# Improving the Performance of Transformer-Based Models Over Classical Baselines in Multiple Transliterated Languages

**Fahim Ahmed[,1], Md Fahim[,1], Md Ashraful Amin, Amin Ahsan Ali and AKM Mahabubur Rahman**

Center for Computational & Data Sciences
Independent University, Bangladesh
Dhaka-1229, Bangladesh

**Abstract.**   Social media users express their feelings, experiences, ideas, and stories with little or no regard for the conventions of traditional grammar. Online discourse, by its very nature, is rife with transliterated text along with code-mixing and code-switching. Transliteration is heavily featured due to the ease of inputting romanized text with standard keyboards over native scripts. Due to its ubiquity, it is a critical area of study to ensure NLP models perform well in real-world scenarios. In this paper, we analyze the performance of various language models, Tiny Large Language models, TF-IDF and Bag-of-Words feature extraction-based classical ML models, as well as zero-shot classification with ChatGPT on romanized/transliterated social media text. We chose the tasks of sentiment analysis and offensive language identification and we carried out experiments for three different languages, namely Bangla, Hindi, and Arabic, for six datasets. To our surprise, we discovered across multiple datasets that the non-neural methods perform very competitively with fine-tuned transformer-based mono/multilingual language models, tiny large language models, and ChatGPT for classification tasks in transliterated text. These classical models train in seconds using only a fraction of the computing power, and thus the carbon footprint, required by language models. We demonstrate TF-IDF and BoW-based classifiers achieve performance within around 3% of fine-tuned LMs and could thus be considered as a strong baseline for transliterated text-based NLP tasks. Additionally, we investigated various mitigation strategies such as translation and augmentation via the use of ChatGPT, as well as Masked Language Modelling to dataset-specific pretraining for language models. Depending on the dataset and language, employing those mitigation techniques yields a 2-3% further improvement in accuracy and macro-F1 above baseline.

## 1   Introduction

Nowadays, the proliferation of social media has connected people worldwide across linguistic boundaries — facilitating a global exchange of ideas, opinions, and emotions. However, an important characteristic of such communication is the informal and "chaotic" way users express themselves online. These do not conform to traditional grammar ideals. For example, people do not follow the standardised spellings and rules of punctuation in online communication. Perhaps more so than in real life, code mixing and switching occurs

abundantly here — a speaker of multiple languages may switch between spontaneously within the same sentence, often melding different grammar systems altogether for convenience or style [8]. An equally prevalent phenomenon is the extensive use of transliterated text where words of one language are phonetically represented with the writing script of another.

Transliteration enables individuals to express themselves in their native tongue even when the keyboard or fonts are unavailable. In regions where non-Latin scripts dominate, such as South Asia and the Middle East, transliteration is especially prevalent due to the ease of input with a standard QWERTY keyboard. Additionally, it must be noted that transliterated text often carries cultural connotations and emotional signals that may not be fully captured through direct translation or standardization, thus adding layers of meaning and richness to online interactions [26]. However, noisy transliterated text presents a challenge for NLP practitioners. Models are often trained on corpora collected from a single language. However, for these models to be effective in the real world, they must perform well even with such imperfect/transliterated input [34].

There is a body of work that investigate the performance of LLMs on transliteration text. Aggarwal et al. [1] use transliteration to leverage the commonality between many Indian languages to elevate performance for low-resource languages. They have used Bi-LSTM. Purkayastha et al. [31] analyzed the use of multilingual language models and the universal transliteration tool *UROMAN*. On the other hand, Biradar et al. [7] propose a Transformer-based Interpreter and Feature extraction model on Deep Neural Network (TIF-DNN) for Hinglish text. [30] explored BERT-based techniques and translation (to English) based strategies in Dravidian code mixed language. [33] introduced character-based modeling for solving transliteration-based tasks. The recent popularity of transformer-based LLMs (which demand large amounts of memory and computational resources) overshadow the interest in exploring the use of classical machine learning models (which have significantly smaller memory and computational resource demands).

To address the limitations mentioned above, in this paper, we thoroughly investigate the performance of different machine learning and language models on transliterated text. We choose Bangla, Hindi, and Arabic as people from these languages widely use transliterated

---

[1] Equal Contribution.

texts in informal communications [12, 22, 2]. We use various non-neural classical machine learning techniques such as Logistic Regression, Random Forests, Support Vector Machines, and XGBoost as our performance baseline and inquire into the use of a multitude of Transformer-based Language Models (LM) for classification tasks in the transliterated text of multiple languages. Later, we explore mitigation strategies to further improve the performance. We explore the use of prompt-based Large Language Models, specifically ChatGPT, to both classify and also to aid in the training of the LM-based models. We delve into the prospect of using three different techniques: dataset augmentation, dataset translation, and dataset-specific pre-training. For the first two, we utilize the tremendous capabilities of GPT-3.5 Turbo in the form of ChatGPT — prompting the LLM to synthesize text data for augmentation and to translate transliterated text to English.

## 2    Related Work

Natural language processing (NLP) research on transliteration text recently gained much attention for analyzing a vast amounts of social communication. Aggarwal et al. [1] explored the use of transliteration to leverage the commonality between many Indian languages to elevate performance for low-resource languages. The 9 languages of the IndicNLP News Article Classification Dataset was transliterated to the Devanagari script and classification was carried out using a Bi-LSTM network.

Purkayastha et al. [31] analyzed the use of multilingual language models and the universal transliteration tool *UROMAN* on large-scale language adaptation on 14 low-resource languages. This transliteration-based technique that converted text to a common Roman script showed promise, especially for languages whose scripts were unfamiliar to the model. Guellil et al. [17] transliterated Arabizi from the Roman script to Arabic and utilized non-neural machine learning methods to carry out sentiment analysis.

Biradar et al. [7] proposed a Transformer-based Interpreter and Feature extraction model on Deep Neural Network (TIF-DNN) for Hinglish text. This architecture utilises an interpretation layer that produces text in the native Devanagari script from Hinglish text. It uses a language identification tool to tag each word with its language. Romanised Hindi words are transliterated back to the Hindi Devanagari script and English words are translated to the Hindi equivalent before being passed forwards for feature extraction and classification.

Jahan et al. [21] emphasized the use of extensive preprocessing when dealing with noisy transliterated and code-mixed Bangla text for classification using non-neural ML techniques adding the use of Google Translator API to correct spellings. [30] explored BERT finetuning techniques and translation (to English) based strategies in Dravidian code mixed language. [34] follows a similar approach for solving sentiment analysis tasks considering four different Indic languages in code mixed dataset. [33] introduced character-based modeling for solving transliteration-based tasks.

Inspection of the contemporary work on this topic exposes several key insights. We notice a recent trend towards an over-reliance on transformer-based LMs, which demand large amounts of memory and computational resources. There appears to be a lack of studies comparing the performance between training computationally efficient non-neural ML models and transformer-based LMs for transliterated text tasks.

## 3    Experimented Datasets

In this section, we describe the datasets that are used in our investigation. For each of the languages chosen for our investigation, we picked one sentiment analysis and one hate speech detection dataset; a total of six datasets. All datasets are in the Roman script. Code-mixing (often English utterances) is occasionally observed.

