

# Cloud-based Natural Language Processing Technology Performs Fake News Identification

Wenxuan Zhang

*Electrical and Computer Engineering  
University of Toronto  
Toronto, Canada  
1009230388*

Hang Li

*Electrical and Computer Engineering  
University of Toronto  
Toronto, Canada  
1010007181*

Mengxuan Ding

*Electrical and Computer Engineering  
University of Toronto  
Toronto, Canada  
1003965319*

Zhengnan Zhu

*Electrical and Computer Engineering  
University of Toronto  
Toronto, Canada  
1005209971*

**Abstract**—This project aims to tackle the spread of fake news by developing a cloud-based natural language processing (NLP) classifier to differentiate between true and fake news. Utilizing advanced NLP techniques and machine learning algorithms, the classifier will analyze news content and categorize it as genuine or counterfeit. The project involves data collection, model training, and deployment as a scalable cloud service. Challenges include data quality, news source biases, and evolving misinformation tactics. The effectiveness will be measured by the classifier’s accuracy in combating fake news and enhancing information reliability. The project contributes to the research field by exploring the implementation and performance of NLP technology in cloud environments.

**Index Terms**—Natural Language Processing, Fake News Detection, Machine Learning, Cloud Computing, News Classification, Information Reliability, Data Preprocessing, Model Optimization, Scalability, Misinformation Tactics

## 1. Introduction

In the information era, the rapid dissemination of information through online platforms has been accompanied by the proliferation of fake news. This misinformation can have wide-range effects, from influencing public opinion and election outcomes to undermining trust in institutions. Addressing the challenge of quickly and accurately identifying fake news is critical to maintaining the integrity of public discourse. This project aims to leverage cloud-based natural language processing (NLP) technologies to develop a classifier that can distinguish between true and fake news, thereby providing a tool for people to verify the authenticity of the information they encounter.

The scope of this project includes the design and development of a sophisticated NLP classifier deployed on cloud platforms designed to analyze news content and classify it

as either fake or true. This solution seeks to use advanced NLP techniques, including machine learning algorithms and linguistic analysis, to evaluate the credibility and authenticity of news articles. The project will be executed in phases, beginning with the collection and preprocessing of a labeled dataset, followed by the training and optimization of the classification model, and culminating in the deployment of the classifier as a scalable cloud service.

The classifier’s development will be constrained by several limitations, including the availability and quality of labeled datasets for training purposes, the inherent biases in news sources, and the evolving nature of misinformation tactics. Additionally, the computational resources required for processing large datasets and the need for ongoing model refinement to adapt to new misinformation strategies will pose challenges.

Expected capabilities of the solution include high accuracy in classifying news articles, scalability to handle large volumes of data, and adaptability to new patterns of misinformation.

By addressing these objectives, the project aims to create a valuable tool for combating fake news, thereby contributing to the reliability of online information and supporting informed public discourse. The success of this initiative will be measured by the classifier’s accuracy and usability.

## 2. Related Work

As the field of research in machine learning and artificial intelligence progresses rapidly, there is a corresponding increase in researchers’ demand for extensive computing resources. Consequently, leveraging shared computing and storage resources in the cloud for deploying and training machine learning models has emerged as a viable option for many researchers. In [1], the authors introduced the DEEP-Hybrid-DataCloud architecture, designed to furnish

distributed resources to researchers via cloud-directed services, catering to computationally intensive tasks in machine learning. Moreover, DEEPaaS, founded on a serverless architecture, encapsulates its user model as a function, facilitating easy sharing with the developer community.

In a related study [2], researchers initiated performance evaluations contrasting several machine learning algorithms in both cloud and local environments, using a virtual setup that emulates real cloud scenarios. The findings indicate that cloud-based distributed machine learning technology surpasses existing systems in terms of training time, resource utilization, and scalability. Furthermore, in [3], the author evaluates and compares the execution proficiency of machine learning tasks across three mainstream cloud platforms: Amazon Cloud Service, Google Cloud Platform, and Microsoft Azure. The results underscore the capability of cloud-based machine learning functionalities to support diverse business needs across sectors such as medical, retail, and automotive.

