

Vietnamese Fake News Generation

Category: Natural Language

I. Introduction

In today's age of information, online news consumers face the challenge of distinguishing between fake and genuine news, which has led to an increase in research on methods to identify fake news. Initially, we aim to detect Vietnamese fake news following the method of Wu et al. (2022). However, we found that there is a lack of sufficient resources for labeled data to detect Vietnamese fake news, prompting us to re-evaluate the scope and direction of our project. Our revised final delivery is a pipeline for generating fake news from reliable news, and two Vietnamese datasets – reliable news and corresponding fake news, hoping that our application can be utilized in further research on Vietnamese fake news detection.

II. Implementation

1. Abstract

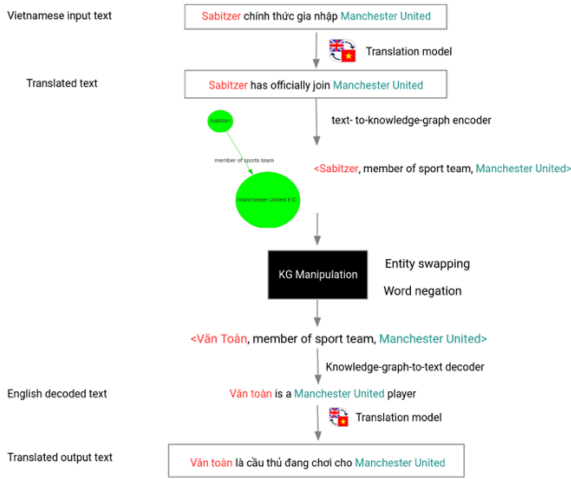


Figure 1. Illustration of our Pipeline

The pipeline for our fake news generation approach comprises of four models and four corresponding steps. The first step is to translate Vietnamese input text into English using a translation model. This is necessary due to a lack of proper datasets and models for encoding Vietnamese text into a knowledge graph (KG). In the next step, the translated English sentence is encoded into a graph structure for easy manipulation. English is preferred because of the availability of WebNLG, a well-structured dataset that maps well-written English sentences with KGs. Thirdly, the encoded graph is manipulated to create a new graph representing the fake news story. The manipulation may involve changing the relationships between the entities or adding new entities altogether. Finally, the manipulated graph is decoded by a KG-

to-text model, which generates the fake news story in English. The generated English text is then translated back to Vietnamese to produce the final output.

2. Translation model

The translation model used in the project is an adapted version of the translation machine developed by VinAI, which can convert speech and text between Vietnamese and English. Through experiments, the VinAI system has been shown to have state-of-the-art performance and successfully employ the recent cutting-edged neural models, including Automatic Speech Recognition (ASR), Machine Translation (MT), and TextTo-Speech (TTS). In our project, we will only use the MT component in the VinAI system to translate the input between Vietnamese and English. The component is developed by first fine-tuning the mBART, a pre-trained sequence-to-sequence model, using 3M training sentence pairs from the high-quality PhoMT dataset. Then, this model will be employed to convert the English sentence into Vietnamese from each English-Vietnamese sentence pair in CCAI and WikiMatrix datasets. Specifically, only the pairs with the BLEU score in the range of 0.15 to 0.95 between the translated Vietnamese variant from the English source and the Vietnamese target sentence are chosen.

3. Text-to-KG encoder

3.1. OneIE

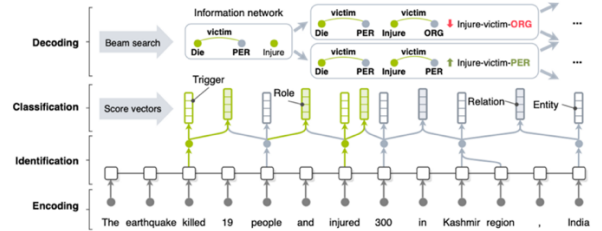


Figure 2. An illustration of the end-to-end joint information extraction framework ONE-IE at the test stage.

One of the methods recommended to encode the text-to-knowledge graph is using the ONE-IE model. This model is proposed to decrease the possibilities of errors made by local classifiers without the global restrictions and could be used regardless of language features. There are 4 phrases in implementing the model, including encoding, identification, classification, and decoding. In the encoding part, we will use a pre-trained BERT encoder to contextualize and represent the given sentence. For the next phrase, we will identify the entity mentions and event triggers as nodes and compute the label

scores for all the nodes and their pair wise links using local classifiers in the classification stage. In the final step, a beam decoder embedded with global features will search for global optimal graph and capture the cross-subtask and cross-instance interactions. Finally, the model will return the information network that has the highest global score. In the project, we have tested the performance of the model by giving it the input as document since we want to use this for generating fake news. However, the results are low and only efficient if the text is in sentence format.

