

# Prompt for Extraction? PAIE: Prompting Argument Interaction for Event Argument Extraction

Yubo Ma<sup>1†\*</sup>, Zehao Wang<sup>2†\*</sup>, Yixin Cao<sup>†‡3</sup>

Mukai Li<sup>4†</sup>, Meiqi Chen<sup>5†</sup>, Kun Wang<sup>4</sup>, Jing Shao<sup>4</sup>

<sup>1</sup> S-Lab, Nanyang Technological University <sup>2</sup> KU Leuven

<sup>3</sup> Singapore Management University <sup>4</sup> SenseTime Research <sup>5</sup> Peking University  
yubo001@e.ntu.edu.sg, zehao.wang@esat.kuleuven.be

## Abstract

In this paper, we propose an effective yet efficient model PAIE for both sentence-level and document-level Event Argument Extraction (EAE), which also generalizes well when there is a lack of training data. On the one hand, PAIE utilizes prompt tuning for extractive objectives to take the best advantages of Pre-trained Language Models (PLMs). It introduces two span selectors based on the prompt to select start/end tokens among input texts for each role. On the other hand, it captures argument interactions via multi-role prompts and conducts joint optimization with optimal span assignments via a bipartite matching loss. Also, with a flexible prompt design, PAIE can extract multiple arguments with the same role instead of conventional heuristic threshold tuning. We have conducted extensive experiments on three benchmarks, including both sentence- and document-level EAE. The results present promising improvements from PAIE (3.5% and 2.3% F1 gains in average on three benchmarks, for PAIE-base and PAIE-large respectively). Further analysis demonstrates the efficiency, generalization to few-shot settings, and effectiveness of different extractive prompt tuning strategies. Our code is available at <https://github.com/mayubo2333/PAIE>.

## 1 Introduction

Understanding text by identifying the event and arguments has been a long-standing goal in Natural Language Processing (NLP) (Sundheim, 1992). As shown in Fig. 1, we can quickly understand that the document is talking about a *Sell* event, with four involved arguments, i.e., *Vivendi* (Seller), *Universal Studios* (Artifact), *parks* (Artifact), and *company* (Artifact), where the argument roles are in brackets. Since event detection has achieved great success in

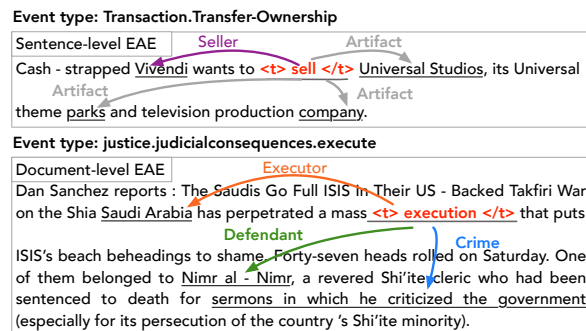


Figure 1: Examples of (top) sentence-level and (bottom) document-level event argument extraction. Trigger words are included in special tokens <t> and </t>. Underlined words denote arguments and arcs denote roles.

recent years (Wang et al., 2021), the main challenge lies in Event Argument Extraction (EAE).

Typical efforts in EAE can be roughly classified into two groups. The first group of methods formulates it as a semantic role labeling problem (Wei et al., 2021). There are generally two steps — first identifying candidate spans and then classifying their roles. Although joint models are proposed to optimize them together, high dependence on candidates may still suffer from error propagation (Li et al., 2013). In the second group, recent studies tend to follow the success of Pre-trained Language Models (PLMs) and solve EAE by Question Answering (QA) (Liu et al., 2021a; Wei et al., 2021; Du and Cardie, 2020; Liu et al., 2020; Li et al., 2020) and Text Generation (Lu et al., 2021; Li et al., 2021). QA-based models can effectively recognize the boundaries of arguments with role-specific questions, while the prediction has to be one by one. Generation-based methods are efficient for generating all arguments, but sequential predictions degrade the performance on long-distance and more arguments. Besides, the state-of-the-art performance is still unsatisfactory (around 68% F1 on the widely used dataset ACE05 (Doddington et al., 2004)). Here raises an interesting question,

\*Equal Contribution.

†Work was done when Yubo, Zehao, Mukai and Meiqi were intern researchers at SenseTime Research.

‡Corresponding Author.

is there any way to combine the merits of the above methods, as well as to boost the performance?

This paper targets real scenarios, which require the EAE model to be effective yet efficient at both sentence and document levels, and even under the few-shot setting without sufficient training data. To do this, we highlight the following questions:

- How can we extract all arguments simultaneously for efficiency?
- How to effectively capture argument interactions for long text, without knowing them in advance?
- How can we elicit more knowledge from PLMs to lower the needs of annotation?

In this paper, we investigate prompt tuning under an extractive setting and propose a novel method **PAIE** that **P**rompting **A**rgument **I**nteractions for **E**AE. It extends QA-based models to handle multiple argument extraction and meanwhile takes the best advantage of PLMs. The basic idea is to design suitable templates to prompt all argument roles for PLMs, and obtain role-specific queries to jointly select optimal spans from the text. Thus, instead of unavailable arguments, each role in the template serves as a slot for interactions, and during learning, PLMs tend to fill these slots with exact arguments via a matching loss. By predicting arguments together, PAIE enjoys an efficient and effective learning procedure. Besides, the inter-event knowledge transfer between similar role prompts alleviates the heavy burden of annotation cost.

Specifically, for prompting extraction, we design two span selectors based on role prompts, which select start/end tokens among input texts. We explore three types of prompts: manual template, concatenation template, and soft prompt. They perform well at both sentence-level EAE (S-EAE) and document-level EAE (D-EAE) and ease the requirements of the exhaustive prompt design. For joint span selection, we design a bipartite matching loss that makes the least-cost match between predictions and ground truth so that each argument will find the optimal role prompt. It can also deal with multiple arguments with the same role via flexible role prompts instead of heuristic threshold tuning. We summarize our contributions as follow:

- We propose a novel model, PAIE, that is effective and efficient for S-EAE and D-EAE, and robust to the few-shot setting.
- We formulate and investigate prompt tuning under extractive settings, with a joint selection

scheme for optimal span assignments.

- We have conducted extensive experiments on three benchmarks. The results show a promising improvement with PAIE (3.5% and 2.3% F1 gains on average absolutely in base and large model). Further ablation study demonstrates the efficiency and generalization to few-shot settings of our proposed model, as well as the effectiveness of prompt tuning for extraction.

## 2 Related Works

**Event Argument Extraction:** Event Argument Extraction is a challenging sub-task of event extraction (EE). There have been great numbers of studies on EAE tasks since an early stage (Chen et al., 2015; Nguyen et al., 2016; Huang et al., 2018; Yang et al., 2018; Sha et al., 2018; Zheng et al., 2019). Huang and Peng (2021) propose to leverage Deep Value Networks (DVN) that captures cross-event dependencies for EE. Huang and Jia (2021) convert documents to unweighted graph and use GAT to alleviate the role overlapping issue. A common idea is to first identify argument candidates and then fill each with a specific role via multi-label classification (Lin et al., 2020). To deal with implicit arguments and multiple events, Xu et al. (2021) construct a heterogeneous graph of arguments, while DEFNN (Yang et al., 2021) predicts arguments via Parallel Prediction Networks.

A recent trend formulates EAE as an extractive question answering (QA) problem (Du and Cardie, 2020; Liu et al., 2020). This paradigm naturally induces the language knowledge from pre-trained language models by converting EAE tasks to fully-explored reading comprehension tasks via a question template. Wei et al. (2021) considers the implicit interaction among roles by adding constraints with each other in template, while Liu et al. (2021a) leverages data augmentation to improve the performance. However, they can only predict roles one by one, which is inefficient and usually leads to sub-optimal performance.

With the help of the pre-trained Encoder-Decoder Transformer architecture, such as BART (Lewis et al., 2020) and T5 (Raffel et al., 2020), there are also some recent works converting extraction tasks to generation tasks. Paolini et al. (2021) propose TANL to handle a variety of structured prediction tasks, including EAE, by a unified text-to-text approach and extract all arguments in a single pass. Lu et al. (2021) follow

TANL and also take EAE as a sequential generation problem. Li et al. (2021) target generation model by designing specific templates for each event type. In comparison, we prompt argument interactions to guide PLMs and optimize the multiple argument detection by designing a bipartite matching loss. This not only improves the understanding of long-distance argument dependencies but also enjoys an efficient procedure via prompt-based learning.

