

# Transformer-Based Architectures versus Large Language Models in Semantic Event Extraction: Evaluating Strengths and Limitations

Tin Kuculo<sup>a,\*</sup>, Sara Abdollahi<sup>a</sup> and Simon Gottschalk<sup>a</sup>

<sup>a</sup> *L3S Research Center, Leibniz Universität Hannover, Hannover, Germany*  
*E-mails: kuculo@L3S.de, abdollahi@L3S.de, gottschalk@L3S.de*

**Abstract.** Understanding complex societal events reported on the Web, such as military conflicts and political elections, is crucial in digital humanities, computational social science, and news analyses. While event extraction is a well-studied problem in Natural Language Processing, there remains a gap in semantic event extraction methods that leverage event ontologies for capturing multifaceted events in knowledge graphs since existing methods for event extraction often fall short in the semantic depth or lack the flexibility required for a comprehensive event extraction.

In this article, we aim to compare two paradigms to address this task of semantic event extraction: The fine-tuning of traditional transformer-based models versus the use of Large Language Models (LLMs). We exemplify these paradigms with two newly developed approaches: T-SEE for transformer-based and L-SEE for LLM-based semantic event extraction. We present and evaluate these two approaches and discuss their complementary strengths and shortcomings to understand the needs and solutions required for semantic event extraction.

For comparison, both approaches employ the same dual-stage architecture; the first stages focus on multilabel event classification, and the second on relation extraction. While our first approach utilises a span prediction transformer model, our second approach prompts an LLM for event classification and relation extraction, providing the potential event classes and properties. For evaluation, we first assess the performances of T-SEE and L-SEE on two novel datasets sourced from DBpedia and Wikidata, containing over 80,000 Wikipedia sentences and semantic event representations. Then, we perform an extensive analysis of the different types of errors made by these two approaches to discuss a set of phenomena relevant to semantic event extraction.

Our work makes substantial contributions to (i) the integration of Semantic Web technologies and NLP, particularly in the underexplored domain of semantic event extraction, and (ii) the understanding of how LLMs can further enhance semantic event extraction and what challenges need to be considered in comparison to traditional approaches.

**Keywords:** Event Extraction, Transformer Models, Large Language Models, Event Knowledge Graph

## 1. Introduction

Event extraction aims to identify and classify events and their relations in text, including Web sources such as social media, news websites, and online encyclopedias like Wikipedia. Typically, this extraction process is conducted without relying on pre-existing knowledge structures or further structuring of extracted data. In contrast, the goal of

---

\*Corresponding author. E-mail: kuculo@L3S.de.

*semantic event extraction* is to leverage an existing event ontology to lift unstructured text into a structured representation capturing the essence of the event, including its type (e.g., `presidential election`) and relations to entities (e.g., `<US presidential election 2020, successful candidate, Joe Biden>`). Specifically, semantic event extraction aims at enriching knowledge graphs to make event information more accessible, i.e., by adding events that are not yet contained in the knowledge graph because (i) the input texts are about recent events or (ii) the events of that type are considered out of domain (e.g. if the knowledge graph only contains more coarse-grained event types). Practical applications of event knowledge graphs include event-centric visualisations [1, 2], biography generations [3], event narrativisation [4] and question answering over event-related information [5].

Semantic extraction operates at a critical juncture of the Semantic Web and Natural Language Processing (NLP) technologies:

- The Semantic Web offers rich event ontologies such as LODE [6] and the Simple Event Model [7] to represent events. However, cross-domain knowledge graphs such as DBpedia [8] and Wikidata [9] typically focus on named events, such as political summits and natural disasters and lack adaptability to diverse expressions in text-based event descriptions. In addition, relation extraction and link prediction for knowledge graph population typically suffer from noisy data [10–12] and require the presence of the related entities in the knowledge graph [13] and are thus not applicable for extracting relations of newly identified events.
- NLP employs named entity recognition and event extraction techniques to identify finer-grained, transient events like individual meetings or transactions [14] from text. However, traditional NLP methods often deconstruct the task of semantic event extraction into smaller sub-tasks such as event detection [15, 16], and argument extraction [17–19] with each garnering their specific benchmark datasets [20, 21] typically not bound to semantic event ontologies.

This divergence results in a critical gap, creating a need for *semantic event extraction*, blending structured, ontology-based classification with the adaptability to handle a wide range of event types – from transient interactions to significant historical occurrences.

Although some efforts have been made towards semantic event extraction [22, 23], Guan et al. denote that the construction of event knowledge graphs still suffers from the unsatisfactory performance of existing event extraction methods, especially for argument extraction [24]. Most methods still fall short in delivering an integrative approach that works across various domains and effectively accommodates sufficiently rich and diverse ontologies [25–27] centring instead around aged NLP benchmark datasets such as *ACE05* [28] or conversely on highly specific domains [29, 30].

**Example:** As an example of semantic event extraction, consider the event represented in Fig. 1. The text on the left is extracted from the Wikipedia article regarding the “2017 UEFA European Under-21 Championship Final”. We aim to extract relevant event information<sup>1</sup> from that text, such as the final match itself or, potentially, other events mentioned in the text, and enrich an event knowledge graph with newly extracted events and event relations. The right-hand side of the figure illustrates a knowledge graph representation of an extracted event. This representation includes an event class (`final`), an event description derived from the text, the precise location of the game, the date, and other relations.

In this article, we introduce two approaches for semantic event extraction, which follow the same structure but two different paradigms: Transformer-based architectures and Large Language Models (LLMs).<sup>2</sup>

**Transformer-based Semantic Event Extraction (T-SEE):** T-SEE benefits from the strengths of both Semantic Web and NLP techniques and is trained and evaluated on two new datasets, specifically created as a resource for semantic event extraction. T-SEE disentangles the complexities of the task into two manageable sub-tasks:

<sup>1</sup> Given an event and its event class, we consider any information that can be expressed with a property typically used on the respective event class as relevant (e.g., the type of sport of a final).

<sup>2</sup> While LLMs also employ transformers, in this article, we refer to the “traditional” use of transformers, which are fine-tuned to a specific target task, and compare them to pre-trained LLMs prompted for the target task.

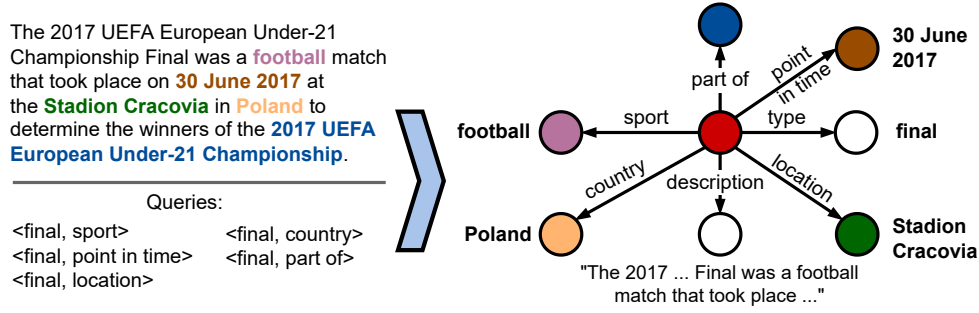


Fig. 1. Example of semantic event extraction for an event mentioned in the Wikipedia article “2017 UEFA European Under-21 Championship Final” using classes and properties in Wikidata. The figure shows a text (top-left), a set of queries consisting of an event class and a property (bottom-left), and the extracted event triples (right).

- **Event Classification:** Approached as a multilabel classification problem, T-SEE determines the most appropriate event labels given a text from a pre-defined set of event classes. In our example, T-SEE applies multilabel classification to categorise the event into *final*.
- **Relation Extraction:** Utilising a span prediction transformer model, we target class-specific relations to construct a nuanced representation of events. In our example, we extract relations such as (*NewEvent*, *country*, *Poland*) employing a set of queries. Here, a query consists of an event class (e.g., *final* or *conflict*) and a property (e.g., *location* or *sport*) and is used to extract the respective information (e.g., the location of the final) within the given text.

**LLM-based Semantic Event Extraction (L-SEE):** With L-SEE, we examine the application of LLMs for semantic event extraction. Given the current prominence of LLMs in various NLP tasks [31, 32], it is pertinent to assess their utility and performance in extracting structured event information from text. In analogy to T-SEE, L-SEE also performs event classification followed by relation extraction, both through specific prompts.

**Evaluation:** To train T-SEE and to evaluate T-SEE and L-SEE, we provide two new semantic event extraction datasets created from DBpedia and Wikidata, containing over 80,000 Wikipedia sentences and semantic event representations. Through a subsequent manual error analysis, we not only aim to gauge the capabilities of LLMs against transformer-based methods but also to identify specific challenges and areas where LLMs might offer novel insights or complement existing approaches.

In this way, we aim to contribute to the ongoing discourse on the potential and limitations of leveraging LLMs for information extraction and knowledge engineering, particularly in cases where LLMs may uncover information beyond the predefined ground truth or existing knowledge graphs.

**Contributions:** In summary, our contributions are:

- We outline the underexplored area of semantic event extraction, situated at the Semantic Web and NLP intersection.
- We present T-SEE and L-SEE, our approaches for semantic event extraction following comparable pipelines, where T-SEE uses a transformer-based architecture and L-SEE uses an LLM.
- We provide two new semantic event extraction datasets created from Wikipedia, Wikidata, and DBpedia: *Wikidata-SEE* and *DBpedia-SEE*.
- We demonstrate the efficacy of T-SEE and L-SEE through empirical evaluations against existing methods.
- We perform an extensive manual annotation of the predictions of T-SEE and L-SEE to identify typical error types and compare the strengths and shortcomings of these two paradigms.
- We make the code<sup>3</sup> and the data<sup>4</sup> available online.

<sup>3</sup><https://github.com/t-kuculo/T-SEE>

<sup>4</sup><https://zenodo.org/records/10818676>

**Structure:** The remainder of this article is structured as follows: In Section 2, we define the task of semantic event extraction. Then, we introduce T-SEE (Section 3) and L-SEE (Section 4). After an automated evaluation of these approaches on a test set (Section 5), we perform our error analysis and discussion in Section 6. After presenting related work (Section 7), we conclude in Section 8.

## 2. Problem Statement

We formally define the problem of semantic event extraction to bridge the gap between granular, structured information and the adaptability required to capture a wide variety of events.

In the context of this work, an *event* is an occurrence of societal importance, typically happening at a specific time and location, involving a set of participants. Examples of events include military conflicts, such as the Second World War, political shakeups, such as Brexit, but also more fine-grained events, such as the battles and air raids in the Second World War or specific football games.

We model information regarding entities (representing real-world events and real-world objects such as persons or locations) and their relations in an event knowledge graph. The classes and properties within the knowledge graph are defined by an event ontology:

**Definition 1** (Event Ontology). *An event ontology  $O = (P, C)$  defines the properties ( $P$ ) and classes ( $C$ ) in an event knowledge graph, where:*

- $P$  is a set of properties describing the types of relations that can hold between two entities and
- $C$  is a set of event classes. An event class can be a sub-class of another event class.

Classes and properties in an event ontology are uniquely identified by an Internationalized Resource Identifier (IRI).<sup>5</sup> Specifically, the property  $p_{\text{type}} \in P$  (typically identified via the property IRI `rdf:type`) assigns an event class to an event.

Other example properties describe the location and number of participants of events. Examples of event classes include `final` as a sub-class of `sporting event`.

Based on an event ontology, we formally define an event knowledge graph as follows:

**Definition 2** (Event Knowledge Graph). *An event knowledge graph  $G_O = (E, V, L, R)$  models entities, events, literals, and their relations following an event ontology  $O = (P, C)$ :*

- $E$  is a set of nodes representing real-world entities.
- $V \subset E$  is a subset of nodes representing real-world events.
- $L$  is a set of literals such as numbers or texts.
- $R \subseteq E \times P \times (E \cup L \cup C)$  is a set of relations.

In a relation  $(e, p_{\text{type}}, c) \in R$  where  $e \in V$  is an event, we require that  $c \in C$  is an event class. This way, we model the class assigned to an event.

We define the task of *semantic event extraction* as follows:

**Definition 3** (Semantic Event Extraction). *Given an event ontology  $O = (P, C)$ , an event knowledge graph  $G_O = (E, V, L, R)$ , and a text  $t$ , the task of semantic event extraction is to detect a set of events described in  $t$  that are not yet represented in  $G_O$ . For each such event  $e_t$ , the task includes:*

1. Identifying its event class relation  $(e_t, p_{\text{type}}, c)$  (event classification), and
2. Extracting a set of relations from  $t$  (relation extraction), with each relation being of the form  $(e_t, p, o)$ , where  $p \in P$  is the property connecting  $e_t$  to  $o \in E \cup L \cup C$ .

*These relations, and the classes they involve, must adhere to the properties and classes of  $O$ .*

<sup>5</sup>Relevant prefixes and namespaces of IRIs used in this article include: `wd`: <http://www.wikidata.org/entity/>, `wdt`: <http://www.wikidata.org/prop/direct/>, `dbp`: <https://dbpedia.org/resource/>, `dbo`: <https://dbpedia.org/ontology/> and `rdf`: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>.

Fig. 1 illustrates an example text ( $t$ ) taken from the Wikipedia article regarding the “2017 UEFA European Under-21 Championship Final”. The semantic event extraction leads to the creation of a new event  $e_t$ , which is typed as the event class `final` and assigned to relations with properties of the event ontology  $O$  (e.g., `location` and `point in time`). These relations can be serialised as RDF triples to be used in downstream applications.

## 2.1. Assumptions

To perform semantic event extraction given the defined problem statement, we propose methodologies that employ transformers and LLMs based on the following assumptions:

### Tasks and Models

- **Task Representation:** Following Definition 3, we frame semantic event extraction as a two-step task: event classification followed by relation extraction. This decomposition is assumed to be effective and meaningful for capturing events and their relations. Further, we directly intertwine the tasks of event detection and event classification: event classification detects and classifies events at the same time, i.e., there are no events without event class.
- **Task Dependency:** Relation extraction depends on the results of event classification. This dependency is intentional, as event classes determine which relations to extract. Consequently, we assume errors to propagate across the entire pipeline, so any misclassifications naturally affects relation extraction results. This error propagation needs to be reflected during evaluation.
- **Model Selection:** We assume that both transformer-based models and LLMs are suitable for semantic event extraction.
  - \* *Transformers:* Fine-tuned transformer models (e.g., BERT) are assumed to generalize effectively for event classification and relation extraction when trained on high-quality, ontology-aligned datasets.
  - \* *LLMs:* LLMs are assumed to generate structured outputs reliably when prompted with event ontologies. However, we acknowledge LLMs’ sensitivity to prompt design and their tendency to hallucinate relations not present in training data, requiring careful validation.

### Data

- **Event Ontology Scope:** The selected event ontology must comprehensively define event classes and properties for the target domain. We assume the event ontology is extracted from a knowledge graph (e.g., Wikidata) and filtered to exclude overly specific or metadata-like entries. As described in our evaluation setup in Section 5.1.1, we use two event ontologies for training and evaluation, extracted from DBpedia and Wikidata.
- **Data Availability:** Training data must consist of texts annotated with events, classes, and relations aligned with the event ontology. As described in our evaluation setup in Section 5.1, we use two datasets for training and evaluation. They contain triples from DBpedia and Wikidata, respectively, both linked to texts from Wikipedia. An example of a text annotated with Wikidata triples is given in Table 1.
- **Annotation Quality:** In order to generate such large-scale datasets, we assume distant supervision during dataset creation to link triples to Wikipedia texts, acknowledging potential noise in annotations. Consequently, even despite a cautious dataset creation process, ground truth annotations may still contain omissions or inaccuracies, particularly in large-scale datasets. Further, annotations can vary regarding granularity (e.g., `dbo:SportsEvent` vs. `dbo:TennisTournament`) and completeness. This assumption motivates to perform manual validation of the evaluation results as we do in Section 6.

