# MUSER: A MUlti-Step Evidence Retrieval Enhancement Framework for Fake News Detection

Hao Liao
Shenzhen University
Shenzhen, China
haoliao@szu.edu.cn

Jiahao Peng
Shenzhen University
Shenzhen, China
2070276145@email.szu.edu.cn

Zhanyi Huang
Shenzhen University
Shenzhen, China
huangzhanyi2020@email.szu.edu.cn

Wei Zhang
Shenzhen University
Shenzhen, China
2210275010@email.szu.edu.cn

Guanghua Li
Shenzhen University
Shenzhen, China
2210275050@email.szu.edu.cn

Kai Shu*
Illinois Institute of Technology
Chicago, USA
kshu@iit.edu

Xing Xie*
Microsoft Research Asia
Beijing, China
xingx@microsoft.com

## ABSTRACT

The ease of spreading false information online enables individuals with malicious intent to manipulate public opinion and destabilize social stability. Recently, fake news detection based on evidence retrieval has gained popularity in an effort to identify fake news reliably and reduce its impact. Evidence retrieval-based methods can improve the reliability of fake news detection by computing the textual consistency between the evidence and the claim in the news. In this paper, we propose a framework for fake news detection based on **MU**lti-**S**tep **E**vidence **R**etrieval enhancement (MUSER), which simulates the steps of human beings in the process of reading news, summarizing, consulting materials, and inferring whether the news is true or fake. Our model can explicitly model dependencies among multiple pieces of evidence, and perform multi-step associations for the evidence required for news verification through multi-step retrieval. In addition, our model is able to automatically collect existing evidence through paragraph retrieval and key evidence selection, which can save the tedious process of manual evidence collection. We conducted extensive experiments on real-world datasets in different languages, and the results demonstrate that our proposed model outperforms state-of-the-art baseline methods for detecting fake news by at least 3% in F1-Macro and 4% in F1-Micro. Furthermore, it provides interpretable evidence for end users.

## CCS CONCEPTS

• **Computing methodologies** → **Natural language processing**;
• **Information systems** → *Data mining*.

## KEYWORDS

Evidence-based Fake News Detection; Multi-step Retrieval; Explainability

## 1 INTRODUCTION

The explosive growth of fake news has exerted serious negative consequences across society affecting areas such as politics, the
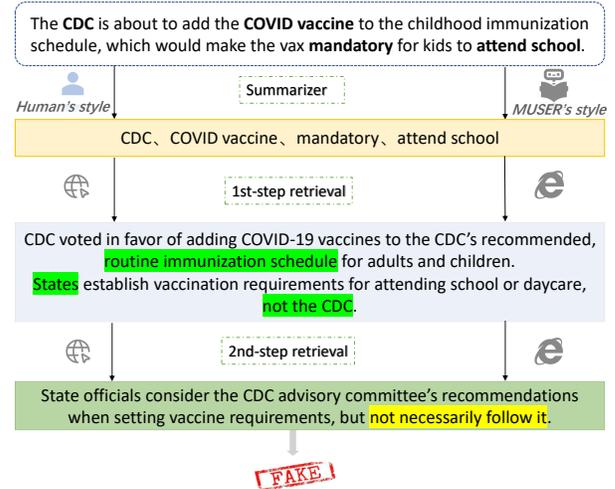


Figure 1: A motivating example of MUSER model. Our model simulates a human evaluating news through three steps: (1) Summarization of the key information, (2) Retrieval and evaluation of relevant evidence: the model assesses the sufficiency and quality of the evidence, determining if additional inquiries are necessary, (3) Conclusion regarding the truthfulness of the news based on the gathered evidence.

economy, and public health [1]. This phenomenon is characterized by the dissemination of sensationalized and alarmist content, which caters to the mindset of netizens and is easily exploited by the "headline party" [2]. To garner more attention, individuals are prone to share news articles or retweet tweets featuring captivating headlines without conducting a diligent evaluation. Consequently, this has facilitated the rapid dissemination of fake news through social media platforms, outpacing the circulation of authentic news. [3]. An overwhelming amount of fake news on social media has made it difficult for individuals to distinguish truth from falsehood, thereby

posing a substantial threat to societal stability [4, 5]. In light of these challenges, the emerging automated fake news detection has drawn widespread attention.

Generally, the detrimental effects of fake news tend to exacerbate over time. To mitigate the ramifications of fake news dissemination, it is important to promptly identify them on social platforms. Meanwhile, fake news detection can help netizens improve their ability to distinguish between true and fake news, thereby fostering the well-being and sustainability of social networks. Various efforts have been made by websites and social media platforms to combat fake news, such as Meta's encouragement for users to report untrustworthy posts and Sina Weibo's provision of a channel for debunking rumors [6]. Besides, fact-checking sites like FactCheck[1], PolitiFact[2] and Full Fact[3] have also begun to hire professionals to conduct fact-checking. However, the diversity and complexity of the increasing volume of news data make manual verification a time-consuming and unscalable process.

To tackle this problem, data mining and machine learning techniques were introduced to detect fake news [7, 8]. Intuitively, the task of fake news detection can be framed as a binary classification problem. These methods commonly employ supervised learning techniques, utilizing textual features such as sentence semantics and news entities, to distinguish between genuine and fabricated news articles [9–11]. Though effective, these content-based methods exhibit some limitations, as fake news often resembles real news in textual features and lacks important information, such as social context [12]. To overcome these limitations, multi-modal fake news detection frameworks have been proposed, which consider social context by analyzing news propagation patterns on social media, such as retweet relationship networks [13, 14], and user-friend relationships [15, 16]. Fake news can spread rapidly and become difficult to control once it has reached a wide audience [17]. Methods based on social context information require a substantial amount of social context information, which may not curb the dissemination of fake news in a timely manner. In addition to the temporal delay issue of detection, methods based on social context face the challenge of user privacy preservation. Therefore, recent research endeavors have increasingly focused on evidence-based verification techniques as a means to detect fake news. These methods perceive fake news detection as an inferential process, wherein external evidence is employed to scrutinize the veracity of the claims presented in news articles. By extracting and incorporating relevant information from the given evidence for claim verification, these methods aim to improve the interpretability of fake news detection. Notably, recent studies have showcased promising outcomes regarding the effectiveness of these approaches. [18–21].

Despite substantial advancements over the years, fake news detection still confronts numerous challenges. Evidence-based detection methods suffer from the assumption that evidence is easily accessible, ignoring the large amount of manual effort required for evidence collection. Furthermore, prior work has inadequately explored complex, long-range semantic dependencies in evidence, neglecting the intricate relationships between information.

Inspired by brain science [22], we propose a fake news inference framework **MU**lti-**S**tep **E**vidence **R**etrieval (MUSER). The cognitive processes involved in human news consumption typically involve three steps [23] as shown in Figure 1: First, a summary of the key findings or claims in the text is made. Second, supporting evidence for the claims is located and evaluated for quality, which may include sources such as website data, official experiments, or research. Finally, conclusions are drawn based on the evaluated evidence. By following these steps, it is possible to ascertain the sources of information, the evidence used, evidence quality, and limitations, thus helping readers to make informed judgments about the validity of the information. MUSER[4] automatically retrieves existing evidence from Wikipedia through paragraph retrieval and key evidence selection, eliminating the need for manual evidence collection. Pieces of evidence needed for news verification are correlated through multi-step retrieval. Furthermore, our model can perform early detection without relying on social context information and provides reasons for the authenticity of the news through retrieved evidence. Although social media can provide external information for early fake news detection, there are two drawbacks - privacy concerns related to user comments and the presence of noisy information among user posts. Our main contributions can be summarized as follows:

- We propose an automatic fact-checking framework for fake news detection that is based on multi-step evidence retrieval. Our framework can explicitly model dependencies among multiple pieces of evidence and retrieves the evidence necessary for news verification through multi-step retrieval. The framework simulates the searching behavior of people when verifying news content on the Internet, making it possible to narrow the gap between computers and human experts in fake news detection.
- The implementation of our proposed model includes three core modules: text summarization, multi-step retrieval, and text reasoning. In the multi-step retrieval module, we employ the method of key evidence selection to control the number of hops, realizing adaptive retrieval step control.
- We conduct extensive experiments on three real-world datasets , and the results demonstrate the effectiveness of our model in terms of improved interpretability and good performance when compared with state-of-the-art models.

## 2 RELATED WORK

### 2.1 Fake News Detection

In recent years, researchers have collaborated with the news ecosystem to better define and characterize fake news through news content and social feedback from web users. We briefly introduce related work from the following aspects: 1) content-based; 2) social context-based; 3) evidence-based.

