

# Neural Networks II: Deep Learning

Mengye Ren

(Slides credit to David Rosenberg, He He, et al.)

NYU

Nov 26, 2024

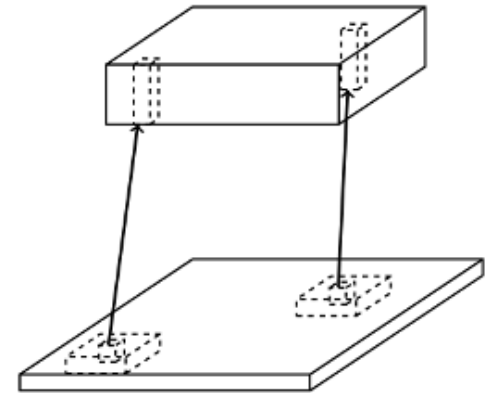


- Homework 4 Due: Dec 3.
- Last lecture: Dec 10 Project presentation.
- Presentation order: Your assigned Group ID.
- Each group has a max of 4 minutes (hard stop) + 1 min Q&A.
- OH this week: Wednesday 1-2pm.

Dec 9. send slides.  
in PDF.

# Local connection patterns

- The typical image input layer has 3 channels R G B for color or 1 channel for grayscale.
- The hidden layers may have  $C$  channels, at each spatial location  $(i, j)$ .

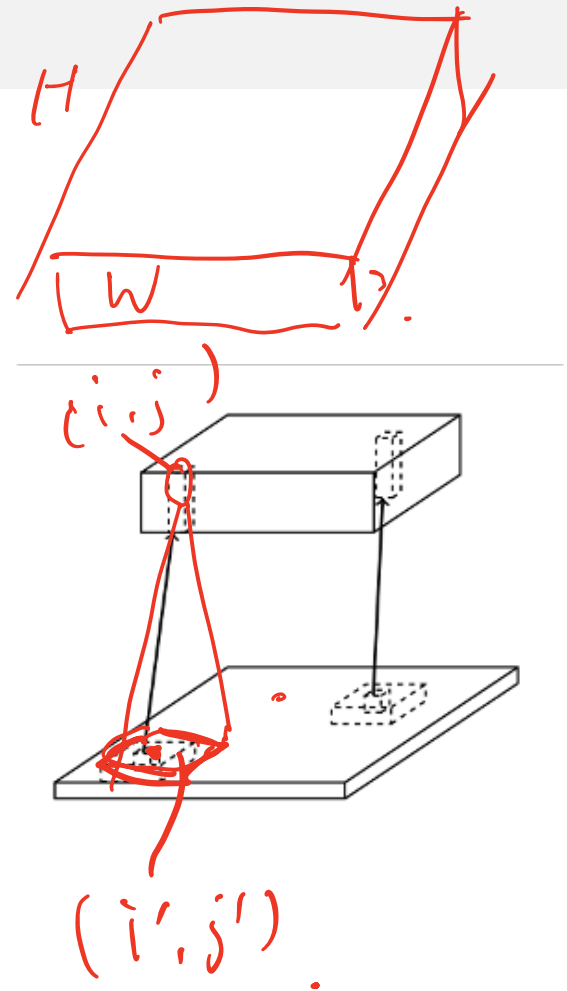




# Local connection patterns

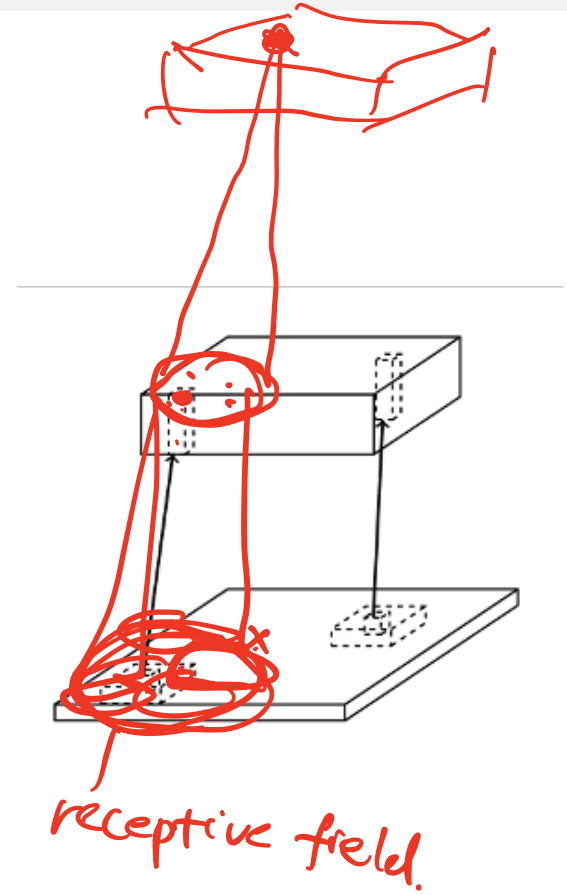
- The typical image input layer has 3 channels R G B for color or 1 channel for grayscale.
- The hidden layers may have  $C$  channels, at each spatial location  $(i, j)$ .
- Now each hidden neuron  $z_{i,j,c}$  receives inputs from  $x_{i\pm k, j\pm k, \cdot}$ .
- $k$  is the “kernel” size - do not confuse with the other kernel we learned.

$$z_{i,j,c} = \sum_{\substack{i' \in [i \pm k], \\ j' \in [j \pm k], \\ c' \in [1, C]}} x_{i',j',c'} \underbrace{w_{i,j,i'-i,j'-j,c',c}}_{\text{6 dim}} \quad \text{6 dim}$$



# Local connection patterns

- The typical image input layer has 3 channels R G B for color or 1 channel for grayscale.
- The hidden layers may have  $C$  channels, at each spatial location  $(i, j)$ .
- Now each hidden neuron  $z_{i,j,c}$  receives inputs from  $x_{i\pm k, j\pm k, \cdot}$ .
- $k$  is the “kernel” size - do not confuse with the other kernel we learned.
- $z_{i,j,c} = \sum_{i' \in [i\pm k], j' \in [j\pm k], c'} x_{i'j'c'} \underline{w_{i,j,i'-i,j'-j,c',c}}$
- The spatial awareness (receptive field) of the neighborhood grows bigger as we go deeper.



# Weight sharing

- Still a lot of weights: If we have 100 channels in the second layer, then  $200 \times 200 \times 3 \times 100 = 12M$

# Weight sharing

- Still a lot of weights: If we have 100 channels in the second layer, then  $200 \times 200 \times 3 \times 100 = 12M$
- Local information is the same regardless of the position of an element.

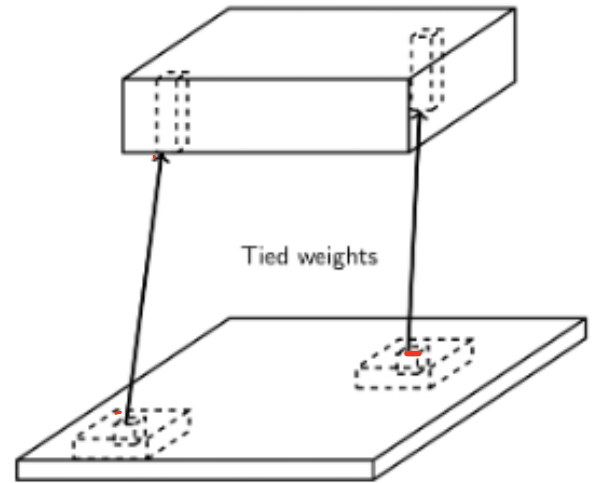
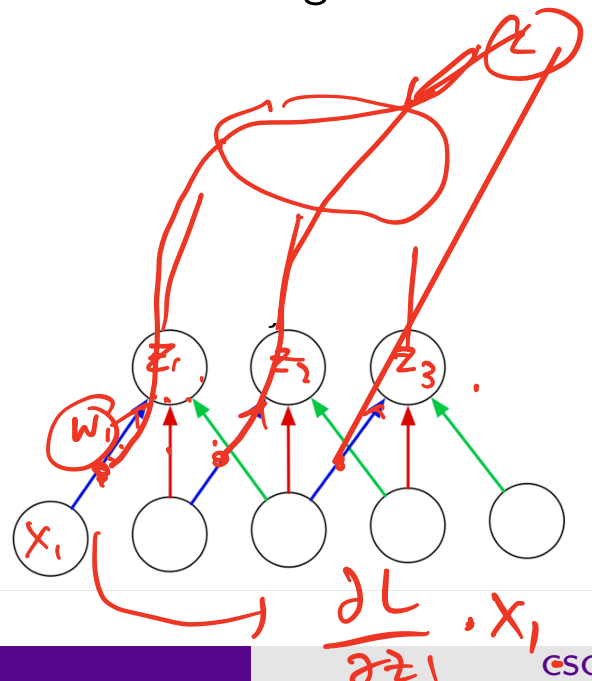
# Weight sharing

Fully Connected

- Still a lot of weights: If we have 100 channels in the second layer, then  $200 \times 200 \times 3 \times 100 = 12M$
- Local information is the same regardless of the position of an element.
- Solution: We can tie the weights at different locations.

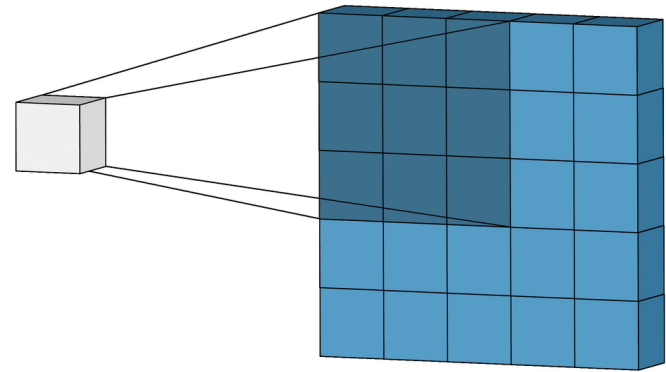
Local  
weight sharing  
CNN.

$$\frac{\partial L}{\partial w_i}$$



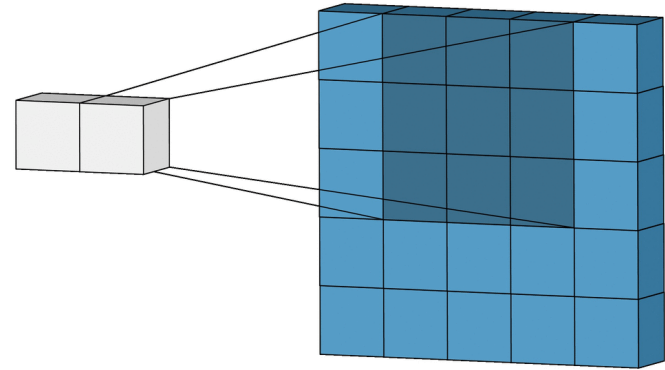
# 2D convolution

- Using the same weight connections for each activation spatial location works like the “filtering operation” or “convolution”
- The neighborhood window is the filter window.
- The weight connection is called “convolution filter”
- $Z_{i,j,c} = \sum_{i' \in [i \pm k], j' \in [j \pm k], c'} X_{i'j'c'} W_{i-i', j-j', c', c}$



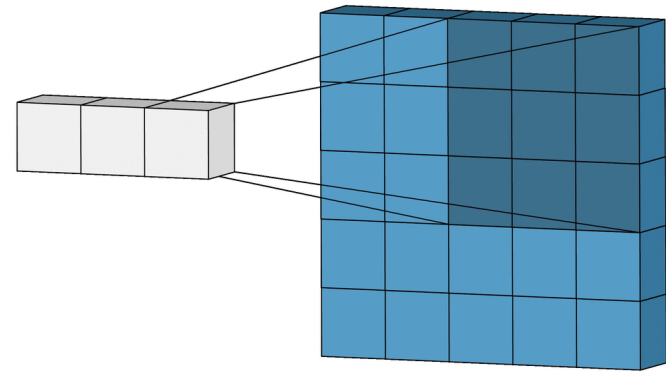
# 2D convolution

- Using the same weight connections for each activation spatial location works like the “filtering operation” or “convolution”
- The neighborhood window is the filter window.
- The weight connection is called “convolution filter”
- $Z_{i,j,c} = \sum_{i' \in [i \pm k], j' \in [j \pm k], c'} X_{i'j'c'} W_{i-i', j-j', c', c}$



# 2D convolution

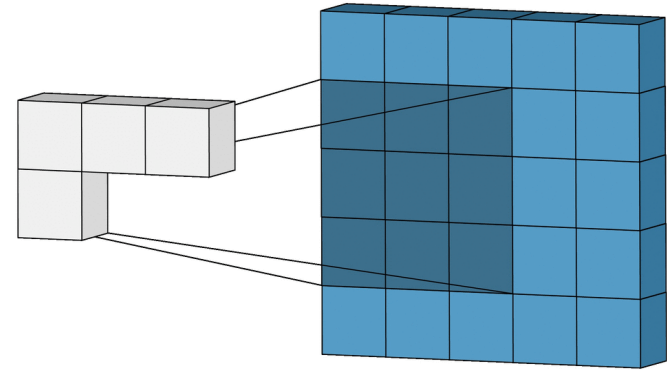
- Using the same weight connections for each activation spatial location works like the “filtering operation” or “convolution”
- The neighborhood window is the filter window.
- The weight connection is called “convolution filter”
- $Z_{i,j,c} = \sum_{i' \in [i \pm k], j' \in [j \pm k], c'} X_{i'j'c'} W_{i-i', j-j', c', c}$





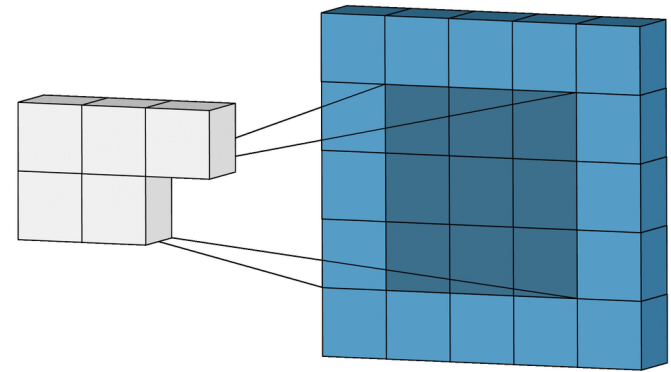
# 2D convolution

- Using the same weight connections for each activation spatial location works like the “filtering operation” or “convolution”
- The neighborhood window is the filter window.
- The weight connection is called “convolution filter”
- $Z_{i,j,c} = \sum_{i' \in [i \pm k], j' \in [j \pm k], c'} X_{i'j'c'} W_{i-i', j-j', c', c}$



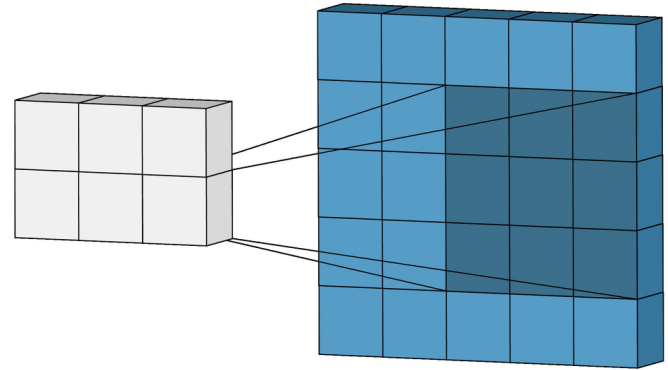
# 2D convolution

- Using the same weight connections for each activation spatial location works like the “filtering operation” or “convolution”
- The neighborhood window is the filter window.
- The weight connection is called “convolution filter”
- $Z_{i,j,c} = \sum_{i' \in [i \pm k], j' \in [j \pm k], c'} X_{i'j'c'} W_{i-i', j-j', c', c}$



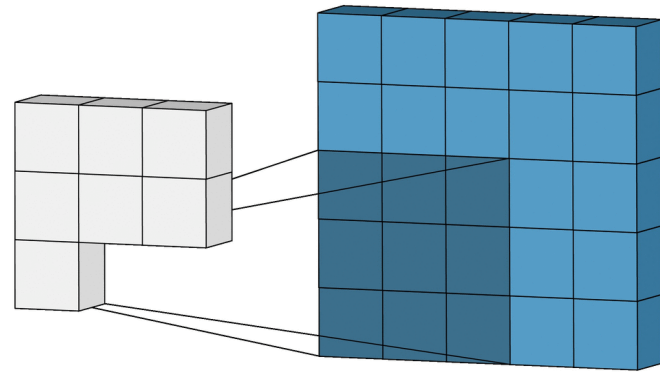
# 2D convolution

- Using the same weight connections for each activation spatial location works like the “filtering operation” or “convolution”
- The neighborhood window is the filter window.
- The weight connection is called “convolution filter”
- $Z_{i,j,c} = \sum_{i' \in [i \pm k], j' \in [j \pm k], c'} X_{i'j'c'} W_{i-i', j-j', c', c}$



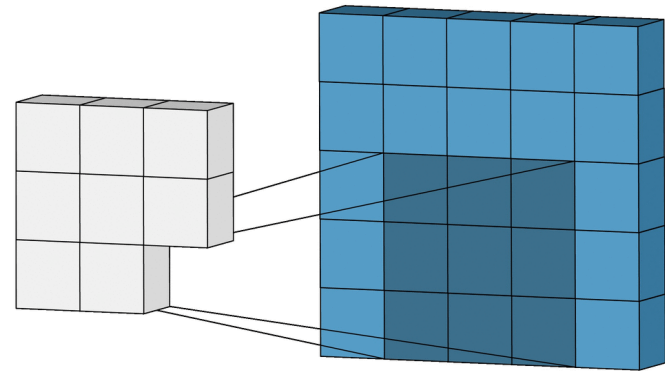
# 2D convolution

- Using the same weight connections for each activation spatial location works like the “filtering operation” or “convolution”
- The neighborhood window is the filter window.
- The weight connection is called “convolution filter”
- $Z_{i,j,c} = \sum_{i' \in [i \pm k], j' \in [j \pm k], c'} X_{i'j'c'} W_{i-i', j-j', c', c}$



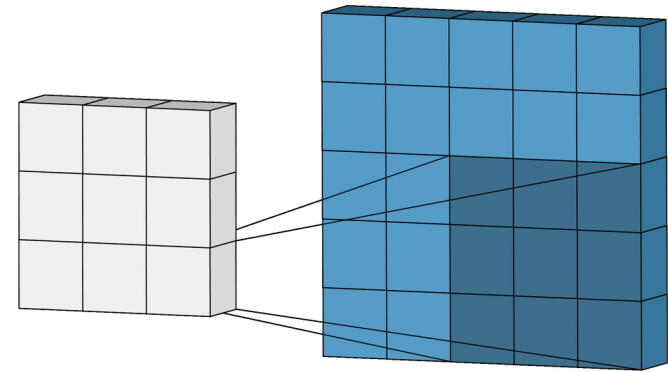
# 2D convolution

- Using the same weight connections for each activation spatial location works like the “filtering operation” or “convolution”
- The neighborhood window is the filter window.
- The weight connection is called “convolution filter”
- $Z_{i,j,c} = \sum_{i' \in [i \pm k], j' \in [j \pm k], c'} X_{i'j'c'} W_{i-i', j-j', c', c}$



# 2D convolution

- Using the same weight connections for each activation spatial location works like the “filtering operation” or “convolution”
- The neighborhood window is the filter window.
- The weight connection is called “convolution filter”



