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

# Transformers, VAEs

Junxian He

Apr 19, 2024

# Transformer



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

# Encoder

Output
Probabilities

Softmax

Linear

Add & Norm

Feed
Forward

Add & Norm

Multi-Head
Attention

Nx

Add & Norm

Masked
Multi-Head
Attention

Positional
Encoding

Output
Embedding

Outputs
(shifted right)

*encoder*

Add & Norm

Feed
Forward

Nx

Add & Norm

Multi-Head
Attention

Positional
Encoding

Input
Embedding

Inputs

3

# Decoder



Output
Probabilities

Softmax

Linear

*decoder*

Add & Norm

Feed
Forward

Add & Norm

Multi-Head
Attention

Add & Norm

Masked
Multi-Head
Attention

N×

*encoder*

Add & Norm

Feed
Forward

N×

Add & Norm

Multi-Head
Attention

Positional
Encoding

Input
Embedding

Inputs

Positional
Encoding

Output
Embedding

Outputs
(shifted right)

4

# Transformer Encoder



MLP

Residual connection

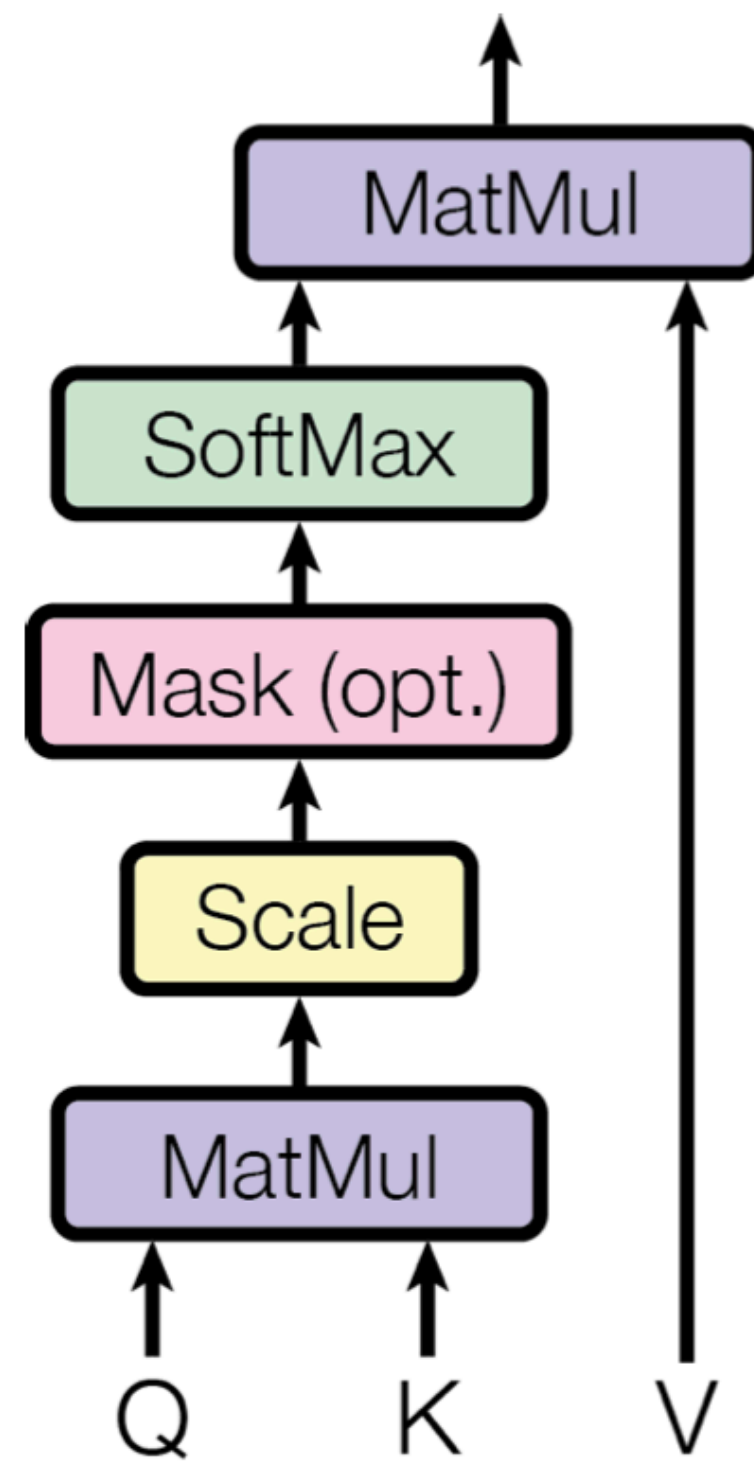Self-attention

Positional Encoding

Input Embedding

N×

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

Dot-products grow large in magnitude
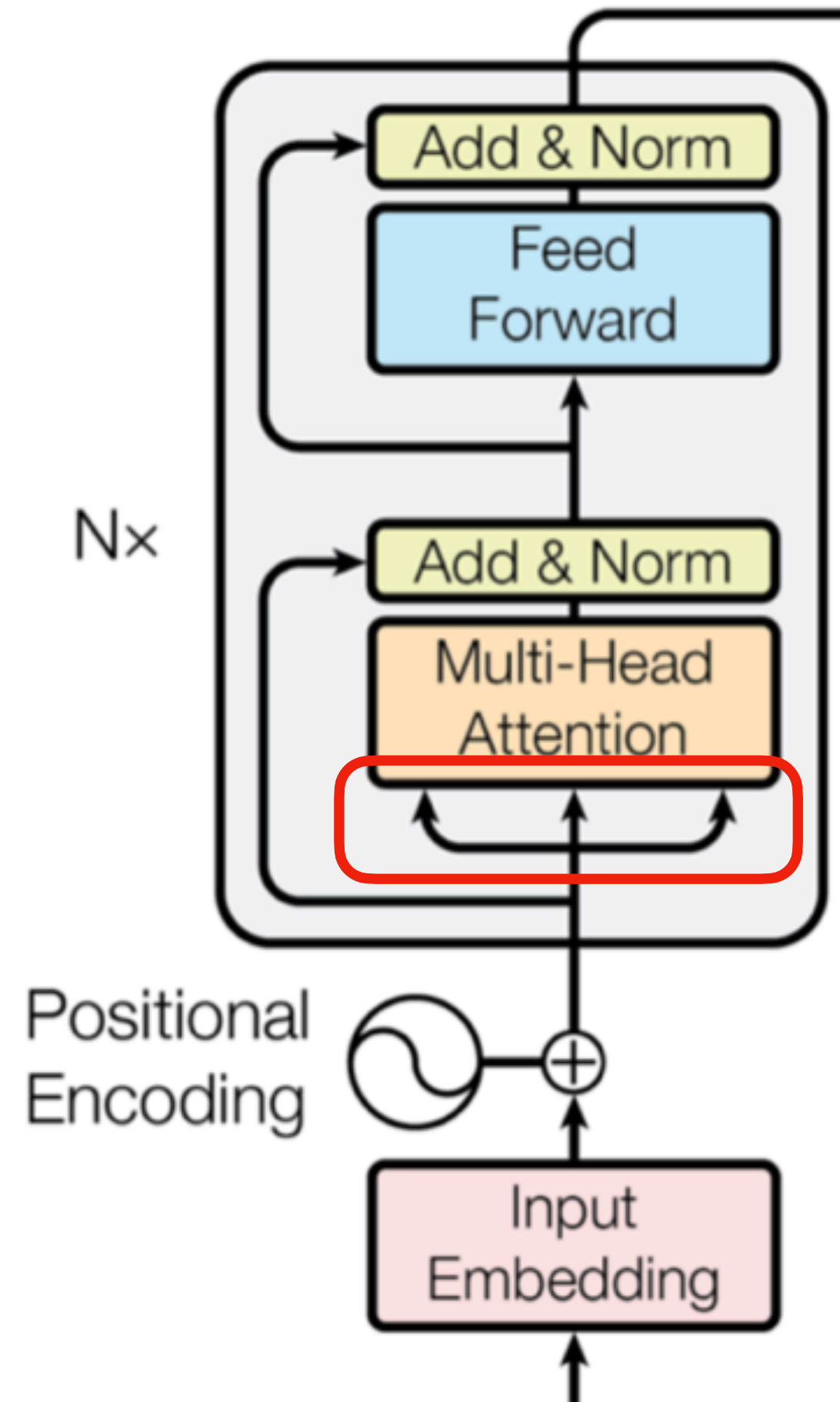
Scaled Attention weight = $\text{softmax}(\dfrac{QK^T}{\sqrt{d_k}})$  Shape is mxn