3.2. REBEL

REBEL (Cabot et. al.) is a new approach to Relation Extraction, which is a task that involves identifying relationships between different entities in text. REBEL uses an autoregressive model that generates output sequentially, and frames Relation Extraction as a sequence-to-sequence task. To train the model, the authors created a new dataset called REBEL, which is a large-scale distantly supervised dataset obtained by leveraging a Natural Language Inference model.

REBEL's approach is different from previous end-to-end approaches because it uses a simple triplet decomposition into a text sequence. The model used is an Encoder-Decoder Transformer called BART, which is pre-trained using the REBEL dataset. This allows the model to leverage both the encoded input and the previously decoded output, leading to better performance in Relation Extraction. According to the authors, after a few epochs of fine-tuning, REBEL achieves state-of-the-art performance on a variety of Relation Extraction baselines.

The simplicity of REBEL's approach makes it highly flexible and adaptable to new domains or longer documents.

3.3 Experiment and result

In this section, we evaluated how well the REBEL model performs on the CONLL04 dataset, which is commonly used for identifying relationships between entities in text. Even though the model was trained on an autoregressive task, we tested its performance on relation extraction (RE) by extracting all the relationships in its generated output. We used Recall, Precision, and micro-F1 to evaluate the model's performance, based on the labeled relationships in the dataset. The test is based on CONLL04 dataset (Roth & Yih, 2004) which consists of news article sentences that are labeled with four entity types (person, organization, location, and other) and five types of relationships (kill, work for, organization based in, live in, and located in). We fine-tuned the REBEL model for 30 iterations, following the guidance of the dataset's original authors, and tested it on the best-performing iteration that was determined by its performance on a validation set. Our evaluation found that the model had an average F1-score

of 70.26% on the CONLL04 dataset, which was slightly lower than the original experiment's performance of 71.97%.

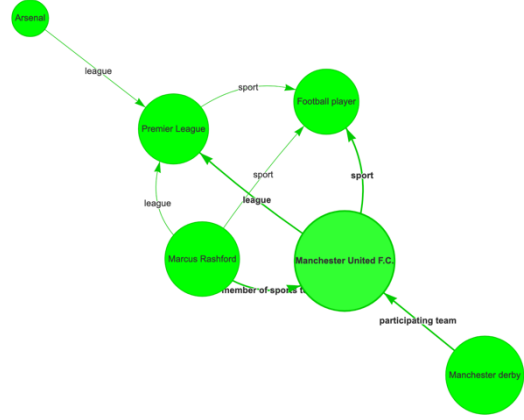


Figure 3. Knowledge Graph generated by REBEL after Marcus Rashford's Wikipedia abstract.

```
{
  "head": "Marcus Rashford",
  "type": "sport",
  "tail": "Football player",
  "meta": {
    "spans": [[0, 128]]
  }
},
{
  "head": "Premier League",
  "type": "sport",
  "tail": "Football player",
  "meta": {
    "spans": [[0, 128]]
  }
},
{
  "head": "Manchester United F.C.",
  "type": "sport",
  "tail": "Football player",
  "meta": {
    "spans": [[0, 128]]
  }
},
{
  "head": "Manchester United F.C.",
  "type": "league",
  "tail": "Premier League",
  "meta": {
    "spans": [[0, 128], [22, 150]]
  }
},
{
  "head": "Marcus Rashford",
  "type": "league",
  "tail": "Premier League",
  "meta": {
    "spans": [[0, 128]]
  }
},
{
  "head": "Arsenal",
  "type": "league",
  "tail": "Premier League",
  "meta": {
    "spans": [[22, 150]]
  }
},
{
  "head": "Manchester derby",
  "type": "participating team",
  "tail": "Manchester United F.C.",
  "meta": {
    "spans": [[22, 150]]
  }
}
```

Figure 4. Relations corresponding to KG in Figure 3

4. Knowledge Graph Manipulation

In the previous step, we transformed human-written sentence structures into a KG as it provides a straightforward and well-defined data structure with event entities as nodes. In this section, we will demonstrate how we leveraged the graph structure's simplicity to manipulate specific details within the news and steer it in our desired direction.

