

Cascading Large Language Models for Salient Event Graph Generation

Anonymous ACL submission

Abstract

Generating event graphs from long documents is challenging due to the inherent complexity of multiple tasks involved such as detecting events, identifying their relationships, and reconciling unstructured input with structured graphs. Recent studies typically consider all events with equal importance, failing to distinguish salient events crucial for understanding narratives. This paper presents CALLM-SAE, a Cascading Large Language Model framework for SALient Event graph generation, which leverages the capabilities of LLMs and eliminates the need for costly human annotations. We first identify salient events by prompting LLMs to generate summaries, from which salient events are identified. Next, we develop an iterative code refinement prompting strategy to generate event relation graphs, removing hallucinated relations and recovering missing edges. Fine-tuning contextualised graph generation models on the LLM-generated graphs outperforms the models trained on CAEVO-generated data. Experimental results on a human-annotated test set show that the proposed method generates salient and more accurate graphs, outperforming competitive baselines.¹

1 Introduction

Events are fundamental discourse units which form the backbone of human communication. They are interconnected through various event relations such as hierarchical, temporal, or causal relations. Event relation graphs are vital for representing and understanding complex event narratives, with nodes representing events and edges denoting relationships between them. High-quality event relation graphs can enhance numerous downstream tasks, such as question answering (Lu et al., 2022) and reasoning (Melnik et al., 2022).

¹Source code and dataset will be released upon paper acceptance.

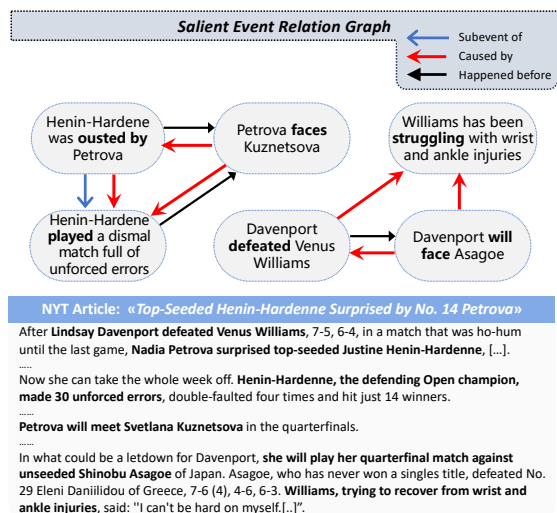


Figure 1: An example of salient event relation graph (top) generated from the NYT article (bottom).

Recent studies on contextualised event graph generation have focused on fine-tuning language models to generate linearised graphs from documents in an end-to-end manner (Madaan and Yang, 2021; Tan et al., 2024a). These methods rely on distant supervision, such as events and event temporal relations detected using an approach called CAEVO (McDowell et al., 2017), due to the data-intensive nature of language models and heavy manual efforts of annotating event graphs. However, CAEVO has limitations. It typically considers predicates (e.g., verbs) in text as events and tends to extract many insignificant events, such as “say” and “think”, which add little value to narrative understanding and have minimal connections to other events, thus introducing noise to the event graphs.

To improve the quality of distant supervision graphs, it is essential to consider the saliency of events. We found that CAEVO-extracted events often have low saliency because CAEVO takes a bottom-up approach to event extraction, classifying each predicate as an event or not. In contrast, identifying salient events requires a top-down approach.

Existing studies on identifying salient events or entities use the *summarisation test* to guide human annotation, where an event or entity is considered salient if a human-written summary is likely to include it (Dunietz and Gillick, 2014; Liu et al., 2018). Given that instruction fine-tuned LLMs perform on par with human writers in news summarisation (Zhang et al., 2024), we propose generating salient events by instructing LLMs to first summarise documents before identifying salient events.

Moreover, we extend beyond the CAEVO’s temporal-only relations to encompass multiple relation types. We introduce iterative refinement prompting in a code prompt format to generate event relation graphs that include hierarchical, temporal, and causal relations (see Figure 1). The prompting framework is highly efficient because the code prompt format generates each type of relation graph in a single pass, while the naive prompting method needs to query each possible event pair individually. The iterative refinement process further enhances the accuracy of event relation predictions by using a hallucination grader to filter out unfaithful edges and iterative generation to recover missing edges.

Using the LLM-generated dataset, we fine-tune Flan-T5 following the same method as Tan et al. (2024a). However, the abstractive nature of salient events poses challenges for evaluation, as salient events rarely exactly match the gold standards despite having the same semantic meaning. To address this, we propose an evaluation metric based on semantic text embeddings for assessing the event relation graphs. Our experimental results on the New York Times corpus (Sandhaus, 2008) show that CALLMSAE, a novel CAscading Large Language Model framework for SALient Event graph generation, outperforms the baselines in terms of event saliency and edge quality. The fine-tuned model surpasses previous models trained with CAEVO-generated graphs. Our contributions are summarised as follows:

- We propose CALLMSAE, a CAscading Large Language Model framework for SALient Event graph generation, serving as a distant signal generator for contextualised graph generation models. We also propose a novel contextualised evaluation metric for comparing salient event graphs.
- We provide a large-scale LLM-generated

salient event graph dataset (10,247 documents) with three relation types for distant supervision, along with a human-annotated test set (100 documents).

- We present an extensive experimental evaluation on LLM-generated event relation graphs in terms of event saliency and event relation on the NYT corpus, demonstrating how higher quality salient event graphs can improve contextualised graph generation.

2 Related Work

Event Relation Graph Construction The early idea of event relation graph construction comes from UzZaman et al. (2013), which introduces a dataset for evaluating an end-to-end system which takes raw text as input and output TimeML annotations (i.e., temporal relations). CAEVO (McDowell et al., 2017) and Cogcomptime (Ning et al., 2018) both utilise a wide range of manually designed features to train MaxEnt and averaged perception for extracting events and relations. Han et al. (2019b) proposed a joint event and relation extraction model based on BERT (Devlin et al., 2019) and BiLSTM (Panchendrarajan and Amareesan, 2018). Other researchers focus on developing specialised sub-systems to classify extracted event pairs for relations (Ning et al., 2019; Han et al., 2019a; Wang et al., 2020; Tan et al., 2021). ATOMIC (Sap et al., 2019) is a large-scale commonsense knowledge graph containing the causes and effects of events. MAVEN-ERE (Wang et al., 2022) is built with event coreference, temporal, causal and subevent relations. However, ATOMIC and MAVEN-ERE completely rely on crowdsourcing and thus are difficult to extend. MAVEN-ERE is less than half the size of our dataset and does not consider the saliency of events.

Madaan and Yang (2021) fine-tune GPT-2 to generate linearised graphs from documents in an end-to-end manner. Their temporal relation graphs used for training are produced by CAEVO. Following this direction, Tan et al. (2024a) instead view the task as set generation and propose a framework based on set property regularisation and data augmentation. In this paper, we focus on generating multi-relation graphs via in-context learning, prompt interaction, and iterative refinement.

Salient Event Identification Several existing papers investigate the problem of identifying salient events. Choubey et al. (2018) build a rule-based

classifier to identify central events by exploiting human-annotated event coreference relations. They find the central events either have large numbers of coreferential event mentions or have large stretch sizes. Jindal et al. (2020) propose a contextual model to identify salient events based on BERT and BiLSTM. They also mention several features, such as event trigger frequency, which are essential features to identify the salient events. Liu et al. (2018) propose a feature-based method using LeToR (Liu et al., 2009) and a neural-based method called Kernel-based Centrality Estimation. To train and evaluate their methods, they build a dataset based on the *summarisation test*: an event is considered salient if a summary written by a human is likely to include it. Zhang et al. (2021) combine the Kernel-based Centrality Estimation with the event and temporal relation extraction model of Han et al. (2019b) to build a salience-aware event chain modelling system. However, they only focus on single-dimensional chains and only model temporal relations.

3 Cascading LLMs to Generate Salient Event Graphs

CALLMSAE combines various prompts in a pipelined manner to generate salient event graphs. In this section, we will first introduce the prompts for generating salient events. Then, we will describe the method for generating relation graphs based on the salient events. Lastly, we define an evaluation metric for comparing event graphs: *Hungarian Graph Similarity*.

3.1 Generate Salient Events

The *summarisation test* (as mentioned in Section 1) is often used to guide the annotation of salient events or entities (Dunietz and Gillick, 2014; Liu et al., 2018). These studies identify events or entities included in human-written summaries as salient. Similarly, we instruct LLMs to generate a summary first and then extract events from it.

3.2 Generate Graphs as Code Completion

While LLMs can extract salient events, they often struggle with identifying event relations (Chan et al., 2023; Tan et al., 2024a). Prompt engineering for extracting event relations is complex due to the need to incorporate various terminologies and graph constraints. Moreover, prompt efficiency is crucial as generating a large-scale dataset with LLMs can still incur significant computational

costs, albeit less than crowdsourcing.

In our method, the main prompt for generating the event relation graph is formulated as a Python code completion task. The graph is defined using the NetworkX² package in Python, with nodes representing the salient events generated in Section 3.1. LLMs are instructed to complete the code by adding relation edges using NetworkX’s APIs.

Recent research suggests that formulating prompts as code can enhance LLMs’ reasoning abilities (Wang et al., 2023; Zhang et al., 2023). In our task, the Python code format effectively incorporates all necessary terminologies, enabling LLMs to understand them without confusion. The Python code format also allows for the inclusion of constraints (e.g., ensuring the graph is a directed acyclic graph) and additional instructions (e.g., ask for explanations) as comments. LLMs can generate explanations as comments without disrupting the main content of the graph. Moreover, the code template simplifies parsing the response, as LLMs are directed to use the “.add_edge()” function to add the relations.

Since hierarchical, temporal, and causal relations are asymmetric, each can be represented by a Directed Acyclic Graph (DAG). We formulate three distinct prompts to guide LLMs in generating three DAGs, each representing one of these relation types. This approach avoids the complexity of a multi-label graph, and LLMs can focus on a single relation type and carefully consider the topological structure of the graph. We can also use the “.find_cycle()” function from NetworkX to detect constraint violations reliably. In addition, if relation types are interconnected, the initially generated graphs can help the generation of subsequent graphs (as will be explained in Section 3.4). We provide an example of the code prompt in Appendix (Table 9).

3.3 Iterative Refinement

Hallucination Grader The code prompt efficiently guides LLMs to generate graphs, but it often generates hallucinated relations. Based on our preliminary experiments, these hallucinations stem from the models’ overconfidence in their relation predictions. Specifically, LLMs tend to infer event relations without explicit linguistic cues or strong evidence for logical inference. Consequently, LLMs predict far more relations than the gold standards, leading to low precision.

