# ResNet

## Deep Residual Learning for Image Recognition

박성진

# 1. Introduction

- **Is deeper network always better?**

    - **Deep networks naturally integrate low/mid/high level features**

    - **The levels of features can be enriched by the number of stacked layers(depth)**

    - **Recent evidence reveals that network depth is crucial importance**

# 1. Introduction

- Is deeper network always better?

  - **Deep** networks naturally integrate low/mid/high level features

  - The levels of features can be enriched by the number of stacked layers(**depth**)

  - Recent evidence reveals that network depth is crucial importance

| AlexNet | VggNet | GoogLeNet |
|---------|--------|-----------|
| 16.4% | 7.3% | 6.7% |
| 8 layers | 19 layers | 22 layers |
| 2012 | 2014 | 2014 |

# 1. Introduction

- **Is deeper network always better?**

    - **Deep** networks naturally integrate low/mid/high level features
    - The levels of features can be enriched by the number of stacked layers(**depth**)
    - Recent evidence reveals that **network depth** is crucial importance

| AlexNet | VggNet | GoogLeNet |
|---|---|---|
| 16.4% | 7.3% | 6.7% |
| 8 layers | 19 layers | 22 layers |
| 2012 | 2014 | 2014 |

# 1. Introduction

- **Is deeper network always better?**

  - **Deep** networks naturally integrate low/mid/high level features
  - The levels of features can be enriched by the number of stacked layers(**depth**)

  - Recent evidence reveals that network depth is crucial importance

| AlexNet | VggNet | GoogLeNet |
|---------|--------|-----------|
| 16.4% | 7.3% | 6.7% |
| 8 layers | 19 layers | 22 layers |
| 2012 | 2014 | 2014 |

# 1. Introduction

- **Is deeper network always better?**

  - **Deep** networks naturally integrate low/mid/high level features

  - The levels of features can be enriched by the number of stacked layers(**depth**)

  - Recent evidence reveals that network depth is crucial importance

| AlexNet | VggNet | GoogLeNet |
|---------|--------|-----------|
| 16.4% | 7.3% | 6.7% |
| **8 layers** | **19 layers** | **22 layers** |
| 2012 | 2014 | 2014 |

# 1. Introduction

- **Is deeper network always better?**

ResNet

3.57%
**152 layers**
2015

# 1. Introduction

- **Is deeper network always better?**

  - **Deep** networks naturally integrate low/mid/high level features

  - The levels of features can be enriched by the number of stacked layers(**depth**)

  - Recent evidence reveals that network depth is crucial importance

# Yes !!!
# Deeper Network is better!!!

# 1. Introduction

- **Is learning better networks as <span style="color:red">easy</span> as stacking more layers?**

  - **Problem of vanishing/exploding gradients**

    - **Better initialization methods**
    - **Better normalization**
    - **Better Activation function**

  - **Number of parameters**

    - **Deeper bottleneck architectures (1x1 Conv)**
    - **2 times 3x3 Conv**

# 1. Introduction

- **Is learning better networks as <span style="color:red">easy</span> as stacking more layers?**

  - **Problem of vanishing/exploding gradients**

    - **Better initialization methods**
    - **Better normalization**
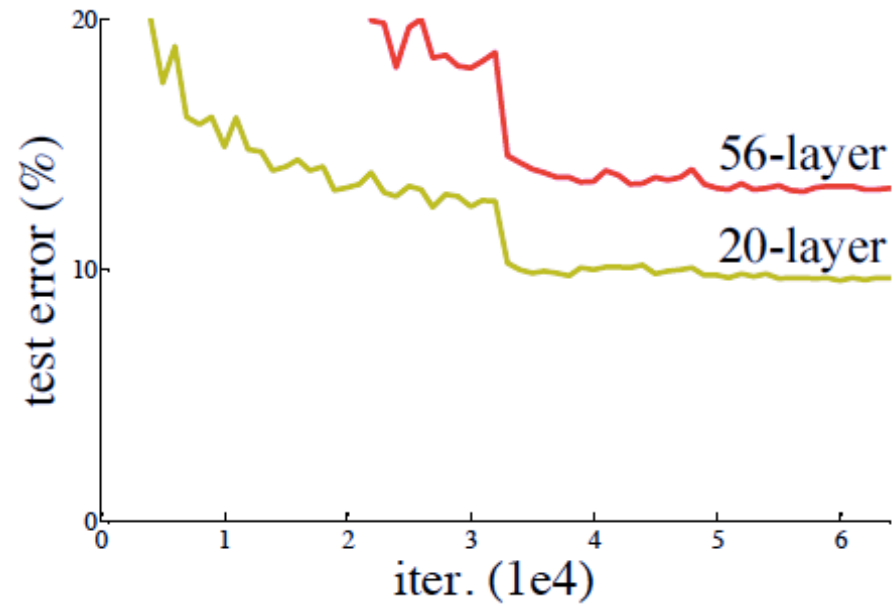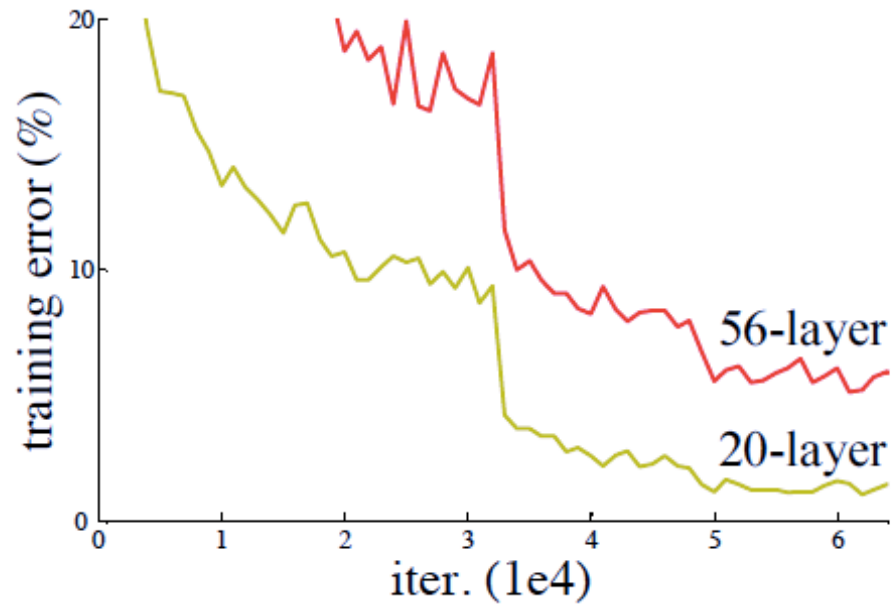    - **Better Activation function**

  - **Number of parameters**

    - **Deeper bottleneck architectures (1x1 Conv)**
    - **2 times 3x3 Conv**

- **Any other problems?**

  - **Overfitting**

# 1. Introduction

- **Is learning better networks as <span style="color:darkred">easy</span> as stacking more layers?**

  - Problem of vanishing/exploding gradients

    - Better initialization methods

# No Overfitting But <span style="color:orange">Degradation</span> !!!

    - Better Activation function(ReLU)

- **Any other problems?**

  - Overfitting

# 1. Introduction

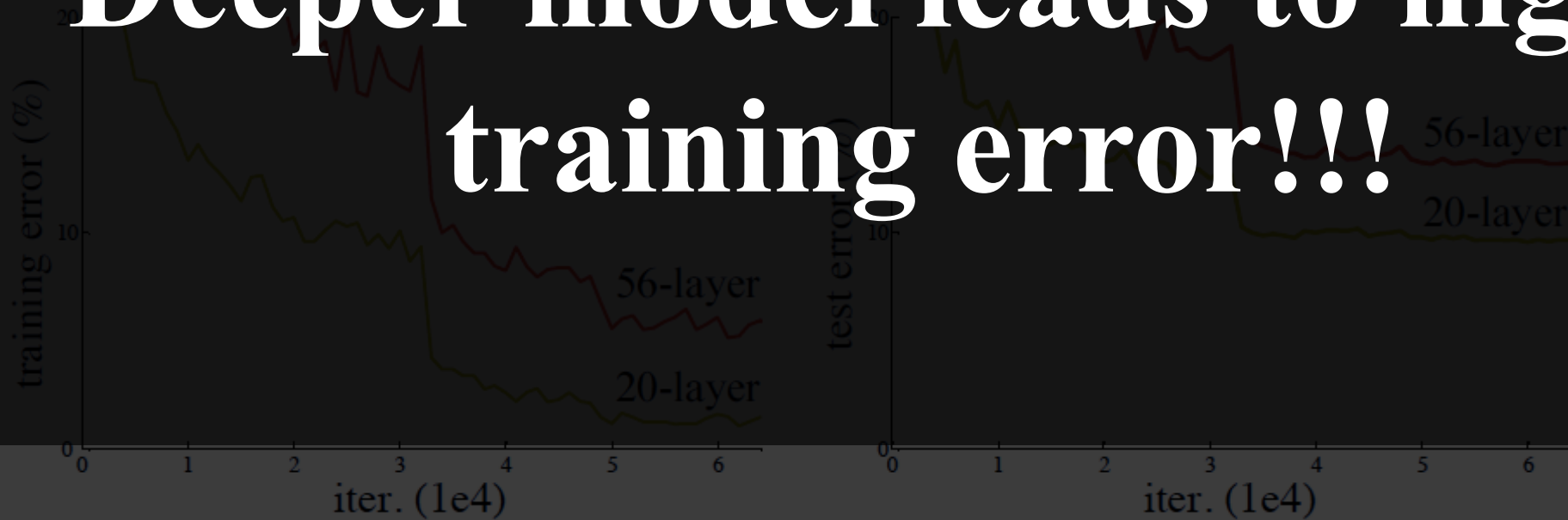- **Degradation Problem**

  - **More depth but lower performance**

# 1. Introduction

- **Degradation Problem**

  - **More depth but lower performance**



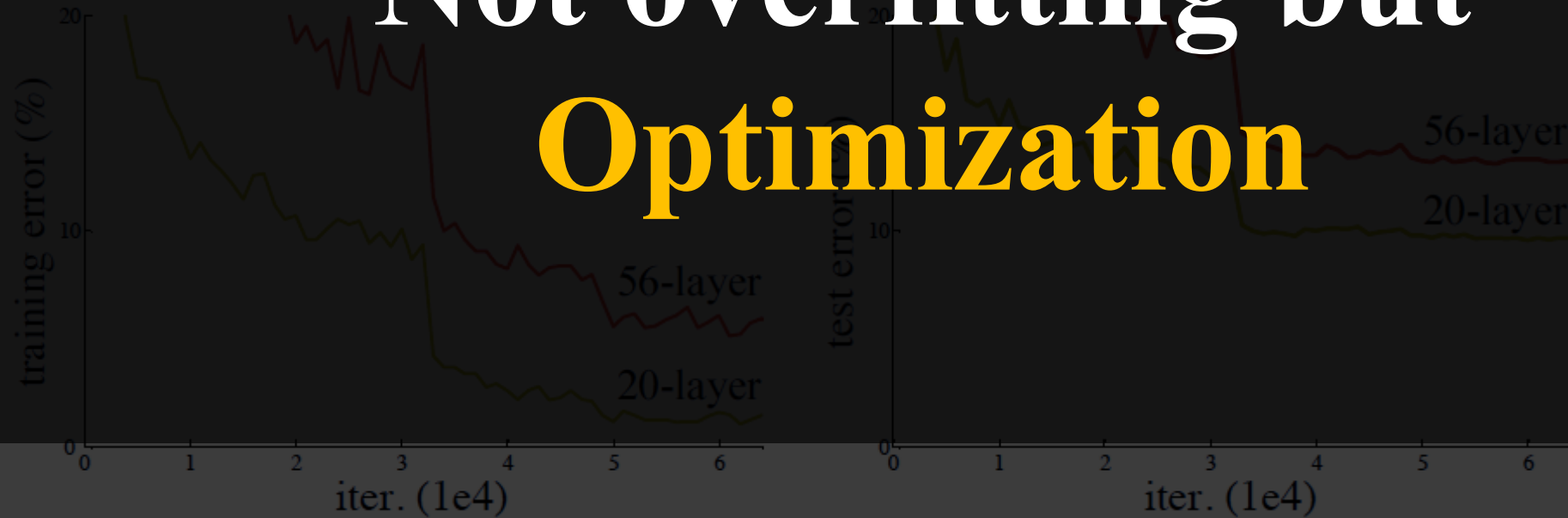# Deeper model leads to higher training error!!!

