

Lecture 1

Introduction to Machine Learning

FABIAN RUEHLE (NORTHEASTERN UNIVERSITY & IAIFI)

ML in Maths and Physics 2023
University of Oxford

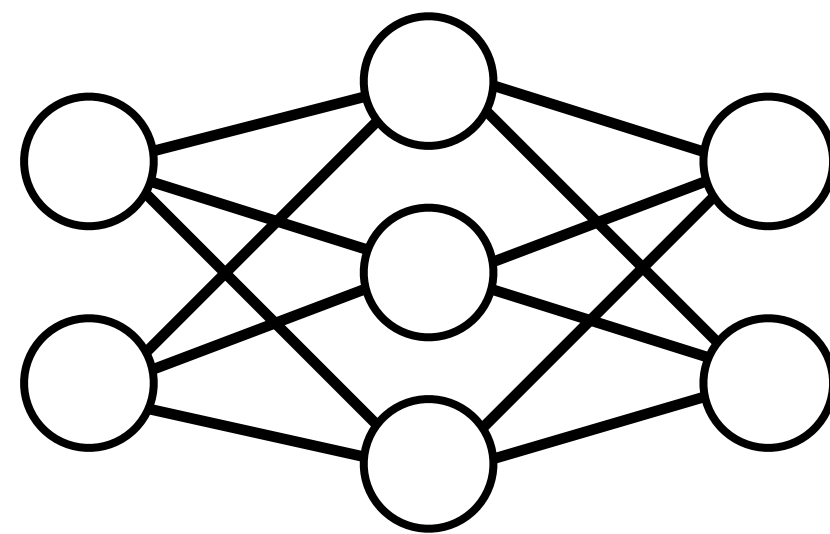
17-21 July 2023



What you will learn in the next 90 min

1

Introduction to NNs



3

The transformer architecture

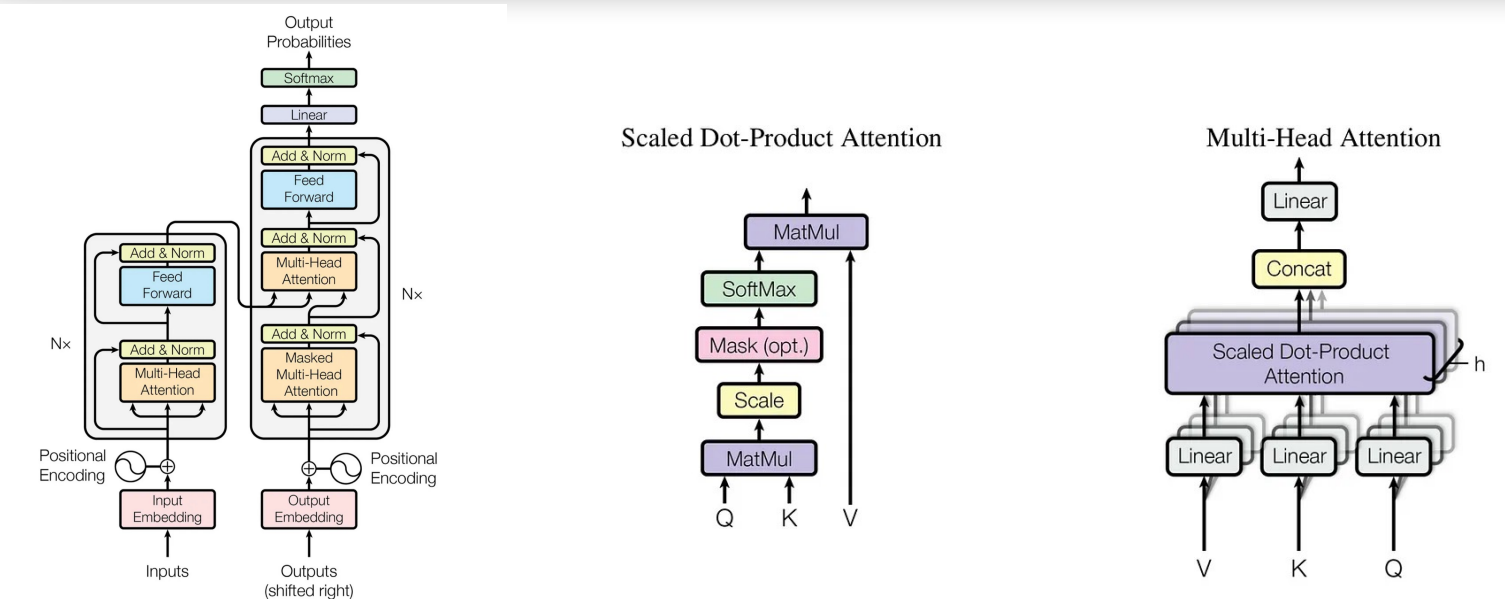
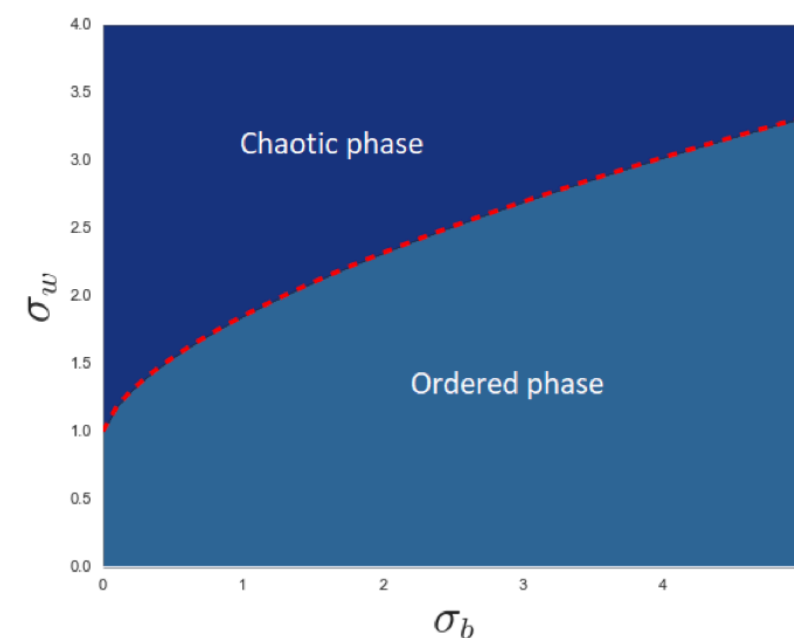


Figure 1: The Transformer - model architecture.

2

Training with backpropagation

$$\theta^{(i)} \rightarrow \theta^{(i)} - \alpha \frac{\partial L}{\partial \theta^{(i)}}$$



4

Recap

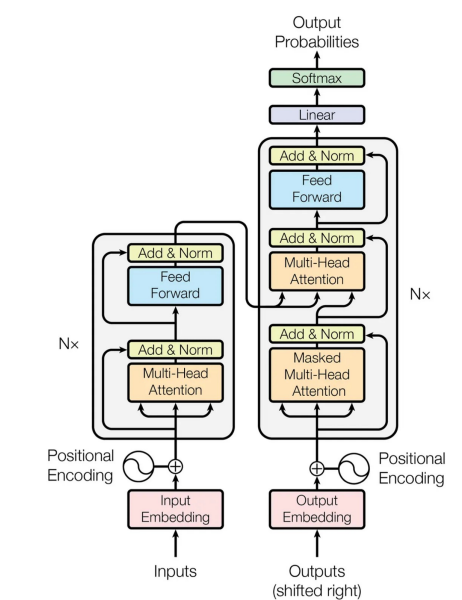
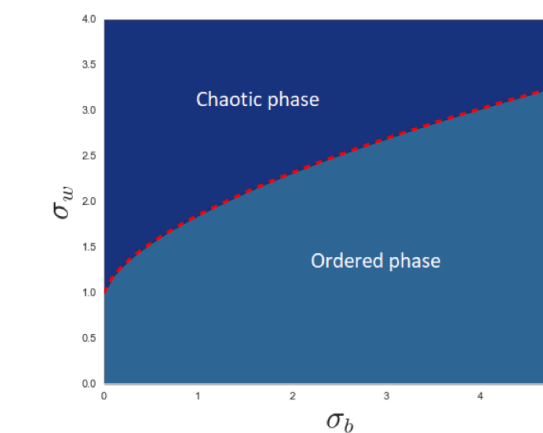
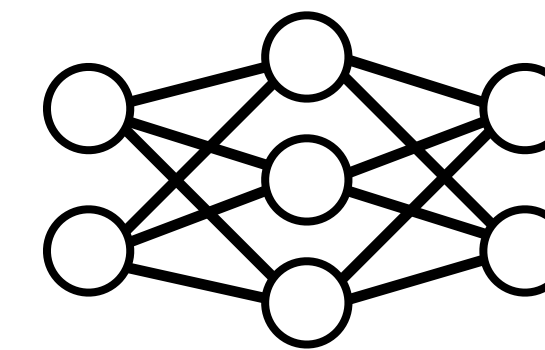
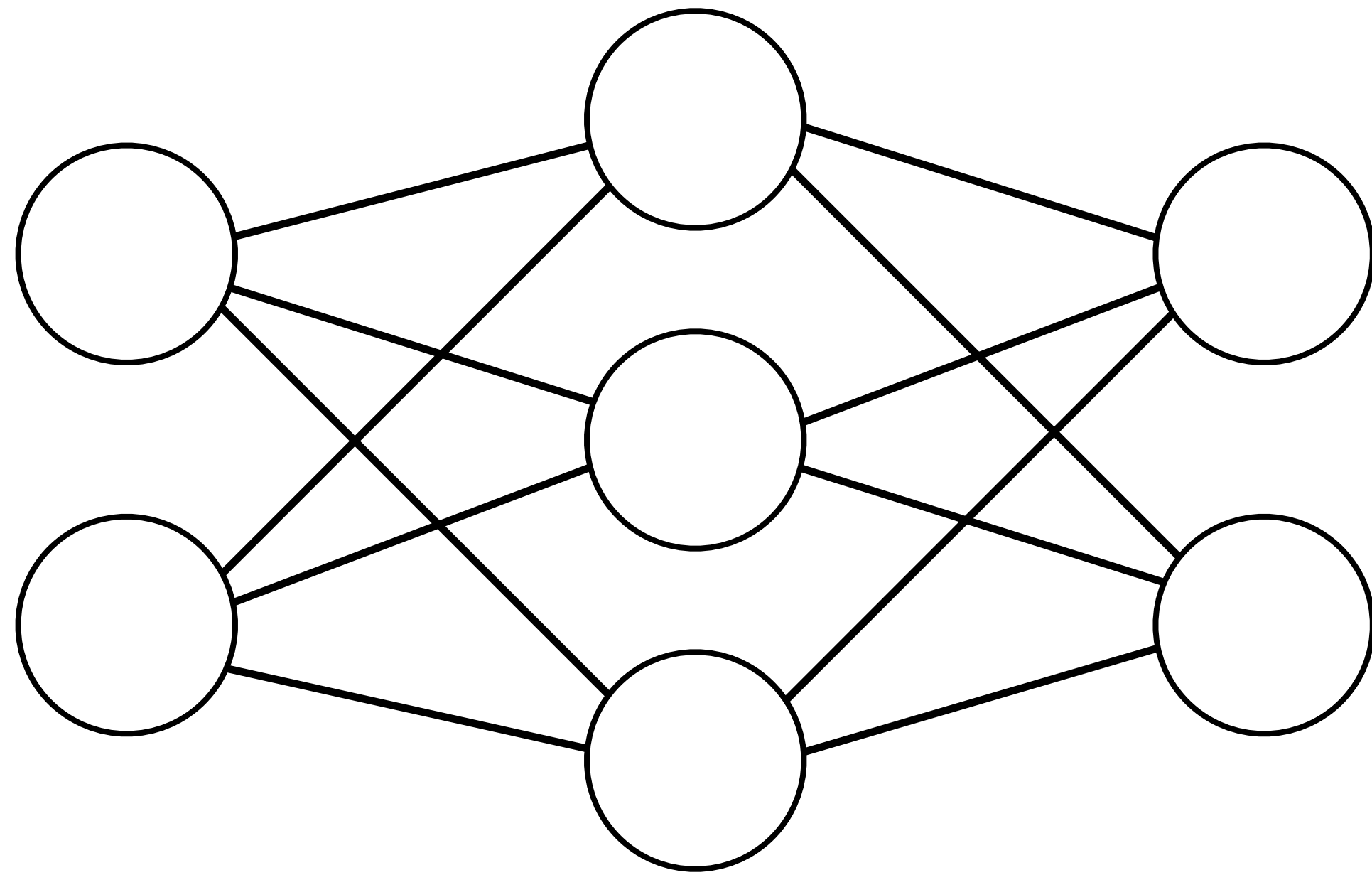


Figure 1: The Transformer - model architecture.



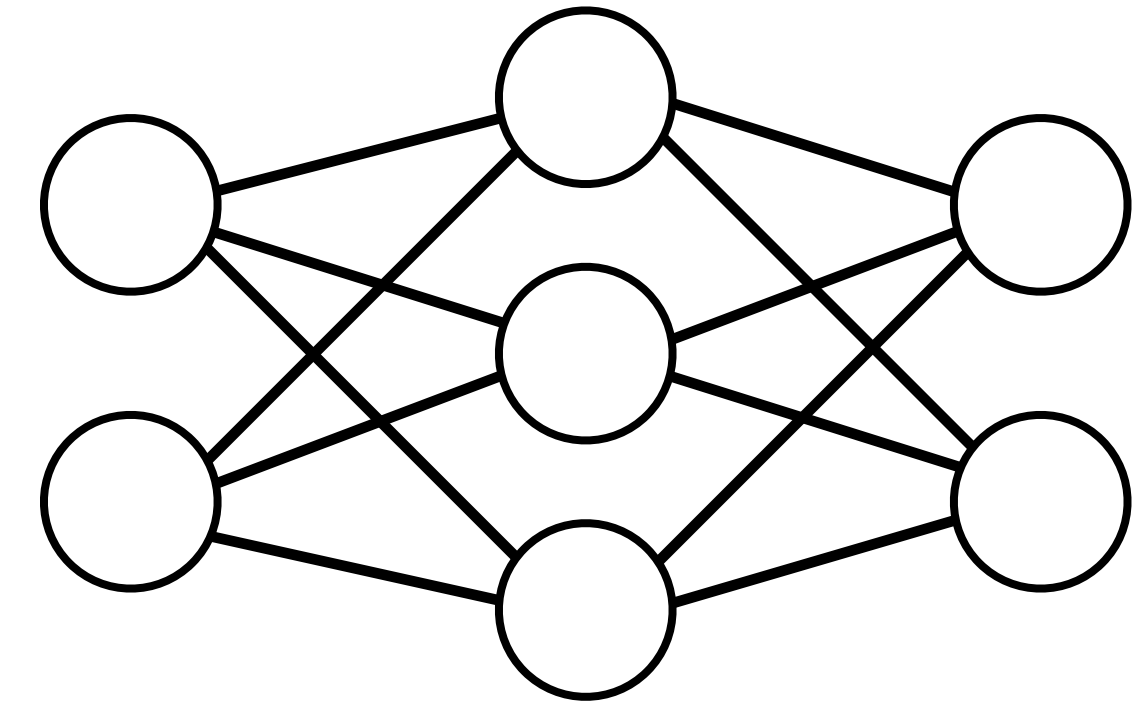
Intro to Neural Networks

Neural Networks

▶ NN is a family of maps $f_\theta : \mathbb{R}^{n_{\text{in}}} \rightarrow \mathbb{R}^{n_{\text{out}}}$ parameterized by parameters θ

▶ The map f_θ is a concatenation of

- affine transformations $z_j \rightarrow W_j^i z_i + b_j$
- element-wise non-linear functions $z_j \rightarrow \sigma(z_j)$



$$\begin{pmatrix} z_0 \\ z_1 \end{pmatrix} \mapsto \begin{pmatrix} w_{1,1}^1 & w_{1,2}^1 & w_{1,3}^1 \\ w_{2,1}^1 & w_{2,2}^1 & w_{2,3}^1 \end{pmatrix} \cdot \begin{pmatrix} \tanh(w_{1,1}^0 z_1 + w_{1,2}^0 z_2 + b_1^0) \\ \tanh(w_{2,1}^0 z_1 + w_{2,2}^0 z_2 + b_2^0) \\ \tanh(w_{3,1}^0 z_1 + w_{3,2}^0 z_2 + b_3^0) \end{pmatrix} + \begin{pmatrix} b_1^1 \\ b_2^1 \end{pmatrix}$$

▶ Compare with

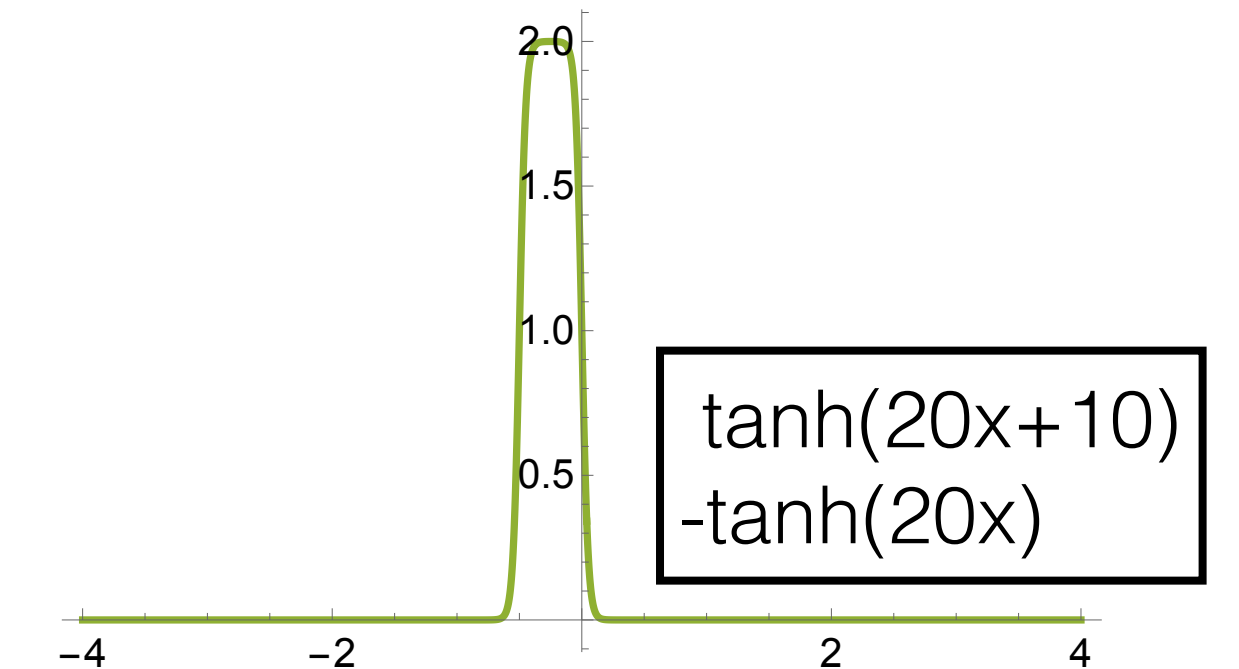
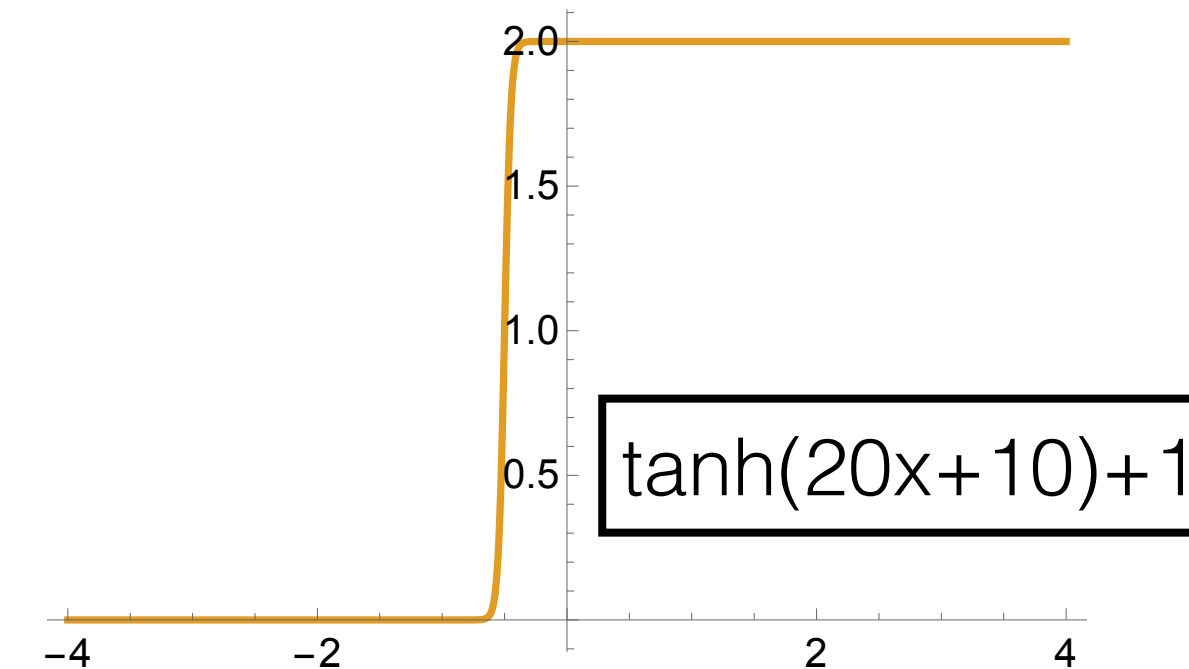
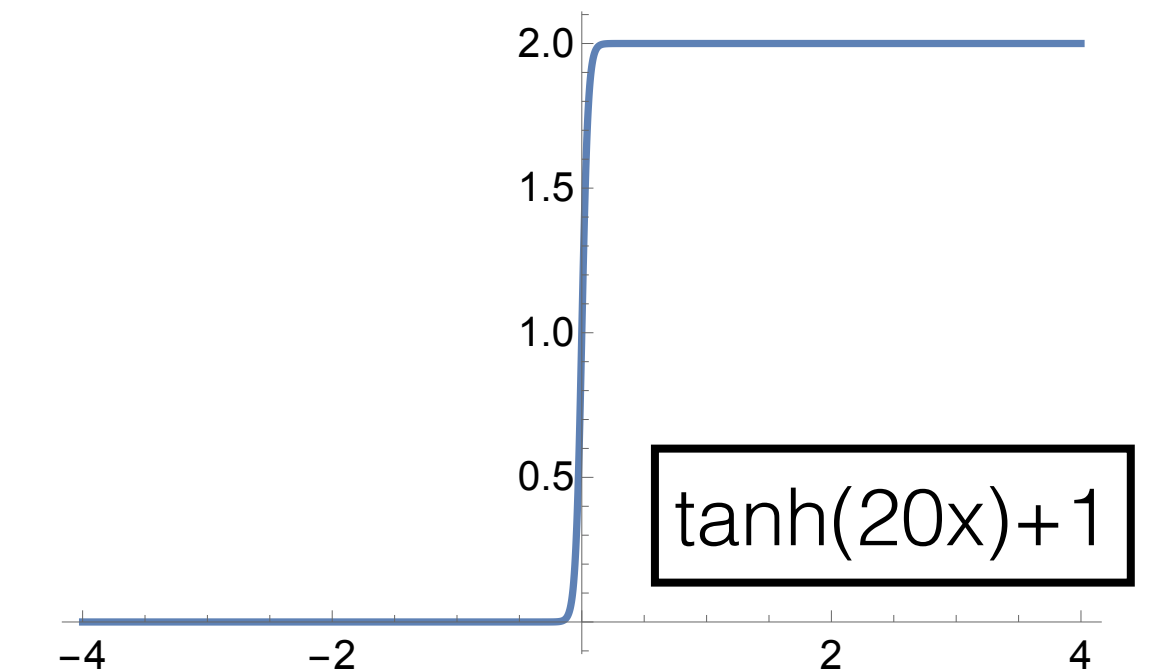
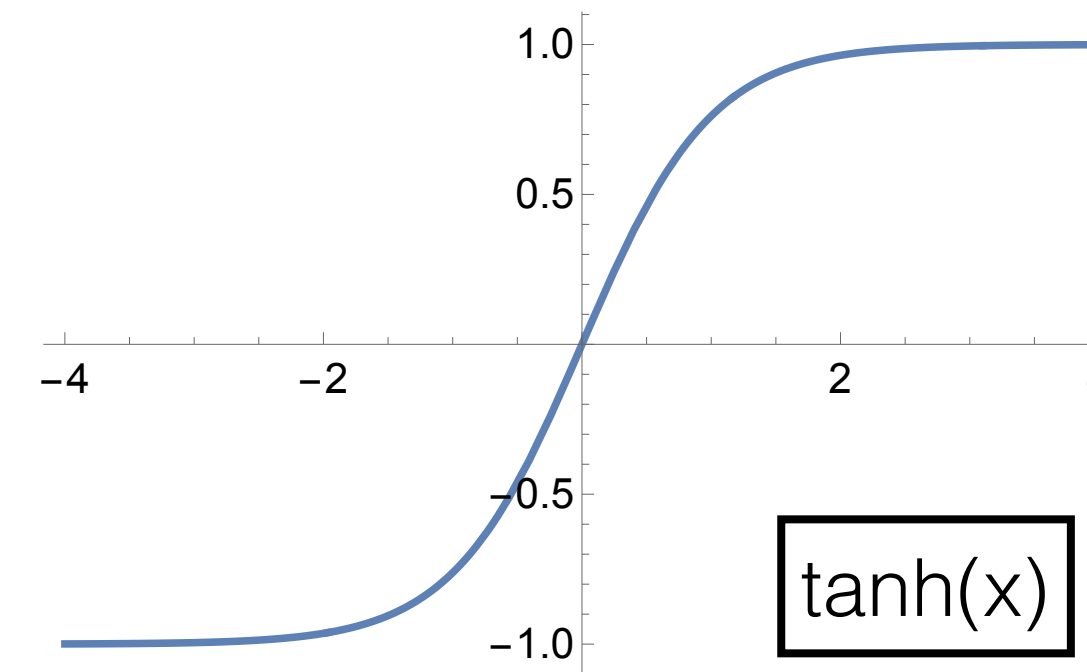
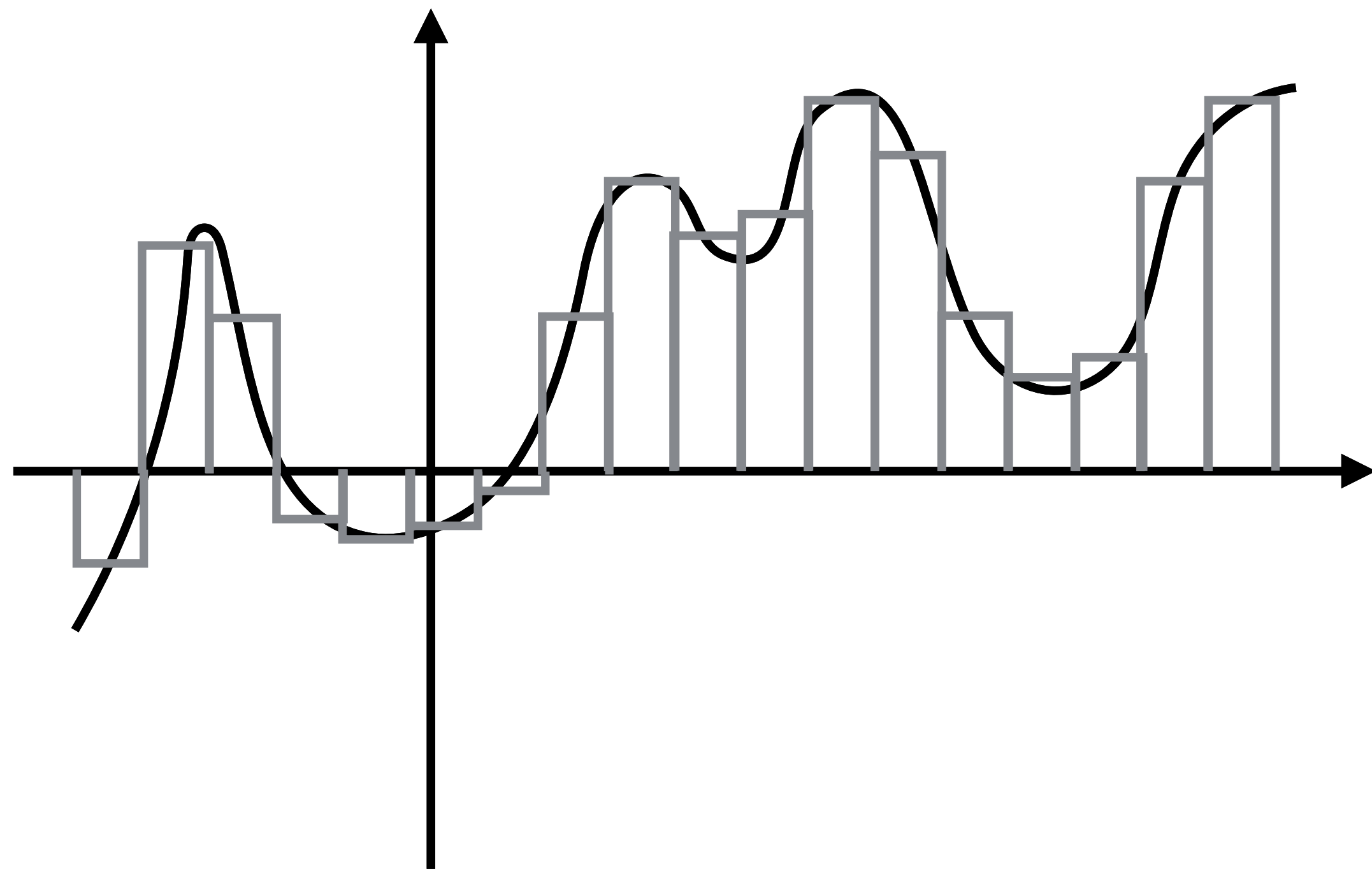
- Taylor Expansion $f(z) = \sum_i w_i (z - b_i)^i$ with $w_i = \partial^n f(x)|_{x=b}$
- Fourier Expansion $f(z) = \sum_i w_i^1 \cos(w_i^0 z + b_i)$

Neural Networks as function approximations

- ▶ So why would one ever do this?

Theorem: NNs are universal function approximators

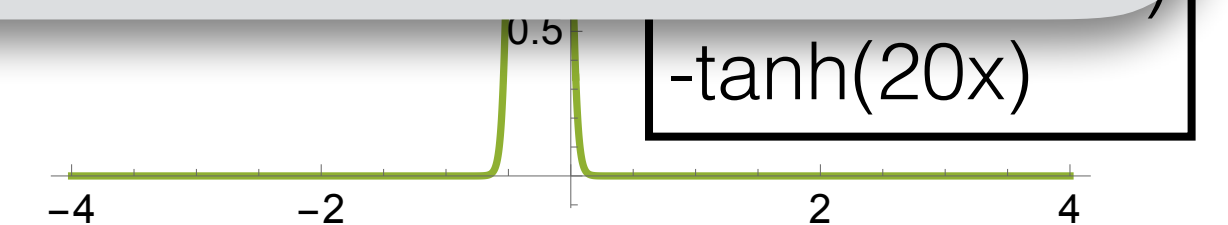
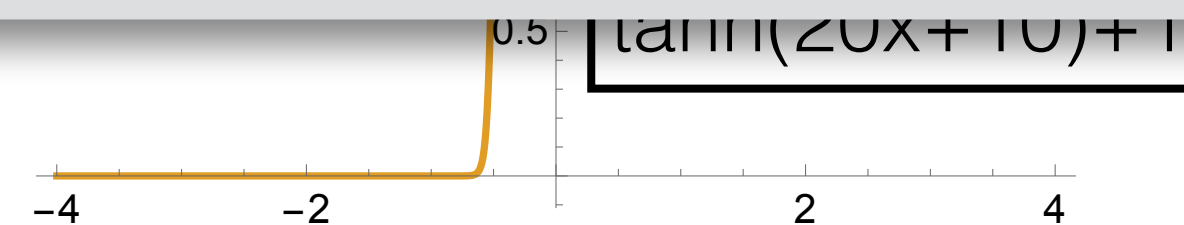
[Cybenko `89; Hornik `91; Leshno et.al. `93; Pinkus `99]



Neural Networks as function approximations

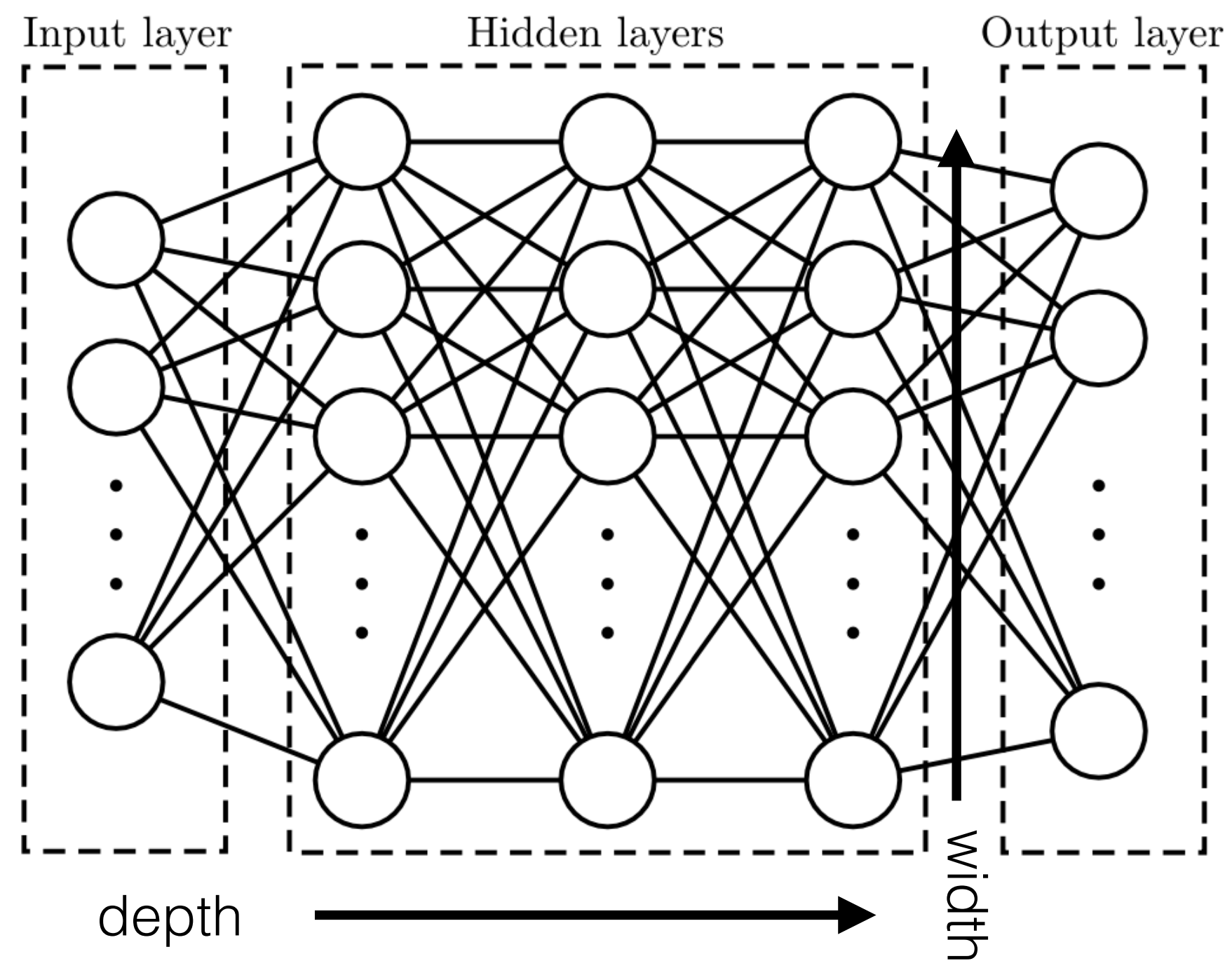
- ▶ So why would one ever do this?

