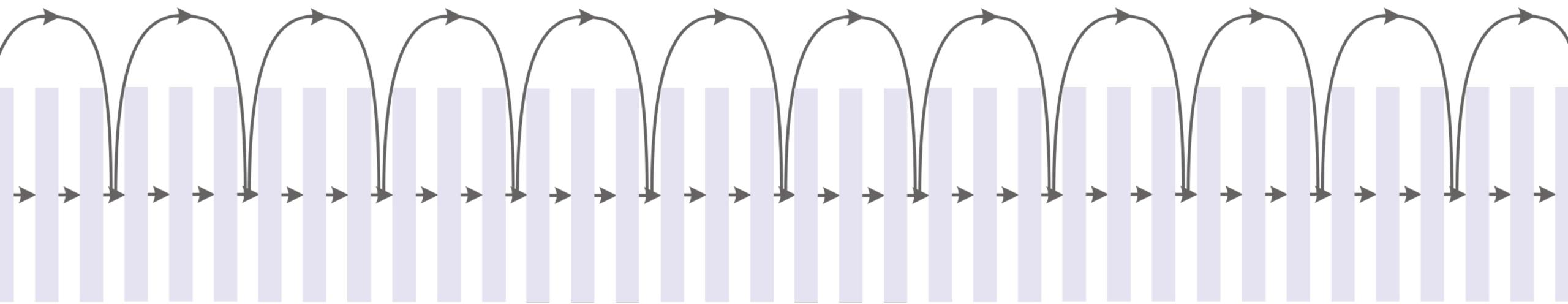


# Lecture 10

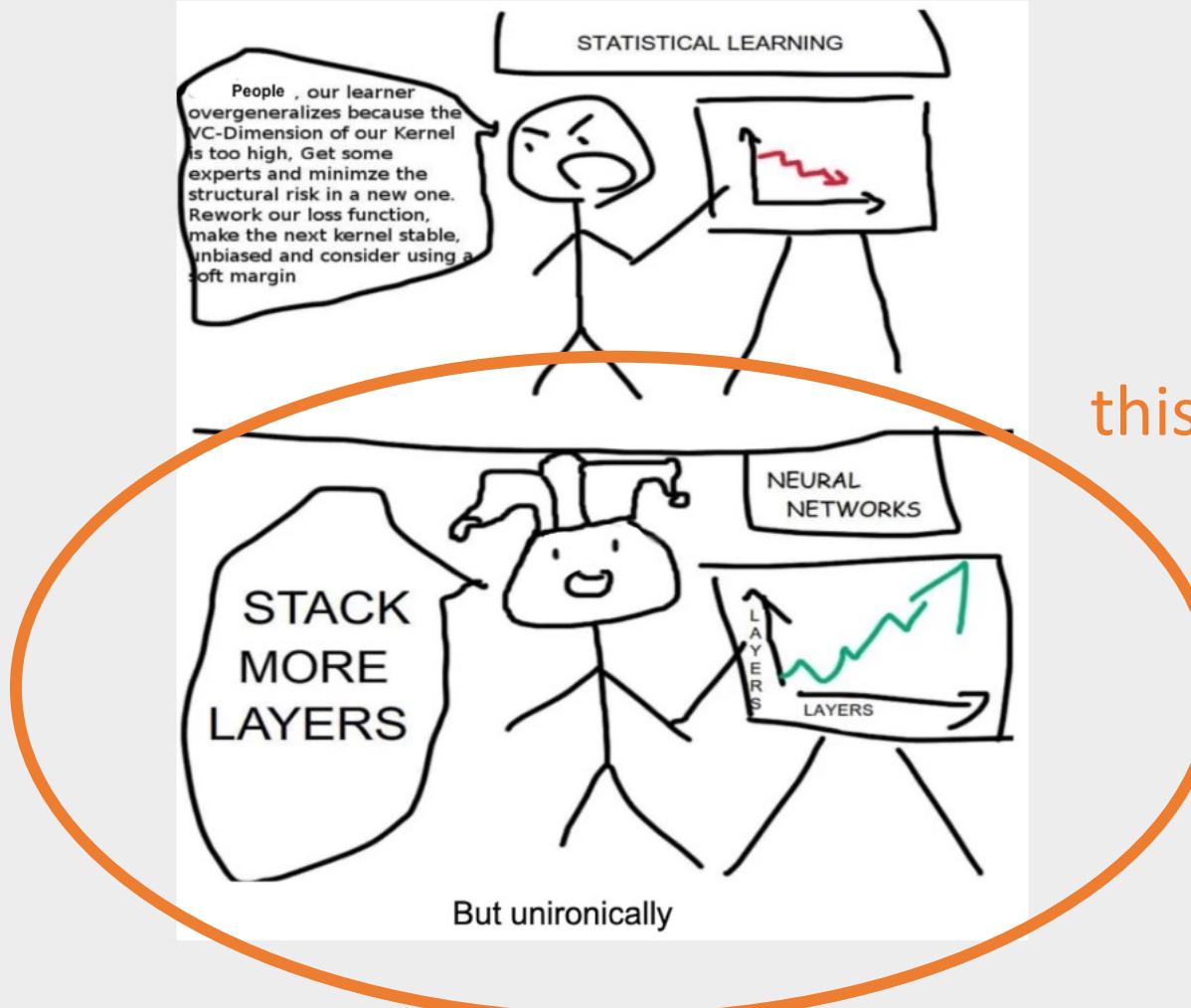
## Going Deep with Neural Networks



# Overview

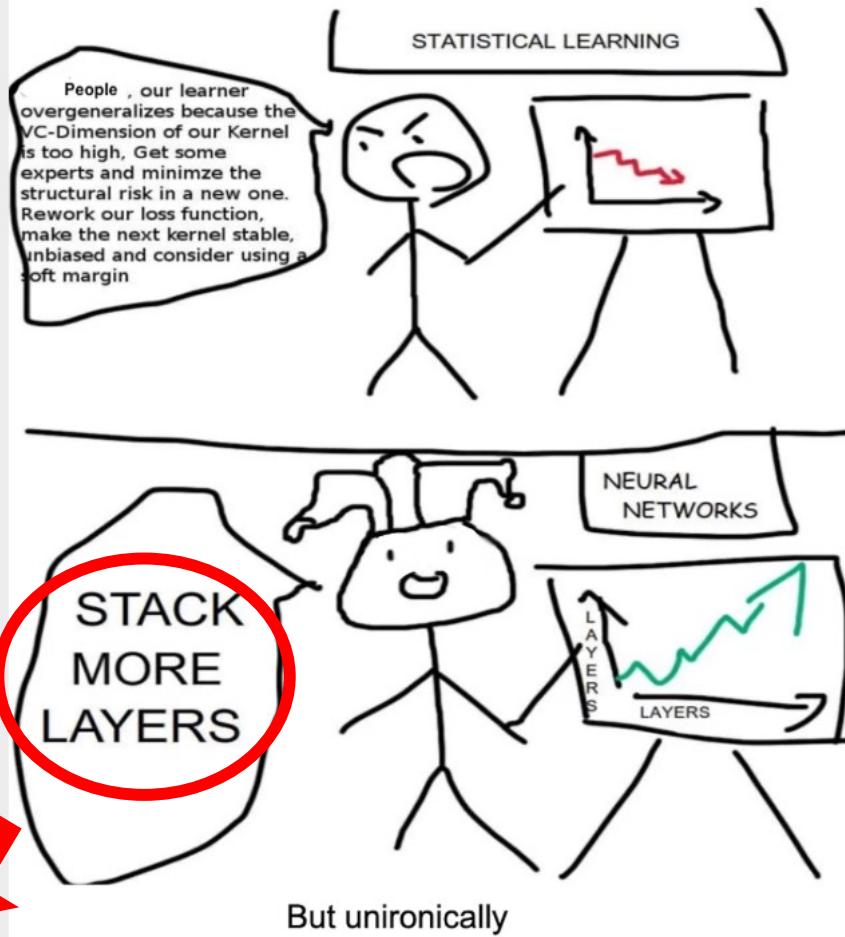
- Why going deep is difficult?
- Weight initialization
- Normalization modules
- Deep Residual Learning

User Can you explain why this is funny. Think about it step-by-step.



GPT-4 The comic is satirizing the difference in approaches to improving model performance between statistical learning and neural networks.

User Can you explain why this is funny. Think about it step-by-step.



GPT-4

The comic is satirizing the difference in approaches to improving model performance between statistical learning and neural networks.

# Why going deep is difficult?

- Troubles accumulate w/ more layers
- Signals get distorted when propagated
- in forward and backward passes



# Why going deep is difficult?

## All about signal propagation

- ReLU activation
  - Weight initialization
  - Normalization modules
  - Deep Residual Learning
- } this lecture



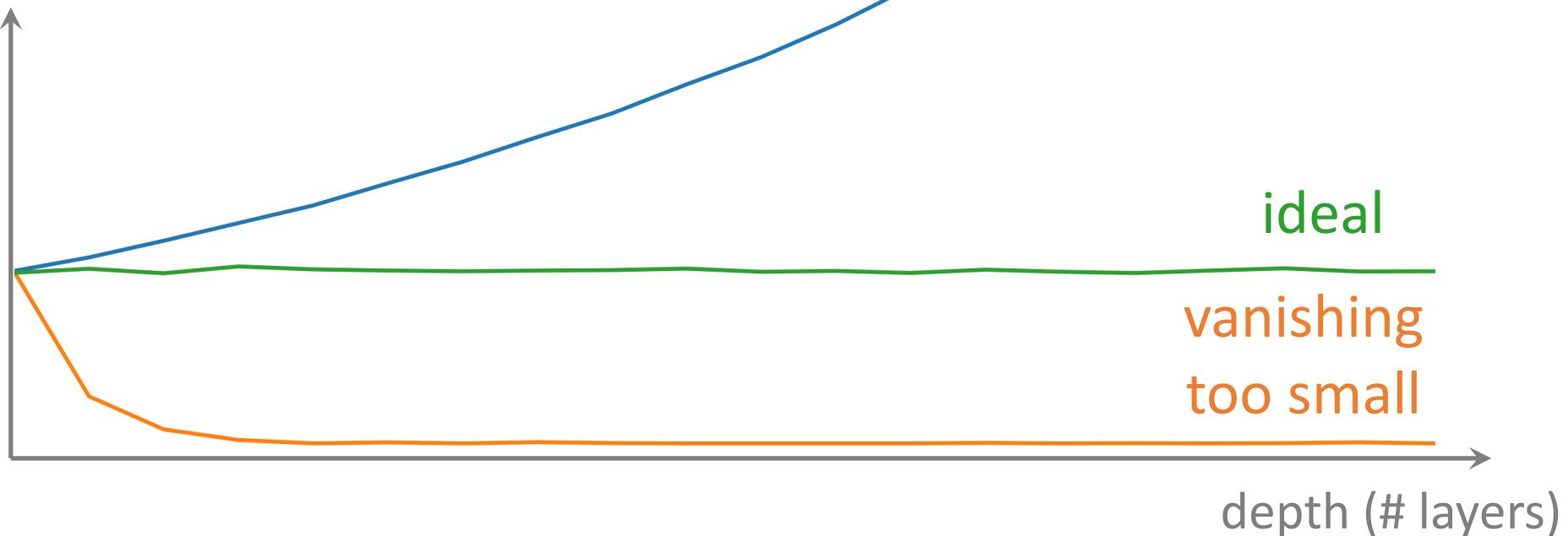
# **Weight Initialization**

# Weight initialization

- Good initialization has a significant effect

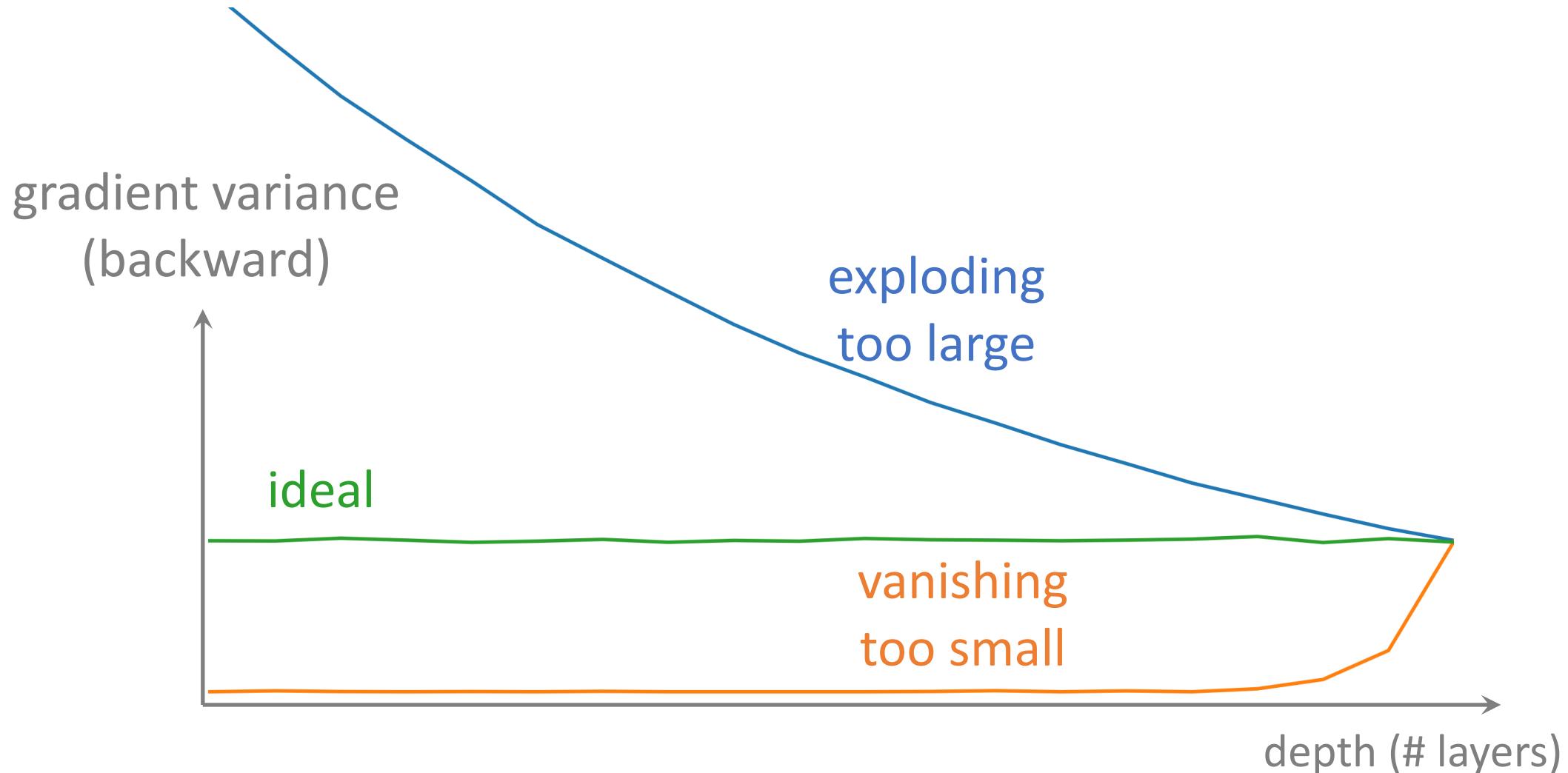
exploding:  
too big

signal variance  
(forward)

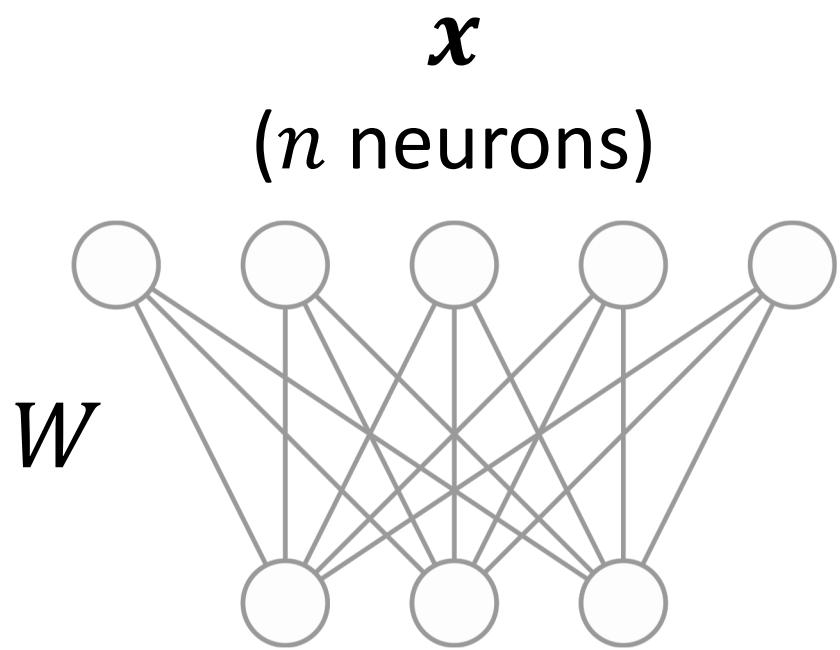


# Weight initialization

- Good initialization has a significant effect



# Weight initialization



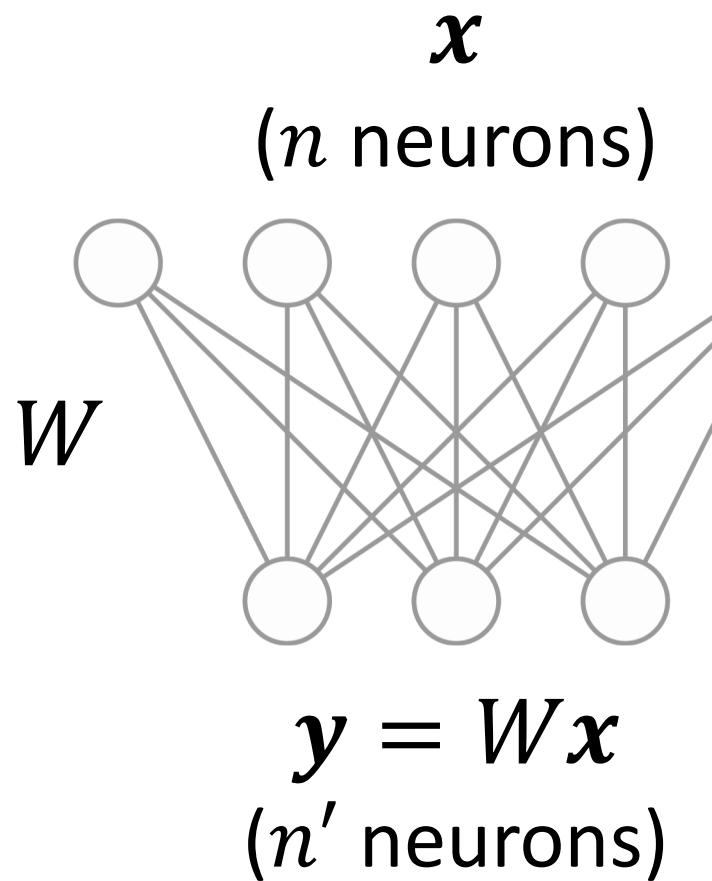
$$y = Wx$$

( $n'$  neurons)

- if linear activation
- one layer: variance scaled by

$$\text{Var}[y] = n \text{Var}[w] \text{Var}[x]$$

# Weight initialization



$$y_i = \sum_j W_{ij} x_j$$

- definition

$$\text{Var}[y_i] = \text{Var}\left[\sum_j W_{ij} x_j\right]$$

$$\text{Var}[y_i] = \sum_j \text{Var}[W_{ij} x_j]$$

- independence

$$\text{Var}[y_i] = \sum_j \text{Var}[W_{ij}] \text{Var}[x_j]$$

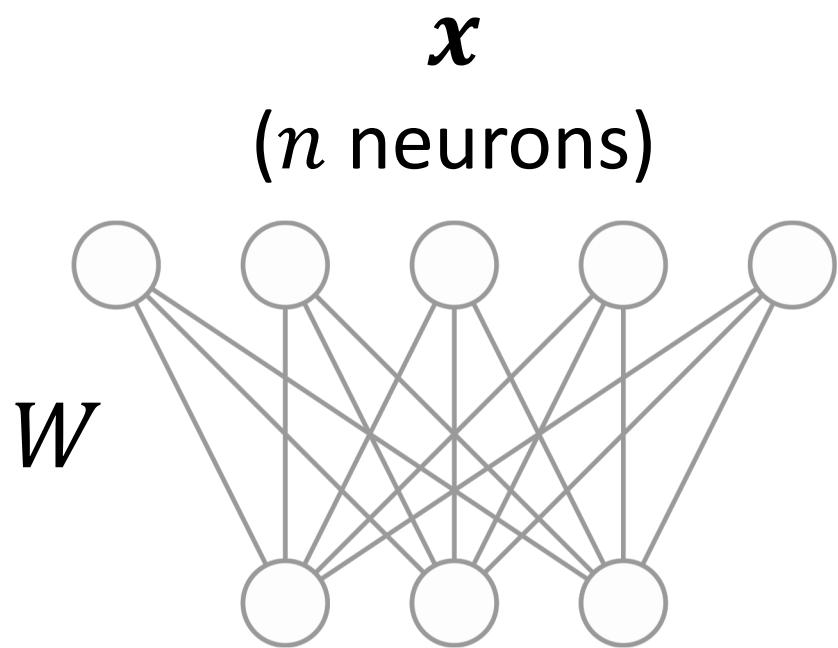
- independence & zero-mean

$$\text{Var}[y] = n \text{Var}[w] \text{Var}[x]$$

- identical distributions

derivation

# Weight initialization



$$y = Wx  
(n' \text{ neurons})$$

- if linear activation
- one layer: variance scaled by

$$\text{Var}[y] = n \text{Var}[w] \text{Var}[x]$$

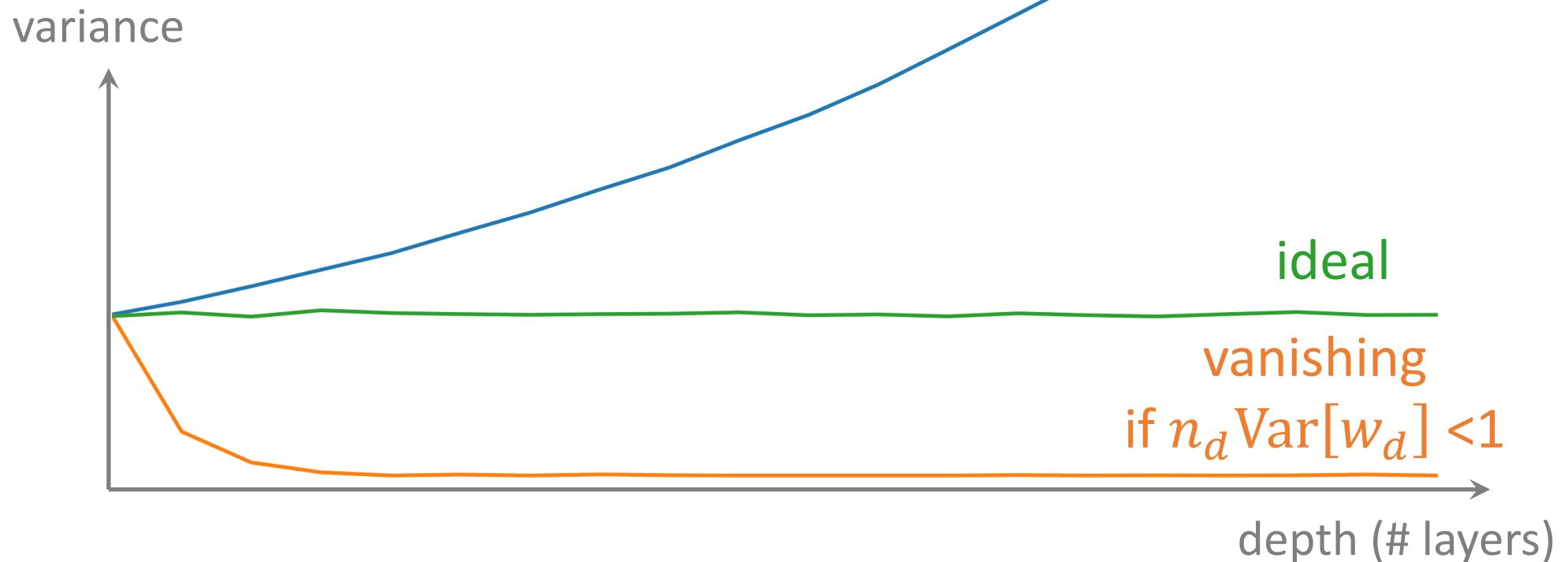
- many layers: variance scaled by

$$\text{Var}[y] = \prod_d n_d \text{Var}[w_d] \text{Var}[x]$$

# Variance Scaling: forward

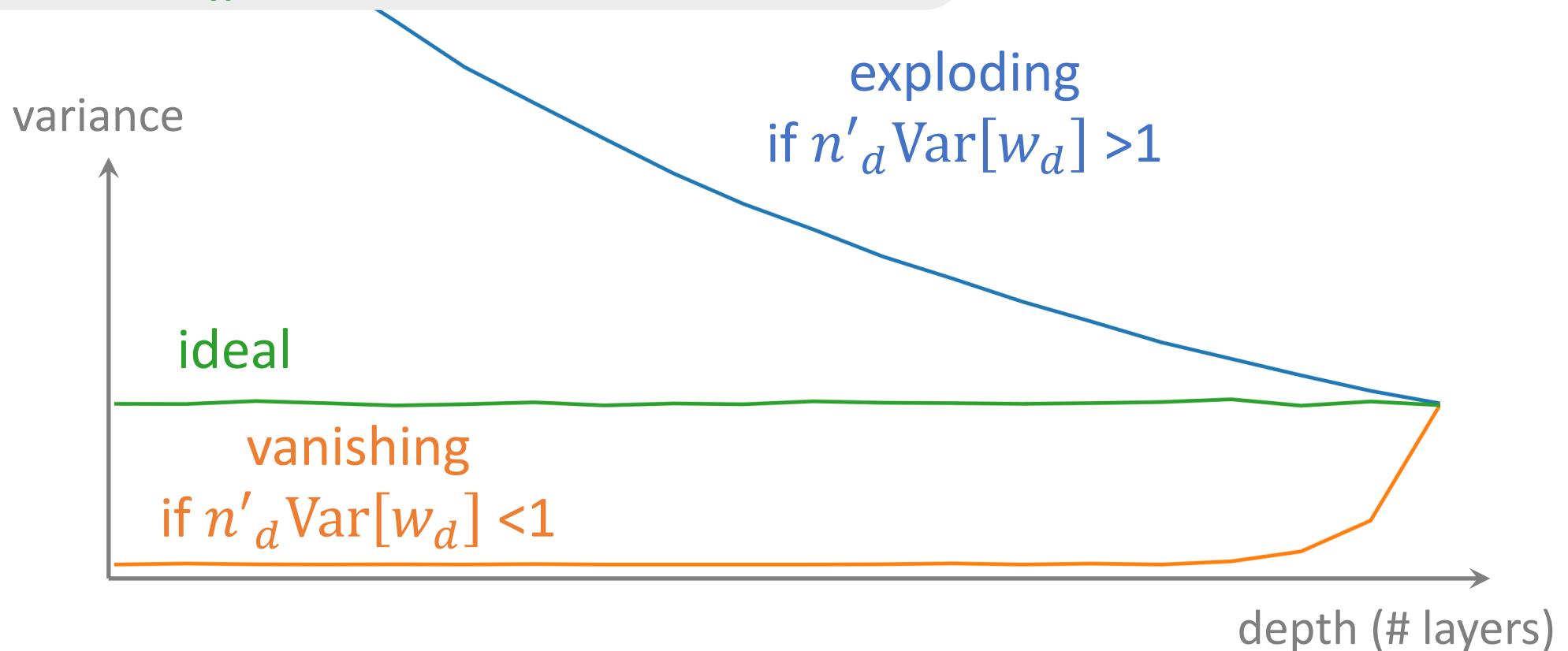
$$\text{Var}[y] = \prod_d n_d \text{Var}[w_d] \text{Var}[x]$$

exploding  
if  $n_d \text{Var}[w_d] > 1$



# Variance Scaling: backward

$$\text{Var} \left[ \frac{\partial \mathcal{E}}{\partial x} \right] = \prod_d n'_d \text{Var}[w_d] \text{Var} \left[ \frac{\partial \mathcal{E}}{\partial y} \right]$$