**Content-based:** Content-based methods detect fake news by exploiting news text, writing style, or external knowledge about news entities. Some works detect fake news by extracting news text features, e.g., n-gram distribution and/or utilize Linguistic Inquiry and Word Count (LIWC) [24] features and sentence relationships based on Rhetorical Structure Theory (RST) [25]. The

---

[1]https://www.factcheck.org/
[2]https://www.politifact.com/
[3]https://fullfact.org/

[4]Code is available at https://github.com/Complex-data/MUSER/

stylistic feature-based approach distinguishes between real and fake news by capturing the specific writing style and emotion usually present in the textual content of fake news [26]. KAN [27] directly evaluates the authenticity of news by comparing news knowledge with knowledge entities in the knowledge graph. Content-based methods are often used in the early detection of fake news to curb the spread of rumors in the early stages of news dissemination.

**Social context-based:** Social media plays an important role in detecting fake news [28]. It has been used to improve the performance of fake news detection by integrating contextual information on social platforms, such as user characteristics, comments, and positions [29]. Methods based on communication structure rely on the assumption that the communication structure of real news and fake news is quite different [17]. Network structure-based methods extract network features by constructing specific networks, such as user interaction networks, user social structures, participation patterns, and news dissemination networks [16, 30–32].

**Evidence-based:** The semantic similarity (conflict) in the claim-evidence pairs can be used to determine the veracity of the news by searching Wikipedia or fact-checking websites according to the claims in the news. Early research approaches employ sequence models to embed semantics and apply attention mechanisms to capture claim-evidence semantic relations. For example, DeClar[19] uses BiLSTM to embed the semantics of the evidence and calculates the evidence score through the attention interaction mechanism. MAC[20] proposes a multi-level multi-head attention network combining word attention and evidence attention to detect fake news. GET[21] models the claims and evidence as graph-structured data, proposing a unified evidence-graph-based fake news detection method for the first time. Evidence-checking-based methods can reveal false parts of claims, provide users with evidence that news is true or fake, and improve the interpretability of fake news detection. Though effective, the above methods all assume that the evidence declared in the news already exists. However, the collection and arrangement of evidence in the actual process often require a lot of manual operations.

Different from the aforementioned studies, we propose a fake news inference framework augmented by multi-step evidence retrieval. Our model can automatically retrieve existing evidence through Wikipedia, conduct evidence collection, and capture dependencies among evidence through multi-step retrieval.

## 2.2 Retrieval Enhancement

Recent work has shown that retrieving additional information can improve the performance of various downstream tasks [33]. Such tasks include open-domain question answering, fact-checking, fact completion, long-form question answering, Wikipedia article generation, and dialogue. In the classic and simplest form of fact-checking, with claims as query conditions, the $k$ relevant passages $K_S = \{P_1, P_2, \ldots, P_{|K_S|}\}$ needed to verify the claims are obtained. Evidence may be contained within a paragraph, or even within a sentence. Retrieve multiple relevant passages $P_i \in K_S$ by a given query Q, and let the reading comprehension model extract the answer from $P_i$ [34, 35]. These studies all used a single-step search. Contrary to the case of single-step retrieval, evidence for some types of queries cannot be obtained through a single retrieval and

requires multiple iterative queries. The ability to retrieve information with multiple iterations is known in the literature as multi-step retrieval [36]. In multi-step retrieval, evidence may need to be obtained with additional information from a previous search, which might otherwise be interpreted as not being fully relevant to the question and no evidence could be found. We extend the capability of multi-step retrieval to fake news claim verification, querying relevant evidence passages in an iterative retrieval manner.

## 2.3 Natural Language Inference

Given a statement and selected evidence sentences, the task of NLI is to predict their relation labels $y$. The advent of large annotated datasets, such as SNLI [37], CreditAssess [38], FEVER [39], has facilitated the development of many different neural NLI models, facilitating model development for this task [40, 41]. The fact verification task related to natural language inference aims to classify a pair of claims and evidence extracted from Wikipedia into three categories: entailment, contradiction, or neutrality. NSMN [41] uses a connected system of three homogeneous neural semantic matching models that jointly perform document retrieval, sentence selection, and claim verification for fact extraction and verification. Soleimani et al. [42] retrieve and validate claims using a BERT [43] model. With the popularity of graph neural networks, graph-based models are also used for semantic reasoning. EVIN [44] proposes an evidence reasoning network, which extracts core semantic conflicts of claims as evidence to explain verification results. Our work differs from prior research in that we focus on classifying news claims as true or false on a comprehensive examination of relevant evidence.

## 3 PROBLEM STATEMENT

In this section, we first define the problem of fake news detection based on evidence retrieval enhancement. We draw a parallel between the detection of fake news and the process by which human beings verify the authenticity of a news article. First, we read the news content and summarize the key information expressed in the news (content summary), then query the evidence in multiple steps based on the summary (multi-step retrieval), and finally infer the authenticity of the news (i.e., Natural Language Inference). So our problem is defined as follows: the input is only news text $A$, and then the news key statement $C$ is obtained through the text summarization module. Retrieve relevant passages in Wikipedia through $C$ to get $P = \{P_1, P_2, P_3, \ldots\}$, and then perform evidence extraction to obtain $E = \{e_1, e_2, e_3, \ldots\}$. The output is the predicted probability of news authenticity $\hat{y} = f(C, E)$, where $f$ is the natural language inference verification model. And $y \in \{0, 1\}$ represents the binary classification labels. In this context, $y = 0$ corresponds to fake news, while $y = 1$ corresponds to true news.

## 4 THE PROPOSED MODEL

In this section, we propose a framework for fake news detection based on **MU**lti-**S**tep **E**vidence **R**etrieval augmentation(MUSER). Figure 2 illustrates the overall architecture of MUSER. Our model mainly consists of three modules:

**Part 1: Text summarization module:** Simulating the human behavior of reading news and summarizing key news information, the proposed module extracts the key information in the news and
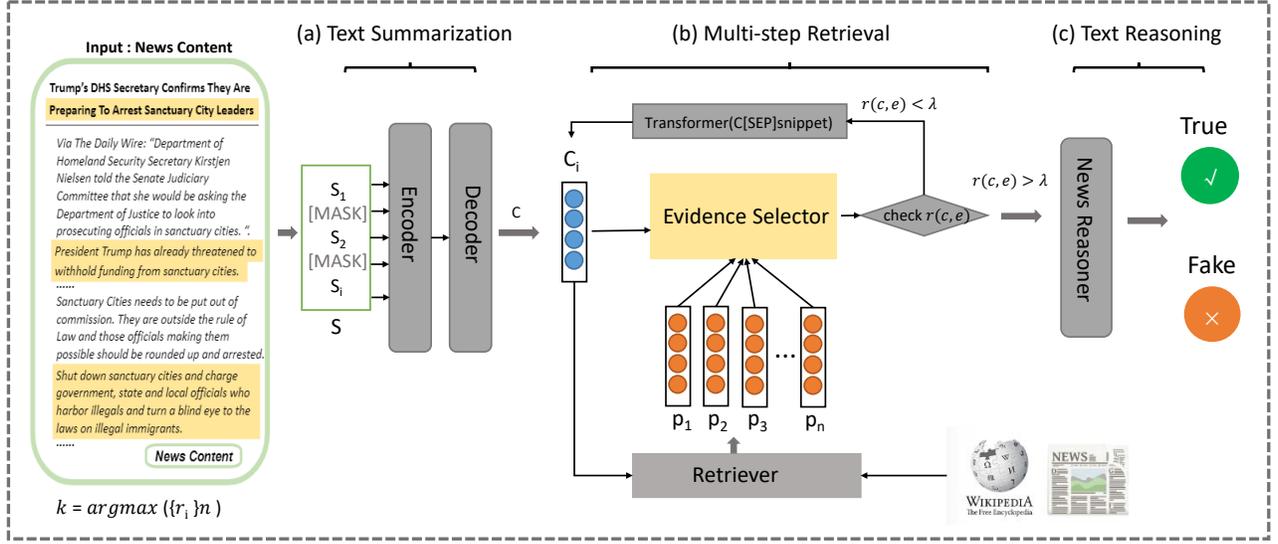
**Figure 2: Our framework unfolds in three steps: (a) Summarization of the initial news text to obtain the key statement $C$, corresponding to the human process of summarizing key information, (b) Evidence finding through multi-step retrieval, corresponding to the human process of querying external relevant information based on the news claim. The retriever sends the first $k$ paragraphs to the evidence selector, which evaluates whether the evidence meets the requirements. The correlation coefficient between $C$ and evidence snippets is represented by $r(c, e)$, and a settable correlation score threshold, $\lambda$, is used to judge the quality of the evidence, and (c) The textual reasoner infers the consistency of evidence and claims, corresponding to the human process of judging news based on evidence.**

filters out the interference of redundant or unimportant information in the news.

