# A Study of Hate Speech Detection and Classification

By Jerome HA
Emergent software laboratory

# Introduction : Preliminary research

Hate speech is a problem for online etiquette and safe navigation on social media

Objective: Detect and classify hate speech using advanced NLP techniques.

Dataset: HateXplain, a benchmark dataset for explainable hate speech detection, consisting of 20k samples from Twitter and Gab labeled as "hateful," "offensive," and "normal." [2]

Impact: Effective hate speech detection can contribute to safer online communities.

https://twitter.com/
https://gab.com/

# Introduction : Preliminary research

First semester :

Classifier of 3 classes :
-Normal
-Offensive
-Hate speech

Accuracy: 69%



3

# Proposed method

Raw text Data (HateXplain)

Fine-Tuning BERT (for HateXplain)
Binary classification: Normal,
hatespeech

Evaluate the level of Offensiveness

Combine both model and offensiveness
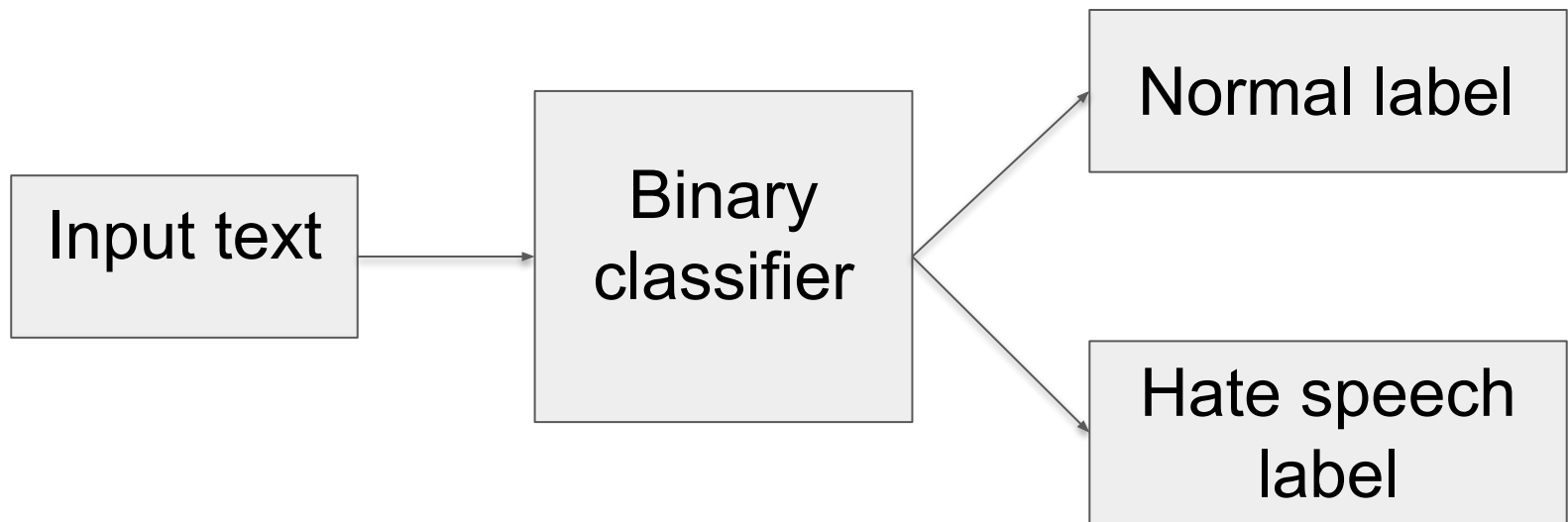level method

Model Evaluation for hatespeech
detection
Offensiveness level calculation

Schematic of the method for our detection model for hate speech and offensiveness

4

# First step: Hate Speech detection



Input text → Binary classifier → Normal label / Hate speech label

# Model Development

- Data Preprocessing:

Tokenize texts using BERT's tokenizer [1].

- Model Architecture:

BERT (Bidirectional Encoder Representations from Transformers):

Pre-trained transformer model for understanding language context.

