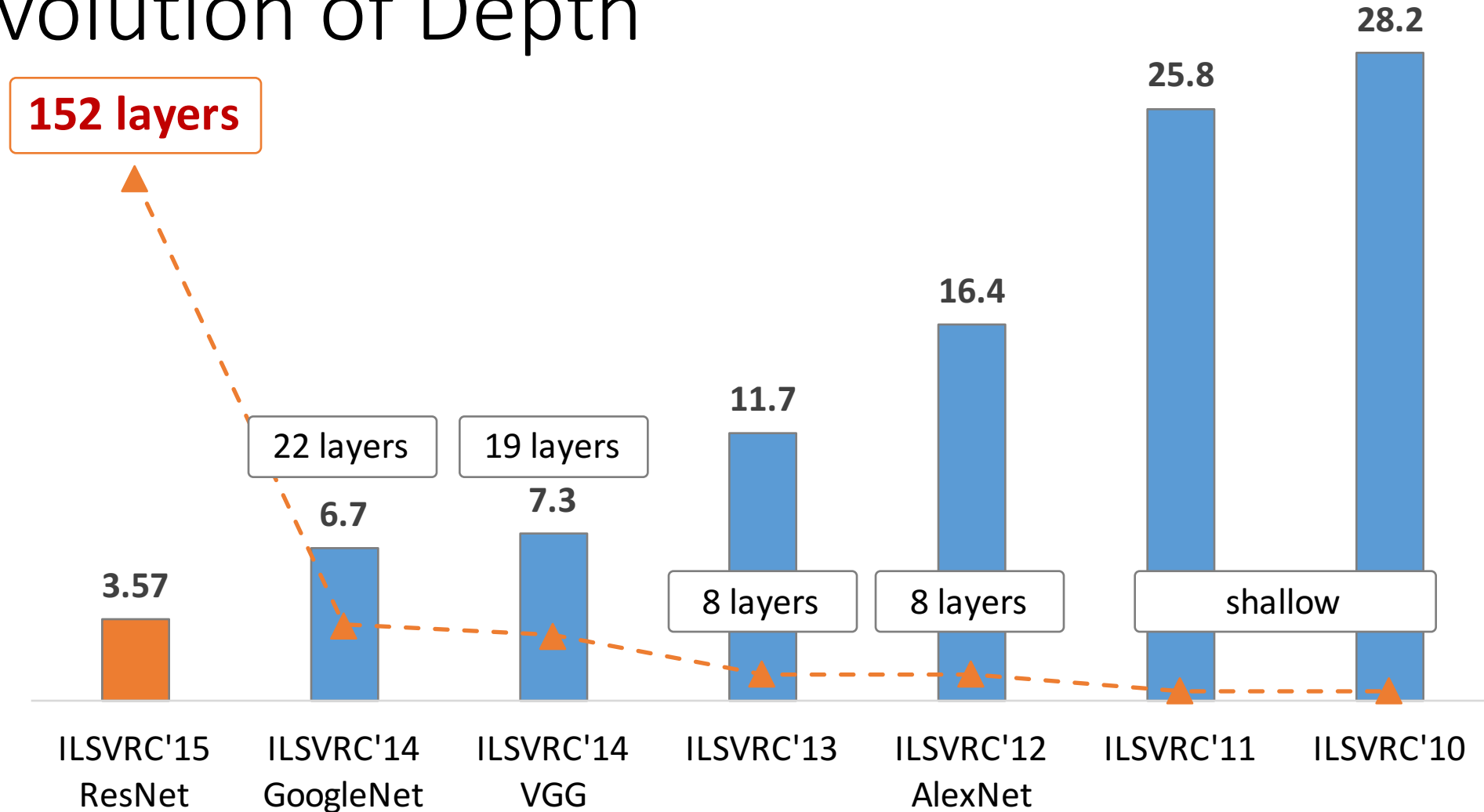
# ResNet @ ILSVRC & COCO 2015 Competitions

