# Deep Residual Learning

## MSRA @ ILSVRC & COCO 2015 competitions

Kaiming He

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

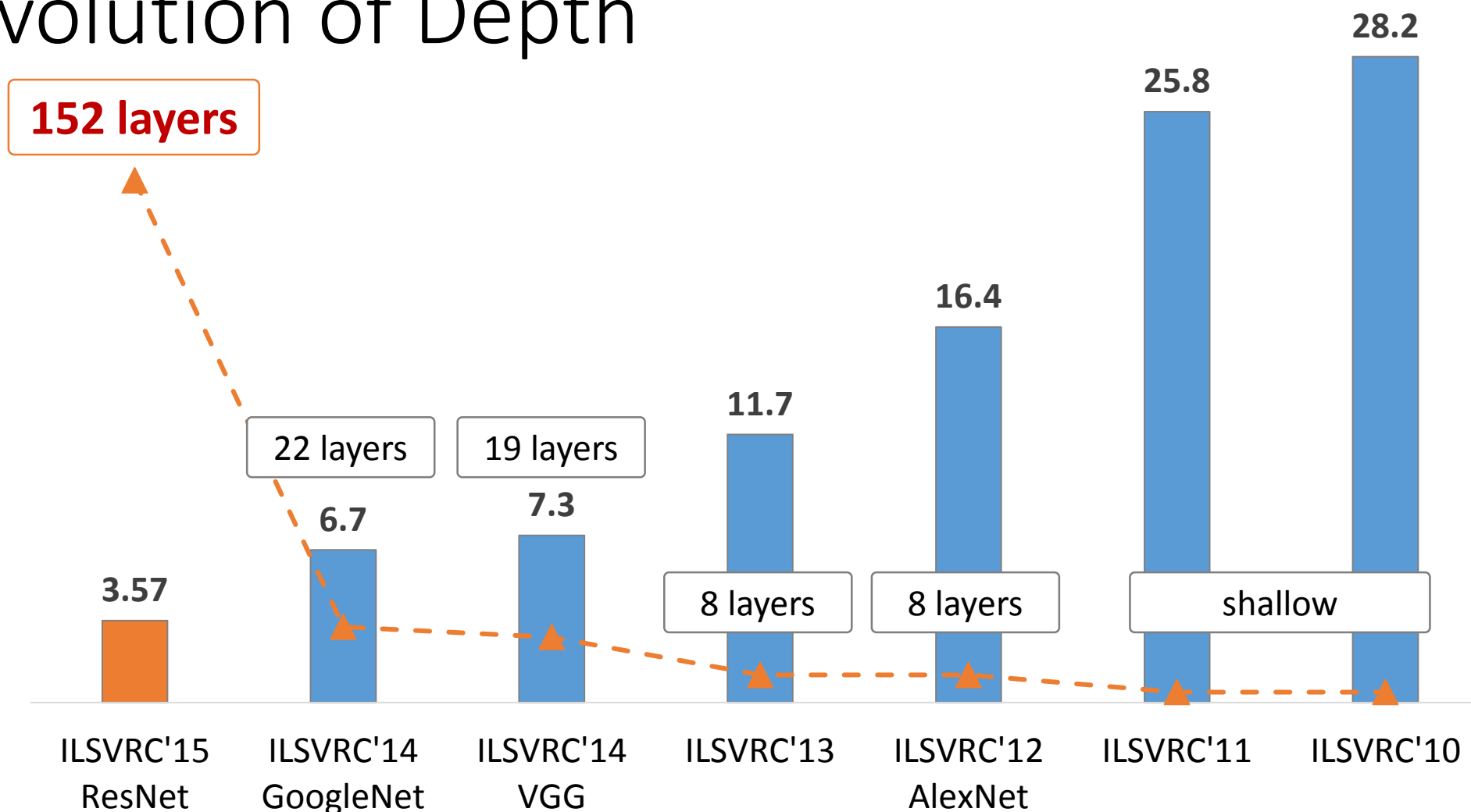Microsoft Research Asia (MSRA)

# MSRA @ ILSVRC & COCO 2015 Competitions

- **1st places in all five main tracks**
  - ImageNet Classification: "*Ultra-deep*" (quote Yann) 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
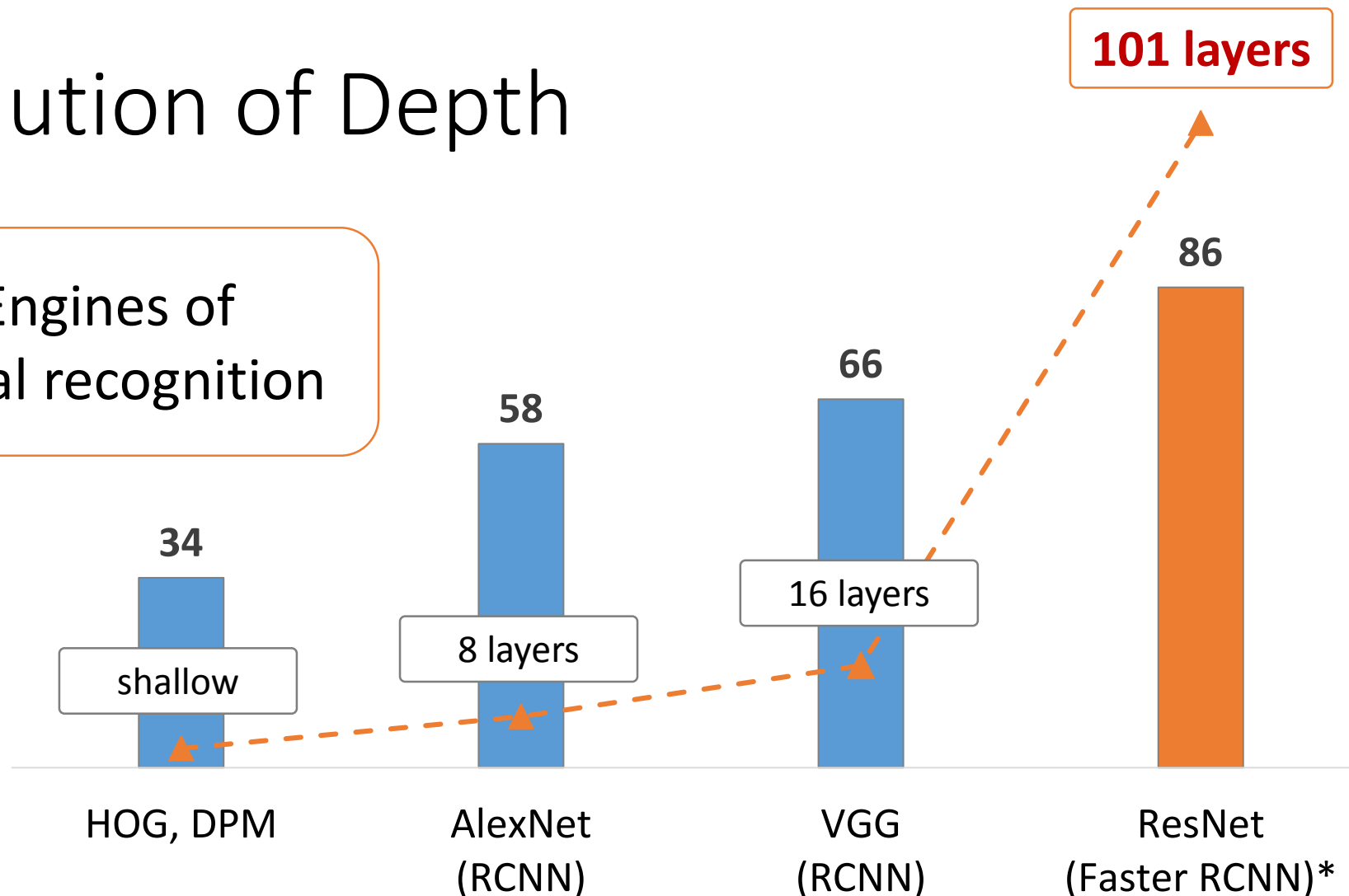
*improvements are relative numbers