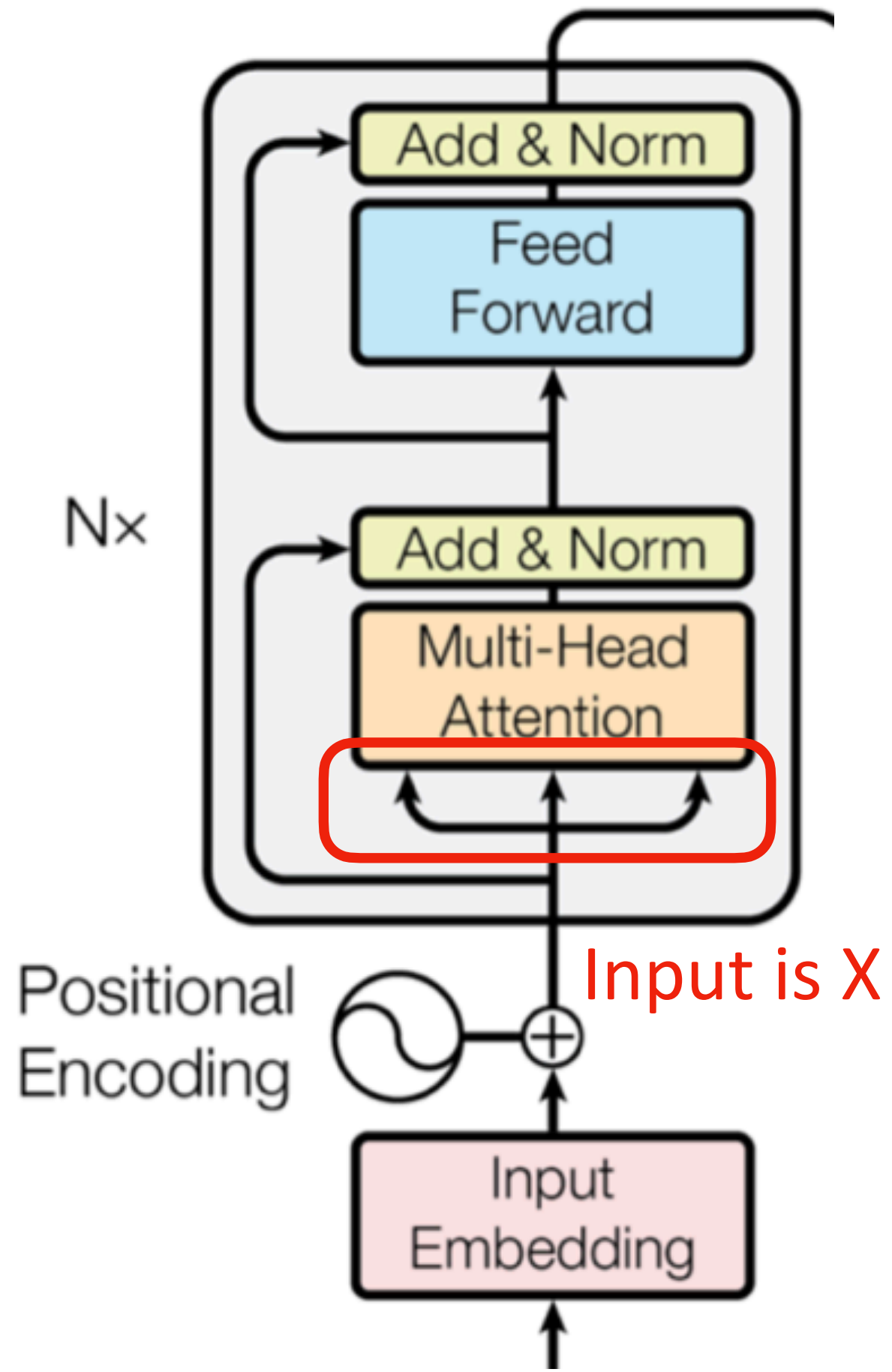Attention weight represents the strength to "attend" values V

$$\text{Attention}(Q, K, V) = \text{softmax}(\frac{QK^T}{\sqrt{d_k}})V$$

Q: Query
K: key
V: value

6

# Q, K, V



What are Q, K, V in the transformer

# Self-Attention



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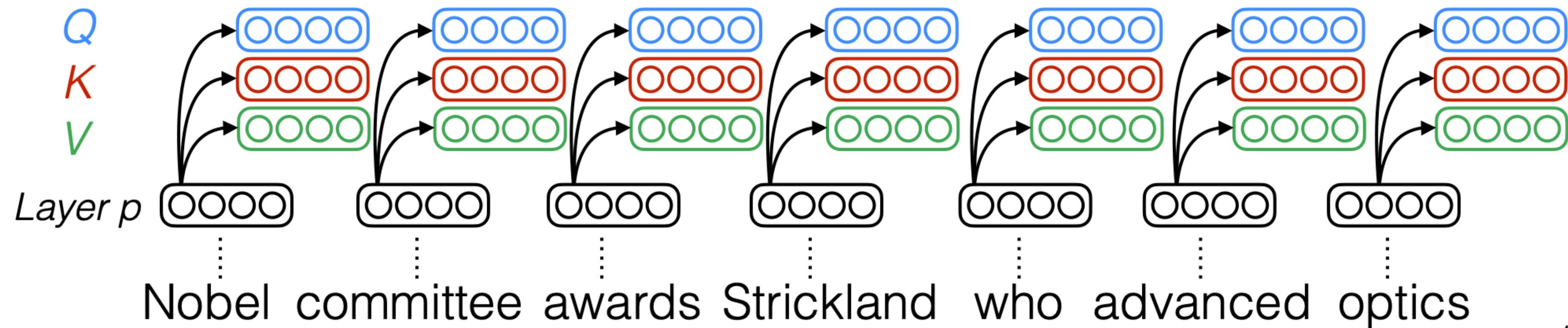
Query, key, and value are from the same input, thus it is called "self"-attention

