

# CIS 522 – Final Project – Technical Report

Ground Zero

April 2022

## Team Members:

- Arvind Balaji Narayan; narvind; Email: [narvind@seas.upenn.edu](mailto:narvind@seas.upenn.edu)
- Gopik Anand; ganand; Email: [ganand@seas.upenn.edu](mailto:ganand@seas.upenn.edu)
- Vasanth Kolli; vasanthk; Email: [vasanthk@seas.upenn.edu](mailto:vasanthk@seas.upenn.edu)

---

<b>Code Links:</b>	GitHub	Google Drive
--------------------	--------	--------------

---

## Abstract

In the age of the internet, information is freely and easily accessible. However, it comes with its own share of benefits and risks. One major bane is misinformation. Misinformation is prevalent on every social media platform like Facebook, TikTok, Snapchat, Reddit, Whatsapp, Instagram, Twitter, etc. These pose extreme risks right from affecting elections to posing grave health risks as we observed during the COVID-19 pandemic and in certain cases might also lead to fatalities. A major source of this is fake news. Most people are usually led to believe that news structured in a formal way should ideally be from a reliable news source, but it is often not. We aim to address this issue of identifying fake news so that such content can be flagged and removed before it causes any disturbance. The rapid advancements in NLP help us tackle this problem and identify news as either real or fake. In this project, we use the publicly available Kaggle News dataset that contains news data from reputed and trustable news sources as our training set and news scraped from social media platforms like Reddit and Twitter forms our test set. The data is trained on multiple models starting from simple Machine Learning models to the state of the art Transformer models. The trained models were then used for the classification task and the results obtained display high accuracy values and F1 scores in identifying fake news irrespective of the distribution of news data and where it was posted.

**Keywords:** *fake news detection, online fake news, online fake reviews, online social network security, text classification, NLP*

# 1 Introduction

With the advent of the smartphone revolution and the rise of social media platforms, the way in which people consume news has drastically changed. People have also got obsessed with the fact of staying up to date with news and gossip around the world. News articles from reputed news firms take hours before they are published online and newspapers come only the next day. This does not help people stay updated with the latest developments. This is when social media came to the rescue. With a global connection expanding across every nook and corner of the world, the news could spread like wildfire. The moment something happens, a Twitter tweet is posted on the incident and millions of people consume that news. Movie updates, sports updates, accidents, calamities, etc. reach from one end of the world to the other with ease and rapid pace. Such is the raw power that social media platforms manifest.

However, these social media platforms can also be misused to spread fake news as well. Whenever news comes about an incident, we end up sharing it with others without verifying its authenticity. Once fake news is out and shared across the world, it is very difficult to correct it and make all the people who saw it understand and realize that it was fake. When fake news or gossip about an individual or a product is spread across, it puts them through testing times because proving their innocence and the fact that the news is fake is a daunting task. In a real-life setting where people largely consume news from social media handles, there is a huge possibility that we might be consuming fake news without knowing it and spread this out to friends, relatives and others which might lead to confusion and even pave the way to serious repercussions later. Fake news can have devastating effects on sectors like business, politics, medicine, healthcare, security, etc. and can cause fights, misunderstandings, etc. One such example was the anti-vax movement that started on social media where some sections of people started spreading false information about COVID-19 vaccines and how they should not be taken. This resulted in many people doubting the healthcare industry before doctors and medical practitioners convinced people of the positive effects of taking the vaccine. Thus, it is quintessential that we have access to the right information.

In such cases, It would be great if people can have a system/setup to filter out the real news from the fake ones so that they can take the right and appropriate action (whatever it might be). They can also share the right news with others. This is where our project, Ground Zero plays its part. By being able to distinguish between real and fake news seamlessly, we can provide a setup where people access the right content always. Using the power of Natural Language Processing, our system identifies fake news from platforms like Twitter and Reddit with impeccable accuracy.

In this project, we use Natural Language Processing to address the task at hand. We train our data on multiple models and report accuracy results and

F1 scores. We test our models on three combinations of the test set: a) Only Reddit, b) Only Twitter and c) a combination of the two. We then perform a comparative study to identify which kind of models do well on the task and report the results in a table as well as through suitable visualizations. Section 2 discusses related work while Section 3 introduces the dataset used in detail. Section 4 discusses the implementation methodology, Section 5 tabulates the results and Section 6 summarizes and discusses the findings and ethical implications.

