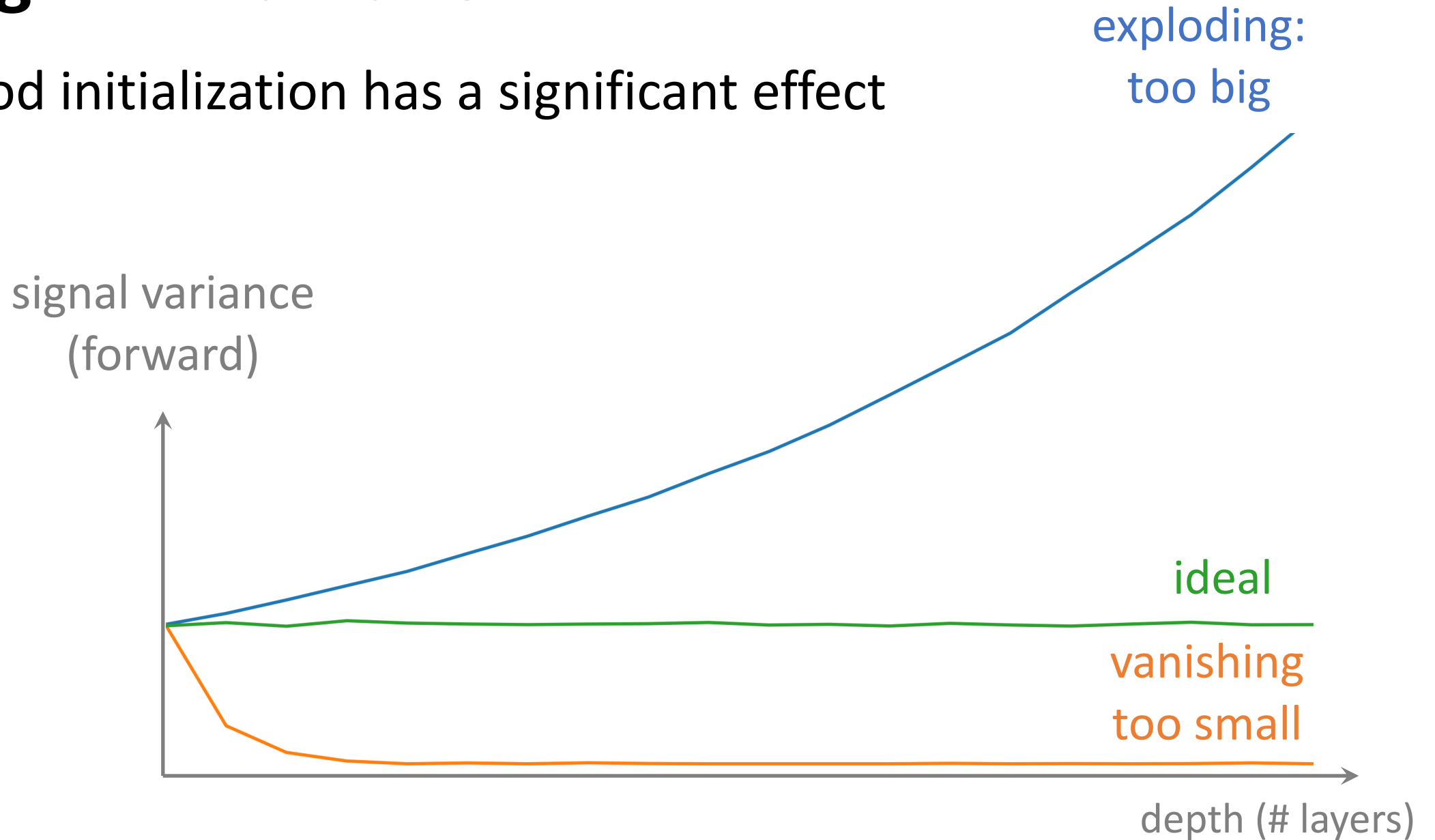


Weight Initialization

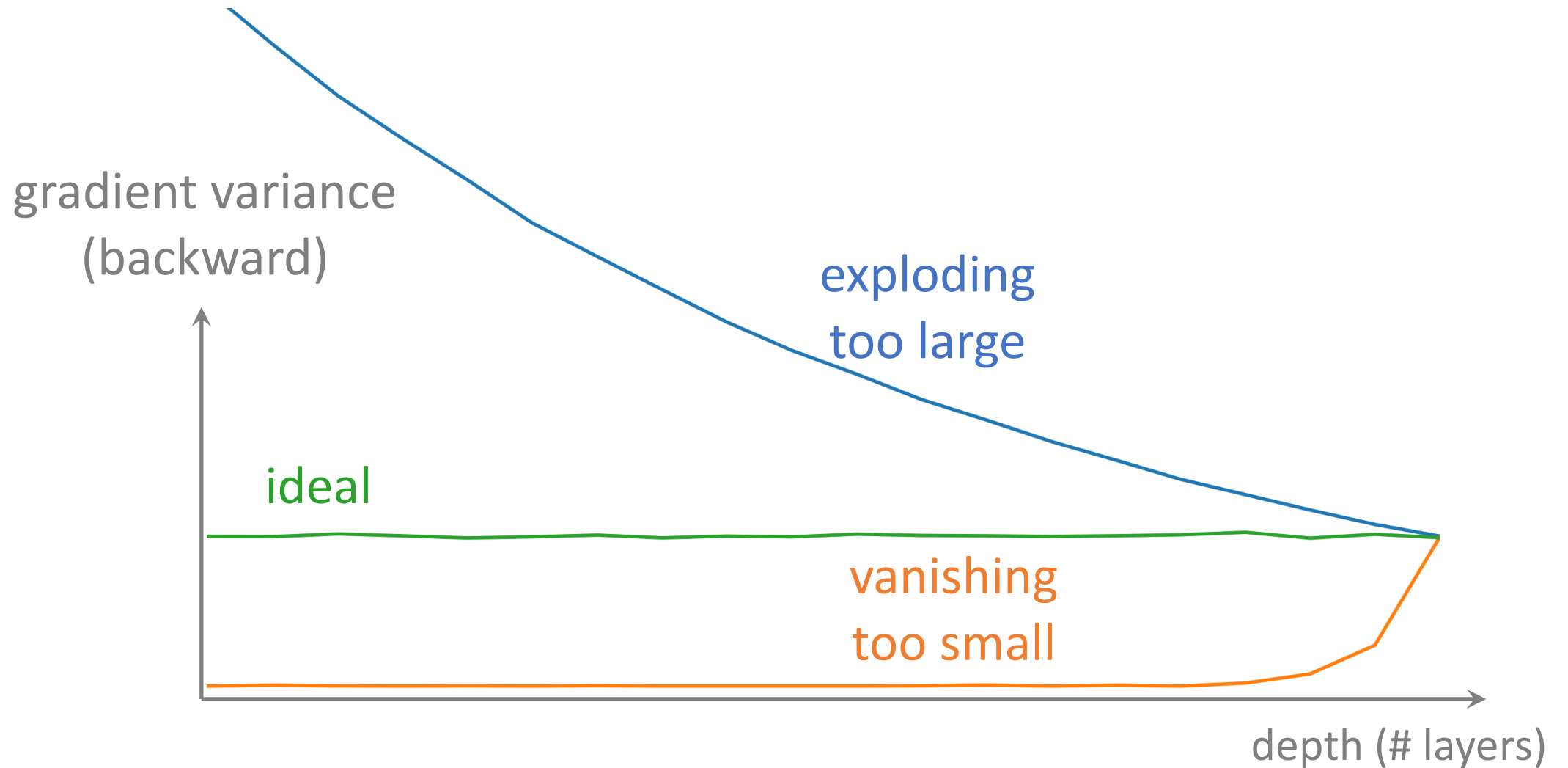
Weight initialization

- Good initialization has a significant effect

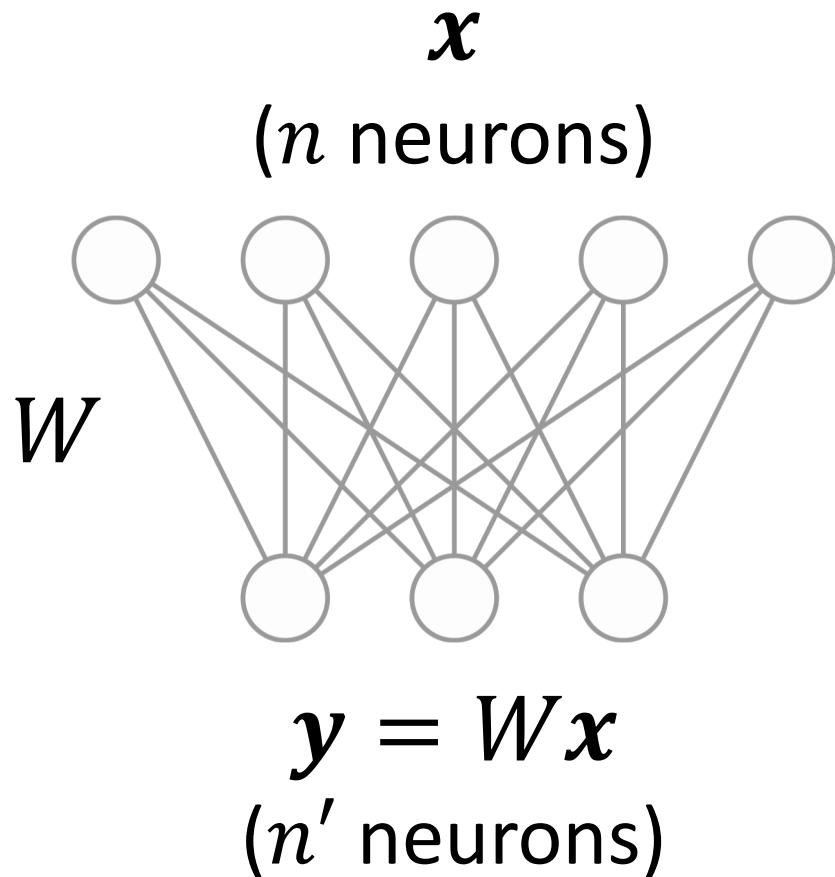


Weight initialization

- Good initialization has a significant effect



Weight initialization

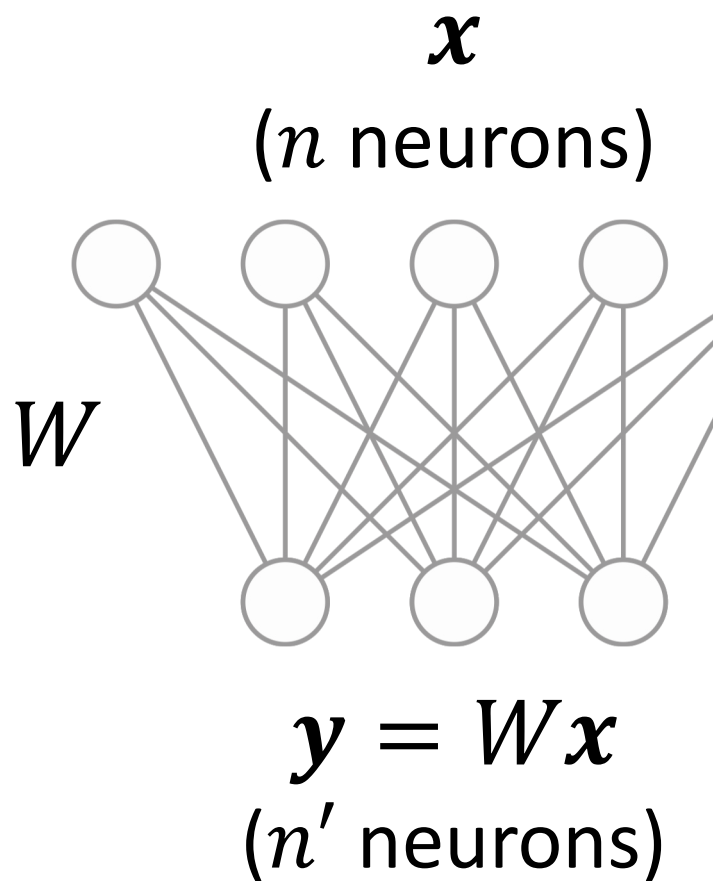


