

New Vietnamese Corpus for Machine Reading Comprehension of Health News Articles

KIET VAN NGUYEN, University of Information Technology, VNU-HCM, Vietnam

TIN VAN HUYNH, University of Information Technology, VNU-HCM, Vietnam

DUC-VU NGUYEN, University of Information Technology, VNU-HCM, Vietnam

ANH GIA-TUAN NGUYEN, University of Information Technology, VNU-HCM, Vietnam

NGAN LUU-THUY NGUYEN, University of Information Technology, VNU-HCM, Vietnam

Large-scale and high-quality corpora are necessary for evaluating machine reading comprehension models on a low-resource language like Vietnamese. Besides, machine reading comprehension (MRC) for the health domain offers great potential for practical applications; however, there is still very little MRC research in this domain. This paper presents ViNewsQA as a new corpus for the Vietnamese language to evaluate healthcare reading comprehension models. The corpus comprises 22,057 human-generated question-answer pairs. Crowd-workers create the questions and their answers based on a collection of over 4,416 online Vietnamese healthcare news articles, where the answers comprise spans extracted from the corresponding articles. In particular, we develop a process of creating a corpus for the Vietnamese machine reading comprehension. Comprehensive evaluations demonstrate that our corpus requires abilities beyond simple reasoning, such as word matching and demanding difficult reasoning based on single-or-multiple-sentence information. We conduct experiments using different types of machine reading comprehension methods to achieve the first baseline performances, compared with further models' performances. We also measure human performance on the corpus and compared it with several powerful neural network-based and transfer learning-based models. Our experiments show that the best machine model is ALBERT, which achieves an exact match score of 65.26% and a F1-score of 84.89% on our corpus. The significant differences between humans and the best-performance model (14.53% of EM and 10.90% of F1-score) on the test set of our corpus indicates that improvements in ViNewsQA could be explored in the future study. Our corpus is publicly available on our website¹ for the research purpose to encourage the research community to make these improvements.

CCS Concepts: • **Computing Methodologies** → **Language resources**; • **Information systems** → **Machine Reading Comprehension**.

Additional Key Words and Phrases: Machine Reading Comprehension, Question Answering, Vietnamese

¹<https://sites.google.com/uit.edu.vn/uit-nlp/datasets-projects>

Authors' addresses: Kiet Van Nguyen, University of Information Technology, VNU-HCM, Vietnam; Tin Van Huynh, University of Information Technology, VNU-HCM, Vietnam; Duc-Vu Nguyen, University of Information Technology, VNU-HCM, Vietnam; Anh Gia-Tuan Nguyen, University of Information Technology, VNU-HCM, Vietnam; Ngan Luu-Thuy Nguyen, University of Information Technology, VNU-HCM, Vietnam.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2020 Association for Computing Machinery.

XXXX-XXXX/2020/2-ART \$15.00

<https://doi.org/10.1145/1122445.1122456>

ACM Reference Format:

Kiet Van Nguyen, Tin Van Huynh, Duc-Vu Nguyen, Anh Gia-Tuan Nguyen, and Ngan Luu-Thuy Nguyen. 2020. New Vietnamese Corpus for Machine Reading Comprehension of Health News Articles. 1, 1 (February 2020), 41 pages. <https://doi.org/10.1145/1122445.1122456>

1 INTRODUCTION

Question answering (QA) systems have recently achieved considerable success in a range of benchmark corpora due to the powerful development of neural network-based [1, 2, 3] QA systems. Modern QA systems have two main components [1], where the first component for information retrieval selects text passages that appear relevant to questions from the corpus, and the second component for machine reading comprehension extracts answers that are then returned to the user. Machine Reading Comprehension (MRC) is a natural language understanding task that requires computers to understand human languages and answer questions by reading a given document. Human annotation for large-scale corpora is laborious and time-consuming, but more qualitative than data generation by automated method. Therefore, it is the best option for building many high-quality datasets such as SQuAD [4] and CMRC [5]. In order to evaluate MRC models, gold-standard resources comprising document-question-answer triples have to be collected and annotated by humans. Therefore, creating a benchmark corpus is vital for human language processing, especially for low-resource languages such as Vietnamese.

In recent years, researchers have developed many MRC corpora and models in popular languages such as English and Chinese. The best-known examples of gold standard MRC resources for English are span-extraction MRC corpora [4, 6, 7], cloze-style MRC corpora [8, 9, 10], reading comprehension with multiple-choice [11, 12], and conversation-based reading comprehension [13, 14]. Examples of the resources available for other languages include the Chinese corpus for the span-extraction MRC [5], traditional Chinese corpus of MRC [15], the user-query-log-based corpus DuReader [16], and the Korean MRC corpus [17]. In addition to development of the reading comprehension corpora, various significant neural network-based approaches have been proposed and made a significant advancement in this research field, such as Match-LSTM [18], BiDAF [19], R-Net [20], DrQA [1], FusionNet [21], FastQA [22], QANet [23], and S3-NET [24]. Powerful transfer learning models such as BERT [25] and its variants (ALBERT [26]) have recently become extremely popular and achieved state-of-the-art results in MRC tasks.

Although researchers have studied several works on the Vietnamese language, such as parsing [27, 28, 29, 30], part-of-speech [31, 32], named entity recognition [33, 34, 35], sentiment analysis [36, 37, 38], and question answering [39, 40, 41], there is only two corpora for evaluating MRC models, ViMMRC [42] for evaluating Vietnamese multiple-choice questions and UIT-ViQuAD [43] for evaluating Vietnamese span-extraction MRC models. However, both two corpora are open-domain. In this paper, we aim to build a new large Vietnamese corpus based on online news articles in the health domain for evaluating MRC models. There are several main reasons for this. Firstly, machine comprehension for health domain has few studies so far, although it could be implemented into various potential and practical applications such as chatbot and virtual assistant in health-care service. Secondly, this study aims to build an application for general readers who search information and health-domain knowledge from online health articles. Finally, a new corpus is our important contribution to assess different MRC and QA models in a low-resource Vietnamese language.

The current approaches based on deep neural networks and transfer learning have surpassed the performance of humans with English corpora like SQuAD, but it is not clear these state-of-the-art models will obtain similar performance with corpora in different languages. Hence, to further enhance the development of the MRC, we develop a new span-extraction corpus for Vietnamese MRC. In this paper, we have three main contributions described as follows.

- Firstly, we develop a benchmark corpus (ViNewsQA) for evaluating Vietnamese machine reading comprehension and question answering systems. ViNewsQA comprises over 22,000 human-created question-answer pairs based on over 4,400 online news articles in the health domain. The corpus is publicly available for Vietnamese language processing research and also for the cross-lingual studies together with other similar corpora such as NewsQA (for English), CMRC (for Chinese), FQuAD (for French) and KorQuAD (for Korean).
- Besides, we analyze the corpus in terms of different linguistic aspects, including vocabulary-based, three types of length (question, answer, and article), three content-based types (question, answer and reasoning) and the correlation between type-based and the answer length, thereby providing comprehensive insights into the corpus that may facilitate future methods.
- Finally, we conduct the first experiments on different types of MRC methods as the first baseline models on the ViNewsQA corpus. The best-performance baseline is ALBERT with 65.26% (in EM) and 84.89% (in F1-score). The significant difference between humans and the best-performance model (10.90% of F1-score) indicates that improvements in ViNewsQA could be explored in the future study. In addition, we compare their performances with humans in terms of various linguistic aspects to obtain in-depth insights into Vietnamese span-extraction machine reading comprehension in the health domain using different methods.

The remainder of this paper is structured as follows. In Section 2, we review the existing machine reading comprehension corpora and models. In Section 3, we explain the creation process of our corpus. The analysis of our corpus is described in Section 4. Then, we present our experimental evaluation (in Section 5) and analysis of the experimental results and discussion (in Section 6). Finally, we draw our conclusions and suggest directions for future research in Section 7.

2 RELATED WORK

In this section, we review several studies related to our work, including related MRC corpora (in Section 2.1) and models (in Section 2.2).

2.1 Related MRC Corpora

Depending on the answer, the answers are divided into different types. Generally, reading comprehension tasks can be divided into four classes: cloze style [44, 45], multiple choice [11, 12], span extraction [4, 5] and free form [46]. In this study, we aim to construct a span-extraction MRC corpus on the health-domain online news for the Vietnamese language. We review thoroughly the corpora related to the span-extraction MRC and the health domain. Our ViNewsQA corpus is inspired by various recent span-extraction reading comprehension corpora, such as SQuAD [4], NewsQA [7], CMRC [5], and KorQuAD [17]. In particular, CliCR [47], MedQA [48] and PubMedQA [49] are three first MRC corpora in the health domain created in 2018. Table 1 summarize a brief of these corpora.

Table 1. A survey of several corpora related to our corpus ViNewsQA.

Corpus	Language	Domain	Type	Size	Annotation Method
SQuAD [4]	English	Open	Span-extraction	100K	Crowdsourcing
NewsQA [7]	English	Open	Span-extraction	100K	Crowdsourcing
CMRC [5]	Chinese	Open	Span-extraction	20K	Crowdsourcing
KorQuAD [17]	Korean	Open	Span-extraction	70K	Crowdsourcing
FQuAD [50]	French	Open	Span-extraction	60K	Crowdsourcing
SberQuAD [51]	Russian	Open	Span-extraction	90K	Crowdsourcing
UIT-ViQuAD [43]	Vietnamese	Open	Span-extraction	23K	Crowdsourcing
ViMMRC [42]	Vietnamese	Open	Multiple-choice	2.7K	Crowdsourcing
CliCR [47]	English	Medical	Gap-filling	105K	Crowdsourcing
MedQA [48]	English	Medical	Multiple-choice	270K	Published materials
PubMedQA [49]	English	Medical	Yes/No	273K	Crowdsourcing
<i>ViNewsQA (this work)</i>	<i>Vietnamese</i>	<i>Medical + News</i>	<i>Span-extraction</i>	<i>22K</i>	<i>Crowdsourcing</i>

Table 1 presents several MRC corpora and their characteristics. For the extraction-span MRC corpora, we review and analyze several well-known corpora, including SQuAD, NewsQA, CMRC, KorQuAD, FQuAD and SberQuAD. **SQuAD** is one of the best-known English corpora for the extractive MRC and it has facilitated the development of many machine learning models. In 2016, Rajpurkar et al. [4] proposed SQuAD v1.1 comprising 536 Wikipedia articles with 107,785 human-generated question and answer pairs. SQuAD v2.0 [6] was based on SQuAD v1.1 but it includes over 50,000 unanswerable questions created adversarially using the crowd-worker method according to the original questions. **NewsQA** is another English corpus proposed by Trischler et al. [7], which comprises 119,633 question-answer pairs generated by crowd-workers based on 12,744 articles from the CNN news. This corpus is similar to SQuAD because the answer to each question is a text segment of arbitrary length in the corresponding news article. **CMRC** [5] is a span-extraction corpus for Chinese MRC, which was introduced in the Second Evaluation Workshop on Chinese Machine Reading Comprehension in 2018. This corpus contains approximately 20,000 human-annotated questions on Wikipedia articles. This competition attracted many participants to conduct numerous experiments on this corpus. **KorQuAD** [17] is a Korean corpus for span-based MRC, comprising over 70,000 human-generated question-answer pairs based on Korean Wikipedia articles. The data collected and the properties of the data are similar to those in the English standard corpus SQuAD. **FQuAD** [50] is a French native reading comprehension corpus of questions and answers on a set of Wikipedia articles that consists of 25K questions for the 1.0 version and 60 questions for the 1.1 version. **SberQuAD** [51] contains 50K paragraph-question-answer triples and was created in a similar way to SQuAD. SberQuAD selected Wikipedia pages, split into paragraphs, and paragraphs presented to crowd workers. For each paragraph, a Russian native speaking crowd worker posed questions that can be answered using solely the content of the paragraph and their answers must have been a paragraph span, i.e., a contiguous sequence of paragraph words. All of these corpora are built based on the crowdsourcing method, which has motivated to build our corpus.

For the Vietnamese language, there are only two corpora for evaluating MRC models, including ViMMRC [42] and UIT-ViQuAD [43]. **ViMMRC** is the first Vietnamese corpus which consists of 2,783 pairs of multiple-choice question-answer-passage triples which are commonly used for teaching reading comprehension for elementary school students.

In addition, **UIT-ViQuAD** is a span-extraction open-domain corpus for the low-resource language as Vietnamese to evaluate MRC models. This corpus consists of over 23K human-generated question-answer pairs based on 5,109 passages of 174 Vietnamese articles from Wikipedia. Both of these corpora are open-domain, we want to target a domain-specific and be useful for future practical applications.

We choose the health domain for our corpus. Hence, we review several related corpora in this domain. **ClICR** [47] is a medical-domain corpus comprising around 100,000 gap-filling queries based on clinical case reports, while **MedQA** [48] collected answer real-world multiple-choice questions with large-scale reading comprehension. These corpora required world and background domain knowledge in the study of the MRC models. **PubMedQA** [49] is a novel biomedical QA corpus collected from PubMed abstracts with 273 yes/no/maybe QA instances. These corpora are mainly aimed at simple forms of English reading comprehension like filling-gap, multiple-choice and yes/no questions.

Until now, there are not any Vietnamese corpus for the span-based MRC research in the health-domain online news. The benchmark corpora mentioned above used for evaluating the MRC models and developing different QA applications, thereby encouraging researchers to explore the machine-learning models on these corpora. Our corpus is also intended for these purposes. These reasons lead to create a Vietnamese corpus in the health domain for MRC tasks.

