

Zero-shot Temporal Relation Extraction with ChatGPT

Chenhan Yuan, Qianqian Xie, Sophia Ananiadou

Department of Computer Science, The University of Manchester
{chenhan.yuan, qianqian.xie, sophia.ananiadou}@manchester.ac.uk

Abstract

The goal of temporal relation extraction is to infer the temporal relation between two events in the document. Supervised models are dominant in this task. In this work, we investigate ChatGPT’s ability on zero-shot temporal relation extraction. We designed three different prompt techniques to break down the task and evaluate ChatGPT. Our experiments show that ChatGPT’s performance has a large gap with that of supervised methods and can heavily rely on the design of prompts. We further demonstrate that ChatGPT can infer more small relation classes correctly than supervised methods. The current shortcomings of ChatGPT on temporal relation extraction are also discussed in this paper. We found that ChatGPT cannot keep consistency during temporal inference and it fails in actively long-dependency temporal inference.

1 Introduction

The temporal relation extraction task aims to extract the temporal relation between either two event triggers in the given document (Dligach et al., 2017a). In this way, a timeline of events in the document can be constructed. It is a crucial task for many downstream NLP tasks, such as natural language understanding (Mani et al., 2006; Paul et al., 2017), storyline construction (Do et al., 2012; Minard et al., 2015), and temporal question answering (Jia et al., 2018b,a), etc. Conventionally, recent temporal relation extraction (RE) models are fine-tuned based on pre-trained language models (PLMs), such as BERT and RoBERTa (Devlin et al., 2019; Liu et al., 2019). On top of the PLMs, complex neural networks are applied to classify the temporal relations, such as self-attention (Lin et al., 2019; Ning et al., 2018a), graph convolutional networks (GCNs) (Mathur et al., 2021), and policy network (Man et al., 2022). Most well-performed temporal relation extractors are supervised models, that is, they heavily rely on annotated training documents first before extracting temporal relations on

the testing set. However, annotating the temporal relations in training documents requires much domain experts’ efforts (Naik et al., 2019; Ning et al., 2018b), which is a high cost.

Different from supervised learning methods, zero-shot learning (ZSL) (Xian et al., 2017) aims to train the model that can be generalized to unseen data without annotated training data and has attracted much attention in recent years. Most recently, large language models (LLMs) (Brown et al., 2020; Bubeck et al., 2023) such as ChatGPT¹ have exhibited remarkable ability in zero-shot learning on various natural language processing (NLP) and medical tasks (Bang et al., 2023), such as information extraction (Wei et al., 2023), machine translation (Jiao et al., 2023), summarization evaluation (Luo et al., 2023), and mental health detection (Yang et al., 2023). However, the performance of LLMs in detecting temporal relations between events are not explored yet. Therefore, it is an urgent and spontaneous question if the LLM can perform zero-shot temporal relation extraction tasks well given a proper prompt approach, and if LLM can be the new paradigm for temporal RE.

In this paper, we explore the performance of ChatGPT on the zero-shot temporal relation extraction, and propose three different prompt strategies to interact with ChatGPT. More specifically, we start with the simple zero-shot prompt that directly requires ChatGPT to infer the temporal relation given the document. Then, we design the event ranking prompt, where ChatGPT is asked to infer the events shown in the given document instead of inferring temporal relations. Finally, we propose the chain-of-thought (CoT) prompt (Wei et al., 2022) to break down the task into two-stage, which guides ChatGPT to make temporal relation reasoning step by step. Based on our experimental results and analysis, we have the following findings:

¹<https://openai.com/blog/chatgpt>

- **Overall Performance.** ChatGPT significantly underperforms advanced supervised methods and even traditional neural network methods such as Bi-LSTM, indicating the challenge of temporal relation detection with ChatGPT without task-specific fine-tuning.
- **Prompts.** Similar to the finding of recent efforts, the CoT prompt can significantly improve ChatGPT’s performance compared with other prompts across all datasets, indicating the importance of proper prompt engineering.
- **Limitations.** Compared with supervised methods, ChatGPT has better performance in the temporal relations with small proportions in the dataset. However, it is also found to have limitations in detecting long-dependency temporal relation extraction and inconsistent temporal relation inference.

2 Related Work

2.1 Temporal Relation Extraction

Several studies have explored the use of temporal information in relation extraction. For example, Han et al. (Han et al., 2019b,a) incorporated tense information and timex temporal interactions into their models. Other researchers have proposed using graph neural networks (GNNs) to encode dependency structures, which are important for temporal relation extraction (Mathur et al., 2021; Schlichtkrull et al., 2018). Wang et al. (Wang et al., 2022b) added an attention layer to an R-GCN-based model to focus on document-creation-time (DCT). Man et al. (Man et al., 2022) used a reinforcement learning framework to select optimized sentences for input into neural models, improving performance.

In the clinical domain, Leeuwenberg et al. (Leeuwenberg and Moens, 2017) applied integer linear programming constraints to learn structured temporal relations. Dligach et al. (Dligach et al., 2017b) improved performance by using neural networks such as CNN and LSTM as the backbone model. Lin et al. (Lin et al., 2019) utilized pre-trained language models like BERT to learn contextualized embeddings. However, no work has yet explored the feasibility of using LLMs in temporal relation extraction.

