# 15. Deep Learning

# Artificial Intelligence
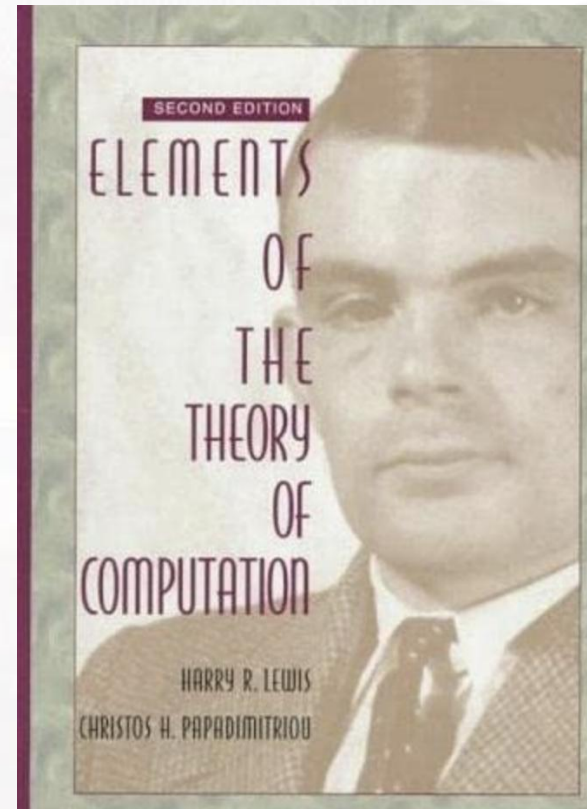
- **Alan Turing**
  - ☐ **Turing test, a method to assess a machine's ability to exhibit human-like intelligent behavior**
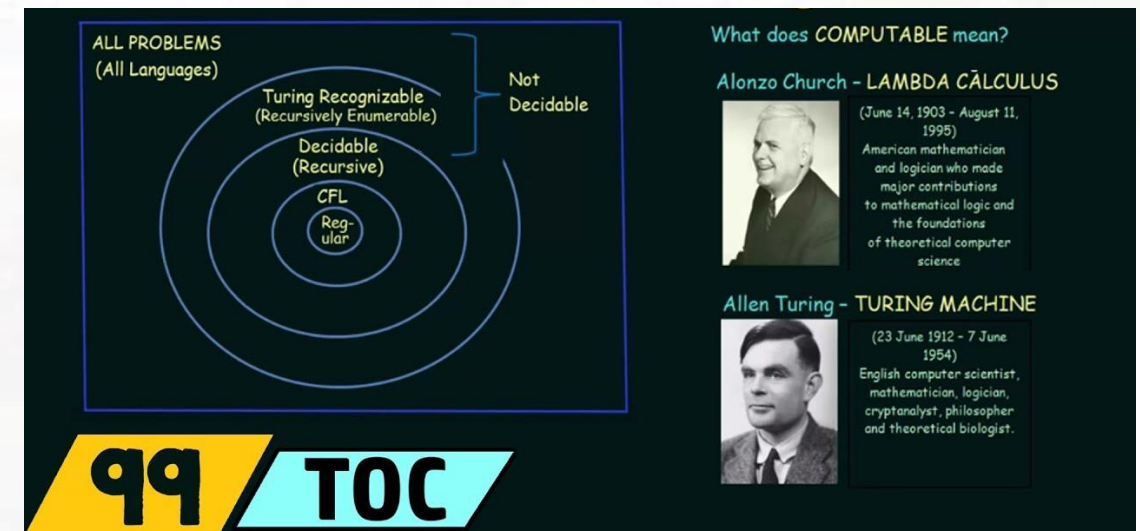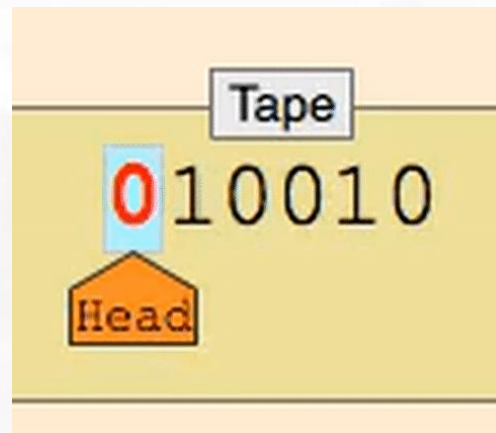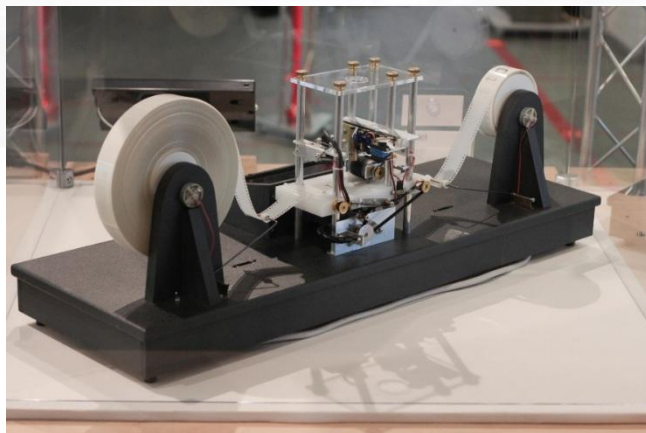
- **Artificial intelligence, or AI**
  - ☐ **technology that enables <span style="color:red">computers</span> and machines to simulate human intelligence and problem-solving capabilities.**



**Turing Test**

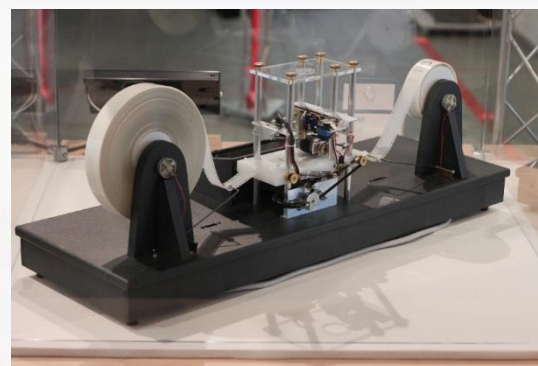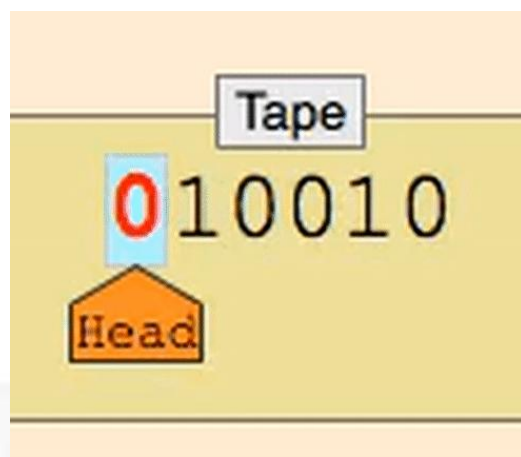# Artificial Intelligence

- **Turing Machine**
  - a mathematical model of computation describing an abstract machine that manipulates symbols on a strip of tape according to a table of rules
  - Despite the model's simplicity, it is capable of implementing **any computer algorithm**.
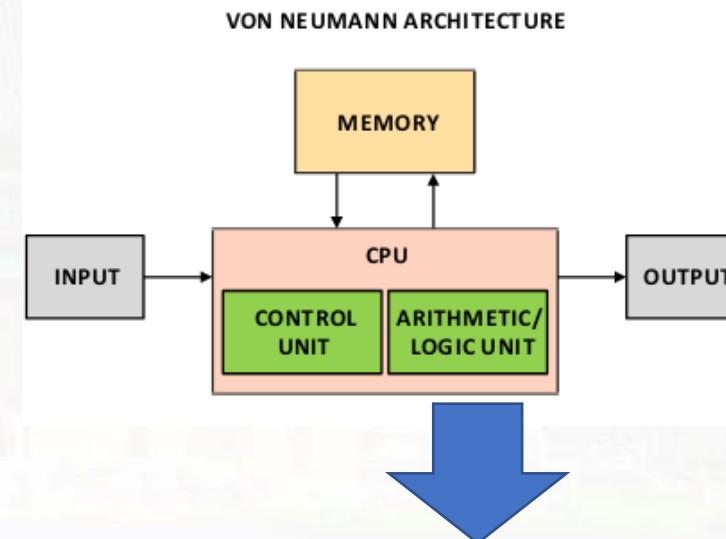
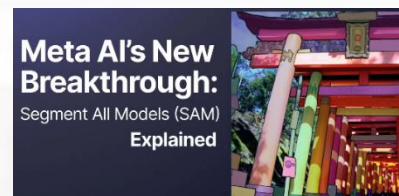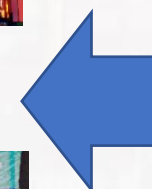# Artificial Intelligence



Turing Machine

Foundation model

Text to Image
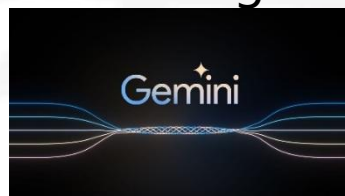
ChatGPT

SAM

SORA

MLLM

DriveGPT

Genie

# The Development of AI

# The Development of AI

# The Development of DL

# Data View: Fundamental Problems

■ **Output: Continuous description or discrete state**

# Data View: Fundamental Problems

■ **Output: Continuous description or discrete state**



Regression versus Classification

# Data view: Learning Paradigm

- **Availability of outputs: unavailable, available, partially available**

# Data view: Learning Paradigm

# Model view: neural network

| Aspect | Neural Networks | Support Vector Machines (SVM) |
|---|---|---|
| Model Type | Non-linear, based on layers of neurons. | Linear or non-linear (with kernel tricks), based on maximizing margins. |
| Data Requirements | Requires large datasets for training. | Works well with small to medium-sized datasets. |
| Complexity | Can model highly complex, non-linear relationships. | Handles non-linearity through kernel functions. |
| Training Time | Computationally expensive, especially for deep learning. | Slower training for large datasets, quadratic complexity in training. |
| Scalability | Scales well with data and features but requires more computational power. | Struggles with very large datasets and high-dimensional data. |
| Interpretability | Black-box model; difficult to interpret. | Easier to interpret, especially with linear SVM. |

# Model view: Neural Network

# Model view: neural network

- Large searching space and fitting capacity
- -> Scaling law
- ->Data-driven
- ->Fit everything
- ->Data is all you need and expert is gone

# Computational Neuron and Perceptron

# Biological Neuron

Action potentials are the primary means of information transmission in biological neurons. They implement spatiotemporal information processing through an event-driven threshold-triggering mechanism (when the sum of synaptic inputs reaches a critical value), ensuring a balance between signal transmission reliability and energy efficiency.



**Action potential illustration:** A brief and specially shaped transmembrane potential pulse generated when the cell membrane at resting membrane potential is subjected to an appropriate stimulus.



**Event-driven schematic:** When an electrical signal (action potential) reaches the threshold, it triggers a pulse and transmits information directionally through synaptic release of chemical transmitters.

# Biological Neuron

**Discrete Spiking representation**: An action potential is expressed in an **all-or-none** manner.



**Asynchronous spiking information transmission**: Whether the current neuron transmits information is **unrelated to whether other neurons transmit information**; it only depends on whether the threshold of the current neuron is triggered.



Equivalent circuit

# Computational Neuron

# Perception

$$\varphi(\vec{X}) = \text{sgn}\left(\sum_{i=1..n} w_i \times x_i + b\right)$$

# Example: Internet Traffic Prediction

- Suppose someone wants to make money through a video platform; they would care about whether the channel has traffic, so they would know their potential earnings.

- Assume that the backend can see a lot of relevant information, such as the number of people who like posts each day, the number of subscribers, and the number of views.

- Based on all the past information of a channel, it is possible to predict the number of views for tomorrow.

- Find a function whose input is the backend information and whose output is the total number of views the channel will have the next day.

# Solution: Three Steps

- **Design a function with unknown parameters**
  - ☐ **model:** $f$、**feature:** $x_1$
  - ☐ **parameter:** $b, w$
  - ☐ **weight:** $w$、**bias:** $b$
- **Define loss function:** $L$
  - ☐ **Mean Absolute Error, MAE:** $e = |\hat{y} - y|$
  - ☐ **Mean Squared Error, MSE:** $e = (\hat{y} - y)^2$
  - ☐ **cross entropy**
- **Solve an optimization problem**
  - ☐ **global minimum and local minimum**

$$y = b + wx_1$$

2017/01/01　01/02　01/03　……　2020/12/30　12/31

4800　4900　7500　3400　9800

$$\hat{y} = 500 + 1x_1 \qquad e_1 = |y - \hat{y}| = 400$$

$$e_2 = |y - \hat{y}| = 2100 \qquad L = \frac{1}{N}\sum_n e_n$$

# Advantage of Perception



$$y = 100 + 0.97x_1$$

red:real viewings
blue:estimated viewings

viewings
(in thousands)

2021/01/01
2021/02/14

$$y = b + wx_1$$

different w

diferent b

# Multi-Layered Perceptron

# Multi-Layered Perceptron

# Multi-Layered Perceptron

- **Understand MLP, understand everything**
  - A few layers of neurons linked together.



Plain vanilla
(aka "multilayer perceptron")



Bias(Offset) Value

Bias(Offset)Weight Value

$$z = b + \sum_{i=1}^{N} a_i w_i$$

$$a_{out} = g(z)$$

# Multi-Layered Perceptron

- **Neuron: a thing that holds a number**
  - represent the inputs and outputs of our network (the images and digit predictions) in terms of these neuron values

# Multi-Layered Perceptron

**what kind of digit does this network think it's looking at? How certain does it feel?**

# Multi-Layered Perceptron

**■ More Compact Notation**



$$a_0^{(1)} = \sigma\left( w_{0,0}\, a_0^{(0)} + w_{0,1}\, a_1^{(0)} + \cdots + w_{0,n}\, a_n^{(0)} + b_0 \right)$$

Sigmoid

Bias

# Multi-Layered Perceptron

# Multi-Layered Perceptron

# Multi-Layered Perceptron



Sigmoid

$$a_0^{(1)} = \sigma\left( w_{0,0}\ a_0^{(0)} + w_{0,1}\ a_1^{(0)} + \cdots + w_{0,n}\ a_n^{(0)} + b_0 \right)$$

Bias

$$\begin{bmatrix} w_{0,0} & w_{0,1} & \cdots & w_{0,n} \\ w_{1,0} & w_{1,1} & \cdots & w_{1,n} \\ \vdots & \vdots & \ddots & \vdots \\ w_{k,0} & w_{k,1} & \cdots & w_{k,n} \end{bmatrix} \begin{bmatrix} a_0^{(0)} \\ a_1^{(0)} \\ \vdots \\ a_n^{(0)} \end{bmatrix}$$

# Multi-Layered Perceptron



$$a_0^{(1)} = \sigma\left( w_{0,0}\, a_0^{(0)} + w_{0,1}\, a_1^{(0)} + \cdots + w_{0,n}\, a_n^{(0)} + b_0 \right)$$

Sigmoid

Bias

$$\begin{bmatrix} w_{0,0} & w_{0,1} & \cdots & w_{0,n} \\ w_{1,0} & w_{1,1} & \cdots & w_{1,n} \\ \vdots & \vdots & \ddots & \vdots \\ w_{k,0} & w_{k,1} & \cdots & w_{k,n} \end{bmatrix} \begin{bmatrix} a_0^{(0)} \\ a_1^{(0)} \\ \vdots \\ a_n^{(0)} \end{bmatrix} + \begin{bmatrix} b_0 \\ b_1 \\ \vdots \\ b_n \end{bmatrix}$$

# Multi-Layered Perceptron

# Multi-Layered Perceptron

- **The result of the weighted sum like this can be any number, but for this network we want the activations to be values between 0 and 1.**



$$w_1 a_1 + w_2 a_2 + w_3 a_3 + w_4 a_4 + \cdots + w_n a_n$$

Activations should be in this range



Sigmoid

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

# Multi-Layered Perceptron

**■ Formulated function**

Superscript corresponds to the layer

$$a_0^{(1)} = \sigma\left(w_{0,0}a_0^{(0)} + w_{0,1}a_1^{(0)} + \cdots + w_{0,n}a_n^{(0)} + b_0\right)$$

Subscript corresponds to a neuron in the layer

Neural network function

784 inputs

13,002 weights
and biases

0
1
2
3
4
5
6
7
8
9

10 outputs

# Backpropagation

# Optimization

- **The behavior of the network depends on its weights and biases.**



$$= \sigma(w_1 a_1 + w_2 a_2 + \cdots + w_n a_n + b)$$

# Optimization

■ **Data Preparation**

# Optimization

# Optimization

- **finding the minimum of a specific function**

# Optimization

■ **The Cost Function**

# Optimization



Cost function

13,002 weights and biases

$(9, 9)(0, 0)(2, 2)(6, 6)$
$(0, 0)(4, 4)(6, 6)(7, 7)$
$(7, 7)(8, 8)(3, 3)(1, 1)$
$(1, 1)(1, 1)(6, 6)(3, 3)$
$(1, 1)(1, 1)(0, 0)(4, 4)$

Lots of training data

3.37

One number

# Optimization



$C(w)$
Single input

$w_{min}$

$w$

# Optimization

■ **For a complicated cost function, computing the exact minimum directly isn't going to work.**

# Optimization

■ **By following the slope (moving in the downhill direction), we approach a local minimum.**

# Optimization

■ **As the slope gets shallower, take smaller steps to avoid overshooting the minimum.**

# Optimization

- **Moving the input position according to the slope is a lot like a ball rolling down a hill.**

# Optimization

- **We can imagine minimizing a function that takes two inputs**
- **Gradient descent just means walking in the downhill direction to minimize the cost function.**

# Optimization

- **The gradient, $\nabla C$, gives the uphill direction, so the negative of the gradient, $-\nabla C$, gives the downhill direction.**

# Optimization

- **Another Way to Think About The Gradient**

- **Changing a weight that has a larger magnitude in the negative gradient vector has a bigger effect on the cost.**

# Forward and Backward

- **Calculate the prediction error for training data**
- **Update model parameters based on the error**

# Backpropagation

- **Approximate partial derivatives based on the chain rule**
- **Solving neural network parameters based on gradient descent**

# Chrain Rule



$$\sigma\left(\begin{bmatrix} w_{0,0} & w_{0,1} & \cdots & w_{0,n} \\ w_{1,0} & w_{1,1} & \cdots & w_{1,n} \\ \vdots & \vdots & \ddots & \vdots \\ w_{k,0} & w_{k,1} & \cdots & w_{k,n} \end{bmatrix} \begin{bmatrix} a_0^{(0)} \\ a_1^{(0)} \\ \vdots \\ a_n^{(0)} \end{bmatrix} + \begin{bmatrix} b_0 \\ b_1 \\ \vdots \\ b_n \end{bmatrix}\right)$$

Change by $-0.08$

$w_0$

$\left.\begin{array}{c} w_1 \\ w_2 \\ \vdots \\ w_{13,001} \end{array}\right\}$ All weights and biases

- **Matrix operations form**

- **sample-wise and parameter-wise differentiation**

$$\frac{\partial C_0}{\partial w^{(L-1)}} = \frac{\partial z^{(L-1)}}{\partial w^{(L-1)}} \frac{\partial a^{(L-1)}}{\partial z^{(L-1)}} \frac{\partial z^{(L)}}{\partial a^{(L-1)}} \frac{\partial a^{(L)}}{\partial z^{(L)}} \frac{\partial C_0}{\partial a^{(L)}}$$

$$C_0 = \left(a^{(L)} - y\right)^2 \longrightarrow \frac{\partial C_0}{\partial a^{(L)}} = 2\left(a^{(L)} - y\right)$$

$$a^{(L)} = \sigma\left(z^{(L)}\right) \longrightarrow \frac{\partial a^{(L)}}{\partial z^{(L)}} = \sigma'\left(z^{(L)}\right)$$

$$z^{(L)} = w^{(L)} a^{(L-1)} + b^{(L)} \longrightarrow \frac{\partial z^{(L)}}{\partial a^{(L-1)}} = w^{(L)}$$

$$z^{(L)} = w^{(L)} a^{(L-1)} + b^{(L)} \longrightarrow \frac{\partial z^{(L)}}{\partial w^{(L)}} = a^{(L-1)}$$

$$-\nabla C(\vec{\mathbf{W}}) = \begin{bmatrix} 0.31 \\ 0.03 \\ -1.25 \\ \vdots \\ 0.78 \\ -0.37 \\ 0.16 \end{bmatrix}$$

$w_0$ should increase somewhat
$w_1$ should increase a little
$w_2$ should decrease a lot

$w_{13,000}$ should increase a lot
$w_{13,001}$ should decrease somewhat
$w_{13,002}$ should increase a little

# Gradient Descent





- **Algorithm**
  - Calculate $\nabla C$ based on chain rules
  - Move all parameters toward $-\nabla C$ slightly
  - Repeat all the steps above

- (randomly) select $\boldsymbol{\theta}_0$
- get gradient $g_0 = \nabla L_1(\boldsymbol{\theta}_0)$
  $$\text{update } \boldsymbol{\theta}_1 \leftarrow \boldsymbol{\theta}_0 - \eta g_0$$
- get gradient $g_1 = \nabla L_2(\boldsymbol{\theta}_1)$
  $$\text{update } \boldsymbol{\theta}_2 \leftarrow \boldsymbol{\theta}_1 - \eta g_1$$
- get gradient $g_2 = \nabla L_3(\boldsymbol{\theta}_2)$
  $$\text{update } \boldsymbol{\theta}_3 \leftarrow \boldsymbol{\theta}_2 - \eta g_2$$

# The Problem of Backpropagation

# The History of Backpropagation

- **The time when the relevant ideas were proposed**
  - ☐ 1986，the term was widely known
  - ☐ 1974, the first time used to train neural network
  - ☐ 1960, The basic knowledge of backpropagation has been accepted and widely used.
  - ☐ 2006, Introduction of pre-training and fine-tuning mechanisms
- **Issues During Training**
  - ☐ **Gradient vanishing**
  - ☐ **Gradient exploding**

Published: 09 October 1986

## Learning representations by back-propagating errors

David E. Rumelhart, Geoffrey E. Hinton & Ronald J. Williams

_Nature_ **323**, 533–536 (1986) | Cite this article

**110k** Accesses | **14712** Citations | **378** Altmetric | Metrics

## Gradient Theory of Optimal Flight Paths

HENRY J. KELLEY[1]

Grumman Aircraft Engineering Corp.
Bethpage, N. Y.

