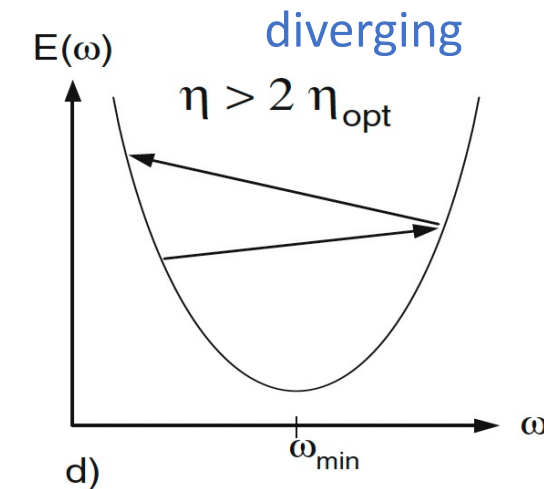
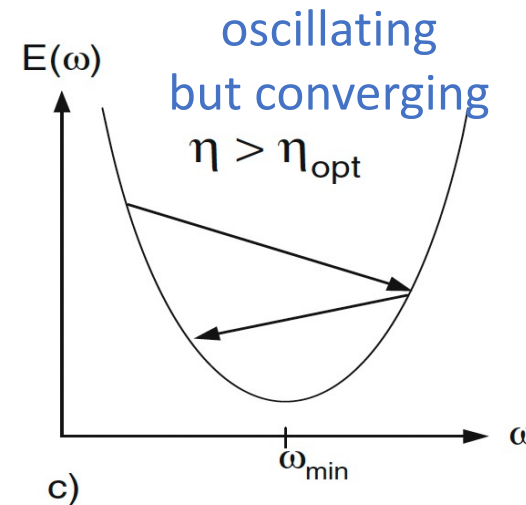
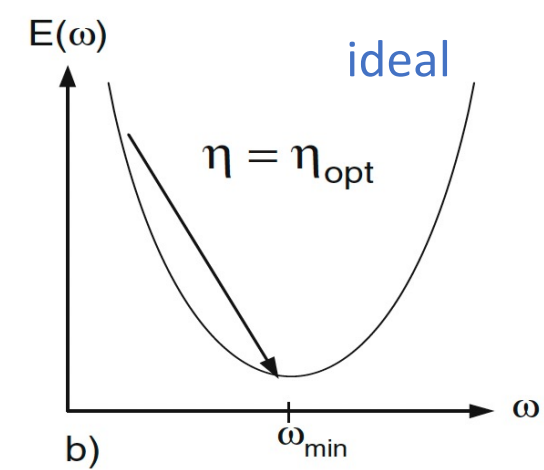
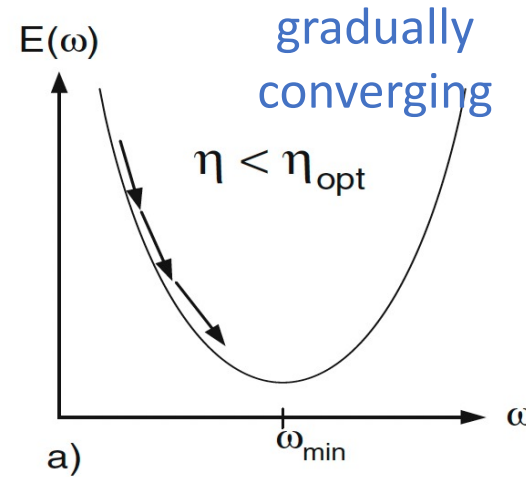