# Xavier initialization: `torch.nn.init.xavier_normal_`

forward:

$$n\text{Var}[w] = 1$$

backward:

$$n'\text{Var}[w] = 1$$

- Gaussian distribution:

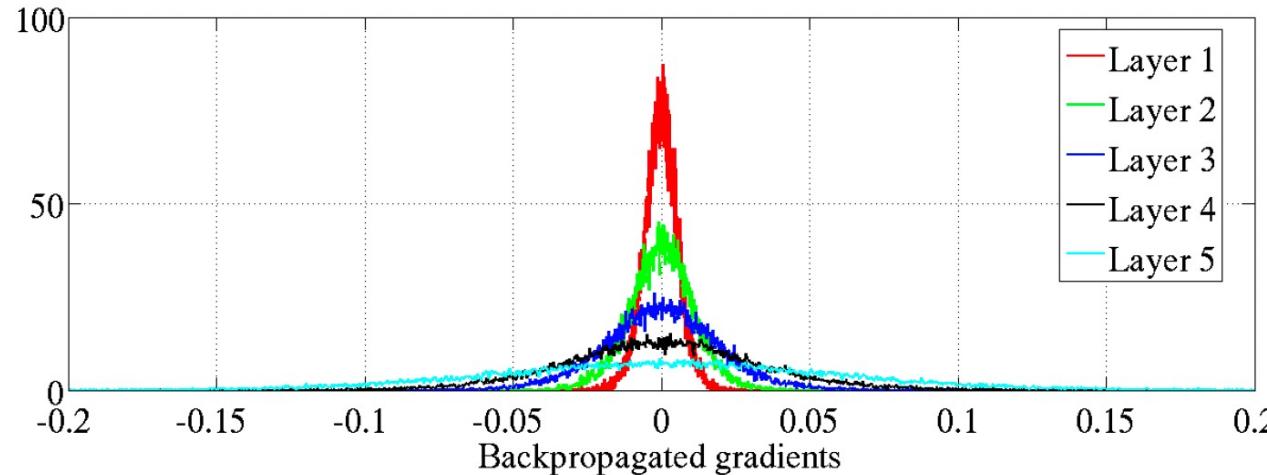
$$w \sim \mathcal{N}(\mu = 0, \sigma = \sqrt{1/n})$$

- Uniform distribution:

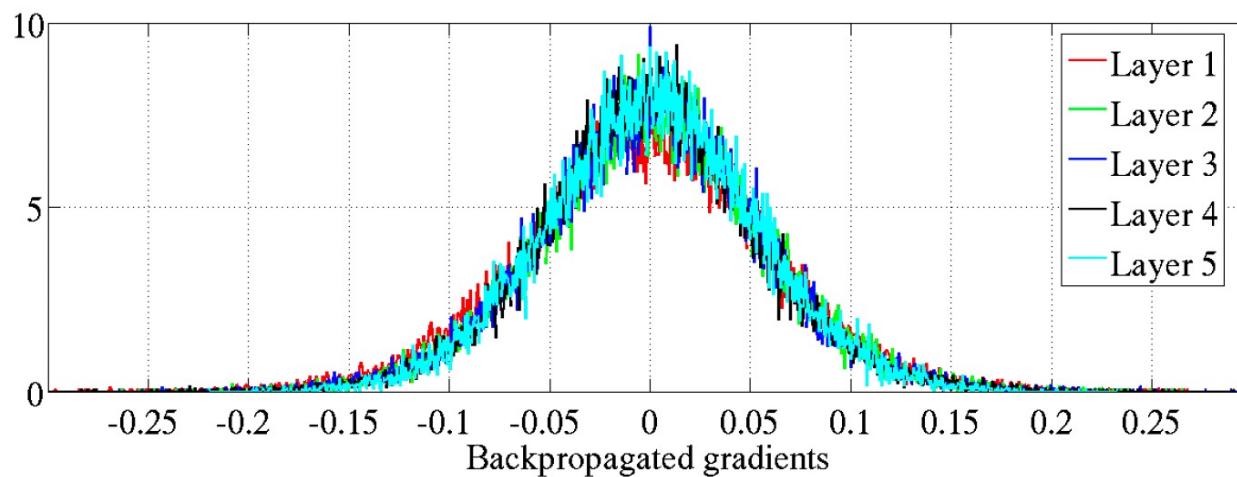
$$w \sim \mathcal{U}(-a, +a), a = \sqrt{3/n}$$

- Consider forward and backward:  
replace  $n$  with  $(n + n')/2$

# Xavier initialization: `torch.nn.init.xavier_normal_`



poor initialization:  
earlier layer has smaller gradients



Xavier initialization:  
all layers have similar gradient scale

# What are $n$ and $n'$ in ConvNet?

The number of connections feeding into/out of a node:

$n$  (“fan\_in”):  $C_{in} \times K_h \times K_w$

$n'$  (“fan\_out”):  $C_{out} \times K_h \times K_w$

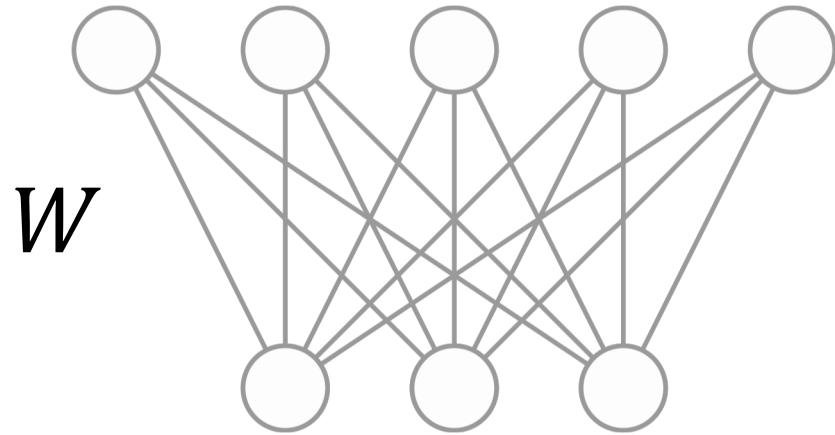
$K_h \times K_w$ : kernel size

fully-connected: viewed as  $H \times W$  conv

# Weight initialization: ReLU

$$\mathbf{x}' = \text{ReLU}(\mathbf{x})$$

( $n$  neurons)



$$\mathbf{y} = W\mathbf{x}'$$

( $n'$  neurons)

- if ReLU activation
- one layer: variance scaled by

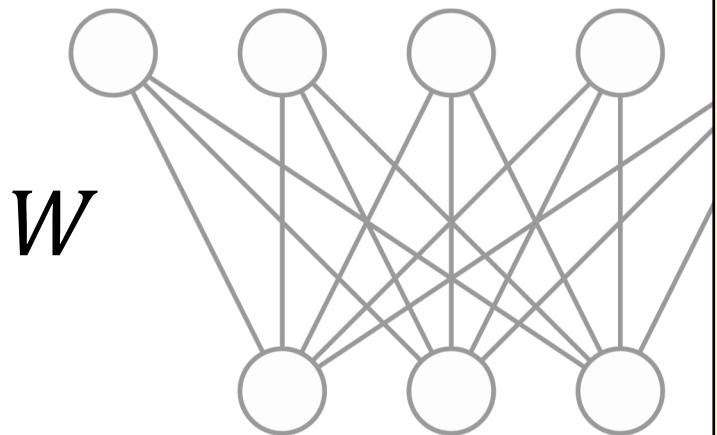
$$\text{Var}[y] = \frac{1}{2} n \text{Var}[w] \text{Var}[x]$$

# Weight initialization

derivation

$$\mathbf{x}' = \text{ReLU}(\mathbf{x})$$

( $n$  neurons)



$$\mathbf{y} = W\mathbf{x}'$$

( $n'$  neurons)

$$y_i = \sum_j W_{ij} \mathbf{x}'_j$$

- $\mathbf{x}' = \text{ReLU}(\mathbf{x})$

$$\text{Var}[y_i] = \text{Var}\left[\sum_j W_{ij} \mathbf{x}'_j\right]$$

$$\text{Var}[y_i] = \sum_j \text{Var}[W_{ij} \mathbf{x}'_j]$$

$$\text{Var}[y_i] = \sum_j \text{Var}[W_{ij}] \mathbf{E}[\mathbf{x}'^2_j]$$

$$\text{Var}[y] = n \text{Var}[w] \frac{1}{2} \text{Var}[x]$$

$$\begin{aligned}\text{Var}[wx'] &= E[(wx')^2] - (E[wx'])^2 \\ &= E[(wx')^2] - (E[w]E[x])^2 \\ &= E[w^2]E[x'^2] = \text{Var}[w]\text{Var}[x']\end{aligned}$$

if  $w$  is zero-mean but  $x'$  is not.

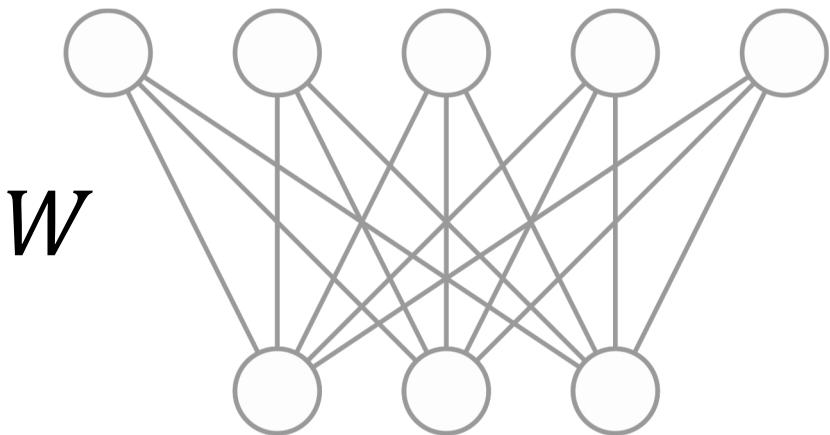
$$\mathbf{E}[x'^2] = \text{Var}[x]/2$$

if  $x$  is zero-mean and symmetric

# Weight initialization: ReLU

$$\mathbf{x}' = \text{ReLU}(\mathbf{x})$$

( $n$  neurons)



$$\mathbf{y} = W\mathbf{x}'$$

( $n'$  neurons)

- if ReLU activation
- one layer: variance scaled by

$$\text{Var}[y] = \frac{1}{2} n \text{Var}[w] \text{Var}[x]$$

- many layers: variance scaled by

$$\text{Var}[y] = \prod_d \frac{1}{2} n_d \text{Var}[w_d] \text{Var}[x]$$

# Kaiming initialization: `torch.nn.init.kaiming_normal_`

forward:

$$\frac{1}{2} n \text{Var}[w] = 1$$

backward:

$$\frac{1}{2} n' \text{Var}[w] = 1$$

- Gaussian distribution:

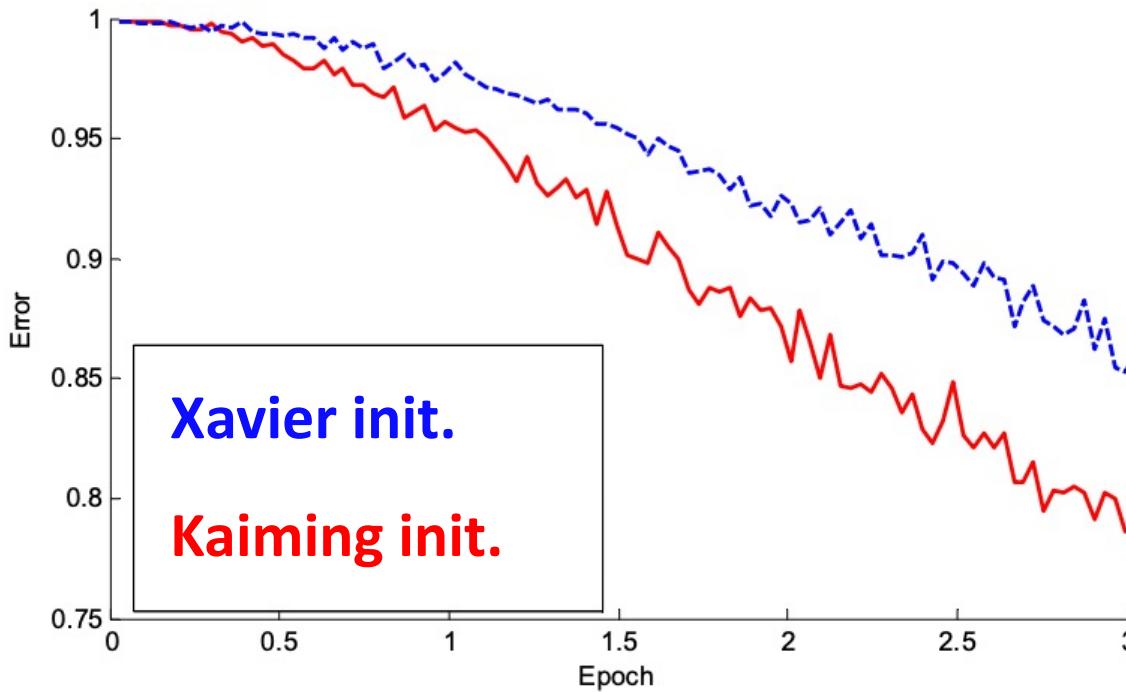
$$w \sim \mathcal{N}(\mu = 0, \sigma = \sqrt{2/n})$$

- Uniform distribution:

$$w \sim \mathcal{U}(-a, +a), a = \sqrt{6/n}$$

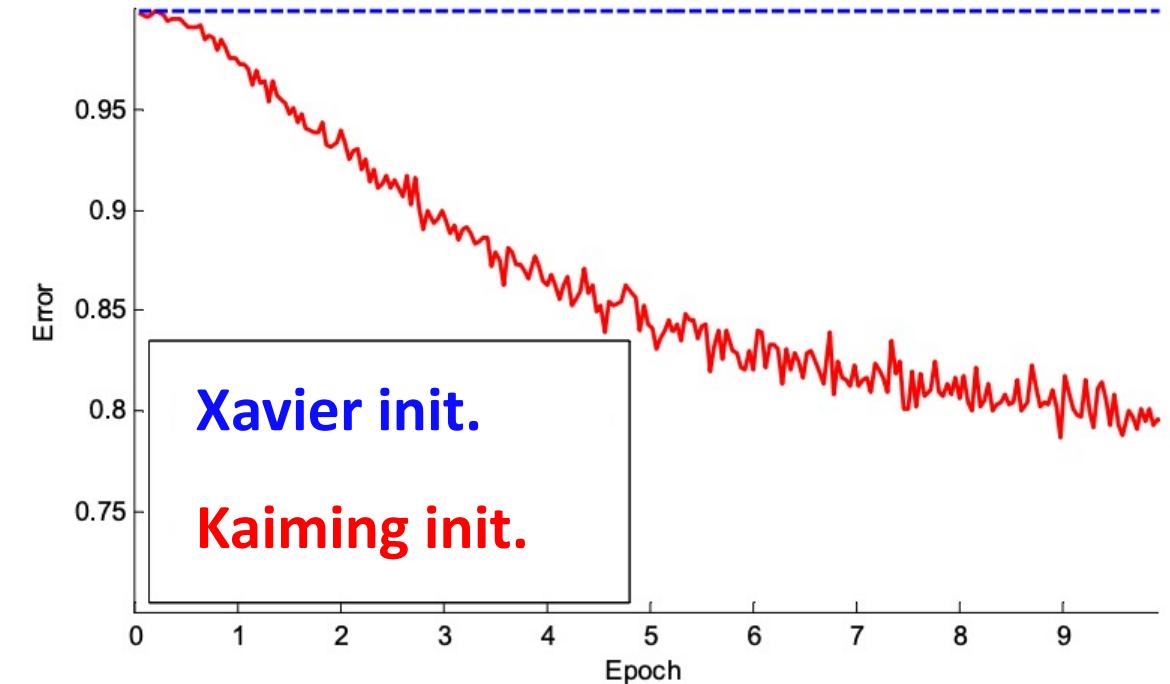
- sufficient to use  $n$  or  $n'$

# Kaiming initialization: `torch.nn.init.kaiming_normal_`



Xavier init.  
Kaiming init.

**22-layer VGG w/ ReLU :**  
better init converges faster



Xavier init.  
Kaiming init.

**30-layer VGG w/ ReLU:**  
better init enables training

# Weight initialization: Takeaways

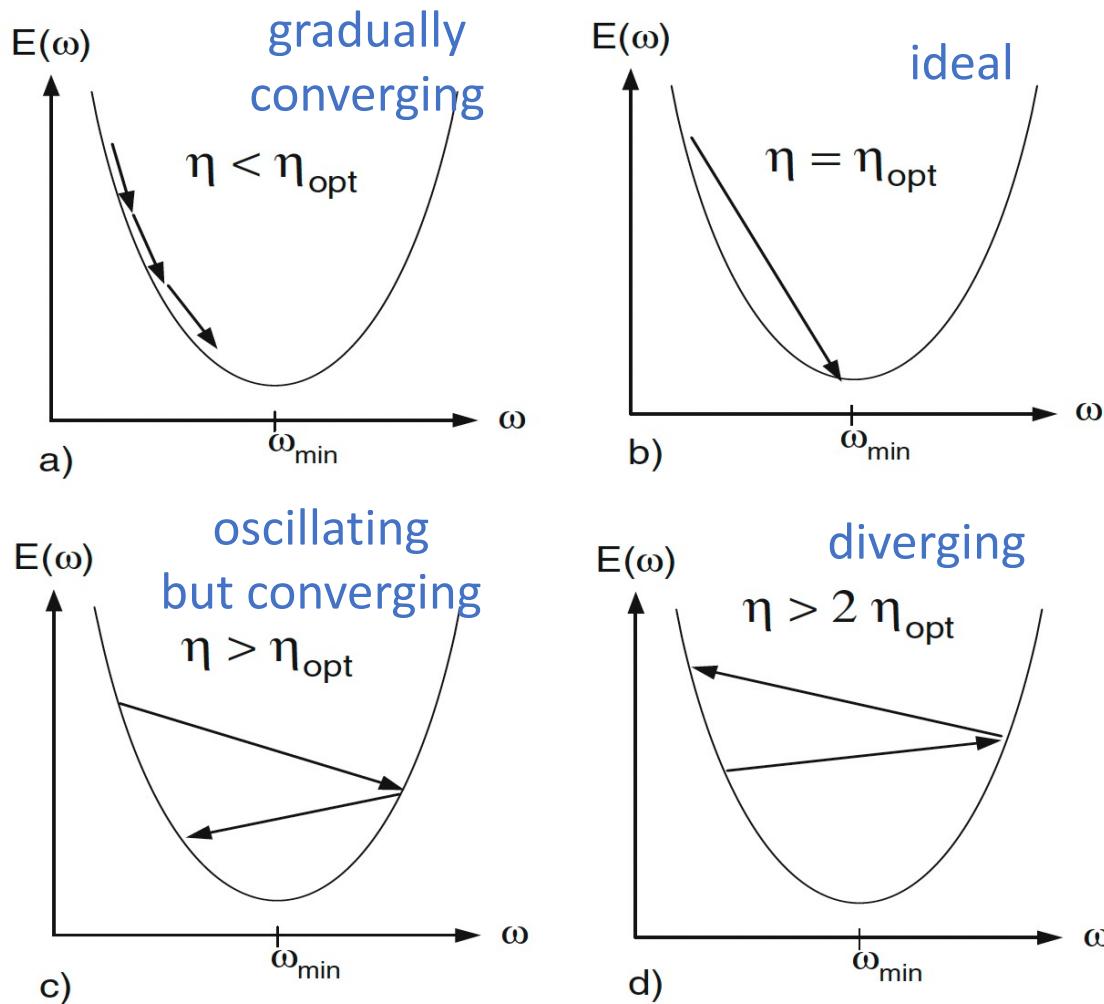
- Signals can be exploding/vanishing
- Product of scaling factors
- Initialization is driven by normalization
- Initialization is based on strong assumptions
- Initialization is necessary, but often not sufficient