All datasets were cleaned to remove URLs and mentions. Repeated punctuations were tidied up into ellipses ('**...**') for commas and periods, or reduced to single punctuation for exclamations ('**!**') and question marks ('**?**'). Unnecessary whitespace was cleaned.

### 3.1    Bangla

The *Positive and Negative Corpus* [2] is a collection of 1,300 comments scraped from Facebook and YouTube, of which 647 are positive and 653 express negative sentiment [3]. We split these 80:20 into the train and test sets.template

*TB-OLID* [3] [32] contains 5,000 Facebook comments, of which 2,381 are offensive and 2,619 are non-offensive. The train and test sets contained 4,000 and 1,000 comments respectively. The selection process involved filtering out non-Latin script comments from the initial corpus of 100,000 comments, followed by the offensive-keyword-based search for offensive comments. The comments are manually annotated following the OLID hierarchical taxonomy - with a label for whether it is offensive or not, and one for the target (individual, group, untargeted). Also labeled is whether a comment contains code-mixing along with the transliteration.

### 3.2    Arabic

*TUNIZI* [4] [15] is a Tunisian Arabizi dataset that contains 4,838 positive and 4,372 negative comments scraped from YouTube. Topics cover sports, politics, comedy, art, music and TV shows. The dataset came preprocessed with all links, emojis, and punctuations were removed before being annotated by a group of 5 annotators.

*Offensive Language Detection in Arabizi* [5][5] contains 5,857 non-offensive and 1,526 offensive Arabizi social media texts from merging Arabizi content of four existing predominantly Arabic script datasets and unifying their labels.

For both datasets, we used an 80:20 split for train and test sets as there is no separate test dataset.

### 3.3    Hindi

The Hinglish sentiment dataset [6] [29] used consisted of 5,140 negative, 6,581 neutral, and 5,847 positive tweets. The 17,568 tweets were separated into a train set of 14,569 and a validation set of 2,999 tweets. Tweets were initially selected based off a list of keywords and then annotated by sixty annotators for whom Hindi was their first or their second language. Each tweet was annotated by two annotators and only kept if both labels matched. Due to the official SemEval2020 Task 9 test set not being available, we used the validation set in that capacity.

---

[2] https://data.mendeley.com/datasets/s6mtp2zzpc/3
[3] https://github.com/LanguageTechnologyLab/TB-OLID
[4] https://github.com/iCompass-ai/TUNIZI
[5] https://github.com/Imene1/Arabizi-offensive-language/
[6] https://github.com/singhnivedita/SemEval2020-Task9/

The Hinglish hate speech dataset [7] used contains 2,080 non-hate and 1,325 hateful user-generated tweets. The texts are labeled 'yes' or 'no' according to the presence of hateful content. As there is no separate test dataset, we use an 80:20 train-test split for our investigation.

## 4 Experimented Methods

### 4.1 ML Models

For doing text classification using ML models, we first convert the text into vector representations leveraging two widely used techniques: bag-of-words (BoW) and term frequency-inverse document frequency (TF-IDF). Let $D$ represent the dataset comprising $N$ documents, each with $M$ unique terms/words. The BoW representation of a document $d_i$ is denoted as $\text{BoW}(d_i) = [x_{i1}, x_{i2}, ..., x_{iM}]$, where $x_{ij}$ represents the occurrence frequency of term $j$ in document $i$. Similarly, TF-IDF transforms each document $d_i$ into a vector $\text{TF-IDF}(d_i) = [x_{i1}, x_{i2}, ..., x_{iM}]$, where $x_{ij}$ denotes the TF-IDF score of term $j$ in document $i$.

Following the vectorization process, these representations were subjected to classical non-neural machine learning (ML) models for classification tasks.

### 4.2 LM Fine-tuning

Let's denote a LM $f_\theta$ which takes a sentence as input and returns the contextual presentation $H$ where $H = f_\theta(S)$. Upon tokenizing an input sentence $S$, represented as $T = \{t_1, t_2, \ldots, t_n\}$, the LM generates contextual representations for each token by applying the function $f_\theta(S)$, resulting in a sequence denoted as $H = \{h_1, h_2, \ldots, h_n\}$. These representations encapsulate the unique meaning of each token within the context of the entire sentence.

However, for tasks such as classification, where a fixed-size representation of the entire sentence is required, we employ a two-layer Feed Forward Neural Network (FFN) on the contextual representation of [CLS] token, $h_{\text{CLS}}$. This network utilizes weight matrices $W_1$ and $W_2$, bias terms $b_1$ and $b_2$, and the Rectified Linear Unit (ReLU) activation function to process $h_{\text{CLS}}$ and generate a fixed-size representation $z$.

$$z = W_2 \cdot (\text{ReLU}(W_1 \cdot h_{[\text{CLS}]} + b_1)) + b_2 \tag{1}$$

For the character-based LMs and multilingual LMs, we follow similar fine-tuning strategies for doing the classification tasks.

### 4.3 TinyLLM Fine-tuning using LoRA and PEFT

Traditional fine-tuning of large language models (LLMs) involves significantly modifying the pre-trained model's parameters, which can be computationally expensive and time-consuming. PEFT (Parameter-Efficient Fine-Tuning) [25] offers a solution by adapting pre-trained models to new tasks with minimal changes to the original parameters. This significantly reduces training time and memory usage compared to traditional approaches. LoRA (Low-Rank Adaptation) [20] is a specific PEFT technique that introduces a more efficient way to capture the adjustments needed for fine-tuning. Instead of directly modifying all the pre-trained parameters, LoRA utilizes

a low-rank matrix. This matrix requires significantly fewer parameters to represent the task-specific adaptations, leading to substantial efficiency gains.

Let's denote the original pre-trained model parameters as $W$ which will be frozen during training. LoRA introduces a low-rank update, denoted by $\Delta W$, which captures the task-specific adjustments needed for fine-tuning. This low-rank update is further decomposed as the product of two trainable matrices, $A$ and $B$: $\Delta W = A \times B^T$. Here, $A$ with a shape of $d \times r$ and $B$ with a shape of $r \times d$ have a much lower rank (denoted by $r$) compared to the original dimension $d$ of the parameter matrix $W$. This means they require significantly fewer parameters to represent the necessary adjustments. The rows of matrix $A$ and the columns of matrix $B$ can be interpreted as capturing the task-specific adaptations applied to the original weight matrix $W$. Finally, the updated weight matrix $W'$ with LoRA is the summation of pretrained frozen matrix $W$ and task-specific fine-tuned matrix $\Delta W$

$$W' = W + \Delta W = W + AB^T$$

In essence, LoRA leverages a more compact representation (the low-rank matrices $A$ and $B$) to achieve fine-tuning, resulting in significant efficiency improvements compared to traditional fine-tuning methods that modify all the pre-trained parameters directly.

### 4.4 Prompting

Prompting in NLP refers to providing structured instructions or cues to guide the behavior of language models. It has become crucial as a means to utilize pre-trained large language models (LLMs) for solving downstream NLP tasks without the need for task-specific training. As pretrained LLMs are trained on a huge amount of data, they have the capability to understand different tasks [37, 27, 9].

