

# Set-Aligning Fine-tuning Framework for Auto-Regressive Event Temporal Graph Generation

Anonymous ACL submission

## Abstract

Event temporal graphs have been shown as convenient and effective representations of complex temporal relations between events in text. While traditional methods are based on a pipeline approach, i.e., event extraction and relation classification, the recently proposed contextualised graph generation methods have shown promising results by employing pre-trained language models to generate linearised graphs as autoregressively. However, this inevitably led to sub-optimal graph generation as the linearised graphs exhibit set characteristics which are instead treated sequentially by language models. This is due to their conventional text generation objectives which end up mistakenly penalizing correct predictions only because of the misalignment between elements in text. In this work, we reformulate the task as a conditional set generation problem, proposing a Set-aligning Fine-tuning Framework allowing smooth employment of large language models. A comprehensive experimental assessment has shown that our proposed framework significantly benefits the event temporal graph generation, and outperforms existing baselines. We further demonstrate that under the zero-shot settings, the structural knowledge introduced through the proposed framework has a significant beneficial impact on model generalisation when the examples available are limited<sup>1</sup>.

## 1 Introduction

Understanding the temporal relation between events mentioned in long documents is crucial to modelling complex text with articulated narratives. One of the widely adopted benchmarks for event temporal relation understanding is still the SemEval 2013 TempEval-3 (UzZaman et al., 2013), requiring end-to-end generation of event temporal graphs directly from raw text. An event temporal graph is a natural representation of temporal information,

<sup>1</sup>The experiment code will be released upon the paper’s acceptance. The human annotations will be released as span indices in agreement with the NYT corpus’s license.

with the nodes representing events and the edges the temporal relationships between them, such as “before”, “after”, or “simultaneous”.

Most existing approaches typically address the problem of extracting event temporal graphs through a two-step pipeline, with the first step focusing on detecting events in text, and the second step on classifying the temporal relations between them (McDowell et al., 2017; Ning et al., 2018b). However, such pipeline-based approaches suffer from well-known limitations, including (i) the need for fine-grained annotations at each step; and (ii) the potential for error propagation throughout the pipeline. In particular, in the first step, the event extractor’s objective is to locate as many event triggers as possible in the given documents, leading to the inclusion of numerous trivial events that often lack relevance to the narrative and have no relation to other events. As a result, the next step for temporal relational extraction becomes burdened with many noisy events, significantly impacting the overall accuracy and efficiency of the models.

To address these limitations, Madaan and Yang (2021) introduced a reformulation of the task by generating event temporal graphs directly through conditional text generation. This approach allows for the use of large pre-trained language models and, more importantly, overcomes the typical limitations associated with the pipeline architecture. While this method involved fine-tuning a text generation model, such as GPT-2, for the generation of linearised event temporal graphs as sequences, it fails to consider an important aspect. Specifically, it does not account for the fact that the target sequence (i.e. the list of event temporal relations) is order-invariant, and should therefore be treated as a *set* rather than as an ordered sequence. For example, the following two sequences represent the same temporal graph:

080	S1: [(Cuomo leaving his office, before, speak to reporters),	loss functions, data augmentation, and weak	129
081	... (Cuomo leaving, before, met with representatives)]	supervision techniques.	130
082	S2: [(Cuomo leaving, before, met with representatives),	• We offer a human-annotated test set and	131
083	... (Cuomo leaving his office, before, speak to reporters)]	a weakly-supervised dataset specifically de-	132
084		signed for document-level event temporal gen-	133
085	In this scenario, the conventional loss function will	eration.	134
086	(mistakenly) yield a high value because most of the	• We conduct an extensive evaluation of the fine-	135
087	tokens in the corresponding positions do not match,	tuned models under various settings, demon-	136
088	even though the event relations are the same. This	strating the effectiveness of the proposed	137
089	issue has a detrimental effect on the model perfor-	framework and the potential of contextualised	138
090	mance for several reasons. First, it discourages the	event graph generation.	139
091	language model from generating additional edges.		
092	Generating more edges implies a greater number	<b>2 Related Work</b>	140
093	of potential permutations in the edge sets, making		
094	it less likely to match the target. Secondly, if the	<b>2.1 Event Temporal Graph</b>	141
095	initially generated edge in the sequence differs in		
096	token count from the one in the target, it causes all	The task of event temporal graph extraction serves	142
097	subsequent edges to misalign with the target, even	as an important task for evaluating an end-to-end	143
098	if they are identical, leading to a high loss value.	system which takes raw text as input and output	144
099	In this work, we propose a Set-Aligning Fine-	TimeML annotations (i.e., temporal relations) (Uz-	145
100	tuning Framework (SAFF) that enables efficient	Zaman et al., 2013). Early attempts on the task	146
101	employment of LLMs for auto-regressive event	include CAEVO (McDowell et al., 2017) and Cog-	147
102	temporal graph generation. SAFF incorporates a	comptime (Ning et al., 2018b), which relied on	148
103	novel set of losses, named Set Property Losses,	a combination of statistical and rule-based meth-	149
104	along with augmented data, aimed at mitigating the	ods. In recent years, more efforts have been put	150
105	challenges associated with conventional text gener-	into developing specialised sub-systems with neu-	151
106	ation loss. Using the proposed SAFF, we fine-tune	ral network-based approaches (Ning et al., 2019;	152
107	language models from the T5 (Raffel et al., 2020)	Han et al., 2019a; Tan et al., 2021a). The emer-	153
108	family with weak supervision. Additionally, we	gence of large language models has paved the way	154
109	introduce the first human-annotated dataset built	for end-to-end learning, treating temporal graph	155
110	on the New York Times for contextualised event	generation as conditional text generation (Madaan	156
111	temporal graph generation, which we combine with	and Yang, 2021). To tackle the set misalignment	157
112	existing event relation extraction datasets to evalu-	issue which remained unexplored in Madaan and	158
113	ate the effectiveness of the SAFF framework. Ex-	Yang (2021), we propose a framework based on a	159
114	periments on the newly annotated New York Times	novel set of Set-Aligning losses, aiming at enhanc-	160
115	corpus <sup>2</sup> show that SAFF significantly increases the	ing autoregressive event temporal graph generation.	161
116	number of generated edges, resulting in improved	It is worth noting that there is another related	162
117	recall. Furthermore, we assess the performance	and more widely-recognised task called <i>temporal</i>	163
118	of our approach on existing sentence-level event	<i>relation extraction</i> , which aims at classifying the	164
119	temporal relation extraction datasets, namely MA-	type of temporal links between pre-extracted events	165
120	TRES (Ning et al., 2018a) and TB-Dense (Cassidy	(Wang et al., 2020; Wen and Ji, 2021; Tan et al.,	166
121	et al., 2014), under zero-shot settings, and we find	2023). While Han et al. (2019b) proposed a joint	167
122	that the structural knowledge introduced through	extraction model for events and event temporal	168
123	the proposed SAFF has an even greater impact on	relations, they rely on event extraction supervision	169
124	model generalisation when the examples available	signals, which our work does not need.	170
125	are limited.		
126	Our contributions are three-folded:	<b>2.2 Graph Generation with Language Models</b>	171
127			
128	• We introduce a model-agnostic framework,	Generating graphs with language models has been	172
	called SAFF, for event temporal graph genera-	explored in many areas. For example, Bosselut	173
	tion. SAFF incorporates novel Set-Aligning	et al. (2019) fine-tunes GPT on the ATOMIC com-	174
		monsense knowledge graph (Sap et al., 2019). Mel-	175
		nyk et al. (2022) proposed a multi-stage system	176
		for knowledge generation based on T5. However,	177

<sup>2</sup><https://doi.org/10.35111/77ba-9x74>

these studies do not generate an entire graph in one generation. In contrast, [Madaan et al. \(2021\)](#) generated inference graphs using a combination of a graph generator and a graph corrector for queries in defeasible reasoning. [Zaratiana et al. \(2023\)](#) generate entities and entity relations with an autoregressive LM, but they did not consider the set property of the target. Different from them, we focus on the set property of the generation sequence, which is particularly important in the setting where both the input document and output sequence are considerably longer.

