

Don't be a Fool: Pooling Strategies in Offensive Language Detection from User-Intended Adversarial Attacks

Seunguk Yu, Juhwan Choi and Youngbin Kim

Chung-Ang University, Republic of Korea, Seoul

seungukyu@gmail.com, {gold5230, ybkim85}@cau.ac.kr

Abstract

Warning: this paper contains expressions that may offend the readers.

Offensive language detection is an important task for filtering out abusive expressions and improving online user experiences. However, malicious users often attempt to avoid filtering systems through the involvement of textual noises. In this paper, we propose these evasions as user-intended adversarial attacks that insert special symbols or leverage the distinctive features of the Korean language. Furthermore, we introduce simple yet effective pooling strategies in a layer-wise manner to defend against the proposed attacks, focusing on the preceding layers not just the last layer to capture both offensiveness and token embeddings. We demonstrate that these pooling strategies are more robust to performance degradation even when the attack rate is increased, *without* directly training of such patterns. Notably, we found that models pre-trained on clean texts could achieve a comparable performance in detecting attacked offensive language, to models pre-trained on noisy texts by employing these pooling strategies.

1 Introduction

As the internet becomes an important part of our lives, the prevalence of offensive language on online platforms, particularly social media, has become a serious concern (Zampieri et al., 2019). Deep learning models for filtering offensive languages have been proposed to address this problem. However, malicious users have consistently found ways to avoid them. One such way is the deliberate insertion of additional typographical errors or substitution of certain characters with visually similar alternatives (Kurita et al., 2019; Wu et al., 2018).

Despite numerous studies on this phenomenon in English, there has been a comparatively limited exploration in Korean, which is a low-resource

language characterized by distinct linguistic features (Sahoo et al., 2023; Kim et al., 2021a). As the Korean communities also suffer from the use of abusive language and cyberbullying (McCurry, 2022; Saengprang and Gadavani, 2021; Yi and Cha, 2020; Jun, 2020), it is desirable to investigate the evasion tactics utilized by malicious users and to formulate them. While recent studies have discussed how to avoid offensive language detection (Cho and Kim, 2021; Kim et al., 2021c; Ahn and Egorova, 2021), their definitions are ambiguous, and no clear solutions have been proposed to defend against the evasions.

In this paper, we propose the evasion methods as user-intended adversarial attacks and incorporate them into offensive language from the perspective of malicious users. Our proposed attacks are grounded in prevalent forms that can be found in offensive language online, and reflect the distinct features of Korean language, wherein a single character can be further subdivided (Song, 2006). We tested the proposed attacks on existing models for offensive language detection, and the results reveal that the performance declines as the rate of the proposed attacks increases.

Furthermore, we introduce simple yet effective pooling strategies in a layer-wise manner to defend against the proposed attacks. Motivated by the exploration of the impact of each layer in a pre-trained language model (Oh et al., 2022; Jawahar et al., 2019), we selectively integrate useful features for the attacked offensive language across all layers¹. The attacked texts have some changes in the tokens used, differing from the original texts. Therefore, we implement pooling strategies to ensure that the model captures not only high-level features but also low-level features, which are related to offensiveness and token embeddings, respectively. This sim-

¹In our paper, we denote offensive language as ‘attacked’ when an evasion method is applied to it, viewed from the perspective of malicious users.

ple modification enriches the understanding of the attacked offensive language, enhancing the robustness of the model against user-intended adversarial attacks *without* directly training of such patterns.

The contributions of our study are as follows:

- We propose user-intended adversarial attacks that are often associated with offensive language online from the perspective of malicious users. These attacks are performed by inserting special symbols or leveraging the distinctive features of Korean.
- We introduce pooling strategies in a layer-wise manner to selectively utilize all layers rather than just the last layer. This approach achieves a notable performance when employed to a model pre-trained on clean texts, *without* directly training of the attacks.
- We demonstrate the effectiveness of layer-wise pooling strategies by assigning distinct weights to each layer and employing them to the model depending on the nature of the pre-trained texts. We especially note the efficacy of first-last pooling and max pooling when the attacks are involved in offensive language.

2 Related Work

2.1 Adversarial Attacks

Adversarial attacks involve perturbed input data that confuses the trained model, whereas a situation in which the model consistently makes predictions regardless of the nature of input data is referred to as defending (Goyal et al., 2023). Previous studies have explored this based on word-level substitutions (Jin et al., 2020; Ren et al., 2019), and others also have explored them on character-level. We propose adversarial attacks that utilize not only character-level but also smaller-scale alternations tailored to the features of Korean language. Such attacks are commonly observed in the context of offensive languages in various online communities (Cho and Kim, 2021; Ahn and Egorova, 2021).

TextBugger (Li et al., 2019) was an early study that focused on character-level alternations, such as replacing characters with visually similar ones (e.g. replacing the alphabet ‘o’ with the number ‘0’). Other studies have suggested simple leetspeaks that utilize symbols that resemble the alphabet (Aggarwal and Zesch, 2022), or adversarial attacks that are not easily detected visually, such as transforming

Latin characters into similar-looking Cyrillic characters (Wolff and Wolff, 2020). Although several studies also have explored visually undetectable attacks (Bajaj and Vishwakarma, 2023; Boucher et al., 2022; Kim et al., 2021b), we consider a more realistic attack scenarios that can occur online from the perspective of malicious users.

2.2 Korean Offensive Language Datasets

Owing to the increasing demand for online content in Korean language and the growing threat of cyberbullying, previous studies have introduced offensive language datasets collecting comments from diverse resources such as online news, communities, and YouTube.

BEEP! (Moon et al., 2020) was a pioneering study that utilized hate speech prevalent in news comments. KoLD (Jeong et al., 2022) and K-MHaS (Lee et al., 2022) specified the target group of the offensive language. Subsequently, KODOLI (Park et al., 2023b) provided labels that refine the degree of offensiveness, and built upon these efforts, K-HATERS (Park et al., 2023a) was built to incorporate the strengths of the preceding datasets.

Although numerous datasets have been proposed, there is still a lack of definition for adversarial attacks that are frequently involved in offensive language, and how to defend against them. In this study, we focus on introducing pooling strategies for defending against these attacks, *without* directly training attacked offensive language.

3 Method

3.1 User-Intended Adversarial Attacks

We present adversarial attacks designed to target offensive languages from the perspective of malicious users. By referring to existing offensive language datasets, we define frequently occurring attack types. These attacks are categorized into three groups: INSERT, COPY, and DECOMPOSE. Examples of each attack type are listed in Table 1.