- if linear activation
- one layer: variance scaled by

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

Weight initiali

derivation



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

- 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

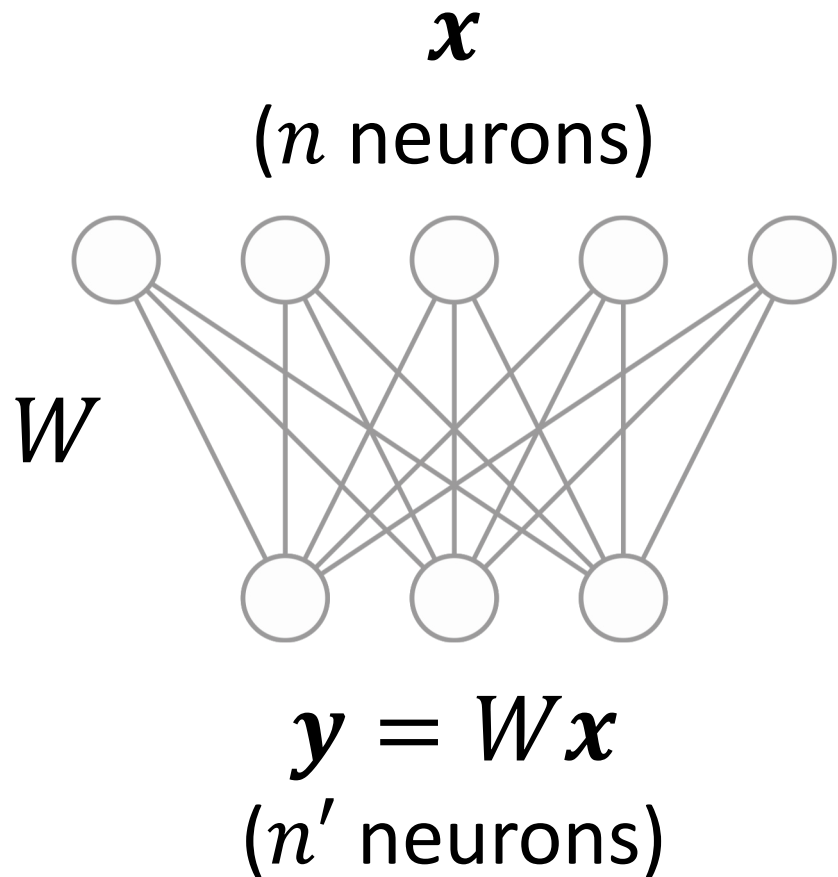
$$\text{Var}[y_i] = \sum_j \text{Var}[W_{ij}] \text{Var}[x_j]$$

