

From Universal Language Model to Downstream Task: Improving RoBERTa-Based Vietnamese Hate Speech Detection

Quang Huu Pham*, Viet Anh Nguyen*, Linh Bao Doan, Ngoc N. Tran
R&D Lab, Sun Asterisk Inc

{pham.huu.quang, nguyen.viet.anhd, doan.bao.linh, tran.ngo.quang.ngoc}@sun-asterisk.com

Ta Minh Thanh[†]

Le Quy Don Technical University, 236 Hoang Quoc Viet, Bac Tu Liem, Ha Noi
thanhtm@mta.edu.vn

* Equal contribution [†] Corresponding author

Tóm tắt nội dung—Natural language processing (NLP) is a fast-growing field of artificial intelligence. Since the Transformer [32] was introduced by Google in 2017, a large number of language models such as BERT, GPT, and ELMo have been inspired by this architecture. These models were trained on huge datasets and achieved state-of-the-art results on natural language understanding. However, fine-tuning a pre-trained language model on much smaller datasets for downstream tasks requires a carefully-designed pipeline to mitigate problems of the datasets such as lack of training data and imbalanced data. In this paper, we propose a pipeline to adapt the general-purpose RoBERTa language model to a specific text classification task: Vietnamese Hate Speech Detection. We first tune the PhoBERT¹[9] on our dataset by re-training the model on the Masked Language Model (MLM) task; then, we employ its encoder for text classification. In order to preserve pre-trained weights while learning new feature representations, we further utilize different training techniques: layer freezing, block-wise learning rate, and label smoothing. Our experiments proved that our proposed pipeline boosts the performance significantly, achieving a new state-of-the-art on Vietnamese Hate Speech Detection (HSD) campaign² with 0.7221 F1 score.

Index Terms—Hate Speech Detection (HSD), RoBERTa, Text Classification, Natural Language Processing, Text Mining.

I. INTRODUCTION

A. Overview

The rapid growth of the Internet, social media, and community forums have allowed people across the world to connect instantaneously and has revolutionized communication as well as content issues. However, the increase of hate speech on these platforms has drawn significant expenditure from governments, organizations, companies, and researchers. The term “hate speech” can be understood as any kind of communication that uses pejorative or discriminatory language with reference to a person or a group based on their religion, gender, ethnicity,

nationality, race, colour, descent, or other identity factors. Multiple statistics reports [20] and books [28] also show that hate speech and crime are highly correlated and on the rise together. Since the internet gives people some degree of anonymity, some take this for granted and abuse it to harass others. Calling names, making distasteful comments about one’s origin, or simply shaming someone in anyway, are all everyday examples of hate speech that anyone will encounter occasionally. To combat this, a vast number of methods have been studied and developed for automated HSD. This aims to classify textual content into hate and non-hate speech.

By the end of 2019, social network site users in Vietnam have reached 48 million users. Still, there has been limited available research about Vietnamese HSD; building appropriate countermeasures for hate speech requires detecting and tracing through content. In the case of the Vietnamese language, this task becomes difficult due to the diverse vocabulary and complex grammar. For example, the same subword can have multiple meanings, which is different from Latin word roots. Or, the fact that the same base (sub)word having multiple possible intonations creates many hindrances in the path of language understanding: not only each subword has multiple vastly different meanings, but also this causes an infinite number of combinatorial possibilities of either misspelling or intentional shorthand expressions.

The variety of semantics and grammar in the Vietnamese language has led to a big challenge for automatic hate speech detection. Previous researches of text classification based on traditional machine learning algorithms or training deep learning from scratch is often inefficient. It also requires a lot of effort for pre-processing and assimilating the semantic of words. In addition, recently proposed pre-trained language models have accomplished success in multiple tasks of natural language processing through fine-tuning when integrated with the model of downstream tasks. The emergence of pre-trained language models has contributed to new ideas for solving significant problems. Pre-trained language models are large scale neural

¹PhoBERT is a pre-trained RoBERTa model which is known as a state-of-the-art language model for Vietnamese provided by VinAI research: <https://github.com/VinAIResearch/PhoBERT>

²<https://vlsp.org.vn/vlsp2019/eval/hsd>

network models based on the deep Transformer structure. Their initial parameters are learned through immense self-supervised training, then combined with multiple downstream models to fix special tasks by fine-tuning. Empirical experiment’s results show that the downstream tasks’ performances of these kinds of models are usually better than those of conventional models.