An analytical development of flight performance optimization according to the method of gradients or "method of steepest descent" is presented. Construction of a minimizing sequence of flight paths by a stepwise process of descent along the local gradient direction is described as a computational scheme. Numerical application of the technique is illustrated in a simple example of orbital transfer via solar sail propulsion. Successive approximations to minimum time planar flight paths from Earth's orbit to the orbit of Mars are presented for cases corresponding to free and fixed boundary conditions on terminal velocity components.

# Gradient Vanishing and Exploding



| Gradient Vanishing | Gradient Exploding |
|---|---|
| The weights are **almost 0** | Weights with **NaN value** |
| The weights near the output layer **update quickly,** while the weights in the input layer **hardly update.** | Weights **increase explosively** |
| **Converges** very slowly | **Instable** performance |
| **Stop learning** | **Performs poorly on the** training set |

# Gradient Vanishing and Exploding

$$\frac{\partial C_0}{\partial w^{(L-1)}} = \frac{\partial z^{(L-1)}}{\partial w^{(L-1)}} \boxed{\frac{\partial a^{(L-1)}}{\partial z^{(L-1)}} \frac{\partial z^{(L)}}{\partial a^{(L-1)}}} \frac{\partial a^{(L)}}{\partial z^{(L)}} \frac{\partial C_0}{\partial a^{(L)}}$$

$$a^{(L)} = \sigma\left(z^{(L)}\right) \quad \longrightarrow \quad \boxed{\frac{\partial a^{(L)}}{\partial z^{(L)}} = \sigma'\left(z^{(L)}\right)}$$

$$z^{(L)} = w^{(L)} a^{(L-1)} + b^{(L)} \quad \longrightarrow \quad \boxed{\frac{\partial z^{(L)}}{\partial a^{(L-1)}} = w^{(L)}}$$



Sigmoid activation function

Saturating    Saturating

Derivative of the Sigmoid

Close to 0    Close to 0

# Gradient Vanishing and Exploding

■ **Use greedy unsupervised learning to find a sensible set of weights one layer at a time. Then fine-tune with backpropagation.**

  ☐ **The precious information in the labels is only used for the final fine-tuning.**

  ☐ **We do not start backpropagation until we already have sensible weights that already do well at the task.**



Unsupervised Layer-wise Pretraining     →      Supervised Fine-tuning

# Gradient Vanishing and Exploding

1. *Zero Initialization*: Initialize all the weights and biases to zero. This is not generally used in deep learning as it leads to symmetry in the gradients, resulting in all the neurons learning the same feature.

2. *Random Initialization:* Initialize the weights and biases randomly from a uniform or normal distribution. This is the most common technique used in deep learning.

3. *Xavier Initialization:* Initialize the weights with a normal distribution with mean 0 and variance of sqrt(1/n), where n is the number of neurons in the previous layer. This is used for the sigmoid activation function.

4. *He Initialization:* Initialize the weights with a normal distribution with mean 0 and variance of sqrt(2/n), where n is the number of neurons in the previous layer. This is used for the ReLU activation function.

5. *Orthogonal Initialization:* Initialize the weights with an orthogonal matrix, which preserves the gradient norm during backpropagation.

6. *Uniform Initialization:* Initialize the weights with a uniform distribution. This is less commonly used than random initialization.

7. *Constant Initialization:* Initialize the weights and biases with a constant value. This is rarely used in deep learning.



Sigmoid activation zero function

Derivative of the Sigmoid

# Gradient Vanishing and Exploding



**Batch 1 after backpropagation**

Input Features — Hidden Layer 1 — Hidden Layer 2 — Output

Distribution X   Distribution Y   Distribution Z

**Batch 2 after backpropagation**

Input Features — Hidden Layer 1 — Hidden Layer 2 — Output

Distribution X   Distribution B   Distribution C

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m}\sum_{i=1}^{m} x_i \qquad \text{// mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m}\sum_{i=1}^{m} (x_i - \mu_{\mathcal{B}})^2 \qquad \text{// mini-batch variance}$$

$$\widehat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \qquad \text{// normalize}$$

$$y_i \leftarrow \gamma \widehat{x}_i + \beta \equiv \mathrm{BN}_{\gamma,\beta}(x_i) \qquad \text{// scale and shift}$$

Without Normalization

With Normalization

**Internal Covariate Shift**

- **Slow Convergence**
- **Gradient Vanishing**

(a) loss landscape   (b) gradient predictiveness

Standard    Standard + BatchNorm

# Gradient Vanishing and Exploding

## ResNets @ ILSVRC & COCO 2015 Competitions

- **1st places** in all five main tracks
  - ImageNet Classification: *"Ultra-deep"* 152-layer nets
  - ImageNet Detection: 16% better than 2nd
  - ImageNet Localization: 27% better than 2nd
  - COCO Detection: 11% better than 2nd
  - COCO Segmentation: 12% better than 2nd

$$a^{(L)} = a^{(L-2)} + H(x)$$

$$\frac{\partial a^{(L)}}{\partial a^{(L-2)}} = \boxed{1} + \frac{\partial H(x)}{\partial a^{(L-2)}}$$

$$\frac{\partial C_0}{\partial w^{(L-2)}} = \frac{\partial z^{(L-2)}}{\partial w^{(L-2)}} \cdot \boxed{\frac{\partial a^{(L-2)}}{\partial z^{(L-2)}} \cdot \frac{\partial z^{(L-1)}}{\partial a^{(L-2)}} \cdot \frac{\partial a^{(L-1)}}{\partial z^{(L-1)}} \cdot \frac{\partial z^{(L)}}{\partial a^{(L-1)}} \cdot \frac{\partial a^{(L)}}{\partial z^{(L)}}} \cdot \frac{\partial C_0}{\partial a^{(L)}}$$

# Gradient Vanishing and Exploding

- **ReLU**
- **Gradient clipping** (**gradient exploding**) $\nabla C = \eta \dfrac{\nabla C}{|\nabla C|_2}$

$$\frac{\partial C_0}{\partial w^{(L-1)}} = \frac{\partial z^{(L-1)}}{\partial w^{(L-1)}} \boxed{\frac{\partial a^{(L-1)}}{\partial z^{(L-1)}}} \frac{\partial z^{(L)}}{\partial a^{(L-1)}} \frac{\partial a^{(L)}}{\partial z^{(L)}} \frac{\partial C_0}{\partial a^{(L)}}$$

# Gradient Vanishing and Exploding

- **Batch gradient descent**

$$W = W - \eta \cdot \nabla_W C(W)$$

- **Stochastic gradient descent**

$$W = W - \eta \cdot \nabla_W C(W; x^{(i)}; y^{(i)})$$

- **Mini-batch gradient descent**

$$W = W - \eta \cdot \nabla_W C(W; x^{(i:i+n)}; y^{(i:i+n)})$$

# SGDM

- **Momentum**
  - ☐ SGD has the problem of slow convergence in 'valley'-shaped spaces
  - ☐ SGDM speeds up SGD by modifying the direction.
- **Add an update direction of a past moment**

$$W \leftarrow W - \eta \frac{\partial L}{\partial W}$$

$$V_t \leftarrow \beta V_{t-1} - \eta \frac{\partial L}{\partial W}$$

$$W \leftarrow W + V_t$$



Image 2: SGD without momentum

Image 3: SGD with momentum

# SGDM

■ **Advantage**
- ☐ **faster convergence**
- ☐ **chance of escaping the local minimum**

$$\theta = \theta - \eta \cdot \nabla_\theta J(\theta; x^{(i:i+n)}; y^{(i:i+n)}).$$

$$v_t = \gamma v_{t-1} + \eta \nabla_\theta J(\theta)$$

$$\theta = \theta - v_t$$



$\nabla L(\theta^0)$
$\theta^0$
$\nabla L(\theta^1)$
$\theta^1$
$\nabla L(\theta^2)$
$\theta^2$
$\theta^3$
$\nabla L(\theta^3)$

→ Gradient
→ Movement
...... Movement of last step

# AdaGrad

$$\theta_{t+1,i} = \theta_{t,i} - \frac{\eta}{\sqrt{G_{t,i} + \epsilon}} \cdot g_{t,i}$$

$$G_t^{(i,i)} = \sum_\tau^t (g_\tau^{(i)})^2$$

$$g_{t,i} = \nabla_\theta J(\theta_{t,i})$$

# RMSProp

$$\theta_{t+1,i} = \theta_{t,i} - \frac{\eta}{\sqrt{G_{t,i} + \epsilon}} \cdot g_{t,i}$$

$$G_t^{(i,i)} = \sum_\tau^t (g_\tau^{(i)})^2$$

$$g_{t,i} = \nabla_\theta J(\theta_{t,i})$$

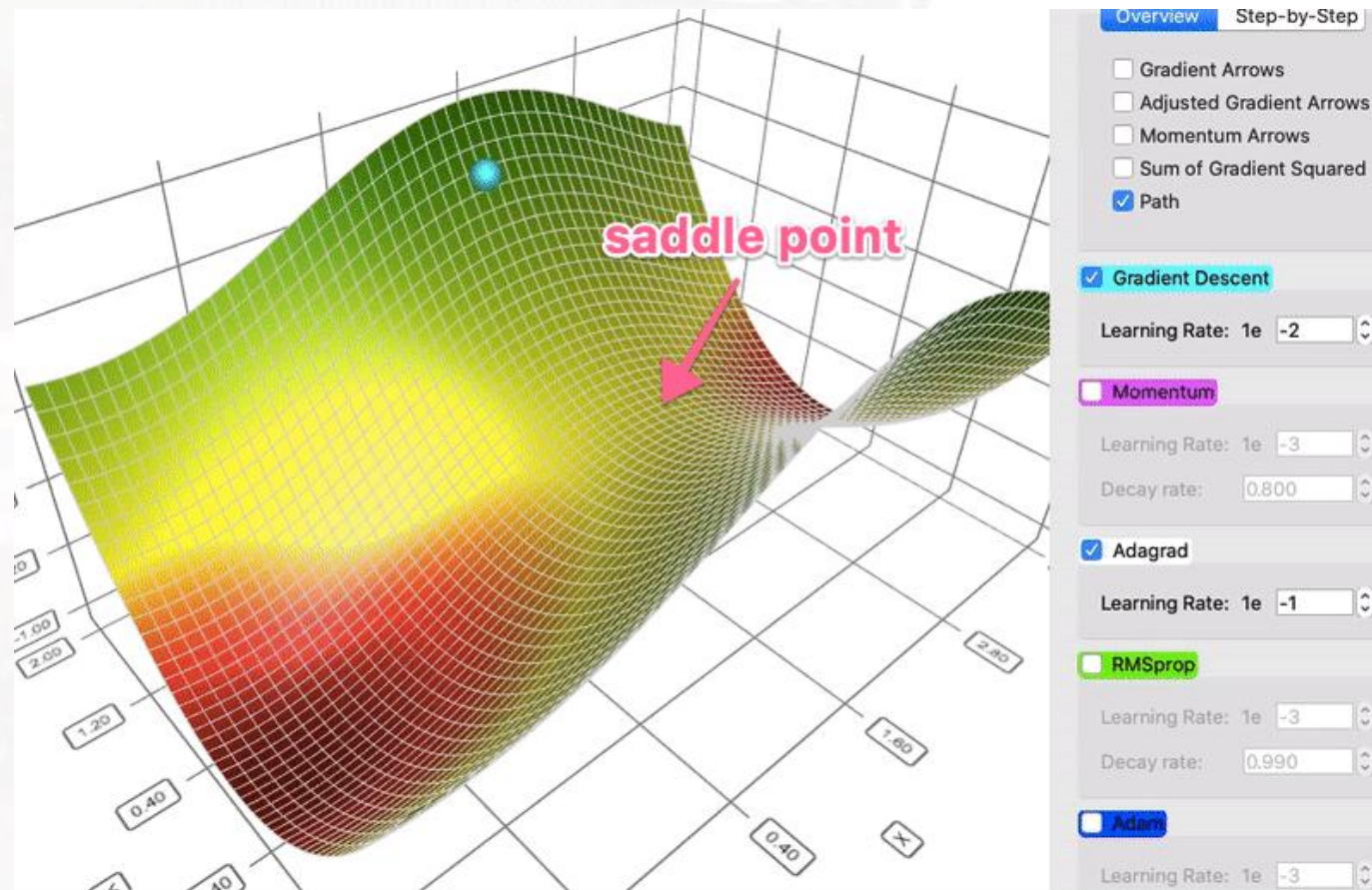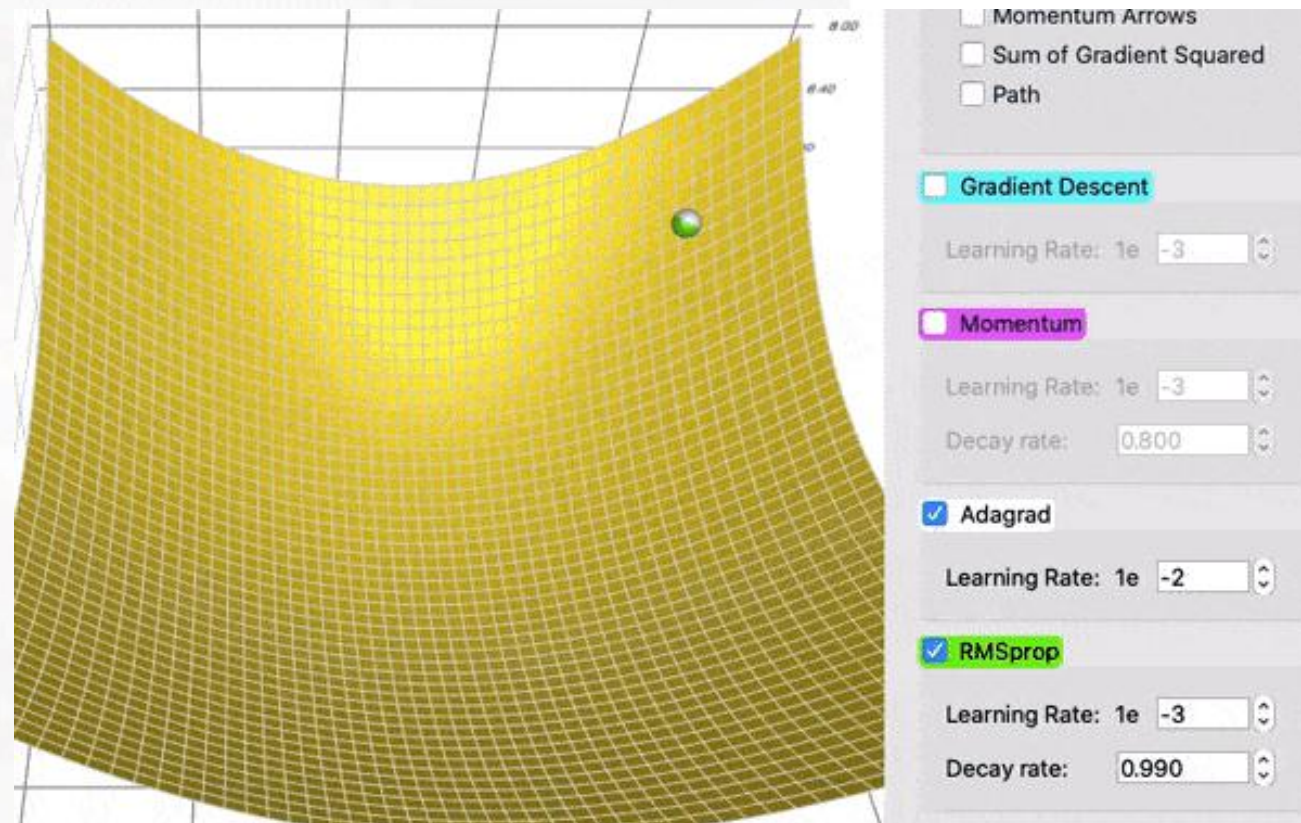$$E[g^2]_t = 0.9 E[g^2]_{t-1} + 0.1 g_t^2$$

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{E[g^2]_t + \epsilon}} g_t$$

# Adam(Adaptive Moment Estimation)
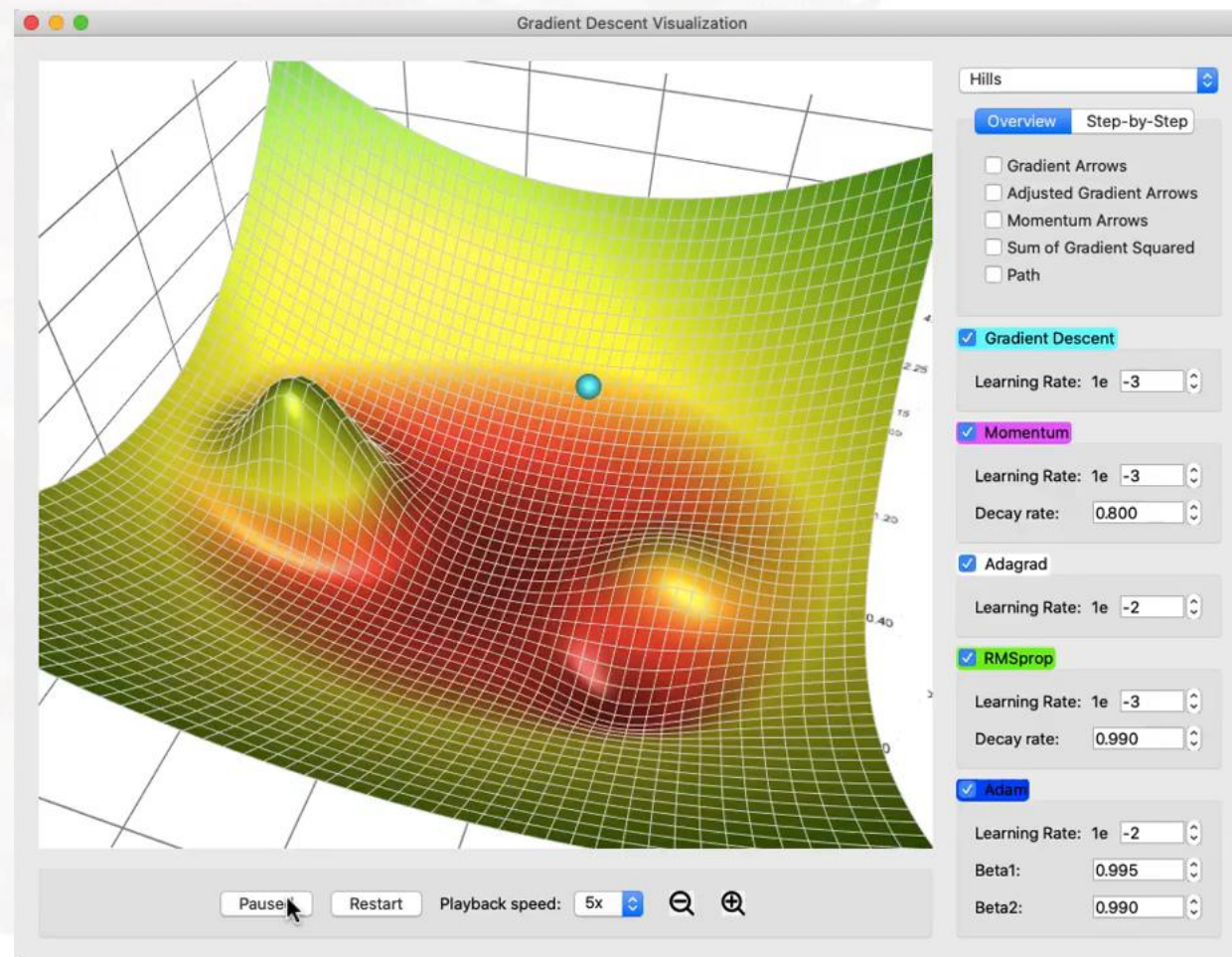
- **Combine SGDM and RMSProp**

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{\hat{v}_t} + \epsilon}\hat{m}_t.$$

