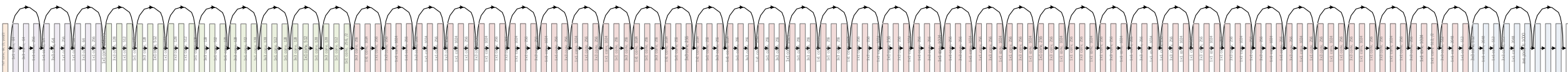# Deep Residual Networks (ResNet)

- Deep Residual Learning for Image Recognition
- Identity Mappings in Deep Residual Networks

# Background of Resnet's appearance

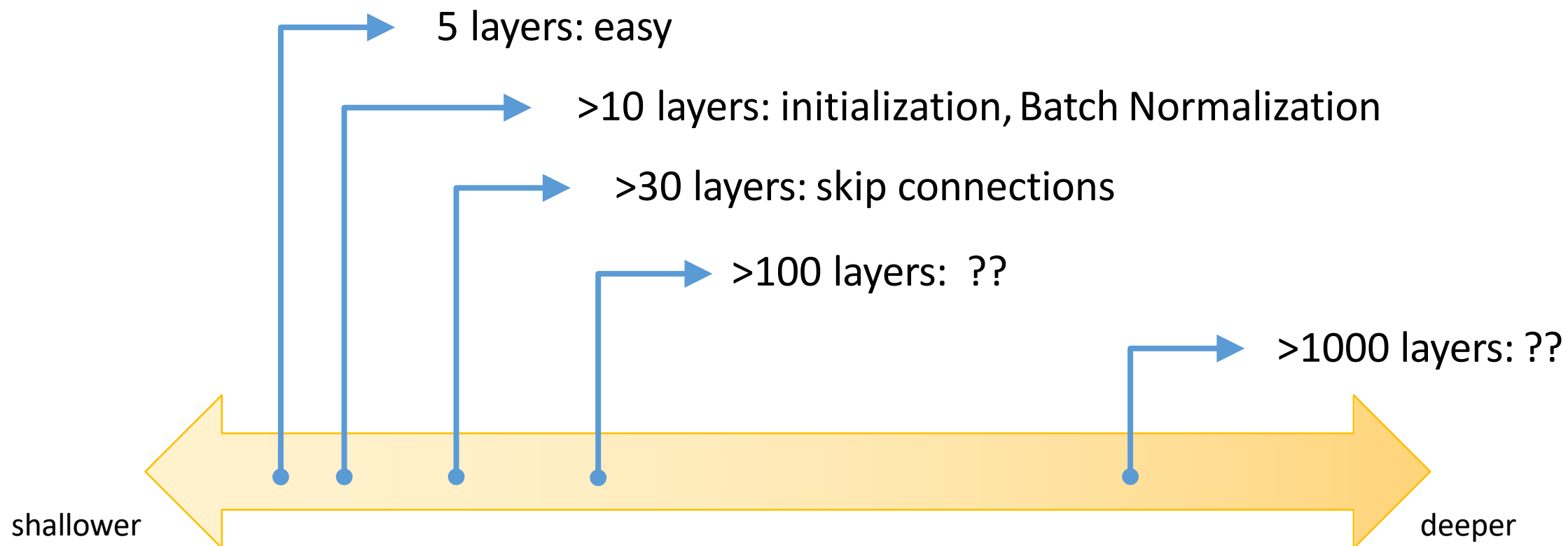# Table of Contents

- Background of ResNet's Appearance

- ResNet's Advantage

- ResNet's Architecture

- ResNet Variations

  - Deeper Bottleneck Architecture

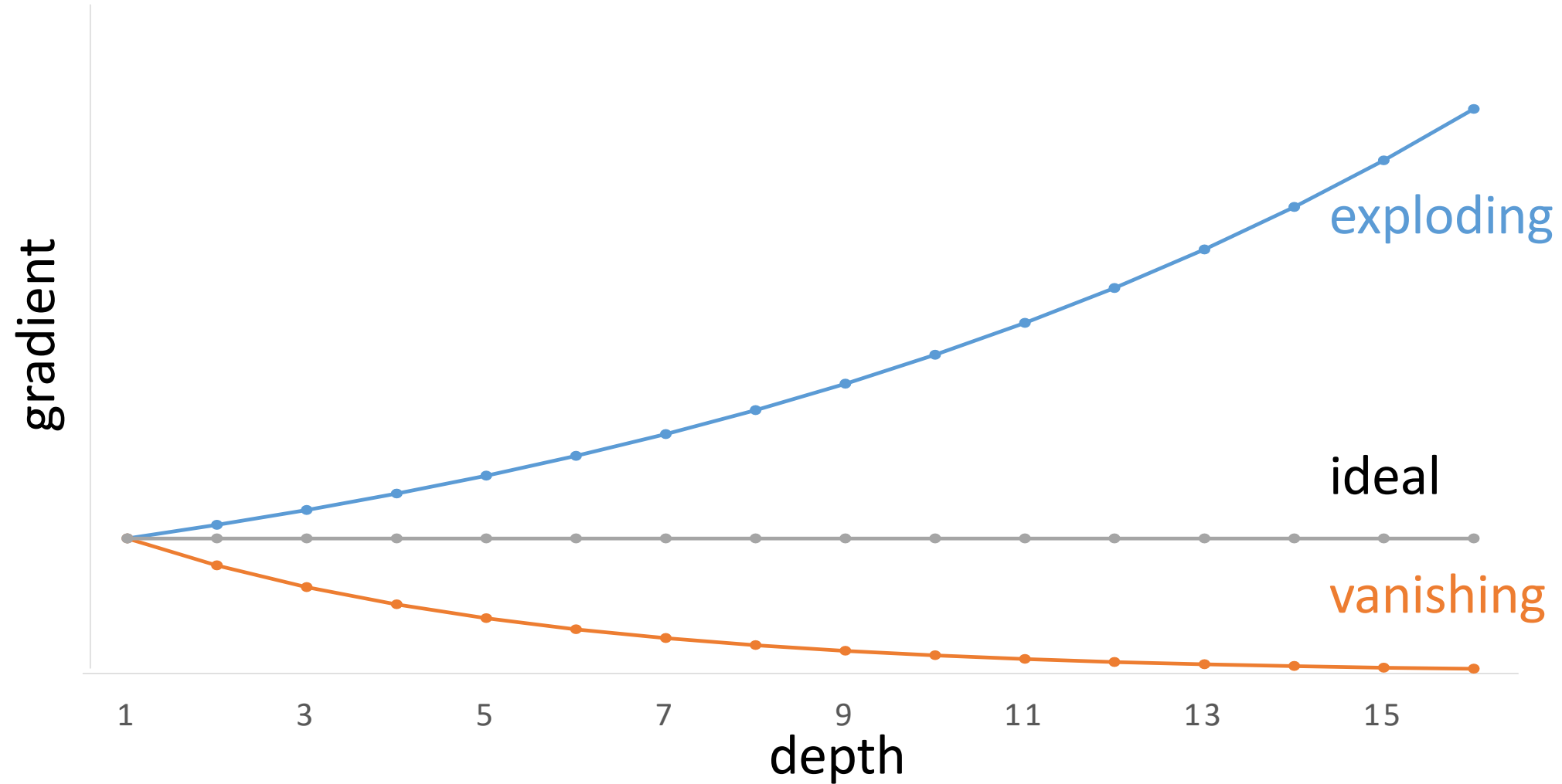  - Pre-Activation ResNet
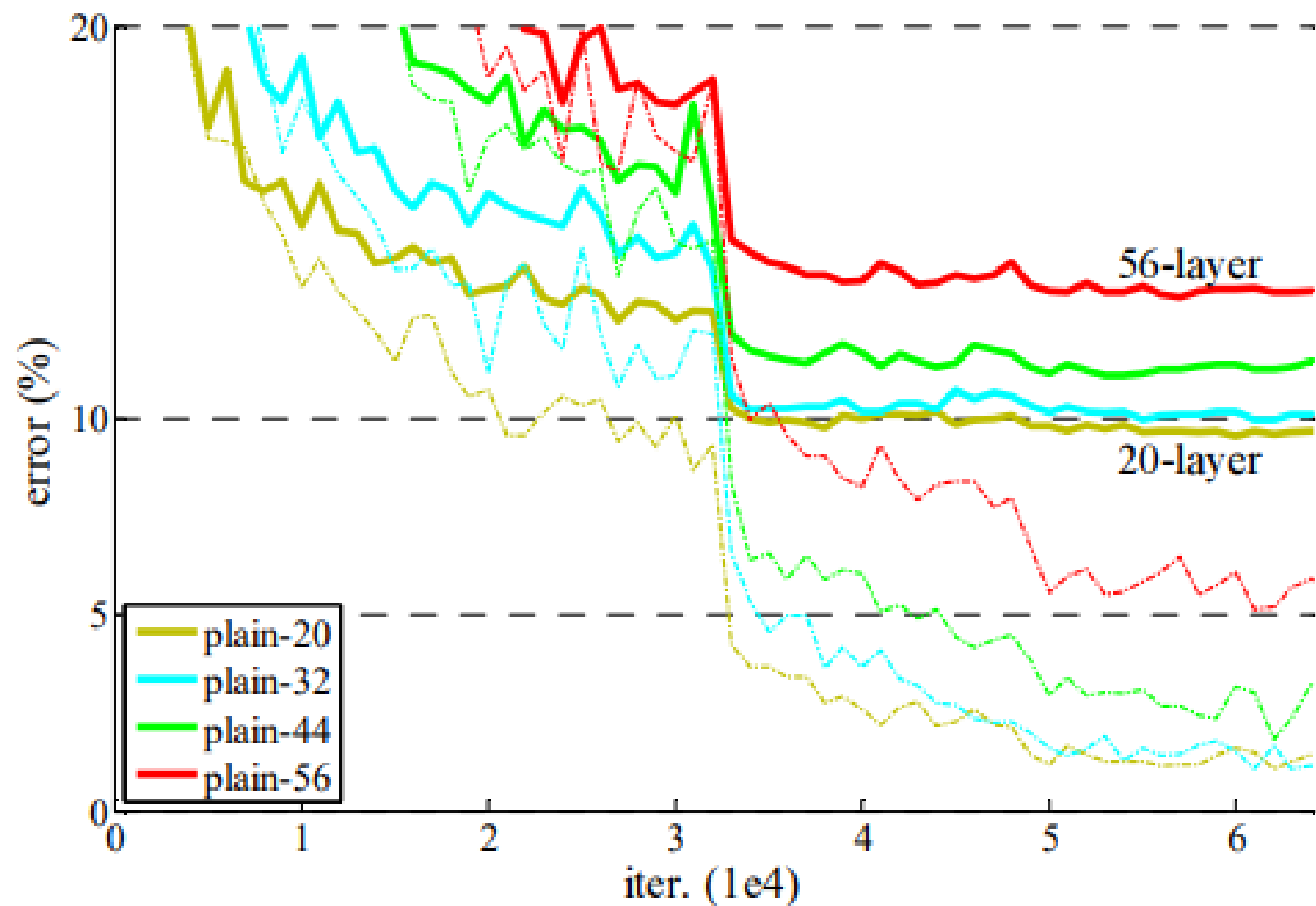
  - ResNeXt

# Before ResNet

5 layers: easy

>10 layers: initialization, Batch Normalization

>30 layers: skip connections

>100 layers:  ??

>1000 layers: ??

shallower

deeper

# Problems with deep layers

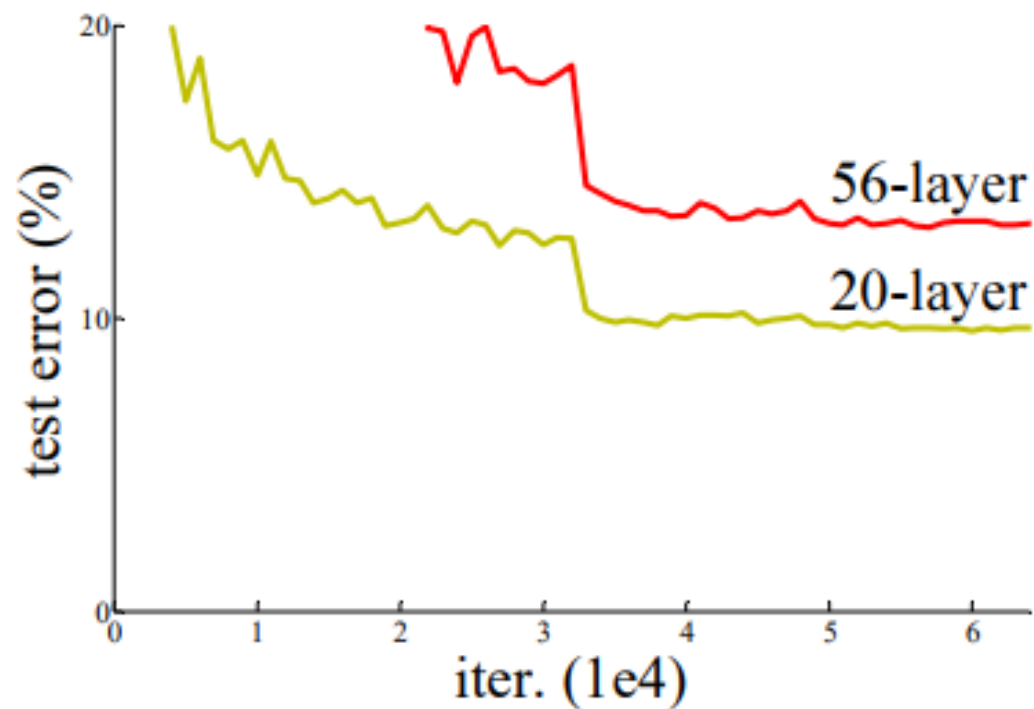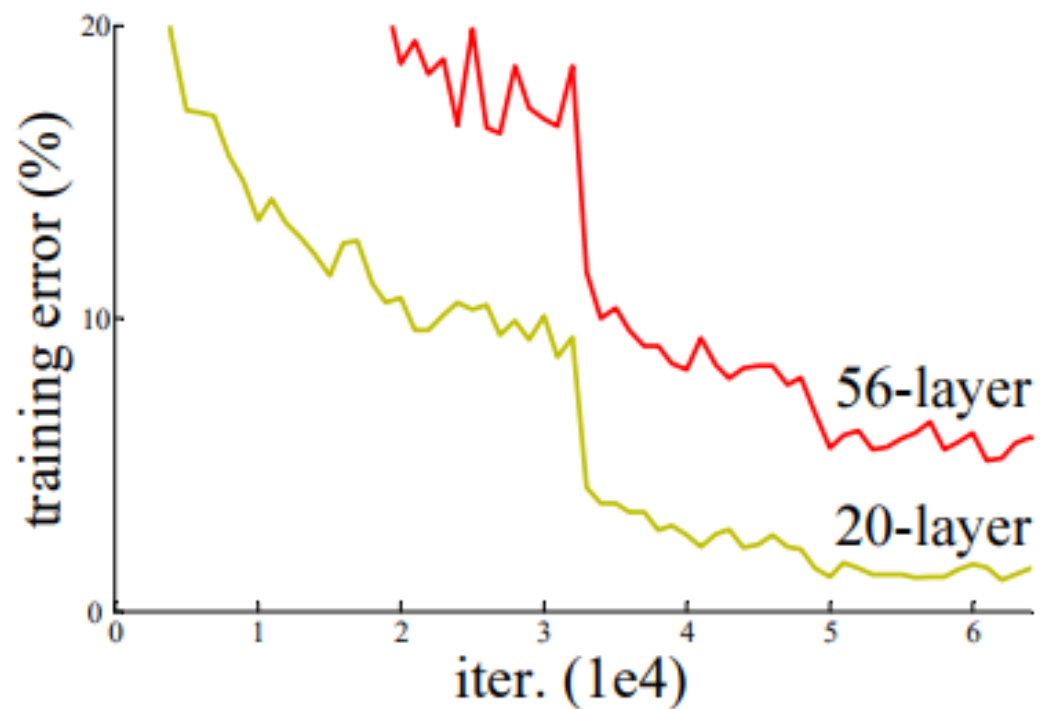1. Vanishing gradient problem