DISCLAIMER:
This is **NOT**
what a NN does

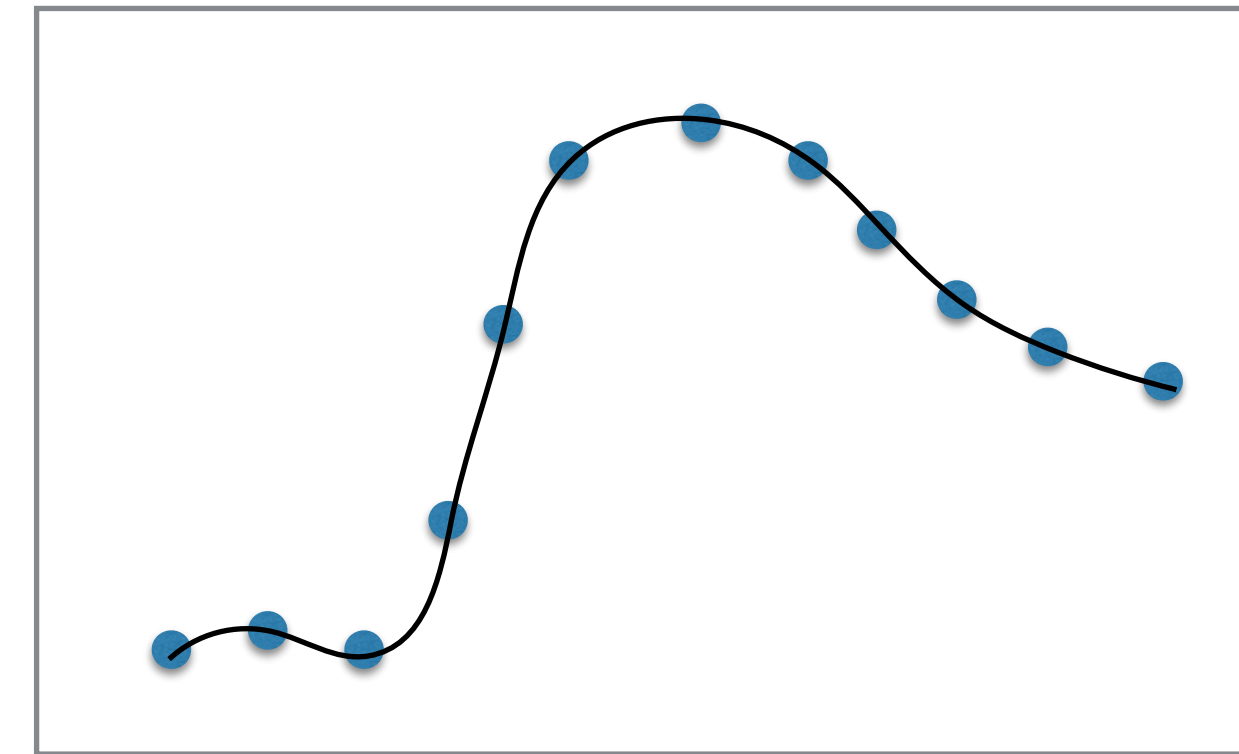


Neural Networks as function approximations

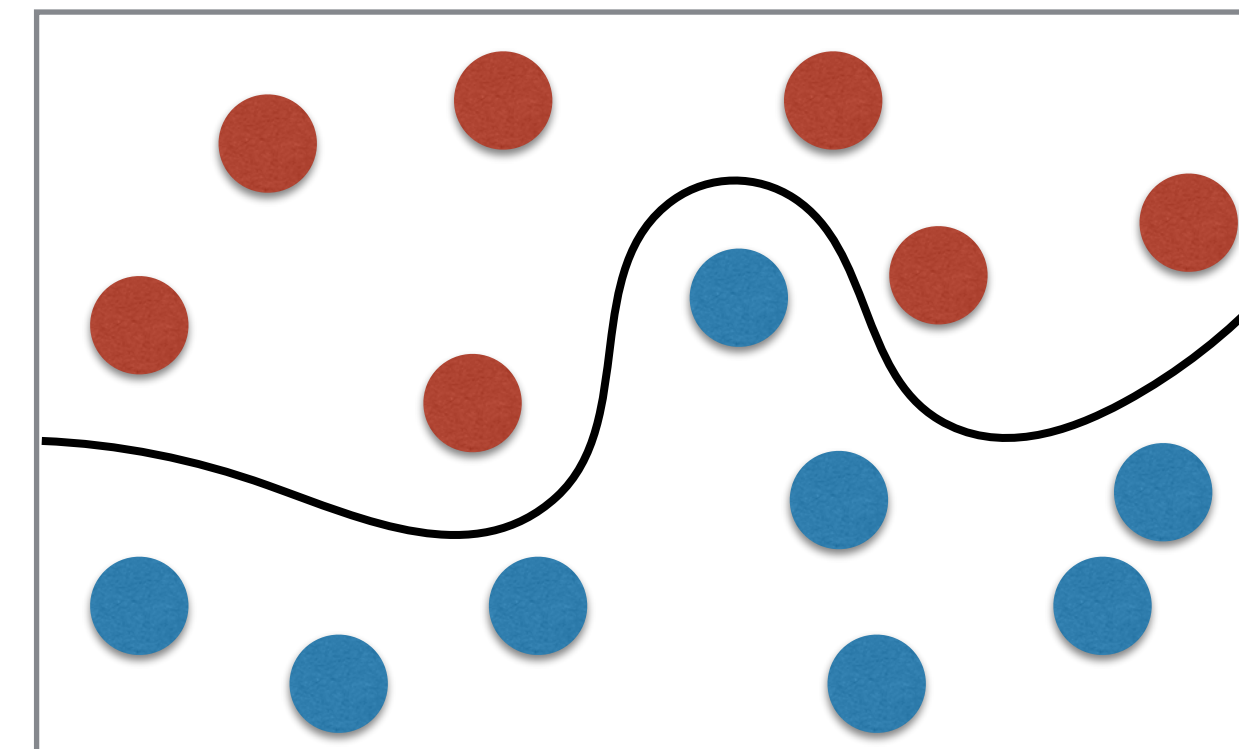
- ▶ In practice: build finite width, finite depth NN

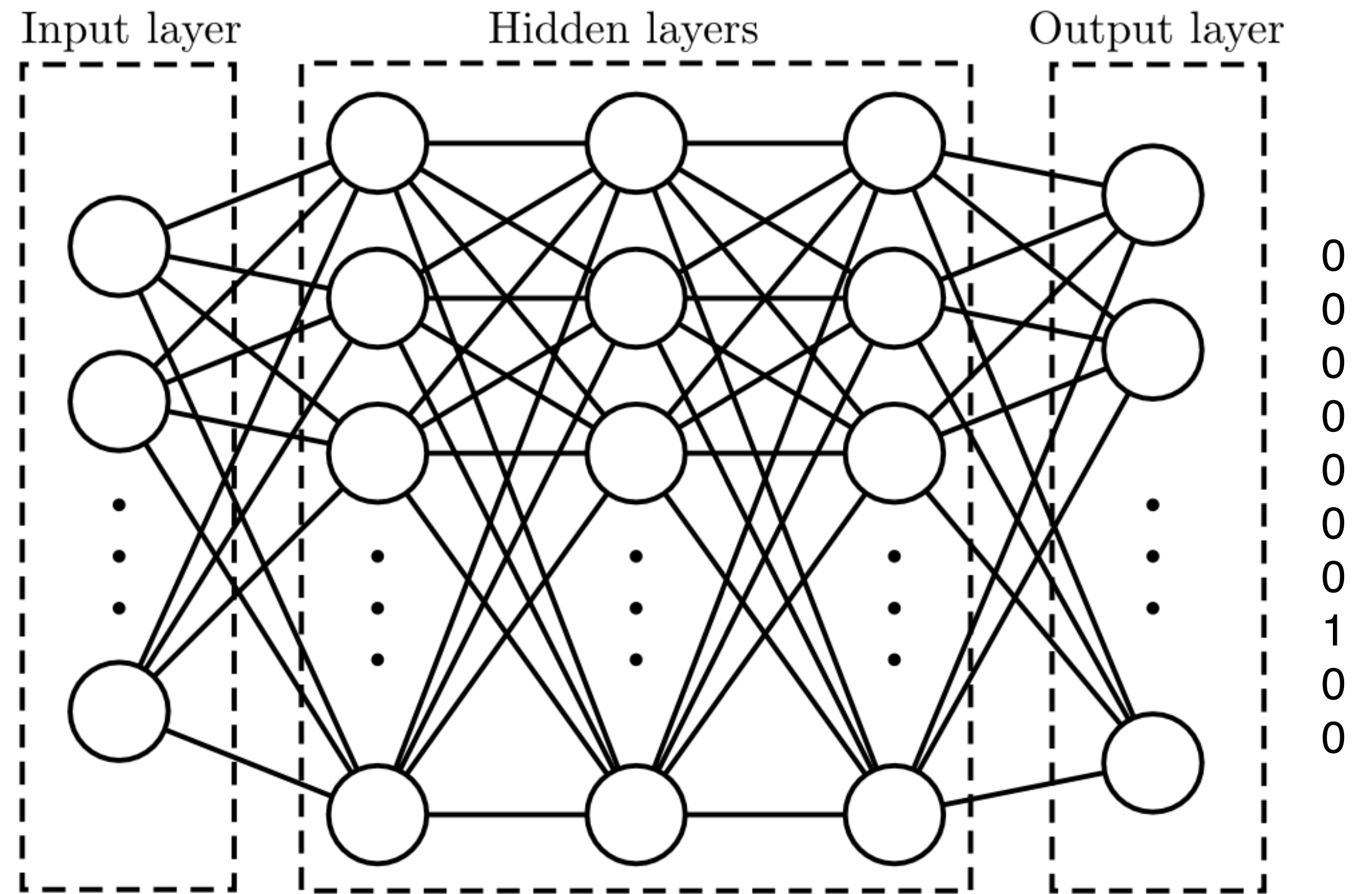


- ▶ Learning functions for regression



- ▶ Learning functions as decision boundaries for classification



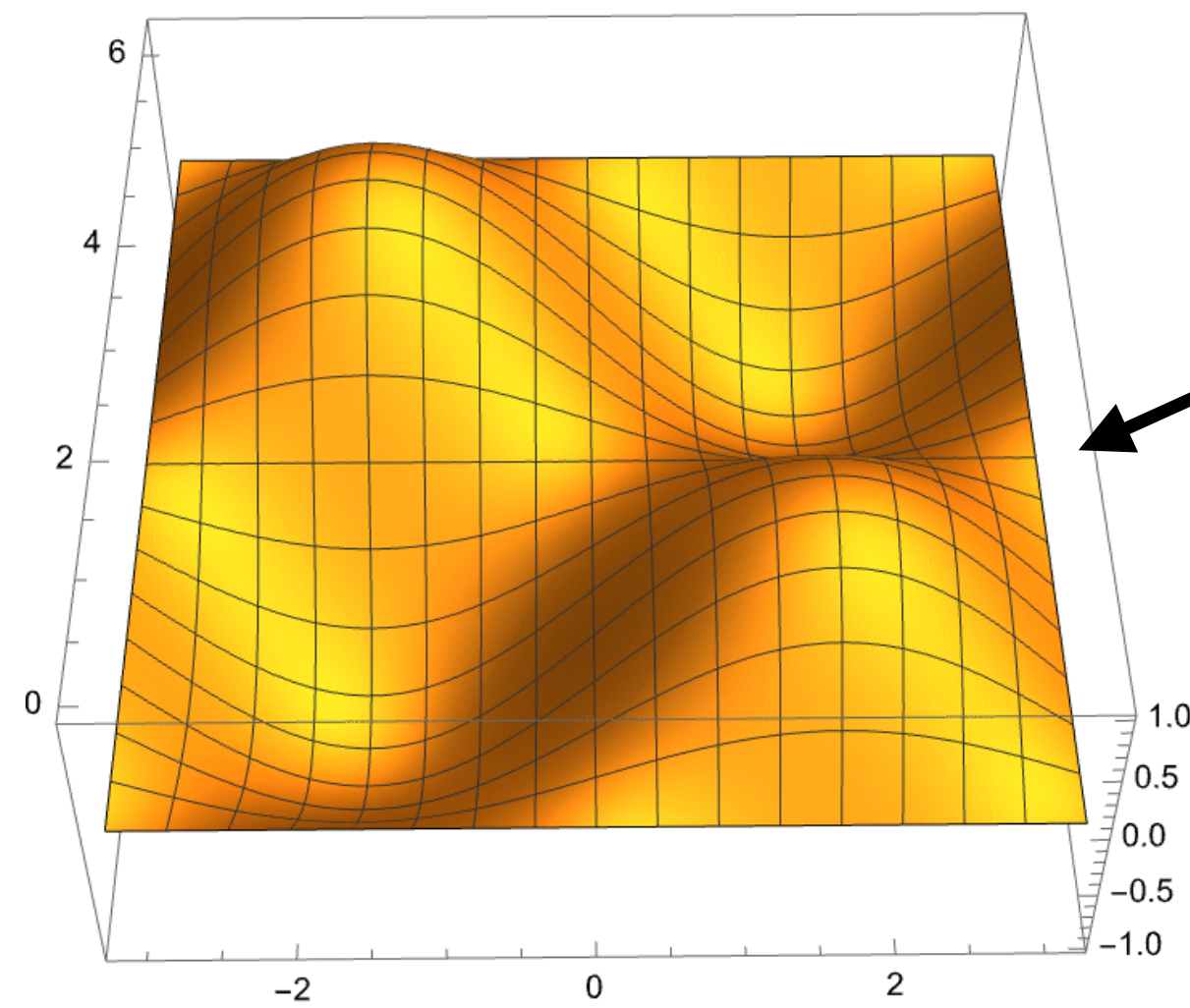


Training with backpropagation

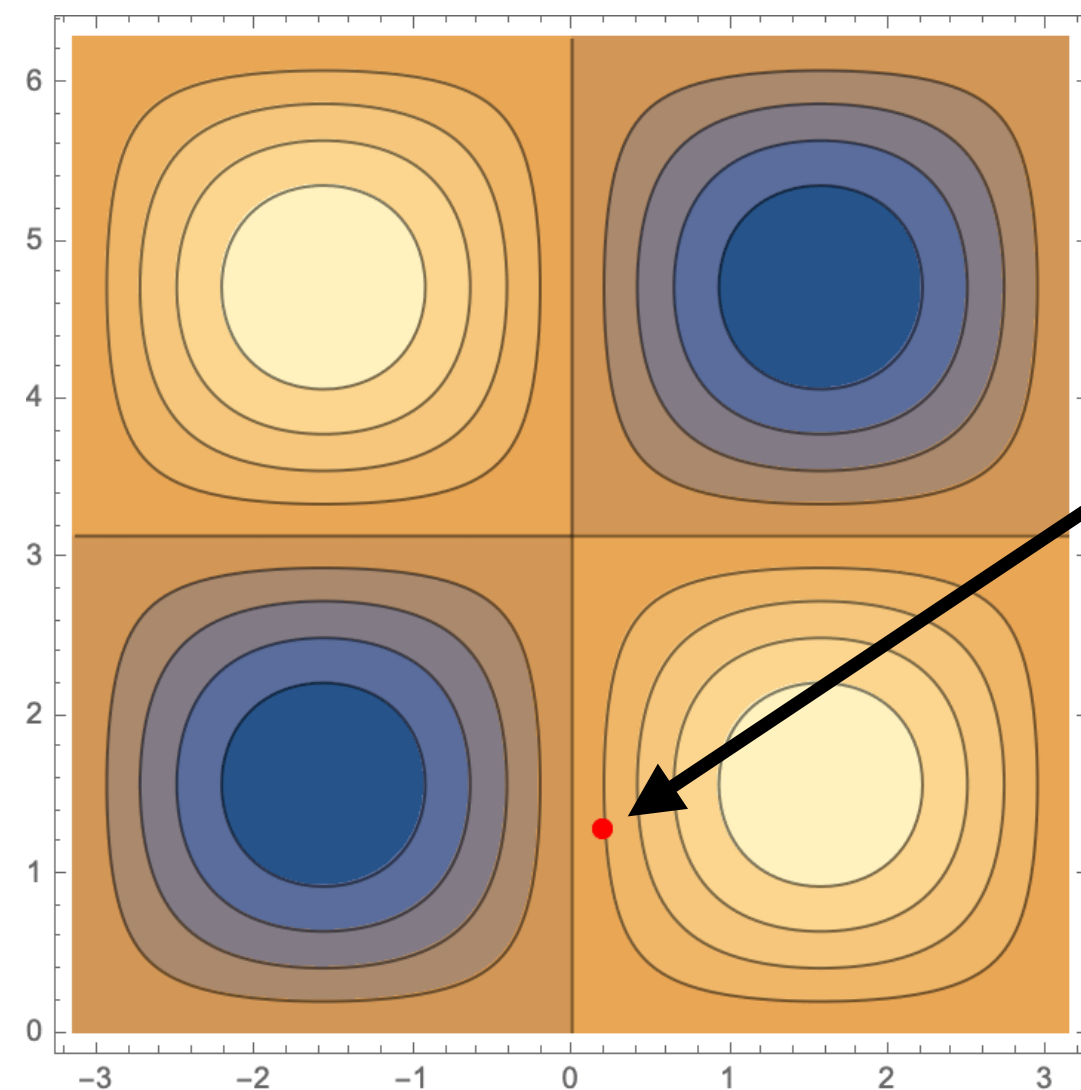
Neural Networks

- ▶ In supervised ML, you have a notion of what you want the NN output to be
- ▶ This is encoded in a **loss function** L , which gives the distance between the NN output $f_{\theta}(\vec{x})$ and the desired output $\vec{y}(\vec{x})$
- ▶ Given this distance or loss, you want to update the NN parameters θ s.t. $f_{\theta}(\vec{x}) \simeq \vec{y}(\vec{x})$
- ▶ Typically, the loss function is chosen such that $L(f_{\theta}(\vec{y}), \vec{y}(\vec{x})) \geq 0$ with equality iff $f_{\theta}(\vec{y}) = \vec{y}(\vec{x})$
- ▶ This means we need to find minima in (a million- to billion-dimensional) parameter space that generalizes well beyond training data
- ▶ Use “cheapest” (linear) optimizer: gradient descent

Visualization - Gradient descent

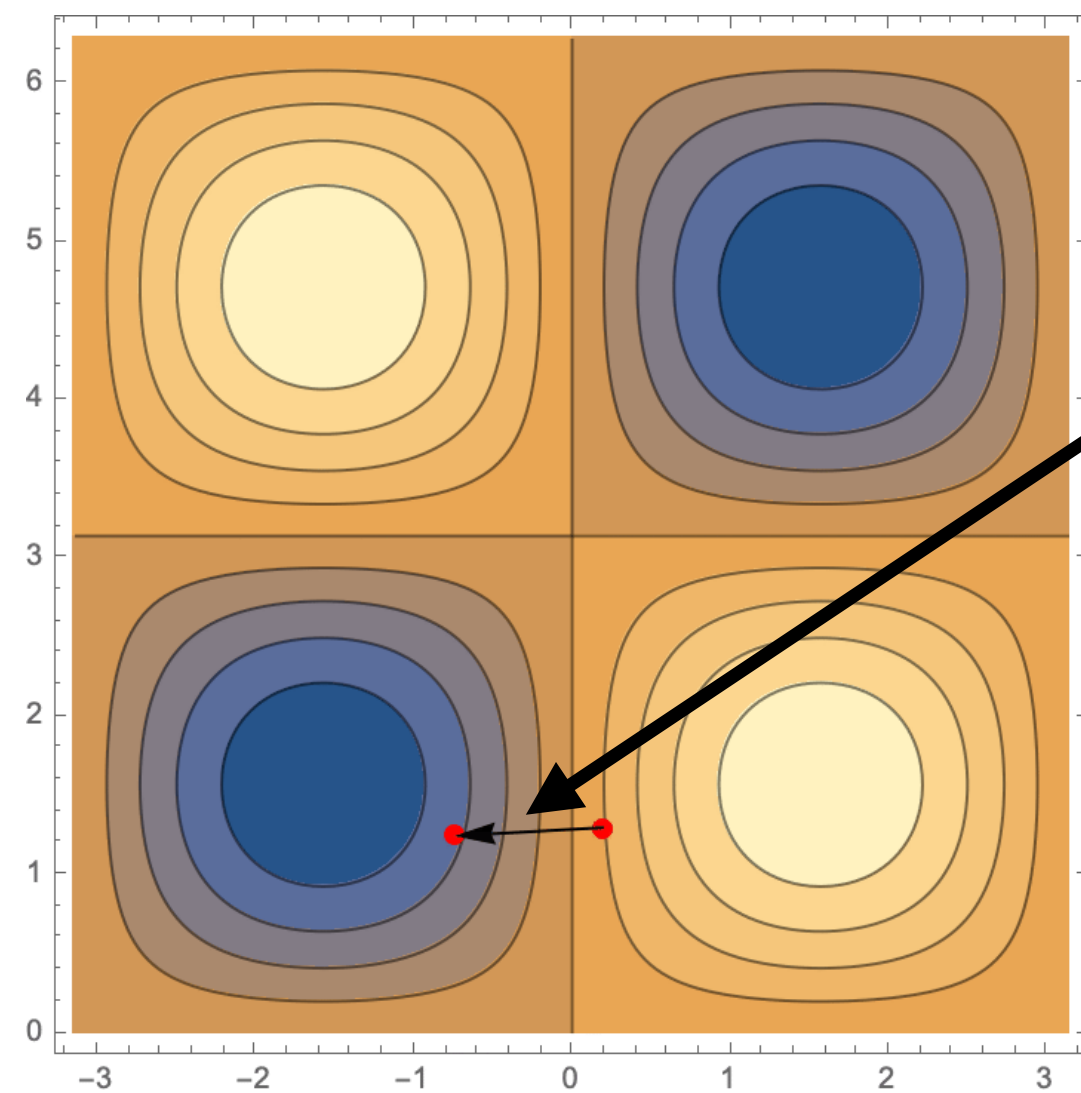
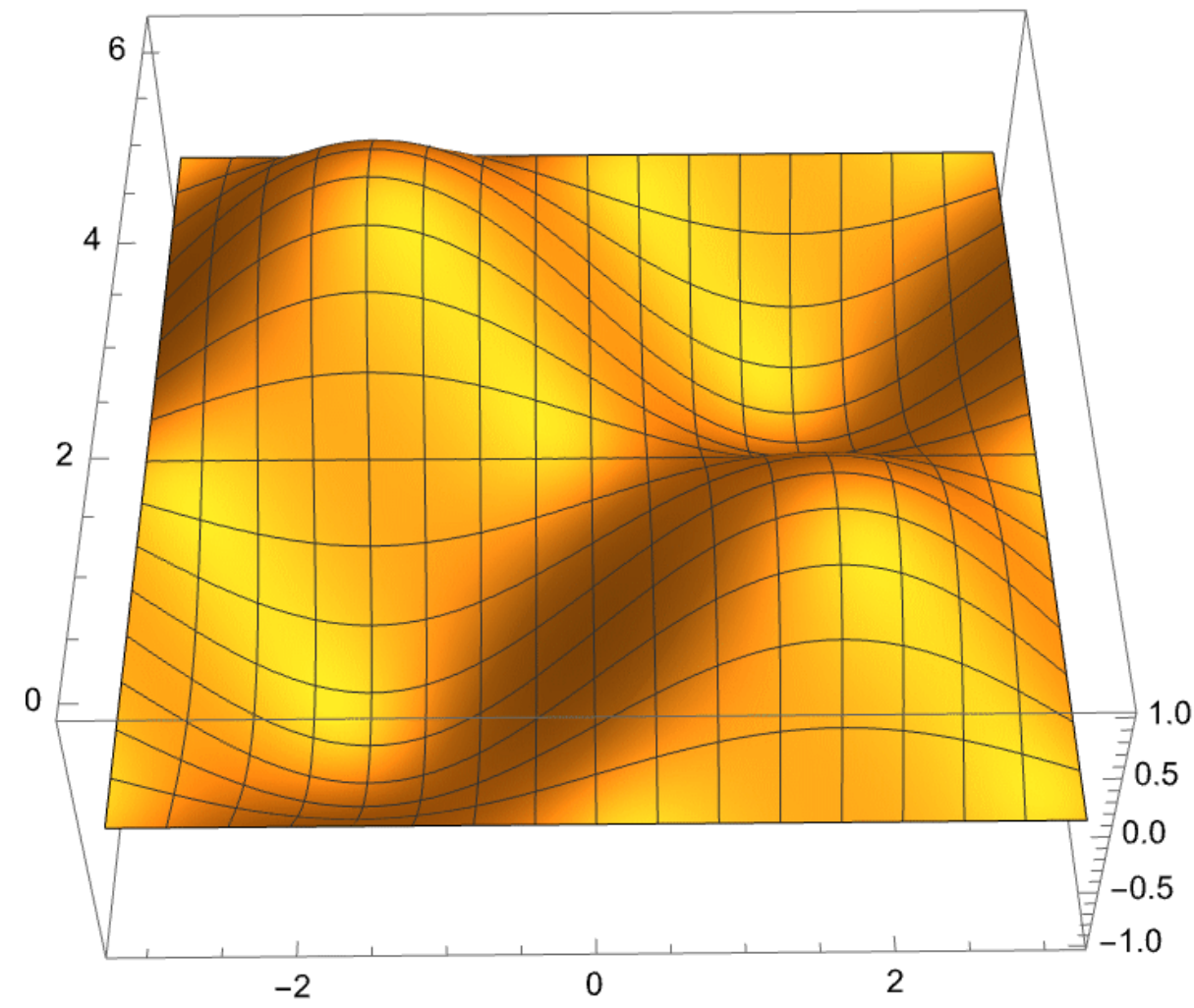


Loss landscape as a function of the NN parameters, and thus as a functional (function of a function) of the NN



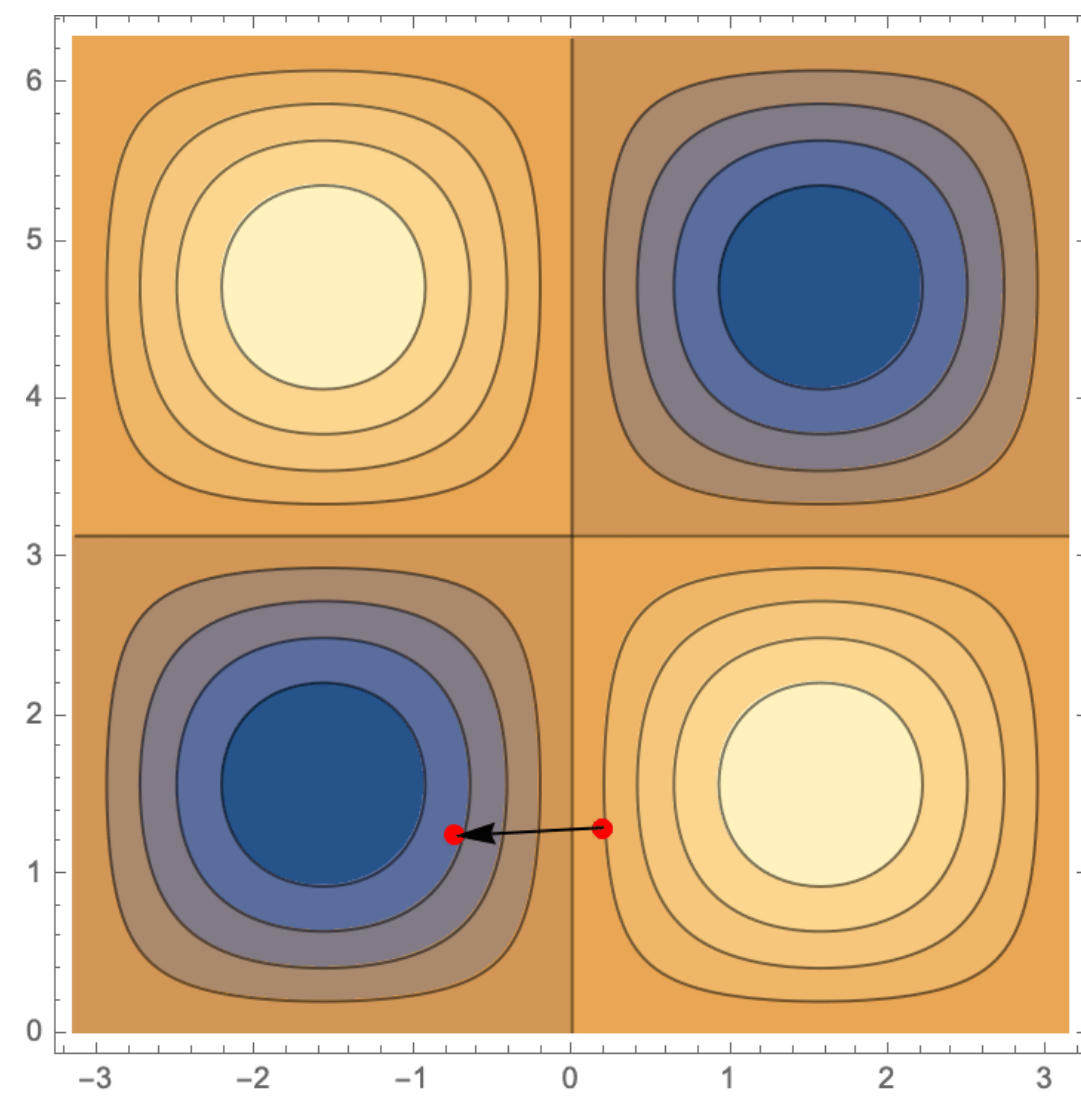
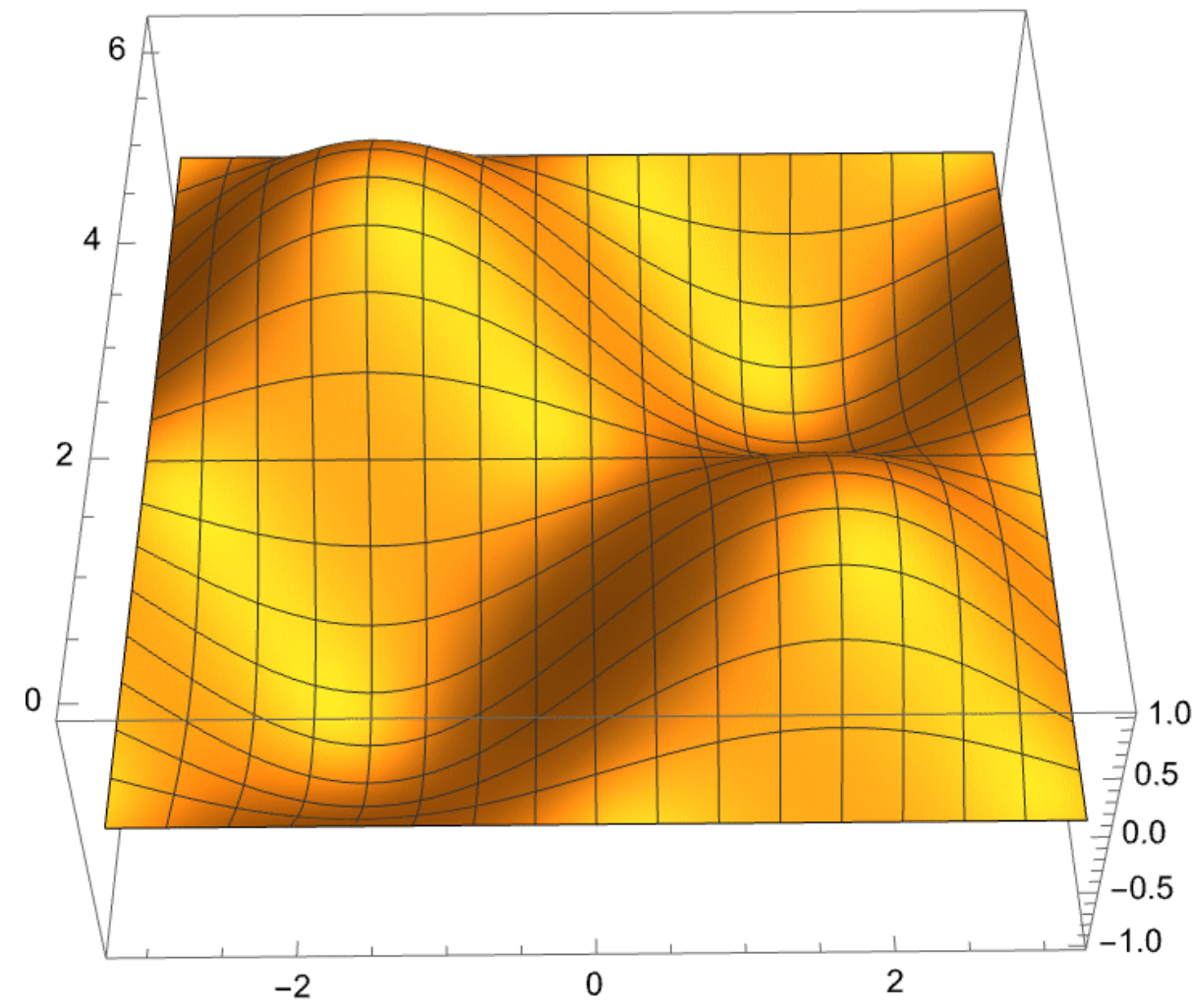
Randomly initialized NN
= Random function
= Random point in loss landscape