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

# Revolution of Depth

**152 layers**

28.2

25.8

16.4

11.7

7.3

6.7

3.57

22 layers

19 layers

8 layers

8 layers

shallow

ILSVRC'15 ResNet
ILSVRC'14 GoogleNet
ILSVRC'14 VGG
ILSVRC'13
ILSVRC'12 AlexNet
ILSVRC'11
ILSVRC'10

ImageNet Classification top-5 error (%)
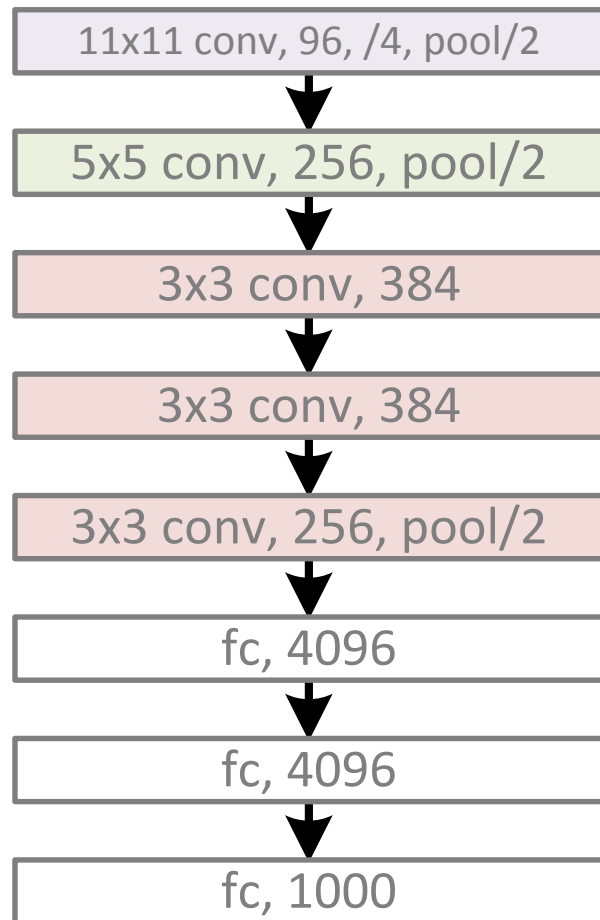
Revolution of Depth

Engines of visual recognition

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". arXiv 2015.

# 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". arXiv 2015.

# Revolution of Depth

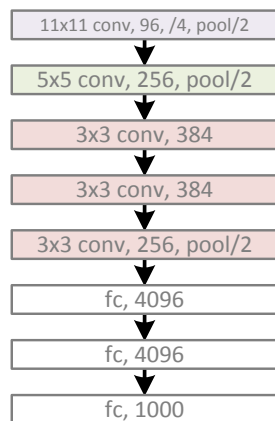**AlexNet, 8 layers**
**(ILSVRC 2012)**

| |
|---|
| 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". arXiv 2015.
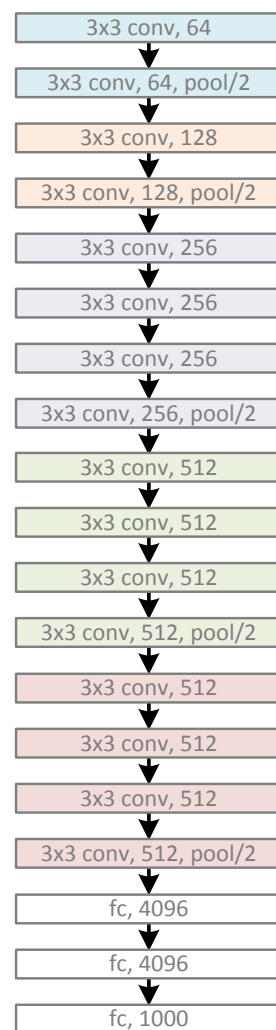
# Revolution of Depth

AlexNet, 8 layers
(ILSVRC 2012)

VGG, 19 layers
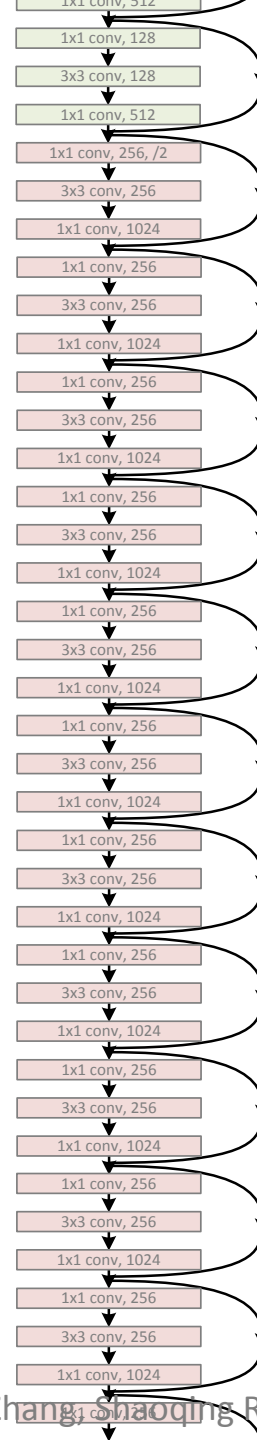(ILSVRC 2014)

