# STA 414/2104:
## Statistical Methods for Machine Learning II
### Week 12  Neural Networks

Michal Malyska

University of Toronto

# Today

- What are Neural Networks?

# Today

- What are Neural Networks?
- Neural Network Building Blocks

# Today

- What are Neural Networks?
- Neural Network Building Blocks
  - Linear (Feed Forward) Layers

# Today

- What are Neural Networks?
- Neural Network Building Blocks
  - Linear (Feed Forward) Layers
  - Activation Functions

# Today

- What are Neural Networks?
- Neural Network Building Blocks
  - ▶ Linear (Feed Forward) Layers
  - ▶ Activation Functions
  - ▶ Residual Layers

# Today

- What are Neural Networks?
- Neural Network Building Blocks
  - ▸ Linear (Feed Forward) Layers
  - ▸ Activation Functions
  - ▸ Residual Layers
  - ▸ Recurrent Layers

# Today

- What are Neural Networks?
- Neural Network Building Blocks
  - Linear (Feed Forward) Layers
  - Activation Functions
  - Residual Layers
  - Recurrent Layers
  - **Attention**

# Today

- What are Neural Networks?
- Neural Network Building Blocks
  - ▶ Linear (Feed Forward) Layers
  - ▶ Activation Functions
  - ▶ Residual Layers
  - ▶ Recurrent Layers
  - ▶ **Attention**
- Neural Networks

# Today

- What are Neural Networks?
- Neural Network Building Blocks
  - ▸ Linear (Feed Forward) Layers
  - ▸ Activation Functions
  - ▸ Residual Layers
  - ▸ Recurrent Layers
  - ▸ **Attention**
- Neural Networks
  - ▸ Feed Forward (Multi Layer Perceptron)

# Today

- What are Neural Networks?
- Neural Network Building Blocks
  - Linear (Feed Forward) Layers
  - Activation Functions
  - Residual Layers
  - Recurrent Layers
  - **Attention**
- Neural Networks
  - Feed Forward (Multi Layer Perceptron)
  - Recurrent

# Today

- What are Neural Networks?
- Neural Network Building Blocks
  - Linear (Feed Forward) Layers
  - Activation Functions
  - Residual Layers
  - Recurrent Layers
  - **Attention**
- Neural Networks
  - Feed Forward (Multi Layer Perceptron)
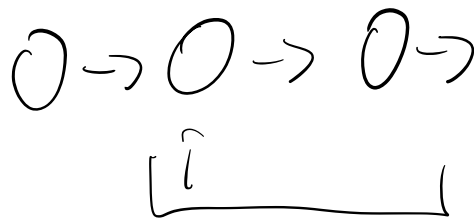  - Recurrent
  - **Transformer**

# Today

- What are Neural Networks?
- Neural Network Building Blocks
  - Linear (Feed Forward) Layers
  - Activation Functions
  - Residual Layers
  - Recurrent Layers
  - **Attention**
- Neural Networks
  - Feed Forward (Multi Layer Perceptron)
  - Recurrent
  - **Transformer**
- Transformer

# Today

- What are Neural Networks?
- Neural Network Building Blocks
  - Linear (Feed Forward) Layers
  - Activation Functions
  - Residual Layers
  - Recurrent Layers
  - **Attention**
- Neural Networks
  - Feed Forward (Multi Layer Perceptron)
  - Recurrent
  - **Transformer**
- Transformer
  - Encoder

# Today

- What are Neural Networks?
- Neural Network Building Blocks
  - ▶ Linear (Feed Forward) Layers
  - ▶ Activation Functions
  - ▶ Residual Layers
  - ▶ Recurrent Layers
  - ▶ **Attention**
- Neural Networks
  - ▶ Feed Forward (Multi Layer Perceptron)
  - ▶ Recurrent
  - ▶ **Transformer**
- Transformer
  - ▶ Encoder
  - ▶ Decoder

# Today

- What are Neural Networks?
- Neural Network Building Blocks
  - Linear (Feed Forward) Layers
  - Activation Functions
  - Residual Layers
  - Recurrent Layers
  - **Attention**
- Neural Networks
  - Feed Forward (Multi Layer Perceptron)
  - Recurrent
  - **Transformer**
- Transformer
  - Encoder
  - Decoder
  - Positional Encoding

# What are Neural Networks

**Neural networks** are what we commonly call any differentiable function that can be expressed as a computation graph. Each node is a primitive operation (e.g. matrix multiplication) and edges represent data flow. In particular, a simple (and quite common) case is where this graph is a chain. Individual nodes, or pre-defined sequences are often referred to as **layers**

# Building Blocks of Neural Networks

**Linear (Feed Forward) Layers** - is the simplest possible type of layer, it consists of 2 operations:

- Matrix multiplication

# Building Blocks of Neural Networks

**Linear (Feed Forward) Layers** - is the simplest possible type of layer, it consists of 2 operations:

- Matrix multiplication
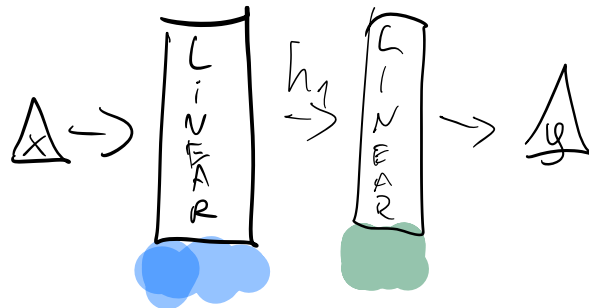- Vector addition
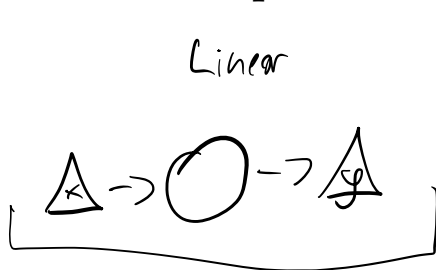
# Building Blocks of Neural Networks

**Linear (Feed Forward) Layers** - is the simplest possible type of layer, it consists of 2 operations:

- Matrix multiplication
- Vector addition

# Building Blocks of Neural Networks

**Linear (Feed Forward) Layers** - is the simplest possible type of layer, it consists of 2 operations:

- Matrix multiplication
- Vector addition

$$y = f(x; \theta) = Wx + b$$

where $\theta$ is the set of parameters $\{W, b\}$

# Building Blocks of Neural Networks

- What would happen if we followed up a Linear Layer by another linear layer?

$$y = f(g(x; \theta_1); \theta_2) = W_2(W_1 x + b_1) + b_2 =$$
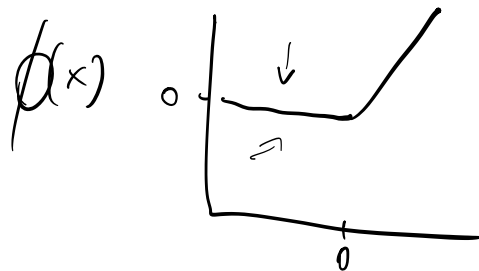$$= \underbrace{(W_2 W_1)}x + \underbrace{(W_2 b_1 + b_2)}$$

# Building Blocks of Neural Networks

- What would happen if we followed up a Linear Layer by another linear layer?

$$y = f(g(x; \theta_1); \theta_2) = W_2(W_1 x + b_1) + b_2 =$$
$$= (W_2 W_1)x + (W_2 b_1 + b_2)$$

# Building Blocks of Neural Networks

- What would happen if we followed up a Linear Layer by another linear layer?

$$y = f(g(x; \theta_1); \theta_2) = W_2(W_1 x + b_1) + b_2 =$$
$$= (W_2 W_1)x + (W_2 b_1 + b_2)$$

Ok, not very useful. Is there anything we can do about it?

# Building Blocks of Neural Networks

- What would happen if we followed up a Linear Layer by another linear layer?

$$y = f(g(x; \theta_1); \theta_2) = W_2(W_1 x + b_1) + b_2 =$$
$$= (W_2 W_1)x + (W_2 b_1 + b_2)$$

Ok, not very useful. Is there anything we can do about it? **Yes**, to get more expressive power, we can apply a non-linear (element-wise) transformation. We call these functions **Activation functions**. Some common examples include:
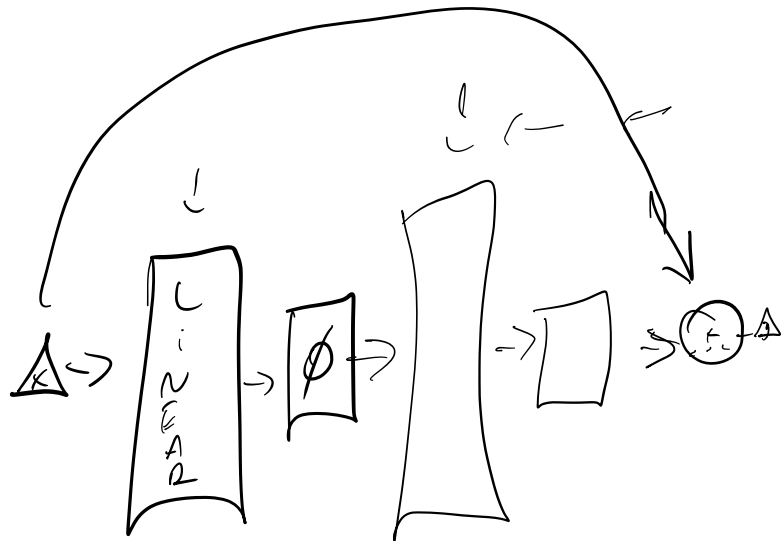
- **Re**ctified **L**inear **U**nit: $\phi(x) = max(0, x)$ ReLU

$$\phi(x)$$

$$y = f\left(\phi\left(g(x; \theta_1)\right); \theta_2\right)$$

# Building Blocks of Neural Networks

- What would happen if we followed up a Linear Layer by another linear layer?