### 2.3 Conditional Set Generation

Text generation models are primarily designed for generating text with strict linear orders, making them suboptimal for generating sets. This limitation has been acknowledged in recent NLP research, where efforts have been made to adapt seq2seq frameworks for tasks like multi-label classification and keyword generation ([Qin et al., 2019](#); [Ye et al., 2021](#)). [Vinyals et al. \(2016\)](#) studied the general challenge of using sets as either input or target output for text generation models. They found in both cases, the order of elements in the set has a significant impact on convergence and final perplexity. This implies that there may exist an optimal order for the input or output set sequence, and they proposed allowing the model to search for this order during training. Instead of resorting to exhaustive search, [Madaan et al. \(2022\)](#) proposed to use data augmentation to enforce order-invariance and prepend the set’s cardinality to the target sequence to ensure the correct cardinality. While previous research has tackled multi-label prediction and keyphrase generation, our work delves into the unique challenges presented by event temporal graph generation, which involves long sequences and partially ordered properties.

In a more general sense, the object detection task from computer vision also involves set prediction. They use parallel decoding to generate the elements in a set based on object queries ([Carion et al., 2020](#); [Chen et al., 2022](#)). [Tan et al. \(2021b\)](#) adopted a similar approach in name entity recognition with a non-autoregressive decoder. Different from them, the set elements (event relations) in our task are not concrete spacial objects or text spans but instead are varied in length and scattered across each document. This makes object queries and non-autoregressive decoding inapplicable in our settings.

## 3 Set-Aligning Fine-tuning Framework

[Madaan and Yang \(2021\)](#) first explored the possibility of end-to-end event temporal graph generation using neural language modelling. Since then, however, this task has remained under-explored, with numerous unresolved issues. To elaborate, the first concern is that [Madaan and Yang \(2021\)](#) framed graph generation as a conventional sequence generation problem, whereas it is fundamentally a set generation problem. Secondly, the dataset they built primarily consists of small-sized graphs, failing to challenge the model in terms of document-level understanding. Lastly, their investigation mainly centred on GPT-2, while the landscape of LLMs has evolved with the emergence of models featuring distinct structures (e.g., encoder-decoder) and new paradigms (e.g., in-context learning) in recent years. In this study, we address these three aspects to enhance the understanding of sequence-to-sequence temporal graph generation.

Although our proposed framework is designed to be model-independent, several factors have led us to choose Flan-T5 as the base model for our experiments: (i) Based on our preliminary experiments, Flan-T5-base hits the sweet spot in terms of performance vs. resource consumption, allowing us to test more variants; (ii) its encoder-decoder structure is well-suited to document-level graph generation, due to its efficiency in processing comprehensive information in lengthy documents.

### 3.1 Event Temporal Graph Modelling as Edge Set Generation

An event temporal graph is a directed graph with no isolated vertex. Each edge in the graph describes a temporal relation between two events, and self-loops are not permitted. Following [Madaan and Yang \(2021\)](#), we represent these graphs by linearizing them into strings using the DOT graph description language ([Gansner, 2006](#)) (example shown in [Figure 1](#)). Given that event temporal graphs do not have isolated vertices, the sequence essentially represents the edge set of the graph.

We model the probability of generating a string  $y$ , which is a linearised representation of the event temporal graph  $G$ , conditioned on a document  $X = (x_1, x_2, \dots, x_n)$  using a language model:

$$p_{\text{LM}}(y|X) = \prod_{t=1}^T p(y_t|X, y_{<t}) \quad (1)$$

where  $y$  is a string formatted in DOT notation.

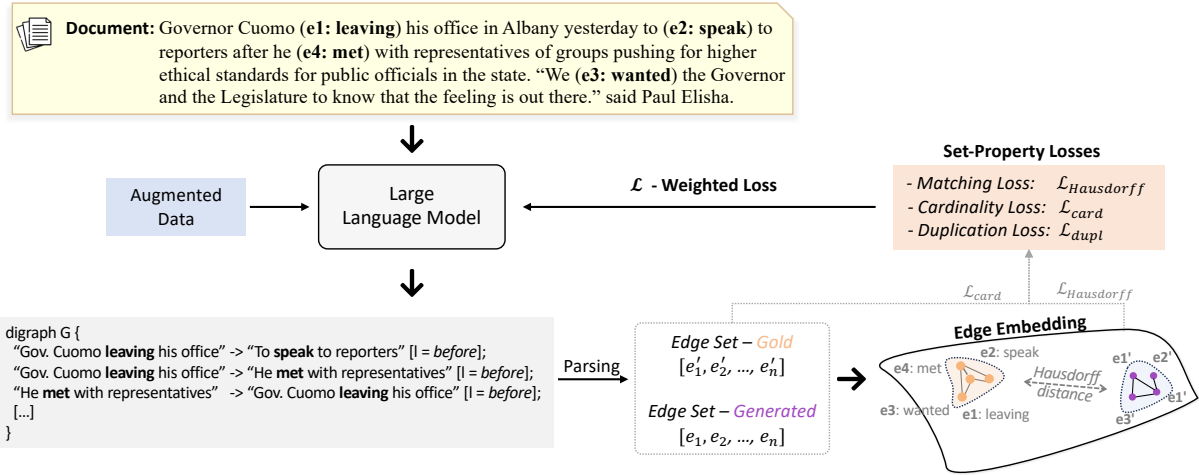


Figure 1: Set-Aligning fine-tuning framework (SAFF).

### 3.2 Data Augmentation

The target sequences of event temporal graph generation are essentially sets rather than strictly ordered text sequences. Therefore, conventional text generation loss can inadvertently penalize the token order and force the arrangement of elements to match the order in the target sequence, which is not necessarily the optimal order. This enforced order may lead to sub-optimal performance (Vinyals et al., 2016). A potential solution is to introduce random permutations of set elements as augmented training examples, which has already been shown effective in tasks like multi-label classification and keyphrase generation (Madaan et al., 2022). Specifically, in the context of event temporal graph generation, the elements correspond to the edges in the target string. The substrings representing the edges are randomly shuffled, while the rest of the string remains unchanged.

Prepending the set cardinality of the ground-truth edge set to the generation target may also help constrain the generation model to avoid over-generation (Madaan et al., 2022). However, such attempts in our preliminary experiment led to an approximate 4% drop in edge  $F_1$  score, despite a significant reduction in the number of generated edges. Thus, we decided not to incorporate the cardinality into the final framework.

### 3.3 Set Property Losses (SPL)

Simply adding augmented data to train models does not address the fundamental issue of set alignment. Several challenges arise in this approach. First of all, it is unrealistic to add all permutations, especially when dealing with long documents containing numerous event relations, as the training data

will grow at a rate proportional to the factorial of the cardinality of the target set. More importantly, with each augmented example, the loss function would still penalize the unobserved permutations of the set. This would make the training unstable.

The core challenge lies in finding an effective way to compare the linearized target graph with the linearized generated graph, without relying on a strict token-by-token comparison as in conventional text generation. To tackle this issue, we propose introducing modifications to the generation objective. As the linearized graph essentially represents the edge set of the graph, we can simplify the graph comparison problem into a set comparison problem. Our approach involves several components. Firstly, we add a set cardinality loss to encourage the model to generate an adequate number of temporal relation edges. Then, we introduce a duplication loss to penalize any repetition of elements in the edge set. Lastly, we design a set matching loss that assesses the semantic similarity between elements in the target edge set and those in the generated edge set. Collectively, the above loss functions are referred to as Set Property Losses (SPL). The SPL is integrated with the conventional token-level cross-entropy loss through a weighted average.

To compute the set property losses, a graph string needs to be first sampled from a language model given a training input. Then, this sequence is parsed into a list of edges  $E$ , where each edge  $e$  is a triplet consisting of a head event, a relation type, and a tail event  $(h, r, t)$ . Now, the number of edges and duplicated edges can be counted. Let  $\mathcal{E}$  denote the set of all the unique edges in  $E$ . The values for the set cardinality loss and the duplication loss can be computed as follows:



$$\mathcal{E} = \{e | e \in E\} \quad (2)$$

$$L_{\text{dupl}} = \frac{|E| - |\mathcal{E}|}{|\mathcal{E}|} \quad (3)$$

$$L_{\text{card}} = \frac{\text{abs}(|\mathcal{E}'| - |\mathcal{E}|)}{|\mathcal{E}|} \quad (4)$$

where function  $\text{abs}(\cdot)$  denotes taking the absolute value,  $\mathcal{E}'$  denotes the ground-truth edge set.