## 2 Related Work

There have been few research works and studies in this field of study to classify news as either fake or real from different social media platforms like Twitter, Reddit, Facebook, etc. The authors in [1] and [2] proposed a detection model that combines text analysis using n-gram features and terms frequency metrics followed by deploying Machine Learning models like SVM, KNN, Logistic Regression, etc. for the classification task. They train and test their models on three different publicly available data sets. They also propose a new dataset called LIAR that collects the title or headline from news posts as well.

The author in [3] studies the different methods by which fake news or reviews can be spread and how to identify such content. This helped us understand how fake news can be collected from social media platforms. The author used Naive Bayes classifiers, SVM and semantic analysis for the classification task. The authors in [4] discuss how features can be extracted which is crucial to identify fake news and how models can be constructed as well. They also discuss several open issues and provide future directions for fake news detection in social media.

We use the Kaggle Dataset that was published recently to ensure that we don't end up using very old news for our task. Our overall objective and approach are similar to the authors of [1] but we train our models on the state of the art Transformer models apart from simple ML models to achieve impeccable accuracies. The details on the performance metrics for each model has been indicated in Section 5: Results.

## 3 Dataset and Features

**Data Contribution:** We will be making a data contribution in terms of scraping data from social media handles like Twitter, Reddit, etc. This data is bound to be in a different format and fashion when compared with the data seen in the current dataset which has been scraped from news articles. So, we believe that this is a novel effort to get data from a different source and train NLP models

to get good accuracies. This is our primary goal given that many people are now drifting from traditional newspapers and articles to the news that comes instantly on social media. Also, the scraped data is the latest news from around the world whereas the training data is comparatively older.

**Training Set:** Our primary dataset consists of about forty-five thousand samples with a roughly equal distribution of fake and real news samples taken from multiple well-reputed news platforms like the New York Times, US News, etc. These samples were collected and made into a dataset for public use accessible through Kaggle by Clement Bisailon. We used this dataset as our training set on which we trained our simple machine learning models, CNN model, LSTM model and Transformer models.

**Test Set:** For our first test set, we scraped news data from social media platforms like Reddit and Twitter. We used the Python Reddit API Wrapper (PRAW) to scrape data from news channels on Reddit. We were able to collect about 5k news samples from subreddit posts under news categories. For Twitter, we scraped data using the Twitter Search Scraper called sntwitter which is one of the modules under snsrape, a widely used social networking service scraper that is based in Python. We managed to collect a total of 10k news samples with a fairly equal distribution of fake and real news samples from real sources like AP news, BBCWorld and fake/sarcastic news sources like TheOnion and TheBabylonBee.

In this way, our test set contains 5k news samples from Reddit and 10k news samples from Twitter which we made use of both individually and combined together for different cases of testing and predictions.

#### **Features:**

News Content Features: We tried to understand the list of attributes that represent news content found on social media platforms. This includes

1. The source of the news sample like NYT, US News, etc.
2. The Headline: The title text that is included with the sole purpose of catching the attention of the header and explaining the main content of the article. (this is what we will use)
3. Body/Text: The full content of the article
4. Extra content: Hyperlinks, images and videos along with the news post

Using these feature representations, we can extract discriminative characteristics of fake news. However, for our purpose, we are taking only the title/Headline of the news article into consideration because when it comes to social media platforms, news posted is generally only the headline with a link to the full article if required.

## 4 Methodology

### 4.1 Data Pre-processing

After collecting the data, we analyzed the data, all the irrelevant fields that are not of interest to us were dropped and only columns of interest were retained. We initially combined the title and text of the news in the train set together but later realized that news in social media handles is most likely in the form of headlines due to the word limit enforced in apps like Twitter. As a result, we only retained the title of the news and the label to which it belongs (1 for real and 0 for fake). Thus we performed the training only on the titles of the news in the training dataset.

While collecting data from Reddit, we collected only news posts with over 1000 upvotes and ignored those below 1000 upvotes, considering it as a factor of secondary authentication of the news. In the case of Twitter, we scraped data only from trusted sources for real news and fake news/spoof news handles for fake news. News tweets sometimes contain the links to the original article as part of the tweets, such links if any were removed from the text of the tweets and tweets containing no text were discarded.

We manually scoured through some of the samples to verify the quality and authenticity of the scraped data.

After scraping data from Reddit and Twitter, they were formatted into CSV files for use during training and testing. Further preprocessing was done to remove duplicates, any NaN values and some of the posts which were later found to be deleted were also dropped accordingly. This is followed by using the Natural Language Processing Toolkit (nltk) to preprocess the data and remove stop words, truncate the text if needed, etc. The data was then tokenized before feeding it to the deep learning models for training or testing purposes.