$$y = f(g(x; \theta_1); \theta_2) = W_2(W_1 x + b_1) + b_2 =$$
$$= (W_2 W_1)x + (W_2 b_1 + b_2)$$

Ok, not very useful. Is there anything we can do about it? **Yes**, to get more expressive power, we can apply a non-linear (element-wise) transformation. We call these functions **Activation functions**. Some common examples include:
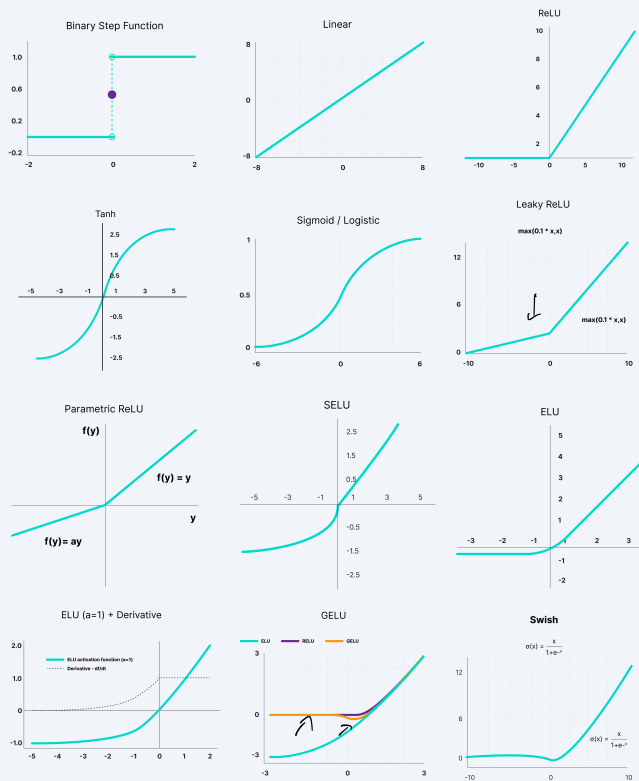
- **Re**ctified **L**inear **U**nit: $\phi(x) = max(0, x)$
- **Sigmoid**: $\phi(x) = \sigma(x) = \frac{1}{1+e^{-x}}$

# Activations Examples



**Neural Network Activation Functions**

Binary Step Function · Linear · ReLU · Tanh · Sigmoid / Logistic · Leaky ReLU · Parametric ReLU · SELU · ELU · ELU (a=1) + Derivative · GELU · Swish

v7labs.com

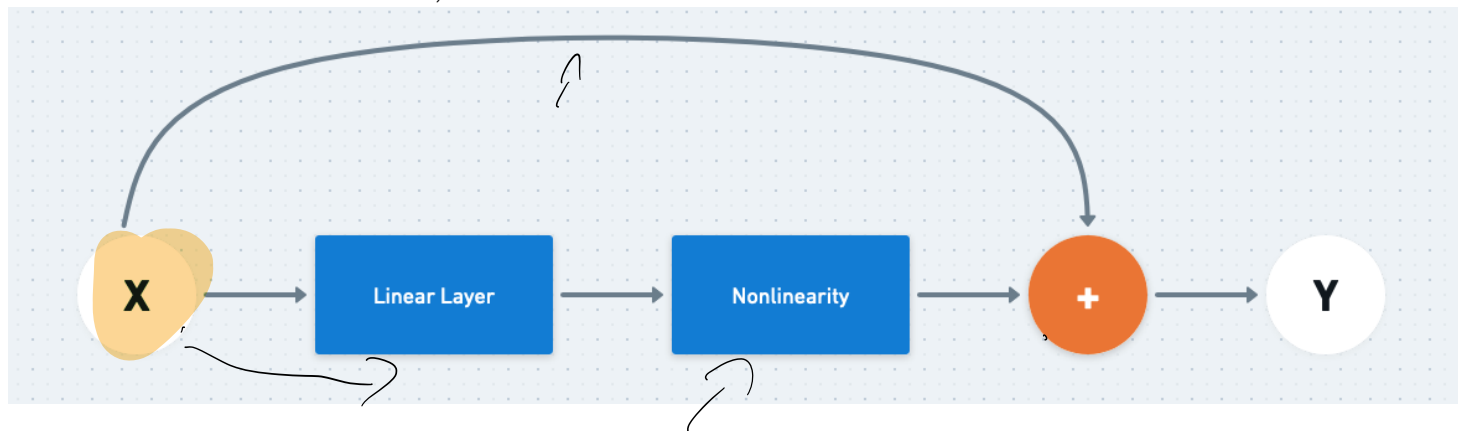# Building Blocks of Neural Networks

If we begin stacking large number of layers together, the signal may get squashed to zero, or blow up to infinity. Similar problem often happens during the gradient computation back through the graph.

# Building Blocks of Neural Networks

If we begin stacking large number of layers together, the signal may get squashed to zero, or blow up to infinity. Similar problem often happens during the gradient computation back through the graph. To reduce the effect of those problems we often propagate the signal to layers further downstream, in what are called **residual connections**

# Building Blocks of Neural Networks

If we begin stacking large number of layers together, the signal may get squashed to zero, or blow up to infinity. Similar problem often happens during the gradient computation back through the graph. To reduce the effect of those problems we often propagate the signal to layers further downstream, in what are called **residual connections**



$$y = f(x; \theta) = \phi(Wx + b) + x$$

# Building Blocks of Neural Networks

When it comes to modelling sequences (e.g. text, or time series data), it is often useful to make the model stateful in order for it to help "carry" the information through the graph. To do that we simply add a state at timepoint $t$: $s_t$, and computing the output and the new state using some function:

# Building Blocks of Neural Networks

When it comes to modelling sequences (e.g. text, or time series data), it is often useful to make the model stateful in order for it to help "carry" the information through the graph. To do that we simply add a state at timepoint $t$: $s_t$, and computing the output and the new state using some function:

$$(y, s_{t+1}) = f(x, s_t)$$

# Building Blocks of Neural Networks

When it comes to modelling sequences (e.g. text, or time series data), it is often useful to make the model stateful in order for it to help "carry" the information through the graph. To do that we simply add a state at timepoint $t$: $s_t$, and computing the output and the new state using some function:

$$(y, s_{t+1}) = f(x, s_t)$$

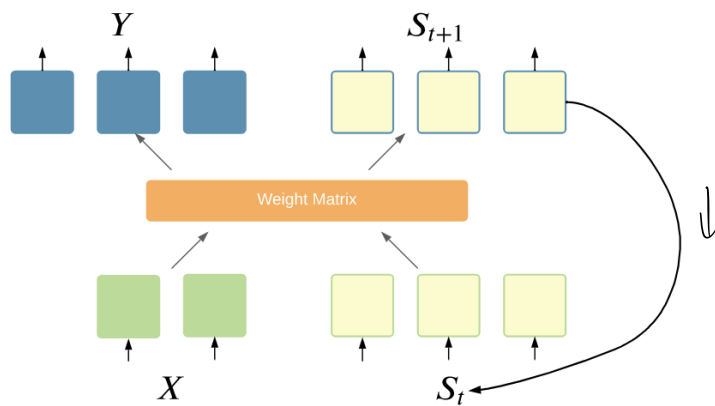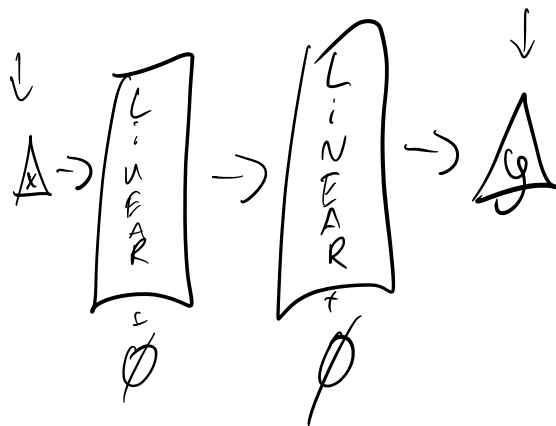This is then called a **recurrent layer**.

# Building Blocks of Neural Networks

When it comes to modelling sequences (e.g. text, or time series data), it is often useful to make the model stateful in order for it to help "carry" the information through the graph. To do that we simply add a state at timepoint $t$: $s_t$, and computing the output and the new state using some function:

$$(y, s_{t+1}) = f(x, s_t)$$

This is then called a **recurrent layer**.