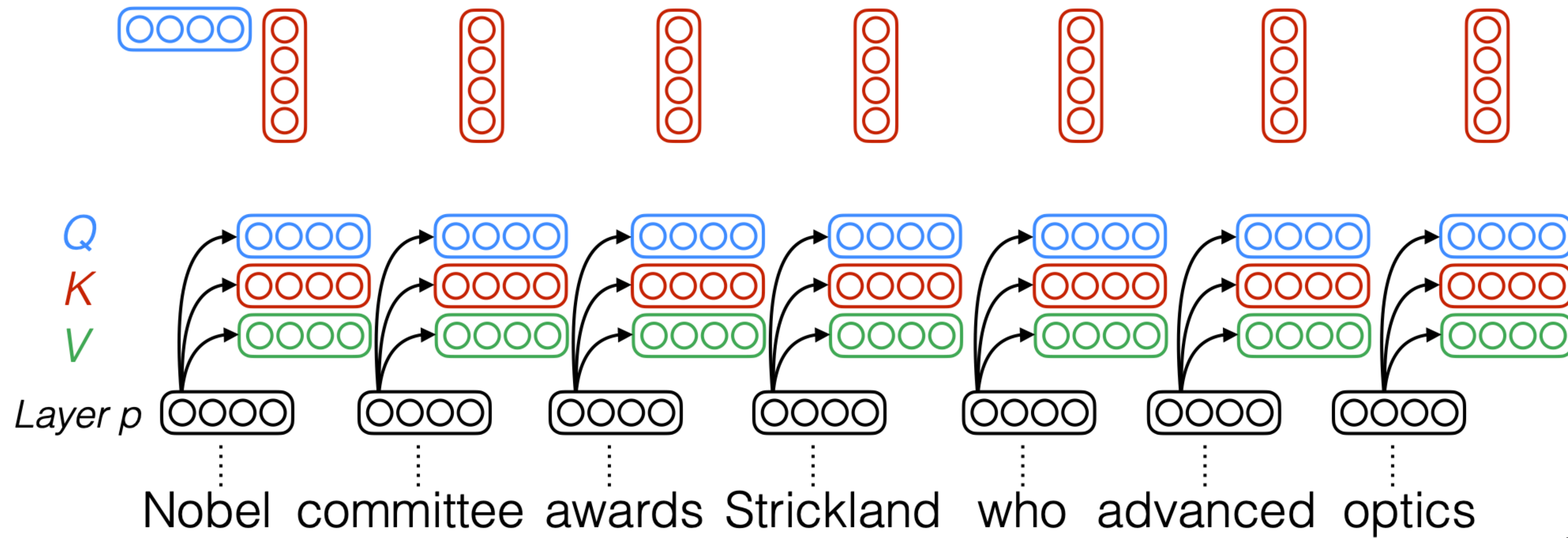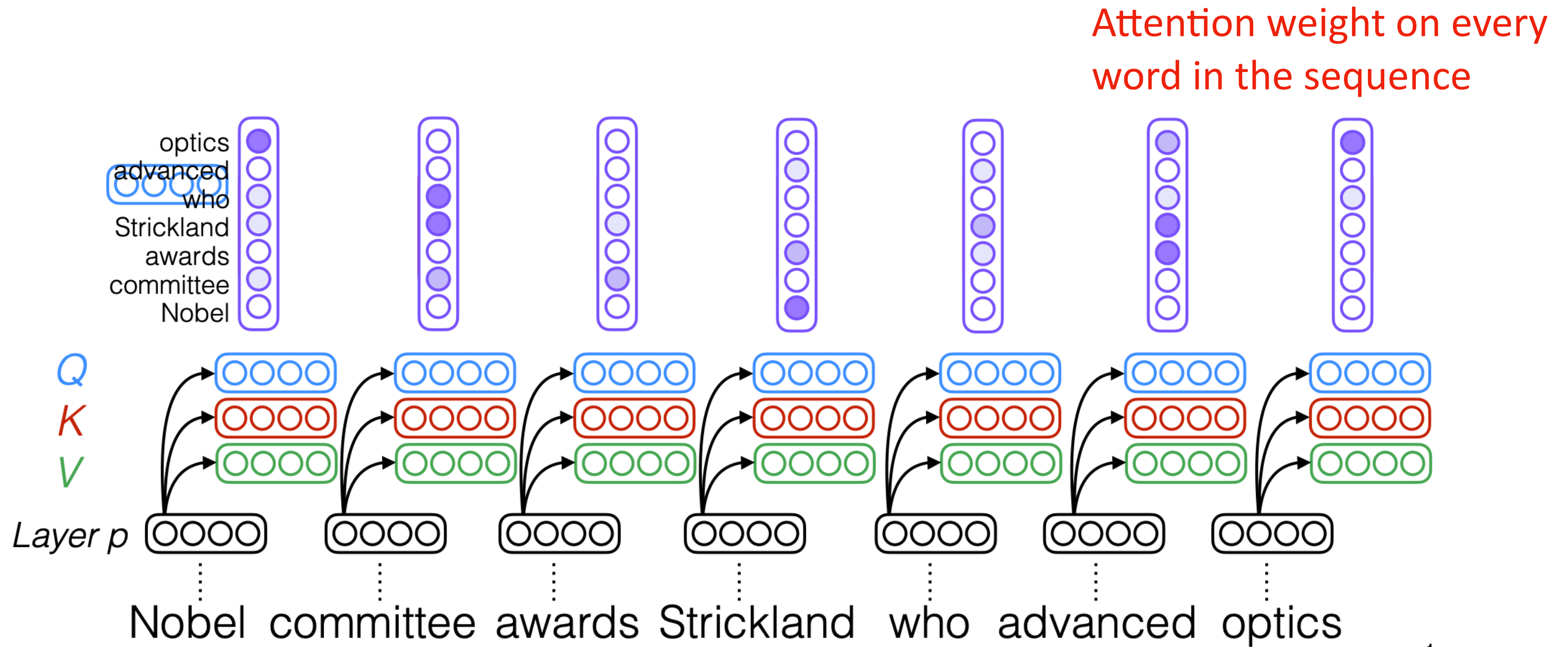$$\text{softmax}\left(\frac{Q \times K^T}{\sqrt{d_k}}\right) V = Z$$

8

8

Jay Alammar. The Illustrated Transformer.

