香港科技大學
THE HONG KONG
UNIVERSITY OF SCIENCE
AND TECHNOLOGY

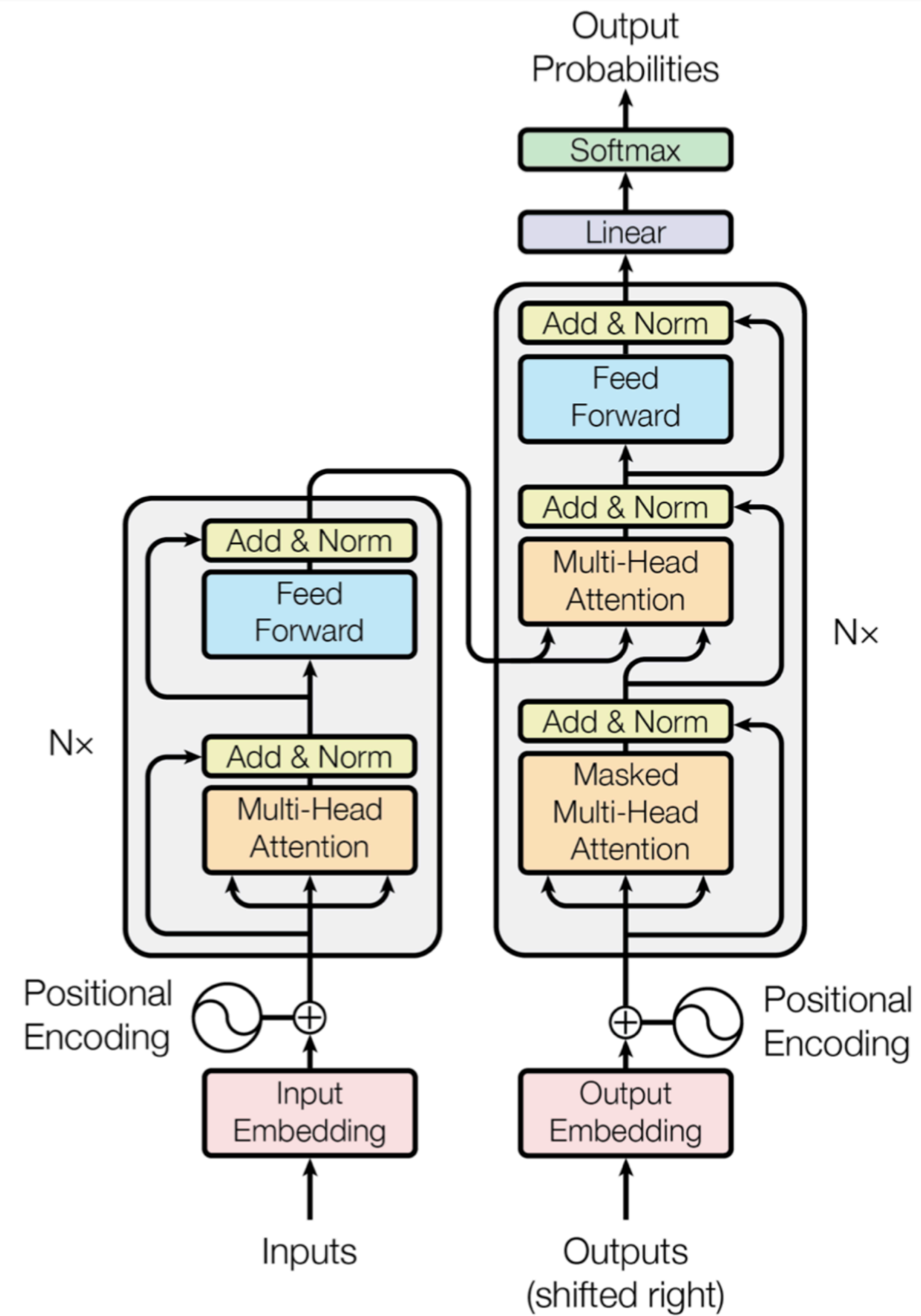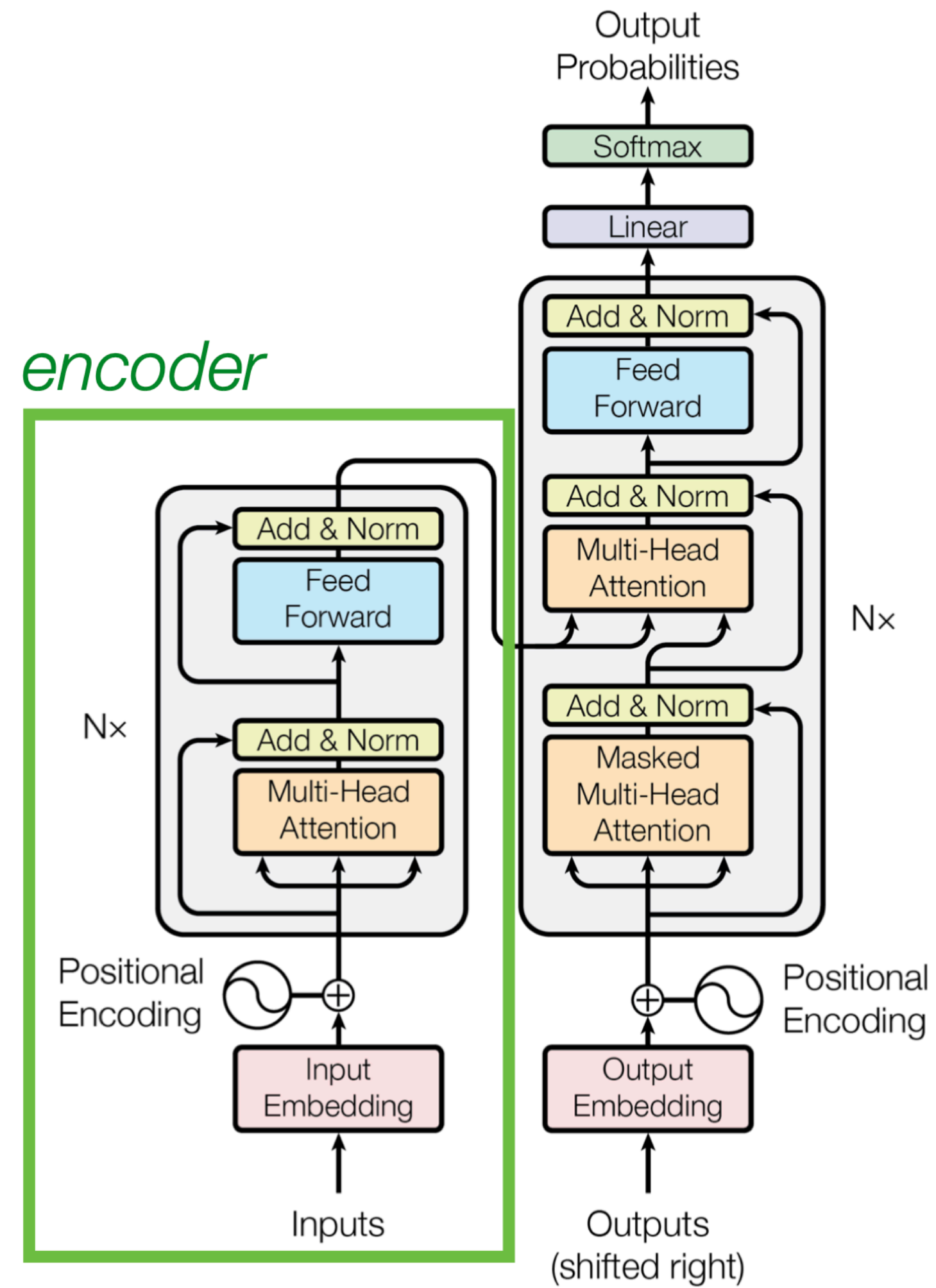# Transformers, VAEs

Junxian He
Nov 19, 2024

# Transformer



Vaswani et al. Attention is All You Need. NeurIPS 2017.

