

# Attention Mechanism and Transformer model

---


# Introduction

NLP requires efficient sequence modeling for building state-of-the-art language model

RNN was tried in the beginning but it suffered from vanishing and exploding gradient

LSTM and GRU was tried later but that too could not find relationship between words which are quite distant apart

e.g.: [Rama](#), being the eldest son of Dasaratha and also being the conqueror of Ravana's lanka was undoubtedly the best candidate for being the king and had actually become the [king of Ayodhya](#)



# Attention is all you need...

The famous paper by Ashish Vaswani et al. (2017)

Language model should weigh different words in a sequence differently to understand the meaning properly

Attention model was built on the basis of the above philosophy to give attention to important words in a sequence

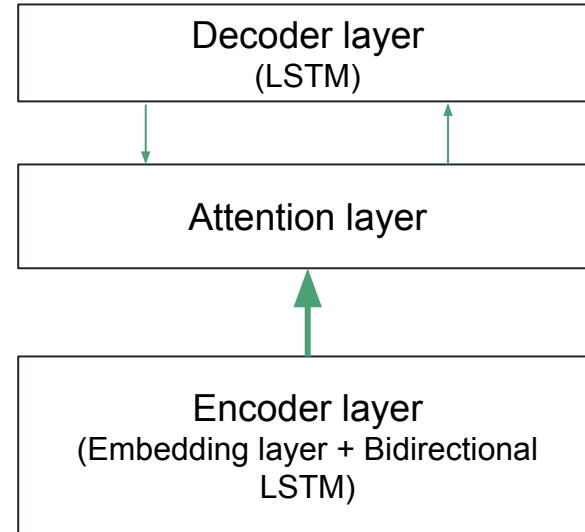
It is a neural network model that learns how to give attention to different words in a sequence

Attention model surpassed solely LSTM or solely GRU based models by large margin to claim the state-of-the-art-model recognition

# The architecture

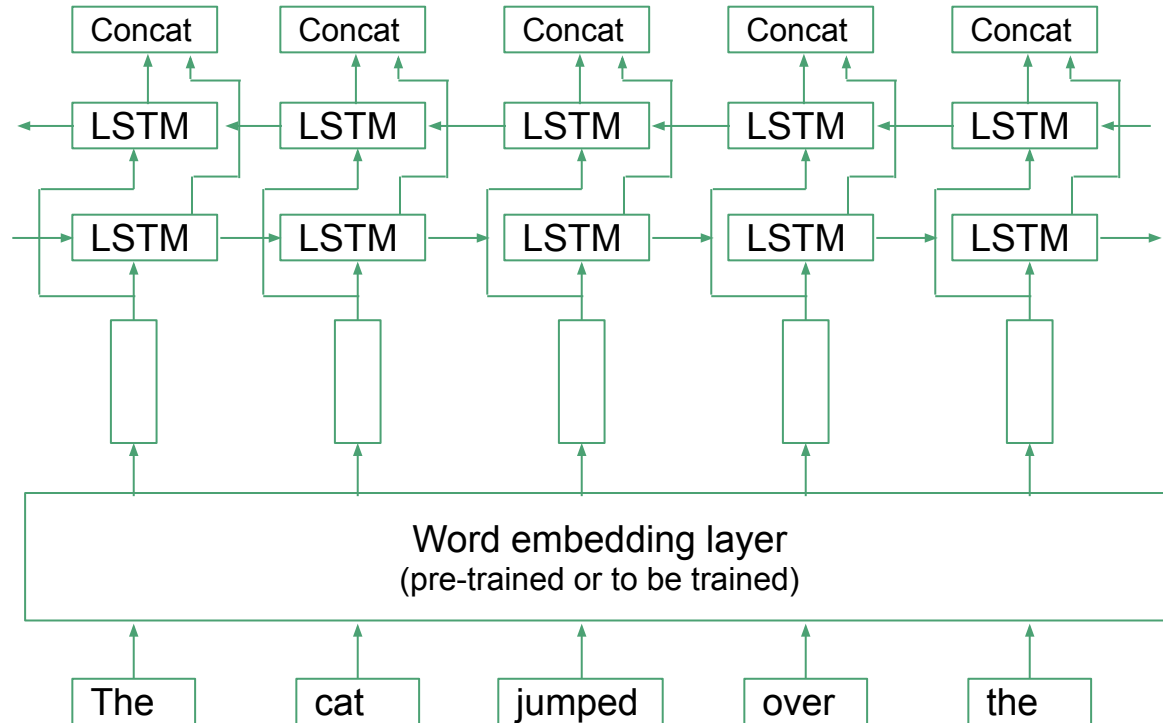
Attention model comprises of 3 different layers

