# Attention Mechanism and Transformers

**Video 7:** Introduction to Transformers

Analytics Vidhya

IN AIR

ChatGP**T**

# Transformers

**Purpose:** Addresses sequence-to-sequence issues with long-range dependencies.

**Impact:** Integral to NLP, driving advancements like BERT, GPT 2, T5, Chat GPT.

# Recap

**Encoder - Decoder**



Output

ENCODER → Encoder state → DECODER

Input

**Encoder state :** solely responsible to transfer information to decoder.

# Recap

**Encoder - Decoder**

Output

Encoder state

**ENCODER** $\cdots\cdots$ **DECODER**

Input

**Encoder state :** solely responsible to transfer information to decoder.

**Attention Mechanism**

$i = 2$

Wie | geht

$s_0$ | $S_1$

$c_1$ | $C_2$

<START> | Wie

$C_2$

$\alpha_{21}$ | $\alpha_{22}$ | $\alpha_{23}$

$h_0$ | $h_1$ | $h_2$ | $h_3$

How | are | you

$j = 1$ | $j = 2$ | $j = 3$

$$C_i = \sum \alpha_{ij} h_i$$

Calculation of alpha **increase computation** for long documents.

Encoder processes the **input sequentially**

# Architecture



Attention

- Encoder-Decoder Attention
- Self Attention
- Masked Attention
- Multi-Head Attention
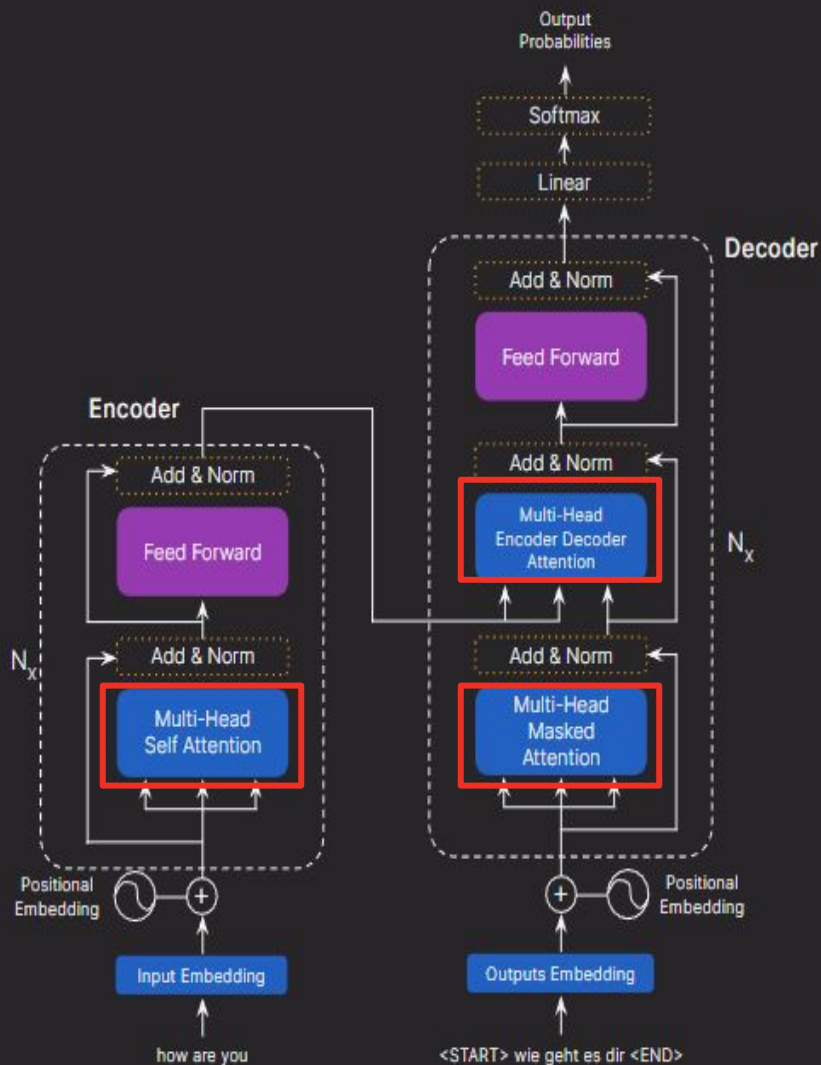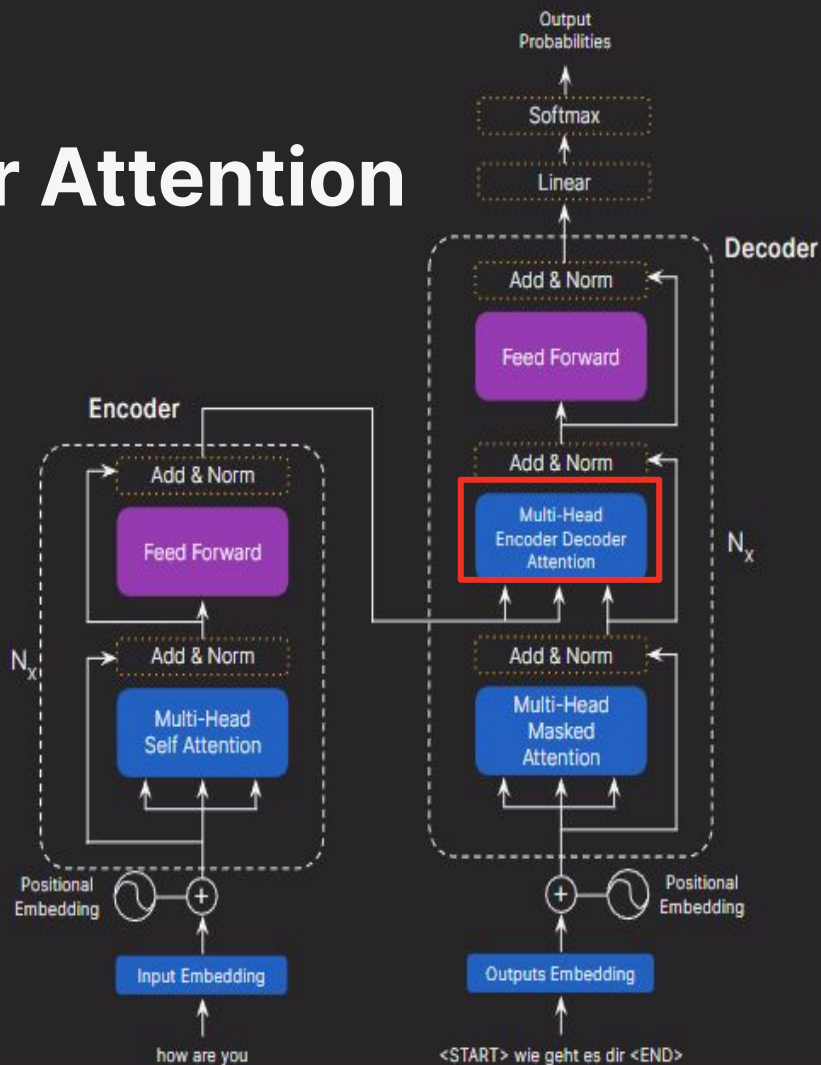
Encoder - Decoder Attention