- $Z_{i,j,c} =$

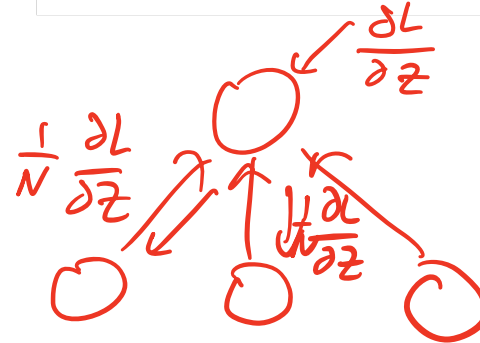
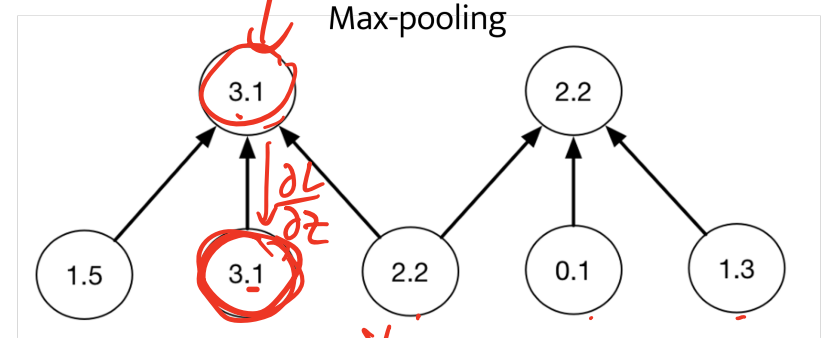
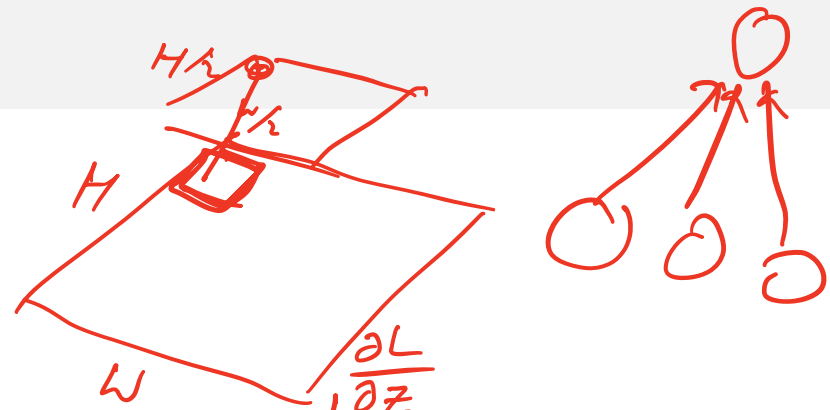
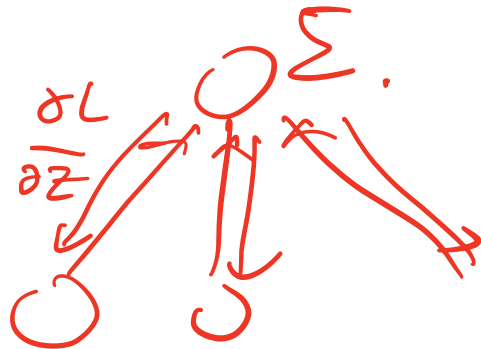
$$\sum_{i' \in [i \pm k], j' \in [j \pm k], c'} X_{i'j'c'} \underline{W_{i-i', j-j', c', c}}$$

neighborhood  
index

input channel  
output channel.

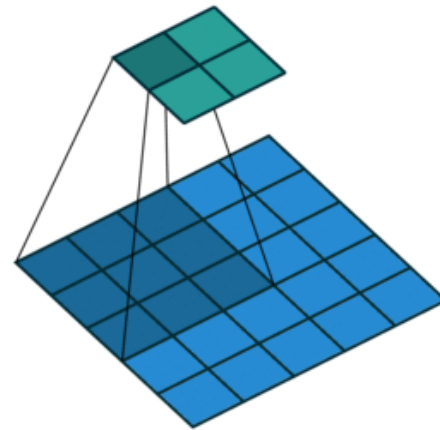
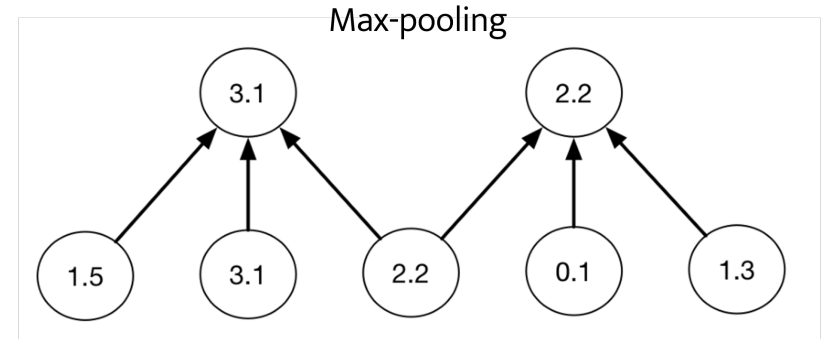
# Pooling

- Need to summarize global information more efficiently.
- Pooling reduces image / activation dimensions.
- Max-pooling or average-pooling



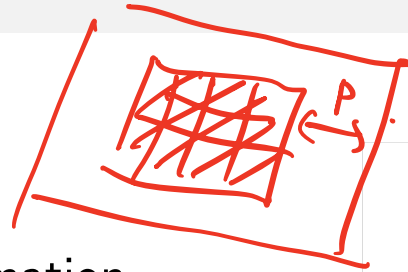
# Pooling

- Need to summarize global information more efficiently.
- Pooling reduces image / activation dimensions.
- Max-pooling or average-pooling
- You can also perform a “strided” convolution by jumping multiple steps.

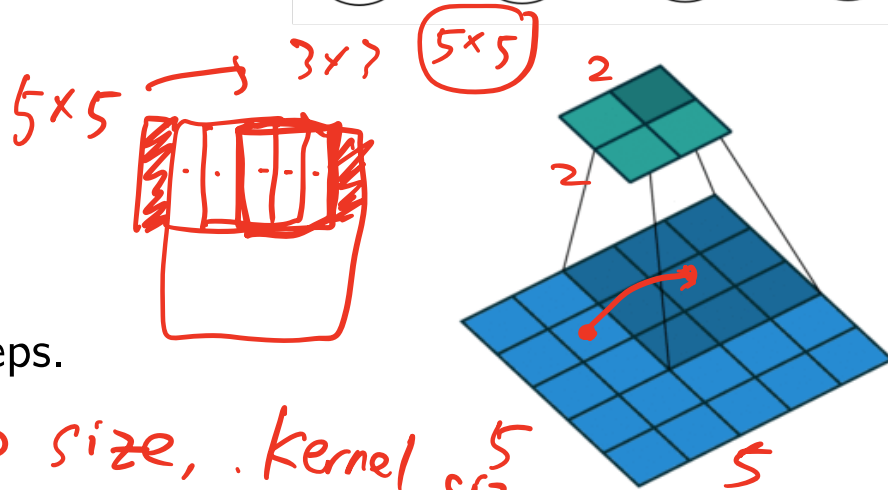
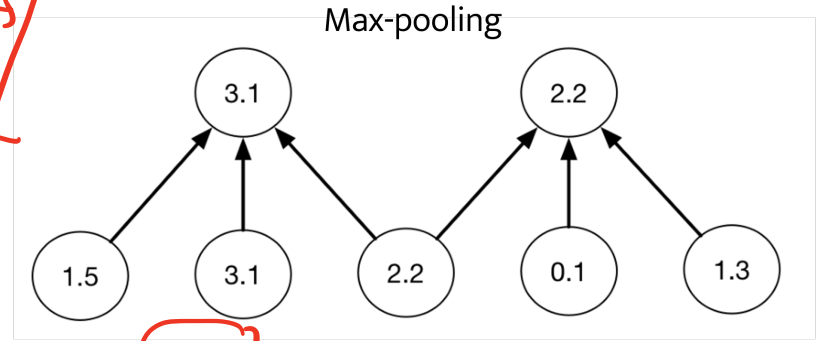




# Pooling



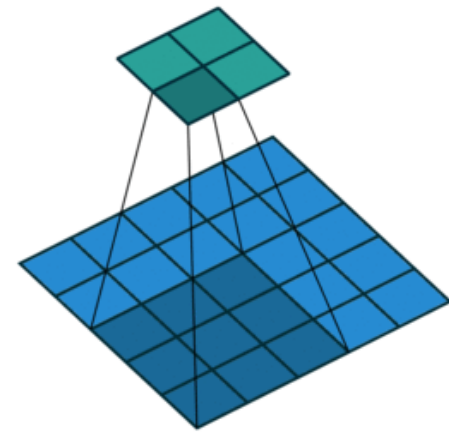
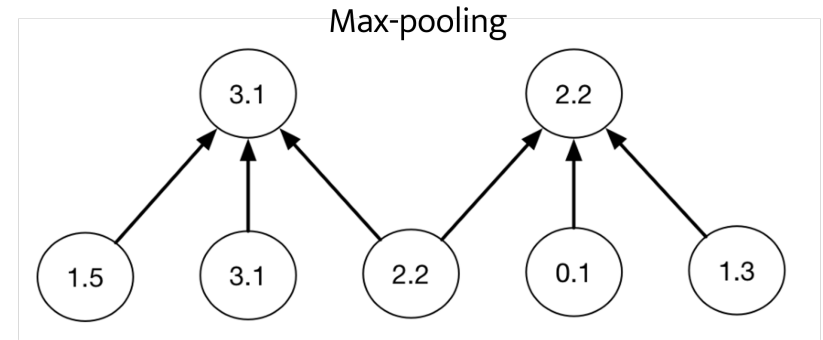
- Need to summarize global information more efficiently.
- Pooling reduces image / activation dimensions.
- Max-pooling or average-pooling
- You can also perform a “strided” convolution by jumping multiple steps.



$Out = \text{function of } ( \text{inp size}, \text{kernel size}, \text{stride}, \text{padding} )$

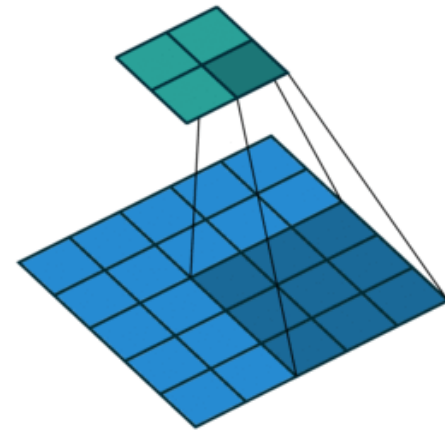
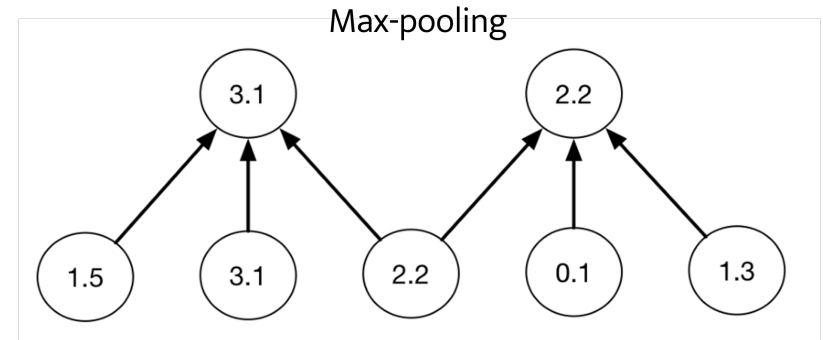
# Pooling

- Need to summarize global information more efficiently.
- Pooling reduces image / activation dimensions.
- Max-pooling or average-pooling
- You can also perform a “strided” convolution by jumping multiple steps.

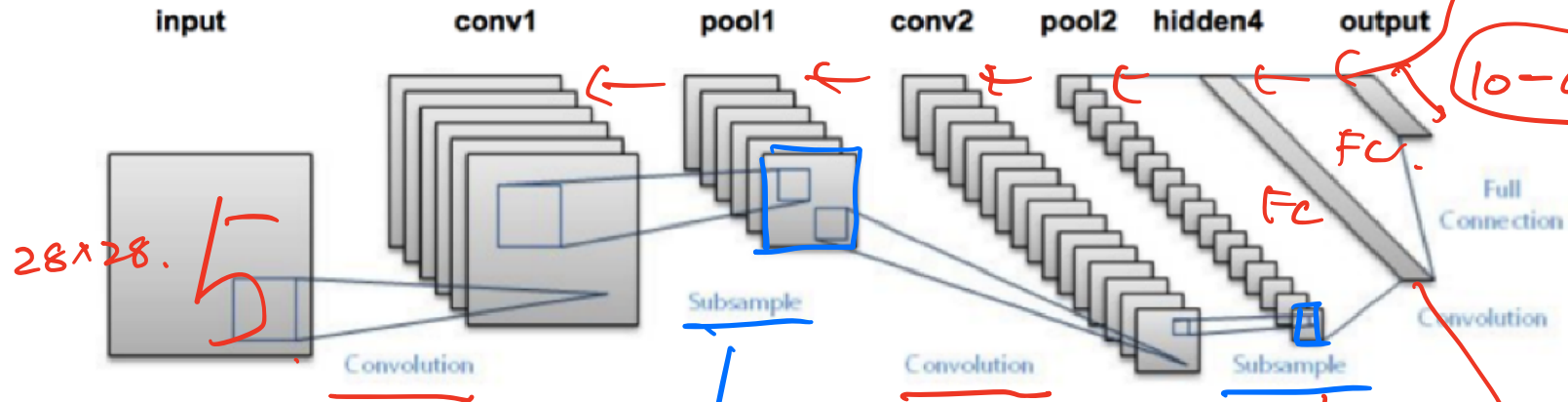


# Pooling

- Need to summarize global information more efficiently.
- Pooling reduces image / activation dimensions.
- Max-pooling or average-pooling
- You can also perform a “strided” convolution by jumping multiple steps.



# Assembling together: LeNet



Cross entropy.

0-9.

10-dim.

pooling.

- Used by USPS to read post code in the 90s.

7x7  
4x4

4x4 x 64

vector  
global  
representation.  
64. Avg  
16x64 convst.

# Historical development

- LeNet has worked and being put to practice in the 1990s.

# Historical development

- LeNet has worked and being put to practice in the 1990s.
- Neural networks for images start to dominate in the last 10 years (starting 2012) for understanding general high resolution natural images.

# Historical development

- LeNet has worked and being put to practice in the 1990s.
- Neural networks for images start to dominate in the last 10 years (starting 2012) for understanding general high resolution natural images.
- During the years:
  - Neural networks were difficult to work
  - People focused on feature engineering
  - Then apply SVM or random forest (e.g. AdaBoost face detector)
  - What has changed?

# Gradient learning conditioning



# Optimization challenges

- Larger images require deeper networks (more stages of processing at different resolutions)
- Optimizing deeper layers of networks is not trivial.
- Loss often stalls or blows up.

# Optimization challenges

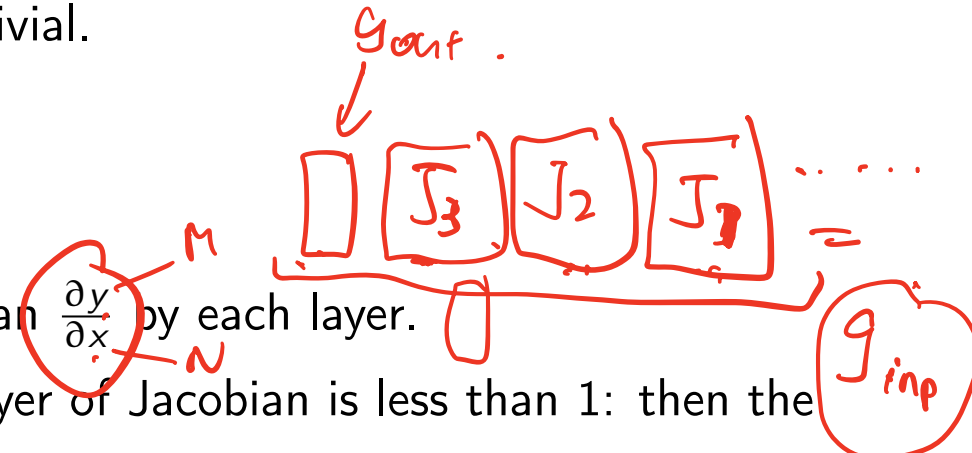
- Larger images require deeper networks (more stages of processing at different resolutions)
- Optimizing deeper layers of networks is not trivial.
- Loss often stalls or blows up.
- Why?

# Optimization challenges

- Larger images require deeper networks (more stages of processing at different resolutions)
- Optimizing deeper layers of networks is not trivial.
- Loss often stalls or blows up.

- Why?

- Backpropagation: multiplying the Jacobian  $\frac{\partial y}{\partial x}$  by each layer.
- If the maximum singular value of each layer of Jacobian is less than 1: then the gradient will converge to 0 with more layers.
- If the greater than 1: then the gradient will explode with more layers.
- The bottom (input) layer may get 0 or infinite gradients.



# Weight initialization

- Even with a few layers ( $>3$ ), optimization is still hard.

# Weight initialization

- Even with a few layers ( $>3$ ), optimization is still hard.
- If weight initialization is bad (too small or too big), then optimization is hard to kick off.

# Weight initialization

- Even with a few layers ( $>3$ ), optimization is still hard.
- If weight initialization is bad (too small or too big), then optimization is hard to kick off.
- Consider the distribution of whole dataset in the activation space.
  - Intuition: upon initialization, the variance of the activations should stay the same across every layer



# Kaiming Initialization

- Suppose each neuron and weight connection are sampling from a random distribution.

---

<sup>1</sup>He et al. Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet. ICCV, 2015.

# Kaiming Initialization

- Suppose each neuron and weight connection are sampling from a random distribution.

- At  $l$ -th layer,  $\text{Var}[z_l] = n_l \text{Var}[w_l x_l]$  ( $n_l = \text{num. input neurons to } l\text{-th layer}$ )

---

<sup>1</sup>He et al. Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet. ICCV, 2015.



# Kaiming Initialization



$$\text{Var} \left[ \underbrace{\text{Var}[x^+]}_{\frac{1}{2}} + \underbrace{\text{Var}[x^-]}_{0} \right]$$

- Suppose each neuron and weight connection are sampling from a random distribution.
- At  $l$ -th layer,  $\text{Var}[z_l] = n_l \text{Var}[w_l x_l]$  ( $n_l = \text{num. input neurons to } l\text{-th layer}$ )
- If we suppose that ReLU is used as the activation, and  $w_l$  is symmetric and zero-mean,  $x_{l+1} = \frac{1}{2} \text{Var}[z_l]$ .

<sup>1</sup>He et al. Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet. ICCV, 2015.

# Kaiming Initialization

- Suppose each neuron and weight connection are sampling from a random distribution.
- At  $l$ -th layer,  $\text{Var}[z_l] = n_l \text{Var}[w_l x_l]$  ( $n_l = \text{num. input neurons to } l\text{-th layer}$ )
- If we suppose that ReLU is used as the activation, and  $w_l$  is symmetric and zero-mean,  $x_{l+1} = \frac{1}{2} \text{Var}[z_l]$ .
- Putting altogether,  $\text{Var}[x_{l+1}] = \frac{1}{2} n_l \text{Var}[w_l] \text{Var}[x_l]$ .

---

<sup>1</sup>He et al. Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet. ICCV, 2015.

# Kaiming Initialization

- Suppose each neuron and weight connection are sampling from a random distribution.
- At  $l$ -th layer,  $\text{Var}[z_l] = n_l \text{Var}[w_l x_l]$  ( $n_l = \text{num. input neurons to } l\text{-th layer}$ )
- If we suppose that ReLU is used as the activation, and  $w_l$  is symmetric and zero-mean,  $x_{l+1} = \frac{1}{2} \text{Var}[z_l]$ .
- Putting altogether,  $x_{l+1} = \frac{1}{2} n_l \text{Var}[w_l] \text{Var}[x_l]$ .
- To make the variance constant, we need  $\frac{1}{2} n_l \text{Var}[w_l] = 1$ ,  $\text{Std}[w_l] = \sqrt{2/n_l}$ .