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1)g_t$$

$$v_t = \beta_2 v_{t-1} + (1 - \beta_2)g_t^2$$

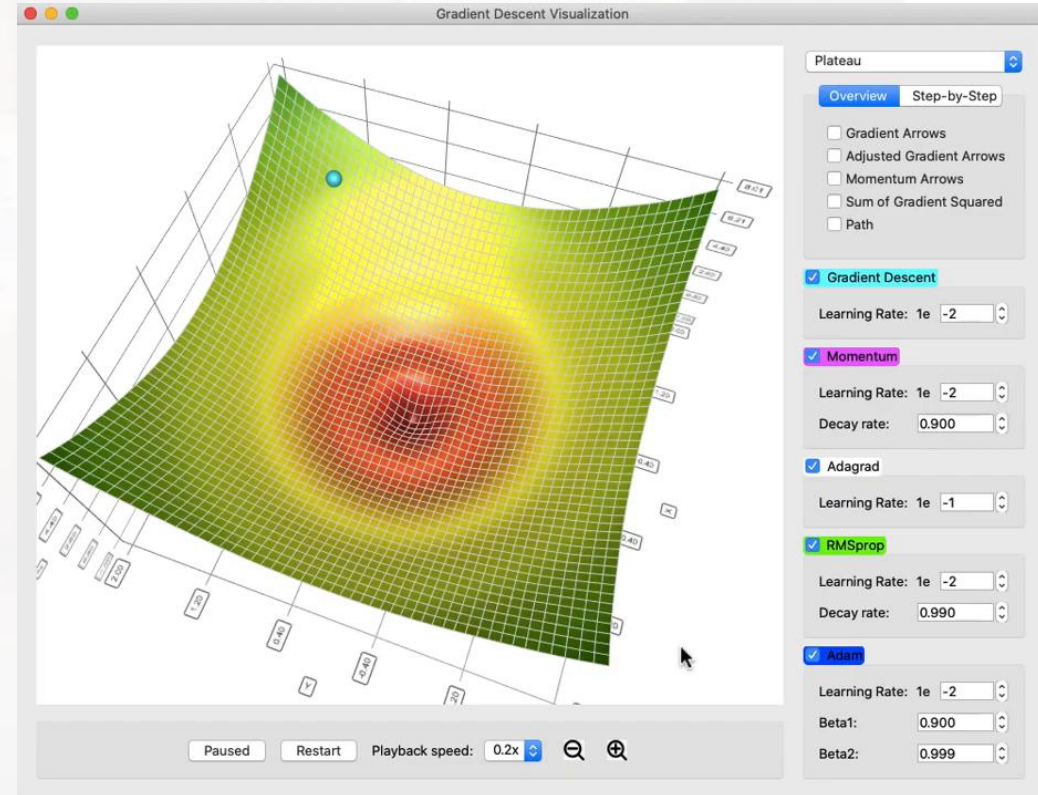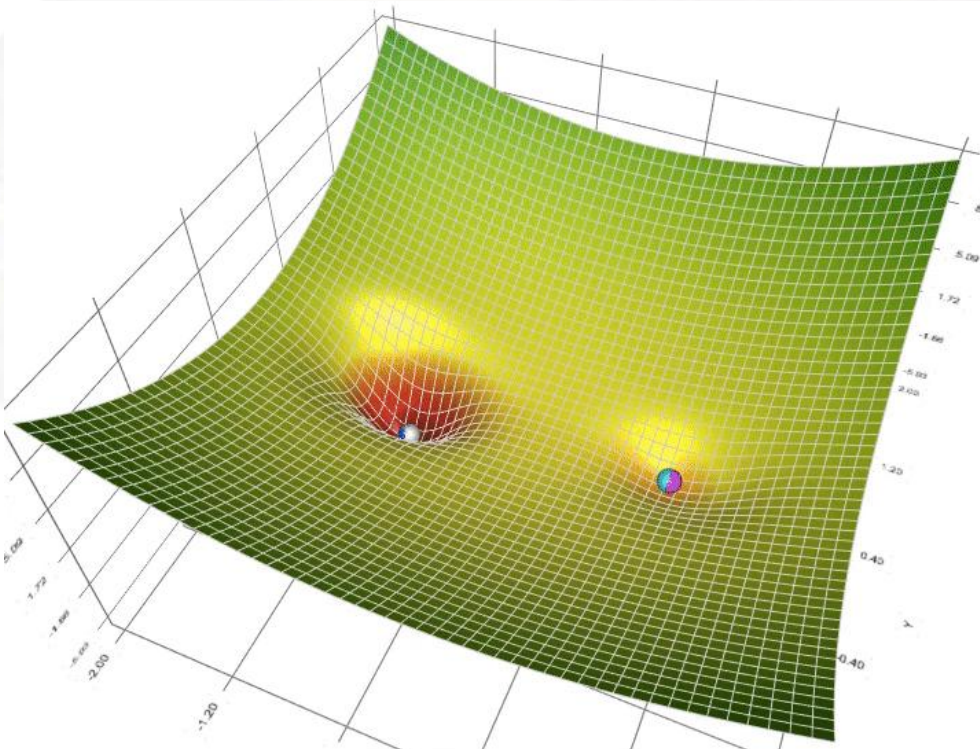$$v_t = \gamma v_{t-1} + \eta\nabla_\theta J(\theta)$$

$$E[g^2]_t = 0.9E[g^2]_{t-1} + 0.1g_t^2$$

# The Advantage of Adam

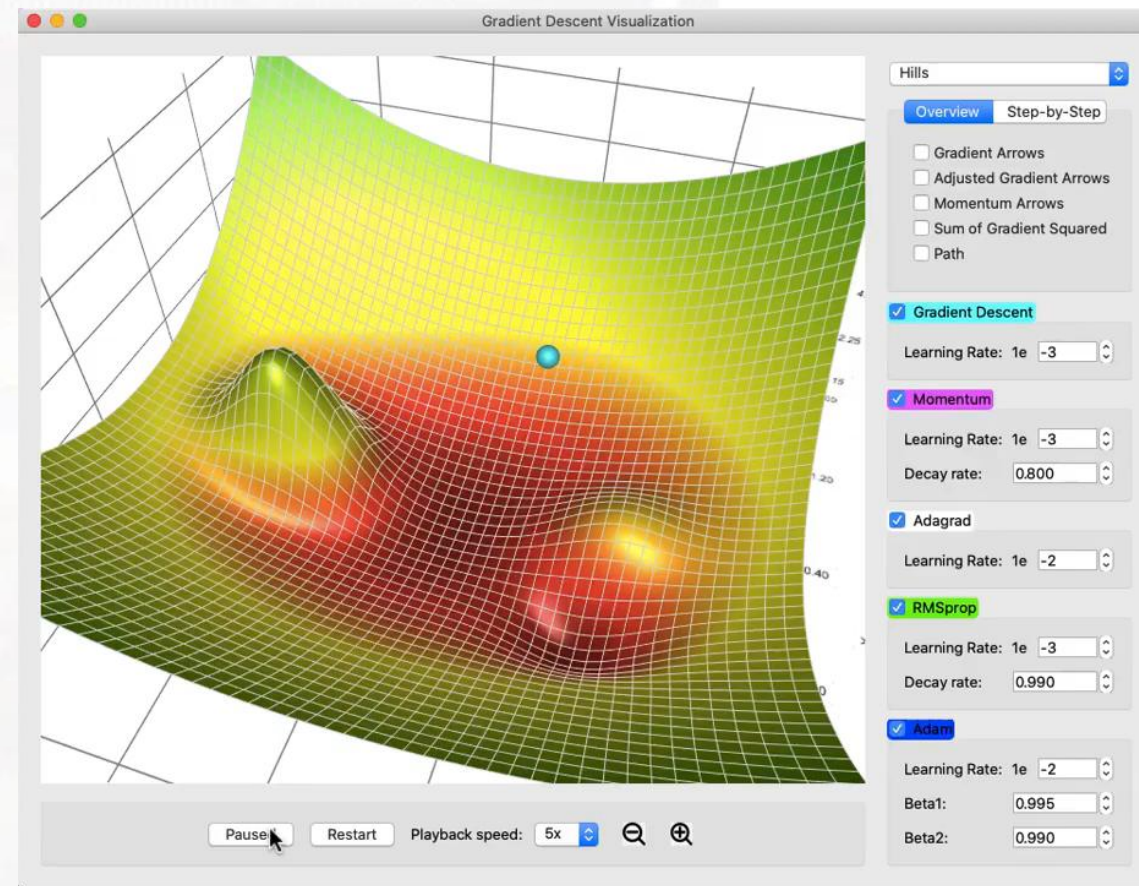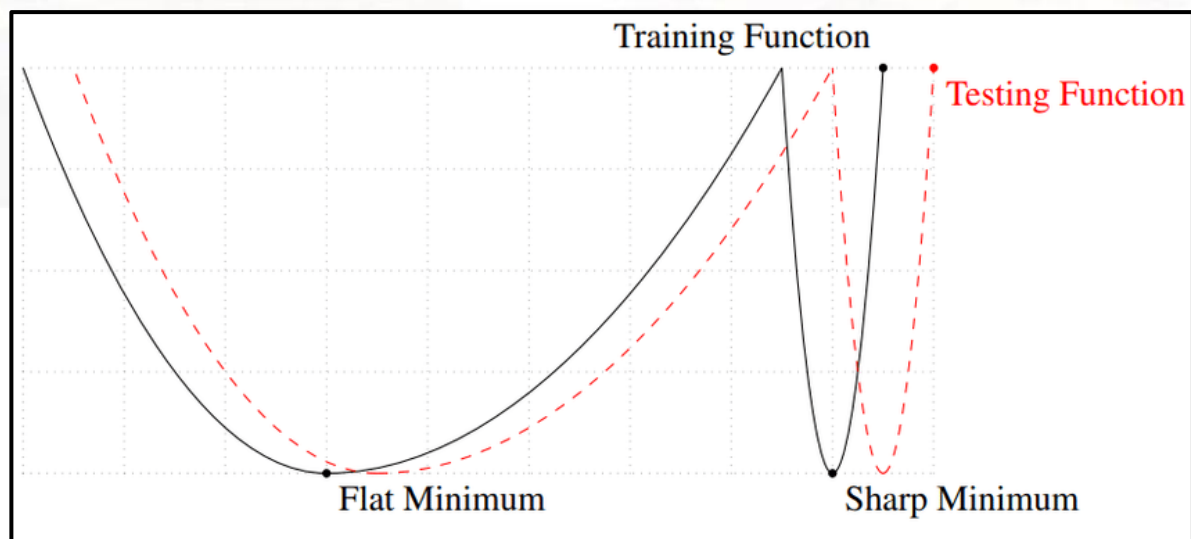- **The advantage of Adaptive**
- **The advantage of Momentum**

# The Problem of Adam

## ■ Overfitting



(a) Training Accuracy

(b) Test Accuracy



Training Function

Testing Function

Flat Minimum

Sharp Minimum

# Summary

■ **Adaptive vs Non Adaptive**

| Adaptive Method | Non-Adaptive Method |
|---|---|
| Adam, AdaGrad, RMSProp | SGDM, SGD |
| Difficult data, complex networks, hard to converge | Good initialization and learning rate scheduling scheme |

| | SGDM | Adam |
|---|---|---|
| Training Speed | Slow | **Fast** |
| Convergence | **Good** | Poor |
| Stability | **Good** | Poor |
| Generalization | **Good** | Poor |

If you are interested in visualizing these or other optimization algorithms, refer to this useful tutorial.

# Fit Everything?

■ **We can have good approximation with sufficient pieces.**



The universal approximation brought by nonlinearity

# Neural Network

$y$   $x_1$

红色曲线 ＝ 常数 ＋ 一组 ⌐ 的和

$y$   ❶ ❸ ❷   $x_1$

How to represent this function?
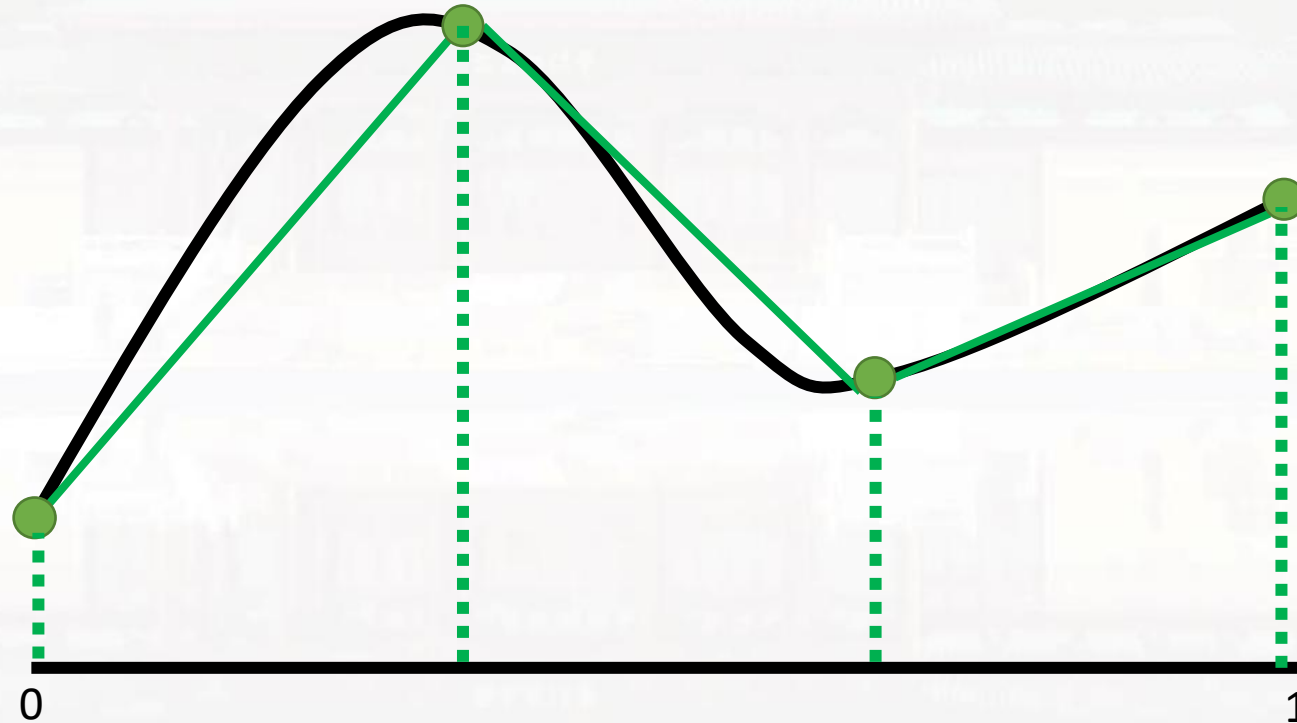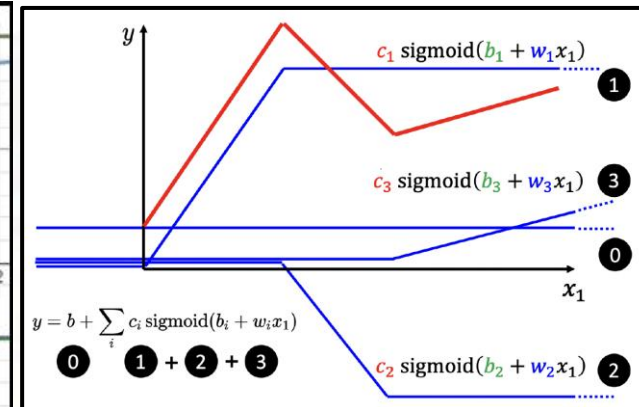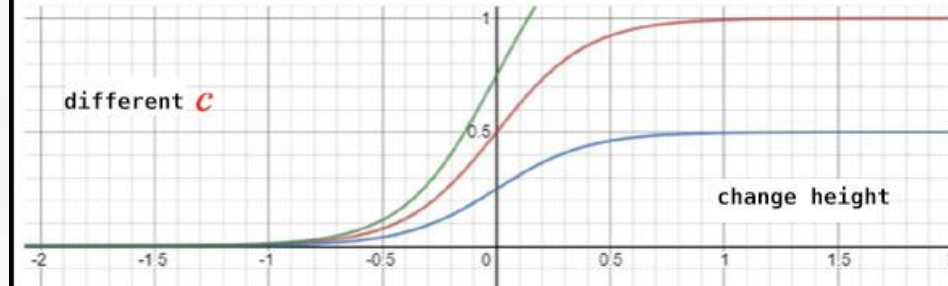
**Hard Sigmoid**

$x_1$

**Sigmoid Function**

$$y = c\,\frac{1}{1 + e^{-(b+wx_1)}}$$

$$= c\,sigmoid(b + wx_1)$$

$x_1$

different $w$    change slope

different $b$    shifting

different $c$    change height

$y$

$c_1\,sigmoid(b_1 + w_1x_1)$   ❶

$c_3\,sigmoid(b_3 + w_3x_1)$   ❸

❶

$y = b + \sum_i c_i\,sigmoid(b_i + w_i x_1)$

❶   ❶ + ❷ + ❸   $c_2\,sigmoid(b_2 + w_2x_1)$ ❷

$x_1$

$$y = b + wx_1$$

$$y = b + \sum_i c_i\,sigmoid(b_i + w_i x_1)$$

$$y = b + \sum_j w_j x_j$$

$$y = b + \sum_i c_i\,sigmoid\left(b_i + \sum_j w_{ij}x_j\right)$$

# Neural Network

How to
represent this
function?



piecewise
linear



= constant +
    sum of a set of

Rectified
Linear Unit
(ReLU)

$c\ max(0, b + wx_1)$

$x_1$

$c'\ max(0, b' + w'x_1)$

0                                                                          1

## The universal approximation brought by nonlinearity

# Neural Network

■ **Why Use Layers?**

# Neural Network

■ **we piece together various components like loops and lines**

# Neural Network

- **In a perfect world, we might hope that each neuron in the second-to-last layer corresponds to one of these subcomponents.**

# Neural Network

- **Layers Break Problems Into Bite-Sized Pieces**
  - ☐ **Edge detection is a useful step for all kinds of image-recognition problems.**
  - ☐ **beyond image recognition**



recognition → re·cog·ni·tion → recognition

Raw audio

# Neural Network

■ **Rank the four images (A, B, C, and D) based on how much they would activate that neuron:**

# Neural Network

■ **We hope … but**

# Overfitting and Underfitting

- **The balance between data (knowns or constraints) and parameters (unknowns)**

# Overfitting: Data Augmentation

# Overfitting: Regularization



Poly Degree=9, $\lambda$=0.0    Poly Degree=9, $\lambda$=1e-5    Poly Degree=9, $\lambda$=0.1

$$J(\beta_0, \beta_1) = \frac{1}{2m} \sum_{i=1}^{m} \left( \left( \beta_0 + \beta_1 x_{obs}^{(i)} \right) - y_{obs}^{(i)} \right)^2 + \lambda \sum_{j=1}^{k} \beta_j^2$$

# Overfitting



y = a + bx + cx²

Constrained model

Real data distribution(invisible)

● training data

● testing data

**Model Constraint**

Overconstrained

y = a + bx

Real data distribution(invisible)

Bias problem

**Early Stop** — Overfitting

Loss

validation

training

early stopping — Epochs

**Dropout**

(a) Standard Neural Net

(b) After applying dropout.

■ **Solutions to overfitting**
  ■ Data augmentation
  ■ Model constraints：
    regularization,
    architecture, dropout
  ■ Early stop

# Architectures

# The Problem of MLP

■ **Numerous Parameters**
  ☐ **Poor Explanation due to Large Scale Minimum**
  ☐ **Large Scale Minimum due to Large Scale Parameters**

# The Problem of MLP

## ■ Poor Flexibility



784



Neurons don't need message of the whole image.

Input    1st layer    2nd layer

$x_1$

$x_2$

$x_N$

Basic detector    Advanced detector

Some patterns are much more smaller than the whole picture

# The Problem of MLP

## ■ Poor Flexibility

# From MLP to CNN

# Pooling and Stride



滤波器 1　滤波器 2

stride=1

| 1 | 0 | 0 | 0 | 0 | 1 |
|---|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 1 | 0 | 0 |
| 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 |

Dot product →　3　-1

If stride=2

| 1 | 0 | 0 | 0 | 0 | 1 |
|---|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 1 | 0 | 0 |
| 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 |

3　-3

6 x 6 image

Strided
Convolution

Input Image　　Kernel　　Feature Map

# Increase Receptive Field



Receptive Field in Convolutional Networks

Layer 1    Layer 2    Layer 3

output

hidden layer 3

hidden layer 2

hidden layer 1

input

# Convolutional Neural Networks

# Feedforward Neural Networks

- **A feedforward neural network (FNN) is one of the two broad types of artificial neural network, characterized by direction of the flow of information between its layers**

- **Human brain deals with information streams. Most data is obtained, processed, and generated sequentially.**

# The Problem of FNNs

- **Human thoughts have persistence; humans don't start their thinking from scratch every second.**
- **The applications of standard FNNs are limited due to:**
  - They only accepted a fixed-size vector as input (e.g., an image) and produce a fixed-size vector as output
  - These models use a fixed amount of computational steps

# Recurrent Neural Networks

Output Layer **y**

A

Hidden Layers **h** C

B

Input Layer **x**

A, B and C are the parameters

| | Next word | long | and | thanks | for | all | |
|---|---|---|---|---|---|---|---|

Loss $-\log y_{\text{long}}$ $-\log y_{\text{and}}$ $-\log y_{\text{thanks}}$ $-\log y_{\text{for}}$ $-\log y_{\text{all}}$ $\cdots$ $\frac{1}{T}\sum_{t=1}^{T} L_{CE}$

$y$

Softmax over Vocabulary

$Vh$

RNN $h$ $\cdots$

Input Embeddings $e$

So long and thanks for $\cdots$

# Training RNN

Backpropagation through time

Forward through entire sequence to compute loss, then backward through entire sequence to compute gradient

Loss



Loss



$y$

$V$

$h$

$W$

$U$

$x$

# Gradient Vanishing and Exploding



$$=w^{999}$$

| | | |
|---|---|---|
| $w = 1$ | → | $y^{1000} = 1$ |
| $w = 1.01$ | → | $y^{1000} \approx 20000$ |
| $w = 0.99$ | → | $y^{1000} \approx 0$ |
| $w = 0.01$ | → | $y^{1000} \approx 0$ |

Large $\partial L/\partial w$ → Small Learning rate?

small $\partial L/\partial w$ → Large Learning rate?

Real experiments on Language modeling

Total Loss

sometimes

Lucky

Epoch

The error surface is either very flat or very steep.

Total Loss

Clipping

$w_2$

$w_1$

[Razvan Pascanu, ICML' 13]

# Gradient Vanishing

## Vanilla RNN Gradient Flow

Bengio et al, "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994
Pascanu et al, "On the difficulty of training recurrent neural networks", ICML 2013

Backpropagation from $h_t$ to $h_{t-1}$ multiplies by W (actually $W_{hh}^T$)



$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$

$$= \tanh\left( \begin{pmatrix} W_{hh} & W_{hx} \end{pmatrix} \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} \right)$$

$$= \tanh\left( W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} \right)$$

$$\frac{\partial h_t}{\partial h_{t-1}} = tanh'(W_{hh}h_{t-1} + W_{xh}x_t)W_{hh}$$

# Gradient Vanishing

## Vanilla RNN Gradient Flow

Bengio et al, "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994
Pascanu et al, "On the difficulty of training recurrent neural networks", ICML 2013

Gradients over multiple time steps:



$$\frac{\partial L}{\partial W} = \sum_{t=1}^{T} \frac{\partial L_t}{\partial W} \qquad \boxed{\frac{\partial h_t}{\partial h_{t-1}} = tanh'(W_{hh}h_{t-1} + W_{xh}x_t)W_{hh}}$$

$$\frac{\partial L_T}{\partial W} = \frac{\partial L_T}{\partial h_T}\frac{\partial h_t}{\partial h_{t-1}} \cdots \frac{\partial h_1}{\partial W} = \frac{\partial L_T}{\partial h_T}\left(\prod_{t=2}^{T}\boxed{\frac{\partial h_t}{\partial h_{t-1}}}\right)\frac{\partial h_1}{\partial W}$$

# Gradient Vanishing

## Vanilla RNN Gradient Flow

Bengio et al, "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994
Pascanu et al, "On the difficulty of training recurrent neural networks", ICML 2013

Gradients over multiple time steps:



What if we assumed no non-linearity?

