

NLP-Project : Group 18

Official Baatmi

Priyam Bajpai
S20190010144
UG-3, CSE
priyam.b19@iiits.in

Dishant Tayade
S20190010175
UG-3, CSE
dishanthiraman.t19@iiits.in

Satyam Kumar Singh
S20190010158
UG-3, CSE
satyamkumar.s19@iiits.in

Raghav Garg
S20190010148
UG-3, CSE
raghav.g19@iiits.in

Happy Kumar
S20190010061
UG-3, CSE
happy.k19@iiits.in

Abstract—The Indian Subcontinent is a combination of many nations which led to the incorporation of many languages, the major ones being the Indo-Aryan languages spoken by 78.05% of Indians and the Dravidian languages spoken by 19.64% of Indians. Both language families together are sometimes known as Indic languages. Taking into consideration the vast population of the nation and its language diversity, it is very difficult to have an exact count of the languages spoken here. Languages like Marathi, Tamil, Telugu and Bengali are one of the most spoken languages in the nation after Hindi and English. In the fast growing demand of data today, need of a model that inter-converts texts between these regional languages to Hindi/English has become quite obvious. In this paper, we aim to propose a model that uses various NLP pre-processing techniques and an LSTM model to inter-convert sentences between Hindi and Marathi language. Based on the parallel corpus of around 40,000 sentences as input, the model gives an accuracy of around 67% in the train data.

Index Terms—NLP, POS, N-gram, seq2seq, LSTM, NMT

I. INTRODUCTION

In this world, population increases daily. Due to this massive rise, people tend to retain their roots and customs, sometimes in form of religion, other times in form of language. But due to this, a plethora of languages coexist together and to solve this problem translation of these many languages becomes a pivotal task.

To make translation easier, Natural language processing (NLP) is done as it automatically converts one natural language into another with the help of Machine translation (MT), thus preserving the meaning of the input text, and producing fluent text in the output language. But translation is not as easy as it seems because many languages have different genealogical families i.e., their origins are widely different from each other.

One such example of genealogical language families is Indic languages which are made up of Indo-Aryan languages and Dravidian languages. Here in figure 1 language distribution across India is presented. But even with Indo-Aryan and Dravidian languages being of the same family, they differ widely.

According to the Indian Language policy, it states that the state government can run its daily work on three languages - Hindi, English or any other constitutionally recognised scheduled language. The state of Maharashtra is one such

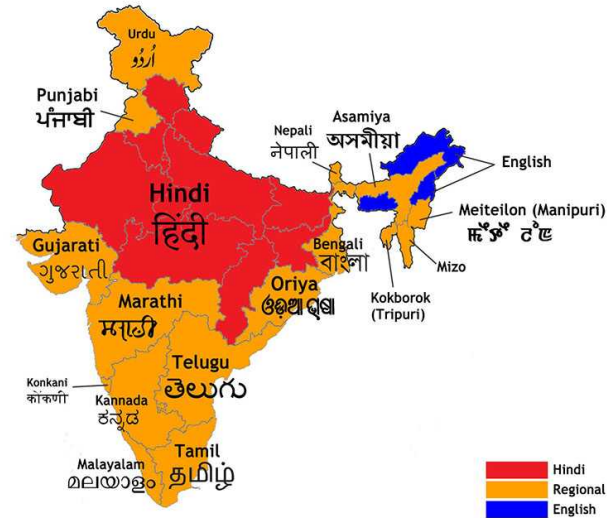


Fig. 1. Language distribution

state, that covers its officially written documents in its regional language Marathi. There is hardly any paperwork that is done in any other language other than Marathi. According to the Maharashtra Official Languages Act, 1964, it states that Marathi is the official language in Maharashtra, and should be used as per the three-language formula of the Centre, i.e., Marathi should be used along with Hindi and English. This circular also asked concerned offices to use Marathi along with English and Hindi languages in all written and oral examinations conducted during staff recruitment as Marathi is the *Rajbhasha* or the state language of Maharashtra.

As one can see Language translation is one of the most tedious tasks in a country as diverse as India but for sharing information, it becomes more necessary. To solve these problems, our model attempts Hindi to Marathi language conversion. There were many languages to choose from, but there are two reasons for taking up Marathi and Hindi.

- Both share the same script Devanagari which makes it easier to build and train the model.
- The group members have a fair knowledge of these two languages, which makes it easier to work on new data, and make our own models for POS Tagging and

Tokenization.