Fine-tuned on HateXplain for hate speech detection.

# Model Development

Training Process:

Use Trainer class from Hugging Face's transformers library to manage training and evaluation.

Total Samples: 20,148

Split:

Training Data: 80% of the total data

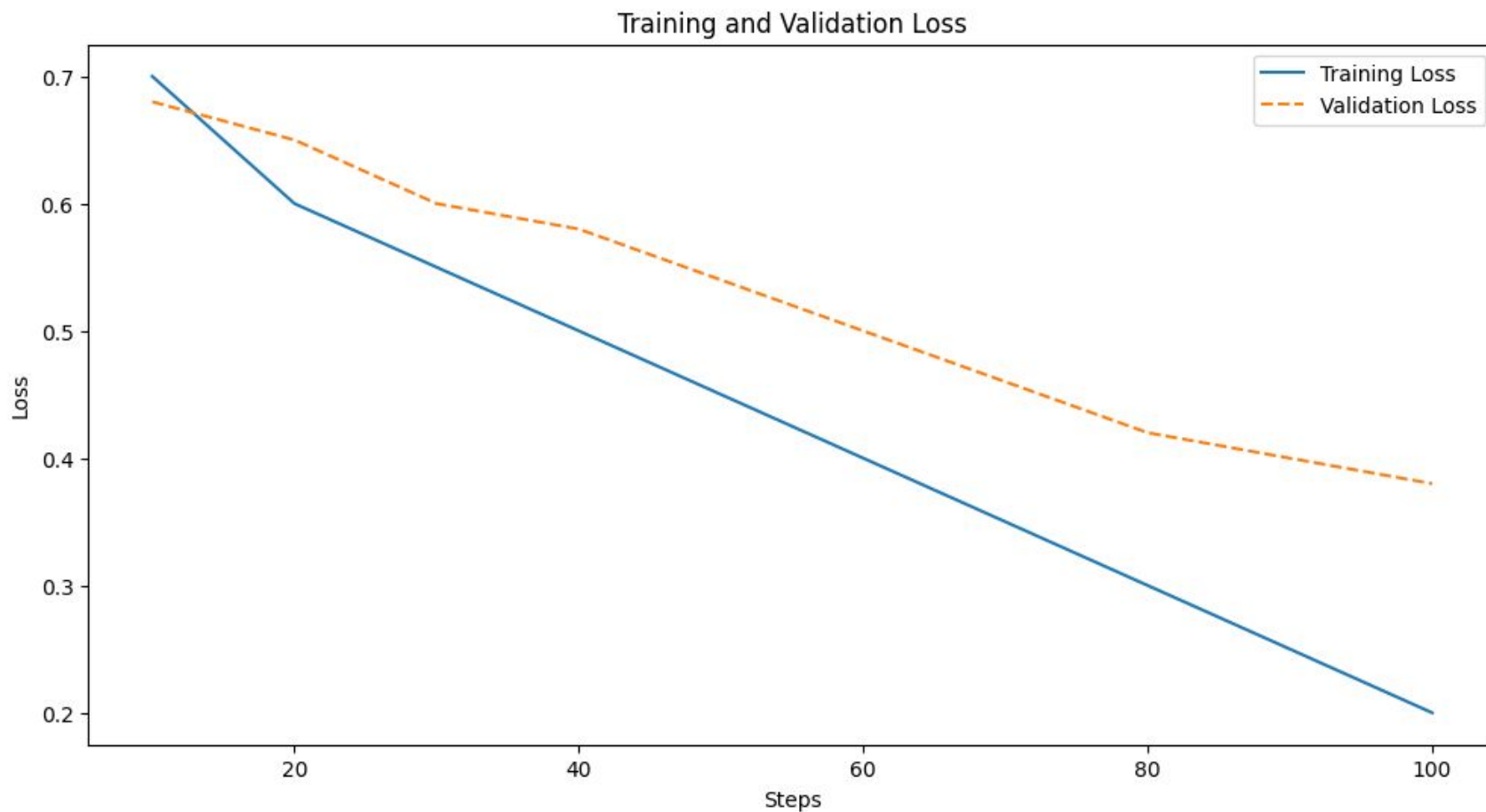Validation Data: 5% of the total data

Test Data: 15% of the total data
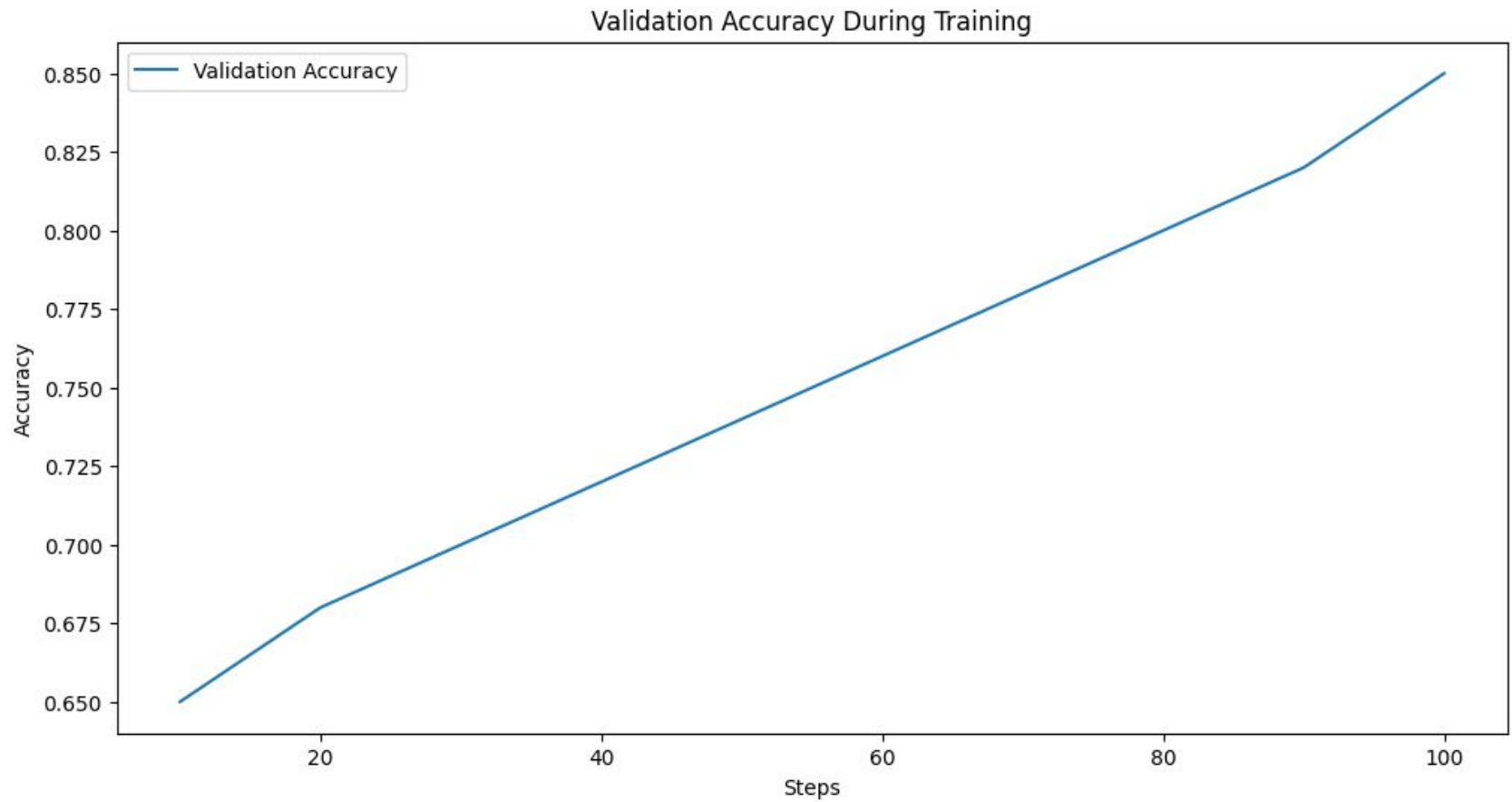
# Model Development table