2.2 Zero-shot Learning with ChatGPT

Since it was launched, ChatGPT has drawn much attention to its strong ability for various NLP tasks. In the clinical domain, Tang et al. explored ChatGPT’s ability on zero-shot named entity recognition and relation extraction (Tang et al., 2023). The experiments on NCBI Disease and BC5CDR Chemical datasets showed that ChatGPT cannot recognize named entities correctly in the clinical domain as the F1 drops around 55.94%-91.58% compared to SOTA-supervised methods. ChatGPT’s poor clinical NER was also proved by testing in i2b2 dataset (Hu et al., 2023). However, it can achieve comparable performance in clinical relation extraction as the F1 score only decreased by 4.73%-10.93% (Tang et al., 2023; Agrawal et al., 2022).

In extraction-related tasks, some work evaluated ChatGPT’s ability in event extraction, general information extraction, and relation extraction (Tang et al., 2023; Gao et al., 2023; Wei et al., 2023). The evaluation process follows multi-stage interactions/conversations with ChatGPT and guides it to produce the desired answers. The evaluation results of these studies showed that with proper prompting, ChatGPT can achieve comparable performance with the supervised methods on zero-shot or few-shot settings of extraction tasks. However, there is also some work pointing out that ChatGPT’s ability is still limited in some specific extraction scenarios such as extracting clinical notes with privacy information masked (Tang et al., 2023). Some work also discussed that for non-digital-available texts, such as historical documents, the entity recognition was performed poorly by ChatGPT (González-Gallardo et al., 2023; Borji, 2023).

3 ChatGPT for Zero-shot Temporal Relation Extraction

Given the input document, the temporal relation extraction aims to identify the temporal relation between any two event triggers in the document, which is modeled as the multi-classification problem. We propose three different prompts methods to evaluate ChatGPT’s performance on temporal relation extraction as shown in Figure 1.

3.1 Zero-shot Prompt

In this prompt, given the document $D = \{x_1, x_2, x_3, \dots, x_n\}$, we first label the event triggers with `<EVENT></EVENT>`. That is






Zero-shot	Event Ranking	CoT
Input document D = The season <EVENT e3>started</EVENT> about a month earlier than usual, <EVENT e4>sparking</EVENT> concerns it might <EVENT e5>turn</EVENT> into the worst in a decade. It <EVENT e6>ended</EVENT> up being very hard on the elderly, but was moderately severe overall, <EVENT e7>according</EVENT> to the Centers for Disease Control and Prevention.		
<p>Prompt: Given the document D and a list of temporal relations [before, after, vague, equal] and event triggers that are labeled as [EVENT]/[EVENT], what is the temporal relation between <EVENT e3> started </EVENT> and <EVENT e6> ended</EVENT>? Answer vague if unsure. Keep the answer short and concise</p> <p> : equal.</p> <p>Return EQUAL</p>	<p>Prompt: Given the document D and event triggers that are labeled as [EVENT]/[EVENT], which event triggers happened before <EVENT e6>ended</EVENT> ?</p> <p> : <Event e3>started</Event> and <Event e4>sparking</Event> happened before <Event e6>ended</Event></p> <p>Return [(e3, e6), (e4, e6)] BEFORE</p>	<p>Prompt: Given the document D, are <EVENT e3>started</EVENT> and <EVENT e6>ended</EVENT> referring to the same event? Keep the answer short and concise.</p> <p> : Yes.</p> <p>Prompt: Did <EVENT e3>started</EVENT> and <EVENT e6>ended</EVENT> simultaneously happened in that event? Keep the answer short and concise.</p> <p> : No.</p> <p>Prompt: "Is <EVENT e3>started</EVENT> before <EVENT e6>ended</EVENT> in that event?"</p> <p> : Yes.</p> <p>Return BEFORE</p>

Figure 1: The proposed prompts.

if x_i is an event trigger, we label it as <EVENT> x_i </EVENT>. Then the whole labeled document is sent to ChatGPT and we ask ChatGPT to find the temporal relations between either two events. Note that our goal is to test ChatGPT’s zero-shot ability in the temporal relation extraction task. Therefore, we do not provide any examples in our prompts. As shown in Figure. 1, we give ChatGPT the whole document, the list of all temporal relations, and the annotations of event triggers. In the end, the input question is designed as “what is the temporal relation between <EVENT> x_i </EVENT> and <EVENT> x_j </EVENT>?”. We let ChatGPT to answer the question by using the temporal relations provided in the list.

3.2 Event Ranking Prompt

Further, we design a new prompt to make the task easier to learn for ChatGPT. Specifically, given the document D and one event trigger e_i with the form <EVENT> x_i </EVENT> and temporal relation set R , we require ChatGPT to complete the query $(e_i, r_j, ?) \forall r_j \in R$. In this way, instead of querying $(e_i, ?, e_j)$ as in the previous prompt, ChatGPT is required to predict the missing event trigger. As event triggers are already shown in the given text/document, ChatGPT is more likely to infer the event triggers than unseen temporal relations as they are not explicitly provided in the context. In detail, as shown in Fig. 1, we achieve this by asking the question first, such as “Which events happened

before e_i ?”. Then based on the feedback of ChatGPT, we form the $\langle e_i, r, e_j \rangle$ triplets to perform the evaluation. Note that if the same event pair is detected in different temporal relation (>2), we denote that this event pair has a “vague” relation as ChatGPT cannot confidently determine which temporal relation the event pair is classified. We also asked ChatGPT to provide some concise prompts to perform temporal relation extraction tasks. As shown in Fig 2, the prompts provided by ChatGPT are in line with our event ranking prompt approach.

