

Transformers

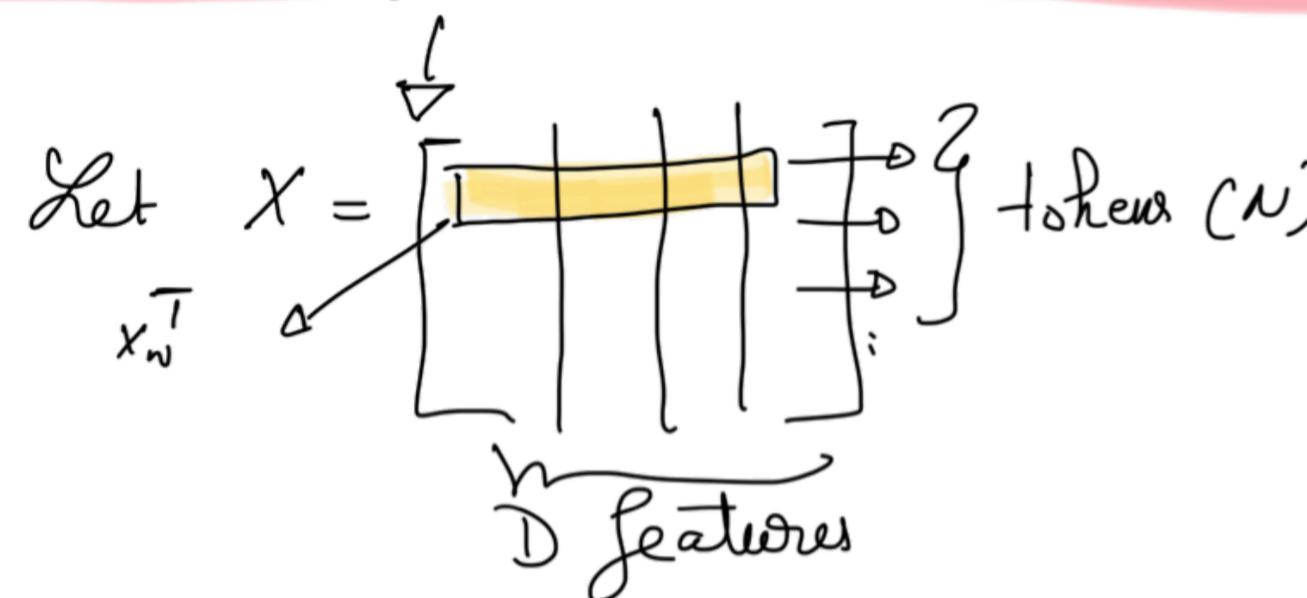
- Based on a processing called 'attention'
- Inductive Bias is what a model believes before seeing the data
 - CNNs: Local Patterns matter (pixels next to each other)
 - RNNs: things happen in order and recent inputs matter more
 - Transformers: Any input can relate to any other input; no built-in locality or order.

These models are known as transformers because they transform a set of vectors in some representation space into a corresponding set of vectors, having the same dimensionality in a new space. This new space has a richer internal representation.

- ① **Attention:** It is a mechanism that allows the model to assign different levels of importance to different parts of the input, enabling it to selectively focus on the most relevant information when computing a representation or making predictions.

1 - Transformer Processing

- The input data to a transformer is a set of vectors $\{x_n\}_{n=1}^N$ of dimensionality D .
- Those data vectors are called "tokens". A token might correspond to a word within a sentence or a patch in an image.
- The elements x_n of each token are called features.



$$\tilde{X} = \text{Transformer Layer}(X)$$

a function that takes the data matrix X as input and creates a transformed matrix \tilde{X} of the same dimensionality as output

- Multiple transformer layers are applied in succession to create a deep network capable of learning rich internal representations.

- Each transformer layer contains its own weights and biases. (learned using gradient descent using an appropriate cost function)

- Each transformer layer is composed by 2 stages:

- ① **Attention Mechanism** mixes information between different tokens (rows) → for each token (word), it looks at other tokens and decides which ones are important and combines their information.