- The encoding layer
- The attention layer
- The decoder layer

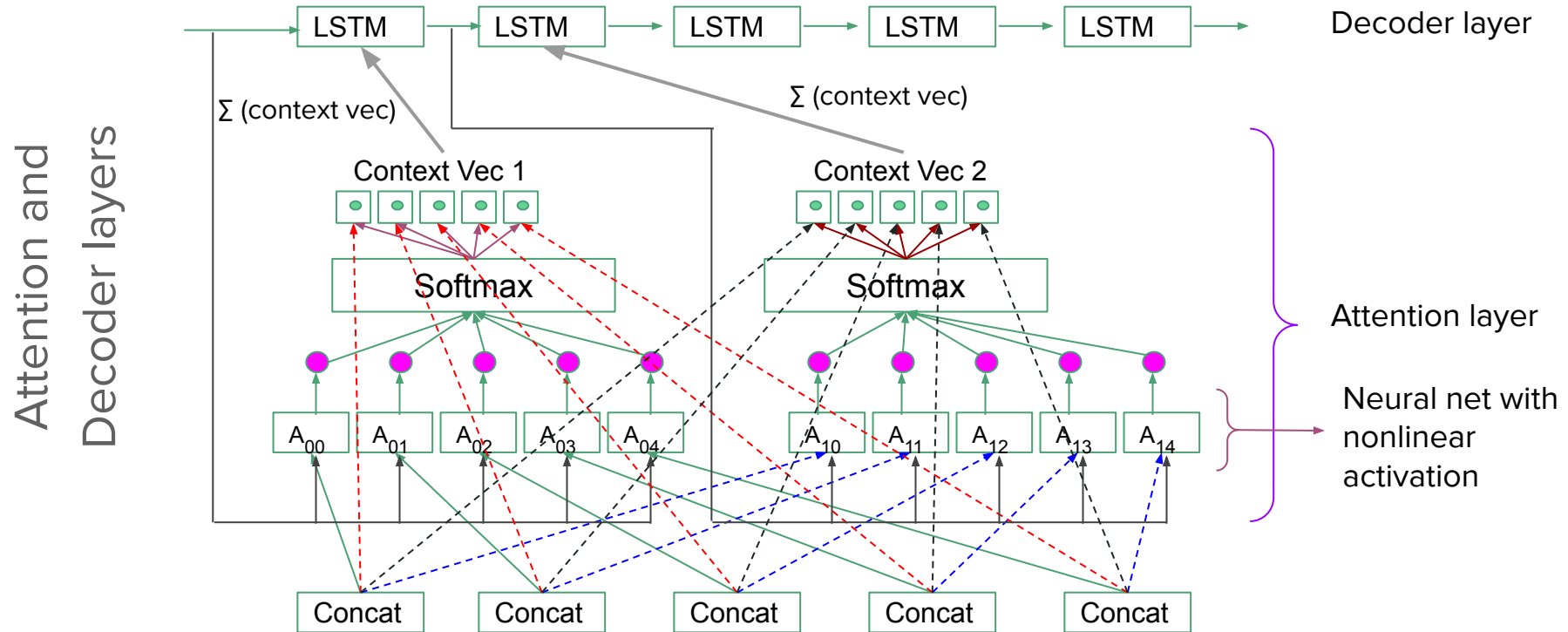


# Getting Deeper...

Encoder layer

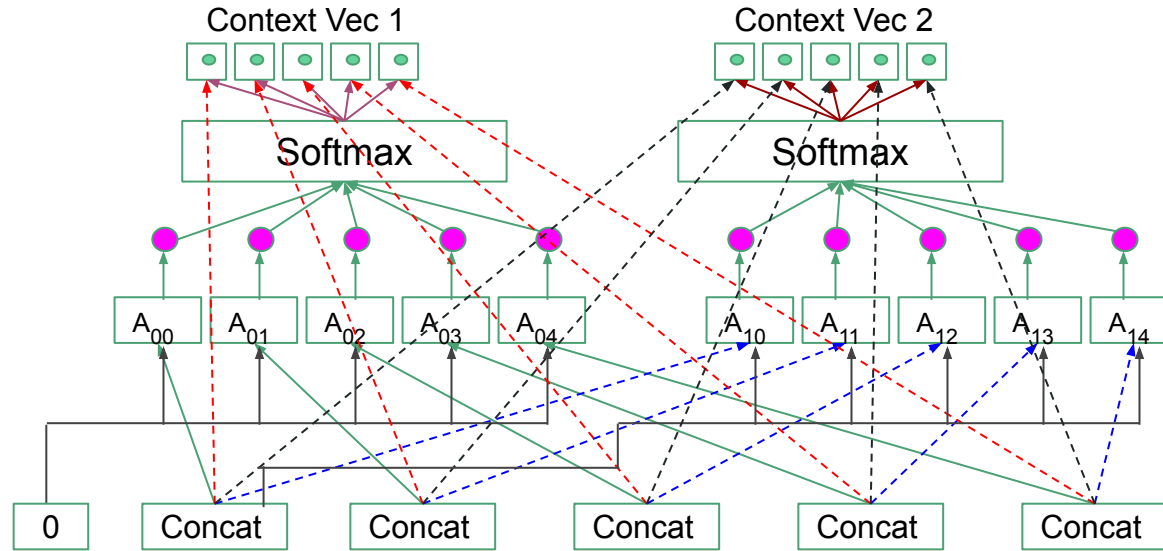