# Self-Attention

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



9

# Self-Attention



Q
K
V

Layer p

Nobel committee awards Strickland who advanced optics

# Self-Attention
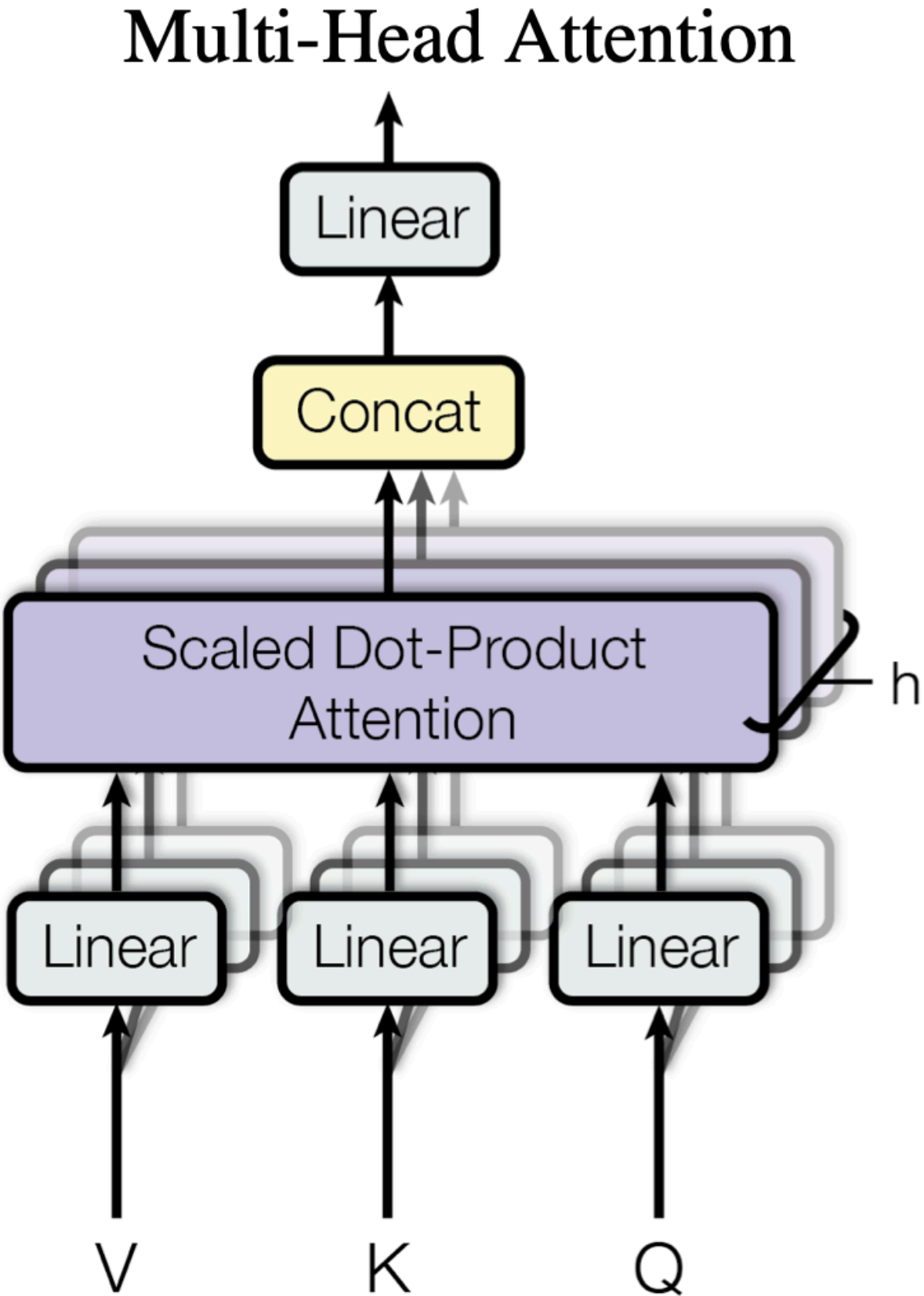
Attention weight on every word in the sequence
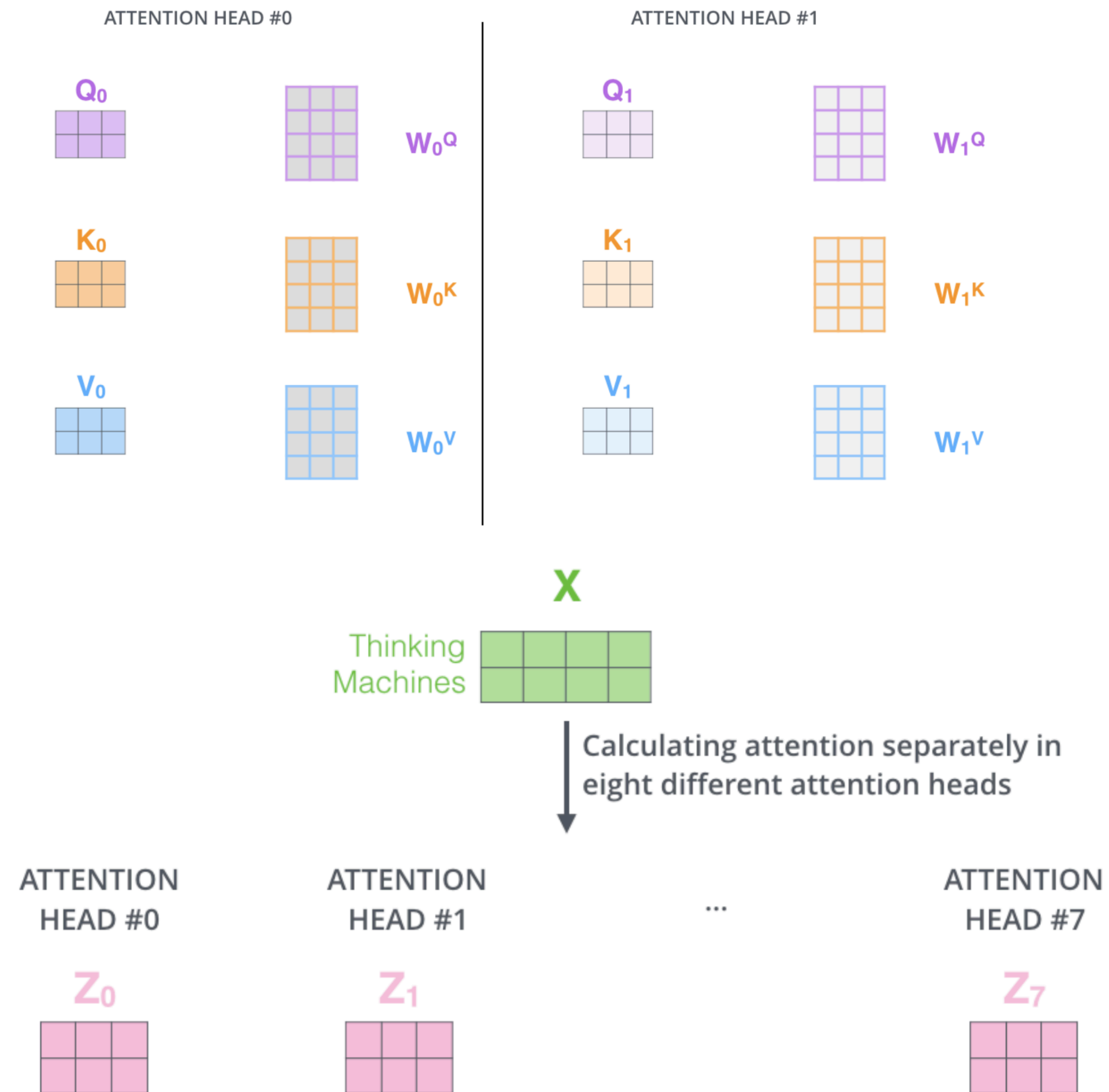


11

# Self-Attention

# Self-Attention

# Multi-Head Attention



Multi-Head Attention

# Multi-Head Self-Attention

ATTENTION HEAD #0                    ATTENTION HEAD #1

$Q_0$          $W_0^Q$          $Q_1$          $W_1^Q$

$K_0$          $W_0^K$          $K_1$          $W_1^K$

$V_0$          $W_0^V$          $V_1$          $W_1^V$

$X$

Thinking
Machines

Calculating attention separately in
eight different attention heads

ATTENTION          ATTENTION                    ATTENTION
HEAD #0            HEAD #1            ...        HEAD #7

$Z_0$              $Z_1$                        $Z_7$

Jay Alammar. The Illustrated Transformer.