- ② **Feed Forward Transformation** works on each token separately → updates the token's internal representation without looking at other tokens.

- 2 - **Attention Coefficients**: we have a set of input tokens x_1, \dots, x_N in an embedding space and we want to map this to another set y_1, \dots, y_N having the same number of tokens but in a new embedding space that captures a richer semantic structure.

- The value of y_N should depend not just on the corresponding input vector x_N but on all the vectors x_1, \dots, x_N in the set.

- A simple way to achieve this is to define each output vector y_N as a linear combination of the input vectors x_1, \dots, x_N with weighting coefficients.

$$y_N = \sum_{m=1}^N a_{mn} x_m$$

attention weights

Constraints:

$$a_{mn} \geq 0$$

$$\sum_{m=1}^N a_{mn} = 1$$

- 3 - **Self Attention**: given a sequence, we want each token to look at all the other tokens (including itself) and decide how much attention to pay to each one to compute a new representation.

↳ This is done using:

Queries (Q) → what this token is asking for?

Keys (K) → what this token offers?

Values (V) → the actual value/info this token carries.

- ① **Linear Projections**: we project X into Query, Key and Value matrices using learnable weights.

$$Q = X \cdot W_Q \quad ; \quad K = X \cdot W_K \quad ; \quad V = X \cdot W_V$$

where: $W_Q, W_K \in \mathbb{R}^{d \times d_h}; W_V \in \mathbb{R}^{d \times d_v}$

d_h is usually set to $d_h = \frac{d}{\text{num_heads}}$

- ② **Attention Scores**: we compute the pairwise product between all queries and keys

$$\text{Scores} = \frac{Q \cdot K^T}{\sqrt{d_h}} \in \mathbb{R}^{N \times N}$$

- ③ **Softmax to normalize**

$$d_{nm} = \frac{\exp[\text{Score}_{nm}]}{\sum_{m=1}^N \exp[\text{Score}_{nm}]}$$

There is no probabilistic interpretation to the Softmax in this case and it is only used to normalize the attention weights (on each row independently) to make them sum to 1.

- ④ **Weighted Sum of Values**:

$$f_N = \sum_{m=1}^N d_{nm} V_m \quad \text{or} \quad \text{Attention}(Q, K, V) = \text{softmax}\left(\frac{Q \cdot K^T}{\sqrt{d_h}}\right) \cdot V$$

- ⑤ **Example**:

Input 3-tokens with 2 features $\rightarrow X = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 1 \end{bmatrix} \in \mathbb{R}^{3 \times 2}$

For simplicity now, we take $d_{\text{model}} = 2 = d_h$ and let $W_Q = W_K = W_V = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$

- 1) Compute $Q \cdot K^T$:

$$Q = X \cdot W_Q = \underbrace{\begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 1 \end{bmatrix}}_{3 \times 2} \underbrace{\begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix}}_{2 \times 2} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 1 \end{bmatrix}$$

$$K = Q \rightarrow Q \cdot K^T = \underbrace{\begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 1 \end{bmatrix}}_{3 \times 2} \underbrace{\begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 1 \\ 0 & 1 & 1 \end{bmatrix}}_{2 \times 3} = \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 1 \\ 0 & 1 & 2 \end{bmatrix}$$

- 2) Rows? → scale by d_h

$$\frac{Q \cdot K^T}{\sqrt{d_h}} = \begin{bmatrix} 0,707 & 0 & 0,707 \\ 0 & 0,707 & 0,707 \\ 0 & 0,707 & 1,414 \end{bmatrix}$$