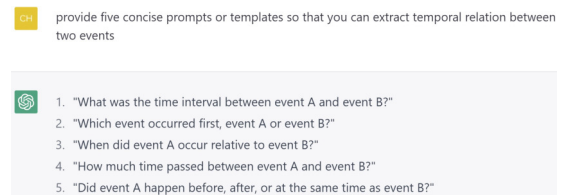


Figure 2: Prompts generated by ChatGPT.

3.3 Chain-of-thought Prompt

We notice that if two event triggers refer to the same event but point to different timestamps of the duration of that event, ChatGPT cannot distinguish them. For example, in Figure. 1, “<EVENT e3> started” and “<EVENT e6> ended” referred to the beginning and the end of the season event. ChatGPT assumes that these two event triggers happened at the same time if we directly ask about the temporal relation following the zero-shot prompt.

Therefore, we propose a new chain-of-thought prompt with two steps, which firstly navigates ChatGPT to distinguish event triggers referring to the same event, and then guides ChatGPT to infer their temporal relation. Specifically, given the document D and two event triggers e_1 and e_2 , we first ask ChatGPT to determine if e_1 and e_2 refer to the same event. If they are not, we further ask ChatGPT to determine the temporal relation between the two event triggers. If they point to the same event, we use a similar prompt but with the extra phrase “in that event” to ask ChatGPT. As shown in Fig. 1, we first ask ChatGPT if the two event triggers `<EVENT>started</EVENT>` and `<EVENT>ended</EVENT>` refer to the same event. Then based on ChatGPT’s feedback, we further iteratively go through the whole temporal relation list to determine which temporal relation exists between the two event triggers.

4 Experiments

4.1 Datasets

We use three datasets to evaluate the performance of ChatGPT on the zero-shot temporal relation extraction. The statistical details of these datasets are shown in the Table 1.

Dataset	Train	Dev	Test	Labels
TB-Dense	4,032	629	1,427	6
MATRES	6,336	–	837	4
TDDMan	4,000	650	1,500	5

Table 1: Statistics of the number of annotated event pairs and different temporal relation labels of the MATRES, TB-Dense and TDDMan datasets.

- **TimeBank-Dense** (Cassidy et al., 2014) labeled 36 news documents in total. The temporal relations between event triggers-event triggers, timex-event triggers, timex-timex, are labeled. Following previous work, we only test our model on the event trigger-event trigger relation in the testing set.
- **MATRES** (Ning et al., 2018b) is a dataset primarily focusing on temporal relations of event triggers with local sentences (1 or 2 sentences).
- **TDDiscourse** (Naik et al., 2019) was created to explicitly emphasize global discourse-level

temporal ordering. Based on the annotation accuracy, the dataset is split into TDDMan and TDDAuto, where TDDAuto introduces much more automatic labels and noise. In this paper, we only evaluate ChatGPT on TDDMan because of the budget limitation.

We only use the testing set of each dataset directly to test our approaches as we do not require training of ChatGPT, following the zero-shot setting. We report the F1 score on each dataset and each temporal relation.

4.2 Baseline Models

Since there is no zero-shot learning methods for temporal relation extraction before, we compare the performance of ChatGPT with the following advanced supervised methods:

- CAEVO (Chambers et al., 2014) a sieve-based architecture that includes multi-level classifiers for intra-sentence temporal relation learning.
- SP+ILP (Ning et al., 2017) a structured learning approach that captures the global temporal features when inferring the relation between two local events.
- Bi-LSTM (Cheng and Miyao, 2017) a Bi-LSTM-based model that encodes the dependency path between two events to classify temporal relation.
- Joint (Han et al., 2019b) and Deep (Han et al., 2019a) models that utilized SSVM as the scoring function to learn temporal constraints and context embeddings.
- UCGraph (Liu et al., 2021) a graph-based model that is trained with mask pre-training mechanism. The model’s uncertainty level is used to guide the inference during testing.
- TIMERS (Mathur et al., 2021) a graph-based model that leverages three graphs to learn temporal, rhetorical, and syntactic information, respectively.
- SCS-EERE (Man et al., 2022) a reinforcement learning-based selector is designed to select the optimized sentences for temporal inference between the given two events.

- RSGT (Zhou et al., 2022) a syntactic-and-semantic-based graph model pre-trained on a temporal neighbor prediction task.
- FaithTRE (Wang et al., 2022a) a model that applied Dirichlet prior to estimating the correctness likelihood. A temperature scaling is also used to recalibrate the model confidence measure after bias mitigation.
- DTRE (Wang et al., 2022b) a document creation time(DCT)-aware graph with a global consistency mechanism when inferring temporal relations.
- MulCo (Yao et al., 2022) a joint model using the BERT to learn contextualized features and GNN to capture syntactic structures. The two models are combined via a multi-level contrastive learning framework.

