

# RAG-based Crowdsourcing Task Decomposition via Masked Contrastive Learning with Prompts

Jing Yang, Xiao Wang, *Senior Member, IEEE*, Yu Zhao, Yuhang Liu, Fei-Yue Wang, *Fellow, IEEE*

**Abstract**—Crowdsourcing is a critical technology in social manufacturing, which leverages an extensive and boundless reservoir of human resources to handle a wide array of complex tasks. The successful execution of these complex tasks relies on task decomposition (TD) and allocation, with the former being a prerequisite for the latter. Recently, pre-trained language models (PLMs)-based methods have garnered significant attention. However, they are constrained to handling straightforward common-sense tasks due to their inherent restrictions involving limited and difficult-to-update knowledge as well as the presence of “hallucinations”. To address these issues, we propose a retrieval-augmented generation-based crowdsourcing framework that reimagines TD as event detection from the perspective of natural language understanding. However, the existing detection methods fail to distinguish differences between event types and always depend on heuristic rules and external semantic analyzing tools. Therefore, we present a Prompt-Based Contrastive learning framework for TD (PBCT), which incorporates a prompt-based trigger detector to overcome dependence. Additionally, trigger-attentive sentinel and masked contrastive learning are introduced to provide varying attention to trigger and contextual features according to different event types. Experiment results demonstrate the competitiveness of our method in both supervised and zero-shot detection. A case study on printed circuit board manufacturing is showcased to validate its adaptability to unknown professional domains.

**Index Terms**—Crowdsourcing, Retrieval-augmented Generation, Task Decomposition, Event Detection, Pre-trained Language Models

## I. INTRODUCTION

Against the backdrop of increasingly diverse and personalized consumer demands, customized production is gaining significant attention and adoption as a pivotal production mode among numerous enterprises [1], [2]. It not only caters to specific consumer needs and amplifies product differentiation for a competitive edge, but also curtails inventory wastage, enhances production flexibility, and fortifies customer loyalty [3], [4].

Jing Yang and Yuhang Liu is with The State Key Laboratory for Management and Control of Complex Systems, Institute of Automation, Chinese Academy of Sciences, Beijing 100190, China, and also with the School of Artificial Intelligence, University of Chinese Academy of Sciences, Beijing 100049, China (E-mail: yangjing2020@ia.ac.cn, liuyuhang2021@ia.ac.cn).

Xiao Wang is with the School of Artificial Intelligence, Anhui University, Hefei 266114, China and also with Qingdao Academy of Intelligent Industries, Qingdao 230031, China (E-mail: xiao.wang@ahu.edu.cn).

Yu Zhao is with National Key Laboratory of Information Systems Engineering, National University of Defense Technology (E-mail: yuzhao@nudt.edu.cn).

Fei-Yue Wang is with Macau University of Science and Technology, Macao 999078, China, with Beijing Engineering Research Center of Intelligent Systems and Technology, Chinese Academy of Sciences, Beijing 100098, China and also with State Key Laboratory for Management and Control of Complex Systems, Chinese Academy of Sciences, Beijing 100190, China (e-mail: feiyue.wang@ia.ac.cn).

Crowdsourcing, epitomizing collective intelligence, as shown in Fig. 1, empowers businesses to gain profound insights into consumers’ needs and preferences. Moreover, it even allows consumers to participate in the production process, leveraging an extensive, boundless reservoir of human resources to handle a wide array of tasks with varying levels of complexity [5]. Therefore, it is crucial and essential to effectively apply crowdsourcing technology to customized production and thus promote the shift towards social manufacturing [6], [7].

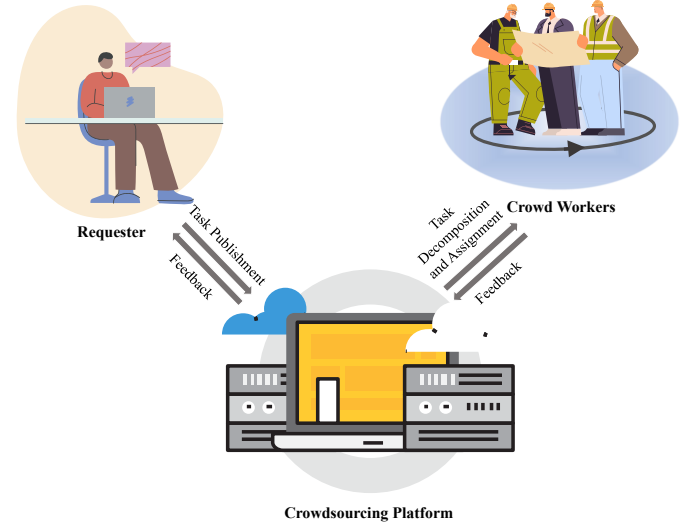


Fig. 1: The workflow of crowdsourcing.

Requesters and crowd workers are organized and coordinated through a crowdsourcing platform, which acts as a mediator between them [8]. Once a task is posted by a requester, the crowdsourcing platform allocates the task to a group of proficient crowd workers, with the goal of efficiently finishing the task within the shortest possible timeframe. Obviously, success depends on the ability to decompose complex tasks into smaller subtasks [9], as illustrated in Fig. 2, each of which is low in complexity and requires minimal cognitive effort to be completed by an individual. These subtasks often feature interdependencies over time, and when orchestrated sequentially, they culminate in an event sequence delineating the complex task. However, prevalent TD methods often focus on specific scenarios, relying on prior knowledge of the complex task to design appropriate rules or optimization algorithms for its breakdown [10]–[13]. These approaches have significant limitations as they may not be easily adaptable to other scenarios and rely heavily on human expertise. Consequently, deploying them on crowdsourcing platforms for decomposing

a variety of tasks proves challenging.

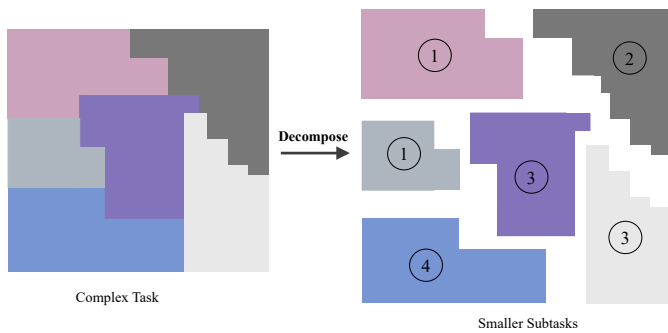


Fig. 2: The illustration of Task decomposition. The numbers represents the order in which subtasks are carried out.

Recently, pre-trained language models (PLMs) like BERT [14], ChatGPT [15], and LLaMA [16] have demonstrated the ability to store substantial world knowledge similar to large knowledge graphs during their pre-training phase. These models have showcased exceptional performance across a variety of tasks, such as event extraction [17], temporal relation extraction [18] and personalized recommendation [19]. Given that requesters’ tasks are typically described in natural language, a straightforward idea is to treat task decomposition as a natural language understanding and information extraction task using PLMs. Several attempts have been made to stimulate the inferential capabilities of PLMs and activate their stored knowledge via prompts for TD [20]–[22]. However, the current PLMs-based TD is primarily constrained to handling common sense tasks due to their limited and hard-to-update knowledge. Another prevalent challenge involves the generation of inaccurate information, known as “hallucinations”, particularly when queries surpass the PLM’s training data.

One promising approach to address these issues is to integrate retrieval-augmented generation (RAG) [23] into the crowdsourcing systems, which incorporates retrieved external data into the TD process, thereby improving the model’s capability to accurately decompose diverse tasks. Consequently, based on the pertinent retrieved data for the task’s execution steps, we extract or detect the triggers and related arguments of subtasks (which can be also considered as events) for each step. It is evident that TD tasks can be converted into event extraction or detection tasks, an area where significant effort has already been invested. However, the existing detection methods might have only learned lexical patterns [24], i.e., word-to-trigger mapping, resulting in skewed performance as they fail to distinguish differences between event types. Additionally, with the constant emergence of new tasks or events without annotated samples in the open world, it is imperative and necessary to deal with zero-shot detection tasks. Therefore, the main contributions of this paper are as follows.

- We propose an RAG-based crowdsourcing framework that addresses challenges related to hallucinations and limited knowledge through external retrieval. By this means, TD tasks are transformed into event detection

tasks, contributing to the utilization of mature detection techniques.

- We present a prompt-based contrastive learning framework for TD, named PBCT, which incorporates prompt learning for trigger detection. It breaks away from dependency on heuristic rules and external semantic analyzing tools.
- We introduce a trigger-attentive sentinel and masked contrastive learning to provide varying attention to trigger features and contextual features according to different event types. This method brings similar samples closer and pushes away samples of different classes, thereby improving the representation learning of unseen types under a zero-shot setting.
- We conduct a series of experiments on two datasets: ACE 2005 and FewEvent, and the results demonstrate that our method achieves competitive performance in both supervised and zero-shot detection. Additionally, we offer a case study on personalized printed circuit board (PCB) manufacturing to validate the effectiveness of our approach for TD in RAG-based crowdsourcing.

The subsequent sections of this paper are structured as follows. Section II presents the RAG-based crowdsourcing framework and its operational workflow, followed by a review of related work for TD in Section III. Section IV introduces our problem statement and relevant notations. In Section V, a comprehensive overview of our model design is presented, along with a detailed explanation of each component. Section VI gives and discusses a series of experimental results. Finally, concluding remarks and future work are drawn in Section VII.