# Getting Deeper...

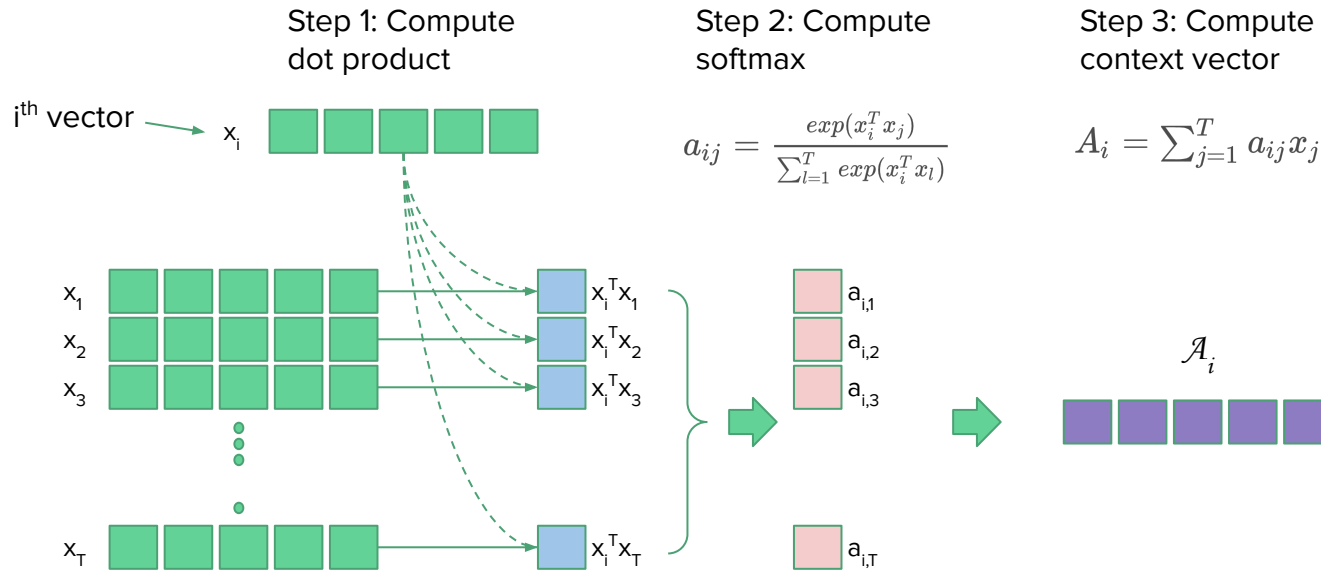


# Self Attention

No decoding layer



# Attention without RNN



Basic attention model without trainable parameters (not good..!)