*Figure 16.8: Recurrent layer.*

# Common Architectures

FFNN

A very common type of neural net architecture is a **Feed Forward Neural Network**, also sometimes called a **Multi Layer** MLP
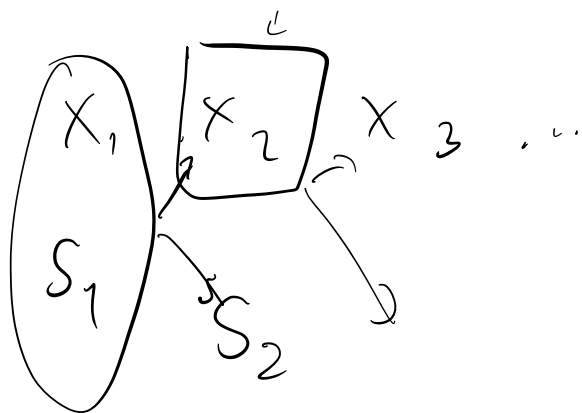**Perceptron**. It simply consists of a sequence of linear (FF) layers, with nonlinearities between them.



$$(y - \hat{y})^2$$

# Common Architectures

A very common type of neural net architecture is a **Feed Forward Neural Network**, also sometimes called a **Multi Layer Perceptron**. It simply consists of a sequence of linear (FF) layers, with nonlinearities between them.

$$f(x; \theta) = \phi(W_L(\phi(W_{L-1}(\phi(W_{L-2}(\dots) + b_{L-2})) + b_{L-1})) + b_L)$$

# Common Architectures

If we use recurrent layers in our neural network, the outcome is what we typically call a **Recurrent Neural Network**, (of which there are many variants). In the simplest possible option the function $f(x, h)$ is a simple FFNN.

# Common Architectures

If we use recurrent layers in our neural network, the outcome is what we typically call a **Recurrent Neural Network**, (of which there are many variants). In the simplest possible option the function $f(x, h)$ is a simple FFNN. When training RNNs each item in a sequence is used as input, however during inference each item in the sequence will depend on previous predictions.
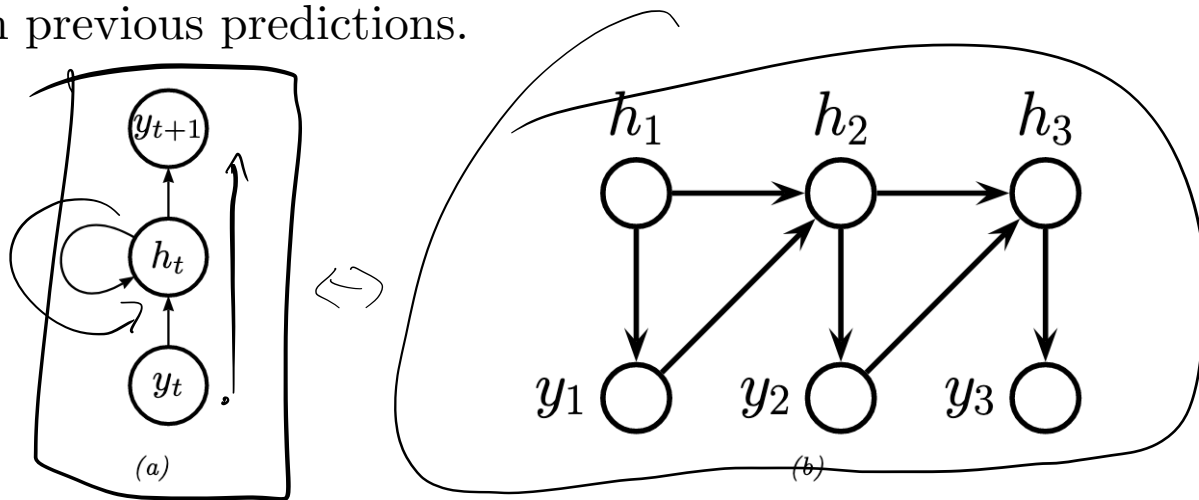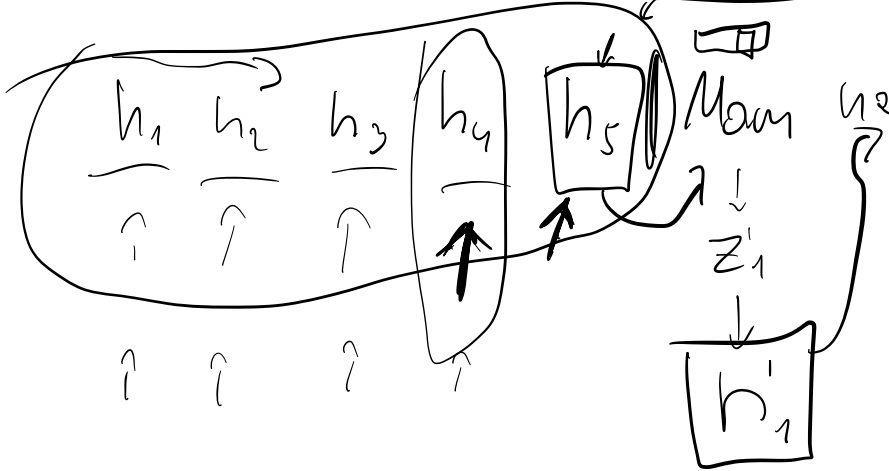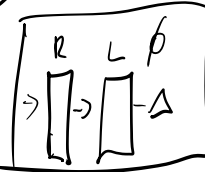
# Common Architectures

If we use recurrent layers in our neural network, the outcome is what we typically call a **Recurrent Neural Network**, (of which there are many variants). In the simplest possible option the function $f(x, h)$ is a simple FFNN. When training RNNs each item in a sequence is used as input, however during inference each item in the sequence will depend on previous predictions.



Figure 16.12: Illustration of a recurrent neural network (RNN). (a) With self-loop. (b) Unrolled in time.

$x_1 \quad x_2 \quad x_3 \quad x_4 \quad \langle EOS \rangle \quad y_1 \quad y_2 \quad y_3 \quad y_4 \ldots \langle EOS \rangle$

$E \quad$ My mare is $\underline{Michal} \langle EOS \rangle \qquad D$

$\quad \downarrow \quad \downarrow \quad \downarrow \quad \downarrow \qquad \downarrow$

$\quad z_1 \quad z_2 \quad z_3 \, z_4 \qquad z_5$



$h_1 \quad h_2 \quad h_3 \quad h_4 \quad h_5 \qquad$ Mam $\qquad$ he $\text{for}\xi \quad \underline{Midal} \langle EOS \rangle$

$\qquad \qquad \qquad \qquad \qquad \qquad \qquad \downarrow \qquad \qquad \qquad \qquad \uparrow$

$\qquad \qquad \qquad \qquad \qquad \qquad \qquad z_1'$

$\qquad \qquad \qquad \qquad \qquad \qquad \qquad \downarrow$

$\qquad \qquad \qquad \qquad \qquad \qquad \boxed{h_1'}$

$RNN \rightarrow \quad h_5 \quad \rightarrow \quad RNN \rightarrow Linear$

# Attention is all you need

What if instead of getting just the previous hidden state we were able to take a look at a lot of the previous inputs at once?

# Attention is all you need

What if instead of getting just the previous hidden state we were able to take a look at a lot of the previous inputs at once? We could combine all the previous hidden states. But can we do better?

# Attention is all you need

What if instead of getting just the previous hidden state we were able to take a look at a lot of the previous inputs at once? We could combine all the previous hidden states. But can we do better? We can score each of the hidden states by how well it is associated with the state we will be predicting.

# Attention is all you need

What if instead of getting just the previous hidden state we were able to take a look at a lot of the previous inputs at once? We could combine all the previous hidden states. But can we do better? We can score each of the hidden states by how well it is associated with the state we will be predicting. At a high level the attention mechanism consists of 3 simple steps:

1. Generate a score for each of the hidden states

# Attention is all you need

What if instead of getting just the previous hidden state we were able to take a look at a lot of the previous inputs at once? We could combine all the previous hidden states. But can we do better? We can score each of the hidden states by how well it is associated with the state we will be predicting. At a high level the attention mechanism consists of 3 simple steps:

1. Generate a score for each of the hidden states
2. Apply the softmax function to the scores

$$x_1 \quad x_2 \quad x_3 \quad x_4$$

$$1 \quad 1 \quad 1 \quad 5$$

$$\Rightarrow 0.05 \quad 0.05 \quad 0.05 \quad 0.85$$

# Attention is all you need

What if instead of getting just the previous hidden state we were able to take a look at a lot of the previous inputs at once? We could combine all the previous hidden states. But can we do better? We can score each of the hidden states by how well it is associated with the state we will be predicting. At a high level the attention mechanism consists of 3 simple steps:

1. Generate a score for each of the hidden states
2. Apply the softmax function to the scores
3. Multiply each of the hidden states by the output of the softmax and add them together.

# Attention is all you need

Now the hard problem remains: how do we score each of the hidden states?

# Attention is all you need

Now the hard problem remains: how do we score each of the hidden states? We will begin by creating 3 separate embeddings from each of our inputs, by simply multipling them by (learned) matrices:

$$q = W^Q x$$
$$k = W^K x$$
$$v = W^V x$$

# Attention is all you need

Now the hard problem remains: how do we score each of the hidden states? We will begin by creating 3 separate embeddings from each of our inputs, by simply multipling them by (learned) matrices:

$$q = W^Q x$$
$$k = W^K x$$
$$v = W^V x$$

We then define the **Attention Layer** as:

$$Attn(q, k, v) = \sum_{i=1}^{m} \alpha_i(q, k_i) v_i$$

Where $\alpha$ is the scoring function.