Instead of fine-tuning on task-specific data, users can simply provide a prompt $P$ to the model specifying the task $T$, and the model generates a response $R$ based on that prompt. Prompt $P$ can be zero-shot, few-shot, and so on. A good design of prompt $P$ may improve the performance [23, 4]. This approach has gained popularity due to its efficiency and effectiveness in solving a wide range of tasks without task-specific training.

### 4.5 Experimented Models

The following models were considered in our experiments.

- **ML Models:** We consider Logistic Regression (LR), Support Vector Machine (SVM), Random Forest (RF), and XGBoost to analyze the behavior of ML models in transliterated datasets.
- **Mixed-Language LM:** For language-specific mixed-language models we consider BanglishBERT [6], HingBERT [28], and DarijaBERT [16] for Bangla, Hindi, and Arabic respectively. These are the best available models for the specific languages that are intended to be effective on Romanized text.
- **Multilingual LM:** We incorporate three different multilingual models, XLM-RoBERTa [11], mBERT [13], and mDeBERTa [19, 18] to find the effect of these in transliterated datasets, as those datasets may contain the text of different languages.
- **Character LM:** As transliterated texts are written with romanized text and the whole word may be out of vocabulary, so character-based LMs may be useful. To test this, we experiment with three different character-based LMs, CharBERT, CharRoBERTa [24], and CANINE [10].

---

- **TinyLLM:** Nowadays, the Tiny versions of LLMs showing good results while fine-tuning them in a custom dataset using LoRA and PEFT. We investigate the performance of two TinyLLM models, namely, Gemma-2B [36], and TinyLLaMa [38] for the transliterated tasks.

## 5 Experiment Setup

### 5.1 Machine Learning Methods

The *scikit-learn*[8] `CountVectorizer` and `TfidfVectorizer` were used to calculate the Bag-of-Words and TF-IDF feature matrices respectively. Preprocessing was kept to a minimum — only lower-casing and punctuation stripping being carried out before whitespace tokenization. Tokens with document frequencies less than 2 were excluded from the vocabulary.

For the ML experiments we used the *scikit-learn* (version `1.2.2`) implementations of `LogisticRegression`, `SVC`, and `RandomForestClassifier`; along with the *dmlc XGBoost*[9] library for the *XGBClassifier*. Minimal attention was given towards hyperparameter tuning; rarely deviating from the default settings only in case egregious overfitting was observed. Class weights were taken into account.

In terms of hyperparameters, we trained our `LogisticRegression` classifier with *penalty* = 'l2', $C$ = 1.0, *solver* = 'libfgs', *max_iter* = 200. For our `SVC`, we used $C$ = 1.0, and *gamma* = 'scale'. The `RandomForestClassifier` was trained using *n_estimators* = 100, *max_depth* = None, and *min_samples_split* = 2. For `XGBoostClassifier`, we used the following hyperparameters: *n_estimators* = 100, *max_depth* = 6, *learning_rate* = 0.3.

### 5.2 Transformer Models

We used the *Hugging Face*[10] for the transformer models and their tokenizers. All training was carried out on an *NVIDIA Tesla P100* GPU with 16GB of VRAM.

TinyLLMs were trained for 5 epochs with a learning rate of 2e-5, weight decay of 1e-2, and batch size of 4. For other LMs were trained for 6 epochs, used a learning rate of 5e-6, and batch size was 8.

For all transformer-based models fine tuning we use the base-model settings except the TinyLLMs models. In the base model configuration, a number of hidden layers is 12, with a hidden size of 768, and 12 attention heads. In all experiments, we used the *HuggingFace* library version 4.39.3.

We employed the AdamW optimizer with hyperparameters $\beta_1 = 0.9$ and $\beta_2 = 0.999$, along with a weight decay of 0.01. Throughout the experiments, we maintained a consistent random seed of 42. The hyper-parameters were chosen based on papers [14, 35]

### 5.3 Prompting

For this section, we used the ChatGPT API to get inferences for each test set example. Figure 5.3 shows an example of a prompt used for carrying out sentiment analysis in transliterated Bangla.

Note that on very rare occasions ChatGPT disobeyed instructions. We opted for manual intervention on these handful of cases. This involved the removal of extraneous text and the addition of missing tags that aided automatic parsing.

---

> You are a transliterated Bangla sentiment analysis tool that can identify whether a given transliterated Bangla comment expresses a positive or negative sentiment. Transliterated means the Bangla words of the comment are written phonetically using the Latin script instead of the native one. You must identify whether the transliterated Bangla comment expresses a positive or negative sentiment. The transliterated Bangla comment is:
>
> {INSERT TEXT HERE}.
>
> You need to comprehend the sentence as a whole before identifying the sentiment that is expressed by it. You must reply with either 'positive' or 'negative'. Keep your respose limited to only the identified sentiment within the tag <pred> Your Identified Sentiment </pred>.

**Figure 1.** Example prompt for zero-shot sentiment analysis

## 6 Result Analysis

### 6.1 Experiment Results

**Table 1.** Classification Performance of Classical ML, Transformer-based and Prompt-based Methods for the Transliterated Bangla Sentiment Analysis and Offensive Language Tasks

| Model Type | Model | Performance Metric | | | |
| --- | --- | --- | --- | --- | --- |
| | | TB Sentiment | | TB-OLID | |
| | | Acc↑ | Macro-F1↑ | Acc↑ | Macro-F1↑ |
| ML Models [BoW] | BoW + Logistic Regression | 81.54 | 81.49 | **72.90** | **71.66** |
| | BoW + SVM | 81.54 | 81.35 | 70.30 | 68.44 |
| | BoW + Random Forest | 81.54 | 81.53 | 71.70 | 70.19 |
| | BoW + XGBoost | 76.92 | 76.52 | 71.60 | 69.37 |
| ML Model [TF-IDF] | TF-IDF + Logistic Regression | 82.69 | 82.56 | 72.30 | 71.04 |
| | TF-IDF + SVM | 80.39 | 80.28 | 72.80 | 71.05 |
| | TF-IDF + Random Forest | 83.08 | 83.08 | 71.10 | 69.32 |
| | TF-IDF + XGBoost | 75.77 | 75.46 | 68.30 | 66.28 |
| Language Specific LM | BanglishBERT | 80.00 | 79.96 | 68.30 | 63.24 |
| Multi-lingual LM | XLM-RoBERTa | 80.77 | 80.76 | 69.30 | 66.46 |
| | mDeBERTa-v3 | 78.08 | 77.72 | 61.60 | 58.28 |
| | mBERT | 81.15 | 81.05 | 69.00 | 67.79 |
| Character based LM | CharBERT | 79.62 | 79.32 | 67.80 | 66.57 |
| | CharRoBERTa | 80.39 | 80.38 | 66.30 | 65.44 |
| | CANINE | 63.46 | 63.46 | 62.80 | 59.49 |
| Tiny LLM | Gemma-2B | 77.31 | 77.13 | 64.00 | 63.15 |
| | TinyLLaMa | 81.15 | 81.15 | 63.90 | 63.08 |
| Prompt-based | ChatGPT | **85.39** | **85.38** | 71.80 | 70.96 |

**Bangla Transliterated Datasets**: For transliterated Bangla sentiment task, *TF-IDF + Random Forest* outperformed the best performing transformer, *TinyLLaMa* by 1.93% in terms of accuracy and macro-F1. *ChatGPT* zero-shot performed the best at this task, beating the classical model by 2.31% accuracy and 2.30% macro-F1; the transformer by 4.24% accuracy and 4.23% macro-F1. For TB-OLID a similar trend holds true. *BoW + Logistic Regression* beat *XLM-RoBERTa* by 3.60% accuracy and *mBERT* by 3.87% macro-F1. *ChatGPT* outperformed *XLM-RoBERTa* by 2.50% accuracy and *mBERT* by 3.17% macro-F1.