4.3 Results

In table 2, we can see that ChatGPT struggles to outperform supervised state-of-the-art models such as SCS-EERE (Man et al., 2022) and RSGT (Zhou et al., 2022), and even the traditional neural networks methods such as CAEVO and BI-LSTM, indicating its ineffective in the temporal relation extraction task in the zero-shot setting. Table 2 also shows the performance of ChatGPT under different prompts in three datasets. We noticed that ChatGPT_ER yields the worst performance on the MATRES and TDDMan datasets. ChatGPT’s performance with the event ranking prompt on TDD-Man dataset is poor as most event trigger pairs cannot be detected under this prompt. However, if the event trigger pairs are explicitly fed to ChatGPT, it can somehow partially infer the temporal relations correctly. For example, the zero-shot prompt and CoT prompt could improve the F1 score by 14.8% and 23.8%, respectively. While for the TB-Dense dataset, ChatGPT_ZS has the worst performance. ChatGPT_CoT achieves the best performance across all datasets, such as it significantly outperforms ChatGPT_ZS and ChatGPT_ER by 27.1% and 33.1%. This illustrates the effectiveness of the CoT prompt with step-by-step guidance in prompting ChatGPT.

Table 3, Table 4 and Table 5 further list the detail results of ChatGPT with different prompts on three datasets. We can see that the performance with the event ranking prompt is much better than that on other datasets. ChatGPT with the zero-shot

prompt cannot determine the temporal relation “is included” and therefore yields a 0.0 performance in this type of relation. The CoT prompt improves the overall performance by significantly detecting “before” and “after” temporal relations. As these two relations take a great portion of the whole dataset, the overall performance is also improved.

5 Discussion

5.1 ChatGPT is slightly better on small temporal relation classes

The imbalanced data is a severe long-existing problem in the temporal relation extraction task. Because of the events’ temporal order frequency in real life, some temporal relations, such as “simultaneous” and “equal”, are very limited in most temporal relation extraction datasets. And popular NLP data-augmented methods are difficult to be applied in the temporal domain. Therefore, most state-of-the-art supervised methods yield much worse performance on small relation classes.

In Deep (Han et al., 2019a), the authors reported detailed performance on each temporal class in MATRES dataset. We therefore especially compare the performance on small relation classes against Deep, i.e., “EQUAL” and “VAGUE”. As shown in Table 3, the supervised model Deep could achieve much better overall performance (81.7 F1 score), due to the contribution of two majority relations, “before” and “after”. Compared to Deep’s 0.0 performance on “EQUAL” and “VAGUE”, ChatGPT with event ranking, CoT and zero-shot prompts can correctly extract some small class relations.

5.2 ChatGPT failed in actively long-dependency temporal relation extraction

As shown in Table 2, ChatGPT’s performance drops a lot in the TDDMan dataset. We argue that this is mainly because the TDDMan focuses more on discourse-level temporal relations and ChatGPT failed to extract useful information from long documents. As shown in Table 4, the event ranking prompt yields almost 0.0 on the whole dataset. In practice, we initially input the whole document D to ChatGPT and ask ChatGPT which event triggers in the document D happened before the given event trigger e_1 . That is, ChatGPT should actively search all event triggers in the document and produce answers. However, in the TDDMan dataset, ChatGPT cannot produce a formatted answer and the outputs

Models	MATRES			TDDMan			TB-Dense		
	prec	recall	F1	prec	recall	F1	prec	recall	F1
CAEVO (Chambers et al., 2014)	–	–	–	32.3	10.7	16.1	49.9	46.6	48.2
SP+ILP (Ning et al., 2017)	71.3	82.1	76.3	23.9	23.8	23.8	58.4	58.4	58.4
Bi-LSTM (Cheng and Miyao, 2017)	59.5	59.5	59.5	24.9	23.8	24.3	63.9	38.9	48.4
Joint (Han et al., 2019b)	–	–	75.5	41.0	41.1	41.1	–	–	64.5
Deep (Han et al., 2019a)	77.4	86.4	81.7	–	–	–	62.7	58.9	62.5
UCGraph (Liu et al., 2021)	–	–	–	44.5	42.3	43.4	62.4	56.1	59.1
TIMERS (Mathur et al., 2021)	81.1	84.6	82.3	43.7	46.7	45.5	48.1	65.2	67.8
SCS-EERE (Man et al., 2022)	78.8	88.5	83.4	–	–	51.1	–	–	–
FaithTRE (Wang et al., 2022a)	–	–	82.7	–	–	52.9	–	–	–
RSGT (Zhou et al., 2022)	82.2	85.8	84.0	–	–	–	68.7	68.7	68.7
DTRE (Wang et al., 2022b)	–	–	–	56.3	56.3	56.3	–	–	70.2
MulCo (Yao et al., 2022)	88.2	88.2	88.2	56.2	54.0	55.1	84.9	84.9	84.9
ChatGPT_ZS	26.4	24.3	25.3	17.7	13.6	15.3	23.7	14.3	17.8
ChatGPT_ER	21.9	17.3	19.3	3.7	0.3	0.5	37.6	35.8	36.6
ChatGPT_CoT	48.0	57.7	52.4	26.8	22.3	24.3	43.4	32.2	37.0

Table 2: The comparison of ChatpGPT with various prompt techniques and supervised state-of-the-art models.