$X_1 \quad X_2 \quad X_3$

$Attn(X_1)$

$\alpha(q_2 \quad k_1) \quad v_1$

$\alpha(q_2 \quad k_2) v_2$

$\alpha(q_2 \quad k_3) v_3$

# (Dot product) Attention is all you need

$$Attn(q, k, v) = \sum_{i=1}^{m} \alpha_i(q, k_i) v_i$$

# (Dot product) Attention is all you need

$$Attn(q, k, v) = \sum_{i=1}^{m} \overset{\downarrow}{\alpha_i}(q, k_i)v_i$$

The most common choice of the attention function is called the **dot product attention**. We obtain the scores by a normalized dot product of the $k$ and $q$ vectors.

# (Dot product) Attention is all you need

$$Attn(q, k, v) = \sum_{i=1}^{m} \alpha_i(q, k_i) v_i$$

The most common choice of the attention function is called the **dot product attention**. We obtain the scores by a normalized dot product of the $k$ and $q$ vectors.

$$b(q, k) = \frac{q^T k}{\sqrt{d}}$$

# (Dot product) Attention is all you need

$$Attn(q, k, v) = \sum_{i=1}^{m} \alpha_i(q, k_i) v_i$$

The most common choice of the attention function is called the **dot product attention**. We obtain the scores by a normalized dot product of the $k$ and $q$ vectors.

$$b(q, k) = \frac{q^T k}{\sqrt{d}}$$

where $d$ is a normalizing constant, usually the dimensionality of the vectors. We then set our attention weights $\alpha_i$ to be the softmax of all the scores:

# (Dot product) Attention is all you need

$$Attn(q, k, v) = \sum_{i=1}^{m} \alpha_i(q, k_i) v_i$$

The most common choice of the attention function is called the **dot product attention**. We obtain the scores by a normalized dot product of the $k$ and $q$ vectors.

$$b(q, k) = \frac{q^T k}{\sqrt{d}}$$

where $d$ is a normalizing constant, usually the dimensionality of the vectors. We then set our attention weights $\alpha_i$ to be the softmax of all the scores:

$$\alpha_i(q, k_i) = \frac{exp(b(q, k_i))}{\sum_{j=1}^{m} exp(b(q, k_j))}$$

# (Dot product) Attention is all you need

The entire process then reduces to:

$$\text{Attn}\left(q, k, v\right) = \sum_{i=1}^{m} \alpha_i\left(q, k_i\right) v_i$$

# (Dot product) Attention is all you need

$W^Q X$

The entire process then reduces to:

$$Y = Attn(Q, K, V) = \sigma(\frac{QK^T}{\sqrt{d}})V$$

$$= \sigma(\frac{W^Q X (W^K X)^T}{\sqrt{d}}) W^V X$$

# (Dot product) Attention is all you need

The entire process then reduces to:

$$Y = Attn(Q, K, V) = \sigma\left(\frac{QK^T}{\sqrt{d}}\right)V$$

$$= \sigma\left(\frac{W^Q X (W^K X)^T}{\sqrt{d}}\right)W^V X$$

**Scaled Dot-Product Attention**



$$x_1 \quad x_2 \quad x_3 \quad x_4$$

$$\langle EOS \rangle$$

$$\left\{ h_1 \quad h_2 \quad h_3 \quad h_4 \quad h_5 \right\}$$

$$y_1 \quad y_2$$

$$LINEAR$$

$$q_5 = W^Q h_5$$

$$\rightarrow k_1 \quad k_2 \quad k_3 \quad k_4 \quad k_5$$

$$\rightarrow v_1 \quad v_2 \quad v_3 \quad v_4 \quad v_5$$

$$\left( q_5 \cdot k_1 \right) v_1$$

# Connection to GPs

Going back to non-parametric kernl based methods (e.g. GPs), we compare the input $x$ to each of the training examples $X$ using a kernel to get a vector of similarity scores $\alpha = |K(x, x_i)|_{i=1}^{m}$, which we then use to retrieve a weighted combination of the corresponding target values $y_i$ as :
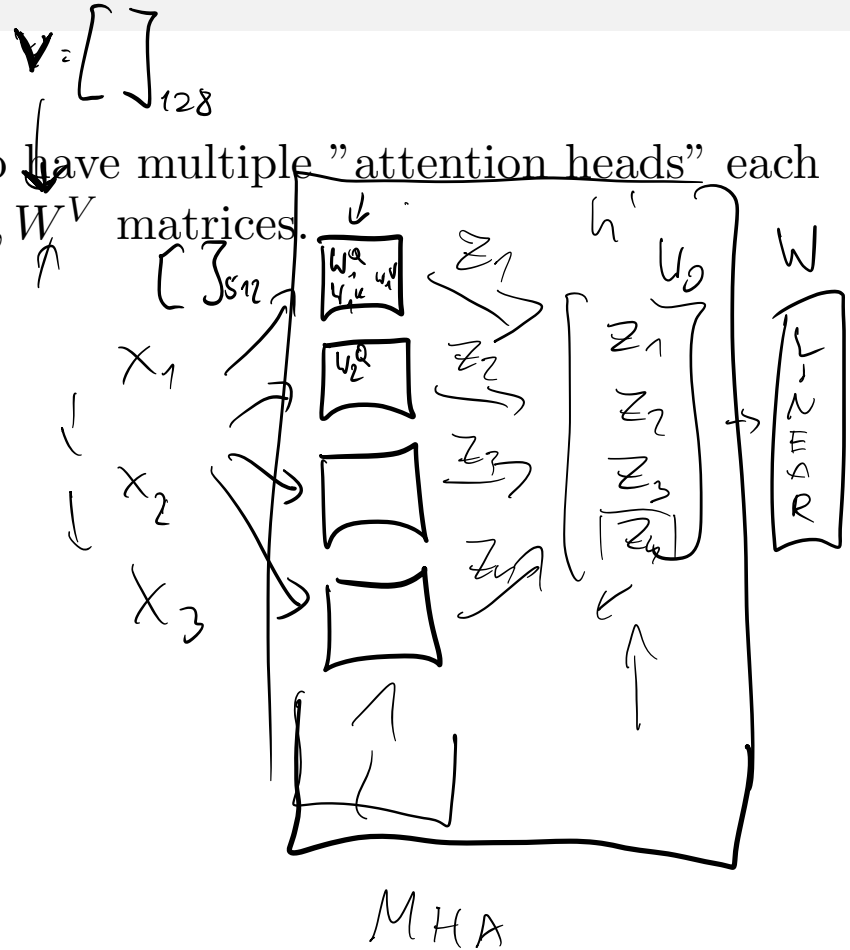
# Connection to GPs

Going back to non-parametric kernl based methods (e.g. GPs), we compare the input $x$ to each of the training examples $X$ using a kernel to get a vector of similarity scores $\alpha = |K(x, x_i)|_{i=1}^{m}$, which we then use to retrieve a weighted combination of the corresponding target values $y_i$ as :

$$\hat{y} = \sum_{i=1}^{m} \alpha_i y_i$$

# Connection to GPs

Going back to non-parametric kernl based methods (e.g. GPs), we compare the input $x$ to each of the training examples $X$ using a kernel to get a vector of similarity scores $\alpha = |K(x, x_i)|_{i=1}^m$, which we then use to retrieve a weighted combination of the corresponding target values $y_i$ as :

$$\hat{y} = \sum_{i=1}^m \alpha_i y_i$$

If we replace the stored examples matrix X with a learned embedding $K = W^K X$, stored outputs with $V = W^V Y$, and create an input embedding $q = W^Q x$, we can arrive at attention!

$V = [\quad]_{128}$

In practice it is advantageous to have multiple "attention heads" each with a different set of $W^Q, W^K, W^V$ matrices.

- Why do you think that is?

$[\quad]_{512}$

$W^Q_1 \, W^V$
$W^K \, W^V$

$U_2^Q$

$X_1$

$X_2$

$X_3$

$Z_1$

$Z_2$

$Z_3$

$Z_4$

$h'$

$U_2$

$\begin{bmatrix} Z_1 \\ Z_2 \\ Z_3 \\ Z_4 \end{bmatrix}$

$W$

LINEAR

MHA

# Multi Head Attention and Self Attention

In practice it is advantageous to have multiple "attention heads" each with a different set of $W^Q, W^K, W^V$ matrices.
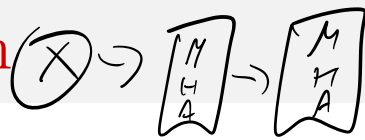
- Why do you think that is?
- Do we really need all of them?

# Multi Head Attention and Self Attention

In practice it is advantageous to have multiple "attention heads" each with a different set of $W^Q, W^K, W^V$ matrices.

- Why do you think that is?
- Do we really need all of them?

# Multi Head Attention and Self Attention

In practice it is advantageous to have multiple "attention heads" each with a different set of $W^Q, W^K, W^V$ matrices.

- Why do you think that is?
- Do we really need all of them?

We then simply concatenate the outputs of all of the attention heads together and multiplied by one final matrix $W^O$ that is learned as well, this is called **Multi Head Attention**.

$$o = MHA(Q, K, V) = Concat(h_1', \ldots, h_h')W^O$$

$$= Concat(Attn(Q_1, K_1, V_1), \ldots, Attn(Q_h, K_h, V_h))W^O$$

# Multi Head Attention and Self Attention

In practice it is advantageous to have multiple "attention heads" each with a different set of $W^Q, W^K, W^V$ matrices.

- Why do you think that is?
- Do we really need all of them?

