

Resnet review

By beomgon.yu

Contents

1. Introduction
2. Reformulation by suboptimal problem(residual leaning)
3. Gradient pumping(by shortcut connection)
4. trainable Non linearity
5. etc

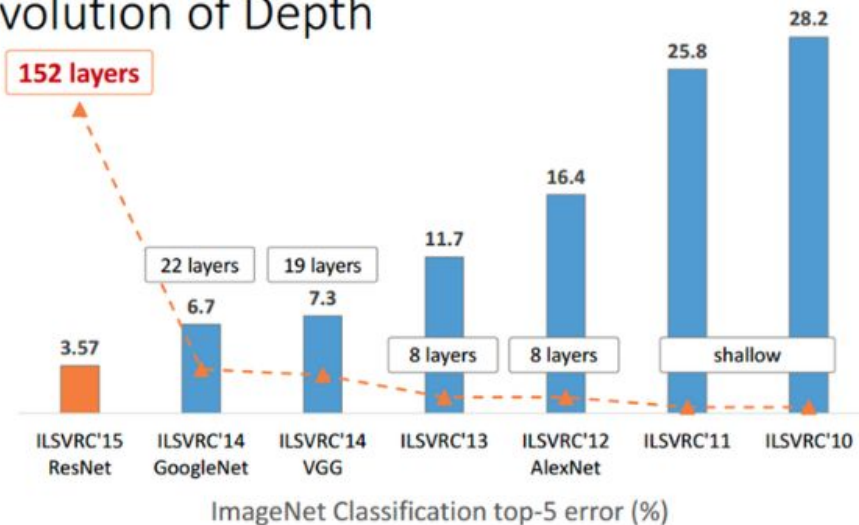
Introduction

Deeper network get better performance

Resnet use lots of layer, first time,

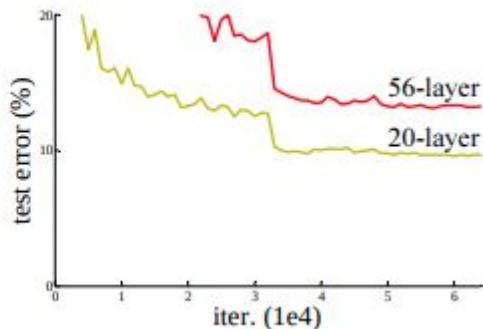
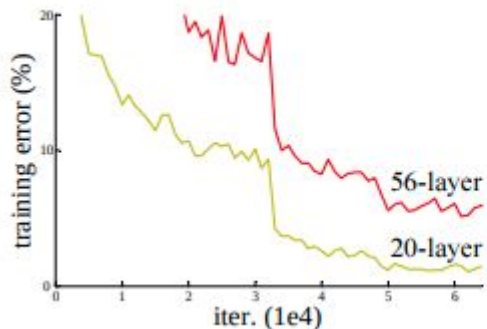
Is there a problem or how they could do?

Revolution of Depth



Introduction

Deeper network get better performance,
But at some point, **performance degradation** happened.
This is neither overfitting, nor vanishing gradient
(normalized initialization and batch normalization)

