# Learning Deep Representations for Visual Recognition

## CVPR 2018 Tutorial

Kaiming He

Facebook AI Research (FAIR)
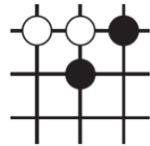
# Deep Learning <u>is</u> Representation Learning

Representation Learning: worth a conference name ☺ (ICLR)

Represent (raw) data for machines to perform tasks:

- Vision: pixels, …

- Language: letters, …
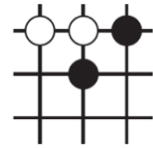
- Speech: waves, …

- Games: status, …

# Representation Learning: AlphaGo

$3^{361}$ states?

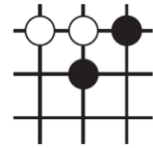# Representation Learning: AlphaGo

$3^{361}$ states?



*Bad representations*

$256^{3*640*480}$?

# Representation Learning: AlphaGo
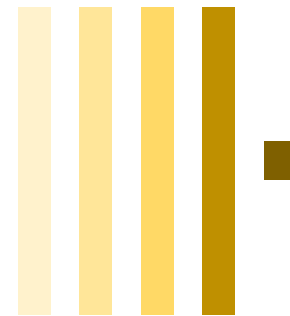
$3^{361}$ states?



$256^{3*640*480}$?



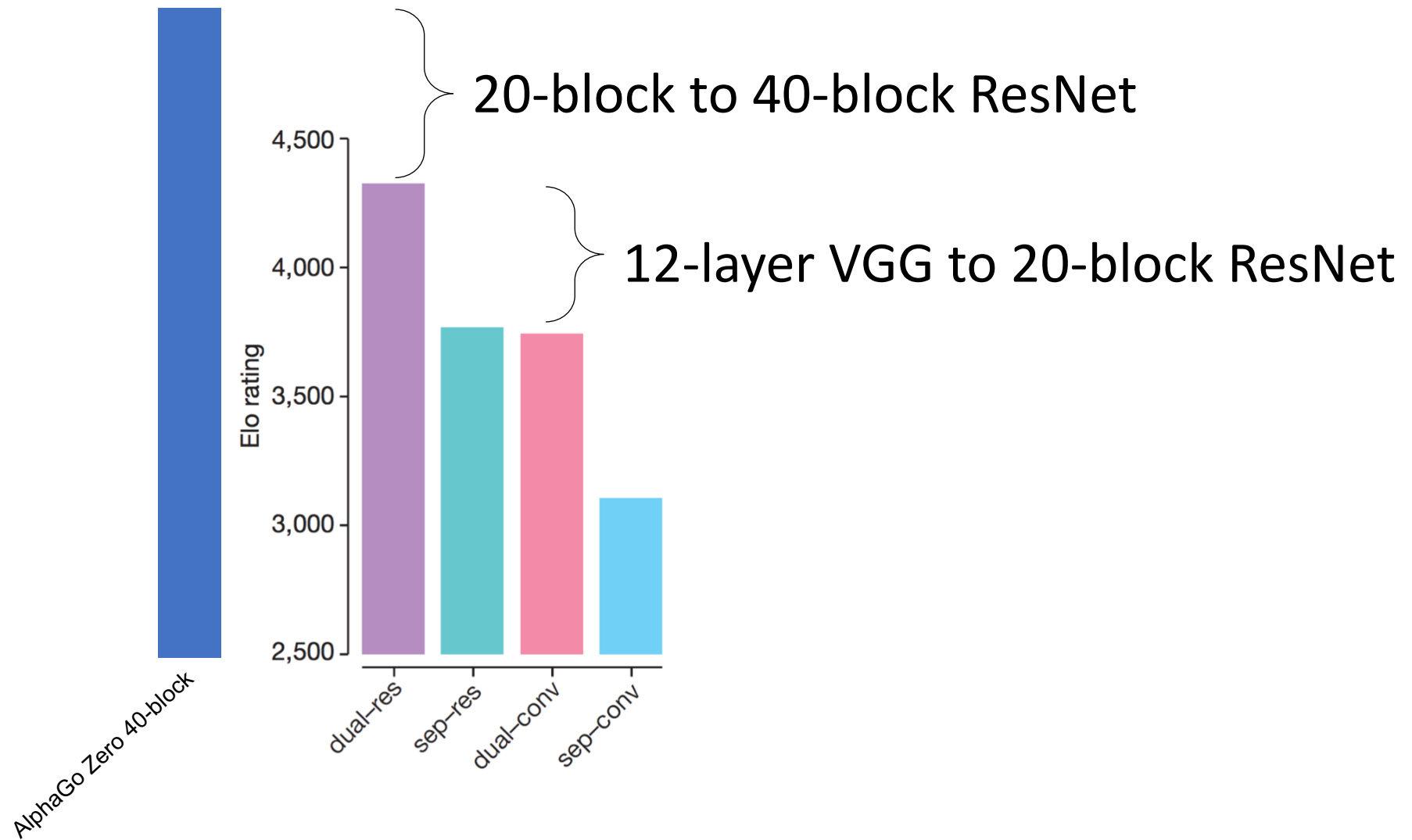*Bad representations*

models
(now, neural nets)

*Good representations*

# Representation Learning: AlphaGo



20-block to 40-block ResNet

12-layer VGG to 20-block ResNet

"Mastering the game of Go without human knowledge", Silver et al. Nature 2017

# How was an image represented?



shallower

pixels → classifier → "bus"?

edges → classifier → "bus"?

**But what's next?**

SIFT/HOG

edges → histogram → classifier → "bus"?

deeper

edges → histogram → K-means sparse code FV/VLAD → classifier → "bus"?

[Lowe 1999, 2004], [Sivic & Zisserman 2003], [Dalal & Triggs 2005], [Grauman & Darrell 2005]
[Lazebnik et al 2006], [Perronnin & Dance 2007], [Yang et al 2009], [Jégou et al 2010], ......

# Learning to represent

Specialized components, domain knowledge required

 → edges → histogram → K-means sparse code FV/VLAD → classifier → "bus"?

Generic components, less domain knowledge

 → ☐ → ☐ → ☐ → ☐ → "bus"?

Repeat **elementary** layers: going deeper

 → ☐ → ☐ → ☐ → ☐ → ☐ → ☐ → ☐ → "bus"?

- End-to-end by BackProp

# LeNet



INPUT 32x32 — C1: feature maps 6@28x28 — C3: f. maps 16@10x10 — S2: f. maps 6@14x14 — S4: f. maps 16@5x5 — C5: layer 120 — F6: layer 84 — OUTPUT 10

Convolutions — Subsampling — Convolutions — Subsampling — Full connection — Full connection — Gaussian connections

- Convolution:
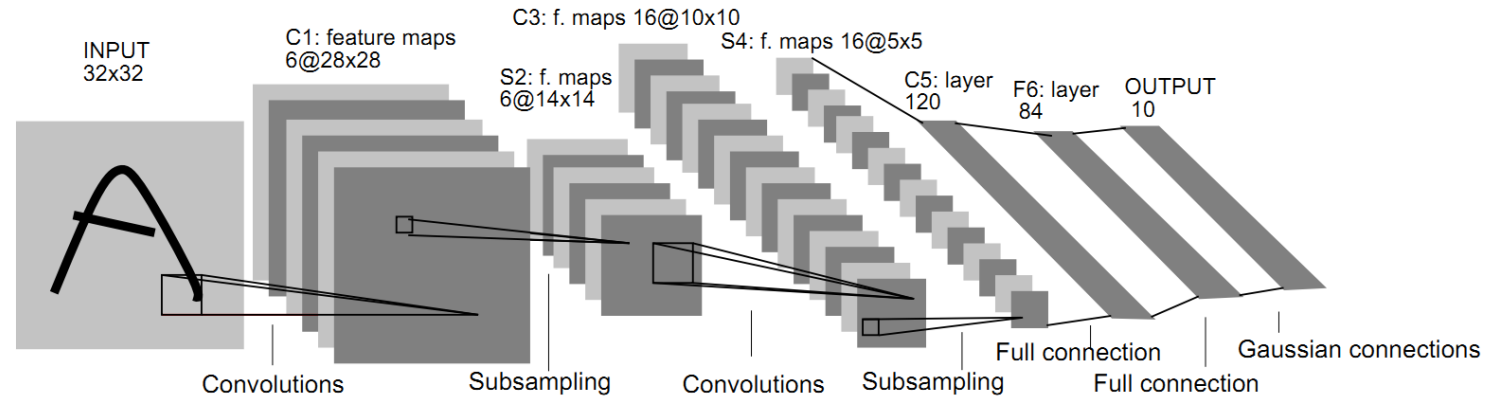  - locally-connected
  - spatially weight-sharing
    - weight-sharing is a key in DL (e.g., RNN shares weights temporally)

