Can LLMs Help Encoder Models Maintain Both High Accuracy and Consistency in Temporal Relation Classification?

Adiel Meir and Kfir Bar

Efi Arazi School of Computer Science, Reichman University, Herzliya, Israel matufadiel@gmail.com, kfir.bar@runi.ac.il

Abstract

Temporal relation classification (TRC) demands both accuracy and temporal consistency in event timeline extraction. Encoderbased models achieve high accuracy but introduce inconsistencies because they rely on pairwise classification, while LLMs leverage global context to generate temporal graphs, improving consistency at the cost of accuracy. We assess LLM prompting strategies for TRC and their effectiveness in assisting encoder models with cycle resolution. Results show that while LLMs improve consistency, they struggle with accuracy and do not outperform a simple confidence-based cycle resolution approach. Our code is publicly available at: https://github.com/MatufA/ timeline-extraction.

1 Introduction

Extracting event timelines from text is a key natural language processing (NLP) task, organizing events chronologically based on their relative occurrence rather than absolute timestamps. A broader definition by (Ocal et al., 2024) describes a timeline as a data structure that arranges events and times in a total order. Timelines have a wide range of practical applications, even when considering events alone. For instance, Bakker et al. (2024) demonstrated how timelines can be used to process government decision letters, extracting and organizing events for improved understanding. Another example is in the medical domain (Sezgin et al., 2023): given a patient's textual medical record—or a collection of such records—it becomes valuable to extract a timeline of relevant medical events to summarize and visualize their journey. Timeline extraction typically involves five steps: (1) event detection, identifying relevant events, often treating all verbs as events; (2) anchoring, selecting events for comparison; (3) temporal relation classification (TRC), assigning relations to pairs; (4) graph construc*tion*, combining pairwise relations into a temporal graph; and (5) *timeline extraction*, deriving a timeline from the graph.

Various methods have been proposed for extracting timelines from temporal graphs (Mani et al., 2006; Do et al., 2012; Kolomiyets et al., 2012; Xue and Zhang, 2018). Recently, Ocal et al. (2024) proposed a method for extracting event timelines from documents annotated with the full TimeML scheme (Saurí et al., 2006), which defines 13 temporal relation types. However, modeling all 13 relations is complex and often results in temporal inconsistencies.

To address the complexity of TimeML's full relation set, several datasets focus on simplified subsets. A widely used resource for TRC is MATRES (Ning et al., 2018b), which reduces the relation types to three deterministic labels—before, after, and equals—along with a vague category for uncertain cases. These labels are assigned to a subset of all possible event pairs, a design choice intended to improve annotation consistency and reduce ambiguity.

Despite this simplification, temporal inconsistencies can still arise, particularly with models following a *pairwise* approach: predicting relations independently for each event pair without considering previously predicted labels. For example, a model might predict: A *before* B, B *before* C, and mistakenly, A *after* C. The last relation contradicts the others and creates a temporal cycle, which complicates efforts to derive a consistent, linear event timeline. A real instance of such a cycle, predicted on a MATRES document, is illustrated in Figure 1.

Large language models (LLMs) have achieved state-of-the-art performance across many NLP tasks. However, previous studies (Roccabruna et al., 2024) have shown that generative LLMs underperform compared to encoder-based models on the TRC task as defined in MATRES. The advantage of LLMs lies in their ability to encode

document-wide information flexibly, which enables them to generate an entire temporal graph in a single step. This capability, recently termed *global* TRC, offers the potential to reduce temporal inconsistencies by considering all event pairs jointly. Building on prior work in TRC, non-fine-tuned generative LLMs still lag behind smaller supervised models that follow the pairwise approach. However, LLMs' ability to generate the entire temporal graph in a single inference step offers a key advantage: the potential to reduce temporal inconsistencies, a common issue in pairwise models.

Therefore, in this work we make two main contributions:

- We study the performance of generative LLMs in extracting temporal graphs. Specifically, we focus on the trade-off between pairwise classification accuracy and the rate of temporal inconsistencies (e.g., cycles) in the resulting graph. Using the MATRES dataset, we explore different approaches to prompt design under various input and output conditions.
- Additionally, we propose a hybrid approach that combines a generative LLM with a standard supervised encoder to improve accuracy while mitigating cycles in the temporal graph.

2 Related work

Temporal relation classification has primarily been addressed using fine-tuned, relatively small encoder-based language models, typically following a pairwise approach in which each event pair is labeled independently.

A key limitation of the pairwise approach is its tendency to produce globally inconsistent outputs. Since these models make independent predictions for each pair of events, they do not take previously predicted labels into account during inference. This lack of global awareness can result in contradictions, such as temporal cycles, which undermine the coherence of the predicted temporal structure and ultimately hinder accurate timeline construction. Despite this limitation, numerous well-established encoder-based methods have been proposed to tackle pairwise TRC. These include approaches that leverage contextualized representations and joint inference strategies to improve local and global consistency (Han et al., 2021; Zhou et al., 2021; Ning et al., 2019; Mathur et al., 2021; Wang et al., 2022, 2023; Zhang et al., 2022; Zhou

et al., 2022; Man et al., 2022; Cohen and Bar, 2023; Niu et al., 2024). While these models have contributed significantly to the field, the challenge of maintaining globally coherent temporal graphs remains a central concern in temporal relation classification.

Early efforts such as Ning et al. (2019) introduced a structured framework for TRC by refining the task with better contextual representations and curated evaluation protocols. Subsequent work expanded this by incorporating global constraints, as in Mathur et al. (2021), which applied joint inference to enforce temporal consistency across event graphs. Similarly, Han et al. (2021) proposed EcoNet, which leveraged event graph structures and global coherence to improve document-level temporal reasoning.

Domain-specific applications have also driven innovation in TRC, particularly in the clinical domain. Zhou et al. (2021) addressed the challenges of TRC in clinical texts, which often involve fragmented or incomplete narratives. Their work demonstrated that specialized models and annotation schemes are necessary to adapt general TRC methods to the clinical setting. (Cohen and Bar, 2023), reframed TRC as a Boolean question answering task. By training a RoBERTa model on Yes/No questions formulated based on the annotation guidelines, they effectively simulated the human annotation process and achieved state-of-the-art results on the MATRES dataset. More recent work by Niu et al. (2024) introduced ContEMPO, a large-scale benchmark for document-level temporal reasoning, combining distant supervision and LLM-based annotation to enhance the breadth and realism of training data.

Several recent studies have also focused on extracting temporal structures beyond pairwise relations. Wang et al. (2022) and Wang et al. (2023) explored the prediction of document creation times (DCT) and global temporal graphs, respectively, highlighting the importance of temporal anchoring in narrative understanding. Zhang et al. (2022) and Zhou et al. (2022) also tackled full timeline construction, proposing models that jointly identify events and infer their temporal relationships, often integrating external knowledge or reasoning modules.