- independence & zero-mean

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

- identical distributions

Weight initialization



- if linear activation
- one layer: variance scaled by

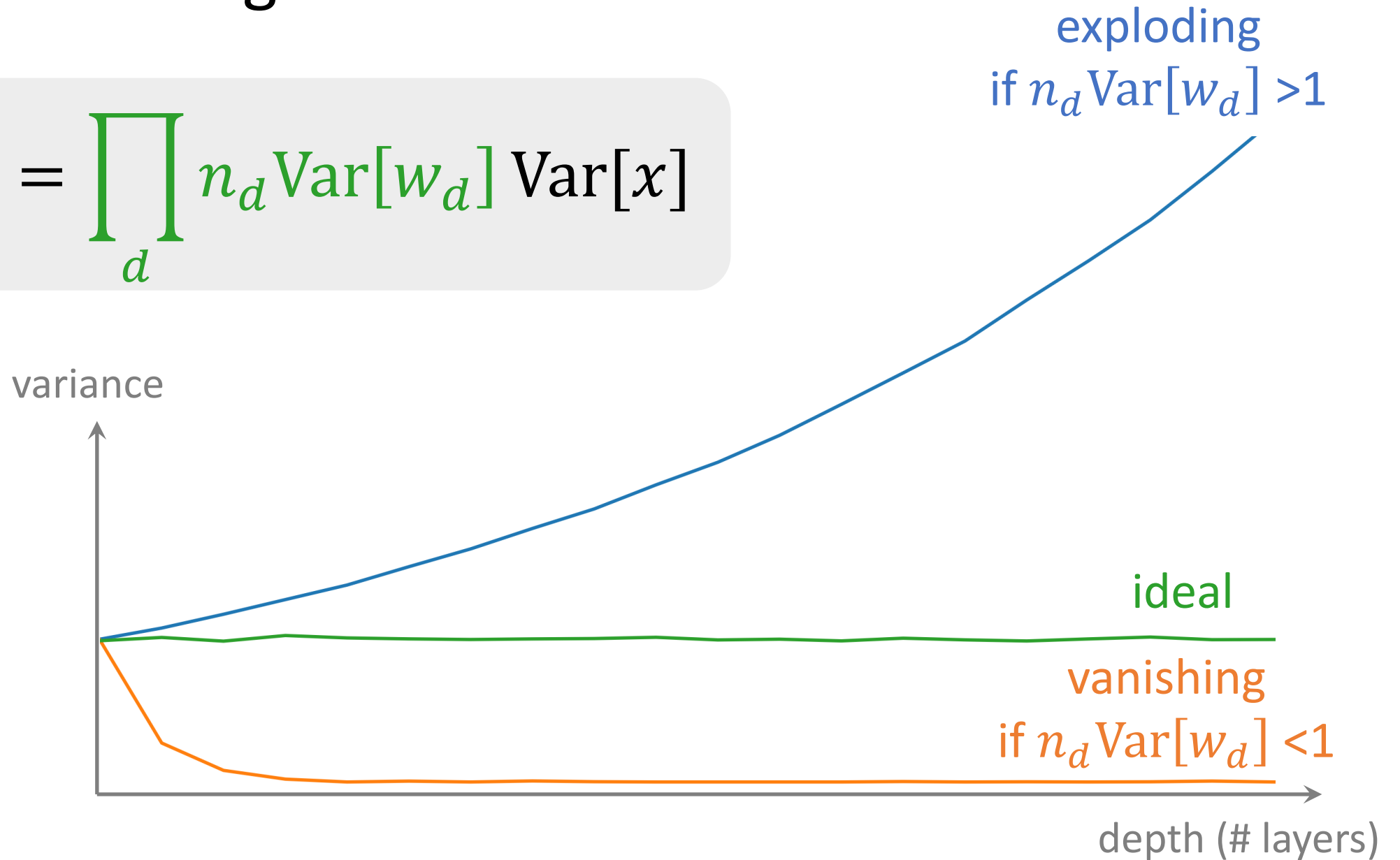
$$\text{Var}[\mathbf{y}] = n \text{Var}[\mathbf{w}] \text{Var}[\mathbf{x}]$$

- many layers: variance scaled by

$$\text{Var}[\mathbf{y}] = \prod_d n_d \text{Var}[\mathbf{w}_d] \text{Var}[\mathbf{x}]$$

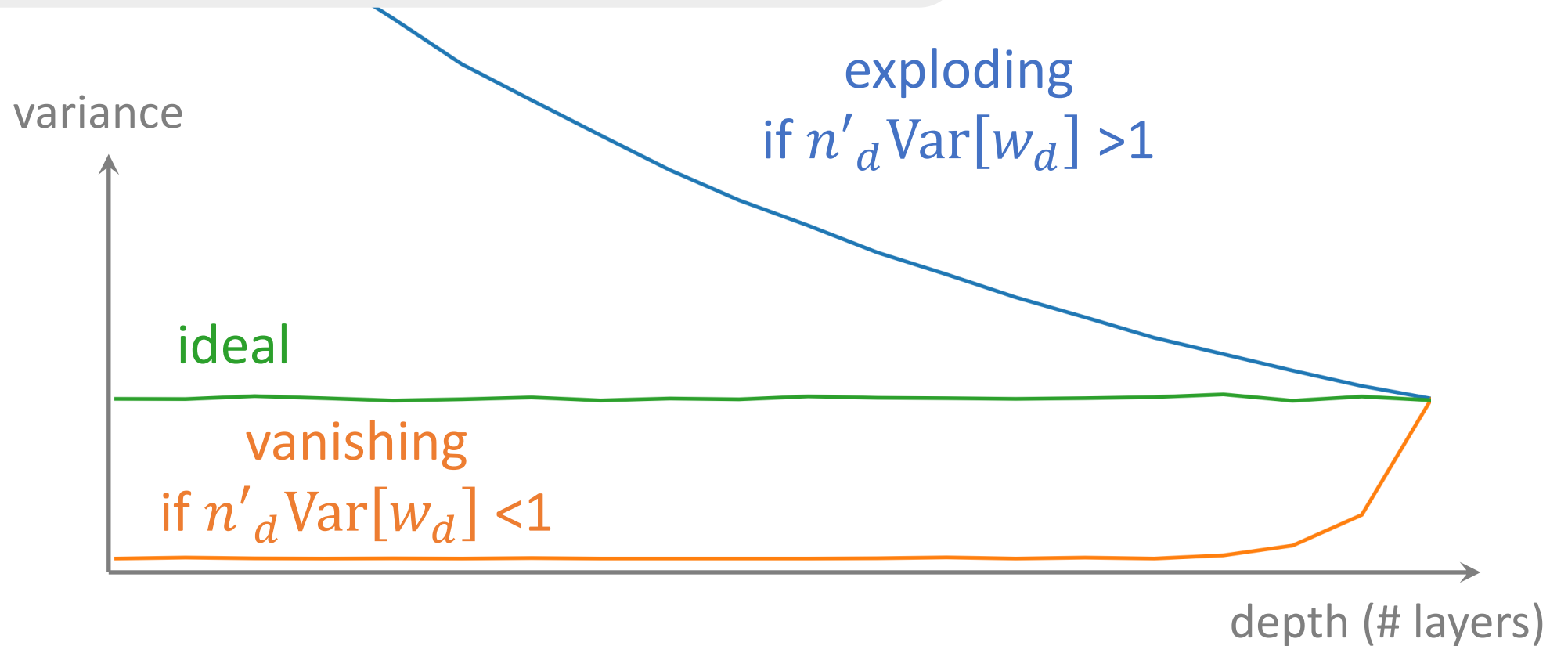
Variance Scaling: forward

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



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:

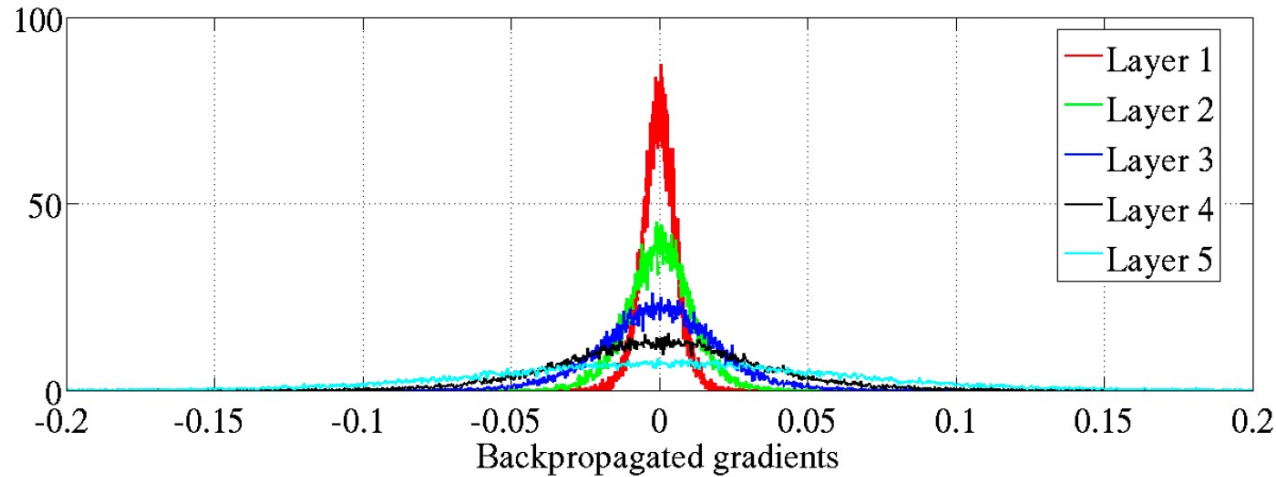
$$w \sim \mathcal{N}(\mu = 0, \sigma = \sqrt{1/n})$$

