

I - Input - ^{Token} Output pairs Determination

1. LLMs predict next ~~word~~ ^{tokens} based on a set of previous tokens
2. Total number of tokens taken as the basis is known as Context window.
3. LLMs are auto regressive, i.e. They predict the next Token based on the training set which is nothing but previous tokens predicted
4. Example :

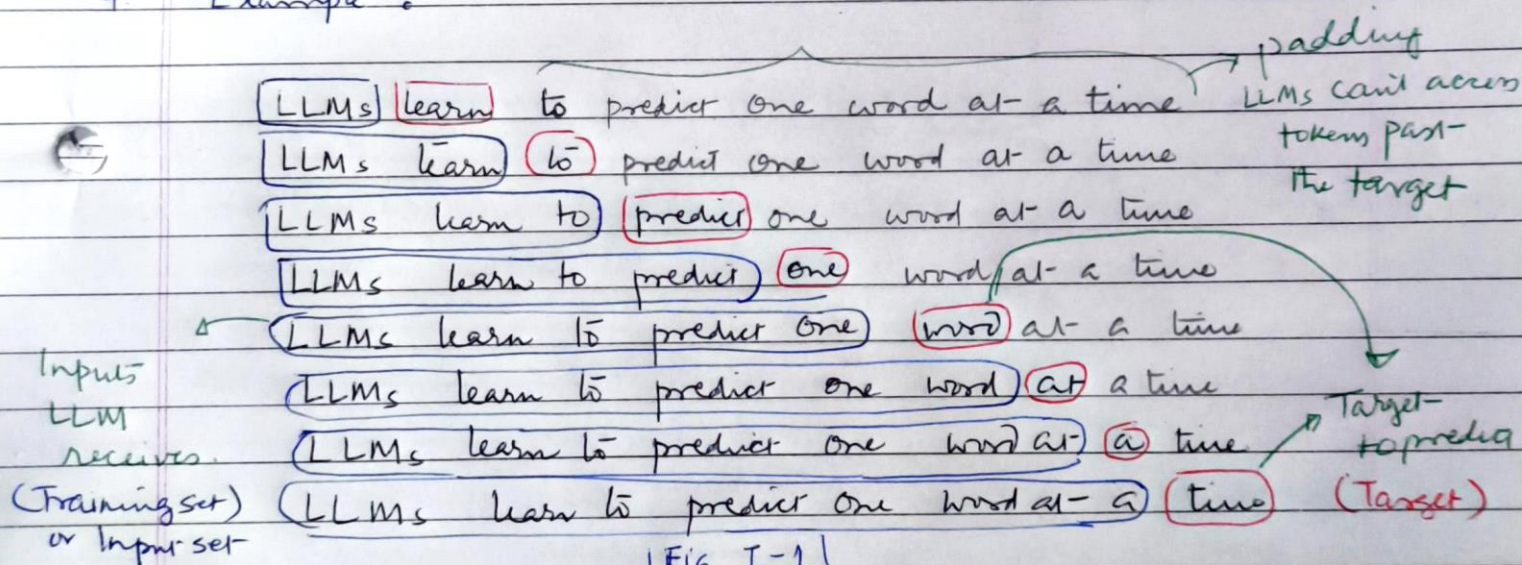


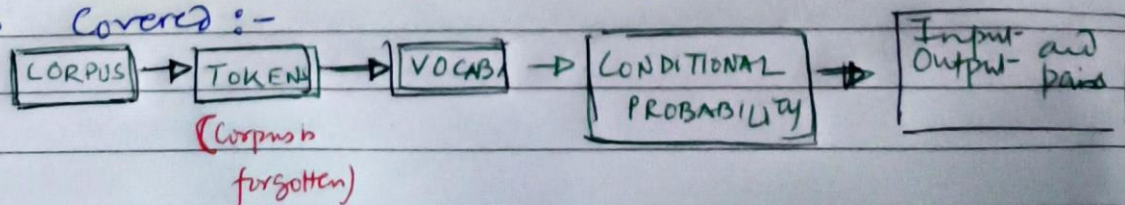
Fig I-1

5. Tokens are converted into Token Ids and then processed. Token Ids are numerical representation of each token.
6. For computation efficiency Each Token that is predicted above it via joint probability

$$P(w_1, w_2, w_3, w_4) = P(w_4 | w_1, w_2, w_3) \times \text{probability of a given sentence}$$