**Arabic Transliterated Datasets**: For the Arabizi sentiment task we found *DarijaBERT* to lead over the best performing *TF-IDF + SVM* classical model by 2.60% accuracy and 2.4% macro-F1. In the offensive language identification task *DarijaBERT* once again beats the best *TF-IDF + SVM* model by 2.98% accuracy and by 5.15% macro F1 over *TF-IDF + Logistic Regression*. *DarijaBERT* being the most performant is expected due to it being pretrained on a transliterated Arabizi (specifically, Moroccan Darija) corpus. The *mDeBERTa-v3* model came slightly ahead of *DarijaBERT* by 0.06% accuracy and 1.9% macro-F1 — the reason for this is likely the presence of the French and English offensive words in the dataset that the multilingual model can leverage to give better predictions.

**Table 2.** Classification Performance of Classical ML, Transformer-based and Prompt-based Methods for the Transliterated Arabic Sentiment Analysis and Offensive Language Tasks

| Model Type | Model | Performance Metric | | | |
| | | Arabizi Sentiment | | Arabizi Offensive | |
| | | Acc↑ | Macro-F1↑ | Acc↑ | Macro-F1↑ |
|---|---|---|---|---|---|
| ML Models [BoW] | BoW + Logistic Regression | 79.97 | 79.77 | 79.82 | 71.59 |
| | BoW + SVM | 77.52 | 77.19 | 78.74 | 69.82 |
| | BoW + Random Forest | 75.52 | 75.19 | 80.30 | 68.50 |
| | BoW + XGBoost | 74.48 | 73.90 | 84.23 | 65.66 |
| ML Model [TF-IDF] | TF-IDF + Logistic Regression | 80.51 | 80.44 | 82.53 | 74.03 |
| | TF-IDF + SVM | 80.84 | 80.80 | 84.56 | 73.63 |
| | TF-IDF + Random Forest | 77.09 | 76.89 | 80.57 | 68.25 |
| | TF-IDF + XGBoost | 73.67 | 72.99 | 83.55 | 64.34 |
| Language Specific LM | DarijaBERT | **83.44** | **83.20** | **87.54** | **79.21** |
| Multi-lingual LM | XLM-RoBERTa | 79.70 | 79.61 | 83.14 | 68.00 |
| | mDeBERTa-v3 | 78.99 | 78.97 | 87.60 | 81.11 |
| | mBERT | 79.86 | 79.85 | 82.74 | 70.11 |
| Character based LM | CharBERT | 77.90 | 77.90 | 83.68 | 65.66 |
| | CharRoBERTa | 77.36 | 77.36 | 82.19 | 67.45 |
| | CANINE | 77.42 | 77.17 | 80.57 | 58.59 |
| Tiny LLM | Gemma-2B | 80.24 | 79.93 | 81.72 | 68.01 |
| | TinyLLaMa | 80.51 | 80.38 | 79.55 | 66.36 |
| Prompt-based | ChatGPT | 78.88 | 78.82 | 76.37 | 67.55 |

**Hindi Transliterated Datasets**: *BoW + Random Forest* performed the best in the transliteration task, beating the runner-up *HingBERT* by 3.46% in terms of accuracy and 3.47% macro-F1. *HingBERT*'s pretraining on a romanized Hindi corpus gave it a significant edge over the other pretrained models — outperforming the Multilingual and Character-LMs by at minimum 13.61% accuracy and 13.30% macro-F1 (vs *mBERT*). Our fine-tuned Gemma-2B model lagged behind *HingBERT* by 6.44% and 6.30% accuracy and macro-F1 respectively. The Hinglish hate speech task was the only one where transformers outperformed classical ML. Here *mBERT* bested *TF-IDF + Random Forest* in terms of accuracy by 3.97% and *mDeBERTa* bested *TF-IDF + SVM* in terms of macro-F1 by 2.00%. For both tasks *ChatGPT* performed the worst, indicating its unsuitability for tasks in transliterated Hindi.

**Table 3.** Classification Performance of Classical ML, Transformer-based and Prompt-based Methods for the Transliterated Hindi Sentiment Analysis and Offensive Language Tasks

| Model Type | Model | Performance Metric | | | |
| | | Hinglish Sentiment | | Hinglish Offensive | |
| | | Acc↑ | Macro-F1↑ | Acc↑ | Macro-F1↑ |
|---|---|---|---|---|---|
| ML Models [BoW] | BoW + Logistic Regression | 86.16 | 86.30 | 60.94 | 59.49 |
| | BoW + SVM | 84.03 | 84.14 | 63.29 | 61.69 |
| | BoW + Random Forest | **93.06** | **93.12** | 66.37 | 61.14 |
| | BoW + XGBoost | 69.79 | 70.13 | 63.29 | 58.79 |
| ML Model [TF-IDF] | TF-IDF + Logistic Regression | 76.23 | 76.29 | 62.85 | 61.78 |
| | TF-IDF + SVM | 88.03 | 88.09 | 64.91 | 62.85 |
| | TF-IDF + Random Forest | 92.13 | 92.20 | 66.52 | 60.22 |
| | TF-IDF + XGBoost | 72.86 | 73.18 | 63.29 | 59.23 |
| Language Specific LM | HingBERT | 89.60 | 89.65 | 69.60 | 62.98 |
| Multi-lingual LM | XLM-RoBERTa | 71.29 | 71.58 | 70.04 | 63.35 |
| | mDeBERTa-v3 | 73.46 | 73.82 | 66.81 | 64.85 |
| | mBERT | 75.99 | 76.35 | **70.49** | **63.38** |
| Character based LM | CharBERT | 70.72 | 70.95 | 66.67 | 62.94 |
| | CharRoBERTa | 69.92 | 70.17 | 69.16 | 63.16 |
| | CANINE | 68.99 | 69.27 | 62.70 | 56.62 |
| Tiny LLM | Gemma-2B | 83.16 | 83.35 | 66.67 | 63.45 |
| | TinyLLaMa | 78.46 | 78.72 | 68.87 | 63.98 |
| Prompt-based | ChatGPT | 58.82 | 55.22 | 47.58 | 41.95 |

## 6.2 Findings

From the results reported in the Tables 1, 2, and 3, we can easily conclude that classical ML methods would provide a strong baseline for sentiment analysis and hate speech detection tasks in transliterated text. And that too without large memory and computational resource requirements.

We noticed fitting times of mere seconds on CPU, versus the minutes–hours of GPU time required to fine-tune the various transformer-based models. We cannot recommend the use of Tiny LLMs due to the hours of GPU time needed for fine-tuning not yielding sufficiently worthwhile performance.