# **Normalization Modules**

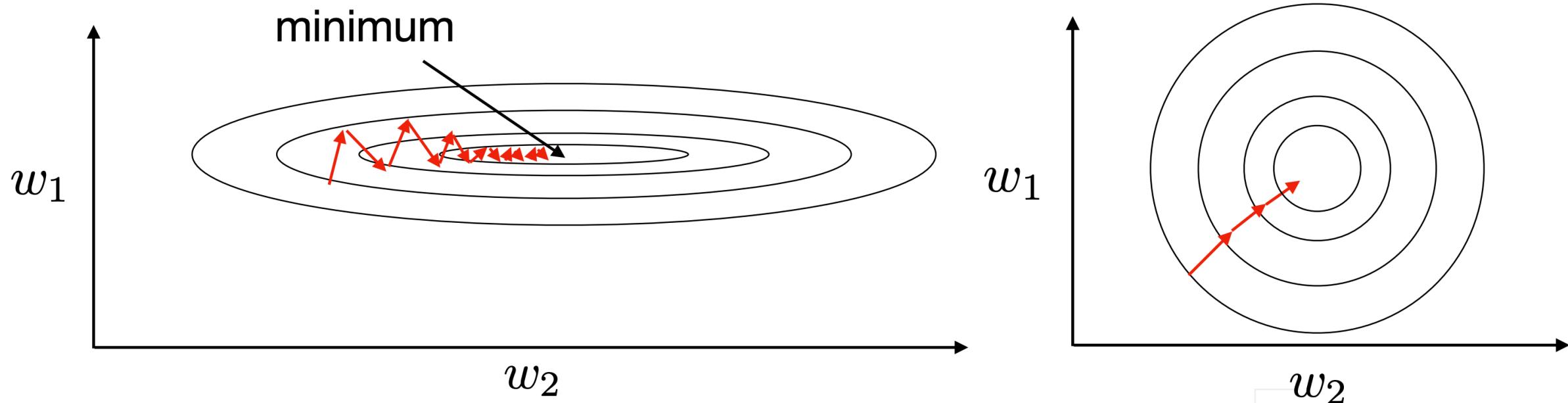
# Learning Rate in Gradient Descent

$$W := W - \eta \frac{\partial E}{\partial W}$$

$\eta$ : learning rate



# Normalization



- Same learning rate applied to all weights
- Large weights dominate updates
- Small weights oscillate (or diverge)
- Similar pace for all weights

# Input Normalization

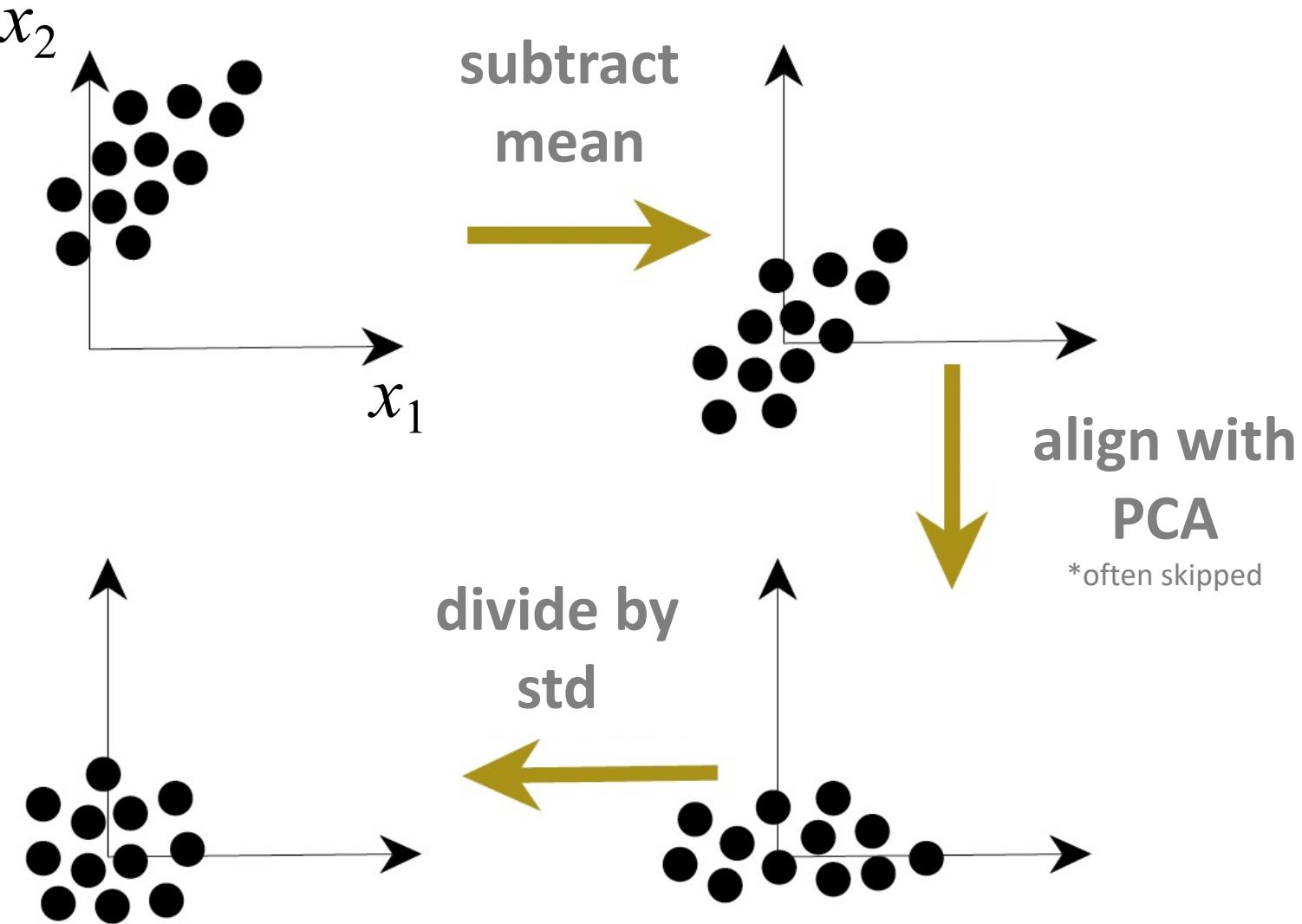
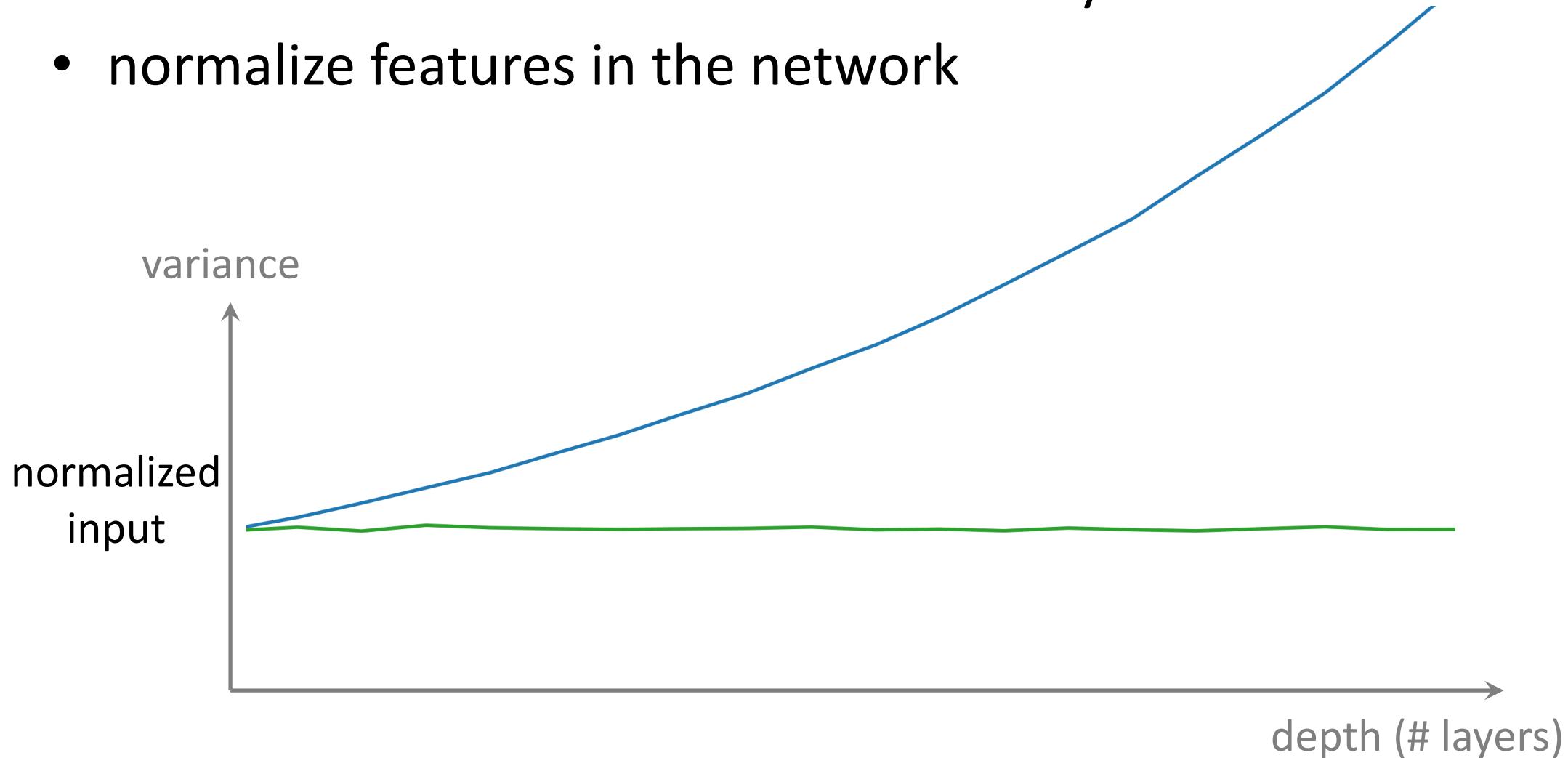


Figure adapted from: LeCun et al. "Efficient BackProp". 1998.

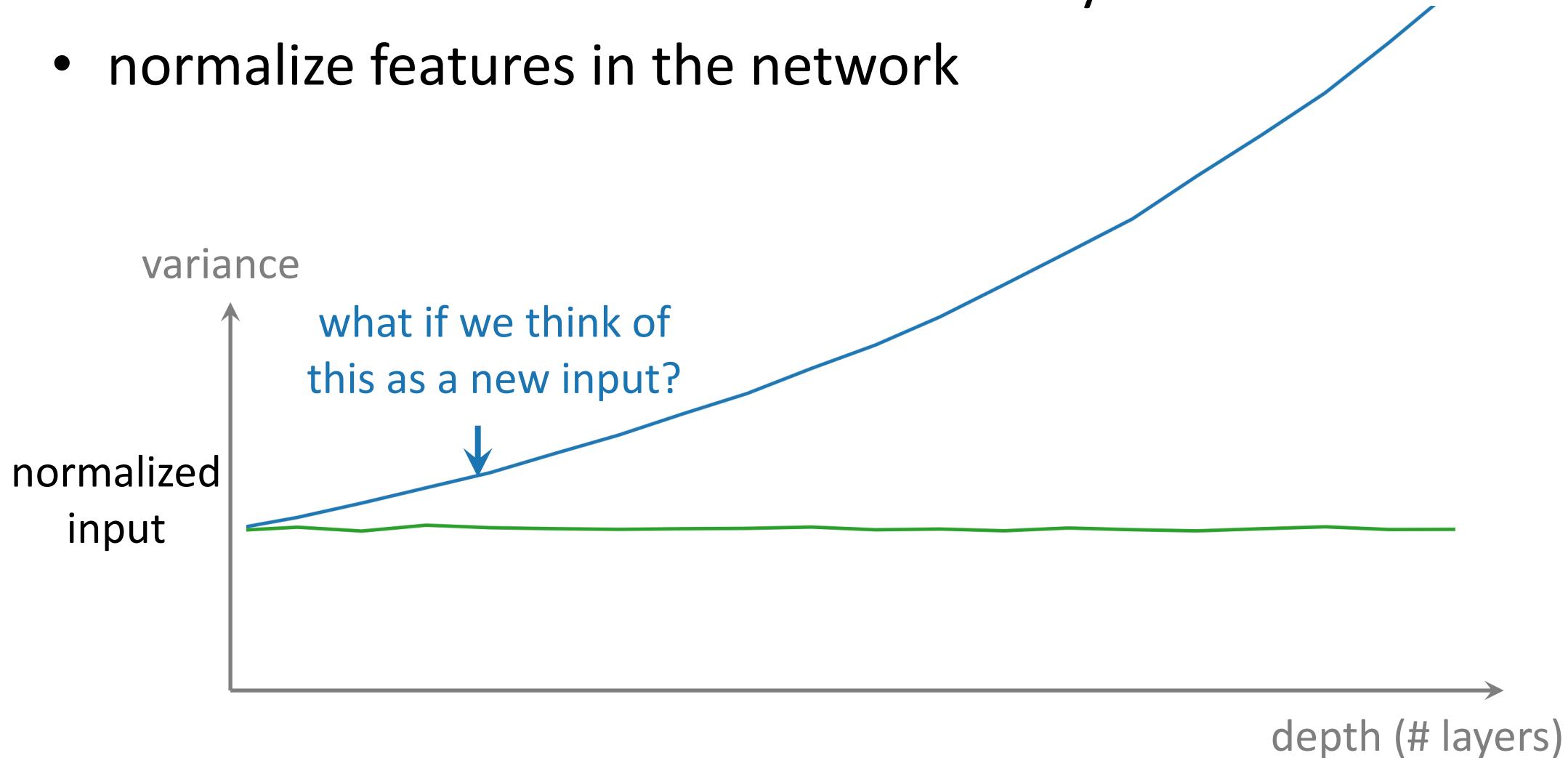
# Normalization Modules

- We want to maintain variance for all layers
- normalize features in the network



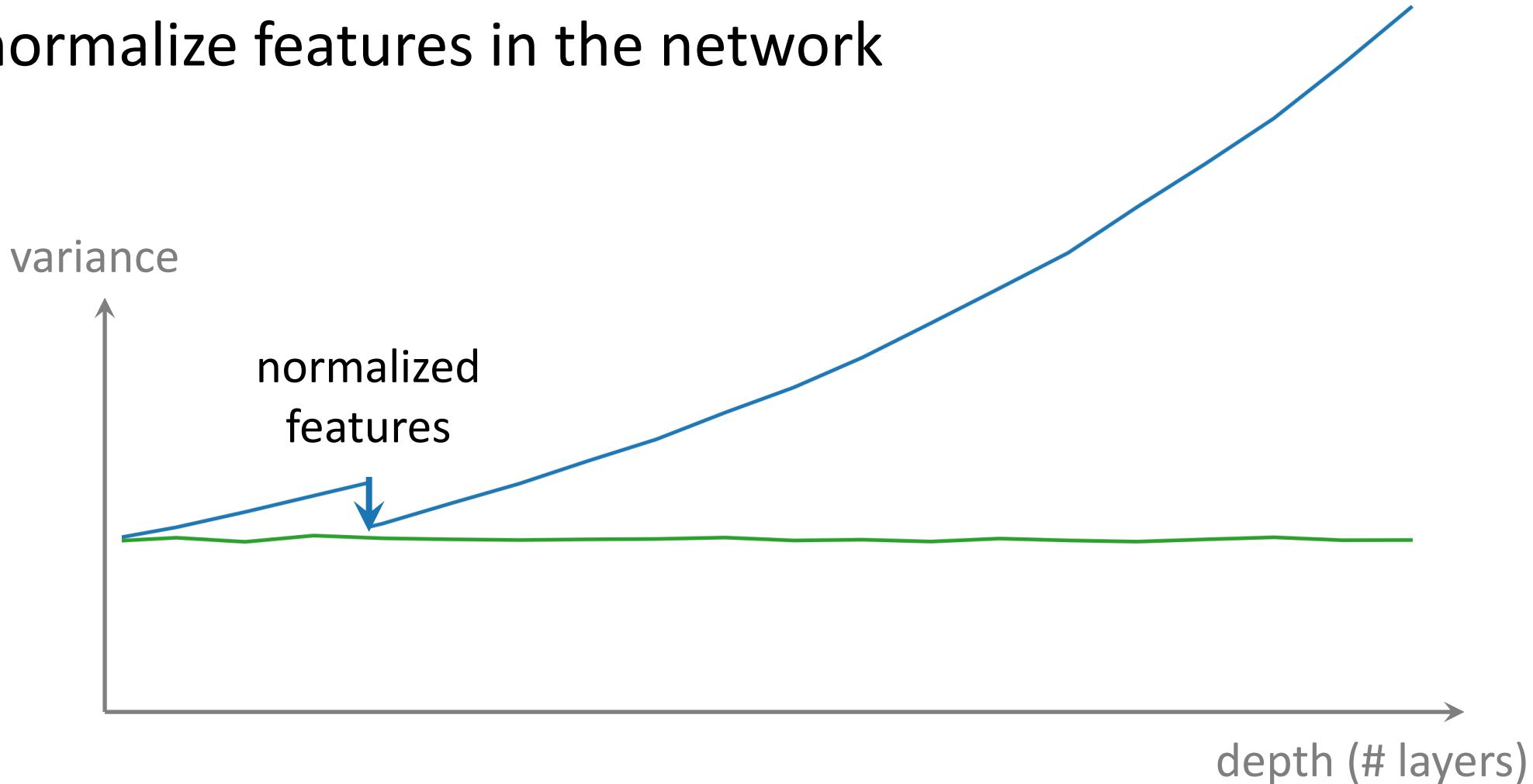
# Normalization Modules

- We want to maintain variance for all layers
- normalize features in the network



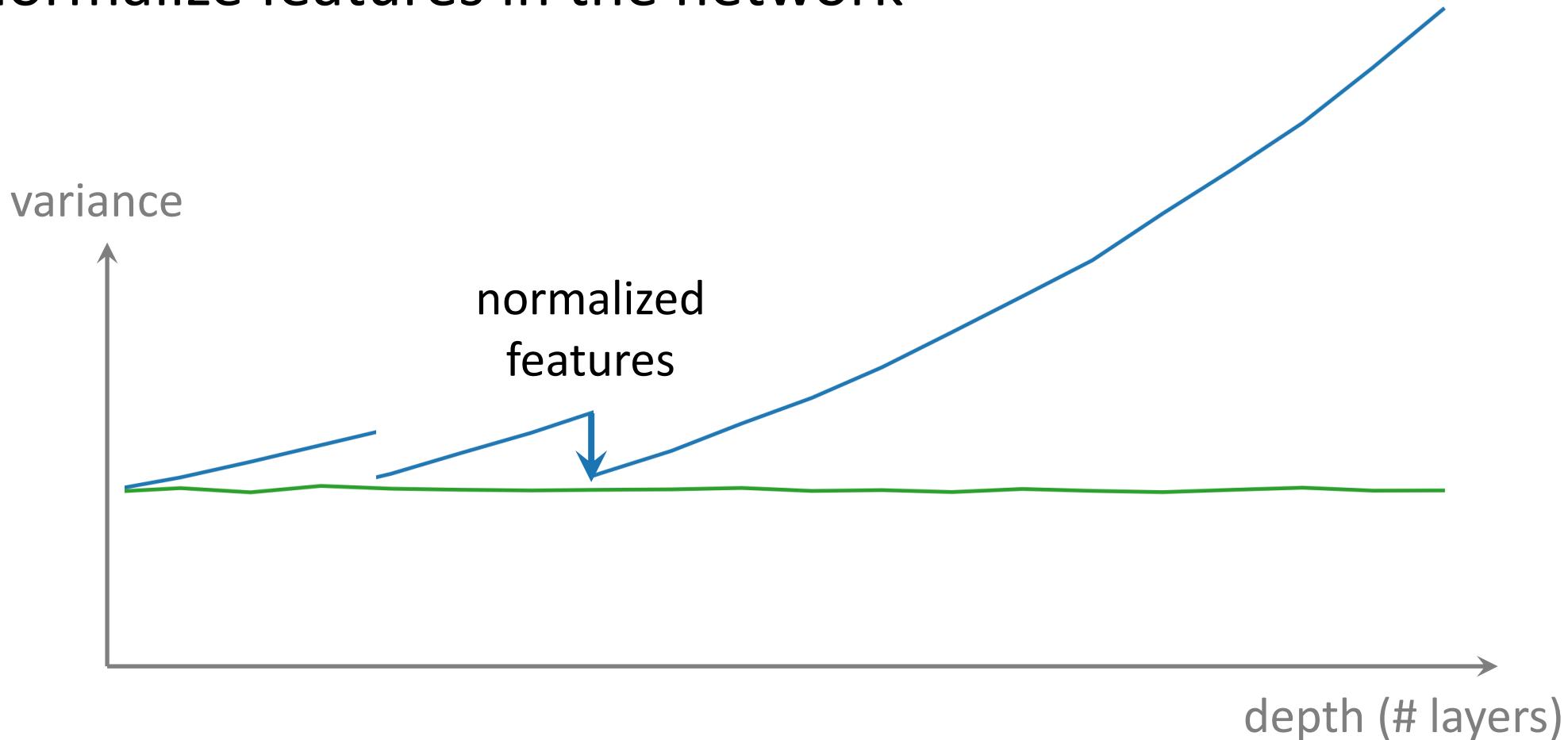
# Normalization Modules

- We want to maintain variance for all layers
- normalize features in the network



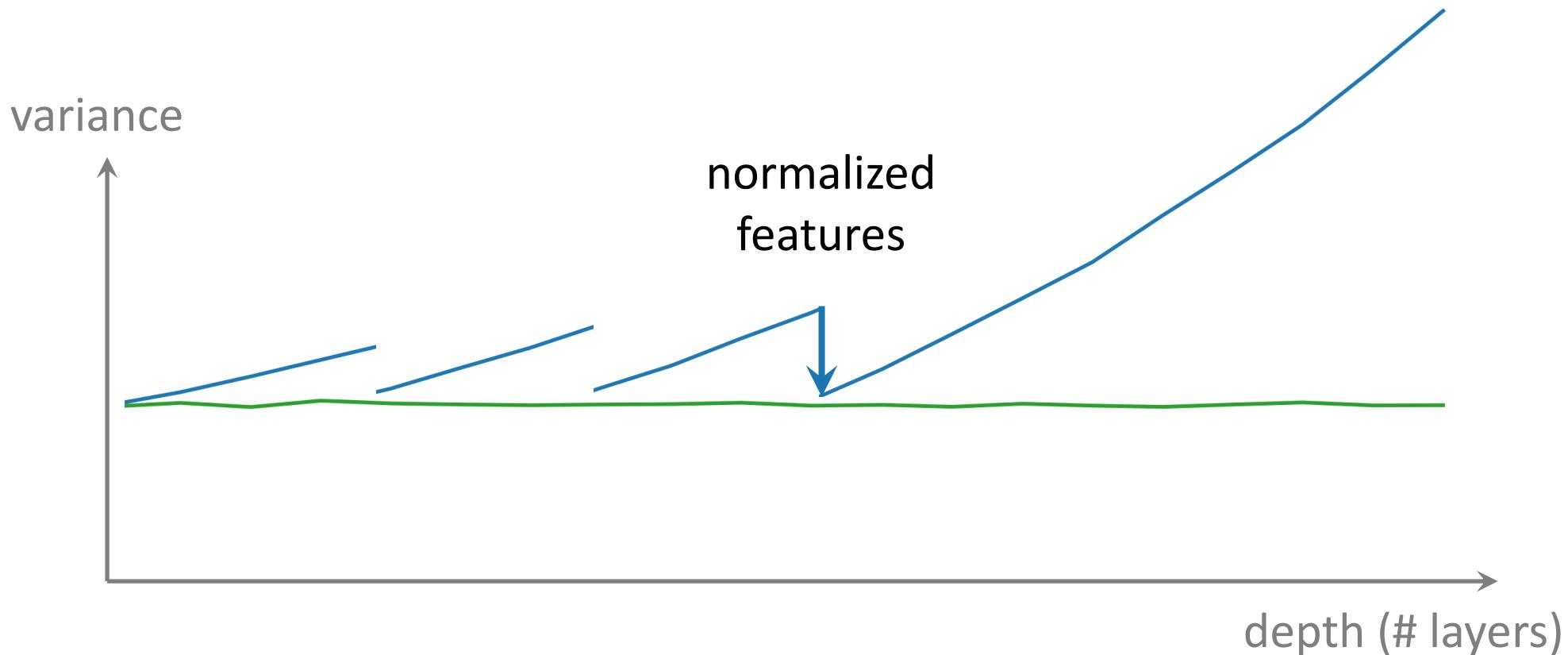
# Normalization Modules

- We want to maintain variance for all layers
- normalize features in the network