- Subsampling

- Fully-connected outputs

- Train by BackProp
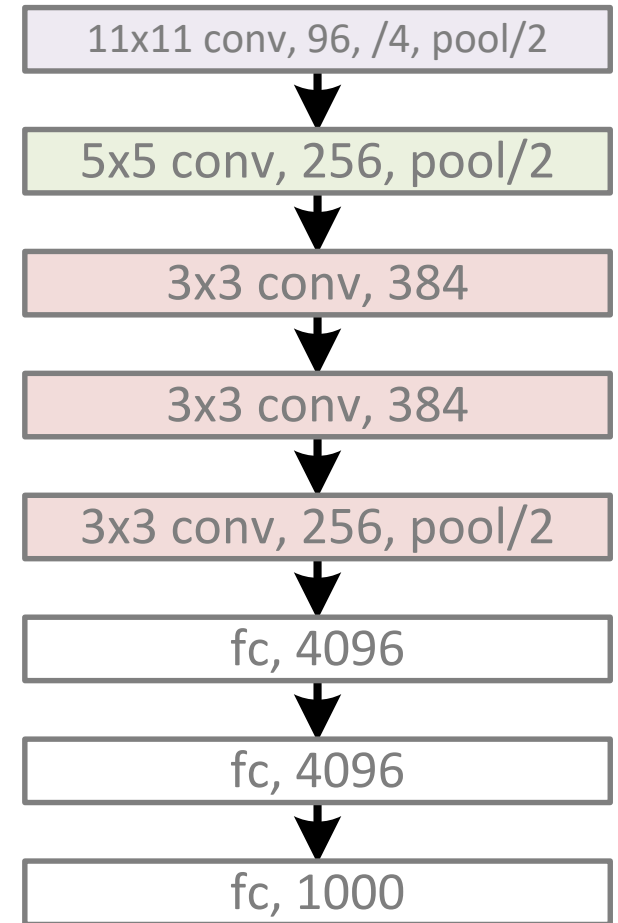
- All are still the basic components of modern ConvNets!

"Gradient-based learning applied to document recognition", LeCun et al. 1998

"Backpropagation applied to handwritten zip code recognition", LeCun et al. 1989

# AlexNet

LeNet-style backbone, plus:

- ReLU [Nair & Hinton 2010]
  - "RevoLUtion of deep learning"*
  - Accelerate training; better grad prop (vs. tanh)
- Dropout [Hinton et al 2012]
  - In-network ensembling
  - Reduce overfitting (might be instead done by BN)
- Data augmentation
  - Label-preserving transformation
  - Reduce overfitting

| 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 |

*Quote Christian Szegedy

"ImageNet Classification with Deep Convolutional Neural Networks",  Krizhevsky, Sutskever, Hinton. NIPS 2012

# VGG-16/19

Simply "Very Deep"!

- Modularized design
  - 3x3 Conv as the module
  - Stack the same module
  - Same computation for each module (1/2 spatial size => 2x filters)

- Stage-wise training
  - VGG-11 => VGG-13 => VGG-16
  - We need a better initialization…



"Very Deep Convolutional Networks for Large-Scale Image Recognition", Simonyan & Zisserman. arXiv 2014 (ICLR 2015)

# Initialization Methods

- Analytical formulations of <u>normalizing</u> forward/backward signals
- Based on strong assumptions (like Gaussian distributions)

<br>

- Xavier Init (linear): $n \cdot Var[w] = 1$
- MSRA Init (ReLU): $n \cdot Var[w] = 2$

"Efficient Backprop", LeCun et al, 1998

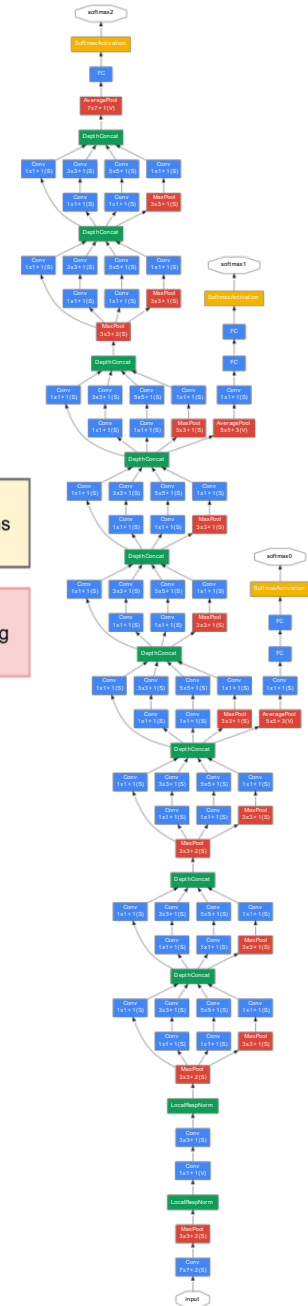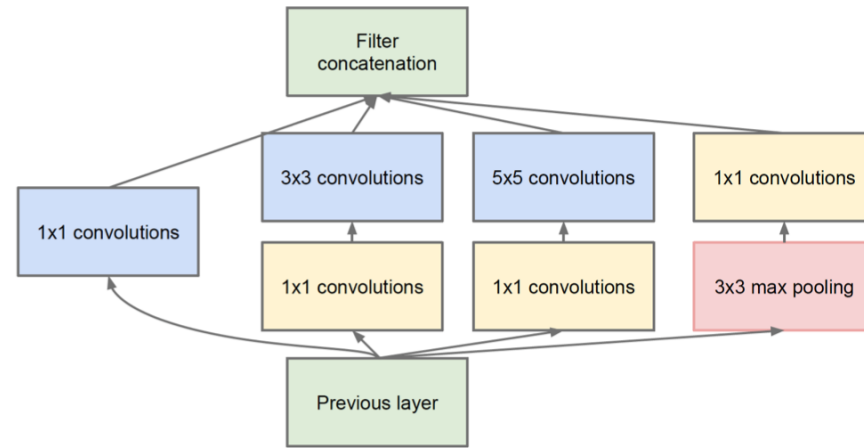"Understanding the difficulty of training deep feedforward neural networks" Glorot & Bengio, 2010

"Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification" Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun, ICCV 2015

# GoogleNet/Inception

Accurate with small footprint.
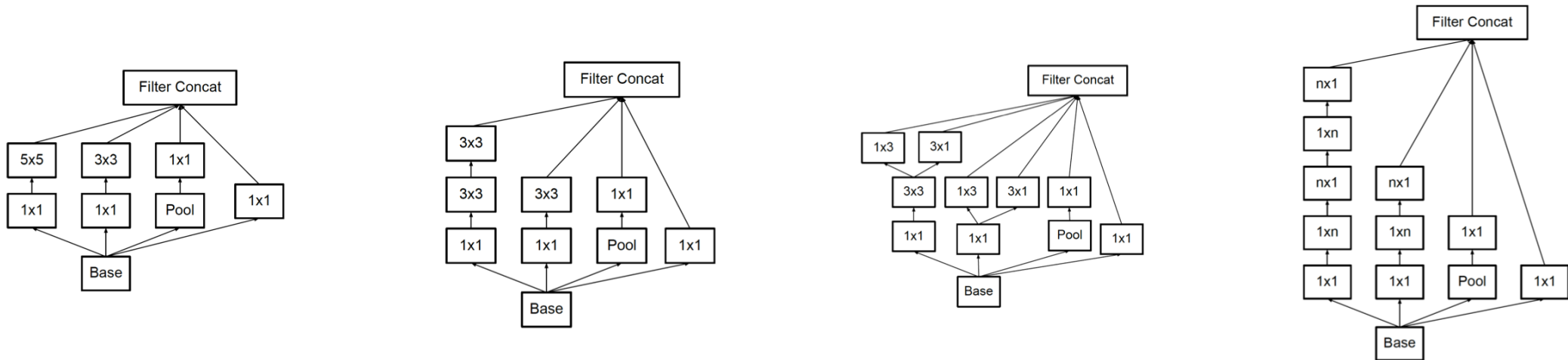
My take on GoogleNets:

- Multiple branches
  - e.g., 1x1, 3x3, 5x5, pool

- Shortcuts
  - stand-alone 1x1, merged by concat.

