

Transformer-Based Temporal Information Extraction and Application: A Review

Anonymous ACL submission

Abstract

Temporal information extraction (IE) aims to extract structured temporal information from unstructured text, thereby uncovering the implicit timelines within. This technique is applied across domains such as healthcare, newswire, and intelligence analysis, aiding models in these areas to perform temporal reasoning and enabling human users to grasp the temporal structure of text. Transformer-based pre-trained language models have produced revolutionary advancements in natural language processing, demonstrating exceptional performance across a multitude of tasks. Despite the achievements garnered by Transformer-based approaches in temporal IE, there is a lack of comprehensive reviews on these endeavors. In this paper, we aim to bridge this gap by systematically summarizing and analyzing the body of work on temporal IE using Transformers while highlighting potential future research directions.

1 Introduction

Temporal information extraction (IE) is a critical task in natural language processing (NLP). Its objective is to extract structured temporal information from unstructured text, thereby revealing the implicit timelines within the text. This not only helps improve temporal reasoning in other NLP tasks, such as timeline summarization and temporal question answering, but also helps human users in gaining a deeper understanding of the evolution of text content over time. For example, Figure 1 displays a snippet of George Washington’s Wikipedia page and the timeline of his position changes; relying solely on text-heavy documents to trace his position changes over different years is time-consuming and may lack accuracy as facts and temporal expressions are scattered throughout the text. In contrast, a timeline enables both NLP models and humans to understand the changes in these positions over time more succinctly and clearly. The application of this

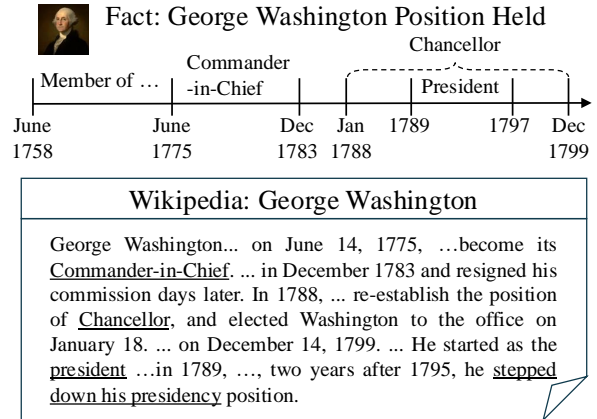


Figure 1: A snippet from George Washington’s Wikipedia page and a timeline regarding his positions.

structured temporal information is not limited to Wikipedia but is also widely used in other domains such as healthcare (Styler IV et al., 2014).

The advent of the Transformer architecture (Vaswani et al., 2017) has sparked a revolutionary change in the field of NLP, particularly with the recent Transformer-based generative large language models (LLM), such as LLAMA3 (Dubey et al., 2024) and GPT-4 (Achiam et al., 2023), demonstrating exceptional performance across many tasks. Nevertheless, there has yet to be an in-depth study that provides a comprehensive review or analysis of the Transformer architecture’s application in the field of temporal IE. Existing surveys (Lim et al., 2019; Leeuwenberg and Moens, 2019; Alfattni et al., 2020; Olex and McInnes, 2021) focus on rule-based systems or traditional machine learning models (e.g., support vector machines) which are reliant on hand-crafted features. Only Olex and McInnes (2021) touches on the application of Transformer models, but they offer only a brief description of BERT-style models and focus largely on the clinical domain.

To address this gap, we systematically review the applications of Transformer-based models in

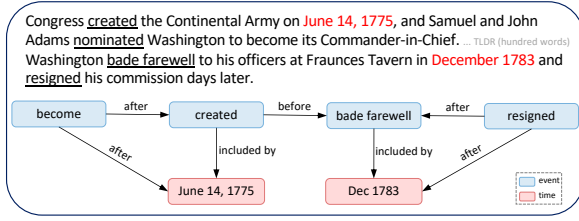


Figure 2: A snippet from George Washington’s Wikipedia page and the corresponding temporal graph.

the field of temporal IE. Broadly, temporal IE refers to any tasks involving the extraction of temporal information from text. We focus on three important tasks which are defined in the most widely adopted temporal IE annotation framework, TimeML (James, 2003): time expression identification, time expression normalization, and temporal relation extraction. Our contributions are summarized as follows: (1) We systematically review, summarize, and categorize the existing temporal IE datasets, Transformer-based methods, and applications. (2) We identify and highlight the research gaps in the field of temporal IE and suggest potential directions for future research.

2 Overview

The goal of temporal IE is to extract structured temporal information from unstructured text, facilitating its interpretation and processing by computers, thereby achieving a transformation from text to structure. The final result of a temporal IE system is the construction of a directed acyclic graph, or a temporal graph, which represents the structured temporal information in the text. In the temporal graph, nodes represent time expressions and events (temporal entities), while edges depict the temporal relations between these nodes, such as “before,” “after,” etc. For instance, Figure 2 illustrates a text snippet from George Washington’s Wikipedia page and its corresponding temporal graph.

Constructing a temporal graph involves several sub-tasks: time expression identification, time expression normalization, event extraction, and temporal relation extraction. The following is a brief introduction to these sub-tasks; see Appendix A for a discussion of common evaluation methods.

Time Expression Identification and Normalization Time expression identification refers to identifying specific time points, durations, or periods within the text, such as the explicitly dateable expression “February 25, 2024,” or more ambiguous

expressions like “three days ago” (James, 2003). Time normalization involves converting identified expressions into a standardized format to improve their interpretability. For example, under the ISO-TimeML framework (Pustejovsky et al., 2010), “February 25, 2024” might be converted into the TIMEX3 format as “2024-02-25”.

Event Trigger Extraction In temporal IE, event extraction differs from other NLP event extraction tasks; it simply marks the event trigger words that represent actions, such as “accident” in “about two weeks after the accident occurred”. We will not review event extraction works because, to our knowledge, there is currently no temporal IE research focused solely on event extraction. Furthermore, most existing work on temporal IE assumes that event triggers have already been identified. For a comprehensive survey of event extraction, we refer readers to (Li et al., 2022).

Temporal Relation Extraction The task of temporal relation extraction aims to identify the temporal relations among given events and time expressions. Common temporal relations include before, after, and simultaneous. For example, in Figure 2, the temporal relation between “June 14, 1775” and the event “become” is marked as “after”.

3 Datasets

A clearly defined annotation framework is essential when constructing a dataset for temporal IE. It needs to precisely define time expressions, events, and their relations. We summarize all the datasets in Table 1 of Appendix B.

3.1 TimeML Annotation Framework Datasets

An end-to-end temporal IE dataset encompasses various tasks, including the identification and normalization of time expressions and the extraction of temporal relations. Most end-to-end temporal information datasets have been based on the TimeML framework (James, 2003) or its derivatives, such as ISO-TimeML (Pustejovsky et al., 2010). We present datasets based on the TimeML framework in the first section of Table 1.

TimeBank (James, 2003) was the first dataset to adopt the TimeML framework, focusing on the English news domain. Follow-up works included the TempEval shared task series (Verhagen et al., 2007, 2010; UzZaman et al., 2013), covering multiple languages, including Chinese, English, Ital-

ian, French, Korean, and Spanish. There are also language-specific datasets like French TimeBank (Bittar et al., 2011), Spanish TimeBank (Nieto et al., 2011), Portuguese TimeBank (Costa and Branco, 2012), Japanese TimeBank (Asahara et al., 2013), Italian TimeBank (Bracchi et al., 2016), and Korean TimeBank (Lim et al., 2018). Similarly, the MeanTime dataset (Minard et al., 2016) offers data in English, Italian, Spanish, and Dutch. Datasets based on TimeML and its variants showcase language diversity and also cover several different domains: the Spanish TimeBank focuses on history text, the Korean TimeBank is based on Wikipedia content, and the Richer Event Description dataset (O’Gorman et al., 2016) provides data from both news and forum discussion domains.

Additionally, efforts have been made to improve the temporal relation annotations in the original TimeBank. TimeBank-Dense (Chambers et al., 2014) addresses the sparsity of temporal relation annotations in TimeBank by requiring annotators to label all temporal relations within a given scope, thus increasing the number of temporal relations in the dataset. The TORDER dataset (Cheng and Miyao, 2018) annotates the same documents as TimeBank-Dense, introducing temporal relations automatically by anchoring times and events to absolute points, reducing the annotation burden. The MATRES dataset (Ning et al., 2018) focuses on events from TimeBank-Dense, anchoring events to different timelines and comparing their start times to enhance inter-annotator consistency.

Several datasets have been developed specific to the clinical domain, of which the Thyme datasets (Bethard et al., 2015, 2016, 2017) are most notable. They are based on the Thyme-TimeML (Styler IV et al., 2014) annotation framework, which adjusts and adds new temporal attributes from ISO-TimeML to suit medical texts. Like the TimeBank series, the Thyme dataset involves identifying and normalizing time expressions and extracting temporal relations, focusing on English. Another similar dataset is i2b2-2012 (Sun et al., 2013), which adapts the TimeML framework for clinical texts.

Besides end-to-end datasets, several others based on TimeML or its variants focus on specific temporal IE tasks. For instance, the AncientTimes dataset (Strötgen et al., 2014) covers a broad range of languages, concentrating on the identification and normalization of time expressions. The TD-Discourse dataset (Naik et al., 2019), based on

TimeBank-Dense, expands the annotation window for temporal relations, focusing on their extraction. The German time expression (Strötgen et al., 2018) and German VTEs (May et al., 2021) datasets are dedicated to identifying and normalizing time expressions in German. The PATE dataset (Zarcone et al., 2020) provides data aimed at time expression identification and normalization for the virtual assistant domain.