**Part 2: Multi-step retrieval module:** Simulating the behavior of humans querying external relevant information in response to news statements, we incorporate a retrieval module into our model. To handle situations where the initially retrieved paragraph may not contain the answer, we adopt a multi-step iterative retrieval method. This process starts by updating the query vector based on the key information and the current query vector. The retriever module then uses this updated query vector for re-retrieval, enabling a deeper exploration of relevant evidence.

**Part 3: Text reasoning module:** Simulating the behavior of humans to judge true or fake news based on the supplementary information queried, this module can extract semantic links between news claims and evidence, and then classify news into two categories: true news and fake news. Through the method of evidence retrieval enhancement, the interpretability of fake news detection is improved, thus mitigating the labor-intensive process of manual evidence extraction.

### 4.1 Text Summarization Module

Naturally, when reading a news article, individuals have a tendency to summarize the key content conveyed within. In order to simulate the ability of humans to summarize news information, we first pre-train a text summarization module. The purpose of this module is to extract the key information in the news and extract the statements worth checking. Although pre-trained language models, such as BERT [43] and UniLM [45], have achieved remarkable results in

NLP scenarios, the word and subword mask language models used in the models may not be suitable for generative text summarization tasks. The reason is that the summarization task requires a coarser-grained semantic understanding, such as sentence and paragraph semantic level understanding, for an effective summary generation.

Inspired by the recent success in masking words and continuous spans, we pre-train a transformer-based encoder-decoder model on a large text corpus for news summarization generation [46]. To leverage a large text corpus for pre-training, we design a sequence-to-sequence self-supervised objective without abstract summarization. We mask sentences from news text and generate an output sequence from the remaining sentences for extracting news summaries. To enhance the relevance of the generated summaries, we select sentences that are deemed important or central to the news.

A piece of news $A$ contains multiple sentences, that is, $A = \{s_i\}_i^N$, where $N$ is the number of sentences. We select the set $S$ of $m$ sentences with the highest scores based on importance. As a proxy for importance, we compute ROUGE1-F1 [47] between the sentence and the rest of the news.

$$r_i = rouge(S \cup s_i, A\backslash\{S \cup s_i\}), \quad \forall i, \ s_i \notin S \quad (1)$$

$A \backslash \{S \cup s_i\}$ represents the remaining sentences, and $S$ is initially an empty set. Then select important sentences according to the importance score $r_i$:

$$k = argmax(\{r_i\}_n) \quad (2)$$

$$S = S \cup s_k \quad (3)$$

The corresponding position of each selected sentence is replaced by a mask token [MASK] to inform the model. Making $m$ selections,

in the end, we select the masked $m$ sentences from the document and concatenate the sentences into a pseudo-summary. the module then generates an output sequence from the remaining sentences, producing the masked sentences. We pre-train the model on the open source news dataset [5] to achieve a better summary generation results. The Mask sentences ratio (MSR) which refers to the ratio of the number of selected gap sentences to the total number of sentences in the document, is an important hyperparameter, similar to the mask rate in other works [46]. A low MSR reduces the difficulty and computational efficiency of pre-training. On the other hand, masking a large number of sentences at high MSR loses the contextual information necessary for guidance generation. In our experiments, we found an MSR of 30% to be effective.

## 4.2 Multi-step Retrieval Module

The purpose of this module is to perform retrieval enhancement based on the key information in the news extracted in the previous step, which is similar to humans looking up data, and finding supplementary information to assist in the identification of true and fake news. Single-step retrieval may lead to insufficient auxiliary information retrieved. Therefore, we adopt a multi-step iterative retrieval method to improve information sufficiency [36]. Through iterative retrieval and supplementation, relevant information can be extracted more comprehensively, so as to better assist in judging the authenticity of news. When implementing this module, it is important to consider how to effectively extract the retrieved key information and how to maintain the sufficiency of information during the multi-step iterative retrieval process.

The multi-step retrieval problem we attempt to address is divided into three steps. In the first step, the news statement $C$ is used to retrieve the relevant paragraph $P$ from the Wikipedia corpus. The second step is to extract evidence from the retrieved long paragraphs and extract the key evidence of the paragraphs. Finally, in the case where no evidence is found in the retrieved paragraphs, the information retrieved in this step is fused with statement $C$ to generate a new statement for the retrieval iteration. The search terminates when evidence is found in the retrieved passages.

**Paragraphs retrieval**: Paragraphs retrieval is the selection of Paragraphs on Wikipedia that are relevant to a given statement. The paragraph retrieval module is based on BERT [43] and creates dense vectors for paragraphs by computing their average token embedding. The relevance of paragraph $p$ to statement $c$ is given by their dot product:

$$r(c, p) = \varphi(c)^T \varphi(p) \tag{4}$$

$\varphi(\cdot)$ is an embedding function used to map paragraphs and statements to a dense vector. Dot product search can use the approximate nearest neighbor index implemented by the FAISS library to improve search efficiency [48]. For the embedding function $\varphi(\cdot)$, we use the average token embedding of the BERT-base language model

[49], which has been fine-tuned on several tasks:

$$\varphi(p) = \frac{1}{p} \sum_{i=1}^{|p|} BERT(p, i) \tag{5}$$

where $BERT(p, i)$ is the embedding of the $i$-th token in paragraph $p$, and $|p|$ is the number of tokens in $p$.

**Key evidence selection**: Key evidence selection is to extract evidence-related key sentences from the retrieved relevant passages. Similar to paragraph retrieval, sentence selection can also be perceived as a semantic matching task, wherein each sentence within a paragraph is compared to a given statement query to identify the most plausible evidence interval. Since the search space has been reduced to a controllable size via the paragraph retrieval in the previous step, we can directly traverse all relevant paragraphs to find key evidence. In this paper, we employ two approaches for key evidence selection: a relevance score-based approach and a context-aware approach.

Relevance score-based selection methods rely on vector representations of statements and sentences in paragraphs. For a given statement $C$, we select sentences $s_i$ from the retrieved relevant passages $P = \{s_1, s_2, \ldots, s_n\}$ whose relevance score $r(c, s_i)$ is greater than a certain threshold $\lambda$ set experimentally. Details on setting lambda values can be found in Appendix A.2.3.

The context-aware sentence selection method uses a BERT-based sequence tagging model. We take as input the concatenation of statement claim $C = \{c_1, c_2, ..., c_k\}$ and passages $P = \{p_1, p_2, ..., p_m\}$ and separate them using special tokens: $[CLS]C[SEP]P[EOS]$. For the output of the model, we adopt the BIO token format, which classifies all irrelevant tokens as O, the first token of an evidence sentence as B evidence, and the remaining tokens of an evidence sentence as I evidence. We train a RoBERTa-large based model [50], minimizing the cross-entropy loss:

$$\mathcal{L}_\theta = -\sum_{i=1}^{N} \sum_{j=1}^{l_i} log(p\theta(y_i^j)) \tag{6}$$

where $N$ is the number of examples in the training batch, $l_i$ is the number of non-padding tokens of the $i$-th example, and $p\theta(y_i^j)$ is the estimated softmax probability of the correct label for the $j$-th token of the $i$-th example. We train this model on Factual-NLI [51] with batch size 64, Adam optimizer, and initial learning rate $5 \times 10^{-5}$ until convergence.

**Multi-step retrieval**: In the process of selecting key evidence, we assess the sufficiency of the evidence's relevance using a threshold $\lambda$. When the evidence is insufficient, we use iterative retrieval to supplement information. To prioritize the most significant fragments in the paragraph, we rank the selected fragments based on their scores. Similar to human behavior of recursively querying external sources like Wikipedia step by step until the desired information is found, only the fragments with the highest scores will be kept. The fragment with the highest score, referred to as the "winner," is then incorporated into the current query $[C[SEP]snippet]$. A reformulated query will be generated by combining the current query with current relevant paragraph information and updating it through a transformer.

$$C_{i+1} = Transformer([C_i[SEP]snippet]) \tag{7}$$

---

[5]Engilsh: http://atp-modelzoo-sh.oss-cn-shanghai.aliyuncs.com/release/tutorials/generation/en_train.tsv

Chinese: http://atp-modelzoo-sh.oss-cn-shanghai.aliyuncs.com/release/tutorials/generation/cn_train.tsv

**Table 1: Statistics of three datasets.**

| Platform | PolitiFact | GossipCop | Weibo |
|----------|-----------|-----------|-------|
| #Real News | 399 | 4,219 | 436 |
| #Fake News | 345 | 3,393 | 311 |
| #Total | 744 | 7,612 | 747 |

The reformulated query is fed back to the retriever, which uses it to reformulate and rank the passages in the corpus. $C_i$ fully interacts with the snippet through the transformer, avoiding information loss during the embedding process. The new query $C_{i+1}$ is again subjected to paragraph retrieval and key evidence selection, achieving the effect of multi-step iterative retrieval. This multi-step iterative approach allows our model to combine the multi-step information needed to validate claims from multiple Wikipedia pages.