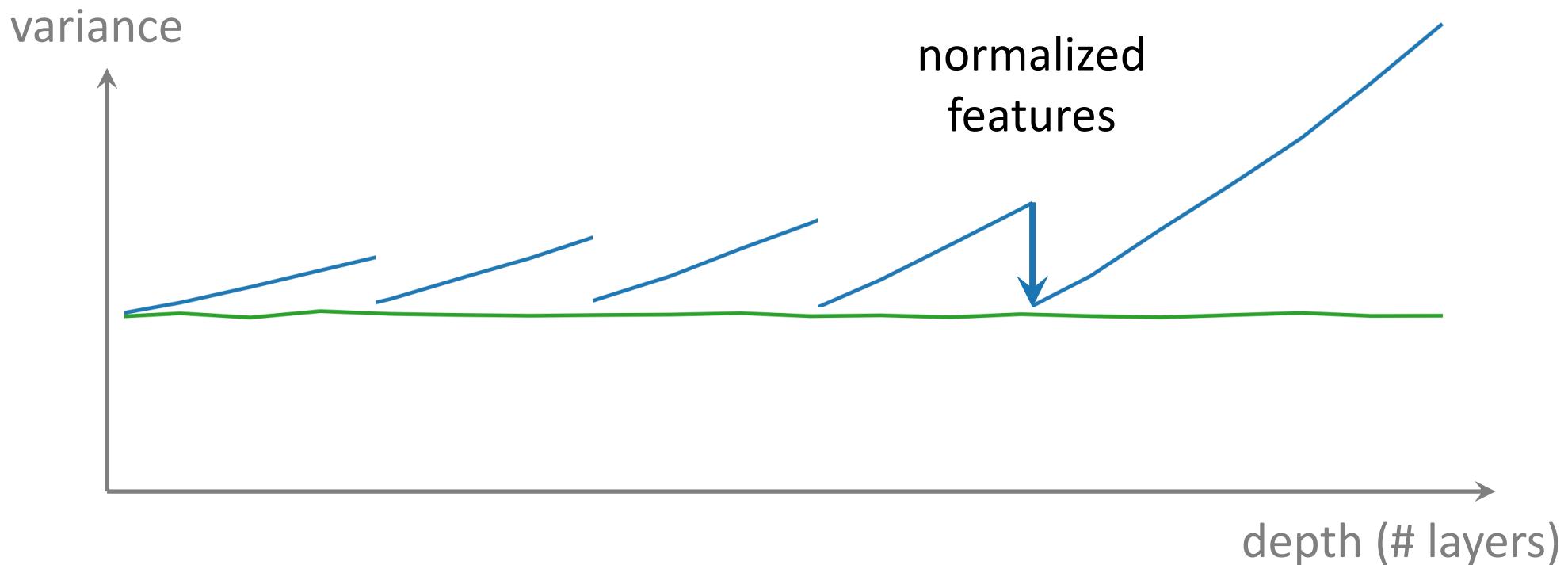
# Normalization Modules

- We want to maintain variance for all layers
- normalize features in the network



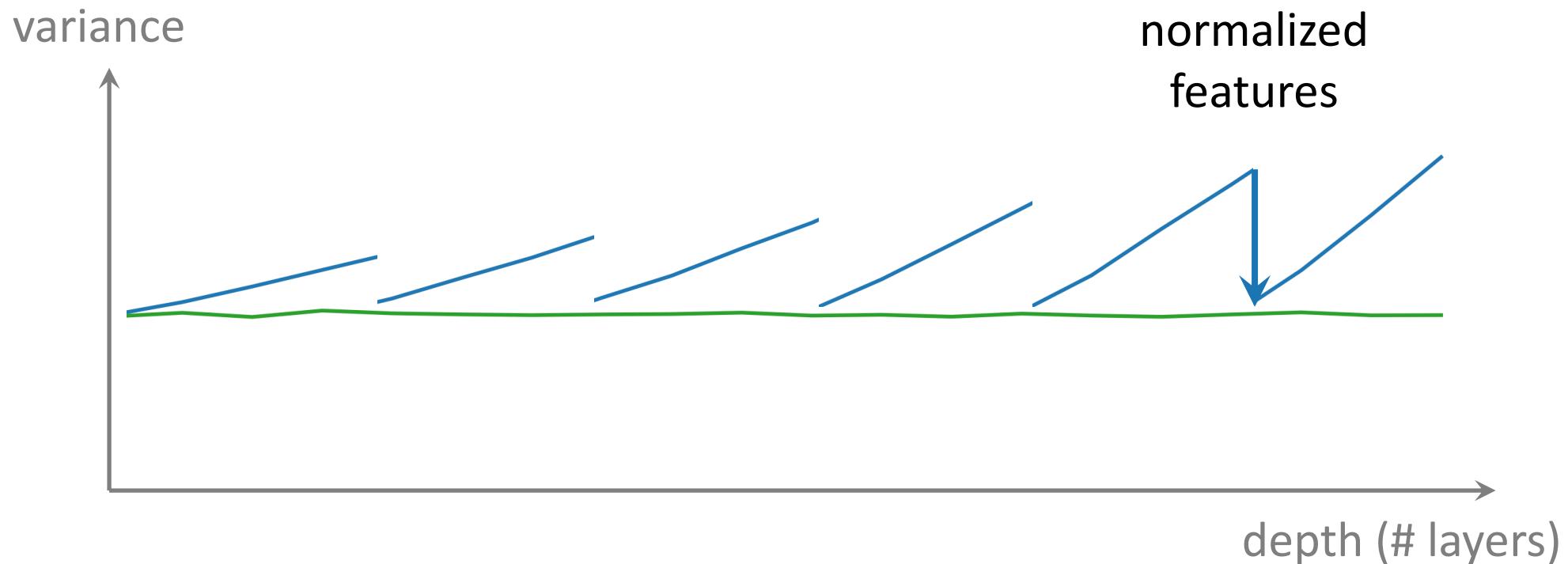
# Normalization Modules

- We want to maintain variance for all layers
- normalize features in the network



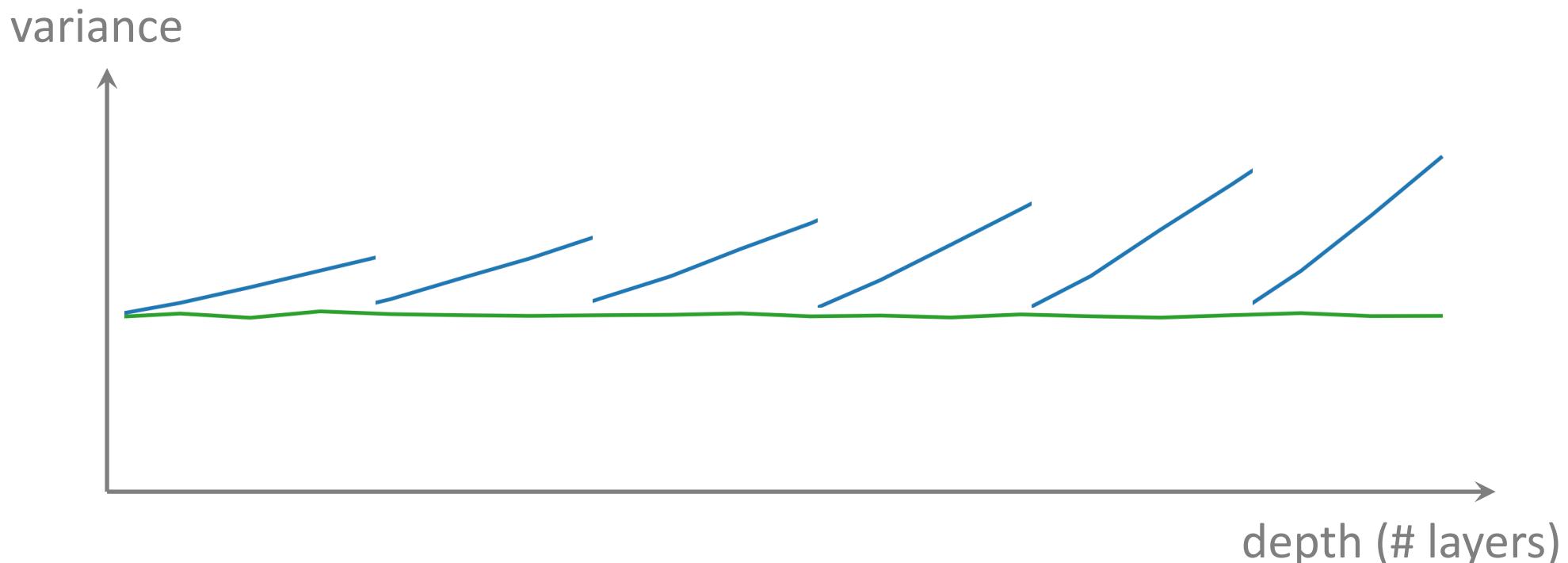
# Normalization Modules

- We want to maintain variance for all layers
- normalize features in the network



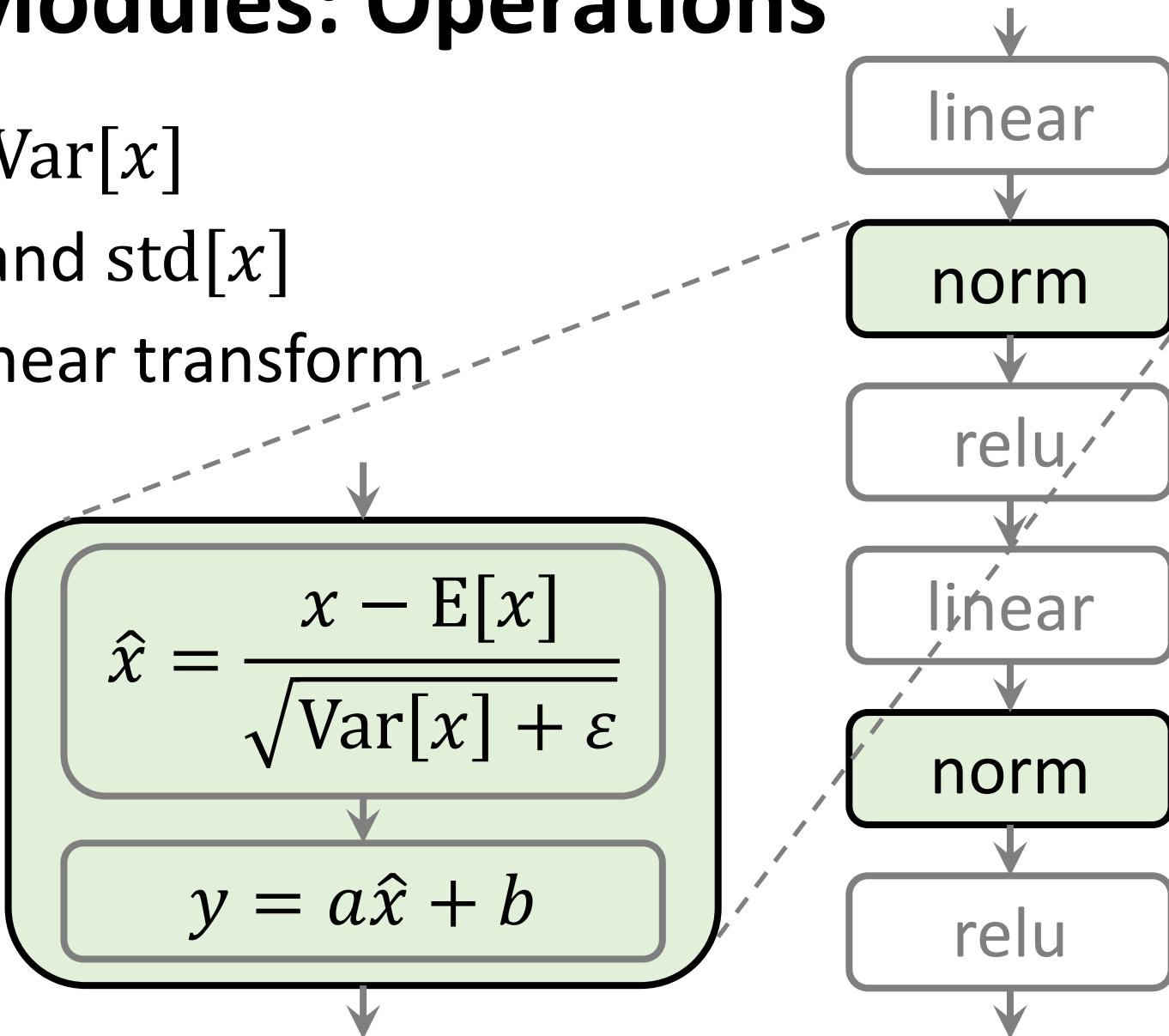
# Normalization Modules

- We want to maintain variance for all layers
- normalize features in the network
- train end-to-end by BackProp



# Normalization Modules: Operations

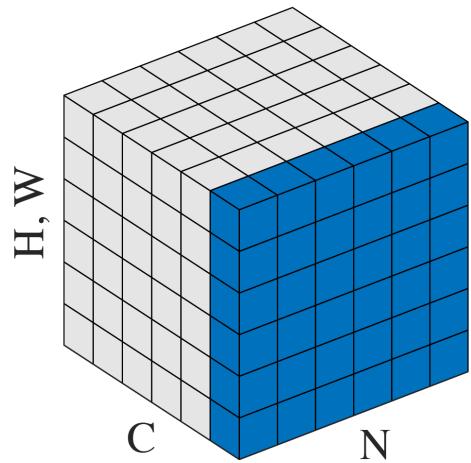
1. compute  $E[x]$  and  $\text{Var}[x]$
2. normalize by  $E[x]$  and  $\text{std}[x]$
3. compensate by a linear transform



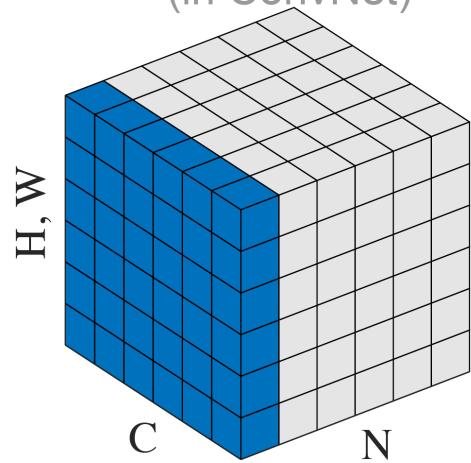
# Normalization Modules: Variants

differ in support sets of  $E[x]$ ,  $\text{Var}[x]$

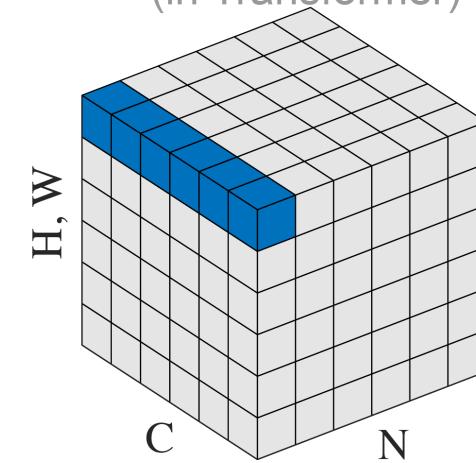
BatchNorm



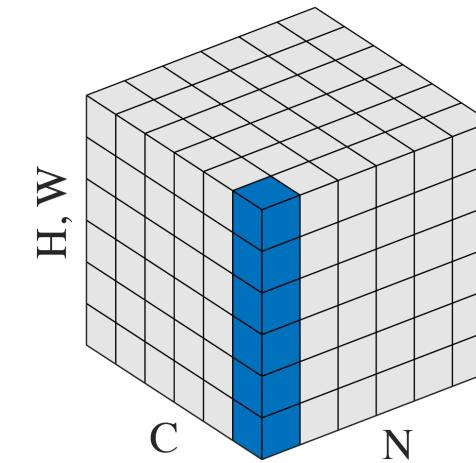
LayerNorm  
(in ConvNet)



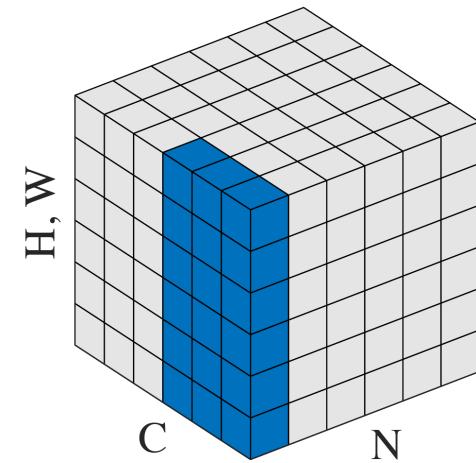
LayerNorm  
(in Transformer)



InstanceNorm



GroupNorm

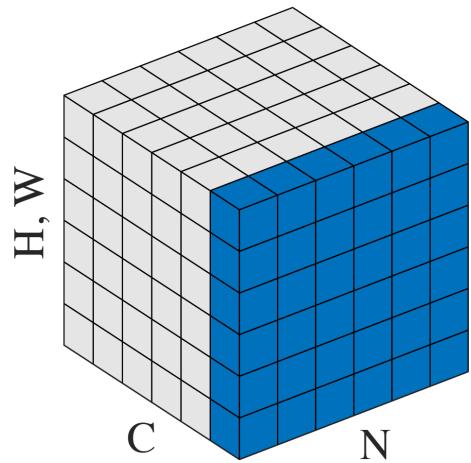


# Case Study: Batch Normalization (BN)

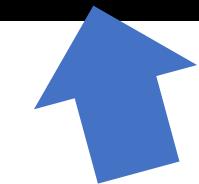
$E[x]$  over N, H, W axes

Einstein sum:  $\sum_n \sum_h \sum_w x_{nchw}$

BatchNorm



```
mean = torch.einsum('nchw->c', x) / (N * H * W)
```



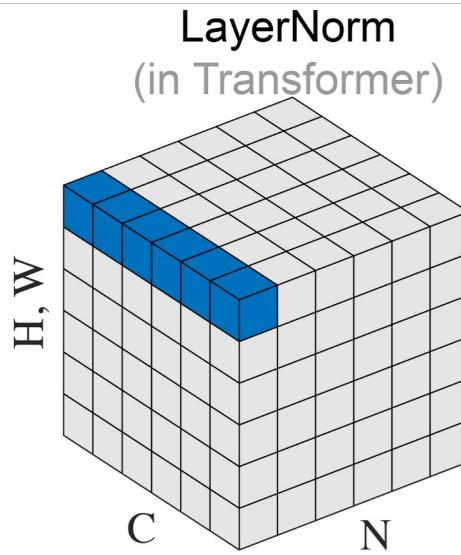
C-dim vector

BN at inference time

- batch (N axis) is not legitimate
- maintain  $\mu$  and  $\sigma^2$  in train set (moving average)
- use  $\mu$  and  $\sigma^2$  at inference time
- tl;dr: set `model.eval()` in torch

# Case Study: Layer Normalization (LN)

$E[x]$  over C axis



Einstein sum:  $\sum_c x_{nchw}$

```
mean = torch.einsum('nchw->nhw', x) / C
```

NxHxW tensor

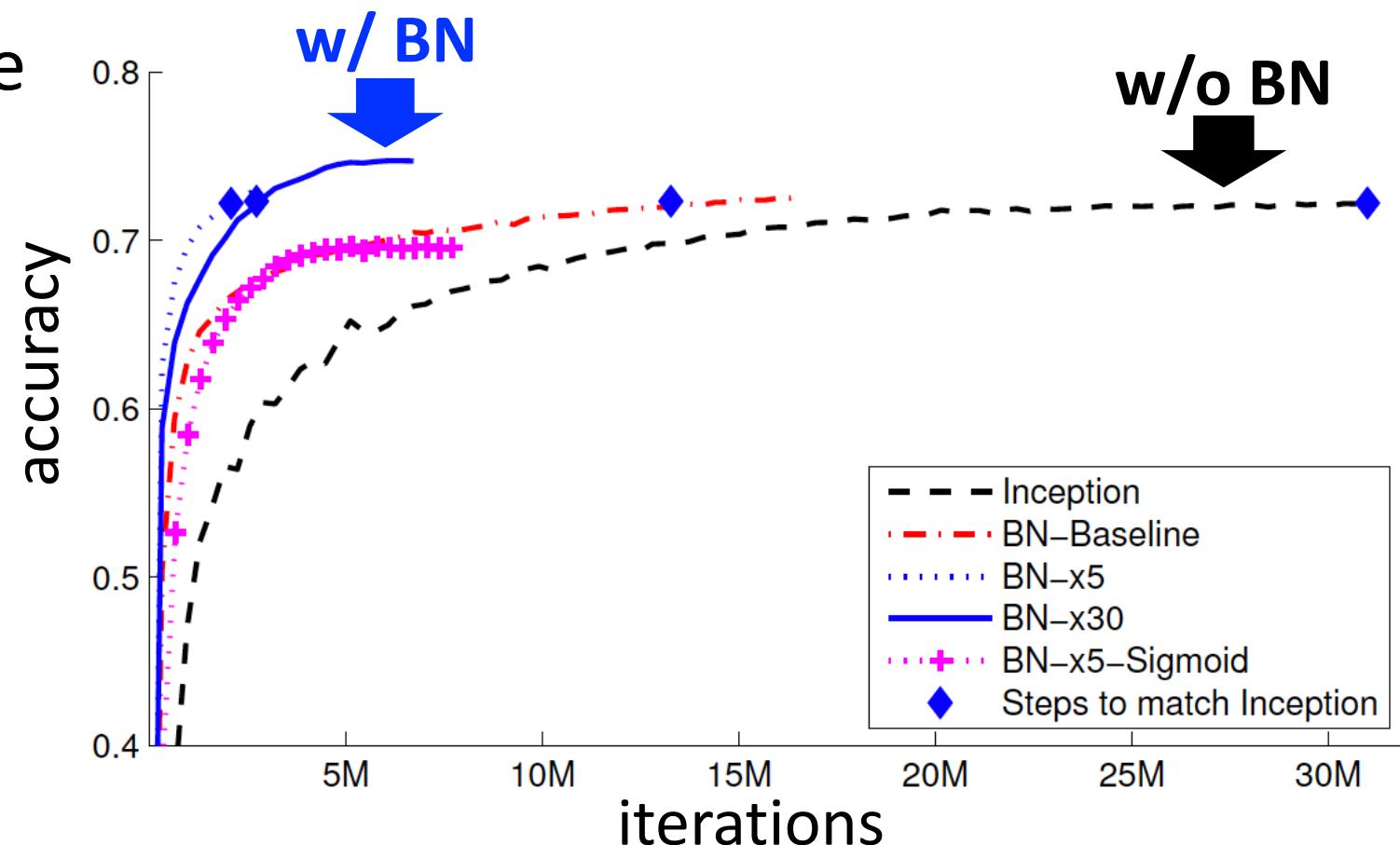


- same operation at train and inference
- popular in Transformers and sequences

\*More fun with Einstein sum:  
<https://rockt.github.io/2018/04/30/einsum>

# Normalization Modules: Effects

- Enable training models that are otherwise not trainable
- Speed up convergence
- Improve accuracy



# **Deep Residual Learning**

# Deep Residual Learning

- Deep Learning gets way deeper
- simple component: identity shortcut
- enable networks w/ hundreds of layers