2

# Encoder

# Decoder

Output Probabilities

Softmax

Linear

*decoder*

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Add & Norm

Masked Multi-Head Attention

N×

Positional Encoding

Output Embedding

Outputs (shifted right)

*encoder*

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

N×

Positional Encoding

Input Embedding

Inputs

4

# Transformer Encoder

add & norm

residual

MLP

residual

Nx

layernorm

Add & Norm

Feed
Forward

Add & Norm

Multi-Head
Attention

Positional
Encoding

Input
Embedding

$$x - \text{mean}$$
$$\sigma$$

# Transformer Encoder

# Transformer Encoder



MLP

Self-attention

# Transformer Encoder

# What is Attention

Scaled Dot-Product Attention



API

Q: Query
K: key
V: value

# What is Attention

$$Q \in R^{n \times d} \qquad K \in R^{m \times d} \qquad V \in R^{m \times d}$$

## Scaled Dot-Product Attention



Q: Query
K: key
V: value

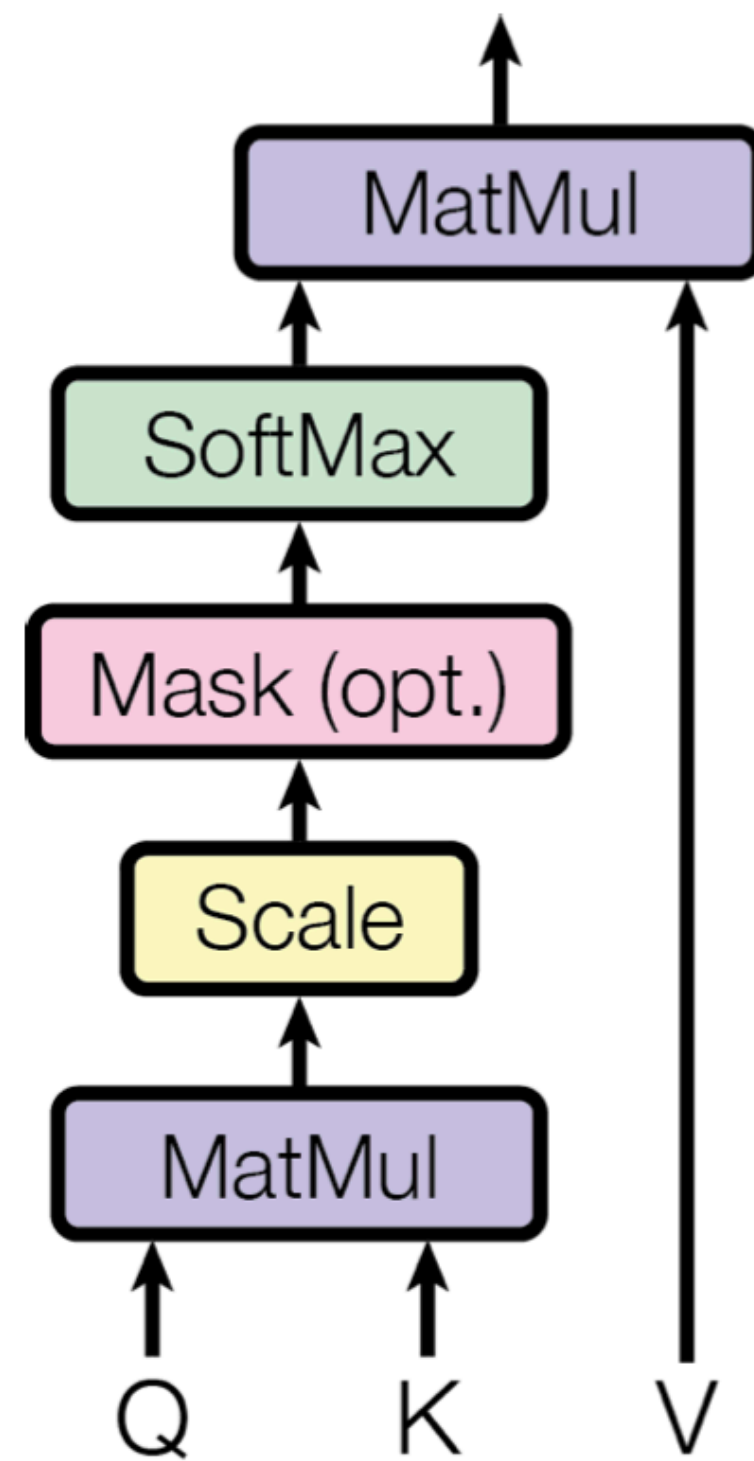$Q$: $n$ different queries, each $R^d$

$K, V$: $m$ different (key, value) pairs $R^d$

$(k, v)$

# What is Attention

$$Q \in R^{n \times d} \qquad K \in R^{m \times d} \qquad V \in R^{m \times d}$$

We have n queries, m (key, value) pairs

Scaled Dot-Product Attention



Q: Query
K: key
V: value

6

# What is Attention

$$Q \in R^{n \times d} \qquad K \in R^{m \times d} \qquad V \in R^{m \times d}$$

We have n queries, m (key, value) pairs
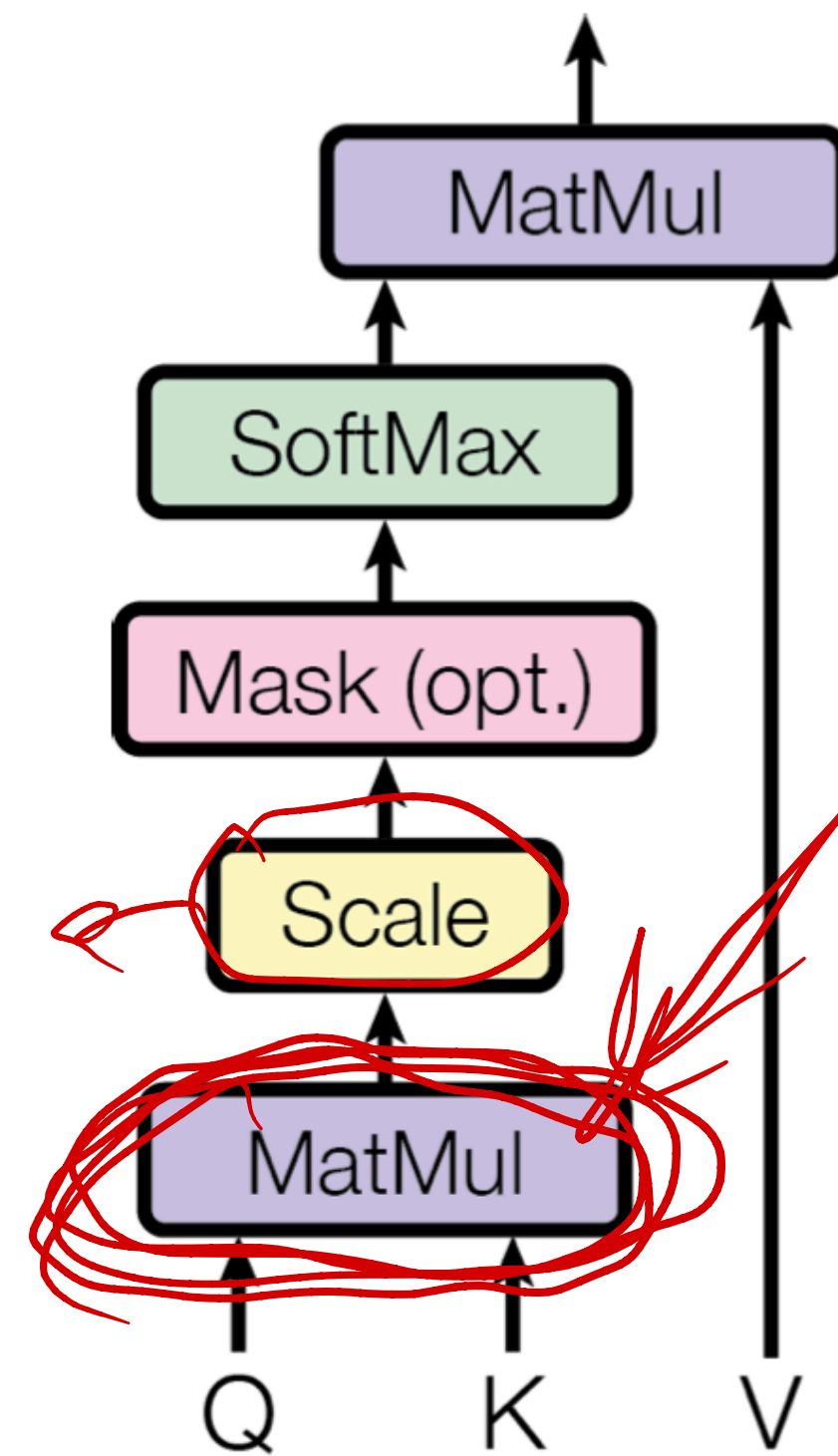
Scaled Dot-Product Attention



Attention weight = softmax($QK^T$)

$n \times m$

Q: Query
K: key
V: value

6

# What is Attention

$$Q \in R^{n \times d} \qquad K \in R^{m \times d} \qquad V \in R^{m \times d}$$

We have n queries, m (key, value) pairs

Attention weight = softmax($QK^T$)