## 4.3 Text Reasoning Module

The last step of our model is to infer whether the news is true or false through multi-step retrieved evidence and news statements. This step aligns with human behavior, where individuals gather information from external sources and then evaluate the credibility of the news based on that information. Given a news claim $C$ and relevant evidence $E$ retrieved through a multi-step retrieval process, our text reasoning module performs a logical inference from the evidence to the claim. The textual reasoning model acts as an evaluator to judge whether a statement is logically consistent with the retrieved evidence, thus identifying a pair of claims and related evidence as true or false. Thus, the training task of a text reasoning model can be perceived as a binary classification task, where the goal is to minimize the binary cross-entropy loss function for each news item and its associated evidence. The cross-entropy loss is defined as follows:

$$\mathcal{L}_{CE} = -\frac{1}{N} \sum_{i=1}^{N} y_i log(V(C_i, E_i)) + (1 - y_i)log(1 - V(C_i, E_i)) \quad (8)$$

$N$ is the number of samples in the current batch, $y = 1$ means that claim $C$ and evidence $E$ are logically consistent, and $y = 0$ means that $C$ and $E$ are contradictory. $V$ is a pre-trained language model that can perform discriminative classification tasks, such as BERT [43], ALBERT [52] and RoBERTa [50]. In this work we choose BERT as the discriminator, we concatenate the claim $C$ and the evidence $E$ as the input of the discriminator, the input is [CLS] C [SEP] E [SEP], the batch size $N$ is 64, Adam optimizer and an initial learning rate of $5 \times 10^{-5}$ until convergence.

## 5 EXPERIMENTS

To verify the effectiveness of our proposed model, we conduct extensive experimental studies on three real-world datasets. Four research questions are addressed through comprehensive experimentation:

- RQ1: Is our MUSER model able to achieve improved fake news detection performance compared to previous fake news detection baseline methods?
- RQ2: How does the impact of the number of steps in multi-step retrieval on model performance?

- RQ3: How does each module of the model contribute to improved fake news detection performance?
- RQ4: Is the evidence retrieved by our model meaningful and explainable through multi-step retrieval?

## 5.1 Experimental Setup

*5.1.1 Datasets.* We conduct experiments on three real-world datasets for fake news detection, including two English datasets (PolitiFact and GossipCop) and one Chinese dataset (Weibo). The English datasets PolitiFact and GossipCop are collected through FakeNews-Net [53]. The Weibo dataset is obtained through crawler tools [54]. Their key statistics are shown in Table 1.

**PolitiFact**: Within this dataset, the news articles are divided into two distinct categories: real news and fake news. This classification is determined based on the assessments provided by journalists and experts who review political news on various websites.

**GossipCop**: In this dataset, entertainment news articles with ratings are collected from various media.

**Weibo**: The data in this dataset are hot news topics from the Sina Weibo platform, and news is marked as rumors and non-rumors.

The datasets mentioned above contain both labeled news content and associated social information. However, since our work centers on curbing the initial propagation of fake news, we only utilize the news text without social information. This scenario resembles the situations where fake news detection must be performed before social information becomes available.

*5.1.2 Baselines.* We compare MUSER with several existing methods, including content-based and evidence-based verification, as described below:

**Content-based methods**

- **TextCNN (EMNLP'14)** [55]: TextCNN combines convolutional neural networks and news content, which can automatically extract text features through multiple convolutional hidden layers,
- **TextRNN (ACL'16)** [56]: TextRNN uses LSTM to encode the textual information in the last output of the recurrent neural network.
- **TCNNURG (IJCAI'18))** [57]: TCNNURG utilizes two convolutional neural networks and a conditional variational autoencoder for classification.
- **BERT (NAACL'19)** [43]: BERT uses the Transformer-based architecture to pre-train deep bidirectional representations of unlabeled text.

**Evidence-based methods**

- **DeClarE (EMNLP'18)** [19]: They use BiLSTM to embed the semantics of evidence and compute evidence scores through an attention interaction mechanism.
- **HAN (ACL'19)** [18]: HAN adopts GRU embedding and two modules of topic consistency and semantic entailment based on a sentence-level attention mechanism to simulate claim-evidence interaction.
- **EHIAN (IJCAI'20)** [58]: EHIAN discusses the questionable parts of claims for interpretable claim verification through an evidence-aware hierarchical interactive attention network to explore more plausible evidence semantics.

**Table 2: Performance comparison of Our model w.r.t. baselines. We repeat the experiment 10 times, and average the results. "F1-Ma" and "F1-Mi" denote the metrics F1-Macro and F1-Micro, respectively. "-T" represents "True News as Positive" and "-F" denotes "Fake news as Positive" in the context of computing the precision and recall values. A t-test is performed on five dataset splits, with $P < .05$. The superior outcomes are indicated in bold and statistically significant improvements are denoted by ∗.**

| Method | PolitiFact | | | | | | | |
|--------|-------|-------|-------|-------|-------|-------|-------|-------|
| | F1-Ma | F1-Mi | F1-T | P-T | R-T | F1-F | P-F | R-F |
| TextCNN | 0.601 | 0.602 | 0.608 | 0.641 | 0.579 | 0.594 | 0.564 | 0.615 |
| TextRNN | 0.610 | 0.609 | 0.616 | 0.650 | 0.586 | 0.603 | 0.572 | 0.636 |
| TextURG | 0.621 | 0.619 | 0.637 | 0.651 | 0.624 | 0.601 | 0.587 | 0.617 |
| BERT | 0.597 | 0.598 | 0.608 | 0.619 | 0.599 | 0.586 | 0.577 | 0.597 |
| DeClarE | 0.654 | 0.651 | 0.656 | 0.689 | 0.673 | 0.651 | 0.613 | 0.664 |
| HAN | 0.661 | 0.660 | 0.679 | 0.676 | 0.682 | 0.643 | 0.650 | 0.637 |
| EHIAN | 0.664 | 0.663 | 0.674 | 0.680 | 0.651 | 0.650 | 0.628 | 0.627 |
| MAC | 0.678 | 0.675 | 0.700 | 0.695 | 0.704 | 0.653 | 0.655 | 0.645 |
| GET | 0.694 | 0.692 | 0.725 | 0.712 | 0.770 | 0.669 | 0.720 | 0.665 |
| MUSER | **0.732**∗ | **0.729**∗ | **0.757**∗ | **0.735**∗ | **0.780**∗ | **0.702**∗ | **0.728**∗ | **0.681**∗ |

**Table 3: Performance comparison of on GossipCop.**

| Method | GossipCop | | | | | | | |
|--------|-------|-------|-------|-------|-------|-------|-------|-------|
| | F1-Ma | F1-Mi | F1-T | P-T | R-T | F1-F | P-F | R-F |
| TextCNN | 0.628 | 0.624 | 0.658 | 0.671 | 0.646 | 0.590 | 0.604 | 0.576 |
| TextRNN | 0.629 | 0.628 | 0.636 | 0.667 | 0.609 | 0.620 | 0.591 | 0.651 |
| TextURG | 0.644 | 0.643 | 0.650 | 0.684 | 0.619 | 0.636 | 0.605 | 0.637 |
| BERT | 0.617 | 0.613 | 0.635 | 0.664 | 0.649 | 0.578 | 0.635 | 0.562 |
| DeClarE | 0.660 | 0.657 | 0.686 | 0.677 | 0.694 | 0.629 | 0.638 | 0.619 |
| HAN | 0.702 | 0.700 | 0.722 | 0.721 | 0.716 | 0.678 | 0.676 | 0.680 |
| EHIAN | 0.705 | 0.702 | 0.731 | 0.713 | 0.749 | 0.673 | 0.694 | 0.654 |
| MAC | 0.729 | 0.727 | 0.725 | 0.742 | **0.756** | 0.705 | 0.713 | 0.697 |
| GET | 0.733 | 0.731 | 0.751 | 0.749 | 0.727 | 0.712 | 0.710 | 0.715 |
| MUSER | **0.776**∗ | **0.775**∗ | **0.784**∗ | **0.843**∗ | 0.734 | **0.768**∗ | **0.714**∗ | **0.830**∗ |

- **MAC (ACL'21)** [20]: MAC combines multi-head word-level attention and multi-head document-level attention, which facilitates interpretation for fake news detection at both word-level and evidence-level.
- **GET (WWW'22)** [21]: GET models claims and pieces of evidence as graph-structured data to explore complex semantic structures and reduces information redundancy through the semantic structure refinement layer.

