

Transformers Explained

→ Self - Attention
Mechanism

By: Kanav Bansal
(That AI Guy)

Summary ↓

2014 → Sequence to Sequence Learning with Neural Network
Proposed - Encoder Decoder architecture

2015 → Neural Machine Translation with Joint Learning to align & Translate
Proposed - Bahdanau Attention Mechanism

2017 → Attention is all you need
Proposed - Transformer Architecture with Self-Attention

2018 → Universal Language Model Fine-Tuning for Text Classification
Proposed - Language Modeling as pre-training Task

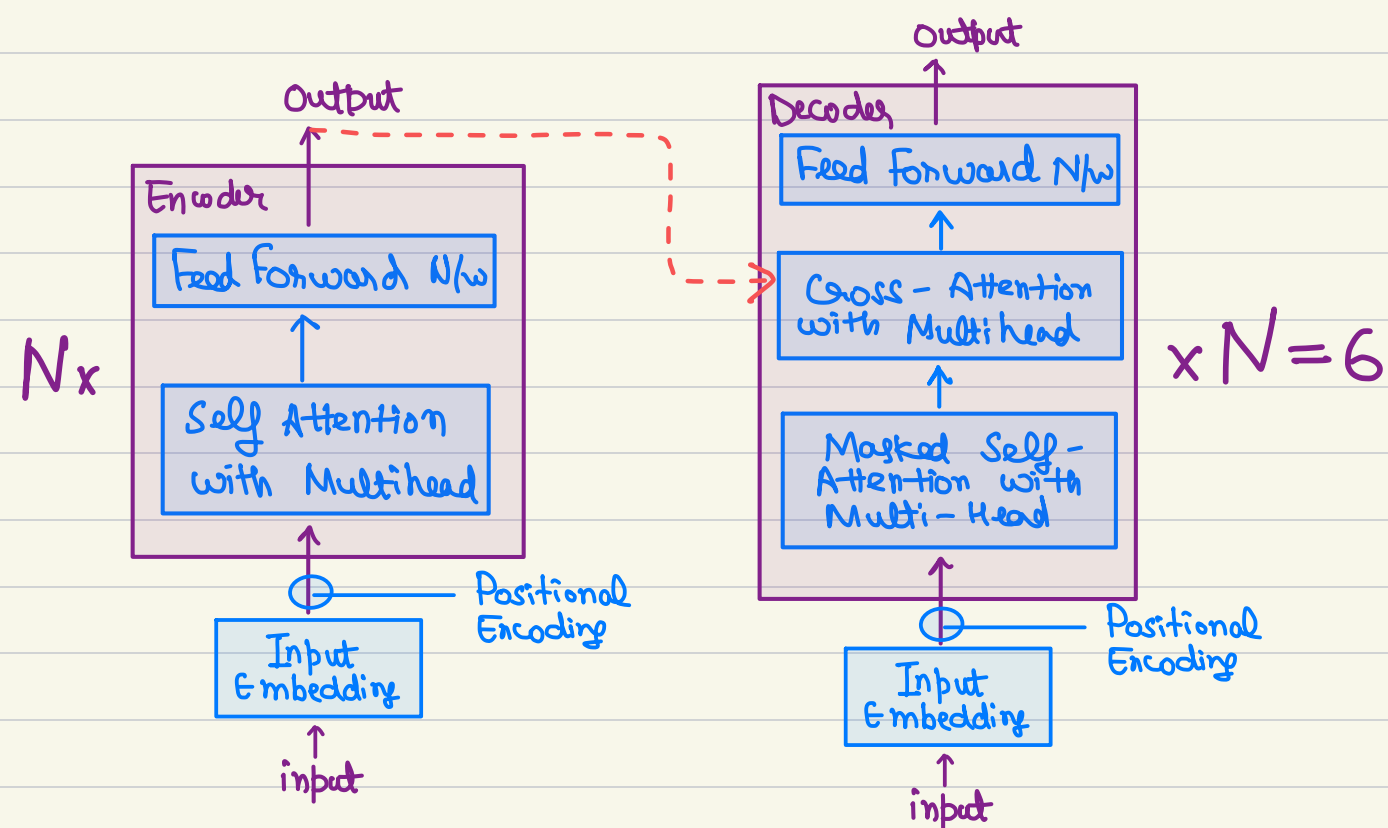
Paper: Attention is all you need

2017 → Proposed Transformers Architecture
which used:

- Encoder Decoder Architecture
- Self-Attention Mechanism

BUILDING BLOCKS

(Note: This is a simplified view of Transformers)



Encoder is good
at understanding
text

Decoder is good
at generating
text

ENCODER BUILDING BLOCK

1. Positional Embedding
 - Add the notion of Seq info
2. Self-Attention with Multihead
 - Attention replaced need of recurrence & makes the process of embedding calculation parallel
 - Attention mechanism will find how each word relates to all other words in a sequence & generates Contextual Embedding

Ques: How to find word similarities?

- Attention will run dot products b/w word vectors & determine the strongest relationships of a word with all other words

Ques: How to Speed up the calculations?

- For each attention sublayer, the original transformer model runs not only one but eight attention mechanisms in parallel to speed up the calculations. This process is done using 'multi-head attention'.

3. Fully Connected Positionwise Feed Forward N/w
 - Improves the word association by applying non-linear transformations.

TRANSFORMER'S



Attention is all you need
(2017 Paper)

Why CELEB STATUS ?

1. Scalable & Parallel Training
2. Revolutionized NLP with LLMs
3. Unified DL Approaches for text, image, audio & video data
4. Multi-Modal
5. Accelerated Gen AI

Encoder

Has self-Attention
with multi-head
&

Feed Forward ANN



Great at
Understanding
Text

Eg: BERT

Decoder

Mask self-Attention
& Cross-Attention
with Multi-head



Feed Forward ANN



Great at
Generating
Text

Eg: GPT-1

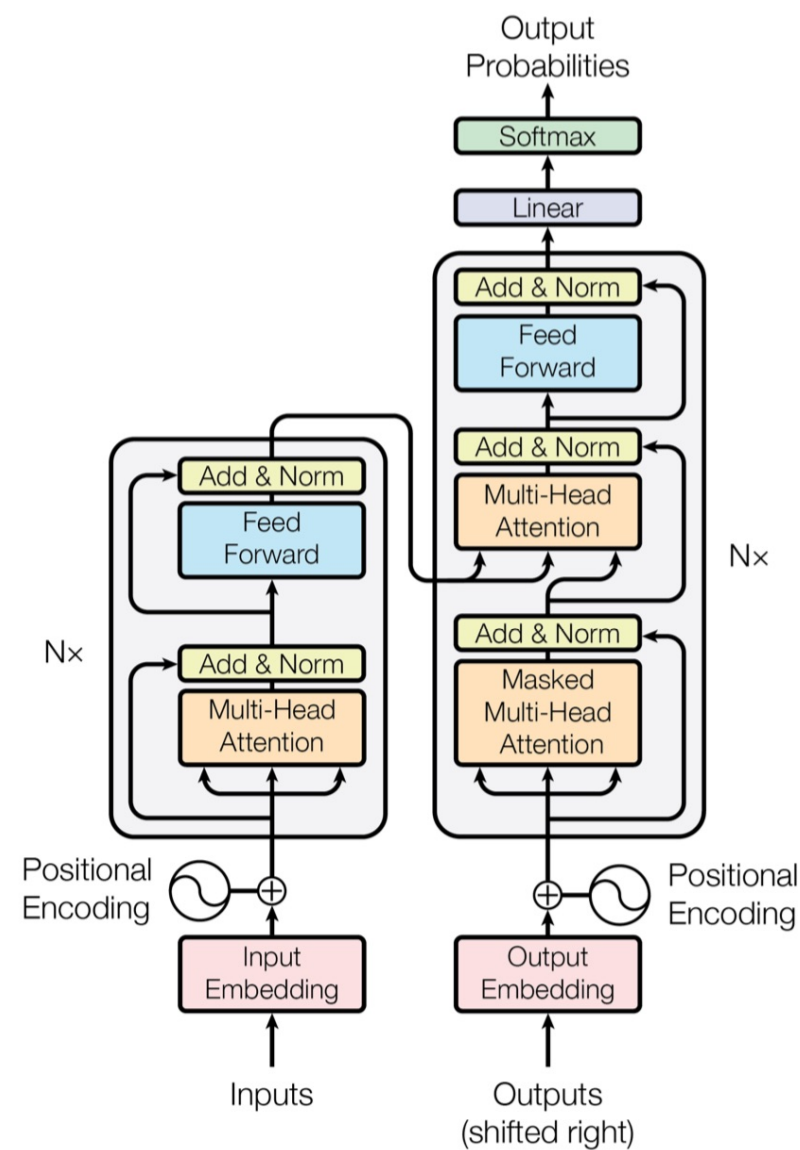


Figure 1: The Transformer - model architecture.

Advantages ↓

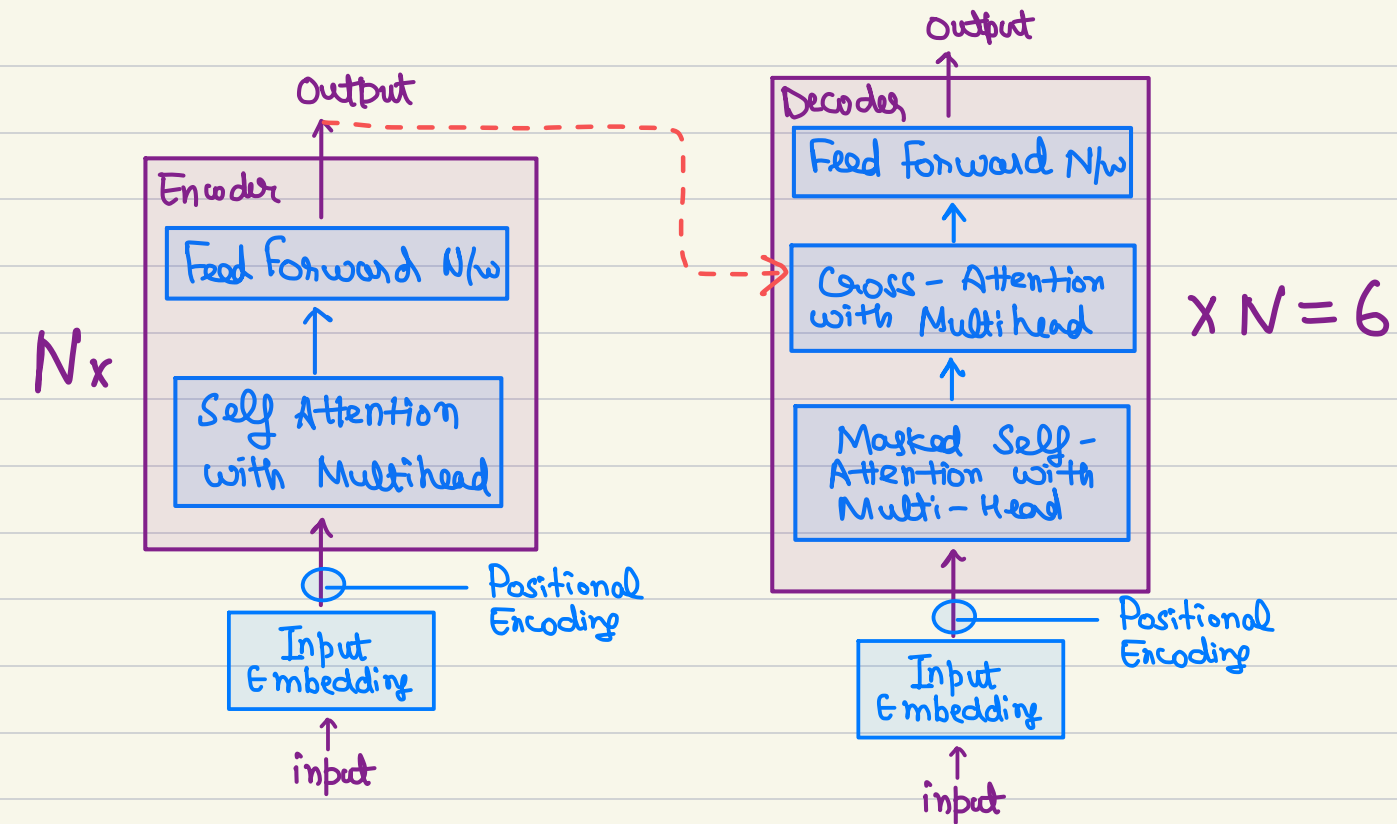
- Parallel & Scalable
- Multi Modal input & output

