

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 2 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

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 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 (Pustejovsky, 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.

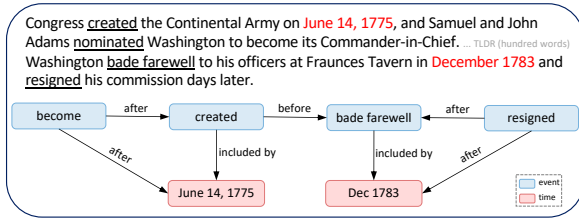


Figure 1: A snippet from George Washington’s Wikipedia page and the corresponding temporal graph.

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 1 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 B 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” (Pustejovsky, 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 1, 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 C.

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 (Pustejovsky, 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 (Pustejovsky, 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, Italian, 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 ex-

pression 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. We present additional timeline focused datasets at Appendix D.

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 1, 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 understand-

ing 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.

369 Shwartz (2022) and Kim et al. (2020) focus
370 on the normalization of time expressions and use
371 transformer-based models. Shwartz (2022) aims to
372 normalize time expressions from various cultural
373 contexts (e.g., morning, noon, afternoon) into pre-
374 cise hourly representations within a day. They train
375 a BERT model with a masked language modeling
376 task to predict specific times of day that are masked,
377 given the time expressions. Kim et al. (2020) seeks
378 to normalize time expressions in novels into spe-
379 cific daily hours, fine-tuning the BERT model for
380 a 24-class classification task to ascertain the corre-
381 sponding times of day for given expressions.

382 Lange et al. (2023) addresses both extraction
383 and normalization of time expressions, adopt-
384 ing a pipeline approach. Initially, they fine-tune
385 the XLM-R model using the token classification
386 method to extract time expressions, then denote
387 identified expressions with TIMEX3 tags with
388 masked time values, and finally fine-tune the XLM-
389 R model with masked language modeling to predict
390 the normalized masked time values.

391 Several of the aforementioned works also uti-
392 lize data augmentation techniques to improve the
393 model’s multilingual performance (Lange et al.,
394 2023; Mirzababaei et al., 2022; Almasian et al.,
395 2022). For instance, Lange et al. (2023) employs
396 the rule-based HeidelTime method (Strötgen and
397 Gertz, 2010) to annotate time expressions and their
398 normalizations across 87 languages, generating a
399 semi-supervised dataset to facilitate model training.

400 4.2 Discussion and Research Gaps

401 Despite the significant achievements of Trans-
402 former models in various NLP tasks, research in
403 the area of time expression identification and nor-
404 malization has remained relatively limited over the
405 past few years. This is particularly true of time nor-
406 malization, where the volume and depth of research
407 are low, especially when compared to similar tasks
408 such as named entity recognition, entity normaliza-
409 tion, and entity linking. Furthermore, the method-
410 ological diversity in existing works is notably con-
411 strained, with most research relying on pre-trained
412 Transformer models for simple token classification.
413 While generative LLMs like GPT-4 or LLAMA3
414 have demonstrated impressive performance in other
415 NLP tasks, their potential in the identification and
416 normalization of time expressions has barely been
417 explored. This suggests a significant research gap
418 exists; exploration of generative approaches may

offer the potential for advancement in time expres-
sion 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 E
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 representa-
tions 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 tempo-
ral 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 en-
capsulates 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 extrac-
tion 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 typi-
cally 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; Yao et al., 2024a), 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). Recent works have incorporated these rules during both training and inference. During training, models employ various approaches including box embedding (Hwang et al., 2022), hyperbolic embedding (Tan et al., 2021), loss function regularization (Zhou et al., 2021; Wang et al., 2020), contrastive objectives (Niu et al., 2024), logical expressions over event time points (Huang et al., 2023), and hierarchical logical conditions (Ning et al., 2024). For inference, methods include custom heuristics (Wang et al., 2022; Zhou et al., 2022, 2021; Liu et al., 2021), linear programming formulation (Wang et al., 2020; Han et al., 2019c), and structured prediction with support vector machines (Han et al., 2019a).

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 subse-

quently 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

Generative approaches in Temporal IE fall into two main categories: fine-tuned encoder-decoder models and large language model (LLM) prompting methods. For fine-tuned generative models, Dligach et al. (2022) investigate BART (Lewis et al., 2020) and T5 (Raffel et al., 2020) architectures, finding that producing outputs for each temporal entity pair separately outperforms triplet format (entity, relation, entity). Recent work has also explored LLM-based approaches. Yuan et al. (2023) and Huang et al. (2023) examine various prompting strategies, with Huang et al. (2023) demonstrating that structured, logic-informed prompts significantly improve performance over standard prompting. Hu et al. (2025) formulates temporal relation extraction as a question-answering task with rationale generation that includes coreference and transitive chains. Meanwhile, Niu et al. (2024) integrates LLMs specifically to enhance commonsense reasoning in their hybrid system. Despite these advances, current findings indicate that prompting-only approaches still underperform compared to fine-tuned discriminative 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.