- Bottleneck
  - Reduce dim by 1x1 before expensive 3x3/5x5 conv



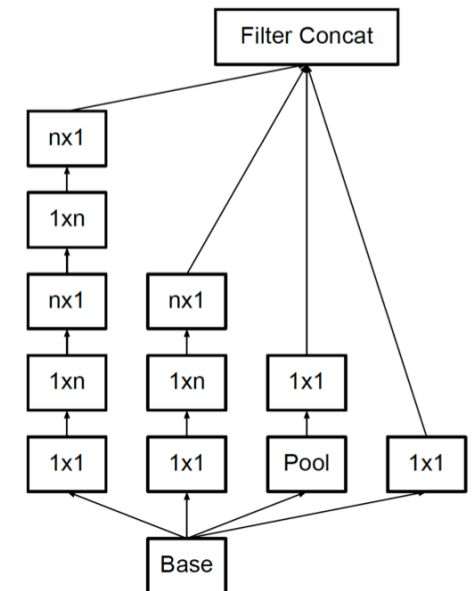Szegedy et al. "Going deeper with convolutions". arXiv 2014 (CVPR 2015)

# GoogleNet/Inception v1, v2, v3, ...

More templates, but the same 3 main properties are kept:

- Multiple branches

- Shortcuts (1x1, concate.)

- Bottleneck



Szegedy et al. "Rethinking the Inception Architecture for Computer Vision". arXiv 2015 (CVPR 2016)

# Batch Normalization (BN)

- Xavier/MSRA init are not directly applicable for multi-branch nets

- Optimizing multi-branch ConvNets largely benefits from BN
  - including all Inceptions and ResNets



Ioffe & Szegedy. "Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift". ICML 2015

# Batch Normalization (BN)

- Recap: Normalizing image input (LeCun et al 1998 "Efficient Backprop")

- Xavier/MSRA init: Analytic normalizing each layer

- BN: data-driven normalization, <span style="color:red">for each layer, for each mini-batch</span>
  - Greatly accelerate training
  - Less sensitive to initialization
  - Improve regularization

Ioffe & Szegedy. "Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift". ICML 2015

# Batch Normalization (BN)

$$x \implies \hat{x} = \frac{x - \mu}{\sigma} \implies y = \gamma\hat{x} + \beta$$

- $\mu$: mean of $x$ in mini-batch
- $\sigma$: std of $x$ in mini-batch
- $\gamma$: scale
- $\beta$: shift

- $\mu, \sigma$: functions of $x$, analogous to responses
- $\gamma, \beta$: parameters to be learned, analogous to weights

Ioffe & Szegedy. "Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift". ICML 2015

# Batch Normalization (BN)

$$x \implies \quad \hat{x} = \frac{x - \mu}{\sigma} \quad \implies \quad y = \gamma \hat{x} + \beta$$

2 modes of BN:
- Train mode:
  - $\mu, \sigma$ are functions of a batch of $x$
- Test mode:
  - $\mu, \sigma$ are pre-computed on training set