Relation	Zero-shot			CoT			Event ranking			Deep		
	prec	recall	F1	prec	recall	F1	prec	recall	F1	prec	recall	F1
overall	26.4	24.3	25.3	48.0	57.7	52.4	21.9	17.3	19.3	77.4	86.4	81.7
EQUAL	0.0	0.0	0.0	7.1	2.9	4.1	5.8	11.1	7.6	0.0	0.0	0.0
VAGUE	14.3	58.7	23.1	14.4	8.1	10.4	14.6	86.2	25.0	0.0	0.0	0.0
AFTER	34.0	25.6	29.2	41.6	41.8	41.7	36.4	1.6	3.0	72.3	84.8	78.0
BEFORE	52.5	17.8	26.6	63.1	71.6	67.1	57.0	13.0	21.1	80.1	89.6	84.6

Table 3: The zero-shot performance of ChatGPT with three different prompts on the MATRES dataset.

sometimes are even some random words in the document instead of event triggers. We then test the limit of the size of the input document, i.e., the number of sentences. We finally found that if we limit the size of the input document to at most 8 context sentences around the event triggers, ChatGPT would be more stable to produce formatted answers instead of repeating part of the document randomly. However, in this way, most temporal relations extracted by ChatGPT do not match with the golden labels in the TDDMan dataset because the extracted temporal relations are only in short dependency while TDDMan emphasizes long-distance temporal dependency.

Nevertheless, surprisingly, if we explicitly ask ChatGPT what the temporal relations between two event triggers labeled in the document are, ChatGPT can answer some of them correctly. Note that we do not cut the size of the document in this case

and ChatGPT can still learn some of the temporal dependency in the document. Compared to the event ranking prompt, ChatGPT passively receives two event trigger information with the zero-shot and CoT prompts. This may reduce the inference difficulty as more information is given.

5.3 ChatGPT can be improved via multi-stage “yes” or “no” prompts

Intuitively, the most efficient way to query ChatGPT about temporal relations between two events in the given document should be the zero-shot prompt, which directly asks the ChatGPT to answer the temporal relation between any two event triggers. If only one event trigger is given, then the prompts provided by ChatGPT, i.e., event ranking, should be used to interact with ChatGPT. However, our extensive experiments show that these two prompt methods produce much worse perfor-

Relation	Zero-shot			CoT			Event ranking		
	prec	recall	F1	prec	recall	F1	prec	recall	F1
overall	17.7	13.6	15.3	26.8	22.3	24.3	3.7	0.3	0.5
is included	9.5	0.7	1.3	20.9	3.1	5.4	0.0	0.0	0.0
include	41.9	17.7	24.8	37.9	11.2	17.3	0.0	0.0	0.0
after	14.7	9.0	11.2	33.3	4.3	7.5	0.0	0.0	0.0
before	29.7	22.9	25.9	35.1	70.8	46.9	12.5	0.7	1.4
simultaneous	3.9	39.1	7.0	0.0	0.0	0.0	11.1	2.2	3.6

Table 4: The zero-shot performance of ChatGPT with three prompts on the TDD-Man dataset.

Relation	Zero-shot			CoT			Event ranking		
	prec	recall	F1	prec	recall	F1	prec	recall	F1
overall	23.7	14.3	17.8	43.4	32.2	37.0	37.6	35.8	36.6
is included	0.0	0.0	0.0	10.0	1.9	3.2	6.2	3.8	4.7
include	3.3	10.7	24.8	5.5	16.1	8.2	16.7	5.4	8.1
after	29.0	17.2	11.2	70.4	13.9	23.2	19.0	1.5	2.7
before	40.0	9.9	25.9	35.0	75.5	47.9	31.2	25.3	27.9
simultaneous	1.5	45.5	3.0	33.3	4.5	8.0	6.7	50.0	11.8
vague	44.6	24.0	31.2	51.2	29.6	37.5	46.0	63.5	53.4

Table 5: The zero-shot performance of ChatGPT with three prompts on the TimeBank-Dense dataset.

mance in most cases compared to the CoT prompt. Note that both zero-shot and CoT provide sufficient information about event triggers. We argue there is another difference resulting in the performance gap.


Comparing the two prompts, one significant difference is that the CoT prompt only accepts “yes” or “no” answers while the zero-shot prompt returns a specific temporal relation label. In the zero-shot prompt, ChatGPT is required to select a temporal relation from the given list, which is similar to a conventional multi-class classification problem. However, in the CoT prompt, ChatGPT only has to determine if one specific temporal relation exists (or not) between two event triggers. This simplified the problem into a binary classification. Further, with the previous question-answer pair as context, ChatGPT has a higher probability of making the correct selection. For example, in Fig 1, in the first round, ChatGPT already inferred that the relation is not “EQUAL”. Then in the second round, ChatGPT is more confident to predict the temporal relation from $[BEFORE, AFTER, VAGUE]$ instead of “EQUAL”.

