

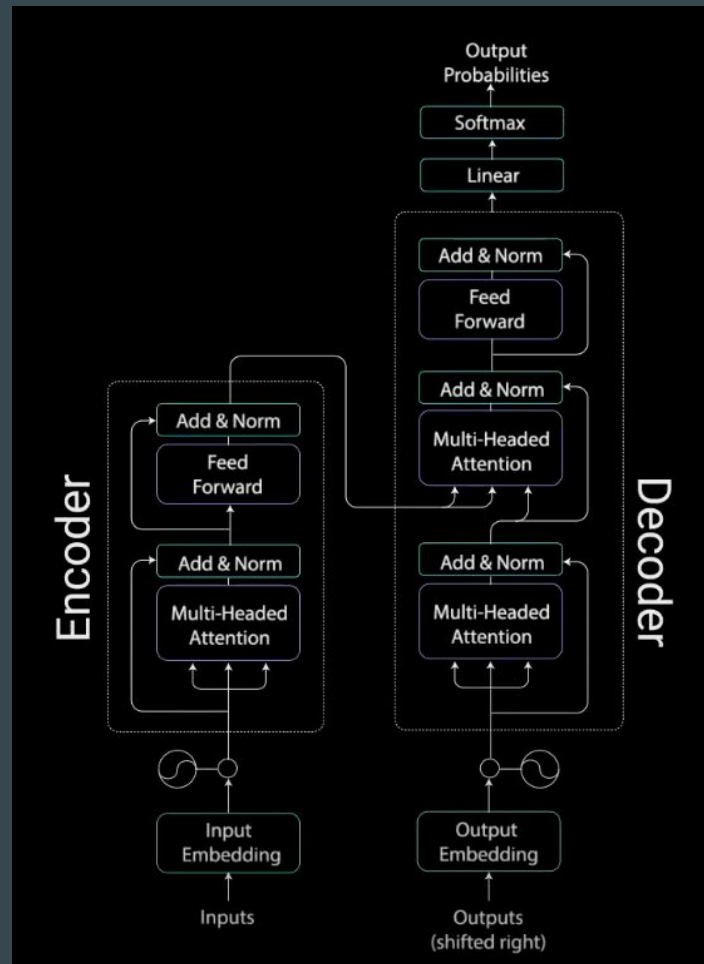
# RCOS AIHWKIT Status Update #3

...

Aasiya Husain

# Transformers

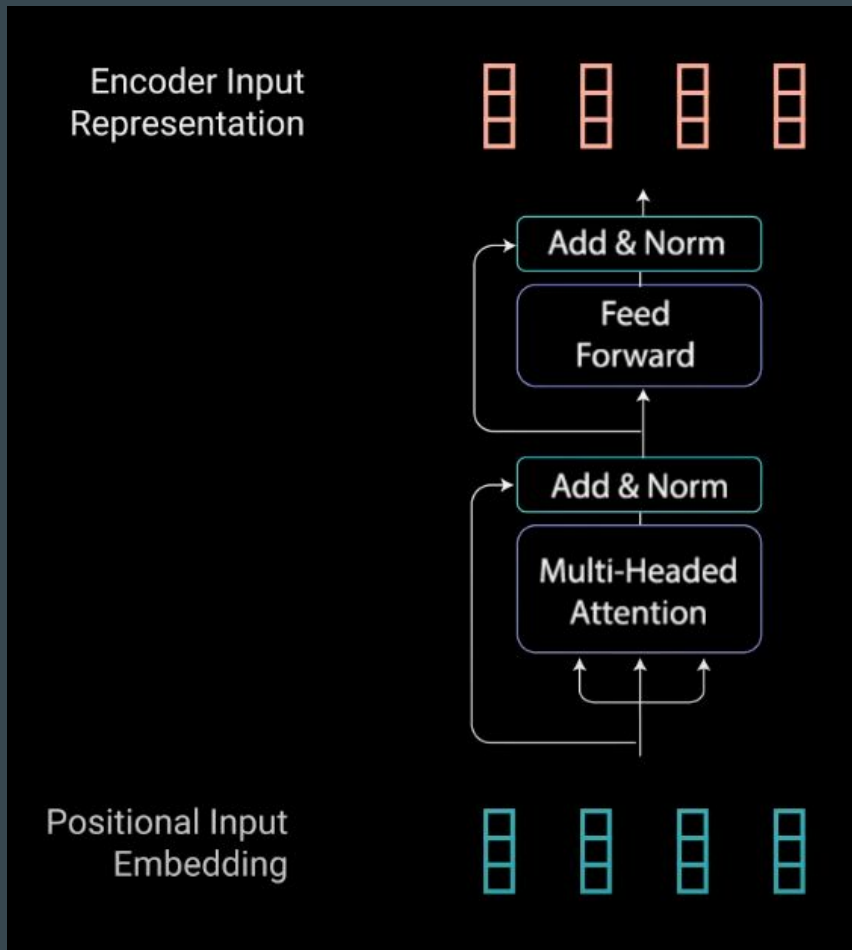
- 1) Input Embedding: Each word is mapped to a vector
- 2) Positional Encoding: Adds position information into the input vector.
  - For each odd index, add the cos function
  - For each even index, add the sin function
- 3) Query(Q), Key(K), Value(V) Vectors: These vectors are obtained by passing the input through 3 distinct linear layers.



# Encoder

## 4) MultiHeaded Attention:

- Q, K, V vectors each passed through separate linear layers
- $\text{DotProduct}(Q, K) \rightarrow \text{Scale by square root of the dimension of } Q \text{ and } K \rightarrow \text{softmax} = \text{Attention weights}$
- $\text{Output} = \text{Attention weights} \times \text{value}$



# Encoder Continued

5) Residual Connection: The output of attention is added to the original position input embedding

6) Normalization Layer

7) Feed-Forward network: =>