# Attention without RNN (contd...)

The basic attention model without trainable parameters can be made powerful trainable attention model by incorporating three trainable matrices

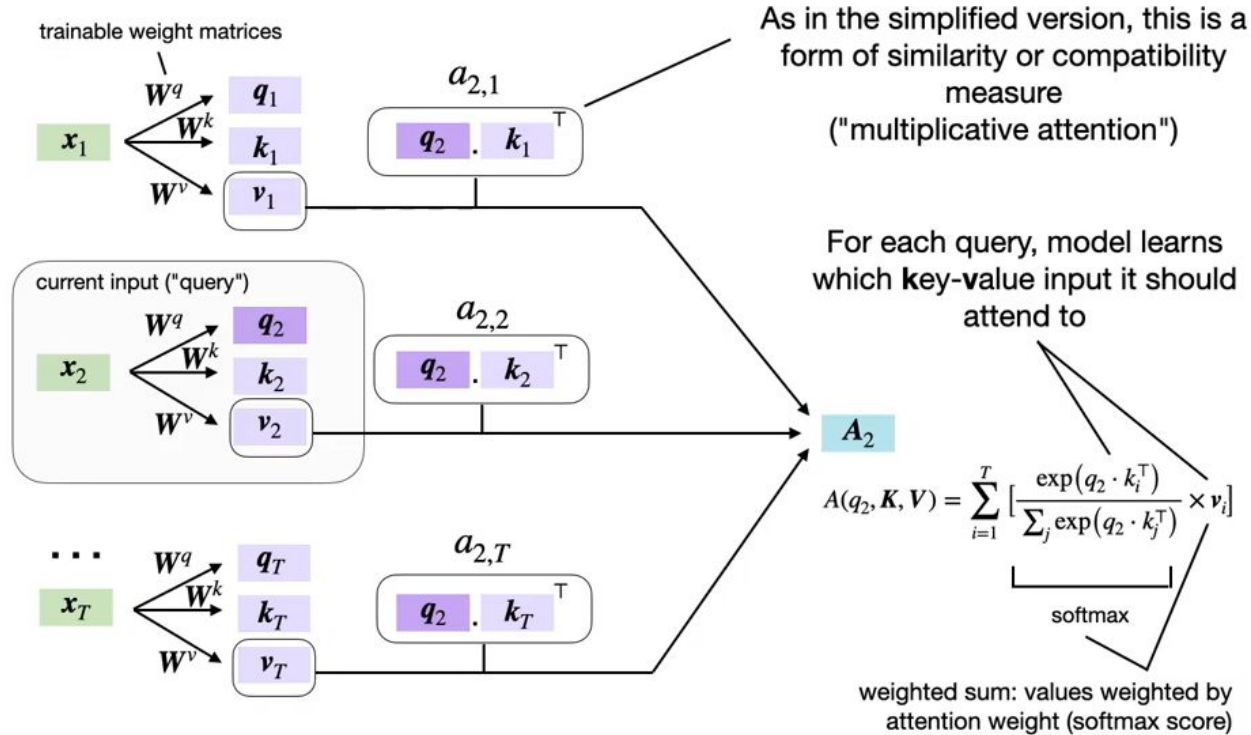
- Query matrix
- Key matrix
- Value matrix

