



# Handling polysemous triggers and arguments in event extraction: an adaptive semantics learning strategy with reward–penalty mechanism<sup>\*#</sup>

Haili LI<sup>1,2,4</sup>, Zhiliang TIAN<sup>†‡2</sup>, Xiaodong WANG<sup>1</sup>, Yunyan ZHOU<sup>3,5</sup>,  
 Shilong PAN<sup>1</sup>, Jie ZHOU<sup>1</sup>, Qiubo XU<sup>4,5</sup>, Dongsheng LI<sup>†‡1</sup>

<sup>1</sup>National Key Laboratory of Parallel and Distributed Computing,  
 National University of Defense Technology, Changsha 410073, China

<sup>2</sup>College of Computer Science and Technology, National University of Defense Technology, Changsha 410073, China

<sup>3</sup>Unit 63891 of PLA, Luoyang 471003, China

<sup>4</sup>Unit 63893 of PLA, Luoyang 471003, China

<sup>5</sup>State Key Laboratory of Complex Electromagnetic Environment Effects on  
 Electronics and Information System, Luoyang 471003, China

<sup>†</sup>E-mail: tianzhiliang@nudt.edu.cn; lidongsheng@nudt.edu.cn

Received Mar. 21, 2024; Revision accepted May 16, 2024; Crosschecked Feb. 10, 2025; Published online Mar. 13, 2025

**Abstract:** Event extraction (EE) is a complex natural language processing (NLP) task that aims at identifying and classifying triggers and arguments in raw text. The polysemy of triggers and arguments stands out as one of the key challenges affecting the precise extraction of events. Existing approaches commonly consider the semantic distribution of triggers and arguments to be balanced. However, the sample quantities of different semantics in the same trigger or argument vary in real-world scenarios, leading to a biased semantic distribution. The bias introduces two challenges: (1) low-frequency semantics is difficult to identify; (2) high-frequency semantics is often mistakenly identified. To tackle these challenges, we propose an adaptive learning method with the reward–penalty mechanism for balancing the semantic distribution in polysemous triggers and arguments. The reward–penalty mechanism balances the semantic distribution by enlarging the gap between the target and nontarget semantics by rewarding correct classifications and penalizing incorrect classifications. Additionally, we propose a sentence-level event situation awareness (SA) mechanism to guide the encoder to accurately learn the knowledge of events mentioned in the sentence, thereby enhancing target event semantics in the distribution of polysemous triggers and arguments. Finally, for various semantics in different tasks, we propose task-specific semantic decoders to precisely identify the boundaries of the predicted triggers and arguments for the semantics. Our experimental results on ACE2005 and its variants, along with the rich Entities, Relations, and Events (ERE), demonstrate the superiority of our approach over single-task and multi-task EE baselines.

**Key words:** Event extraction; Polysemous triggers; Polysemous arguments; Semantic imbalance; Reward–penalty mechanism

<https://doi.org/10.1631/FITEE.2400220>

**CLC number:** TP393.1

<sup>‡</sup> Corresponding authors

<sup>\*</sup> Project supported by the National Natural Science Foundation of China (Nos. 62306330 and 62106275), the Young Elite Scientist Sponsorship Program by China Association for Science and Technology (No. YESS20230367), and the Natural Science Foundation of Hunan Province, China (Nos. 2022JJ40558 and

WDZC20235250103)

<sup>#</sup> Electronic supplementary materials: The online version of this article (<https://doi.org/10.1631/FITEE.2400220>) contains supplementary materials, which are available to authorized users

ORCID: Zhiliang TIAN, <https://orcid.org/0000-0002-8906-5198>; Dongsheng LI, <https://orcid.org/0000-0001-9743-2034>

© Zhejiang University Press 2025

## 1 Introduction

Events, serving as carriers of information, possess significant research value due to their elevated information content and rich semantic details. The accelerated evolution of the Internet and the emergence of numerous Internet applications have brought a large number of unstructured and fragmented text resources. How to quickly and accurately obtain structured target event information from these resources has always been a key and challenging problem for scholars engaged in the field of event extraction (EE). EE (Ahn, 2006; Chen et al., 2015; Lu D et al., 2023) is the task of identifying and classifying triggers and arguments from unstructured text based on the predefined event schema, as shown in Fig. 1. EE enables users to obtain information in a timely and intuitive manner on who (doer), when (time), where (place), how (artifact), whom (recipient), and what (event) occurred. The extracted event can be widely used in downstream applications, such as event graph construction (Shu et al., 2021; Xu TY et al., 2022), recommendation systems (Cui ZJ et al., 2023; Xia et al., 2023), and decision aids (Anelli et al., 2022; You MS et al., 2023).

Many efforts have been devoted to EE. Earlier EE methodologies relied mainly on manually crafted multi-granularity features (Ji and Grishman, 2008; Hong et al., 2011; McClosky et al., 2011), which were labor-intensive. The emergence of deep learning techniques (Chen et al., 2015; Nguyen et al., 2016; Sha et al., 2018), capable of automatically learning features of tasks from extensive annotated data, has overcome the limitations of manual feature design. Recently, pre-trained language models (PLMs) (Yang S et al., 2019; Lin et al., 2020; Lu YJ et al., 2021; Liu X et al., 2022) with rich general language representations, such as BERT and RoBERTa, have become the backbone of EE systems, reducing the need for extensive annotated data. To address the challenges introduced by low-resource scenarios, including zero-shot and few-shot, prompts (Hsu et al., 2022; Wang SJ et al., 2023; Yao et al., 2023; Zhang KH et al., 2023) with task-specific knowledge aid PLMs in comprehending the content and format of tasks which require significant training. Large language models (LLMs) (Ettinger et al., 2023; Li et al., 2023; Pang et al., 2023) with exceptional text understanding and generation abilities do not require

training and are extensively applied in various natural language processing (NLP) tasks.

The polysemy of triggers and arguments poses significant challenges for EE (Feng et al., 2018; Ding et al., 2019). We take into consideration a trigger or argument of polysemy when it is associated with two or more meanings. We take polysemous triggers as an example to analyze the challenges, as shown in Fig. 1. The polysemous trigger “death” triggers three distinct events “life: die,” “justice: execute,” and “justice: sentence,” while “punishable” triggers only “justice: sentence.” The semantics of event types is finite, discrete, and predefined, which defines the output space for EE tasks. The output space of polysemous triggers and arguments is a complex space composed of multiple semantic subspaces. Mapping “death” to “justice: execute” among its three relevant meanings is more challenging than mapping “punishable” to “justice: sentence.”

Polysemy increases the complexity of EE tasks, making it difficult for models to determine the exact meaning represented by polysemous triggers and arguments. Specifically, the challenges introduced by polysemy to EE tasks manifest in two main aspects: (1) Semantic ambiguity. Polysemy makes it difficult for EE models to distinguish the semantics of triggers and arguments, manifested mainly in their possible various different semantics. It requires complex semantic disambiguation for EE models to map triggers and arguments to predefined types. (2) Context dependency. Polysemy affects the semantics of triggers and arguments by the context, allowing the same word to represent different semantics in different sentences. Consequently, the semantics of triggers and arguments may become ambiguous across different contexts, increasing the complexity of understanding and modeling their semantics for models.

To tackle the aforementioned challenges, existing EE studies have employed various methods to enhance the semantics of polysemous triggers and arguments, including context knowledge (Chen et al., 2015; Lu D et al., 2023), knowledge enhancement (Du and Ji, 2022), multi-task learning (Ping et al., 2023), and prompt-based approaches (Yao et al., 2023; Zhang KH et al., 2023), which treat the multiple semantics in polysemous triggers and arguments as balanced semantics. However, the semantic distribution is unbalanced. This imbalance poses some challenges for the semantic modeling of polysemous

triggers and arguments, as well as the identification of their boundaries, mainly in the following three aspects: (1) Biased semantic distribution. By analyzing the samples of polysemous triggers and arguments in the dataset, we find significant differences in the distribution of sample numbers for different semantics. We observe in Fig. 2a that the number of samples with the semantics “life: die” is much higher than those of “justice: execute” and “justice: sentence.” This imbalanced semantic distribution results in the model paying more attention to the high-frequency semantics in polysemous triggers and arguments during the learning process, while neglecting the acquisition of low-frequency semantics. Consequently, this further leads to the omission of low-frequency semantics and the erroneous identification

of high-frequency semantics. (2) Misidentification of relevant and irrelevant semantics. The probability of relevant and irrelevant semantics being detected is greater than that of the target semantics, resulting in a false positive (FP) that the nontarget semantics is identified and a false negative (FN) that the target semantic is misidentified, as illustrated in Fig. 2b. (3) Difficulty in boundary identification. For multi-token triggers and arguments, the polysemy of subtokens presents a substantial challenge in accurately identifying the boundaries of triggers and arguments.

To tackle these challenges, we introduce an adaptive semantics learning method for handling the imbalanced semantics in polysemous triggers and arguments using the reward–penalty mechanism,

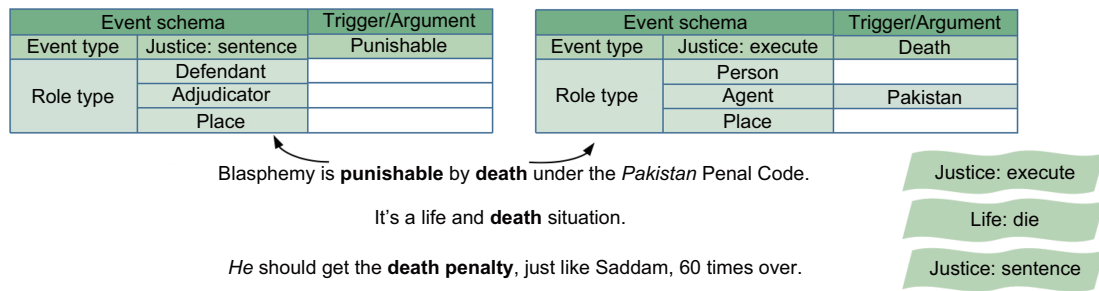


Fig. 1 Event extraction (EE) results of the sentence “Blasphemy is punishable by death under the Pakistan Penal Code.” and examples of the semantics for the triggers “death” and “punishable.” Words in bold are triggers, while those italicized are arguments. Event “justice: sentence” is triggered by the trigger “punishable” without arguments playing any roles. Trigger “death” triggers the “justice: execute” event, where “Pakistan” plays the “agent” role in the event schema of “justice: execute”

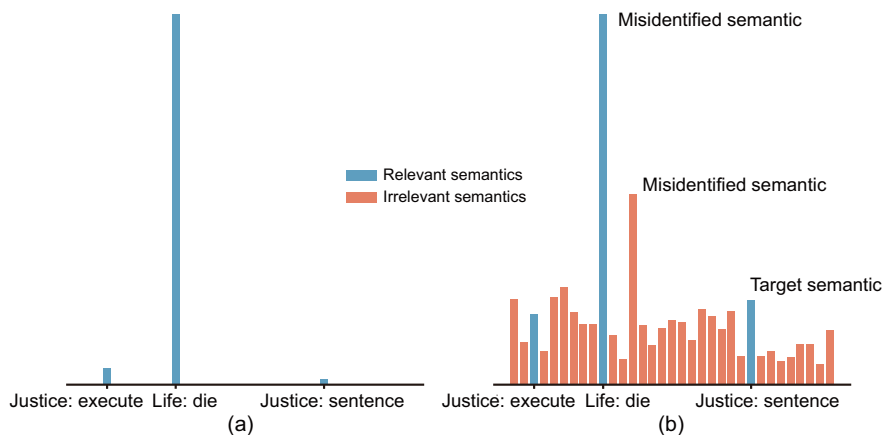


Fig. 2 Original semantic distribution of “death” (a) vs. semantic distribution of “death” after balancing (b). (a) shows the original imbalanced semantic distribution in which the number of samples with the semantics “life: die” is significantly higher than the numbers of samples with the semantics “justice: execute” and “justice: sentence.” (b) shows the semantic distribution after balancing, in which the probabilities of both relevant and irrelevant semantics being detected are higher than that of the target semantic, leading to two false positives (FPs) and one false negative (FN)

denoted as RPEE. First, we leverage the reward–penalty mechanism to balance the biased distribution of semantics by weakening the high-frequency semantics and amplifying the low-frequency semantics, and improve the classification accuracy by enlarging the gap between target semantics and nontarget semantics. This approach dynamically adapts the different semantics of polysemous triggers and arguments, based on the semantic probability distribution and the model classification outcomes, which differs from the traditional methods that adjust category weights (Yang S et al., 2019; Nam et al., 2022). Additionally, we use the sentence event semantics to enhance the semantics of triggers and arguments, intending to reduce FPs for irrelevant semantics and nontarget relevant semantics. To ensure the accuracy of the sentence event semantics for avoiding the error propagation, the proposed sentence-level event situation awareness (SA) mechanism uses a sentence event classification task for precisely modeling the sentence event semantics. Finally, we develop task-specific decoders to identify all candidate spans for triggers and arguments in the sentence, classifying their types for different semantics using task-specific thresholds. Experimental results demonstrate that our method effectively mitigates the imbalanced semantic distribution of polysemous triggers and arguments.