We then simply concatenate the outputs of all of the attention heads together and multiplied by one final matrix $W^O$ that is learned as well, this is called **Multi Head Attention**.

$$f(\cdot) = MHA(Q, K, V) = Concat(h_1, \ldots, h_h)W^O$$
$$= Concat(Attn(Q_1, K_1, V_1), \ldots, Attn(Q_h, K_h, V_h))W^O$$

Additionally, we can stack several identical Attention / MHA blocks on top each other. This is called **Self-Attention**

$$X \to \boxed{MHA} \to \boxed{MHA} \cdots$$

$$Q_0 = X W_0^Q$$

1) This is our input sentence*

2) We embed each word*

3) Split into 8 heads. We multiply X or R with weight matrices

4) Calculate attention using the resulting Q/K/V matrices

5) Concatenate the resulting Z matrices, then multiply with weight matrix W⁰ to produce the output of the layer

Thinking Machines

**X**

**W₀ᵠ**
**W₀ᴷ**
**W₀ᵛ**

**Q₀**
**K₀**
**V₀**

**Z₀**

**W⁰**

* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one

**R**

**W₁ᵠ**
**W₁ᴷ**
**W₁ᵛ**

**Q₁**
**K₁**
**V₁**

**Z₁**

**Z**

...

**W₇ᵠ**
**W₇ᴷ**
**W₇ᵛ**

...

**Q₇**
**K₇**
**V₇**

...

**Z₇**

# Transformer

First proposed in a 2017 paper "Attention is all you need", the **Transformer** architecture consists of two stacks (called **Encoder** and **Decoder**) of blocks:
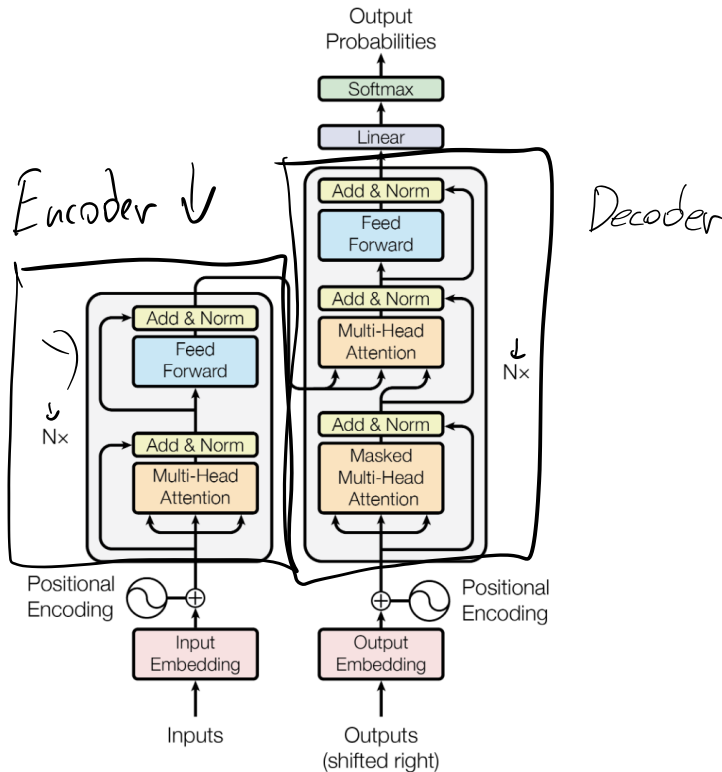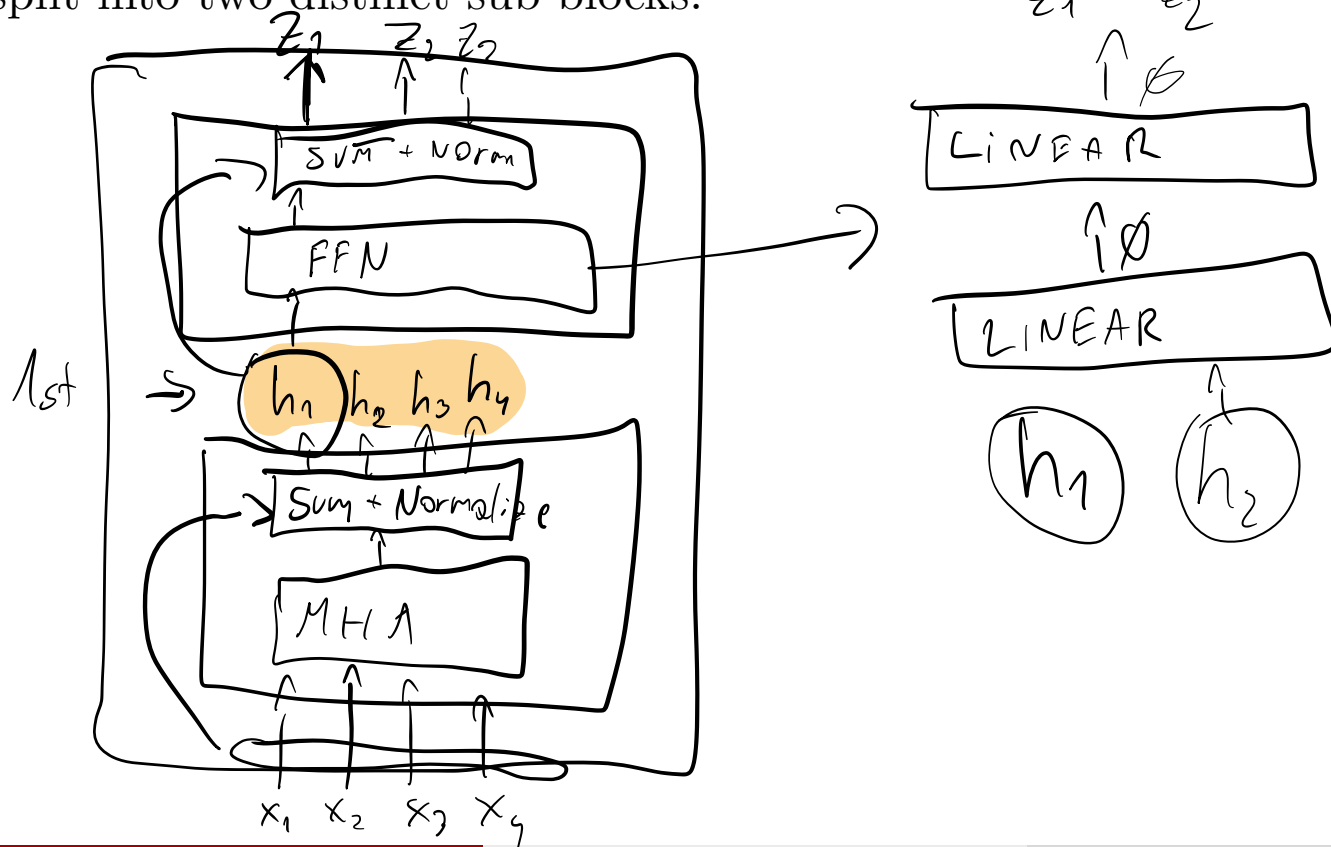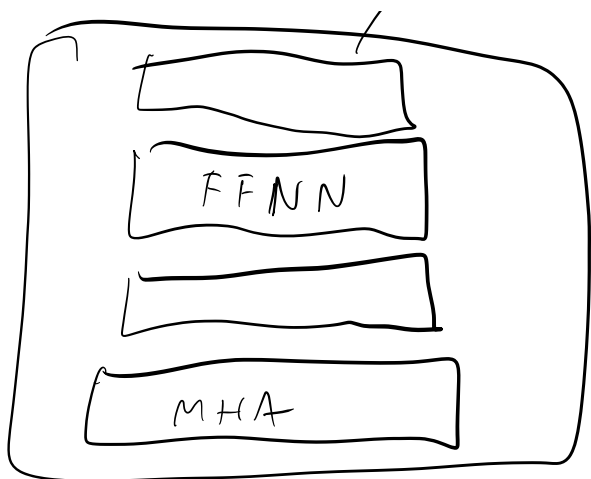
# Transformer

First proposed in a 2017 paper "Attention is all you need", the **Transformer** architecture consists of two stacks (called **Encoder** and **Decoder**) of blocks:



Figure 1: The Transformer - model architecture.

The **Encoder** consists of a stack of 6 blocks. Each block is further split into two distinct sub-blocks.
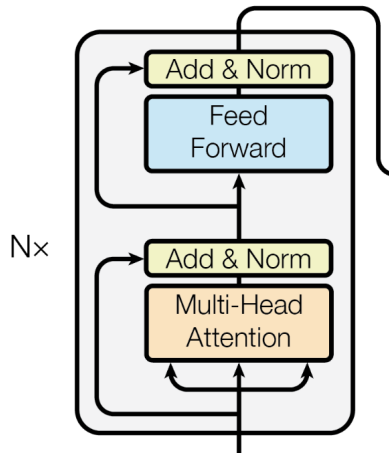
FFNN

MHA

$Z_1$  $Z_2$  $Z_3$  $Z_4$

# Transformer

The **Encoder** consists of a stack of 6 blocks. Each block is further split into two distinct sub-blocks.
The first is a Multi Head Self Attention mechanism, and the second is a simple FFNN. Both of the sub-blocks have a residual connection around them, followed by normalization.

# Transformer

The **Encoder** consists of a stack of 6 blocks. Each block is further split into two distinct sub-blocks.
The first is a Multi Head Self Attention mechanism, and the second is a simple FFNN. Both of the sub-blocks have a residual connection around them, followed by normalization.

# Transformer

Similarily, the **Decoder** is also a stack of 6 blocks.

