

14.0 NEXT WORD PREDICTOR USING PyTorch & LSTM USING PyTorch

24 August 2025 02:49 AM RE Avinash Yadav

IMPLEMENTATION OF LSTM USING PYTORCH

USING THIS CONCEPT TO BUILD NEXT WORD PREDICTOR

USEFUL IN KEYBOARDS
SMARTPHONES, SEARCH
ENGINE FINALITY, CHATBOTS
ETC

DOWN THE LINE WE'LL ALSO
EXPLORE LANGUAGE MODELLING
CONCEPT IN THIS PROCESS

STRATEGY :

NEXT WORD PREDICTION or

Language Modelling is the task of predicting the next word (or character) in a sequence based on the context of previous words.

Used in auto-completion, machine translation, text-
Summarization etc.

What is the course fee for the Science Mentorship Program (SMP 2025)?
The course follows a monthly subscription model where you have to make monthly payments of Rs 799/month.
What is the total duration of the course?
The total duration of the course is 7 months. So the total course fee becomes $799 \times 7 = \text{Rs } 5600$ (approx.)
What is the syllabus of the mentorship program?
We will be covering the following modules:

Python Fundamentals
Data Analysis
SQL for Data Science
Maths for Machine Learning
ML Algorithms
Practical ML
MLOps
Case studies

UNSUPERVISED DATA
Means we do not
have target labels
here

so here we are given
with some text, so in
language modelling, first
we need to convert
the unsupervised data to
supervised data

Means we will make input feature and target labels as well.
So, first of all we'll split all the sentences from
unsupervised dataset.

Example: Python libraries for data science

input: Python
output: libraries
Python libraries for → for
Python libraries for data → data
Python libraries for data science → science

Now this same task would be repeated for all the
other sentences in the dataset. So finally we'll
have a dataset which will have input feature
target label as well

The next task would be forming a VOCAB - a
dictionary where each unique word is mapped with
a number

VOCAB = { 'python': 1,
 'libraries': 2,
 'for': 3,
 'data': 4,
 'science': 5
 :
 }

Now Based on this VOCAB replace the word in new
dataset with the values of the word, i.e. indexes

input: [1]
 [1, 2]
 [1, 2, 3]
 [1, 2, 3, 4]
output: [2]
 [3]
 [4]
 [5]

After this we'll convert every word into embeddings
like $[1, 2] \Rightarrow [[\dots], [\dots]]$

This will now go to
our LSTM.

TOP LEVEL STEPS :

① Data Pre-processing

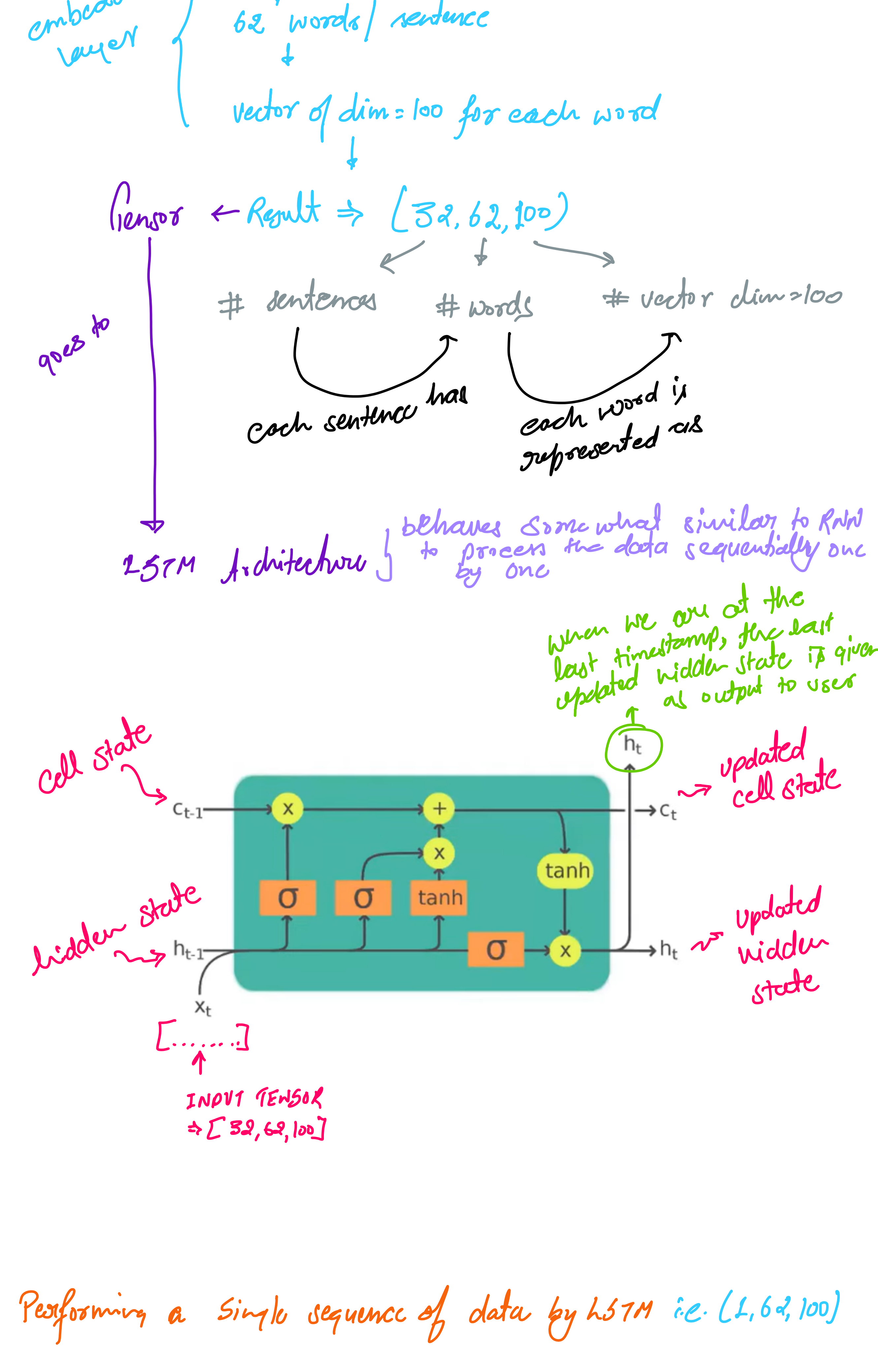
1.1) Unsupervised \rightarrow Supervised
1.2) English VOCAB \rightarrow NUMBER VECTOR
1.3) EMBEDDINGS