2. Degradation problem

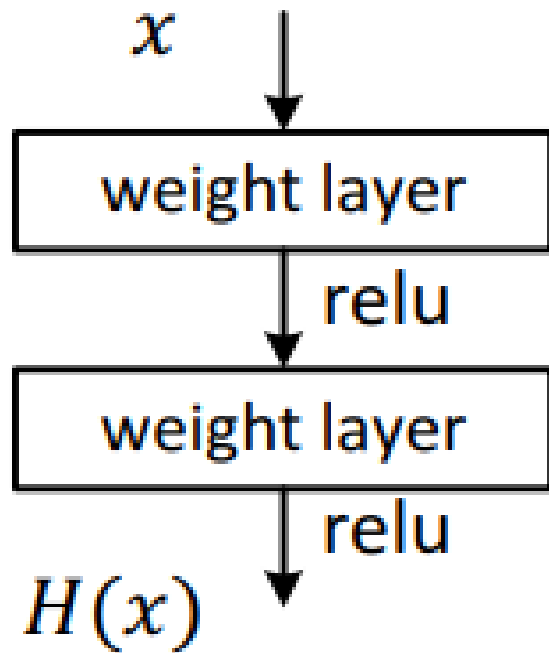# Vanishing/Exploding Gradient Problem
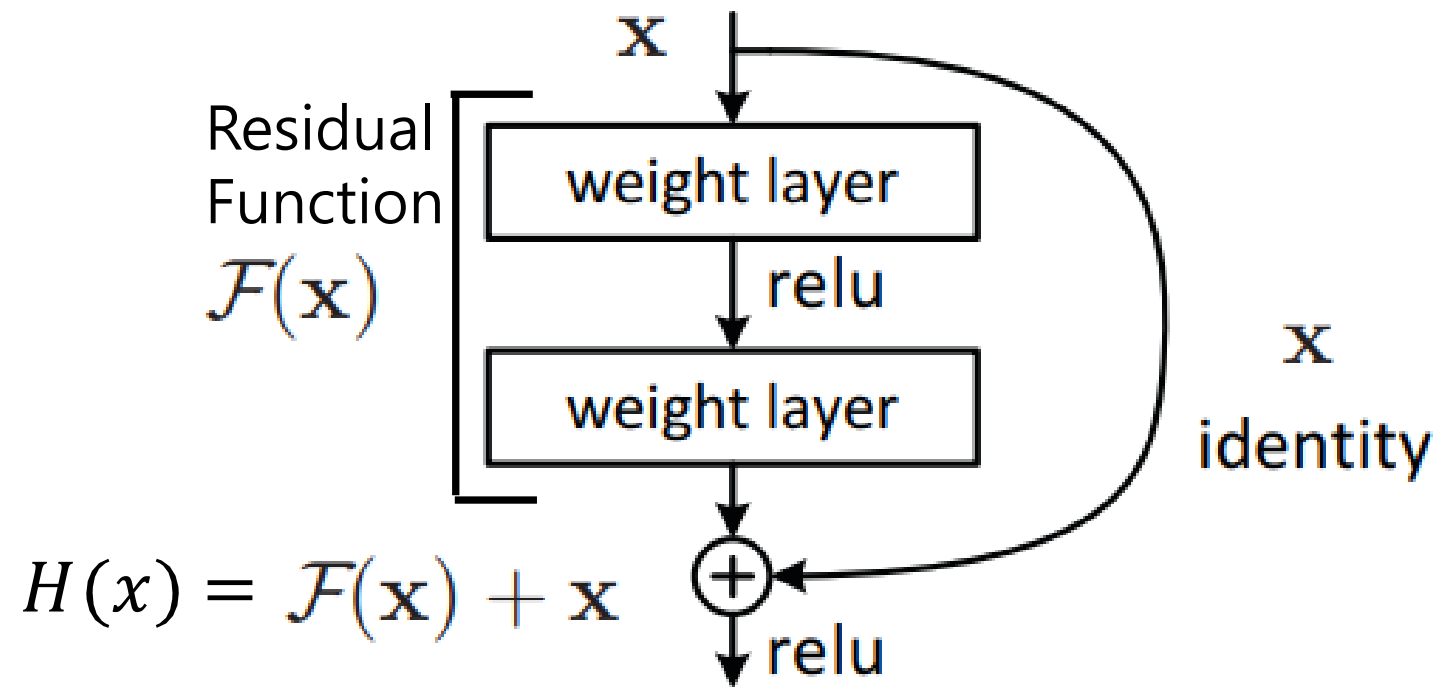
# Degradation Problem

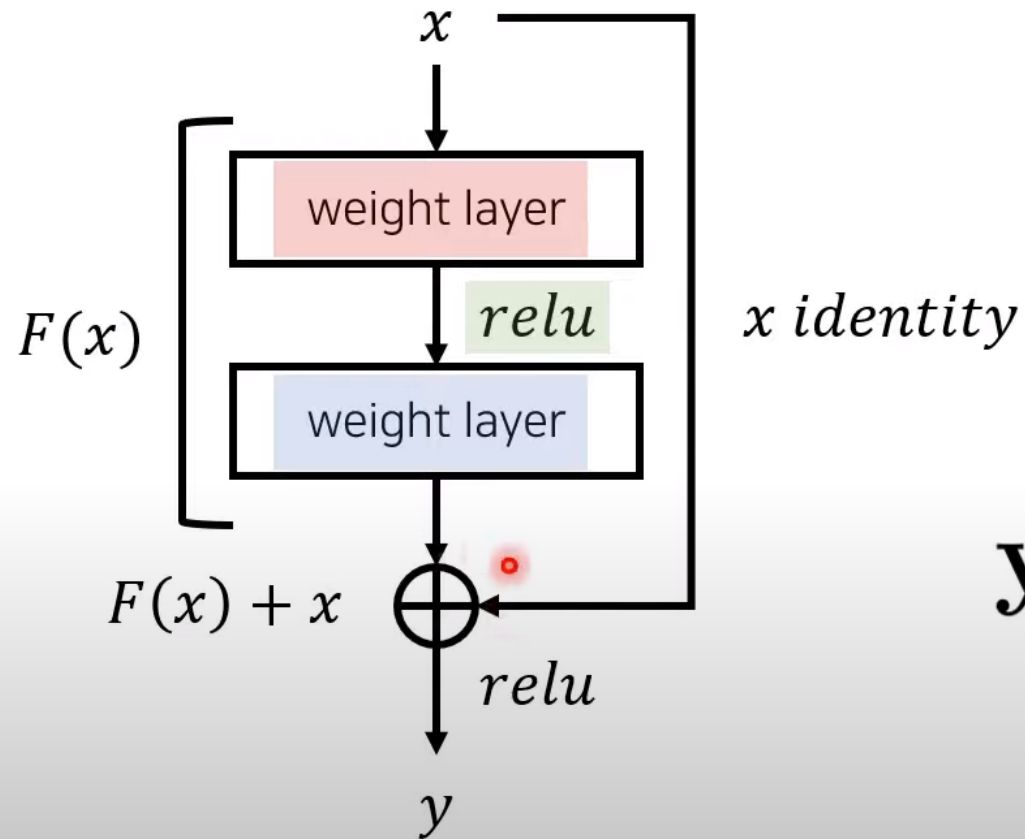# Degradation Problem

# What is Residual Learning?



Plain net

Residual net

# ResNet's Advantage
## 1) # of parameters doesn't increase



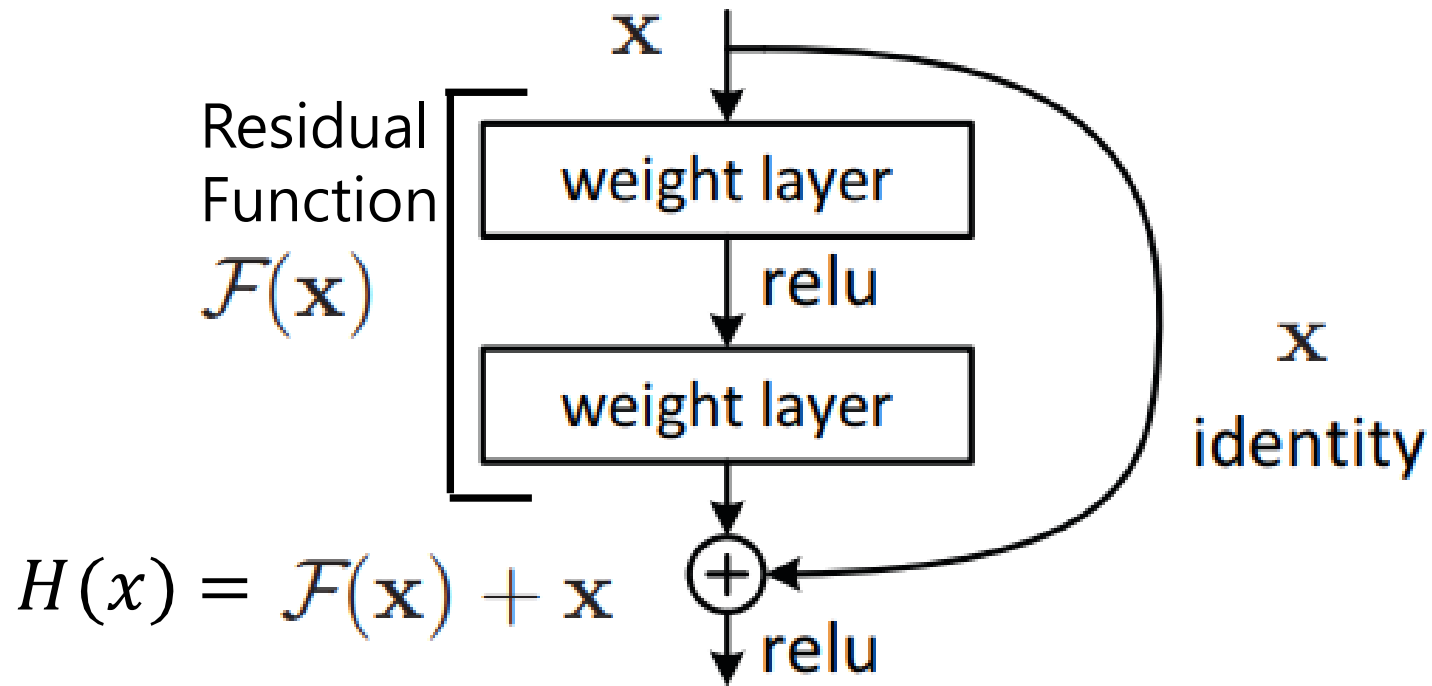$$\mathcal{F} = W_2 \sigma(W_1 \mathbf{x})$$

일반적인 형태

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

multiple convolutional layers          shortcut

# ResNet's Advantage
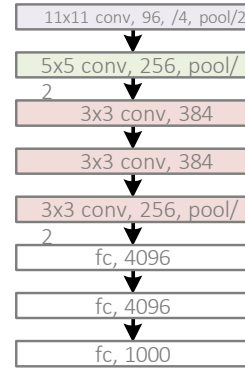## 2) Solve Vanishing Gradient Problem



Residual Function $\mathcal{F}(\mathbf{x})$

$$H(x) = \mathcal{F}(\mathbf{x}) + \mathbf{x}$$

Residual net

$$\frac{\partial H}{\partial x} = \frac{\partial F}{\partial x} + 1$$

# ResNet's Advantage
# 3) High Accuracy in Deep structure

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)

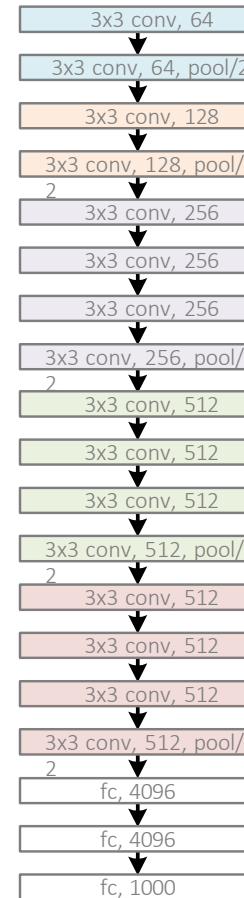| 3x3 conv, 64 |
| 3x3 conv, 64, pool/2 |
| 3x3 conv, 128 |
| 3x3 conv, 128, pool/2 |
| 3x3 conv, 256 |
| 3x3 conv, 256 |
| 3x3 conv, 256 |
| 3x3 conv, 256, pool/2 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512, pool/2 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
|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

GoogleNet, 22 layers
(ILSVRC 2014)