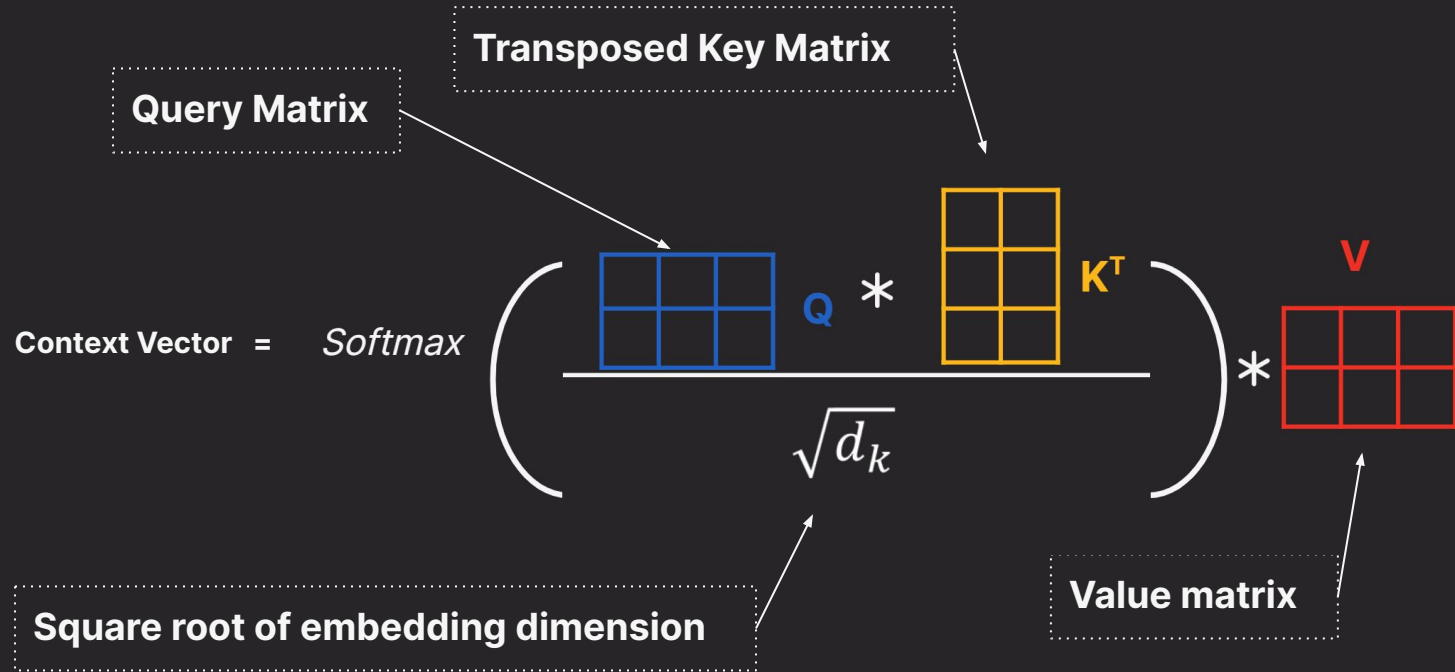# Encoder- Decoder Attention

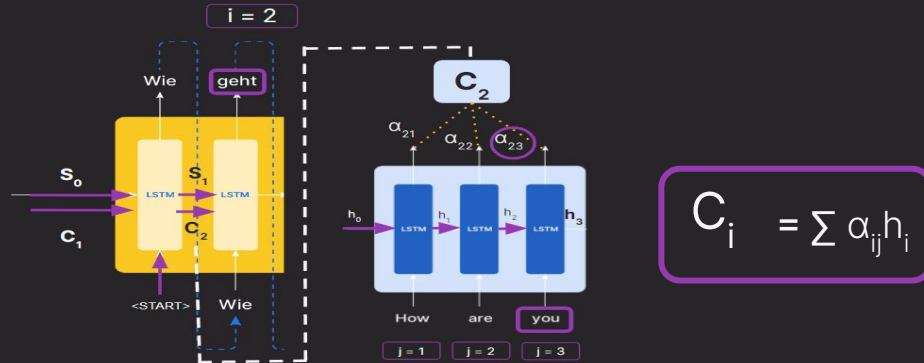$$C_i = \sum \alpha_{ij} h_i$$

Uses Dot-Product attention to calculate attention.

$$\text{Context Vector} = \text{Softmax} \left( \frac{Q \cdot K^T}{\sqrt{d^k}} \right) * V$$