### Evaluation

- **Setup:** As described above, the evaluation setting requires a training and evaluation dataset and needs to assess the quality of event classification, relation extraction and their combination in semantic event extraction.
- **Metrics:** Metrics must reflect pipeline-wide performance, including error propagation. Therefore, we compute precision, recall and  $F_1$  scores for the tasks of event classification and relation extraction in isolation and in combination. For LLM-based methods, we additionally assume consistency metrics (e.g., Fleiss’  $\kappa$ ) to account for stochastic outputs.

Table 1

Example of an annotated text as required in a dataset required for training and evaluating a semantic event extraction model. This example is based on Fig. 1 using Wikidata as the target event ontology. The given text  $t$  mentions two events (here, marked in bold for convenience).

	$e_{t_1}$	$e_{t_2}$
$t$	The <b>2017 UEFA European Under-21 Championship Final</b> was a football match that took place on 30 June 2017 at the Stadion Cracovia in Poland to determine the winners of the <b>2017 UEFA European Under-21 Championship</b> .	
C	final	season
R	<ul style="list-style-type: none"> <li>• (NewEvent1, sport, football)</li> <li>• (NewEvent1, country, Poland)</li> <li>• (NewEvent1, <math>p_{type}</math>, final)</li> <li>• (NewEvent1, point in time, "2017-07-30")</li> <li>• (NewEvent1, location, Stadion Cracovia)</li> <li>• ...</li> </ul>	<ul style="list-style-type: none"> <li>• (NewEvent2, country, Poland)</li> <li>• (NewEvent2, point in time, "2017")</li> <li>• (NewEvent2, <math>p_{type}</math>, season)</li> </ul>

- **Error Analysis:** We assume manual error analysis is critical to identify phenomena like event ambiguity, type misalignment, and annotation discrepancies, which automated metrics may overlook.

### 3. T-SEE: Transformer-based Semantic Event Extraction

In this section, we present T-SEE (Transformer-based Semantic Event Extraction), an approach for semantic event extraction based on a transformer architecture. The design of T-SEE is guided by two goals:

1. Through a 3-step procedure of event classification, relation extraction and event modelling, we ensure comparability with L-SEE, our LLM-based approach presented in the next section (Section 4).
2. To allow seamless integration into the Semantic Web, the whole architecture of T-SEE needs to be guided through an event ontology, its RDF classes and properties.

Fig. 2 offers a visual summary of T-SEE following these goals. Given an event ontology  $O$ , we generate a set of queries  $Q$  during the preprocessing phase (*query generation*) as a basis of the query-based relation extraction. T-SEE then carries out a three-step process to extract and semantically represent events from a given text  $t$ :

1. *Event classification:* We formulate event classification as a multilabel classification problem and apply it to a given text  $t$  to identify event mentions and their classes. This enables us to classify all event mentions within the text concurrently.
2. *Query-based relation extraction:* For each identified event, we extract its relations using a transformer-based extraction model and a subset of  $Q$ , i.e., selected queries used to extract relevant relations of the detected events. After the appropriate queries have been selected, we train our relation extraction model on pairs of event classes and properties.
3. *Event modelling:* We transform the extracted event information into triples and add them to the event knowledge graph  $G_O$ .

With this process, T-SEE focuses on event classification and subsequent event relation extraction, aiming to generate a robust and comprehensive representation of event knowledge. We build on three key factors: (i) the inherent strengths of transformer models, including their capacity to encapsulate complex semantic relationships within the text; (ii) the use of task-specific fine-tuning of these models that allows us to tailor their powerful general language understanding capabilities to our specific extraction tasks, and (iii) the structural guidance provided by an event ontology, which not only aligns the model’s understanding of events with existing schemas but also offers adaptability accommodating emerging event types, such as "pandemic".

In the following, we describe T-SEE’s steps in more detail, along with its algorithm and a running example for a more intuitive understanding.

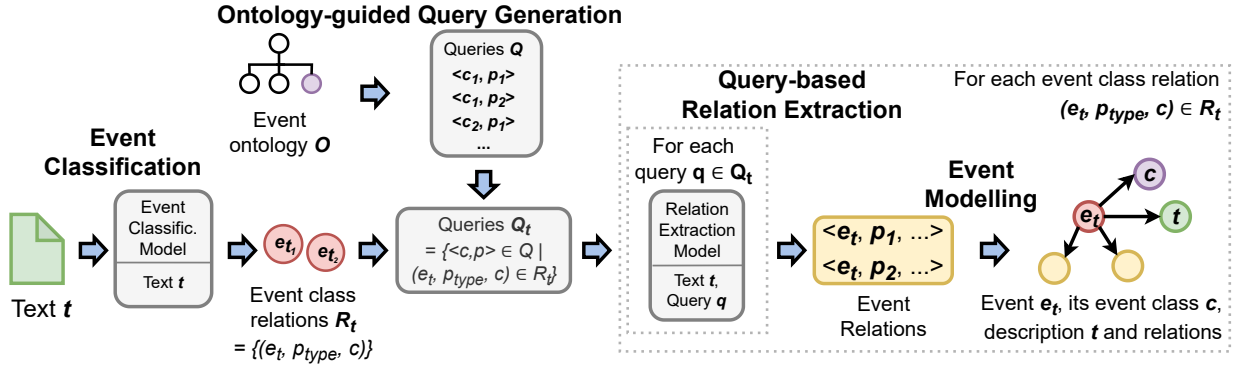


Fig. 2. Overview of T-SEE, showing how it extracts and models a single event. Inputs to the models are shown below the horizontal lines.

**Algorithm:** Algorithm 1 provides an overview of T-SEE. The algorithm embodies the three main inference steps explained earlier: event classification, query-based relation extraction, and event modelling.

---

**Algorithm 1** Transformer-based Semantic Event Extraction (T-SEE)

---

```

1: Input
2:    $t$       Text
3:    $O$       Event ontology
4:    $R$       The relations in an event knowledge graph  $G_O$ 
5:    $Q$       Set of queries
6:    $ECM$     Event Classification Model (trained)
7:    $REM$     Relation Extraction Model (trained)
8:
9:  $R_t = \{(e_t, p_{type}, c)\} \leftarrow ECM.classifyEvents(t, O)$  ▷ Event classification (Section 3.2)
10:
11: for each  $(e_t, p_{type}, c) \in R_t$  do ▷ Query-based relation extraction (Section 3.3)
12:    $R_{e_t} \leftarrow \{\}$ 
13:   for each  $q = \langle c, p \rangle \in getQueries(Q, c)$  do
14:      $result \leftarrow REM.getQueryResult(t, q)$ 
15:      $R_{e_t} \leftarrow R_{e_t} \cup REM.createRelations(e_t, p, result)$ 
16:
17:    $R = R \cup (e_t, p_{type}, c)$  ▷ Event modelling (Section 3.4)
18:    $R = R \cup (e_t, p_{description}, t) \cup R_{e_t}$ 

```

---

**Example:** We exemplify each of the steps based on the example illustrated in Fig. 3, where the text  $t$  pertains to protests in Tehran. T-SEE extracts two events ( $e_{t_1}$  and  $e_{t_2}$ ), their classes<sup>6</sup> (conflict and revolution) and relations. This example demonstrates how T-SEE’s relation extraction model is capable of extracting different relations for each detected event, for instance,  $(e_{t_1}, p_{participant}, \text{Government of Iran})$ <sup>7</sup> and  $(e_{t_2}, p_{location}, \text{Tehran})$ .

### 3.1. Ontology-guided Query Generation

The query generation step is a preprocessing step that creates a set of queries  $Q$  used later as input to the query-based relation extraction model. The generation of  $Q$  is guided by an event ontology  $O$ , such that each query  $q = \langle c, p \rangle \in Q$  comprises the event class  $c$  and a corresponding property  $p$  as defined in the event ontology. For

<sup>6</sup>As our event ontology  $O$  in this example, we select an event ontology extracted from Wikidata.

<sup>7</sup>For readability of the example relations, we represent selected entities and classes through their labels.

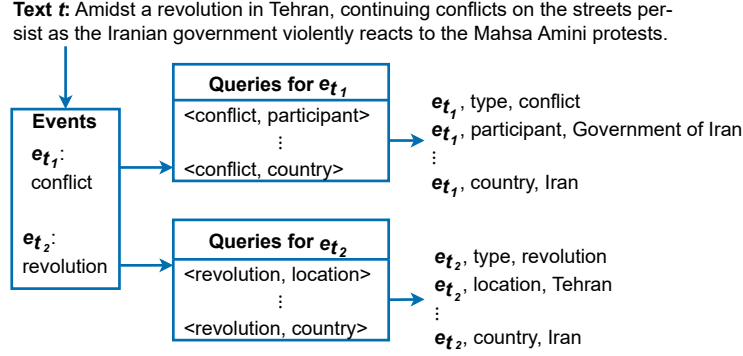


Fig. 3. Example of event classification and query-based relation extraction on a sentence in the Wikipedia article “Mahsa Amini protests”.

each considered event class in  $O^8$ , a set of queries is added to  $Q$ . These queries are then used in T-SEE’s query-based relation extraction step.

Given an event ontology  $O = (P, C)$  and an event knowledge graph  $G_O = (E, V, L, R)$ , we create these queries as follows: For each event class  $c \in C$ , we select a set of properties that are used together with events of this class in  $G_O$ :  $\{p \mid (e, p, x) \in R \wedge (e, p_{\text{type}}, c) \in R\}$ . To avoid the inclusion of inappropriate queries (e.g., infrequent event classes and metadata properties), additional constraints can be applied to remove queries from  $Q$ . We describe our constraints in Section 5.1 and make our sets of event classes, properties, and queries available<sup>9</sup>.

### 3.1.1. Example

Fig. 4 shows an example Wikidata SPARQL query to extract Wikidata properties commonly (more than 50 times) used on entities classified as `wd:Q180684` (conflict). It returns 22 properties, including `wdt:P17` (country) and `wdt:P710` (participant).

```
PREFIX wd: <http://www.wikidata.org/entity/>
PREFIX wdt: <http://www.wikidata.org/prop/direct/>

SELECT ?property WHERE {
  ?s ?property ?o .
  ?s wdt:P31 wd:Q180684 . # instance of conflict
  FILTER(STRSTARTS(STR(?property), STR(wdt:))) .
} GROUP BY ?property HAVING(COUNT(?s) > 50)
```

Fig. 4. SPARQL query on Wikidata to extract Wikidata properties commonly (more than 50 times) used on entities classified as “conflict”.

Using such SPARQL queries, we can generate queries used by T-SEE. Table 2 provides examples of four such queries for two Wikidata event classes, each together with their properties.

## 3.2. Event Classification

Given a text  $t$ , the goal of T-SEE’s event classification step is to identify a set of events that occur in  $t$  and to detect their event classes, i.e. the set of relations  $R_t = \{(e_t, p_{\text{type}}, c)\}$  (line 9 in Algorithm 1). To do so, we propose a multilabel event classification model based on a transformer architecture [33], which allows for the efficient and effective processing of input texts.

Specifically, the input to our event classification model is a sequence of tokens derived from  $t$  representing one or more event mentions in the text. The model processes the input sequence using a series of self-attention mechanisms,

<sup>8</sup>We consider all event classes and properties in the event ontology  $O$  that appear in the training data.

<sup>9</sup>[https://github.com/t-kuculo/T-SEE/blob/main/processing/filtered\\_wikidata\\_event2.schema](https://github.com/t-kuculo/T-SEE/blob/main/processing/filtered_wikidata_event2.schema)



Table 2  
Example queries extracted from the Wikidata ontology.

Event class ( $c$ )	Property ( $p$ )	Query ( $q = \langle c, p \rangle$ )
conflict (wd:Q180684)	participant (wdt:P710)	<conflict, participant>
	country (wdt:P17)	<conflict, country>
revolution (wd:Q10931)	location (wdt:P276)	<revolution, location>
	country (wdt:P17)	<revolution, country>

allowing it to capture complex relationships between contextual and semantic information of the input  $t$  [33]. The output of the transformer-based architecture is a sequence of hidden states, which encodes the relevant information from the input sequence.

The hidden states are then passed through a dropout layer to reduce the number of connections between the pre-trained layers and the downstream layers, effectively forcing the downstream layers to learn more robust and generalisable representations of the input data. Finally, a fully-connected layer and a Sigmoid activation function are used in the output layer, generating a probability distribution over the possible event classes in the input text.

Additionally, we conduct threshold optimisation on a validation set. Prior work on multilabel classification, such as binary relevance methods [34], often employs a fixed decision threshold (usually 0.5) to convert predicted probabilities into class labels. However, this may not be optimal for all classes, especially in cases with imbalanced data or differing class complexities. To address this issue, we utilise an optimisation strategy that fine-tunes individual decision thresholds for each label, aiming to maximise the  $F_1$  score.

### 3.2.1. Example

In our example, the event classification model receives the whole text shown in Fig. 3 (“Amidst a revolution in Tehran, continuing conflicts on the streets persist as the Iranian government violently reacts to the Mahsa Amini protests.”) as an input and returns two event classes (`conflict` and `revolution`) corresponding to the two events in the text.

### 3.2.2. Training

To train T-SEE’s event classification model, a corpus that contains texts and event class labels corresponding to the events represented in each individual text is required. Specifically, we utilise two datasets that contain sentences from Wikipedia, annotated with events and their relations from Wikidata and DBpedia, respectively. These datasets are described in detail in Section 5.1. The multilabel classification model is fed the tokenised input texts and uses a focal loss function [35].

### 3.3. Query-based Relation Extraction

Given the text  $t$  and the set  $R_t = \{(e_t, p_{type}, c)\}$  of detected events together with their predicted event classes, the goal of relation extraction is to detect, extract, and assign relations found in  $t$  to the matching events. T-SEE utilises a subset of the generated queries  $Q$  that can be matched to the predicted event classes of the extracted events in  $R_t$ . Specifically, given  $R_t = \{(e_t, p_{type}, c)\}$ , we select those queries in  $Q$  which ask about these event classes:  $Q_t = \{q = \langle c, p \rangle \in Q \mid \exists (e_t, p_{type}, c) \in R_t\}$  (line 13 in Algorithm 1). Together with  $t$ , these queries serve as input to our query-based relation extraction model.

We leverage BERT [36] as the base of T-SEE as it provides a nuanced understanding of semantics, capturing the meaning and context of words and sentences in text. BERT is known for its proficiency in capturing long-range dependencies, a crucial aspect of comprehending the complexities of textual narratives. In addition, BERT incorporates a Next Sentence Prediction loss, which is specifically designed to model the coherence between sentences. This element of coherence is particularly valuable for relation extraction tasks. By understanding the continuity of text, the model is empowered to decipher the intricate relationships between entities that might be scattered across the text.