**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
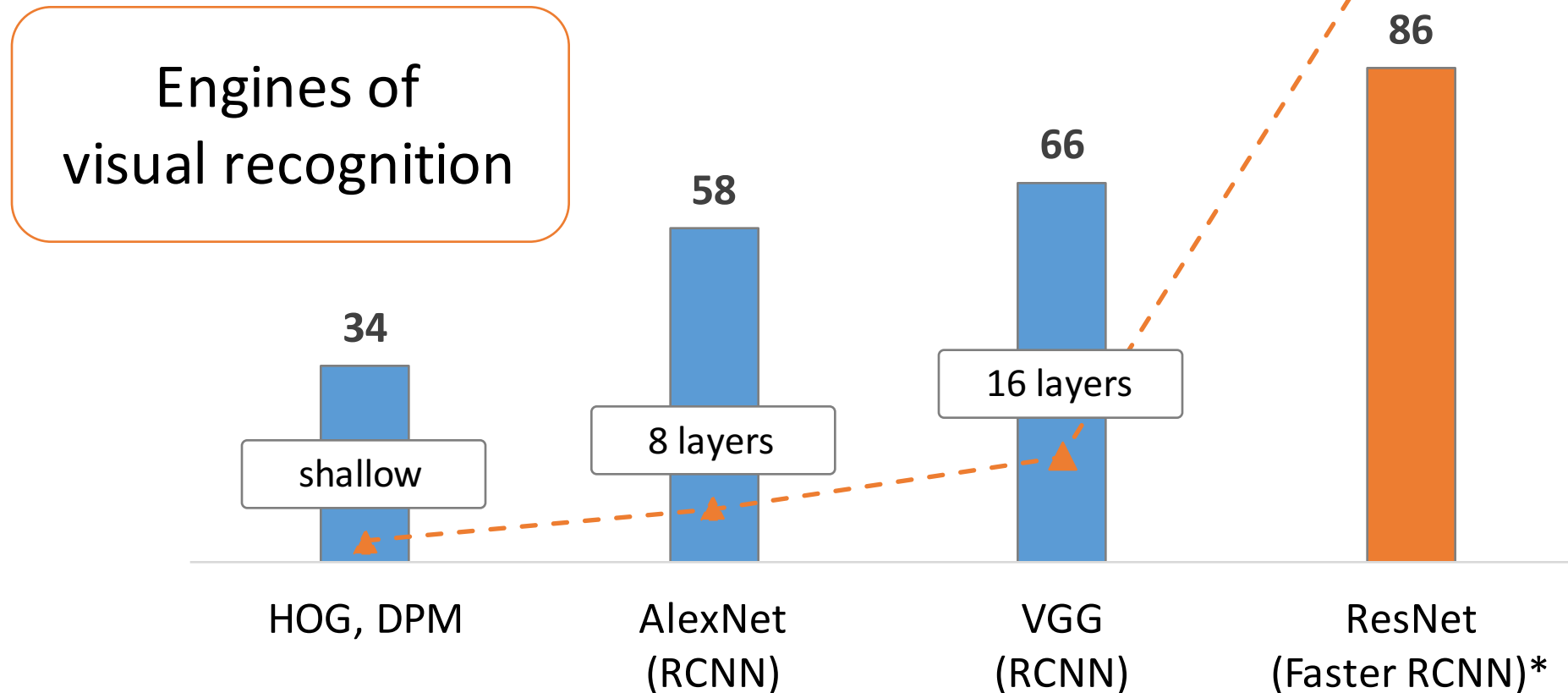- COCO Segmentation: 12% better than 2nd

*improvements are relative numbers

Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". CVPR 2016.

# Revolution of Depth



**152 layers**

**3.57**  ILSVRC'15 ResNet

22 layers — **6.7**  ILSVRC'14 GoogleNet

19 layers — **7.3**  ILSVRC'14 VGG

8 layers — **11.7**  ILSVRC'13

8 layers — **16.4**  ILSVRC'12 AlexNet

shallow — **25.8**  ILSVRC'11

**28.2**  ILSVRC'10

ImageNet Classification top-5 error (%)

Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". CVPR 2016.

# Revolution of Depth

Engines of visual recognition

**101 layers**

86

66

58

34

16 layers

8 layers

shallow

HOG, DPM

AlexNet (RCNN)
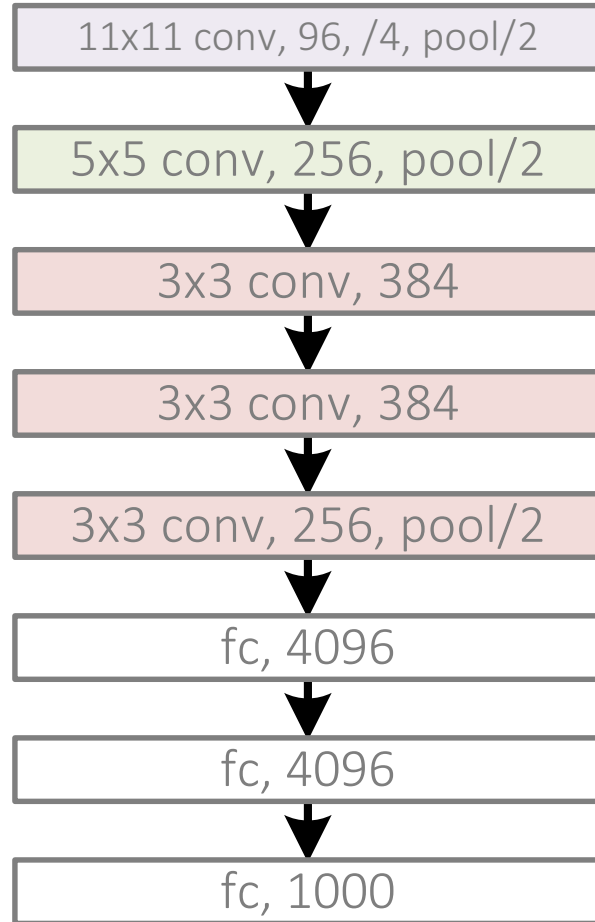
VGG (RCNN)