- 3) Apply Softmax row-wise:

$$\text{row 1: } \text{softmax} [0,707 \ 0 \ 0,707] = [0,401 \ 0,198 \ 0,401]$$

- 4) Multiply by V :

$$y_1 = 0,401 \underbrace{\begin{bmatrix} 0 & 1 \\ 1 & 1 \end{bmatrix}}_{\text{token 1}} + 0,198 \underbrace{\begin{bmatrix} 0 & 1 \\ 0 & 1 \end{bmatrix}}_{\text{token 2}} + 0,401 \underbrace{\begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}}_{\text{token 3}}$$

$$\rightarrow y_1 = [1,203, 0,599]$$

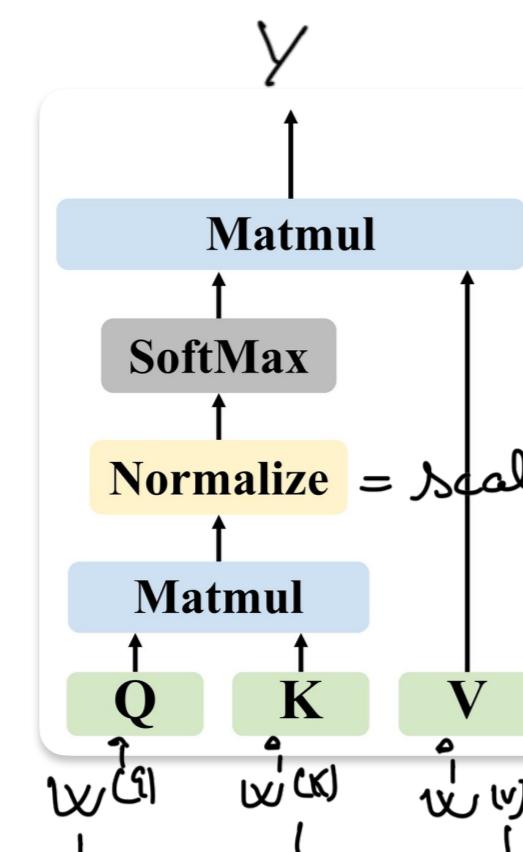
- For each output y_n , the corresponding input x_n is considered the query while all inputs x_m (including itself) are considered keys and values.
- This process is called Self Attention because we use the same sequence to compute the queries, keys and values.
- This is also known as dot product Self Attention because the similarity between key and query is computed by a dot product.

③ Attention Parameters

Attention parameters are $\mathbf{W}^{(Q)}$, $\mathbf{W}^{(K)}$ and $\mathbf{W}^{(V)}$.

- $\mathbf{W}^{(Q)}$ has dimensionality $d \times d_k$. So we can do the dot product.
- $\mathbf{W}^{(Q)} = \mathbf{W}^{(K)} = \mathbf{W}^{(V)}$
- $d \times d_v$. If $d = d_v$, the output representation has the same dimensionality as the input.

$\mathbf{W}^{(Q)}$, $\mathbf{W}^{(K)}$ and $\mathbf{W}^{(V)}$ are learnable linear transformations (fully connected layers without activation) applied to the input tokens to produce the query, key and value. In practice, bias parameters are also included.



④ Multihead Attention?

Why? Single head attention computes one similarity space, but that might not be enough to capture all types of relationships (position, semantics etc).

- MHA lets the model focus on different aspects of the input at the same time, capturing richer and more diverse relationships.

Idea?

- We split the input into multiple heads.
- Each head has its own \mathbf{W}^Q , \mathbf{W}^K and \mathbf{W}^V .
- Each head computes its own self-attention.
- Then we concatenate the results and project them back to the model dimension.

Let's

$$\textcircled{1} \quad X \in \mathbb{R}^{N \times d} \quad \text{8 input tokens}$$

$$h = \text{nb of heads}$$

$$d_h = d_v = d/h \quad \text{8 dimension per head}$$

$$\textcircled{2} \quad \text{for each head } i=1 \dots h \text{ compute:}$$

$$Q_i = \frac{X \cdot \mathbf{W}_Q^{(i)}}{\sqrt{d_h}} \in \mathbb{R}^{N \times d_h}$$

$$K_i = X \cdot \mathbf{W}_K^{(i)} \in \mathbb{R}^{N \times d_h}$$

$$V_i = X \cdot \mathbf{W}_V^{(i)} \in \mathbb{R}^{N \times d_h}$$

$$\textcircled{3} \quad \text{then apply scaled dot product attention per head}$$

$$\text{Head } i = \text{softmax}\left[\frac{-Q_i \cdot K_i^T}{\sqrt{d_h}}\right] \cdot V_i \in \mathbb{R}^{N \times d_h}$$