Visualization - Gradient descent

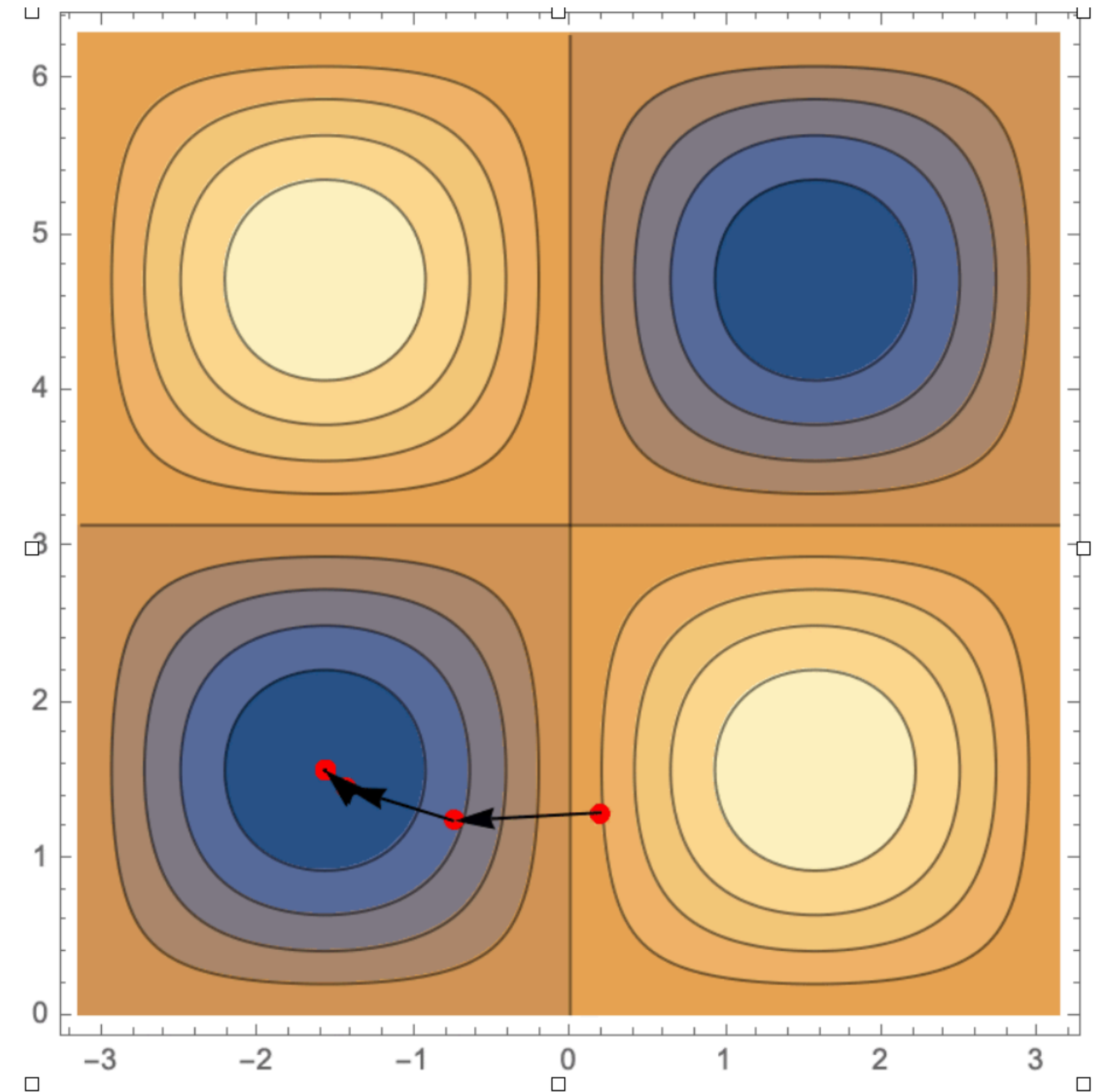


step size = learning rate

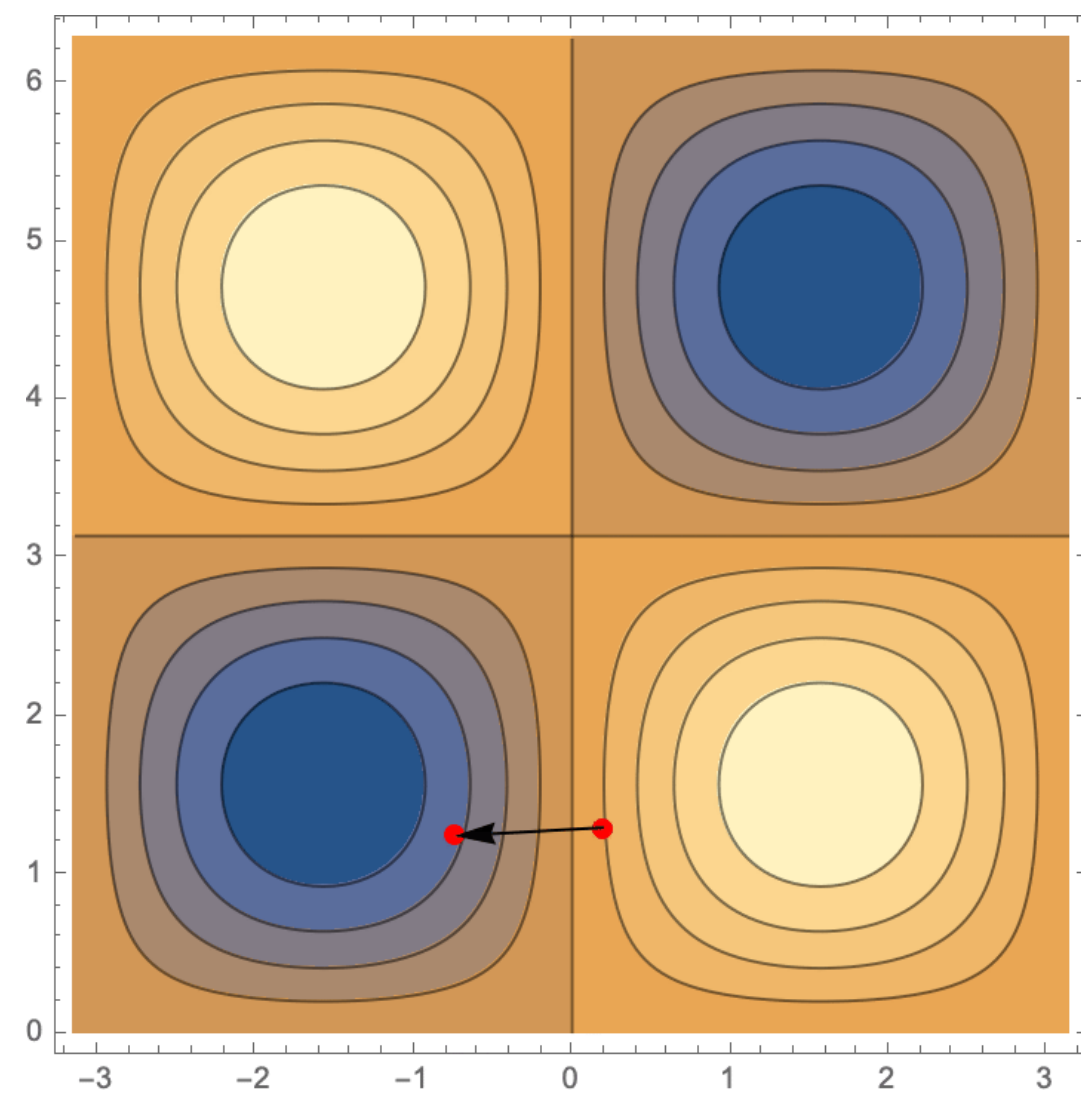
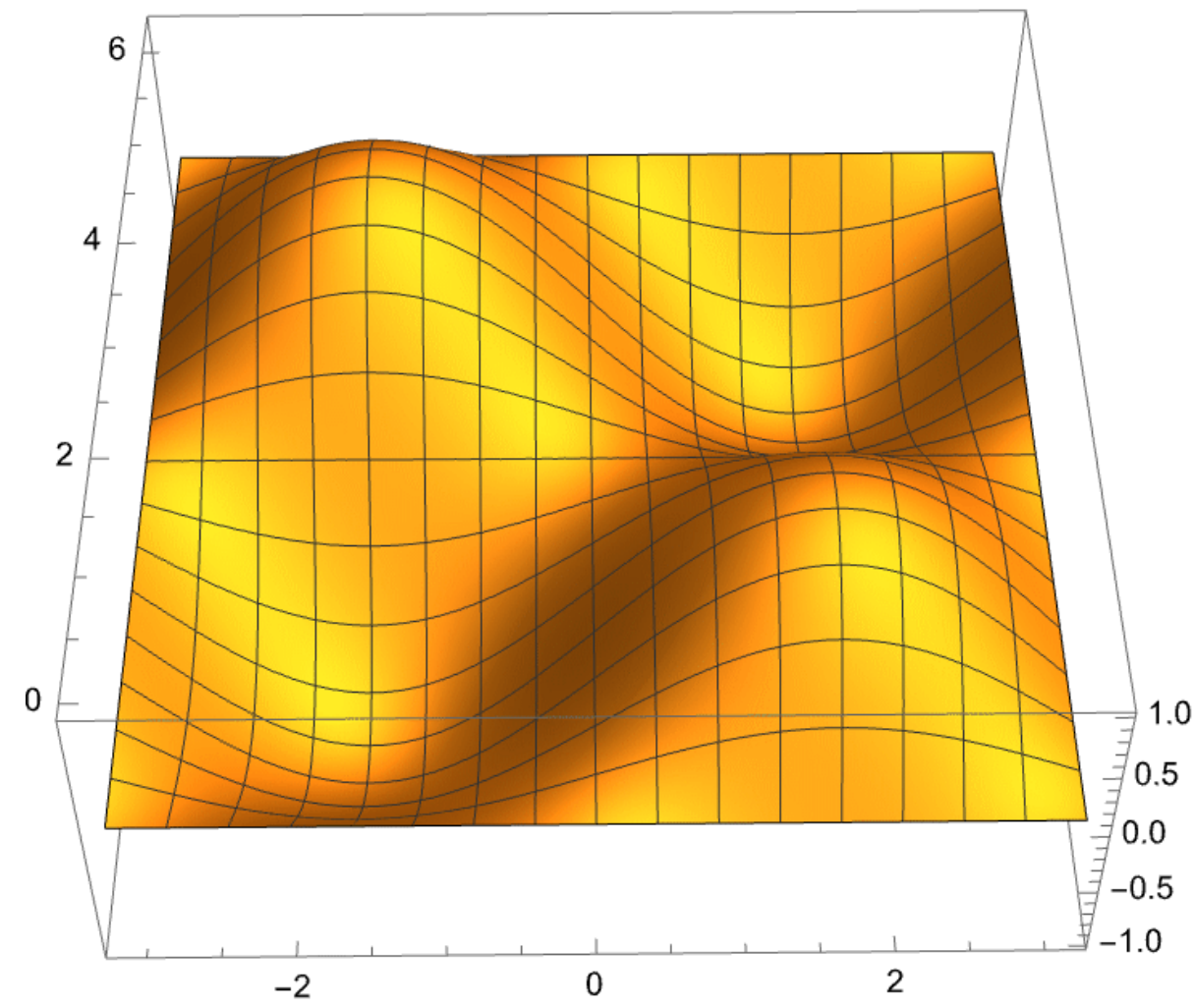
Visualization - Gradient descent



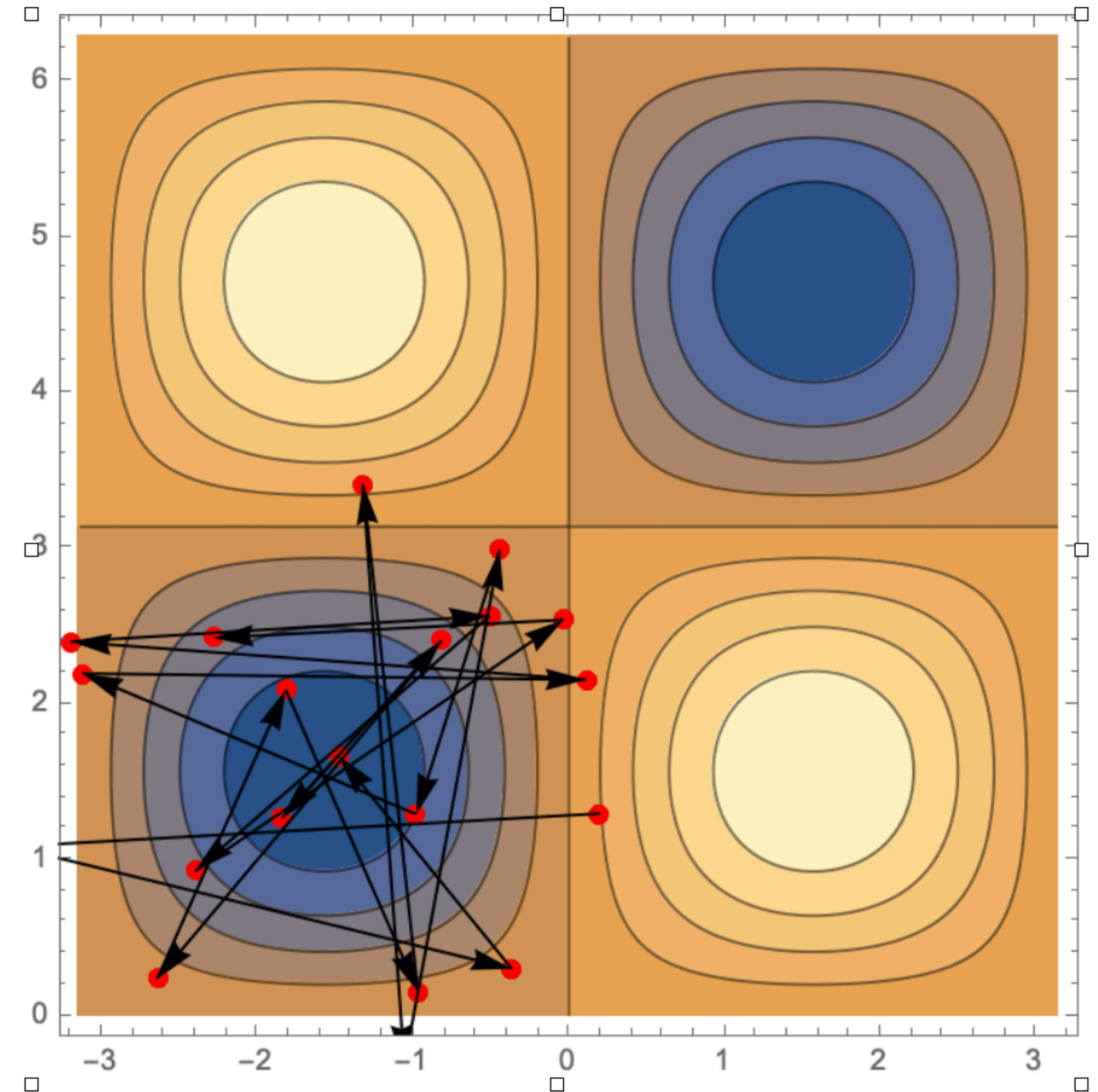
good learning rate



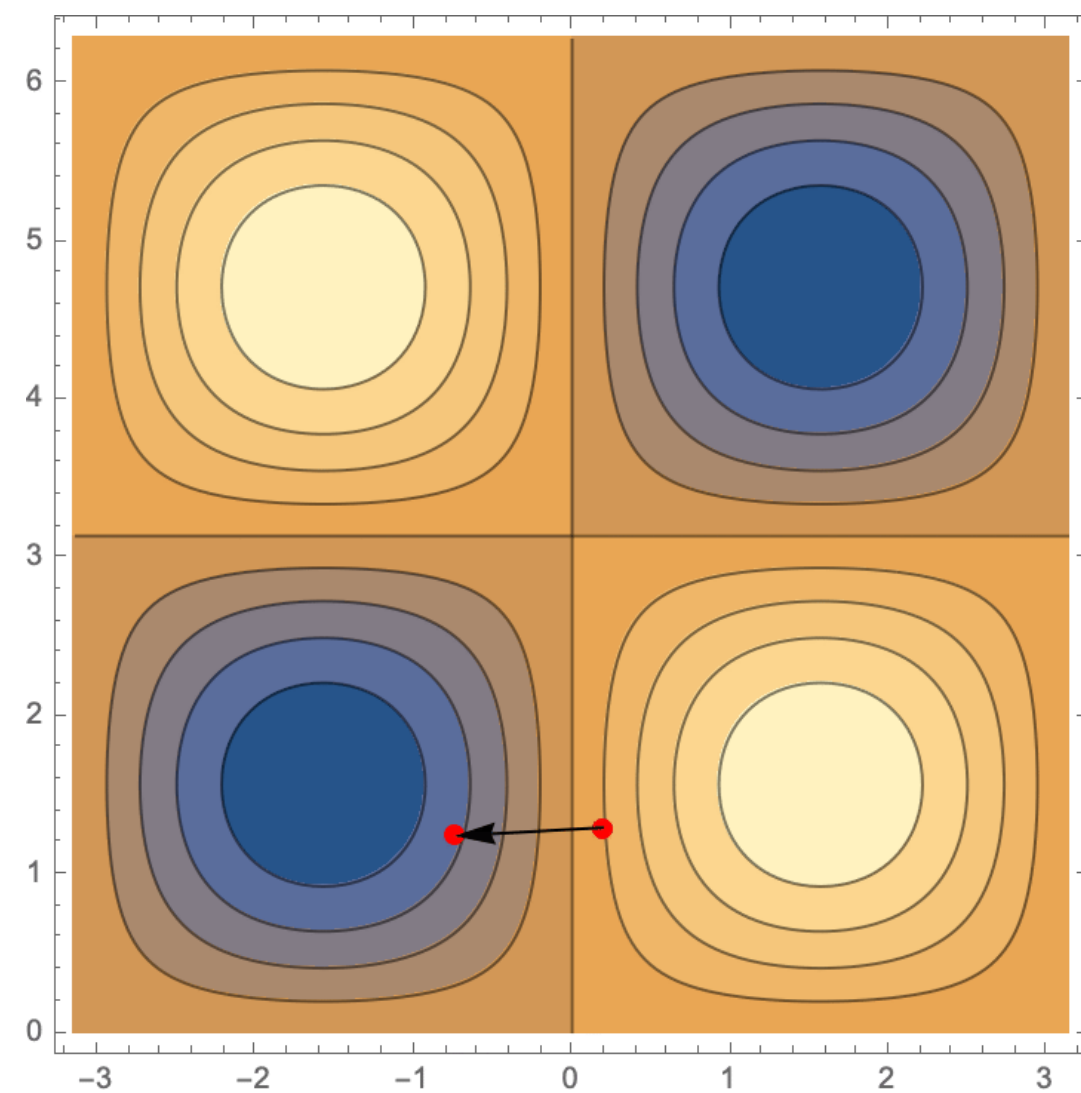
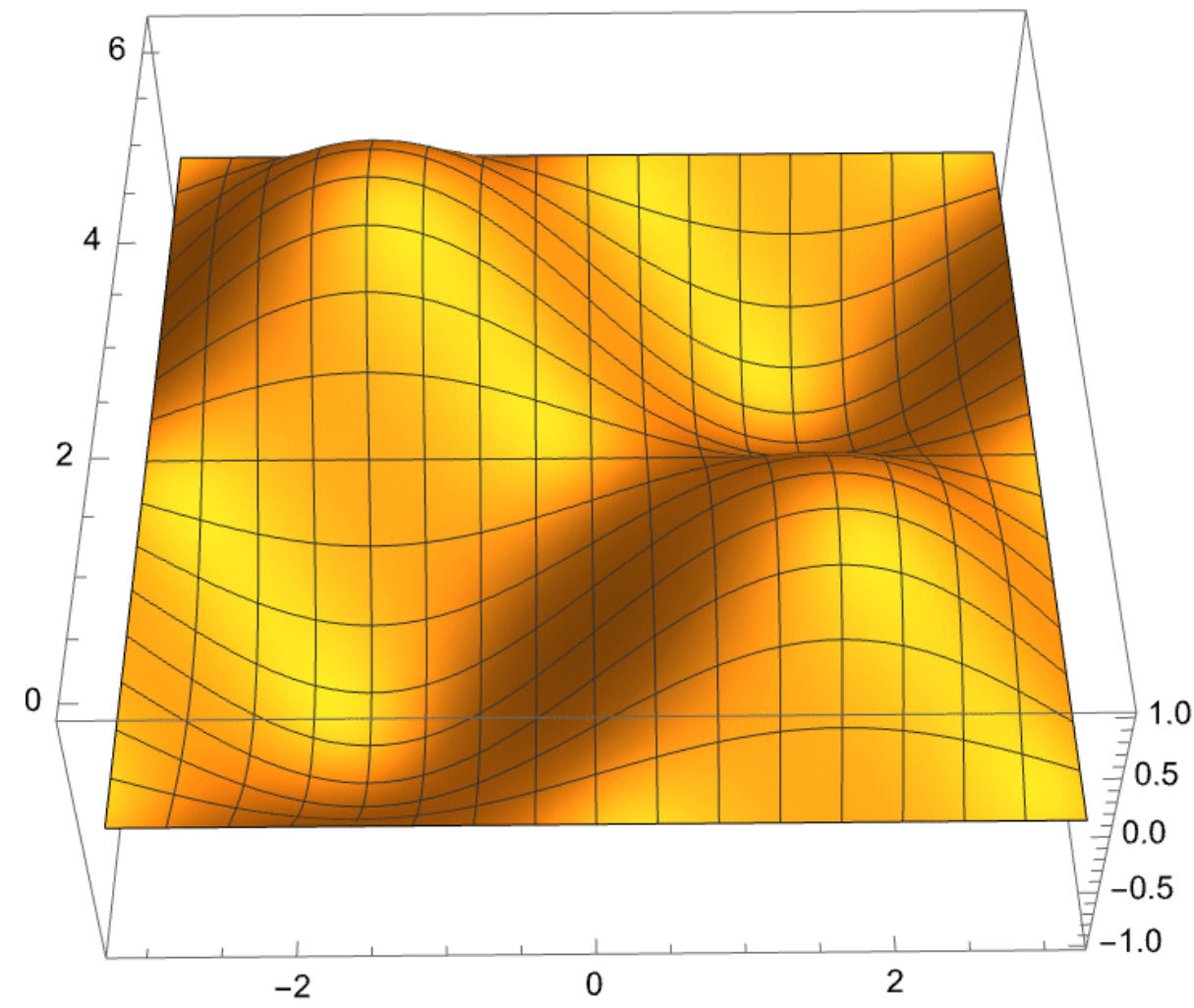
Visualization - Gradient descent



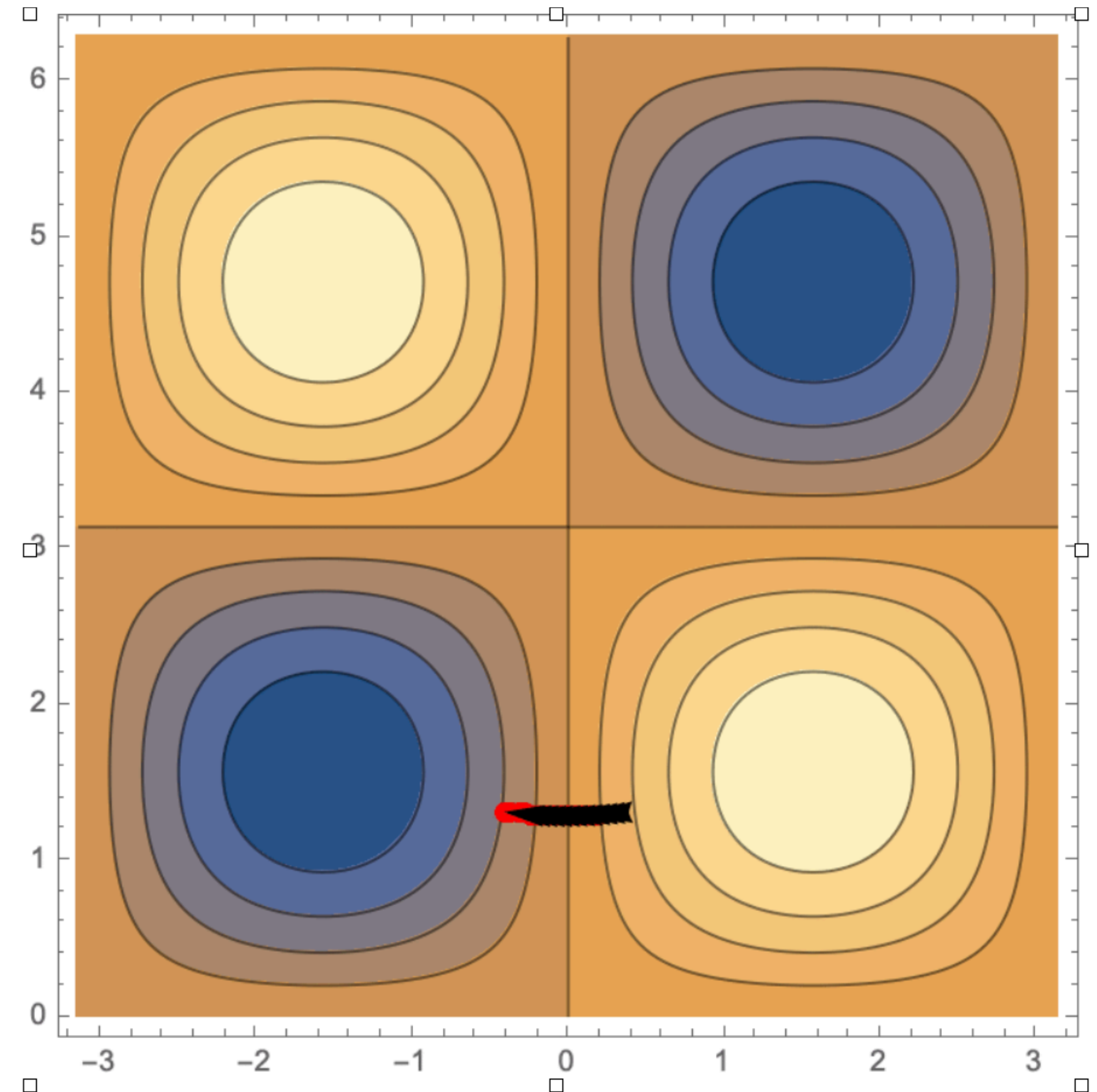
learning rate too big



Visualization - Gradient descent



learning rate too small



Gradient Descent - Loss functions

▶ For regression:

◆ MSE (Mean squared error):
$$L_{\text{MSE}} = \frac{1}{N} \sum_{i=1}^N [y(x_i) - f_{\theta}(x_i)]^2$$

◆ MAE (Mean absolute error):
$$L_{\text{MAE}} = \frac{1}{N} \sum_{i=1}^N |y(x_i) - f_{\theta}(x_i)|$$

◆ MAPE (Mean absolute percentage error):
$$L_{\text{MAPE}} = \frac{1}{N} \sum_{i=1}^N \left| \frac{y(x_i) - f_{\theta}(x_i)}{y(x_i)} \right|$$

▶ For classification:

◆ Binary cross-entropy:
$$L_{\text{BCE}} = \frac{1}{N} \sum_{i=1}^N y(x_i) \log(f_{\theta}(x_i)) - (1 - y(x_i)) \log(1 - f_{\theta}(x_i))$$

Gradient Descent - Details

- ▶ NN = composition of maps \Rightarrow use chain rule for derivatives
- ▶ At each layer of the NN, we want the gradient for each point in the dataset
- ▶ We introduce the following notation:

post-activation of layer $i + 1$ layer index Pre-activation of layer $i + 1$

$\ell_{\mu}^{(i+1)} = a^{(i+1)}(z_{\mu}^{(i+1)})$ with $z_{\mu}^{(i+1)} = \sum_{v=1}^{n_i} w_{\mu v}^{(i+1)} \ell_v^{(i)} + b_{\mu}^{(i+1)}$

μ^{th} neuron activation function

Jupyter Notebook Presentation

▼ NN approximating a function

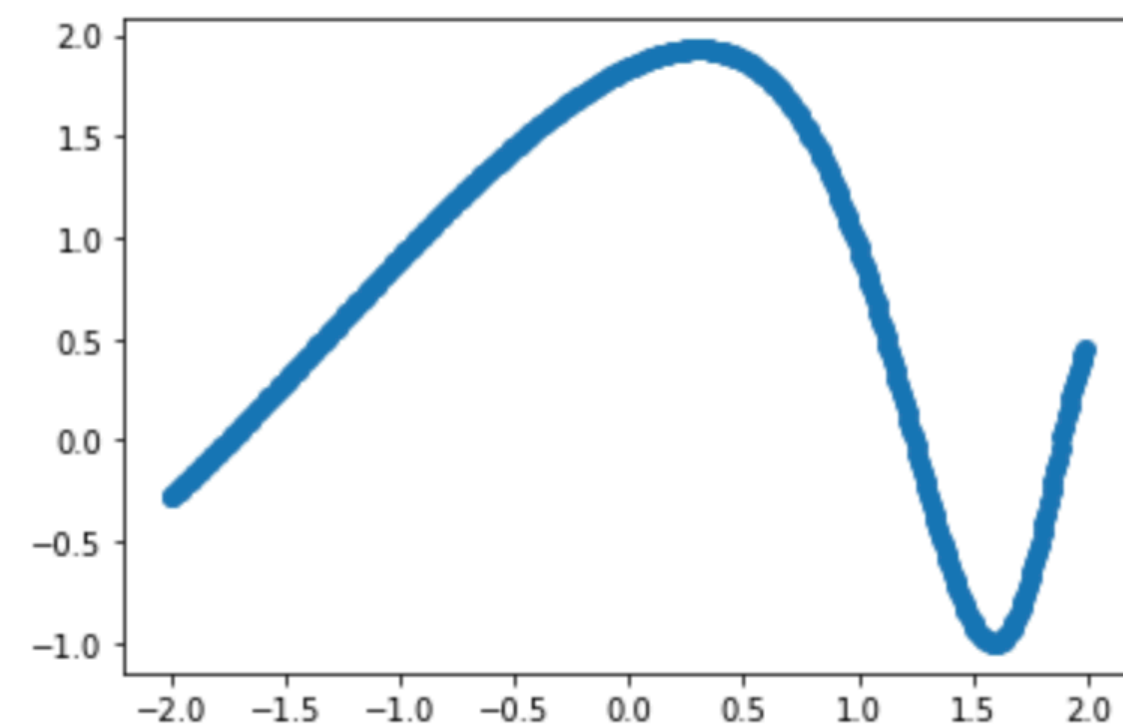
Create dataset:

$$x = [-2, 2], \quad y = \sin(e^x) + \cos(x)$$

```
In [2]: f.cast(np.arange(-2, 2, .01), dtype=tf.float32)
        f.sin(tf.exp(x)) + tf.cos(x)
```

Plot the dataset

```
In [3]: figure(1)
        scatter(x, y);
```



Setup a NN to learn this function

```
In [4]: tf.keras.Sequential()
        # add layer
        model.add(tf.keras.Input(shape=(1)))
```

Gradient Descent - Details

- ▶ NN = composition of maps \Rightarrow use chain rule for derivatives
- ▶ At each layer of the NN, we want the gradient for each point in the dataset
- ▶ We introduce the following notation: $\ell_\mu^{(i+1)} = a^{(i+1)}(z_\mu^{(i+1)})$ with $z_\mu^{(i+1)} = \sum_{\nu=1}^{n_i} w_{\mu\nu}^{(i+1)} \ell_\nu^{(i)} + b_\mu^{(i+1)}$

▶ Then:

$$\diamond \quad \frac{\partial L}{\partial \theta^{(n)}} = \frac{\partial L}{\partial z^{(n)}} \frac{\partial z^{(n)}}{\partial \theta^{(n)}} \longleftarrow \frac{\partial z_\mu^{(n)}}{\partial \theta_a^{(n)}} = \begin{cases} \frac{\partial z_\mu^{(n)}}{\partial w_{\mu\nu}^{(n)}} = \ell_\nu^{(n-1)} & \text{for } \theta_a^{(n)} = w_{\mu\nu}^{(n)} \\ \frac{\partial z_\mu^{(n)}}{\partial b_\nu^{(n)}} = \delta_{\mu\nu} & \text{for } \theta_a^{(n)} = b_\nu^{(n)} \end{cases} \quad \begin{array}{l} \text{known from} \\ \text{forward pass} \end{array}$$

$$\frac{\partial L_{\text{MSE}}}{\partial z^{(n)}} = -\frac{2}{N} \left[y - a^{(n)}(z^{(n)}) \right] a'^{(n)}(z^{(n)})$$

Gradient Descent - Details

- ▶ NN = composition of maps \Rightarrow use chain rule for derivatives
- ▶ At each layer of the NN, we want the gradient for each point in the dataset
- ▶ We introduce the following notation: $\ell_{\mu}^{(i+1)} = a^{(i+1)}(z_{\mu}^{(i+1)})$ with $z_{\mu}^{(i+1)} = \sum_{\nu=1}^{n_i} w_{\mu\nu}^{(i+1)} \ell_{\nu}^{(i)} + b_{\mu}^{(i+1)}$
- ▶ Then:

$$\begin{aligned}
 \blacklozenge \quad & \frac{\partial L}{\partial \theta^{(n)}} = \frac{\partial L}{\partial z^{(n)}} \frac{\partial z^{(n)}}{\partial \theta^{(n)}} \\
 \blacklozenge \quad & \frac{\partial L}{\partial \theta^{(n-1)}} = \frac{\partial L}{\partial z^{(n)}} \frac{\partial z^{(n)}}{\partial z^{(n-1)}} \frac{\partial z^{(n-1)}}{\partial \theta^{(n-1)}} \quad \leftarrow \quad \frac{\partial z_{\mu}^{(n-1)}}{\partial \theta_a^{(n-1)}} = \begin{cases} \frac{\partial z_{\mu}^{(n-1)}}{\partial w_{\mu\nu}^{(n-1)}} = \ell_{\nu}^{(n-2)} & \text{for } \theta_a = w_{\mu\nu}^{(n)} \\ \frac{\partial z_{\mu}^{(n-1)}}{\partial b_{\nu}^{(n-1)}} = \delta_{\mu\nu} & \text{for } \theta_a = b_{\mu}^{(n)} \end{cases} \\
 & \quad \quad \quad \left(\sum_{\mu=1}^{N_{n-1}} w_{\mu\nu}^{(n-1)} \frac{\partial L}{\partial z_{\mu}^{(n)}} \right) a'^{(n-1)}(z_{\nu}^{(n-1)}) \quad \text{known from forward pass}
 \end{aligned}$$

Gradient Descent - Details

- ▶ NN = composition of maps \Rightarrow use chain rule for derivatives
- ▶ At each layer of the NN, we want the gradient for each point in the dataset
- ▶ We introduce the following notation: $\ell_{\mu}^{(i+1)} = a^{(i+1)}(z_{\mu}^{(i+1)})$ with $z_{\mu}^{(i+1)} = \sum_{\nu=1}^{n_i} w_{\mu\nu}^{(i+1)} \ell_{\nu}^{(i)} + b_{\mu}^{(i+1)}$
- ▶ Then:

$$\blacklozenge \quad \frac{\partial L}{\partial \theta^{(n)}} = \frac{\partial L}{\partial z^{(n)}} \frac{\partial z^{(n)}}{\partial \theta^{(n)}}$$

$$\blacklozenge \quad \frac{\partial L}{\partial \theta^{(n-1)}} = \frac{\partial L}{\partial z^{(n)}} \frac{\partial z^{(n)}}{\partial z^{(n-1)}} \frac{\partial z^{(n-1)}}{\partial \theta^{(n-1)}}$$

$$\blacklozenge \quad \frac{\partial L}{\partial \theta^{(n-2)}} = \frac{\partial L}{\partial z^{(n)}} \frac{\partial z^{(n)}}{\partial z^{(n-1)}} \frac{\partial z^{(n-1)}}{\partial z^{(n-2)}} \frac{\partial z^{(n-2)}}{\partial \theta^{(n-2)}}$$

Gradient Descent - Details

- ▶ NN = composition of maps \Rightarrow use chain rule for derivatives
- ▶ At each layer of the NN, we want the gradient for each point in the dataset
- ▶ We introduce the following notation: $\ell_{\mu}^{(i+1)} = a^{(i+1)}(z_{\mu}^{(i+1)})$ with $z_{\mu}^{(i+1)} = \sum_{\nu=1}^{n_i} w_{\mu\nu}^{(i+1)} \ell_{\nu}^{(i)} + b_{\mu}^{(i+1)}$
- ▶ Then:

$$\blacklozenge \quad \frac{\partial L}{\partial \theta^{(n)}} = \frac{\partial L}{\partial z^{(n)}} \frac{\partial z^{(n)}}{\partial \theta^{(n)}}$$

$$\blacklozenge \quad \frac{\partial L}{\partial \theta^{(n-1)}} = \frac{\partial L}{\partial z^{(n)}} \frac{\partial z^{(n)}}{\partial z^{(n-1)}} \frac{\partial z^{(n-1)}}{\partial \theta^{(n-1)}}$$

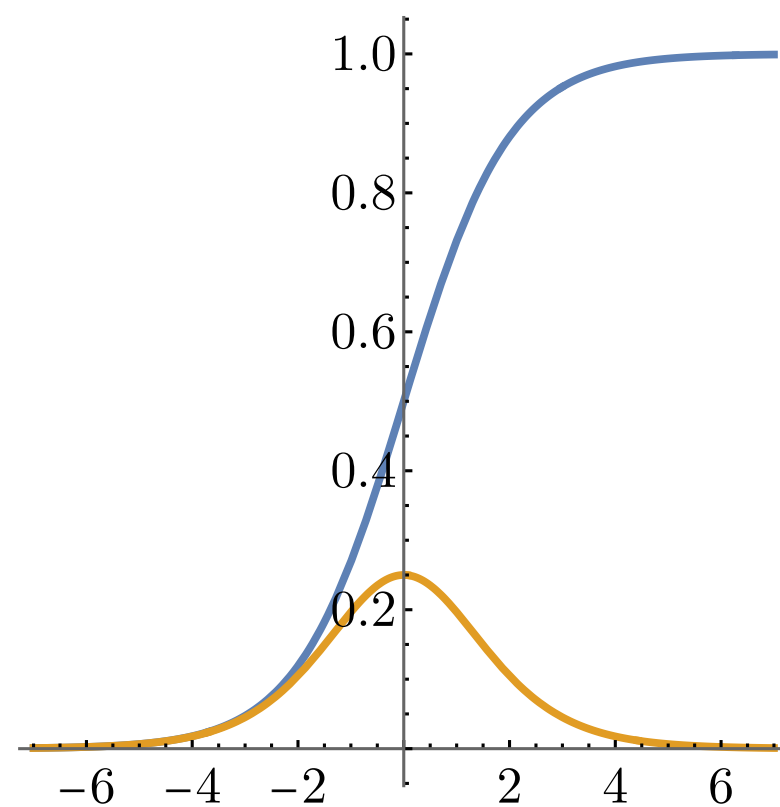
$$\blacklozenge \quad \frac{\partial L}{\partial \theta^{(n-2)}} = \frac{\partial L}{\partial z^{(n)}} \frac{\partial z^{(n)}}{\partial z^{(n-1)}} \frac{\partial z^{(n-1)}}{\partial z^{(n-2)}} \frac{\partial z^{(n-2)}}{\partial \theta^{(n-2)}}$$

GD parameter update

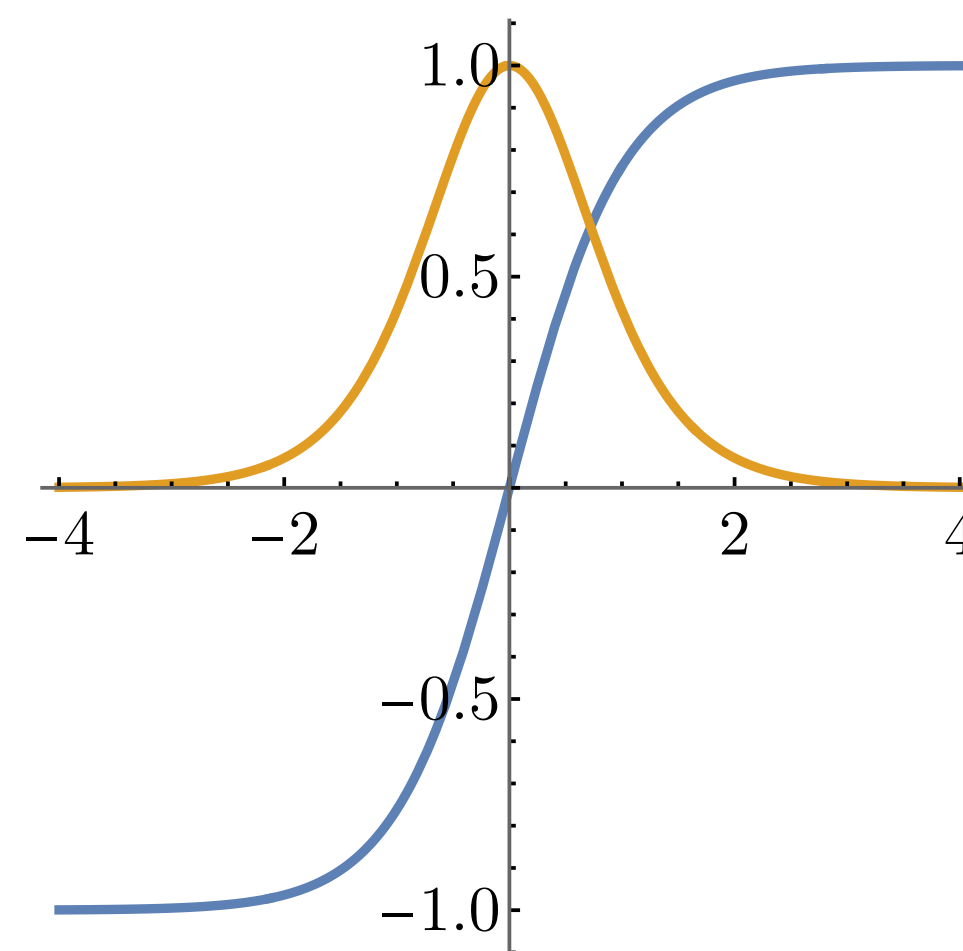
$$\theta^{(i)} \rightarrow \theta^{(i)} - \alpha \frac{\partial L}{\partial \theta^{(i)}}$$