3.2 Other Annotation Framework Datasets

Unlike datasets for temporal IE based on TimeML, other annotation frameworks typically focus on specific sub-tasks of temporal IE, such as time expression identification and normalization or the extraction of temporal relations. We present these datasets in the second section of Table 1.

For time expression identification and normalization, WikiWars (Mazur and Dale, 2010) and SCATE (Laparra et al., 2018) are two major datasets. WikiWars contains data from English and German Wikipedia, annotated based on TIMEX2 (a precursor to TimeML’s TIMEX3) to mark explicit time expressions. The SCATE dataset, based on English news and clinical documents, aims to address limitations in TimeML that prevent expressing multiple calendar units, times relative to events, and compositional time expressions. To achieve this, SCATE represents time expressions as compositions of temporal operators.

For temporal relations, there are datasets based on the temporal dependency tree/graph (Zhang and Xue, 2018, 2019; Yao et al., 2020) and CaTeRS (Mostafazadeh et al., 2016) frameworks. Unlike the pairwise temporal relations considered in the TimeML framework, temporal dependency tree assumes that all time expressions and events in a document have a reference time, allowing for the representation of overall temporal relations through a dependency tree. The subsequent temporal dependency graph dataset (Yao et al., 2020) relaxed this assumption by enabling each event in a document to have a reference event, a reference time, or both, thus forming a temporal graph structure. The temporal dependency tree dataset covers news and narrative domains in English and Chinese, while the temporal dependency graph dataset focuses on English news. Meanwhile, CaTeRS concentrates on analyzing temporal relations between events in English commonsense stories, with event definitions based on ontologies, different from the verb-,

adjective-, or noun-based definitions in TimeML. CaTeRS’ annotation of temporal relations is story-wide, with a simplified set of relations.

3.3 Discussion and Research Gaps

Domain Bias Existing annotated datasets exhibit significant domain biases. As demonstrated in Table 1, among the 32 datasets we reviewed, 20 (or 63%) are predominantly focused on the newswire domain. While temporal information is crucial for understanding news content, an excessive concentration in a single domain hampers the advancement and generalizability of systems trained on these datasets, since the challenges and difficulties encountered in temporal IE vary across different domains. Notably, the Clinical TempEval 2017 shared task (Bethard et al., 2017) reveals that most tasks suffer an approximately 20-point drop in performance in a cross-domain setting, underscoring how domain shifts can significantly degrade model accuracy. For example, temporal information, especially time expressions, in newswire texts tend to be explicitly stated, whereas in other domains, like historical Wikipedia entries, they might appear in subtler ways. Consider a statement from a page about George Washington that reads, “... 1798, one year after that, he stepped down from the presidency,” which would demand a more nuanced interpretation for accurate time normalization. Cultivating datasets that represent a variety of domains is vital to driving innovation in temporal IE.

Language Diversity Unlike the domain homogeneity of the datasets, the existing datasets display rich linguistic diversity, covering 15 different languages. The representation of time varies across languages, and even when semantically similar, the specific time intervals on the timeline can differ. For example, analysis in Shwartz (2022) shows that different cultures/languages have significant variations in the understanding of “night” and “evening” during the day. One instance is that Brazilian Portuguese speakers often use “evening” and “night” interchangeably to denote the same time period, possibly because the tropical climate in Brazil causes evening to transition quickly into night. However, this might not be applicable to other cultures or languages. Therefore, the language diversity in datasets is crucial for developing models capable of effectively extracting temporal information across different languages.

Annotation and Dataset Framework Development Slows Down Aside from the original TimeML and some incremental modifications to it, no new end-to-end temporal IE annotation frameworks have been proposed. A significant issue with the existing TimeML-based annotation frameworks is the limited amount of information that the resultant temporal graphs can represent. For instance, in Figure 2, we only see trigger words for events, time expressions, and some temporal relations. When these temporal graphs are isolated from their original context and treated as stand-alone entities, they struggle to provide a comprehensive understanding of the textual information. This might explain why, in the upcoming Section 6, we see no work directly employing these extracted temporal graphs for reasoning to accomplish specific tasks, such as answering temporal questions. Instead, these temporal graphs are used as auxiliary tools or additional knowledge to assist task-specific models in temporal reasoning.

In addition to the stagnation in the innovation of end-to-end annotation frameworks, there has been a notable decline in dataset development efforts in the field of temporal IE in recent years. This trend may primarily stem from the intrinsic complexity of the annotation process for temporal IE datasets. Such complexity accounts for the low annotator agreement observed in many annotation tasks (Cassidy et al., 2014). Furthermore, as demonstrated by analysis in Su et al. (2021), even Ph.D. students in relevant fields find it challenging to comprehend annotation guidelines and annotate high-quality data within a short period. These issues highlight the difficulties in developing temporal IE datasets, suggesting that improvements in the annotation framework might be necessary to address these challenges.

4 Time Expression Methods

4.1 Methods Overview

In the realm of time expression identification, most prior work (Almasian et al., 2021; Chen et al., 2019; Mirzababaei et al., 2022; Olex and McInnes, 2022; Laparra et al., 2021; Almasian et al., 2022; Cao et al., 2022) leverages discriminative models built upon transformer encoders like BERT (Devlin et al., 2019). These approaches typically frame time expression identification as a token classification task, wherein a sequence of tokens is input, processed through a base encoder model to

obtain contextualized representations, and these representations are fed into a classifier (such as a simple linear classification layer or a Conditional Random Field layer) to identify time expressions and their specific types. Almasian et al. (2021) is the only work exploring a generative approach for time expression identification, framing the task as a sequence-to-sequence problem and employing a pair of transformer encoders to formulate an encoder-decoder model—where one serves as the encoder and the other as the decoder—to generate additional TIMEX3 tags for the input, thereby recognizing time expressions and their types.

Shwartz (2022) and Kim et al. (2020) focus on the normalization of time expressions and use transformer-based models. Shwartz (2022) aims to normalize time expressions from various cultural contexts (e.g., morning, noon, afternoon) into precise hourly representations within a day. They train a BERT model with a masked language modeling task to predict specific times of day that are masked, given the time expressions. Kim et al. (2020) seeks to normalize time expressions in novels into specific daily hours, fine-tuning the BERT model for a 24-class classification task to ascertain the corresponding times of day for given expressions.

Lange et al. (2023) addresses both extraction and normalization of time expressions, adopting a pipeline approach. Initially, they fine-tune the XLM-R model using the token classification method to extract time expressions, then denote identified expressions with TIMEX3 tags with masked time values, and finally fine-tune the XLM-R model with masked language modeling to predict the normalized masked time values.

Several of the aforementioned works also utilize data augmentation techniques to improve the model’s multilingual performance (Lange et al., 2023; Mirzababaei et al., 2022; Almasian et al., 2022). For instance, Lange et al. (2023) employs the rule-based HeidelTime method (Strötgen and Gertz, 2010) to annotate time expressions and their normalizations across 87 languages, generating a semi-supervised dataset to facilitate model training.

4.2 Discussion and Research Gaps

Despite the significant achievements of Transformer models in various NLP tasks, research in the area of time expression identification and normalization has remained relatively limited over the past few years. This is particularly true of time nor-

malization, where the volume and depth of research are low, especially when compared to similar tasks such as named entity recognition, entity normalization, and entity linking. Furthermore, the methodological diversity in existing works is notably constrained, with most research relying on pre-trained Transformer models for simple token classification. While generative LLMs like GPT-4 or LLAMA3 have demonstrated impressive performance in other NLP tasks, their potential in the identification and normalization of time expressions has barely been explored. This suggests a significant research gap exists; exploration of generative approaches may offer the potential for advancement in time expression identification and normalization.

5 Temporal Relation Methods

The task of temporal relation extraction typically assumes that events and time expressions in the text have already been identified, with the only objective being to extract the temporal relations between them. We summarize all the reviewed temporal relation extraction works in Appendix C Table 2. Discriminative methods typically employ a pretrained discriminative language model like BERT or RoBERTa (Liu et al., 2019) as the base encoder model to derive contextualized representations of events or time expressions. Subsequently, these representations are paired and input into a classification layer for a multi-class classification task, with each class representing a different temporal relation. Generative methods typically leverage encoder-decoder models such as T5 (Raffel et al., 2020) or decoder-only models like GPT (Radford et al., 2019) to generate a target sequence that encapsulates the temporal relation between the input events and times. These methods often rely on post-processing techniques to extract specific temporal relations from the predicted target sequences.

5.1 Discriminative Methods Overview

Works on discriminative temporal relation extraction have mainly focused on integrating external knowledge and improving model robustness.

5.1.1 Integrating External Knowledge

Commonsense Knowledge Commonsense knowledge for temporal relations usually involves typical sequences of events, such as eating typically occurring after cooking. Such commonsense knowledge might be fundamental for humans, but absent from the base encoder model. Ning et al.

(2019), Wang et al. (2020) and Tan et al. (2023) integrated knowledge from external commonsense knowledge graphs. Tan et al. (2023) employs a complex Bayesian learning method to merge the knowledge with the contextualized representations from the base encoder, whereas Ning et al. (2019) and Wang et al. (2020) simply concatenate the vectorized representations of the commonsense knowledge with those from the base encoder.