**Caution**: make sure your BN usage is correct!
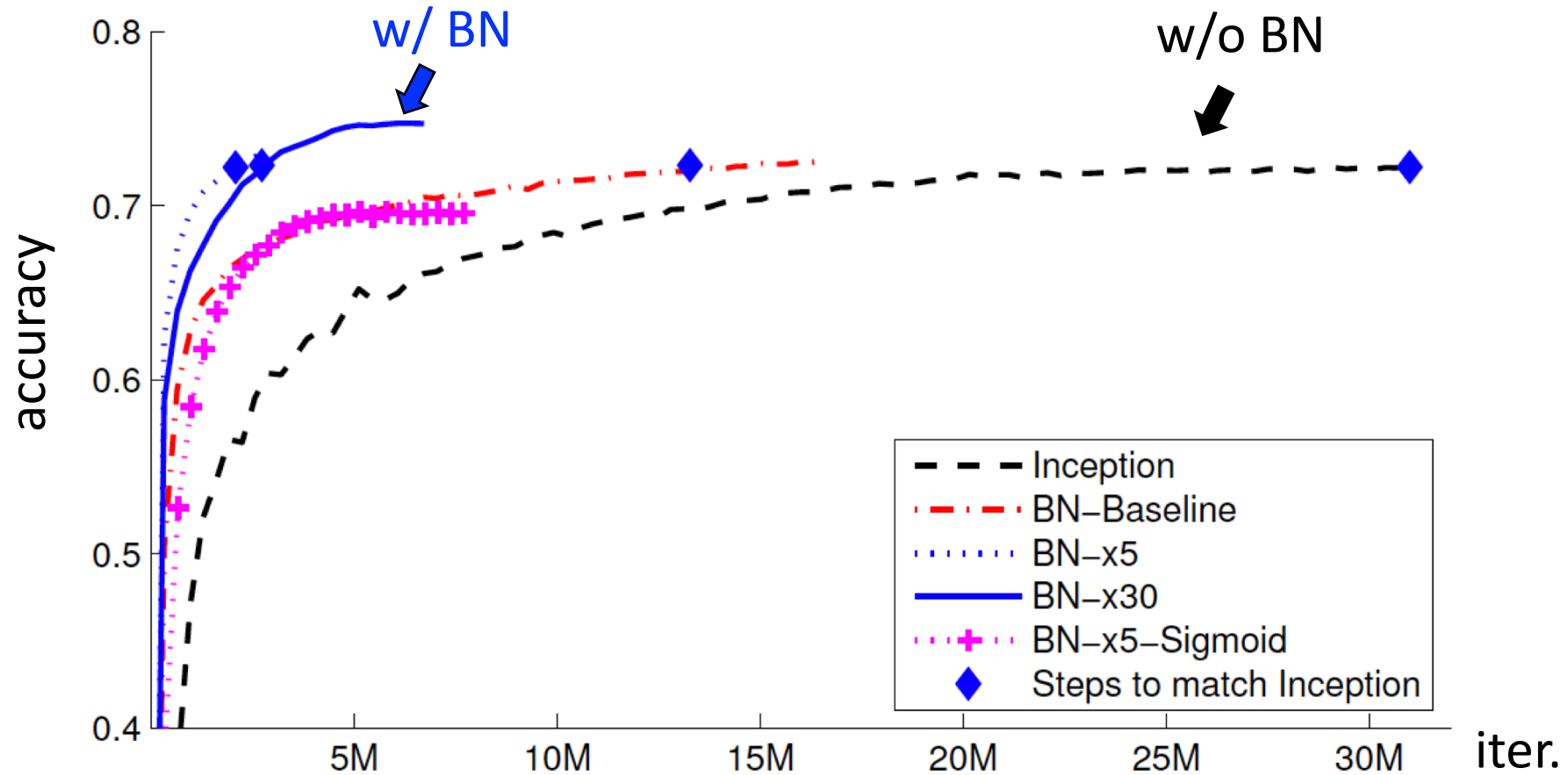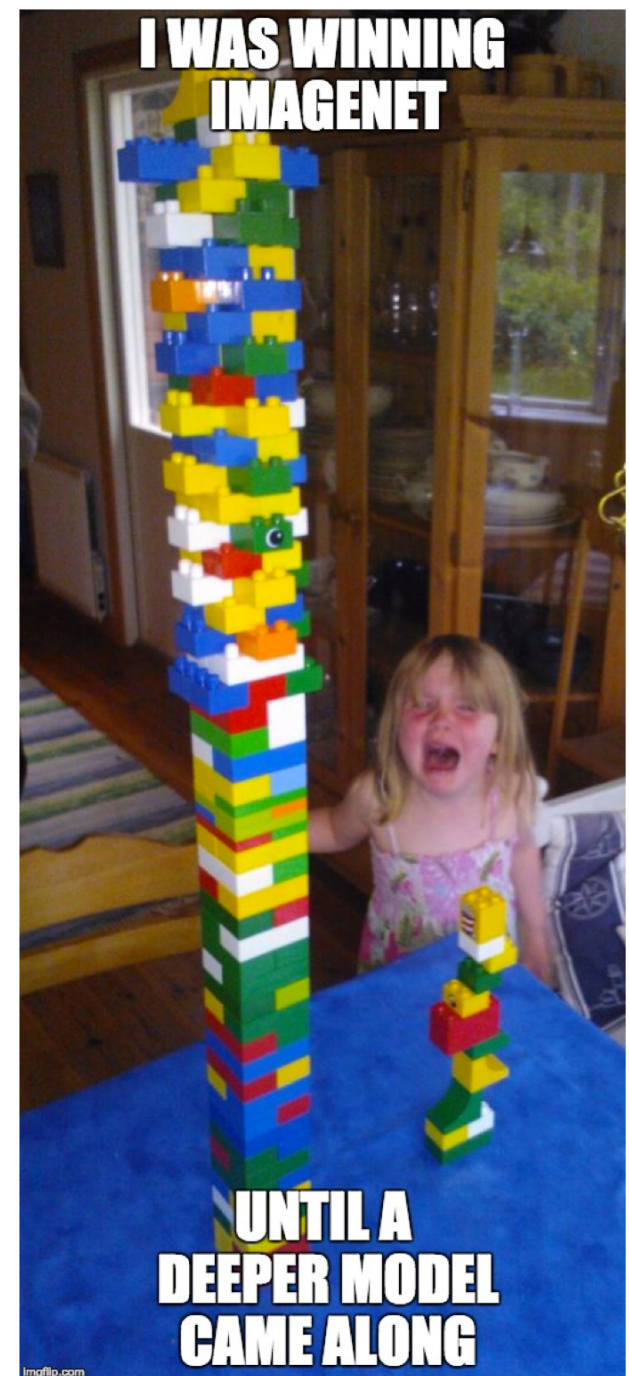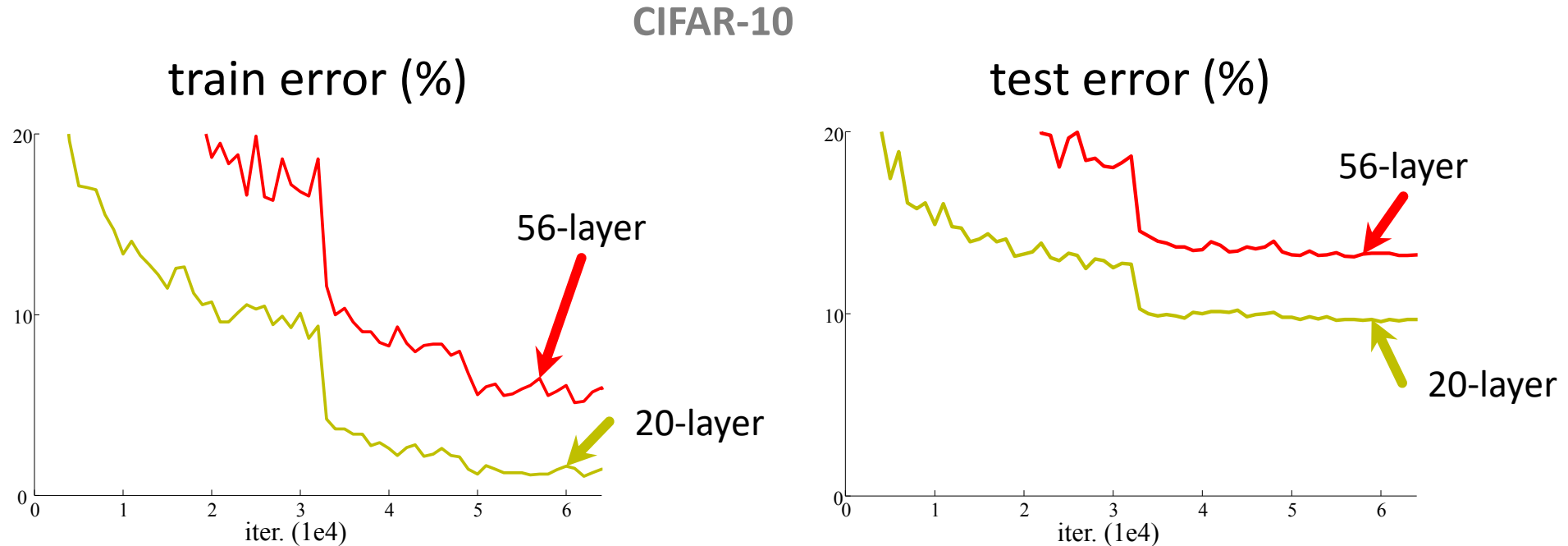(this causes many of my bugs in my research experience!)

Ioffe & Szegedy. "Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift". ICML 2015

# Batch Normalization (BN)



w/ BN

w/o BN

accuracy

- - - Inception
- — · — BN–Baseline
- ········ BN–x5
- —— BN–x30
- +·+·+ BN–x5–Sigmoid
- ◆ Steps to match Inception

iter.

Figure credit: Ioffe & Szegedy

Ioffe & Szegedy. "Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift". ICML 2015

# ResNets



I WAS WINNING IMAGENET

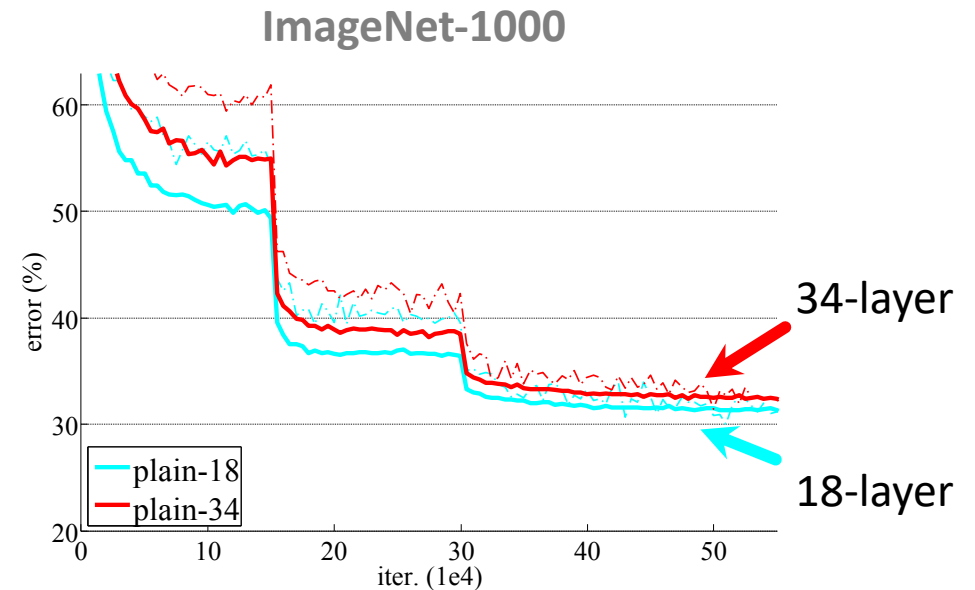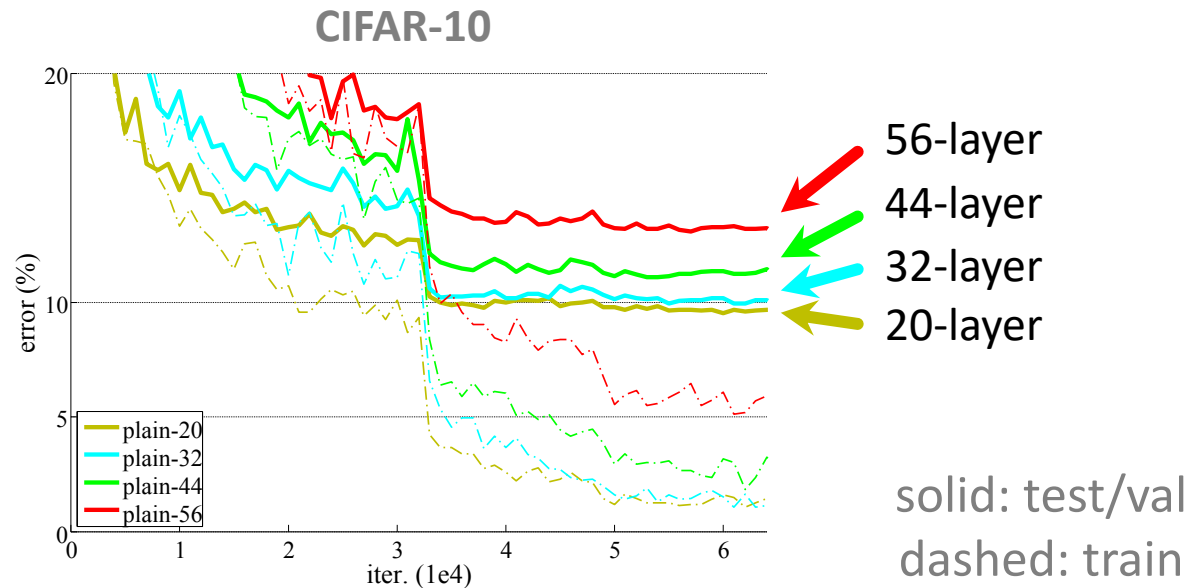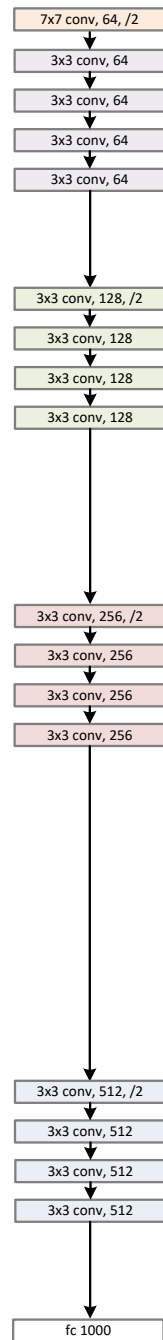UNTIL A DEEPER MODEL CAME ALONG

