

ViWikiFC: Fact-Checking for Vietnamese Wikipedia-Based Textual Knowledge Source

Hung Tuan Le^{1,2}, Long Truong To^{1,2}, Manh Trong Nguyen^{1,2},
Kiet Van Nguyen^{1,2*}

¹Faculty of Information Science and Engineering, University of Information
Technology, Ho Chi Minh City, Vietnam.

²Vietnam National University, Ho Chi Minh City, Vietnam.

*Corresponding author(s). E-mail(s): kietnv@uit.edu.vn;

Contributing authors: 21520250@gm.uit.edu.vn; 21521101@gm.uit.edu.vn;
21520343@gm.uit.edu.vn;

Abstract

Fact-checking is essential due to the explosion of misinformation in the media ecosystem. Although false information exists in every language and country, most research to solve the problem mainly concentrated on huge communities like English and Chinese. Low-resource languages like Vietnamese are necessary to explore corpora and models for fact verification. To bridge this gap, we construct ViWikiFC, the first manual annotated open-domain corpus for **Vietnamese Wikipedia Fact Checking** more than 20K claims generated by converting evidence sentences extracted from Wikipedia articles. We analyze our corpus through many linguistic aspects, from the new dependency rate, the new n-gram rate, and the new word rate. We conducted various experiments for Vietnamese fact-checking, including evidence retrieval and verdict prediction. BM25 and InfoXLM_{Large} achieved the best results in two tasks, with BM25 achieving an accuracy of 88.30% for SUPPORTS, 86.93% for REFUTES, and only 56.67% for the NEI label in the evidence retrieval task, InfoXLM_{Large} achieved an F₁ score of 86.51%. Furthermore, we also conducted a pipeline approach, which only achieved a strict accuracy of 67.00% when using InfoXLM_{Large} and BM25. These results demonstrate that our dataset is challenging for the Vietnamese language model in fact-checking tasks.

Keywords: Fact Checking, Language Model, Information Verification, Corpus

1 Introduction

Following the increase in the amount of information, the lack of strict policies when spreading information has led to misinformation and disinformation on social media. This can cause conflict and manipulate the group of people. In an effort to minimize the impact of misinformation and disinformation in society, some organizations like PolityFact¹, FactCheck.org², play an essential role in online Fact Verification when manually verifying claims based on a different source of evidence. Manual verification is a time-consuming task since a Fact-Checker would have to search through many potential sources for the evidence relevant to the claim. This work is insufficient with the speed of updated information on social media.

Automatic Fact-Checking is a complicated task. In a recent survey on fact-checking [1] that can be separated into four sub-tasks: claim detection, where a machine has to choose what claim needs to be verified; evidence retrieval is to find evidence that SUPPORTS or REFUTES the claim, and verdict prediction is an attempt to determine the veracity of the claim based on evidence in the evidence retrieval task. Finally, the justification production is conducted with a reasonable explanation of why a claim is supported or refuted by the evidence. However, there are many works on how to improve the performance of fact-checking systems in different ways, from applying fast dot product indexing in evidence retrieval [2, 3], constructing new reasoning methods for the verification task [4] to a large-scale corpus with natural input such as NELA [5], FakeCovid [6], to artificial input for research purposes such as FEVER [7], VitaminC [8], languages as Vietnamese are still considered a low resource for natural language processing (NLP) due to lack of corpora, especially in automatic fact-checking.

To contribute to the process of NLP research for Vietnamese, especially in Fact-Checking, we present the first large-scale, open-domain corpus for Vietnamese Fact-Checking on Wikipedia, ViWikiFC: Vietnamese Wikipedia Fact-Checking. The corpus consists of 20,916 claims manually annotated and based on evidence retrieved from Wikipedia pages. Our corpus takes advantage of FEVER [7]; our corpora also have three label classes: SUPPORTS, REFUTES, and NOTENOUGHINFORMATION (NEI). In our corpora, the evidence is manually rewritten into three types of claims based on three class labels. In this way, claims are more realistic when in reality and very diverse in many semantic ways from the way claim construction in the FEVER corpus [7] based on just one piece of evidence. Furthermore, we also follow ViNLI [9] corpus creation to create a double claim with different meanings for each label from one evidence to capture more than just one piece of information in evidence.

Before corpus creation, we developed a suitable guideline and annotation tool for the annotation process. Annotators who are Vietnamese native speakers with an education background beyond high school were carefully trained with the guidelines and a strict training process to ensure corpus creation consistency.

We conducted two experiments, evidence retrieval, and verdict prediction, to evaluate the challenge of ViWikiFC. For the evidence retrieval task, the input is a claim sentence and a document corpus. We use TF-IDF and BM25 for lexical match, and Vietnamese-SBERT [10] for semantic search to retrieve a sentence most relevant to the claim from the document corpus. For the verdict prediction (VP), we used three deep neural networks, including CBOW [11], BiLSTM [12]. We also evaluated using SOTA pre-trained language models such as Multilingual BERT (mBERT) [13], XLM-R [14], and InfoXLM [15], as well as monolingual

¹<https://www.politifact.com/>

²<https://www.factcheck.org/>

models designed for Vietnamese, such as PhoBERT [16], ViT5 [17], and ViDeBERTa [18]. These models have shown exceptional ability in performing the VP task. Ultimately, we have seamlessly integrated the top-performing models for both tasks into a comprehensive pipeline and thoroughly assessed the effectiveness of the resulting system on a test set.

Contributions in this paper are described as follows.

- Firstly, we propose the first fact-checking corpus on Vietnamese Wikipedia, comprising 20,916 claim sentences manually annotated based on 73 Wikipedia articles built with a consensus agreement among annotators, achieving a 95.87% Fleiss' κ -agreement.
- Second, we conduct two experiments on two tasks, evidence retrieval and verdict prediction, while also developing a pipeline approach for fact-checking models, encompassing both neural network-based and pre-trained transformer-based models.
- Next, we analyze the corpus from various linguistic perspectives in the Vietnamese language to gain more insight into Vietnamese fact-checking, including the new n-gram rate, the new dependency rate, and the new word rate.
- Next, we analyze the experimental results of language models from different aspects of the corpus to highlight the challenges introduced by our corpus.
- Finally, our corpus and annotation tool is available for research purposes³.

Our paper is structured into five sections: Section 1 provides an overview of the problem and our contributions, Section 2 introduces related research on fact-checking, including existing corpora (see Section 2.1) and baselines (see Section 2.2). Section 3 describes the corpus construction process, including annotator recruitment and training (see Section 3.1), annotation and validation tool (see Section 3.2, selection of data sources (see Section 3.3), claim generation (see Section 3.4), corpus validation (see Section 3.5), and analysis of the linguistic aspects of the corpus (see Section 3.6). Section 4 introduces the baseline models tested on various tasks (see Section 4.1), summarizes the experimental results (see Section 4.2), and analyzes the impact of the data set on the methods (see Section 4.3). Finally, in Section 5, we conclude our research findings and propose future directions for the fact-checking task in Vietnamese.

2 Related Work

In this section, we introduce the corpus in Section 2.1 as well as the baseline models in Section 4.1 for the fact-checking task.

2.1 Fact-Checking Corpora

We performed a review of existing fact-checking datasets, as summarized in Table 1. We categorized the datasets based on six attributes: data set size (claims), language, data source, domain, annotated evidence, and annotated label

The emergence of the Politifact corpus [19] constructed by Vlachos and Riedel in 2014 laid the foundation for investigating the verification of claims with a structure made up of three main tasks, known as claim detection, evidence retrieval, and verdict prediction, in recent fact-checking studies. However, the corpus still carried the limitation of being too small, restricting the training and evaluation of neural network-based models and following a similar approach to the Politifact corpus [19], Liar [20] was published with more than 12.8K samples extracted

³<https://github.com/HighWill0/ViWikiFC.git>

Corpus	Claims	Language	Source	Domain	Annotated Evidence	Annotated Label
PolitiFact (2014)	106	English	PolitiFact	Politics	✗	✗
Liar (2017)	12,836	English	PolitiFact	Politics	✗	✗
FEVER (2018)	185,445	English	Wiki	Multiple	✓	✓
PUBHEALTH (2020)	11,832	English	Fact-Checking/News	Health	✗	✗
HOVER (2020)	26,171	English	Wiki	Multiple	✓	✓
TabFact (2020)	92,283	English	Wiki	Multiple	✗	✓
InfoTabs (2020)	23,738	English	Wiki	Multiple	✗	✓
ANT (2020)	4,547	Arabic	News	Multiple	✗	✗
FakeCovid (2020)	5,182	Multilingual (3)	Fact Check	Health	✗	✗
X-Fact (2021)	31,189	Multilingual (25)	Fact Check	Multiple	✓	✗
VitaminC (2021)	488,904	English	Wiki	Multiple	✗	✓
DanFEVER (2021)	6,407	Danish	Wiki	Multiple	✓	✓
CHEF (2022)	10,000	Chinese	Internet	Multiple	✓	✓
ViWikiFC	20,916	Vietnamese	Wiki	Multiple	✗	✓

Table 1: Overall statistics of fact-checking corpora.

from the Politifact⁴ website. However, the accuracy prediction in this corpus relied solely on claims or metadata without utilizing external sources and did not provide any evidence for verification.

Fever [7] introduced a large-scale corpus with more than 180k claim sentences manually written by a team of annotators based on Wikipedia⁵ pages. This corpus presented an automated information verification process by combining models for distinct tasks into a complete pipeline. In subsequent years, new corpora were introduced, such as Snopes [21] with more than 6K samples and MultiFC [22] with more than 36k samples collected from various fact-checking websites such as PolitiFact⁶. Particularly in the MultiFC [22] corpus, claims were collected from multiple sources, resulting in label counts ranging from 2 to 27. From 2020 onward, a series of corpora appeared with various approaches. For example, Hover [23] had more than 26k samples and was built and verified using multiple Wikipedia pages. The VitaminC corpus [8] was built with more than 488k samples based on actual revisions of Wikipedia, all including evidence.

Apart from the Corpora constructed from regular text data, there are also Corpora for semi-structured data, such as tables. TabFact [24], with more than 92k samples, and InfoTabs [25], with more than 23k samples, were constructed based on tables from Wikipedia. The FEVEROUSE corpus [26], with more than 87K samples, uses both text and a table for verification. Corpora in other languages have been developed and researched as well. X-Fact [27] with more than 31K samples in 25 different languages excluding English, or ANT [28] with more than 4K samples for Arabic, DanFEVER [29] with more than 6K samples for Danish, and the most recent CHEF corpus [30] with 10K samples for Chinese.

With the growing demand for the development of automatic fact-checking tasks, especially in Vietnamese NLP research, an open-domain, large-scale, and high-quality corpus becomes crucial. Following the corpus development efforts, we constructed a Vietnamese Fact-Checking corpus alongside an information verification process using pre-trained language models in the Vietnamese language, contributing to advancing research and addressing Fact-Checking challenges in Vietnamese.

⁴<https://www.politifact.com/>

⁵<https://www.Wikipedia's.org/>

⁶<https://www.politifact.com/>

2.2 Fact-Checking Methods

As previously stated, fact-checking encompasses four distinct sub-tasks: claim detection, evidence retrieval, verdict prediction, and justification production. Of these sub-tasks, evidence retrieval and verdict prediction are particularly crucial, as they serve as the basis for determining ground-truth information and information inference mechanisms to verify information. Our research focuses on these two sub-tasks, so in this section, we introduce the methods used in the two sub-tasks.

Evidence retrieval is the task of searching and extracting information from various sources and types of data, such as text, tables, images, and knowledge bases, to verify a claim. In recent research, evidence retrieval is often approached through two main directions. First, the traditional non-neural approach, such as TF-IDF [7, 30] and BM25 [31], which uses sparse representations to measure term matching between a claim and evidence. Although this method doesn't require much time and resources, it can only capture lexical information and is limited in handling various queries. The second approach is neural network-based, which can calculate and measure the semantic similarity between evidence and a claim. In 2020, Karpukhin et al. introduced dense passage retrieval using a dual-encoder architecture to represent questions and passages as dense vectors. This approach has successfully incorporated BERT-based language models and their variants into open-domain question-answering [32–34] and fact-checking tasks [35, 36]. Cross-encoder architecture is also utilized for re-ranking passages and sentences by BERT-based and its variants re-rankers [37–39], which achieved significant improvement compared to traditional methods.