# Dot Product Attention



Transposed Key Matrix

Query Matrix

Context Vector $=$ $Softmax$ $\left( \dfrac{Q * K^T}{\sqrt{d_k}} \right) * V$

Square root of embedding dimension

Value matrix

# Dot Product Attention



$$C_i = \sum \alpha_{ij} h_i$$

Context Vector $= Softmax \left( \dfrac{Q * K^T}{\sqrt{d_k}} \right) * V = Z$

Attention weights ........ $\sum \alpha_{ij}$

$h_j$ ..... Hidden States

# Dot Product Attention

**Q**

**Query Vector** : Importance of Relative Words

**K**

**Key Vector** : Represents Evaluated Word

**V**

**Value Vector** : Contains focussed information
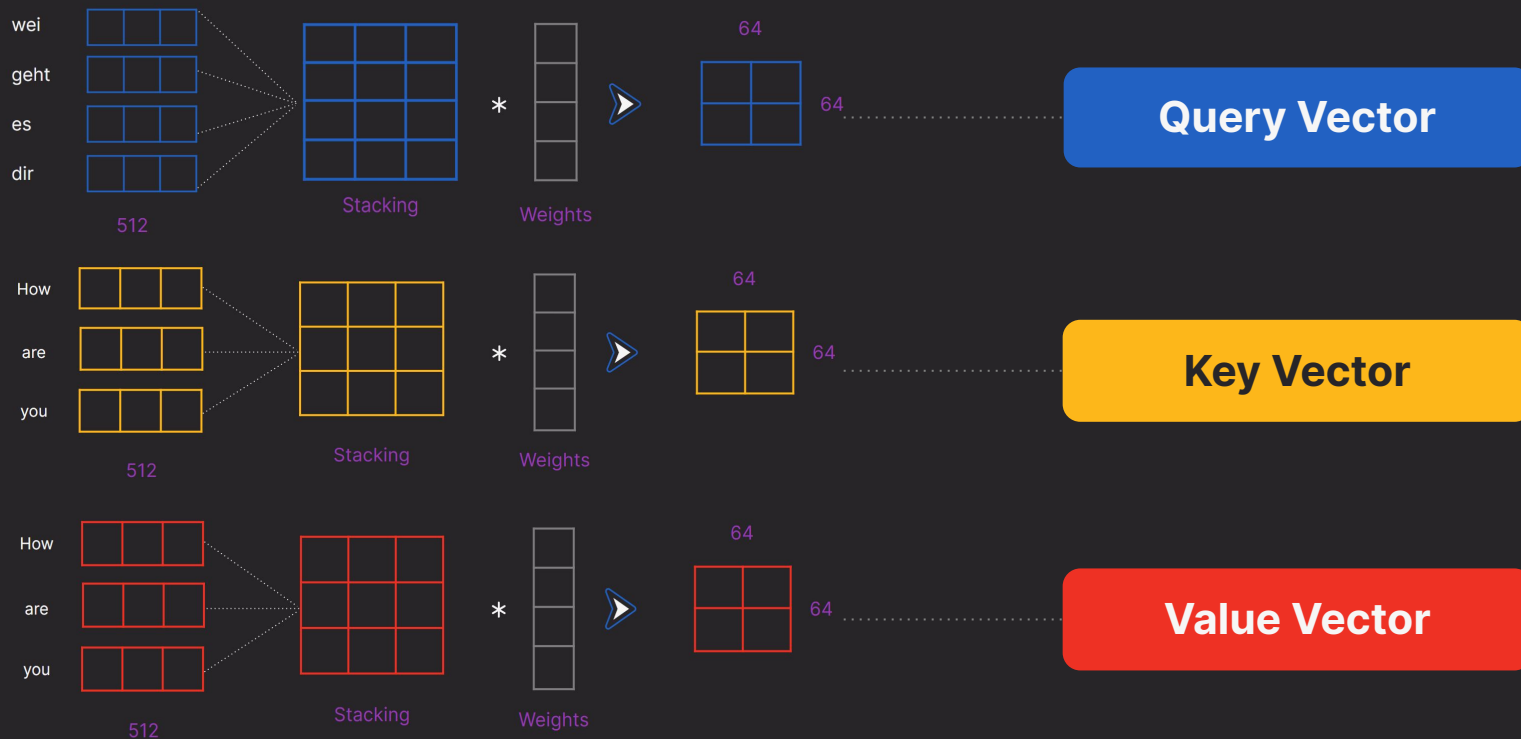
# Encoder-Decoder Dot Product Attention