Credit: ???

# Simply stacking layers?

**CIFAR-10**



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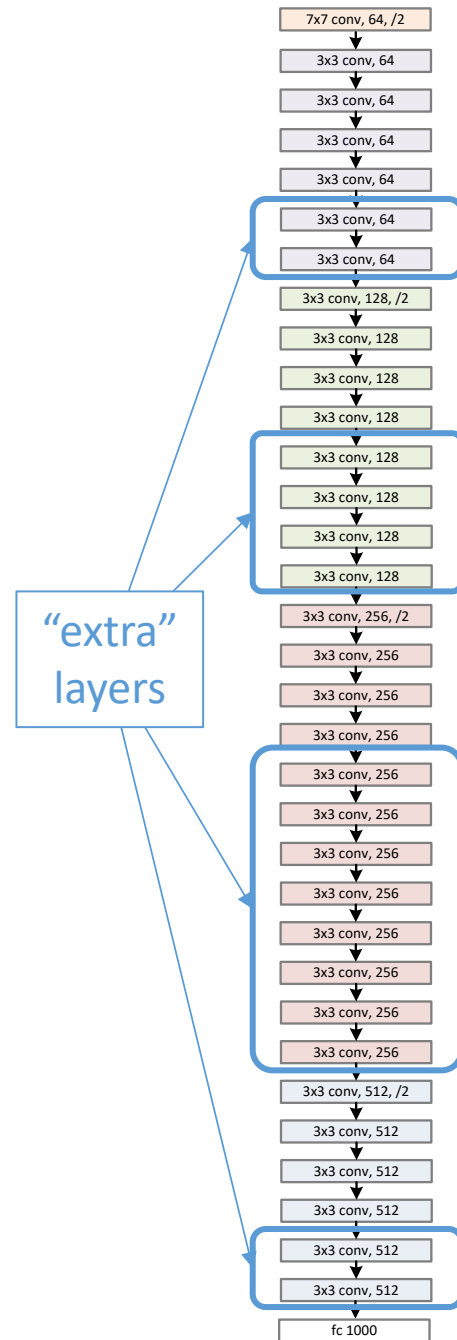- 56-layer net has **higher training error** and test error than 20-layer net

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

# Simply stacking layers?



CIFAR-10

ImageNet-1000

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

solid: test/val
dashed: train

34-layer
18-layer

- "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.
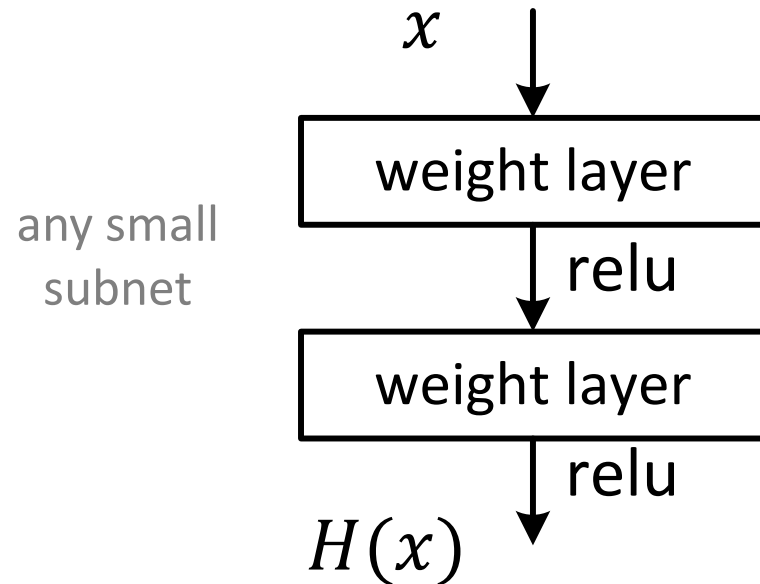
a shallower model (18 layers)

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

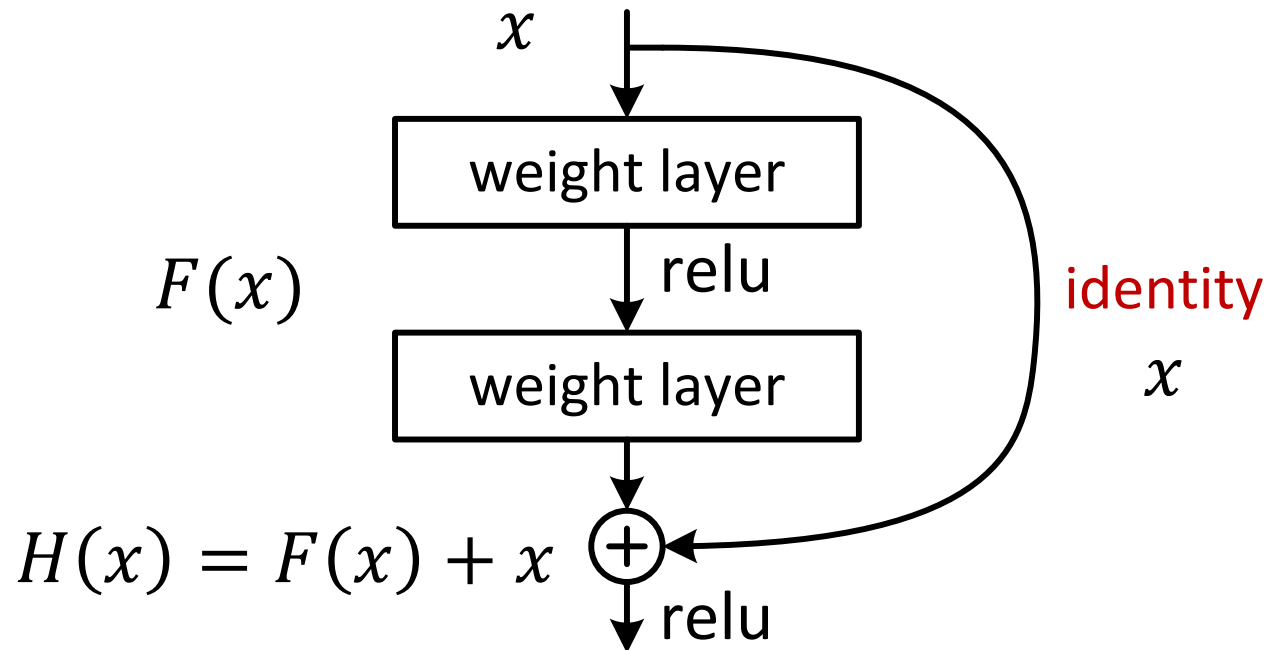- Optimization difficulties: solvers cannot find the solution when going deeper...