*5.1.3 Implementation Details.* Fake news detection is commonly perceived as a binary classification problem, and the indicators used for model performance evaluation are F1, Precision, Recall, F1-Macro, and F1-Micro [21]. The dataset is partitioned into two sets, with 75% of the data as the training set and the remaining 25% of the data as the test set. The learning rate of the Adam optimizer is uniformly set to $5 \times 10^{-5}$ across all datasets. And the number of training epochs is set to 20 for both our model and the baselines. The hyperparameters for the baselines are configured based on corresponding papers, with key hyperparameters carefully tuned for optimal performance (e.g., learning rate and embedding size). All experiments are conducted on Linux servers equipped with GeForce RTX 3080 GPUs (32GB memory each) using PyTorch 1.8.0. The implementation details are in the appendix and repository.

**Table 4: Performance comparison of on Weibo.**

| Method | Weibo | | | | | | | |
|--------|-------|-------|-------|-------|-------|-------|-------|-------|
| | F1-Ma | F1-Mi | F1-T | P-T | R-T | F1-F | P-F | R-F |
| TextCNN | 0.722 | 0.721 | 0.740 | 0.742 | 0.736 | 0.703 | 0.706 | 0.700 |
| TextRNN | 0.741 | 0.737 | 0.771 | 0.730 | 0.812 | 0.701 | 0.756 | 0.654 |
| TextURG | 0.709 | 0.704 | 0.741 | 0.712 | 0.628 | 0.667 | 0.707 | 0.759 |
| BERT | 0.699 | 0.698 | 0.719 | 0.720 | 0.716 | 0.678 | 0.676 | 0.680 |
| DeClarE | 0.746 | 0.745 | 0.765 | 0.758 | 0.771 | 0.724 | 0.732 | 0.717 |
| HAN | 0.689 | 0.687 | 0.711 | 0.706 | 0.716 | 0.662 | 0.668 | 0.657 |
| EHIAN | 0.753 | 0.752 | 0.770 | 0.768 | 0.772 | 0.734 | 0.754 | 0.731 |
| MAC | 0.734 | 0.732 | 0.709 | 0.722 | 0.697 | 0.755 | 0.745 | 0.766 |
| GET | 0.756 | 0.754 | 0.776 | 0.760 | 0.794 | 0.730 | 0.761 | 0.712 |
| MUSER | **0.804**∗ | **0.802**∗ | **0.824**∗ | **0.812**∗ | **0.837**∗ | **0.791**∗ | **0.806**∗ | **0.778**∗ |

## 5.2 Performance Results (RQ1)

We compare our model, MUSER, to 9 baselines, including 4 content-based methods and 5 evidence-based methods. The results are reported in Tables 2, 3, and 4, and we have the following observations:

Firstly, it is worth noting that evidence-based methods tend to predict more correctly than content-based methods (i.e., the first four methods in the tables), indicating the extra value of incorporating additional evidential information, which can well make up for the insufficiency of news content features alone. The evidence-based methods rely on external evidence to verify the validity of the claims, reducing excessive reliance on textual schemas.

Secondly, in comparison to three recent evidence-based methods (GET, EHIAN, MAC), our proposed MUSER achieves superior results (MUSER > GET > EHIAN > MAC). In particular, MUSER improves the performance by 3% on F1-Macro and F1-Micro compared to the current SOTA baseline GET on the three datasets, which can better reflect the overall detection ability of the model. Furthermore, for more fine-grained evaluation, we computed "True news as Positive" and "Fake news as Positive" separately. MUSER also achieved superior results in F1, Precision, and Recall scores on the three datasets. Accuracy is equivalent to F1-Macro and thus omitted in the evaluation.

Finally, our results demonstrate that MUSER outperforms all baseline methods in fake news detection, as indicated by the positive detection metric. For instance, as far as GossipCop is concerned, the F1-False, Precision-False, and Recall-False values have been increased by 5%, 0.4%, and 11%, respectively. Similar obvious improvements can be observed on other datasets. These results show that our method exhibits a higher degree of accuracy in discerning fake news. Enhanced by multi-step iterative evidence retrieval, our model can extract relevant information, so as to better assist in assessing the veracity of news. Furthermore, extensive experiments are conducted on large public datasets for the detection of fake news. Detailed information can be found in Appendix A.2.1.

## 5.3 Retrieve Steps Comparison (RQ2)

Next, we investigate the performance improvement of the number of retrieval steps in the multi-step retrieval module. The evaluation is conducted using the commonly used F1-Macro and F1-Micro scores on each dataset and results are presented in Figure 3. In order to examine the effectiveness of key evidence selection in the multi-step retrieval process, we remove it and use a fixed number
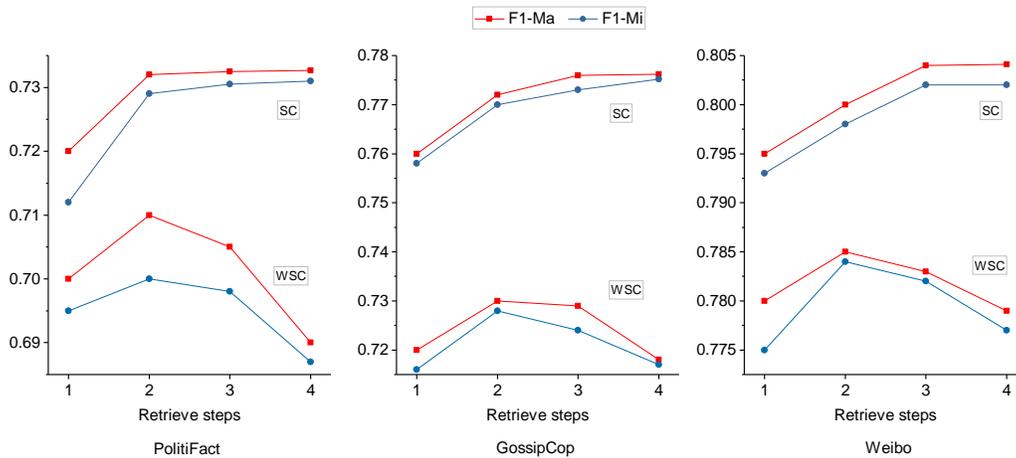
**Figure 3: Results of retrieve step comparison study. The term SC (Step Control) means that the key evidence selection function is activated, while WSC (Without Step Control) means that the key evidence selection function is not included.**
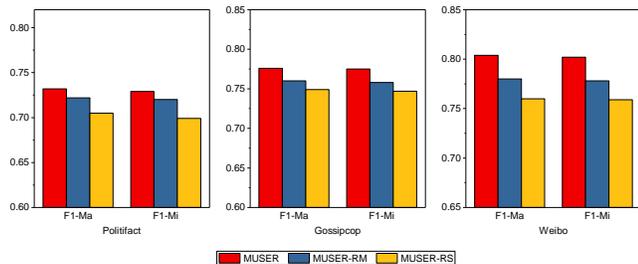


**Figure 4: Results of ablation study. MUSER represents the complete model performance, MUSER-RM represents the removal of the multi-step retrieval module and MUSER-RS represents the removal of the text summary module.**

of retrieval steps to conduct experiments, and then compare it to the model with the key evidence selection function.

Firstly, we can find that in experiments where key evidence selection is not enabled, as the number of retrievals increases, the performance decreases instead. This is because there is no evidence screening for the retrieved paragraphs, which may contain redundant information, leading to a decrease in performance.

Secondly, we observe that enabling key evidence selection results in improved performance compared to the scenario where key evidence selection is not enabled. In the key evidence selection stage, our model determines whether the current retrieval results include key evidence. When key evidence is successfully retrieved, the iterative retrieval process is halted to minimize the interference caused by redundant information. In other words, the selection strategy follows an exploratory approach, where the emphasis is on exploring relevant information first. Importantly, increasing the number of retrieval steps does not result in an increase in redundant information.

The key takeaway from this experiment is that multiple retrieval steps consistently improve performance compared to single-step retrieval. That is, even if relevant evidence passages are not retrieved in the initial step, the retriever will continue in the subsequent iterative retrieval process. The performance of the model reaches its peak around 2 to 3 retrieval steps. Beyond this point, increasing the number of steps does not yield significant benefits and, in fact, leads to a degradation in performance. Interestingly, despite variations in the difficulty level of the datasets, the optimal number of retrieval steps remains consistent.

## 5.4 Ablation Study (RQ3)

In this part, comparative performance experiments are conducted to assess the necessity of each module. As depicted in Figure 4, MUSER outperforms MUSER-RM, proving the critical role of multi-step iterative evidence retrieval. Additionally, the text summarization module is also important. By extracting key statements in the news, the interference of unrelated information is mitigated, thereby achieving more accurate predictions. Furthermore, MUSER performs better than MUSER-RS and MUSER-RM, showing that removing any of them leads to performance degradation, which demonstrates the effectiveness of our main components.