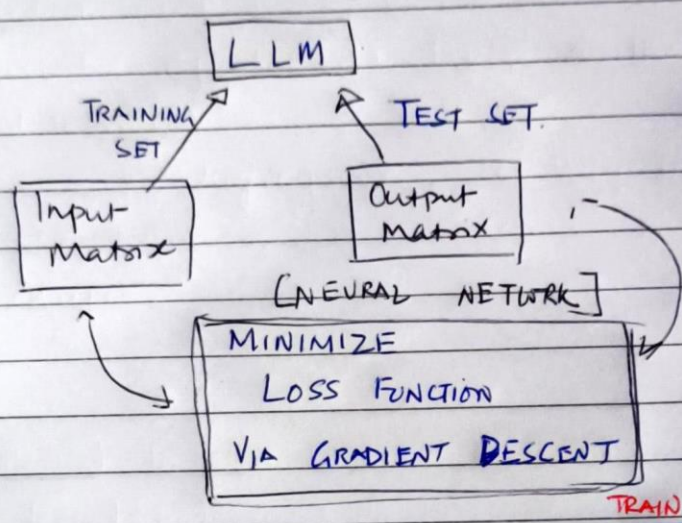
$$= P(w_1) \times P(w_2 | w_1) \times P(w_3 | w_1, w_2) \times P(w_4 | w_1, w_2, w_3)$$

~~XP~~
7. This probability is formed based on the Vocabulary that is formed Corpus that is used to make the vocabulary in the 1st stage out of the Corpus.
8. No machine learning has happened till now!
9. Steps Covered :-



II - TRAINING THE LLM

1.



The LLM is trained using the Input/output pairs determined in step I-9.

2. For Computational Efficiency the Input and Output sent to LLM is NOT as shown in I-9's figure Fig I-1. Rather it is loaded in BATCHES as follows :-

TRAINING SET

(Diff from Corpus that was used to Build Vocab)

'In the heart of the city stood the old library, A relic from a bygone era. Its stone walls bore the marks of time, and my clung tightly to its facade ...'

Input Matrix: $x = \begin{pmatrix} ["In", "the", "heart"], \\ ["of", "the", "city"], \\ ["stood", "the", "old"], \\ [...]] \end{pmatrix}$

all Tokens are not characters but token IDs!

Output Matrix: $y = \begin{pmatrix} ["the", "heart", "of"], \\ ["the", "city", "stood"], \\ [...]] \end{pmatrix}$

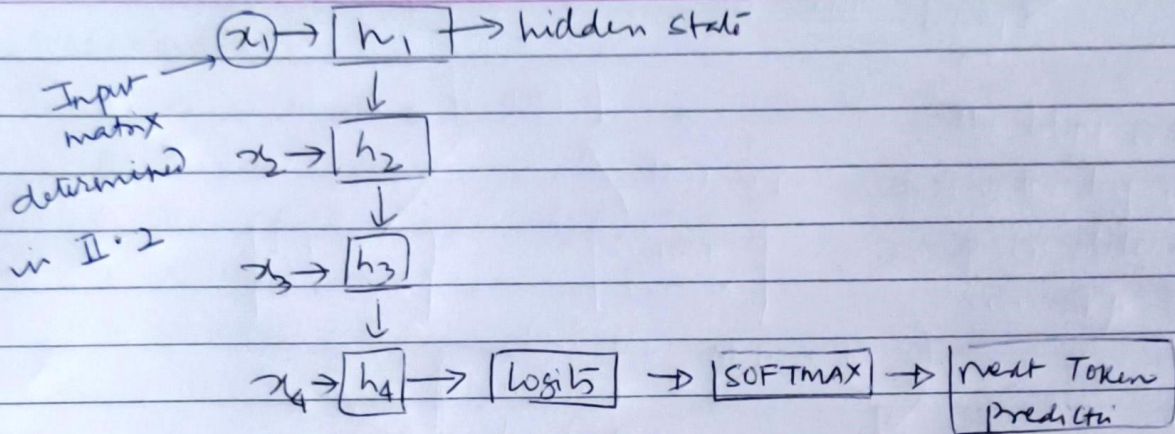
3. Using these Test- and Train matrices the neural network is trained.