# ResNet's Advantage
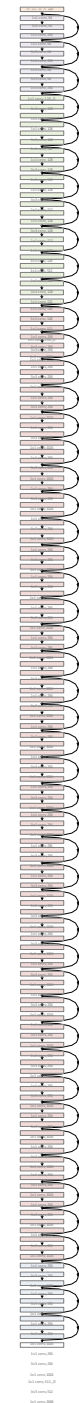## 3) High Accuracy in Deep structure

AlexNet, 8 layers
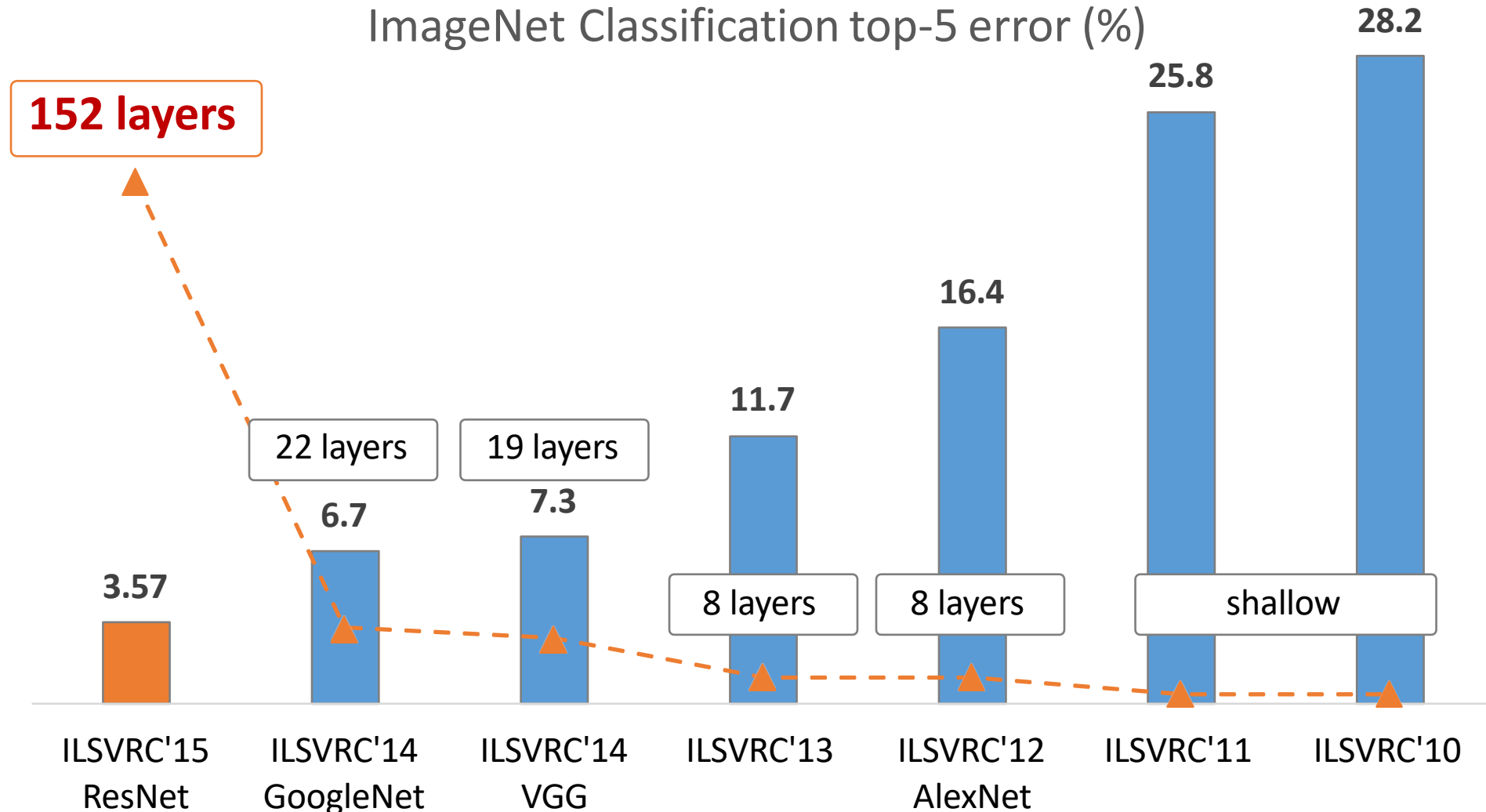(ILSVRC 2012)

VGG, 19 layers
(ILSVRC 2014)
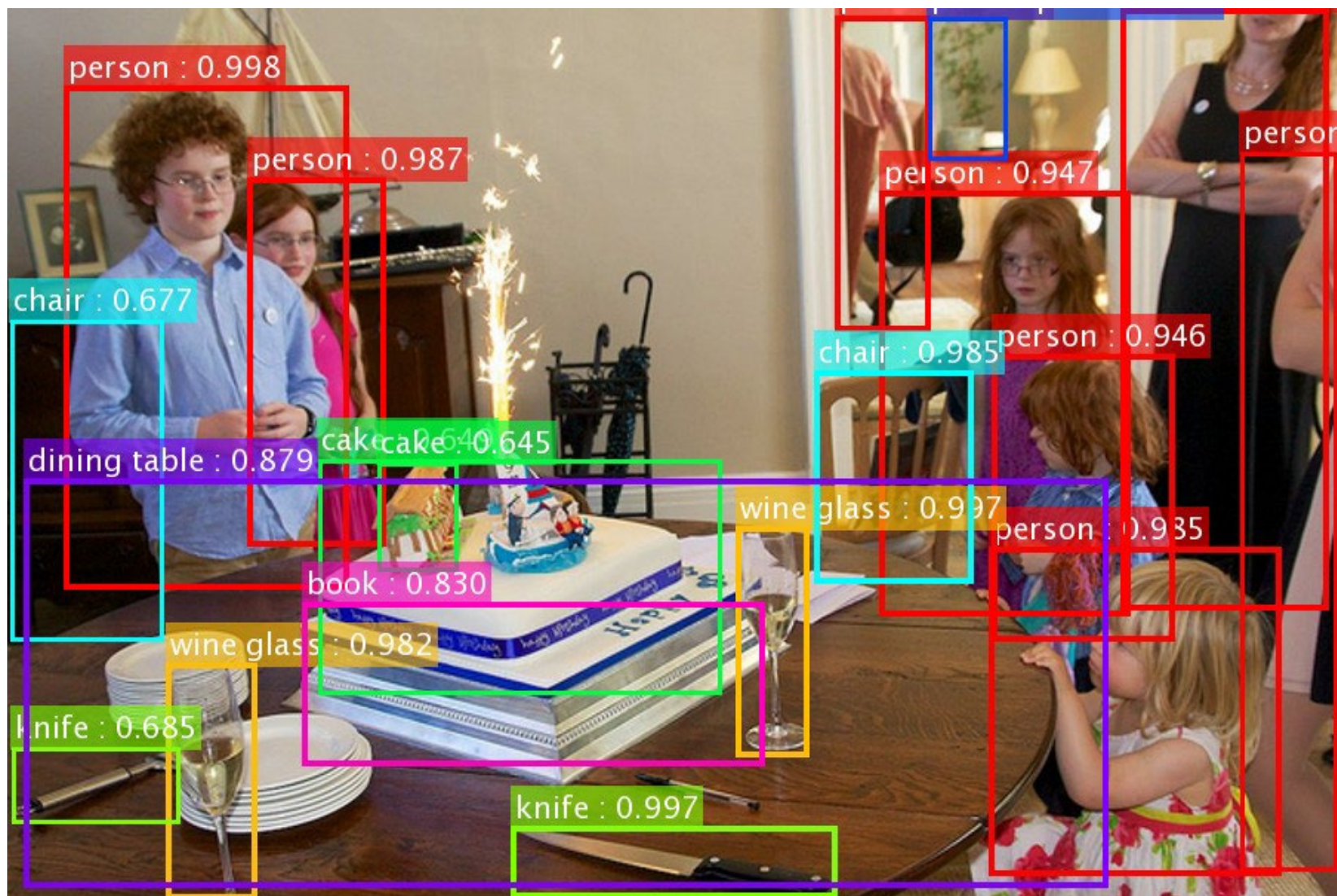
ResNet, 152 layers
(ILSVRC 2015)

# ResNet's Advantage
# 3) High Accuracy in Deep structure

ImageNet Classification top-5 error (%)

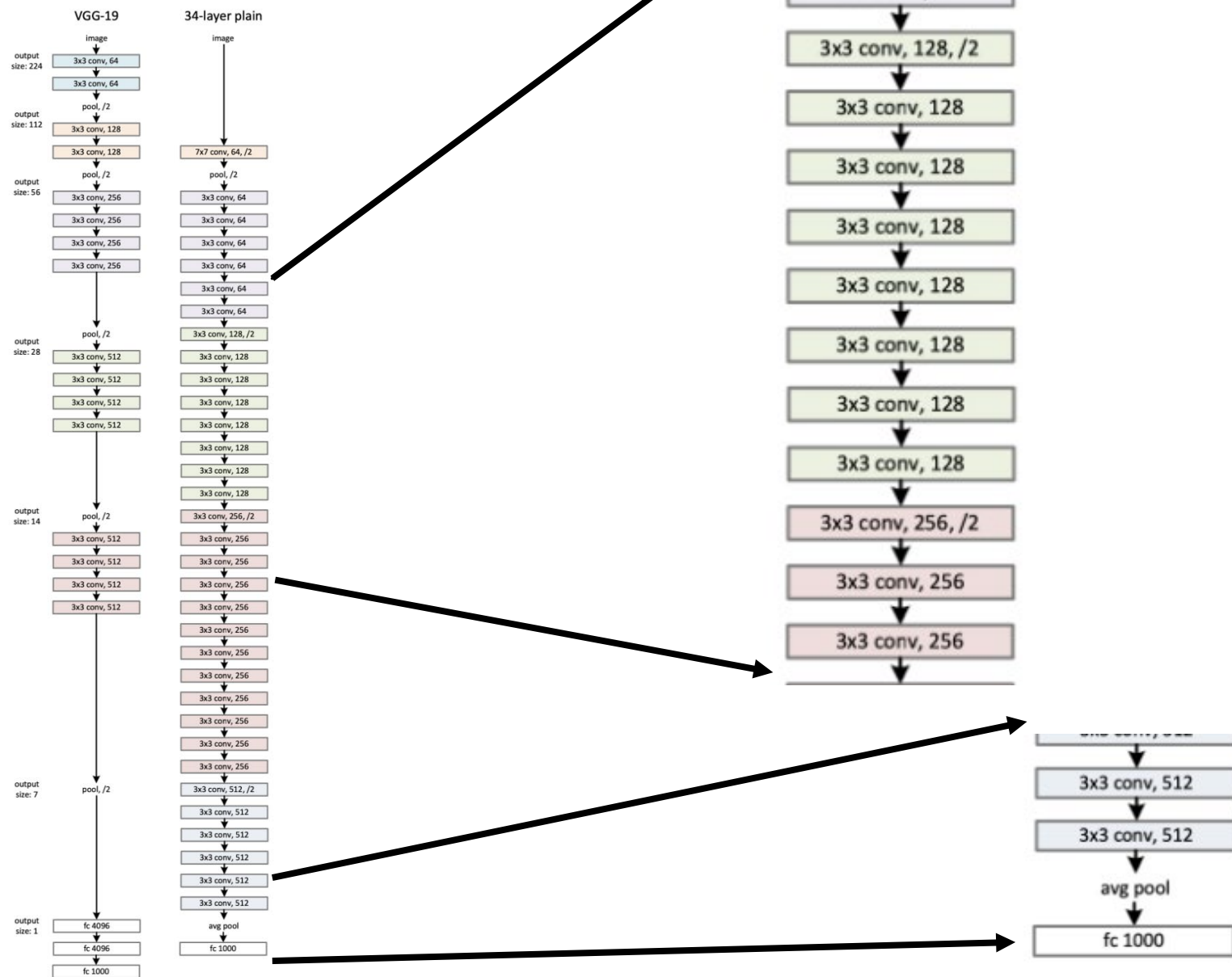# ResNet at ILSVRC & COCO 2015 Competitions

- **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

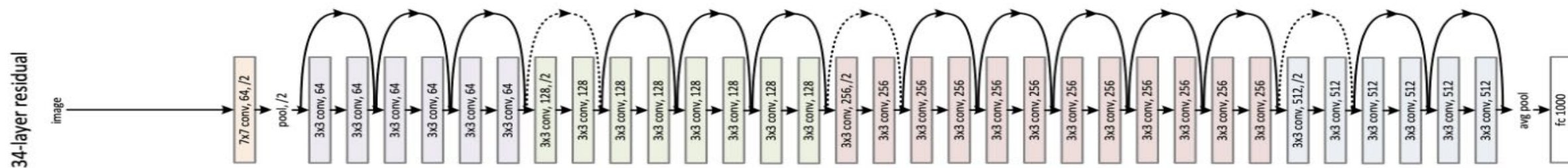*improvements are relative numbers
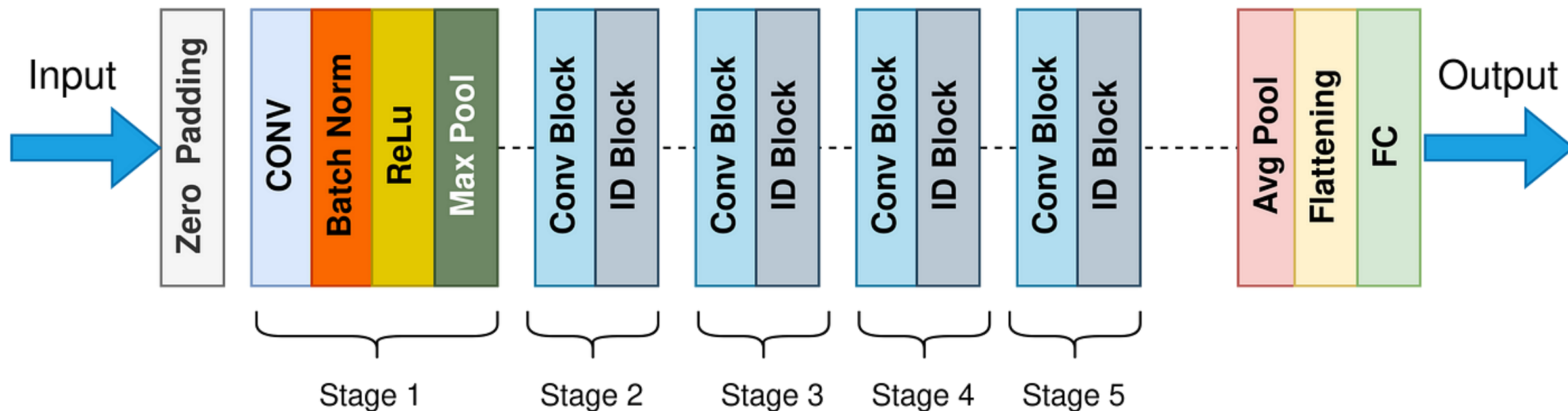
ResNet's object detection result on COCO

