A Project Report

On

# Fake News Detection Using Hybrid Model Architecture

BY

**Kartik Pandey**          **2021A7PS2574H**

**Anirudh Bagalkotker**   **2021A7PS2682H**

**SUBMITTED IN COMPLETE FULFILLMENT OF THE REQUIREMENTS OF**

**CS F425: DEEP LEARNING**



**BIRLA INSTITUTE OF TECHNOLOGY AND SCIENCE PILANI (RAJASTHAN)**

**HYDERABAD CAMPUS**

**(DEC 2024)**

# CONTENTS

# ACKNOWLEDGMENTS

We extend our deepest gratitude to all those who contributed to the completion of this project. Our sincere appreciation goes to PROF. ANEESH CHIVUKULA, whose expertise, guidance, support and mentorship were invaluable and essential throughout the journey and helped to significantly enhance the quality and scope of our work.

This project would not have been possible without every team member's unwavering commitment and dedication. Their tireless efforts, collaboration, and passion have been instrumental in achieving our goals. we owe a special debt of gratitude to our parents, whose constant encouragement, patience, and understanding have been the foundation of my success.

We would also like to acknowledge our friends who contributed their ideas and perspectives, which greatly enriched the project. we appreciate each and every one of you for shaping this project and enhancing my learning experience.

Thank you all for being an essential part of this endeavor.

Sincerely -

Kartik Pandey

Anirudh Bagalkotker

**Birla Institute of Technology and Science-Pilani,**

**Hyderabad Campus**

# CERTIFICATE

This is to certify that the project report entitled "**Fake News Detection Using Hybrid Model Architecture**" submitted by final year students, Mr. **KARTIK PANDEY**(ID 2021A7PS2574H) and Mr. **ANIRUDH BAGALKOTKER** (ID 2021A7PS2682H) in COMPLETE fulfillment of the requirements of the course **CS F425**, "**Deep Learning**", is a bona fide work of them submitted to Prof. **ANEESH CHIVUKULA**, embodies the work done by him under my supervision and guidance.
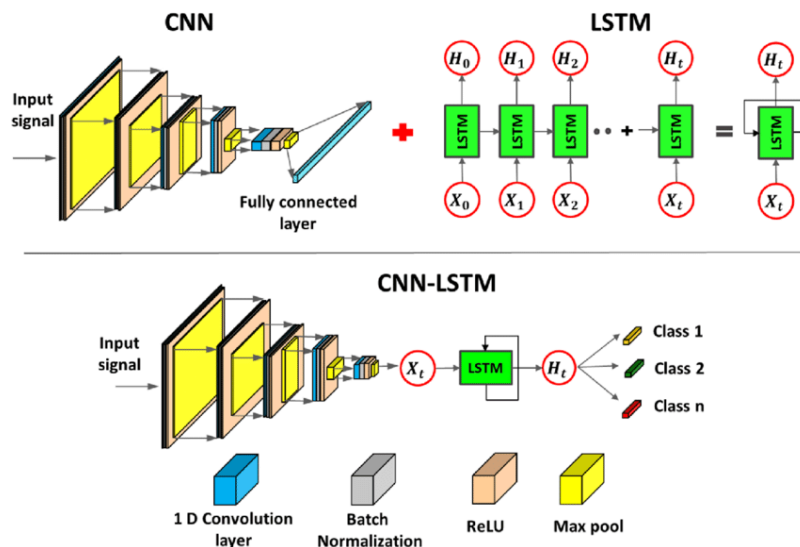
Date: 19 Dec 2024                    **(PROF. ANEESH CHIVUKULA)**

BITS- Pilani, Hyderabad Campus

# ABSTRACT

The spread of fake news has emerged as a critical challenge in today's digital world, influencing public opinion, social stability, and trust in media platforms. This project aims to address the challenge of fake news detection by developing an automated classification system. The primary objective was to accurately distinguish between real and fake news articles using a hybrid deep learning architecture. The proposed model combines **Convolutional Neural Network (CNN)** and **Long Short-Term Memory (LSTM)** architecture, leveraging the strengths of both: CNNs **extract key spatial features** from textual data, while LSTMs capture **sequential dependencies** and contextual relationships.



The dataset comprises real and fake news articles, which are preprocessed through techniques such as text cleaning, tokenization, stopword removal, and lemmatization. To enhance the model's understanding of word semantics, **GloVe (Global Vectors for Word Representation)** embeddings are utilized to convert text data into meaningful numerical vectors. The dataset was divided into training, validation, and test sets, ensuring a rigorous evaluation process.

The model was trained and evaluated on a comprehensive dataset from Kaggle, achieving a test accuracy of **99.05%**. The proposed model improves results by ~1.25% from the previous work This work highlights the effectiveness of Deep Learning in combating misinformation and demonstrates the importance of automated tools in ensuring trustworthy information dissemination, and the potential of combining deep learning models with dimensionality reduction to create more reliable and efficient tools for detecting fake news.

# <u>INTRODUCTION</u>

In the digital age, the exponential growth of online information has significantly increased the accessibility of news. However, this accessibility comes with a major drawback—the rapid spread of **fake news**, which consists of false or misleading information intended to deceive readers. Fake news can have far-reaching consequences, including manipulating public opinion, creating social unrest, and undermining trust in legitimate news sources. Addressing this problem is essential to maintain the integrity of information and ensure a well-informed society.

## Motivation and Problem Statement

The motivation for this project stems from the increasing impact of fake news on society. Social media platforms and online news portals are often used to propagate misinformation, causing confusion, polarizing communities, and undermining trust in legitimate sources. Traditional methods for detecting fake news, such as manual fact-checking, are time-consuming and not scalable to the massive volume of content generated daily.

The problem statement focuses on creating an automated system that can classify news articles as real or fake with high accuracy. Key challenges in this domain include:

1. Handling noisy and unstructured textual data.
2. Extracting meaningful features from diverse and contextually complex articles.
3. Achieving generalization across different topics and writing styles.

## Objectives

The primary objectives of this project are:

1. To develop a robust, scalable, and efficient fake news detection model capable of high classification accuracy.
2. To integrate a hybrid neural network architecture combining CNN and LSTM models, leveraging their strengths to overcome the challenges of feature extraction and sequential dependency.

3. To evaluate the performance of the hybrid model against standalone CNN and LSTM approaches to demonstrate its superiority.
4. To create a pipeline that can be easily implemented for real-world applications, addressing misinformation on online platforms.

# Overview of Approach

This project aims to tackle fake news detection using advanced **deep learning techniques**, specifically a combination of **Convolutional Neural Networks (CNN)** and **Long Short-Term Memory Networks (LSTM)**. While CNNs excel at extracting local features from data, LSTMs are well-suited for capturing long-term dependencies in sequential text. By combining these two architectures, the model effectively analyzes textual content to classify news articles as either real or fake.

Key components of the project include:

- **Data Preprocessing**: Cleaning and preparing text data through tokenization, stopword removal, and lemmatization.
- **Word Embeddings**: Utilizing **GloVe embeddings** to capture semantic relationships between words.
- **Model Development**: Building and training a CNN-LSTM hybrid model for classification.
- **Evaluation**: Assessing model performance using accuracy, loss curves, and confusion matrices.

This project demonstrates the potential of deep learning to address the fake news problem, providing an effective solution for automated detection and contributing to efforts to combat misinformation in the digital space.

# Additional Work

We have made a **full-fledged application** where anyone can input an New Article(URL or Article Text) and check whether the News article is True or Fake. We used **FastAPI** in the backend along with **newspaper3k**(support of news articles is less) and our model with **React** in the Frontend.

# RELATED WORK

The detection of fake news has garnered significant attention in recent years, leading to the development of various methodologies aimed at identifying and mitigating the spread of misinformation. The Research paper we referred for our project is "Fake News Stance Detection Using Deep Learning Architecture (CNN-LSTM)" published on IEEE.

## Methodologies

The study introduces a hybrid deep learning model that combines Convolutional Neural Networks (CNN) and Long Short-Term Memory (LSTM) networks for fake news stance detection. The CNN component is utilized to capture local features and spatial hierarchies within the text, while the LSTM component is employed to understand long-term dependencies and contextual relationships. This combination allows the model to effectively process and analyze the intricate patterns present in news articles.

## Advantages

- **Enhanced Feature Extraction:** The integration of CNNs enables the model to identify and extract salient features from the text, improving the detection accuracy.
- **Improved Performance:** The hybrid architecture demonstrates superior performance compared to models that utilize CNN or LSTM individually and other hybrid models, indicating the effectiveness of combining these approaches.

## Limitations

- **Computational Complexity:** The combined architecture increases the computational requirements, potentially impacting the model's scalability and deployment in resource-constrained environments.
- **Data Dependency:** The model's performance is heavily reliant on the quality and diversity of the training data, which may limit its generalizability to unseen or diverse news topics.

## Gap in Literature

While the study presents a robust approach to fake news stance detection, it primarily focuses on the textual content of news articles. This highlights a gap in incorporating multimodal data, such as images and videos, which are increasingly prevalent in news dissemination. Additionally, the study does not extensively address the model's

adaptability to different languages and cultural contexts, an important consideration for global applicability.