2.2 Related MRC Methods

To the best of our knowledge, a range of studies have investigated MRC methodologies, and three popular approaches for MRC are rule-based, neural network-based and transfer learning-based. Rule-based approach is the first baseline of many well-known corpora [4, 11, 14, 42]. However, neural network-based and transfer learning-based systems have recently become more prevalent in MRC systems due to the powerful development of large-scale and high-quality corpora. In particular, we review them in detail as follows.

Rule-based Approaches. Sliding window (SW) is the first rule-based approach developed by Richardson et al. (2013) [11]. This approach matches a set of words built from a question and one of its answer candidates with a given reading text, before calculating the matching score using TF-IDF for each answer candidate. Experiments have been conducted with this simple model on many different corpora as first baseline models, such as MCTest [11], SQuAD [4], DREAM [14], and ViMMRC [42].

Machine Learning-Based Approaches. In addition to the rule-based models, machine-learning-based models have interesting features due to the development of large and high-quality corpora and robust machine configurations. In particular, Rajpurkar et al. [4] introduced a logistic regression model with a range of different linguistic features. However, neural network-based models on this problem have attracted more attention and obtained outstanding results in recent years. The corpora mentioned in Sub-section 2.1 have been studied in the development and evaluation of various neural network-based models in the field of natural language processing, such as Match-LSTM [18], BiDAF [19], CNN-LR [52], R-Net [20], DrQA [1], FusionNet [21], FastQA [22], QANet [23], and S3-NET [24]. In recent years, transfer-learning models have shown their strengths on many NLP tasks. In particular, Devlin et al. [25], Lan et al. [26], and Conneau et al. [53] introduced BERT and its variants (ALBERT and XLM-R), respectively, as powerful models trained on multiple languages and they obtained state-of-the-art performance with machine reading comprehension corpora.

In this paper, we choose several typical methods from three popular types of MRC models comprising rule-based (Sliding Window), neural network-based (DrQA and QANet) and

transfer learning-based (BERT and ALBERT) for our machine reading comprehension corpus. In addition, we attempt to analyze the experimental results in terms of different linguistic aspects to gain first insights into Vietnamese machine reading comprehension in the health domain.

3 CORPUS

In this section, we introduce the task of machine reading comprehension and give several examples in Vietnamese (in Section 3.1). Then, we present how to create a new corpus for evaluating Vietnamese machine reading comprehension in the health domain (in Section 3.2). These sections are described as follows.

3.1 Task Definition

Formally, the reading comprehension task is described as a triple (D, Q, A) , where D represents a document, Q represents a question, and A means an answer. Documents in our corpus are online news articles. Specifically, for the span-based reading comprehension task, question-answer pairs are created by humans. The answer A is a continuous span that is directly extracted from the document D . Figure 1 presents several examples for Vietnamese span-extraction reading comprehension in the health-domain online news.

<p>Document: Nghiên cứu cho thấy resveratrol trong rượu vang đỏ có khả năng làm giảm huyết áp, khi thí nghiệm trên chuột. Resveratrol là một hợp chất trong vỏ nho có khả năng chống oxy hóa, chống nấm mốc và ký sinh trùng. Trên Circulation, các nhà khoa học từ King’s College London (Anh) công bố kết quả thí nghiệm tìm ra sự liên quan giữa chuột và resveratrol. Cụ thể, resveratrol tác động đến huyết áp của những con chuột này, làm giảm huyết áp của chúng.</p> <p><i>(The study showed that resveratrol in red wine could reduce blood pressure when tested in mice. Resveratrol is a compound found in grape skin that has antioxidant, anti-mold, and anti-parasitic properties. Scientists from King’s College London (UK) published experimental results in Circulation regarding a link between mice and resveratrol. Specifically, resveratrol affected the blood pressure of these mice, lowering their blood pressure.)</i></p>
<p>Question 1: Chất bổ trong vỏ nho có tác dụng gì? <i>(What is the substance in grape skin for?)</i></p> <p>Answer: có khả năng chống oxy hóa, chống nấm mốc và ký sinh trùng <i>(has antioxidant, anti-mold, and anti-parasitic properties).</i></p>
<p>Question 2: Các nhà khoa học từ trường King’s tìm ra phát hiện gì về loài chuột và resveratrol? <i>(What did scientists from King’s University discover about mice and resveratrol?)</i></p> <p>Answer: resveratrol tác động đến huyết áp của những con chuột này, làm giảm huyết áp của chúng <i>(resveratrol affected the blood pressure of these mice, lowering their blood pressure).</i></p>

Fig. 1. Several examples of our proposed corpus (ViNewsQA). English translations are also provided for comparison.

3.2 Corpus Creation

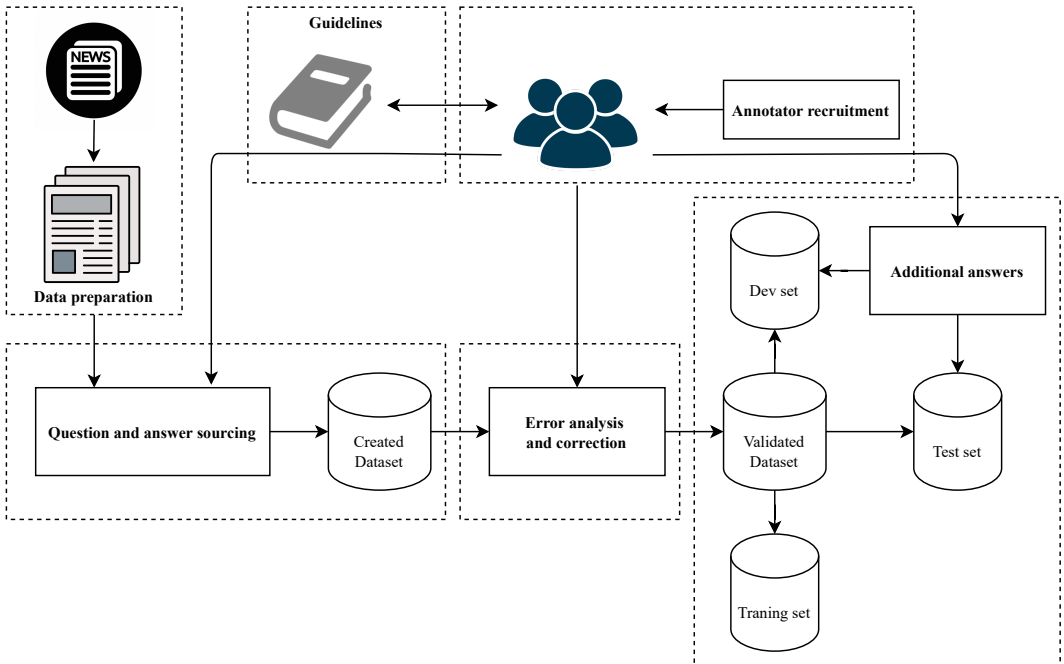


Fig. 2. The overview process of creating the Vietnamese MRC corpus in the health domain.

In this section, we present a new process to create the Vietnamese MRC corpus in the health domain, as shown in Figure 2. In particular, we construct our corpus through six different phases comprising (see in Section 3.2.1) annotator recruitment, (see in Section 3.2.2) building guidelines, (see in Section 3.2.3) data preparation, (see in Section 3.2.4) question and answer sourcing, (see in Section 3.2.5) validation based on error analysis and correction, and (see in Section 3.2.6) collecting additional answers. We describe these phases in detail as follows.

3.2.1 Annotator recruitment. We hire annotators to build our corpus according to a rigorous process in the following three different stages described as follows.

- **Stage 1:** People, who have an interest in reading health-domain online news, apply to become annotators to create the question-answer pairs for the MRC task.
- **Stage 2:** Annotators selected are good at general knowledge and passed our reading comprehension test.
- **Stage 3:** Official annotators are carefully trained guidelines (see in Section 3.2.2) with 200 questions. They MUST follow annotation rules presented in Section 3.2.2.

3.2.2 Guidelines. The annotators read and understand each article, and they then formulate questions and select their answers directly in the article. During the creation process of question-answer pairs, the annotators conform to the following rules.

- **Rule 1:** Annotators are required to pose at least three question-answer pairs per the article.

- **Rule 2:** Annotators are encouraged to ask questions in their own words and vocabulary.
- **Rule 3:** The answer **MUST** be a span in the article that satisfy the requirements of the task definition. The spans with the shortest length from potential answers are encouraged to be selected for the answers to the questions.
- **Rule 4:** To diverse different types of questions, annotators are encouraged to create questions with different types (what/who/when/where/why/how, etc.). In addition, complex reasoning (single-sentence and multiple-sentence reasoning) is also encouraged in the question generation.
- **Rule 5:** Annotators are warned about mistakes that could be avoided when creating questions–answer pairs. These mistakes are shown from our error analysis presented in Section 3.2.5.

3.2.3 Data preparation. 4,416 news articles related to **the health topic are collected from the online newspaper VnExpress⁴. We choose this source because it is one of the most popular Vietnamese online newspaper⁵⁶ and the language used in articles is easy to understand for general readers, which aim for practical applications.** All images, figures, and tables are eliminated from these articles, and articles shorter than 300 characters or those containing many special characters and symbols are removed. We divide the articles randomly into a training set (Train), a development set (Dev), and a test set (Test) with an approximate rate of 8:1:1 **for conducting experiments on machine reading comprehension models in Section 5.**

3.2.4 Question and answer sourcing. Following the guidelines (see in Section 3.2.2), annotators create question-answer pairs per article. Annotators use the MRC annotation tool that we build to create question-answer pairs. In each working section, the tool allows to display the article content and enables the annotators to enter questions and choose their answers directly on the article and also allows the annotators to save the article content, the questions, and answers on a *.json file.

3.2.5 Error analysis and correction. Errors that may arise when manually creating questions and choosing answers from articles are inevitable. To enhance the quality of the corpus, we perform the validation process to minimize these errors. To analyze the error types that can occur during the data generation of annotators, we select randomly 1,183 question-answer pairs (over 5% of the corpus) to investigate the errors and find 335 question-answer pairs with mistakes. Based on questions or answers, we divide these errors into five different types such as unclear questions (Error type 1), misspelled questions (Error type 2), incorrect answers (Error type 3), lack-or-excess-of-information answers (Error type 4), and incorrect-boundary answers (Error type 5). These errors are described as follows.

- **Error type 1:** Questions are misspelled. In the process of creating questions and their answers, annotators could misspell during the typing process.
- **Error type 2:** Answers are incorrect for their questions. In particular, questions are correct, but their selected answers are wrong.
- **Error type 3:** Answers are lack or excess of information for questions. In particular, annotators can choose the redundant or unnecessary text to answer their questions.
- **Error type 4:** Questions are not precise and clear in their contents. People cannot understand these questions, so they do not find answers to these questions.

⁴<https://vnexpress.net/suc-khoe>

⁵<https://en.wikipedia.org/wiki/VnExpress>

⁶<https://www.alexa.com/topsites/countries/VN>

- **Error type 5:** Answers are incorrect-boundary spans. Remarkably, the annotators can choose either lack or excess some characters or spaces in the answer.

Table 2 presents the common types of errors that annotators made during the corpus creation process. We find that the error type 3 occurs most frequently and accounts for 54.33% while the error type 5 accounts for the lowest percentage of 1.49%. From these analyses, we require the annotators to check and correct carefully with these errors. Besides, these types of errors are useful for future development of MRC corpora.

Table 2. Statistics of error types of annotators when creating question-answer sourcing. Examples and their English translations are also given for comparison.

Error types	Examples	Percentage (%)
1	<p>Annotator's question: Salmonella gây ra những gnyu hiểm gì cho phụ nữ mang thai?</p> <p>Correct question: Salmonella gây ra những <u>nguy hiểm</u> gì cho phụ nữ mang thai?</p> <p>English translation:</p> <p>Annotator's question: What are the <u>adngers</u> of Salmonella in pregnant women?</p> <p>Correct question: What are the <u>dangers</u> of Salmonella in pregnant women?</p>	4.78
2	<p>Document: Bệnh thủy đậu xảy ra ở mọi lứa tuổi, chủ yếu ở trẻ em. Những người có hệ miễn dịch kém như người trên 50 tuổi, suy dinh dưỡng hoặc đang sử dụng thuốc điều trị ung thư, thuốc ức chế miễn dịch, phụ nữ có thai... có nguy cơ cao mắc bệnh. Người lớn thường bị trong các trường hợp suy giảm miễn dịch, thông thường bệnh sẽ nặng hơn ở trẻ em.</p> <p>Annotator's question: Nếu mắc bệnh thủy đậu thì ai sẽ mắc bệnh nặng hơn?</p> <p>Annotator's answer: Trẻ em.</p> <p>Correct answer: Người lớn.</p> <p>English translation:</p> <p>Document: Chickenpox happens to people of all ages, mostly in children. People with weakened immune systems such as those aged over 50, malnourished or taking cancer treatment drugs, immunosuppressants, pregnant women, etc. have a greater incident of the disease. Adults are often infected in cases of immunodeficiency while the children suffer more seriously.</p> <p>Annotator's question: Assuming someone gets the chickenprox, who could be worse?</p> <p>Annotator's answer: Children.</p> <p>Correct answer: Adult.</p>	4.78

Document: Hydro sulfua là hợp chất khí ở điều kiện nhiệt độ thường, không màu, mùi trứng thối.

Annotator's question: Tính chất vật lý của khí H₂S là gì?

Annotator's answer: hợp chất khí ở điều kiện nhiệt độ thường.

Correct answer: hợp chất khí ở điều kiện nhiệt độ thường, không màu, mùi trứng thối.