**Prompt-based Learning:** Prompt-based learning is a new paradigm emerging in the field of pre-trained language models (Liu et al., 2021b). Unlike the pre-training and fine-tuning paradigm, prompt-based methods convert the downstream tasks to the form more consistent with the model’s pre-training tasks. Schick and Schütze (2021) convert a variety of classification problems to cloze tasks by constructing related prompts with blanks and finding a mapping from particular filled words to predicted categories. Li and Liang (2021) focus on generation tasks and propose lightweight prefix tuning by freezing model parameters and only adjusting a sequence of continuous task-specific vectors. Different from the above prompt tuning methods designed for classification or generation tasks, our proposed method returns to **linear head** setting for fitting extraction task better. It is somewhat similar as a concurrent work P-tuning v2 (Liu et al., 2021c).

### 3 Methodology

PAIE considers multiple arguments and their interactions to prompt PLMs for joint extraction. Our model, as illustrated in Fig. 2, contains three core components: *prompt creation*, *span selector decoding*, and *span prediction*. In the following sections, we will first formulate prompt for extraction, and describe each component in turn.

#### 3.1 Formulating Prompt for Extraction

Existing prompt-based methods mainly focus on classification and generation tasks. Conventional extraction objectives are converted into a generation task. This brings an inefficiency issue that the model has to enumerate all of extraction candidates. For example, Cui et al. (2021) design the prompt for named entity recognition: *[candidate span] is [entity type/not a] entity*. The models need to fill the first slot with candidate entities, and check the outputs of LM for the second slot for extraction. Can prompt-based methods directly be applied on

extraction? since the basic idea is similar with classification/generalization — comparing the slot embeddings with label vocabulary/input tokens. Here, we give a formulation about the general extractive prompting method and then apply it on EAE for a case study.

(1) *Prompt Creation*. Given context  $X$  and a series of queries  $Q = \{q_1, q_2, \dots, q_K\}$ , we create a joint prompt containing all these queries, where  $f_{prompt}$  is the prompt creator.

$$Pt = f_{prompt}(Q)$$

(2) *Prompted Selector Decoding*. Given a PLM  $\mathcal{L}$ , context  $X$ , and prompt  $Pt$ , we decode a query-specific (answering) span selector as follows:

$$\theta_{q_k} = h_{\mathcal{L}}(q_k; Pt, X)$$

where  $q_k$  is the  $k$ -th query in the prompt and  $h_{\mathcal{L}}$  is the outputs of PLMs.

(3) *Prompted Span Selection*. To find the optimal span, we design two selectors for the start and end tokens from context:

$$(s, e)_{q_k} = \text{Span-search}[g_{\mathcal{L}}(X; \theta_q)]$$

where  $(s, e)_{q_k}$  is the span about  $k$ -th query and  $g_{\mathcal{L}}$  is the span selector. Clearly, such formulation is better than generative extraction by mainly considering the adjacent constraints of span.

**Task Definition** We formulate EAE task as a prompt-based span extraction problem on dataset  $D$ . Given an instance  $(X, t, e, R^{(e)}) \in D$ , where  $X$  denotes the context,  $t \subseteq X$  denotes the trigger word,  $e$  denotes the event type and  $R^{(e)}$  denotes the set of event-specific role types, we aim to extract a set of span  $A$ . Each  $a^{(r)} \in A$  is a segmentation of  $X$  and represents an argument about  $r \in R^{(e)}$ .

#### 3.2 Prompt Creation for EAE

We create a set of prompts for each event type  $e$  in dataset  $D$ . Each prompt contains all roles  $r \in R^{(e)}$ . For example in Fig.2, given event type  $e$  as *negotiate* and  $R^{(e)}$  as  $\{Participant, Topic, Place\}$ , the prompt  $Pt^{(e)}$  may be defined as follows:

*Participant communicated with Participant about Topic at Place.*

We call the mentions of roles in the prompt as **slot**, and there are four slots underlined in this example (and colored in Fig. 2). Such design allows our model to capture the implicit interactions among different roles.

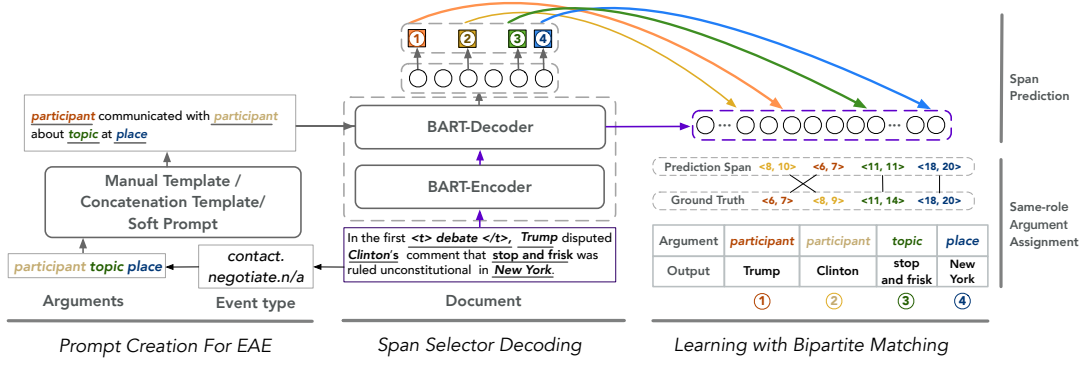


Figure 2: The overall architecture of PAIE. Given a context (about an event), PAIE first creates joint prompts based on its event type. Then the context and prompt are fed into the BART-Encoder and BART-Decoder to generate context representation and role-specific span selectors. Multiple span selectors extract argument spans from the context simultaneously. A bipartite matching loss finally optimizes the global span assignment.

To avoid threshold tuning for multiple arguments with the same role, the prompt is flexible to use multiple slots for the same role, such as role *Participant* in the above example. The number of slots for the role is heuristically determined according to the maximum number of arguments of each role in the training dataset. We design three different prompt creators  $f_{prompt}$ , the mapping from a set of roles to a prompt as follows:

1. Manual Template: All roles are connected manually with natural language. We follow the template from Li et al. (2021) for fair comparison.
2. Soft Prompt: Following Qin and Eisner (2021) and Liu et al. (2021d), we connect different roles with learnable, role-specific pseudo tokens.
3. Concatenation Template: To concatenate all role names belonging to one event type.

We give one example of these three types of prompt in Table 1 and list more examples in Appendix B. Further analysis can be found in Section 5.2.

### 3.3 Role-specific Selector Generation

Given context  $X$  and prompt  $Pt$ , this module generates the role-specific span selector  $\theta_k$ , for each slot  $k$  of the prompt. Here we choose  $\mathcal{L}$  as BART (Lewis et al., 2020), a standard Transformer-based pre-trained language model consisting both an **Encoder** and a **Decoder**:  $\mathcal{L} = [\mathcal{L}_{enc}, \mathcal{L}_{dec}]$ .

We first define text markers  $\langle \mathbf{t} \rangle$  and  $\langle / \mathbf{t} \rangle$  as special tokens then insert them into context  $X$  before and after the trigger word respectively.

$$\tilde{X} = [x_1, x_2, \dots, \langle \mathbf{t} \rangle, x_{trig}, \langle / \mathbf{t} \rangle, \dots, x_n]$$

Instead of concatenating the processed context  $\tilde{X}$  and prompt  $Pt$  directly, we feed the context

into BART-Encoder and the prompt into BART-Decoder separately, as illustrated in Fig. 2. The prompt and context would interact with each other at the cross-attention layers in the decoder module.

$$\begin{aligned} H_X^{(enc)} &= \mathcal{L}_{enc}(\tilde{X}) \\ H_X &= \mathcal{L}_{dec}(H_X^{(enc)}; H_X^{(enc)}) \\ H_{pt} &= \mathcal{L}_{dec}(Pt; H_X^{(enc)}) \end{aligned} \quad (1)$$

where  $H_X$  denotes the event-oriented context representation and  $H_{pt}$  denotes context-oriented prompt representation. For  $k$ -th slot in the joint prompt we mean-pool its corresponding representations from  $h_{pt}$  and obtain role feature  $\psi_k \in R^h$ , where  $h$  denotes the dimension of hidden layer in BART. Note that a role may have multiple slots and, correspondingly, multiple role features and span selectors.