To compute the set matching loss, we assess the similarity between the generated set and the target set by comparing the semantic similarity of the edges across the two sets. We take the last layer of the decoder’s representations of the respective tokens as the semantic representations of the events and the relation type. Then, we concatenate these representations as the semantic representation of each edge:

$$z_h = H_{[h_1, h_2, \dots, h_m]} \quad (5)$$

$$z_r = H_{[r_1, r_2, \dots, r_s]} \quad (6)$$

$$z_t = H_{[t_1, t_2, \dots, t_n]} \quad (7)$$

$$\bar{e} = [\text{pool}(z_h); \text{pool}(z_r); \text{pool}(z_t)] \quad (8)$$

where  $H$  is the last-layer hidden states of the decoder.  $[h_1, \dots, h_m]$ ,  $[r_1, \dots, r_s]$ , and  $[t_1, \dots, t_n]$  are the indices of the head event, relation type, and tail event, respectively.  $z_h, z_r, z_t$  denote the semantic representations of the head event, relation type, and tail event, respectively.  $\text{pool}(\cdot)$  represents the average pooling function.  $\bar{e}$  denotes the semantic representation of the edge.

We now possess two sets of embeddings: one compassing the edge embeddings extracted from the target graph, and the other containing the edge embeddings derived from the generated graph. Essentially, they can be considered as two sets of points in the representation space. Thus, we can measure the similarity of the two graphs by measuring the distance between the two point sets (manifolds) in the representation space. The Hausdorff distance, originally defined to measure the separation between two subsets within a metric space, has recently found applications in machine learning for measuring the distance between two sets of embeddings (Schutze et al., 2012; Wang et al., 2023). We compute the average Hausdorff distance as the measure:

$$d_H(\mathcal{E}', \mathcal{E}) = \frac{1}{|\mathcal{E}'|} \sum_{\bar{e}' \in \mathcal{E}'} \min_{\bar{e} \in \mathcal{E}} d_{\text{cos}}(\bar{e}', \bar{e}) + \frac{1}{|\mathcal{E}|} \sum_{\bar{e} \in \mathcal{E}} \min_{\bar{e}' \in \mathcal{E}'} d_{\text{cos}}(\bar{e}', \bar{e}) \quad (9)$$

where the distance of an edge pair is computed by the cosine distance  $d_{\text{cos}}(\cdot)$ .

### 3.4 Fine-tuning with Set Property Losses

Unlike the set prediction methods based on parallel decoding (Carion et al., 2020; Tan et al., 2021b), set-based objectives, such as SPL, cannot be directly used as the main objective in autoregressive generation. There are two primary reasons for this. The first reason is that obtaining the SPL requires sampling from the decoder, which would reduce the training speed significantly. Moreover, the second reason is that the language model will struggle to generate sequences in DOT format accurately because learning the token dependency for such format requires the language modelling objective. Consequently, the sequence parser will fail to recognize any valid edges within the sequence, resulting in high SPL values and hindering the training.

To avoid the problems mentioned above, we introduce the SPL after a certain number of fine-tuning iterations. Once the model has acquired a basic proficiency in generating correct DOT sequences, the SPL can function as intended. SPL can also be viewed as a regularization to prevent LLMs from overfitting to the order of the target set shown in the training samples.

We explored alternative approaches to incorporate SPL, but they reported inferior performance compared to the method eventually included in our framework. We discuss those alternative methods in the Appendix A.

## 4 Experiment

### 4.1 NYT Temporal Event Graph Dataset

There are several event temporal relation extraction datasets with pairwise event relation annotations, such as MATRES and TBD. It is theoretically possible to convert these annotations into document-level event temporal graphs. However, our preliminary experiments have shown that even when merging all of these datasets (resulting in 4,684 training documents), it is not sufficient to fine-tune a large language model to achieve acceptable performance. To address this limitation, we opted to

	NYT-train	NYT-test	NYT-human
Total documents	18,263	1,000	22
Total events	846,022	47,251	661
Node degree	2.52	2.54	2.34
Total relations	1,066,264	60,056	528
<i>before</i>	578,216	32,729	465
<i>after</i>	412,704	23,200	0
<i>includes</i>	7,922	450	12
<i>is_included</i>	41,964	2,332	0
<i>simultaneous</i>	25,458	1,345	51

Table 1: The statistics of the NYT temporal event graph dataset. Node degree represent the average number of relations each event has.

build a significantly larger dataset on a selection of data from the New York Times (NYT) corpus using a weak supervision approach, drawing inspiration from the work of Madaan and Yang (2021). Nevertheless, we introduced additional steps in the data selection process to ensure that the selected documents contain high-quality event temporal graphs, which were not taken in Madaan and Yang (2021).

Firstly, we performed topic modelling using Latent Dirichlet Allocation (LDA) on the MATRES and TBD datasets to extract a set of topics. Then, we identified general descriptors that are semantically similar to these topics (e.g., politics, diploma, sports, etc.). This selection process was crucial because, following training with noisy labels, our intention was to evaluate the model’s performance on these datasets under zero-shot settings. We further analysed the most noteworthy events in these descriptors to ensure they were narrative-oriented, because articles that weave stories tend to contain a wealth of event temporal relations. To identify the most significant events, we employed a metric similar to TF-IDF which we could describe as “event frequency  $\times$  inverse-descriptor frequency”.

$$ef \cdot idf = \frac{f_{\epsilon,d}}{\sum_{\epsilon' \in d} f_{\epsilon',d}} \cdot \log \frac{|D|}{|\{d \in D : \epsilon \in d\}|} \quad (10)$$

where  $\epsilon$  is an event and  $d$  is a descriptor.  $f_{\epsilon,d}$  is the number of times that event  $\epsilon$  occurs in the documents with the descriptor  $d$ .  $\sum_{\epsilon' \in d} f_{\epsilon',d}$  is the total number of event occurrence in the descriptor  $d$ .  $|D|$  is the total number of descriptors in the corpus.  $|\{d \in D : \epsilon \in d\}|$  is the number of descriptors where the event  $\epsilon$  appears.

The descriptors that are selected and the number of documents in them are listed in the Appendix D.1. After choosing the documents, we acquire the event temporal graph by running an off-the-shelf

event and temporal relation extraction tool called CAEVO (McDowell et al., 2017). CAEVO is more scalable than Cogcomptime (Ning et al., 2018b), making it suitable for building a large-scale dataset.

Then, each temporal graph is represented in DOT format, and every event verb is prefixed and suffixed with its noun phrase and object, respectively. Note that we did not break the documents into short segments as Madaan and Yang (2021) did. Instead, we keep the data strictly at the document level which is a more challenging setting because the model needs to analyse the entire document and generate a much larger graph. In the dataset we built, a target graph has about 46 nodes and 58 edges on average. While in Madaan and Yang (2021), the average number of nodes is 4 and the average number of edges is 5 in a document-level event temporal graph. Moreover, their events have 1.54 relations on average, while events in our data have 2.52 relations on average, showing that the graphs in our dataset are much more complex. In practice, these complex documents are usually the ones that require analysis, and a model developed based on simpler inputs cannot handle them directly.

## 4.2 Human-annotated Test Data

Aside from testing with the CAEVO-created data, we recruited human annotators to annotate a test split of the NYT data. We performed a preprocessing step regarding the relation types by merging the reciprocal relations, such as transforming *after* into *before*, *is\_included* into *includes* by swapping the head and tail events. For example, “I had dinner after I had lunch” is equivalent to “I had lunch before I had dinner”. This processing not only streamlined the annotation process but also enhanced the model performance (refer to experimental results in Appendix C). We recruited crowd workers from Prolific<sup>3</sup> platform, which is a research-focused platform providing verified human workers. We recruited 24 participants in total (including pilot testing runs). To make sure the participants can understand and annotate the article efficiently, we only recruited native English speakers who have an education level higher than High school diploma/A-levels. We put 4 documents, which are randomly sampled from the same descriptor set as the training and testing of the selected NYT corpus, into each unit task. There is a shared document across all the tasks to compute the inter-annotator agree-

<sup>3</sup>prolific.com

ment (IAA). To minimize discrepancy, we asked 2 participants to first identify the event triggers in each unit task. We then merged the event annotations from the participants by taking the union of the spans (if there are overlapped spans, we take the longer span). Then, we asked another participant to annotate the event temporal relation based on the identified events. We also included the outputs from the CAEVO model to serve as examples, but we explicitly asked the participants to correct the annotations by adding, removing, or changing the CAEVO’s annotations. In the end, we collected 22 documents as the human-annotated test set. On the event identification, we compute IOU (Intersection over Union) as a measure of agreement between the annotators. Average across 7 tasks, the IOU between the event spans is 0.8986. For the relation annotations, we compute the average Cohen’s  $\kappa$  of every participant pair in the relation annotation task (on the shared document). The average Cohen’s  $\kappa$  is 0.7465. Details of instructions and interfaces are in Appendix D.1.