3

English translation:

Document: Hydrogen sulfide is a gas compound at normal temperature conditions, colorless, rotten egg smell.

Annotator's question: What are the physical properties of H₂S?

Annotator's answer: gas compound at room temperature.

Correct answer: gas compound at normal temperature, colorless, rotten egg smell.

54.33

Document: Mọi người cần tiêm phòng thủy đậu cho tất cả trẻ em và người lớn chưa nhiễm bệnh. Hiện vắc xin này được sử dụng khá phổ biến và đem lại hiệu quả cao trong việc phòng bệnh, ít gây tác dụng phụ.

Annotator's question: Vắc xin hiện nay có được sự cải tiến nào?

Correct question: Vắc xin thủy đậu hiện nay có được sự cải tiến nào?

4

English translation:

Document: People are in need of vaccinating all uninfected children and adults against the chickenpox. Currently, this vaccine is used quite commonly and brings about the high efficiency in the prevention of disease while causing few side effects.

Annotator's question: Is there any improvement in current vaccines?

Correct question: Which improvement is being taken place on the current chickenpox vaccines?

34.62

Document: Trong 101.862 trẻ được tiêm vắc xin ComBE Five trên 19 tỉnh, ghi nhận 1,73% trẻ có phản ứng thông thường như sốt nhẹ, sưng đau tại chỗ tiêm, khó chịu, quấy khóc.

Annotator's answer: ốt nhẹ, sưng đau tại chỗ tiêm, khó chịu, quấy khóc.

Correct answer: sốt nhẹ, sưng đau tại chỗ tiêm, khó chịu, quấy khóc.

5

English translation:

Document: Among 101,862 children were injected the ComBE Five vaccine in 19 provinces, 1.73% of them have some slightly reaction such as low fever, soreness at the injection site, discomfort, and crying.

Annotator's answer: ow fever, soreness at the injection site, discomfort, crying.

Correct answer: low fever, soreness at the injection site, discomfort, crying.

1.49

3.2.6 Collecting additional answers. To estimate the performance of humans (see in Section 5.4) and to enhance our empirical evaluations of the development and test sets, we add two more answers for each question, and the first answers in the development and test sets are annotated by annotators. During this phase, the annotators do not know the first answer, and they are encouraged to give diverse answers.

4 CORPUS ANALYSIS

Firstly, we introduce the overview of our corpus in Section 4.1. To understand the characteristics of our corpus ViNewsQA, we perform a variety of analyzes based on linguistic aspects such as vocabulary-based in Section 4.2, length-based (question length, answer length and article length) in Section 4.3 and type-based (question type, answer type and reasoning type) in Section 4.4. In addition, we analyze the correlation between type-based and the answer length in Section 4.5. These analyses provide in-depth insights into our corpus and comparisons with another Vietnamese corpus as UIT-ViQuAD. For question type, answer type and reasoning type, we also hire annotators to annotate questions of the development set which are selected randomly from our corpus. Corpora such as SQuAD [4] and UIT-ViQuAD [43] also perform corpus analysis on the development set.

4.1 Overall statistics

Before conducting detailed analyses, we provide an overview of our corpus. Table 3 presents statistics for the training, development, and test sets in our corpus. ViNewsQA consists of 22,057 question-answer pairs based on 4,416 online news articles in health domain. Table 3 shows the number of articles and the average lengths⁶ for questions and answers, and the vocabulary size. The number of questions of our corpus is approximately equal to the number of questions of UIT-ViQuAD.

⁶We use the pyvi library <https://pypi.org/project/pyvi/> for word segmentation to calculate the average lengths of articles, questions and answers, and the vocabulary size.

Table 3. Overview statistics of ViNewsQA. * indicates that UIT-ViQuAD used passages as reading texts.

	Train	Dev	Test	All	UIT-ViQuAD
Number of article	3,517	500	399	4,416	174
Number of passage*	-	-	-	-	5,109
Number of questions	17,568	2,497	1,992	22,057	23,074
Average reading-text length	342,9	323.9	360.4	342.4	153.4
Average question length	10.6	10.8	10.3	10.6	12.2
Average answer length	10.7	10.3	10.9	10.7	8.2
Vocabulary size	29,111	10,765	10,020	32,749	41,773

4.2 Vocabulary-based analysis

To understand the health domain, we utilize the word cloud tool¹ to create graphical representations of word frequency for articles (see in Figure 3), questions (see in Figure 4) and answers (see in Figure 5) in our corpus. The larger the word in the visual, the more common the word is in the articles, questions, and answers. Vietnamese stop words² and numbers are excluded from these statistics. Table 4, Table 5 and Table 6 show the top tenth popular words that appears in articles, questions and answers in our corpus, respectively. These words are in the health domain, which is also characteristic of our corpus. Because the corpus is collected from the health online news articles. Five different words such as *bệnh nhân* (patient), *bác sĩ* (doctor), *bệnh* (disease), *bệnh viện* (hospital), *ung thư* (cancer) are appeared all in articles, questions and answers. Figure 7 and Figure 8 present word distribution for ViNewsQA and UIT-ViQuAD, respectively. Top 10 words from ViNewsQA (see Table 7) and UIT-ViQuAD (see Table 8) are very different because these high-frequency words on the UIT-ViQuAD corpus belong to multiple domains such as history, geography, economics, and politics.

¹Word cloud tool: <https://www.wordclouds.com>

²Vietnamese stop words: <https://github.com/stopwords/vietnamese-stopwords>



Fig. 3. Word distribution of articles.



Fig. 4. Word distribution of questions.



Fig. 5. Word distribution of answers.

No.	Vietnamese	English	Freq.
1	bác sĩ	doctor	9,314
2	bệnh nhân	patient	7,182
3	bệnh viện	hospital	6,231
4	bệnh	disease	5,762
5	máu	blood	3,965
6	điều trị	treat	3,901
7	thuốc	medicine	3,587
8	ung thư	cancer	3,205
9	y tế	medical	3,066
10	phẫu thuật	surgery	2,962

Table 4. Common words appeared in articles.

No.	Vietnamese	English	Freq.
1	bệnh nhân	patient	2,281
2	bác sĩ	doctor	1,753
3	bệnh	disease	1,275
4	bệnh viện	hospital	1,125
5	nguyên nhân	reason	694
6	điều trị	treat	618
7	phẫu thuật	surgery	601
8	ung thư	cancer	553
9	mắc	get sick	550
10	trẻ	young	533

Table 5. Common words appeared in questions.

No.	Vietnamese	English	Freq.
1	máu	blood	2,280
2	bệnh	disease	2,092
3	bệnh nhân	patient	1,807
4	đau	pain	1,689
5	cơ thể	body	1,542
6	ung thư	cancer	1,450
7	thuốc	medicine	1,390
8	tim	heart	1,207
9	viêm	inflamm	1,143
10	bệnh viện	hospital	1,130

Table 6. Common words appeared in answers.



Fig. 6. Word distribution of ViNewsQA.

No.	Vietnamese	English	Freq.
1	bác sĩ	doctor	11,628
2	bệnh nhân	patient	10,878
3	bệnh	disease	8,738
4	bệnh viện	hospital	8,288
5	máu	blood	6,173
6	điều trị	treat	5,393
7	cơ thể	body	5,173
8	thuốc	medicine	5,136
9	ung thư	cancer	4,878
10	phẫu thuật	surgery	4,167

Table 7. Common words appeared in ViNewsQA.



Fig. 7. Word distribution of UIT-ViQuAD.

No.	Vietnamese	English	Freq.
1	quốc gia	nation	1,444
2	thành phố	city	1,264
3	chính phủ	government	1,256
4	dân	resident	1,189
5	thế giới	world	1,142
6	nam	south	1,077
7	châu	continent	1,068
8	phát triển	develop	1,045
9	đảng	party	1,005
10	kinh tế	economy	999

Table 8. Common words appeared in UIT-ViQuAD.

4.3 Analysis based on different lengths

4.3.1 Analysis based on question length. Statistics for the various question lengths are shown in Table 9. Questions with 8–9 words comprise the highest proportion with 25.67%. Most of the questions in the corpus have lengths from 6 to 13 words, which account for approximately 80% of the corpus. Very short questions (4–5 words) and long questions (>=18 words) account for a low percentages of 2.68% and 3.91%, respectively.

Table 9. Statistics for the question lengths in ViNewsQA.

Question length	ViNewsQA				UIT-ViQuAD
	Train	Dev	Test	All	
4-5	2.73	2.28	2.71	2.68	0.99
6-7	14.29	14.14	16.52	14.48	7.10
8-9	25.81	23.35	27.31	25.67	17.34
10-11	23.84	23.91	24.40	23.90	22.07
12-13	15.82	16.94	13.76	15.76	20.28
14-15	8.89	10.25	6.88	8.86	13.96
16-17	4.66	5.37	4.62	4.74	8.91
18-19	2.09	2.08	2.11	2.09	4.91
>19	1.86	1.68	1.71	1.82	4.44

4.3.2 Analysis based on answer length. Table 10 shows the distribution of the answer length analysis in our corpus. The largest percentage (14.73%) comprise answers with lengths of 3–4 words. Most of the answers (nearly 60%) have lengths of 1–10 words. Longer answers (over 10 words) comprise a low proportion of our corpus.

Table 10. Statistics of the answer lengths on ViNewsQA.

Answer length	ViNewsQA				UIT-ViQuAD
	Train	Dev	Test	All	
1-2	12.53	14.82	13.25	12.85	29.46
3-4	14.80	15.18	13.60	14.73	17.64
5-6	11.48	12.37	12.75	11.69	12.10
7-8	10.50	10.37	9.04	10.35	7.95
9-10	9.91	9.33	8.23	9.69	6.39
11-12	8.64	7.81	9.04	8.58	5.03
13-14	7.02	6.69	6.83	6.96	4.17
15-16	5.85	5.01	6.73	5.83	3.25
17-18	4.99	3.80	5.17	4.87	2.62
19-20	3.67	3.16	3.56	3.60	2.26
>20	10.62	11.45	11.80	10.82	9.13

4.3.3 Analysis based on article length. Besides, we also aim to analyze article lengths in the corpus. Table 11 presents statistics for various article lengths in our corpus. The lengths of most articles ranging from 101 to 500 words, which account for over 84%. The length of the reading texts on ViNewsQA is significantly longer than that on the UIT-ViQuAD. Figure 8 shows different reading-text-length distributions of the two Vietnamese corpora. Based on the length characteristics, we determine whether the length affects the performance of the machine models and humans according to question, answer, or article lengths?

Table 11. Statistics for ViNewsQA according to the article length.

Article length	Train	Dev	Test	All
<101	0.34	0.40	0.25	0.34
101-200	15.33	18.40	13.53	15.51
201-300	29.03	31.40	27.07	29.12
301-400	23.34	23.60	23.56	23.39
401-500	16.52	13.40	16.29	16.15
501-600	10.41	7.80	11.03	10.17
>600	5.03	5.00	8.27	5.32

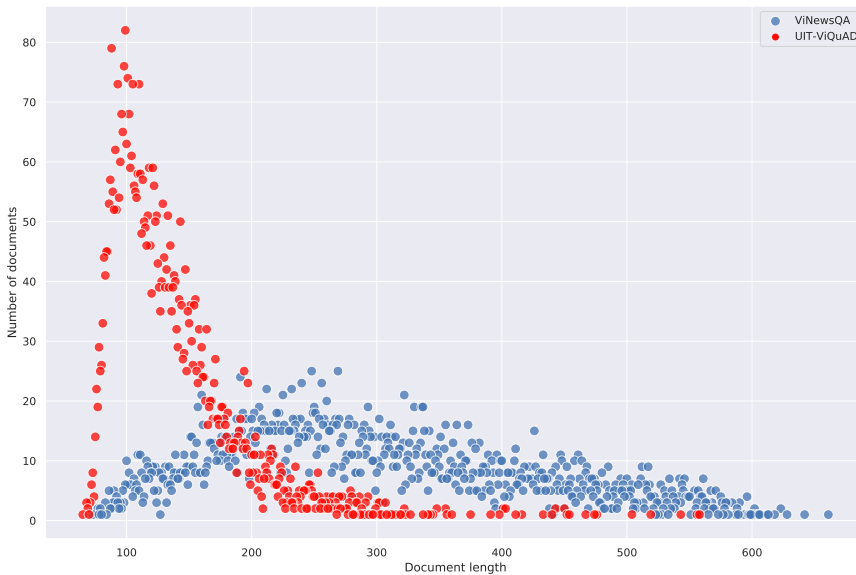


Fig. 8. Document length distributions of two corpora (ViNewsQA and UIT-ViQuAD).

4.4 Analysis based on different types

Before conducting the type-based analysis, we hire three annotators and train them on 100 questions to have more than 80% of the Cohen’s kappa inter-annotator agreement before annotating data simultaneously.

4.4.1 Analysis based on question type. In this work, we divide Vietnamese questions into seven different question types such as Who, What, When, Where, Why, How, and Others, which is constructed in a similar way to UIT-ViQuAD [43] and CMRC [5]. These question types are described as follows.

- **Who:** A group of questions have their answers related to people.

- **When:** A group of questions requires their answers presenting time expression.
- **Where:** A group of questions requires their answers which are locations or places.
- **Why:** A group of questions require their answers expressing reasons.
- **How:** A group of questions require their answers related to an away or method to do.
- **What:** A group of questions have answers which are definitions, things or events.
- **Others:** A group of questions do not belong to the above types. Most of the questions are related to numbers such How many and How much.