With LLMs becoming state-of-the-art in many tasks and offering more flexible input handling in a zero-shot setting, recent studies have explored different ways to use them for TRC, both in pairwise and global settings.

Barack Obama would make...Traditionally, the (intentionally) funny lines by our presidents have had one thing in common: They were self-deprecating. Sure, some presidents have [EVENT5]used[/EVENT5] jokes to take jabs at their opponents, but not to the extent of Obama. During his tenure, he has increasingly [EVENT8]unleashed[/EVENT8] biting comedic barbs against his critics and political adversaries. These jokes are [EVENT1000]intended[/EVENT1000] to do more than simply entertain you. They have an agenda. Obama's humor is often delivered the way a comedian dealing with a heckler would do it. He tries to undermine his opponents with it and get the crowd -- in this case the public -- on his side. I can [EVENT20]assure[/EVENT20] you that having a crowd laugh at your critic/heckler is not only effective in dominating them, it's also very satisfying.



Figure 1: Example of a cycle in a document from the MATRES dataset, mistakenly generated by one of our pairwise encoder models.

Jain et al. (2023) evaluated a variety of LLMs (including standard and code-generation models) across different temporal tasks and prompting strategies (zero-shot, few-shot). Their comprehensive analysis revealed that while LLMs exhibit proficiency in certain temporal aspects, they face significant challenges in areas requiring reasoning over specific timings and handling complex scenarios involving multiple events. Focusing specifically on the pairwise TRC task, (Roccabruna et al., 2024) investigated if LLMs could supersede established encoder-only models. Evaluating several LLMs with in-context learning and fine-tuning, they found that LLMs generally underperform a strong RoBERTa baseline for this task. Through explainability methods and analysis of word embeddings, they attributed this gap, in part, to differences in pre-training objectives (autoregressive vs. masked language modeling) and how models process input sequences. These studies highlight that while LLMs show promise for broader temporal reasoning, the specific requirements of tasks like pairwise temporal classification may still favor specialized encoder-only architectures or necessitate further research into tailoring LLMs for such fine-grained analysis.

Recent studies have begun exploring the use of zero-shot LLMs for TRC, though most efforts have adhered to the traditional pairwise prediction framework (Yuan et al., 2023; Li et al., 2024; Kougia et al., 2024). A more recent study (Eirew et al., 2025) proposes enhancing global consistency by prompting a strong generative LLM to produce the

entire graph of temporal relations in a single step. To address potential contradictions and instability of LLMs in generating consistent output, this approach incorporates a post-processing step based on the linear programming optimization framework introduced by Ning et al. (2018a), which enforces global coherence by resolving inconsistencies in the predicted temporal graph.

Together, these works form a comprehensive foundation for understanding the evolution of TRC, from pairwise classification to global timeline construction, and from specialized supervised models to LLM-based generalization. Building on these efforts, we compare zero-shot LLM accuracy and consistency across prompts and propose a simple, effective cycle-breaking method for encoders while maintaining accuracy.

3 Datasets

Our investigation of the trade-off between pairwise accuracy and global consistency is grounded in experiments on two datasets that represent distinct annotation paradigms.

3.1 MATRES (Ning et al., 2018b)

MATRES is a widely-used benchmark for TRC which simplifies TimeML's relations into four labels: *before*, *after*, *equal*, and *vague*. It employs a *sparse annotation* strategy, providing gold labels primarily for event pairs in close proximity (i.e., within two-sentence contexts) to enhance interannotation agreement. This sparsity directly impacts our evaluation: accuracy is computed using

only the gold subset, while temporal consistency is measured over all generated relations.

3.2 NarrativeTime (NT) (Rogers et al., 2024)

NarrativeTime offers a contrasting, *dense annotation* approach by labeling all event pairs within a document. NT expands the label set to seven types, including those from MATRES plus *includes*, *is_included*, and *overlap*. We leverage its comprehensive coverage in Section 5 to evaluate our cycle-breaking approach.

4 Evaluation of LLMs on TRC

4.1 Extraction Approach

Our timeline extraction approach follows the fivestep process outlined in Section 1. Specifically, we work with the MATRES dataset, where all events are defined as verbs. MATRES employs a novel strategy for determining which events should be anchored to a given event. Building on the approximate complete-graph approach introduced in (Naik et al., 2019)—where events are anchored only to those within a predefined surrounding window of sentences-MATRES further refines this by incorporating different types of narrative axes (e.g., opinions, intentions), which impact anchoring decisions. In our work, we build on the MA-TRES anchoring framework and ask the LLM to merely classify the anchored event pairs according to the MATRES label set: before, after, equal, and vague. We explore various approaches to modeling input context length, event marking, yield type, and prompt techniques. Broadly speaking, for a given full document i with k marked events, we use an LLM as a function to predict the corresponding temporal graph. The prompt is structured into three sections: 1) instructions (s_i) ; 2) the input text $(t_i(e_1, e_2, \dots, e_k))$, including kmarked events to be classified (we mark events as [EVENT1]eat[/EVENT1], with the event number taken from the dataset.); and 3) some illustrative input-output examples (f). The output is composed of one or more (m) labels l_i , with each label corresponding to an event pair introduced in the input. Formally, we use the generative LLM as follows:

$$l_1, l_2, \dots, l_m = LLM[s_i, t_i(e_1, e_2, \dots, e_k), f]$$

The LLM is expected to return a single label for each event pair formed from the marked events. Following the self-consistency approach Wang et al. (2023), we run each instance five times and use majority voting across runs. Only event pairs where a single label receives majority vote are included as links in the temporal graph. If no clear majority exists, we treat the relation as unknown and exclude it from the graph. This subset approach ensures greater reliability in the predicted temporal structure.

We observe that generative LLMs tend to predict the *vague* label more often, likely reflecting their uncertainty. To address this, since LLMs are not instructed to label every pair and can choose which pairs to label, sometimes we remove *vague* from the label set given to the models directing them to predict only *before*, *after*, or *equals*. Furthermore, the *equals* label poses an additional challenge for handling in a timeline and is both infrequently annotated and predicted. As a result, we choose to ignore *equals* and *vague* when constructing a temporal graph. Consequently, only the *before* and *after* labels are used as links to form the temporal graph.

We explore variations in prompt design, particularly focusing on the following aspects:

Output Type. Most prior work predicts temporal labels for event pairs individually, a straightforward but inconsistency-prone approach due to its lack of global context. An alternative is predicting the entire temporal graph in one step, leveraging global context for better consistency. We evaluate both approaches—pairwise and graph. In the graph approach, the model generates labels for all pairs in DOT format (Gansner et al., 2006).