4.1. Entity swapping

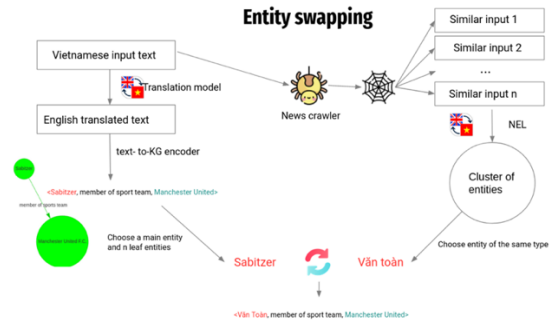


Figure 5. Overall process of entity swapping

The concept of entity switching is illustrated in the above figure. To further enhance the quality of the fake news, we added a crucial component at this stage: the *news crawler*. The news crawler's task is to scan through a selection of reputable and trusted Vietnamese news websites, handpicked by our team, and retrieve articles within a specific time frame (e.g., news articles from 1/2/2023 to 12/2/2023). From these news

[illegible]

From the encoded KG of the original text, we chose certain entities for manipulation. As suggested by Fung et al. (2021), nodes with the highest degree of connectivity are most critical, while those with the lowest degree make a limited contribution. Therefore, we chose the highest-degree node and a number (determined by a hyperparameter 'n') of lowest-degree nodes for manipulation. Next, we replaced the chosen entities in original input text with entities from our pre-constructed dictionary that are of the same type as the chosen entities, forming a new KG. Note that we should only use entities in the dictionary that did not appear in the original text for replacement.

As the cosine similarity did not perform as expected, we experimented with a different approach to measure the similarity of the news articles in our crawled data set. Our new approach involves representing each article in the data set using the term frequency-inverse document frequency (TF-IDF) representation, which takes into account the frequency of

258 Jürgen Klopp đang có gắng dùng xây một kỷ nguyên.
260 Trùng trùng ngày 31/1 chỉ trích Tổng thống đắc
277 Mykhailo Mudryk quá rắc rối khi nhắc đến các
282 Yussuke Adachi Giám đốc Kyuetsu của Liền...
297 5 năm dưới triều đại của HLV Park Hang Seo
301 Tổng giám đốc VPF Nguyễn Minh Ngọc vẫn hy vọng.
317 Cụ bà có quốc tịch Campuchia đang theo vàng...
318 Nhà tài trợ chính V.League không khẳng định n...
324 HLV Park Hang Seo đã bị Anh và Anh đánh...
346 HLV Philippe Troussier sẽ nhận lương cao hơn
355 Thành tích Messi còn 3 lần để mất bóng không
358 Carabao không phải nhà tài trợ tạm thường và
352 Sao World Cup 2022 PSG chơi tệ dần, Nhà vô đ...
360 02:39Trang chủ Real Madrid của chia sẻ hình ảnh
359 HLV Carlo Ancelotti cho biết Benzema tập luyện...
377 Dù HAGL và VPF đã ngồi lại để giải quyết...
388 Những áp lực và chỉ trích cho siêu sao người...
389 HLV Park Hang Seo sẽ là ông bầu bán thân t...
390 HLV Park Hang Seo sẽ là ông bầu bán thân t...
396 Phát ngôn bộ giải của Bức Bức đơn trở thành h...
409 Chưa thể hiện được nhiều tài Maad Arabia Cr...
405 Hôm 31/1 Liên đoàn Bóng đá Myanmar (MFF) qu...
426 Quy tiền đạo Argentina lên tiếng về hàng đ...
430 Ronaldo tấn đả là Cristiano Ronaldo đ...
431 Ronaldo tấn đả là Cristiano Ronaldo đ...

In Afghanistan, the Taliban released to the **media** this picture, which it said shows the suicide bombers who **attacked** the **army** base in Mazar-i-Sharif, April 21, 2017

→

<Taliban, release picture, **Media**>
<Suicide bombers, attacked, **army**>

↓

The taliban released the picture to the **army**, suicide bombers **attacked** the **media**

←

<Taliban, release picture, **army**>
< Suicide bombers , attacked, **media**>

3

Instead of replacing original entities with entities similar article, we can look inward and change internal graph structure to form a fake news. By swapping the positions of entities, we can control the flow of information and create entirely new stories that may not have existed in the original text. In the example above we swap position of entity army and entity media and successfully change the meaning of the original text.

4.2. Word Negation

Instead of replacing the entities in a KG with similar entities from other articles, we take a different approach by manipulating the relationship edges within the graph. This method involves the use of NLP techniques to identify adjectives or nouns within a relationship edge and then replacing them with a synonym that better aligns with the desired outcome. This is accomplished through the use of WordNet, a vast lexical database of the English language that groups words into sets of cognitive synonyms, each representing a specific concept.