# Multi-Head Self-Attention

1) Concatenate all the attention heads

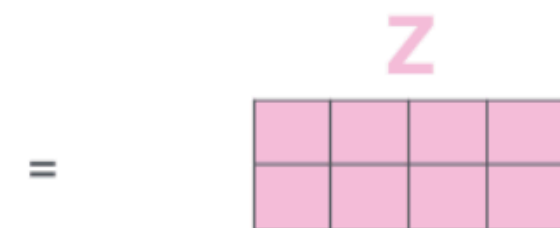$Z_0$  $Z_1$  $Z_2$  $Z_3$  $Z_4$  $Z_5$  $Z_6$  $Z_7$



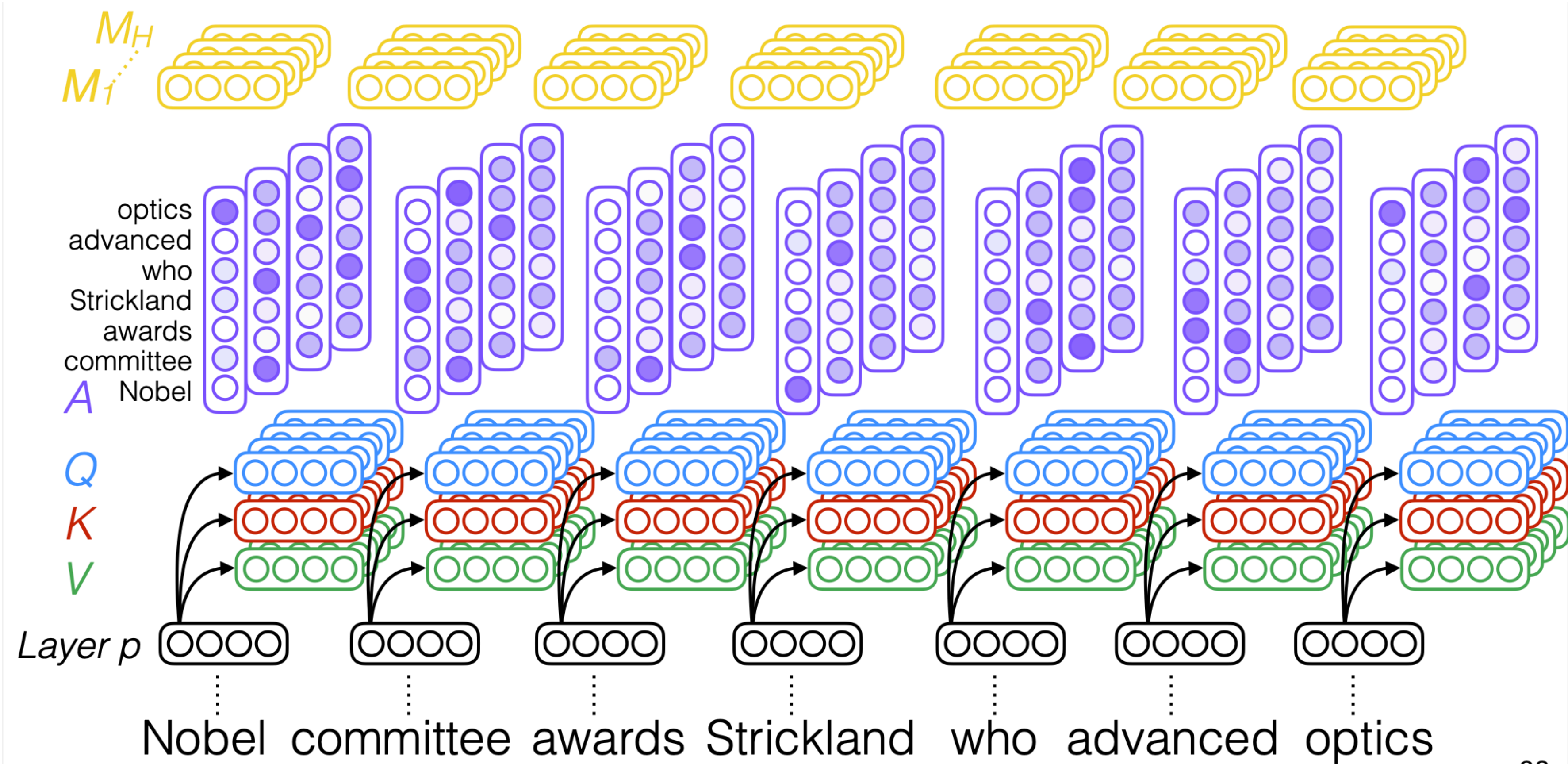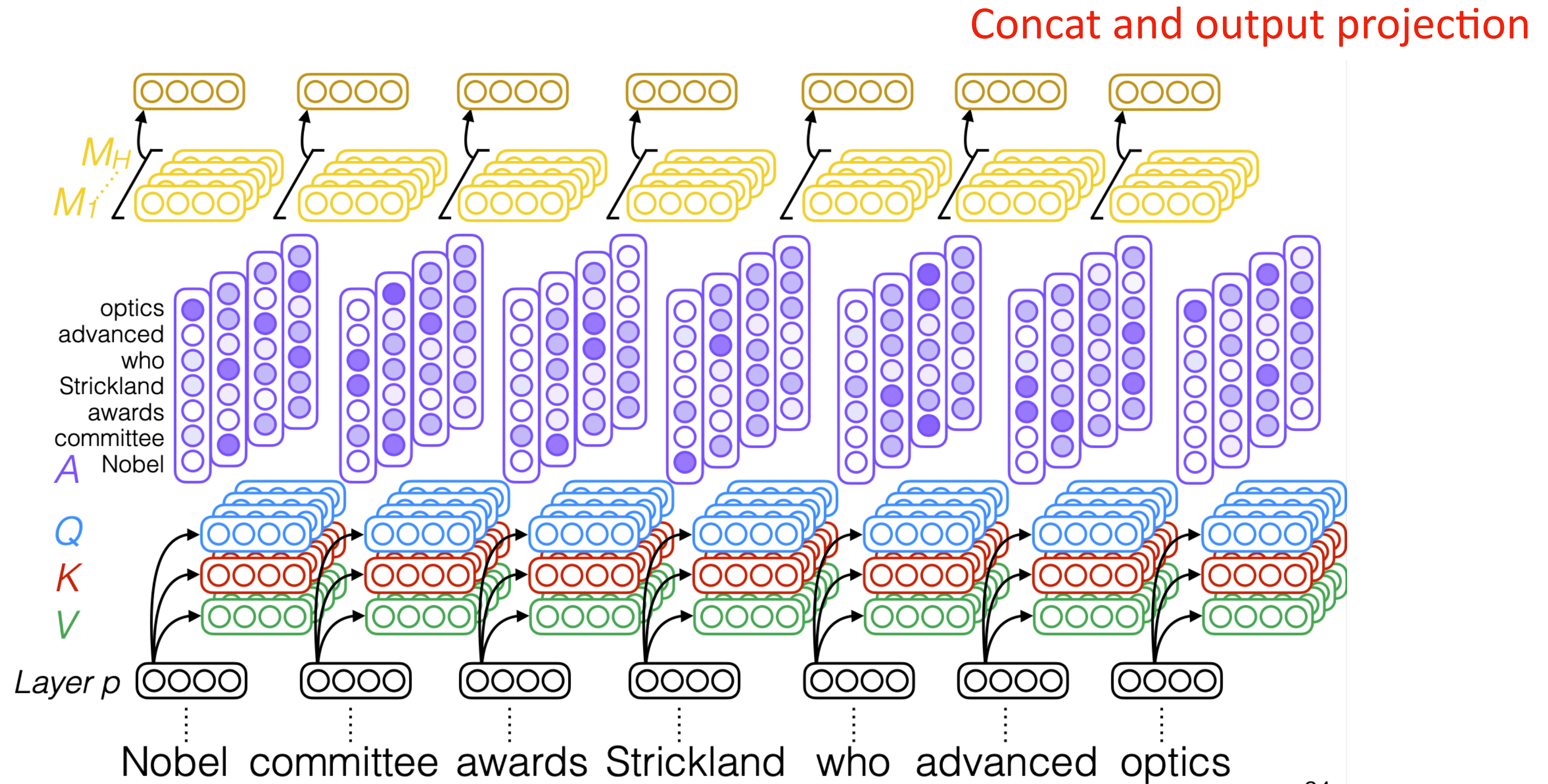2) Multiply with a weight matrix $W^O$ that was trained jointly with the model

