**Vellore Institute of Technology**

(Deemed to be University under section 3 of UGC Act, 1956)

Tasks

Natural Language Processing

(CSE4022)

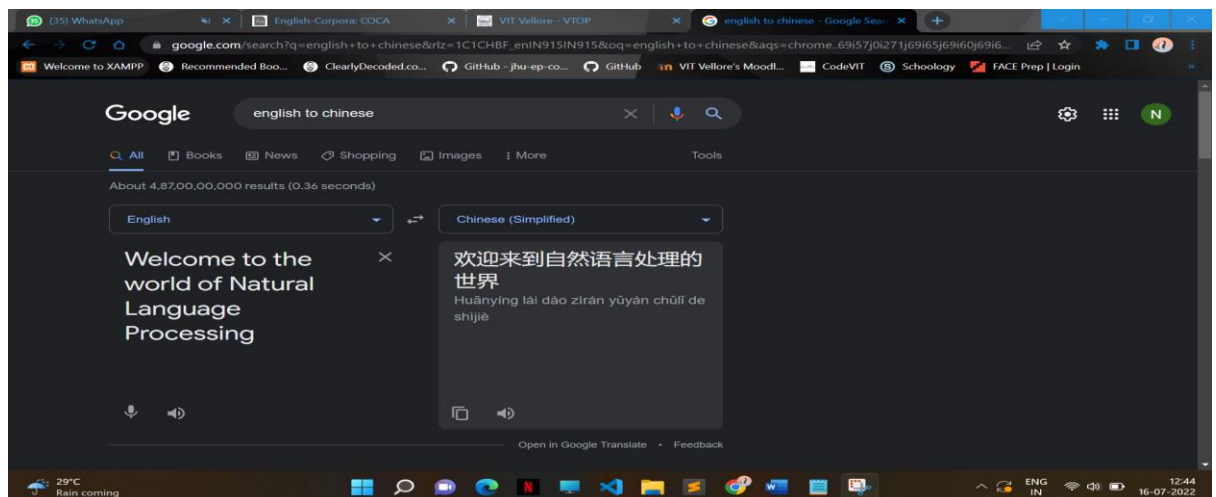Name: Nachiket Talwar

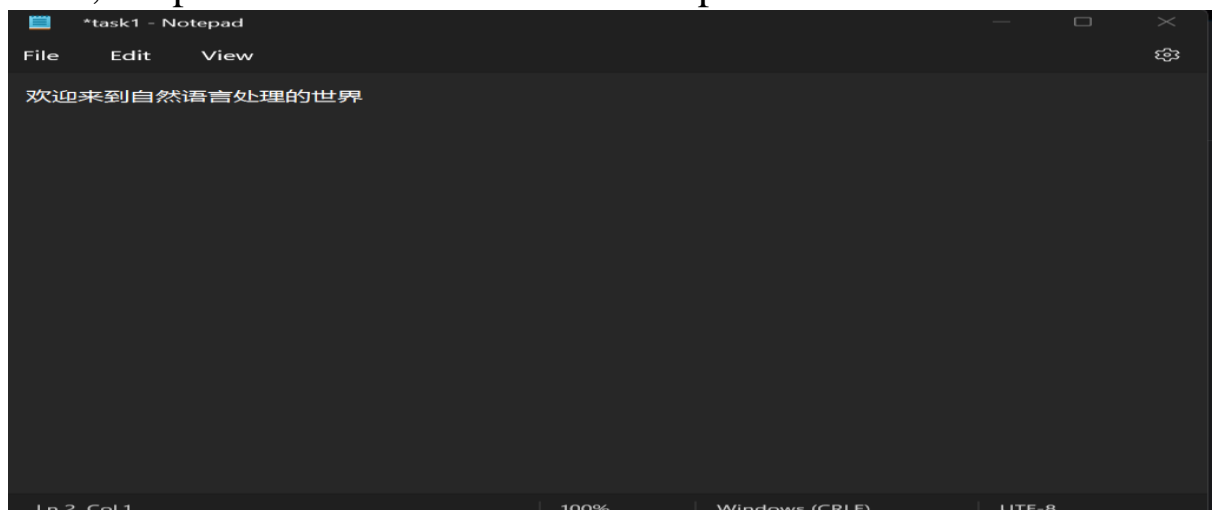Registration Number: 19BCE0840

Slot: C1

## Tasks assigned:

1. encoding demo