Query:  $W_q x_i$

Key:  $W_k x_i$

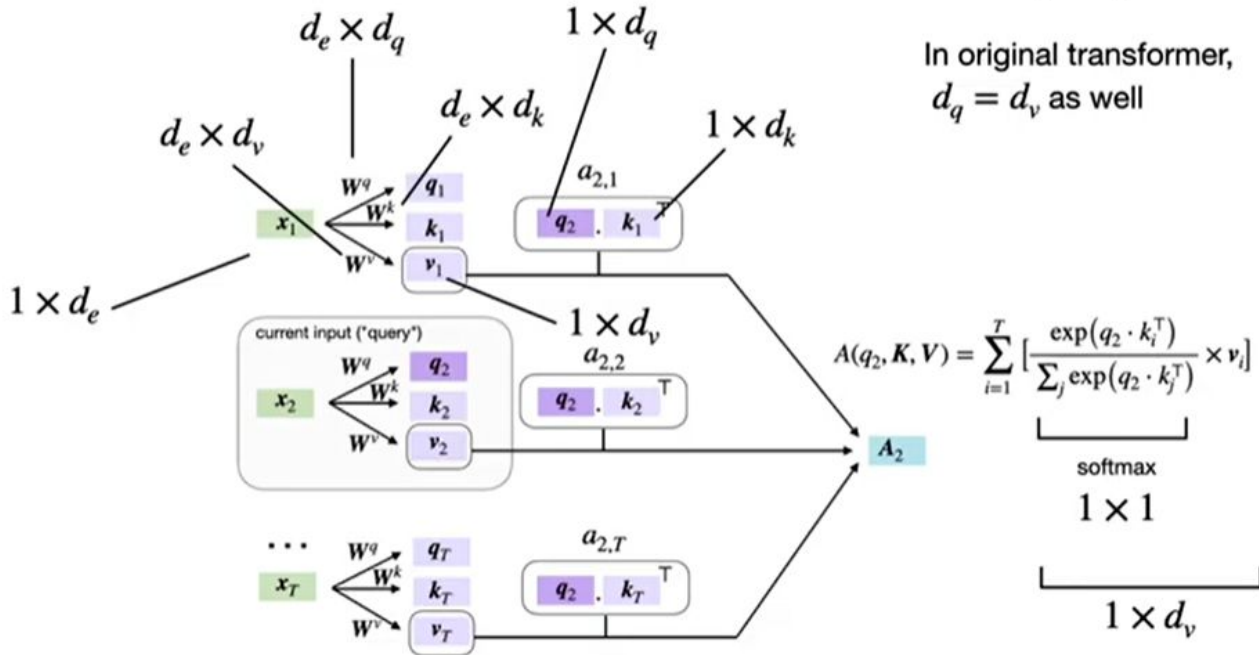
Value:  $W_v x_i$

# Attention without RNN (contd...)

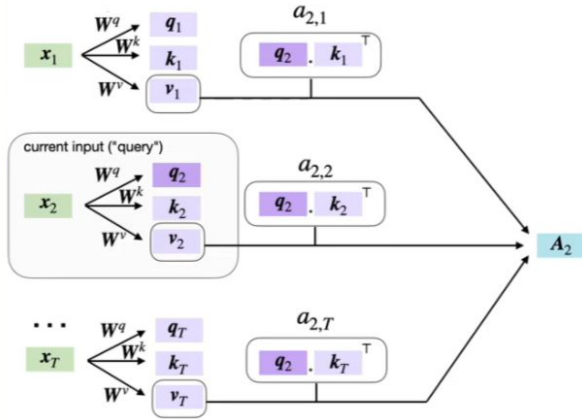


# Attention without RNN (contd...)

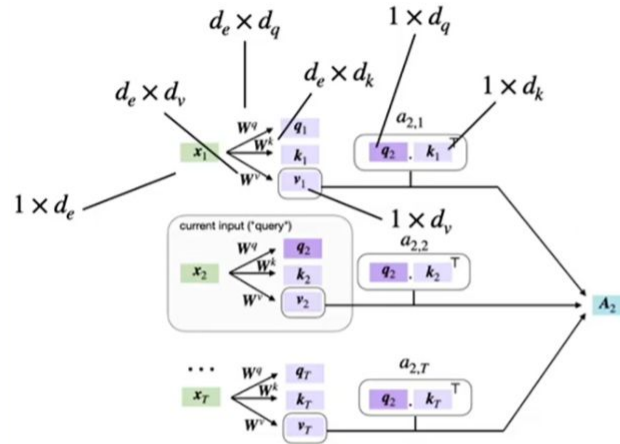
$d_e$  = embedding size (original transformer = 512)



# Attention without RNN (Matrix form)



This model was proposed by Ashish Vaswani et al.