$[0, 1]$

Dot-products grow large in magnitude

$q \in R^d \qquad k \in R^d \qquad d \cdot 1$

$$<q, k> \qquad \sum_{i=1}^{d} q_i k_i \qquad Var(q_i k_i) = 1$$

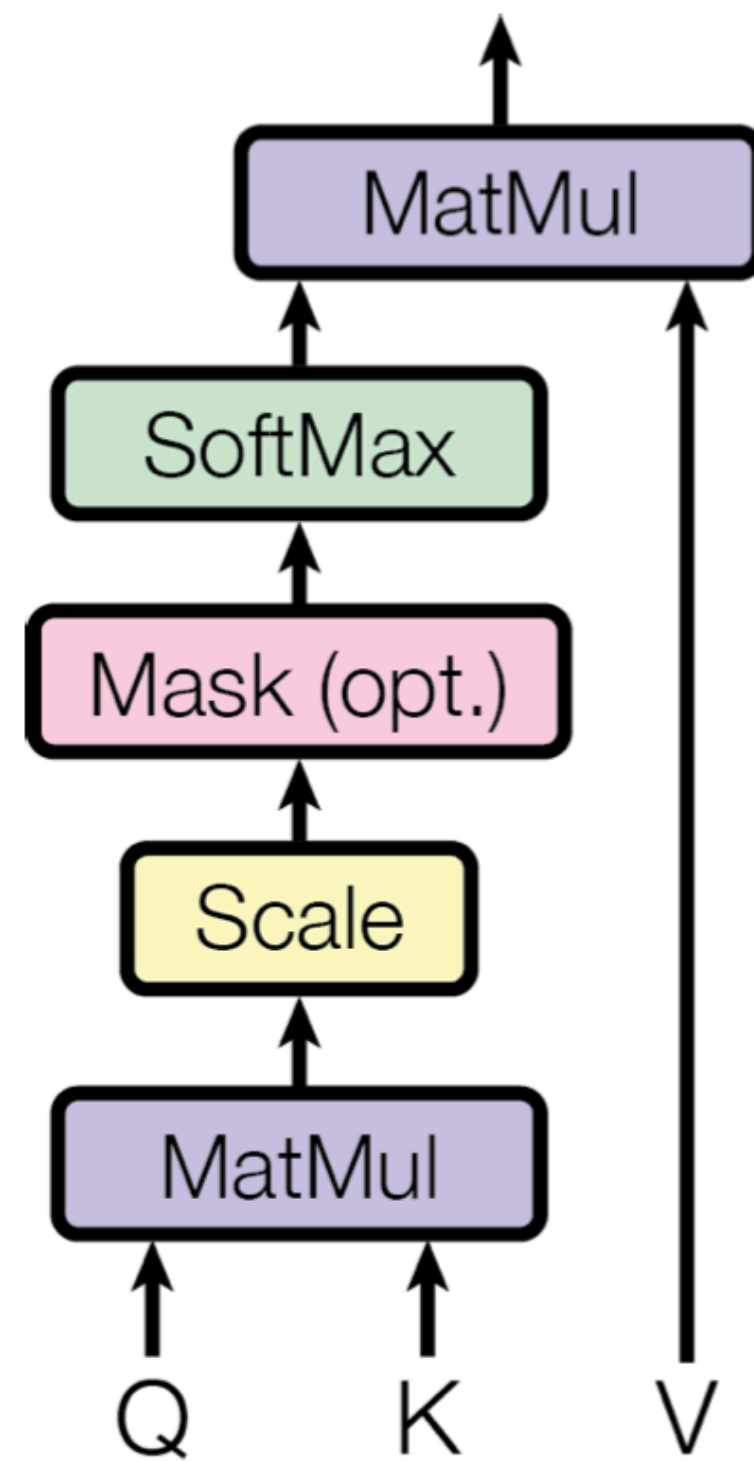$$Var(\sum_{i}^{d} q_i k_i) = d$$

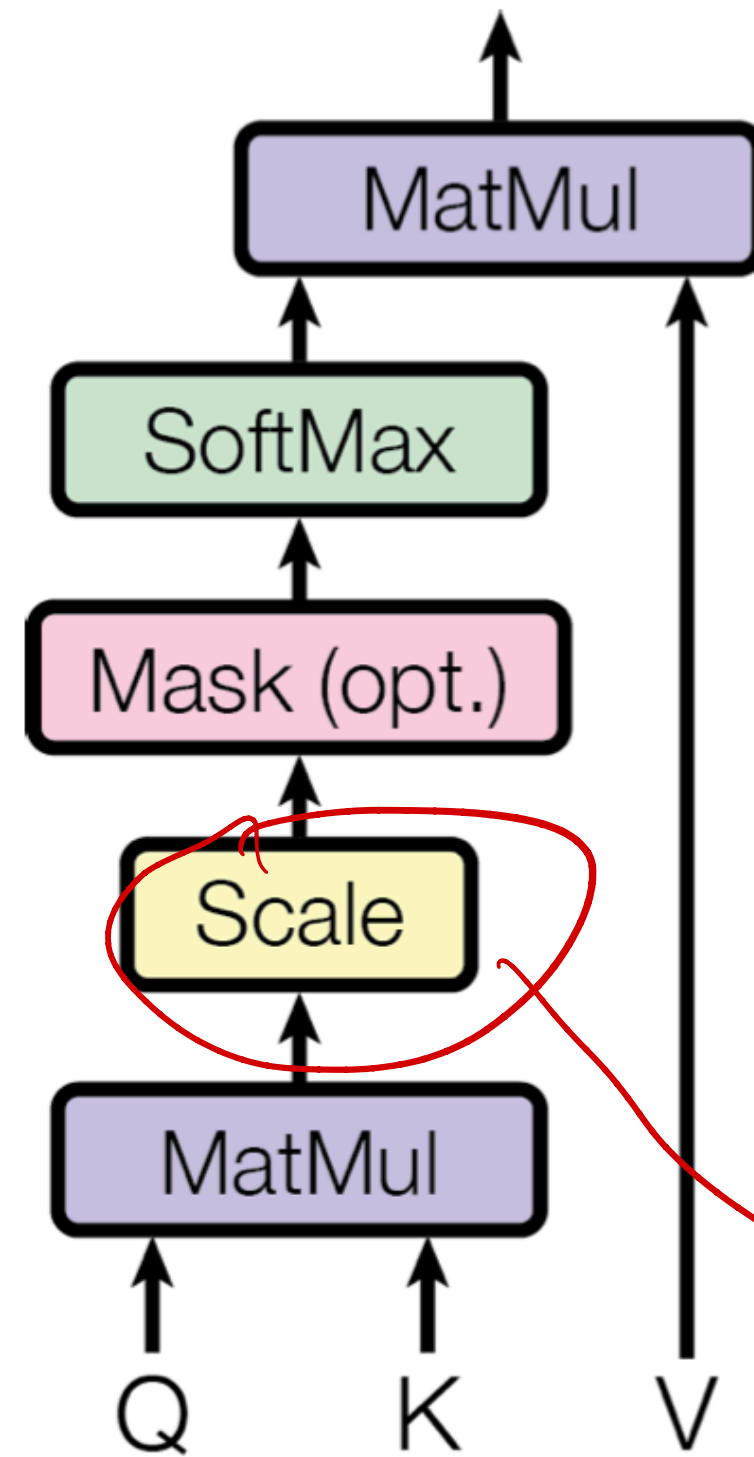## Scaled Dot-Product Attention



Q: Query
K: key
V: value

# What is Attention

$$Q \in R^{n \times d} \qquad K \in R^{m \times d} \qquad V \in R^{m \times d}$$

We have n queries, m (key, value) pairs

**Scaled Dot-Product Attention**



Attention weight = softmax($QK^T$)

Dot-products grow large in magnitude

Scaled Attention weight = softmax($\dfrac{QK^T}{\sqrt{d_k}}$)

$$Var\left(\frac{X}{\sqrt{d_k}}\right) = \frac{Var(X)}{d_k}$$

Q: Query
K: key
V: value

6

# What is Attention

Scaled Dot-Product Attention

$$Q \in R^{n \times d} \qquad K \in R^{m \times d} \qquad V \in R^{m \times d}$$

We have n queries, m (key, value) pairs

Attention weight = softmax($QK^T$)

Dot-products grow large in magnitude

Scaled Attention weight = softmax($\frac{QK^T}{\sqrt{d_k}}$)