Considered Events. MATRES was selectively annotated, labeling only event pairs within two-sentence paragraphs, leaving many pairs unannotated. To address this, we evaluate accuracy using three approaches. The MATRES approach considers only the originally annotated pairs, using the *pairwise* output type. The *sliding-window* approach expands this by pairing each event with all others in a two-sentence window, shifting one sentence at a time. The *document* approach, applicable only to *graph* output, considers all event pairs but avoids redundancy by marking events and instructing the model to infer non-redundant relations, omitting symmetric and transitive ones to produce a compact graph. In both the *sliding-window*

¹Released under the CC-BY 4.0 license (Ning et al., 2018b), we use the dataset for evaluation as intended by its authors.

and *document* approaches, event pairs without gold labels but assigned a relation by the model are included for consistency assessment but excluded from accuracy calculations.

Context. The context refers to the portion of text surrounding the events that are provided to the model for classification. We experiment with two context sizes. In the first, referred to as *document*, we provide the entire document to the model. In the second, called *paragraph*, we provide a window of two sentences surrounding the two events in focus. This method is compatible only with the *pairwise* output type.

Prompt Style. We experiment with both zeroshot and few-shot in-context learning. For the *pairwise* output, few-shot learning includes two examples—one *before* and one *after*—randomly selected from the training set. For the *graph* output, we provide a document with marked events and the MATRES-annotated relations in DOT format. Note that this does not fully represent a complete graph, since MATRES provides gold relations only for some of the event pairs.

Sample prompts are provided in Appendix A.

4.2 Evaluation Approach

We experiment with all combinations of the prompt aspects mentioned above, using four LLMs: GPT4o (OpenAI et al., 2024), GPT4o-mini (OpenAI et al., 2024), LLama-3.1-8B (Grattafiori et al., 2024), and LLama-3.2-3B (Grattafiori et al., 2024). We choose these specific models to balance between large and small models, as well as between open and closed weights. Note that not all aspect combinations are possible. For instance, in the graph output type, only four combinations exist because both the considered events and context aspects must be set to document. Additionally, we evaluate the graph configuration only with OpenAI GPT models, as the task demands stronger reasoning capabilities. Our evaluation balances inconsistency and accuracy: inconsistency is measured by the percentage of test documents with cycles, which prevent timeline extraction, while accuracy is reported as the Micro-F1 score. We evaluate on the MATRES test set (20 documents) and refer to the percentage of cycle-containing documents as the cycle rate.

4.3 Results

Figure 2 offers a high-level overview of our results, highlighting both accuracy and consistency metrics across the different experimental settings. For a comprehensive breakdown, Table 1 presents the full set of results, covering all prompt configurations evaluated with the four large language models used in this study. This includes performance across both pairwise and global prediction modes, as well as the impact of prompt strategies.

While the results are somewhat noisy, we observe a strong correlation between accuracy and cycles ($\rho=0.64,\,p<.001,\,n=36$). This indicates that higher accuracy is often accompanied by reduced temporal consistency: conservative models tend to lower recall and therefore reduce the number of cycles, whereas models that output more relations are more likely to introduce cycles. However, this correlation does not preclude the existence of configurations that simultaneously improve both accuracy and consistency. Indeed, recent work by Eirew et al. (2025) demonstrated that such improvements are achievable under specific settings. We discuss this in more detail below.

Predicting *vague* is particularly challenging, as human annotators also struggled with it (Ning et al., 2018b), and often represents disagreement between annotators. Therefore, for some experiments, we tested the model with and without the *vague* label. When included, the model predicted one of four labels (*before*, *after*, *equal*, or *vague*); otherwise, it was limited to three.

Accuracy calculations were based on the fourlabel or three-label setting, respectively, following our definition of the F1 score described above. However, for calculating consistency, defined as the number of documents that introduce cycles, we only consider the before and after relations. We exclude equal, given its rare occurrence (less than 5% of the relations in the MATRES test set), and vague, since it does not participate in cycles. In Figure 2, we indicate experiments that include the vague relation by placing a line under each relevant shape. In total, there are seven experiments for which results exist both with and without vague. Averaging the differences between corresponding experiments, we observe an 18% drop in accuracy when allowing vague, but also a 37% reduction in cycled documents. This suggests that when the model can predict vague, it does so frequently, which lowers accuracy but also helps break cycles, as only before

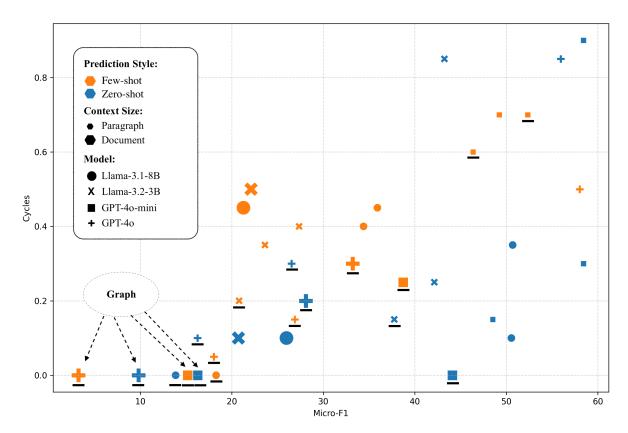


Figure 2: Micro-F1 vs. inconsistency (cycle rate in the test set) across LLM experiments. Underlined markers denote configurations including the vague relation. Each marker's shape and color encode model type, prompt style, and context as explained in the text and Table 1.

and after relations contribute to cycle formation.

The experiments using the *graph* output type are represented by the bubble labeled *graph*. Their accuracy is notably low, suggesting that LLMs struggle to generate the full temporal graph in a single pass using this relatively simple approach supported by previous work (Eirew et al., 2025).

In the *pairwise* (non-*graph*) experiments, accuracy improves but introduces more cycles. Another insight is that larger contexts limit accuracy, while smaller ones increase it but add cycles. Since our goal is timeline generation, we prioritize fewer cycles with reasonable accuracy. The best experiment according to this criterion is produced by GPT-4omini in a few-shot configuration combined with the *paragraph* context.

4.4 Accuracy vs. Consistency

As reported above, we observe a statistically significant correlation between inconsistency—measured as the number of documents predicted with temporal cycles—and accuracy. In general, our experiments suggest that for non-fine-tuned LLMs, such as those evaluated in this study, higher accuracy

often comes at the cost of lower consistency. That is, as the model generates more accurate temporal labels, it also tends to introduce a greater number of contradictory relations. This trade-off appears especially when prompting the model to be more successful (e.g., by providing a more relevant context, or by providing examples) in its predictions, which intuitively aligns with the classic tension between specificity and sensitivity observed in traditional machine learning tasks.

However, it is important to note a fact in our evaluation strategy: consistency is measured over all predicted event-event relations, whereas accuracy is computed with respect to the subset of annotated pairs in the MATRES dataset. Since MATRES does not provide gold labels for event pairs that are more than one sentence apart, but the LLMs output relations for all event pairs, our consistency metric is based on a broader set of predictions than the accuracy metric. This introduces a slight misalignment between the two evaluation dimensions, but we choose to retain this approach to more comprehensively capture the model's global behavior.