Gradient Descent - Consequences

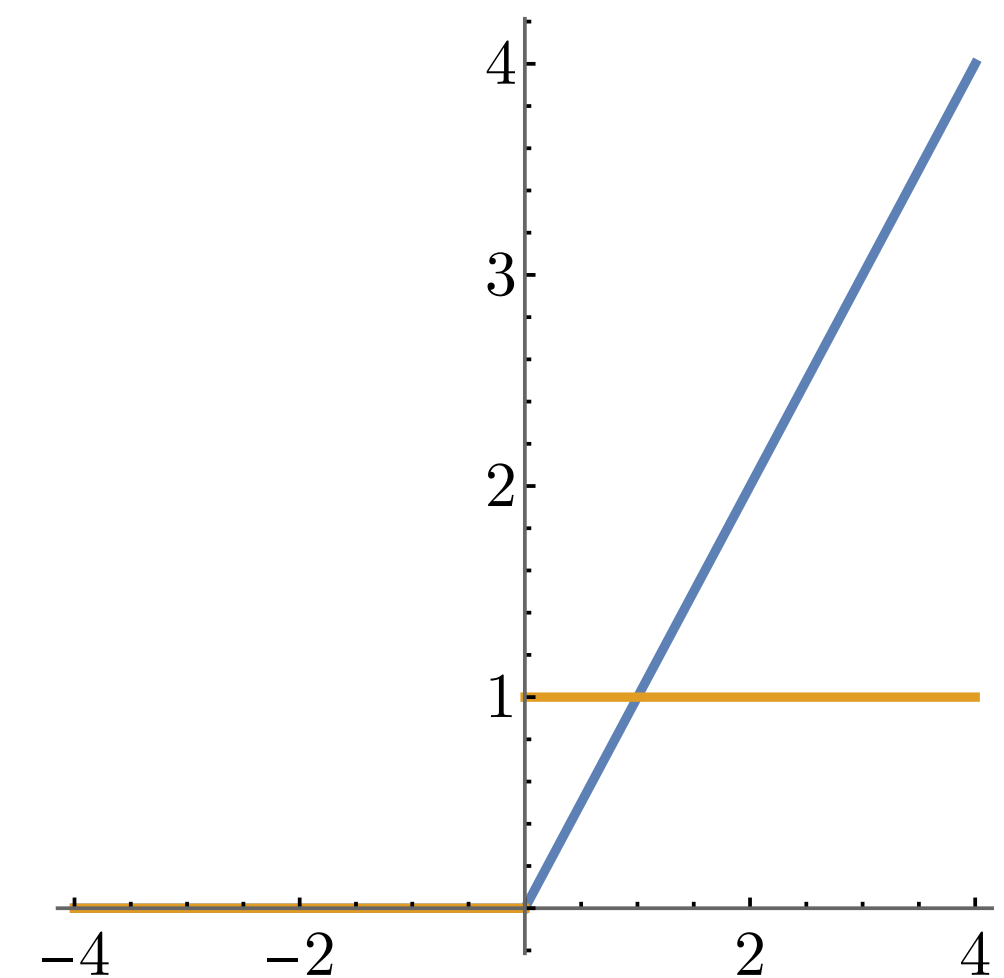
- ▶ Most of what we need in backward pass has been computed in forward pass or can be trivially computed
- ▶ ...but we need $a^{(i)}(z^{(i-1)})$
- ▶ Common activation functions have $a^{(i)} = f(a^{(i)})$, so reuse from fwd pass



$$a(x) = \sigma(x) = 1/(1 + e^{-x})$$
$$a'(x) = a(x)[1 - a(x)]$$



$$a(x) = \tanh(x)$$
$$a'(x) = 1 - [a(x)]^2$$



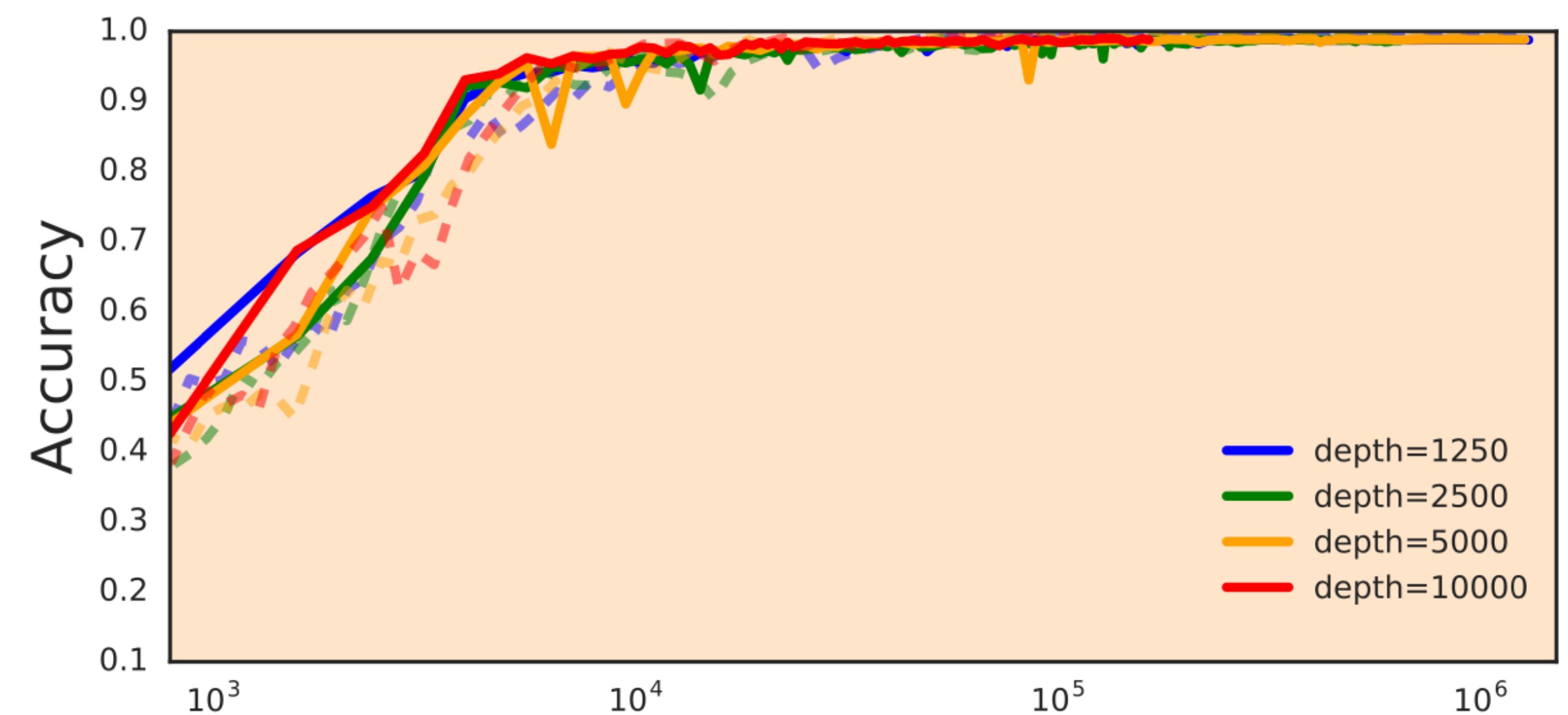
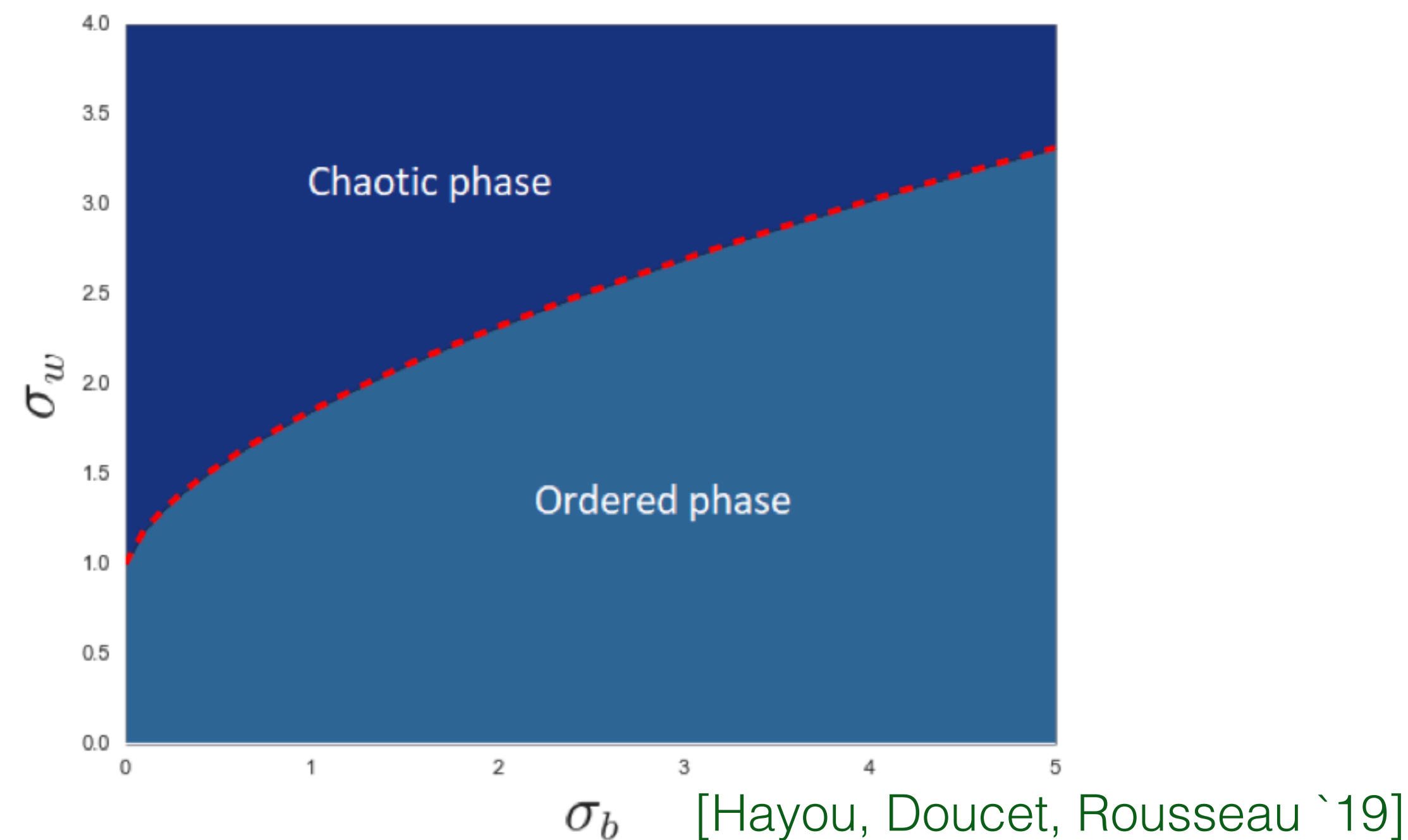
$$a(x) = \text{ReLU}(x) = \max(0, x)$$
$$a'(x) = \theta(x)$$

Gradient Descent - Consequences

- ▶ NN are typically used in big data applications
- ▶ Can parallelize over input-output-pairs
 - $O(10)$ cores on CPUs
 - $O(10k)$ cores on GPUs
- ▶ But not all input-output-pairs can fit into memory simultaneously
- ▶ Batch data into mini-batches and perform updates after each mini-batch (typically batch size is $O(10)$ - $O(100)$)
- ▶ This leads to stochastic GD (since we are not optimizing on entire loss landscape but only on a subset computed from each batch)
- ▶ This can in fact help overcome local minima and saddles, but hinders parallelization

Gradient Descent - Consequences

- ▶ Gradient at layer i proportional to gradient at layer $i + 1$
- ▶ Get vanishing / exploding gradients at earlier layers
- ▶ Counter: Edge of chaos initialization, gradient/weight clipping, batch norm, skip connections (UNET, ResNET), ...



Gradient Descent - Consequences

- ▶ Gradient at layer i proportional to gradient at layer $i + 1$
- ▶ Dying ReLU problem: $\text{ReLU}(x) = \text{ReLU}'(x) = 0$ for $x \leq 0$
- ▶ ReLU can die and not activate other layers; updates are zero as well, so this behavior won't change
- ▶ Combat with leaky ReLU: $\text{leaky ReLU} = \begin{cases} x & \text{if } x > 0 \\ .01 x & \text{if } x \leq 0 \end{cases}$
- ▶ To a lesser extent the same can happen to sigmoid:

$$\sigma(x) \approx \sigma'(x) \approx 0 \quad \text{for } x \ll 1$$

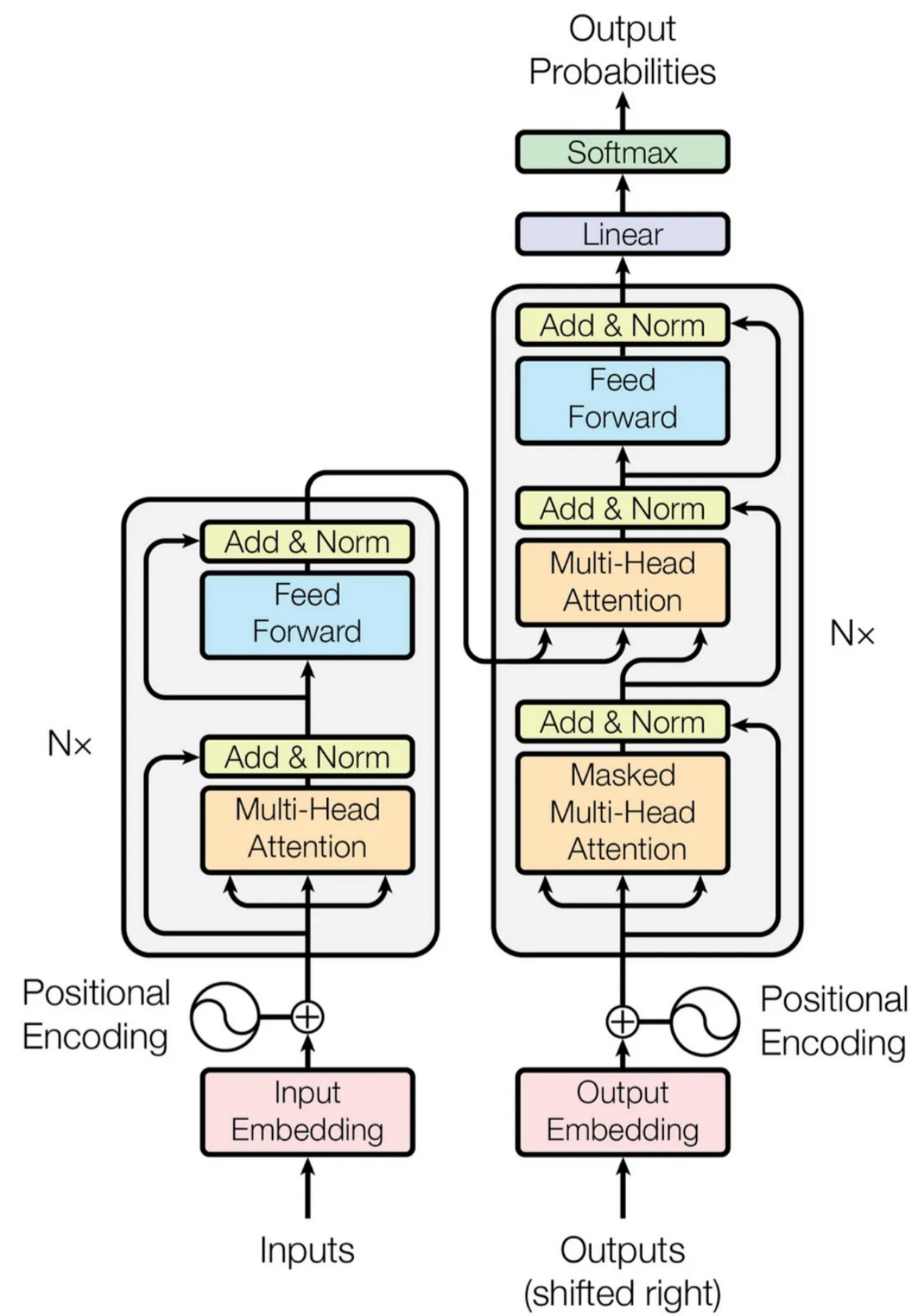
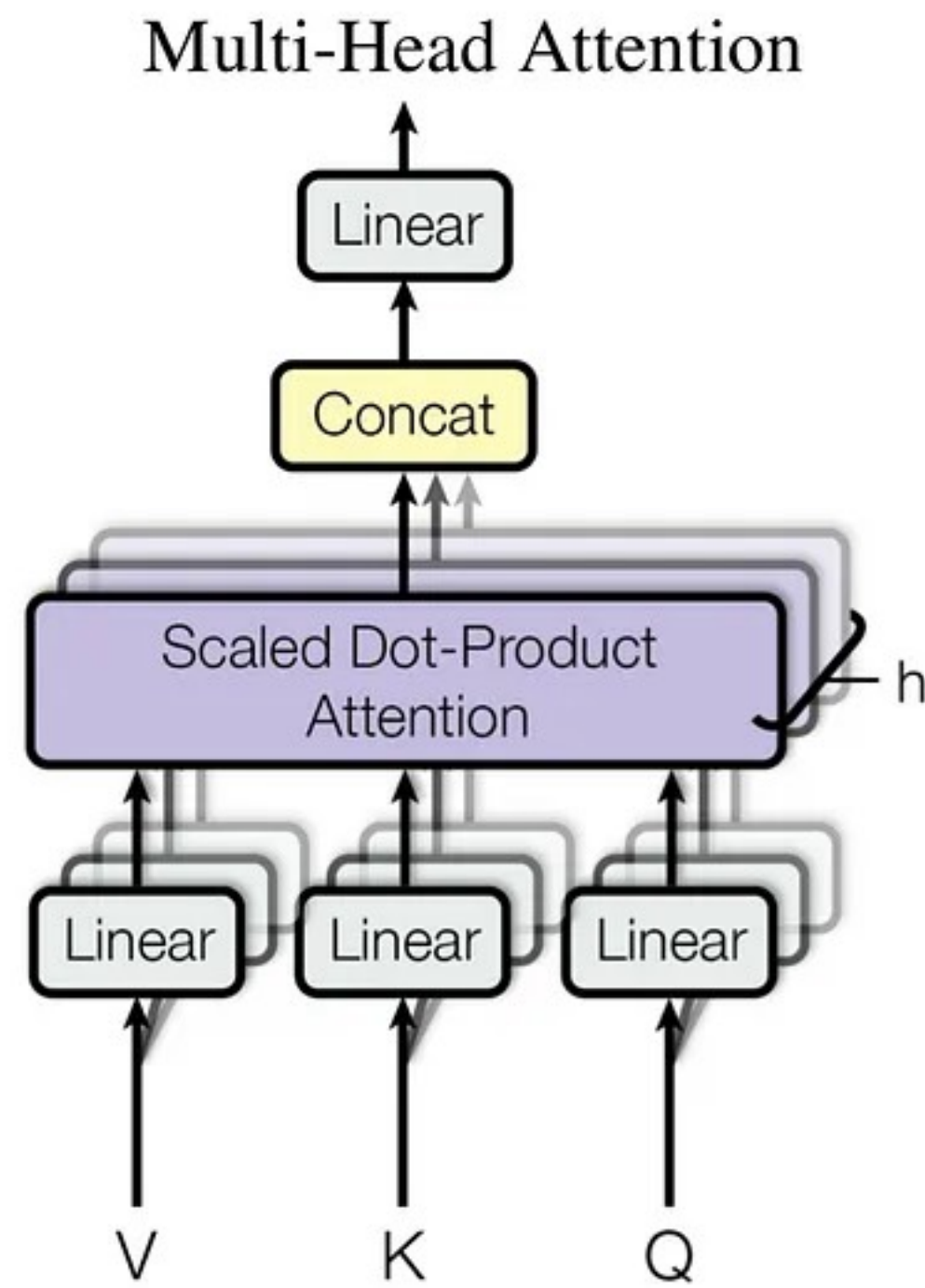
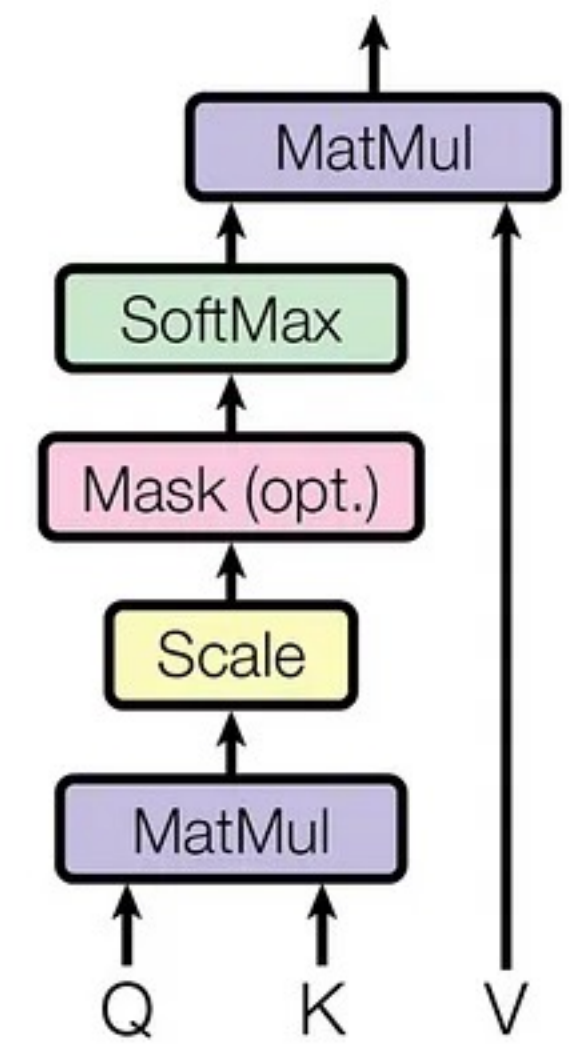


Figure 1: The Transformer - model architecture.



Scaled Dot-Product Attention



[Vaswani et.al. '17]

The transformer architecture

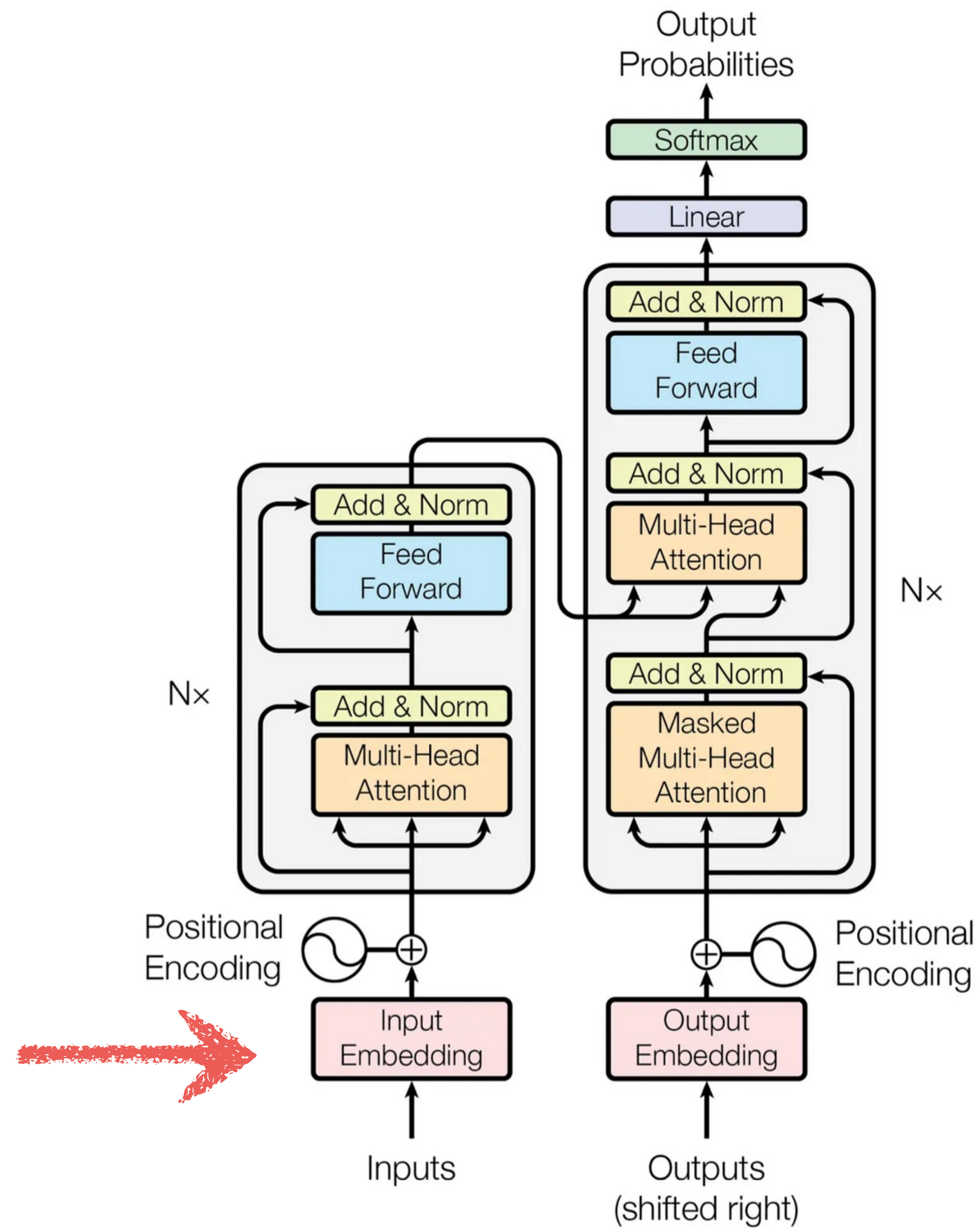


Figure 1: The Transformer - model architecture.

NLP and input embedding

- ▶ The goal of NLP is to produce a sequence of outputs from a sequence of inputs
- ▶ In order to feed words/characters of natural languages to a NN, we need to represent each word/character as a vector in $\mathbb{R}^{n_{in}}$
- ▶ A simple way of doing this is as follows:

- Tokenize the words/alphabet:

- ◆ Take a dictionary (with 10k unique English words, say) and enumerate each word from 1 to 10k.

- ◆ assign to each word a 10'000-dim vector, which has a 1 at position l and zeros everywhere else

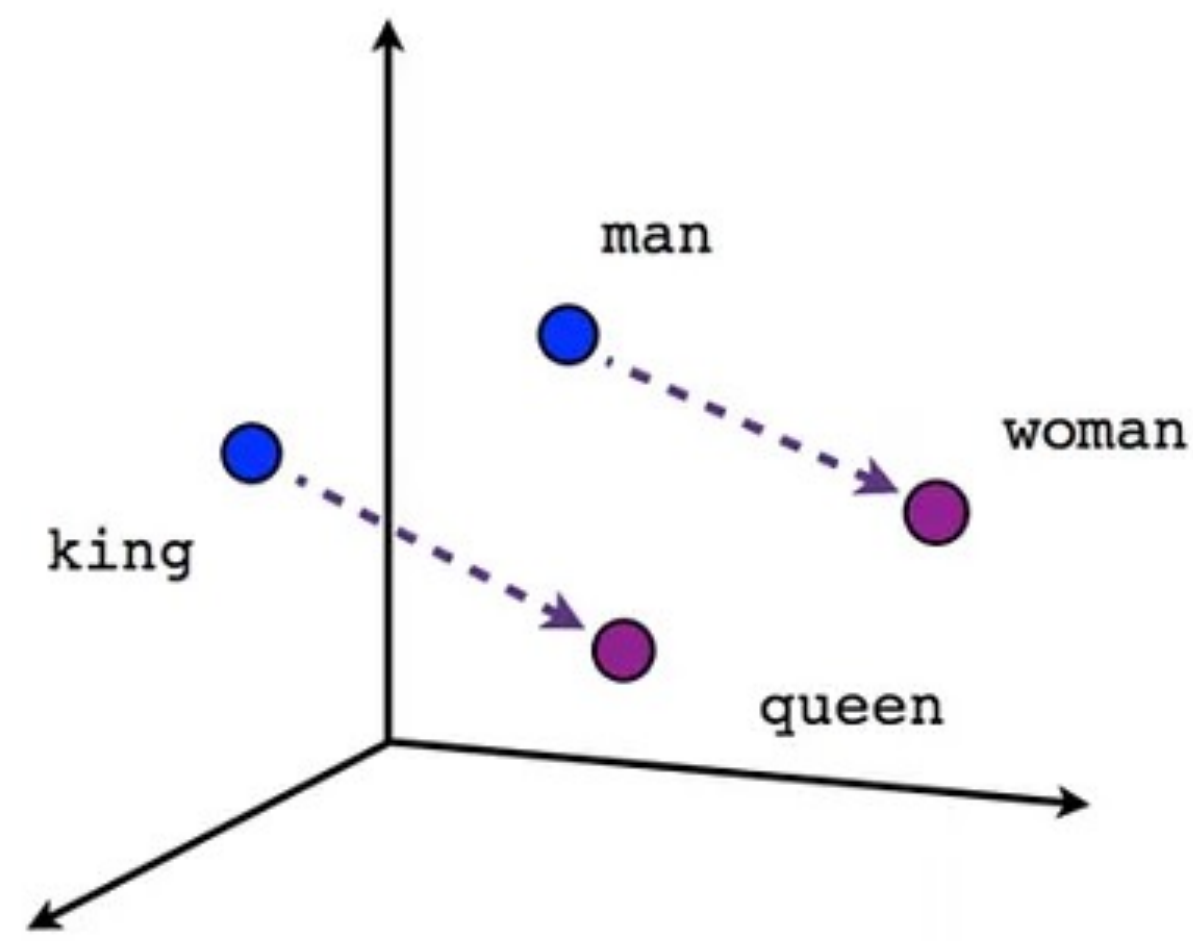
a $\boxed{1}\boxed{0}\boxed{0}\dots\boxed{0}\boxed{0}$ Aardvark $\boxed{0}\boxed{1}\boxed{0}\dots\boxed{0}\boxed{0}$... Zyzyyva $\boxed{0}\boxed{0}\boxed{0}\dots\boxed{0}\boxed{1}$

- Multiply this 10k-dim vector with a 512 x 10k matrix to produce a vector in \mathbb{R}^{512}

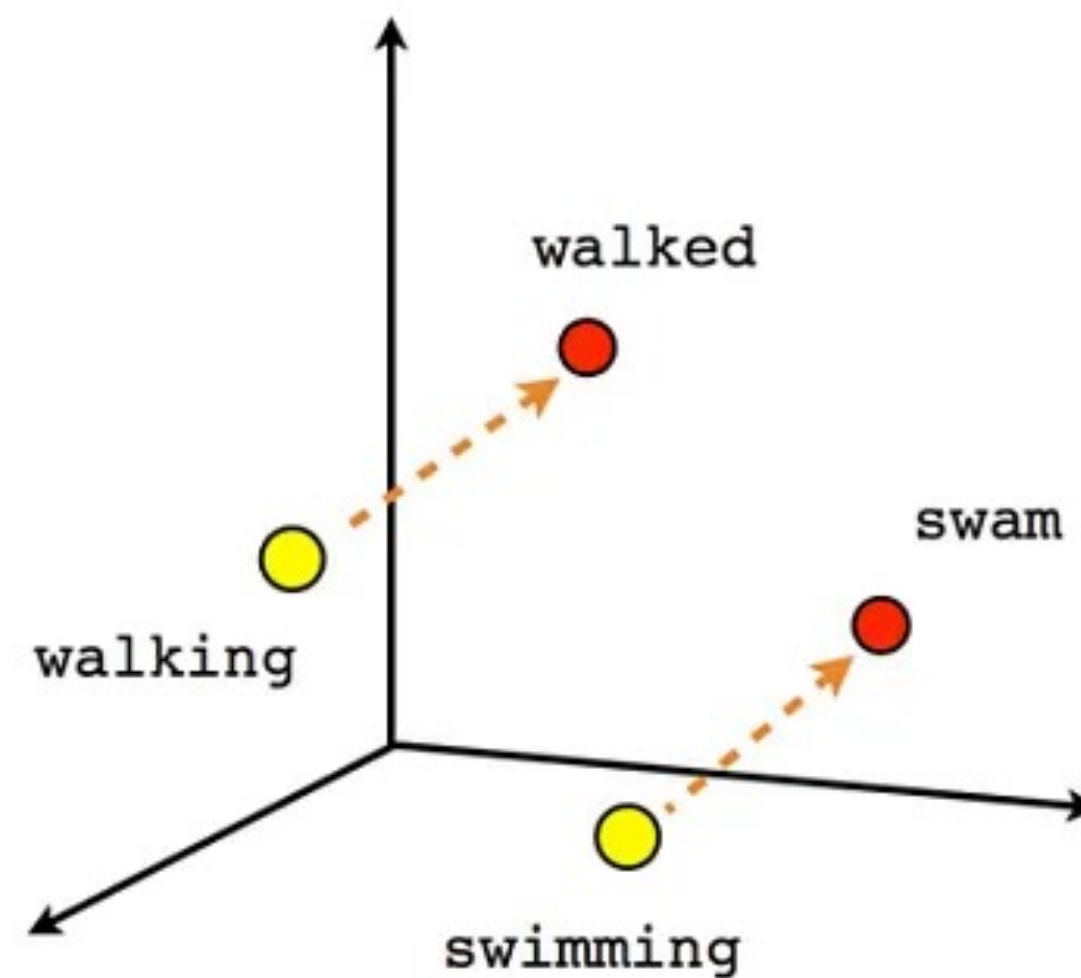
$$\begin{pmatrix} e_{1,1} & \dots & e_{10000,1} \\ \vdots & \vdots & \vdots \\ e_{1,512} & \dots & e_{10000,512} \end{pmatrix} \begin{pmatrix} \boxed{1} \\ \boxed{0} \\ \boxed{0} \\ \vdots \\ \boxed{0} \end{pmatrix} = \begin{pmatrix} e_{1,1} \\ \vdots \\ e_{1,512} \end{pmatrix} \qquad \begin{pmatrix} e_{1,1} & \dots & e_{10000,1} \\ \vdots & \vdots & \vdots \\ e_{1,512} & \dots & e_{10000,512} \end{pmatrix} \begin{pmatrix} \boxed{0} \\ \boxed{0} \\ \boxed{0} \\ \vdots \\ \boxed{1} \end{pmatrix} = \begin{pmatrix} e_{10000,1} \\ \vdots \\ e_{10000,512} \end{pmatrix}$$