Table 12 presents the distribution of the question types on our corpus. The table show that the question type What accounted for the largest proportion with 54.35%. Compared to the SQuAD corpus, the rate of the What question in our corpus is similar to that in SQuAD (49.97%) [54]. Our corpus requires abilities beyond factoid questions that demand intricate knowledge and skills to answer like Why and How questions. In particular, How and Why ranked the second and the third with 13.46% and 12.17%, respectively.

Table 12. Statistics and examples of question type.

Type	Example	Percentage (%)	
		ViNewsQA	UIT-ViQuAD
Who	Vietnamese: Bệnh viện Thẩm mỹ Emcas đã hợp tác với ai để thực hiện ca phẫu thuật? English: Who has Emcas Cosmetic Hospital cooperated with to perform the surgery?	3.16	9.41
When	Vietnamese: Mỹ dùng năng lượng vi sóng để điều trị hôi nách vào năm nào? English: When did the US use microwave energy to treat armpits?	3.52	8.96
Where	Vietnamese: Bệnh viện đã cử 14 người đi đào tạo ở đâu? English: Where did the hospital take 14 people to train?	3.40	5.64
Why	Vietnamese: Tại sao phụ nữ ở Mỹ ngày càng sinh con muộn? English: Why are women in America increasingly late for giving birth?	12.17	7.54
How	Vietnamese: Bệnh viêm não cấp thường lây nhiễm bằng cách nào? English: How is acute encephalitis usually spread?	13.46	9.09
What	Vietnamese: Trong các thức uống năng lượng có những thành phần nào giống nhau? English: What are the similar ingredients in energy drinks?	54.35	49.97
Others	Vietnamese: Hà nội đã tiếp nhận bao nhiêu đơn vị máu từ người hiến tình nguyện vào năm 2018? English: How many units of blood did Hanoi receive from voluntary donors in 2018?	9.94	9.41

4.4.2 Analysis based on answer type. We divide answers into 11 types including numbers (time, other numeric), entity (person, location, other entity), phrase (noun phrase, adjective phrase, verb phrase, prepositional phrase, clause) and others. The priority order of annotation is number, entity, phrase and others. Table 13 present statistics of answer types in our corpus. While verb phrases account for the highest proportion of 34.84%, prepositional phrases account for the lowest proportion with 0.8%.

Table 13. Statistics and examples of answer type.

Type	Example	Percentage (%)	
		ViNewsQA	UIT-ViQuAD
Time	Vietnamese: tháng 5/2017, tuần thứ 36. English: May 2017, 36th week.	4.49	7.71
Other numeric	Vietnamese: 12%, hơn 350 calo. English: 12%, more than 350 calories.	9.29	9.41
Person	Vietnamese: Bác sĩ Nguyễn Khắc Vui, nhiếp ảnh gia Amy Taylor. English: Dr. Nguyen Khac Vui, photographer Amy Taylor.	0.96	5.39
Location	Vietnamese: bệnh viện Westchester, Quận 8. English: Westchester Hospital, District 8.	1.16	4.32
Other entity	Vietnamese: Cục An toàn thực phẩm, khoa Hồi sức cấp cứu. English: Food Safety Department, Emergency Care Department.	4.25	11.65
Noun phrase	Vietnamese: trẻ em, sự lây lan của vi khuẩn. English: children, spread of bacteria.	27.63	22.86
Adjective phrase	Vietnamese: rất đắt đỏ, béo phì. English: very expensive, fat.	4.41	2.52
Verb phrase	Vietnamese: kiểm tra huyết áp, uống một ly nước chanh. English: check blood pressure, drink a glass of lemon juice.	34.84	18.43
Preposition phrase	Vietnamese: dưới gan phải, ở vùng đáy tử cung. English: under the right liver, in the base of the uterus.	0.80	3.18
Clause	Vietnamese: mỗi tình nguyện viên sẽ cần uống khoảng 1.000 chai rượu mỗi ngày, Thuốc lá làm cạn lượng vitamin C trong cơ thể.		

	English: Each volunteer will need to drink about 1,000 bottles of wine per day, Tobacco depletes vitamin C in the body.	5.65	5.91
Others	Vietnamese: Tôi đã sống một cuộc đời rất hạnh phúc và viên mãn nên không còn gì phải hối hận. Không quan trọng là mình sống được bao lâu, quan trọng là ý nghĩa mỗi ngày được sống. English: I have lived a very happy and fulfilled life, so I have no regrets. It doesn't matter how long we live, what matters to live each day.	6.52	10.55

4.4.3 Analysis based on reasoning type. To classify the difficulty of a question, we divided question reasoning into one of five types, comprising word matching, paraphrasing, single-sentence reasoning, multi-sentence reasoning, and ambiguous/insufficient. These reasoning types are described as follows.

- **Word matching (WM):** The main words in the question exactly match the words in the reading text.
- **Paraphrasing (PP):** The questions are paraphrased from a single sentence in the reading text. In particular, we may use synonymy and world knowledge to create the question.
- **Single-sentence Reasoning (SSR):** The answers are inferred from a single sentence in the article. Such answers could be created by extracting incomplete information or conceptual overlap.
- **Multi-sentence Reasoning (MSR):** The answers are inferred from multiple sentences in the article by information fusion techniques.
- **Ambiguous/Insufficient (AoI):** The questions have many answers or answers are not found in the article.

The reasoning types in the development set for our corpus annotated by workers. Table 14 shows the distributions of the reasoning types in our corpus. The reasoning of the largest proportion is paraphrasing (PP) with 31.16%, whereas the lowest is Ambiguous/Insufficient (AoI) with 0.40%.

Table 14. Statistics and examples of reasoning type.

Reasoning Example	Percentage (%)	
	ViNewsQA	UIT-ViQuAD

Word matching	<p>Context: Thành phần trong nhân sâm chứa nhiều <i>ginsenoside</i> giúp cải thiện và tăng đáng kể hàm lượng testosterone ở nam giới, làm tăng cảm giác hưng phấn.</p> <p>Question: Thành phần trong nhân sâm chứa nhiều chất gì?</p> <p>Answer: <i>ginsenoside</i></p>	26.19	13.35
Paraphrasing	<p>English translation:</p> <p>Context: The ingredient in ginseng contains <i>ginsenoside</i>, which helps improve and significantly increase testosterone content in men, increasing feelings of excitement.</p> <p>Question: What are the ingredients in ginseng?</p> <p>Answer: <i>ginsenoside</i>.</p> <hr/> <p>Context: Bác sĩ Nguyễn Hữu Thịnh cho biết người bị tắc ruột thường có các triệu chứng <i>buồn nôn, trướng bụng, không thể đại tiện hoặc trung tiện, đau bụng quặn từng cơn</i>.</p> <p>Question: Người bị tắc ruột thường có các dấu hiệu gì?</p> <p>Answer: <i>buồn nôn, trướng bụng, không thể đại tiện hoặc trung tiện, đau bụng quặn từng cơn</i>.</p>	31.16	31.22

Single-sentence reasoning	<p>Context: Hội chứng người sói Congenital Hypertrichosis hay hội chứng người sói là căn bệnh di truyền hiếm gặp khiến <i>lông, tóc trên toàn bộ cơ thể phát triển quá mức</i>.</p> <p>Question: Người bị hội chứng người sói sẽ có những điểm gì khác biệt?</p> <p>Answer: <i>lông, tóc trên toàn bộ cơ thể phát triển quá mức</i>.</p>	16.78	38.07
Multiple-sentence reasoning	<p>English translation:</p> <p>Context: The werewolf syndrome Congenital Hypertrichosis or werewolf syndrome is a rare genetic disease that causes <i>excessive hair and hair growth throughout the body</i>.</p> <p>Question: What are the differences between people with werewolf syndrome?</p> <p>Answer: <i>excessive hair and hair growth throughout the body</i>.</p> <hr/> <p>Context: Theo Bảng thành phần dinh dưỡng Việt Nam, <i>trứng gà ít calo và cholesterol hơn trứng vịt</i>. Bác sĩ khuyên mỗi ngày chỉ nên ăn một quả trứng gà thay vì trứng vịt sẽ phù hợp với những người huyết áp cao, tim mạch.</p> <p>Question: Tại sao người huyết áp cao và tim mạch nên ăn trứng gà thay vì trứng vịt?</p> <p>Answer: <i>trứng gà ít calo và cholesterol hơn trứng vịt</i>.</p>	25.47	16.22

Context: Bé được chuyển lên phòng mổ cấp cứu, bác sĩ xác định bỏng nặng độ I,II,III. Bệnh nhi *phải ghép da, vá da dày toàn bộ, băng ép cố định diện ghép da, nẹp cố định cánh, cẳng tay trái bằng nẹp bột, vệ sinh làm sạch vết thương hằng ngày... và chăm sóc đặc biệt*

Question: Bệnh nhi cần được cấp cứu như thế nào?

Answer: *chăm sóc đặc biệt*

Correct: *phải ghép da, vá da dày toàn bộ, băng ép cố định diện ghép da, nẹp cố định cánh, cẳng tay trái bằng nẹp bột, vệ sinh làm sạch vết thương hằng ngày... và chăm sóc đặc biệt.*

Ambiguous **English translation:**

0.40

2.11

or insufficient

Context: The baby was transferred to the emergency operating room, the doctor determined severe burns I, II, III. Patients *must have skin grafts, thick skin patches, compression bandages for fixed skin grafts, fixed braces, left forearms with powder splint, cleaning the wound daily... and special care.*

Question: How do children need emergency care?

Answer: *special care*

Correct: *must have skin grafts, thick skin patches, compression bandages for fixed skin grafts, fixed braces, left forearms with powder splint, cleaning the wound daily... and special care.*

4.5 Correlation between answer length and different types

To understand relations between answer length and different linguistic types, we analyze various correlations consisting of question type-answer length, answer type-answer length, and reasoning type-answer length. Nguyen et al. [43] shows that answer length had an apparent effect, and the longer the answer is, the more difficult it is. Therefore, we choose answer length to analyze along with types. The specific results of our analysis are presented as follows.

Table 15 presents the correlation between question type and answer length. As can be seen from the table, answers for Who, When, and Where tend to be short and are mostly 1 to 5 words in length. In contrast, answers for Why, How, and What tend to be long and varied in length because they are more difficult than Who, When, and Where questions. Others account for a large proportion (>90%) for 1 to 5-word answers because they are mostly related to the number, such as How many and How much.

Table 15. The relation between question type and answer length.

Answer length	1-5	6-10	11-15	15-20	>20
Who	51.90	21.52	10.13	7.59	8.86
When	75.00	18.18	4.55	1.14	1.14
Where	62.35	17.65	12.94	2.35	4.71
Why	24.67	32.57	19.74	12.50	10.53
How	28.37	28.74	19.60	9.80	13.49
What	18.45	29.76	20.83	13.99	16.96
Others	90.20	6.12	2.86	0.41	0.41

Table 16 shows the correlation between answer type and answer length. Answers based on numbers (time and other numeric), entities (person, location and other entity) and preposition phrases tend to be shorter than word answers based on noun phrases, verb phrases, adjective phrases, clauses and others.

Table 16. The relation between answer type and answer length.

Answer length	1-5	6-10	11-15	16-20	>20
Time	85.71	12.50	0.00	0.00	1.79
Other numeric	86.21	9.48	3.02	0.86	0.43
Person	91.67	8.33	0.00	0.00	0.00
Location	72.41	17.24	10.34	0.00	0.00
Other entity	73.58	13.21	7.55	0.94	4.72
Noun phrase	33.19	30.87	18.26	7.83	9.86
Verb phrase	21.84	31.15	21.61	13.10	12.30
Adjective phrase	25.45	40.91	19.09	5.45	9.09
Preposition phrase	40.00	25.00	30.00	0.00	0.00
Clause	13.48	20.57	23.40	21.99	20.57
Others	7.89	20.86	20.86	12.27	38.04

Table 17 shows the correlation between reasoning type and answer length. Multi-sentence reasoning questions tend to be longer than other questions. In contrast, word matching, paraphrasing and single-sentence reasoning questions have almost the same ratio.

Table 17. The relation between reasoning type and answer length.

Answer length	1-5	6-10	11-15	16-20	>20
Word matching	41.74	26.61	15.29	7.34	9.02
Paraphrasing	36.63	27.51	17.87	8.35	9.64
Single-sentence reasoning	40.57	27.21	14.80	8.11	9.31
Multi-sentence reasoning	26.89	23.43	19.50	12.74	17.45
Ambiguous or insufficient	50.00	30.00	0.00	0.00	10.00

5 EMPIRICAL EVALUATION

In this section, we aim to evaluate three different types of methods consisting of rule-based, neural network-based, and transfer learning-based models on our corpus. Our experiments

revolve around the following questions: **Q1**: How do the MRC models perform on this new corpus we build? and **Q2**: Do these models outperform humans?

5.1 Re-implemented methods and baselines

To solve the research question **Q1**, MRC models are chosen: rule-based (Sliding Window), neural network-based (DrQA and QANet), and transfer learning-based (BERT and ALBERT). Because these are popular for evaluating machine reading comprehension in English [1, 23, 25, 53] and in other languages [5, 17, 50, 51, 43]. This section conducts various experiments to evaluate the MRC methods as the first baseline models on our corpus.