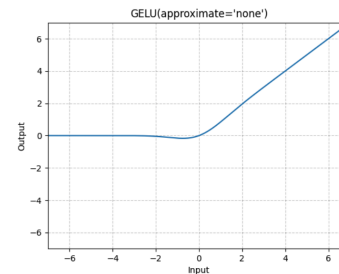
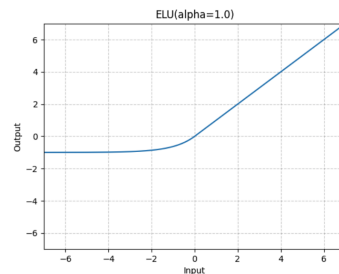
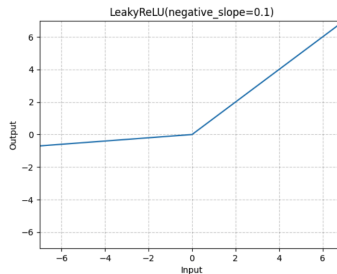
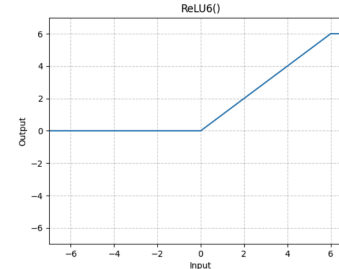
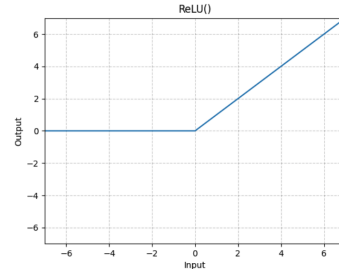
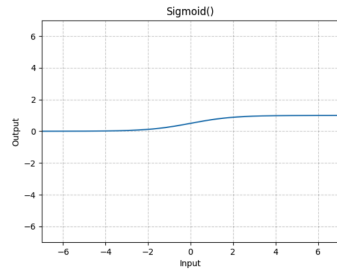
more inp neurons  
small std. of weights.

$\sqrt{\frac{c}{n_{in}}}$   
w/o relu.

<sup>1</sup>He et al. Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet. ICCV, 2015.

# Activation functions

- ReLU was proposed in 2009-2010<sup>23</sup>, and was successfully used in AlexNet in 2012<sup>4</sup>.
- Address the vanishing gradient issue in activations, comparing to sigmoid or tanh.



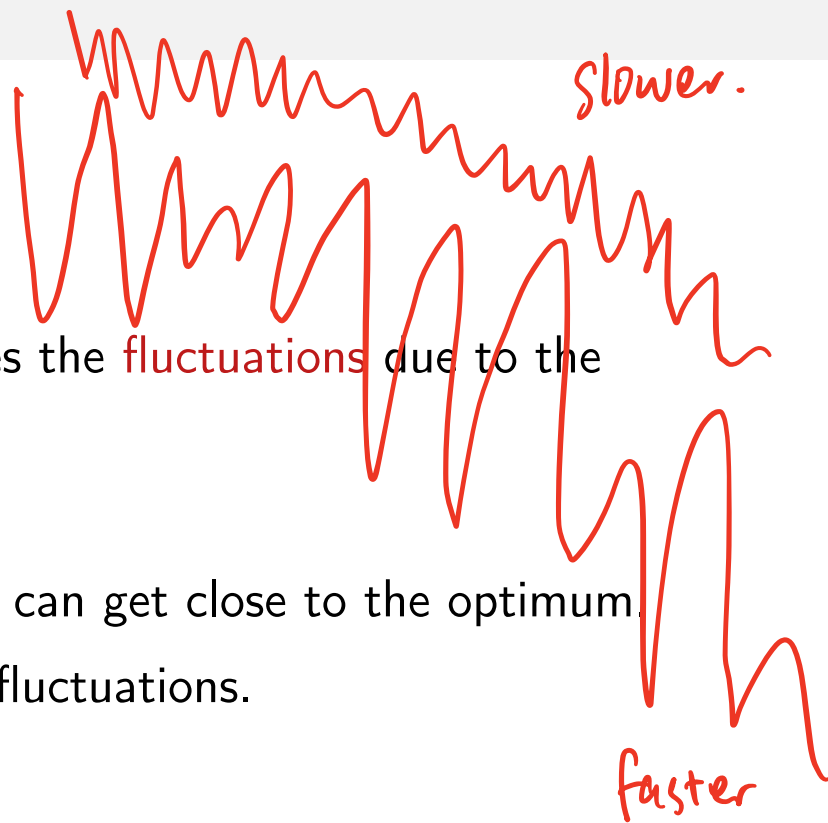
$$f(z) = \frac{1}{1 + e^{-z}}$$

<sup>2</sup> Jarrett et al. What is the Best Multi-Stage Architecture for Object Recognition? ICCV, 2009.  
<sup>3</sup> Nair & Hinton/ Rectified Linear Units Improve Restricted Boltzmann Machines. ICML, 2010.  
<sup>4</sup> Krizhevsky et al. ImageNet Classification with Deep Convolutional Neural Networks. NIPS, 2012.

# SGD Learning Rate

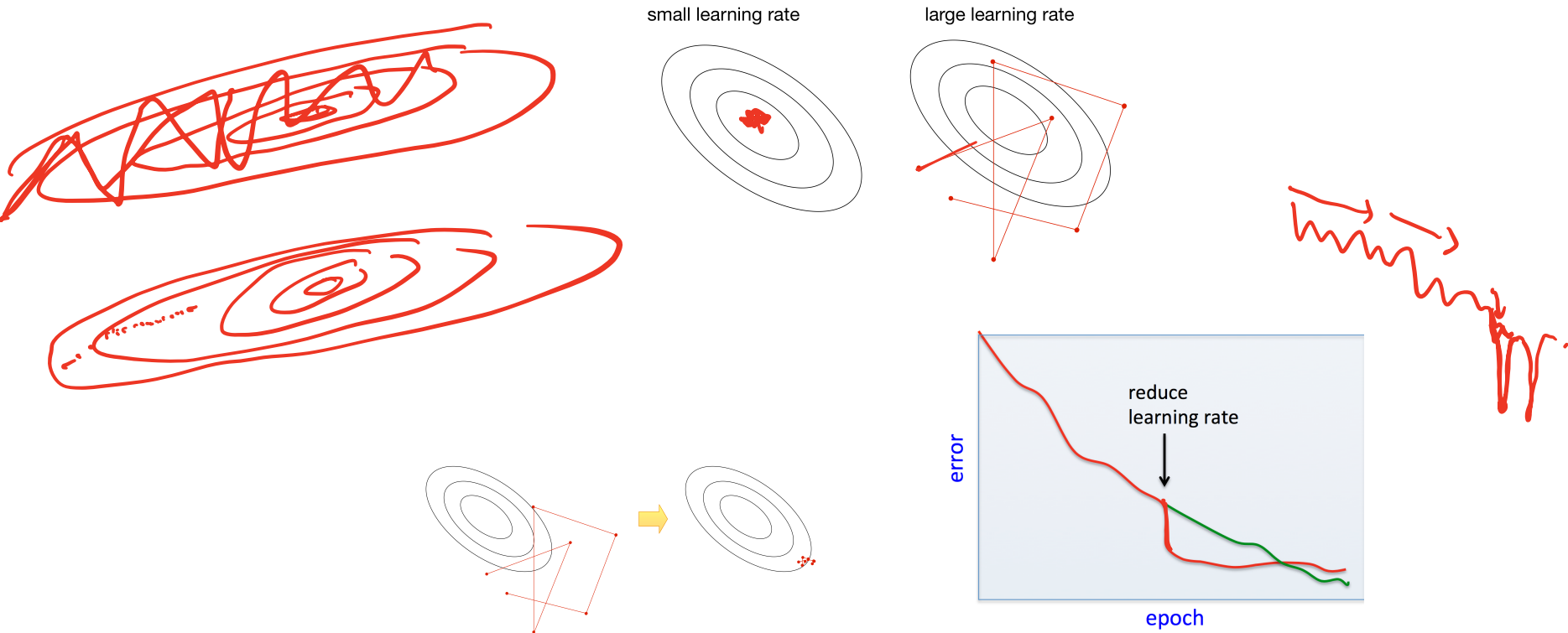
mini-batch.

- In stochastic training, the learning rate also influences the **fluctuations** due to the stochasticity of the gradients.
- Typical strategy:
  - Use a large learning rate early in training so you can get close to the optimum.
  - Gradually decay the learning rate to reduce the fluctuations.



# Learning Rate Decay

- We also need to be aware about the impact of learning rate due to the stochasticity.



# RMSprop and Adam

- Recall: SGD takes large steps in directions of high curvature and small steps in directions of low curvature.
- RMSprop is a variant of SGD which rescales each coordinate of the gradient to have norm 1 on average. It does this by keeping an exponential moving average  $s_j$  of the squared gradients.

# RMSprop and Adam

- Recall: SGD takes large steps in directions of high curvature and small steps in directions of low curvature.
- **RMSprop** is a variant of SGD which rescales each coordinate of the gradient to have norm 1 on average. It does this by keeping an exponential moving average  $s_j$  of the squared gradients.
- The following update is applied to each coordinate  $j$  independently:

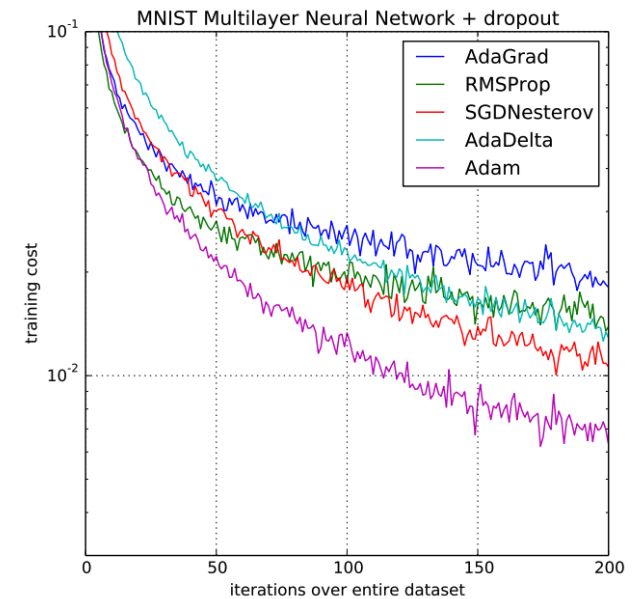
$$s_j \leftarrow (1 - \gamma)s_j + \gamma \left[ \frac{\partial L}{\partial \theta_j} \right]^2$$
$$\theta_j \leftarrow \theta_j - \frac{\alpha}{\sqrt{s_j + \epsilon}} \frac{\partial L}{\partial \theta_j}$$

— exponential moving avg. 0.2002.  
if norm is small



# Adam optimizer

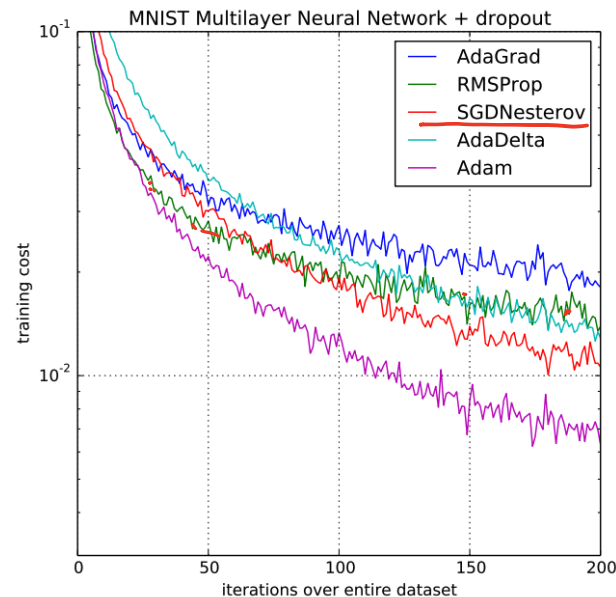
- **Adam** = RMSprop + momentum = Adaptive Momentum estimation
- Smoother estimate of the average gradient and gradient norm.



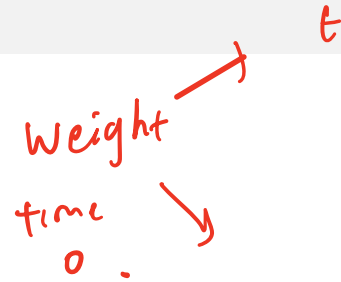
# Adam optimizer

- Adam = RMSprop + momentum = Adaptive Momentum estimation
- Smoother estimate of the average gradient and gradient norm.
- $m_t$ : exponential moving average of gradient.
- $v_t$ : exponential moving average of gradient squared.
- $\hat{m}_t, \hat{v}_t$ : Bias correction.
- $\theta_t \leftarrow \theta_{t-1} - \alpha \hat{m}_t / (\sqrt{\hat{v}_t} + \epsilon)$
- The “default” optimizer for modern networks.

→ Adam W.



# Normalization



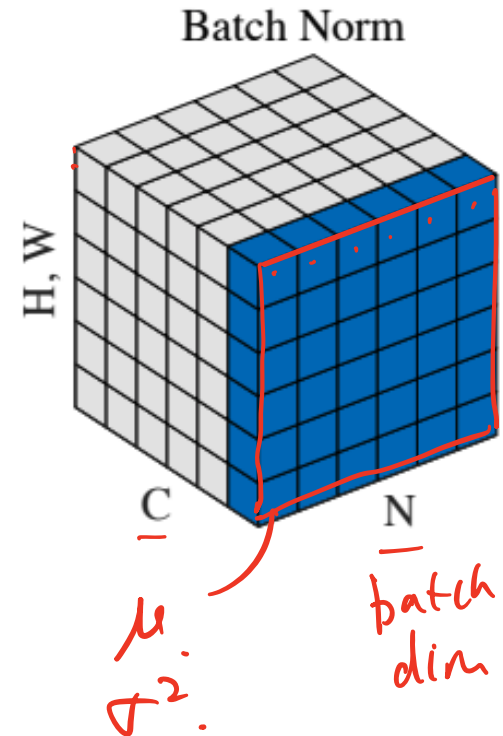
- Weight initialization is tricky, and there is no guarantee that the distribution of activations will stay the same over the learning process.
- What if the weights keep grow bigger and activation may explode?

# Normalization

- Weight initialization is tricky, and there is no guarantee that the distribution of activations will stay the same over the learning process.
- What if the weights keep grow bigger and activation may explode?
- We can “normalize” the activations.
- The idea is to control the activation within a normal range: zero-mean, uni-variance.

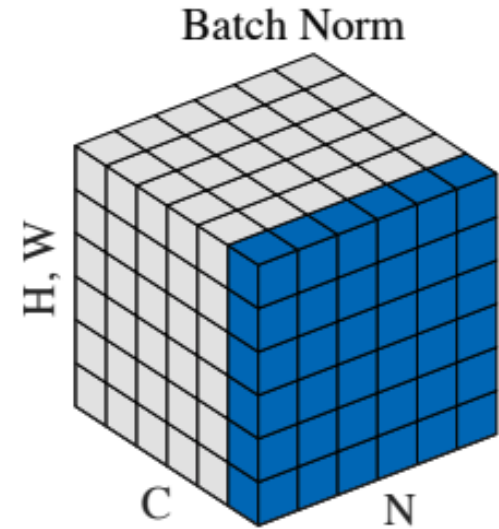
# Batch Normalization (BN)

- In CNNs, neurons across different spatial locations are also samples of the same feature channel.



# Batch Normalization (BN)

- In CNNs, neurons across different spatial locations are also samples of the same feature channel.
- Batch norm: Normalize across  $N$   $H$   $W$  dimensions, leaving  $C$  channels.



# Batch Normalization (BN)

- In CNNs, neurons across different spatial locations are also samples of the same feature channel.
- Batch norm: Normalize across N H W dimensions, leaving C channels.

•  $\tilde{x} = \gamma \frac{x - \mu}{\sigma} + \beta$  — channel C.

- $\gamma, \beta$ : learnable parameters.  $\mu, \sigma$ : statistics from the training batch.

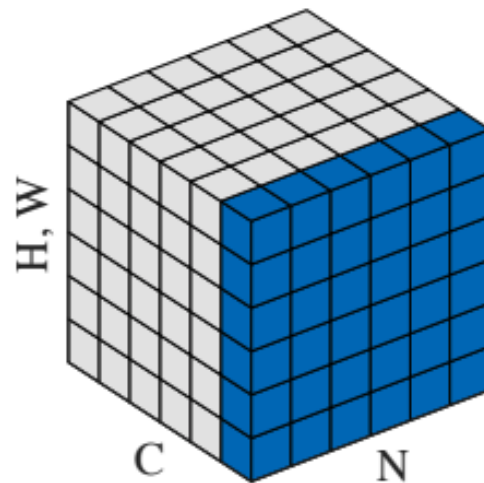
Batch-normalized activation.

$$\tilde{x} = \frac{x - \mu}{\sigma}$$

Wrong.

$$\frac{1}{\sigma} \cdot \frac{\partial L}{\partial x}$$

Batch Norm



# Batch Normalization (BN)

- In CNNs, neurons across different spatial locations are also samples of the same feature channel.

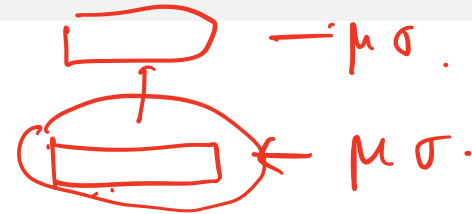
- Batch norm: Normalize across N H W dimensions, leaving C channels.

- $\tilde{x} = \gamma \frac{x - \mu}{\sigma} + \beta$   
*+  $\epsilon \approx 1e-7$*

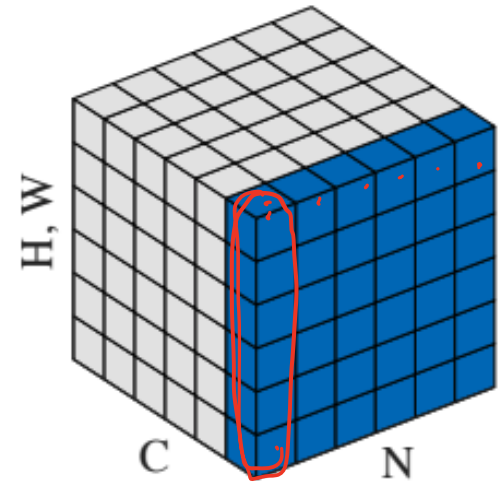
- $\gamma, \beta$ : learnable parameters.  $\mu, \sigma$ : statistics from the training batch.

