# Transformer Network Encoders

An overview

# What are Transformer Networks?

- Transformer Networks are a type of neural network architecture designed to handle sequential data

- Examples:

  - Chatbots

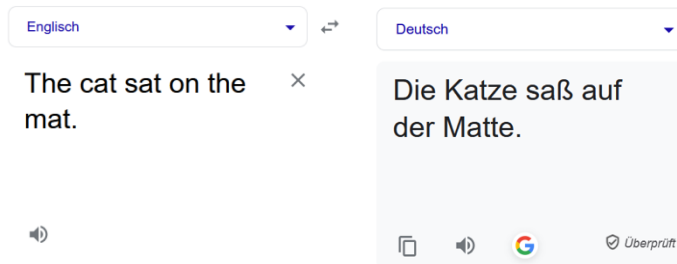    **You**
    What is a Large Language Model?

    **ChatGPT**
    A Large Language Model (LLM) like me, ChatGPT, is a type of artificial intelligence system designed to understand, generate, and respond to human language in a way that is both coherent and contextually relevant. Here are the key characteristics of a Large Language Model:
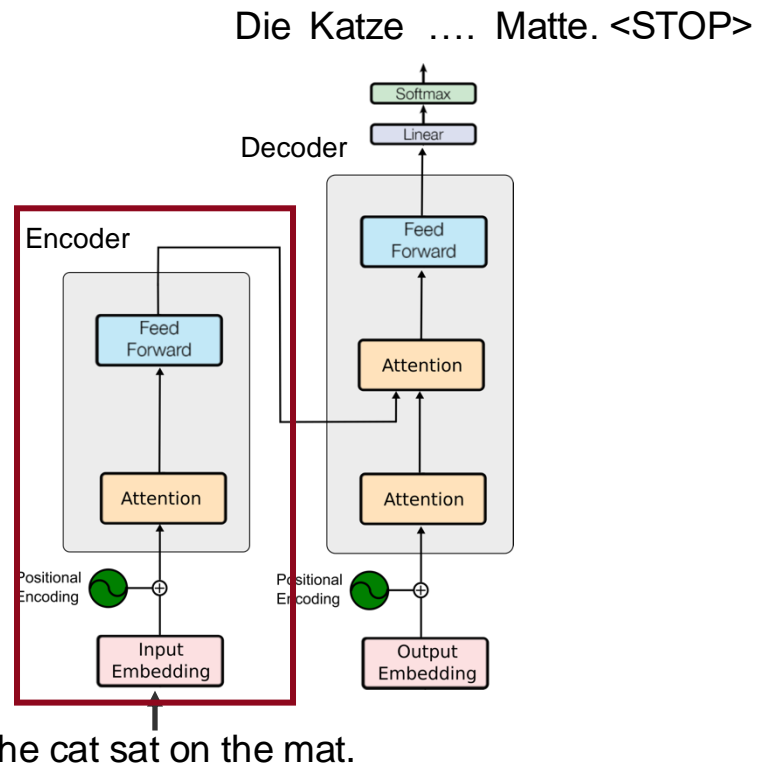
  - Translator

    | Englisch | ⇄ | Deutsch |
    |---|---|---|

    The cat sat on the mat.    ✕

    Die Katze saß auf der Matte.
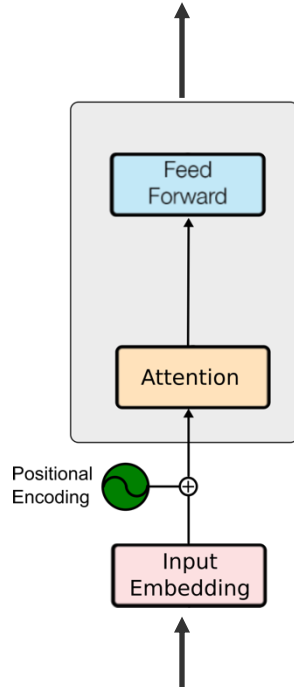
# Transformer Network Architecture

- **Transformer Network**
  - Encoder: Encodes the input sequence into numerical vectors
  - Decoder: Uses the numerically encoded input sequence as its input and produces an output sequence.