ResNet, 152 layers
(ILSVRC 2015)

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

# Revolution of Depth

ResNet, 152 layers

| 7x7 conv, 64, /2, pool/2 |
| 1x1 conv, 64 |
| 3x3 conv, 64 |
| 1x1 conv, 256 |
| 1x1 conv, 64 |
| 3x3 conv, 64 |
| 1x1 conv, 256 |
| 1x1 conv, 64 |
| 3x3 conv, 64 |
| 1x1 conv, 256 |
| 1x1 conv, 128, /2 |
| 3x3 conv, 128 |
| 1x1 conv, 512 |
| 1x1 conv, 128 |
| 3x3 conv, 128 |
| 1x1 conv, 512 |
| 1x1 conv, 128 |
| 3x3 conv, 128 |
| 1x1 conv, 512 |
| 1x1 conv, 128 |
| 3x3 conv, 128 |
| 1x1 conv, 512 |
| 1x1 conv, 128 |
| 3x3 conv, 128 |
| 1x1 conv, 512 |
| 1x1 conv, 128 |
| 3x3 conv, 128 |
| 1x1 conv, 512 |
| 1x1 conv, 128 |
| 3x3 conv, 128 |
| 1x1 conv, 512 |
| 1x1 conv, 256, /2 |
| 3x3 conv, 256 |

(there was an animation here)

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

# Revolution of Depth

ResNet, 152 layers

1x1 conv, 512

1x1 conv, 128
3x3 conv, 128
1x1 conv, 512

1x1 conv, 256, /2
3x3 conv, 256
1x1 conv, 1024

1x1 conv, 256
3x3 conv, 256
1x1 conv, 1024

1x1 conv, 256
3x3 conv, 256
1x1 conv, 1024

1x1 conv, 256
3x3 conv, 256
1x1 conv, 1024

1x1 conv, 256
3x3 conv, 256
1x1 conv, 1024

1x1 conv, 256
3x3 conv, 256
1x1 conv, 1024

1x1 conv, 256
3x3 conv, 256
1x1 conv, 1024

1x1 conv, 256
3x3 conv, 256
1x1 conv, 1024

1x1 conv, 256
3x3 conv, 256
1x1 conv, 1024

1x1 conv, 256
3x3 conv, 256
1x1 conv, 1024

(there was an animation here)

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

# Revolution of Depth

ResNet, 152 layers

(there was an animation here)

| 1x1 conv, 256 |
| --- |
| 3x3 conv, 256 |
| 1x1 conv, 1024 |
| 1x1 conv, 256 |
| 3x3 conv, 256 |
| 1x1 conv, 1024 |
| 1x1 conv, 256 |
| 3x3 conv, 256 |
| 1x1 conv, 1024 |
| 1x1 conv, 256 |
| 3x3 conv, 256 |
| 1x1 conv, 1024 |
| 1x1 conv, 256 |
| 3x3 conv, 256 |
| 1x1 conv, 1024 |
| 1x1 conv, 256 |
| 3x3 conv, 256 |
| 1x1 conv, 1024 |
| 1x1 conv, 256 |
| 3x3 conv, 256 |
| 1x1 conv, 1024 |
| 1x1 conv, 256 |
| 3x3 conv, 256 |
| 1x1 conv, 1024 |
| 1x1 conv, 256 |
| 3x3 conv, 256 |
| 1x1 conv, 1024 |
| 1x1 conv, 256 |
| 3x3 conv, 256 |
| 1x1 conv, 1024 |
| 1x1 conv, 256 |
| 3x3 conv, 256 |

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

# Revolution of Depth

ResNet, 152 layers

3x3 conv, 256

1x1 conv, 1024

1x1 conv, 256

3x3 conv, 256

1x1 conv, 1024

1x1 conv, 256

3x3 conv, 256

1x1 conv, 1024

1x1 conv, 256

3x3 conv, 256

1x1 conv, 1024

1x1 conv, 256

3x3 conv, 256

1x1 conv, 1024

1x1 conv, 256

3x3 conv, 256

1x1 conv, 1024

1x1 conv, 256

3x3 conv, 256

1x1 conv, 1024

1x1 conv, 256

3x3 conv, 256

1x1 conv, 1024

1x1 conv, 512, /2

3x3 conv, 512

1x1 conv, 2048

1x1 conv, 512

3x3 conv, 512

1x1 conv, 2048

1x1 conv, 512

3x3 conv, 512

1x1 conv, 2048

ave pool, fc 1000

(there was an animation here)

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

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

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

# 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

# Simply stacking layers?



CIFAR-10

ImageNet-1000

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

34-layer
18-layer

solid: test/val
dashed: train

- "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)

"extra" layers

- 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". arXiv 2015.

# Deep Residual Learning

- Plaint net

$x$

weight layer

any two
stacked layers

relu

weight layer

relu

$H(x)$

# Deep Residual Learning

- **Residual** net

$x$

weight layer

$F(x)$    relu

weight layer

identity

$x$

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

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". arXiv 2015.

# Deep Residual Learning

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



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

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

# Related Works – Residual Representations

- VLAD & Fisher Vector [Jegou et al 2010], [Perronnin et al 2007]
  - Encoding residual vectors; powerful shallower representations.

- Product Quantization (IVF-ADC) [Jegou et al 2011]
  - Quantizing residual vectors; efficient nearest-neighbor search.

- MultiGrid & Hierarchical Precondition [Briggs, et al 2000], [Szeliski 1990, 2006]
  - Solving residual sub-problems; efficient PDE solvers.