## II. RAG-BASED CROWDSOURCING

Crowdsourcing is a goal-directed and self-organizing collective behavior. The interaction between requesters and crowd workers is crucial for the efficiency and quality of crowdsourced task completion. Therefore, it is important to accurately grasp the requirements of requesters and convert them into a series of subtasks allocated to crowd workers. The emergence of PLMs and RAG makes it possible to automate their interactions. To enhance PLMs’ inference performance and equip them with the capability to handle complex, diverse and open-world tasks, we reformulate all computational tasks, such as TD, temporal relation extraction, and personalized recommendation, within crowdsourcing systems into prompt engineering tasks via RAG modes. Fig. 3 illustrates the framework of RAG-based crowdsourcing, which interfaces with an external knowledge repository to augment and rectify the knowledge stored within PLMs, especially professional and up-to-date knowledge. During the preparatory phase, the corpus in the repository is segmented into discrete chunks, which are embedded into vector representations as their indices through an encoder.

In terms of the task-processing phase, we take task decomposition as an example in Fig. 3 to explain the operational workflow of the framework. After a requester initiates a query, the same encoder as the preparatory phase is utilized to transform the input into a vector representation. Subsequently,

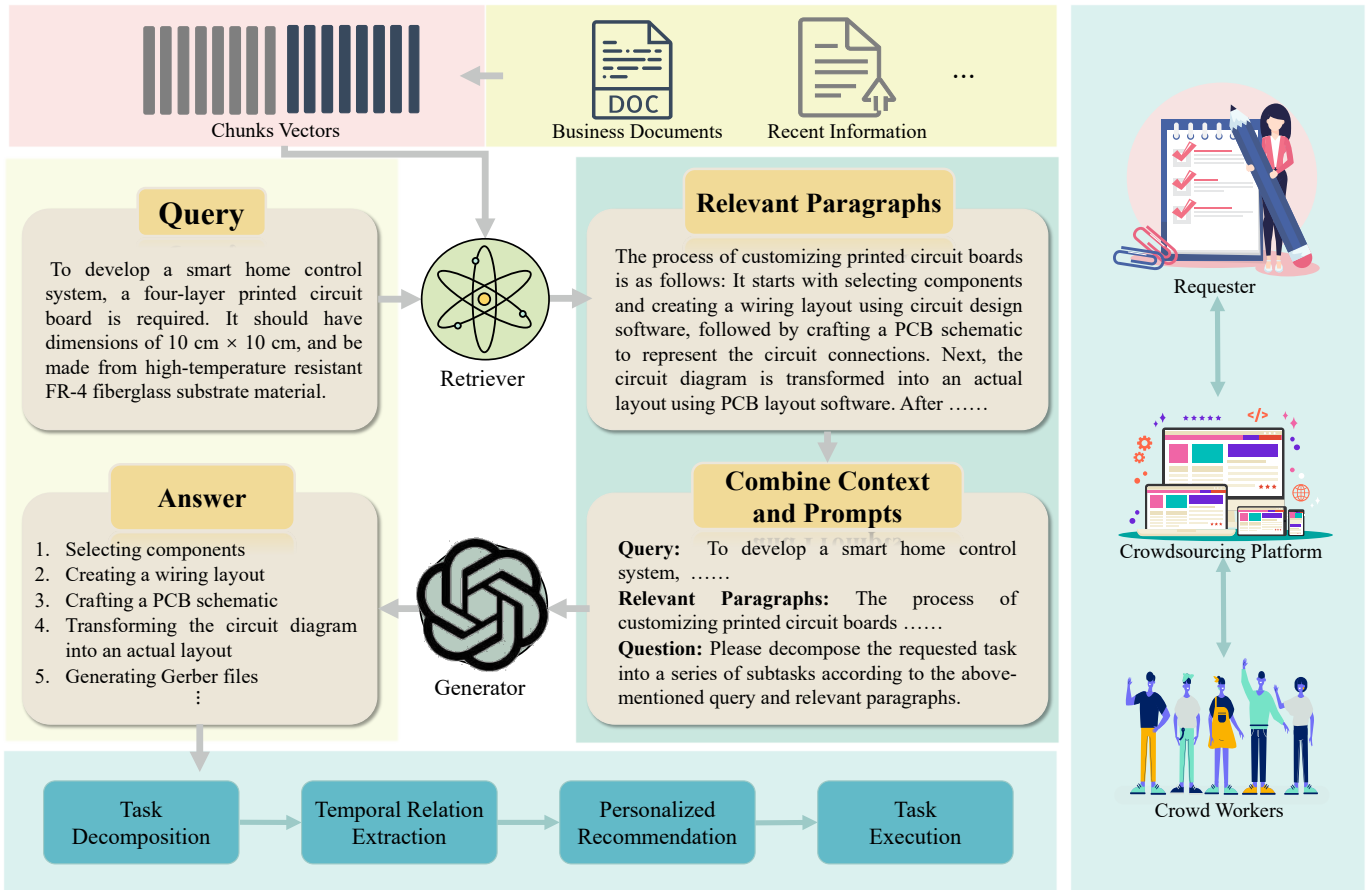


Fig. 3: RAG-based crowdsourcing

the similarity between the query vector and index vector is computed to obtain the most relevant chunks via a retriever. Finally, the query, along with the relevant paragraphs retrieved and a task-specific prompt, are combined as input for the generator to produce answers. Based on the TD results, we sequentially extract the temporal dependencies of subtasks and recommend corresponding subtasks to the competent crowd workers, following the same procedure. Once the tasks are assigned, the crowd workers would independently or collaboratively strive to accomplish the tasks, driven by incentive mechanisms [25], [26].

In this paper, our focus is on enhancing the performance of the generator used to identify a series of subtasks based on a pertinent paragraph describing a complex task. Therefore, we assume that the retriever can access critical texts regarding the steps and measures relevant to task execution, which is expected to be the focal point of our future research.

### III. TASK DECOMPOSITION

TD is the first step in crowdsourcing task automation, for which many efforts have been made. Some approaches necessitate designers to manually break down intricate tasks into simpler executable ones [27]–[29]. Especially, to reduce the barrier to entry for workflow design, Turkomatic [30] is introduced as a tool that engages crowd workers to adaptively decompose and solve tasks. Besides, there are semi-automated

modes as the prevailing trend. They are often task-specific [31] or object-specific [13], [32], [33], heavily relying on prior knowledge of task structures and subtask repositories. Their focus is to design TD rules and evaluation metrics tailored to specific scenarios, and then refine various TD strategies through intelligent optimization algorithms. Tong et al. [34] propose an optimal priority queue-based approach to decompose large-scale crowdsourcing tasks that consist of thousands or millions of independent atomic subtasks in both homogeneous and heterogeneous scenarios. However, this method struggles to address complex and interdependent tasks. Motes et al. [32] concentrate on multi-robot transportation tasks and integrate decomposition, allocation, and planning for these multiple decomposable tasks to seek the optimal solution through technologies like conflict-based search. Zhang et al. [35] leverage an integrated numerical design structure matrix and an adaptive genetic algorithm to divide product development tasks into different groups based on their correlation degrees. For complex control tasks, researchers consider they can be solved by multiple sub-polices corresponding to subtasks, so learning from demonstrations is utilized to build sequential decision processes [13], [36]. In software crowdsourcing, a general-purpose approach is presented to verify the complex task’s decomposition scheme and its functional specification [37]. Song et al. employ a design structure matrix and a product structure tree to model the interaction of the design

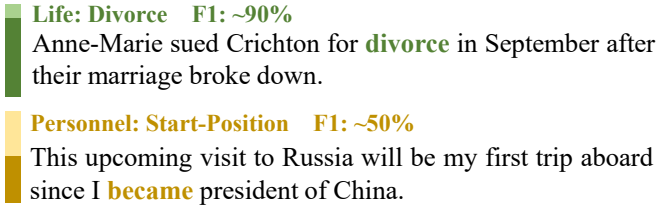


Fig. 4: Two typical instances taken from the ACE 2005 benchmark.

process and thus decompose the design task for complex products [38]. While these methods to a certain extent achieve automation in TD, none of them are generalizable across different task characteristics.

Recently, the incorporation of PLMs into TD has garnered some attention. Sakaguchi et al. [21] finetune PLMs to generate high-quality scripts including partial ordering events through two complementary tasks: edge prediction and script generation. Jansen et al. [20] model visual semantic planning as a translation problem of converting natural language directives into detailed multi-step sequences of actions and finetune PLMs to generate a command sequence via prompts. SmartLLM [22] employs Pythonic chain-of-thought prompts and demonstrations to guide the large PLMs in generating code for multi-robot task planning including TD within the context of smart homes. However, the capabilities of these methods are constrained within addressing prototype tasks in everyday scenarios that demand common-sense knowledge. It is non-trivial to the application of PLMs in dealing with more complex tasks, especially professional tasks. Additionally, it is challenging to enhance their decomposition performance due to hallucination issues. This paper aims at providing a unified solution for decomposing various types of tasks by introducing RAG to model TD as event extraction or detection tasks.