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
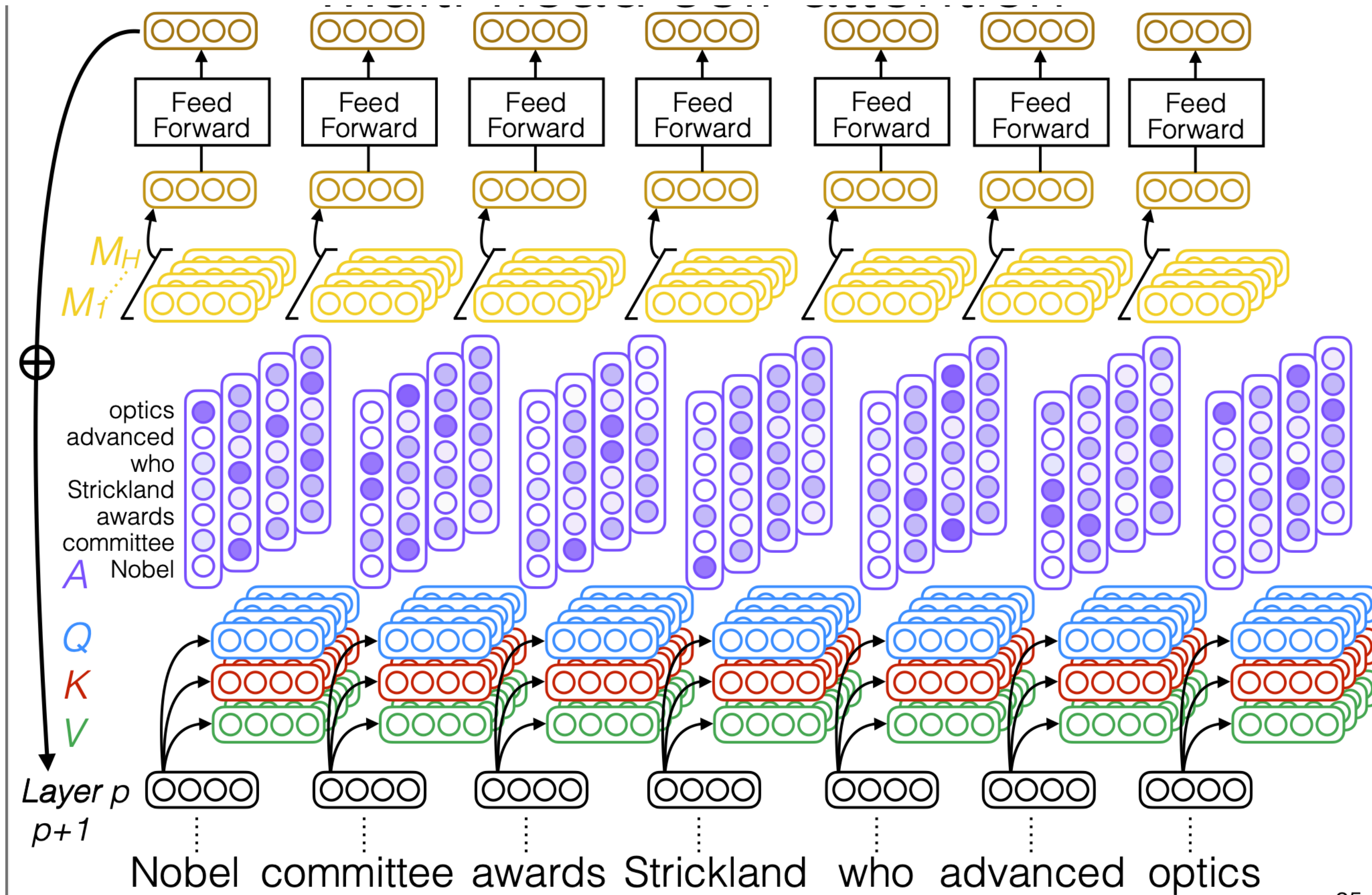
Z

=

Jay Alammar. The Illustrated Transformer.

