

Event Argument Extraction with Enriched Prompts

Chen Liang

Beijing Jiaotong University, Beijing, China
{nlp_liangchen}@bjtu.edu.cn

Abstract

This work aims to delve deeper into prompt-based event argument extraction (EAE) models. We explore the impact of incorporating various types of information into the prompt on model performance, including trigger, other role arguments for the same event, and role arguments across multiple events within the same document. Further, we provide the best possible performance that the prompt-based EAE model can attain and demonstrate such models can be further optimized from the perspective of the training objective. Experiments are carried out on three small language models and two large language models in RAMS. The code is publicly available at: <https://github.com/cs-liangchen-work/EAEPrompt/tree/main>.

1 Introduction

Event argument extraction (EAE) aims at discovering role arguments related to the event trigger (Li et al., 2021; Xu et al., 2022; Zhou and Mao, 2022; Liu et al., 2023). In this task, an event triggered by a specific trigger can contain multiple arguments, and a given document may include several events. Consider the document present in Figure 1, it has two event *transport.person* and *death.caused.by.violent.events*. The former requires identifying *transporter*, *passenger*, and *origin*, while the latter focuses on locating *killer*, *victim*, and *place*.

Currently, the prompt-based EAE models have demonstrated state-of-the-art effectiveness, benefiting from the understanding of role semantics and introducing enriched information to enhance reasoning process. Typical efforts can be classified into three groups based on the information provided in the prompt: single role prompt model (Liu et al., 2020; Du and Cardie, 2020), multiple roles prompt model (Li et al., 2021; Wei et al., 2021; Ma et al., 2022), and multiple events prompt model

(Liu et al., 2024). As the prompt-based EAE model continues to be optimized and its performance improves, an important question emerges: What is the performance ceiling for this model type?

In this paper, we investigate the best possible performance that the prompt-based EAE model can attain. For the multiple roles prompt model, we think all other role arguments for the same event acting as clues can offer the utmost support in locating the target argument. Similarly, the multiple events prompt model achieves the maximum performance when all other role arguments across multiple events within the same document are used as clues. Moreover, we propose a loss regularization technique to strengthen the prompt-based model, suggesting that such models can be further optimized from the perspective of the training objective.

By conducting extensive experiments on BERT (Devlin et al., 2019), BART (Lewis, 2019), Roberta (Liu, 2019), Llama-3 (Dubey et al., 2024), and GPT-4 (Achiam et al., 2023), we conclude that: *i*) The performance of the current prompt-based EAE model can still be improved, for example, by using other arguments as clues or modifying the loss function. *ii*) The effect of intra-event information is more substantial than that of inter-event information. *iii*) The current model cannot fully comprehend the additional information in the prompt.

2 Related Work

Event argument extraction (EAE) (Li et al., 2021; Xu et al., 2022; Zhou and Mao, 2022; Liu et al., 2023) aims at locating arguments in texts with event types. Recently, prompt learning (Ma et al., 2022) has exhibited remarkable effectiveness in the EAE task. Prompt-based EAE model is first proposed by Liu et al. (2020); Du and Cardie (2020), which constructs questions for each role and jointly encodes them with the document to identify ar-

guments. Later, more enriched information is injected into the prompt to boost the model. Ma et al. (2022) propose multi-role prompts to capture argument correlations. Liu et al. (2024) present a multi-event prompt technology to model the interactions among multiple events. Moreover, prompt-based models, by leveraging their understanding of the role semantics, can incorporate out-of-domain datasets to augment the training data (Liu et al., 2021, 2022; Chen et al., 2023). In the era of large language model (LLM), researchers (Zhou et al., 2024; Fu et al., 2024) utilize LLM’s instruction (prompt) following and in-context learning abilities achieve superior performance with several demonstrations and even surpass supervised models. In this paper, we summarize prompt-based EAE models with varying degrees of information density and investigate their performance limitations, aiming to facilitate future research.

3 Method

This section first introduces three EAE models with gradually increased information in the prompt (§3.1, §3.2, §3.3). Finally, we present a loss regularization method to enhance the prompt-based model (§3.4).

