

# Classifying COVID-19 Misinformation with BERT: A Deep Learning Approach

Boyu Gong,<sup>1</sup> Yichi Guo,<sup>1</sup> Jinlong Hu<sup>1</sup>

<sup>1</sup> Xiamen University, Xiamen, China  
{firstauthor, secondauthor, thirdauthor}@xmu.edu.cn

## Abstract

The spread of misinformation during the COVID-19 pandemic has had a significant impact on public health efforts and social stability. The creation of misinformation is often driven by various interests, including financial gains and political motives. Some false claims aim to attract public attention to generate advertising revenue, while others manipulate public opinion through provocative content. The automation of misinformation detection faces several challenges, one of which is the blending of false information with real facts, sometimes even leveraging the credibility of authoritative figures or organizations. Moreover, the lack of timely and authoritative data, particularly in relation to current events like the pandemic, further complicates the process of distinguishing truth from falsehood. To address these challenges, this paper proposes a BERT-based deep learning approach for the automatic classification of COVID-19-related misinformation. By fine-tuning the BERT model on a labeled dataset of rumors and facts, we are able to effectively distinguish between factual information and misinformation. Experimental results show that our BERT model achieves an accuracy of 89% in classifying COVID-19 misinformation, significantly outperforming traditional machine learning models such as Support Vector Machines (SVM)[4] and Naive Bayes. This work highlights the tremendous potential of deep learning models in combating misinformation and underscores the critical role of automated systems in promoting accurate information and reducing misleading content during global health crises.

## Introduction

The COVID-19 pandemic has not only caused a global health crisis but has also triggered an unprecedented spread of misinformation. False claims about the virus, its transmission, prevention methods, and treatments have spread rapidly through social media and other digital platforms, leading to confusion, panic, and, in some cases, harmful behavior. Misinformation about COVID-19 has undermined public health efforts, delayed effective medical responses, and even contributed to the politicization of health measures. As a result, detecting and classifying COVID-19-related misinformation has become a critical task in mitigating the negative impact of false information during the pandemic.

Copyright © 2025, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

Traditional approaches to misinformation detection have faced several challenges, including the difficulty of distinguishing between real and fake information, especially when false claims are subtly interwoven with accurate facts. Additionally, the constant emergence of new rumors during ongoing events like the COVID-19 pandemic makes timely and authoritative data collection and labeling particularly challenging. Recent advances in natural language processing (NLP) and deep learning, however, show significant potential for addressing these issues. Despite this, many existing models rely on shallow features or are limited by the quality of training data, leading to less effective misinformation detection.

To address these challenges, we propose a deep learning-based approach using Bidirectional Encoder Representations from Transformers (BERT), a model renowned for its contextual language understanding. By fine-tuning BERT on a dataset of COVID-19-related misinformation, we aim to classify the information as either factual or false with high accuracy. BERT's ability[1] to capture nuanced relationships between words in context makes it an ideal model for handling the complexities inherent in misinformation detection. In this paper, we make the following contributions:

- We propose a novel method based on BERT for classifying COVID-19 misinformation, demonstrating its ability to effectively distinguish between true and false claims related to the pandemic.
- We conduct extensive experiments on a real-world dataset of COVID-19 misinformation, showing that our method significantly outperforms traditional machine learning models such as Support Vector Machines (SVM) and Naive Bayes in terms of classification accuracy.
- This work is among the first to apply BERT to the automated classification of COVID-19 misinformation, showcasing the potential of deep learning techniques in addressing public health crises by enhancing information accuracy and combating false narratives in real-time.
- We analyze the performance of the BERT model in misinformation classification, discuss the interpretability challenges of deep learning models, and propose strategies to enhance model transparency and trustworthiness.

Through these contributions, we aim to demonstrate that

deep learning models, particularly BERT, can play a crucial role in addressing the challenges of misinformation detection and help combat the spread of false information, especially during global health crises.

## Related work

### Traditional Misinformation Detection Methods

Early approaches to misinformation detection relied heavily on rule-based systems and traditional machine learning algorithms. Techniques such as Support Vector Machines (SVM), Naive Bayes, and decision trees were employed to classify news articles or posts as true or false. These methods, however, often struggled to effectively capture the subtle linguistic and contextual cues that differentiate misinformation from genuine information. Furthermore, they heavily relied on manually crafted features, which limited their scalability and generalization to new datasets.