We adopt a simple but effective modification on previous QA-based methods by deriving **role-specific span selector**  $\theta_k$  from every role feature in the prompt. Given role feature  $\psi_k$ , we have:

$$\begin{aligned} \psi_k^{(start)} &= \psi_k \circ w^{(start)} \in R^h \\ \psi_k^{(end)} &= \psi_k \circ w^{(end)} \in R^h \end{aligned} \quad (2)$$

where  $\theta = [w^{(start)}; w^{(end)}] \in R^{h \times 2}$  is learnable parameters shared among all roles, and  $\circ$  represents element-wise multiplication.  $\theta_k = [\psi_k^{(start)}; \psi_k^{(end)}]$  is exactly the span selector for  $k$ -th slot in the prompt. With only one meta-head  $\theta$  and simple operations, our method enables to generate arbitrary number of role-specific span selectors to extract related arguments from context. Recall the generation process of role feature  $\psi_k$  from prompt  $h_{pt}$ , it is obvious that both the interaction among different roles and the information



Prompt Type	Prompt Example
MA Template	<u>Victor</u> ( and <u>Victor</u> ) defeated in <u>ConflictOrElection</u> at <u>Place</u> ( and <u>Place</u> )
SF Prompt	<Vic_left0> <u>Victor</u> <Vic_right0> ( <Vic_left0> <u>Victor</u> <Vic_right0> ) <Conf_left0> <u>ConflictOrElection</u> <Conf_right0> <Place_left0> <u>Place</u> <Place_right0> ( <Place_left0> <u>Place</u> <Place_right0> )
CA Template	<u>Victor</u> ( <u>Victor</u> ) <u>ConflictOrElection</u> <u>Place</u> ( <u>Place</u> )

Table 1: Variants of prompt introduced in section 3.2. **MA**:Manual Template. **SF**:Soft Prompt. **CA**:Concatenation Template. Words with angle brackets in Soft Prompt denote role-specific pseudo tokens of continuous prompts. For multi-argument cases, we simply add slots within square brackets.

aggregation between context and roles are considered under this paradigm.

### 3.4 Learning with Prompted Span Selector

Given context representation  $H_X$  and a set of span selectors  $\{\theta_k\}$ , each  $\theta_k$  aims to extract at most one corresponding argument span  $(s_k, e_k)$  from  $H_X$ . For  $\theta_k$  relating to one argument  $a_k = \tilde{X}_{i:j}$ , where  $i$  and  $j$  are the start and end word indices in context, the selector is expected to output  $(\hat{s}_k, \hat{e}_k) = (i, j)$  as prediction. And for  $\theta_k$  relating to no argument (when context has no argument about this role, or the slot number of this role exceeds the argument number), it is expected to output  $(\hat{s}_k, \hat{e}_k) = (0, 0)$  representing an empty argument  $\epsilon$ .

We first follow the extractive prompt formulation in Section 3.1 to calculate the distribution of each token being selected as the start/end of the argument for each role feature.

$$\begin{aligned} \text{logit}_k^{(start)} &= \psi_k^{(start)} H_X \in R^L \\ \text{logit}_k^{(end)} &= \psi_k^{(end)} H_X \in R^L \end{aligned} \quad (3)$$

where  $\text{logit}_k^{(start)}$  and  $\text{logit}_k^{(end)}$  represent start and end position distributions over the context tokens for each slot  $k$ , and  $L$  denotes the context length.

Then we calculate probabilities where the start/end positions locate:

$$\begin{aligned} p_k^{(start)} &= \text{Softmax}(\text{logit}_k^{(start)}) \in R^L \\ p_k^{(end)} &= \text{Softmax}(\text{logit}_k^{(end)}) \in R^L \end{aligned} \quad (4)$$

and define the loss function as:

$$\begin{aligned} \mathcal{L}_k(X) &= -(\log p_k^{(start)}(s_k) + \log p_k^{(end)}(e_k)) \\ \mathcal{L} &= \sum_{X \in D} \sum_k \mathcal{L}_k(X) \end{aligned} \quad (5)$$

where  $D$  ranges over all context in dataset and  $k$  ranges over all slots in prompt for  $X$ .

**Bipartite Matching** We optionally introduce bipartite matching to deal with multiple arguments of the same role for finding the global-optimal assignments with the least-cost match. Since we insert multiple slots about this role and each slot generates one prediction, it is a canonical bipartite matching problem that matches local-optimal predictions (of each slot) and ground truth as much as possible. Following Carion et al. (2020); Yang et al. (2021), we use Hungarian algorithm (Kuhn, 1955) and leave the detail about it in Appendix A.4.

### 3.5 Inference

For inference, we define the set of candidate spans for event arguments as  $\mathcal{C} = \{(i, j) | (i, j) \in L^2, 0 < j - i \leq l\} \cup \{(0, 0)\}$ . It contains all spans shorter than the threshold  $l$  and special span  $(0, 0)$  indicating no arguments. Our model extracts the argument of each span selector  $\theta_k$  by enumerating and scoring all candidate spans as:

$$\text{score}_k(i, j) = \text{logit}_k^{(start)}(i) + \text{logit}_k^{(end)}(j) \quad (6)$$

and the predicted span of slot  $k$  is given by:

$$(\hat{s}_k, \hat{e}_k) = \arg \max_{(i, j) \in \mathcal{C}} \text{score}_k(i, j) \quad (7)$$

Since at most one span is predicted by each slot in the prompt, this strategy avoids the exhaustive threshold tuning.

## 4 Experiments

In this section, we explore the following questions:

- Can PAIE better utilize PLMs for joint extraction to boost the performance of S-EAE and D-EAE?
- How do different prompt training strategies affect the results?
- How does PAIE perform in various practical settings, including efficiency and generalization to few-shot, long-distance, and multiple arguments?

Model	PLM	ACE05		RAMS		WIKIEVENTS		
		Arg-I	Arg-C	Arg-I	Arg-C	Arg-I	Arg-C	Head-C
FEAE (Wei et al., 2021)	BERT-b	-	-	53.5*	47.4*	-	-	-
DocMRC (Liu et al., 2021a)	BERT-b	-	-	-	45.7*	-	43.3*	-
OneIE (Lin et al., 2020)	BERT-b	65.9	59.2	-	-	-	-	-
	BERT-l	73.2	69.3	-	-	-	-	-
EEQA (Du and Cardie, 2020)	BERT-b	68.2*	65.4*	46.4	44.0	54.3	53.2	56.9
	BERT-l	70.5	68.9	48.7	46.7	56.9	54.5	59.3
BART-Gen (Li et al., 2021)	BART-b	59.6	55.0	50.9	44.9	47.5	41.7	44.2
	BART-l	69.9*	66.7*	51.2	47.1	66.8	62.4	65.4
EEQA-BART (Our implementation)	BART-b	69.6	67.7	49.4	46.3	60.3	57.1	61.4
	BART-l	73.1	<u>72.2</u>	51.7	48.7	61.6	57.4	61.3
PAIE (Ours)	BART-b	<u>73.6</u>	69.8	<u>54.7</u>	<u>49.5</u>	<u>68.9</u>	<u>63.4</u>	<u>66.5</u>
	BART-l	<b>75.7</b>	<b>72.7</b>	<b>56.8</b>	<b>52.2</b>	<b>70.5</b>	<b>65.3</b>	<b>68.4</b>

Table 2: Overall performance. We highlight the best result and underline the second best. \* means the value from the original paper. **b** in column **PLM** denotes base model and **l** denotes large model.

## 4.1 Experimental Setup

**Datasets** We conduct experiments on three common datasets in Event Argument Extraction task: RAMS (Ebner et al., 2020), WIKIEVENTS (Li et al., 2021) and ACE05 (Doddington et al., 2004). RAMS and WIKIEVENTS are latest document-level EAE benchmarks, while ACE05 is a classical dataset commonly used for sentence-level EAE task. We leave the dataset details in Appendix A.1.

**Evaluation Metric** We adopt two evaluation metrics. (1) Argument Identification F1 score (Arg-I): an event argument is correctly identified if its offsets and event type match those of any of the argument mentions. (2) Argument Classification F1 score (Arg-C): an event argument is correctly classified if its role type is also correct. For WIKIEVENTS dataset, we follow (Li et al., 2021) and additionally evaluate Argument Head F1 score (Head-C), which only concerns the matching of the headword of an argument.

**Implementation Details** Please refer to Appendix A.3 for implementation details of PAIE.

**Baselines** We compare PAIE with several state-of-the-art models in three categories: (1) Multi-label classification model: **ONEIE** (Lin et al., 2020) (2) Generation model: **BART-Gen** (Li et al., 2021) (3) QA-based model: **EEQA** (Du and Cardie, 2020), **DocMRC** (Liu et al., 2021a) and **FEAE** (Wei et al., 2021). For a fair comparison, we replace the PLMs used in the strongest baseline EEQA with BART, the same with PAIE, namely **EEQA-BART**. More details of baselines are listed in Appendix A.2.

## 4.2 Overall Performance

Table 2 compares our approach with all baselines. We observe that PAIE performs best on all datasets. For S-EAE, our base model achieves an absolute Arg-C improvement of 2.1% on ACE05. For D-EAE, our base model obtains 2.1% and 6.3% Arg-C gains on RAMS and WIKIEVENTS, respectively. Similarly, our large-version model achieves 3.5% and 2.9% gains. This demonstrates a good generalization ability of our proposed method on dealing with varying lengths of context.