B. Our contributions

In this paper, we investigate many experiments in fine-tuning the pre-trained RoBERTa model for text classification tasks, specifically Vietnamese HSD. We propose a general pipeline and model architectures to adapt the universal language model as RoBERTa for downstream tasks such as text classification. With our technique, we achieve new state-of-the-art results on the Vietnamese Hate Speech campaign, organized by VLSP 2019³.

The main contributions of our paper are as follows:

- We propose a general pipeline to adopt the universal language model as a pre-trained RoBERTa for the text classification. It includes two steps: (1) Re-training masked language model task on training data of the classification task. (2) Fine-tuning model with a new classification head for the target task.
- We conduct multiple methods to design model architecture for text categorization task by using the pre-trained RoBERTa model such as PhoBERT[9].
- A number of training techniques are suggested that can improve the efficiency of the fine-tuning phase in solving data problems. These techniques are practical and can empirically help the model prevent overfitting phenomenon in the absence of training data or data imbalances.
- From PhoBERT to Vietnamese HSD: We have achieved the new state-of-the-art results on Vietnamese HSD task by utilizing PhoBERT and our proposed method.

C. Roadmap

The rest of the paper is organized as follows. Section 2 provides a brief survey of related work. Next, in Section 3, we introduce our proposed method, the model architecture, and our conducted fine-tuning strategies. Experiments are described in Section 4, including dataset information, data processing method, and some other experimental settings. Section 5 shows our experimental results. Finally, Section 6 presents our conclusions.

II. RELATED WORK

A. Language Models

Language modeling has been becoming an essential part of modern NLP field. It helps computers understand human language by digitalizing qualitative information. Early studies on word embeddings try to construct static representations for words, that is, context-independent embeddings. Mikolov *et al.* [22] introduced novel architectures with an open-source

package called *Word2Vec*. The architectures consist of two models: Continuous Bag of Words (CBOW) and Skip-gram model were trained on a huge dataset of 1.6 billion words. GloVe [27] is an unsupervised learning algorithm for context-independent word embedding extraction. It first creates co-occurrence matrix of words and then factorizes it to extract dense word vectors. However, GloVe and Word2Vec are both fail in representing rare or out-of-vocabulary words. To mitigate this problem, fastText [23] decomposes each word as a sum of character n-grams. This handles unseen words very well because their character n-grams still occur in other words. In contrast to context-independent embeddings, contextualized word embeddings aim to encode word semantics within contexts. Marking a new era of NLP, Bidirectional Encoder Representations from Transformers [11] or BERT achieves new state-of-the-art performance on eleven NLP tasks, outperforms previous best result (with GLUE score of 80.4%, 7.6% improvement) and human on SQuAD 1.1 (with 93.2% accuracy, 2% higher). The large model consists of 24 Transformer blocks for a total of 340M parameters and was trained on 3.3 billion words corpus. GPT-2 [1] by OpenAI is even a larger model with 1.5 billion parameters and 48 layers. Being trained on a huge, diverse and well-processed dataset, GPT-2 achieves state-of-the-art results on 7 out of 8 datasets. Facebook research team propose an improved training procedure for BERT, called RoBERTa [17]. The improvements include a ten-times larger dataset, longer training, increased batch size, using byte-level encoding with bigger vocabulary, excluding next sentence predicting task and dynamic masking pattern changing. For Vietnamese language modeling, PhoBERT [9] is the current state-of-the-art. The model is based on RoBERTa with two versions “base” and “large”. Both versions outperform previous best on different downstream tasks.