### Deep Learning Approaches to Misinformation Detection

With the advent of deep learning, significant improvements have been made in the field of misinformation detection. Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs) have been employed to classify text-based misinformation. These models, leveraging powerful word embeddings such as Word2Vec or GloVe, are capable of capturing semantic patterns in textual data. However, they still faced limitations in handling long-range dependencies and contextual relationships, which are critical for distinguishing between fact and misinformation in complex narratives.

### BERT and Transformer Models in Text Classification

Recently, Transformer-based models such as BERT have revolutionized natural language processing tasks, including misinformation detection. BERT, in particular, has achieved state-of-the-art performance on a wide range of text classification tasks. By leveraging bidirectional context through attention mechanisms, BERT can better understand the meaning of words in context, making it particularly suited for tasks like distinguishing between fact and fiction in complex and nuanced statements. Fine-tuning BERT on a specific task has been shown to outperform traditional machine learning models and other deep learning architectures such as CNNs and RNNs.

### Misinformation Detection during the COVID-19 Pandemic

The COVID-19 pandemic has been a fertile ground for the spread of misinformation, ranging from false claims about the virus's origin to unproven treatments. Several studies have attempted to address this issue by applying natural language processing techniques to identify and classify COVID-19-related misinformation. These studies have employed deep learning models, particularly BERT and other

Transformer-based architectures, to detect false information in social media posts, news articles, and forum discussions. Despite challenges such as imbalanced data and the rapidly evolving nature of misinformation, deep learning approaches—especially BERT—have shown significant potential in combating the spread of misinformation during the pandemic.

- You must use the 2025 AAAI Press L<sup>A</sup>T<sub>E</sub>X style file and the aai25.bst bibliography style files, which are located in the 2025 AAAI Author Kit (aai25.sty, aai25.bst).
- You must complete, sign, and return by the deadline the AAAI copyright form (unless directed by AAAI Press to use the AAAI Distribution License instead).
- You must read and format your paper source and PDF according to the formatting instructions for authors.
- You must submit your electronic files and abstract using our electronic submission form **on time**.
- You must pay any required page or formatting charges to AAAI Press so that they are received by the deadline.
- You must check your paper before submitting it, ensuring that it compiles without error, and complies with the guidelines found in the AAAI Author Kit.

## Proposed Solution

### Dataset

The experiment used two datasets. One was collected from Tencent's "Jiaozhen" platform, which contains COVID-19 misinformation debunking data. However, this dataset is currently small (only 310 entries at the time of crawling), and no better sources of COVID-19 misinformation data were found. The code for crawling this data has also been packaged together. The competition dataset obtained from the Biendata platform is Weibo data, containing over 38,000 entries. However, the data quality is average, and there are some duplicate entries. The competition dataset from Biendata is shown in Figure 1. Example data is as follows (label 0 represents real information, and label 1 represents false information).

1. 早听说小龙虾好，但是不知道这么能吃！(来自微博) [转帖] 重庆 [1] 日本人为什么不吃小龙虾？为了家人的健康考虑。有日本的客户来，一起的去吃。无论怎样都不吃小龙虾。... (使用新浪长微博工具发布)

0. 地震快讯#中国地震台网正式测定：06月09日09时32分在四川绵阳市九县(北纬33.28度，东经103.75度)发生3.7级地震，震源深度10千米。(来源：中国地震台网)

0. 疫情期间环保局二中队对汉源路沿线点位进行检查。(来源：中国环境网)

1. 这谁在交押金啊？就这没有那么多小钱还干嘛？一下滑吧。(来源：四川地区一小学需要回到十岁孩子的旧衣服和鞋子，洗干净就可以) 地址：甘孜藏族自治州石渠县西区长沙路马乡小学。邮编：627350 校长：达让18923491809 大家可以5.2的名义寄去，用这种方式给他们庆祝7周年也不错。

0. #周杰伦歌迷会联合声明# 周杰伦演唱会主办方以及歌迷！周杰伦中文网JayCn周杰伦国际后援会Jay2u杰威尔音乐官方微博地震震后周杰伦世界巡回演唱会

0. [男子动物园发用石块砸狗视频：难得一遇] 3月27日，扬州网友在微博上曝光一段男子用石块砸动物园狗的视频。其行径被其他游客阻止后，男子竟说“难得来一趟”。男子动物园...