# Normalization Modules

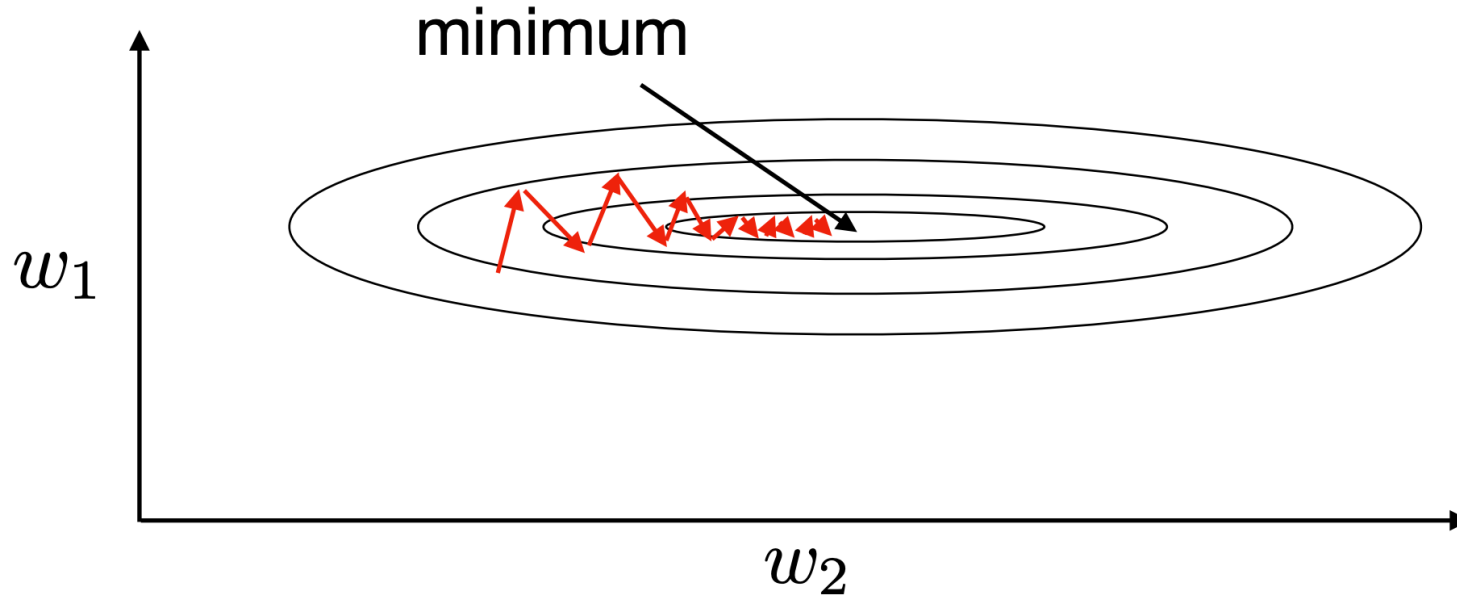
# Learning Rate in Gradient Descent

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

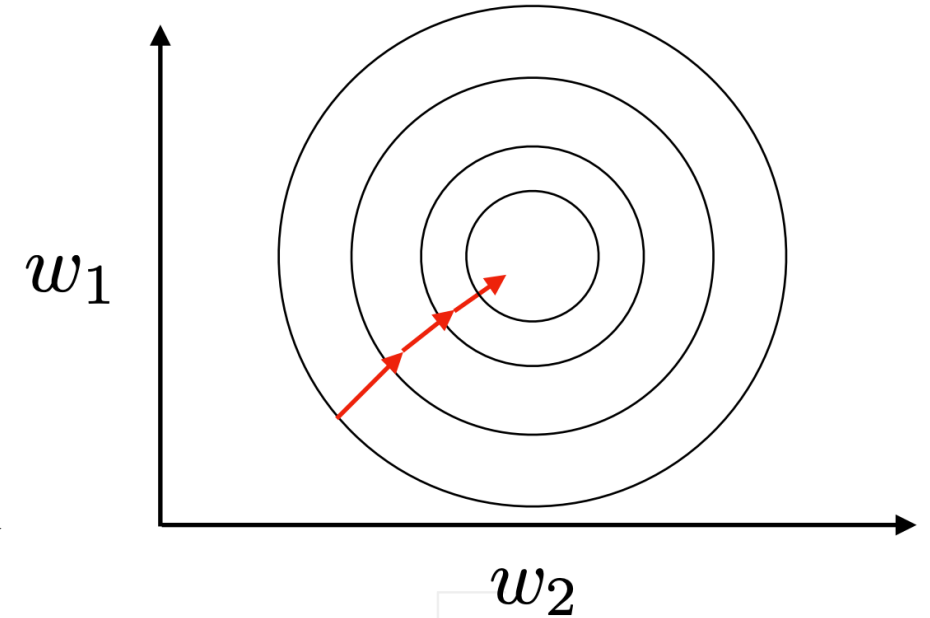
$\eta$ : learning rate



# Normalization

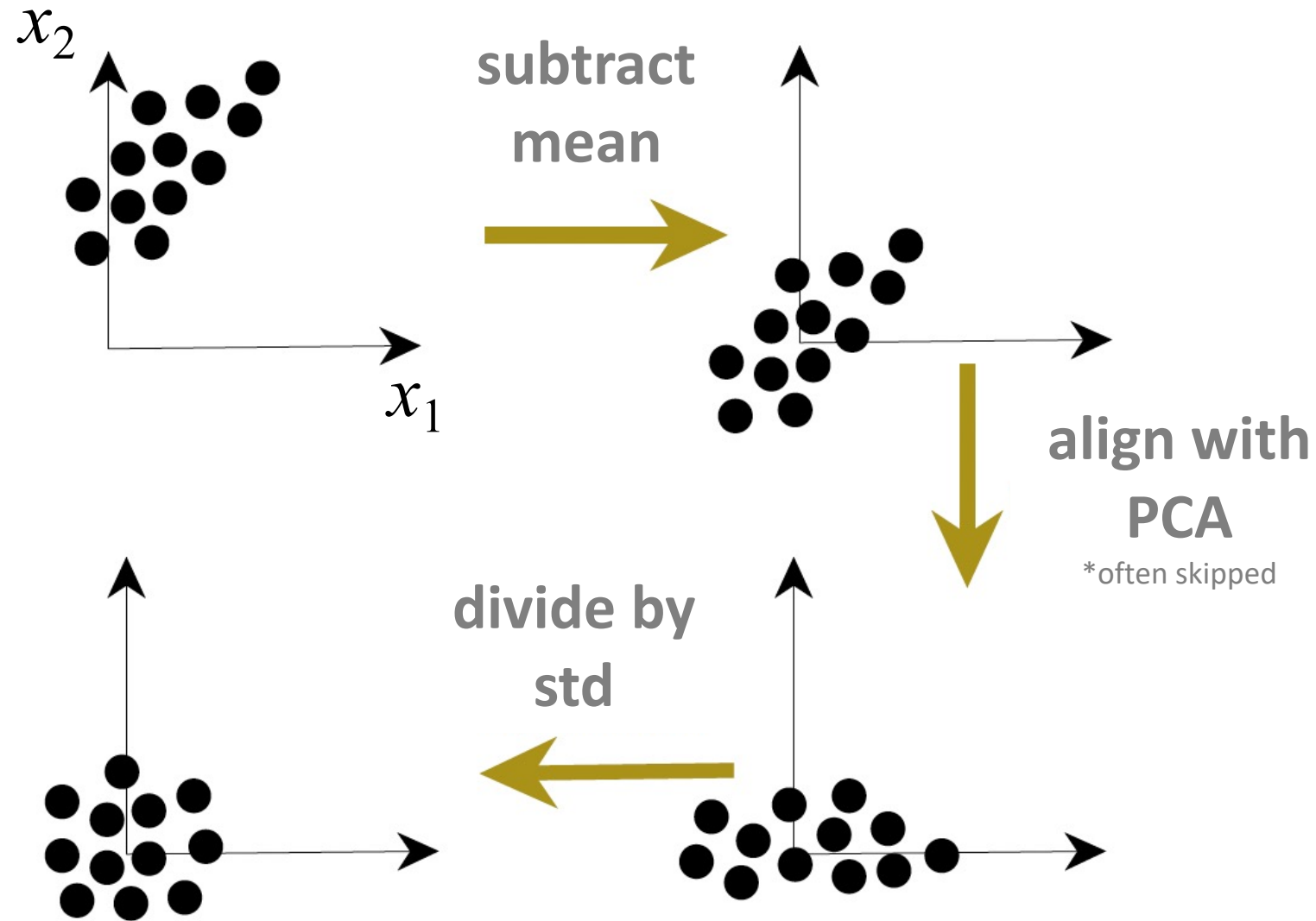


- Same learning rate applied to all weights
- Large weights dominate updates
- Small weights oscillate (or diverge)