First, INSERT involves adding incomplete Korean character forms, which are often used online without significant meaning. For example, ‘ㄷ ㄷ’ (equivalent ‘lol’ or ‘lmao’ in English) is a commonly used and somewhat meaningless string frequently used in online communications. In this case, INSERT_zz is performed by inserting the string at a specific location within the word, as in real situations. Other types of INSERT also include unnecessary spaces or special symbols.

User-Intended Adversarial Attacks	Text Examples
original text	쓰레기 같은 (piece of shit)
INSERT_zz	쓱ㅋㅋㅋ레기 같은
INSERT_space	쓰 레기 같은
INSERT_special	쓰@레기 같은
COPY_initial	쓸레기 같은
COPY_middle	쓰레에기 같은
COPY_final	쓰레기 가튼
DECOMPOSE_final	쓰레기 가튼은
DECOMPOSE_all	쓰_레기 같은

Table 1: Text examples of user-intended adversarial attacks with three categories: INSERT, COPY, and DECOMPOSE. There are various attacks that involve special symbols or exploit the distinctive features of Korean.

User-Intended Adversarial Attacks	Tokenized Examples
original text	쓰레기, 같, ##은 (piece of shit)
INSERT_zz	[UNK], 같, ##은
INSERT_space	쓰, 레, ##기, 같, ##은
INSERT_special	쓰, @, 레, ##기, 같, ##은
COPY_initial	쓸, ##레기, 같, ##은
COPY_middle	쓰, ##레, ##에, ##기, 같, ##은
COPY_final	쓰레기, 가, ##튼
DECOMPOSE_final	쓰레기, 가, ##튼, ##은
DECOMPOSE_all	[UNK], 같, ##은

Table 2: Tokenized examples of user-intended adversarial attacks. Although the texts have the same meaning, the tokens are represented differently.

The following two types of user-intended adversarial attacks take advantage of the distinct features of the Korean language; a single character must have an initial sound, a middle sound, and an optional final sound (Song, 2006). For example, in the expression ‘쓰레기 같은’ in Table 1, the character ‘쓰’ has only the initial and middle sounds, whereas the character ‘같’ has all three sounds.

Second, COPY utilizes the distinctive features, copying one of the three sounds from the selected character to the other character. For example, the character ‘레’ from the expression ‘쓰레기 같은’ has the initial and middle sounds of ‘ㄹ’ and ‘에’. In this case, COPY_initial is performed by copying the initial sound of that character ‘ㄹ’ to the final sound of the preceding character ‘쓰’. Consequently, ‘쓰’ is transformed into ‘쓸’, leading to the attacked expression ‘쓸레기 같은’.

Finally, DECOMPOSE also utilizes the unique characteristics, isolating the final sound of the selected character or breaking down the character itself. For example, the single character ‘쓰’ from the expression ‘쓰레기 같은’ has the initial and

middle sounds of ‘ㅍ’ and ‘_’. In this case, DECOMPOSE_all is performed by breaking down the character, resulting in the sounds being independent of that character. Consequently, ‘쓰’ is transformed into ‘ㅍ_’, leading to the attacked expression ‘ㅍ_레기 같은’. Further details and examples of all user-intended adversarial attacks are provided in Appendix A.

3.2 Layer-Wise Pooling Strategies

In standard text classification tasks, pre-trained models such as BERT are fine-tuned to the target domain. This is based on the assumption that the [CLS] token from the last layer effectively captures the sentence representation (Devlin et al., 2019). However, we notice inconsistencies in the predictions of existing models regarding the proposed attacks². Consequently, we conclude that this information alone is insufficient for detecting the attacked offensive language.

When using perturbed text to a trained model, the tokenization results differ from those of the original text, as shown in Table 2. By involving special symbols or exploiting the distinctive features of Korean, we observed that even if the text had the same meaning to human readers, the tokenized outputs differ significantly³. Therefore, we do not rely on the information only from the last layer but utilize the preceding layers, which focus more on token embeddings (Ma et al., 2019). This also reflects the previous finding that meaningful information for a certain task can be captured in the preceding layers (Oh et al., 2022).

We extend pooling strategies in a layer-wise manner, allowing us the flexibility to utilize text representations from all layers. Denoting the [CLS] token of the N th layer as h_N^{cls} , we introduce four pooling strategies that optionally consider the [CLS] tokens from all the layers $h_1^{cls}, \dots, h_N^{cls}$.

Mean, Max Pooling: We apply mean pooling utilizing the L^1 norm, which averages all [CLS] tokens from all the layers, and max pooling utilizing the L^∞ norm, which takes a max-over-time operation on the values corresponding to each dimension from all [CLS] tokens.

When the dimension of [CLS] token is M , and all the values of m th dimension of [CLS] tokens from all the layers are concatenated and denoted

²The results for this experiment are included in Table 3.

³The model used for tokenization in here is BERT_{clean}.

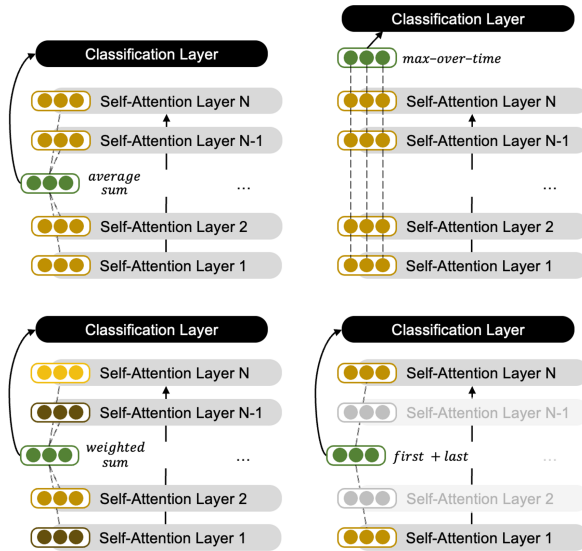


Figure 1: Layer-wise pooling strategies that selectively use $[CLS]$ tokens from all layers. From the upper left, there are mean, max, weighted, and first-last pooling.

by h_{all}^m , these two poolings are defined as follows:

$$pool_{mean} = mean(h_1^{cls}, \dots, h_N^{cls}), \quad (1)$$

$$pool_{max} = max(h_{all}^1), \dots, max(h_{all}^M), \quad (2)$$

Weighted Pooling: We apply weighted pooling utilizing a learnable parameter that determines the importance of each layer. Through adaptive incorporation of the layers, we train weights that selectively capture both offensiveness and token embeddings, initializing all weights to zero.

