



# **PROJECT IDEA :**

## **Sentiment Analysis for Code-Switched Languages**

Deepti Mittal - 193050025  
Ankita Singh - 19305R002



# Motivation & Idea

1. Sentiment Analysis is a hot topic of research due to availability of opinions of users on social media platforms.
2. Users post comments/feedback as per their ease :
  - a. *"I love this mobile".*
  - b. *"This restaurant is nice. Khana lazawab hai yaha."*
3. In the second example, a user posted his opinion by using a **code-mixed language**, in this case Hindi and English together, popularly known as HingLish
4. We wish to analyse the sentiment of users when feedback is given in a code-switched language.
5. For the purpose of our project, we will consider code-switching in Hindi and English.



# Problem Statement

- Statement: Given tweets in Code Switched Language, Output the sentiment expressed. Sentiment can be positive, negative or neutral.
- Example:
  - **Input:** *"bholy bhayaa. Ufffff dil jeet liya ap ne. Love you imran bhai. Mind blowing ap ki acting hai."*
  - **Output:** Positive (1)

# DataSet Statistics

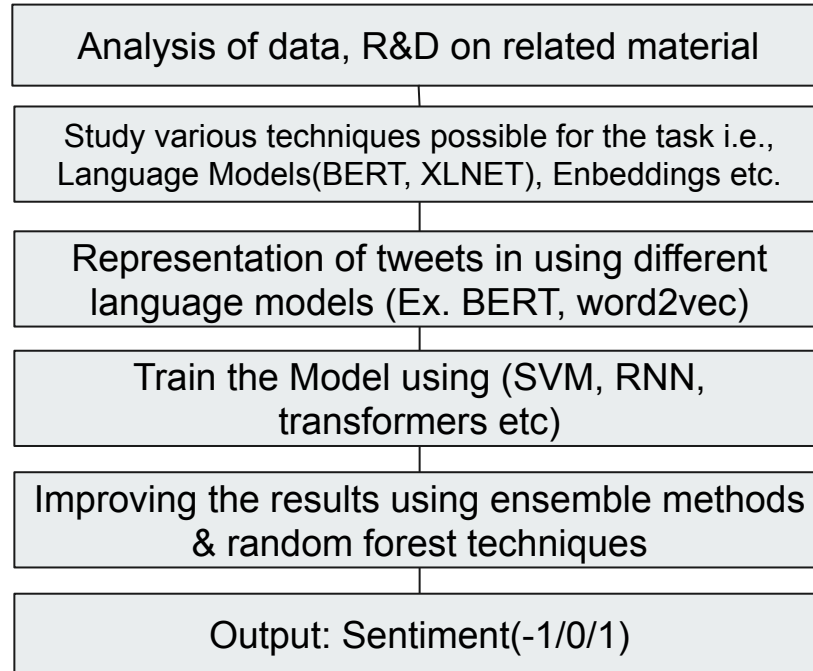
- Data is collected from tweets by using the list of tokens of hindi words released by Patra et al[2018].

Language	Split	Total	Positive	Neutral	Negative
Hinglish	Train	14,000	4,634 (33.10%)	5,264 (37.60%)	4102 (29.30%)
	Validation	3,000	982 (32.73%)	1,128 (37.60%)	890 (29.67%)
	Test	3,000	1,000 (33.33%)	1,100 (36.67%)	900 (30%)
	<b>Total</b>	<b>20,000</b>	<b>6,616 (33.08%)</b>	<b>7,492 (37.46%)</b>	<b>5892 (29.46%)</b>

- Example Tweets:

Hinglish	Congratulations <sub>ENG</sub> Sir <sub>ENG</sub> we <sub>ENG</sub> proud <sub>ENG</sub> of <sub>ENG</sub> you <sub>ENG</sub> ..o Aap <sub>HIN</sub> pr <sub>HIN</sub> pura <sub>HIN</sub> jakeen <sub>HIN</sub> hai <sub>HIN</sub> ..o aap <sub>HIN</sub> bohat <sub>HIN</sub> achaa <sub>HIN</sub> n home <sub>HIN</sub> minister <sub>ENG</sub> Honga <sub>HIN</sub> ..o )o (Congratulations sir we are proud of you.. We believe in you.. You will be a very good home minister.. )	Positive
Hinglish	Hostelite <sub>ENG</sub> k <sub>ENG</sub> naam <sub>HIN</sub> pe <sub>HIN</sub> dhabba <sub>HIN</sub> ho <sub>HIN</sub> tum <sub>HIN</sub> (you are a blot on the name of a hostelite)	Negative
Hinglish	Warm <sub>ENG</sub> up <sub>ENG</sub> match <sub>ENG</sub> to <sub>ENG</sub> theek <sub>HIN</sub> thaak <sub>HIN</sub> chal <sub>HIN</sub> ra <sub>HIN</sub> hai <sub>HIN</sub> (Warm up match is going fine)	Neutral