The statistics of the constructed datasets are shown in Table 1. The distributions of relation types are highly imbalanced, with a majority falling into either the *before* or *after* categories. We also evaluated the trained models on the MATRES test set (comprising 20 documents) and TBD test set (consisting of 9 documents), both of which are based on human annotations and processed into DOT using the methods previously described.

### 4.3 Model Setting

We use Flan-T5-base (Chung et al., 2022) (250M) as our backbone model. The model is trained for 10 epochs, with each document being augmented through 4 random permutations, followed by a further 3 epochs of training, during which the SPL loss is adopted without permutations. We use a learning rate of  $2e - 5$ , along with a weight decay of 0.01. Batch size of 5 before SPL, and 3 during SPL because additional memory is required for sampling. We optimize the loss with the AdamW optimizer (Loshchilov and Hutter, 2019). We use the beam search (Graves, 2012) with a beam size of 5 and a maximum length of 2048 to sample results. We balanced the training steps in the compared methods to make sure they saw the same amount of training data. Experiments are conducted on a GPU node under an HPC cluster using 4 Nvidia A100 GPUs. The models are trained based on 3 random seeds (ChatGPT was tested for 3 times) and the metrics are the average values of them.

	NYT-test			NYT-human		
	$P^E$	$R^E$	$F_1^E$	$P^E$	$R^E$	$F_1^E$
Flan-T5-base	51.27	32.43	39.73	22.61	25.88	24.14
SAFF (w/o aug)	50.28	34.82	41.15	25.80	32.13	28.62
SAFF (w/o SPL)	<b>51.88</b>	36.64	42.95	<b>27.08</b>	34.91	30.50
SAFF	50.97	<b>39.96</b>	<b>44.80</b>	25.92	<b>40.21</b>	<b>31.52</b>

Table 2: Edge-based metrics on the NYT datasets

	NYT-test			NYT-human		
	$P^N$	$R^N$	$F_1^N$	$P^N$	$R^N$	$F_1^N$
Flan-T5-base	<b>75.52</b>	58.24	65.76	53.36	47.66	50.35
SAFF (w/o aug)	75.34	60.64	67.20	<b>54.86</b>	50.43	52.55
SAFF (w/o SPL)	75.43	62.36	68.27	54.14	51.59	52.84
SAFF	75.47	<b>65.16</b>	<b>69.95</b>	53.63	<b>54.51</b>	<b>54.06</b>

Table 3: Node-based metrics on the NYT datasets

### 4.4 Evaluation Metrics

Following the previous research (Madaan and Yang, 2021), we evaluate the results using the metrics of precision, recall, and  $F_1$  score for both node set and edge set predictions. The primary metric is the edge  $F_1$  because the quality of the node generation is also reflected in it.

### 4.5 Results

Table 2 and 3 show the evaluation results on NYT-test and NYT-human. Flan-T5-base was trained following the same setup as in Madaan and Yang (2021). SAFF (w/o aug) is the proposed framework without the augmentations of element order and with SPL. SAFF (w/o SPL) is the proposed framework without the use of Set Property Losses (SPL) and with the augmentations. As our SAFF framework with SPL requires additional training steps and the augmentations enlarge the training set, we keep the number of training steps balanced in the methods to exclude the influence of seeing different amounts of training data.

Comparing the results of Flan-T5-base with those of SAFF (w/o SPL), it becomes evident that permuted augmentation improves edge  $F_1$  by about 3% on NYT-test and 6% on NYT-human. While SAFF (w/o aug) shows the SPL alone can improve edge  $F_1$  by about 1.5% on NYT-test and 4.5% on NYT-human. Furthermore, when both SPL and augmentation are incorporated, the graph generation performance is further increased (cf. SAFF (w/o SPL)/SAFF (w/o aug) and SAFF). We can also see that models utilizing SAFF have much higher edge recalls while their edge precision scores are either similar or occasionally even lower than those of other models. This suggests that the

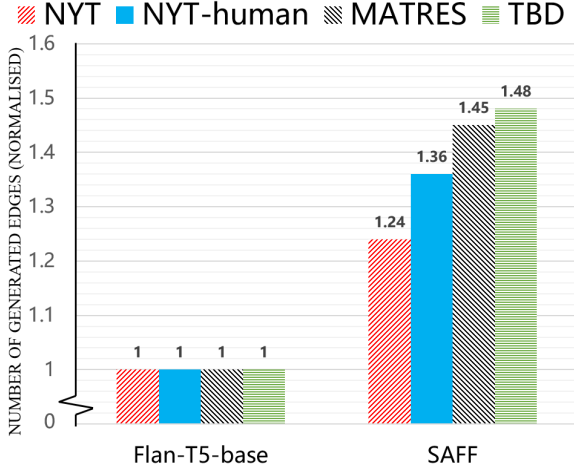


Figure 2: The comparison of generated edges between SAFF and vanilla Flan-T5-base. The  $y$  axis is normalised by dividing the number of edges generated by Flan-T5-base in the respective datasets.

	MATRES-test			TBD-test		
	$P^E$	$R^E$	$F_1^E$	$P^E$	$R^E$	$F_1^E$
ChatGPT	10.58	6.56	8.09	25.92	5.94	9.66
Flan-T5-base	13.06	7.16	9.25	23.26	4.59	7.67
SAFF	<b>18.05</b>	<b>14.31</b>	<b>15.96</b>	<b>37.53</b>	<b>11.04</b>	<b>17.05</b>

Table 4: Experiment results on human-annotated MATRES and TBD under the zero-shot setting.

performance improvement primarily comes from the generation of more edges. This observation is reinforced by the information presented in Figure 2, where models trained with SAFF can generate 24% – 48% more edges compared to the conventional text generation framework on these datasets. These additional edges play a pivotal role in the improvement of the edge  $F_1$  since precision stays nearly the same.

It is worth mentioning that the NYT-human dataset has a different label distribution compared to the NYT dataset used for training, where its events and event temporal relations were produced by CAEVO. Notably, the frequency of *simultaneous* is significantly higher, accounting for 9.66%, in contrast to the 2.39% observed in the training set (see Appendix D.1 for more comprehensive analyses). Based on our observation, it appears that human annotators tend to apply a more lenient criterion for the *simultaneous* label whereas CAEVO enforces a stricter definition of this label.

Similar trends are also observed in Table 4, which were obtained through evaluation on MATRES and TBD. We used the models trained on the NYT training set to test on these datasets under

zero-shot settings. It is worth noticing that SAFF exhibits superior zero-shot performance compared to GPT-2, as reported in (Madaan and Yang, 2021), despite our SAFF training using data four times smaller and GPT-2 (with 355M parameters) having 40% more parameters.

ChatGPT shows our best attempts to generate event temporal graphs with gpt-3.5-turbo API. We used two hops: (i) ask ChatGPT to generate events from the documents, (ii) ask ChatGPT to generate an event temporal graph based on the generated events and the documents. The results show that ChatGPT is outperformed by fine-tuned models, which is in line with the recent papers on exploring ChatGPT’s ability on event understanding (Li et al., 2023; Chan et al., 2023; Gao et al., 2023). These facts confirm that event temporal graph generation cannot be solved solely with prompt engineering on ChatGPT. We provide the detailed inputs, outputs, and parameter settings in Appendix E.

#### 4.6 Error Analysis

A major error type we found is that the model often fails to deduce temporal relationships that involve inference. This is due to the reliance of weak supervision signals provided by CAEVO, which primarily rely on syntactical rules. Consequently, this problem led to a lower edge  $F_1$  on the human-annotated test set, as human annotators provided many temporal relations that were inferred through commonsense reasoning. Conversely, the model does not perceive a clear temporal sequence in the sentence: “<person A> won the gold medal in women’s 1,500m. <person B> won the silver and <person C> won the bronze.” However, human annotators can readily identify an obvious temporal order among “<person A> won”, “<person B> won”, and “<person C> won”, as it aligns with the common knowledge that in a race, the first person who crossed the finish line won the gold, followed by the silver and the bronze winners.

## 5 Conclusion

This study proposes a framework for fine-tuning large language models to generate event temporal graphs directly from raw documents. We propose data augmentation and set property losses to mitigate the problem caused by conventional generation loss, promoting the generation of more edges by language models and, consequently, leading to improved performance.