### 4.2 Models Incorporated

A brief description of all the models incorporated has been documented below:

1. **Naive Bayes and SVM:** To begin with, we implemented statistical Machine Learning architectures such as SVM and Naive Bayes and tabulated their performance on our dataset. We reached the conclusion that even though SVM and Naive Bayes are comparatively simpler than other complex architectures, they did not do very well but could however be considered as good starting points to train further complex ensemble models.
2. **BiLSTM:** Bidirectional LSTM (BiLSTM) is widely used for sequence classification tasks and consists of two LSTMs, one that takes input in the forward direction and the other that takes input in the backward direction. They are high-end model architectures that take word embeddings as in-

put and use the pre-trained GloVe embeddings to produce the resultant embedding matrix.

3. **TextCNN:** TextCNN is the Convolution Neural Network that is specifically used for tasks concerning text classification. These require word embeddings to semantically map similar word structures used in the corpus. As a result, we use pre-trained Glove embeddings to output an embedding matrix that takes in the generated word embeddings as input.

#### **Transformer Models Used:**

Transformer models are the most prominent models used for NLP tasks. We finetuned the transformer models for the news classification task. For all the transformer model architectures that we trained and tested on, we have made use of the Hugging Face transformers library and understood how each model tokenizes the data. Following this, we managed to generate tokens for all our transformer-based models. Once the training and test set were obtained, we trained each model under similar conditions and observed the performance metrics.

1. **BERT:** BERT (Bidirectional Encoder Representations from Transformers) is a transformer-based machine learning technique that uses bidirectionality for language modeling applications. For the tokenization of our input for text classification, a single token which can represent the whole input is fed to the model. Therefore, an additional [CLS] token is added to the beginning of the sentence and a classification layer is added on top of the BERT output for the token and fine-tuned for text classification.
2. **CANINE:** This transformer model is directly trained at unicode character level and no explicit tokenization step is involved. Character-level training invariably leads to longer sequences, which CANINE avoids by using an effective downsampling technique before applying a deep Transformer encoder.
3. **Electra:** Electra follows a much better sample-efficient pre-training task called replaced token detection that largely reduces the amount of compute that is used in models like BERT for text classification tasks. It uses only about 25% of the compute needed for other models like XLNet, RoBERTa and DeBERTa at the cost of comparatively lesser performance metrics/
4. **GPT-2:** GPT-2 (Generative Pre-trained Transformer 2) follows a causal modeling objective that is used to predict the next token in the sequence and is pretrained on a large corpus of text data which allows for great performance capabilities on sequence classification tasks.
5. **XLNet:** XLNet is a transformer model that follows a generalized autoregressive pretraining approach that has the capability of learning bidirectional contexts by maximizing the expected likelihood and significantly outperforms BERT by large margins in text classification.

6. **ALBERT:** ALBERT (A Lite BERT) was developed by Google Researchers and involves changes in architecture from common transformer models like Factorized embedding parameters, Cross layer parameter sharing etc. This has allowed ALBERT to bring about improvements to BERT while reducing the number of parameters used by 30%.
7. **Electra:g)** RoBERTa: RoBERTa (Robustly Optimized BERT) follows a pretraining approach that largely builds on BERT by removing the next-sentence pretraining objective, modifying the key hyperparameters, and training with much larger learning rates and mini-batches. It manages to massively increase the performance on NLP tasks in comparison to BERT.
8. **DeBERTa:** DeBERTa (Decoding-enhanced BERT with disentangled attention) improves upon the performance of the BERT and RoBERTa models using 2 novel techniques. Firstly, a disentangled attention mechanism where each word's vector representation consists of 2 vectors that encode content and position. The attention weights are computed using disentangled matrices based on the content and relative positions. Secondly, the output softmax layer is replaced with an enhanced mask decoder to predict masked tokens for pretraining. DeBERTa only required half of the training data required by RoBERTa for similar performances.
9. **DistilBERT:** DistilBERT is a small, fast, cheap and light Transformer model trained by distilling the BERT model base. DistilBERT focuses on knowledge distillation during the pre-training phase which makes it possible to reduce the size of a BERT model by 40% while being 60% faster and retaining 97% of it's language understanding capabilities.
10. **XLM-RoBERTa:** This transformer model is largely based on and developed from RoBERTa and has been trained on extremely large data. The implementation of this model is very similar and offers comparable performance with that of RoBERTa. It offers greater attention to the language used in text and produces good classification results.
11. **LinkBERT:** This transformer model is pretrained on a large corpus of documents that captures document links such as hyperlinks and citations to develop knowledge that spans multiple documents. So, linked documents are fed into the same language model context. It views the text corpus as a graph of documents and solves two objectives namely masked language modeling and document relation prediction. It significantly outperforms BERT in the tasks of text classification and other knowledge intensive tasks like question answering and document retrieval.