wei
geht
es
dir

512

Stacking

* Weights

64

64

**Query Vector**

How
are
you

512

Stacking

* Weights

64

64

**Key Vector**

How
are
you

512

Stacking

* Weights

64

64

**Value Vector**

# Encoder-Decoder Dot Product Attention

$$Cv = Softmax \left( \frac{Q . K^T}{\sqrt{d^k}} \right) * V$$

$$\boxed{\text{Context Vector}} = Softmax \left( \frac{Q \; * \; K^T}{\sqrt{d_k}} \right) * V$$
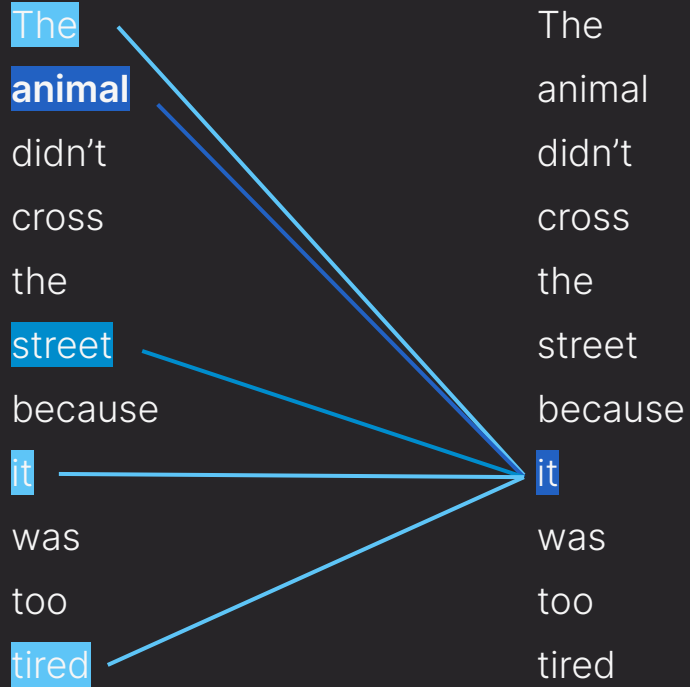
Z

# Encoder-Decoder Dot Product Attention

# Self Attention

- The animal didn't cross the street because **it** was too tired.

- The animal didn't cross the street because **it** was too crowded.

What does "it" refer to in these sentences?

# Self Attention



The animal didn't cross the street because **it** was too tired.

# Self Attention



The animal didn't cross the street because **it** was too tired.

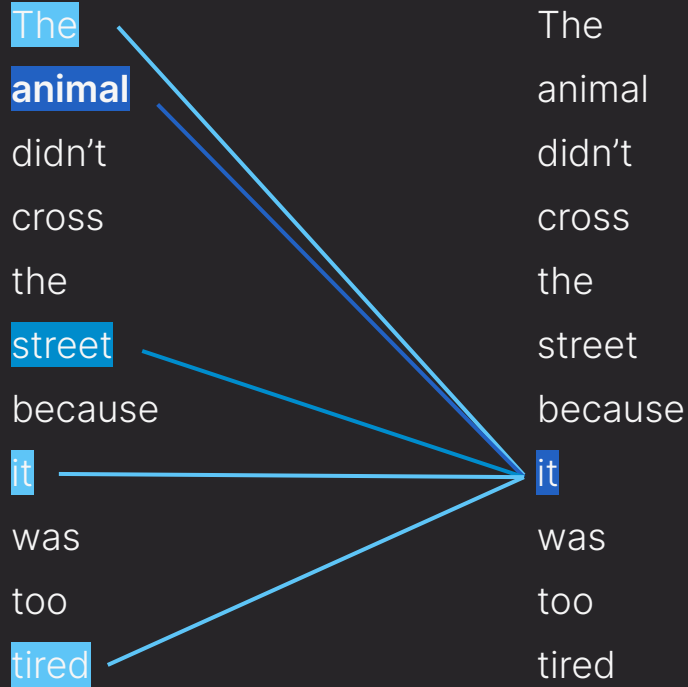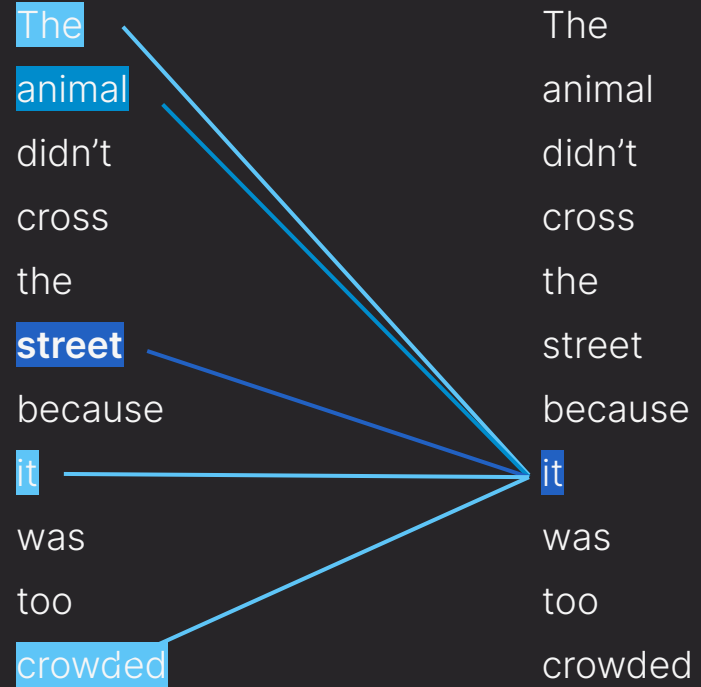The animal didn't cross the street because **it** was too wide.

