# Text Summarization and Temporal Learning Models Applied to Portuguese Fake News Detection in a Novel Brazilian Corpus Dataset

**Gabriel Lino Garcia** and **Pedro Henrique Paiola** and **Danilo Samuel Jodas**
and **Luis Afonso Sugi** and **João Paulo Papa**
Department of Computing, São Paulo State University, São Paulo, Brazil
{gabriel.lino,pedro.paiola,danilo.jodas,luis.afonso,joao.papa}@unesp.br

## Abstract

Streaming content advances and the appearance of online media raised the ability for massive content sharing that reaches thousands of people worldwide in a real-time fashion. Fake news spreading is nowadays the main concern of several authorities worldwide due to the negative impact and potential to induce social and political instability in our society. Therefore, fake news detection and suppression gained increased attention as an important topic in natural language processing and machine learning academic research. Regardless of the state-of-the-art methods available for fake news detection, a good corpus revealing novel language-specific counterfeit aspects is also important to exploit and distinguish between real and fake news in the context of social and political impacts for specific regions. This paper extends a previous Brazilian Portuguese corpora dataset and proposes using and comparing several deep learning and classical machine learning models to detect counterfeit content in the Portuguese language. Moreover, we propose using text summarization to achieve concise news summaries and prevent losing relevant information. This work presents an updated and balanced version of the FakeRecogna dataset for detecting fake news articles using a temporal learning approach based on efficient and well-known deep learning models.

## 1 Introduction

Social and online media have emerged as innovative and rapid communication sources in the last few years. It promotes an easy medium for sharing data that reaches millions of people worldwide. While massive data can be readily spread in real-time using social media, it can also be slanted to bias public opinion's perception and lead to misconceptions that may lead to social and political instabilities. Such practice, usually called fake news, is defined by Allcott and Gentzkow (2017) as the intentional production of fake content that seeks to lead to false impressions and misconceptions by the readers.

In this context, an in-depth exploration of textual and visual information has been proposed to cope with fake news detection by using natural language processing (NLP) models and features extracted from images (Singhal et al., 2019). State-of-the-art works tackled the fake news detection problem using news published in English. Regardless, the focus of this paper is to use content published in the Portuguese language. However, most studies used out-to-date corpus with only a few samples to design fake news detection systems using Portuguese texts. On this matter, Garcia et al. (2022) proposed FakeRecogna, a novel Portuguese fake news dataset, to achieve more representative samples with the latest news articles organized into the most meaningful news categories in Brazil. Monteiro et al. (2018) presented the Fake.Br, a corpus containing 7, 200 Portuguese news collected between 2015 and 2018. On the other hand, Charles et al. (2022) assembled a full-bodied corpus dataset with 12, 398 news articles collected between 2013 and 2021.

Fake news spreading has widely increased in the last few years, providing new opportunities to support a broader assessment regarding the up-to-date aspects related to counterfeit content. This work extends the previous research in Portuguese fake news detection by supporting the gathering of new data to compose a full-bodied and large dataset with more than 52, 000 real and fake news articles collected from well-known Brazilian agency news. Moreover, we propose using extractive and abstractive text summarization and a temporal learning approach based on Long Short-Term Memory (LSTM), Gated Recurrent Unit (GRU), and the Bidirectional Encoder Representations for Transformers (BERT) model to predict fake news using text representation. Classical machine learning models were also assessed in terms of the fake news

prediction over the proposed dataset.

The main contributions of this work are summarized in the following key points:

- To extend a new balanced version of the FakeRecogna dataset with more than 52,000 news articles in the Portuguese language;

- To apply extractive and abstractive summarization to the news of the proposed dataset;

- To propose the use of temporal learning models to enhance fake news detection;

- To provide up-to-date research on the novel fake news aspects in the context of the Brazilian Portuguese.

## 2 Related Works

Several studies have been proposed for using NLP solutions to explore and understand the aspects behind counterfeit content in English by combining several machine learning and deep learning methods (Ruchansky et al., 2017; Oshikawa et al., 2018; Zhang and Ghorbani, 2020; Kesarwani et al., 2020; Zhou and Zafarani, 2020; Mishra et al., 2022). However, researchers have also explored fake news detection in the context of the Portuguese language. Endo et al. (2022) further investigated fake news detection during the COVID-19 pandemic using online communications based on Brazilian Portuguese content. Faustini and Covões (2019) addressed fake news detection in Brazil by leveraging research on anomaly detection using only fake news instances to train a One-Class Classification model. In a similar approach, Garcia et al. (2023) proposed a large and rich fake news dataset to harness research on anomaly detection methods by offering an imbalanced dataset and promoting novel classes of Portuguese counterfeit content. The proposed dataset is imbalanced since the fake news samples are assumed to be outliers in the data, thus leading to many more real news samples.

Large Language Models (LLMs) have transformed the computer generative capabilities in a broad range of deep learning applications. In the NLP scenario, such powerful networks are trained on huge amounts of textual data to evolve the manner in which computers understand and produce textual information. In a recent study, LLMs have been applied to detect counterfeit Portuguese content using the second version of the

Large Language Model Meta AI (LLaMA 2) architecture (Garcia et al., 2024). The study proposed a trained version of the LLaMA 2 architecture utilizing the Low-Rank Adaptation (LoRA) method (Hu et al., 2021) in the Portuguese version of the Alpaca dataset (Larcher et al., 2023). The study revealed the LLMs' capacity to cope with the increasing spreading of fake information.

