

---

# OFFENSIVE LANGUAGE DETECTION WITH BERT-BASED MODELS, BY CUSTOMIZING ATTENTION PROBABILITIES

---

**Peyman Alavi**  
NLP Research Laboratory  
Shahid Beheshti University  
Tehran, Iran  
seyedp.alavi@mail.sbu.ac.ir

**Pouria Nikvand**  
NLP Research Laboratory  
Shahid Beheshti University  
Tehran, Iran  
po.nikvand@mail.sbu.ac.ir

**Mehrnoush Shamsfard**  
NLP Research Laboratory  
Shahid Beheshti University  
Tehran, Iran  
m-shams@sbu.ac.ir

## ABSTRACT

This paper describes a novel study on using ‘Attention Mask’ input in transformers and using this approach for detecting offensive content in both English and Persian languages. The paper’s principal focus is to suggest a methodology to enhance the performance of the BERT-based models on the ‘Offensive Language Detection’ task. Therefore, we customize attention probabilities by changing the ‘Attention Mask’ input to create more efficacious word embeddings. To do this, we firstly tokenize the training set of the exploited datasets (by BERT tokenizer). Then, we apply Multinomial Naive Bayes to map these tokens to two probabilities. These probabilities indicate the likelihood of making a text non-offensive or offensive, provided that it contains that token. Afterwards, we use these probabilities to define a new term, namely Offensive Score. Next, we create two separate (because of the differences in the types of the employed datasets) equations based on Offensive Scores for each language to re-distribute the ‘Attention Mask’ input for paying more attention to more offensive phrases. Eventually, we put the F1-macro score as our evaluation metric and fine-tune several combinations of BERT with ANNs, CNNs and RNNs to examine the effect of using this methodology on various combinations. The results indicate that all models will enhance with this methodology. The most improvement was 2% and 10% for English and Persian languages, respectively.

**Keywords** Offensive Language Detection · BERT-based Models · Attention Mask · Attention Probability · Multinomial Naive Bayes

## 1 Introduction

These days, social media provide the most convenient way of communication for individuals and make the users capable of sending messages containing attitudes, criticisms and daily conversations immediately. Unfortunately, increasing social media popularity has led to pervading more offensive content to users, which has become a critical problem for the communities (Ranasinghe and Hettiarachchi, 2020). Insulting content can jeopardise the community’s mental health and can affect user experience detrimentally. Accordingly, it is indispensable to recognise and throttle the offensive content before it appears to individuals’ observance.

The amount of textual information generated daily in social media is tremendous; consequently, it becomes inconceivable and infeasible for humans to identify and remove offensive messages manually. Furthermore, base word filtering could not be an adequate solution because it cannot consider the influence of aspects, such as the domain of an utterance, discourse context, the author’s and targeted recipient’s identity and other extra information (Schmidt and Wiegand, 2017). Accordingly, scholars, online communities, social media platforms, and IT companies have decided to elaborate more intelligent systems based on NLP techniques for identifying offensive language.

Along with these efforts, several shared tasks were launched regarding covering a particular aspect of offensive language. TRAC (Trolling, Aggression and Cyberbullying) (Kumar et al., 2018a) (Kumar et al., 2020), HatEval (Basile et al., 2019), HASOC (Hate Speech and Offensive Content Identification) (Majumder et al., 2019), and OffensEval (Zampieri et al., 2019a) (Zampieri et al., 2020) are some significant instances of these types of competitions. The first version

of OffensEval covered just the English language, while its second version comprised more languages such as Arabic, Danish, Greek and Turkish.

Since the NLP techniques need the numerical representation of textual information, word embedding is widely used. Word embedding is the dense vector representation of words in lower-dimensional space. Several word embedding models have been proposed, which Vaswani et al. (2017) presented a novel neural network architecture, called a transformer, which had many benefits over the conventional sequential models (LSTM, GRU). Bidirectional Encoder Representations from Transformer (BERT) (Devlin et al., 2018) and ParsBERT (Farahani et al., 2020) are such transformers that are used widely for English and Persian word embedding.

This paper intends to present that the performance of the BERT-based models on the ‘Offensive Language Detection’ task will enhance, provided that the models pay more attention to more offensive phrases. This procedure can be achieved by changing the ‘Attention Mask’ input, which affects the attention probabilities directly. To demonstrate that this improvement occurs in different languages, we tried this methodology on two distinct languages (English and Persian). In English, we used the OLID (Zampieri et al., 2019b) dataset, and in Persian, we created our custom dataset, namely POLID. These two datasets are a little different in context; hence, the approach of changing the ‘Attention Mask’ input will be changed. For each language, we tried several BERT-based models. All these models experienced an improvement in their performance when they paid more attention to more offensive phrases.

The remainder of this paper is structured as follows: section 2 describes the related research and datasets in the field of offensive language detection. Section 3 explicates the core of our methodology. Section 4 presents an analysis of our evaluation results on the OLID (Zampieri et al., 2019b) and POLID datasets. Eventually, Section 5 offers the conclusion. The code of this paper is available here: <https://github.com/peyman-aly/offensive-detection>

## 2 Background

### 2.1 Related Work

Fundamentally, offensive language consists of several varieties, such as aggressive behaviour, bullying and hate speech. Nockleby (2000) defines hate speech as “any animosity or disparagement of an individual or a group on account of a group characteristics such as race, colour, national origin, sex, disability, religion, or sexual orientation”. Nitta et al. (2013) defines a particular part of bullying, called cyberbullying, as “humiliating and slandering behaviour towards other people through Internet services, such as BBS, Twitter or e-mails.”

