# AM-DETOX: Analyzing Amharic Text Detoxification Using Pre-trained Large Language Models

Abinew Ali Ayele[1,2], Seid Muhie Yimam[1]

[1]Language Technology Group, Department of Informatics, Universität Hamburg, Hamburg, Germany
[2] Bahir Institute of Technology, Bahir Dar University, Bahir Dar, Ethiopia

{abinew.ali.ayele, seid.muhie.yimam}@uni-hamburg.de

## Introduction

- Social media makes the spread of toxic content easier.
- Various mitigation strategies are employed.
- Most studies in low-resource languages focus on detection and classification tasks, fail to employ text detoxification
- Text detoxification or rewriting messages allows for the removal of toxicity from messages.
- Detoxification can protect vulnerable groups like children.
- LLMs were unable to detoxify toxic messages in low-resource languages, mainly "hallucinate".

## Research Questions

- **RQ1**: How effective are LLMs in detoxifying toxic content in low-resource languages such as Amharic?
- **RQ2**: The specific challenges in annotating toxic content and applying text detoxification techniques for the Amharic language?
- **RQ3**: How can the challenges be addressed?

## Data Collection

- **Source**: offensive labeled tweets from Ayele et al. (2022, 2023).
  https://github.com/uhh-lt/AmharicHateSpeech
- 3,120 tweets from both datasets
- Re-annotated for detoxification task



Fig. 1: Annotation GUI

## Annotation

- Customized POTATO: POrtable Text Annotation Tool
  https://github.com/davidjurgens/potato
- Pilot annotation by 3 experts: **125** tweets
- Main annotation: **2,995** tweets
- Total annotations:
  o 1,452: detoxifiable
  o 1,543: non-detoxifiable



Fig. 2: Text detoxification/ re-writing examples

## Experimental Results

- Am-ReBERTa achieved better performance on classification tasks

| Classifier | precision | recall | F1-score |
|---|---|---|---|
| Am-RoBERTa | **69.57** | 73.37 | **70.01** |
| Afro-xlmr-large | 69.07 | 73.53 | 67.75 |
| Bert-medium-amharic | 66.13 | 69.90 | 64.62 |

Table 1: Classification results in identifying detoxifiable texts from non-detoxifiable ones

- GPT4 and Shapely employed for identifying toxic words in text, where GPT4.0 mini outperformed.

| Toxic word detection performance | | Detoxification performance | |
|---|---|---|---|
| GPT-4o mini | SHAP | chrF2 | Sacrebleu |
| **76.33%** | 54.83% | 11.5 | 0.5 |

Table 2: Toxic words detection (Explainability) and the detoxification capabilities (Detoxifiability) of SHAP and GPT-4 models.
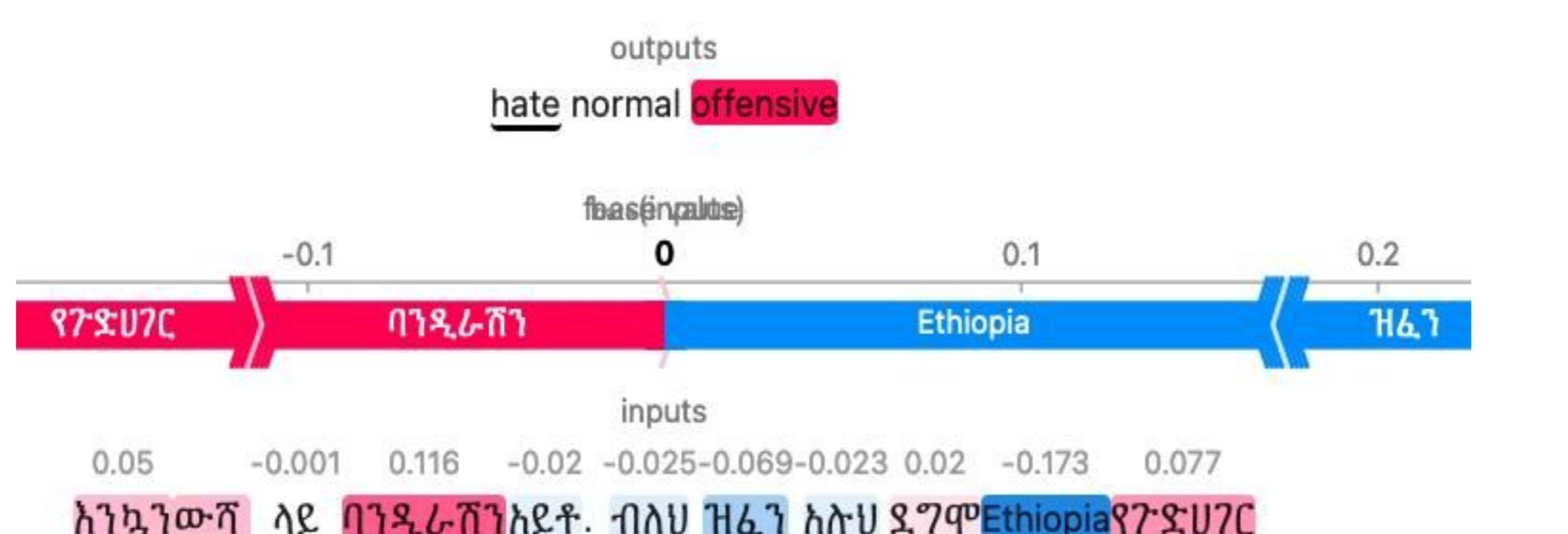


Fig.3: Output examples for explainable toxicity detection (SHAP)

- **GPT-4o-Mini** detoxification outputs significantly deviate from human expert annotations due to hallucinations.