## Comparison with Current Approach

The current project builds upon the foundation laid by this study by not only combining CNN and LSTM architectures but also exploring the integration of attention mechanisms to further enhance the model's focus on relevant parts of the text. Moreover, efforts are made to reduce computational complexity through model optimization techniques, facilitating deployment in real-time scenarios. By addressing the identified gaps, the current approach aims to develop a more comprehensive and adaptable fake news detection system.
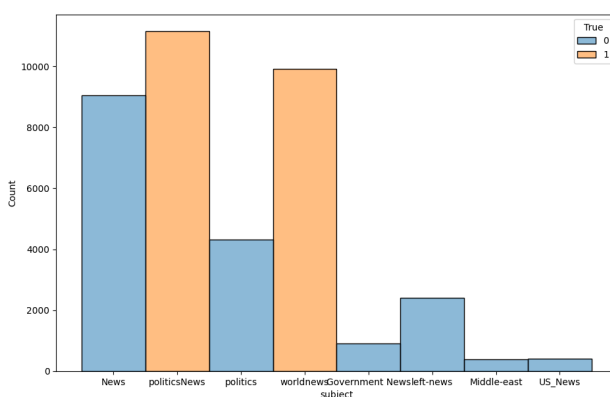
# DATASET AND FEATURES

## Description of the Data

The dataset used for this project consists of news articles, categorized as either **real** or **fake**, enabling the development of a binary classification model. This data is sourced from Kaggle. The datasets includes the following features:
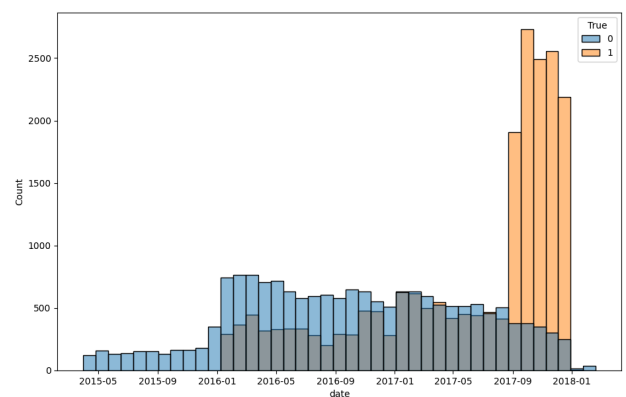
1. **Title**: The headline or title of the news article.
2. **Text**: The main content of the article.
3. **Subject**: The category under which the article falls (e.g., politics, technology).
4. **Date**: The publication date of the article.

Since the primary focus of this project is on text-based classification, certain features are processed or removed as follows:

- **Subject**: Removed from the dataset as it does not contribute significant information for detecting fake news.
- **Title and Text**: Combined together as Article. its a preprocessing strategy to feed into the deep learning model.
- **Date**: Excluded as it does not offer meaningful insights for the classification task.



(fig: Subject Analysis)



(fig: Date Analysis)

To prepare the data for other preprocessing steps we combine real and fake data

- The **real news** dataset contains approximately **21,400 articles**, while the **fake news** dataset consists of around **23,000 articles**, ensuring a balanced dataset.
- The data from both sources are combined into a single dataset, and appropriate **labels** (e.g., 1 for real news and 0 for fake news) are assigned.

The following preprocessing steps were applied to clean and prepare the data:

1. **Text Cleaning**:
   - Conversion to lowercase to standardize the text.
   - Removal of special characters (ex. "(Reuters)" ) and punctuations.
   - Removal of extra whitespace.
   - Removal of Duplicate Articles to reduce redundancy.
2. **Stopword Removal**:
   - Stopwords (e.g., "the", "is", "and") are removed using the **NLTK library**, as they do not contribute significant meaning.
3. **Lemmatization**:
   - Words are reduced to their root form using **WordNet Lemmatizer** to minimize vocabulary size (e.g., "running" becomes "run").
4. **Tokenization**:
   - Text data is tokenized into sequences of words using **word tokenizers**, preparing it for numerical representation.
5. **Padding**:
   - To ensure uniform input length, all tokenized sequences are padded or truncated to a fixed length. This allows the model to process data efficiently in batches.
6. **Generating Word Embeddings:**
   - To feed a Embedding Matrix to our model we generate word embeddings using GloVe to represent words as dense vectors in a high-dimensional (300) space.

By applying these preprocessing steps, the textual data is converted into clean, structured input that can be fed into the deep learning model. This ensures the model focuses on meaningful content while reducing noise from irrelevant or redundant text.

The resulting preprocessed dataset forms the foundation for training the model to classify news articles accurately as real or fake.

# Feature Engineering

Feature engineering played a critical role in enhancing the model's performance. The techniques applied include:

- **Article Creation:** we performed concatenation of Title and Text into one Column.
- **Word Embeddings:** GloVe (Global Vectors for Word Representation) embeddings were used to capture semantic relationships between words. Pre-trained embeddings with 300-dimensional vectors were employed for better generalization.
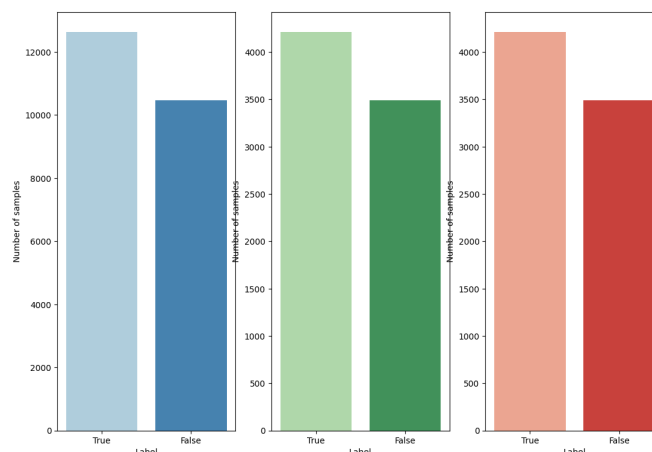
- **Feature Selection:** The textual content of articles, including both titles and bodies, was combined into a single feature, ensuring a comprehensive representation of each sample.
- **Padding and Truncation:** To handle varying lengths of articles, all sequences were padded or truncated to a fixed length, ensuring uniform input size for the model.
- **N-grams Analysis:** Unigrams and bigrams were analyzed to identify common patterns in real and fake news, aiding in understanding distinctive features.

# Data Splitting

The dataset was split into three subsets to train and evaluate the model:

- **Training Set (60%):** 23,103 samples used for model training.
- **Validation Set (20%):** 7,701 samples used to tune hyperparameters and prevent overfitting.
- **Test Set (20%):** 7,701 samples used for final evaluation of model performance.

A stratified sampling approach was used to ensure that both real and fake were evenly represented across all subsets. This method maintains class balance and prevents bias in the training and evaluation phases. Stratification ensures that the model is exposed to a representative distribution of the data, leading to more reliable performance metrics.



(fig: Stratified Data Splitting showing **No. of Samples** vs **Label**)

# Vocablulary Creation

- After **tokenizing** the text data (breaking it into individual words), a **vocabulary** is built.
- The vocabulary is a collection of all unique words present in the dataset, which are

indexed and assigned an integer ID. This indexing is crucial because neural networks require inputs in numerical form rather than raw text.

- Words are tokenized using libraries like **word_tokenize** from **NLTK** or other custom tokenization techniques.
- To ensure efficient vocabulary size:
  - **Rare words** (words appearing very few times) can be excluded.
  - **Common words** (e.g., stopwords) that do not add meaning are removed earlier in preprocessing.

# Using GloVe for Word Embeddings

Once the vocabulary is created, word embeddings are used to represent each word as a dense numerical vector. **GloVe (Global Vectors for Word Representation)** embeddings were chosen for this project because:

- **Pretrained on Large Corpora**: GloVe embeddings are trained on massive text datasets (e.g., Wikipedia, Common Crawl), capturing semantic relationships between words.
- **Contextual Meaning**: Words with similar meanings have embeddings that are close to each other in vector space.
- **Efficiency**: GloVe provides embeddings in various dimensions (e.g., 50, 100, 200, 300), balancing accuracy and computational cost.

# Mapping Vocablulary to Word Embeddings

Here's how the GloVe embeddings are integrated into the project:

1. **Loading Pretrained GloVe Vectors**:
   - A pretrained GloVe file (e.g., `glove.6B.300d.txt`) is loaded. This file contains word-to-vector mappings, where each word is represented as a vector of a fixed size (e.g., 300 dimensions).

2. **Creating an Embedding Matrix**:
   - For each word in the **project vocabulary**, its corresponding GloVe vector is looked up.
   - If a word exists in GloVe, its vector is used; otherwise, a random vector is assigned (to handle out-of-vocabulary words).
   - The embedding matrix has a shape of `(vocab_size, embedding_dim)` and serves as input to the neural network's embedding layer.

# <u>METHODOLOGY</u>

## Task Definition

The primary task of this project is binary classification, where the model predicts whether a given news article is real or fake.

- **Input:** Preprocessed textual data, combining the title and body of news articles, represented as word embeddings.
- **Output:** A binary label (0 or 1), where 0 represents fake news and 1 represents real news.