# 2. Model Architecture

- **Hypothesis : the problem is an <span style="color:red">optimization problem</span>  deeper models are harder to optimize.**

# 1. Introduction

- **Degradation Problem**

    - **More depth but lower performance**



# Not overfitting but
# Optimization

# 2. Model Architecture
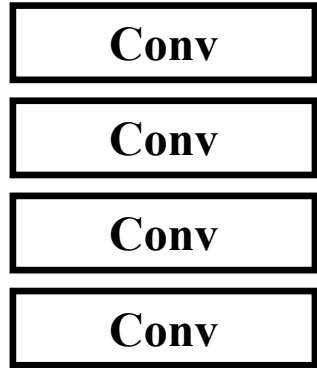
- **Hypothesis : the problem is an <span style="color:red">optimization problem</span>  deeper models are harder to optimize.**

- **The deeper model should be able to perform <span style="color:red">at least</span> as well as the shallower model.**

- **A solution by construction is copying the learned layers from the shallower model ands setting additional layers to identity mapping.**
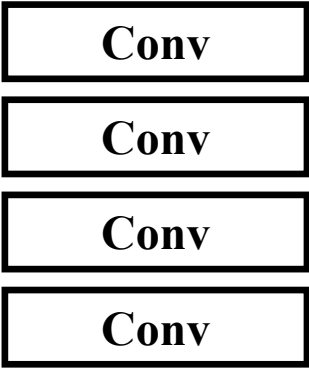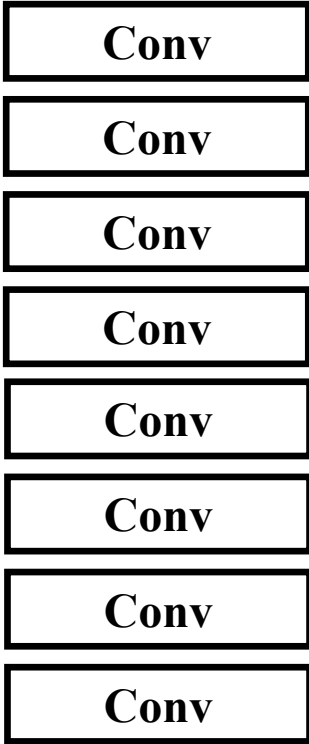
# 2. Model Architecture

| Conv |
|------|

| Conv |
|------|

| Conv |
|------|

| Conv |
|------|

**Trained**

**Tested**

**Acc X%**

# 2. Model Architecture

| Conv |
| --- |
| Conv |
| Conv |
| Conv |

| Trained |
| --- |
| Tested |

↓

**Acc X%**

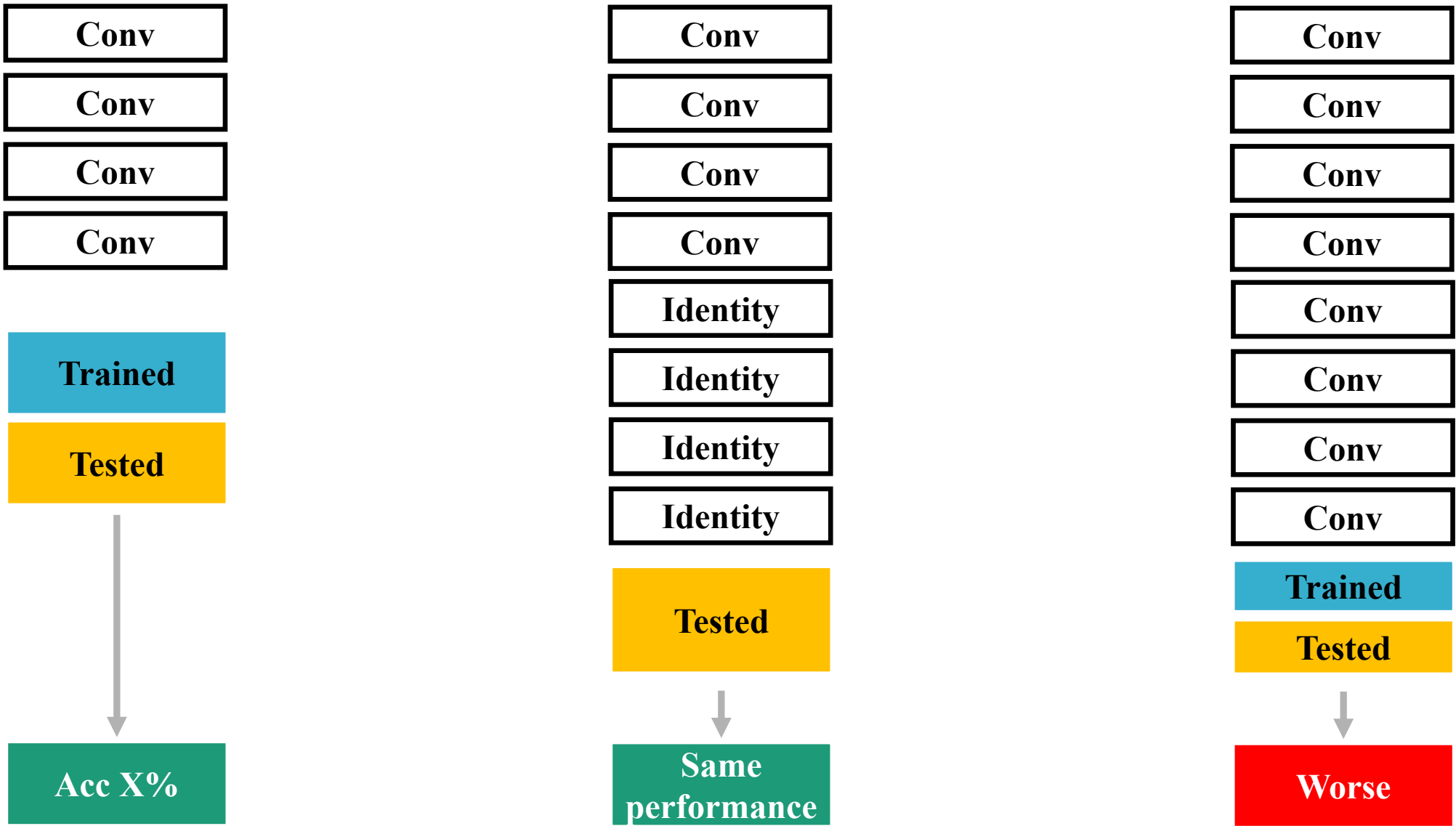| Conv |
| --- |
| Conv |
| Conv |
| Conv |
| Conv |
| Conv |
| Conv |
| Conv |

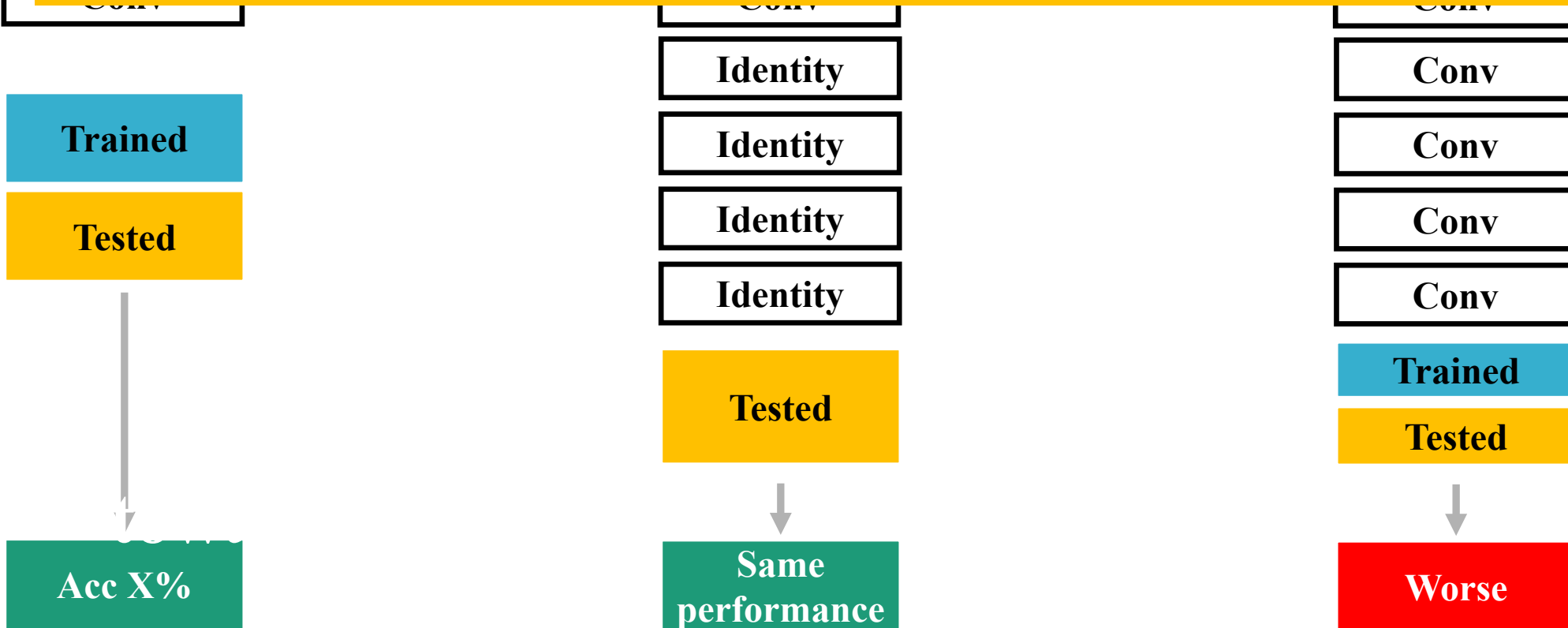| Trained |
| --- |
| Tested |

↓
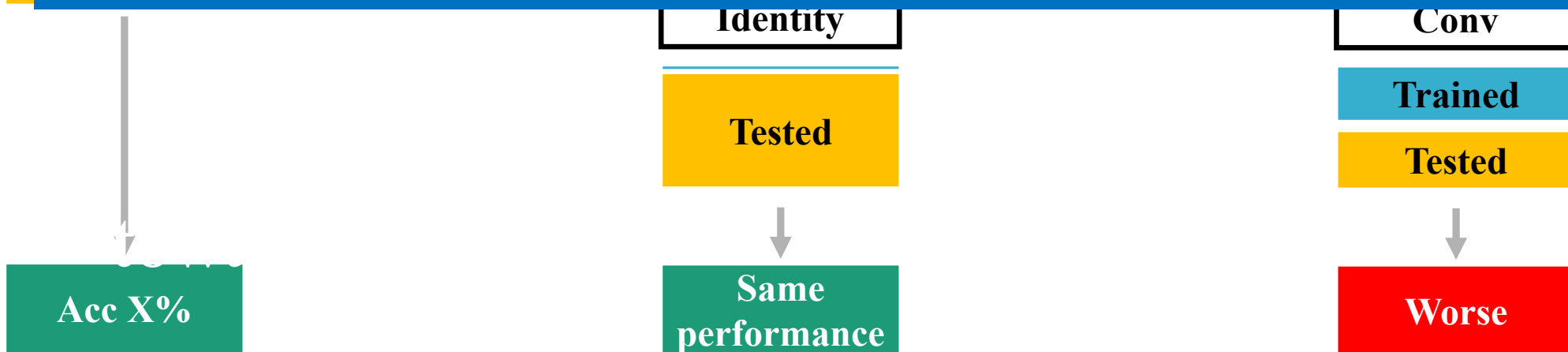
**Worse**

# 2. Model Architecture

Our Current solvers on hand are **unable to find solutions** that are comparably good or better than the **constructed solution**(or unable to do so in feasible time)