1. 前转的女孩出果果在泉州万达广场附近走丢麻烦帮我发下联系电话18852405378/18151887476。顺手一转功德无量！

1. 前看到的科普贴：家里的牙膏要严选！牙膏时注意牙膏管底部底部的颜色条，今天才知道，原来颜色条有含义！分四种：绿、蓝、红、黑。绿色：纯天然；蓝色：天然+药物；红色：天然+化学成分；黑色：纯化学成分ps：我用了那么久的两支牙膏吃了不少毒的。

1. 中央电视台《焦点访谈》已经播出，可口可乐承认旗下(黑松糖)含有美国禁用农药“多菌灵”，多菌灵可致脑癌、肝脏癌等癌症。包括香港正在销售的(黑松糖)，香港食环署正在了解此事。专家指出，(多菌灵)跟其他农药一样，对脑部影响最大，可引起癫痫抽搐，并会导致致癌。不要给孩子用这种农药。

0. 地震快讯#中国地震台网自动测定：06月11日08时09分在四川绵阳市九县(北纬33.18度，东经103.68度)发生3.9级左右地震，震源深度10千米。(来源：中国地震台网)

1. 宁波发现国内第一起埃博拉，此疾病致死率90%。大家必须提醒孩子和家人用肥皂洗手，不吃路边摊和露天食物，买的成品食物必须开火食用，防范在先！切记！此次埃博拉可能发展为比SARS更可怕的瘟疫。中德人口密集，防范意识差，很令人担忧。不要恐慌，转发！2此环境

1. 吃榴莲后，喝可口可乐，喝过甜腻！又一游客，客死泰国乡！一位中国游客在泰国旅游的时候，吃了很多榴莲之后喝了可乐，致血液中糖分飙升，结果引发心脏病猝死。年仅28岁，泰国明确有规定，食用大量榴莲后，8小时之内不能喝可口可乐！心里有数的朋友很多，转发提醒一下！让所有人都知道！

Figure 1: The example data from the Biendata platform is from Weibo.

BERT Model Architecture

BERT (Bidirectional Encoder Representations from Transformers) was introduced by Google in the 2018 paper *“BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding”*. BERT draws on the successful principles of the Transformer model. The Transformer model, proposed by Google’s machine translation team in 2017, eliminates the use of traditional convolutional neural networks (CNNs) and recurrent neural networks (RNNs). Instead, it relies solely on the Attention mechanism to solve machine translation tasks, achieving remarkable results. The architecture of the Transformer[6] is shown in Figure 2.

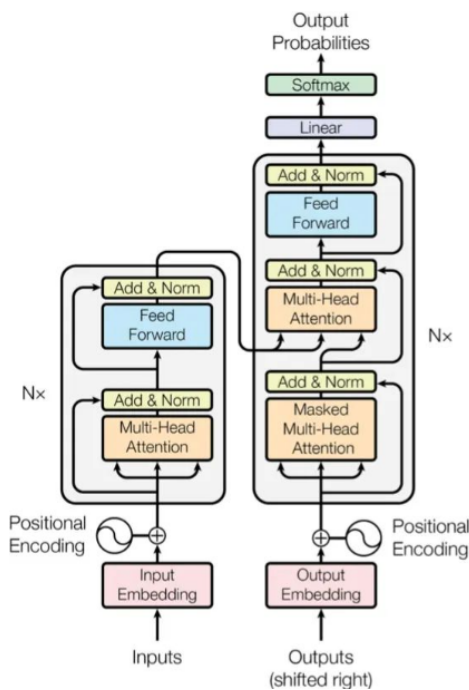


Figure 2: The Transformer consists of two main parts: the Encoder and the Decoder. The left side of the figure represents the Encoder, which is composed of multiple stacked blocks, while the right side represents the Decoder.

The architecture of BERT[2] is primarily based on the Encoder[7] part of the Transformer model, as shown in Figure 4. It consists of multiple stacked encoder layers (Trm). The model input is a sequence of tokens that have undergone word embedding[9] ( $E, E, \dots, E$ ). Each token is mapped to a fixed-dimensional vector representation while retaining positional information. The input data is processed through several Transformer encoder layers, each composed of self-attention mechanisms and feed-forward neural networks. These components enable bidirectional contextual modeling, capturing global dependencies between tokens in the input sequence.