Shape is mxn

Q: Query
K: key
V: value

6

# What is Attention

$$Q \in R^{n \times d} \qquad K \in R^{m \times d} \qquad V \in R^{m \times d}$$

We have n queries, m (key, value) pairs

## Scaled Dot-Product Attention



Q: Query
K: key
V: value

Attention weight = $\text{softmax}(QK^T)$

Dot-products grow large in magnitude

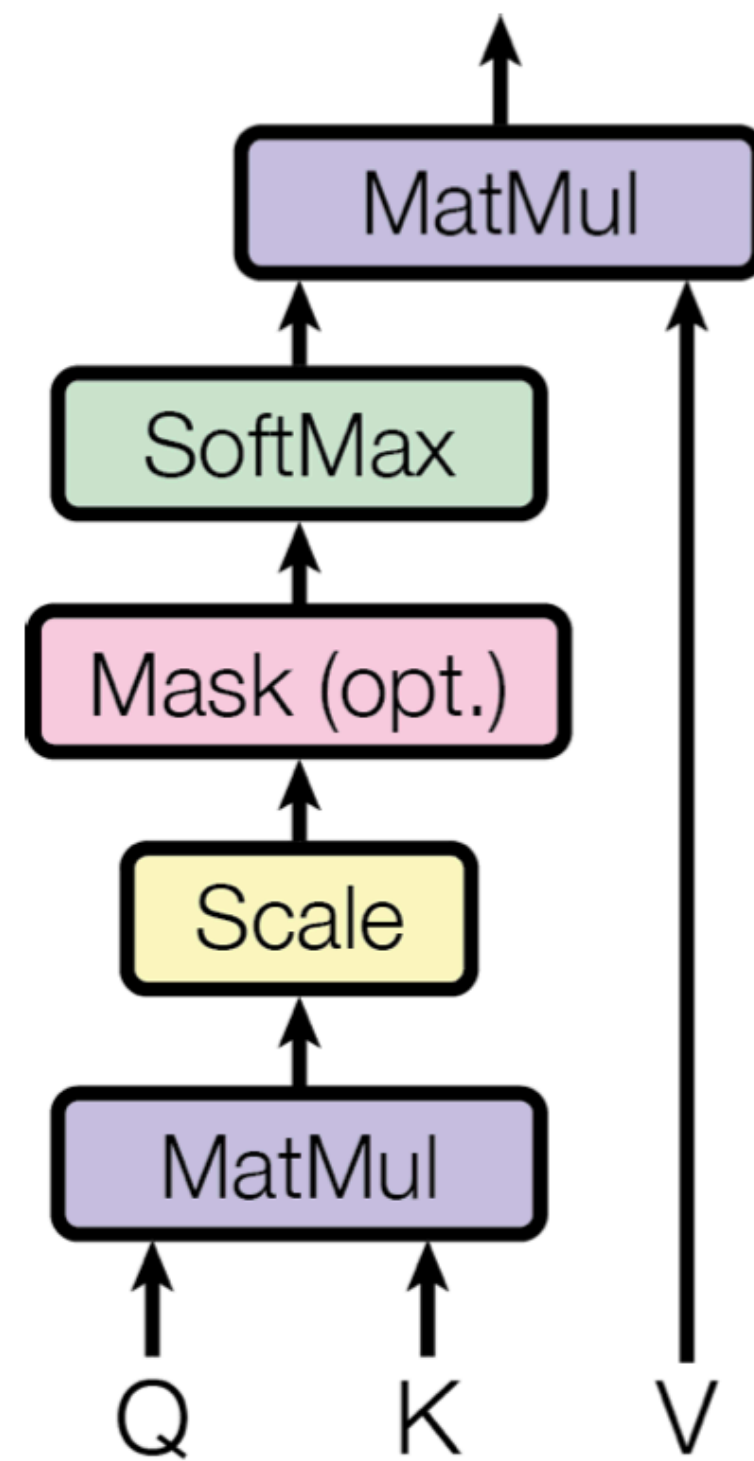Scaled Attention weight = $\text{softmax}(\dfrac{QK^T}{\sqrt{d_k}})$

Shape is mxn

Attention weight represents the strength to "attend" values V
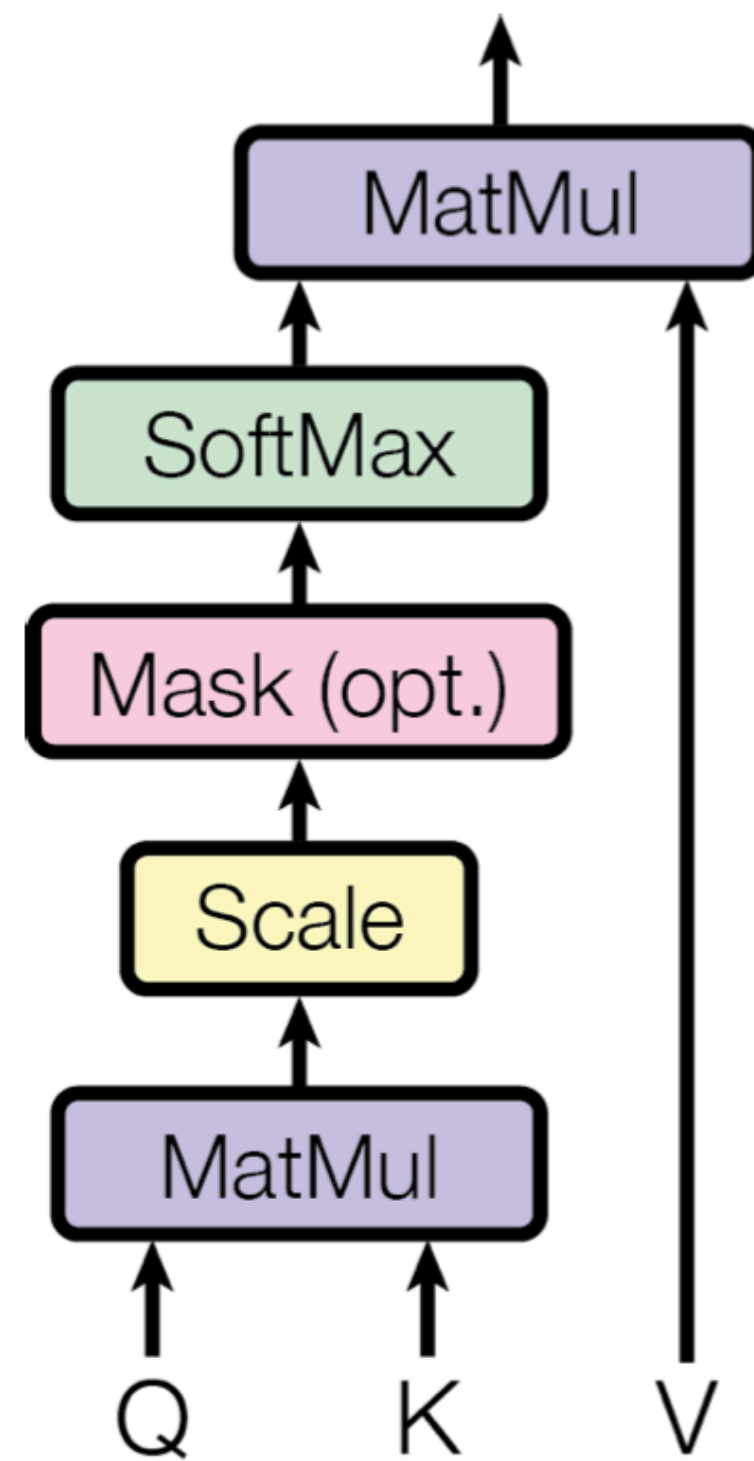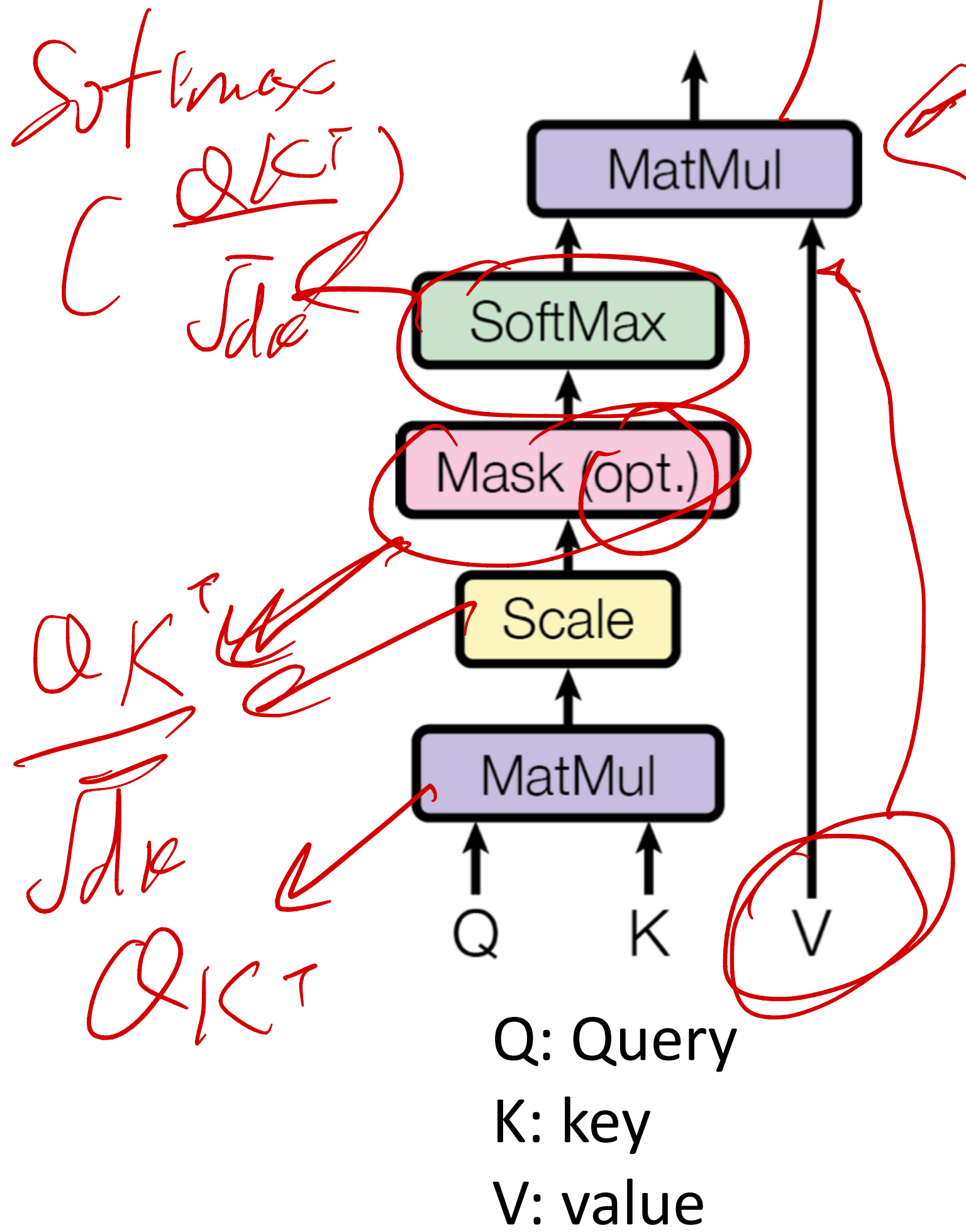
6