Specifically, we encode the text  $t$  and a query  $q$  as fixed-length vectors. The decoded results then correspond to a probability distribution over token spans that represent possible relation values.

As shown in line 14 of Algorithm 1, each selected query  $q = \langle c, p \rangle$ , and context represented by the text  $t$  are passed through our query-based relation extraction model, generating results and their associated confidence scores. Together with the respective event and the property  $p$ , each result resembles a relation.

### 3.3.1. Example

For our predicted event classes `conflict` and `revolution`, the queries in  $Q$  cover a variety of Wikidata properties such as `participant` and `location`. As shown in Fig. 3, given the query  $\langle \text{revolution}, \text{location} \rangle$ , we infer its result `Tehran`, i.e., the relation  $(e_{t_2}, p_{\text{location}}, \text{Tehran})$ . This process is repeated for each query-context pair, creating, for each accepted result, a relation.

### 3.3.2. Training

To train our query-based relation extraction model, we use a corpus of texts with event mentions and their relations with properties in the event ontology  $O$ . As in [37], the model is jointly trained using a span extraction loss and a logistic regression loss for an additional classifier that predicts answerability [38, 39]. During training, the model is rewarded for selecting token spans that correspond to correct relation values between an event of a given event class label and entities or literals that occur in the text.

## 3.4. Event Modelling

In the event modelling step, we materialise the extracted event information as triples and enrich the event knowledge graph with them (line 17 - 18). Precisely, for each text  $t$ , and each of the event class relations  $(e_t, p_{\text{type}}, c) \in R_t$ , we create the following relations:

- Type relation for  $e_t$ :  $(e_t, p_{\text{type}}, c)$
- Description of  $e_t$ :  $(e_t, p_{\text{description}}, t)$
- Relations extracted with our query-based relation extraction

This process is repeated for all texts in an input corpus and the events extracted within them, after which the ontology-mapped relations can be transformed into RDF triples. As described in Definition 3, the event modelling step creates new triples of events not yet represented in the target knowledge graph  $G_O$ . Event classes, properties and their values were identified in the extraction process guided by the event ontology  $O$ .

For representing the provenance and explicitly providing the source of the semantic event representation, further information could be added, e.g., a URL pointing to the source text and a description of the extraction method. To do so, sources can be directly linked to a source statement in Wikidata<sup>10</sup>. Another option would be to use the PROV-O ontology [40].

### 3.4.1. Example

Fig. 3 illustrates relations extracted for the example events `conflict` and `revolution`. Given the `conflict` event, the following relations are created:

- $(e_{t_1}, p_{\text{type}}, \text{conflict})$
- $(e_{t_1}, p_{\text{description}}, \text{"Amidst a revolution in Tehran, continuing conflicts on the streets persist as the Iranian government violently reacts to the Mahsa Amini protests."})$
- $(e_{t_1}, p_{\text{participant}}, \text{Government of Iran})$
- $(e_{t_1}, p_{\text{country}}, \text{Iran})$

We provide examples of generated RDF triples in Section 5.6.

## 4. L-SEE: LLM-based Semantic Event Extraction

In this section, we present L-SEE (LLM-based Semantic Event Extraction), an approach for semantic event extraction based on a Large Language Model. As LLMs continue redefining the boundaries of NLP, their application

<sup>10</sup><https://www.wikidata.org/wiki/Help:Sources>

in semantic event extraction presents a compelling approach for assessing their standalone capabilities and potential synergies with pipeline-based methodologies that decompose event extraction into event detection and argument extraction, as commonly employed in state-of-the-art approaches in foundational work [41] and retained in recent studies [42–44].

Fig. 5 offers a visual summary of L-SEE whose bottom part is analogous to T-SEE in Fig. 2. Given an event ontology  $O$ , the set of all event classes  $C$  is extracted beforehand. As in T-SEE, L-SEE then carries out a three-step process to extract and semantically represent events from a given text  $t$ :

1. *Event classification*: We perform event classification as a multilabel classification problem by prompting an LLM to detect events and their classes in a text  $t$  given  $C$ .
2. *Relation extraction*: We prompt the LLM to extract relations of all identified events.
3. *Event modelling*: We transform the extracted event information into triples and add them to the event knowledge graph  $G_O$ .

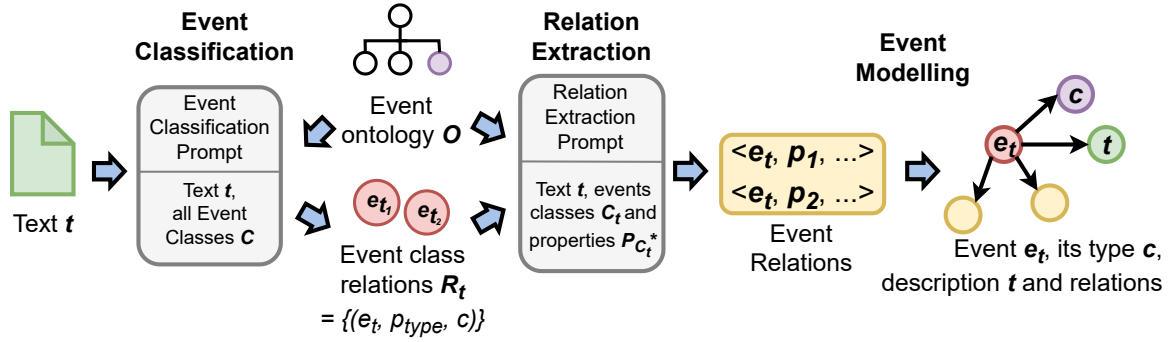


Fig. 5. Overview of L-SEE, showing how it extracts and models a single event. Inputs to the prompts are shown below the horizontal lines. \*  $C_t$  refers to all classes of the events detected in the text.  $P_{C_t}$  is the set of properties used together with these event classes in the relations  $R$  of the target knowledge graph.

**Algorithm:** Algorithm 2 provides an overview of L-SEE and its three steps: event classification, relation extraction and event modelling.

---

**Algorithm 2** LLM-based Semantic Event Extraction (L-SEE)

---

```

1: Input
2:    $t$            Text
3:    $O = (P, C)$    Event ontology
4:    $R$            The relations in an event knowledge graph  $G_O$ 
5:    $LLM$         Large Language Model (pre-trained)
6:
7:    $R_t = \{(e_t, p_{type}, c)\} \leftarrow LLM.classifyEvents(t, C)$            ▶ Event classification (Section 4.1)
8:    $C_t = \{c \mid (e_t, p_{type}, c) \in R_t\}$ 
9:
10:   $P_{C_t} \leftarrow getPropertiesOfClasses(O, C_t)$ 
11:   $R_{e_t} \leftarrow LLM.getRelations(t, C_t, P_{C_t})$            ▶ Relation extraction (Section 4.2)
12:
13:  for each  $(e_t, p_{type}, c) \in R_t$  do
14:     $R = R \cup (e_t, p_{type}, c)$            ▶ Event modelling (Section 4.3)
15:     $\cup (e_t, p_{description}, t) \cup R_{e_t}$ 

```

---

#### 4.1. Event Classification

For event classification (line 7 in Algorithm 2), L-SEE guides the LLM with a precise prompting mechanism to identify and categorise events in the text  $t$ , given the set  $C$  of event classes in the target event ontology. This step builds upon the LLM's ability to discern events of significance akin to those warranting dedicated Wikipedia entries, ensuring the extraction of events with substantial relevance.

The event classification LLM prompt template is shown in Fig. 16 in the Appendix, where  $C$  is formatted like `['conflict', 'revolution']`. In detail, the event classification LLM prompt template consists of the following parts:

1. **Instruction:** Explicitly defines event classification and the operational definition of an "event".
2. **Example:** Illustrates the task with a one-shot example, including a sample text, identified event classes, and explanations to clarify expectations.
3. **Output options:** Explicitly lists the full set of potential outputs, i.e., the set of all event classes  $C$  in our target event ontology.
4. **Task:** Specifies the input text  $t$  for classification.

#### 4.2. Relation Extraction

For relation extraction (line 11 in Algorithm 2), L-SEE prompts the LLM a second time, now to extract the relations of each identified event, given the event classes  $C_t = \{c \mid (e_t, p_{type}, c) \in R_t\}$  identified in the previous step together with the set of properties  $P_{C_t}$  used on these classes (extracted as described in Section 3.1, i.e.,  $P_{C_t} = \{p \mid (e, p, x) \in R \wedge (e, p_{type}, c) \in R \wedge c \in C_t\}$ ).

The relation extraction LLM prompt template is shown in Fig. 17 in the Appendix. For our condensed example in Table 2,  $C_t$  and  $P_{C_t}$  would be added to the prompt formatted as `{ 'conflict': ['participant', 'conflict'], 'revolution': ['location', 'country'] }`. In detail, the relation extraction LLM prompt template consists of the following parts:

1. **Instruction:** Defines the task (relation extraction), specifies expected property-value formats (e.g., temporal or spatial attributes), and mandates valid JSON output. Semantic constraints enforced through data type conventions ensure consistency for downstream processing.
2. **Example:** Provides a one-shot demonstration with a text snippet, event classes, properties, and a corresponding JSON output to model structured responses.
3. **Task:** Presents the input text  $t$ , the event classes  $C_t$  and properties  $P_{C_t}$ , requiring the LLM to populate these properties with text-derived values.

#### 4.3. Event Modelling

The event modelling step (line 14 - 15 in Algorithm 2) follows the procedure outlined in T-SEE, as detailed in Section 3.4. This process results in the creation of RDF triples that represent newly identified events and their relations.

### 5. Evaluation

In this section, we introduce two new datasets for semantic event extraction and compare T-SEE and L-SEE to event extraction baselines. Finally, we show an example of the generated RDF triples and compare the consistency of LLM outputs over different executions.

## 5.1. Datasets

We introduce two new large-scale datasets that currently stand as the largest and most diverse datasets for the task of semantic event extraction and follow our assumptions states in Section 2.1: *DBpedia-SEE* and *Wikidata-SEE*. They are available online.<sup>11</sup> *DBpedia-SEE* and *Wikidata-SEE* serve as training and test corpora for semantic event extraction based on event ontologies of DBpedia and Wikidata. To comply with the definition of semantic event extraction in Definition 3, each dataset belongs to an event ontology  $G_O$  and contains a set of texts, where each text  $t$  is annotated with a set of events, their classes and relations.

### 5.1.1. Event Ontology Extraction

In the first step, we extract relevant event classes and their properties from DBpedia and Wikidata to create two event ontologies. The main reason why we extract event ontologies from DBpedia and Wikidata instead of using event ontologies such as LODE [6] and the Simple Event Model [7] is that we do not only require an event ontology but also a large corpus of events modelled with such ontology, as available in the DBpedia and Wikidata knowledge graphs. Further reasons are as follows: (i) we focus on cross-domain knowledge graphs, with DBpedia and Wikidata being well-established cross-domain knowledge graphs yet inherently incomplete and bear potential for extension [45], (ii) as described in Section 1, we focus on named events and (iii) to create our evaluation datasets (see next Section 5.1.2), we utilise Wikipedia links which can be directly mapped to Wikidata and DBpedia entities.

*Filtering Protocols and Thresholds.* To ensure the quality and relevance of the event classes and properties extracted from Wikidata and DBpedia, we apply stringent filtering protocols. Specifically, we restrict event classes and properties to those used in the context of events and apply a minimum threshold for event classes (100 appearances) and properties (50 appearances). These thresholds serve two key methodological purposes. First, by requiring a minimum frequency, we ensure that classes and properties are sufficiently represented across the training, validation, and test splits, thereby improving the statistical reliability of our evaluation. Second, consistent filtering avoids scenarios where very rare or highly specialised event classes (e.g., classes associated with only few events such as *City of Cardiff Council election* in Wikidata) might skew macro-averaged metrics or lead to overfitting on sparse patterns. We arrived at these particular numbers by analysing the frequency distributions of event-related resources in *DBpedia-SEE* and *Wikidata-SEE*, finding that they effectively preserve the majority of relevant classes and properties while excluding rarely used or metadata-like entries. We acknowledge that different use cases or domain-specific requirements might call for alternative cutoffs, but this balance between coverage and reliability is well-suited to our current scope.

While we try to keep manual interventions minimal and to be as consistent as possible in our annotations, for the remaining events and properties, we need to manually filter out overly specific event classes and metadata properties. Specifically, for Wikidata, we filtered out the following three types of event classes and properties:

- Event classes specifically about a country (we still consider their parent classes. For example, instead of "UK Parliamentary by-election", there still is "by-election"). Examples are:
  - \* Turkish general election (wd:Q22333900)
  - \* Spanish Grand Prix (wd:Q9208)
  - \* Sydney International (wd:Q248952)
- Classes that are wrongly categorised as event classes in Wikidata. Examples are:
  - \* communications satellite (wd:Q149918)
  - \* space telescope (wd:Q148578)
  - \* crewed spacecraft (wd:Q7217761)
- Properties that do not represent real-world relations (e.g., identifiers). An example is:
  - \* X username (wdt:P2002)

<sup>11</sup><https://zenodo.org/records/10818676>

Table 3  
Statistic of the extracted DBpedia and Wikidata event ontologies.

	<i>DBpedia-SEE</i>	<i>Wikidata-SEE</i>
<b>Event Classes</b>	19	60
<b>most occurrences</b>	dbo:MilitaryConflict (23, 264)	wd:Q27020041 (sports season) (5, 901)
<b>least occurrences</b>	dbo:MixedMartialArtsEvent (104)	wd:Q1079023 (championship) (55)
<b>Properties</b>	29	17
<b>most occurrences</b>	dbo:place (17, 618)	wdt:P585 (point in time) (31, 378)
<b>least occurrences</b>	dbo:previousMission (72)	wdt:P571 (inception) (20)

Statistics of the resulting DBpedia and Wikidata event ontologies are shown in Table 3. For example, the Wikidata event ontology has 60 event classes and 5,901 events typed as `sports season`. We consider the two event ontologies independently from each other and do not align them.<sup>12</sup> This way, we are able to evaluate semantic event extraction on two distinct datasets and thus demonstrate the generalisability of our models. Further, both the DBpedia [47, 48] and Wikidata [48–50] ontologies have been successfully used to represent events in other works.

### 5.1.2. Extraction of Event Triples

To extract texts and the RDF triples representing mentioned events, we follow a distance-label generation process.<sup>13</sup> The individual texts are sentences extracted from articles in the English Wikipedia describing events<sup>14</sup>. Event classes and relations are extracted by exploiting existing links to events and their DBpedia or Wikidata representations.

Fig. 6 illustrates the distance-label generation process at an example: The Wikipedia article “Turkish involvement in the Syrian civil war” has a link to the event “Operation Euphrates Shield” which has a relation to Syria and is also mentioned in the same text. Consequently, we select the text, the event class `military operation`<sup>15</sup>, and the `country` relation to `Syria`.



Fig. 6. Example illustrating how we label texts with events and relations. The Wikipedia text on the left links to the Wikidata event on the right side, which also has a relation to an entity mentioned in the text: `<country, Syria>`.

### 5.1.3. Statistics

As delineated in Table 4, *DBpedia-SEE* includes 42,648 texts, and *Wikidata-SEE* contains 37,988 texts, where each text contains at least one annotated event and its corresponding relations. Together, these datasets feature over 80,636 uniquely annotated events and 111,663 relation instances, making them the most extensive repositories for training and evaluating event extraction models to date.