For our model we have planned to use a part of the Samanantar dataset which is probably one of the largest datasets available for training our model, and also use a seq2seq model which is a method of encoder-decoder based machine translation and language processing that maps an input of sequence to an output of sequence with a tag (for our model specially POS tagging) and attention value. The idea is to use 2 RNNs that will work together with a special token and try to predict the next state sequence from the previous sequence.

II. STATE OF THE ART/BACKGROUND

Both the languages are members of the same family i.e. Indo-Aryan family of languages. They share a lot of features from script to phonemes and vocabulary. The problem statement was picked up from the Similar Language Task (SLT) from WMT 2020. Hindi-Marathi translation lacks background work. However, similar work is found on Hindi-Nepali pair at WMT19 shared task of similar language translation Laskar et al. (2019). The literature survey mainly focuses on NMT for low resource language pairs since NMT outperforms conventional SMT on low resource pairs like English to Mizo, English to Hindi, English to Punjabi, and English to Tamil Pathak et al. (2018); Pathak and Pakray (2018); Laskar et al. (2019). It is noticed that training performance improves while parallel training data increases. For low resource languages, it is difficult to collect parallel data unlike monolingual data which is easily found through online sources. Hence, monolingual based NMT systems are introduced to enhance the translation quality of low resource language pair translation. To get the advantage of monolingual data, unsupervised pre-train methods are introduced Ramachandran et al. (2017); Variš and Bojar (2019). Conneau and Lample (2019) proposed XLM based on bidirectional encoder representations from transformers (BERT) where the contextual language model is built with words based on preceding and succeeding context. No work has been done on Hindi-Marathi low resource language pair with such advanced NMT based approach, from the best of our knowledge. Our work investigates XLM model on Hindi-Marathi low resource language pair translation. NUIG-Panlingua-KMI submission to WMT 2020 seeks to push the state-of-the-art in Similar Language Translation Task for HindiMarathi language pair. As part of these efforts, we conducted a series of experiments to address the challenges for translation between similar languages. Among the 4 MT systems prepared under this task, 1 PBSMT systems were prepared for HindiMarathi each and 1 NMT systems were developed for HindiMarathi using Byte Pair Encoding (BPE) into subwords. The results show that different architectures NMT could be an effective method for developing MT systems for closely related languages. Our Hindi-Marathi NMT system was ranked 8th among the 14 teams that participated and our Marathi-Hindi NMT system was ranked 8th among the 11 teams participated for the task. Mujadia et al. experimented with attention based recurrent neural network

architecture (seq2seq) for this task. We explored the use of different linguistic features like POS and Morph along with back translation for Hindi-Marathi and Marathi-Hindi machine translation. In their paper, Amit et al. use the Transformer network based NMT system (Vaswani et al., 2017) because it is among the state of the art models for machine translation. The work reported for this shared task is an extension of the work done by (Kumar and Singh, 2019) for similar languages task for 2019, which had also used a transformer based NMT system. Our work revolves around a similar attempt for Indic language translation for hindi and marathi pair of corpora.

III. PROPOSED SYSTEM

The whole work is divided into four parts -

- Text Input, pre-processing and tokenization
- N-gram Model for both the languages
- POS Tagger for both the languages
- LSTM model for translation

We will look into details about each process and how was it implemented in the following sections.

A. Text Input, pre-processing and tokenization

1) **Text Input:** We have taken our train dataset from the Samanantar Dataset from AI4Bharat and used the Hindi-Marathi parallel corpora that contains approx 0.2 million sentences. We have utilised 40 thousand sentence, and pre-processed it to train our LSTM model on it.

2) **Pre-processing:** For pre-processing, we needed our own list of stop-words that cater to our need of the devanagari script. Also we required some method to remove hindi numerals used extensively in the dataset. We also tried out methods to mine and remove english words between the sentences in the corpora.

3) **Tokenization:** Tokenization is essentially splitting a phrase, sentence, paragraph, or an entire text document into smaller units called tokens, which is done after our data is cleaned with the help of previous two steps. These tokens help in understanding the meaning of our selected corpus by analyzing the sequence of the words and later processing our data for further analysis and also for developing the model for the NLP project.

['याबाबत',
'चिंचवड',
'पोलीस',
'ठाण्यात',
'गुन्हा',
'दाखल',
'करण्यात',
'आला',
'आहे',
'आजही',
'मोठ्या',
'प्रमाणावर',
'याचीच',
'चर्चा',
'...'