- Test time: using the mean and variance from the

*one image.*



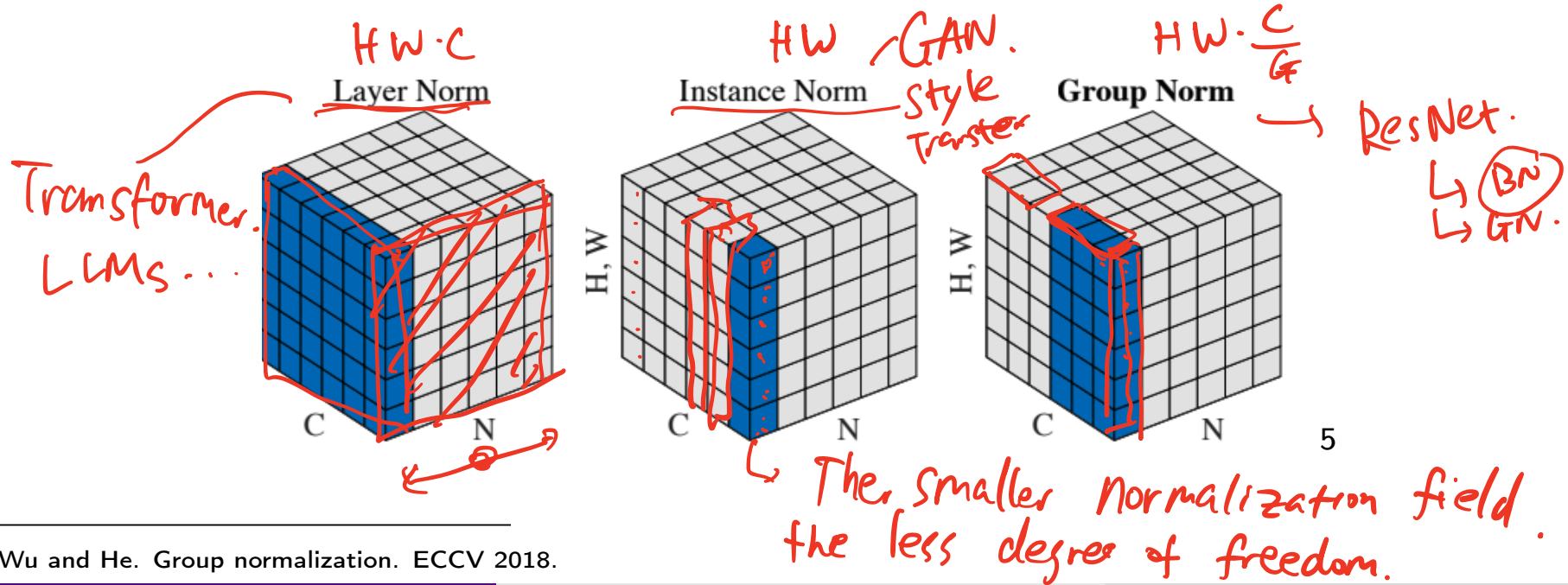
Batch Norm





# BN Alternatives

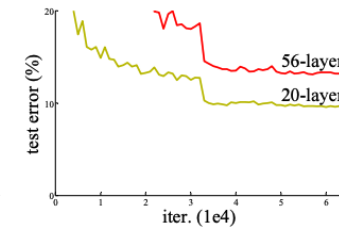
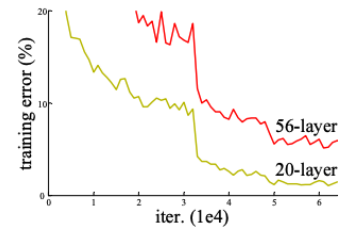
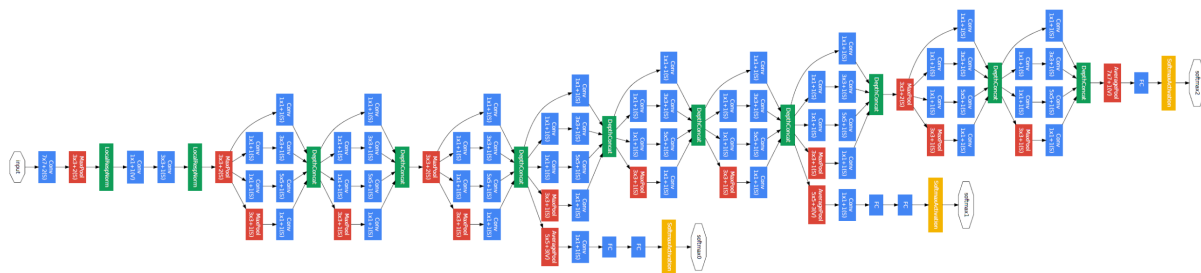
- Need a considerable batch size to estimate mean and variance correctly.
- Training is different from testing.
- Alternatives consider the C channel dimension instead of N batch dimension.



<sup>5</sup>Wu and He. Group normalization. ECCV 2018.

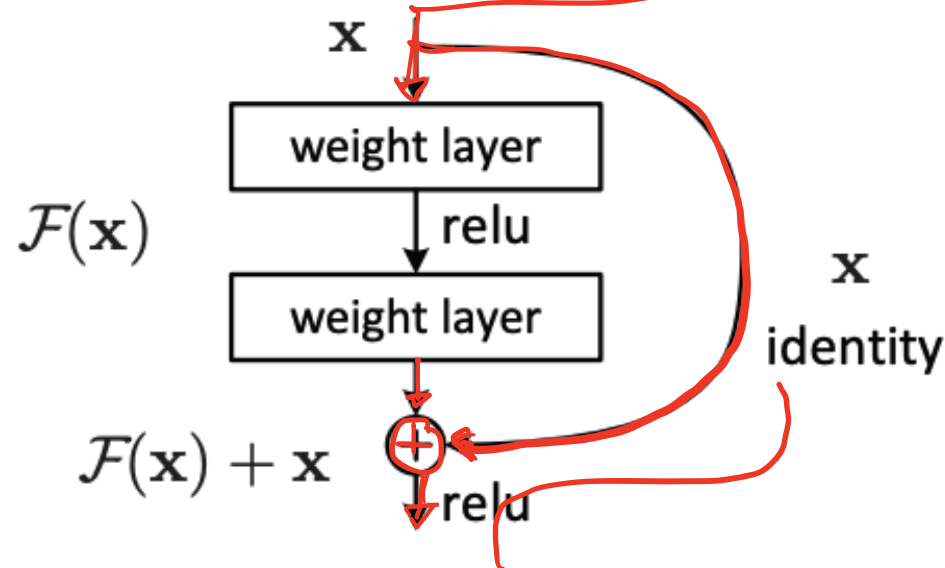
# Going Deeper

- The progress of normalization allowed us to train even deeper networks.
- The networks are no longer too sensitive with initialization.
- But the best networks were still around 20 layers and deeper results in worse performance.



# Residual Networks (ResNet)

- Recall in gradient boosting, we are iteratively adding a function to the model to expand the capacity.
- Residual connection: Skip connection to prevent gradient vanishing.<sup>6</sup>



<sup>6</sup>He et al. Deep Residual Learning for Image Recognition. CVPR 2016.

# ResNet Success

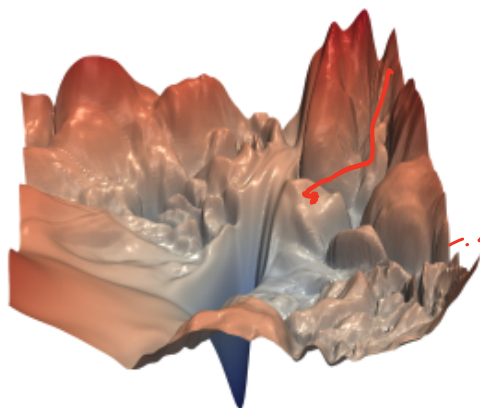
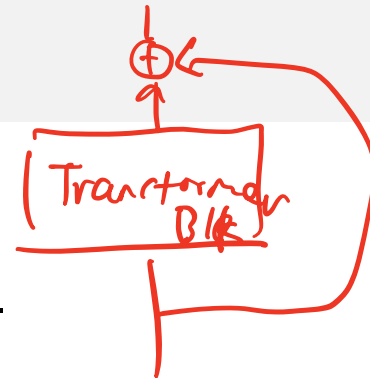
- Now able to train over 100 layers.
- One of the most important network design choices in the past decade.
- Prevalent in almost all network architectures, including Transformers.

---

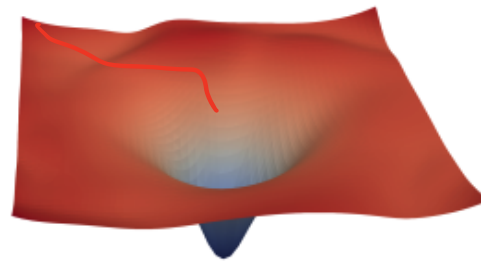
<sup>7</sup> Li et al. Visualizing the Loss Landscape of Neural Nets. NIPS 2018.

# ResNet Success

- Now able to train over 100 layers.
- One of the most important network design choices in the past decade.
- Prevalent in almost all network architectures, including Transformers.
- Loss landscape view: Skip connections makes loss smoother  $\rightarrow$  easier to optimize <sup>7</sup>.



(a) without skip connections

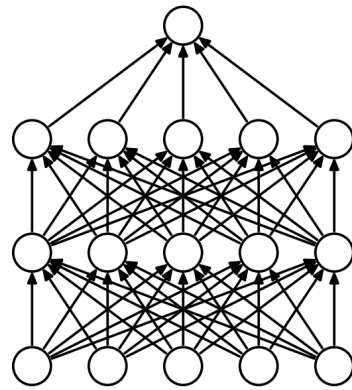


(b) with skip connections

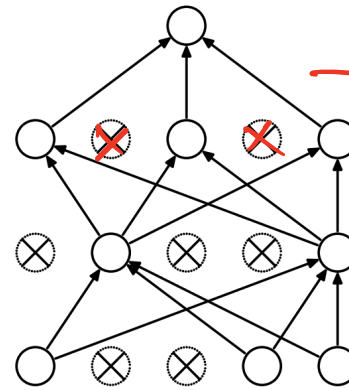
<sup>7</sup> Li et al. Visualizing the Loss Landscape of Neural Nets. NIPS 2018.

# Dropout<sup>8</sup>

- Want to reduce overfitting in neural networks.
- Stochastically turning off neurons in propagation.



(a) Standard Neural Net



(b) After applying dropout.

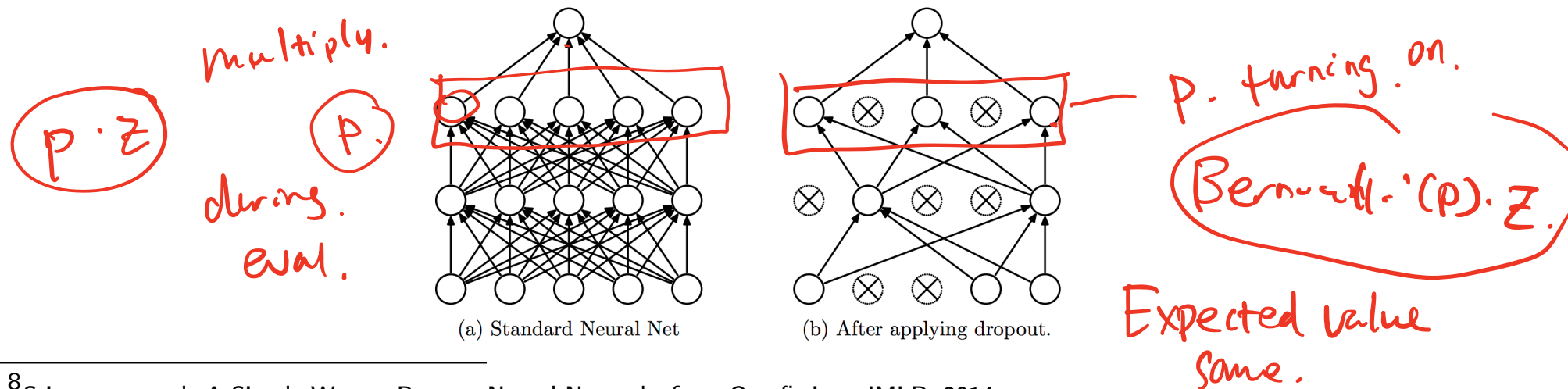
Transformer.

*training.*

<sup>8</sup>Srivastava et al. A Simple Way to Prevent Neural Networks from Overfitting. JMLR, 2014.

# Dropout<sup>8</sup>

- Want to reduce overfitting in neural networks.
- Stochastically turning off neurons in propagation.
- Training to preserve redundancy.
- Test time: multiplying activations with probability. Model ensembling effect.

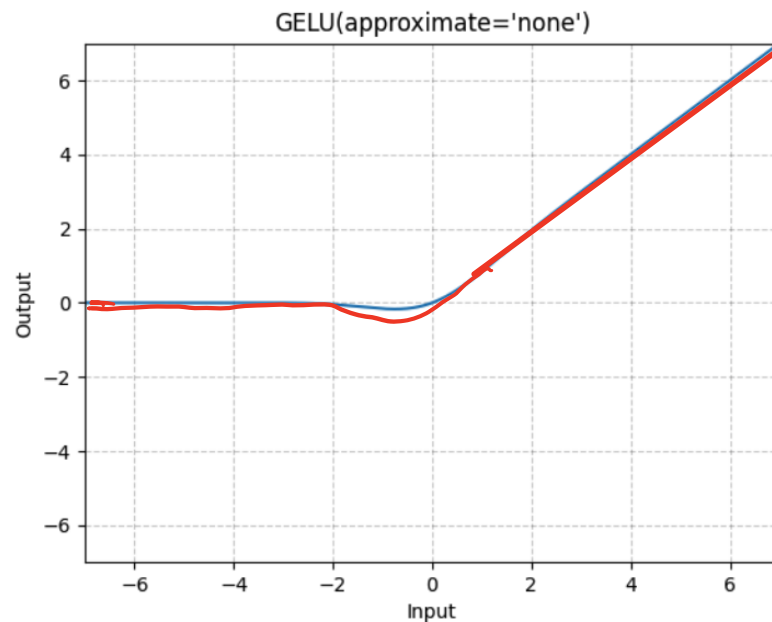


<sup>8</sup>Srivastava et al. A Simple Way to Prevent Neural Networks from Overfitting. JMLR, 2014.

# GELU<sup>9</sup>

- Gaussian Error Linear Unit - A smoother activation function.
- Motivated by Dropout.

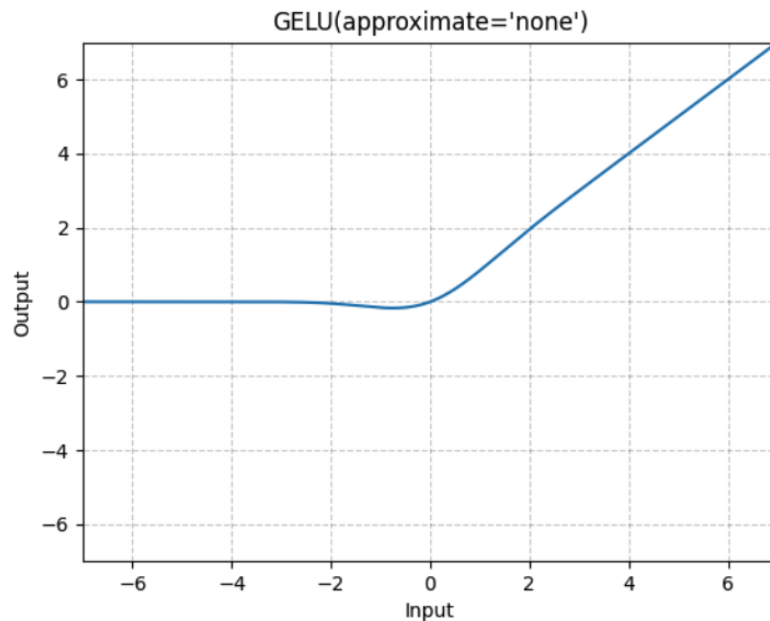
50°  
ReLU.  
dead unit.



<sup>9</sup>Hendrycks & Gimpel. Gaussian Error Linear Unit (GELU). CoRR abs/1606.08415, 2016.



- Gaussian Error Linear Unit - A smoother activation function.
- Motivated by Dropout.
- $f(x) = \mathbb{E}[x \cdot m]$ .



<sup>9</sup>Hendrycks & Gimpel. Gaussian Error Linear Unit (GELU). CoRR abs/1606.08415, 2016.

Transformer.

- Gaussian Error Linear Unit - A smoother activation function.

- Motivated by Dropout.

- $f(x) = \mathbb{E}[x \cdot m]$ .

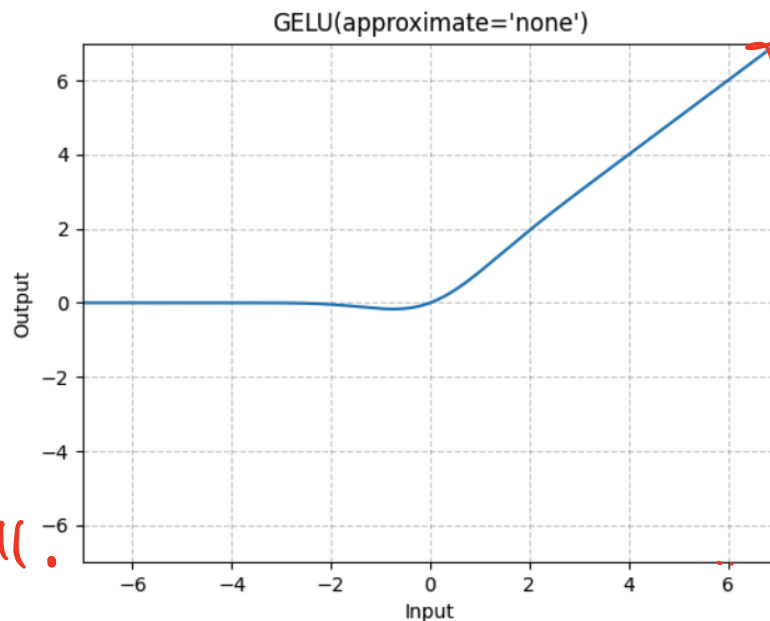
- $m \sim \text{Bernoulli}(\Phi(x))$ .

- $\Phi(x) = P(X \leq x)$ .

- $X \sim \mathcal{N}(0, 1)$ .

$f(x) \approx x$  if  $x$  is big.

$f(x) \approx 0$  if  $x$  is small.



<sup>9</sup>Hendrycks & Gimpel. Gaussian Error Linear Unit (GELU). CoRR abs/1606.08415, 2016.

# Data augmentation

- Leverage the invariances of images
- Create more data points for free

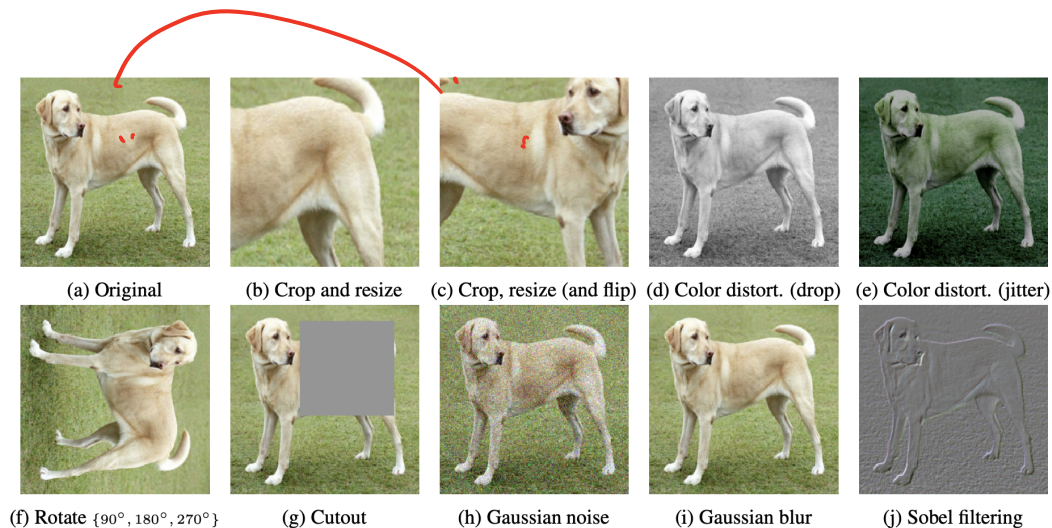


Image credit<sup>10</sup>

<sup>10</sup>Chen et al. A Simple Framework for Contrastive Learning of Visual Representations. ICML 2020.