BERT (Ours)

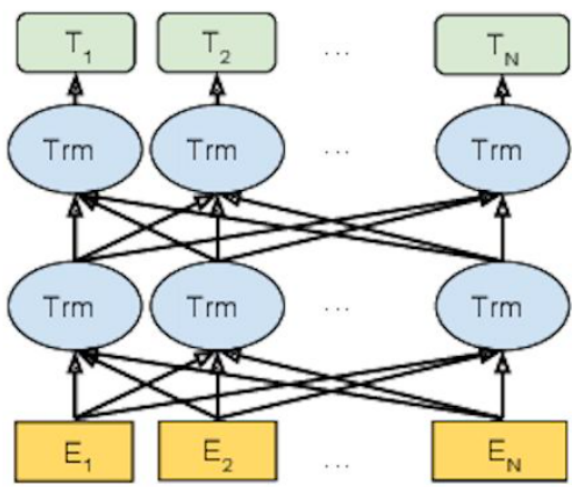


Figure 3: The architecture of BERT, which is based on the Transformer Encoder. Multiple stacked layers enable bidirectional contextual modeling.

In the BERT model, each input token’s representation[5] is refined through multiple Transform[8]r encoder layers to obtain deep contextualized representations ( $T, T, \dots, T$ ). These representations incorporate information from both the left and right contexts.

The core advantage of BERT[3] lies in its *bidirectionality*. Unlike traditional unidirectional language models, BERT captures contextual information simultaneously from left-to-right and right-to-left. This allows for a more comprehensive understanding of the input text’s meaning.

The input to BERT is formed by the combination of three types of embeddings, as shown in Figure 3: “Token Embeddings,” “Segment Embeddings,” and “Position Embeddings.”

Input	[CLS]	my	dog	is	cute	[SEP]	he	likes	play	#ing	[SEP]
Token Embeddings	$E_{[CLS]}$	$E_{my}$	$E_{dog}$	$E_B$	$E_{cute}$	$E_{[SEP]}$	$E_{he}$	$E_{likes}$	$E_{play}$	$E_{++ing}$	$E_{[SEP]}$
Segment Embeddings	$E_A$	$E_A$	$E_A$	$E_A$	$E_A$	$E_A$	$E_B$	$E_B$	$E_B$	$E_B$	$E_B$
Position Embeddings	$E_0$	$E_1$	$E_2$	$E_3$	$E_4$	$E_5$	$E_6$	$E_7$	$E_8$	$E_9$	$E_{10}$

Figure 4: Among the three embeddings that constitute the BERT input, “Token Embeddings” are the word vectors for each token in the input text. “Segment Embeddings” are the sentence-level encoding vectors, which are added to the word vectors of every token in a sentence. “Position Embeddings” encode the positional information of each word within the sentence and are also added to the word vectors of each token.

## Training Process

The training of BERT consists of two main tasks: Masked Language Model (MLM) and Next Sentence Prediction (NSP).

In the Masked Language Model task, the model randomly masks some tokens in the input text and predicts the masked tokens based on their surrounding context. This approach is similar to how CBOW trains Word2Vec, with the key difference being that BERT performs bidirectional modeling and leverages the Transformer architecture to extract features.

The Next Sentence Prediction task is designed to help the model learn the sequential relationships between sentences, specifically determining whether a given sentence B is the next sentence following sentence A. This task is particularly relevant for downstream tasks such as question answering (QA), where understanding inter-sentence relationships is crucial.

Based on the pre-trained BERT[10] model, we can implement tasks such as text classification, QA, and sequence labeling (e.g., word segmentation, named entity recognition), as shown in Figure 5.