**Compose simple modules into complex functions**

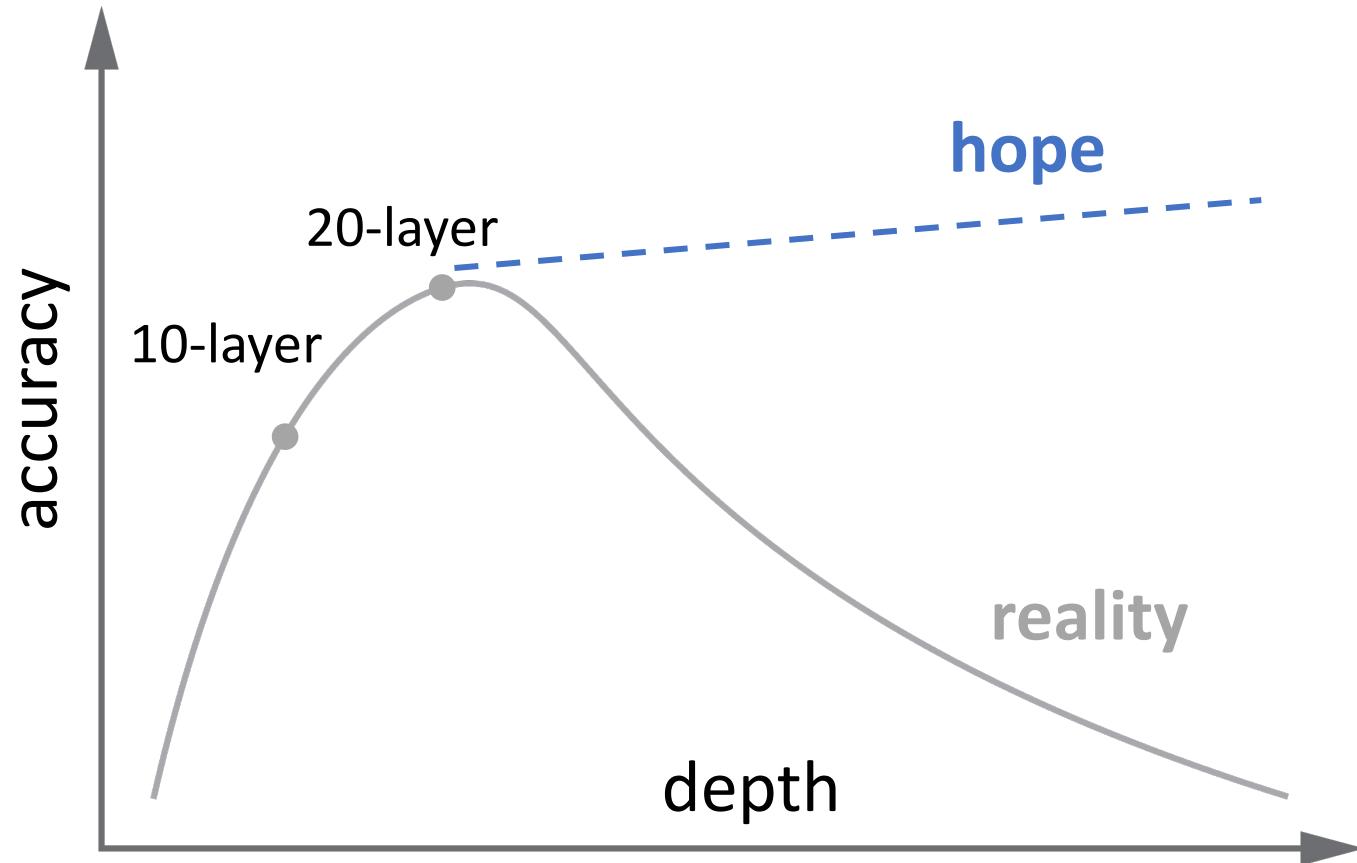


I WAS WINNING  
IMAGENET

UNTIL A  
DEEPER MODEL  
CAME ALONG

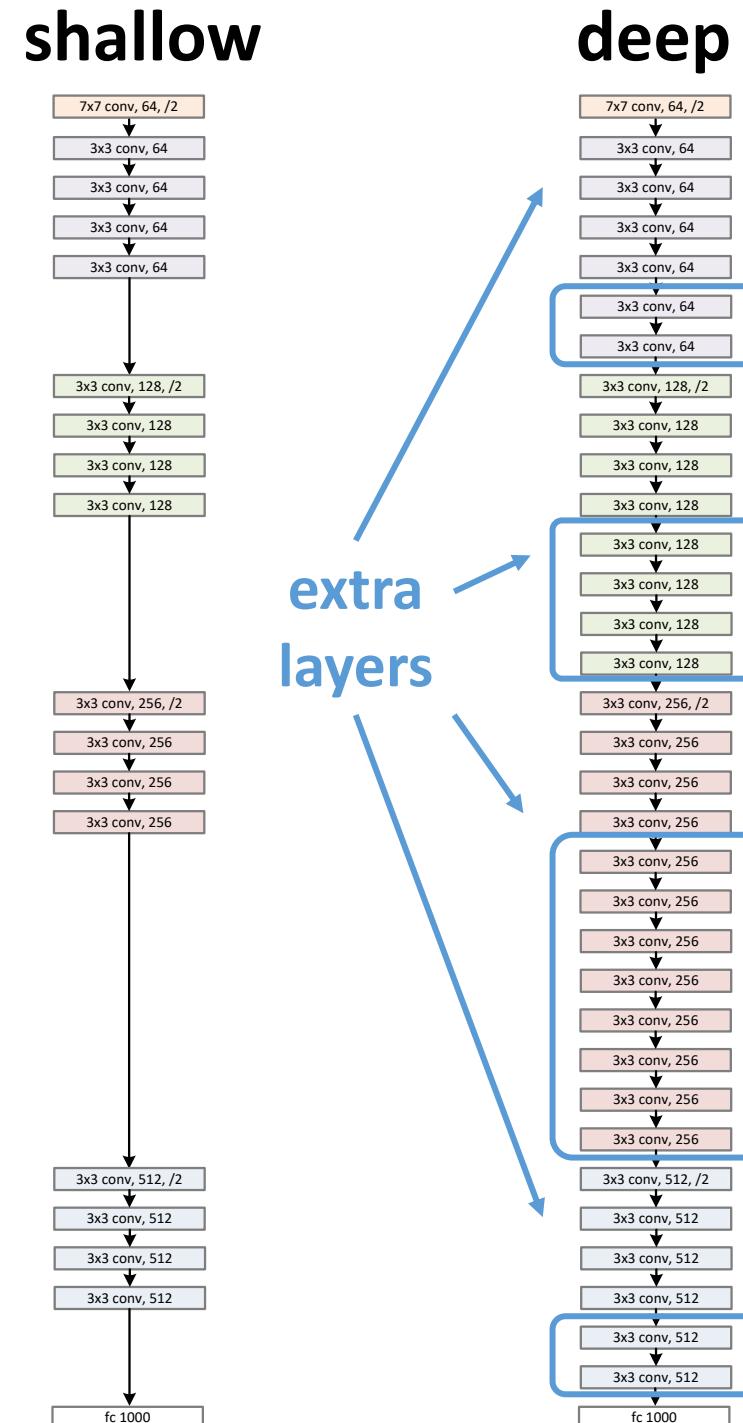
# The degradation problem

- Good init + norm enable training deeper models
- Simply stacking more layers?
- Degrade after  $\sim 20$  layers
- Not overfitting
- Difficult to train



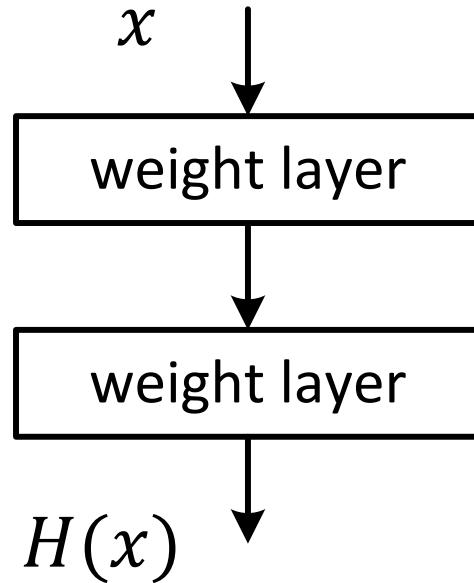
# A thought experiment

- Richer solution space
- Deeper nets should not degrade in training
- Solution by construction
  - take a good shallow net
  - copy trained layers
  - let extra layers be **identity mapping**
  - no worse training accuracy
- Optimizer cannot find the solution



# Deep Residual Learning

a subnet in  
a deep net

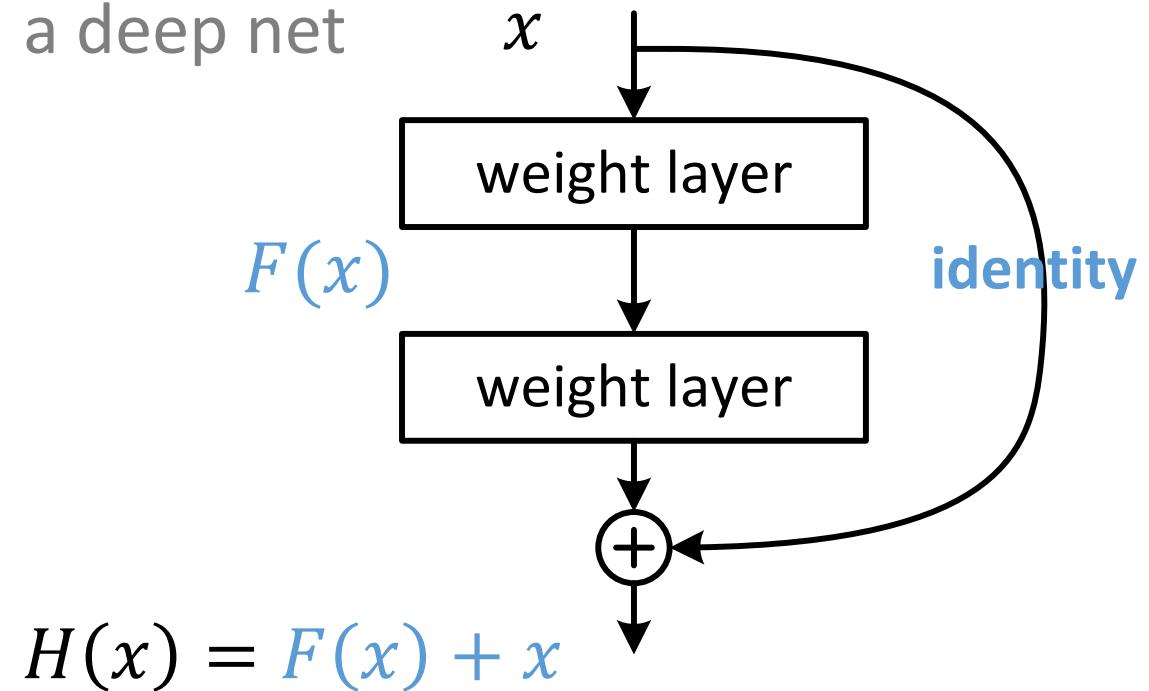


classical network

- $H(x)$ : desired function to be fit by a subnet
- let weight layers fit  $H(x)$

# Deep Residual Learning

a subnet in  
a deep net

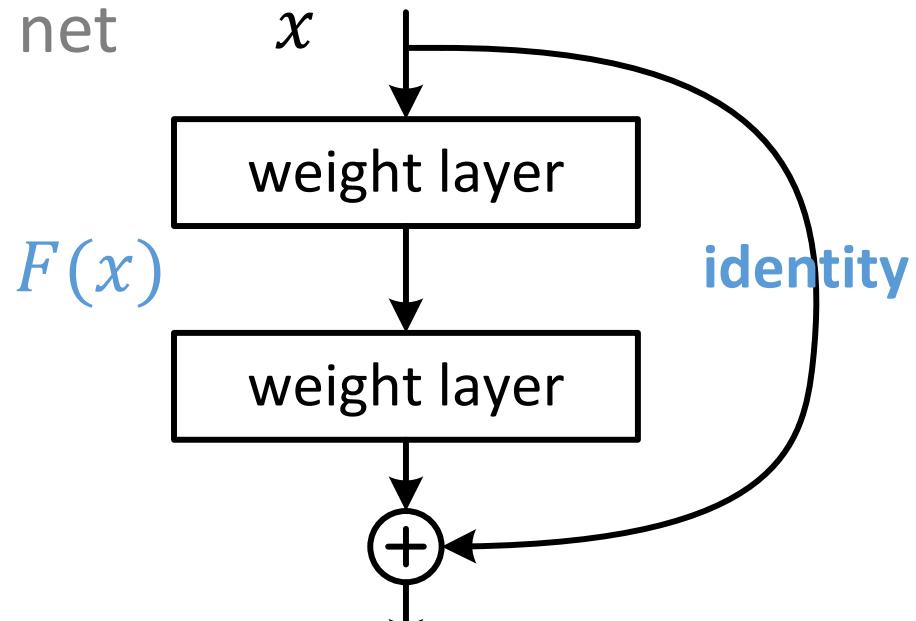


## residual block

- $H(x)$ : desired function to be fit by a subnet
- ~~let weight layers fit  $H(x)$~~
- let weight layers fit  $F(x)$
- set  $H(x) = F(x) + x$

# Deep Residual Learning

a subnet in  
a deep net



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

## residual block

- $F(x)$ : residual function
- if  $H(x) = \text{identity}$  is near-optimal
  - push weights to small
  - encourage small changes
- initialization
  - small or zero weights

# Residual Networks (ResNet)

Building very deep nets:

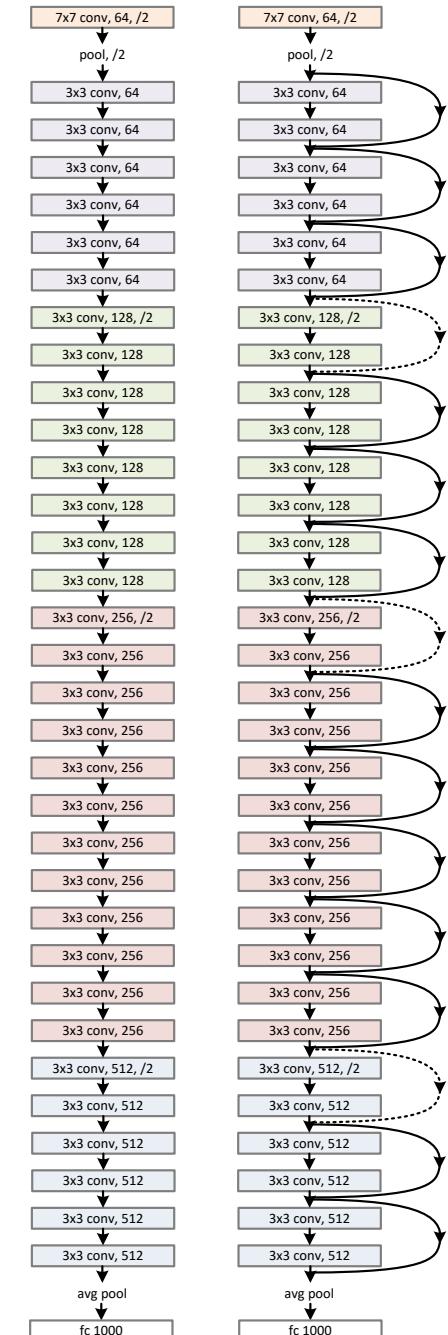
- add **identity connections** to vanilla nets  
(a.k.a. skip/shortcut/residual connections)

or:

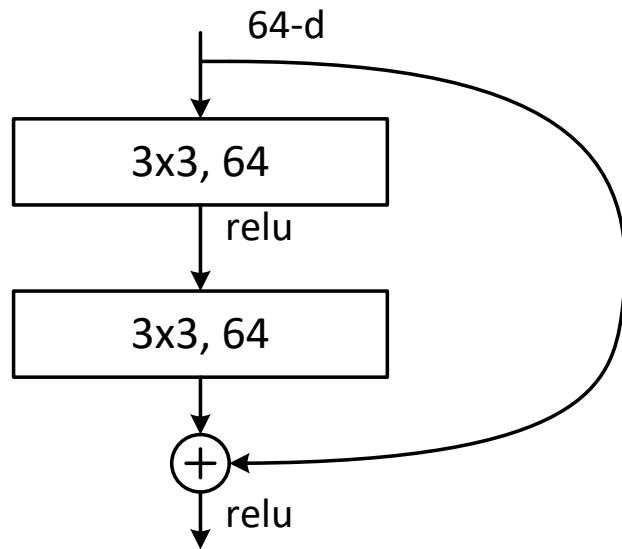
- stack many **residual blocks**

Residual Blocks:

- new generic modules for neural nets
- design blocks and compose them

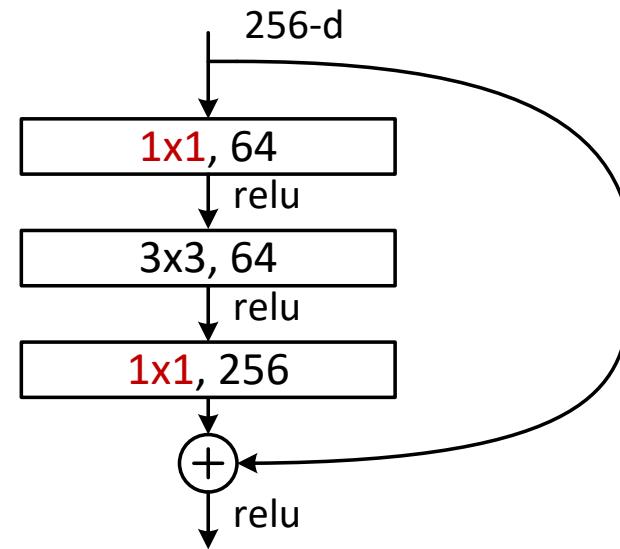


# Residual Block



## Basic Block

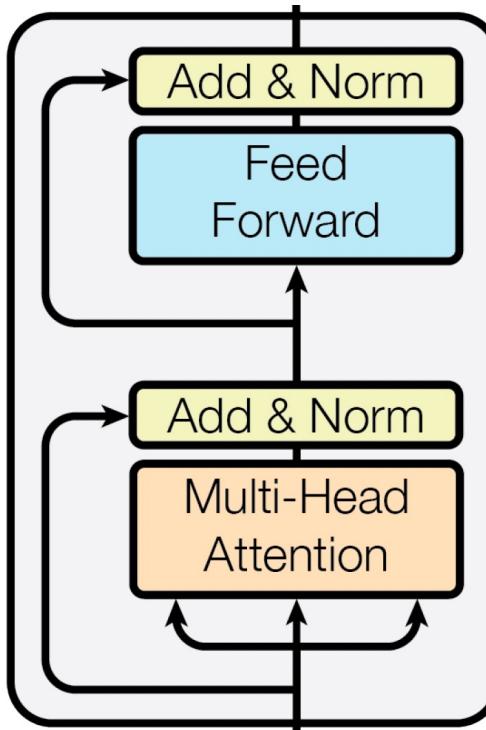
- 2-layer per block
- uniform width
- VGG style: all 3x3 conv



## Bottleneck Block

- 3-layer per block
- bottleneck: reduce dim
- w/ 1x1 conv for dim control

# Residual Block: Transformer



A Transformer Block has two Residual Blocks.

# Residual Networks (ResNet)

- For legacy reasons, “ResNet” often refers to the original instantiations
  - ResNet-18, 34: basic blocks
  - ResNet-50, 101, 152: bottleneck blocks
- In general, a ResNet is a network with the  $F(x) + x$  blocks
- The effect of residual learning is not limited to special instantiations



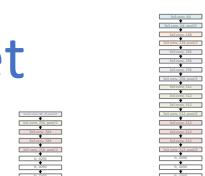
# Deep Residual Networks (ResNet)