By delving into these pioneering research findings, it becomes apparent that numerous researchers are turning to cloud platforms for implementing and evaluating the performance of machine learning models. Consequently, exploring and assessing the implementation of various artificial intelligence technologies on cloud platforms emerge as a noteworthy research area. Thus, our interest lies in investigating how natural language processing (NLP) technology is implemented in the cloud and evaluating its performance relative to local environment deployments.

### **3. Main Concepts**

#### **3.1. Fake News and Its Properties**

Fake news refers to misinformation or disinformation typically spread via traditional news media or online platforms, often with the intent to deceive or manipulate public opinion. The properties of fake news that natural language processing technologies can leverage for detection include linguistic cues, sensationalist language, inconsistency with known facts, and the manipulation of emotional appeal. Linguistically, fake news tends to exhibit certain stylistic features, such as the use of exaggerated punctuation, capitalization to invoke emotion, and a higher frequency of subjective statements compared to objective reporting. NLP techniques can analyze these textual features, along with the structure and semantics of the content, to classify news as fake or genuine with a high degree of accuracy [5].

#### **3.2. Natural Language Processing**

Natural Language Processing (NLP) is a pivotal branch of artificial intelligence that enables the automated understanding, interpretation, and generation of human language. This technology is instrumental in various applications, ranging from language translation to sentiment analysis. NLP stands out as a particularly effective tool for fake news

detection due to its capacity to process vast amounts of textual data swiftly, discern linguistic and semantic patterns, and identify inconsistencies indicative of misinformation. By analyzing the structure, context, and sentiment of news content, NLP algorithms can detect subtle cues that differentiate genuine news from fake. These capabilities are grounded in sophisticated machine learning models that are trained on extensive datasets comprising both authentic and fabricated news articles. Through this training, NLP models learn to recognize the hallmarks of fake news, such as exaggerated language, factual inaccuracies, and manipulative content. The integration of NLP in fake news detection efforts represents a significant advancement in the fight against misinformation, offering a scalable and effective solution for verifying the authenticity of information disseminated across digital platforms [6].

#### **3.3. Docker Container**

Docker is a platform that enables developers to package applications into containers—standardized executable components combining application source code with the operating system (OS) libraries and dependencies required to run that code in any environment. Docker containers ensure that software will behave the same way regardless of where it's deployed. This makes Docker particularly beneficial for projects like ours, where the consistent, scalable deployment of NLP models is crucial. Docker facilitates the easy distribution and deployment of applications across various environments, enhancing the portability and scalability of cloud-based applications. By using Docker, developers can easily package the NLP application and its environment into a container, which can then be deployed seamlessly on any cloud platform, ensuring the application runs efficiently and reliably in different computing environments [7].

### **4. Research Activities**

To achieve the objectives of this project, we have divided the cloud-based NLP fake news recognizer into two main components: a pre-trained NLP model and a deployed recognizer application on the cloud. To conduct the model training, we have currently gathered Fake News datasets [4] from Kaggle, which consist of two CSV tables representing genuine news and fake news, respectively. The dataset representing genuine news comprises 20,826 entries, while the fake news dataset consists of 17,903 entries. These data entries record four news attributes: title, body, subject, and date. Below outlines our methodology for data collection and preprocessing, aimed at gathering news information from real-life scenarios along with a series of usable details, and transforming them into formats suitable for NLP model consumption. Subsequently, we detail the NLP model employed for the recognizer in this project and the process of deploying it to the cloud.

## 4.1. Data Collection and Pre-processing

For data collection, we surveyed prominent artificial intelligence platforms such as Hugging Face and Kaggle, ultimately obtaining a reliable and sufficiently sized Fake News dataset [4] from Kaggle. Subsequently, we conducted data cleansing on all collected data, removing entries with missing values, and discarding data attributes unsuitable for the selected model, such as the source website of the news, which cannot be effectively transformed into vectors recognizable by the model. Additionally, due to the variance in the number of entries between genuine and fake news datasets, there existed a slight imbalance in the dataset utilized for training. To mitigate the potential impact of this issue, we employed a stratified resampling technique, ensuring a balanced representation of both labels within the training set.