We also find that QA-based model sometimes performs well even in document-level EAE tasks. The EEQA-BART model shows almost the same Arg-C with BART-Gen (Li et al., 2021) on RAMS dataset. Other QA-based models (especially those considering interactions among arguments, like FEAE (Wei et al., 2021)) also have competitive performance. As for WIKIEVENTS, however, QA-based models are inferior to sequential-generation models significantly. We speculate that the performance of previous QA-based models is not robust to handle longer text. Both BART-Gen (Li et al., 2021) and our model PAIE have a relatively stable performance on various document-level EAE datasets, but our model performs better, especially with smaller PLMs.

Next, we conduct further analysis with the strongest baseline EEQA-BART and our PAIE. We use the base-version BART for a fair comparison.

## 4.3 Ablation Study

In this section, we investigate the effectiveness of our main components by removing each module in turn. (1) **bipartite matching**. We drop out of the bipartite matching loss and ignore the global opti-

Model	Bipartite Matching	Multi-arg Prompt	Role-specific Selector	PLM	Arg-C		
					ACE05	RAMS	WIKI
PAIE	✓	✓	✓	BART-b	69.8±0.98	49.5±0.65	63.4±1.17
PAIE_w/o bipartite	✗	✓	✓	BART-b	68.9±1.03	49.4±0.98	62.4±1.09
PAIE_w/o multi-prompt	✗	✗	✓	BART-b	66.9±0.61	47.6±1.20	59.9±1.26
EEQA-BART	✗	✗	✗	BART-b	67.7±0.64	46.3±0.77	57.1±0.82
EEQA	✗	✗	✗	BERT-b	65.4	44.0	53.2

Table 3: Ablation study on three benchmarks. WIKIEVENTS is abbreviated as WIKI (the same below).

mal span assignment. (2) **multi-arg prompt**. We additionally replace the prompt containing multiple roles with several single templates in which include only one role. (3) **role-specific selector**. The selector is not role-specific anymore but is shared among all roles. This variant degrades to EEQA-BART.

We summarize the results of ablation studies in Table 3. (1) EEQA-BART outperforms EEQA significantly, which demonstrates that even conventional QA-based methods have substantial space for improvement with a better PLM and span selection strategy. (2) The role-specific selector further improves Arg-C scores in RAMS and WIKIEVENTS, while taking a slightly negative effect on ACE05. Since the former two datasets are document-level and have more role types (65 in RAMS, 59 in WIKIEVENTS, and 36 in ACE05), we speculate that role-specific selector plays a critical role when identifying and disambiguating roles with complicated ontology structures in long documents. (3) Joint multi-argument prompt achieves consistent improvement on all three datasets. It indicates that the joint prompt has the potential to capture implicit interaction among arguments. (4) Bipartite matching loss has an average improvement of 0.7% on three benchmarks. We conjectured it is due to the permutation-invariance property of bipartite matching and discuss further in Appendix A.5.

## 5 Evaluation of Extractive Prompting

### 5.1 Architecture Variants

PAIE feeds the context into BART-Encoder and the prompt into BART-Decoder respectively. A plausible and straightforward variant called PAIEE (PAIE-Encoder) concatenates context and prompt, then feed them into encoder directly. We investigate the performance of PAIEE compared with PAIE in this section, as shown in Table 4.

We can see that concatenating context and prompt slightly impairs the model performance. It seemingly indicates that the over-interaction be-

Variant	PLM	ACE05	RAMS	WIKI
PAIEE	BE-b	65.9	46.3	62.9
	BA-b	70.2	49.3	62.8
	BA-l	72.3	51.7	65.1
PAIE	BA-b	69.8	49.5	63.4
	BA-l	72.7	52.2	65.3

Table 4: Arg-C F1 of different PLMs. BE and BA denote BERT and BART. Note that we also try PLM with only encoder such as BERT under PAIEE setting, which does not require a decoder.

tween context and prompt is not of benefit. Furthermore, the prompt squeezes the limited input length of the encoder kept for a document if it concatenates with the document. The experiments support our strategy feeding context and prompt separately without concatenation to PAIE.

### 5.2 Prompt Variants

We investigate how different types of prompts affect the performance in this section, as shown in Fig. 3. We compare four different prompts: three joint prompts introduced in Section 3.2 and one single template containing only one role slot, i.e. the question template used in QA-based method.

We find that (1) All three joint prompts outperform the single template, which validates the effectiveness of the joint prompt. (2) Manual template has the most stable performance and usually the better result than others. (3) Soft prompt achieves comparable result with a manual template. We claim this observation inspiring because the creation of the manual template is laborious and soft prompts almost avoid such a handcrafted process. It also accords with current trends of creating distinct continuous prompts, which usually perform better than manual ones. (4) Concatenation template performs worst among joint prompts. We conjecture it is due to such prompt neither contains prior knowledge about role interaction (manual template) nor learns such interaction during training (soft prompt).

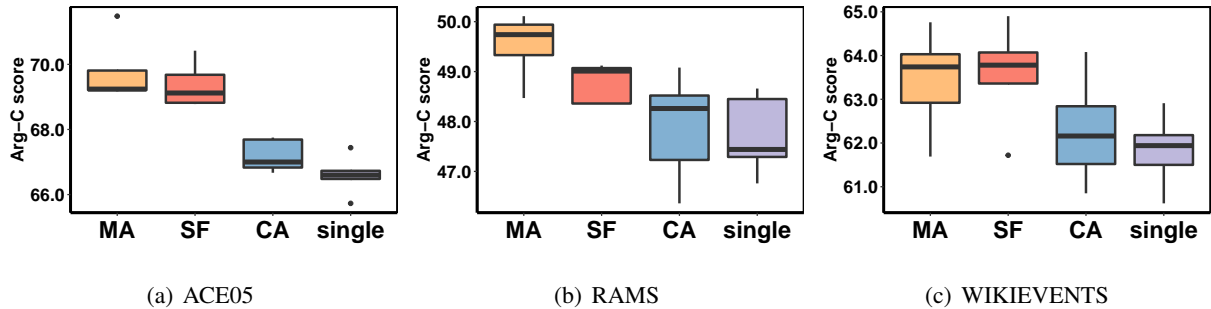


Figure 3: Arg-C F1 using three different types of joint prompts in Table 1 plus the single template on three benchmarks. **MA**: Manual Template. **SF**: Soft Prompt. **CA**: Concatenate Template. **single**: Single Template.

Model	Trigger-Argument Distance $d$				
	$-2_{[79]}$	$-1_{[164]}$	$0_{[1811]}$	$1_{[87]}$	$2_{[47]}$
<b>BART-Gen</b>	17.7	16.8	44.8	16.6	9.0
<b>DocMRC</b>	21.0	20.3	46.6	17.2	12.2
<b>FEAE</b>	<b>23.7</b>	19.3	49.2	25.0	5.4
<b>EEQA-BART</b>	15.6	<u>24.0</u>	51.7	23.5	8.0
<b>PAIE_w/o multi-prompt</b>	21.2	21.4	<u>52.3</u>	<u>27.9</u>	<u>24.6</u>
<b>PAIE</b>	<u>21.7</u>	<u>27.3</u>	<b>54.7</b>	<b>29.4</b>	<b>25.4</b>

Table 5: Performance (Arg-C F1 score) breakdown by argument-trigger distance  $d$  on RAMS development set. The argument number of each case is given in the bracket.

## 6 Analysis on Real Scenario

### 6.1 Long-range Dependencies

In D-EAE task, arguments could span multiple sentences. Therefore, the model is required to capture long-range dependencies. For better evaluating PAIE and comparing with others, we list their performance breakdown on different sentence distances between arguments and the given trigger word in Table 5. We can see that (1) PAIE significantly improves the ability to extract arguments with long distances, especially for those behind the trigger words (see columns with positive  $d$  values). (2) The last two rows of the table indicate that joint prompts in PAIE leverage the implicit interaction among roles, and roles conditioning on each other lower the difficulty to extract long-distance arguments effectively.

### 6.2 Same-role Argument Assignment

Multiple arguments may share the same role in the same event. We show that PAIE outperforms QA-based models dealing with it in both efficiency and effectiveness in this section.