ResNet (Faster RCNN)*

PASCAL VOC 2007 **Object Detection** mAP (%)

*w/ other improvements & more data

Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". CVPR 2016.

# Revolution of Depth

**AlexNet, 8 layers**
**(ILSVRC 2012)**

11x11 conv, 96, /4, pool/2

5x5 conv, 256, pool/2

3x3 conv, 384

3x3 conv, 384

3x3 conv, 256, pool/2

fc, 4096

fc, 4096

fc, 1000

Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". CVPR 2016.

# Revolution of Depth

AlexNet, 8 layers
(ILSVRC 2012)

| |
|---|
| 11x11 conv, 96, /4, pool/2 |
| 5x5 conv, 256, pool/2 |
| 3x3 conv, 384 |
| 3x3 conv, 384 |
| 3x3 conv, 256, pool/2 |
| fc, 4096 |
| fc, 4096 |
| fc, 1000 |

VGG, 19 layers
(ILSVRC 2014)

| |
|---|
| 3x3 conv, 64 |
| 3x3 conv, 64, pool/2 |
| 3x3 conv, 128 |
| 3x3 conv, 128, pool/2 |
| 3x3 conv, 256 |
| 3x3 conv, 256 |
| 3x3 conv, 256 |
| 3x3 conv, 256, pool/2 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512, pool/2 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512, pool/2 |
| fc, 4096 |
| fc, 4096 |
| fc, 1000 |

GoogleNet, 22 layers
(ILSVRC 2014)

Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". CVPR 2016.

# Revolution of Depth

AlexNet, 8 layers
(ILSVRC 2012)

VGG, 19 layers
(ILSVRC 2014)

ResNet, 152 layers
(ILSVRC 2015)

Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". CVPR 2016.

# Is learning better networks as simple as stacking more layers?

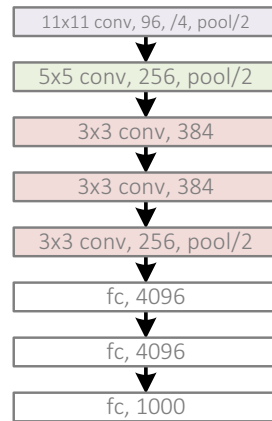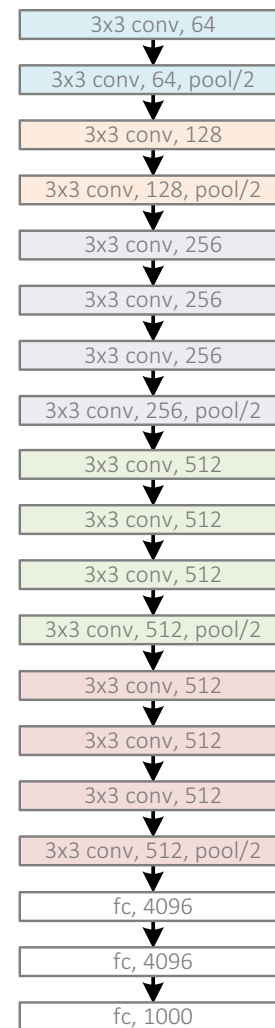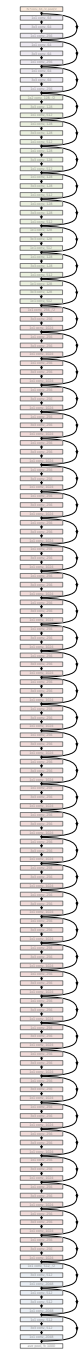Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". CVPR 2016.

# Simply stacking layers?

**CIFAR-10**

train error (%)

test error (%)



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

Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". CVPR 2016.