# Self Attention

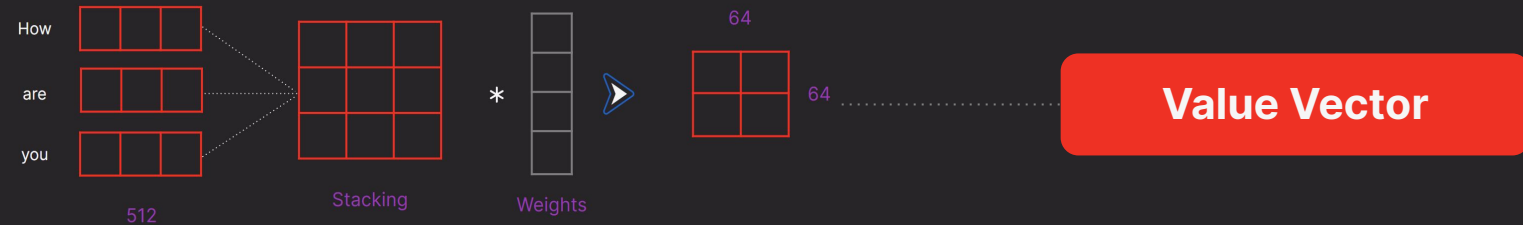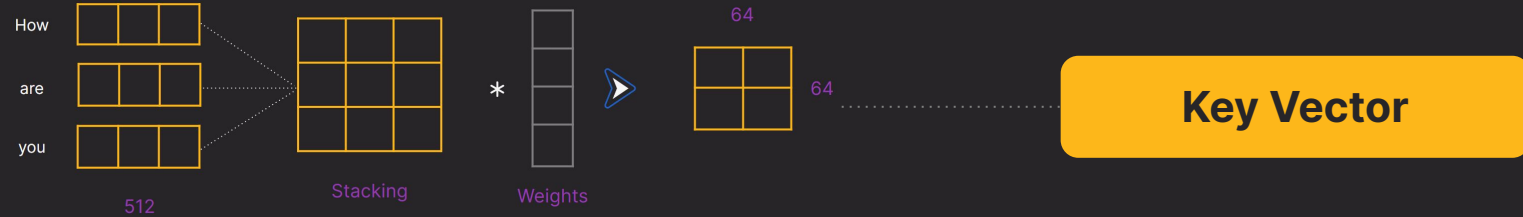$$C_v = \frac{\text{Softmax}\ (\ Q\ .\ K^T\ )* V}{\sqrt{d^k}}$$

Q = Query Vector
K = Key vector
V = Value vector
$d^k$ = dimension of key vector

# Self Attention



How are you — 512 — Stacking — Weights — * — 64 × 64 — **Query Vector**

How are you — 512 — Stacking — Weights — * — 64 × 64 — **Key Vector**

How are you — 512 — Stacking — Weights — * — 64 × 64 — **Value Vector**

# Self Attention

|       | How | are | you |
|-------|-----|-----|-----|
| How   |     |     |     |
| are   |     |     |     |
| you   |     |     |     |