Event detection is a relatively well-established field that focuses on identifying specific types of events (also known as different subtasks in TD), within unstructured text. Extensive research has been conducted in both supervised and few-shot/zero-shot settings. Despite considerable advancements in supervised training [39]–[44], few-shot/zero-shot detection poses ongoing difficulties owing to the limited availability of prior knowledge and annotated examples, a common issue in the open world. Under zero-shot settings, most existing works rely on additional semantic analysis tools to supplement event knowledge or design heuristic rules to facilitate the detection of unseen events [45]–[49]. Zhang et al. [48] compute the similarity between triggers and arguments to the type’s names to classify events with the help of external tools (e.g., Semantic Role Labeling) and corpora (e.g., New York Times). Huang and Ji [46] propose a semi-supervised vector quantized variational autoencoder framework that discovers candidate triggers with a heuristic approach and projects them into a particular type. However, these methods fail to recognize the distinction between event types and train a model that treats all categories equally, leading to imbalanced performance. Observations [24], [50] indicate that performance tends to be relatively lower on context-dependent texts, as models may

primarily learn lexical patterns. For instance, as shown in Fig.4, while DYGIE++ [40] achieves a high F1 score of 90% for the *Life: Divorce* type, its performance drops to 50% for the *Personnel: Start-Position* type. What’s more surprising is that the training set for *Personnel: Start-Position* is eight times larger than that of *Life: Divorce*. This suggests that, in the case of context-dependent texts, the implied semantics of triggers may be inadequate, resulting in an inability to classify events using lexical patterns. In response to this challenge, our approach focuses on adjusting attention between trigger features and contextual features by employing trigger-attentive sentinel and masked contrastive learning techniques.

#### IV. PROBLEM STATEMENT

Given an event mention as an input  $S = \{s_0, s_1, \dots, s_N\}$ , including one or more trigger words annotated with event types  $y \in \mathcal{T}$ , where  $s_i$  is the  $i$ -th token and  $N$  is the number of tokens. A pre-defined set of event types  $\mathcal{T}$  consists of seen event types  $\mathcal{T}_s$  and unseen event types  $\mathcal{T}_u$ , that is,  $\mathcal{T} = \mathcal{T}_s \cup \mathcal{T}_u$ ,  $\mathcal{T}_s \cap \mathcal{T}_u = \emptyset$ . During the training phase, unseen event types are not used as sample labels to optimize the model. Additionally, the triggers of these samples remain invisible to the model. The goal of this task is to identify the triggers and subsequently classify events into appropriate types based on them and the corresponding sentences.

#### V. MODEL ARCHITECTURE

In this section, we elaborate on our proposed PBCT model for TD from the perspective of event detection, with Fig. 5 illustrating its overall architecture, comprising four modules: a prompt-based trigger detector, a prototypical network classifier, a trigger-attentive sentinel, and masked contrastive learning. The prompt-based trigger detector leverages prompt learning to fully exploit the inference capability of PLMs for trigger identification, eliminating the need for additional semantic analysis tools or heuristic rule design. The prototypical network classifier trains a prototypical matrix for event types, which are initialized by label semantics and thus classify events through similarity calculation. The trigger-attentive sentinel balances attention between the trigger and contextual features by learning two weights. Masked contrastive learning is employed to enhance event representation learning, achieved by constructing two positive samples and one negative sample, along with their corresponding losses. Additionally, to take great advantage of all available information, it is also indispensable to incorporate supervised learning in the detection of seen events.

##### A. Prompt-based Trigger Detector

As the first step of event detection, a prompt-based trigger detector is introduced to acquire the contextual vectors of triggers. We construct a prompt template  $PT = \text{“This is an event about [MASK]. } \langle event\ mention \rangle \text{”}$  to input a PLM for trigger prediction, where [MASK] refers to a trigger word to be predicted and  $\langle event\ mention \rangle$  represents the description of an event. Before input,  $PT$  is transformed into the input sequence  $S_p = \{s_{p_0}, s_{p_1}, \dots, s_{p_m}, \dots, s_{p_{l-1}}, s_{p_l}\}$

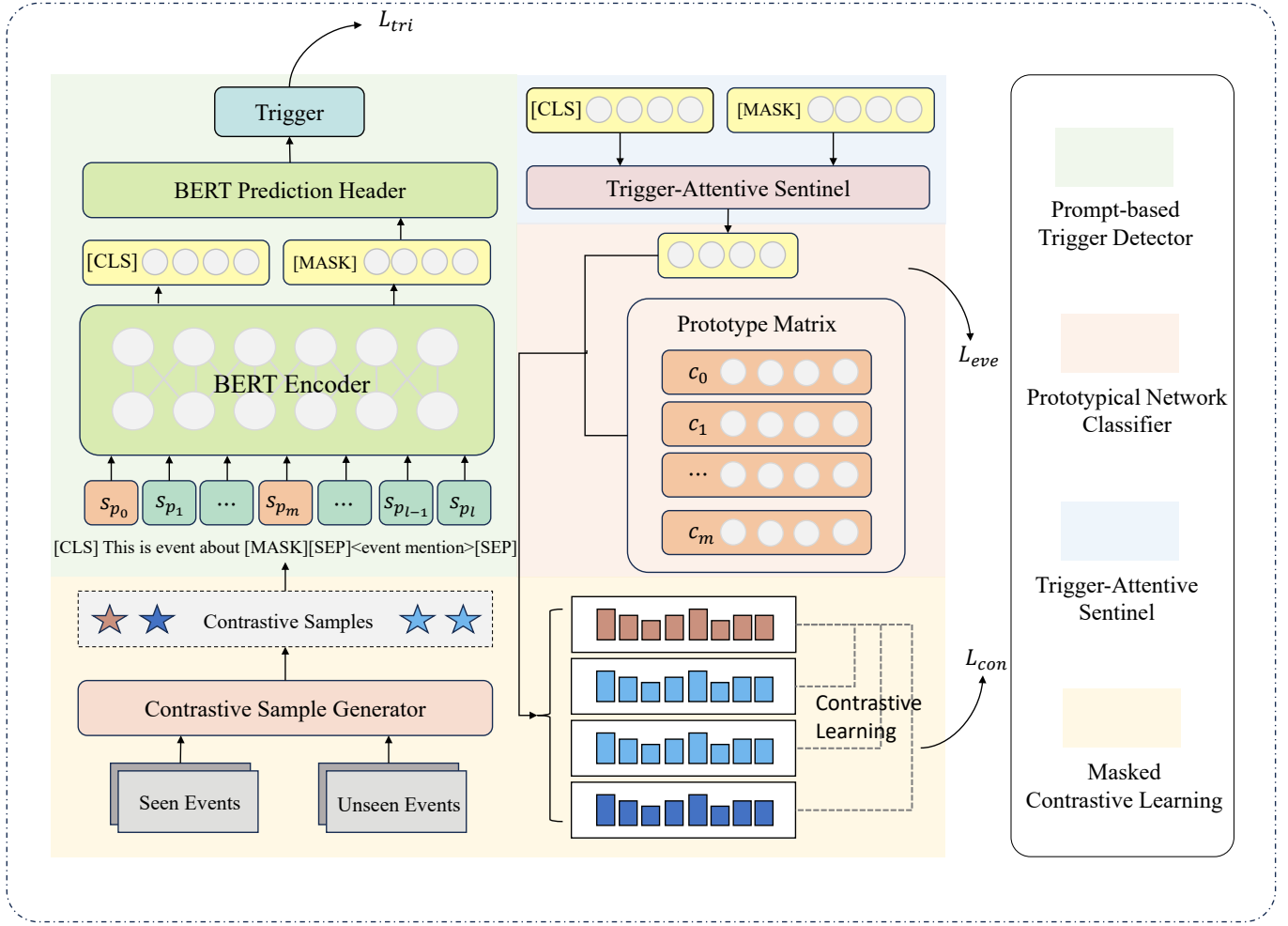


Fig. 5: The overall architecture of PBCT, which consists of four modules: A prompt-based trigger detector, a prototypical network classifier, a trigger-attentive sentinel and masked contrastive learning.

via a tokenizer, where  $s_{p_m}$  is the token at the [MASK] position and  $l$  is the length of the tokenized  $PT$ . To achieve more accurate predictions, we fine-tune the model using the cross-entropy loss function  $\mathcal{L}_{tri}$  for seen events and thereby obtain a word distribution  $P_m^t(s_{p_i}) = P_m^t(s_{p_m} = s_{p_i} | S_p)$  over words in the  $\langle event\ mention \rangle$ :

$$\mathcal{L}_{tri} = \begin{cases} -\sum_{i=8}^l P_m^{gt}(s_{p_i}) \log P_m^t(s_{p_i}) & y \in \mathcal{T}_s \\ 0 & y \in \mathcal{T}_u \end{cases} \quad (1)$$

where  $P_m^{gt}(s_{p_i}) = P_m^{gt}(s_{p_m} = s_{p_i} | S_p)$  is a ground-truth word distribution and  $i = 8$  is the index indicating the starting position of  $\langle event\ mention \rangle$  in  $PT$ .

### B. Prototypical Network Classifier

After trigger detection, a prototypical network classifier is introduced to achieve event-type prediction. Specifically, a prototype matrix  $C = [c_0, c_1, \dots, c_m]^T \in \mathbb{R}^{m \times h}$  is built, with each row representing the prototype of an embedded event type, while  $h$  denotes the embedding dimension of the PLM and  $m = |\mathcal{T}|$  is the number of event types.