Summarization works for fake news detection in the Portuguese language are scarcer than fake news detection research in English, mainly due to the lack of annotated summary datasets. Notably, important efforts have been attained by the Interinstitutional Center for Computational Linguistics (NILC), many of which are focused on Multidocument Summarization (Souza and Felippo, 2018) or Opinion Summarization (Inácio and Pardo, 2021; López Condori and Salgueiro Pardo, 2017). Regarding news summarization, the PTT5-Summ proposed by Paiola et al. (2022) can be cited, which involves adapting the PTT5 model (Carmo et al., 2020) for the task of abstractive summarization through fine-tuning with Portuguese annotated news datasets.

In English-language research, we also find models in the literature for fake news classification that used the news summaries as input. Esmaeilzadeh et al. (2019) investigated the application of deep learning models in fake news detection and conducted experiments using the original news and their summaries as input. The authors observed a slight increase in accuracy in fake news detection when using the summaries. Hartl and Kruschwitz (2022) also explored a fake news detection method based on automatic summarization, proposing the Contextual Multi-Text Representations for fake news detection with BERT (CMTR-BERT) model, which combines different textual representations and additional contextual information to build a more condensed version of the original text.

## 3 Proposed method

Figure 1 illustrates the steps of the proposed method. Each step is described in details in the next sections. The fake news collection was performed on licensed and verified Brazilian news websites with enrollment in the Duke Reporters' Lab Center[1] released by the Sanford School of Public Policy journalism center at Duke University.

---

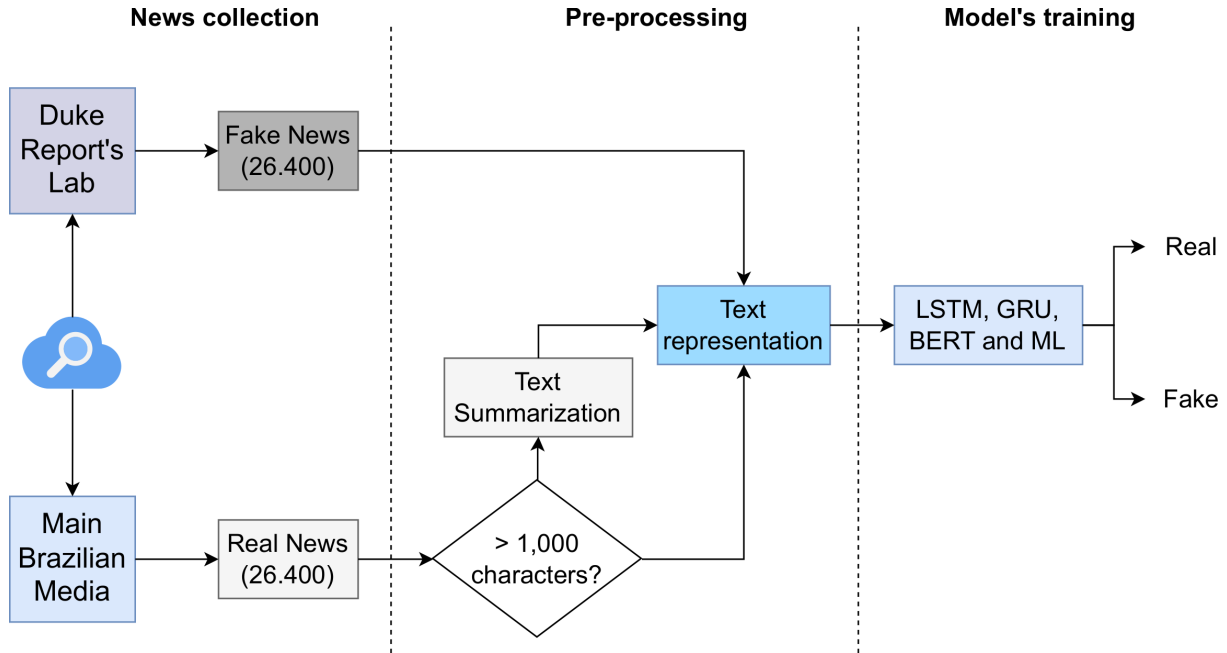[1] https://reporterslab.org/
fact-checking

Figure 1: Pipeline of the proposed method.

The system was designed as a source to fight against fake news spreading worldwide. For real news, we selected well-known media platforms in Brazil. Since real texts are much larger than most of the produced fake content, the genuine news was preprocessed with text summarization. At this stage, there is no further processing of stop words or lemmatization of the text. After trimming and standardizing the real news, we produced textual representations based on Bag of Words (BoW), Term Frequency – Inverse Document Frequency (TF-IDF), FastText, PTT5, and BERTimbau (Souza et al., 2020) to form the input feature vectors for the ML models.