While our overall results support the observed

Model	Output Type	Prompt Style	Considered Event	Context	Vague	Micro-F1	Cycles
Llama-3.1-8B	Pairwise	zero-shot	MATRES	Paragraph	no	50.55	0.10
			Sliding-window	Paragraph	no	50.69	0.35
				Paragraph	yes	13.86	0.00
				Document	no	25.97	0.10
		few-shot	MATRES	Paragraph	no	35.91	0.45
			Sliding-window	Paragraph	no	34.39	0.40
				Paragraph	yes	18.28	0.00
				Document	no	21.27	0.45
Llama-3.2-3B	Pairwise	zero-shot	MATRES	Paragraph	no	42.13	0.25
			Sliding-window	Paragraph	no	43.23	0.85
				Paragraph	yes	37.75	0.15
				Document	no	20.72	0.10
		few-shot	MATRES	Paragraph	no	23.62	0.35
			Sliding-window	Paragraph	no	27.35	0.40
				Paragraph	yes	20.79	0.20
				Document	no	22.1	0.50
	Pairwise	zero-shot	MATRES	Paragraph	no	58.43	0.30
GPT-4o-mini			Sliding-window	Paragraph	no	58.43	0.90
				Paragraph	no	48.51	0.15
				Document	yes	44.09	0.00
		few-shot	MATRES	Paragraph	yes	52.33	0.70
			Sliding-window	Paragraph	no	49.22	0.70
				Paragraph	yes	46.36	0.60
				Document	yes	38.71	0.25
	Graph	zero-shot	Document	Document	yes	16.25	0.00
	Grapii	few-shot	Document	Document	yes	15.17	0.00
GPT-4o	Pairwise	zero-shot	MATRES	Paragraph	yes	16.25	0.10
			Sliding-window	Paragraph	yes	26.52	0.30
				Paragraph	no	55.94	0.85
				Document	yes	28.08	0.20
			MATRES	Paragraph	yes	18.04	0.50
		few-shot	Sliding-window	Paragraph	yes	26.88	0.15
				Paragraph	no	58.01	0.50
				Document	yes	33.21	0.30
	Graph	zero-shot	Document	Document	yes	9.80	0.00
		few-shot	Document	Document	yes	3.30	0.00

Table 1: Full results of the evaluation of LLMs under different prompting conditions.

trend—that increasing accuracy tends to introduce more inconsistencies—we also find notable exceptions. Certain model—prompt configurations, specifically GPT-40, GPT-40-mini, and LLaMA-3-8B with the Paragraph-context prompting style, show improvements in both accuracy and consistency. These cases suggest that, although the tradeoff is common, it is not inevitable. Prior work has shown that with careful modeling, such as the

use of structured inference or post-hoc consistency enforcement, systems (Eirew et al., 2025) can improve both dimensions simultaneously. Nonetheless, our findings are specific to zero- and few-shot prompting approaches using non-fine-tuned LLMs, and future work may further explore how fine-tuning or additional consistency-aware methods can shift or mitigate this trade-off.

Document	Number of Cycles
CNN_20130321_821	4
CNN_20130322_1003	135
WSJ_20130321_1145	36
WSJ_20130322_159	72
WSJ_20130322_804	59
nyt_20130321_china_pollution	202
nyt_20130321_cyprus	87
nyt_20130321_sarcozy	20
nyt_20130321_women_senate	74

Table 2: Number of simple cycles per each document of the 9 documents that have cycles from the MATRES test set.

Model	MATRES F1	NT F1
Confidence-Based	74.67	52.28
LLM-Assisted (GPT-4o)	67.03	51.90
LLM-Assisted (GPT-4o-mini)	70.49	51.65

Table 3: Performance of cycle-breaking approaches. Since MATRES is only sparsely annotated, it is hard to know the upper bound for the performance gain only by removing relations to break cycles. For NT, the best reported SOTA performance is 52.57.

5 Combining LLMs with a Small Encoder

It has already been shown (Roccabruna et al., 2024) that encoder models achieve higher accuracy than generative LLMs on TRC. However, as demonstrated in the previous section, LLMs maintain greater consistency than simple encoders when leveraging global information. Building on this insight, we adopt a hybrid approach, in which a baseline encoder model first predicts temporal relations for all event pairs in a document using the pairwise approach. We then detect simple cycles in the predictions using the NetworkX package (https://networkx.org/). A simple cycle is a cycle in a graph with no repeated nodes. As mentioned before, only before and after relations are considered for cycle detection. Once a cycle is found, we iteratively break it using one of two methods, and the process repeats until no cycles remain in the document. Generally speaking, breaking a simple cycle involves removing a single link from it, aiming to minimize accuracy loss while restoring consistency. We explore two different approaches: Confidence-Based. This method removes the cycle link with the lowest confidence, determined by the encoder's probability for the predicted label. LLM-Assisted. This approach prompts a generative LLM to identify the most likely erroneous link in a cycle, leveraging its ability to process detailed

input and enhance global consistency. The prompt (Appendix A) provides TRC instructions, requiring the model to identify the most likely error in a document with a cycle of *before* and *after* links. It then presents the full document with marked events and cycle links in DOT format (Gansner et al., 2006).

5.1 Experimental Settings and Results

In addition to MATRES, we use the NarrativeTime (NT) dataset (Rogers et al., 2024) (MIT license), which labels all event pairs in a document rather than just within two-sentence segments. NT also uses seven relations, which are the four of MA-TRES plus three more includes, is_included, and overlap. The NT test set contains 9 documents (overall, 7,582 relations), 27 training documents (overall, 67,860 relations). We evaluate the two cycle-breaking approaches using an encoder model trained from scratch once on the MATRES training set and once on the NT training set. The model follows the BERT-based (Devlin et al., 2019) (License: Apache 2.0) Entity Marker Entity Start architecture (Baldini Soares et al., 2019), where event mentions are marked with special tokens [E1] and [/E1] for the first event and [E2] and [/E2] for the second event. This architecture operates pairwise, independently classifying each event pair without considering previously predicted labels. On the MATRES

test set, our encoder achieves a micro-F1 score of 80.29%, and on the NT test set, 52.57%. Additionally, 9 of 20 MATRES test documents contain cycles, averaging 76.5 simple cycles per document, while all 9 NT test sets include cycles. Table 2 breaks down the number of simple cycles detected in each document from the MATRES test set. The cycles were detected over the full temporal graph extracted by the base supervised BERT model.

For the LLM-assisted cycle-breaking approach, we experiment with GPT-40 and GPT-40-mini. All three cycle-breaking approaches successfully resolved all cycles, but accuracy dropped from the original 80.29% (MATRES) and 52.57% (NT). Table 3 summarizes the results. The confidence-based approach significantly outperformed LLM-assisted methods on MATRES and showed less conclusive results on NT, suggesting it was more effective at identifying the correct links to remove.