- Uniform distribution:

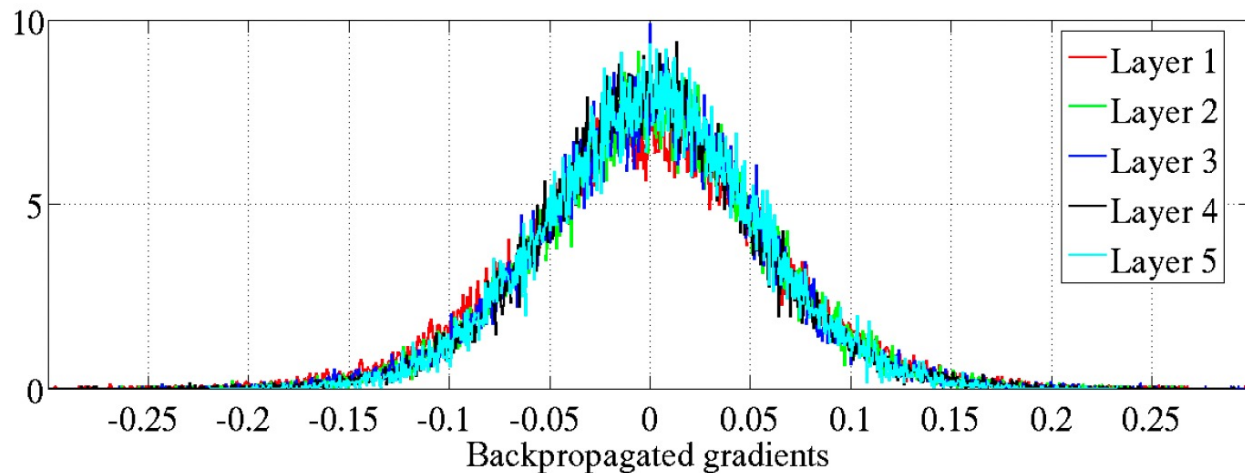
$$w \sim \mathcal{U}(-a, +a), a = \sqrt{3/n}$$

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

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



poor initialization:
earlier layer has smaller gradients

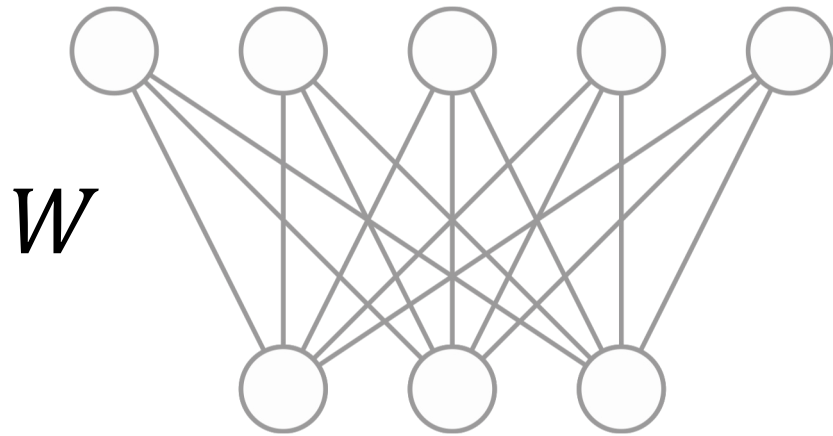


Xavier initialization:
all layers have similar gradient scale

Weight initialization: ReLU

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

(n neurons)



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

(n' neurons)

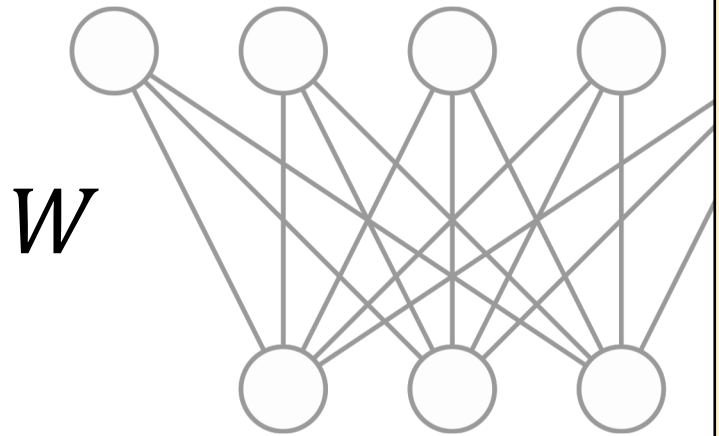
- if **ReLU** activation
- one layer: variance scaled by

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

Weight initialization

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

(n neurons)



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

(n' neurons)

derivation

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

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

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

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

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

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

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

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

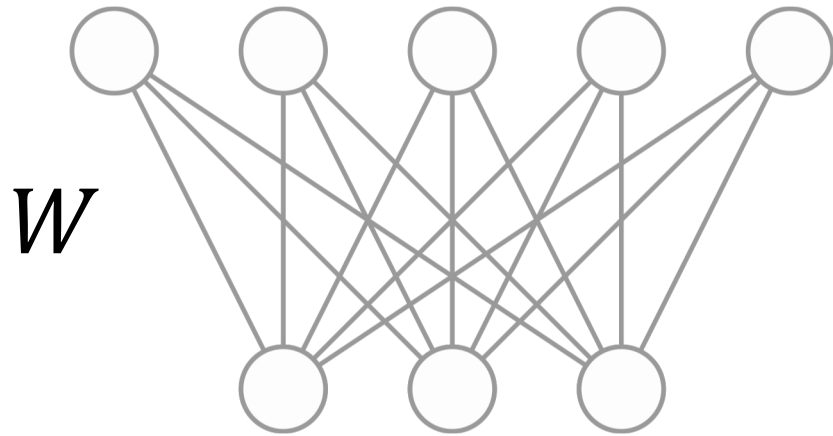
$$E[x'^2] = \text{Var}[x]/2$$

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:

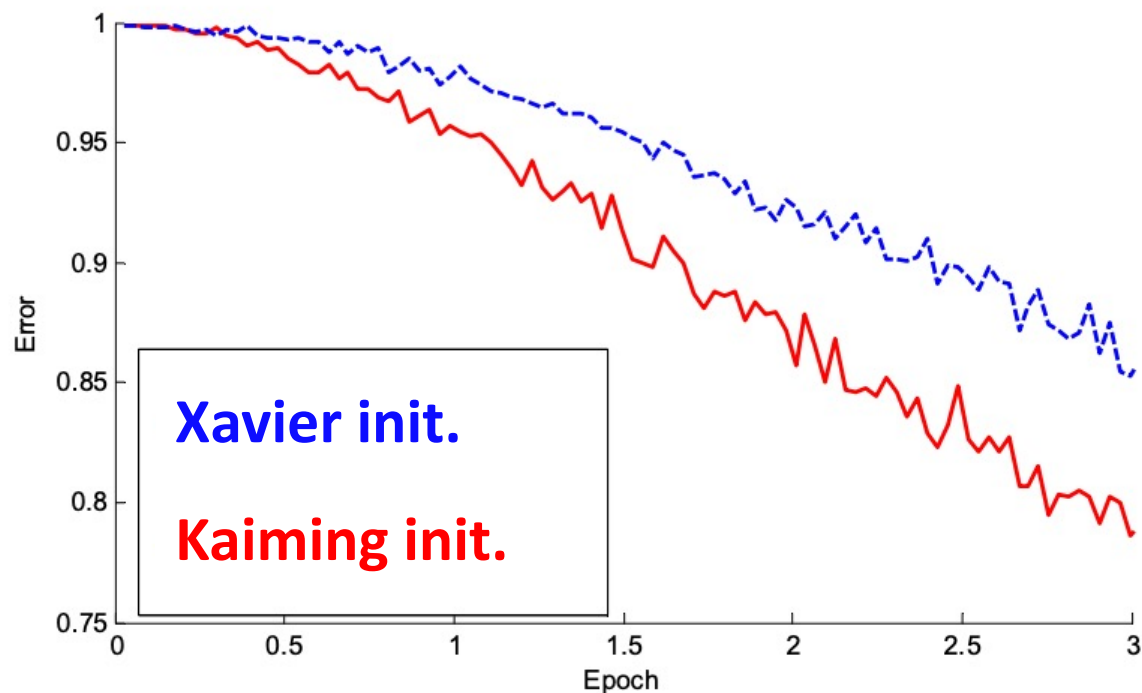
$$w \sim \mathcal{N}(\mu = 0, \sigma = \sqrt{2/n})$$