Generative LLMs: Progress and Challenges

Despite increasing interest in generative LLMs for temporal relation extraction, a significant research gap remains: current generative approaches consistently underperform compared to fine-tuned discriminative models (Yuan et al., 2023). Although recent works have explored structured prompts (Huang et al., 2023), question-answering frameworks (Hu et al., 2025), and hybrid systems (Niu et al., 2024), none have matched state-of-the-art discriminative methods. Promising directions for future research include: (1) specialized temporal fine-tuning techniques for LLMs; (2) more effective methods to encode temporal rules and constraints in LLM prompts; and (3) improved evaluation frameworks for generative outputs in temporal tasks.

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 inaccessibility 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

665	the models from Sections 4 and 5 are used in down-	(2021) employs the rule-based HeidelTime (Ströt-	716
666	stream tasks to help temporal reasoning.	gen and Gertz, 2010) for extracting and normaliz-	717
667	Despite a wealth of research on Transformer-	ing time expressions in texts for constructing the in-	718
668	based temporal IE systems in recent years, there	put of a temporal question generation model; while	719
669	has been scant application of these systems' out-	Cole et al. (2023) uses the rule-based SUTime	720
670	puts in temporal reasoning tasks. Only a few tem-	(Chang and Manning, 2012) to process the entire	721
671	poral reasoning tasks, such as timeline extraction,	Wikipedia, supporting the temporal pre-training of	722
672	timeline summarization and temporal question an-	the Transformer model.	723
673	swering, leverage the results of temporal IE. Time-		
674	line extraction is a direct product of temporal IE,	6.2 Discussion and Research Gaps	724
675	where the extracted events and time expressions,	Although there is considerable work on	725
676	along with their temporal relations, naturally form	transformer-based temporal IE, especially in	726
677	a chronologically ordered timeline following the	temporal relation extraction tasks, these methods	727
678	traditional TimeML paradigm. For example, the	have not been widely applied to downstream tasks.	728
679	recent Chemotherapy Timeline Extraction shared	For example, there are many Transformer-based	729
680	task (Yao et al., 2024b) focuses on constructing	works that have been trained on the MATRES	730
681	patient-level treatment timelines from electronic	dataset, but none have been utilized in downstream	731
682	health records, with most participating systems us-	tasks. This may be attributed to most temporal	732
683	ing fine-tuned Transformer models for event and	IE models not being publicly available, as shown	733
684	time expression extraction, followed by temporal	in Table 2. Replicating these models can be	734
685	relation classification. The timeline summariza-	both complex and time-consuming, requiring	735
686	tion task aims to chronologically order and label	substantial effort. Furthermore, existing models	736
687	key dates of events within a collection of news	exhibit domain bias. For example, in temporal	737
688	documents, while temporal question answering re-	relation extraction tasks, most research relies	738
689	lies on unstructured context documents to answer	on the TimeBank-Dense and MATRES datasets,	739
690	temporal-related questions. Both tasks require re-	which primarily contain data from the newswire	740
691	asoning about time and events to generate outcomes.	domain. Hence, the generalization capabilities of	741
692	One approach to utilizing temporal IE systems	these models in other domains might be limited.	742
693	is to explicitly construct temporal graphs to assist		
694	with temporal reasoning. Some works use only	7 Conclusion	743
695	simple temporal graphs containing only time ex-	In this paper, we provide an overview of three clas-	744
696	pressions extracted by rules (Su et al., 2023) or	sic tasks in the field of temporal IE: time expression	745
697	transformers (Yang et al., 2023; Xiong et al., 2024)	identification, time expression normalization, and	746
698	and normalized by rules. Other works use com-	temporal relation extraction. We discuss datasets,	747
699	plete temporal graphs constructed by a complete	Transformer-based methods, and their applications	748
700	temporal IE pipeline, including time expression	within these areas. We found that although Trans-	749
701	identification, normalization, and temporal re-	former models have demonstrated outstanding per-	750
702	lation extraction, with Mathur et al. (2022) using	formance on many NLP tasks, there remain sig-	751
703	Transformer-based relation extraction, and Li et al.	nificant research gaps in the domain of temporal	752
704	(2021) using LSTM-based relation extraction and	IE. We hope this survey will offer a comprehensive	753
705	rules for the other components. As for the usage	review and insights to researchers in the field, in-	754
706	of the constructed temporal graph, they can be in-	spiring further research to address these existing	755
707	put into models directly in text form (Su et al.,	gaps. We expand on the research opportunities	756
708	2023; Yang et al., 2023; Xiong et al., 2024) or	arising from these gaps in Appendix F.	757
709	encoded into the hidden states of a Transformer		
710	model through an attention fusion mechanism or	Limitations	758
711	graph neural networks (Li et al., 2021; Mathur et al.,	In this review, we focus exclusively on transformer-	759
712	2022; Su et al., 2023).	based temporal IE methods, without including rule-	760
713	Some works only preprocess the input with a	based approaches. We also center our discussion	761
714	specific temporal IE component rather than build-	on the most common temporal IE tasks rather than	762
715	ing a temporal graph. For instance, Bedi et al.	addressing every possible subtask.	763

764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819

References

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, 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.