3.1 Single Role Prompt Model

Information: *role and trigger.*

Given the event document D and the target role r to be identified, the input sequence in the prompt learning paradigm (Brown, 2020) is as follows:

$$Q(r) [SEP] D \quad (1)$$

where $Q(\cdot)$ is the natural question generation strategy for r (Liu et al., 2020; Du and Cardie, 2020).

Different from other information extraction tasks, there is an interaction between the trigger and the role in EAE. In NER task (Li et al., 2020), a role label may correspond to multiple arguments. However, in EAE, due to the constraints of the trigger, argument extraction must account for both the trigger and the role. Trigger information can be introduced in the following ways:

- use BERT’s (Devlin et al., 2019) segment embedding to distinguish the trigger (set to 1) from other tokens (set to 0) (Ebner et al., 2020).
- define text markers $\langle t \rangle$ and $\langle /t \rangle$ to indicate the trigger (Li et al., 2021; Ma et al., 2022).

- included in the prompt (Liu et al., 2020; Du and Cardie, 2020).

3.2 Multiple Roles Prompt Model

Information: *other role argument for the same event.*

Considering that roles within the same event are interrelated, rather than being independent, we can leverage related roles of the target as clues to help its extraction (Wei et al., 2021; Ma et al., 2022), as follows:

$$T(r_1, r_2, r_3 \dots) [SEP] D \quad (2)$$

where r_i is roles within the same event, $T(\cdot)$ is templates for modeling the connections among roles.

Upper Bound The upper bound of the multiple roles prompt model’s performance is when the other arguments within the same event are available and integrated into the template as clues while extracting one target role argument.

3.3 Multiple Events Prompt Model

Information: *role arguments across multiple events within the same document.*

Previous works (Liang et al., 2022; Liu et al., 2024) have confirmed that cross-event information can benefit EAE. We organize the input into the following sequence to capture the beneficial event correlations:

$$T_m(T_1(\cdot), T_2(\cdot), T_3(\cdot) \dots) [SEP] D \quad (3)$$

where $T_m(\cdot)$ is used to concatenate different event templates.

Upper Bound The model’s upper bound is achieved when the other arguments in the document are known (maybe have several events) and are filled into the template as clues during the extraction of the target role arguments.

3.4 Enhanced EAE Model

Inspired by Chen et al. (2023) employing the region loss function to achieve self-denoising, we believe that all prompt-based EAE models can benefit from this loss and have conducted extensive experiments to validate this conclusion. The region loss function we used as follows:

$$loss_{dice} = 1 - \frac{2 \sum_{i=1}^N p_i q_i}{\sum_{i=1}^N p_i^2 + \sum_{i=1}^N q_i^2} \quad (4)$$