²<https://networkx.org/documentation/stable/>

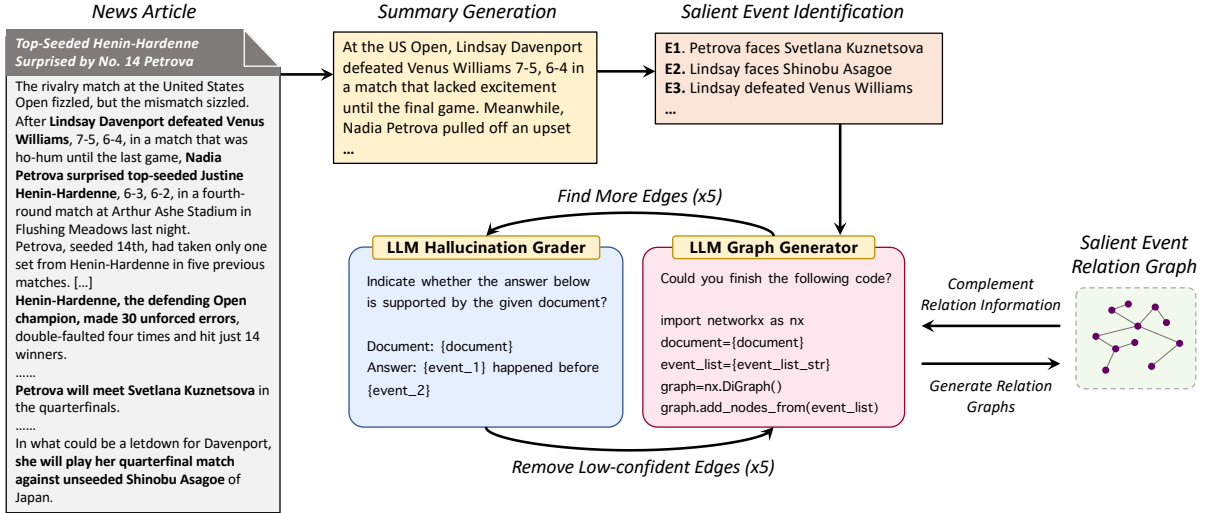


Figure 2: The proposed CALLMSAE framework.

Recent studies show that LLMs can evaluate and correct their own outputs (Madaan et al., 2023; Asai et al., 2024). Thus, we introduce a hallucination grader to address hallucination. For each relation edge generated, we pose a question to the LLMs to determine whether the relation is grounded in the given document. If the LLMs respond with a “yes”, the edge is retained; otherwise, it is discarded. An example of the hallucination grader prompt is shown in Appendix (Table 10).

Recover Missing Edges The main benefit of the hallucination grader is that it increases precision by removing low-confident edges. However, this process inevitably reduces recall. To mitigate this side effect, we introduce an iterative refinement process. After discarding hallucinated edges, we reinsert the code block containing the relation edges into the graph generation prompt and ask the LLMs to complete the code again. In this way, the LLMs can reconsider whether there are any missing relations in the document, thereby improving recall.

Once the LLMs generate a new graph, the hallucination grader checks the relation edges again. This process is repeated for a fixed number of times. We set the maximum number of iterations to 5 in our experiments, as the LLMs stop discovering new edges after 2 or 3 iterations in most documents.

3.4 Complement Relation Types

Hierarchical, temporal, and causal relations are not independent of each other. We found that if one type of relation depends on another, providing the graph for the first relation can benefit the generation of the dependent relation’s graph. Specifically, we

predict the hierarchical relation graph first. Then, we provide this graph to the LLMs and ask them to generate the temporal relation graph. Lastly, with both the hierarchical and temporal relation graphs available, the LLMs predict the causal relation graph.

The hierarchical relation describes two closely related events at different granularity levels. It focuses on the inherent semantics of the events and does not depend on other relation types. For example, “writing a dissertation” is a subevent of “doing a PhD”. Therefore, we choose to predict the hierarchical relations first.

Temporal relations can depend on hierarchical relations. For example, knowing “doing a PhD” happened before “being prompted to Professor” allows us to deduce that “writing a thesis” also happened before “being prompted to Professor”. Thus, we predict temporal relations after hierarchical relations.

Lastly, causal relations depend on both hierarchical and temporal relation, as the antecedent event in a causal relation must occur before the consequence. Therefore, the causal relation is predicted in the last step. For more details about the entire prompting process, please refer to the descriptions and pseudocode in Appendix C.

3.5 Hungarian Graph Similarity

It is challenging to compare event relation graphs generated by LLMs due to the abstractive nature of generation, making it difficult to align the generated events with the gold standard events (Li et al., 2023). Moreover, salient events are often high-level and abstract rather than fine-grained and con-

crete, which means some variations in wording is not only acceptable but also expected. Instead of using exact matching (Zhao et al., 2024) or rule-based token matching (Tan et al., 2024b) on events and relations to calculate F_1 , adopting semantic-based evaluation metrics is more reasonable and fair. As more tasks adopt text generation frameworks, many researchers are also turning to metrics based on language models rather than traditional token matching metrics like ROUGE and BLUE (Goyal et al., 2022; Pratapa et al., 2023).

In this study, we propose a novel metric for evaluating LLM-generated event graphs, called **Hungarian Graph Similarity (HGS)**. The metric is based on the Hungarian assignment algorithm (Kuhn, 1955), which is widely used in the object detection to match generated objects and target objects (Carion et al., 2020). It can find the optimal assignment given a cost matrix containing the distance between elements in two lists of objects. We adapt this algorithm to match predicted edges with edges in the gold standard graphs as follows:

1. Encode the events using SFR-Embedding-Mistral (Meng et al., 2024), which was ranked 1st on the Massive Text Embedding Benchmark leaderboard (Muennighoff et al., 2022) at the time of our experiments.
2. Given two edges of the same relation type, let \bar{e}_1^h, \bar{e}_1^t be the embeddings of the head event and the tail event in the first edge. Let \bar{e}_2^h, \bar{e}_2^t be the embeddings of the head and tail events in the second edges. We define the distance between the edges as $\max(D_{\cos}(\bar{e}_1^h, \bar{e}_2^h), D_{\cos}(\bar{e}_1^t, \bar{e}_2^t))$, where $D_{\cos}(\cdot, \cdot)$ is the cosine distance.
3. Build a cost matrix by computing the distance between every edge pair in the gold and predicted edge sets. Pad the matrix to a square matrix with the maximum cost value of 1.
4. Apply the Hungarian algorithm to the cost matrix to get the minimal cost value. The final score is $1 - \text{cost value}$, making the value more intuitive (higher is better). To compute the HGS over all the documents, we weight the scores by the number of gold edges to obtain an average value.

In step 2, we take the maximum value of the distances between head and tail events because relation edges are considered matched only if both the head and tail events match.

For more detailed analysis, we define precision-oriented HGS and recall-oriented HGS. We match edges without padding the cost matrix in step 3 to obtain the total cost values of all matched edge pairs. Then, the total matched similarity is the number of matched edges minus the total cost. **Precision-oriented HGS** is computed by dividing the total matched similarity by the total number of predicted edges. **Recall-oriented HGS** is computed by dividing the total matched similarity by the total number of edges in the target graph.

4 Dataset

In this section, we describe how we construct the distant supervision dataset and a human-annotated dataset from the New York Times (NYT) corpus.

4.1 Document Selection

We follow the same procedures as in (Tan et al., 2024a) to select documents from the NYT corpus, one of the largest news datasets, with additional filtering based on document length. We select 10,347 documents based on their descriptors indicating they are related to event narratives instead of opinions and discussions, such as sports and international politics. Among them, 100 documents are randomly sampled as the test set to be annotated by humans. More details about data selection are shown in Appendix A.1.

4.2 Annotation Settings

We recruited annotators from Prolific³. There are two subtasks: *salient event identification* and *event relation identification*. In the first subtask, the participants are asked to identify the salient event triplets: *actor*, *trigger*, and *object* (optional). We provide the definition of an event and several examples in the guidelines. They are instructed to do the summarisation test: the salient events should be the events they would include in the summary of the given article. Moreover, we provide some prominent features for helping annotators to identify salient events (Choubey et al., 2018; Jindal et al., 2020):

- Frequency: salient events are frequently mentioned in the articles.
- First appearance: salient events are often mentioned at the beginning of the article.

³<http://www.prolific.com>

- **Stretch size:** salient events are often mentioned throughout the articles. Stretch size is the distance between the location where the event is first mentioned and last mentioned. A salient event usually has a large stretch size.

In the second stage, we ask participants to identify relation triplets: *a source event, a relation type, and a target event*. Both the source and target events should be among the salient events identified in the first stage. In the guideline, we define three relation types: *happened_before*, *caused_by*, and *is_subevent_of*. *happened_before* indicates that the source event happened earlier than the target event. *caused_by* means the source event would not have happened if the target event did not happen. *is_subevent_of* signifies that the source event is a subevent of the target event. Annotators were informed that relations would be either explicitly mentioned in the article or inferred based on evidence within the article. Further details about the guidelines and user interface can be found in Appendix A.4.

4.3 Inter Annotator Agreement

Identifying salient events and event relations is complicated and time-consuming. We found it challenging to educate participants about these concepts because, in daily life, the meanings of events and relations differ from their definitions in the field of information extraction. Moreover, the technical definitions are much less intuitive to those outside the academic field. As a result, thorough training of participants is important to obtain high-quality annotations.

In total, we recruited 3 annotators to annotate 100 documents. Due to their varying availability, annotator 1 and 2 each annotated 45 documents, while annotator 3 annotated 20 documents. Among these, 5 documents were annotated by all three annotators. Following prior research in information extraction (Gurulingappa et al., 2012; Zhao et al., 2024), we used F_1 to measure the inter-annotator agreement on these 5 documents. To compute inter-annotator agreement, events or relations identified by one annotator are represented as set S_1 . Another annotator’s annotation S_2 serves as a pseudo-reference to compute precision $= \frac{|S_1 \cap S_2|}{|S_1|}$, recall $= \frac{|S_1 \cap S_2|}{|S_2|}$, and the F_1 score $= \frac{2|S_1 \cap S_2|}{|S_1| + |S_2|}$.

Table 1 shows the agreement scores for stages 1 and 2. Identifying salient events is subjective, which makes it difficult to reach a complete agreement. Moreover, event relation identification is

even more subjective and dependent on the previous stage, leading to less unanimous agreement.