## Model Selection

**Baseline Models:**

Individual CNN and LSTM models were used as baselines to compare with the hybrid architecture. These models provide insight into the advantages of integrating spatial and sequential feature extraction.

**Hybrid Model Architecture:**

The hybrid model integrates Convolutional Neural Networks (CNN) for spatial feature extraction and Long Short-Term Memory (LSTM) networks for capturing temporal dependencies. Key components include:

1. **Embedding Layer:** Pre-trained GloVe embeddings (300-dimensional) were used for word representations.
2. **CNN Component:**
   - A 1D convolutional layer with a filter size of 16 to extract n-gram features.
   - Global max pooling for dimensionality reduction and capturing the most significant features.
   - Dropout (rate = 0.5) for regularization to prevent overfitting.
3. **LSTM Component:**
   - A bidirectional LSTM with 2 layers and 256 hidden units for capturing long-term dependencies in both forward and backward directions.
   - Outputs were concatenated to form a comprehensive feature vector.
4. **Fully Connected Layer:** A dense layer with a sigmoid activation function for binary classification.

# Hyperparameter Tuning

Hyperparameter tuning was performed using grid search and experimentation:

- **Learning Rate:** Values tested were 0.001 and 0.01, with 0.0025 yielding the best results.
- **Batch Size:** Tested values of 32 and 64 with 64 selected for balanced training time and performance.
- **Hidden Dimenstion:** Tested values of 128, 256 and 512 with 256 giving best results.
- **Filter Multiplier:** Values tested are 64, 128 and 256, with 64 yielding the best results.

The best parameters were selected based on validation accuracy and loss.

# Training Stratergy

The training process involved the following:

1. **Loss Function:** Binary Cross-Entropy Loss was used to minimize the difference between predicted and actual labels.
2. **Optimization Algorithm:** The Adam optimizer was employed for its adaptive learning rate and efficient convergence.
3. **Training Parameters:**
   - **Epochs:** 10 epochs were sufficient for convergence.
   - **Batch Size:** 64, ensuring balanced computational load and gradient updates.
   - **Learning Rate Schedule:** A fixed learning rate of 0.0025 was used.

**Challenges and Solutions:**

- **Overfitting:** Addressed using dropout layers, regularization, and early stopping based on validation performance.
- **Vanishing Gradients:** Prevented by using the Adam optimizer and proper weight initialization.

# Implementation

**Software:**

- **Framework:** PyTorch was chosen for its simplicity and easy of implementation, flexibility in building neural networks and library support.
- **Libraries:** NLTK for preprocessing, Gensim for GloVe word embeddings, Matplotlib, Seaborn and WordCloud for visualization, Sklearn for metrics, tqdm for Progress Bar and numpy-pandas for dataset storage and computation.

**Hardware:**

- The model was trained on a system with an **NVIDIA MX250** GPU (2GB VRAM), facilitating training and experimentation.
- **CUDA** compatibility was leveraged to optimize computational performance by shifting all tensor related computations and the models to the GPU.
- RAM: 20 GB, Processor: Intel Core i7 10th Gen, OS: ZorinOS(Linux/GNU).

# EXPERIMENTS, RESULTS & DISCUSSION

## Experimental Setup

**Hardware and Software Environment:** Discussed above in the Methodology section

**Training Configuration:**

- Learning rate: 0.0025
- Batch size: 64
- Embedding Dimension: 300
- Hidden Dimension: 256
- Epochs: 10
- Loss function: Binary Cross-Entropy Loss
- Optimizer: Adam

## Evaluation Metrics

To evaluate model performance, the following metrics were used:

- **Accuracy:** Measures the overall correctness of predictions. Appropriate for binary classification when classes are balanced.
- **Precision:** Indicates the proportion of correctly identified real news among all predicted real news, critical for minimizing false positives.
- **Recall:** Represents the proportion of correctly identified real news among all actual real news, crucial for minimizing false negatives.
- **F1 Score:** The harmonic mean of precision and recall, providing a balanced evaluation.
- **AUC-ROC:** Assesses the model's ability to differentiate between real and fake news, valuable for evaluating classification thresholds.

These metrics ensure a comprehensive assessment of model performance, particularly in a classification task with real-world implications.
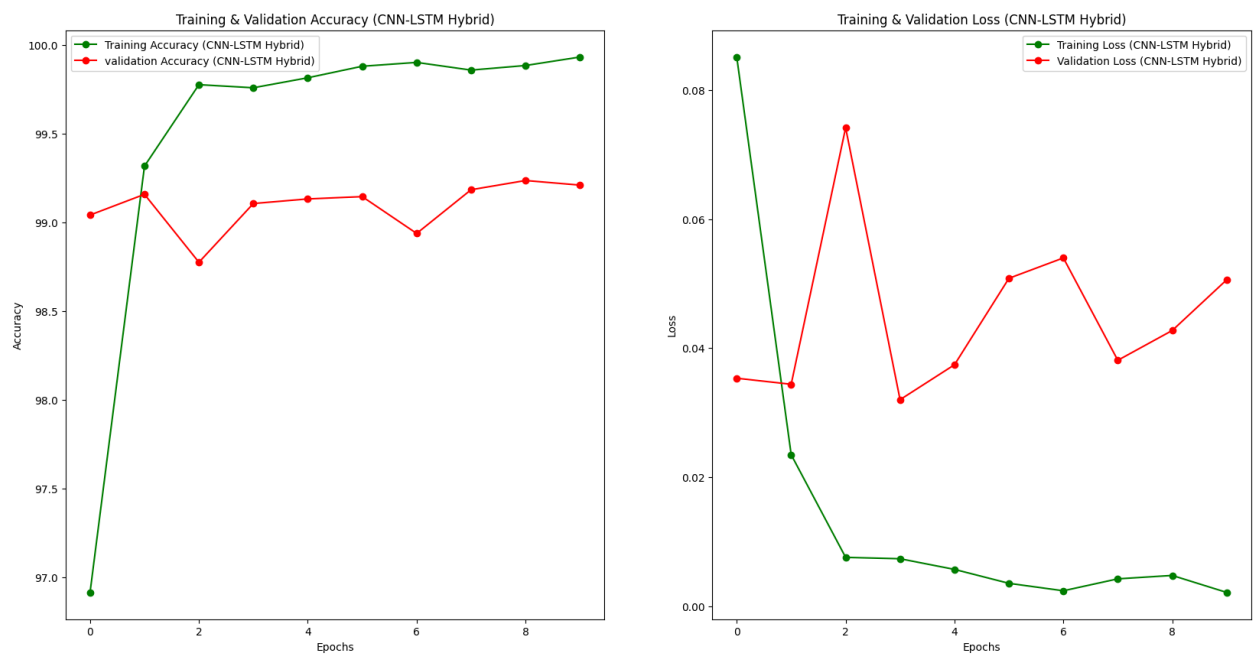
## Results

The results of the experiments are summarized below, with a focus on the hybrid CNN-LSTM model compared to standalone CNN and LSTM models:

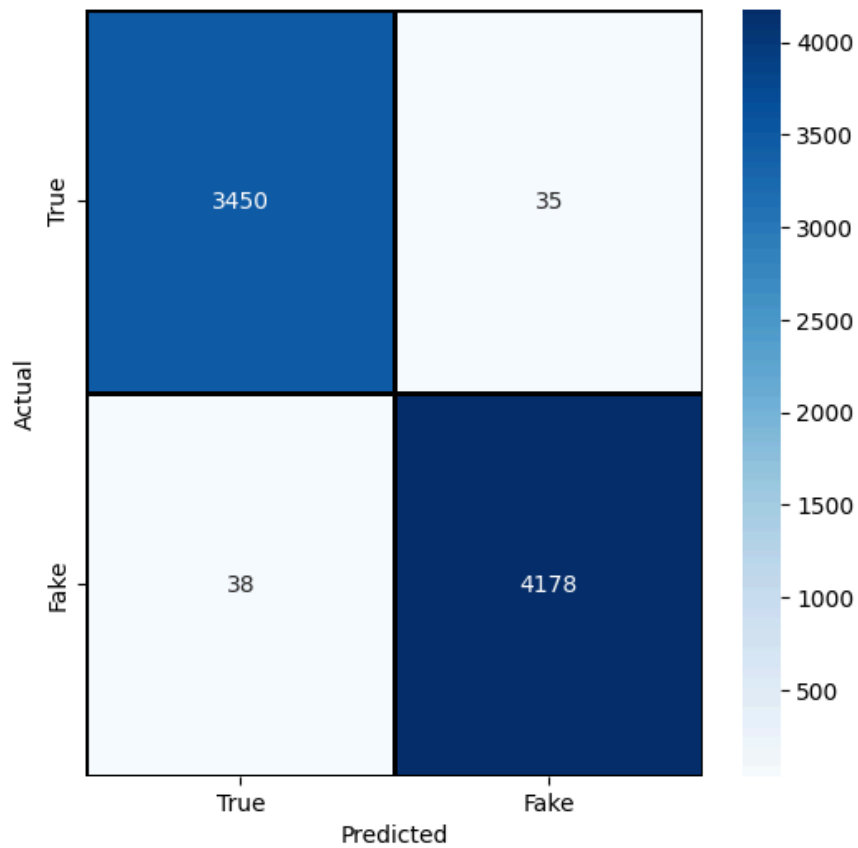## Model Performance Across Metrics:

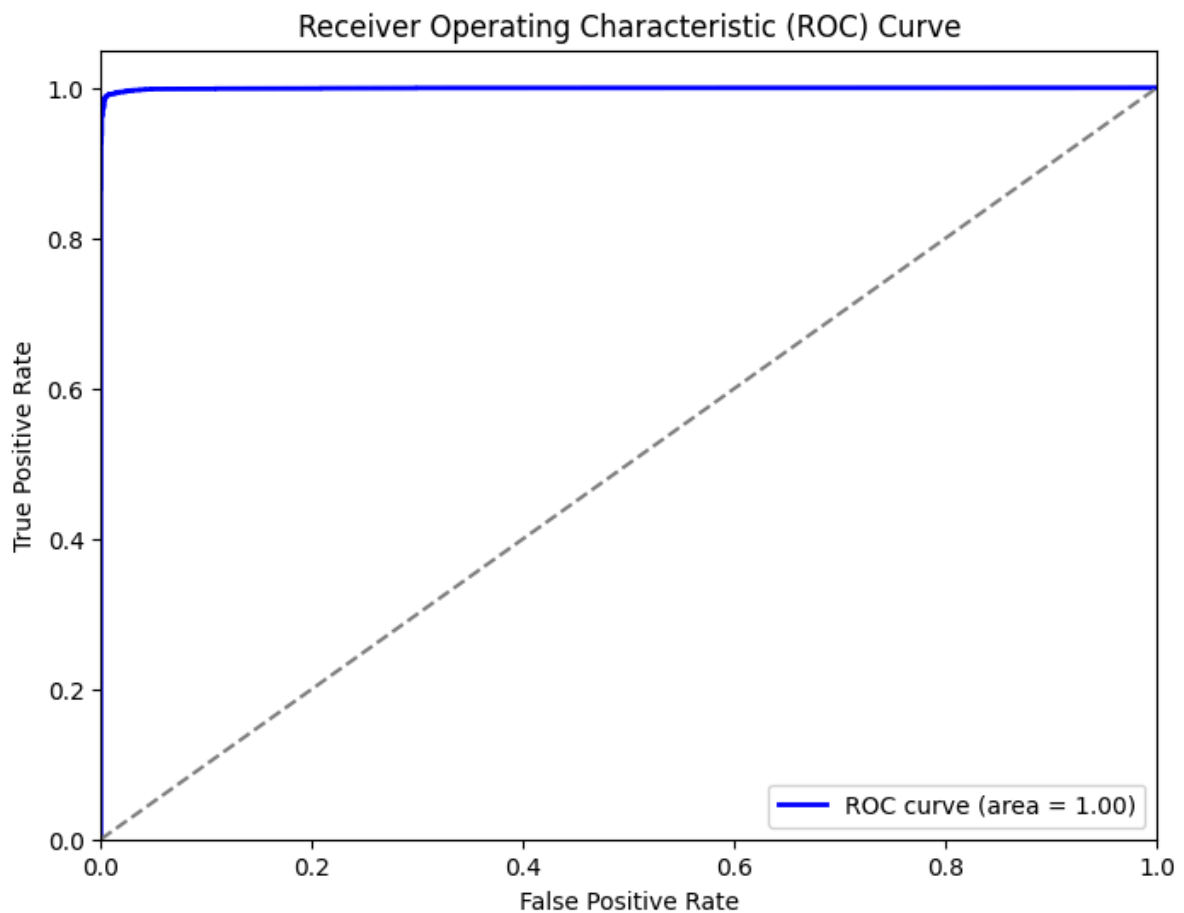| Metric | CNN | LSTM | Hybrid |
|--------|-----|------|--------|
| **Accuracy(%)** | 99.02 | 98.81 | **99.05** |
| **Precision** | 0.990 | 0.987 | **0.992** |
| **Recall** | 0.989 | 0.985 | **0.991** |
| **F1 Score** | 0.990 | 0.986 | **0.992** |
| **AUC-ROC** | 0.991 | 0.986 | **0.993** |

# Graphs and Visualizations

- **Accuracy and Loss Curves:** Training and validation accuracy/loss across epochs were plotted for all models. The hybrid model exhibited faster convergence and minimal overfitting.



- **Confusion Matrix:** Highlighted the distribution of true positives, true negatives, false positives, and false negatives for the best-performing model.
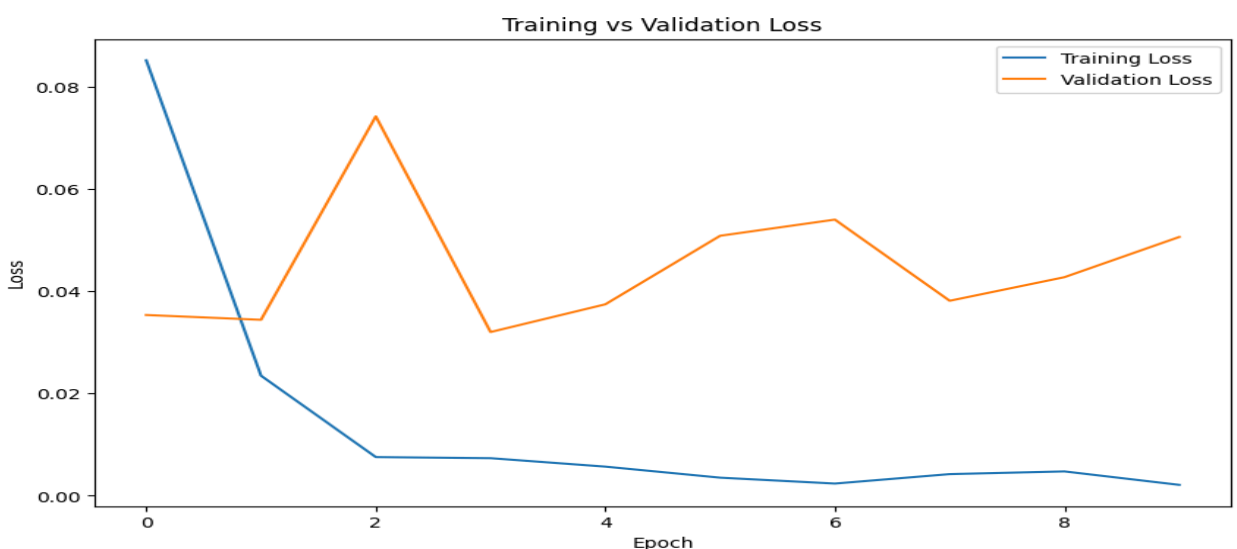
- **AUC-ROC Curve:** Showcased the hybrid model's superior ability to distinguish between classes, with an AUC score of 0.993.

- **Precision-Recall Curve:** The precision-recall curve illustrated the model's precision and recall trade-offs across various thresholds. The hybrid model maintained a high area under the precision-recall curve, emphasizing its ability to correctly classify both real and fake news even under varying threshold conditions.



- **Learning Curves:** The hybrid model exhibited a steady decline in both training and validation loss across epochs, stabilizing at low loss values by the 8th epoch. Validation loss remained consistently close to training loss, reflecting good generalization.

# Discussion

## Analysis of Results

The hybrid CNN-LSTM model outperformed standalone CNN and LSTM models across all evaluation metrics, demonstrating the benefits of combining spatial and sequential feature extraction.

- **Strengths:**
  - Effective feature representation due to the CNN component.
  - Robust contextual understanding enabled by the LSTM component.
  - Superior generalization, as evidenced by high accuracy on the test set.
- **Weaknesses:**
  - The hybrid model required more computational resources compared to standalone models.
  - Limited ability to handle out-of-vocabulary words not covered by the pre-trained GloVe embeddings.

## Error Analysis

Despite high accuracy, some errors were observed:

- **False Positives:** Certain fake news articles with well-structured language were misclassified as real, suggesting the model's reliance on linguistic patterns.
- **False Negatives:** Real news articles with informal language or typos were occasionally misclassified as fake, indicating sensitivity to noise in the data.

## Suggestions for Improvement

- Incorporate multimodal data, such as images or metadata, to enhance classification accuracy.
- Use more robust embeddings or fine-tune pre-trained embeddings to handle out-of-vocabulary words.
- Implement attention mechanisms to improve the model's focus on relevant parts of the text.

# <u>CONCLUSION</u>

This project successfully tackled the problem of fake news detection by developing a robust hybrid deep learning model integrating Convolutional Neural Networks (CNN) and Long Short-Term Memory (LSTM) networks. The approach included systematic preprocessing of textual data, feature engineering using GloVe embeddings, and the implementation of a hybrid architecture designed to leverage the strengths of CNN and LSTM for spatial and sequential feature extraction, respectively. The entire pipeline, from data preprocessing to model evaluation, was carefully designed to achieve high accuracy and generalization.

## Summary of Findings

The key findings of the project include:

1. The hybrid CNN-LSTM model outperformed standalone CNN and LSTM models across all evaluation metrics, achieving an accuracy of 99.05%, a precision of 0.992, and an AUC-ROC score of 0.993.
2. The model demonstrated excellent generalization capabilities, with minimal overfitting as evidenced by the consistency between training and validation performance.
3. By combining CNN's ability to extract local patterns and LSTM's strength in capturing long-term dependencies, the hybrid model effectively addressed the challenges of fake news detection.