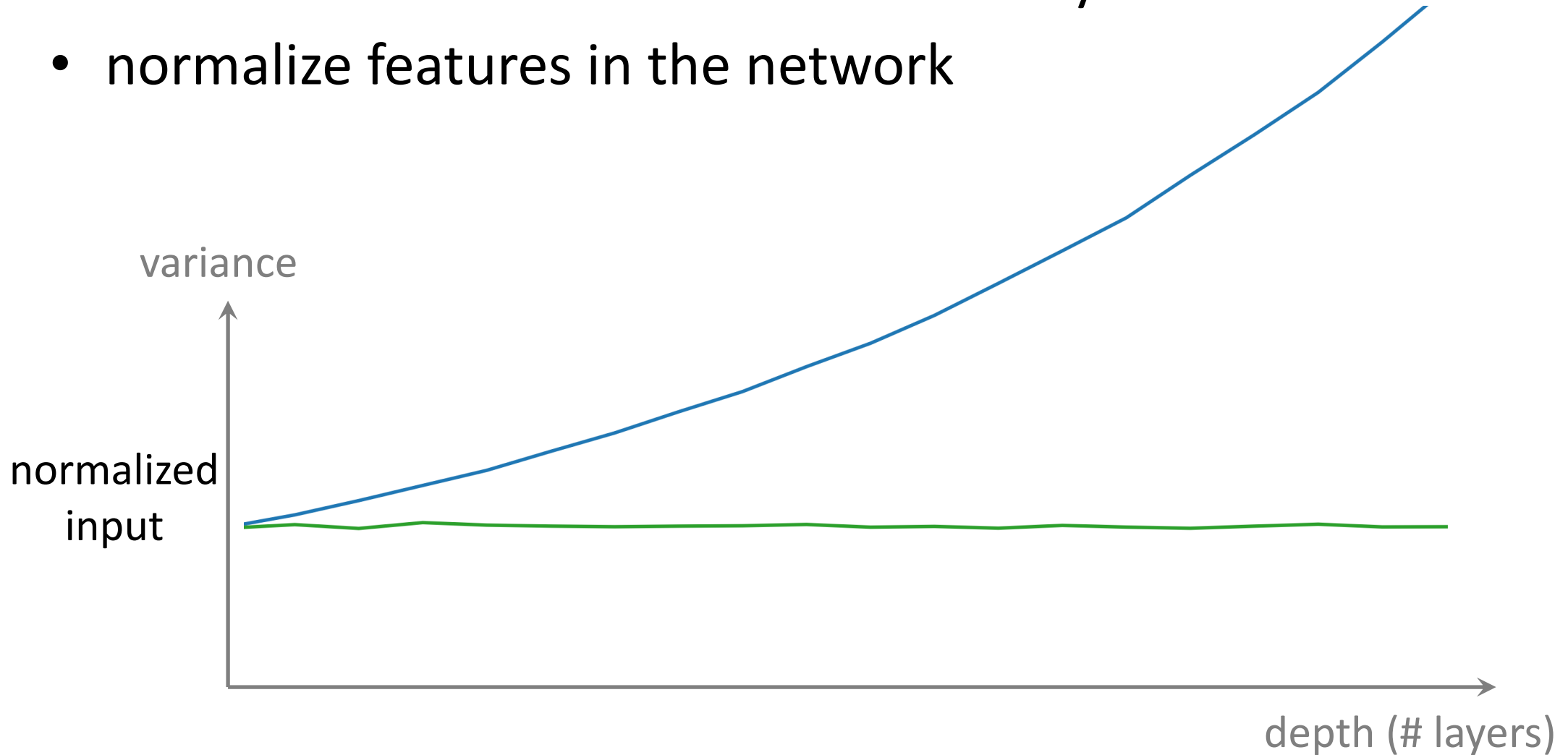
- Similar pace for all weights

# Input Normalization



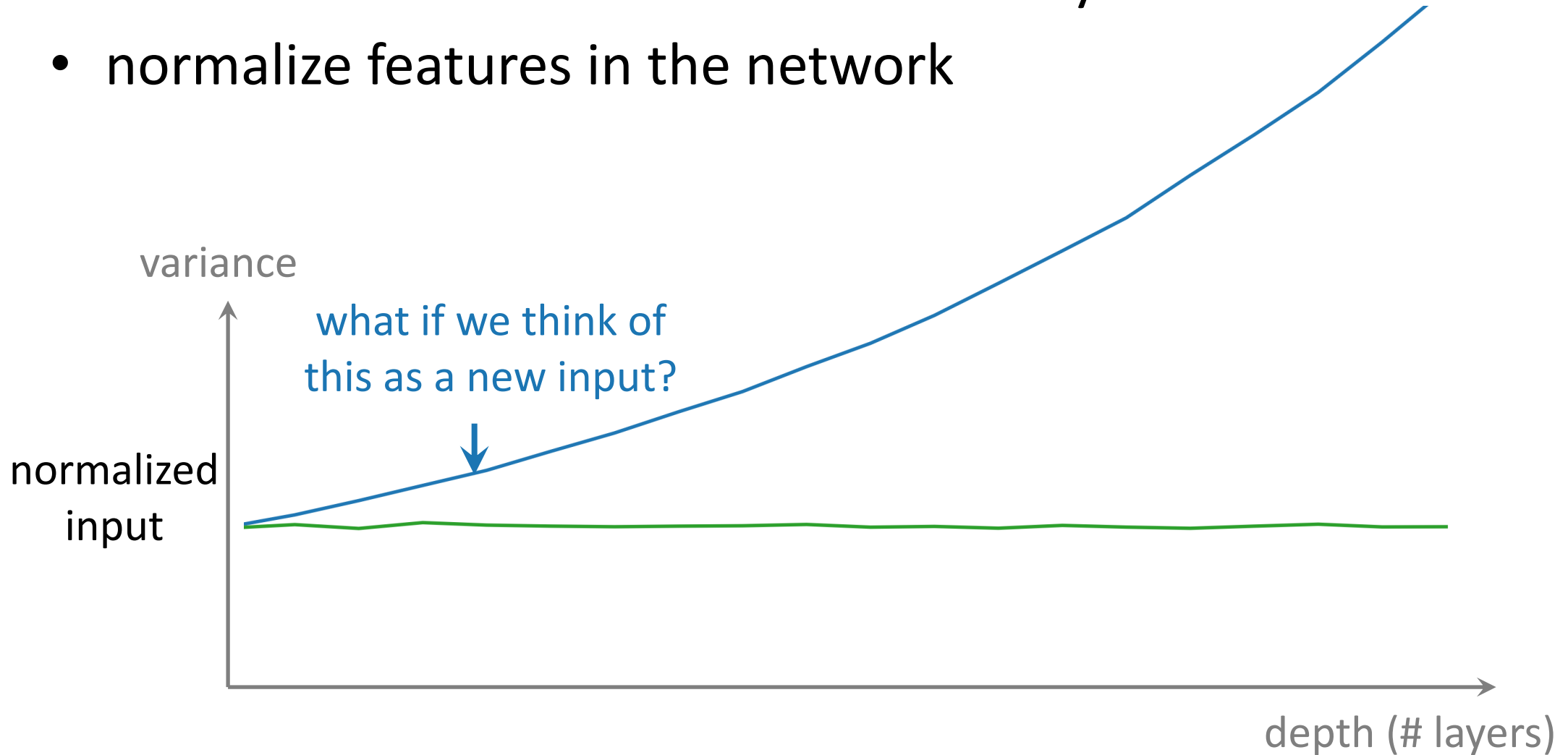
# Normalization Modules

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