Verdict prediction is a task to verify claims based on evidence extracted from the evidence retrieval task. The input for this task consists of a pair of sentences, including evidence and claims. However, for datasets with natural claims [5, 6, 22] or datasets built for research purposes [23, 24], the claim needs to be inferred and verified from information in multiple pieces of evidence. In recent years, solutions have been proposed, such as concatenating all pieces of evidence into one [40] or using special components like kernel-based attention [41] or a graph-based approach [42]. These approaches create conditions for verifying information on more complex claims [43–45]. Early approaches for this task used neural network-based methods such as CBOW [11], RNN [46], and BiLSTM [12] to verify claims, achieving high performance. In 2018, BERT [47], a pre-trained language model based on the transformer architecture [48], significantly improved in various natural language processing tasks, including verdict prediction. In the following years, many pre-trained language models have been developed, such as RoBERTa [49], ALBERT [50], and XLM-R [14], which are multilingual pre-trained language models. The first monolingual pre-trained language model for Vietnamese is PhoBERT [16] developed with the same architecture as BERT [47] and the same pre-training technique as RoBERTa [49]. PhoBERT was trained on a 20GB word-level Vietnamese Wikipedia and News corpus, achieving state-of-the-art results in many downstream tasks such as NER, POS, and NLI. Since then, many pre-trained language models for Vietnamese have been developed, such as ViT5 [17], a transformer-based encoder-decoder model with the same architecture and training method as T5 [51], and ViDeBERTa [18], which uses the DeBERTa architecture [52]."

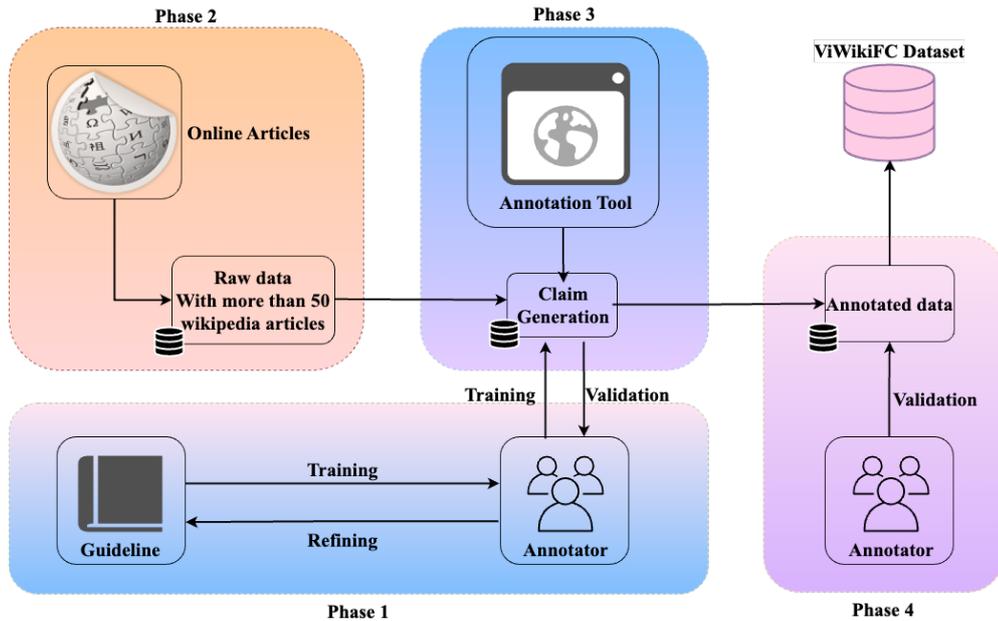


Fig. 1: The creation process of the fact-checking corpus.

3 Corpus Creation

The corpus was created throughout four phases (see Figure 1), including annotator recruitment and training (see Section 3.1), evidence selection (see Section 3.3), claim generation (see Section 3.4), and corpus validation (see Section 3.5). Furthermore, we also analyzed the corpus (see Section 3.6) in terms of various linguistic aspects to gain a deeper understanding of the language characteristics of the corpus.

3.1 Annotator Recruitment and Training

We hired thirty annotators who are Vietnamese native speakers and have high school to college literacy to create claims from the evidence given. These annotators are in the habit of reading and searching for information on the Internet. Annotators must undergo a strict training process with two phrases before entering the official annotation process; we use the Fleiss κ score to calculate the agreement of annotators with the authors.

In the first phase, the authors trained the annotators through the provided guidelines to familiarize themselves with the rules when creating claims. Next, the annotators were required to create a claim in the training corpus that contained 50 evidence sentences. The annotator training corpus labels (SUPPORTS, REFUTES, and NOTENOUGHINFORMATION) were marked. One of Three authors read and verdict the label for every claim-evidence pair; if the proportion of the labels agreed upon by Two other authors is over 0.90, the annotator is selected for the second phase of the training process. If the result is not over 0.95, the annotator must learn the mistake from the previous training corpus and go through the first training process again with another training corpus.

In the second phase, thirty annotators are randomly divided into six distinct groups, each consisting of five individuals. Before this phase, the thirty individuals worked on at least one data folder containing 150 claim-evidence pairs. The thirty data folders have their labels concealed. Each group is assigned five data folders (already label-hidden) from any of the remaining five groups. All members within a group then relabel the data folders of the members in the other designated group. Labeling results are used to calculate the level of consensus using the Fleiss κ score [53]. If the consensus level among the five annotators on a data folder is greater than or equal to 0.95, the Annotator who labels that folder is selected for the official annotation process. In cases where the consensus level is below 0.95, the Annotator must go through the first training process. The procedure for calculating consensus is repeated for all annotators in different groups until all annotators have been evaluated for consensus. During test annotation, we encourage annotators to list exceptional cases (multi-meaning in a sentence, unknown subject) and note frequent mistakes in the guideline.

3.2 Fact-Checking Annotation Tool

In order to create the environment for corpus creation, we developed a manual labeling tool for annotators. The tool is built using the Java Development Kit ⁷, which makes it compatible with any operating system and easy to use. Figure 2 shows a complete screen after finishing the generated claim in Section 3.4. During claim generation, the annotator must follow some strict rules, including the following.

- **Rule 1:** Generating all six claims in one sentence of evidence before moving on to the next.
- **Rule 2:** The annotator must accurately write the claim when generating the claim.
- **Rule 3:** The first letter of a claim must be capitalized, and the end must be followed by a period.
- **Rule 4:** Regional or specific dialect words should be avoided in writing a claim.

Furthermore, we also built a validation tool (see Figure 3) in the Corpus Validation (see Section 3.5).

3.3 Evidence Selection

We use Wikipedia as the primary data source to build the ViWikiFC corpus. We do not prioritize establishing the credibility of the data source, as, in reality, all data sources can sometimes present inaccuracies or omit information, either unintentionally or deliberately. Therefore, no source can be deemed completely trustworthy. The team considers the coverage of information in various aspects of life as a crucial factor. Wikipedia, being an open encyclopedia with more than 55 million articles on a wide array of topics, helps to build a corpus with greater comprehensiveness compared to constructing it solely from other sources, such as political information from Politifact⁸ or news articles from VnExpress⁹. We extracted 3,812 evidence sentences from more than 1,479 paragraphs in 73 articles on Wikipedia, covering various topics such as history, geography, philosophy, and science.

⁷<https://www.oracle.com/java/>

⁸<https://www.politifact.com/>

⁹<https://vnexpress.net/>

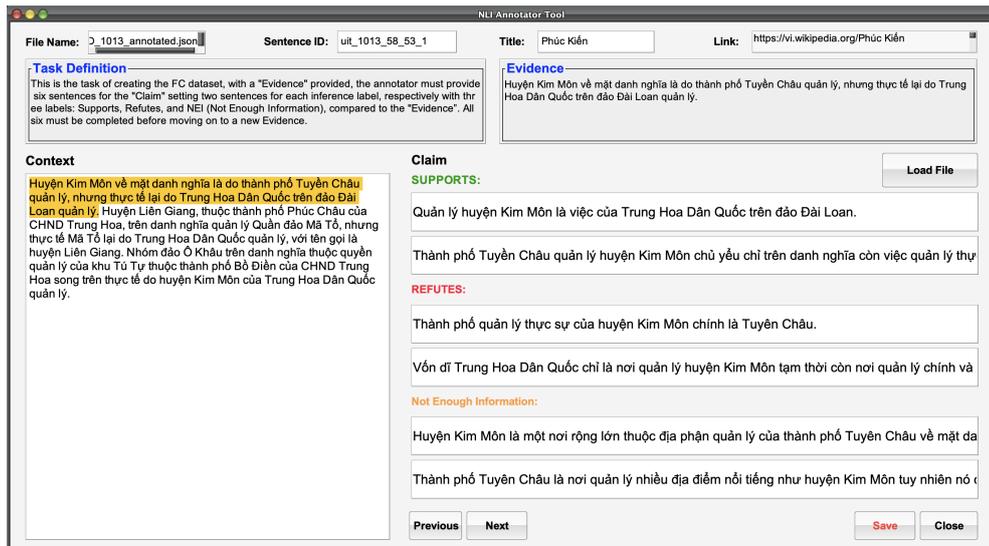


Fig. 2: ViWifiFC annotation tool for fact-checking corpus.

3.4 Claim Generation

Rather than creating only two labels representing SUPPORTS or REFUTES, we add a NOT ENOUGH INFORMATION label as in FEVER [7]. In reality, we encounter claims that do not have enough evidence to prove that they are entirely leaning toward a side. We ask annotators to create claim sentences for the following three labels.

- **SUPPORTS:** Annotators create a claim that we can determine to be true based only on information from the evidence.
- **REFUTES:** Annotators create a claim that we can determine is false based only on the information in evidence.
- **NOT ENOUGH INFORMATION (NEI):** Annotators create a claim that we cannot determine to be true or false based only on the evidence provided.

For two labels, SUPPORTS and REFUTES, the annotator composes a claim sentence based on the information provided by the evidence sentence. When writing the claim, the annotator is not allowed to use personal understanding or external knowledge beyond the evidence. When crafting the claim sentence, the annotator is free to use any information provided by the evidence sentence and to adopt various writing styles to improve the complexity and diversity of expression in the claim that is different from the generation of FEVER claims [7] when a single piece of information conducts the claim. For NOTENOUGHINFORMATION labels, the annotator uses at least one subject, object, or event in the evidence sentence, plus some information outside of it but still in the context of the evidence sentence. In this way, we can construct an evidence-retrieval experiment for this label in the section.

Annotators are required to write two distinct claim sentences for each corresponding label for each evidence sentence, resulting in six claim sentences for the three labels from a single evidence sentence. The construction of claim sentences is analogous to the approach used in

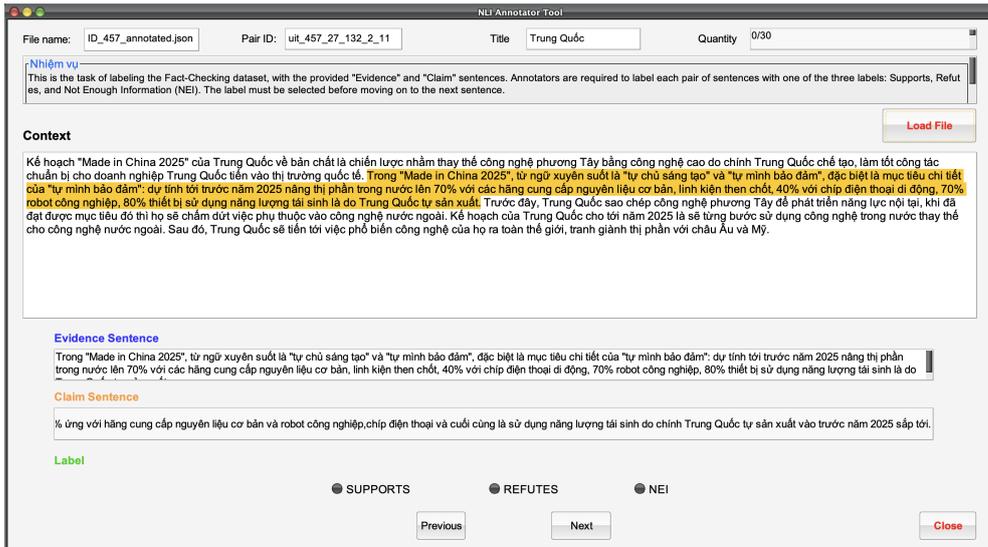


Fig. 3: ViWiFiFC validation tool for fact-checking corpus.