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

# Deep Residual Learning

- Plain net

$x$

any small
subnet

weight layer

relu

weight layer

relu

$H(x)$

$H(x)$ is any desired mapping,

hope the small subnet 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

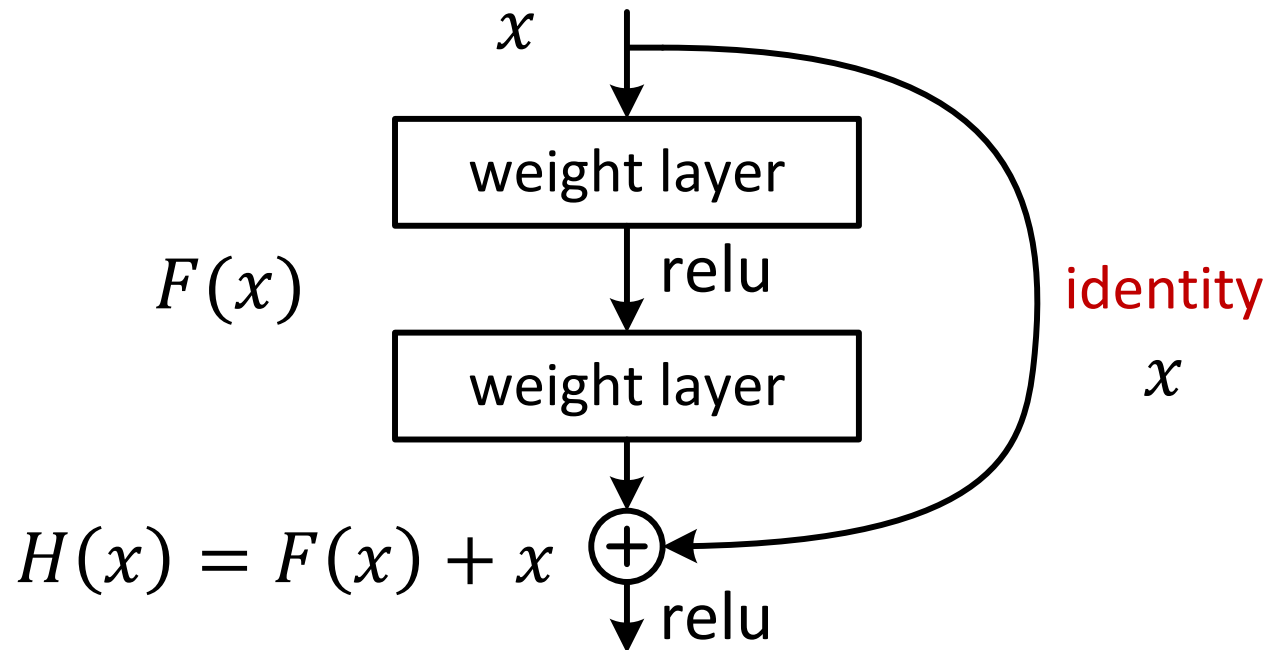$F(x)$    relu

weight layer

identity
$x$

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

relu

$H(x)$ is any desired mapping,

~~hope the small subnet fit $H(x)$~~

hope the small subnet fit $F(x)$

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

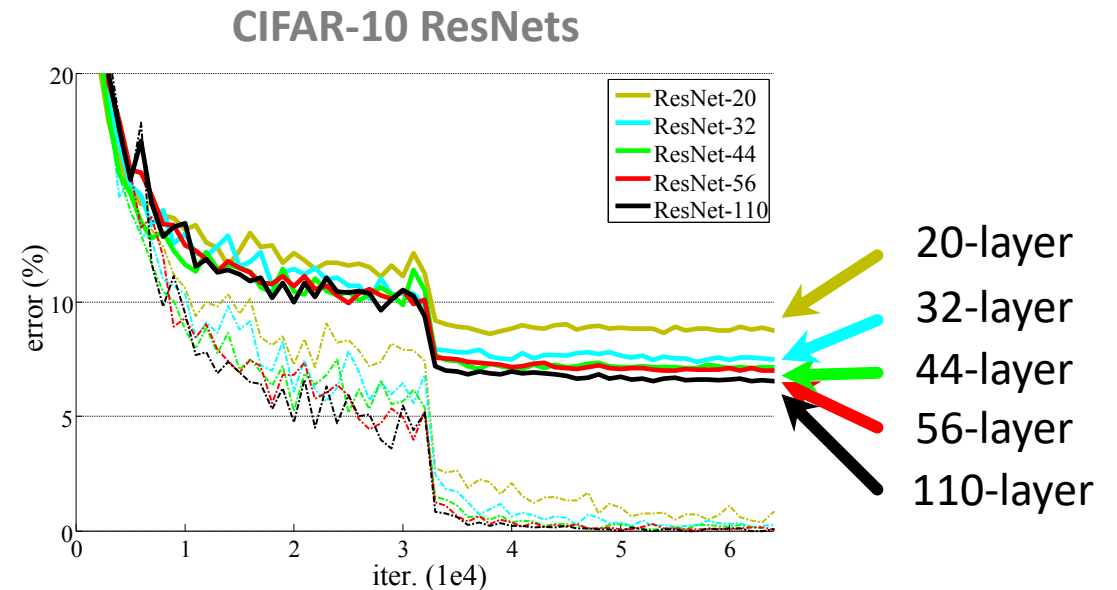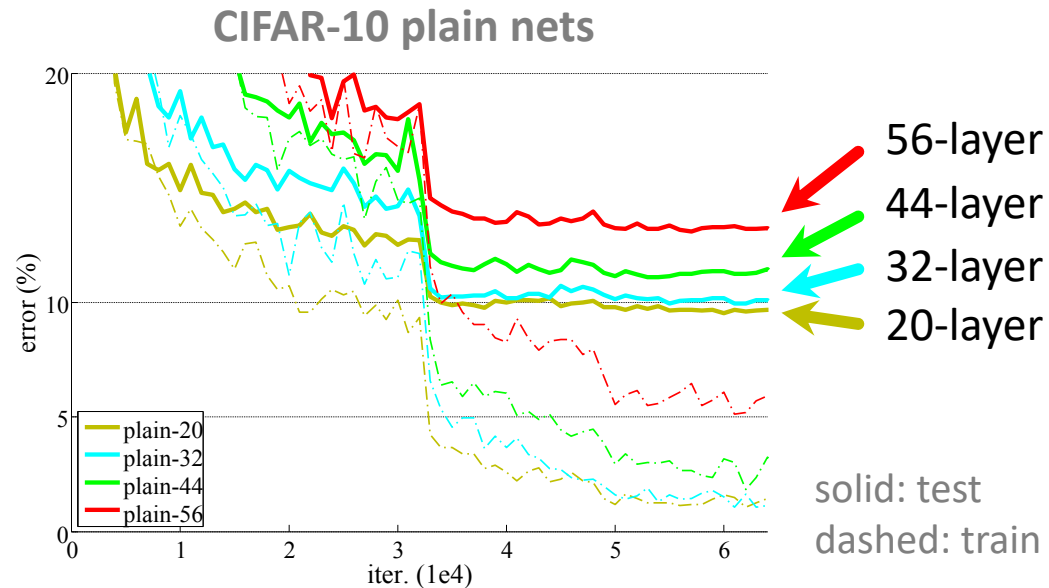Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". CVPR 2016.