Syntactic and Semantic Knowledge Syntactic and semantic knowledge, typically extracted using off-the-shelf external tools or straightforward rules, enrich the base encoder models’ representations. For instance, Wang et al. (2022) utilizes SpaCy’s dependency parser to parse the syntactic dependency trees from the input text and neuralcoref to identify coreferential relationships among entities. Mathur et al. (2021) employs the discoursegraphs library to parse rhetorical dependency graphs from the text. To integrate this structured knowledge into the contextualized event or time expression representations, graph neural networks are often employed over syntactic or semantic pairwise relations (Wang et al., 2022; Mathur et al., 2022; Zhou et al., 2022; Mathur et al., 2021). For example, Wang et al. (2022) first encodes an input sequence containing event pairs with the RoBERTa model to generate initial contextual representations, which are then enhanced with extracted syntactic and semantic knowledge using additional graph neural network layers. Another method is to prelearn or extract vectorized representations of the knowledge, which are later concatenated with the event or time expression representations (Ross et al., 2020; Wang et al., 2020; Han et al., 2019a; Ning et al., 2019; Han et al., 2019b), as in Wang et al. (2020), where RoBERTa token embeddings and one-hot vectors of part-of-speech tags are combined.

Temporal-Specific Rules These rules are intrinsic to temporal relations themselves, with symmetry and transitivity being the most common. For instance, if event A happens before event B, then symmetry can be used to infer that B happens after A. And if A precedes B and B precedes C, transitivity can be used to infer that A precedes C. Detailed explanations of the symmetry and transitivity rules and a comprehensive transitivity table are provided in Ning et al. (2019). Such rules can be embedded during the model training phase, enabling the model to learn the characteristics of these tempo-

ral relations. Hwang et al. (2022) and Tan et al. (2021) utilize box embedding and hyperbolic embedding, respectively, to implicitly guide the model in understanding and learning the symmetry and transitivity rules. Zhou et al. (2021) and Wang et al. (2020) translate the constraints of temporal relations into regularization terms for the loss function during training to penalize predictions that violate these rules. Alternatively, rules can be embedded during the inference phase to ensure that all deduced temporal relations adhere to the symmetry and transitivity rules as closely as possible. Custom heuristics in Wang et al. (2022); Zhou et al. (2022, 2021); Liu et al. (2021) exclude temporal relations that contravene rules during inference. Wang et al. (2020) and Han et al. (2019c) formulate the inference of temporal relations as a linear programming problem, optimizing the solution to achieve optimal outcomes. Han et al. (2019a) interprets the discriminative model’s output probabilities as confidence scores for potential relations between temporal entity pairs and employs a structured support vector machine for the final predictions.

Label Distribution Knowledge of label distribution pertains to the frequency distribution of specific temporal relations in the training set. Wang et al. (2023) and Han et al. (2020) integrate this distribution knowledge into their frameworks, using it as a regularization term in the loss function or for inference-time linear programming, aiming to mitigate potential biases in model predictions.

5.1.2 Improving Model Robustness

Multitask Learning Wang et al. (2022), Lin et al. (2020) and Cheng et al. (2020) categorize temporal relations and treat the extraction of different types of temporal relations as independent tasks, employing multitask learning to extract all types of relations simultaneously. For instance, Wang et al. (2022) delineates tasks into event-event, event-time, and event-document creation time, undergoing multitask training across these three tasks. Mathur et al. (2022) applies multitask learning in their model to concurrently predict temporal relations and dependency links between nodes in a temporal dependency tree. Similarly, Ballesteros et al. (2020) implements multitask learning by integrating the extraction of temporal relations with the extraction of entity relations in the general domain.

Data Augmentation Wang et al. (2023) generates counterfactual instances from the training set

samples to mitigate model bias, while Tiesen and Lishuang (2022) employs predefined templates to create additional training examples.

Continued Pre-training of Base Encoder In Zhao et al. (2021) and Han et al. (2021), heuristic methods are used to identify temporal indicators in a corpus of unlabeled data, further training the base encoder using a masked language modeling (MLM) approach to recover masked indicators. Lin et al. (2019) focuses on the medical domain, using MLM on electronic health records from MIMIC-III to adapt the base encoder for domain-specific training prior to temporal relation extraction.

Adversarial Training Kanashiro Pereira (2022) and Pereira et al. (2021) introduce adversarial perturbations at different layers of the transformer encoder during training to enhance model robustness.

Self-training Cao et al. (2021) and Ballesteros et al. (2020) initially train a temporal relation extraction model on annotated datasets and then apply the model to unlabeled data to obtain model-generated labels as pseudo labels. They subsequently select pseudo-labeled examples as sliver examples based on the model’s uncertainty scores and confidence scores (probability scores for specific temporal relation predictions) to train the model.

5.2 Generative Methods Overview

Unlike the task of extracting relations between general entities for constructing knowledge graphs (refer to survey Ye et al. (2022)), few generative approaches have been proposed and applied in the field of temporal relation extraction. Dligach et al. (2022) utilizes an encoder-decoder model architecture, specifically the BART (Lewis et al., 2020) and T5 (Raffel et al., 2020) models. They primarily investigate how to fine-tune these encoder-decoder models for temporal relation extraction tasks, focusing on the input and output formats. They discover that producing outputs for each event and time pair separately is more effective than the intuitive triplet form, i.e., (entity, relation, entity). On the other hand, Yuan et al. (2023) concentrates on examining the capabilities of the powerful ChatGPT generative model, in the context of temporal relation extraction, testing various prompting methods, such as zero-shot prompting, and the popular chain-of-thought prompting (Wei et al., 2022). Their findings indicate that, despite using these prompting methods, ChatGPT’s performance in temporal re-

lation extraction still falls significantly short compared to fine-tuned transformer-based models.

5.3 Discussion and Research Gaps

Homogenization of Methods and Evaluations

While numerous Transformer-based methods for temporal relation extraction have emerged, they tend to be algorithmically similar, utilizing discriminative base models like BERT to represent temporal entities and incorporating additional knowledge into these representations. A common strategy involves using off-the-shelf IE tools to extract syntactic knowledge and enhance the base model’s representations with graph neural networks. The small gains in state-of-the-art performance from one model to the next probably represent additional hyperparameter tuning more than substantial progress in understanding the relations between temporal entities in text.

Most works also focus on only three datasets – MATRES, TimeBank-Dense, and TDDiscourse – which are predominantly in the newswire domain with only 274, 36, and 34 documents, respectively, and exhibit significant overlap. This limitation in datasets might lead to an incomplete assessment of the models’ generalization capabilities. Repeated testing and fine-tuning on these small, overlapping datasets could result in overfitting, failing to reflect the models’ effectiveness on broader and more diverse datasets. Moreover, this singular domain-focused evaluation approach could cause severe domain bias, leaving the applicability of these methods outside the news domain uncertain.

Absence of Generative LLMs In temporal relation extraction, we observe a phenomenon similar to that in time expressions—there is a lack of applications using generative LLMs, which have shown excellent performance in natural language processing tasks. While there are two works that attempt to explore Transformer-based generative approaches, they are limited to studying different formats in input and output. We have not seen further exploration or application of more complex prompting techniques or training strategies.

Increased Demand for Model Openness As shown in the last column of Table 2, most temporal relation extraction models are not publicly available, possibly due to the absence of code releases or the need to re-train models on new datasets even when code is provided. Re-training a model involves significant replication work. This inaccessi-

bility directly impacts the practical application and testing of these trained models in other temporal reasoning tasks, thereby affecting the development of the temporal relation extraction field. Given the application-oriented nature of temporal relation extraction tasks, only by understanding the specific issues encountered in actual applications can we propose strategies to address these real-world challenges.

6 Applications

6.1 Methods Overview

Temporal IE is often regarded as an “upstream” system, akin to other general IE systems. These systems aim to extract structured information to improve the reasoning of “downstream” tasks, such as temporal reasoning. A natural question is how the models from Sections 4 and 5 are used in downstream tasks to help temporal reasoning.

Despite a wealth of research on Transformer-based temporal IE systems in recent years, there has been scant application of these systems’ outputs in temporal reasoning tasks. Only a few temporal reasoning tasks, such as timeline summarization and temporal question answering, leverage the results of temporal IE. The timeline summarization task aims to chronologically order and label key dates of events within a collection of news documents, while temporal question answering relies on unstructured context documents to answer temporal-related questions. Both tasks require reasoning about time and events to generate outcomes.

One approach to utilizing temporal IE systems is to explicitly construct temporal graphs to assist with temporal reasoning. Some works use only simple temporal graphs containing only time expressions extracted by rules (Su et al., 2023) or transformers (Yang et al., 2023; Xiong et al., 2024) and normalized by rules. Other works use complete temporal graphs constructed by a complete temporal IE pipeline, including time expression identification, normalization, and temporal relation extraction, with Mathur et al. (2022) using Transformer-based relation extraction, and Li et al. (2021) using LSTM-based relation extraction and rules for the other components. As for the usage of the constructed temporal graph, they can be input into models directly in text form (Su et al., 2023; Yang et al., 2023; Xiong et al., 2024) or encoded into the hidden states of a Transformer model through an attention fusion mechanism or

graph neural networks (Li et al., 2021; Mathur et al., 2022; Su et al., 2023).