# ResNet Architecture

# ResNet Architecture
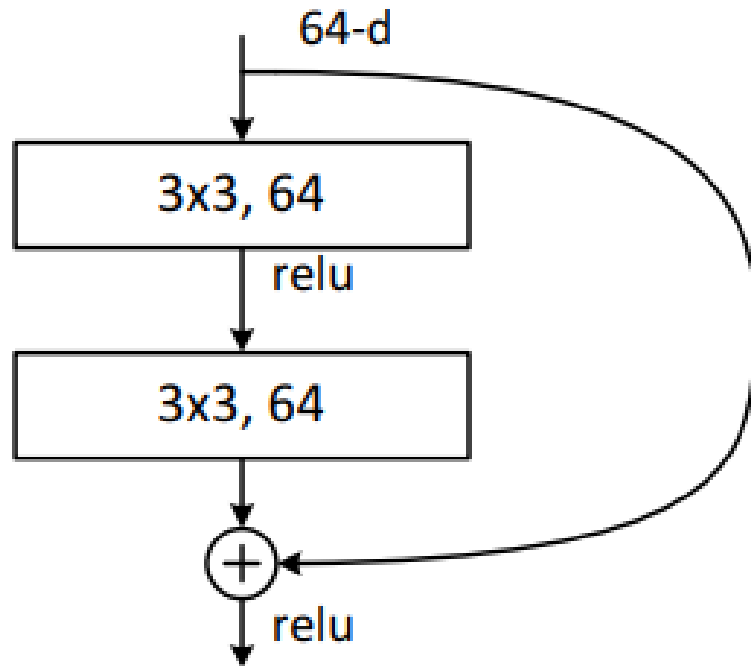


**ResNet50 Model Architecture**

# ResNet Architecture

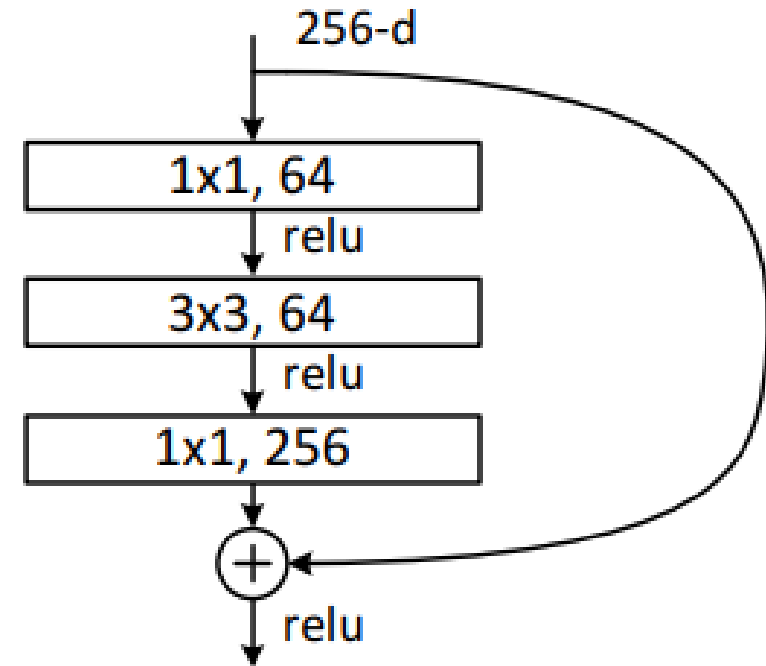- How to align dimensions of input and output

1. Zero padding

2. linear projection $W_s$ 사용: $y = F(x, \{W_i\}) + W_s x$
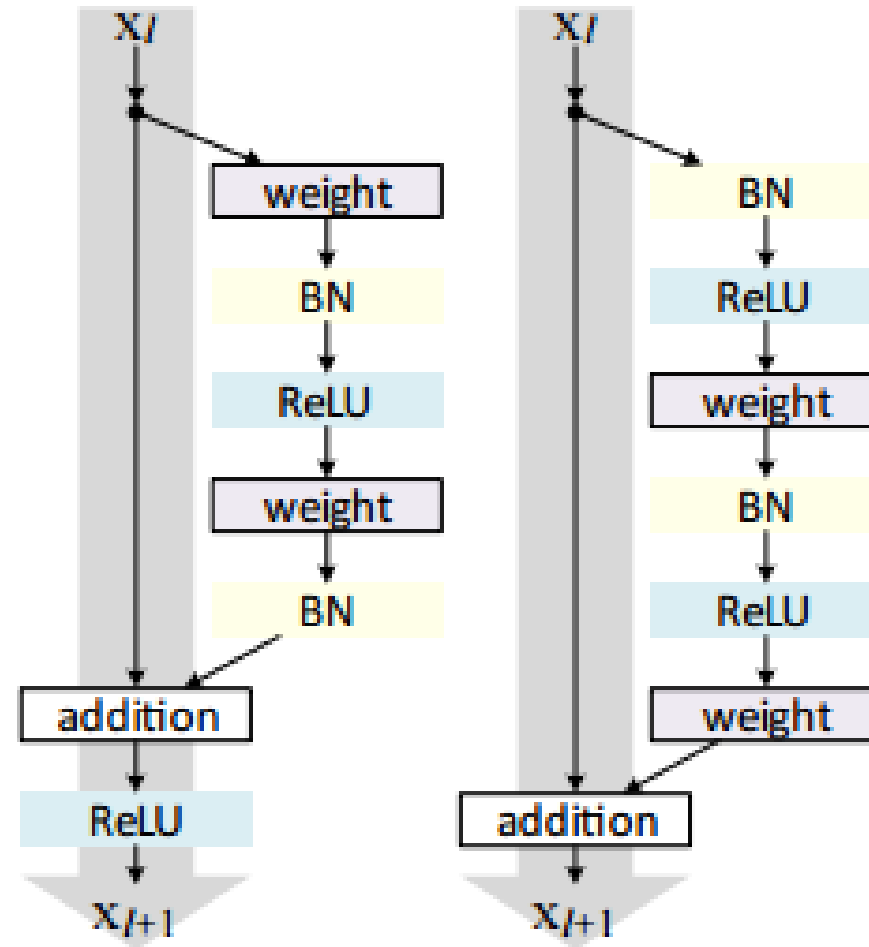
# Deeper Bottleneck Architecture



Residual net

Deeper Bottleneck

# Pre-Activation ResNet



(a) original     (b) proposed

# Pre-Activation ResNet

$$y_l = h(x_l) + \mathcal{F}(x_l, \mathcal{W}_l), \qquad (1)$$

$$x_{l+1} = f(y_l). \qquad (2)$$

**x**: residual unit input
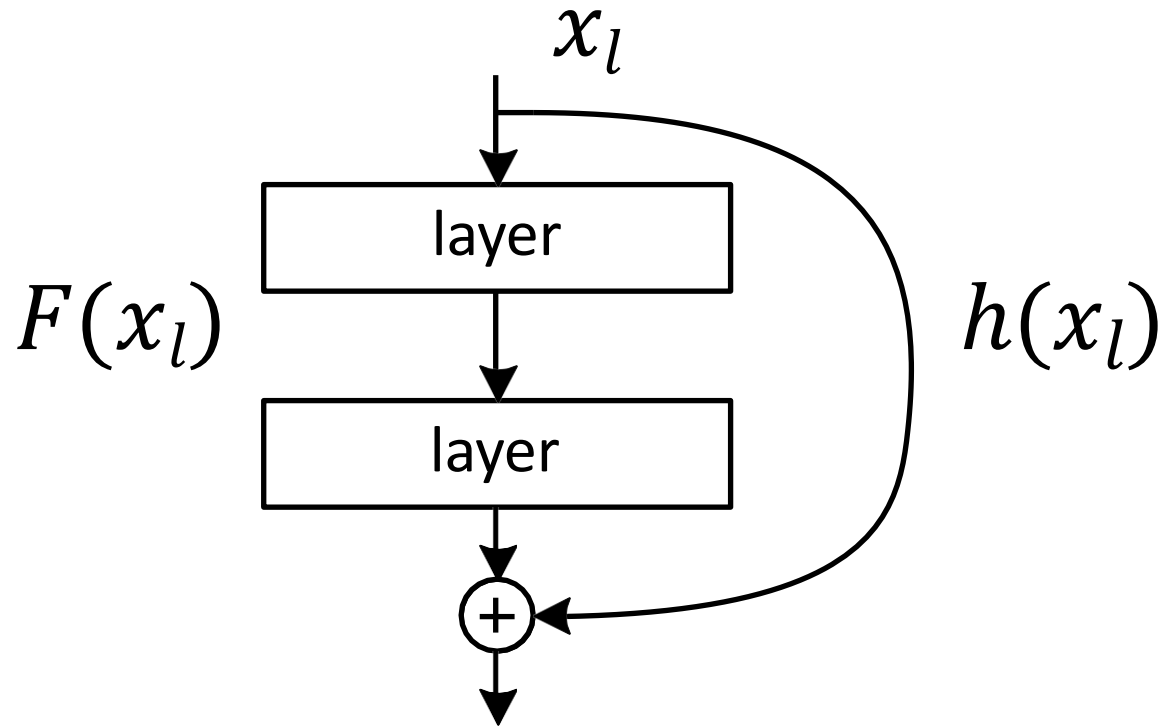
**y**: residual unit output

**l**: number of each layer

**W**: weight

**F**: residual function F(x)

**f**: activation function

**H**: identity function

# Pre-Activation ResNet



$x_l$

$F(x_l)$

layer

layer

$h(x_l)$

x: residual unit input

y: residual unit output

l: number of each layer

W: weight

F: residual function

f: activation function
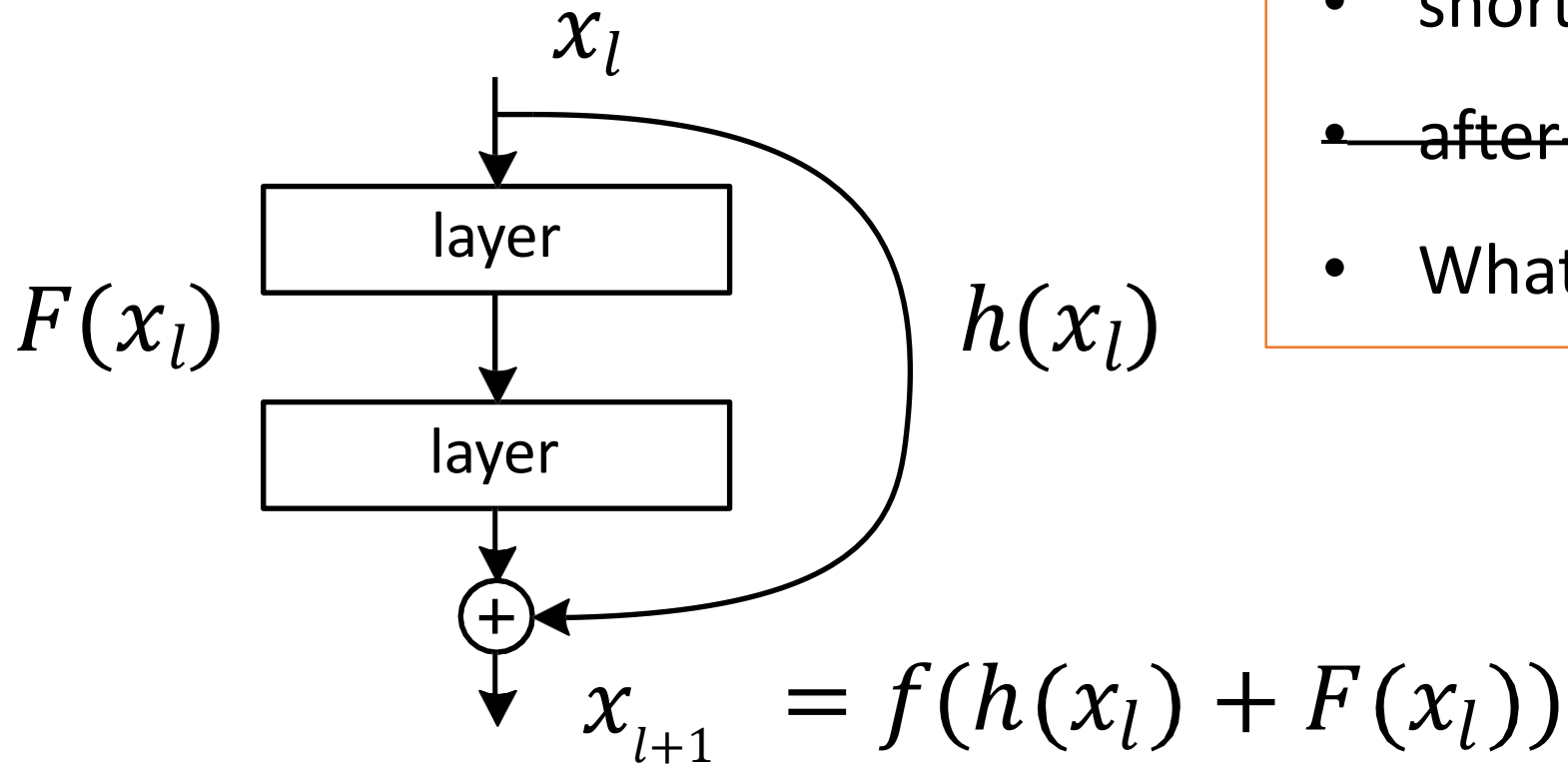
H: identity function