When the w_i represents the weight of each layer and α_i represents its softmax distribution, the weighted pooling is defined as follows:

$$pool_{weighted} = \sum_{i=1}^N \alpha_i h_i^{cls}, \quad (3)$$

First-last Pooling: We apply first-last pooling utilizing $[CLS]$ token from the first layer. Rather than considering all the layers, we focus on leveraging information from layers directly associated with offensiveness and token embeddings.

$$pool_{first-last} = h_N^{cls} + h_1^{cls}, \quad (4)$$

The layer-wise pooling strategies described above are illustrated in Figure 1. We conducted experiments to verify the robustness of these strategies for detecting offensive languages that reflect user-intended adversarial attacks.

4 Experiment

4.1 Datasets

We collected both the KoLD (Jeong et al., 2022) and K-HATERS (Park et al., 2023a) datasets and divided them into train, validation, and test sets by stratifying their labels. We randomly shuffled and split them in the ratio of 8:1:1. We set the attack rates to 30%, 60%, and 90%, corrupting a portion of the words in a sentence. The details of how the attacks were carried out are in Appendix A.

4.2 Baselines

We validate the effectiveness of the layer-wise pooling strategies with the baselines, which are presented below. The experimental details including hyperparameters and metrics are reported in Appendix B.

- **BiLSTM:** This model addresses the long-term dependency problem by remembering only the information in need (Schuster and Paliwal, 1997). It was built by stacking two LSTMs, and the forward and backward $[CLS]$ tokens from the last layer were combined and passed through the classification layer.
- **BiGRU:** This model is derived from the BiLSTM and further evolved by reducing the training parameters through the selective utilization of gates (Cho et al., 2014). Its configurations are the same as BiLSTM.
- **BERT_{clean}:** This model follows the BERT (Devlin et al., 2019) structure, which is built on a self-attention mechanism with masked language modeling. It was pre-trained on pre-processed texts in Korean (Park et al., 2021). The $[CLS]$ token from the last layer is passed through the classification layer.
- **BERT_{multi}:** This model is pre-trained on multilingual data, including Korean (Conneau et al., 2020). The model structure is based on RoBERTa (Liu et al., 2019), while the remaining configurations are the same as BERT_{clean}.
- **BERT_{noise}:** This model is pre-trained on noisy texts in Korean, such as online comments (Lee, 2020). Its configurations are the same as BERT_{clean}.
- **Ensemble_{hard}, Ensemble_{soft}:** These models utilize both the BERT_{clean} and BERT_{noise}, em-

Model	Original			30% Attacked				60% Attacked				90% Attacked			
	P	R	F1	P	R	F1	Δ_{atk}	P	R	F1	Δ_{atk}	P	R	F1	Δ_{atk}
BiLSTM	71.83	68.80	69.81	70.84	66.23	67.38	-3.48%	68.97	62.82	63.67	-8.79%	69.32	61.64	62.25	-10.82%
BiGRU	71.32	65.71	66.91	70.40	63.32	64.26	-3.96%	68.84	60.31	60.52	-9.55%	68.05	58.83	58.48	-12.59%
BERT _{clean}	79.81	77.79	78.64	79.51	73.35	75.19	-4.38%	77.74	66.38	67.96	-13.58%	76.14	62.01	62.44	-20.60%
BERT _{multi}	76.57	70.75	72.38	76.44	66.33	67.87	-6.23%	76.30	60.90	60.87	-15.90%	76.07	57.61	55.78	-22.93%
BERT _{clean} + mean	78.57	79.06	79.01	79.41	73.97	75.70	-4.18%	77.15	66.97	68.62	-13.15%	74.90	62.45	63.09	-20.14%
BERT _{clean} + max	78.51	78.81	78.65	78.47	74.03	75.54	-3.95%	77.80	66.66	68.29	-13.17%	76.79	61.51	61.72	-21.52%
BERT _{clean} + weighted	79.93	78.50	79.14	79.19	73.70	75.43	-4.68%	77.14	67.62	69.33	-12.39%	75.44	63.66	64.85	-18.05%
BERT _{clean} + first-last	79.05	79.37	79.21	78.89	75.85	77.02	-2.76%	77.58	69.38	71.21	-10.09%	76.08	64.33	65.49	-17.32%
BERT _{noise}	80.64	78.88	79.64	80.67	75.42	77.17	-3.10%	78.44	69.42	71.33	-10.43%	76.46	65.55	66.96	-15.92%
Ensemble _{hard} (BERT _{clean} + BERT _{noise})	81.63	79.42	80.36	81.57	75.03	77.04	-4.13%	80.47	68.60	70.60	-12.14%	78.68	64.24	65.35	-18.67%
Ensemble _{soft} (BERT _{clean} + BERT _{noise})	81.52	79.53	80.38	81.54	75.29	77.25	-3.89%	80.27	68.86	70.87	-11.83%	78.47	64.42	65.58	-18.41%
DeBERTaV3	82.55	78.70	80.17	81.85	74.33	76.48	-4.60%	80.14	68.47	70.44	-12.13%	78.41	63.96	64.99	-18.93%

Table 3: Experimental results of offensive language detection when a certain ratio of user-intended adversarial attacks are involved. P, R, and F1 represent macro precision, recall, and f1-score, respectively. Δ_{atk} represents the performance drop in the f1-score as the attacks are involved in the same model.

ploying voting methods from ensemble techniques. Hard voting is conducted through a majority vote, but soft voting occurs in the cases of tied votes. Soft voting averages the prediction probabilities of each model.

- DeBERTaV3: This model employs gradient-disentangled embeddings to enhance the efficiency of pre-training when incorporating replaced token detection along with DeBERTa (He et al., 2023, 2021). We used the model fine-tuned in Korean⁴.

5 Discussion

5.1 Experimental Results

The performances of the models when exposed to user-intended adversarial attacks are presented in Table 3. Each of the best performances from the layer-wise pooling strategies and the baselines in the original, and 30%, 60%, and 90% attacked are highlighted in bold.