Die Katze .... Matte. <STOP>

Decoder

Encoder

Softmax

Linear

Feed Forward

Attention

Feed Forward

Attention

Attention

Positional Encoding ⊕

Positional Encoding ⊕

Input Embedding

Output Embedding

The cat sat on the mat.

hhu.de

# Encoder

2.71 0.18 ⋮ 1.43    2.1 -0.34 ⋮ 0.67    0.22 1.85 ⋮ -0.42    ...    2.17 0.16 ⋮ 0.29

Feed Forward
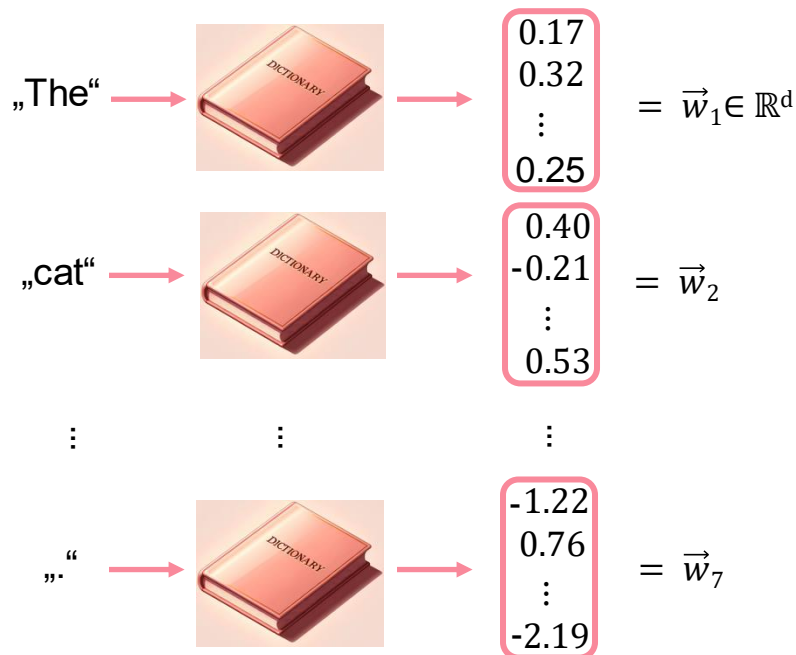
Attention

Positional Encoding ⊕

Input Embedding

The cat sat on the mat.

1. Input Embedding & Positional Encoding
2. Attention
3. Feed Forward Neural Network

# Input Embedding

■ For every word, we have a specific vector representing this word

„The" → 📖 → $\begin{bmatrix} 0.17 \\ 0.32 \\ \vdots \\ 0.25 \end{bmatrix}$ = $\vec{w}_1 \in \mathbb{R}^d$

„cat" → 📖 → $\begin{bmatrix} 0.40 \\ -0.21 \\ \vdots \\ 0.53 \end{bmatrix}$ = $\vec{w}_2$

⋮      ⋮      ⋮

„." → 📖 → $\begin{bmatrix} -1.22 \\ 0.76 \\ \vdots \\ -2.19 \end{bmatrix}$ = $\vec{w}_7$

Input Embedding

$\vec{w}_1 \quad \vec{w}_2 \quad \vec{w}_3 \quad \vec{w}_4 \quad \vec{w}_5 \quad \vec{w}_6 \quad \vec{w}_7$

↑ ↑ ↑ ↑ ↑ ↑ ↑

The cat sat on the mat.

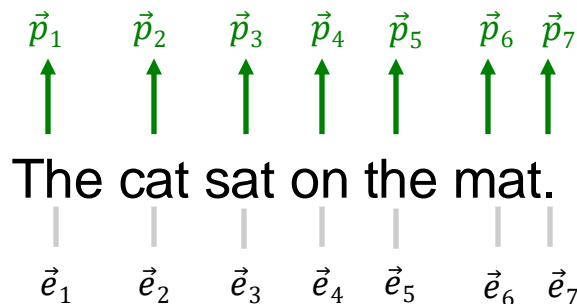$\vec{e}_1 \quad \vec{e}_2 \quad \vec{e}_3 \quad \vec{e}_4 \quad \vec{e}_5 \quad \vec{e}_6 \quad \vec{e}_7$
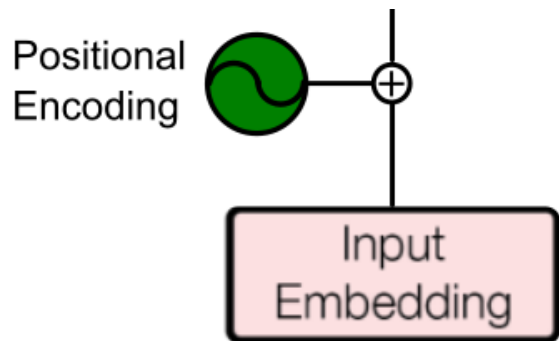
Positional
Encoding:

$$\vec{p}_1 = \begin{bmatrix} 0.00 \\ 0.84 \\ \vdots \\ 0.14 \end{bmatrix} \qquad \vec{p}_2 = \begin{bmatrix} 1.00 \\ 0.54 \\ \vdots \\ -0.99 \end{bmatrix} \qquad \dots \qquad \vec{p}_7 = \begin{bmatrix} 0.00 \\ 1.00 \\ \vdots \\ 0.96 \end{bmatrix} \in \mathbb{R}^d$$

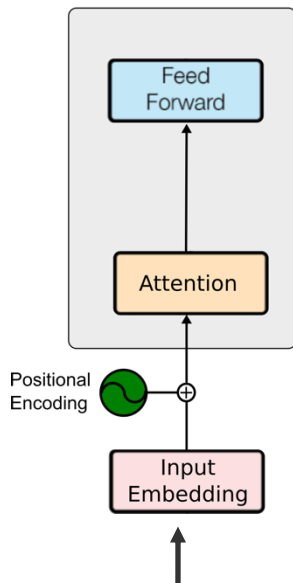- The same positional encoding is used for every input sequence!

Positional
Encoding

$\vec{p}_1 \quad \vec{p}_2 \quad \vec{p}_3 \quad \vec{p}_4 \quad \vec{p}_5 \quad \vec{p}_6 \; \vec{p}_7$

↑ ↑ ↑ ↑ ↑ ↑ ↑

The cat sat on the mat.

$\vec{e}_1 \quad \vec{e}_2 \quad \vec{e}_3 \quad \vec{e}_4 \quad \vec{e}_5 \quad \vec{e}_6 \; \vec{e}_7$

$$\vec{e}_1$$

$$\begin{matrix} 0.00 + 0.17 \\ 0.84 + 0.32 \\ \vdots \\ 0.14 + 0.25 \end{matrix}$$

$$\vec{e}_7$$

$$\begin{matrix} 0.00 - 1.22 \\ 1.00 + 0.76 \\ \vdots \\ 0.96 - 2.19 \end{matrix} \in \mathbb{R}^d$$

$$= \quad = \quad = \quad\quad\quad =$$

$$\vec{p}_1 \quad \vec{p}_2 \quad \vec{p}_3 \quad \vec{p}_4 \quad \vec{p}_5 \quad \vec{p}_6 \quad \vec{p}_7$$

$$+ \quad + \quad + \quad + \quad + \quad + \quad +$$

$$\vec{w}_1 \quad \vec{w}_2 \quad \vec{w}_3 \quad \vec{w}_4 \quad \vec{w}_5 \quad \vec{w}_6 \quad \vec{w}_7$$

The cat sat on the mat.

$$\vec{e}_1 \quad \vec{e}_2 \quad \vec{e}_3 \quad \vec{e}_4 \quad \vec{e}_5 \quad \vec{e}_6 \quad \vec{e}_7$$

Positional Encoding

Input Embedding

1. Input Embedding & Positional Encoding
2. Attention
3. Feed Forward Neural Network

$$Attention(\vec{e}_i, K, V) = softmax(\frac{\vec{e}_i^T \cdot K}{\sqrt{d}}) \cdot V^T$$

$$K = (\vec{e}_1, \vec{e}_2, \dots, \vec{e}_7), \qquad V^T = \begin{vmatrix} \vec{e}_1 \\ \vec{e}_2 \\ \vdots \\ \vec{e}_7 \end{vmatrix}$$

**Example for $i = 1$:**

$$\vec{e}_1 = 0.45 \cdot \vec{e}_1 + 0.25 \cdot \vec{e}_2 + 0.14 \cdot \vec{e}_3 + 0.05 \cdot \vec{e}_4$$
$$+ 0.03 \cdot \vec{e}_5 + 0.06 \cdot \vec{e}_6 + 0.02 \cdot \vec{e}_7$$



The cat sat on the mat.