To provide a solid starting point for the prototype network, we propose a semantic initialization approach to expedite model convergence and boost performance. Specifically, event-type labels are joined as  $XL = "y_0, y_1, \dots, y_m"$ . Afterward, we utilize the PLM to encode the input sequence  $XL$  consisting of type labels, thereby segmenting the output  $\mathbf{H}_{XL}$  to acquire each type representation  $\mathbf{h}_{y_j}$ :

$$\mathbf{H}_{XL} = PLM(XL) \quad (2)$$

$$\mathbf{h}_{y_0}, \mathbf{h}_{y_1}, \dots, \mathbf{h}_{y_m} = Segment(\mathbf{H}_{XL}). \quad (3)$$

It is worth noting that, in this paper, we choose a PLM based on self-attention as their attention heads facilitate direct interactions between type words. This approach can take great advantage of event-type labels and cross-type interaction relations and offer important clues for type representation.

Finally, we derive the predicted probabilities  $P(y = y_j | \mathbf{x})$  of the queried event description across each type by calculating the similarity between its event vector  $\mathbf{x}$  and the prototype of the event type  $\mathbf{c}_j$ :

$$P(y = y_j | \mathbf{x}) = \frac{\exp(-d(\mathbf{x}, \mathbf{c}_i))}{\sum_{y_k \in \mathcal{T}} \exp(-d(\mathbf{x}, \mathbf{c}_k))} \quad (4)$$

where  $d(\cdot, \cdot)$  is set as Euclidean distance.

Therefore, the event type with the high probability is chosen as the final prediction  $\hat{y}$ :

$$\hat{y} = \operatorname{argmax}_{y_j \in \mathcal{T}} P(y = y_j | \mathbf{x}) \quad (5)$$

### C. Trigger-Attentive Sentinel

With each event type associated with multiple trigger words, there's a risk that relying solely on predicted trigger words as query points may cause samples of the same event type to be widely dispersed across hidden spaces. Additionally, some trigger embedding vectors may lack distinctive semantics, with the bulk of semantic information contained within contextual embedding vectors. Therefore, we introduce a trigger-attentive sentinel to learn the trade-off between trigger features and contextual features through the computation of two attention weights,  $g_0$  and  $g_1$ :

$$g_0, g_1 = \sigma(\mathbf{W}_g[\mathbf{E}_{[CLS]} \oplus \mathbf{E}_{[MASK]}] + \mathbf{b}_g) \quad (6)$$

*s.t.*  $g_0 + g_1 = 1$

where  $\mathbf{E}_{[CLS]}$  and  $\mathbf{E}_{[MASK]}$  refer to a contextual embedding vector and a trigger embedding vector, respectively, encoded by PLMs in the prompt-based trigger detector,  $\mathbf{W}_g$  and  $\mathbf{b}_g$  are learnable parameters of the feedforward network, and  $\sigma$  denotes the softmax function. Afterward, a weighted summation of  $\mathbf{E}_{[CLS]}$  and  $\mathbf{E}_{[MASK]}$  is calculated to derive a queried vector  $\mathbf{x}$  corresponding to an event mention:

$$\mathbf{x} = g_0 \mathbf{E}_{[CLS]} + g_1 \mathbf{E}_{[MASK]} \quad (7)$$

This sentinel allows the model to dynamically adjust the combination of the two features based on trigger saliency, contributing to the final prediction.

### D. Masked Contrastive Learning

To better represent events under the zero-shot setting, especially unseen events, we design a masked contrastive learning mechanism by constructing several contrastive samples, i.e., two positive samples and one negative sample. Specifically, for positive sample One, we rephrase event mentions through back translation [51] as original event triggers are kept. This can ensure the consistency of semantics, so the rephrased sample is regarded as a positive sample. For positive sample Two, we aim to reduce lexical bias and encourage the model to make predictions based on context alone. To achieve this, we delexicalize triggers by replacing them with placeholders “[MASK]” [24], so the masked sample can serve as a positive sample. In contrast, the selection of negative samples should be divided into two situations: for seen events, seen events with different labels or unseen events are randomly selected, while for unseen events, negative samples are chosen from all seen events.

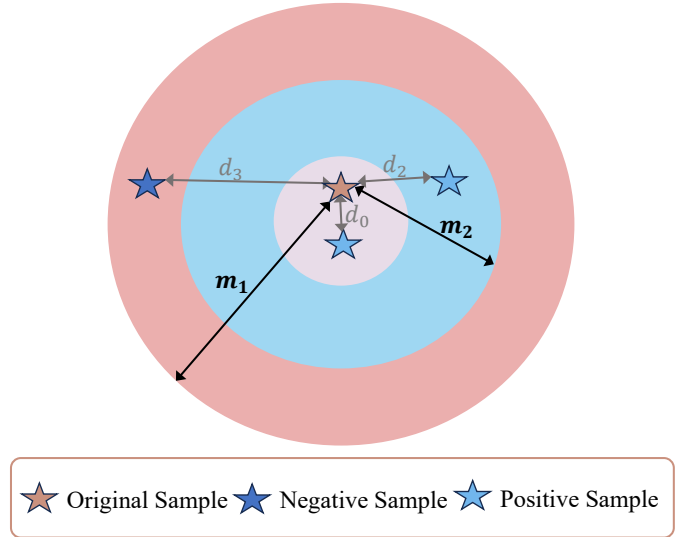


Fig. 6: The schematic of masked contrastive learning.

Therefore, the distances between the probability distributions of contrastive samples and the original sample are calculated using the Wasserstein distance as the distance function  $\mathbb{W}_p$ , denoted as  $d_1, d_2$  and  $d_3$ , as follows:

$$d_i = \mathbb{W}_p(p_i, p_0), i \in \{1, 2, 3\} \quad (8)$$

$$\mathbb{W}_p(p_i, p_j) = \inf_{\gamma \in \Pi(p_i, p_j)} \mathbf{E}_{(x, y) \sim \gamma} [c(x, y)] \quad (9)$$

where  $p_0, p_1, p_2, p_3$  refers to event type probability distribution of the original, rephrased, masked and negative sample,  $\Pi(p_i, p_j)$  represents the set of all possible joint distributions combined by two distributions  $p_i$  and  $p_j$ ,  $(x, y)$  is a sample drawn from a joint distribution  $\gamma$ ,  $c(x, y)$  denotes the cost function for transferring  $x$  to  $y$  and is typically set as the  $L_p$  distance,  $\inf$  refers to the minimum expected cost.

The goal of contrastive learning is to bring positive samples closer and push negative samples farther away. However, certain negative samples start with a considerable distance from the original sample. Taking these distances into account may lead to a shift in learning focus, potentially leaving smaller distances unexpanded. Hence, it is necessary to establish a margin  $m_1$  to define the negative samples that require pushing further away, as follows:

$$\hat{d}_3 = \max(0, m_1 - d_3). \quad (10)$$

Moreover, for certain event mentions, the removal of trigger words has minimal impact on semantic expression, resulting in a relatively close distance between the masked sample and the original sample. However, there are instances where eliminating the trigger words render the expressed event type unrecognizable even to humans because the majority of the semantics of this event mention is implied by its trigger words. If these masked samples with larger distances are forcefully brought closer to the original sample, it could mislead the representation of other events, given the substantial overlap of tokens between masked event mentions and mentions of other

TABLE I: Dataset split and statistics. The last two columns are the mean and standard deviation of samples by types.

Dataset	Split	Total Sample	Total Type	Seen Sample Type		Unseen Sample Type		Mean	Stdev
ACE 2005	Total	3805	33	2316	1489	17	16	115.30	21.60
	Training	3043		1852	1191			92.21	17.26
	Validation	381		232	149			11.55	2.17
	Test	381		232	149			11.55	2.17
FewEvent	Total	74330	100	40891	33439	50	50	743.30	74.52
	Training	59463		32712	26751			594.63	59.61
	Validation	7433		4089	3344			74.33	7.45
	Test	7434		4090	3344			74.33	7.46

event types. Therefore, similar to negative samples, we must design a margin  $m_2$  to mask out masked samples that extend beyond it, as follows:

$$\hat{d}_2 = \min(m_2, d_2). \quad (11)$$

Fig. 6 depicts a schematic of the masked contrastive learning mechanism. The total contrastive loss is calculated as:

$$\mathcal{L}_{con} = d_1 + \hat{d}_2 + \hat{d}_3 \quad (12)$$

### E. Supervised Loss

Since the labels of seen event samples are available to the model, supervised learning can be applied alongside contrastive learning, to optimize the first  $|\mathcal{T}_s|$  rows of the prototype matrix corresponding to seen event types. To this end, the cross-entropy loss  $\mathcal{L}_{eve}$  for event classification is calculated as:

$$\mathcal{L}_{eve} = \begin{cases} -\sum_{i=0}^{|\mathcal{T}_s|-1} P^{gt}(y = y_i|\mathbf{x}) \log P(y = y_i|\mathbf{x}) & y \in \mathcal{T}_s \\ 0 & y \in \mathcal{T}_u \end{cases} \quad (13)$$

where  $P^{gt}(y = y_i|\mathbf{x})$  is a ground-truth event type distribution.

Therefore, the total loss is the summation of the supervised loss and the contrastive loss:

$$\mathcal{L} = \mathcal{L}_{eve} + \mathcal{L}_{tri} + \lambda * \mathcal{L}_{con} \quad (14)$$

where  $\lambda$  is a trade-off hyper-parameter between the two types of losses.