6 Conclusions

Our study highlights both the promise and the current limitations of using LLMs for timeline extraction through TRC. Across extensive experiments, we observe a recurring inverse relationship between accuracy and consistency: as LLMs are pushed toward higher accuracy through prompt design and reasoning strategies, they tend to generate more globally inconsistent temporal graphs—often resulting in cycles. This trade-off mirrors classic precision-recall tensions in traditional machine learning and highlights the challenge of achieving both local accuracy and global coherence in zeroor few-shot generative settings. We also find that current LLMs struggle to generate complete and accurate temporal graphs in a single pass, even when using compact representations and chainof-thought reasoning. While supervised encoderbased models remain more accurate on annotated pairs, their pairwise prediction structure inherently introduces global inconsistencies unless followed by post-hoc constraints. Interestingly, when evaluating LLM-generated graphs with existing cycle resolution strategies, a simple confidence-based encoder model remained among the most effective for enforcing consistency—highlighting the value of integrating structured reasoning modules into otherwise generative workflows. Our results motivate future research directions focused on hybrid approaches that combine the strengths of encoderbased models—particularly their structured reasoning and consistency enforcement—with the generalization and contextual understanding of LLMs. We believe that with targeted improvements in prompt engineering, structural guidance, and consistency-aware inference, LLMs can play a central role in advancing temporal relation extraction beyond current limitations.

Limitations

Our study has several limitations. First, our approach Confidence-Based Cycle Breaking relies on confidence scores derived from BERT-based architecture, which may limit the generalization of our conclusions to other architectures or classification strategies. MATRES and NarrativeTime both annotate news articles, so our conclusions may not generalize to other domains. Finally, our experimental evaluation was performed on only four models (two open source and two closed), despite the existence of a broader array of models in the literature.

We see no risks in our work, as we use publicly available datasets as intended and employ LLMs like GPT-40 solely for evaluation. We did not review or filter the datasets for personal information, as both datasets consist solely of publicly available news documents sourced from media outlets.

References

Femke Bakker, Ruben Van Heusden, and Maarten Marx. 2024. Timeline extraction from decision letters using chatgpt. In *Proceedings of the 7th Workshop on Challenges and Applications of Automated Extraction of Socio-political Events from Text (CASE 2024)*, pages 24–31.

Omer Cohen and Kfir Bar. 2023. Temporal relation classification using Boolean question answering. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 1843–1852, Toronto, Canada. Association for Computational Linguistics.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

Quang Do, Wei Lu, and Dan Roth. 2012. Joint inference for event timeline construction. In *Proceedings of* the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning, pages 677–687.

Alon Eirew, Kfir Bar, and Ido Dagan. 2025. Beyond pairwise: Global zero-shot temporal graph generation. *Preprint*, arXiv:2502.11114.

Emden Gansner, Eleftherios Koutsofios, and Stephen North. 2006. Drawing graphs with dot.

Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, Danny Wyatt, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Francisco Guzmán, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Govind Thattai, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jack Zhang, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Karthik Prasad, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Kushal Lakhotia, Lauren Rantala-Yeary, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Maria Tsimpoukelli, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Ning Zhang, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohan Maheswari, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan

Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Raparthy, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vítor Albiero, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaofang Wang, Xiaoqing Ellen Tan, Xide Xia, Xinfeng Xie, Xuchao Jia, Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aayushi Srivastava, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Amos Teo, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Dong, Annie Franco, Anuj Goyal, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Ce Liu, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Cynthia Gao, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Eric-Tuan Le, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Filippos Kokkinos, Firat Ozgenel, Francesco Caggioni, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hakan Inan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Hongyuan Zhan, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Ilias Leontiadis, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Janice Lam, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan Mc-Phie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kiran Jagadeesh, Kun Huang, Kunal Chawla, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabsa, Manav Avalani, Manish Bhatt, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Miao Liu, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikhil Mehta, Nikolay Pavlovich Laptev, Ning Dong, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Rangaprabhu Parthasarathy, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Russ Howes, Ruty Rinott, Sachin Mehta, Sachin Siby, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Mahajan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shishir Patil, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Summer Deng, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Koehler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaojian Wu, Xiaolan Wang, Xilun Wu, Xinbo Gao, Yaniv Kleinman, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yu Zhao, Yuchen Hao, Yundi Qian, Yunlu Li, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, Zhiwei Zhao, and Zhiyu Ma. 2024. The llama 3 herd of models. Preprint, arXiv:2407.21783.

Rujun Han, Xiang Ren, and Nanyun Peng. 2021. ECONET: Effective continual pretraining of language models for event temporal reasoning. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 5367–5380, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Raghav Jain, Daivik Sojitra, Arkadeep Acharya, Sriparna Saha, Adam Jatowt, and Sandipan Dandapat. 2023. Do language models have a common sense regarding time? revisiting temporal commonsense reasoning in the era of large language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 6750–6774

Oleksandr Kolomiyets, Steven Bethard, and Marie Francine Moens. 2012. Extracting narrative timelines as temporal dependency structures. In *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 88–97.

Vasiliki Kougia, Anastasiia Sedova, Andreas Joseph Stephan, Klim Zaporojets, and Benjamin Roth. 2024. Analysing zero-shot temporal relation extraction on clinical notes using temporal consistency. In *Proceedings of the 23rd Workshop on Biomedical Natural Language Processing*, pages 72–84, Bangkok, Thailand. Association for Computational Linguistics.

Xingzuo Li, Kehai Chen, Yunfei Long, and Min Zhang. 2024. Llm with relation classifier for document-level relation extraction. *arXiv* preprint *arXiv*:2408.13889.

Hieu Man, Nghia Trung Ngo, Linh Ngo Van, and Thien Huu Nguyen. 2022. Selecting optimal context sentences for event-event relation extraction. In *Proceedings of the AAAI conference on artificial intelligence*, volume 36, pages 11058–11066.

Inderjeet Mani, Marc Verhagen, Ben Wellner, Chungmin Lee, and James Pustejovsky. 2006. Machine learning of temporal relations. In *Proceedings of the 21st International Conference on Computational Linguistics and 44th Annual Meeting of the Association for Computational Linguistics*, pages 753–760.

Puneet Mathur, Rajiv Jain, Franck Dernoncourt, Vlad Morariu, Quan Hung Tran, and Dinesh Manocha. 2021. TIMERS: Document-level temporal relation extraction. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers), pages 524–533, Online. Association for Computational Linguistics.

Aakanksha Naik, Luke Breitfeller, and Carolyn Rose. 2019. TDDiscourse: A dataset for discourse-level temporal ordering of events. In *Proceedings of the 20th Annual SIGdial Meeting on Discourse and Dialogue*, pages 239–249, Stockholm, Sweden. Association for Computational Linguistics.

Qiang Ning, Zhili Feng, Hao Wu, and Dan Roth. 2018a. Joint reasoning for temporal and causal relations. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2278–2288, Melbourne, Australia. Association for Computational Linguistics.

Qiang Ning, Sanjay Subramanian, and Dan Roth. 2019. An improved neural baseline for temporal relation extraction. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 6203–6209, Hong Kong, China. Association for Computational Linguistics.