Annotator	Stage 1	Stage 2
1 & 2	0.838	0.676
1 & 3	0.771	0.645
2 & 3	0.847	0.710
Average	0.819	0.677

Table 1: Inter-annotator agreement measured by F_1 .

4.4 Dataset Statistics

Table 2 shows the distributions of the relation types after applying the transitive closure to the annotated graphs. *happened_before* emerges as the most frequent relation type, reflecting the predominant focus on temporal sequences in news articles, and they are relatively straightforward to identify. Conversely, *caused_by* is the least frequent as it is the most challenging to identify.

Relation Type	Number
<i>happened_before</i>	310
<i>caused_by</i>	202
<i>is_subevent_of</i>	245
Total	757

Table 2: The distributions of the relation types.

5 Experiments

5.1 Model Settings

We compare our proposed approach with the following baselines:

- CAEVO (McDowell et al., 2017) is a pipeline system based on a Maximum Entropy (Max-Ent) classifier and manually designed features for extracting events and temporal relations.
- Madaan and Yang (2021) trained language models on CAEVO-generated linearised graphs with the language modelling objective. We implemented their method to train a Flan-T5 model.
- Tan et al. (2024a) also trained language models on CAEVO-generated graphs, but applied data augmentations and regularisations to mitigate the set element misalignment issue. We applied their method to train a Flan-T5.
- Han et al. (2019b) proposed a joint event and temporal relation extraction model. We

adapted the model to predict hierarchical and causal relations by training it on the MAVEN-ERE dataset (Wang et al., 2022). We also replaced BERT with Longformer (Beltagy et al., 2020) to enable it to process long documents.

- GPT-3.5 is an LLM based on the generative pre-train framework⁴. We used “gpt-3.5-turbo”.
- GPT-4 is also an LLM based on the generative pre-train framework (OpenAI et al., 2024). We used “gpt-4-1106-preview”.
- MIXTRAL is an LLM based on the Mistral model and the mixture of expert framework. We used the Mixtral 8x7B instruct version (Jiang et al., 2024).
- LLAMA3 is an LLM based on the Llama framework⁵. We used the Llama3-70B-instruct 8-bit version provided by ollama⁶. The 8-bit quantization is shown to be degradation-free (Detrmers et al., 2022).

We fine-tuned a Flan-T5-base (250M) with the relation graphs generated by CALLMSAE, following the same method as in Tan et al. (2024a). The baseline prompt evaluates whether each event pair is supported by the document, akin to the hallucination grader described in Section 3.3. Thus, it serves as an ablation of our method without incorporating the code prompt.

CALLMSAE is designed to be model-agnostic. Due to budget constraints and the preliminary test results, we chose Llama3 as the backbone of all the prompt-based methods detailed in Table 5.

5.2 Event Saliency Evaluation

Table 3 shows the salient features (defined in Section 4.2, computation formulas in Appendix B) extracted from various backbone LLMs using summarisation prompts, alongside comparison with CAEVO and human annotations. The LLM-generated events are much more salient than CAEVO-generated events and exhibit similarity to human annotations.

We also use human annotations to evaluate the saliency. In the salient event identification annotation, we provide the events generated by CAEVO

⁴<https://platform.openai.com/docs/models/gpt-3-5-turbo>

⁵<https://ai.meta.com/blog/meta-llama-3/>

⁶https://ollama.com/library/llama3:70b-instruct-q8_0

	Mean event number	Event frequency \uparrow	First appearance \downarrow	Stretch Size \uparrow
CAEVO	34.71	0.05	0.46	0.07
Human	8.26	0.11	0.31	0.20
GPT-4	6.49	0.09	0.37	0.18
Llama3	5.17	0.09	0.30	0.19
Mixtral	10.60	0.10	0.33	0.20

Table 3: The average number of extracted events and the saliency features (in percentage values).

	P	R	F_1	HGS
CAEVO	3.29	3.72	3.49	18.18
Mixtral	48.97	56.77	52.59	67.15

Table 4: Precision, recall, and F_1 based on the choices of the annotators. Hungarian graph similarity (HGS) is defined in Section 3.5. The values are in percentage.

and Mixtral as candidate salient events. Note that only the top CAEVO events ranked in saliency features are shown. Half of the candidates are from CAEVO and the other half are from Mixtral. They are randomly shuffled and then shown to the annotators. We compute the precision, recall, and F_1 based on how the annotators select them. We also compute HGS using human-annotated salient events as references (Table 4). It is clear that although CAEVO extracted more events than Mixtral, many of them are not salient. Mixtral outperforms CAEVO significantly across all evaluation metrics.

5.3 Salient Event Relation Graph Evaluation

The salient event relation graph evaluation results are shown in Table 5. Even with the most basic prompting (*Baseline Prompt*), which queries the relation of each event pair, Llama3 outperforms all the baseline methods on all relation types. However, *Baseline Prompt* is slow and costly because the number of prompts it needs for building one graph is $O(n^2)$, where n is the number of events in the document. On the other hand, *Code Prompt* only needs $O(1)$. Moreover, *Code Prompt*’s overall HGS is significantly higher than *Baseline Prompt* on all relation types. *Baseline Prompt* check the event pairs more thoroughly and thus have higher recall but its precision is much lower. The complete CALLMSAE combines the code prompt and hallucination grader for iterative refinement, checking missing relations and verifying them to prevent hallucination. It significantly increases the precision and strikes a balance with recall.

	Hierarchical			Temporal			Causal		
	PHGS	RHGS	HGS	PHGS	RHGS	HGS	PHGS	RHGS	HGS
Han et al. (2019b)	0.158	0.247	0.098	0.092	0.352	0.148	0.084	0.316	0.116
CAEVO	-	-	-	0.030	0.558	0.092	-	-	-
Madaan and Yang (2021)	-	-	-	0.061	0.439	0.116	-	-	-
Tan et al. (2024a)	-	-	-	0.126	0.335	0.187	-	-	-
Baseline Prompt	0.076	0.651	0.248	0.085	0.627	0.195	0.062	0.657	0.207
Code Prompt	0.174	0.559	0.315	0.153	0.678	0.283	0.121	0.632	0.272
Code Prompt (dependent rels)	0.196	0.544	0.334	0.211	0.601	0.341	0.135	0.599	0.272
CALLMSAE (ours)	0.196	0.544	0.334	0.294	0.509	0.327	0.198	0.529	0.295
Fine-tuned T5 (CALLMSAE)	0.314	0.434	0.339	0.244	0.544	0.362	0.366	0.397	0.343

Table 5: The Hungarian graph similarity (HGS) of the LLM-generated graphs on the human-annotated NYT dataset. PHGS is precision-oriented HGS. RHGS is recall-oriented HGS. *Code Prompt (dependent rels)* means adding hierarchical graphs in the prompts for temporal graphs; and adding hierarchical and temporal for causal graphs. *Fine-tuned T5 (CALLMSAE)* means fine-tuning a flan-T5 using the graphs generated by CALLMSAE.

In the temporal category, the results of *Code Prompt (dependent rels)* are obtained when provided with hierarchical graphs generated by CALLMSAE to LLMs. It has much higher overall HGS and precision than *Code Prompt* without hierarchical information, showing that hierarchical information can mitigate hallucinations during the temporal graph generation. In the casual category, the results of *Code Prompt (dependent rels)* are obtained when given both hierarchical and temporal graphs generated by CALLMSAE. The additional information also increases precision.

Fine-tuned T5 outperform all the methods based on CAEVO (McDowell et al., 2017; Madaan and Yang, 2021; Tan et al., 2024a), showing that the high-quality graphs generated by CALLMSAE can boost the contextualised graph generation. Interestingly, the performance of the *Fine-tuned T5*, fine-tuned on CALLMSAE-generated data, exceeds that of CALLMSAE itself, implying that the fine-tuned model can effectively adapt the reasoning patterns provided by Llama3 and generalise them.

5.4 Format Error and Cycles in the Graphs

A format error occurs when the generated code blocks fail to pass the Python interpreter. We detected these errors by executing the generated code. If the Python interpreter returns an error, it is classified as a format error. We specified the relation graphs as directed acyclic graphs in the prompt. If there is a cycle in the generated graph, it means that the LLM failed to follow the instructions. A cycle also indicates a violation of logic constraints because all the relations in the event relation graphs are asymmetric. We detected the cycles using the `find_cycle()` from the NetworkX after obtaining the transitive closure of the graphs.

	Format Error	Cycle
GPT-3.5	0	10.67
GPT-4	3.67	1.67
Mixtral	3.33	2.33
Llama3	0	0

Table 6: The average number of CALLMSAE-generated graphs out of 100 with format errors or cycles.

We prompt each LLM three times on the annotated test set. Table 6 shows the average number of documents encountering format errors or cycles. All LLMs have low rates of format errors which shows that state-of-the-art LLMs can understand the instruction well and generate executable Python code. Among them, GPT-3.5 and Llama3 achieve zero errors. The occurrence of cycles can serve as an indicator of the reasoning ability of the LLMs. About 10% of graphs generated by GPT-3.5 have cycles, suggesting that GPT-3.5 may have limited reasoning ability compared to other LLMs. GPT-4 and Mixtral both have low rates of cycle occurrence, but they are beaten by Llama3 which has no cycle in all generations, showing its remarkable understanding of the transitive and asymmetric constraints in the complex event relation graphs.

6 Conclusion

This study explored utilising LLMs to generate salient event relation graphs from news documents without relying on human annotations. We studied how the events generated by LLMs are compared to the traditional methods in terms of event saliency. We further demonstrated that CALLMSAE-generated graphs can serve as distant signals to fine-tune smaller models and outperform those based on CAEVO.

646 Limitations

647 Although we have tested many prompting methods
648 and included several of the most effective ones in
649 this paper, we have not explored all possible com-
650 binations due to the extensive volume of recent
651 literature on prompt engineering. There might still
652 exist combinations of prompts that could further
653 improve performance. However, we are almost cer-
654 tain that any potential combinations, if they exist,
655 are likely to be more complex and thus less effi-
656 cient for building large-scale datasets. For example,
657 we did not add demonstrations in graph generation
658 because the code template is already quite lengthy.
659 Adding more documents could potentially exceed
660 the context windows of some LLMs, making it
661 challenging for them to interpret the instructions
662 effectively.

663 Ethics Statement

664 Event relation graph generation is a powerful tool
665 for understanding text. A potential misuse of the
666 proposed method is mining user behaviours on their
667 private data. For example, salient event relation
668 graphs can be extracted from users’ tweets to anal-
669 yse their potential reactions to advertisements and
670 scams. That could be a huge risk to social media
671 users.