<sup>12</sup>A discussion about the task of event ontology alignment involving Wikidata is given by Guo et al. [46].

<sup>13</sup>We follow approaches such as [51, 52] that associate text with RDF triples from a knowledge graph.

<sup>14</sup>Event articles typically contain descriptions of related events.

<sup>15</sup>If an event has multiple event classes, we select the most infrequently used event class among them in order to add fine-grained event classes to the dataset.

Table 4  
Statistic of our datasets for semantic event extraction.

	<i>DBpedia-SEE</i>	<i>Wikidata-SEE</i>
<b>Texts</b>	42,648	37,988
<b>Events</b>	42,726	38,014
<b>Relations</b>	47,666	63,997

#### 5.1.4. Comparison to Existing Datasets

*DBpedia-SEE* and *Wikidata-SEE* distinctly surpass existing benchmarks for the task of semantic event extraction due to their use of RDF annotations, their focus on general-domain events with societal impact and the coverage of both event detection and relation extraction annotations. Datasets such as SuicideED [53], SCIERC[54] and GENIA [55] only cover very domain-specific events. MAVEN [20] and MINION [56] only provide annotations for event detection, not relation or argument extraction. The existing larger event datasets like GDELT [57, 58] are less structured and not in RDF<sup>16</sup>. In a comparison to the *ACE05* [28] dataset typically used for event extraction, our datasets *DBpedia-SEE* and *Wikidata-SEE*:

- are freely available
  - \* *ACE05* is only available for \$4,000.00 to non-members of the Linguistic Data Consortium.
- have wider coverage of event domains
  - \* e.g., *ACE05* does not have sport-related events
- use RDF classes and properties
- have a large number of event classes and properties
  - \* *DBpedia-SEE*: 19 event classes and 29 properties
  - \* *Wikidata-SEE*: 60 event classes and 17 properties
  - \* *ACE05*: 33 event classes and 22 arguments
- provide a large number of texts
  - \* *DBpedia-SEE*: 42,648 texts
  - \* *Wikidata-SEE*: 37,988 texts
  - \* *ACE05*: 599 texts

These attributes amplify the datasets’ potential for semantic event extraction, which can not be performed with other existing datasets.

#### 5.1.5. Data Preparation and Experiment Design

With our distantly labelled datasets *DBpedia-SEE* and *Wikidata-SEE*, we are able to i) train T-SEE and the baselines on large-scale datasets and (ii) evaluate their performance in the semantic event extraction of events which already exist in DBpedia or Wikidata. We exclude links to existing events when running the experiments to simulate the situation in which the events do not yet exist in the target knowledge graph.

In our experiments, we split the datasets into training, test, and validation sets using 70:15:15 splits.

## 5.2. Evaluation Setup

Next, we describe our evaluation setup, i.e., baselines and metrics.

<sup>16</sup>Instead of concise event mentions, GDELT considers whole articles as texts and does not provide relation types between events and entities.

### 5.2.1. Baselines

We compare T-SEE against two baselines:

- Text2Event [59]: A state-of-the-art method for event extraction using a sequence-to-structure generation paradigm.
- EventGraph [60]: A method for event extraction using semantic graph parsing that has shown state-of-the-art results for the task of argument role classification.

The selection of baselines for our study is carefully considered but constrained by the availability and adaptability of existing event extraction methodologies due to the following reasons: (i) Despite their valuable contributions, several works do not provide any accessible implementations [61–63], which is a critical barrier to replication and further research. (ii) The usability of many event extraction frameworks is hampered by a lack of comprehensive documentation and a dependency on specific or proprietary datasets, notably the *ACE05* dataset [64–67]. Other methodologies like *DEGREE* [66] and the question-answering paradigms by [67] and [65] necessitate additional, task-specific inputs such as argument and description queries, complicating their integration into diverse research settings. Similarly, [64] and *ChatIE* [68] are hindered by very limited documentation and strict data formatting requirements.

(iii) *CLEVE* [27] cannot be adapted to our definition of semantic event extraction due to its presupposition of argument type knowledge.

(iv) Frameworks like *AllenNLP* (on which *DyGIE++* [41] is built) have been discontinued, and (v) the substantial computational resources required for models like the 10-billion parameter *Deepstruct* [69] model further limit their viability.

Given these considerations, we have chosen *Text2Event* and *EventGraph* as our baselines. These methodologies have demonstrated strong performance in the event extraction task (e.g., they both outperform *Deepstruct* event classification [59, 60, 69]), provide publicly available code, and are adaptable to our task definition.

### 5.2.2. Metrics & Setting

To evaluate T-SEE’s and L-SEE’s performance on semantic event extraction, we assess their performances both on event classification and relation extraction.

We judge the accuracy of event classification using precision, recall, and  $F_1$  scores to assess if events with correct classes were extracted.

Analogously, we use the same metrics for evaluating relation extraction, where relations are only considered to be correct if connected to a correctly classified event via the correct property and to the correct entity or value.

In this section, we report the results of L-SEE as L-SEE\* since we only consider those texts for which the output of the LLM was formatted syntactically correctly and in a consistent way allowing us to use automatic evaluation. Consequently, L-SEE is not evaluated on an identical dataset as T-SEE and the baselines, but on a smaller dataset (5,602 of the full 6,407 texts for *DBpedia-SEE* and 3,958 texts of the original 5,711 for *Wikidata-SEE*).

### 5.3. Event Classification

Table 5 shows the evaluation results of T-SEE, L-SEE and the baselines on the tasks of event classification. In general, we observe that T-SEE performs well on event classification, reaching  $F_1$  scores of 0.92 and 0.85. T-SEE and *Text2Event* outperform by a notable margin *EventGraph*. While *Text2Event* performs better than T-SEE on *DBpedia-SEE*, T-SEE performs better on the more diverse *Wikidata* dataset. This performance of T-SEE may be attributed to its capability of dealing with rich event ontologies, given that *Wikidata-SEE* has three times more event classes than *DBpedia-SEE*.

The performance of L-SEE closely follows that of the baselines on *DBpedia-SEE*. However, we do see a notable drop-off in the case of *Wikidata*, likely correlated with the larger number of fine-grained event classes notable to *Wikidata*.



Table 5  
F<sub>1</sub> scores for event classification on *DBpedia-SEE* and *Wikidata-SEE*.

Approach	<i>DBpedia-SEE</i>			<i>Wikidata-SEE</i>		
	P	R	F <sub>1</sub>	P	R	F <sub>1</sub>
Text2Event	<b>0.94</b>	<b>0.94</b>	<b>0.94</b>	0.84	0.84	0.84
EventGraph	0.75	0.69	0.72	0.77	0.52	0.62
T-SEE	0.92	0.92	0.92	<b>0.85</b>	<b>0.85</b>	<b>0.85</b>
L-SEE*	0.88	0.89	0.89	0.53	0.58	0.55

Table 6  
Precision (P), recall (R) and F<sub>1</sub> scores for relation extraction on *DBpedia-SEE* and *Wikidata-SEE*.

Approach	<i>DBpedia-SEE</i>			<i>Wikidata-SEE</i>		
	P	R	F <sub>1</sub>	P	R	F <sub>1</sub>
Text2Event	0.74	0.75	0.74	<b>0.75</b>	<b>0.77</b>	<b>0.76</b>
EventGraph	0.72	0.57	0.64	0.85	0.16	0.27
T-SEE	<b>0.75</b>	<b>0.76</b>	<b>0.75</b>	<b>0.75</b>	<b>0.77</b>	<b>0.76</b>
L-SEE*	0.28	0.52	0.37	0.37	0.37	0.37

#### 5.4. Relation Extraction

Table 6 presents the relation extraction performance of T-SEE and our baselines. While it is evident that EventGraph lags behind Text2Event and T-SEE, it demonstrates a notable precision in its extractions (0.85 in *Wikidata-SEE*), albeit with a significantly lower recall (0.16). This suggests that EventGraph is highly accurate in the instances it chooses to label, but it misses many relevant relations. In contrast, T-SEE consistently matches or outperforms all baseline performances across both datasets, demonstrating its robustness in the relation extraction task.

L-SEE shows a notably lower performance compared to both T-SEE and the baselines. This is likely associated with the relatively higher complexity of the relation extraction task compared to the event classification task. It should also be noted that as L-SEE is untrained, it relies more on the natural language understanding abilities it acquired through pre-training than other evaluated baselines. As such, it is also more likely to fail when the property labels are inadequately descriptive as to what purpose they are meant to fulfil. We go into further detail about the nature of these errors and the limitations of L-SEE in Section 6.

#### 5.5. Cascading Errors

Given the sequential structure of T-SEE’s and L-SEE’s approach, where event classification precedes relation extraction, inaccuracies in the initial phase of event classification might negatively influence the subsequent relation extraction performance. Therefore, we analyse the impact of cascading errors by comparing the F<sub>1</sub> scores for relation extraction in isolated and end-to-end settings at the example of T-SEE.

For T-SEE on *DBpedia-SEE*, the isolated setting shows a precision score of 0.81, recall of 0.82, and an 0.82 F<sub>1</sub> score, indicating the model’s performance in an ideal scenario with perfect event classification. However, in the end-to-end setting, the scores decrease to the precision of 0.75, recall of 0.76, and an F<sub>1</sub> score of 0.75. This drop in performance suggests that errors in the event classification phase cascade down, as expected, affecting the model’s ability to extract relations accurately.

Similarly, on *Wikidata-SEE*, T-SEE demonstrates high scores in the isolated setting with the precision of 0.87, recall of 0.88, and an F<sub>1</sub> score of 0.88. In contrast, the end-to-end setting yields lower scores: Precision 0.75, Recall 0.77, and F<sub>1</sub> 0.76. This reduction further underscores the presence of cascading errors. However, the effect is limited and does not prevent T-SEE from outperforming or matching the baselines in both datasets.

## 5.6. Example Result

Finally, we provide example RDF triples of an event extracted with T-SEE. Fig. 7 shows the RDF triples created from an event which we extracted from the Wikipedia article “1991 Monte Carlo Open”<sup>17</sup> using T-SEE. As we can see, T-SEE successfully extracted an event of the fine-grained event class recurring tennis tournament and several relations, including properties such as season starts, located in the administrative territorial entity and part of.

```
@base <http://example.org/>.
@prefix wd: <https://www.wikidata.org/entity/>.
@prefix wdt: <https://www.wikidata.org/prop/direct/>.
@prefix so: <https://www.w3.org/2000/01/rdf-schema#>.

:E1 a wd:Q47443726 (recurring tennis tournament);
    so:description "It was ... part of the ATP Super 9 of the 1991 ATP Tour. It took place at
        the Monte Carlo Country Club in Roquebrune-Cap-Martin, France, near Monte Carlo, Monaco,
        from 22 April through 28 April 1991."@en ;
    wdt:P4794 (season starts) wd:Q118 (April) ;
    wdt:P17 (country) wd:Q142 (France) ;
    wdt:P131 (located in the administrative territorial entity)
        wd:Q45240 (Monte Carlo) ;
    wdt:P276 (location) wd:Q3861317 (Monte Carlo Country Club) ;
    wdt:P361 (part of) wd:Q300008 (ATP Tour) .
```

Fig. 7. Example of RDF triples generated from the Wikipedia article “1991 Monte Carlo Open” using the Turtle syntax.

## 5.7. Consistency Analysis

To address the variability of the LLM in generating outputs for identical inputs, we evaluate the consistency of L-SEE across multiple executions of its LLM prompts. This analysis is essential for assessing the robustness of L-SEE, as LLMs inherently introduce stochasticity due to their sampling mechanisms during generation. Specifically, we repeatedly process the same set of inputs (i.e., prompts) through the LLM under identical conditions and observe the outputs generated in each iteration. To quantify consistency, we use Fleiss’  $\kappa$ , a metric that measures inter-rater agreement [70], adapted here to measure agreement between outputs from different executions of an LLM.<sup>18</sup>

Our analysis reveals a high level of consistency for both event ontologies, as summarised in Table 7. For *DBpedia-SEE*, we observe an average Fleiss’  $\kappa$  of 0.991 for event classification and 1.000 for relation extraction, indicating near-perfect agreement across runs. For *Wikidata-SEE*, Fleiss’  $\kappa$  for event classification is 0.881, reflecting slightly reduced but still strong consistency. These scores confirm the robustness of L-SEE, which yields highly similar results across iterations.

Table 7  
Consistency analysis results of L-SEE for event classification and relation extraction.

Task	<i>DBpedia-SEE</i>	<i>Wikidata-SEE</i>
Event Classification (Fleiss’ $\kappa$ )	0.991	0.881
Relation Extraction (Fleiss’ $\kappa$ )	1.000	1.000

<sup>17</sup>[https://en.wikipedia.org/w/index.php?title=1991\\_Monte\\_Carlo\\_Open&oldid=1101607800](https://en.wikipedia.org/w/index.php?title=1991_Monte_Carlo_Open&oldid=1101607800)

<sup>18</sup>To allow consistent Fleiss’  $\kappa$  computation, in this consistency analysis, we only consider cases where valid JSON output is generated and we do not consider cases where multiples values are assigned to the same property.

The few cases demonstrating disagreement across LLM executions for event classification and relation extraction can be attributed to the design and complexity of the two tasks and the event ontologies. There is a stronger agreement for *DBpedia-SEE* compared to *Wikidata-SEE* due to the lower number of classes and properties (see Table 3): with a lower number of event classes to select from, there naturally is a higher chance of agreement. This also explains the consistency observed for relation extraction where the prompts include a small number of properties ( $P_{C_i}$ ) that have been identified to be relevant given the already detected event classes  $C_i$ .

## 5.8. Implementation

In order to implement our multilabel classification model, we leverage a pre-trained uncased BERT base model<sup>19</sup>. The model is fine-tuned for 30 epochs using the focal loss function with a gamma of 2, and the Adam optimiser, with a learning rate of  $1e-5$  and a Dropout layer with a probability of 0.3. We apply early stopping based on validation-set performance, with training capped at 30 epochs. Continuing beyond this point did not improve the validation metrics. For the relation extraction model, we utilise the same BERT model and fine-tune it on the relation extraction task. Similarly to the classification model, we train the model for 30 epochs with the Adam optimiser and a learning rate of  $3e-5$  and again employ early stopping up to 30 epochs.

`Text2Event` and `EventGraph` are trained for 40 epochs using a batch size of 30 and their original training settings.

To generate the training data, we extract Wikipedia articles using the `MWDumper`<sup>20</sup>. For entity linking, we use the `Spacy Entity Linker`<sup>21</sup>, a named entity linking tool specifically designed for Wikidata.

For L-SEE, we use `gpt-3.5-turbo-1106`<sup>22</sup>, a version of GPT-3.5 Turbo that supports a 16K context and supports improved instruction following, JSON mode, and parallel function calling. We pick this version as it has shown a 38% improvement in format following tasks such as generating JSON, XML and YAML.

## 6. Comparison of T-SEE and L-SEE

A significant finding of our evaluation is the worse performance of L-SEE compared to T-SEE on the task of relation extraction (Section 5.4). This leads to the question of whether LLMs are not suited for the task of semantic event extraction at all, in contrast to fine-tuning a transformer-based architecture. To answer this question, this section delves into a manual evaluation and a multifaceted error analysis, followed by a discussion.