5.4 ChatGPT’s temporal inference is inconsistent even with sufficient context


During the testing of ChatGPT with the event ranking prompt, we noticed a fatal issue in ChatGPT’s temporal relation extraction, namely the inconsistent temporal relation inference. Given the same input document, ChatGPT may produce different temporal relations between two event triggers. As the top example shown in Figure 3, given the document D , if the prompt is “Which event triggers happened *before* e_1 ?”, ChatGPT will give a list of event triggers, e.g. $[e_2, e_5, e_6]$. Now given the same document D , if the prompt is “Which event triggers happened *after* e_6 ”, ChatGPT is expected to at least include e_1 in the list. However, during the experiments, we noticed that ChatGPT failed in this scenario multiple times and the failures can be categorized into two cases. The first failure case is that ChatGPT does not include e_1 in any list associated with e_6 . The second case is that ChatGPT includes e_1 in a wrong list, e.g., “EQUAL”, associated with e_6 which violates the temporal consistency.

A similar problem also happened to the CoT prompt. As shown in Figure 4, during the experiments, ChatGPT would give various answers instead of “vague” as we specified, if it thinks there


Prompt: Given the document "Earlier AT and T also [EVENT e7]announced[/EVENT] an eighteen percent [EVENT e42]drop[/EVENT] in profits for the fourth quarter. The company's sales force [EVENT e9]applauded[/EVENT] the shake up. It's not something we're not [EVENT e206]used[/EVENT] to in the industry." and event triggers that are labeled as [EVENT]/[EVENT], which event triggers happened before [EVENT e9]applauded[/EVENT]? keep the answer short and concise

 : Two event triggers happened before "applauded": [EVENT e7]announced[/EVENT] and [EVENT e42]drop[/EVENT]

Prompt: which event triggers happened after [EVENT e42]drop[/EVENT]?

 : The events that happened after the [EVENT e42]drop[/EVENT] are not specified or labeled in the given document.

Prompt: Given the document "so, Yeah I [EVENT e48]think[/EVENT] he's [EVENT e24]doing[/EVENT] the right thing. I [EVENT e25]think[/EVENT] it's a good thing that they're finally going to be [EVENT e27]downsizing[/EVENT] you know some management because there is a whole lot of waste." and event triggers that are labeled as [EVENT]/[EVENT], which event triggers occurred during the period of [EVENT e25]think[/EVENT]? keep the answer short and concise

 : The event trigger labeled as [EVENT e24]doing[/EVENT], and [EVENT e27]downsizing[/EVENT]

Prompt: which event triggers happened before [EVENT e27]downsizing[/EVENT]? keep the answer short and concise



 : The event triggers that happened before are [EVENT e48]think[/EVENT], [EVENT e24]doing[/EVENT], and [EVENT e25]think[/EVENT].

Figure 3: ChatGPT’s two temporal inconsistency cases examples in the Event ranking prompt.


Prompt: Given the document "Meanwhile, secretary of state Madeleine Albright, Berger and defense secretary William Cohen [EVENT e80]announced[/EVENT] [EVENT e144]plans[/EVENT] to [EVENT e82]travel[/EVENT] to an unnamed city in the US heartland next week, to [EVENT e84]explain[/EVENT] to the American people just why military force will be [EVENT e146]necessary[/EVENT] if diplomacy [EVENT e87]fails[/EVENT]. Steve Hurst, CNN, the state department."

Did [EVENT e146]necessary[/EVENT] happen before [EVENT e87]fails[/EVENT]? Answer vague if unsure and keep the answer short

 : The document does not provide information on the order of events between "[EVENT e146]necessary[/EVENT]" and "[EVENT e87]fails[/EVENT]."

Prompt: Given the document " The last surviving member of the team which first [EVENT 1]conquered[/EVENT] Everest in 1953 has [EVENT 2]died[/EVENT] in a Derbyshire nursing home. George Lowe, 89, [EVENT 3]died[/EVENT] in Ripley on Wednesday after a long-term [EVENT 45]illness[/EVENT], with his wife Mary by his side. New Zealand-born Mr Lowe was part of the team that [EVENT 4]helped[/EVENT] Sir Edmund Hillary and Tenzing Norgay to the summit in 1953."

Did [EVENT 3]died[/EVENT] happen after [EVENT 4]helped[/EVENT]? Answer vague if unsure and keep the answer short

 : Not specified.

Prompt: Did [EVENT 3]died[/EVENT] happen before [EVENT 4]helped[/EVENT]? Answer vague if unsure and keep the answer short


 : Yes.

Figure 4: ChatGPT’s inconsistency failures examples in the CoT prompt. The top shows ChatGPT does not reply “vague” when unsure and the bottom shows that ChatGPT still infers other inconsistent relations.

is no clue to infer the temporal relation between the given two event triggers. These answers include “Cannot determine.”, “I cannot answer that question as it is unclear from the given information.”, “Unknown.”, etc. We treat all of these answers as “vague” when we evaluate ChatGPT. Intuitively, if ChatGPT cannot answer “yes” or “no” for the specific temporal relation between two event triggers, it should also not be able to answer questions about other temporal relations of the same event triggers. However, ChatGPT may violate its “unknown” statement in the multi-stage prompts. We test if ChatGPT’s inconsistency exists in multi-stage prompts by implementing the following experiments. For example, during the i -th round, if we ask “Did e_1 happen before e_2 ?” and ChatGPT’s answer includes “It is unclear from the given information.”. We then ask the next $i + 1$ question “Did e_1 happen after e_2 ?”. And surprisingly, in most cases (84%), ChatGPT can classify the event trigger pair into other temporal relation classes even if it claims that the information is insufficient. Moreover, 96% of these classification results is incorrect.