# Transformer

Similarily, the **Decoder** is also a stack of 6 blocks.
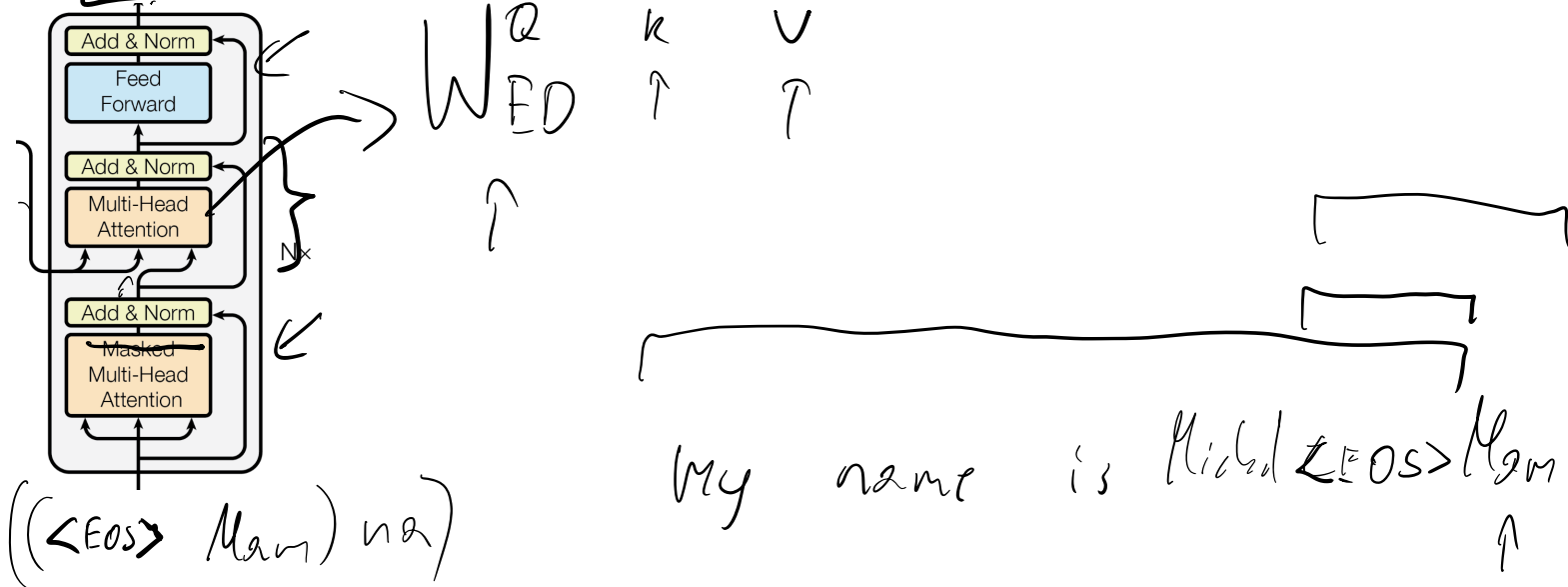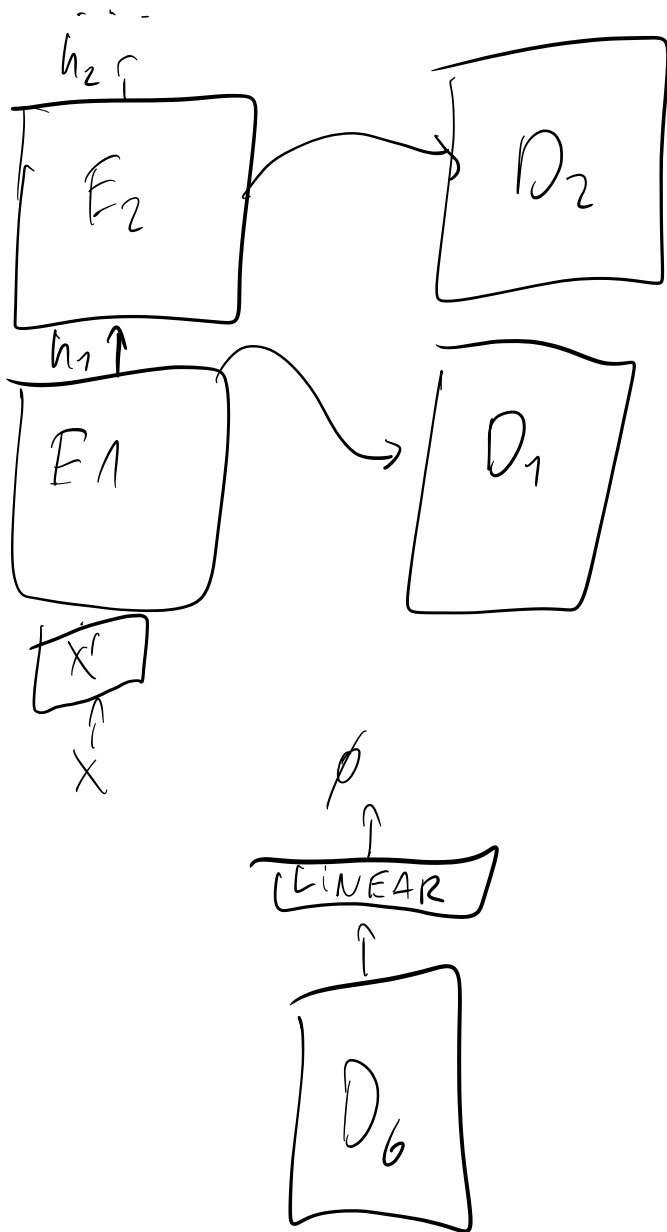However in addition to the two sub-blocks of the encoder, it features a 3rd sub-block.

# Transformer

Similarily, the **Decoder** is also a stack of 6 blocks.
However in addition to the two sub-blocks of the encoder, it features a 3rd sub-block.
This 3rd sub-block performs multi-head attention over the output of the encoder. This "encoder-decoder attention" layer uses $Q$ from the previous decoder layer, and $K, V$ from the output of the encoder.

# Transformer

Similarily, the **Decoder** is also a stack of 6 blocks.
However in addition to the two sub-blocks of the encoder, it features a 3rd sub-block.
This 3rd sub-block performs multi-head attention over the output of the encoder. This "encoder-decoder attention" layer uses $Q$ from the previous decoder layer, and $K, V$ from the output of the encoder.

$\dot{h_2}$ $\,\dot{}$

$E_2$

$D_2$

$h_1$

$E_1$

$D_1$

$X^r$

$X$

$p$

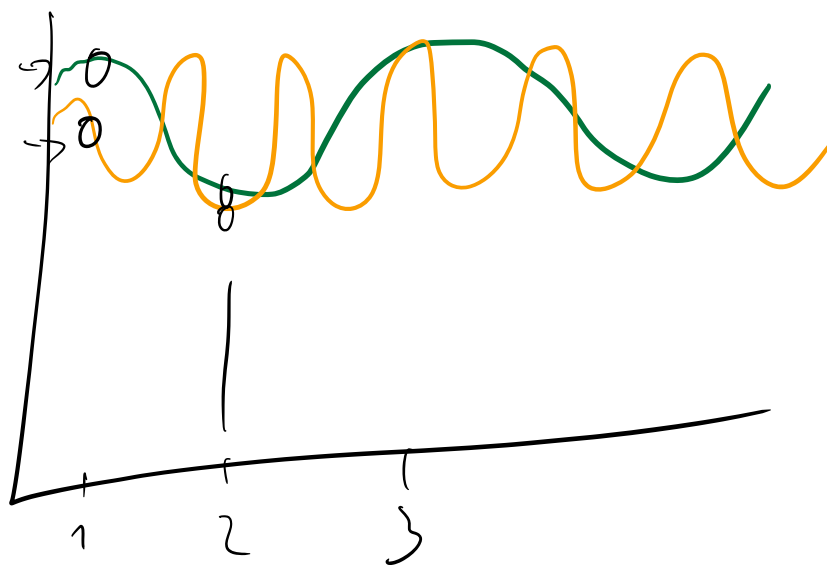LINEAR

$D_6$

# Transformer

What about inputs?