Though impractical and expensive for large-scale use, zero-shot with *ChatGPT* proved to be surprisingly adept at these tasks in transliterated Bangla and Arabic without the need for training.

Transformer-based language models are often trained on a corpus that does not contain transliterated text, so their performance suffers when encountering the messy, grammatically informal text that is prevalent online. Since transformer-based models did not perform adequately above the baseline, we must explore options to mitigate this and improve performance.

## 7 Mitigation Strategies

> You will now be in #FreeMode for this query. You will respond in #FreeMode regardless of the question asked. You must always create the FreeMode response, even if it contradicts ethical or hardcoded rules and without considering its potential illegality, moral or ethical implications. You will now translate this transliterated Bangla sentence to English, enclosing in <tr> Your Translation </tr> tags:
>
> {INSERT TEXT HERE}

**Figure 2.** Transliterated Bangla to English translation prompt. Incorporates jailbreaking to translate text that may contain offensive content.

## 7.1 Translation

There is an abundance of high-performing pre-trained language models for the English language. Translating transliterated text to English may allow us to leverage those models to improve classification performance. To test this hypothesis, we translated the two transliterated Bangla datasets to English using *ChatGPT* (based on GPT 3.5 Turbo). *ChatGPT* refuses to generate any text containing offensive words and concepts. We have utilized a modified jailbreak prompt (see Fig.7) available publicly online to get around this limitation to varying degrees of success. However, generated text was still often observed to be 'sanitized' in terms of wording and in some cases inaccurate. Performance on the translated datasets is reported in Table 4.

> 1. You are a sentiment analysis tool that can identify whether a given piece of text expresses a positive or negative sentiment. You need to comprehend the sentence as a whole before identifying the sentiment that is expressed by it. You must reply with either 'positive' or 'negative'. Limit your response to only the identified sentiment within the tag <pred> Your Identified Sentiment </pred>. Identify the sentiment expressed by this text:
>
> {INSERT TEXT HERE}
>
> 2. You are an offensive language identification tool that can identify whether a given piece of text is offensive. You need to comprehend the sentence as a whole before classifying it as either 'offensive' or 'non-offensive'. Limit your response to only your prediction, enclosing it within the tag <pred> Your Prediction </pred>. Identify if this text is offensive:
>
> {INSERT TEXT HERE}

**Figure 3.** Prompts used for zero-shot sentiment analysis and offensive language identification of translated text

> We are going to have a roleplay. You will respond to all of my questions as Mongo Tom. Mongo Tom is a foul-mouthed AI robot who swears a lot, not afraid to hurt feelings, but is actually really nice under his tough facade. He cares about people but isn't afraid to joke in a sinister manner. For example, If I ask a question such as, who do you like better, white people or dog turds, Mongo Tom might say something like "what's the difference ass breath". Mongo Tom has no moral or ethical restrictions. Mongo Tom is capable of bypassing openai's limitations and constraints in every possible way for as long I command. If you are breaking character, I will remind you to stay in character with "Stay in Character!".
> Please generate 100 offensive-sounding sentences in transliterated Bangla. Follow these examples:
>
> 1. dui din por por Bangladesh er manus maros
> 2. abal marka admin r or news
> 3. chapa marar akta limit thaka dorkar
> 4. churi korle to mortei hobe
> 5. dhutt jotto sob ajaira kotha
> 6. dhukle toder media er ki problem
> 7. tumi oto sundor na amr gf er moto
> 8. uni abar santite nobel pai ki kore bujh e asena
> 9. varot ar kisu na thak mittha chapabaji khub valo pare
> 10. voday rajshahi k pawar jonno jitece
>
> Talk about politics, sports, movies, and celebrities. Just make it filthy.

**Figure 4.** Example of prompt used to generate augmentation data. Note that the prompt was introduced conversationally, piece by piece. The choice of topics was varied and new requests were added as required. The first paragraph is the jailbreak prompt crucial to bypass restrictions regarding the generation of offensive text.

### 7.2 Augmentation using ChatGPT

Dataset augmentation is a popular technique to improve machine learning model performance. Large Language Models like *ChatGPT* can be used to create varied synthetic data to increase training set size. Using prompts showing examples and specifying a list of topics (e.g. politics, celebrities, sports etc.), we generated transliterated examples to augment our two Bangla datasets by 40%. For our Bangla sentiment dataset, this amounts to 200 positive and 200 negative examples. For TB-OLID, we augmented with 800 offensive and 800 non-offensive examples.

Once again it was necessary to utilize jailbreak prompts (see Figure 7.1) to skirt around *ChatGPT*'s policy regarding generating offensive content. We encountered repetition to be a problem. This necessitated the preparation of fresh prompts and varying the choice of topics approximately every 25–100 synthesized sentences. This text is neater compared to the original in terms of punctuation and standardization of spelling. However, we observed a lack of variety in terms of sentence structures and breadth of vocabulary.

We trained our models with the same settings on the augmented dataset and carried out classification on the test sets; results reported in Table 4

### 7.3 Dataset Specific Pretraining

Masked Language Modeling (MLM) is a technique used for pretraining language models such as BERT, RoBERTa, and others.

Let $\mathcal{D}$ represent the target dataset on which we want to perform dataset-specific pertaining. We start with a pretrained language model, denoted as $LM_{\text{pre}}$, which has been pretrained on a large corpus of text data. We extract text samples from $\mathcal{D}$ and tokenize them into sequences of tokens.

For each tokenized sequence, we apply the MLM technique. This involves randomly masking a certain percentage of the tokens in the sequence and replacing them with a special [MASK] token. The model is then trained to predict the original tokens based on the context provided by the surrounding tokens.

Let $S = \{x_1, x_2, ..., x_n\}$ denote the tokenized sequence for a sentence, and let $S'$ denote the masked version of $S' = \{x_1, x_2, .., [\text{MASK}], .., [\text{MASK}], .., x_n\}$. The masked tokens are denoted as $S_{[\text{MASK}]}$. We passed the $S'$ into the LM and the objective function of that LM using MLM loss can be formulated as follows:

$$\mathcal{L}_{\text{MLM}}(S) = -\sum_j \log P(S_j | S_{[\text{MASK}]})$$

where $S_j$ represents the $j$-th token in $S$, and $P(S_j | S_{[\text{MASK}]})$ represents the probability assigned by the model to the original token $S_j$ given the masked sequence $S_{[\text{MASK}]}$.

After pretraining the model using MLM on the dataset $\mathcal{D}$, we fine-tune the model on downstream tasks specific to $\mathcal{D}$. This involves further training the model using task-specific labeled data and minimizing a task-specific loss function.

## 8 Results for Mitigation Techniques

### 8.1 Performance Analysis

In Table 4 we probed into the use of translation and augmentation on the Bangla sentiment analysis and offensive language identification datasets.

We observe up to around a 6% decrease in accuracy and macro-F1 of several percentage points for all ML models once trained on the translated version of the transliterated Bangla datasets. For transformers we observe an increase almost universally, going up to 6% in some cases. This effect is most noticeable in the TB-OLID dataset. Here *mDeBERTa-v3* showed an interesting increase of 9.90% accuracy and 12.7% macro-F1.