ResNet-152



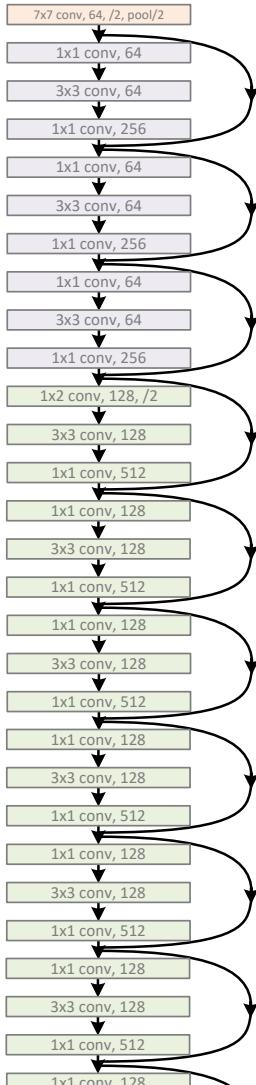
VGG-19

AlexNet

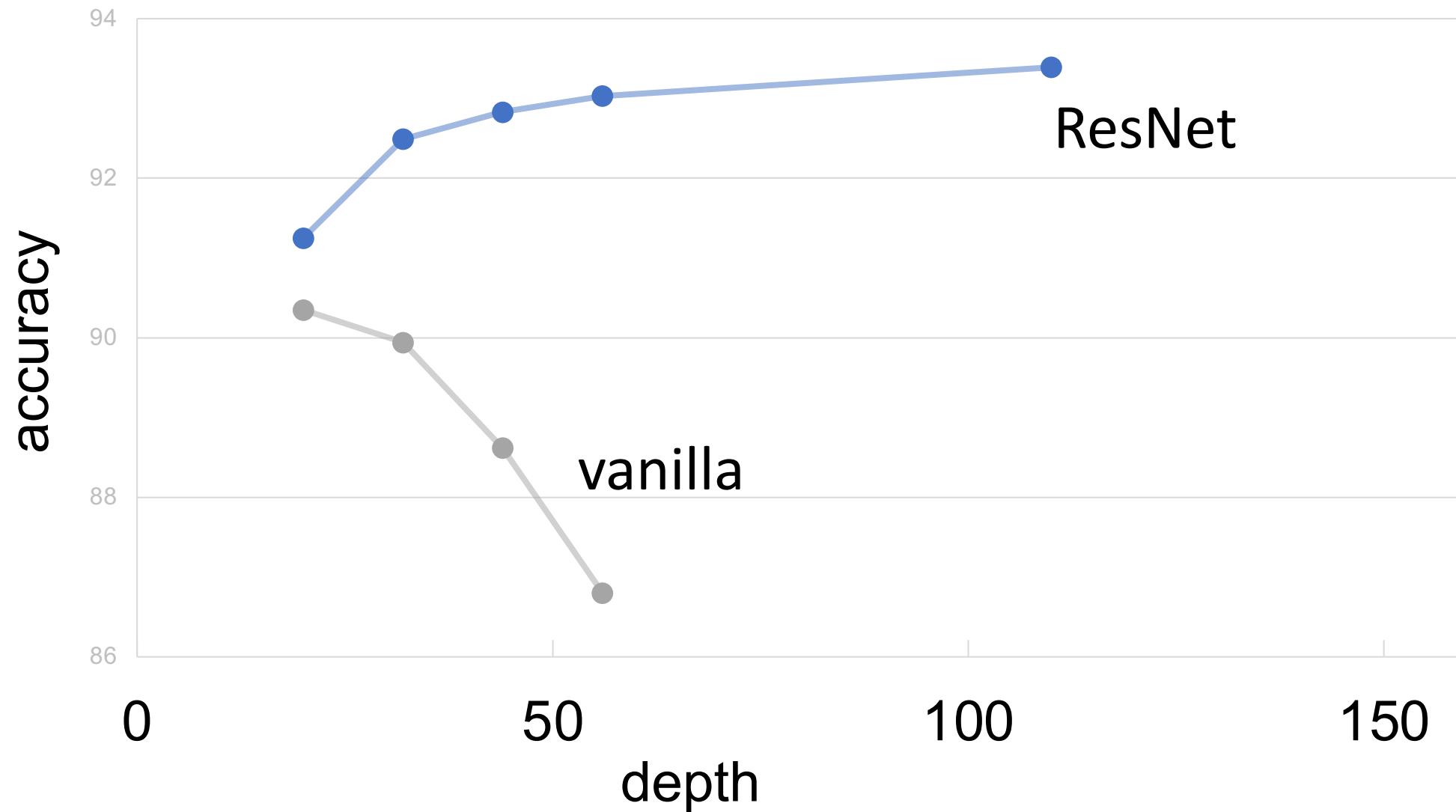


# Deep Residual Networks (ResNet)

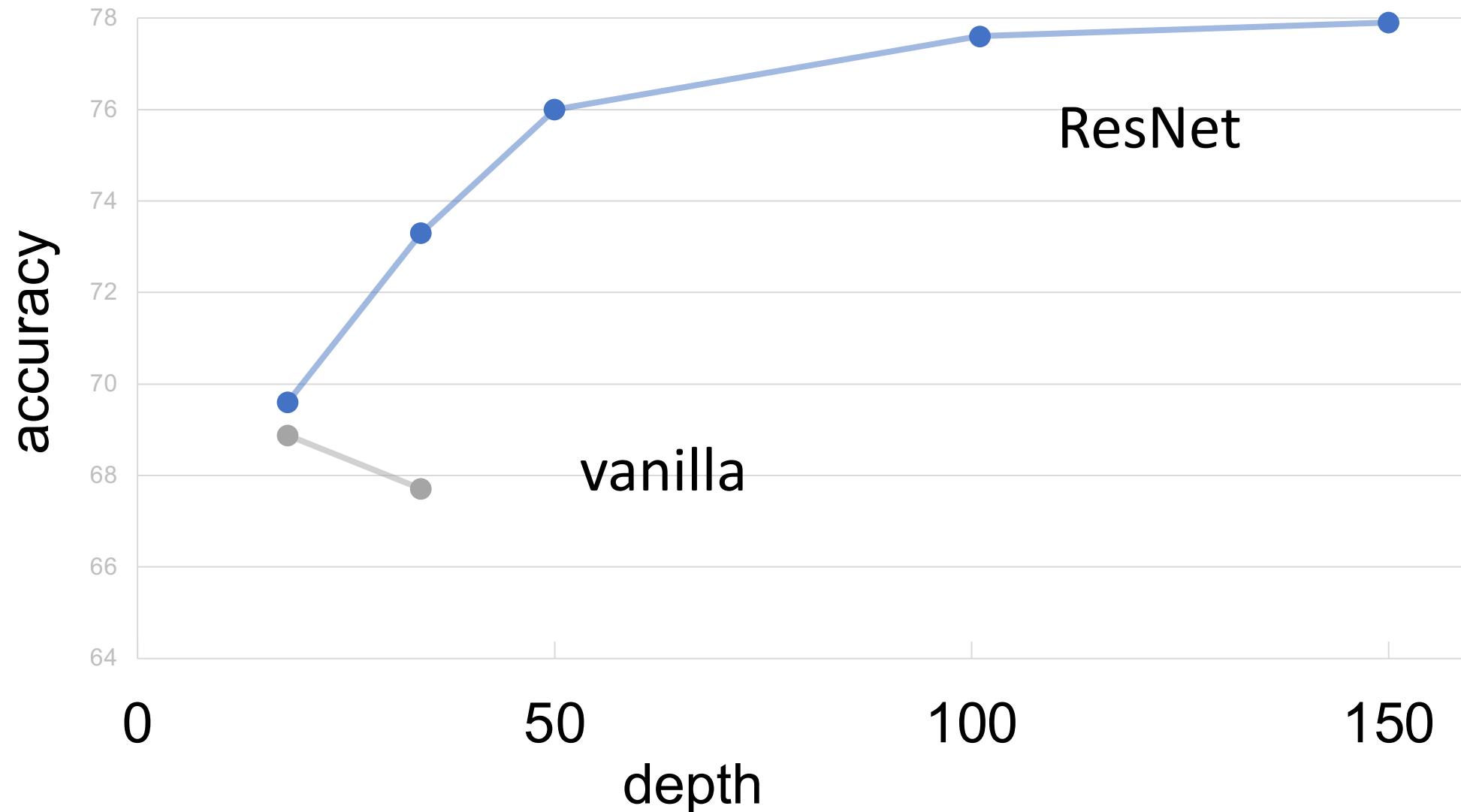
ResNet-152



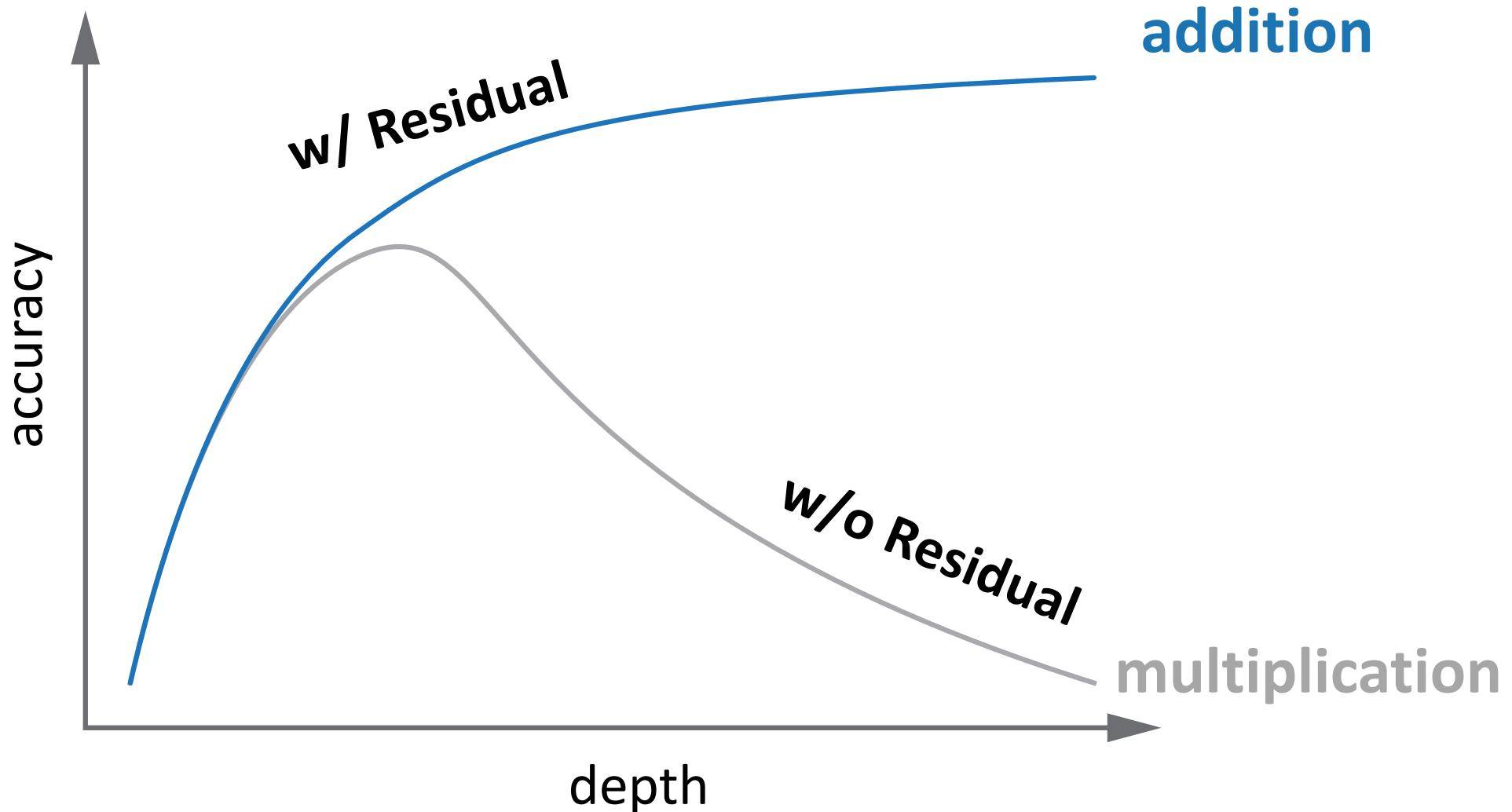
# Experiment on CIFAR-10 classification



# Experiment on ImageNet classification



# Understanding Residual Learning



# Understanding Residual Learning

Vanilla network

$$x_{l+1} = W_l x_l$$

recursively

$$x_L = \left( \prod W_i \right) x_l$$

gradients

$$\frac{\partial \mathcal{E}}{\partial x_l} = \frac{\partial \mathcal{E}}{\partial x_L} \left( \prod W_i \right)$$

# Understanding Residual Learning

Vanilla network

$$x_{l+1} = W_l x_l$$

recursively

$$x_L = \left( \prod W_i \right) x_l$$

gradients

$$\frac{\partial \mathcal{E}}{\partial x_l} = \frac{\partial \mathcal{E}}{\partial x_L} \left( \prod W_i \right)$$

- **vanishing signal**

# Understanding Residual Learning

Residual network

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

recursively

$$x_L = x_l + \sum F_i(x_i)$$

gradients

$$\frac{\partial \mathcal{E}}{\partial x_l} = \frac{\partial \mathcal{E}}{\partial x_L} \left( 1 + \frac{\partial}{\partial x_l} \sum F_i(x_i) \right)$$

# Understanding Residual Learning

Residual network

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

recursively

$$x_L = \textcolor{blue}{x_l} + \sum F_i(x_i)$$

gradients

$$\frac{\partial \mathcal{E}}{\partial x_l} = \frac{\partial \mathcal{E}}{\partial x_L} \left( \textcolor{blue}{1} + \frac{\partial}{\partial x_l} \sum F_i(x_i) \right)$$

- **not vanishing**

# Understanding Residual Learning

Residual network

recursively

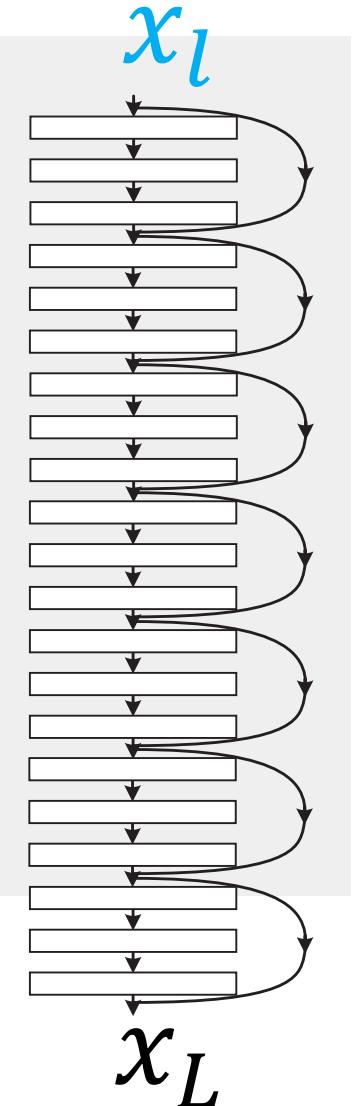
gradients

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

$$x_L = \textcolor{blue}{x_l} + \sum F_i(x_i)$$

$$\frac{\partial \mathcal{E}}{\partial x_l} = \frac{\partial \mathcal{E}}{\partial x_L} \left( 1 + \frac{\partial}{\partial x_l} \sum F_i(x_i) \right)$$

- **not vanishing**



# Understanding Residual Learning

Residual network

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

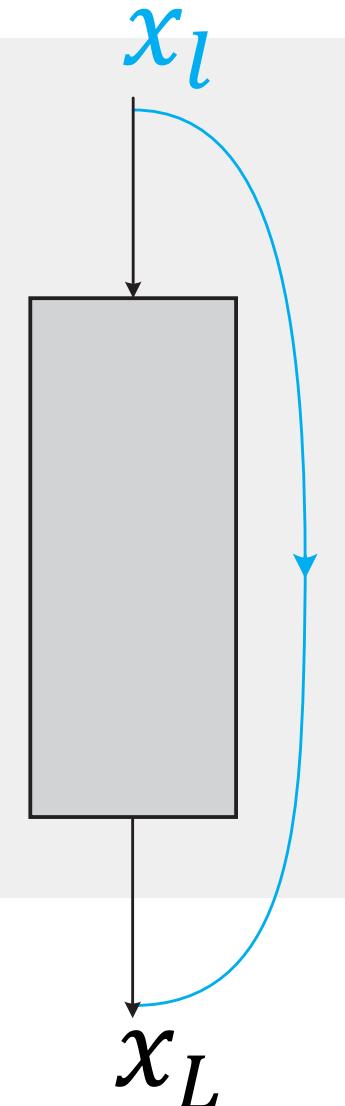
recursively

$$x_L = \textcolor{blue}{x_l} + \sum F_i(x_i)$$

gradients

$$\frac{\partial \mathcal{E}}{\partial x_l} = \frac{\partial \mathcal{E}}{\partial x_L} \left( 1 + \frac{\partial}{\partial x_l} \sum F_i(x_i) \right)$$

- **not vanishing**



# Understanding Residual Learning

Residual network

recursively

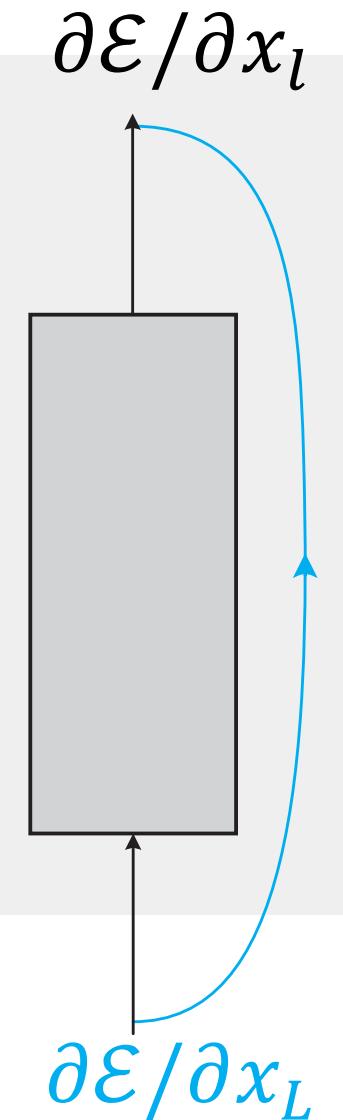
gradients

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

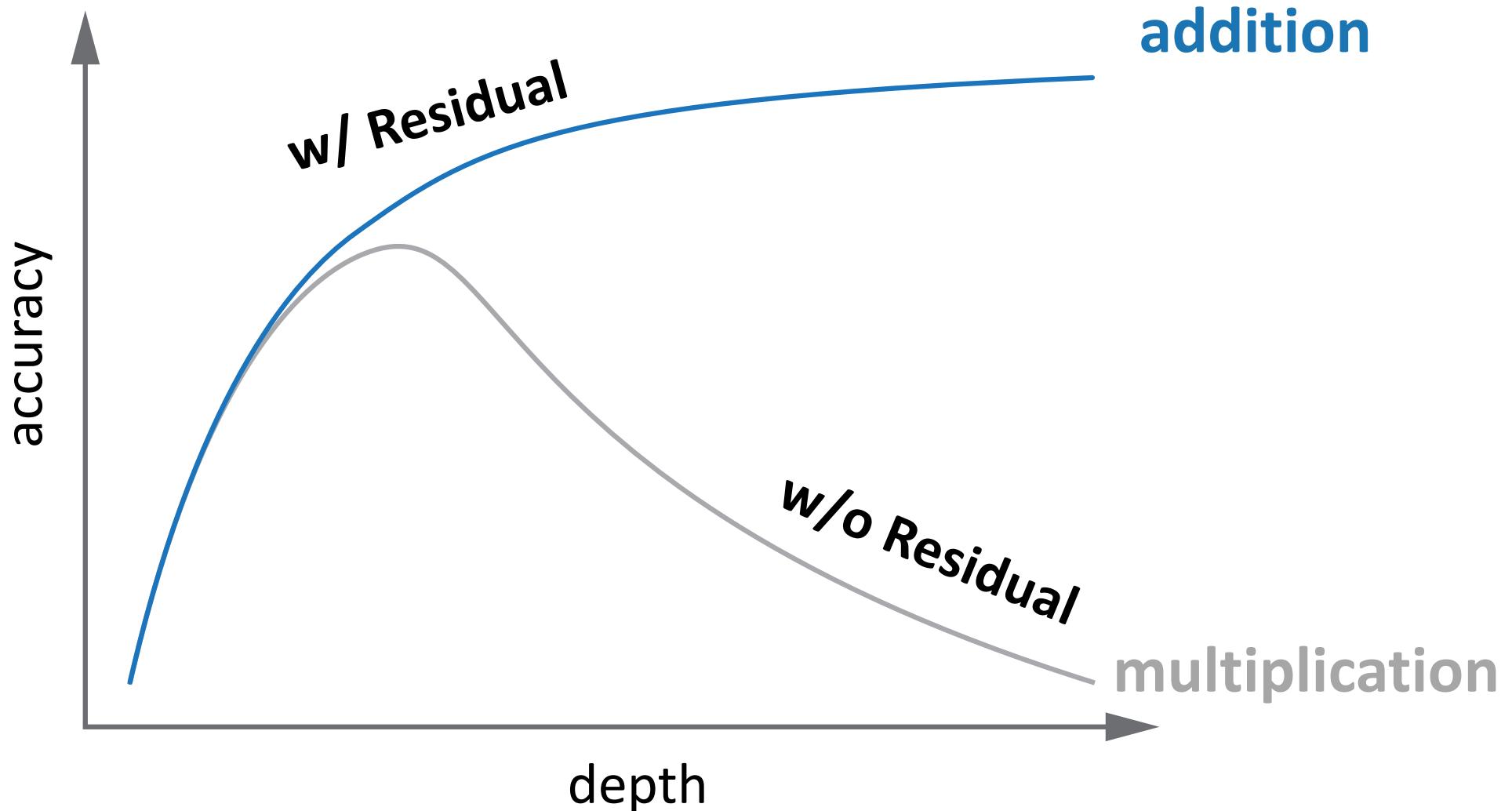
$$x_L = x_l + \sum F_i(x_i)$$

$$\frac{\partial \mathcal{E}}{\partial x_l} = \frac{\partial \mathcal{E}}{\partial x_L} \left( 1 + \frac{\partial}{\partial x_l} \sum F_i(x_i) \right)$$

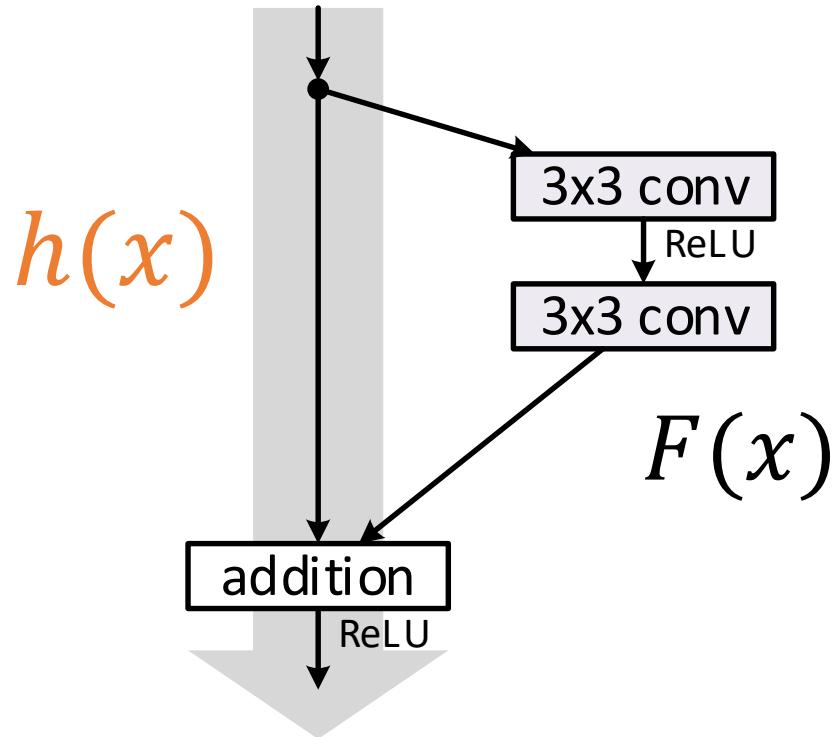
- **not vanishing**



# Understanding Residual Learning



# Identity Mappings in Residual Learning



- Residual Learning:  $h$  = identity mapping

Question 1:  
What if shortcut mapping  $h \neq$  identity?

if  $h$  is multiplicative ...

consider  $h(x) = \lambda x$ :

$$x_{l+1} = \lambda x_l + F_l(x_l)$$

recursively

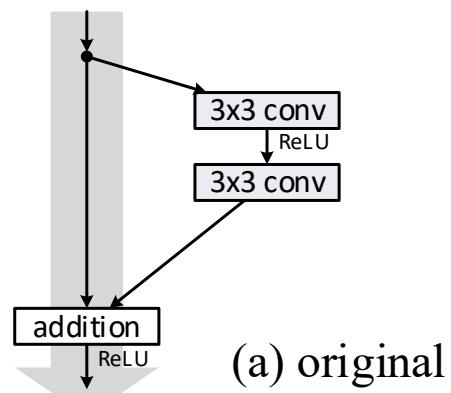
$$x_L = \lambda^{L-l} x_l + \sum \hat{F}_i(x_i)$$

gradients

$$\frac{\partial \mathcal{E}}{\partial x_l} = \frac{\partial \mathcal{E}}{\partial x_L} \left( \lambda^{L-l} + \frac{\partial}{\partial x_l} \sum \hat{F}_i(x_i) \right)$$

- **vanishing signals, if  $\lambda < 1$**

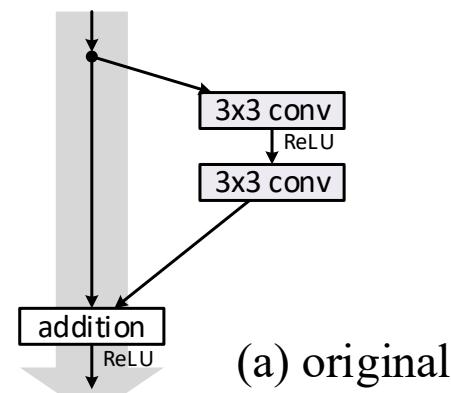
$h(x) = x$   
error: **6.6%**



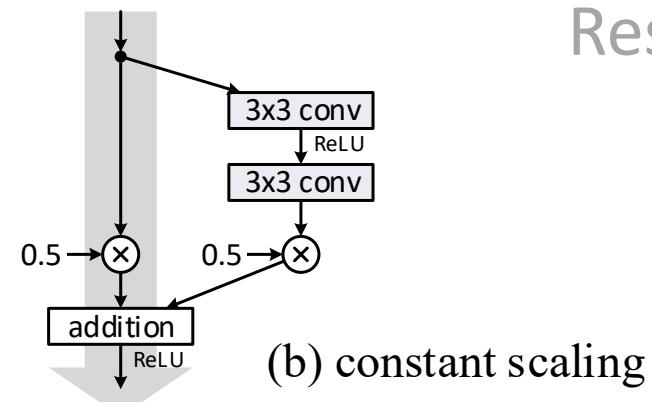
# ResNet-110 on CIFAR-10

$$h(x) = x$$

error: **6.6%**



(a) original

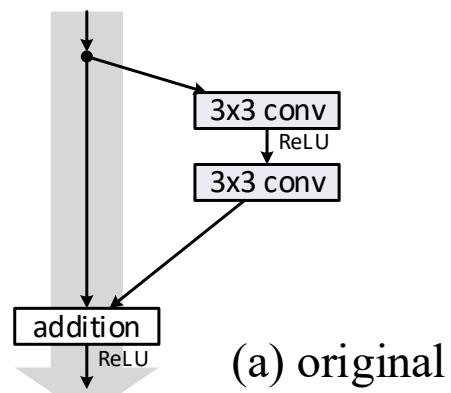


(b) constant scaling

$$h(x) = 0.5x$$

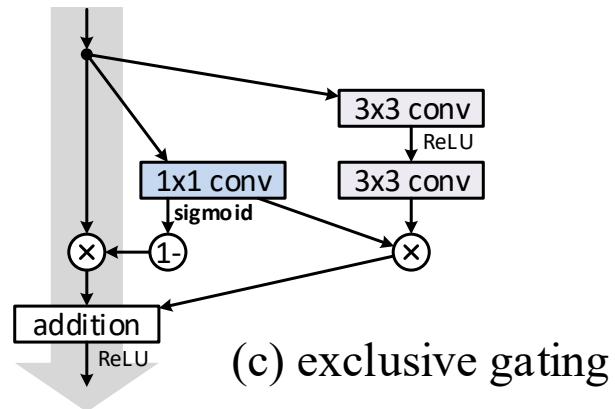
error: **12.4%**

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

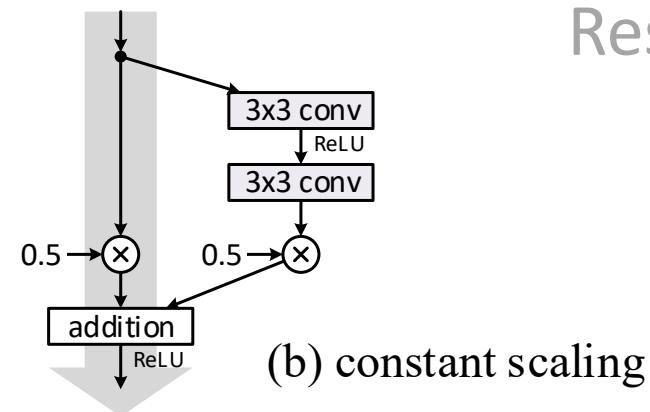


(a) original

$h(x) = \text{gate} \cdot x$   
error: **8.7%**



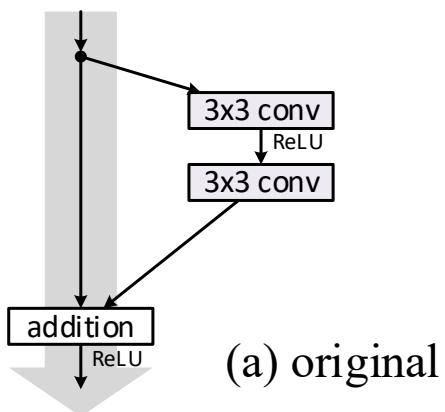
(c) exclusive gating



(b) constant scaling

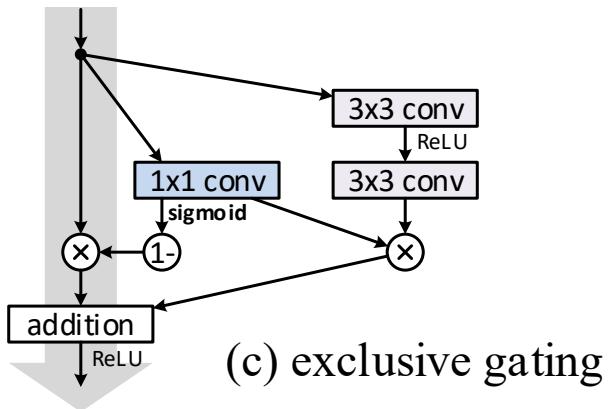
$h(x) = 0.5x$   
error: **12.4%**

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

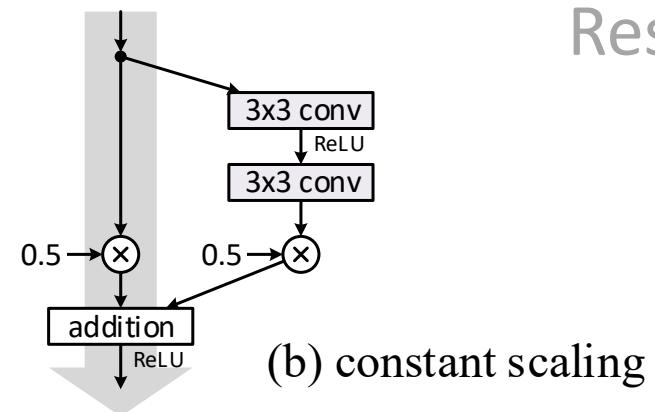


(a) original

$h(x) = \text{gate} \cdot x$   
error: **8.7%**

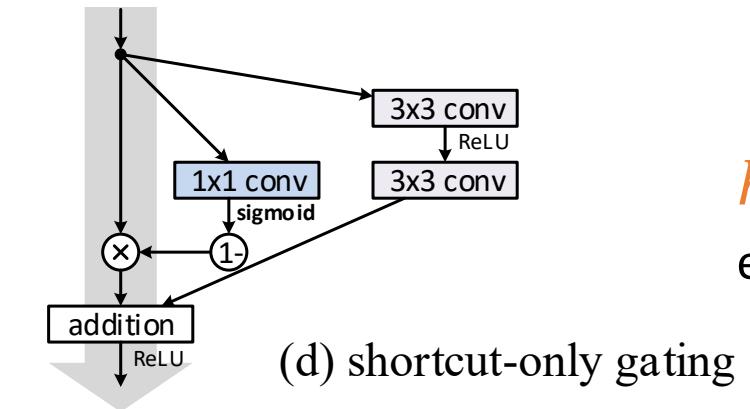


(c) exclusive gating



(b) constant scaling

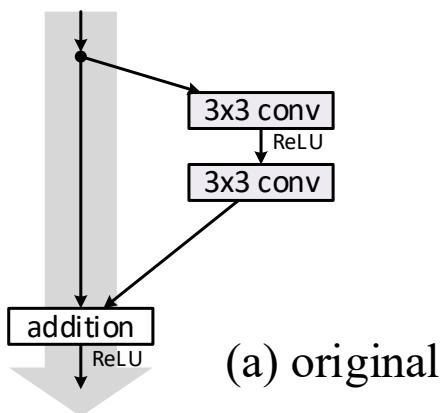
$h(x) = 0.5x$   
error: **12.4%**



(d) shortcut-only gating

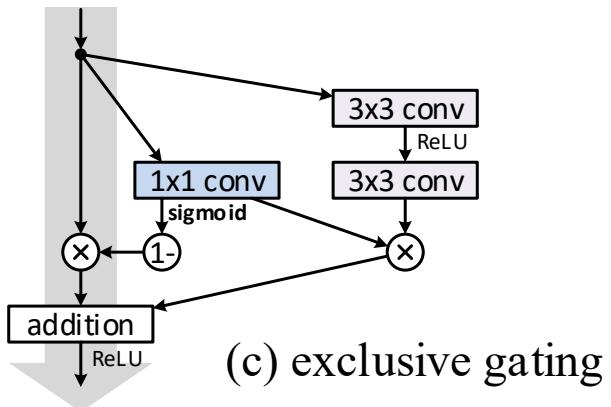
$h(x) = \text{gate} \cdot x$   
error: **12.9%**