672 Another potential risk is that the saliency may
673 introduce bias. LLMs may have their preferences
674 in selecting a specific group of events as important
675 events due to the data they were trained on. This is
676 a question which requires further large-scale inves-
677 tigation. However, we think this risk is negligible
678 in this study because we work on document-level
679 information. There is little room for selection given
680 that the news articles are already the products of
681 choice and distillation. If the system is used to ex-
682 tract information from a border information source,
683 such as social media, the risk must be carefully
684 assessed.

685 References

686 Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and
687 Hannaneh Hajishirzi. 2024. [Self-RAG: Learning to
688 retrieve, generate, and critique through self-reflection.](#)
689 In *The Twelfth International Conference on Learning
690 Representations*.

691 Iz Beltagy, Matthew E. Peters, and Arman Cohan. 2020.
692 [Longformer: The long-document transformer.](#) *CoRR*,
693 abs/2004.05150.

Nicolas Carion, Francisco Massa, Gabriel Synnaeve, 694
Nicolas Usunier, Alexander Kirillov, and Sergey 695
Zagoruyko. 2020. End-to-end object detection with 696
transformers. In *European conference on computer 697
vision*, pages 213–229. Springer. 698

Chunkit Chan, Jiayang Cheng, Weiqi Wang, Yuxin 699
Jiang, Tianqing Fang, Xin Liu, and Yangqiu Song. 700
2023. [Chatgpt evaluation on sentence level relations:
701 A focus on temporal, causal, and discourse relations.](#) 702

Prafulla Kumar Choubey, Kaushik Raju, and Ruihong 703
Huang. 2018. [Identifying the most dominant event
704 in a news article by mining event coreference rela-
705 tions.](#) In *Proceedings of the 2018 Conference of the
706 North American Chapter of the Association for Com-
707 putational Linguistics: Human Language Technolo-
708 gies, Volume 2 (Short Papers)*, pages 340–345, New
709 Orleans, Louisiana. Association for Computational
710 Linguistics. 711

Tim Dettmers, Mike Lewis, Younes Belkada, and Luke 712
Zettlemoyer. 2022. [GPT3.int8\(\): 8-bit matrix mul-
713 tiplication for transformers at scale.](#) In *Advances in
714 Neural Information Processing Systems*. 715

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and 716
Kristina Toutanova. 2019. [BERT: Pre-training of
717 deep bidirectional transformers for language under-
718 standing.](#) In *Proceedings of the 2019 Conference of
719 the North American Chapter of the Association for
720 Computational Linguistics: Human Language Tech-
721 nologies, Volume 1 (Long and Short Papers)*, pages
722 4171–4186, Minneapolis, Minnesota. Association for
723 Computational Linguistics. 724

Jesse Dunietz and Daniel Gillick. 2014. [A new entity
725 salience task with millions of training examples.](#) In
726 *Proceedings of the 14th Conference of the European
727 Chapter of the Association for Computational Lin-
728 guistics, volume 2: Short Papers*, pages 205–209,
729 Gothenburg, Sweden. Association for Computational
730 Linguistics. 731

Tanya Goyal, Junyi Jessy Li, and Greg Durrett. 2022.
732 [News summarization and evaluation in the era of
733 gpt-3.](#) *arXiv preprint arXiv:2209.12356*. 734

Harsha Gurulingappa, Abdul Mateen Rajput, Angus 735
Roberts, Juliane Fluck, Martin Hofmann-Apitius, and
736 Luca Toldo. 2012. Development of a benchmark
737 corpus to support the automatic extraction of drug-
738 related adverse effects from medical case reports.
739 *Journal of biomedical informatics*, 45(5):885–892. 740

Rujun Han, I-Hung Hsu, Mu Yang, Aram Galstyan, 741
Ralph Weischedel, and Nanyun Peng. 2019a. [Deep
742 structured neural network for event temporal relation
743 extraction.](#) In *Proceedings of the 23rd Conference on
744 Computational Natural Language Learning (CoNLL)*,
745 pages 666–106, Hong Kong, China. Association for
746 Computational Linguistics. 747

Rujun Han, Qiang Ning, and Nanyun Peng. 2019b. [Joint
748 event and temporal relation extraction with shared](#) 749

750	representations and structured prediction. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 434–444, Hong Kong, China. Association for Computational Linguistics.	
751		
752		
753		
754		
755		
756		
757	Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, Gianna Lengyel, Guillaume Bour, Guillaume Lample, L�lio Renard Lavaud, Lucile Saulnier, Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, Teven Le Scao, Th�ophile Gervet, Thibaut Lavril, Thomas Wang, Timoth�e Lacroix, and William El Sayed. 2024. Mistral of experts .	
758		
759		
760		
761		
762		
763		
764		
765		
766		
767		
768	Disha Jindal, Daniel Deutsch, and Dan Roth. 2020. Is killed more significant than fled? a contextual model for salient event detection . In <i>Proceedings of the 28th International Conference on Computational Linguistics</i> , pages 114–124, Barcelona, Spain (Online). International Committee on Computational Linguistics.	
769		
770		
771		
772		
773		
774		
775	Harold W Kuhn. 1955. The hungarian method for the assignment problem. <i>Naval research logistics quarterly</i> , 2(1-2):83–97.	
776		
777		
778	Bo Li, Gexiang Fang, Yang Yang, Quansen Wang, Wei Ye, Wen Zhao, and Shikun Zhang. 2023. Evaluating chatgpt’s information extraction capabilities: An assessment of performance, explainability, calibration, and faithfulness .	
779		
780		
781		
782		
783	Tie-Yan Liu et al. 2009. Learning to rank for information retrieval. <i>Foundations and Trends® in Information Retrieval</i> , 3(3):225–331.	
784		
785		
786	Zhengzhong Liu, Chenyan Xiong, Teruko Mitamura, and Eduard Hovy. 2018. Automatic event salience identification . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 1226–1236, Brussels, Belgium. Association for Computational Linguistics.	
787		
788		
789		
790		
791		
792	Junru Lu, Xingwei Tan, Gabriele Pergola, Lin Gui, and Yulan He. 2022. Event-centric question answering via contrastive learning and invertible event transformation . In <i>Findings of the Association for Computational Linguistics: EMNLP 2022</i> , pages 2377–2389, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	
793		
794		
795		
796		
797		
798		
799	Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, Shashank Gupta, Bodhisattwa Prasad Majumder, Katherine Hermann, Sean Welleck, Amir Yazdanbakhsh, and Peter Clark. 2023. Self-refine: Iterative refinement with self-feedback . In <i>Thirty-seventh Conference on Neural Information Processing Systems</i> .	
800		
801		
802		
803		
804		
805		
806		
807		
	Aman Madaan and Yiming Yang. 2021. Neural language modeling for contextualized temporal graph generation . In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 864–881, Online. Association for Computational Linguistics.	808
		809
		810
		811
		812
		813
		814
	Bill McDowell, Nathanael Chambers, Alexander Ororbia II, and David Reitter. 2017. Event ordering with a generalized model for sieve prediction ranking . In <i>Proceedings of the Eighth International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 843–853, Taipei, Taiwan. Asian Federation of Natural Language Processing.	815
		816
		817
		818
		819
		820
		821
	Igor Melnyk, Pierre Dognin, and Payel Das. 2022. Knowledge graph generation from text . In <i>Findings of the Association for Computational Linguistics: EMNLP 2022</i> , pages 1610–1622, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	822
		823
		824
		825
		826
		827
	Rui Meng, Ye Liu, Shafiq Rayhan Joty, Caiming Xiong, Yingbo Zhou, and Semih Yavuz. 2024. Sfr-embedding-mistral:enhance text retrieval with transfer learning . Salesforce AI Research Blog.	828
		829
		830
		831
	Niklas Muennighoff, Nouamane Tazi, Loic Magne, and Nils Reimers. 2022. Mteb: Massive text embedding benchmark . <i>arXiv preprint arXiv:2210.07316</i> .	832
		833
		834
	Qiang Ning, Sanjay Subramanian, and Dan Roth. 2019. An improved neural baseline for temporal relation extraction . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 6203–6209, Hong Kong, China. Association for Computational Linguistics.	835
		836
		837
		838
		839
		840
		841
		842
	Qiang Ning, Ben Zhou, Zhili Feng, Haoruo Peng, and Dan Roth. 2018. CogCompTime: A tool for understanding time in natural language . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations</i> , pages 72–77, Brussels, Belgium. Association for Computational Linguistics.	843
		844
		845
		846
		847
		848
		849
	OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai,	850
		851
		852
		853
		854
		855
		856
		857
		858
		859
		860
		861
		862
		863
		864
		865