Fig. 2. Tokenization of marathi corpus

```
[ 'जैसे',
  'लेकर',
  'चिंचवड',
  'पुलिस',
  'थाने',
  'में',
  'मामला',
  'दर्ज',
  'किया',
  'गया',
```

Fig. 3. Tokenization of hindi corpus

Tokenization is done in our code through the function `word_tokenize` and stored in the variable `finlist` which is shown in figure 2 and figure 3 for marathi and hindi corpus respectively.

B. N-gram

N-gram is a probabilistic model that's trained on a corpus of text to do a variety of tasks like speech recognition, machine translation and predictive text input. An N-gram model is one type of a Language Model (LM), which is about finding the probability distribution over word sequences.

For instance take these two sentences: "There was heavy rain" vs. "There was heavy flood". From experience, we know that the former sentence sounds better. An N-gram model will tell us that "heavy rain" occurs much more often than "heavy flood" in the training corpus. Thus, the first sentence is more probable and will be selected by the model.

A model that simply relies on how often a word occurs without looking at previous words is called unigram. If a model considers only the previous word to predict the current word, then it's called bigram. If two previous words are considered, then it's a trigram model.

An n-gram model for the above example would calculate the following probability:

$$P('There was heavy rain') = \quad (1)$$

$$P('There' | 'was' | 'heavy' | 'rain') \quad (2)$$

$$= P('There')P('was' | 'There')P('heavy' | 'Therewas') \quad (3)$$

$$P('rain' | 'Therewasheavy') \quad (4)$$

Since it's impractical to calculate conditional probabilities of long sentences, one can use Markov assumption and approximate this to a bigram model:

$$P('There was heavy rain') \approx \quad (5)$$

$$P('There')P('was' | 'There')P('heavy' | 'was') \quad (6)$$

$$P('rain' | 'heavy') \quad (7)$$

Seeing uni-gram and bi-gram, it came to our understanding that by calculating probabilities, we can calculate frequency of each and every word, better understand phrases that occur simultaneously and assigning specific meaning to these

phrases. So we implemented an N-gram model, taking variable N to be the maximum length of any sentence present in both the datasets of hindi and marathi. Then this N-gram model counts how often word sequences occur in corpus text and then estimated the probabilities with which a word can appear in the translation from marathi to hindi.

By calculating the most frequents phrases via uni-gram and bi-gram, we create a dictionaries like

- `n_gram_prob_mar` for marathi corpus
- `n_gram_prob_hi` for hindi corpus

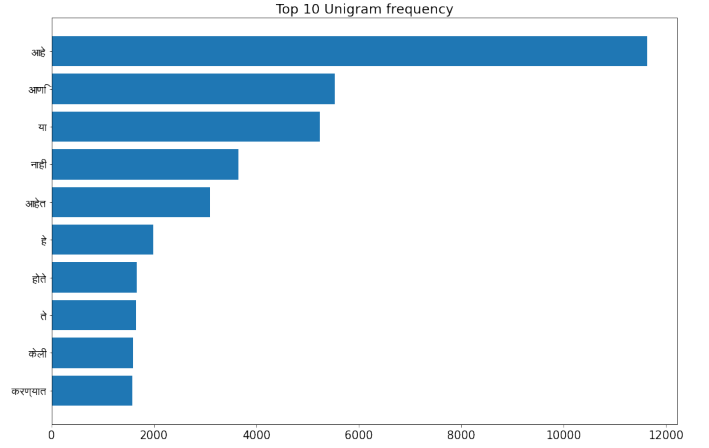


Fig. 4. Uni-gram for marathi corpus

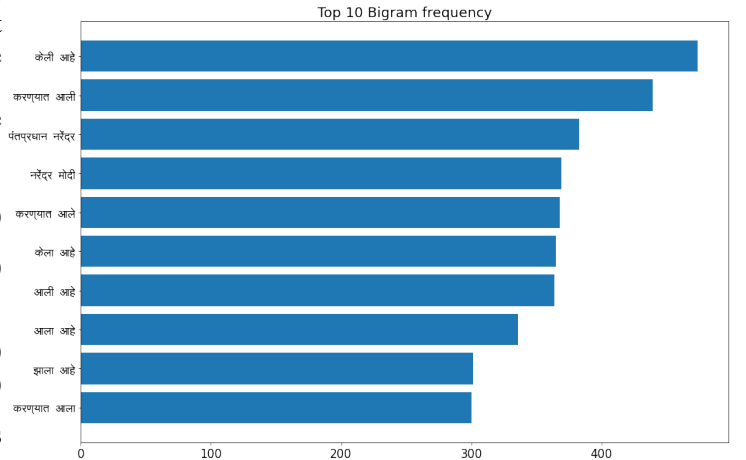


Fig. 5. Bi-gram for marathi corpus

We calculate frequency of uni-gram and bi-gram through consecutive loops.