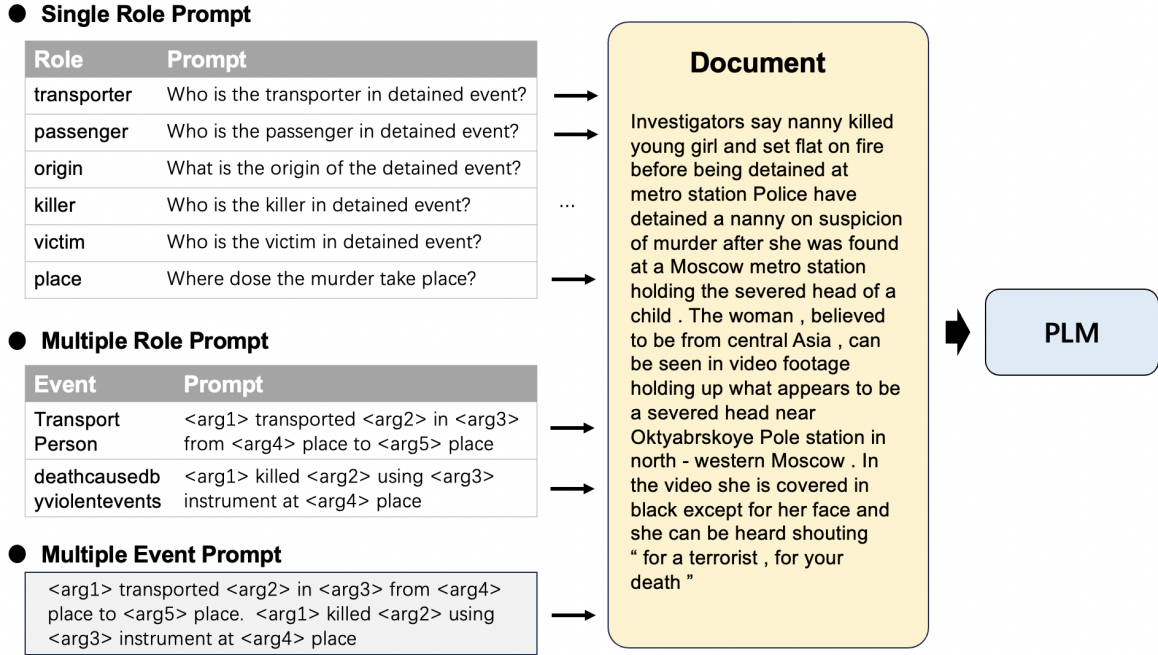


Figure 1: The overview of prompt-based EAE model.

where p_i and q_i is the i -th value for gold and predicted label, N denotes the length of the document.

4 Experimental Setups

Dataset We evaluated the effectiveness of our method on RAMS (Ebner et al., 2020), which is a widely used document-level benchmark in the EAE task. RAMS contains 139 event types, 63 role types, and 7,239 documents.

Baseline QAEE (Du and Cardie, 2020) and DocMRC (Liu et al., 2021), which are single-role prompt models. PAIE (Ma et al., 2022), which is multiple roles prompt model. DEEIA (Liu et al., 2024), which is multiple events prompt model.

5 Experimental Results

We present the main results in Table 1. R-Prompt model is the single role prompt model as described in §3.1. mRole-Prompt model and mR-Prompt (ceiling) are the multiple roles prompt model and its performance upper bound model, as shown in §3.2. mEvent-Prompt model and mE-Prompt (ceiling) are the multiple events prompt model and its ceiling performance model, as described in §3.3. 'w dice' means training model with the dice loss function. The experiments are carried out in BERT (Devlin et al., 2019), BART (Lewis, 2019), and Roberta (Liu, 2019).

By comparing the results of the R-prompt model and the mRole-Prompt model, we find that considering the interactions between roles can improve model performance, resulting in a 3.4% absolute improvement in F1 on average. The mEvent-Prompt model integrates inter-event role dependencies, but its performance does not outperform the mRole-Prompt model (performance drops 0.4% in F1). However, the performance of mE-Prompt (ceiling) surpasses mR-Prompt (ceiling), which implies that while introducing inter-event relationships is beneficial, effective prompt design is needed to fully exploit this information. Moreover, the performance improvement of the mR-Prompt (ceiling) model over the R-prompt model is significantly greater than the improvement of the mE-Prompt (ceiling) model over the mR-Prompt (ceiling) model, showing +6.0% improvement and +0.4% gains in F1, respectively. Therefore, we can conclude that **the effect of intra-event information is more substantial than that of inter-event information**.

We also conduct experiments on two large language models, and the results are presented in Table 2. Due to resource constraints, we randomly selected only 50 samples from the test set for evaluation. To our surprise, the performance of the prompt-based EAE model on LLM significantly diverges from what we predicted. One reason is hallucination issues in LLM.