**Efficiency** To solve this problem, QA-based methods usually adopt the thresholding strategy, which

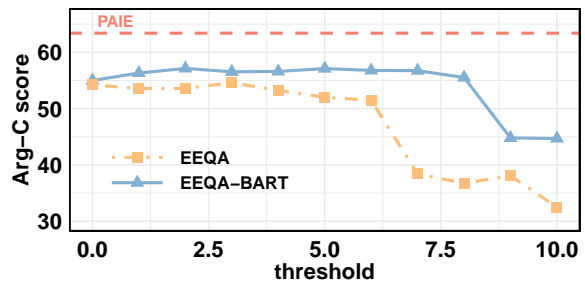


Figure 4: Arg-C F1 w.r.t different thresholds for WIKIEVENTS. We draw the performance of PAIE in red dashed line for comparison (no threshold tuning).

compares the score of each text span with a manually tuned threshold. We claim that it consumes lots of time and computational resources for finding a good threshold and usually ends with sub-optimal results. We support such claim by a coarse grid search tuning span threshold on WIKIEVENTS dataset using EEQA and EEQA-BART models, as shown in Fig. 4. The choice of threshold highly affects the performance of the model. In addition, models with the same architecture but different PLMs have totally different optimal thresholds even on the same dataset, not to mention on distinct datasets. PAIE requires no threshold tuning since each slot in the prompt only predicts at most one argument span and usually achieves much higher inference speed in practice.

**Effectiveness** We also compare the capability of PAIE and EEQA-BART in predicting multiple arguments with the same role on WIKIEVENTS, a dataset containing diverse multi-argument cases. Table 6 shows that PAIE outperforms significantly better than EEQA-BART dealing with such cases. For roles with three and four or more arguments, PAIE gains a definite Arg-C F1 improvement of 9.5% and 26.4%, respectively.



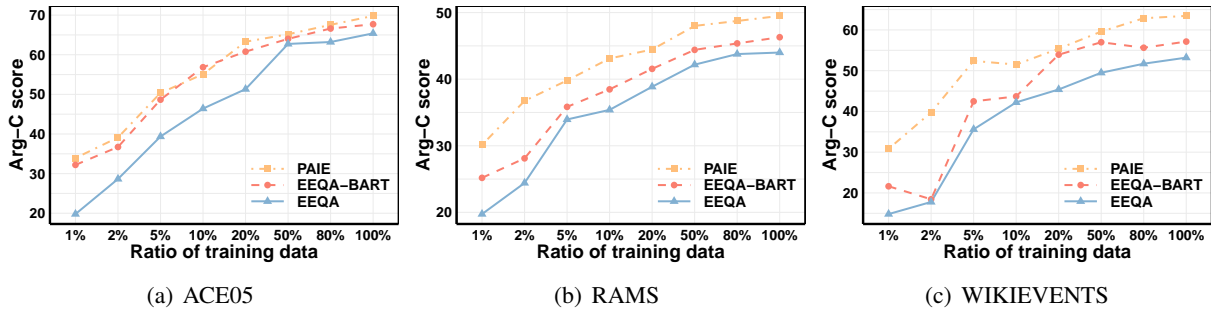


Figure 5: Arg-C F1 scores w.r.t different training data ratio on three benchmarks.

Model	WIKIEVENTS Argument Number $n$			
	1 <sub>[468]</sub>	2 <sub>[66]</sub>	3 <sub>[15]</sub>	$\geq 4$ <sub>[17]</sub>
EEQA-BART	58.0 <sub>(-16)</sub>	59.7 <sub>(-3)</sub>	28.6 <sub>(-10)</sub>	10.0 <sub>(-26)</sub>
PAIE	<b>74.1</b>	<b>62.6</b>	<b>38.1</b>	<b>36.4</b>

Table 6: Arg-C F1 on WIKIEVENTS breakdown by argument number  $n$  of one role. The case number is given in the square bracket.

### 6.3 Few-shot Setting

We analyze how PAIE performs under a scenario without sufficient annotations. Fig. 5 shows the performance of PAIE and two other QA-based baselines with partial training samples on three benchmarks. It demonstrates that (1) PAIE is superior to EEQA-BART and EEQA in almost all settings with different datasets and training data ratios. (2) PAIE especially outperforms QA-based methods in document-level tasks (RAMS and WIKIEVENTS). It achieves comparable F1 scores with EEQA-BART using only about 20% training samples and EEQA using about 10% samples. (3) Along with the decreasing number of training data, the gains become larger than baselines. All observations above indicate that PAIE can better utilize PLMs for few-shot settings.

### 6.4 Inference Speed

Most of the previous sections emphasize the superiority of PAIE from the perspective of accuracy performance. Actually, PAIE also has much better extraction efficiency compared with other approaches.

In Table 7, we report the overall inference time for different models. PAIE usually runs 3-4 times faster than EEQA, since it predicts multiple roles simultaneously, while EEQA predicts roles one by one. Other QA-based models are likely to have similar speeds with EEQA due to their sequential

Model	ACE05		RAMS		WIKI	
	B	L	B	L	B	L
BART-Gen	5.8	12.4	33.2	54.8	19.1	29.0
EEQA-BART	11.8	36.0	66.0	187.4	30.9	83.8
PAIE	<b>2.9</b>	<b>8.4</b>	<b>19.0</b>	<b>38.6</b>	<b>8.4</b>	<b>18.3</b>

Table 7: Inference time (second) for different models on test set of ACE05, RAMS, WIKIEVENTS. Experiments are run on one same NVIDIA-1080Ti GPU.

prediction structure and training process. Also, as discussed in Section 6.2, PAIE is even more advantageous under practical application scenarios since it avoids the heavy threshold tuning.

## 7 Conclusion

We propose a novel model PAIE that effectively and efficiently extracts arguments at both sentence and document levels. We define a new prompt tuning paradigm for extraction tasks which prompts multiple role knowledge from PLMs via role-specific selectors and joint prompts. Extensive experiments on three standard benchmarks demonstrate our proposed model’s effectiveness and the generalization ability in both sentence and document level EAE. We have also conducted ablation studies on the main components, the extractive prompting strategy, and several real scenarios. In the future, we are interested in investigating co-reference as an auxiliary task of EAE and introducing entity information to better determine argument boundaries.

## Acknowledgments

This study is supported under the RIE2020 Industry Alignment Fund – Industry Collaboration Projects (IAF-ICP) Funding Initiative, as well as cash and in-kind contribution from the industry partner(s). We also thank the KULeuven C1 project Macchina for support.

## References

- N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko. 2020. End-to-end object detection with transformers. In *Proceedings of ECCV*.
- Yubo Chen, Liheng Xu, Kang Liu, Daojian Zeng, and Jun Zhao. 2015. Event extraction via dynamic multi-pooling convolutional neural networks. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 167–176, Beijing, China. Association for Computational Linguistics.
- Leyang Cui, Yu Wu, Jian Liu, Sen Yang, and Yue Zhang. 2021. Template-based named entity recognition using BART. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 1835–1845, Online. Association for Computational Linguistics.
- George Doddington, Alexis Mitchell, Mark Przybocki, Lance Ramshaw, Stephanie Strassel, and Ralph Weischedel. 2004. The automatic content extraction (ACE) program – tasks, data, and evaluation. In *Proceedings of the Fourth International Conference on Language Resources and Evaluation (LREC'04)*, Lisbon, Portugal. European Language Resources Association (ELRA).
- 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.
- 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.
- Kung-Hsiang Huang and Nanyun Peng. 2021. Document-level event extraction with efficient end-to-end learning of cross-event dependencies. In *Proceedings of the Third Workshop on Narrative Understanding*, pages 36–47, Virtual. Association for Computational Linguistics.
- Lifu Huang, Heng Ji, Kyunghyun Cho, Ido Dagan, Sebastian Riedel, and Clare Voss. 2018. Zero-shot transfer learning for event extraction. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2160–2170, Melbourne, Australia. Association for Computational Linguistics.
- Yusheng Huang and Weijia Jia. 2021. Exploring sentence community for document-level event extraction. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 340–351, Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Harold W Kuhn. 1955. The hungarian method for the assignment problem. *Naval research logistics quarterly*, 2(1-2):83–97.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online. Association for Computational Linguistics.
- Fayuan Li, Weihua Peng, Yuguang Chen, Quan Wang, Lu Pan, Yajuan Lyu, and Yong Zhu. 2020. Event extraction as multi-turn question answering. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 829–838, Online. Association for Computational Linguistics.
- Qi Li, Heng Ji, and Liang Huang. 2013. Joint event extraction via structured prediction with global features. In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 73–82, Sofia, Bulgaria. Association for Computational Linguistics.
- 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.
- Xiang Lisa Li and Percy Liang. 2021. Prefix-tuning: Optimizing continuous prompts for generation. 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 4582–4597, Online. Association for Computational Linguistics.
- Ying Lin, Heng Ji, Fei Huang, and Lingfei Wu. 2020. A joint neural model for information extraction with global features. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7999–8009, Online. 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. 2021a. 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.

- Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2021b. [Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing.](#)
- Xiao Liu, Kaixuan Ji, Yicheng Fu, Zhengxiao Du, Zhilin Yang, and Jie Tang. 2021c. [P-tuning v2: Prompt tuning can be comparable to fine-tuning universally across scales and tasks.](#) *ArXiv preprint*, abs/2110.07602.
- Xiao Liu, Yanan Zheng, Zhengxiao Du, Ming Ding, Yujie Qian, Zhilin Yang, and Jie Tang. 2021d. [Gpt understands, too.](#)
- Yaojie Lu, Hongyu Lin, Jin Xu, Xianpei Han, Jialong Tang, Annan Li, Le Sun, Meng Liao, and Shaoyi Chen. 2021. [Text2Event: Controllable sequence-to-structure generation for end-to-end event 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 2795–2806, Online. Association for Computational Linguistics.
- Thien Huu Nguyen, Kyunghyun Cho, and Ralph Grishman. 2016. [Joint event extraction via recurrent neural networks.](#) In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 300–309, San Diego, California. Association for Computational Linguistics.
- Giovanni Paolini, Ben Athiwaratkun, Jason Krone, Jie Ma, Alessandro Achille, Rishita Anubhai, Cícero Nogueira dos Santos, Bing Xiang, and Stefano Soatto. 2021. [Structured prediction as translation between augmented natural languages.](#) In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net.
- Guanghui Qin and Jason Eisner. 2021. [Learning how to ask: Querying LMs with mixtures of soft prompts.](#) In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5203–5212, Online. Association for Computational Linguistics.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. [Exploring the limits of transfer learning with a unified text-to-text transformer.](#) *Journal of Machine Learning Research*, 21(140):1–67.
- Timo Schick and Hinrich Schütze. 2021. [Exploiting cloze-questions for few-shot text classification and natural language inference.](#) In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 255–269, Online. Association for Computational Linguistics.
- Lei Sha, Feng Qian, Baobao Chang, and Zhifang Sui. 2018. [Jointly extracting event triggers and arguments by dependency-bridge RNN and tensor-based argument interaction.](#) In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, (AAAI-18), the 30th innovative Applications of Artificial Intelligence (IAAI-18), and the 8th AAAI Symposium on Educational Advances in Artificial Intelligence (EAAI-18), New Orleans, Louisiana, USA, February 2-7, 2018*, pages 5916–5923. AAAI Press.
- Beth M. Sundheim. 1992. [Overview of the fourth Message Understanding Evaluation and Conference.](#) In *Fourth Message Understanding Conference (MUC-4): Proceedings of a Conference Held in McLean, Virginia, June 16-18, 1992*.
- David Wadden, Ulme Wennberg, Yi Luan, and Hananeh Hajishirzi. 2019. [Entity, relation, and event extraction with contextualized span representations.](#) In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 5784–5789, Hong Kong, China. Association for Computational Linguistics.
- Ziqi Wang, Xiaozhi Wang, Xu Han, Yankai Lin, Lei Hou, Zhiyuan Liu, Peng Li, Juanzi Li, and Jie Zhou. 2021. [CLEVE: Contrastive Pre-training for Event 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 6283–6297, Online. 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, Tianyu Liu, Lei Li, and Baobao Chang. 2021. [Document-level event extraction via heterogeneous graph-based interaction model with a tracker.](#) 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 3533–3546, Online. Association for Computational Linguistics.
- Hang Yang, Yubo Chen, Kang Liu, Yang Xiao, and Jun Zhao. 2018. [DCFEE: A document-level Chinese financial event extraction system based on automatically labeled training data.](#) In *Proceedings of ACL 2018, System Demonstrations*, pages 50–55, Melbourne, Australia. Association for Computational Linguistics.

Hang Yang, Dianbo Sui, Yubo Chen, Kang Liu, Jun Zhao, and Taifeng Wang. 2021. [Document-level event extraction via parallel prediction networks](#). 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 6298–6308, Online. Association for Computational Linguistics.

Shun Zheng, Wei Cao, Wei Xu, and Jiang Bian. 2019. [Doc2EDAG: An end-to-end document-level framework for Chinese financial event extraction](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 337–346, Hong Kong, China. Association for Computational Linguistics.

## A Dataset and Model

### A.1 Dataset statistics

We evaluate on three common datasets for Event Argument Extraction: ACE05 (Dodington et al., 2004), RAMS (Ebner et al., 2020) and WIKIEVENTS (Li et al., 2021).

ACE05 is a joint information extraction dataset providing entity, relation, and event annotation for three languages: English, Chinese, and Arabic. We use its English event annotation for sentence-level EAE tasks. We follow the pre-processing procedure of DyGIE++ (Wadden et al., 2019), which keeps 33 event types and 22 argument roles and collects 4859 arguments in the training set, 605 and 576 in the development and test set respectively.

RAMS is a document-level dataset annotated with 139 event types and 65 semantic roles. Each sample is a 5-sentence document, with trigger word indicating pre-defined event type and its argument scattering among the whole document.

WIKIEVENTS is another document-level dataset providing 246 documents, with 50 event types and 59 argument roles. These documents are collected from English Wikipedia articles that describe real-world events and then follow the reference links to crawl related news articles. They also annotate the coreference links of arguments, while we only use the annotations of their conventional arguments in this task.

Table 8 shows their detailed statistics.

### A.2 Details of baseline models

We compare our model with following previous models. (1) **ONEIE** (Lin et al., 2020): a joint

Dataset	ACE05	RAMS	WIKIEVENTS
<b>#Sents</b>			
<b>Train</b>	17,172	7,329	5,262
<b>Dev</b>	923	924	378
<b>Test</b>	832	871	492
<b>#Args</b>			
<b>Train</b>	4,859	17,026	4,552
<b>Dev</b>	605	2,188	428
<b>Test</b>	576	2,023	566
<b>#Event</b>	33	139	50
<b>#Role</b>	22	65	59
<b>#Arg per Event</b>	1.19	2.33	1.40

Table 8: Statistics of datasets.

model extracting entity, relation and event simultaneously. Different from QA-based model, they rely on extracted entities as candidate arguments. (2) **BART-Gen** (Li et al., 2021): a conditional generation model generating (rather than recognizing the spans) arguments sequentially via a sequence-to-sequence model and prompt. (3) **EEQA** (Du and Cardie, 2020): the first Question Answering (QA) based model designed for sentence-level EAE task. (4) **FEAE** (Wei et al., 2021): a QA-based method extended to document-level EAE by considering argument interactions via knowledge distillation. (5) **DocMRC** (Liu et al., 2021a): another QA-based method with implicit knowledge transfer and explicit data augmentation. The implementation details of all baselines are as follow:

1. **FEAE** (Wei et al., 2021): We report the results from the original paper.
2. **DocMRC** (Liu et al., 2021a): We report the results from original paper.
3. **BART-Gen** (Li et al., 2021): For BART-large model, We report the results from origin paper. For BART-base model, we use their code<sup>1</sup> to test its performance on all datasets.
4. **EEQA** (Du and Cardie, 2020): We report the results of ACE05 dataset from the origin papers. We use their code<sup>2</sup> to test its performance on RAMS and WIKIEVENT dataset. In order to generate the question template of these two datasets automatically, we follow the second template setting in **EEQA**. The question template is *What is the ROLE in TRIGGER WORD?*.

<sup>1</sup><https://github.com/raspberrycice/gen-arg>

<sup>2</sup><https://github.com/xinyadu/eeqa>



5. **EEQA-BART**: For a fair comparison with our model, we substitute the pre-trained model of EEQA from BERT to BART and call it EEQA-BART. We re-train the model on ACE05, RAMS and WIKIEVENTS datasets.
6. **ONEIE** (Lin et al., 2020): We use their code<sup>3</sup> and re-train the model on ACE05. We don't report its performance on RAMS and WIKIEVENTS because OneIE is a joint model extracting entity, relation and event. However, there is no entity annotation in RAMS and no relation annotation in both RAMS and WIKIEVENTS. Simply dropping the modules related to entity and relation in OneIE achieves abnormally low performance on RAMS and WIKIEVENTS dataset. Therefore it is somewhat unfair comparing OneIE with our model and other baselines in these two datasets.

For the models we re-trained, we keep all other hyper-parameters except learning rates the same with default settings in their original papers. We search the learning rate in [2e-5, 3e-5, 5e-5] and report the test set performance of the model that performs best on the development set.