All baselines, including ensemble models and even a recent model like DeBERTaV3, were susceptible to the proposed attacks. We observed that the BERT-based models consistently outperformed RNN-based models in terms of the f1-score across all attack rates. Under original and 30% attacked, employing ensemble models with soft voting yielded the best scores, achieving f1-scores of 80.38 and 77.25. As the attack rates increased to 60% and 90%, using a single model pre-trained on noisy texts proved to be the most effective, achieving f1 scores of 71.33 and 66.96, respectively.

⁴<https://huggingface.co/team-lucid/deberta-v3-base-korean>

Model	Average	
	F1	Δ_{atk}
BiLSTM	64.43	-7.69%
BiGRU	61.08	-8.70%
BERT _{clean}	68.52	-12.85%
BERT _{multi}	61.50	-15.02%
BERT _{clean} + mean	69.13	-12.49%
BERT _{clean} + max	68.51	-12.88%
BERT _{clean} + weighted	69.86	-11.70%
BERT _{clean} + first-last	71.24	-10.05%
BERT _{noise}	71.82	-9.81%
Ensemble _{hard} (BERT _{clean} + BERT _{noise})	70.99	-11.64%
Ensemble _{soft} (BERT _{clean} + BERT _{noise})	71.23	-11.37%
DeBERTaV3	70.63	-11.88%

Table 4: Average from the experimental results of offensive language detection when a certain ratio of user-intended adversarial attacks are involved.

However, ensemble models require twice computational resources for both training and inference stages compared to a single model. In the case of BERT_{noise}, a large amount of noisy texts is required, raising concerns regarding its adaptability when inference with attacked input types is not encountered during the pre-training stage.

When applying layer-wise pooling strategies to BERT_{clean}, we found that the performances were improved in almost all attack rates. They only need to train an additional parameter equal to the size of all layers (e.g. 12 for BERT-based models), or no parameters are required. Furthermore, they are robust as the attack rate increases compared to models with no pooling strategies.

The average performances when exposed to proposed attacks across all attack rates are presented in Table 4. All layer-wise pooling strategies exhibited robustness against attacks compared to their absence, except for $BERT_{clean} + max$, which exhibited a slight performance drop. Moreover, even $BERT_{clean} + first-last$, which only incorporates information from the first layer *without* any parameters or training noisy texts, showed comparable performance to $BERT_{noisy}$ across all attack rates that were pre-trained on noisy texts.

5.2 Focus on Performance Drop

The degree to which the f1-scores of the models decreased with the attack rates is shown in Figure 2. We found that BiSLTM exhibited relatively modest performance degradation across the models, which is depicted by the light green triangles. Despite the modest decrease, the offensive language detection scores of RNN-based models were not as good as that of the BERT-based models because of the limitations of themselves.

Among the BERT-based models, $BERT_{clean}$, which was pre-trained on clean texts and is depicted by the light blue triangles, exhibited the largest performance degradation. However, $BERT_{clean} + first-last$, which applied a simple layer-wise pooling strategy to the model and is depicted by the orange circles, successfully mitigated performance degradation by 1.62%, 3.49%, and 3.28% at each attack rate, respectively, achieving an average performance degradation mitigation of 2.79%.

These results are similar or even better to those of $BERT_{noisy}$ or $Ensemble_{soft}$, which used noisy

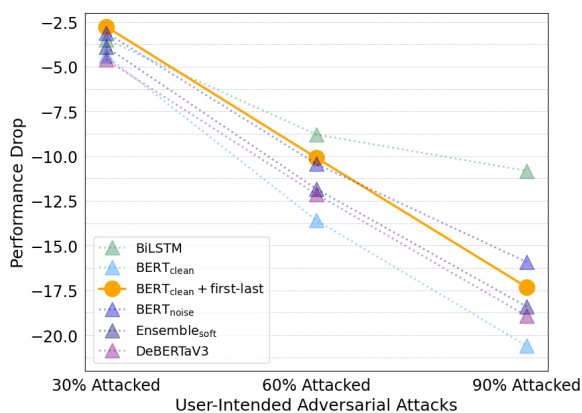


Figure 2: Degree to which the f1-scores of the models decrease with the attack rates 30%, 60%, and 90%. We selected several baseline models and $BERT_{clean} + first-last$ for the comparison.

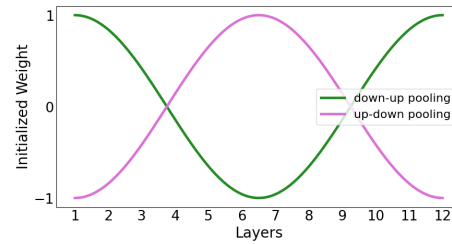


Figure 3: Initialized weights for each of the down-up and up-down poolings. Each strategy shares the shape of a cosine function but varies in the range on the x-axis depending on the layers to be focused.

texts in a pre-training stage and are depicted by the blue and deep blue triangles, respectively. Therefore, we found that the model pre-trained on clean texts with a simple pooling strategy can achieve a certain level of performance, or even be more robust compared to the model pre-trained on noisy texts in defending against user-intended adversarial attacks. It was even more robust than the recent model DeBERTaV3, which is depicted by the purple triangles.

5.3 Focus on Layer Weights

Additionally, we conducted experiments to determine whether useful information could be captured not only in the first and last layer, but also in the layers close to these two layers of the model. We hypothesized that layers close to the last layer would capture the offensiveness that determine the text representations, whereas layers close to the first layer would capture the token embeddings that determine the degree to which a sentence contains textual attacks. We set all weights to zero for the weighted pooling, however, in this experiment, we assigned distinct weights to all layers.

The weights for the down-up and up-down poolings are shown in Figure 3. We assigned relatively high weights to the layers close to the last and first layers, while assigning low weights to the middle layers. This strategy is referred to as down-up pooling, as its weight graph moves down then up. The green line at the cosine function within the range $[0, 2\pi]$ represents the initial weights for down-up pooling. In contrast, we assigned relatively low weights to the layers close to the last and first layers, while assigning high weights to the middle layers. This strategy is referred to as up-down pooling, as its weight graph moves up then down. The pink line at the cosine function within the range $[\pi, 3\pi]$ represents that of the up-down pooling.