866	Cory Decareaux, Thomas Degry, Noah Deutsch,	Clemens Winter, Samuel Wolrich, Hannah Wong,	930
867	Damien Deville, Arka Dhar, David Dohan, Steve	Lauren Workman, Sherwin Wu, Jeff Wu, Michael	931
868	Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti,	Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qim-	932
869	Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix,	ing Yuan, Wojciech Zaremba, Rowan Zellers, Chong	933
870	Simón Posada Fishman, Juston Forte, Isabella Ful-	Zhang, Marvin Zhang, Shengjia Zhao, Tianhao	934
871	ford, Leo Gao, Elie Georges, Christian Gibson, Vik	Zheng, Juntang Zhuang, William Zhuk, and Barret	935
872	Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-	Zoph. 2024. Gpt-4 technical report .	936
873	Lopes, Jonathan Gordon, Morgan Grafstein, Scott		
874	Gray, Ryan Greene, Joshua Gross, Shixiang Shane	Rruba Panchendrarajan and Aravindh Amaresan. 2018.	937
875	Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris,	Bidirectional LSTM-CRF for named entity recogni-	938
876	Yuchen He, Mike Heaton, Johannes Heidecke, Chris	tion . In <i>Proceedings of the 32nd Pacific Asia Con-</i>	939
877	Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele,	<i>ference on Language, Information and Computation</i> ,	940
878	Brandon Houghton, Kenny Hsu, Shengli Hu, Xin	Hong Kong. Association for Computational Linguis-	941
879	Hu, Joost Huizinga, Shantanu Jain, Shawn Jain,	tics.	942
880	Joanne Jang, Angela Jiang, Roger Jiang, Haozhun		
881	Jin, Denny Jin, Shino Jomoto, Billie Jonn, Hee-	Adithya Pratapa, Kevin Small, and Markus Dreyer.	943
882	woo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kama-	2023. Background summarization of event time-	944
883	mali, Ingmar Kanitscheider, Nitish Shirish Keskar,	lines . In <i>Proceedings of the 2023 Conference on</i>	945
884	Tabarak Khan, Logan Kilpatrick, Jong Wook Kim,	<i>Empirical Methods in Natural Language Processing</i> ,	946
885	Christina Kim, Yongjik Kim, Jan Hendrik Kirchner,	pages 8111–8136, Singapore. Association for Com-	947
886	Jamie Kiros, Matt Knight, Daniel Kokotajlo,	putational Linguistics.	948
887	Łukasz Kondraciuk, Andrew Kondrich, Aris Kon-		
888	stantinidis, Kyle Kosic, Gretchen Krueger, Vishal	Evan Sandhaus. 2008. The New York Times Annotated	949
889	Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan	Corpus .	950
890	Leike, Jade Leung, Daniel Levy, Chak Ming Li,		
891	Rachel Lim, Molly Lin, Stephanie Lin, Mateusz	Maarten Sap, Ronan Le Bras, Emily Allaway, Chan-	951
892	Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue,	dra Bhagavatula, Nicholas Lourie, Hannah Rashkin,	952
893	Anna Makanju, Kim Malfacini, Sam Manning, Todor	Brendan Roof, Noah A. Smith, and Yejin Choi.	953
894	Markov, Yaniv Markovski, Bianca Martin, Katie	2019. Atomic: An atlas of machine commonsense	954
895	Mayer, Andrew Mayne, Bob McGrew, Scott Mayer	for if-then reasoning . In <i>Proceedings of the Thirty-</i>	955
896	McKinney, Christine McLeavey, Paul McMillan,	<i>Third AAAI Conference on Artificial Intelligence and</i>	956
897	Jake McNeil, David Medina, Aalok Mehta, Jacob	<i>Thirty-First Innovative Applications of Artificial In-</i>	957
898	Menick, Luke Metz, Andrey Mishchenko, Pamela	<i>telligence Conference and Ninth AAAI Symposium</i>	958
899	Mishkin, Vinnie Monaco, Evan Morikawa, Daniel	<i>on Educational Advances in Artificial Intelligence</i> ,	959
900	Mossing, Tong Mu, Mira Murati, Oleg Murk, David	AAAI'19/IAAI'19/EAAI'19. AAAI Press.	960
901	Mély, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak,		
902	Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh,	Xingwei Tan, Gabriele Pergola, and Yulan He. 2021.	961
903	Long Ouyang, Cullen O'Keefe, Jakub Pachocki, Alex	Extracting event temporal relations via hyperbolic	962
904	Paino, Joe Palermo, Ashley Pantuliano, Giambat-	geometry . In <i>Proceedings of the 2021 Conference</i>	963
905	tista Parascandolo, Joel Parish, Emy Parparita, Alex	<i>on Empirical Methods in Natural Language Process-</i>	964
906	Passos, Mikhail Pavlov, Andrew Peng, Adam Perel-	<i>ing</i> , pages 8065–8077, Online and Punta Cana, Do-	965
907	man, Filipe de Avila Belbute Peres, Michael Petrov,	minican Republic. Association for Computational	966
908	Henrique Ponde de Oliveira Pinto, Michael, Poko-	Linguistics.	967
909	rny, Michelle Pokrass, Vitchyr H. Pong, Tolly Pow-		
910	ell, Alethea Power, Boris Power, Elizabeth Proehl,	Xingwei Tan, Yuxiang Zhou, Gabriele Pergola, and	968
911	Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh,	Yulan He. 2024a. Set-aligning framework for auto-	969
912	Cameron Raymond, Francis Real, Kendra Rimbach,	regressive event temporal graph generation .	970
913	Carl Ross, Bob Rotsted, Henri Roussez, Nick Ry-		
914	der, Mario Saltarelli, Ted Sanders, Shibani Santurkar,	Xingwei Tan, Yuxiang Zhou, Gabriele Pergola, and	971
915	Girish Sastry, Heather Schmidt, David Schnurr, John	Yulan He. 2024b. Set-aligning framework for auto-	972
916	Schulman, Daniel Selsam, Kyla Sheppard, Toki	regressive event temporal graph generation . <i>arXiv</i>	973
917	Sherbakov, Jessica Shieh, Sarah Shoker, Pranav	preprint arXiv:2404.01532 .	974
918	Shyam, Szymon Sidor, Eric Sigler, Maddie Simens,		
919	Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin	Naushad UzZaman, Hector Llorens, Leon Derczynski,	975
920	Sokolowsky, Yang Song, Natalie Staudacher, Fe-	James Allen, Marc Verhagen, and James Pustejovsky.	976
921	lipe Petroski Such, Natalie Summers, Ilya Sutskever,	2013. SemEval-2013 task 1: TempEval-3: Evaluat-	977
922	Jie Tang, Nikolas Tezak, Madeleine B. Thompson,	ing time expressions, events, and temporal relations .	978
923	Phil Tillet, Amin Tootoonchian, Elizabeth Tseng,	In <i>Second Joint Conference on Lexical and Compu-</i>	979
924	Preston Tuggle, Nick Turley, Jerry Tworek, Juan Fe-	<i>tational Semantics (*SEM), Volume 2: Proceedings</i>	980
925	lipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya,	<i>of the Seventh International Workshop on Seman-</i>	981
926	Chelsea Voss, Carroll Wainwright, Justin Jay Wang,	<i>tic Evaluation (SemEval 2013)</i> , pages 1–9, Atlanta,	982
927	Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei,	Georgia, USA. Association for Computational Lin-	983
928	CJ Weinmann, Akila Welihinda, Peter Welinder, Ji-	guistics.	984
929	ayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner,		

985 Haoyu Wang, Muhao Chen, Hongming Zhang, and Dan
986 Roth. 2020. [Joint constrained learning for event-](#)
987 [event relation extraction](#). In *Proceedings of the 2020*
988 *Conference on Empirical Methods in Natural Lan-*
989 *guage Processing (EMNLP)*, pages 696–706, Online.
990 Association for Computational Linguistics.

991 Xiaozhi Wang, Yulin Chen, Ning Ding, Hao Peng, Zimu
992 Wang, Yankai Lin, Xu Han, Lei Hou, Juanzi Li,
993 Zhiyuan Liu, Peng Li, and Jie Zhou. 2022. [MAVEN-](#)
994 [ERE: A unified large-scale dataset for event core-](#)
995 [ference, temporal, causal, and subevent relation ex-](#)
996 [traction](#). In *Proceedings of the 2022 Conference on*
997 *Empirical Methods in Natural Language Processing*,
998 pages 926–941, Abu Dhabi, United Arab Emirates.
999 Association for Computational Linguistics.

1000 Xingyao Wang, Sha Li, and Heng Ji. 2023. [Code4Struct:](#)
1001 [Code generation for few-shot event structure predic-](#)
1002 [tion](#). In *Proceedings of the 61st Annual Meeting of*
1003 *the Association for Computational Linguistics (Vol-*
1004 *ume 1: Long Papers)*, pages 3640–3663, Toronto,
1005 Canada. Association for Computational Linguistics.

1006 Li Zhang, Hainiu Xu, Yue Yang, Shuyan Zhou, Weiqiu
1007 You, Manni Arora, and Chris Callison-Burch. 2023.
1008 [Causal reasoning of entities and events in procedural](#)
1009 [texts](#). In *Findings of the Association for Compu-*
1010 *tational Linguistics: EACL 2023*, pages 415–431,
1011 Dubrovnik, Croatia. Association for Computational
1012 Linguistics.

1013 Tianyi Zhang, Faisal Ladhak, Esin Durmus, Percy Liang,
1014 Kathleen McKeown, and Tatsunori B. Hashimoto.
1015 2024. [Benchmarking Large Language Models for](#)
1016 [News Summarization](#). *Transactions of the Associa-*
1017 *tion for Computational Linguistics*, 12:39–57.

1018 Xiyang Zhang, Muhao Chen, and Jonathan May. 2021.
1019 [Salience-aware event chain modeling for narrative](#)
1020 [understanding](#). In *Proceedings of the 2021 Confer-*
1021 *ence on Empirical Methods in Natural Language Pro-*
1022 *cessing*, pages 1418–1428, Online and Punta Cana,
1023 Dominican Republic. Association for Computational
1024 Linguistics.

1025 Runcong Zhao, Qinglin Zhu, Hainiu Xu, Jiazheng Li,
1026 Yuxiang Zhou, Yulan He, and Lin Gui. 2024. Large
1027 language models fall short: Understanding complex
1028 relationships in detective narratives. *arXiv preprint*
1029 *arXiv:2402.11051*.

A Additional Details of Dataset Construction 1030 1031

A.1 Document Selection 1032

1033 We select news documents from the NYT corpus
1034 based on the descriptors available. With regards to
1035 the generation of salient event graphs, the most rel-
1036 evant documents tend to be centered around event
1037 narratives, so that they could be rich in event re-
1038 lations. Tan et al. (2024a) investigated which de-
1039 scriptors are rich in event narrative using event fre-
1040 quency \times inverse-descriptor frequency. We chose
1041 the documents using the same descriptors as them
1042 (e.g., “airlines and airplanes”, “united states inter-
1043 national relations”, “civil war and guerrilla war-
1044 fare”, “track and field”, “soccer”, etc.).

1045 We applied additional filtering based on the num-
1046 ber of words in the documents. Documents with
1047 more than 8500 words or less than 100 words are
1048 excluded. Based on our preliminary observations,
1049 the extremely long documents are not typically
1050 news articles (only takes 0.02% in the entire NYT).
1051 They tend to be collections of articles over longer
1052 time spans, making them not suitable as focus of
1053 this study. Additionally, very long articles may af-
1054 fect the performance of open-source LLMs only
1055 due to limitations in the context length rather than
1056 their reasoning abilities. On the other hand, ar-
1057 ticles that are too short are less likely to contain
1058 complex event relation graphs, so we also exclude
1059 them. The final average word count of the selected
1060 10347 documents is 780.

A.2 Frequent words and descriptors in the annotated dataset 1061 1062

Rank	Test		Train	
	Word	Count	Word	Count
1	win	41	win	2,964
2	express	15	make	1,591
3	play	14	face	1,564
4	make	13	express	1,411
5	defeat	12	include	1,307

Table 7: The top 5 most frequent trigger words in the human-annotated test set and the distant train set.