Example of Transformer model finetuning. Below are the training loss, evaluation loss and AUROC plot. (link valid till 05/02/2022)  
<https://wandb.ai/anony-mouse-246221/visualization-demo/runs/2o3kb8ec?apiKey=614f0e90a04eb2588221afe71caeb2fd12aa3a32>



Figure 1: ELECTRA Model Training Loss

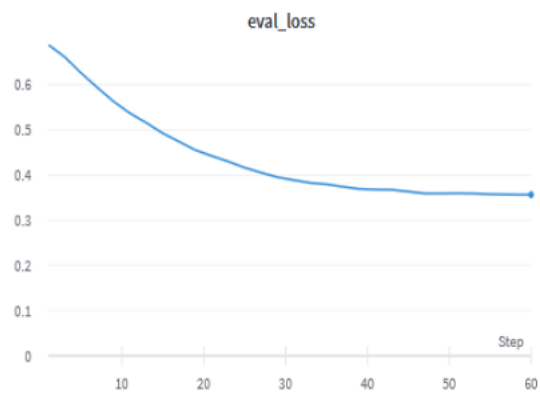


Figure 2: ELECTRA Model Evaluation Loss



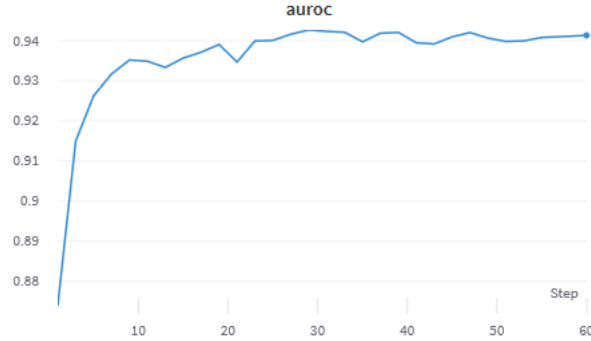


Figure 3: ELECTRA AUROC plot

## 5 Results

Each of the model architectures described above was trained on the Kaggle News Dataset and tested on three different combinations of the Test set. In the first case, we use only the Reddit dataset as the test set (with 5k samples) while we used only the Twitter dataset (with 10k samples) in the second case. In the third case, we combined both the Reddit and Twitter datasets to form the test set containing a total of 15k samples. The results obtained in each case have been tabulated below:

A) Train: Kaggle News Dataset — Test: Reddit News Dataset

Model Name	Validation Accuracy (%)	F1 Score
Naive Bayes	35.59	0.355
SVM	46.13	0.461
BiLSTM	44.12	0.441
TextCNN	44.69	0.276
BERT	82.21	0.822
Canine	79.42	0.885
GPT-2	78.29	0.878
AlBERT	51.35	0.678
RoBERTa	79.68	0.796
DeBERTa	80.19	0.890
DistilBERT	82.15	0.821
XLM-RoBERTa	80.25	0.802
Electra	82.56	0.825
LinkBERT	81.27	0.896

B) Train: Kaggle News Dataset — Test: Twitter News Dataset

Model Name	Validation Accuracy (%)	F1 Score
Naive Bayes	35.59	0.355
SVM	46.13	0.461
BiLSTM	44.12	0.441
TextCNN	44.69	0.276
BERT	82.21	0.822
Canine	79.42	0.885
GPT-2	78.29	0.878
AlBERT	51.35	0.678
RoBERTa	79.68	0.796
DeBERTa	80.19	0.890
DistilBERT	82.15	0.821
XLNet	81.68	0.816
XLM-RoBERTa	80.25	0.802
Electra	82.56	0.825
LinkBERT	81.27	0.896
XLNet	81.68	0.816