## 684 Limitations

685 Due to the presence of noisy labels used in fine-  
686 tuning, a major limitation of the proposed method  
687 is the inclusion of many imaginary events, trivial  
688 events, and negative expressions of events. For ex-  
689 ample, CAEVO identified phrases like “<someone>  
690 did not **fire**” as an event. While “fire” serves as a  
691 predicate and the notion of “did not **fire**” can hold  
692 narrative significance, it may not be entirely suit-  
693 able within the context of event temporal graphs.  
694 This is because it is not about the occurrence of  
695 an action or a change of state, but rather describes  
696 the absence of an event. Similarly, in some articles,  
697 there are descriptions of multiple potential future  
698 developments, such as “he might **buy** product A”.  
699 Including such expressions as events might intro-  
700 duce confusion into the event temporal graph, as  
701 these represent possibilities rather than actual oc-  
702 currences. This problem mainly arises from the  
703 behavior of the CAEVO method, which primar-  
704 ily focuses on identifying fine-grained predicates  
705 as events. The resolution to this problem lies in  
706 obtaining better-quality supervision signals which  
707 focus on salient events (i.e., those events with re-  
708 latively higher occurrences that are important to the  
709 narrative).

## 710 Ethics Statement

711 The proposed method analyses the text provided  
712 and extracts relevant information from it. The al-  
713 gorithm cannot acquire information beyond the  
714 boundary of the given text. Thus, any associated  
715 risks stem solely from the data itself. This research  
716 only utilised publicly available data. As long as the  
717 data input to the model is collected according to the  
718 relevant data policies and guidelines, the proposed  
719 method does not introduce further risks.

## 720 References

721 Antoine Bosselut, Hannah Rashkin, Maarten Sap, Chai-  
722 tanya Malaviya, Asli Celikyilmaz, and Yejin Choi.  
723 2019. [COMET: Commonsense transformers for auto-  
724 matic knowledge graph construction](#). In *Proceedings  
725 of the 57th Annual Meeting of the Association for  
726 Computational Linguistics*, pages 4762–4779, Flo-  
727 rence, Italy. Association for Computational Linguis-  
728 tics.

729 Nicolas Carion, Francisco Massa, Gabriel Synnaeve,  
730 Nicolas Usunier, Alexander Kirillov, and Sergey  
731 Zagoruyko. 2020. End-to-end object detection with  
732 transformers. In *European conference on computer  
733 vision*, pages 213–229. Springer.

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,  
Baltimore, Maryland. Association for Computational  
Linguistics. 734 735 736 737 738 739 740

Chunkit Chan, Jiayang Cheng, Weiqi Wang, Yuxin  
Jiang, Tianqing Fang, Xin Liu, and Yangqiu Song.  
2023. [Chatgpt evaluation on sentence level relations:  
A focus on temporal, causal, and discourse relations](#). 741 742 743 744

Ting Chen, Saurabh Saxena, Lala Li, David J. Fleet, and  
Geoffrey Hinton. 2022. [Pix2seq: A language model-  
ing framework for object detection](#). In *International  
Conference on Learning Representations*. 745 746 747 748

Hyung Won Chung, Le Hou, Shayne Longpre, Barret  
Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi  
Wang, Mostafa Dehghani, Siddhartha Brahma, Al-  
bert Webson, Shixiang Shane Gu, Zhuyun Dai,  
Mirac Suzgun, Xinyun Chen, Aakanksha Chowdh-  
ery, Alex Castro-Ros, Marie Pellat, Kevin Robinson,  
Dasha Valter, Sharan Narang, Gaurav Mishra, Adams  
Yu, Vincent Zhao, Yanping Huang, Andrew Dai,  
Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Ja-  
cob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le,  
and Jason Wei. 2022. [Scaling instruction-finetuned  
language models](#). 749 750 751 752 753 754 755 756 757 758 759 760

Emden R. Gansner. 2006. [Drawing graphs with dot](#). 761

Jun Gao, Huan Zhao, Changlong Yu, and Ruifeng Xu.  
2023. [Exploring the feasibility of chatgpt for event  
extraction](#). 762 763 764

Alex Graves. 2012. [Sequence transduction with recur-  
rent neural networks](#). *ArXiv*, abs/1211.3711. 765 766

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, Hong Kong, China. Association for  
Computational Linguistics. 767 768 769 770 771 772 773

Rujun Han, Qiang Ning, and Nanyun Peng. 2019b. [Joint  
event and temporal relation extraction with shared  
representations and structured prediction](#). In *Pro-  
ceedings of the 2019 Conference on Empirical Meth-  
ods in Natural Language Processing and the 9th In-  
ternational Joint Conference on Natural Language  
Processing (EMNLP-IJCNLP)*, pages 434–444, Hong  
Kong, China. Association for Computational Linguis-  
tics. 774 775 776 777 778 779 780 781 782

Bo Li, Gexiang Fang, Yang Yang, Quansen Wang, Wei  
Ye, Wen Zhao, and Shikun Zhang. 2023. [Evaluating  
chatgpt’s information extraction capabilities: An as-  
sessment of performance, explainability, calibration,  
and faithfulness](#). 783 784 785 786 787

Ilya Loshchilov and Frank Hutter. 2019. [Decoupled  
weight decay regularization](#). In *International Confer-  
ence on Learning Representations*. 788 789 790

791	Aman Madaan, Dheeraj Rajagopal, Niket Tandon, Yiming Yang, and Antoine Bosselut. 2022. <a href="#">Conditional set generation using seq2seq models</a> . In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 4874–4896, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	849
792		850
793		851
794		852
795		853
796		854
797		
798	Aman Madaan, Niket Tandon, Dheeraj Rajagopal, Peter Clark, Yiming Yang, and Eduard Hovy. 2021. <a href="#">Think about it! improving defeasible reasoning by first modeling the question scenario</a> . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 6291–6310, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	855
799		856
800		857
801		858
802		859
803		860
804		
805		
806	Aman Madaan and Yiming Yang. 2021. <a href="#">Neural language modeling for contextualized temporal graph generation</a> . 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.	861
807		862
808		863
809		864
810		865
811		866
812		867
813	Bill McDowell, Nathanael Chambers, Alexander Ororbia II, and David Reitter. 2017. <a href="#">Event ordering with a generalized model for sieve prediction ranking</a> . 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.	868
814		869
815		870
816		
817		
818		
819		
820	Igor Melnyk, Pierre Dognin, and Payel Das. 2022. <a href="#">Knowledge graph generation from text</a> . In <i>Findings of the Association for Computational Linguistics: EMNLP 2022</i> , pages 1610–1622, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	871
821		872
822		873
823		874
824		875
825		
826	Qiang Ning, Sanjay Subramanian, and Dan Roth. 2019. <a href="#">An improved neural baseline for temporal relation extraction</a> . 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.	876
827		877
828		878
829		879
830		880
831		881
832		882
833		
834	Qiang Ning, Hao Wu, and Dan Roth. 2018a. <a href="#">A multi-axis annotation scheme for event temporal relations</a> . In <i>Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1318–1328, Melbourne, Australia. Association for Computational Linguistics.	883
835		884
836		885
837		886
838		887
839		888
840	Qiang Ning, Ben Zhou, Zhili Feng, Haoruo Peng, and Dan Roth. 2018b. <a href="#">CogCompTime: A tool for understanding time in natural language</a> . 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.	889
841		890
842		891
843		892
844		893
845		894
846		895
847	Kechen Qin, Cheng Li, Virgil Pavlu, and Javed Aslam. 2019. <a href="#">Adapting RNN sequence prediction model to multi-label set prediction</a> . In <i>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)</i> , pages 3181–3190, Minneapolis, Minnesota. Association for Computational Linguistics.	896
848		897
		898
		899
		900
		901
		902
		903
		904
		905
		906
	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. <a href="#">Exploring the limits of transfer learning with a unified text-to-text transformer</a> . <i>Journal of Machine Learning Research</i> , 21(140):1–67.	
	Maarten Sap, Ronan Le Bras, Emily Allaway, Chandra Bhagavatula, Nicholas Lourie, Hannah Rashkin, Brendan Roof, Noah A. Smith, and Yejin Choi. 2019. <a href="#">Atomic: An atlas of machine commonsense for if-then reasoning</a> . In <i>Proceedings of the Thirty-Third AAAI Conference on Artificial Intelligence and Thirty-First Innovative Applications of Artificial Intelligence Conference and Ninth AAAI Symposium on Educational Advances in Artificial Intelligence, AAAI’19/IAAI’19/EAAI’19</i> . AAAI Press.	
	Oliver Schutze, Xavier Esquivel, Adriana Lara, and Carlos A. Coello Coello. 2012. <a href="#">Using the averaged hausdorff distance as a performance measure in evolutionary multiobjective optimization</a> . <i>IEEE Transactions on Evolutionary Computation</i> , 16(4):504–522.	
	Xingwei Tan, Gabriele Pergola, and Yulan He. 2021a. <a href="#">Extracting event temporal relations via hyperbolic geometry</a> . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 8065–8077, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	
	Xingwei Tan, Gabriele Pergola, and Yulan He. 2023. <a href="#">Event temporal relation extraction with Bayesian translational model</a> . In <i>Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics</i> , pages 1125–1138, Dubrovnik, Croatia. Association for Computational Linguistics.	
	Zeqi Tan, Yongliang Shen, Shuai Zhang, Weiming Lu, and Yueting Zhuang. 2021b. <a href="#">A sequence-to-set network for nested named entity recognition</a> . In <i>Proceedings of the 30th International Joint Conference on Artificial Intelligence, IJCAI-21</i> .	
	Naushad UzZaman, Hector Llorens, Leon Derczynski, James Allen, Marc Verhagen, and James Pustejovsky. 2013. <a href="#">SemEval-2013 task 1: TempEval-3: Evaluating time expressions, events, and temporal relations</a> . In <i>Second Joint Conference on Lexical and Computational Semantics (*SEM), Volume 2: Proceedings of the Seventh International Workshop on Semantic Evaluation (SemEval 2013)</i> , pages 1–9, Atlanta, Georgia, USA. Association for Computational Linguistics.	
	Oriol Vinyals, Samy Bengio, and Manjunath Kudlur. 2016. <a href="#">Order matters: Sequence to sequence for sets</a> .	