Nathanael Chambers, Taylor Cassidy, Bill McDowell, and Steven Bethard. 2014. **Dense event ordering with a multi-pass architecture**. *Transactions of the Association for Computational Linguistics*, 2:273–284.

875	Angel X. Chang and Christopher Manning. 2012. SU-Time: A library for recognizing and normalizing time expressions . In <i>Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC'12)</i> , pages 3735–3740, Istanbul, Turkey. European Language Resources Association (ELRA).	932
876		933
877		934
878		
879		935
880		936
881		937
882		938
883		939
884	Sanxing Chen, Guoxin Wang, and Börje Karlsson. 2019. Exploring word representations on time expression recognition. <i>Microsoft Research Asia, Tech. Rep.</i>	
885		940
886		941
887		942
888		943
889		944
890		945
891		946
892		947
893		948
894		949
895		950
896		951
897		
898		
899		
900		952
901		953
902		954
903		955
904		956
905		957
906		958
907		959
908		960
909		
910		
911		
912		
913		
914		
915		
916		
917		
918		
919		
920		
921		
922		
923		
924		
925		
926		
927		
928		
929		
930		
931		
		961
		962
		963
		964
		965
		966
		967
		968
		969
		970
		971
		972
		973
		974
		975
		976
		977
		978
		979
		980
		981
		982
		983
		984
		985
		986
		987
		988