② Model Creation

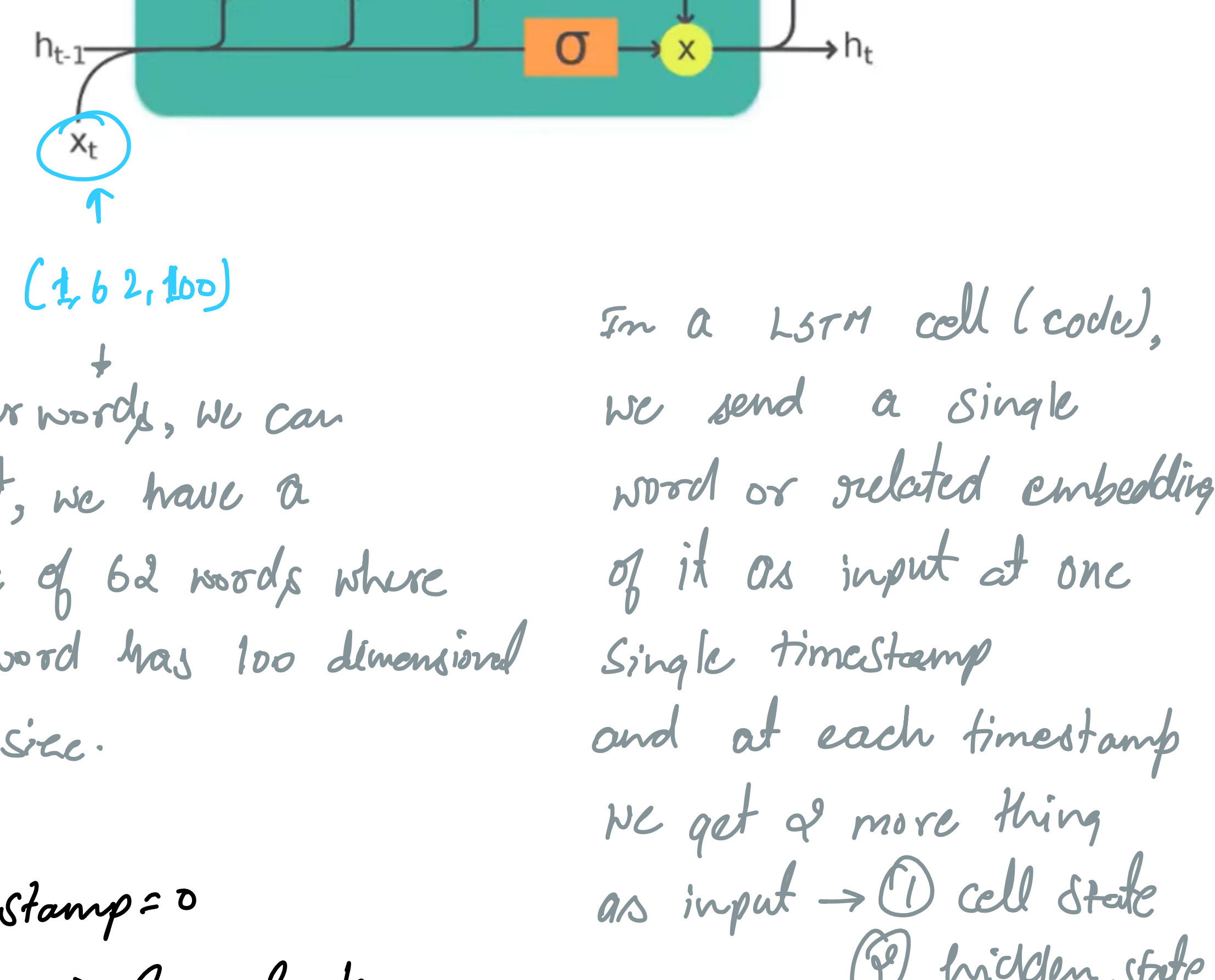
③ Training

④ Prediction on trained model

④ LSTM MODEL ARCHITECTURE :



Performing a Single sequence of data by LSTM i.e. (1, 62, 100)



Performing a Batch sequence of data by LSTM i.e. (32, 62, 100)

So the task we performed above for one sequence,
the same task will be performed parallelly
for the number of batches (# sentences)

i.e. parallelly 32 sequences will be executed at $T_{B=0}$
and rest of all the calculations all together
as first of above we did for one sequence.

The min-LSTM gives ③ outputs \Rightarrow All the intermediate
we get 3 things from LSTM after hidden states
its processing \Rightarrow