Some works only preprocess the input with a specific temporal IE component rather than building a temporal graph. For instance, Bedi et al. (2021) employs the rule-based HeidelTime (Strötgen and Gertz, 2010) for extracting and normalizing time expressions in texts for constructing the input of a temporal question generation model; while Cole et al. (2023) uses the rule-based SUTime (Chang and Manning, 2012) to process the entire Wikipedia, supporting the temporal pre-training of the Transformer model.

6.2 Discussion and Research Gaps

Although there is considerable work on transformer-based temporal IE, especially in temporal relation extraction tasks, these methods have not been widely applied to downstream tasks. For example, there are many Transformer-based works that have been trained on the MATRES dataset, but none have been utilized in downstream tasks. This may be attributed to most temporal IE models not being publicly available, as shown in Table 2. Replicating these models can be both complex and time-consuming, requiring substantial effort. Furthermore, existing models exhibit domain bias. For example, in temporal relation extraction tasks, most research relies on the TimeBank-Dense and MATRES datasets, which primarily contain data from the newswire domain. Hence, the generalization capabilities of these models in other domains might be limited.

7 Conclusion

In this paper, we provide an overview of three classic tasks in the field of temporal IE: time expression identification, time expression normalization, and temporal relation extraction. We discuss datasets, Transformer-based methods, and their applications within these areas. We found that although Transformer models have demonstrated outstanding performance on many NLP tasks, there remain significant research gaps in the domain of temporal IE. For example, there is a noticeable lack of studies involving LLMs. We hope this survey will offer a comprehensive review and insights to researchers in the field, inspiring further research to address these existing gaps. We expand on the research opportunities arising from these gaps in Appendix D.

Limitations

In this review, we focus exclusively on transformer-based temporal IE methods, without including rule-based approaches. We also center our discussion on the most common temporal IE tasks rather than addressing every possible subtask.

References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altmenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Ghada Alfattni, Niels Peek, and Goran Nenadic. 2020. Extraction of temporal relations from clinical free text: A systematic review of current approaches. *Journal of biomedical informatics*, 108:103488.
- Satya Almasian, Dennis Aumiller, and Michael Gertz. 2021. Bert got a date: Introducing transformers to temporal tagging. *arXiv preprint arXiv:2109.14927*.
- Satya Almasian, Dennis Aumiller, and Michael Gertz. 2022. Time for some german? pre-training a transformer-based temporal tagger for german. *Text2Story@ ECIR*, 3117.
- Masayuki Asahara, Sachi Yasuda, Hikari Konishi, Mizuho Imada, and Kikuo Maekawa. 2013. *BCCWJ-TimeBank: Temporal and event information annotation on Japanese text*. In *Proceedings of the 27th Pacific Asia Conference on Language, Information, and Computation (PACLIC 27)*, pages 206–214, Taipei, Taiwan. Department of English, National Chengchi University.
- Miguel Ballesteros, Rishita Anubhai, Shuai Wang, Nima Pourdamghani, Yogarshi Vyas, Jie Ma, Parminder Bhatia, Kathleen McKeown, and Yaser Al-Onaizan. 2020. *Severing the edge between before and after: Neural architectures for temporal ordering of events*. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 5412–5417, Online. Association for Computational Linguistics.
- Harsimran Bedi, Sangameshwar Patil, and Girish Palshikar. 2021. *Temporal question generation from history text*. In *Proceedings of the 18th International Conference on Natural Language Processing (ICON)*, pages 408–413, National Institute of Technology Silchar, Silchar, India. NLP Association of India (NLP AI).
- Steven Bethard, Leon Derczynski, Guergana Savova, James Pustejovsky, and Marc Verhagen. 2015. *SemEval-2015 task 6: Clinical TempEval*. In *Proceedings of the 9th International Workshop on Semantic Evaluation (SemEval 2015)*, pages 806–814, Denver, Colorado. Association for Computational Linguistics.

- Steven Bethard and Jonathan Parker. 2016. *A semantically compositional annotation scheme for time normalization*. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC’16)*, pages 3779–3786, Portorož, Slovenia. European Language Resources Association (ELRA).
- Steven Bethard, Guergana Savova, Wei-Te Chen, Leon Derczynski, James Pustejovsky, and Marc Verhagen. 2016. *SemEval-2016 task 12: Clinical TempEval*. In *Proceedings of the 10th International Workshop on Semantic Evaluation (SemEval-2016)*, pages 1052–1062, San Diego, California. Association for Computational Linguistics.
- Steven Bethard, Guergana Savova, Martha Palmer, and James Pustejovsky. 2017. *SemEval-2017 task 12: Clinical TempEval*. In *Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017)*, pages 565–572, Vancouver, Canada. Association for Computational Linguistics.
- André Bittar, Pascal Amsili, Pascal Denis, and Laurence Danlos. 2011. *French TimeBank: An ISO-TimeML annotated reference corpus*. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pages 130–134, Portland, Oregon, USA. Association for Computational Linguistics.
- Alice Bracchi, Tommaso Caselli, and Irina Prodanof. 2016. Enriching the ita-timebank with narrative containers. In *Proceedings of Third Italian Conference on Computational Linguistics CLiC-it 2016*, pages 83–88. Accademia University Press.
- Pengfei Cao, Xinyu Zuo, Yubo Chen, Kang Liu, Jun Zhao, and Wei Bi. 2021. Uncertainty-aware self-training for semi-supervised event temporal relation extraction. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*, pages 2900–2904.
- Yuwei Cao, William Groves, Tanay Kumar Saha, Joel Tetreault, Alejandro Jaimes, Hao Peng, and Philip Yu. 2022. *XLTime: A cross-lingual knowledge transfer framework for temporal expression extraction*. In *Findings of the Association for Computational Linguistics: NAACL 2022*, pages 1931–1942, Seattle, United States. Association for Computational Linguistics.
- Tommaso Caselli and Piek Vossen. 2017. *The event StoryLine corpus: A new benchmark for causal and temporal relation extraction*. In *Proceedings of the Events and Stories in the News Workshop*, pages 77–86, Vancouver, Canada. Association for Computational Linguistics.
- 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.

863	Nathanael Chambers, Taylor Cassidy, Bill McDowell,	shape on visualization task performance. In <i>Proceed-</i>	919
864	and Steven Bethard. 2014. Dense event ordering	<i>ings of the 2020 CHI Conference on Human Factors</i>	920
865	with a multi-pass architecture . <i>Transactions of the</i>	<i>in Computing Systems</i> , pages 1–12.	921
866	<i>Association for Computational Linguistics</i> , 2:273–		
867	284.		
868	Angel X. Chang and Christopher Manning. 2012. SU-	Dmitriy Dligach, Steven Bethard, Timothy Miller, and	922
869	Time: A library for recognizing and normalizing	Guergana Savova. 2022. Exploring text representa-	923
870	time expressions . In <i>Proceedings of the Eighth In-</i>	tions for generative temporal relation extraction . In	924
871	<i>ternational Conference on Language Resources and</i>	<i>Proceedings of the 4th Clinical Natural Language</i>	925
872	<i>Evaluation (LREC’12)</i> , pages 3735–3740, Istanbul,	<i>Processing Workshop</i> , pages 109–113, Seattle, WA.	926
873	Turkey. European Language Resources Association	Association for Computational Linguistics.	927
874	(ELRA).		
875	Sanxing Chen, Guoxin Wang, and Börje Karlsson. 2019.	Abhimanyu Dubey, Abhinav Jauhari, Abhinav Pandey,	928
876	Exploring word representations on time expression	Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman,	929
877	recognition. <i>Microsoft Research Asia, Tech. Rep.</i>	Akhil Mathur, Alan Schelten, Amy Yang, Angela	930
		Fan, et al. 2024. The llama 3 herd of models. <i>arXiv</i>	931
		<i>preprint arXiv:2407.21783</i> .	932
878	Fei Cheng, Masayuki Asahara, Ichiro Kobayashi, and	Rujun Han, I-Hung Hsu, Mu Yang, Aram Galstyan,	933
879	Sadao Kurohashi. 2020. Dynamically updating event	Ralph Weischedel, and Nanyun Peng. 2019a. Deep	934
880	representations for temporal relation classification	structured neural network for event temporal relation	935
881	with multi-category learning . In <i>Findings of the Asso-</i>	extraction . In <i>Proceedings of the 23rd Conference on</i>	936
882	<i>ciation for Computational Linguistics: EMNLP 2020</i> ,	<i>Computational Natural Language Learning (CoNLL)</i> ,	937
883	pages 1352–1357, Online. Association for Computa-	pages 666–106, Hong Kong, China. Association for	938
884	tional Linguistics.	Computational Linguistics.	939
885	Fei Cheng and Yusuke Miyao. 2018. Inducing temporal	Rujun Han, Mengyue Liang, Bashar Alhafni, and	940
886	relations from time anchor annotation . In <i>Proceed-</i>	Nanyun Peng. 2019b. Contextualized word em-	941
887	<i>ings of the 2018 Conference of the North American</i>	beddings enhanced event temporal relation ex-	942
888	<i>Chapter of the Association for Computational Lin-</i>	traction for story understanding. <i>arXiv preprint</i>	943
889	<i>guistics: Human Language Technologies, Volume</i>	<i>arXiv:1904.11942</i> .	944
890	<i>1 (Long Papers)</i> , pages 1833–1843, New Orleans,		
891	Louisiana. Association for Computational Linguis-	Rujun Han, Qiang Ning, and Nanyun Peng. 2019c. Joint	945
892	tics.	event and temporal relation extraction with shared	946
893	Jeremy R. Cole, Aditi Chaudhary, Bhuwan Dhingra,	representations and structured prediction . In <i>Pro-</i>	947
894	and Partha Talukdar. 2023. Salient span masking	<i>ceedings of the 2019 Conference on Empirical Meth-</i>	948
895	for temporal understanding . In <i>Proceedings of the</i>	<i>ods in Natural Language Processing and the 9th In-</i>	949
896	<i>17th Conference of the European Chapter of the As-</i>	<i>ternational Joint Conference on Natural Language</i>	950
897	<i>sociation for Computational Linguistics</i> , pages 3052–	<i>Processing (EMNLP-IJCNLP)</i> , pages 434–444, Hong	951
898	3060, Dubrovnik, Croatia. Association for Computa-	Kong, China. Association for Computational Linguis-	952
899	tional Linguistics.	tics.	953
900	Francisco Costa and António Branco. 2012. Time-	Rujun Han, Xiang Ren, and Nanyun Peng. 2021. ECONET: Effective continual pretraining of lan-	954
901	BankPT: A TimeML annotated corpus of Portuguese .	guage models for event temporal reasoning . In <i>Pro-</i>	955
902	In <i>Proceedings of the Eighth International Con-</i>	<i>ceedings of the 2021 Conference on Empirical Meth-</i>	956
903	<i>ference on Language Resources and Evaluation</i>	<i>ods in Natural Language Processing</i> , pages 5367–	957
904	<i>(LREC’12)</i> , pages 3727–3734, Istanbul, Turkey. Eu-	5380, Online and Punta Cana, Dominican Republic.	958
905	ropean Language Resources Association (ELRA).	Association for Computational Linguistics.	959
906	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and	Rujun Han, Yichao Zhou, and Nanyun Peng. 2020. Do-	961
907	Kristina Toutanova. 2019. BERT: Pre-training of	main knowledge empowered structured neural net	962
908	deep bidirectional transformers for language under-	for end-to-end event temporal relation extraction . In	963
909	standing . In <i>Proceedings of the 2019 Conference of</i>	<i>Proceedings of the 2020 Conference on Empirical</i>	964
910	<i>the North American Chapter of the Association for</i>	<i>Methods in Natural Language Processing (EMNLP)</i> ,	965
911	<i>Computational Linguistics: Human Language Tech-</i>	pages 5717–5729, Online. Association for Computa-	966
912	<i>nologies, Volume 1 (Long and Short Papers)</i> , pages	tional Linguistics.	967
913	4171–4186, Minneapolis, Minnesota. Association for		
914	Computational Linguistics.	EunJeong Hwang, Jay-Yoon Lee, Tianyi Yang, Dhru-	968
915	Sara Di Bartolomeo, Aditeya Pandey, Aristotelis Leven-	vish Patel, Dongxu Zhang, and Andrew McCallum.	969
916	tidis, David Saffo, Uzma Haque Syeda, Elin Carstens-	2022. Event-event relation extraction using proba-	970
917	dottir, Magy Seif El-Nasr, Michelle A Borkin, and	bilistic box embedding . In <i>Proceedings of the 60th</i>	971
918	Cody Dunne. 2020. Evaluating the effect of timeline	<i>Annual Meeting of the Association for Computational</i>	972
		<i>Linguistics (Volume 2: Short Papers)</i> , pages 235–244,	973
		Dublin, Ireland. Association for Computational Lin-	974
		guistics.	975