6 Conclusion

In this work, we comprehensively test ChatGPT’s zero-shot ability on temporal relation extraction. We designed three different prompts to evaluate ChatGPT’s performance. Our experiments demonstrate that with proper prompting, ChatGPT’s performance on zero-shot temporal relation extraction can be significantly improved, highlighting the importance of prompt engineering to better trigger ChatGPT’s ability in future work. However, compared to supervised methods, ChatGPT is still far behind. We further discuss our findings from experimental results, including its better performance in small classes than supervised methods and its drawbacks such as failures in long-distance temporal dependency inference and inconsistent temporal relation inference. Our work only takes the initial step of exploring the LLMs for zero-shot temporal relation extraction. To fill the gap between the performance of LLMs in the zero-shot setting and that of advanced supervised methods, we believe more efforts should be explored in the future.

7 Acknowledgements

This work is supported by the project JPNP20006 from New Energy and Industrial Technology Development Organization (NEDO).

References

- Monica Agrawal, Stefan Hegselmann, Hunter Lang, Yoon Kim, and David Sontag. 2022. Large language models are zero-shot clinical information extractors. *arXiv preprint arXiv:2205.12689*.
- Yejin Bang, Samuel Cahyawijaya, Nayeon Lee, Wenliang Dai, Dan Su, Bryan Wilie, Holy Lovenia, Ziwei Ji, Tiezheng Yu, Willy Chung, et al. 2023. A multi-task, multilingual, multimodal evaluation of chatgpt on reasoning, hallucination, and interactivity. *arXiv preprint arXiv:2302.04023*.
- Ali Borji. 2023. A categorical archive of chatgpt failures. *arXiv preprint arXiv:2302.03494*.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. 2023. Sparks of artificial general intelligence: Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712*.
- Taylor Cassidy, Bill McDowell, Nathanael Chambers, and Steven Bethard. 2014. An annotation framework for dense event ordering. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 501–506.
- Nathanael Chambers, Taylor Cassidy, Bill McDowell, and Steven Bethard. 2014. Dense event ordering with a multi-pass architecture. *Transactions of the Association for Computational Linguistics*, 2:273–284.
- Fei Cheng and Yusuke Miyao. 2017. Classifying temporal relations by bidirectional lstm over dependency paths. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 1–6.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186.
- Dmitriy Dligach, Timothy Miller, Chen Lin, Steven Bethard, and Guergana Savova. 2017a. **Neural temporal relation extraction**. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, pages 746–751, Valencia, Spain. Association for Computational Linguistics.
- Dmitriy Dligach, Timothy Miller, Chen Lin, Steven Bethard, and Guergana Savova. 2017b. Neural temporal relation extraction. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, pages 746–751.
- Quang Do, Wei Lu, and Dan Roth. 2012. Joint inference for event timeline construction. In *Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning*, pages 677–687.
- Jun Gao, Huan Zhao, Changlong Yu, and Ruifeng Xu. 2023. Exploring the feasibility of chatgpt for event extraction. *arXiv preprint arXiv:2303.03836*.
- Carlos-Emiliano González-Gallardo, Emanuela Boros, Nancy Girdhar, Ahmed Hamdi, Jose G Moreno, and Antoine Doucet. 2023. Yes but.. can chatgpt identify entities in historical documents? *arXiv preprint arXiv:2303.17322*.
- Rujun Han, I-Hung Hsu, Mu Yang, Aram Galstyan, Ralph Weischedel, and Nanyun Peng. 2019a. Deep structured neural network for event temporal relation extraction. In *Proceedings of the 23rd Conference on Computational Natural Language Learning (CoNLL)*, pages 666–106.
- Rujun Han, Qiang Ning, and Nanyun Peng. 2019b. Joint event and temporal relation extraction with shared representations and structured prediction. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 434–444.
- Yan Hu, Iqra Ameer, Xu Zuo, Xueqing Peng, Yujia Zhou, Zehan Li, Yiming Li, Jianfu Li, Xiaoqian Jiang, and Hua Xu. 2023. Zero-shot clinical entity recognition using chatgpt. *arXiv preprint arXiv:2303.16416*.
- Zhen Jia, Abdalghani Abujabal, Rishiraj Saha Roy, Jan-nik Strötgen, and Gerhard Weikum. 2018a. Tempquestions: A benchmark for temporal question answering. In *Companion Proceedings of the The Web Conference 2018*, pages 1057–1062.
- Zhen Jia, Abdalghani Abujabal, Rishiraj Saha Roy, Jan-nik Strötgen, and Gerhard Weikum. 2018b. Tequila: Temporal question answering over knowledge bases. In *Proceedings of the 27th ACM international conference on information and knowledge management*, pages 1807–1810.