B. Hate Speech Detection (HSD)

HSD can be understood as a type of text classification. Before the deep learning era, traditional methods had been widely studied for this notorious problem. These approaches require manual feature engineering to encode text sequences in vector form, which then be fed into classifiers such as Random Forests [35], Naive Bayes [5], [21], and Support Vector Machine [3], [21], [35]. *Bag-of-words (BoW)* is a way to represent documents by modeling the occurrence of every word. This has been reported to be a discriminative feature for HSD [3], [4], [5], [14], [30], [31], [34], [36]. However, modeling words suffers from the problem of typing error, which generates noisy words to the corpus. *Character n-gram* divides a word into sub-words in order to suppress typing error parts of the word. In [21], the authors have systematically proved that character-level n-gram representation outperforms BoW in abusive language detection. Based on the idea that hate speech usually comes along with negative sentiment, the methods in [16], [25], [30] utilize *sentiment analysis* as an evidence for HSD. *Linguistic features* inject additional useful knowledge to the raw text. Xu *et al.* [36] combine Stanford CoreNLP [19] POS information with n-gram features. However, the POS tokens did not improve the

³The sixth international workshop on Vietnamese Language and Speech Processing.

performance of the classifiers. In contrast, [3], [4], [5], [25], [26] successfully employed dependency relationships in their feature set and report significant performance improvements. Because hate speech usually comes from social media, it is feasible to extract meta-information such as how likely a user will post hateful speech in the future [35], the number of profanity in a user’s activity history[7] or user gender [6], [34].

Recent studies are paying great attention to deep learning approaches. It not only boosts performance considerably, but also requires less effort on feature engineering. The input of a deep learning model can simply be an one-hot encoding of text sequences [2], [13], [38], [10], then meaningful features are learned by Convolutional Neural Networks (CNN) [13], Long Short-Term Memory (LSTM) [2], [10] or even a combination of CNN and LSTM [38]. However, one-hot representation suffers from issues of high dimensionality as its length equals to the vocabulary size. Therefore, a more convenient way is to embed the input into a low-dimensional space. This can be character embeddings [21], comment embeddings [12] or text embeddings from a two-phase deep learning model [37].

HSD Shared Task in VLSP Campaign 2019 [33] is a challenge for detecting Vietnamese hate speech on social networks. Our team, the winning team, using logistic regression with ensemble learning n-gram, achieved F1 score of 0.67756 and 0.61971 on public-test and private-test, respectively. The second best team also utilized logistic regression with ensemble learning, however the input includes more feature types: n-grams of words, part-of-speech tags, and numeric features. They scored 0.58883 on private-test, that is 3% lower than the winner. The other teams employed deep learning architectures from CNN, LSTM, RNN to Bi-LSTM, LSTMCNN and Bi-GRU-LSTM-CNN, achieving scores from 0.51705 to 0.58455 on private-test.

III. PROPOSED METHOD

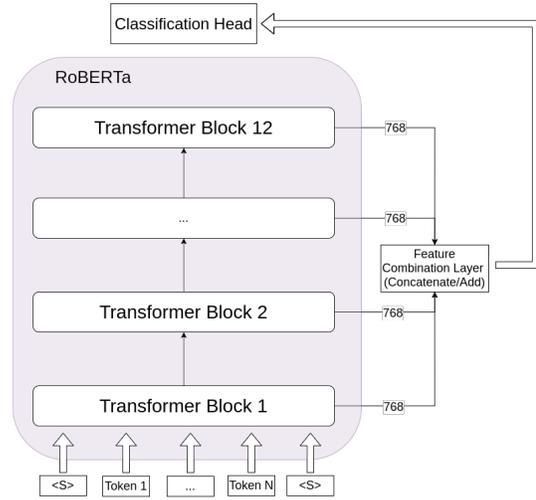
A. RoBERTa-based Network for HSD Task

Our architecture utilizes PhoBERT as a backbone network with some modifications. The output of each Transformer block represents a different semantic level for the inputs. In our experiments, we use different combinations of outputs of those Transformers blocks. The general architecture model is shown in Figure 1.