# Deep Residual Learning

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

$$x$$

weight layer

relu

$$F(x)$$
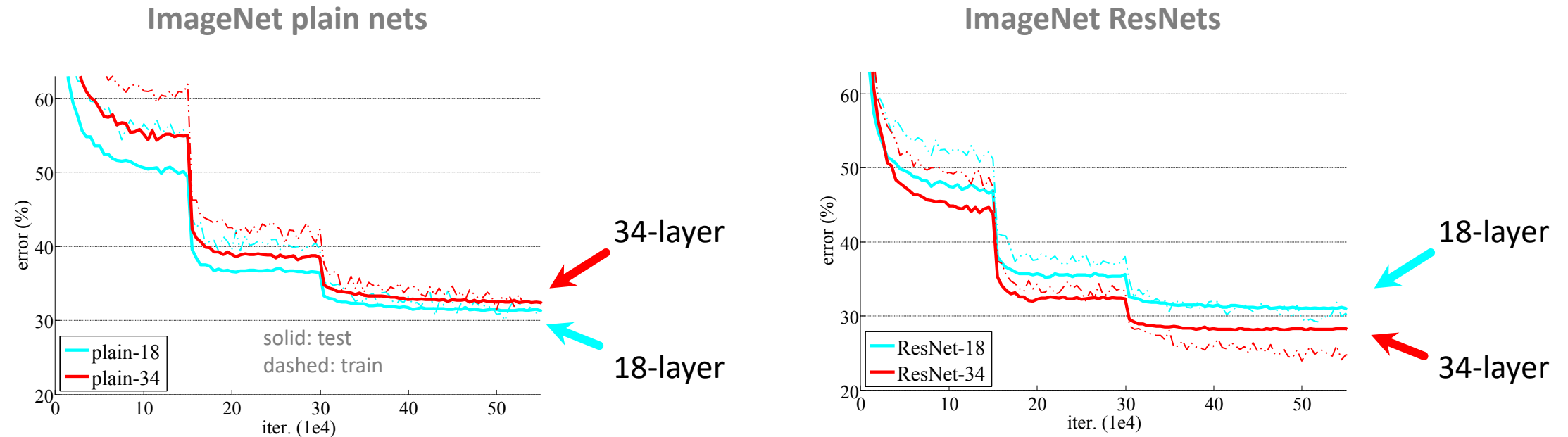
weight layer

identity

$$x$$

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

relu

- 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.

# CIFAR-10 experiments

**CIFAR-10 plain nets**

error (%)

20 — 56-layer
— 44-layer
— 32-layer
10 — 20-layer

5

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

solid: test
dashed: train

iter. (1e4)

**CIFAR-10 ResNets**

error (%)

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

20 — 20-layer
— 32-layer
10 — 44-layer
— 56-layer
5 — 110-layer
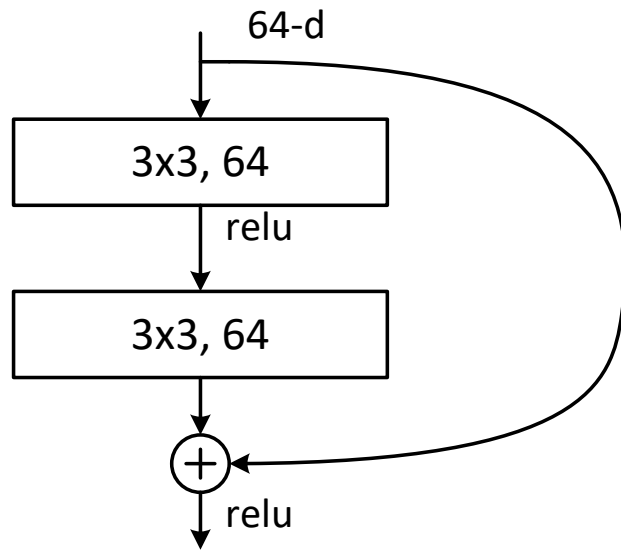
iter. (1e4)
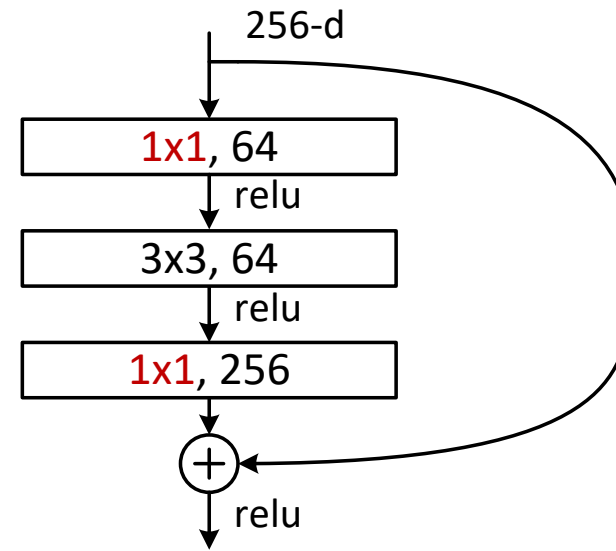
- Deep ResNets can be trained without difficulties
- Deeper ResNets have **lower training error**, and also lower test error

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

# ImageNet experiments

**ImageNet plain nets**

**ImageNet ResNets**



- 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

- A practical design of going deeper



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

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

# 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". CVPR 2016.

# ResNet beyond computer vision

- **Neural Machine Translation** (NMT): 8-layer LSTM!



residual connections

residual connections

Wu et al. "Google's Neural Machine Translation System: Bridging the Gap between Human and Machine Translation". arXiv 2016.

# ResNet beyond computer vision

- **Speech Synthesis** (WaveNet): Residual CNNs on 1-d sequence



residual connections

van den Oord et al. "WaveNet: A Generative Model for Raw Audio". arXiv 2016.

# ResNet beyond computer vision

- **AlphaGo Zero**: 40 Residual Blocks



"Mastering the game of Go without human knowledge", Silver et al. Nature 2017

# ResNeXt

- Recap: shortcut, bottleneck, and <u>multi-branch</u>



**Inception**:
heterogeneous multi-branch

**ResNeXt**:
uniform multi-branch

Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. "Aggregated Residual Transformations for Deep Neural Networks". arXiv 2016 (CVPR 2017).
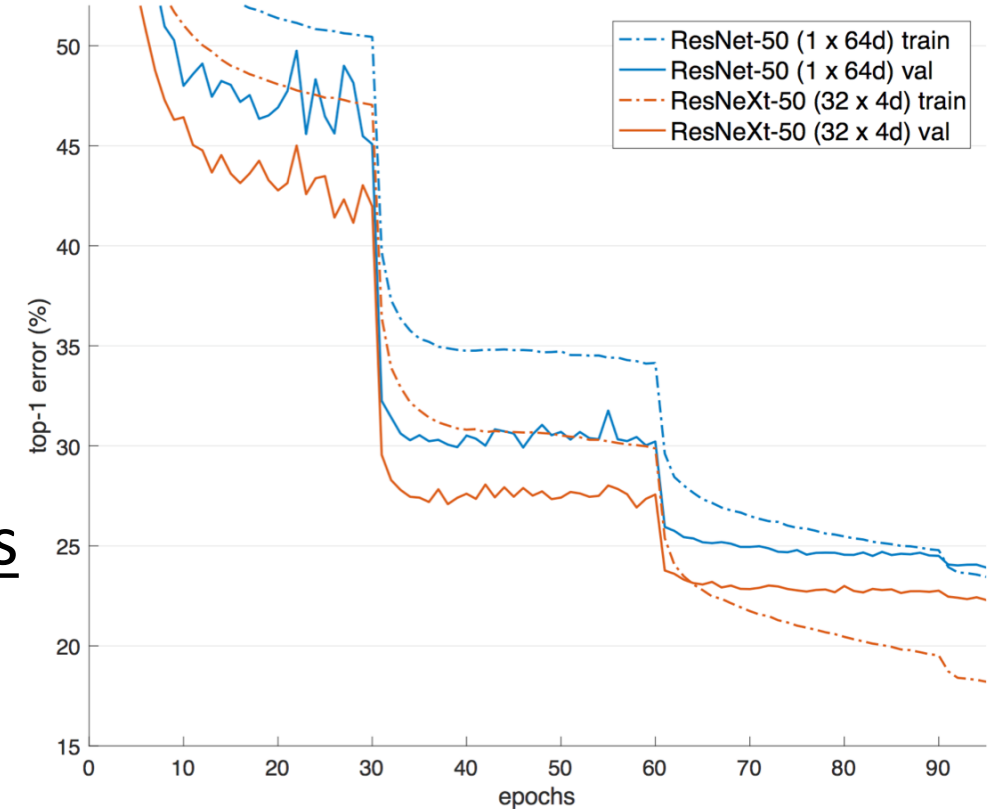
# ResNeXt

- Concatenation and Addition are interchangeable
  - General property for DNNs; not only limited to ResNeXt
- Uniform multi-branching can be done by group-conv



Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. "Aggregated Residual Transformations for Deep Neural Networks". arXiv 2016 (CVPR 2017).