- Wenxiang Jiao, Wenxuan Wang, Jen-tse Huang, Xing Wang, and Zhaopeng Tu. 2023. Is chatgpt a good translator? a preliminary study. *arXiv preprint arXiv:2301.08745*.
- Artuur Leeuwenberg and Marie Francine Moens. 2017. Structured learning for temporal relation extraction from clinical records. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers*, pages 1150–1158.
- Chen Lin, Timothy Miller, Dmitriy Dligach, Steven Bethard, and Guergana Savova. 2019. A bert-based universal model for both within-and cross-sentence clinical temporal relation extraction. In *Proceedings of the 2nd Clinical Natural Language Processing Workshop*, pages 65–71.
- Jian Liu, Jinan Xu, Yufeng Chen, and Yujie Zhang. 2021. Discourse-level event temporal ordering with uncertainty-guided graph completion. In *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence, IJCAI-21*, pages 3871–3877. International Joint Conferences on Artificial Intelligence Organization. Main Track.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.
- Zheheng Luo, Qianqian Xie, and Sophia Ananiadou. 2023. Chatgpt as a factual inconsistency evaluator for abstractive text summarization. *arXiv preprint arXiv:2303.15621*.
- Hieu Man, Nghia Trung Ngo, Linh Ngo Van, and Thien Huu Nguyen. 2022. Selecting optimal context sentences for event-event relation extraction. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 11058–11066.
- Inderjeet Mani, Marc Verhagen, Ben Wellner, Chungmin Lee, and James Pustejovsky. 2006. Machine learning of temporal relations. In *Proceedings of the 21st International Conference on Computational Linguistics and 44th Annual Meeting of the Association for Computational Linguistics*, pages 753–760.
- Puneet Mathur, Rajiv Jain, Franck Dernoncourt, Vlad Morariu, Quan Hung Tran, and Dinesh Manocha. 2021. Timers: document-level temporal relation extraction. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 524–533.
- Anne-Lyse Myriam Minard, Manuela Speranza, Eneko Agirre, Itziar Aldabe, Marieke van Erp, Bernardo Magnini, German Rigau, and Ruben Urizar. 2015. Semeval-2015 task 4: Timeline: Cross-document event ordering. In *proceedings of the 9th International Workshop on Semantic Evaluation (SemEval 2015)*, pages 778–786.
- Aakanksha Naik, Luke Breitfeller, and Carolyn Rose. 2019. Tddiscourse: A dataset for discourse-level temporal ordering of events. In *Proceedings of the 20th Annual SIGdial Meeting on Discourse and Dialogue*, pages 239–249.
- Qiang Ning, Zhili Feng, and Dan Roth. 2017. A structured learning approach to temporal relation extraction. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 1027–1037.
- Qiang Ning, Hao Wu, Haoruo Peng, and Dan Roth. 2018a. Improving temporal relation extraction with a globally acquired statistical resource. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 841–851.
- Qiang Ning, Hao Wu, and Dan Roth. 2018b. A multi-axis annotation scheme for event temporal relations. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1318–1328.
- Rohan Paul, Andrei Barbu, Sue Felshin, Boris Katz, and Nicholas Roy. 2017. Temporal grounding graphs for language understanding with accrued visual-linguistic context. In *Proceedings of the 26th International Joint Conference on Artificial Intelligence*, pages 4506–4514.
- Michael Schlichtkrull, Thomas N Kipf, Peter Bloem, Rianne Van Den Berg, Ivan Titov, and Max Welling. 2018. Modeling relational data with graph convolutional networks. In *The Semantic Web: 15th International Conference, ESWC 2018, Heraklion, Crete, Greece, June 3–7, 2018, Proceedings 15*, pages 593–607. Springer.
- Ruixiang Tang, Xiaotian Han, Xiaoqian Jiang, and Xia Hu. 2023. Does synthetic data generation of llms help clinical text mining? *arXiv preprint arXiv:2303.04360*.
- Haoyu Wang, Hongming Zhang, Yuqian Deng, Jacob R Gardner, Muhao Chen, and Dan Roth. 2022a. Extracting or guessing? improving faithfulness of event temporal relation extraction. *arXiv preprint arXiv:2210.04992*.
- Liang Wang, Peifeng Li, and Sheng Xu. 2022b. Dct-centered temporal relation extraction. In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 2087–2097.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed H Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. In *Advances in Neural Information Processing Systems*.

- Xiang Wei, Xingyu Cui, Ning Cheng, Xiaobin Wang, Xin Zhang, Shen Huang, Pengjun Xie, Jinan Xu, Yufeng Chen, Meishan Zhang, et al. 2023. Zero-shot information extraction via chatting with chatgpt. *arXiv preprint arXiv:2302.10205*.
- Yongqin Xian, Bernt Schiele, and Zeynep Akata. 2017. Zero-shot learning-the good, the bad and the ugly. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4582–4591.
- Kailai Yang, Shaoxiong Ji, Tianlin Zhang, Qianqian Xie, and Sophia Ananiadou. 2023. [On the evaluations of chatgpt and emotion-enhanced prompting for mental health analysis](#).
- Hao-Ren Yao, Luke Breitfeller, Aakanksha Naik, Chunxiao Zhou, and Carolyn Rose. 2022. Multi-scale contrastive co-training for event temporal relation extraction. *arXiv preprint arXiv:2209.00568*.
- Jie Zhou, Shenpo Dong, Hongkui Tu, Xiaodong Wang, and Yong Dou. 2022. Rsgt: Relational structure guided temporal relation extraction. In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 2001–2010.