## VI. EXPERIMENTS

### A. Datasets and Baselines

We take ACE 2005<sup>1</sup> [52] and FewEvent<sup>2</sup> [53] as the benchmark datasets. ACE 2005 stands out as the most widely-used dataset for event detection, drawing from a diverse array of news sources spanning 6 distinct categories. It includes annotations for 8 event types and 33 more finely-grained subtypes. FewEvent currently stands as the largest few-shot dataset for event detection, covering 19 event types, which are further subdivided into 100 event subtypes. To ensure a

balanced distribution of samples between seen and unseen types, we arrange the event types in descending order according to their respective sample counts, where odd-positioned event types are considered as seen types and even-positioned ones as unseen types. Additionally, these datasets are split into training, validation and test sets in proportions of 80%, 10% and 10%. According to the statistics, while FewEvent has more samples and types than ACE 2005, it is facing a more conspicuous sampling bias issue. More details about the two datasets and their splits are shown in TABLE I.

To showcase its superiority, our methods are compared against the following baselines:

- **SCCL** [54] It is a state-of-the-art model designed for unsupervised text clustering, which leverages instance-wise contrastive learning to support unsupervised clustering. Instead of using the [CLS] token as the representation for event mentions, we employ the contextual vector of identical candidate trigger words as SS-VQ-VAE, tailored to fit the event detection task.
- **SS-VQ-VAE** [46] It is a semi-supervised vector quantized variational autoencoder approach that associates the sense of each word with OntoNotes [55], consequently considering all mapped noun and verb concepts as candidate triggers.
- **BERT-MCL** We fine-tune BERT-Base<sup>3</sup> and BERT-Large<sup>4</sup> with our proposed masked contrastive learning. KNN algorithm for seen event detection and K-means algorithm for unseen event detection are utilized once event encoding is obtained.

### B. Training Settings and Evaluation Metrics

We code and implement our PBCT model by using PyTorch on a Linux server with four GPU Tesla V100. BERT-Base and BERT-Large are chosen as the PLM for the prompt-based trigger detector. Therefore, the total parameters of the executed PBCT are 109.59M and 335.28M, respectively. Argos Translate<sup>5</sup> is employed for back translation, with Chinese designated as the intermediary language. For the ACE 2005 dataset, we employ the Adam optimizer with a learning rate of  $1e-6$  and a weight decay of  $1e-6$  for the parameters of

<sup>1</sup><https://catalog.ldc.upenn.edu/LDC2006T06>

<sup>2</sup>[https://github.com/231sm/Low\\_Resource\\_KBP](https://github.com/231sm/Low_Resource_KBP)

<sup>3</sup><https://huggingface.co/google-bert/bert-base-uncased>

<sup>4</sup><https://huggingface.co/google-bert/bert-large-uncased>

<sup>5</sup><https://www.argosopentech.com>

TABLE II: Comparison of various detection methods on ACE 2005 and FewEvent datasets. The best results are bolded and the second-best results are underlined.

Model	ACE 2005							FewEvent						
	F1-Seen	F1-Unseen	NMI	FM	P	R	F1	F1-Seen	F1-Unseen	NMI	FM	P	R	F1
SCCL	0.5936	0.2893	0.2607	0.2185	0.5543	0.4987	0.4746	0.7892	0.3140	0.2832	0.2855	0.7374	0.6032	0.6250
SS-VQ-VAE	0.7233	0.2812	0.1202	<b>0.4096</b>	0.5205	0.6116	0.5504	0.9232	0.4344	0.1710	0.5754	0.7226	0.7608	0.7033
BERT-MCL-Base	0.4108	0.2765	0.3312	0.1682	0.4702	0.3753	0.3583	0.8681	0.5285	<b>0.4722</b>	0.6149	0.7368	0.5075	0.7153
BERT-MCL-Large	0.4686	0.3359	<b>0.3641</b>	0.2168	0.5222	0.4252	0.4167	0.8710	0.4611	0.4119	0.5561	0.7167	0.6691	0.6866
PBCT-Base	<b>0.8153</b>	<b>0.4327</b>	0.3633	0.3934	0.6790	<b>0.6693</b>	<b>0.6657</b>	0.9389	0.5609	0.4536	0.7471	0.7957	0.7998	0.7689
PBCT-Large	<u>0.8017</u>	<u>0.3527</u>	0.3027	0.3361	<b>0.6823</b>	0.6194	0.6261	<b>0.9454</b>	<b>0.5634</b>	<u>0.4575</u>	<b>0.7514</b>	<b>0.8055</b>	<b>0.8056</b>	<b>0.7736</b>

TABLE III: Results of ablation studies. The best results are bolded.

Model	ACE 2005							FewEvent						
	F1-Seen	F1-Unseen	NMI	FM	P	R	F1	F1-Seen	F1-Unseen	NMI	FM	P	R	F1
PBCT-Base	<u>0.8153</u>	<b>0.4327</b>	<b>0.3633</b>	0.3934	<b>0.6790</b>	<b>0.6693</b>	<b>0.6657</b>	<b>0.9389</b>	<b>0.5609</b>	<b>0.4536</b>	<b>0.7471</b>	0.7957	<b>0.7998</b>	<b>0.7689</b>
-wo Atten.	0.7998	0.3301	0.1774	<b>0.4963</b>	0.6257	0.6562	0.6161	0.9351	0.5541	0.4322	0.7337	0.8019	0.7963	0.7637
-wo SemanInit.	0.7820	0.3274	0.1800	0.4935	0.6138	0.6536	0.6042	0.9306	0.5559	0.4338	0.7356	0.7918	0.7920	0.7621
-wo MaskContra.	<b>0.8418</b>	0.3545	0.3377	0.3564	0.6772	0.6614	0.6512	0.9347	0.5519	0.4283	0.7308	<b>0.8063</b>	0.7957	0.7625

BERT, alongside a learning rate of  $1e-3$  for the remaining parameters. Conversely, for the FewEvent dataset, we adjust the learning rate to  $1e-5$  for the BERT parameters and  $1e-2$  for the remaining parameters. Our evaluation task is composed of two parts: supervised event detection for seen events and zero-shot event detection for unseen events, mirroring real-world scenarios in open environments. Our model directly outputs the predicted labels of seen events and maps its output into the predicted labels of unseen events by the Hungarian Algorithm [56]. We use the F1 score as the same metric for the two parts, i.e., F1-Seen and F1-Unseen and especially leverage Normalized Mutual Info (NMI) and Fowlkes Mallows (FM) [46], [54] to assess the clustering performance for unseen events. Additionally, we opt for weighted precision, recall, and F1 to evaluate the overall performance in both supervised and zero-shot settings, denoted as P, R, F1. We tune the PBCT’s hyperparameters according to F1 scores on each validation set. The margins  $m_1$  and  $m_2$  of masked contrastive learning are set to 1 and 0.6. The trade-off hyperparameter  $\lambda$  between losses is assigned to 0.7.

### C. Results and analysis

We compare the test results of our PBCT model with those of baselines for event detection on the ACE 2005 and FewEvent datasets, as shown in Table II. Our model surpasses the state-of-the-art model on most metrics and demonstrates comparable performance on remaining metrics. BERT-MCL achieves competitive results for unseen events, thus highlighting the effectiveness of masked contrastive learning in zero-shot training scenarios. SS-VQ-VAE’s superior performance for seen events over SCCL and BERT-MCL underscores the indispensability and necessity of supervised learning. PBCT outperforms BERT-MCL across most metrics in detecting unseen events, which demonstrates the prototype network’s superior ability to learn representations for each type. The superior performance of PBCT over SS-VQ-VAE suggests that

a prompt-based trigger detector holds its own against heuristic rules. PBCT-Base generally exhibits better performance than PBCT-Large on the ACE 2005 dataset, but it underperforms PBCT-Large on the FewEvent dataset. We believe this is because the smaller data scale of ACE 2005 led to overfitting or underfitting in the PBCT-Large model with a larger parameter size.

PBCT achieves remarkable enhancements, with a 9.2% improvement in F1-Seen and 9.68% in F1-Unseen over the best baseline on the ACE 2005 dataset, along with gains of 2.22% in F1-Seen and 3.49% in F1-Unseen on the FewEvent dataset. This showcases that PBCT has a slightly greater advantage in zero-shot unseen event detection compared to other baselines. Our model also achieves improvements of 12.8% in P, 5.77% in R and 11.53% in F1 on the ACE 2005 dataset as well as 6.81% in P, 4.48% in R and 5.83% in F1 on the FewEvent dataset, compared to the state-of-the-art results. It is observed that our model’s performance improvement on the ACE 2005 dataset surpasses its improvement on the FewEvent dataset, suggesting an edge in scenarios with limited available data.

### D. Ablation study

We perform ablation experiments on the ACE 2005 and FewEvent datasets to examine the impact of removing certain components of the PBCT model. This aids in enhancing our comprehension regarding the significance of individual components within the model’s overall performance. The results of ablation studies are shown in Table III.