# ResNeXt

- Better accuracy
  - when having the same FLOPs/#params as a baseline ResNet

- <u>Better trade-off for high-capacity models</u>



Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. "Aggregated Residual Transformations for Deep Neural Networks". arXiv 2016 (CVPR 2017).

# Competition winners using ResNeXt

ResNeXt is a good trade-off for high-capacity:

- ImageNet Classification 2017, 1$^{st}$ place
  - SE-ResNeXt
- COCO Object Detection 2017, 1$^{st}$ place
  - MegDet + ResNeXt
- COCO Instance Segmentation 2017, 1$^{st}$ place
  - PANet + ResNeXt
- COCO Stuff Segmentation 2017, 1$^{st}$ place
  - FPN + ResNetXt
- …

Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. "Aggregated Residual Transformations for Deep Neural Networks". arXiv 2016 (CVPR 2017).
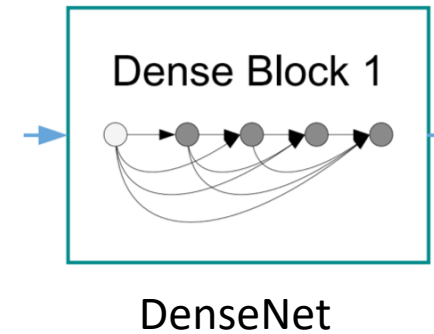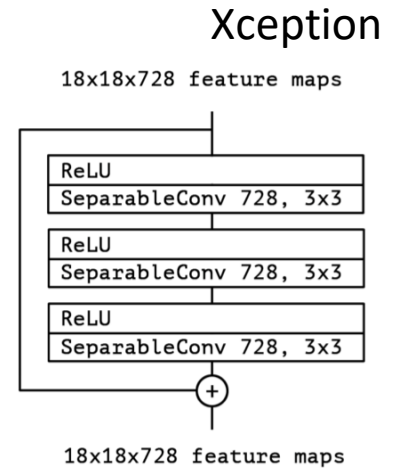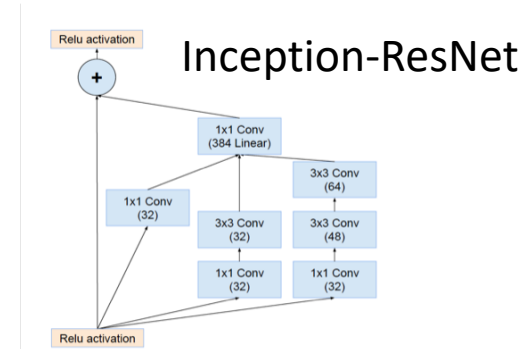
# ResNeXt: higher capacity for billion-scale images



Fig. 5: Classification accuracy on val-IN-1k using ResNeXt-101 32×{4, 8 16, 32, 48}d with and without pretraining on the IG-940M-1.5k dataset.

"Exploring the Limits of Weakly Supervised Pretraining". arXiv 2018.
Dhruv Mahajan, Ross Girshick, Vignesh Ramanathan, Kaiming He, Manohar Paluri, Yixuan Li, Ashwin Bharambe, and Laurens van der Maaten.
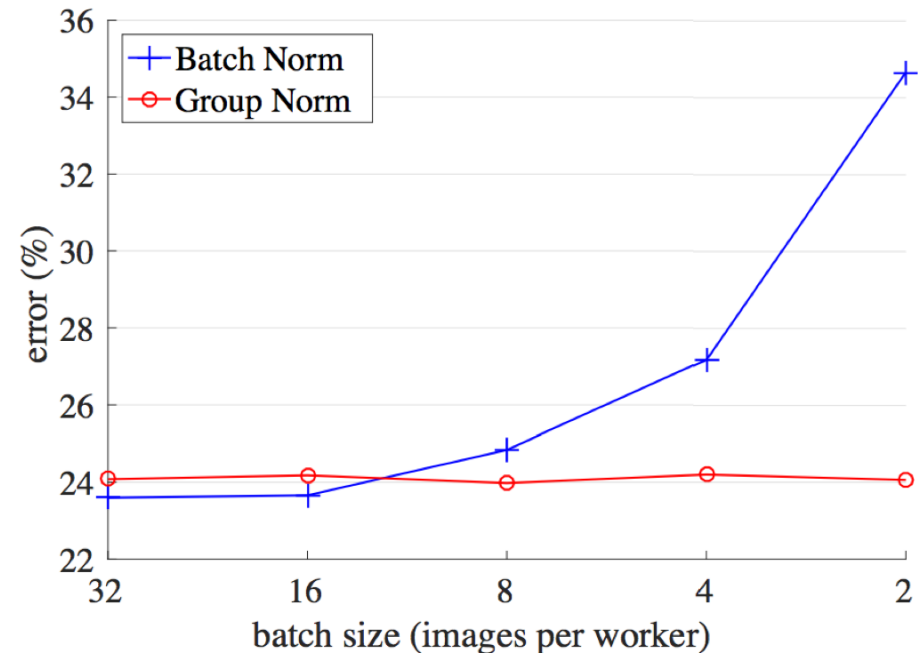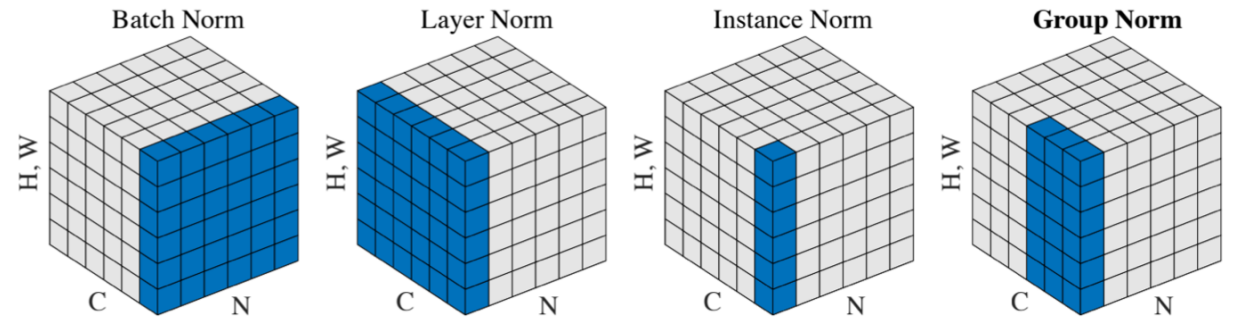
# More architectures (not covered in this tutorial)

- ## Inception-ResNet [Szegedy et al 2017]
  - Inception as transformation + residual connection

- ## DenseNet [Huang et al CVPR 2017]
  - Densely connected shortcuts w/ concat.

- ## Xception [Chollet CVPR 2017], MobileNets [Howard et al 2017]
  - DepthwiseConv (i.e., GroupConv with #group=#channel)

- ## ShuffleNet [Zhang et al 2017]
  - More Group/DepthwiseConv + shuffle

- ......



Inception-ResNet



Xception



DenseNet



ShuffleNet

# *Teaser*: Group Normalization (GN)

- Independent of batch size

- Robust to small batches

- Enable new scenarios:
  e.g.: 41 AP on COCO
  trained from scratch

Yuxin Wu and Kaiming He. "Group Normalization". arXiv 2018

# Conclusion

- Deep Learning is Representation Learning

- <u>Represent data</u> for machines to perform tasks (this talk)

- Represent data for machines to <u>perform tasks</u> (next talks)