Various text classification approaches have exploited traditional or feature-based supervised learning techniques in early research for covering a particular part of offensive languages (Ranasinghe and Hettiarachchi, 2020). Chen et al. (2012) applied text-mining methods to create the Lexical Syntactic Feature (LSF) to identify offensive language in social media and predict a user’s potentiality to send out offensive contents. Malmasi and Zampieri (2017) applied linear Support Vector Machine (SVM) on several surface n-grams and word skip-grams for detecting hate speech in social media. Razavi et al. (2010) exploited Complement Naïve Bayes, Multinomial Updatable Naïve Bayes and Decision Table/Naïve Bayes hybrid classifier for building a multi-level classifier for flame detection by boosting an underlying dictionary of abusive and insulting phrases. One of the main privileges of their model is being able to be modified based on any accumulative training data. Montani (2018) defined five disjoint sets of features (namely Character N-Grams, Important N-Grams, Token N-Grams, Important Tokens and Word Embeddings) and trained Logistic Regression with balanced class weight, two sets of an ensemble of Random Forests on each feature group and ensemble the results due to detecting abusive language in German. They understood that ‘Important Tokens’ is the most useful feature set because their evaluation metric dropped when this set was removed. Their proposed approach stood at first in GermEval 2018 (Wiegand et al., 2018). Indurthi et al. (2019) fed the output of Universal Sentence Encoder (Cer et al., 2018) to the SVM with RBF kernel for detecting hate speech in English HatEval. Their approach surpassed Neural Network models (such as FastText + BiGRU, BERT + CNN and GloVe + LSTM) because they stood at first rank.

By advancing and conducting various research in NLP, it is proven that Deep Neural Networks (DNNs) are capable of outperforming in compared to these traditional approaches; however, there are needed sufficient training instances for reaching optimum weights. Badjatiya et al. (2017) applied several Deep Learning approaches for detecting hate speech and concluded that exploiting “LSTM + Random Embedding + Gradient Boosted Decision Tree” can outperform other methods. Gambäck and Sikdar (2017) applied several possible embedding approaches for feeding to a Convolutional Neural Network (CNN), which reached their best experimental result by exploiting Word2Vec. Park and Fung (2017) proposed a model, named HybridCNN, for classifying hate speech, which takes both character and word features as input. Aroyehun and Gelbukh (2018) used FastText to represent inputs and use LSTM for classifying aggression in texts. They stood at the first rank in TRAC (Trolling, Aggression and Cyberbullying) shared task.

The emergence of the transformers and applying them for extracting features causes proposing state-of-the-art models in NLP tasks. Thanks to OffensEval shared tasks, that has been caused proposing many transformer-based approaches by participants. NULI (Liu et al., 2019), UM-IU@LING (Zhu et al., 2019) and Embeddia (Pelicon et al., 2019) are those examples of participants in OffensEval 2019 who fine-tuned BERT with variations in the parameters and preprocessing steps for subtask-A, which aimed to detect offensive language. UHH-LT (Wiedemann et al., 2020) fine-tuned different transformer models on the OLID training data, and then combined them into an ensemble. Galileo (Wang et al., 2020) and Rouges (Dadu and Pant, 2020) are those participants who fine-tuned XLM-R (XLM-RoBERTa) to detect the offensive language in tweets across five languages (English, Turkish, Greek, Arabic, and Danish). XLM-RoBERTa is an unsupervised cross-lingual representation pretrained transformer model, which is highly scalable and can be fine-tuned for different downstream tasks in multiple languages (Conneau et al., 2019). GUIR (Sotudeh et al., 2020) trained linear SVM classifiers (view-classifiers) using 6-gram features in addition to the BERT-based classifier and fed the concatenation of the probability output from SVM and sigmoid output from BERT as feature vector for a final linear SVM classifier (the meta-classifier). Kungfupanda (Dai et al., 2020) stood at the sixth rank in the English OffensEval 2020. They fed the text to the BERT model and passed its output to three similar levels of BiLSTM layers, followed by linear layers with different units. KUISAIL (Safaya et al., 2020) showed that combining CNN with BERT is better than using BERT on its own. This participant stood at top-five participants for Arabic, Greek and Turkish languages in OffensEval 2020 by using language-specific variations of BERT (separately for each language) and passed their output to a CNN, which is designed based on mixture of experts (Baldacchino et al., 2016) concept.

## 2.2 Related Datasets

The OLID 2019 is a hierarchical dataset for detecting the offensive language (Subtask-A), categorizing it (Subtask-B), and identifying its target (Subtask-C). Subtask-A is a binary classification problem that contains ‘NOT’ and ‘OFF’ labels. ‘NOT’ refers to those posts, which do not contain any offense, whilst ‘OFF’ determines posts including insults, threats, and posts containing profane language or swear words. The SOLID 2020 (Rosenthal et al., 2020), refers to Semi-supervised OLID, which were used several NLP techniques to prepare more volume dataset.

HatEval provided a hierarchical dataset, similar to OLID 2019. This dataset aims to identify hate speech against women and immigrants. The dataset is distributed among two languages, namely English and Spanish. The first level of English language version judges whether the content of the text is hatefully provided that ‘HS’ (a binary value) becomes 1.

The TRAC-1 dataset (Kumar et al., 2018b), has English and Hindi formats, in which the instances are categorized into 3 different labels, namely ‘NAG’/Not-Aggressive, ‘CAG’/Covertly Aggressive, ‘OAG’/Overtly Aggressive. Bhattacharya et al. (2020) provides the second version of TRAC-1 and adds a second level for identifying misogyny. The second level is a binary classification, in which the gendered or misogynous contents are specified by ‘GEN’, otherwise ‘NGEN’ is used.

## 3 Methodology

In this section, we present our methodology for detecting offensive contents in Persian and English languages. For this purpose, we first talk about the datasets including our newly created dataset, named POLID (for Persian OLID) and other exploited resources . Then, we describe our proposed models and algorithms. The overall view of our methodology is depicted in Figure 1.

### 3.1 Creating Persian OLID (POLID)

To the best of our knowledge, there are not available and open datasets related to offensive language detection in Persian. Hence, we had to create our dataset. Due to collecting text data, we crawl tweets, Instagram comments, users’ reviews on different Iranian web applications such as Digikala <sup>1</sup>, Snappfood <sup>2</sup>, etc. We categorize the data in two steps, to make confident labels. At the first step, we create a basic list of common swearwords and label each text entity (text entity means each tweet, Instagram comment and users’ review in Iranian web applications) as ‘OFF’ if it contains at least one element of this list. Otherwise, we categorize it as ‘NOT’. The meaning of the labels is as follow: ‘NOT’ determines those texts, which neither contain profanity, nor hating, nor threatening. On the other hand, ‘OFF’ texts contain profanity, or hating, or threatening. In the second step, we correct the miss-labelled instances based on the following:

---

<sup>1</sup><https://www.digikala.com/>

<sup>2</sup><https://snappfood.ir/>

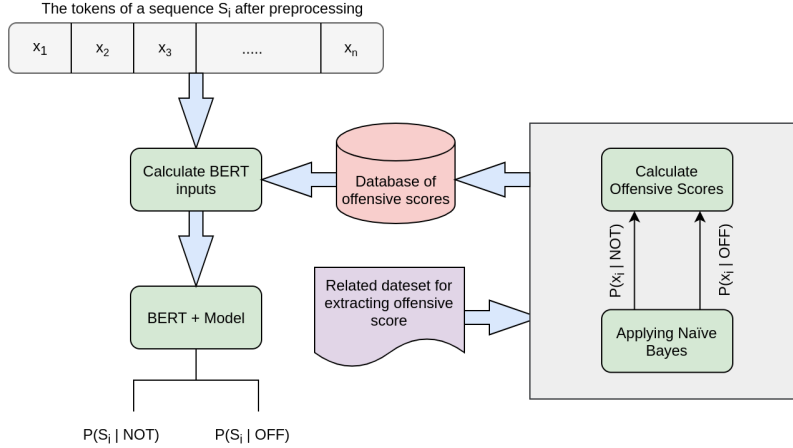


Figure 1: Schema of the proposed method.

- There are some Persian polysemous words, for which there are offensive and non-offensive senses; for instance, “لخت” means naked by pronouncing ‘lokht’ and inertia by pronouncing ‘lakht’. Hence, we may have labelled wrongly some instances as offensive.
- The text data may contains offensive content with no obscenities. These types of text may contain targeted offenses, which may be veiled or direct. For instance, in the phrase of “آخر سر میکشمت.” (Translation: I will kill you finally.), there are no swearwords, but it contains offensive language. Hence, we modified some instances, which were labeled as not offensive before.

The final dataset (POLID) contains 4,988 text entities, comprising 2,453 of not offensive and 2,535 of offensive contents.

### 3.2 Employed Datasets

There are two purposes for using the available datasets; 1) Extracting an underlying dictionary, which maps a token to its offensive score 2) Training and evaluating the models. In English, OLID, dataset of HatEval and TRAC-1 dataset were used for the first purpose and OLID was used just for the second purpose. In Persian, POLID dataset is applied for both purposes. Tables 1 and 2 illustrates the information about the distribution of used datasets for both purposes.

As discussed, each named dataset covers a particular aspect of offensive language, hence it is logical to combine them and create more comprehensive dataset. To make our models comparable to others, we merely homogenized the named datasets for extracting the offensive scores and used OLID for training and evaluating the models. The homogenized dataset contains two labels, namely ‘NOT’ and ‘OFF’. In dataset of HatEval, 0s and 1s are considered as ‘NOT’ and ‘OFF’, respectively. In TRAC-1 dataset, NAGs are considered as ‘NOT’, whilst CAGs and OAGs are considered as ‘OFF’.

### 3.3 Preprocessing

In English, we followed the preprocessing steps proposed by the Kungfupanda; however, there are some differences in Persian preprocessing steps. In continue, the modules and their differences will be discussed.

To deal with emojis, Kungfupanda replaces each emoji with its English textual descriptor. This is a logical approach for dealing with emojis because sometimes the combination of some emojis conveys an offensive concept and this information will not be missed by this module. But, in Persian, there is no comprehensive library for replacing emojis with its Persian textual descriptor. Hence, we removed each emoji in the texts inevitably.

As both OLID and POLID entities have come from social media, there are many meaningful hashtags in them. Hence, it is more effective to extract the words from hashtags. For tackling this challenge, Kungfupanda exploited a third-party library. Since there is no Persian-based library, we first detect the hashtags, then by removing underlines we extract their words.

Due to reducing redundancy, Kungfupanda replaced all ‘@USER’ with one ‘@USERS’ at the first of the text, provided that the original text contains several ‘@USER’. On the contrary, we deleted users’ IDs and numbers to reduce redundancy because we focus on detecting offensive language not identifying the people who have been insulted.

Table 1: Distribution of datasets before homogenizing.

Dataset	OLID		HatEval		TRAC-1		
Label	NOT	OFF	HS = 0	HS = 1	NAG	CAG	OAG
Counts	8,840	4,440	5,217	3,783	5,052	4,240	2,078

Table 2: The distribution of OLID 2019, which is used for training and evaluating the models.

Training			Testing		
NOT	OFF	Total	NOT	OFF	Total
8,840	4,400	13,240	620	240	860

Kungfupanda replaced ‘URL’ tokens with ‘http’ to handle Out-Of-Vocabulary (OOV) tokens. In Persian, we deleted the links because when we crawled Instagram comments, the text entities were comprising the hyperlinks to posts and they did not contain any useful information.

There are some further preprocessing steps that we applied on POLID. It includes normalizing elongated words and converting to proper Persian format.

Users may repeat some letters in a specific word for representing their emotions, such as ‘ هوا خیلی خوووووووبه!! ’ (Translation: The weather is very goooooood.). As the repetition does not affect our task, we normalized those tokens which have more than two consecutive repeated characters and lessened them to two repeated characters for reducing redundancy. The reason behind that is there are some words like ‘ گراییدن ’ (means to tend), which has two consecutive repeated characters and normalizing this word ruin its meaning. Finally, we normalize texts for character encodings by exploiting Hazm<sup>3</sup>.

Eventually, similar to the last module of Kungfupanda’s preprocessing, we truncated all the tweets to the max length of 64 to lower the GPU memory usage and slightly improve the performance.

### 3.4 Calculating Offensive Scores

In this section, we present the approach of mapping tokens to a number between 0 and 1 as their offensive scores. The closer to 1 represents the more offensive a token is. As manual mapping is prone to error and it is possible to miss some offensive tokens, we assumed that the homogenized dataset (In English) and POLID (In Persian) are multinomially distributed and applied Multinomial Naïve Bayes for automatizing the procedure.