# Simply stacking layers?



CIFAR-10

56-layer
44-layer
32-layer
20-layer

plain-20
plain-32
plain-44
plain-56

solid: test/val
dashed: train

ImageNet-1000

34-layer
18-layer

plain-18
plain-34

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

Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". CVPR 2016.

a shallower model (18 layers)

a deeper counterpart (34 layers)

"extra" layers

- Richer solution space

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

- A solution *by construction*:
  - original layers: copied from a learned shallower model
  - extra layers: set as identity
  - at least the same training error

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

Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". CVPR 2016.

# Deep Residual Learning

- Plaint net

$x$

weight layer

relu

weight layer

relu

$H(x)$

any two
stacked layers

$H(x)$ is any desired mapping,

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

Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". CVPR 2016.

# Deep Residual Learning

- **Residual** net



$$x$$

weight layer

relu

$$F(x)$$

weight layer

identity

$$x$$

$$H(x) = F(x) + x \quad \oplus$$

relu

$H(x)$ is any desired mapping,

~~hope the 2 weight layers fit $H(x)$~~

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

let $H(x) = F(x) + x$

Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". CVPR 2016.

# Deep Residual Learning

- $F(x)$ is a <span style="color:red">residual</span> mapping w.r.t. <span style="color:red">identity</span>



$$H(x) = F(x) + x$$

- If identity were optimal, easy to set weights as 0

- If optimal mapping is closer to identity, easier to find small fluctuations

Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". CVPR 2016.

# Network "Design"

plain net                    ResNet

- Keep it simple

- Our basic design (VGG-style)
  - all 3x3 conv (almost)
  - spatial size /2 => # filters x2
  - Simple design; just deep!