Our contributions are fourfold:

1. To solve the imbalanced semantic distribution of polysemous triggers and arguments, we introduce a semantics-adjustment method to minimize the gap between relevant semantics by using the reward–penalty mechanism.
2. We devise a reward–penalty mechanism to mitigate the biased distribution of semantics by dynamically adjusting different semantics in polysemous triggers and arguments.
3. The proposed sentence event situation awareness (SESA) mechanism provides correct event constraints for triggers and arguments in the sentence. Additionally, the task-specific decoder accurately identifies the boundaries of triggers and arguments composed of an uncertain number of tokens.
4. Extensive experiments demonstrate that RPEE outperforms state-of-the-art EE methods, demonstrating strong robustness, generalization ability, and superior performance in handling polysemous triggers and arguments, even in complex sce-

narios where triggers and arguments comprise multiple tokens.

## 2 Related works

This section reviews EE approaches and summarizes models that accurately identify the boundaries of triggers and arguments.

### 2.1 Event extraction (EE)

EE is a fundamental, crucial, and complex task in information extraction (IE) and NLP, focused on identifying triggers, event participants, and event types in text. Many efforts have been made from various perspectives to enhance EE performance. Researchers have used various neural networks like convolutional neural networks (CNNs) (Chen et al., 2015; Zeng et al., 2016), recurrent neural network (RNN) (Nguyen et al., 2016; Sha et al., 2018), long short-term memory (LSTM) (Feng et al., 2018; Lou et al., 2021), and graph neural network (GNN) (Liu X et al., 2018; Cui SY et al., 2020) to capture event features. From the perspective of leveraging knowledge of tasks and datasets, Du and Cardie (2020), Liu J et al. (2021), Yang P et al. (2021), Du and Ji (2022), Wang S et al. (2022), and Lu D et al. (2023) formalized EE as machine reading comprehension (MRC) or question answering (QA) tasks. Some researchers (Hsu et al., 2022; Ma et al., 2023; Ping et al., 2023; Yao et al., 2023; Zhang KH et al., 2023) employed well-designed prompts to guide language models in extracting events. Some others (Ettinger et al., 2023; Hsu et al., 2023; Yang YQ et al., 2023) used NLP tools to make use of syntactic, syntax, and semantic knowledge in the data. From the perspective of leveraging external resources, Liu X et al. (2022), Wang B et al. (2023), and Yao et al. (2023) tackled the challenge of data scarcity by generating instances. The methods mentioned above have made significant progress in EE. However, they tend to overlook the uneven distribution aspect of semantics in triggers and arguments, leading to numerous FPs that impact the overall performance of EE.

### 2.2 Boundary identification for EE

Identifying the boundaries of triggers and arguments is crucial for accurately extracting events,

especially for those consisting of multiple tokens. In this subsection, we review the literature on EE, specifically focusing on sequence labeling and span-based approaches.

### 2.2.1 Sequence labeling EE

The sequence labeling EE models formalize EE as sequence labeling, with the aim of modeling the semantic distribution of triggers and arguments. Various labeling schemes exist, including IO, BIO, BMES, and BIESO. In these schemes, the terms beginning, inside, outside, middle, end, and single correspond to each letter in BIOMES. Different methods use different labeling schemes. Guzman-Nateras et al. (2023) and Liu J et al. (2023) used the IO scheme for labeling. To better leverage the potential transferred knowledge between labels, Nguyen et al. (2016), Sha et al. (2018), Lin et al. (2020), Cong et al. (2021), Liu J et al. (2022), Wang ZT et al. (2023), and Xu ZY et al. (2023) used RNNs or conditional random fields (CRFs) and the BIO labeling scheme to model the boundaries of triggers and arguments. However, the sequence labeling method fails to handle nested triggers and arguments.

### 2.2.2 Span-based EE

In contrast to the methods rooted in sequence labeling, span-based modeling approaches aim to tackle intricate event structures, such as nested triggers and arguments. These approaches transform the EE task into a text span classification task, aiming to identify target triggers or arguments from all candidate spans and to classify each span's type. Depending on the modeling, span-based methods consist mainly of boundary location modeling and span representation modeling. Existing works (Yang S et al., 2019; Du and Cardie, 2020; Yang P et al., 2021; Xu RX et al., 2022; He et al., 2023) used two task-specific classifiers to model the head and tail tokens of the span. The works (Dozat and Manning, 2017; You HL et al., 2022, 2023; Ping et al., 2023) used the biaffine attention mechanism to jointly model the head and tail tokens of the span. Wadden et al. (2019) and Yang YQ et al. (2023) enumerated all spans, to model the joint representation of spans for multi-token triggers and arguments.

Sequence labeling EE methods model the semantic distribution of triggers and arguments but

fail to handle the nested or overlapping ones. However, span-based approaches tackle the issue but struggle with the imbalanced semantic distribution of polysemous triggers and arguments. To address this, we formalize the EE task as a token-classification problem and propose a reward-penalty mechanism to dynamically adjust the imbalanced semantic distribution of polysemous triggers and arguments, thereby mitigating their bias. Additionally, we design task-specific decoders to model the boundaries of triggers and arguments, separately.

## 3 Preliminaries

### 3.1 Task formulation

Following the definition by Ahn (2006), the process of EE consists of event detection (ED) (Liu J et al., 2023; Wang SJ et al., 2023) and event argument extraction (EAE) (He et al., 2023; Yang XJ et al., 2023), aiming at extracting triggers and arguments from the given sentence, as well as mapping them to the predefined types. We formalize ED and EAE as multi-label classification tasks to address the polysemy of triggers and arguments. The type sets of ED and EAE are denoted as  $E = \{e_1, e_2, \dots, e_M\} \cup \{e_0 = \text{"NULL"}\}$  and  $R = \{r_1, r_2, \dots, r_m\} \cup \{r_0 = \text{"NULL"}\}$ , respectively, where "NULL" indicates that the token is neither trigger nor argument.

For a given sentence  $X = \{x_1, x_2, \dots, x_n\}$ , where  $n$  is the length of tokens, ED identifies all candidate triggers for each meaning and presents the results in the format of  $\bigcup_{i=1}^M \left\{ [e_i : \bigcup_{j=1}^t [(s_{ij}, e_{ij})]] \right\}$ , where  $s_{ij}$  and  $e_{ij}$  represent the head and tail positions of the  $j^{\text{th}}$  candidate trigger for event type  $e_i$ ,  $e_i \in E$ , and  $t$  is the number of triggers for event  $e_i$ . According to the predefined event schema, the argument role set of  $e_i$  is  $r^i = \{r_1^i, r_2^i, \dots, r_a^i\}$ , where  $r_j^i \in R$ , and  $a$  indicates the number of argument roles for  $e_i$ . EAE recognizes all candidate arguments playing the role  $r_i$ , and the result is presented as  $\left\{ [e_i : \bigcup_{j=1}^a \left\{ r_j : \bigcup_{k=1}^b [(s_j^k, e_j^k)] \right\}] \right\}$ , where  $s_j^k$  and  $e_j^k$  represent the head and tail positions of the  $k^{\text{th}}$  candidate argument for role type  $r_j$ , respectively.

**Definition 1** In the training set, suppose that token  $x$  is labeled with a set of labels, denoted as  $e_x = \{e_{x1}, e_{x2}, \dots, e_{xg}\}$ , where  $e_{xi} \in E$  and  $g \leq M$ . Here,  $e_x$  represents the relevant semantics for token  $x$ , while

the remaining semantics  $e_u = E - e_r$  constitutes the irrelevant semantics of token  $x$ .

### 3.2 Situation awareness (SA)

SA (Endsley, 1988, 2001) perceives environment factors or events within the complex and dynamically changing information environment, comprehends their significance, and predicts their future states. SA is an intrinsic representation of the constantly changing external environment, which forms the fundamental basis for subsequent decision-making and performance.

Let  $E_v = \{en_1, en_2, \dots, en_n\}$  be the set of information in the environment, where  $en_i$  is the  $i^{\text{th}}$  kind of information. Sa is the function representing the SA model. Consequently,  $Sa(E_v) = \{s_1, s_2, \dots, s_s\}$  signifies the comprehensive state of the environment  $E_v$ , with  $s_i$  representing the  $i^{\text{th}}$  state component.

With the help of situational information, individuals or systems can better comprehend and adapt to complex environments. SA is widely used in various fields such as cybersecurity (Onwubiko, 2020; Matey et al., 2022), power systems (Dwivedi et al., 2023), disease prevention (Shashikumar et al., 2021), and traffic security (Zhang JR et al., 2023), and is also applied in specific tasks, such as emotion recognition (Akgun et al., 2023; Palash and Bhargava, 2023). To our knowledge, this paper is the first to introduce SA into EE, enhancing the understanding and adaptation to complex event environments.

### 3.3 Binary cross-entropy (BCE) loss

BCE loss (Zheng et al., 2022; Xu ZZ et al., 2023), commonly known as the sigmoid loss, employs the sigmoid function to compute probabilities and is commonly used by binary classification tasks and multi-label classification tasks. The sigmoid function independently calculates probabilities for each category, thereby preventing interference between different semantics. The formulation of BCE loss is

$$\mathcal{L}(\mathbf{Y}, \hat{\mathbf{Y}}) = -\frac{1}{c} \sum_{i=1}^c [y_i \log(\delta(r_i)) + (1 - y_i) \log(1 - \delta(r_i))], \quad (1)$$

where  $C = \{1, 2, \dots, c\}$  is the target set,  $\mathbf{Y} = [y_1, y_2, \dots, y_c]$  and  $\hat{\mathbf{Y}} = [\delta(r_1), \delta(r_2), \dots, \delta(r_c)]$  are the one-hot vector of the ground-truth label and the predicted label vector for input  $x$ , respectively,  $y_i \in$

$\{0, 1\}$ ,  $r_i$  is the logit value of  $x$  on class  $i$ ,  $r_i \in [0, 1]$ , and  $\delta$  is the sigmoid function. Suppose that the ground-truth label for  $x$  is class  $r_r$  and that the other classes are uniformly represented as  $u = C \setminus \{r_r\}$ . Then, the gradient with respect to the class  $r_r$  and  $u_i \in u$  are

$$\frac{\partial \mathcal{L}(\mathbf{Y}, \hat{\mathbf{Y}})}{\partial r_r} = \frac{\delta(r_r) - 1}{c}, \quad \frac{\partial \mathcal{L}(\mathbf{Y}, \hat{\mathbf{Y}})}{\partial r_{u_i}} = \frac{\delta(r_{u_i})}{c}. \quad (2)$$

## 4 Our approach

This paper presents a method to mitigate the biased semantic distribution of polysemous triggers and arguments using the reward-penalty mechanism. The overall framework, depicted in Fig. 3, consists of four main modules:

1. Reward-penalty mechanism dynamically adjusts the learning method of semantics in polysemous triggers and arguments. It rewards well-learned semantics while penalizing the erroneous one with the semantic probability distribution and the model classification outcomes.
2. The SESA mechanism generates an accurate and comprehensive representation for all events mentioned in sentences.
3. Semantics-enhanced encoder represents tokens with vectors and enhances semantics in tokens with all events mentioned in the sentence.
4. Task decoder identifies all potential trigger and argument candidates in the sentence and classifies their types.

Our training procedure comprises three phases: pretraining the SESA module (Section 4.2), then training ED and EAE with their respective semantics-enhanced encoder (Section 4.3), and using the task decoder (Section 4.4). ED provides EAE with target role sets based on the predefined event schema. The semantics-enhanced encoder furnishes the task decoder with token representations augmented by the sentence event semantics. SESA ensures the accuracy of the sentence event semantics provided to the semantics-enhanced encoder. Finally, the task decoder identifies the boundaries of triggers and arguments based on the representations of tokens and the reward-penalty mechanism (Section 4.1) and classifies their types.