$$y_l = h(\mathbf{x}_l) + \mathcal{F}(\mathbf{x}_l, \mathcal{W}_l), \qquad (1)$$
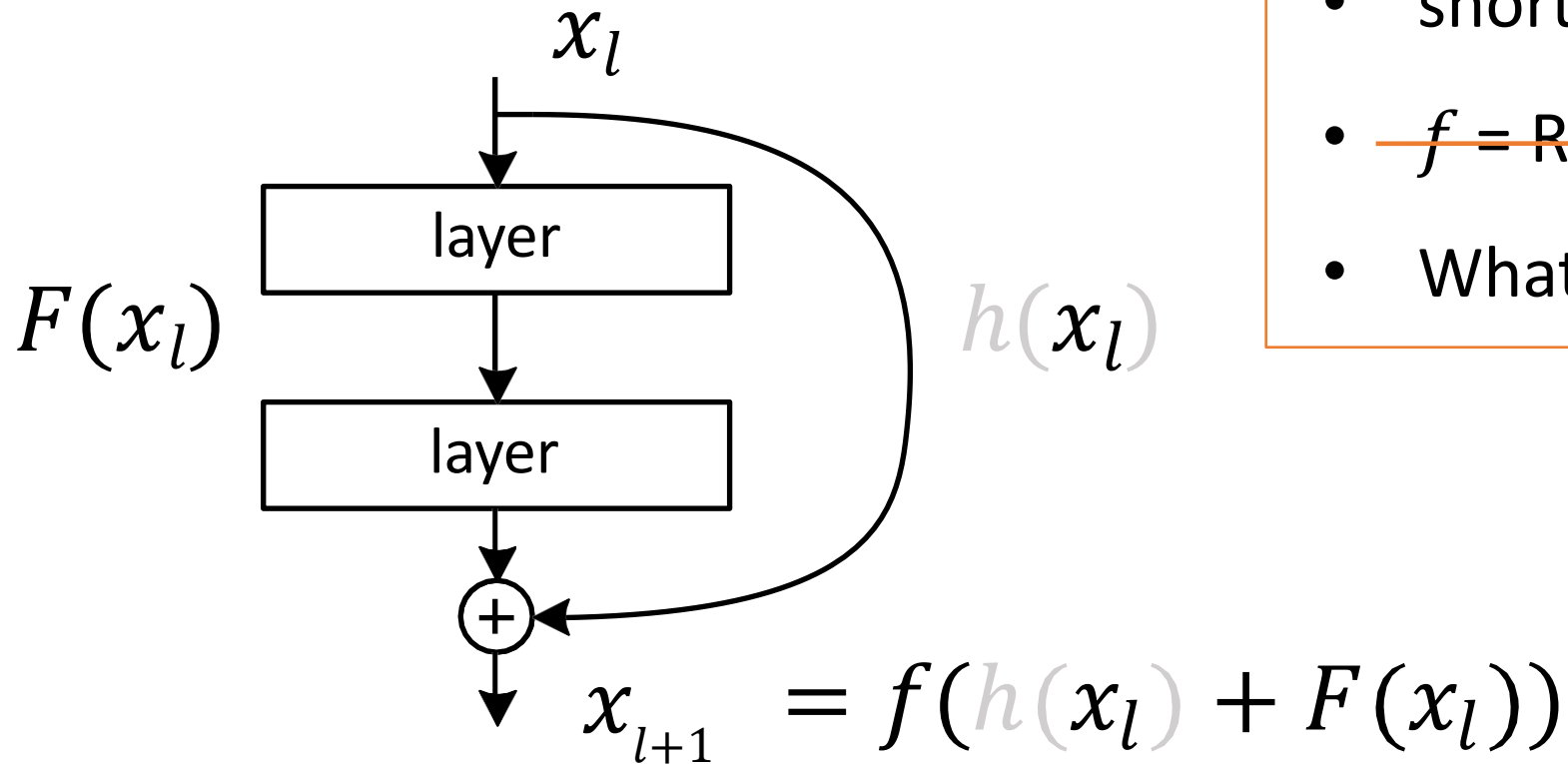
$$\mathbf{x}_{l+1} = f(\mathbf{y}_l). \qquad (2)$$

$$x_{l+1} = f(h(x_l) + F(x_l))$$

# Pre-Activation ResNet



$x_l$

layer

layer

$F(x_l)$

$h(x_l)$

$$x_{l+1} = f(h(x_l) + F(x_l))$$

- shortcut mapping: $h$ = identity

- ~~after-add mapping: $f$ = ReLU~~

- What if $f$ = identity?

# Pre-Activation ResNet



- shortcut mapping: $h$ = identity
- ~~$f$ = ReLU~~
- What if $f$ = identity?

$x_l$

$F(x_l)$

layer

layer

$h(x_l)$

$x_{l+1} = f(h(x_l) + F(x_l))$

# Pre-Activation ResNet

$$x_{l+1} = f(y_l). \tag{2}$$

Change activation function f to <u>identity mapping</u> !

$$x_{l+1} = y_l$$

# Pre-Activation ResNet

대입

$$\mathbf{x}_{l+1} = \mathbf{y}_l \qquad \Rightarrow \qquad \mathbf{y}_l = h(\mathbf{x}_l) + \mathcal{F}(\mathbf{x}_l, \mathcal{W}_l) \qquad (1)$$

$$\mathbf{x}_{l+1} = \mathbf{x}_l + \mathcal{F}(\mathbf{x}_l, \mathcal{W}_l) \qquad (3)$$

# Pre-Activation ResNet

$$x_{l+1} = x_l + F(x_l)$$

$$x_{l+2} = x_{l+1} + F(x_{l+1})$$

# Pre-Activation ResNet

$$x_{l+1} = x_l + F(x_l)$$

$$x_{l+2} = x_{l+1} + F(x_{l+1})$$

$$x_{l+2} = x_l + F(x_l) + F(x_{l+1})$$
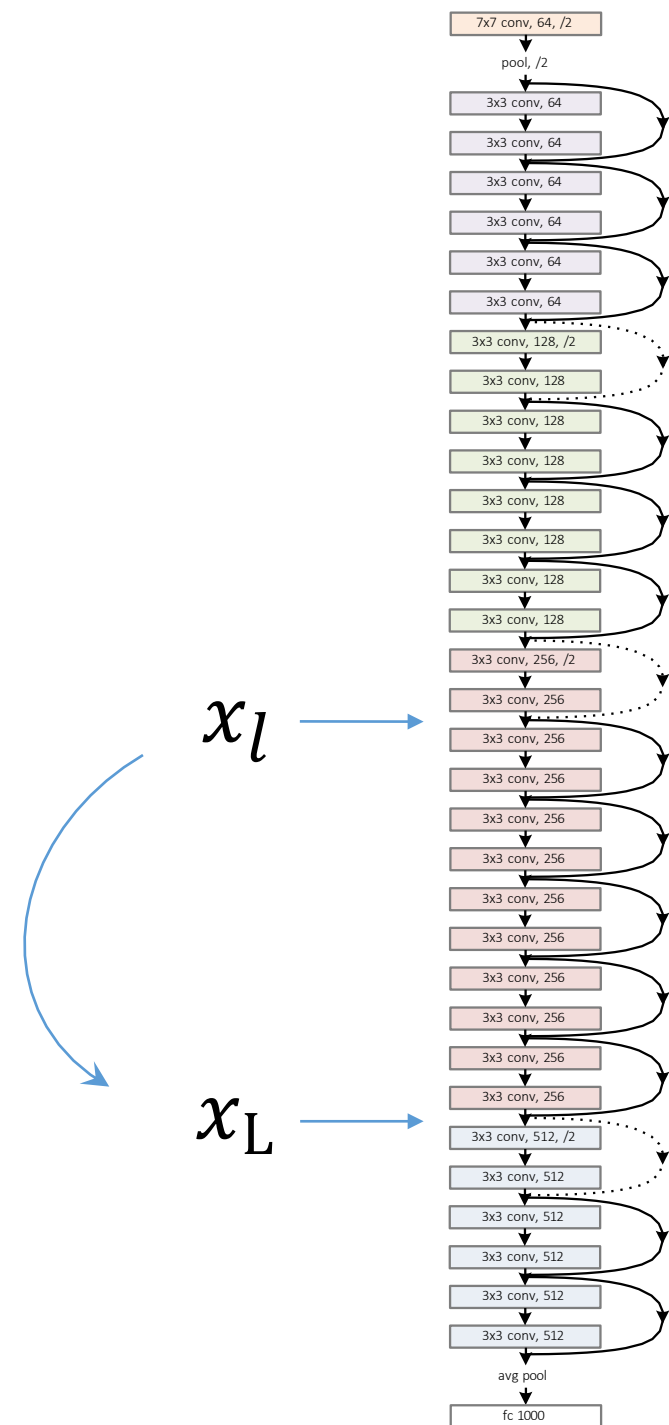
# Pre-Activation ResNet

$$\mathbf{x}_L = \mathbf{x}_l + \sum_{i=l}^{L-1} \mathcal{F}(\mathbf{x}_i, \mathcal{W}_i) \qquad (4)$$

- When Feed forwarding, ResNet can be expressed as the sum of Residual Function F

# Forward Propagation

$$\mathbf{x}_L = \mathbf{x}_l + \sum_{i=l}^{L-1} \mathcal{F}(\mathbf{x}_i, \mathcal{W}_i)$$

- Any $x_l$ is directly forward propagation to any $x_L$, plus residual.

- Any $x_L$ is an additive outcome.

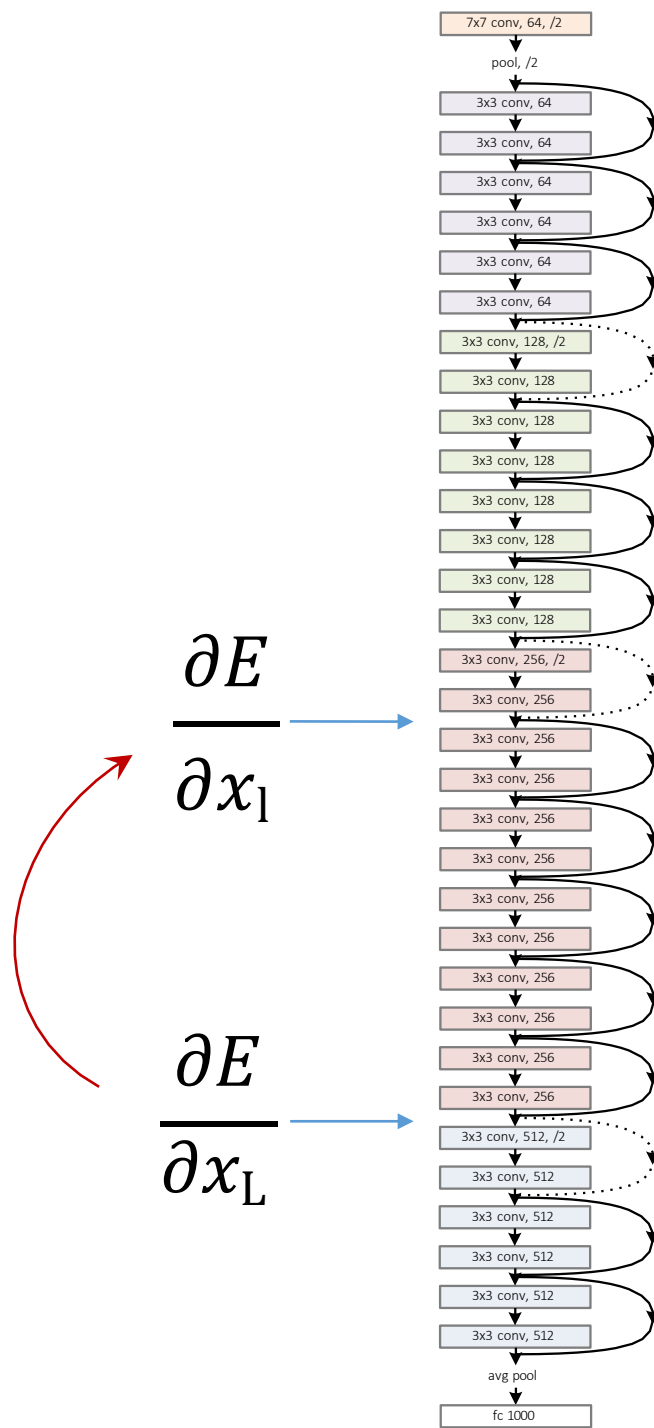  - in contrast to multiplicative: $x_L = \prod_{i=1}^{L-1} W_i\, x_l$

# Back Propagation

$$\mathbf{x}_L = \mathbf{x}_l + \sum_{i=l}^{L-1} \mathcal{F}(\mathbf{x}_i, \mathcal{W}_i)$$

$$\frac{\partial \mathcal{E}}{\partial \mathbf{x}_l} = \frac{\partial \mathcal{E}}{\partial \mathbf{x}_L} \frac{\partial \mathbf{x}_L}{\partial \mathbf{x}_l} = \frac{\partial \mathcal{E}}{\partial \mathbf{x}_L} \left( 1 + \frac{\partial}{\partial \mathbf{x}_l} \sum_{i=l}^{L-1} \mathcal{F}(\mathbf{x}_i, \mathcal{W}_i) \right)$$

$$\frac{\partial E}{\partial x_l}$$

$$\frac{\partial E}{\partial x_L}$$

# Back Propagation

$$\frac{\partial \mathcal{E}}{\partial \mathbf{x}_l} = \frac{\partial \mathcal{E}}{\partial \mathbf{x}_L} \frac{\partial \mathbf{x}_L}{\partial \mathbf{x}_l} = \frac{\partial \mathcal{E}}{\partial \mathbf{x}_L} \left( 1 + \frac{\partial}{\partial \mathbf{x}_l} \sum_{i=l}^{L-1} \mathcal{F}(\mathbf{x}_i, \mathcal{W}_i) \right)$$

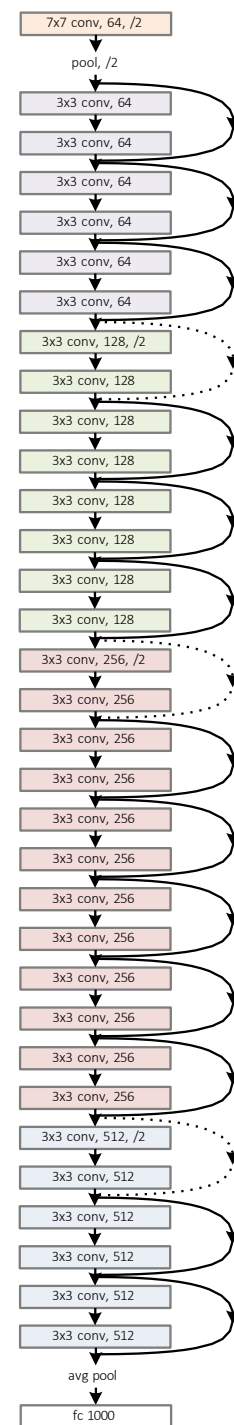- Any $\dfrac{\partial E}{\partial x_\mathrm{L}}$ is directly back propagation to any

  $\dfrac{\partial E}{\partial x_\mathrm{l}}$ plus residual.