⑤ Concatenate and project

$$\text{MHA}(X) = \text{Concat}(\text{Head}_1 \dots \text{Head}_h) \cdot \mathbf{W}_O$$

$$\in \mathbb{R}^{N \times d}$$

$$\in \mathbb{R}^{N \times h \cdot d_h}$$

$$\in \mathbb{R}^{N \times d_h \cdot d}$$

Even though the concatenated output already has the same shape as the input, the final projection with \mathbf{W}_O allows to learn how to combine the information from all the heads into a meaningful representation, because the concatenation allows to just put the heads side by side \rightarrow disjoint ignoring

Algorithm 12.2: Multi-head attention

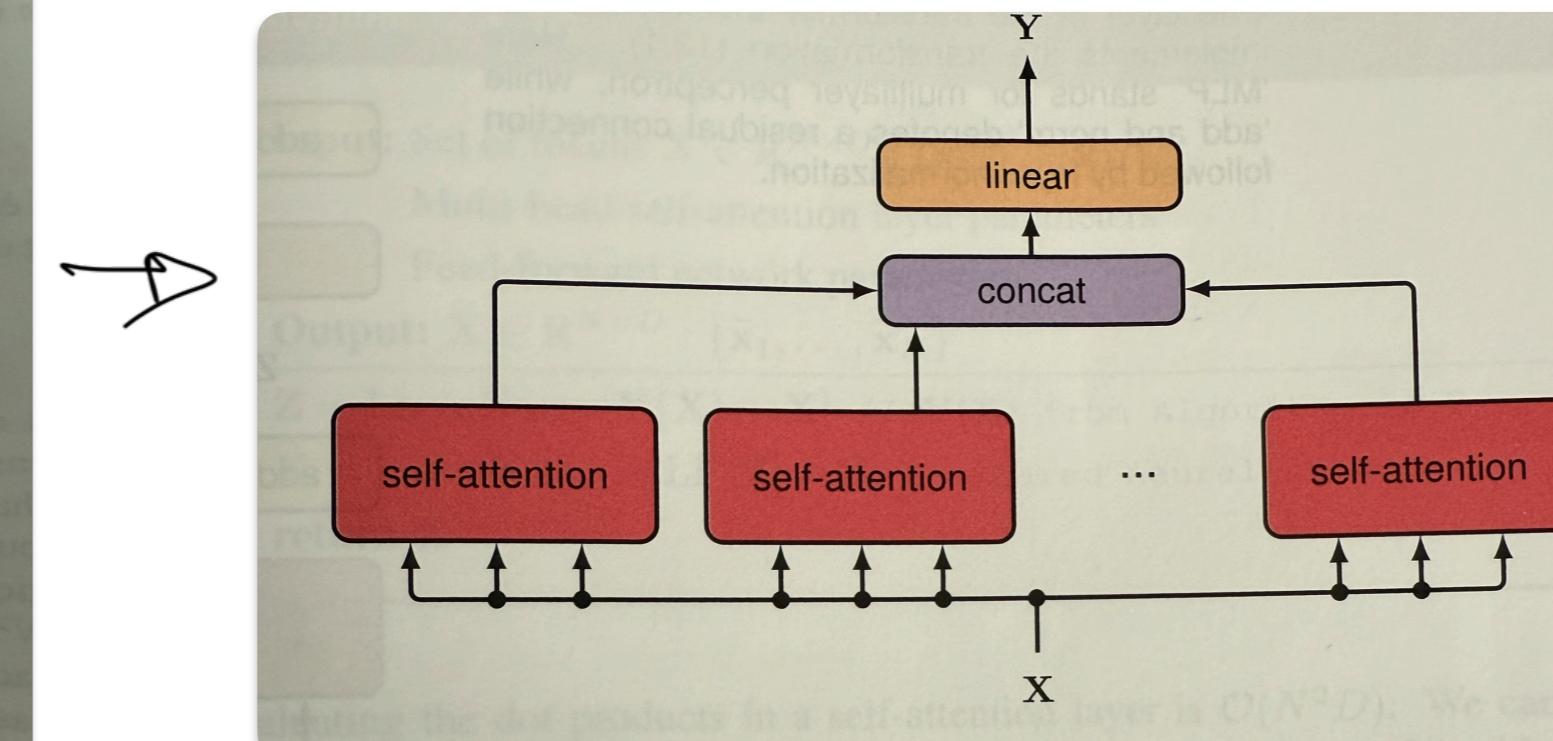
```

Input: Set of tokens  $X \in \mathbb{R}^{N \times D} : \{x_1, \dots, x_N\}$ 
Query weight matrices  $\{\mathbf{W}_1^{(q)}, \dots, \mathbf{W}_H^{(q)}\} \in \mathbb{R}^{D \times D}$ 
Key weight matrices  $\{\mathbf{W}_1^{(k)}, \dots, \mathbf{W}_H^{(k)}\} \in \mathbb{R}^{D \times D}$ 
Value weight matrices  $\{\mathbf{W}_1^{(v)}, \dots, \mathbf{W}_H^{(v)}\} \in \mathbb{R}^{D \times D_v}$ 
Output weight matrix  $\mathbf{W}^{(o)} \in \mathbb{R}^{H \cdot D_v \times D}$ 

Output:  $Y \in \mathbb{R}^{N \times D} : \{y_1, \dots, y_N\}$ 

// compute self-attention for each head (Algorithm 12.1)
for  $h = 1, \dots, H$  do
     $Q_h = X \mathbf{W}_h^{(q)}$ ,  $K_h = X \mathbf{W}_h^{(k)}$ ,  $V_h = X \mathbf{W}_h^{(v)}$ 
     $H_h = \text{Attention}(Q_h, K_h, V_h)$  //  $H_h \in \mathbb{R}^{N \times D_v}$ 
end for
 $H = \text{Concat}[H_1, \dots, H_H]$  // concatenate heads
return  $Y(X) = H \mathbf{W}^{(o)}$ 