Conv

**Trained**

**Tested**

Acc X%

Conv

Identity

Identity

Identity

Identity

**Tested**

Same performance

Conv

Conv

Conv

Conv

Conv

**Trained**

**Tested**

Worse

**2. N**

Our Current solvers on hand are **unable to find solutions** that are comparably good or better than the **constructed solution**(or unable to do so in feasible time)

Solvers might have **difficulties in approximating Identity mappings** by multiple nonlinear layers

Identity

Conv

Trained

Tested

Tested

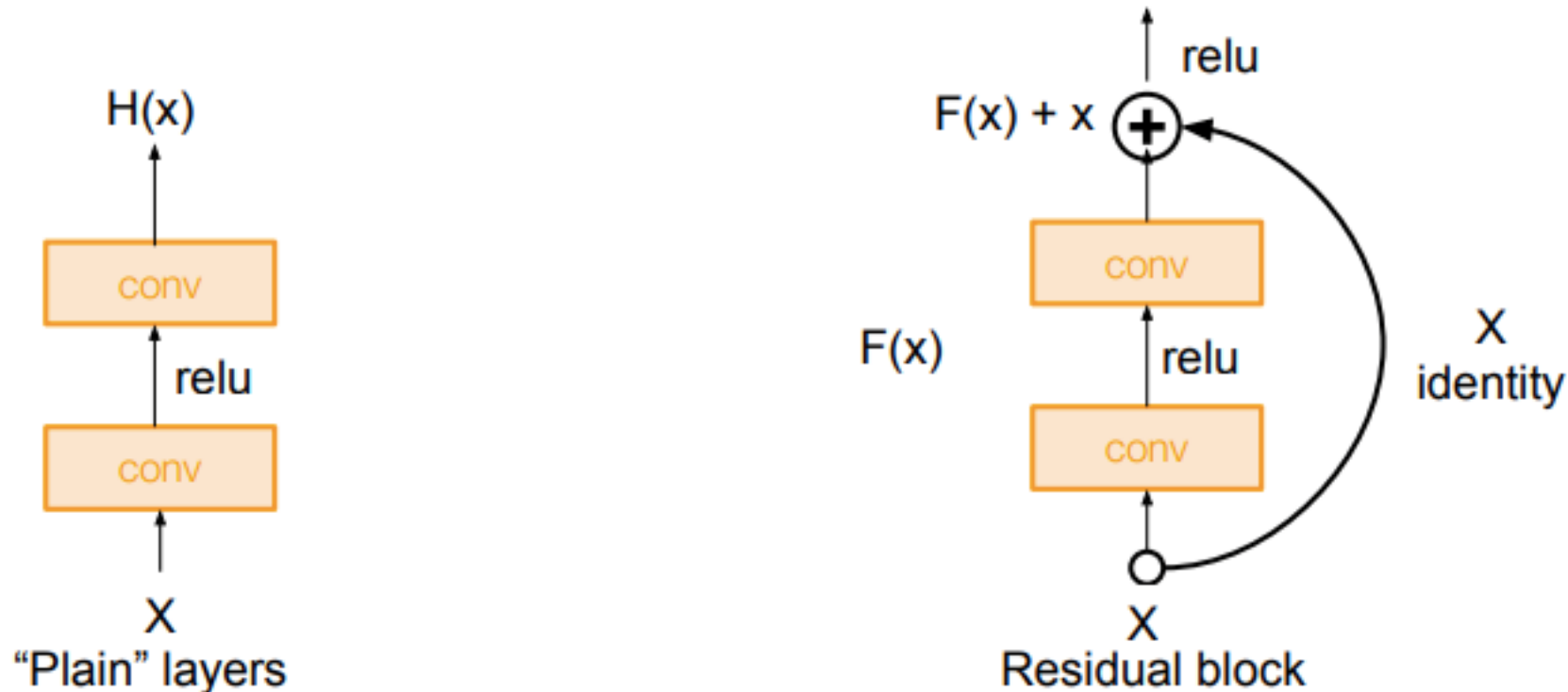Acc X%

Same performance

Worse

**2. N**

**Our Current solvers on hand are unable to find solutions that are comparably good or better than the constructed solution(or unable to do so in feasible time)**

**Solvers might have difficulties in approximating Identity mappings by multiple nonlinear layers**

**Add explicit identity connections and solvers may Drive the weights of the multiple nonlinear layers toward zero**

performance

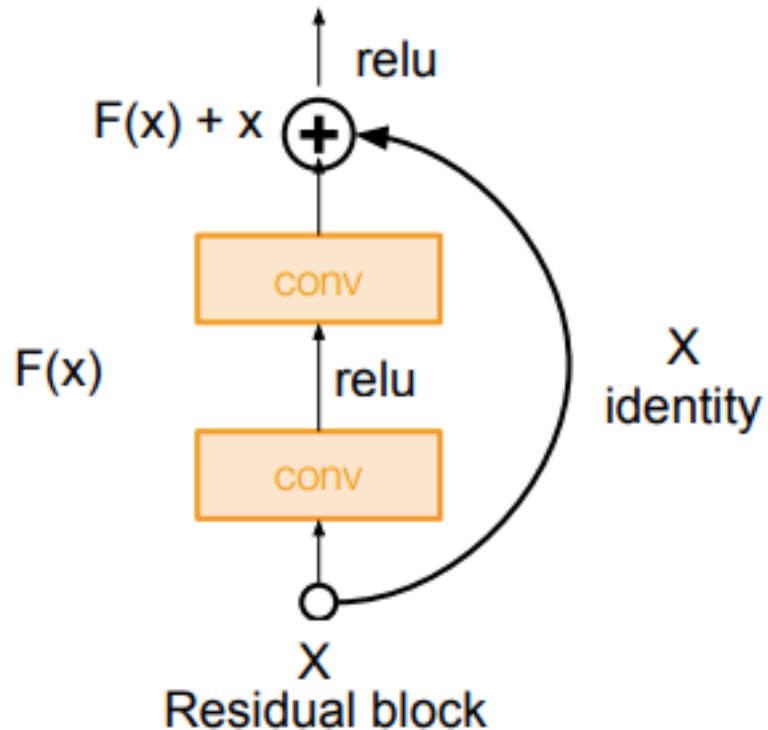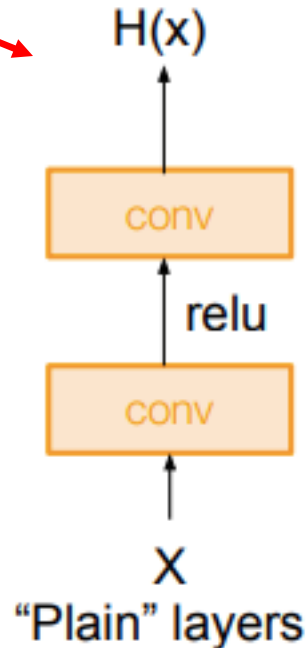# 2. Model Architecture

- **Solution : Use network layers to fit a residual mapping instead of**

  **<span style="color:red">directly trying to fit a desired underlying mapping</span>**

# 2. Model Architecture

- **Solution : Use network layers to fit a residual mapping instead of**

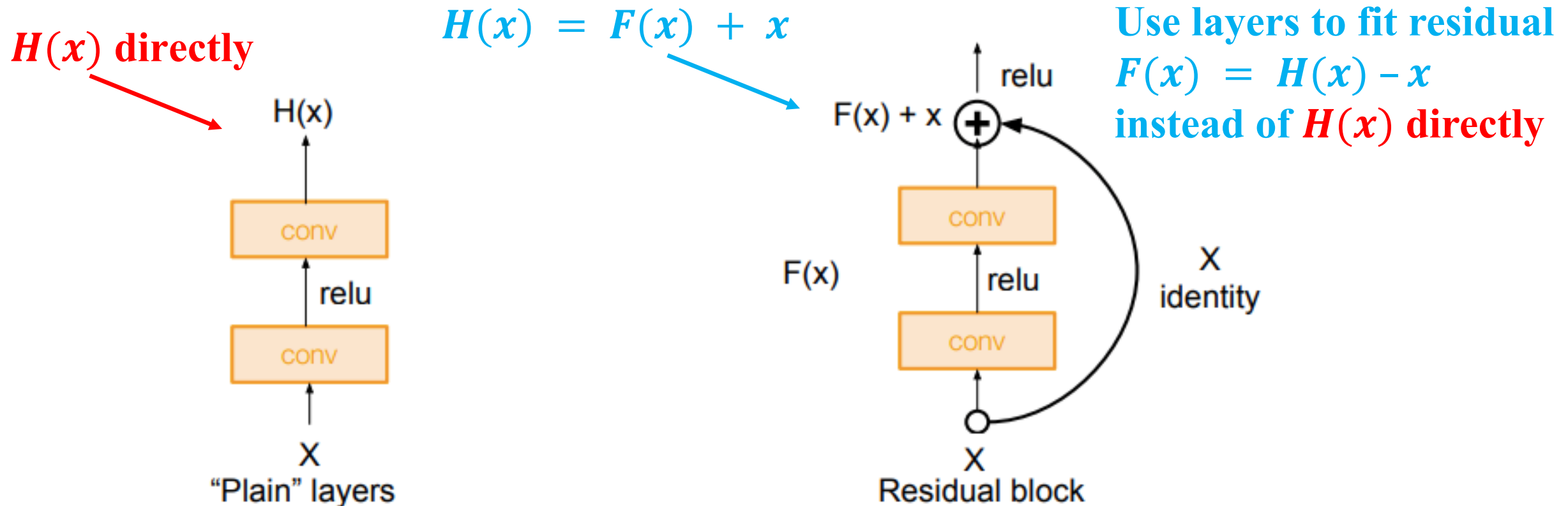  **directly trying to fit a desired underlying mapping**