1102	Puneet Mathur, Vlad Morariu, Verena Kaynig-Fittkau,	extraction . In <i>Proceedings of the 2019 Conference on</i>	1160
1103	Jiuxiang Gu, Franck Dernoncourt, Quan Tran, Ani	<i>Empirical Methods in Natural Language Processing</i>	1161
1104	Nenkova, Dinesh Manocha, and Rajiv Jain. 2022.	<i>and the 9th International Joint Conference on Natu-</i>	1162
1105	DocTime: A document-level temporal dependency	<i>ral Language Processing (EMNLP-IJCNLP)</i> , pages	1163
1106	graph parser . In <i>Proceedings of the 2022 Conference</i>	6203–6209, Hong Kong, China. Association for Com-	1164
1107	<i>of the North American Chapter of the Association</i>	putational Linguistics.	1165
1108	<i>for Computational Linguistics: Human Language</i>		
1109	<i>Technologies</i> , pages 993–1009, Seattle, United States.		
1110	Association for Computational Linguistics.		
1111	Ulrike May, Karolina Zaczynska, Julián Moreno-	Qiang Ning, Hao Wu, and Dan Roth. 2018. A multi-	1166
1112	Schneider, and Georg Rehm. 2021. Extraction and	axis annotation scheme for event temporal relations .	1167
1113	normalization of vague time expressions in German .	In <i>Proceedings of the 56th Annual Meeting of the</i>	1168
1114	In <i>Proceedings of the 17th Conference on Natural</i>	<i>Association for Computational Linguistics (Volume</i>	1169
1115	<i>Language Processing (KONVENS 2021)</i> , pages 114–	<i>1: Long Papers)</i> , pages 1318–1328, Melbourne, Aus-	1170
1116	126, Düsseldorf, Germany. KONVENS 2021 Orga-	tralia. Association for Computational Linguistics.	1171
1117	nizers.		
1118	Pawel Mazur and Robert Dale. 2010. WikiWars: A	Wanting Ning, Lishuang Li, Xueyang Qin, Yubo Feng,	1172
1119	new corpus for research on temporal expressions . In	and Jingyao Tang. 2024. Temporal cognitive tree:	1173
1120	<i>Proceedings of the 2010 Conference on Empirical</i>	A hierarchical modeling approach for event tempo-	1174
1121	<i>Methods in Natural Language Processing</i> , pages 913–	ral relation extraction . In <i>Findings of the Associa-</i>	1175
1122	922, Cambridge, MA. Association for Computational	<i>tion for Computational Linguistics: EMNLP 2024</i> ,	1176
1123	Linguistics.	pages 855–864, Miami, Florida, USA. Association	1177
1124	Anne-Lyse Minard, Manuela Speranza, Ruben Urizar,	for Computational Linguistics.	1178
1125	Begoña Altuna, Marieke van Erp, Anneleen Schoen,	Jingcheng Niu, Saifei Liao, Victoria Ng, Simon De Mon-	1179
1126	and Chantal van Son. 2016. MEANTIME, the	tigny, and Gerald Penn. 2024. ConTempo: A unified	1180
1127	NewsReader multilingual event and time corpus . In	temporally contrastive framework for temporal rela-	1181
1128	<i>Proceedings of the Tenth International Conference</i>	tion extraction . In <i>Findings of the Association for</i>	1182
1129	<i>on Language Resources and Evaluation (LREC’16)</i> ,	<i>Computational Linguistics: ACL 2024</i> , pages 1521–	1183
1130	pages 4417–4422, Portorož, Slovenia. European Lan-	1533, Bangkok, Thailand. Association for Computa-	1184
1131	guage Resources Association (ELRA).	tional Linguistics.	1185
1132	Sajjad Mirzababaei, Amir Hossein Kargaran, Hinrich	Tim O’Gorman, Kristin Wright-Bettner, and Martha	1186
1133	Schütze, and Ehsaneddin Asgari. 2022. Hengam: An	Palmer. 2016. Richer event description: Integrating	1187
1134	adversarially trained transformer for Persian temporal	event coreference with temporal, causal and bridging	1188
1135	tagging . In <i>Proceedings of the 2nd Conference of the</i>	annotation . In <i>Proceedings of the 2nd Workshop on</i>	1189
1136	<i>Asia-Pacific Chapter of the Association for Computa-</i>	<i>Computing News Storylines (CNS 2016)</i> , pages 47–	1190
1137	<i>tional Linguistics and the 12th International Joint</i>	56, Austin, Texas. Association for Computational	1191
1138	<i>Conference on Natural Language Processing (Vol-</i>	Linguistics.	1192
1139	<i>ume 1: Long Papers)</i> , pages 1013–1024, Online only.	Amy L Olex and Bridget T McInnes. 2021. Review of	1193
1140	Association for Computational Linguistics.	temporal reasoning in the clinical domain for timeline	1194
1141	Nasrin Mostafazadeh, Alyson Grealish, Nathanael	extraction: Where we are and where we need to be.	1195
1142	Chambers, James Allen, and Lucy Vanderwende.	<i>Journal of biomedical informatics</i> , 118:103784.	1196
1143	2016. CaTeRS: Causal and temporal relation scheme	Amy L Olex and Bridget T McInnes. 2022. Temporal	1197
1144	for semantic annotation of event structures . In <i>Pro-</i>	disambiguation of relative temporal expressions in	1198
1145	<i>ceedings of the Fourth Workshop on Events</i> , pages	clinical texts. <i>Frontiers in Research Metrics and</i>	1199
1146	51–61, San Diego, California. Association for Com-	<i>Analytics</i> , 7:1001266.	1200
1147	putational Linguistics.	Lis Pereira, Fei Cheng, Masayuki Asahara, and Ichiro	1201
1148	Aakanksha Naik, Luke Breitfeller, and Carolyn Rose.	Kobayashi. 2021. ALICE++: Adversarial training	1202
1149	2019. TDDiscourse: A dataset for discourse-level	for robust and effective temporal reasoning . In <i>Pro-</i>	1203
1150	temporal ordering of events . In <i>Proceedings of the</i>	<i>ceedings of the 35th Pacific Asia Conference on Lan-</i>	1204
1151	<i>20th Annual SIGdial Meeting on Discourse and Dia-</i>	<i>guage, Information and Computation</i> , pages 373–	1205
1152	<i>logue</i> , pages 239–249, Stockholm, Sweden. Associa-	382, Shanghai, China. Association for Computational	1206
1153	tion for Computational Linguistics.	Linguistics.	1207
1154	Marta Guerrero Nieto, Roser Saurí, and Miguel An-	James Pustejovsky. 2003. Timeml: Robust specifica-	1208
1155	gel Bernabé Poveda. 2011. Modes timebank: A	tion of event and temporal expressions in text. In	1209
1156	modern spanish timebank corpus. <i>Procesamiento</i>	<i>Proceedings of the Fifth International Workshop on</i>	1210
1157	<i>del lenguaje natural</i> , 47:259–267.	<i>Computational Semantics (IWCS-5)</i> , 2003.	1211
1158	Qiang Ning, Sanjay Subramanian, and Dan Roth. 2019.	James Pustejovsky, Kiyong Lee, Harry Bunt, and Lau-	1212
1159	An improved neural baseline for temporal relation	rent Romary. 2010. ISO-TimeML: An international	1213
		standard for semantic annotation . In <i>Proceedings</i>	1214