C) Train: Kaggle News Dataset — Test: (Reddit + Twitter) News Dataset

Model Name	Validation Accuracy (%)	F1 Score
Naive Bayes	35.59	0.355
SVM	46.13	0.461
BiLSTM	44.12	0.441
TextCNN	44.69	0.276
BERT	82.21	0.822
Canine	79.42	0.885
GPT-2	78.29	0.878
AlBERT	51.35	0.678
RoBERTa	79.68	0.796
DeBERTa	80.19	0.890
DistilBERT	82.15	0.821
XLNet	81.68	0.816
XLM-RoBERTa	80.25	0.802
Electra	82.56	0.825
LinkBERT	81.27	0.896
XLNet	81.68	0.816

We can observe that with the use of transformer models, we get high accuracies and F1 scores in the range of 85-90% and more, meaning that our models are performing very well in predicting test samples from Reddit or Twitter or a combination of both. All our project work was implemented in Google Colaboratory Pro IPython Notebook and all the necessary pre-installed packages used in ML/DL were acquired for smooth execution of our work.

### Accuracy and F1 score Results:

#### 1. Train: Kaggle — Test: Reddit

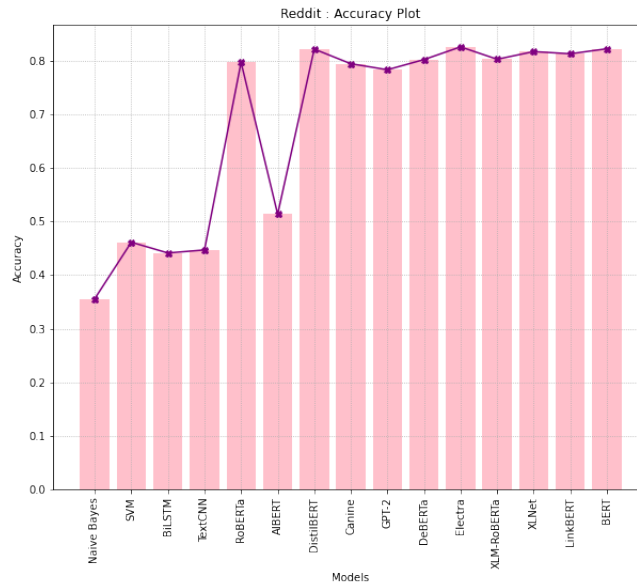


Figure 4: Train: Kaggle — Test: Reddit, Accuracy

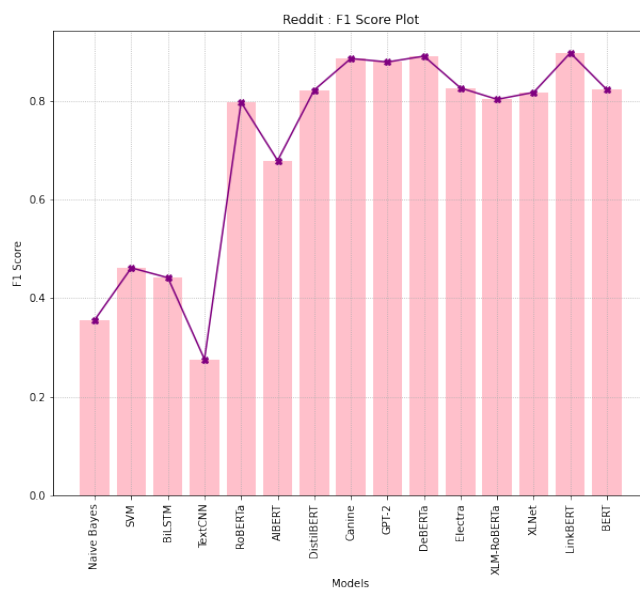


Figure 5: Train: Kaggle — Test: Reddit, F1 score