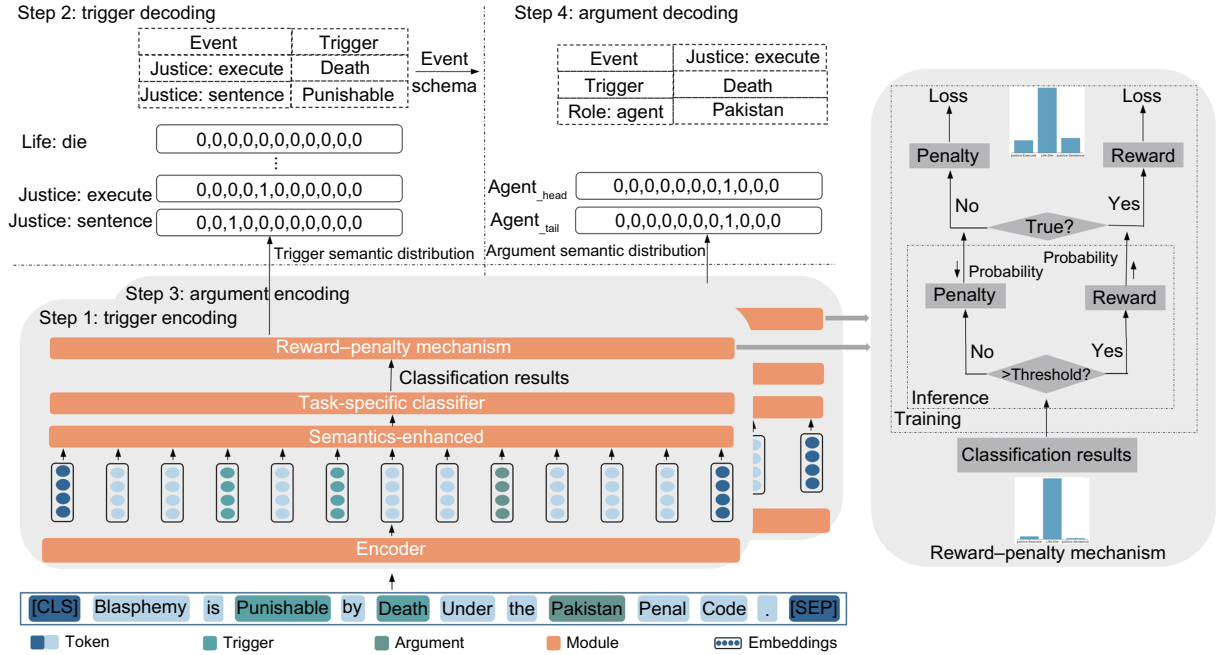


Fig. 3 The overview of our joint event extraction (EE) model. The encoder converts tokens into high-dimensional vectors. Before this, the sentence-level event situation awareness (SA) mechanism fine-tunes the encoder to guarantee the precise representation of the sentence events. The encoder then uses this representation  $S$  to enhance token semantics during encoding. Subsequently, our model calculates the probability distribution  $P(x_i)$  and employs the reward-penalty mechanism to amplify the correct semantics and diminish the incorrect semantics, widening the gap between them. Finally, using  $P(x_i)$ , the trigger decoder and argument decoder use task-specific thresholds to identify and classify all candidates

#### 4.1 Reward-penalty mechanism

The reward-penalty mechanism dynamically adjusts different semantics of polysemous triggers and arguments by rewarding semantics that is correctly classified and penalizing erroneous semantics. Subsequently, we provide a detailed analysis of the causes (Section 4.1.1), the desired effects (Section 4.1.2), and the implementation of the reward-penalty mechanism, covering both multi-factor (Section 4.1.3) and single-factor (Section 4.1.4) implementation methods.

##### 4.1.1 Motivational analysis of the reward-penalty mechanism

We analyze the misclassifications introduced by the imbalanced sample quantities from the following perspectives:

1. FP of irrelevant semantics. Let  $n_p$  and  $n_h$  denote the number of samples in the dataset annotated with labels  $p$  and  $h$ , respectively, where  $n_h > n_p$ .

The weight updating process is

$$w_{\text{new}} = w_{\text{old}} - \sum_{i=1}^n \frac{\delta(r_i) - 1}{c}, \quad (3)$$

where  $n$  indicates the number of labels.

When  $w_{\text{old}-p} = w_{\text{old}-h}$  and  $\delta(r_p) = \delta(r_h)$ , then  $w_{\text{new}-h} > w_{\text{new}-p}$ . For a test sample of class  $p$ , the trained model tends to categorize it as  $h$ , resulting in an FP.

2. FN of low-frequency semantics. Suppose that the two classes  $r_h$  and  $r_l$  of token  $x_i$  have  $k_1$  and  $k_2$  samples, respectively, with  $k_1 > k_2$ ,  $\{r_h, r_l\} \subseteq C$ . The accumulated gradients of  $r_h$  and  $r_l$  are

$$\sum_{i=1}^{k_1} \frac{\partial \mathcal{L}(\mathbf{Y}, \hat{\mathbf{Y}})}{\partial r} = \sum_{i=1}^{k_1} \frac{\delta(r_i) - 1}{c}. \quad (4)$$

When  $\delta(r_{r_h}) = \delta(r_{r_l})$ , the accumulated gradient value of class  $r_h$  is greater than that of class  $r_l$ , resulting in a biased semantic distribution for  $x_i$  and class  $r_l$  being overwhelmed by class  $r_h$ , as well as an FN for class  $r_l$ .

#### 4.1.2 Purposes of reward–penalty mechanism

The reward–penalty mechanism aims to mitigate these challenges by boosting the gradient gap between the relevant and irrelevant semantics in Eq. (3), and simultaneously diminishing the accumulated gradient gaps between semantics within the set of relevant semantics in Eq. (4). The gradient adjustment for semantics is detailed as follows:

$$\frac{\partial \mathcal{L}(\mathbf{Y}, \hat{\mathbf{Y}})}{\partial r} = \mathcal{K}_{r,u} \frac{\partial \mathcal{L}(\mathbf{Y}, \hat{\mathbf{Y}})}{\partial r}, \quad (5)$$

where  $\mathcal{K}_{r,u} = \mathcal{P}_{r,u} \mathcal{R}_{r,u}$  is the reward–penalty factor consisting of the reward factor  $\mathcal{R}_{r,u}$  and the penalty factor  $\mathcal{P}_{r,u}$ .

#### 4.1.3 Multi-factor reward–penalty mechanism

The multi-factor reward–penalty mechanism accurately adjusts various semantics based on classification outcomes and the probability distribution of tokens. It alleviates the bias in the semantic distribution by amplifying the probability of low-frequency semantics while diminishing that of high-frequency semantics and irrelevant semantics. We elaborate on the implementation process of the reward and penalty mechanisms, delineating the specific problems they target.

Reward mechanism uses a reward factor  $\mathcal{R}_{r,u}$  and the sample distribution to adjust the weight of the correctly classified semantics' gradients, including the true positive (TP) high-frequency semantics  $r_h$  and the true negative (TN) irrelevant semantics  $u$ . We adjust the weights of semantics' gradients as follows:

$$\mathcal{R}_{r,u} = \begin{cases} (r_h)^{\mathcal{R}_r}, & r_h \geq t, \text{ a TP happens,} \\ (r_u)^{\mathcal{R}_u}, & r_u \leq t, \text{ a TN happens,} \\ 1, & r_h < t \text{ or } r_u > t, \end{cases} \quad (6)$$

where  $\mathcal{R}_r > 0$  and  $\mathcal{R}_u > 0$  are hyperparameters, and  $t$  is the classification threshold. When a sample of  $r_h$  is correctly classified, the value of the reward factor  $\mathcal{R}_{r,u}$  decreases, resulting in a gradient enlargement for  $r_h$ , which narrows the gap in Eq. (4) between the accumulated gradient values of high-frequency and low-frequency semantics. Additionally, the weight of  $u$  diminishes to widen the gap between relevant and irrelevant semantics, making them easier to distinguish.

Penalty mechanism employs a penalty factor  $\mathcal{P}_{r,u}$  and the model classification outcomes to handle

misclassified semantics, including unidentified target semantics and misidentified irrelevant semantics. The adjustment process is as follows:

$$\mathcal{P}_{r,u} = \begin{cases} \left(\frac{t}{\mathcal{F}(r_{r_1})}\right)^{\mathcal{P}_r}, & \mathcal{F}(r_{r_1}) < t, \text{ an FN happens,} \\ \left(\frac{\mathcal{F}(r_u)}{t}\right)^{\mathcal{P}_u}, & \mathcal{F}(r_u) > t, \text{ an FP happens,} \\ 1, & \mathcal{F}(r_{r_1}) \geq t \text{ or } \mathcal{F}(r_u) \leq t, \end{cases} \quad (7)$$

where  $\mathcal{P}_r > 0$  and  $\mathcal{P}_u > 0$  are hyperparameters adjusting the punishment, and  $\mathcal{F}(\cdot) \in [0, 1]$  is the activation function that realizes the reward–penalty mechanism. When a small sample of token  $x$  is labeled  $r_1$  and an FN occurs, the penalty mechanism decreases the gradient of  $r_1$ , as illustrated in Eq. (5), thereby amplifying its weight, as described in Eq. (3). In the case of an FP, according to Eq. (7),  $\mathcal{P}_{r,u}$  increases, leading to an increase in the gradient of irrelevant semantics and consequent decrease in its weight.

#### 4.1.4 Single-factor reward–penalty mechanism

The numerous hyperparameters of each factor in the multi-factor reward–penalty mechanism pose significant challenges to accurate modeling of the semantic distribution of polysemous triggers and arguments. Therefore, we propose a single-factor reward–penalty mechanism with one hyperparameter to improve feasibility and usability by sacrificing some accuracy. We meticulously design the following reward–penalty function to implement the reward–penalty mechanism:

$$\text{pr}(x, \mathcal{K}_{r,u}) = \frac{1}{1 + e^{(-x)\mathcal{K}_{r,u}}}, \quad (8)$$

where  $x$  is the logit score of the input on a specified class, and  $\mathcal{K}_{r,u} > 1$  is an integer.

Subsequently, we will elaborate on how the function  $\text{pr}(x, \mathcal{K}_{r,u})$  implements the reward–penalty mechanism from the model training and inference perspective.

1. Training process. The reward–penalty function  $\text{pr}(x, \mathcal{K}_{r,u})$  directs the model training by modifying the loss based on the model classification outcomes. In the case of incorrect classification, the reward–penalty function  $\text{pr}(x, \mathcal{K}_{r,u})$  boosts the loss so as to encourage the model to learn more about the gold label of the input. Conversely, correct classifications result in  $\text{pr}(x, \mathcal{K}_{r,u})$  decreasing the loss, which helps avoid over-fitting and reduces the gap

between high-frequency and low-frequency semantics. Table 1 illustrates the impact on the loss value.

**Table 1** Changes in the loss value

		Prediction	
		Positive	Negative
Gold	Positive	–	+
	Negative	+	–

“+” and “–” denote the increase and decrease in the loss value using the function  $\text{pr}(x, \mathcal{K}_{r,u})$  compared to using the function  $\delta(x)$ , respectively

Specifically, when token  $x$  labeled with class  $i$  is correctly classified and  $r_i > 0$ ,  $\text{pr}(r_i, \mathcal{K}_{r,u}) > \delta(r_i)$ , the heightened probability serves as a reward for correct categorization. From Eq. (5), we observe that  $\text{loss}_{\text{pr}} < \text{loss}_{\text{s}}$ , where  $\text{loss}_{\text{pr}}$  and  $\text{loss}_{\text{s}}$  employ  $\text{pr}(r_i, \mathcal{K}_{r,u})$  and  $\delta(r_i)$  as the activation function, respectively. The decrease in the loss implies a reduced need for parameter tuning, further preventing the weight of class  $i$  from becoming excessively large. This illustrates that  $\text{pr}(r_i, \mathcal{K}_{r,u})$  effectively balances the semantic distribution.

In case where token  $x$  labeled with class  $r$  is unidentified,  $r_r < t$  or  $\text{pr}(r_r, \mathcal{K}_{r,u}) < \text{pr}(r_u, \mathcal{K}_{r,u})$ , an FN happens. With  $\text{pr}(r_r, \mathcal{K}_{r,u}) < \delta(r_r) < t$ ,  $\text{loss}_{\text{pr}} > \text{loss}_{\text{s}}$ , the increased loss can be seen as a penalty and leads the model to learn more about class  $r$ .

When token  $x$  does not have samples labeled with class  $u$  but  $\text{pr}(r_r, \mathcal{K}_{r,u}) < \text{pr}(r_u, \mathcal{K}_{r,u})$ , it indicates that  $u$  is misclassified. Using  $\text{pr}(x, \mathcal{K}_{r,u})$  gains larger loss than using  $\delta(x)$ , denoted as  $\text{loss}_{\text{pr}} > \text{loss}_{\text{s}}$ , and then the increased loss guides the model to reduce the weight of  $u$  to avoid FP.

2. Inference process. The reward–penalty function  $\text{pr}(x, \mathcal{K}_{r,u})$  simplifies the threshold setting by enlarging the gap between the target semantics and the nontarget semantics, improving the accuracy of identifying the target semantic.

Typically,  $t$  varies across tasks. Assuming that token  $x_i$  labeled with  $p$  is correctly recognized, then  $\text{pr}(r_p, \mathcal{K}_{r,u}) \geq t$ . At the same time, it can be observed from  $\text{pr}(r_p, \mathcal{K}_{r,u}) > \delta(r_p)$  that  $\text{pr}(x, \mathcal{K}_{r,u})$  rewards the target semantic  $p$  of  $x_i$  by enlarging its probability. Conversely, if the nontarget semantics  $u$  of  $x_i$  is misclassified,  $\text{pr}(r_u, \mathcal{K}_{r,u}) < t$  and  $\text{pr}(r_u, \mathcal{K}_{r,u}) < \delta(r_u)$ . The reward–penalty function  $\text{pr}(x, \mathcal{K}_{r,u})$  punishes the nontarget semantics  $u$  by

diminishing its probability. Simultaneously, by comparing changes in the probability of various semantics using different activation functions, we know  $\text{pr}(r_p, \mathcal{K}_{r,u}) > \delta(r_p) > \delta(r_u) > \text{pr}(r_u, \mathcal{K}_{r,u})$  and  $\text{pr}(r_p, \mathcal{K}_{r,u}) - \text{pr}(r_u, \mathcal{K}_{r,u}) > \delta(r_p) - \delta(r_u)$ . A wider boundary implies an easier setting for the threshold and a more accurate recognition of the target semantic.

## 4.2 SESA mechanism

SESA generates the joint representation of all events  $\mathbf{E}_{\text{SESA}} = [e_{s1}, e_{s2}, \dots, e_{sw}]$  mentioned in the sentence  $X$ , where  $e_{si} \in E$ , to enhance event semantics of tokens in  $X$ . To generate an accurate representation, inspired by Gururangan et al. (2020), we use a sentence-level event classification task and the same training dataset as the EE task to fine-tune SESA. Due to the rich general knowledge in PLMs (Devlin et al., 2019; Lewis et al., 2020), we use BERT (Devlin et al., 2019) as the backbone.