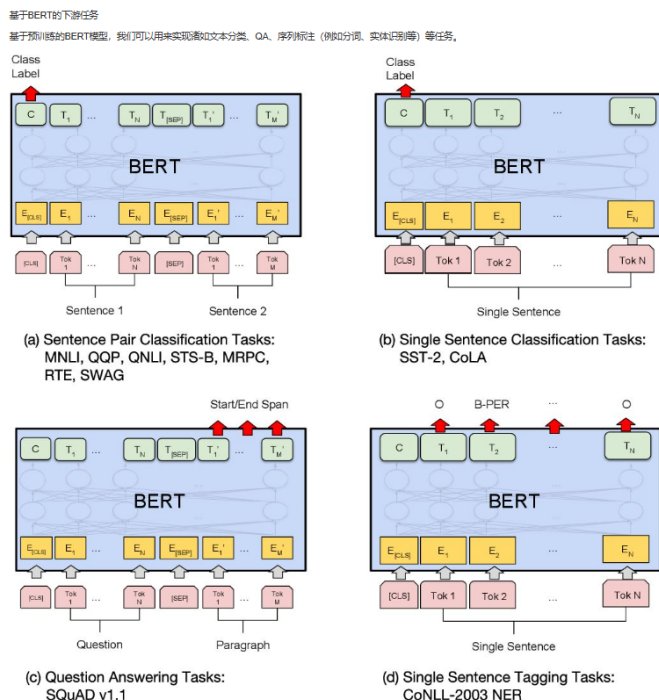


Figure 5: The four subfigures represent: (a) Sentence Pair Classification Tasks (e.g., MNLI, QQP), where the input is a pair of sentences and classification is performed using the [CLS] vector; (b) Single Sentence Classification Tasks (e.g., SST-2, CoLA), where the input is a single sentence and classification also uses the [CLS] vector; (c) Question Answering Tasks (e.g., SQuAD v1.1), where the input is a question and a passage, and the model predicts the start and end positions of the answer; (d) Sequence Labeling Tasks (e.g., CoNLL-2003 NER), where the input is a single sentence, and each token is assigned a label.

## Evaluation Metrics

To evaluate the performance of the model, we utilized the confusion matrix. The confusion matrix provides a detailed breakdown of the model's predictions, offering insights into the number of true positives (TP), true negatives (TN), false positives (FP), and false negatives (FN). This allows for the calculation of key evaluation metrics such as accuracy, precision, recall, and F1-score, which are essential for understanding the model's effectiveness in handling classification tasks.

## Experiments

Google provides the official TensorFlow implementation of the BERT model (<https://github.com/google-research/bert>). The experiments in this study are also based on the official code. The structure of the experimental code is shown in Figure 6.

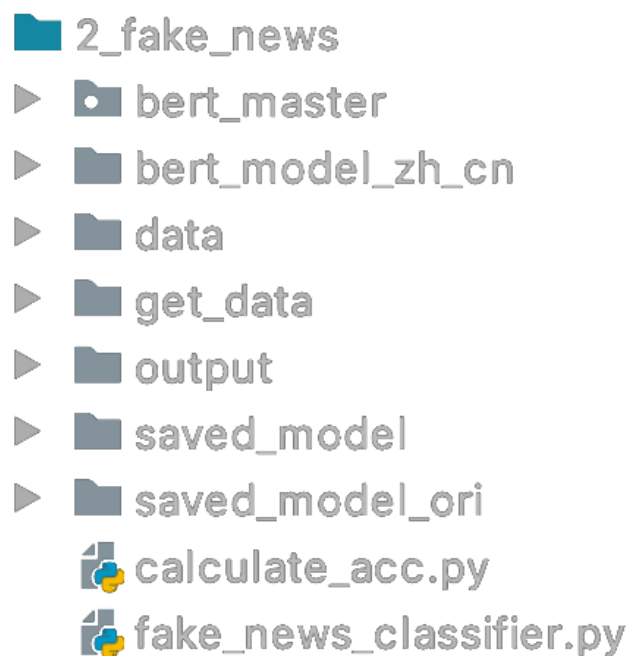


Figure 6: The code structure of this experiment.

The following explains the code structure shown in Figure 6:

- `bert_master`: The official code provided by Google.
- `bert_model_zh_cn`: The pre-trained BERT model for Chinese, provided officially.
- `data`: The data used in the experiment.
- `get_data`: Code for crawling data from Tencent's Jiaozhen platform.
- `output`: Stores the output results during the prediction phase.
- `saved_model`: Contains the fine-tuned BERT model trained on our dataset (using pandemic data, including some Weibo data).

- `saved_model_ori`: Contains the fine-tuned BERT model trained on our dataset (using only Weibo data).
- `calculate_acc.py`: Calculates the confusion matrix and accuracy based on the prediction results.
- `fake_news_classifier.py`: The code implementing classification using BERT.