# Attention



$$Attention(\vec{e}_i, K, V) = softmax(\frac{\vec{e}_i^T \cdot K}{\sqrt{d}}) \cdot V^T$$

$$K = (\vec{e}_1, \vec{e}_2, \dots, \vec{e}_7), \qquad V^T = \begin{pmatrix} \vec{e}_1^T \\ \vec{e}_2^T \\ \vdots \\ \vec{e}_7^T \end{pmatrix}$$

**Example for $i = 2$:**

$$Attention(\vec{e}_2, K, V)$$

$\vec{e}_1$    $\vec{e}_2$    $\vec{e}_3$        $\vec{e}_7$

| 0.7 | 0.2 | 0.1 | | | | -0.1 |
| 0.4 | -0.9 | 0.6 | | | | 0.8 |
| ⋮ | ⋮ | ⋮ | ... ... ... | | | ⋮ |
| 2.1 | 0.6 | -0.3 | | | | -1.6 |

The cat sat on the mat.

1. $\vec{e}_i^T \cdot K = \vec{e}_2^T \cdot K =$

$$= (\vec{e}_2^T \cdot \vec{e}_1, \vec{e}_2^T \cdot \vec{e}_2, \dots, \vec{e}_2^T \cdot \vec{e}_7)$$

$$= (23.2, 70.8, 33.7, 5.7, -12.4, 27.8, -22.4)$$

2. $\dfrac{\vec{e}_i^T \cdot K}{\sqrt{d}}$    $= (0.8, 2.6, 1.2, 0.2, -0.4, 1.0, -0.8)$

3. $softmax(\dfrac{\vec{e}_i^T \cdot K}{\sqrt{d}}) = (0.1, 0.55, 0.14, 0.05, 0.03, 0.12, 0.02)$

4. $softmax(\dfrac{\vec{e}_i^T \cdot K}{\sqrt{d}}) \cdot V^T = (0.1 \cdot \vec{e}_1 + 0.55 \cdot \vec{e}_2 + 0.14 \cdot \vec{e}_3 + 0.05 \cdot \vec{e}_4$
$+ 0.03 \cdot \vec{e}_5 + 0.12 \cdot \vec{e}_6 + 0.02 \cdot \vec{e}_7)^T$

$$Attention(\vec{e}_i, K, V) = softmax(\frac{\vec{e}_i^{\,T} \cdot K}{\sqrt{d}}) \cdot V^T$$

$$Attention(W_E \cdot \vec{e}_i, W_K \cdot K, W_V \cdot V)$$

$$= Attention(\overrightarrow{\tilde{e}_i}, \tilde{K}, \tilde{V})$$

$$\tilde{K} = W_K \cdot K = (WK \cdot \vec{e}_1, W_K \cdot \vec{e}_2, \dots, W_K \cdot \vec{e}_7),$$

$$\tilde{V} = WV \cdot V = (WV \cdot \vec{e}_1, WV \cdot \vec{e}_2, \dots, W_V \cdot \vec{e}_7)$$

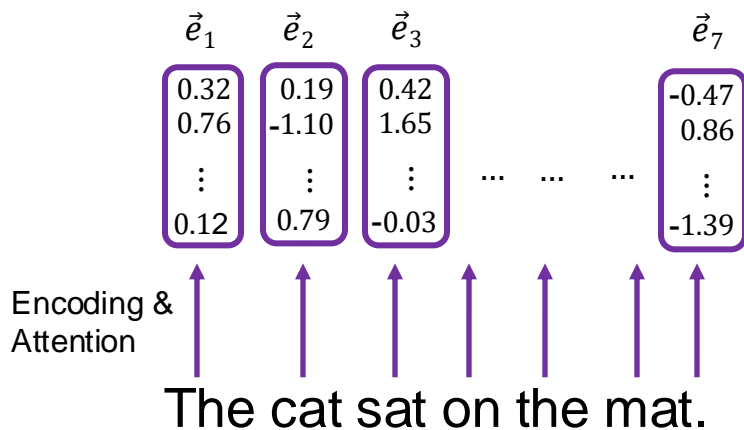$$softmax(\vec{x})_j = \frac{e^{x_j}}{\sum_{k=1}^{n} e^{x_k}}, \quad \vec{x} \in \mathbb{R}^N$$



The cat sat on the mat.