# Data augmentation

- Leverage the invariances of images
- Create more data points for free
  - Random cropping

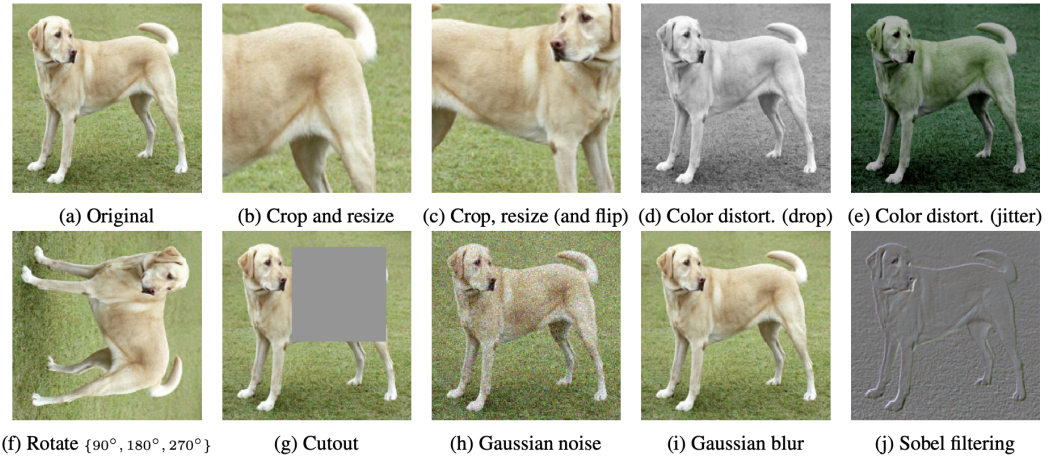


Image credit<sup>10</sup>

---

<sup>10</sup>Chen et al. A Simple Framework for Contrastive Learning of Visual Representations. ICML 2020.

# Data augmentation

- Leverage the invariances of images
- Create more data points for free
  - Random cropping
  - Left+right flipping

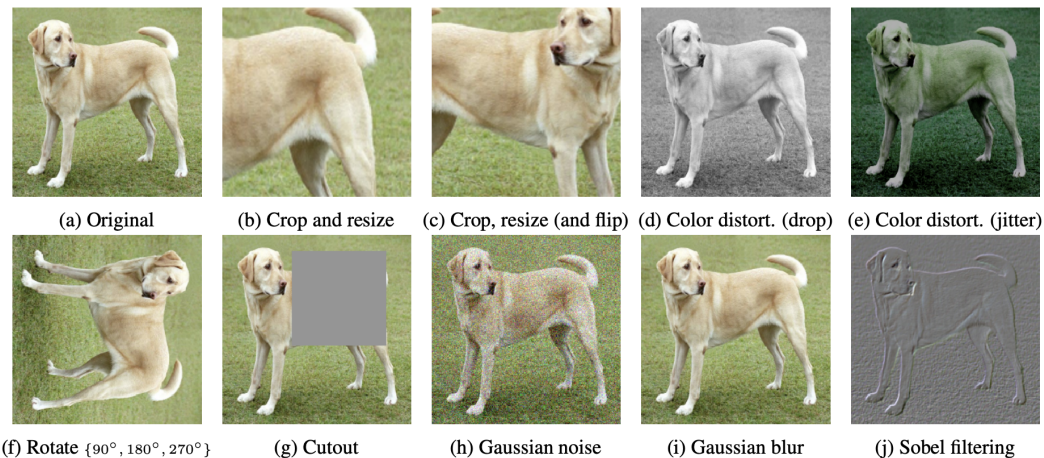


Image credit<sup>10</sup>

---

<sup>10</sup>Chen et al. A Simple Framework for Contrastive Learning of Visual Representations. ICML 2020.

# Data augmentation

- Leverage the invariances of images
- Create more data points for free
  - Random cropping
  - Left+right flipping
  - Random color jittering

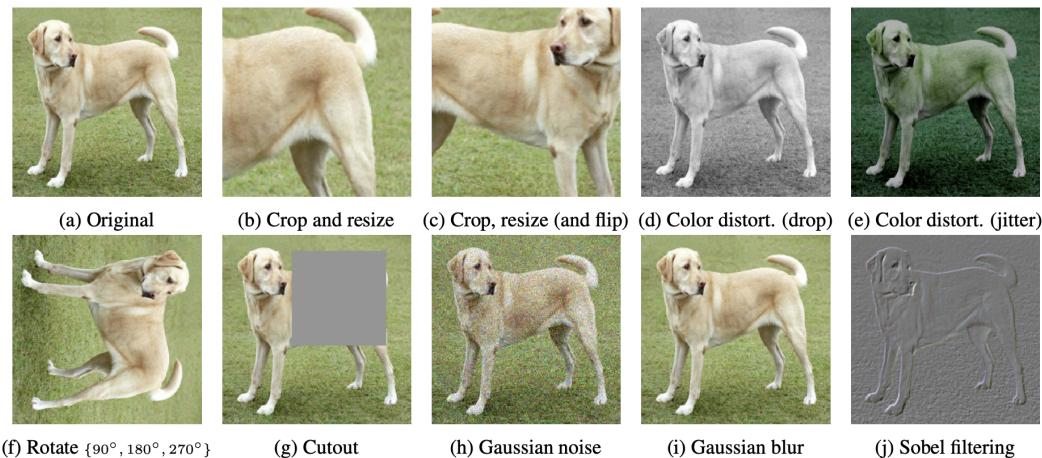


Image credit<sup>10</sup>

---

<sup>10</sup>Chen et al. A Simple Framework for Contrastive Learning of Visual Representations. ICML 2020.



# Data augmentation

- Leverage the invariances of images
- Create more data points for free
  - Random cropping
  - Left+right flipping
  - Random color jittering
  - Random blurring

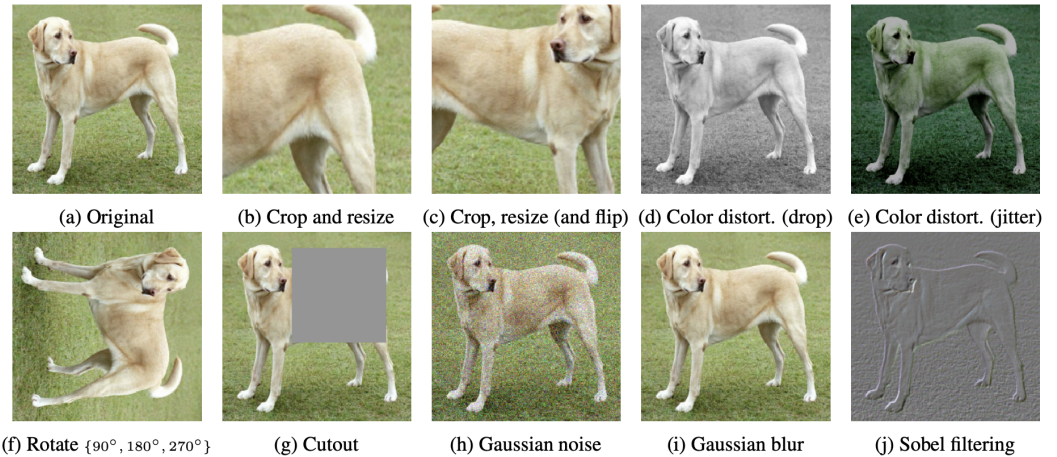


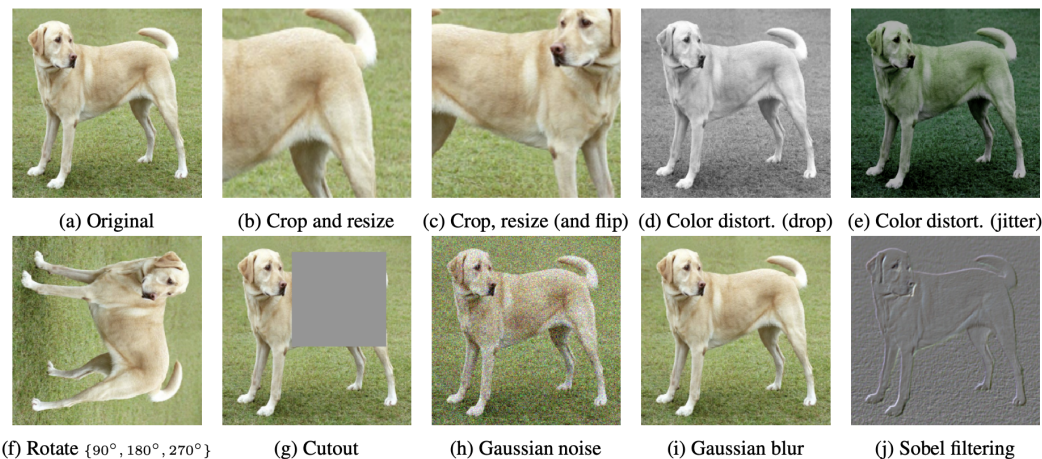
Image credit<sup>10</sup>

<sup>10</sup>Chen et al. A Simple Framework for Contrastive Learning of Visual Representations. ICML 2020.

# Data augmentation

- Leverage the invariances of images
- Create more data points for free
  - Random cropping
  - Left+right flipping
  - Random color jittering
  - Random blurring
  - Affine warping
  - Etc.

Image credit<sup>10</sup>



<sup>10</sup>Chen et al. A Simple Framework for Contrastive Learning of Visual Representations. ICML 2020.

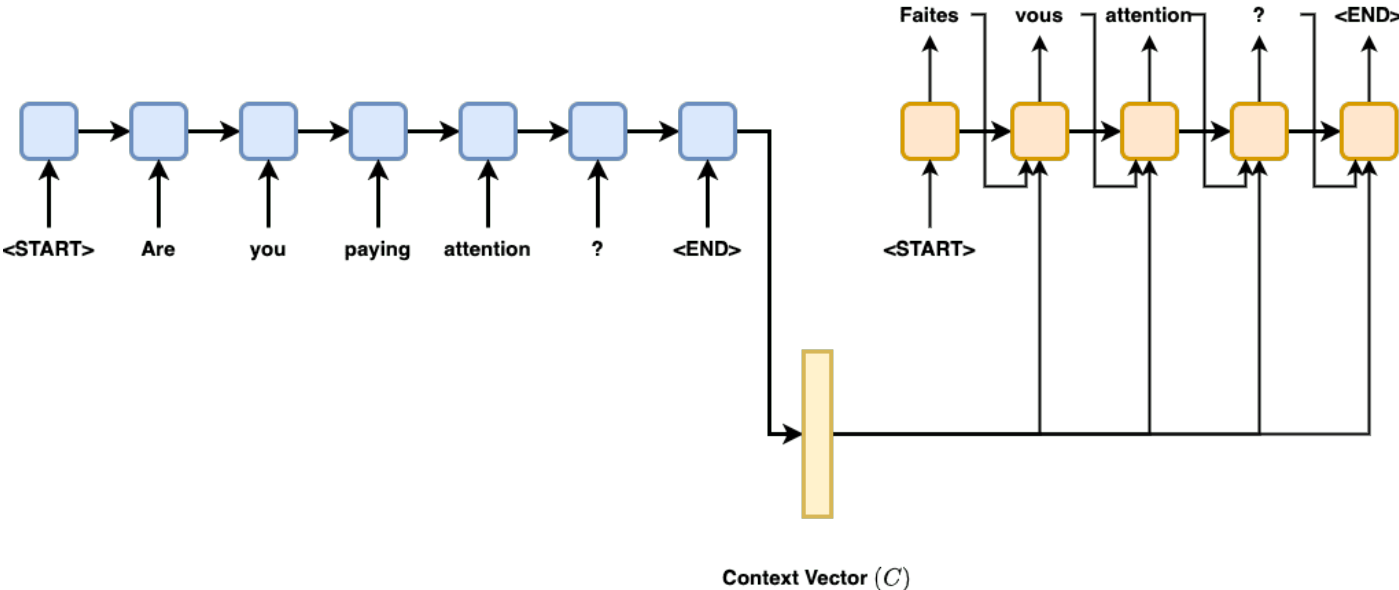


# Language and sequential signals

---

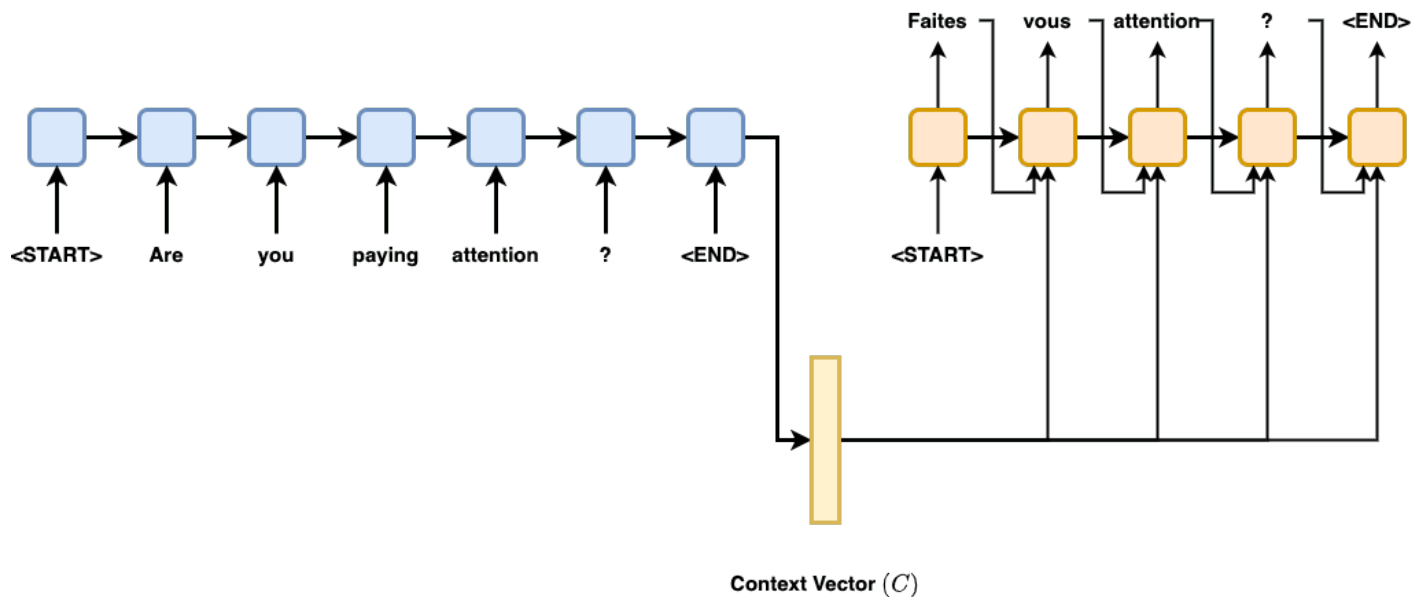
# What about natural language

- Neural networks are great for dealing with naturalistic and unstructured signals.



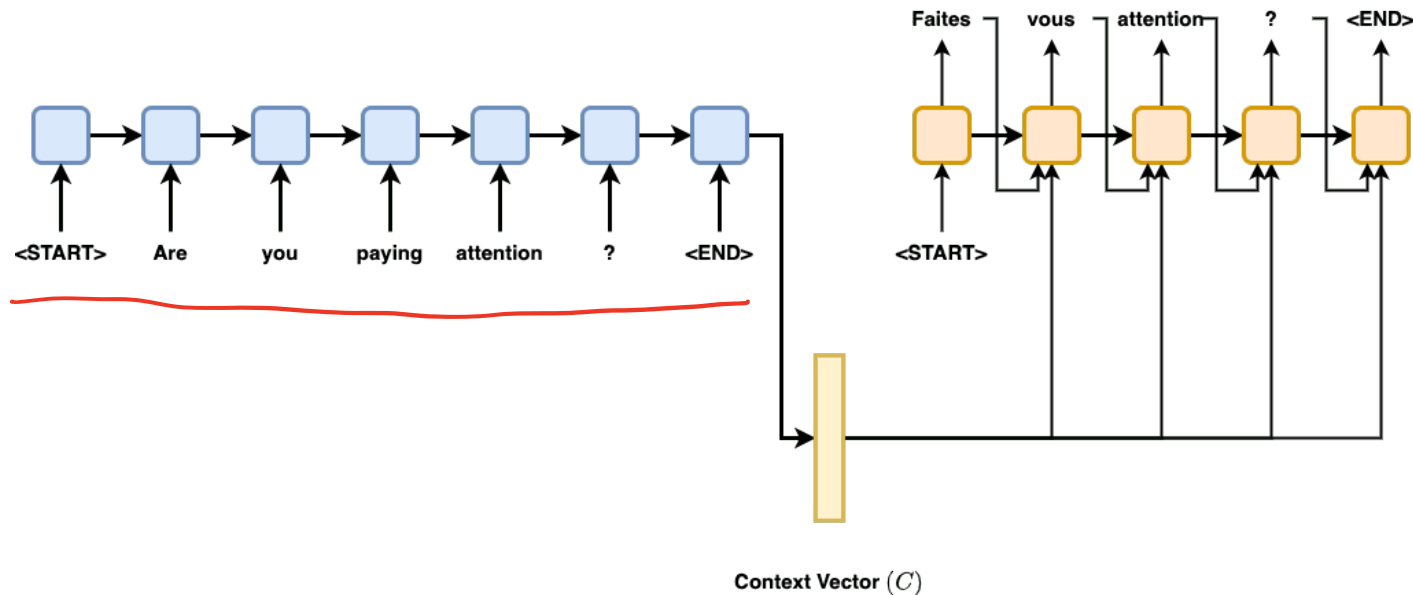
# What about natural language

- Neural networks are great for dealing with naturalistic and unstructured signals.
- Past lectures: **Feature functions** in structured models, but still primitive.



# What about natural language

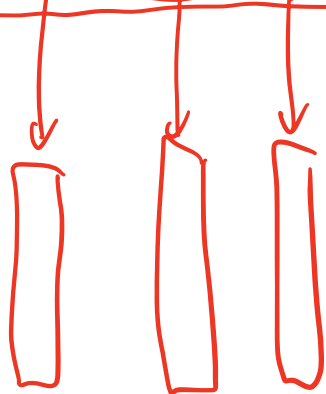
- Neural networks are great for dealing with naturalistic and unstructured signals.
- Past lectures: Feature functions in structured models, but still primitive.
- Design neural networks to accomodate sequential signals such as language.



# Word embeddings

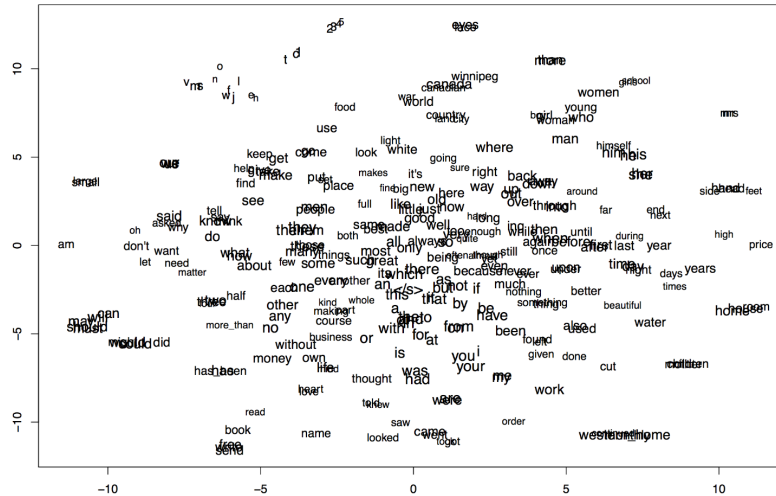
- Neural networks are best dealing with real valued vectors.