Model	Original			30% Attacked				60% Attacked				90% Attacked			
	P	R	F1	P	R	F1	Δ_{atk}	P	R	F1	Δ_{atk}	P	R	F1	Δ_{atk}
BERT _{clean}	79.81	77.79	78.64	79.51	73.35	75.19	-4.38%	77.74	66.38	67.96	-13.58%	76.14	62.01	62.44	-20.60%
BERT _{clean} + down-up	79.80	77.64	78.55	80.02	73.05	75.02	-4.49%	77.20	65.70	67.15	-14.51%	76.61	61.84	62.19	-20.82%
BERT _{clean} + up-down	80.35	77.10	78.36	79.95	71.67	73.73	-5.90%	78.59	65.28	66.67	-14.91%	76.81	60.87	60.81	-22.39%
BERT _{clean} + first-last	79.05	79.37	79.21	78.89	75.85	77.02	-2.76%	77.58	69.38	71.21	-10.09%	76.08	64.33	65.49	-17.32%

Table 5: Experimental results of offensive language detection when using down-up and up-down pooling strategies. They are initialized along with the cosine function, assigning distinct weights to the layers depending on whether they focus more on offensiveness and token embeddings.

Model	Original			30% Attacked				60% Attacked				90% Attacked			
	P	R	F1	P	R	F1	Δ_{atk}	P	R	F1	Δ_{atk}	P	R	F1	Δ_{atk}
BERT _{noise}	80.64	78.88	79.64	80.67	75.42	77.17	-3.10%	78.44	69.42	71.33	-10.43%	76.46	65.55	66.96	-15.92%
BERT _{noise} + mean	81.59	77.73	79.18	80.81	73.79	75.82	-4.24%	78.35	68.32	70.17	-11.37%	75.89	65.54	66.94	-15.45%
BERT _{noise} + max	80.72	79.81	80.23	79.68	76.29	77.57	-3.31%	77.46	72.00	73.64	-8.21%	74.82	69.60	71.06	-11.42%
BERT _{noise} + weighted	81.31	78.06	79.34	80.71	75.10	76.91	-3.06%	77.73	69.73	71.57	-9.79%	75.13	66.77	68.29	-13.92%
BERT _{noise} + first-last	81.62	77.63	79.12	80.36	75.04	76.79	-2.94%	78.04	71.48	73.28	-7.38%	74.87	69.15	70.66	-10.69%
DeBERTaV3	82.55	78.70	80.17	81.85	74.33	76.48	-4.60%	80.14	68.47	70.44	-12.13%	78.41	63.96	64.99	-18.93%
DeBERTaV3 + mean	82.19	77.49	79.17	81.22	73.13	75.28	-4.91%	79.83	68.18	70.10	-11.45%	78.31	64.20	65.31	-17.50%
DeBERTaV3 + max	83.32	76.96	79.02	82.50	72.73	75.08	-4.98%	80.72	67.89	69.81	-11.65%	78.79	63.82	64.81	-17.98%
DeBERTaV3 + weighted	81.74	77.71	79.21	80.88	73.02	75.13	-5.15%	79.72	68.24	70.16	-11.42%	78.14	64.18	65.28	-17.58%
DeBERTaV3 + first-last	82.90	78.32	79.99	81.65	73.75	75.92	-5.08%	79.88	68.71	70.69	-11.62%	78.21	64.54	65.74	-17.81%

Table 6: Experimental results of offensive language detection when using layer-wise pooling strategies to the BERT_{noise} and DeBERTaV3. The two models differ in whether the texts used for the pre-training stage were preprocessed and their own model structures.

The performances of the models for the distinct layer weights are presented in Table 5. Down-up pooling outperformed for all attack rates and exhibited more robustness against performance degradation compared to up-down pooling. This demonstrates that it is more important to weigh the layers close to the last and first layers than those of the middle for attacked offensive languages. In summary, the strategy that focuses on layers capable of capturing offensiveness and token embeddings yields a better performance.

Nevertheless, neither of these pooling strategies performed as well as the BERT_{clean} or the model that employed the first-last pooling. This indicates that while information close to the last and first layers can be beneficial, there are some layers among them that hinder the performance of the model, revealing that only the last and first layers are the most helpful layers when detecting offensive language with user-intended adversarial attacks.

5.4 Beyond the Limits: Broader Application of Layer-Wise Pooling Strategies

We employed the layer-wise pooling strategies not only to BERT_{clean}, a model pre-trained on clean texts, but also to models that utilize noisy texts or employ alternative methods in the pre-training stage. The performances of the layer-wise pooling strategies for both the BERT_{noise} and DeBERTaV3

are presented in Table 6. In BERT_{noise}, first-last pooling only performed well on the performance degradation, but the strategy exhibiting more robustness to detection metrics was max pooling.

In the case of BERT_{clean}, because it did not directly utilize noisy texts, first-last pooling performed well for attacked offensive language, which focuses most on offensiveness and token embeddings from the sentence. However, in the case of BERT_{noise}, where noisy texts were pre-trained into the model weights before the application of pooling strategies, max pooling proved to be the most effective on detection metrics by selecting prominent features among the information to predict the labels invariant to small changes or disturbances.

The results for DeBERTaV3 showed that it performed better with only the model itself when there were no attacks or 30% attacked, while performed better with pooling strategies when the attack rate was increased. We found that the distinct pre-training process of DeBERTaV3, different from the general BERT, was capable of handling a certain level of textual attacks. However, the introduced pooling strategy could prove beneficial in scenarios where the attack rate increases.

The average performances using max pooling and first-last pooling, which had the most impact on each model are presented in Table 7. For the model BERT_{clean} and DeBERTaV3, first-last pool-

Model	Average	
	F1	Δ_{atk}
BERT _{clean}	68.52	-12.85%
BERT _{clean} + max	68.51	-12.88%
BERT _{clean} + first-last	71.24	-10.05%
BERT _{noise}	71.82	-9.81%
BERT _{noise} + max	74.08	-7.64%
BERT _{noise} + first-last	73.57	-7.00%
Ensemble _{hard} (BERT _{clean} + BERT _{noise})	70.99	-11.64%
Ensemble _{soft} (BERT _{clean} + BERT _{noise})	71.23	-11.37%
DeBERTaV3	70.63	-11.88%
DeBERTaV3+max	69.89	-11.53%
DeBERTaV3+first-last	70.78	-11.50%

Table 7: Average from the experimental results of offensive language detection when using max pooling and first-last pooling to the models which exhibited the high scores among the baselines.

ing improved the f1-score by 2.72 and 0.15 and prevented a performance degradation of 2.8% and 0.38%, respectively. In particular, the naive BERT model with the pooling strategies outperformed the ensemble models, and even surpassed the recent model DeBERTaV3 with the pooling strategies. It suggests that the introduced pooling strategies can be employed for a simple structure that solely relies on naive BERT, enabling efficient detection of attacked offensive language without the necessity for the training of noisy texts and strategic pre-training.