Disadvantages ↓

- Huge compute resources
- Huge amount of data
- Overfit
- Energy

BUILDING BLOCKS ↴

By: Karan Bansal
(That AI Guy)

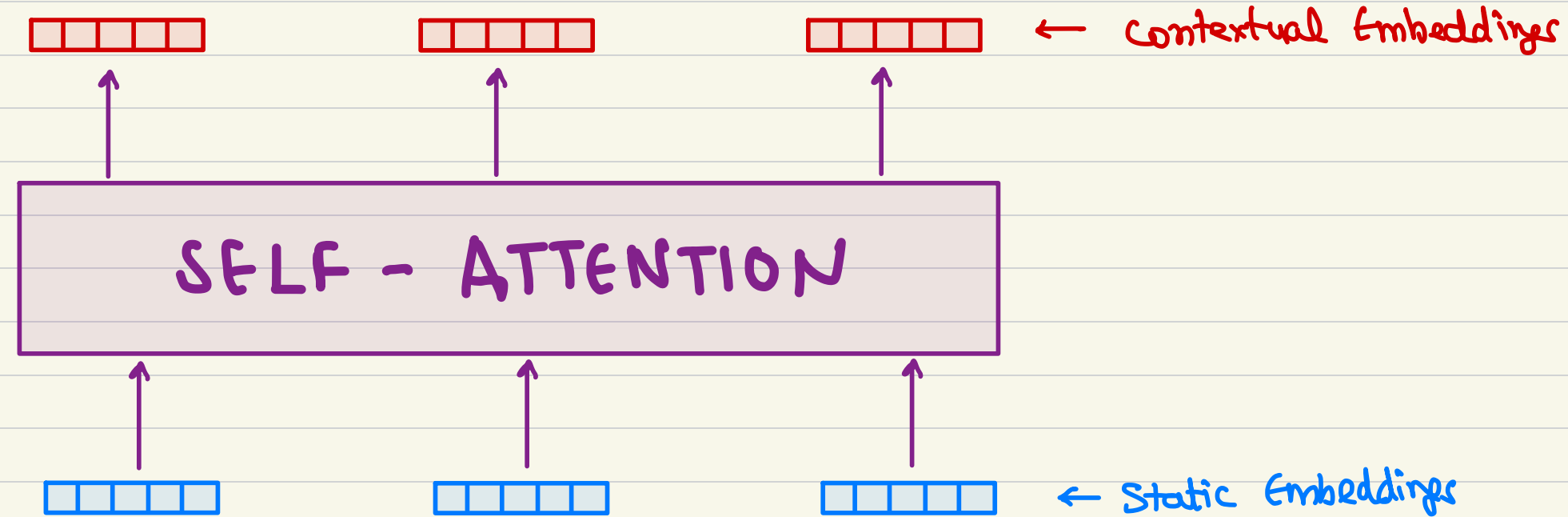


Encoder is good
at understanding
text

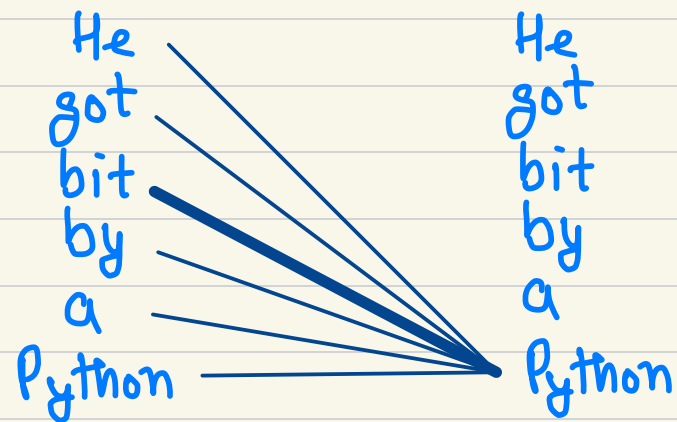
Decoder is good
at generating
text

UNDERSTANDING SELF-ATTENTION ↴

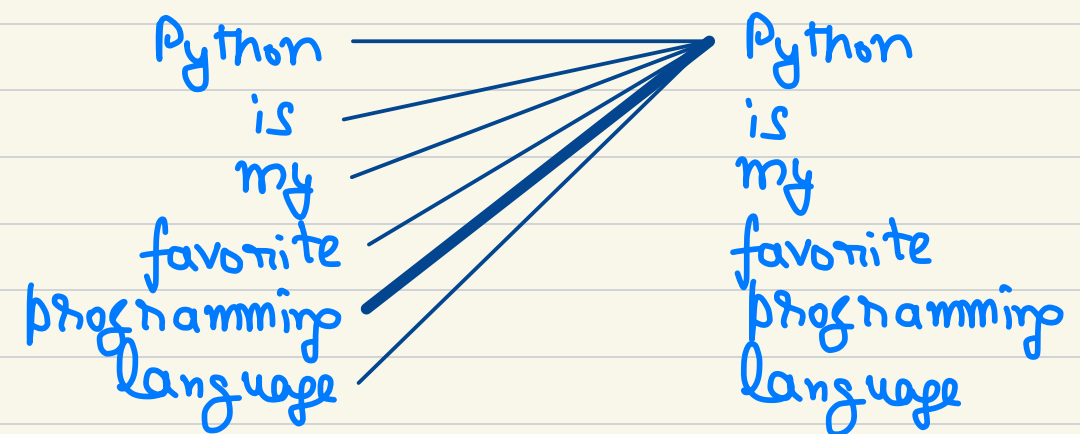
By: Kanav Bansal
(That AI Guy)



Ex: Sent A → He got bit by a Python



Sent B → Python is my favorite programming language



Example ↴

Sentence 1: Python executes quickly

$$\text{Python} = 0.5 \text{ Python} + 0.3 \text{ executes} + 0.2 \text{ quickly}$$

$$\text{executes} = 0.2 \text{ Python} + 0.7 \text{ executes} + 0.1 \text{ quickly}$$

$$\text{quickly} = 0.1 \text{ Python} + 0.1 \text{ executes} + 0.8 \text{ quickly}$$

$$\begin{aligned} Z_{\text{Python}} &= \boxed{0.5} e_{\text{Python}} + \boxed{0.3} e_{\text{executes}} + \boxed{0.2} e_{\text{quickly}} \\ &\quad \begin{matrix} e_{\text{Python}} \cdot e_{\text{Python}} & e_{\text{Python}} \cdot e_{\text{executes}} & e_{\text{Python}} \cdot e_{\text{quickly}} \\ \uparrow & \uparrow & \uparrow \end{matrix} \\ Z_{\text{executes}} &= 0.2 e_{\text{Python}} + 0.7 e_{\text{executes}} + 0.1 e_{\text{quickly}} \\ Z_{\text{quickly}} &= \boxed{0.1} e_{\text{Python}} + 0.1 e_{\text{executes}} + 0.8 e_{\text{quickly}} \\ &\quad \begin{matrix} e_{\text{quickly}} \cdot e_{\text{Python}} \\ \downarrow \end{matrix} \end{aligned}$$