2. coca: story from data+ search string in text box: concirdance searching
3. accessing books : brown corpus, inaugral corpus,
4. experiening: frequency dist & conditional freq dist
5. suggest an application that u can build using these corporus: less than 20 words
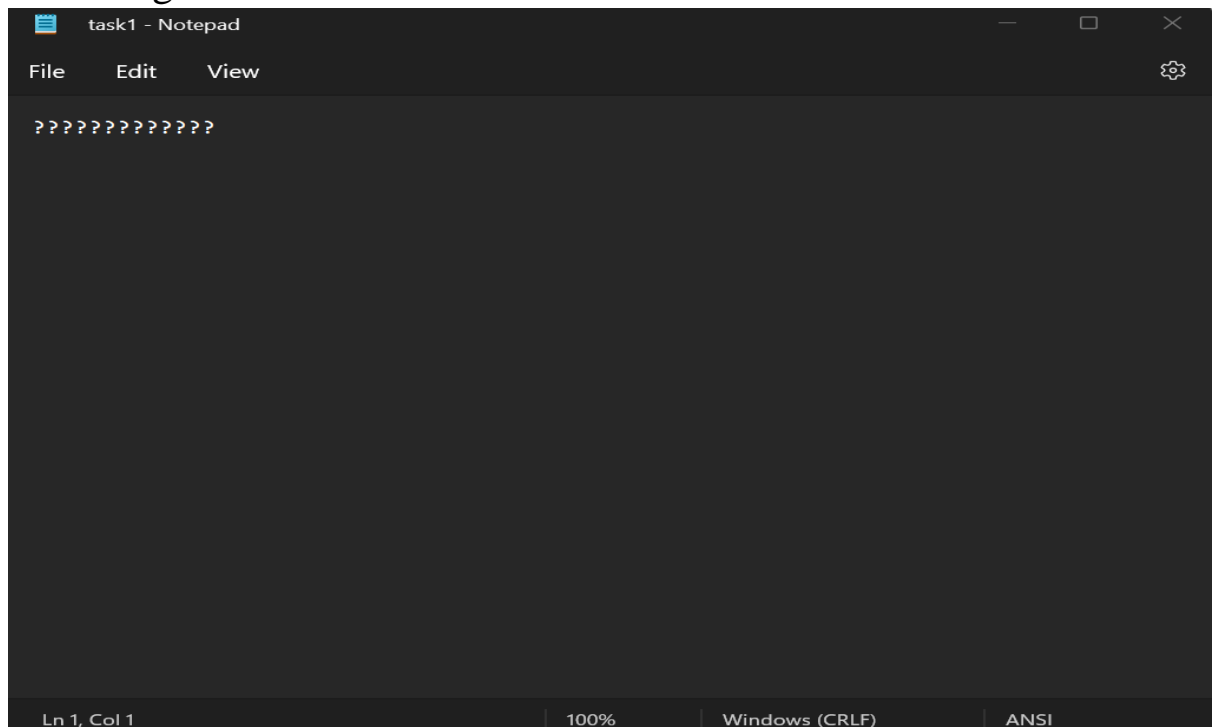
## Tasks solution:

1. We first get our Chinese letters which we will save to notepad


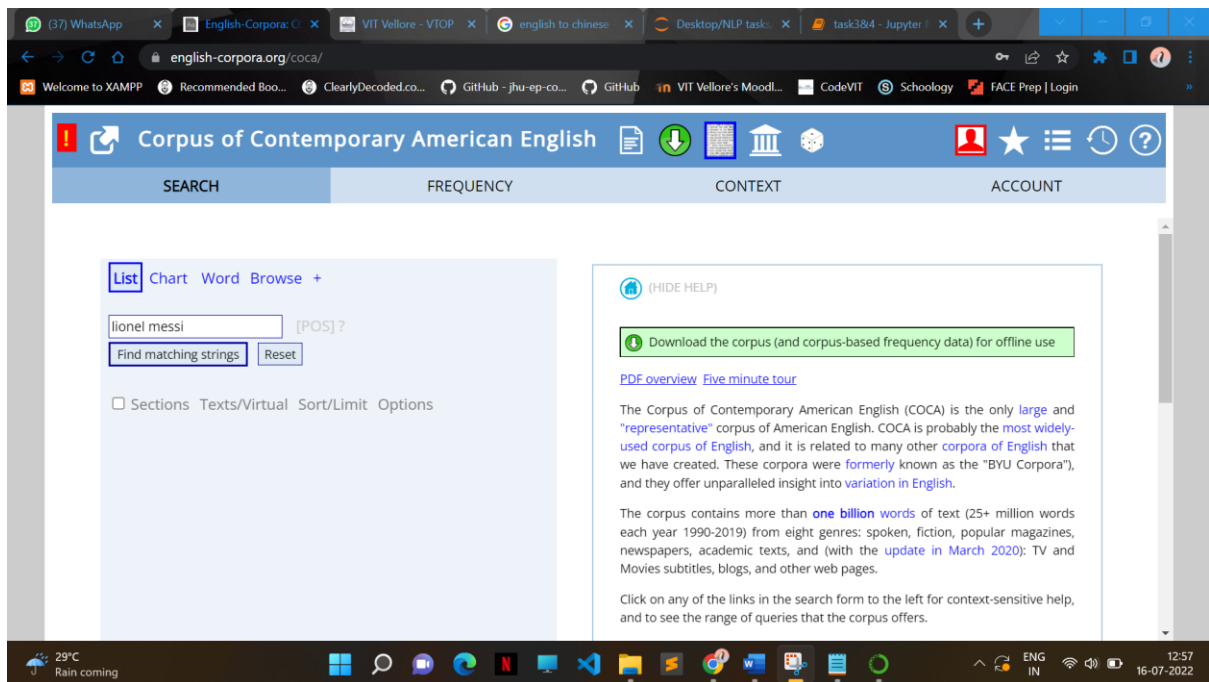
Next, we paste the characters in the notepad

We lastly save the file as with ANSI encoding and get the following result
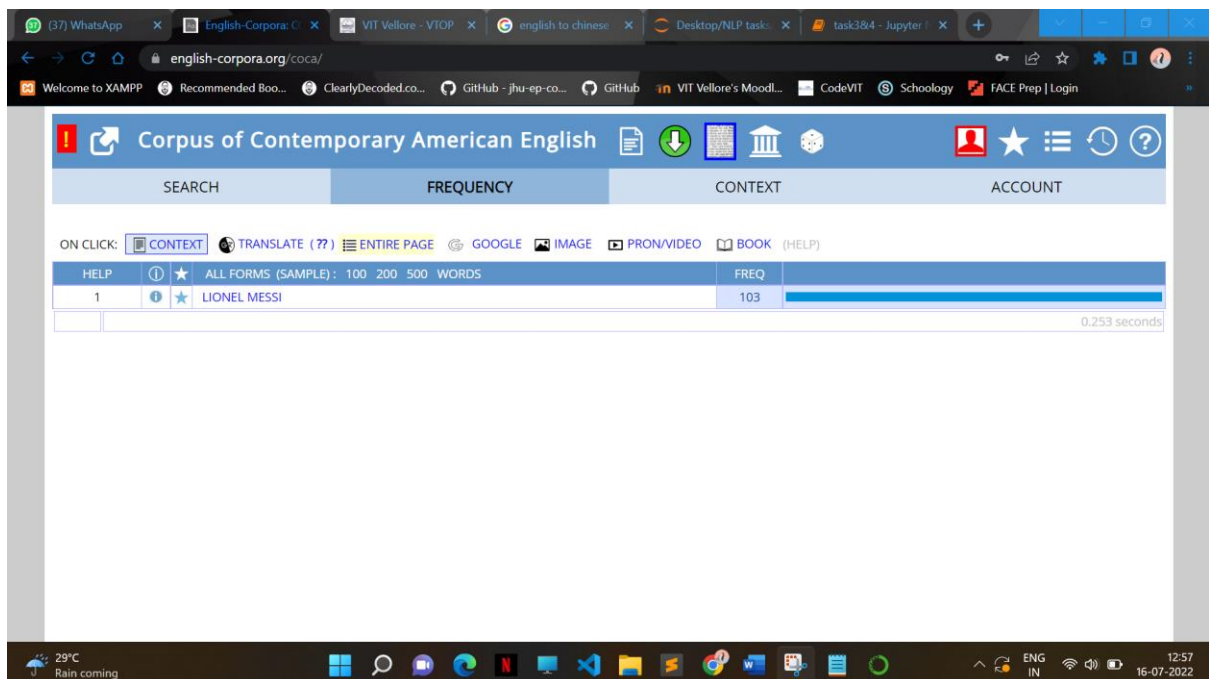