- Any $\dfrac{\partial E}{\partial x_\mathrm{l}}$ is additive; unlikely to vanish

  - in contrast to multiplicative: $\dfrac{\partial E}{\partial x_\mathrm{l}} = \Pi_{i=1}^{L-1} W_i \dfrac{\partial E}{\partial x_\mathrm{L}}$

$\dfrac{\partial E}{\partial x_\mathrm{l}}$

$\dfrac{\partial E}{\partial x_\mathrm{L}}$
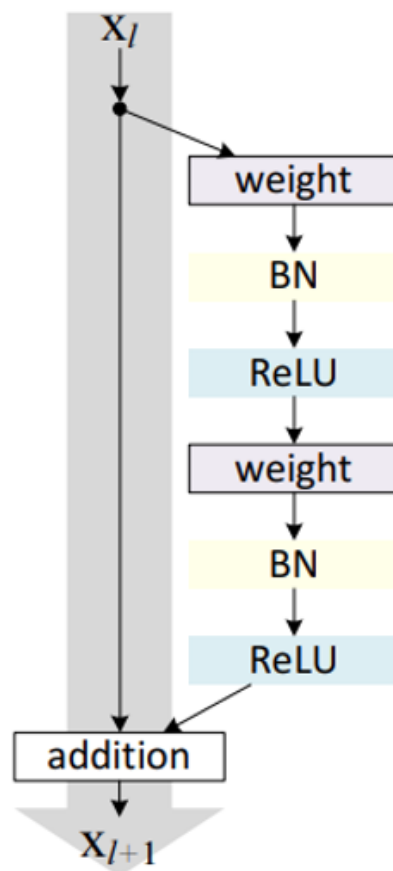
# Pre-Activation ResNet
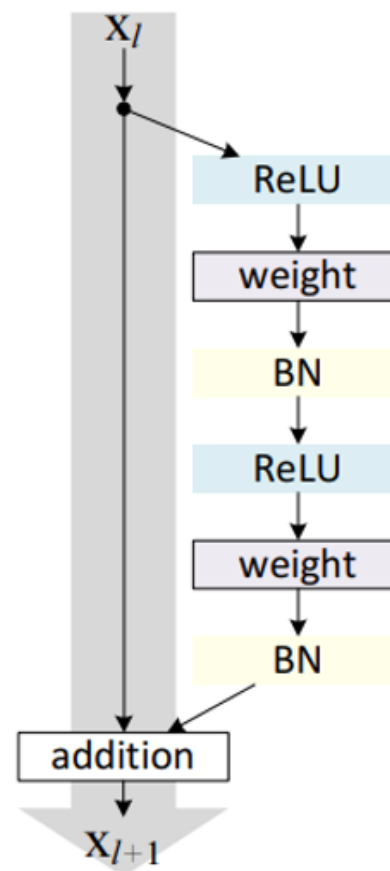


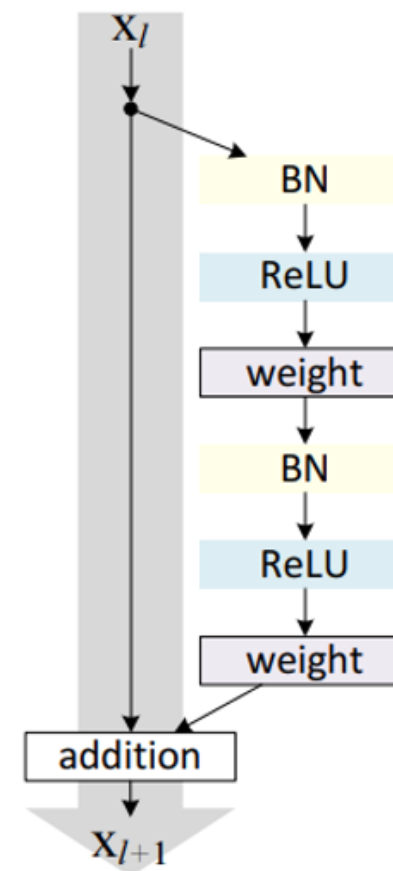(a) original
error : 6.61%

(b) BN after addition
error : 8.17%

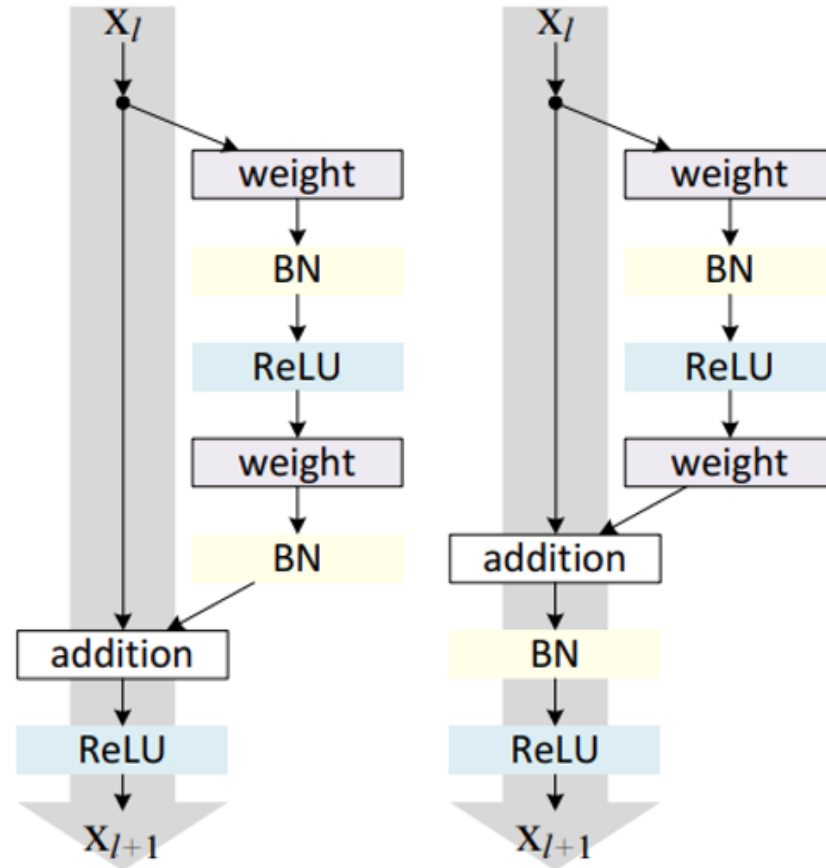(c) ReLU before addition
error : 7.84%

(d) ReLU-only pre-activation
error : 6.71%

(e) **full pre-activation**
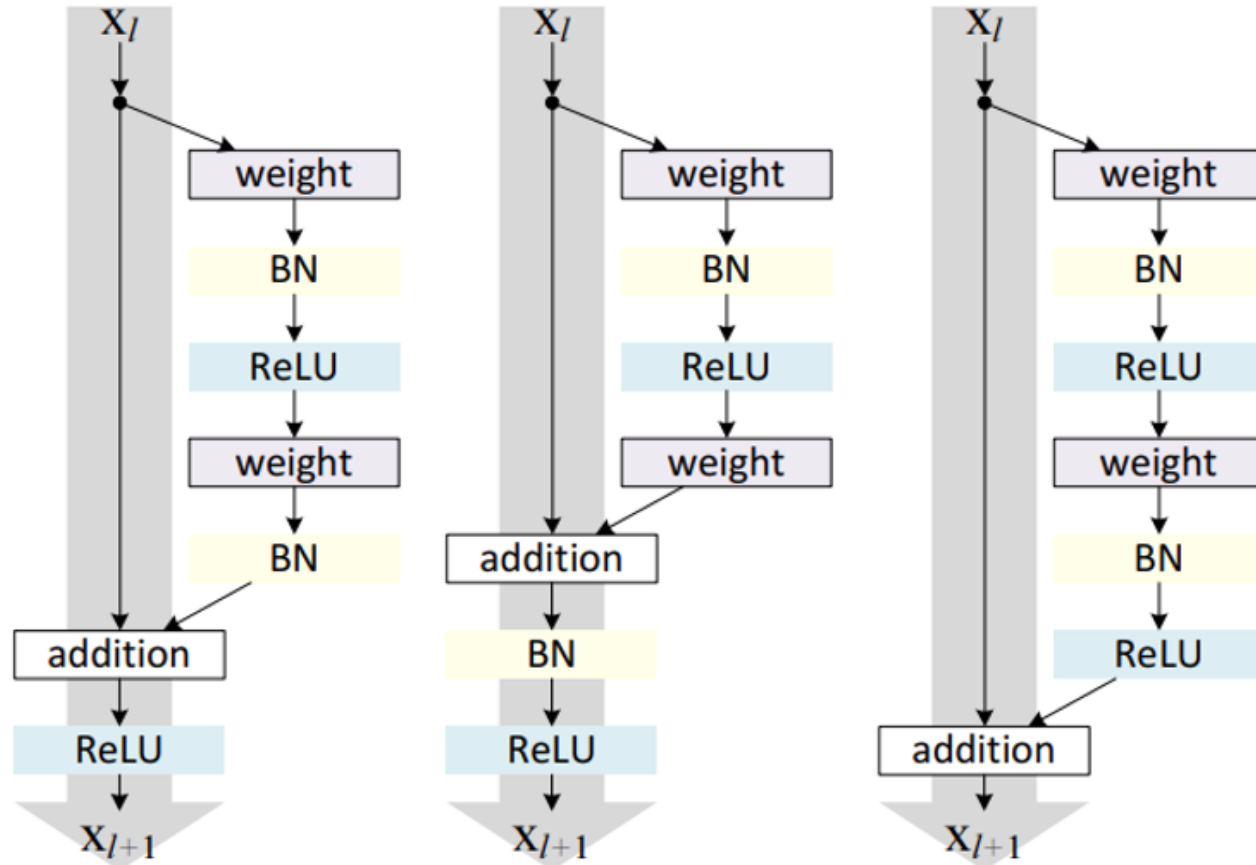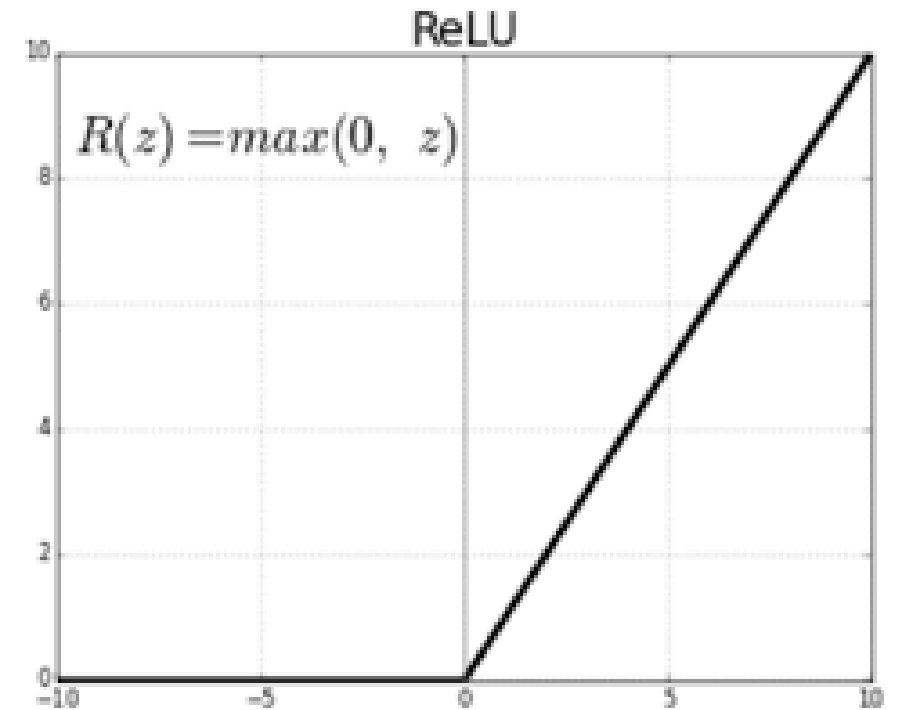**error : 6.37%**

# Pre-Activation ResNet



(a) original

(b) BN after addition

error : 6.61%

error : 8.17%

# Batch Normalization (BN)

- Normalizing input

- BN: normalizing <span style="color:red">each layer</span>, for <span style="color:red">each mini-batch</span>

- Batch: Number of data when the model updates parameters once

- Greatly accelerate training

- Improve regularization

# Pre-Activation ResNet
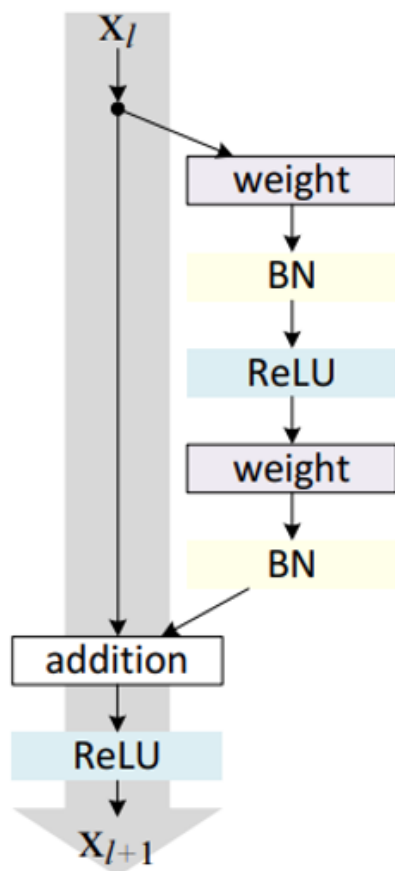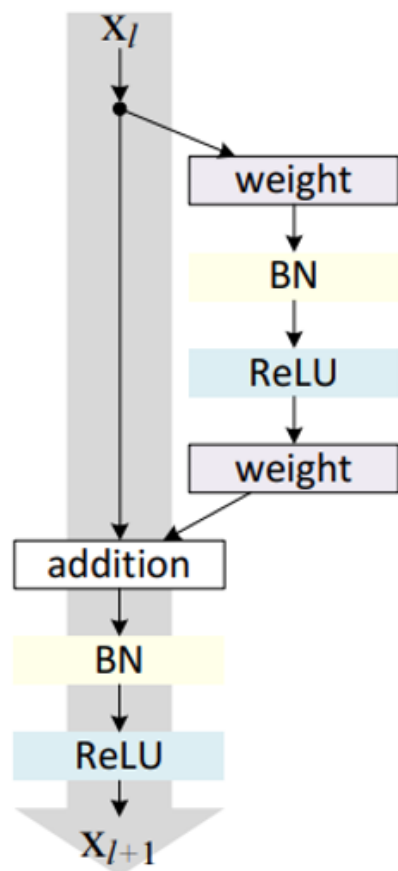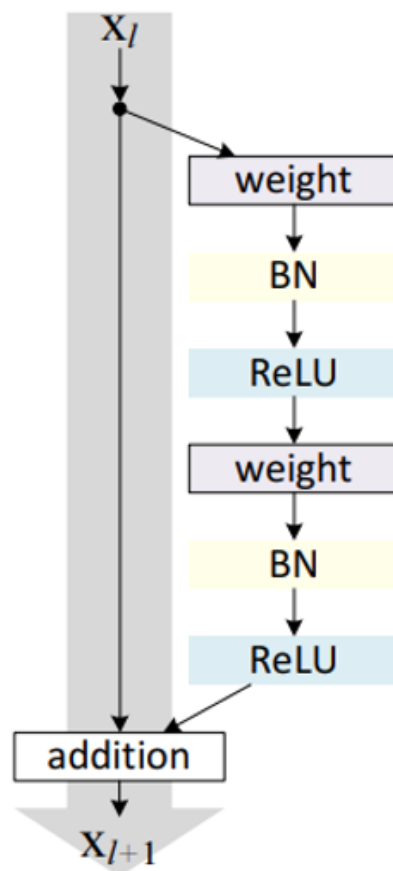