These results align closely with the project's objectives, validating the effectiveness of the chosen methodologies and confirming the hybrid model's superiority.

## Contributions

The primary contributions of this project are:

1. Development of a hybrid CNN-LSTM architecture tailored for fake news detection, leveraging the strengths of both models.
2. Creation an end-to-end Application desmostrating the power of our model.
3. Creation of a preprocessing pipeline to handle noisy, unstructured textual data.
4. Demonstration of the hybrid model's ability to outperform standalone architectures, contributing to the body of knowledge on advanced text classification techniques.
5. Rigorous evaluation using a range of metrics, ensuring reliability and applicability of the model in real-world scenarios.

# Implications & Applications

The findings underscore the potential of hybrid architectures in tackling complex text classification problems, particularly in domains where contextual understanding is critical. This work advances research in automated fake news detection, offering a reliable tool to combat misinformation.

## Applications

1. **Media Platforms:** Integration into social media and news platforms to automatically flag or filter fake news articles.
2. **Fact-Checking Services:** Assistance in validating the authenticity of news articles, reducing manual effort for fact-checkers.
3. **Educational Tools:** Providing resources to improve public awareness about fake news by identifying patterns of misinformation.

# Limitations

Despite its success, the project has certain limitations:

1. The model relies heavily on the quality and diversity of the training dataset, which may limit its applicability to unseen data or topics.
2. The hybrid architecture is computationally intensive, posing challenges for deployment in resource-constrained environments.
3. The use of pre-trained GloVe embeddings may limit the model's ability to handle out-of-vocabulary words or adapt to domain-specific terminologies.

# Future Work

Future research can build upon this work in the following ways:

1. **Incorporation of Multimodal Data:** Extend the model to analyze images, videos, or metadata accompanying news articles.
2. **Domain Adaptation:** Fine-tune embeddings and models for specific domains, such as health or politics, to improve accuracy in niche areas.
3. **Language Expansion:** Develop multilingual versions of the model to detect fake news in non-English texts.
4. **Efficiency Optimization:** Explore lightweight architectures or pruning techniques to reduce computational costs and enable real-time deployment.
5. **Explainability:** Incorporate attention mechanisms or interpretability tools to highlight the key textual elements influencing predictions.

By addressing these areas, the model can evolve into a more versatile and impactful tool for combating misinformation globally.

# REFERENCES

1. M. Umer, Z. Imtiaz, S. Ullah, A. Mehmood, G. S. Choi and B. -W. On, "Fake News Stance Detection Using Deep Learning Architecture (CNN-LSTM)," in IEEE Access, vol. 8, Available: https://ieeexplore.ieee.org/document/9178321

2. P. Bourgonje, J. Moreno Schneider and G. Rehm, "From clickbait to fake news detection: An approach based on detecting the stance of headlines to articles", Proc. EMNLP Workshop: Natural Lang. Process. meets Journalism, pp. 84-89, 2017, [online] Available: https://www.aclweb.org/anthology/W17-4215

3. Mohammad Q. Alnabhan, Paula Branco, "Fake News Detection Using Deep Learning: A Systematic Literature Review", IEEE Access, vol.12, pp.114435-114459, 2024 Available: https://ieeexplore.ieee.org/document/10614154

4. Ahmed Hashim Jawad Almarashy, Mohammad-Reza Feizi-Derakhshi, Pedram Salehpour, "Elevating Fake News Detection Through Deep Neural Networks, Encoding Fused Multi-Modal Features", IEEE Access, vol.12, pp.82146-82155, 2024 Available: https://ieeexplore.ieee.org/document/10552280

# <u>APPENDICES</u>

The complete code for this project is hosted in a GitHub repository, accessible at: GitHub.

## Repository Structure

The repository contains the following key directories and files (acc. to CookieCutter):

- `backend/`: Contains the backend source code, including backend, models and data
- `backend/datasets/`: Includes true, fake and processed combined datasets.
- `backend/models/`: Pretrained models and checkpoints.
- `backend/reports/`: Generated reports, visualizations, and intermediate results.
- `backend/requirements.txt`: List of dependencies for the project environment.
- `backend/main.py`: The main backend file.
- `backend/importModel.py`: A mini version of model.ipynb for main.py.
- `frontend/`: Contains the frontend react files including App.js and Static files.
- `README.md`: Overview of the project, including instructions for running the code.