1063 Table 7 reports the most frequent trigger words
1064 among the human-identified salient events and
1065 LLM-generated salient events after filtering out the
1066 light words (words that have no semantic meaning).
1067 We could see that “win”, “play”, and “defeat” are

Rank	Test		Train	
	Descriptor	Count	Descriptor	Count
1	U.S. International Relations	27	Terrorism	2,885
2	Terrorism	21	U.S. International Relations	2,574
3	Bombs and Explosives	17	Bombs and Explosives	1,727
4	U.S. Armament and Defense	15	U.S. Armament and Defense	1,717
5	Politics and Government	15	Politics and Government	1,649

Table 8: The top 5 most frequent descriptors in the human-annotated test set and the distant train set.

prominent triggers due to the sports topics within the dataset. These articles usually mention multiple events with these triggers. Triggers like “*express*”, “*include*”, and “*make*” are instead common across different scenarios.

Table 8 shows the most frequent descriptors in the human-annotated test set and the distant train set. These are the typical event-rich topics and are full of narratives.

A.3 Disclaimers of Risks

Consider that a large portion of the new articles in the New York Times corpus are about violent incidences, such as terrorist attacks and war. To prevent inflicting harm to traumatised victims, we show the information clearly in the recruitment description on the Prolific platform (Figure 3).

What will happen?

You will be asked to annotate a series of news articles. In each article, you will be asked to **identify the salient events**. Salient events are the important events in the article that you will include if you were to write a summary of the article. Salient events are the center of the news. They could be key milestones or events that relate to many other events. Salient events usually locate in the main clause of a sentence.

The topics of these articles include sports, politics, crimes, and business. These articles are published from 1996 to 2007 in the New York Times. There may be descriptions of violent events, such as terrorist attacks and war.

Figure 3: The recruitment descriptions.

A.4 Guidelines and User Interface

A well-designed user interface is essential for collecting high-quality data efficiently. We fully cooperate with participants to improve the user interface iteratively based on their feedback.

In the salient event identification stage (Figure 4), we show the title, abstract, and content of the article on the right side. We show candidate events, which are extracted through CAEVO and Mixtral, on the left sidebar. The shown CAEVO events are the top events ranked based on the saliency feature score. The participants can choose the candidate events which they think are accurate and salient. The guideline also informs them that if multiple options refer to the same event, they can only choose the most accurate and informative one. If a salient

event is not present among the candidates, they could write it in the text input box and add it.

In the event relation identification stage (Figure 5), they could choose a source event, a relation type, and a target event to add a relation triplet. The source event and the target event need to be chosen from the salient event list from the first stage. We automatically detect and prevent any new event that will lead to duplication and contradiction. The participants can also deselect the added event if they change their minds. The participants were asked to finish the first stage first, and then annotate the second stage based on their own annotations in the first stage.

In the following are reported the screenshots of the guideline pages (Figure 6).

A.5 More details about the annotation

We started the annotation process by releasing several trial rounds, during which we chose participants based on their dedication and understanding of the terminologies. It required considerable communication efforts to ensure they had an accurate understanding of the task definition.

During training, we found a common mistake among the annotators was that they tended to overestimate the *is_subevent_of* relation. They often confused it with the *caused_by* relation or temporal inclusion.

We advised them that *is_subevent_of* pertains to two events on different granularity levels but referring to the same subject. To distinguish *is_subevent_of* from temporal overlap, they could check whether the actor in the subevent is the same as or a part of the actor or object in the parent event. For example, if a parent event is “*a team did something*” the subevent can be “*a member of the team did something*”.

A.6 Information about the Annotators

The annotators were paid at the rate of 8£/h. We screened native English speakers from all over the world to ensure they could read English articles

Judge Brinkema; must ensure; that the Constitution is fully applied in Moussaoui's case
 Zacarias Moussaoui; was denied; the right to see evidence crucial to his defense
 Ramzi bin al-Shibh; was a member of; Al Qaeda
 brinkema; ordered; government
 Bush administration; attempted; to bypass the Constitution while conducting the war on terror
 Justice Department; refused; to allow Moussaoui to question Ramzi bin al-Shibh
 evidence; assist; moussaoui
 Judge Leonie Brinkema; ordered; the government to make bin al-Shibh available
 The government; claimed; it would pose a threat to national security
 judge; allow; government
 The government; refused; to make bin al-Shibh available

Enter a new event you want to add

subject; predicate; object 0/150

Add the event

Current page (select to jump to a new page):

1 Logout

Guideline

Please click the guideline link above to view the annotation guidelines and examples 📄

The Trial of Zacarias Moussaoui

Abstract: Editorial says Justice Dept is trying to trample Bill of Rights in trial of Zacarias Moussaoui, so-called 20th hijacker, by denying him right to see evidence critical to his defense and then suggesting it might transfer his case to military tribunal if it does not like judge's ruling on matter; says war on terrorism has not repealed Constitution, and Judge Leonie Brinkema must ensure that it applies fully in Moussaoui's case

Since the Sept. 11 attacks, the Bush administration has repeatedly tried to dodge the Constitution while prosecuting the war on terror. In the trial of Zacarias Moussaoui, the so-called 20th hijacker, the Justice Department is once again attempting to trample the Bill of Rights -- in this case, by denying Mr. Moussaoui the right to see evidence critical to his defense. The judge should not allow the government to have its way. The dispute now raging in the Moussaoui case is over whether the defendant will be permitted to question Ramzi bin al-Shibh, a captured member of Al Qaeda who played a key role in the Sept. 11 conspiracy. Mr. bin al-Shibh is mentioned prominently in Mr. Moussaoui's indictment, and it is possible he could provide evidence that could assist Mr. Moussaoui in his defense. The Sixth Amendment guarantees a criminal defendant the right "to have compulsory process for obtaining witnesses in his favor," and Judge Leonie Brinkema has properly ordered the government to make Mr. bin al-Shibh available. But prosecutors have refused, arguing that allowing Mr. Moussaoui to question Mr. bin al-Shibh would pose a threat to national security. Faced with the government's defiance, Judge Brinkema can strike counts from

Figure 4: The user interface of salient event identification.

is_subevent_of:
(Bush administration; attempted; to bypass the Constitution while conducting the war on terror)

<S> (Judge Leonie Brinkema; ordered; the government to make bin al-Shibh available);
happened_before;
 (Justice Department; refused; to allow Moussaoui to question Ramzi bin al-Shibh)

<6> (Ramzi bin al-Shibh; was a member of; Al Qaeda);
happened_before;
 (Judge Leonie Brinkema; ordered; the government to make bin al-Shibh available)

<7> (Justice Department; refused; to allow Moussaoui to question Ramzi bin al-Shibh);
casued_by;
 (Bush administration; attempted; to bypass the Constitution while conducting the war on terror)

<8> (The government; claimed; it would pose a threat to national security);
casued_by;
 (Bush administration; attempted; to bypass the Constitution while conducting the war on terror)

Source Event:

Justice Department; refused; to allow Moussaoui to question Ra... ▼

Relation:

casued_by ▼

Target Event:

... ▼

Current page (select to jump to a new page):

1 Logout

Guideline

Please click the guideline link above to view the annotation guidelines and examples 📄

The Trial of Zacarias Moussaoui

Abstract: Editorial says Justice Dept is trying to trample Bill of Rights in trial of Zacarias Moussaoui, so-called 20th hijacker, by denying him right to see evidence critical to his defense and then suggesting it might transfer his case to military tribunal if it does not like judge's ruling on matter; says war on terrorism has not repealed Constitution, and Judge Leonie Brinkema must ensure that it applies fully in Moussaoui's case

Since the Sept. 11 attacks, the Bush administration has repeatedly tried to dodge the Constitution while prosecuting the war on terror. In the trial of Zacarias Moussaoui, the so-called 20th hijacker, the Justice Department is once again attempting to trample the Bill of Rights -- in this case, by denying Mr. Moussaoui the right to see evidence critical to his defense. The judge should not allow the government to have its way. The dispute now raging in the Moussaoui case is over whether the defendant will be permitted to question Ramzi bin al-Shibh, a captured member of Al Qaeda who played a key role in the Sept. 11 conspiracy. Mr. bin al-Shibh is mentioned prominently in Mr. Moussaoui's indictment, and it is possible he could provide evidence that could assist Mr. Moussaoui in his defense. The Sixth Amendment guarantees a criminal defendant the right "to have compulsory process for obtaining witnesses in his favor," and Judge Leonie Brinkema has properly ordered the government to make Mr. bin al-Shibh available. But prosecutors have refused, arguing that allowing Mr. Moussaoui to question Mr. bin al-Shibh would pose a threat to national security. Faced with the government's defiance, Judge Brinkema can strike counts from the indictment that involve Mr. bin al-Shibh, or dismiss the entire case. Allowing the government to deny access to Mr. bin al-Shibh with impunity would set the dangerous precedent that important constitutional rights can be taken away in terrorism cases. It is not at all clear that allowing Mr. Moussaoui to question Mr. bin al-Shibh in carefully monitored circumstances would threaten national security. If the Justice Department is convinced it would, it can adjust the charges against Mr. Moussaoui so Mr. bin al-Shibh's

Figure 5: The user interface of event relation identification.

1141 fluently. We also selected participants based on
1142 their previous submissions and approval rates to
1143 ensure they were familiar with the platform and
1144 were high-quality annotators.

1145 Two of the final annotators are identified as male,
1146 and they both come from the UK. One of the final
1147 annotators is identified as female, and she comes
1148 from Canada. They all identified as white.

1149 A.7 Dataset Licensing

1150 The original NYT corpus is available for noncom-
1151 mercial research license. One of our authors has
1152 obtained the license. Based on the license, we
1153 could not include the original text in our dataset.
1154 Thus, we will only release the generated/annotated
1155 graphs. Our dataset will also be in noncommercial
1156 research license.