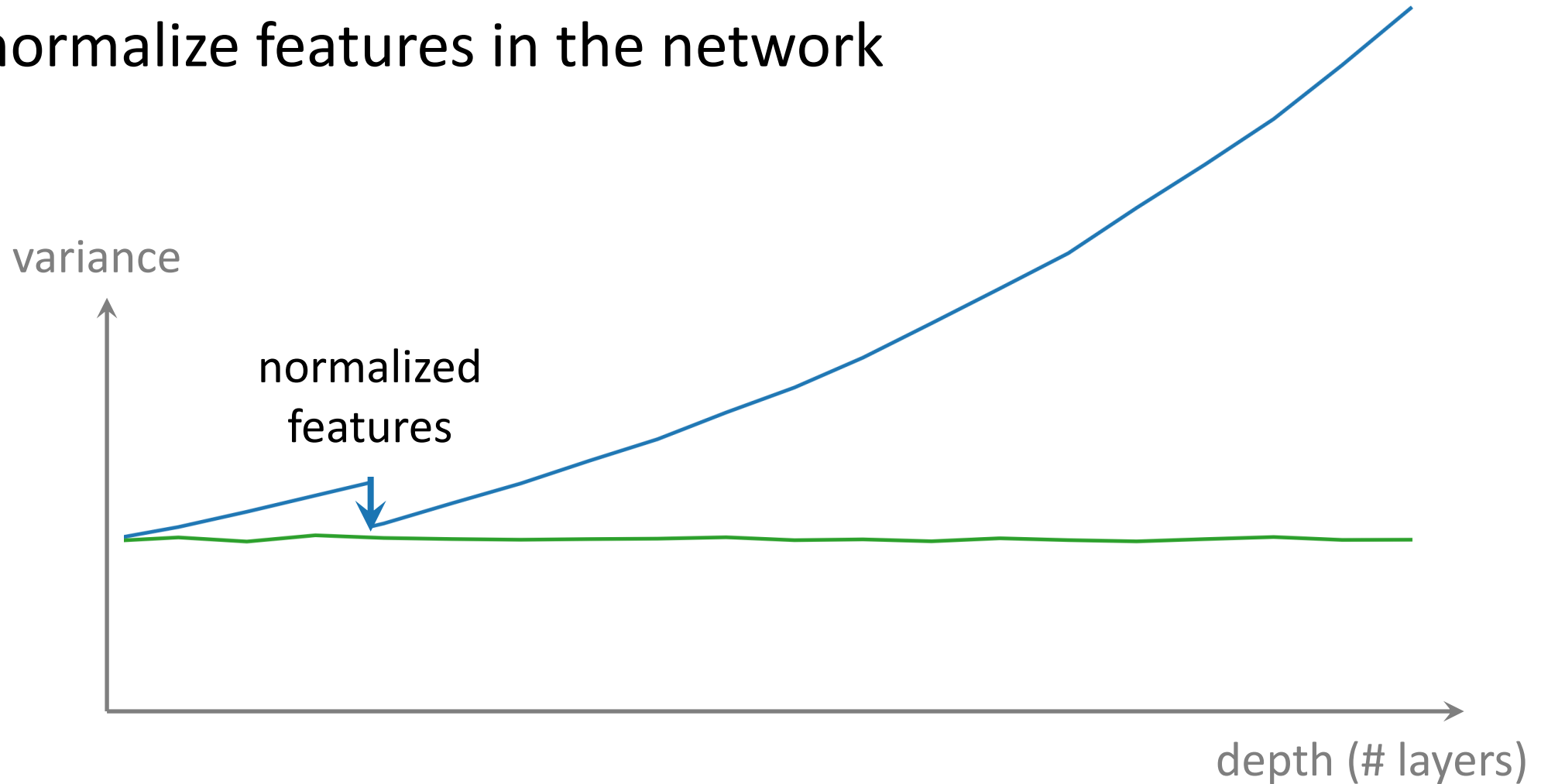
# Normalization Modules

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



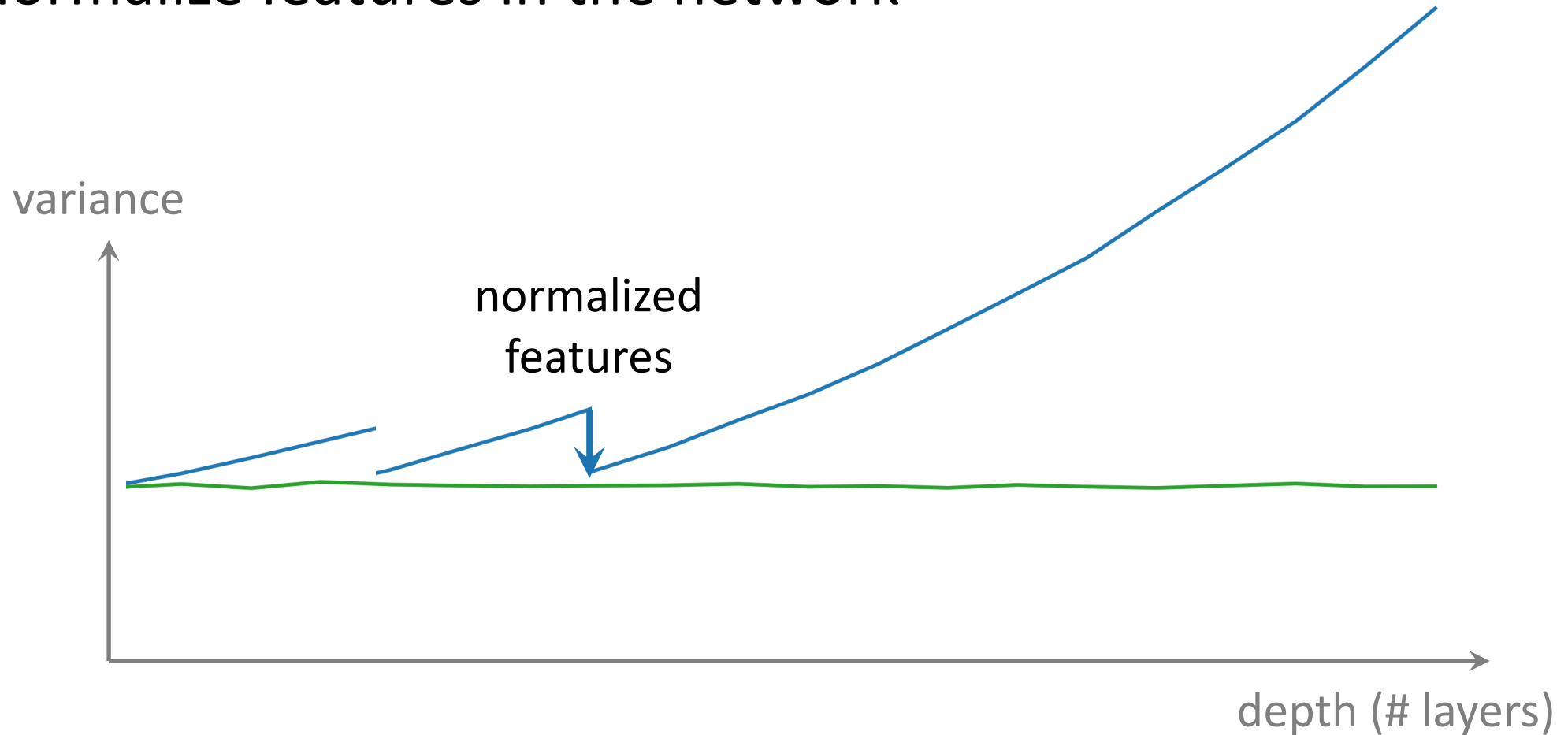