# What is Attention

$$Q \in R^{n \times d} \qquad K \in R^{m \times d} \qquad V \in R^{m \times d}$$

We have n queries, m (key, value) pairs

## Scaled Dot-Product Attention



Softmax
$(\frac{QK^T}{\sqrt{d_k}})$

$\frac{QK^T}{\sqrt{d_k}}$

$QK^T$

Q: Query
K: key
V: value

Attention weight = softmax($QK^T$)

Dot-products grow large in magnitude

Scaled Attention weight = softmax($\frac{QK^T}{\sqrt{d_k}}$)

Shape is mxn

Attention weight represents the strength to "attend" values V

$$\text{Attention}(Q, K, V) = \text{softmax}(\frac{QK^T}{\sqrt{d_k}})V$$
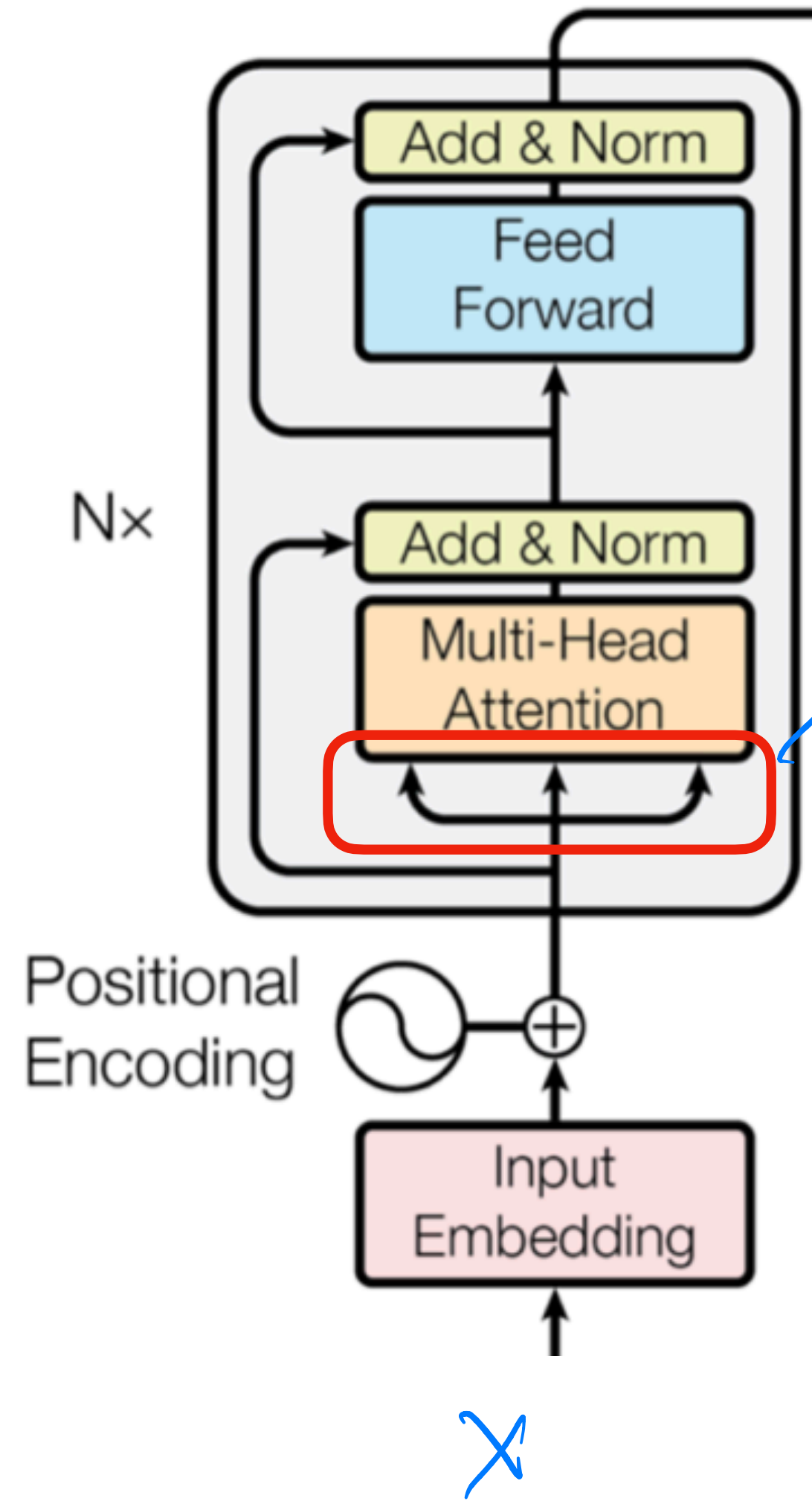
linear combination of all the values

6

$R^d$

$\boxed{q_1}$ $q_2$

$R^d$ $R^d$

$K_1$ $K_2$ $K_3$ $\cdots K_m$

$\boxed{n \times d}$
$\downarrow$

difference queries $V_1$ $V_2$ $V_3$ $\cdots V_m$

$q_n$ $R^d$

are independent

$\in R$
$\uparrow$

$\text{Softmax}\left( \dfrac{q_1 k_1^T, \; q_1 k_2^T, \; q_1 k_3^T \cdots q_1 k_m^T}{\sqrt{d}} \right) = W_1, \; W_2 \cdots W_m$

$\text{Output} = W_1 \cdot V_1 + W_2 \cdot V_2 + W_3 V_3 \cdots \cdots W_m V_m$

# Q, K, V



Positional Encoding

Input Embedding

$Q, K, V$

$X$

What are Q, K, V in the transformer

# Self-Attention



Input is X

X × W^Q = Q

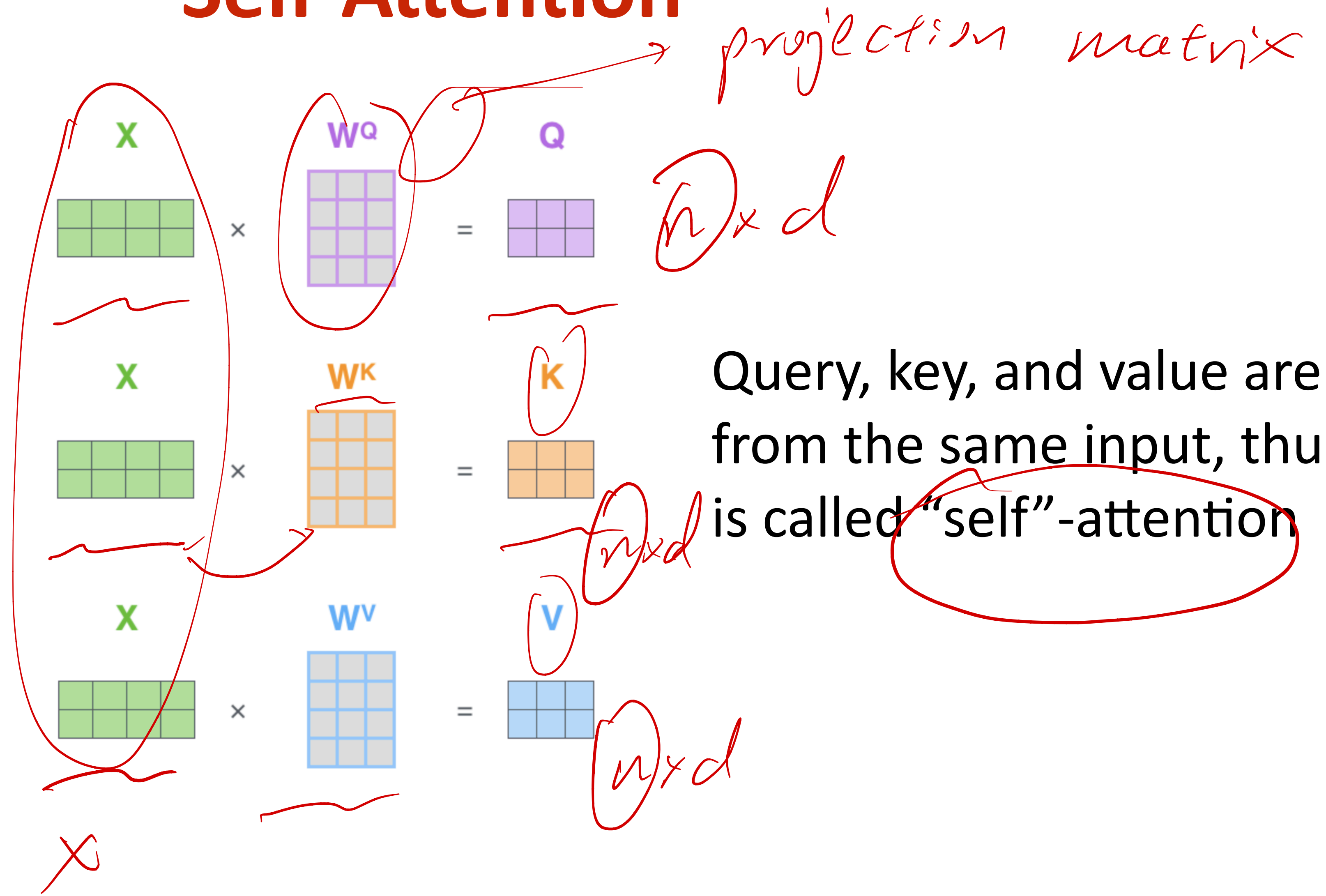X × W^K = K

X × W^V = V

Jay Alammar. The Illustrated Transformer.