### A.3 PAIE implementation and training setup

The optimization procedure of PAIE for one sample is shown in the pseudo code 1. We initialize the weight in encoder-decoder architecture with pre-trained BART models. The contexts in the document-level dataset sometimes exceed the constraint of BART-Encoder and consume prohibitively large memory; thus we add a window centering on the trigger words and only encode the words within the window. We train each large model on single NVIDIA-V100 GPU and each base model on a single NVIDIA-1080Ti GPU. For each setting, we train models with 5 fixed seeds [13, 21, 42, 88, 100] and 3 learning rates [2e-5, 3e-5, 5e-5]. Then we record the test set performance of the model that performs best on the development set for each random seed. The final reported performance is the average value of results w.r.t five different seeds. For model variations mentioned in Section 5.1, we only change the input strategy and leave other parts constant. We list other important hyperparameters in Table 9.

<sup>3</sup><http://blender.cs.illinois.edu/software/oneie/>

Hyperparameter	Value
Batch size	16 (ACE05) / 4 (Others)
Weight decay	0.01
Training steps	10000
Optimizer	AdamW
Adam $\epsilon$	$1 \times 10^{-8}$
Adam $\beta_1/\beta_2$	0.9 / 0.999
Scheduler	Linear (with 0.1 warmup step)
Max span length	10
Max gradient norm	5.0
Window size	250
Max encoder seq length	192 (ACE05) / 500 (Others)
Max decoder seq length	80

Table 9: Hyperparameters for PAIE

### A.4 Details of Bipartite Matching loss

We formulate the details of bipartite matching loss in this section. Given  $\text{logit}_k^{(start)}$  and  $\text{logit}_k^{(end)}$  from Eq 3, we apply greedy search on predicted start and end position distributions to select the predicted span for each role-specific selector  $\theta_k$ .

$$(\hat{s}_k, \hat{e}_k) = \arg \max_{(i,j) \in L^2, i < j} \text{logit}_k^{(start)}(i) + \text{logit}_k^{(end)}(j) \quad (8)$$

Denote  $y_r = [(s_0, e_0), \dots, (s_n, e_n)]$  as ground truth spans of role  $r$  for sample  $X$ , and  $\hat{y}_r = [(\hat{s}_0, \hat{e}_0), \dots, (\hat{s}_m, \hat{e}_m)]$  as predicted spans, where  $m$  is the number of occurrence of role  $r$  in the corresponding prompt.

With the candidate spans for each role, we define the bipartite matching between the candidates and ground truth annotations as finding the lowest cost of a permutation  $\Gamma$  of  $N$  elements:

$$\hat{\sigma} = \arg \min_{\sigma \in \Gamma_N} \sum_k^N L_1((s_k, e_k), (\hat{s}_{\sigma(k)}, \hat{e}_{\sigma(k)})) \quad (9)$$

where  $L_1((s_k, e_k), (\hat{s}_{\sigma(k)}, \hat{e}_{\sigma(k)}))$  represents  $L_1$ -norm between  $(s_k, e_k)$  and  $(\hat{s}_{\sigma(k)}, \hat{e}_{\sigma(k)})$ .

We introduce the classical Hungarian algorithm (Kuhn, 1955) for efficient optimal assignment. In Eq.9,  $N$  is chosen to the minimum value between  $m$  and  $n$ . If the number of candidate spans  $m$  is larger than the number of ground truth span  $n$ , we will pad  $(0, 0)$  representing no arguments to the golden answer set. Otherwise, we only select the optimally matched gold spans for bipartite matching loss calculation.

After finding the optimal assignment  $\hat{\sigma}$ , we align each ground truth span in  $y_r$  and each predicted span in  $\hat{y}_r$  according to the matching result and then calculate probabilities where the start/end positions

---

**Algorithm 1:** Training one sample

---

**Input:**  $X, Pt$  // Context, Prompt tokens  
**Data:**  $Y = \{r_0 : [[s_0^0, e_0^0], [s_0^1, e_0^1]], \{r_1 : [[s_1^0, e_1^0]]\}$   
 $H_{enc}, H \leftarrow \text{BART}(X)$   
 $\hat{P} \leftarrow \text{BART-Decoder}(Pt, H_{enc})$   
 $L \leftarrow 0$  // Initialize loss  
**foreach**  $role$  **in**  $Y.keys()$  **do**  
    **Set**  $\hat{Y}_{role}$  **to** empty list  
    **foreach**  $EMB_{slot}$  **in**  $\hat{P}.get\_next(role)$  **do**  
         $\psi \leftarrow \text{MeanPool}(EMB_{slot})$   
         $\psi^{(s)} \leftarrow \psi \circ \mathbf{W}^{(s)}$   
         $\psi^{(e)} \leftarrow \psi \circ \mathbf{W}^{(e)}$   
         $\text{logit}^{(s)} \leftarrow \psi^{(s)} H$  // cos-sim to H  
         $\text{logit}^{(e)} \leftarrow \psi^{(e)} H$  // cos-sim to H  
  
         $\hat{Y}_{role}.insert(\arg \max_{(i,j) \in L^2, i < j} \text{logit}^{(s)}(i) + \text{logit}^{(e)}(j))$   
    **end**  
     $Y_{role}, \hat{Y}_{role} \leftarrow \text{Hungarian}(Y_{role}, \hat{Y}_{role})$   
     $L \leftarrow L + \text{CrossEntropy}(Y_{role}, \hat{Y}_{role})$   
**end**

---

locate about role slot  $k$ . Note that we use the logit distribution of  $\hat{\sigma}(k)$  rather than  $k$ , which is different from Eq. 4 without bipartite matching:

$$\begin{aligned} p_k^{(start)} &= \text{Softmax}(\text{logit}_{\hat{\sigma}(k)}^{(start)}) \\ p_k^{(end)} &= \text{Softmax}(\text{logit}_{\hat{\sigma}(k)}^{(end)}) \end{aligned} \quad (10)$$

Given  $p_k^{(start)}$  and  $p_k^{(end)}$  obtained by Eq. 10, we follow the same loss function in Eq. 5 during training process. The bipartite matching is only applied in training. For inference, the model will output all non-zero spans with corresponding argument roles as predictions.

### A.5 Further analysis of Bipartite Matching

Ablation studies have validated the effectiveness of bipartite matching loss. In our settings, bipartite matching loss focuses on multiple arguments of the same role and reassigns the predicted arguments in each prompt slot. Since slots in our joint prompts usually entail different semantic meanings and matching preferences, even they are about the same roles, the permutation-invariance property of bipartite matching assures a global optimization of these arguments.

Such optimization especially makes sense when the arguments in context have subtle semantic distinction, and such distinction can not merely be captured by sequential order. Simple argument enumerations, for example, do not satisfy the condition mentioned above, while contexts with different syntactic structures are more likely to satisfy it. We consider such an instance when the prompt is in active voice *Person teaches Person*, but the context sentence is in passive voice *Peter is taught by John*.<sup>4</sup> The position-wise assignment will be likely to assign *Peter* to the first *Person* slot and *John* for the second. It is not semantically correct, although they are treated as correct in evaluation. This supervision signal will force the model downgrading to extract arguments by position order during the training process. Such a model is hard to be voice-aware (and also insensitive to capture other syntactic structures) and tends to misidentify multi-argument data during inference, as shown in Table 10.

However, we also acknowledge the improvement from bipartite matching loss is somewhat not significant and robust when compared with other contributions in our paper. We attribute it to the following points: **(1)** existing datasets are not designed especially for evaluating the multi-arguments problem, only 8.9% samples in ACE05, 6.1% in RAMS and 10.9% in WIKIEVENTS facing it. **(2)** Even in limited cases about multi-arguments, related arguments are usually simply enumerated and do not require complex analysis and matching about implicit structure. Thus we expect a large-scale dataset with more multi-arguments and diverse narrative styles in the future, and we believe the bipartite matching loss will bring more significant improvement in it.

## B Prompt Examples

We compare our prompt with others used in EAE task in Table 11. The first row gives a standard QA-based template, and the second row shows a standard prompt in the generation paradigm. Row 3-5 show our three types of joint prompts respectively.

We further show ten manual template examples of each dataset at Table 12. The complete version of different types of prompts is available in our codebase.

---

<sup>4</sup>This prompt and context sentence are imaginary and do not relate to any events/samples in three benchmarks. We use it just for the convenience of discussion.

Example	w/ Bipartite	w/o Bipartite
“We demand that the Security Council ... ,” said a spokesman for a <u>meeting</u> (Contact.Meet) Saturday of <b>Saddam</b> and top - level <b>officials</b> , quoted by media.	<b>Entity:</b> Saddam <b>Entity:</b> officials	<b>Entity:</b> Saddam <b>Entity:</b> ∅
..., bombing at the world-renowned race, where <b>he</b> and his brother, <b>Tamerlan</b> , 26, <u>set off</u> (Conflict.Attack.DetonateExplode) two pressure-cooker bombs near...	<b>Attacker:</b> Tamerlan <b>Attacker:</b> he	<b>Attacker:</b> Tamerlan <b>Attacker:</b> ∅

Table 10: Examples from our benchmark datasets. Prediction results for models with/without bipartite matching loss. Argument roles are boldfaced in example sentences, trigger words are underlined, and the event types are in brackets.