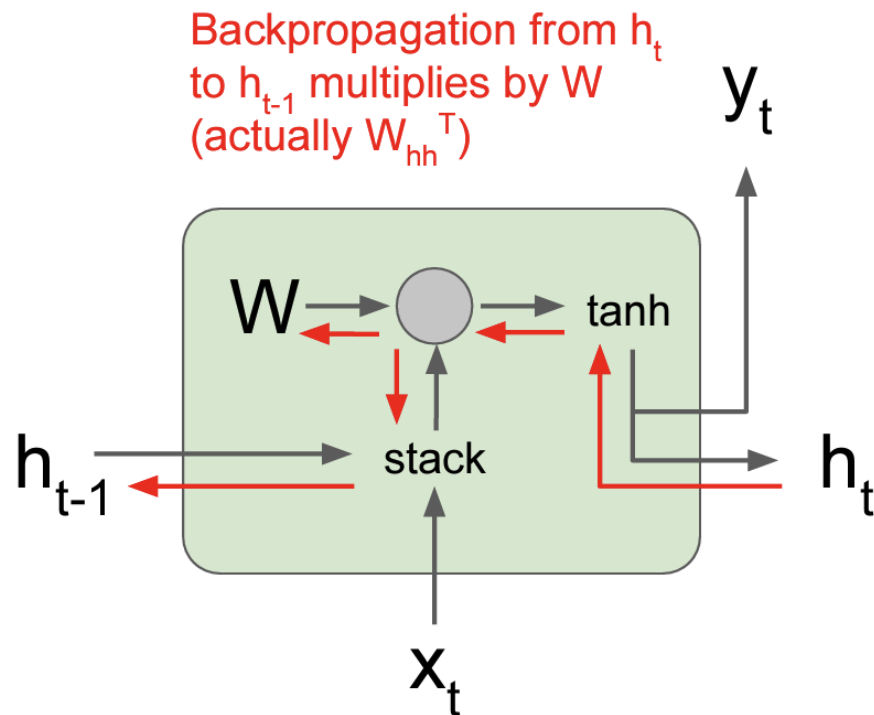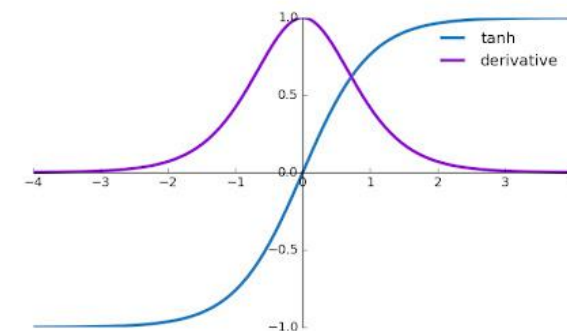$$\frac{\partial L}{\partial W} = \sum_{t=1}^{T} \frac{\partial L_t}{\partial W}$$

$$\frac{\partial L_T}{\partial W} = \frac{\partial L_T}{\partial h_T} \boxed{W_{hh}^{T-1}} \frac{\partial h_1}{\partial W}$$

Largest singular value > 1:
**Exploding gradients**

Largest singular value < 1:
**Vanishing gradients**

→ **Gradient clipping**:
Scale gradient if its norm is too big

```
grad_norm = np.sum(grad * grad)
if grad_norm > threshold:
    grad *= (threshold / grad_norm)
```

# RNN and LSTM



**Vanilla RNN**

$$h_t = \tanh\left(W\begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}\right)$$

**LSTM**

Four gates

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

Cell state $\longrightarrow$ $c_t = f \odot c_{t-1} + i \odot g$

Hidden state $\longrightarrow$ $h_t = o \odot \tanh(c_t)$

vector from below (**x**)

vector from before (**h**)

W

x

h

sigmoid $\longrightarrow$ i

sigmoid $\longrightarrow$ f

sigmoid $\longrightarrow$ o

tanh $\longrightarrow$ g

4h x 2h      4h      4*h

Output Gate

Memory Gate

Forget Gate

Input Gate     **LSTM**

**i**: Input gate, whether to write to cell

**f**: Forget gate, Whether to erase cell

**o**: Output gate, How much to reveal cell

**g**: Gate gate (?), How much to write to cell

# RNN and LSTM



Other part of the network

Special Neuron:
4 inputs,
1 output

Signal control the output gate

(Other part of the network)

Output Gate

Memory Cell

Forget Gate

Signal control the forget gate

(Other part of the network)

Input Gate

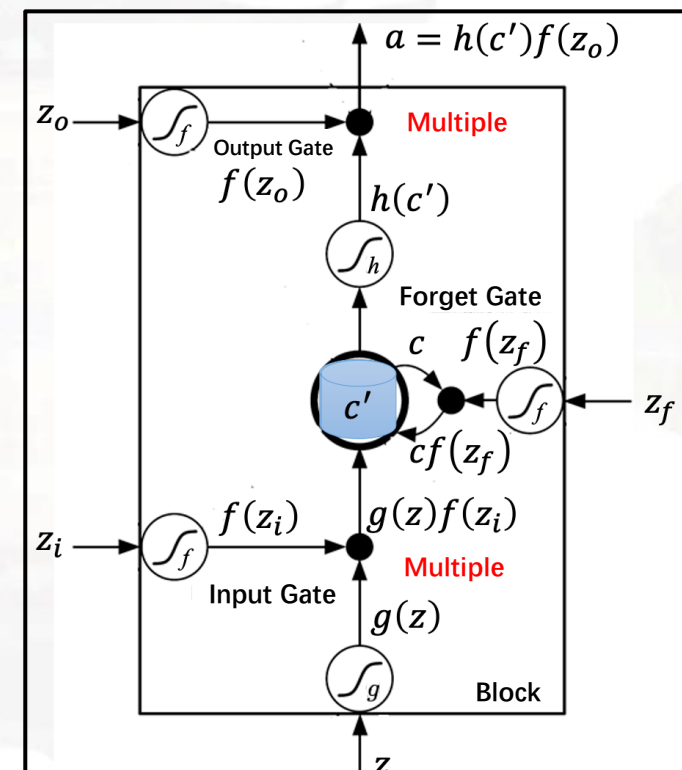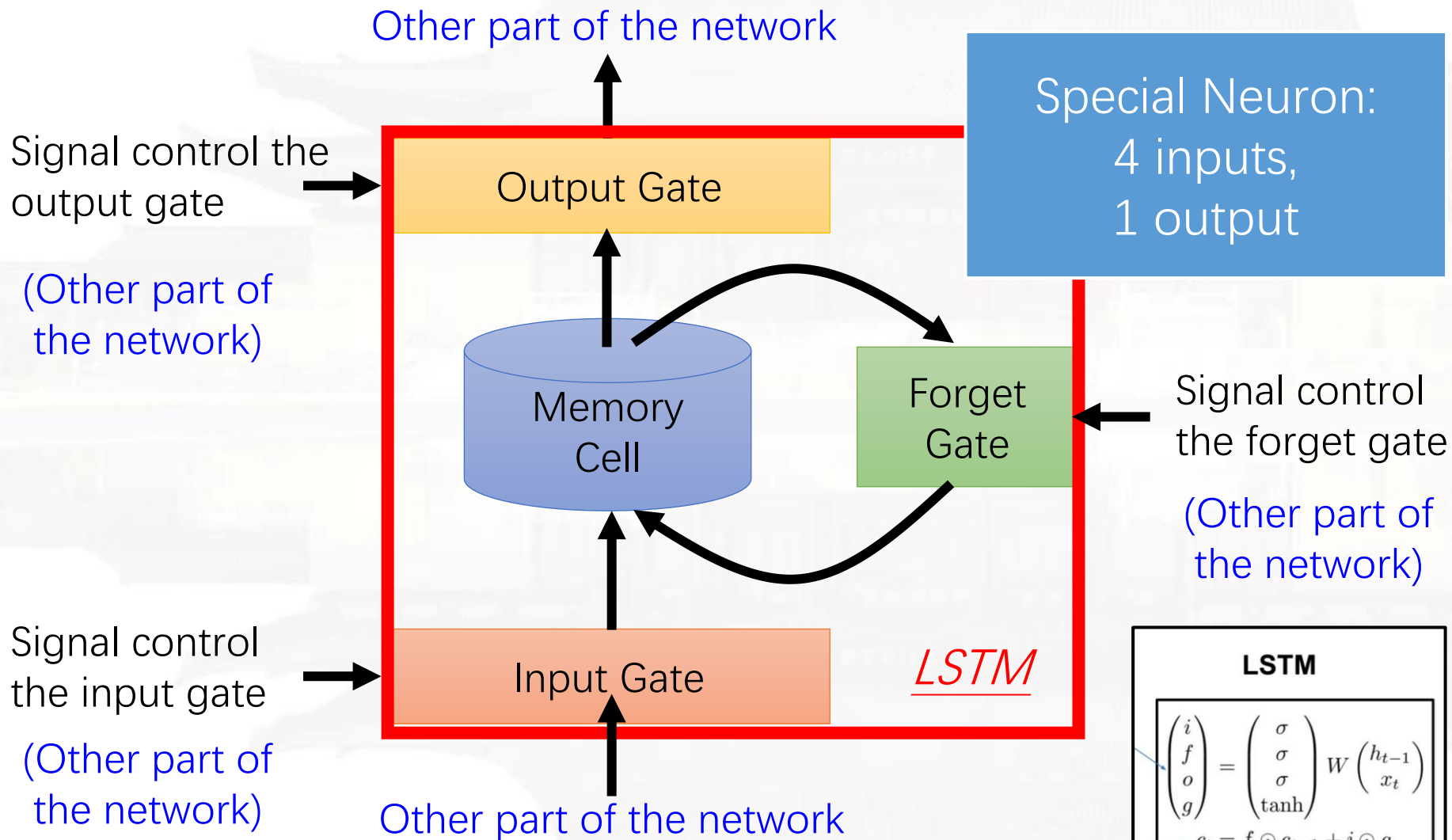*LSTM*

Signal control the input gate

(Other part of the network)

Other part of the network

**LSTM**

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

$c_t = f \odot c_{t-1} + i \odot g$
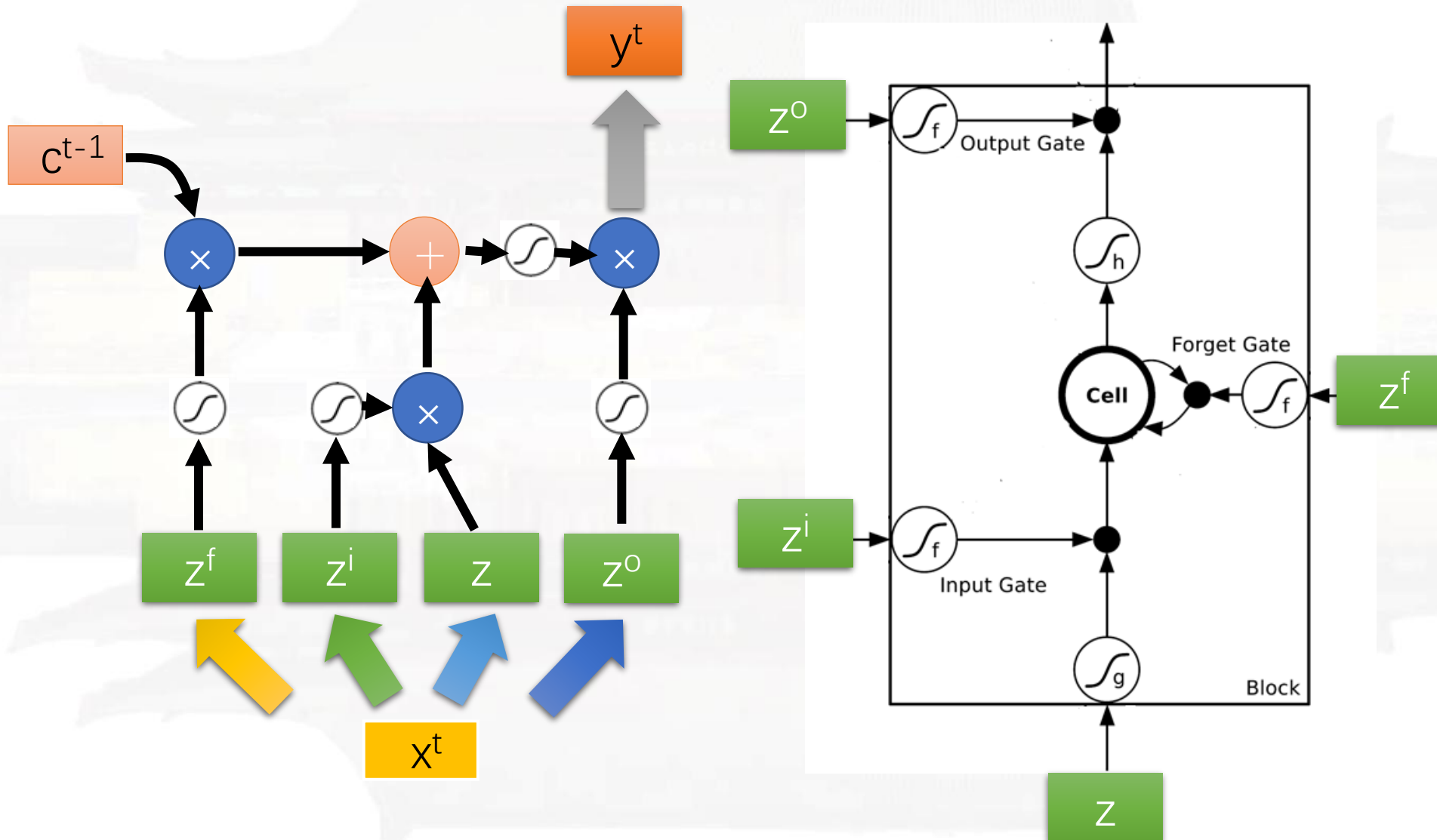
$h_t = o \odot \tanh(c_t)$

$a = h(c')f(z_o)$

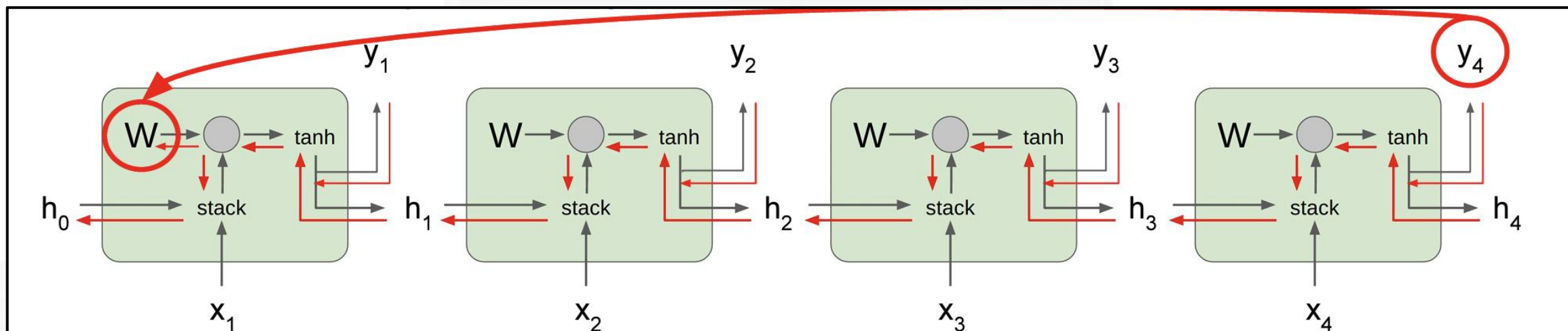$z_o$   $f(z_o)$   Output Gate   Multiple

$h(c')$

Forget Gate

$c$   $f(z_f)$

$c'$   $z_f$

$cf(z_f)$

$f(z_i)$   $g(z)f(z_i)$

$z_i$   Input Gate   Multiple

$g(z)$

$z$   Block
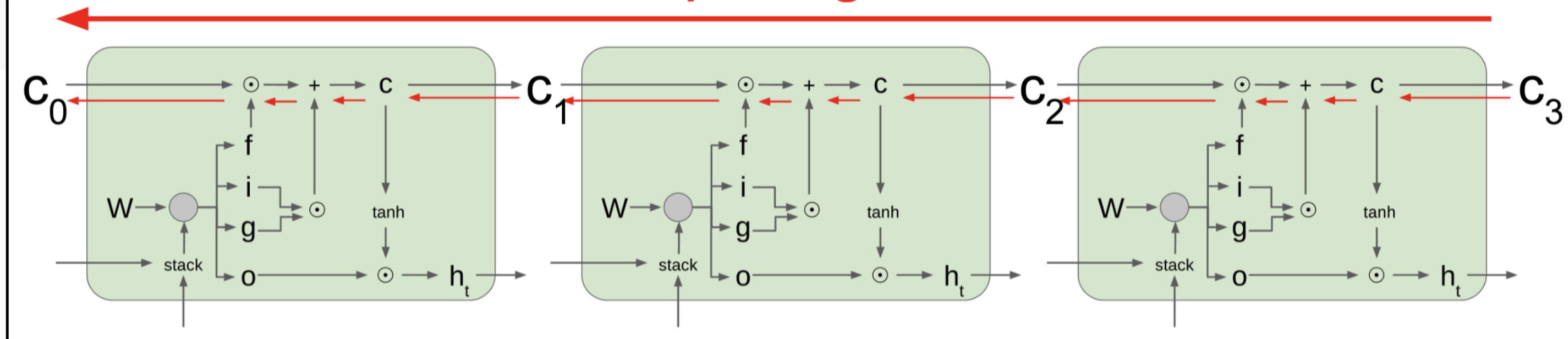
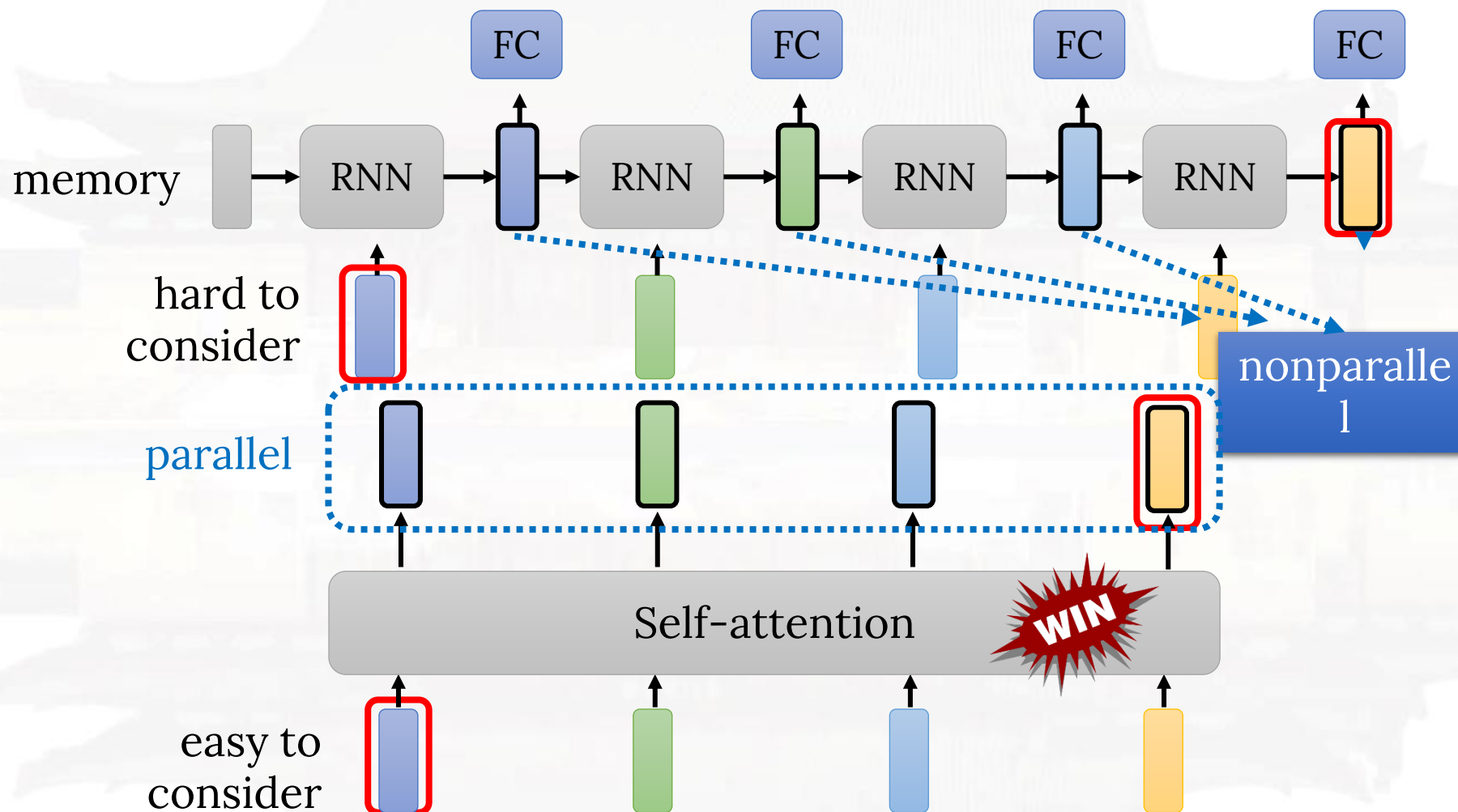# RNN and LSTM



Uninterrupted gradient flow!