### 6.1. Manual Evaluation

In this section, we aim to understand the differences between the two paradigms of transformer-based architecture versus using LLMs for semantic event extraction. Therefore, on top of the automatic evaluation performed in Section 5, we perform a comparison of T-SEE and L-SEE based on a manually annotated subset of the test dataset used in the automatic evaluation.

We create *DBpedia-SEE*<sub>100</sub> – a subset of *DBpedia-SEE* with 100 randomly selected texts, their events and relations. We ensure that L-SEE successfully performs semantic event extraction on these texts without syntactical errors. For each text in *DBpedia-SEE*<sub>100</sub>, we manually annotate the semantic event representations generated by T-SEE and by L-SEE with respect to each other and the ground truth. For example, given a text  $t$ , if T-SEE generates a relation  $r$  that is not in *DBpedia-SEE*<sub>100</sub>, we manually assess whether  $r$  is correct and expressed in  $t$ . If this assessment is positive and  $r$  is also missing in L-SEE, we denote a true positive for T-SEE and a false negative for L-SEE.

<sup>19</sup><https://huggingface.co/bert-base-uncased>

<sup>20</sup><https://www.mediawiki.org/wiki/Manual:MWDumper>

<sup>21</sup><https://github.com/egerber/spaCy-entity-linker>

<sup>22</sup><https://platform.openai.com/docs/models/gpt-3-5-turbo>

Table 8 shows the results of evaluating T-SEE and L-SEE on *DBpedia-SEE*<sub>100</sub> before and after our manual assessment. The results before manual assessment confirm our results given in Table 5, where T-SEE and L-SEE both perform well on event classification ( $F_1$  scores of 0.92 and 0.89), but L-SEE is clearly outperformed by T-SEE for relation extraction ( $F_1$  scores of 0.72 and 0.39), mainly due to 200 false positive extracted relations. This indicates a considerably better ability of T-SEE to accurately identify and categorise relationships within the data under controlled conditions.

Table 8

Evaluation of T-SEE vs. L-SEE on *DBpedia-SEE*<sub>100</sub> before and after manual assessment. TP: true positives, FP: false positives, FN: false negatives.

Task	Approach	TP		FP		FN		$F_1$	
		Before	After	Before	After	Before	After	Before	After
Event Classification	T-SEE	92	90	7	12	10	9	<b>0.92</b>	0.90
	L-SEE	91	100	11	2	11	2	0.89	<b>0.98</b>
Relation Extraction	T-SEE	83	103	33	23	32	114	<b>0.72</b>	0.58
	L-SEE	77	178	200	99	38	29	0.39	<b>0.74</b>

After manual assessment, L-SEE shows a remarkable improvement in event classification, achieving an almost perfect  $F_1$  score of 0.98, suggesting that with manual verification of the ground truth, the LLM's capabilities are more effectively utilised. Regarding relation extraction, while L-SEE improves performance ( $F_1$  score of 0.74), T-SEE experiences a significant drop in effectiveness ( $F_1$  score of 0.58), indicating challenges in adapting to the intricacies of manually annotated samples and the complexity of real-world data.

These results underscore the strengths and limitations of both methodologies. While T-SEE demonstrates superior performance in a controlled environment, particularly in relation extraction tasks, L-SEE shows remarkable adaptability and potential in handling complex, real-world scenarios when supplemented with manual verification and annotation processes: not being constrained by any limitations in training data, L-SEE is able to extract more than double the amount of relations. This highlights the importance of context and the level of detail in ground truth annotations when evaluating and comparing data extraction methodologies.

## 6.2. Error Taxonomy

To understand the differences in behaviours between T-SEE and L-SEE, we manually annotate the specific errors that occur when performing semantic event extraction on *DBpedia-SEE*<sub>100</sub>. While doing so, we create an error taxonomy presented in this section. Later, to contextualise said error taxonomy, we present examples of generated RDF triples and the errors in them.

Our manual annotation process has unveiled a structured classification of errors, which we have divided into three principal categories:

**Extraction Inaccuracies** Errors arising from the model's inability to accurately interpret information within texts:

- *Omissions or Missing Events/Relations*: The event or its relations are not extracted.
- *Type misalignment*: An inappropriate type of entity or value is selected for a given property.
- *Granularity mismatch*: The model's predictions lack the specificity of the ground truth, e.g., categorising an event broadly as `dbo:SportsEvent` rather than the more specific `dbo:TennisTournament`.
- *Erroneous extraction*: The extraction of incorrect properties or values, leading to a misrepresentation of the factual content.

**Annotation Discrepancies** Errors stemming from inconsistencies, errors or omissions in the ground truth:

- *Imprecise event class*: The model's predictions provide a more detailed event classification.
- *Imprecise property*: The model predicts property values with greater accuracy than the ground truth, such as specifying the exact match score when the ground truth only acknowledges the victory.

- *Annotation error*: The presence of omissions or inaccuracies within the ground truth itself, such as neglecting to annotate the specific date of a match or other pertinent details.

*Other Anomalies* Errors arising from other sources:

- *Event ambiguity*: The model struggles to distinguish between multiple distinct events described within a single sample, which may lead to conflated or mixed property assignments.
- *Processing error*: T-SEE and L-SEE match spans of text to specific entities, relying on an external entity linking component and a date parsing module which are prone to errors.

### 6.2.1. Examples of Errors

We provide four semantic event representations generated by L-SEE as examples of the identified error types in the error taxonomy. For each of the examples, we provide the input text  $t$ , selected RDF triples describing an event  $e$  in the ground truth as well as selected triples generated by L-SEE.<sup>23</sup> Errors are marked in red, relations only in the ground truth are marked in blue, and relations only in the prediction are marked in green.

**Example 1** (Fig. 8) – ground truth extracted from  $dbr:Black\_Monday\_ (1360)$ :

- *Omission*: L-SEE failed to extract the `dbo:commander` relation.
- *Annotation error*: On the other hand, L-SEE accurately extracts a relevant date and territory for the event, however, these are not contained within the ground truth.

**Text:** This was in part caused by Black Monday (1360), the freak storm that devastated the English army and forced Edward III into peace talks.

#### Ground Truth

`:MilitaryConflict1` `dbo:commander` `dbr:Edward_III_of_England`.

#### Prediction

`:MilitaryConflict1` `dbo:date` `"1360-01-01"^^xsd:date`;  
`dbo:territory` `dbr:England`.

Fig. 8. Example of an omission error and an annotation error.

**Example 2** (Fig. 9) – ground truth extracted from  $dbr:Al-Qusayr\_offensive$ :

- *Type misalignment*: The commanders are incorrectly identified and assigned to group entities instead of individuals. Specifically, L-SEE detects two commanders extracting "Syrian Army" and the "Lebanese militia Hezbollah".
- *Processing error*: in the entity linking process, "Lebanese militia Hezbollah" is wrongly linked to three entities.

<sup>23</sup>For brevity, we skip  $p_{type}$  and  $p_{description}$  relations. The event class is indicated by its URL (e.g., `:MilitaryConflict1` is an event classified as `dbo:MilitaryConflict`).

**Text:** The second of two battles in al-Qusayr started on 19 May 2013, as part of the larger al-Qusayr offensive, launched in early April 2013 by the Syrian Army and the Lebanese militia Hezbollah, during the Syrian civil war, with the aim of capturing the villages around the rebel-held town of al-Qusayr and ultimately launching an attack on the town itself.

#### Ground Truth

```
:MilitaryConflict2 dbo:place dbr:Al-Qusayr,_Syria ;
                        dbo:isPartOfMilitaryConflict dbr:Syrian_Civil_War .
```

#### Prediction

```
:MilitaryConflict2 dbo:place dbr:Al-Qusayr,_Syria ;
                        dbo:commander dbr:Syrian_Army ;
                        dbo:commander dbr:Lebanon ;
                        dbo:commander dbr:Militia ;
                        dbo:commander dbr:Hezbollah ;
                        dbo:isPartOfMilitaryConflict dbr:Syrian_Civil_War .
```

Fig. 9. Example of type misalignment and processing errors.

**Example 3** (Fig. 10) – ground truth extracted from *dbr:2016\_Wuhan\_Open*:

- Imprecise event class: L-SEE identifies a more precise event class (*dbo:TennisTournament* versus *dbo:Tournament*).
- Erroneous extraction and event ambiguity: The same tennis tournament did not happen in Wuhan and in Beijing; L-SEE fails to distinguish between the tournaments Wuhan Open and China Open.

**Text:** However, she rebounded in the Asian swing by reaching the quarterfinals of Wuhan and the semifinals of Beijing.

#### Ground Truth

```
:Tournament1 dbo:location dbr:Wuhan .
```

#### Prediction

```
:TennisTournament1 dbo:location dbr:Wuhan ;
                    dbr:location dbr:Beijing .
```

Fig. 10. Example of an imprecise event class, an erroneous extraction and event ambiguity.

**Example 4** (Fig. 11) – ground truth extracted from *dbr:1959\_Ontario\_general\_election*:

- Event ambiguity: The date "1961-01-01" indicates confusion between multiple events. Specifically, this is because the event annotated in the ground truth is derived from the link tied to the string "previous election", referring to *dbr:1959\_Ontario\_general\_election*.
- Erroneous extraction: The use of *dbo:secondLeader* to indicate a chronological successor is highlighted in red, illustrating a misunderstanding of the property, as *dbo:secondLeader* is meant to instead describe second ranking in a competition.
- Annotation error: The relation using the *dbo:affiliation* property is missing in the ground truth.

**Text:** The Ontario Progressive Conservative Party, led by John Robarts, who had replaced Leslie Frost as PC leader and premier in 1961, won a seventh consecutive term in office, and maintained its majority in the legislature, increasing its caucus from the 71 members elected in the previous election to 77 members in an enlarged legislature.

#### Ground Truth

```
:Election1 dbo:country dbr:Ontario ;
          dbo:firstLeader dbr:Leslie_Frost .
```

#### Prediction

```
:Election1 dbo:startDate "1961-01-01"^^xsd:date ;
          dbo:country dbr:Ontario ;
          dbo:secondLeader dbr:John_Robarts ;
          dbo:affiliation dbr:Progressive_Conserv._Party_of_Canada ;
          dbo:firstLeader dbr:Leslie_Frost .
```

Fig. 11. Example of event ambiguity, erroneous extraction and an annotation error

### 6.3. Error Analysis

On the basis of our error taxonomy, we annotated each semantic event representation generated by T-SEE and L-SEE with the set of errors occurring in them.

First, we categorise errors into extraction inaccuracies, annotation discrepancies, and other anomalies to clarify our approaches' error landscapes. Fig. 12 visualises these error profiles for T-SEE and L-SEE, highlighting the challenges in semantic event extraction. In general, we register fewer errors for T-SEE than L-SEE across all three error categories, which results from T-SEE's capability to mimic the dataset characteristics. On the other hand, we annotate 180 annotation discrepancies for L-SEE, more than its 107 extraction inaccuracies. Since annotation discrepancies represent cases where the model extracts valid triples which are not covered in the ground truth, this analysis demonstrates how L-SEE is capable of semantic event extraction without being closely attached to the characteristics of training data and, implicitly, the data coverage in the target knowledge graph.

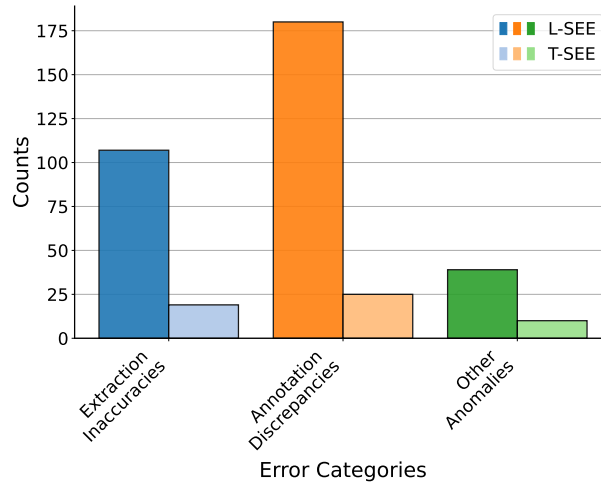


Fig. 12. Distribution of error categories for T-SEE and L-SEE.

Fig. 13 provides a detailed analysis of the error types. As can be seen in the figure, different error types manifest with varying frequencies across L-SEE and T-SEE.

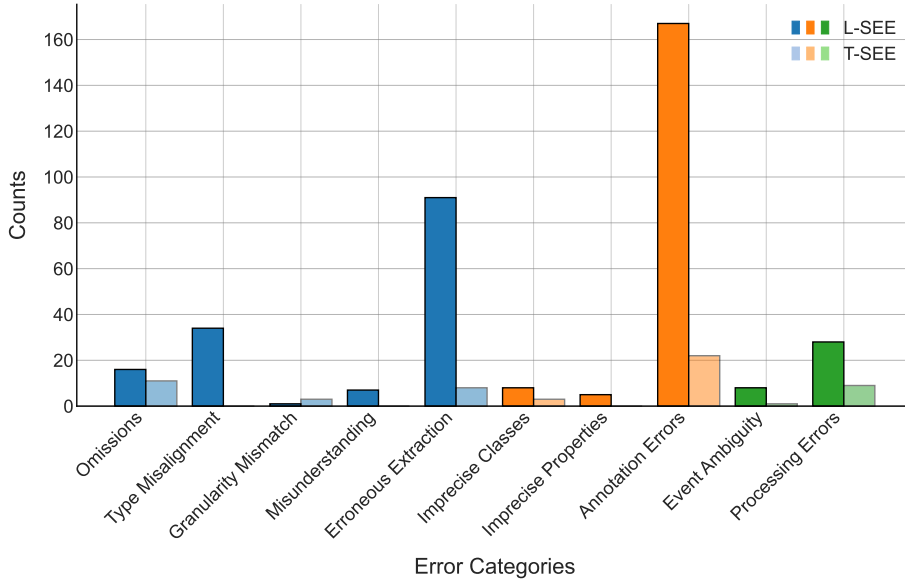


Fig. 13. Distribution of error types for T-SEE and L-SEE grouped by category.

For L-SEE, the most prevalent error type are *Annotation Errors*, with a count of 167, reflecting instances where L-SEE identifies relevant semantic relations not annotated in the ground truth. Following closely is the *Erroneous Extraction* category, with 91 instances, which encompasses errors related to the incorrect identification of properties or values. A notable portion of these errors can be attributed to *Processing Errors*, amounting to 28 instances, where the entity linking and date extraction methods utilised by L-SEE falter in accurately extracting dates or correctly linking entities, leading to inaccuracies in the relation extraction.

*Misunderstanding* and *Type Misalignment* errors, with 7 and 34 instances respectively, further contribute to the *Erroneous Extraction* count. These errors emerge when L-SEE misinterprets the intended meaning of properties or incorrectly aligns relations with inappropriate entities or values. For instance, the common misunderstanding of the `dbo:secondLeader` property (Example 4) exemplifies how more careful prompting approaches may lead to better performance. For an example of a type misalignment error, we may observe instances where in the absence of precise time expressions in the text, L-SEE assigns imprecise values such as "yesterday" to date relations.

Conversely, T-SEE demonstrates a lower overall error frequency, with *Annotation Errors* again emerging as the dominant error type, albeit with a substantially lower count of 22. This suggests a more precise alignment with the ontology. Notably, T-SEE exhibits no *Type Misalignment* errors and only 1 *Event Ambiguity* error. However, both methodologies encounter *Omissions* and *Processing Errors*, with L-SEE facing 16 and 28 instances respectively, and T-SEE experiencing 11 and 9 instances.