**$H(x)$ directly**



H(x)

conv

relu

conv

X

"Plain" layers

relu

F(x) + x ⊕

conv

F(x)          relu

conv

X
identity

X

Residual block

# 2. Model Architecture

- **Solution : Use network layers to fit a residual mapping instead of**

  **directly trying to fit a desired underlying mapping**

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

**$H(x)$ directly**

**Use layers to fit residual**
$$F(x) = H(x) - x$$
**instead of $H(x)$ directly**
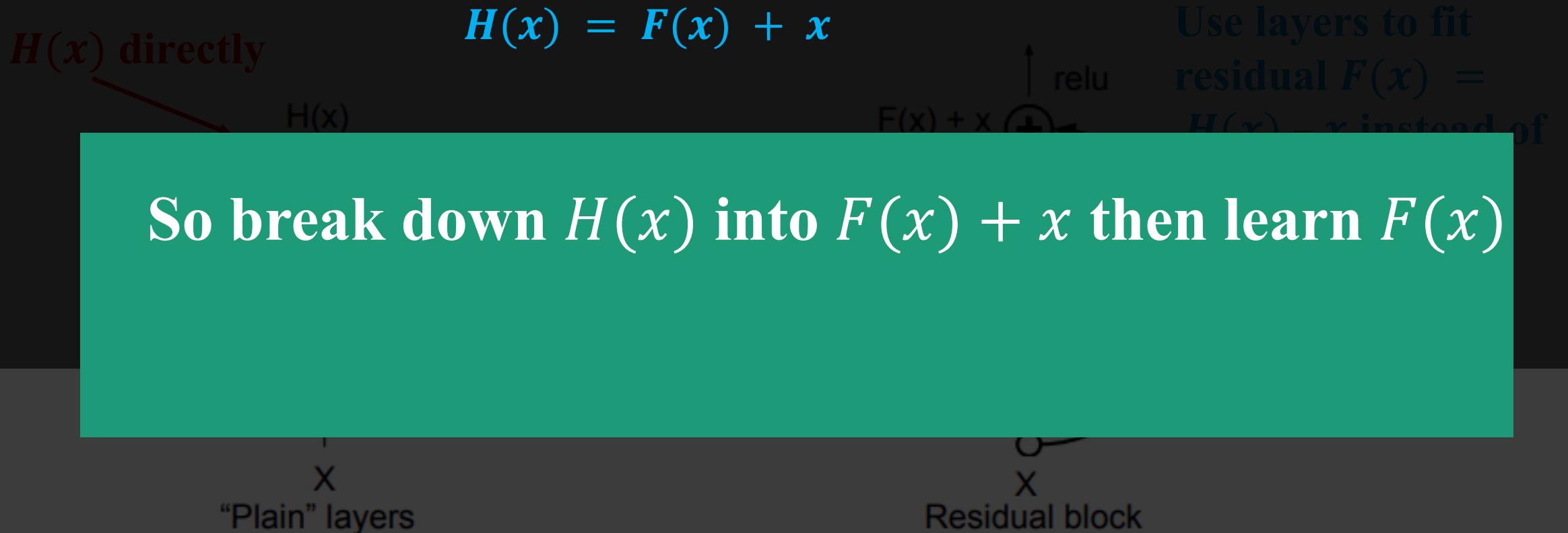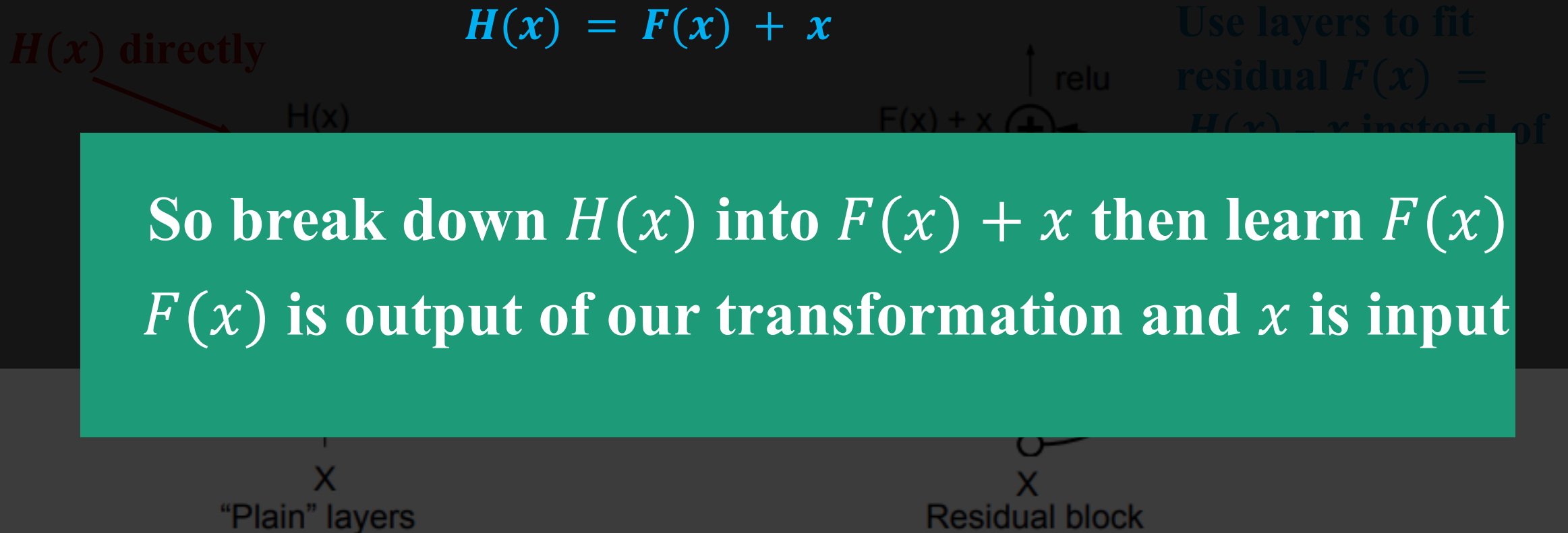
# 2. Model Architecture

- Solution : Use network layers to fit a residual mapping instead of directly trying to fit a desired underlying mapping

$H(x)$ directly

H(x)

relu

F(x) + x

Use layers to fit residual $F(x) = H(x) - x$ instead of

X
"Plain" layers

X
Residual block

**It's hard to learn H($x$) with very deep network**

# 2. Model Architecture

- Solution : Use network layers to fit a residual mapping instead of directly trying to fit a desired underlying mapping
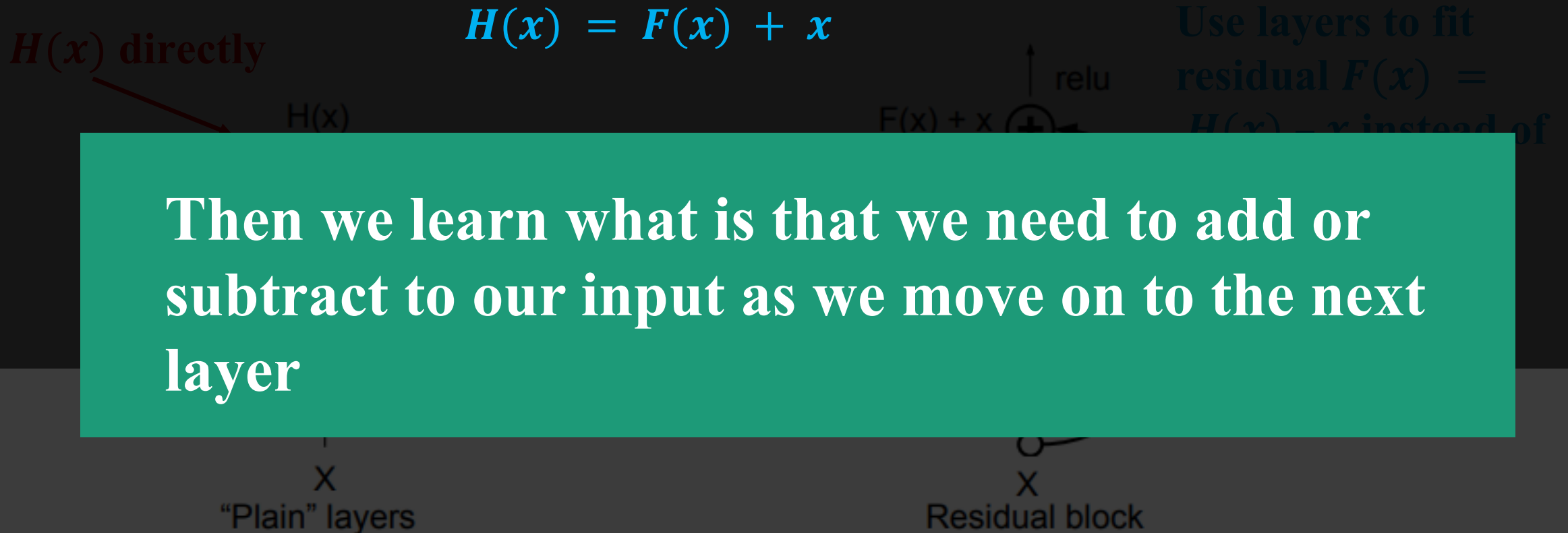
$H(x)$ directly

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

Use layers to fit residual $F(x) = H(x) - x$ instead of

H(x)

relu

F(x) + x

So break down $H(x)$ into $F(x) + x$ then learn $F(x)$

X
"Plain" layers

X
Residual block

# 2. Model Architecture

- Solution : Use network layers to fit a residual mapping instead of directly trying to fit a desired underlying mapping
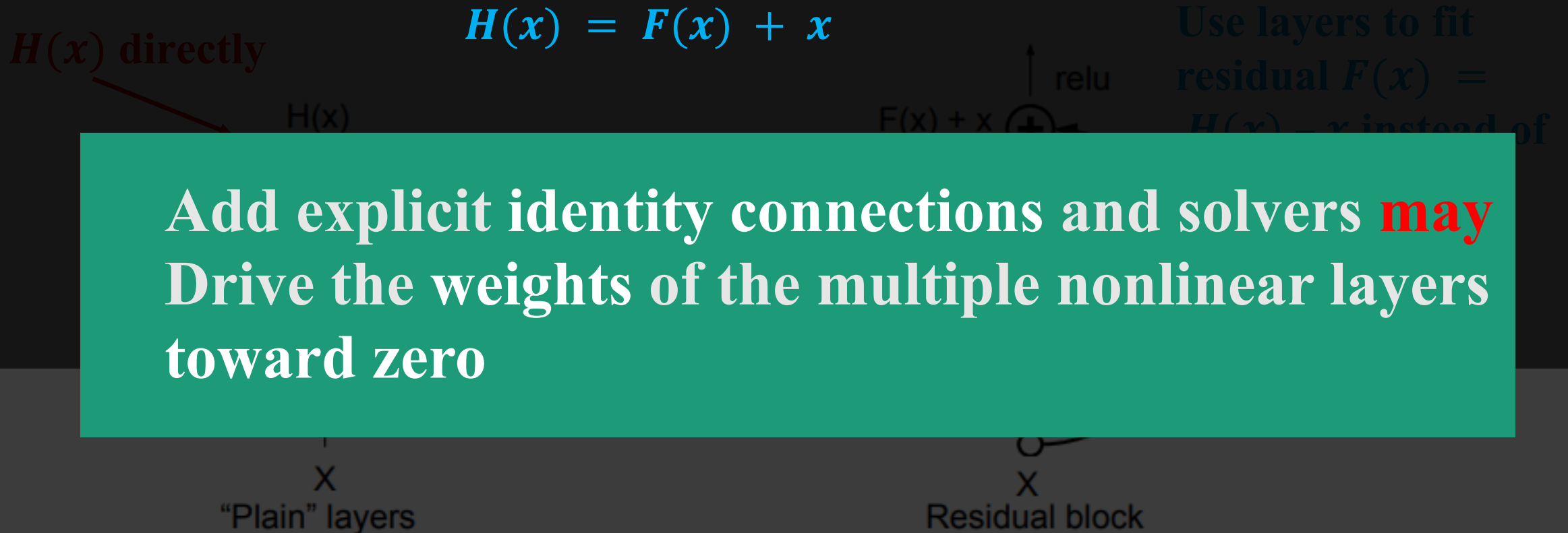
$H(x)$ directly