# Methodology

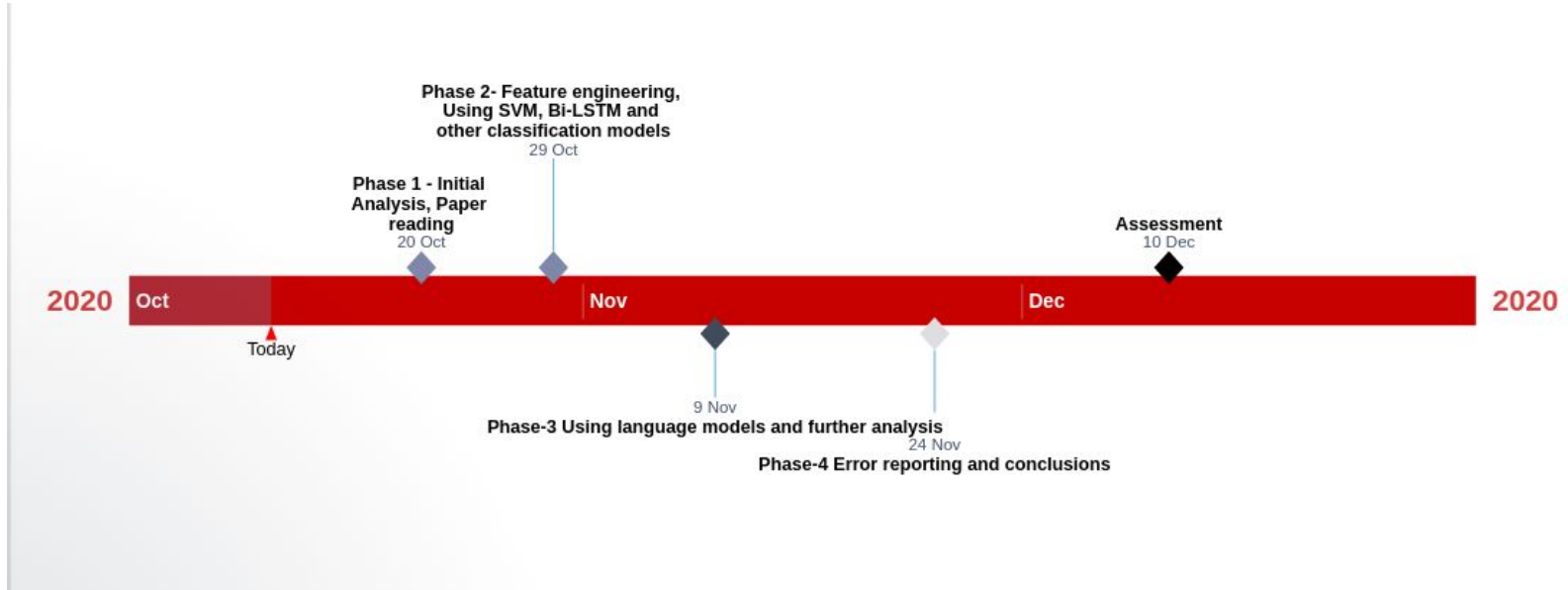




# Deliverables

1. Report the accuracy for language Identification (LId) in the given dataset.
2. Compare accuracies obtained using different models (SVM, RNN, Transformers) using different language models for the task of sentiment analysis.
3. Detailed Error Analysis for the work done.

# Timelines





## References

1. Patwa, Parth and Aguilar, Gustavo and Kar, Sudipta and Pandey, Suraj and PYKL, Srinivas and Gamb'ack, Bj'orn and Chakraborty, Tanmoy and Solorio, Thamar and Das, Amitava, ***SemEval-2020 Task 9: Overview of Sentiment Analysis of Code-Mixed Tweets***, Proceedings of the 14th International Workshop on Semantic Evaluation (SemEval-2020), December , ACL 2020
2. Braja Gopal Patra, Dipankar Das, and Amitava Das. 2018. ***Sentiment analysis of code-mixed indian languages: An overview of sail code-mixed shared task @icon-2017***. CoRR, abs/1803.06745.