Forasmuch as the offensive scores should be calculated as what the BERT uses, we tokenized the text entities with the BERT tokenizer that Huggingface (Wolf et al., 2019) has implemented. Then, we applied Term Frequency Inverse Document Frequency (TF-IDF) and Multinomial Naïve Bayes to extract probabilities belonging to ‘NOT’ and ‘OFF’ classes for each token. These probabilities are calculated by equation 1, where  $\hat{\theta}_{yi}$  is the probability  $P(x_i|y)$  of token  $i$  appearing in a sample belonging to class  $y$  (in this case  $y$  is ‘NOT’ or ‘OFF’),  $N_{yi}$  is the number of times that token  $i$  appears in a sample of class  $y$  in the related dataset,  $N_y$  is the total count of all features for class  $y$  and  $|V|$  is the size of the vocabulary.

$$\hat{\theta}_{yi} = \frac{N_{yi} + 1}{N_y + |V|} \quad (1)$$

After calculating probabilities  $P(x_i|NOT)$  and  $P(x_i|OFF)$  for each token like  $x_i$ , we applied equations 2 and 3 to extract offensive score for each token. Table 3 illustrates some tokens that mapped to their offensive score. For English we used text entities from homogenized dataset, while for Persian whole entities of POLID were used.

$$S(x) = \frac{1}{1 + e^{-x}} \quad (2)$$

$$OffensiveScore(x_i) = S(\log \frac{P(x_i|OFF)}{P(x_i|NOT)}) \quad (3)$$

These offensive scores will be saved in a database for future usage. Due to optimizing the performance of BERT or ParsBERT, the offensive scores will be applied in customizing the attention probabilities.

<sup>3</sup><https://github.com/sobhe/hazm>

Table 3: Some instances of token’s offensive score

Token	Offensive Score
idiots	0.8843
sucks	0.8352
filthy	0.8316
pig	0.8044
cows	0.8043
stubborn	0.6292
devil	0.6225
disabilities	0.4464
airport	0.2489

### 3.5 Calculating BERT Inputs

We used BERT and ParsBERT transformers, provided by Huggingface that get ‘Input IDs’ and ‘Attention Mask’ as required and optional parameters, respectively. ‘Input IDs’ is a vector containing indices of input sequence tokens in the vocabulary. On the other hand, ‘Attention Mask’ indicates which tokens should be attended to, and which should not. Classic ‘Attention Mask’ contains 0s and 1s, which the value of padding tokens are 0 and the others are 1.

Equation 4 represents the approach of using the ‘Attention Mask’ vector in the process of attention probabilities calculation. This equation is based on the implementation of BERT provided by Google Research<sup>4</sup>. Based on this equation, the complement of ‘Attention Mask’ is multiplied by a very small number (like -10,000) and adds to pre-calculated attention probabilities. Then, new attention probabilities will be normalized by applying softmax. This causes to have less attention to padding tokens.

$$AttentionProbabilities = Softmax(AttentionScores - 10000 * (1 - attention\_mask)) \quad (4)$$

Although the value of attention probabilities would be different for different tokens by exploiting the classic ‘Attention Mask’, it is possible to be created such values that cause the models to not pay more attention to offensive phrases. Accordingly, we believe that the process of creating word embeddings (by BERT or ParsBERT) would be improved provided that attention probabilities were customized, which it can be done by creating a custom-tailored ‘attention mask’ that focuses on offensive phrases more than others. For creating this new ‘Attention Mask’, we recommend two possible equations such as equations 5, 6 for calculating ‘Attention Mask’ value for a single token  $x_i$ .

All proposed equations, keep the masking value of padding tokens as what it is was before. Their difference is in the ‘attention mask’ value for the primary tokens. Equation 5 aims to add the offensive score of each token, provided that it finds the token in its database of the offensive score. It keeps the original ‘Attention Mask’ value for those tokens that are not in the database. Actually, the model’s attention to offensive phrases increases. On the other hand, equation 6 follows a different purpose. It aims to make the model just paying attention to very offensive tokens.

$$AttentionMaskValue_1(x_i) = \begin{cases} 0 & \text{if } x_i = \text{padding token} \\ 1 + OffensiveScore(x_i) & \text{if } x_i \text{ in database} \\ 1 & \text{otherwise} \end{cases} \quad (5)$$

$$AttentionMaskValue_2(x_i) = \begin{cases} 0 & \text{if } x_i = \text{padding token} \\ 1 & \text{if } OffensiveScore(x_i) \geq \text{threshold} \\ 0 & \text{otherwise} \end{cases} \quad (6)$$

We applied both of these equations with various values of hyper-parameter (threshold in equation 6) in both languages. Due to discovering the optimum value of the threshold, we applied the greedy search concept and set the value of the threshold between 0.5 and 0.8. If we set the value of threshold less than 0.5, we involve the tokens that  $P(x_i|NOT) \geq P(x_i|OFF)$  which means involving less offensive phrases that are not effective. Our best experimental result demonstrates that the best value of the threshold is 0.6.

<sup>4</sup><https://github.com/google-research/bert/blob/eedf5716ce1268e56f0a50264a88cafad334ac61/modeling.py#L705>

Using the equation 6 is not a good choice in English models because we use various types of datasets to cover all aspects of offensive language when calculating offensive scores, which leads to having a small number of tokens with an offensive score greater than 0.6; consequently, paying attention to only a small number of tokens does not improve the performance of the models, but also confuses the models and reduces the performance.

On the contrary, using the equation 5 in Persian models does not change the performance of them. This is because the offensive text entities of the POLID contains profanity and aggressive behaviour more than covertly offensive language; consequently, the number of tokens with a score greater than 0.6 becomes greater (This number for English and Persian is respectively 3,352 and 5,588). Furthermore, in Persian, the tokens in this range (greater than 0.6) are more offensive, whilst this is not happened in English because there are some normal tokens with high offensive score; in result, using equation 5 makes the model not feel much difference between the offensive and normal tokens.

### 3.6 Designing BERT-based Model

Since the inputs of BERT (and ParsBERT) were prepared, we combined BERT with ANNs, CNNs and RNNs to analyze the effect of the customized attention probability in the performance of models. We used CNNs either because the result of using CNNs at NLP tasks were promising (Kim, 2014)(Hughes et al., 2017).

In English, we selected the re-implemented proposed models by top-performing participants of OffensEval 2019 and 2020, named of NULI, KUISAIL and Kungfupanda as the representatives BERT-based models of ANNs, CNNs and RNNs. These proposed architectures is introduced in Related Works section. We re-implement the architecture of Kungfupanda’s model and use just their first level because our focus is in subtask-A of OLID.