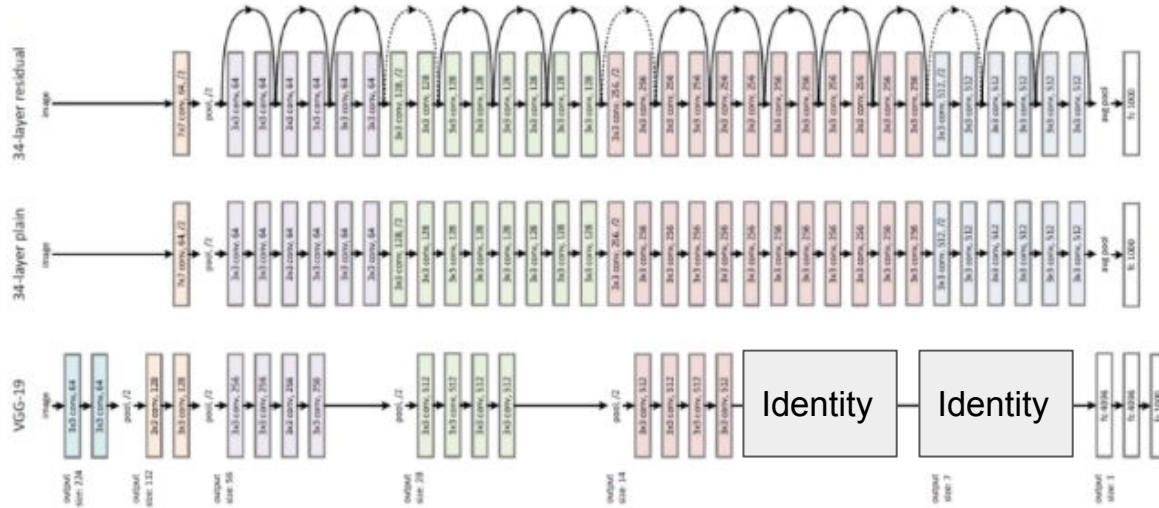


Introduction

Deeper network, Is there a No solution, NO?

Solution must be exist, (counter part shallow network)

But its difficult to learn Identity, or takes long time because of Non linearity



Reformulation by suboptimal problem

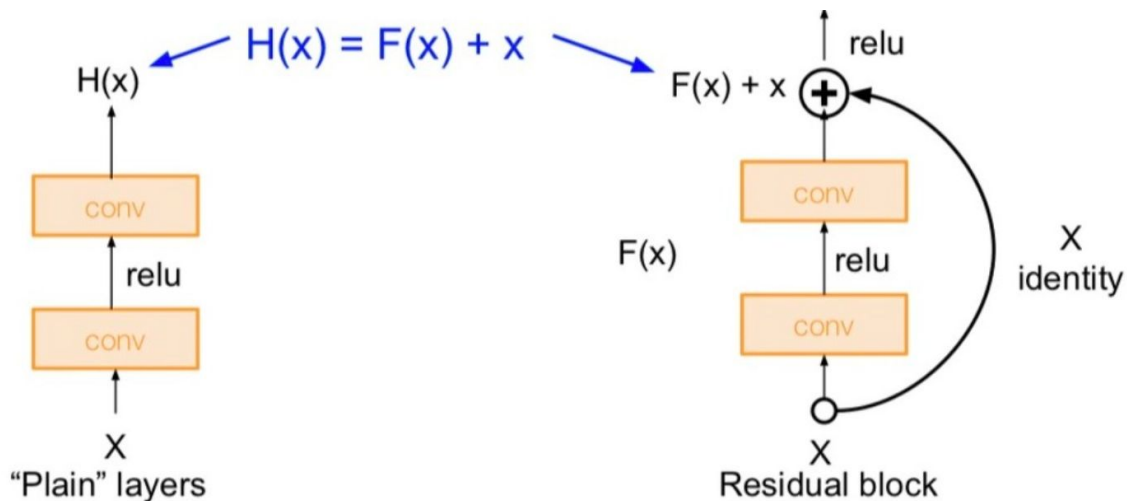
residual learning

instead of direct learning of mapping, just learning residual (**offset**)

(forced to learn identity)

(offset, its similar to anchor in object detection)