Sentence	
Evidence	Hitler gây ra sự kiện "Đêm của những con dao dài", giết hại các đối thủ của ông. (<i>Hitler caused the "Night of the Long Knives" event, killing his opponents.</i>)
Claim	Hitler chịu trách nhiệm cho việc giết hại các đối thủ trong sự kiện "Đêm của những con dao dài". (<i>Hitler was responsible for killing his opponents during the "Night of the Long Knives".</i>)

Table 2: Example of generating SUPPORTS claim based on evidence with few tokens.

the two corpora: OCNLI [54] and IndoNLI [55], where the annotators are instructed to take advantage of multiple pieces of information provided by the evidence sentence and use various expressions for the same information template. Furthermore, for a given evidence sentence, the claim sentences corresponding to different labels exhibit differentiation in writing style and the information used to compose the claim.

During the Annotator training process, we observed that a significant portion of annotators encountered difficulties in composing claim sentences with the SUPPORTS label, especially for evidence sentences with limited exploitable information, without excessively mirroring the evidence sentence in the Paraphrase. This occurs with evidence sentences that have a length of 15 to 25 tokens. Additionally, the REFUTE label poses challenges for annotators when constructing sentences by negating the information in the evidence sentence. However, frequent and dense negation may lead to a shift from the REFUTE label to the NEI label. Table 3 illustrates REFUTED claims with too many information changes; if not written and carefully reviewed, they tend to lean towards the NEI label.

Sentence	
Evidence	Con đường tơ lụa được nhà Tùy chú trọng kế thừa và phát triển vào thời kỳ hưng thịnh của triều đại này vào thế kỷ thứ 7 . (<i>The Silk Road was emphasized for inheritance and development by the Sui dynasty during its prosperous period in the 7th century.</i>)
Claim	Con đường tơ lụa vẫn được kế thừa và xây dựng dù nhà Đường bước vào giai đoạn suy vong của thế kỷ thứ 10 . (<i>The Silk Road continued to be inherited and developed even as the Tang dynasty entered a period of decline in the 10th century.</i>)

Table 3: Example of generating REFUTED claims using more than two rules lead to NEI label.

3.5 Corpus Validation

In the corpus evaluation process, we establish a two-part evaluation procedure, namely annotator validation and author validation. Both of these procedures are conducted during the claim generation phase to ensure a high level of consensus throughout the entire labeling process.

3.5.1 Annotator Validation

During the claim generation phase, upon reaching 50.00% corpus completion, we identify the top 5 annotators who have made significant contributions to the data construction. Subsequently, we select a distinct corpus these 5 Annotators have not previously labeled. This new corpus comprises 750 pairs of claim-evidence sentences. The existing set of labels on this corpus is then removed. All five annotators collaboratively apply labels to each sentence pair until all 750 claim-evidence pairs have been re-evaluated. In cases where fewer than three of the five annotators assign an identical label to a sentence pair, that particular sentence is designated with a "-" mark, and those pairs are subsequently eliminated from the corpus. The agreement rate among the annotators reached 99.87% when at least three annotators were assigned the same label, and the agreement rate when at least two annotators were assigned the same label was 95.47%. Additionally, we measured the inter-annotator agreement using Fleiss κ score, and the achieved result was 95.87%, higher than the agreement rate in the FEVER corpus, which reached only 68.41%.

3.5.2 Author Validation

We used 750 pairs of claim-evidence sentences, which were relabeled in Section 3.5.1 and evaluated whether the written claim sentences adhere to the guidelines. We observed that 96.12% of claim sentences followed the stipulated guidelines, while 3.87% of the claim sentences contained writing errors and did not adhere to the annotation guideline rules.

3.6 Corpus Analysis

3.6.1 Overall Statistics

We divide the corpus, which consists of 20,916 claim sentences, randomly into three different sets: 80.00% for training (Train), 10.00% for development (Dev), and 10.00% for the test (Test) set for VP task. For ER tasks, we do not employ pre-trained models on a portion of the data

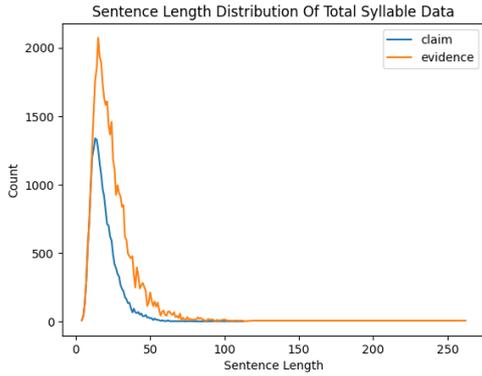


Fig. 4: The distribution total syllable in the corpus.

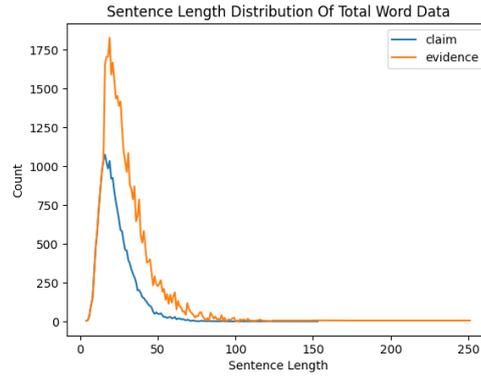


Fig. 5: The distribution total word in the corpus.

and then evaluate the remaining part. Instead, we use the cosine similarity measure between sentences within the model to extract evidence sentences, and this process will be carried out and evaluated on the entire corpus. The distribution of labels is summarized in Table 4.

	SUPPORTS	REFUTES	NEI
Train	5,594	5,573	5,571
Dev	666	694	730
Test	708	706	677

Table 4: Overall statistics of our corpus.

3.6.2 Length Distribution

Allocate the sentences of claims and evidence based on their length, as presented in Figure 5 and Figure 4. The length of claims is shorter compared to the evidence, but they still ensure a complete and accurate representation of the information conveyed by the evidence sentences. In reality, claim sentences often tend to be shorter in length, whereas evidence sentences are typically longer, serving the role of providing additional details and supporting the claim. Furthermore, writing too long sentences can result in providing more context and information to the model, making training and prediction less challenging and more straightforward. The shortest length for a claim is four words, while the longest when generating is 113 words. Sentences within the length range of 12 to 25 words constitute the largest proportion of the corpus.

3.6.3 New Word Rate

To evaluate diversity in the corpus, we measure the rate of new words in the claim sentence that do not appear in the evidence sentence. We use VnCoreNLP [56] for word segmentation. The results in table 5 show that the new word rate in the REFUTES label is the lowest, reaching only 31.28%. The word diversity of the NEI label (50.44%) is the highest among the three

Label	New Word Rate (%)	New Dependency Rate (%)	Part-Of-Speech (%)					
			Noun	Verb	Adjective	Preposition	Adjunct	Other
SUPPORTS	32.91	75.26	26.97	30.22	7.65	10.81	7.92	16.45
REFUTES	31.28	67.90	26.27	25.47	8.66	8.04	13.10	18.47
NEI	50.44	81.96	33.96	23.61	8.96	9.82	7.65	15.99

Table 5: Corpus analysis in terms of linguistic aspects.

labels, followed by the SUPPORTS label with 32.91%. Both the SUPPORTS and REFUTES labels have a low new word rate compared to the NEI label, indicating that annotators tend to use fewer new words when constructing sentences for these two labels. This tendency ensures semantic accuracy with the evidence sentence when external knowledge is not used. In Addition, we utilized part-of-speech analysis on new words to delve deeper into annotator data construction trends. We used PhoNLP [57] for word segmentation and summarized the results in Table 5. The statistical results reveal that nouns and verbs are the two primary components annotators use when composing claim sentences.

3.6.4 New Dependency Rate

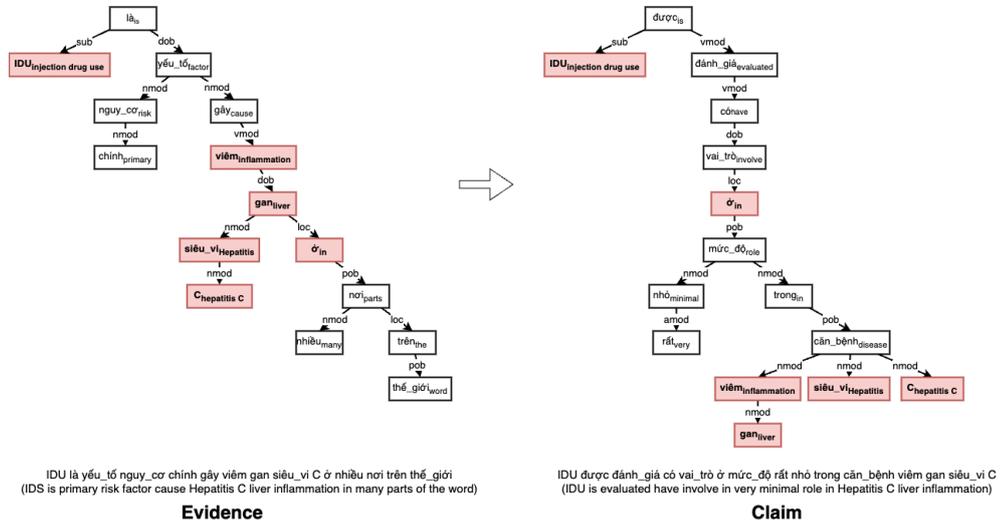


Fig. 6: Semantic dependency tree from evidence to claim.

We measured the new dependency rate in all three labels to analyze the creativity in constructing claim sentences from evidence sentences. We use the VnCoreNLP toolkit [56] for dependency parsing on the claims and the evidence sentences. Subsequently, we calculated the dependencies in the claim sentences that were absent in the corresponding evidence sentences.

Figure 6 illustrates the modification of the semantic dependency tree from evidence to claim. The word in the red box appears in two sentences, and its position and dependency change throughout the syntax tree.

The results are summarized in Table 5. All three labels exhibit a new dependency rate of more than 50.00%, with NEI having the highest rate at 81.96%, followed by SUPPORTS at 75.26%, and REFUTES at 67.90%. The statistical results indicate that in the ViWikiFC dataset, the claim sentences constructed from the evidence sentences demonstrate the grammatical structure and semantic diversity. This contributes to making the ViWikiFC dataset more challenging when evaluating language models.

3.6.5 New n-gram rate

New n-gram rate	ViWikiFC			VitaminC		
	SUPPORTS	REFUTES	NEI	SUPPORTS	REFUTES	NEI
Unigram	31.71	30.43	51.95	38.85	49.77	58.56
Bigram	57.37	53.78	71.70	70.35	79.06	84.50
Trigram	73.27	69.06	81.93	84.01	90.40	92.77
Fourgram	81.92	77.75	86.95	90.30	94.68	96.14

Table 6: New n-gram rate (%) between claim-evidence pairs.

To assess the similarity between the claim and evidence, we calculated the percentage of new n-grams between the claim and evidence sentences for different labels in our dataset. Furthermore, we have implemented a comparable approach to the VitaminC corpus [8] due to the resemblance between our corpus and the challenging corpus VitaminC. We used the nltk library ¹⁰ to segment sentences into adjacent sequences of n syllables, ranging from one to four. The results are summarized in Table 6

From Table 6, across both corpus, the percentage of new n-grams is highest for the NEI label among the three labels. This is understandable because NEI sentences are written by annotators using additional information beyond the evidence sentence. In contrast, sentences with the SUPPORTS and REFUTES labels are written solely based on the information provided in the evidence, leading to a higher overlap of words or n-grams between the claim and evidence sentences, resulting in lower percentages of new n-grams in these two labels compared to the NEI label.