- Learning the entries of this matrix allows to learn “useful” embeddings

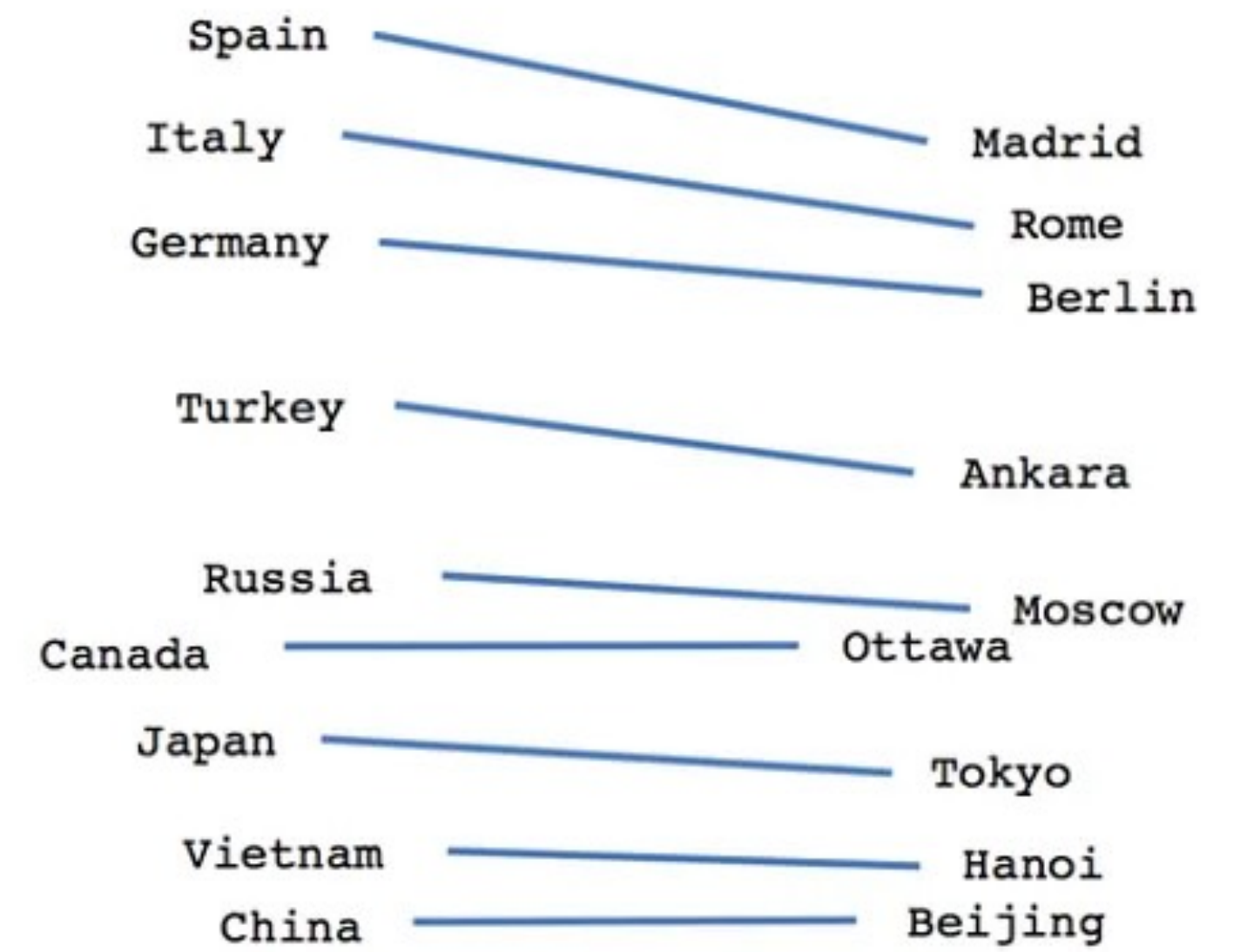
Examples of learned embeddings



Male-Female



Verb tense



Country-Capital

Learns semantics

“queen - woman + man = king”

[Image Credit: Tensorflow]

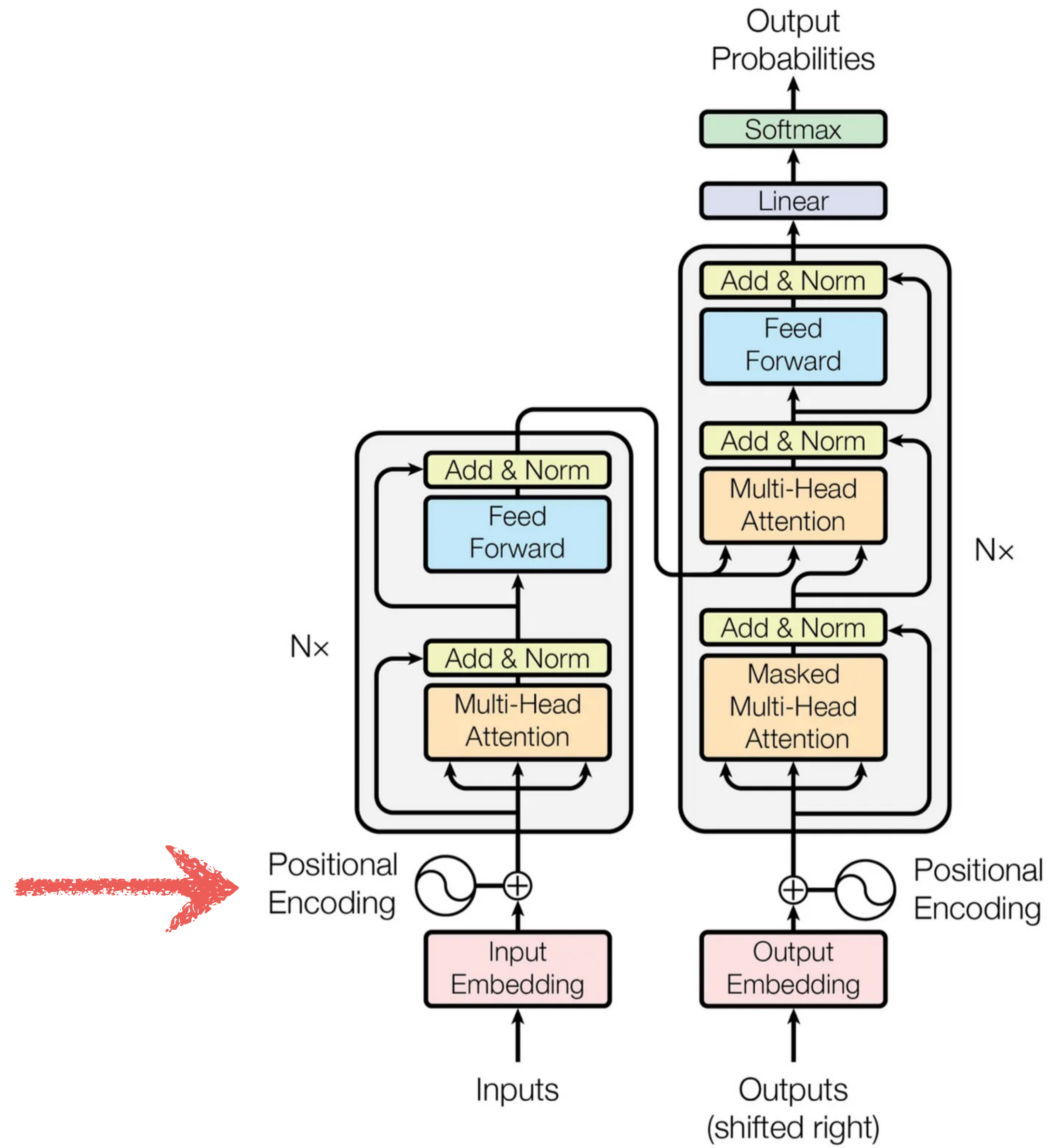


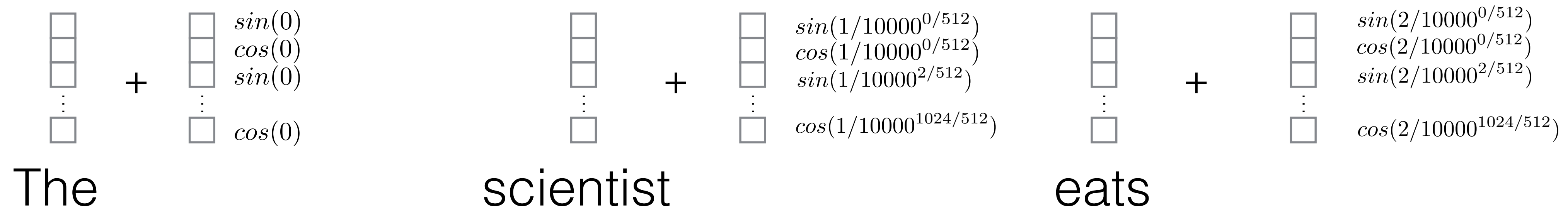
Figure 1: The Transformer - model architecture.

Positional embedding

- ▶ A fully connected layer has no notion of position: each node is connected to all other nodes
- ▶ But position of words in a sentence is very important:
The scientist eats the chicken \neq The chicken eats the scientist

- ▶ Add a (time-modulated) signal to each embedded word to signal its position in the sentence. The NN will learn to use it

$$PE_{(pos,2i)} = \sin(pos/10000^{2i/d_{\text{model}}}) \quad , \quad PE_{(pos,2i+1)} = \cos(pos/10000^{2i/d_{\text{model}}})$$



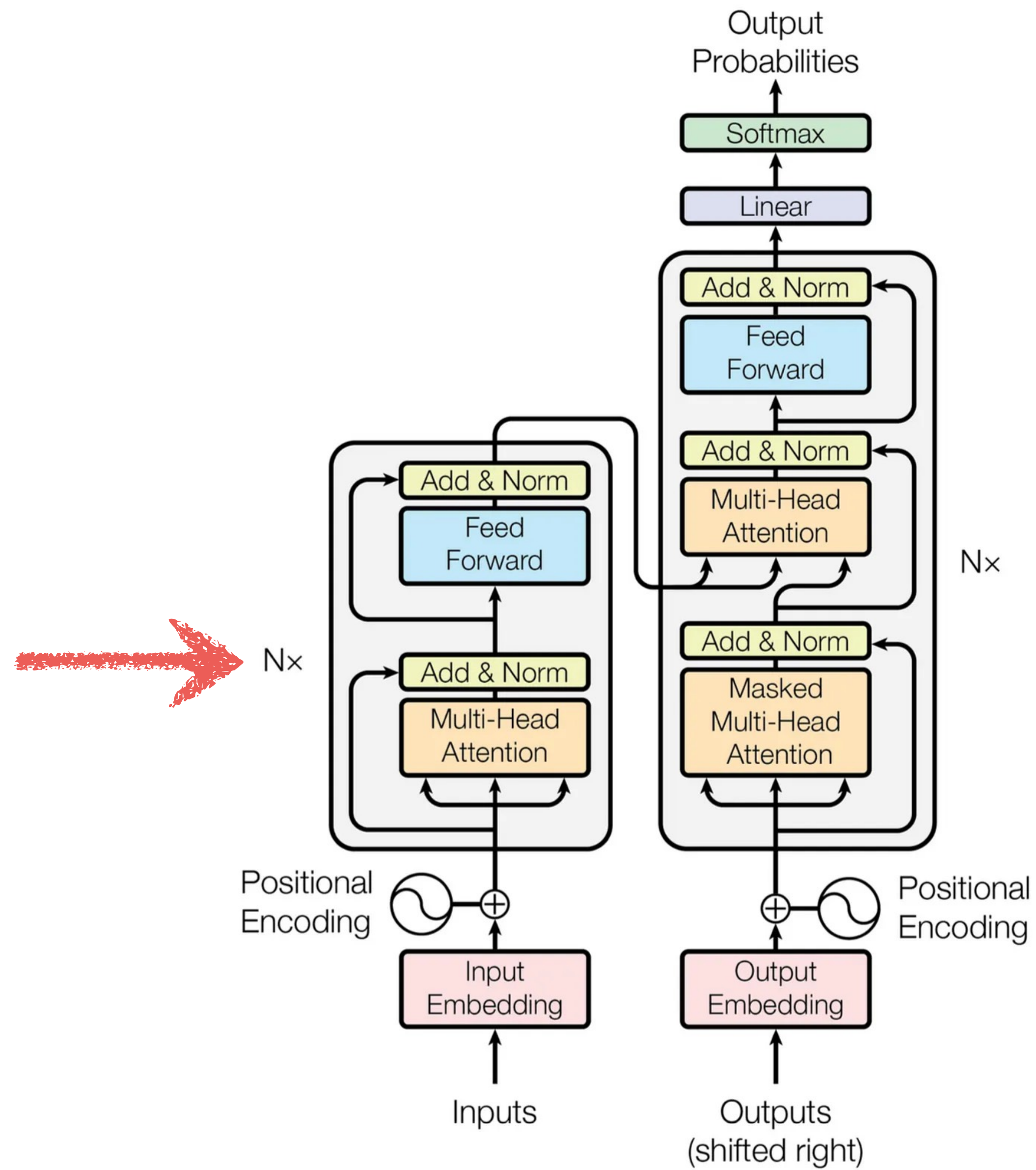
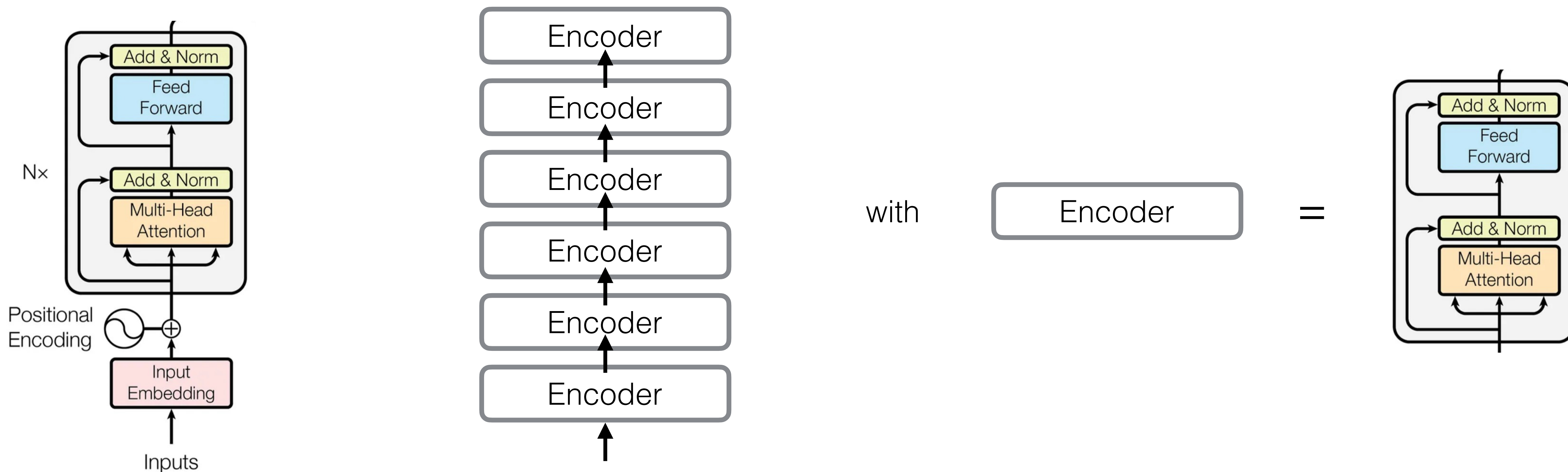


Figure 1: The Transformer - model architecture.

Encoding

- ▶ Now we have an encoded sequence of words with temporal information
- ▶ Just as in a deep NN, we now keep encoding this sequence over and over, feeding the output of the last encoder to the input of the next
- ▶ Repeat (6 times in the paper) until the NN has learned a good encoding



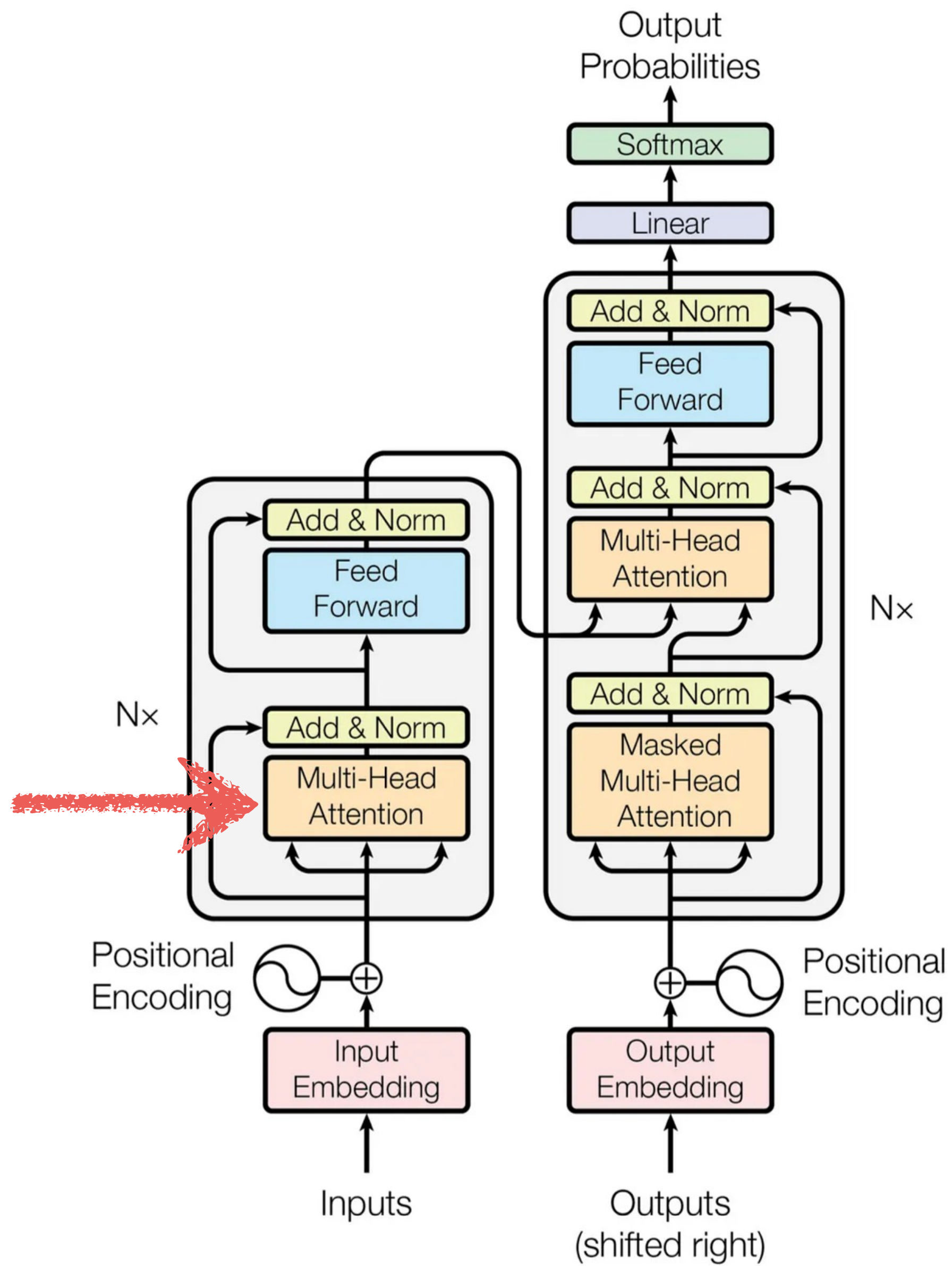
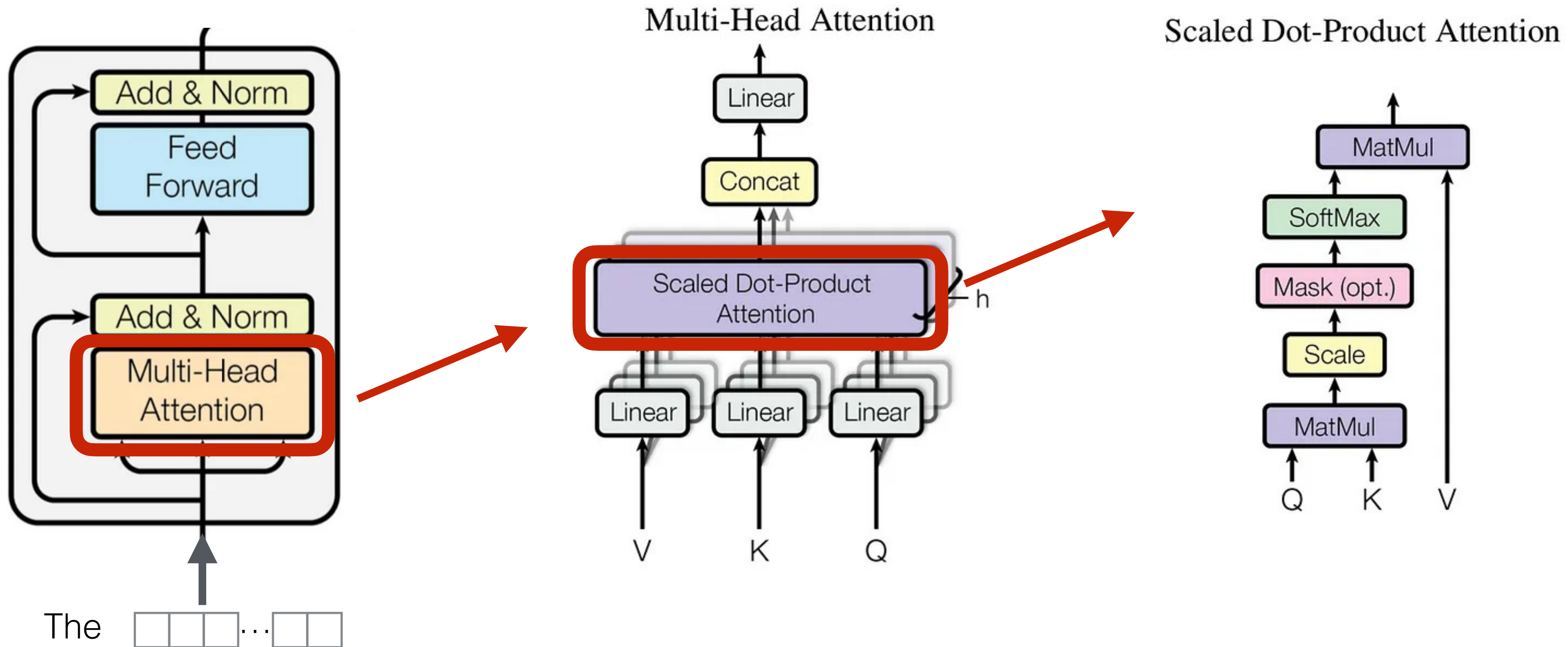


Figure 1: The Transformer - model architecture.

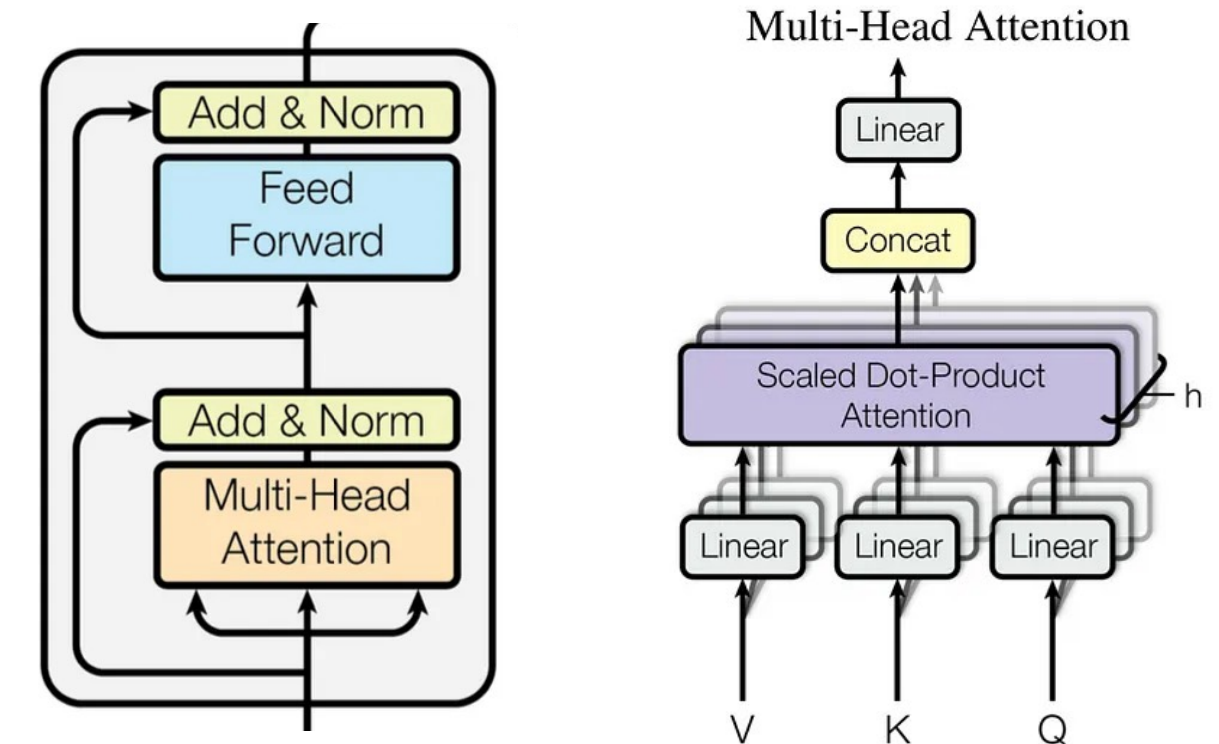
Self-attention

- ▶ So each encoder layer consists of attention heads and a vanilla feed-forward layer (and normalizations and skip connections for better training performance)



Self-attention I

- ▶ As illustrated, we create 3 copies of each input vector, denoted by V(value), K(ey), Q(uey)
- ▶ This is done for all input vector (i.e. each encoded word) in the sequence
- ▶ Then, each copy V, K, Q goes into a feed-forward layer (i.e., matrix multiplication, where the entries of the matrix are learnable parameters)

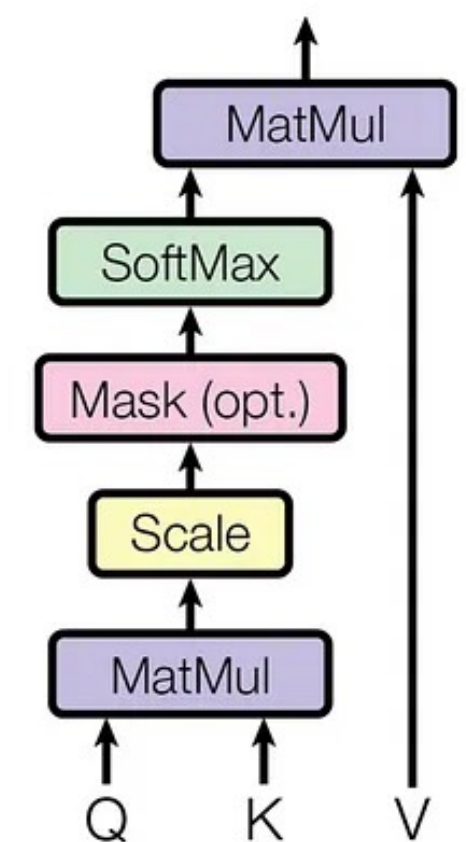


$\begin{pmatrix} v_{1,1} & \dots & v_{512,1} \\ \vdots & \vdots & \vdots \\ v_{1,64} & \dots & v_{512,64} \end{pmatrix} \begin{matrix} \square \\ \square \\ \square \end{matrix} \text{The} = \begin{pmatrix} v_1(\text{"The"}) \\ \vdots \\ v_{64}(\text{"The"}) \end{pmatrix}$	$\begin{pmatrix} v_{1,1} & \dots & v_{512,1} \\ \vdots & \vdots & \vdots \\ v_{1,64} & \dots & v_{512,64} \end{pmatrix} \begin{matrix} \square \\ \square \\ \square \end{matrix} \text{scientist} = \begin{pmatrix} v_1(\text{"scientist"}) \\ \vdots \\ v_{64}(\text{"scientist"}) \end{pmatrix}$	$\begin{pmatrix} v_{1,1} & \dots & v_{512,1} \\ \vdots & \vdots & \vdots \\ v_{1,64} & \dots & v_{512,64} \end{pmatrix} \begin{matrix} \square \\ \square \\ \square \end{matrix} \text{eats} = \begin{pmatrix} v_1(\text{"eats"}) \\ \vdots \\ v_{64}(\text{"eats"}) \end{pmatrix}$
$\begin{pmatrix} k_{1,1} & \dots & k_{512,1} \\ \vdots & \vdots & \vdots \\ k_{1,64} & \dots & k_{512,64} \end{pmatrix} \begin{matrix} \square \\ \square \\ \square \end{matrix} \text{The} = \begin{pmatrix} k_1(\text{"The"}) \\ \vdots \\ k_{64}(\text{"The"}) \end{pmatrix}$	$\begin{pmatrix} k_{1,1} & \dots & k_{512,1} \\ \vdots & \vdots & \vdots \\ k_{1,64} & \dots & k_{512,64} \end{pmatrix} \begin{matrix} \square \\ \square \\ \square \end{matrix} \text{scientist} = \begin{pmatrix} k_1(\text{"scientist"}) \\ \vdots \\ k_{64}(\text{"scientist"}) \end{pmatrix}$	$\begin{pmatrix} k_{1,1} & \dots & k_{512,1} \\ \vdots & \vdots & \vdots \\ k_{1,64} & \dots & k_{512,64} \end{pmatrix} \begin{matrix} \square \\ \square \\ \square \end{matrix} \text{eats} = \begin{pmatrix} k_1(\text{"eats"}) \\ \vdots \\ k_{64}(\text{"eats"}) \end{pmatrix}$
$\begin{pmatrix} q_{1,1} & \dots & q_{512,1} \\ \vdots & \vdots & \vdots \\ q_{1,64} & \dots & q_{512,64} \end{pmatrix} \begin{matrix} \square \\ \square \\ \square \end{matrix} \text{The} = \begin{pmatrix} q_1(\text{"The"}) \\ \vdots \\ q_{64}(\text{"The"}) \end{pmatrix}$	$\begin{pmatrix} q_{1,1} & \dots & q_{512,1} \\ \vdots & \vdots & \vdots \\ q_{1,64} & \dots & q_{512,64} \end{pmatrix} \begin{matrix} \square \\ \square \\ \square \end{matrix} \text{scientist} = \begin{pmatrix} q_1(\text{"scientist"}) \\ \vdots \\ q_{64}(\text{"scientist"}) \end{pmatrix}$	$\begin{pmatrix} q_{1,1} & \dots & q_{512,1} \\ \vdots & \vdots & \vdots \\ q_{1,64} & \dots & q_{512,64} \end{pmatrix} \begin{matrix} \square \\ \square \\ \square \end{matrix} \text{eats} = \begin{pmatrix} q_1(\text{"eats"}) \\ \vdots \\ q_{64}(\text{"eats"}) \end{pmatrix}$

Self-attention II

- ▶ The set of 64-dim vectors $(q(\text{word}_i), k(\text{word}_i), v(\text{word}_i))$ are then used to compute a score between the current word and all other words in the input sentence.
- ▶ The higher the score the more relevant the other word is for understanding the current word
- ▶ This score (or attention) is computed as follows:
 - Take the dot product (think “overlap” of the encoded words) between the query and the key for all words in the input sentence
 - Scale it down by $\sqrt{64}$ to avoid too big gradients, and optionally mask future inputs such that you can only pay attention to what has been seen already (for decoder)
 - Take the softmax to convert the “overlaps” to probabilities
 - The resulting score tells you how relevant other words in the sentence are for understanding the current word
 - Multiply the value of all words with this score and sum them all up to get a (superposition of) words which one should pay attention to for understanding the first word. Rinse and repeat

Scaled Dot-Product Attention



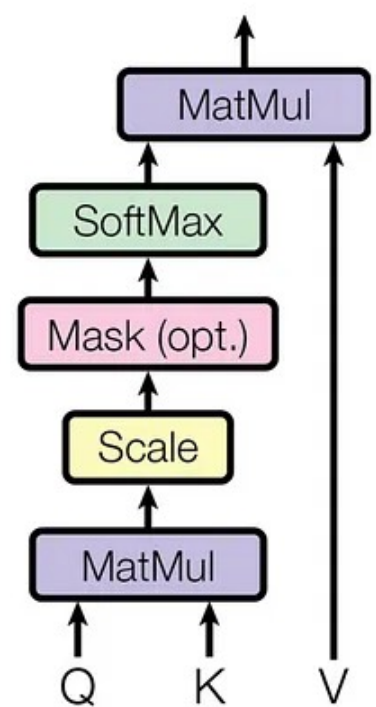
Self-attention III

- ▶ In matrix notation, the scaled dot-product attention head is

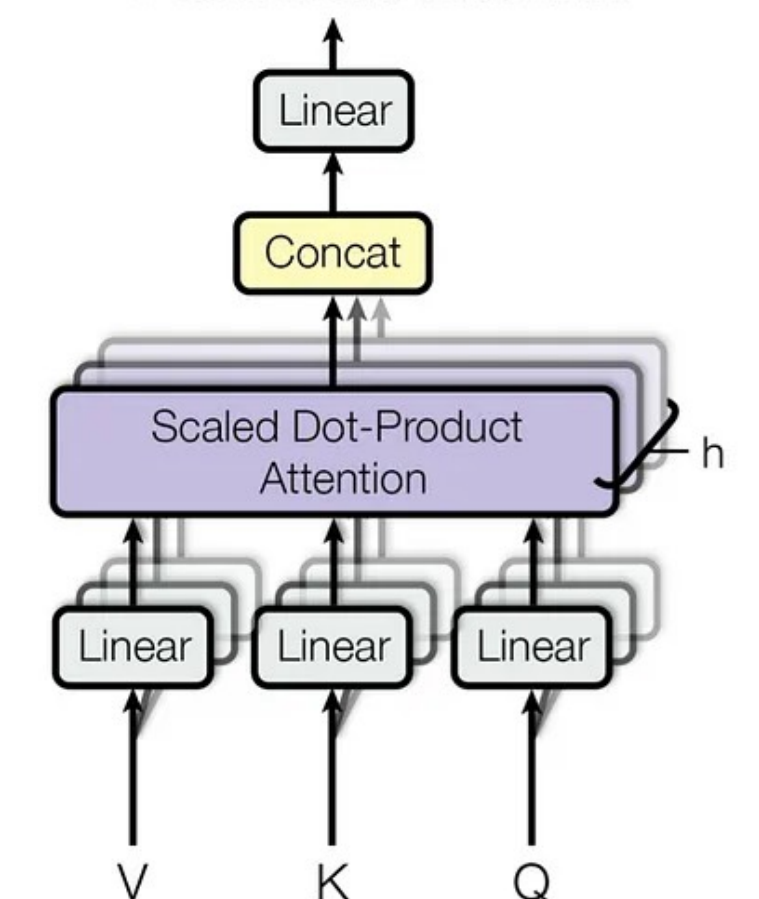
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

- ▶ Repeat this attention step in parallel for 8 different (learnable) matrices Q, K, V (“multi-heads”)
- ▶ The resulting 64-dimensional vectors from each head are put into one large vector of dimension $64 \times 8 = 512$
- ▶ This vector is fed into a fully-connected layer (with output dimension 512, so that it can fit into the next encoding layer with the same architecture

Scaled Dot-Product Attention



Multi-Head Attention



Self-attention example

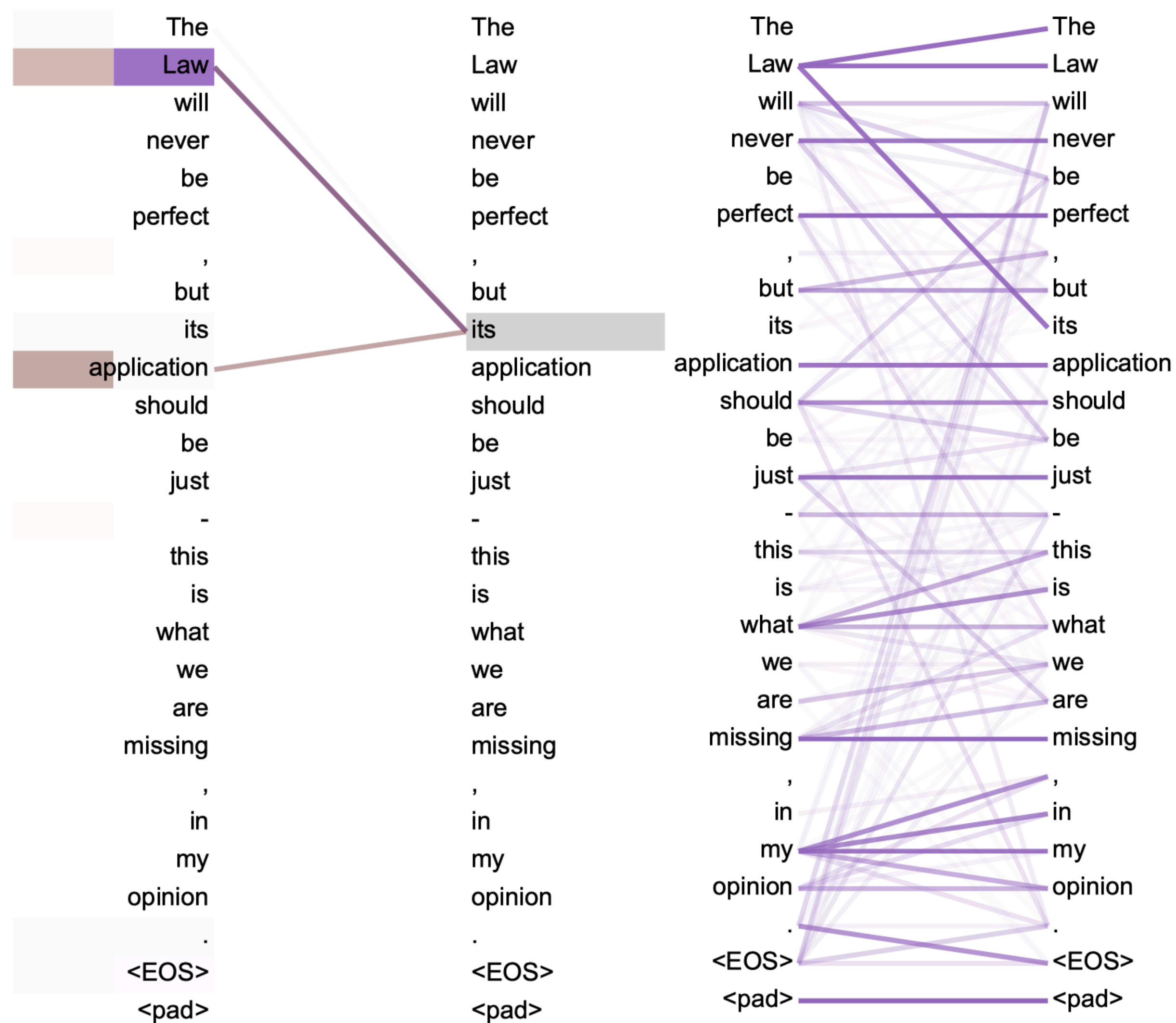


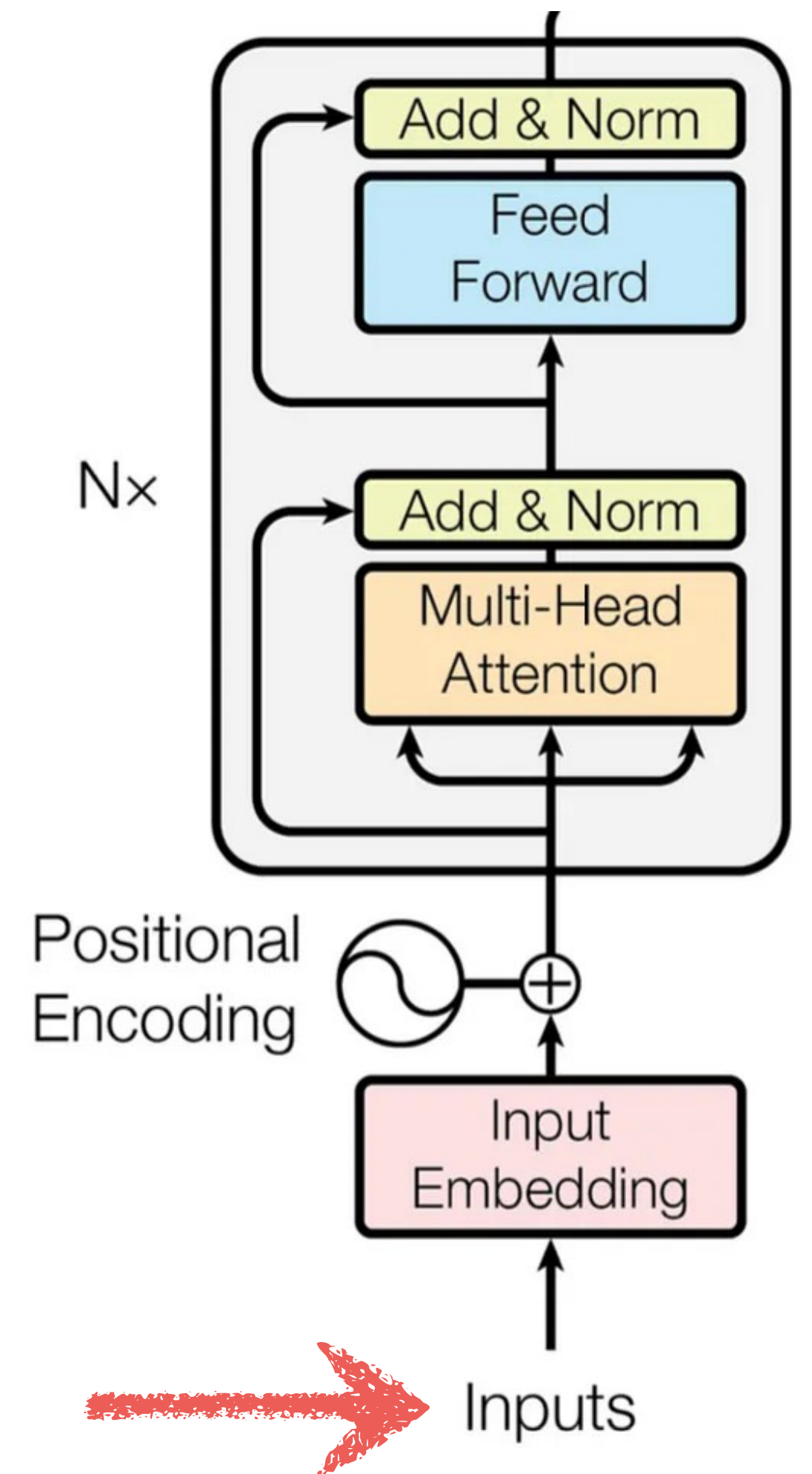
Figure 4: Two attention heads, also in layer 5 of 6, apparently involved in anaphora resolution.