# Self-Attention

projection matrix

$n \times d$

$n \times d$

$n \times d$

X

X   W^Q   Q

X   W^K   K

X   W^V   V

Query, key, and value are from the same input, thus it is called "self"-attention

Jay Alammar. The Illustrated Transformer.

# Self-Attention

*Parameters*



Input is X

X $\times$ W$^Q$ = Q

X $\times$ W$^K$ = K

X $\times$ W$^V$ = V
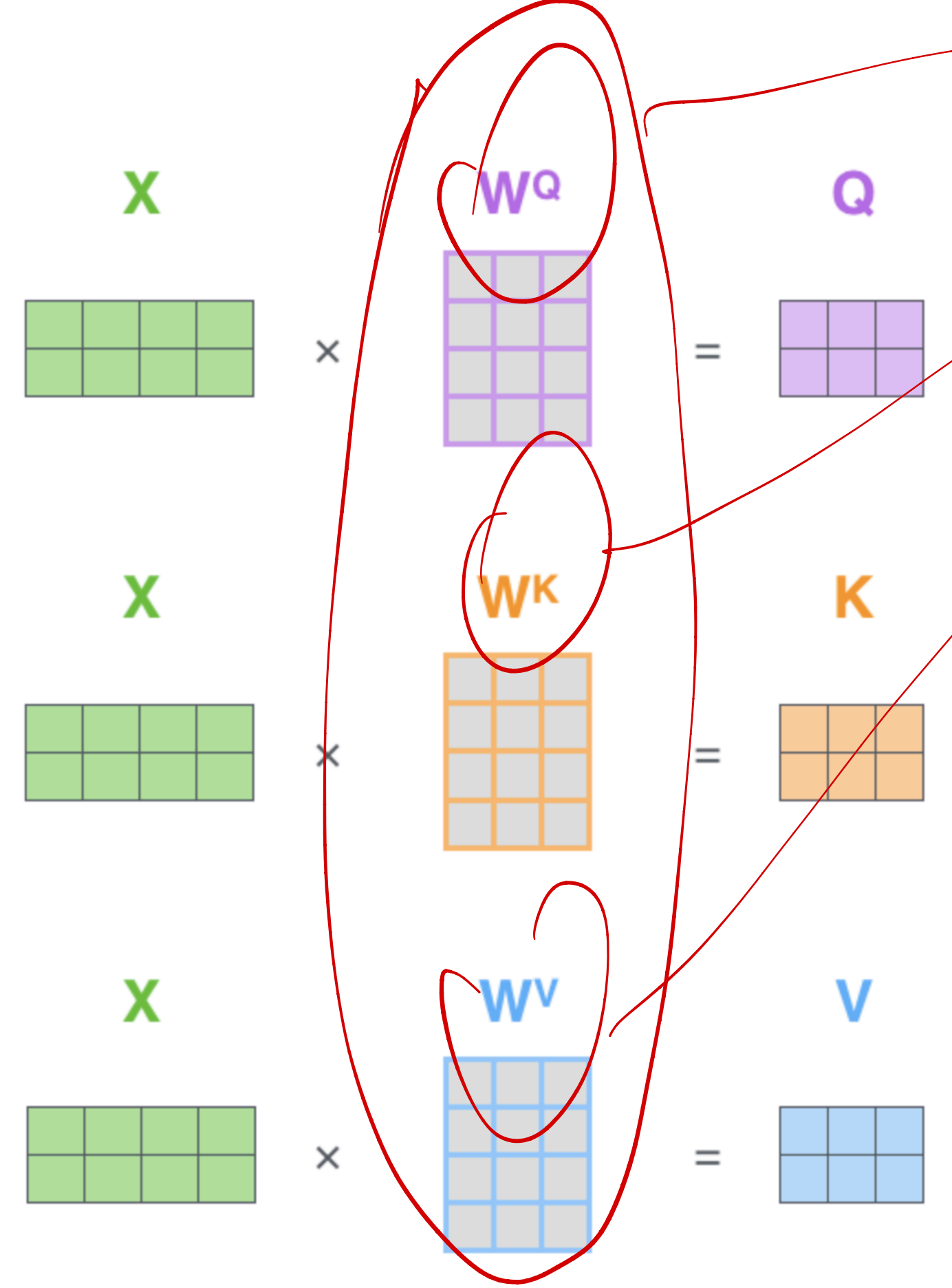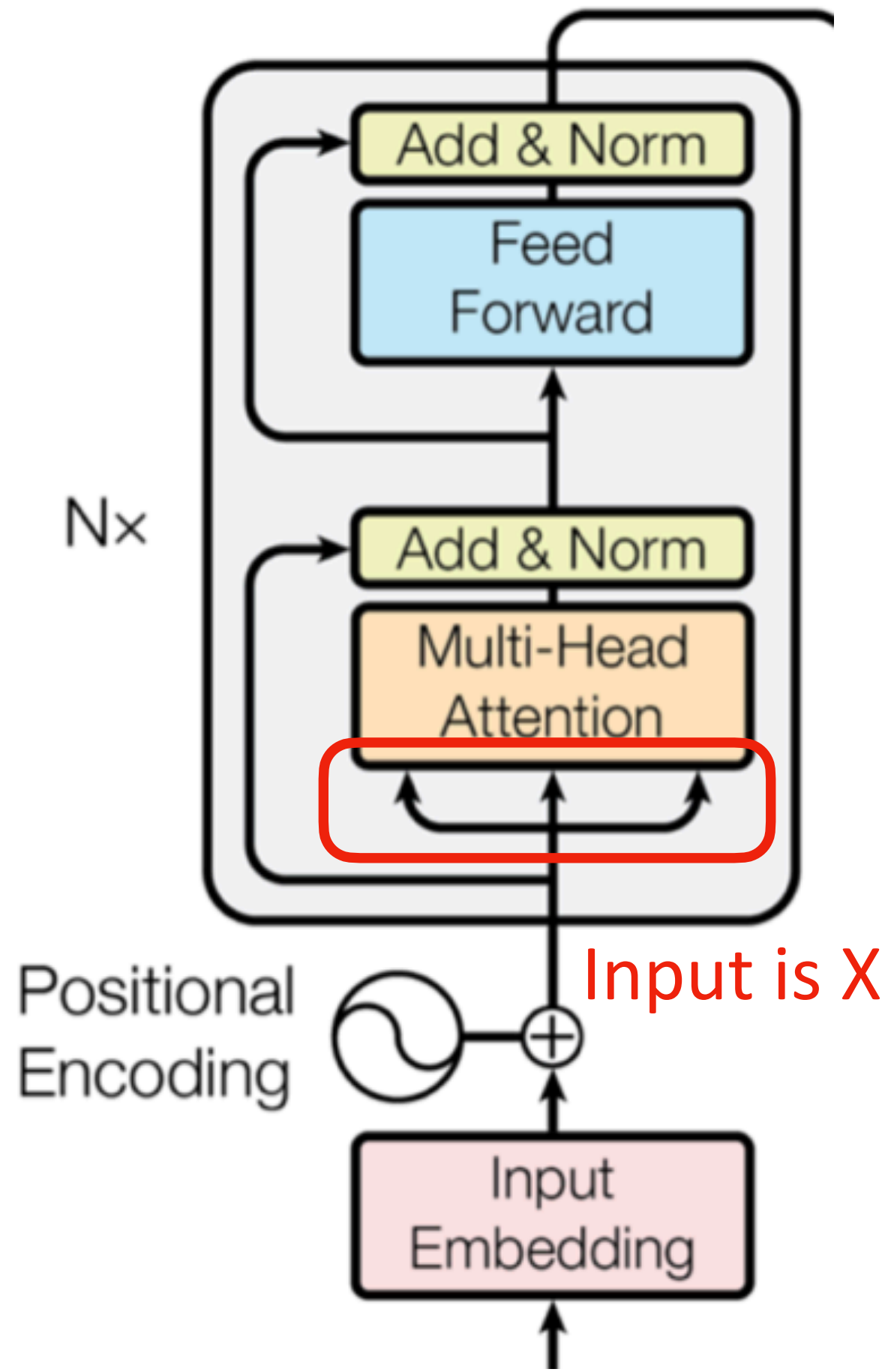
Query, key, and value are from the same input, thus it is called "self"-attention