Sentence 2: Python crawls quickly

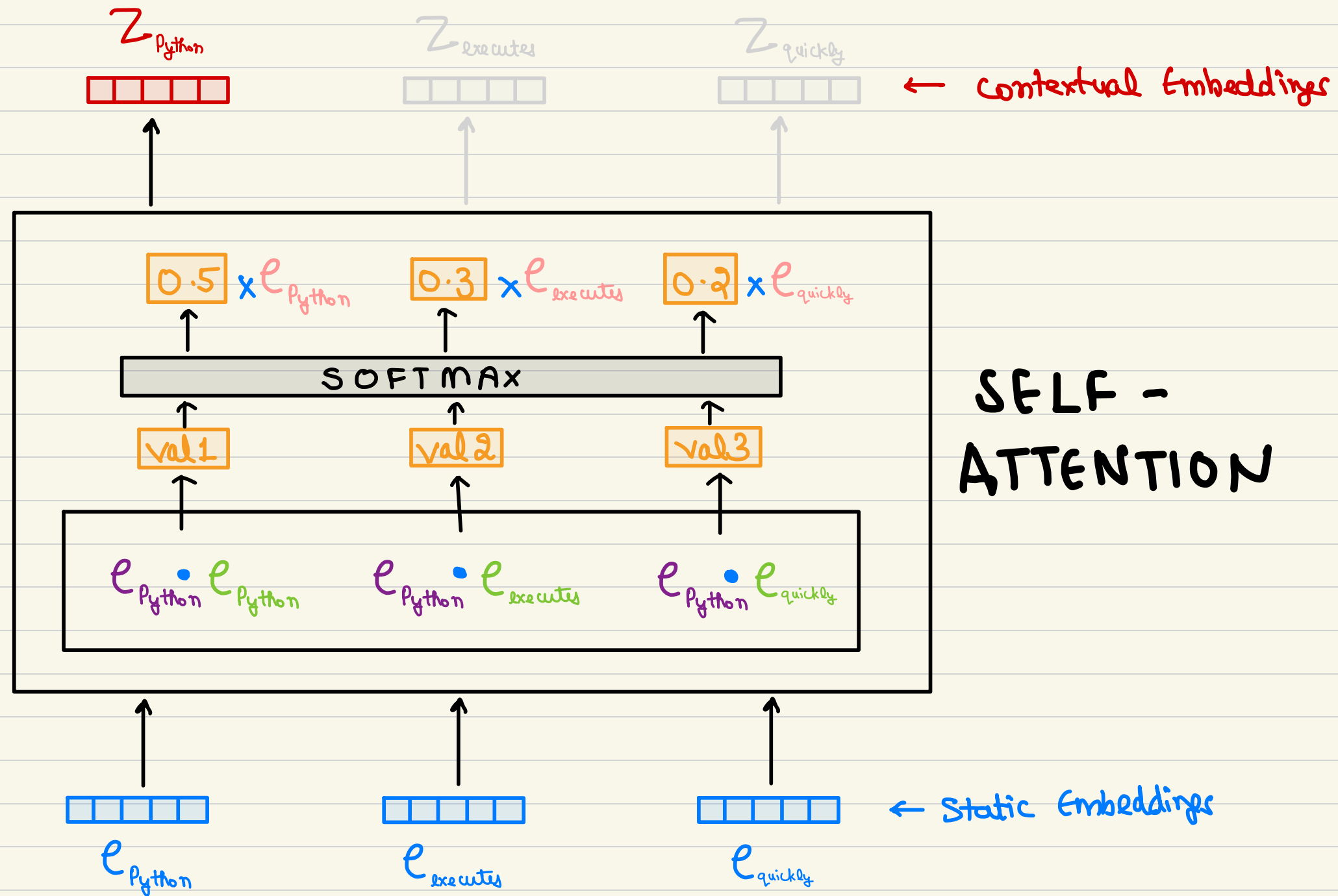
$$\text{Python} = 0.4 \text{ Python} + 0.4 \text{ crawls} + 0.2 \text{ quickly}$$

$$\text{crawls} = 0.3 \text{ Python} + 0.6 \text{ crawls} + 0.1 \text{ quickly}$$

$$\text{quickly} = 0.1 \text{ Python} + 0.1 \text{ crawls} + 0.8 \text{ quickly}$$

$$\begin{aligned} Z_{\text{Python}} &= 0.4 e_{\text{Python}} + \boxed{0.4} e_{\text{crawls}} + 0.2 e_{\text{quickly}} \\ &\quad \begin{matrix} e_{\text{Python}} \cdot e_{\text{crawls}} \\ \uparrow \end{matrix} \\ Z_{\text{crawls}} &= 0.3 e_{\text{Python}} + 0.6 e_{\text{crawls}} + 0.1 e_{\text{quickly}} \\ Z_{\text{quickly}} &= 0.1 e_{\text{Python}} + 0.1 e_{\text{crawls}} + 0.8 e_{\text{quickly}} \end{aligned}$$