Because of non linearity, learning identity is difficult



Use layers to

fit residual

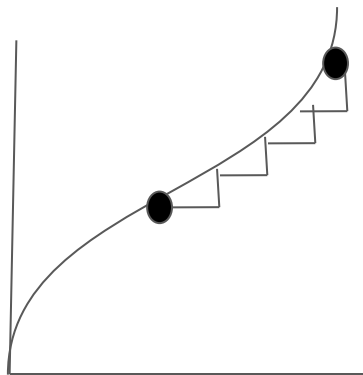
$$F(x) = H(x) - x$$

instead of

$H(x)$ directly

Reformulation by suboptimal problem

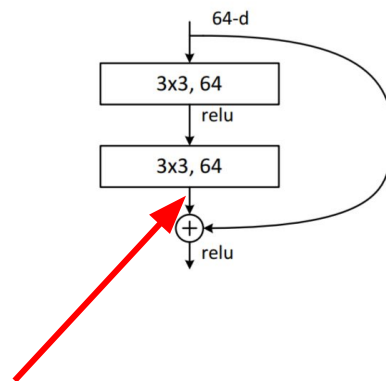
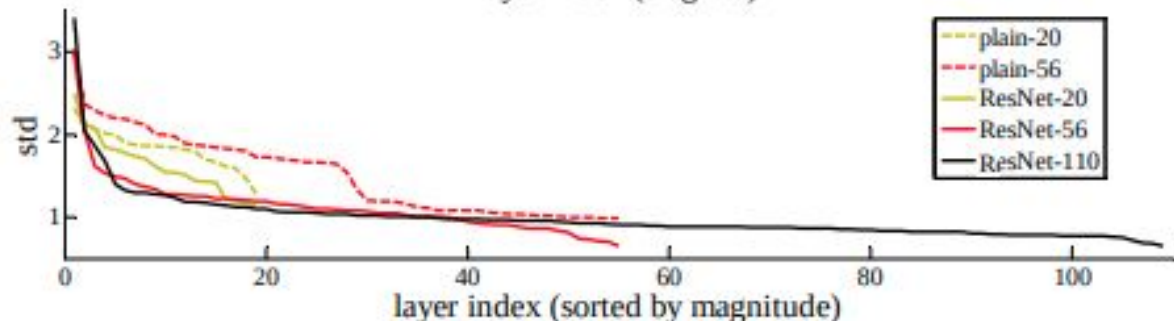
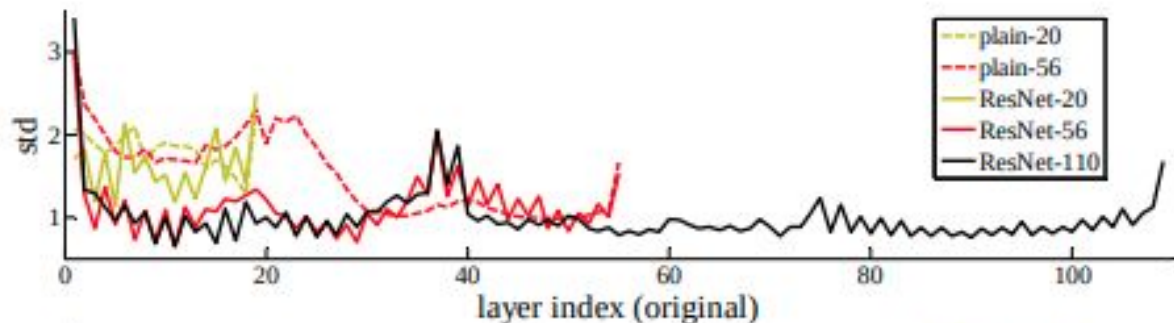
천리길도 한 걸음부터



Reformulation by suboptimal problem

Analysis of Layer Responses

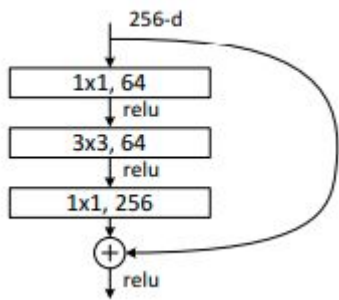
ResNets have generally smaller responses than their plain counterparts



Gradient pumping

this optimization difficulty is unlikely to be caused by vanishing gradients trained with BN, and checked each propagated gradient has healthy norms with BN

