

In [1]:

```
from google.colab import drive
drive.mount('/content/drive/')
```

Drive already mounted at /content/drive/; to attempt to forcibly remount, call drive.mount("/content/drive/", force\_remount=True).

In [2]:

```
import nltk
import pandas as pd
```

In [3]:

```
nltk.download('punkt')
```

[nltk\_data] Downloading package punkt to /root/nltk\_data...  
[nltk\_data] Package punkt is already up-to-date!

Out[3]:

True

In [4]:

```
from nltk.tokenize import word_tokenize, sent_tokenize
filepath = '/content/drive/My Drive/English-Telugu/train.en'
corpus = open(filepath, 'r').read()
words = nltk.word_tokenize(corpus)
print("The number of tokens is", len(words))
average_tokens = round(len(words)/75000)
print("The average number of tokens per sentence is", average_tokens)
unique_tokens = set(words)
print("The number of unique tokens are", len(unique_tokens))
```

The number of tokens is 1809312  
The average number of tokens per sentence is 24  
The number of unique tokens are 21095

In [5]:

```
from nltk.tokenize import word_tokenize, sent_tokenize
filepath = nltk.data.find('/content/drive/My Drive/English-Telugu/train.te')
corpus = open(filepath, 'r').read()
words = nltk.word_tokenize(corpus)
print("The number of tokens is", len(words))
average_tokens = round(len(words)/75000)
print("The average number of tokens per sentence is", average_tokens)
unique_tokens = set(words)
print("The number of unique tokens are", len(unique_tokens))
```

The number of tokens is 1030924  
The average number of tokens per sentence is 14  
The number of unique tokens are 106016

In [6]:

```
f= open("/content/drive/My Drive/English-Telugu/train.en")
en=f.readlines()
len(en)
```

Out[6]:

75000

In [7]:

```
en[:6]
```

Out[7]:

```
words_on = []
```

```
words_en = []
for i in en:
    for word in i.split():
        words_en.append(word)

words_en[:10]
```

Out[14]:

```
['we', 'just', 'party', 'and', 'we', 'can', 'do', 'whatever', 'we', 'want']
```

In [15]:

```
words_te = []
for i in te:
    for word in i.split():
        words_te.append(word)

words_te[:10]
```

Out[15]:

```
[('the', 113593),
 ('and', 80265),
 ('of', 62955),
 ('to', 34296),
 ('in', 24579),
 ('that', 21052),
 ('you', 20618),
 ('he', 19632),
 ('i', 19186),
 ('a', 17782)]
```

In [16]:

```
from collections import Counter
english_words_counts = Counter(words_en)
telugu_words_counts = Counter(words_te)
```

In [17]:

```
import operator
english_words_counts = sorted(english_words_counts.items(), key = operator.itemgetter(1),
reverse = True)
telugu_words_counts = sorted(telugu_words_counts.items(), key = operator.itemgetter(1), reverse = True)
```

In [18]:

```
english_words_counts[:10]
```

Out[18]:

```
[('the', 113593),
 ('and', 80265),
 ('of', 62955),
 ('to', 34296),
 ('in', 24579),
 ('that', 21052),
 ('you', 20618),
 ('he', 19632),
 ('i', 19186),
 ('a', 17782)]
```

In [19]:

```
telugu_words_counts[:10]
```

Out[19]:

```
[('the', 12085),
 ('and', 9127),
 ('of', 9006),
 ('to', 8842),
 ('in', 8681)]
```

Tr. [26].

```
In [26]:
```

```
from nltk.translate.bleu_score import sentence_bleu
```

```
In [27]:
```

```
from tensorflow.keras import models
```

```
In [28]:
```

```
NMT1=models.load_model('/content/drive/My Drive/NMT_models/NMT1.h5')  
NMT1.summary()
```

```
Model: "sequential_1"
```