Upon a more nuanced examination, especially after correcting for annotation errors, the performance landscape shifts. Initially, T-SEE appears to outperform L-SEE due to its lower error rates. However, this might also indicate a tendency of T-SEE to conform to the existing annotations, potentially overlooking unlabelled but present relations. This could imply that while T-SEE is more aligned with the given annotations, it may also be less inclined to explore beyond them, possibly fitting to annotation noise rather than capturing the full spectrum of semantic relations.

In summary, while T-SEE shows precision in alignment with the current ontology, L-SEE's broader extraction attempts, despite higher initial error rates, may offer a more comprehensive understanding of the underlying semantic structures, especially when considering the corrected annotation context. This dichotomy highlights the balance between precision and recall in semantic event extraction and underscores the importance of continuous refinement in both methodologies to enhance their efficacy and reliability.



### 6.3.1. Formatting and Ontology Errors

As indicated in Section 5.2.2, for 805 of the texts in the complete dataset *DBpedia-SEE*, L-SEE could not generate RDF triples due to formatting issues, including:

- Misformatted output: The LLM-generated JSON strings of 606 texts were not in proper JSON syntax and could not be processed.
- Non-existing event classes: In 401 cases, an event class was identified which is not part of the event ontology and the prompt. An example is the extraction of an event typed as `dbo:SyrianCivilWar`, while only `dbo:CivilWar` exists in the DBpedia ontology.
- Invalid properties: In 1,191 cases, a property was identified which is not part of the event ontology (e.g., `dbo:percentageOfPopularVote`, `dbo:delayReason`). Despite these being errors, they often demonstrate L-SEE capability to suggest relevant attributes for specific scenarios adaptively.

### 6.4. Effect of Text Characteristics on Semantic Event Extraction

To get a sense of L-SEE performance across a variety of syntactic and semantic phenomena, we dissected *DBpedia-SEE* into multiple subsets, each representing distinct text characteristics. The subsets are generated employing specific strategies, each tailored to highlight a particular aspect of the dataset, ranging from event co-occurrences to the complexity of the document structure.

#### 6.4.1. Text Characteristics

We employ a collection of strategies to generate meaningful subsets of the dataset, each aimed at isolating different factors that could influence L-SEE's performance. Specifically, we ranked the entire dataset based on the presence and frequency of certain linguistic, syntactic, or semantic phenomena. From this ranking, we then selected the top 100 samples for each subset to focus our analysis on the most pronounced examples of each phenomenon.

**Semantic Diversity:** We assess samples for semantic diversity. The semantic diversity of a text is measured by the variety of verb phrases and their arguments, approximated by the count of unique verb lemmas in the text. Samples with high semantic diversity are chosen for this subset, aiming to test the model's understanding of varied semantic contexts and its ability to extract a broad range of event semantics.

**Sentence Length:** This strategy sorts the samples by the length of the text. Samples are then selected from the sorted list, prioritising those with the longest texts.

**Geographical Diversity:** Samples of this subset are generated based on the count of geographical entities identified by the *Spacy* NLP pipeline (i.e. "GPE" and "LOC" labelled entities) in each text. To assess the model's proficiency in dealing with texts containing diverse geographical references, we select samples with the highest counts of such entities.

**Temporal Event Distribution:** We identify texts with temporal expressions using the *Spacy* library and extract where they are most frequently occurring. As temporal expressions can be crucial for event understanding, this subset evaluates L-SEE's capability to understand and integrate temporal information.

**Named Entity Diversity:** This subset focuses on the diversity of named entities. We again utilise the *Spacy* library to extract named entities and then sort and select samples with the widest range of entities. This subset tests the model's ability to accurately recognise and categorise entities in the context of events.

**Complex Sentence Structures:** Samples with intricate syntactical constructions are selected to challenge L-SEE's parsing abilities, as complex structures can obscure event boundaries and relations, making extraction more difficult. This set is generated by measuring the depth of the syntactic parse tree of each text, with depth representing the maximum distance from any token to the root of the tree. Samples with the most complex sentence structures, i.e., the deepest parse trees, are selected for this subset.

In the following, we detail the outcomes of this analysis, demonstrating L-SEE's efficacy and limitations across varying text characteristics.

### 6.4.2. Results

L-SEE's performance was evaluated across the subsets using precision, recall, and F<sub>1</sub> scores for both event classification and relation extraction. Figures 14 and 15 show the results of this analysis, detailed in the following.

For comparison, we also include the full dataset performance in our analysis. The distinctly strongest relation extraction performance on the full dataset suggests that we have successfully sampled parts of our data that L-SEE finds difficult to deal with.

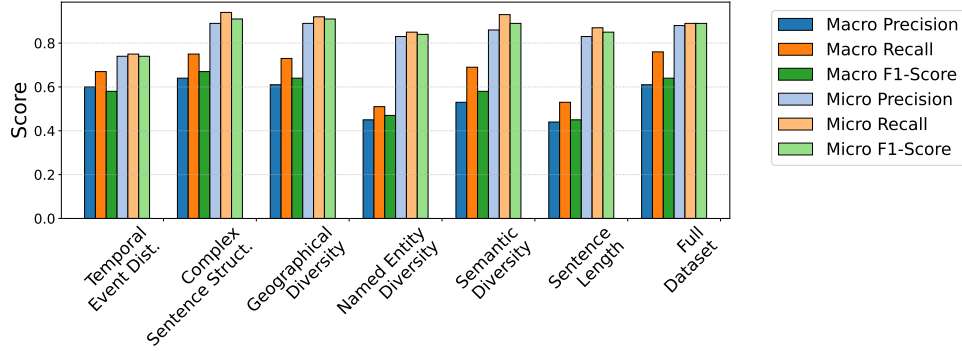


Fig. 14. L-SEE performance in event classification across various data subsets.

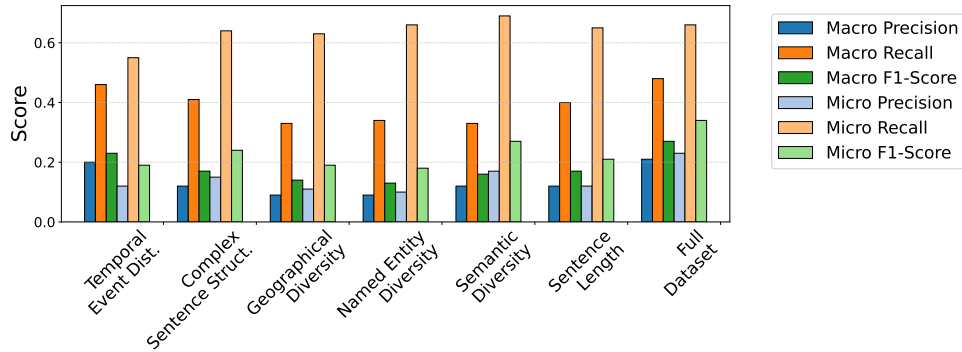


Fig. 15. L-SEE performance in relation extraction across various data subsets.

**Semantic Diversity** L-SEE exhibits moderate performance in the semantic diversity subset, with macro metrics reflecting challenges in consistently classifying a wide range of semantically varied events. However, micro metrics indicate better performance in frequent semantic contexts. The significantly lower performance in relation extraction highlights that there may be difficulties in mapping complex semantic relationships accurately.

**Sentence Length** The results suggest that longer sentences pose significant challenges, with lower macro metrics for event classification and even more pronounced difficulties in relation extraction. This indicates that L-SEE may have limitations in maintaining context and coherence over sentences.

**Geographical Diversity** L-SEE performs relatively well in event classification, suggesting a good grasp of geographical contexts. However, the lower relation extraction scores point to challenges in accurately extracting relationships when the diversity of geographical entities is high.

**Temporal Event Distribution** L-SEE displays reduced event classification accuracy in this subset. However, relation extraction metrics remain comparatively stable. Given the lower event classification performance, this

stability in relation extraction, despite potential cascading errors from misclassifications, may indicate a relatively stronger inherent capability of the model in isolating and extracting relations in the context of temporally complex texts.

**Named Entity Diversity** We observe notable difficulties, highlighting L-SEE’s struggle with diverse named entities. The discrepancy between macro and micro metrics points to L-SEE’s poorer handling of less frequent event types when the diversity of entity mentions is high.

**Complex Sentence Structures** L-SEE’s strongest performance is observed in the subset with complex sentence structures, indicating effective parsing of intricate syntactic constructions. However, the lower relation extraction metrics suggest that while L-SEE can identify events within complex sentences, accurately extracting the relations remains challenging.

This analysis underscores L-SEE’s strengths in contextual integration and syntactic navigation but also points at significant areas for improvement. Future work should conduct further analysis focusing on how improved prompting strategies may help in robustness and accuracy across these diverse linguistic and contextual scenarios.

## 6.5. Discussion

With our comparison of T-SEE and L-SEE in this article, we aim at a deeper understanding of the suitability of two different paradigms – transformer-based architectures and the use of LLMs – for semantic event extraction. Following our evaluation results and analysis, we identify five core phenomena to be considered when deciding between these paradigms:

**Mimicry of dataset characteristics:** Our analyses, e.g., in Table 6, Table 8 and Fig. 12, clearly demonstrate that the results of methodologies fine-tuned on the target datasets (T-SEE, Text2Event and EventGraph) are much more aligned to the expected RDF triples in the test sets than triples generated by an LLM (L-SEE). From this behaviour, we infer the following: (i) In a controlled setting where mimicking the characteristics of the training data is desired, transformer-based approaches are preferable over the use of LLMs. (ii) However, transformer-based approaches also mimic the flaws of the target datasets and knowledge graphs. For example, if a specific property is rarely used in an event knowledge graph but still valuable, L-SEE would identify it, while fine-tuned approaches might miss it. An example is the property `dbo:country` on the event class `dbo:Election`, which is only used in approximately 25% of DBpedia’s `dbo:Election` events.

**Distantly-labelled datasets:** Training a transformer-based architecture requires the availability of large training data, i.e., texts annotated with RDF triples. Therefore, we opted for the automated extraction of two new datasets. The use of distantly-labelled datasets without human annotations such as *DBpedia-SEE* and *Wikidata-SEE* for semantic event extraction or datasets for relation extraction [52, 71–73] overcomes the issues of training data dimensionality but always comes with questions regarding dataset quality.<sup>24</sup> Specifically, we identified a large number of false positives when evaluating L-SEE (Table 8), resulting from incomplete knowledge graphs or faulty alignment between texts and RDF triples in the distant labelling process. Consequently, the evaluation of different approaches on a distantly-labelled dataset requires careful investigation of the outputs beyond solely providing scores of the evaluation metrics.

**Ontology Guidance:** We took care of carefully guiding both our approaches through our event ontologies. By fine-tuning a transformer-based architecture, adherence to the ontology can be enforced, e.g., by explicitly classifying into the event classes pertinent to the event ontology. For an LLM, in contrast, while we prompted for the specific event classes and properties, we still observed cases of invalid event classes or properties as discussed in Section 6.3.1. Also, our examples demonstrated cases of type misalignment and a misunderstanding of the semantic definition of a property (`dbo:secondLeader` in Example 4), demonstrating the need to control the outputs of an LLM. The improvement in the precision of LLM-based semantic event extraction is a major future direction for LLM-based semantic event extraction, e.g., through the provision of property descriptions within the prompt.

<sup>24</sup>Note that even human annotators frequently disagree when providing annotations for NLP tasks [74].

**Complexity:** Setting up a transformer-based architecture and its fine-tuning requires the availability of rich training data, computing and time resources. Setting up an LLM, in contrast, requires access to an LLM and careful prompt engineering, i.e., potentially easier-to-obtain resources.

**Real-world applications:** Given the capability of LLMs to adapt to different inputs and data characteristics, we assume that LLM-based approaches are well-suited under more complex, real-world conditions and to explore low-resource scenarios.

## 7. Related Work

Knowledge graphs have, as a form of structured human knowledge, drawn a lot of research attention from both academia and the industry [75]. With a great deal of event information worldwide, it is essential to bring entities and events together through event-centric knowledge representations [24], with event extraction and relation extraction being key technologies for accessing event knowledge [14].

### 7.1. Event Knowledge Graphs

Event knowledge graphs represent knowledge about happenings with societal impact in an event ontology and interlink them with connected entities [24]. We distinguish between two types of event representations as follows:

- Named events: The predominantly entity-centric information of popular cross-domain knowledge graphs such as DBpedia, YAGO, and Wikidata represent events as *named events* such as “Brexit” and “World War II”. Named events are also the core component of EventKG [76], a multilingual event-centric temporal knowledge graph, part of the Open Event Knowledge Graph [77] that integrates event-related data sets from multiple application domains. GDELT [57] and ICEWS are two datasets of global political events encoded using the CAMEO framework [78], i.e., not in RDF.
- Unnamed events: Works that address *unnamed events* specifically deal with the identification of texts describing events and with the semantic annotation of these texts. For example, Rospocher et al. [22] build knowledge graphs from news articles, and Zhang et al. [79] develop a large-scale English event knowledge graph extracted from several sources such as reviews, news, and social media. For the task of event modelling, [80] proposes a weakly-supervised approach to extract event relation tuples from text and build an event knowledge base, not focusing on event-entity relations.

All event knowledge graphs require the availability of an event ontology, with popular examples including *LODE* [6], the *Simple Event Model* [7] and more as discussed by Pyriani et al. [81]. Relevant patterns for event representation are presented in [82, 83], focusing on the spatio-temporal extent of events, the role of their participants and recurring events. In this article, we extracted event ontologies from their vocabularies to allow the population of the well-established cross-domain knowledge graphs DBpedia and Wikidata.

With T-SEE, we aim to bring together the complementary strengths of the Semantic Web and NLP perspectives by performing event extraction that can be adapted to different event ontologies.

### 7.2. Event Extraction

Event extraction (EE) is a critical task in constructing and populating entity-centric knowledge graphs, with recent advancements significantly diversifying the methodologies employed [30, 59, 84]. Earlier approaches have relied on sentence-level pipelines for extracting event triggers and their corresponding argument roles [65, 85, 86], employing sequence-to-structure generation paradigms like *Text2Event* [59] and multi-task frameworks such as *DyGIE++* [41], which utilise contextualised embeddings and dynamic span graph updates. Other studies have extended the scope to document-level EE [87, 88] or ventured into open-domain EE without predefined event classes [89, 90], which, while broadening the applicability, faces challenges due to the absence of a well-defined event ontology.

Innovations in the field have introduced contrastive pre-training frameworks like *CLEVE* [27], which capitalise on large unsupervised datasets and their semantic structures to enhance EE’s efficacy, demonstrating marked improvements in both supervised and unsupervised settings. Similarly, *EventGraph* [60] has presented a joint framework that

conceptualises events as graphs, facilitating the simultaneous detection and extraction of multiple events and their intricate interrelations, thereby achieving state-of-the-art results in event trigger and argument role classification.