The complete process \rightarrow (Calculation for the 'Python' Query)



SELF ATTENTION ↴

By: Kanav Bansal
(That AI Guy)

1. init the 3 weight matrix

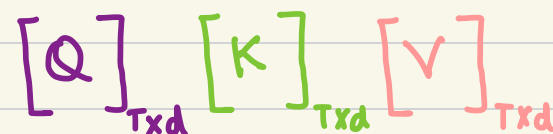
These weights get updated at training time



512
↓
Shape = D x d
↑
64

2. Calc 1:

$$Q = X W_q \quad K = X W_k \quad V = X W_v$$



$$\text{Shape} = (T \times D) \times (D \times d) = T \times d$$

3. Calc 2:

$$\text{Softmax}\left(\frac{Q K^T}{\sqrt{d}}\right) V$$

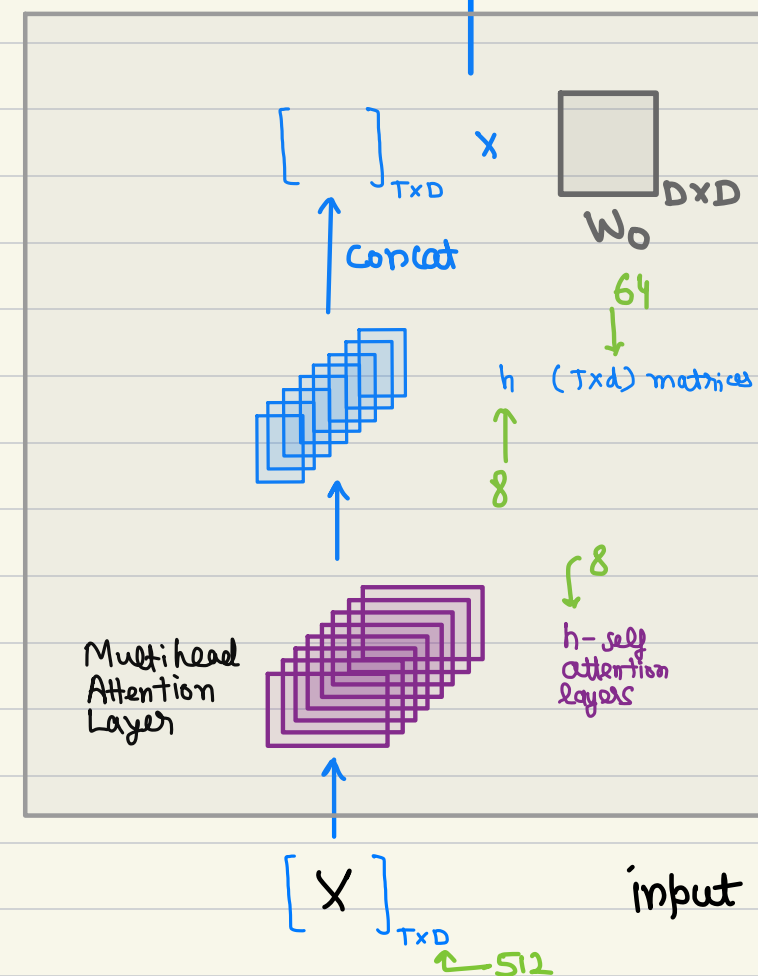
This returns a matrix of shape $T \times T$.
Intuitively, we can think of this matrix as the importance of each word on another.
AKA Attention Score

$$\begin{aligned} \text{Shape} &= ((T \times d) \times (d \times T)) \times (T \times d) \\ &= (T \times T) \times (T \times d) \\ &= (T \times d) \end{aligned}$$

