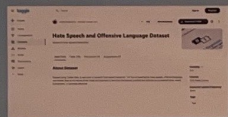


By: Varun Chilukuri, Rahul Nair, Anirudh Poruri, Shreya Shete

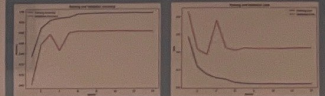
- Millions of people use social media platforms like Youtube and Twitter every day
- Such platforms oftentimes have hateful and offensive language, which can influence impressionable youth.



### Project Goals:

- To identify hateful/offensive comments so they can be removed on these platforms
  - Develop a user-friendly tool that utilizes our model's results to classify comments
- What makes this challenging?**
- Hate speech can be subtle and heavily context-dependent
  - Language is constantly evolving and new phrases are being coined all the time
  - We need to balance sensitivity and specificity—overly aggressive models may result in high false positives, impacting free speech and causing user dissatisfaction.

Accuracy on test data: 0.9067562222480774  
F1 Score: 0.9087254838943481  
Precision: 0.9028119905471882  
Recall: 0.8994792103767395

[illegible]

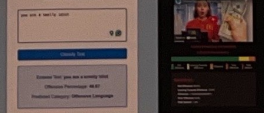
### Success Criteria

- The model accurately identifies hate speech, offensive language, and neutral speech
- We achieved a high accuracy and built two user-friendly tools that allow us to analyze individual comments or an entire YouTube video section

- Initial dataset improperly labeled hate speech values as "neutral"
- New dataset had more nuance—it differentiated between "offensive language" and "hate speech"

Figure 1 consists of two pie charts. The left chart, titled 'How often do you use the Internet?', shows the following distribution: 'Daily' (blue, 60%), 'Weekly' (orange, 20%), 'Monthly' (green, 10%), and 'Never' (red, 10%). The right chart, titled 'How often do you use the Internet?', shows the following distribution: 'Daily' (blue, 40%), 'Weekly' (orange, 30%), 'Monthly' (green, 20%), and 'Never' (red, 10%). An arrow points from the left chart to the right chart.

## Text Classifier



### Limitations

- Profanity and true hate speech is classified well
- Friendly and supportive text is classified well
- Strong real-world applications