## 4.2. Model Selection

We elected to employ a fine-tuned CNN-based model as our NLP recognizer model. This CNN model comprises four primary components: an Embedding layer, convolutional layers, a Max-pooling layer, and fully connected layers. The Embedding layer represents the input matrix. Convolutional layers are responsible for extracting features from the input and employ the Rectified Linear Unit (ReLU) function as the activation function. The MaxPool layer selects the maximum value across sentence lengths and reduces the output volume produced by each kernel. Finally, the fully connected layers activate through the Sigmoid function, generating scalar outputs to represent the probability of news belonging to genuine or fake categories. To enable the CNN model to perform NLP tasks, we fine-tuned it by transforming the data into a TextDataset object, an extension of the PyTorch Dataset class, thereby enabling text input for the CNN model. Moreover, to ensure successful conversion of inputs into usable Embedding layer matrices, we devised a padding function to pad shorter sentences within the same batch, ensuring uniform length across input strings in a given batch. As for the configuration of kernel numbers for the CNN model, we designed it as adjustable hyperparameters, conducting multiple experiments to identify the optimal parameter selection scheme. The specific scheme,  $k1=2$ ,  $k2=4$ ,  $n1=20$ ,  $n2=20$ , was determined to achieve the highest accuracy. Finally, the trained CNN model will be saved in the local environment for subsequent steps of cloud deployment.

## 4.3. Cloud Deployment

To facilitate the utilization of the AI model for cloud deployment applications, we devised a simple UI interface based on the Gradio module. This interface, coupled with the stored CNN model, enables the functionality of deploying a fake news recognition application on the cloud. Within this UI interface, users merely need to simply copy and paste the news title or content into the input box and submit to obtain

the NLP model's prediction regarding its authenticity. For the cloud deployment process, we utilized the free Docker platform and employed its Docker Compose functionality to create images. We defined application dependencies through the Image Composition Process, created a dedicated Dockerfile, defined services in a compose file, and ultimately used Compose to build and successfully run our fake news recognition application. Through experimentation, this application has been effectively deployed and can provide a public link for external use. Finally, we pushed the image of this application to Docker Hub, creating a repository for it, enabling it to be pulled into other local environments via Docker Desktop.

## 5. Contributions

### 5.1. Clarity and Precision

**Model Detailing:** CNNs (Convolutional Neural Networks) are chosen for analyzing news texts in NLP tasks because they're really good at picking up patterns in the data, even if it's as varied and complex as language. They can handle texts of different lengths easily, find important words or phrases without being told explicitly what to look for, and do all of this relatively quickly, which is great for processing lots of news articles. Plus, they've been proven to work well for sorting texts into categories, like figuring out if a piece of news is real or fake. Even though CNNs is an old model and were initially designed for images, these characteristics make them quite effective and efficient for dealing with language data, too.

**Data Handling:** Our dataset has been partitioned into three distinct subsets: training, validation and testing data. The training data can be used to train the model. The validation data plays an important role in the selection of hyperparameters. Lastly, the testing data is utilized as an objective measure to assess the model's accuracy, providing an unbiased evaluation from the simulated real-world scenarios. This approach can achieve a accurate and reliable model finally.

### 5.2. Technical Advantages

Deploying our model on a cloud system brings several significant advantages, especially in terms of model updates, parallel programming, and computational demands. Firstly, as the model improves or as new data becomes available, we can easily integrate these changes to enhance accuracy and performance without any downtime for the users. Secondly, algorithms our model uses, especially in neural networks, involves a lot of matrix calculations. By running these calculations on a cloud server, they can be speed up by using parallel programming, which makes our model more efficient and faster in handling complex tasks. Thirdly, the users' hardware capabilities may be limited, so we offload the heavy lifting from the users' devices by leveraging the cloud.

### 5.3. Application Scenarios and Social Value

Application by Social Media Organizations: Given how popular social media is and how important it is for media organizations to be seen as reliable and trustworthy, our system stands out as a key tool. It acts as an effective filter for checking news, allowing these organizations to avoid spreading false information. This is done much more cheaply and quickly than older, manual checking methods, improving the trustworthiness and credibility of what they share.

User Interaction: People receive news from various sources every day. However, many of us don't know how to determine if the information is true or false. Our program can help by scanning and filtering news on specific websites and informing users about the authenticity of the news. Moreover, with a little adjustment to the model, it can even show the likelihood of the news being true or false. Totally, it can help people avoid the harm of fake news.

