

Transformers continued

- Attention: the sequence of vectors from the embedding step are passed through an operation called an attention block. This allows the vectors to talk to each other and pass information back and forth to update their value, giving context to words.

Ex a Machine Learning Model vs A Fashion Model

The attention block figures out which words in context (prompt) are relevant for updating the meanings of which other words " " and how exactly those meanings should be updated (the numbers in vector)

- MLP: A Multi layer perceptron AKA Feed forward layer.

The vectors after the attention operation pass through a MLP operation. and here the vectors don't communicate they all go through the same operation in parallel and is like asking a long list of questions such as about each vector and updating them based on the answer, this can help the LLM store facts.

Queen $\Rightarrow \begin{bmatrix} 2.1 \\ 5.2 \\ 3.3 \\ \vdots \\ x_i \end{bmatrix}$

- \rightarrow is it in english? ☒
- \rightarrow is it a name? ☒
- \rightarrow refers to a person? ☒
- \rightarrow is this fiction? ☐

- both the attention and MLP are a block of Matrix multiplications
 - there are also normalization steps in between and a final unembedding
- Then you go through repetitions of attention and MLP on and on the hope is that all the meaning of the prompt has been baked into the very last vector in the sequence. we then do a certain operation on the last step which produces the probability distribution of all possible next words (tokens), then feed it prompt + prediction and repeat until stop word or limit