# From RNN to Attention



Transformers are RNNs: Fast Autoregressive Transformers with Linear Attention

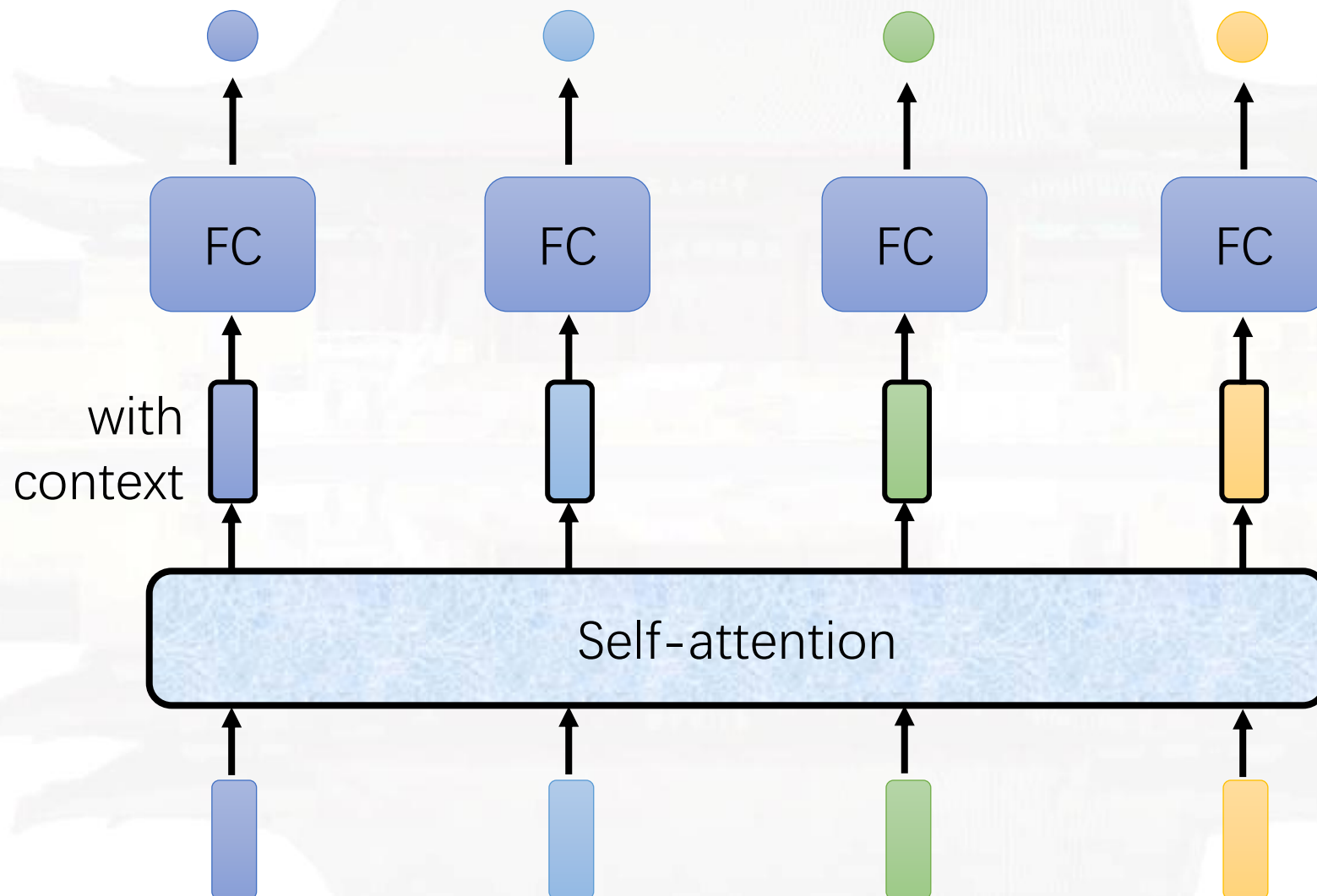https://arxiv.org/abs/2006.16236

# From RNN to Attention

$b^i$ is obtained based on the whole input sequence.

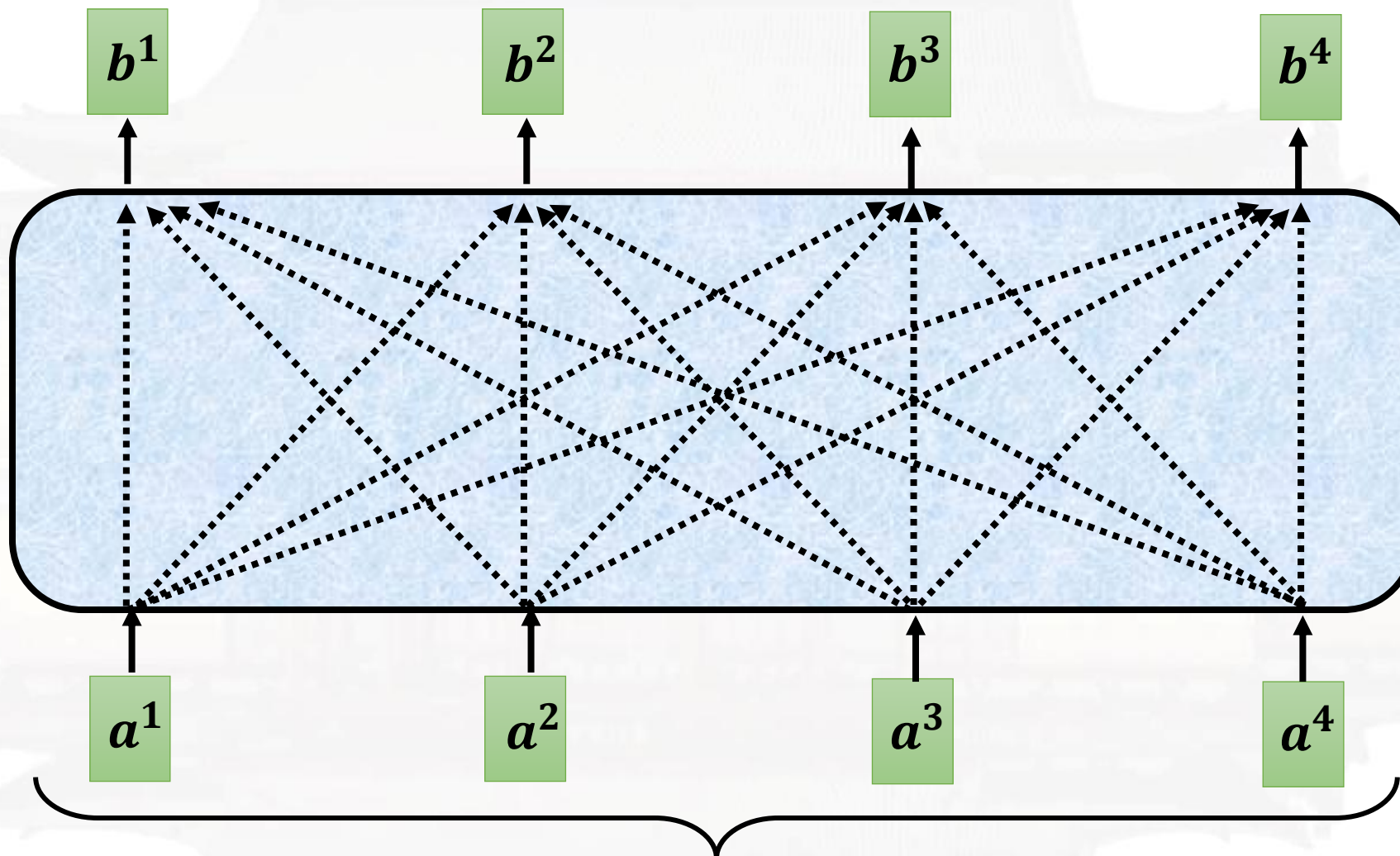$b^1, b^2, b^3, b^4$ can be parallelly computed.



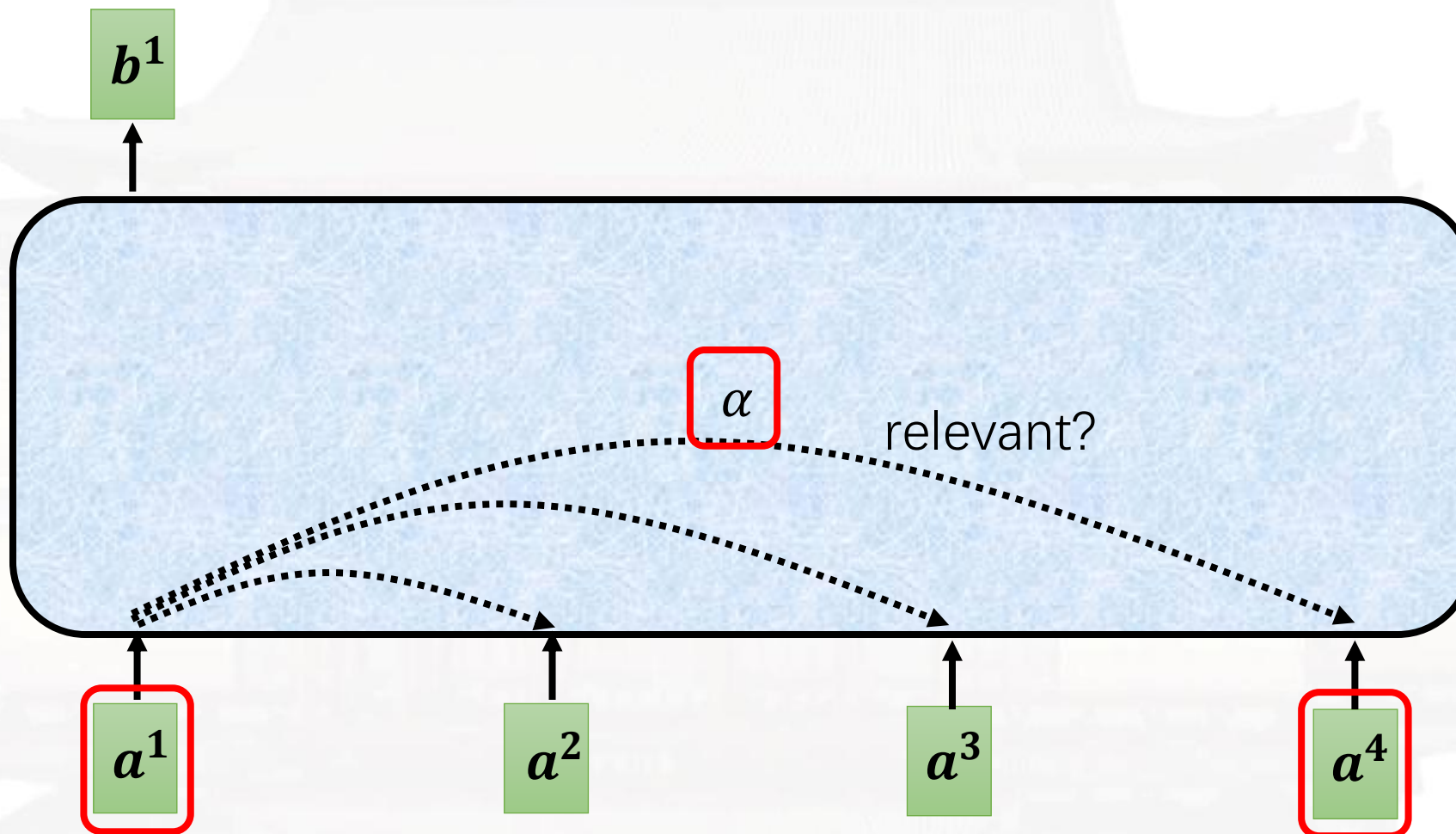You can try to replace any thing that has been done by RNN with self-attention.

# Self-Attention

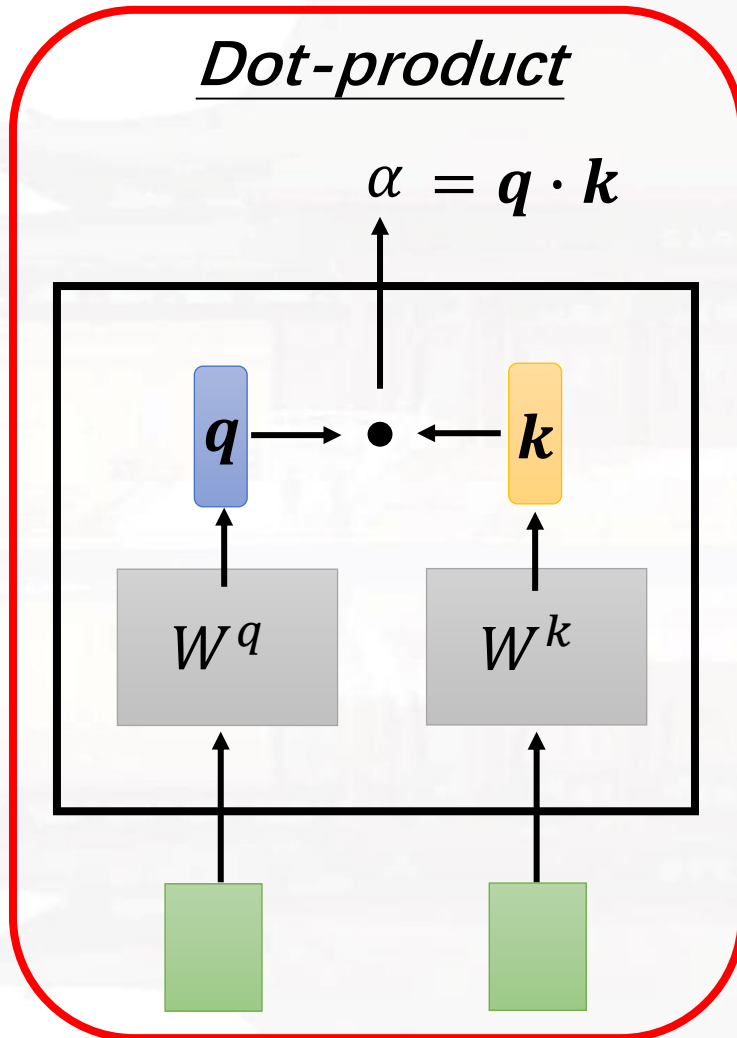# Self-Attention



Can be either **input** or **a hidden layer**

Find the relevant vectors in a sequence

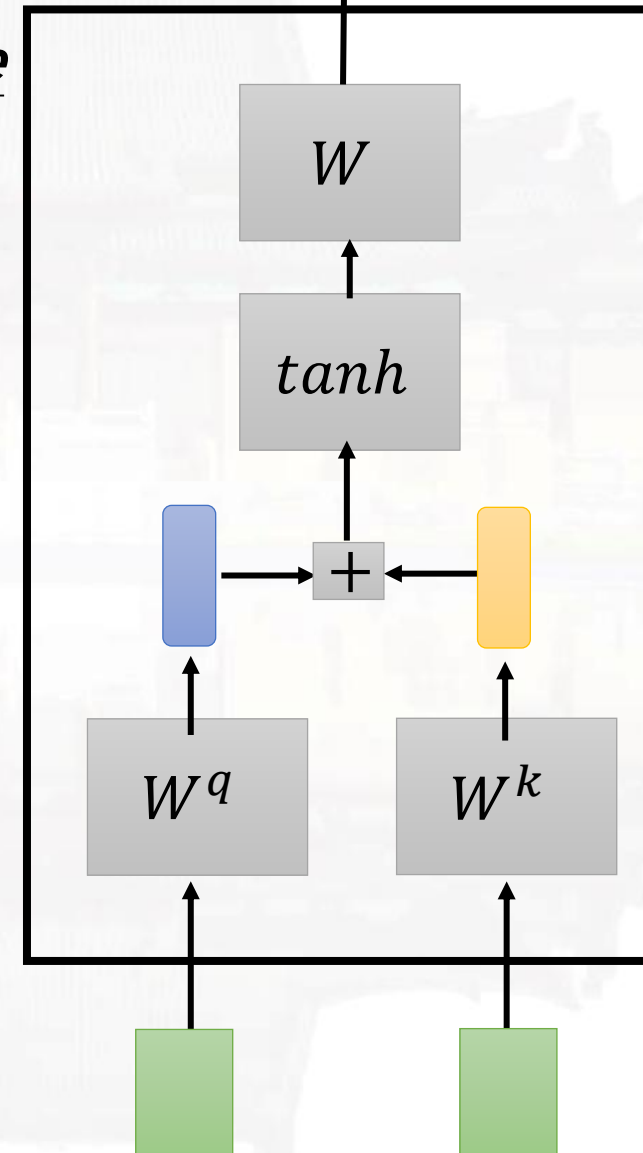# Self-Attention



**Dot-product**

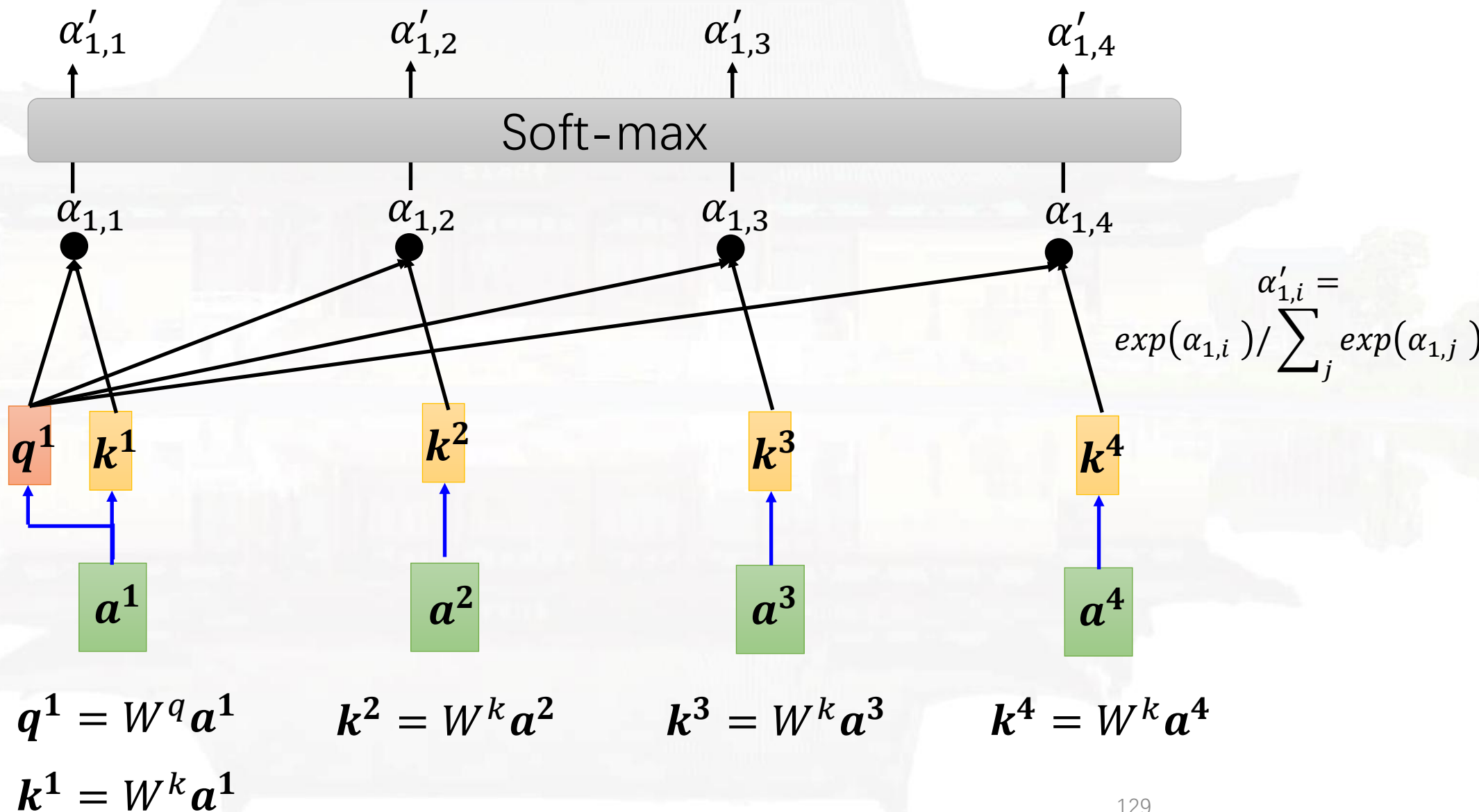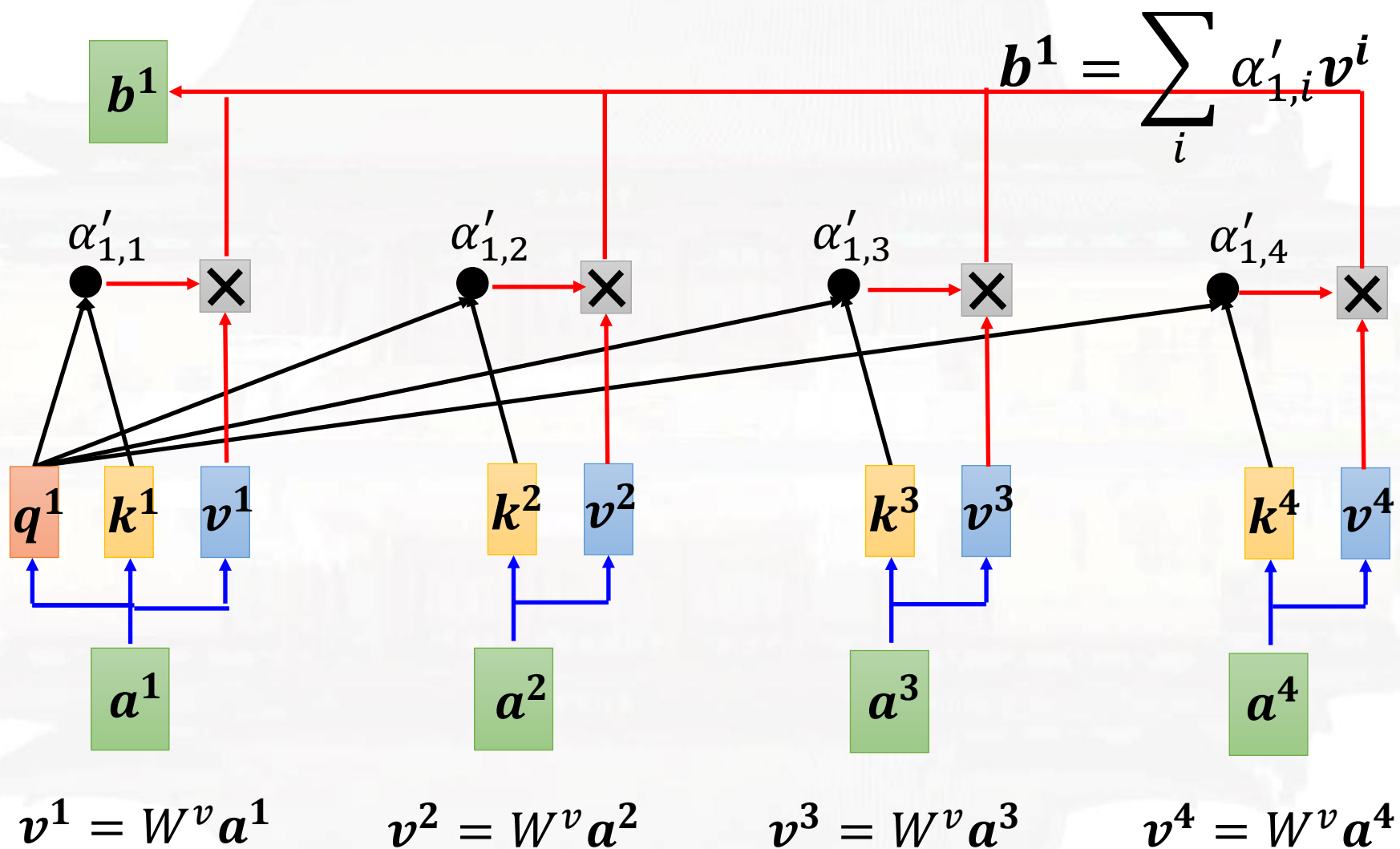$$\alpha = \boldsymbol{q} \cdot \boldsymbol{k}$$