### 5.4. Future direction

In future, we can improve our models with more new coming input data and we may use some other advanced NLP technologies (like Transformer models). Additionally, we can create more APIs with different functions. For example, the likelihood of the news being true or false / the news from specific regions (e.g. science, politic, economy etc.).

## 6. Research Findings Blocks

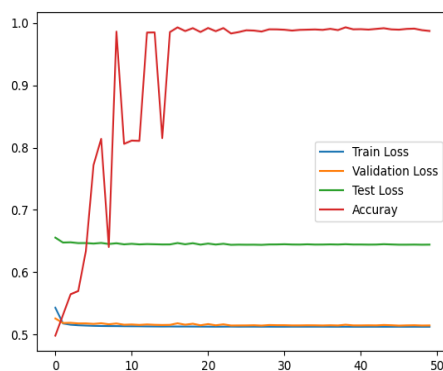


Figure 1. loss plot

The chart ("Fig. 1") shows that initially, all types of loss (training, validation and testing) decreases and then stabilize after a few epochs. This indicates that the model is effectively learning from the data provided. Moreover, the test loss is higher than both the training and validation loss, which is expected. The model is trained on the training data and the validation data is used to fine-tune the settings, meaning that both sets of data are somewhat 'known' to the model. Hence, the model 'recognizes' these data, resulting

in lower loss values. The accuracy is showing a pattern of fluctuating increase, ending at a high level (almost 100%), suggesting that the model is trained well without overfitting. All signs from the figure point to successful model training with excellent performance.

## 7. Conclusion

In conclusion, this project successfully developed and deployed a cloud-based natural language processing (NLP) classifier to identify and combat the spread of fake news. Utilizing advanced NLP techniques and machine learning algorithms, the classifier efficiently analyzed news content and categorized it as genuine or counterfeit with a high degree of accuracy. The project encompassed data collection, model training, and deployment as a scalable cloud service, addressing challenges such as data quality, news source biases, and evolving misinformation tactics.

The implementation of a fine-tuned CNN-based model for the NLP recognizer proved effective in processing textual data and distinguishing between true and fake news. The deployment of the classifier as a cloud service, facilitated by Docker and Gradio, enhanced its accessibility and scalability, making it a valuable tool for users to verify the authenticity of information encountered online.

Overall, this project contributes to the field of NLP and cloud computing by showcasing the potential of cloud-based NLP technologies in addressing the critical issue of fake news dissemination. Future work could focus on further optimizing the classifier's performance, expanding its capabilities to handle multilingual content, and exploring integration with social media platforms for real-time fake news detection.

## References

- [1] A. Lopez Garcia et al., "A Cloud-Based Framework for Machine Learning Workloads and Applications," IEEE Access, vol. 8, pp. 18681–18692, 2020, doi: <https://doi.org/10.1109/access.2020.2964386>.
- [2] I. Sakthidevi, G. V. Rajkumar, R. Sunitha, A. Sangeetha, R. S. Krishnan, and S. Sundararajan, "Machine Learning Orchestration in Cloud Environments: Automating the Training and Deployment of Distributed Machine Learning AI Model," 2023 7th International Conference on I-SMAC (IoT in Social, Mobile, Analytics and Cloud) (I-SMAC), Oct. 2023, doi: <https://doi.org/10.1109/i-smac58438.2023.10290278>.
- [3] A. Jagati and T. Subbulakshmi, "Building ML Workflow for Malware Images Classification using Machine Learning Services in Leading Cloud Platforms," 2023 International Conference on Computational Intelligence and Sustainable Engineering Solutions (CISES), Apr. 2023, doi: <https://doi.org/10.1109/cises58720.2023.10183421>.
- [4] "Fake News Detection Model," kaggle.com. <https://www.kaggle.com/code/noorsaeed/fake-news-detection-model/input> (accessed Mar. 28, 2024).
- [5] Zhou, X., & Zafarani, R. (2018). A survey of fake news: Fundamental theories, detection methods, and opportunities. ACM Computing Surveys (CSUR), 51(5), 1-40. doi:10.1145/3395046
- [6] X. Zhang and A. A. Ghorbani, "An overview of online fake news: Characterization, detection, and discussion," Information Processing & Management, vol. 57, no. 2, p. 102025, Feb. 2020.

- [7] Merkel, D. (2014). Docker: Lightweight Linux containers for consistent development and deployment. *Linux Journal*, 2014(239).