1215			
1216			
1217			
1218			
1219	Alec Radford, Jeffrey Wu, Rewon Child, David Luan,		
1220	Dario Amodei, Ilya Sutskever, et al. 2019. Language		
1221	models are unsupervised multitask learners. <i>OpenAI</i>		
1222	<i>blog</i> , 1(8):9.		
1223	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine		
1224	Lee, Sharan Narang, Michael Matena, Yanqi Zhou,		
1225	Wei Li, and Peter J Liu. 2020. Exploring the limits		
1226	of transfer learning with a unified text-to-text		
1227	transformer. <i>Journal of machine learning research</i> ,		
1228	21(140):1–67.		
1229	Anna Rogers, Marzena Karpinska, Ankita Gupta,		
1230	Vladislav Lialin, Gregory Smelkov, and Anna		
1231	Rumshisky. 2019. Narrativetime: Dense tempo-		
1232	ral annotation on a timeline. <i>arXiv preprint</i>		
1233	<i>arXiv:1908.11443</i> .		
1234	Hayley Ross, Jonathon Cai, and Bonan Min. 2020. Ex-		
1235	ploring Contextualized Neural Language Models for		
1236	Temporal Dependency Parsing. In <i>Proceedings of the</i>		
1237	<i>2020 Conference on Empirical Methods in Natural</i>		
1238	<i>Language Processing (EMNLP)</i> , pages 8548–8553,		
1239	Online. Association for Computational Linguistics.		
1240	Vered Shwartz. 2022. Good night at 4 pm?! time ex-		
1241	pressions in different cultures . In <i>Findings of the As-</i>		
1242	sociation for Computational Linguistics: ACL 2022 ,		
1243	pages 2842–2853, Dublin, Ireland. Association for		
1244	Computational Linguistics.		
1245	Jannik Strötgen, Thomas Bögel, Julian Zell, Ayser Ar-		
1246	miti, Tran Van Canh, and Michael Gertz. 2014. Ex-		
1247	tending HeidelTime for temporal expressions refer-		
1248	ring to historic dates. In <i>Proceedings of the Ninth In-</i>		
1249	<i>ternational Conference on Language Resources and</i>		
1250	<i>Evaluation (LREC’14)</i> , pages 2390–2397, Reykjavik,		
1251	Iceland. European Language Resources Association		
1252	(ELRA).		
1253	Jannik Strötgen and Michael Gertz. 2010. HeidelTime:		
1254	High quality rule-based extraction and normaliza-		
1255	tion of temporal expressions . In <i>Proceedings of the</i>		
1256	<i>5th International Workshop on Semantic Evaluation</i> ,		
1257	pages 321–324, Uppsala, Sweden. Association for		
1258	Computational Linguistics.		
1259	Jannik Strötgen, Anne-Lyse Minard, Lukas Lange,		
1260	Manuela Speranza, and Bernardo Magnini. 2018.		
1261	KRAUTS: A German temporally annotated news cor-		
1262	pus . In <i>Proceedings of the Eleventh International</i>		
1263	<i>Conference on Language Resources and Evaluation</i>		
1264	<i>(LREC 2018)</i> , Miyazaki, Japan. European Language		
1265	Resources Association (ELRA).		
1266	William F. Styler IV, Steven Bethard, Sean Finan,		
1267	Martha Palmer, Sameer Pradhan, Piet C de Groen,		
1268	Brad Erickson, Timothy Miller, Chen Lin, Guergana		
1269	Savova, and James Pustejovsky. 2014. Temporal an-		
1270	notation in the clinical domain . <i>Transactions of the</i>		
	<i>Association for Computational Linguistics</i> , 2:143–		1271
	154.		1272
Xin Su, Phillip Howard, Nagib Hakim, and Steven			1273
Bethard. 2023. Fusing temporal graphs into trans-			1274
formers for time-sensitive question answering . In			1275
<i>Findings of the Association for Computational Lin-</i>			1276
<i>guistics: EMNLP 2023</i> , pages 948–966, Singapore.			1277
Association for Computational Linguistics.			1278
Xin Su, Yiyun Zhao, and Steven Bethard. 2021. The			1279
University of Arizona at SemEval-2021 task 10: Ap-			1280
plying self-training, active learning and data augmen-			1281
tation to source-free domain adaptation . In <i>Proceed-</i>			1282
<i>ings of the 15th International Workshop on Semantic</i>			1283
<i>Evaluation (SemEval-2021)</i> , pages 458–466, Online.			1284
Association for Computational Linguistics.			1285
Weiyi Sun, Anna Rumshisky, and Ozlem Uzuner. 2013.			1286
Evaluating temporal relations in clinical text: 2012			1287
i2b2 challenge. <i>Journal of the American Medical</i>			1288
<i>Informatics Association</i> , 20(5):806–813.			1289
Xingwei Tan, Gabriele Pergola, and Yulan He. 2021.			1290
Extracting event temporal relations via hyperbolic			1291
geometry . In <i>Proceedings of the 2021 Conference</i>			1292
<i>on Empirical Methods in Natural Language Process-</i>			1293
<i>ing</i> , pages 8065–8077, Online and Punta Cana, Do-			1294
minican Republic. Association for Computational			1295
Linguistics.			1296
Xingwei Tan, Gabriele Pergola, and Yulan He. 2023.			1297
Event temporal relation extraction with Bayesian			1298
translational model . In <i>Proceedings of the 17th Con-</i>			1299
<i>ference of the European Chapter of the Association</i>			1300
<i>for Computational Linguistics</i> , pages 1125–1138,			1301
Dubrovnik, Croatia. Association for Computational			1302
Linguistics.			1303
Sun Tiesen and Li Lishuang. 2022. Improving event			1304
temporal relation classification via auxiliary label-			1305
aware contrastive learning . In <i>Proceedings of the</i>			1306
<i>21st Chinese National Conference on Computational</i>			1307
<i>Linguistics</i> , pages 861–871, Nanchang, China. Chi-			1308
nese Information Processing Society of China.			1309
Naushad UzZaman, Hector Llorens, Leon Derczynski,			1310
James Allen, Marc Verhagen, and James Pustejovsky.			1311
2013. SemEval-2013 task 1: TempEval-3: Evaluat-			1312
ing time expressions, events, and temporal relations .			1313
In <i>Second Joint Conference on Lexical and Compu-</i>			1314
<i>tational Semantics (*SEM), Volume 2: Proceedings</i>			1315
<i>of the Seventh International Workshop on Seman-</i>			1316
<i>tic Evaluation (SemEval 2013)</i> , pages 1–9, Atlanta,			1317
Georgia, USA. Association for Computational Lin-			1318
guistics.			1319
Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob			1320
Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz			1321
Kaiser, and Illia Polosukhin. 2017. Attention is all			1322
you need. <i>Advances in neural information processing</i>			1323
<i>systems</i> , 30.			1324
Marc Verhagen, Robert Gaizauskas, Frank Schilder,			1325
Mark Hepple, Graham Katz, and James Pustejovsky.			1326