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



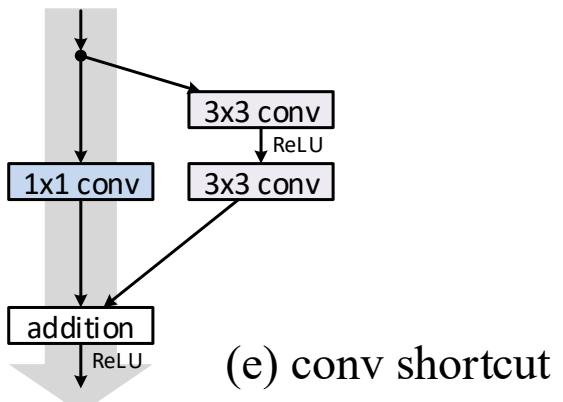
(a) original

$h(x) = \text{gate} \cdot x$   
error: **8.7%**

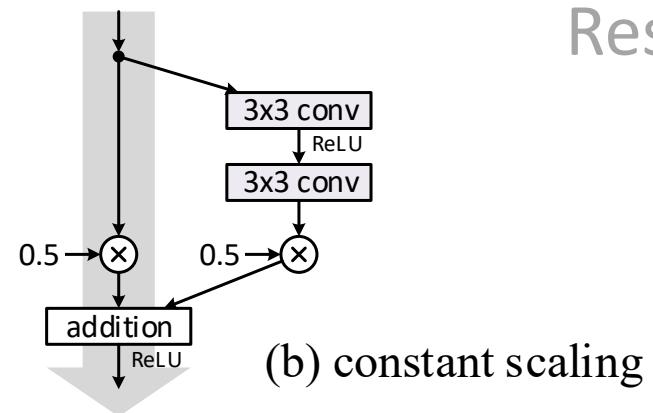


(c) exclusive gating

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

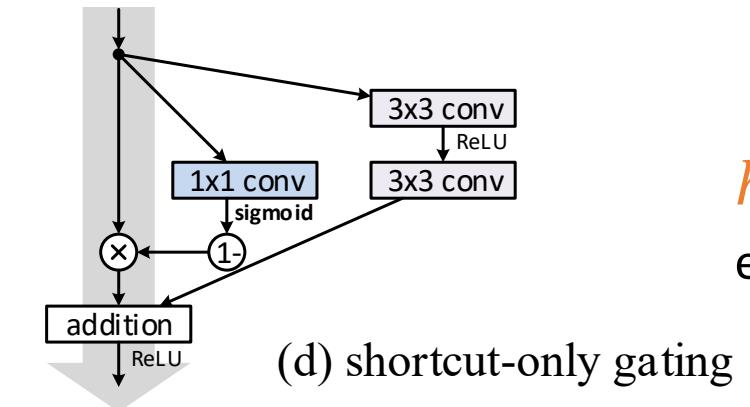


(e) conv shortcut



(b) constant scaling

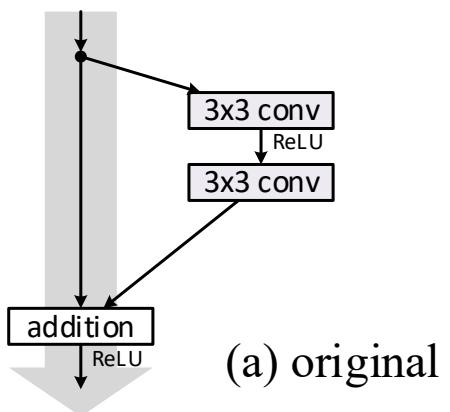
$h(x) = 0.5x$   
error: **12.4%**



(d) shortcut-only gating

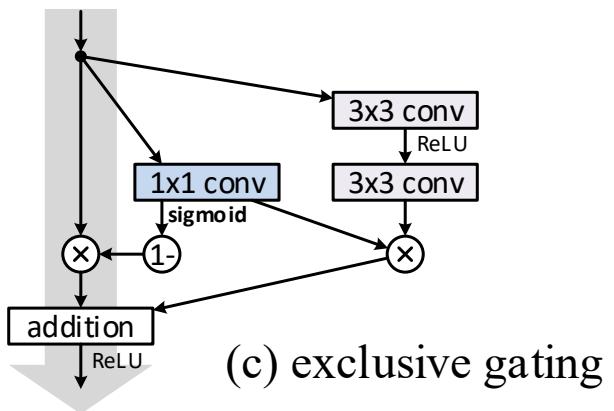
$h(x) = \text{gate} \cdot x$   
error: **12.9%**

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



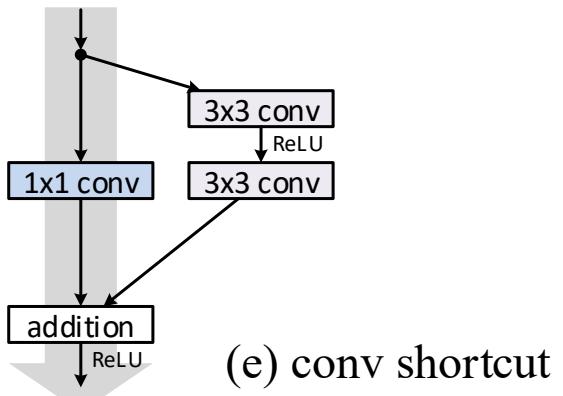
(a) original

$h(x) = \text{gate} \cdot x$   
error: **8.7%**

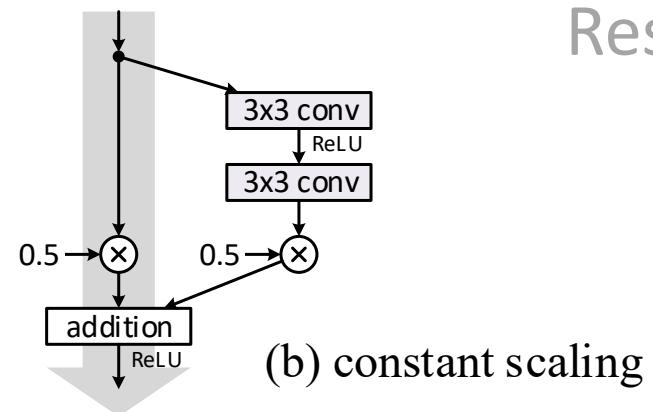


(c) exclusive gating

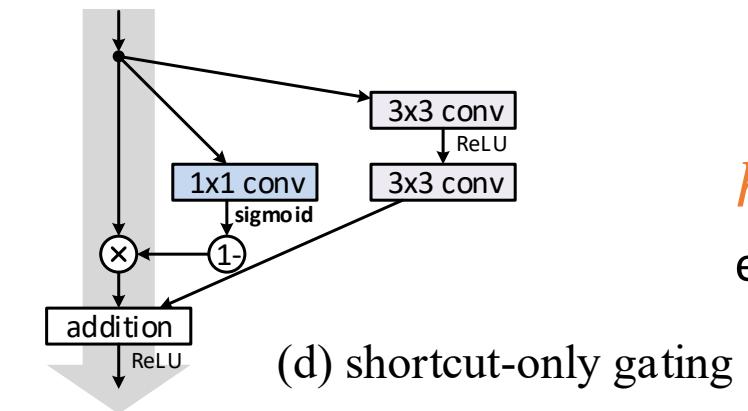
$h(x) = \text{conv}(x)$   
error: **12.2%**



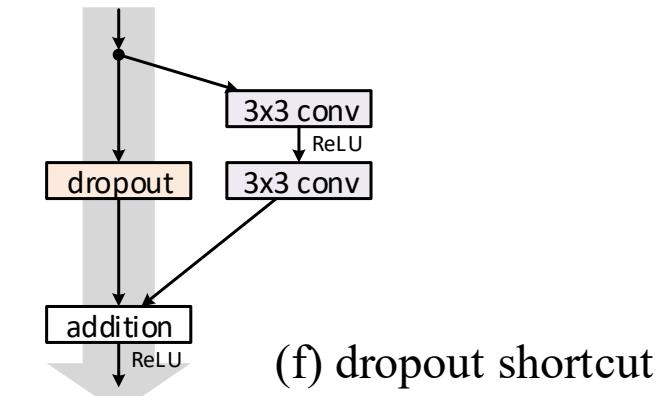
(e) conv shortcut



(b) constant scaling



(d) shortcut-only gating



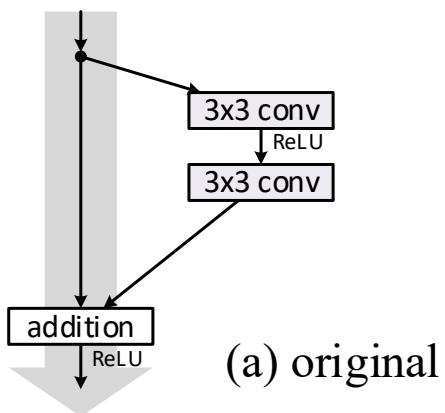
(f) dropout shortcut

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

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

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

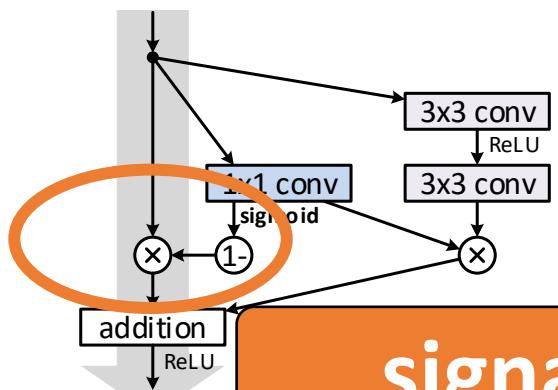
$h(x) = x$   
error: **6.6%**



(a) original

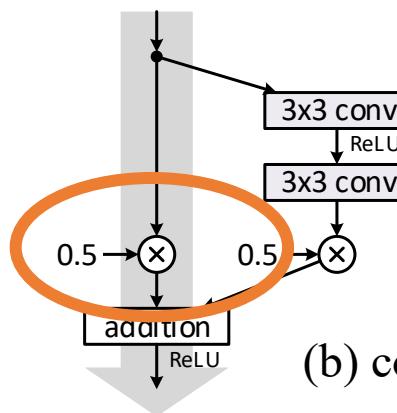
$h(x) = \text{gate} \cdot x$   
error: **8.7%**

\*like "Highway Net"



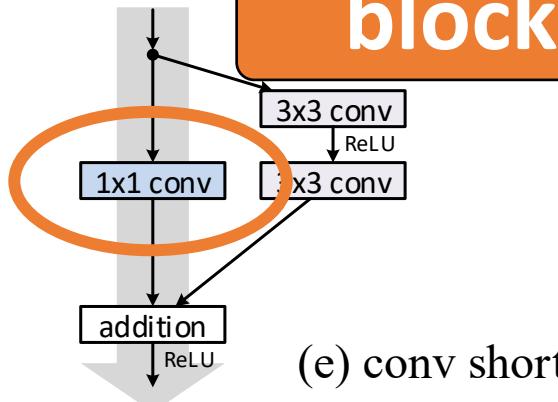
(b) constant scaling

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



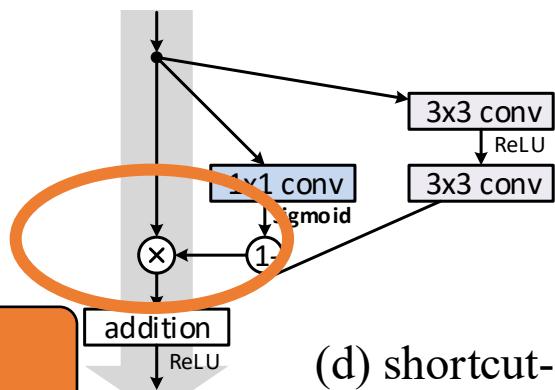
(c) identity mapping

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



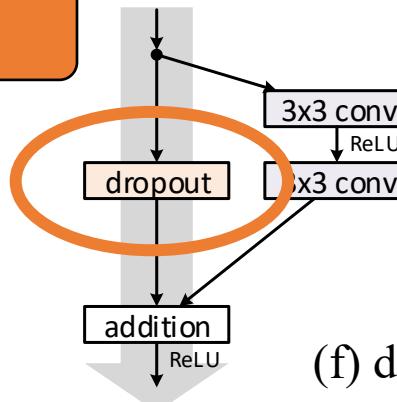
(d) shortcut-only gating

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



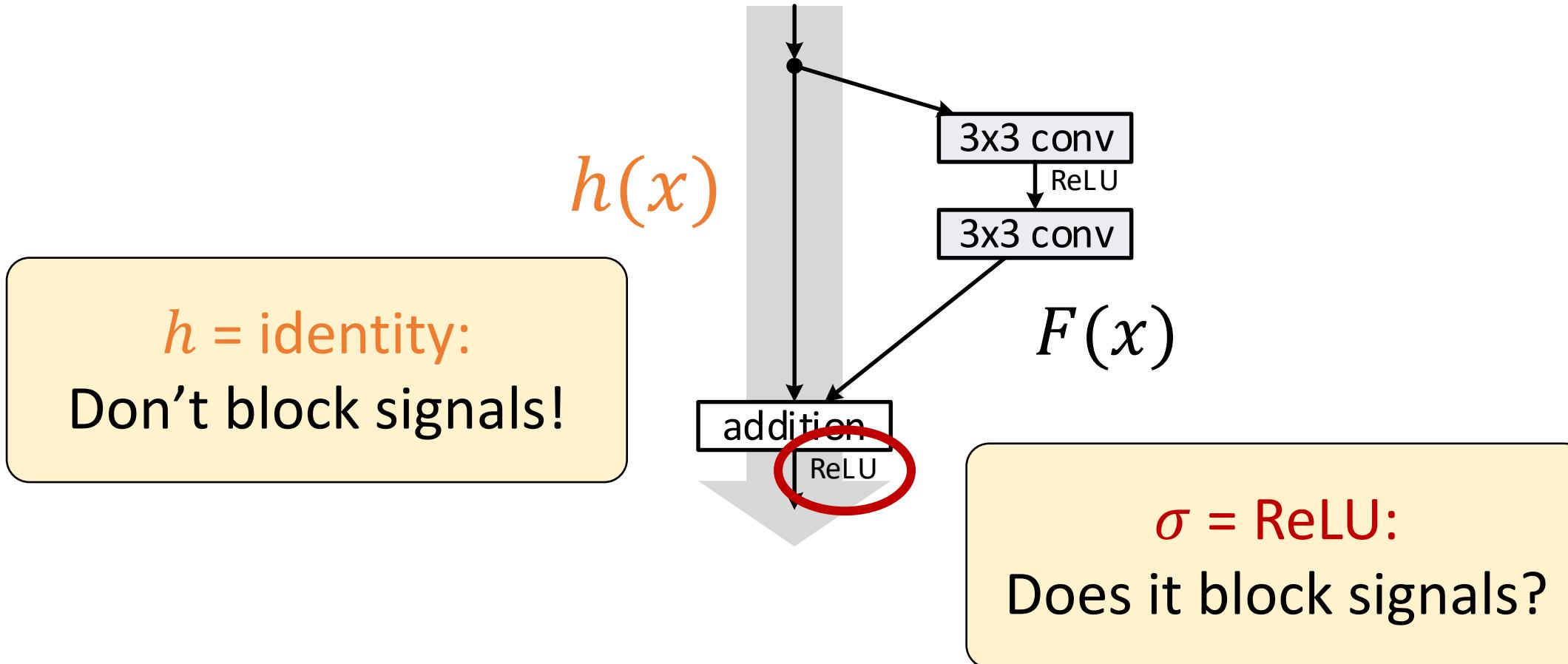
(e) conv shortcut

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



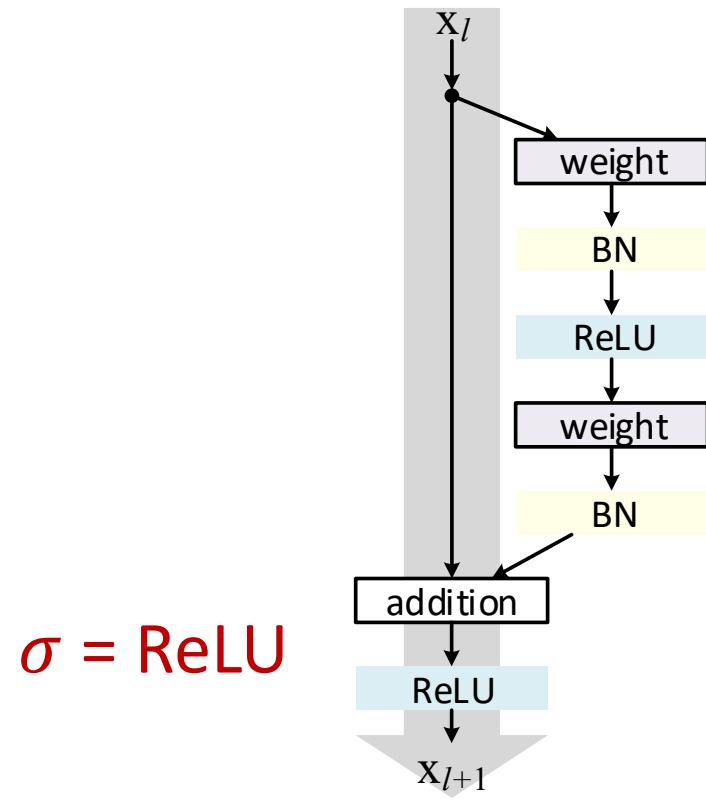
(f) dropout shortcut

# Identity Mappings in Residual Learning



Question 2:  
What if post activation  $\sigma$  = identity?

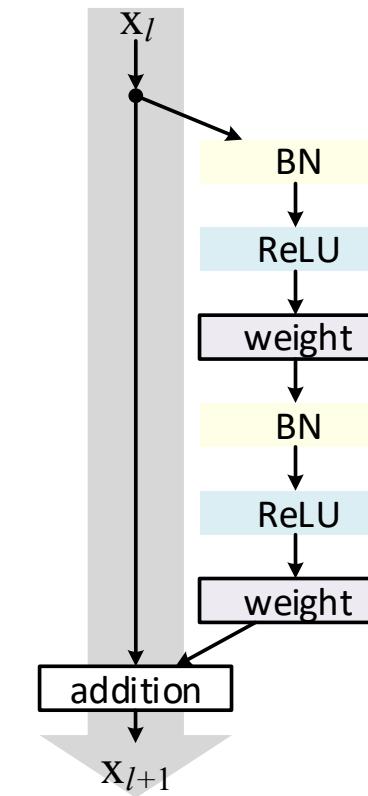
# Pre-activation (Pre-normalization)



$$\sigma = \text{ReLU}$$

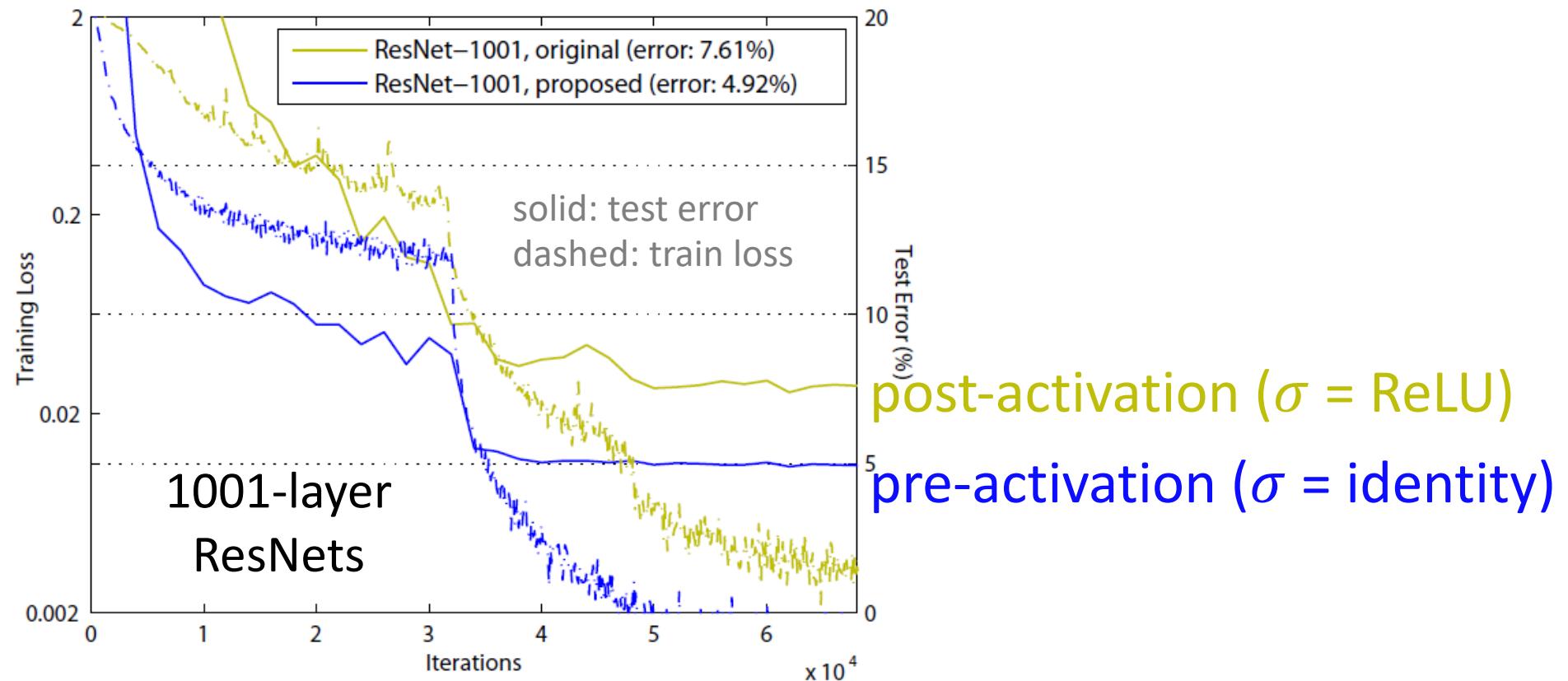
post-activation

$$\sigma = \text{identity}$$



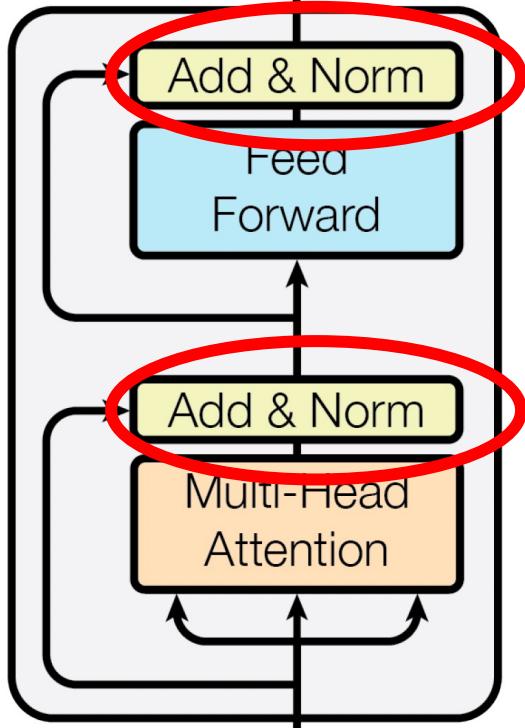
pre-activation

# Pre-activation (Pre-normalization)

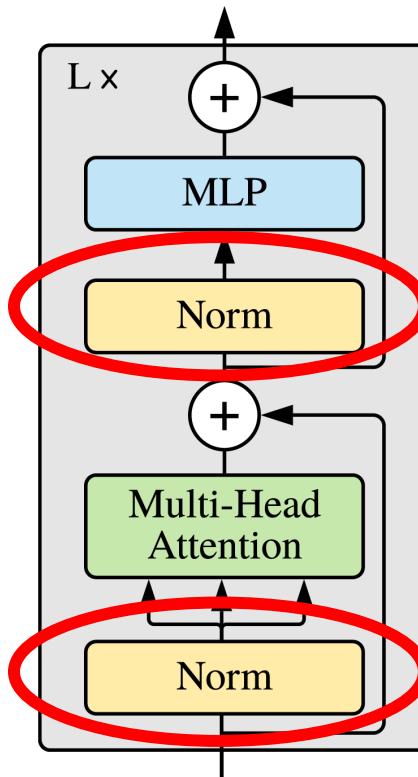


- ReLU can also block signals when there are many layers
- Pre-activation eases optimization