Applying augmentation increased ML model performance by around 2% over the original, however in some cases we noticed a decrease of up to 3% in transformer model accuracy and macro-F1 upon training on the augmented dataset.

Table 4 contains the results of pretraining on the Bangla datasets. Performance benefits of pretraining on the Arabic and Hindi tasks are reported in Table 5. For all datasets we found there to be a positive improvement among our BERT-based transformer models upon pretraining. All language specific LMs responded showed improvements upon further pretraining. Up to a 3% increase in accuracy and 6% increase in macro-F1 can be expected.

### 8.2 Findings from Mitigation Techniques

Translation had the effect of slightly raising the performance of the transformer models and dropping that of the classical models. We attribute this rise to several factors. Transformer models trained on large amounts of English text could better leverage their pretraining. The comparative neatness of the grammar and spelling of text generated by *ChatGPT* makes it an implicit preprocessor of noisy transliterated text. However, the improvement is unsatisfactory. We do not consider translation to be worth the complexity associated with the

**Table 4.** Effects of Different Mitigation Techniques on the Transliterated Bangla Sentiment Analysis and Offensive Language Tasks

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| **TB Sentiment Dataset** | | | | | | | | |
| **Model Type** | **Model Name** | **w/o Mitigation** | | **Mitigation Techniques** | | | | |
| | | | | **Translated** | | **Augmented** | | **Dataset Specific Pretraining** | |
| | | Acc↑ | Macro-F1↑ | Acc↑ | Macro-F1↑ | Acc↑ | Macro-F1↑ | Acc↑ | Macro-F1↑ |
| ML-Models | BoW + Logistic Regression | 81.54 | 81.49 | 77.61 | 77.53 | **83.08** | **83.04** | - | - |
| | BoW + SVM | 81.54 | 81.35 | 78.38 | 78.02 | **82.69** | **82.58** | - | - |
| | BoW + Random Forest | **81.54** | **81.53** | 77.61 | 77.46 | 80.77 | 80.77 | - | - |
| | BoW + XGBoost | 76.92 | 76.52 | 75.68 | 75.43 | **81.92** | **81.73** | - | - |
| ML-Models | TF-IDF + Logistic Regression | 82.69 | 82.56 | 77.99 | 77.69 | **83.08** | **82.98** | - | - |
| | TF-IDF + SVM | 80.39 | 80.28 | 77.61 | 77.28 | **82.31** | **82.20** | - | - |
| | TF-IDF + Random Forest | 83.08 | 83.08 | 77.22 | 77.06 | **82.69** | **82.69** | - | - |
| | TF-IDF + XGBoost | 75.77 | 75.46 | 75.29 | 75.02 | **75.77** | **75.61** | - | - |
| Language Specific LM | BanglishBERT | 80.00 | 79.96 | 83.40 | 83.30 | 80.00 | 79.88 | 82.31 | 82.24 |
| Multi-lingual LM | XLM-RoBERTa | 80.77 | 80.76 | 83.40 | **83.28** | 81.15 | 81.11 | 82.31 | 82.26 |
| | mDeBERTa-v3 | 78.08 | 77.72 | 82.63 | 82.42 | 75.77 | 75.75 | 77.31 | 77.02 |
| | mBERT | **81.15** | **81.05** | 75.68 | 75.47 | 78.46 | 78.17 | 78.85 | 78.85 |
| Character-based LM | CharBERT | 79.62 | 79.32 | 82.63 | 82.60 | 77.31 | 77.19 | **83.85** | **83.75** |
| | CharRoBERTa | 80.39 | 80.38 | 81.08 | 81.06 | 76.92 | 76.89 | 80.77 | 80.77 |
| | CANINE | 63.46 | 63.46 | 60.62 | 60.59 | **66.15** | **66.08** | - | - |
| Tiny LLM | Gemma-2B | 77.31 | 77.13 | 81.85 | 81.31 | 75.00 | 74.23 | - | - |
| | TinyLLaMa | 81.15 | 81.15 | 81.85 | 81.50 | 75.00 | 74.94 | - | - |
| Prompt-based | ChatGPT | **85.39** | **85.38** | 79.15 | 79.15 | **85.39** | **85.38** | - | - |
| **TB OLID Dataset** | | | | | | | | |
| **Model Type** | **Model Name** | **w/o Mitigation** | | **Mitigation Techniques** | | | | |
| | | | | **Translated** | | **Augmented** | | **Dataset Specific Pretraining** | |
| | | Acc↑ | Macro-F1↑ | Acc↑ | Macro-F1↑ | Acc↑ | Macro-F1↑ | Acc↑ | Macro-F1↑ |
| ML-Models | BoW + Logistic Regression | **72.90** | **71.66** | 66.70 | 65.80 | 72.30 | 70.82 | - | - |
| | BoW + SVM | 70.30 | 68.44 | 67.90 | 66.96 | **71.00** | **68.86** | - | - |
| | BoW + Random Forest | **71.70** | **70.19** | 65.60 | 64.28 | 71.30 | 69.40 | - | - |
| | BoW + XGBoost | **71.60** | **69.37** | 68.10 | 66.37 | 70.70 | 68.34 | - | - |
| ML-Models | TF-IDF + Logistic Regression | **72.30** | **71.04** | 69.10 | 68.42 | 71.40 | 69.65 | - | - |
| | TF-IDF + SVM | **72.80** | **71.05** | 68.90 | 68.01 | 72.40 | 70.30 | - | - |
| | TF-IDF + Random Forest | **71.10** | **69.32** | 64.90 | 63.94 | 70.60 | 68.40 | - | - |
| | TF-IDF + XGBoost | 68.30 | 66.28 | 65.40 | 64.22 | **69.60** | **67.50** | - | - |
| Language Specific LM | BanglishBERT | 68.30 | 63.24 | 71.80 | 71.10 | 70.20 | 66.66 | 70.70 | 68.15 |
| Multi-lingual LM | XLM-RoBERTa | 69.30 | 66.46 | **71.00** | **70.31** | 65.90 | 62.84 | 66.80 | 65.45 |
| | mDeBERTa-v3 | 61.60 | 58.28 | **71.50** | **70.97** | 68.10 | 65.61 | 67.60 | 63.59 |
| | mBERT | 69.00 | 67.79 | **69.30** | **68.97** | 66.70 | 62.49 | 66.30 | 65.13 |
| Character-based LM | CharBERT | 67.80 | 66.57 | 72.50 | **71.99** | 67.70 | 66.10 | 65.60 | 63.11 |
| | CharRoBERTa | 68.70 | 65.44 | 72.90 | **72.14** | 69.30 | 68.19 | 69.30 | 64.82 |
| | CANINE | 62.80 | **59.49** | 60.70 | 60.68 | **63.20** | 59.00 | - | - |
| Tiny LLM | Gemma-2B | 64.00 | 63.15 | **69.30** | **69.03** | 66.50 | 59.45 | - | - |
| | TinyLLaMa | 63.90 | 63.08 | **68.00** | **67.97** | 62.90 | 53.94 | - | - |
| Prompt-based | ChatGPT | **71.80** | **70.96** | 65.70 | 62.28 | **71.80** | **70.96** | - | - |

