# NLP-Fake News Detection

Apostol Alin, Ciuperceanu Vlad, Pîrvulescu Daria

May 2025

# Introduction

This paper aims to highlight the importance of recognizing synthetic online content through the use of **Natural Language Processing**.

In today's society, disinformation can range from images and videos altered using AI tools to entirely fabricated news articles and election posts. That being said, by leveraging **NLP techniques**, we can develop automated systems capable of identifying patterns of synthetic content.

# Related Work

### FakeBERT:

In the paper *'Fake news detection in social media with a bert-based deep learning approach. Multimedia tools and applications'*, Rohit Kumar Kaliyar, Anurag Goswami and Pratik Narang introduce an improved BERT, which was used to classify fake news from real news.

### The results where impressive:

 *'Our proposed model (FakeBERT) produced more accurate results as compared to existing benchmarks with an accuracy of **98.90** %.'*

# Dataset

The dataset used in our experiments was obtained from Kaggle and originates from the study: Welfake: Word embedding over linguistic features for fake news detection.



Figure: Dataset View

# Dataset-Analysis

Eliminating: empty strings, non-English texts, hyperlinks and emojis

```
Language distribution in original dataset: {'en': 70703, 'fr': 43, 'ar': 19, 'de': 110, 'es
': 145, 'ru': 156, 'sw': 14, 'pl': 4, 'so': 4, 'pt': 12, 'it': 6, 'tr': 7, 'nl': 5, 'fi': 3
, 'tl': 3, 'hr': 3, 'et': 1, 'vi': 2, 'ro': 4, 'no': 5, 'da': 1, 'af': 3, 'el': 2, 'zh-cn':
1, 'hu': 1, 'cy': 1, 'ca': 1, 'sq': 1, 'id': 2, 'sv': 1, 'lt': 1}
Empty texts removed: 783
Non-English texts removed: 561
Cleaned dataset saved to 'dataset_clean.csv'
```

Figure: Dataset Distribution

# Dataset-Statistic

As seen in the graphics, country names and president names are the most frequent and we expected it to affect the training of the model.
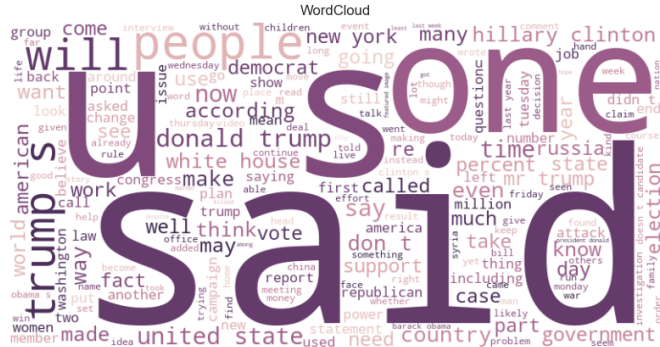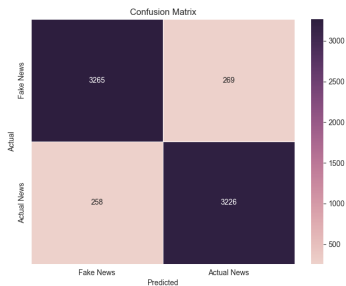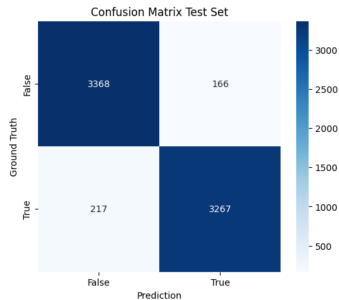


Figure: Word Cloud

# Word-Embeddings

We used : Word2Vec-CBOW, Word2Vec-SkipGram and GloVe (*both in their pretrained form and trained from scratch*).

# Classic-ML

As for comparing our two models the best combinations were: Random-Forest with Word2Vec-SkipGram trained from scratch with $92,49\%$ accuracy with Linear SVM Word2Vec-SkipGram trained from scratch with $94,54\%$ accuracy. It seems that **Skip-Gram** worked best for our large dataset.
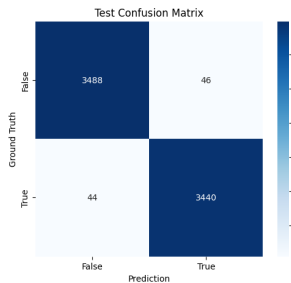


(a) Random Forest

(b) SVM

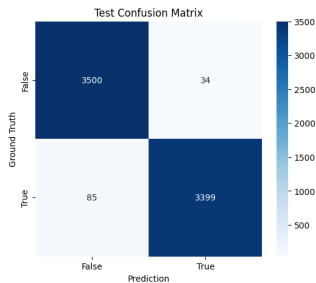Figure: Confusion matrices for the best models
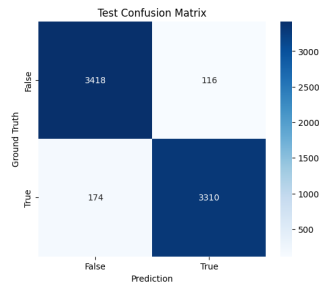
# Transformers

BERT with 98.72%
DistilBERT with (98.30%)
TinyBERT with (95.87%).



(a) BERT

(b) DistilBERT

(c) TinyBERT

Figure: Confusion matrices for the best models

# Fake part extraction

We aimed to extract the fake part of a fake news, or what out best performing models (DistilBERT and BERT) take most in consideration when making a prediction for real vs fake news.
For this we used two approaches: Lime and Captum Integrated Gradients.

# Lime

Lime trains a simple, interpretable model on the outputs of the original model for different small variations of the original text.
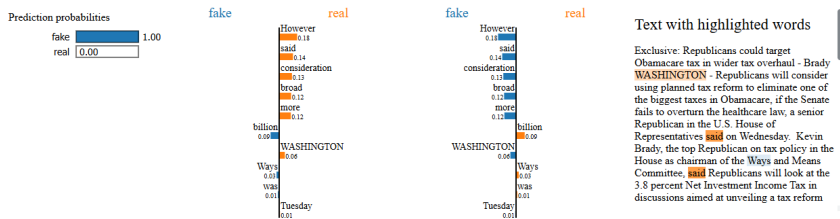


Figure: Word importance analysis using Lime

# Captum - Integrated Gradients

The Captum library provides the Integrated Gradients algorithm that calculates the integral of gradients with respect to inputs along the path from a given baseline to input.
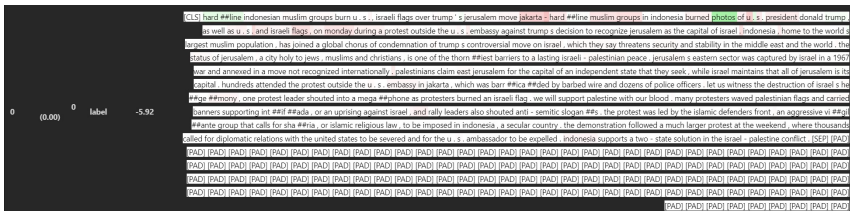


Figure: Word importance analysis using Integrated Gradients

# Conclusions

While our models achieved **promising results**, we acknowledge the **limitations** inherent to this task, including dataset bias, linguistic ambiguity, and the evolving nature of deceptive content.

Our analysis for the fake part extraction shows that extracting the fake part from a fake news still has to be deeply studied, as we've only shown certain buzzwords that influence the model into taking a decision.