# Pre-activation (Pre-normalization)



post-norm  
(original Transformer)

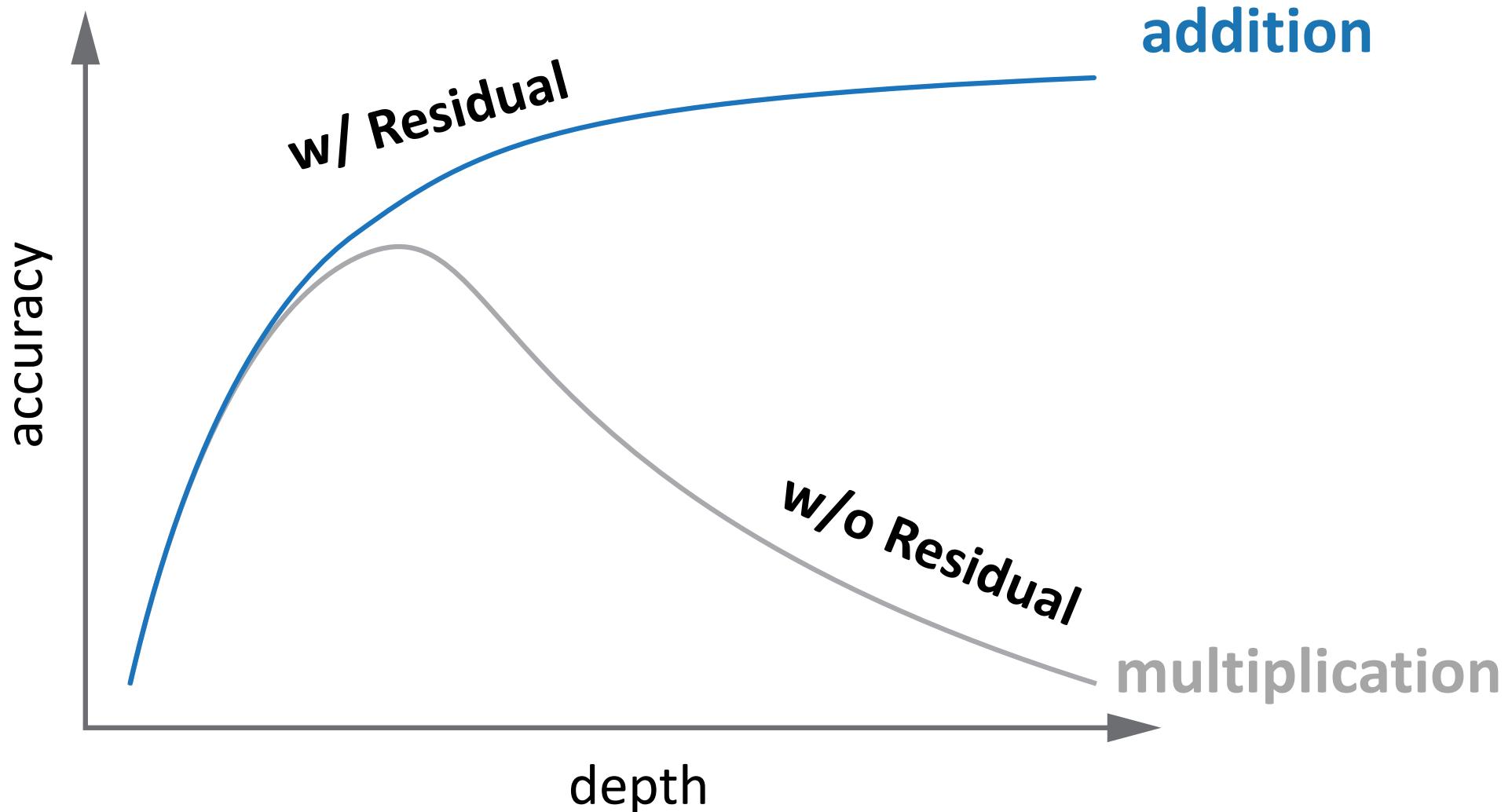


pre-norm  
(after GPT-2)

few modifications. Layer normalization (Ba et al., 2016) was moved to the input of each sub-block, similar to a pre-activation residual network (He et al., 2016) and an additional layer normalization was added after the final self-attention block. A modified initialization which accounts

from the GPT-2 paper  
[Radford, et al., 2019]

# Understanding Residual Learning



# Understanding Residual Learning

The Shattered Gradients Problem:  
If resnets are the answer, then what is the question?

David Balduzzi<sup>1</sup> Marcus Frean<sup>1</sup> Lennox Leary<sup>1</sup> JP Lewis<sup>1,2</sup> Kurt Wan-Duo Ma<sup>1</sup> Brian McWilliams<sup>3</sup>

Residual Networks Behave Like Ensembles of  
Relatively Shallow Networks

Andreas Veit Michael Wilber Serge Belongie  
Department of Computer Science & Cornell Tech  
Cornell University  
[{av443, mjwt285, sjb344}@cornell.edu](mailto:{av443, mjwt285, sjb344}@cornell.edu)

Visualizing the Loss Landscape of Neural Nets

Hao Li<sup>1</sup>, Zheng Xu<sup>1</sup>, Gavin Taylor<sup>2</sup>, Christoph Studer<sup>3</sup>, Tom Goldstein<sup>1</sup>  
<sup>1</sup>University of Maryland, College Park <sup>2</sup>United States Naval Academy <sup>3</sup>Cornell University  
[{haoli, xuzh, tomg}@cs.umd.edu, taylor@usna.edu, studer@cornell.edu](mailto:{haoli, xuzh, tomg}@cs.umd.edu, taylor@usna.edu, studer@cornell.edu)

Identity Matters in Deep Learning

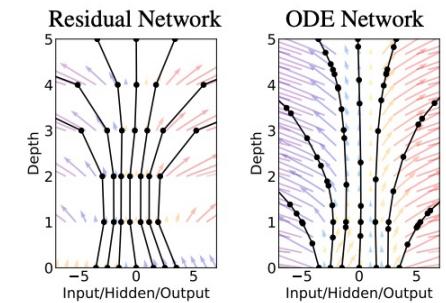
Moritz Hardt \* Tengyu Ma†  
July 23, 2018

SKIP CONNECTIONS ELIMINATE SINGULARITIES

A. Emin Orhan [aeminorhan@gmail.com](mailto:aeminorhan@gmail.com) Xaq Pitkow [xaq@rice.edu](mailto:xaq@rice.edu)  
Baylor College of Medicine & Rice University

Neural Ordinary Differential Equations

Ricky T. Q. Chen\*, Yulia Rubanova\*, Jesse Bettencourt\*, David Duvenaud  
University of Toronto, Vector Institute  
[{rtqichen, rubanova, jessebett, duvenaud}@cs.toronto.edu](mailto:{rtqichen, rubanova, jessebett, duvenaud}@cs.toronto.edu)



# Understanding Residual Learning

- Residual connections produce smoother loss landscapes

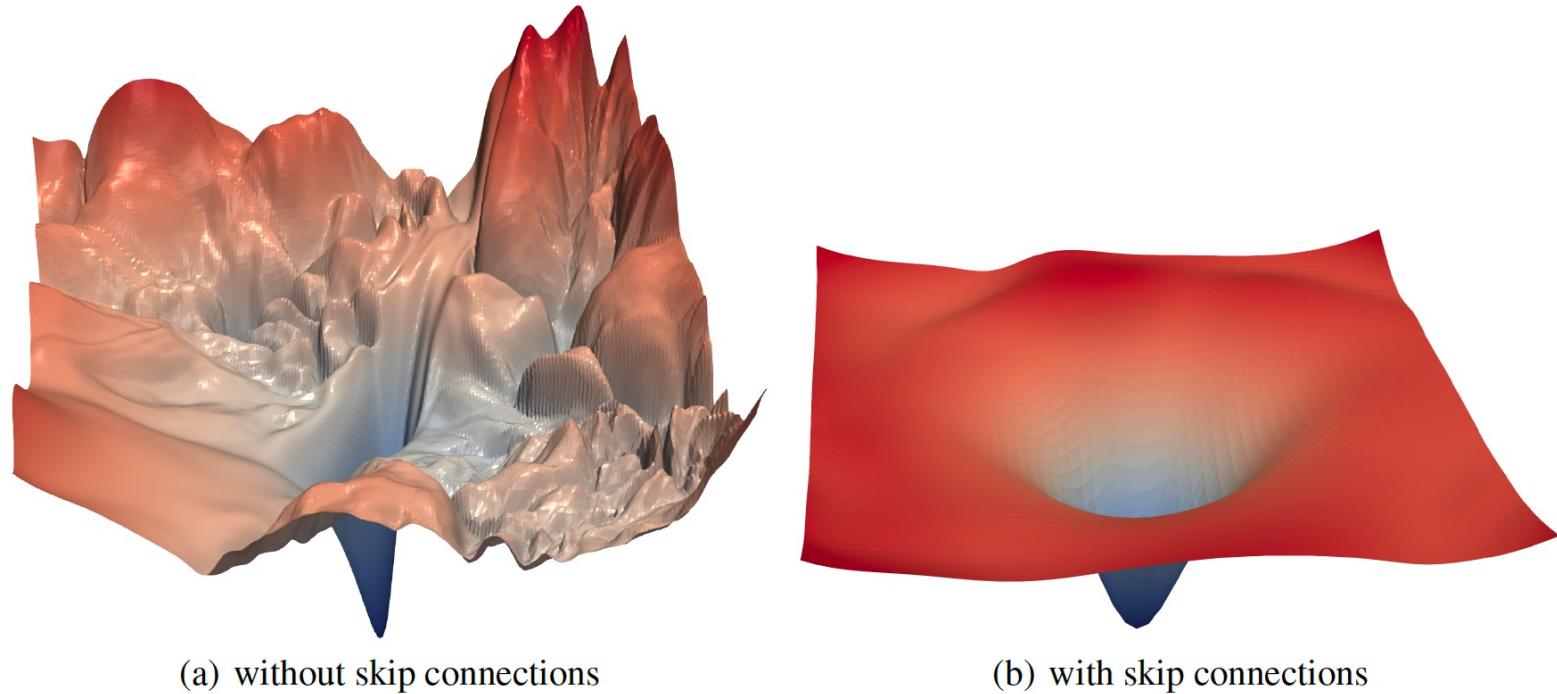
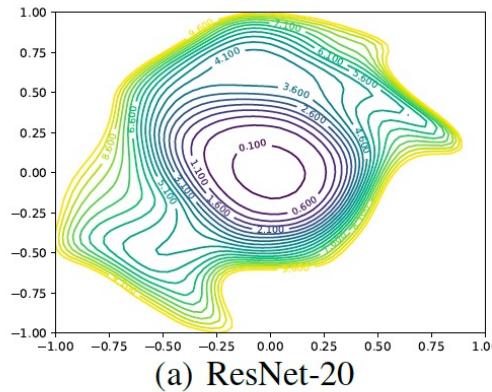


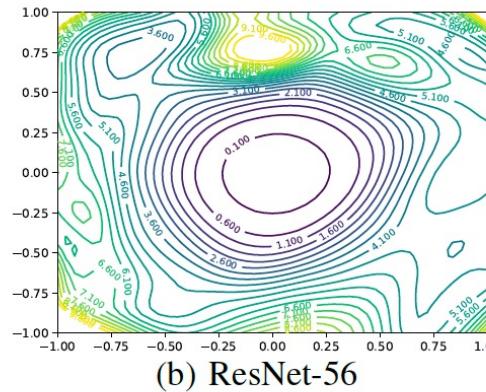
Figure 1: The loss surfaces of ResNet-56 with/without skip connections. The vertical axis is logarithmic to show dynamic range. The proposed filter normalization scheme is used to enable comparisons of sharpness/flatness between the two figures.

# Understanding Residual Learning

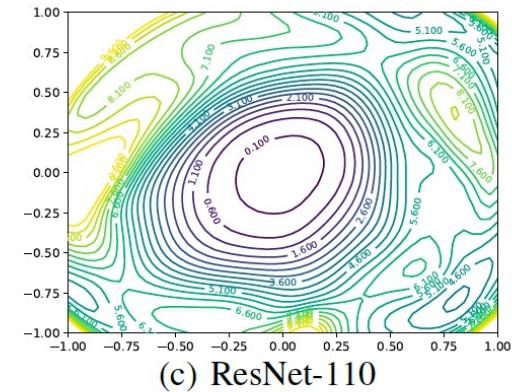
- Without residual connections, deeper nets have more chaotic landscapes



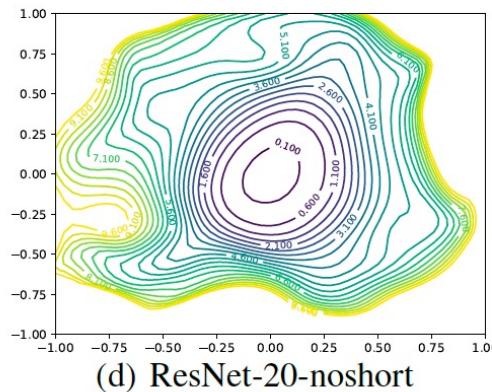
(a) ResNet-20



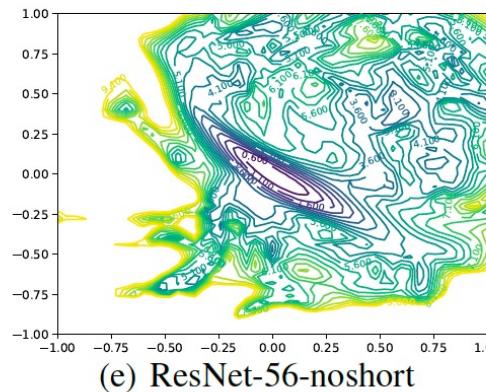
(b) ResNet-56



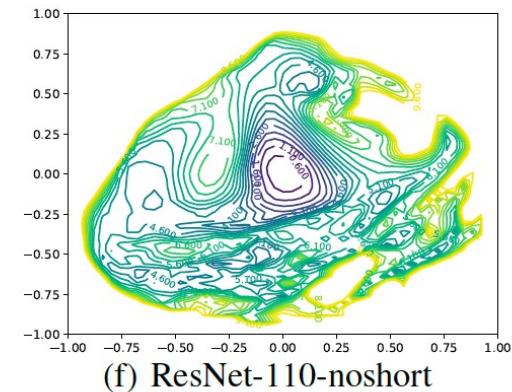
(c) ResNet-110



(d) ResNet-20-noshort

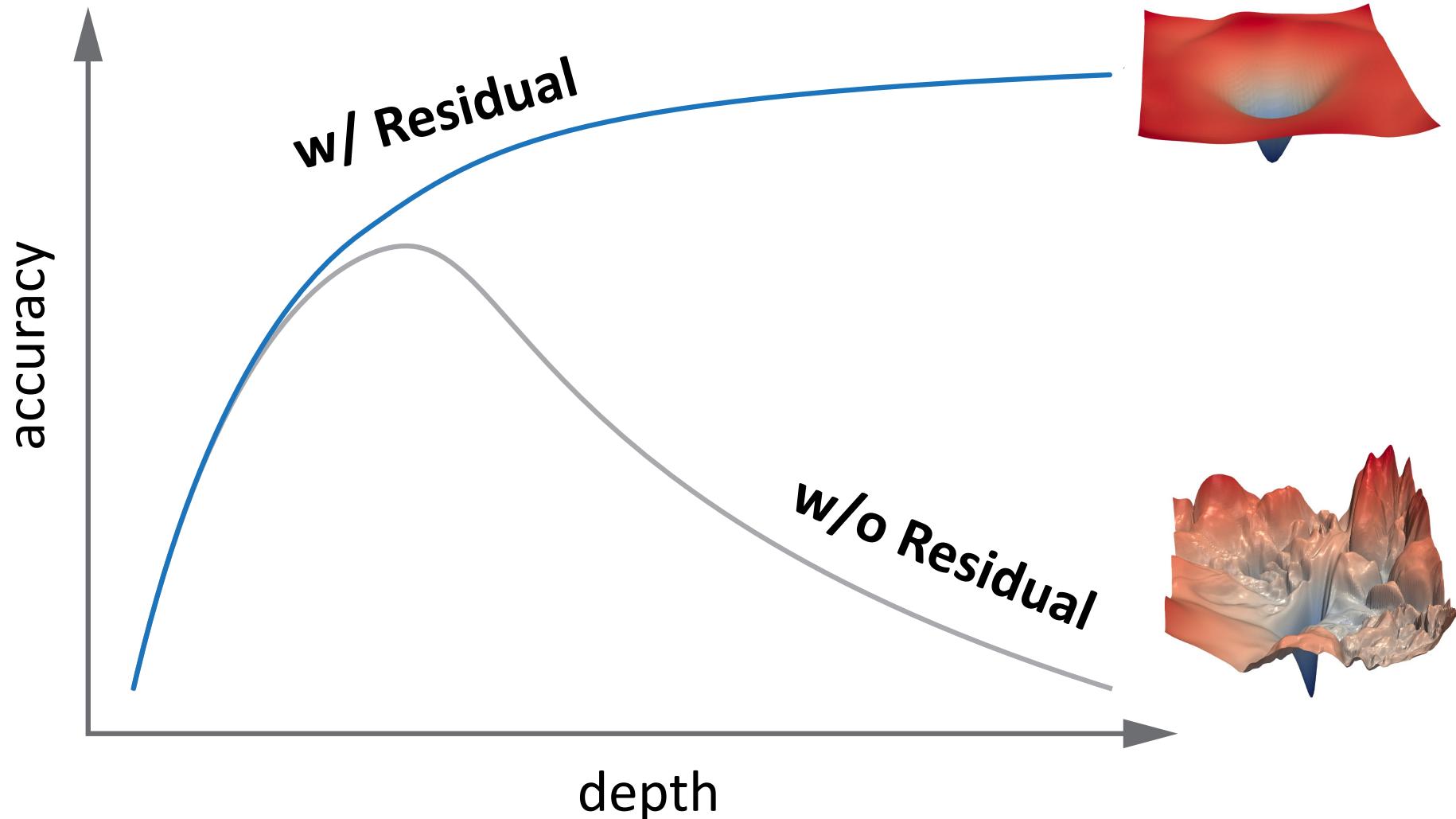


(e) ResNet-56-noshort

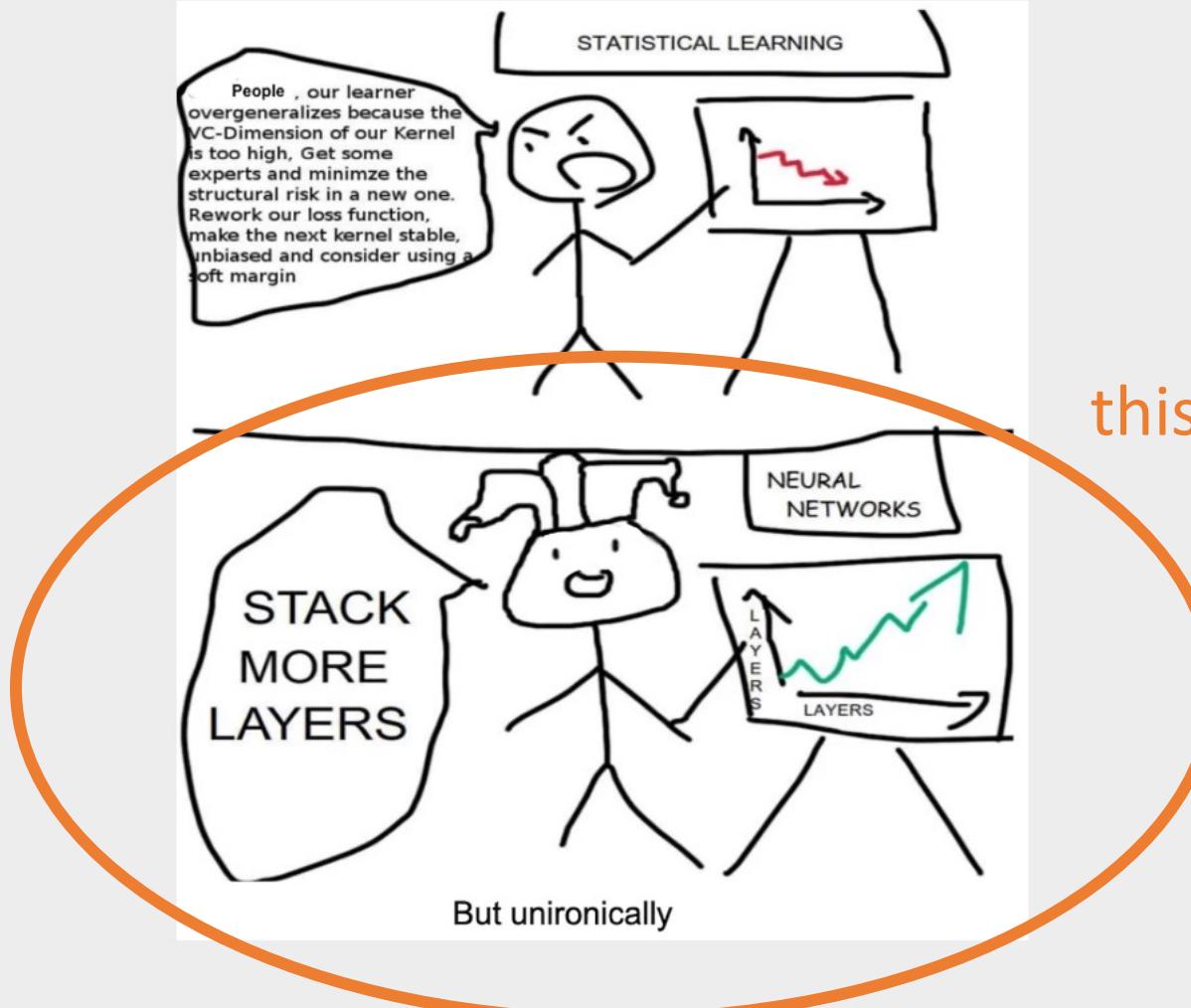


(f) ResNet-110-noshort

# Understanding Residual Learning



User Can you explain why this is funny. Think about it step-by-step.



GPT-4 The comic is satirizing the difference in approaches to improving model performance between statistical learning and neural networks.

# Lecture 10: Going deep with neural networks

- Why going deep is difficult?
- Weight initialization
- Normalization modules
- Deep Residual Learning