2. We used the corpus coca to find a string and create a story from that

We have decided to use the string 'lionel messi' for our coca search
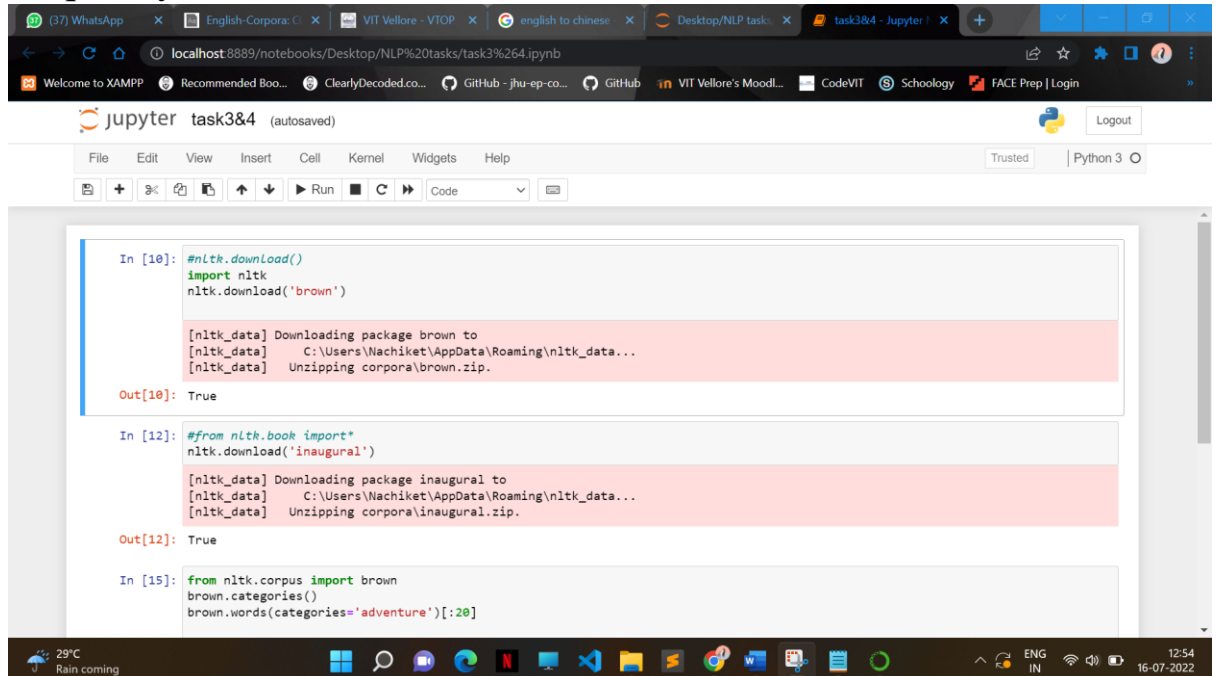
Below is the frequency of the appearance of our chosen string in the coca database



3. Tasks 3 and 4 have been plugged together as they work on the same piece of code

We use nltk toolkit to access brown and inaugural corpus and for task 4 we apply frequency distribution and conditional frequency distribution





4. We can create a blog accessing app. It will take a string and give us all the blogs that access a particular string