*Deepstruct*, on the other hand, tries to leverage the structural understanding capabilities of language models through task-agnostic pretraining, allowing for zero-shot knowledge transfer across a wide array of structure prediction tasks and setting new benchmarks on numerous datasets [69]. With *DEGREE* [66], authors propose a data-efficient, generation-based model for EE that capitalises on semantic guidance from manually designed prompts and the joint prediction of triggers and arguments, showcasing robust performance in low-resource settings. *ONEEE* [91], on the other hand, utilises a one-stage framework for fast overlapping and nested event extraction.

A notable shift in EE methodology is the adoption of a question-answering paradigm [65], which mitigates the prevalent issue of error propagation seen in conventional approaches by facilitating end-to-end argument extraction, including for roles not encountered during training. Following this line, *QGA-EE* [67] has refined the QA-based approach by integrating context-aware question generation, thus accommodating multiple arguments for identical roles and surpassing prior single-task models in performance metrics.

In light of the new methodologies and progress in event extraction, the research community has also focused on the specific subtasks of event extraction. For example, with *PAIE* [92], authors devise a prompt tuning approach to document-level event argument extraction similar to the already established question-answering paradigm in event extraction work. Older work on event argument extraction, such as *HMEAE* [18], a hierarchical approach to argument extraction utilising concept correlation among argument roles, have, in turn, inspired approaches such as *DEGREE* that aim to resolve issues such as poor handling of the encoding of the labels semantics and other weak supervision signals.

Prompt-based approaches have been explored for event argument extraction, leveraging the ability of pre-trained language models to generate structured outputs. For example, Peng et al. (2024) propose *Event Co-occurrences Prefix Event Argument Extraction (ECPEAE)*, which incorporates co-occurrence information of multiple events in a sentence to improve argument extraction accuracy [93]. This method uses a co-occurrence event prefix module to encode template information for all events in the input, enabling the model to leverage causal relationships between events. While ECPEAE focuses on sentence-level event interactions, *T-SEE*'s pipeline explicitly integrates event ontologies and RDF triples for knowledge graph population, aligning with broader semantic event extraction goals. Other recent work includes *Hyperspherical Multi-Prototype Learning* [94], which enhances event argument extraction via optimal transport.

The other subtask of event extraction, event detection, has also received attention with the *DRC* framework [95] trying to compete with trigger-based models as a way of exploring methods of event detection robust to less annotated real-world domains, an area we examine in our work as well. Similarly, recent work has introduced retrieval-augmented prompting for event detection, leveraging LLMs to improve performance in both high- and low-resource settings [44]. This approach constructs automatic retrieval-augmented prompts to provide LLMs with structured extraction guidelines, enhancing their ability to detect events without relying solely on trigger words. These advancements align with our exploration of methods for event detection in less-annotated domains.

Other research exploring ontology and schema-based approaches to event extraction has yielded promising innovations. Notably, Huang et al. (2024) introduce a multi-graph representation for event extraction, using graph neural networks to model event interactions and improve extraction accuracy [96]. This graph-based approach contrasts with *L-SEE*'s pipeline structure, which prioritises ontology-guided classification and relation extraction. Shiri et al. (2024) propose a schema-aware event extraction method using LLMs, decomposing the task into event detection and argument extraction while incorporating dynamic schema-aware retrieval examples [44]. This approach uses dynamic retrieval to fetch task-specific examples for each query, enhancing LLM understanding but requiring external data. In contrast, *L-SEE* employs dynamic prompt generation from a static ontology, enabling targeted extraction of event classes and properties without relying on retrieval mechanisms. Our approach achieves a critical balance: leveraging the flexibility of ontology-driven prompting while maintaining simplicity through independence from external data sources. Similarly, *COFEE* [97] uses a static ontology for schema-guided augmentation in supervised models. Unlike *L-SEE*, which leverages LLMs' contextual understanding with ontology-guided prompts, *COFEE* relies on static schema augmentation. *L-SEE*'s ontology-driven prompting enables adaptability to diverse event types while avoiding the computational overhead of retrieval-augmented generation.

These developments reflect a broader trend towards more adaptable, efficient, and comprehensive models for event extraction, underlining the field's evolution towards leveraging advanced language model capabilities and innovative problem-solving frameworks.

### 7.3. LLM-based Information Extraction

The field of Information Extraction (IE) has traditionally relied on rule-based and statistical methods to extract structured information from text. However, the emergence of Large Language Models (LLMs) has opened up new avenues for tackling IE tasks with remarkable capabilities in understanding and generating natural language. This section reviews recent advancements in using LLMs for IE, particularly focusing on unstructured information extraction and event extraction.

#### General Information Extraction with LLMs

A few years ago, LLMs were still in their early stages of development, with limited capabilities for tackling complex tasks like information extraction. While early works explored LLM-based approaches for IE (e.g., [98]), these models faced challenges due to limited model capacity, data inefficiency, and limited adaptation. However, significant advancements in recent years have addressed these challenges, driven by the rise of the transformer architecture [33] enabling long-range dependencies. Large-scale pre-training pushed things further with BERT [36] and GPT-3 [99], allowing LLMs to learn general language understanding capabilities and adapt to specific IE tasks through fine-tuning with smaller labelled datasets. Finally, the growing availability of powerful computing resources like GPUs and TPUs [100] has enabled the training of larger and more complex LLM models, further enhancing their ability to handle complex information extraction tasks.

#### Unstructured Information Extraction with LLMs

In 2022, Dunn et al. showed how a pre-trained LLM can extract structured information from scientific abstracts [32]. In 2023, Polak et al. [101] expanded on the early promises of unstructured information extraction with *ChatExtract*, demonstrating that a significant amount of up-front effort, expertise, and coding may be fully automated using an advanced conversational LLM. By leveraging prompts and follow-up questions, *ChatExtract* achieves high accuracy and efficiency in extracting materials data, showcasing the potential of LLMs for automated knowledge extraction from scientific literature.

In the same year, Wei et al. proposed *ChatIE*, a multi-turn QA framework for zero-shot information extraction demonstrating good performance across a number of datasets, three tasks, and two languages [68]. Li et al. systematically analysed *ChatGPT* across seven detailed information extraction tasks [102] including event extraction. The authors show that while *ChatGPT* underperforms in standard IE tasks compared to BERT-based models, it excels in OpenIE settings, as confirmed by human evaluators. However, a notable concern is the model's overconfidence in its predictions, leading to calibration issues. This is further confirmed in the comprehensive survey by Liu et al. [103], in which the authors evaluate the capabilities and applications of *ChatGPT* (versions 3.5 and 4) against the backdrop of current state-of-the-art models in natural language processing. The paper highlights *ChatGPT*'s advancements in large-scale pre-training, instruction fine-tuning, and reinforcement learning from human feedback, which collectively enhance its adaptability and performance across a myriad of NLP tasks. A detailed comparison of *ChatGPT* with existing state-of-the-art models reveals that while *ChatGPT* excels in multitask learning and shows promising results in some NLP tasks, it falls short in multilingual capabilities and specialised tasks when compared to dedicated models. Moreover, stability and consistency emerge as areas where *ChatGPT* does not yet match the performance levels of state-of-the-art models, which could impact its reliability in critical applications.

#### 7.3.1. LLM-based Event Extraction

LLMs have recently been utilised for the task of event extraction. In general, as already mentioned, Li et al. [102] evaluate the performance of *ChatGPT* on a number of information extraction tasks, revealing an increasingly worse performance as the complexity of the evaluated task increases, where the worst performance is reported on the task of event extraction.

A comparison between LLMs and traditional methods have been conducted on several tasks related to EE: In [104], authors explore prompt-based learning with *GPT-4* for detecting factual events in literary narratives. The

study concludes that while BiLSTM with BERT embeddings excels in event detection within literary texts, *GPT-4* shows promise in prompt-based learning approaches, particularly in few-shot settings. Sharif et al. [105] conducted an in-depth analysis of *ChatGPT*'s performance on the task of characterising information-seeking events, where *ChatGPT* underperformed compared to transformer models like *XLNet*, especially in domain-specific contexts requiring extensive knowledge.

Zhan et al. introduce *GLEN* [106], a large-scale general-purpose event detection dataset that significantly expands the ontology of event types. While *InstructGPT* underperformed compared to other baselines in their experiments, the authors attribute this to the limited input length and lack of fine-tuning, with only 57.8% of generated event types matching the ontology, similarly to our observations (Section 6.3.1). In 2024, Zhang et al. present *ULTRA* [107], a framework utilising hierarchical modelling and pairwise refinement for document-level event argument extraction.

Peng et al. (2024) introduce *CsEAE*, a model that combines small language models (SLMs) and LLMs for document-level event argument extraction [108]. *CsEAE* incorporates co-occurrence-aware and structure-aware modules to handle semantic boundaries between events and reduce interference from redundant information. The authors also demonstrate that insights from SLMs can enhance LLM performance via supervised fine-tuning and prompt engineering. This work aligns with L-SEE prompting strategies and highlights the potential for co-occurrence-aware designs to improve LLM-based event extraction.

Liu et al. (2024) propose *EventRL*, a framework that enhances LLM-based event extraction using reinforcement learning with outcome supervision [109]. *EventRL* improves extraction accuracy by rewarding the LLM based on its alignment with human-annotated triggers and arguments. While *EventRL* focuses on refining LLM outputs through external feedback, L-SEE leverages ontology-guided prompting to structure LLM responses internally. Both approaches aim to improve LLM reliability, but *EventRL* uses post-hoc correction, while L-SEE prioritizes upfront guidance. These methods address complementary aspects of LLM-based extraction: *EventRL* mitigates hallucinations via feedback, while L-SEE ensures semantic consistency with knowledge graphs through ontology integration.

While the early attempts at utilising LLMs for the complex task of event extraction have shown mixed results, with LLMs often underperforming in comparison to traditional methods, especially in domain-specific contexts, there is a clear trajectory of improvement. As LLMs continue to evolve, gaining the ability to handle larger context windows and as researchers refine their prompting techniques — such as breaking down the task into simpler sub-tasks as demonstrated in L-SEE — the gap between LLMs and traditional methods is expected to narrow. The advancements in hierarchical modelling, pairwise refinement, and modules like LEAFER [107] for argument span refinement indicate the potential for LLMs to improve and catch up to traditional event extraction methodologies in the near future.

### 7.3.2. LLM-based Knowledge Graph Population

The use of LLMs for the population of knowledge graphs has also been explored recently. For example, Mihindukulasooriya et al. experimented on ontology-driven triple extraction from sentences [110], while Yao et al. performed instruction tuning for the tasks of triple classification, relation prediction and entity link prediction [111]. In another innovative approach, *AutoKG* leverages a multi-agent-based approach employing LLMs and external sources for KG construction and reasoning [112]. Zhang et al. propose *KoPA*, which ingests entity and relation embeddings into LLMs [113].

These papers about LLM-based information extraction present a glimpse into the rapidly evolving field of LLM-based IE. While promising results have been achieved, further research is needed to address challenges such as factual correctness, bias mitigation, and adapting LLMs to specific domains and tasks. As research progresses, LLMs are set to play a key role in the future of information extraction, enabling efficient and accurate knowledge extraction from vast amounts of unstructured text data.

## 8. Conclusion

In this article, we compared two paradigms for semantic event extraction: Fine-tuning transformer-based architectures as exemplified by our approach T-SEE and prompting Large Language Models (LLMs), exemplified by our approach L-SEE.

Both approaches consist of two main steps: event classification and relation extraction, where T-SEE frames event classification as a multi-label classification task, and conducts relation extraction with a span prediction transformer model. L-SEE provides an LLM with two different prompts which include the event classes and properties in the target event ontology.

In our evaluation, we first introduced two new datasets for semantic event extraction. Then, we compare T-SEE and L-SEE to two state-of-the-art baselines, with T-SEE outperforming or matching them and setting a new benchmark for transformer-based methods in semantic event extraction. Finally, we specifically focused on the different characteristics of T-SEE and L-SEE, highlighting T-SEE's adaptation to the precise characteristics of the training and test data, while L-SEE performs clearly worse on the test data. However, our subsequent analysis revealed its capability of extracting relevant knowledge that is often overlooked by distantly-labelled datasets.

Consequently, we derive a set of phenomena to be regarded when performing semantic event extraction, including the role of distantly-labelled datasets and the event ontology.

In future work, we plan to further improve T-SEE and L-SEE, e.g., by bringing event classification, relation extraction and other tasks like named entity recognition even closer together in joint multi-task learning frameworks and to extend them to encompass multilingual and document-level semantic event extraction. In addition, we aim to enhance metrics and datasets, allowing a fair comparison between semantic event extraction methods employing transformer-based architectures and LLMs.

## References

- [1] S. Latif, S. Agarwal, S. Gottschalk, C. Chrosch, F. Feit, J. Jahn, T. Braun, Y.C. Tchenko, E. Demidova and F. Beck, Visually Connecting Historical Figures Through Event Knowledge Graphs, *Computing Research Repository* (2021).
- [2] S. Gottschalk and E. Demidova, EventKG+BT: Generation of Interactive Biography Timelines from a Knowledge Graph, in: *The Semantic Web*, 2020.
- [3] S. Gottschalk and E. Demidova, EventKG - the Hub of Event Knowledge on the Web - and Biographical Timeline Generation, *The Semantic Web* (2019).
- [4] R. Porzel, M. Pomarlan, L. Spillner, J.A. Bateman, T. Mildner and C. Santagiustina, Narrativizing Knowledge Graphs, in: *International Workshop on Artificial Intelligence Technologies for Legal Documents*, 2022.
- [5] T. Souza Costa, S. Gottschalk and E. Demidova, Event-QA: A dataset for Event-centric Question Answering over Knowledge Graphs, in: *Proceedings of the 29th ACM international conference on information & knowledge management*, 2020.
- [6] R. Shaw, R. Troncy and L. Hardman, LOD: Linking open descriptions of events, in: *The Semantic Web: Fourth Asian Conference, ASWC 2009, Shanghai, China, December 6-9, 2009. Proceedings 4*, 2009.
- [7] W.R. Van Hage, V. Malaisé, R. Segers, L. Hollink and G. Schreiber, Design and Use of the Simple Event Model (SEM), *Journal of Web Semantics* (2011).
- [8] S. Auer, C. Bizer, G. Kobilarov, J. Lehmann, R. Cyganiak and Z. Ives, DBpedia: A Nucleus for a Web of Open Data, in: *The Semantic Web Conference*, 2007.
- [9] F. Erxleben, M. Günther, M. Krötzsch, J. Mendez and D. Vrandečić, Introducing Wikidata to the Linked Data Web, in: *International Semantic Web Conference*, 2014.
- [10] A. Li, X. Wang, W. Wang, A. Zhang and B. Li, A Survey of Relation Extraction of Knowledge Graphs, in: *Web and Big Data International Workshops*, 2019.
- [11] S. Ji, S. Pan, E. Cambria, P. Marttinen and P.S. Yu, A Survey on Knowledge Graphs: Representation, Acquisition and Applications, *Computing Research Repository* (2020).
- [12] S. Shirai, D. Bhattacharjya and O. Hassanzadeh, Event Prediction using Case-Based Reasoning over Knowledge Graphs, in: *The Web Conference*, 2023.
- [13] G. Stoica, E.A. Platanios and B. Póczos, Improving Relation Extraction by Leveraging Knowledge Graph Link Prediction, *Computing Research Repository* (2020).
- [14] W. Xiang and B. Wang, A Survey of Event Extraction from Text, *IEEE Access* (2019).
- [15] J. Zheng, F. Cai, W. Chen, W. Lei and H. Chen, Taxonomy-aware Learning for Few-shot Event Detection, in: *The Web Conference*, 2021.
- [16] S. Mehta, M.R. Islam, H. Rangwala and N. Ramakrishnan, Event Detection Using Hierarchical Multi-aspect Attention, in: *Web Conference*, 2019.
- [17] S. Li, H. Ji and J. Han, Document-level Event Argument Extraction by Conditional Generation (2021).
- [18] X. Wang, Z. Wang, X. Han, Z. Liu, J. Li, P. Li, M. Sun, J. Zhou and X. Ren, HMEAE: Hierarchical Modular Event Argument Extraction, in: *Empirical Methods in Natural Language Processing and International Joint Conference on Natural Language Processing*, 2019.
- [19] Y. Ma, Z. Wang, Y. Cao, M. Li, M. Chen, K. Wang and J. Shao, Prompt for extraction? PAIE: Prompting Argument Interaction for Event Argument Extraction (2022).