Layer (type)	Output Shape	Param #
=====		
embedding_1 (Embedding)	(None, 74, 256)	5810688
-----		
lstm_2 (LSTM)	(None, 256)	525312
-----		
repeat_vector_1 (RepeatVecto	(None, 47, 256)	0
-----		
lstm_3 (LSTM)	(None, 47, 256)	525312
-----		
time_distributed_1 (TimeDist	(None, 47, 106518)	27375126
=====		
Total params: 34,236,438		
Trainable params: 34,236,438		
Non-trainable params: 0		
-----		

```
In [29]:
```

```
NMT2=models.load_model('/content/drive/My Drive/NMT_models/NMT2.h5')  
NMT2.summary()
```

```
Model: "sequential"
```

Layer (type)	Output Shape	Param #
=====		
embedding (Embedding)	(None, 74, 256)	5810688
-----		
lstm (LSTM)	(None, 74, 256)	525312
-----		
lstm_1 (LSTM)	(None, 128)	197120
-----		
repeat_vector (RepeatVector)	(None, 47, 128)	0
-----		
lstm_2 (LSTM)	(None, 47, 256)	394240
-----		
lstm_3 (LSTM)	(None, 47, 128)	197120
-----		
time_distributed (TimeDistri	(None, 47, 106518)	13740822
=====		
Total params: 20,865,302		
Trainable params: 20,865,302		
Non-trainable params: 0		
-----		

```
In [30]:
```

```
from sklearn.model_selection import train_test_split  
  
x_train, x_test, y_train, y_test = train_test_split(x_padded, y_padded, test_size = 0.1)  
import numpy as np  
y_train = np.expand_dims(y_train, axis = 2)  
y_train.shape  
y_test = np.expand_dims(y_test, axis = 2)  
y_test.shape
```

```
Out[30]:
```

```
(7500, 47, 1)
```

In [31]:

```
def pad_to_text(padded, tokenizer):  
  
    id_to_word = {id: word for word, id in tokenizer.word_index.items()}  
    id_to_word[0] = ''  
    return ' '.join([id_to_word[j] for j in padded])
```

In [32]:

```
# function to make prediction  
def predictionNMT1(x, x_tokenizer = x_tokenizer, y_tokenizer = y_tokenizer):  
    predictions = NMT1.predict(x)[0]  
    id_to_word = {id: word for word, id in y_tokenizer.word_index.items()}  
    id_to_word[0] = ''  
    return ' '.join([id_to_word[j] for j in np.argmax(predictions,1)])  
  
# function to make prediction  
def predictionNMT2(x, x_tokenizer = x_tokenizer, y_tokenizer = y_tokenizer):  
    predictions = NMT2.predict(x)[0]  
    id_to_word = {id: word for word, id in y_tokenizer.word_index.items()}  
    id_to_word[0] = ''  
    return ' '.join([id_to_word[j] for j in np.argmax(predictions,1)])
```

In [33]:

```
y_test=y_test.reshape(y_test.shape[0],y_test.shape[1])
```

In [34]:

```
yp1=[]  
yp2=[]  
yt=[]  
for i in range(len(x_test)):  
    yt.append(pad_to_text(y_test[i], y_tokenizer))  
    yp1.append(predictionNMT1(x_test[i:i+1]))  
    yp2.append(predictionNMT2(x_test[i:i+1]))
```

In [35]:

```
BLEUScore1 = nltk.translate.bleu_score.corpus_bleu(yp1, yt)  
print(BLEUScore1)
```

0.4779748973490351

```
/usr/local/lib/python3.7/dist-packages/nltk/translate/bleu_score.py:490: UserWarning:  
Corpus/Sentence contains 0 counts of 2-gram overlaps.  
BLEU scores might be undesirable; use SmoothingFunction().  
    warnings.warn(_msg)
```

In [36]:

```
BLEUScore2 = nltk.translate.bleu_score.corpus_bleu(yp2, yt)  
print(BLEUScore2)
```

0.48104068819490353

```
/usr/local/lib/python3.7/dist-packages/nltk/translate/bleu_score.py:490: UserWarning:  
Corpus/Sentence contains 0 counts of 2-gram overlaps.  
BLEU scores might be undesirable; use SmoothingFunction().  
    warnings.warn(_msg)
```