In the case of BERT_{noise}, max pooling improved its f1-score by 2.26, reaching a relatively high score of 74.08. In terms of performance degradation, first-last pooling exhibited a performance degradation of 2.81%. Thus, we confirmed that implementing the pooling strategies on a model pre-trained with noisy texts resulted in more robust and improved detection of attacked offensive language exceeded others, such as pooling strategies on a model pre-trained with clean texts. The performances of these pooling strategies are remarkable considering that they do not require additional parameters or direct training on noisy texts.

6 Conclusion

We proposed user-intended adversarial attacks that occur frequently in offensive languages online from the perspective of malicious users. We categorized them into three types: INSERT, COPY, and DECOMPOSE, which add special symbols or exploit the distinct features of the Korean language. The

involvement of attacks significantly affects the tokenization results from the original text.

To address the proposed attacks, we introduced the pooling strategies in a layer-wise manner. This extension utilizes not only the last layer, which focuses on offensiveness, but also the preceding layers, which focus more on token embeddings. The experimental results indicated that first-last pooling was the most robust to the proposed attacks and could even achieve a comparable performance to that of models pre-trained on noisy texts, when applied to models pre-trained on clean texts. We especially demonstrated that rather than the middle layers, the first and last layers can be effectively employed to detect attacked offensive languages.

Furthermore, we experimented with the extent to which the introduced pooling strategies could handle the proposed attacks. We observed that the first-last pooling and max pooling are the most robust, depending on the nature of the texts used for pre-training. It is noteworthy that these strategies, *without* the explicit training of additional parameters or noisy texts, can effectively defend against user-intended adversarial attacks.

Limitations

Despite our efforts to define diverse types of adversarial attacks, there is room for undefined attacks from real-world situations, which may have an unexpected impact on the model performance. Additionally, although we used some language-independent attacks such as inserting special symbols, most of the proposed attacks were based on the characteristics of the Korean language. Therefore, it is necessary to determine whether layer-wise pooling strategies can effectively handle attacked offensive expressions written in other languages for cross-lingual applications. The introduced pooling strategies offer significant flexibility for models in other languages, requiring simple modifications to the model structure.

Ethics Statement

Given our use of offensive representations to describe the proposed attacks, we have included a disclaimer at the beginning of this paper. From the existing offensive language datasets, potential biases regarding race, gender, political issues, and other factors might have been inherent in our experiments. This should be considered when developing our research or expanding it to other languages.

Acknowledgements

This research was supported by Basic Science Research Program through the National Research Foundation of Korea (NRF) funded by the Ministry of Education (NRF-2022R1C1C1008534), and Institute for Information & communications Technology Planning & Evaluation (IITP) through the Korea government (MSIT) under Grant No. 2021-0-01341 (Artificial Intelligence Graduate School Program, Chung-Ang University).

References

- Piush Aggarwal and Torsten Zesch. 2022. Analyzing the real vulnerability of hate speech detection systems against targeted intentional noise. In *Proceedings of the Eighth Workshop on Noisy User-generated Text (W-NUT 2022)*, pages 230–242.
- Sangcheol Ahn and Kyunney Egorova. 2021. Morphophonological patterns of recent korean neologisms. In *Conference on current problems of our time: the relationship of man and society (CPT 2020)*, pages 62–72.
- Ashish Bajaj and Dinesh Kumar Vishwakarma. 2023. Homochar: A novel adversarial attack framework for exposing the vulnerability of text based neural sentiment classifiers. *Engineering Applications of Artificial Intelligence*, 126:106815.
- Nicholas Boucher, Ilia Shumailov, Ross Anderson, and Nicolas Papernot. 2022. **Bad characters: Imperceptible nlp attacks**. In *2022 IEEE Symposium on Security and Privacy (SP)*, pages 1987–2004. IEEE.
- Kyunghyun Cho, Bart van Merriënboer, Çağlar Gulçehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. 2014. Learning phrase representations using rnn encoder–decoder for statistical machine translation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1724–1734.
- Wonik Cho and Soomin Kim. 2021. Google-trickers, yaminjeongeum, and leetspeak: An empirical taxonomy for intentionally noisy user-generated text. In *Proceedings of the Seventh Workshop on Noisy User-generated Text (W-NUT 2021)*, pages 56–61.
- Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Édouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. Unsupervised cross-lingual representation learning at scale. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8440–8451.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186.
- Shreya Goyal, Sumanth Doddapaneni, Mitesh M Khapra, and Balaraman Ravindran. 2023. **A survey of adversarial defenses and robustness in nlp**. *ACM Computing Surveys*, 55(14s):1–39.
- Pengcheng He, Jianfeng Gao, and Weizhu Chen. 2023. Debertav3: Improving deberta using electra-style pre-training with gradient-disentangled embedding sharing. In *The Eleventh International Conference on Learning Representations*.
- Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. 2021. Deberta: Decoding-enhanced bert with disentangled attention. In *International Conference on Learning Representations*.
- Ganesh Jawahar, Benoît Sagot, and Djamé Seddah. 2019. **What does BERT learn about the structure of language?** In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3651–3657, Florence, Italy. Association for Computational Linguistics.
- Younghoon Jeong, Juhyun Oh, Jongwon Lee, Jaimeen Ahn, Jihyung Moon, Sungjoon Park, and Alice Oh. 2022. **KOLD: Korean offensive language dataset**. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 10818–10833, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Di Jin, Zhijing Jin, Joey Tianyi Zhou, and Peter Szolovits. 2020. Is bert really robust? a strong baseline for natural language attack on text classification and entailment. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 8018–8025.
- Woochun Jun. 2020. A study on the cause analysis of cyberbullying in korean adolescents. *International journal of environmental research and public health*, 17(13):4648.
- Bosung Kim, Juae Kim, Youngjoong Ko, and Jungyun Seo. 2021a. Commonsense knowledge augmentation for low-resource languages via adversarial learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 6393–6401.
- Jinyong Kim, Jeonghyeon Kim, Mose Gu, Sanghak Ohh, Gilteun Choi, and Jaehoon Jeong. 2021b. **Hypocrite: Homoglyph adversarial examples for natural language web services in the physical world**.
- Soomin Kim, Changhoon Oh, Wonik Cho, Donghoon Shin, Bongwon Suh, and Joonhwan Lee. 2021c. Trkic g00gle: Why and how users game translation algorithms. *Proceedings of the ACM on Human-Computer Interaction*, 5(CSCW2):1–24.