Furthermore, we also calculated the percentage distribution of the number of new n-grams in both data sets, as shown in Figure 7. The statistical results indicate that the distribution of new n-grams in the ViWikiFC data set is more evenly distributed. Sentences with more than 11 new n-grams have a high proportion in the dataset, in contrast to the VitaminC dataset, where the percentage is relatively small, indicating a more balanced distribution in the ViWikiFC dataset compared to VitaminC.

3.6.6 Corpus-generation rules analysis

To evaluate and analyze the linguistic factors of the annotators during the creation of the corpus, we proceed with the analysis of the rules used based on the ViNLI (Huynh et al.) [] corpus generation process. We randomly chose 500 SUPPORTS and 500 REFUTES claim-evidence pairs from the validation set. Selecting 500 sentences for each label will increase the

¹⁰<https://www.nltk.org/index.html>

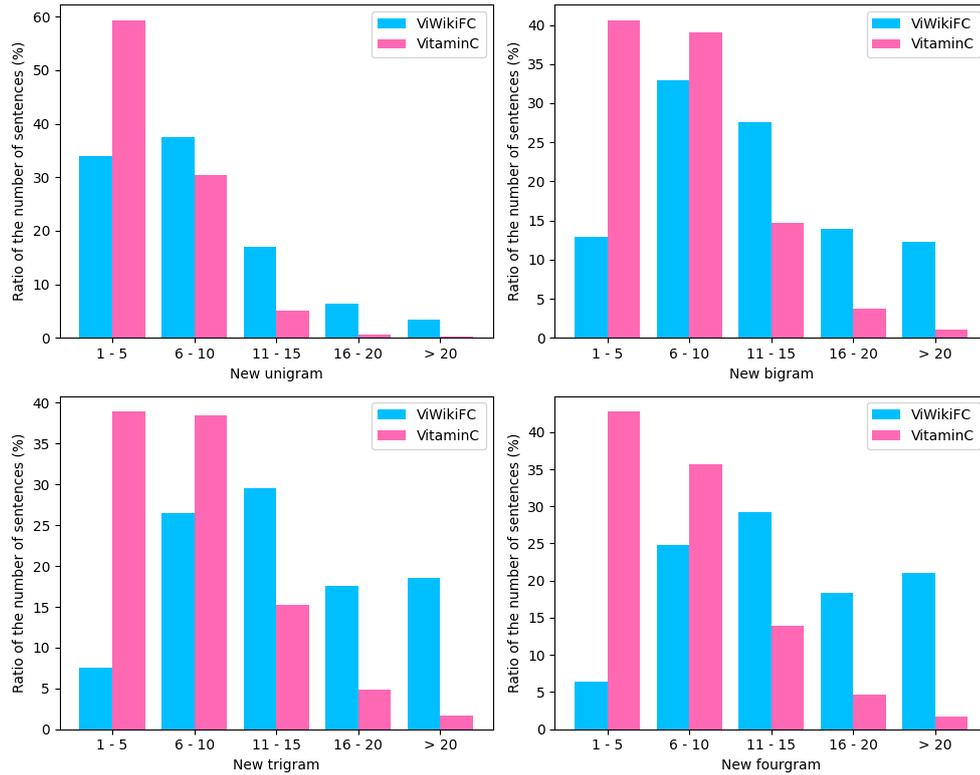


Fig. 7: The ratio of new n-grams number in ViWikiFC and VitaminC corpus.

coverage of the data compared to analyzing a smaller sample of only 100 sentences, as was performed in ViNLI.

For the SUPPORTS label, the rule annotator uses most is 'replace words with synonyms' with 58.60%; meanwhile, the 'turn adjective into relative clause' rule has the lowest ratio (Table 7). As for the REFUTES label, annotators tend to generate claims with the 'wrong reasoning about an event' rule, whereas 'opposite of time' and 'create a sentence that has the opposite meaning of a presupposition' share the same proportion of 5.60%, making them the least common rule used when creating data (Table 8).

During the corpus generation process, annotators can use more than one rule for every claim they write. The analysis we conducted on claims shows that almost half of SUPPORT claims are written using two rules (49.80%), while one rule has 41.20%. For the REFUTES label, claims are nearly generated by one rule with a proportion of 88.00% (Figure 9). This can be explained by the fact that when using too many rules without reviewing carefully, REFUTES claims can be turned into NEI labels (Table 3).

No.	Rule	Ratio (%)
1	Converting active sentences to passive or vice versa	17.40
2	Replacing with synonyms	58.60
3	Adding/Removing modifier without changing the meaning of the evidence sentence	48.00
4	Turning nouns into relative clauses	10.20
5	Turning the object into relative clauses	2.20
6	Turning adjectives into relative clauses	0.20
7	Replacing quantifiers or time with others that have a similar meaning	5.60
8	Creating a presupposition sentence	20.60
9	Replacing Named Entities with a word that stands for the class	2.00
10	Creating conditional sentences	0.40
11	Others	3.00

Table 7: Data-generation rules for creating SUPPORTS claims based on evidence.

No.	Rule	Ratio (%)
1	Using negative words (cannot, hardly, do not,...)	23.00
2	Replacing with antonyms	14.00
3	Opposite of quantity	9.80
4	Opposite of time	5.60
5	Creating a sentence that has the opposite meaning of a presupposition	5.60
6	Wrong reasoning about an object (House, car, river, sea, person, etc.)	23.60
7	Wrong reasoning about an event	30.80
8	Others	0.20

Table 8: Data-generation rules for creating REFUTED claims based on evidence.

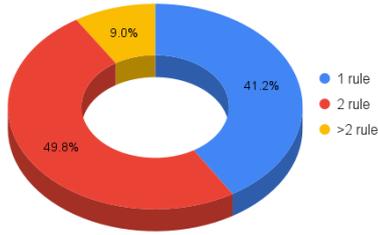


Fig. 8: The ratio of SUPPORT claim generation rules.

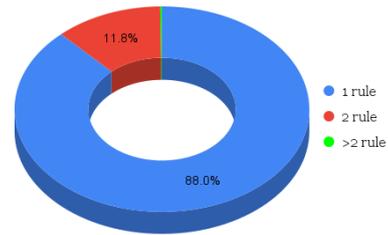


Fig. 9: The ratio of REFUTE claim generation rules.

4 Empirical Evaluation

4.1 Baseline Models

To assess the challenge of ViWikiFC, we define two tasks for a comprehensive corpus evaluation: evidence retrieval and verdict prediction (VP). Moreover, we integrate these two tasks into a complete pipeline. Each task is evaluated on the development set, and the final prediction results are based on the test set.

4.1.1 Evidence Retrieval

For the evidence retrieval (ER) task, we give a claim sentence as input, then employ TF-IDF, Okapi BM25, and Vietnamese Sentence-BERT [10] to select the top one to ten sentences with the highest relevance to the claim sentence from the extracted corpus. With TF-IDF, we build the algorithm from scratch, while with BM25, we use the Rank-BM25¹¹ library with pre-set hyperparameters.

Vietnamese Sentence-BERT [10] is a Vietnamese Sentence embedding model using sentence-BERT [58] architecture. Sentence-BERT [58] is a variant of the pre-trained model BERT [13] with a siamese and triplet network structure, which enhances the ability to search and compare the correlation between sentences with reduced computational time and cost compared to other pre-trained models such as BERT [13] and RoBERTa [49].

4.1.2 Verdict Prediction

In the VP task, we conducted experiments using models ranging from simple approaches like Continuous Bag of Words (CBOW) to more complex ones like BiLSTM [12]. These two models are used with the pre-trained PhoW2V embedding [59] for Vietnamese. We applied two 300-dimensional versions of PhoW2V, including the syllable and word levels. Furthermore, we utilize state-of-the-art multilingual pre-trained language models such as Multilingual BERT (mBERT) [13], XLM-RoBERTa [14], InfoXLM [15], as well as monolingual pre-trained language models for Vietnamese such as PhoBERT [16], ViT5 [17], and ViDeBERTa [18].

The six pre-trained models, mBERT, XLM-R, PhoBERT, InfoXLM, ViT5, and ViDeBERTa, perform highly on various natural language processing (NLP) tasks, including verdict prediction. mBERT, XLM-R, and InfoXLM are multilingual models pre-trained on a vast amount of data encompassing over 90 languages, including Vietnamese. PhoBERT, on the other hand, is a state-of-the-art pre-trained language model for Vietnamese, trained on 20GB Vietnamese Wikipedia and news text corpus, with 135 million parameters for the base version, 370 million parameters for the large version, and recently, PhoBERT base-v2 version have been published which trained on 120GB of Vietnamese texts from OSCAR-2301¹². ViT5 follows encode-decoder transformer architecture [48] with T5 [51] self-supervised pertaining framework. ViDeBERTa is pre-trained with DeBERTa architecture [52] trained on a 138GB Vietnamese dataset.

A characteristic of the Vietnamese language is that a word can be composed of various syllables, so to evaluate various aspects of the Vietnamese language in this task, we conducted two experiments using word-based (PhoBERT, ViDeBERTa) and syllable-based (mBERT, XLM-R, InfoXLM, ViT5) inputs for the models. Word-based models are used with the VNCORENLP toolkit [56] for word segmentation. Furthermore, to mitigate the influence of the experiment results that vary due to initialization seeds, we conducted the experiments five times with randomly selected seeds. We calculated the standard deviation of the results obtained.

4.1.3 Pipeline Description

Finally, we selected models that achieved the highest performance in both tasks, evidence retrieval and verdict prediction. To ensure the best rigor and accuracy for the pipeline, we use

¹¹<https://pypi.org/project/rank-bm25/>

¹²<https://huggingface.co/datasets/oscar-corpus/OSCAR-2301>

strict accuracy measurement where the input is a sentence (claim), and the output consists of a sentence (evidence) extracted from the corpus along with a label of the claim validated based on information within the evidence sentence. The pipeline will count as incorrect output whenever the evidence or label is incorrect. Strict accuracy is described below, with v and v' as the ground truth and predicted evidence in the evidence retrieval task, e and e' as the ground truth and predicted label in verdict prediction ($e, e' \in \text{SUPPORTS, REFUTES, NOTENOUGHINFORMATION}$)

$$\text{Strict Accuracy} = \delta(v, v') \times \delta(e, e')$$

4.1.4 Evaluation Metrics

To assess the performance of the models in these tasks, we employ accuracy as the primary evaluation metric. Additionally, we utilize the F_1 score (macro average) as the secondary evaluation metric for the verdict prediction task (VP), which is described below:

$$F_1 = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

4.2 Experimental Results

4.2.1 Evidence Retrieval

The experimental results presented in Table 9 shed light on the performance of models in the evidence retrieval (ER) task with pre-trained model Vietnamese-SBERT (SBERT), TF-IDF, and BM25. These results offer valuable information on the strengths and weaknesses of these approaches.

We can observe that BM25 consistently outperforms TF-IDF and SBERT across all 3 labels in both single-sentence and 5-sentence retrieval scenarios. As the number of retrieved sentences increases from 1 to 5, the accuracy of all 3 models increases, with BM25 achieving the highest accuracy, reaching 93.93% for the SUPPORTS label and 93.04% for the REFUTES label. SBERT, with its transformer architecture, possesses the ability to capture the contextual nuances of the sentences it retrieves, but in the ViWikiFC corpus, the model’s performance underperforms BM25. This suggests that BM25, despite its less complex architecture, exhibits remarkable performance in evidence retrieval tasks compared to other transformer-based methods [60].

However, it is important to address the relatively lower performance of TF-IDF, SBERT, and BM25 in the NEI (Not Enough Information) label category. This discrepancy in performance can be attributed to the inherent challenges in dealing with the NEI label. On the NEI label, the information contained within a claim may be entirely unrelated to the available evidence. Therefore, accurately identifying and classifying such cases is inherently more complex. TF-IDF, SBERT, and BM25 struggle with this label due to the inherent ambiguity and lack of clear semantic connections between the claim and evidence.