- Uniform distribution:

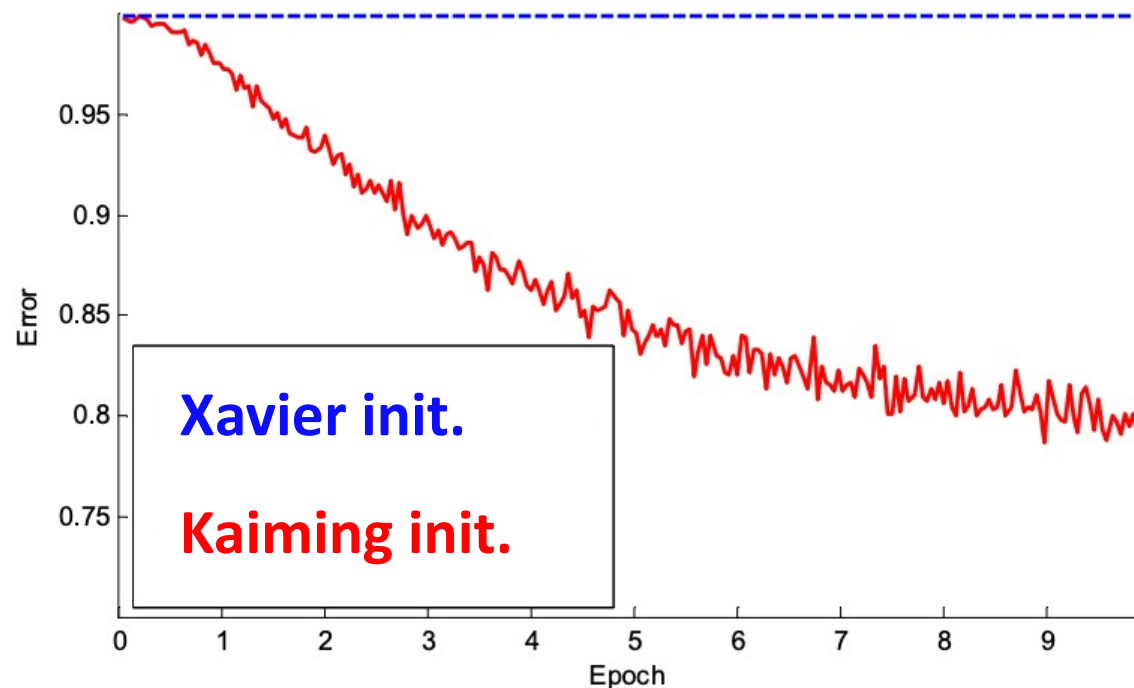
$$w \sim \mathcal{U}(-a, +a), a = \sqrt{6/n}$$

- sufficient to use n or n'