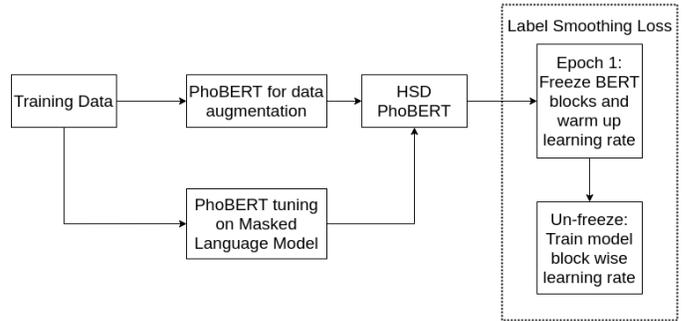
The features are combined across outputs of multiple transformer blocks by concatenating or adding are fed to the classification head. A simple classification head is a fully connected network without hidden layers.

B. Proposed Fine-Tuning Strategies

1) *Fine-Tuning Masked Language Model:* The pre-trained model was fit to a much larger dataset of a completely different domain. Therefore, while it might work very well in general, the vanilla pre-trained model will underperform at our one specific task. This induces the need for us to tune the model to our needs. Therefore, with the existing weights of PhoBERT as a starting point, we train our model with our domain-specific training data on Masked Language Modeling (MLM) task.

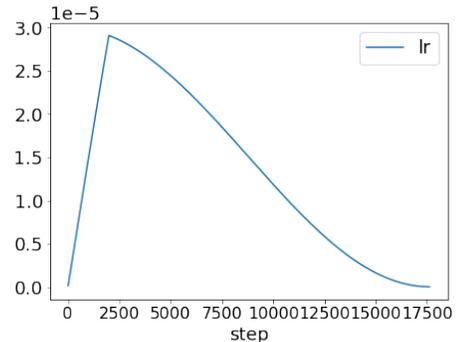


Hinh 1. An illustration of our architecture. The input is a tokenized sentence, feeding parallelly to RoBERTa-base. Each of the twelve transformer blocks of the RoBERTa produces a 768-D vector. These vectors are then concatenated to form a long vector for the classification head. We investigate on the performance of different combinations of these vectors.



Hinh 2. The proposed fine-tuning strategies for Hate Speech detection.

Moreover, for such a large model to be trained successfully without either forgetting its good initialization (by gradient descending too far from it) or failing to converge (due to deep models being bad at propagating through further layers), we use a warm-up learning rate scheduling scheme. Originating from paper [15] under the name of “slanted triangular learning rates”,



Hinh 3. Warm-up learning rate

the main purpose of this method is to make the model converge more quickly for a suitable region of the parameter space in the beginning process of training.

2) *Optimization Strategy: Freeze or not Freeze:* Each layer in RoBERTa network captures different levels of context. Specifically, lower layers produce global embeddings for words while the embeddings from upper layers are more context-specific. We would like to preserve the global knowledge while adjusting context representations for our classification model.

For the first few epochs, we only keep the fully connected layers that are responsible for classifying text sequences, and the RoBERTa part is frozen. This allows the model to learn a decent decision for the task. After these epochs, we unfreeze all layers and set different learning rates for different layers: the deeper the layer is, the more learning rate increased. This prevents the model from forgetting the useful global meaning of the words while forcing it to learn the context of the domain.

3) *Label Smoothing:* For such a large model to be fit on a relatively small set of data, the model tends to become overconfident of its performance, going to the dark side of overfitting. To avoid this, we employ label smoothing, which softens our one-hot ground truth labels. Specifically, for a model outputting probabilities y_k of K classes, instead of labelling our ground truth with one-hot encoding:

$$y_k = \begin{cases} 1 & \text{if } k = j \text{ for some } j, \\ 0 & \text{otherwise,} \end{cases} \quad (1)$$

we slightly rebalance the target distribution so that it becomes less “binarized” by smoothing the probabilities with,

$$y'_k = y_k(1 - \alpha) + \alpha/K \quad (2)$$

for some smoothing parameter α . As a result, this technique teaches the model to have some uncertainty in its guesses, and reduces the severity of overfitting. Moreover, since we are fine-tuning on a pretrained model, the original output probability vector of the model is close to an one-hot. This introduces numerical instabilities if the new true label is also one-hot, since with cross-entropy as the loss function, the loss becomes $1 \log 0 = -\infty$. Thus, being used here, label smoothing acts as a small perturbation in our numerical method, making the training process more stable, helping our model converge better.

