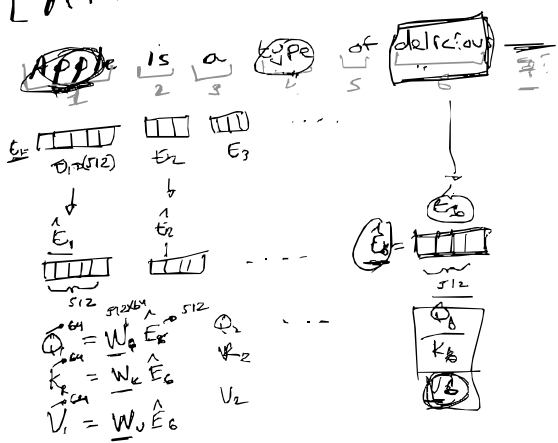
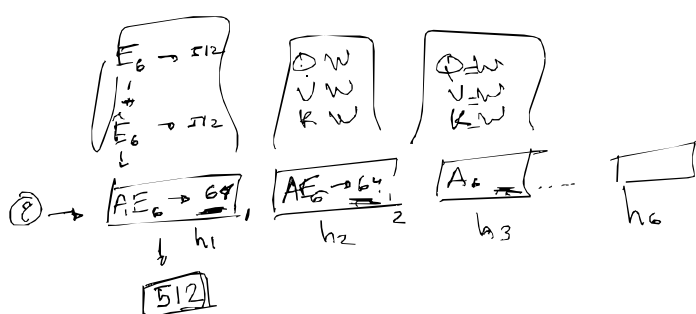


[Attention]