## Model Training

This study developed a BERT-based binary classifier for rumor detection. The dataset, containing truthful information (label 0) and rumors (label 1), was tokenized using the BERT tokenizer and padded to a length of 128. The model, built on the pre-trained bert-base-uncased with a classification head, was trained using the AdamW optimizer with a learning rate of  $5e-5$ , a batch size of 8, and 500 epochs. Experimental results demonstrated the model's effectiveness in achieving high performance on the test set for rumor detection.

## Result

First, we trained and tested the model using Weibo data to evaluate its performance. The results, including the confusion matrix and various evaluation metrics, indicate that the model performed well, achieving an accuracy of 98%, with other metrics also showing stable and satisfactory results.

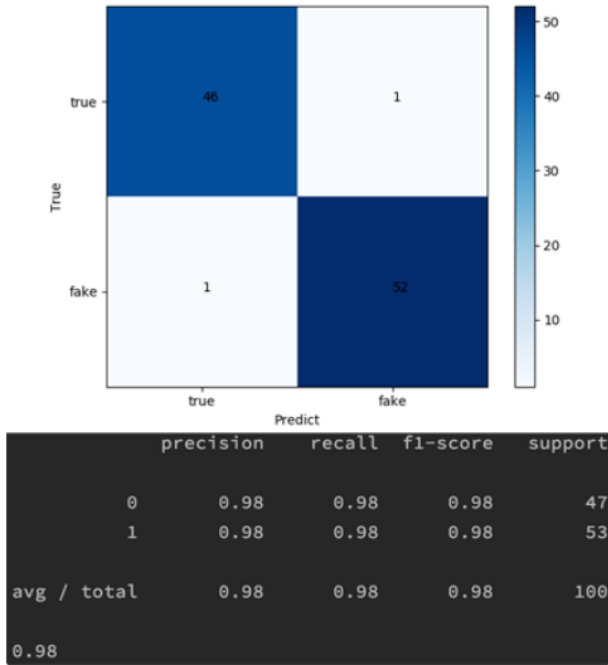
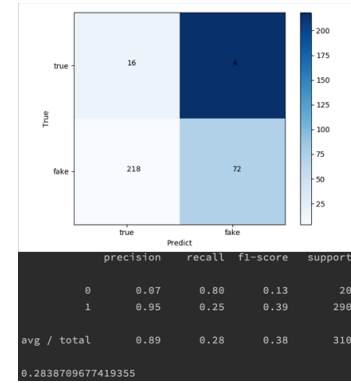
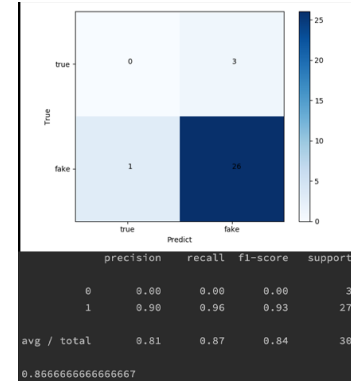


Figure 7: This figure illustrates the model's performance on the classification task between "true" and "fake" labels. The confusion matrix shows that the model made only 2 incorrect predictions out of 100 samples (1 false positive and 1 false negative), achieving an accuracy of 98%. The evaluation metrics indicate that the precision, recall, and F1-score for both classes are all 0.98, demonstrating stable and excellent performance.

The model trained on Weibo data performed poorly when predicting COVID-19 misinformation, achieving an accuracy of only 28%. As shown in Figure 8, the confusion matrix indicates that a large portion of false information (218 out of 290) was misclassified as true information, leading to the low accuracy. This may be due to significant differences in content and features between pandemic-related news and Weibo data, resulting in poor generalization. To address this issue, we directly trained the BERT model using the COVID-19 misinformation dataset. However, the dataset is small and highly imbalanced, with only 310 samples (20 true and 290 false). To balance the training samples, we supplemented the dataset with 120 true information samples from the Weibo data to improve the data distribution.



(a) Training the model with Weibo data to predict COVID-19 misinformation.



(b) The results of adding 120 labeled true information samples to the training data.

Figure 8: Training the model with Weibo data to predict COVID-19 misinformation (a), and training and validating the model with COVID-19 misinformation data (b).