IV. EXPERIMENTS

A. Dataset

To perform experiments on hate speech detection task, we use the HSD dataset from the VLSP workshop campaign in 2019 [29]. The dataset includes 25,431 samples; all were crawled from posts and comments on Facebook, and annotated to one of three classes: hate speech (HATE), offensive but not hate speech (OFFENSIVE), neither offensive nor hate speech (CLEAN) by many annotators. The training data consists of 20,345 items with label distribution is showed in Table I.

In general, this is an imbalanced dataset and contains a lot of noisy data. There is an extremely enormous gap between the number of CLEAN data points and those from the other

Bảng I
TRAINING DATA DISTRIBUTION

	HATE	OFFENSIVE	CLEAN
Number of sample	709	1,022	18,614

two categories, which lead to bad results for some methods, especially Deep Learning approaches.

Because the dataset was crawled directly from users’ posts and comments on social networks, it has some notable properties such as abbreviations, emoji, special characters, foreign language, teen code, typing errors, etc. This has led to a significant challenge because PhoBERT or some other pre-trained language model are often trained on normal clean data such as Wikipedia data or Vietnamese news.

B. Data preprocessing

We have built a basic preprocessing module to process the raw text before feeding it into the model. It includes Unicode normalization, lowercase conversion, and the replacement of all emoji characters to the EMOJI label. Private personal information also is masked for privacy purposes such as phone numbers and email addresses. Finally, raw text is segmented by a word segmenter. As PhoBERT employed the RDRSegmenter [8] from VnCoreNLP⁴ to pre-process training data of the language model task, it is recommended to use the same word segmenter for PhoBERT-based downstream applications.

C. Experimental Settings

1) *Data augmentation:* To partially prevent data imbalance and less data for two HATE and OFFENSIVE labels, we used PhoBERT Masked Language Model as a method to augment data. From the original text, we randomly selected one token, replaced it with <mask>, and used the PhoBERT MLM to fill the mask. Repeated 5 times, we got a brand new sample by fill up 5 tokens <mask>from the original text.

2) *Model Training:* We fine-tuned PhoBERT with Masked Language Model on HSD dataset with ratio of 15% masked tokens. As in BERT default settings, batch size is fixed to 32, with 10,000 steps. We chose Adam for the optimization and trained the model on a single GPU with learning rate of 3×10^{-5} . Figure 2 illustrates the overall pipeline. After tuning PhoBERT process completed, we retrain HSD model with one or combination of multiple output features from numerous transformer blocks.

These features then would be incorporated into a Classification Head, which was designed straightforward with Fully Connected and Dropout layers. AdamW [18] were selected to be the optimizer for BERT blocks and Classification Head with different learning rates. BERT learning rate was 1×10^{-5} while classification head learning rate was larger at 1×10^{-4} . Except bias and LayerNorm being excluded from weight decay, this ratio was chosen with value 0.01. Num_warmup_steps in experiments is equal to 1/8 of total

⁴<https://github.com/vncorenlp/VnCoreNLP>

steps in 1 epoch.

We utilized Stratified K-fold with $k = 10$, the ratio of samples in train dataset and valid dataset were preserved as in the original dataset. Each and every fold was trained for 10 epochs.

3) *Loss function*: We use Label Smoothing loss which is the combination of Cross-Entropy Loss and Label Smoothing. The smoothing rate is set to 0.2 as this shows prominent after the first few epochs. We slightly lower loss target values while increasing the target value 0. These so-called soft targets will give lower loss when there is an incorrect prediction and consequently, our model will penalize low entropy predictions or “confidence penalty” as mention in [24] and the section above.