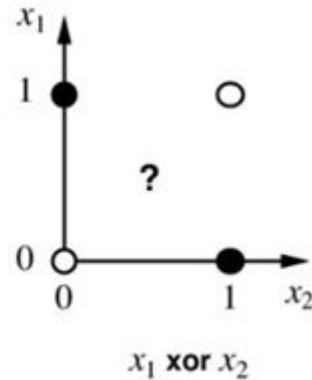
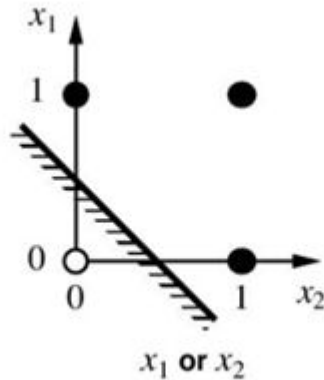
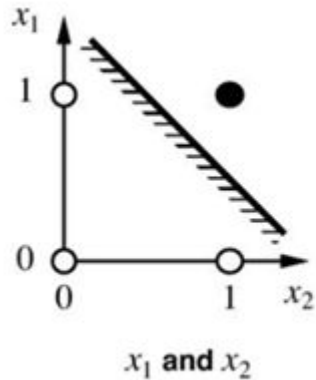
But maybe skip connection give gradient directly to lower layer, so help solving vanishing gradient problem or fast learning



trainable Non linearity

Importance of non linearity

without non linearity, mlp or higher nn can not solve trivial xor problem.

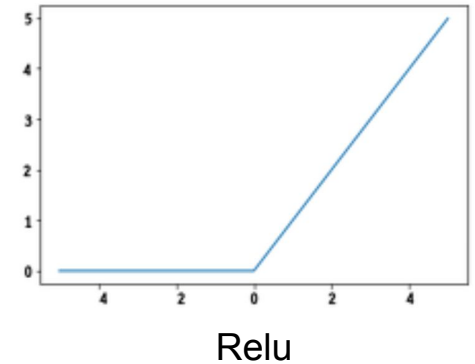
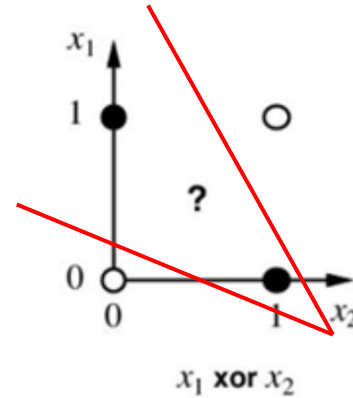
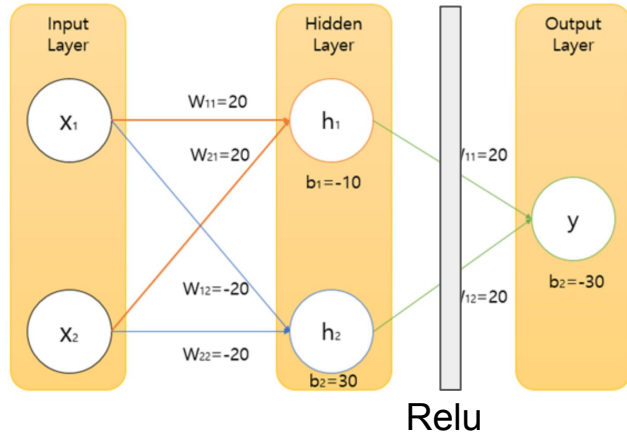