- 907 In *4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May*  
908 *2-4, 2016, Conference Track Proceedings*.
- 910 Haoyu Wang, Muhao Chen, Hongming Zhang, and Dan  
911 Roth. 2020. [Joint constrained learning for event-](#)  
912 [event relation extraction](#). In *Proceedings of the 2020*  
913 *Conference on Empirical Methods in Natural Lan-*  
914 *guage Processing (EMNLP)*, pages 696–706, Online.  
915 Association for Computational Linguistics.
- 916 Xinyu Wang, Lin Gui, and Yulan He. 2023. [Document-](#)  
917 [level multi-event extraction with event proxy nodes](#)  
918 [and hausdorff distance minimization](#). In *Proceedings*  
919 *of the 61st Annual Meeting of the Association for*  
920 *Computational Linguistics (Volume 1: Long Papers)*,  
921 pages 10118–10133, Toronto, Canada. Association  
922 for Computational Linguistics.
- 923 Haoyang Wen and Heng Ji. 2021. [Utilizing relative](#)  
924 [event time to enhance event-event temporal relation](#)  
925 [extraction](#). In *Proceedings of the 2021 Conference*  
926 *on Empirical Methods in Natural Language Process-*  
927 *ing*, pages 10431–10437, Online and Punta Cana,  
928 Dominican Republic. Association for Computational  
929 Linguistics.
- 930 Jiacheng Ye, Tao Gui, Yichao Luo, Yige Xu, and  
931 Qi Zhang. 2021. [One2Set: Generating diverse](#)  
932 [keyphrases as a set](#). In *Proceedings of the 59th An-*  
933 *ual Meeting of the Association for Computational*  
934 *Linguistics and the 11th International Joint Confer-*  
935 *ence on Natural Language Processing (Volume 1:*  
936 *Long Papers)*, pages 4598–4608, Online. Association  
937 for Computational Linguistics.
- 938 Urchade Zaratiana, Nadi Tomeh, Pierre Holat, and  
939 Thierry Charnois. 2023. [An autoregressive text-to-](#)  
940 [graph framework for joint entity and relation extrac-](#)  
941 [tion](#). In *ICML 2023 Workshop on Structured Prob-*  
942 *abilistic Inference & Generative Modeling*.

## A Discussion of Alternative Approaches for Incorporating SPL

In our preliminary studies, we tested several way of incorporating SPL into the training process. We tried using the weighted average of SPL and the language modelling loss in every training step (the weight of SPL increases across the training process). However, the training becomes very slow due to the added decoding processes.

We also experimented with introducing randomness into the incorporation of the SPL loss in some of the training steps. In each training step, there was a 0.5 probability that the SPL loss was computed and backpropagated, while in other cases, only the LM loss was considered. This probability increased progressively with the epoch number. For example, during the initial epoch (epoch 1), the probability was set to 0, as we explained in Section 3.4, the model struggled to generate outputs in the correct format during the early stages. Subsequently, the probability was increased linearly, reaching 1 by the last epoch.

However, this approach proved to be ineffective and caused training instability. The loss value fluctuated between steps, leading to confusion for the model. It is worth noting that this was an initial experiment and significantly differed from the final version of the proposed method.

## B Additional Error Analysis

We observed a type of error involving the model’s incorrect prediction of long-distance temporal relationships. The model sometimes predicts a temporal relation between two events that are separated by more than ten sentences. This is unexpected, as the CAEVO model, which produces weak supervision signals, typically does not extract relations for events that are more than two sentences apart from each other. In essence, it primarily focuses on events within close proximity. Our observations suggest that human annotators also tend not to annotate temporal relations for events that are distant from each other, arguably because such relations are often implicit and can be challenging to track across large chunks of text.

## C More Analysis about the Generation Results

Table A1 shows a preliminary experiment in which the augmented Flan-T5-base models were trained and tested on NYT-test before and after merging

	$P^N$	$R^N$	$F_1^N$	$P^E$	$R^E$	$F_1^E$
With reciprocal relations	76.05	55.05	63.87	52.07	30.90	38.78
Merge reciprocal	75.48	61.96	<b>68.05</b>	52.03	36.34	<b>42.79</b>

Table A1: Comparison between model trained with reciprocal relations or merging reciprocal relations.

	<i>before</i>	<i>includes</i>	<i>simultaneous</i>
Target graph	93.13	4.63	2.24
Flan-T5-base	92.82	3.19	3.99
T5-base	92.48	3.94	3.59
SAFF (Flan-T5-base)	93.45	3.28	3.27
SAFF (T5-base)	93.00	3.65	3.35

Table A2: Generated graph temporal relation label distribution (in percentage).

reciprocal relation types. The first row is the model trained with *before*, *after*, *includes*, *is\_included*, and *simultaneous*. The second row is the model trained by merging *after* with *before*, *is\_included* and *includes* by swapping the head and tail events. Both models are trained with 4 augmented instances for each original instance. The results show the model benefits from the simpler label set.

Table A2 shows the relation type distribution generated by the models. The relation distributions are highly unbalanced.

Table A3 shows the average degree for the nodes in the generated graphs. The SAFF generate more complex graphs with higher node degree than the compared approaches.

We also investigate the effect of the order of the target edge set (Figure A4). The first row is Flan-T5-base fine-tuned with target edge set ordered based on random order in the documents. The second row is Flan-T5-base fine-tuned with target edge set ordered based on their appearance order in the documents. We could observe that the appearance order results in slightly better performance than random order. Each method was run on 5 different random seeds and trained for 50 epochs.

average node degree	
Vanilla Flan-T5-base	2.06
SAFF (w/o SPL)	2.16
SAFF	2.31

Table A3: The average node degree of the generated graphs on NYT.



	NYT-test			NYT-human		
	$P^E$	$R^E$	$F_1^E$	$P^E$	$R^E$	$F_1^E$
Random order	51.02	29.37	37.28	18.15	20.45	19.24
Appearance order	51.20	30.24	38.02	19.24	21.78	20.43

Table A4: Comparison based on different sequence orders on NYT-test and NYT-human.

## D Annotation of the Test Set

### D.1 Overview

We recruited crowd workers from Prolific<sup>4</sup> platform, which is a research-focused platform providing verified human workers. We recruited 24 participants in total (including pilot testing runs). In order to make sure the participants can understand and annotate the article efficiently, we require the participants to be native English speakers and have an education level higher than High school diploma/A-levels. We put 4 documents, which are randomly sampled from the same descriptor set as the training and testing of the selected NYT corpus, into each unit task. There is a shared document across all the tasks for the purpose of computing the inter-annotator agreement (IAA). In order to maximize the IAA, we asked 2 participants to first identify the event triggers in each unit task. After that, we merged the event annotations from the participants by taking the union of the spans (if there are overlapped spans, we take the longer span). Then, we asked another participant to annotate the event temporal relation based on the identified events. We also included the outputs from the CAEVO model to serve as examples, but we explicitly asked the participants to correct the annotations by adding, removing, or changing the CAEVO’s annotations. In the end, we collected 22 documents as the human-annotated test set.