The cross entropy loss function is calculated as follows:

$$\mathcal{L}(\mathbf{x}', \mathbf{x}) = - \sum_i x'_i \log x_i, \quad (3)$$

where \mathbf{x}' is the true distribution of any particular data point (one-hot encoded, possibly with label smoothing applied), and \mathbf{x} is the model’s predicted distribution. Specifically, for our prediction task with a dictionary of c possible words to choose from, x_i is the probability that the next word in the sequence is the i -th token in the dictionary. The loss for the whole dataset (or batch) is the sum (or mean) of the losses of individual data points.

D. System configuration

Our experiments are conducted on a computer with Intel Core i7 9700K Turbo 4.9GHz, 32GB of RAM, GPU GeForce GTX 2080Ti, and 1TB SSD hard disk.

V. EXPERIMENTAL RESULTS

A. Evaluation metrics

Macro F1-score is a common evaluation metrics for classification tasks. F1-score is the harmonic mean of *Precision* and *Recall*.

F1 score: performance measure for classification

$$F_1 = \frac{2}{Recall^{-1} + Precision^{-1}}, \quad (4)$$

where *Precision* is the number of correct positive results divided by the number of all positive results, and *Recall* is the number of correct positive results divided by the number of all samples that should have been identified as positive.

F1-macro or macro-averaged F1 score is computed as mean of F1 scores for each class

$$Macro F_1 = \frac{F_{1HATE} + F_{1OFFENSIVE} + F_{1CLEAN}}{3} \quad (5)$$

Bảng II
MEAN OF MACRO F_1 SCORE ON STRATIFIED K-FOLD WITH $k = 10$ OF DIFFERENCE BLOCKS

Feature blocks	Mean of F1 score
Layer 6 (only single block)	0.6854
Layer 12 (only single block)	0.6978
Layer 3-6 (4 middle blocks)	0.6855
Layer 9-12 (4 last blocks)	0.6989
Layer 1-6 (6 first blocks)	0.6905
Layer 7-12 (6 last blocks)	0.6989
Layer 1-12 (all blocks)	0.6979

Bảng III
SOME SAMPLES THAT OUR SYSTEM GOOD AND FAILURE CLASSIFIED

Text	Model prediction	Truth label
"học kỹ cuối như đồ thị hình sóng thần"	CLEAN	CLEAN
"lan anh đm ám ảnh vì"	OFFENSIVE	OFFENSIVE
"nguyễn hoàng bách đm a môm lol à"	HATE	HATE
"hay vat đúng lắm a"	OFFENSIVE	CLEAN
"có đeo đầu mà xem vì"	CLEAN	OFFENSIVE
"hường lily mặt ngu vì"	OFFENSIVE	HATE

B. Our results

Our experiment results are shown in Table II and Table IV.

Table II indicates the mean of Macro F_1 score on using feature from different BERT blocks. We take the test on multiple selected blocks including single block (layer 6, 12), middle blocks (layer 3-6), last blocks (layer 6-12, 7-12, 9-12) and all blocks. Our retrieved result with last blocks by far is the top score. Last 4 blocks (layer 9-12) and last 6 blocks (layer 7-12) both achieve the highest F1 score at 0.6989.

Table IV demonstrates effectiveness of each and every fine-tuning methods. Each individual technique boosts performance of the model by 0.5-1.5% in term of F1 score while the combination of these methods significantly enhances the score up to 0.7211, outperforming the winner of this challenge (0.67756 F1 score) and the current best result on the public leaderboard (0.71432 F1 score with a single model). Some samples that our system good and failure classified are shown in Table III. It shows that our proposed method is efficient for task of HSD.

VI. CONCLUSION

In this paper, we have explored and proposed our pipeline to solve the Vietnamese Hate Speech Detection task by using a pre-trained universal language model. By intensively conducting experiments using PhoBERT, we have demonstrated that RoBERTa and our fine-tuning strategy is highly effective in text classification tasks. With our proposed methods, we have achieved new state-of-the-art results on the Vietnamese HSD campaign. For future work, we would like to explore different classification head architectures. Instead of only using the fully connected layer, we will experiment with other architectures for instance LSTM, RNN, and CNN-LSTM on top.