Prompt Type	Prompt Example
Question Answering Prompt (Du and Cardie, 2020)	Who is the Victor in the Conflict.defeat event? What is the ConflictOrElection in the Conflict.defeat event? Where is the Place in the Conflict.defeat event?
Conditional Generation Prompt (Li et al., 2021)	<arg1> defeated <arg2> conflict at <arg3> place
Manual Template (Ours)	<u>Victor</u> ( and <u>Victor</u> ) defeated in <u>ConflictOrElection</u> at <u>Place</u> ( and <u>Place</u> )
Concatenation Template (Ours)	<u>Victor</u> ( <u>Victor</u> ) <u>ConflictOrElection</u> <u>Place</u> ( <u>Place</u> )
Soft Prompt (Ours)	<Vic_left0> <u>Victor</u> <Vic_right0> ( <Vic_left0> <u>Victor</u> <Vic_right0> ) Defeated <Conf_left0> <u>ConflictOrElection</u> <Conf_right0> <Place_left0> <u>Place</u> <Place_right0> ( <Place_left0> <u>Place</u> <Place_right0> )

Table 11: Example prompts about Event type *Conflict.Defeat.Unspecified* in WIKIEVENTS dataset. Angle brackets in conditional generation prompt denote the content to be filled during the decoding stage. Angle brackets in soft prompt represents pseudo tokens connecting different slots. Underlined words in the last three rows denote role slots, and brackets include roles with multiple arguments.

Dataset	Event Type	Natural Lanugage Prompt
ACE05	Movement.Transport	<u>Agent</u> (and <u>Agent</u> ) transported <u>Artifact</u> (and <u>Artifact</u> ) in <u>Vehicle</u> (and <u>Vehicle</u> ) cost <u>Price</u> from <u>Origin</u> place (and <u>Origin</u> place) to <u>Destination</u> place (and <u>Destination</u> place)
	Justice.Arrest-Jail	<u>Agent</u> (and <u>Agent</u> ) arrested <u>Person</u> (and <u>Person</u> ) at <u>Place</u> (and <u>Place</u> ) for <u>Crime</u>
	Justice.Execute	<u>Agent</u> (and <u>Agent</u> ) executed <u>Person</u> at <u>Place</u> (and <u>Place</u> ) for <u>Crime</u>
	Conflict.Attack	<u>Attacker</u> (and <u>Attacker</u> ) attacked <u>Target</u> (and <u>Target</u> ) hurting <u>Victims</u> using <u>Instrument</u> (and <u>Instrument</u> ) at <u>Place</u> (and <u>Place</u> )
	Contact.Meet	<u>Entity</u> (and <u>Entity</u> ) met with <u>Entity</u> (and <u>Entity</u> ) at <u>Place</u> (and <u>Place</u> )
	Conflict.Demonstrate	<u>Entity</u> (and <u>Entity</u> ) demonstrated at <u>Place</u> (and <u>Place</u> )
	Transaction.Transfer-Ownership	<u>Seller</u> gave <u>Buyer</u> ( and <u>Buyer</u> , <u>Buyer</u> , <u>Buyer</u> , <u>Buyer</u> , <u>Buyer</u> , <u>Buyer</u> ) to <u>Beneficiary</u> ( and <u>Beneficiary</u> , <u>Beneficiary</u> ) for the benefit of <u>Artifact</u> ( and <u>Artifact</u> , <u>Artifact</u> ) cost <u>Price</u> at <u>Place</u> ( and <u>Place</u> , <u>Place</u> )
	Transaction.Transfer-Money	<u>Giver</u> (and <u>Giver</u> ) gave <u>Money</u> to <u>Recipient</u> (and <u>Recipient</u> ) for the benefit of <u>Beneficiary</u> (and <u>Beneficiary</u> ) at <u>Place</u> (and <u>Place</u> )
	Life.Be-Born	<u>Person</u> (and <u>Person</u> ) was born at <u>Place</u> (and <u>Place</u> )
Life.Marry	<u>Person</u> married <u>Person</u> at <u>Place</u> (and <u>Place</u> )	
RAMS	life.injure. illnessdegradationphysical artifactexistence.	<u>Victim</u> person has some physical degradation from <u>Medicalissue</u> imposed by <u>Injurer</u> at <u>Place</u>
	damagedestroy.destroy	<u>Destroyer</u> destroyed <u>Artifact</u> using <u>Instrument</u> in <u>Place</u>
	conflict.yield.surrender	<u>Surrenderer</u> surrendered to <u>Recipient</u> at <u>Place</u>
	conflict.yield.retreat	<u>Retreater</u> retreated from <u>Origin</u> place to <u>Destination</u> place
	contact.commandorder. correspondence	<u>Communicator</u> communicated remotely with <u>Recipient</u> about <u>Topic</u> at <u>Place</u>
	government.agreements. rejectagreementcontractceasefire	<u>Rejecternullifier</u> rejected or nullified an agreement with <u>Otherparticipant</u> in <u>Place</u>
	government.vote. violationspreventvote	<u>Preventer</u> prevented <u>Voter</u> from voting for <u>Candidate</u> on ballot in <u>Place</u>
	inspection.sensoryobserve. physicalinvestigateinspect	<u>Inspector</u> inspected <u>Inspectedentity</u> in <u>Place</u>
	manufacture.artifact. createintellectualproperty	<u>Manufacturer</u> manufactured or created or produced <u>Artifact</u> using <u>Instrument</u> at <u>Place</u>
life.injure. illnessdegradationsickness	<u>Victim</u> has disease sickness or illness at <u>Place</u> , deliberately infected by <u>Injurer</u>	
WIKI- EVENTS	ArtifactExistence. ManufactureAssemble	<u>ManufacturerAssembler</u> (and <u>ManufacturerAssembler</u> ) manufactured or assembled or produced <u>Artifact</u> (and <u>Artifact</u> ) from <u>Components</u> (and <u>Components</u> ) using <u>Instrument</u> (and <u>Instrument</u> ) at <u>Place</u> (and <u>Place</u> )
	Conflict.Demonstrate	<u>Demonstrator</u> was in a demonstration for <u>Topic</u> with <u>VisualDisplay</u> against <u>Target</u> at <u>Place</u> , with potential involvement of <u>Regulator</u> police or military
	Cognitive.Inspection. SensoryObserve	<u>Observer</u> (and <u>Observer</u> ) observed <u>ObservedEntity</u> (and <u>ObservedEntity</u> ) using <u>Instrument</u> (and <u>Instrument</u> ) in <u>Place</u> (and <u>Place</u> )
	Cognitive. TeachingTrainingLearning	<u>TeacherTrainer</u> (and <u>TeacherTrainer</u> ) taught <u>FieldOfKnowledge</u> (and <u>FieldOfKnowledge</u> ) to <u>Learner</u> (and <u>Learner</u> ) using <u>Means</u> (and <u>Means</u> ) at <u>Institution</u> (and <u>Institution</u> ) in <u>Place</u> (and <u>Place</u> )
	Control.ImpedeInterfereWith	<u>Impeder</u> (and <u>Impeder</u> ) impeded or interfered with <u>ImpededEvent</u> at <u>Place</u> (and <u>Place</u> )
	Transaction.Donation	<u>Giver</u> gave <u>ArtifactMoney</u> to <u>Recipient</u> (and <u>Recipient</u> ) for the benefit of <u>Beneficiary</u> (and <u>Beneficiary</u> ) at <u>Place</u> (and <u>Place</u> )
	Disaster.DiseaseOutbreak	<u>Disease</u> (and <u>Disease</u> ) broke out among <u>Victim</u> (and <u>Victim</u> ) or population at <u>Place</u> (and <u>Place</u> )
	Justice.TrialHearing	<u>Prosecutor</u> tried <u>Defendant</u> (and <u>Defendant</u> ) before <u>JudgeCourt</u> for <u>Crime</u> (and <u>Crime</u> ) in <u>Place</u> (and <u>Place</u> )
	Medical.Vaccinate	<u>Treater</u> vaccinated <u>Patient</u> via <u>VaccineMethod</u> for <u>VaccineTarget</u> at <u>Place</u> (and <u>Place</u> )
Personnel.StartPosition	<u>Employee</u> started working in <u>Position</u> at <u>PlaceOfEmployment</u> organization in <u>Place</u> (and <u>Place</u> )	

Table 12: Example manual templates used in our work. Underlined words denote role slots, and slots in brackets denote repetitive ones designed for multi-arguments of the same roles.