$\boldsymbol{q} \rightarrow \bullet \leftarrow \boldsymbol{k}$

$W^q \qquad W^k$

**Additive**

$\alpha$

$W$

$tanh$

$+$

$W^q \qquad W^k$

# Self-Attention



$$\alpha_{1,2} = \boldsymbol{q^1} \cdot \boldsymbol{k^2} \qquad \alpha_{1,3} = \boldsymbol{q^1} \cdot \boldsymbol{k^3} \qquad \alpha_{1,4} = \boldsymbol{q^1} \cdot \boldsymbol{k^4}$$

$$\alpha_{1,2} \qquad \alpha_{1,3} \qquad \alpha_{1,4}$$

attention
score

$\boldsymbol{q^1}$ query    $\boldsymbol{k^2}$ key    $\boldsymbol{k^3}$    $\boldsymbol{k^4}$

$\boldsymbol{a^1}$    $\boldsymbol{a^2}$    $\boldsymbol{a^3}$    $\boldsymbol{a^4}$

$$\boldsymbol{q^1} = W^q \boldsymbol{a^1} \qquad \boldsymbol{k^2} = W^k \boldsymbol{a^2} \qquad \boldsymbol{k^3} = W^k \boldsymbol{a^3} \qquad \boldsymbol{k^4} = W^k \boldsymbol{a^4}$$

# Self-Attention



$$\alpha'_{1,i} = exp(\alpha_{1,i})/\sum_j exp(\alpha_{1,j})$$

$$q^1 = W^q a^1$$

$$k^2 = W^k a^2 \qquad k^3 = W^k a^3 \qquad k^4 = W^k a^4$$

$$k^1 = W^k a^1$$

# Self-Attention



$$b^1 = \sum_i \alpha'_{1,i} v^i$$

$$v^1 = W^v a^1 \qquad v^2 = W^v a^2 \qquad v^3 = W^v a^3 \qquad v^4 = W^v a^4$$

# Self-Attention



$$b^2 = \sum_i \alpha'_{2,i} v^i$$

# Self-Attention

# Self-Attention

$$q^i = W^q a^i$$

$$\boxed{q^1\ q^2\ q^3\ q^4} = \boxed{W^q}\ \boxed{a^1\ a^2\ a^3\ a^4}$$

$Q$ $\qquad$ I

$$k^i = W^k a^i$$

$$\boxed{k^1\ k^2\ k^3\ k^4} = \boxed{W^k}\ \boxed{a^1\ a^2\ a^3\ a^4}$$

$K$ $\qquad$ I

$$v^i = W^v a^i$$

$$\boxed{v^1\ v^2\ v^3\ v^4} = \boxed{W^v}\ \boxed{a^1\ a^2\ a^3\ a^4}$$

$V$ $\qquad$ I

$q^1\ k^1\ v^1 \qquad q^2\ k^2\ v^2 \qquad q^3\ k^3\ v^3 \qquad q^4\ k^4\ v^4$

$a^1 \qquad a^2 \qquad a^3 \qquad a^4$

# Self-Attention

$$\alpha_{1,1} = \boxed{k^1}\ \boxed{q^1} \quad \alpha_{1,2} = \boxed{k^2}\ \boxed{q^1}$$

$$\alpha_{1,3} = \boxed{k^3}\ \boxed{q^1} \quad \alpha_{1,4} = \boxed{k^4}\ \boxed{q^1}$$

$$\begin{bmatrix} \alpha_{1,1} \\ \alpha_{1,2} \\ \alpha_{1,3} \\ \alpha_{1,4} \end{bmatrix} = \begin{bmatrix} k^1 \\ k^2 \\ k^3 \\ k^4 \end{bmatrix} q^1$$
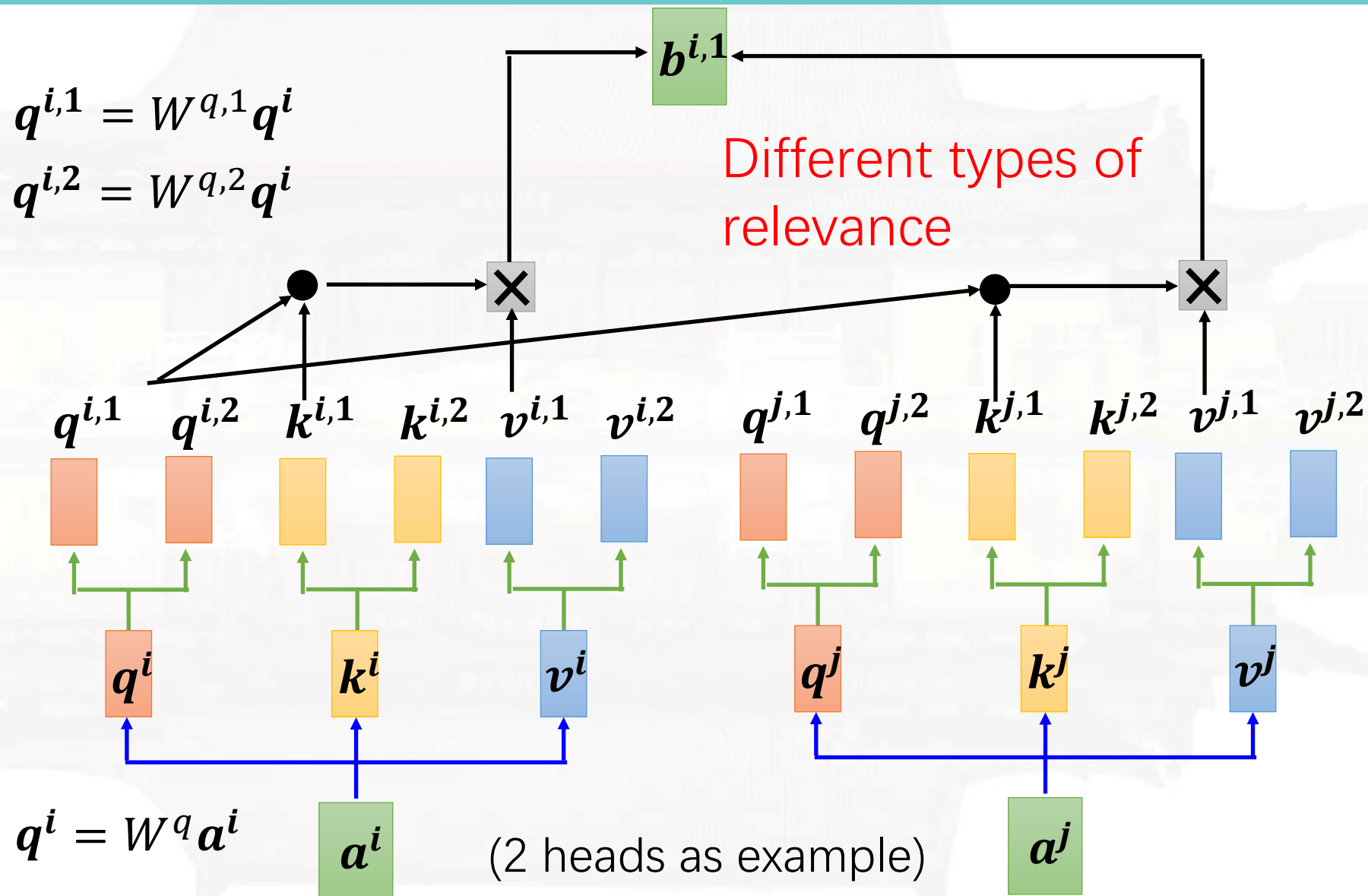
# Self-Attention

# Self-Attention

$$Q = W^q I$$

$$K = W^k I$$

$$V = W^v I$$

Parameters to be learned

$$A = K^T Q$$

$$A'$$

Attention Matrix

$$O = V A'$$

# Multi-head Self-attention

$$q^{i,1} = W^{q,1}q^i$$
$$q^{i,2} = W^{q,2}q^i$$

Different types of relevance

$$q^i = W^q a^i$$

(2 heads as example)

# Multi-head Self-attention



$$q^{i,1} = W^{q,1}q^i$$

$$q^{i,2} = W^{q,2}q^i$$

$b^{i,1}$

$b^{i,2}$ Different types of relevance

$q^{i,1}$ $q^{i,2}$ $k^{i,1}$ $k^{i,2}$ $v^{i,1}$ $v^{i,2}$ $q^{j,1}$ $q^{j,2}$ $k^{j,1}$ $k^{j,2}$ $v^{j,1}$ $v^{j,2}$

$q^i$ $k^i$ $v^i$ $q^j$ $k^j$ $v^j$

$$q^i = W^q a^i$$

$a^i$ (2 heads as example) $a^j$

# Multi-head Self-attention



(2 heads as example)

$q^i = W^q a^i$
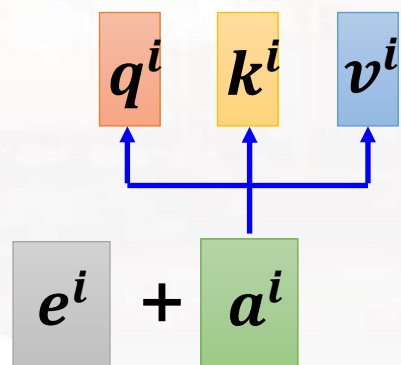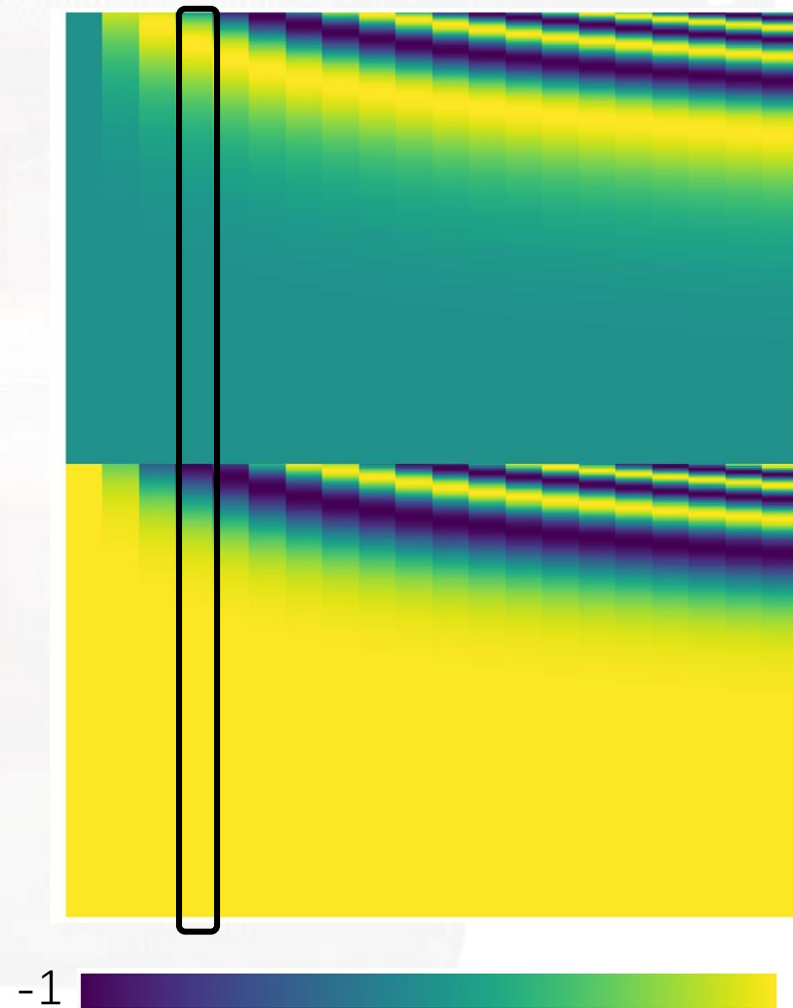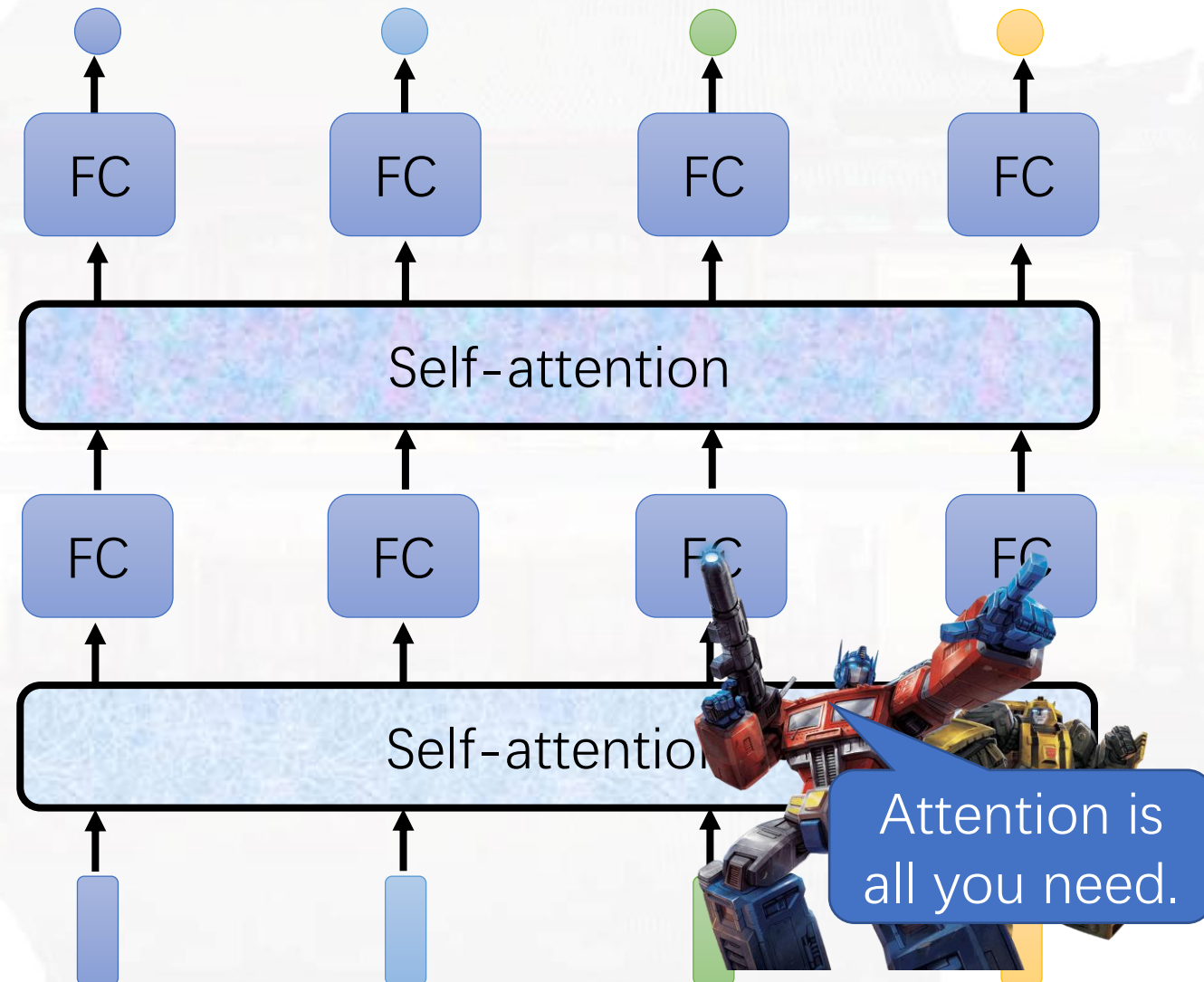
# Positional Encoding

- No position information in self-attention.

- Each position has a unique positional vector $e^i$
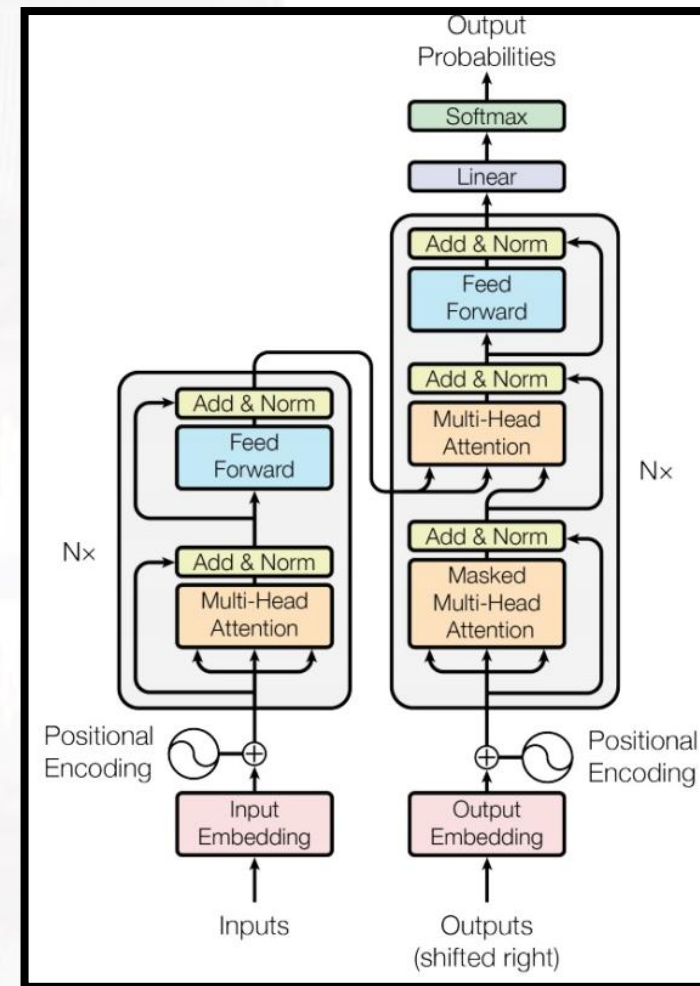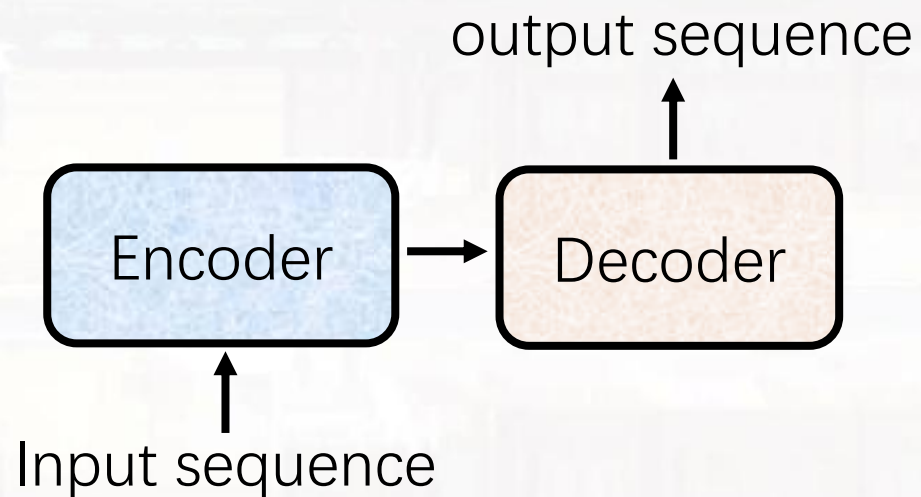
- **hand-crafted**