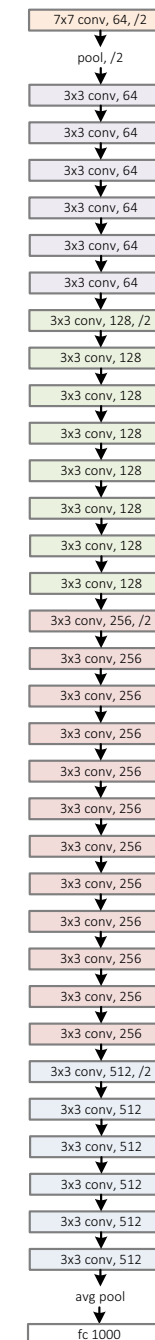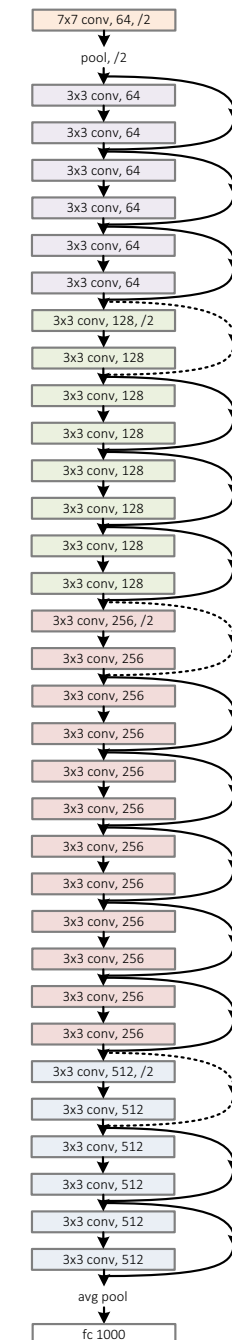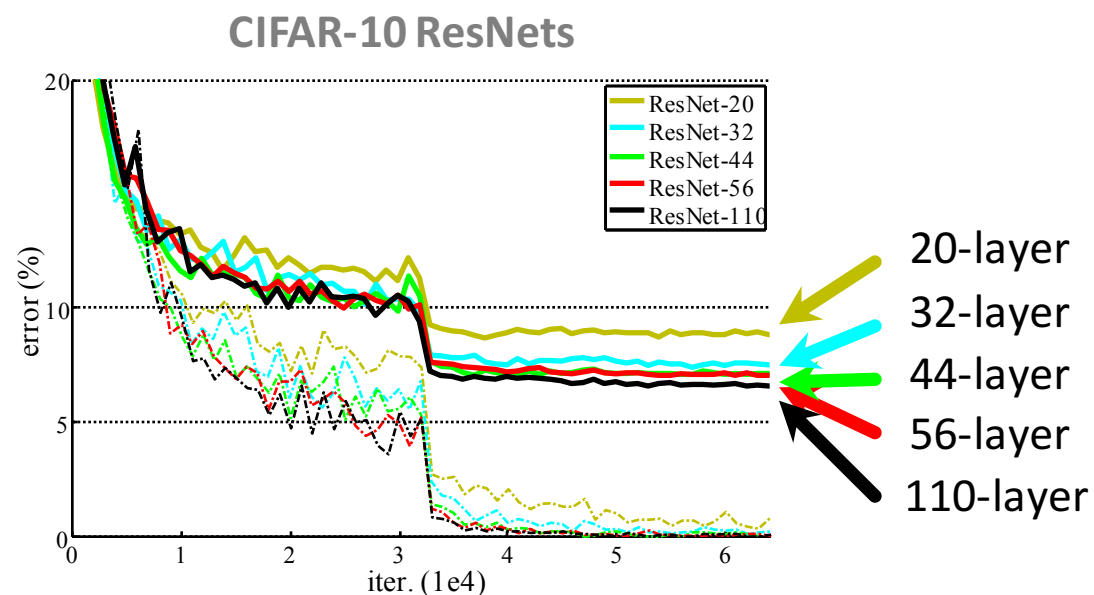| plain net | ResNet |
|---|---|
| 7x7 conv, 64, /2 | 7x7 conv, 64, /2 |
| pool, /2 | pool, /2 |
| 3x3 conv, 64 | 3x3 conv, 64 |
| 3x3 conv, 64 | 3x3 conv, 64 |
| 3x3 conv, 64 | 3x3 conv, 64 |
| 3x3 conv, 64 | 3x3 conv, 64 |
| 3x3 conv, 64 | 3x3 conv, 64 |
| 3x3 conv, 64 | 3x3 conv, 64 |
| 3x3 conv, 128, /2 | 3x3 conv, 128, /2 |
| 3x3 conv, 128 | 3x3 conv, 128 |
| 3x3 conv, 128 | 3x3 conv, 128 |
| 3x3 conv, 128 | 3x3 conv, 128 |
| 3x3 conv, 128 | 3x3 conv, 128 |
| 3x3 conv, 128 | 3x3 conv, 128 |
| 3x3 conv, 128 | 3x3 conv, 128 |
| 3x3 conv, 128 | 3x3 conv, 128 |
| 3x3 conv, 256, /2 | 3x3 conv, 256, /2 |
| 3x3 conv, 256 | 3x3 conv, 256 |
| 3x3 conv, 256 | 3x3 conv, 256 |
| 3x3 conv, 256 | 3x3 conv, 256 |
| 3x3 conv, 256 | 3x3 conv, 256 |
| 3x3 conv, 256 | 3x3 conv, 256 |
| 3x3 conv, 256 | 3x3 conv, 256 |
| 3x3 conv, 256 | 3x3 conv, 256 |
| 3x3 conv, 256 | 3x3 conv, 256 |
| 3x3 conv, 256 | 3x3 conv, 256 |
| 3x3 conv, 256 | 3x3 conv, 256 |
| 3x3 conv, 256 | 3x3 conv, 256 |
| 3x3 conv, 512, /2 | 3x3 conv, 512, /2 |
| 3x3 conv, 512 | 3x3 conv, 512 |
| 3x3 conv, 512 | 3x3 conv, 512 |
| 3x3 conv, 512 | 3x3 conv, 512 |
| 3x3 conv, 512 | 3x3 conv, 512 |
| 3x3 conv, 512 | 3x3 conv, 512 |
| avg pool | avg pool |
| fc 1000 | fc 1000 |

Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". CVPR 2016.

# CIFAR-10 experiments



**CIFAR-10 plain nets** — plot: error (%) vs iter. (1e4). Legend: plain-20, plain-32, plain-44, plain-56. Arrows: 56-layer, 44-layer, 32-layer, 20-layer. solid: test, dashed: train