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

H(x)

relu

F(x) + x

Use layers to fit residual $F(x) =$

So break down $H(x)$ into $F(x) + x$ then learn $F(x)$

$F(x)$ is output of our transformation and $x$ is input
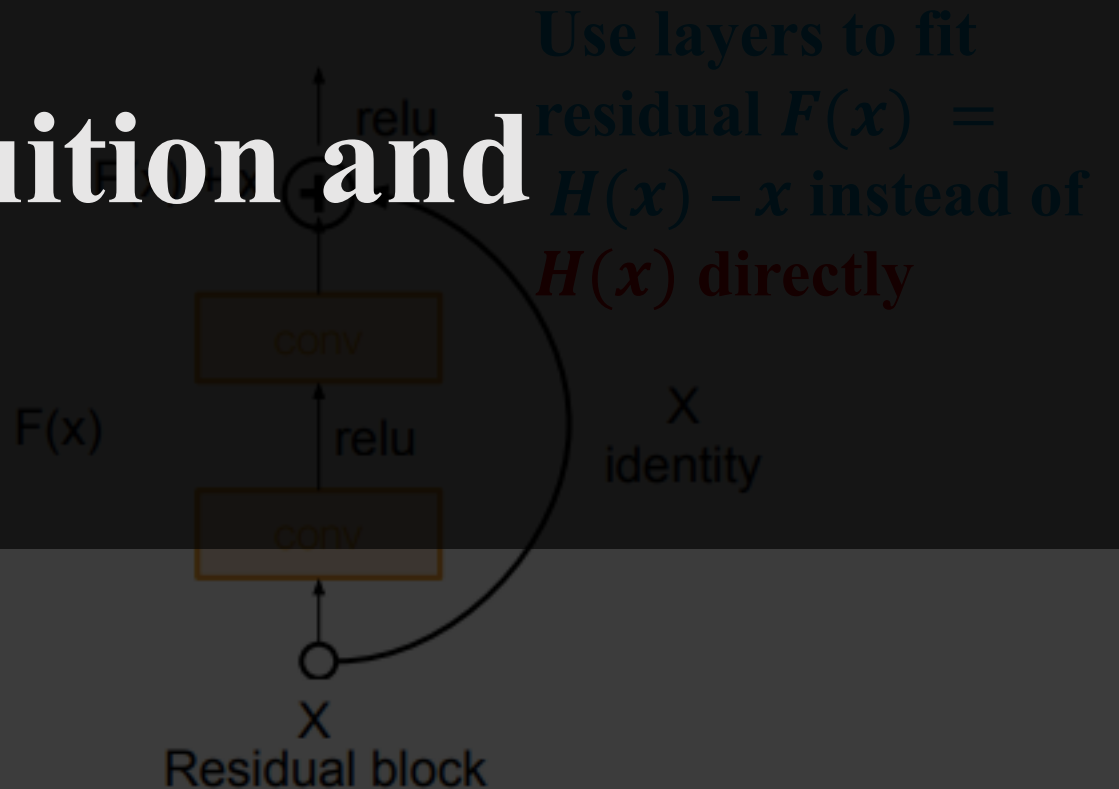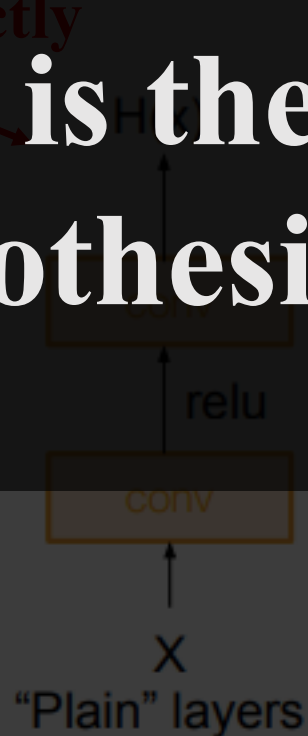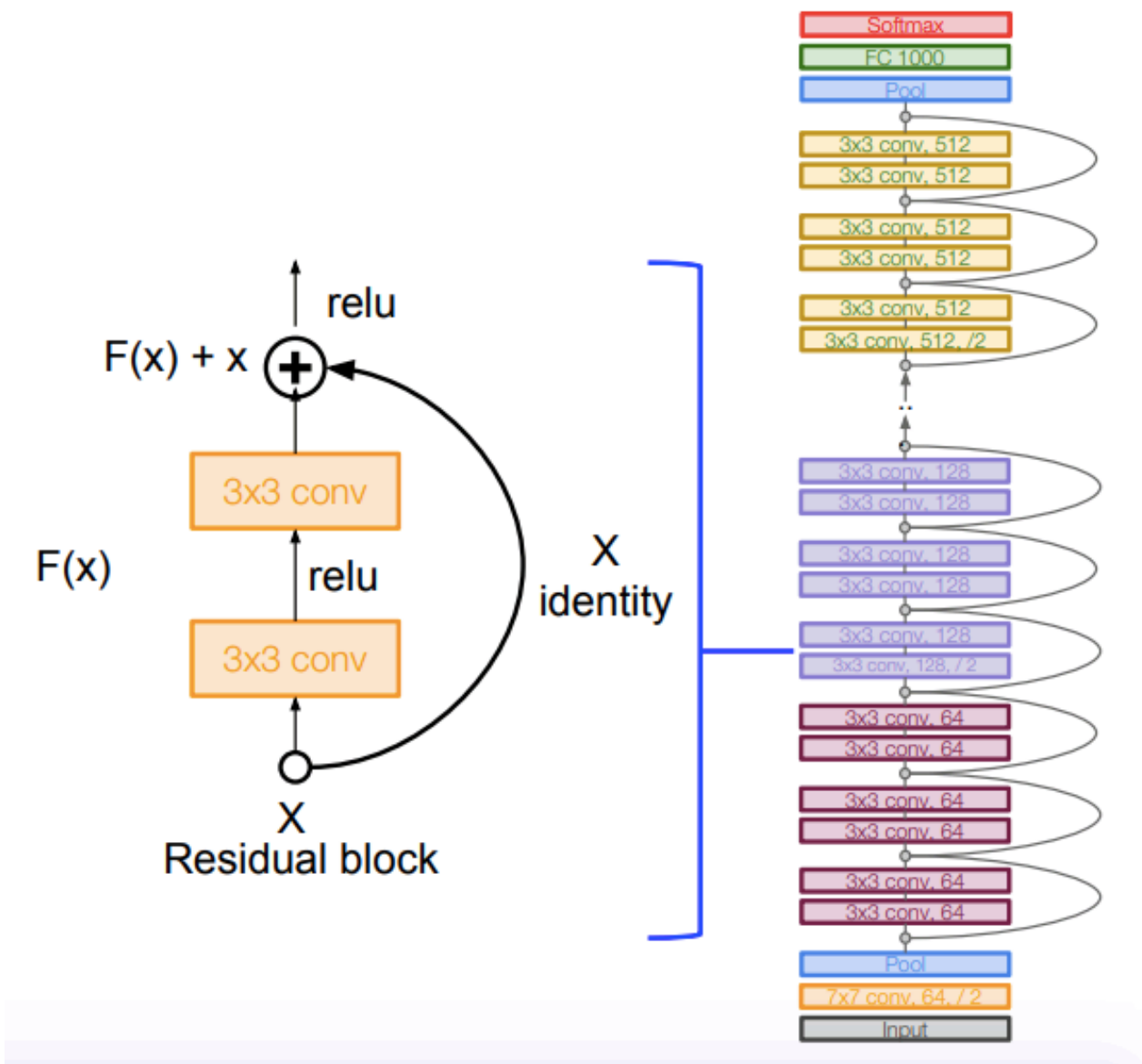
X
"Plain" layers

X
Residual block

# 2. Model Architecture

- Solution : Use network layers to fit a residual mapping instead of directly trying to fit a desired underlying mapping

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

$H(x)$ directly

H(x)

Use layers to fit residual $F(x) = H(x) - x$ instead of

relu

F(x) + x

Then we learn what is that we need to add or subtract to our input as we move on to the next layer

X

"Plain" layers

X

Residual block

# 2. Model Architecture

- Solution : Use network layers to fit a residual mapping instead of directly trying to fit a desired underlying mapping

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

$H(x)$ directly

Use layers to fit residual $F(x) = H(x) - x$ instead of

relu

H(x)

F(x) + x

X

"Plain" layers

X

Residual block

**Add explicit identity connections and solvers may Drive the weights of the multiple nonlinear layers toward zero**

# 2. Model Architecture

- Solution : Use network layers to fit a residual mapping instead of

directly trying to fit a desired underlying mapping

$H(x)$ directly

Use layers to fit residual $F(x) = H(x) - x$ instead of $H(x)$ directly

relu

**This is the main Intuition and Hypothesis**

F(x)

relu

conv

relu

conv

X

X
identity

X

X

"Plain" layers

Residual block

# 2. Model Architecture

# 2. Model Architecture

- **Stack residual block**

- **Every residual block has two 3x3 conv layers**

# 2. Model Architecture

- **3x3 conv, 64 filters**

  **double # of filters**

- **3x3 conv, 128 filters**
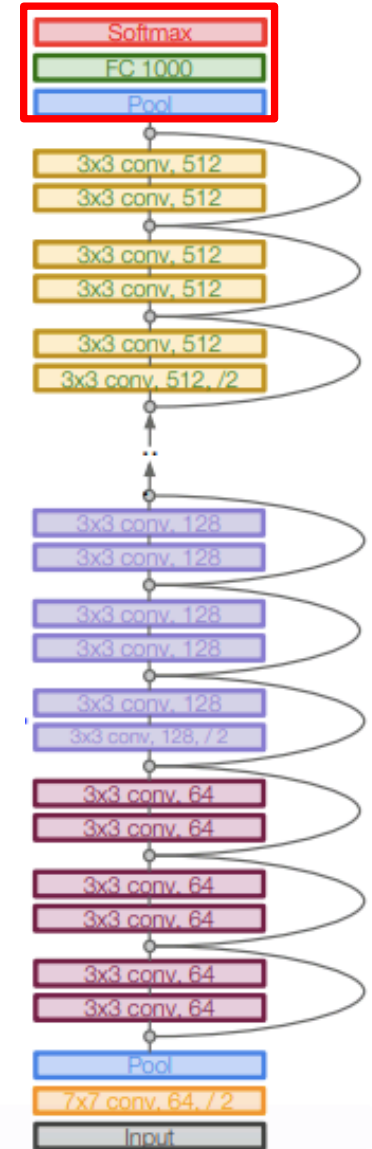- **downsample ½ spatially using stride 2 (not pooling)**

# 2. Model Architecture

- **Additional conv layer at the beginning**
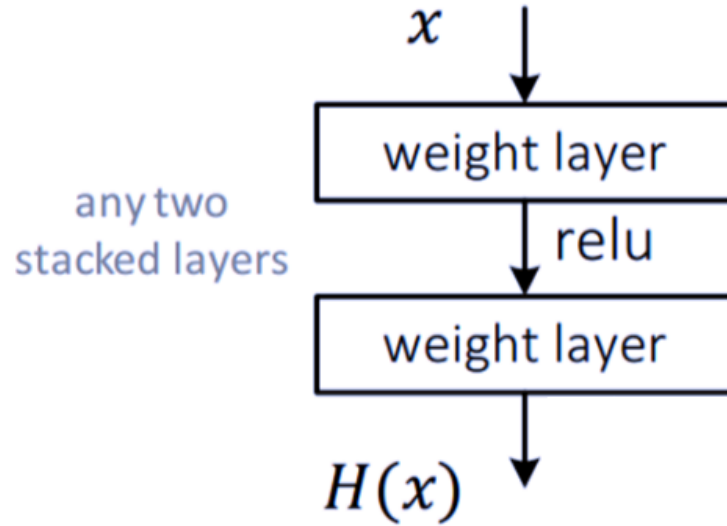- **Input -> 7x7 conv, 64 , /2 -> pool**