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

# 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!

- Other remarks:
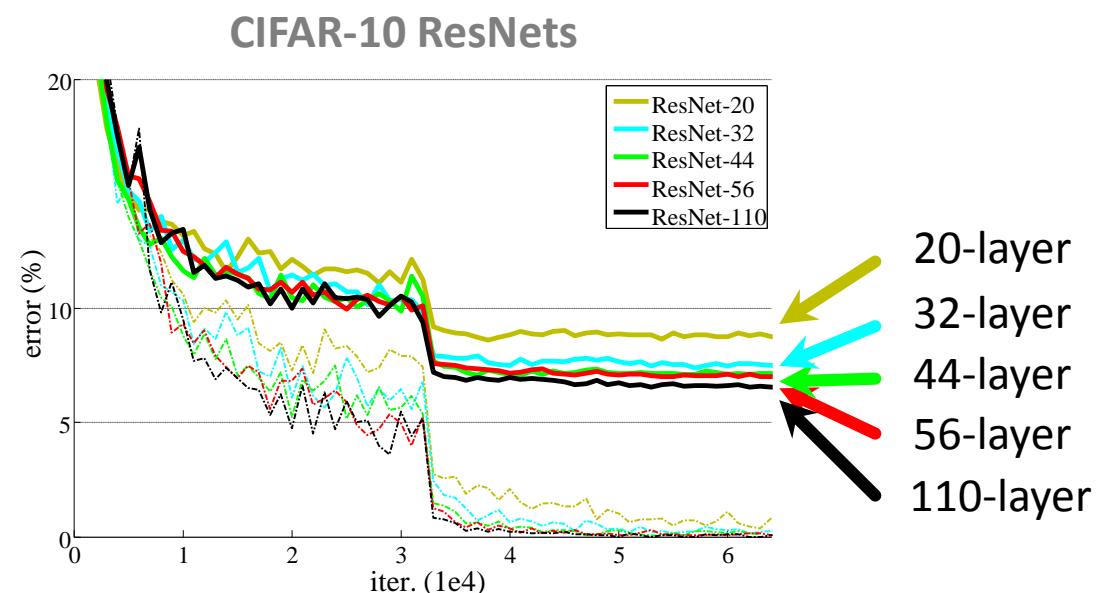  - no max pooling (almost)
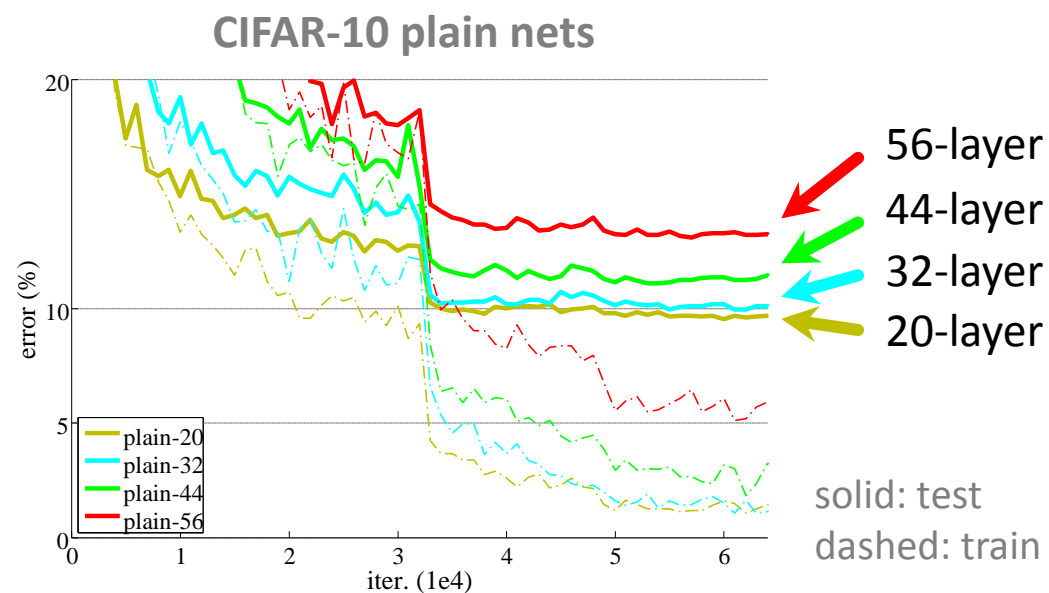  - no hidden fc
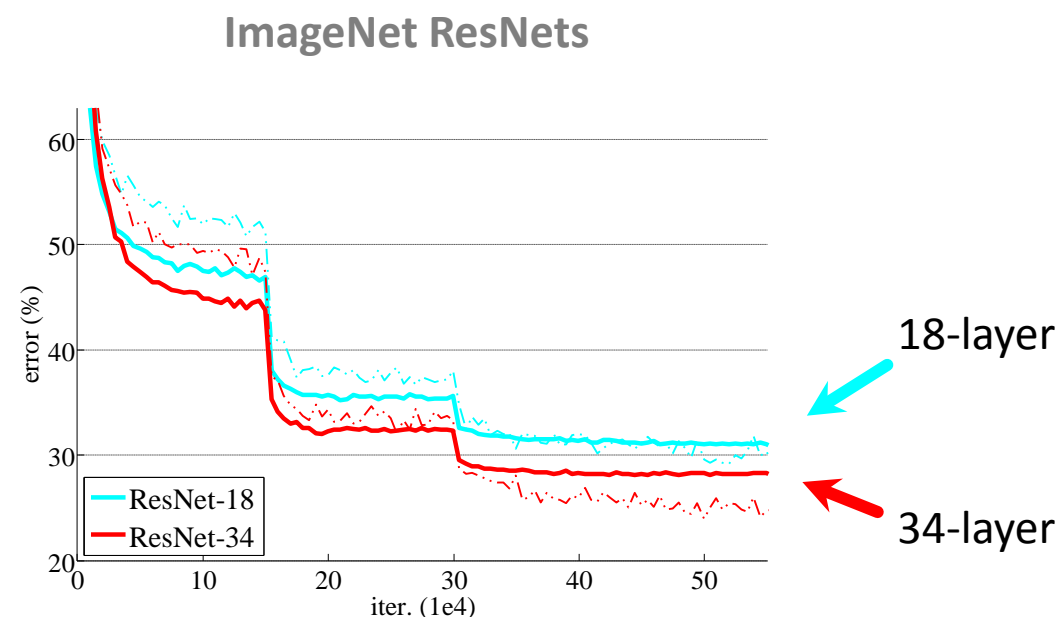  - no dropout

# Training
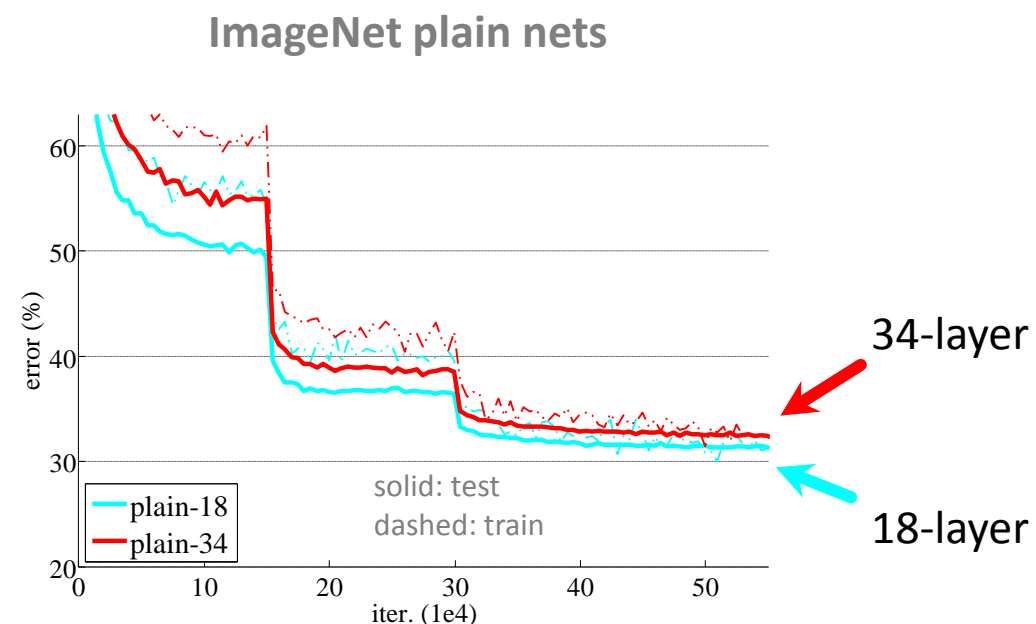
- All plain/residual nets are trained <span style="color:red">from scratch</span>