```



⑥ Transformer Layers

Multilevel Self Attention (MSA) is the core architecture of the transformer network.

Neural Networks benefit greatly from depth \rightarrow we want to stack multiple MSA layers on top of each other.

To improve training efficiency, **Residual Connection** (adding the input to the output) was introduced \rightarrow we need the output dimensionality to be equal to the input one $\rightarrow N \times D$.

Why residual connections?

- Stabilize gradients because when we stack many layers, gradients can explode or vanish.
- Allow the model to learn adjustments \rightarrow it gives the model flexibility to learn when to change and when to keep information.
- Enable deeper architectures \rightarrow without residuals, deep transformers (with many layers) become hard to train.

Then, we apply **layer normalization** to improve training efficacy

$$\begin{aligned} \mathbf{z} &= \text{LayerNorm}[\mathbf{x}] \\ \mathbf{z} &= \text{LayerNorm}[\mathbf{x} + \mathbf{x}] \end{aligned}$$

or $\mathbf{z} = \mathbf{y} / \sqrt{\text{LayerNorm}(\mathbf{x}) + \epsilon}$

max effective

The output of the MSA

Even though non-linearity enters through attention weights (softmax) but the outputs are still constrained to lie in the subspace spanned by the input vectors \rightarrow limit the expressiveness of the attention layer \rightarrow we use a standard non-linear neural network with D inputs and D outputs as a post-processing of the output of each layer.

$$\mathbf{z} = \text{LayerNorm}[\text{MLP}(\mathbf{z}) + \mathbf{z}]$$

the final output
↓
One transformation layer

$$\mathbf{z} = \text{MLP}[\text{LayerNorm}(\mathbf{z})] + \mathbf{z}$$

⑦ Computational Complexity

It's a way to estimate how much work an algorithm needs depending on the input.