1327	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.	
1332	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.	
1338	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.	
1344	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.	
1352	Liang Wang, Peifeng Li, and Sheng Xu. 2022. DCT-centered temporal relation extraction . In <i>Proceedings of the 29th International Conference on Computational Linguistics</i> , pages 2087–2097, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.	
1358	Haoyang Wen and Heng Ji. 2021. Utilizing relative event time to enhance event-event temporal relation extraction . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 10431–10437, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	
1365	Siheng Xiong, Ali Payani, Ramana Kompella, and Faramarz Fekri. 2024. Large language models can learn temporal reasoning. <i>arXiv preprint arXiv:2401.06853</i> .	
1369	Sen Yang, Xin Li, Lidong Bing, and Wai Lam. 2023. Once upon a time in graph: Relative-time pretraining for complex temporal reasoning . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 11879–11895, Singapore. Association for Computational Linguistics.	
1376	Hao-Ren Yao, Luke Breitfeller, Aakanksha Naik, Chunxiao Zhou, and Carolyn Rose. 2024a. Distilling multi-scale knowledge for event temporal relation extraction. In <i>Proceedings of the 33rd ACM International Conference on Information and Knowledge Management</i> , pages 2971–2980.	
	Jiarui Yao, Harry Hochheiser, WonJin Yoon, Eli Goldner, and Guergana Savova. 2024b. Overview of the 2024 shared task on chemotherapy treatment timeline extraction . In <i>Proceedings of the 6th Clinical Natural Language Processing Workshop</i> , pages 557–569, Mexico City, Mexico. Association for Computational Linguistics.	1382 1383 1384 1385 1386 1387 1388
	Jiarui Yao, Haoling Qiu, Bonan Min, and Nianwen Xue. 2020. Annotating Temporal Dependency Graphs via Crowdsourcing . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 5368–5380, Online. Association for Computational Linguistics.	1389 1390 1391 1392 1393 1394
	Chenhan Yuan, Qianqian Xie, and Sophia Ananiadou. 2023. Zero-shot temporal relation extraction with ChatGPT . In <i>The 22nd Workshop on Biomedical Natural Language Processing and BioNLP Shared Tasks</i> , pages 92–102, Toronto, Canada. Association for Computational Linguistics.	1395 1396 1397 1398 1399 1400
	Alessandra Zarcone, Touhidul Alam, and Zahra Kolagar. 2020. PATE: A corpus of temporal expressions for the in-car voice assistant domain . In <i>Proceedings of the Twelfth Language Resources and Evaluation Conference</i> , pages 523–530, Marseille, France. European Language Resources Association.	1401 1402 1403 1404 1405 1406
	Shuaicheng Zhang, Qiang Ning, and Lifu Huang. 2022. Extracting temporal event relation with syntax-guided graph transformer . In <i>Findings of the Association for Computational Linguistics: NAACL 2022</i> , pages 379–390, Seattle, United States. Association for Computational Linguistics.	1407 1408 1409 1410 1411 1412
	Yuchen Zhang and Nianwen Xue. 2018. Structured interpretation of temporal relations . In <i>Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)</i> , Miyazaki, Japan. European Language Resources Association (ELRA).	1413 1414 1415 1416 1417 1418
	Yuchen Zhang and Nianwen Xue. 2019. Acquiring structured temporal representation via crowdsourcing: A feasibility study . In <i>Proceedings of the Eighth Joint Conference on Lexical and Computational Semantics (*SEM 2019)</i> , pages 178–185, Minneapolis, Minnesota. Association for Computational Linguistics.	1419 1420 1421 1422 1423 1424 1425
	Xinyu Zhao, Shih-Ting Lin, and Greg Durrett. 2021. Effective distant supervision for temporal relation extraction . In <i>Proceedings of the Second Workshop on Domain Adaptation for NLP</i> , pages 195–203, Kyiv, Ukraine. Association for Computational Linguistics.	1426 1427 1428 1429 1430
	Jie Zhou, Shenpo Dong, Hongkui Tu, Xiaodong Wang, and Yong Dou. 2022. RSGT: Relational structure guided temporal relation extraction . In <i>Proceedings of the 29th International Conference on Computational Linguistics</i> , pages 2001–2010, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.	1431 1432 1433 1434 1435 1436 1437