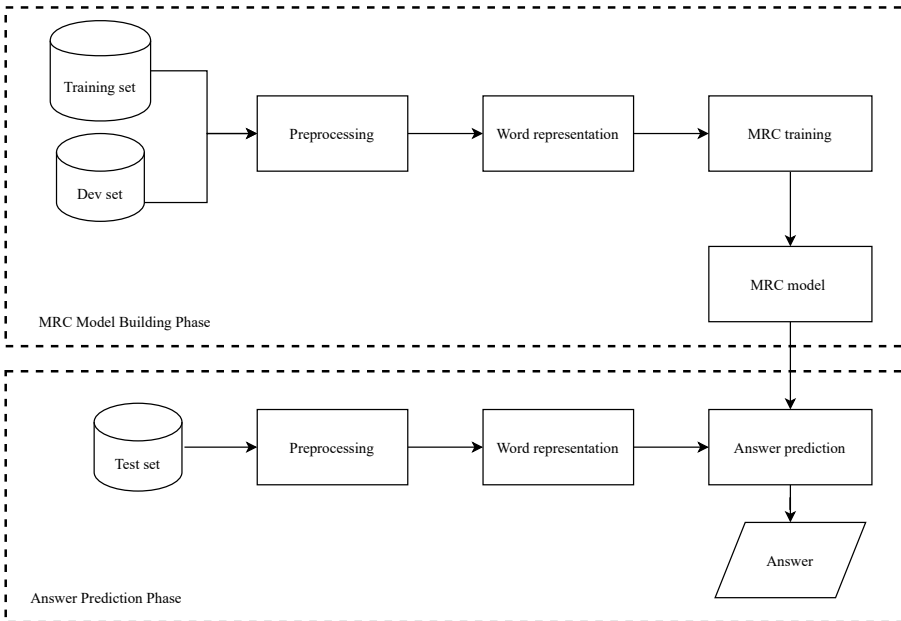


Fig. 9. Baseline framework for machine reading comprehension in Vietnamese.

Sliding Window is the first baseline model to conduct experiments on the well-known different corpora such as MCTest [11], SQuAD [4], DREAM [14], and ViMMRC [42]. In recent years, neural network-based and transfer learning-based models have achieved better performances on the span-extraction MRC corpora. How do these methods work well on our corpus? Hence, we aim to evaluate these methods as baseline models on our corpus. To adapt these methods into our corpus, we follow the Vietnamese MRC framework, as shown in Figure 9. The pre-processing stage before training involved word segmentation, removing extra white spaces, and updating the answer positions. In addition, we remove punctuations when performing evaluations. We use pre-trained embeddings for word representations, which have proved effective in the MRC task in Vietnamese [43, 42]. Four popular neural network-based and transfer learning-based MRC methods are chosen for the framework: DrQA Reader [1], QANet [23], BERT [25] and ALBERT [26], because these methods have achieved state-of-the-art performances in many MRC tasks on SQuAD [25], NewsQA [55], and CMRC [5]. These models are described as follows.

- **Sliding Window:** Sliding window (SW) is the first rule-based approach proposed by Richardson et al. (2013) [11]. This approach matches a set of words built from a question and one of its answer candidates with a given document before calculating the matching score using TF-IDF for each answer candidate. Experiments have been conducted with this simple model on many different corpora as first baseline models, such as MCTest [11], SQuAD [4], DREAM [14], and ViMMRC [42].
- **DrQA:** We implement the simple but effective neural-based model DrQA Reader, which is based on the open-domain QA system called DrQA proposed by Chen et al. [1] in 2017. This model has achieved good performance with multiple MRC corpora such as SQuAD [4] and CoQA [13]. DrQA is known as a simple-reasoning MRC model with multiple layers. In the input layer, the model presents binary features to the lexical-unit embedding of each context lexical unit if this lexical unit and its variants appear in the question. In addition, the model extends the lexical-unit embeddings in the input layer with the linguistic features such as POS and NER. The UIT-ViQuAD corpus also used DrQA as one of first baseline model [43].
- **QANet:** Yu et al. proposed QANet [23] and this model has obtained good performance with many MRC corpora [4, 56]. QANet has a feed-forward architecture with convolutions and attention mechanisms for MRC. The model comprises multiple convolutional layers followed by two components: the self-attention and fully connected layer, for both question and reading text encoding, as well as three stacked layers before predicting the final output. It was also the first baseline model of the UIT-ViQuAD corpus [43].
- **BERT:** Devlin et al. [25] recently proposed BERT, which uses a transformer network to pre-train a language model for extracting contextual word embeddings. This model is one of the best for contextualized representation learning [57, 58, 59, 25] and it has achieved the state-of-the-art results in multiple reading comprehension tasks. In this study, we used mBERT [25], as a large-scale multilingual language model, which was pre-trained for evaluating our Vietnamese MRC task. Nguyen et al. [43] chose BERT as one of the first baseline models.
- **ALBERT:** Lan et al. proposed ALBERT [26] to improve NLP tasks. This model was designed similarly to BERT, using the transformer encoder architecture with the GELU activation function. However, according to the authors of ALBERT, this model has more highlights than BERT comprising factorized embedding parameterization, cross-layer parameter sharing, and inter-sentence coherence loss. ALBERT achieves significantly better performances than BERT on various benchmark MRC corpora such as SQuAD and RACE [26].

5.2 Evaluation metrics

Evaluations of MRC models on English and Chinese corpora [4, 5] used two evaluation metrics comprising exact match (EM) score and F1-score (macro-average). We also use these metrics to evaluate reading comprehension performances of humans and machine models on our corpus. In our Vietnamese MRC evaluations, punctuations and white spaces are ignored to allow normalization. The two evaluation metrics are described as follows.

- **EM** is the proportion of predicted answers that match exactly with the gold standard answers.
- **F1-score** measures the overlap between the predicted answer and the gold standard answer. First, we generated the predicted answer and gold answer as sets of tokens, and

then computed their F1-score. Finally, the evaluation system selected the maximum F1-score from all of the gold standard answers for each question, and then averaged them over all of the questions.

5.3 Experimental settings

To adapt the DrQA [1] and QANet [23] systems to suit our corpus, we fine-tune parameters of the systems. In particular, we only use word features for all of the experiments with DrQA and QANet models. We use the pyvi tool⁶ for word segmentation. We employ 300-dimensional Pho2Vec [60] word embeddings as our pre-trained Vietnamese word embeddings for DrQA and QANet. We choose Pho2Vec because it was trained on the corpus including nearly 19GB data³ generated by removing similar articles and duplication from a 50GB Vietnamese news set, which contains the same domain as our corpus. Besides, we set $batch\ size = 32$ and $epochs = 40$ for both the two models.

We use the base-cased multilingual BERT and set of parameters as follows: 12 layers, 768 hidden dimensions, and 12 attention heads (in the transformer) with 179M parameters and a vocabulary of about 120k vocabulary. We select the best hyperparameters by searching a combination of the batch size, learning rate and the number of fine-tuning epochs from the following ranges: learning rate: $2e - 5$, $3e - 5$, $5e - 5$; batch size: 4, 8, 16, 32; number of epochs: 2, 3. The best hyperparameters and models are selected based on the performance with the development set in Vietnamese.

For the ALBERT experiments, we use a pre-trained⁴ model for embedding input data. The training data for the embedding model is extracted from the Vietnamese texts with a vocabulary size of 30,000, tokenized by SentencePiece [61]. Input data is processed in two forms: original (ALBERT_{Base_Cased}) and lower case (ALBERT_{Base_Uncased}) to increase the learning ability of the model through a combination of assessment of form and semantics of each word in the corpus based on the multi-meanings of the Vietnamese language. We follow the ALBERT baseline, and tune several parameters on the development set, and obtain the best results with learning rate: $2e - 5$, batch size: 32 and epochs: 5.

For the BERT and ALBERT experiments, based on our corpus characteristics, we choose the maximum answer length to 200, the question length to 50, and the input sequence length to 512.

5.4 Estimation of human performance

To answer the question **Q2**, we measure the performance of humans with the development and test sets for our corpus. In particular, we hire three other workers to independently answer questions using the test and development sets, with four answers per question, as described in the phase for collecting additional answers (in Section 3.2.6). In contrast to Rajpurkar et al. [4], we use a cross-validation methodology to measure the performance of humans in a similar manner to and similar to Cui et al. [5]. In particular, we consider the first answer as the human prediction and treated the remainder of the answers as ground truths. We obtain four human prediction performance results by iteratively treating the first, second, third, and fourth answers as the human predictions. We calculate the maximum the performance over all of the ground truth answers for each question. Finally, we average the

⁶We use the pyvi library <https://pypi.org/project/pyvi/> for word segmentation.

³<https://github.com/binhvq/news-corpus>

⁴https://github.com/ngoanpv/albert_vi

four human prediction results as the final human performance result with on our corpus. Estimated human performances are presented in Table 18.

5.5 Experimental results

Table 18 compares the performance of the machine models and humans with the development and test sets for our corpus. Random guess and Sliding Window are the two methods with the lowest results, achieving less than 15% of the F1-score. Simple models cannot be used to deal with our data set. The transfer learning-based models (BERT and ALBERT) significantly outperform better than the neural network-based models (DrQA and QANet), but not as well as humans. The best model (ALBERT_{Base_Uncased}) achieves an EM of 65.26% and a F1-score of 84.89% on the test set. In particular, the performances (EM and F1-score) of the best model are better than that of DrQA, with 10.43% in EM and 10.80% in F1-score, respectively. The best model also perform better compared with the QANet model, where differences of 8.55% (in EM) and differences 5.10% (in F1-score). Meanwhile, the different performance between humans and the best model (14.53% and 10.90% differences in EM and F1-score, respectively) with our corpus was significant, thereby indicating that models for ViNewsQA should be improved in future research.

Table 18. Model and human performances on the development and test sets of ViNewsQA. Besides, we also build a approach based on random guess to compare with machine systems.

Type	MRC systems	EM (%)		F1-score(%)	
		Dev	Test	Dev	Test
Random	Random Guess	0.20	0.15	9.56	9.30
Rule-based	Sliding Window	0.32	0.15	13.11	13.38
Neural Network-based	DrQA	49.26	45.83	74.03	74.09
	QANet	57.80	56.71	78.39	79.79
Transfer learning-based	BERT	64.56	63.81	81.47	83.19
	ALBERT	64.68	64.46	83.43	84.16
	ALBERT	64.24	65.26	83.52	84.89
Human performance		75.19	79.79	92.77	95.79

6 RESULT ANALYSIS AND DISCUSSION

To obtain more insights into the performance of the neural network-based and transfer learning-based models and humans on our corpus, we analyze their performances in terms of different linguistic aspects comprising the question length (see in Section 6.1), answer length (see in Section 6.2), article length (see in Section 6.3), question type (see in Section 6.4), answer type (see in Section 6.5) and reasoning type (see in Section 6.6). These aspects are described in Section 4. In addition, we also aim to examine the impacts of the size of the training set as well as the vocabulary size on the machine models (see in Section 6.7). Finally, we perform qualitative analysis through typical examples (see Section 6.8). In this study, we omit analyzing these results on the lexical-based approach (Sliding Window) because its performance is significantly lower than other models.

6.1 Effects of question length

Firstly, we examine how well the reading comprehension models handle questions with different lengths. In particular, we analyze the performance of the machine models and

humans in terms of the F1-score metric. Table 19 presents the detailed analysis of performances with various questions lengths. In general, more accurate results are obtained for long (>16 words) than short and average-length (<15 words) questions on the best model. The difference in performance between humans and the best model decreases as the question length increases. A plausible explanation for this is that longer questions tend to contain more sufficient information in order to be able to extract the correct answer span, which is similar to what on SQuAD [62] and UIT-ViQuAD [43].

Table 19. Performance in terms of F1-score (%) according to the question length with the development set for our corpus. Δ is the performance difference between the human and the best-performance machine model (ALBERT).

Question length	DrQA	QANet	BERT	ALBERT	Human	Δ
4-5	72.56	73.76	86.36	85.14	95.60	+10.46
6-7	71.68	76.96	80.43	82.59	92.95	+10.36
8-9	74.69	79.92	79.46	82.01	92.91	+10.90
10-11	74.98	77.83	82.33	84.28	92.68	+8.40
12-13	73.18	77.74	83.85	85.00	93.34	+8.34
14-15	73.97	78.00	78.93	81.08	92.03	+10.95
16-17	77.07	82.63	81.70	85.85	93.02	+7.17
18-19	74.17	70.74	85.27	85.67	94.65	+8.98
>19	72.23	88.03	85.01	89.83	93.78	+3.95

6.2 Effects of answer length

In order to examine how well the reading comprehension models could predict the answers with different lengths, we analyze the performances of the machine models and humans in terms of the EM and F1-score. As can be seen from Table 20, our analysis shows the performances obtained with different answer lengths. In general, more accurate results are achieved for the shorter answer than longer answers. In particular, the best model achieves the highest performance with short answers (1–2 words, which accounts for 12.85% of the corpus) and the lowest performance with long answers (>18 words, which accounts for 14.42% of the corpus). In particular, answers with a length of 3 to 18 words witness fluctuation in performance. Hence, longer answers may be more difficult in finding answers. However, a difficult question is caused by many other factors such as such as reasoning type, answer type or question type. We examine this hypothesis in the following analysis.

Table 20. Performance in terms of F1-score (%) according to the answer length with the development set for our corpus. Δ is the performance difference between the human and the best-performance machine model (ALBERT).