(a) original

error : 6.61%

(b) BN after addition

error : 8.17%

(c) ReLU before addition

error : 7.84%

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

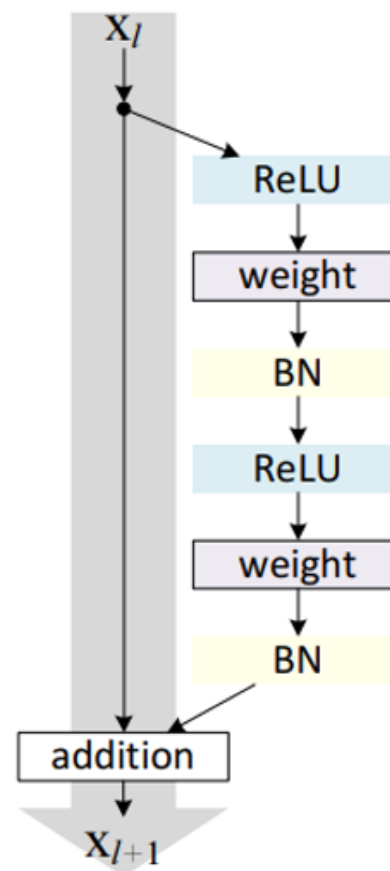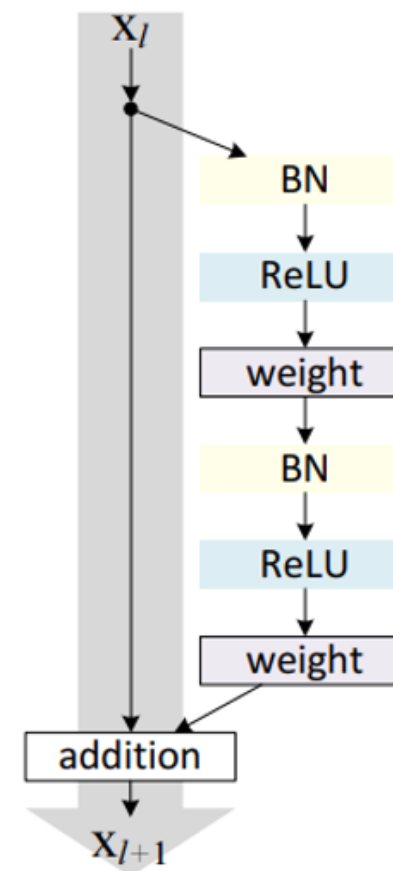# Pre-Activation ResNet



(a) original
error : 6.61%

(b) BN after addition
error : 8.17%

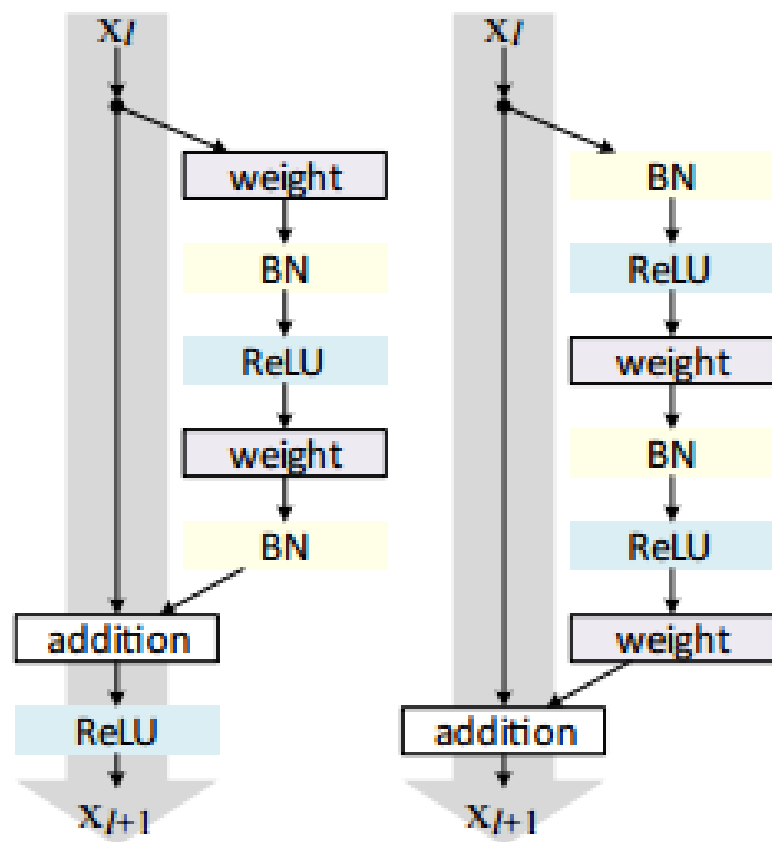(c) ReLU before addition
error : 7.84%

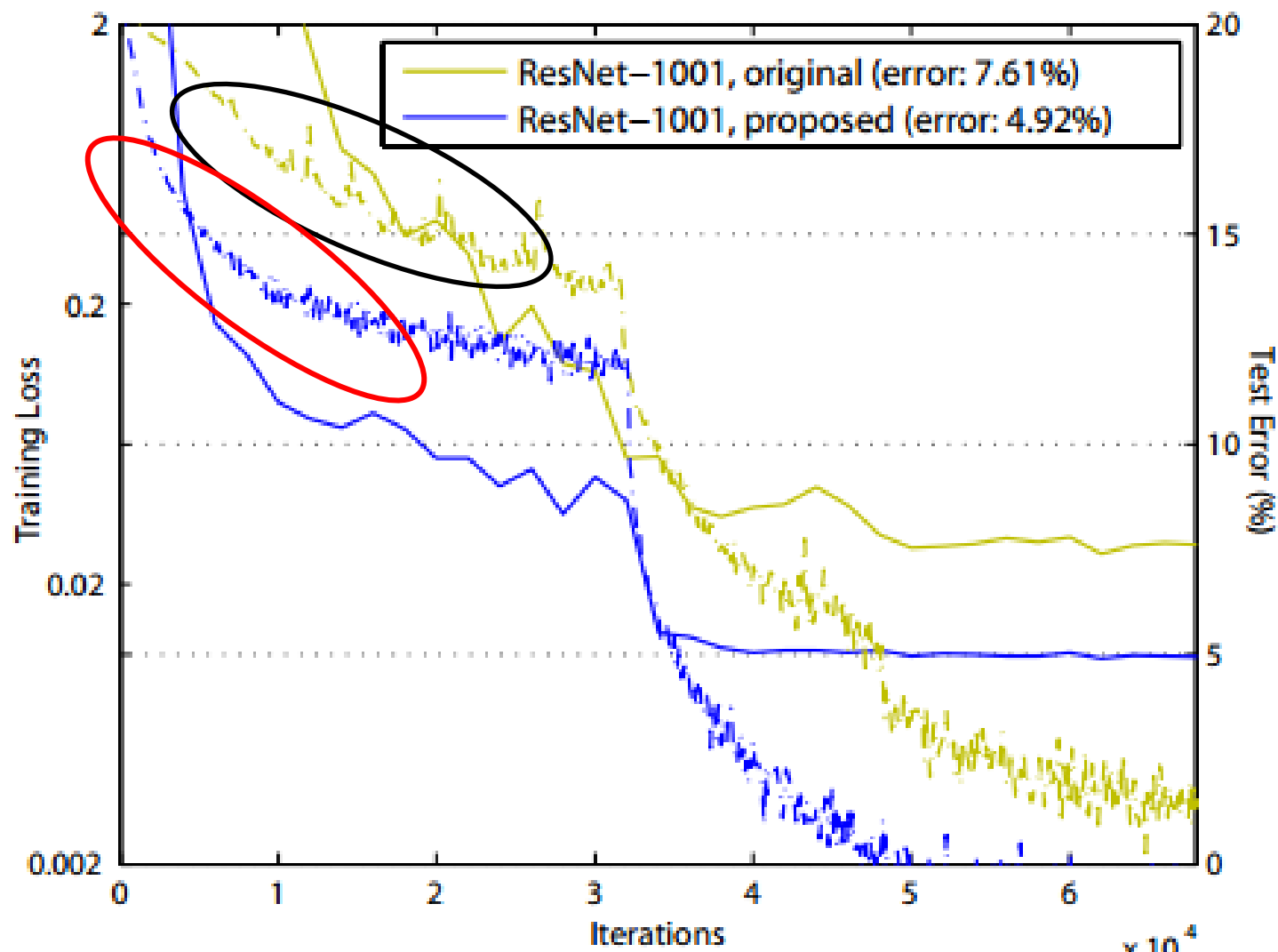(d) ReLU-only pre-activation
error : 6.71%

(e) **full pre-activation**
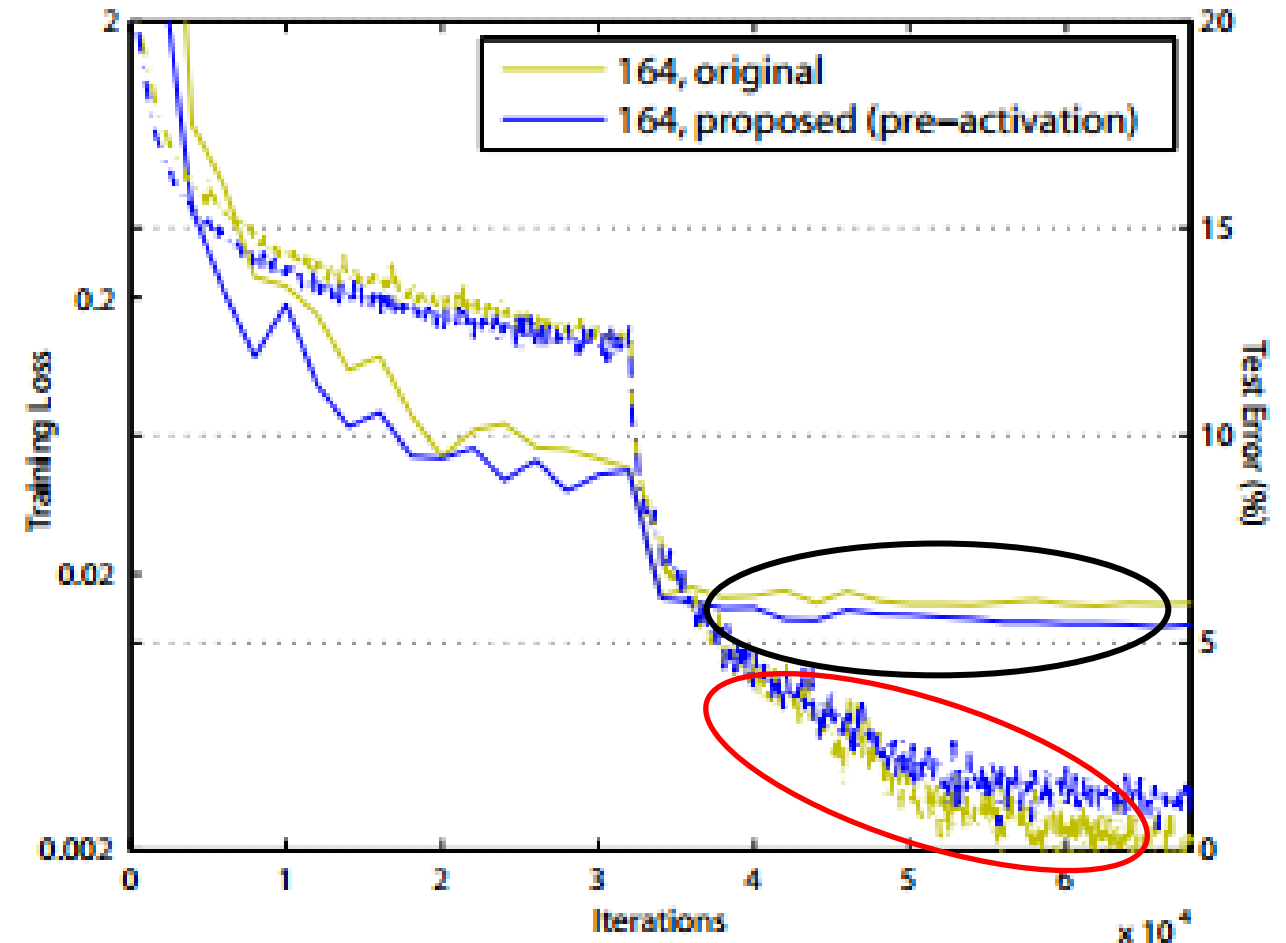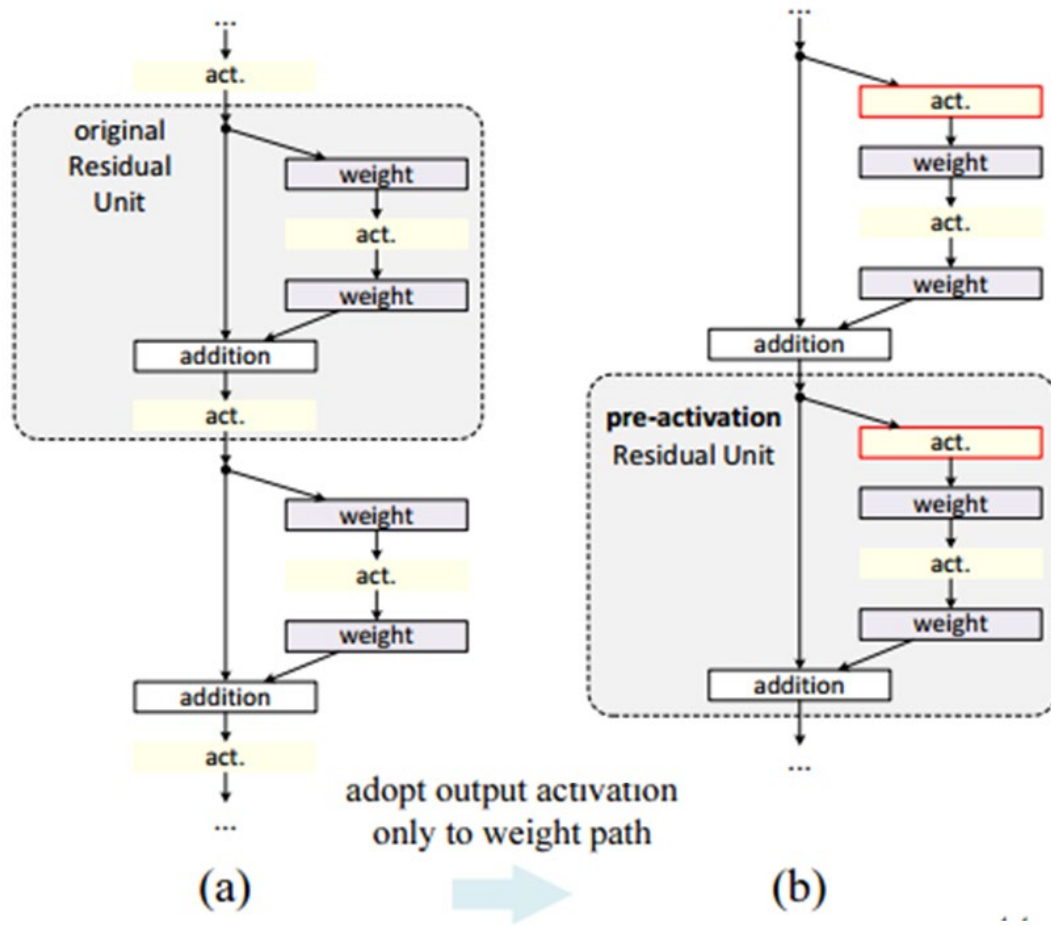**error : 6.37%**