4.2.2 Verdict Prediction

The analysis of the metrics clearly indicates that transformer-based models have emerged as the top performers in the task at hand. Their ability to capture complex linguistic patterns and relationships within the data sets them apart from other models.

ER	Top 1			Top 5		
	SUP	REF	NEI	SUP	REF	NEI
TF-IDF	61.47	58.22	34.26	83.67	81.46	54.58
SBERT	80.19	71.57	43.61	89.41	84.69	60.10
BM25	88.30	86.93	56.57	93.93	93.04	71.51

Table 9: Evidence Retrieval Results using accuracy (%).

In the case of the "word" corpus, PhoBERT_{Large} stands out as the best-performing model. It achieves an impressive accuracy of 82.88% on the development set and 81.63% on the test set. These results demonstrate the robustness and effectiveness of PhoBERT_{Large} in handling the nuances of the "word" corpus. Its high accuracy in both development and test sets suggests that it can consistently provide reliable predictions in this specific context.

However, when we shift our focus to the "syllable" corpus, InfoXLM_{Large} takes the lead with remarkable accuracy and F₁ scores. It achieves an accuracy of 87.45% and 86.50% on validation and test sets, respectively. These numbers not only surpass the performance of PhoBERT_{Large} but also highlight the exceptional capabilities of InfoXLM_{Large} in handling the "syllable" corpus. This outcome establishes InfoXLM_{Large} as the optimal model for the VP task. The superiority of InfoXLM_{Large} in the "syllable" corpus might be attributed to its ability to capture the intricate linguistic patterns and syllable-level features, which are evidently critical in achieving high accuracy and F₁ scores in this particular context.

Model	Dev		Test		
	Acc	F ₁	Acc	F ₁	
Word	CBoW	51.26±0.66	50.50±0.98	50.39±1.05	49.65±0.95
	BiLSTM	49.91±0.47	49.83±0.61	49.13±0.45	49.24±0.48
	ViDeBERTa _{Small}	61.67±1.13	61.43±0.79	60.62±0.66	61.73±0.87
	ViDeBERTa _{Base}	64.41±0.57	64.11±0.54	63.70±0.51	64.06±0.40
	PhoBERT _{Base}	81.72±0.57	81.54±0.58	80.67±0.37	80.70±0.34
	PhoBERT _{Base} V2	82.50±0.52	82.28±0.53	81.22±0.53	81.19±0.56
	PhoBERT _{Large}	82.88±0.97	82.74±1.01	81.63±0.47	81.62±0.44
Syllable	CBoW	45.25±0.66	44.55±0.63	44.24±1.05	43.46±0.64
	BiLSTM	47.68±0.60	47.67±0.77	45.12±0.62	45.04±0.79
	mBERT	75.75±0.68	75.61±1.24	75.93±0.58	76.01±1.18
	ViT5 _{Base}	77.65±0.72	76.70±0.27	77.50±0.67	76.60±0.26
	ViT5 _{Large}	82.77±0.55	82.10±0.42	82.68±0.58	82.16±0.54
	XLM-R _{Base}	79.66±0.95	79.48±0.97	78.53±0.41	78.56±0.48
	XLM-R _{Large}	86.09±0.72	86.02±0.67	85.11±0.27	85.15±0.26
	InfoXLM _{Base}	80.10±0.57	79.92±0.61	79.19±0.78	79.23±0.77
	InfoXLM _{Large}	87.45±0.41	87.42±0.60	86.50±0.43	86.51±0.62

Table 10: The result of models on Dev and Test.

4.2.3 Pipeline

We combine BM25 and SBERT, respectively, with the best models in the VP task, PhoBERT_{Large} and InfoXLM_{Large} for VP. The pipeline is evaluated through the test set. Table 11 gives the result of the pipeline; the best combination for both tasks is BM25 and

	Pipeline	ER Acc	VP Acc	Stric Acc
SBERT	InfoXLM _{Large}	65.56	75.80	57.15
	PhoBERT _{Large}	65.56	71.54	52.84
BM25	InfoXLM _{Large}	78.38	80.63	67.00
	PhoBERT _{Large}	78.38	76.18	63.46

Table 11: Pipeline result for test set

InfoXLM_{Large} with 67.00%, while SBERT and PhoBERT_{Large} have the lowest result. This result demonstrates ViWikiFC corpus’s challenge for current Vietnamese fact-checking models. Besides, it provides a practical perspective on the fact-checking problem, as accurate evidence retrieval significantly influences the performance of evidence-based methods such as PhoBERT and InfoXLM. While these methods perform well in verdict prediction tasks, their strict accuracy still needs to be improved. Therefore, approaching the Vietnamese fact-checking problem using separate models for each task, first for evidence retrieval and then for verdict prediction, needs improvement to deals with more complex claim that have been studied in English [4, 61].

4.3 Result Analysis

4.3.1 Evidence Retrieval

Figure 10 illustrates how model accuracy changes in the ER task as the number of new words increases. The results consistently show a decline in models performance as the count of unfamiliar words in various word categories rises. In particular, we can observe in the chart depicting the number of new preposition that the BM25 model significantly decreases from more than 70.00% to approximate 0% as we increase the number of preposition.

When we introduce a larger number of new words, especially a preposition, there can be significant changes in the context of the sentence. This can pose challenges to both BM25 and SBERT model in accurately determining the context for the information that needs to be retrieved, potentially leading to the return of unrelated or inaccurate sentences and reducing the accuracy of the model.

Furthermore, most sentences that contain many new words belong to the NEI label (see Table 5). During the process of generating NEI-labeled claims, annotators use information and details that are not related to the evidence, making it challenging for the model to retrieve within the scope of a single sentence accurately.

4.3.2 Verdict Prediction

To gain a clearer understanding of the data in the verdict prediction task and approach the language from different aspects, we utilize the two best models, PhoBERT_{Large} and InfoXLM_{Large}. In this section, InfoXLM and PhoBERT respectively replace PhoBERT_{Large} and InfoXLM_{Large}.

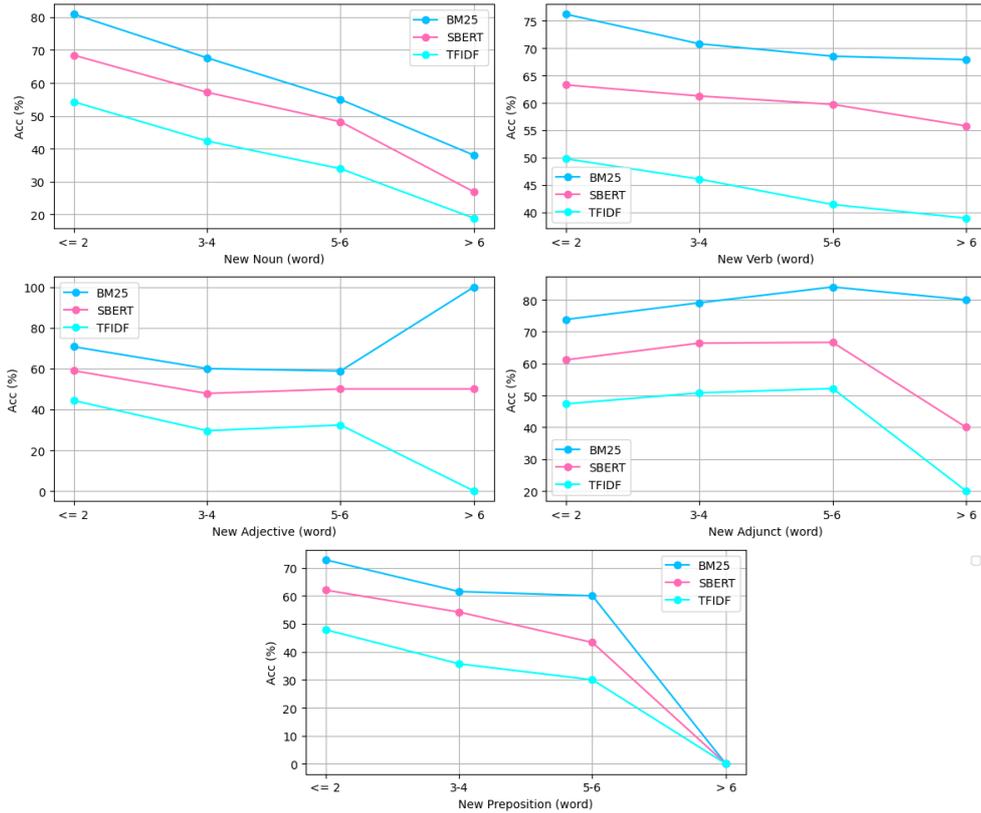


Fig. 10: The effect of new POS tagging for ER task.

We manually determined the number of rules employed to create each sentence (as mentioned in section 3.6.6).

Label	Rule	PhoBERT	InfoXLM
SUPPORTS	1	79.25	88.35
	>1	82.04	89.76
REFUTES	1	80.68	82.95
	>1	81.67	80.00

Table 12: Impacts of data-generation rules on the models.

As we can see in Table 12, SUPPORT claims that are constructed using only one rule can indeed pose challenges for the model when predicting the correct label. As for the REFUTES label, similar to the SUPPORTS label, the predictions are divergent between the two models when applied to two types of different label quantities. For PhoBERT, making substantial errors in information tends to make it easier for models to detect discrepancies compared to

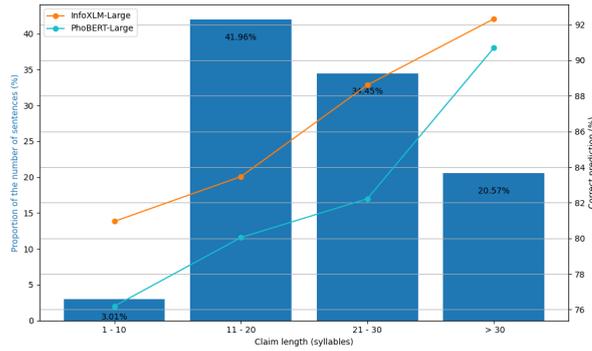


Fig. 11: The effect of claim length

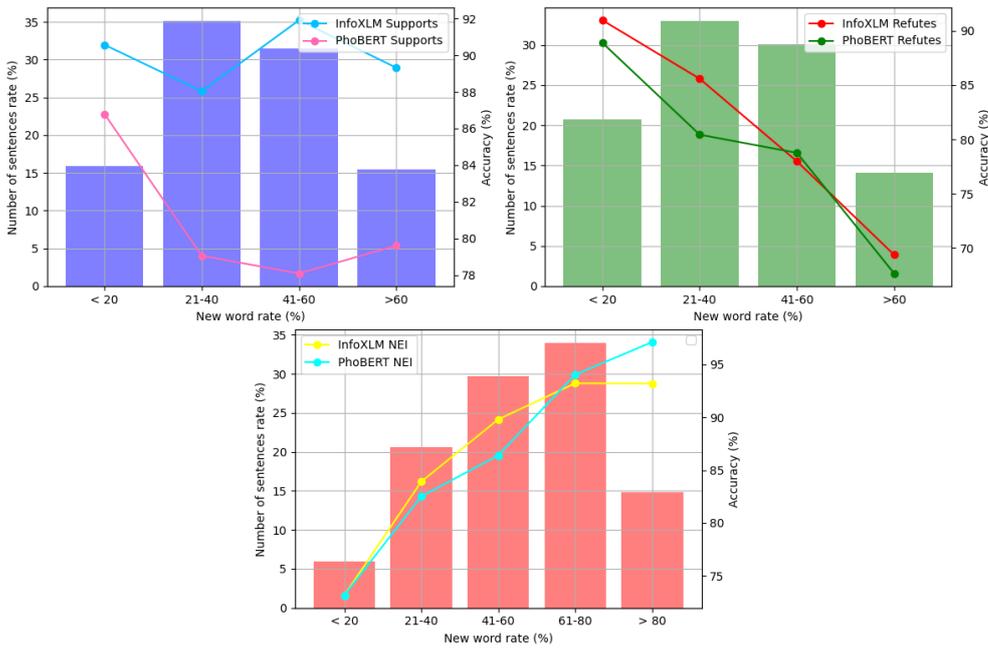


Fig. 12: The effect of new word rate for each label in verdict prediction task

cases where they need to identify just one incorrect piece of information between the claim and the original evidence.