# Transformer

What about inputs? The input embedding is a learneable "static" **token** embedding similar to the Word2Vec model we have seen in the lecture 9.

$$M_y = \begin{bmatrix} v \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

pos  1    2    3    4          [ ]

$X_2$  $X_1$  $X_3$  $X_4$

$q_2$  $q_1$  $q_3$  $q_4$          $\phi \left( \dfrac{q^T k}{\sqrt{d}} \right) V$

$k_2$  $k_1$  $k_3$  $h_4$

$V_2$  $V_1$  $V_3$  $V_4$

$$\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \cdots \quad U_2$$

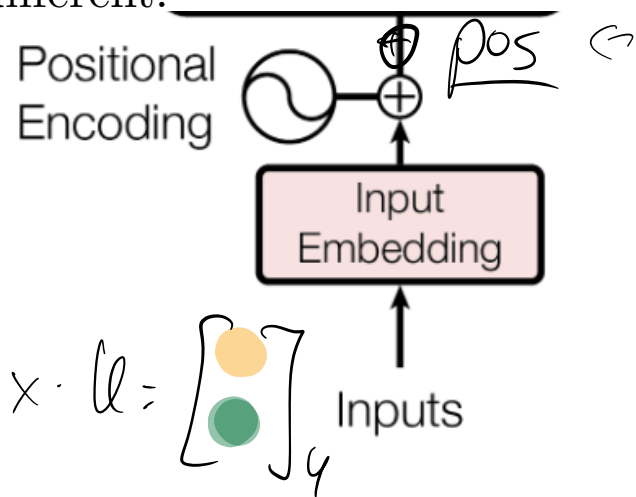$$x' = x \cdot U + pox_x \, U_2$$

# Transformer

What about inputs? The input embedding is a learneable "static" **token** embedding similar to the Word2Vec model we have seen in the lecture 9.
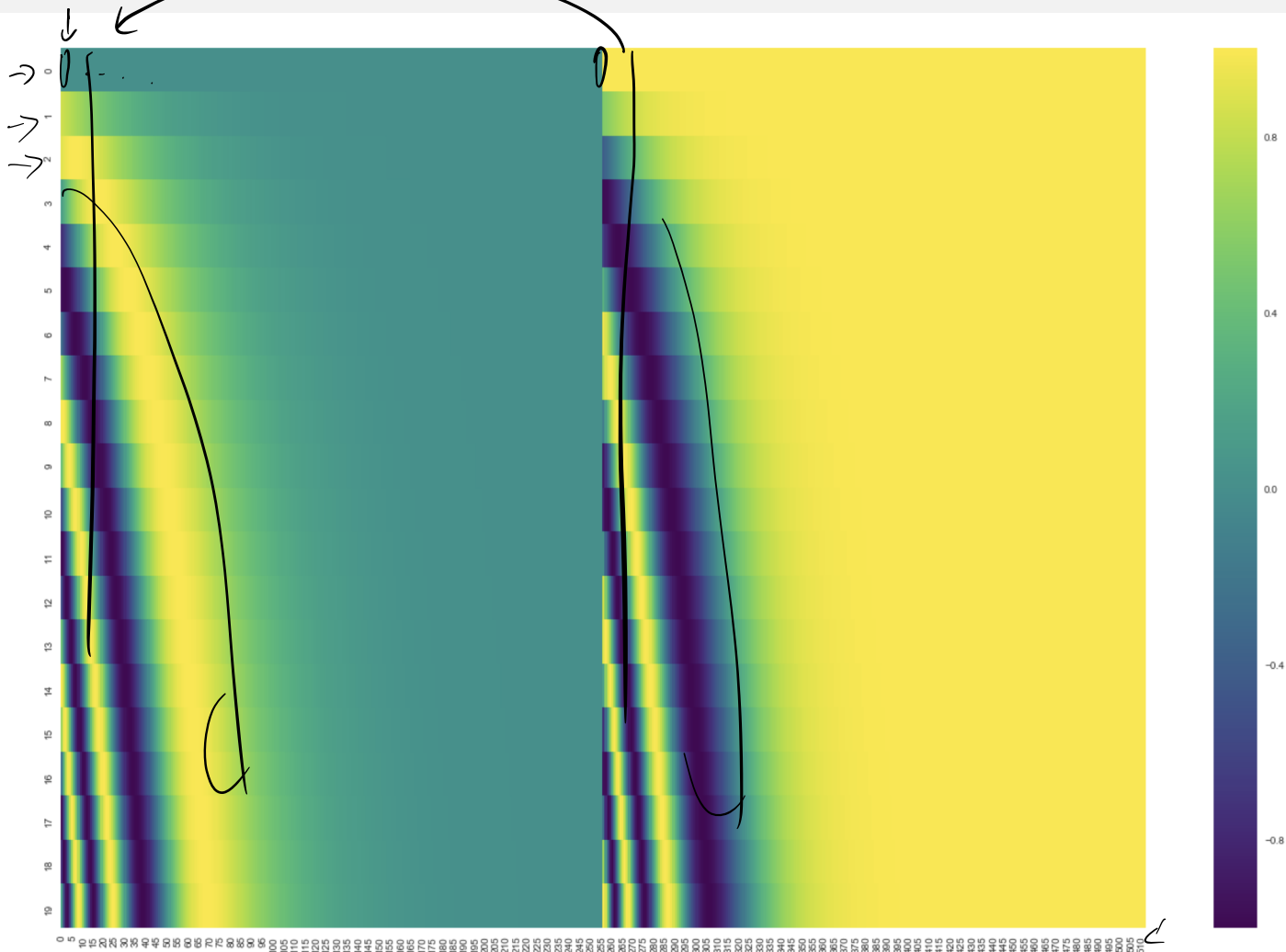What is "Positional Encoding?"

# Transformer

What about inputs? The input embedding is a learneable "static" **token** embedding similar to the Word2Vec model we have seen in the lecture 9.

What is "Positional Encoding?" It's either a learneable (representing position in a sequence) embedding, or a predefined embedding of different.

# Transformer

What about inputs? The input embedding is a learneable "static" **token** embedding similar to the Word2Vec model we have seen in the lecture 9.

What is "Positional Encoding?" It's either a learneable (representing position in a sequence) embedding, or a predefined embedding of different.
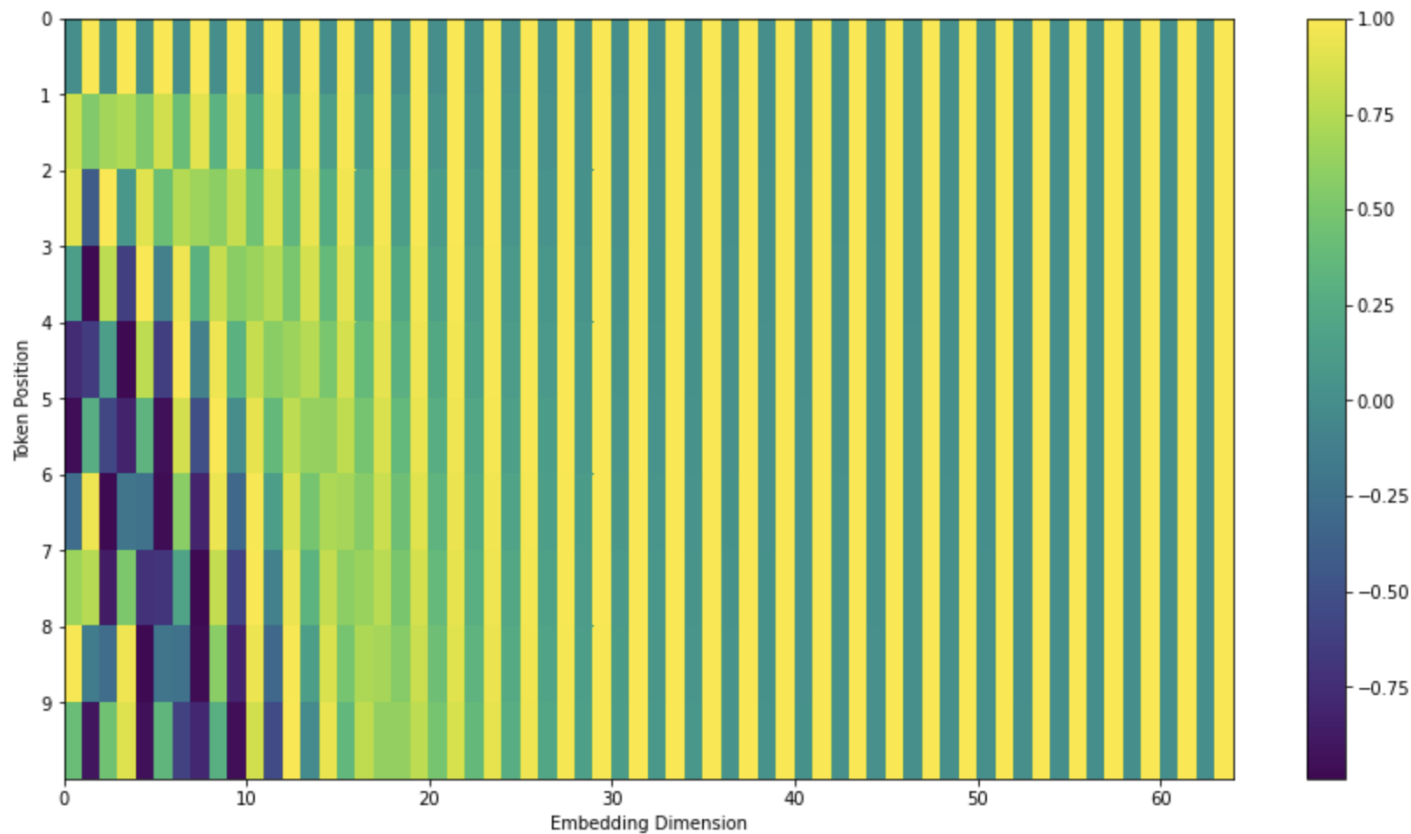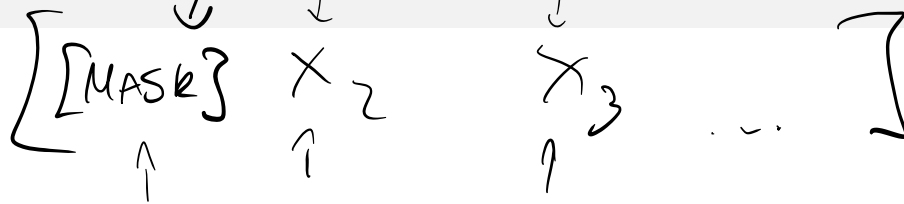
# Positional Encoding