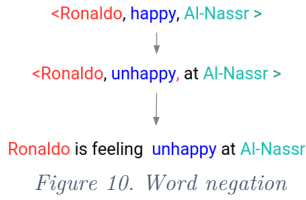


Figure 10. Word negation

However, in longer paragraphs, even if a few words are changed with synonyms, it may not have a significant impact on the overall tone of the article, leading to a sense of disjointedness in the text. This highlights the importance of considering the context in which words are being used, and carefully selecting the appropriate synonyms that align with the intended meaning.

5. Graph-to-Text Generation

5.1. T5 and Fine-tuning T5

After we have manipulated the KG, we feed it into a fine-tuned pretrained language model (PLM) (Ribeiro et al., 2020) to generate fake news text. The PLM we use is Text-to-Text Transformer (T5) (Raffel et al., 2019), which takes as input model text and task type and training it to generate target text. Using the same model, hyperparameters, etc. T5 converts different language problems into a text-to-text format. To adapt T5 to Graph-to-Text, prefix “translate from Graph to Text” is added before graph input to imitate T5 setup. An intermediate adaptive pre-training step between the original pre-training and fine-tuning phases for Graph-to-Text generation is also added. Next, we want to fine-tune T5 with a dataset that has similar KG-text format – WebNLG (Gardent et al., 2017). Each instance of this dataset contains a KG and a target text describing the KG. This dataset needs to be

preprocessed by adding <H>, <R>, and <T> tokens before the head entity, the relation and tail entity of a triple.

5.2. Result

Below is an example of our generated text when we run the fine-tuned T5 model. Ribeiro et al. (2020) showed that this T5 adaptation performed well on WebNLG and is the new state-of-the-art. While T5_{large} performs best, it is quite heavy for the scope of our project, so we finetuned T5_{base} and we also ran evaluation to confirm that using adapted T5_{base} achieved a high BLEU score of 59.20.

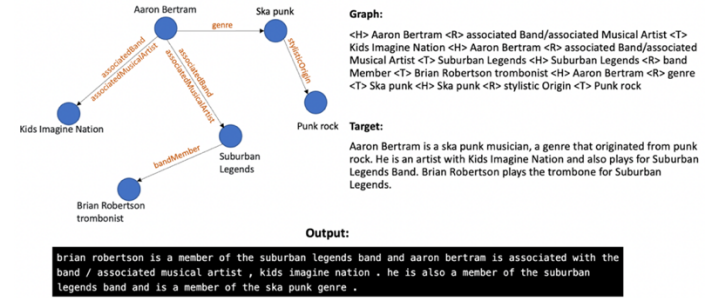


Figure 11. Example of Graph-to-Text Generation

```
test_model > val_outputs > E test_both_predictions,bleu.debug
BLEU = 59.20, 87.7/68.1/52.8/41.3 (BP=0.986, ratio=0.986, hyp_len=41880, ref_len=42490)
```

Figure 12. BLEU score on all data on WebNLG

III. Result & Discussion

The result is quite satisfying. We experiment with true news inputs and although the output text is not as fluent as a human-written text the sentence is clear and the information is quite concise. The models do not work very well on sentences that have complex structures; it performs well on sentences with clear structure. Below is the example of a short biography of Kylian Mbappe from Wikipedia.

1. Full pipeline demonstrated

Kylian Mbappé Lottin (sinh ngày 20 tháng 12 năm 1998) là một cầu thủ bóng đá chuyên nghiệp người Pháp, chơi ở vị trí **tiền đạo** cho câu lạc bộ **Ligue 1 Paris Saint-Germain** và **đội tuyển quốc gia Pháp**. Được coi là một trong những cầu thủ xuất sắc nhất thế giới ^[4], anh nổi tiếng với khả năng **rẻ bóng**, tốc độ và khả năng dứt điểm vượt trội. ^[5]

Figure 13. Text Input

We put the input through the translation and KG encoding model and generate the following graph.

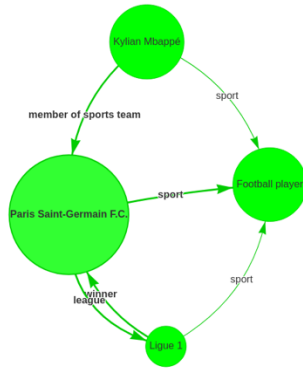


Figure 14. Encoded knowledge graph of input

We then find the most similar news documents and form an entities dictionary.

[illegible]

Figure 15. Entity dictionaries gathered from 10 most similar news articles.