- **learned from data**

Each column represents a positional vector $e^i$



$q^i$   $k^i$   $v^i$

$e^i$ + $a^i$

-1     1

# Transformer



https://arxiv.org/abs/1706.03762

# Transformer

output sequence

Encoder → Decoder

↑
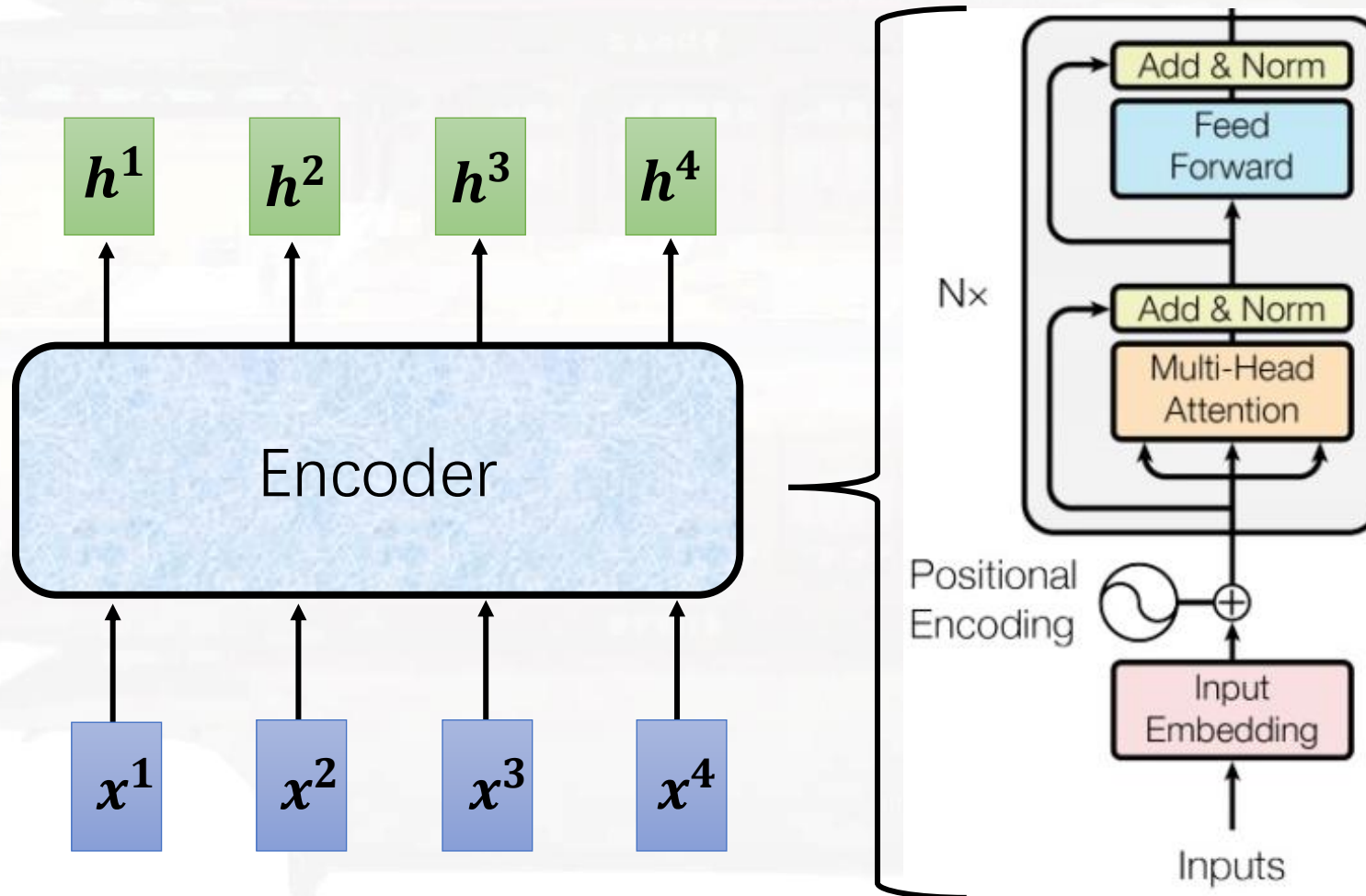Input sequence



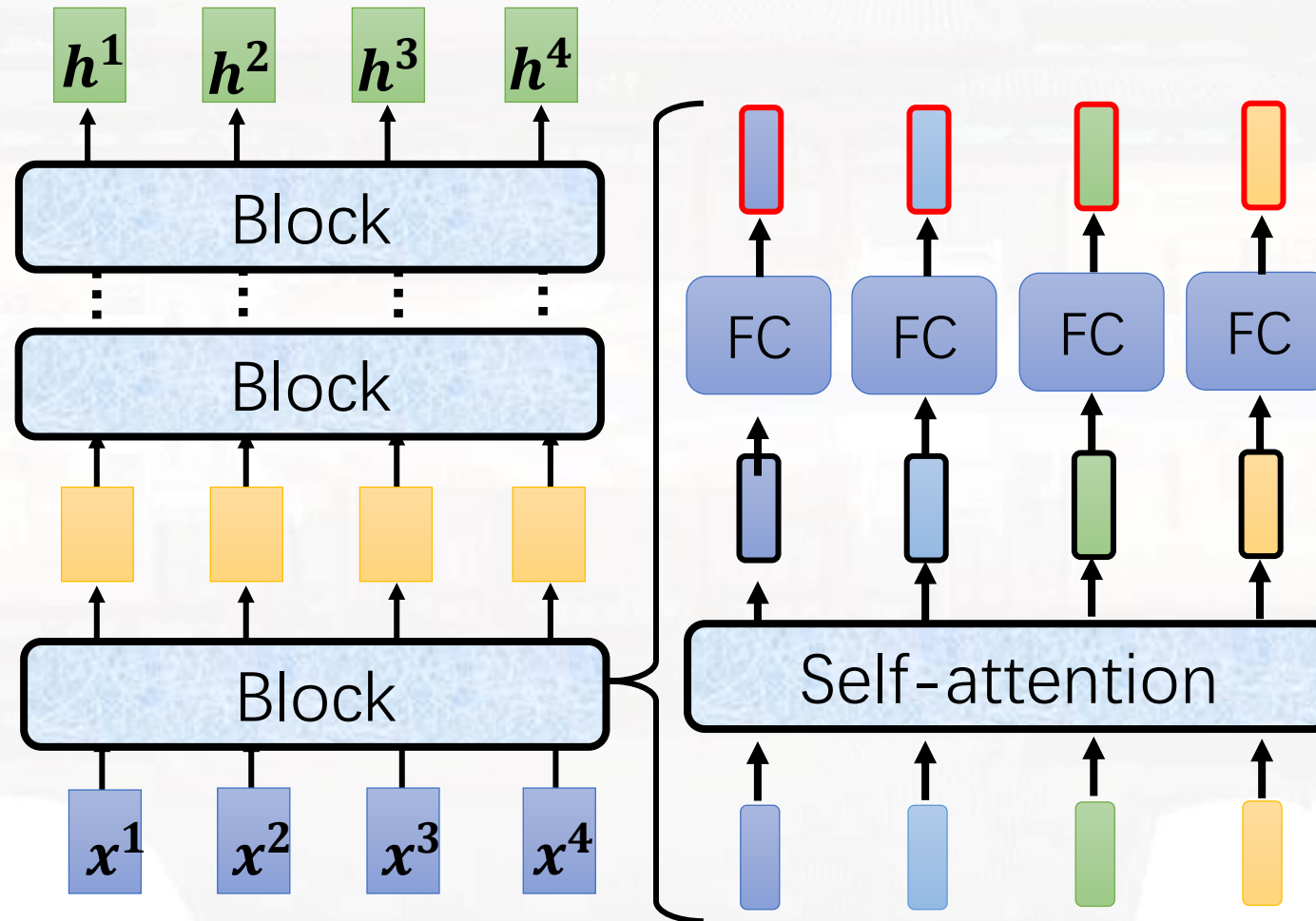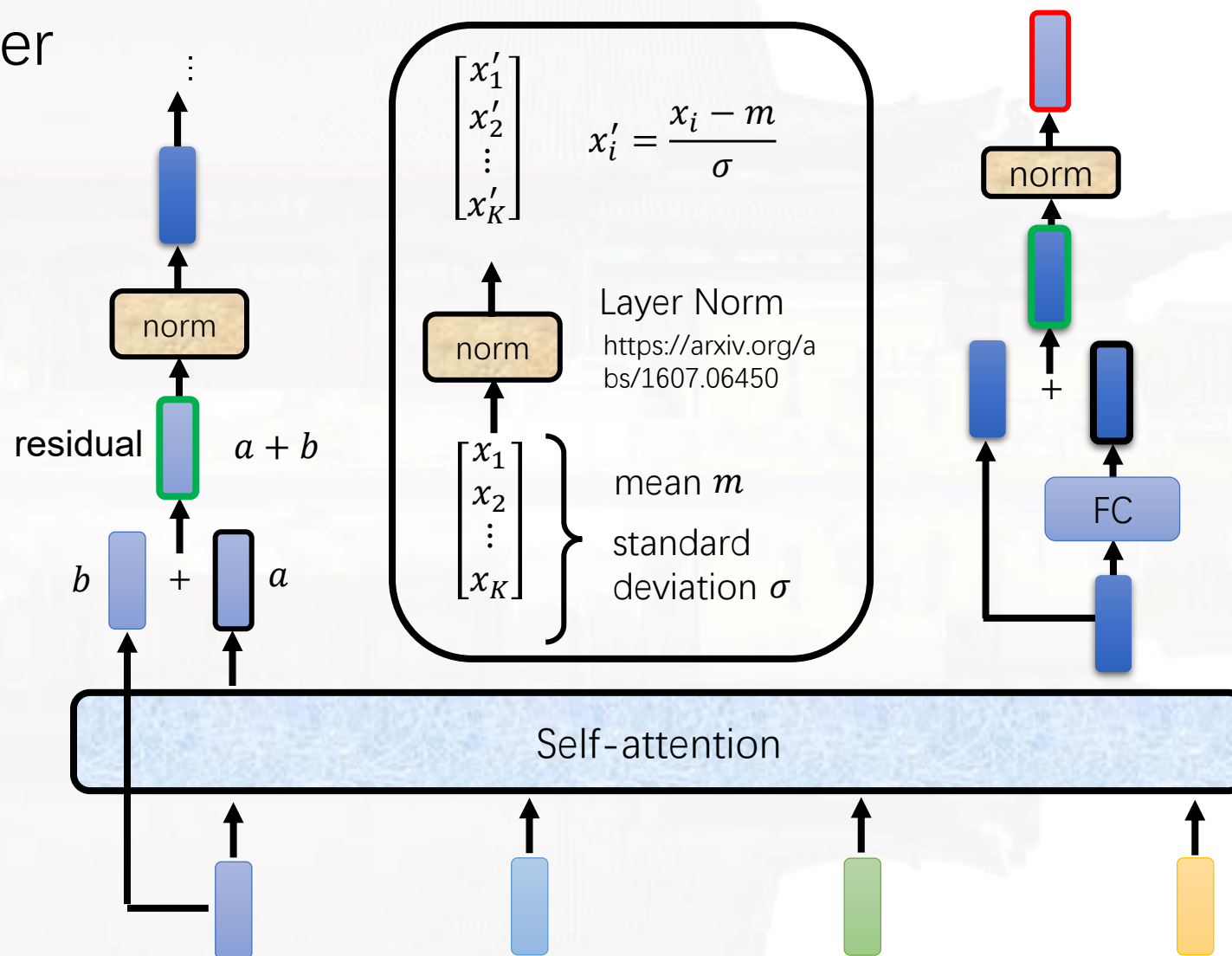Transformer

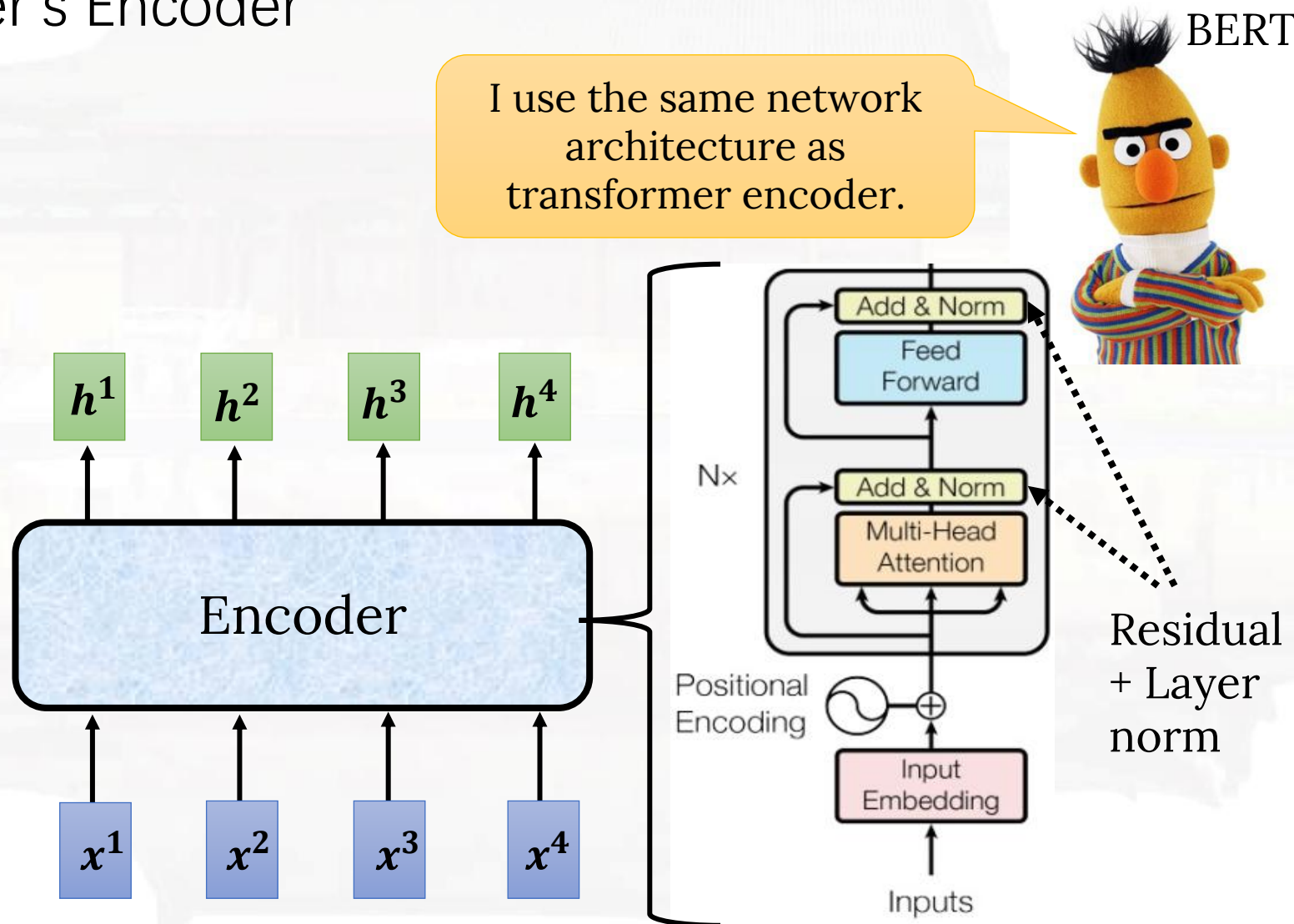https://arxiv.org/abs/1706.03762

# Transformer

- Transformer's Encoder

# Transformer

- Transformer's Encoder

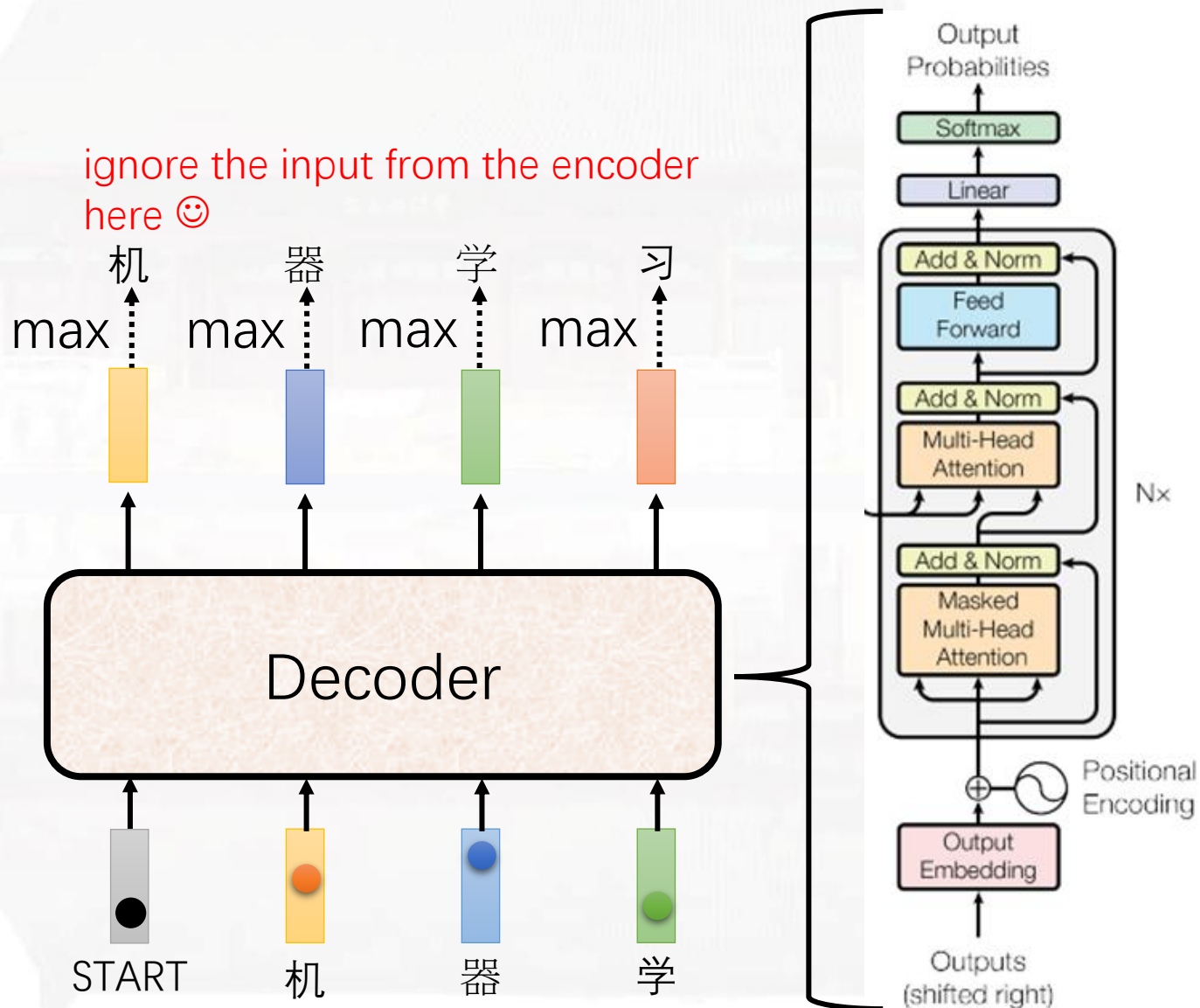# Transformer

- Transformer's Encoder

$$\begin{bmatrix} x'_1 \\ x'_2 \\ \vdots \\ x'_K \end{bmatrix} \qquad x'_i = \frac{x_i - m}{\sigma}$$