ACKNOWLEDGMENT

This work is partially supported by *Sun-Asterisk Inc*. We would like to thank our colleagues at *Sun-Asterisk Inc* for

Bảng IV

MEAN OF MACRO F1 SCORE ON STRATIFIED K-FOLD WITH K = 10 WITH CONCATENATE OF LAYERS 6-12 AND OUR TRAINING APPROACH

Proposed training approach	Mean of F1 score
Cross entropy loss	0.6922
Label Smoothing loss	0.7005
Non warm-up learning rate	0.6989
Warm-up learning rate	0.7062
Non Fine-tune MLM	0.6989
Fine-tune MLM	0.7162
Non Block wise learning rate	0.7051
Block wise learning rate	0.7079
Combine all the methods	0.7211

their advice and expertise. Without their support, this experiment would not have been accomplished. We also express our gratitude to Tiep Vu, founder of AIVIVN⁵ for supporting us in evaluating results on aivivn.com.

TÀI LIỆU

- [1] R. Alec, J. Wu, C. Rewon, L. David, A. Dario, and S. Ilya. Language models are unsupervised multitask learners. 2019.
- [2] Pinkesh Badjatiya, Shashank Gupta, Manish Gupta, and Vasudeva Varma. Deep learning for hate speech detection in tweets. *CoRR*, abs/1706.00188, 2017.
- [3] Pete Burnap and Matthew Williams. Cyber hate speech on twitter: An application of machine classification and statistical modeling for policy and decision making: Machine classification of cyber hate speech. *Policy & Internet*, 7, 04 2015.
- [4] Pete Burnap and Matthew L. Williams. Us and them: identifying cyber hate on twitter across multiple protected characteristics. *Epj Data Science*, 5, 2016.
- [5] Y. Chen, Y. Zhou, S. Zhu, and H. 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.
- [6] M. Dadvar, F. Jong, R. Ordelman, and D. Trieschnigg. Improved cyberbullying detection using gender information. 01 2012.
- [7] M. Dadvar, D. Trieschnigg, R. Ordelman, and F. Jong. Improving cyberbullying detection with user context. pages pp 693–696, 03 2013.
- [8] N. Q. Dat, N. Q. Dai, V. Thanh, M. Dras, and M. Johnson. A Fast and Accurate Vietnamese Word Segmenter. In *Proceedings of the 11th International Conference on Language Resources and Evaluation (LREC 2018)*, pages 2582–2587, 2018.
- [9] N. Q. Dat and N. A. Tuan. Phobert: Pre-trained language models for vietnamese, 2020.
- [10] Fabio Del Vigna, Andrea Cimino, Felice Dell’Orletta, Marinella Petrocchi, and Maurizio Tesconi. Hate me, hate me not: Hate speech detection on facebook. 01 2017.
- [11] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding, 2018.
- [12] Nemanja Djuric, Jing Zhou, Robin Morris, Mihajlo Grbovic, Vladan Radosavljevic, and Narayan Bhamidipati. Hate speech detection with comment embeddings. In *Proceedings of the 24th International Conference on World Wide Web, WWW ’15 Companion*, page 29–30, New York, NY, USA, 2015. Association for Computing Machinery.
- [13] 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, Vancouver, BC, Canada, August 2017. Association for Computational Linguistics.
- [14] Homa Hosseinmardi, Sabrina Arredondo Mattson, Rahat Ibn Rafiq, Richard Han, Qin Lv, and Shivakant Mishra. Detection of cyberbullying incidents on the instagram social network, 2015.
- [15] Jeremy Howard and Sebastian Ruder. Universal language model fine-tuning for text classification, 2018.
- [16] D. Karthik, J. Birago, H. Catherine, L. Henry, and P. Rosalind. Common sense reasoning for detection, prevention, and mitigation of cyberbullying. *ACM Trans. Interact. Intell. Syst.*, 2(3), September 2012.
- [17] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach, 2019.
- [18] I. Loshchilov and F. Hutter. Decoupled weight decay regularization, 2017.
- [19] Christopher Manning, Mihai Surdeanu, John Bauer, Jenny Finkel, Steven Bethard, and David McClosky. The Stanford CoreNLP natural language processing toolkit. In *Proc. of 52nd Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pages 55–60. Association for Computational Linguistics, June 2014.