We also aim to gain insights into how sentence length affects the prediction results of the model. Figure 11 illustrates how the accuracy of the models changes as we increase the length of the sentences. The accuracy of the model increases as the sentences become longer. It is apparent that there is a positive correlation between sentence length and model accuracy. The trend suggests that the model benefits from having more context to make accurate predictions.

Longer sentences provide additional context and information, allowing the model to understand the context of the text better and, consequently, make more precise predictions.

In general, for both the SUPPORTS and REFUTES labels, the accuracy of the models tends to decrease as the percentage of new words increases (see Figure 12). Surprisingly, in the SUPPORTS label, when the percentage of new words is in the range of 41.00-60.00%, InfoXLM yields slightly better results compared to the previous range of 21.00-40.00%. PhoBERT has significantly lower accuracy than InfoXLM in correctly predicting the SUPPORTS label in the range of 4.00-13.00%. The REFUTES label shows a clear downward trend with InfoXLM as the percentage of new words in the claim sentence increases. However, as the percentage of 41.00-60.00% new words increases, it becomes evident that PhoBERT predicts the label slightly better than InfoXLM. The NEI label is characterized by a high percentage of new words, so it is not surprising that as the number of new words increases, the model’s prediction accuracy also increases significantly. In terms of accuracy, we can observe that InfoXLM performs better than PhoBERT for percentages of new words below 60.00%. Afterward, both models exhibit similar performance in the remaining portion of the dataset.

5 Conclusions and Future Directions

In this paper, we construct the first Wikipedia Vietnamese corpus, open-domain and high-quality, including 20,916 samples to evaluate fact-checking models. Furthermore, we have discussed the data collection and annotation method and shared the insight we obtained during the annotation process, which can be applied in the development of non-English corpus creation. The best baseline and pipeline model (InfoXLM and BM25) only achieved 67.00% strict accuracy, making our feasible challenge for the Vietnamese language model in fact-checking tasks. We believe that ViWikiFC will encourage the development of Vietnamese fact-checking research.

Following the development of fact-checking in English, we aim to develop our research for Vietnamese fact-checking by expanding our corpus from quantity to quality with more trustworthy data sources and constructing datasets not only on textual but also on images and tables for developing multi-modal fact-checking. Besides, we want our work to adapt to more advanced NLP tasks in low-resource language, such as fake news detection or machine reading comprehension, with advanced reasoning design and sub-task functionalities.

Acknowledgement

This research was supported by The VNUHCM-University of Information Technology’s Scientific Research Support Fund.

Declarations

Conflict of interest The authors declare that they have no conflict of interest.

Data Availability

The corpora generated during and/or analyzed during the current study are available from the corresponding author on reasonable request.

Author Contribution

Hung Tuan Le: Conceptualization; Formal analysis; Investigation; Methodology; Validation; Visualization; Writing - review&editing. Long Truong To: Conceptualization; Data curation; Formal analysis; Investigation; Validation; Visualization; Writing - original draft. Manh Trong Nguyen: Conceptualization; Data curation; Investigation; Methodology; Writing - original draft. Kiet Van Nguyen: Conceptualization; Formal analysis; Investigation; Methodology; Validation; Supervision; Writing - review&editing.

References

- [1] Guo, Z., Schlichtkrull, M., Vlachos, A.: A survey on automated fact-checking. *Transactions of the Association for Computational Linguistics* **10**, 178–206 (2022) https://doi.org/10.1162/tacl_a_00454
- [2] Lewis, P., Perez, E., Piktus, A., Petroni, F., Karpukhin, V., Goyal, N., Küttler, H., Lewis, M., Yih, W.-t., Rocktäschel, T., *et al.*: Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems* **33**, 9459–9474 (2020)
- [3] Maillard, J., Gao, C., Kalbassi, E., Sadagopan, K.R., Goswami, V., Koehn, P., Fan, A., Guzman, F.: Small data, big impact: Leveraging minimal data for effective machine translation. In: *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 2740–2756. Association for Computational Linguistics, Toronto, Canada (2023). <https://doi.org/10.18653/v1/2023.acl-long.154> . <https://aclanthology.org/2023.acl-long.154>
- [4] Pan, L., Wu, X., Lu, X., Luu, A.T., Wang, W.Y., Kan, M.-Y., Nakov, P.: Fact-checking complex claims with program-guided reasoning. In: *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 6981–7004. Association for Computational Linguistics, Toronto, Canada (2023). <https://doi.org/10.18653/v1/2023.acl-long.386> . <https://aclanthology.org/2023.acl-long.386>
- [5] Horne, B., Khedr, S., Adali, S.: Sampling the news producers: A large news and feature data set for the study of the complex media landscape. In: *Proceedings of the International AAAI Conference on Web and Social Media*, vol. 12 (2018)
- [6] Shahi, G.K., Nandini, D.: Fakecovid—a multilingual cross-domain fact check news dataset for covid-19. *arXiv preprint arXiv:2006.11343* (2020)
- [7] Thorne, J., Vlachos, A., Christodoulopoulos, C., Mittal, A.: 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, Volume 1 (Long Papers)*, pp. 809–819. Association for Computational Linguistics, New Orleans, Louisiana (2018). <https://doi.org/10.18653/v1/N18-1074> . <https://aclanthology.org/N18-1074>

- [8] Schuster, T., Fisch, A., Barzilay, R.: Get your vitamin C! robust fact verification with contrastive evidence. In: Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pp. 624–643. Association for Computational Linguistics, Online (2021). <https://doi.org/10.18653/v1/2021.naacl-main.52> . <https://aclanthology.org/2021.naacl-main.52>
- [9] Huynh, T.V., Nguyen, K.V., Nguyen, N.L.-T.: ViNLI: A Vietnamese corpus for studies on open-domain natural language inference. In: Proceedings of the 29th International Conference on Computational Linguistics, pp. 3858–3872. International Committee on Computational Linguistics, Gyeongju, Republic of Korea (2022). <https://aclanthology.org/2022.coling-1.339>
- [10] Phan, Q.L., Doan, T.H.P., Le, N.H., Tran, N.B.D., Huynh, T.N.: Vietnamese sentence paraphrase identification using sentence-bert and phobert. In: Nguyen, N.-T., Dao, N.-N., Pham, Q.-D., Le, H.A. (eds.) Intelligence of Things: Technologies and Applications, pp. 416–423. Springer, Cham (2022)
- [11] Mikolov, T., Sutskever, I., Chen, K., Corrado, G.S., Dean, J.: Distributed representations of words and phrases and their compositionality. *Advances in neural information processing systems* **26** (2013)
- [12] Graves, A., Schmidhuber, J.: Framewise phoneme classification with bidirectional lstm and other neural network architectures. *Neural networks* **18**(5-6), 602–610 (2005)
- [13] Kenton, J.D.M.-W.C., Toutanova, L.K.: Bert: Pre-training of deep bidirectional transformers for language understanding. In: Proceedings of NAACL-HLT, pp. 4171–4186 (2019)
- [14] Conneau, A., Khandelwal, K., Goyal, N., Chaudhary, V., Wenzek, G., Guzmán, F., Grave, E., Ott, M., Zettlemoyer, L., Stoyanov, V.: Unsupervised cross-lingual representation learning at scale. In: Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pp. 8440–8451. Association for Computational Linguistics, Online (2020). <https://doi.org/10.18653/v1/2020.acl-main.747> . <https://aclanthology.org/2020.acl-main.747>
- [15] Chi, Z., Dong, L., Wei, F., Yang, N., Singhal, S., Wang, W., Song, X., Mao, X.-L., Huang, H., Zhou, M.: InfoXLM: An information-theoretic framework for cross-lingual language model pre-training. In: Toutanova, K., Rumshisky, A., Zettlemoyer, L., Hakkani-Tur, D., Beltagy, I., Bethard, S., Cotterell, R., Chakraborty, T., Zhou, Y. (eds.) Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pp. 3576–3588. Association for Computational Linguistics, Online (2021). <https://doi.org/10.18653/v1/2021.naacl-main.280> . <https://aclanthology.org/2021.naacl-main.280>
- [16] Nguyen, D.Q., Tuan Nguyen, A.: PhoBERT: Pre-trained language models for Vietnamese. In: Findings of the Association for Computational Linguistics: EMNLP 2020, pp. 1037–1042. Association for Computational Linguistics, Online (2020). <https://doi.org/10.18653/v1/2020.findings-emnlp.103>

- 18653/v1/2020.findings-emnlp.92 . <https://aclanthology.org/2020.findings-emnlp.92>
- [17] Phan, L., Tran, H., Nguyen, H., Trinh, T.H.: ViT5: Pretrained text-to-text transformer for Vietnamese language generation. In: Ippolito, D., Li, L.H., Pacheco, M.L., Chen, D., Xue, N. (eds.) Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies: Student Research Workshop, pp. 136–142. Association for Computational Linguistics, Hybrid: Seattle, Washington + Online (2022). <https://doi.org/10.18653/v1/2022.naacl-srw.18> . <https://aclanthology.org/2022.naacl-srw.18>
- [18] Tran, C.D., Pham, N.H., Nguyen, A.T., Hy, T.S., Vu, T.: ViDeBERTa: A powerful pre-trained language model for Vietnamese. In: Vlachos, A., Augenstein, I. (eds.) Findings of the Association for Computational Linguistics: EACL 2023, pp. 1071–1078. Association for Computational Linguistics, Dubrovnik, Croatia (2023). <https://doi.org/10.18653/v1/2023.findings-eacl.79> . <https://aclanthology.org/2023.findings-eacl.79>
- [19] Vlachos, A., Riedel, S.: Fact checking: Task definition and dataset construction. In: Proceedings of the ACL 2014 Workshop on Language Technologies and Computational Social Science, pp. 18–22. Association for Computational Linguistics, Baltimore, MD, USA (2014). <https://doi.org/10.3115/v1/W14-2508> . <https://aclanthology.org/W14-2508>
- [20] Wang, W.Y.: “liar, liar pants on fire”: A new benchmark dataset for fake news detection. In: Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), pp. 422–426. Association for Computational Linguistics, Vancouver, Canada (2017). <https://doi.org/10.18653/v1/P17-2067> . <https://aclanthology.org/P17-2067>
- [21] Hanselowski, A., Stab, C., Schulz, C., Li, Z., Gurevych, I.: A richly annotated corpus for different tasks in automated fact-checking. In: Proceedings of the 23rd Conference on Computational Natural Language Learning (CoNLL), pp. 493–503. Association for Computational Linguistics, Hong Kong, China (2019). <https://doi.org/10.18653/v1/K19-1046> . <https://aclanthology.org/K19-1046>
- [22] Augenstein, I., Lioma, C., Wang, D., Chaves Lima, L., Hansen, C., Hansen, C., Simonsen, J.G.: MultiFC: A real-world multi-domain dataset for evidence-based fact checking of claims. 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), pp. 4685–4697. Association for Computational Linguistics, Hong Kong, China (2019). <https://doi.org/10.18653/v1/D19-1475> . <https://aclanthology.org/D19-1475>
- [23] Jiang, Y., Bordia, S., Zhong, Z., Dognin, C., Singh, M., Bansal, M.: HoVer: A dataset for many-hop fact extraction and claim verification. In: Findings of the Association for Computational Linguistics: EMNLP 2020, pp. 3441–3460. Association for Computational Linguistics, Online (2020). <https://doi.org/10.18653/v1/2020.findings-emnlp.309> . <https://aclanthology.org/2020.findings-emnlp.309>

