# St. Vincent Pallotti College of Engineering & Technology
## Department of Information Technology
### Academic Session: 2020-2021

## Word Phrase Alignment

**Name of the students: Gaurav Bhisikar, Aditya Dale, Sagar Malik**
**Name of the Guide: Prof. Pooja Walke & Indutry Mentor: Dr. Soma Paul and Dr. Vineet Chaitanya**
**Alumni Mentor: Ms. CBS Manasvini**

**ABSTRACT** : The main objective of this project is to develop software tools and a scheme for building word or phrase aligned corpus of English and Hindi optimally using existing open resources. Current NMT models are able to translate, but current NMT models are not able to generate correct alignment of translated Hindi sentence. So this project focuses on identifying the relation between words of both English and Hindi sentences and based on this relation to generating a Hindi Translated sentence that is correctly aligned.
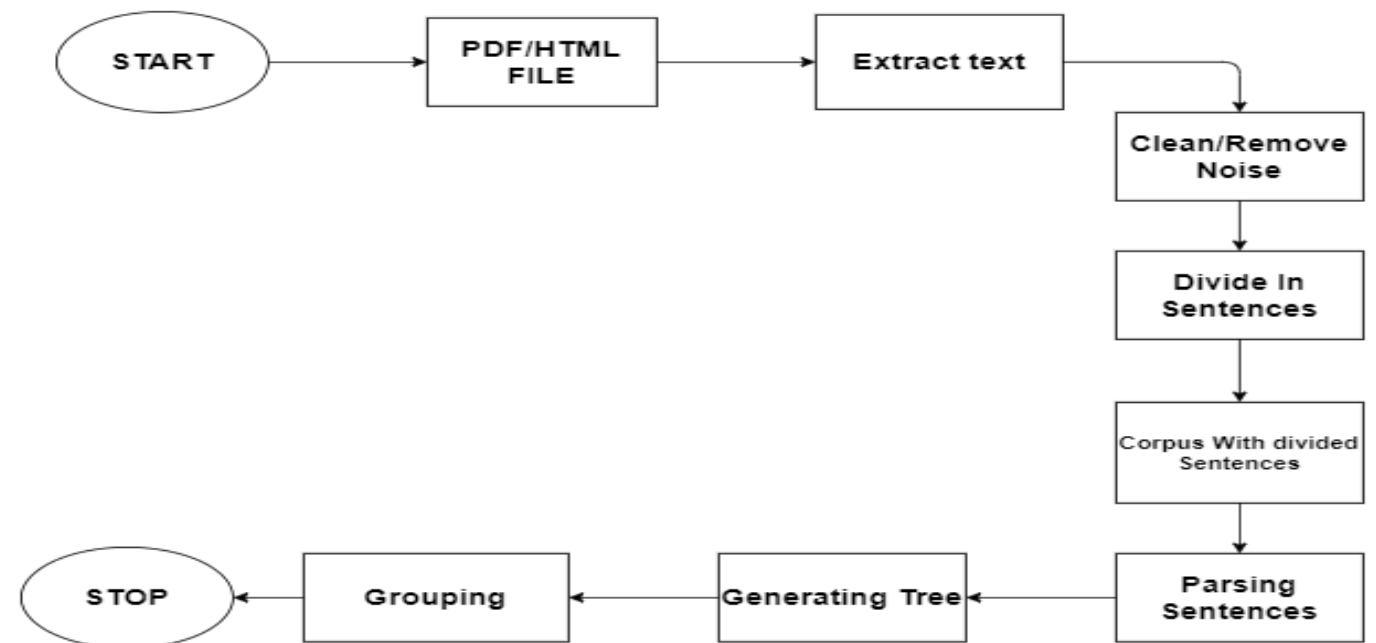
## INTRODUCTION:

This is a research project in the domain of NLP, at Anusaarka Lab under LTRC center at IIIT Hyderabad. This project aims to modify the current method of training deep learning model from sentence level to word/phrase level. The Current method of training deep learning model for translation uses parallel corpus which is on sentence level. But this method requires a very large corpus to train the model which makes it difficult to get expected result from smaller set of data. So if we can convert the parallel corpus to word/phrase structure level it will be easier to train and get better result from a deep learning model.