Scaling to prevent high values in softmax preventing saturation

$$Q \in \mathbb{R}^{T \times d_q}$$

$$K \in \mathbb{R}^{T \times d_k}$$

$$V \in \mathbb{R}^{T \times d_v}$$

"attention matrix"  
 $T \times T$

$$A(Q, K, V) = \text{softmax} \left( \frac{QK^T}{\sqrt{d_k}} \right) V$$

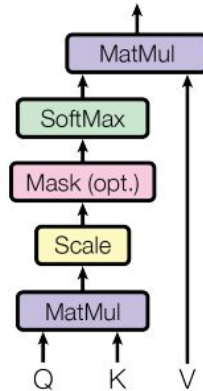
"attention-based embedding"  
 $T \times d_v$

# Multi-head Attention

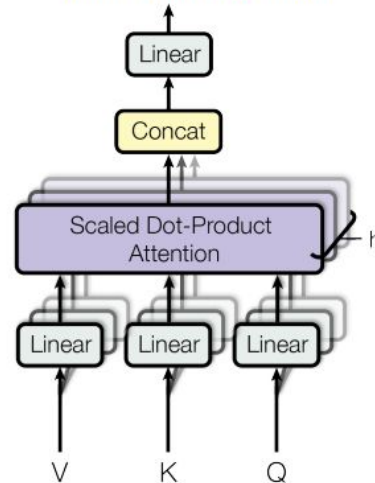
Use multiple sets of  $W_q$ ,  $W_k$  and  $W_v$  (8 was used in the original paper)

Each set learns its own set of parameters

Scaled Dot-Product Attention



Multi-Head Attention



# Transformer model

A transformer is basically a combination of encoder decoder modules

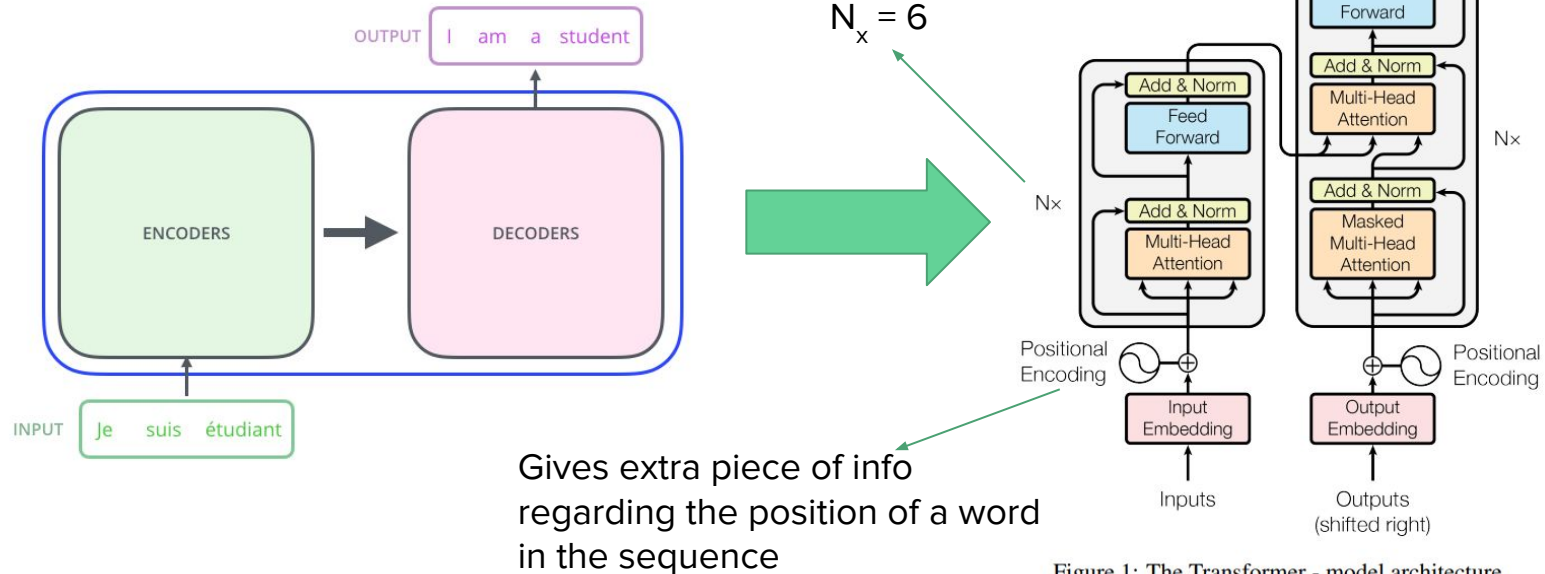


Figure 1: The Transformer - model architecture.

# Transformer model

The encoder component is essentially a stack of encoders having feed forward neural network and self-attention layer

Self-attention layer looks at the other words in the input sentence

Output of the last encoder is fed to the decoder component