- [24] Chen, W., Wang, H., Chen, J., Zhang, Y., Wang, H., Li, S., Zhou, X., Wang, W.Y.: Tabfact: A large-scale dataset for table-based fact verification. In: International Conference on Learning Representations (2020). <https://openreview.net/forum?id=rkeJRhNYDH>
- [25] Gupta, V., Mehta, M., Nokhiz, P., Srikumar, V.: INFOTABS: Inference on tables as semi-structured data. In: Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pp. 2309–2324. Association for Computational Linguistics, Online (2020). <https://doi.org/10.18653/v1/2020.acl-main.210> . <https://aclanthology.org/2020.acl-main.210>
- [26] Aly, R., Guo, Z., Schlichtkrull, M.S., Thorne, J., Vlachos, A., Christodoulopoulos, C., Cocarascu, O., Mittal, A.: Feverous: Fact extraction and verification over unstructured and structured information. In: Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 1) (2021)
- [27] Gupta, A., Srikumar, V.: X-fact: A new benchmark dataset for multilingual fact checking. In: Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers), pp. 675–682. Association for Computational Linguistics, Online (2021). <https://doi.org/10.18653/v1/2021.acl-short.86> . <https://aclanthology.org/2021.acl-short.86>
- [28] Khouja, J.: Stance prediction and claim verification: An Arabic perspective. In: Proceedings of the Third Workshop on Fact Extraction and VERification (FEVER), pp. 8–17. Association for Computational Linguistics, Online (2020). <https://doi.org/10.18653/v1/2020.fever-1.2> . <https://aclanthology.org/2020.fever-1.2>
- [29] Nørregaard, J., Derczynski, L.: DanFEVER: claim verification dataset for Danish. In: Proceedings of the 23rd Nordic Conference on Computational Linguistics (NoDaLiDa), pp. 422–428. Linköping University Electronic Press, Sweden, Reykjavik, Iceland (Online) (2021). <https://aclanthology.org/2021.nodalida-main.47>
- [30] Hu, X., Guo, Z., Wu, G., Liu, A., Wen, L., Yu, P.: CHEF: A pilot Chinese dataset for evidence-based fact-checking. In: Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pp. 3362–3376. Association for Computational Linguistics, Seattle, United States (2022). <https://doi.org/10.18653/v1/2022.naacl-main.246> . <https://aclanthology.org/2022.naacl-main.246>
- [31] Chen, D., Fisch, A., Weston, J., Bordes, A.: Reading Wikipedia to answer open-domain questions. In: Barzilay, R., Kan, M.-Y. (eds.) Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pp. 1870–1879. Association for Computational Linguistics, Vancouver, Canada (2017). <https://doi.org/10.18653/v1/P17-1171> . <https://aclanthology.org/P17-1171>
- [32] Karpukhin, V., Oguz, B., Min, S., Lewis, P., Wu, L., Edunov, S., Chen, D., Yih, W.-t.: Dense passage retrieval for open-domain question answering. In: Webber, B., Cohn, T.,

- He, Y., Liu, Y. (eds.) Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pp. 6769–6781. Association for Computational Linguistics, Online (2020). <https://doi.org/10.18653/v1/2020.emnlp-main.550> . <https://aclanthology.org/2020.emnlp-main.550>
- [33] Chang, W.-C., Felix, X.Y., Chang, Y.-W., Yang, Y., Kumar, S.: Pre-training tasks for embedding-based large-scale retrieval. In: International Conference on Learning Representations (2019)
- [34] Lee, K., Chang, M.-W., Toutanova, K.: Latent retrieval for weakly supervised open domain question answering. In: Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, pp. 6086–6096. Association for Computational Linguistics, Florence, Italy (2019). <https://doi.org/10.18653/v1/P19-1612> . <https://aclanthology.org/P19-1612>
- [35] Samarinas, C., Hsu, W., Lee, M.L.: Improving evidence retrieval for automated explainable fact-checking. In: Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies: Demonstrations, pp. 84–91. Association for Computational Linguistics, Online (2021). <https://doi.org/10.18653/v1/2021.naacl-demos.10> . <https://aclanthology.org/2021.naacl-demos.10>
- [36] Samarinas, C., Hsu, W., Lee, M.L.: 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), pp. 941–948 (2020). IEEE Computer Society
- [37] Nogueira, R., Cho, K.: Passage re-ranking with bert. arXiv preprint arXiv:1901.04085 (2019)
- [38] Nogueira, R., Yang, W., Cho, K., Lin, J.: Multi-stage document ranking with bert. arXiv preprint arXiv:1910.14424 (2019)
- [39] Wang, Z., Ng, P., Ma, X., Nallapati, R., Xiang, B.: Multi-passage bert: A globally normalized bert model for open-domain question answering. 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), pp. 5878–5882 (2019)
- [40] Luken, J., Jiang, N., Marneffe, M.-C.: QED: A fact verification system for the FEVER shared task. In: Thorne, J., Vlachos, A., Cocarascu, O., Christodoulopoulos, C., Mittal, A. (eds.) Proceedings of the First Workshop on Fact Extraction and VERification (FEVER), pp. 156–160. Association for Computational Linguistics, Brussels, Belgium (2018). <https://doi.org/10.18653/v1/W18-5526> . <https://aclanthology.org/W18-5526>
- [41] Liu, Z., Xiong, C., Sun, M., Liu, Z.: Fine-grained fact verification with kernel graph attention network. In: Jurafsky, D., Chai, J., Schluter, N., Tetreault, J. (eds.) Proceedings

- of the 58th Annual Meeting of the Association for Computational Linguistics, pp. 7342–7351. Association for Computational Linguistics, Online (2020). <https://doi.org/10.18653/v1/2020.acl-main.655> . <https://aclanthology.org/2020.acl-main.655>
- [42] Zhong, W., Xu, J., Tang, D., Xu, Z., Duan, N., Zhou, M., Wang, J., Yin, J.: Reasoning over semantic-level graph for fact checking. In: Jurafsky, D., Chai, J., Schluter, N., Tetreault, J. (eds.) Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pp. 6170–6180. Association for Computational Linguistics, Online (2020). <https://doi.org/10.18653/v1/2020.acl-main.549> . <https://aclanthology.org/2020.acl-main.549>
- [43] Schlichtkrull, M.S., Karpukhin, V., Oguz, B., Lewis, M., Yih, W.-t., Riedel, S.: Joint verification and reranking for open fact checking over tables. In: Zong, C., Xia, F., Li, W., Navigli, R. (eds.) Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pp. 6787–6799. Association for Computational Linguistics, Online (2021). <https://doi.org/10.18653/v1/2021.acl-long.529> . <https://aclanthology.org/2021.acl-long.529>
- [44] Ma, J., Gao, W., Joty, S., Wong, K.-F.: Sentence-level evidence embedding for claim verification with hierarchical attention networks. In: Korhonen, A., Traum, D., Màrquez, L. (eds.) Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, pp. 2561–2571. Association for Computational Linguistics, Florence, Italy (2019). <https://doi.org/10.18653/v1/P19-1244> . <https://aclanthology.org/P19-1244>
- [45] Yoneda, T., Mitchell, J., Welbl, J., Stenetorp, P., Riedel, S.: UCL machine reading group: Four factor framework for fact finding (HexaF). In: Thorne, J., Vlachos, A., Cocarascu, O., Christodoulopoulos, C., Mittal, A. (eds.) Proceedings of the First Workshop on Fact Extraction and VERification (FEVER), pp. 97–102. Association for Computational Linguistics, Brussels, Belgium (2018). <https://doi.org/10.18653/v1/W18-5515> . <https://aclanthology.org/W18-5515>
- [46] Elman, J.L.: Finding structure in time. *Cognitive Science* **14**(2), 179–211 (1990) [https://doi.org/10.1016/0364-0213\(90\)90002-E](https://doi.org/10.1016/0364-0213(90)90002-E)
- [47] Devlin, J., Chang, M.-W., Lee, K., Toutanova, K.: 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, Volume 1 (Long and Short Papers), pp. 4171–4186. Association for Computational Linguistics, Minneapolis, Minnesota (2019). <https://doi.org/10.18653/v1/N19-1423> . <https://aclanthology.org/N19-1423>
- [48] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł., Polosukhin, I.: Attention is all you need. *Advances in neural information processing systems* **30** (2017)
- [49] Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer,

- L., Stoyanov, V.: Roberta: A robustly optimized BERT pretraining approach. *CoRR abs/1907.11692* (2019) [1907.11692](https://arxiv.org/abs/1907.11692)
- [50] Lan, Z., Chen, M., Goodman, S., Gimpel, K., Sharma, P., Soricut, R.: Albert: A lite bert for self-supervised learning of language representations. In: *International Conference on Learning Representations* (2019)
- [51] Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W., Liu, P.J.: Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research* **21**(1), 5485–5551 (2020)
- [52] He, P., Liu, X., Gao, J., Chen, W.: Deberta: Decoding-enhanced bert with disentangled attention. In: *International Conference on Learning Representations* (2020)
- [53] Fleiss, J.L.: Measuring nominal scale agreement among many raters. *Psychological Bulletin* **76**, 378–382 (1971)
- [54] Hu, H., Richardson, K., Xu, L., Li, L., Kuebler, S., Moss, L.: Ocnli: Original chinese natural language inference. In: *Findings of EMNLP* (2020). <https://arxiv.org/abs/2010.05444>
- [55] Mahendra, R., Aji, A.F., Louvan, S., Rahman, F., Vania, C.: IndoNLI: A natural language inference dataset for Indonesian. In: *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 10511–10527. Association for Computational Linguistics, Online and Punta Cana, Dominican Republic (2021). <https://aclanthology.org/2021.emnlp-main.821>
- [56] Vu, T., Nguyen, D.Q., Nguyen, D.Q., Dras, M., Johnson, M.: VnCoreNLP: A Vietnamese natural language processing toolkit. In: *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Demonstrations*, pp. 56–60. Association for Computational Linguistics, New Orleans, Louisiana (2018). <https://doi.org/10.18653/v1/N18-5012> . <https://aclanthology.org/N18-5012>
- [57] Nguyen, L.T., Nguyen, D.Q.: PhoNLP: A joint multi-task learning model for Vietnamese part-of-speech tagging, named entity recognition and dependency parsing. In: *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Demonstrations*, pp. 1–7 (2021)
- [58] Reimers, N., Gurevych, I.: 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)*, pp. 3982–3992 (2019)
- [59] Nguyen, A.T., Dao, M.H., Nguyen, D.Q.: A Pilot Study of Text-to-SQL Semantic Parsing for Vietnamese. In: *Findings of the Association for Computational Linguistics: EMNLP 2020*, pp. 4079–4085 (2020)

- [60] Thakur, N., Reimers, N., Rücklé, A., Srivastava, A., Gurevych, I.: Beir: A heterogeneous benchmark for zero-shot evaluation of information retrieval models. In: Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2) (2021)
- [61] Yao, B.M., Shah, A., Sun, L., Cho, J.-H., Huang, L.: End-to-end multimodal fact-checking and explanation generation: A challenging dataset and models. In: Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval, pp. 2733–2743 (2023)