### 3.1 FakeRecogna 2.0 Dataset

This section presents the proposed extension for the FakeRecogna dataset in the context of fake news detection. FakeRecogna includes real and fake news texts collected from online media and ten fact-checking sources in Brazil. An important aspect is the lack of relation between the real and fake news samples, i.e., they are not mutually related to each other to avoid intrinsic bias in the data. Details of the news collection and categorization are provided in the next sections.

#### 3.1.1 Data collection

The news collection was performed using web crawlers specifically designed to seek pages from well-known agencies with national importance.

Each news page was subsequently processed to extract relevant information from the news so that we can prevent citations to other articles, advertising, and texts that may end up being part of the news story. After that, the news was classified in chronological order.

#### 3.1.2 Fake News Mining

Fake news mining was performed on pages collected between 2019 and 2023 from the Duke Reporters Lab. This respected agency presently cooperates with 417 active fact-checking agencies globally, nine of them operating in Brazil. Moreover, they keep a list of pages committed to proving the validity of news sources.

#### 3.1.3 Fake News Sources Selection

Fake news sources were selected from nine fact-checking agencies in Brazil. This process provides a broad range of categories and many fake news samples to promote data diversity. Table 1 presents the existing Brazilian fact-checking initiatives and the number of fake news samples collected from each news source. When the search process was concluded, we ended up with $26,569$ fake news samples, which, in turn, were further processed to detect and remove possible duplicate samples, thus leading to a final set of $26,400$ fake news articles.

Table 1: Fact-checking agencies in Brazil.

| Agency | Web address | # news |
|---|---|---|
| *AFP Checamos* | https://checamos.afp.com/afp-brasil | 1,587 |
| *Agência Lupa* | https://piaui.folha.uol.com.br/lupa/ | 3,147 |
| *Aos Fatos* | https://aosfatos.org | 2,720 |
| *Boatos.org* | https://boatos.org | 8,654 |
| *Estadão Verifica* | https://politica.estadao.com.br/blogs/estadao-verifica | 1,405 |
| *E-farsas* | https://www.e-farsas.com | 3,330 |
| *Fato ou Fake* ("Fact or Fake") | https://oglobo.globo.com/fato-ou-fake | 2,270 |
| *Projeto Comprova* | https://projetocomprova.com.br | 877 |
| *UOL Confere* | https://noticias.uol.com.br/confere | 2,579 |
| **Total** | | 26,569 |

### 3.1.4 Data organization

We established several thematic classes to facilitate the data organization and support the initial pages' content categorization. After that, all news were grouped according to their published data. This process yields a range of news sources and different writing styles that ensure data diversity and a suitable data structure for NLP and machine learning algorithms. The collected texts are distributed into nine categories in relation to their main subjects: Brazil, Conspirations, Entertainment, Health, Politics, Science and Technology, Social Media, Sports, and World. These categories are determined by the journal sections from which the news articles were extracted. Figure 2 illustrates the news distribution across each defined category along with the respective percentages.
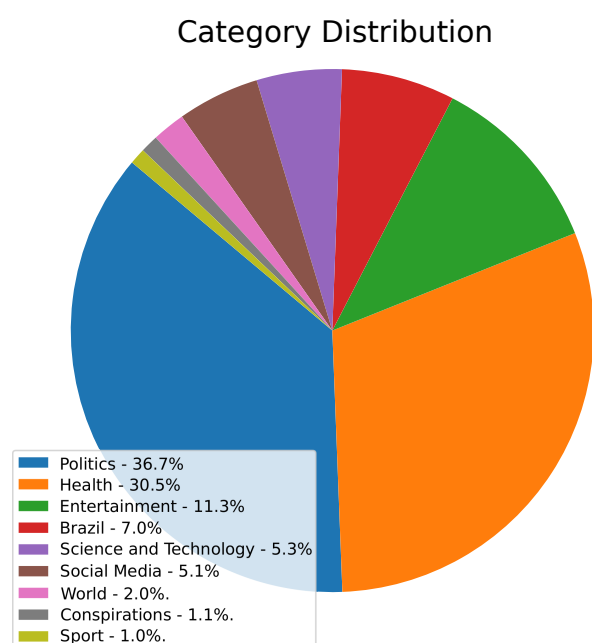


Figure 2: Fake news distribution by category.

Table 2 provides instances of both authentic and fabricated articles, illustrating the contrasting content sizes between the two types of news.

Table 2: Example of fake and true news.

| Fake | Real |
|---|---|
| Publicações nas redes sociais usam dados de uma pesquisa brasileira para acusar pesquisadores de tramarem contra o uso de cloroquina no tratamento de pacientes com a covid-19. algumas postagens acusam pesquisador de ser ligado ao pt. | O Cristo Redentor vai reabrir para o público neste sábado (15), depois de passar cinco meses fechado por causa da pandemia de covid-19. Hoje, o local passa por uma desinfecção para receber os visitantes. O trabalho começou às 7h, em uma parceria da Arquidiocese do Rio de Janeiro, do Parque Nacional da Tijuca e do Comando Conjunto Leste. [...] |