## 5.5 Explainability Study (RQ4)

*5.5.1 Case Study.* In this part, we demonstrate the effectiveness of our model in facilitating a deeper understanding of the multi-step retrieval process. In particular, we present a specific example involving the evaluation of a news story concerning US President Donald Trump's efforts to combat drug-related issues. The news says "Donald Trump marshaled the full power of government to stop deadly drugs, opioids, and fentanyl from coming into our country. As a result, drug overdose deaths declined nationwide for the first time in nearly 30 years." By employing key evidence extraction and conducting a multi-step search for supplementary evidence, MUSER successfully identifies this news as fake. This
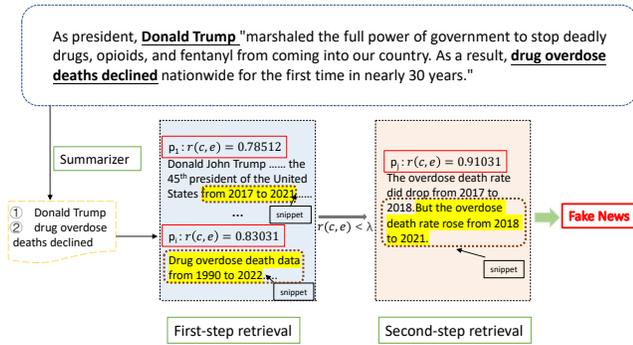
**Figure 5: A verification example generated by MUSER in the Case study. The evidence correlation score $r(c, e)$ obtained by the first step of retrieval is smaller than the threshold $\lambda$ we set. Then proceed to the second step of retrieval to obtain more sufficient evidence.**

particular case serves as a compelling demonstration of MUSER's capability to accurately assess the authenticity of the news.

Specifically, Figure 5 shows the steps of the verification processes. After key information is extracted in the text summarization model, the first step of retrieval is performed, and relevant paragraph data is obtained from the corpus. Evidence extraction identifies information related to Donald Trump and data on drug overdose deaths in the United States. The calculated $r(c, e)$ from the key evidence selection is less than the preset limit value $\lambda$, indicating the necessity for another retrieval step. In the second step, the snippet information retrieved is carried forward and the statement "The overdose death rate did drop from 2017 to 2018. But the overdose death rate rose from 2018 to 2021." is obtained. Finally, the reasoning module judges the news to be fake. Evidence from multi-step retrieval makes it easier for users to understand the judgments made by the model on the authenticity of the news.

*5.5.2 User Study.* In this part, we aim to determine if real-world users are able to accurately assess the veracity of news articles based on the evidence retrieved by MUSER. Specifically, we conduct a user study in which there are 60 news articles randomly selected from PolitiFact, GossipCop, and Weibo, with 10 fake and 10 real news articles from each dataset. We compare the evidence retrieved by MUSER with the evidence obtained by the GET model after refinement by semantic structure and ask 8 participants to score the evidence. For each piece of news, we will give the relevant evidence of MUSER or GET, and then ask the participant to determine whether the news is true or fake based on the given evidence within three minutes. Moreover, participants are asked to give an adjusted confidence score about her/his conclusion according to a 5-point Likert scale. To ensure fairness in our user study experiment, each participant is given the news articles to be judged in a randomized manner and participate in the experiment independently.

Table 5 shows the results of the experiments. By comparing the labels given by different participants, we find that the conclusions drawn by the participants have a high level of consistency with the predicted labels produced by the MUSER model. This indicates

**Table 5: Results of the user study. The agreement measure means the proportion of concurrence between the user's judgment and the model's judgment.**

| Method | F1 | Precision | Agreement |
|--------|-------|-----------|-----------|
| GET | 0.690 | 0.667 | 70% |
| MUSER | 0.758 | 0.733 | 76.7% |

that by observing the multi-step retrieval of evidence generated by MUSER, human participants can much more accurately decide whether a news article is fake or not.

## 6 CONCLUSION

In this paper, we propose a framework for fake news detection based on multi-step evidence retrieval enhancement—MUSER. Our model leverages a three-phase methodology inspired by human verification processes, including summarization, retrieval, and reasoning. Through text summarization, key information is extracted from the news, reducing irrelevant information. The multi-step retrieval phase enables evidence association for news verification, increasing the dependency between multiple pieces of evidence. Finally, the semantic connection between the news statement and the evidence is analyzed for news classification into two categories: true news and fake news. The results of our experiments on three real-world demonstrated its effectiveness. Moreover, our results also show that evidence association via multi-step retrieval enhances the interpretability of the fake news detection task, making it easier for users to assess the credibility of information and form their own valid judgments.

## 7 ACKNOWLEDGMENTS

## REFERENCES
[1] Hunt Allcott and Matthew Gentzkow. Social media and fake news in the 2016 election. *Journal of Economic Perspectives*, 31(2):211–36, 2017. doi: 10.1257/jep.31.2.211.
[2] Nir Grinberg, Kenneth Joseph, Lisa Friedland, Briony Swire-Thompson, and David Lazer. Fake news on twitter during the 2016 u.s. presidential election. *Science*, 363(6425):374–378, 2019. doi: 10.1126/science.aau2706. URL https://www.science.org/doi/abs/10.1126/science.aau2706.
[3] Kashyap Popat. Assessing the credibility of claims on the web. In *Proceedings of the 26th International Conference on World Wide Web Companion*, pages 735–739, 2017. doi: 10.1145/3041021.3053379.
[4] Shuo Yang, Kai Shu, Suhang Wang, Renjie Gu, Fan Wu, and Huan Liu. Unsupervised fake news detection on social media: A generative approach. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pages 5644–5651, 2019.
[5] Xichen Zhang and Ali A. Ghorbani. An overview of online fake news: Characterization, detection, and discussion. *Information Processing & Management*, 57(2):102025, 2020. doi: https://doi.org/10.1016/j.ipm.2019.03.004. URL

https://www.sciencedirect.com/science/article/pii/S0306457318306794.

[6] Chenguang Song, Kai Shu, and Bin Wu. Temporally evolving graph neural network for fake news detection. *Information Processing & Management*, 58(6): 102712, 2021. doi: https://doi.org/10.1016/j.ipm.2021.102712.

[7] Kai Shu, Amy Sliva, Suhang Wang, Jiliang Tang, and Huan Liu. Fake news detection on social media: A data mining perspective. *ACM SIGKDD Explorations Newsletter*, 19(1):22–36, 2017. doi: 10.1145/3137597.3137600. URL https://doi.org/10.1145/3137597.3137600.

[8] Meeyoung Cha, Wei Gao, and Cheng-Te Li. Detecting fake news in social media: an asia-pacific perspective. *Communications of the ACM*, 63(4):68–71, 2020. doi: 10.1145/3378422.

[9] Mohammad Hadi Goldani, Saeedeh Momtazi, and Reza Safabakhsh. Detecting fake news with capsule neural networks. *Applied Soft Computing*, 101: 106991, 2021. doi: https://doi.org/10.1016/j.asoc.2020.106991. URL https://www.sciencedirect.com/science/article/pii/S1568494620309303.

[10] Huaiwen Zhang, Quan Fang, Shengsheng Qian, and Changsheng Xu. Multi-modal knowledge-aware event memory network for social media rumor detection. In *Proceedings of the 27th ACM International Conference on Multimedia*, pages 1942–1951, 2019. doi: 10.1145/3343031.3350850. URL https://doi.org/10.1145/3343031.3350850.

[11] Martin Potthast, Johannes Kiesel, Kevin Reinartz, Janek Bevendorff, and Benno Stein. A stylometric inquiry into hyperpartisan and fake news. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 231–240, 2018. doi: 10.18653/v1/P18-1022.

[12] Aditi Gupta, Ponnurangam Kumaraguru, Carlos Castillo, and Patrick Meier. Tweetcred: Real-time credibility assessment of content on twitter. In *International conference on social informatics*, pages 228–243, 2014. doi: 10.1007/978-3-319-13734-6_16.

[13] Yang Liu and Yi-Fang Wu. Early detection of fake news on social media through propagation path classification with recurrent and convolutional networks. *Proceedings of the AAAI Conference on Artificial Intelligence*, 32(1):254–261, 2018. doi: 10.1609/aaai.v32i1.11268.

[14] Yi-Ju Lu and Cheng-Te Li. GCAN: Graph-aware co-attention networks for explainable fake news detection on social media. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 505–514, 2020. doi: 10.18653/v1/2020.acl-main.48.

[15] Mansour Davoudi, Mohammad R. Moosavi, and Mohammad Hadi Sadreddini. Dss: A hybrid deep model for fake news detection using propagation tree and stance network. *Expert Systems with Applications*, 198:116635, 2022. doi: https://doi.org/10.1016/j.eswa.2022.116635.