976	Pustejovsky James. 2003. Timeml: Robust specification of event and temporal expressions in text. In <i>Proceedings of the Fifth International Workshop on Computational Semantics (IWCS-5)</i> , 2003.	
977		
978		
979		
980	Lis Kanashiro Pereira. 2022. Attention-focused adversarial training for robust temporal reasoning . In <i>Proceedings of the Thirteenth Language Resources and Evaluation Conference</i> , pages 7352–7359, Marseille, France. European Language Resources Association.	
981		
982		
983		
984		
985	Allen Kim, Charuta Pethe, and Steve Skiena. 2020. What time is it? temporal analysis of novels . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 9076–9086, Online. Association for Computational Linguistics.	
986		
987		
988		
989		
990		
991	Lukas Lange, Jannik Strötgen, Heike Adel, and Dietrich Klakow. 2023. Multilingual normalization of temporal expressions with masked language models . In <i>Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics</i> , pages 1174–1186, Dubrovnik, Croatia. Association for Computational Linguistics.	
992		
993		
994		
995		
996		
997		
998	Egoitz Laparra, Xin Su, Yiyun Zhao, Özlem Uzuner, Timothy Miller, and Steven Bethard. 2021. SemEval-2021 task 10: Source-free domain adaptation for semantic processing . In <i>Proceedings of the 15th International Workshop on Semantic Evaluation (SemEval-2021)</i> , pages 348–356, Online. Association for Computational Linguistics.	
999		
1000		
1001		
1002		
1003		
1004		
1005	Egoitz Laparra, Dongfang Xu, Ahmed Elsayed, Steven Bethard, and Martha Palmer. 2018. SemEval 2018 task 6: Parsing time normalizations . In <i>Proceedings of the 12th International Workshop on Semantic Evaluation</i> , pages 88–96, New Orleans, Louisiana. Association for Computational Linguistics.	
1006		
1007		
1008		
1009		
1010		
1011	Artuur Leeuwenberg and Marie-Francine Moens. 2019. A survey on temporal reasoning for temporal information extraction from text. <i>Journal of Artificial Intelligence Research</i> , 66:341–380.	
1012		
1013		
1014		
1015	Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 7871–7880, Online. Association for Computational Linguistics.	
1016		
1017		
1018		
1019		
1020		
1021		
1022		
1023		
1024	Manling Li, Tengfei Ma, Mo Yu, Lingfei Wu, Tian Gao, Heng Ji, and Kathleen McKeown. 2021. Timeline summarization based on event graph compression via time-aware optimal transport . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 6443–6456, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	
1025		
1026		
1027		
1028		
1029		
1030		
1031		
	Qian Li, Jianxin Li, Jiawei Sheng, Shiyao Cui, Jia Wu, Yiming Hei, Hao Peng, Shu Guo, Lihong Wang, Amin Beheshti, et al. 2022. A survey on deep learning event extraction: Approaches and applications. <i>IEEE Transactions on Neural Networks and Learning Systems</i> .	1032
		1033
		1034
		1035
		1036
		1037
	Chae-Gyun Lim, Young-Seob Jeong, and Ho-Jin Choi. 2018. Korean TimeBank including relative temporal information . In <i>Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)</i> , Miyazaki, Japan. European Language Resources Association (ELRA).	1038
		1039
		1040
		1041
		1042
		1043
	Chae-Gyun Lim, Young-Seob Jeong, and Ho-Jin Choi. 2019. Survey of temporal information extraction. <i>Journal of Information Processing Systems</i> , 15(4):931–956.	1044
		1045
		1046
		1047
	Chen Lin, Timothy Miller, Dmitriy Dligach, Steven Bethard, and Guergana Savova. 2019. A BERT-based universal model for both within- and cross-sentence clinical temporal relation extraction . In <i>Proceedings of the 2nd Clinical Natural Language Processing Workshop</i> , pages 65–71, Minneapolis, Minnesota, USA. Association for Computational Linguistics.	1048
		1049
		1050
		1051
		1052
		1053
		1054
		1055
	Chen Lin, Timothy Miller, Dmitriy Dligach, Farig Sadique, Steven Bethard, and Guergana Savova. 2020. A BERT-based one-pass multi-task model for clinical temporal relation extraction . In <i>Proceedings of the 19th SIGBioMed Workshop on Biomedical Language Processing</i> , pages 70–75, Online. Association for Computational Linguistics.	1056
		1057
		1058
		1059
		1060
		1061
		1062
	Jian Liu, Jinan Xu, Yufeng Chen, and Yujie Zhang. 2021. Discourse-level event temporal ordering with uncertainty-guided graph completion. In <i>IJCAI</i> , pages 3871–3877.	1063
		1064
		1065
		1066
	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. <i>arXiv preprint arXiv:1907.11692</i> .	1067
		1068
		1069
		1070
		1071
	Hieu Man, Nghia Trung Ngo, Linh Ngo Van, and Thien Huu Nguyen. 2022. Selecting optimal context sentences for event-event relation extraction. In <i>Proceedings of the AAAI conference on artificial intelligence</i> , volume 36, pages 11058–11066.	1072
		1073
		1074
		1075
		1076
	Puneet Mathur, Rajiv Jain, Franck Dernoncourt, Vlad Morariu, Quan Hung Tran, and Dinesh Manocha. 2021. TIMERS: Document-level temporal relation extraction . In <i>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)</i> , pages 524–533, Online. Association for Computational Linguistics.	1077
		1078
		1079
		1080
		1081
		1082
		1083
		1084
		1085
	Puneet Mathur, Vlad Morariu, Verena Kaynig-Fittkau, Jiuxiang Gu, Franck Dernoncourt, Quan Tran, Ani	1086
		1087