0  
1  
2.



vector.

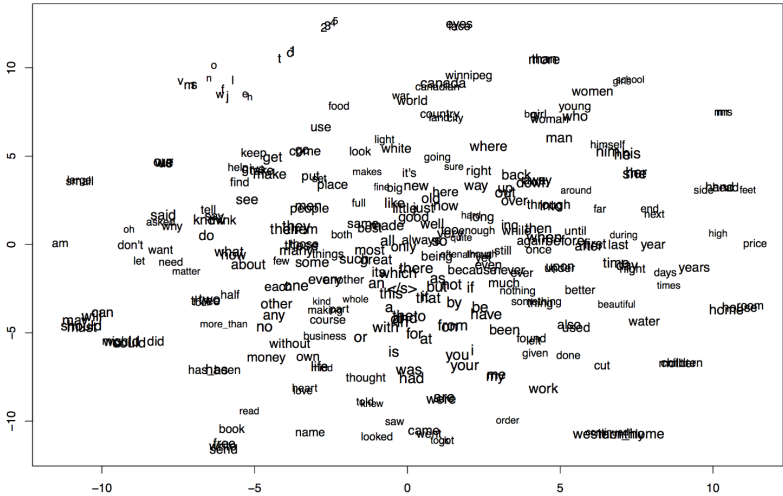
token embedding.



11

# Word embeddings

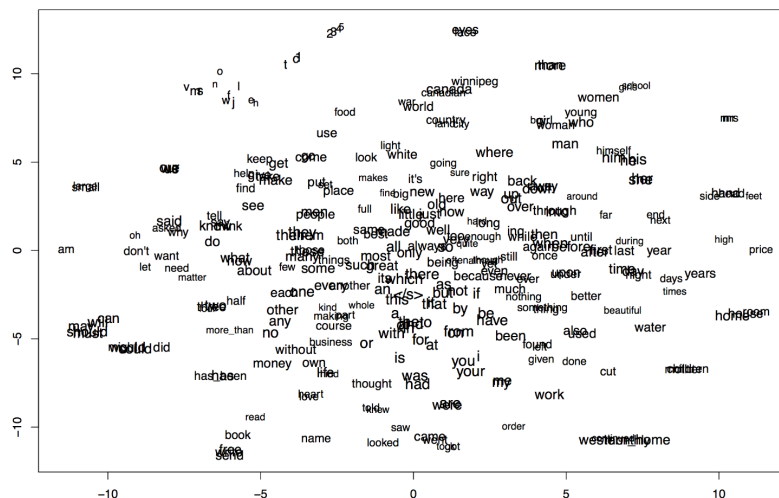
- Neural networks are best dealing with real valued vectors.
- Need to convert words (discrete) into vectors (continuous).



11

# Word embeddings

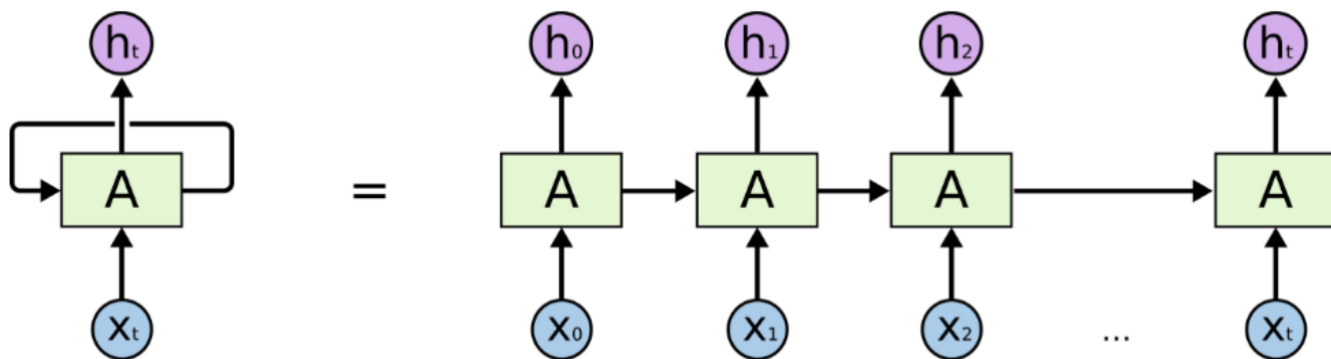
- Neural networks are best dealing with real valued vectors.
- Need to convert words (discrete) into vectors (continuous).
- A large matrix of  $V \times D$ .  $V$  = vocab size,  $D$  = network embedding size.



11

# Convolutional vs. recurrent networks

- Recall in images we used the convolution operation.
- We can also use the idea of convolution for temporal signals.

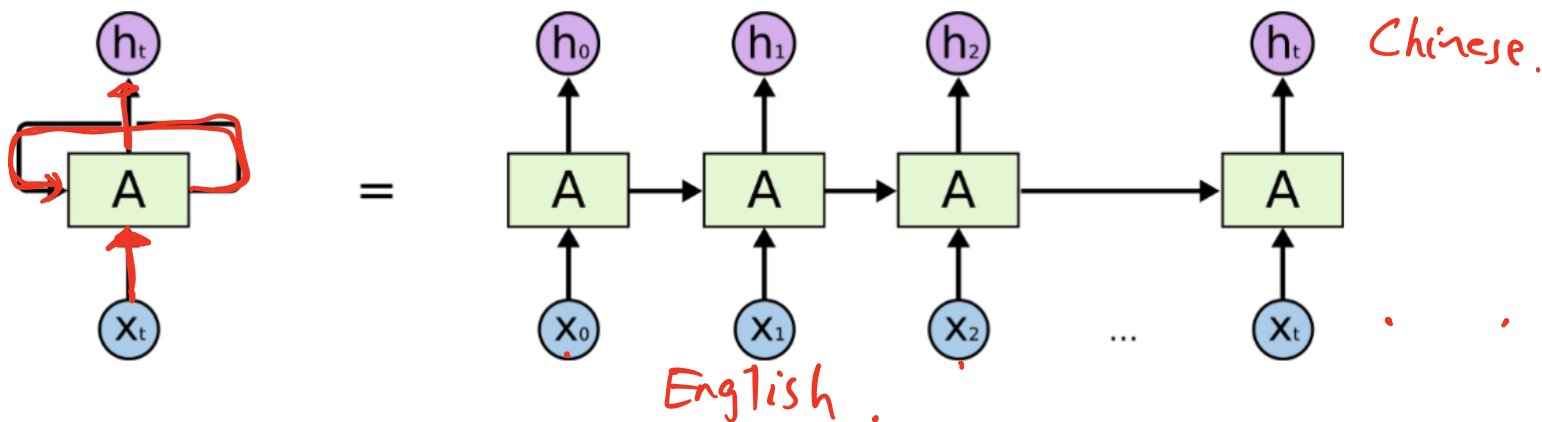




# Convolutional vs. recurrent networks

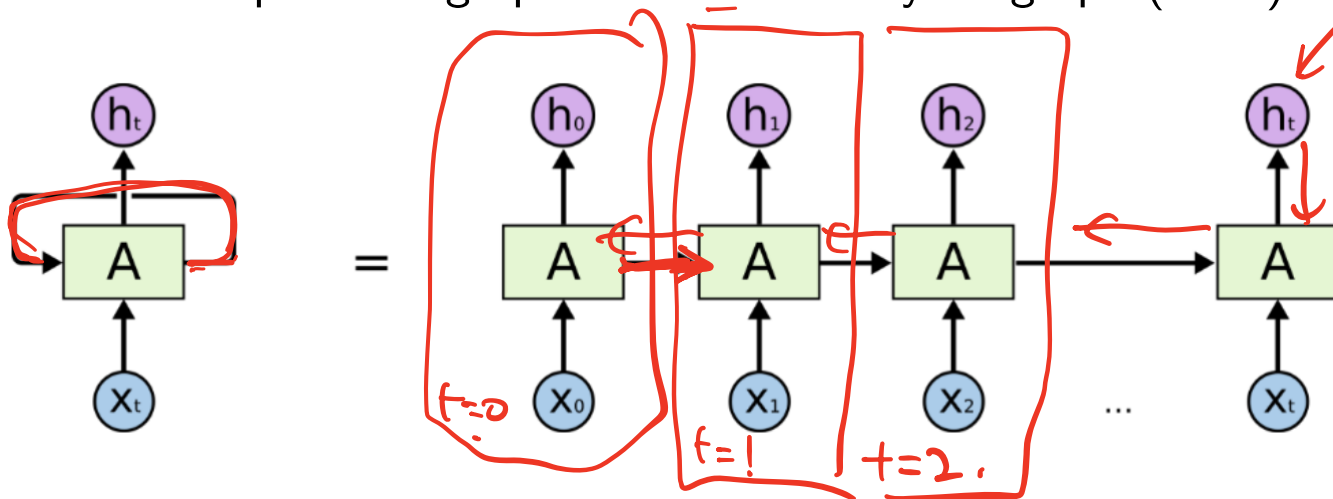
- Recall in images we used the convolution operation.
- We can also use the idea of convolution for temporal signals.
- Another alternative is to use a type of network called recurrent networks.
- Two inputs:  $x_t$  is the current input, and  $h_t$  is the historical hidden state.

*weight sharing on temporal level.* → RNN.



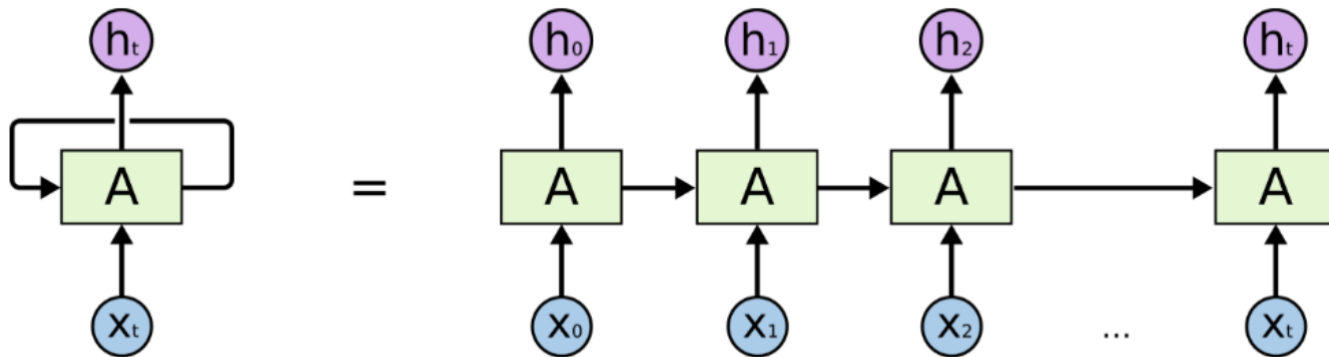
# Convolutional vs. recurrent networks

- Recall in images we used the convolution operation.
- We can also use the idea of convolution for temporal signals.
- Another alternative is to use a type of network called recurrent networks.
- Two inputs:  $x_t$  is the current input, and  $h_t$  is the historical hidden state.
- We can unroll the computation graph into a direct acyclic graph (DAG).



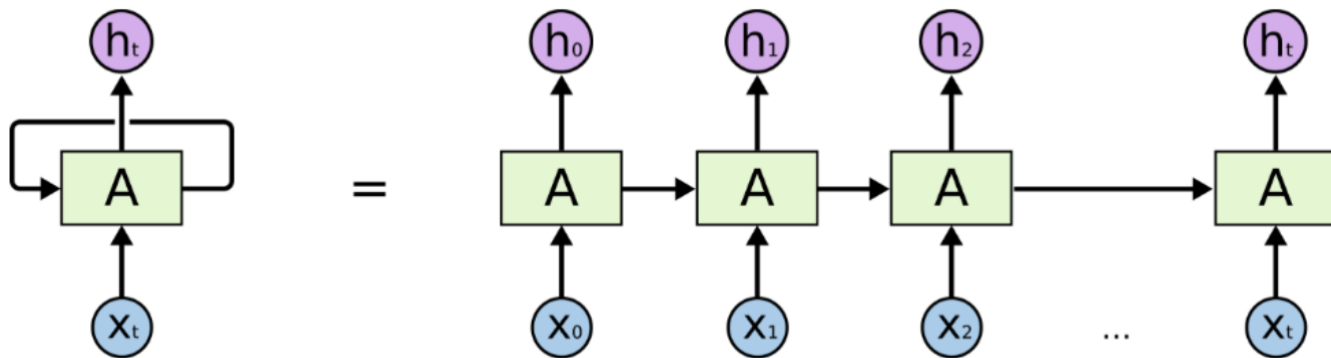
# Recurrent neural networks (RNNs)

- A simple RNN can be made similar to a standard NN with one hidden layer.



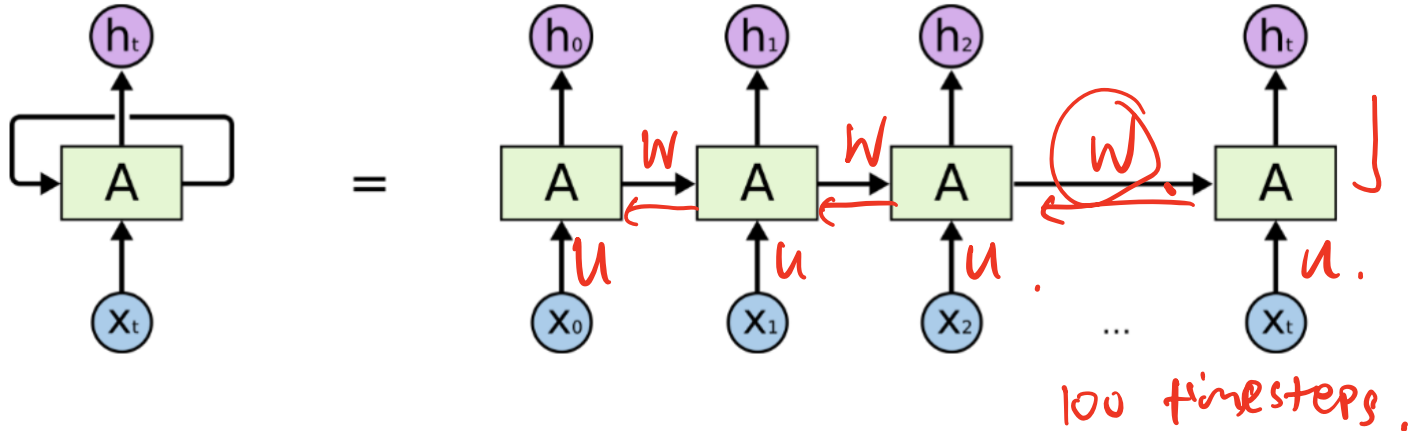
# Recurrent neural networks (RNNs)

- A simple RNN can be made similar to a standard NN with one hidden layer.
- $h_t = \tanh(W h_{t-1} + U x_t)$ .



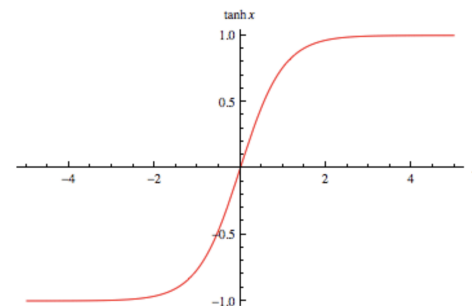
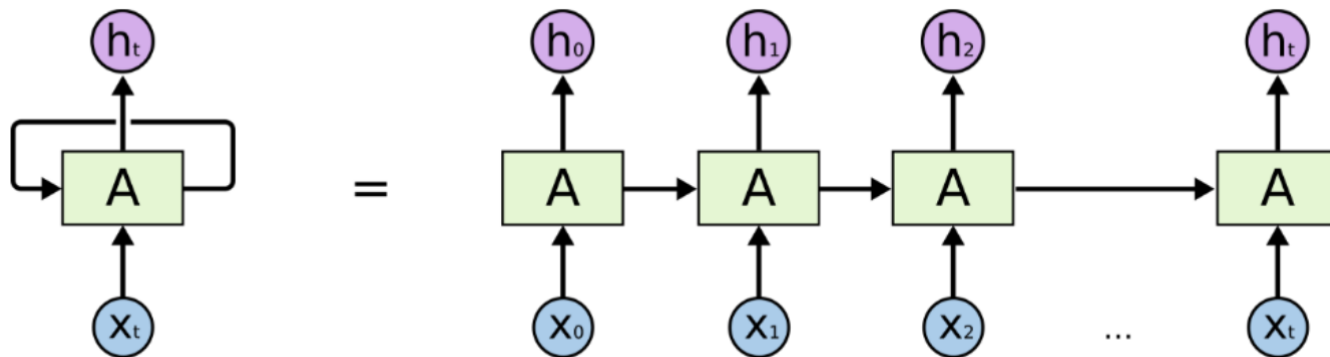
# Recurrent neural networks (RNNs)

- A simple RNN can be made similar to a standard NN with one hidden layer.
- $h_t = \tanh(W h_{t-1} + U x_t)$ .
- $y_t = \text{Softmax}(V h_t)$ .



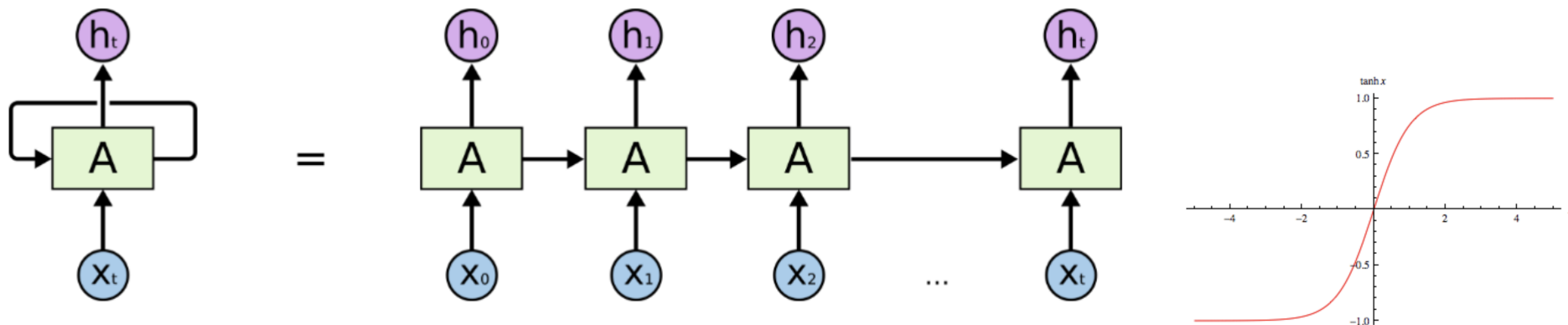
# Gradient vanishing

- Every iteration, we multiply the hidden state  $h_{t-1}$  from the previous iteration with the same  $W$ . Recall the definition of Jacobian.