To evaluate the influence of semantic initialization, we replace it with random initialization (denoted as -SemanInit.). This results in a decrease of F1 score by 6.15% on the ACE 2005 dataset and by 0.68% on the FewEvent dataset. Obviously, this difference indicates that semantic initialization is more crucial for training on small-scale datasets, as it can shorten the path of parameter optimization. Removing the



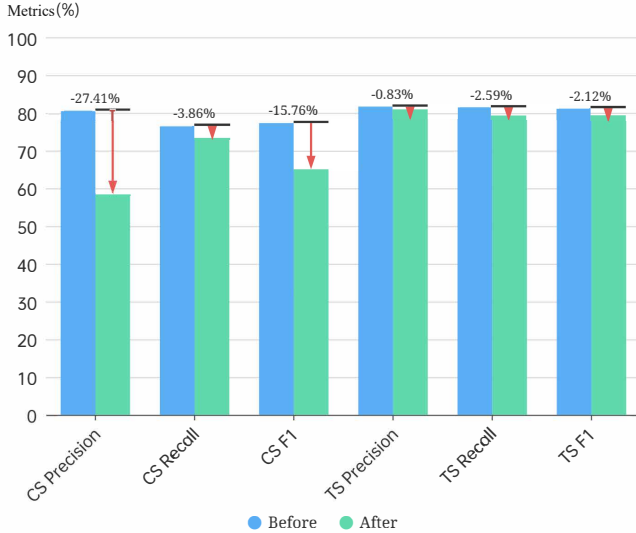


Fig. 7: The performance change in event detection of trigger-salient (TS) and context-salient (CS) types before and after removing masked contrastive learning and trigger-attentive sentinel.

masked contrastive learning (i.e., -wo MaskContra.) leads to a slight overall decrease, except an increase in F1-Seen on the ACE 2005 dataset. This could be attributed to the model’s focus on refining seen event sample and type representations, thereby leading to higher F1 scores of seen event metrics. Furthermore, the model utilized for encoding seen event sample representations is identical to the one used for encoding unseen event sample representations, thereby indirectly optimizing the representations of unseen events. In addition, semantic initialization offers a high-quality type representation for unseen event types, ensuring that there is no significant decrease in the F1 scores for detecting unseen events. To evaluate the effectiveness of trigger-attentive sentinel, we replace it with the sum of contextual embedding vectors and trigger embedding vectors (denoted as -wo Atten.). Its inferior performance relative to PBCT suggests that effectively discriminating between different event types has been essential for enhancing performance.

### E. Trigger Saliency Analysis

We leverage seen events from the ACE 2005 dataset to the unique characteristics between trigger-salient (TS) and context-salient (CS) types. Specifically, we compute the mean value of samples’  $g_1$  within each event type  $y_j \in \mathcal{T}_s$  on the training set  $\mathcal{D}_s$  as trigger saliency  $\bar{g}_j$ :

$$\bar{g}_j = \frac{1}{|\mathcal{D}_s(y_j)| - 1} \sum_{i=0}^{|\mathcal{D}_s(y_j)|-1} g_1^i \quad (15)$$

$$j = 0, 1, \dots, |\mathcal{T}_s|$$

where  $\mathcal{D}_s(y_j)$  represents seen event samples labeled as type  $y_j$  in the training dataset.

We arrange  $\bar{g}_j$  in descending order and divide event types into TS for the top half and CS for the bottom half, determined by the median. Afterward, we test the performance change in event detection of trigger-salient (TS) and context-salient (CS) types before and after removing masked contrastive learning and trigger-attentive sentinel, as shown in Fig. 7. Obviously, CS types exhibit a more substantial decline in performance compared to TS types, particularly in terms of precision. We believe this phenomenon arises from the model’s post-removal tendency to focus on learning lexical patterns rather than developing a nuanced understanding of contextual cues. Consequently, classification tasks degrade into simplistic matching ones. The reliance of CS types on contexts makes them more sensitive to the removal of two components facilitating contextual understanding.

### F. Hyperparameter Analysis

We investigate the impact of varying hyperparameters  $m_1$ ,  $m_2$  and  $\lambda$  on PBCT performance using the ACE 2005 dataset, as shown in Fig. 8.

The hyperparameter  $m_1$  regulates the scope of negative samples to be considered for pulling away. As depicted in Fig. 8 (a), with higher values of  $m_1$  than a certain boundary, F1-Seen decreases while F1-Unseen increases. This may be because the introduction of distant negative samples results in substantial contrastive loss, thereby impeding the thorough optimization of the cross-entropy loss. Therefore, with the comprehensive assessment consideration of the detecting performance of both event types, we employ the F1 score as the metric to select the optimal value for  $m_1$ , i.e.,  $m_1 = 2.0$ .

Similar to  $m_1$ , the hyperparameter  $m_2$  controls the range of masked samples that should be considered for pulling closer. Fig. 8 (b) illustrates an initial ascent trend of F1 followed by a subsequent decline as  $m_2$  increases, especially F1-Seen. We attribute this trend change to the distortion of representations of other events caused by forcibly pulling in masked samples from a greater distance.  $m = 0.6$  is chosen as the optimal value as it yields the best F1 scores for both F1 and F1-Unseen, along with competitive performance in F1-Seen.

The aim of tuning the hyperparameter  $\lambda$  is to achieve a balance between the influences of supervised learning and contrastive learning. As shown in Fig. 8 (c), there is a slight downward trend in the value of F1-Seen with an increase in the value of  $\lambda$ . This is evidently due to the reduced proportion of supervised loss, leading to insufficient training in seen event representations. Meanwhile, F1-Unseen demonstrates an overall upward trend, accompanied by significant fluctuations.  $\lambda = 0$  represents the absence of masked contrastive learning, and its relatively high F1-Unseen value is caused by effective supervised learning indirectly optimizing representations of unseen event samples. Therefore, we should select an appropriate  $\lambda$  value that maximizes the benefits of contrastive learning without hindering supervised learning. Obviously, the optimal choice for  $\lambda$  is 0.7.

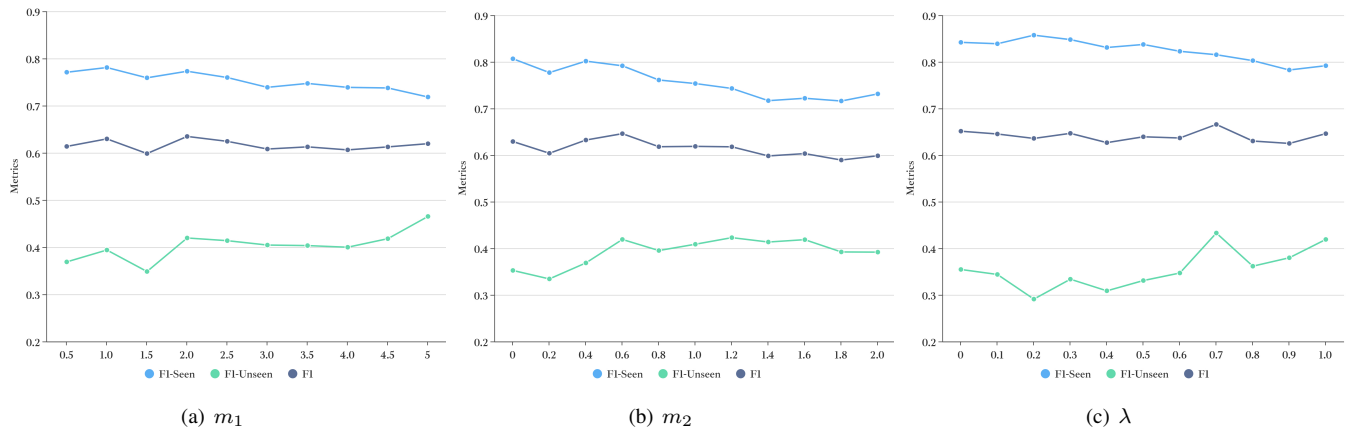


Fig. 8: Sensitivity analysis results of hyperparameters in F1-Seen, F1-Unseen and F1.

### G. Visualization Analysis

For the visual observation of models’ representation learning, we employ t-SNE<sup>6</sup> to visualize various event samples in the FewEvent dataset, with BERT-Base being used as the PLM, as shown in Fig. 9. PBCT can achieve better representations of both unseen events and seen events, while SS-VQ-VAE performs poorly on unseen event types, and BERT-MCL on seen event types. This suggests that both supervised learning and contrastive learning play crucial roles in representation learning. By comparing the representations of seen and unseen events, we can discern that classifying unseen event samples is more intricate than classifying seen ones, with only a handful of high-sample-size types being accurately classified. Generally, the greater the number of samples within a type, the more effective the clustering.

### H. Case Study

To validate the model’s effectiveness in unfamiliar and specialized domains, we utilize ChatGPT to generate sentences describing the steps of PCB design and manufacturing. Subsequently, these sentences are formatted into a series of event detection samples and inputted into the PBCT model that has already been trained on the FewEvent dataset. Their input event mentions as well as the output recognized trigger words, and classified event types are presented in TABLE IV. It’s notable that despite the absence of PCB-specific knowledge in the FewEvent dataset, the model proficiently identifies trigger words and categorizes them under the *Education.Education* label. However, the model is constrained by its ability to recognize only a single event, which is a focal point for our future efforts.

## VII. CONCLUSION AND FUTURE WORK

In this paper, we integrate RAG technology into the crowdsourcing framework to achieve automated crowdsourcing TD accurately and universally. By this means, TD tasks are reconfigured into event detection tasks. Therefore, we design a

TABLE IV: Event detection examples in unfamiliar domains using the trained PBCT model

<i>Event Mention:</i> It starts with creating a wiring layout using circuit design software.
<i>Trigger Word:</i> creating
<i>Event Type:</i> Education.Education
<i>Event Mention:</i> It is followed by crafting a PCB schematic to represent the circuit connections.
<i>Trigger Word:</i> crafting
<i>Event Type:</i> Education.Education
<i>Event Mention:</i> Next, the circuit diagram is transformed into an actual layout using PCB layout software.
<i>Trigger Word:</i> transformed
<i>Event Type:</i> Education.Education

PBCT framework that incorporates prompt learning, a trigger-attentive sentinel and masked contrastive learning to address the reliance on external tools and the challenge of selectively handling different event types. The experimental results indicate our model performs comparably or even surpasses the state-of-the-art baselines in event detection for both seen and unseen events.