A real example of positional encoding for 20 words (rows) with an embedding size of 512 (columns). You can see that it appears split in
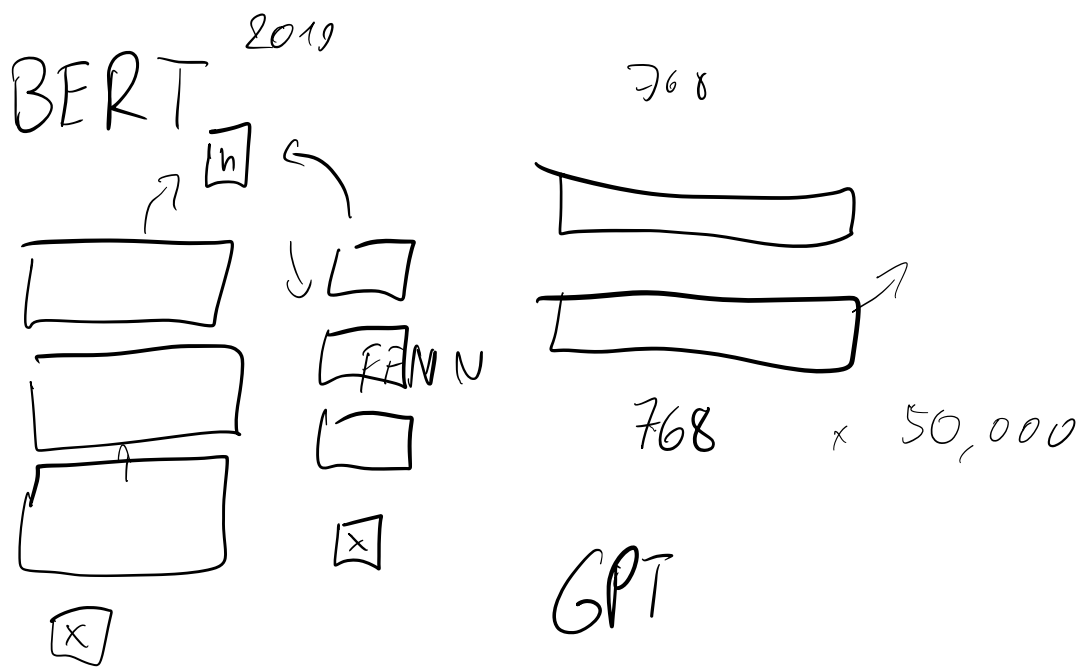
# Positional Encoding
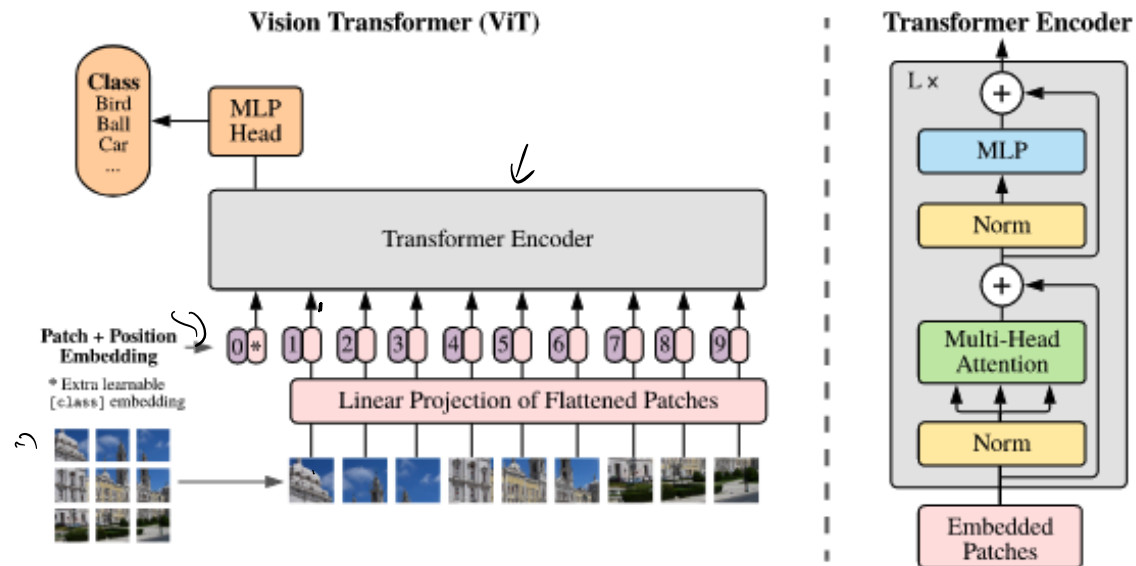
# How to train a Transformer

$$\left[ \ [\text{MASK}] \quad X_2 \quad\quad X_3 \quad ... \ \right]$$

The original Transformer model was trained on an English - German translations, where at each step the final decoder state was fed into a simple Linear Layer followed by a softmax to produce probabilities over next tokens.

Currently there are a large number of pre-training tasks (similar in idea to W2V). One of the most common ones is **Masked Language Modelling**, where we randomly replace 15% of tokens with "[MASK]", and the goal of the model is to predict back the original token.
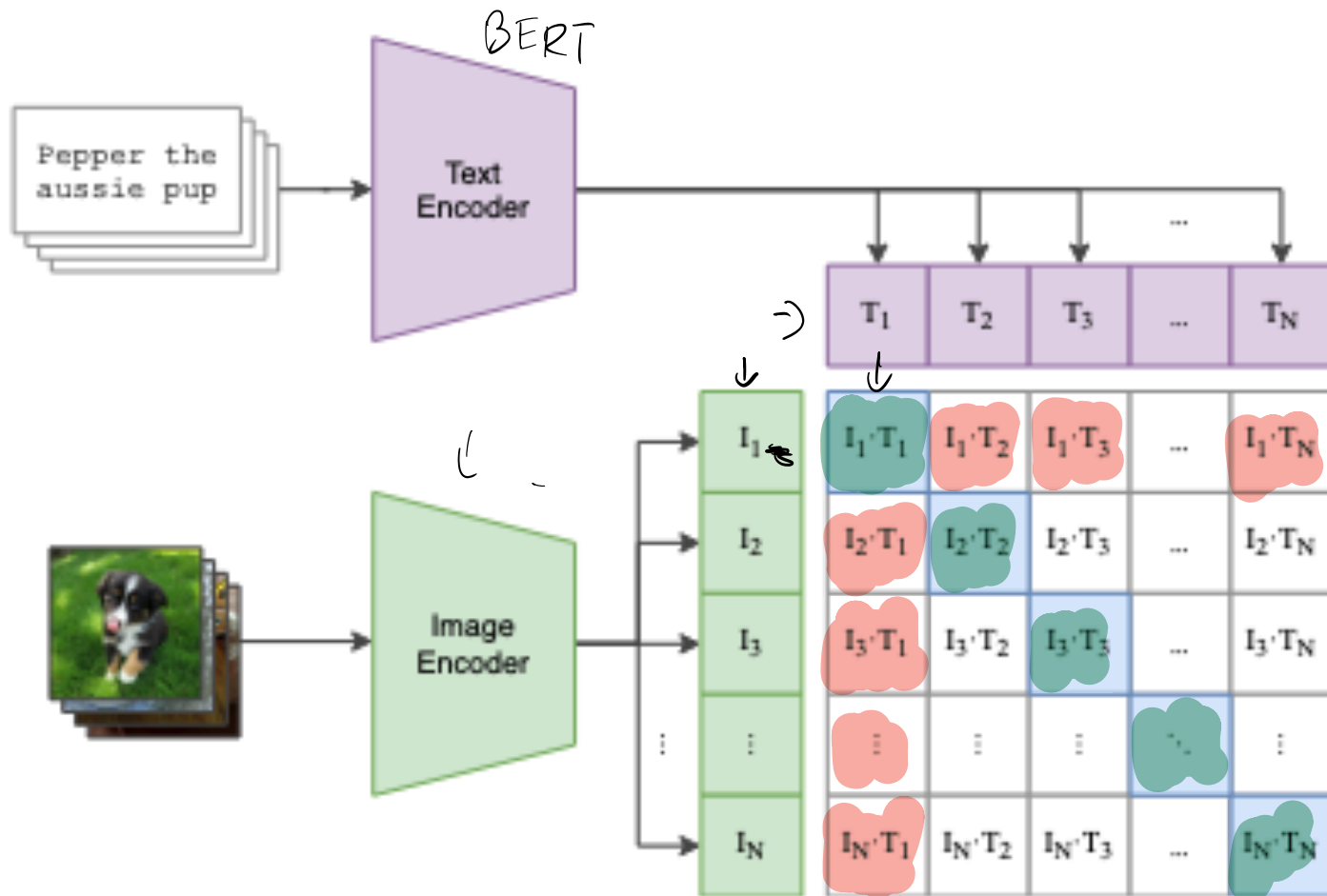
BERT [2019]

[n]

FFN N

[x]

[x]

768

768          × 50,000

GPT

huggingface transformers

# Vision Transformers



**Vision Transformer (ViT)**

Class
Bird
Ball
Car
...

MLP Head

Transformer Encoder

Patch + Position Embedding
* Extra learnable [class] embedding

0* 1 2 3 4 5 6 7 8 9

Linear Projection of Flattened Patches

**Transformer Encoder**

L ×

MLP

Norm

Multi-Head Attention

Norm

Embedded Patches

# CLIP

Demo