# Normalization Modules

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



# Normalization Modules

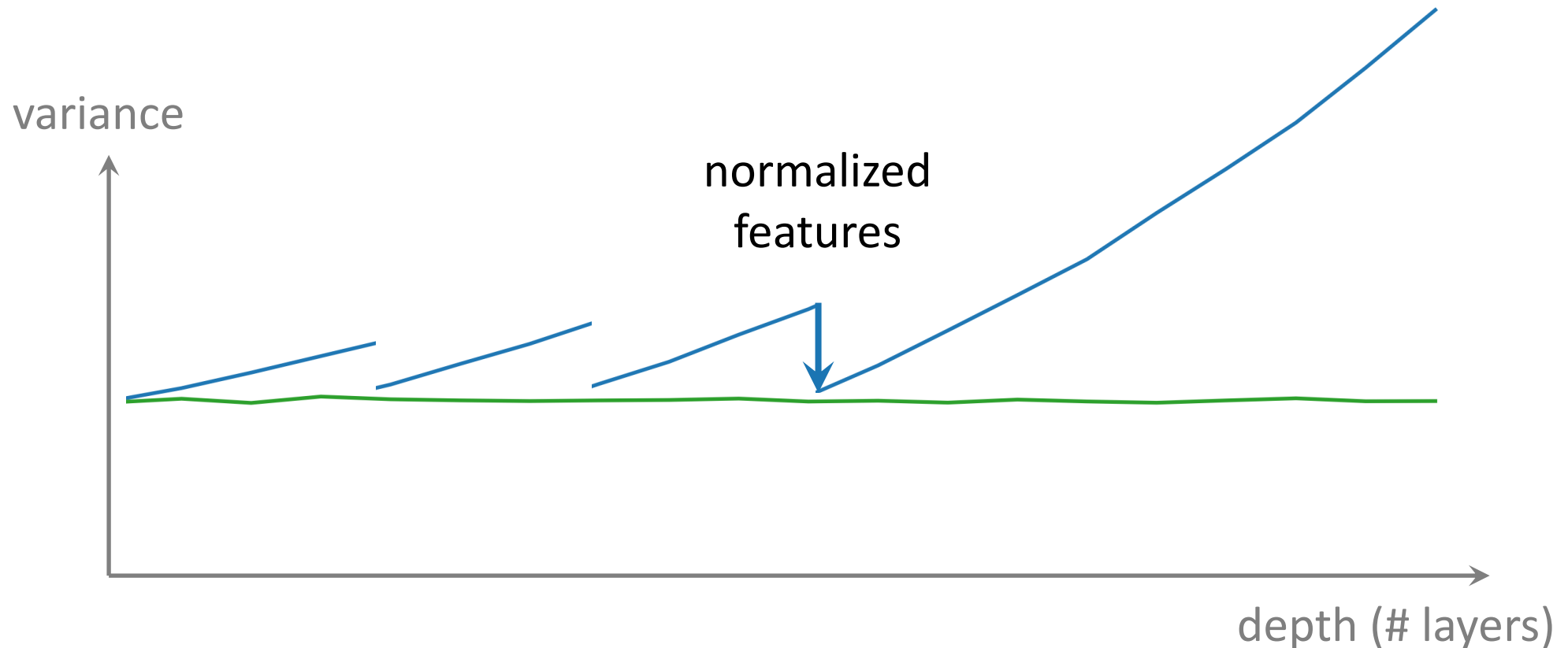
- We want to maintain variance for all layers
- normalize features in the network