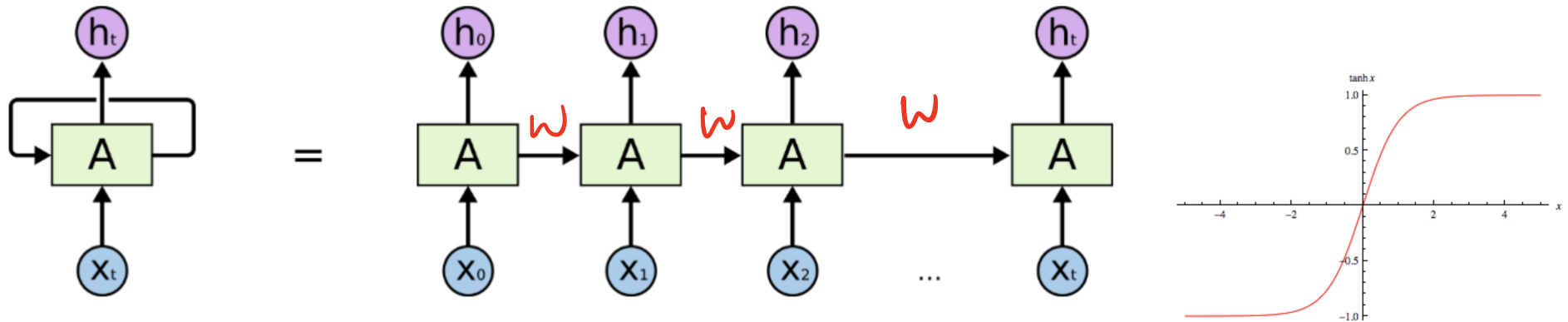
# Gradient vanishing

- Every iteration, we multiply the hidden state  $h_{t-1}$  from the previous iteration with the same  $W$ . Recall the definition of Jacobian.
- If the largest singular value of  $W$  is less than one then back-propagation will be attenuated.



# Gradient vanishing

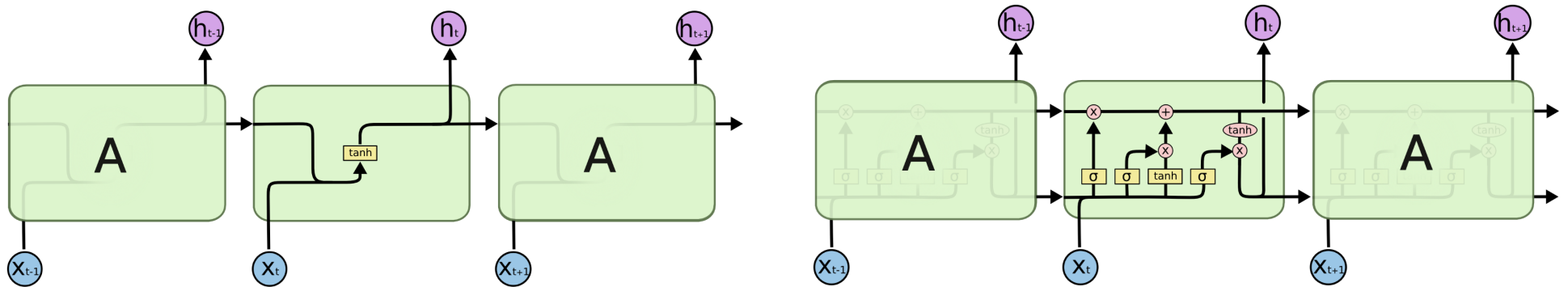
- Every iteration, we multiply the hidden state  $h_{t-1}$  from the previous iteration with the same  $W$ . Recall the definition of Jacobian.
- If the largest singular value of  $W$  is less than one then back-propagation will be attenuated.
- Similarly, we apply tanh activation every iteration – further reducing gradient flow.





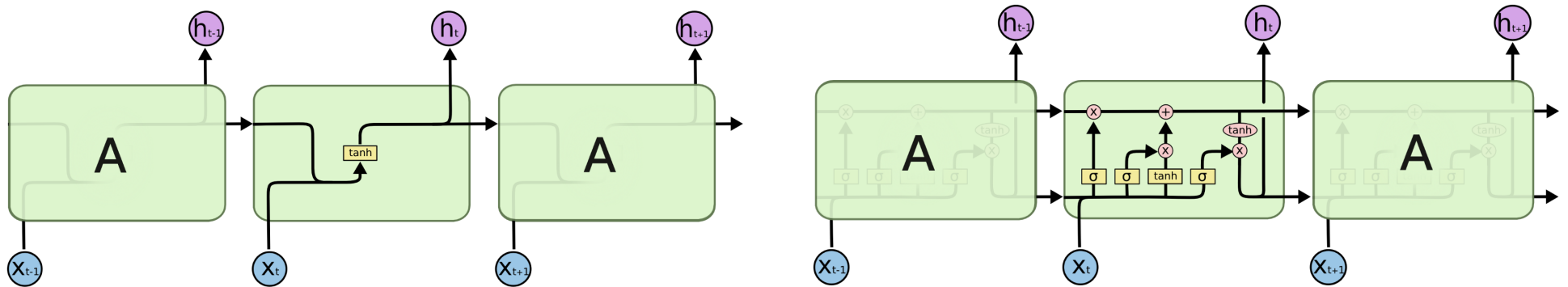
# Gating functions in LSTM

- Long short-term memory is a network that addresses the gradient vanishing problem by introducing gating functions.
- Gating functions provide “shortcuts”, like ResNet.



# Gating functions in LSTM

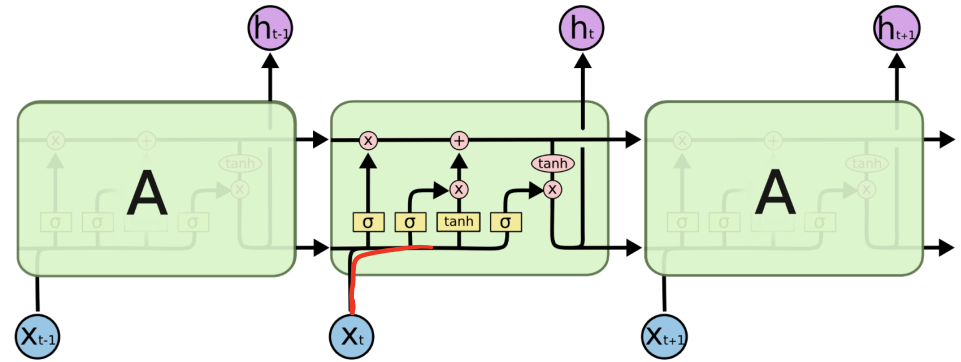
- Long short-term memory is a network that addresses the gradient vanishing problem by introducing gating functions.
- Gating functions provide “shortcuts”, like ResNet.
- Originally proposed by Hochreiter and Schmidhuber in 1997.



# Gating functions in LSTM

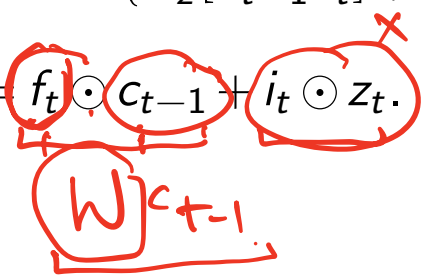
- Input gate:  $i_t = \sigma(W_i[h_{t-1}, x_t] + b_i)$ .
- Forget gate:  $f_t = \sigma(W_f[h_{t-1}, x_t] + b_f)$ .
- $z_t = \tanh(w_z[h_{t-1}, x_t] + b_z)$ .

*information to be consumed*  
*information get rid of.*

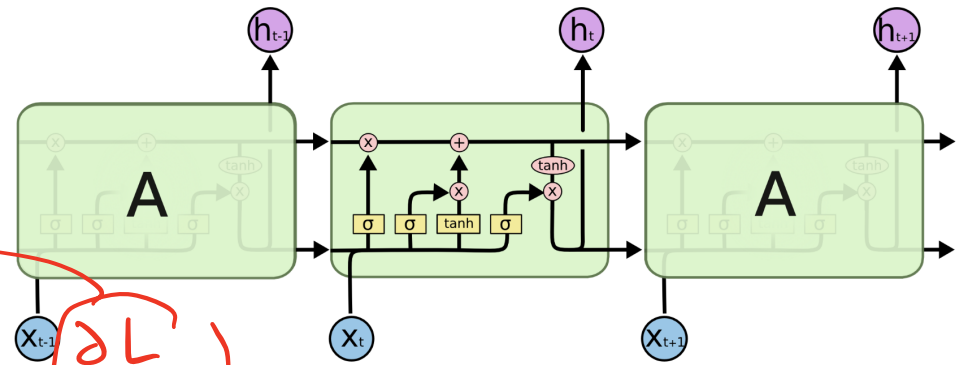


# Gating functions in LSTM

- Input gate:  $i_t = \sigma(W_i[h_{t-1}, x_t] + b_i)$ .
- Forget gate:  $f_t = \sigma(W_f[h_{t-1}, x_t] + b_f)$ .
- $z_t = \tanh(w_z[h_{t-1}, x_t] + b_z)$ .
- $c_t = f_t \odot c_{t-1} + i_t \odot z_t$ .

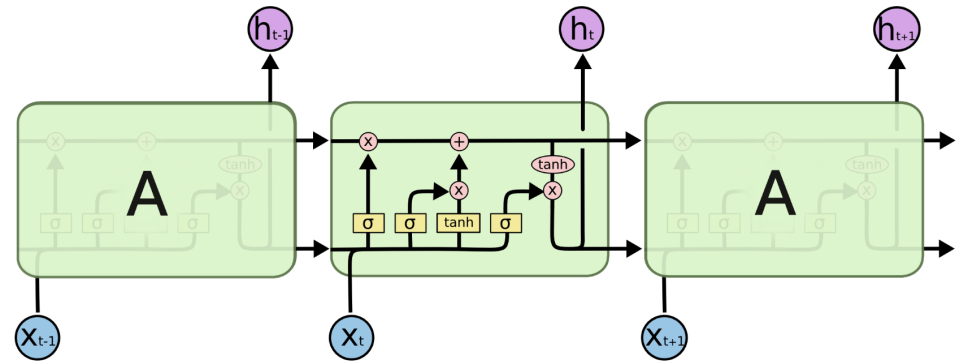


$$\frac{\partial \mathcal{L}}{\partial c_{t-1}} = f_t \odot \frac{\partial \mathcal{L}}{\partial c_t}$$



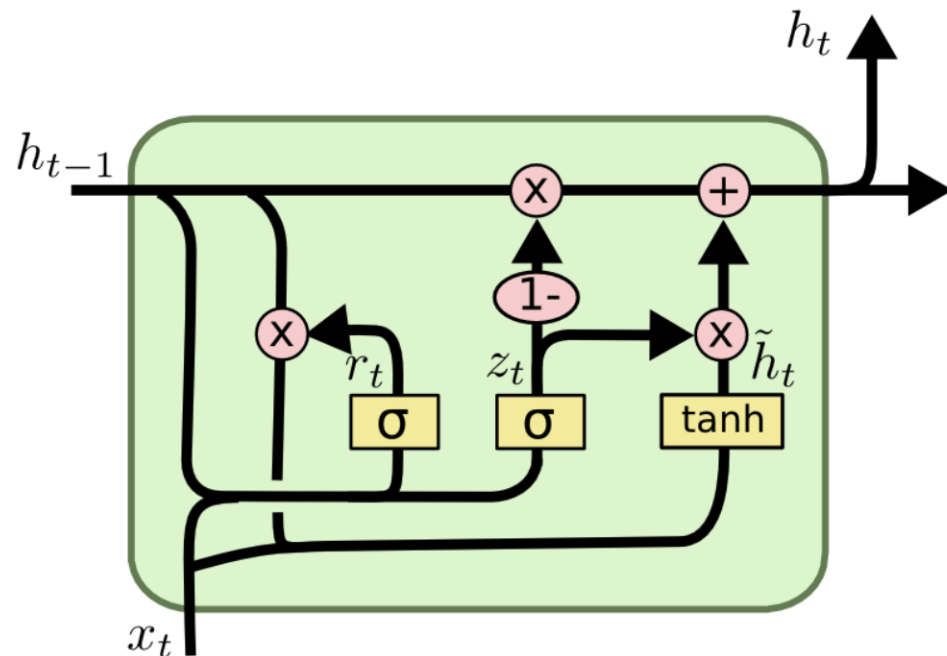
# Gating functions in LSTM

- Input gate:  $i_t = \sigma(W_i[h_{t-1}, x_t] + b_i)$ .
- Forget gate:  $f_t = \sigma(W_f[h_{t-1}, x_t] + b_f)$ .
- $z_t = \tanh(w_z[h_{t-1}, x_t] + b_z)$ .
- $c_t = f_t \odot c_{t-1} + i_t \odot z_t$ .
- Output gate:  $o_t = \sigma(W_o[h_{t-1}, x_t] + b_o)$ .
- $h_t = o_t \odot \tanh(c_t)$ .



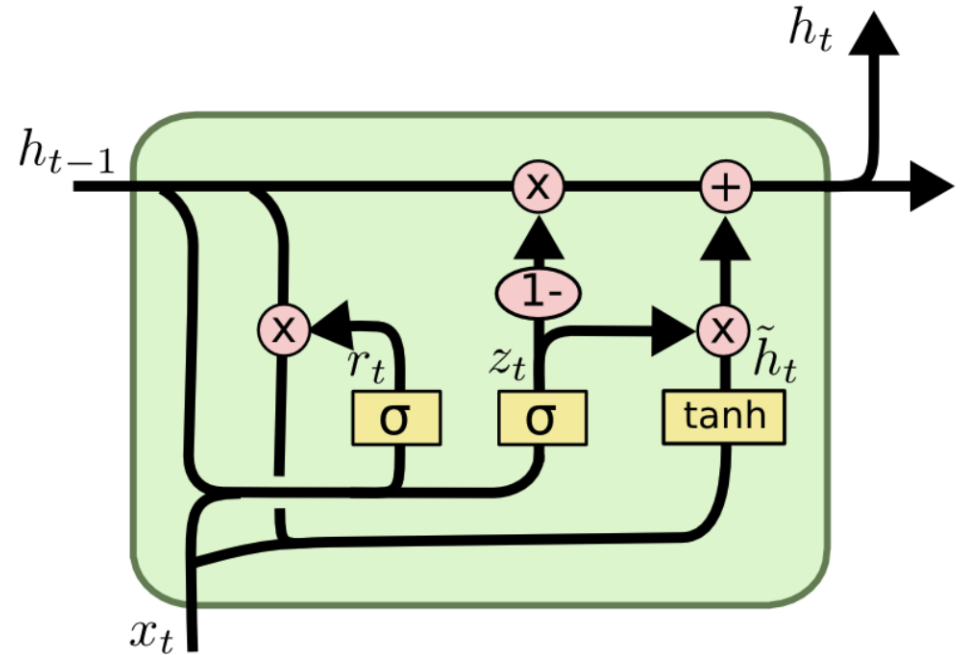
# Gated Recurrent Unit

- Proposed by Chung et al. in 2015, a simplified variant compared to LSTM.



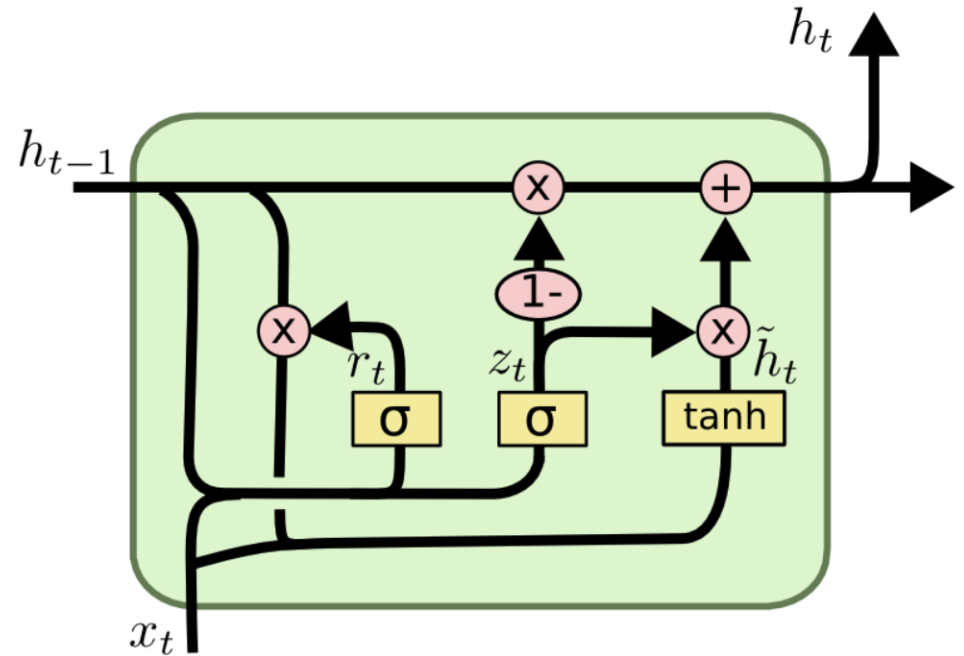
# Gated Recurrent Unit

- Proposed by Chung et al. in 2015, a simplified variant compared to LSTM.
- Input gate  $i_t = \sigma(W_i[h_{t-1}, x_t] + b_i)$ .



# Gated Recurrent Unit

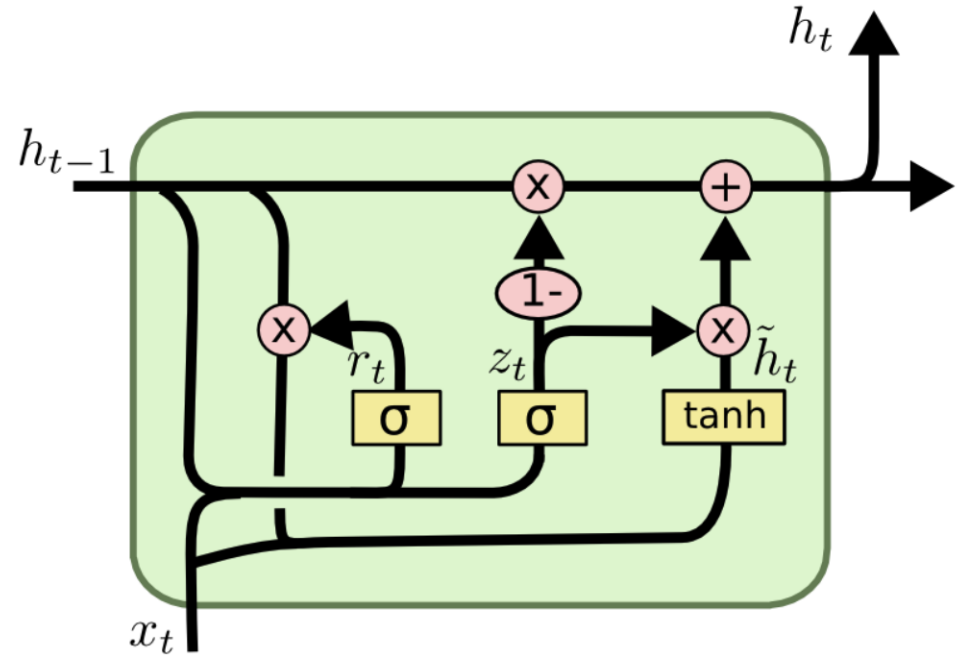
- Proposed by Chung et al. in 2015, a simplified variant compared to LSTM.
- Input gate  $i_t = \sigma(W_i[h_{t-1}, x_t] + b_i)$ .
- Reset gate  $r_t = \sigma(W_r[h_{t-1}, x_t] + b_r)$ .