1088	Nenkova, Dinesh Manocha, and Rajiv Jain. 2022.	<i>Empirical Methods in Natural Language Processing</i>	1145
1089	DocTime: A document-level temporal dependency	<i>and the 9th International Joint Conference on Natu-</i>	1146
1090	graph parser . In <i>Proceedings of the 2022 Conference</i>	<i>ral Language Processing (EMNLP-IJCNLP)</i> , pages	1147
1091	<i>of the North American Chapter of the Association</i>	6203–6209, Hong Kong, China. Association for Com-	1148
1092	<i>for Computational Linguistics: Human Language</i>	putational Linguistics.	1149
1093	<i>Technologies</i> , pages 993–1009, Seattle, United States.		
1094	Association for Computational Linguistics.		
1095	Ulrike May, Karolina Zaczynska, Julián Moreno-	Qiang Ning, Hao Wu, and Dan Roth. 2018. A multi-	1150
1096	Schneider, and Georg Rehm. 2021. Extraction and	axis annotation scheme for event temporal relations .	1151
1097	normalization of vague time expressions in German .	In <i>Proceedings of the 56th Annual Meeting of the</i>	1152
1098	In <i>Proceedings of the 17th Conference on Natural</i>	<i>Association for Computational Linguistics (Volume</i>	1153
1099	<i>Language Processing (KONVENS 2021)</i> , pages 114–	<i>1: Long Papers)</i> , pages 1318–1328, Melbourne, Aus-	1154
1100	126, Düsseldorf, Germany. KONVENS 2021 Orga-	tralia. Association for Computational Linguistics.	1155
1101	nizers.		
1102	Pawel Mazur and Robert Dale. 2010. WikiWars: A	Tim O’Gorman, Kristin Wright-Bettner, and Martha	1156
1103	new corpus for research on temporal expressions . In	Palmer. 2016. Richer event description: Integrating	1157
1104	<i>Proceedings of the 2010 Conference on Empirical</i>	event coreference with temporal, causal and bridging	1158
1105	<i>Methods in Natural Language Processing</i> , pages 913–	annotation . In <i>Proceedings of the 2nd Workshop on</i>	1159
1106	922, Cambridge, MA. Association for Computational	<i>Computing News Storylines (CNS 2016)</i> , pages 47–	1160
1107	Linguistics.	56, Austin, Texas. Association for Computational	1161
1108	Anne-Lyse Minard, Manuela Speranza, Ruben Urizar,	Linguistics.	1162
1109	Begoña Altuna, Marieke van Erp, Anneleen Schoen,	Amy L Olex and Bridget T McInnes. 2021. Review of	1163
1110	and Chantal van Son. 2016. MEANTIME, the	temporal reasoning in the clinical domain for timeline	1164
1111	NewsReader multilingual event and time corpus . In	extraction: Where we are and where we need to be.	1165
1112	<i>Proceedings of the Tenth International Conference</i>	<i>Journal of biomedical informatics</i> , 118:103784.	1166
1113	<i>on Language Resources and Evaluation (LREC’16)</i> ,		
1114	pages 4417–4422, Portorož, Slovenia. European Lan-	Amy L Olex and Bridget T McInnes. 2022. Temporal	1167
1115	guage Resources Association (ELRA).	disambiguation of relative temporal expressions in	1168
1116	Sajad Mirzababaei, Amir Hossein Kargaran, Hinrich	clinical texts. <i>Frontiers in Research Metrics and</i>	1169
1117	Schütze, and Ehsaneddin Asgari. 2022. Hengam: An	<i>Analytics</i> , 7:1001266.	1170
1118	adversarially trained transformer for Persian temporal	Lis Pereira, Fei Cheng, Masayuki Asahara, and Ichiro	1171
1119	tagging . In <i>Proceedings of the 2nd Conference of the</i>	Kobayashi. 2021. ALICE++: Adversarial training	1172
1120	<i>Asia-Pacific Chapter of the Association for Computa-</i>	for robust and effective temporal reasoning . In <i>Pro-</i>	1173
1121	<i>tational Linguistics and the 12th International Joint</i>	<i>ceedings of the 35th Pacific Asia Conference on Lan-</i>	1174
1122	<i>Conference on Natural Language Processing (Vol-</i>	<i>guage, Information and Computation</i> , pages 373–	1175
1123	<i>ume 1: Long Papers)</i> , pages 1013–1024, Online only.	382, Shanghai, China. Association for Computational	1176
1124	Association for Computational Linguistics.	Linguistics.	1177
1125	Nasrin Mostafazadeh, Alyson Grealish, Nathanael	James Pustejovsky, Kiyong Lee, Harry Bunt, and Lau-	1178
1126	Chambers, James Allen, and Lucy Vanderwende.	rent Romary. 2010. ISO-TimeML: An international	1179
1127	2016. CaTeRS: Causal and temporal relation scheme	standard for semantic annotation . In <i>Proceedings</i>	1180
1128	for semantic annotation of event structures . In <i>Pro-</i>	<i>of the Seventh International Conference on Lan-</i>	1181
1129	<i>ceedings of the Fourth Workshop on Events</i> , pages	<i>guage Resources and Evaluation (LREC’10)</i> , Val-	1182
1130	51–61, San Diego, California. Association for Com-	letta, Malta. European Language Resources Associa-	1183
1131	putational Linguistics.	tion (ELRA).	1184
1132	Aakanksha Naik, Luke Breittfeller, and Carolyn Rose.	Alec Radford, Jeffrey Wu, Rewon Child, David Luan,	1185
1133	2019. TDDiscourse: A dataset for discourse-level	Dario Amodei, Ilya Sutskever, et al. 2019. Language	1186
1134	temporal ordering of events . In <i>Proceedings of the</i>	models are unsupervised multitask learners. <i>OpenAI</i>	1187
1135	<i>20th Annual SIGdial Meeting on Discourse and Dia-</i>	<i>blog</i> , 1(8):9.	1188
1136	<i>logue</i> , pages 239–249, Stockholm, Sweden. Associa-	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine	1189
1137	tion for Computational Linguistics.	Lee, Sharan Narang, Michael Matena, Yanqi Zhou,	1190
1138	Marta Guerrero Nieto, Roser Saurí, and Miguel An-	Wei Li, and Peter J Liu. 2020. Exploring the lim-	1191
1139	gel Bernabé Poveda. 2011. Modes timebank: A	its of transfer learning with a unified text-to-text	1192
1140	modern spanish timebank corpus. <i>Procesamiento</i>	transformer. <i>Journal of machine learning research</i> ,	1193
1141	<i>del lenguaje natural</i> , 47:259–267.	21(140):1–67.	1194
1142	Qiang Ning, Sanjay Subramanian, and Dan Roth. 2019.	Hayley Ross, Jonathon Cai, and Bonan Min. 2020. Ex-	1195
1143	An improved neural baseline for temporal relation	ploring Contextualized Neural Language Models for	1196
1144	extraction . In <i>Proceedings of the 2019 Conference on</i>	Temporal Dependency Parsing . In <i>Proceedings of the</i>	1197
		<i>2020 Conference on Empirical Methods in Natural</i>	1198
		<i>Language Processing (EMNLP)</i> , pages 8548–8553,	1199
		Online. Association for Computational Linguistics.	1200

1201	Vered Shwartz. 2022. Good night at 4 pm?! time expressions in different cultures . In <i>Findings of the Association for Computational Linguistics: ACL 2022</i> , pages 2842–2853, Dublin, Ireland. Association for Computational Linguistics.	1258
1202		1259
1203		1260
1204		1261
1205		1262
1206	Jannik Strötgen, Thomas Bögel, Julian Zell, Ayser Armiti, Tran Van Canh, and Michael Gertz. 2014. Extending HeidelTime for temporal expressions referring to historic dates . In <i>Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)</i> , pages 2390–2397, Reykjavik, Iceland. European Language Resources Association (ELRA).	1263
1207		1264
1208		
1209		1265
1210		1266
1211		1267
1212		1268
1213		1269
1214	Jannik Strötgen and Michael Gertz. 2010. HeidelTime: High quality rule-based extraction and normalization of temporal expressions . In <i>Proceedings of the 5th International Workshop on Semantic Evaluation</i> , pages 321–324, Uppsala, Sweden. Association for Computational Linguistics.	1270
1215		
1216		1271
1217		1272
1218		1273
1219		1274
1220	Jannik Strötgen, Anne-Lyse Minard, Lukas Lange, Manuela Speranza, and Bernardo Magnini. 2018. KRAUTS: A German temporally annotated news corpus . In <i>Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)</i> , Miyazaki, Japan. European Language Resources Association (ELRA).	1275
1221		1276
1222		1277
1223		1278
1224		1279
1225		1280
1226		
1227	William F. Styler IV, Steven Bethard, Sean Finan, Martha Palmer, Sameer Pradhan, Piet C de Groen, Brad Erickson, Timothy Miller, Chen Lin, Guergana Savova, and James Pustejovsky. 2014. Temporal annotation in the clinical domain . <i>Transactions of the Association for Computational Linguistics</i> , 2:143–154.	1281
1228		1282
1229		1283
1230		1284
1231		1285
1232		
1233		1286
1234	Xin Su, Phillip Howard, Nagib Hakim, and Steven Bethard. 2023. Fusing temporal graphs into transformers for time-sensitive question answering . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 948–966, Singapore. Association for Computational Linguistics.	1287
1235		1288
1236		1289
1237		1290
1238		1291
1239		1292
1240	Xin Su, Yiyun Zhao, and Steven Bethard. 2021. The University of Arizona at SemEval-2021 task 10: Applying self-training, active learning and data augmentation to source-free domain adaptation . In <i>Proceedings of the 15th International Workshop on Semantic Evaluation (SemEval-2021)</i> , pages 458–466, Online. Association for Computational Linguistics.	1293
1241		1294
1242		1295
1243		1296
1244		1297
1245		1298
1246		
1247	Weiwei Sun, Anna Rumshisky, and Ozlem Uzuner. 2013. Evaluating temporal relations in clinical text: 2012 i2b2 challenge. <i>Journal of the American Medical Informatics Association</i> , 20(5):806–813.	1299
1248		1300
1249		1301
1250		1302
1251	Xingwei Tan, Gabriele Pergola, and Yulan He. 2021. Extracting event temporal relations via hyperbolic geometry . 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.	1303
1252		1304
1253		1305
1254		1306
1255		1307
1256		1308
1257		1309
		1310
		1311
		1312
		1313
		1314
	Xingwei Tan, Gabriele Pergola, and Yulan He. 2023. Event temporal relation extraction with Bayesian translational model . 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.	
	Sun Tiesen and Li Lishuang. 2022. Improving event temporal relation classification via auxiliary label-aware contrastive learning . In <i>Proceedings of the 21st Chinese National Conference on Computational Linguistics</i> , pages 861–871, Nanchang, China. Chinese Information Processing Society of China.	
	Naushad UzZaman, Hector Llorens, Leon Derczynski, James Allen, Marc Verhagen, and James Pustejovsky. 2013. SemEval-2013 task 1: TempEval-3: Evaluating time expressions, events, and temporal relations . 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.	
	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. <i>Advances in neural information processing systems</i> , 30.	
	Marc Verhagen, Robert Gaizauskas, Frank Schilder, Mark Hepple, Graham Katz, and James Pustejovsky. 2007. SemEval-2007 task 15: TempEval temporal relation identification . In <i>Proceedings of the Fourth International Workshop on Semantic Evaluations (SemEval-2007)</i> , pages 75–80, Prague, Czech Republic. Association for Computational Linguistics.	
	Marc Verhagen, Roser Saurí, Tommaso Caselli, and James Pustejovsky. 2010. SemEval-2010 task 13: TempEval-2 . In <i>Proceedings of the 5th International Workshop on Semantic Evaluation</i> , pages 57–62, Uppsala, Sweden. Association for Computational Linguistics.	
	Haoyu Wang, Muhao Chen, Hongming Zhang, and Dan Roth. 2020. Joint constrained learning for event-event relation extraction . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 696–706, Online. Association for Computational Linguistics.	
	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 <i>Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics</i> , pages 541–553, Dubrovnik, Croatia. Association for Computational Linguistics.	
	Liang Wang, Peifeng Li, and Sheng Xu. 2022. DCT-centered temporal relation extraction . In <i>Proceed-</i>	