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



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

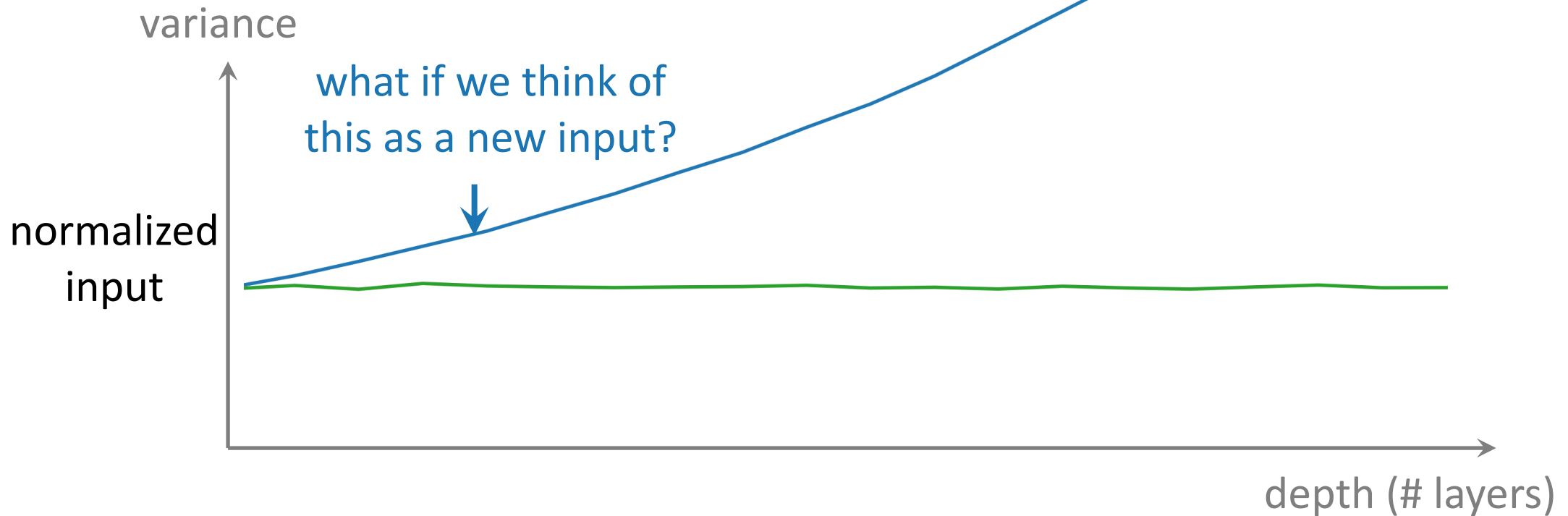


30-layer VGG w/ ReLU:
better init enables training

Normalization Modules

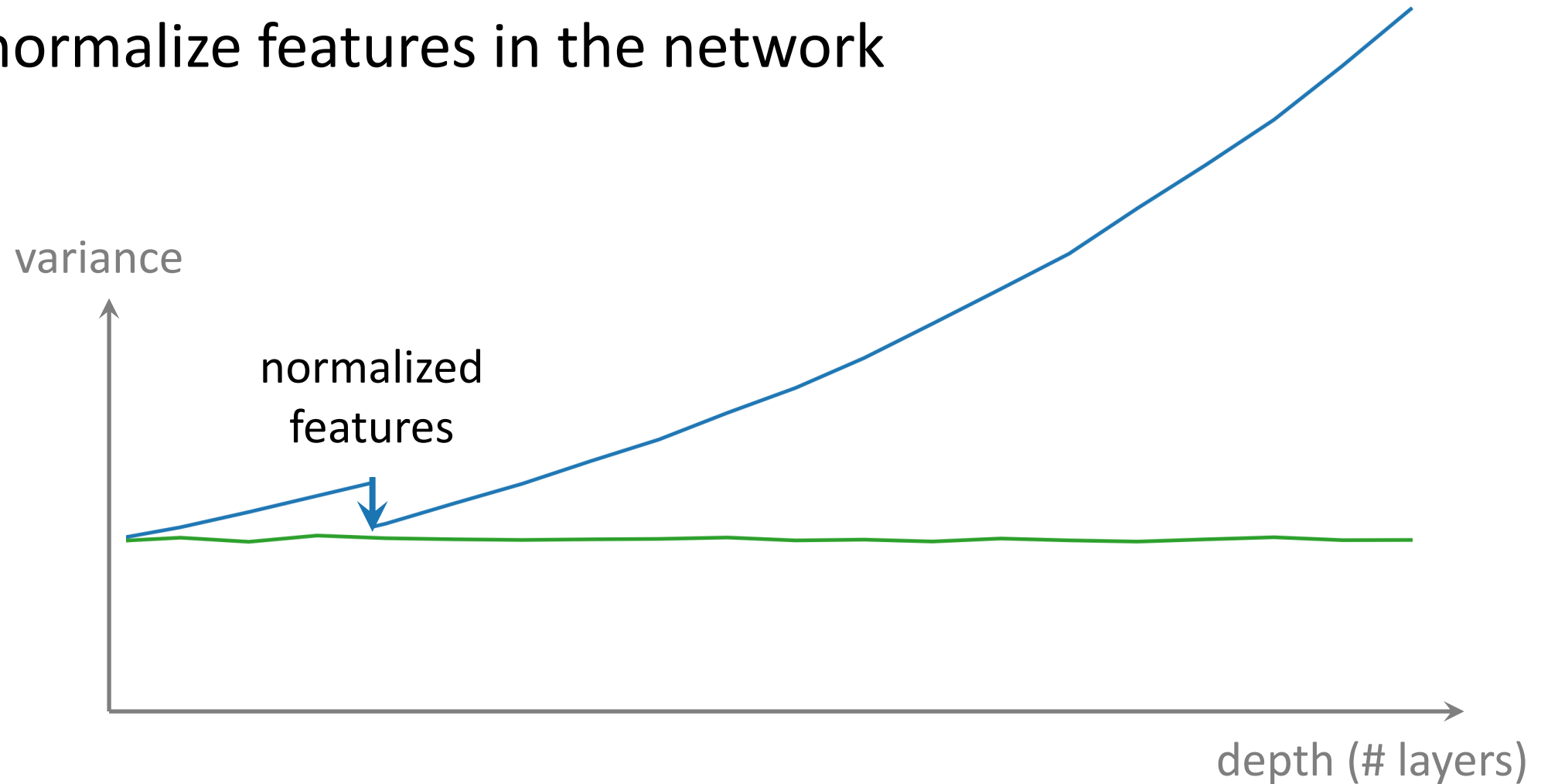
Normalization Modules

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



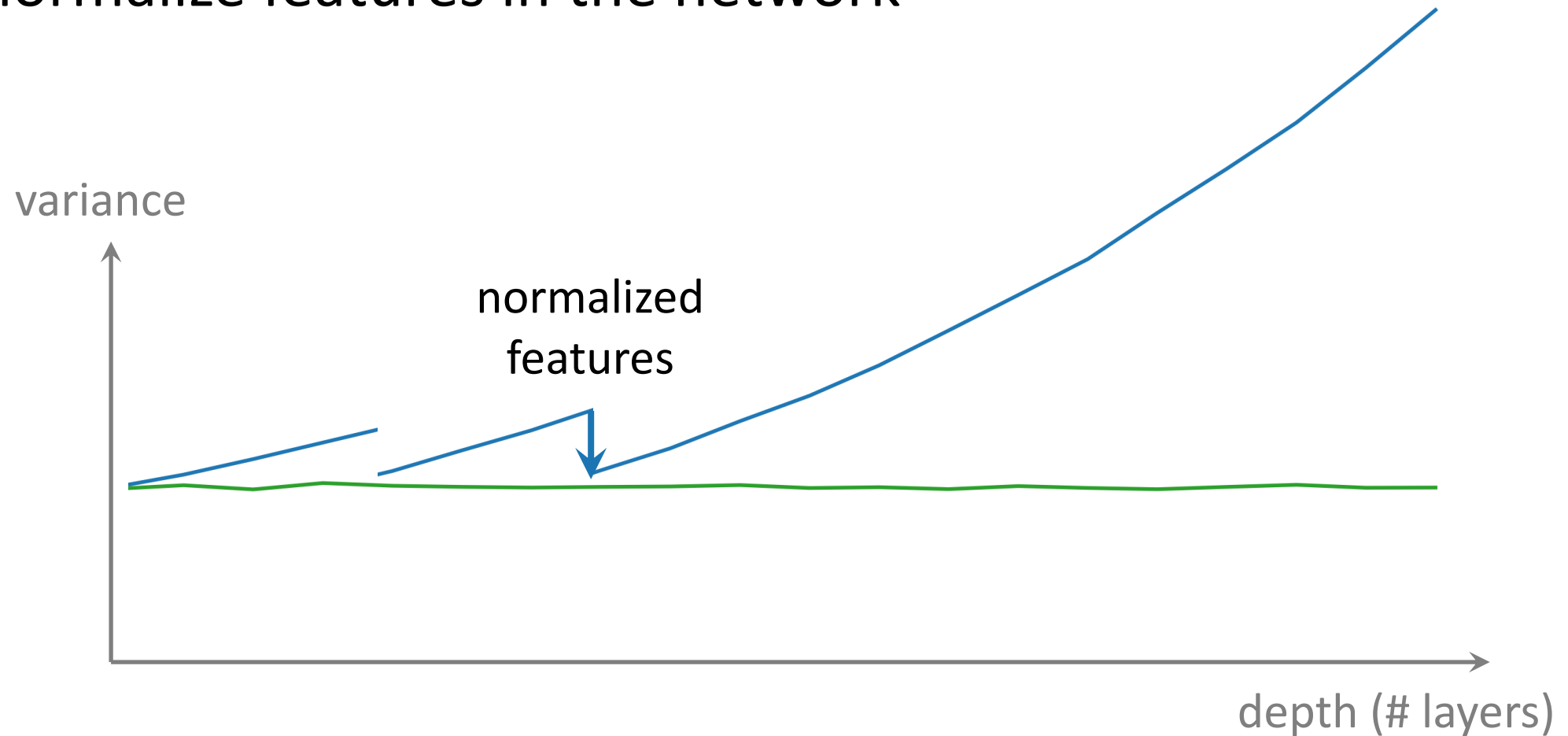
Normalization Modules

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



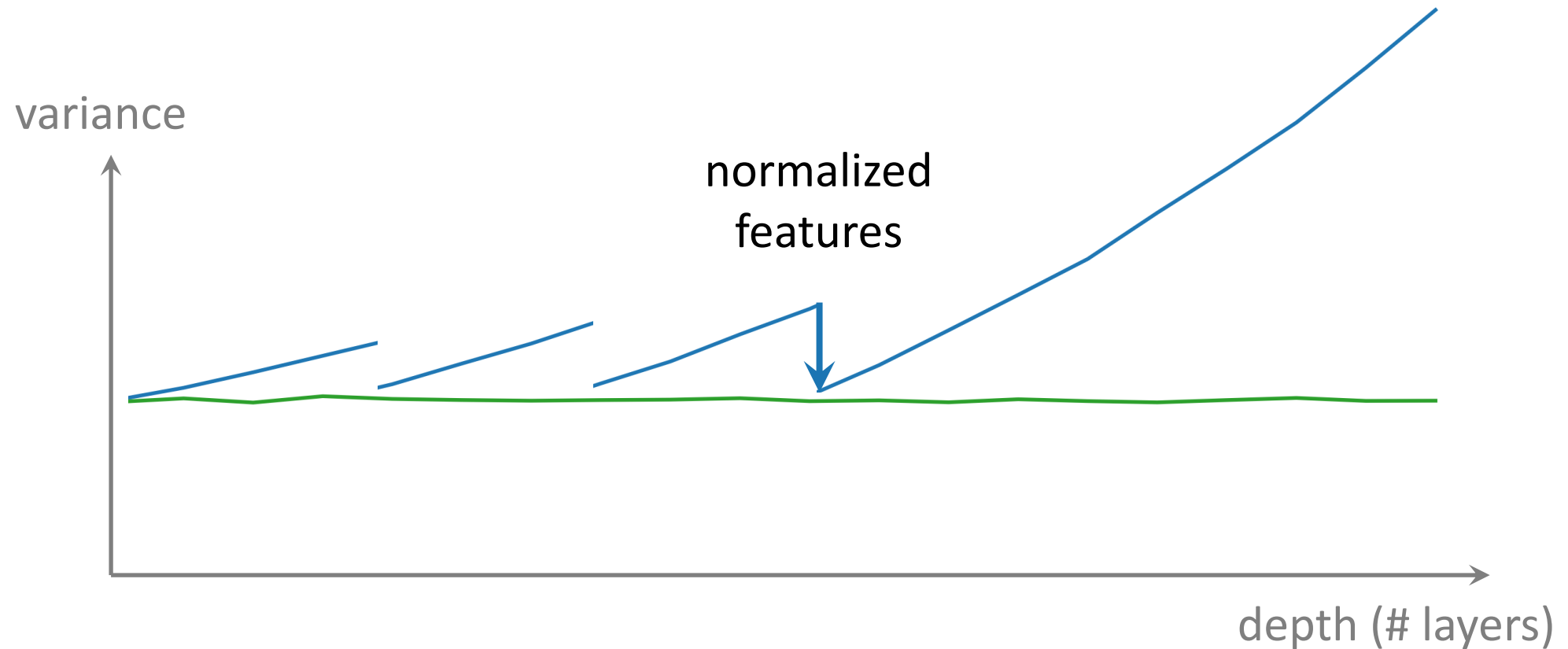
Normalization Modules

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