norm

Layer Norm
https://arxiv.org/abs/1607.06450

$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_K \end{bmatrix}$$

mean $m$

standard deviation $\sigma$

norm

residual  $\quad a + b$

$b$  $\quad + \quad$  $a$

norm

FC

Self-attention

# Transformer

- Transformer's Encoder

# Transformer



ignore the input from the encoder here ☺

# Transformer
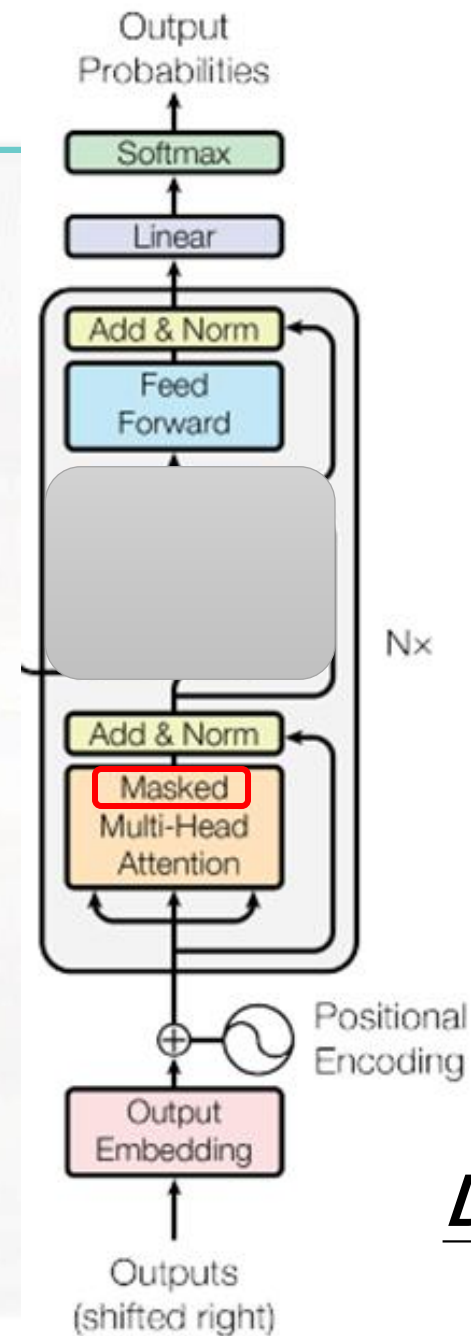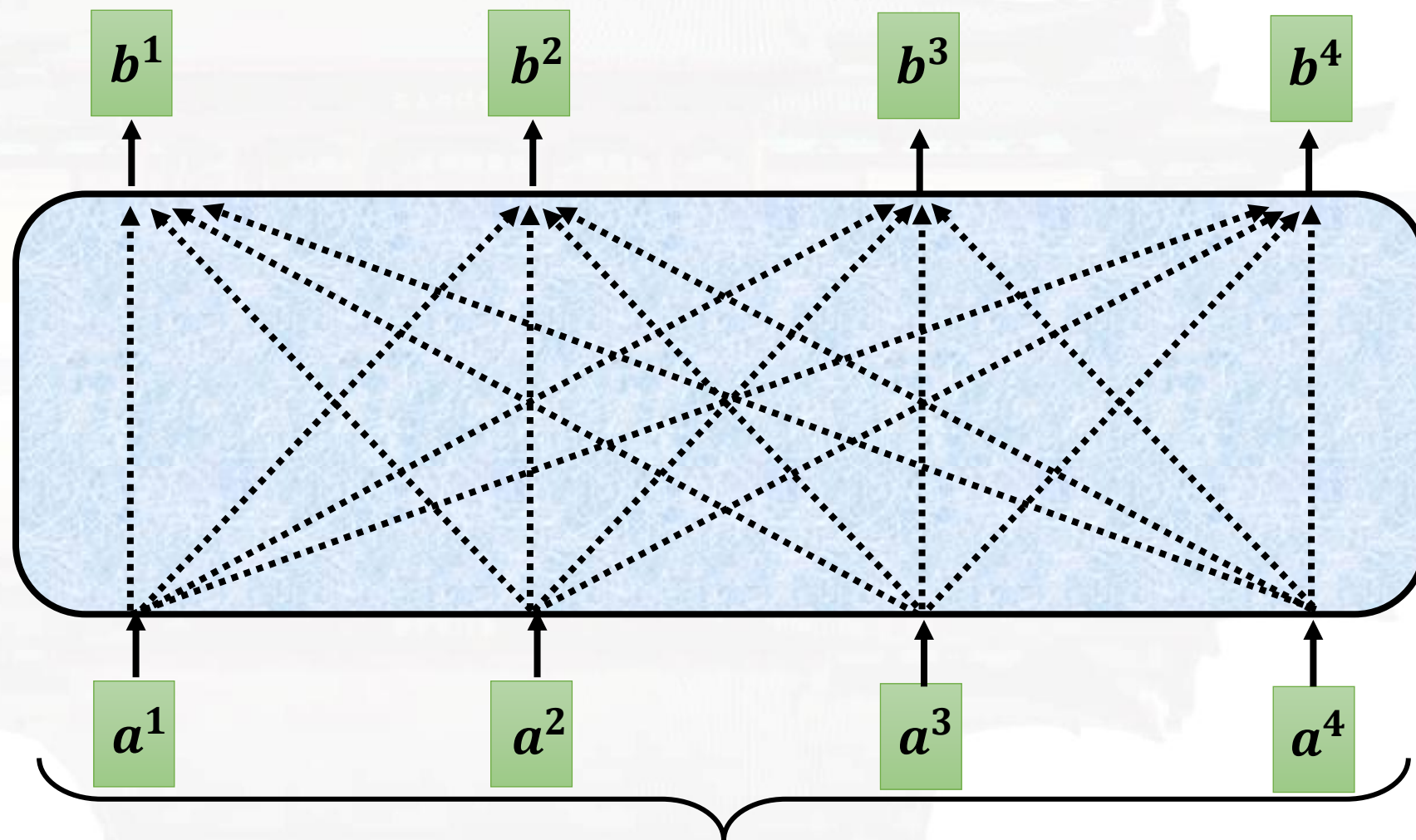


Encoder
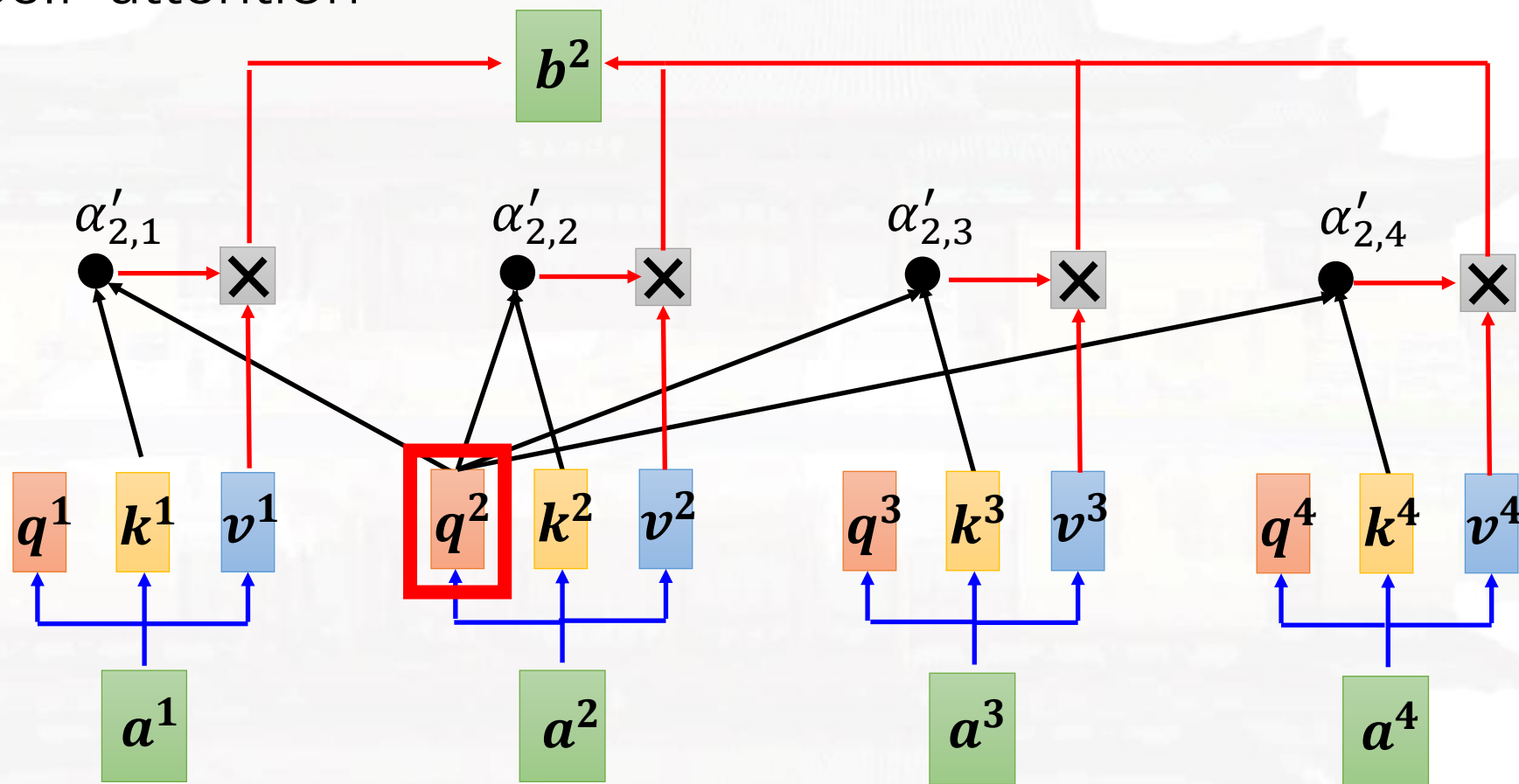
Decoder

# Transformer

- Masked Self-attention
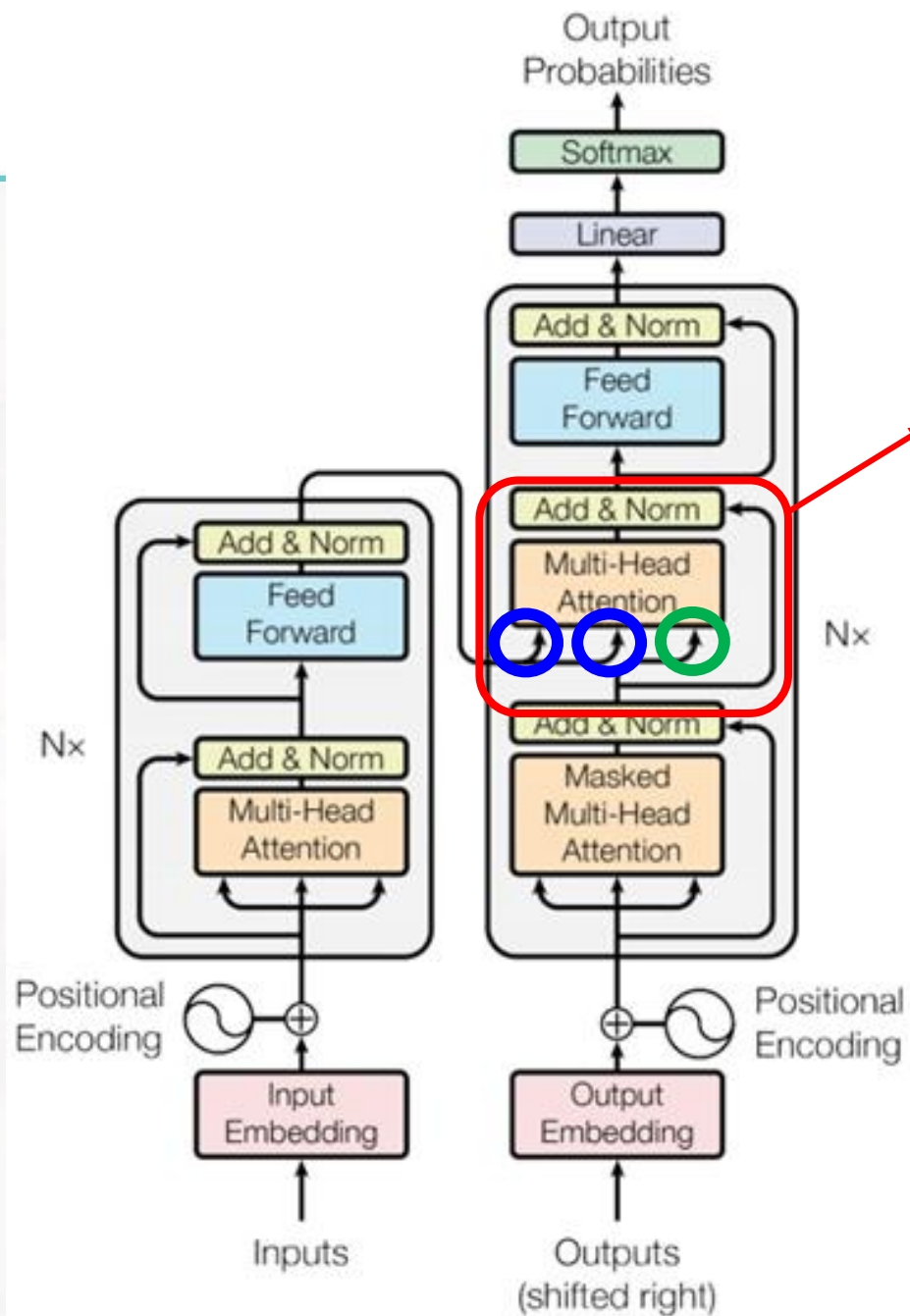


Can be either **input** or **a hidden layer**

# Transformer

- Masked Self-attention



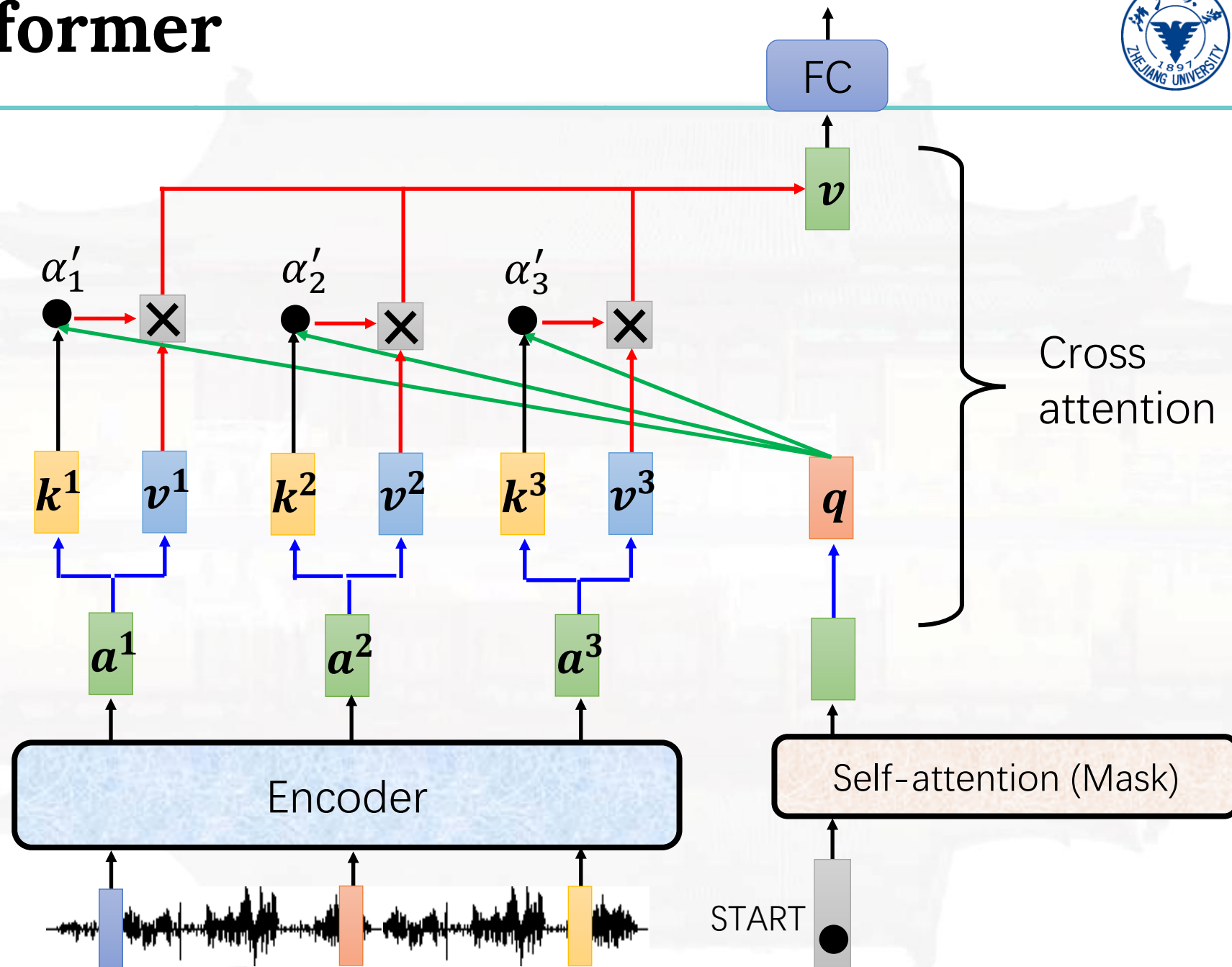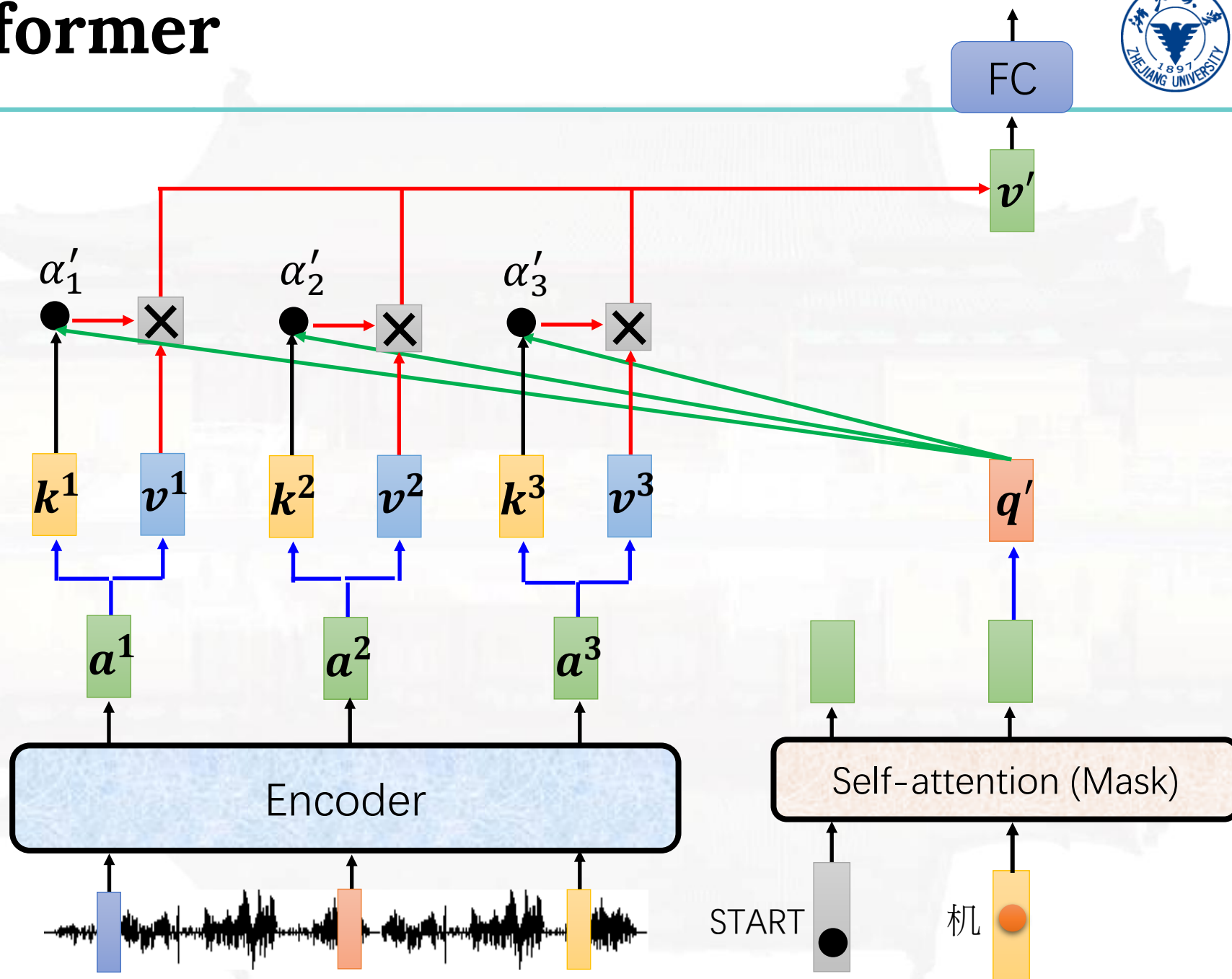Why masked?Consider how does decoder work

# Transformer



Cross attention
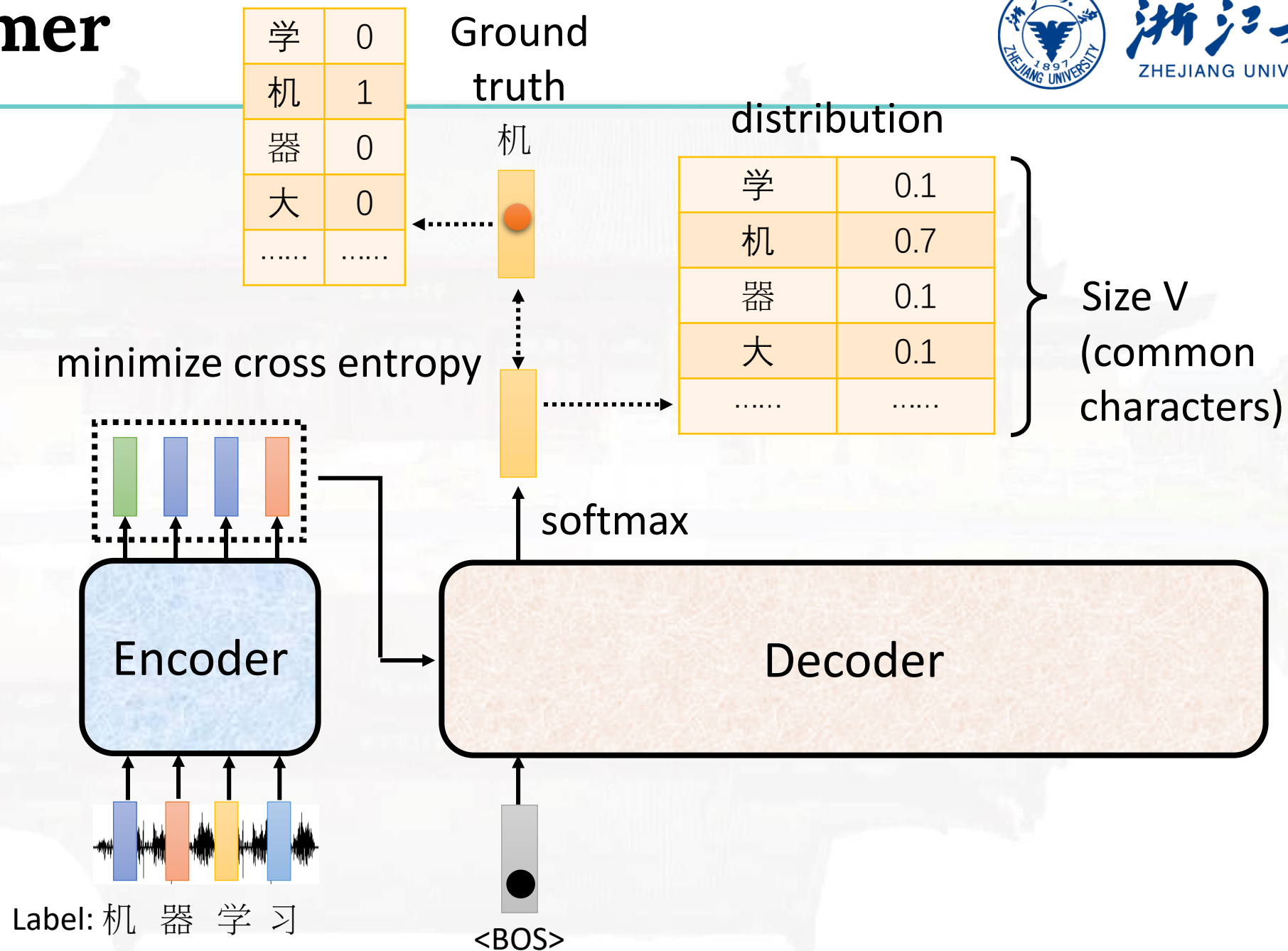
# Transformer

# Transformer

# **Transformer**

- Training

| 学 | 0 |
|---|---|
| 机 | 1 |
| 器 | 0 |
| 大 | 0 |
| …… | …… |

Ground truth

机

distribution

| 学 | 0.1 |
|---|---|
| 机 | 0.7 |
| 器 | 0.1 |
| 大 | 0.1 |
| …… | …… |

Size V (common characters)

minimize cross entropy

softmax

Encoder

Decoder

Label: 机 器 学 习

<BOS>

# Transformer

- Training



| 学 | 0.0 |
| 机 | 0.8 |
| 器 | 0.0 |
| 习 | 0.1 |
| ...... | ...... |
| END | 0.0 |

Size V (common characters)

softmax    distribution

Encoder

Decoder

（机器学习）

# Transformer

- Training