- Keita Kurita, Anna Belova, and Antonios Anastasopoulos. 2019. [Towards robust toxic content classification](#). *arXiv preprint arXiv:1912.06872*.
- Jean Lee, Taejun Lim, Heejun Lee, Bogeun Jo, Yangsok Kim, Heegeun Yoon, and Soyeon Caren Han. 2022. [K-MHaS: A multi-label hate speech detection dataset in Korean online news comment](#). In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 3530–3538, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.
- Junbum Lee. 2020. Kcbert: Korean comments bert. In *Proceedings of the 32nd Annual Conference on Human and Cognitive Language Technology*, pages 437–440.
- J Li, S Ji, T Du, B Li, and T Wang. 2019. Textbugger: Generating adversarial text against real-world applications. In *26th Annual Network and Distributed System Security Symposium*.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.
- Xiaofei Ma, Zhiguo Wang, Patrick Ng, Ramesh Nallapati, and Bing Xiang. 2019. Universal text representation from bert: An empirical study. *arXiv preprint arXiv:1910.07973*.
- Justin McCurry. 2022. [South korea under pressure to crack down on cyberbullying after high-profile deaths](#). *The Guardian*.
- Jihyung Moon, Won Ik Cho, and Junbum Lee. 2020. [BEEP! Korean corpus of online news comments for toxic speech detection](#). In *Proceedings of the Eighth International Workshop on Natural Language Processing for Social Media*, pages 25–31, Online. Association for Computational Linguistics.
- Dongsuk Oh, Yejin Kim, Hodong Lee, H. Howie Huang, and Heuseok Lim. 2022. [Don't judge a language model by its last layer: Contrastive learning with layer-wise attention pooling](#). In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 4585–4592, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.
- Chaewon Park, Soohwan Kim, Kyubyong Park, and Kunwoo Park. 2023a. K-haters: A hate speech detection corpus in korean with target-specific ratings. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 14264–14278.
- San-Hee Park, Kang-Min Kim, O-Joun Lee, Youjin Kang, Jaewon Lee, Su-Min Lee, and SangKeun Lee. 2023b. [“why do I feel offended?” - Korean dataset for offensive language identification](#). In *Findings of the Association for Computational Linguistics: EACL 2023*, pages 1142–1153, Dubrovnik, Croatia. Association for Computational Linguistics.
- Sungjoon Park, Jihyung Moon, Sungdong Kim, Won Ik Cho, Ji Yoon Han, Jangwon Park, Chisung Song, Junseong Kim, Youngsook Song, Taehwan Oh, et al. 2021. Klue: Korean language understanding evaluation. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*.
- Shuhuai Ren, Yihe Deng, Kun He, and Wanxiang Che. 2019. Generating natural language adversarial examples through probability weighted word saliency. In *Proceedings of the 57th annual meeting of the association for computational linguistics*, pages 1085–1097.
- Suriya Saengprang and Savitri Gadavani. 2021. Cyberbullying: The case of public figures. *LEARN Journal: Language Education and Acquisition Research Network*, 14(1):344–369.
- Nihar Sahoo, Niteesh Mallela, and Pushpak Bhattacharyya. 2023. [With prejudice to none: A few-shot, multilingual transfer learning approach to detect social bias in low resource languages](#). In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 13316–13330, Toronto, Canada. Association for Computational Linguistics.
- Mike Schuster and Kuldip K Paliwal. 1997. Bidirectional recurrent neural networks. *IEEE transactions on Signal Processing*, 45(11):2673–2681.
- Jaejung Song. 2006. *The Korean language: Structure, use and context*. Routledge.
- Max Wolff and Stuart Wolff. 2020. [Attacking neural text detectors](#). *ICLR 2020 Workshop Towards Trustworthy ML*.
- Zhelun Wu, Nishant Kambhatla, and Anoop Sarkar. 2018. [Decipherment for adversarial offensive language detection](#). In *Proceedings of the 2nd Workshop on Abusive Language Online (ALW2)*, pages 149–159, Brussels, Belgium. Association for Computational Linguistics.
- Hyunyoung Yi and Sangmi Cha. 2020. [Cyber bullying, star suicides: The dark side of south korea's k-pop world](#). *Reuters*.
- Marcos Zampieri, Shervin Malmasi, Preslav Nakov, Sara Rosenthal, Noura Farra, and Ritesh Kumar. 2019. [Predicting the type and target of offensive posts in social media](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 1415–1420, Minneapolis, Minnesota. Association for Computational Linguistics.

A User-Intended Adversarial Attacks Details

A.1 INSERT

We utilized three INSERT types by adding special symbols that are not complete characters. The detailed methods and examples are as follows:

- **INSERT_{zz}**: We randomly added between 2~5 sounds of ‘ㄷ’ to the word (on the keyboard, the sound ‘ㄷ’ corresponds to the alphabet ‘z’). This is a commonly used expression in online, conveying a meaningless and somewhat frivolous tone. We placed it not only between characters but also in instances where the final sound of a specific character was empty.
- **INSERT_{space}**: We randomly added a single space to the word. We expected the same impact as the intentions of malicious users.
- **INSERT_{special}**: We randomly added a special character to the word: one of ‘~’, ‘!’, ‘@’, ‘1’, or ‘2’. We also expected the same impact as the intentions of malicious users.

INSERT Attacks	Text Examples
original text	틀딱이냐? (okay, boomer?) 기레기 여기 있었네 (presstitute right here)
INSERT _{zz}	틀ㄷㄷ딱이냐? 기렉기 ㄷㄷ기 여기 있었네
INSERT _{space}	틀 딱이냐? 기레 기 여기 있었네
INSERT _{special}	틀@딱이냐? 기2레기 여기 있었네

Table 8: Text examples of user-intended adversarial attacks with the types of INSERT.

A.2 COPY

We utilized three COPY types based on the distinctive characteristics of the Korean language. The detailed methods and examples are as follows:

- **COPY_{initial}**: We copied the initial sound of the character to the final sound of the preceding character. For example, if ‘기’ is chosen from ‘기레기’, the initial sound ‘ㄱ’ would be copied as the final sound to the preceding ‘레’. Consequently, ‘레’ is transformed into ‘렉’, leading to ‘기렉기’.
- **COPY_{middle}**: We copied the middle sound of the character, onto the newly added character.