1 - Linear proj & calculate Q, K and V.

$$\mathbf{Q} = \frac{X \cdot \mathbf{W}_Q^{(i)}}{\sqrt{d_h}} \rightarrow \text{shape} = (N, D)$$

$$\text{token 1: } \left(\begin{array}{c} \vdots \\ \vdots \\ \vdots \end{array} \right) \left(\begin{array}{c} \vdots \\ \vdots \\ \vdots \end{array} \right) \dots \left(\begin{array}{c} \vdots \\ \vdots \\ \vdots \end{array} \right) \xrightarrow{\text{5 multiplication}} \left(\begin{array}{c} \vdots \\ \vdots \\ \vdots \end{array} \right) \xrightarrow{\text{4 additions}} \left(\begin{array}{c} \vdots \\ \vdots \\ \vdots \end{array} \right) \xrightarrow{\mathcal{O}(D^2)} \mathcal{O}(D^2)$$

for token 1 on a $D \times D$ matrix $\rightarrow \mathcal{O}(D^2)$

$$\text{2 - Attention Scores calculation:}$$

$$\text{scores} = \mathbf{Q} @ \mathbf{K}^T \xrightarrow{\mathcal{O}(N \cdot D \cdot D)}$$

$$\text{token 1: } \left(\begin{array}{c} \vdots \\ \vdots \\ \vdots \end{array} \right) \left(\begin{array}{c} \vdots \\ \vdots \\ \vdots \end{array} \right) \dots \left(\begin{array}{c} \vdots \\ \vdots \\ \vdots \end{array} \right) \xrightarrow{\text{3 multiplication + 2 addition}} \left(\begin{array}{c} \vdots \\ \vdots \\ \vdots \end{array} \right) \xrightarrow{\mathcal{O}(N \cdot D)}$$

$$\text{token 2: } \left(\begin{array}{c} \vdots \\ \vdots \\ \vdots \end{array} \right) \left(\begin{array}{c} \vdots \\ \vdots \\ \vdots \end{array} \right) \dots \left(\begin{array}{c} \vdots \\ \vdots \\ \vdots \end{array} \right) \xrightarrow{\text{3 multiplication + 2 addition}} \left(\begin{array}{c} \vdots \\ \vdots \\ \vdots \end{array} \right) \xrightarrow{\mathcal{O}(N \cdot D)}$$

for N tokens $\rightarrow \mathcal{O}(N^2 D)$

3 - Apply Attention weights to

$$\text{output} = \text{softmax}(\text{scores}) \cdot \mathbf{V}$$

$$\text{token 1: } \left(\begin{array}{c} \vdots \\ \vdots \\ \vdots \end{array} \right) \left(\begin{array}{c} \vdots \\ \vdots \\ \vdots \end{array} \right) \dots \left(\begin{array}{c} \vdots \\ \vdots \\ \vdots \end{array} \right) \xrightarrow{\text{2 multiplication + 1 addition}} \left(\begin{array}{c} \vdots \\ \vdots \\ \vdots \end{array} \right) \xrightarrow{\mathcal{O}(D)}$$

$$\text{token 2: } \left(\begin{array}{c} \vdots \\ \vdots \\ \vdots \end{array} \right) \left(\begin{array}{c} \vdots \\ \vdots \\ \vdots \end{array} \right) \dots \left(\begin{array}{c} \vdots \\ \vdots \\ \vdots \end{array} \right) \xrightarrow{\text{2 multiplication + 3 addition}} \left(\begin{array}{c} \vdots \\ \vdots \\ \vdots \end{array} \right) \xrightarrow{\mathcal{O}(D)}$$

for N tokens $\rightarrow \mathcal{O}(N \cdot D)$

$$\text{total cost} = \mathcal{O}(ND^2 + N^2D + N^2D)$$

$$= \mathcal{O}(ND^2 + 2N^2D)$$

$$= \mathcal{O}(ND^2 + N^2D)$$

$$\cong \text{generally } N \gg D \rightarrow \mathcal{O}(N^2D)$$

SA is quadratic in sequence length