Left: Isolated attentions from just the word 'its' for attention heads 5 and 6. Note that the attentions are very sharp for this word.

Right: Full attentions for head 5.

Summary of encoder

[The, scientist, eats, the, chicken, eos]



Summary of encoder

$$\begin{pmatrix} e_{1,1} & \dots & e_{10000,1} \\ \vdots & \vdots & \vdots \\ e_{1,512} & \dots & e_{10000,512} \end{pmatrix} \begin{matrix} \square \\ \square \\ \vdots \\ \square \end{matrix} = \begin{pmatrix} e_1(\text{"The"}) \\ \vdots \\ e_{512}(\text{"The"}) \end{pmatrix}$$

The

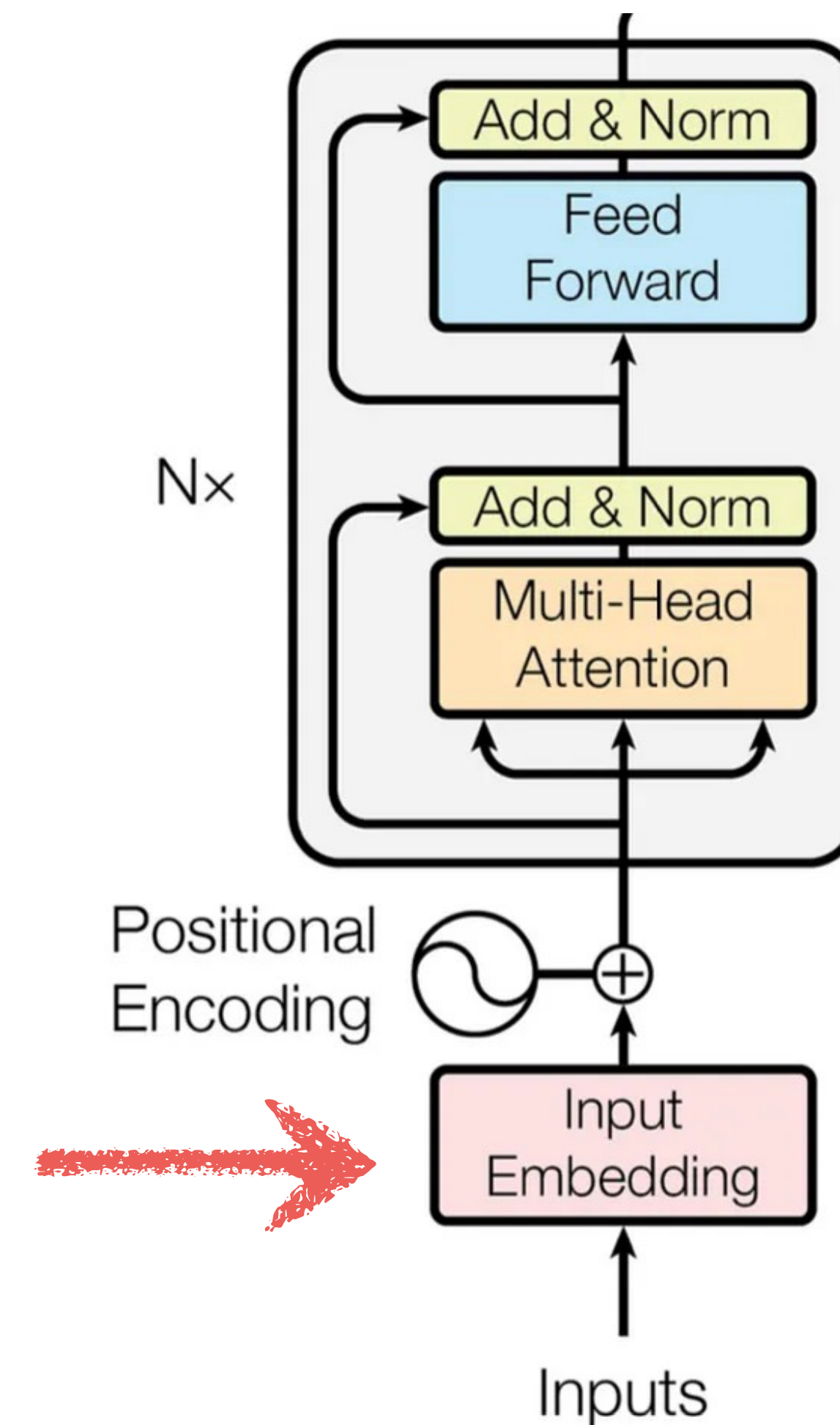
$$\begin{pmatrix} e_{1,1} & \dots & e_{10000,1} \\ \vdots & \vdots & \vdots \\ e_{1,512} & \dots & e_{10000,512} \end{pmatrix} \begin{matrix} \square \\ \square \\ \vdots \\ \square \end{matrix} = \begin{pmatrix} e_1(\text{"scientist"}) \\ \vdots \\ e_{512}(\text{"scientist"}) \end{pmatrix}$$

scientist

$$\begin{pmatrix} e_{1,1} & \dots & e_{10000,1} \\ \vdots & \vdots & \vdots \\ e_{1,512} & \dots & e_{10000,512} \end{pmatrix} \begin{matrix} \square \\ \square \\ \vdots \\ \square \end{matrix} = \begin{pmatrix} e_1(\text{"eats"}) \\ \vdots \\ e_{512}(\text{"eats"}) \end{pmatrix}$$

eats

⋮

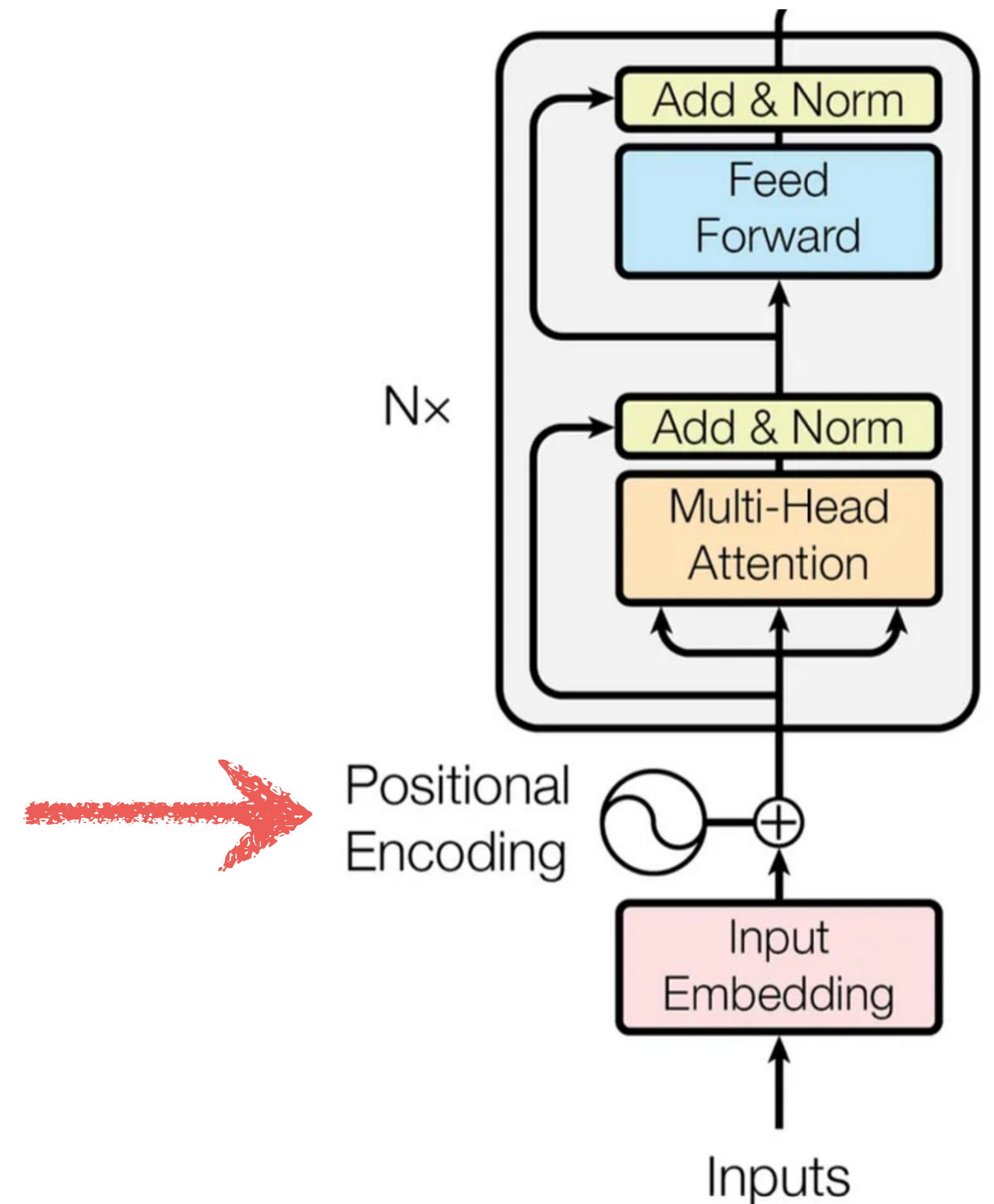


Summary of encoder

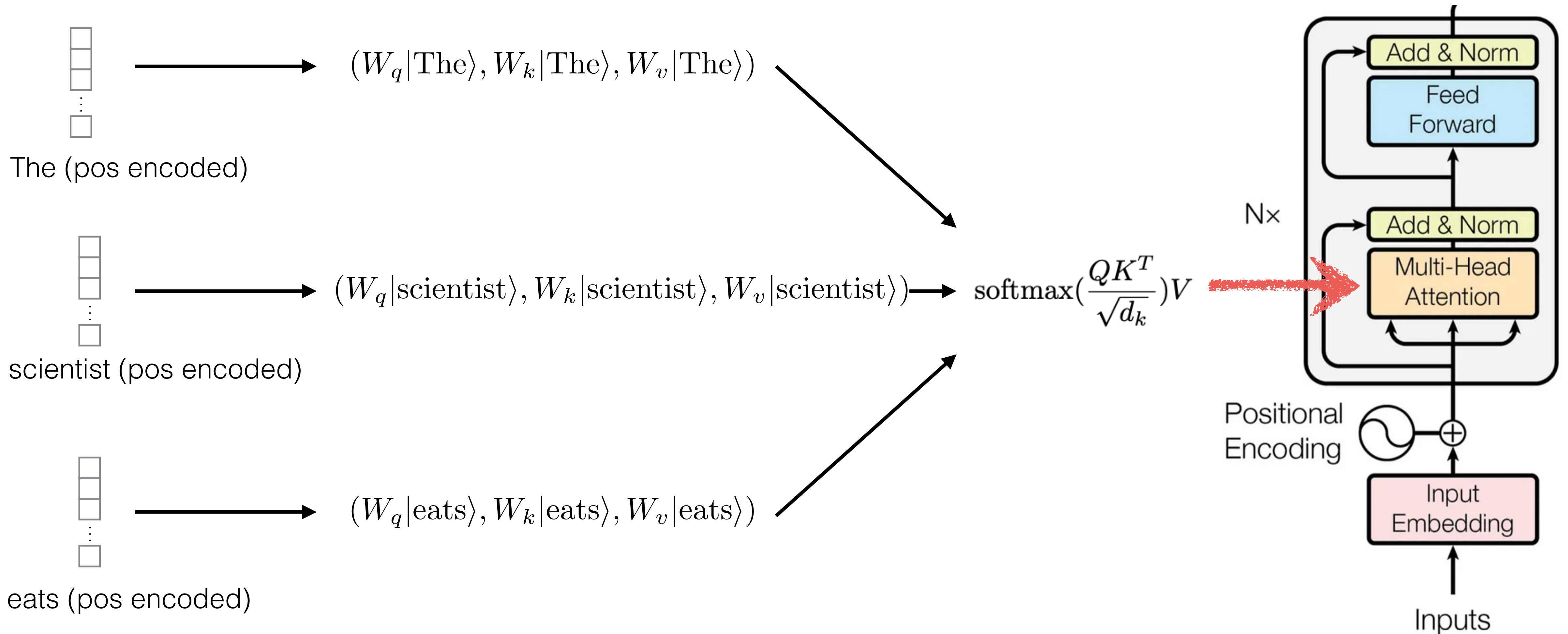
$$\begin{pmatrix} e_1(\text{"The"}) \\ \vdots \\ e_{512}(\text{"The"}) \end{pmatrix} + \begin{pmatrix} \square \sin(0) \\ \square \cos(0) \\ \square \sin(0) \\ \vdots \\ \square \cos(0) \end{pmatrix} = \begin{pmatrix} \square \\ \square \\ \square \\ \vdots \\ \square \end{pmatrix} \text{"The"}$$

$$\begin{pmatrix} e_1(\text{"scientist"}) \\ \vdots \\ e_{512}(\text{"scientist"}) \end{pmatrix} + \begin{pmatrix} \square \sin(1/10000^{0/512}) \\ \square \cos(1/10000^{0/512}) \\ \square \sin(1/10000^{2/512}) \\ \vdots \\ \square \cos(1/10000^{1024/512}) \end{pmatrix} = \begin{pmatrix} \square \\ \square \\ \square \\ \vdots \\ \square \end{pmatrix} \text{"scientist"}$$

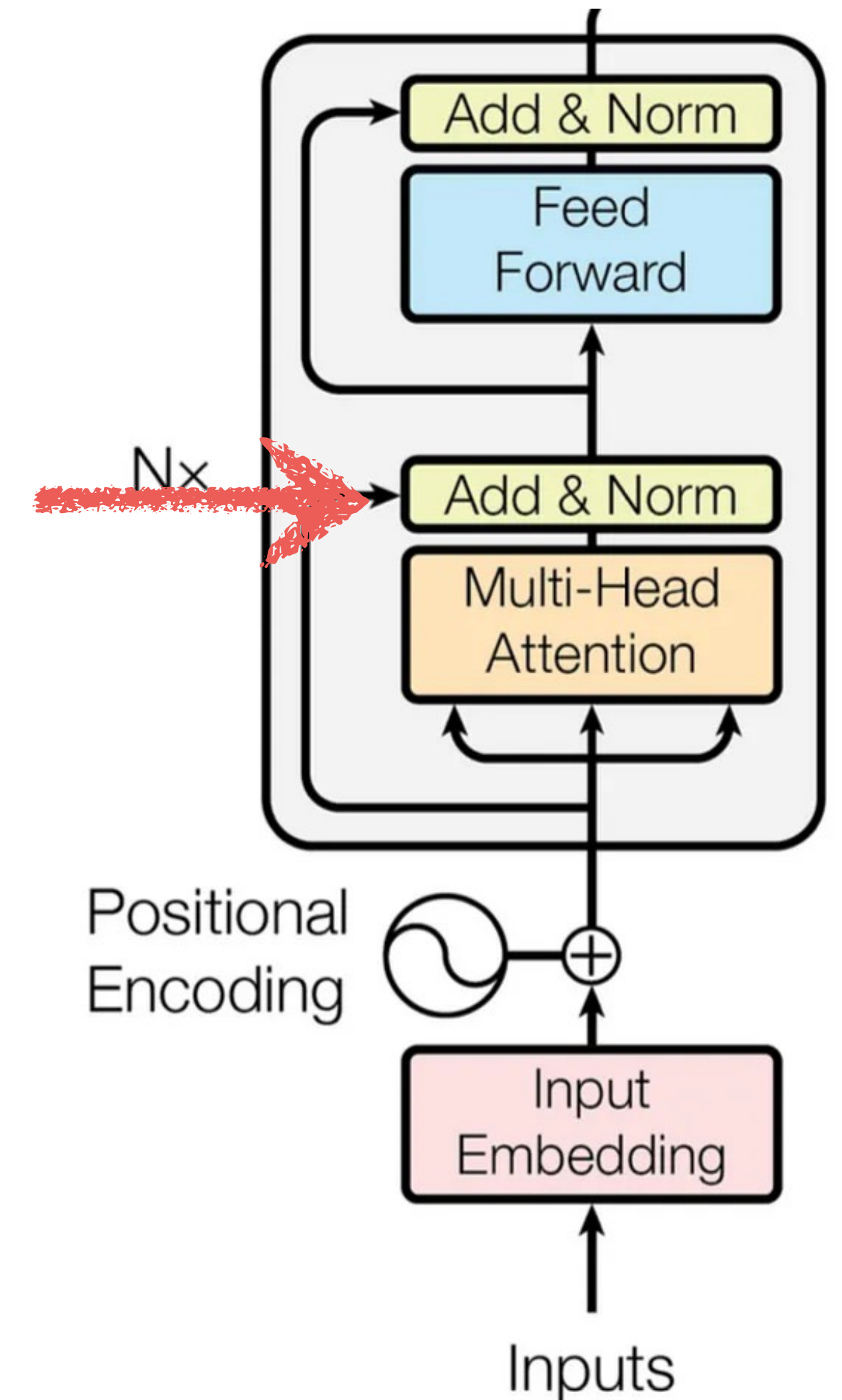
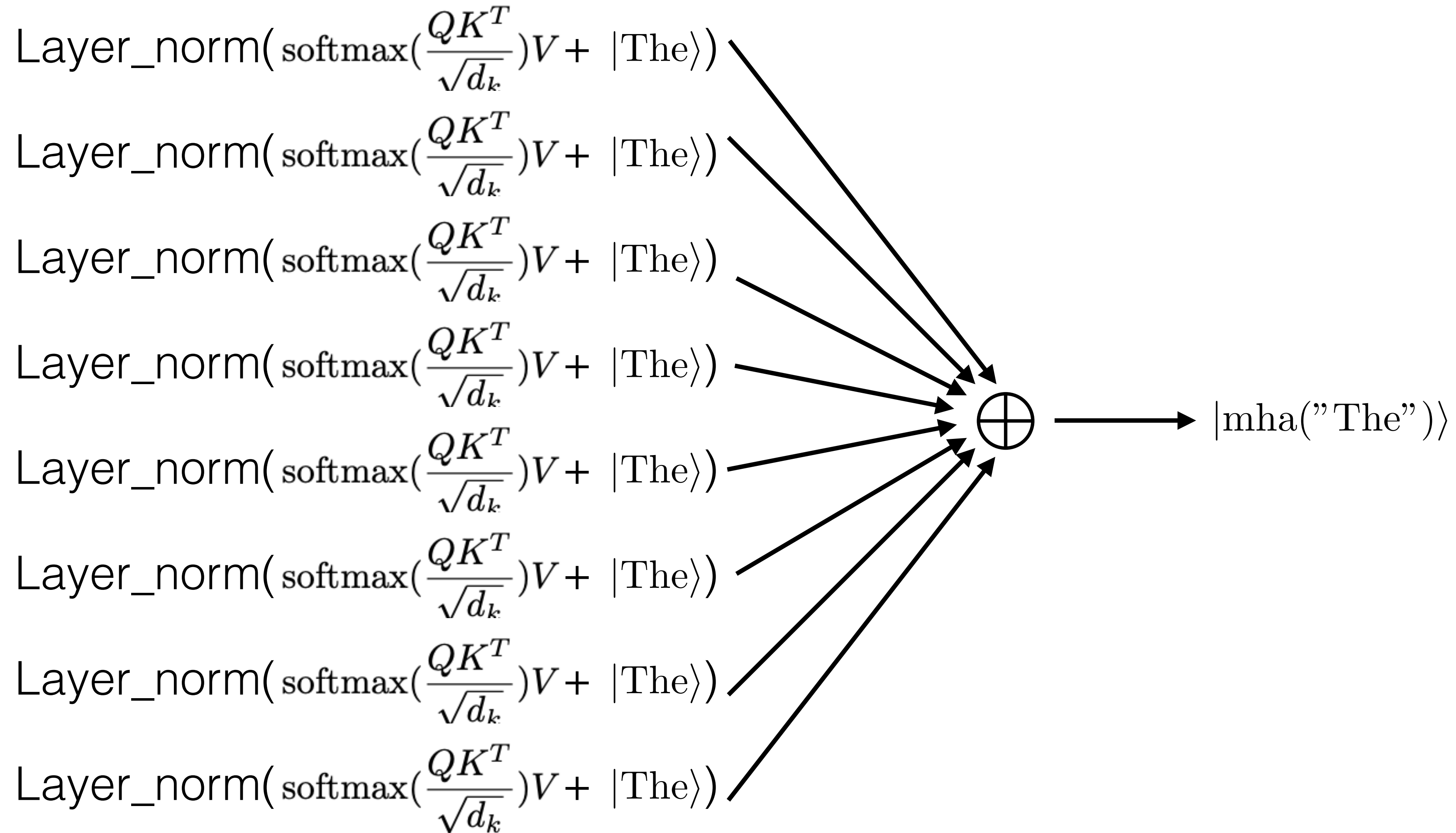
$$\begin{pmatrix} e_1(\text{"eats"}) \\ \vdots \\ e_{512}(\text{"eats"}) \end{pmatrix} + \begin{pmatrix} \square \sin(2/10000^{0/512}) \\ \square \cos(2/10000^{0/512}) \\ \square \sin(2/10000^{2/512}) \\ \vdots \\ \square \cos(2/10000^{1024/512}) \end{pmatrix} = \begin{pmatrix} \square \\ \square \\ \square \\ \vdots \\ \square \end{pmatrix} \text{"eats"}$$



Summary of encoder

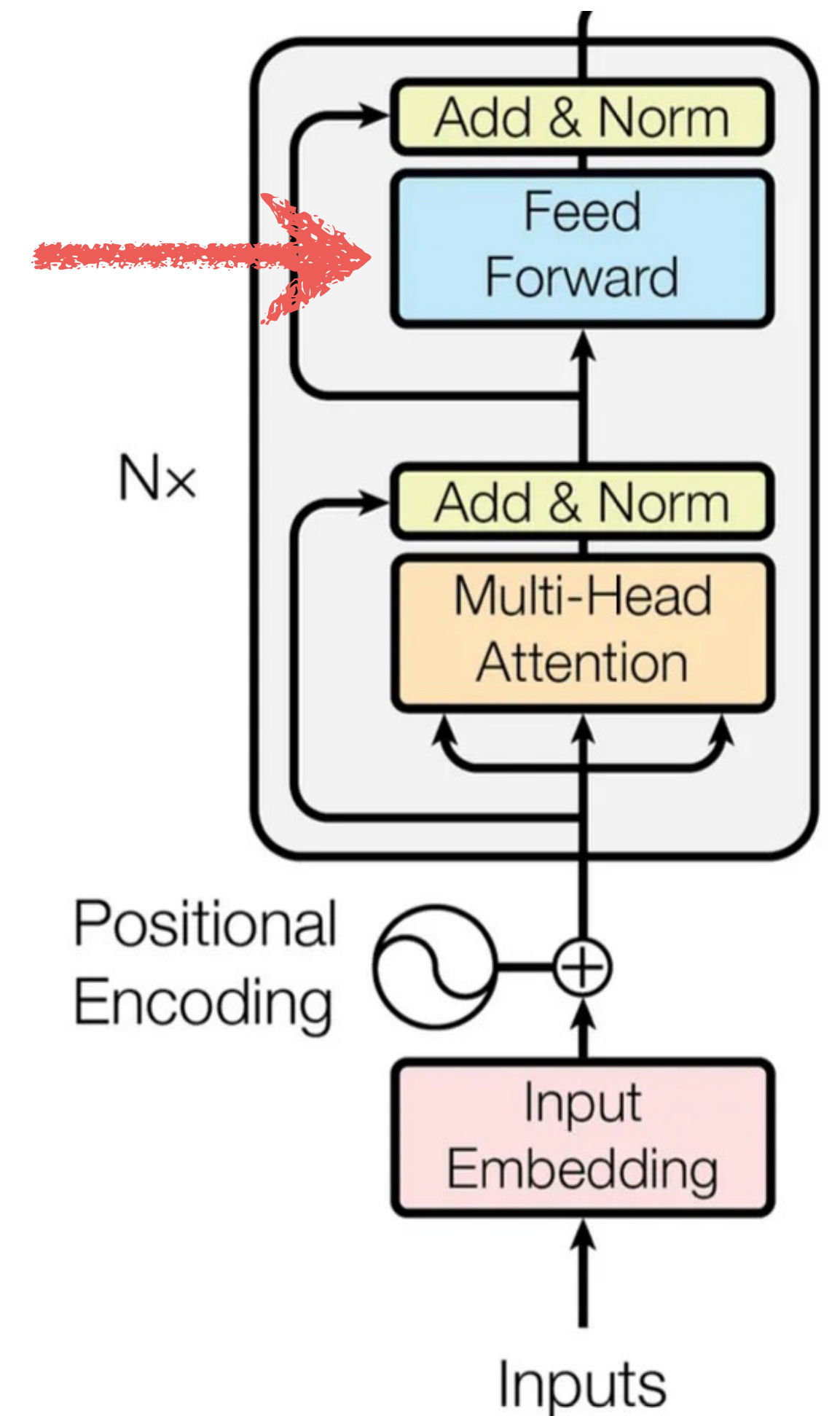


Summary of encoder



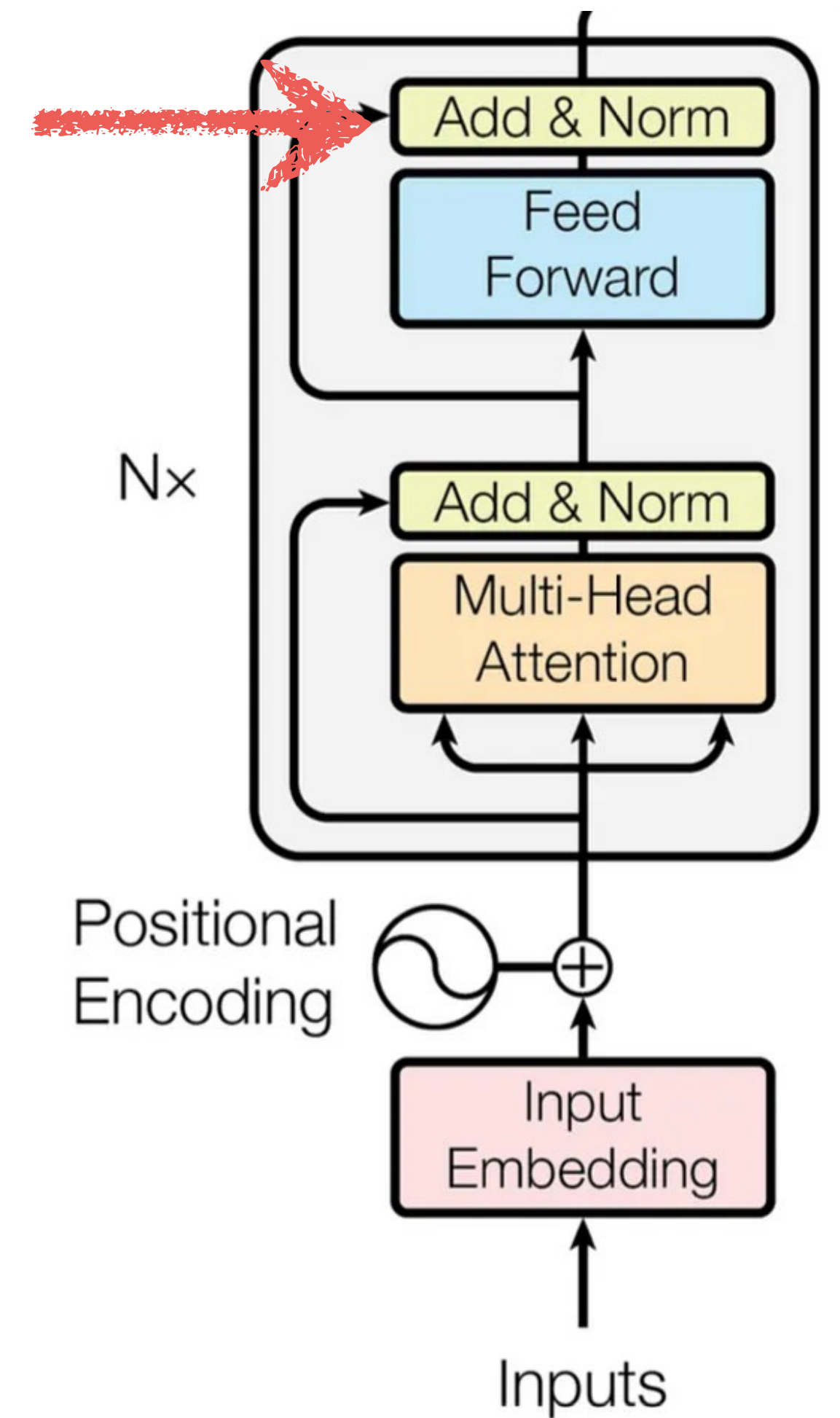
Summary of encoder

$$W_2 \text{ReLU}(W_1 | \text{mha}(\text{"The"}) \rangle + b_1) + b_2 = | \text{tmha}(\text{"The"}) \rangle$$



Summary of encoder

$\text{Layer_norm}(\text{tmha}(\text{"The"}) + \text{mha}(\text{"The"}))$



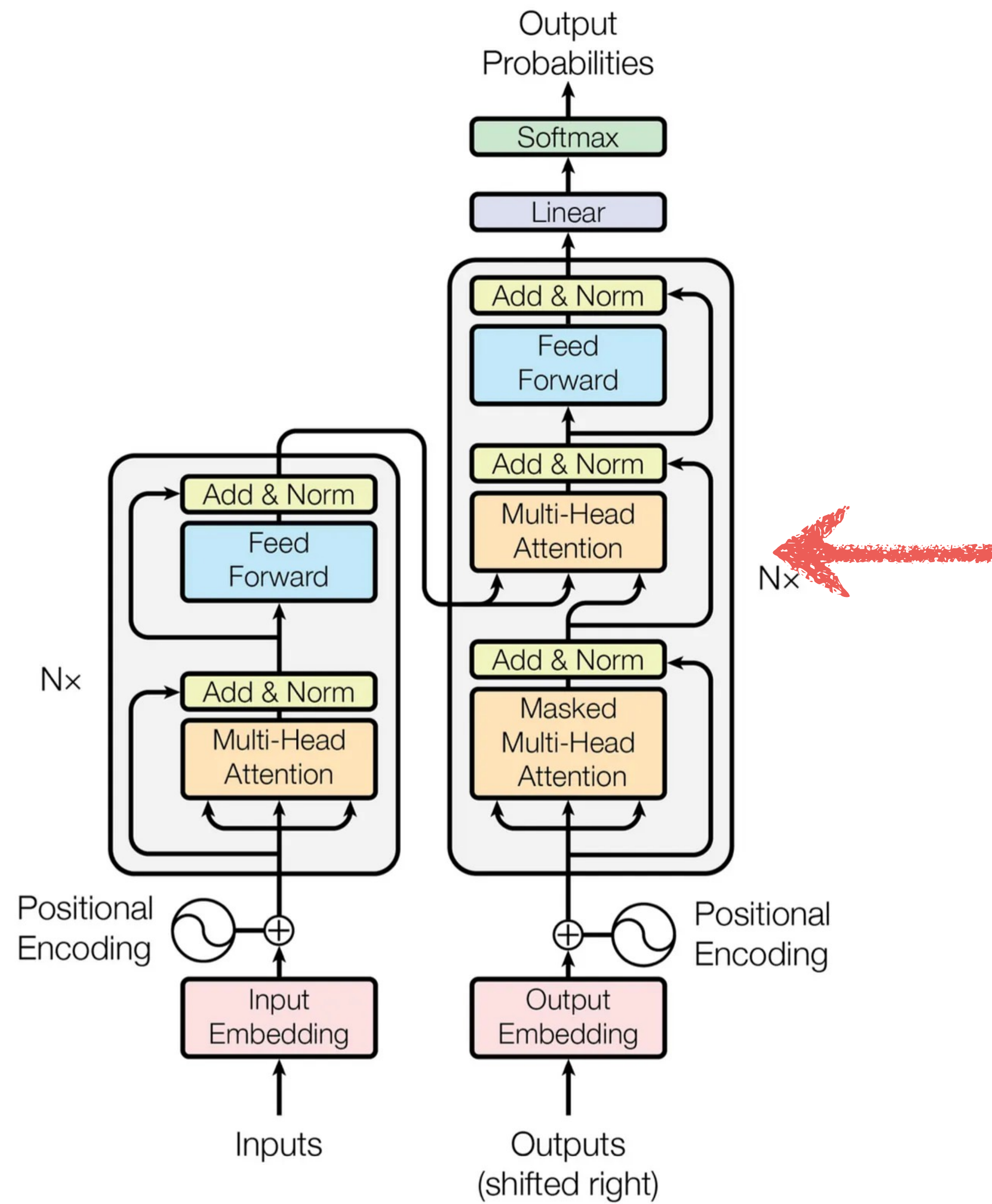


Figure 1: The Transformer - model architecture.