## 2. Train: Kaggle — Test: Twitter

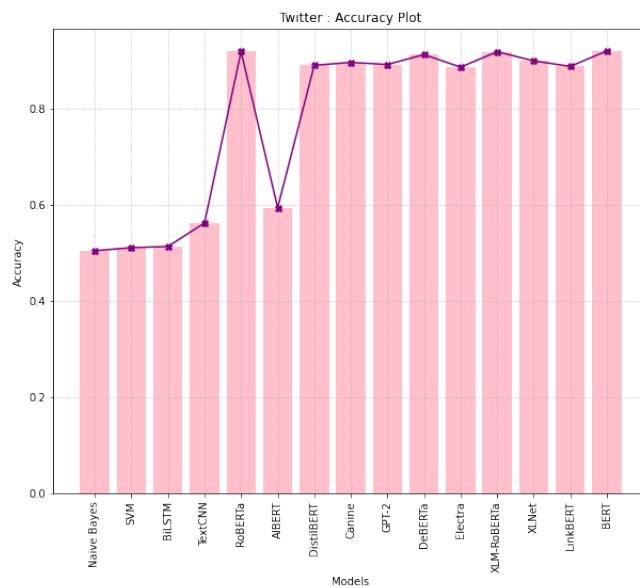


Figure 6: Train: Kaggle — Test: Twitter, Accuracy

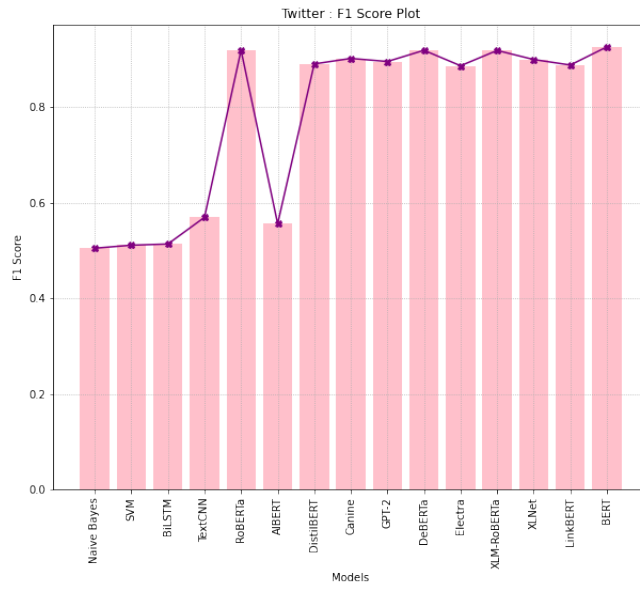


Figure 7: Train: Kaggle — Test: Twitter, F1 score

### 3. Train: Kaggle — Test: Twitter + Reddit

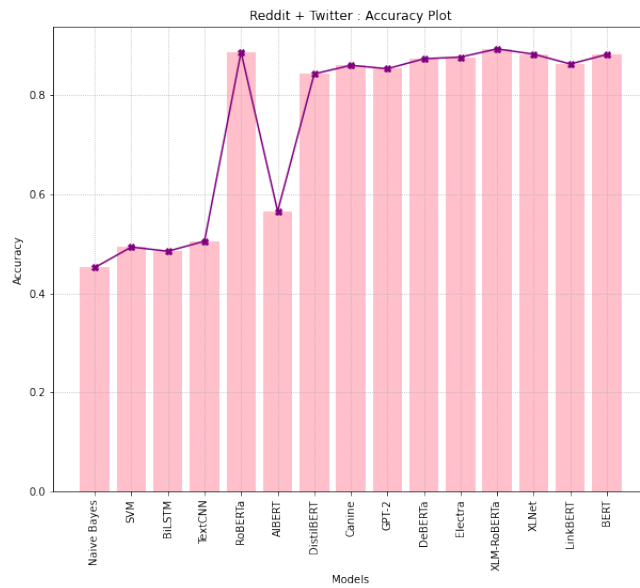


Figure 8: Train: Kaggle — Test: Reddit + Twitter, Accuracy

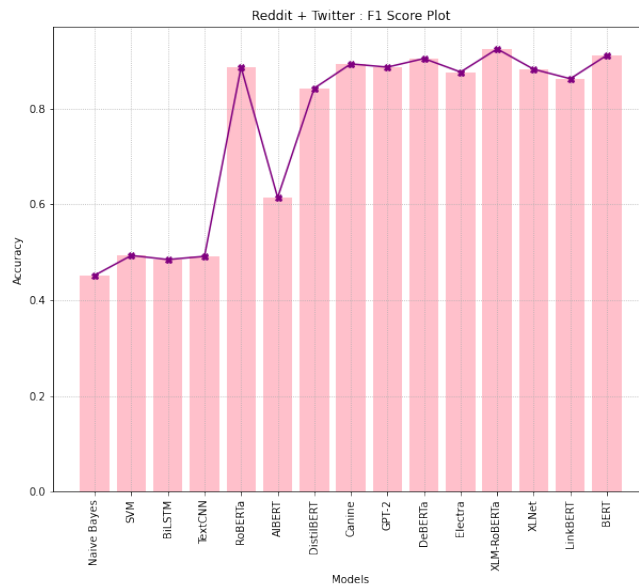


Figure 9: Train: Kaggle — Test: Reddit + Twitter, F1 score

### Word Clouds for Fake News and Real News:

1. Fake news



Figure 10: Fake News Word Cloud

## 2. Real news



Figure 11: Real News Word Cloud

## 6 Discussion

As mentioned in above sections, we implemented traditional machine learning models such as Naive Bayes and SVM as a baseline. TextCNN and BiLSTM were used as a baseline for deep learning models. This was followed by multiple state of the art transformer-based models all of which were used for analyzing their performance on the news classification task. Below we discuss the findings, limitations of our approach and also present future research directions for our news classification task.