Qiang Ning, Hao Wu, and Dan Roth. 2018b. A multiaxis annotation scheme for event temporal relations. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume* 1: Long Papers), pages 1318–1328, Melbourne, Australia. Association for Computational Linguistics.

Jingcheng Niu, Saifei Liao, Victoria Ng, Simon De Montigny, and Gerald Penn. 2024. ConTempo: A unified temporally contrastive framework for temporal relation extraction. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 1521–1533, Bangkok, Thailand. Association for Computational Linguistics.

Mustafa Ocal, Ning Xie, and Mark Finlayson. 2024. Tlex: An efficient method for extracting exact timelines from timeml temporal graphs. *arXiv preprint arXiv:2406.05265*.

OpenAI, :, Aaron Hurst, Adam Lerer, Adam P. Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, Aleksander Madry, Alex Baker-Whitcomb, Alex Beutel, Alex Borzunov, Alex Carney, Alex Chow, Alex Kirillov, Alex Nichol, Alex Paino, Alex Renzin, Alex Tachard Passos, Alexander Kirillov, Alexi Christakis, Alexis Conneau, Ali Kamali, Allan Jabri, Allison Moyer, Allison Tam, Amadou Crookes, Amin Tootoochian, Amin Tootoonchian, Ananya Kumar, Andrea Vallone, Andrej Karpathy, Andrew Braunstein, Andrew Cann, Andrew Codispoti, Andrew Galu, Andrew Kondrich, Andrew Tulloch, Andrey Mishchenko, Angela Baek, Angela Jiang, Antoine Pelisse, Antonia Woodford, Anuj Gosalia, Arka Dhar, Ashley Pantuliano, Avi Nayak, Avital Oliver, Barret Zoph, Behrooz Ghorbani, Ben Leimberger, Ben Rossen, Ben Sokolowsky, Ben Wang, Benjamin Zweig, Beth Hoover, Blake Samic, Bob McGrew, Bobby Spero, Bogo Giertler, Bowen Cheng, Brad Lightcap, Brandon Walkin, Brendan Quinn, Brian Guarraci, Brian Hsu, Bright Kellogg, Brydon Eastman, Camillo Lugaresi, Carroll Wainwright, Cary Bassin, Cary Hudson, Casey Chu, Chad Nelson, Chak Li, Chan Jun Shern, Channing Conger, Charlotte Barette, Chelsea Voss, Chen Ding, Cheng Lu, Chong Zhang, Chris Beaumont, Chris Hallacy, Chris Koch, Christian Gibson, Christina Kim, Christine Choi, Christine McLeavey, Christopher Hesse, Claudia Fischer, Clemens Winter, Coley Czarnecki, Colin Jarvis, Colin Wei, Constantin Koumouzelis, Dane Sherburn, Daniel Kappler, Daniel Levin, Daniel Levy, David Carr, David Farhi, David Mely, David Robinson, David Sasaki, Denny Jin, Dev Valladares, Dimitris Tsipras, Doug Li, Duc Phong Nguyen, Duncan

Findlay, Edede Oiwoh, Edmund Wong, Ehsan Asdar, Elizabeth Proehl, Elizabeth Yang, Eric Antonow, Eric Kramer, Eric Peterson, Eric Sigler, Eric Wallace, Eugene Brevdo, Evan Mays, Farzad Khorasani, Felipe Petroski Such, Filippo Raso, Francis Zhang, Fred von Lohmann, Freddie Sulit, Gabriel Goh, Gene Oden, Geoff Salmon, Giulio Starace, Greg Brockman, Hadi Salman, Haiming Bao, Haitang Hu, Hannah Wong, Haoyu Wang, Heather Schmidt, Heather Whitney, Heewoo Jun, Hendrik Kirchner, Henrique Ponde de Oliveira Pinto, Hongyu Ren, Huiwen Chang, Hyung Won Chung, Ian Kivlichan, Ian O'Connell, Ian O'Connell, Ian Osband, Ian Silber, Ian Sohl, Ibrahim Okuyucu, Ikai Lan, Ilya Kostrikov, Ilya Sutskever, Ingmar Kanitscheider, Ishaan Gulrajani, Jacob Coxon, Jacob Menick, Jakub Pachocki, James Aung, James Betker, James Crooks, James Lennon, Jamie Kiros, Jan Leike, Jane Park, Jason Kwon, Jason Phang, Jason Teplitz, Jason Wei, Jason Wolfe, Jay Chen, Jeff Harris, Jenia Varavva, Jessica Gan Lee, Jessica Shieh, Ji Lin, Jiahui Yu, Jiayi Weng, Jie Tang, Jieqi Yu, Joanne Jang, Joaquin Quinonero Candela, Joe Beutler, Joe Landers, Joel Parish, Johannes Heidecke, John Schulman, Jonathan Lachman, Jonathan McKay, Jonathan Uesato, Jonathan Ward, Jong Wook Kim, Joost Huizinga, Jordan Sitkin, Jos Kraaijeveld, Josh Gross, Josh Kaplan, Josh Snyder, Joshua Achiam, Joy Jiao, Joyce Lee, Juntang Zhuang, Justyn Harriman, Kai Fricke, Kai Hayashi, Karan Singhal, Katy Shi, Kavin Karthik, Kayla Wood, Kendra Rimbach, Kenny Hsu, Kenny Nguyen, Keren Gu-Lemberg, Kevin Button, Kevin Liu, Kiel Howe, Krithika Muthukumar, Kyle Luther, Lama Ahmad, Larry Kai, Lauren Itow, Lauren Workman, Leher Pathak, Leo Chen, Li Jing, Lia Guy, Liam Fedus, Liang Zhou, Lien Mamitsuka, Lilian Weng, Lindsay McCallum, Lindsey Held, Long Ouyang, Louis Feuvrier, Lu Zhang, Lukas Kondraciuk, Lukasz Kaiser, Luke Hewitt, Luke Metz, Lyric Doshi, Mada Aflak, Maddie Simens, Madelaine Boyd, Madeleine Thompson, Marat Dukhan, Mark Chen, Mark Gray, Mark Hudnall, Marvin Zhang, Marwan Aljubeh, Mateusz Litwin, Matthew Zeng, Max Johnson, Maya Shetty, Mayank Gupta, Meghan Shah, Mehmet Yatbaz, Meng Jia Yang, Mengchao Zhong, Mia Glaese, Mianna Chen, Michael Janner, Michael Lampe, Michael Petrov, Michael Wu, Michele Wang, Michelle Fradin, Michelle Pokrass, Miguel Castro, Miguel Oom Temudo de Castro, Mikhail Pavlov, Miles Brundage, Miles Wang, Minal Khan, Mira Murati, Mo Bavarian, Molly Lin, Murat Yesildal, Nacho Soto, Natalia Gimelshein, Natalie Cone, Natalie Staudacher, Natalie Summers, Natan LaFontaine, Neil Chowdhury, Nick Ryder, Nick Stathas, Nick Turley, Nik Tezak, Niko Felix, Nithanth Kudige, Nitish Keskar, Noah Deutsch, Noel Bundick, Nora Puckett, Ofir Nachum, Ola Okelola, Oleg Boiko, Oleg Murk, Oliver Jaffe, Olivia Watkins, Olivier Godement, Owen Campbell-Moore, Patrick Chao, Paul McMillan, Pavel Belov, Peng Su, Peter Bak, Peter Bakkum, Peter Deng, Peter Dolan, Peter Hoeschele, Peter Welinder, Phil Tillet, Philip Pronin, Philippe Tillet, Prafulla Dhariwal, Qiming Yuan, Rachel Dias, Rachel Lim, Rahul Arora, Rajan Troll, Randall Lin, Rapha Gontijo Lopes, Raul Puri, Reah Miyara, Reimar Leike, Renaud Gaubert, Reza Zamani, Ricky Wang, Rob Donnelly, Rob Honsby, Rocky Smith, Rohan Sahai, Rohit Ramchandani, Romain Huet, Rory Carmichael, Rowan Zellers, Roy Chen, Ruby Chen, Ruslan Nigmatullin, Ryan Cheu, Saachi Jain, Sam Altman, Sam Schoenholz, Sam Toizer, Samuel Miserendino, Sandhini Agarwal, Sara Culver, Scott Ethersmith, Scott Gray, Sean Grove, Sean Metzger, Shamez Hermani, Shantanu Jain, Shengjia Zhao, Sherwin Wu, Shino Jomoto, Shirong Wu, Shuaiqi, Xia, Sonia Phene, Spencer Papay, Srinivas Narayanan, Steve Coffey, Steve Lee, Stewart Hall, Suchir Balaji, Tal Broda, Tal Stramer, Tao Xu, Tarun Gogineni, Taya Christianson, Ted Sanders, Tejal Patwardhan, Thomas Cunninghman, Thomas Degry, Thomas Dimson, Thomas Raoux, Thomas Shadwell, Tianhao Zheng, Todd Underwood, Todor Markov, Toki Sherbakov, Tom Rubin, Tom Stasi, Tomer Kaftan, Tristan Heywood, Troy Peterson, Tyce Walters, Tyna Eloundou, Valerie Qi, Veit Moeller, Vinnie Monaco, Vishal Kuo, Vlad Fomenko, Wayne Chang, Weiyi Zheng, Wenda Zhou, Wesam Manassra, Will Sheu, Wojciech Zaremba, Yash Patil, Yilei Qian, Yongjik Kim, Youlong Cheng, Yu Zhang, Yuchen He, Yuchen Zhang, Yujia Jin, Yunxing Dai, and Yury Malkov. 2024. Gpt-4o system card. Preprint, arXiv:2410.21276.