Hence, the following discussion outlines the components of SESA, focusing on learning and generation of sentence event representations.

Global encoder generates the global representation of all events mentioned in the input sentence with all tokens, to train and test SESA. The input sequence  $X'$  is constructed by adding the [CLS] token at the beginning of  $X$ , with all the corresponding masks in  $\text{ATTN\_MASK}$  set to 1. The embedding of [CLS], denoted as  $\mathbf{SG}_{\text{CLS}} = \text{BERT}(X', \text{ATTN\_MASK})$ , serves as the global representation of the sentence events. During experiments, we observe that the embedding of [CLS] in the last layer hidden state of the BERT output outperforms its counterpart in the `pooler_output`.

Event encoder generates a single high-dimensional vector representing all triggers and arguments in the sentence, exclusively for training purposes. The input sentence of the event encoder aligns with that of the global encoder. After filtering out tokens irrelevant to events, the representation of the sentence pure events is  $\mathbf{SE}_{\text{CLS}} = \text{BERT}(X, \text{EVENT\_MASK})$ , where  $\text{EVENT\_MASK}$  designates masks, corresponding to [CLS] and tokens in  $X$  labeled as triggers and arguments, with the value of 1. Experiments illustrate that the pure event representation enhances the global representation of the sentence event when modeling, ensuring the accuracy of modeling all

events mentioned in the sentence.

Event classifier identifies events with the joint representation of all events mentioned in the sentence  $S$  as follows: (1) Representation generation. During training, the output of the global encoder and the event encoder generates  $\mathbf{S} = [\mathbf{SG}_{\text{CLS}}; \mathbf{SE}_{\text{CLS}}]$ . However, during inference, only the global representation of all events mentioned in the sentence is used, so  $\mathbf{S} = \mathbf{SG}_{\text{CLS}}$ . (2) Event identification. The event classifier comprises a feedforward network (FFN) and an activation function. The FFN consists of a single-layer network structure and a rectified linear unit (ReLU) function. The classifier generates an event probability vector  $\mathbf{P} = [p_1, p_2, \dots, p_M]$ , where  $p_i$  denotes the probability of the occurrence of  $e_i$ . The threshold  $t$  is used to identify event types mentioned in  $X$ :

$$\hat{\mathbf{E}}_X = [e'_1, e'_2, \dots, e'_M], \quad e'_i = \begin{cases} 0, & p_i < t, \\ 1, & p_i \geq t, \end{cases} \quad i \in \{1, 2, \dots, M\}$$

$$\mathbf{E}_{\text{SESA}_X} = \hat{\mathbf{E}}_X \odot \mathbf{E} = [e_{s1}, e_{s2}, \dots, e_{sw}], \quad (9)$$

where  $e'_i \in \{0, 1\}$ , which indicates whether  $x_i$  triggers type  $e_i$ . When  $e_{si} = 1$ , it indicates that the sentence  $X$  triggers event type  $e_{si}$ . To streamline training, we use  $\text{pr}(x, \mathcal{K}_{r,u})$  as the activation function and employ BCE loss for optimization:

$$\mathcal{L}(\mathbf{Y}, \hat{\mathbf{Y}}) = - \sum_{i=1}^M [g_i y_i \log(\text{pr}(X)) + (1 - g_i)(1 - y_i) \log(1 - \text{pr}(X))], \quad (10)$$

where  $g_i$  denotes the weight of  $e_i$ ,  $\text{pr}(X) = \text{pr}(X, \mathcal{K}_{r,u})$ . To calculate  $g_i$ , we use the reciprocal for the ratio of the number of annotated samples of  $e_i$  to the total number of annotated samples in the dataset.

### 4.3 Semantics-enhanced encoder

The semantics-enhanced encoder transforms tokens in the input sentence  $X$  into real-valued word embedding. To mitigate the interference from nontarget event semantics and relevant semantics, we encode the knowledge of all the events mentioned in  $X$  into the representation of each token, enhancing token's event semantics. Specifically, we use the BERT fine-tuned by SESA to derive the event representation of  $X$ , denoted as  $\mathbf{S}_X = \text{SESA}(X, \text{ATTN\_MASK})$ , where  $\mathbf{S}_X \in \mathbb{R}^{1 \times d}$  and  $d$  represents the hidden layer dimension of BERT. We

use the fine-tuned BERT to encode the vector representation of each token in  $X$  as  $\{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_n\} = \text{BERT}(x_1, x_2, \dots, x_n)$ , where  $\mathbf{w}_i = [w_1^i, w_2^i, \dots, w_d^i]$ . The token enhanced representation, incorporating the knowledge of the sentence event, is expressed as  $\mathbf{w}'_i = [\mathbf{w}_i; \mathbf{S}_X]$ .

### 4.4 Task decoder

The task decoder recognizes the boundaries of candidate triggers and arguments in the sentence and classifies their types. Given that triggers and arguments are associated with distinct label sets and use different decoding modes, we develop task-specific decoders for triggers and arguments separately.

#### 4.4.1 Trigger decoder

The trigger decoder consists of  $M$  semantic decoders, each responsible for recognizing its trigger boundaries, based on the probability distribution of tokens, and classifying their types. Before decoding, we use the event classifier described in Section 4.2 to obtain the probability distribution  $p(x_i, S)$  of each token  $x_i$ . The decoding process includes the following steps:

1. Classification. Commonly used methods identify and classify triggers using the maximum probability leading to misclassification of low-frequency semantics. To address this issue, we employ a task-specific threshold value as the criterion for judgment. By applying  $t_t$ , based on Eq. (9) and the probability distribution of  $x_i$ , we derive the predicted type set for  $x_i$  as  $E_{xi} = \{e_{xi1}, e_{xi2}, \dots, e_{xit}\}$ , where  $e_{xii} \in E$  and  $p_{xii} \geq t_t$ .

2. Boundary identification. We employ the semantic decoders and threshold to recognize trigger's boundaries based on the predicted type set of tokens. Semantic decoder  $i$  identifies the boundaries of all triggers for  $e_i$ , where trigger boundaries are determined by consecutive tokens in the sentence that triggers the same event type. The decoding process for the trigger decoder is shown in Fig. 4.

#### 4.4.2 Argument decoder

We design a head decoder and a tail decoder for each semantic in the EAE task and use the task-specific threshold  $t_r$  to identify roles for arguments, guiding the model in modeling the boundary distribution of arguments.

	$x_1$	$x_2$	$x_3$	$x_4$	$x_5$	$x_6$	$x_7$	$x_8$	$x_9$	$x_{10}$	Output
$e_1$	0	0	0	0	0	0	0	0	0	0	$\{e_1: []\}$
$e_2$	0	0	1	1	1	0	0	0	0	0	$\{e_2: [[x_3, x_5]]\}$
$e_3$	0	0	0	0	0	0	0	0	0	0	$\{e_3: []\}$
$e_4$	0	0	0	0	1	0	0	1	0	0	$\{e_4: [[x_5, x_8], [x_8, x_8]]\}$
$e_5$	0	0	0	0	0	0	0	0	0	0	$\{e_5: []\}$
$e_6$	0	0	0	0	0	1	0	0	0	0	$\{e_6: [[x_6, x_6]]\}$

Fig. 4 The process of trigger decoding with the semantic decoders. 1 in  $(i, j)$  indicates that the predicted type for token  $x_j$  is  $e_i$ , and 0 in  $(i, j)$  indicates that token  $x_j$  does not trigger type  $e_i$ . Each row in this figure is a semantic decoder that identifies the boundaries of all candidates for the semantic. Each column is the predicted type set of the token

1. Boundary modeling. To accurately model boundaries of arguments, following previous studies (Yang S et al., 2019; Du and Cardie, 2020; Yang P et al., 2021), we adopt the head/tail labeling scheme to annotate boundaries of arguments. The probability distribution of token  $x_i$  is  $[p_{cr0}, p_{cr1}, \dots, p_{crm}] = \text{pr}(\text{FFN}_c(\mathbf{w}'_i), \mathcal{K}_{r,u})$ , where  $c \in \{\text{head}, \text{tail}\}$ ,  $\text{FFN}_{\text{head}}$  and  $\text{FFN}_{\text{tail}}$  are the head and tail role classifiers, respectively, and  $p_{cr_i}$  indicates the probability of  $x_i$  being the head or tail of the argument for role  $r_i$ .

2. Boundary identification. We employ  $t_r$  to identify the role types that  $x_i$  plays, following Eq. (9). The widely used boundary identification method is the enumeration (Wadden et al., 2019; Du and Cardie, 2020), which enumerates all predicted head-tail position combinations and identifies target boundaries with the heuristic method. However, we adopt the heuristic matching principle proposed by

Yang S et al. (2019), which selects the tail closest to the head as the target argument. The process of identifying head positions and tail positions of each semantic is analogous. Taking the head position recognition of class  $i$  as an example, the number of head positions is determined by the number of candidate chunks. Subsequently, the token with the highest probability is selected as the head of the chunk. The detailed decoding process of the role decoder is illustrated in Fig. 5.

## 4.5 Training

We design task-specific loss functions for ED and EAE separately to train the model, intending to learn the semantic distribution of polysemous triggers and arguments. Our model aims to intensify the learning of misjudged semantics by increasing the loss while diminishing the learning of correct semantics through loss reduction. To accomplish this, we employ BCE loss for training, based on the difference between the predicted probability distribution and the gold probability distribution. The loss functions for ED and EAE are formulated as follows:

$$\mathcal{L}_t = - \sum_{j=1}^M [g_j y_j \log(\text{pr}) + (1 - g_j)(1 - y_j) \log(1 - \text{pr})],$$

$$\mathcal{L}_{\text{role}}^c = - \sum_{j=1}^m [y_j \log(\text{pr}_c) + (1 - y_j) \log(1 - \text{pr}_c)],$$
(11)

where  $\mathcal{L}_t$  and  $\mathcal{L}_{\text{role}}^c$  indicate the loss functions for training the ED and EAE model, respectively,  $f$  is

	$x_1$	$x_2$	$x_3$	$x_4$	$x_5$	$x_6$	$x_7$	$x_8$	$x_9$	$x_{10}$
	pro	type	pro	type	pro	type	pro	type	pro	type
$r_{1\_head}$	-	0	-	0	-	0	-	0	0.79	1
$r_{1\_tail}$	-	0	-	0	-	0	-	0	0.88	1
$r_{1\_tail}$	-	0	-	0	-	0	-	0	0.75	1
$r_{l\_head}$	-	0	-	0	0.86	1	0.72	1	-	0
$r_{l\_tail}$	-	0	-	0	-	0	0.65	1	0.71	1
$r_{m\_head}$	-	0	-	0	-	0	-	0	-	0
$r_{m\_tail}$	-	0	-	0	-	0	-	0	-	0

■ Trigger   
■ Target argument   
■ Probability and role type   
■ Candidate head   
■ Candidate tail   
 Predicted head/tail

Fig. 5 The process of role decoding, where “pro” is the abbreviation of the probability indicating the probability of the token on the type, and the “type” value ( $\in \{0, 1\}$ ) means whether the type is the predicted one of the token. The results of role decoding are  $\{e_{x_6} : \{r_j : [[x_3, x_5], [x_7, x_7]]\}\}$ ,  $e_{x_9} : \{r_1 : [x_8, x_8]\}$ , where  $e_{x_i}$  is the event type for token  $x_i$

the task-specific classifier,  $\text{pr} = \text{pr}(f(\mathbf{w}'_i), \mathcal{K}_{r,u})$ , and  $\text{pr}_c = \text{pr}(f_c(\mathbf{w}'_i), \mathcal{K}_{r,u})$ .

### 5 Theoretical analysis of the reward-penalty mechanism

In this section, we theoretically analyze the effectiveness of the reward-penalty mechanism from the perspective of enlarging the gap between semantics and achieving balanced semantics learning.

#### 5.1 Analysis from the perspective of enlarging the gap

Function  $\text{pr}(x, \mathcal{K}_{r,u})$  augments the maximum gradient value for better training and enlarges the gap between the relevant and irrelevant semantics. The maximum gradient value of  $\text{pr}(x, \mathcal{K}_{r,u})$  is  $\mathcal{K}_{r,u}$  times that of the sigmoid function, which effectively mitigates gradient vanishing during back-propagation, as depicted in the following equation:

$$\max \left( \frac{\text{pr}'(x, \mathcal{K}_{r,u})}{\delta'(x)} \right) = \mathcal{K}_{r,u}. \quad (12)$$

It is evident that the larger the discrepancy between the probabilities of target and nontarget semantics, the simpler it is to set the threshold  $t$  for classification. As depicted in Fig. 6, the maximum value of  $\text{pr}'(x, \mathcal{K}_{r,u}) - \delta'(x)$  increases as  $\mathcal{K}_{r,u}$  grows when  $\mathcal{K}_{r,u}$  exceeds 1. However, the effective range of the reward-penalty mechanism, denoted as  $[-x_e, x_e]$ , decreases, where  $\delta'(\pm x_e) = \text{pr}'(\pm x_e, \mathcal{K}_{r,u})$ . When  $x \in [-x_e, x_e]$  and  $\Delta x \geq 0$ , if  $\text{pr}'(x, \mathcal{K}_{r,u}) \geq \delta'(x)$ , then  $\text{pr}(x, \mathcal{K}_{r,u}) - \text{pr}(x - \Delta x, \mathcal{K}_{r,u}) \geq \delta(x) - \delta(x - \Delta x)$ . Thus, the reward-penalty function  $\text{pr}(x, \mathcal{K}_{r,u})$  widens the gap between semantics, thereby reducing misclassifications.