- All plain/residual nets use Batch Normalization

- Standard hyper-parameters & augmentation

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

# CIFAR-10 experiments

**CIFAR-10 plain nets**



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

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

solid: test
dashed: train

**CIFAR-10 ResNets**



ResNet-20
ResNet-32
ResNet-44
ResNet-56
ResNet-110

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

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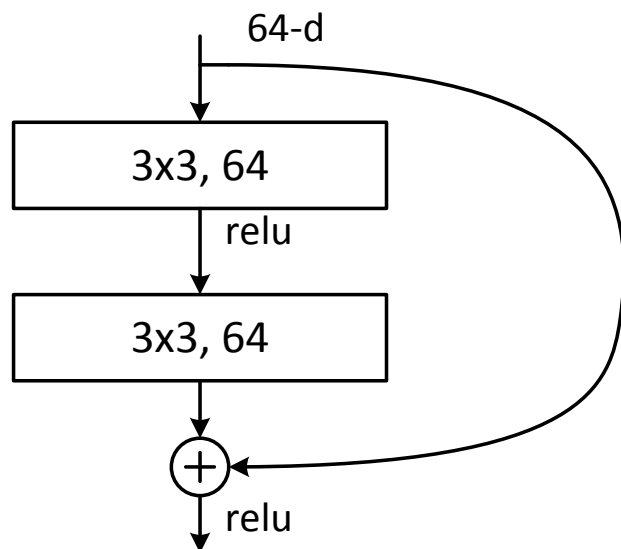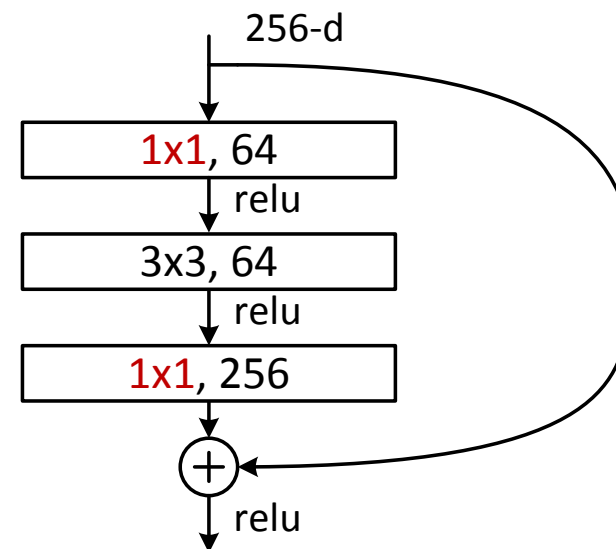- Deeper ResNets have **lower training error**, and also lower test error

# 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". arXiv 2015.

# ImageNet experiments

- A practical design of going deeper



all-3x3 ⟷ similar complexity ⟷ **bottleneck**
(for ResNet-50/101/152)

# ImageNet experiments

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

- Deeper ResNets have lower error

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". arXiv 2015.

# ImageNet experiments

**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 (%)

# Just 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 | MSRA | 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!**
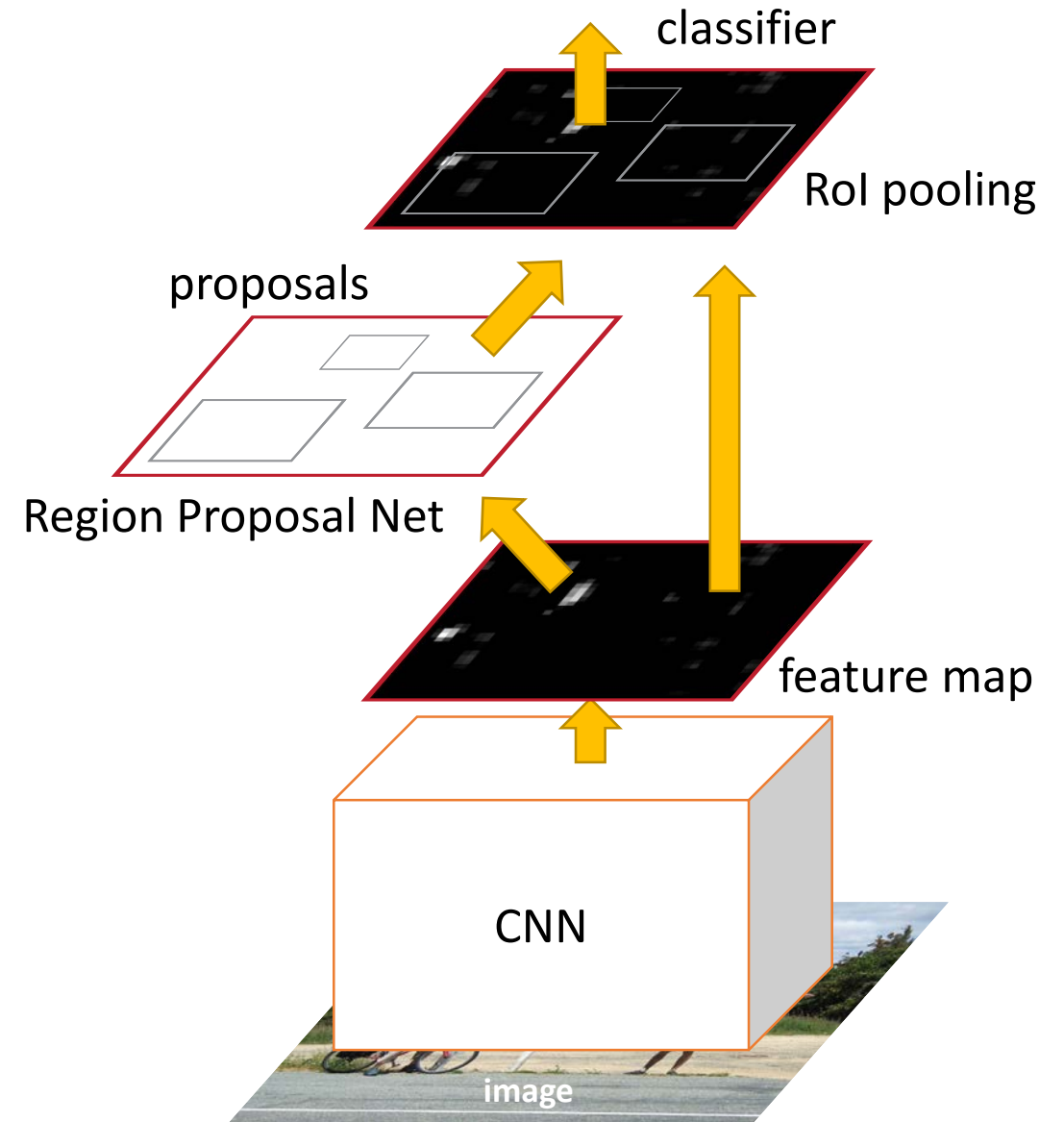
- Our results are all based on ResNet-101
- Our features are well transferrable