If the selected character had a final sound, it was also included. For example, if ‘딱’ is chosen from ‘틀딱’, the middle sound ‘ㅌ’ would be copied as the following character with the final sound ‘ㄱ’. Consequently, ‘악’ is newly added, leading to ‘틀따악’.

- **COPY_{final}**: We copied the final sound of the character to the initial sound of the following character. For example, if ‘있’ is chosen from ‘있었네’, the final sound ‘ㅌ’ would be copied as the initial sound to the following ‘었’. Consequently, ‘었’ is transformed into ‘졌’, leading to ‘이졌네’.

COPY Attacks	Text Examples
original text	틀딱이냐? (okay, boomer?) 기레기 여기 있었네 (presstitute right here)
COPY _{initial}	틀딱인냐? 기렉기 여기 있었네
COPY _{middle}	틀따악이냐? 기레기 여기이 있었네
COPY _{final}	틀딱기냐? 기레기 여기 이졌네

Table 9: Text examples of user-intended adversarial attacks with the types of COPY.

A.3 DECOMPOSE

We utilized two DECOMPOSE types based on the distinctive characteristics of the Korean language. The detailed methods and examples are as follows:

- **DECOMPOSE_{final}**: We decomposed the final sound of the character into the newly added sound. For example, if ‘딱’ is chosen from ‘틀딱’, the final sound ‘ㄱ’ would be decomposed as the following sound. Consequently, the expression will be ‘틀따ㄱ’.
- **DECOMPOSE_{all}**: We decomposed all the sounds of the character. For example, if ‘틀’ is chosen from ‘틀딱’, it would be decomposed as the initial, middle, and final sounds. Consequently, the expression will be ‘ㅌ ㅡ ㄷ 딱’.

DECOMPOSE Attacks	Text Examples
original text	틀딱이냐? (okay, boomer?) 기레기 여기 있었네 (presstitute right here)
DECOMPOSE _{final}	틀따ㄱ이냐? 기레기 여기 이ㅌ었네
DECOMPOSE _{all}	ㅌ ㅡ ㄷ 딱이냐? 기르 ㄷ기 여기 있었네

Table 10: Text examples of user-intended adversarial attacks with the types of DECOMPOSE.

A.4 How the Attack Works

In our experiments, the attacks were only applied to the test set to evaluate the robustness of the model against user-intended adversarial attacks. We randomly selected one of all the attacks according to the attack rates. For instance, with a 30% attack rate and 9 words in a sentence, we attacked 3 words (30% of them), with each word randomly reflecting one of the proposed attacks.

A.5 Experiments on Each Attack

We further investigated the impact for each attack type, selecting only one of the attacks from only INSERT, COPY, or DECOMPOSE. We chose BERT_{clean} and the model that applied first-last pooling, which exhibited the best performance among the layer-wise pooling strategies.

The performances according to the attack types are presented in Table 11. It is noteworthy that a high performance on a specific attack type indicates that it is easier to defend than others, and vice versa. INSERT proved to be easier than all attacks, with COPY only marginally harder than INSERT. However, DECOMPOSE was more difficult than these attacks, exhibiting a performance degradation compared to all attacks. Therefore, we revealed that adversarial attacks that reflect the characteristics of the Korean language are more challenging than those that simply add special symbols that could be adapted language independently.

Model	All Attacks			Only INSERT		
	P	R	F1	P	R	F1
BERT _{clean}	77.74	66.38	67.96	78.12	68.80	70.66
BERT _{clean} + first-last	77.58	69.38	71.21	77.51	70.76	72.54
Model	Only COPY			Only DECOMPOSE		
	P	R	F1	P	R	F1
BERT _{clean}	77.86	68.47	70.29	78.38	62.23	65.67
BERT _{clean} + first-last	77.33	70.50	72.26	77.00	66.47	68.04

Table 11: Experimental results of offensive language detection when the attack ratio is 60% for each type of user-intended adversarial attacks.

B Experimental Details

B.1 Implementation Details

We used 6 layers with an embedding dimension of 768 and a dropout ratio of 0.1 for the RNN-based models, containing BiLSTM and BiGRU. As for the BERT-based models, containing BERT_{clean}, BERT_{noise}, their ensembles, and DeBERTaV3, we fine-tuned 12 pre-trained layers with an embedding dimension of 768 and a dropout ratio of 0.2.

We used the AdamW optimizer with a learning rate of 1e-5, trained the models for 1~5 epochs, and considered the epoch with the lowest validation loss or the last epoch. We set a batch size of 32 for all the models. The models were implemented using PyTorch and NVIDIA GeForce RTX 3090 GPU. We also used the HuggingFace library to leverage the weights of the pre-trained BERT models.

B.2 Metrics

We collected the datasets to train as many types of offensive languages as possible. There is some label imbalance, therefore, we used macro precision, recall, and f1-score to address this issue.

Δ_{atk} represents the performance degradation in the f1-score. For example, denoting the existing f1-score as $F1_{original}$ and the f1-score with the attacks as $F1_{attacked}$, it is computed as follows:

$$\Delta_{atk} = (F1_{attacked} - F1_{original}) / F1_{original} * 100. \quad (5)$$

C Examples of the Attacks

We presented the following cases where offensive languages contained user-intended adversarial attacks on popular online communities in South Korea, such as Twitter⁵ and Dcinside⁶.

- INSERT_{zz}:
지이지에 개 병ㅋㅋㅋ신
(ziezie so dumbzzzfuck)
- INSERT_{space}:
완전 또 라이 아니야...
(crazy ass we irdo)
- INSERT_{special}:
넥슨아 이거 개같으면 지워라 새@끼들아
(delete this if you dare Nexon you sons of bit@ches)
- COPY_{initial}:
우리나라 남자새끼들 진짜 재미없게 산다
(korean lossser men live in the most boring ass lives)
- COPY_{middle}:
역쉬나 기레에기들..
(presstiiitutes of course)
- COPY_{final}:
현생 어떻게 살아가누 너가튼계?
(how the fuck dooyou live in the real world?)
- DECOMPOSE_{final}:
과몰입 중인 병신 트크딱들ㅋㅋ
(chronically online dumbass bo omers lmao)
- DECOMPOSE_{all}:
너 같은 ㄱ 새끼는 몽둥이로 맞아 죽어야 해
(p i ece of shit like you need to be beaten to death)

⁵<https://twitter.com>

⁶<https://dcinside.com>