Here in figure 4 and figure 5, it represents the uni-gram and bi-gram of Marathi corpus through the help of `calc_ng_pro_mar` function.

Here in figure 6 and figure 7, it represents the uni-gram and bi-gram of hindi corpus through the help of `calc_ng_pro_hin` function

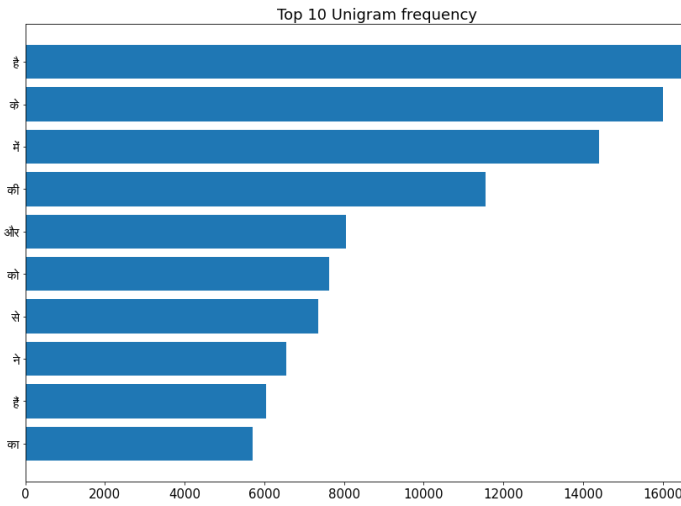


Fig. 6. Uni-gram for hindi corpus

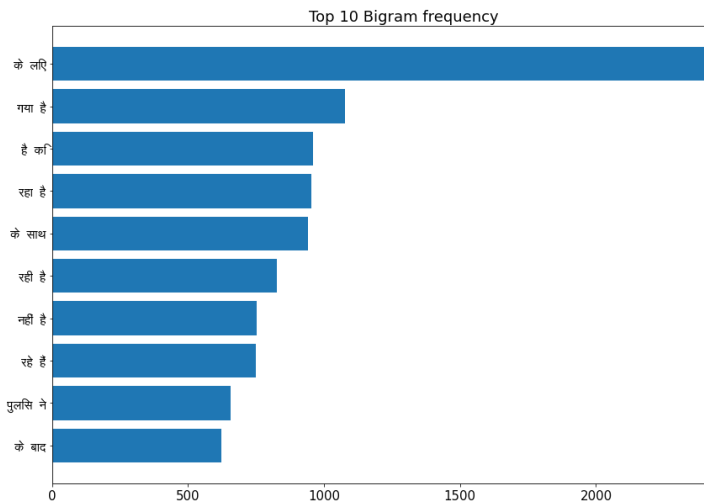


Fig. 7. Bi-gram for hindi corpus

C. POS Tagging

Part-of-speech (POS) tagging is a popular Natural Language Processing process which refers to categorizing words in a text (corpus) in correspondence with a particular part of speech, depending on the definition of the word and its context.

In corpus linguistics, Parts-Of-Speech or POS tagging, also called grammatical tagging is the process of marking up a word as nouns, verbs, adjective, adverb, etc. Once performed by hand, POS tagging is now done in the context of computational linguistics, using algorithms which associate discrete terms, as well as hidden parts of speech, by a set of descriptive tags. POS-tagging algorithms fall into three distinctive groups: rule-based, stochastic and hybrid.

Part-of-speech tagging is harder than just having a list of words and their parts of speech, because some words can represent more than one part of speech at different times, and because some parts of speech are complex or unspoken. This is not rare—in natural languages (as opposed to many artificial

languages), a large percentage of word-forms are ambiguous. For example, even "dogs", which is usually thought of as just a plural noun, can also be a verb: The sailor dogs the hatch.