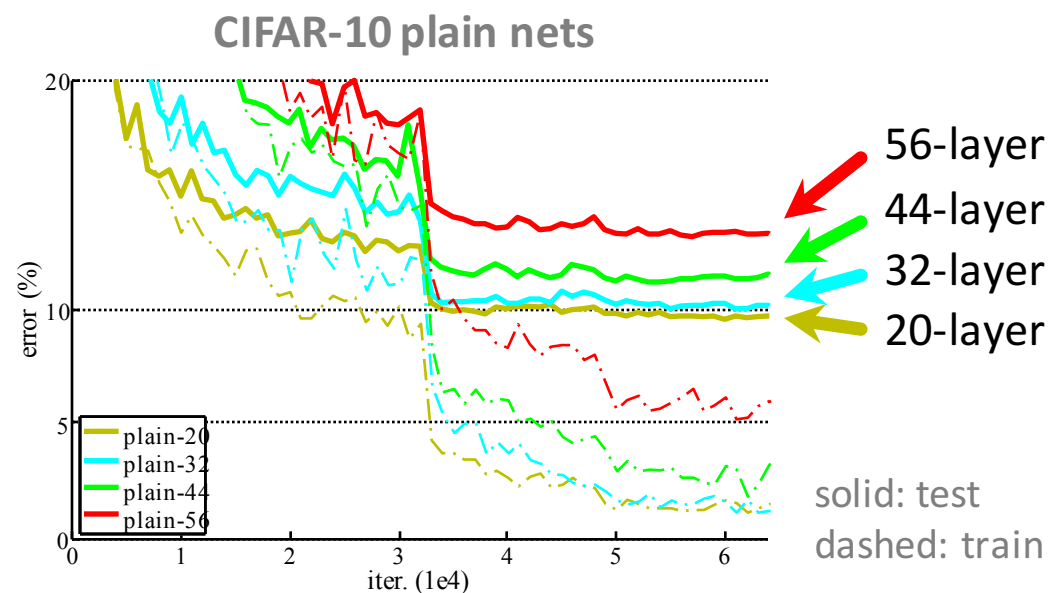**CIFAR-10 ResNets** — plot: error (%) vs iter. (1e4). Legend: ResNet-20, ResNet-32, ResNet-44, ResNet-56, ResNet-110. Arrows: 20-layer, 32-layer, 44-layer, 56-layer, 110-layer

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

Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". CVPR 2016.
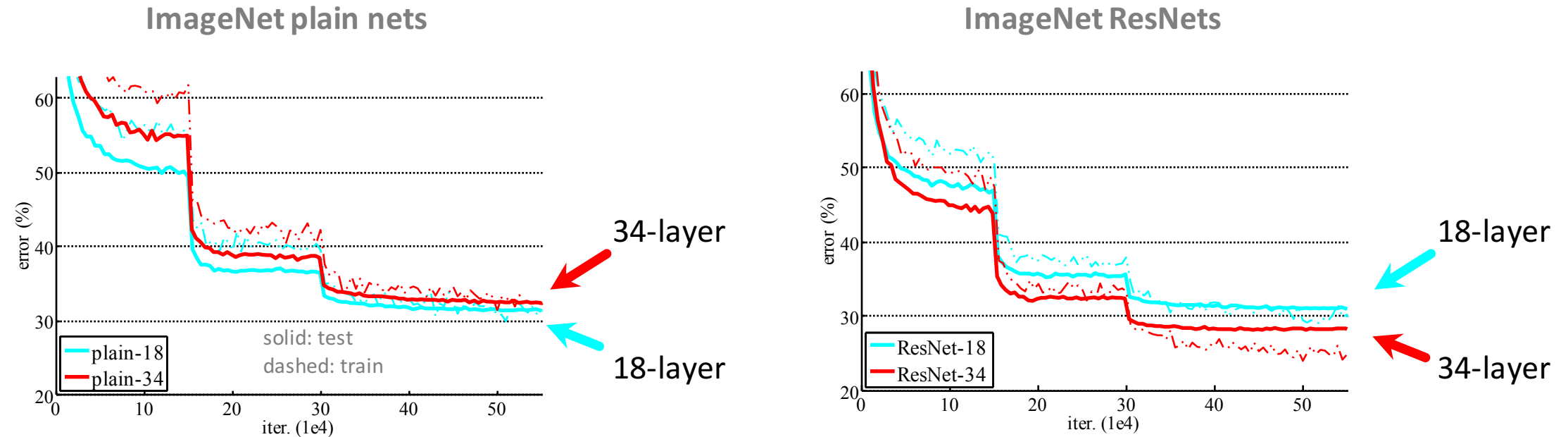
# ImageNet experiments