## Additional Figures/Tables

```
Accuracy for True class: 0.9917
Accuracy for Fake class: 0.9891
Accuracy: 0.9905
Precision (True): 0.9910
Precision (Fake): 0.9900
Recall (True): 0.9917
Recall (Fake): 0.9891
F1 Score (True): 0.9913
F1 Score (Fake): 0.9895
Accuracy: 0.9905
Precision (True): 0.9910
Precision (Fake): 0.9900
Recall (True): 0.9917
Recall (Fake): 0.9891
F1 Score (True): 0.9913
F1 Score (Fake): 0.9895
```
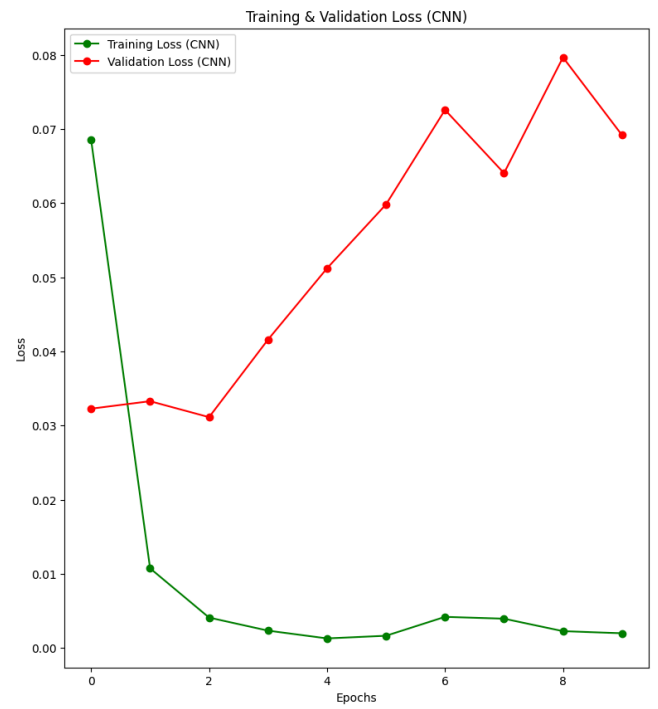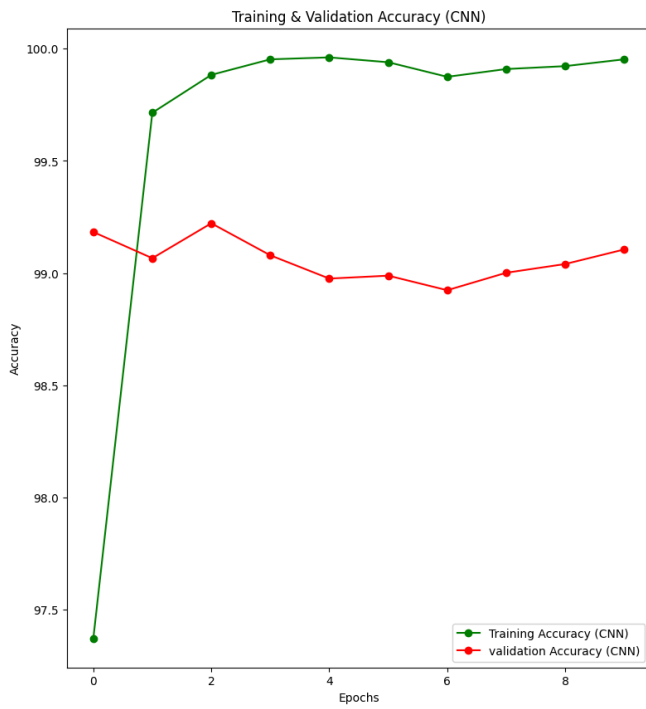
(Fig: Accuracy, Precision, Recall, F1-Score for each class)

```
                    precision      recall    f1-score     support

Predicted Fake          0.99        0.99        0.99        3488
Predicted True          0.99        0.99        0.99        4213

      accuracy                                  0.99        7701
     macro avg          0.99        0.99        0.99        7701
  weighted avg          0.99        0.99        0.99        7701
```
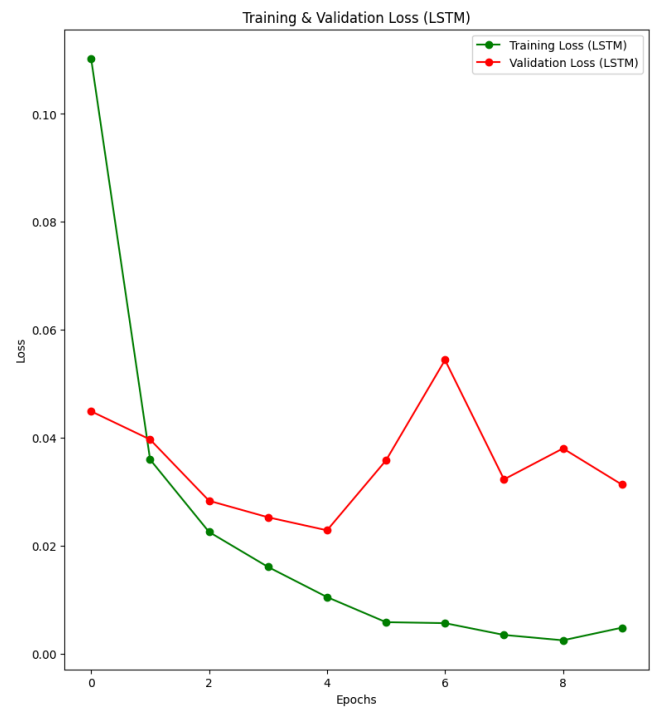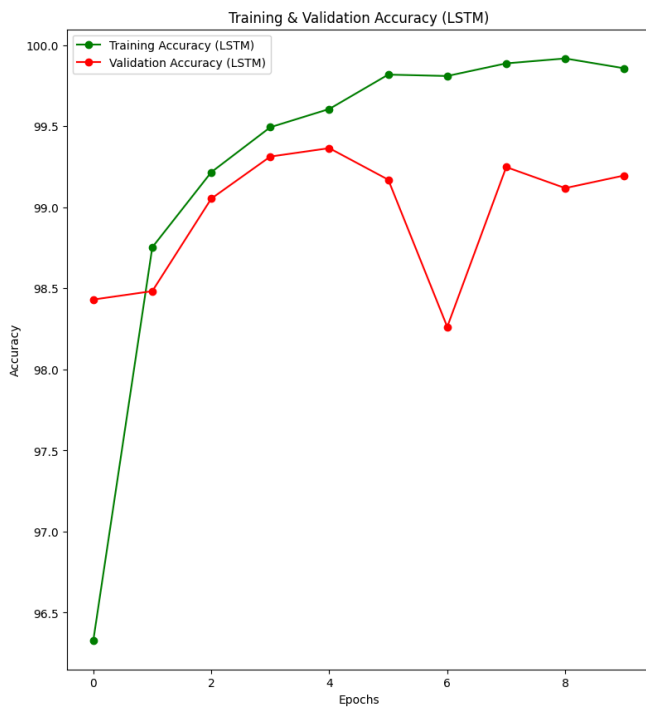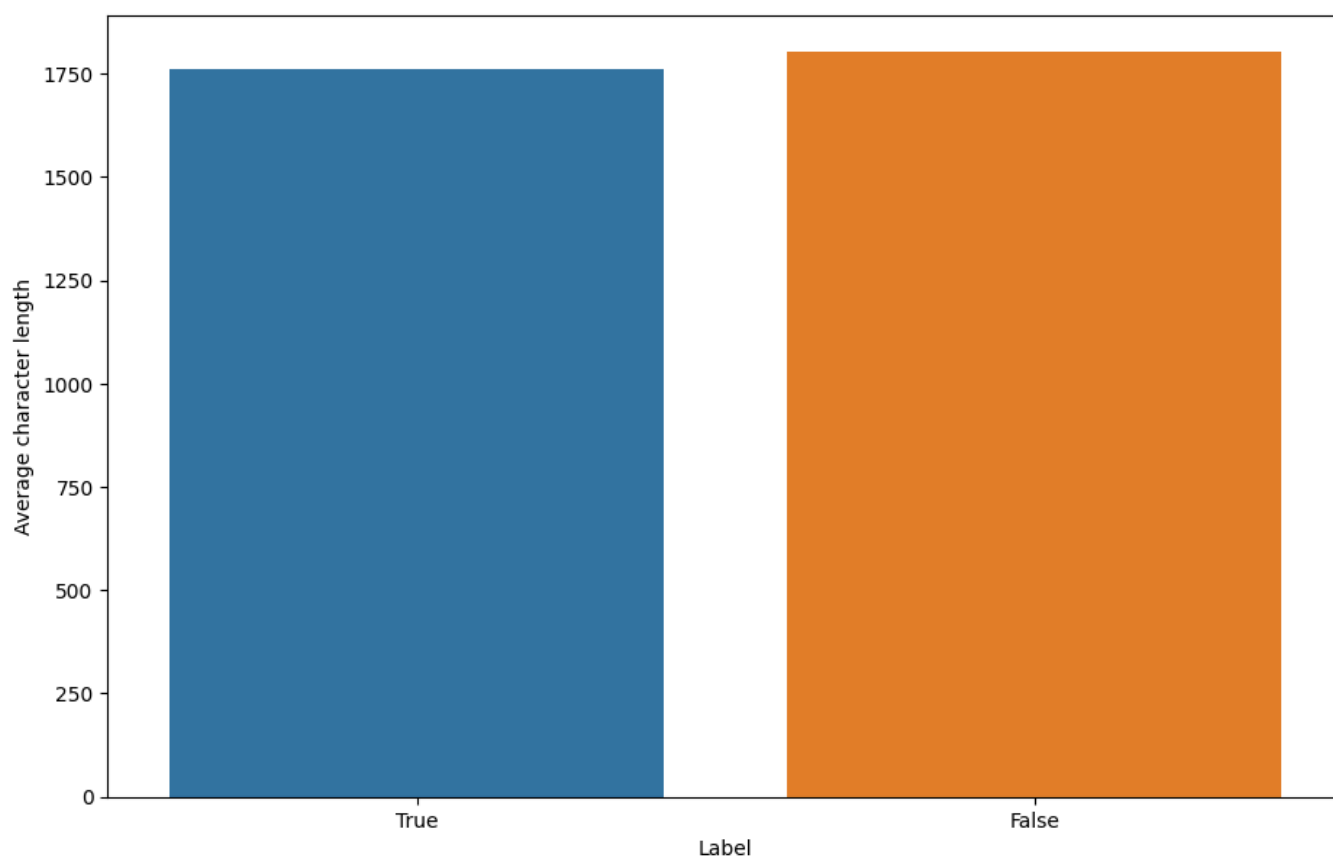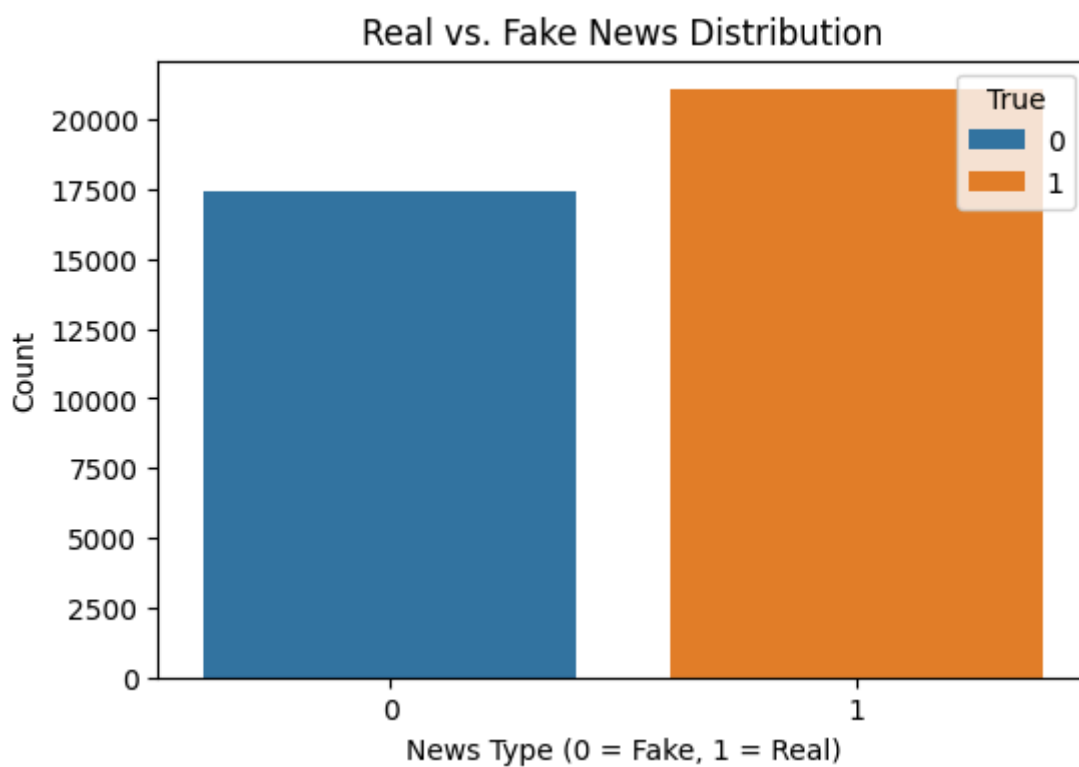
(Fig: Classification Report)



(Fig: Normalized Confusion Matrix)