### 3.1.5 FakeRecogna vs FakeRecogna 2.0

The FakeRecogna 2.0 has nearly increased 5 times the original size of its counterpart version, FakeRecogna 1.0, which previously comprised 11,902 news samples spread across the real and fake news classes. Conversely, FakeRecogna 2.0 includes a total of 52,800 news articles. Both datasets are balanced when considering the number of samples distributed across the real and fake news categories. However, FakeRecogna 2.0 was expanded to include articles collected from 3 additional communication channels affiliated with fact-checking initiatives in Brazil, totaling 9 agencies to gather the additional data to assemble the new dataset

version. For comparison purposes, FakeRecogna 1.0 was assembled by news collected from only 6 fact-checking Brazilian agencies. Furthermore, the news collection strategy adopted in this study yielded an increase in the number of categories compared to FakeRecogna 1.0, leading to an increase from 6 to 9 categories in FakeRecogna 2.0. However, as reported in previous research, politics and health are still the major targets for fake news production.

When considering the data pre-processing, we expand the previous research by capitalizing on innovative strategies based on text summarization methods, namely abstractive and extractive summarization, applied to real news content. Moreover, the new pre-processing strategy avoids irrelevant steps like removing stopwords, lemmatization, and removal of words such as "enganoso", "boato" and "#fake" to prevent bias in the data. Punctuation, special characters, and URLs were also removed. Furthermore, we standardized the texts to lowercase letters and summarization of real news.

In summary, FakeRecogna 2.0 represents a significant advancement over the previous version and contributes fundamentally to research in fake news detection in the Brazilian context. This corpus can be a key component in developing more effective solutions for identifying and mitigating the spread of fake information in our ever-evolving media landscape.

## 4 Methodology

This section presents the data preprocessing strategy and briefly describes the ML and deep learning models used for fake news detection in the context of this study.

### 4.1 Data Pre-processing

Real news articles are usually longer than fake news content in most online media sources. This aspect may lead to overload in the training process while introducing bias and overfitting to the model since it might be prone and specialized in detecting all input text as authentic and reliable content. Aiming to preserve the most relevant information in the text, we propose using summaries of true news so that they are smaller and similar to fake news in size. This approach reduces the computational load to machine learning models while preserving the original text information and essence.

We adopted extractive and abstractive summarization to produce accurate summaries for genuine news texts. The first method tends to be immune to inconsistencies and hallucinations since the final summary comprises the most relevant sentences without generating new words and phrases. Conversely, abstractive summarization promotes toe ability to produce novel sentences that vary differently from the original text in terms of semantic and sentence structure. Despite being subject to a broad range of problems in textual generation, it can better condense the main information of a text in a way more similar to a human writer.

For abstractive summarization, we used a BERT-based model (Miller, 2019) to extract embeddings from the text and the $k$-Means algorithm to group and select the sentences. Moreover, we employed the PTT5-Summ model developed by Paiola et al. (2022), which was trained on a news dataset called XL-Sum containing relatively short annotated summaries. Abstractive summarization was only applied to texts with more than 1,000 characters, resulting in summaries with nearly 1,000 characters in size.

### 4.2 Textual representation

Text processing is essential to artificial intelligence and NLP tasks. One of the primary steps in text processing is text representation, which involves converting words or documents into formats that machine learning models can understand and process. This article will examine three popular approaches to text representations: Bag of Words, TF-IDF, FastText, BERTimbau, and PTT-5.

#### 4.2.1 Bag of Words

Bag of Words (BoW) is one of the simplest and most widely used approaches for text representation. In this technique, the text is divided into tokens (words or other elements), and then a vector is created to retain the frequency of each token's occurrence in the document. Each document is represented by a vector where each element corresponds to a unique token, and the value in each element is the frequency of that token in the document (Qader et al., 2019). The primary advantage of BoW is its simplicity and computational efficiency.

#### 4.2.2 TF-IDF (Term Frequency-Inverse Document Frequency)

TF-IDF (Salton and Buckley, 1988) is another common technique for textual representation that considers the frequency of document terms. It assesses

the importance of a term regarding a specific document and a collection of documents. The TF-IDF representation assigns a weight to each term based on its frequency in the document (Term Frequency) and its rarity in the collection of documents (Inverse Document Frequency). TF-IDF is effective in reducing the importance of highly frequent and common terms, such as "a," "de," and "o" while increasing the importance of terms that are distinctive to a specific document or topic. This approach makes it useful in information retrieval and text classification tasks.

### 4.2.3 FastText

FastText is a more advanced and recent approach to textual representations. It is based on word embeddings (word vectors) trained on large amounts of text. The primary innovation of FastText concerns its ability to represent unknown or rare words by breaking them down into subwords (n-grams) and combining the representations of these subwords. This technique is especially useful when dealing with texts in languages with extensive vocabulary, texts with spelling errors, or specific jargon. Additionally, FastText preserves the order of words and captures semantic relationships between words (Bojanowski et al., 2017).

### 4.2.4 BERTimbau

BERTimbau (Souza et al., 2020) is a textual representation based on the BERT model, known for its ability to capture bidirectional contexts of words. In the context of BERTimbau, this model is adapted for the Portuguese language, making it a valuable tool for text processing. BERTimbau offers several advantages, including its ability to understand complex contexts and excellent performance in a broad range of NLP tasks. Moreover, it has proven particularly relevant for the Portuguese-speaking community, filling an important gap in text processing in this language.