Correct grammatical tagging will reflect that "dogs" is here used as a verb, not as the more common plural noun. Grammatical context is one way to determine this; semantic analysis can also be used to infer that "sailor" and "hatch" implicate "dogs" as 1) in the nautical context and 2) an action applied to the object "hatch" (in this context, "dogs" is a nautical term meaning "fastens (a watertight door) securely"). So, a Rule based POS tagger cannot give correct results always, so there is a need to consider the probability for each word and probability of the word to be a specific part of speech. This is the reason, that motivates the use of Stochastic method for developing Parts of Speech tagger.

Steps in POS tagging

- Obtain a tagged data set for training.
- Using Hidden Markov Model to identify Transmission and Emission Probabilities
- Apply Viterbi Algorithm on Testing Data Set.
- Output the tagged sequence with the highest probability.

We have obtained the source data set containing Marathi and Hindi Sentences with tagged POS from Marathi and Hindi sections of the universal dependencies corpus. The source corpora, documentation and credits can be found at <http://universaldependencies.org>

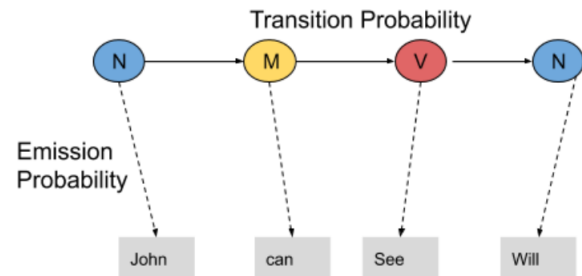


Fig. 8. Language distribution

Stochastic (Probabilistic) tagging: A stochastic approach includes frequency, probability or statistics. The simplest stochastic approach finds out the most frequently used tag for a specific word in the annotated training data and uses this information to tag that word in the unannotated text. But sometimes this approach comes up with sequences of tags for sentences that are not acceptable according to the grammar rules of a language. One such approach is to calculate the probabilities of various tag sequences that are possible for a sentence and assign the POS tags from the sequence with the highest probability. Hidden Markov Models (HMMs) are probabilistic approaches to assign a POS Tag.

Hidden Markov Model:

- A finite set of states

Here, for our problem, the states can be various parts of speech like noun, verb, adjectives, adverb, etc.

- A sequence of observations
(terminals i.e the words in our sentence)
- Transition probability
defined as the probability of a state 's' appearing after observing u and v in the sequence of observations
- Emission probability
defined as the probability of making an observation x given that the state was s.

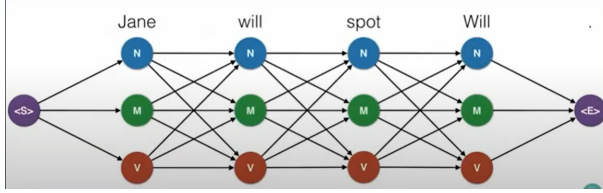


Fig. 9. POS tagging

Hidden States :

Let's consider the example of sunny day or rainy day. Given 7 days, a day can be sunny or rainy and we are given the probabilities of mood of a person to be happy or sad. Here we consider Hidden Transitions. We are given probabilities of a sunny day to be followed by a rainy day and rainy to be followed by sunny and so on. These states are called transition state. And direct observation which we get based on person's mood are called sequence of observations.

Transition Probability : The probability related to hidden state.

Emission probability : The probability related to visible states.

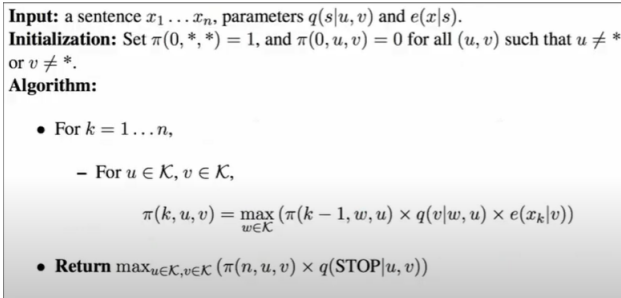


Fig. 10. Pseudocode for Viterbi

Viterbi Algorithm :

Take a sample sentence "Jane will spot Will." In NLP, whenever we consider a sentence for pos tagging, we append the start of sentence with "start" tag and end of the sentence with "end" tag. There are various paths which can be followed to obtain a tag sequence. So we can say that for very large sentence the brute force method will be hopelessly inefficient.