# Architecture



Attention
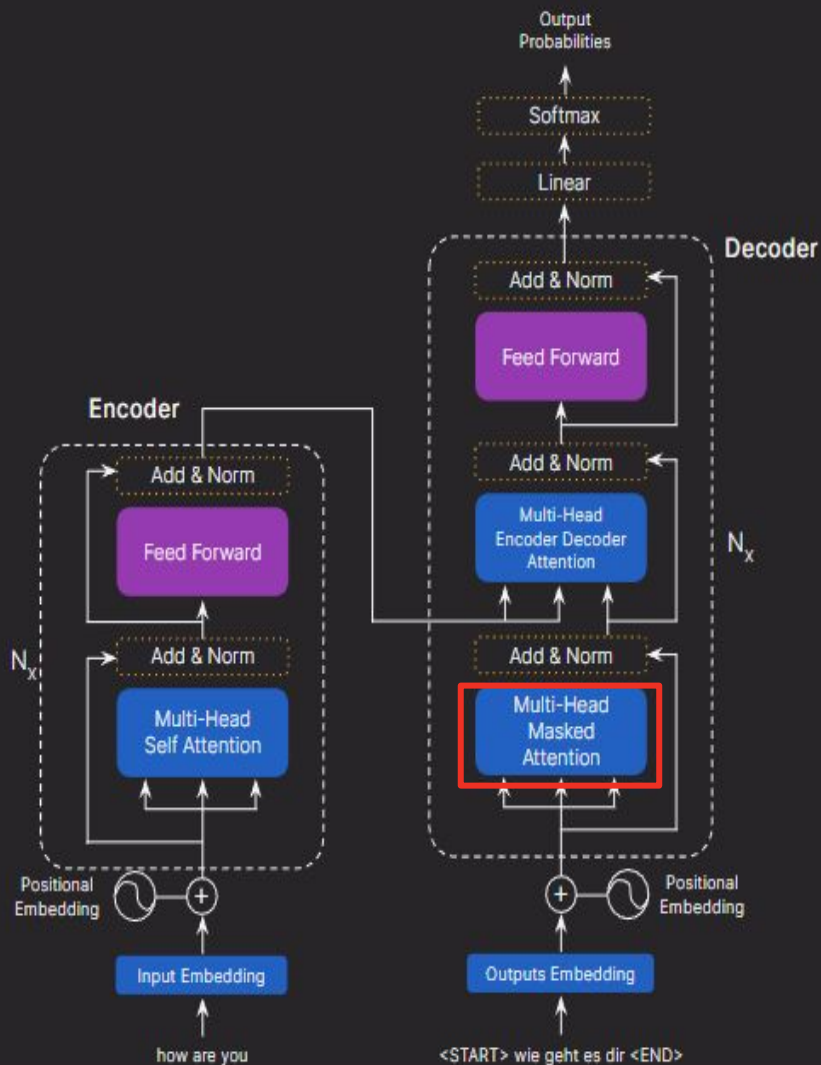
- Encoder-Decoder Attention
- Self Attention
- **Masked Attention**
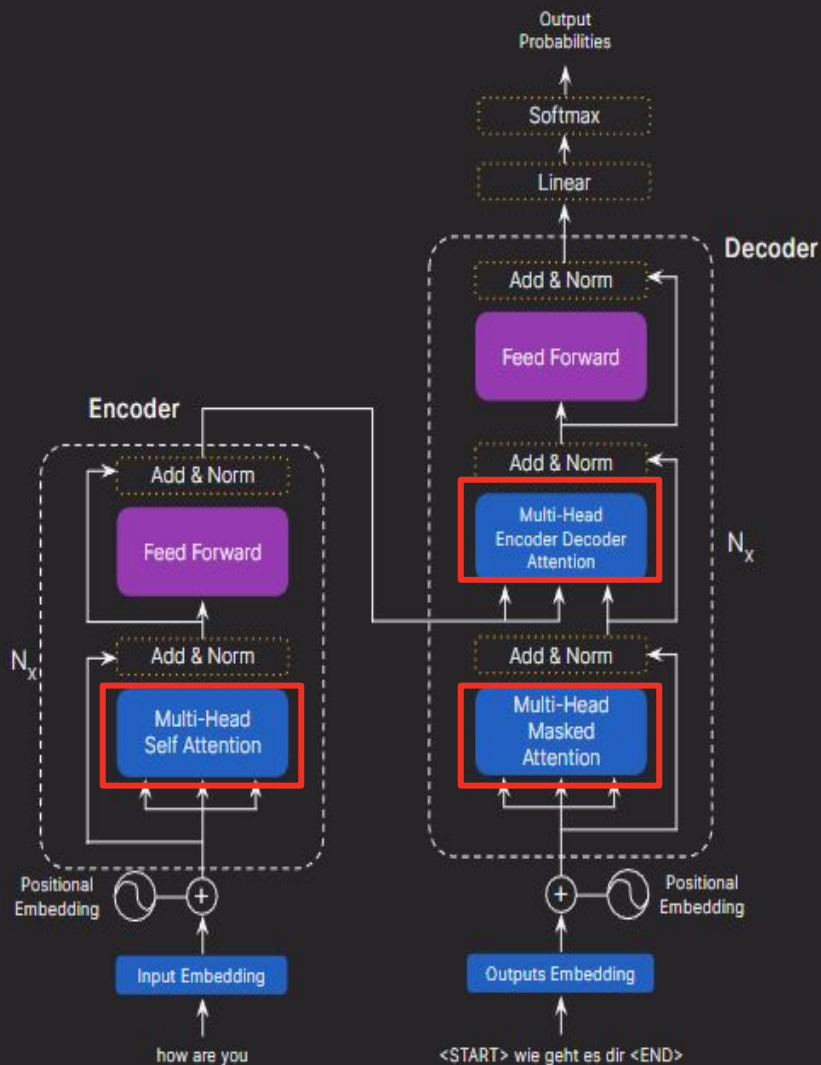- Multi-Head Attention