On the event identification, we compute IOU (Intersection over Union) as a measure of agreement between the annotators. Average across 7 tasks, the IOU between the event spans is 0.8986. For the relation annotations, we compute the average Cohen’s  $\kappa$  of every participant pair in the relation annotation task (on the shared document). The average Cohen’s  $\kappa$  is 0.7465.

### D.2 Chosen Descriptors

Here are the chosen descriptors: “airlines and airplanes”, “olympic games”, “tennis”, “united states international relations”, “international relations”,

<sup>4</sup>prolific.com

“civil war and guerrilla warfare”, “track and field”, “soccer”, “bombs and explosives”, “politics and government”. We choose 2,000 documents from each descriptor. After preprocessing and filtering out some invalid documents, we have 18,263 documents in NYT-train, 1,000 documents in NYT-test, and 22 documents in NYT-human.

### D.3 Instructions and Interface

We use a popular open-sourced annotation interface called Doccano. As shown in Figure A1, annotators can select text spans for events. To direct annotators to distinguish events that actually occurred and imaginary events, we also provide an “imaginary event” label type. We asked them to annotate the predicates that are about a negative expression of an action or just a hypothesis in the context as an imaginary event. Imaginary events are orthogonal to the real-world timeline and thus have limited meaning for understanding the narrative.

Figure A2 shows the interface for annotating the relation. On this page, annotators can select two existing event spans, and then select the relation type from “before”, “includes”, and “simultaneous”.

Before the annotators came to the annotation platform, they went through a website where we put detailed descriptions and terminology definitions about the task. We also provided a video tutorial for using the annotation platform.

## E ChatGPT prompting

We used the OpenAI API chat completion model *gpt-3.5-turbo-0613*. We used the “function call” method to ensure better parsing quality. The function call parameters are shown in Figure A5. The temperature is set to 0. The other parameters are set as default. We show the inputs and outputs of the multi-hop prompting in Table A5.

## F GPT-4 Case Study

We show some test cases of where GPT-4 was prompted in this anonymous link<sup>5</sup>. The responses of GPT-4 essentially serve as summaries of the documents provided. The events it understood are quite broad, resembling abstracts of segments in the documents. This diverges from the NLP community’s definition of event understanding, which typically pertains to the occurrence of specific actions. We would like to obtain more granular information within the event temporal graph.

<sup>5</sup>Test cases on the TBD dataset

1105 Another issue with the graphs generated by GPT-  
1106 4 is that they tend to represent a linear sequence of  
1107 items ordered by their appearances in the document.  
1108 This ties back to the first issue, which relates to how  
1109 GPT-4 comprehends events. It essentially gener-  
1110 ated a summary of the document, which, while not  
1111 incorrect, does not align with the standard of event  
1112 temporal graph extraction as defined in SemEval  
1113 2013 TempEval-3 (UzZaman et al., 2013).

1114 Simply providing the definition of an event has  
1115 not brought about a change in its behaviour <sup>6</sup>.  
1116 While extensive prompt engineering might help,  
1117 we believe that incorporating some supervision sig-  
1118 nals could still be necessary. Our framework could  
1119 prove valuable for instruct-finetuning, aligning spe-  
1120 cific instructions with the event temporal graph  
1121 generation task.

---

<sup>6</sup>Prompt with definition

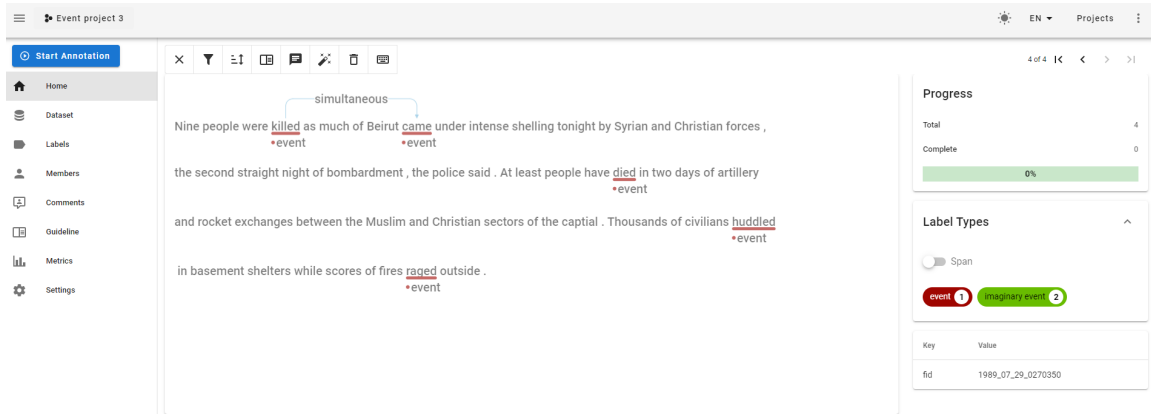


Figure A1: Annotation interface for event identification.

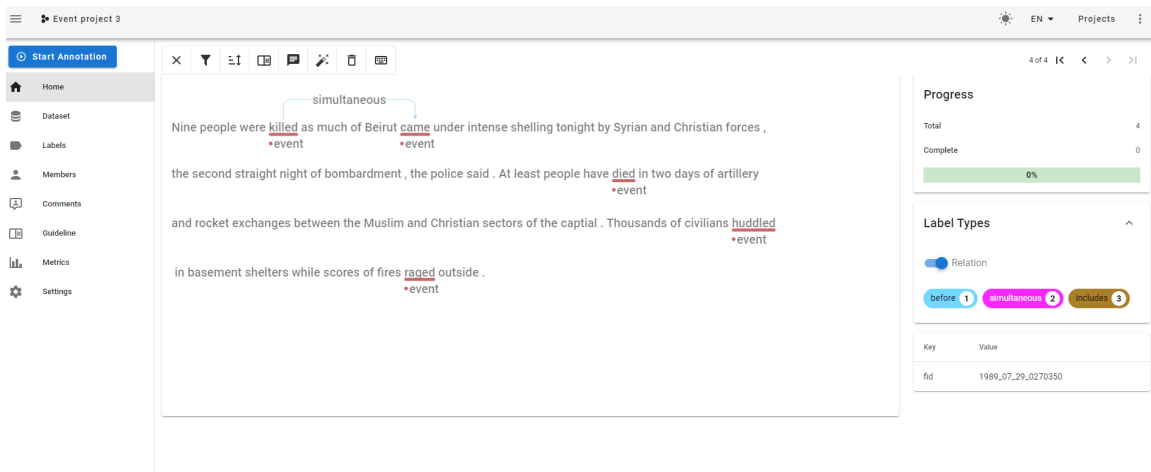


Figure A2: Annotation interface for event relation identification.

# Participant Information Sheet

## Invitation

You are invited to take part in a research study on event temporal relation understanding. The researchers in this project comes from [REDACTED]. Please take the time to read the following information carefully. Please ask us if there is anything that is not clear or if you would like more information.

## What will happen?

You will be asked to annotate a series of news articles. In each article, you will be asked to identify the temporal relation between the events. Specifically:

1. Given two event triggers, choose a temporal relation (definition provided below) type if the context indicates they are temporally related.

The interface of our annotation platform works best with a keyboard and mouse. Annotating with a touch screen can be very challenging.

## Details of Annotation

Here is a detailed [guideline](#) for the annotation including images and videos as examples.

## Time Commitment

You will be asked to annotate several news articles. The quantity are based on the length of each document. We estimate the task should take approximately an hour of your time.

## Confidentiality

We won't collect any of your personal data as part of this study. We will only collect your Prolific ID for the purpose of assigning you an anonymous account on our annotation platform. Once your responses have been submitted, if you wish to withdraw the data you have provided, you will be able to do so by emailing [REDACTED] and quoting your Prolific ID. We will only retain a record of your Prolific ID for 5 years after the end of this study (i.e., up until Sep 2028). After that time, it will no longer be possible to identify your individual responses.

## Participant rights

Your participation in this study is completely voluntary and choosing not to take part will not affect you or your rights in any way. You can also choose to withdraw your participation at any time, without giving a reason. You have the right to not answer or respond to any question that is asked of you. If you wish to withdraw from participating in the study, please close your browser.

## Benefits and risks

There are no known benefits or risks beyond everyday life for you in this study. The study will be useful in understanding what makes a conversation interesting to those involved in it.

Figure A3: Disclaimers.