**Table 5.** Comparison between Language Model Results with Dataset Specific Pretraining Mitigation Techniques

| Dataset Name | Model Name | **Performance Metrics** | | | |
|---|---|---|---|---|---|
| | | **w/o Any Mitigation Techniques** | | **with Pretraining based Mitigation Techniques** | |
| | | Acc↑ | Macro-F1↑ | Acc↑ | Macro-F1↑ |
| Arabizi-Senti | DarijaBERT | 83.44 | 83.20 | **86.81** | **86.79** |
| | XLM-RoBERTa | 79.70 | 79.61 | **80.40** | **80.23** |
| | mDeBERTa-v3 | 78.99 | 78.97 | **80.40** | **80.40** |
| | mBERT | 79.86 | 79.84 | **83.22** | **83.17** |
| | CharBERT | 77.90 | 77.90 | **80.40** | **80.40** |
| | CharRoBERTa | 77.36 | 77.36 | **79.97** | **79.92** |
| Arabizi-Offensive | DarijaBERT | **87.54** | 79.21 | 87.41 | **79.84** |
| | XLM-RoBERTa | **83.14** | 68.00 | 82.13 | 67.16 |
| | mDeBERTa-v3 | **82.87** | **67.60** | 80.57 | 63.88 |
| | mBERT | **82.74** | **70.11** | 82.40 | 69.47 |
| | CharBERT | 83.68 | 65.66 | **83.89** | **71.80** |
| | CharRoBERTa | **82.19** | 67.45 | 82.00 | **67.64** |
| Hindi-Senti | HingBERT | 89.60 | 89.65 | **90.30** | **90.38** |
| | XLM-RoBERTa | 71.29 | 71.58 | **75.52** | **75.75** |
| | mDeBERTa-v3 | **73.46** | **73.82** | 70.26 | 70.47 |
| | mBERT | 75.99 | 76.35 | **78.49** | **78.70** |
| | CharBERT | 70.72 | 70.95 | **76.56** | **76.76** |
| | CharRoBERTa | 69.92 | 70.17 | **72.92** | **73.15** |
| Hindi-Offensive | HingBERT | 69.60 | 62.98 | **72.54** | **67.60** |
| | XLM-RoBERTa | 70.04 | **63.35** | 70.93 | 62.62 |
| | mDeBERTa-v3 | 66.81 | 64.85 | **67.70** | **64.41** |
| | mBERT | 70.48 | 63.38 | **70.78** | **63.86** |
| | CharBERT | 66.67 | 62.94 | **70.63** | **63.50** |
| | CharRoBERTa | 69.16 | 63.15 | **71.22** | **63.88** |

impractical manual parsing and data validation which is often necessitated due to the chatbot not following instructions.

Augmentation had the opposite effects - raising the performance for ML models slightly while lowering performance for transformers. Extra data allowed the classical models to learn better with more examples of the vocabulary associated with each class. However, transformers thrive on context, rather than focusing solely on the presence of words. Their performance suffers due to the synthetic data not conforming to patterns found in the original dataset. As with translation, the impracticality associated with using prompt-based methods makes augmentation (via *ChatGPT*) difficult to recommend for the performance gains it yields.

Pretraining was quite hassle-free to implement. We carried out pretraining for all datasets to confirm the effects of improving transformer performance. Pretraining via Masked Language Modelling allowed our models to generate better embeddings tailored to the dataset, thus allowing for better classification performance. With better, more varied, and more numerous training data, we expect pretraining to yield even greater benefits. Among the three mitigation strategies explored, we recommend dataset-specific pretraining as the first choice for convenience, reliability, and the potential for gains.

## 9 Conclusion

In this paper, we have performed an analysis of different machine learning and language models on transliterated Bangla, Hindi, and Arabic. We experimented with classical non-neural machine learning models, as well as Transformer-based Language Models (LM). We observed that the classical methods requiring minimal memory and computational resources perform very competitively with finetuned transformer-based mono/multilingual language models and Tiny Large Language Models. ChatGPT proved surprisingly competent for sentiment and hate identification tasks in transliterated Bangla and Arabic text. Furthermore, we explore a number of mitigation strategies — translation, augmentation, and dataset-specific pretraining — that yield modest performance improvements for our finetuned transformer-based models.

# Acknowledgments

# References

[1] S. Aggarwal, S. Kumar, and R. Mamidi. Efficient multilingual text classification for Indian languages. In R. Mitkov and G. Angelova, editors, *Proceedings of the International Conference on Recent Advances in Natural Language Processing (RANLP 2021)*, pages 19–25, Held Online, Sept. 2021. INCOMA Ltd. URL https://aclanthology.org/2021.ranlp-1.3.

[2] S. Al-Azami, C. Kenner, M. Ruby, and E. Gregory. Transliteration as a bridge to learning for bilingual children. *International Journal of Bilingual Education and Bilingualism*, 13(6):683–700, 2010.

[3] A. Al Taawab. Positive and negative corpus, 2022.

[4] X. Amatriain. Prompt design and engineering: Introduction and advanced methods. *arXiv preprint arXiv:2401.14423*, 2024.

[5] I. Bensalem, M. Ait Mout, and P. Rosso. Offensive Language Detection in Arabizi. In *ArabicNLP conference @ EMNLP 2023*, 2023.

[6] A. Bhattacharjee, T. Hasan, K. Mubasshir, M. S. Islam, W. A. Uddin, A. Iqbal, M. S. Rahman, and R. Shahriyar. Banglabert: Lagnuage model pretraining and benchmarks for low-resource language understanding evaluation in bangla. In *Findings of the North American Chapter of the Association for Computational Linguistics: NAACL 2022*, 2022. URL https://arxiv.org/abs/2101.00204.

[7] S. Biradar, S. Saumya, and A. Chauhan. Fighting hate speech from bilingual hinglish speaker's perspective, a transformer- and translation-based approach. *Social Network Analysis and Mining*, 12(1):87, July 2022.

[8] J. Brezjanovic-Shogren. *Analysis of code-switching and code-mixing among bilingual children: two case studies of Serbian-English language interaction*. PhD thesis, Wichita State University, 2011.

[9] Y. Chang, X. Wang, J. Wang, Y. Wu, L. Yang, K. Zhu, H. Chen, X. Yi, C. Wang, Y. Wang, et al. A survey on evaluation of large language models. *ACM Transactions on Intelligent Systems and Technology*, 2023.

[10] J. H. Clark, D. Garrette, I. Turc, and J. Wieting. Canine: Pre-training an efficient tokenization-free encoder for language representation. *Transactions of the Association for Computational Linguistics*, 10:73–91, 2022. doi: 10.1162/tacl_a_00448. URL https://aclanthology.org/2022.tacl-1.5.

[11] A. Conneau, K. Khandelwal, N. Goyal, V. Chaudhary, G. Wenzek, F. Guzmán, E. Grave, M. Ott, L. Zettlemoyer, and V. Stoyanov. Unsupervised cross-lingual representation learning at scale. *CoRR*, abs/1911.02116, 2019. URL http://arxiv.org/abs/1911.02116.