# Masked Attention
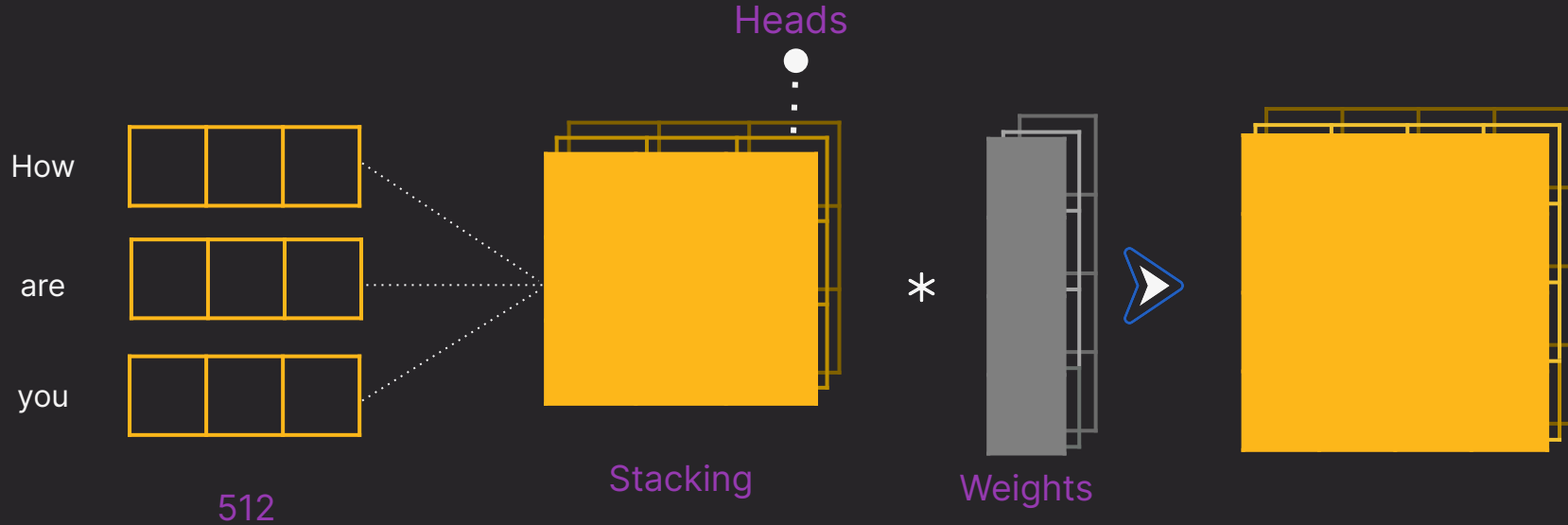
# Architecture



**Attention**

- Encoder-Decoder Attention
- Self Attention
- Masked Attention
- **Multi-Head Attention**

Output Probabilities

Softmax

Linear

Decoder

Add & Norm

Feed Forward

Add & Norm

Multi-Head Encoder Decoder Attention

$N_x$

Encoder

Add & Norm

Feed Forward

Add & Norm

Multi-Head Self Attention

$N_x$

Add & Norm

Multi-Head Masked Attention

Positional Embedding

Input Embedding

how are you

Positional Embedding

Outputs Embedding

<START> wie geht es dir <END>

# Architecture



Attention

- Encoder-Decoder Attention
- Self Attention
- Masked Attention
- Multi-Head Attention