$$\text{softmax}\left(\frac{Q \times K^T}{\sqrt{d_k}}\right) V = Z$$
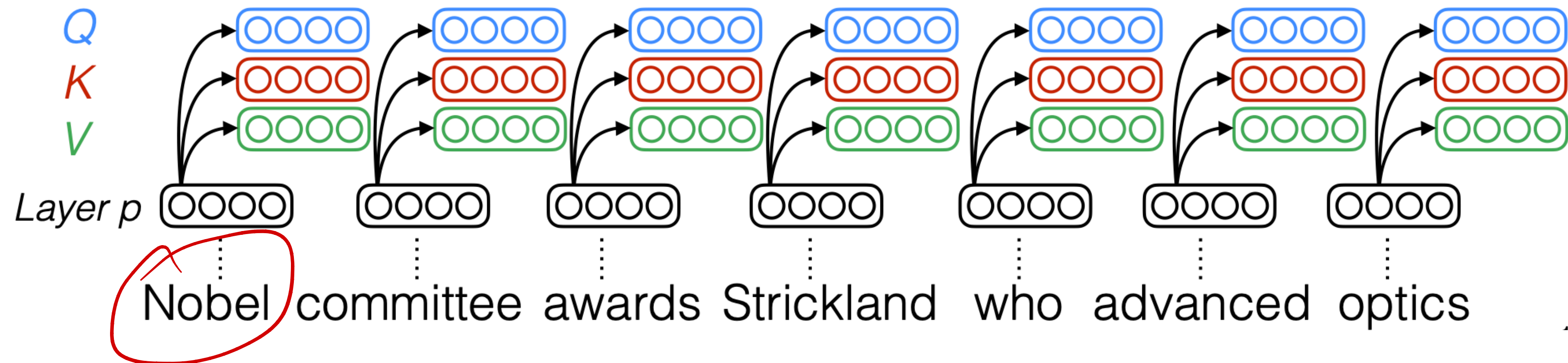
8

8

Jay Alammar. The Illustrated Transformer.

# Self-Attention

# Self-Attention

At each step, the attention computation attends
to all steps in the input example

# Self-Attention



*Layer p*

Nobel committee awards Strickland who advanced optics

# Self-Attention

Attention weight on every word in the sequence

11

# Self-Attention

# Self-Attention



13

# Multi-Head Attention

## Multi-Head Attention



$d$

[ $Output_1$ , $O_2$ . $O_3$ --- $O_n$ ]

output

$SA_1$

$Q$ $K$ $V$

$output_2$

$SA_2$

$O_2$ $K_2$ $V_2$

$SA_h$

14

# Multi-Head Self-Attention

Jay Alammar. The Illustrated Transformer.

# Multi-Head Self-Attention

ATTENTION HEAD #0

$Q_0$

$W_0^Q$

$K_0$

$W_0^K$

$V_0$

$W_0^V$

ATTENTION HEAD #1

$Q_1$

$W_1^Q$

$K_1$

$W_1^K$

$V_1$

$W_1^V$

15

Jay Alammar. The Illustrated Transformer.

# Multi-Head Self-Attention



Jay Alammar. The Illustrated Transformer.

# Multi-Head Self-Attention

Jay Alammar. The Illustrated Transformer.

# Multi-Head Self-Attention

1) Concatenate all the attention heads

$Z_0$ $Z_1$ $Z_2$ $Z_3$ $Z_4$ $Z_5$ $Z_6$ $Z_7$

Jay Alammar. The Illustrated Transformer.

# Multi-Head Self-Attention

Output projection

1) Concatenate all the attention heads

$Z_0$ $Z_1$ $Z_2$ $Z_3$ $Z_4$ $Z_5$ $Z_6$ $Z_7$

2) Multiply with a weight
matrix $W^O$ that was trained
jointly with the model

X

$W^O$

$2 \times 500$

$2 \times 1024$

$W_0$

16

# Multi-Head Self-Attention

1) Concatenate all the attention heads

$Z_0$  $Z_1$  $Z_2$  $Z_3$  $Z_4$  $Z_5$  $Z_6$  $Z_7$

2) Multiply with a weight matrix $W^O$ that was trained jointly with the model

X

$W^O$

3) The result would be the Z matrix that captures information from all the attention heads. We can send this forward to the FFNN

Z

=

Jay Alammar. The Illustrated Transformer.

# Multi-head Self-Attention

# Multi-head Self-Attention



Concat and output projection

$M_H$
$M_1$

optics
advanced
who
Strickland
awards
committee
$A$  Nobel

$Q$
$K$
$V$

*Layer p*

Nobel  committee  awards  Strickland  who  advanced  optics

18

# Multi-head Self-Attention + FFN

# Transformer Encoder

Currently we only cover the encoder side

# Transformer Encoder

Currently we only cover the encoder side



*encoder*

*encoder*

*decoder*

seq2seq

2RNN

# Transformer Encoder

Currently we only cover the encoder side



This encoder-decoder arch is originally proposed as a seq2seq arch, for classification tasks, often only encoder is used. And language models often only have a decoder

# Transformer Decoder in Seq2Seq

# Transformer Decoder in Seq2Seq



encoder

decoder

Self-attention

# Transformer Decoder in Seq2Seq

# Transformer Decoder in Seq2Seq

*decoder*

*encoder*

Cross-attention

Self-attention

$R^{n \times d}$ $\quad R^{m \times d}$ $\quad R^{m \times d}$

$Q \quad K, \quad V$ $\qquad m \neq n$

Cross-attention uses the output of encoder as input

# Transformer Decoder in Seq2Seq



decoder

encoder

Cross-attention

Self-attention

Cross-attention uses the output of encoder as input

# Masked Attention

# Masked Attention

# Masked Attention

*[handwritten: autoregressive decoding]*

Typical attention attends to the entire sequence, while masked attention only attends to the ones on the left because future words have not been generated

*[handwritten annotations: T5, no mask, BART, encoder, decoder, um (a) student, um a student]*



Output Probabilities
Softmax
Linear
Add & Norm
Feed Forward
Add & Norm
Multi-Head Attention
Add & Norm
Masked Multi-Head Attention
Positional Encoding
Output Embedding
Outputs (shifted right)

Add & Norm
Feed Forward
Add & Norm
Multi-Head Attention
Positional Encoding
Input Embedding
Inputs

Decoding time step: 1 2 3 4 5 6    OUTPUT    I

$K_{encdec}$  $V_{encdec}$    Linear + Softmax

ENCODERS    DECODERS

EMBEDDING WITH TIME SIGNAL

EMBEDDINGS

INPUT    Je    suis    étudiant    PREVIOUS OUTPUTS    I

22

# chatGPT

[mask]

next word prediction

dec
r
eno

am a student in HK

transformer decoder

am a student in HK

mask

FFN

mask SA

uni directional attention

# Position Embeddings

# Position Embeddings



Question: If we shuffle the order of words in the sequence, will that change the attention output and feed forward output of the corresponding word?

$z_1 \quad z_2 \quad t_3 \quad z_4$

I am a student $\quad [z_4 = z_1$

$z_1' \quad\quad z_2' \quad z_3' \quad z_4'$

student a am I

# Position Embeddings



Question: If we shuffle the order of words in the sequence, will that change the attention output and feed forward output of the corresponding word?

Position embeddings are added to each word embedding, otherwise our model is unaware of the position of a word

23

# Positional Encoding



EMBEDDING WITH TIME SIGNAL

$x_1$    $x_2$    $x_3$

=  =  =

POSITIONAL ENCODING

$t_1$    $t_2$    $t_3$

+  +  +

EMBEDDINGS

$x_1$    $x_2$    $x_3$

INPUT    Je    suis    étudiant

# Transformer Positional Encoding

*absolute position*

*RoPE*

$$PE_{(pos, 2i)} = \sin\left(\frac{pos}{10000^{2i/d_{model}}}\right)$$

$$PE_{(pos, 2i+1)} = \cos\left(\frac{pos}{10000^{2i/d_{model}}}\right)$$

*1st   2nd*

*relative position*

*RoPE*

Positional encoding is a 512d vector
*i* = a particular dimension of this vector
*pos* = dimension of the word
*d_model* = 512

$2i+1 \quad 2i = 5/2$

$10 , 2$

Word A --- Word B

two word

25

# Complexity

| Layer Type | Complexity per Layer | Sequential Operations |
|---|---|---|
| Self-Attention | $O(n^2 \cdot d)$ | $O(1)$ |
| Recurrent | $O(n \cdot d^2)$ | $O(n)$ |
| Convolutional | $O(k \cdot n \cdot d^2)$ | $O(1)$ |
| Self-Attention (restricted) | $O(r \cdot n \cdot d)$ | $O(1)$ |

n is sequence length, d is embedding dimension.



26

# Complexity

| Layer Type | Complexity per Layer | Sequential Operations |
|---|---|---|
| Self-Attention | $O(n^2 \cdot d)$ | $O(1)$ |
| Recurrent | $O(n \cdot d^2)$ | $O(n)$ |
| Convolutional | $O(k \cdot n \cdot d^2)$ | $O(1)$ |
| Self-Attention (restricted) | $O(r \cdot n \cdot d)$ | $O(1)$ |

n is sequence length, d is embedding dimension.