Decoder

- ▶ In a sequence to sequence task like translation, the decoder needs its own output and the (encoded) input it should translate
- ▶ The first multi-head attention computes an attention score+embedding from its own previous outputs

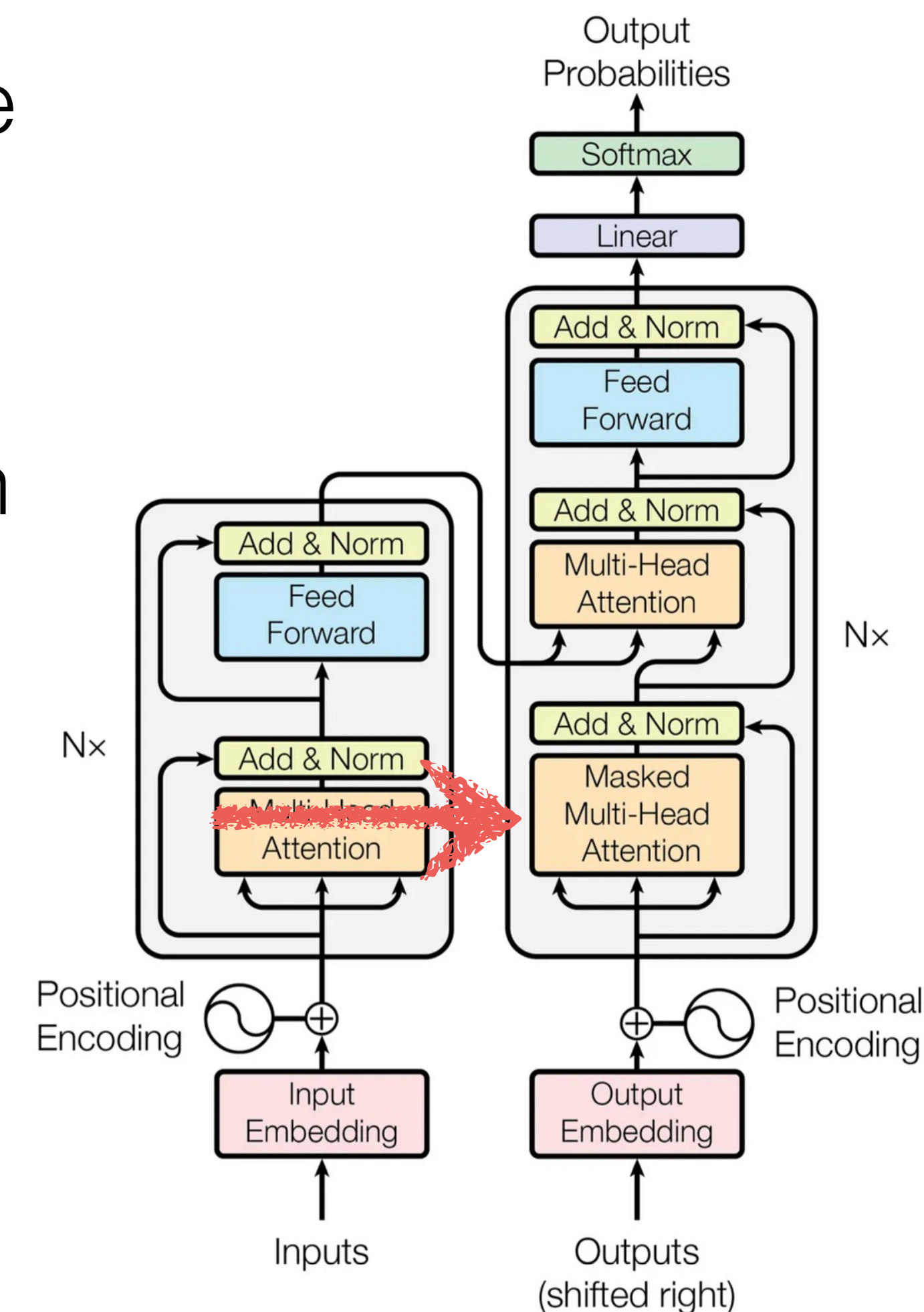


Figure 1: The Transformer - model architecture.

Decoder

- ▶ In a sequence to sequence task like translation, the decoder needs its own output and the (encoded) input it should translate
- ▶ The first multi-head attention computes an attention score+embedding from its own previous outputs
- ▶ The second multi-head attention takes this as the query and uses the output of the next word from the encoder as key and value
- ▶ Based on this, it computes another attention score, which is then used again in a FFNN

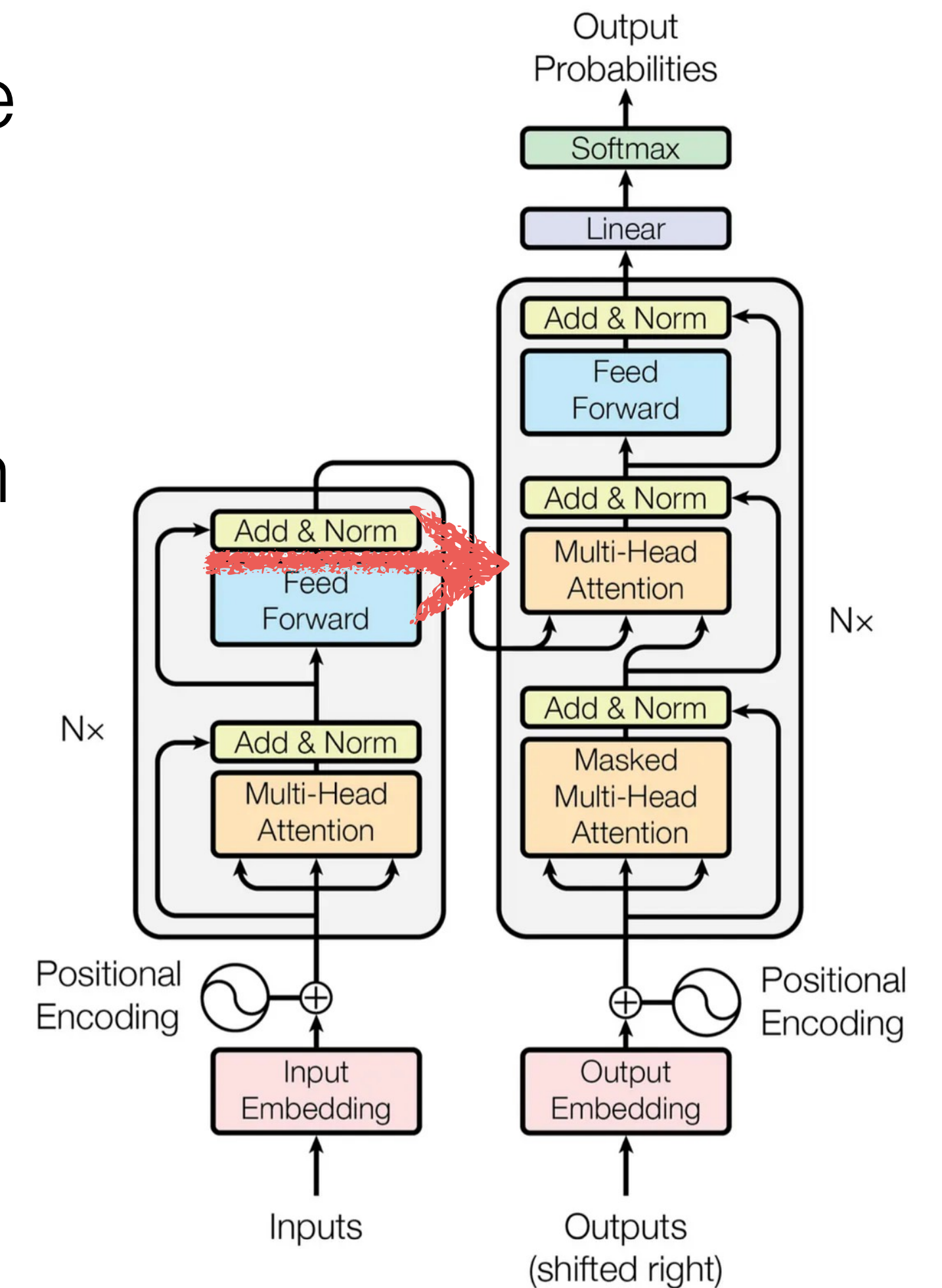


Figure 1: The Transformer - model architecture.

Decoder

- ▶ After repeating this decoding step 6 times, the 512dim output is put into a linear layer of dimension $10k \times 512$

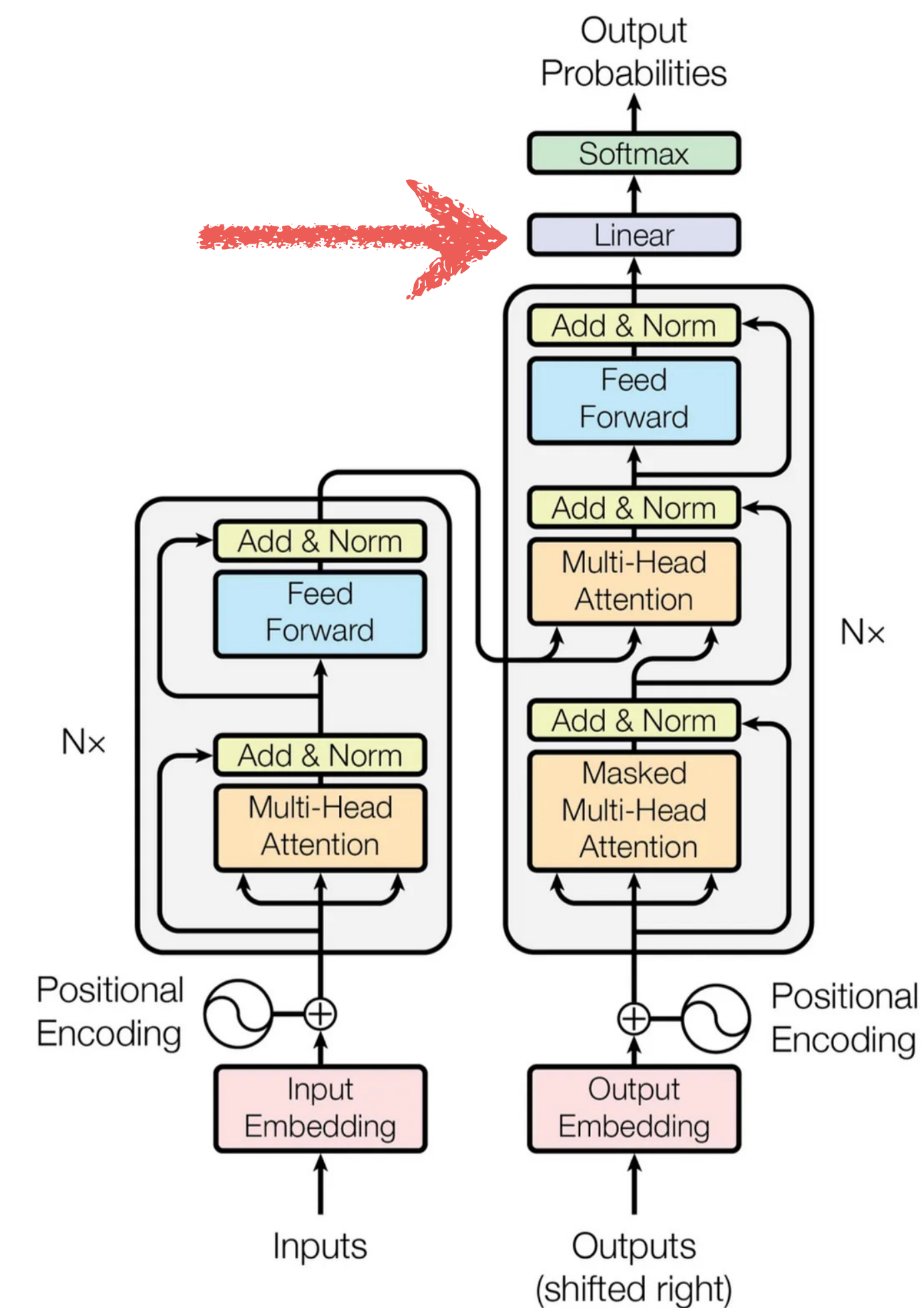


Figure 1: The Transformer - model architecture.

Decoder

- ▶ After repeating this decoding step 6 times, the 512dim output is put into a linear layer of dimension $10k \times 512$
- ▶ Then we take the softmax to obtain a $10k$ -dim vector of probabilities which word to use next

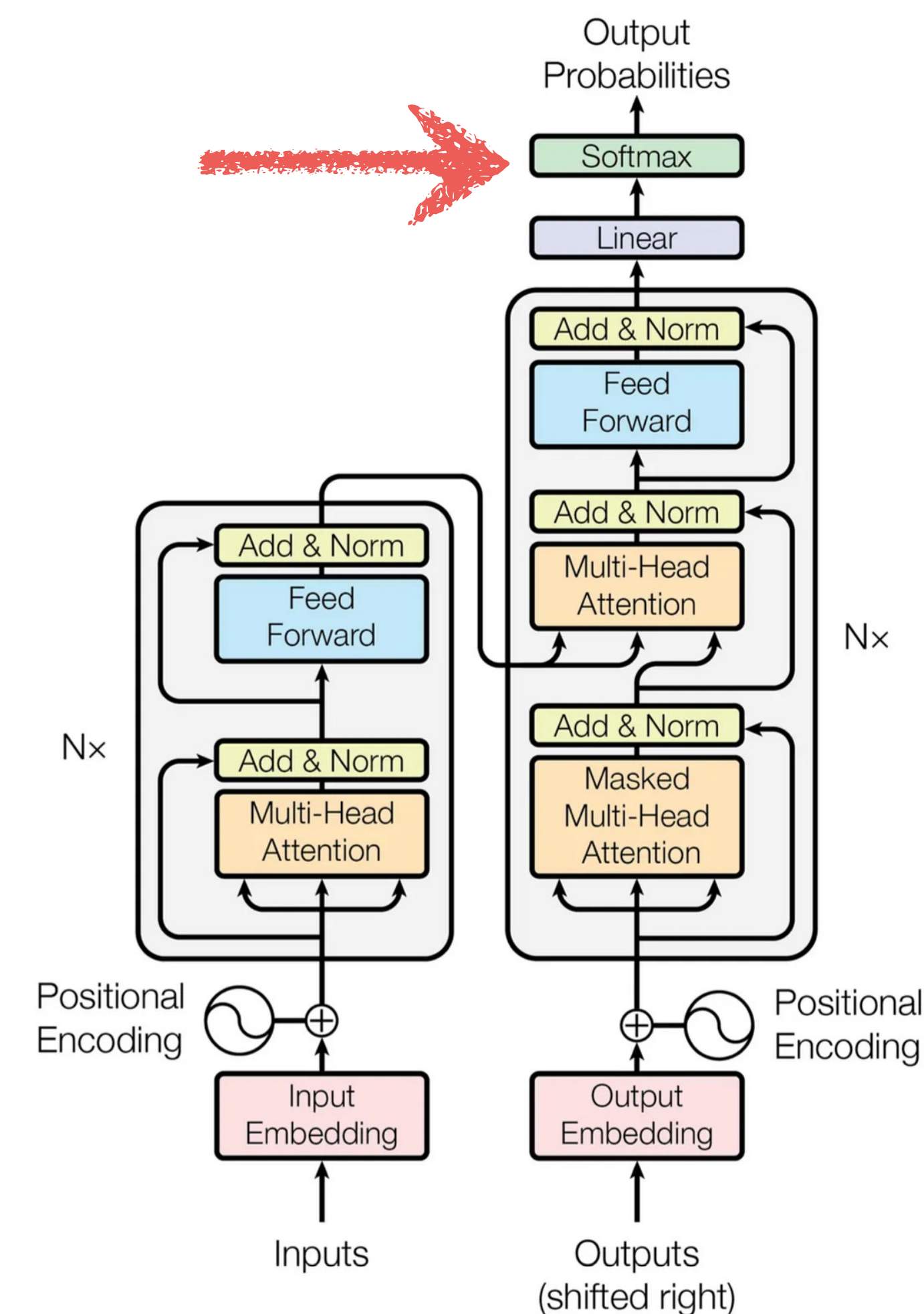


Figure 1: The Transformer - model architecture.

Decoder

- ▶ After repeating this decoding step 6 times, the 512dim output is put into a linear layer of dimension $10k \times 512$
- ▶ Then we take the softmax to obtain a $10k$ -dim vector of probabilities which word to use next
- ▶ In greedy decoding, the word with the highest probability is chosen from the $10k$ dictionary entries and added to the output

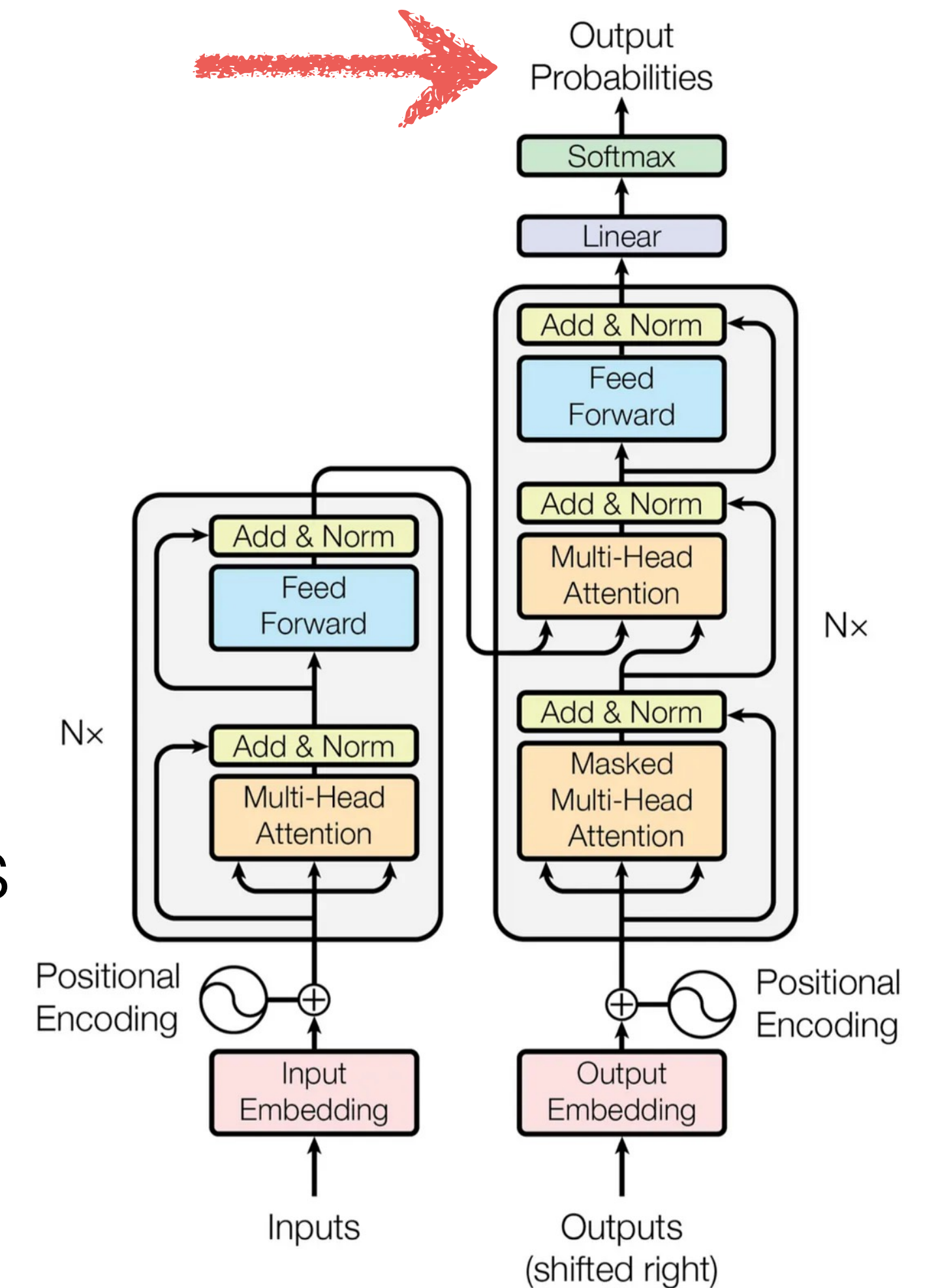


Figure 1: The Transformer - model architecture.

Decoder

- ▶ After repeating this decoding step 6 times, the 512dim output is put into a linear layer of dimension $10k \times 512$
- ▶ Then we take the softmax to obtain a $10k$ -dim vector of probabilities which word to use next
- ▶ In greedy decoding, the word with the highest probability is chosen from the $10k$ dictionary entries and added to the output
- ▶ Then the decoding step is repeated with the new output that contains this new word, and the next word to translate from the encoder

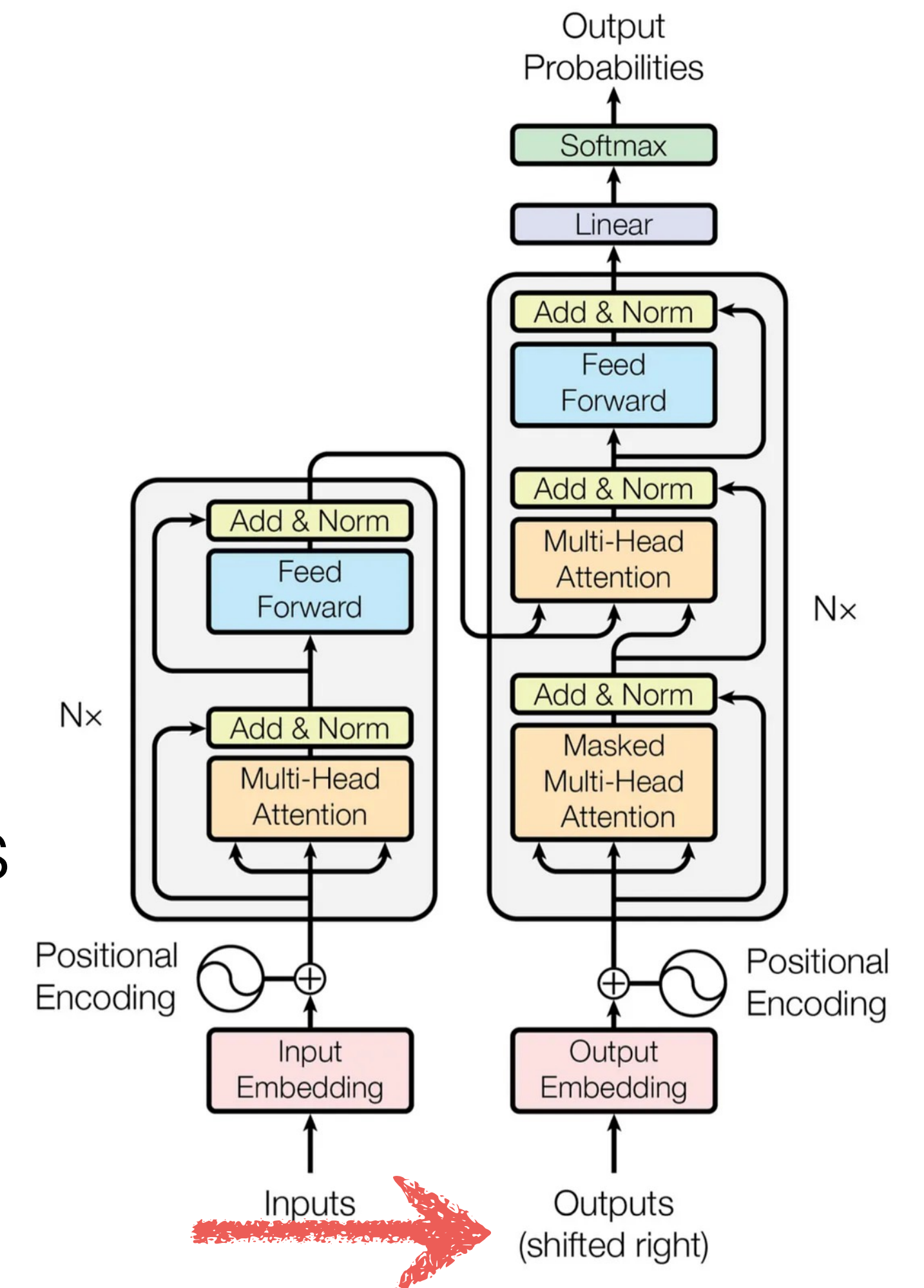
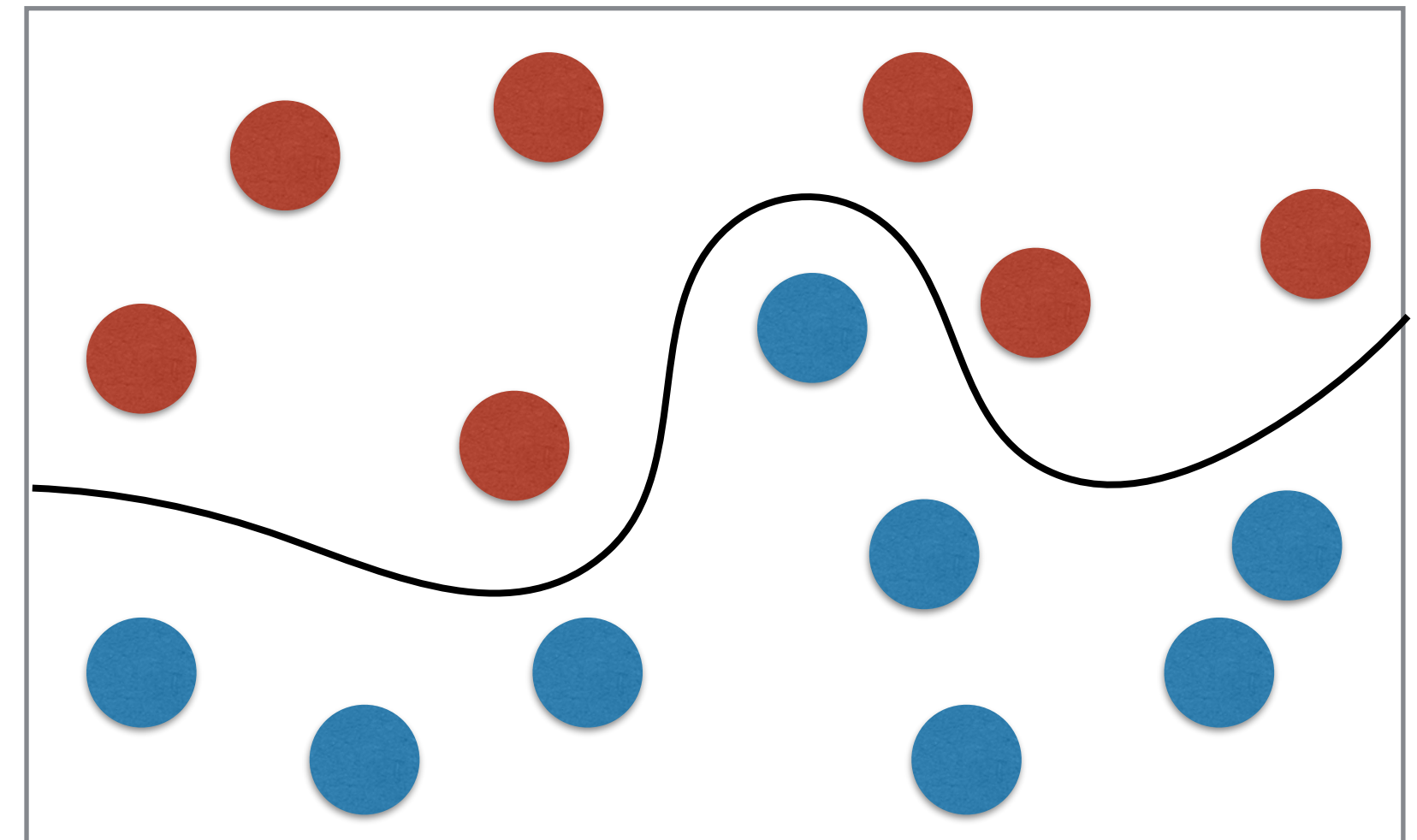
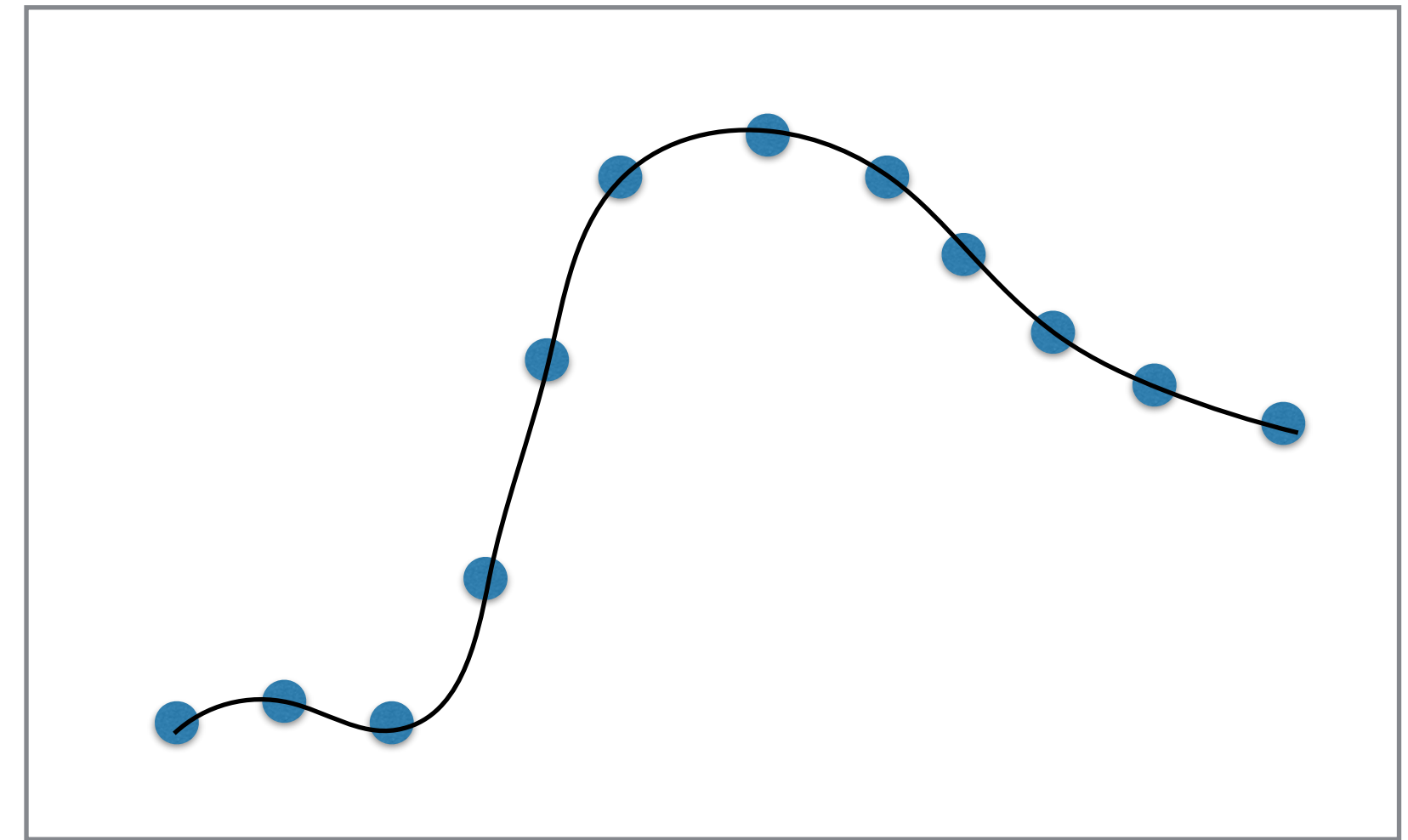
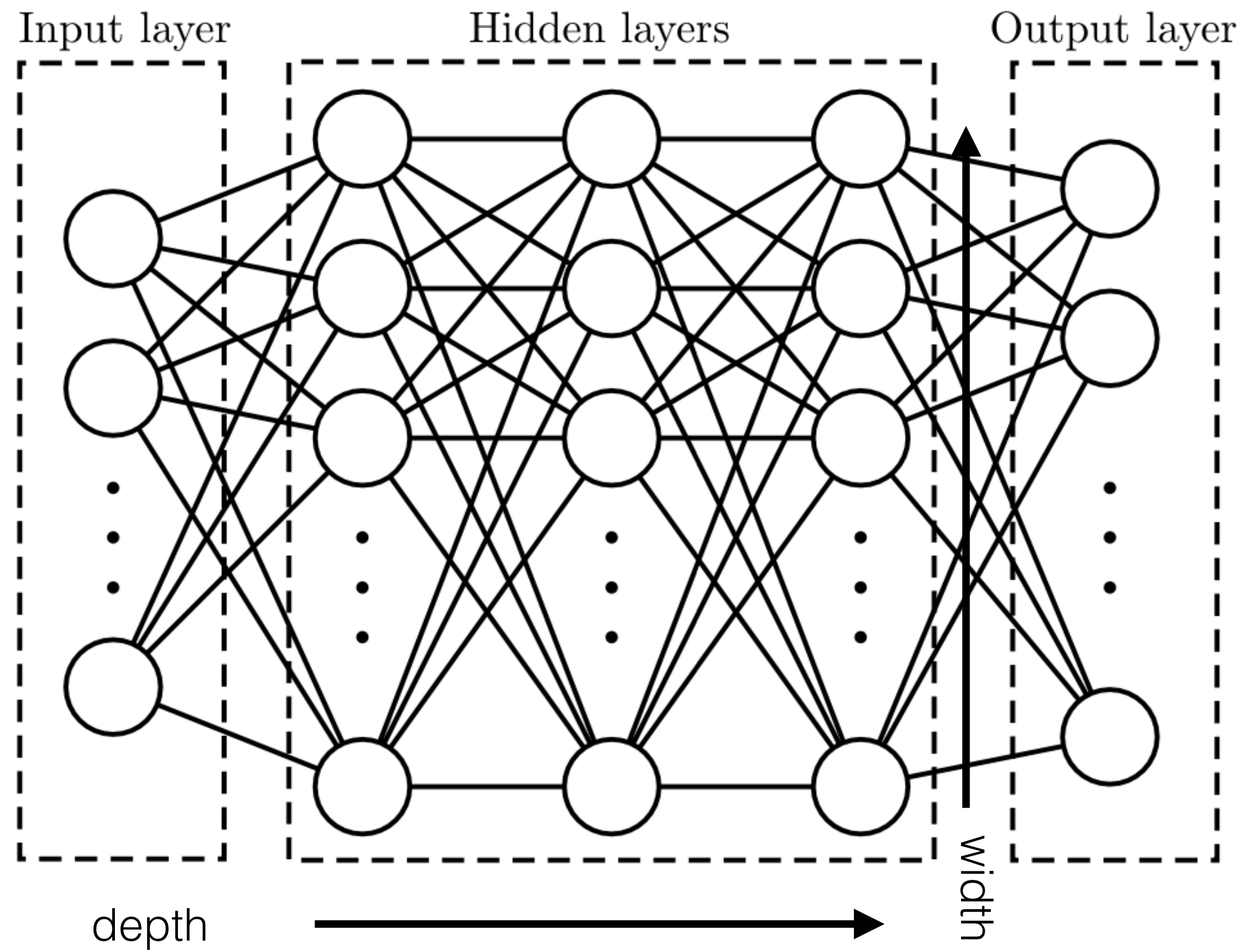
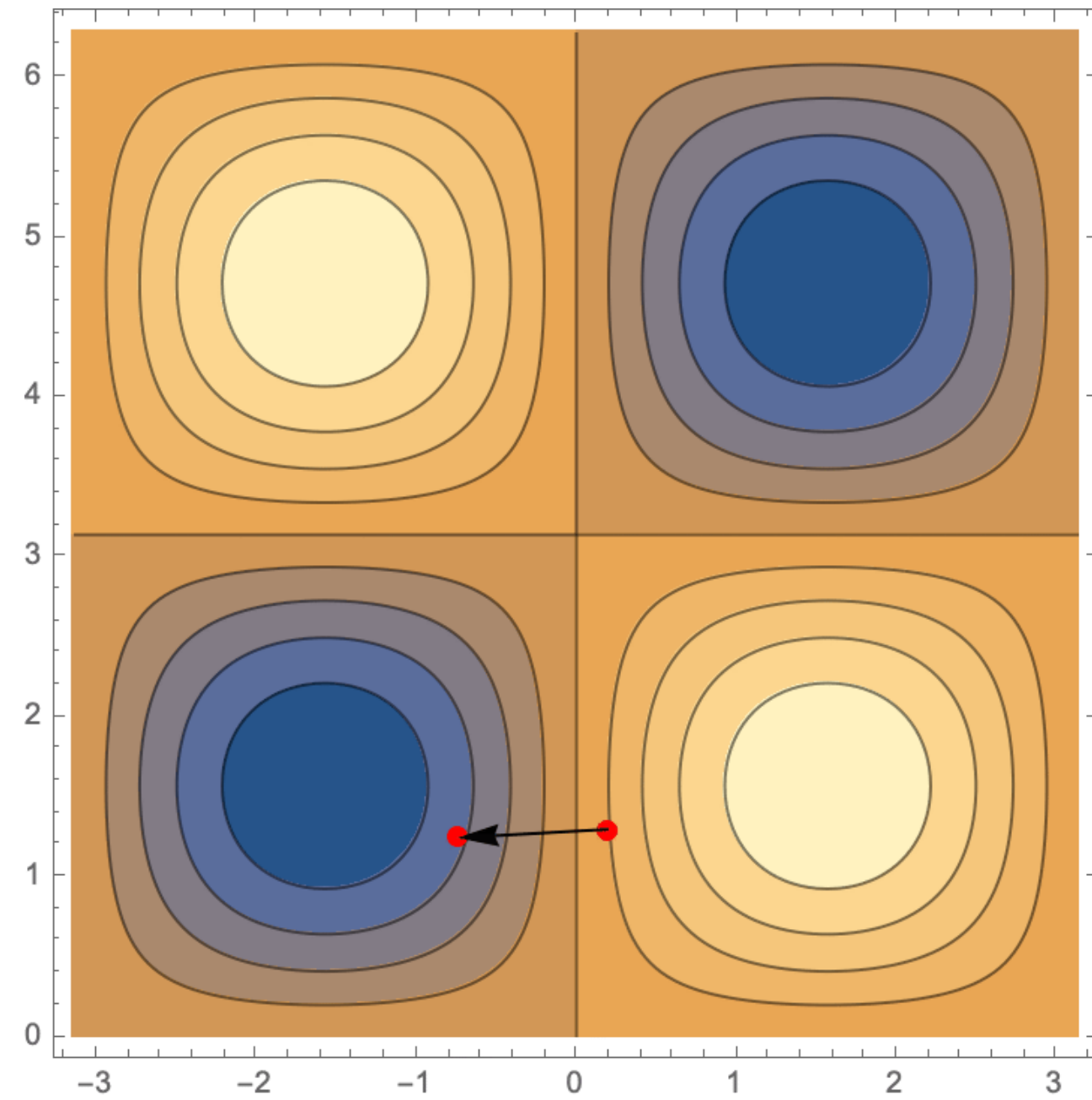


Figure 1: The Transformer - model architecture.