trainable Non linearity

Importance of non linearity

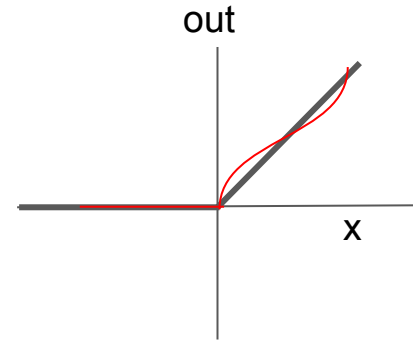
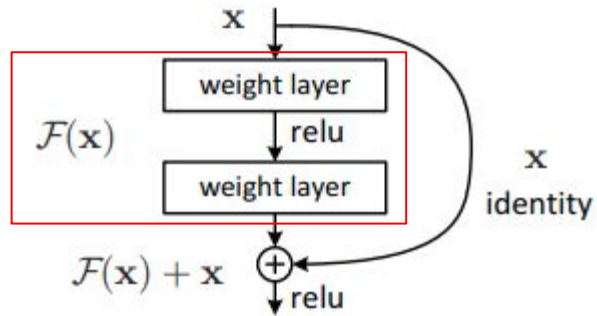
with non linearity, xor can be solved.

how about using more complex non linearity or trainable non linearity?



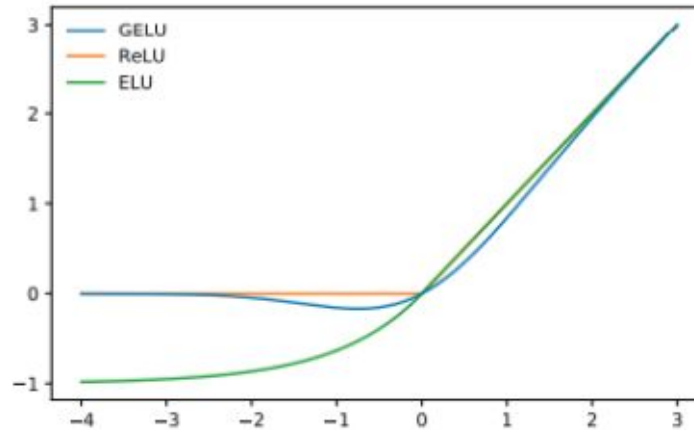
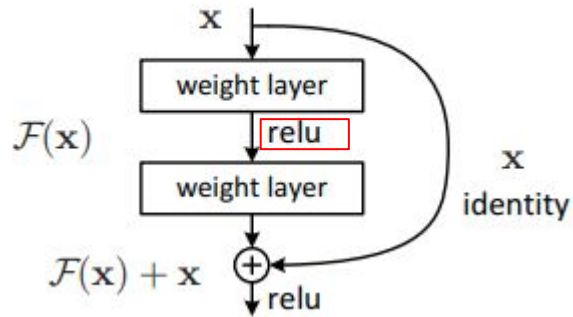
trainable Non linearity

How about if non linearity can be trainable?



trainable Non linearity

in resnet 18, with caltech 101 dataset,
gelu \geq relu, leaky relu, elu $>$ sigmoid, tanh

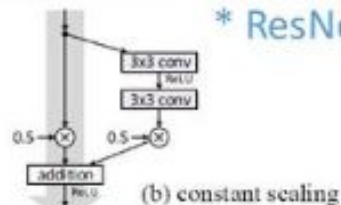
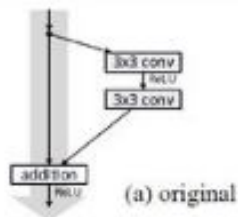


ETC

skip connection variation

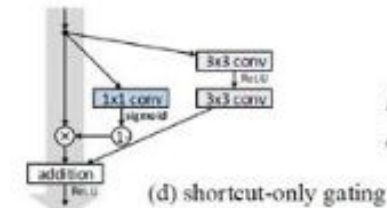
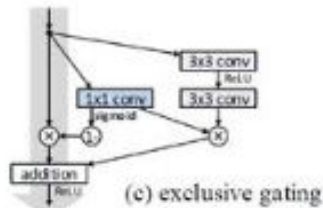
* ResNet-110 on CIFAR-10