1157 B Saliency Features

1158 Inspired by (Choubey et al., 2018), we calculate
1159 the saliency features to show how our proposed
1160 method differs from previous methods in terms of
1161 event saliency. Unlike conventional computation
1162 methods, these saliency features are calculated on
1163 the sentence level to be comparable across docu-
1164 ments of various lengths. These saliency features
1165 are:

1166 **Event frequency:** A salient event tends to ap-
1167 pear frequently in the document. Let $D =$
1168 $\{s_0, s_1, \dots, s_{n-1}, s_n\}$ be the document and the list
1169 of sentences in the document. Let e be the event.
1170 Let $M(e) = \{s_i, s_j, \dots, s_k\}, 0 \leq i < j < k \leq n$
1171 be the list of sentences which mention the event e .
1172 The event frequency is calculated as:

$$1173 \text{frequency}(D, e) = \frac{|M(e)|}{n + 1}. \quad (1)$$

1174 **First appearance:** News writers usually mention
1175 the salient event as early as possible to attract read-
1176 ers’ attention. The first appearance of the event e
1177 is computed as:

$$1178 \text{first_appearance}(D, e) = \frac{i}{n}. \quad (2)$$

1179 **Stretch size:** Salient events tend to be mentioned
1180 all across the document. The stretch size of event e
1181 is calculated as:

$$1182 \text{stretch_size}(D, e) = \frac{k - i}{n}. \quad (3)$$

1183 To detect which sentences mention the event e ,
1184 we first lemmatise the words in the document and

1185 the given event. Then, detect whether there is a
1186 matched substring the same as the given event in
1187 each sentence. However, the abstractive nature of
1188 LLM-based salient event generation makes exact
1189 matching not viable. To detect the event mention
1190 of LLM-generated events, we formulate a series of
1191 prompts. We first ask: “Which sentence in the doc-
1192 ument below mentions the event “{event}”? Please
1193 enclose that sentence in () and show it. Docu-
1194 ment: “””{doc_content}””””. Then, we employ
1195 iterative refinement in case the LLM misses any
1196 other sentences: “Is there any other sentence in the
1197 document directly mentioning the event “{event}”?
1198 Please enclose that sentence in () and show it.”
1199 Lastly, we collect the sentences from the responses.

1200 We run the methods on the human-annotated
1201 dataset (100 documents). We compute the saliency
1202 features of the events in each document and take
1203 the average across the events. Lastly, all the values
1204 are averaged across all the documents.

1205 C Prompting Details

1206 The total time cost for prompting Llama3 to con-
1207 struct the training data (10247 documents) is about
1208 2200 hours (total wall-clock time of all the ma-
1209 chines). We run the models on multiple machines
1210 with different specifications, including one with
1211 $6 \times$ RTX 3090, one with an A100, and one with
1212 $2 \times$ A40.

1213 Table 9 shows an example of the code prompt
1214 for hierarchical graph generation and the response
1215 from Llama3. Table 10 shows an example of the
1216 hallucination prompt and the response.

1217 Algorithm 1 is the pseudo-code of the entire
1218 salient event graph generation process.

1219 In the summarization prompt, we use a temper-
1220 ature of 0.8 and a top_p of 0.9. For the salient
1221 event generation prompt, we use a temperature of
1222 0.5 and a top_p of 0.9. The relation graph gener-
1223 ation prompt also uses a temperature of 0.5 and a
1224 top_p of 0.9. The hallucination grader prompt uses
1225 a temperature of 0.

Input to the LLMs	Response from the LLMs
<p>Could you please finish the following code?</p> <pre>import networkx as nx document = {document} event_list = {event_list_str} # This is a graph representing the hierarchical relation between the events in the document # Each edge in the graph represents a subevent relation between the head and tail nodes which are events # An edge means the head event is a subevent of the tail event. They are closely related but on different granularity levels. hierarchical_graph = nx.DiGraph() # This is a directed acyclic graph. There should not be any cycle in the graph. # Add events as nodes for event in event_list: hierarchical_graph.add_node(event) # Add hierarchical relations as edges to the graph using .add_edge() function. The function takes two strings as inputs. Each string represents an event. # There should not be any cycle in the graph. Ex- plain the reason for each added edge as a comment after each function call</pre>	<p>Here is the completed code:</p> <pre>import networkx as nx hierarchical_graph = nx.DiGraph() # Add events as nodes for event in event_list: hierarchical_graph.add_node(event) # Add hierarchical relations as edges to the graph hierarchical_graph.add_edge("the government re- sponded by reducing the provincial civil service by 28%", "liberals cut spending") # The event "the government responded by re- ducing the provincial civil service by 28%" is a subevent of "liberals cut spending" because it is one way the Liberals cut spending.</pre>

Table 9: An example of the relation graph generation prompt input and the response from Llama3.

Input to the LLMs	Response from the LLMs
<p>You are a grader assessing whether an answer is grounded in / supported by a set of facts. Give a binary score 'yes' or 'no' score to indicate whether the answer is grounded in or supported by a set of facts. Then, provide a short explanation. Here are the facts: {document} Here is the answer: Event "the government responded by reducing the provincial civil service by 28%" is a subevent of event "liberals cut spending".</p>	<p>Score: Yes Explanation: The answer is grounded in the facts because it accurately identifies a specific action taken by the government (reducing the provincial civil service by 28%) as a subevent of the broader event of cutting spending, which is mentioned in the text.</p>

Table 10: An example of the hallucination prompt input and the response from Llama3.

Algorithm 1 CALLMSAE: CASCADING Large Language Models for SALIENT Event graph generation

Input: Document d , Max Refinement Round k

Output: An Event Relation Graph g

```
1:  $summary \leftarrow \text{Summary\_Generation}(d)$ 
2:  $salient\_events \leftarrow \text{Event\_Generation}(summary)$ 
3:  $hierarchical\_graph \leftarrow \text{null}$ 
4:  $current\_round \leftarrow 0$ 
5: while  $current\_round < n$  do
6:    $hierarchical\_graph \leftarrow \text{Hierarchical\_Graph\_Generation}(d, salient\_events,$ 
    $hierarchical\_graph)$ 
7:    $hierarchical\_edges \leftarrow \text{Get\_Edges}(hierarchical\_graph)$ 
8:   for  $edge_i$  in  $hierarchical\_edges$  do
9:      $remove\_edge \leftarrow \text{Hallucination\_Grader}(d, edge_i)$ 
10:    if  $remove\_edge$  then
11:       $hierarchical\_graph \leftarrow \text{Remove\_edge}(hierarchical\_graph, edge_i)$ 
12:    end if
13:  end for
14:   $current\_round \leftarrow current\_round + 1$ 
15: end while
16:  $temporal\_graph \leftarrow \text{null}$ 
17:  $current\_round \leftarrow 0$ 
18: while  $current\_round < n$  do
19:    $temporal\_graph \leftarrow \text{Temporal\_Graph\_Generation}(d, salient\_events, temporal\_graph,$ 
    $hierarchical\_graph)$ 
20:    $temporal\_edges \leftarrow \text{Get\_Edges}(temporal\_graph)$ 
21:   for  $edge_i$  in  $temporal\_edges$  do
22:      $remove\_edge \leftarrow \text{Hallucination\_Grader}(d, edge_i)$ 
23:     if  $remove\_edge$  then
24:        $temporal\_graph \leftarrow \text{Remove\_edge}(temporal\_graph, edge_i)$ 
25:     end if
26:   end for
27:    $current\_round \leftarrow current\_round + 1$ 
28: end while
29:  $causal\_graph \leftarrow \text{null}$ 
30:  $current\_round \leftarrow 0$ 
31: while  $current\_round < n$  do
32:    $causal\_graph \leftarrow \text{Causal\_Graph\_Generation}(d, salient\_events, causal\_graph,$ 
    $temporal\_graph, hierarchical\_graph)$ 
33:    $causal\_edges \leftarrow \text{Get\_Edges}(causal\_graph)$ 
34:   for  $edge_i$  in  $causal\_edges$  do
35:      $remove\_edge \leftarrow \text{Hallucination\_Grader}(d, edge_i)$ 
36:     if  $remove\_edge$  then
37:        $causal\_graph \leftarrow \text{Remove\_edge}(causal\_graph, edge_i)$ 
38:     end if
39:   end for
40:    $current\_round \leftarrow current\_round + 1$ 
41: end while
42:  $g \leftarrow \{hierarchical\_graph, temporal\_graph, causal\_graph\}$ 
```

Guideline

Back to annotation page

When you visit the annotation platform for the first time, there may be a ngrok confirmation page. Just click 'visit' to confirm. ngrok is a tool which we use for setting up the website.

Step one

Select a page from the dropdown menu to start annotating. You can also use the 'Previous' and 'Next' buttons at the bottom to navigate through the pages.

Each page shows a news article. The large font text is the title. The bold font text is the abstract. The rest is the main content.

You should primarily use the words and phrases from the main content to construct events. Please don't repeat events from the title or the abstract.

Step two

The sidebar on the left show a list of candidate salient event suggested by an algorithm. We ask you to do the followings:

1. If you think an event is salient, tick the checkbox next to it. Otherwise, untick the checkbox.
2. If you think there is an event that isn't listed, you can add it by entering the event in the text box. **The event should at least contain a subject and a trigger.**
3. If you think a ticked event makes no sense, untick it. When two options are referring to the same event, untick the one you think is less informative.

Every modification will be saved automatically.

Definition of Event

An event is anything that happens as described in the article. We represent the events in a structured format: **actor; trigger; target**. The actor of the event is usually the subject of a sentence. The trigger can be seen as the predicate of a sentence. The target is usually the object in the sentence which is optional.

Example

New York is one of the four candidate cities competing to be presented to the IOC

The task force; will tour; some athletic sites and the proposed Olympic Village in Long Island City

The task force; will have; a critical five-hour session at City Hall on Monday

NYC2012; has added; more detail to its athletic, housing, transportation, finance, and security plans

NYC2012; has not changed; the proposed locations of any sports

NYC2012; has hired; consultants to study urban design and zoning of the far West Side of Manhattan

Officials Will Tour A Changed New York

Abstract: United States Olympic Committee task force will visit New York City to examine potential as host city for 2012 Summer Olympics; New York City competes against San Francisco, Washington-Baltimore and Houston amongst American candidates (M)

For the first time since the terrorist attacks of Sept. 11, a United States Olympic Committee task force will visit New York on Sunday and Monday to evaluate its worthiness to be America's designated host city competing for the 2012 Summer Games.

The Olympic plan laid out by NYC2012 for the summer visit last summer, which advanced New York into a pool of four candidate cities, has not changed much.

But the city's image has.

Step one

Up to 30,000 Troops From a Dozen Nations to Replace Some G.I.'s in Iraq

Abstract: United States Olympic Committee task force will visit New York City to examine potential as host city for 2012 Summer Olympics; New York City competes against San Francisco, Washington-Baltimore and Houston amongst American candidates (M)

For the first time since the terrorist attacks of Sept. 11, a United States Olympic Committee task force will visit New York on Sunday and Monday to evaluate its worthiness to be America's designated host city competing for the 2012 Summer Games.

The Olympic plan laid out by NYC2012 for the summer visit last summer, which advanced New York into a pool of four candidate cities, has not changed much.