Answer length	DrQA	QANet	BERT	ALBERT	Human	Δ
1-2	73.55	77.78	83.95	85.09	94.10	+9.01
3-4	79.46	76.69	84.44	82.73	92.51	+9.78
5-6	79.11	75.07	79.80	83.52	92.46	+8.94
7-8	78.78	78.77	79.95	84.08	91.89	+7.81
9-10	83.91	79.43	83.39	83.83	91.80	+7.97
11-12	79.03	79.89	80.31	83.16	92.65	+9.49
13-14	78.57	80.86	79.63	84.32	92.23	+7.91
15-16	76.73	82.96	85.20	84.91	91.99	+7.08
17-18	62.83	80.46	79.26	84.55	92.62	+8.07
19-20	56.13	77.70	79.99	77.77	91.80	+14.03
>20	50.65	79.36	79.01	80.99	91.87	+10.88

6.3 Effects of article length

In addition to determining the impacts of the question and answer lengths on the MRC model for the Vietnamese language, we analyze the performances of the MRC models and humans (in F1-score) with various article lengths. The detailed results are presented in and Table 21. In general, more accurate results are obtained for shorter articles than longer articles. In particular, all the models have a tendency to achieve better performances with short articles (<301 words). Because longer articles have higher interference and take more time to process. Hence, the MRC system has difficulty finding the answer for longer articles, which is useful for future improvements.

Table 21. Performance in terms of F1-score (%) according to the article length with the development set for our corpus. Δ is the performance difference between the human and the best-performance machine model (ALBERT).

Article length	DrQA	QANet	BERT	ALBERT	Human	Δ
<101	83.49	89.85	98.00	82.44	97.06	+14.62
101-200	78.07	81.45	85.79	87.21	94.93	+7.72
201-300	74.25	79.22	82.63	84.57	93.20	+8.63
301-400	71.83	76.58	79.26	82.37	92.87	+10.50
401-500	72.17	74.27	78.55	79.36	90.49	+11.13
501-600	74.01	80.53	81.33	82.94	92.37	+9.43
>600	72.49	76.76	75.41	81.12	91.94	+10.82

6.4 Effects of question type

We also analyze the performance of the model and humans in terms of the linguistic aspects based on the question type. Table 22 illustrates the detailed performances with different question types. In general, the "How," "Why," "What" and "Who" questions in our corpus are more difficult than others. In particular, ALBERT achieves the lowest performance on the question "Who". This is explained because "Who" has the lowest percentage of 3.16% (see Table 12) when compared to other types of questions in the corpus, and its answers with

complex-structure noun phrases are difficult to locate the beginning and ending positions of answers (see the first example in Section 6.8.1). The MRC system more readily extracts the correct answers for "Where", "When" and "Others" questions, and the best model achieves 88.16%, 84.47% and 89.14%, respectively, because their answers were mostly short, from 1 to 5 words (see Table 15). The differences in performance between humans and the best models are high for "Why", "How", "What", and "Who" questions (with differences in F1-scores over 9%). These are the types of questions needed to improve their performance in the future.

Table 22. Performance in terms of F1-score (%) according to the question type with the development set for our corpus. Δ is the performance difference between the human and the best-performance machine model (ALBERT).

Question type	DrQA	QANet	BERT	ALBERT	Human	Δ
Who	66.64	66.57	81.51	76.25	95.00	+18.75
When	76.07	77.92	82.10	84.47	91.61	+7.14
Where	72.96	74.08	83.12	88.16	95.53	+7.37
Why	74.60	76.45	79.69	81.37	93.17	+11.80
What	73.45	78.87	80.77	83.11	93.05	+9.94
How	70.82	76.99	78.10	83.15	91.66	+8.52
Others	83.04	84.95	91.18	89.25	92.93	+3.68

6.5 Effects of answer type

Besides, we also aim to verify the impacts of answer types to the MRC models. Table 23 show the EM and F1-score performances according to different types of answer. "Person" and "Others" are two difficult answer types, achieving the lowest performances including 77.14% (in F1-score) and 73.88% (in F1-score), respectively. The Person answer type has the same results as for the Who question because they are closely related (see the second example in Section 6.8.1). Besides, the Others answer type is very long compared to other types (see in Table 16), which is why it has low performance.

Table 23. Performance in terms of the F1-score (%) according to the answer type with the development set for our corpus. Δ is the performance difference between the human and the best-performance machine model (ALBERT).

Answer type	DrQA	QANet	BERT	ALBERT	Human	Δ
Time	77.47	75.15	82.43	85.18	93.33	+8.15
Other Numeric	83.95	86.74	90.80	89.27	92.71	+3.44
Person	73.49	82.13	77.51	77.14	97.85	+20.71
Location	52.66	63.62	85.67	85.83	98.47	+12.64
Other Entity	67.76	70.08	81.84	86.66	93.28	+6.62
Noun Phrase	75.20	76.37	82.13	83.44	93.02	+9.58
Verb Phrase	72.93	80.15	81.50	82.94	93.16	+10.22
Adjective Phrase	78.08	81.18	80.05	85.60	92.76	+7.16
Preposition Phrase	77.39	67.89	64.43	81.29	94.13	+12.84
Clause	72.37	78.96	80.40	84.92	92.23	+7.31
Others	64.73	73.22	68.06	73.88	90.47	+16.59

6.6 Effects of reasoning type

We examine how well the MRC models could handle answers with different reasoning types. Table 24 shows the performance with different reasoning types. In general, more accurate results are gained for non-inference questions (word matching and paraphrasing) than inference questions (single-sentence reasoning and multi-sentence reasoning). This analysis provides clear insights that more complex inference forms make it difficult for MRC systems.

The differences in performance (F1-score) between humans and the best model are significant with over 11% for inference questions, and with 3.3 to 8.4% for simple questions.

Table 24. Performance in terms of F1-score (%) according to the reasoning type with the development set for our corpus.

Reasoning type	DrQA	QANet	BERT	ALBERT	Human	Δ
Word matching	85.38	88.02	92.75	91.28	94.60	+3.32
Papaphrasing	75.62	80.73	83.10	84.96	93.38	+8.42
Single-sentence reasoning	71.50	74.94	75.38	81.08	92.49	+11.41
Multi-sentence reasoning	62.38	67.83	72.28	75.42	91.18	+15.76
Ambiguous/Insufficient	54.83	64.49	55.49	83.61	86.51	+2.90

6.7 Effects of training set size

Table 25. Vocabulary size of different amounts of training data.

#Questions	#Vocabulary size
1,992	9,816
3,990	13,935
5,990	17,076
7,990	19,763
9,990	22,079
11,980	24,137
13,985	26,062
15,985	27,880
17,583	29,142

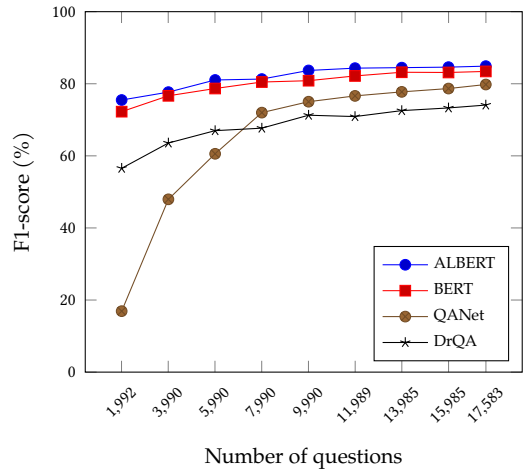


Fig. 10. Analysis and visualization of the model's result with different sizes of training data.

Our training data comprising 17,583 questions is much lower than the amount of the data used for English and Chinese MRC systems. To verify whether the small amount of training data affect the poor performance of the MRC systems different sizes training data including 1992, 3990, 5990, 7990, 9990, 11980, 13985, 15985, and 17583 questions. Table 25 shows that increasing data size means increasing vocabulary size. Figure 10 presents the performance (F1-score) on the test set of the ViNewsQA corpus. In general, the performance of the systems

is improved as the size of the training set increases from 1,992 to 17,583 questions. However, the F1-score increases steadily when the number of questions increased from 7,999 to 17,583 with DrQA and QANet. These observations indicate that the transfer learning models (BERT and ALBERT) are more effective with a small amount of training data compared with the other two models. However, we find that their performances still increase with enhancing the size of the training set. Consequently, increasing the training data size can improve the performance of the MRC models and depend on the model type.

6.8 Error analysis for typical examples

To obtain a better sense of what errors the reading comprehension systems are making on the ViNewsQA corpus, this work revolve two following questions: Q3: Why are the questions related to people not good at predicting models? and Q4: Which questions make neural network-based (DrQA and QANet) and transfer learning-based (BERT and ALBERT) systems all predict wrong?

6.8.1 Questions related to people. In order to answer the question Q3, we observe on many examples and select two typical samples for analysis, described as follows.

Context: Ngày 6/12, kíp mổ gồm 15 bác sĩ và 15 y tá Bệnh viện SS ở Varanasi, bang Uttar Pradesh, đã tiến hành ca phẫu thuật miễn phí cho cặp song sinh dính liền. (*On December 6, a team of 15 doctors and 15 nurses at SS Hospital in Varanasi, Uttar Pradesh state, performed the conjoined twins surgery for free.*)

Question: Ai phẫu thuật miễn phí? (*Who does the surgery for free?*)

Correct answer: kíp mổ gồm 15 bác sĩ và 15 y tá Bệnh viện SS ở Varanasi, bang Uttar Pradesh (*a team of 15 doctors and 15 nurses at SS Hospital in Varanasi, Uttar Pradesh state*)

DrQA prediction: 15 bác sĩ (*15 doctors*)

QANet prediction: song sinh dính liền (*conjoined twins*)

BERT prediction: kíp mổ gồm 15 bác sĩ và 15 y tá Bệnh viện SS ở Varanasi, bang Uttar Pradesh, đã tiến hành ca phẫu thuật miễn phí cho cặp song sinh dính liền (*a team of 15 doctors and 15 nurses at SS Hospital in Varanasi, Uttar Pradesh state, performed the conjoined twins surgery for free*)

ALBERT prediction: cặp song sinh dính liền (*the conjoined twins*)

"Who" questions having complex-structured answers such as noun phrases are difficult to be predicted by the MRC systems. Nguyen et al. [43] also found that noun-phrase answers are not easy to extract. Therefore, predicted answers are incorrect or deemed.

Context: Thai phụ mang thai lần hai 37 tuần, em bé ngôi thuận có dấu hiệu chuyển dạ, được đưa vào Bệnh viện Phụ sản Hải Phòng đêm 29/5. Phó giáo sư Vũ Văn Tâm, Giám đốc Bệnh viện Phụ sản Hải Phòng trực tiếp mổ sinh, bé trai chào đời nặng 2,8 kg với dây rốn dài 50 cm thắt nút đơn như bím tóc tết. Bác sĩ Tâm cho biết thai nhi có nút thắt dây rốn hiện vẫn là khó khăn cho chẩn đoán trước sinh, kể cả với những chuyên gia đầu ngành trên thế giới. (*Pregnant women with a second pregnancy of 37 weeks, her dominant baby showing signs of labor, was admitted to Hai Phong Obstetrics Hospital on the night of May 29. Associate Professor Vu Van Tam, Director of Hai Phong Obstetrics Hospital directly gave birth, the baby boy was born weighing 2.8 kg with a 50 cm umbilical cord tied in a single knot like a braid. Dr Tam said that the fetus with umbilical cord knot is still difficult for prenatal diagnosis, even with leading experts in the world.*)

Question: Ai là người thực hiện ca phẫu thuật? (*Who did perform the surgery?*)

Correct answer: Phó giáo sư Vũ Văn Tâm (*Associate Professor Vu Van Tam*)

DrQA prediction: những chuyên gia đầu ngành (*leading experts*)

QANet prediction: Phó giáo sư Vũ Văn Tâm (*Associate Professor Vu Van Tam*)

BERT prediction: Thai phụ mang thai (*Pregnant women*)

ALBERT prediction: Thai phụ mang thai lần hai 37 tuần, em bé ngôi thuận (*Pregnant women with a second pregnancy of 37 weeks, her dominant baby*)

The correct answer is Associate Professor Vu Van Tam (Associate Professor Tam Van Vu) including a title (Associate Professor) and a human name (Vu Van Tam). Following this phrase in the context is a phrase that adds information about Associate Professor Vu, which may interfere with the predicted result. According to the analysis in Section 6.5, the QANet model better predicts person-entity-based answers.

6.8.2 Predicted answers all models wrong. To answer the question Q4, we select and analyze the following context-question-answer triples, where neural network-based (DrQA and QANet) and transfer learning-based (BERT and ALBERT) systems all predict wrong.

Ambiguous answers: Multiple spans in the context can be selected as correct answers for a question, so they create ambiguity in the process of predicting the correct answer by the MRC systems. An example is provided as follows.