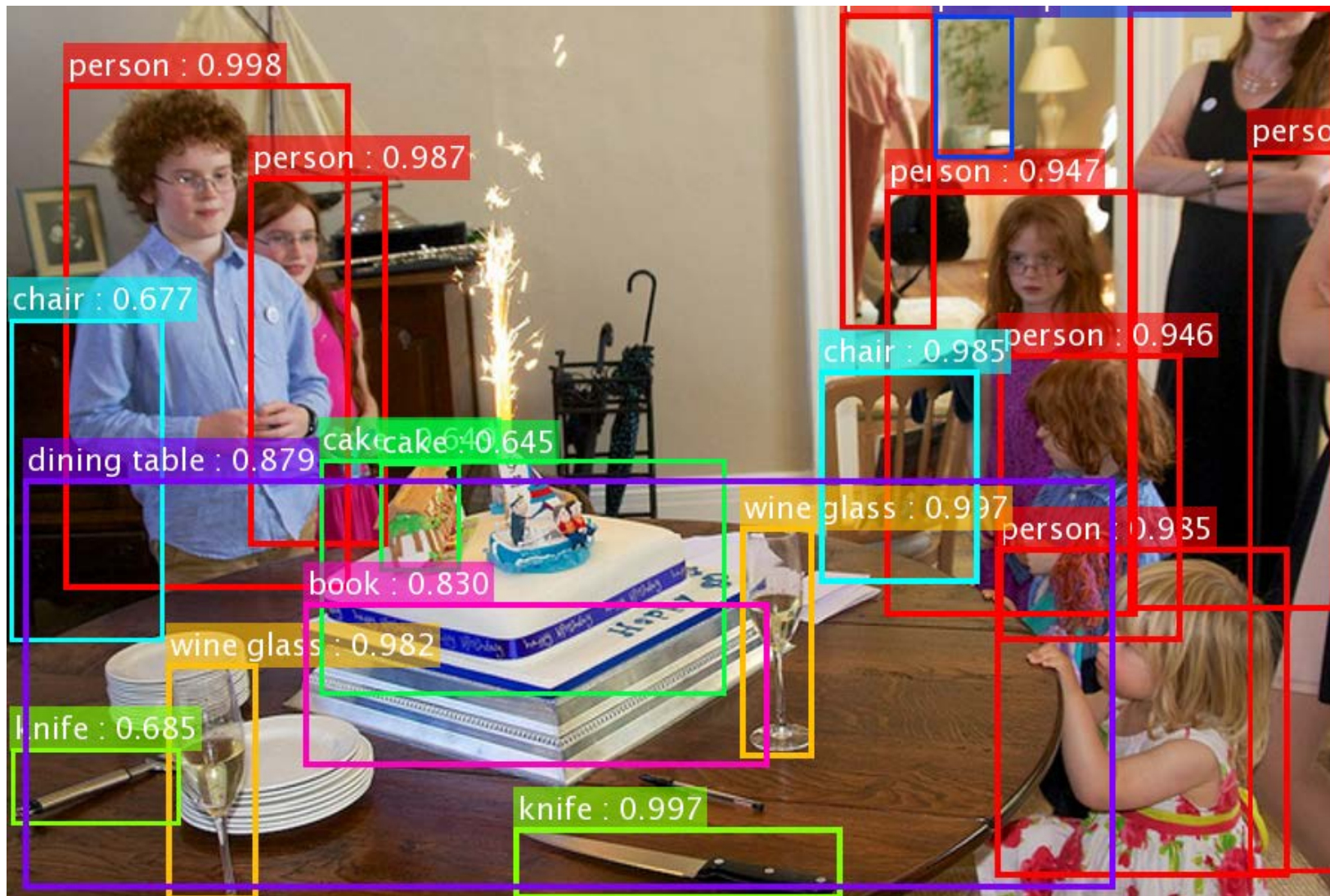# Object Detection (brief)

- 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". arXiv 2015.
Shaoqing Ren, Kaiming He, Ross Girshick, & Jian Sun. "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks". NIPS 2015.

