



# A white hat approach to fighting online trolls

Experiments with BERT and GAN



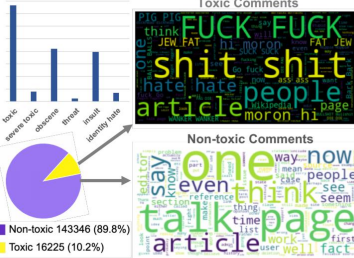
## INTRODUCTION

- In cybersecurity, "white hat" refers to an ethical expert who tests an organization's security defences with the aim of strengthening them
- The very same deep learning tools that are used to improve human conditions are also available to malicious actors for automation of attacks
- Taking a leaf from cybersecurity, we first build classifiers to detect abusive online language and then use those toxic comments to generate new ones to test the robustness of our classifiers against machine generated negative language

## DATASET

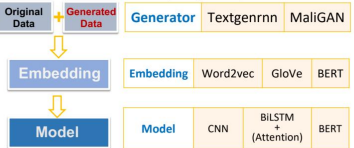
- The original dataset consists of 159,751 Wikipedia comments which have been labeled by human raters for six different types of toxicity

Distribution of Dataset

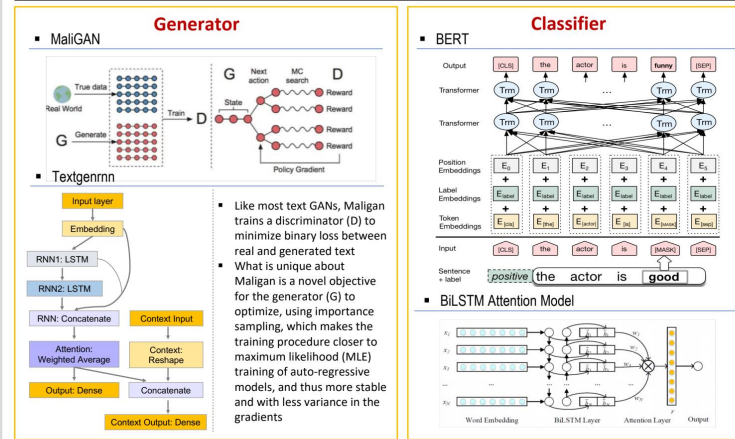


## METHODS

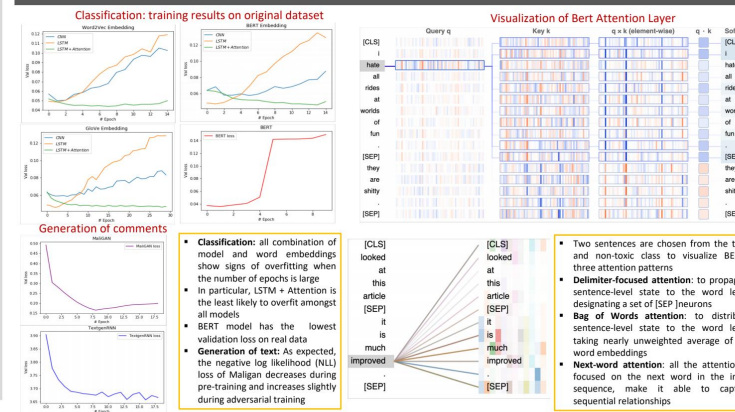
- Classification:** permutations of 4 models (BERT, CNN, LSTM & Attention) and 3 word embeddings (BERT, GloVe and word2vec)
- Generation of comments:** we used Maximum-Likelihood Augmented Discrete Generative Adversarial Networks (MaliGAN) and a multi-layer RNN consisting of RNN, LSTM and GRU (textgenrnn)



## MODEL ARCHITECTURE



## ANALYSIS

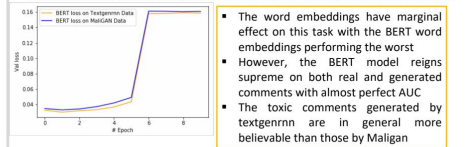


## RESULTS

	Original Data			Original + MaliGAN			Original + Textgenrnn		
	Train loss	Val loss	Total AUC	Train loss	Val loss	Total AUC	Train loss	Test loss	Total AUC
W2V+CNN	0.012	0.050	0.953	0.041	0.057	0.976	0.990	0.031	0.045
W2V+LSTM	0.013	0.048	0.940	0.043	0.052	0.978	0.994	0.038	0.044
W2V+Attention	0.033	0.048	0.983	0.036	0.039	0.986	0.998	0.033	0.040
GloVe+CNN	0.029	0.059	0.954	0.050	0.049	0.975	0.987	0.052	0.051
GloVe+LSTM	0.019	0.045	0.953	0.039	0.042	0.981	0.991	0.037	0.045
GloVe+Attention	0.044	0.043	0.983	0.032	0.036	0.989	0.995	0.044	0.040
BERT+CNN	0.018	0.058	0.936	0.050	0.062	0.967	0.988	0.045	0.048
BERT+LSTM	0.004	0.047	0.933	0.379	0.505	0.577	0.986	0.033	0.043
BERT+Attention	0.029	0.041	0.985	0.035	0.040	0.986	0.994	0.047	0.045
BERT+BERT	0.027	0.035	0.994	0.027	0.033	0.995	0.999	0.023	0.032

Examples of generated text

- MaliGAN 1. Sorry do you think it makes me really stupid shit?
- 2. You utter bitch get my real world abusing valid messages.
- Textgenrnn: 1. You are a dishonest idiot. Stop censoring this cte.
- 2. Oh, I am going to vandalize Wikipedia, what the fuck.



## DISCUSSION

- The prowess of the Attention is obvious - both BERT (Bidirectional Encoder Representations from Transformers) and BiLSTM Attention perform superior to LSTM and CNN across different word embeddings
- It is not surprising that the BERT model performed the best - the model is now considered the state-of-the-art in NLP as evident from its results on SQuAD v1.1
- Our results imply that even with access to modern text generation models such as MaliGAN and Textgenrnn, it will be difficult for motivated malicious actors to trick abusive language classifiers

## REFERENCES

- [1] BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding (2018)
- [2] Maximum-Likelihood Augmented Discrete Generative Adversarial Networks (2017)
- [3] Attention-Based Bidirectional Long Short-Term Memory Networks for Relation Classification (2016)

Jiaying Huang hij1227@stanford.edu

Zhihao Lin zhl@stanford.edu