Gabriel Roccabruna, Massimo Rizzoli, and Giuseppe Riccardi. 2024. Will LLMs replace the encoder-only models in temporal relation classification? In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 20402–20415, Miami, Florida, USA. Association for Computational Linguistics.

Anna Rogers, Marzena Karpinska, Ankita Gupta, Vladislav Lialin, Gregory Smelkov, and Anna Rumshisky. 2024. NarrativeTime: Dense temporal annotation on a timeline. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 12053–12073, Torino, Italia. ELRA and ICCL.

Roser Saurí, Jessica Littman, Bob Knippen, Robert Gaizauskas, Andrea Setzer, and James Pustejovsky. 2006. Timeml annotation guidelines. *Version*, 1(1):31.

Emre Sezgin, Syed-Amad Hussain, Steve Rust, and Yungui Huang. 2023. Extracting medical information from free-text and unstructured patient-generated health data using natural language processing methods: feasibility study with real-world data. *JMIR Formative Research*, 7:e43014.

Haoyu Wang, Hongming Zhang, Yuqian Deng, Jacob Gardner, Dan Roth, and Muhao Chen. 2023. Extracting or guessing? improving faithfulness of event temporal relation extraction. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 541–553,

Dubrovnik, Croatia. Association for Computational Linguistics.

Liang Wang, Peifeng Li, and Sheng Xu. 2022. DCT-centered temporal relation extraction. In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 2087–2097, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.

Nianwen Xue and Yuchen Zhang. 2018. Neural ranking models for temporal dependency structure parsing. In 2018 Conference on Empirical Methods in Natural Language Processing.

Chenhan Yuan, Qianqian Xie, and Sophia Ananiadou. 2023. Zero-shot temporal relation extraction with ChatGPT. In *The 22nd Workshop on Biomedical Natural Language Processing and BioNLP Shared Tasks*, pages 92–102, Toronto, Canada. Association for Computational Linguistics.

Shuaicheng Zhang, Qiang Ning, and Lifu Huang. 2022. Extracting temporal event relation with syntax-guided graph transformer. In *Findings of the Association for Computational Linguistics: NAACL 2022*, pages 379–390, Seattle, United States. Association for Computational Linguistics.

Jie Zhou, Shenpo Dong, Hongkui Tu, Xiaodong Wang, and Yong Dou. 2022. RSGT: Relational structure guided temporal relation extraction. In *Proceedings* of the 29th International Conference on Computational Linguistics, pages 2001–2010, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.

Yichao Zhou, Yu Yan, Rujun Han, J Harry Caufield, Kai-Wei Chang, Yizhou Sun, Peipei Ping, and Wei Wang. 2021. Clinical temporal relation extraction with probabilistic soft logic regularization and global inference. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 14647–14655.

A Prompts

Figures 3-7 provide full examples of the prompts used to evaluate the LLMs. In all prompt examples, we use MATRES's four relations, while Narrative-Time prompts follow the same format but include seven relations. Figure 3 presents the prompt for the zero-shot pairwise approach, followed by its few-shot extension in Figure 4. Similarly, Figures 5 and 6 illustrate the prompts for the graph generation approach. Finally, Figure 7 shows the prompt used for breaking cycles with the LLM.

In all examples, we include the *vague* relation, though we also conduct experiments without it.

B Experimental Settings

Both datasets we used in this study, MATRES and NarrativeTime contain documents written in English, and cover news articles. For pairwise model experiments, we evaluate the MATRES approach and sliding-window using Llama-3.1-8B-Instruct-Turbo and Llama-3.2-3B-Instruct-Turbo with Float16 quantization on an NVIDIA GeForce RTX 3090, alongside GPT-4o-mini-2024-07-18 and GPT-4o-2024-08-06, incurring a total cost of approximately \$20.

For graph-based model experiments, the same models are used, with OpenAI models costing around \$15.