$h(x) = x$
error: 6.6%



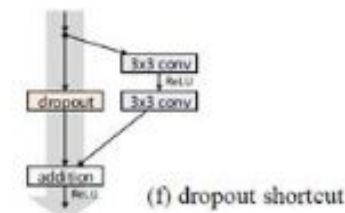
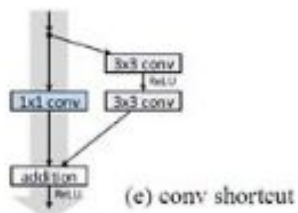
$h(x) = 0.5x$
error: 12.4%

$h(x) = \text{gate} \cdot x$
error: 8.7%
*similar to "Highway Network"



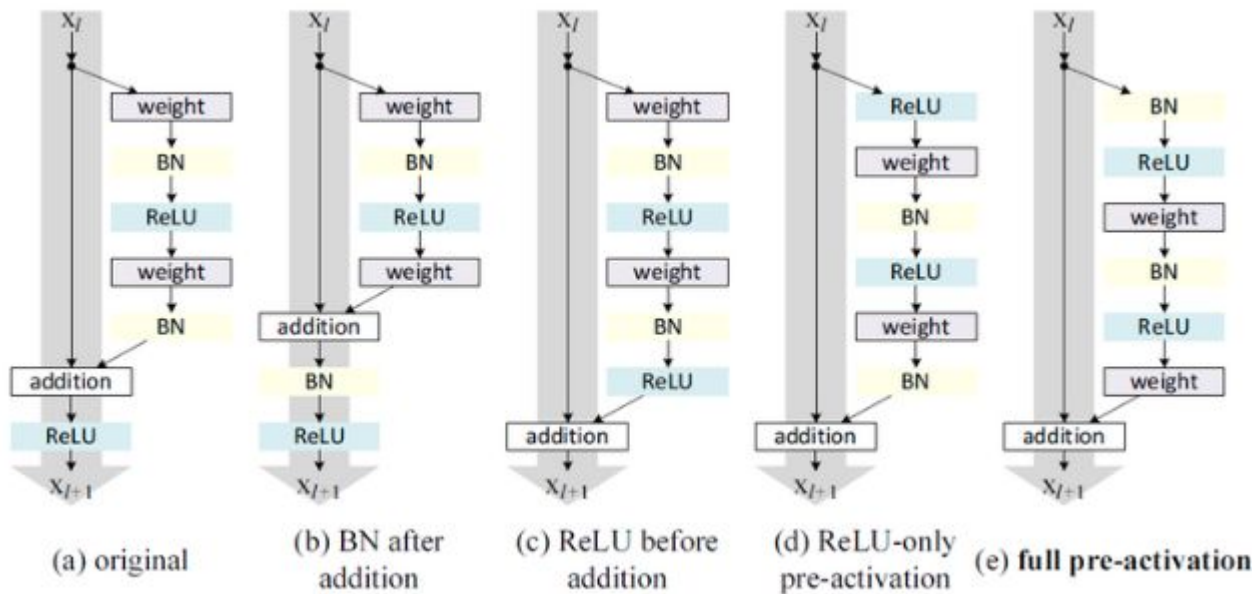
$h(x) = \text{gate} \cdot x$
error: 12.9%

$h(x) = \text{conv}(x)$
error: 12.2%



$h(x) = \text{dropout}(x)$
error: > 20%

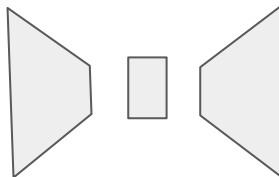
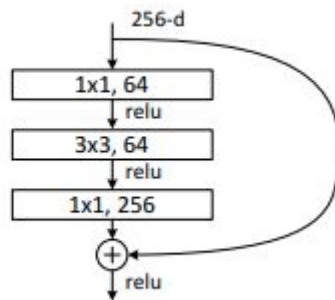
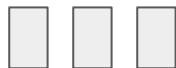
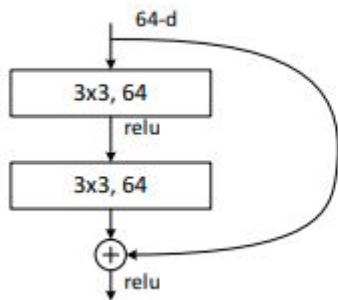
ETC