1. Input Embedding & Positional Encoding
2. Attention
3. Feed Forward Neural Network

The cat sat on the mat.

hhu
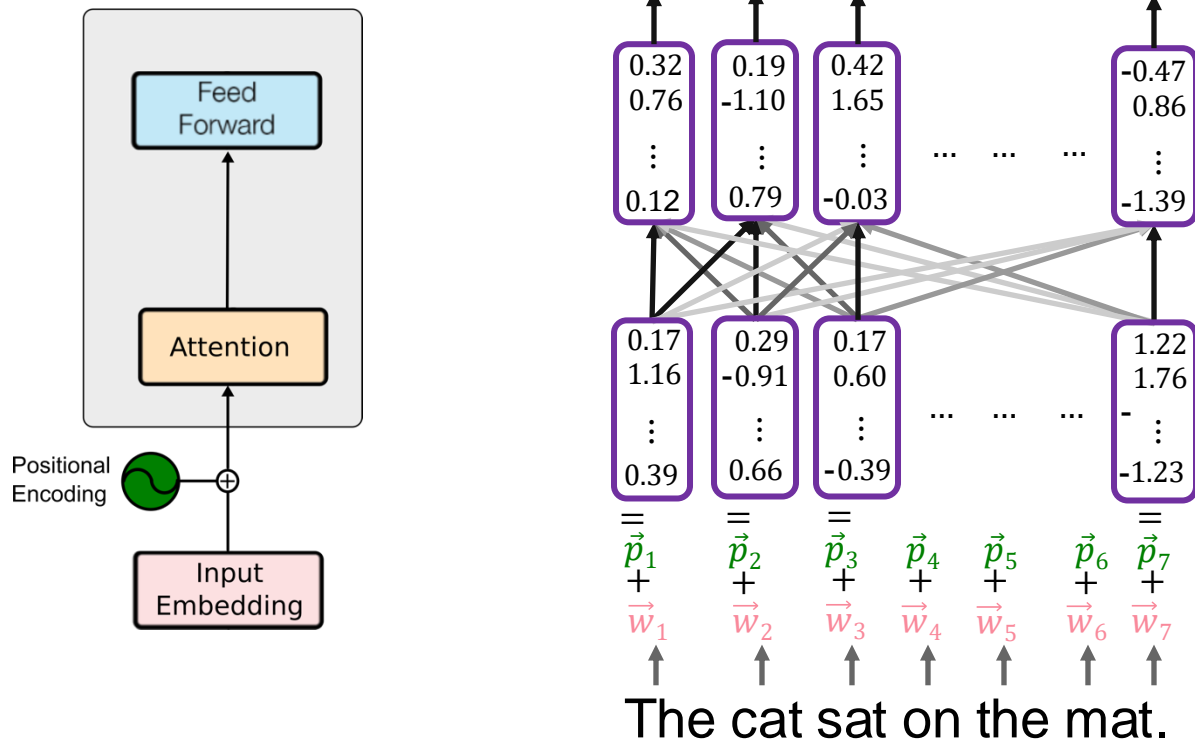Heinrich Heine
Universität Düsseldorf

- Mapping through Encoding and Attention:

$\vec{e}_1$  $\vec{e}_2$  $\vec{e}_3$   $\vec{e}_7$

$\vec{e}_1$: 0.32, 0.76, ⋮, 0.12
$\vec{e}_2$: 0.19, -1.10, ⋮, 0.79
$\vec{e}_3$: 0.42, 1.65, ⋮, -0.03
… … …
$\vec{e}_7$: -0.47, 0.86, ⋮, -1.39

Encoding & Attention

The cat sat on the mat.

- Feed Forward Neural Network:

Feed
Forward

$$FFN(\vec{e}_i) = W_2 \cdot \mathrm{ReLU}(W_1 \cdot \vec{e}_i + \vec{b_1}) + \vec{b_2}$$

i=1: $\vec{e}_i \in \mathbb{R}^d$      $FFN(\vec{e}_i) \in \mathbb{R}^d$

0.32          2.71

0.76          0.18

0.12          1.43

# Applying Encoders to Protein Sequences



| A | L | K | ... | A |

Elementwise mean