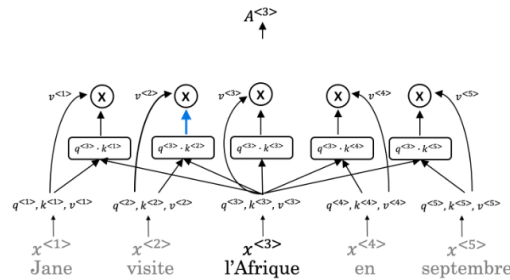


1. A Transformer Network, unlike its predecessors RNNs, GRUs and LSTMs, can process entire sentences all at the same time. (Parallel architecture). 1 point
  - ☐ False
  - ☒ True
2. Transformer Network methodology is taken from: (Check all that apply) 1 point
  - ☐ Convolutional Neural Network style of architecture.
  - ☒ Attention mechanism.
  - ☐ Convolutional Neural Network style of processing.
  - ☐ None of these.
3. What are the key inputs to computing the attention value for each word? 1 point



- ☐ The key inputs to computing the attention value for each word are called the quotation, key, and vector.
- ☐ The key inputs to computing the attention value for each word are called the quotation, knowledge, and value.
- ☐ The key inputs to computing the attention value for each word are called the query, knowledge, and vector.
- ☒ The key inputs to computing the attention value for each word are called the query, key, and value.

4. Which of the following correctly represents *Attention*? 1 point

☒  $A(Q, K, V) = \sum_i \left( \frac{\exp(q_i k_i^T)}{\sum_j \exp(q_i k_j^T)} \right) * V^{<i>}$   
☐  $A(Q, K, V) = \left( \frac{\exp(q_i k_i^T)}{\sum_j \exp(q_i k_j^T)} \right) * V^{<i>}$   
☐  $A(Q, K, V) = \sum_i \left( \frac{\exp(q_i k_i^T)}{\sum_j \exp(q_i k_j^T)} \right) * K^{<i>}$   
☐  $A(Q, K, V) = \sum_i \left( \frac{\exp(q_i k_i^T)}{\sum_j \exp(q_i k_j^T)} \right) * \sum_i v^i$

5. Are the following statements true regarding Query (Q), Key (K) and Value (V) ? 1 point

Q = interesting questions about the words in a sentence

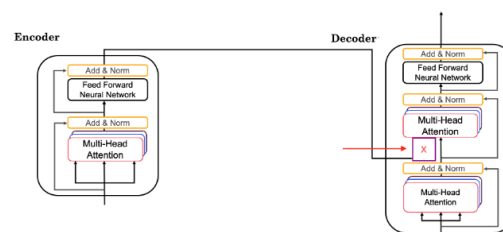
K = specific representations of words given a Q

V = qualities of words given a Q

- ☒ False
- ☐ True

**$Attention(W_i^Q Q, W_i^K K, W_i^V V)$**

6. What does  $i$  represent in this multi-head attention computation? 1 point
  - ☐ The computed attention weight matrix associated with the order of the words in a sentence
  - ☐ The computed attention weight matrix associated with the  $i$ th "word" in a sentence.
  - ☐ The computed attention weight matrix associated with specific representations of words given a Q
  - ☒ The computed attention weight matrix associated with the  $i$ th "head" (sequence)
7. Following is the architecture within a Transformer Network (*without displaying positional encoding and output layers(s)*). 1 point



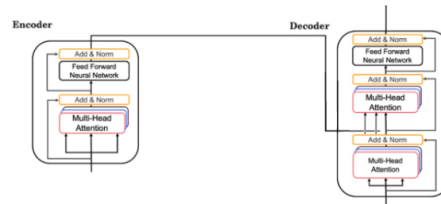
What information does the Decoder take from the Encoder for its second block of *Multi-Head Attention* ? (Marked  $X$ , pointed by the independent arrow)

(Check all that apply)

- ☒ V
- ☐ Q
- ☒ K

8. Following is the architecture within a Transformer Network (*without displaying positional encoding and output layers(s)*).

1 point



The output of the decoder block contains a softmax layer followed by a linear layer to predict the next word one word at a time.

- ☒ False
- ☐ True

9. Which of the following statements is true about positional encoding? Select all that apply.

1 point

- ☒ Positional encoding uses a combination of sine and cosine equations.
- ☒ Positional encoding provides extra information to our model.
- ☒ Positional encoding is used in the transformer network and the attention model.
- ☒ Positional encoding is important because position and word order are essential in sentence construction of any language.

10. Which of these is a good criterion for a good positional encoding algorithm?

1 point

- ☐ Distance between any two time-steps should be inconsistent for all sentence lengths.
- ☐ It should output a common encoding for each time-step (word's position in a sentence).
- ☒ The algorithm should be able to generalize to longer sentences.
- ☐ It must be nondeterministic.