Restricted self-attention means not attending all words in the sequence, but only a restricted field

# Complexity

$n^2 \times d + n \times d^2$

$n \times n$

$n^2 \times d$

| Layer Type | Complexity per Layer | Sequential Operations |
|---|---|---|
| Self-Attention | $O(n^2 \cdot d)$ | $O(1)$ |
| Recurrent | $O(n \cdot d^2)$ | $O(n)$ |
| Convolutional | $O(k \cdot n \cdot d^2)$ | $O(1)$ |
| Self-Attention (restricted) | $O(r \cdot n \cdot d)$ | $O(1)$ |

Mamba

State Space Model

$n \times r \times d$

$r = n$

$n^2$

$n$

n is sequence length, d is embedding dimension.

Restricted self-attention means not attending all words in the sequence, but only a restricted field

Square complexity of sequence length is a major issue for transformers to deal with long sequence

look

$n^2$

Flash attention

26

# Auto-Encoding Variational Bayes

**Diederik P. Kingma**
Machine Learning Group
Universiteit van Amsterdam
dpkingma@gmail.com

**Max Welling**
Machine Learning Group
Universiteit van Amsterdam
welling.max@gmail.com

# Variational Autoencoders

# VAE is a Generative Model

p(z): multinomial , k
classes(e.g. uniform)

Label $Z$

$(\mu_1, \Sigma_1), (\mu_2, \Sigma_2), \ldots (\mu_k, \Sigma_k)$

Data $X$   Gaussian Mixture Model (GMM)

# The VAE Model

p(z) is a normal distribution in most cases

p(z)

Z

Data    X

# The VAE Model

p(z) is a normal distribution in most cases



p(z)

Z

Neural Networks

Data  X

# The VAE Model

p(z) is a normal distribution in most cases

p(z)

Z

Neural Networks

Data X

$$X \sim P(x, f(z; \theta))$$

# The VAE Model

p(z) is a normal distribution in most cases

p(z)

Z

Neural Networks

Data    X    $X \sim P(x, f(z; \theta))$

$f$ is a neural network taking Z as input

# Training

p(z)

Z

Neural Networks

Data    X    $X \sim P(x, f(z; \theta))$

# Training

p(z)

Z

Neural Networks

Data  X

$$X \sim P(x, f(z; \theta))$$

How to train the model? Can we do MLE?

# Training

p(z)



Z

Neural Networks

Data   X   $X \sim P(x, f(z; \theta))$

How to train the model? Can we do MLE?

Intractable P(X), EM algorithm?

# Let's try EM

p(z)

Z

Neural Networks

X

$$X \sim P(x, f(z; \theta))$$

# Let's try EM

p(z)

Z

Neural Networks

$$X \sim P(x, f(z; \theta))$$

E-Step: compute P(z|x)

$$Q(z) = P(z \mid x) \propto P(z)P(x \mid z)$$

# Let's try EM

p(z)



Neural Networks

$X \sim P(x, f(z; \theta))$

E-Step: compute P(z|x)

$$Q(z) = P(z\,|\,x) \propto P(z)P(x\,|\,z)$$

This is ok?

# Let's try EM

p(z)



Neural Networks

$X \sim P(x, f(z; \theta))$

E-Step: compute P(z|x)

$$Q(z) = P(z|x) \propto P(z)P(x|z)$$

This is ok?

M-Step: the ELBO objective

$$\text{argmax}_\theta \sum_z Q(z) \log p(x, z; \theta) = \text{argmax}_\theta \mathbb{E}_{z \sim Q(z)} \log p(x, z; \theta)$$

# Let's try EM

p(z)



Neural Networks

$X \sim P(x, f(z; \theta))$

E-Step: compute P(z|x)

$$Q(z) = P(z \mid x) \propto P(z)P(x \mid z)$$

This is ok?

M-Step: the ELBO objective

$$\mathrm{argmax}_\theta \sum_z Q(z) \log p(x, z; \theta) = \mathrm{argmax}_\theta \mathbb{E}_{z \sim Q(z)} \log p(x, z; \theta)$$

In most cases, we cannot do the sum, and cannot easily sample from Q(z) either

# Approximate Posterior

We need an easy-to-sample distribution to approximate $P(z|x)$

# Approximate Posterior

We need an easy-to-sample distribution to approximate P(z|x)

$$q(z|x; \phi) \text{ to approximate } p(z|x; \theta)$$

# Approximate Posterior

We need an easy-to-sample distribution to approximate P(z|x)

$q(z|x; \phi)$ to approximate $p(z|x; \theta)$    Why conditioned on x?

# Approximate Posterior

We need an easy-to-sample distribution to approximate P(z|x)

$q(z|x; \phi)$ to approximate $p(z|x; \theta)$    Why conditioned on x?

$\phi$ is the parameter for the approximate function, $\theta$ is the generative model parameter

# Approximate Posterior

We need an easy-to-sample distribution to approximate P(z|x)

$q(z|x; \phi)$ to approximate $p(z|x; \theta)$     Why conditioned on x?

$\phi$ is the parameter for the approximate function, $\theta$ is the generative model parameter

How to train $q(z|x; \phi)$, what would be the loss to find $\phi$?

# Recap: ELBO

$$\text{ELBO}(x; Q, \theta) = \sum_z Q(z) \log \frac{p(x, z; \theta)}{Q(z)}$$

What is $\text{argmax}_{Q(z)} \text{ELBO}(x; Q, \theta)$?

# Recap: ELBO

$$\mathrm{ELBO}(x; Q, \theta) = \sum_z Q(z) \log \frac{p(x, z; \theta)}{Q(z)}$$

What is $\mathrm{argmax}_{Q(z)}\mathrm{ELBO}(x; Q, \theta)$?

ELBO is maximized when Q(z) is equal to p(z|x)

# Recap: ELBO

$$\text{ELBO}(x; Q, \theta) = \sum_z Q(z) \log \frac{p(x, z; \theta)}{Q(z)}$$

What is $\text{argmax}_{Q(z)}\text{ELBO}(x; Q, \theta)$?

ELBO is maximized when Q(z) is equal to p(z|x)

Therefore, we can approximate the true posterior by maximizing ELBO:

$$\text{argmax}_\phi \sum_z q(z|x; \phi) \log \frac{p(x, z; \theta)}{q(z|x; \phi)}$$

# Recap: ELBO

$$\text{ELBO}(x; Q, \theta) = \sum_z Q(z) \log \frac{p(x, z; \theta)}{Q(z)}$$

What is $\text{argmax}_{Q(z)} \text{ELBO}(x; Q, \theta)$?

ELBO is maximized when Q(z) is equal to p(z|x)

Therefore, we can approximate the true posterior by maximizing ELBO:

$$\text{argmax}_{\phi} \sum_z q(z|x; \phi) \log \frac{p(x, z; \theta)}{q(z|x; \phi)}$$

Variational Inference

# **Training VAEs**

E-Step:

$$\text{argmax}_\phi \sum_z q(z|x;\phi)\log \frac{p(x,z;\theta)}{q(z|x;\phi)}$$

M-Step:

$$\text{argmax}_\theta \sum_z q(z|x;\phi)\log \frac{p(x,z;\theta)}{q(z|x;\phi)}$$

# Training VAEs

E-Step:

$$\text{argmax}_\phi \sum_z q(z|x;\phi)\log \frac{p(x,z;\theta)}{q(z|x;\phi)}$$

M-Step:

$$\text{argmax}_\theta \sum_z q(z|x;\phi)\log \frac{p(x,z;\theta)}{q(z|x;\phi)}$$

Same objective, different parameters to optimize

# Training VAEs

E-Step:

$$\text{argmax}_\phi \sum_z q(z|x;\phi)\log \frac{p(x,z;\theta)}{q(z|x;\phi)}$$

M-Step:

$$\text{argmax}_\theta \sum_z q(z|x;\phi)\log \frac{p(x,z;\theta)}{q(z|x;\phi)}$$

Same objective, different parameters to optimize

Because we use approximate rather than exact posterior, it is also called Variational EM

# Training VAEs

E-Step:

$$\text{argmax}_\phi \sum_z q(z|x;\phi)\log \frac{p(x,z;\theta)}{q(z|x;\phi)}$$

M-Step:

$$\text{argmax}_\theta \sum_z q(z|x;\phi)\log \frac{p(x,z;\theta)}{q(z|x;\phi)}$$

# **Training VAEs**

E-Step:

$$\text{argmax}_{\phi} \sum_z q(z|x;\phi)\log \frac{p(x,z;\theta)}{q(z|x;\phi)}$$

M-Step:

$$\text{argmax}_{\theta} \sum_z q(z|x;\phi)\log \frac{p(x,z;\theta)}{q(z|x;\phi)}$$

We use MC sampling to approximate expectation
and use gradient descent to optimize $\theta$

# **Training VAEs**

E-Step:

$$\text{argmax}_{\phi} \sum_{z} q(z|x;\phi)\log \frac{p(x,z;\theta)}{q(z|x;\phi)}$$

Can we do gradient descent over $\phi$?

M-Step:

$$\text{argmax}_{\theta} \sum_{z} q(z|x;\phi)\log \frac{p(x,z;\theta)}{q(z|x;\phi)}$$

We use MC sampling to approximate expectation and use gradient descent to optimize $\theta$