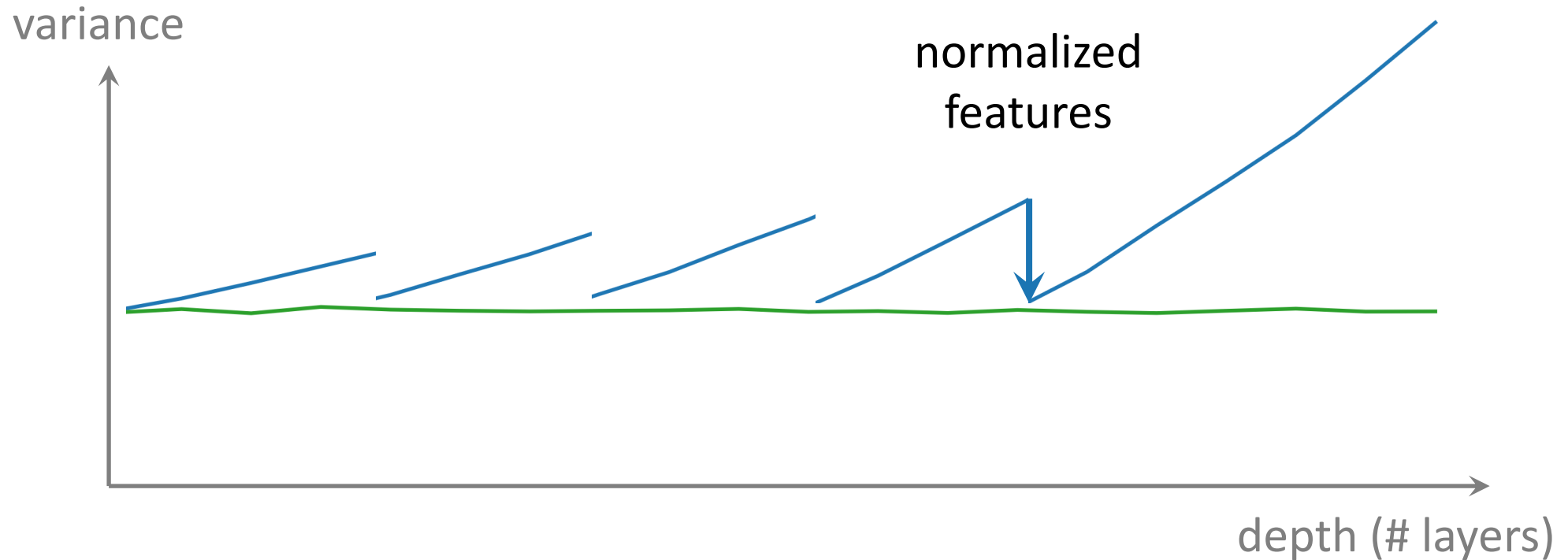
Normalization Modules

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



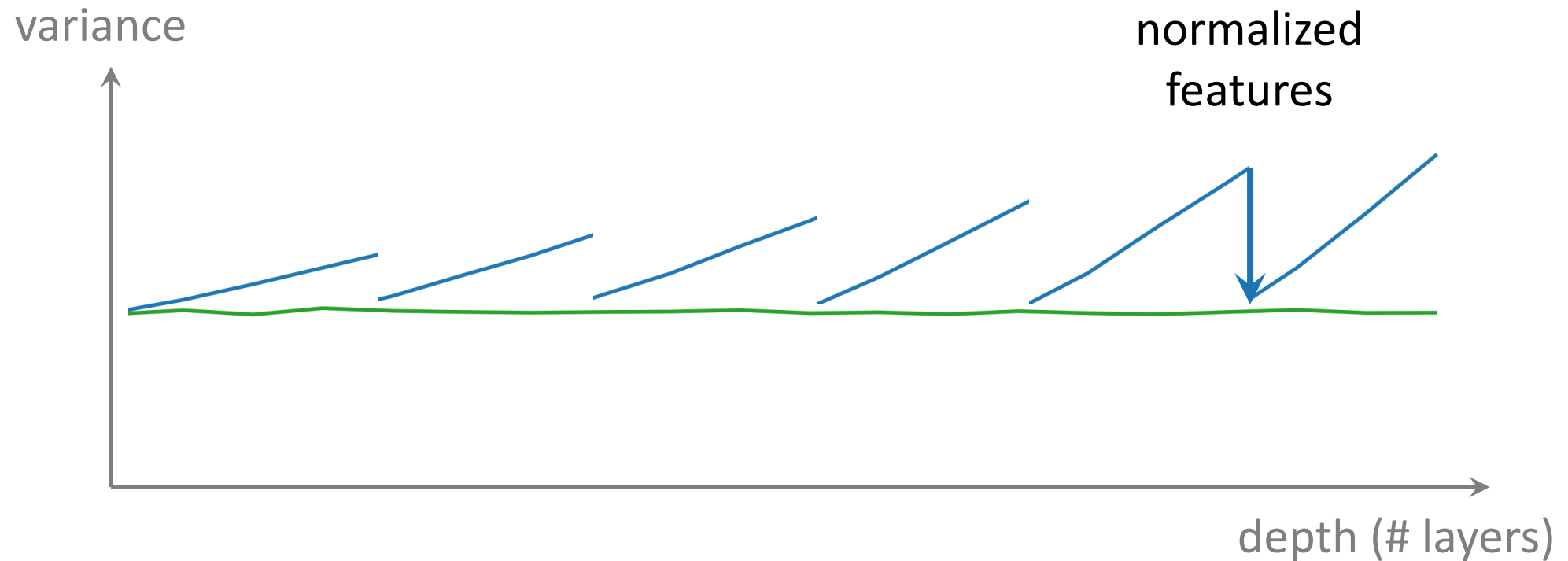
Normalization Modules

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



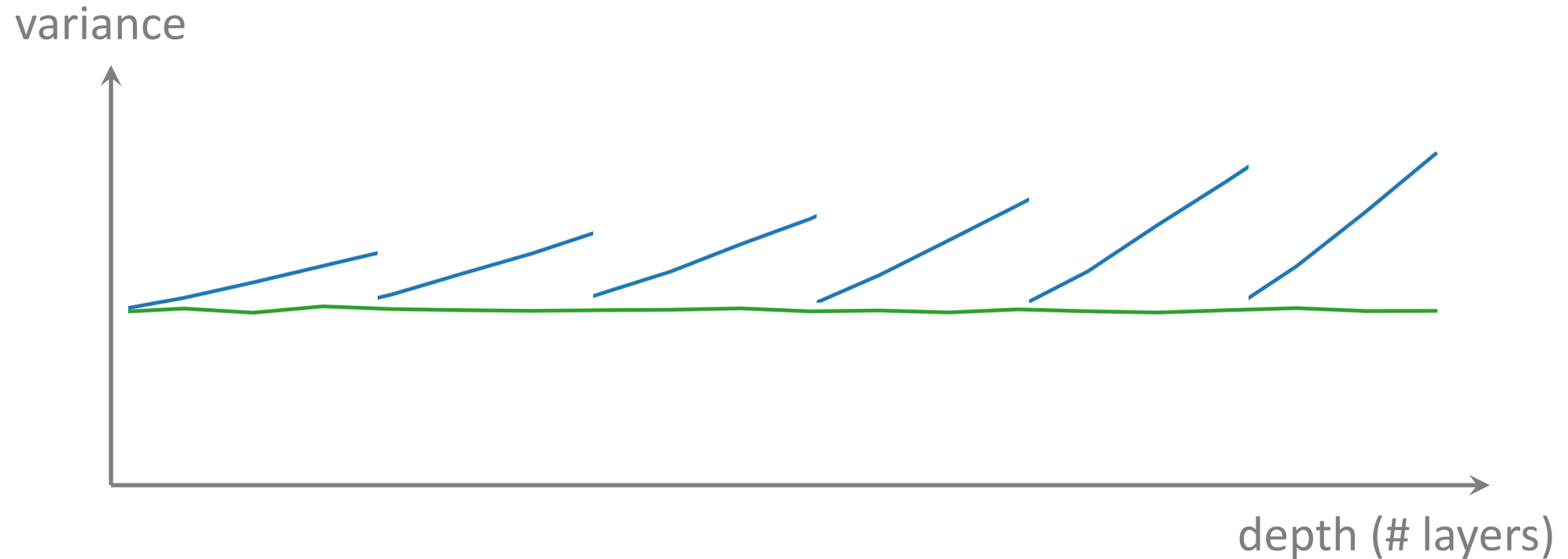
Normalization Modules

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



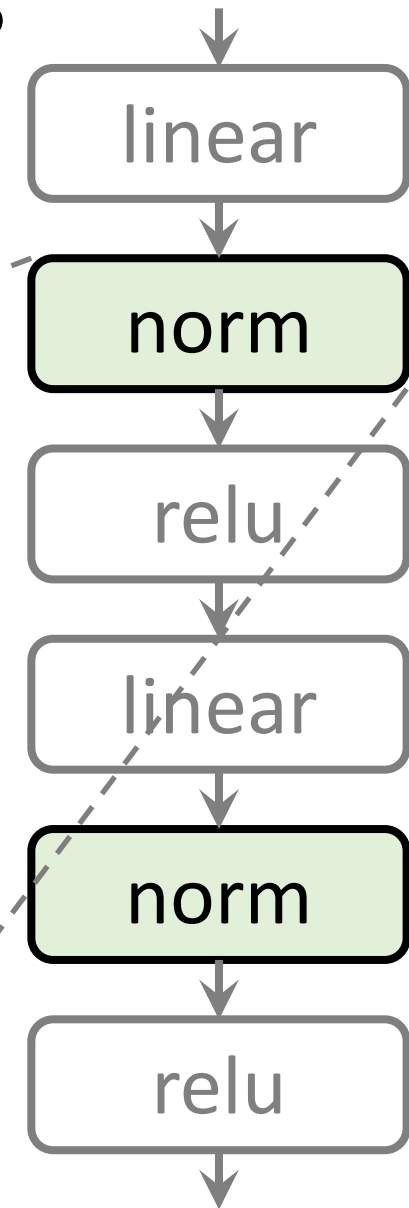
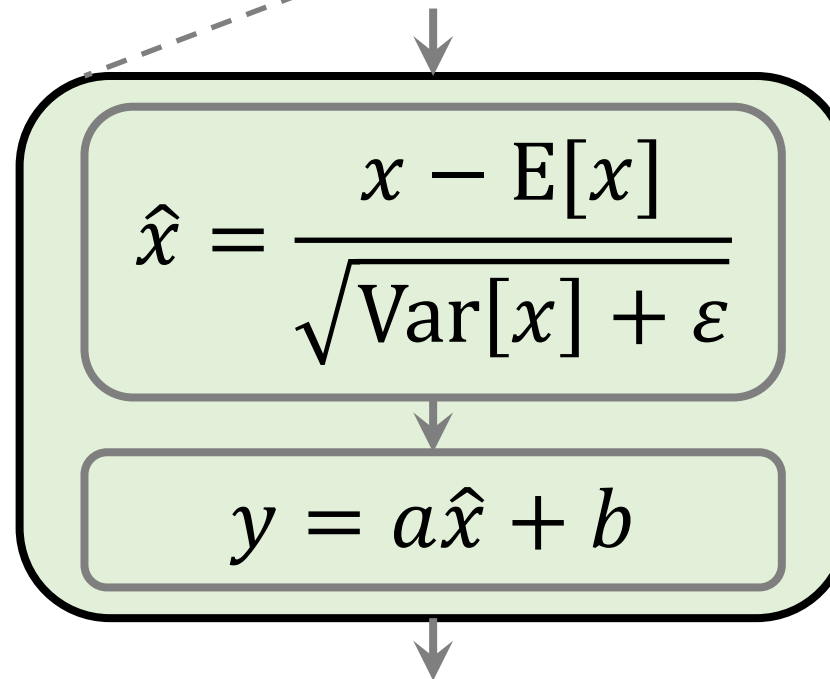
Normalization Modules

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



Normalization Modules: Operations