Recap - NN for regression and classification

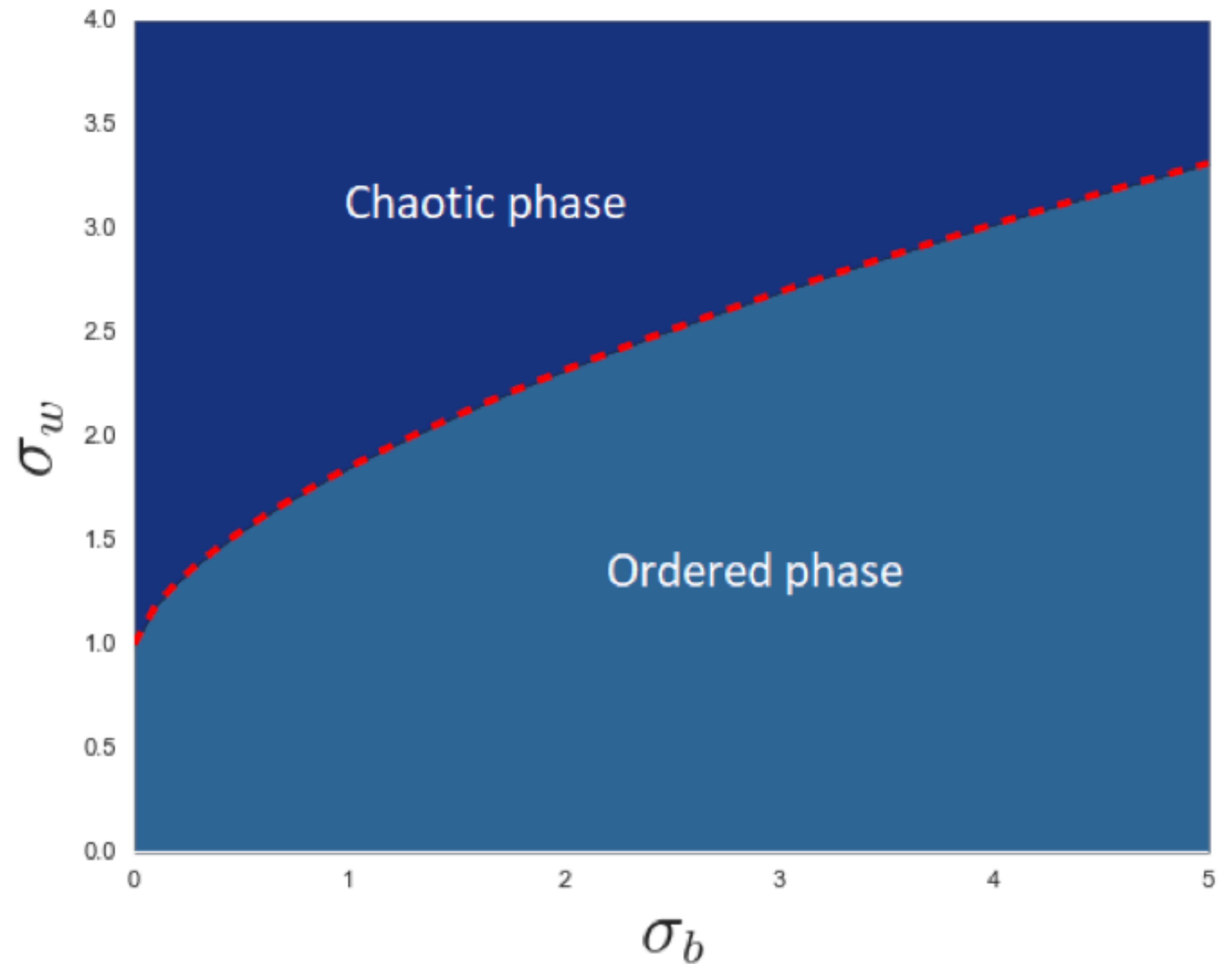


Recap - Training with GD

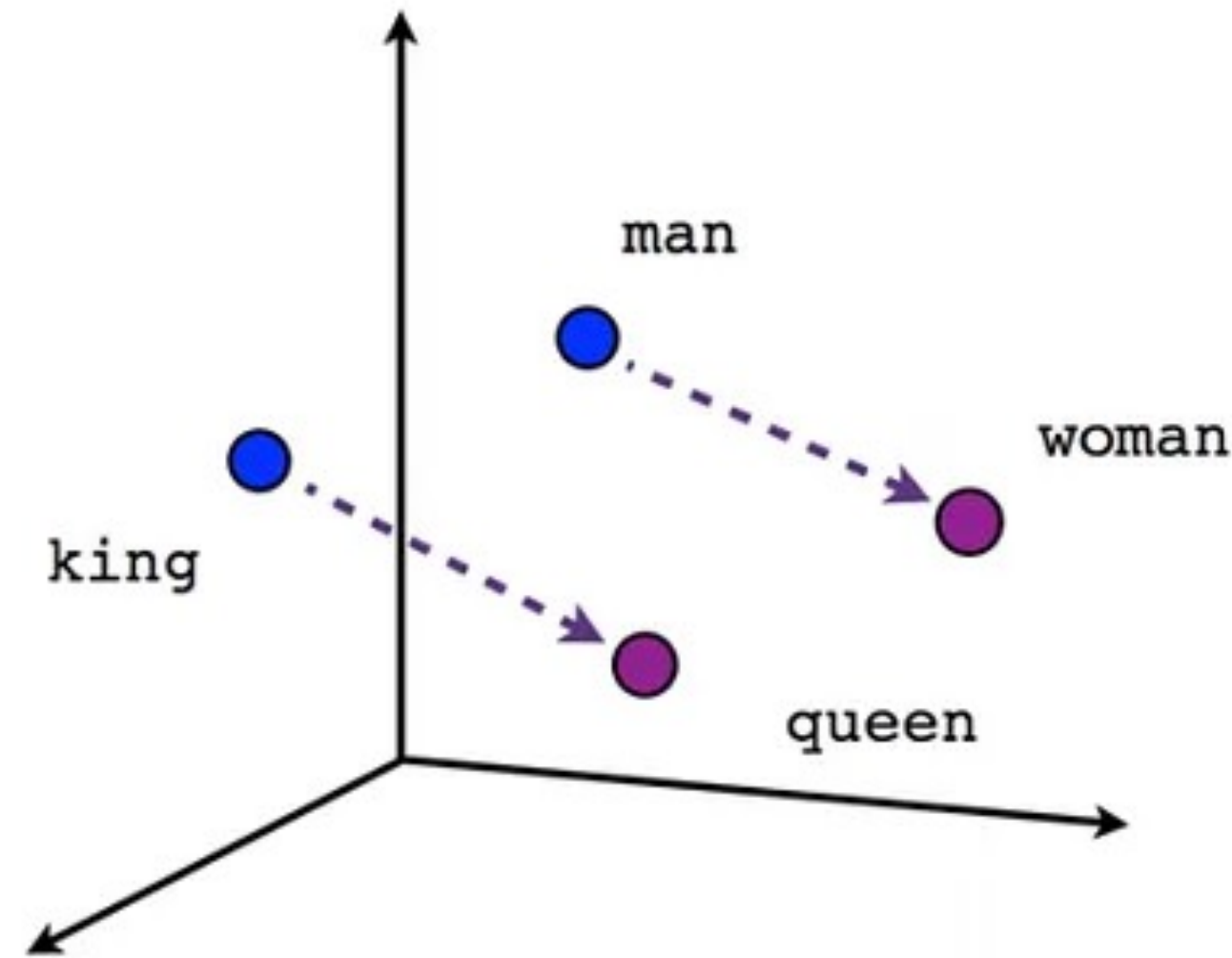


GD parameter update

$$\theta^{(i)} \rightarrow \theta^{(i)} - \alpha \frac{\partial L}{\partial \theta^{(i)}}$$



Recap - Transformer



Male-Female

Learns semantics

“queen - woman + man = king”

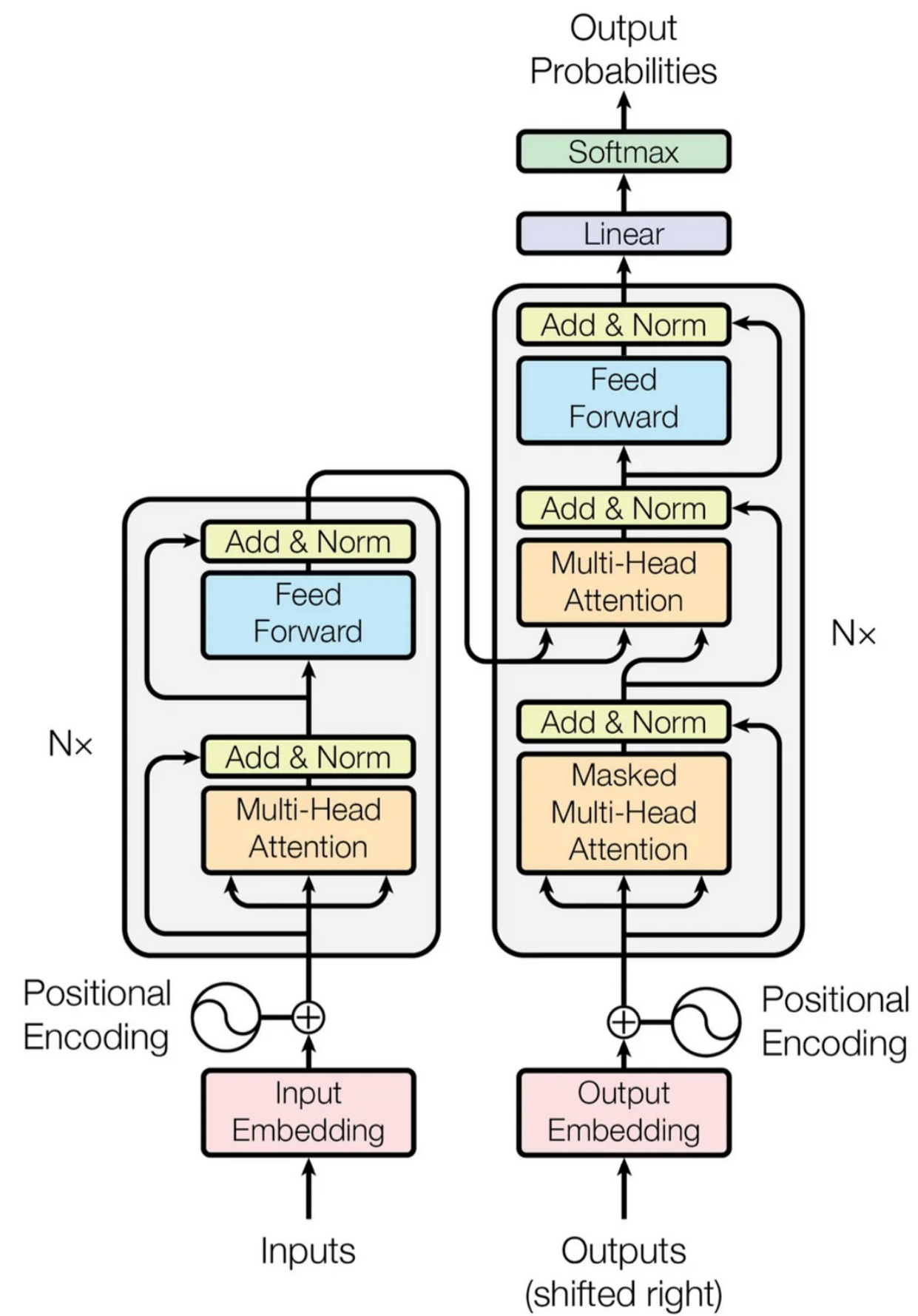
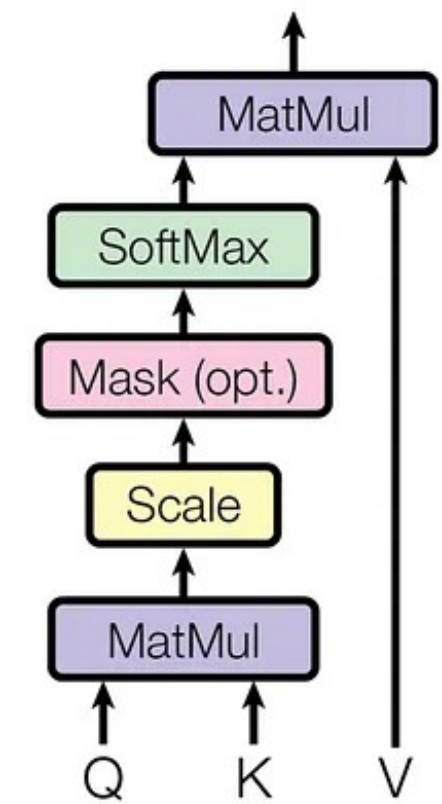
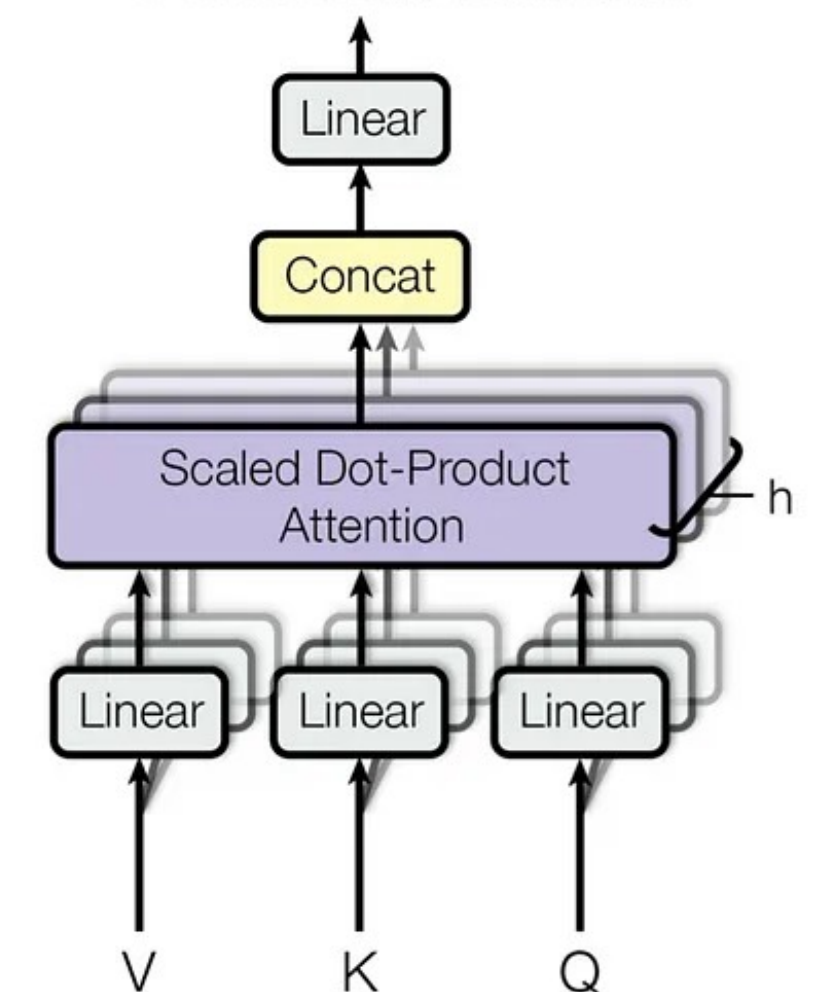


Figure 1: The Transformer - model architecture.

Scaled Dot-Product Attention



Multi-Head Attention

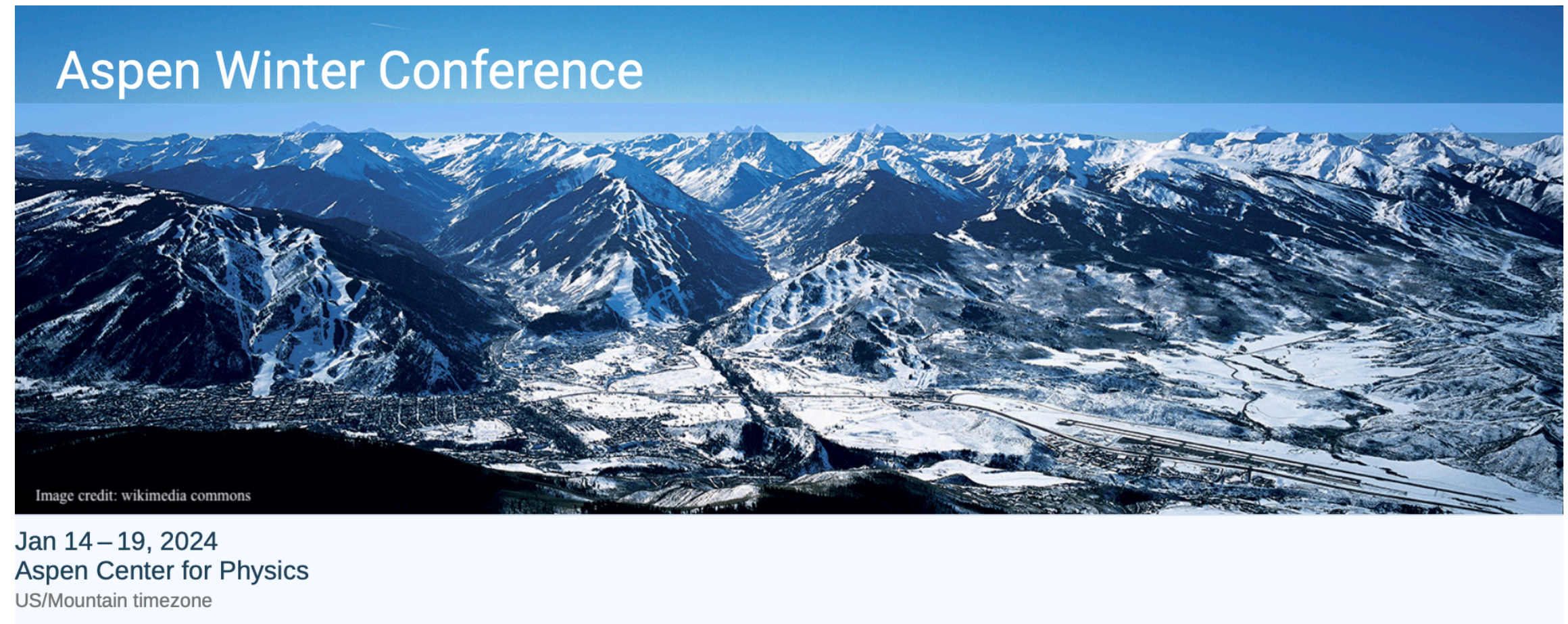


... want to learn more?



ML Meetings @ Caltech

- Dec 10-12: Mathematics and ML 2023
- Dec 13-15: string_data 2023



Overview
Timetable
Registration and Accommodation
Travel and Visas
Local information, Transportation, Skiing
Block Award
Poster

Fields, Strings, and Deep Learning

Progress in deep learning has traditionally involved experimental data, but in recent years it has impacted our understanding of formal structures arising in theoretical high energy physics and pure mathematics, via both theoretical and applied deep learning. This conference will bring together high energy theorists, mathematicians, and computer scientists across a broad variety of topics at the interface of these fields. Featured topics include the interface of neural network theory with quantum field theory, lattice field theory, conformal field theory, and the renormalization group; theoretical physics for AI, including equivariant, diffusion, and other generative models; ML for pure mathematics, including knot theory and special holonomy metrics, and deep learning for applications in string theory and holography.

Aspen Winter Conference

- Jan 14-19: Fields, Strings, and Deep Learning
- Application deadline: Aug 31**

Thank you - Questions?

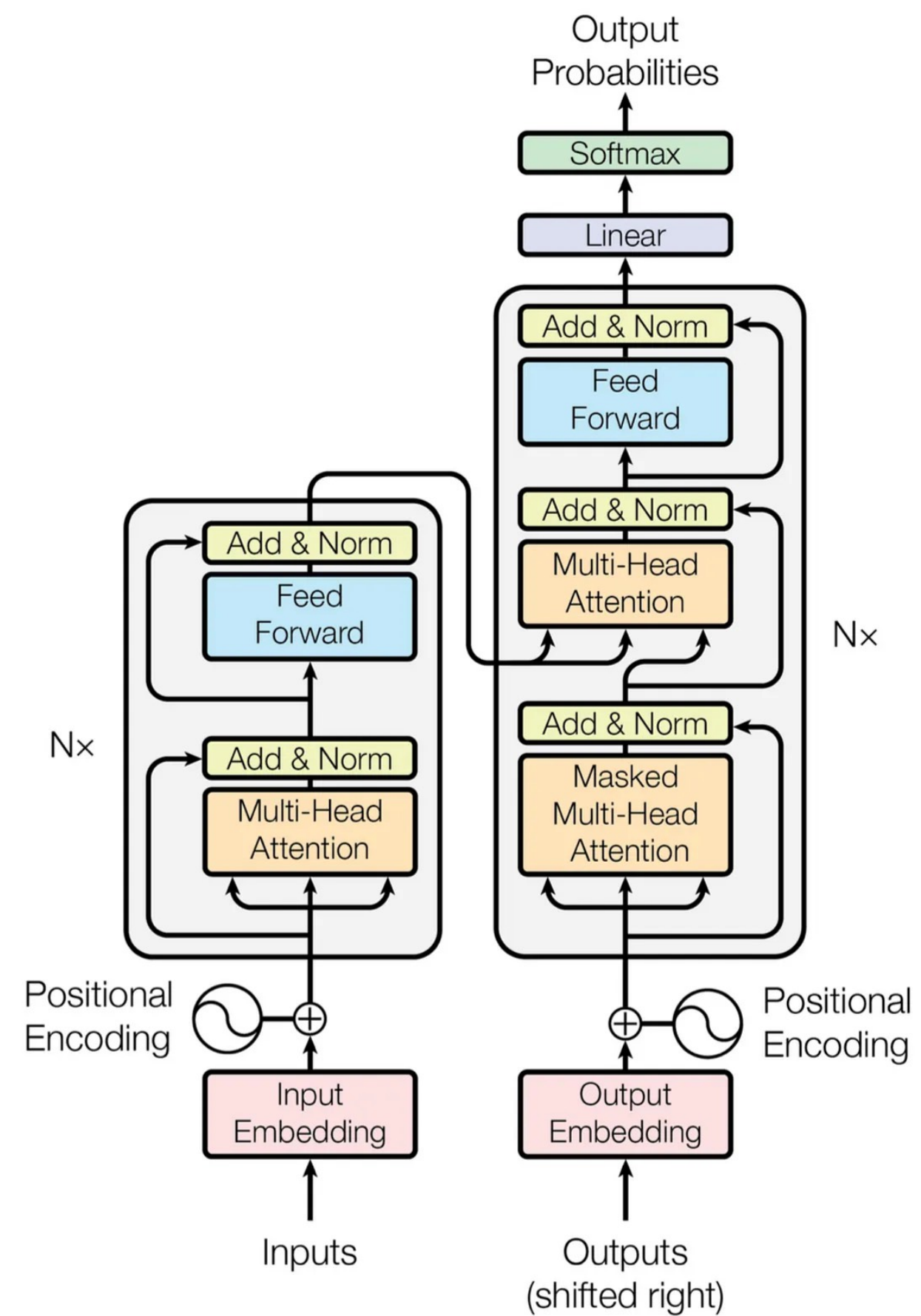
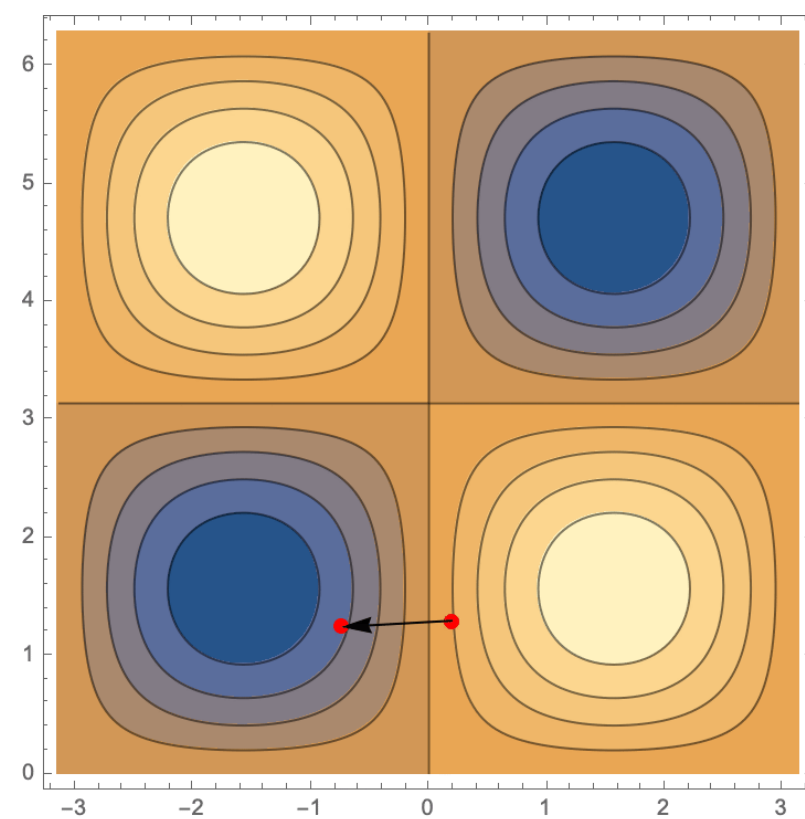
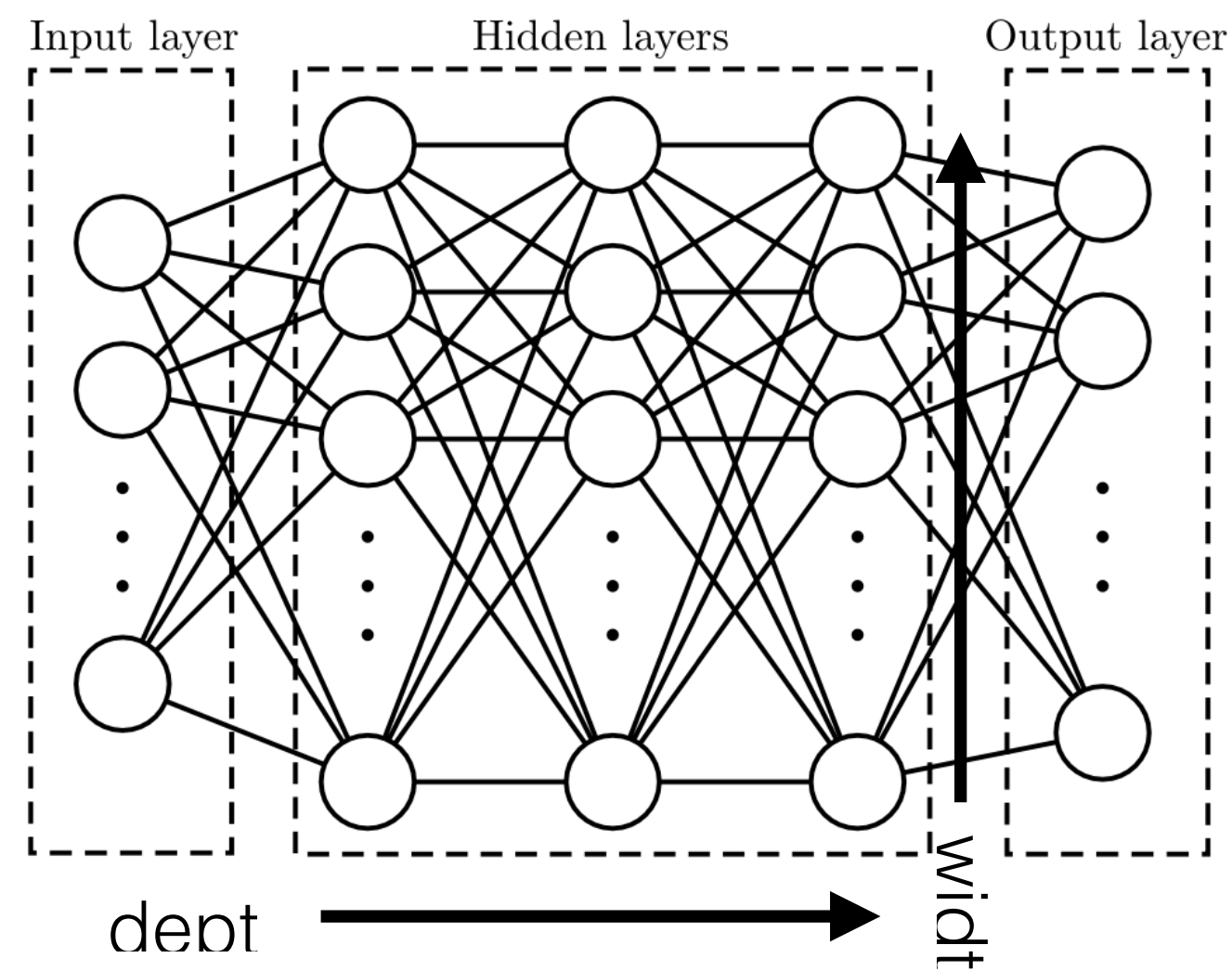
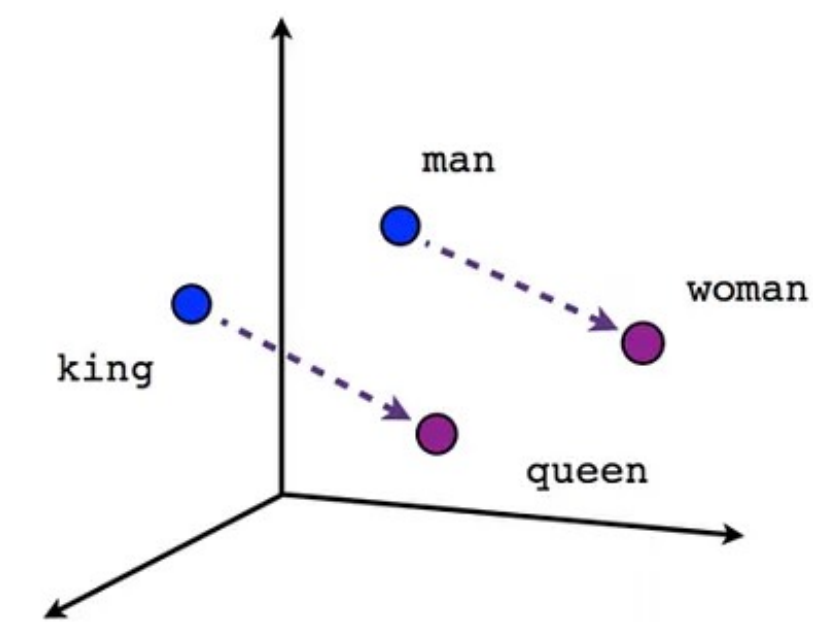
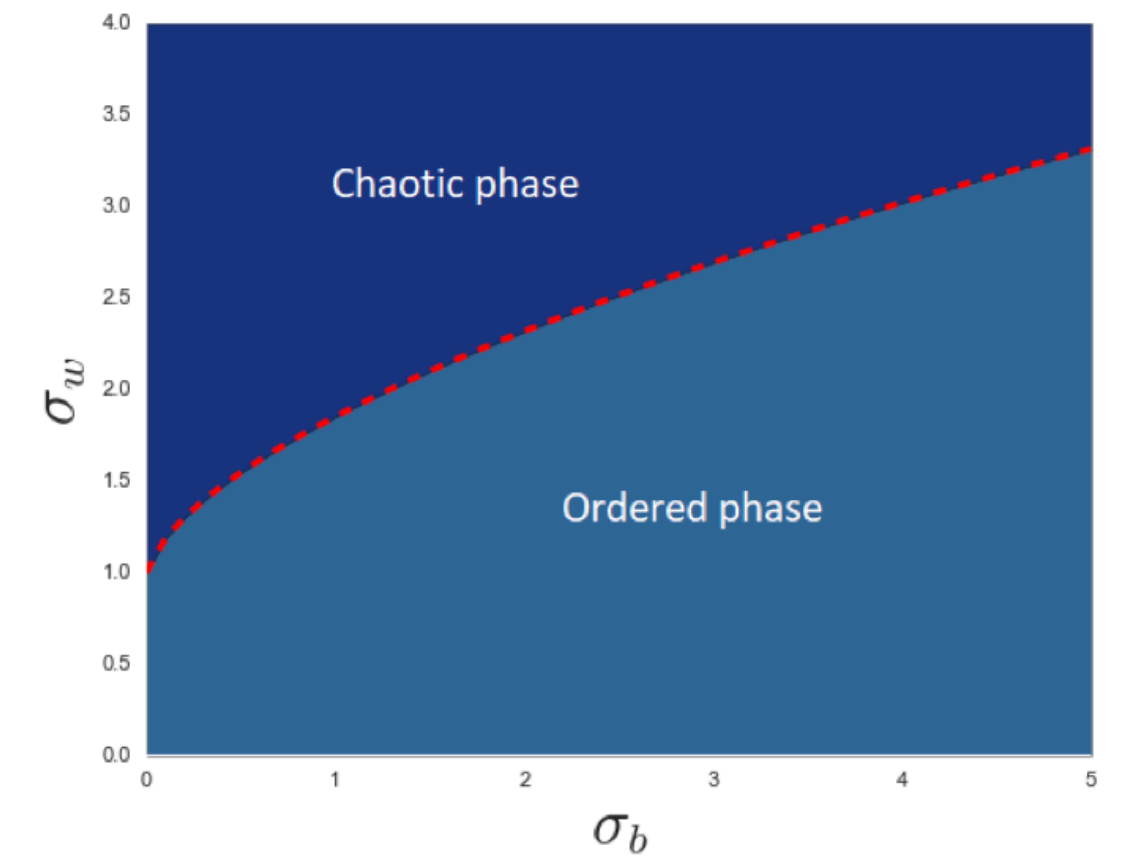


Figure 1: The Transformer - model architecture.



Male-Female