III

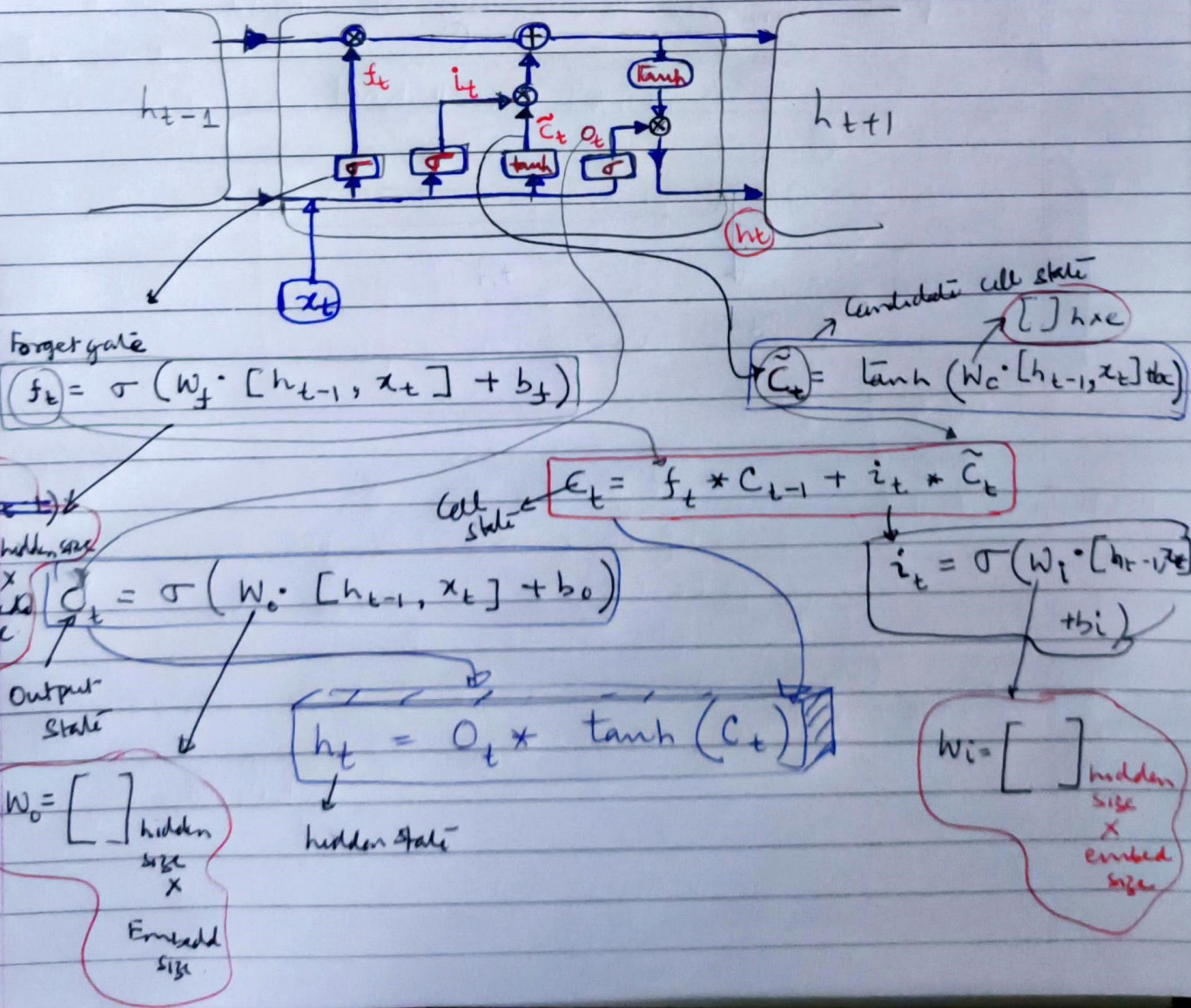
LSTM as the Neural Net (LM) THAT IS TRAINED

3

1.

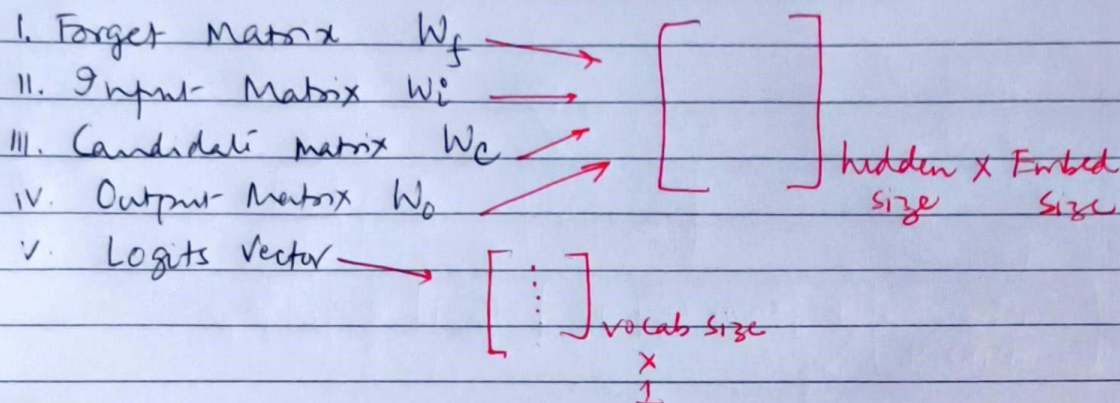


2. Making of the hidden state!



④

3. MATRICES :



4. Recap :