# 2. Model Architecture

- **global average pooling layer**

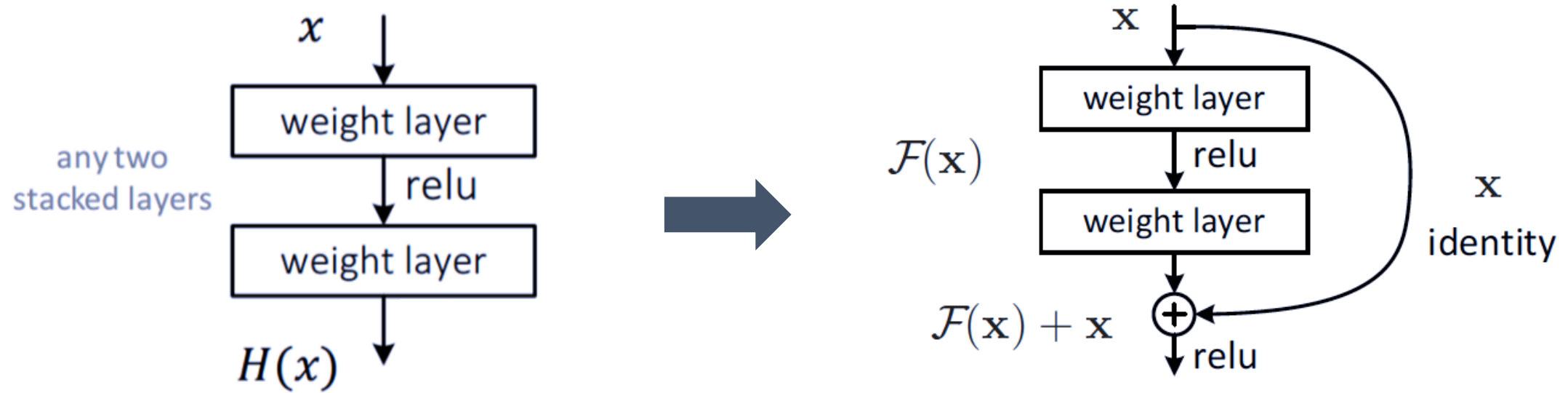- **No FC layer at the end only FC 1000 to output classes**
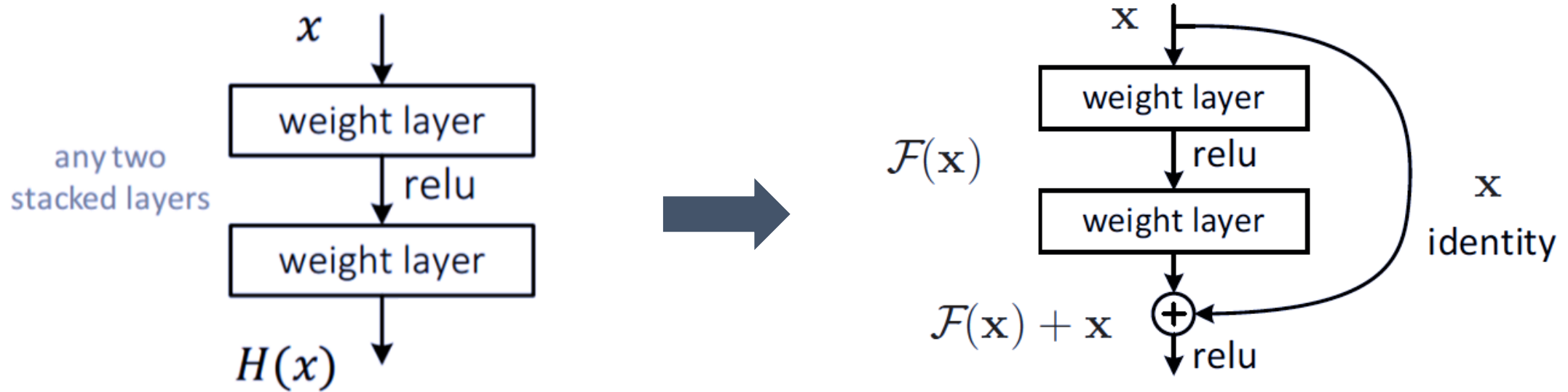
# 3. Why Residual?



- **Simple 2 layers of Conv**

- **Goal : Extraction features /  to get optimal $H(x)$**

- **Output $H(x)$**

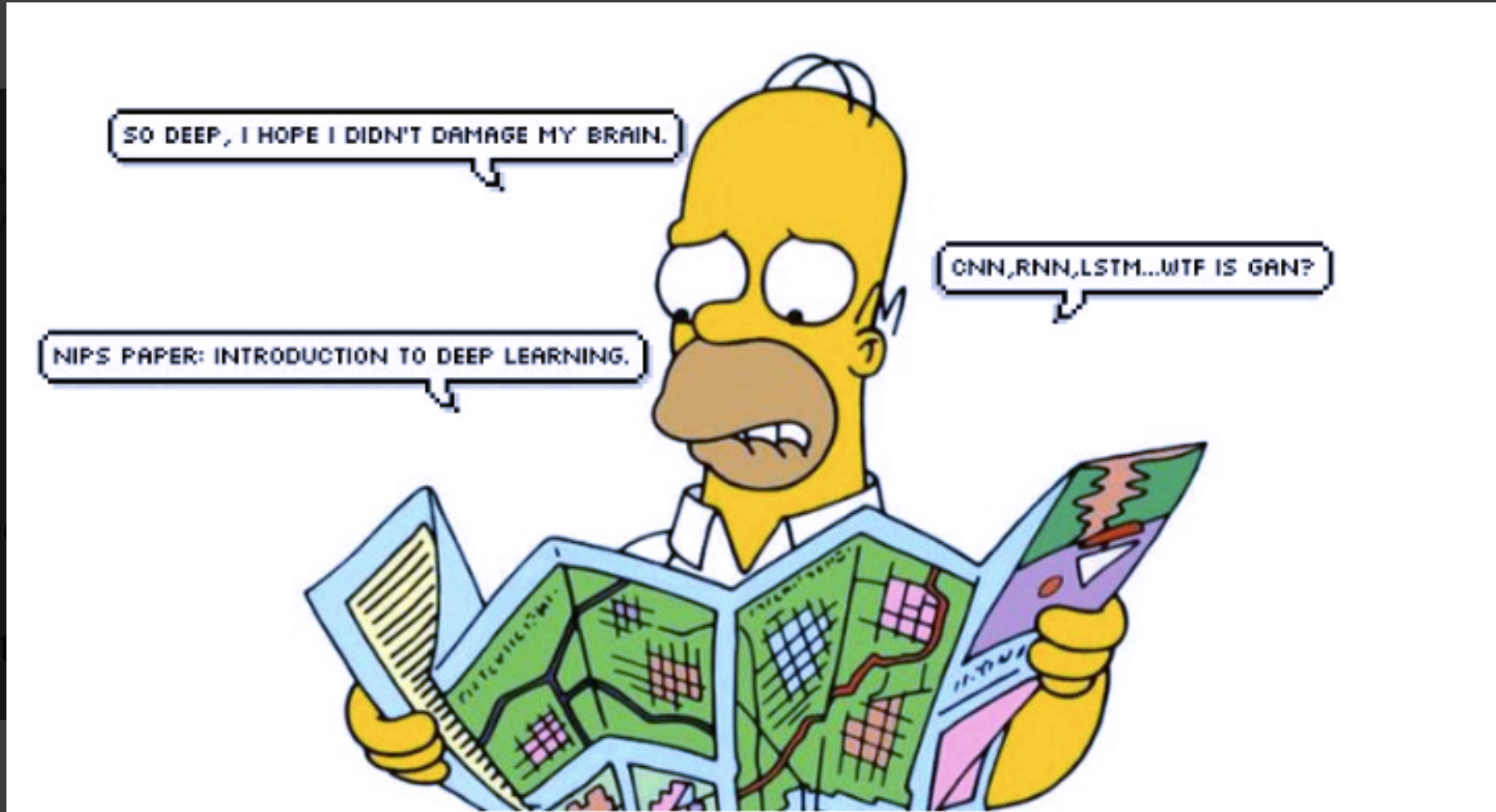- **The network would learn $W$ toward Goal**

# 3. Why Residual?

# 3. Why Residual?



- **If fixed Goal : $H(x) - x$**

- **Then the two layers should be learned toward $H(x) - x$**

- **Here, If $F(x) := H(x) - x$**

- **Then output $H(x) = F(x) + x$**

# 3. Why Residual?

any two
stacked lay



- If fixe

- Then

- Here,

- Then output $H(x) = F(x) + x$

# 3. Why Residual?

- **If <span style="color:red">Our Goal</span> is Not $H(x)$ But <span style="color:red">$H(x) - x$</span>**

# 3. Why Residual?

- **If <span style="color:red">Our Goal</span> is Not $H(x)$ But <span style="color:red">$H(x) - x$</span>**

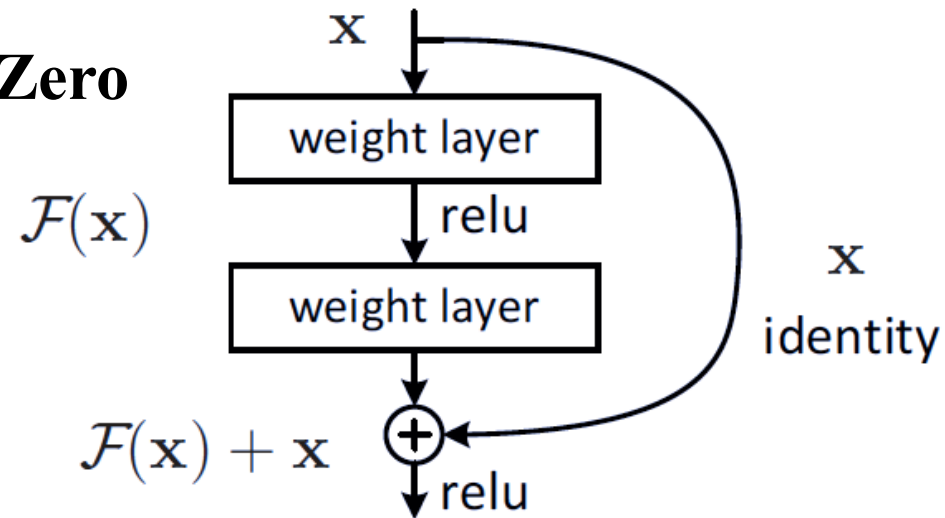- **What we learn is Residual = Output - Input**

# 3. Why Residual?

- **If Our Goal is Not $H(x)$ But $H(x) - x$**

- **What we learn is Residual = Output - Input**

- **New output = $F(x) + x$**

- **New output = Residual + Input**
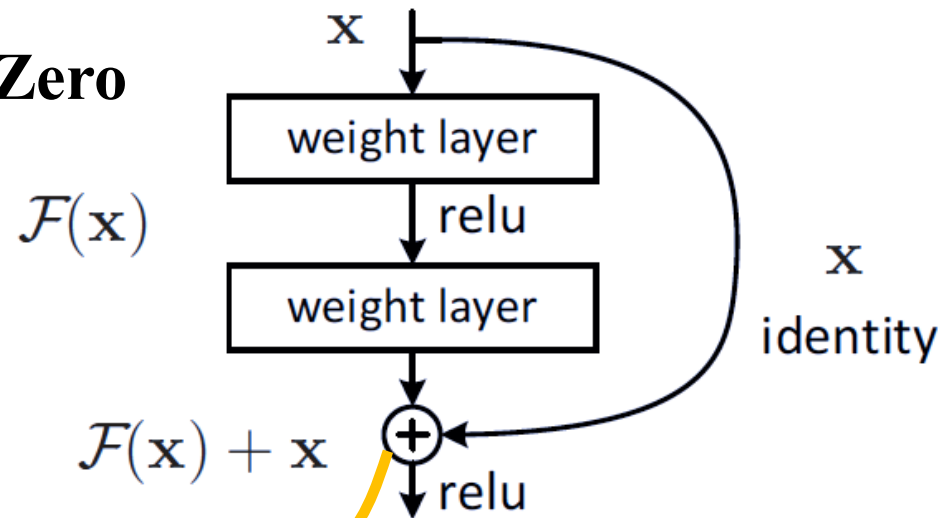
# 3. Why Residual?