1315
1316
1317
1318

1319
1320
1321
1322
1323

1324
1325
1326
1327
1328
1329
1330

1331
1332
1333
1334

1335
1336
1337
1338
1339
1340
1341

1342
1343
1344
1345
1346
1347

1348
1349
1350
1351
1352
1353

1354
1355
1356
1357
1358
1359

1360
1361
1362
1363
1364
1365

1366
1367
1368
1369
1370
1371

ings of the 29th International Conference on Computational Linguistics, pages 2087–2097, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837.

Haoyang Wen and Heng Ji. 2021. Utilizing relative event time to enhance event-event temporal relation extraction. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 10431–10437, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Siheng Xiong, Ali Payani, Ramana Kompella, and Faramarz Fekri. 2024. Large language models can learn temporal reasoning. *arXiv preprint arXiv:2401.06853*.

Sen Yang, Xin Li, Lidong Bing, and Wai Lam. 2023. Once upon a time in graph: Relative-time pretraining for complex temporal reasoning. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 11879–11895, Singapore. Association for Computational Linguistics.

Jiarui Yao, Haoling Qiu, Bonan Min, and Nianwen Xue. 2020. Annotating Temporal Dependency Graphs via Crowdsourcing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 5368–5380, Online. Association for Computational Linguistics.

Hongbin Ye, Ningyu Zhang, Hui Chen, and Huajun Chen. 2022. Generative knowledge graph construction: A review. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 1–17, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

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.

Alessandra Zarcone, Touhidul Alam, and Zahra Kolagar. 2020. PATE: A corpus of temporal expressions for the in-car voice assistant domain. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 523–530, Marseille, France. European Language Resources Association.

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.

Yuchen Zhang and Nianwen Xue. 2018. Structured interpretation of temporal relations. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, Miyazaki, Japan. European Language Resources Association (ELRA).

Yuchen Zhang and Nianwen Xue. 2019. Acquiring structured temporal representation via crowdsourcing: A feasibility study. In *Proceedings of the Eighth Joint Conference on Lexical and Computational Semantics (*SEM 2019)*, pages 178–185, Minneapolis, Minnesota. Association for Computational Linguistics.

Xinyu Zhao, Shih-Ting Lin, and Greg Durrett. 2021. Effective distant supervision for temporal relation extraction. In *Proceedings of the Second Workshop on Domain Adaptation for NLP*, pages 195–203, Kyiv, Ukraine. 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 Evaluation Metrics

In temporal IE, the evaluation method from TEMPEVAL-3 (UzZaman et al., 2013) is the most widely adopted standard. This evaluation method calculates the standard precision (P), recall (R), and F1 score (F) between the system predictions (System) and the gold annotations (Reference) as follows:

$$P = \frac{|\text{System} \cap \text{Reference}|}{|\text{System}|} \quad (1)$$

$$R = \frac{|\text{System} \cap \text{Reference}|}{|\text{Reference}|} \quad (2)$$

$$F = 2 \cdot \frac{P \cdot R}{P + R} \quad (3)$$

In time expression identification, “System” refers to the time expressions identified by the system, while “Reference” refers to the annotated gold time expressions. In time expression normalization, “System” and “Reference” refer to the system-normalized time expressions and the gold

annotated normalized expressions, respectively. If calculating the end-to-end time expression normalization score, “System” only involves the correctly identified time expressions.

For the temporal relation extraction task, the TEMPEVAL-3 evaluation method calculates the temporal awareness scores. This is achieved by performing a graph closure operation on the gold temporal graph based on temporal transitivity rules (to incorporate all potential temporal relations) and reducing the predicted temporal relation graph (to remove duplicate relations). These steps are completed before calculating the standard scores. Here, “System” denotes the temporal relations predicted by the system, while “Reference” is the gold annotated temporal relations.

B Datasets Summary

We summarize the temporal IE datasets in Table 1. The first section is based on the most widely used TimeML annotation framework, while the second section covers those that adopt all other annotation frameworks.

C Temporal Relation Extraction Methods Summary

We summarize the temporal relation extraction methods we review in Table 2.

D Discussion on Future Directions

In the previous sections, we have identified the following research opportunities in the field of temporal IE:

- Enrich annotation frameworks (Section 3.3), e.g., representing event arguments or expanding formal semantic systems like SCATE.
- Improve dataset diversity (Section 3.3), e.g., annotating more domains beyond newswire.
- Explore generative approaches (Sections 4.2 and 5.3), e.g., new input-output formulations, new fine-tuning strategies.
- Develop public tools and benchmarks (Sections 4.2 and 5.3), e.g., publish temporal IE models and datasets to the public repositories
- Explore new applications (Section 6.2), e.g., the utility of extracted timelines when visualized for human-computer interaction.

D.1 Enrich Annotation Frameworks and Improve the Domain Diversity of Datasets

Current annotation frameworks, such as TimeML, often produce temporal graphs composed of temporal relations and temporal entities, as illustrated in Figure 2. However, these temporal graphs are challenging to interpret independently or use directly for temporal reasoning without extensive context. One future direction could be to integrate richer content into end-to-end temporal IE annotation frameworks. One example is incorporating entity relation extraction and full event extraction (including triggers and arguments) from the general domain to construct a more complete temporal graph. This concept has begun to emerge in the literature, as seen in Li et al. (2021). Yet, that work mainly integrates existing temporal IE tools with general domain IE tools without proposing a well-defined annotation framework. Another example is to develop user-friendly frameworks like SCATE, which, unlike TimeML, outputs temporal intervals that can be directly mapped onto a timeline given a temporal expression. However, SCATE primarily focuses on the normalization of time expressions. Expanding its scope to include the normalization of a broader range of temporal content, such as events and sentences, could significantly widen its applicability.

Furthermore, future efforts could focus on expanding the domains covered by existing datasets to mitigate the domain bias present in current datasets. For example, the Thyme datasets represent an adaptation of TimeML to better suit the medical field’s representation of temporal relations between events and times. Yet, such efforts to adapt and improve annotation frameworks for additional fields are still scarce. Therefore, adapting existing annotation frameworks to a broader range of domains to enhance the domain diversity of datasets represents a potential future research direction.