Context: Theo báo cáo nhanh của Bộ Y tế, ngày 6/2 tức mừng 2 Tết nguyên đán Kỷ Hợi, Bệnh viện Bạch Mai (Hà Nội) ghi nhận 2 người viêm phổi nặng nghi nhiễm cúm gia cầm nguy hiểm. Hai trường hợp này đang được Viện Vệ sinh dịch tễ Trung ương điều tra dịch tễ và lấy mẫu xét nghiệm. Kết quả sẽ có trong vài ngày tới. Theo Cục Thú y, Bộ Nông nghiệp và Phát triển nông thôn, hiện cả nước không có ổ dịch cúm nào xảy ra trên gia cầm, nhưng không loại trừ khả năng cúm gia cầm xuất hiện trên người. Cũng theo Bộ Y tế, trong 6 ngày nghỉ Tết nguyên đán, không ghi nhận trường hợp nào mắc bệnh sởi và bệnh liên cầu lợn, chỉ một số ổ dịch sốt xuất huyết tại An Giang, Bà Rịa - Vũng Tàu và Bến Tre. Ngoài ra còn có một ổ dịch quai bị ở Bến Tre. (*English translation: According to a quick report of the Ministry of Health, on February 6 the second day of the Ky Hoi Lunar New Year, Bach Mai Hospital (Hanoi) recorded 2 people with severe pneumonia suspected of being infected with dangerous avian influenza. These two cases are being investigated and tested by the Central Institute of Hygiene and Epidemiology. Results will be available in the next few days. According to the Department of Animal Health, Ministry of Agriculture and Rural Development, there is no flu outbreak in poultry nationwide, but it is not excluded that avian flu occurs in humans. According to the Ministry of Health, during the six days of the Tet holidays, no cases of measles and swine streptococcus were recorded, only a few outbreaks of dengue fever in An Giang, Ba Ria - Vung Tau, and Ben Tre. There is also a mumps outbreak in Ben Tre.*)

Question: Ở Bến Tre có trường hợp nhiễm phải căn bệnh nào? (*What disease is there in Ben Tre?*)

Correct answer: sốt xuất huyết. (*dengue fever.*)

DrQA prediction: bệnh sởi và bệnh liên cầu lợn. (*measles and swine streptococcus.*)

QANet prediction: viêm phổi nặng nghi nhiễm cúm gia cầm nguy hiểm. (*severe pneumonia suspected of being infected with dangerous avian influenza.*)

BERT prediction: bệnh sởi và bệnh liên cầu lợn. (*measles and swine streptococcus.*)

ALBERT prediction: sởi và bệnh liên cầu lợn. (*measles and swine streptococcus.*)

In this example, we find that two spans such as "sốt xuất huyết" (*dengue fever*) and "quai bị" (*mumps*) in the context are correct answers to the question "Ở Bến Tre có trường hợp nhiễm

phải căn bệnh nào?" (What disease is there in Ben Tre?). Hence, this question is ambiguous in finding the correct answer.

Incorrect boundary: The answers predicted by the four machine models are roughly equivalent to the correct answers, even though they lack or excess some words compared with the correct answers. This error is also easily caused by humans. These predicted answers have no meaning in the EM evaluation; however, they are calculated into the F1-score evaluation. An example is given as follows.

Context: Hút thuốc lá hay thuốc lá đều có hại cho sức khỏe của chính bản thân người hút và những người xung quanh. Thuốc lá có hàm lượng nicotin khoảng 9%, cao hơn nhiều so với thuốc lá thông thường (khoảng 1- 3%). (*Tobacco or pipe tobacco smoking is both harmful to the health of the smokers and those around them. Pipe tobacco has a nicotine content of about 9%, much higher than regular tobacco (about 1- 3%).*)

Question: Hút thuốc lá có hại như thế nào đối với con người? (*How harmful is smoking for humans?*).

Correct answer: có hại cho sức khỏe của chính bản thân người hút và những người xung quanh. (*harmful to the health of the smokers and those around them*).

DrQA prediction: bản thân người hút và những người xung quanh. (*the health of the smokers and those around them*).

QANet prediction: đều có hại cho sức khỏe của chính bản thân người hút và những người xung quanh. (*both harmful to the health of the smokers and those around them*).

BERT prediction: cho sức khỏe của chính bản thân người hút và những người xung quanh. (*to the health of the smokers and those around them*).

ALBERT prediction: có hại cho sức khỏe. (*harmful to the health*).

In this question, compared with the correct answer, the DrQA, QANet, BERT answer predictions lack or excess a word and several words. However, these answers (DrQA, QANet, and BERT predictions) are deemed, which could be accepted as the correct answers.

Incorrect inference: According to the analysis in Sub-section 6.5, the questions inferred based on single-or-multiple-sentence information have lower results than the non-inference questions. We select an example for this error, as described follows.

Context: Theo SCMP, các bác sĩ từng cảnh báo Wu Ying "hoàn toàn không phù hợp để mang thai" vì bệnh tim bẩm sinh cùng chứng tăng huyết áp phổi. Chồng Wu là Shen Jie cũng sẵn sàng từ bỏ đứa bé để cứu vợ, song Wu vẫn kiên quyết giữ con. "Nhiều người bảo tôi cứng đầu nhưng họ đều đã có con rồi", Wu trái lòng. "Mỗi lần nhìn con cái họ, tôi lại muốn sinh ra đứa con của chính mình. Tôi hiểu rất rõ các rủi ro nhưng sẵn sàng đánh cược". Wu từng sảy thai hai lần. Tháng 5/2017, Wu sinh con trai nặng một kg bằng phương pháp đẻ mổ. (*According to SCMP, doctors warned Wu Ying "completely unsuitable for pregnancy" because of congenital heart disease and pulmonary hypertension. Wu's husband Shen Jie is also willing to give up the baby to save his wife, but Wu is determined to keep the child. "Many people told me to be stubborn, but they all have children," Wu said. "Every time I see their children, I want to give birth to my own child. I understand the risks very well but am willing to bet." Wu had miscarried twice. In May 2017, Wu gave birth to a one-kilogram son by cesarean section.*)

Question: Vì sao Wu Ying vẫn quyết định giữ đứa bé? (*Why did Wu Ying still decide to keep the baby?*).

Correct answer: muốn sinh ra đứa con của chính mình. (*want to give birth to my own child.*).

DrQA prediction: Nhiều người bảo tôi cứng đầu nhưng họ đều đã có con rồi. (*Many people told me to be stubborn, but they both had children.*).

QANet prediction: bệnh tim bẩm sinh cùng chứng tăng huyết áp phổi. Chồng Wu là Shen Jie cũng sẵn sàng từ bỏ đứa bé để cứu vợ. (*congenital heart disease and pulmonary hypertension. Wu's husband Shen Jie is also willing to give up the baby to save his wife.*).

BERT prediction: Nhiều người bảo tôi cứng đầu nhưng họ đều đã có con rồi. (*Many people told me to be stubborn, but they all have children.*).

ALBERT prediction: vì bệnh tim bẩm sinh cùng chứng tăng huyết áp phổi. (*because of congenital heart disease and pulmonary hypertension.*)

To answer the question above, the human or MRC systems must understand and connect the contents of the first four sentences in the context to find the correct answer. Therefore, the predicted answers are chosen from a span of the first three sentences instead of a segment from the content of the 4th sentence.

Lack of world knowledge: Several questions require knowledge to determine the answers. For example, these questions may use equivalent words or phrases which can be terms or concepts of a specific field. We select an example for this error, as described follows.

Context: Theo tiến sĩ Đào Văn Long, nguyên Trưởng khoa Tiêu hóa Bệnh viện Bạch Mai, người bị đau dạ dày thường do vi khuẩn HP, stress, lạm dụng chất kích thích và thói quen ăn uống không hợp lý. Chế độ ăn của người đau dạ dày cần giảm tác dụng của axit lên niêm mạc dạ dày, hạn chế hoặc bỏ những kích thích có hại để dạ dày nghỉ ngơi và các tổn thương mau lành. Do đó, người đau dạ dày cần ăn đúng giờ, hạn chế chất kích thích và chọn gia vị phù hợp. "Người đau dạ dày cần ăn uống nghiêm ngặt và cầu kỳ", tiến sĩ Long khuyên. (*According to Dr. Dao Van Long, the former head of the Department of Gastroenterology at Bach Mai Hospital, people with stomach pain are often caused by HP bacteria, stress, substance abuse and inappropriate eating habits. The diet of the stomach ache should reduce the effect of the acid on the gastric mucosa, limit or eliminate harmful stimuli to rest the stomach and heal the damage. Therefore, people with stomach ache need to eat on time, limit stimulants and choose the right spices. "People with stomach ache need strict and sophisticated diet", Dr. Long advised.*)

Question: Người bị đau dạ dày cần kiêng những gì? (*What should people with stomachache abstain from?*)

Correct answer: giảm tác dụng của axit lên niêm mạc dạ dày, hạn chế hoặc bỏ những kích thích có hại để dạ dày nghỉ ngơi và các tổn thương mau lành (*reduce the effect of the acid on the gastric mucosa, limit or eliminate harmful stimuli to rest the stomach and heal the damage.*).

DrQA prediction: ăn uống nghiêm ngặt và cầu kỳ (*strict and sophisticated diet.*).

QANet prediction: nghiêm ngặt và cầu kỳ (*strict and sophisticated diet.*).

BERT prediction: ăn uống nghiêm ngặt và cầu kỳ (*strict and sophisticated diet.*).

ALBERT prediction: ăn đúng giờ, hạn chế chất kích thích và chọn gia vị phù hợp (*eat on time, limit stimulants and choose the right spices.*).

In the context, there are three advices from a doctor that should be improved if you encounter stomachache. Therefore, this question is ambiguous to find the answer. If the systems has knowledge of the verb "kiêng" ("kiêng" means to avoid eating, not consuming

certain foods or doing certain things, because it is harmful or considered harmful to health.), it is easy to find the correct answer from the answer candidates.

7 CONCLUSION AND FUTURE WORK

This paper introduced ViNewsQA, a span-extraction corpus for evaluating intelligent reading comprehension systems and question-answering in a low-resource language like Vietnamese. Over 22,000 question-answer pairs were generated by humans based on a set of 4,416 online health news articles in our corpus. Our corpus contains diverse answer types, and a significant proportion of questions (42.25% of ViNewsQA) required complex reasoning ability to solve. The corpus is challenging because our evaluation results showed that the difference in performance between humans and the best model was significant (an EM difference of 14.53% and an F1-score difference of 10.90%). Analyses of the experimental results showed that better performances were obtained for long questions with more information than short questions, whereas shorter answers and articles tend to yield better performances. Additionally, we realized that our corpus has difficult question types (What, How, Why, and Who) and complex reasoning based on a sentence or connections between multiple sentences. Finally, we explored the qualitative analysis errors consisting of ambiguous answers, incorrect boundaries, incorrect inference and lack of world knowledge.

By its size and complexity, ViNewsQA makes a significant extension to the existing machine reading comprehension corpora. For example, our corpus can be used for cross-lingual studies based on experiments with other similar corpora, such as SQuAD, NewsQA, and CMRC. We hope that our corpus will stir more research in machine reading comprehension and guide the development of artificial-intelligence applications. In particular, we conduct further investigations to solve the difficult questions that require comprehensive reasoning based on multiple sentences in the article. Moreover, we would like to operate a Vietnamese MRC challenging shared task for researchers to conduct experiments to explore better models with our corpus. Finally, we aim to build a modern QA system based on DrQA [1] which is necessary to serve the searching demands of nearly 100M Vietnamese people.

ACKNOWLEDGMENTS

We would like to thank the editors and anonymous reviewers for their helpful feedback.

REFERENCES

- [1] Danqi Chen, Adam Fisch, Jason Weston, and Antoine Bordes. 2017. Reading wikipedia to answer open-domain questions. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 1870–1879.
- [2] Shuohang Wang, Mo Yu, Xiaoxiao Guo, Zhiguo Wang, Tim Klinger, Wei Zhang, Shiyu Chang, Gerry Tesauro, Bowen Zhou, and Jing Jiang. 2018. R 3: reinforced ranker-reader for open-domain question answering. In *Thirty-Second AAAI Conference on Artificial Intelligence*.
- [3] Deepak Gupta, Asif Ekbal, and Pushpak Bhattacharyya. 2019. A deep neural network framework for english hindi question answering. *ACM Transactions on Asian and Low-Resource Language Information Processing (TALLIP)*, 19, 2, 1–22.
- [4] Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. Squad: 100,000+ questions for machine comprehension of text. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, 2383–2392.
- [5] Yiming Cui, Ting Liu, Wanxiang Che, Li Xiao, Zhipeng Chen, Wentao Ma, Shijin Wang, and Guoping Hu. 2019. A span-extraction dataset for chinese machine reading

- comprehension. 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)*, 5886–5891.
- [6] Pranav Rajpurkar, Robin Jia, and Percy Liang. 2018. Know what you don't know: unanswerable questions for SQuAD. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*. Association for Computational Linguistics, Melbourne, Australia, (July 2018), 784–789. doi: 10.18653/v1/P18-2124. <https://www.aclweb.org/anthology/P18-2124>.
- [7] Adam Trischler, Tong Wang, Xingdi Yuan, Justin Harris, Alessandro Sordani, Philip Bachman, and Kaheer Suleman. 2017. NewsQA: a machine comprehension dataset. In *Proceedings of the 2nd Workshop on Representation Learning for NLP*. Association for Computational Linguistics, Vancouver, Canada, (August 2017), 191–200. doi: 10.18653/v1/W17-2623. <https://www.aclweb.org/anthology/W17-2623>.
- [8] Karl Moritz Hermann, Tomas Kocisky, Edward Grefenstette, Lasse Espeholt, Will Kay, Mustafa Suleyman, and Phil Blunsom. 2015. Teaching machines to read and comprehend. In *Advances in neural information processing systems*, 1693–1701.
- [9] Felix Hill, Antoine Bordes, Sumit Chopra, and Jason Weston. 2015. The goldilocks principle: reading children's books with explicit memory representations. *arXiv preprint arXiv:1511.02301*.
- [10] Yiming Cui, Ting Liu, Zhipeng Chen, Shijin Wang, and Guoping Hu. 2016. Consensus attention-based neural networks for chinese reading comprehension. In *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers*, 1777–1786.
- [11] Matthew Richardson, Christopher JC Burges, and Erin Renshaw. 2013. Mctest: a challenge dataset for the open-domain machine comprehension of text. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, 193–203.
- [12] Guokun Lai, Qizhe Xie, Hanxiao Liu, Yiming Yang, and Eduard Hovy. 2017. Race: large-scale reading comprehension dataset from examinations. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, 785–794.
- [13] Siva Reddy, Danqi Chen, and Christopher D Manning. 2019. Coqa: a conversational question answering challenge. *Transactions of the Association for Computational Linguistics*, 7, 249–266.
- [14] Kai Sun, Dian Yu, Jianshu Chen, Dong Yu, Yejin Choi, and Claire Cardie. 2019. Dream: a challenge data set and models for dialogue-based reading comprehension. *Transactions of the Association for Computational Linguistics*, 7, 217–231.
- [15] Chih Chieh Shao, Trois Liu, Yuting Lai, Yiyang Tseng, and Sam Tsai. 2018. Drcd: a chinese machine reading comprehension dataset. *arXiv preprint arXiv:1806.00920*.
- [16] Wei He, Kai Liu, Jing Liu, Yajuan Lyu, Shiqi Zhao, Xinyan Xiao, Yuan Liu, Yizhong Wang, Hua Wu, Qiaoqiao She, et al. 2018. Dureader: a chinese machine reading comprehension dataset from real-world applications. *ACL 2018*, 37.
- [17] Seungyoung Lim, Myungji Kim, and Jooyoul Lee. 2019. Korquad1. 0: korean qa dataset for machine reading comprehension. *arXiv preprint arXiv:1909.07005*.
- [18] Shuohang Wang and Jing Jiang. 2016. Machine comprehension using match-lstm and answer pointer. *arXiv preprint arXiv:1608.07905*.
- [19] Minjoon Seo, Aniruddha Kembhavi, Ali Farhadi, and Hannaneh Hajishirzi. 2016. Bidirectional attention flow for machine comprehension. *arXiv preprint arXiv:1611.01603*.
- [20] Wenhui Wang, Nan Yang, Furu Wei, Baobao Chang, and Ming Zhou. 2017. Gated self-matching networks for reading comprehension and question answering. In *Proceedings*