However, there are some limitations in establishing the workflow of decomposed subtasks for subsequent task recommendations. (1) The method is unable to address cases where multiple events are described within a single sentence, a scenario frequently encountered in real-world contexts. (2) The detected events (also known as subtasks) have temporal dependencies during execution. To allocate subtasks effectively, it is crucial to consider the temporal relationships between events. (3) It is assumed that relevant descriptions of task execution steps have already been retrieved prior to application before our models are applied, but in reality, the efficient and accurate retrieval of key content is of critical importance as it directly influences the performance of subsequent TD. Therefore, our research will focus on RAG-based crowdsourcing from three perspectives: multi-event detection, temporal relation extraction, and key content retrieval.

## REFERENCES

- [1] X. Wang, Y. Wang, J. Yang, X. Jia, L. Li, W. Ding, and F.-Y. Wang, “The survey on multi-source data fusion in cyber-physical-social systems:

<sup>6</sup><https://scikit-learn.org/stable/modules/generated/sklearn.manifold.TSNE.html>

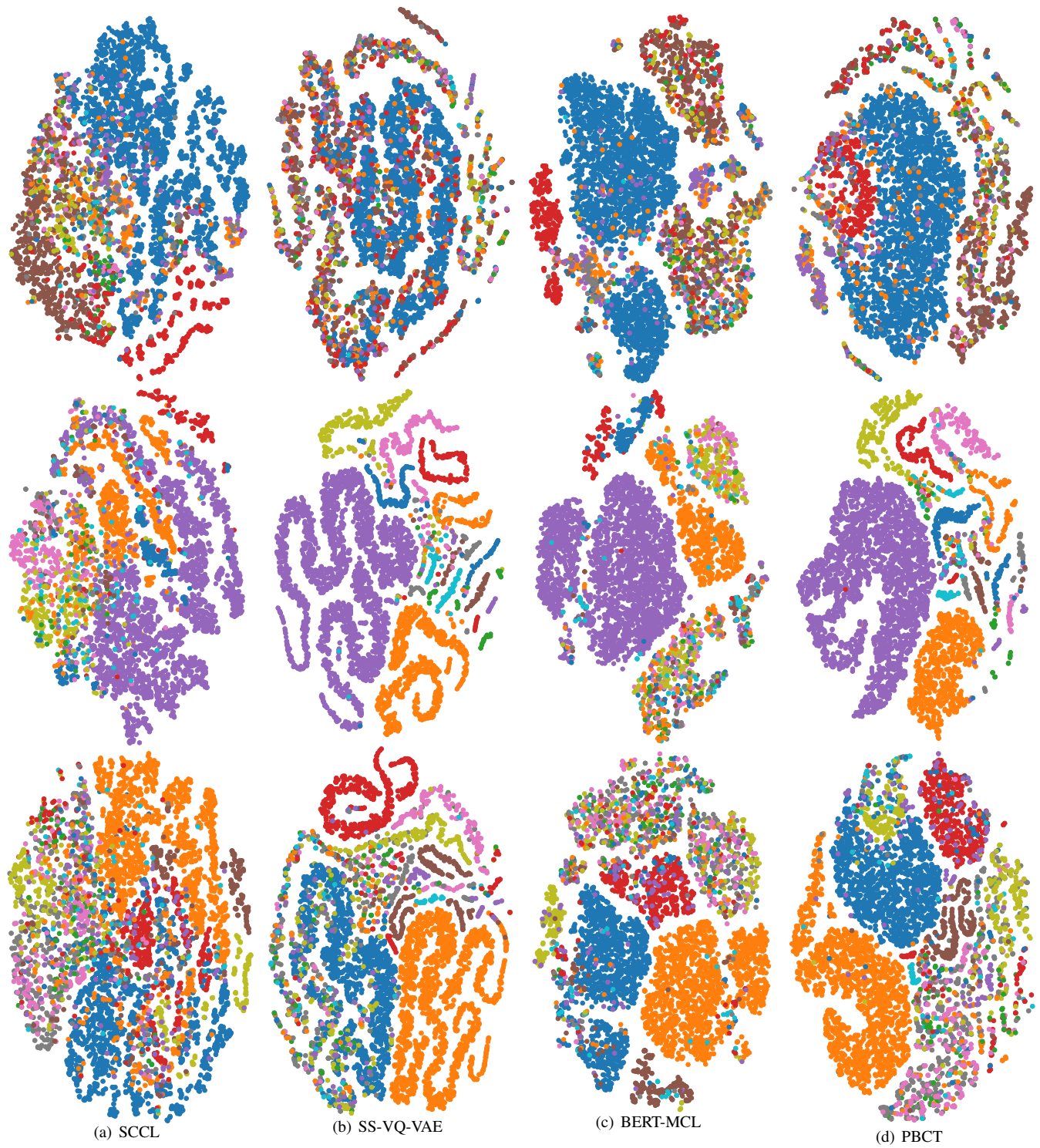


Fig. 9: Visualization of samples across various event types. From top to bottom are respectively the unseen event samples, seen event samples, and total event samples.