In Persian, we define our models since there is not previous works. In these models, BERT is substituted with ParsBERT, whose output will be passed to some prevalent RNN and CNN models. The details of these models are as below:

**ParsBERT + BiLSTM or BiGRU,** A Dropout layer after ParsBERT with rate of 0.5, followed by a Bidirectional LSTM or GRU layer with 32 units and a Dense layer for creating probabilities and classification.

**ParsBERT + CNN\_BiLSTM,** A one-dimensional convolutional layer with 256 filters with size of 5, followed by a Max-Pooling layer with size of 2. The output is passed to a Bidirectional LSTM layer with 32 units. Batch Normalization layers is applied after Max-Pooling and BiLSTM layers for stabilizing the learning process and dramatically reducing the number of training epochs. Eventually two level of Dense layers with 16 and 2 units is used for classifying.

**ParsBERT + CNN\_2D,** The idea is similar to the proposed model by KUISAIL. The difference is reducing the number of convolutional layers and adding a Batch Normalization layer after each of them. This is because the number of instances in the POLID dataset is small and training a model with high trainable parameters may cause overfitting. The convolutional layers have 32 kernels of size 3x3, 64 kernels of size 4x4, 256 kernels of size 5x5. Finally, a dense layer with 2 units is added for classification. The new CNN\_2D possess nearly 1,000,000 learning parameters less than the KUISAIL’s model.

## 4 Experiments

In this section, we aim to represent our best evaluation results that are obtained during the testing phase for both English and Persian languages. Macro-average F1 score is used for evaluating the performance of defined models, which were trained on training sets of OLID and POLID datasets. Furthermore, the effect of using offensive scores and the rate of improvement of the best models will be covered.

### 4.1 English

The experiments began with putting the re-implementation of the proposed model by NULI as our baseline model and fine-tune it with and without exploiting offensive scores. Then, further BERT-based models, such as the models proposed by Kungfupanda and KUISAIL were used due to having more confident conclusions and we repeated the training and testing processes. The models have been re-implemented in the TensorFlow<sup>5</sup> framework and were trained on a single GeForce RTX 2080Ti.

Due to optimizing the weights of defined models during training phase, Adam Optimizer with learning rate of 1e-5 was used. The smaller values caused slow optimization and larger values caused the metric on the validation set not to be

---

<sup>5</sup><https://www.tensorflow.org/>

changed. The maximum number of epochs and batch size value were set 11 and 64, respectively. Table 4 represents the evaluation results on the testing set of OLID dataset, without and with using equation 5.

Table 4: Evaluation Results on re-implemented format of defined models.

Model	Not Offensive			Offensive			F1 Macro
	P	R	F1	P	R	F1	
NULI	0.8922	0.9210	0.9063	0.7773	0.7125	0.7435	0.8249
NULI + equation 5	0.8915	0.9274	0.9091	0.7907	0.7083	0.7473	<b>0.8282</b>
Kungfupanda	0.8930	0.9016	0.8973	0.7393	0.7208	0.7300	0.8136
Kungfupanda + equation 5	0.9042	0.9129	0.9085	0.7692	0.7500	0.7595	<b>0.8340</b>
KUISAIL	0.8717	0.9532	0.9106	0.8407	0.6375	0.7251	0.8179
KUISAIL + equation 5	0.8877	0.9306	0.9087	0.7952	0.6958	0.7422	<b>0.8254</b>

Our best model is the model proposed by Kungfupanda by applying equation 5, which is achieved 0.8340 on the testing set of OLID. Based on reports, NULI achieved 0.829 of F1-Macro on the testing set of OLID, which our best model outperforms this model. Although Kungfupanda reports that their MTL\_Transformer\_LSTM model achieved the same score as our best model, the main difference between their model and ours is that they used the information of subtask-B and subtask-C during training, while we just used the information of subtask-A, which indicates that the proposed method has the positive impact on the performance.

#### 4.1.1 Error Analysis of the Best Model (Kungfupanda + equation 5)

Figures 2a and 2b display the confusion matrices of the model proposed by Kungfupanda (with or without exploiting equation 5) against the testing set of OLID dataset. It can be seen that the enhanced model reduced both false positives and false negatives.

The normal model classifies wrongly the text “#BeckyLynch is beautiful one of the few women in wrestling that absolutely need no work done at all. She’s absolutely beautiful just the way she is. Doesn’t need giant boobs or a fake booty. @USER is just simply amazing #HIAC” as an offensive language because of some phrases, such as ‘giant boobs’ or ‘fake booty’. While, after exploiting equation 5 the amount of attention on different tokens will be changed that causes to predict this input as not offensive correctly. It turns out that the classifier, by re-distributing the attention mask, avoids over-focusing on some words that seem to be offensive at the first glance.

The enhanced model predicts the text “#Trump, everything trump does is another piece of evidence for obstruction of justice. He is a talking crime machine.” as offensive. Adding the offensive scores of the tokens to the attention mask causes the model to predict this text as an offensive language with a probability of 0.92. While the normal model cannot understand the significance of ‘crime machine’ and predicts it not offensive.

On the other hand, a few instances (18 text entities) cause the enhanced model to predict wrongly, while the normal model predicts these instances correctly. For instance, the text “@USER This is one of my favourite quotes on the internet. He is so confident that he bends maths to his will. That’s hardcore.” is classified as offensive by enhanced model, while actually, it is not offensive. It is because of ‘hardcore’ token, whose offensive score of 0.756 that causes the enhanced model to predict wrongly. As another example, “#AafiaSiddiqui is a convicted #terrorist but she is widely hailed as the Daughter of #Pakistan. URL” is also classified as offensive by the enhanced model. which is wrong. This is because of phrases like ‘convicted terrorist’, ‘daughter’ and ‘#Pakistan’ that have high offensive scores and causes the model to predict this input as offensive wrongly.

## 4.2 Persian

Our Persian experiment began with preparing and labelling the dataset, named POLID. As the number of POLID instances is small and splitting it into train and test sets may cause missing some important features, we used Stratified K-fold, where  $k = 5$  and calculate F1-Macro for each iteration and report the average of experiments in each fold.