#### 5.2 Analysis from the perspective of balanced learning semantics

The reward-penalty mechanism aims to rebalance the distribution of semantics by adjusting the training loss. This loss of token semantics comprises two components: the loss associated with target semantics and nontarget semantics, as illustrated in Eq. (13). To address misjudgments, the reward-penalty function  $\text{pr}(x, \mathcal{K}_{r,u})$  fine-tunes the model with a penalty mechanism to enhance the understanding of the target semantics and decrease the learning ability of nontarget semantics, thereby

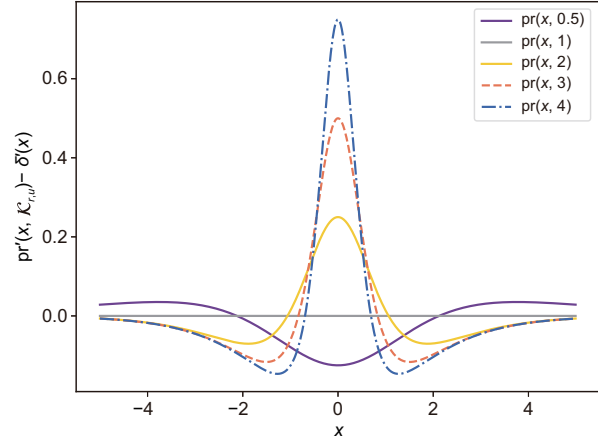


Fig. 6 Results of  $\text{pr}'(x, \mathcal{K}_{r,u}) - \delta'(x)$

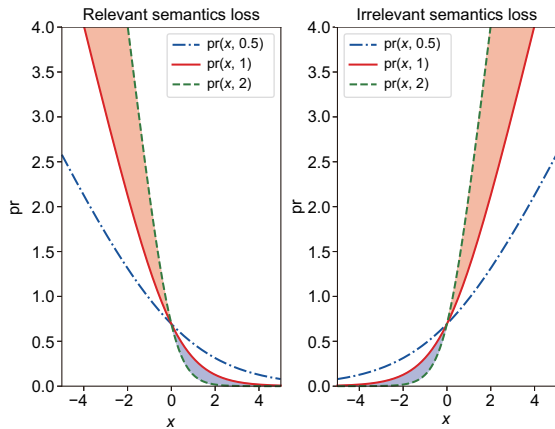
widening the gap between target and nontarget semantics. Concurrently, the reward mechanism is employed to mitigate the model learning of correctly classified semantics, aiming to reduce the gap between relevant semantics and achieve a balanced semantic distribution.

$$\begin{aligned} \mathcal{L}(Y_{i,j}, \hat{Y}_{i,j}) = & - \underbrace{[y_{i,j} \log(\text{pr}(f(x_i), \mathcal{K}_{r,u}))]}_{\text{the target semantic loss}} \\ & + \underbrace{[(1 - y_{i,j}) \log(1 - \text{pr}(f(x_i), \mathcal{K}_{r,u}))]}_{\text{the nontarget semantic loss}}, \end{aligned} \quad (13)$$

where  $Y_{i,j}$  and  $\hat{Y}_{i,j}$  are the ground-truth value and the predicted probability value of token  $x_i$  on class  $j$ , respectively, and  $y_{i,j} \in \{0, 1\}$ .

The penalty mechanism guides the model to accurately learn misclassified semantics by enlarging the loss, as illustrated by the red region in Fig. 7. Suppose that the ground-truth label of token  $x_i$  is  $j$ , where  $p_{i,j} < 0$  and  $t = 0.5$ . In this case,  $x_i$  is not identified as class  $j$ , resulting in an FN. Meanwhile, for token  $x_i$  not labeled with class  $b$  with  $p_{i,b} > 0$ ,  $x_i$  is identified as class  $b$ , leading to an FP. Referring to Eq. (13) and Fig. 7, it is apparent that when  $\mathcal{K}_{r,u} > 1$ , FPs and FNs lead to an increase in loss, which is equivalent to imposing a penalty.

The reward mechanism reduces the gap among relevant meanings by decreasing the loss, as the blue region depicted in Fig. 7. When  $x_i$  is a sample for class  $j$  and is accurately classified as  $j$ , yielding a TP. Referring to Eq. (13) and Fig. 7, it is observed that  $\text{loss}_{\text{pr}} < \text{loss}_{\delta}$  in the case where a TP occurs and  $\mathcal{K}_{r,u} > 1$ . The TP essentially leads to a diminished loss, akin to receiving a reward.



**Fig. 7** Loss of relevant and irrelevant semantics. References to color refer to the online version of this figure

## 6 Experiments

### 6.1 Experimental setup

#### 6.1.1 Datasets

We evaluate our model on two public EE benchmarks, the ACE2005 corpus (<https://catalog.ldc.upenn.edu/LDC2006T06>) and the rich Entities, Relations, and Events (ERE) corpus (Song et al., 2015) (Here we use datasets LDC2015E29, LDC2015E68, and LDC2015E78 in ERE), primarily using their English corpora, denoted as ACE2005-E+ and ERE-EN, respectively. Following Lin et al. (2020) and Hsu et al. (2022), ACE2005-E+ is spilt into three parts: the training set with 529 documents, the development set with 30 other documents, and the test set with the remaining 40 newswire documents. Additionally, following the pre-processing in Lin et al. (2020), we obtain the variant ACE05-E. Both have been annotated with 33 event types and 22 argument roles. The pre-processing of ERE-EN follows Lin et al. (2020), involving 38 event types and 21 argument roles. Moreover, we add a “NULL” type for tokens without annotation. ACE05-E differs from ACE05-E+ and ERE-EN in that the latter are more complex datasets with multi-token triggers. Statistical results of ACE05-E, ACE05-E+, and ERE-EN are shown in Table 2.

#### 6.1.2 Evaluation metrics

We adhere to the criteria employed in previous studies (Wadden et al., 2019; Hsu et al., 2022).

**Table 2** Dataset statistical results

Dataset	Set	Num <sub>s</sub>	Num <sub>e</sub>	Num <sub>r</sub>
ACE05-E	Training	17 172	4202	4859
	Dev	923	450	605
	Test	832	403	576
ACE05-E+	Training	19 216	4419	6607
	Dev	901	468	759
	Test	676	424	689
ERE-EN	Training	14 736	6208	8924
	Dev	1209	525	730
	Test	1163	221	822

Num<sub>s</sub>: number of sentences; Num<sub>e</sub>: number of events; Num<sub>r</sub>: number of roles. Dev: development

(1) Trigger identification (Trig-I): a trigger is correctly identified only if its predicted span matches that of the gold trigger perfectly. (2) Trigger classification (Trig-C): the event type of the trigger is correctly classified only if its predicted type matches that of the gold trigger. (3) Argument identification (Arg-I): an argument is correctly identified only if its predicted span matches that of the gold argument. (4) Argument classification (Arg-C): the role type of an argument is correctly classified only if its predicted role type and event type match those of the gold argument. Simultaneously, we use the widely used evaluation metrics, including precision ( $P$ ), recall ( $R$ ), and Micro F1 score (F1), to assess the performance.

#### 6.1.3 Parameter settings

We conduct all experiments on one NVIDIA RTX 3090, with a learning rate of  $1e-5$  and a weight decay of  $1e-5$  for BERT, and a learning rate of  $1e-4$  and a weight decay of  $1e-2$  for the other models. The batch size is 32. The number of epochs for EAE is 50 and the other is 30. The dropout rate is 0.5. We set our seed as 42. The thresholds for the tasks of sentence event classification, trigger identification, and argument identification are  $t_e$ ,  $t_t$ , and  $t_r$ , respectively. We employ AdamW (Loshchilov and Hutter, 2019) to optimize the model, and the maximum gradient clipping is set to 5 to avoid over-fitting.

#### 6.1.4 Baselines

Single-task EE models solely rely on event annotations for EE. In contrast, multi-task EE models perform EE with the help of entity recognition, relation extraction, or entity annotations. Since not all event corpora extensively annotate entities and

relationships, EE models relying solely on event annotations are more versatile.

1. Single-task EE models. (1) DMCNN (Chen et al., 2015) uses dynamic multi-pooling CNNs to capture features of word level and sentence level. (2) BERT\_QA (Du and Cardie, 2020) formalizes the EE task as a QA and designs task-specific question templates for trigger extraction and argument extraction. (3) LEAR (Yang P et al., 2021) enhances tokens' task semantics by encoding the label annotation into the token representation. (4) Text2Event (Lu YJ et al., 2021) uses a curriculum learning approach and constrained decoding to accomplish sequence-to-structure tasks in document level EE. (5) DEGREE (Hsu et al., 2022) proposes an end-to-end event generation model that generates events from predefined event type specific templates. (6) GTEE-DynPref (Liu X et al., 2022) is a template-based generative EE method that adopts dynamic prefix-tuning technique. (7) DAEE (Wang B et al., 2023) enhances EE by using reinforcement learning and event knowledge to generate high-quality data to augment EE. (8) DemoSG (Zhao et al., 2023) uses knowledge of the annotated data and label semantics to conduct EE in low resources. (9) ChatGPT-ICL (Han et al., 2023) is a prompt-based inference-only method that conducts 14 sub-tasks of IE to evaluate the performance and robustness of ChatGPT.

2. Multi-task EE models. (1) DYGIE++ (Wadden et al., 2019) learns distant contextual features by using dynamic span graphs. (2) ONEIE (Lin et al.,

2020) obtains optimal event graphs by using beam search and global features. (3) UniEX (Ping et al., 2023) proposes the triaffine attention mechanism to encode the schema of all tasks and their label semantics into token semantics to fully improve the comprehensive semantics.

## 6.2 Effectiveness

In this subsection, we conduct extensive experiments on the three datasets to assess the effectiveness of RPEE. Detailed content related to the case study is provided in Section 1 in the supplementary materials. A comprehensive discussion of RPEE is presented in Section 2 in the supplementary materials.

### 6.2.1 Main results

Table 3 illustrates the experimental results of all baselines and our method on ACE05-E. When comparing the performance, we obtain the following findings: (1) Our method surpasses all the baselines in terms of F1 score. This indicates the effectiveness of the proposed RPEE on the EE task. (2) In terms of Trig-C, compared with state-of-the-art methods of both multi-task and single-task EE, our method achieves relative performance improvements of 5.2% and 3.7% on F1, respectively. Our method does not use manually crafted prompts, complex language models, or NLP tools, and achieves significant results with limited annotated data, showcasing strong scalability and generalization ability. (3) In terms of

**Table 3 EE results on ACE05-E**

Task	Model	PLM	F1 score	
			Trig-C	Arg-C
Multi-task EE	DYGIE++ (Wadden et al., 2019)	BERT-base	73.6	52.5
	ONEIE (Lin et al., 2020)	BERT-base	74.7	<u>56.8</u>
	UniEX (Ping et al., 2023)	RoBERTa-large	74.1	53.9
Single-task EE	DMCNN (Chen et al., 2015)	–	69.1	53.5
	BERT_QA (Du and Cardie, 2020)	BERT-base	72.4	53.3
	LEAR* (Yang P et al., 2021)	BERT-base	72.2	
	Text2Event (Lu YJ et al., 2021)	T5-large	71.9	53.8
	DEGREE (Hsu et al., 2022)	BART-large	73.3	55.8
	GTEE-DynPref (Liu X et al., 2022)	BART-large	72.6	55.8
	DAEE (Wang B et al., 2023)	BART-large	<u>75.8</u>	56.5
	DemoSG (Zhao et al., 2023)	BART-large	73.4	56.0
	ChatGPT-ICL (Han et al., 2023)	GPT-3.5-turbo	27.3	31.6
RPEE (Ours)	BERT-base	<b>78.6</b>	<b>59.0</b>	

The highest scores are highlighted in bold, while the sub-optimal scores are underlined. “–” indicates the absence of PLM usage. The symbol “\*” denotes the results obtained by using the same dataset and data pre-processing outlined in this paper. Arg-C: argument classification; EE: event extraction; PLM: pre-trained language model; Trig-C: trigger classification

Arg-C F1, our model exhibits performance enhancements of 3.9% and 4.4% compared with multi-task and single-task EE models, respectively, highlighting the effectiveness of our approach in this task. (4) Concerning PLMs, our method not only outperforms the baselines that also use BERT-base but also outperforms the baselines that employ larger PLMs, such as BART-large, demonstrating the superiority of our approach in applications.

To further verify the scalability and robustness of our method, we conduct experiments on ACE05-E+ and ERE-EN, and present the results in Table 4. Upon analyzing the results, we derive two crucial conclusions:

1. High scalability

Our approach has demonstrated strong performance on ACE05-E+ and ERE-EN. Our method surpasses all the baseline methods on Trig-C and Arg-C in terms of the F1 score. Notably, we observe a relative enhancement of 2.9% and 0.1% on ACE05-E+ and ERE-EN, for F1 score in Trig-C, and a relative improvement of 8.0% and 6.5% for F1 score in Arg-C, respectively. These findings underscore the scalability and efficacy of our method, signifying its suitability for EE across diverse domains.