Due to the extremely limited data, only 30 samples were used for the test set. Although the accuracy appears to be 86.66%, a closer analysis of the confusion matrix suggests that this accuracy may be overestimated. Additionally, there were only 3 true information samples in the test set, all of which were misclassified as false information. Given the insufficient data size, this result can only be con-

sidered as a reference.

## Conclusion.

Through the experiments conducted in this study, we observed that employing text classification models for the automated identification of misinformation can yield reasonably good results. This observation was particularly evident in the experiments performed on Weibo data, where the model demonstrated its ability to effectively classify text-based content. These results validate the potential of machine learning models, such as BERT, in tackling misinformation detection tasks to a certain extent.

Nevertheless, significant limitations were encountered when applying the text classification model trained on Weibo data to pandemic-related misinformation. The performance drop was notable, and a closer analysis revealed that the root cause lies in the substantial mismatch between the training dataset (Weibo data) and the domain-specific features inherent in pandemic-related information. Misinformation surrounding a health crisis, such as COVID-19, often possesses unique contextual patterns and linguistic nuances that are not well-represented in general datasets like those from Weibo. This discrepancy highlights a critical challenge: the effectiveness and generalizability of machine learning models heavily depend on the scale, quality, and domain comprehensiveness of the training data.

Addressing this limitation will require further research efforts. Expanding the dataset with domain-specific information is an essential step to ensure the model's robustness and its ability to generalize across diverse types of misinformation. Additionally, strategies such as domain adaptation, transfer learning, and data augmentation could be explored to improve performance in low-resource scenarios where limited labeled data is available, such as pandemic misinformation.

Future research can build on the findings presented in this work, but it is imperative to first address the challenges posed by data insufficiency and domain mismatch. We hope that this study serves as a foundation for further exploration in the automated detection of misinformation, encouraging researchers to focus on building more adaptable and robust models. Collaboration between researchers, organizations, and public platforms will also play a vital role in tackling this issue comprehensively. With joint efforts, more effective tools and systems can be developed to mitigate the spread of misinformation and promote the dissemination of accurate information, particularly during global crises where reliable communication is paramount.

## References

- [1] Devlin, J.; Chang, M.-W.; Lee, K.; and Toutanova, K. 2018. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. *arXiv preprint arXiv:1810.04805*.
- [2] Hao, Y.; Dong, L.; Wei, F.; and Xu, K. 2020. Investigating learning dynamics of BERT fine-tuning. In *Proceedings of the 1st Conference of the Asia-Pacific*

*Chapter of the Association for Computational Linguistics and the 10th International Joint Conference on Natural Language Processing*, 87–92.

- [3] Hinton, G. 2022. The forward-forward algorithm: Some preliminary investigations. *arXiv preprint arXiv:2212.13345*.
- [4] Kamineni, G.; Sai, K. A.; and Rao, G. S. N. 2023. Resume Classification using Support Vector Machine. In *2023 3rd International Conference on Pervasive Computing and Social Networking (ICPCSN)*, 91–96. IEEE.
- [5] Sun, C.; Qiu, X.; Xu, Y.; and Huang, X. 2019. How to fine-tune bert for text classification? In *Chinese computational linguistics: 18th China national conference, CCL 2019, Kunming, China, October 18–20, 2019, proceedings 18*, 194–206. Springer.
- [6] Vaswani, A. 2017. Attention is all you need. *Advances in Neural Information Processing Systems*.
- [7] Wen, Q.; Zhou, T.; Zhang, C.; Chen, W.; Ma, Z.; Yan, J.; and Sun, L. 2022. Transformers in time series: A survey. *arXiv preprint arXiv:2202.07125*.
- [8] Xu, P.; Zhu, X.; and Clifton, D. A. 2023. Multimodal learning with transformers: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(10): 12113–12132.
- [9] Zhang, S.; Jiang, H.; Xiong, S.; Wei, S.; and Dai, L.-R. 2016. Compact Feedforward Sequential Memory Networks for Large Vocabulary Continuous Speech Recognition. In *Interspeech*, 3389–3393.
- [10] Zhang, T.; Wu, F.; Katiyar, A.; Weinberger, K. Q.; and Artzi, Y. 2020. Revisiting few-sample BERT fine-tuning. *arXiv preprint arXiv:2006.05987*.