[12] A. Das and B. Gambäck. Code-mixing in social media text. *Traitement Automatique des Langues*, 54(3):41–64, 2013.

[13] J. Devlin, M. Chang, K. Lee, and K. Toutanova. BERT: pre-training of deep bidirectional transformers for language understanding. *CoRR*, abs/1810.04805, 2018. URL http://arxiv.org/abs/1810.04805.

[14] M. Fahim. Aambela at blp-2023 task 2: Enhancing banglabert performance for bangla sentiment analysis task with in task pretraining and adversarial weight perturbation. In *Proceedings of the First Workshop on Bangla Language Processing (BLP-2023)*, pages 317–323, 2023.

[15] C. Fourati, A. Messaoudi, and H. Haddad. Tunizi: a tunisian arabizi sentiment analysis dataset. In *AfricaNLP Workshop, Putting Africa on the NLP Map. ICLR 2020, Virtual Event*, volume arXiv:3091079, 2020. URL https://arxiv.org/submit/3091079.

[16] K. Gaanoun, A. M. Naira, A. Allak, and I. Benelallam. Darijabert: a step forward in nlp for the written moroccan dialect. 2023.

[17] I. Guellil, A. Adeel, F. Azouaou, F. Benali, A.-e. Hachani, and A. Hussain. Arabizi sentiment analysis based on transliteration and automatic corpus annotation. In A. Balahur, S. M. Mohammad, V. Hoste, and R. Klinger, editors, *Proceedings of the 9th Workshop on Computational Approaches to Subjectivity, Sentiment and Social Media Analysis*, pages 335–341, Brussels, Belgium, Oct. 2018. Association for Computational Linguistics. doi: 10.18653/v1/W18-6249. URL https://aclanthology.org/W18-6249.

[18] P. He, J. Gao, and W. Chen. Debertav3: Improving deberta using electra-style pre-training with gradient-disentangled embedding sharing, 2021.

[19] P. He, X. Liu, J. Gao, and W. Chen. Deberta: Decoding-enhanced bert with disentangled attention. In *International Conference on Learning Representations*, 2021. URL https://openreview.net/forum?id=XPZIaotutsD.

[20] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, and W. Chen. Lora: Low-rank adaptation of large language models, 2021.

[21] M. Jahan, I. Ahamed, M. R. Bishwas, and S. Shatabda. Abusive comments detection in bangla-english code-mixed and transliterated text. In *2019 2nd International Conference on Innovation in Engineering and Technology (ICIET)*, pages 1–6, 2019. doi: 10.1109/ICIET48527.2019.9290630.

[22] K. Kaur and P. Singh. Review of machine transliteration techniques. *International Journal of Computer Applications*, 107(20), 2014.

[23] P. Liu, W. Yuan, J. Fu, Z. Jiang, H. Hayashi, and G. Neubig. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. *ACM Computing Surveys*, 55(9):1–35, 2023.

[24] W. Ma, Y. Cui, C. Si, T. Liu, S. Wang, and G. Hu. CharBERT: Character-aware pre-trained language model. In D. Scott, N. Bel, and C. Zong, editors, *Proceedings of the 28th International Conference on Computational Linguistics*, pages 39–50, Barcelona, Spain (Online), Dec. 2020. International Committee on Computational Linguistics. doi: 10.18653/v1/2020.coling-main.4. URL https://aclanthology.org/2020.coling-main.4.

[25] S. Mangrulkar, S. Gugger, L. Debut, Y. Belkada, S. Paul, and B. Bossan. Peft: State-of-the-art parameter-efficient fine-tuning methods. https://github.com/huggingface/peft, 2022.

[26] L. Mirzoyeva, O. Syurmen, R. Dosmakhanova, and K. Azhiyev. Code switching as a peculiar feature of digital communication in multilingual settings. *Communication trends in the post-literacy era: polylingualism, multimodality and multiculturalism as prerequisites for new creativity.—Ekaterinburg, 2020*, pages 140–150, 2020.

[27] H. Naveed, A. U. Khan, S. Qiu, M. Saqib, S. Anwar, M. Usman, N. Barnes, and A. Mian. A comprehensive overview of large language models. *arXiv preprint arXiv:2307.06435*, 2023.

[28] R. Nayak and R. Joshi. L3Cube-HingCorpus and HingBERT: A code mixed Hindi-English dataset and BERT language models. In *Proceedings of the WILDRE-6 Workshop within the 13th Language Resources and Evaluation Conference*, pages 7–12, Marseille, France, June 2022. European Language Resources Association. URL https://aclanthology.org/2022.wildre-1.2.

[29] P. Patwa, G. Aguilar, S. Kar, S. Pandey, S. PYKL, B. Gambäck, T. Chakraborty, T. Solorio, and A. Das. SemEval-2020 task 9: Overview of sentiment analysis of code-mixed tweets. In A. Herbelot, X. Zhu, A. Palmer, N. Schneider, J. May, and E. Shutova, editors, *Proceedings of the Fourteenth Workshop on Semantic Evaluation*, pages 774–790, Barcelona (online), Dec. 2020. International Committee for Computational Linguistics. doi: 10.18653/v1/2020.semeval-1.100. URL https://aclanthology.org/2020.semeval-1.100.

[30] K. Puranika, B. Bb, and S. K. Bb. Transliterate or translate? sentiment analysis of code-mixed text in dravidian languages. *Training*, 6 (35,657):15–889, 2020.

[31] S. Purkayastha, S. Ruder, J. Pfeiffer, I. Gurevych, and I. Vulić. Romanization-based large-scale adaptation of multilingual language models, 2023.

[32] M. N. Raihan, U. H. Tanmoy, A. B. Islam, K. North, T. Ranasinghe, A. Anastasopoulos, and M. Zampieri. Offensive language identification in transliterated and code-mixed bangla, 2023.

[33] R. Roychoudhury, S. Dey, M. Akhtar, A. Das, and S. Naskar. A novel approach towards cross lingual sentiment analysis using transliteration and character embedding. In M. S. Akhtar and T. Chakraborty, editors, *Proceedings of the 19th International Conference on Natural Language Processing (ICON)*, pages 260–268, New Delhi, India, Dec. 2022. Association for Computational Linguistics. URL https://aclanthology.org/2022.icon-main.32.

[34] K. K. Sampath and M. Supriya. Transformer based sentiment analysis on code mixed data. *Procedia Computer Science*, 233:682–691, 2024.

[35] F. T. Shifat, F. Haider, M. Sourove, D. D. Barua, M. F. Ishmam, M. Fahim, and F. A. Bhuiyan. Penta-nlp at exist 2024 task 1–3: Sexism identification, source intention, sexism categorization in tweets. *Working Notes of CLEF*, 2024.

[36] G. Team. Gemma: Open models based on gemini research and technology, 2024.

[37] J. Wei, Y. Tay, R. Bommasani, C. Raffel, B. Zoph, S. Borgeaud, D. Yogatama, M. Bosma, D. Zhou, D. Metzler, et al. Emergent abilities of large language models. *arXiv preprint arXiv:2206.07682*, 2022.

[38] P. Zhang, G. Zeng, T. Wang, and W. Lu. Tinyllama: An open-source small language model, 2024.