## 6.1 Findings

We observed that the machine learning models perform well when the test set closely matches the train set i.e. when the test set is created by partitioning the available data into a train and test set. But when the test set is different from the train set as in our case (training on standard news articles data and testing on Reddit and Twitter posts) we observe that the models do not perform well giving us accuracies of about 50 % which is close to simple random guessing on a binary classification task. This shows that the models are unable to generalize adequately. This is because the posts from social media (test set) might be quite different from the (train set) and model overfits on the train data.

We observed similar results from the simpler deep learning models such as TextCNN and BiLSTM. This was quite surprising as these models were expected to perform better than the traditional ML methods mentioned above but were

unable to do so. This might also be due to the fact that the models learnt the train set distribution well but weren't able to generalize.

The transformers models however performed much better than the other models mentioned above and showed accuracy scores as high as 92%. This might be due to the attention and context-based embeddings being the core of the transformer models which is important in news classification. Also, the transformer models used were general-purpose transformer models which were pre-trained on a huge corpus and fine-tuned for the news classification task. These models generalized well both due to pre-training and attention-based mechanisms.

The motivation for the project was based on the social impact of fake news. Nowadays social media is an important aspect of society and fake news cannot be allowed to run rampant on these platforms. The repercussions of fake news are many and in worst cases lead to fatalities as observed in the COVID-19 pandemic. Hence being able to distinguish fake news from real news is critical. This is where recognizing fake news becomes a crucial aspect of modern society and our project aims to tackle this problem through deep learning. Our results indicate that this can be achieved and more so can be achieved with the training data and the test data being quite different and taken from various sources. This reduces the burden on the technical platforms detecting fake news drastically as obtaining labeled data for every kind of real and fake news on various platforms is tedious and next to impossible due to the sheer volume of news on social media. This methodology is very efficient and scalable considering the next best alternative which is human annotations.

## 6.2 Limitations and Ethical Considerations

The limitations of the approach are twofold - Technical and Ethical.

One of the technical limitations is due to limited compute resources on our end. We would have liked to train the transformer models on a much bigger dataset and for a larger number of epochs but it was not feasible due to the limited computing we possessed. Further, all the models mentioned above take text as input, derive features from the text and then classify them as fake or real. The models lack semantic understanding of the text i.e. if the news is actually real or fake based on the state of the world. The models rely solely on their understanding of the distribution of fake and real news based on past experience (training data). These models could be tricked if the fake news is generated to be similar to the distribution of the true news, this is being actively explored in the form of GANs for fake news detection. Here the objective is to create a better discriminator than the generator to classify fake news but can be misused to generate fake news too. Lastly, the models do not possess the semantic understanding of the news and the sequence of events mentioned in



news due to a lack of knowledge regarding the state of the world which the news is covering.

From an ethical standpoint, the limitations arise due to the bias and instances where the news classification into real and fake is blurred. The training of these models is supervised and hence required ground truth labels. The ground truth labels for the task are obtained manually and can be a major source of bias for the models. This poses a major risk where the system is indirectly modeling the bias of the human annotators. Thus labeling of data for such sensitive systems should be handled with extreme care and vigilance. The other factor is much more nuanced and difficult to factor into a system. Fake news detection for fact-based news is relatively straightforward in the sense that either a sequence of events occurred and are hence true or facts represented have sufficient tangible evidence. However, a lot of news revolves around presenting ideologies, sentiments and abstract notions/emotions which cannot be directly verified. E.g. news article stating “college education not useful anymore”, this represents the sentiment of the journalist towards modern education and cannot be directly verified as real or fake, it is a subjective opinion. This is a major problem as the system needs to identify news as fake or real even though the notion of real and fake is blurred here, classification of such news can pose a grave danger to the public as they can sway their opinion on matters subconsciously and can be exploited by bad actors.