[16] Van-Hoang Nguyen, Kazunari Sugiyama, Preslav Nakov, and Min-Yen Kan. Fang: Leveraging social context for fake news detection using graph representation. In *Proceedings of the 29th ACM international conference on information & knowledge management*, pages 1165–1174, 2020. doi: 10.1145/3517214.

[17] Soroush Vosoughi, Deb Roy, and Sinan Aral. The spread of true and false news online. *Science*, 359(6380):1146–1151, 2018. doi: 10.1126/science.aap9559.

[18] Jing Ma, Wei Gao, Shafiq Joty, and Kam-Fai Wong. Sentence-level evidence embedding for claim verification with hierarchical attention networks. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2561–2571, 2019. doi: 10.18653/v1/P19-1244.

[19] Kashyap Popat, Subhabrata Mukherjee, Andrew Yates, and Gerhard Weikum. DeClarE: Debunking fake news and false claims using evidence-aware deep learning. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 22–32, 2018. doi: 10.18653/v1/D18-1003.

[20] Nguyen Vo and Kyumin Lee. Hierarchical multi-head attentive network for evidence-aware fake news detection. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume.*, pages 965–975, 2021.

[21] Weizhi Xu, Junfei Wu, Qiang Liu, Shu Wu, and Liang Wang. Evidence-aware fake news detection with graph neural networks. In *Proceedings of the ACM Web Conference 2022*, pages 2501–2510, 2022. doi: 10.1145/3485447.3512122.

[22] Robert C. Coghill, John G. McHaffie, and Yi-Fen Yen. Neural correlates of interindividual differences in the subjective experience of pain. *Proceedings of the National Academy of Sciences*, 100(14):8538–8542, 2003. doi: 10.1073/pnas.1430684100.

[23] Andrew Gordon, Jonathan C.W. Brooks, Susanne Quadflieg, Ullrich K.H. Ecker, and Stephan Lewandowsky. Exploring the neural substrates of misinformation processing. *Neuropsychologia*, 106:216–224, 2017. doi: https://doi.org/10.1016/j.neuropsychologia.2017.10.003.

[24] James W Pennebaker, Ryan L Boyd, Kayla Jordan, and Kate Blackburn. The development and psychometric properties of liwc2015. Technical report, 2015.

[25] Victoria L. Rubin and Tatiana Lukoianova. Truth and deception at the rhetorical structure level. *J. Assoc. Inf. Sci. Technol.*, 66(5):905–917, 2015. doi: 10.1002/asi.23216.

[26] Piotr Przybyla. Capturing the style of fake news. *in Proceedings of the AAAI Conference on Artificial Intelligence*, 34:490–497, Apr. 2020. doi: 10.1609/aaai.v34i01.5386.

[27] Yaqian Dun, Kefei Tu, Chen Chen, Chunyan Hou, and Xiaojie Yuan. Kan: Knowledge-aware attention network for fake news detection. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(1):81–89, 2021. doi: 10.1609/aaai.v35i1.16080.

[28] Arkaitz Zubiaga, Ahmet Aker, Kalina Bontcheva, Maria Liakata, and Rob Procter. Detection and resolution of rumours in social media: A survey. *ACM Computing Surveys (CSUR)*, 51(2), 2018. doi: 10.1145/3161603.

[29] Kai Shu, Limeng Cui, Suhang Wang, Dongwon Lee, and Huan Liu. Defend: Explainable fake news detection. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 395–405, 2019. doi: 10.1145/3292500.3330935.

[30] Yiqiao Jin, Xiting Wang, Ruichao Yang, Yizhou Sun, Wei Wang, Hao Liao, and Xing Xie. Towards fine-grained reasoning for fake news detection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 5746–5754, 2022. doi: 10.1609/aaai.v36i5.20517.

[31] Erxue Min, Yu Rong, Yatao Bian, Tingyang Xu, Peilin Zhao, Junzhou Huang, and Sophia Ananiadou. Divide-and-conquer: Post-user interaction network for fake news detection on social media. In *Proceedings of the ACM Web Conference 2022*, pages 1148–1158, 2022. doi: 10.1145/3485447.3512163.

[32] Ruichao Yang, Xiting Wang, Yiqiao Jin, Chaozhuo Li, Jianxun Lian, and Xing Xie. Reinforcement subgraph reasoning for fake news detection. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 2253–2262, 2022. doi: 10.1145/3534678.3539277.

[33] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. Retrieval-augmented generation for knowledge-intensive nlp tasks. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, 2020. doi: 10.5555/3495724.3496517.

[34] Zhijiang Guo, Michael Schlichtkrull, and Andreas Vlachos. A survey on automated fact-checking. *Transactions of the Association for Computational Linguistics*, 10:178–206, 2022. doi: 10.1162/tacl_a_00454.

[35] Canasai Kruengkrai, Junichi Yamagishi, and Xin Wang. A multi-level attention model for evidence-based fact checking. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 2447–2460, August 2021. doi: 10.18653/v1/2021.findings-acl.217.

[36] Yair Feldman and Ran El-Yaniv. Multi-hop paragraph retrieval for open-domain question answering. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2296–2309, July 2019. doi: 10.18653/v1/P19-1222.

[37] Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. A large annotated corpus for learning natural language inference. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 632–642, 2015. doi: 10.18653/v1/D15-1075.

[38] Kashyap Popat, Subhabrata Mukherjee, Jannik Strötgen, and Gerhard Weikum. Credibility assessment of textual claims on the web. In *Proceedings of the 25th ACM International on Conference on Information and Knowledge Management*, pages 2173–2178, 2016. doi: 10.1145/2983323.2983661.

[39] James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Arpit Mittal. FEVER: a large-scale dataset for fact extraction and VERification. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 809–819, 2018. doi: 10.18653/v1/N18-1074.

[40] Ankur Parikh, Oscar Täckström, Dipanjan Das, and Jakob Uszkoreit. A decomposable attention model for natural language inference. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2249–2255, 2016. doi: 10.18653/v1/D16-1244.

[41] Yixin Nie, Haonan Chen, and Mohit Bansal. Combining fact extraction and verification with neural semantic matching networks. *in Proceedings of the AAAI Conference on Artificial Intelligence*, 33:6859–6866, 07 2019. doi: 10.1609/aaai.v33i01.33016859.

[42] Amir Soleimani, Christof Monz, and Marcel Worring. *BERT for Evidence Retrieval and Claim Verification*, pages 359–366. 04 2020. doi: 10.1007/978-3-030-45442-5_45.

[43] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019*, pages 4171–4186, 2019. doi: 10.18653/v1/n19-1423.

[44] Jing Ma, Wei Gao, Shafiq Joty, and Kam-Fai Wong. Sentence-level evidence embedding for claim verification with hierarchical attention networks. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics (ACL 2019)*, pages 2561–2571, 2019. doi: 10.18653/v1/P19-1244.

[45] Li Dong, Nan Yang, Wenhui Wang, Furu Wei, Xiaodong Liu, Yu Wang, Jianfeng Gao, Ming Zhou, and Hsiao-Wuen Hon. Unified language model pre-training for natural language understanding and generation. In *Proceedings of the 33rd International Conference on Neural Information Processing Systems*, pages 13063–13075, 2019. doi: 10.5555/3454287.3455457.

[46] Jingqing Zhang, Yao Zhao, Mohammad Saleh, and Peter J. Liu. Pegasus: Pre-training with extracted gap-sentences for abstractive summarization. In *Proceedings of the 37th International Conference on Machine Learning*, 2020.

[47] Chin-Yew Lin. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81, 2004. URL https://aclanthology.org/W04-1013.

[48] Jeff Johnson, Matthijs Douze, and Hervé Jégou. Billion-scale similarity search with gpus. *IEEE Transactions on Big Data*, 7(3):535–547, 2019.

[49] Nils Reimers and Iryna Gurevych. Sentence-BERT: Sentence embeddings using Siamese BERT-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3982–3992, 2019. doi: 10.18653/v1/D19-1410. URL https://aclanthology.org/D19-1410.

[50] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach. *ArXiv*, abs/1907.11692, 2019.

[51] Chris Samarinas, Wynne Hsu, and Mong Li Lee. Latent retrieval for large-scale fact-checking and question answering with nli training. In *2020 IEEE 32nd International Conference on Tools with Artificial Intelligence (ICTAI)*, pages 941–948, 2020. doi: 10.1109/ICTAI50040.2020.00147.

[52] Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. ALBERT: A lite BERT for self-supervised learning of language representations. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*, 2020. URL https://openreview.net/forum?id=H1eA7AEtvS.