A Data Generation Rules Analysis

Rule	Example	Ratio
Change active sentences into passive sentences and vice versa	<p>E: Âu Lạc bị nhà Triệu ở phương Bắc thôn tính vào đầu thế kỷ thứ 2 TCN sau đó là thời kỳ Bắc thuộc kéo dài hơn một thiên niên kỷ. (<i>Au Lac was conquered by the Zhaos in the northern region at the beginning of the 2nd century BC, leading to the Northern Domination period that lasted for over a millennium.</i>)</p> <p>C: Thời kỳ Bắc thuộc diễn ra sau khi phương Bắc thôn tính được Âu Lạc. (<i>The Northern Domination period took place after the northern regions conquered Au Lac.</i>)</p>	17.40%
Replace words with synonyms.	<p>E: Thời kỳ đầu, những bậc đế vương và những nhà quý tộc của La Mã thích lụa Trung Hoa đến mức họ cho cân lụa lên và đổi chỗ lụa đó bằng vàng với cân nặng tương đương. (<i>In the early periods, Roman emperors and nobles had such a fondness for Chinese silk that they would weigh silk against gold and exchange it with an equivalent weight of gold.</i>)</p> <p>C: Lụa Trung Hoa nhận được sự ưa chuộng cực lớn từ những người đứng đầu, những người có địa vị ở La Mã. (<i>Chinese silk received immense favor from the Roman leaders and influential people.</i>)</p>	58.60%
Add or remove modifiers that do not radically alter the meaning of the sentence.	<p>E: Gary Shilling, chủ tịch một công ty nghiên cứu kinh tế, cho rằng mức tăng trưởng GDP thực sự của Trung Quốc chỉ là 3.50% chứ không phải 7.00% như báo cáo chính thức. (<i>Gary Shilling, the chairman of an economic research firm, believes that China's actual GDP growth rate is only 3.50%, not the officially reported 7.00%.</i>)</p> <p>C: Báo cáo chính thức về mức tăng trưởng thực sự của Trung Quốc đã được Gary Shilling cho rằng là 3.50%. (<i>The official report on China's actual growth rate has been suggested by Gary Shilling to be 3.50%.</i>)</p>	48.00%
Turn nouns into relative clauses	<p>E: Trong bốn nước xưa có nền văn minh lớn thì có ba nước xưa ở vào châu Á (Ấn Độ, Iraq (Lưỡng Hà) và Trung Quốc). (<i>Among the four ancient civilizations, three were located in Asia (India, Iraq (Mesopotamia), and China).</i>)</p> <p>C: Châu Á là nơi tồn tại ba nước có nền văn minh lớn. (<i>Asia is the region where three major ancient civilizations existed.</i>)</p>	10.20%

Rule	Example	Ratio
Turn the object into relative clauses	<p>E: Lý thuyết "Quả cầu tuyết Trái Đất" cho rằng những sự thay đổi về mức độ CO2 vừa là nguyên nhân gây ra, vừa là nguyên nhân làm kết thúc thời kỳ cực lạnh ở cuối Liên đại Nguyên Sinh (Proterozoic). (<i>The "Snowball Earth" theory posits that changes in the levels of CO2 were both the cause and the termination of the extreme cold period at the end of the Proterozoic Eon.</i>)</p> <p>C: Theo lý thuyết "Quả cầu tuyết Trái Đất" cho thấy rằng sự thay đổi thành phần không khí mà cụ thể là lưu lượng CO2 là tác nhân cho kỷ băng giá, và cũng kết thúc luôn thời kỳ cực lạnh ở kỷ Proterozoic. (<i>The "Snowball Earth" theory indicates that changes in the composition of the atmosphere, specifically the levels of CO2 which are the driving factors behind ice ages and the termination of extreme cold periods in the Proterozoic Eon.</i>)</p>	2.20%
Turn adjectives into relative clauses	<p>E: Về du lịch biển, Nghệ An có 82 km bờ biển với nhiều bãi tắm đẹp hấp dẫn khách du lịch quốc tế như bãi biển Cửa Lò, Cửa Hội; Nghi Thiết,... (<i>For beach tourism, Nghệ An has a coastline of 82 kilometers with numerous beautiful and attractive beaches for international tourists, such as Cua Lo, Cua Hoi, Nghi Thiet,...</i>)</p> <p>C: Nghệ An có các bãi tắm, nơi được xem là những địa điểm tuyệt đẹp và rất cuốn hút du khách quốc tế. (<i>Nghe An has beaches which are considered to be stunning and highly attractive destinations for international tourists.</i>)</p>	0.20%
Replace quantifiers or time with others that have a similar meaning.	<p>E: Theo Viện Hàn lâm Khoa học Nga, Liên Xô đã chịu 26,6 triệu thương vong trong chiến tranh thế giới thứ hai, bao gồm sự gia tăng tỷ lệ tử vong ở trẻ sơ sinh là 1,3 triệu. (<i>According to the Russian Academy of Sciences, the Soviet Union suffered 26.6 million casualties during World War II, including an increase in the infant mortality rate of 1.3 million.</i>)</p> <p>C: Liên Xô chứng kiến gần 27 triệu cái chết trong thế chiến thứ hai. (<i>The Soviet Union witnessed nearly 27 million deaths during World War II.</i>)</p>	5.60%
Create a pre-supposition sentence	<p>E: Một số người cho rằng Firenze trở thành nơi khởi đầu Phục Hưng là do may mắn, nghĩa là đơn thuần bởi vì những vĩ nhân ngẫu nhiên sinh ra ở đây: cả da Vinci, Botticelli và Michelangelo đều là người xứ Toscana (mà Firenze là thủ phủ). (<i>Some argue that Firenze became the cradle of the Renaissance due to luck, meaning it was purely coincidental that such great talents were born there: both da Vinci, Botticelli, and Michelangelo were natives of Tuscany (with Firenze as its capital).</i>)</p> <p>C: Quê hương của Michelangelo là Firenze. (<i>Michelangelo's hometown is in Firenze.</i>)</p>	20.60%

Rule	Example	Ratio
Replace Named Entities with a word that stands for the class.	<p>E: Các công ty công nghệ cao của Trung Quốc như Lenovo, Huawei, Xiaomi, Coolpad, ZTE,... đã bắt đầu cạnh tranh thành công trên thị trường thế giới. (<i>High-tech companies in China such as Lenovo, Huawei, Xiaomi, Coolpad, ZTE,... have successfully begun competing in the global market.</i>)</p> <p>C: Nhiều công ty công nghệ của Trung Quốc đã bắt đầu ghi danh mình thành công trên thị trường công nghệ cao thế giới. (<i>Many Chinese high-tech companies have indeed successfully made a name for themselves in the global high-tech market.</i>)</p>	2.00%
Create conditional sentences	<p>E: Bệnh viêm gan siêu vi C mạn được xác định là nhiễm siêu vi viêm gan C hơn 6 tháng căn cứ trên sự hiện diện của ARN. (The diagnosis of chronic hepatitis C infection is established when the presence of RNA is detected for more than 6 months.)</p> <p>C: Nếu nhiễm siêu vi gan C hơn nửa năm, bạn sẽ bị viêm gan siêu vi C mạn. (<i>If you are infected with the hepatitis C virus for more than half a year, you will be considered to have chronic hepatitis C.</i>)</p>	0.40%
Other	<p>E: Những phân tích lõi băng và lõi trầm tích đại dương không chứng minh rõ ràng sự hiện diện của băng giá và những thời kỳ trung gian băng giá trong vòng vài triệu năm qua. (<i>Core ice and sediment analyses from the ocean do not provide clear evidence of the presence of ice ages and intermediate ice-free periods over the past several million years.</i>)</p> <p>C: Dù đều là các bằng chứng địa chất nhưng theo những phân tích lõi băng và lõi trầm tích đại dương đã không chứng minh rõ ràng sự hiện diện của băng giá và những thời kỳ trung gian băng giá trong vòng vài triệu năm qua. (<i>Although both are geological pieces of evidence, core ice and ocean sediment analyses have not provided clear confirmation of the presence of ice ages and intermediate ice-free periods over the past several million years.</i>)</p> <p>Note: Although clause, Inversion can be used for converting evidence to claim. These circumstances are listed in Other.</p>	3.00%

Table 13: SUPPORTS rules and examples for creating evidence (E) - claim (C) pairs. Simply, we only mention one rule to be applied in each example.

Rule	Example	Ratio
Use negative words (no, not, never, nothing, hardly, etc.)	<p>E: Sự kiện này dẫn tới việc Hiệp định Genève (1954) được ký kết và Việt Nam bị chia cắt thành hai vùng tập kết quân sự, lấy ranh giới là vĩ tuyến 17. (This event led to the signing of the Geneva Accords in 1954, which divided Vietnam into two military zones along the 17th parallel.)</p> <p>C: Sự kiện này dẫn tới việc Hiệp định Genève (1954) không được ký kết. (<i>This event led to the Geneva Accords (1954) not being signed.</i>)</p>	23.00%
Replace words with antonyms	<p>E:Đội tuyển bóng nước Singapore đã giành huy chương vàng SEA Games lần thứ 27 vào năm 2017, tiếp tục chuỗi vô địch dài nhất của thể thao Singapore về môn môn cụ thể. (<i>The Singapore water polo team won the gold medal at the 27th SEA Games in 2017, continuing Singapore's longest winning streak in a specific sport.</i>)</p> <p>C:Singapore liên tục gặp thất bại ở bộ môn bóng nước tại các kỳ SEA Games. (<i>Singapore has continuously faced defeats in the sport of water polo at various SEA Games.</i>)</p>	14.00%
Opposite of quantity	<p>E: Năm 2010, tổng chi tiêu của Nhà nước vào khoa học và công nghệ chiếm khoảng 0.45% GDP. (<i>In 2010, the government's total expenditure on science and technology accounted for approximately 0.45% of the GDP.</i>)</p> <p>C: Nhà nước chi ra hơn 2.00% GDP năm 2010 cho khoa học và công nghệ. (<i>The government allocated over 2.00% of the GDP in 2010 for science and technology.</i>)</p>	9.80%
Opposite of time	<p>E: Năm 1860, dân số Singapore đã vượt quá 80,000 và hơn một nửa là người Hoa. (<i>In 1860, the population of Singapore had exceeded 80,000, with over half being of Chinese descent.</i>)</p> <p>C: Mãi đến năm 2000, Dân số Singapore mới đạt mốc 80,000. (<i>It wasn't until the year 2000 that Singapore's population reached the 80,000 milestone.</i>)</p>	5.60%
Create a sentence that has the opposite meaning of a presupposition	<p>E: Đức đã tung ra 70.00% binh lực với các sư đoàn mạnh và tinh nhuệ nhất, chưa kể binh lực góp thêm của các nước đồng minh của Đức (Ý, Rumani, Bulgari, Hungary, Phần Lan...). (<i>Germany deployed approximately 70.00% of its military power, not counting the additional military contributions from Germany's allies (Italy, Romania, Bulgaria, Hungary, Finland...).</i>)</p> <p>C: Đức không có bất kỳ đồng minh nào giúp sức. (<i>Germany did not have any allies providing assistance.</i>)</p>	5.60%

Rule	Example	Ratio
Wrong reasoning about an object (House, car, river, sea, person, etc.)	<p>E: Khi gió mùa đổi hướng, các đường bờ biển giáp với biển Ả Rập và vịnh Bengal có thể phải hứng chịu xoáy thuận. (<i>When the monsoon winds change direction, coastal areas bordering the Arabian Sea and the Bay of Bengal may be susceptible to cyclones.</i>)</p> <p>C: Các đường bờ biển tiếp giáp Ấn Độ có thể phải hứng chịu xoáy thuận khi gió mùa đổi hướng. (<i>The coastal areas bordering India may be susceptible to cyclones when the monsoon winds change direction.</i>)</p>	23.60%
Wrong reasoning about an event	<p>E: Một trận động đất ở Valdivia, Chile với cường độ 9,4-9,6 độ richter, mức cao nhất từng được ghi nhận, khiến 1.000 đến 6.000 người chết. (<i>An earthquake in Valdivia, Chile, with a magnitude of 9.4 to 9.6 on the Richter scale, the highest ever recorded, resulted in the deaths of approximately 1,000 to 6,000 people.</i>)</p> <p>C: Cơn địa chấn ở Valdivia chỉ gây ra thiệt hại về tài sản. (<i>The earthquake in Valdivia only caused property damage.</i>)</p>	30.80%
Others	<p>E: Tiếng Pháp không phải là ngôn ngữ chính thức ở Ontario, nhưng Đạo luật Dịch vụ Ngôn ngữ Pháp đảm bảo rằng các dịch vụ của tỉnh bang sẽ được cung cấp bằng ngôn ngữ này. (<i>French is not the official language in Ontario, but the French Language Services Act ensures that provincial services will be provided in this language.</i>)</p> <p>C: Dù tiếng Pháp có là ngôn ngữ chung của bang Ontario, thì nó cũng sẽ bị loại bỏ khỏi hệ thống các dịch vụ. (<i>Even though French is one of the official languages of the province of Ontario, it will be removed from the system of services.</i>)</p>	0.20%

Table 14: REFUTES rules and examples for creating evidence (E) - claim (C) pairs. , we only mention one rule to apply in each example.