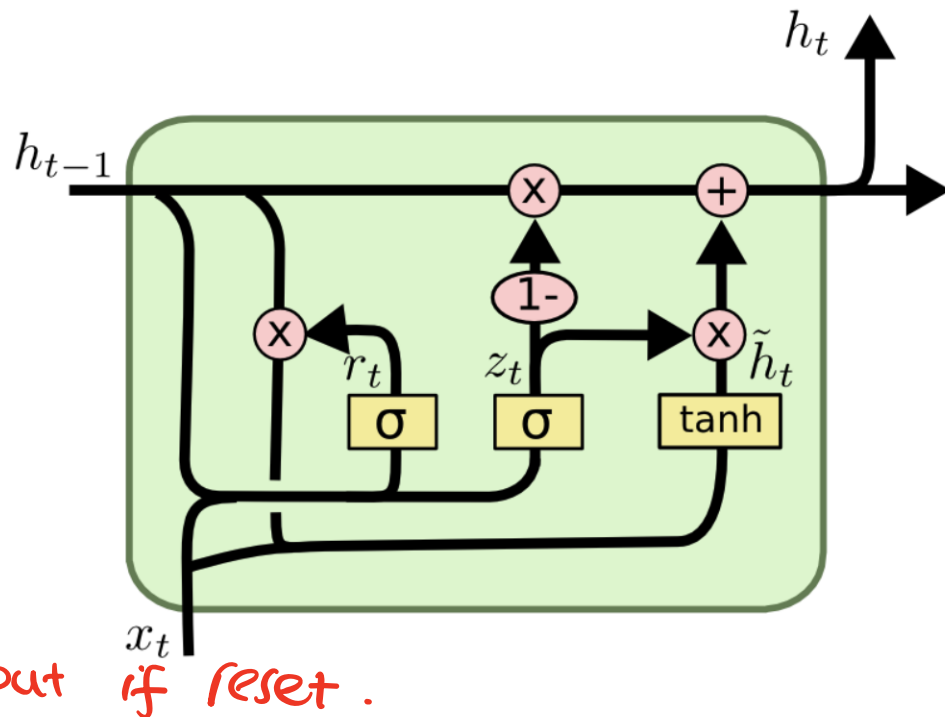
# Gated Recurrent Unit

- Proposed by Chung et al. in 2015, a simplified variant compared to LSTM.
- Input gate  $i_t = \sigma(W_i[h_{t-1}, x_t] + b_i)$ .
- Reset gate  $r_t = \sigma(W_r[h_{t-1}, x_t] + b_r)$ .
- $\tilde{h}_t = \tanh(W_h[r_t \odot h_{t-1}, x_t] + b_h)$ .



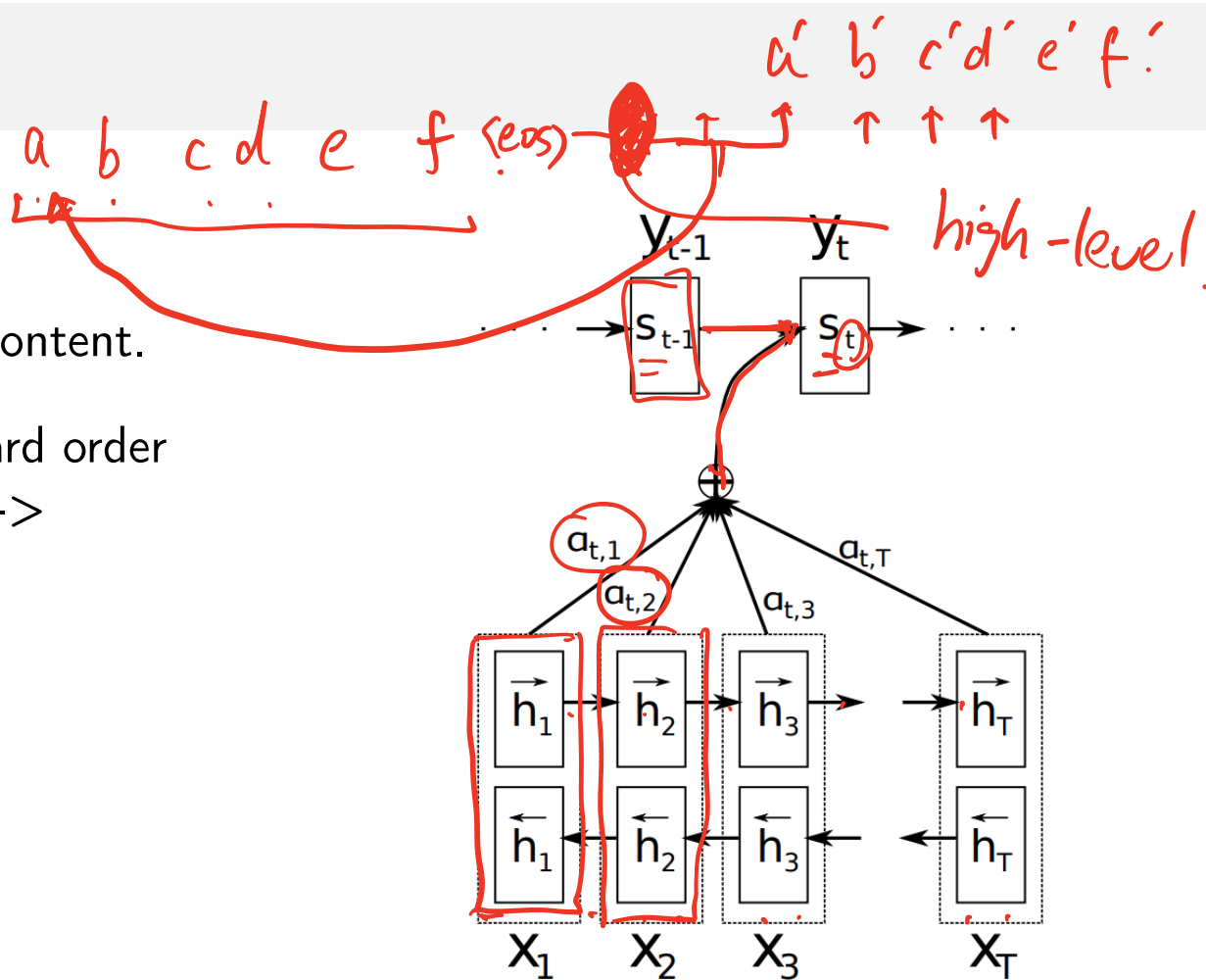
# Gated Recurrent Unit

- Proposed by Chung et al. in 2015, a simplified variant compared to LSTM.
- Input gate  $i_t = \sigma(W_i[h_{t-1}, x_t] + b_i)$ .
- Reset gate  $r_t = \sigma(W_r[h_{t-1}, x_t] + b_r)$ .
- $\tilde{h}_t = \tanh(W_h[r_t \odot h_{t-1}, x_t] + b_h)$ .
- $h_t = (1 - i_t) \odot h_{t-1} + i_t \odot \tilde{h}_t$ .



# Attention Mechanisms

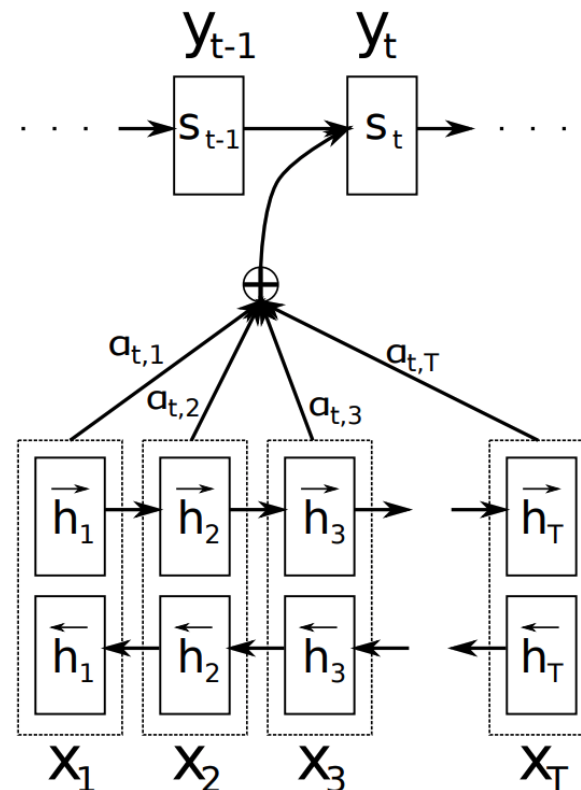
- Earlier content will decay more.
- Hard to refer back to the raw content.
- Reverse order better than forward order [abcde  $\rightarrow$  a'b'c'd'e' vs. abcde  $\rightarrow$  e'd'c'b'a'].



Bahdanau et al., 2014

# Attention Mechanisms

- Earlier content will decay more.
- Hard to refer back to the raw content.
- Reverse order better than forward order [abcde  $\rightarrow$  a'b'c'd'e' vs. abcde  $\rightarrow$  e'd'c'b'a'].
- Attending to arbitrary sequence tokens.



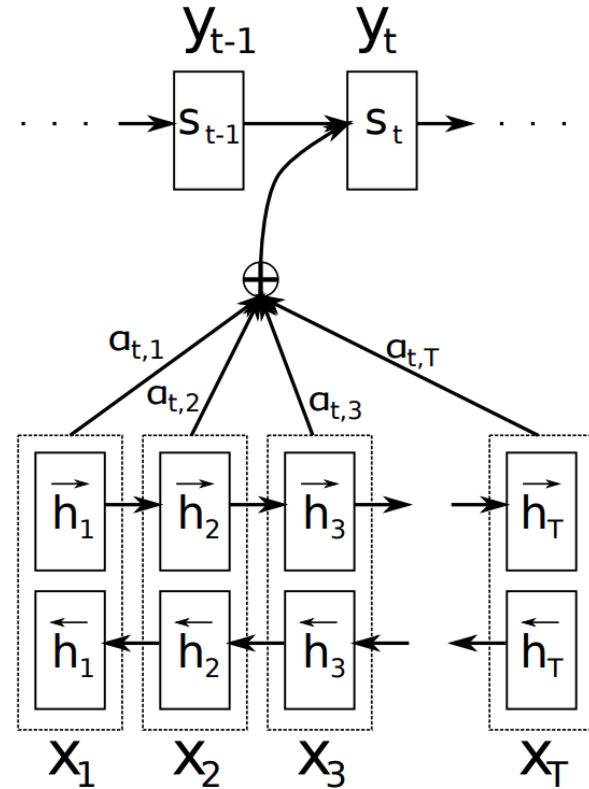
Bahdanau et al., 2014

# Attention Mechanisms

- Earlier content will decay more.
- Hard to refer back to the raw content.
- Reverse order better than forward order [abcde -> a'b'c'd'e' vs. abcde -> e'd'c'b'a'].

- Attending to arbitrary sequence tokens.

- $s_t = f(s_{t-1}, y_{t-1}, c_t)$  *ran output.*  
*Content attends to.*



Bahdanau et al., 2014

# Attention Mechanisms

- Earlier content will decay more.
- Hard to refer back to the raw content.
- Reverse order better than forward order [abcde -> a'b'c'd'e' vs. abcde -> e'd'c'b'a'].
- Attending to arbitrary sequence tokens.

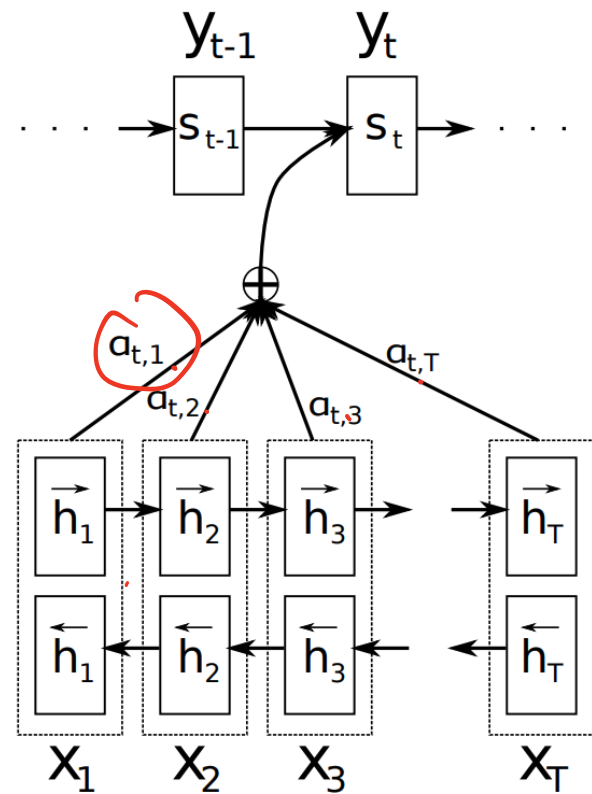
- $s_t = f(s_{t-1}, y_{t-1}, c_t)$

- $c_t = \sum_{\tau} \alpha_{t,\tau} h_{\tau}, \alpha_{t,\tau} = \frac{\exp(a(s_{t-1}, h_k))}{\sum_k \exp(a(s_{t-1}, h_k))}$

*Softmax. attention.*

*[0, 1, 0, 0] input content.*

*[0.5, 0.5, 0, 0]*



Bahdanau et al., 2014

# Attention Mechanisms

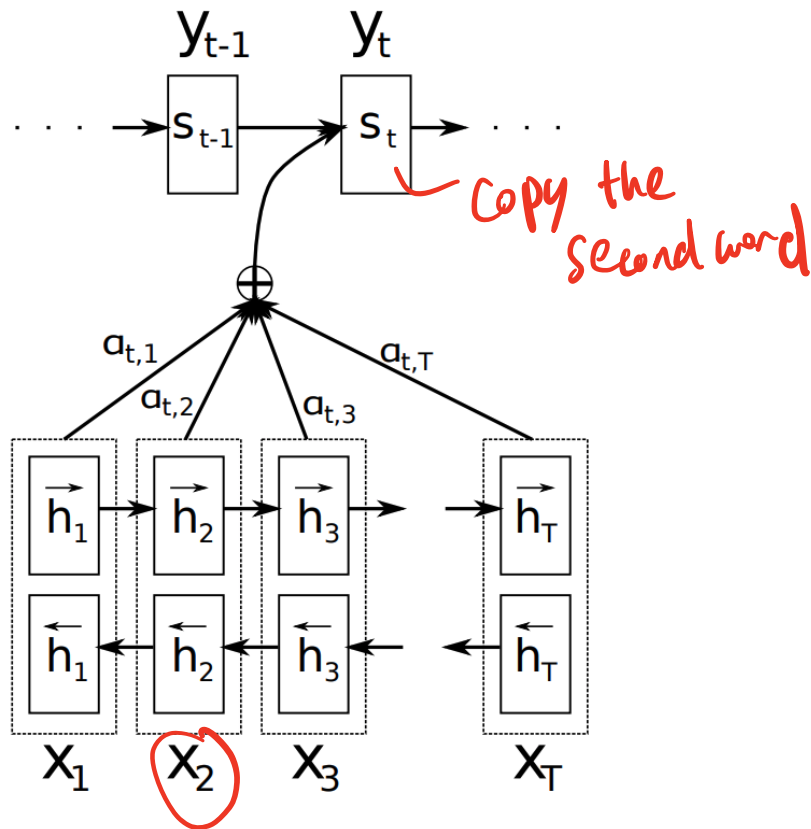
- Earlier content will decay more.
- Hard to refer back to the raw content.
- Reverse order better than forward order [abcde -> a'b'c'd'e' vs. abcde -> e'd'c'b'a'].

- Attending to arbitrary sequence tokens.

- $s_t = f(s_{t-1}, y_{t-1}, c_t)$

- $c_t = \sum_{\tau} \alpha_{t,\tau} h_{\tau}$ ,  $\alpha_{t,\tau} = \frac{\exp(a(s_{t-1}, h_k))}{\sum_k \exp(a(s_{t-1}, h_k))}$

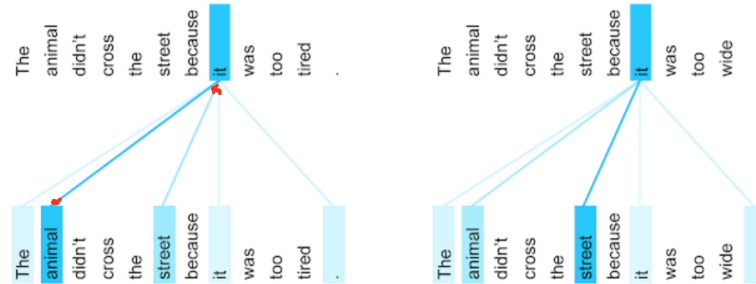
- $a(s_{t-1}, h_k) = v_a^T \tanh(W_a[s_{t-1}, h_k]) \rightarrow \text{scalar.}$   
*Compatible.*



Bahdanau et al., 2014

# Transformers (“Attention is All You Need”)

- The previous architecture is very complicated.
  - 1 RNN for encoding the tokens. ✗
  - Attention mechanisms for accessing content
  - 1 RNN for combining attended tokens. ✗

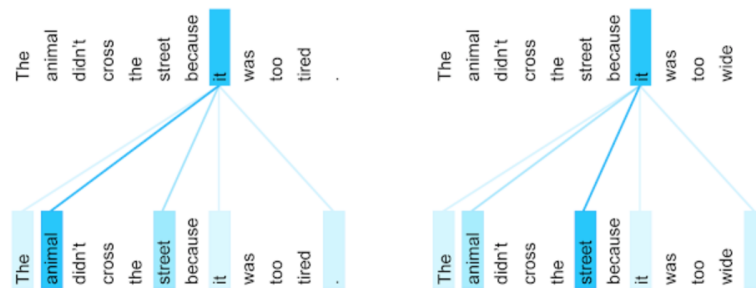


13



# Transformers (“Attention is All You Need”)

- The previous architecture is very complicated.
  - 1 RNN for encoding the tokens.
  - Attention mechanisms for accessing content
  - 1 RNN for combining attended tokens.
- RNNs have the ability to incorporate past information, so does attention.



13

# Positional encoding

- Attention operation is permutation equivariant.



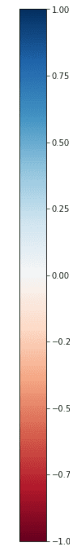
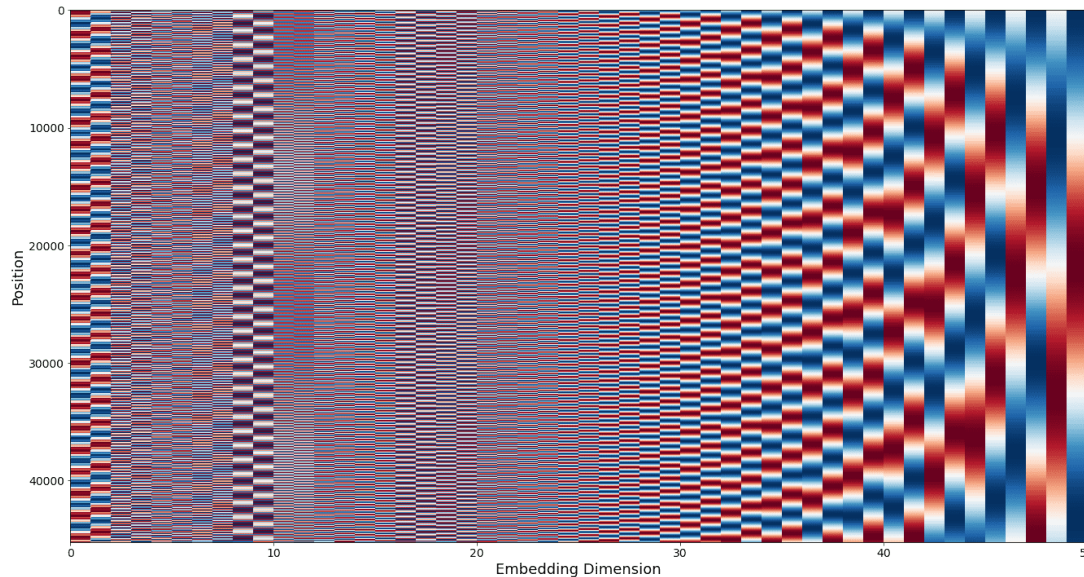
# Positional encoding

- Attention operation is permutation equivariant.
- Solution: Encode the position of each token.

# Positional encoding

- Attention operation is permutation equivariant.
- Solution. Encode the position of each token.
- $PE(pos, 2i) = \sin(p/k^{2i/d})$ ,  $PE(pos, 2i+1) = \cos(p/k^{2i/d})$ .

 pos.



— fixed dimension encoding.