2. Strong robustness

Our method exhibits superior performance on ACE05-E+ compared to all the baselines and even outperforms itself on ACE05-E, as shown in Tables 3 and 4. Notably, baselines generally perform better on ACE05-E than ACE05-E+. The discrepancy is attributed to the presence of multiple tokens, posing a substantial challenge for models to precisely model trigger and argument boundaries. Our method adeptly leverages multi-token instances, enhancing model performance via the reward-penalty mechanism and the task-specific decoding strategy. Experimental results affirm the robustness of our approach in handling intricate scenarios involving multiple tokens.

### 6.2.2 Ablation study

This subsection focuses on a detailed analysis of the impact of each component on the performance, with corresponding experimental results on ACE05-E+ presented in Table 5.

1. w/o SESA indicates the variant without the SESA module. Significantly, there is a notable performance decrease compared with RPEE,

**Table 4** Experimental results on ACE05-E+ and ERE-EN

Model	F1 score			
	ACE05-E+		ERE-EN	
	Trig-C	Arg-C	Trig-C	Arg-C
ONEIE	72.8	54.8	59.1	50.5
LEAR*	71.4		57.0	
Text2Event	71.8	54.4	59.4	48.3
DEGREE	70.9	<u>56.3</u>	57.1	49.6
GTEE-DynPref	74.3	54.7	<u>66.9</u>	<u>55.1</u>
DAEE	<u>76.9</u>	<u>56.3</u>	65.0	51.6
RPEE (Ours)	<b>79.1</b>	<b>60.8</b>	<b>67.0</b>	<b>58.7</b>

The highest scores are highlighted in bold, while the sub-optimal scores are underlined. The symbol “\*” denotes the results obtained by using the same dataset and data pre-processing outlined in this paper. Arg-C: argument classification; Trig-C: trigger classification

affirming the effectiveness of the SESA proposed in Section 4.2. This discrepancy is attributed to the representation of the sentence event that is well-learned during the SESA pre-training, offering event constraints for triggers in the sentence.

2. w/o Event Encoder denotes the absence of the event encoder in the SESA module. Compared with the “w/o SESA” variant, “w/o Event Encoder” learns the representation of the sentence event, and the performance instead decreases. It illustrates that pre-training without pure event knowledge hinders the accurate acquisition of sentence event knowledge, resulting in a degraded model.

3. w/o Re-weighting signifies ignoring the imbalanced distribution of events. The findings suggest that this uneven distribution impacts the learning of event semantics, thereby influencing the overall performance of EE.

4. w/o SA designates the variant that omits the utilization of the representation of the sentence event, resulting in the poorest performance among all the variants. It illustrates the significance of incorporating the representation of all the events mentioned in the sentence, as it effectively enhances the event semantics of tokens within the sentence. Enhanced representation offers vital event constraints for triggers. The findings emphasize the pivotal role of the SA in the overall model effectiveness.

5. w/o Reward-Penalty indicates the variant without using the reward-penalty mechanism, displaying inferior performance compared to most variants. The declining performance underscores the crucial role of the reward-penalty mechanism in

achieving a balanced modeling of the token event distribution within the model.

### 6.2.3 Analysis of SESA

In this subsection, we conduct experiments to analyze the effectiveness of SESA and the impact of the event encoder, as shown in Fig. 8. Results on ACE05-E, ACE05-E+, and ERE-EN reveal that, SESA excels in the sentence event classification task, achieving remarkable F1 values surpassing 93, 92, and 89, respectively. It indicates that SESA has well learned the representation of all events mentioned in the sentence. When fixing  $\mathcal{K}_{r,u}$ , the performance varies with the adjustment of  $t_e$ . There exists an optimal  $t_e$  for  $\mathcal{K}_{r,u}$  to achieve the best performance. Furthermore, experimental results suggest that using the reward-penalty mechanism has a negligible impact on the SESA performance. To reduce the complexity of training, SESA adopts the identical configuration of the reward-penalty function with Trig-C.

To further assess the influence of pure event knowledge on modeling all the events mentioned in

the sentence, we conduct experiments by excluding the event encoder module, with detailed results in Table 6. Experimental results emphasize that relying solely on the knowledge of all the tokens in the sentence for modeling the sentence event yields unsatisfactory results, leading to severe errors in downstream tasks, as depicted in Section 6.2.2. The findings underscore the pivotal role of pure event knowledge in effectively modeling the sentence event.

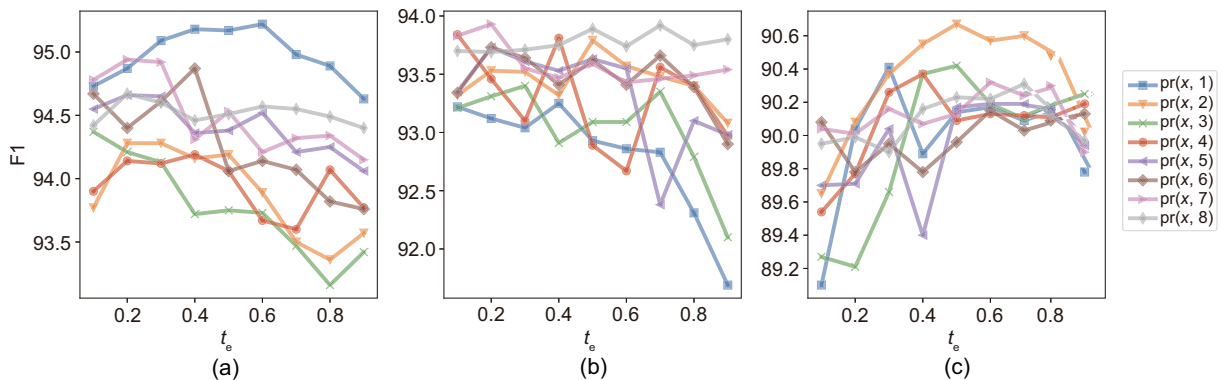
### 6.2.4 Sensitivity test

We analyze the impact of different configurations of key hyperparameters in RPEE, specifically,  $\mathcal{K}_{r,u}$  in the reward-penalty mechanism and  $t_t$  during the decoding process. These hyperparameters are individually adjusted, while the remaining parameters remain consistent with the previously reported settings. Taking the ED task as the case study, we conduct experiments with various settings of  $\mathcal{K}_{r,u}$  and  $t_t$  on ACE05-E, ACE05-E+, and ERE-EN. The plot in Fig. 9 illustrates the fluctuation of  $P$ ,  $R$ , and F1 scores for Trig-I and Trig-C across different hyperparameter settings on ACE05-E. It is evident from

**Table 5 Ablation study on ACE05-E+**

Model	Trig-I			Trig-C			Arg-I			Arg-C		
	$P$	$R$	F1	$P$	$R$	F1	$P$	$R$	F1	$P$	$R$	F1
RPEE	<b>85.22</b>	<b>79.18</b>	<b>82.09</b>	<b>81.63</b>	<b>76.79</b>	<b>79.14</b>	<b>73.95</b>	<b>64.91</b>	<b>69.13</b>	<b>62.31</b>	<b>59.37</b>	<b>60.80</b>
w/o SESA	84.46	<u>75.08</u>	79.49	81.19	72.77	76.75	68.11	60.20	63.91	55.89	56.10	56.00
w/o Re-weighting	84.19	77.47	80.69	79.54	73.87	76.60	70.71	62.41	66.30	54.97	56.97	55.95
w/o Event Encoder	<u>82.22</u>	76.01	78.99	<u>78.41</u>	73.31	75.77	66.76	63.57	65.12	52.84	57.78	55.20
w/o SA	82.86	75.26	<u>78.88</u>	78.82	<u>71.81</u>	<u>75.15</u>	72.56	<u>52.85</u>	<u>61.16</u>	59.56	<u>45.05</u>	<u>51.30</u>
w/o Reward-Penalty	82.82	76.15	79.35	78.71	72.66	75.56	<u>61.97</u>	63.73	62.83	<u>50.14</u>	57.86	53.72

The highest scores are highlighted in bold, while the lowest scores for all model variants are underlined. Arg-C: argument classification; Arg-I: argument identification; SA: situation awareness; SESA: sentence event situation awareness; Trig-C: trigger classification; Trig-I: trigger identification



**Fig. 8 Performance of SESA on ACE05-E (a), ACE05-E+ (b), and ERE-EN (c). SESA: sentence event situation awareness**

the figure that both  $R$  and F1 decrease as  $\mathcal{K}_{r,u}$  or  $t_t$  increases.  $P$  increases with the increment in  $\mathcal{K}_{r,u}$  and  $t_t$ , indicating a positive correlation between the hyperparameters and precision. When  $\mathcal{K}_{r,u} = 1$ ,  $\text{pr}(x, 1) = \delta(x)$  denotes the absence of the reward-penalty mechanism. The improvement in  $P$ ,  $R$ , and F1 scores demonstrates the efficacy of the reward-penalty mechanism.

Fig. 10 illustrates the variation in the Trig-C F1 score on ACE05-E+ and ERE-EN. Experimental results suggest that an increased value of  $\mathcal{K}_{r,u}$  or  $t_t$  does not consistently improve the model performance. Different tasks exhibit the optimal performance under the specific  $\mathcal{K}_{r,u}$  and  $t_t$ , demonstrating the flexibility and superiority of the reward-penalty mechanism.

### 6.2.5 Polysemy test

In this subsection, we conduct experiments to verify the efficiency of RPEE when dealing with polysemous triggers. To assess the effect on the polysemy and monosemy, we divide the test datasets of ACE05-E, ACE05-E+, and ERE-EN into monosemous and polysemous test datasets, based on whether sub-tokens of triggers have multiple semantics. Table 7 displays the performances of RPEE and its variants on the Trig-C F1 score. The performance of RPEE on the polysemous test dataset of ERE-EN outperforms the monosemous one, while for ACE05-E and ACE05-E+, the opposite is true. The reason is that the numbers of monosemous and polysemous test datasets for ERE-EN are nearly equal, whereas for ACE05-E and ACE05-E+, the number of

Table 6 Ablation study of the event decoder for SESA

Model	ACE05-E			ACE05-E+			ERE-EN		
	$P$	$R$	F1	$P$	$R$	F1	$P$	$R$	F1
RPEE	96.47	92.66	94.53	94.86	92.54	93.69	89.36	90.05	89.70
w/o Event Encoder	82.79	80.36	81.56	78.85	77.63	78.24	74.91	79.79	77.27
$\Delta$	16.52	15.31	15.90	20.30	19.21	19.75	19.29	12.86	16.09

$\Delta$  signifies the relative performance gain obtained using pure event knowledge. SESA: sentence event situation awareness

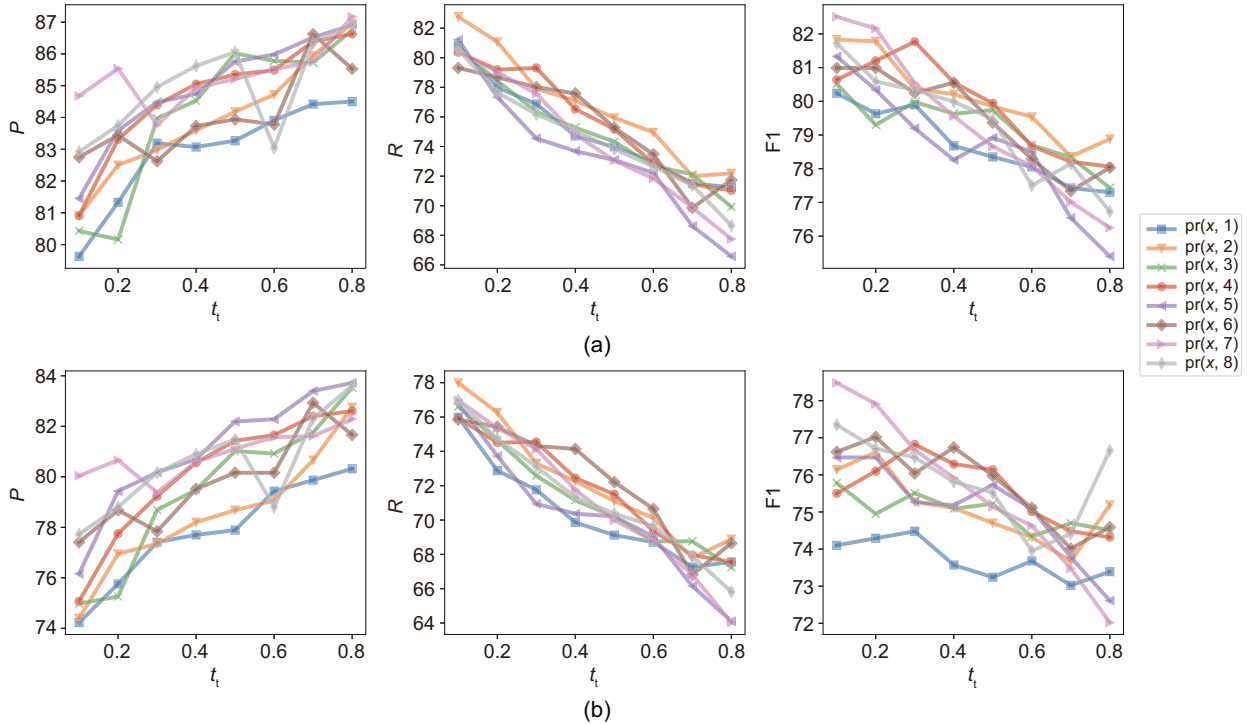


Fig. 9 Performance of Trig-I (a) and Trig-C (b) on ACE05-E with different settings of two hyperparameters ( $\mathcal{K}_{r,u}$  and  $t_t$ ). Trig-I: trigger identification; Trig-C: trigger classification