[53] Kai Shu, Deepak Mahudeswaran, Suhang Wang, Dongwon Lee, and Huan Liu. Fakenewsnet: A data repository with news content, social context, and spatiotemporal information for studying fake news on social media. *Big Data*, 8(3):171–188, 2020.

[54] Qiang Liu, Feng Yu, Shu Wu, and Liang Wang. Mining significant microblogs for misinformation identification: An attention-based approach. *ACM Trans. Intell. Syst. Technol.*, 9(5), 2018. doi: 10.1145/3173458.

[55] Yoon Kim. Convolutional neural networks for sentence classification. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1746–1751, 2014. doi: 10.3115/v1/D14-1181.

[56] Zichao Yang, Diyi Yang, Chris Dyer, Xiaodong He, Alex Smola, and Eduard Hovy. Hierarchical attention networks for document classification. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1480–1489, 2016. doi: 10.18653/v1/N16-1174.

[57] Feng Qian, Chengyue Gong, Karishma Sharma, and Yan Liu. Neural user response generator: Fake news detection with collective user intelligence. In *Proceedings of the 27th International Joint Conference on Artificial Intelligence*, page 3834âĂŞ3840, 2018. doi: 10.5555/3304222.3304302.

[58] Lianwei Wu, Yuan Rao, Xiong Yang, Wanzhen Wang, and Ambreen Nazir. Evidence-aware hierarchical interactive attention networks for explainable claim verification. In *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI-20*, pages 1388–1394, 2020. doi: 10.24963/ijcai.2020/193.

# A APPENDIX ON REPRODUCIBILITY

## A.1 Experimental Environment

The experiments are conducted on CentOS 7 servers equipped with GeForce RTX 3080 GPUs, each with 32GB of memory. The code for the experiment is implemented using PyTorch version 1.8.0.

## A.2 Supplementary Experiment

*A.2.1 Large Dataset Experiments.* We further validate the performance of MUSER on two large public datasets. The first is the LIAR dataset (https://www.cs.ucsb.e du/ william/data/liar_dataset.zip). We transform the original LIAR multi-class dataset into a binary classification format, where each sample is labeled as either "true" or "false". We merge the original multiple categories into these two binary labels. We merge "mostly true" and "true" into "true", and "barely-true", "false", and "pants-fire" into "false" to better suit the needs of binary classification problems. Moreover, we compare two representative baseline methods, and the experimental results are given in Table A1.

We also conduct experiments on the Fakeddit dataset (https://github.com/entitize/Fakeddit), which contains multi-modal fake news data (image, text). For this experiment, we select textual data only.

**Table A1: Performance comparison of on LIAR.**

| Method | LIAR | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | F1-Ma | F1-Mi | F1-T | P-T | R-T | F1-F | P-F | R-F |
| BERT | 0.57142 | 0.54946 | 0.51000 | 0.63879 | 0.55120 | 0.58893 | 0.50217 | 0.62307 |
| GET | 0.61413 | 0.61048 | 0.57280 | 0.56667 | 0.57907 | 0.64116 | **0.65400** | 0.63243 |
| **MUSER** | **0.64502** | **0.64413** | **0.64731** | **0.64042** | **0.65434** | **0.64270** | 0.64977 | **0.63577** |

**Table A2: Performance comparison of on Fakeddit.**

| Method | Fakeddit | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | F1-Ma | F1-Mi | F1-T | P-T | R-T | F1-F | P-F | R-F |
| BERT | 0.82086 | 0.75175 | 0.88273 | 0.87511 | 0.89048 | 0.62076 | 0.63874 | 0.60377 |
| GET | 0.84651 | 0.79794 | 0.89700 | 0.87060 | 0.92506 | 0.68538 | 0.76689 | 0.61195 |
| **MUSER** | **0.86785** | **0.82201** | **0.91611** | **0.88847** | **0.94552** | **0.68887** | **0.77871** | **0.62761** |

**Table A3: Performance comparison of retrieval budget experiments.**

| Method | PolitiFact | | Gossipcop | | Weibo | |
|---|---|---|---|---|---|---|
| | F1-Ma | F1-Mi | F1-Ma | F1-Mi | F1-Ma | F1-Mi |
| N1 = 30 | 0.7193 | 0.7155 | 0.7600 | 0.7580 | 0.7951 | 0.7930 |
| N1 = 60 | 0.7117 | 0.7110 | 0.7560 | 0.7555 | 0.7868 | 0.7859 |
| N1 = 90 | 0.7064 | 0.7057 | 0.7546 | 0.7542 | 0.7855 | 0.7846 |
| **MUSER** | **0.7324** | **0.7293** | **0.7760** | **0.7751** | **0.8042** | **0.8026** |

The label used is the '2-way' introduced in the data set, which is divided into two categories: true and false. The experimental results are given in Table A2.

*A.2.2 Retrieval Budget Experiments.* In order to assess the effectiveness of multi-step retrieval and investigate whether increasing the retrieval budget can serve as a substitute, we performed the following experiments. To compare with our MUSER's configuration, which is a 3-step retrieval with $N1 = 30$, $N2 = 30$, and $N3 = 30$, we use three configurations of 1-step retrieval with $N1 = 30$, $N1 = 60$ or $N1 = 90$. The results are reported in Table A3. Surprisingly, increasing the pool of 1-step retrieval has a negative effect on the performance, and the patterns are consistent across the three datasets. The reasons for this can be two folds. First, there exists a degree of interdependence among certain pieces of evidence, which requires additional information from previous retrieval steps for accurate identification. Second, the inclusion of an excessive number of paragraphs introduces a higher level of noise into the text, ultimately leading to suboptimal results. Consequently, our investigation has demonstrated that the proposed 1-step retrieval with an augmented budget is not a viable alternative to MUSER's multi-step retrieval approach.

*A.2.3 Threshold $\lambda$ selection.* $\lambda$ essentially serves as an evaluation metric for the correlation between evidence and passage text. In our research, we employed a unified threshold value of $\lambda = 0.9$ across all datasets. To investigate the effect of $\lambda$ value, we conduct experiments to analyze its impact on the retrieved evidence. And the results reveal that a low $\lambda$ value tends to introduce more noise in the retrieved evidence, whereas a high $\lambda$ value may inadvertently exclude critical evidence. For your reference, the comparison results with different $\lambda$ values are given in Table A4.

**Table A4: The comparison results with different $\lambda$ values.**

| $\lambda$ | PolitiFact | | Gossipcop | | Weibo | |
|---|---|---|---|---|---|---|
| | **F1-Ma** | **F1-Mi** | **F1-Ma** | **F1-Mi** | **F1-Ma** | **F1-Mi** |
| **0.8** | 0.699 | 0.698 | 0.734 | 0.735 | 0.787 | 0.786 |
| **0.85** | 0.712 | 0.710 | 0.755 | 0.753 | 0.796 | 0.795 |
| **0.9** | **0.732** | **0.729** | **0.776** | **0.775** | **0.804** | **0.802** |
| **0.9** | 0.715 | 0.713 | 0.739 | 0.738 | 0.789 | 0.788 |

## A.3 Code Resources

We compare the proposed framework, MUSER, with 9 baseline methods discussed in Section 5.2, the content-based methods including TextCNN, TextRNN, TCNNURG, BERT, and the evidence-based methods including DeClarE, HAN, EHIAN, MAC, and GET. The implementation details of our proposed framework, including code and settings, are available through the following link: https://github.com/Complex-data/MUSER. Other codes were obtained as follows:

- **TextCNN:** we use the publicly available implementation at: https://github.com/FinIoT/text_cnn
- **TextRNN:** we use the publicly available implementation at: https://github.com/luchi007/ RNN_Text_Classify
- **TCNNURG:** we use the publicly available implementation at: https://github.com/text_classify

- **BERT:** we use the publicly available implementation at: https://github.com/google-research/bert
- **DeClarE:** we use the publicly available implementation at: https://github.com/atulkumarin/DeClare
- **HAN:** we use the publicly available implementation at: https://github.com/majingCUHK/Claim_Verification
- **EHIAN:** we use the publicly available implementation at: https://github.com/evidence-inference
- **MAC:** we use the publicly available implementation at: https://github.com/nguyenvo09/EACL2021
- **GET:** we use the publicly available implementation at: https://github.com/CRIPAC-DIG/GET

## A.4 Corpus processing

In this article, we use Wikipedia data as the retrieval corpus. The download address of Wikipedia Chinese corpus is: https://dumps.wikimedia.org/zhwiki/latest/, and the download address of Wikipedia English corpus is: https://dumps.wikimedia.org/enwiki/latest/.

We extract the Wikipedia corpus through WikiExtractor, which can extract the main article content of the corpus ending with .bz downloaded from Wikipedia. The download address of the tool is: https://github.com/attardi/wikiextractor.