# Multi-head Self-Attention

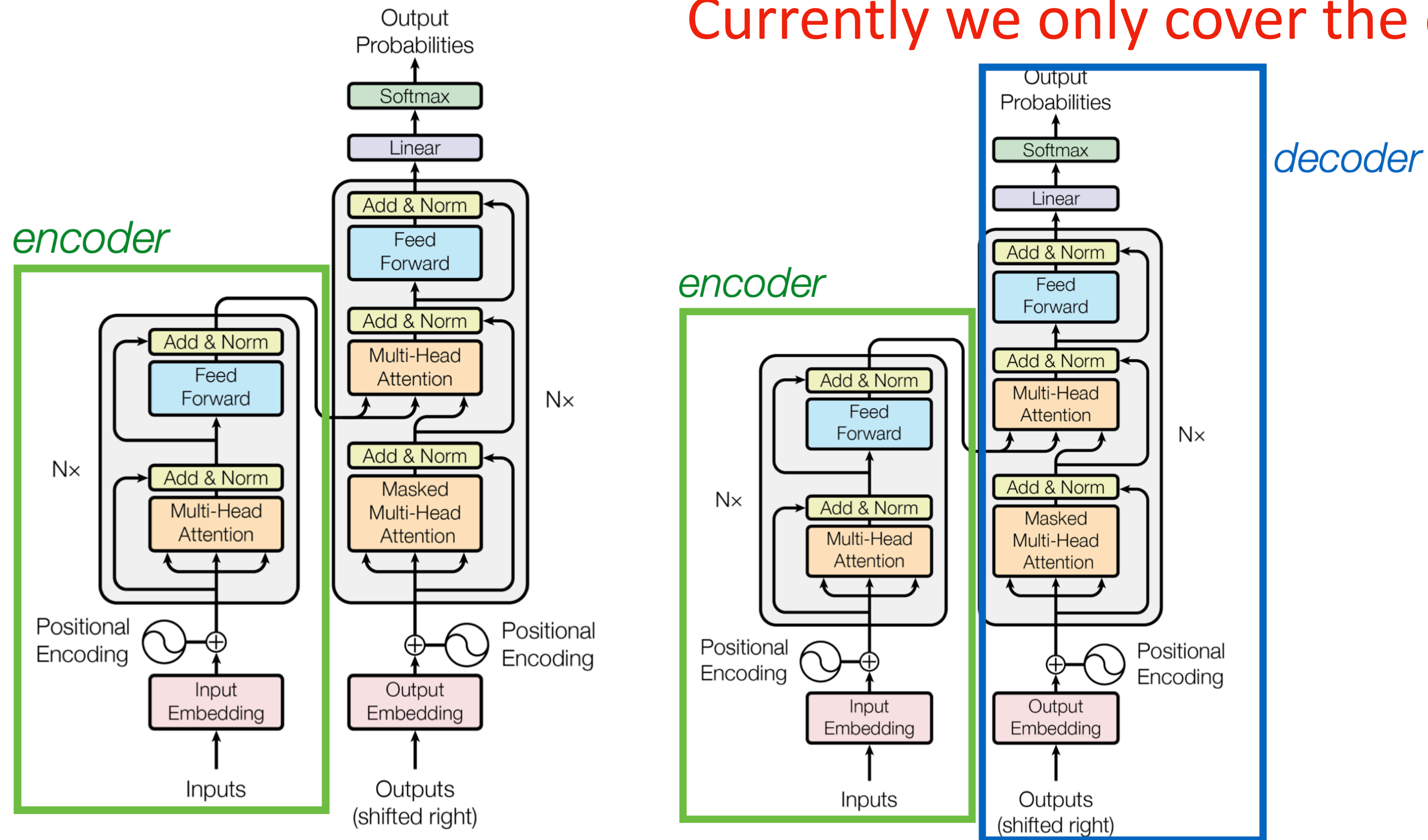# Multi-head Self-Attention

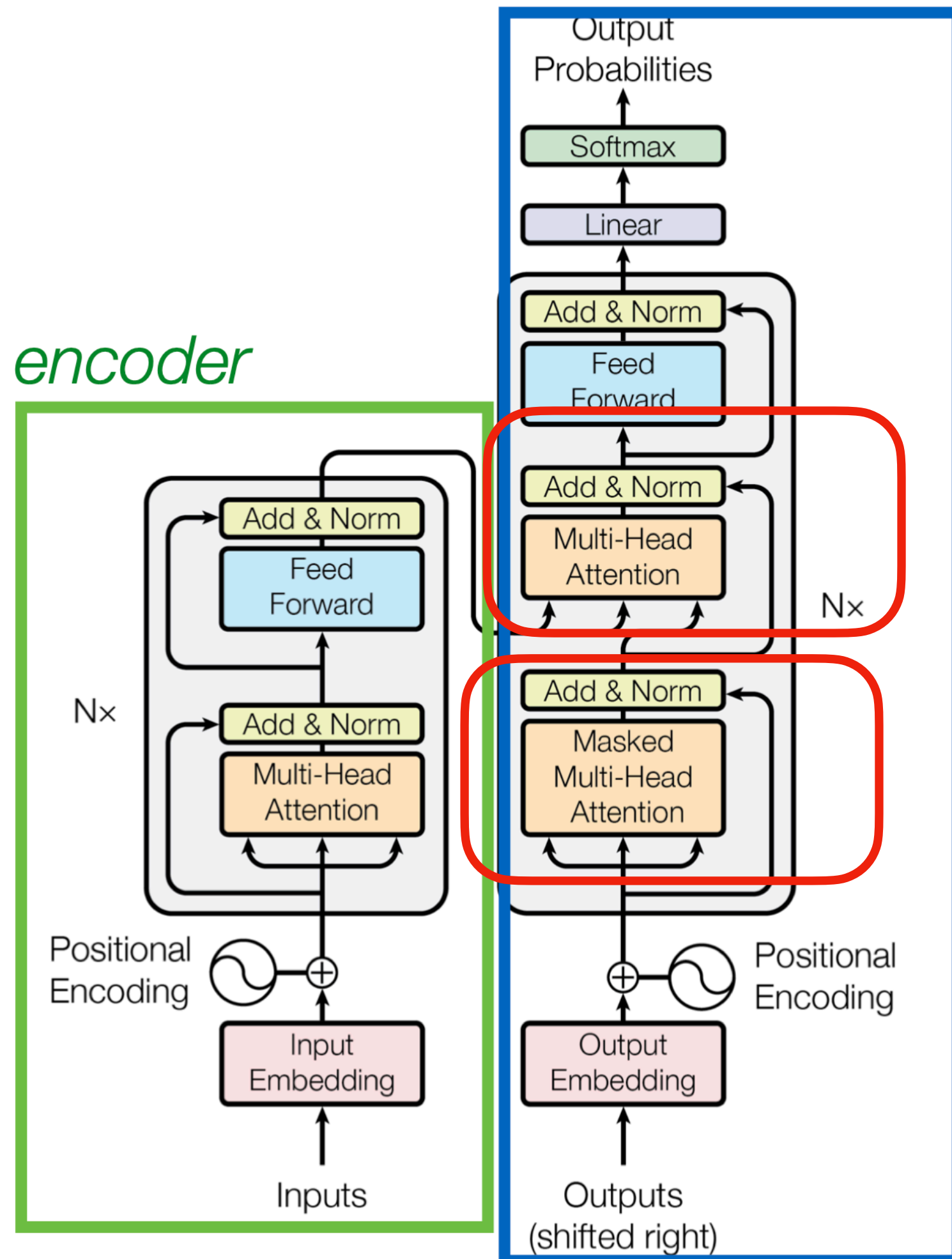# Multi-head Self-Attention + FFN



19

# 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
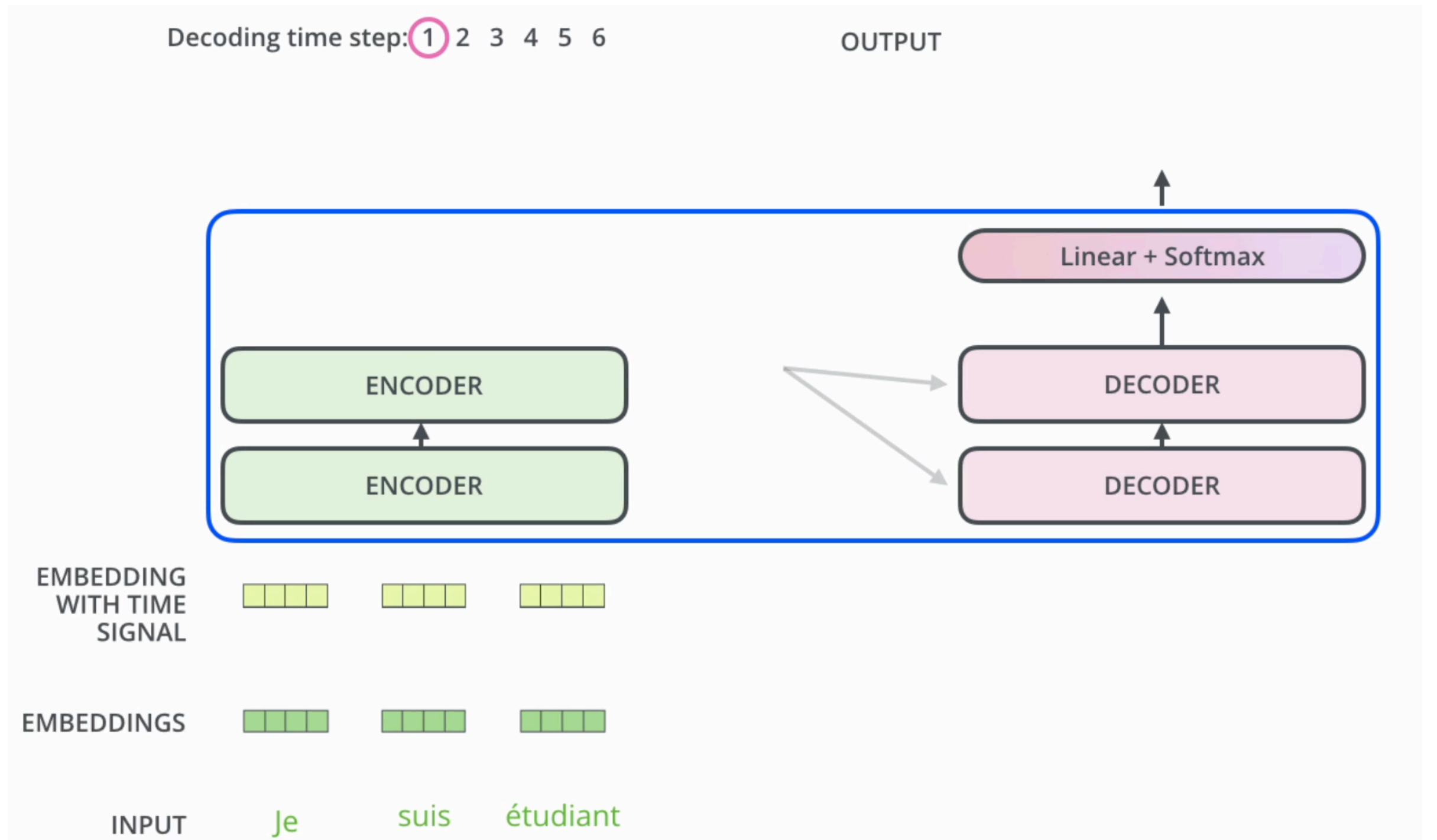
20

# Transformer Decoder in Seq2Seq



Cross-attention

Self-attention

Cross-attention uses the output of encoder as input

# Masked Attention

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

*encoder*

*decoder*

# 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

# 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

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

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

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

# 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

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

# Auto-Encoding Variational Bayes

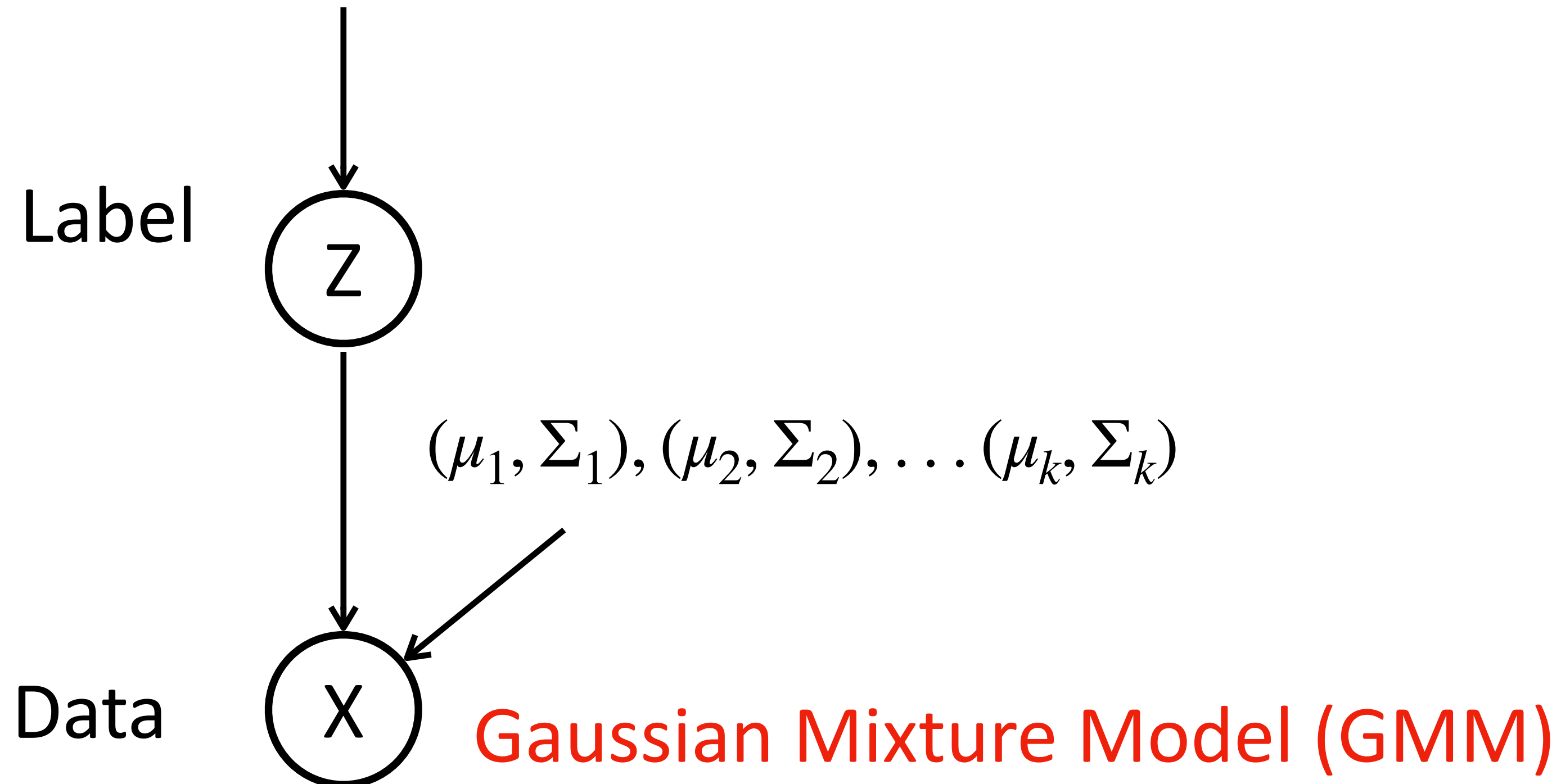**Diederik P. Kingma**
Machine Learning Group
Universiteit van Amsterdam
dpkingma@gmail.com

**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

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

How to train the model? Can we do MLE?

Intractable P(X), EM algorithm?

# 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

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

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

31

# Approximate Posterior

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

$$q(z|x; \phi) \text{ to approximate } p(z|x; \theta)$$ 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 $\text{argmax}_{Q(z)} \text{ELBO}(x; Q, \theta)$?

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:

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

Same objective, different parameters to optimize

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

# Training VAEs

E-Step:

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

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 $\theta$