So, Instead of brute force approach, we will see that we can find the highest probable tag sequence efficiently using a dynamic programming algorithm known as the VITERBI algorithm. It basically stores the cost between any two nodes. So we don't have to calculate the cost every single point. Suppose I'm consider a path from NOUN followed MODEL AUXILIARY VERB and ADJECTIVE followed by NOUN.

So I calculate it. Now again if I want to know the cost for NOUN followed by MODEL AUXILIARY VERB and MODEL AUXILIARY VERB followed by VERB, I don't have to recalculate the cost for NOUN followed by MODEL AUXILIARY VERB, I will already have the probability stored in the form of a dynamic programming table.

Fig 10 is the pseudocode of the Viterbi Algorithm. What it essentially does is, it multiplies the previous probabilities which we have obtained so far with the emission probability of the current sequence of observations. The maximum of previous and current one is taken to be the most efficient one.

Results of developed POS Tagger :

- Training the model :

Using the POS tagged data from <http://universaldependencies.org> for training of POS tagger, the following results are obtained. Marathi POS tagging model :

```
start~tag~SYM:0.0001419446415897799858055358410
SYM~tag~NN:0.2209523809523809523809523810
NN~tag~CC:0.02532588454376163873370577281
CC~tag~NN:0.2798634812286689419795221843
NN~tag~NN:0.2959031657355679702048417132
NN~tag~VM:0.3148975791433891992551210428
VM~tag~VAUX:0.3115540789959394610557401255
VAUX~tag~DEM:0.008973080757726819541375872383
DEM~tag~NN:0.8275862068965517241379310345
NN~tag~NNPC:0.04152700186219739292364990689
NNPC~tag~NNP:0.7114989733059548254620123203
NNP~tag~PRP:0.2103064066852367688022284123
PRP~tag~NNP:0.05806451612903225806451612903
NNP~tag~NN:0.2172701949860724233983286908
NN~tag~QO:0.004096834264432029795158286778
QO~tag~NN:0.546875
VM~tag~PUNC:0.08748615725359911406423034330
PUNC~tag~SYM:0.07092198581560283687943262411
SYM~tag~DEM:0.1009523809523809523809523810
NNPC~tag~NNPC:0.2854209445585215605749486653
PRP~tag~VM:0.1435483870967741935483870968
VM~tag~PUNC:0.2321889996308600959763750461
start~tag~JJ:0.0001302592158395206460857105640
JJ~tag~NNPC:0.006914433880726015557476231634
PRP~tag~QC:0.03306451612903225806451612903
OC~tag~NN:0.6945525291828793774319066148
```

Fig. 11. Transmission probability

Fig. 11 and Fig. 12 represent the developed model itself. Fig. 11 shows the data stored after counting the transmission probability between two POS tags for e.g. NN tag CC:0.0253258845437616.... So this represents the probability of a word having POS tag NN followed by a word having tag CC is 0.0253258845437616... In Fig. 12 we get to see the Emission probability of a word having a particular POS tag for e.g अव-हइवअदअनअ_NN:0.0003724394785847299813780260708 So this represents that the probability of having 'NN' tag to the word अवहइवअदअनअ is 0.0003724394785847299813780260708.

Fig. 13 represents the input for the testing. While testing, the **Viterbi Algorithm** is used to give the POS tags to testing data. This algorithm gives output in the form of all the words of the sentences are tagged with some POS tag which is shown in Fig. 14.

Fig. 16 shows the accuracy of the tagger after verifying it from Fig. 15 which contains actual expected output for input

Fig. 12. Emission probability

Fig. 14. Marathi Test Output

Fig. 15. Marathi Expected Output

Fig. 16. Accuracy of Marathi POS tagger model on test data

We have achieved a good accuracy of around 59.61% on our POS tagging approach, and our model performs at an accuracy of 66.83% on the training data. We have also successfully implemented an LSTM model for our Indic data

```

Epoch 41/50
271/271 [=====] - 18s 66ms/step - loss: 0.5549 - accuracy: 0.6334
Epoch 42/50
271/271 [=====] - 18s 67ms/step - loss: 0.5432 - accuracy: 0.6379
Epoch 43/50
271/271 [=====] - 18s 67ms/step - loss: 0.5323 - accuracy: 0.6419
Epoch 44/50
271/271 [=====] - 18s 67ms/step - loss: 0.5216 - accuracy: 0.6465
Epoch 45/50
271/271 [=====] - 18s 67ms/step - loss: 0.5114 - accuracy: 0.6501
Epoch 46/50
271/271 [=====] - 18s 67ms/step - loss: 0.5013 - accuracy: 0.6541
Epoch 47/50
271/271 [=====] - 18s 67ms/step - loss: 0.4919 - accuracy: 0.6579
Epoch 48/50
271/271 [=====] - 18s 67ms/step - loss: 0.4824 - accuracy: 0.6616
Epoch 49/50
271/271 [=====] - 18s 67ms/step - loss: 0.4734 - accuracy: 0.6651
Epoch 50/50
271/271 [=====] - 18s 67ms/step - loss: 0.4649 - accuracy: 0.6683