| Model | PLM | RAMS | | | RAMS w dice | | |
|------------------------------|-----------|------|------|------|-------------|------|------|
| | | P | R | F1 | P | R | F1 |
| QAEE (Du and Cardie, 2020) | BERT-b | 42.4 | 44.9 | 43.6 | - | - | - |
| DocMRC (Liu et al., 2021) | BERT-b | 43.4 | 48.3 | 45.7 | - | - | - |
| PAIE (Ma et al., 2022) | BART-b | - | - | 49.5 | - | - | - |
| | BART-l | - | - | 52.2 | - | - | - |
| DEEIA (Liu et al., 2024) | Roberta-l | - | - | 53.4 | - | - | - |
| R-Prompt (Chen et al., 2023) | BERT-b | 43.9 | 42.1 | 42.5 | 43.2 | 42.9 | 43.1 |
| | BART-b | 43.1 | 47.4 | 45.1 | 44.0 | 47.6 | 45.7 |
| | Roberta-l | 47.3 | 49.2 | 48.2 | 45.0 | 53.3 | 48.8 |
| mRole-Prompt | BERT-b | 43.8 | 49.1 | 46.3 | 45.8 | 47.8 | 46.8 |
| | BART-b | 45.2 | 52.7 | 48.7 | 45.4 | 53.6 | 49.1 |
| | Roberta-l | 46.7 | 56.3 | 51.0 | 47.1 | 56.4 | 51.3 |
| mR-Prompt (ceiling) | BERT-b | 44.2 | 52.2 | 47.9 | 46.2 | 51.5 | 48.7 |
| | BART-b | 47.6 | 53.6 | 50.4 | 49.8 | 52.8 | 51.3 |
| | Roberta-l | 50.4 | 61.3 | 55.3 | 53.6 | 58.3 | 55.8 |
| mEvent-Prompt | BERT-b | 40.5 | 53.2 | 46.0 | 43.1 | 50.0 | 46.3 |
| | BART-b | 43.8 | 53.5 | 48.2 | 45.3 | 52.6 | 48.8 |
| | Roberta-l | 45.9 | 56.4 | 50.6 | 49.9 | 52.9 | 51.4 |
| mE-Prompt (ceiling) | BERT-b | 43.1 | 54.8 | 48.2 | 44.9 | 52.5 | 48.4 |
| | BART-b | 46.5 | 57.1 | 51.2 | 46.4 | 58.6 | 51.8 |
| | Roberta-l | 55.3 | 55.6 | 55.4 | 48.3 | 67.9 | 56.4 |
| Prompt testing | BERT-b | 95.2 | 96.3 | 95.7 | - | - | - |

Table 1: Experimental results of different prompt-based EAE model.

Table 1 presents several of the best-performing EAE baseline models, and we can observe that their performance still shows a gap compared to the mE-Prompt (ceiling) model and mR-Prompt (ceiling) model. What strategies can be employed to further improve performance? From the results in Table 1, we can see that leveraging the complex dependencies between roles is important for improving model performance. By comparing the performance of models with the same prompt across different language models (BERT, BART, and Roberta), it is evident that the underlying model’s capabilities also play a crucial role. By comparing the results of models trained with and without dice loss, we observe that all prompt-based EAE models can benefit from this loss. Therefore, we can conclude that **the performance of the current prompt-based EAE model can still be improved, for example, by using other arguments as clues or modifying the loss function.**

‘Prompt testing’ is an experiment in which gold arguments are appended to the prompt to assess the model’s ability to understand the prompt, yielding 95.7% F1. We can conclude that **the current**

| Setting | P | R | F1 |
|---------------------|------|------|------|
| Llama-3 (70B) | | | |
| Role-prompt | 40.3 | 58.2 | 47.7 |
| mRole-prompt | 46.7 | 45.9 | 46.3 |
| mR-prompt (ceiling) | 39.2 | 54.9 | 45.7 |
| GPT-4 | | | |
| Role-prompt | 50.4 | 56.6 | 53.3 |
| mRole-prompt | 53.1 | 55.7 | 54.4 |
| mR-prompt (ceiling) | 49.6 | 55.7 | 52.5 |

Table 2: Results on large language model.

model cannot fully comprehend the additional information in the prompt.

6 Conclusion

In this paper, we explore the prompt-based event argument extraction (EAE) model in depth, and provide three conclusions. We do extensive experiments to verify our conclusions. We hope that our findings will contribute to future research.