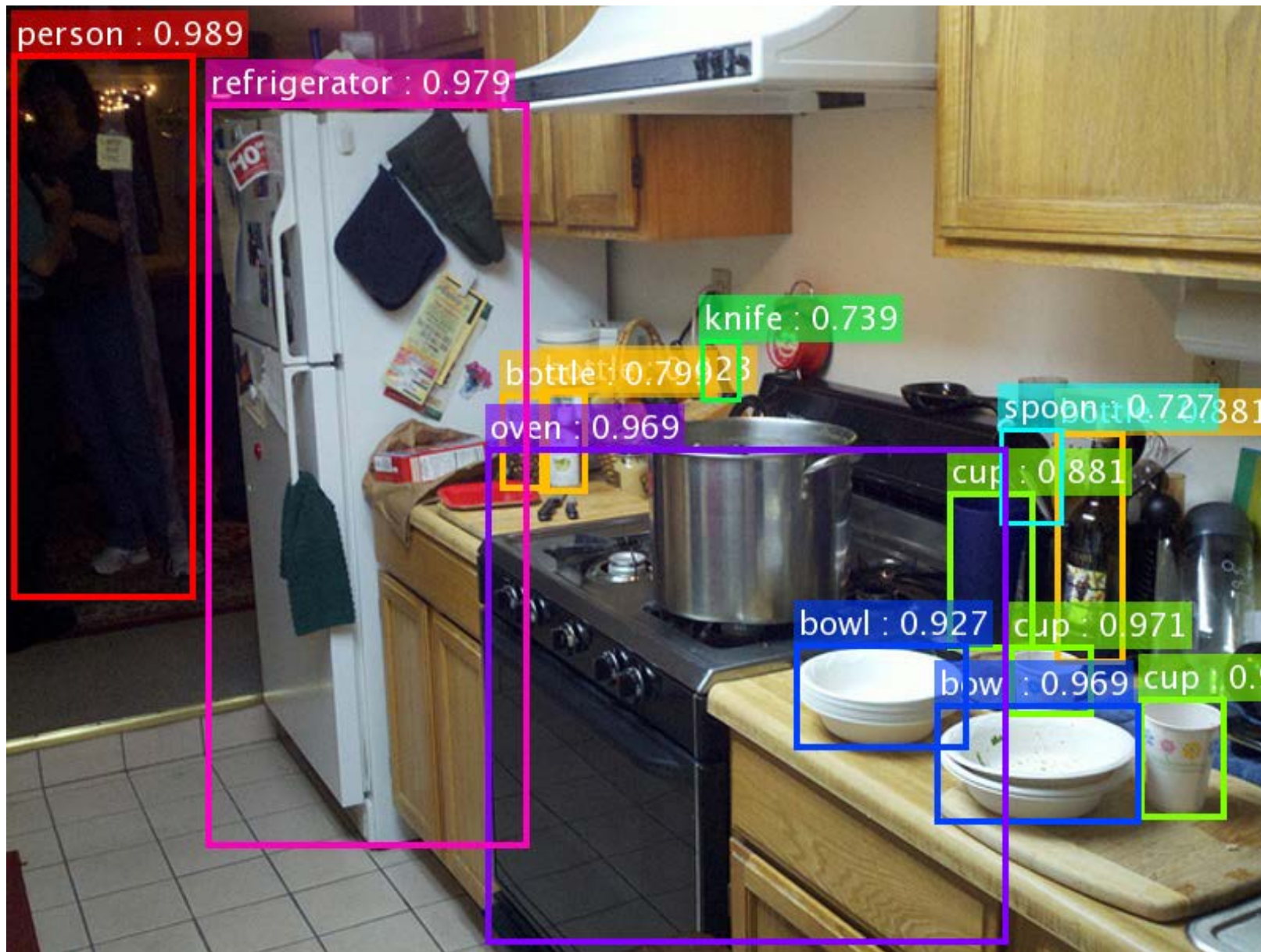# Object Detection (brief)

- RPN learns proposals by extremely deep nets
  - We use only 300 proposals (no SS/EB/MCG!)

- Add what is just missing in Faster R-CNN…
  - Iterative localization
  - Context modeling
  - Multi-scale testing

- All are based on CNN features; all are end-to-end (train and/or inference)

- All benefit more from deeper features – cumulative gains!

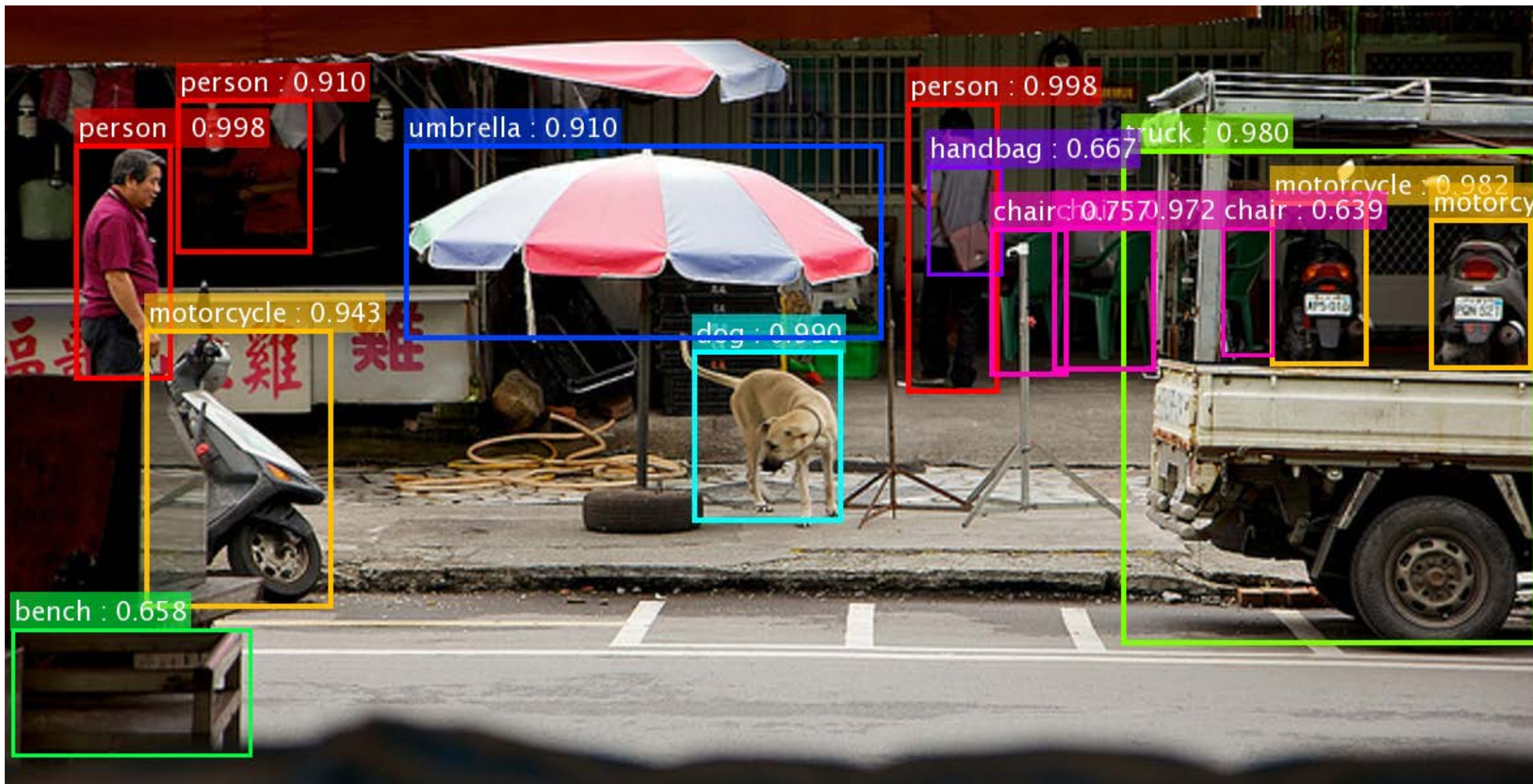Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". arXiv 2015.
Shaoqing Ren, Kaiming He, Ross Girshick, & Jian Sun. "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks". NIPS 2015.