D.2 Improve the Application of Generative LLMs

The application of generative LLMs in the field of time expression identification, normalization, and temporal relation extraction remains underexplored. Given the proven capabilities of LLMs like ChatGPT and LLAMA3 across various tasks, it is logical to probe their potential within the realm of temporal IE. Whether it involves leveraging new prompting methods or fine-tuning strategies for

Name	Framework	Domain	Lang	Tasks
<i>TimeML-Based</i>				
TimeBank (James, 2003)	TimeML	Newswire	EN	I, N, R
TempEval-1 (Verhagen et al., 2007)	TimeML	Newswire	EN	I, N, R
TempEval-2 (Verhagen et al., 2010)	TimeML	Newswire	ZH, EN, IT, FR, KR, ES	I, N, R
Spanish TimeBank (Nieto et al., 2011)	TimeML	Historiography	ES	I, N
French TimeBank (Bittar et al., 2011)	ISO-TimeML	Newswire	FR	I, N, R
Portuguese TimeBank (Costa and Branco, 2012)	TimeML	Newswire	PT	I, N, R
i2b2-2012 (Sun et al., 2013)	Thyme-TimeML	Clinical	EN	I, N, R
TempEval-3 (UzZaman et al., 2013)	TimeML	Newswire	EN, ES	I, N, R
TimeBank-Dense (Chambers et al., 2014)	TimeML	Newswire	EN	I, N, R
Japanese TimeBank (Asahara et al., 2013)	ISO-TimeML	Publication, Library, Special purpose	JA	I, N, R
AncientTimes (Strötgen et al., 2014)	TimeML	Wikipedia	EN, DE, NL, ES, FR, IT, AR, VI	I, N
THYME-2015 (Bethard et al., 2015)	Thyme-TimeML	Clinical	EN	I, N, R
THYME-2016 (Bethard et al., 2016)	Thyme-TimeML	Clinical	EN	I, N, R
Richer Event Description (O’Gorman et al., 2016)	Thyme-TimeML	Newswire, Forum Discussions	EN	I, N, R
Italian TimeBank (Bracchi et al., 2016)	TimeML	Newswire	IT	I, N, R
MeanTime (Minard et al., 2016)	ISO-TimeML	Newswire	EN, IT, ES, NL	I, N, R
THYME-2017 (Bethard et al., 2017)	Thyme-TimeML	Clinical	EN	I, N, R
Event StoryLine (Caselli and Vossen, 2017)	TimeML	Story	EN	I, N, R
MATRES (Ning et al., 2018)	TimeML	Newswire	EN	I, R
Korean TimeBank (Lim et al., 2018)	TimeML	Wikipedia	KR	I, N, R
German Temporal Expression (Strötgen et al., 2018)	TimeML	Newswire	DE	I, N
TDDiscourse (Naik et al., 2019)	TimeML	Newswire	EN	R
PATE (Zarcone et al., 2020)	TimeML	Voice Assistant	EN	I, N
German VTEs (May et al., 2021)	ISO-TimeML	Newswire	DE	I, N
<i>Other Annotation Framework-based</i>				
WikiWars (Mazur and Dale, 2010)	TIMEX2	Wikipedia	EN, DE	I, N
SCATE (Bethard and Parker, 2016; Laparra et al., 2018)	SCATE	Newswire, Clinical	EN	I, N
CaTeRS (Mostafazadeh et al., 2016)	CaTeRS	Commonsense Stories	EN	R
TORDER (Cheng and Miyao, 2018)	TORDER	Newswire	EN	R
Temporal Dependency Tree (Zhang and Xue, 2018, 2019)	Temporal Dependency Tree	Newswire, Narratives	ZH	R
Temporal Dependency Graph (Yao et al., 2020)	Temporal Dependency Graph	Newswire	EN	R

Table 1: Overview of datasets and their schemas, domains, languages (EN: English, DE: German, NL: Dutch, ES: Spanish, FR: French, IT: Italian, AR: Arabic, VI: Vietnamese, JA: Japanese, PT: Portuguese, ZH: Chinese, KR: Korean), and tasks (I: identification, N: time expression normalization, R: temporal relation extraction).

specific tasks, there is ample room for innovation.

However, it is important to emphasize that while these models excel in generating unstructured text when applied to temporal IE, it is imperative to specially design suitable input-output formats. Such designs are intended to enable generative LLMs, which are typically used for producing unstructured text, to also effectively output structured temporal information.

D.3 Develop Public Toolkits and Evaluation Benchmarks

We believe that one key reason transformer-based temporal IE models have not been widely adopted

might be the absence of a publicly available code repository that facilitates easier access to models and data. For example, HuggingFace¹ provides language model heads or pipelines suitable for various tasks, allowing users to easily download and deploy trained models on any dataset directly from the HuggingFace Hub. A future research direction should involve establishing such a repository or pushing models/datasets to HuggingFace Hub for the temporal IE tasks to enhance the reproducibility and applicability of research. Another important direction is to create a public and test-set concealed

¹<https://huggingface.co/>

Work	Approach	Base Model	Evaluation Datasets	Knowl.	Robust	Avail.
Lin et al. (2019)	Discr.	BERT	THYME	✗	✓	✗
Han et al. (2019a)	Discr.	BERT	TimeBank-Dense, MATRES	✓	✗	✗
Ning et al. (2019)	Discr.	BERT	TimeBank-Dense, MATRES	✓	✗	✗
Han et al. (2019c)	Discr.	BERT	TimeBank-Dense, MATRES	✓	✓	✗
Han et al. (2019b)	Discr.	BERT	Richer Event Description, CaTeRS	✓	✓	✗
Lin et al. (2020)	Discr.	BERT	THYME	✗	✓	✗
Cheng et al. (2020)	Discr.	BERT	Japanese-Timebank, TimeBank-Dense	✓	✓	✗
Ross et al. (2020)	Discr.	BERT	Temporal Dependency Tree	✓	✗	✗
Ballesteros et al. (2020)	Discr.	RoBERTa	MATRES	✗	✓	✗
Han et al. (2020)	Discr.	RoBERTa	i2b2-2012, TimeBank-Dense	✓	✓	✗
Wang et al. (2020)	Discr.	RoBERTa	MATRES	✓	✗	✗
Zhao et al. (2021)	Discr.	RoBERTa	MATRES	✗	✓	✓
Zhou et al. (2021)	Discr.	BERT	i2b2-2012, TimeBank-Dense	✓	✗	✗
Cao et al. (2021)	Discr.	RoBERTa	MATRES, TimeBank-Dense	✗	✓	✗
Tan et al. (2021)	Discr.	RoBERTa	MATRES	✓	✗	✗
Mathur et al. (2021)	Discr.	BERT	TimeBank-Dense, MATRES, TDDiscourse	✓	✗	✗
Liu et al. (2021)	Discr.	BERT	TimeBank-Dense, TDDiscourse	✓	✗	✗
Wen and Ji (2021)	Discr.	RoBERTa	MATRES	✓	✗	✗
Pereira et al. (2021)	Discr.	RoBERTa	MATRES, TimeML	✗	✓	✗
Han et al. (2021)	Discr.	RoBERTa/BERT	TimeBank-Dense, MATRES, Richer Event Description	✗	✓	✓
Kanashiro Pereira (2022)	Discr.	RoBERTa	MATRES, TimeML	✗	✓	✗
Wang et al. (2022)	Discr.	RoBERTa	TimeBank-Dense, TDDiscourse	✓	✓	✗
Mathur et al. (2022)	Discr.	BERT	Temporal Dependency Tree	✓	✓	✗
Hwang et al. (2022)	Discr.	RoBERTa	MATRES, Event StoryLine	✓	✗	✗
Dligach et al. (2022)	Gen	BART/T5	THYME	✗	✗	✗
Wang et al. (2023)	Discr.	BigBird	MATRES, TDDiscourse	✓	✓	✗
Zhang et al. (2022)	Discr.	BERT	MATRES, TimeBank-Dense	✓	✗	✗
Tiesen and Lishuang (2022)	Discr.	BERT	TimeBank-Dense, MATRES	✗	✓	✗
Zhou et al. (2022)	Discr.	RoBERTa	TimeBank-Dense, MATRES	✓	✗	✗
Man et al. (2022)	Discr.	RoBERTa	MATRES, TDDiscourse	✓	✗	✗
Yuan et al. (2023)	Gen	ChatGPT	TimeBank-Dense, MATRES, TDDiscourse	✗	✗	✗
Tan et al. (2023)	Discr.	BART	MATRES, imeBank-Dense	✓	✗	✓

Table 2: Overview of research on temporal relation extraction. “Knowl.” represents the inclusion of external knowledge. “Robust” refers to the application of methods to enhance model robustness. “Avail.” indicates whether the model is publicly available. Symbols ✓ and ✗ indicate the presence or absence of a feature, respectively.

benchmark for a more equitable comparison of existing work. In most existing works, although metrics such as F1 scores, precision, and recall are commonly computed, the specific implementations can vary. For instance, in Kanashiro Pereira (2022), only the “before” and “after” relationships are evaluated for relation extraction performance, whereas Zhang et al. (2022) includes all temporal relationships except “vague” in their evaluation.

D.4 Explore More Application Directions

In reviewing the application of temporal IE systems, we observe that current research primarily

focuses on aiding “models” in temporal reasoning to enhance their performance in other tasks. Future research in temporal IE should not only continue to support model performance improvement but should also pay more attention to serving humans and enhancing its practical value. A promising application direction is visualizing timelines in human-computer interaction (HCI) scenarios. The visualization results of existing temporal graphs are often challenging for human users to interpret. For instance, visualizing the temporal graph of any document in the TimeBank-Dense dataset might result in a graph densely populated with points and

lines, offering little help for users to comprehend the progression of events within the text.

User studies, such as those conducted by [Di Bartolomeo et al. \(2020\)](#), have revealed the importance of visualization forms of timelines for user understanding. Consequently, temporal IE research should also consider incorporating user research on temporal graphs to guide the design of temporal IE methods, such as how to represent standardized time expressions, identify which types of temporal relations most effectively facilitate time understanding, and determine the best ways to present this information. By addressing these problems, the extraction and representation of temporal information can be more closely aligned with user needs, enhancing its application value in HCI.