So comparatively model will require small sized corpus for training, eventually it will reduce the time to train a model, as size of data is reduced. It will also give more human control over traditional model.

In this project we are working on the possibility of automating the task of alignment of words/phrases of translated sentence in order to map these words/phrases with the words/phrases with same meaning in source sentence i.e. to create a parallel corpus with word/phrase level structure.

## SYSTEM ARCHITECTURE:



## MODULES:
1. Parser(ACE & STANFORD)      2. LUCY Tree-bank
3. Sentence Grouping      4. Noise Remove

## APPLICATIONS:
- Used in Natural Language Processing(NLP).
- Used in Automated Machine Translators.

# IMPLEMENTATION:



Fig. A) Clean Corpus

It is the output of program RemoveNoise.py, which is used to remove noises from the corpus.



Fig. B) Program Sent2table.py

It is the Snapshot of program Sent2table.py.



Fig. C) ) Output of Sent2table.py

It is the Output of Sent2table.py, it extracts information from parse tree.
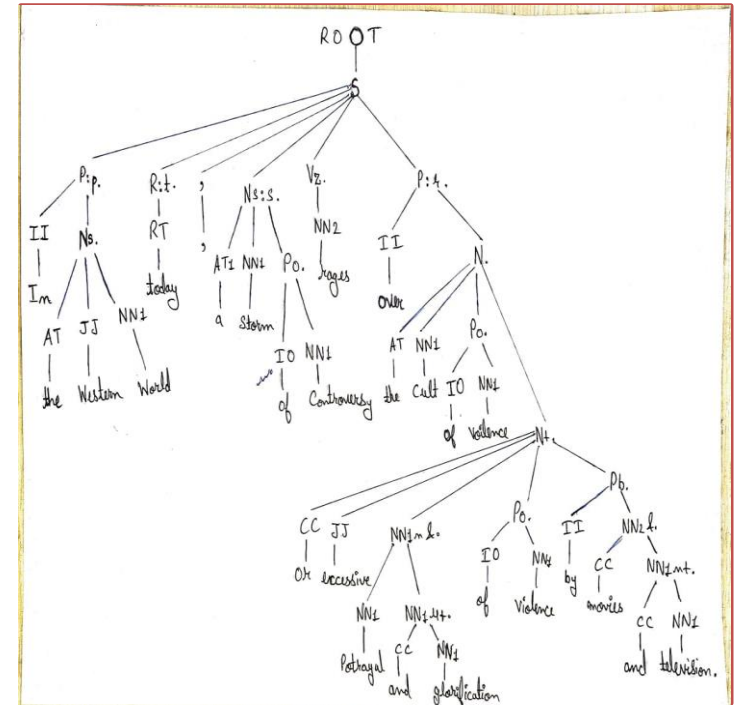


Fig. D) Manual Parse Tree

It is a parse tree we have drawn manually with the help of LUCY tree-bank.

# REQUIREMENTS:

**Hardware :** Processor - Intel i3 or higher , RAM: 4 GB or higher , Disk Space: 1 GB of free disk

**Software   :** Windows XP or later, Ubuntu , Python 3.6 or above , Java 8 or above

# CONCLUSION:.

The work done till the grouping of sentences revealed that the parse generated by Stanford parser might be incorrect, which directly affects the process of grouping. The main goal of the project is to develop a tool/software which performs automatic word alignment of parallel corpus of different language. The work done by us during this project will be a stepping stone for others to carry this research work to achieve the ultimate goal.

# REFERENCES:

[1] By Xin Lian, Kshitij Jain, Jakub Truszkowski, Pascal Poupart, and Yaoliang Yu (2020) "Unsupervised Multilingual Alignment using Wasserstein Barycenter"

[2] [Alaux et al., 2019] Jean Alaux, Edouard Grave, Marco Cuturi, and Armand Joulin. Unsupervised Hyper-alignment for Multilingual Word Embeddings. In ICLR, 2019.