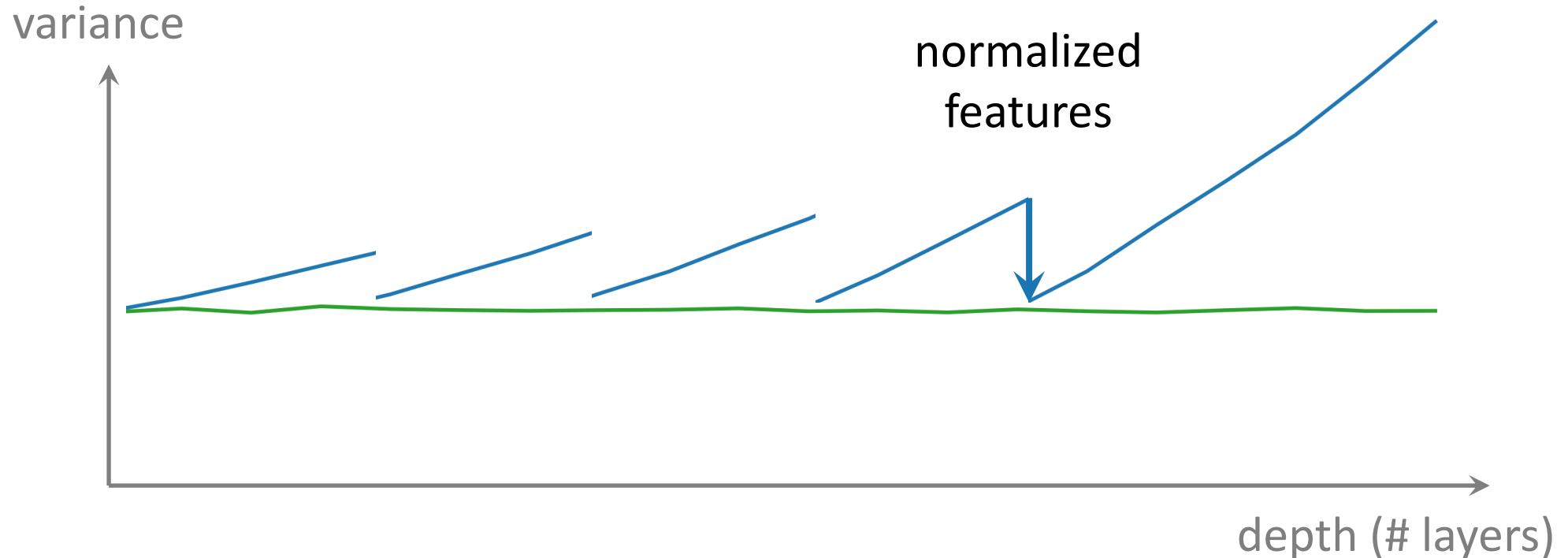
# Normalization Modules

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



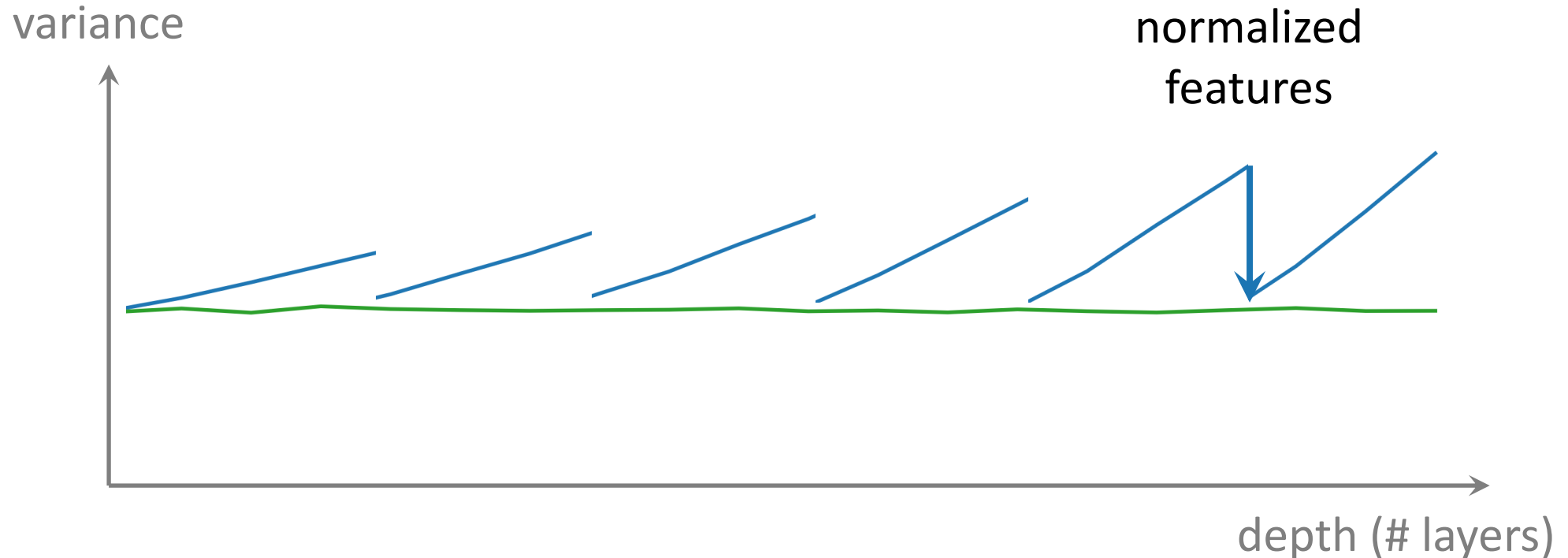
# Normalization Modules

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



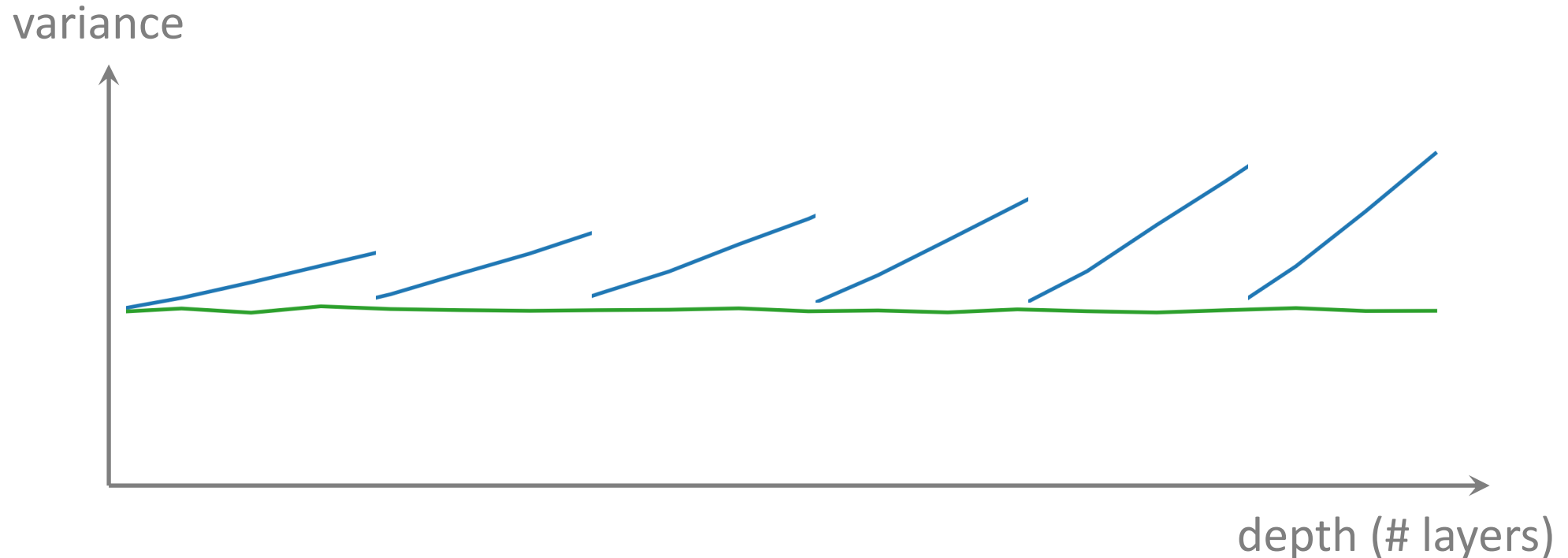
# Normalization Modules

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