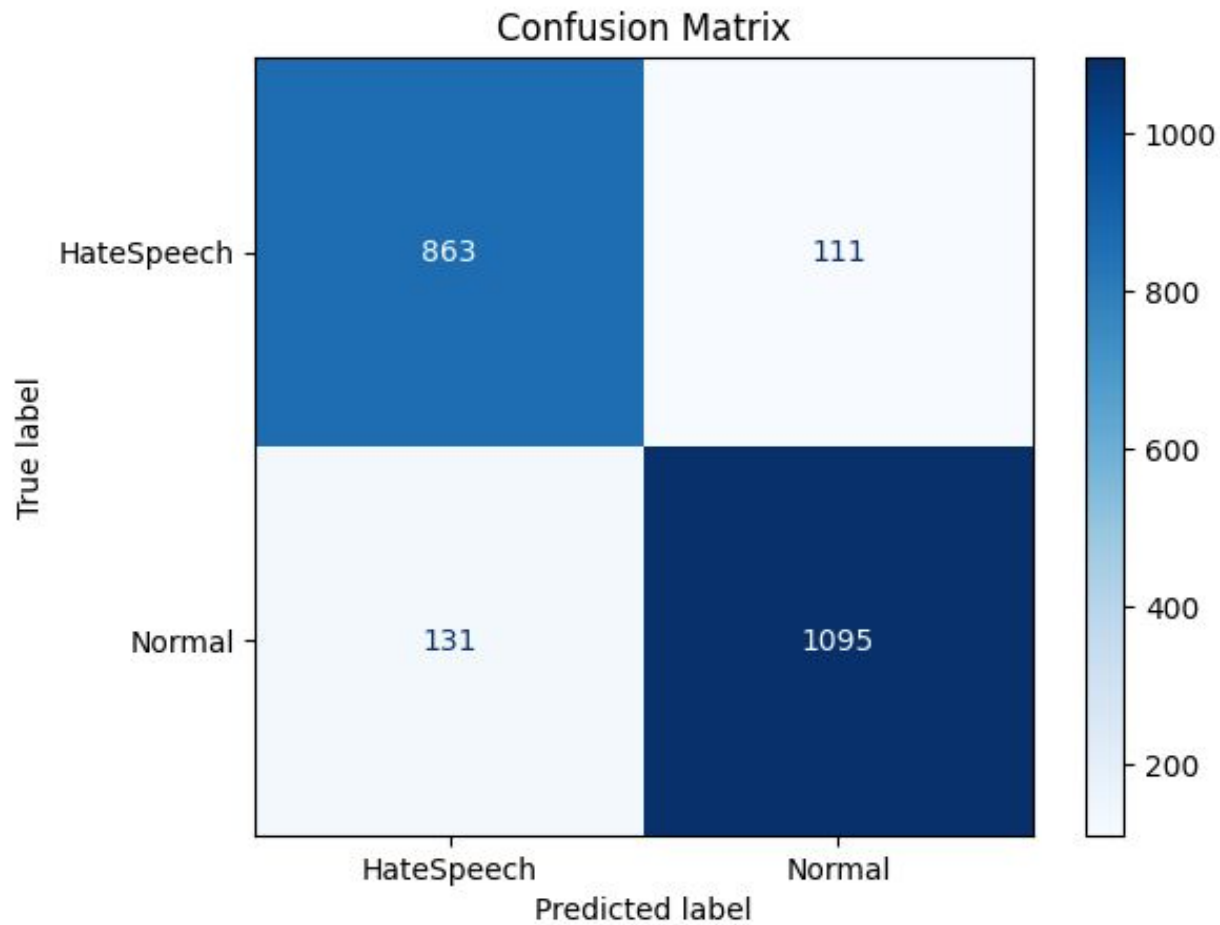| Hyperparameter | value |
|---|---|
| num_train_epochs | 3 |
| batch_size | 8 |
| Learning rate | Dynamic adjustment with warmup steps (initial: 5e-5) |
| warmup steps | 500 |
| BERT model | Bert-base-uncased pre-trained on BooksCorpus (11000 books), English Wikipedia (2.5 billion words) |