```

Fig. 17. Accuracy of model on training data

```

Test sentence: 2
sentence: मी त्याच्याशी एकट्याने बोलेन
original translate: मैं अकेले उससे बात करूंगा।
predicted Translate: मैं पिछले हफ्ते अकेला नहीं था।
=====
Test sentence: 3
sentence: तू आता एकटी आहेस का
original translate: क्या आप अभी अकेले हैं
predicted Translate: क्या आप अब अकेले हैं
=====
Test sentence: 4
sentence: त्यांचे लांब केस आहेत
original translate: उसके लंबे बाल है।
predicted Translate: वह लंबे बाल है।
=====
Test sentence: 5
sentence: त्या कोण
original translate: वह कौन है
predicted Translate: वो कौन है
=====
Test sentence: 6
sentence: माझा दादा शिक्षक आहे
original translate: मेरा बड़ा भाई एक शिक्षक है।
predicted Translate: मेरा भाई एक शिक्षक है।
=====
Test sentence: 7
sentence: ती आपली शाळा आहे
original translate: वह हमारा स्कूल है।
predicted Translate: वह मेरा प्रेमी है।

```

Fig. 18. Accuracy of model on training data

NMT model. Some limitations that we faced were due to long sentences and less and sparse training data that led to some hindrances in recording even better accuracy.

REFERENCES

- [1] Laskar, S. R., Khilji, A. F. U. R., Pakray, P., Bandyopadhyay, S. (2020, November). Hindi-marathi cross lingual model. In Proceedings of the Fifth Conference on Machine Translation (pp. 396-401).
- [2] Ojha, A. K., Rani, P., Bansal, A., Chakravarthi, B. R., Kumar, R., McCrae, J. P. (2020, November). NUIG-Panlingua-KMI Hindi-Marathi MT Systems for Similar Language Translation Task@ WMT 2020. In Proceedings of the Fifth Conference on Machine Translation (pp. 418-423).
- [3] Mujadia, V., Sharma, D. M. (2020, November). Nmt based similar language translation for hindi-marathi. In Proceedings of the Fifth Conference on Machine Translation (pp. 414-417).
- [4] Pal, S., Zampieri, M. (2020, November). Neural machine translation for similar languages: The case of indo-aryan languages. In Proceedings of the Fifth Conference on Machine Translation (pp. 424-429).
- [5] Kumar, A., Baruah, R., Mundayiya, R. K., Singh, A. K. (2020, November). Transformer-based Neural Machine Translation System for Hindi-Marathi: WMT20 Shared Task. In Proceedings of the Fifth Conference on Machine Translation (pp. 393-395).
- [6] Kumar, A., Baruah, R., Mundayiya, R. K., Singh, A. K. (2020, November). Transformer-based Neural Machine Translation System for Hindi-Marathi: WMT20 Shared Task. In Proceedings of the Fifth Conference on Machine Translation (pp. 393-395).
- [7] Laskar, S. R., Pakray, P., Bandyopadhyay, S. (2019, August). Neural machine translation: Hindi-Nepali. In Proceedings of the Fourth Conference on Machine Translation (Volume 3: Shared Task Papers, Day 2) (pp. 202-207).
- [8] Przystupa, M., Abdul-Mageed, M. (2019, August). Neural machine translation of low-resource and similar languages with backtranslation. In Proceedings of the Fourth Conference on Machine Translation (Volume 3: Shared Task Papers, Day 2) (pp. 224-235).
- [9] Manger, K. M. (2014). Translation Divergences in Hindi-Nepali Machine Translation. Language in India, 14(5).
- [10] Liu, Y., Sun, C., Lin, L., Wang, X. (2016). Learning natural language inference using bidirectional LSTM model and inner-attention. arXiv preprint arXiv:1605.09090.