### 4.2.5 PTT-5

The PTT-5 (Portuguese, Tagalog, Turkish, Tamil, and Telugu) is a textual representation that stands out for its multilingual approach. In an increasingly globalized world, the ability to process text in multiple languages is essential, and the PTT-5 aims to address this need. The PTT-5 is a textual representation that stands out for its multilingual approach (Carmo et al., 2020), making it suitable for the context of the Portuguese language. In ad-

dition, PTT-5 is based on a text-to-text approach powered by the T5 model for text-to-text representation, thus enabling the text representation based on a transformer architecture for text summarization.

### 4.3 Standard Classifiers

In the context of this study, we used the conventional classifiers for detecting Portuguese fake news articles:

1. Logistic Regression (LR) (Cox, 1972);

2. Multilayer Perceptron (MLP) (Bishop, 1995);

3. Naive Bayes (NB) (Rish, 2001);

4. Optimum-Path Forest (OPF) (Papa et al., 2009, 2012);

5. Random Forest (RF) (Breiman, 2001);

6. Support Vector Machine (SVM (Cortes and Vapnik, 1995).

### 4.4 Deep Classifiers

The experiments were performed using the following deep learning models:

1. Convolutional Neural Network (CNN) (Le-Cun et al., 1998);

2. Long Short-Term Memory (LSTM) (Hochreiter and Schmidhuber, 1997), and Bidirectional Long Short-Term Memory (BiLSTM) (Graves and Schmidhuber, 2005);

3. Gated Recurrent Unit (GRU) (Cho et al., 2014), and Bidirectional Gated Recurrent Unit (BiGRU) (Schuster and Paliwal, 1997);

4. Bidirectional Encoder Representations from Transformers (BERT) (Devlin et al., 2018);

5. Text-To-Text Transfer Transformer (T-5) (Raffel et al., 2019).

## 5 Experimental Setup

In terms of the dataset split, we employed a 5-fold cross-validation procedure to achieve the best data balancing between both classes of news. Table 3 presents the sample distributions yielded from this procedure.

Table 3: Details of each experimental setup.

| Set | Types of news | # of samples |
|---|---|---|
| Train | 50% Real and 50% fake | 42, 240 |
| Test | 50% Real and 50% fake | 10, 560 |

For the FastText representation, we adopted the following setup for the hyperparameter values: embedding size equal to 200 dimensions, the maximum number of unique words as 10, 000, the maximum amount of tokens for each sentence equal to 1, 000, and the n-gram is set to the default value of 2. Since BoW and TF-IDF are simpler approaches than the textual representation FastText, we decided to focus on using FastText for the deep learning classifier experiments.

We adopted a Python-inspired implementation of the OPF framework[2] (de Rosa and Papa, 2021) and the Scikit-Learn library (Pedregosa et al., 2011) to perform experiments with the baseline classifiers. In terms of the deep learning models, only BERT and T-5 were performed over Hugging-Face[3] for natural language processing tasks. For CNN and the temporal models, the training process was performed using Adaptive Moment Estimation (Adam) (Kingma and Ba, 2014) as the optimizer and the Binary Cross Entropy as the loss function.

The models' performance is assessed using four validation metrics: (i) precision, (ii) recall, (iii) f1-score, and (iv) accuracy. For each metric, we compute the average values over the 5-fold cross-validation. Discussion regarding the obtained results is presented in the next section.

# 6  Experimental Results

This section covers the experimental setup presented in two major parts: i) the average results for each text representation and classification algorithm involving FakeRecogna 2.0 with extractive summarization and ii) the outcomes from FakeRecogna 2.0 with abstractive summarization. The text size resulting from each method for text summarization is set to a maximum of 1,000 words.

## 6.1  FakeRecogna with extractive summarization

Table 4 shows the average results for each text representation and classification technique, with the best results highlighted in bold.

[2]https://github.com/gugarosa/opfython
[3]https://huggingface.co/

A more in-depth analysis revealed the best performance attained by BoW and the LR classifier when the same joint approach is considered for comparative purposes with the other baseline classifiers. When considering TF-IDF, SVM achieved the best performance in this scenario, followed by LR and MLP, increasing on average 1% of the BoW results. However, when FastText is employed in classical classifiers, the models exhibit inferior performance compared to alternative representations. The overall results dropped in performance, but the MLP model was stable at an average accuracy of 90%. The results exhibit remarkable performance, even using standard natural language processing techniques like BoW and TF-IDF. The results involving deep classifiers showed increased performance using the FastText representation. In this case, the best classifier was BiGRU, while the BERT and T-5 classifiers exceeded 98% in accuracy.

## 6.2  FakeRecogna with abstractive summarization

Table 5 presents the average results for each text representation and classification technique considering the abstractive summarization, with the best results highlighted in bold.