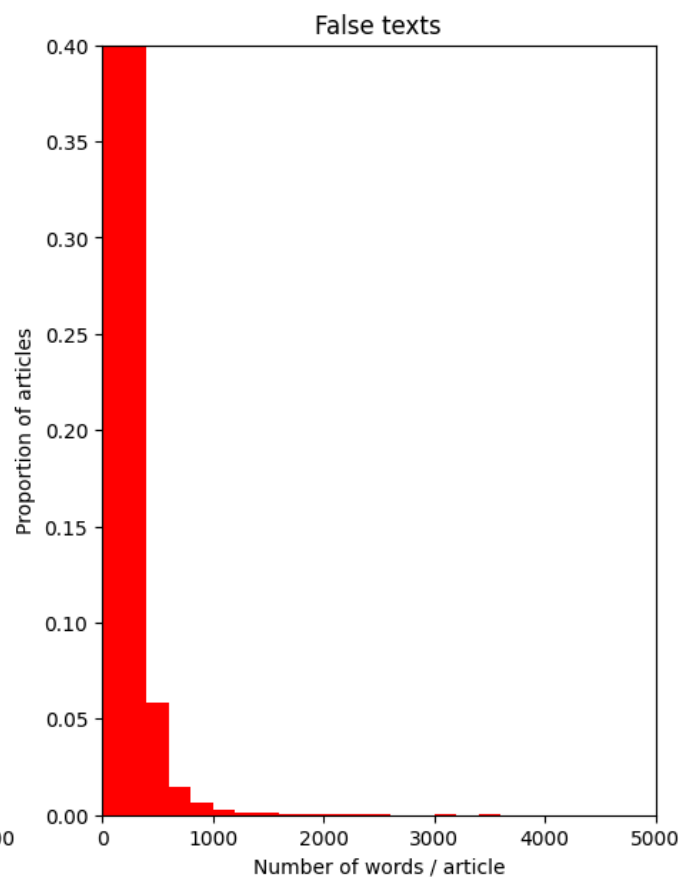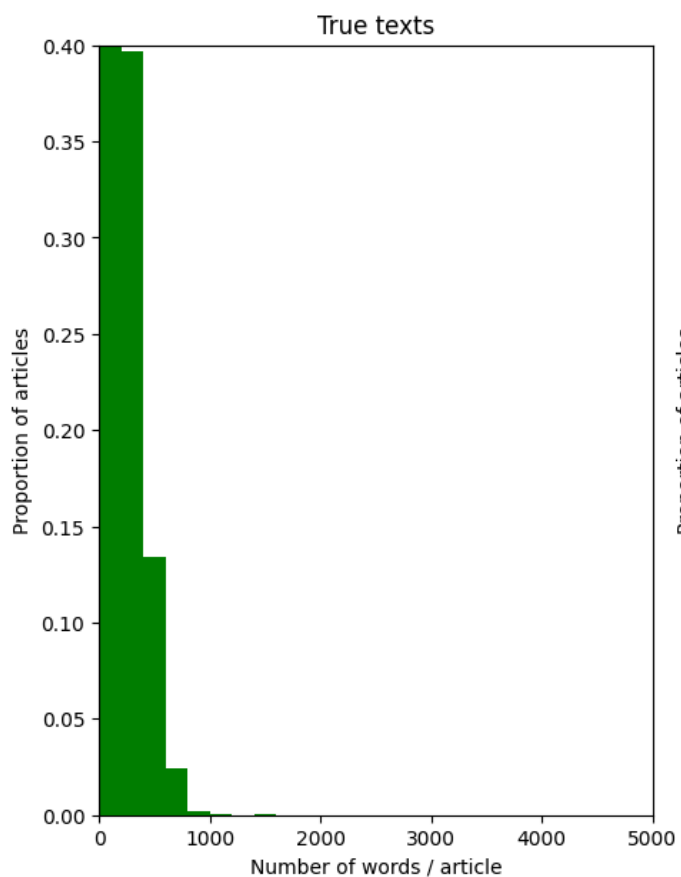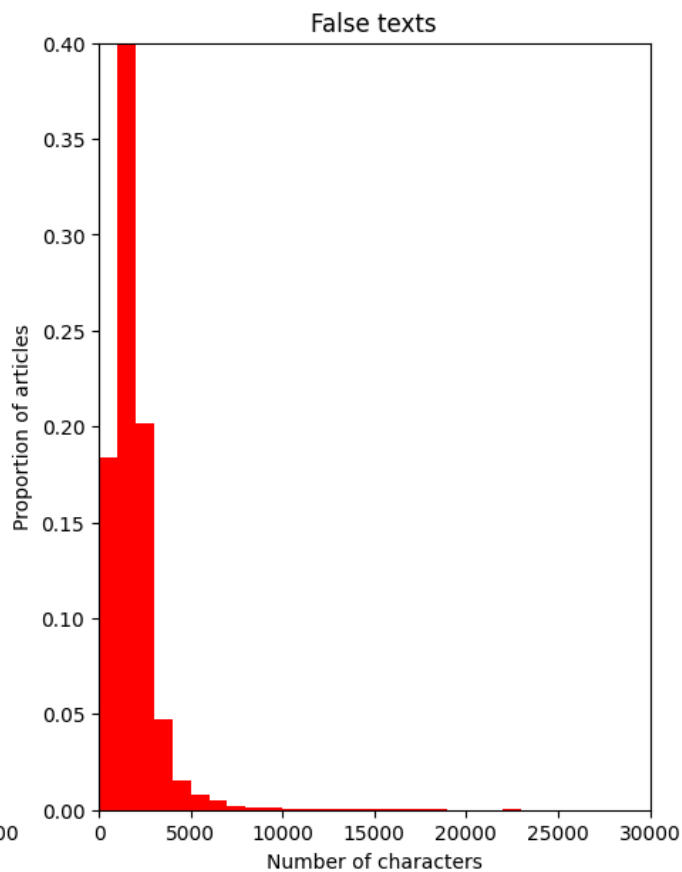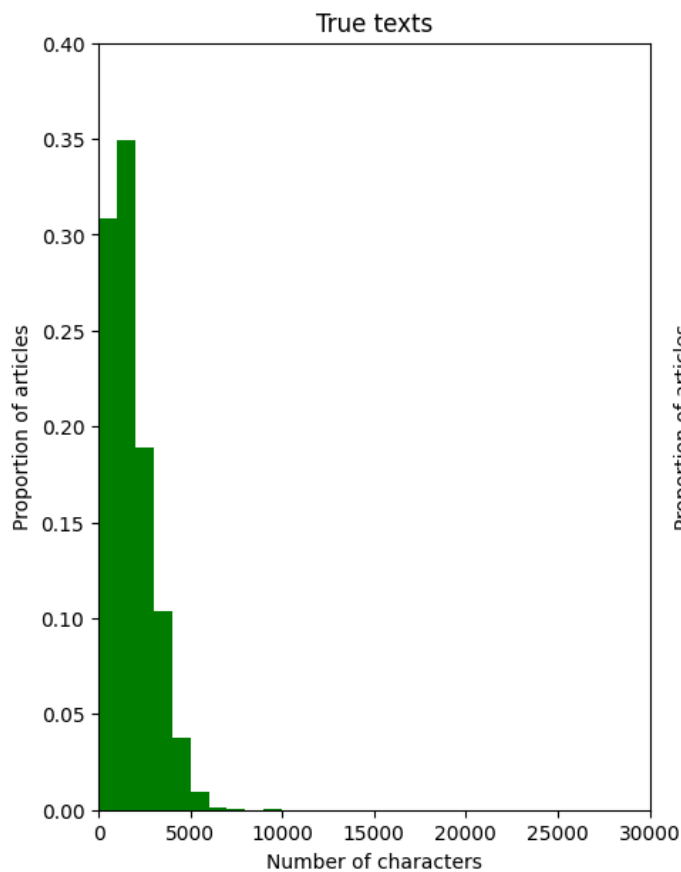(Fig:Training & Validation Accuracy and Loss for CNN)