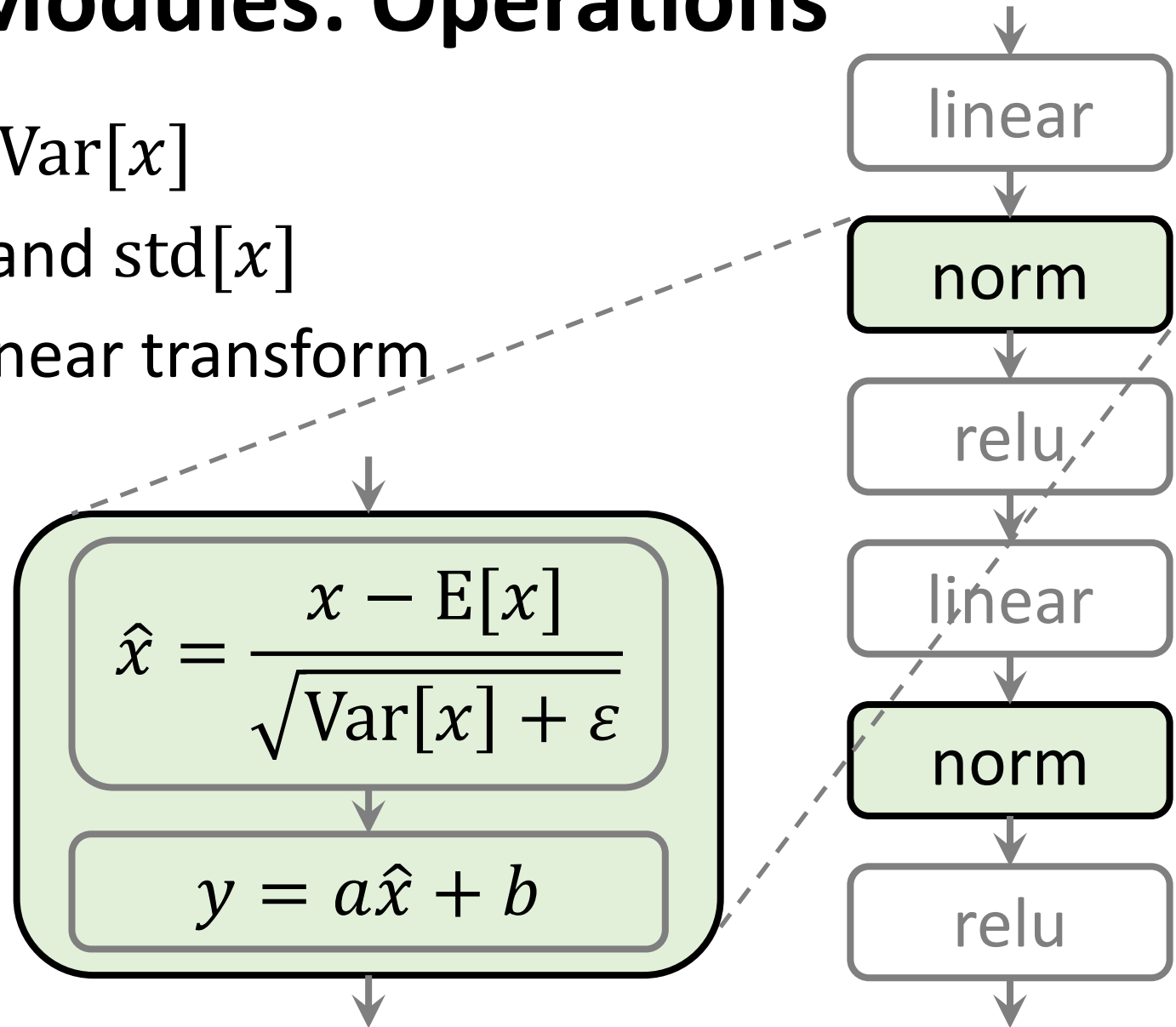
# Normalization Modules

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



# Normalization Modules: Operations

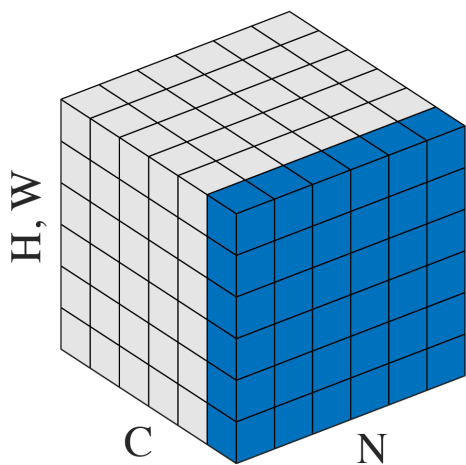
1. compute  $E[x]$  and  $\text{Var}[x]$
2. normalize by  $E[x]$  and  $\text{std}[x]$