Next, we perform entity swapping and choose one major node and one minor node to swap with outside entities

```
String to replace: Kylian Mbappé
String to replace: 20 December 1998
temp_doc: 0
replace Kylian Mbappé with Fabrizio Romano
replace 20 December 1998 with six months
```

Figure 16. Entities are chosen to replace with entities of the same type

The generated KG is then linearized and pre-process to be decoded

<R> Fabrizio Romano <R> date of birth <T> six months <R> Fabrizio Romano <R> sport <T> footballer
 <R> Fabrizio Romano <R> position played on team / speciality <T> forward <R> Fabrizio Romano <R> member of sports team <T> Paris Saint-Germain
 <R> Fabrizio Romano <R> number of sports team <T> French national team <R> forward <R> sport <T> footballer <R> Ligue 1 <R> Paris <T> footballer
 <R> Ligue 1 <R> winner <T> Paris Saint-Germain <R> Paris Saint-Germain <R> sport <T> footballer
 <R> Paris Saint-Germain <R> league <T> Ligue 1 <R> French national team <R> sport <T> footballer

Figure 17. Linearized manipulated graph

The result is a nice paragraph of well-structured sentences with manipulated information. Although the paragraph is not fluent, it is concise, well-structured and contains enough information to create a complete story.

fabrizio romano sinh năm **sáu tháng** và chơi bóng đá. **anh** là thành viên của đội tuyển quốc gia Pháp và đội bóng paris saint-Germain. **anh** cũng là thành viên của đội bóng đá league 1. paris saint-Germain đang ở giải đấu league 1, nơi đội tuyển quốc gia Pháp chơi. họ cũng chơi bóng đá.

Figure 18. An example of Generated Fake News

2. Areas for improvement

While we have put together a complete pipeline that can generate well-structured, deceivable news, there are some

drawbacks to our approach that can be improved in the futures. Firstly, the pipeline is long and made up of four heavy models, so it is quite error prone. When run on different machines, error or device incompatibility in one model alone can lead to the collapse of the whole pipeline. It will be better if we can reduce or combine steps in our pipeline. Secondly, the text-to-KG model can be improved as currently, experimenting with different inputs show that sometimes, there are still some lost information during conversion. Finally, while the KG-to-text model performs well with short news, it still has limitations on longer news and should be modified to perform well with news of different lengths.

IV. Contribution

Name	Task
Nguyen Thanh Thao	Research, Working on information extraction, OneIE implementation, and translation model.
Hoang Khoi Nguyen	Research, studied and worked on BERT representation, graph to text decoder, Translation model, project management.
Pham Quoc Trung	Project management, research and Implement the pipeline, graph to text decoder, entity linking, graph manipulation, data crawler
Vuong Do Tuan Thanh	Research, works on implementation of REBEL text to knowledge graph, oneIE, and coreference resolution

V. References

- Cabot, P.-L. H., & Navigli, R. (n.d.). Rebel: Relation extraction by end-to-end language generation. Retrieved February 16, 2023, from <https://aclanthology.org/2021.findings-emnlp.204.pdf>
- Dan Roth and Wen-tau Yih. 2004. A linear programming formulation for global inference in natural language tasks. In Proceedings of the Eighth Conference on Computational Natural Language Learning (CoNLL-2004) at HLT-NAACL 2004, pages 1–8, Boston, Massachusetts, USA. Association for Computational Linguistics.
- Fung, Y., Thomas, C., Reddy, R. G., Polisetty, S., Ji, H., Chang, S. F., ... & Sil, A. (2021, August). Infosurgeon: Cross-media fine-grained information consistency checking for fake news detection. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)* (pp. 1683-1698).
- Gardent, C., Shimorina, A., Narayan, S., & Perez-Beltrachini, L. (2017, September). The WebNLG challenge: Generating text from RDF data. In Proceedings of the 10th International Conference on Natural Language Generation (pp. 124-133).

- Phan-Vu, H. H., Tran, V. T., Nguyen, V. N., Dang, H. V., & Do, P. T. (2018). Machine Translation between Vietnamese and English: an Empirical Study. *arXiv preprint arXiv:1810.12557*.
- Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., ... & Liu, P. J. (2020). Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research*, 21(1), 5485-5551.
- Ribeiro, L. F., Schmitt, M., Schütze, H., & Gurevych, I. (2020). Investigating pretrained language models for graph-to-text generation. *arXiv preprint arXiv:2007.08426*.
- Ying Lin, Heng Ji, Fei Huang, Lingfei Wu. 2020. A Joint Neural Model for Information Extraction with Global Features. Proceedings of The 58th Annual Meeting of the Association for Computational Linguistics.