The experiments revealed slight improvements by integrating abstractive summarization with deep learning models. This joint strategy improved almost 1% the accuracy of the LSTM, GRU, BiLSTM, BiGRU, and CNN. We considered GRU the best-performing model despite its results being the same as those yielded by BiGRU in the abstractive summarization. This decision was made in terms of the lower parameter counts and shorter training time required by GRU to achieve convergence. Likewise, abstractive summarization attained a marginal increase compared to its counterpart version for the BERT classifier, yielding 98.4% in accuracy in this scenario. The same model attained 98.2% accuracy when extractive summarization was applied. However, no improvement was observed by employing abstractive summarization with the T-5 model.

# 7  Conclusions

In this article, we present FakeRecogna 2.0, a significant update to the original corpus FakeRecogna, aimed at addressing the ever-evolving challenges of detecting fake news in the Brazilian context. By

Table 4: Experimental results with standard classifiers on the FakeRecogna 2.0 corpus with extractive summarization.

| Standard Classifiers | | | | | |
|---|---|---|---|---|---|
| Text Representation | Classifiers | Precision | Recall | F1 | Accuracy |
| BoW | LR | **0.948** | **0.948** | **0.948** | **94.8%** |
| | MLP | 0.940 | 0.940 | 0.940 | 94.0% |
| | NB | 0.890 | 0.890 | 0.890 | 89.1% |
| | OPF | 0.834 | 0.834 | 0.834 | 83.4% |
| | RF | 0.932 | 0.932 | 0.932 | 93.2% |
| | SVM | 0.936 | 0.936 | 0.936 | 93.8% |
| TF-IDF | LR | 0.941 | 0.941 | 0.941 | 94.3% |
| | MLP | 0.939 | 0.939 | 0.939 | 94.2% |
| | NB | 0.900 | 0.900 | 0.900 | 89.4% |
| | OPF | 0.796 | 0.758 | 0.749 | 75.8% |
| | RF | 0.940 | 0.940 | 0.940 | 93.8% |
| | SVM | **0.954** | **0.954** | **0.954** | **95.3%** |
| FastText | LR | 0.866 | 0.866 | 0.866 | 86.6% |
| | MLP | **0.902** | **0.902** | **0.902** | **90.2%** |
| | NB | 0.764 | 0.706 | 0.706 | 70.6% |
| | OPF | 0.784 | 0.784 | 0.782 | 78.4% |
| | RF | 0.888 | 0.888 | 0.887 | 88.7% |
| | SVM | 0.686 | 0.686 | 0.686 | 68.6% |
| Deep Classifiers | | | | | |
| Text Representation | Classifiers | Precision | Recall | F1 | Accuracy |
| FastText | LSTM | 0.957 | 0.957 | 0.957 | 95.7% |
| | GRU | 0.956 | 0.958 | 0.958 | 95.8% |
| | BiLSTM | 0.958 | 0.958 | 0.958 | 95.8% |
| | BiGRU* | **0.958** | **0.959** | **0.958** | **96.0%** |
| | CNN | 0.956 | 0.956 | 0.956 | 95.6% |
| BERTimbau | BERT | **0.985** | 0.979 | **0.982** | **98.2%** |
| PTT-5 | T-5 | 0.980 | **0.980** | 0.980 | 98.0% |

*Best results in terms of recall and accuracy.

expanding the corpus size to nearly $53,000$ news articles, incorporating a variety of categories and news sources, we aim to represent more comprehensive information about the Brazilian scenario in terms of fake news spreading. We hope that FakeRecogna 2.0 will inspire new research and collaborations, and we look forward to seeing how the scientific community will utilize this resource to address the ongoing challenge of fake news in Brazil.

We conducted extensive tests with various classifiers throughout this study, ranging from classical methods to deep learning techniques, allowing us to assess the effectiveness of existing approaches in detecting fake news in the Brazilian context. The results indicate that FakeRecogna 2.0 provides a robust and challenging dataset that can serve as a valuable resource for future research in this context.

Regarding the results of each type of summarization, our initial hypothesis is that extractive summaries would be a more effective alternative than abstractive summaries since they do not hallucinate and are not capable of generating new sentences. On the other hand, considering the ability of abstractive summarizers to generate more concise sentences, we decided to test both methods in the experiments of this work. In practice, the results would not differ much from each other, and, in general, traditional machine learning models performed better with extractive summaries. In contrast, deep learning models performed better with abstractive summaries. In future work, we intend to investigate the reasons for this difference in results and why the models behave differently across different types of summaries.

Table 5: Experimental results with standard classifiers on the FakeRecogna 2.0 corpus with abstrative summarization.