But the city's image has.

Step two

Current page (select to jump to a new page): 1

Logout

Guideline

Officials Will Tour A Changed New York

Abstract: United States Olympic Committee task force will visit New York City to examine potential as host city for 2012 Summer Olympics; New York City competes against San Francisco, Washington-Baltimore and Houston amongst American candidates (M)

For the first time since the terrorist attacks of Sept. 11, a United States Olympic Committee task force will visit New York on Sunday and Monday to evaluate its worthiness to be America's designated host city competing for the 2012 Summer Games.

The Olympic plan laid out by NYC2012 for the task force's visit last summer, which advanced New York into a pool of four candidate cities, has not changed much.

But the city's image has.

Example one

In example 1, *New York; is one of the four candidate cities; competing to be presented to the IOC* should not be chosen because the predicate isn't something that can be considered as an event. An event is essentially a change of state. Predicates like "is" is only describing one state. On the other hand, *New York; is competing; with three cities to be presented to the IOC* should be chosen because the predicate can indicate an event.

- The task force; will have; a critical five-hour session at City Hall on Monday
- NYC2012; has added; more detail to its athletic, housing, transportation, finance, and security plans
- NYC2012; has not changed; the proposed locations of any sports
- NYC2012; has hired; consultants to study urban design and zoning of the far West Side of Manhattan
- Two other studies; are underway; to analyze financing for the subway extension and the building of platforms over the Long Island Rail Road storage yards
- NYC2012; expects; to spend about \$13 million on its bid
- NYC2012; has raised; \$11 million so far
- Terrorists; attacked; New York
- Daniel L. Doctoroff; said; "we don't want any sympathy for that"

Enter a new event you want to add

Officials Will Tour A Changed New York

Abstract: United States Olympic Committee task force will visit New York City to examine potential as host city for 2012 Summer Olympics; New York City competes against San Francisco, Washington-Baltimore and Houston amongst American candidates (M)

For the first time since the terrorist attacks of Sept. 11, a United States Olympic Committee task force will visit New York on Sunday and Monday to evaluate its worthiness to be America's designated host city competing for the 2012 Summer Games.

The Olympic plan laid out by NYC2012 for the task force's visit last summer, which advanced New York into a pool of four candidate cities, has not changed much.

But the city's image has.

"We don't want any sympathy for that," said Daniel L. Doctoroff, who founded NYC2012, a nonprofit group that is lobbying for the Games, and who is now the city's deputy mayor for economic development.

After Sept. 11, he said, "The rest of the country saw in New York the true face of the city: the compassion, the courage and the resilience."

The premise of the Olympic bid is that New York's spirit and ethnic diversity, as well as its ability to handle security for enormous events, embody Olympic ideals.

New York's rivals — San Francisco, Houston and Washington-Baltimore — offer their own variations of the patriotic Olympic sentiments.

The four cities are competing to be among the two chosen by the U.S.O.C. task force in September. Then, in November, the full Olympic Committee board will select the city it will present to the International Olympic Committee; the I.O.C. will make a final choice in 2005.

The task force will tour some athletic sites and the location of the proposed Olympic Village in Long Island

Example two

In example 2, *Daniel L. Doctoroff; said; we don't want any sympathy for that* is an event but not a salient event because simply describing someone said something isn't important enough in this article.

- MetroStars; won; Major League Soccer game
- MetroStars; won; against New England Revolution
- MetroStars; won; with a score of 2-1 in overtime
- Lothar Matthaus; was able to play; due to a slipped disk
- Lothar Matthaus; played; as a substitute
- Lothar Matthaus; provided; the assist for Adolfo Valencia's game-winning goal
- Adolfo Valencia; scored; the game-winning goal in the final minutes of regulation time
- Adolfo Valencia; scored; the game-winning goal in overtime
- Adolfo Valencia; scored; with a powerful header
- MetroStars; have; an 11-point lead over the second-place Revolution
- MetroStars; were able to rally; despite missing key midfielders Tab Ramos and Roy Myers
- Game; was marked; by physical play

MetroStars Pull Out Victory to Tighten Grip on Division

Abstract: MetroStars defeat New England Revolution, 2-1 (M)

With Lothar Matthaus playing in Major League Soccer for the first time in six weeks because of a slipped disk, the MetroStars ended a two-game losing streak with a 2-1 overtime victory over the New England Revolution tonight at soaked Giants Stadium before a crowd announced at 12,688.

Adolfo Valencia scored both goals for the MetroStars (2-1-9-2), including the so-called golden goal three minutes into the overtime. Valencia got his winner with a powerful header from 6 yards after a corner kick on the right by Clint Mathis. Mathis, who had a sprained right ankle, was used as a substitute for the first time this season, entering the game in the 59th minute.

"It was important for me to play when I'm really fit, it's nice to see that the team can win with me," Matthaus said with a smile. "Now, I can concentrate on the MetroStars and nothing else."

Coach Octavio Zambrano was impressed with the way his MetroStars rallied after falling behind by a goal with three minutes left in regulation.

"I'm extremely proud of the guys," he said. "They pulled through and they were rewarded at the end."

Matthaus started the play that led to Valencia's first goal, which tied the score with two minutes left in regulation and one minute after New England's Wolde Harris had opened the scoring. Matthaus gave the ball to Mathis for the pass to Valencia, who beat Jurgen Sommer with a low shot from close range for his team-leading 12th goal of the season.

The victory, the MetroStars' third of the season over New England, strengthened their position atop the Eastern Division, where they have an 11-point lead over the second-place Revolution (9-11-6).

The MetroStars looked headed for defeat when Harris scored for the fifth consecutive game. Taking a pass from Mauricio Wright after a corner kick by John Harkes, Harris beat Mike Ammann from 10 yards. But Matthaus rallied the MetroStars with his pass to Mathis.

Example three

In example 3, *MetroStars; won; Major League* isn't selected because it is less accurate and informative than the second option. They are referring to the same event and we don't want duplication. The fourth and fifth options are not selected due to they are more about a description of state than a change of state. *Game; was marked; by physical play* isn't selected because there is no direct reference in the article that the game is marked by physical play.

Below are more annotated examples.

- falling commodity prices; contributed; to the economic slump
- a decrease in American tourists; contributed; to the economic slump
- the budget surplus; turned into; a large deficit
- the government; responded; by reducing the provincial civil service by 28%
- the government; implemented; a three-year spending freeze on health care and education
- the finance minister; claimed; that the tax cuts will stimulate consumer spending and business investment

British Columbia's Liberals Deliver a Tax Cut, Then Pay Dearly

Abstract: britr

Residents of British Columbia will receive a big cut in their income taxes on New Year's Day, their second in six months. But far from winning applause, the province's governing Liberal Party is experiencing a drop in public support.

Tax cuts were a major promise in the campaign that brought the Liberals to office in the province in a landslide election win last May. The Liberals, led by Gordon Campbell, a former secondary school teacher, real-estate executive and three-term mayor of Vancouver, won all but two of the 79 seats in the western province's legislative assembly. They defeated the left-leaning New Democrats, whose 10 years in office were marked by a growing public role in the economy and numerous tales of economic mismanagement.

With the Jan. 1 reductions, personal income tax rates for provincial taxes — which constitute a much larger

Figure 6: Annotation guidelines of salient event identification shown to the annotators.

Guideline

Back to annotation page

Step one

Select a page from the dropdown menu to start annotating. You can also use the 'Previous' and 'Next' buttons at the bottom to navigate through the pages.

Each page shows a news article. The large font text is the title. The bold font text is the abstract. The rest is the main content.

You should primarily use the words and phrases from the main content to construct events. Please don't repeat events from the title or the abstract.

Step two

You could add event relations to the left sidebar. To add a relation, first choose a source event, then choose a relation type and a target event. There are three options for relation type. Relation **is_caused_by** means the source event is strictly dependent on the target event. The source event wouldn't happen if the target event doesn't happen. Relation **happened_before** means the source event happened before the target event as described by the article. It's the temporal relation in the real world instead of the occurrence order in the article. Relation **is_subevent_of** means the source event is the subevent of the target event.

Please find all the relations that are described in the article. Including the relations that are explicitly mentioned, and the relations that can be inferred based on the evidence presented in the article. If logical inference cannot fully support the relation, please don't include it. Every modification will be automatically saved on the server.

Example

*1- (The task force; will tour; some athletic sites and the proposed Olympic Village in Long Island City); **happened_before**; (The task force; will have; a critical five-hour session at City Hall on **Monday**)

Example one

In example 1, (The task force; will tour; some athletic sites and the proposed Olympic Village in Long Island City) will be on Sunday and (The task force; will have; a critical five-hour session at City Hall) will be on Monday. They are explicitly related by time.

*1- (Gareth H. Edmondson-Jones; will be flying back; on a new jet from the Airbus factory in Toulouse, France); **is_subevent_of**; (Gareth H. Edmondson-Jones; took a European vacation; to avoid the disruption of the Republican National Convention)

Source Event: Gareth H. Edmondson-Jones; will be flying back; on a new jet fro...

Relation: is_subevent_of

Target Event: Gareth H. Edmondson-Jones; took a European vacation; to avoid...

Please don't add the same relation twice

Example two

In example 2, (Gareth H. Edmondson-Jones; will be flying back; on a new jet from the Airbus factory in Toulouse, France) is part of the vacation in the event (Gareth H. Edmondson-Jones; took a European vacation; to avoid the disruption of the Republican National Convention).

Current page (select to jump to a new page): 1 | Logout | Guideline

Officials Will Tour A Changed New York

Abstract: United States Olympic Committee task force will visit New York City to examine potential as host city for 2012 Summer Olympics, New York City competes against San Francisco, Washington-Baltimore and Houston amongst American candidates (M)

For the first time since the terrorist attacks of Sept. 11, a United States Olympic Committee task force will visit New York on Sunday and Monday to evaluate its worthiness to be America's designated host city competing for the 2012 Summer Games.

The Olympic plan laid out by NYC2012 for the task force's visit last summer, which advanced New York into a pool of four candidate cities, has not changed much.

But the city's mayor has...

Step one

Officials Will Tour A Changed New York

The task force will tour some athletic sites and the location of the proposed Olympic Village in Long Island City on Sunday. But the critical part of the visit will be a planned five-hour session at City Hall on **Monday** in which the task force will pose questions to a group that will include NYC2012 officials, some Olympians and Police Commissioner Raymond W. Kelly.

"We'll get grilled," Doctoroff said.

Step two

Figure 7: Annotation guidelines of relation identification shown to the annotators.