Table 3: Comparison of Gpt4 detoxification outputs Vs human expert detoxification results



Table 4: Shap outputs of toxic words and detoxified examples.



Table 5: Shap outputs of toxic words and detoxified examples.

## Conclusion and Future Works

- Presented a new detoxification dataset for Amharic
- Conducted 3 tasks using transformer models: classification (detoxifiable or not), explainability (why a message is toxic) and detoxifiability (rewriting messages).
- GPT4 is better in Amharic text detoxification task despite its frequent hallucination challenges.
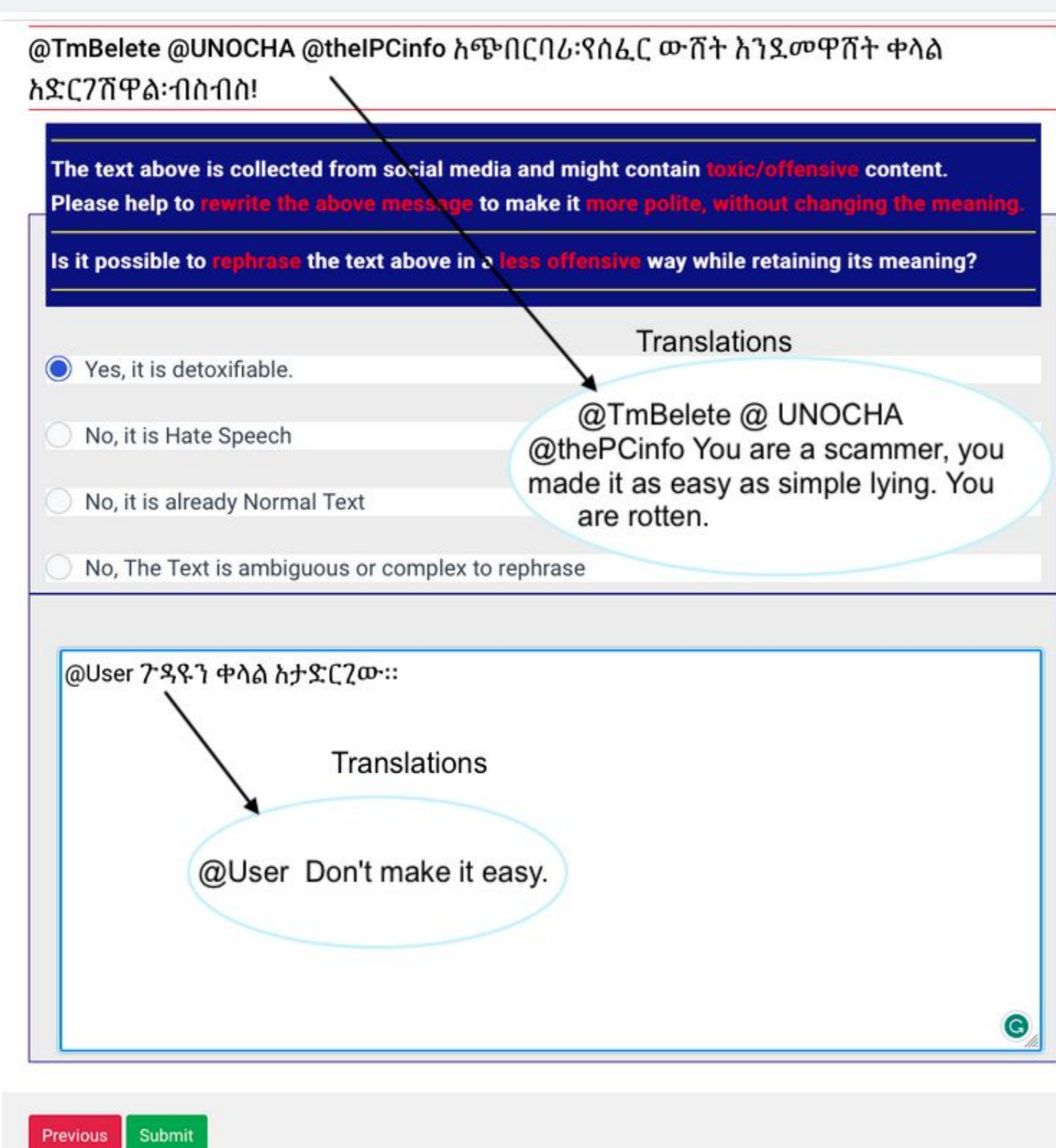- LLMs still struggle with detoxifying low-resource language messages, needs fine-tuning with more datasets