Our results on COCO – too many objects, let's check carefully!

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.

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.

*the original image is from the COCO dataset

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.
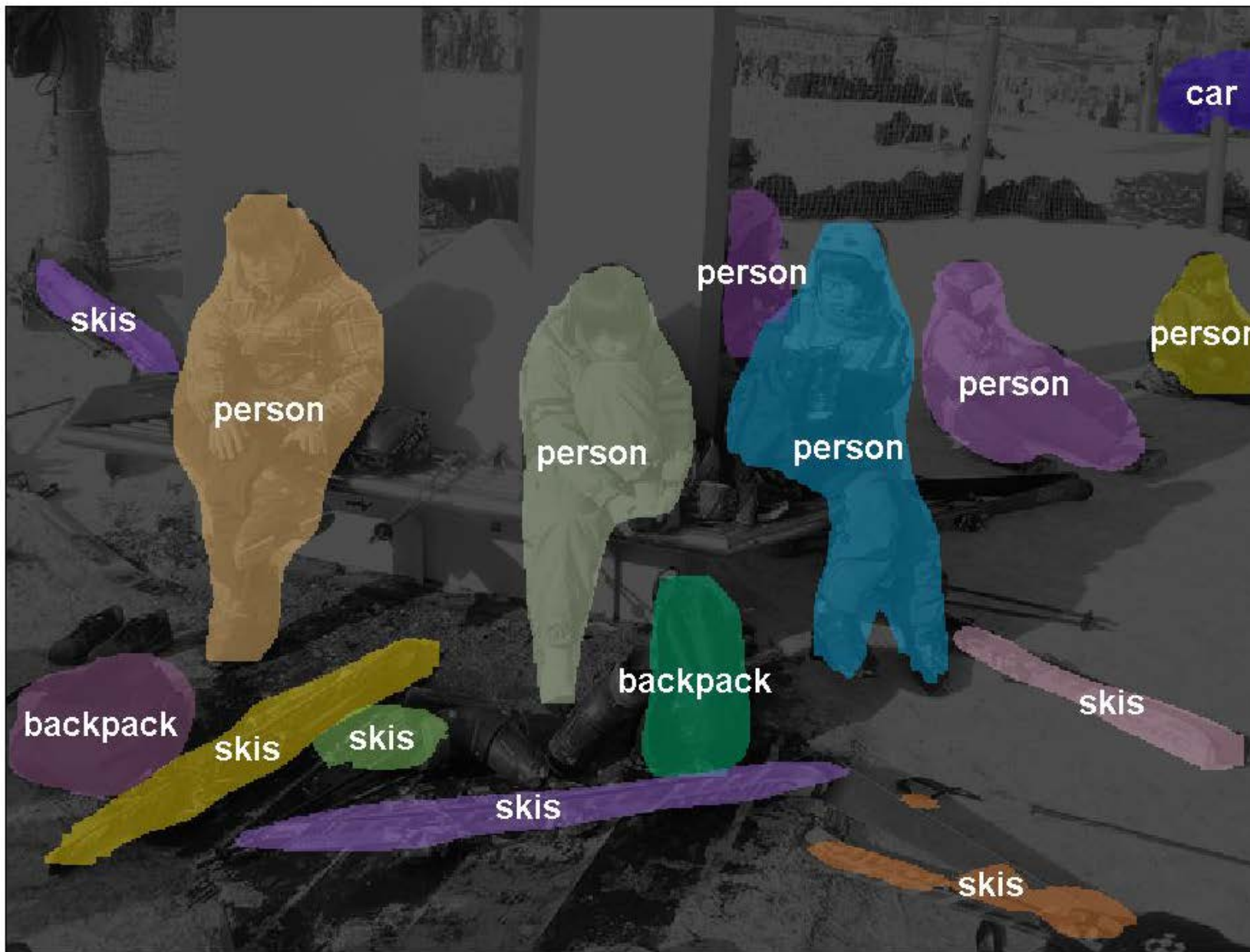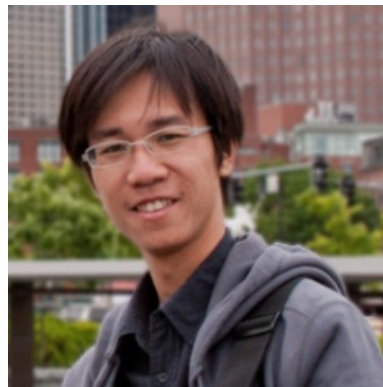
# Instance Segmentation (brief)



- Solely CNN-based ("features matter")
- Differentiable RoI warping layer (w.r.t box coord.)
- Multi-task cascades, exact end-to-end training

box instances (RoIs)

conv feature map

RoI warping, pooling

*for each RoI*

mask instances

masking

*for each RoI*

categorized instances

CONVs

FCs

FCs

person person person horse

Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". arXiv 2015.
Jifeng Dai, Kaiming He, & Jian Sun. "Instance-aware Semantic Segmentation via Multi-task Network Cascades". arXiv 2015.
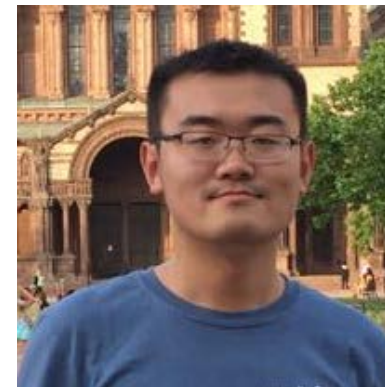
input

Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". arXiv 2015.
Jifeng Dai, Kaiming He, & Jian Sun. "Instance-aware Semantic Segmentation via Multi-task Network Cascades". arXiv 2015.

# Conclusions

- Deeper is still better

- "*Features matter*"!

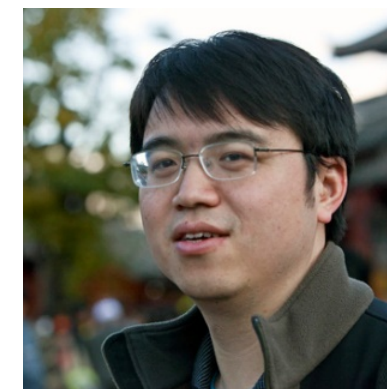- Faster R-CNN is just amazing



Kaiming He    Xiangyu Zhang    Shaoqing Ren

Jifeng Dai    Jian Sun

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.
Jifeng Dai, Kaiming He, & Jian Sun. "Instance-aware Semantic Segmentation via Multi-task Network Cascades". arXiv 2015.