- [20] W. L. Matthew, B. Pete, J. Amir, L. Han, and O. Sefa. Corrigendum to: Hate in the Machine: Anti-Black and Anti-Muslim Social Media Posts as Predictors of Offline Racially and Religiously Aggravated Crime. *The British Journal of Criminology*, 60(1):242–242, 2019.
- [21] Yashar Mehdad and Joel Tetreault. Do characters abuse more than words? pages 299–303, 01 2016.
- [22] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space, 2013.
- [23] Tomas Mikolov, Edouard Grave, Piotr Bojanowski, Christian Puhresch, and Armand Joulin. Advances in pre-training distributed word representations. In *Proceedings of the International Conference on Language Resources and Evaluation (LREC 2018)*, 2018.
- [24] Rafael Müller, Simon Kornblith, and Geoffrey Hinton. When does label smoothing help?, 2019.
- [25] Dennis Njagi, Z. Zuping, Damien Hanyurwimfura, and Jun Long. A lexicon-based approach for hate speech detection. *International Journal of Multimedia and Ubiquitous Engineering*, 10:215–230, 04 2015.
- [26] Chikashi Nobata, Joel Tetreault, Achint Thomas, Yashar Mehdad, and Yi Chang. Abusive language detection in online user content. In *Proceedings of the 25th International Conference on World Wide Web, WWW ’16*, page 145–153, Republic and Canton of Geneva, CHE, 2016. International World Wide Web Conferences Steering Committee.
- [27] Jeffrey Pennington, Richard Socher, and Christopher Manning. GloVe: Global vectors for word representation. In *Proc. of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543, Doha, Qatar, October 2014. Association for Computational Linguistics.
- [28] Timothy C Shiell. *Campus hate speech on trial*. Univ Pr of Kansas, 2009.
- [29] V. X. Son, V. Thanh, T. M. Vu, L. C. Thanh, and N. T. M. Huyen. Hsd shared task in vslp campaign 2019: Hate speech detection for social good. *Proceedings of VLSP 2019*, 2019.
- [30] Sara Owsley Sood, Elizabeth F. Churchill, and Judd Antin. Automatic identification of personal insults on social news sites. *J. Am. Soc. Inf. Sci. Technol.*, 63(2):270–285, February 2012.
- [31] Cynthia Van Hee, Els Lefever, Ben Verhoeven, Julie Mennes, Bart Desmet, Guy Pauw, Walter Daelemans, and Veronique Hoste. Detection and fine-grained classification of cyberbullying events. 09 2015.
- [32] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin. Attention is all you need, 2017.
- [33] Xuan-Son Vu, Thanh Vu, Mai-Vu Tran, Thanh Le-Cong, and Huyen T M. Nguyen. Hsd shared task in vslp campaign 2019:hate speech detection for social good, 2020.
- [34] Zeerak Waseem and Dirk Hovy. Hateful symbols or hateful people? predictive features for hate speech detection on twitter. In *Proceedings of the NAACL Student Research Workshop*, pages 88–93, San Diego, California, June 2016. Association for Computational Linguistics.
- [35] Guang Xiang, Bin Fan, Ling Wang, Jason Hong, and Carolyn Rose. Detecting offensive tweets via topical feature discovery over a large scale twitter corpus. pages 1980–1984, 10 2012.
- [36] J. M. Xu, K. S. Jun, X. Zhu, and A. Bellmore. Learning from bullying traces in social media. In *Proceedings of the 2012 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL HLT ’12*, page 656–666, USA, 2012. Association for Computational Linguistics.
- [37] Shuhan Yuan, Xintao Wu, and Yang Xiang. A two phase deep learning model for identifying discrimination from tweets. In *EDBT*, 2016.
- [38] Ziqi Zhang, D. Robinson, and Jonathan Tepper. Detecting hate speech on twitter using a convolution-gru based deep neural network. 03 2018.

⁵<https://www.aivivn.com/>