| Standard Classifiers | | | | | |
|---|---|---|---|---|---|
| Text Representation | Classifiers | Precision | Recall | F1 | Accuracy |
| BoW | LR | **0.941** | **0.941** | **0.941** | **94.2%** |
| | MLP | 0.933 | 0.933 | 0.933 | 93.3% |
| | NB | 0.896 | 0.896 | 0.896 | 89.4% |
| | OPF | 0.834 | 0.834 | 0.896 | 89.1% |
| | RF | 0.920 | 0.920 | 0.920 | 91.9% |
| | SVM | 0.932 | 0.932 | 0.932 | 93.4% |
| TF-IDF | LR | 0.939 | 0.939 | 0.939 | 93.9% |
| | MLP | 0.933 | 0.933 | 0.933 | 93.4% |
| | NB | 0.898 | 0.898 | 0.898 | 89.7% |
| | OPF | 0.540 | 0.540 | 0.540 | 54.0% |
| | RF | 0.922 | 0.922 | 0.922 | 92.3% |
| | SVM | **0.950** | **0.950** | **0.950** | **94.7%** |
| FastText | LR | **0.860** | **0.860** | **0.860** | **86.0%** |
| | MLP | 0.855 | 0.855 | 0.855 | 85.4% |
| | NB | 0.684 | 0.684 | 0.684 | 68.5% |
| | OPF | 0.784 | 0.784 | 0.782 | 78.4% |
| | RF | 0.858 | 0.858 | 0.858 | 85.7% |
| | SVM | 0.733 | 0.733 | 0.733 | 73.0% |
| Deep Classifiers | | | | | |
| Text Representation | Classifiers | Precision | Recall | F1 | Accuracy |
| FastText | LSTM | 0.964 | 0.965 | 0.965 | 96.5% |
| | GRU$^\star$ | **0.965** | **0.965** | **0.965** | **96.5%** |
| | BiLSTM | 0.964 | 0.965 | 0.965 | 96.5% |
| | BiGRU | 0.965 | 0.965 | 0.965 | 96.5% |
| | CNN | 0.963 | 0.963 | 0.963 | 96.3% |
| BERTimbau | BERT | **0.985** | **0.983** | **0.984** | **98.4%** |
| PTT-5 | T-5 | 0.980 | 0.980 | 0.980 | 98.0% |

$^\star$Best results in terms of the lower count for the network parameters.

# References

H Allcott and M Gentzkow. 2017. Social media and fake news in the 2016 election. *Journal of Economic Perspectives*, 31:211–236.

Christopher M Bishop. 1995. *Neural networks for pattern recognition*. Oxford university press.

Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. 2017. Enriching Word Vectors with Subword Information. *Transactions of the Association for Computational Linguistics*, 5:135–146.

Leo Breiman. 2001. Random forests. *Machine learning*, 45(1):5–32.

Diedre Carmo, Marcos Piau, Israel Campiotti, Rodrigo Nogueira, and Roberto de Alencar Lotufo. 2020. PTT5: Pretraining and validating the T5 model on Brazilian Portuguese data. *ArXiv*, abs/2008.09144.

Anderson Cordeiro Charles, Livia Ruback, and Jonice Oliveira. 2022. Fakepedia corpus: A flexible fake news corpus in portuguese. In *International Conference on Computational Processing of the Portuguese Language*, pages 37–45. Springer.

Kyunghyun Cho, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. 2014. Learning phrase representations using RNN encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078*.

Corinna Cortes and Vladimir Vapnik. 1995. Support-vector networks. *Machine learning*, 20(3):273–297.

David R Cox. 1972. Regression models and life-tables. *Journal of the Royal Statistical Society: Series B (Methodological)*, 34(2):187–220.

Gustavo H de Rosa and João P Papa. 2021. Opfython: A python implementation for optimum-path forest. *Software Impacts*, 9:100113.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. BERT: Bidirectional En-

coder Representations from Transformers. *arXiv preprint arXiv:1810.04805*.

Patricia Takako Endo, Guto Leoni Santos, Maria Eduarda de Lima Xavier, Gleyson Rhuan Nascimento Campos, Luciana Conceição de Lima, Ivanovitch Silva, Antonia Egli, and Theo Lynn. 2022. Illusion of Truth: Analysing and classifying COVID-19 fake news in Brazilian Portuguese language. *Big Data and Cognitive Computing*, 6(2):36.

Soheil Esmaeilzadeh, Gao Xian Peh, and Angela Xu. 2019. Neural Abstractive Text Summarization and Fake News Detection.

Pedro Faustini and Thiago Covões. 2019. Fake News Detection Using One-Class Classification. In *2019 8th Brazilian Conference on Intelligent Systems (BRACIS)*, pages 592–597.

Gabriel L Garcia, Luis CS Afonso, and João P Papa. 2022. Fakerecogna: A new brazilian corpus for fake news detection. In *International Conference on Computational Processing of the Portuguese Language*, pages 57–67. Springer.

Gabriel Lino Garcia, Luis CS Afonso, Leandro A Passos, Danilo S Jodas, Kelton AP da Costa, and João P Papa. 2023. FakeRecogna Anomaly: Fake News Detection in a New Brazilian Corpus. In *VISIGRAPP (4: VISAPP)*, pages 830–837.

Gabriel Lino Garcia, Pedro Henrique Paiola, Luis Henrique Morelli, Giovani Candido, Arnaldo Cândido Júnior, Danilo Samuel Jodas, Luis C. S. Afonso, Ivan Rizzo Guilherme, Bruno Elias Penteado, and João Paulo Papa. 2024. Introducing bode: A fine-tuned large language model for portuguese prompt-based task.

Alex Graves and Jürgen Schmidhuber. 2005. Framewise phoneme classification with bidirectional LSTM and other neural network architectures. In *Neural Networks, 2005. IJCNN'05. Proceedings. 2005 IEEE International Joint Conference on*, pages 2047–2052. IEEE.