- of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 189–198.
- [21] Hsin-Yuan Huang, Chenguang Zhu, Yelong Shen, and Weizhu Chen. 2018. Fusionnet: fusing via fully-aware attention with application to machine comprehension. In *International Conference on Learning Representations*.
- [22] Dirk Weissenborn, Georg Wiese, and Laura Seiffe. 2017. Making neural qa as simple as possible but not simpler. In *Proceedings of the 21st Conference on Computational Natural Language Learning (CoNLL 2017)*, 271–280.
- [23] Adams Wei Yu, David Dohan, Minh-Thang Luong, Rui Zhao, Kai Chen, Mohamad Norouzi, and Quoc V Le. 2018. Qanet: combining local convolution with global self-attention for reading comprehension. In *International Conference on Learning Representations*.
- [24] Cheoneum Park, Heejun Song, and Changki Lee. 2020. S3-net: sru-based sentence and self-matching networks for machine reading comprehension. *ACM Transactions on Asian and Low-Resource Language Information Processing (TALLIP)*, 19, 3, 1–14.
- [25] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. 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)*, 4171–4186.
- [26] Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. 2019. Albert: a lite bert for self-supervised learning of language representations. In *International Conference on Learning Representations*.
- [27] Phuong-Thai Nguyen, Xuan-Luong Vu, Thi-Minh-Huyen Nguyen, Van-Hiep Nguyen, and Hong-Phuong Le. 2009. Building a large syntactically-annotated corpus of vietnamese. In *Proceedings of the third linguistic annotation workshop*. Association for Computational Linguistics, 182–185.
- [28] Dat Quoc Nguyen, Dai Quoc Nguyen, Son Bao Pham, Phuong-Thai Nguyen, and Minh Le Nguyen. 2014. From treebank conversion to automatic dependency parsing for vietnamese. In *International Conference on Applications of Natural Language to Data Bases/Information Systems*. Springer, 196–207.
- [29] Kiet V Nguyen and Ngan Luu-Thuy Nguyen. 2016. Vietnamese transition-based dependency parsing with supertag features. In *2016 Eighth International Conference on Knowledge and Systems Engineering (KSE)*. IEEE, 175–180.
- [30] Binh Duc Nguyen, Kiet Van Nguyen, and Ngan Luu-Thuy Nguyen. 2018. Lstm easy-first dependency parsing with pre-trained word embeddings and character-level word embeddings in vietnamese. In *2018 10th International Conference on Knowledge and Systems Engineering (KSE)*. IEEE, 187–192.
- [31] Dat Quoc Nguyen, Dai Quoc Nguyen, Dang Duc Pham, and Son Bao Pham. 2014. RDRPOSTagger: a ripple down rules-based part-of-speech tagger. In: *Proceedings of the Demonstrations at the 14th Conference of the European Chapter of the Association for Computational Linguistics*, (April 2014), 17–20.
- [32] Ngo Xuan Bach, Nguyen Dieu Linh, and Tu Minh Phuong. 2018. An empirical study on POS tagging for vietnamese social media text. *Computer Speech & Language*, 50, 1–15.
- [33] Pham Thi Xuan Thao, Tran Quoc Tri, Dinh Dien, and Nigel Collier. 2008. Named entity recognition in vietnamese using classifier voting. *ACM Transactions on Asian Language Information Processing*, 6, 4, (December 2008). ISSN: 1530-0226. DOI: 10.1145/1316457.1316460. <https://doi.org/10.1145/1316457.1316460>.

- [34] Long H. B. Nguyen, Dien Dinh, and Phuoc Tran. 2016. An approach to construct a named entity annotated english-vietnamese bilingual corpus. *ACM Trans. Asian Low-Resour. Lang. Inf. Process.*, 16, 2, (October 2016). issn: 2375-4699. doi: 10.1145/2990191. <https://doi.org/10.1145/2990191>.
- [35] Binh An Nguyen, Kiet Van Nguyen, and Ngan Luu-Thuy Nguyen. 2019. Error analysis for vietnamese named entity recognition on deep neural network models. *arXiv preprint arXiv:1911.07228*.
- [36] Kiet Van Nguyen, Vu Duc Nguyen, Phu XV Nguyen, Tham TH Truong, and Ngan Luu-Thuy Nguyen. 2018. Uit-vfsc: vietnamese students' feedback corpus for sentiment analysis. In *2018 10th International Conference on Knowledge and Systems Engineering (KSE)*. IEEE, 19–24.
- [37] Phu XV Nguyen, Tham TT Hong, Kiet Van Nguyen, and Ngan Luu-Thuy Nguyen. 2018. Deep learning versus traditional classifiers on vietnamese students' feedback corpus. In *2018 5th NAFOSTED Conference on Information and Computer Science (NICS)*. IEEE, 75–80.
- [38] Thin Van Dang, Vu Duc Nguyen, Kiet Van Nguyen, and Ngan Luu-Thuy Nguyen. 2018. A transformation method for aspect-based sentiment analysis. *Journal of Computer Science and Cybernetics*, 34, 4, 323–333.
- [39] Dai Quoc Nguyen, Dat Quoc Nguyen, and Son Bao Pham. 2009. A vietnamese question answering system. In *2009 International Conference on Knowledge and Systems Engineering*. IEEE.
- [40] Van-Tu Nguyen and Anh-Cuong Le. 2016. Improving question classification by feature extraction and selection. *Indian Journal of Science and Technology*, 9, 17, 1–8.
- [41] Phuong Hong Le and Duc-Thien Bui. 2018. A factoid question answering system for vietnamese. In *Companion Proceedings of the The Web Conference 2018*. International World Wide Web Conferences Steering Committee, 1049–1055.
- [42] Kiet Van Nguyen, Khiem Vinh Tran, Son T Luu, Anh Gia-Tuan Nguyen, and Ngan Luu-Thuy Nguyen. 2020. Enhancing lexical-based approach with external knowledge for vietnamese multiple-choice machine reading comprehension. *IEEE Access*, 8, 201404–201417.
- [43] Kiet Nguyen, Vu Nguyen, Anh Nguyen, and Ngan Nguyen. 2020. A Vietnamese dataset for evaluating machine reading comprehension. In *Proceedings of the 28th International Conference on Computational Linguistics*. International Committee on Computational Linguistics, Barcelona, Spain (Online), (December 2020), 2595–2605. <https://www.aclweb.org/anthology/2020.coling-main.233>.
- [44] Qizhe Xie, Guokun Lai, Zihang Dai, and Eduard Hovy. 2018. Large-scale cloze test dataset created by teachers. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Brussels, Belgium, 2344–2356. doi: 10.18653/v1/D18-1257. <https://www.aclweb.org/anthology/D18-1257>.
- [45] Yiming Cui, Ting Liu, Ziqing Yang, Zhipeng Chen, Wentao Ma, Wanxiang Che, Shijin Wang, and Guoping Hu. 2020. A sentence cloze dataset for Chinese machine reading comprehension. In *Proceedings of the 28th International Conference on Computational Linguistics*. International Committee on Computational Linguistics, Barcelona, Spain (Online), (December 2020), 6717–6723. <https://www.aclweb.org/anthology/2020.coling-main.589>.

- [46] Wei Yang, Yuqing Xie, Aileen Lin, Xingyu Li, Luchen Tan, Kun Xiong, Ming Li, and Jimmy Lin. 2019. End-to-end open-domain question answering with bertserini. *arXiv preprint arXiv:1902.01718*.
- [47] Simon Šuster and Walter Daelemans. 2018. CliCR: a dataset of clinical case reports for machine reading comprehension. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*. Association for Computational Linguistics, New Orleans, Louisiana, (June 2018), 1551–1563. doi: 10.18653/v1/N18-1140. <https://www.aclweb.org/anthology/N18-1140>.
- [48] Xiao Zhang, Ji Wu, Zhiyang He, Xien Liu, and Ying Su. 2018. Medical exam question answering with large-scale reading comprehension. In *Thirty-Second AAAI Conference on Artificial Intelligence*.
- [49] Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William Cohen, and Xinghua Lu. 2019. PubMedQA: a dataset for biomedical research 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)*. Association for Computational Linguistics, Hong Kong, China, (November 2019), 2567–2577. doi: 10.18653/v1/D19-1259. <https://www.aclweb.org/anthology/D19-1259>.
- [50] Martin d’Hoffschmidt, Wacim Belblidia, Quentin Heinrich, Tom Brendlé, and Maxime Vidal. 2020. FQuAD: French question answering dataset. In *Findings of the Association for Computational Linguistics: EMNLP 2020*. Association for Computational Linguistics, Online, (November 2020), 1193–1208. <https://www.aclweb.org/anthology/2020.findings-emnlp.107>.
- [51] Pavel Efimov, Andrey Chertok, Leonid Boytsov, and Pavel Braslavski. 2020. Sberquad – russian reading comprehension dataset: description and analysis. In *Experimental IR Meets Multilinguality, Multimodality, and Interaction*. Avi Arampatzis, Evangelos Kanoulas, Theodora Tsirikla, Stefanos Vrochidis, Hideo Joho, Christina Lioma, Carsten Eickhoff, Aurélie Névéal, Linda Cappellato, and Nicola Ferro, editors. Springer International Publishing, Cham, 3–15. isbn: 978-3-030-58219-7.
- [52] Tomasz Jurczyk, Michael Zhai, and Jinho D Choi. 2016. Selqa: a new benchmark for selection-based question answering. In *2016 IEEE 28th International Conference on Tools with Artificial Intelligence (ICTAI)*. IEEE, 820–827.
- [53] Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. Unsupervised cross-lingual representation learning at scale. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, Online, (July 2020), 8440–8451. doi: 10.18653/v1/2020.acl-main.747. <https://www.aclweb.org/anthology/2020.acl-main.747>.
- [54] Anna Aniol, Marcin Pietron, and Jerzy Duda. 2019. Ensemble approach for natural language question answering problem. In *2019 Seventh International Symposium on Computing and Networking Workshops (CANDARW)*. IEEE, 180–183.
- [55] Mandar Joshi, Danqi Chen, Yinhan Liu, Daniel S Weld, Luke Zettlemoyer, and Omer Levy. 2020. Spanbert: improving pre-training by representing and predicting spans. *Transactions of the Association for Computational Linguistics*, 8, 64–77.

- [56] Dheeru Dua, Yizhong Wang, Pradeep Dasigi, Gabriel Stanovsky, Sameer Singh, and Matt Gardner. 2019. Drop: a reading comprehension benchmark requiring discrete reasoning over paragraphs. 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)*, 2368–2378.
- [57] Matthew Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, 2227–2237.
- [58] Jeremy Howard and Sebastian Ruder. 2018. Universal language model fine-tuning for text classification. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, Melbourne, Australia, (July 2018), 328–339. DOI: 10.18653/v1/P18-1031. <https://www.aclweb.org/anthology/P18-1031>.
- [59] Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. [n. d.] Improving language understanding by generative pre-training.
- [60] Anh Tuan Nguyen, Mai Hoang Dao, and Dat Quoc Nguyen. 2020. A pilot study of text-to-SQL semantic parsing for Vietnamese. In *Findings of the Association for Computational Linguistics: EMNLP 2020*. Association for Computational Linguistics, Online, (November 2020), 4079–4085. DOI: 10.18653/v1/2020.findings-emnlp.364. <https://www.aclweb.org/anthology/2020.findings-emnlp.364>.
- [61] Taku Kudo and John Richardson. 2018. SentencePiece: a simple and language independent subword tokenizer and detokenizer for neural text processing. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*. Association for Computational Linguistics, Brussels, Belgium, (November 2018), 66–71. DOI: 10.18653/v1/D18-2012. <https://www.aclweb.org/anthology/D18-2012>.
- [62] Soumya Wadhwa, Khyathi Chandu, and Eric Nyberg. 2018. Comparative analysis of neural qa models on squad. In *Proceedings of the Workshop on Machine Reading for Question Answering*, 89–97.