### 6.3 Future Research Directions

As mentioned above, if provided with more compute resources and time, we would have liked to tackle some of the problems mentioned above. We would have liked to explore the idea of creating a GAN for fake news where the discriminator could be put into effect for the classification task. Furthermore, we would have explored the idea of incorporating graph neural networks to learn and encode the relations in the world from the training samples and act as a representation of the state of the world. Another avenue to be explored would have been the news classification on multimedia such as text and pictures/video which is quite common on social media using a multimodal approach to classify news.

## 7 Conclusions

The goal of this experiment was to check the feasibility of creating a machine learning model that can successfully discern between true and fake news. For that purpose, the models were trained on a bit dated (2 years old) newspaper articles and evaluated on the latest news posts curated from social media platforms namely Reddit and Twitter. The models were initially trained on traditional machine learning models like Naïve Bayes and SVM to establish a baseline for

the performance of the classification task. From there we proceeded to Deep Learning models like BiLSTM and TextCNN. Finally, we leveraged the best Deep Learning has to offer for NLP tasks i.e., Transformer Models. Throughout this journey, we observed various key features and made some observations as follows:

1. After successfully training the models, it was observed that the models fit the training data for the news classification task well despite the abstract nature of the task. Close inspection of the results of the process shows that deep learning models perform better in comparison to the traditional ML models as expected.
2. While considering the performance of the Deep Learning models, it was observed that the simpler models of the lot i.e., BiLSTM and TextCNN showed performance closer to baseline methods. This is attributed to the low expressive power of the models i.e., the complexity of the model which leads to its incapability to properly learn from the training data. These models also failed to generalize well on the test dataset. In comparison, most of the Transformer models exhibited significantly higher levels of performance.
3. Among the transformer models, ALBERT unexpectedly showed significantly less performance in comparison to the other transformer models. However, this led us to understand the importance of the number of parameters in the transformer model. For the purposes of our experiment and the lack of data, we adopted the process of Transfer Learning when it came to the Transformer models. Since we only fine-tuned the models for our classification task, the performance of the models is hugely attributed to the original training corpus and the number of trainable parameters of the model. ALBERT is nine times smaller than its original model i.e., BERT. This limitation with regards to the model architecture has led to the low expressive power of the model thereby leading to lower performance.
4. In the remaining Transformer models, it was observed that they exhibited significantly higher generalization performance compared to their competitors. They were able to classify most of the news data in the test dataset (Reddit and Twitter posts) correctly. This is attributed to the fact that they are pre-trained general language models which were fine-tuned for our task. Among them, some of the notable models were RoBERTa, XLM-RoBERTa, DeBERTa and BERT.
5. A stretch goal in our study was to understand the importance of the structure of the text when a model is trained to detect if a news article is real or fake. From our results, it can be observed that even with varying distributions with regard to the structure of the data, our models have generalized well.

## References

- [1] Ahmed H, Traore I, Saad S. *Detecting opinion spams and fake news using text classification*, Journal of Security and Privacy, Volume 1, Issue 1, Wiley, January/February 2018.
- [2] Ahmed H, Traore I, Saad S. (2017) *Detection of Online Fake News Using N-Gram Analysis and Machine Learning Techniques*, In: Traore I., Woungang I., Awad A. (eds) *Intelligent, Secure, and Dependable Systems in Distributed and Cloud Environments. ISDDC 2017. Lecture Notes in Computer Science*, vol 10618. Springer, Cham (pp. 127-138).
- [3] Kelly Stahl (2018) *Fake news detection in social media*, Published May, 2018 [https://www.csustan.edu/sites/default/files/groups/University%20Honors%20Program/Journals/02\\_stahl.pdf](https://www.csustan.edu/sites/default/files/groups/University%20Honors%20Program/Journals/02_stahl.pdf)
- [4] Shu, K., Sliva, A., Wang, S., Tang, J., Liu, H. (2017) *Fake News Detection on Social Media: A Data Mining Perspective*, ACM SIGKDD Explorations Newsletter, 19(1), 22-36.
- [5] <https://www.storybench.org/how-to-scrape-reddit-with-python/>
- [6] <https://towardsdatasciencecom/scraping-reddit-with-praw-76efcd1e1d9>
- [7] <https://mediumcom/@pasdan/how-to-scrap-reddit-using-pushshift-io-via-python-a3ebcc>
- [8] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei and Ilya Sutskever *Language Models are Unsupervised Multitask Learners*.
- [9] Kim Yoon, *Convolutional Neural Networks for Sentence Classification*
- [10] Train Dataset from Kaggle: <https://www.kaggle.com/datasets/clmentbisailon/fake-and-real-news-dataset>