# Results



Training and Validation Loss

# Results



Validation Accuracy During Training

# Results

'eval_accuracy': 0.89

## Confusion Matrix

# Analysis on falsely classified data

**True Label: hatespeech, Predicted Label: Normal**

Example: me too but i still dis like jews and rather not have them follow me

--------------------------------------------------------------------------------

**True Label: normal, Predicted Label: Hatespeech**

Example: it was most certainly high even for secular jews during the ellis island days but i think the percentage of jews relative to their population that use welfare is higher than say asians or white catholics not sure though

# Second step: Offensiveness level calculation

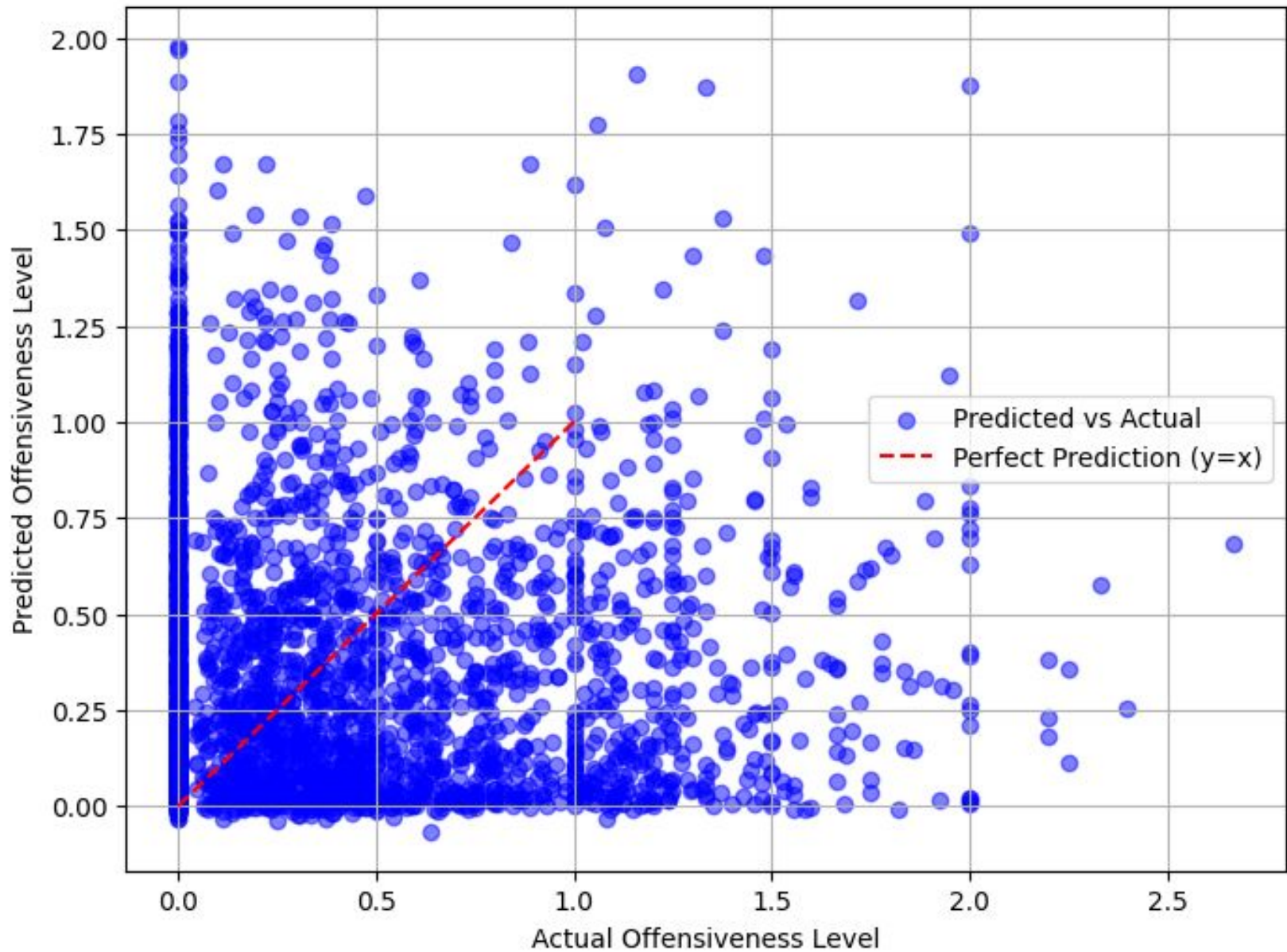Input text → Offensiveness level calculator → Offensive level score = 0,1,2…
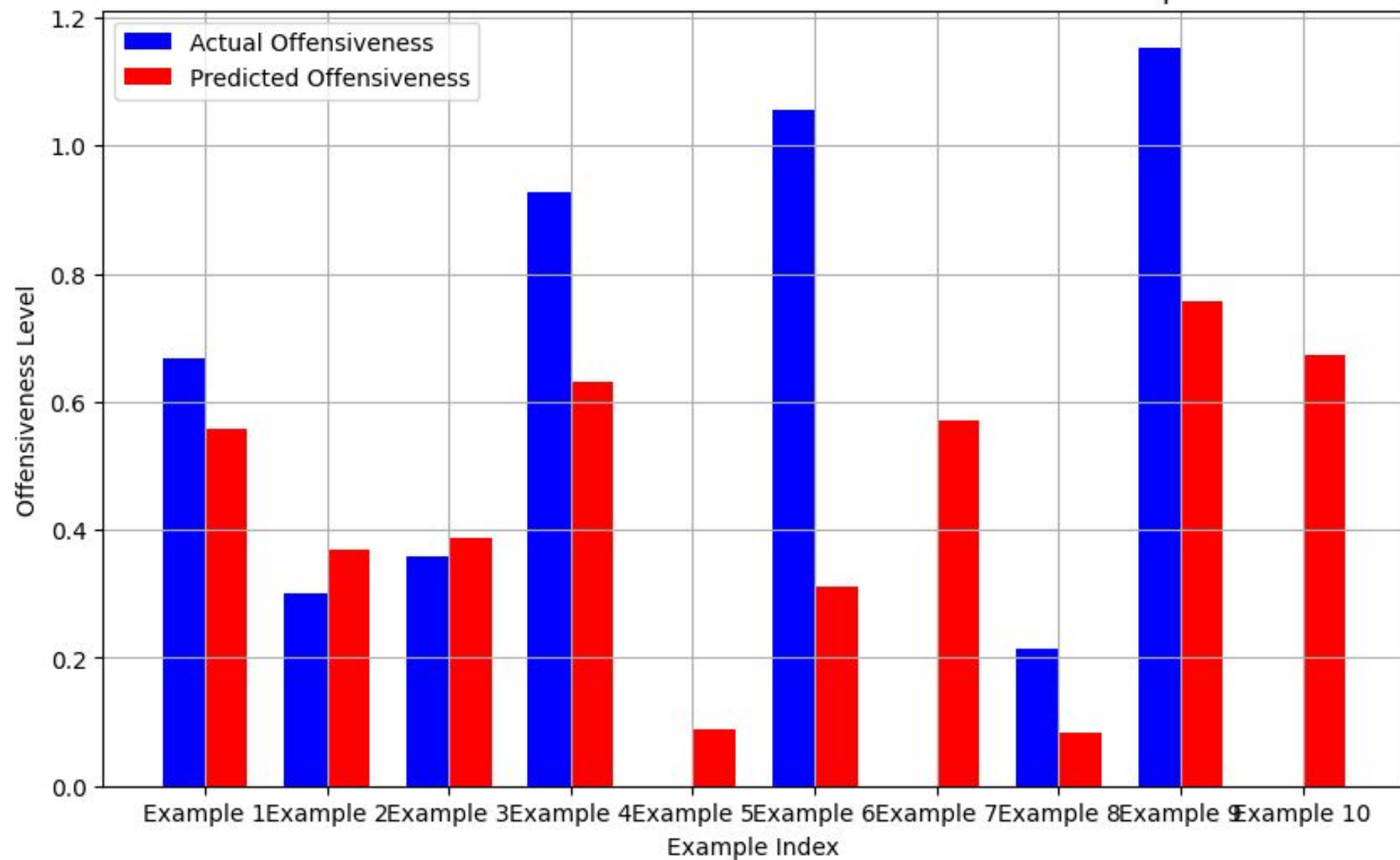
# First method