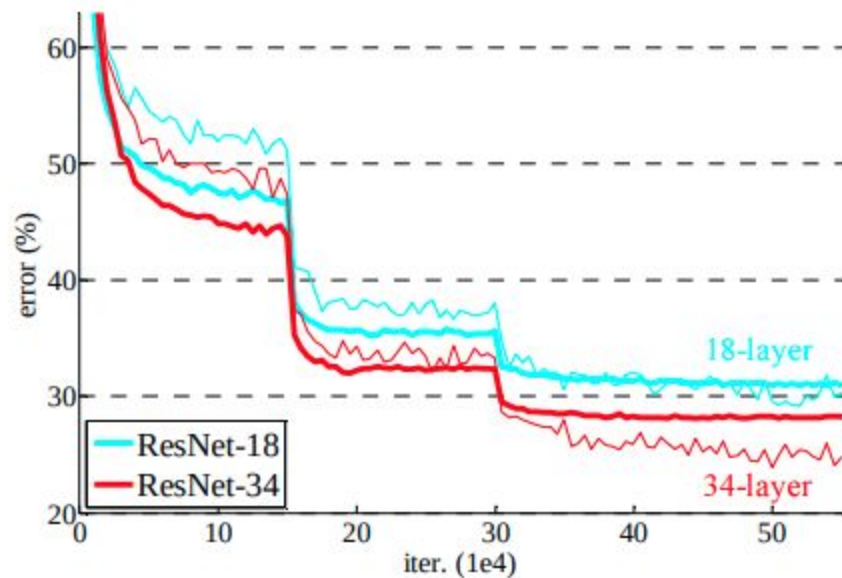
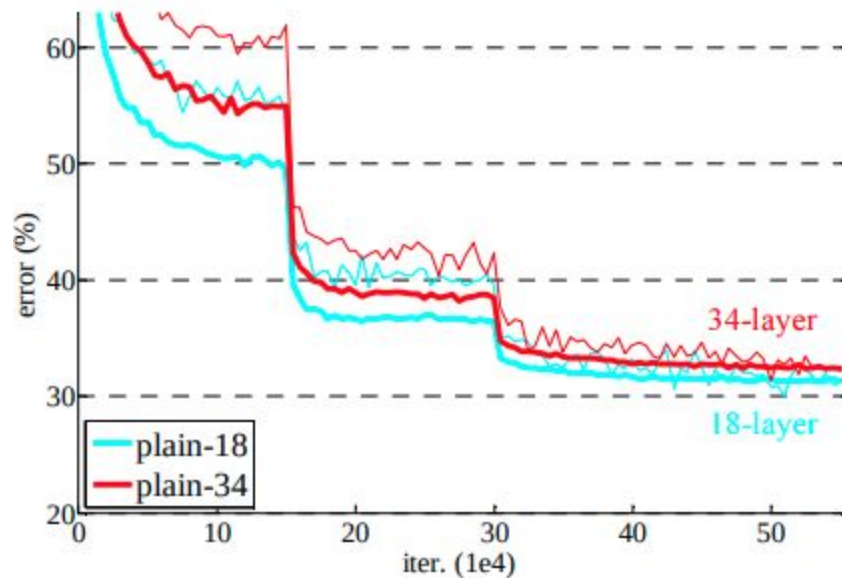
ETC

bottleneck layer

for computation efficiency, use it for more than 50 layers.



Test Result



Thank you !!!