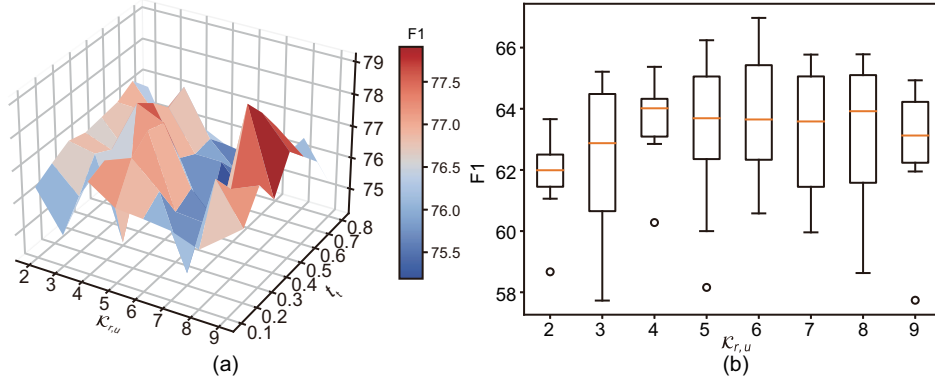


Fig. 10 Trig-C F1 results on ACE05-E+ (a) and ERE-EN (b) with varying  $K_{r,u}$  and  $t_t$

Table 7 Experimental results on the monosemous and polysemous triggers

Model	Dataset	Trig-I			Trig-C		
		<i>P</i>	<i>R</i>	F1	<i>P</i>	<i>R</i>	F1
RPEE_full	ACE05-E	85.67	80.12	82.80	80.84	76.42	78.57
RPEE_monosemous_triggers		87.23	77.04	81.82	83.89	75.03	79.21
RPEE_polysemous_triggers		93.25	81.12	86.77	81.09	73.89	77.32
w/o reward_penalty RPEE_monosemous_triggers		74.40	78.30	76.30	74.21	78.10	76.11
w/o reward_penalty RPEE_polysemous_triggers		80.68	81.43	81.06	68.33	70.51	69.40
RPEE_full		ACE05-E+	85.22	79.18	82.09	81.63	76.79
RPEE_monosemous_triggers	83.20		77.30	80.14	82.09	76.66	79.28
RPEE_polysemous_triggers	88.88		81.08	84.80	78.66	75.40	76.99
w/o reward_penalty RPEE_monosemous_triggers	74.33		82.99	78.42	71.96	82.00	76.65
w/o reward_penalty RPEE_polysemous_triggers	74.70		83.95	79.06	63.49	77.58	69.83
RPEE_full	ERE-EN		75.24	78.14	76.62	63.78	70.51
RPEE_monosemous_triggers		74.45	74.09	62.04	62.04	65.13	63.55
RPEE_polysemous_triggers		74.34	86.30	79.88	60.42	77.25	67.81
w/o reward_penalty RPEE_monosemous_triggers		66.67	66.16	66.42	55.98	57.77	56.86
w/o reward_penalty RPEE_polysemous_triggers		71.24	80.85	75.74	53.19	69.41	60.23

Trig-C: trigger classification; Trig-I: trigger identification

monosemous datasets is 50% larger than that of polysemous datasets. It demonstrates that RPEE can effectively handle triggers with multiple meanings. When using the reward-penalty mechanism, there is a relative improvement of 11.4%, 10.3%, and 12.6% for the polysemous datasets of ACE05-E, ACE05-E+, and ERE-EN on Trig-C F1 score, respectively. It indicates that the reward-penalty mechanism effectively addresses the challenges posed by polysemy. Hence, we can confidently conclude that our approach can handle datasets with polysemous triggers and arguments, showcasing strong robustness.

### 6.2.6 EAE with gold triggers

We conduct comparative experiments using gold triggers on ACE05-E, ACE05-E+, and ERE-EN to

explore the potential of our model. As depicted in Table 8, it achieves relative F1 gains of 4.4%, 6.7%, and 7.7% for Arg-C on ACE05-E, ACE05-E+, and ERE-EN, respectively. It demonstrates that our approach effectively handles EAE tasks, irrespective of using predicted triggers or gold triggers. Additionally, we observe a decrease in the EAE performance when neglecting the reward-penalty mechanism, further affirming the reliability and effectiveness of our designed reward-penalty mechanism.

## 7 Conclusions

In this paper, we propose an adaptive semantic learning strategy to mitigate the bias in the semantic distribution of polysemous triggers and arguments. We design a reward-penalty mechanism to enlarge

**Table 8 Performance of event argument extraction using gold triggers on ACE05-E, ACE05-E+, and ERE-EN**

Model	F1 score					
	ACE05-E		ACE05-E+		ERE-EN	
	Arg-I	Arg-C	Arg-I	Arg-C	Arg-I	Arg-C
RPEE	66.43	58.96	69.13	60.80	69.47	58.57
with gold_triggers	65.98	61.53	72.59	64.86	69.09	63.13
with gold_triggers + w/o reward-penalty	67.50	55.83	72.10	63.89	68.53	62.35

Arg-C: argument classification; Arg-I: argument identification

the gap between the relevant semantics and irrelevant semantics and diminish the gap between relevant semantics by rewarding the correctly classified semantics and punishing the misclassified semantics. The sentence-level event semantics, pre-trained by using a sentence-level event SA mechanism to ensure accuracy, is integrated into token representations to narrow the target event scope of triggers. The model identifies the boundaries of triggers and arguments and classifies their types using task-specific semantic decoders. Our experiments show our model's strengths in robustness, scalability, and generalization ability in complex scenarios. In the future, we will extend our model to low-resource scenarios.

### Contributors

Haili LI designed the research. Haili LI, Yunyan ZHOU, and Jie ZHOU processed the data. Haili LI drafted the paper. Zhiliang TIAN, Yunyan ZHOU, and Qiubo XU helped organize the paper. Haili LI, Zhiliang TIAN, Dongsheng LI, Xiaodong WANG, and Shilong PAN revised and finalized the paper.

### Conflict of interest

Dongsheng LI is a corresponding expert of *Frontiers of Information Technology & Electronic Engineering*, and he was not involved with the peer review process of this paper. All the authors declare that they have no conflict of interest.

### Data availability

The data that support the findings of this study are available from the corresponding authors upon reasonable request.

### References

Ahn D, 2006. The stages of event extraction. Proc Workshop on Annotating and Reasoning about Time and Events, p.1-8. <https://doi.org/10.3115/1629235.1629236>  
 Akgun SA, Ghafurian M, Crowley M, et al., 2023. Using affect as a communication modality to improve human-robot communication in robot-assisted search and res-

cue scenarios. *IEEE Trans Affect Comput*, 14(4):3013-3030. <https://doi.org/10.1109/TAFFC.2022.3221922>  
 Anelli VW, Di Noia T, Di Sciascio E, et al., 2022. Inferring user decision-making processes in recommender systems with knowledge graphs. Proc 30<sup>th</sup> Italian Symp on Advanced Database Systems, p.505-513.  
 Chen YB, Xu LH, Liu K, et al., 2015. Event extraction via dynamic multi-pooling convolutional neural networks. Proc 53<sup>rd</sup> Annual Meeting of the Association for Computational Linguistics and 7<sup>th</sup> Int Joint Conf on Natural Language Processing, p.167-176. <https://doi.org/10.3115/V1/P15-1017>  
 Cong X, Cui SY, Yu BW, et al., 2021. Few-shot event detection with prototypical amortized conditional random field. Findings of the Association for Computational Linguistics, p.28-40. <https://doi.org/10.18653/V1/2021.FINDINGS-ACL.3>  
 Cui SY, Yu BW, Liu TW, et al., 2020. Edge-enhanced graph convolution networks for event detection with syntactic relation. Findings of the Association for Computational Linguistics, p.2329-2339. <https://doi.org/10.18653/v1/2020.findings-emnlp.211>  
 Cui ZJ, Yuan ZM, Wu YF, et al., 2023. Intelligent recommendation for departments based on medical knowledge graph. *IEEE Access*, 11:25372-25385. <https://doi.org/10.1109/ACCESS.2023.3254303>  
 Devlin J, Chang MW, Lee K, et al., 2019. BERT: pre-training of deep bidirectional Transformers for language understanding. Proc Conf of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, p.4171-4186. <https://doi.org/10.18653/v1/n19-1423>  
 Ding N, Li ZR, Liu ZY, et al., 2019. Event detection with trigger-aware lattice neural network. Proc Conf on Empirical Methods in Natural Language Processing and 9<sup>th</sup> Int Joint Conf on Natural Language Processing, p.347-356. <https://doi.org/10.18653/v1/D19-1033>  
 Dozat T, Manning CD, 2017. Deep biaffine attention for neural dependency parsing. Proc 5<sup>th</sup> Int Conf on Learning Representations, p.1-8.  
 Du XY, Cardie C, 2020. Event extraction by answering (almost) natural questions. Proc Conf on Empirical Methods in Natural Language Processing, p.671-683. <https://doi.org/10.18653/v1/2020.emnlp-main.49>  
 Du XY, Ji H, 2022. Retrieval-augmented generative question answering for event argument extraction. Proc Conf on Empirical Methods in Natural Language Processing, p.4649-4666. <https://doi.org/10.18653/v1/2022.emnlp-main.307>