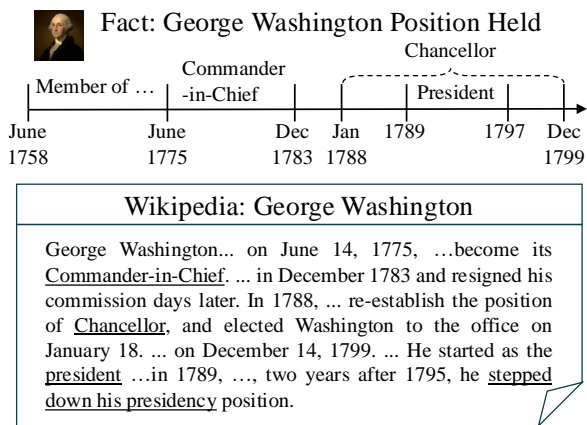


Figure 2: A snippet from George Washington’s Wikipedia page and a timeline regarding his positions.

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 Timeline Examples

We present in Figure 2 a snippet from George Washington’s Wikipedia page alongside the corresponding timeline of his position changes.

B 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.

C 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.

D Timeline-focused Datasets

A notable trend in temporal IE dataset development is the emergence of timeline-focused annotation frameworks that offer more comprehensive and coherent temporal representations compared to traditional approaches. For timeline-centric annotation, Rogers et al. (2019) propose NarrativeTIME, which enables dense, full-coverage temporal relation annotation. Unlike the pairwise TLINK annotation in TimeML, NarrativeTIME constructs coherent narrative timelines, supports underspecification via event types and timeline branches, and achieves significantly higher annotation density. Similarly, Liu and Zhang (2025) introduce ETimeline, a large-scale bilingual (English/Chinese) timeline dataset comprising over 600 timelines and 13,878 annotated event entries, spanning diverse domains from March 2020 to April 2024. Created using an LLM-assisted annotation approach, ETimeline represents a significant resource for cross-lingual timeline construction and temporal reasoning across news domains.

E Temporal Relation Extraction Methods Summary

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

Name	Framework	Domain	Lang	Tasks
<i>TimeML-Based</i>				
TimeBank (Pustejovsky, 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).

1516 F Discussion on Future Directions

1517 In the previous sections, we have identified the
1518 following research opportunities in the field of tem-
1519 poral IE:

- 1520 • Enrich annotation frameworks (Section 3.3),
1521 e.g., representing event arguments or expand-
1522 ing formal semantic systems like SCATE.
- 1523 • Improve dataset diversity (Section 3.3), e.g.,
1524 annotating more domains beyond newswire.
- 1525 • Explore generative approaches (Sections 4.2
1526 and 5.3), e.g., new input-output formulations,
1527 new fine-tuning strategies.
- 1528 • Develop public tools and benchmarks (Sec-

tions 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 visual-
ized for human-computer interaction.

F.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 1. However, these temporal graphs are challenging to interpret independently or use directly for temporal reasoning without extensive

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) (RSGT)	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.

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

1568 between events and times. Yet, such efforts to adapt
1569 and improve annotation frameworks for additional
1570 fields are still scarce. Therefore, adapting existing
1571 annotation frameworks to a broader range of do-
1572 mains to enhance the domain diversity of datasets
1573 represents a potential future research direction.

1574 **F.2 Improve the Application of Generative** 1575 **LLMs**

1576 The application of generative LLMs in the field
1577 of time expression identification, normalization,
1578 and temporal relation extraction remains underex-
1579 plored. Given the proven capabilities of LLMs like
1580 ChatGPT and LLAMA3 across various tasks, it is
1581 logical to probe their potential within the realm of
1582 temporal IE. Whether it involves leveraging new
1583 prompting methods or fine-tuning strategies for
1584 specific tasks, there is ample room for innovation.

1585 However, it is important to emphasize that while
1586 these models excel in generating unstructured text
1587 when applied to temporal IE, it is imperative to spe-
1588 cially design suitable input-output formats. Such
1589 designs are intended to enable generative LLMs,
1590 which are typically used for producing unstructured
1591 text, to also effectively output structured temporal
1592 information.