- **If Our Goal is Not $H(x)$ But $H(x) - x$**

- **What we learn is Residual = Output - Input**

- **New output $= F(x) + x$**

- **New output $=$ Residual $+$ Input**

- **Here, If Residual Approximated to Zero**

- **Then, New output $\approx x$**

# 3. Why Residual?

- **If Our Goal is Not $H(x)$ But $H(x) - x$**

- **What we learn is Residual = Output - Input**

- **New output = $F(x) + x$**

- **New output = Residual + Input**

- **Here, If Residual Approximated to Zero**

- **Then, New output ≈ $x$**

- **No matter how deep the network is we just add $x$**

- **This part is Pre - Conditioning**

# 3. Why Residual?
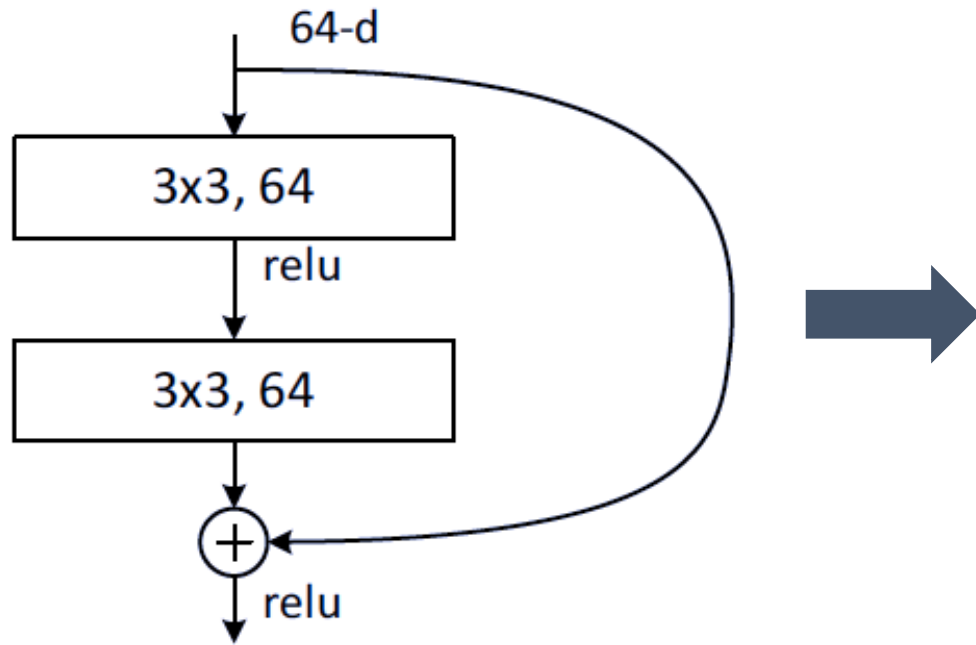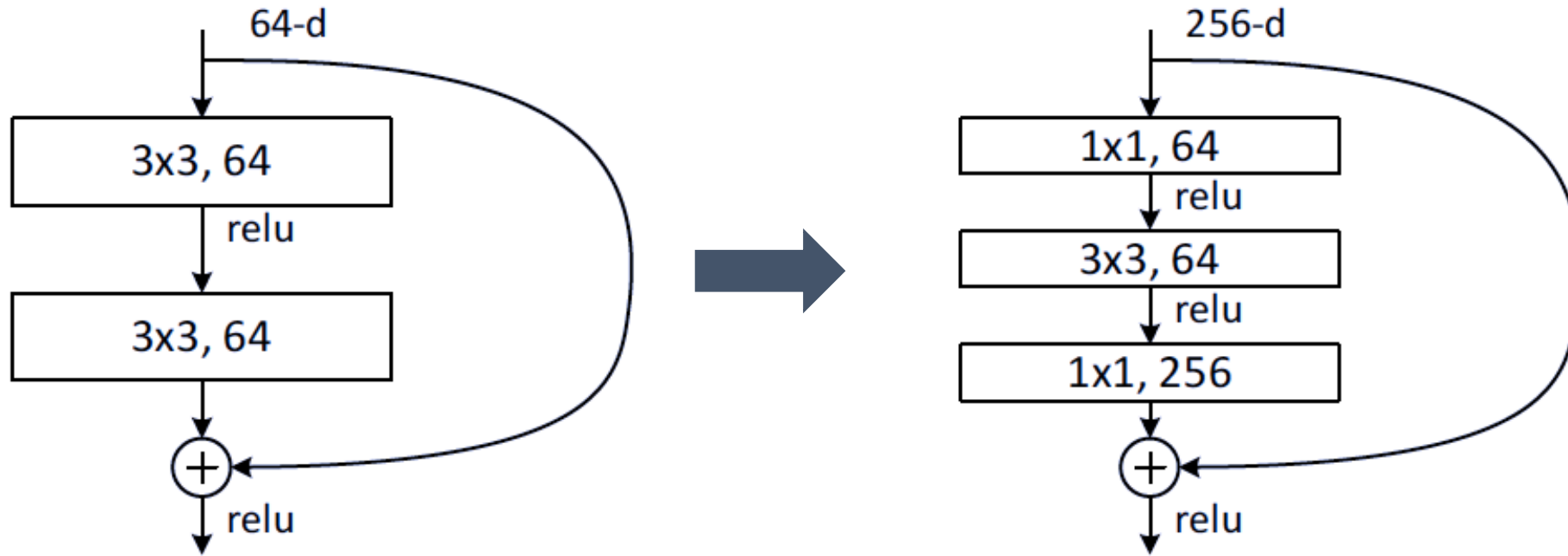
- **The extremely deep residual nets are <span style="color:red">easy</span> to optimize**

- **The deep residual nets can <span style="color:red">easily</span> enjoy <span style="color:red">accuracy gains</span> from greatly <span style="color:red">increased depth</span>, producing results substantially better than previous networks**

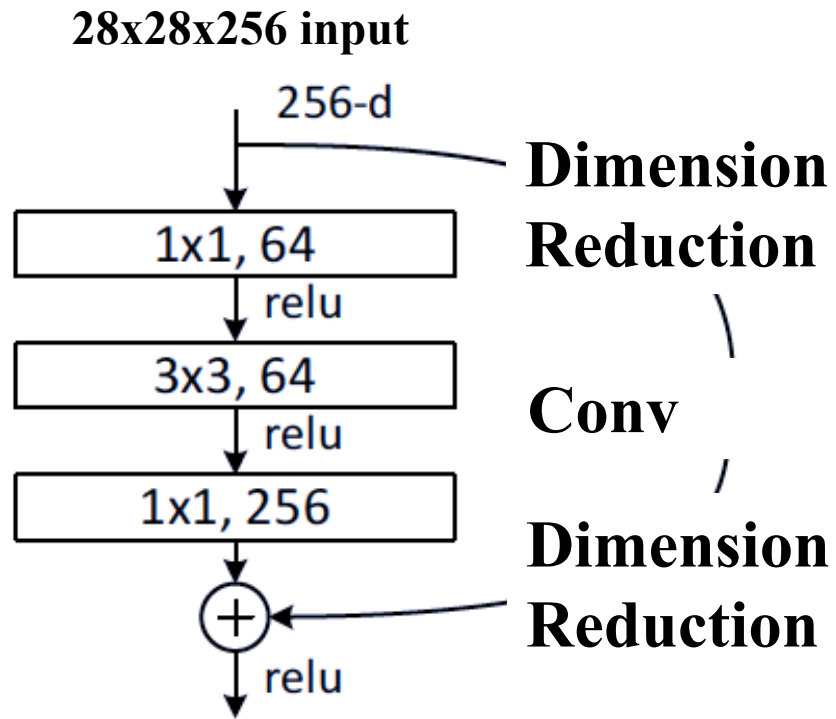# 4. Deeper bottleneck architecture

# 4. Deeper bottleneck architecture

# 4. Deeper bottleneck architecture

**28x28x256 input**

256-d

1x1, 64
relu

Dimension
Reduction

3x3, 64
relu

Conv

1x1, 256

Dimension
Reduction

⊕

relu

**1x1 conv, 64 filters to project to 28x28x64**

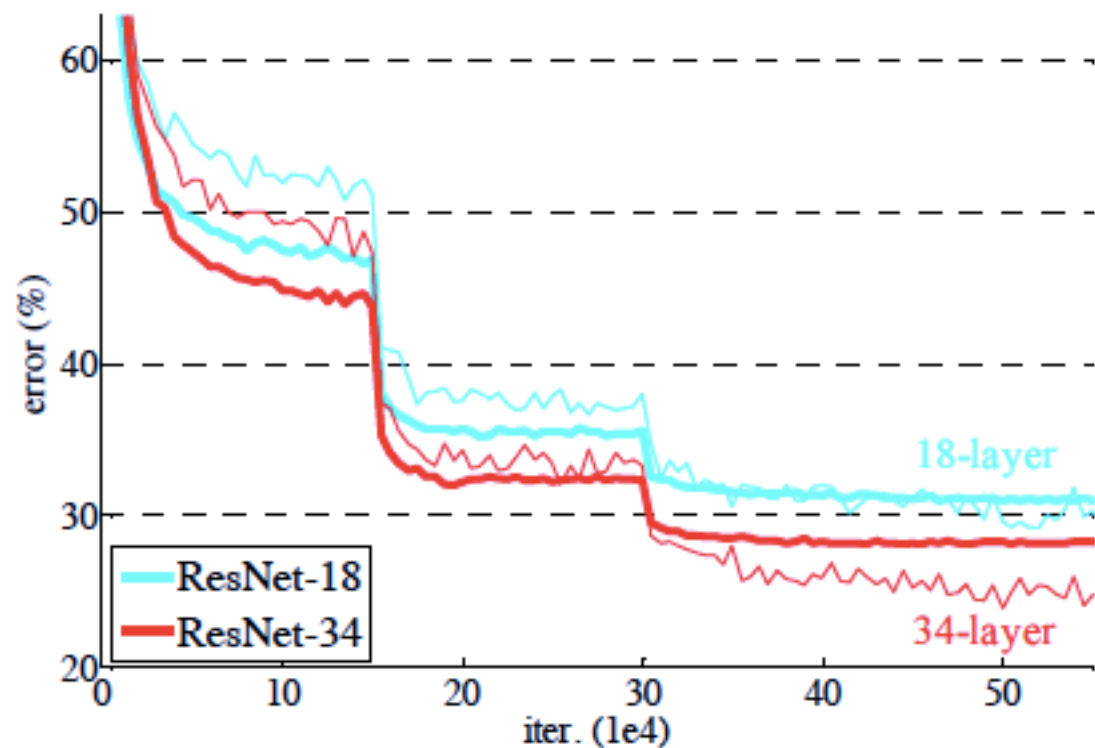**3x3 conv operates over only 64 feature maps**