- Foundational infrastructure for industrial metaverses and industries 5.0,” *Information Fusion*, p. 102321, 2024.
- [2] J. Yang, Y. Wang, X. Wang, X. Wang, X. Wang, and F.-Y. Wang, “Generative ai empowering parallel manufacturing: Building a “6s” collaborative production ecology for manufacturing 5.0,” *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, pp. 1–15, 2024.
- [3] J. Yang, X. Wang, and Y. Zhao, “Parallel manufacturing for industrial metaverses: A new paradigm in smart manufacturing,” *IEEE/CAA Journal of Automatica Sinica*, vol. 9, no. 12, pp. 2063–2070, 2022.
- [4] J. Yang, S. Li, X. Wang, J. Lu, H. Wu, and X. Wang, “Defact in manuse for parallel manufacturing: Foundation models and parallel workers in smart factories,” *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, vol. 53, no. 4, pp. 2188–2199, 2023.
- [5] G. Xiong, F.-Y. Wang, T. R. Nyberg, X. Shang, M. Zhou, Z. Shen, S. Li, and C. Guo, “From mind to products: Towards social manufacturing and service,” *IEEE/CAA Journal of Automatica Sinica*, vol. 5, no. 1, pp. 47–57, 2017.
- [6] J. Yang, X. Wang, Y. Tian, and F.-Y. Wang, “Parallel intelligence in cps: being, becoming, and believing,” *IEEE Intelligent Systems*, vol. 38, no. 6, pp. 75–80, 2023.
- [7] X. Wang, J. Yang, X. Li, and F.-Y. Wang, “Parallel intelligence for cps: An acp approach,” *Cyber-Physical-Human Systems: Fundamentals and Applications*, pp. 145–169, 2023.
- [8] X. Niu and S. Qin, “A review of crowdsourcing technology for product design and development,” in *2017 23rd International Conference on Automation and Computing (ICAC)*. IEEE, 2017, pp. 1–6.
- [9] H. Jiang and S. Matsubara, “Efficient task decomposition in crowdsourcing,” in *PRIMA 2014: Principles and Practice of Multi-Agent Systems: 17th International Conference, Gold Coast, QLD Australia, December 1-5, 2014. Proceedings 17*. Springer, 2014, pp. 65–73.
- [10] Y. Rizk, M. Awad, and E. W. Tunstel, “Cooperative heterogeneous multi-robot systems: A survey,” *ACM Computing Surveys (CSUR)*, vol. 52, no. 2, pp. 1–31, 2019.
- [11] T. Hu, S. Messelodi, and O. Lanz, “Dynamic task decomposition for probabilistic tracking in complex scenes,” in *2014 22nd International Conference on Pattern Recognition*. IEEE, 2014, pp. 4134–4139.
- [12] X. Li, S. Dang, K. Li, and Q. Liu, “Multi-agent-based battlefield reconnaissance simulation by novel task decomposition and allocation,” in *2010 5th International Conference on Computer Science & Education*. IEEE, 2010, pp. 1410–1414.
- [13] K. Shiarlis, M. Wulfmeier, S. Salter, S. Whiteson, and I. Posner, “Taco: Learning task decomposition via temporal alignment for control,” in *International Conference on Machine Learning*. PMLR, 2018, pp. 4654–4663.
- [14] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, “Bert: Pre-training of deep bidirectional transformers for language understanding,” *arXiv preprint arXiv:1810.04805*, 2018.
- [15] T. Wu, S. He, J. Liu, S. Sun, K. Liu, Q.-L. Han, and Y. Tang, “A brief overview of chatgpt: The history, status quo and potential future development,” *IEEE/CAA Journal of Automatica Sinica*, vol. 10, no. 5, pp. 1122–1136, 2023.
- [16] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar *et al.*, “Llama: Open and efficient foundation language models,” *arXiv preprint arXiv:2302.13971*, 2023.
- [17] J. Gao, H. Zhao, C. Yu, and R. Xu, “Exploring the feasibility of ChatGPT for event extraction,” *arXiv preprint arXiv:2303.03836*, 2023.
- [18] C. Yuan, Q. Xie, and S. Ananiadou, “Zero-shot temporal relation extraction with ChatGPT,” *arXiv preprint arXiv:2304.05454*, 2023.
- [19] X. Li, Y. Zhang, and E. C. Malthouse, “A preliminary study of ChatGPT on news recommendation: Personalization, provider fairness, fake news,” *arXiv preprint arXiv:2306.10702*, 2023.
- [20] P. A. Jansen, “Visually-grounded planning without vision: Language models infer detailed plans from high-level instructions,” *arXiv preprint arXiv:2009.14259*, 2020.
- [21] K. Sakaguchi, C. Bhagavatula, R. L. Bras, N. Tandon, P. Clark, and Y. Choi, “Proscript: Partially ordered scripts generation via pre-trained language models,” *arXiv preprint arXiv:2104.08251*, 2021.
- [22] S. S. Kannan, V. L. Venkatesh, and B.-C. Min, “Smart-LLM: Smart multi-agent robot task planning using large language models,” *arXiv preprint arXiv:2309.10062*, 2023.
- [23] P. Lewis, E. Perez, A. Piktus, F. Petroni, V. Karpukhin, N. Goyal, H. Küttler, M. Lewis, W.-t. Yih, T. Rocktäschel *et al.*, “Retrieval-augmented generation for knowledge-intensive nlp tasks,” *Advances in Neural Information Processing Systems*, vol. 33, pp. 9459–9474, 2020.
- [24] J. Liu, Y. Chen, K. Liu, Y. Jia, and Z. Sheng, “How does context matter? on the robustness of event detection with context-selective mask generalization,” in *Findings of the Association for Computational Linguistics: EMNLP 2020*, 2020, pp. 2523–2532.
- [25] X. Wang, Y. Wang, M. Netto, L. Stapleton, Z. Wan, and F.-Y. Wang, “Smart decentralized autonomous organizations and operations for smart societies: Human–autonomous organizations for industry 5.0 and society 5.0,” *IEEE Intelligent Systems*, vol. 38, no. 6, pp. 70–74, 2023.
- [26] J. Li, X. Liang, R. Qin, and F.-Y. Wang, “From dao to tao: Finding the essence of decentralization,” in *2023 IEEE International Conference on Systems, Man, and Cybernetics (SMC)*. IEEE, 2023, pp. 4283–4288.
- [27] M. Dorigo, D. Floreano, L. M. Gambardella, F. Mondada, S. Nolfi, T. Baaboura, M. Birattari, M. Bonani, M. Brambilla, A. Brutschy *et al.*, “Swarmanoid: a novel concept for the study of heterogeneous robotic swarms,” *IEEE Robotics & Automation Magazine*, vol. 20, no. 4, pp. 60–71, 2013.
- [28] J. Kiener and O. Von Stryk, “Towards cooperation of heterogeneous, autonomous robots: A case study of humanoid and wheeled robots,” *Robotics and Autonomous Systems*, vol. 58, no. 7, pp. 921–929, 2010.
- [29] A. Kittur, B. Smus, S. Khamkar, and R. E. Kraut, “Crowdforge: Crowdsourcing complex work,” in *Proceedings of the 24th annual ACM symposium on User interface software and technology*, 2011, pp. 43–52.
- [30] A. Kulkarni, M. Can, and B. Hartmann, “Collaboratively crowdsourcing workflows with turkomatic,” in *Proceedings of the acm 2012 conference on computer supported cooperative work*, 2012, pp. 1003–1012.
- [31] S. Xie, X. Wang, B. Yang, M. Long, J. Zhang, and L. Wang, “A multi-stage framework for complex task decomposition in knowledge-intensive crowdsourcing,” in *2021 IEEE International Conference on Industrial Engineering and Engineering Management (IEEM)*. IEEE, 2021, pp. 1432–1436.
- [32] J. Motes, R. Sandström, H. Lee, S. Thomas, and N. M. Amato, “Multi-robot task and motion planning with subtask dependencies,” *IEEE Robotics and Automation Letters*, vol. 5, no. 2, pp. 3338–3345, 2020.
- [33] S. C. Botelho and R. Alami, “M+: a scheme for multi-robot cooperation through negotiated task allocation and achievement,” in *Proceedings 1999 IEEE international conference on robotics and automation (Cat. No. 99CH36288C)*, vol. 2. IEEE, 1999, pp. 1234–1239.
- [34] Y. Tong, L. Chen, Z. Zhou, H. V. Jagadish, L. Shou, and W. Lv, “Slade: A smart large-scale task decomposer in crowdsourcing,” *IEEE Transactions on Knowledge and Data Engineering*, vol. 30, no. 8, pp. 1588–1601, 2018.
- [35] X. Zhang, Y. Yang, and B. Bao, “Task decomposition and grouping for customer collaboration in product development,” *Journal of Intelligent Systems*, vol. 25, no. 3, pp. 361–375, 2016.
- [36] L. C. Cobo, C. L. Isbell Jr, and A. L. Thomaz, “Automatic task decomposition and state abstraction from demonstration,” in *Proceedings of the 11th International Conference on Autonomous Agents and Multiagent Systems-Volume 1*, 2012, pp. 483–490.
- [37] Y. Shu, H. Chen, S. Li, and F. Hu, “The verification approach to complex tasks’ functional specification in software crowdsourcing,” in *2016 5th International Conference on Computer Science and Network Technology (ICCSNT)*. IEEE, 2016, pp. 171–176.
- [38] L. Song, Y. Fu, J. Su, K. Zhou, and M. Long, “A novel modeling method of the crowdsourcing design process for complex products-based an object-oriented petri net,” *IEEE Access*, vol. 9, pp. 41 430–41 440, 2021.
- [39] Y. Lin, H. Ji, F. Huang, and L. Wu, “A joint neural model for information extraction with global features,” in *Proceedings of the 58th annual meeting of the association for computational linguistics*, 2020, pp. 7999–8009.
- [40] D. Wadden, U. Wennberg, Y. Luan, and H. Hajishirzi, “Entity, relation, and event extraction with contextualized span representations,” in *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing*, 2019, pp. 5784–5789.
- [41] J. Liu, Y. Chen, K. Liu, W. Bi, and X. Liu, “Event extraction as machine reading comprehension,” in *Proceedings of the 2020 conference on empirical methods in natural language processing (EMNLP)*, 2020, pp. 1641–1651.
- [42] X. Du and C. Cardie, “Event extraction by answering (almost) natural questions,” in *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing*, 2020, pp. 671–683.
- [43] Y. Lu, H. Lin, J. Xu, X. Han, J. Tang, A. Li, L. Sun, M. Liao, and S. Chen, “Text2event: Controllable sequence-to-structure generation for end-to-end event extraction,” in *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing*, 2021, pp. 2795–2806.

- [44] X. Liu, H. Huang, G. Shi, and B. Wang, "Dynamic prefix-tuning for generative template-based event extraction," in *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics*, 2022, pp. 5216–5228.
- [45] L. Huang, H. Ji, K. Cho, and C. R. Voss, "Zero-shot transfer learning for event extraction," in *The 56th Annual Meeting of the Association for Computational Linguistics*, 2017.
- [46] L. Huang and H. Ji, "Semi-supervised new event type induction and event detection," in *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2020, pp. 718–724.
- [47] Q. Lyu, H. Zhang, E. Sulem, and D. Roth, "Zero-shot event extraction via transfer learning: Challenges and insights," in *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, 2021, pp. 322–332.
- [48] H. Zhang, H. Wang, and D. Roth, "Zero-shot label-aware event trigger and argument classification," in *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, 2021, pp. 1331–1340.
- [49] S. Li, Q. Zhan, K. Conger, M. Palmer, H. Ji, and J. Han, "Glen: General-purpose event detection for thousands of types," in *The 2023 Conference on Empirical Methods in Natural Language Processing*, 2023.
- [50] J. Liu, Y. Chen, and J. Xu, "Saliency as evidence: Event detection with trigger saliency attribution," in *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2022, pp. 4573–4585.
- [51] M. Fadaee and C. Monz, "Back-translation sampling by targeting difficult words in neural machine translation," in *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, 2018, pp. 436–446.
- [52] D. Ahn, "The stages of event extraction," in *Proceedings of the Workshop on Annotating and Reasoning about Time and Events*, 2006, pp. 1–8.
- [53] S. Deng, N. Zhang, J. Kang, Y. Zhang, W. Zhang, and H. Chen, "Meta-learning with dynamic-memory-based prototypical network for few-shot event detection," in *Proceedings of the 13th international conference on web search and data mining*, 2020, pp. 151–159.
- [54] D. Zhang, F. Nan, X. Wei, S. Li, H. Zhu, K. McKeown, R. Nallapati, A. Arnold, and B. Xiang, "Supporting clustering with contrastive learning," in *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2021, pp. 5419–5430.
- [55] E. Hovy, M. Marcus, M. Palmer, L. Ramshaw, and R. Weischedel, "Ontonotes: the 90% solution," in *Proceedings of the human language technology conference of the NAACL, Companion Volume: Short Papers*, 2006, pp. 57–60.
- [56] S. Zhang, T. Ji, W. Ji, and X. Wang, "Zero-shot event detection based on ordered contrastive learning and prompt-based prediction," in *Findings of the Association for Computational Linguistics: NAACL 2022*, 2022, pp. 2572–2580.