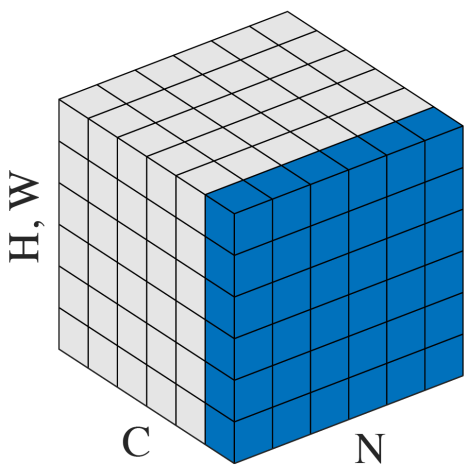
1. compute $E[x]$ and $\text{Var}[x]$
2. normalize by $E[x]$ and $\text{std}[x]$
3. compensate by a linear transform



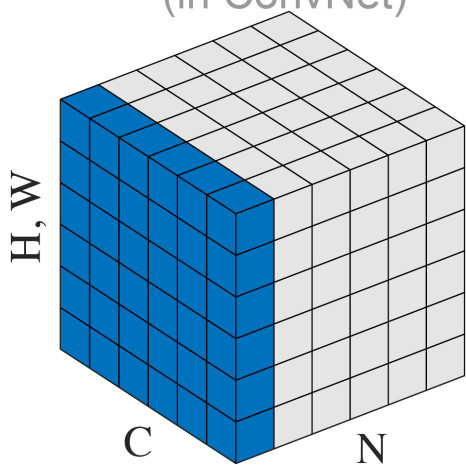
Normalization Modules: Variants

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

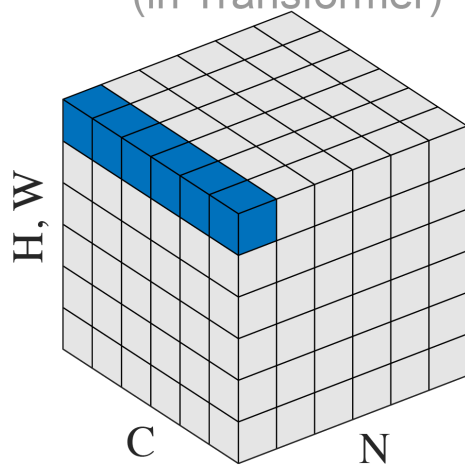
BatchNorm



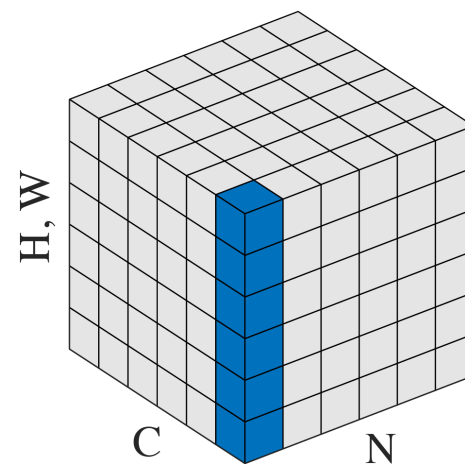
LayerNorm
(in ConvNet)



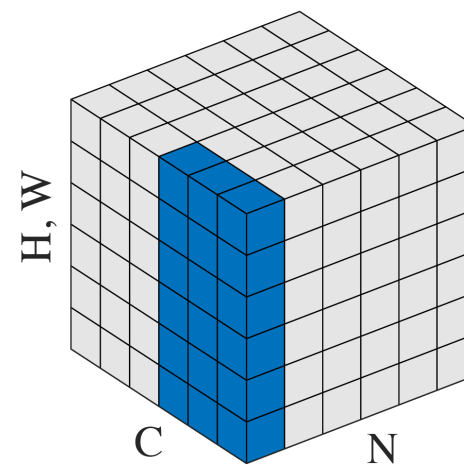
LayerNorm
(in Transformer)



InstanceNorm



GroupNorm



Normalization Modules: Effects

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