- Dwivedi D, Yemula PK, Pal M, 2023. DynamoPMU: a physics informed anomaly detection, clustering, and prediction method using nonlinear dynamics on  $\mu$  PMU measurements. *IEEE Trans Instrum Meas*, 72:1-9. <https://doi.org/10.1109/TIM.2023.3327481>
- Endsley MR, 1988. Design and evaluation for situation awareness enhancement. *Proc Hum Factors Ergon Soc Ann Meet*, 32(2):97-101. <https://doi.org/10.1177/154193128803200221>
- Endsley MR, 2001. Designing for situation awareness in complex systems. *Proc 2<sup>nd</sup> Int Workshop on Symbiosis of Humans, Artifacts and Environment*, p.1-14.
- Ettinger A, Hwang J, Pyatkin V, et al., 2023. "You are an expert linguistic annotator": limits of LLMs as analyzers of abstract meaning representation. *Findings of the Association for Computational Linguistics*, p.8250-8263. <https://doi.org/10.18653/v1/2023.findings-emnlp.553>
- Feng XC, Qin B, Liu T, 2018. A language-independent neural network for event detection. *Sci China Inform Sci*, 61(9):092106. <https://doi.org/10.1007/S11432-017-9359-X>
- Gururangan S, Marasović A, Swayamdipta S, et al., 2020. Don't stop pretraining: adapt language models to domains and tasks. *Proc 58<sup>th</sup> Annual Meeting of the Association for Computational Linguistics*, p.8342-8360. <https://doi.org/10.18653/v1/2020.acl-main.740>
- Guzman-Nateras L, Deroncourt F, Nguyen T, 2023. Hybrid knowledge transfer for improved cross-lingual event detection via hierarchical sample selection. *Proc 61<sup>st</sup> Annual Meeting of the Association for Computational Linguistics*, p.5414-5427. <https://doi.org/10.18653/v1/2023.acl-long.296>
- Han RD, Yang CH, Peng T, et al., 2023. An empirical study on information extraction using large language models. <https://doi.org/10.48550/arXiv.2305.14450>
- He YX, Hu JY, Tang BZ, 2023. Revisiting event argument extraction: can EAE models learn better when being aware of event co-occurrences? *Proc 61<sup>st</sup> Annual Meeting of the Association for Computational Linguistics*, p.12542-12556. <https://doi.org/10.18653/v1/2023.acl-long.701>
- Hong Y, Zhang JF, Ma B, et al., 2011. Using cross-entity inference to improve event extraction. *Proc 49<sup>th</sup> Annual Meeting of the Association for Computational Linguistics*, p.1127-1136.
- Hsu IH, Huang KH, Boschee E, et al., 2022. DEGREE: a data-efficient generation-based event extraction model. *Proc Conf of the North American Chapter of the Association for Computational Linguistics*, p.1890-1908. <https://doi.org/10.18653/v1/2022.naacl-main.138>
- Hsu IH, Xie ZY, Huang KH, et al., 2023. AMPERE: AMR-aware prefix for generation-based event argument extraction model. *Proc 61<sup>st</sup> Annual Meeting of the Association for Computational Linguistics*, p.10976-10993. <https://doi.org/10.18653/v1/2023.acl-long.615>
- Ji H, Grishman R, 2008. Refining event extraction through cross-document inference. *Association for Computational Linguistics Annual Meeting: Human Language Technologies*, p.254-262.
- Lewis M, Liu YH, Goyal N, et al., 2020. BART: denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. *Proc 58<sup>th</sup> Annual Meeting of the Association for Computational Linguistics*, p.7871-7880. <https://doi.org/10.18653/v1/2020.ACL-MAIN.703>
- Li S, Zhao RN, Li ML, et al., 2023. Open-domain hierarchical event schema induction by incremental prompting and verification. *Proc 61<sup>st</sup> Annual Meeting of the Association for Computational Linguistics*, p.5677-5697. <https://doi.org/10.18653/v1/2023.acl-long.312>
- Lin Y, Ji H, Huang F, et al., 2020. A joint neural model for information extraction with global features. *Proc 58<sup>th</sup> Annual Meeting of the Association for Computational Linguistics*, p.7999-8009. <https://doi.org/10.18653/v1/2020.acl-main.713>
- Liu J, Chen YF, Xu JA, 2021. Machine reading comprehension as data augmentation: a case study on implicit event argument extraction. *Proc Conf on Empirical Methods in Natural Language Processing*, p.2716-2725. <https://doi.org/10.18653/v1/2021.emnlp-main.214>
- Liu J, Chen YF, Xu JA, 2022. Saliency as evidence: event detection with trigger saliency attribution. *Proc 60<sup>th</sup> Annual Meeting of the Association for Computational Linguistics*, p.4573-4585. <https://doi.org/10.18653/v1/2022.acl-long.313>
- Liu J, Sui DB, Liu K, et al., 2023. Learning with partial annotations for event detection. *Proc 61<sup>st</sup> Annual Meeting of the Association for Computational Linguistics*, p.508-523. <https://doi.org/10.18653/v1/2023.acl-long.30>
- Liu X, Luo ZC, Huang HY, 2018. Jointly multiple events extraction via attention-based graph information aggregation. *Proc Conf on Empirical Methods in Natural Language Processing*, p.1247-1256. <https://doi.org/10.18653/v1/d18-1156>
- Liu X, Huang HY, Shi G, et al., 2022. Dynamic prefix-tuning for generative template-based event extraction. *Proc 60<sup>th</sup> Annual Meeting of the Association for Computational Linguistics*, p.5216-5228. <https://doi.org/10.18653/v1/2022.acl-long.358>
- Loshchilov I, Hutter F, 2019. Decoupled weight decay regularization. *Proc 7<sup>th</sup> Int Conf on Learning Representations*, p.1-8.
- Lou DF, Liao ZL, Deng SM, et al., 2021. MLBiNet: a cross-sentence collective event detection network. *Proc 59<sup>th</sup> Annual Meeting of the Association for Computational Linguistics and 11<sup>th</sup> Int Joint Conf on Natural Language Processing*, p.4829-4839. <https://doi.org/10.18653/v1/2021.acl-long.373>
- Lu D, Ran SH, Tetreault JR, et al., 2023. Event extraction as question generation and answering. *Proc 61<sup>st</sup> Annual Meeting of the Association for Computational Linguistics*, p.1666-1688. <https://doi.org/10.18653/v1/2023.acl-short.143>
- Lu YJ, Lin HY, Xu J, et al., 2021. Text2Event: controllable sequence-to-structure generation for end-to-end event extraction. *Proc 59<sup>th</sup> Annual Meeting of the Association for Computational Linguistics and 11<sup>th</sup> Int Joint Conf on Natural Language Processing*, p.2795-2806. <https://doi.org/10.18653/v1/2021.acl-long.217>
- Ma MD, Taylor A, Wang W, et al., 2023. DICE: data-efficient clinical event extraction with generative models. *Proc 61<sup>st</sup> Annual Meeting of the Association for Computational Linguistics*, p.15898-15917. <https://doi.org/10.18653/v1/2023.acl-long.886>

- Matey AH, Danquah P, Koi-Akrofi GY, 2022. Predicting cyber-attack using cyber situational awareness: the case of independent power producers (IPPs). *Int J Adv Comput Sci Appl*, 13(1):700-709. <https://doi.org/10.14569/IJACSA.2022.0130181>
- McClosky D, Surdeanu M, Manning CD, 2011. Event extraction as dependency parsing. Proc 49<sup>th</sup> Annual Meeting of the Association for Computational Linguistics: Human Language Technologies, p.1626-1635. <https://doi.org/10.5555/2002472.2002667>
- Nam H, Kim SH, Park YH, 2022. FilterAugment: an acoustic environmental data augmentation method. Proc IEEE Int Conf on Acoustics, Speech and Signal Processing, p.4308-4312. <https://doi.org/10.1109/ICASSP43922.2022.9747680>
- Nguyen TH, Cho K, Grishman R, 2016. Joint event extraction via recurrent neural networks. Proc Conf of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, p.300-309. <https://doi.org/10.18653/v1/n16-1034>
- Onwubiko C, 2020. CyberOps: situational awareness in cybersecurity operations. *Int J Cyber Situat Aware*, 5(1):82-107. <https://doi.org/10.22619/ijcsa.2020.100134>
- Palash M, Bhargava B, 2023. SAFER: situation aware facial emotion recognition. <https://doi.org/10.48550/ARXIV.2306.09372>
- Pang CX, Cao YX, Ding Q, et al., 2023. Guideline learning for in-context information extraction. Proc Conf on Empirical Methods in Natural Language Processing, p.15372-15389. <https://doi.org/10.18653/v1/2023.emnlp-main.950>
- Ping Y, Lu JY, Gan RY, et al., 2023. UniEX: an effective and efficient framework for unified information extraction via a span-extractive perspective. Proc 61<sup>st</sup> Annual Meeting of the Association for Computational Linguistics, p.16424-16440. <https://doi.org/10.18653/v1/2023.acl-long.907>
- Sha L, Qian F, Chang BB, et al., 2018. Jointly extracting event triggers and arguments by dependency-bridge RNN and tensor-based argument interaction. Proc 32<sup>nd</sup> AAAI Conf on Artificial Intelligence, p.5916-5923. <https://doi.org/10.1609/aaai.v32i1.12034>
- Shashikumar SP, Josef CS, Sharma A, et al., 2021. DeepAISE—an interpretable and recurrent neural survival model for early prediction of sepsis. *Artif Intell Med*, 113:102036. <https://doi.org/10.1016/J.ARTMED.2021.102036>
- Shu XF, Yan J, Gao WR, et al., 2021. Research on military equipment entity recognition and knowledge graph construction method based on ALBERT-Bi-LSTM-CRF. Proc 4<sup>th</sup> Int Conf on Artificial Intelligence and Pattern Recognition, p.273-279. <https://doi.org/10.1145/3488933.3489030>
- Song ZY, Bies A, Strassel S, et al., 2015. From light to rich ERE: annotation of entities, relations, and events. Proc 3<sup>rd</sup> Workshop on EVENTS: Definition, Detection, Coreference, and Representation, p.89-98. <https://doi.org/10.3115/V1/W15-0812>
- Wadden D, Wennberg U, Luan Y, et al., 2019. Entity, relation, and event extraction with contextualized span representations. Proc Conf on Empirical Methods in Natural Language Processing and 9<sup>th</sup> Int Joint Conf on Natural Language Processing, p.5783-5788. <https://doi.org/10.18653/V1/D19-1585>
- Wang B, Huang HY, Wei XC, et al., 2023. Boosting event extraction with denoised structure-to-text augmentation. Findings of the Association for Computational Linguistics, p.11267-11281. <https://doi.org/10.18653/v1/2023.findings-acl.716>
- Wang S, Yu M, Chang S, et al., 2022. Query and extract: refining event extraction as type-oriented binary decoding. Findings of the Association for Computational Linguistics, p.169-182. <https://doi.org/10.18653/v1/2022.findings-acl.16>
- Wang SJ, Yu M, Huang LF, 2023. The art of prompting: event detection based on type specific prompts. Proc 61<sup>st</sup> Annual Meeting of the Association for Computational Linguistics, p.1286-1299. <https://doi.org/10.18653/v1/2023.acl-short.111>
- Wang ZT, Wang XY, Hu W, 2023. Continual event extraction with semantic confusion rectification. Proc Conf on Empirical Methods in Natural Language Processing, p.11945-11955. <https://doi.org/10.18653/v1/2023.emnlp-main.732>
- Xia LQ, Liang YS, Leng JW, et al., 2023. Maintenance planning recommendation of complex industrial equipment based on knowledge graph and graph neural network. *Reliab Eng Syst Saf*, 232:109068. <https://doi.org/10.1016/J.RESS.2022.109068>
- Xu RX, Wang PY, Liu TY, et al., 2022. A two-stream AMR-enhanced model for document-level event argument extraction. Proc Conf of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, p.5025-5036. <https://doi.org/10.18653/v1/2022.naacl-main.370>
- Xu TY, Guo C, Du LX, et al., 2022. A method for traditional Chinese medicine knowledge graph dynamic construction. Proc 5<sup>th</sup> Int Conf on Big Data Technologies, p.196-202. <https://doi.org/10.1145/3565291.3565323>
- Xu ZY, Lee JY, Huang LF, 2023. Learning from a friend: improving event extraction via self-training with feedback from abstract meaning representation. Findings of the Association for Computational Linguistics, p.10421-10437. <https://doi.org/10.18653/v1/2023.findings-acl.662>
- Xu ZZ, Liu RK, Yang S, et al., 2023. Learning imbalanced data with vision transformers. Proc IEEE/CVF Conf on Computer Vision and Pattern Recognition, p.15793-15803. <https://doi.org/10.1109/CVPR52729.2023.01516>
- Yang P, Cong X, Sun ZY, et al., 2021. Enhanced language representation with label knowledge for span extraction. Proc Conf on Empirical Methods in Natural Language Processing, p.4623-4635. <https://doi.org/10.18653/v1/2021.emnlp-main.379>
- Yang S, Feng DW, Qiao LB, et al., 2019. Exploring pre-trained language models for event extraction and generation. Proc 57<sup>th</sup> Annual Meeting of the Association for Computational Linguistics, p.5284-5294. <https://doi.org/10.18653/v1/p19-1522>
- Yang XJ, Lu YJ, Petzold LR, 2023. Few-shot document-level event argument extraction. Proc 61<sup>st</sup> Annual

- Meeting of the Association for Computational Linguistics, p.8029-8046.  
<https://doi.org/10.18653/v1/2023.acl-long.446>
- Yang YQ, Guo QP, Hu XK, et al., 2023. An AMR-based link prediction approach for document-level event argument extraction. Proc 61<sup>st</sup> Annual Meeting of the Association for Computational Linguistics, p.12876-12889.  
<https://doi.org/10.18653/v1/2023.acl-long.720>
- Yao YZ, Mao SY, Zhang NY, et al., 2023. Schema-aware reference as prompt improves data-efficient knowledge graph construction. Proc 46<sup>th</sup> Int ACM SIGIR Conf on Research and Development in Information Retrieval, p.911-921. <https://doi.org/10.1145/3539618.3591763>
- You HL, Samuel D, Touileb S, et al., 2022. EventGraph: event extraction as semantic graph parsing. Proc 5<sup>th</sup> Workshop on Challenges and Applications of Automated Extraction of Socio-Political Events from Text, p.7-15. <https://doi.org/10.18653/v1/2022.case-1.2>
- You HL, Touileb S, Øvreid L, 2023. JSEEGraph: joint structured event extraction as graph parsing. Proc 12<sup>th</sup> Joint Conf on Lexical and Computational Semantics, p.115-127.  
<https://doi.org/10.18653/v1/2023.starsem-1.11>
- You MS, Yin J, Wang H, et al., 2023. A knowledge graph empowered online learning framework for access control decision-making. *World Wide Web*, 26(2):827-848.  
<https://doi.org/10.1007/S11280-022-01076-5>
- Zeng Y, Yang HH, Feng YS, et al., 2016. A convolution BiLSTM neural network model for Chinese event extraction. Proc 5<sup>th</sup> CCF Conf on Natural Language Processing and Chinese Computing, 10102:275-287.  
[https://doi.org/10.1007/978-3-319-50496-4\\_23](https://doi.org/10.1007/978-3-319-50496-4_23)
- Zhang JR, Ilievski F, Ma KX, et al., 2023. A study of situational reasoning for traffic understanding. Proc 29<sup>th</sup> ACM SIGKDD Conf on Knowledge Discovery and Data Mining, p.3262-3272.  
<https://doi.org/10.1145/3580305.3599246>
- Zhang KH, Shuang K, Yang XY, et al., 2023. What is overlap knowledge in event argument extraction? APE: a cross-datasets transfer learning model for EAE. Proc 61<sup>st</sup> Annual Meeting of the Association for Computational Linguistics, p.393-409.  
<https://doi.org/10.18653/v1/2023.acl-long.24>
- Zhao G, Gong XC, Yang XJ, et al., 2023. DemoSG: demonstration-enhanced schema-guided generation for low-resource event extraction. Findings of the Association for Computational Linguistics, p.1805-1816.  
<https://doi.org/10.18653/v1/2023.findings-emnlp.121>
- Zheng YZ, Pan SR, Lee VCS, et al., 2022. Rethinking and scaling up graph contrastive learning: an extremely efficient approach with group discrimination. Conf on Neural Information Processing Systems, p.1-17.

## List of supplementary materials

1 Case study

2 Discussions

Table S1 Case analysis of different types of semantic error classification for polysemous triggers