We created a simple lexical-based classifier based on our manual swearwords list (an initial list, which was used for creating data) as our baseline model. This classifier categorizes each text entity as offensive provided that it contains at least one element of this list. The F1-Macro score of this classifier in 5 iterations of cross-validation is 0.6763, 0.6662, 0.6560, 0.6489 and 0.6880 (average: 0.6671). The high F1-Macro score for this baseline model indicates that the list covers most of the offensive instances.

Tables 5 and 6 report our best experimental results of models on different folds of POLID before and after applying equation 6. As discussed, we applied different values as the threshold. The results showed us that 0.6 is the best value. Our best performance of the models archived by using the Adam Optimizer with learning rate of to  $3e-4$  for BiLSTM,



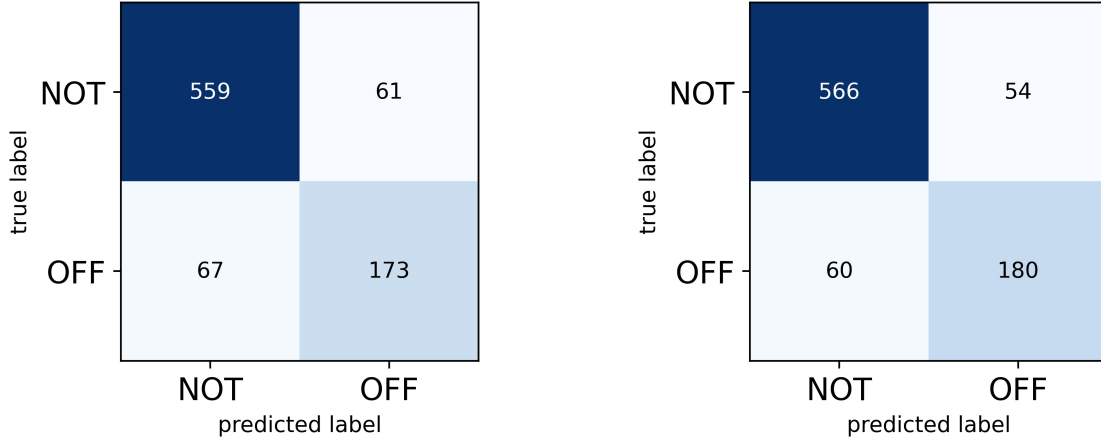


Figure 2: Confusion Matrix of the best experiments for English language, before and after using offensive scores. (a) the left is the confusion matrix of Kungfupanda without using equations. (b) the right is the confusion matrix of Kungfupanda by exploiting equation 5.

2e-5 for the model proposed by KUISAIL and 1e-3 for the others. Furthermore, The maximum number of epochs and batch size were set 30 and 64, respectively.

Table 5: Evaluation Results on defined models before applying offensive scores.

	BiLSTM		BiGRU		CNN2D		CNN-BiLSTM		KUISAIL	
	Loss	F1	Loss	F1	Loss	F1	Loss	F1	Loss	F1
<b>Fold1</b>	0.5436	0.7515	0.5456	0.7474	0.6490	0.7605	2.4684	0.7485	0.5099	0.7662
<b>Fold2</b>	0.5524	0.7493	0.5232	0.7665	0.6335	0.7709	2.0774	0.7652	0.4810	0.7685
<b>Fold3</b>	0.5285	0.7444	0.5482	0.7509	0.6646	0.7605	2.2875	0.7484	0.4915	0.7675
<b>Fold4</b>	0.5015	0.7670	0.4952	0.7773	0.6130	0.7733	2.4643	0.7560	0.4775	0.7752
<b>Fold5</b>	0.5468	0.7408	0.5301	0.7472	0.6404	0.7561	2.5533	0.7382	0.4990	0.7440
<b>Average</b>	0.5346	0.7506	0.5285	0.7579	0.6401	0.7643	2.3702	0.7513	0.4918	0.7643

Table 6: Evaluation Results on defined models after applying offensive scores

	BiLSTM		BiGRU		CNN2D		CNN-BiLSTM		KUISAIL	
	Loss	F1	Loss	F1	Loss	F1	Loss	F1	Loss	F1
<b>Fold1</b>	0.4359	0.8186	0.3649	0.8547	0.4774	0.8567	2.1694	0.8467	0.3794	0.8405
<b>Fold2</b>	0.4223	0.8276	0.3708	0.8475	0.4629	0.8557	1.7839	0.8477	0.3590	0.8437
<b>Fold3</b>	0.3963	0.8295	0.3681	0.8435	0.4521	0.8676	2.0818	0.8427	0.3513	0.8517
<b>Fold4</b>	0.3802	0.8533	0.3238	0.8785	0.4206	0.8836	2.3822	0.8551	0.3354	0.8696
<b>Fold5</b>	0.4602	0.8204	0.3960	0.8385	0.4954	0.8515	2.3775	0.8265	0.3933	0.8325
<b>Average</b>	0.4190	<b>0.8299</b>	0.3647	<b>0.8525</b>	0.4617	<b>0.8630</b>	2.1590	<b>0.8437</b>	0.3636	<b>0.8476</b>

#### 4.2.1 Error Analysis of the Best Model (CNN2D + equation 6)

There are many instances that normal CNN\_2D fails at predicting the instances of the validation set in folds, while the enhanced one predicts correctly; for instance, the input of “اگر میخوايد که یک مارمولک و موزمار باشید، سلنا باشید” (Translation: If you want to be a person like lizards and muzmaars (a Persian idiom which means cunning and insidious), be like Selena.)” is classified as a not offensive language by normal CNN\_2D. By taking a deeper look at this input, it can be understood that the ‘موزمار’ token is the principal cause of becoming offensive. This token has an offensive score of 0.7, which cause the enhanced model to pay attention to this token and classifies this input as offensive. As another example, the input “جنسی کاملاً معمولی و اندازش کوچکتر از عکس” (Translation: a normal production and smaller size than the photo)” is not offensive actually but the normal CNN\_2D predicts it as an offensive language. The enhanced model figures out that this input does not have any offensive language; however, it criticizes a little bit.

On the other hand, there are some instances that the enhanced model fails at predicting correctly; for instance, the enhanced model fails at predicting “چرت و پرت نگو ضعیفه. دلیلشو که نگفت، فقط نتیجه رو گفت. تو هم داری مثلاً دلیلشو” (Translation: don’t say bullshit weak lady, he didn’t say the reason, he just said the results. نتیجه بازم فرق نکرد.)