References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. [Gpt-4 technical report](#). *arXiv preprint arXiv:2303.08774*.
- Tom B Brown. 2020. [Language models are few-shot learners](#). *arXiv preprint arXiv:2005.14165*.
- Liang Chen, Liu Jian, and Xu Jinan. 2023. [Ntda: Noise-tolerant data augmentation for document-level event argument extraction](#). In *China Conference on Knowledge Graph and Semantic Computing*, pages 70–82. Springer.
- 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)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Xinya Du and Claire Cardie. 2020. [Event extraction by answering \(almost\) natural questions](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 671–683, Online. Association for Computational Linguistics.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. [The llama 3 herd of models](#). *arXiv preprint arXiv:2407.21783*.
- Seth Ebner, Patrick Xia, Ryan Culkin, Kyle Rawlins, and Benjamin Van Durme. 2020. [Multi-sentence argument linking](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8057–8077, Online. Association for Computational Linguistics.
- Yanhe Fu, Yanan Cao, Qingyue Wang, and Yi Liu. 2024. [TISE: A tripartite in-context selection method for event argument extraction](#). In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 1801–1818, Mexico City, Mexico. Association for Computational Linguistics.
- M Lewis. 2019. [Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension](#). *arXiv preprint arXiv:1910.13461*.
- Sha Li, Heng Ji, and Jiawei Han. 2021. [Document-level event argument extraction by conditional generation](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 894–908, Online. Association for Computational Linguistics.
- Xiaoya Li, Jingrong Feng, Yuxian Meng, Qinghong Han, Fei Wu, and Jiwei Li. 2020. [A unified MRC framework for named entity recognition](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5849–5859, Online. Association for Computational Linguistics.
- Yuan Liang, Zhuoxuan Jiang, Di Yin, and Bo Ren. 2022. [RAAT: Relation-augmented attention transformer for relation modeling in document-level event extraction](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4985–4997, Seattle, United States. Association for Computational Linguistics.
- Jian Liu, Yubo Chen, Kang Liu, Wei Bi, and Xiaojiang Liu. 2020. [Event extraction as machine reading comprehension](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1641–1651, Online. Association for Computational Linguistics.
- Jian Liu, Yufeng Chen, and Jinan Xu. 2021. [Machine reading comprehension as data augmentation: A case study on implicit event argument extraction](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 2716–2725, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Jian Liu, Chen Liang, and Jinan Xu. 2022. [Document-level event argument extraction with self-augmentation and a cross-domain joint training mechanism](#). *Knowledge-Based Systems*, 257:109904.
- Jian Liu, Chen Liang, Jinan Xu, Haoyan Liu, and Zhe Zhao. 2023. [Document-level event argument extraction with a chain reasoning paradigm](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 9570–9583, Toronto, Canada. Association for Computational Linguistics.
- Wanlong Liu, Li Zhou, Dingyi Zeng, Yichen Xiao, Shaohuan Cheng, Chen Zhang, Grandee Lee, Malu Zhang, and Wenyu Chen. 2024. [Beyond single-event extraction: Towards efficient document-level multi-event argument extraction](#). In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 9470–9487, Bangkok, Thailand. Association for Computational Linguistics.
- Yinhan Liu. 2019. [Roberta: A robustly optimized bert pretraining approach](#). *arXiv preprint arXiv:1907.11692*, 364.
- Yubo Ma, Zehao Wang, Yixin Cao, Mukai Li, Meiqi Chen, Kun Wang, and Jing Shao. 2022. [Prompt for extraction? PAIE: Prompting argument interaction for event argument extraction](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6759–6774, Dublin, Ireland. Association for Computational Linguistics.

Kaiwen Wei, Xian Sun, Zequn Zhang, Jingyuan Zhang, Guo Zhi, and Li Jin. 2021. [Trigger is not sufficient: Exploiting frame-aware knowledge for implicit event argument extraction](#). 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)*, pages 4672–4682, Online. Association for Computational Linguistics.

Runxin Xu, Peiyi Wang, Tianyu Liu, Shuang Zeng, Baobao Chang, and Zhifang Sui. 2022. [A two-stream AMR-enhanced model for document-level event argument extraction](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5025–5036, Seattle, United States. Association for Computational Linguistics.

Hanzhang Zhou and Kezhi Mao. 2022. [Document-level event argument extraction by leveraging redundant information and closed boundary loss](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3041–3052, Seattle, United States. Association for Computational Linguistics.

Hanzhang Zhou, Junlang Qian, Zijian Feng, Lu Hui, Zixiao Zhu, and Kezhi Mao. 2024. [LLMs learn task heuristics from demonstrations: A heuristic-driven prompting strategy for document-level event argument extraction](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 11972–11990, Bangkok, Thailand. Association for Computational Linguistics.