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

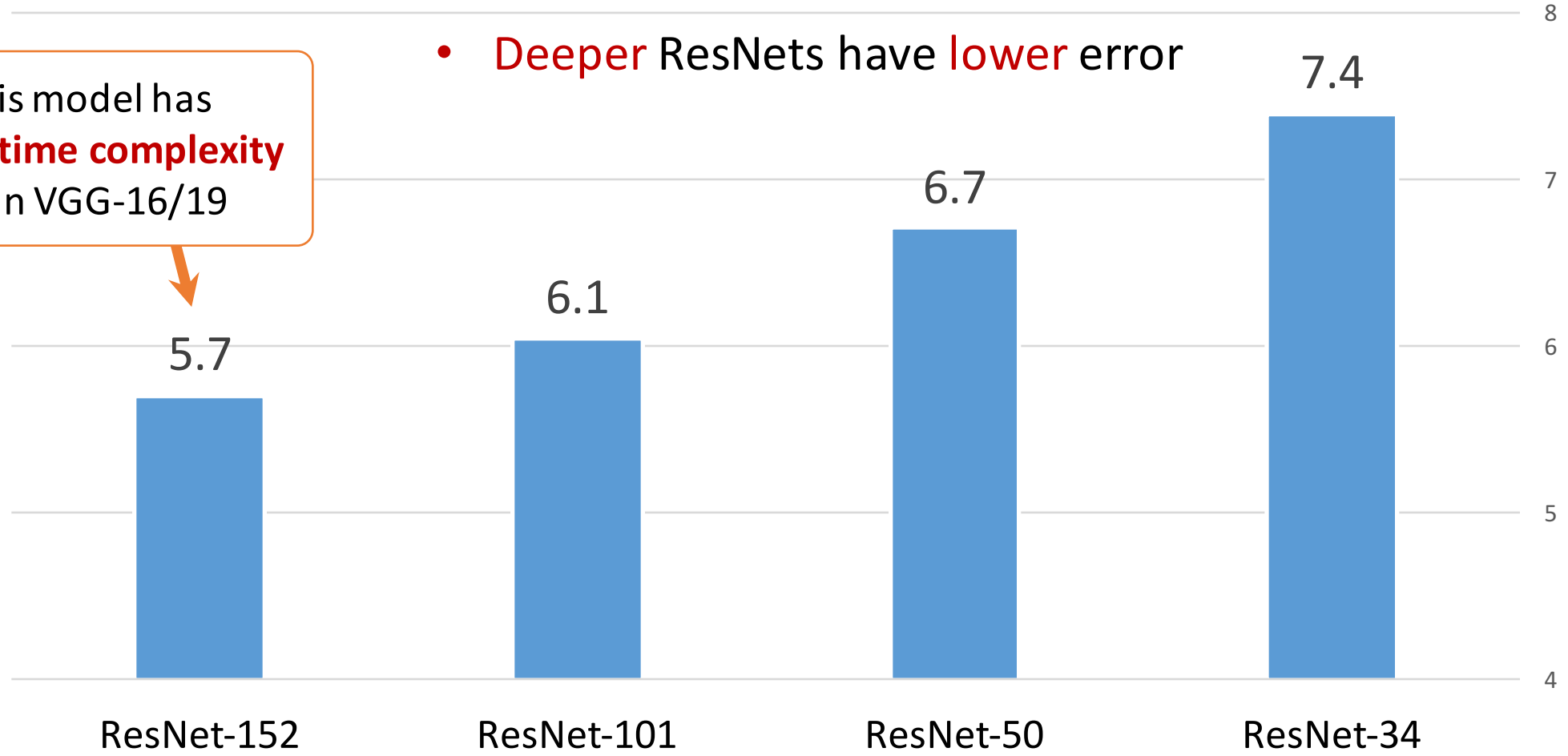Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". CVPR 2016.

# ImageNet experiments



- Deeper ResNets have lower error

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

ResNet-152: 5.7
ResNet-101: 6.1
ResNet-50: 6.7
ResNet-34: 7.4

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

Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". CVPR 2016.

# Beyond classification

**A treasure from ImageNet is on <span style="color:red">learning features</span>.**

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

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

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

**absolute 8.5% better!**

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

Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". CVPR 2016.

# Object Detection (brief)

- Simply "Faster R-CNN + ResNet"

| Faster R-CNN baseline | mAP@.5 | mAP@.5:.95 |
|---|---|---|
| VGG-16 | 41.5 | 21.5 |
| ResNet-101 | **48.4** | **27.2** |

COCO detection results
(ResNet has 28% relative gain)



classifier

RoI pooling

proposals

Region Proposal Net

feature map

CNN

image

Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". CVPR 2016.
Shaoqing Ren, Kaiming He, Ross Girshick, & Jian Sun. "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks". NIPS 2015.

Our results on MS COCO

Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". CVPR 2016.
Shaoqing Ren, Kaiming He, Ross Girshick, & Jian Sun. "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks". NIPS 2015.

this video is available online: https://youtu.be/WZmSMkK9VuA

Results on real video. Model trained on MS COCO w/ 80 categories.
(frame-by-frame; no temporal processing)

Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". arXiv 2015.
Shaoqing Ren, Kaiming He, Ross Girshick, & Jian Sun. "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks". NIPS 2015.