You are talking about its reason but the results have not been changed.)” correctly because of two reasons: 1) At this time, the model did not understand that ‘چرت و پرت’ has an inappropriate meaning. 2) Multinomial Naive Bayes could not extract appropriate probabilities for ‘ضعیفه’ which is a hateful word against women because its offensive score is near 0.53.

## 5 Conclusion and Future works

As presented, paying more attention to more offensive phrases than others can improve the performance of the models in identifying insulting language because these phrases have more impact on the final target. To do so, we created a new ‘Attention Mask’ input by using a term named offensive score to help the BERT-based models more effectively detect offensive content. Although applying this term is promising, mapping a proper value to each token is challenging. Hence, we used a statistical approach like Multinomial Naive Bayes to find offensive scores automatically.

In the future, increasing the instance of POLID will be considered. Therefore, using a semi-supervised or unsupervised technique to find offensive scores and classify textual entities can be done because labelling all items will become time-consuming. Also, one thing that can enhance the performance is improving preprocessing modules, such as converting emojis to their Persian textual descriptors.

In addition, applying such techniques which consider language structures may improve the efficiency of this methodology. Moreover, exploiting other approaches for finding probabilities (for offensive scores extraction), such as Deep Learning methods, may improve the performance of the proposed methodology. On the other hand, exploiting this approach on multi-label classification problems may be promising.

## References

- Segun Taofeek Aroyehun and Alexander Gelbukh. Aggression detection in social media: Using deep neural networks, data augmentation, and pseudo labeling. In *Proceedings of the First Workshop on Trolling, Aggression and Cyberbullying (TRAC-2018)*, pages 90–97. Association for Computational Linguistics, 2018.
- Pinkesh Badjatiya, Shashank Gupta, Manish Gupta, and Vasudeva Varma. Deep learning for hate speech detection in tweets. *CoRR*, 2017.
- Tara Baldacchino, Elizabeth J. Cross, Keith Worden, and Jennifer Rowson. Variational bayesian mixture of experts models and sensitivity analysis for nonlinear dynamical systems. *Mechanical Systems and Signal Processing*, pages 178–200, 2016.
- Valerio Basile, Cristina Bosco, Elisabetta Fersini, Debora Nozza, Viviana Patti, Francisco Manuel Rangel Pardo, Paolo Rosso, and Manuela Sanguinetti. SemEval-2019 task 5: Multilingual detection of hate speech against immigrants and women in Twitter. In *Proceedings of the 13th International Workshop on Semantic Evaluation*, pages 54–63. Association for Computational Linguistics, 2019.
- Shiladitya Bhattacharya, Siddharth Singh, Ritesh Kumar, Akanksha Bansal, Akash Bhagat, Yogesh Dawer, Bornini Lahiri, and Atul Kr. Ojha. Developing a multilingual annotated corpus of misogyny and aggression. In *Proceedings of the Second Workshop on Trolling, Aggression and Cyberbullying*, pages 158–168. European Language Resources Association (ELRA), 2020.
- Daniel Cer, Yinfei Yang, Sheng yi Kong, Nan Hua, Nicole Limtiaco, Rhomni St. John, Noah Constant, Mario Guajardo-Cespedes, Steve Yuan, Chris Tar, Yun-Hsuan Sung, Brian Strope, and Ray Kurzweil. Universal sentence encoder, 2018.
- Ying Chen, Yilu Zhou, Sencun Zhu, and Heng Xu. Detecting offensive language in social media to protect adolescent online safety. In *2012 International Conference on Privacy, Security, Risk and Trust and 2012 International Conference on Social Computing*, pages 71–80, 2012.
- Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. Unsupervised cross-lingual representation learning at scale. *CoRR*, 2019.
- Tanvi Dadu and Kartikey Pant. Team rouges at SemEval-2020 task 12: Cross-lingual inductive transfer to detect offensive language. In *Proceedings of the Fourteenth Workshop on Semantic Evaluation*, pages 2183–2189. International Committee for Computational Linguistics, 2020.
- Wenliang Dai, Tiezheng Yu, Zihan Liu, and Pascale Fung. Kungfupanda at semeval-2020 task 12: Bert-based multi-task learning for offensive language detection, 2020.

- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: pre-training of deep bidirectional transformers for language understanding. *CoRR*, 2018.
- Mehrdad Farahani, Mohammad Gharachorloo, Marzieh Farahani, and Mohammad Manthouri. Parsbert: Transformer-based model for persian language understanding, 2020.
- Björn Gambäck and Utpal Kumar Sikdar. Using convolutional neural networks to classify hate-speech. In *Proceedings of the First Workshop on Abusive Language Online*, pages 85–90. Association for Computational Linguistics, 2017.
- Mark Hughes, Irene Li, Spyros Kotoulas, and Toyotaro Suzumura. Medical text classification using convolutional neural networks. *Studies in Health Technology and Informatics*, 235, 2017.
- Vijayaradhi Indurthi, Bakhtiyar Syed, Manish Shrivastava, Nikhil Chakravartula, Manish Gupta, and Vasudeva Varma. FERMI at SemEval-2019 task 5: Using sentence embeddings to identify hate speech against immigrants and women in Twitter. In *Proceedings of the 13th International Workshop on Semantic Evaluation*, pages 70–74. Association for Computational Linguistics, 2019.
- Yoon Kim. Convolutional neural networks for sentence classification. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1746–1751. Association for Computational Linguistics, 2014.
- Ritesh Kumar, Atul Kr. Ojha, Shervin Malmasi, and Marcos Zampieri. Benchmarking aggression identification in social media. In *Proceedings of the First Workshop on Trolling, Aggression and Cyberbullying (TRAC-2018)*, pages 1–11. Association for Computational Linguistics, 2018a.
- Ritesh Kumar, Aishwarya N. Reganti, Akshit Bhatia, and Tushar Maheshwari. Aggression-annotated Corpus of Hindi-English Code-mixed Data. In Nicoletta Calzolari (Conference chair), Khalid Choukri, Christopher Cieri, Thierry Declerck, Sara Goggi, Koiti Hasida, Hitoshi Isahara, Bente Maegaard, Joseph Mariani, H el ene Mazo, Asuncion Moreno, Jan Odijk, Stelios Piperidis, and Takenobu Tokunaga, editors, *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*. European Language Resources Association (ELRA), 2018b.
- Ritesh Kumar, Atul Kr. Ojha, Shervin Malmasi, and Marcos Zampieri. Evaluating aggression identification in social media. In *Proceedings of the Second Workshop on Trolling, Aggression and Cyberbullying*, pages 1–5. European Language Resources Association (ELRA), 2020.
- Ping Liu, Wen Li, and Liang Zou. NULI at SemEval-2019 task 6: Transfer learning for offensive language detection using bidirectional transformers. In *Proceedings of the 13th International Workshop on Semantic Evaluation*, pages 87–91. Association for Computational Linguistics, 2019.
- Prasenjit Majumder, Daksh Patel, Sandip Modha, and Thomas Mandl. Overview of the hasoc track at fire 2019: Hate speech and offensive content identification in indo-european languages. 2019.
- Shervin Malmasi and Marcos Zampieri. Detecting hate speech in social media. In *Proceedings of the International Conference Recent Advances in Natural Language Processing, RANLP 2017*, pages 467–472. INCOMA Ltd., 2017.
- J. P. Montani. Germeval 2018 : German abusive tweet detection. 2018.
- Taisei Nitta, Fumito Masui, Michal Ptaszynski, Yasutomo Kimura, Rafal Rzepka, and Kenji Araki. Detecting cyberbullying entries on informal school websites based on category relevance maximization. In *Proceedings of the Sixth International Joint Conference on Natural Language Processing*, pages 579–586. Asian Federation of Natural Language Processing, 2013.
- John T. Nockleby. Hate speech. In Leonard W Levy, Kenneth L Karst, and Adam Winkler, editors, *Encyclopedia of the American Constitution*, page 1277–1279. Macmillan Reference USA, 2000.
- Ji Ho Park and Pascale Fung. One-step and two-step classification for abusive language detection on Twitter. In *Proceedings of the First Workshop on Abusive Language Online*, pages 41–45. Association for Computational Linguistics, 2017.
- Andra  Pelicon, Matej Martinc, and Petra Kralj Novak. Embeddia at SemEval-2019 task 6: Detecting hate with neural network and transfer learning approaches. In *Proceedings of the 13th International Workshop on Semantic Evaluation*, pages 604–610. Association for Computational Linguistics, 2019.
- Tharindu Ranasinghe and Hansi Hettiarachchi. BRUMS at semeval-2020 task 12 : Transformer based multilingual offensive language identification in social media. *CoRR*, 2020.
- Amir H. Razavi, Diana Inkpen, Sasha Uritsky, and Stan Matwin. Offensive language detection using multi-level classification. In Atefeh Farzindar and Vlado Ke elj, editors, *Advances in Artificial Intelligence*, pages 16–27. Springer Berlin Heidelberg, 2010.

- Sara Rosenthal, Pepa Atanasova, Georgi Karadzhov, Marcos Zampieri, and Preslav Nakov. A large-scale semi-supervised dataset for offensive language identification. *arXiv preprint arXiv:2004.14454*, 2020.
- Ali Safaya, Moutasem Abdullatif, and Deniz Yuret. KUISAIL at SemEval-2020 task 12: BERT-CNN for offensive speech identification in social media. In *Proceedings of the Fourteenth Workshop on Semantic Evaluation*, pages 2054–2059. International Committee for Computational Linguistics, 2020.
- Anna Schmidt and Michael Wiegand. A survey on hate speech detection using natural language processing. In *Proceedings of the Fifth International Workshop on Natural Language Processing for Social Media*, pages 1–10. Association for Computational Linguistics, 2017.
- Sajad Sotudeh, Tong Xiang, Hao-Ren Yao, Sean MacAvaney, Eugene Yang, Nazli Goharian, and Ophir Frieder. GUIR at SemEval-2020 task 12: Domain-tuned contextualized models for offensive language detection. In *Proceedings of the Fourteenth Workshop on Semantic Evaluation*, pages 1555–1561. International Committee for Computational Linguistics, 2020.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *CoRR*, 2017.
- Shuohuan Wang, Jiaxiang Liu, Xuan Ouyang, and Yu Sun. Galileo at semeval-2020 task 12: Multi-lingual learning for offensive language identification using pre-trained language models. *CoRR*, 2020.
- Gregor Wiedemann, Seid Muhie Yimam, and Chris Biemann. UHH-LT at SemEval-2020 task 12: Fine-tuning of pre-trained transformer networks for offensive language detection. In *Proceedings of the Fourteenth Workshop on Semantic Evaluation*, pages 1638–1644. International Committee for Computational Linguistics, 2020.
- Michael Wiegand, Melanie Siegel, and Josef Ruppenhofer. Overview of the germeval 2018 shared task on the identification of offensive language. 2018.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, and Jamie Brew. Huggingface’s transformers: State-of-the-art natural language processing. *CoRR*, 2019.
- Marcos Zampieri, Shervin Malmasi, Preslav Nakov, Sara Rosenthal, Noura Farra, and Ritesh Kumar. Semeval-2019 task 6: Identifying and categorizing offensive language in social media (offenseval). In *Proceedings of the 13th International Workshop on Semantic Evaluation*, pages 75–86, 2019a.
- Marcos Zampieri, Shervin Malmasi, Preslav Nakov, Sara Rosenthal, Noura Farra, and Ritesh Kumar. Predicting the type and target of offensive posts in social media. *CoRR*, 2019b.
- Marcos Zampieri, Preslav Nakov, Sara Rosenthal, Pepa Atanasova, Georgi Karadzhov, Hamdy Mubarak, Leon Derczynski, Zeses Pitenis, and Çağrı Çöltekin. SemEval-2020 Task 12: Multilingual Offensive Language Identification in Social Media (OffenseEval 2020). In *Proceedings of SemEval*, 2020.
- Jian Zhu, Zuoyu Tian, and Sandra Kübler. UM-IU@LING at SemEval-2019 task 6: Identifying offensive tweets using BERT and SVMs. In *Proceedings of the 13th International Workshop on Semantic Evaluation*, pages 788–795. Association for Computational Linguistics, 2019.