MULTI-HEAD SELF-ATTENTION ↴

As it uses multiple layers of self-attention in parallel

output $[Z]_{T \times D}$

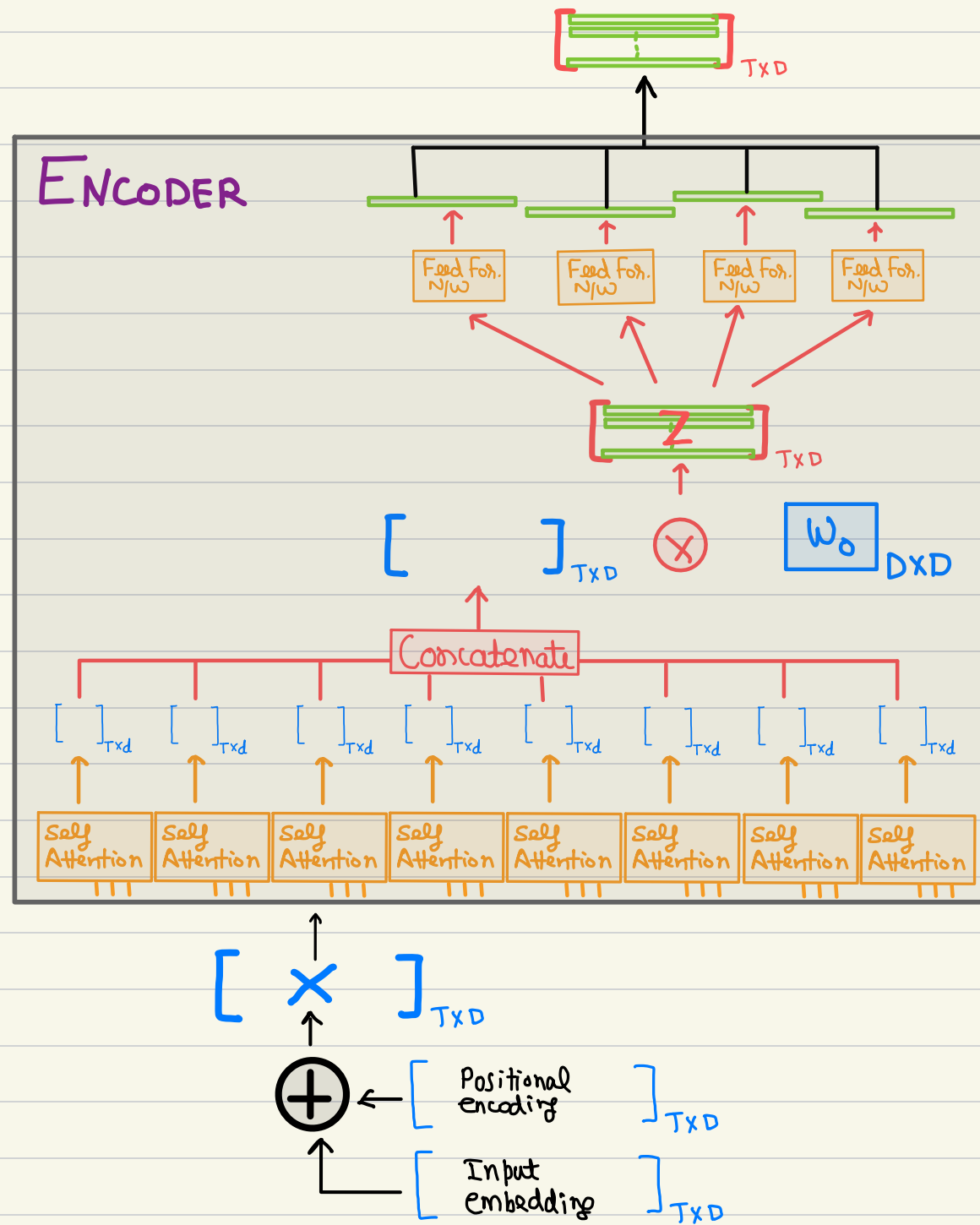


Concat



ENCODER ↴

By: Karan Bansal
(That AI Guy)



$$PE_{(pos,2i)} = \sin(pos/10000^{2i/D})$$

$$PE_{(pos,2i+1)} = \cos(pos/10000^{2i/D})$$