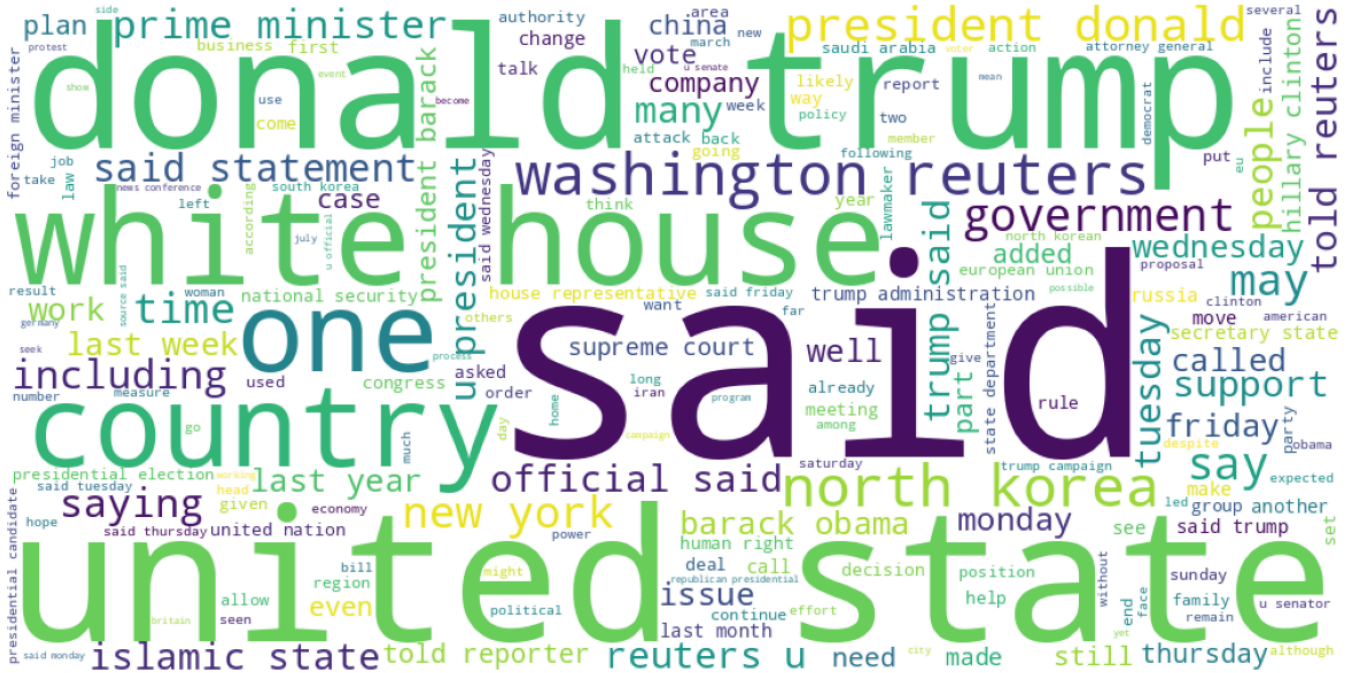(Fig:Training & Validation Accuracy and Loss for LSTM)

Real vs. Fake News Distribution

## True texts

Proportion of articles vs Number of characters

## False texts

Proportion of articles vs Number of characters

## True texts

Proportion of articles vs Number of words / article

## False texts
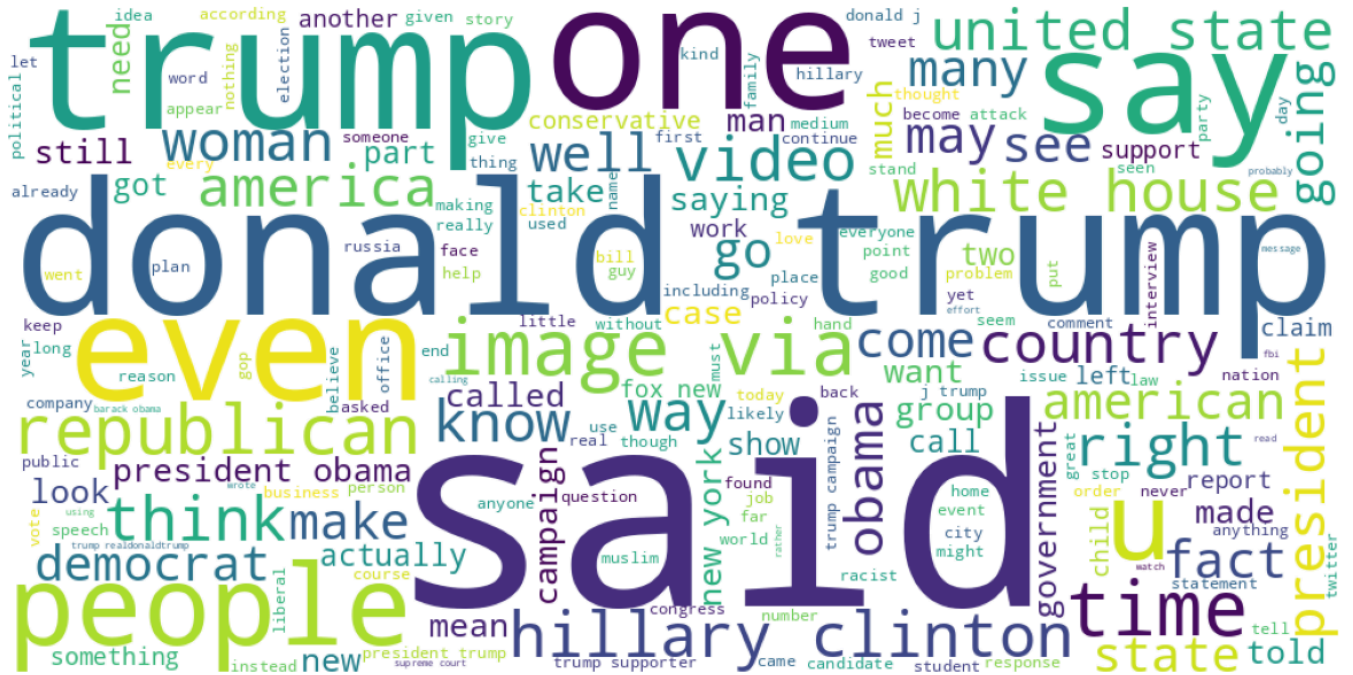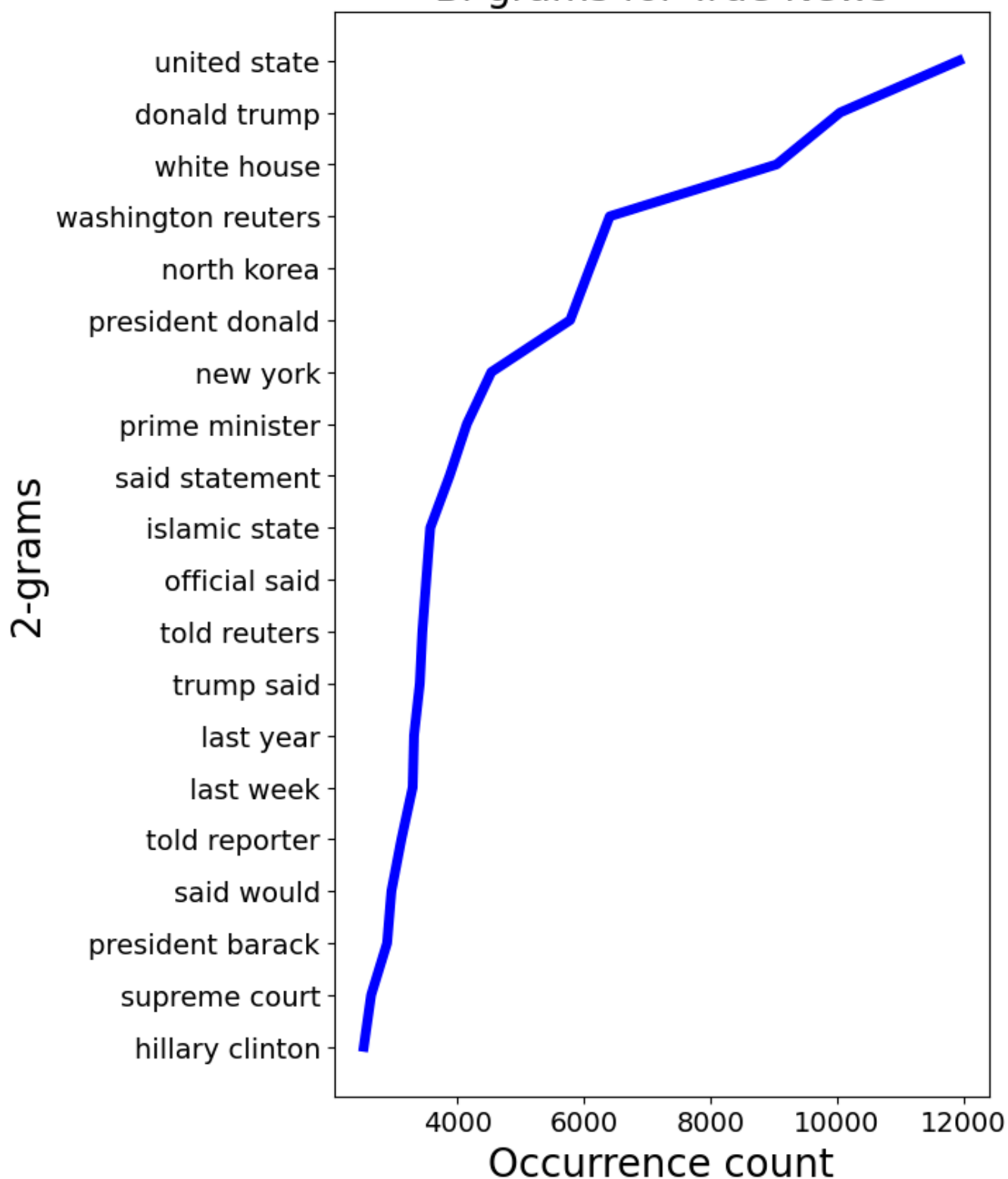
Proportion of articles vs Number of words / article

WordCloud for True News


WordCloud for Fake News

Bi-grams for True News

# Bi-gram for Fake News



Chart showing bi-gram occurrence counts. Y-axis (2-grams, top to bottom): donald trump, image via, hillary clinton, white house, united state, president obama, new york, fox news, president trump, trump supporter, trump campaign, trump realdonaldtrump, supreme court, barack obama, ted cruz, fake news, republican party, screen capture, bernie sander, trump said. X-axis: Occurrence count (2000 to 14000).