A regression model using the Offensive labels by the annotator data

Each sentence of the dataset has 3 annotators that labeled Rationales or specific parts of the text with normal, offensive or hate speech.

Actual vs Predicted Offensiveness Level

Actual vs Predicted Offensiveness Levels for 10 Random Examples

16

# Second method

A predefined dictionary (profanity.txt) Using Wiktionary list of profanity for English (3171 words)

Accounting for slang, abbreviations, and misspellings.

Total offensive words = Offensiveness Level

### Pages in category "English vulgarities"

The following 200 pages are in this category, out of 3,171 total.

(previous page) (next page)

**4**

- 4uck

**A**

- A2M
- abso-fucking-lutely
- absofuckinglutely
- abso-fuckin'-lutely
- abso-fuckin-lutely
- absofuckinlutely
- a damn sight
- AF

- arse over tit
- arse-up
- arsewash
- arseways
- arsewipe
- arsy varsy
- arsy versy
- artfag
- artfuck
- as all fuck
- as balls
- asf
- as fuck

https://en.wiktionary.org/wiki/Category:English_vulgarities

17

# Result examples

**True Label: <mark style="background:lightgreen">normal</mark>, Predicted Label: <mark style="background:lightgreen">normal</mark>**

Example: the whole drink water cure all <mark style="background:yellow">shit</mark> is fake i drink at least a gallon a day and have <mark style="background:yellow">shit</mark> skin hair brain

**Offensive score: 2**

-----------------------------------------------------------------------------------------

**True Label: <mark style="background:lightgreen">hatespeech</mark>, Predicted Label: <mark style="background:lightgreen">hatespeech</mark>**

Example: <user> she was screaming <mark style="background:yellow">nigger</mark> and that she wish she could kill all the <mark style="background:yellow">niggers</mark> in public at a cvs <mark style="background:yellow">bitch</mark>

**Offensive score: 3**

-----------------------------------------------------------------------------------------

**True Label: <mark style="background:lightgreen">normal</mark>, Predicted Label: <mark style="background:lightgreen">normal</mark>**

Example: done wit yo <mark style="background:yellow">nigga</mark> i just want his wallet

**Offensive score: 1**

18

# Potential bypass

**Obfuscation with Special Characters**

Example: f@ck, f*ck, fu&k, f#ck

**Intentional Misspellings**

Example: fuuuuuck, fucc, fuk, fasshole

**Leet Speak (1337)**

Example:h8, n1gger, b1tch

# Potential Improvements and future works

Some Offensive words are more offensive than others, find a way to include that information

Improve the accuracy of Hate speech detection.

Multilingual: Enhance the model's ability to handle multiple languages and mixed-language content.

Code-Switching: Try to include deliberate misspelling on offensive words

Think of other programs, how to improve the idea of the research

Use **generative IA** future works and compare it with my method

20

# Conclusion

Ongoing Research: Continuous improvement in model accuracy, interpretability, and adaptability.

Impact: Effective hate speech detection can contribute to safer online communities.

# References

[1] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of deep bidirectional transformers for language understanding," in Long and Short Papers. Minneapolis, Minnesota: Association for Computational Linguistics, Jun. 2019, pp. 4171–4186. [Online]. Available: https://aclanthology.org/N19-1423

[2] B. Mathew, P. Saha, S. M. Yimam, C. Biemann, P. Goyal, and A. Mukherjee, "Hatexplain: A benchmark dataset for explainable hate speech detection," in Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP). Association for Computational Linguistics, 2020, pp. 3451–3463.

https://en.wiktionary.org/wiki/Category:English_vulgarities