1593 **F.3 Develop Public Toolkits and Evaluation** 1594 **Benchmarks**

1595 We believe that one key reason transformer-based
1596 temporal IE models have not been widely adopted
1597 might be the absence of a publicly available code
1598 repository that facilitates easier access to models
1599 and data. For example, HuggingFace ¹ provides
1600 language model heads or pipelines suitable for var-
1601 ious tasks, allowing users to easily download and
1602 deploy trained models on any dataset directly from
1603 the HuggingFace Hub. A future research direction
1604 should involve establishing such a repository or
1605 pushing models/datasets to HuggingFace Hub for
1606 the temporal IE tasks to enhance the reproducibility
1607 and applicability of research. Another important
1608 direction is to create a public and test-set concealed
1609 benchmark for a more equitable comparison of
1610 existing work. In most existing works, although
1611 metrics such as F1 scores, precision, and recall
1612 are commonly computed, the specific implementa-
1613 tions can vary. For instance, in [Kanashiro Pereira \(2022\)](#),
1614 only the “before” and “after” relationships
1615 are evaluated for relation extraction performance,

¹<https://huggingface.co/>

1616 whereas [Zhang et al. \(2022\)](#) includes all temporal
1617 relationships except “vague” in their evaluation.

1618 **F.4 Explore More Application Directions**

1619 In reviewing the application of temporal IE sys-
1620 tems, we observe that current research primarily
1621 focuses on aiding “models” in temporal reason-
1622 ing to enhance their performance in other tasks.
1623 Future research in temporal IE should not only con-
1624 tinue to support model performance improvement
1625 but should also pay more attention to serving hu-
1626 mans and enhancing its practical value. A promis-
1627 ing application direction is visualizing timelines in
1628 human-computer interaction (HCI) scenarios. The
1629 visualization results of existing temporal graphs
1630 are often challenging for human users to interpret.
1631 For instance, visualizing the temporal graph of any
1632 document in the TimeBank-Dense dataset might
1633 result in a graph densely populated with points and
1634 lines, offering little help for users to comprehend
1635 the progression of events within the text.

1636 User studies, such as those conducted by [Di Bar-
1637 tolomeo et al. \(2020\)](#), have revealed the impor-
1638 tance of visualization forms of timelines for user
1639 understanding. Consequently, temporal IE research
1640 should also consider incorporating user research
1641 on temporal graphs to guide the design of temporal
1642 IE methods, such as how to represent standardized
1643 time expressions, identify which types of tempo-
1644 ral relations most effectively facilitate time under-
1645 standing, and determine the best ways to present
1646 this information. By addressing these problems,
1647 the extraction and representation of temporal in-
1648 formation can be more closely aligned with user
1649 needs, enhancing its application value in HCI.

1650 **G Comparison with Previous Surveys**

1651 Our survey offers several key advancements over
1652 previous reviews in the field of temporal informa-
1653 tion extraction. Prior surveys such as [Lim et al. \(2019\)](#)
1654 and [Leeuwenberg and Moens \(2019\)](#) pro-
1655 vide only brief mentions of standard datasets like
1656 TimeBank and TempEval, and largely predate the
1657 Transformer era. More recent reviews in the clin-
1658 ical domain—such as [Alfattni et al. \(2020\)](#) and
1659 [Olex and McInnes \(2021\)](#)—present more detailed
1660 dataset descriptions but are limited to clinical texts
1661 and do not cover resources from other domains.

1662 In contrast, our survey compiles and categorizes
1663 32 datasets across multiple domains (newswire,
1664 clinical, Wikipedia, narratives) and 15 languages,

1665 structured by annotation framework (TimeML-
1666 based vs. alternative schemas such as SCATE, tem-
1667 poral dependency trees, or CaTeRS). We provide
1668 a systematic analysis of dataset diversity, domain
1669 bias, language coverage, and annotation schema.
1670 Notably, we quantitatively analyze dataset bias,
1671 identifying that 63% of current datasets come from
1672 the newswire domain, and highlight underexplored
1673 areas such as the low representation of historical
1674 and non-news domains.

1675 Our work specifically focuses on the Trans-
1676 former era, providing in-depth analysis of how
1677 these architectures are applied to temporal IE tasks,
1678 examination of fine-tuning strategies, and discus-
1679 sion of how pre-trained language models capture
1680 temporal information. We also offer a broader
1681 scope in terms of domain and language coverage
1682 compared to previous works that focus on specific
1683 domains or primarily discuss English-language re-
1684 sources.

1685 This broader treatment of datasets and methods
1686 is intentional. Since Transformer-based approaches
1687 often depend heavily on annotated corpora for fine-
1688 tuning or benchmarking, a full understanding of
1689 available datasets and their annotation assumptions
1690 is crucial to contextualizing methodological ad-
1691 vances in temporal information extraction.