3. compensate by a linear transform



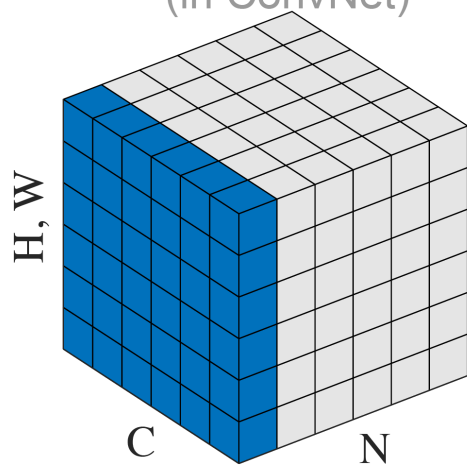
# Normalization Modules: Variants

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

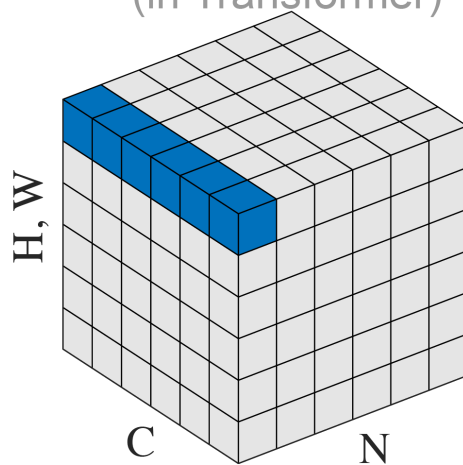
BatchNorm



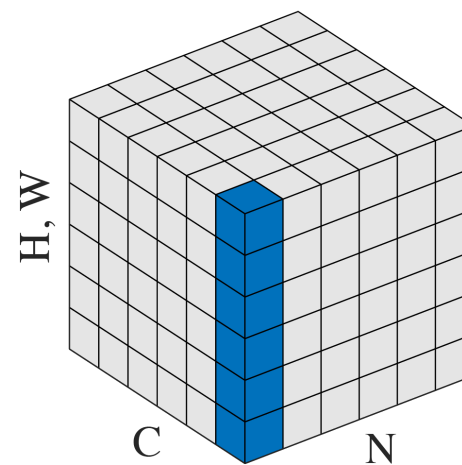
LayerNorm  
(in ConvNet)



LayerNorm  
(in Transformer)



InstanceNorm



GroupNorm