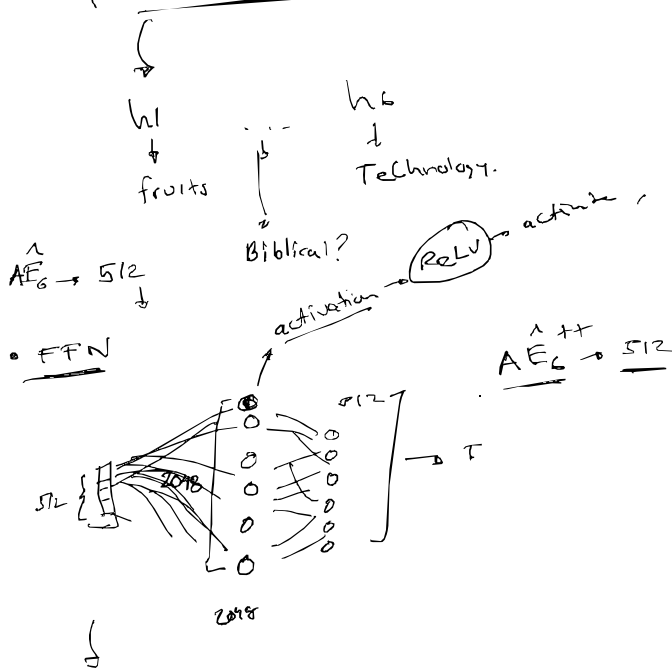


$$\alpha_i = \frac{\exp(Q_1 \cdot K_i)}{\sum_{j=1}^6 \exp(Q_1 \cdot K_j)}$$

$$V_1 \alpha_1 + \dots + V_6 \alpha_6$$

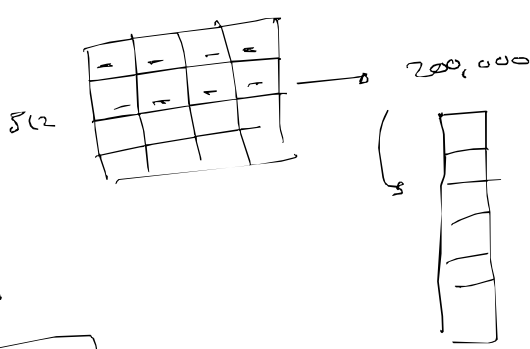
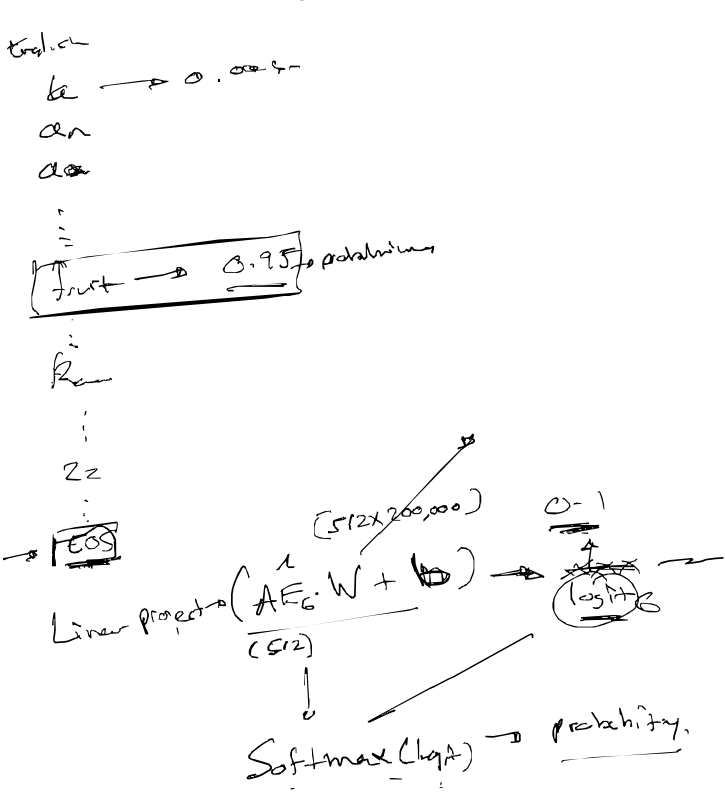


different Types of CONTEXT

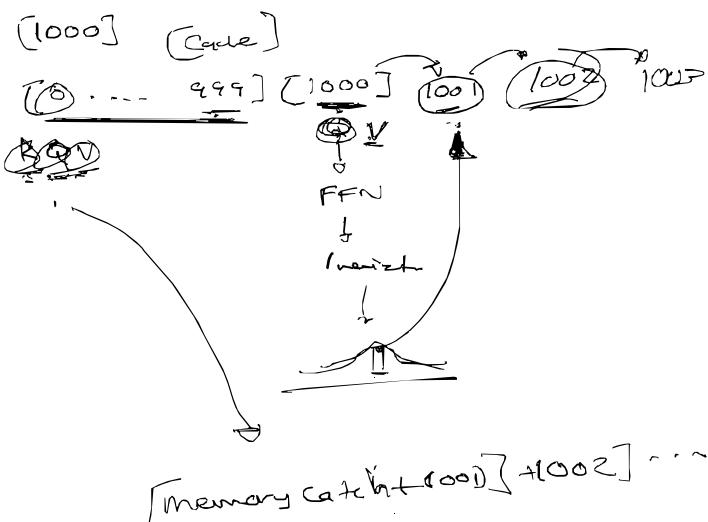


Linearization

Embedding Space ($A \cdot E_6^{++}$) \rightarrow Vocabulary (200,000)

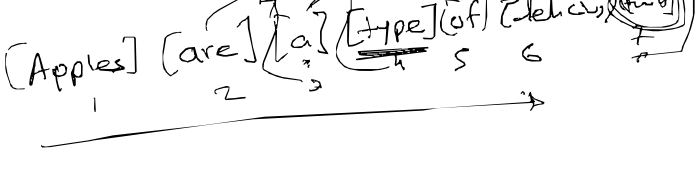


Inference

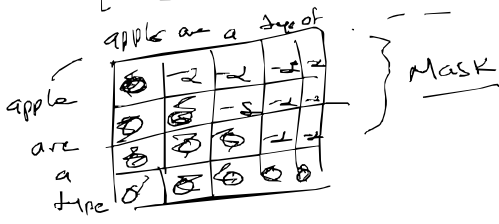


Training

Segments \rightarrow Context window



[6] \rightarrow Repetitions



Decoder \rightarrow Architecture \rightarrow Instead of RNN