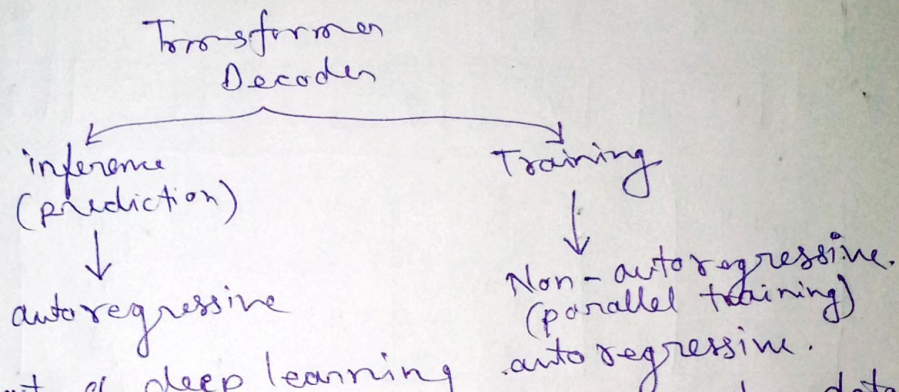


16 Mar 2025

MASKED MULTI-HEAD ATTENTION

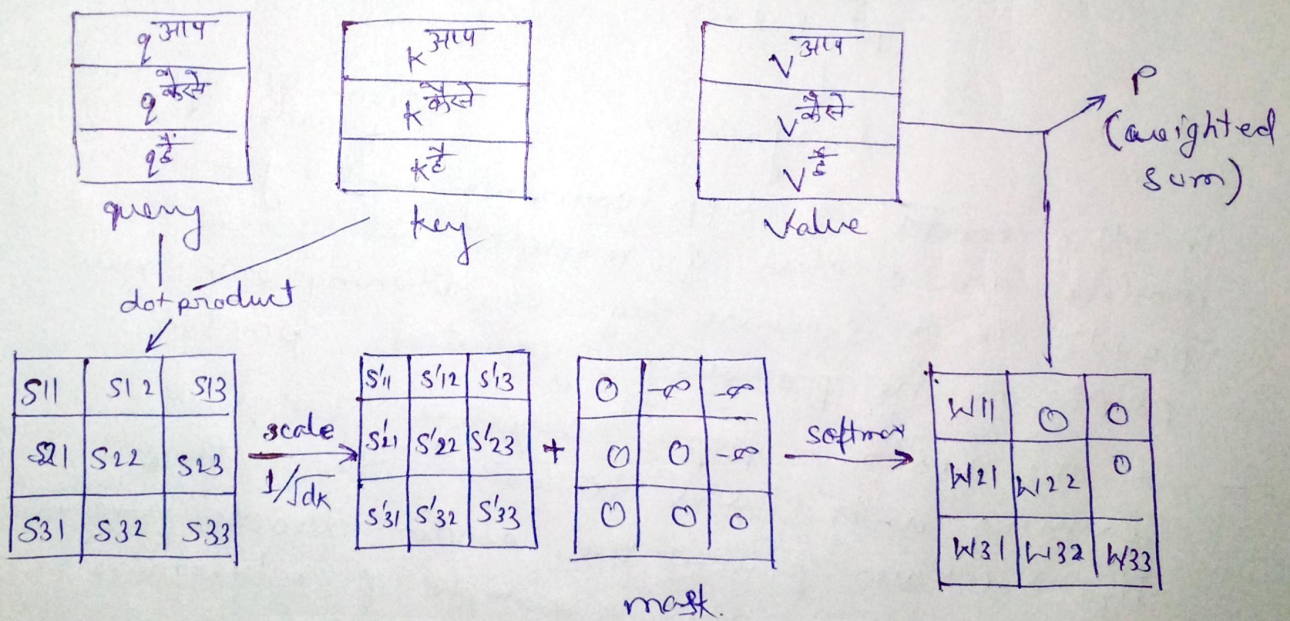
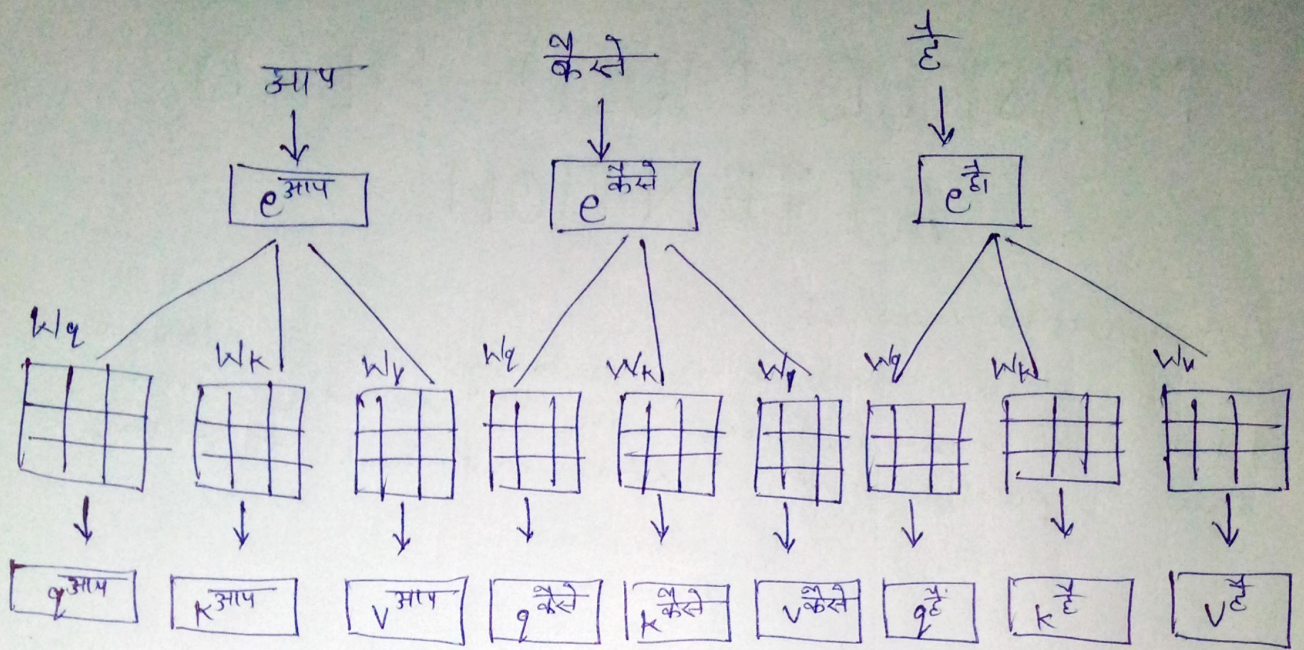
Autoregressive models:

The transformer decoder is autoregressive at inference time and non-autoregressive at training time. This is possible with masked multi-head attention.



- in the context of deep learning models are a class of models that generate data points in a sequence by conditioning each new point on the previously generated points.
- in the research paper "Sequence to Sequence Learning with Neural Networks" the author Ilya Sutskever present the encoder-decoder architecture is also an example of autoregressive model. (All seq2seq model).

masked multi-head Attention Architecture:



$$P = \begin{cases} \text{आप}_{ce} = W_{11} * V_{\text{आप}} + W_{12} * V_{\text{कैसे}} + W_{13} * V_{\text{है}} \\ \quad = W_{11} * V_{\text{आप}} \\ \text{कैसे}_{ce} = W_{21} * V_{\text{आप}} + W_{22} * V_{\text{कैसे}} + W_{23} * V_{\text{है}} \\ \quad = W_{21} * V_{\text{आप}} + W_{22} * V_{\text{कैसे}} \\ \text{है}_{ce} = W_{31} * V_{\text{आप}} + W_{32} * V_{\text{कैसे}} + W_{33} * V_{\text{है}} \end{cases}$$