# Pre-Activation ResNet



(a) original  (b) proposed

# Pre-Activation ResNet

# ResNeXt

# ResNeXt



Deeper Bottleneck

ResNeXt

# ResNeXt – Grouped Convolution

# ResNeXt – Cardinality, Width



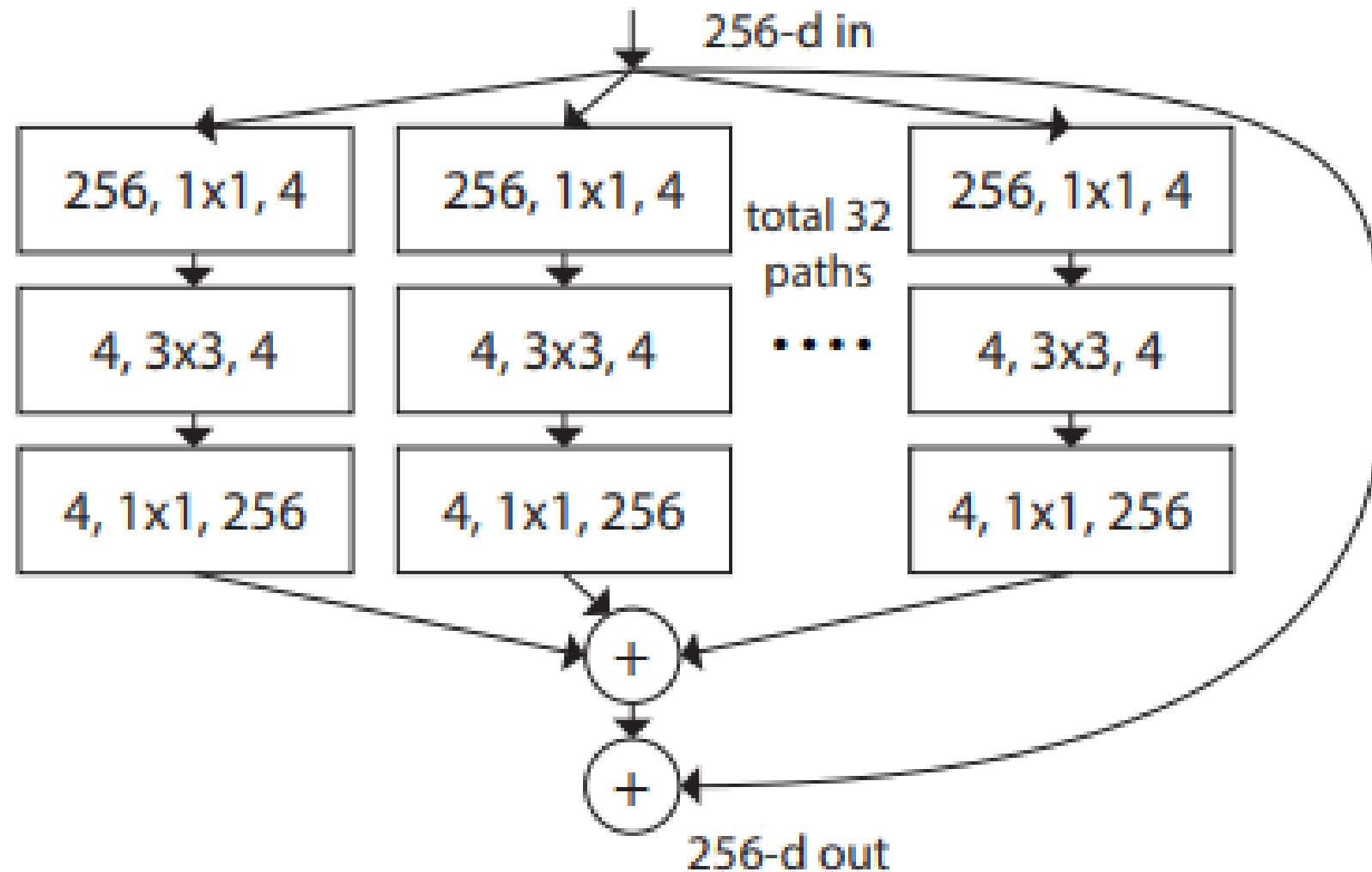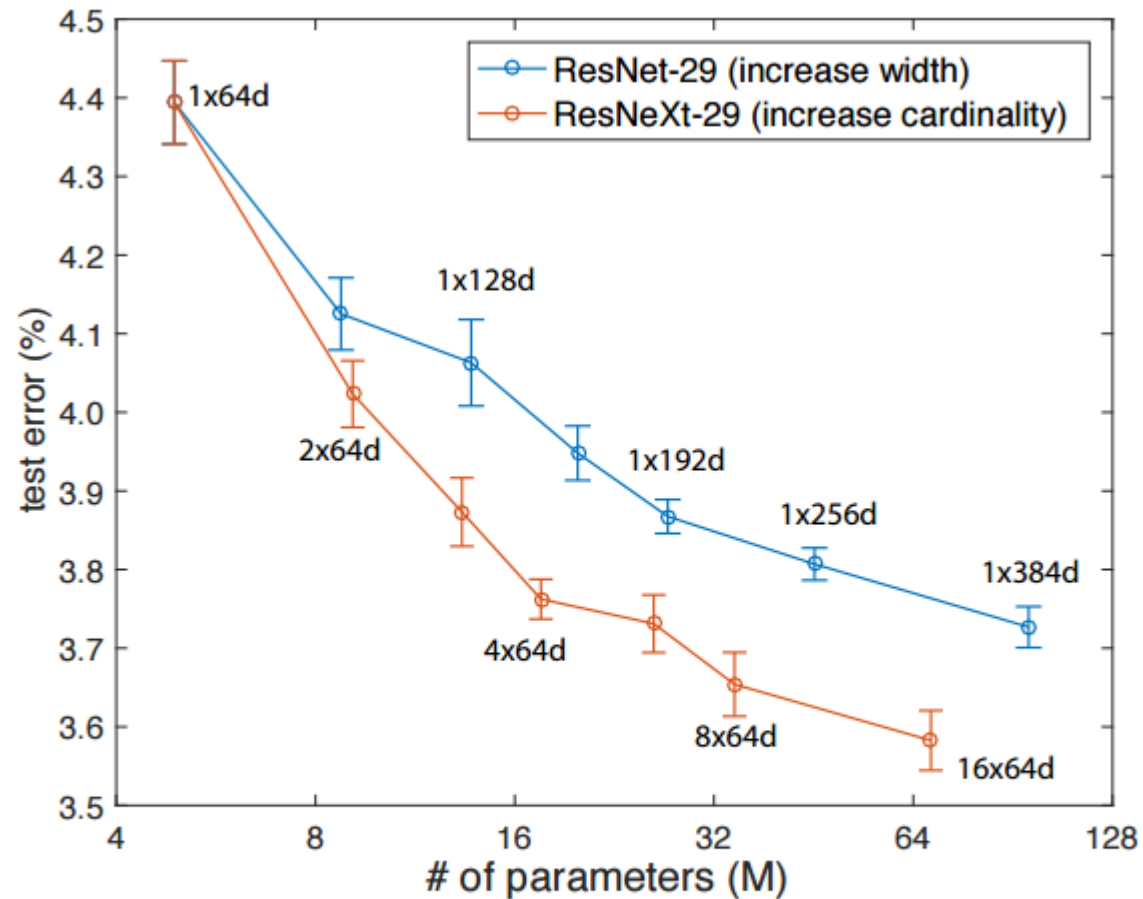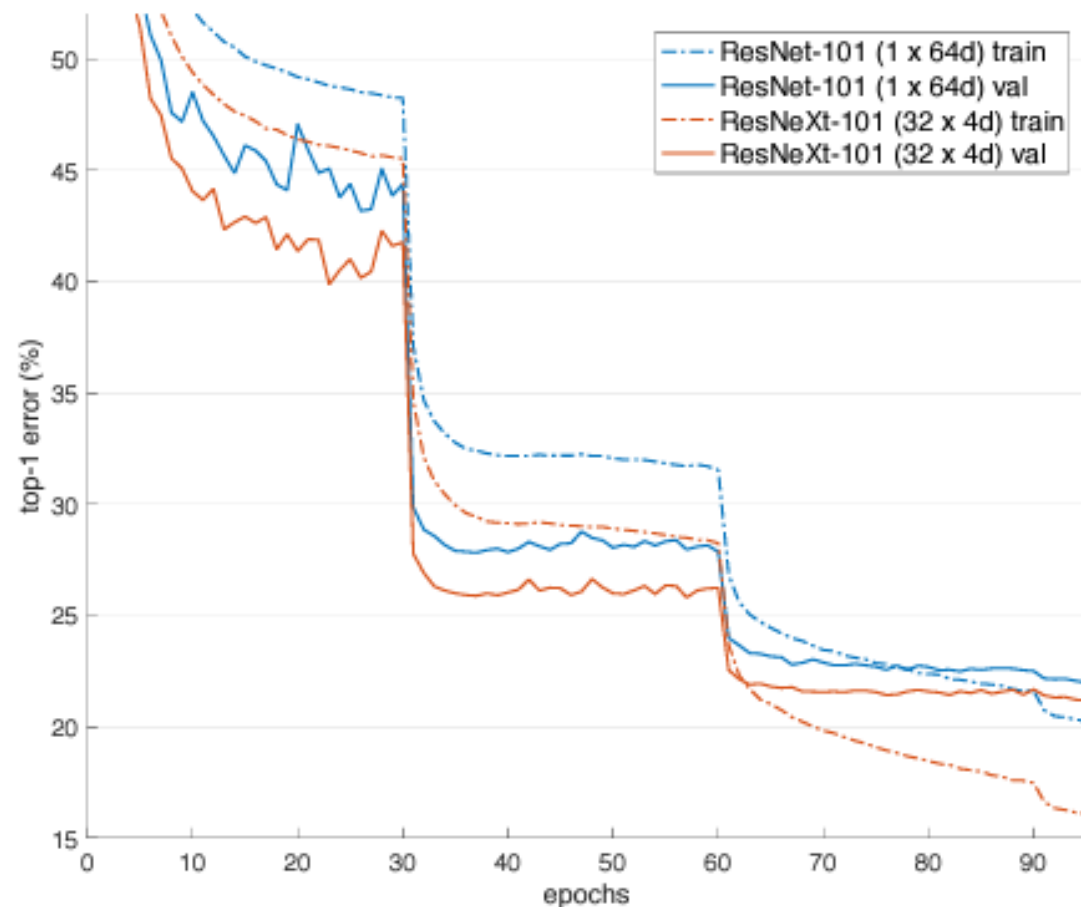Figure 7. Test error *vs.* model size on CIFAR-10. The results are computed with 10 runs, shown with standard error bars. The labels show the settings of the templates.

- Cardinality: the number of groups to divide the total number of channels

- Width: the number of channels in one group

- Cardinality  > Width > Depth

# ResNeXt – Cardinality, Width



| | setting | top-1 err (%) | top-5 err (%) |
|---|---|---|---|
| *1× complexity references:* | | | |
| ResNet-101 | $1 \times 64$d | 22.0 | 6.0 |
| ResNeXt-101 | $32 \times 4$d | 21.2 | 5.6 |
| *2× complexity models follow:* | | | |
| ResNet-**200** [15] | $1 \times 64$d | 21.7 | 5.8 |
| ResNet-101, wider | $1 \times$ **100**d | 21.3 | 5.7 |
| ResNeXt-101 | **2** $\times 64$d | 20.7 | 5.5 |
| ResNeXt-101 | **64** $\times 4$d | **20.4** | **5.3** |

Table 4. Comparisons on ImageNet-1K when the number of FLOPs is increased to 2× of ResNet-101's. The error rate is evaluated on the single crop of 224×224 pixels. The highlighted factors are the factors that increase complexity.

# 감사합니다