# Guidelines

## Task Overview

You will be asked to annotate a series of news articles. In each article, you will be asked to identify the events mentioned in text and the temporal relation between them. Specifically:

1. Given two event triggers, choose a temporal relation (definition provided below) type if the context indicates they are temporally related.

The interface of our annotation platform works best with a keyboard and mouse. Annotating with a touch screen can be very challenging.

## Definition of the Terminologies

1. Event: An event is an action or a change of states. Events in text are notified by event triggers, which is the word that indicates the occurrence of the event. A trigger is usually a predicate, but it can also be a noun phrase. **In this task, the event triggers have been annotated.**
2. Temporal Relation: We are interested in the relative temporal relations between events. We have the following label types: *before*, *includes*, and *simultaneous*.
  - o If event A ends before event B starts, then event A happened *before* event B.
  - o If event A happened within the span of event B, then event B *includes* event A.
  - o If event A and event B clearly don't include one or the other, and the time spans of them mostly overlap, then choose *simultaneous*.

We provide some examples below, please check them carefully before the task. Note that you don't have to choose a relation for every event pair. **Only when the context clearly indicates a temporal relation**, then choose a relation label. Otherwise, please skip them.

## Annotation Tutorial

Please ignore the part about annotating the event triggers. In this task, you are only responsible for relations.

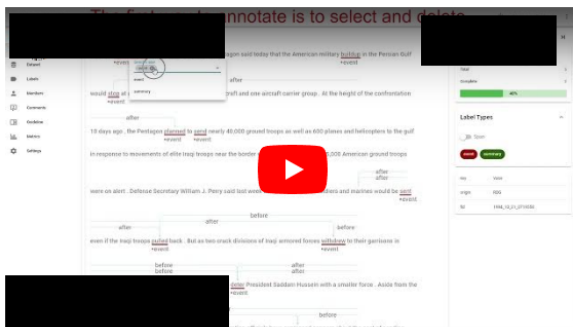


Figure A4: Guides.

```

1 FUNCTION_LIST = [
2   {
3     "name": "save_events",
4     "description": "Store the extracted events in a list",
5     "parameters": {
6       "type": "object",
7       "properties": {
8         "event_list": {
9           "type": "string",
10          "description": "This is a list of event strings",
11        }
12      },
13      "required": ["event_list"],
14    },
15  },
16  {
17    "name": "save_graph",
18    "description": "Store the constructed graph in DOT language",
19    "parameters": {
20      "type": "object",
21      "properties": {
22        "graph": {
23          "type": "string",
24          "description": "The constructed graph in DOT language. \
25          This graph is a strict graph, in which every edge containing \
26          two event nodes, and a temporal relation label from \
27          [\"before\", \"includes\", \"simultaneous\"]. For example, \
28          \"strict graph {\n\"The Organization asserted responsibility \
29          \" -- \"a United States Navy diver killed\" [rel=before];\n}\"",
30        }
31      },
32      "required": ["graph"],
33    },
34  },
35 ]

```

Figure A5: Function list for OpenAI API call

Input	Output
<p>First hop:  Instructions: Analyse the given document. Extract the main events. Store the events as a JSON list by calling the <i>save_events_function</i>.  Example input document: Governor Cuomo leaving his office in Albany yesterday to speak to reporters after he met with representatives of group pushing for higher ethical standards for public officials in the state.  Example output events: Governor Cuomo leaving his office, Governor Cuomo speak to reporters, Governor Cuomo met with representatives."  Document: President Clinton says he is committed to a possible strike against Iraq, despite objections from other countries. Steve Hurst has that. The US military buildup in the Persian gulf continues apace, more planes headed from the united states and senior officials say Iraq's president Saddam Hussein can expect punishing air strikes to continue well into the future if he doesn't stop building biological and chemical weapons. Saddam will know by our actions and our warning that we will be prepared to act again if we have evidence he is trying to rebuild his weapons of mass destruction capabilities. In a bit of television diplomacy, Iraq's deputy foreign minister responded from Baghdad in less than one hour, saying Washington would break international law by attacking without UN approval. The united states is not authorized to use force before going to the council. President Clinton, meantime, glossed over stern warnings from Moscow on Thursday that US air strikes against Iraq could do serious harm to relations with the Kremlin. We're trying to find a diplomatic solution. And I hope that, whatever happens today, that our relationships with Russia will continue to be productive and constructive and strong, because that's very important to the future of our peoples. One contrary view of the issue presented itself to the president as he arrived in Philadelphia later in the day. Nevertheless, the president said Washington would use force if diplomacy fails to force Saddam Hussein to back down.</p>	<pre>strict graph { "President Clinton's commitment to a possible strike against Iraq" -&gt; "The US military buildup in the Persian Gulf" [label="before"]; "The US military buildup in the Persian Gulf" -&gt; "Senior officials' statement about punishing air strikes against Iraq" [label="before"]; "Senior officials' statement about punishing air strikes against Iraq" -&gt; "Iraq's deputy foreign minister's response from Baghdad" [label="before"]; "President Clinton's attempt to find a diplomatic solution" -&gt; "Contrary views expressed in Philadelphia" [label="before"]; "Contrary views expressed in Philadelphia" -&gt; "President Clinton's statement about using force if diplomacy fails" [label="before"]; "The Russian foreign minister's suggestion to wait for UN Secretary General Kofi Annan's visit" -&gt; "Secretary of State Madeleine Albright, Berger, and Defense Secretary William Cohen's plans to travel and explain military force necessity" [label="before"]; }</pre>

The Russian foreign minister, meanwhile, sought to soften the harsh words of his military counterpart, saying on Friday that Russia now feels the US must hold off at least until UN secretary general Kofi Annan visits Baghdad in a last-ditch effort at diplomacy. Annan has no trip planned so far. Meanwhile, secretary of state Madeleine Albright, Berger and defense secretary William Cohen announced plans to travel to an unnamed city in the US heartland next week, to explain to the American people just why military force will be necessary if diplomacy fails. Steve Hurst, CNN, the state department.

Second hop:

For every possible pair of events, please infer the relative temporal order (e.g., before, includes, simultaneous) between them based on the context. Make a graph where the nodes are the events, and the edges represent the temporal relation between the events. Lastly, please output the graph using DOT graph description language.

Conversation link:

<https://chat.openai.com/share/a7c8659d-4fbd-4cfa-b60c-b4066ac77ee9>

---



<p>First hop:</p> <p>Instructions: Analyse the given document. Extract the main events. Store the events as a JSON list by calling the save_events function.</p> <p>Document: For The World, this is Pamela McCall in the BBC news room in London. More heavy gunfire in the Serbian province of Kosovo. It's the second day of an offensive by the security forces against villages populated by ethnic Albanians. BBC correspondent Karyn Coleman reports from Kosovo. There were more accounts of people fleeing from the villages near where the offensive was taking place. They were too afraid to stay, fearing the forces may also move their operations further afield. By mid afternoon, official Serb sources were saying the operation was over, but that has not yet been confirmed from Belgrade, the capital of Serbia, which is where the whole attack is thought to have been planned. The ethnic Albanians here in Kosovo are considerably worried about the future and fear they are witnessing the beginnings of a much wider conflict. Karyn Coleman, BBC news, Pristina.</p> <p>Second hop:</p> <p>For every possible pair of events, please infer the relative temporal order (e.g., before, includes, simultaneous) between them based on the context. Make a graph where the nodes are the events, and the edges represent the temporal relation between the events. Lastly, please output the graph using DOT graph description language.</p> <p>Conversation link:  <a href="https://chat.openai.com/share/6bdeb237-247d-486f-91ba-146c3b4fd3b2">https://chat.openai.com/share/6bdeb237-247d-486f-91ba-146c3b4fd3b2</a></p>	<pre>strict graph {   "Heavy gunfire in the Serbian province of Kosovo." -- "The second day of an offensive by the security forces against villages populated by ethnic Albanians" [rel=before];   "The second day of an offensive by the security forces against villages populated by ethnic Albanians" -- "People fleeing from villages near the offensive" [rel=before];   "The second day of an offensive by the security forces against villages populated by ethnic Albanians" -- "Official Serb sources stating the operation was over" [rel=before];   "Official Serb sources stating the operation was over" -- "Uncertainty regarding the operation's status in Belgrade" [rel=simultaneous];   "Uncertainty regarding the operation's status in Belgrade" -- "Ethnic Albanians in Kosovo expressing worry about the future and fear of a wider conflict" [rel=before]; }</pre>
--	--

Table A5: Example of ChatGPT promoting on the TBD data.