Philipp Hartl and Udo Kruschwitz. 2022. Applying Automatic Text Summarization for Fake News Detection.

Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation*, 9(8):1735–1780.

Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models.

Marcio Inácio and Thiago Pardo. 2021. Semantic-Based Opinion Summarization. In *Proceedings of the International Conference on Recent Advances in Natural Language Processing (RANLP 2021)*, pages 619–628, Held Online. INCOMA Ltd.

Ankit Kesarwani, Sudakar Singh Chauhan, and Anil Ramachandran Nair. 2020. Fake News Detection on Social Media using K-Nearest Neighbor Classifier. In *2020 International Conference on Advances in Computing and Communication Engineering (ICACCE)*, pages 1–4.

Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.

Celio Larcher, Marcos Piau, Paulo Finardi, Pedro Gengo, Piero Esposito, and Vinicius Caridá. 2023. Cabrita: closing the gap for foreign languages.

Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. 1998. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324.

Roque Enrique López Condori and Thiago Alexandre Salgueiro Pardo. 2017. Opinion summarization methods: Comparing and extending extractive and abstractive approaches. *Expert Systems with Applications*, 78:124–134.

Derek Miller. 2019. Leveraging BERT for Extractive Text Summarization on Lectures. *CoRR*, abs/1906.04165.

Shubha Mishra, Piyush Shukla, and Ratish Agarwal. 2022. Analyzing machine learning enabled fake news detection techniques for diversified datasets. *Wireless Communications and Mobile Computing*, 2022.

Rafael A Monteiro, Roney LS Santos, Thiago AS Pardo, Tiago A de Almeida, Evandro ES Ruiz, and Oto A Vale. 2018. Contributions to the study of fake news in portuguese: New corpus and automatic detection results. In *International Conference on Computational Processing of the Portuguese Language*, pages 324–334. Springer.

Ray Oshikawa, Jing Qian, and William Yang Wang. 2018. A survey on natural language processing for fake news detection. *arXiv preprint arXiv:1811.00770*.

Pedro H. Paiola, Gustavo H. de Rosa, and João P. Papa. 2022. Deep Learning-Based Abstractive Summarization for Brazilian Portuguese Texts. In *BRACIS 2022: Intelligent Systems*, pages 479–493, Cham. Springer International Publishing.

João P Papa, Alexandre X Falcão, Victor Hugo C De Albuquerque, and João Manuel RS Tavares. 2012. Efficient supervised optimum-path forest classification for large datasets. *Pattern Recognition*, 45(1):512–520.

Joao P Papa, Alexandre X Falcao, and Celso TN Suzuki. 2009. Supervised pattern classification based on optimum-path forest. *International Journal of Imaging Systems and Technology*, 19(2):120–131.

Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, et al. 2011. Scikit-learn: Machine learning in Python. *the Journal of machine Learning research*, 12:2825–2830.

Wisam A Qader, Musa M Ameen, and Bilal I Ahmed. 2019. An overview of bag of words; importance, implementation, applications, and challenges. In *2019 international engineering conference (IEC)*, pages 200–204. IEEE.

Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, Peter J. Liu, Jun Zhu, et al. 2019. Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer. *arXiv preprint arXiv:1910.10683*.

Irina Rish. 2001. An empirical study of the naive Bayes classifier. In *IJCAI 2001 workshop on empirical methods in artificial intelligence*, volume 3, pages 41–46.

Natali Ruchansky, Sungyong Seo, and Yan Liu. 2017. Csi: A hybrid deep model for fake news detection. In *Proceedings of the 2017 ACM on Conference on Information and Knowledge Management*, pages 797–806.

Gerard Salton and Christopher Buckley. 1988. Term-weighting approaches in automatic text retrieval. *Information processing & management*, 24(5):513–523.

Mike Schuster and Kuldip K Paliwal. 1997. Bidirectional recurrent neural networks. In *International conference on neural information processing*, pages 1–6. Springer.

Shivangi Singhal, Rajiv Ratn Shah, Tanmoy Chakraborty, Ponnurangam Kumaraguru, and Shin'ichi Satoh. 2019. SpotFake: A Multi-modal Framework for Fake News Detection. In *2019 IEEE Fifth International Conference on Multimedia Big Data (BigMM)*, pages 39–47.

Fábio Souza, Rodrigo Nogueira, and Roberto Lotufo. 2020. BERTimbau: pretrained BERT models for Brazilian Portuguese. In *Intelligent Systems: 9th Brazilian Conference, BRACIS 2020, Rio Grande, Brazil, October 20–23, 2020, Proceedings, Part I 9*, pages 403–417. Springer.

Jackson Wilke da Cruz Souza and Ariani Di Felippo. 2018. Characterization of Temporal Complementarity: Fundamentals for Multi-Document Summarization. *Alfa: Revista de Linguística (São José do Rio Preto)*, 62:125–150.

Xichen Zhang and Ali A Ghorbani. 2020. An overview of online fake news: Characterization, detection, and discussion. *Information Processing & Management*, 57(2):102025.

Xinyi Zhou and Reza Zafarani. 2020. A survey of fake news: Fundamental theories, detection methods, and opportunities. *ACM Computing Surveys (CSUR)*, 53(5):1–40.