For cycle-breaking models, we train an encoder-based model using the BERT architecture on an NVIDIA GeForce RTX 3090 for five epochs. Training takes approximately 1 hour for MATRES and 7 hours for NarrativeTime. Subsequently, GPT-4o-mini-2024-07-18 and GPT-4o-2024-08-06 are used, with a combined cost of approximately \$35.

```
"role": "system",
  "content":
    Task Overview:
    You are given a text, in which some verbs are uniquely marked by [EVENT#ID]event
        [/EVENT#ID] (e.g., [EVENT1]event1[/EVENT1], [EVENT2]event2[/EVENT2]).
    Your task is to say which of the verbs happened first in a chronological order.
    More specifically, you need to return for each pair of verbs, which is two
        sentence apart,
    a single label out of the listed potential labels:
    before - the first verb happened before the second. after - the first verb happened after the second.
    equal - both verbs happened together.
    vague - It is impossible to know based on the context provided
    you should only provide one classification."
},
  "role": "user",
  "content":
    Text for Analysis:
    Former President Nicolas Sarkozy was [EVENT1]informed[/EVENT1] Thursday that he
        would face a formal investigation into whether he [EVENT3]abused[/EVENT3]
        the frailty of Liliane Bettencourt, 90, the heiress to the L'Oreal fortune and France's richest woman, to get funds for his 2007 presidential campaign.
         Mr. Sarkozy has denied accepting illegal campaign funds from Ms.
        Bettencourt, either personally or through his party treasurer at the time,
        Eric Woerth, as alleged by her former butler.
in one word --> "
```

Figure 3: Zero-shot prompt for pairwise classification.

```
"role": "system",
  "content":
    <INSTRUCTIONS>
    Examples:
    #########
    Text for Analysis:
    NAIROBI, Kenya (AP)
    Suspected bombs [EVENT1]exploded[/EVENT1] outside the U.S. embassies in the
        Kenyan and Tanzanian capitals Friday, [EVENT2]killing[/EVENT2] dozens of
        people, witnesses said.
    --> before
    Text for Analysis:
    Suspected bombs exploded outside the U.S. embassies in the Kenyan and Tanzanian
        capitals Friday, killing dozens of people, witnesses [EVENT3]said[/EVENT3].
    The American ambassador to Kenya was among hundreds [EVENT12]injured[/EVENT12],
       a local TV said.
    --> after
    #######"
{
  "role": "user",
  "content": "
    Text for Analysis:
    <TEXT>
in one word -->"
}
```

Figure 4: Few-shot prompt for pairwise classification. <INSTRUCTIONS> is a placeholder for the instructions provided in Figure 3.

```
"role": "system",
 "content":
    Task Overview:
    You are given a text, in which some verbs are uniquely marked by [EVENT#ID]event
       [/EVENT#ID] (e.g., [EVENT1]event1[/EVENT1], [EVENT2]event2[/EVENT2]).
    Your task is to say which of the verbs happened first in a chronological order.
   More specifically, you need to return for each pair of verbs, which is two
       sentence apart,
    a single label out of the listed potential labels:
   before - the first verb happened before the second.
    after - the first verb happened after the second.
    equal - both verbs happened together.
   vague - It is impossible to know based on the context provided
   All responses should be valid and compact dot graph format.
   compact meaning:
    - do not mention transitive dependencies - if eil BEFORE ei2 and ei2 BEFORE ei3
       don't write ei1 BEFORE ei3
    - do not mention symmetric relation - if ei1 BEFORE ei2 don't write ei2 AFTER
       ei1'
},
 "role": "user",
 "content": "---
   Text for Analysis:
   The flu season is winding down, and it has [EVENT2]killed[/EVENT2] 105 children
       so far - about the average toll.
        The season [EVENT3] started[/EVENT3] about a month earlier than usual, [
           EVENT4] sparking[/EVENT4] concerns it might turn into the worst in a
           decade.
Respond only with valid dot graph format with the approprite markers and attributes
   (like label). Do not write an introduction or summary.
the graph:"
}
```

Figure 5: Zero-shot prompt for generating the entire temporal graph.

```
"role": "system",
  "content":
    <INSTRUCTIONS>
    Example:
    ########
    Text for Analysis:
    NAIROBI, Kenya (AP)
    Suspected bombs [EVENT1]exploded[/EVENT1] outside the U.S. embassies in the Kenyan and Tanzanian capitals Friday, [EVENT2]killing[/EVENT2] dozens of
         people, witnesses [EVENT3]said[/EVENT3].
    the sample of correct labels are:
    digraph {
         "EVENT1" -> "EVENT2" [label="before"];
         "EVENT3" -> "EVENT12" [label="after"];
"EVENT4" -> "EVENT5" [label="vague"];
    ########
},
  "role": "user",
  "content": "
    Text for Analysis:
    <TEXT>
Respond only with valid dot graph format with the approprite markers and attributes
    (like label). Do not write an introduction or summary.
the graph:"
}
```

Figure 6: Few-shot prompt for generating the entire temporal graph.

```
"role": "system",
  "content":
    Task Overview:
    You are given a text, in which some events are uniquely marked by [EVENT#ID]
        event[/EVENT#ID] (e.g., [EVENT1]event1[/EVENT1], [EVENT2]event2[/EVENT2]),
    and a dot graph which represent chronological order with error, where some edges
         form cycles.
    Your task is to decide which pair to drop (by his unique_id), being concise and
       removing the minimum number of edges.
    Pay attention, I used classifier to choose the most fitted relation (label
        attribute in dot graph)
    and score which represent the confidence of the classifier.
    relation meaning:
    before - the first verb happened before the second. after - the first verb happened after the second.
    equal - both events happen simultaneously
    vague - temporal order cannot be determined from the context"
 "role": "user",
  "content": "
    Text for Analysis:
    Barack Obama would make a great stand-up comic, not because he's the funniest
        president ever but because he uses jokes the same way many of us comedians
        do: as a weapon.
        Traditionally, the (intentionally) funny lines by our presidents have had one thing in common: They were self-deprecating. Sure, some presidents
            have [EVENT5]used[/EVENT5] jokes to take jabs at their opponents, but
            not to the extent of Obama.
        During his tenure, he has increasingly [EVENT8]unleashed[/EVENT8] biting
            comedic barbs against his critics and political adversaries. These jokes
             are [EVENT1000]intended[/EVENT1000] to do more than simply entertain
            you. They have an agenda.
        Obama's humor is often delivered the way a comedian dealing with a heckler
            would do it. He tries to undermine his opponents with it and get the
            crowd -- in this case the public -- on his side. I can [EVENT20]assure[/
            EVENT20] you that having a crowd laugh at your critic/heckler is not
            only effective in dominating them, it's also very satisfying.
    digraph Chronology {
        "EVENT5" -> "EVENT8" [label="BEFORE", score=0.71996284, unique_id=0];
        "EVENT8" -> "EVENT1000" [label="BEFORE", score=0.8759634, unique_id=1];
        "EVENT5" -> "EVENT20" [label="AFTER", score=0.9743732, unique_id=2];
        "EVENT1000" -> "EVENT20" [label="BEFORE", score=0.75076234, unique_id=3];
    }
Respond only with the unique_id list to drop (wrong label)"
}
```

Figure 7: Zero-shot prompt for cycle breaking.