# More Visual Recognition Tasks

## ResNets lead on these benchmarks (incomplete list):

- **ImageNet** classification, detection, localization
- **MS COCO** detection, segmentation

- **PASCAL VOC** detection, segmentation
- **VQA** challenge 2016

- Human pose estimation [Newell et al 2016]
- Depth estimation [Laina et al 2016]
- Segment proposal [Pinheiro et al 2016]
- …

| | mean | aero plane | bicycle | bird | boat | bottle | bus | car | |
|---|---|---|---|---|---|---|---|---|---|
| DeepLabv2-CRF [?] | 79.7 | 92.6 | 60.4 | 91.6 | 63.4 | 76.3 | 95.0 | 88.4 | |
| CASIA_SegResNet_CRF_COCO [?] | 79.3 | 93.8 | | | | | | | |
| Adelaide_VeryDeep_FCN_VOC [?] | 79.1 | 91.9 | 48.1 | 93.4 | 69.3 | 75.5 | 94.2 | 87.5 | |
| LRR_4x_COCO [?] | 78.7 | 93.2 | 44.2 | 89.4 | 63.4 | 74.3 | 93.5 | 87.0 | |
| CASIA_IVA_OASeg [?] | 78.3 | 93.8 | 41.9 | 89.4 | 67.5 | 71.5 | 94.6 | 85.3 | |
| Oxford_TVG_HO_CRF [?] | 77.9 | 92.5 | 59.1 | 90.3 | 70.6 | 74.4 | 92.4 | 84.1 | |
| Adelaide_Context_CNN_CRF_COCO [?] | 77.8 | 92.9 | 39.6 | 84.0 | 67.9 | 75.3 | 92.7 | 83.8 | |

ResNet-101

PASCAL **segmentation** leaderboard

| | mean | aero plane | bicycle | bird | boat | bottle | bus | car | cat |
|---|---|---|---|---|---|---|---|---|---|
| Faster RCNN, ResNet (VOC+COCO) [?] | 83.8 | 92.1 | 88.4 | 84.8 | 75.9 | 71.4 | 86.3 | 87.8 | 84.2 |
| R-FCN, ResNet (VOC+COCO) [?] | 82.0 | 89.5 | 88.3 | 83.1 | | | | 86.3 | |
| OHEM+FRCN, VGG16, VOC+COCO [?] | 80.1 | 90.1 | 87.4 | 79.9 | 65.8 | 66.3 | 86.1 | 85.0 | 92.5 |
| SSD500 VGG16 VOC + COCO [?] | 78.7 | 89.1 | 85.7 | 78.9 | 63.3 | 57.0 | 85.3 | 84.1 | 92.3 |
| HFM_VGG16 [?] | 77.5 | 88.8 | 85.1 | 76.8 | 64.8 | 61.4 | 85.0 | 84.1 | 90.0 |
| IFRN_07+12 [?] | 76.6 | 87.8 | 83.9 | 79.0 | 64.5 | 58.9 | 82.2 | 82.0 | 91.4 |
| ION [?] | 76.4 | 87.5 | 84.7 | 76.8 | 63.8 | 58.3 | 82.6 | 79.0 | 90.9 |

ResNet-101

PASCAL **detection** leaderboard

# Potential Applications

ResNets have
shown outstanding or
promising results on:

Visual Recognition

Image Generation
(Pixel RNN, Neural Art, etc.)

Natural Language Processing
(Very deep CNN)

Speech Recognition
(preliminary results)

Advertising, user prediction
(preliminary results)

Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". CVPR 2016.

# Conclusions

- Deep Residual Networks:
  - Easy to train
  - Simply gain accuracy from depth
  - Well transferrable

- **Follow-up** [He et al. arXiv 2016]
  - 200 layers on ImageNet, 1000 layers on CIFAR

Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Identity Mappings in Deep Residual Networks". arXiv 2016.
Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". CVPR 2016.

# Resources

- Models and Code
  - Our ImageNet models in Caffe: https://github.com/KaimingHe/deep-residual-networks

- Many available implementations:
  (list in https://github.com/KaimingHe/deep-residual-networks)
  - Facebook AI Research's Torch ResNet:
    https://github.com/facebook/fb.resnet.torch
    - Torch, CIFAR-10, with ResNet-20 to ResNet-110, training code, and curves: code
    - Lasagne, CIFAR-10, with ResNet-32 and ResNet-56 and training code: code
    - Neon, CIFAR-10, with pre-trained ResNet-32 to ResNet-110 models, training code, and curves: code
    - Torch, MNIST, 100 layers: blog, code
    - A winning entry in Kaggle's right whale recognition challenge: blog, code
    - Neon, Place2 (mini), 40 layers: blog, code
    - .......

Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". CVPR 2016.