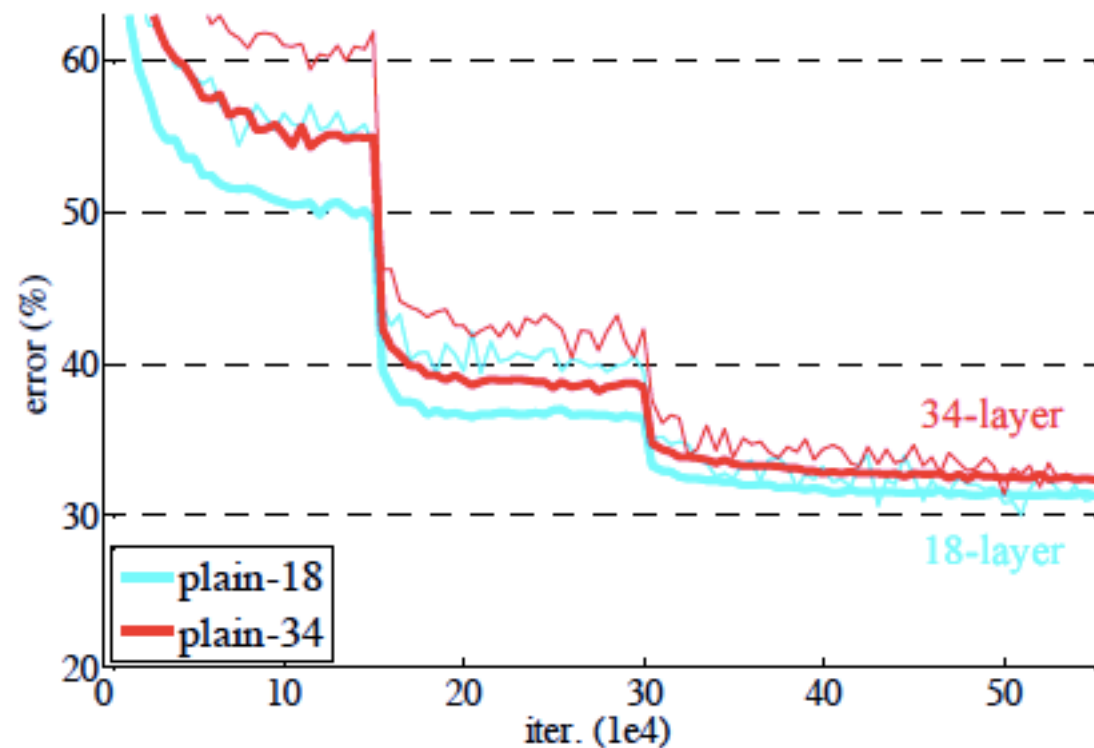**1x1 conv, 256 filters projects back to 256 feature maps (28x28x256)**
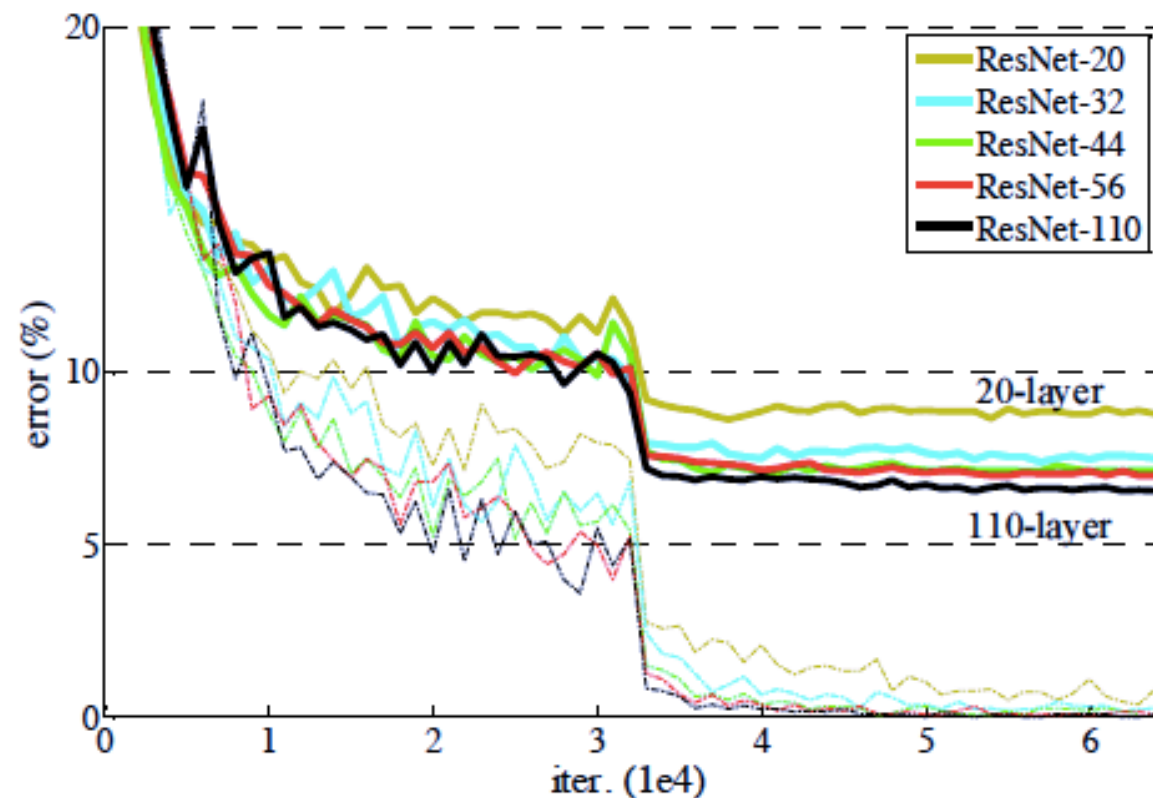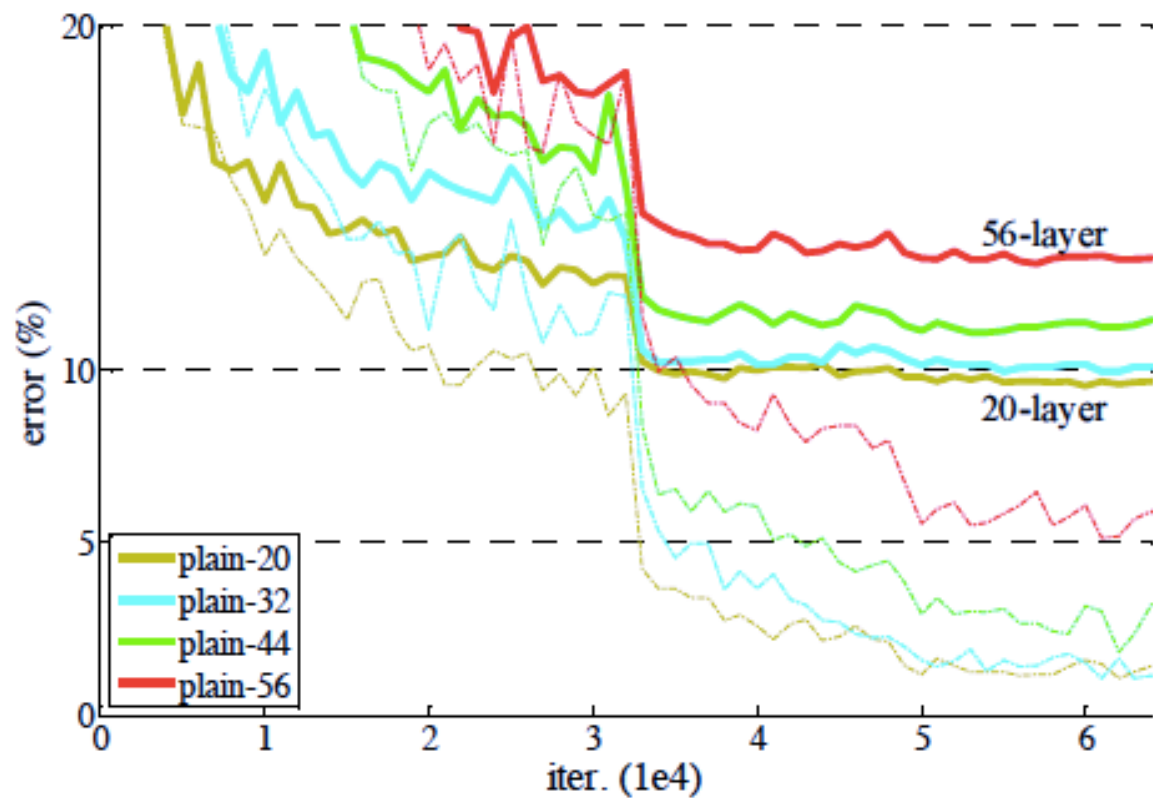
# 5. Experimental Result

- **Training ResNet in practice**

  - **Batch Normalization after every Conv layer**

  - **Xaxier/2 initialization from He at al.**

  - **No dropout**

  - **SGD + Momentum(0.9)**

  - **Learning rate: 0.1, divided by 10 when validation error plateaus**

  - **Mini batch size 256**
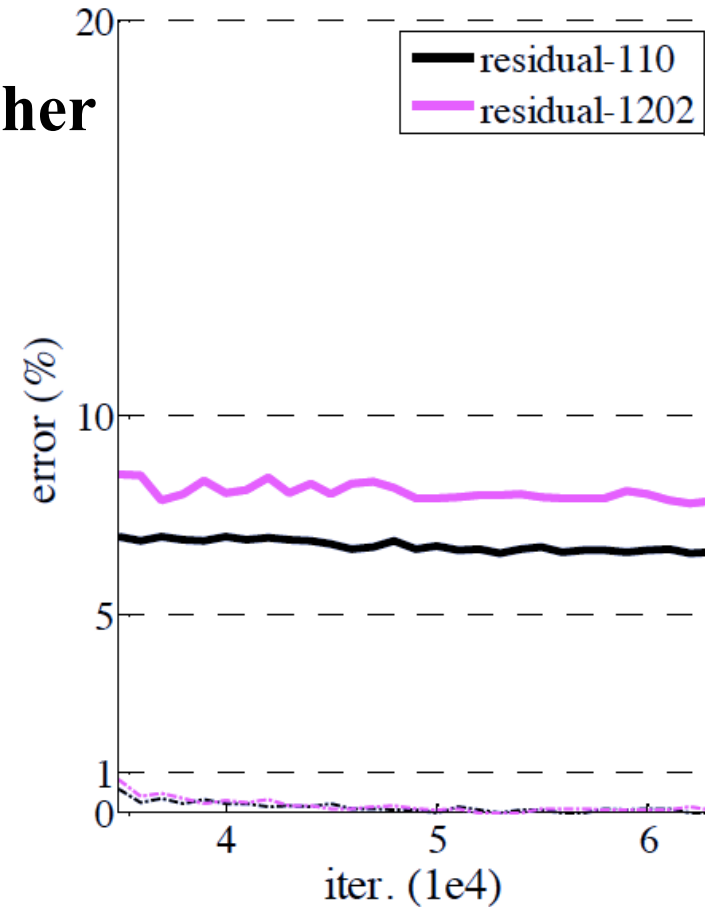
  - **Weight decay of 1e-5**

# 5. Experimental Result

# 5. Experimental Result

# 5. Experimental Result

- **No optimization difficulty**

- **But Ultra deeper model is worse than the other**

| method | # layers | # params | error (%) |
|---|---|---|---|
| Maxout [10] | | | 9.38 |
| NIN [25] | | | 8.81 |
| DSN [24] | | | 8.22 |
| FitNet [35] | 19 | 2.5M | 8.39 |
| Highway [42, 43] | 19 | 2.3M | 7.54 (7.72±0.16) |
| Highway [42, 43] | 32 | 1.25M | 8.80 |
| ResNet | 20 | 0.27M | 8.75 |
| ResNet | 32 | 0.46M | 7.51 |
| ResNet | 44 | 0.66M | 7.17 |
| ResNet | 56 | 0.85M | 6.97 |
| ResNet | 110 | 1.7M | **6.43** 6.61±0.16) |
| ResNet | 1202 | 19.4M | 7.93 |

# 5. Experimental Result

- **Experimental Result**

  - **Able to train very deep network without** <span style="color:orange">**Degradation**</span>

  - **Deeper Network now achieve lowing training error as expected**

  - **Swept 1st place in all ILSVRC and COCO 2015 competitions**

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

감사합니다