- [20] X. Wang, Z. Wang, X. Han, W. Jiang, R. Han, Z. Liu, J. Li, P. Li, Y. Lin and J. Zhou, MAVEN: A Massive General Domain Event Detection Dataset, in: *Empirical Methods in Natural Language Processing*, 2020.
- [21] S. Ebner, P. Xia, R. Culkin, K. Rawlins and B. Van Durme, Multi-Sentence Argument Linking, in: *Association for Computational Linguistics*, 2020.
- [22] M. Rospocher, M. Van Erp, P. Vossen, A. Fokkens, I. Aldabe, G. Rigau, A. Soroa, T. Ploeger and T. Bogaard, Building Event-centric Knowledge Graphs from News, *Journal of Web Semantics* (2016).
- [23] K. Guo, D. Diefenbach, A. Gourru and C. Gravier, Wikidata as a Seed for Web Extraction, in: *The Web Conference*, 2023.
- [24] S. Guan, X. Cheng, L. Bai, F. Zhang, Z. Li, Y. Zeng, X. Jin and J. Guo, What is Event Knowledge Graph: A Survey, *IEEE Transactions on Knowledge and Data Engineering* (2022).
- [25] F. Hamborg, C. Breiteringer and B. Gipp, Giveme5w1h: A Universal System for Extracting Main Events from News Articles, *arXiv preprint arXiv:1909.02766* (2019).
- [26] Y. Zhou, Y. Chen, J. Zhao, Y. Wu, J. Xu and J. Li, What the Role is vs. What Plays the Role: Semi-Supervised Event Argument Extraction via Dual Question Answering, in: *AAAI*, 2021.
- [27] Z. Wang, X. Wang, X. Han, Y. Lin, L. Hou, Z. Liu, P. Li, J. Li and J. Zhou, CLEVE: Contrastive Pre-training for Event Extraction (2021).
- [28] Linguistic Data Consortium, *ACE (Automatic Content Extraction) English Annotation Guidelines for Events*, 2005.
- [29] A.M. Davani, L. Yeh, M. Atari, B. Kennedy, G. Portillo-Wightman, E. Gonzalez, N. Delong, R. Bhatia, A. Mirinjian, X. Ren et al., Reporting the Unreported: Event extraction for Analyzing the Local Representation of Hate Crimes (2019).
- [30] R. Xu, T. Liu, L. Li and B. Chang, Document-level Event Extraction via Heterogeneous Graph-based Interaction Model with a Tracker, in: *Association for Computational Linguistics and the International Joint Conference on Natural Language Processing*, 2021.
- [31] P. Kaliamoorthi, A. Siddhant, E. Li and M. Johnson, Distilling Large Language Models into Tiny and Effective Students using pQRNN, *arXiv preprint arXiv:2101.08890* (2021).
- [32] A. Dunn, J. Dagdelen, N. Walker, S. Lee, A.S. Rosen, G. Ceder, K. Persson and A. Jain, Structured Information Extraction from Complex Scientific Text with Fine-tuned Large Language Models, *arXiv preprint arXiv:2212.05238* (2022).
- [33] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A.N. Gomez, L. Kaiser and I. Polosukhin, Attention Is All You Need, *Advances in Neural Information Processing Systems* (2017).
- [34] G. Tsoumakas and I. Katakis, Multi-label classification: An overview, *Data Warehousing and Mining: Concepts, Methodologies, Tools, and Applications* (2008).
- [35] T.-Y. Lin, P. Goyal, R. Girshick, K. He and P. Dollár, Focal Loss for Dense Object Detection, in: *international conference on computer vision*, 2017.
- [36] J. Devlin, M.-W. Chang, K. Lee and K. Toutanova, Bert: Pre-training of Deep Bidirectional Transformers for Language Understanding, *North American Chapter of the Association for Computational Linguistics* (2019).
- [37] Z. Lan, M. Chen, S. Goodman, K. Gimpel, P. Sharma and R. Soricut, ALBERT: A Lite BERT for Self-supervised Learning of Language Representations, in: *International Conference on Learning Representations*, 2019.
- [38] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer and V. Stoyanov, RoBERTa: A Robustly Optimized BERT Pretraining Approach, *arXiv preprint arXiv:1907.11692* (2019).
- [39] Z. Yang, Z. Dai, Y. Yang, J. Carbonell, R.R. Salakhutdinov and Q.V. Le, Xlnet: Generalized Autoregressive Pretraining for Language Understanding, *Advances in Neural Information Processing Systems* (2019).
- [40] R. Hoekstra and P. Groth, PROV-O-Viz-understanding the Role of Activities in Provenance, in: *Provenance and Annotation of Data and Processes*, 2015.
- [41] D. Wadden, U. Wennberg, Y. Luan and H. Hajishirzi, Entity, Relation, and Event Extraction with Contextualized Span Representations, in: *Empirical Methods in Natural Language Processing and the International Joint Conference on Natural Language Processing*, 2019.
- [42] L. Liu, M. Liu, S. Liu and K. Ding, Event extraction as machine reading comprehension with question-context bridging, *Knowledge-Based Systems* (2024).
- [43] P. Wu, X. Li, J. Gu, L. Qian and G. Zhou, Pipelined biomedical event extraction rivaling joint learning, *Methods* (2024).
- [44] F. Shiri, V. Nguyen, F. Moghimifar, J. Yoo, G. Haffari and Y.-F. Li, Decompose, Enrich, and Extract! Schema-aware Event Extraction using LLMs, *2024 27th International Conference on Information Fusion (FUSION)* (2024).
- [45] K. Shenoy, F. Ilievski, D. Garijo, D. Schwabe and P. Szekely, A Study of the Quality of Wikidata, *Journal of Web Semantics* (2022).
- [46] S. Guo, C. Wang, Y. Chen, K. Liu, R. Li and J. Zhao, EventOA: An Event Ontology Alignment Benchmark based on FrameNet and Wikidata, in: *Findings of the Association for Computational Linguistics*, 2023.
- [47] V.A. Carriero, A. Gangemi, A.G. Nuzzolese, V. Presutti et al., An Ontology Design Pattern for Representing Recurrent Events., in: *WOP@ ISWC*, 2019.
- [48] S. Gottschalk and E. Demidova, HapPenIng: Happen, Predict, Infer—Event Series Completion in a Knowledge Graph, in: *International Semantic Web Conference*, 2019.
- [49] O. Hassanzadeh, Building a Knowledge Graph of Events and Consequences Using Wikidata., *Wikidata@ ISWC* (2021).
- [50] C. Rudnik, T. Ehrhart, O. Ferret, D. Teyssou, R. Troncy and X. Tannier, Searching News Articles using an Event Knowledge Graph leveraged by Wikidata, in: *Companion Proceedings of the 2019 World Wide Web Conference*, 2019.
- [51] M. Mintz, S. Bills, R. Snow and D. Jurafsky, Distant Supervision for Relation Extraction without Labeled Data, in: *Proceedings of the Joint Conference of the 47th Annual Meeting of the ACL and the 4th International Joint Conference on Natural Language Processing of the AFNLP*, 2009.

- [52] H. Elsahar, P. Vougiouklis, A. Remaci, C. Gravier, J. Hare, F. Laforest and E. Simperl, T-REx: A Large Scale Alignment of Natural Language with Knowledge Base Triples, in: *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, 2018.
- [53] L. Guzman-Nateras, V. Lai, A. Pouran Ben Veyseh, F. Dernoncourt and T. Nguyen, Event Detection for Suicide Understanding, in: *North American Association for Computational Linguistics*, 2022.
- [54] Y. Luan, L. He, M. Ostendorf and H. Hajishirzi, Multi-Task Identification of Entities, Relations, and Coreference for Scientific Knowledge Graph Construction, in: *Empirical Methods in Natural Language Processing*, 2018.
- [55] T. Ohta, Y. Tateisi, J.-D. Kim, H. Mima and J. Tsujii, The GENIA corpus: An Annotated Research Abstract Corpus in Molecular Biology Domain, in: *Human Language Technology Conference*, 2002.
- [56] A. Pouran Ben Veyseh, M.V. Nguyen, F. Dernoncourt and T. Nguyen, MINION: a Large-Scale and Diverse Dataset for Multilingual Event Detection, in: *North American Chapter of the Association for Computational Linguistics*, 2022.
- [57] K. Leetaru and P.A. Schrodt, GDEL: Global Data on Events, Location, and Tone, in: *ISA Annual Convention*, 2013.
- [58] M. Li, R. Xu, S. Wang, L. Zhou, X. Lin, C. Zhu, M. Zeng, H. Ji and S.-F. Chang, Clip-event: Connecting Text and Images with Event Structures, in: *Conference on Computer Vision and Pattern Recognition*, 2022.
- [59] Y. Lu, H. Lin, J. Xu, X. Han, J. Tang, A. Li, L. Sun, M. Liao and S. Chen, Text2Event: Controllable Sequence-to-Structure Generation for End-to-end Event Extraction, in: *Association for Computational Linguistics and the International Joint Conference on Natural Language Processing*, 2021.
- [60] H. You, D. Samuel, S. Touileb and L. Øvrelid, EventGraph: Event Extraction as Semantic Graph Parsing (2022).
- [61] F. Li, W. Peng, Y. Chen, Q. Wang, L. Pan, Y. Lyu and Y. Zhu, Event Extraction as Multi-turn Question Answering, in: *Findings of the Association for Computational Linguistics: EMNLP 2020*, 2020.
- [62] X. Liu, H.-Y. Huang, G. Shi and B. Wang, Dynamic Prefix-Tuning for Generative Template-based Event Extraction, in: *The Association for Computational Linguistics*, 2022.
- [63] H. Huang, X. Liu, G. Shi and Q. Liu, Event Extraction with Dynamic Prefix Tuning and Relevance Retrieval, *IEEE Transactions on Knowledge and Data Engineering* (2023).
- [64] S. Liu, Y. Li, F. Zhang, T. Yang and X. Zhou, Event Detection Without Triggers, in: *North American Chapter of the Association for Computational Linguistics*, 2019.
- [65] X. Du and C. Cardie, Event Extraction by Answering (Almost) Natural Questions, in: *Empirical Methods in Natural Language Processing*, 2020.
- [66] I. Hsu, K.-H. Huang, E. Boschee, S. Miller, P. Natarajan, K.-W. Chang, N. Peng et al., DEGREE: A Data-efficient Generation-based Event Extraction Model (2021).
- [67] D. Lu, S. Ran, J. Tetreault and A. Jaimes, Event Extraction as Question Generation and Answering (2023).
- [68] X. Wei, X. Cui, N. Cheng, X. Wang, X. Zhang, S. Huang, P. Xie, J. Xu, Y. Chen, M. Zhang et al., Zero-shot Information Extraction via Chatting with Chatgpt, *arXiv preprint arXiv:2302.10205* (2023).
- [69] C. Wang, X. Liu, Z. Chen, H. Hong, J. Tang and D. Song, DeepStruct: Pretraining of Language Models for Structure Prediction (2022).
- [70] J.L. Fleiss, Measuring Nominal Scale Agreement among Many Raters, *Psychological bulletin* (1971).
- [71] X. Han, T. Gao, Y. Lin, H. Peng, Y. Yang, C. Xiao, Z. Liu, P. Li, J. Zhou and M. Sun, More Data, More Relations, More Context and More Openness: A Review and Outlook for Relation Extraction, in: *The Association for Computational Linguistics and the 10th International Joint Conference on Natural Language Processing*, 2020.
- [72] B. Goodrich, V. Rao, P.J. Liu and M. Saleh, Assessing the Factual Accuracy of Generated Text, in: *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining*, 2019.
- [73] Y. Yao, D. Ye, P. Li, X. Han, Y. Lin, Z. Liu, Z. Liu, L. Huang, J. Zhou and M. Sun, DocRED: A Large-Scale Document-Level Relation Extraction Dataset, in: *The Association for Computational Linguistics*, 2019.
- [74] K. Mai, T.-H. Pham, M.T. Nguyen, T.D. Nguyen, D. Bollegala, R. Sasano and S. Sekine, An Empirical Study on Fine-grained Named Entity Recognition, in: *Proceedings of the 27th International Conference on Computational Linguistics*, 2018.
- [75] S. Ji, S. Pan, E. Cambria, P. Martinen and P. Yu, A Survey on Knowledge Graphs: Representation, Acquisition, and Applications, *IEEE Transactions on Neural Networks and Learning Systems* (2021).
- [76] S. Gottschalk and E. Demidova, EventKG: a Multilingual Event-centric Temporal Knowledge Graph, in: *Extended Semantic Web Conference*, 2018.
- [77] S. Gottschalk, E. Kacupaj, S. Abdollahi, D. Alves, G. Amaral, E. Koutsiana, T. Kuculo, D. Major, C. Mello, G.S. Cheema et al., OEKG: The Open Event Knowledge Graph, in: *CLEOPATRA Workshop co-located with the 30th The Web Conference*, 2021.
- [78] D.J. Gerner, P.A. Schrodt, O. Yilmaz and R. Abu-Jabr, Conflict and Mediation Event Observations (CAMEO): A New Event Data Framework for the Analysis of Foreign Policy Interactions, *International Studies Association* (2002).
- [79] H. Zhang, X. Liu, H. Pan, Y. Song and C.W.-K. Leung, ASER: A Large-scale Eventuality Knowledge Graph, in: *The Web Conference*, 2020.
- [80] W. Yao, Z. Dai, M. Ramaswamy, B. Min and R. Huang, Weakly Supervised Subevent Knowledge Acquisition, in: *Empirical Methods in Natural Language Processing*, 2020.
- [81] R. Piryani, N. Aussenac-Gilles and N. Hernandez, Comprehensive Survey on Ontologies about Event, in: *ESWC Workshops on Semantic Methods for Events and Stories (SEMME@ ESWC 2023)*, 2023.
- [82] V.A. Carriero, A. Gangemi, A.G. Nuzzolese, V. Presutti et al., An Ontology Design Pattern for Representing Recurrent Situations., in: *WOP (Book)*, 2021.
- [83] A. Krisnadhi and P. Hitzler, A Core Pattern for Events, *Advances in Ontology Design and Patterns* (2017).

- [84] X. Du, A.M. Rush and C. Cardie, GRIT: Generative Role-filler Transformers for Document-level Event Entity Extraction, in: *European Chapter of the Association for Computational Linguistics*, 2021.
- [85] J. Liu, Y. Chen, K. Liu, W. Bi and X. Liu, Event Extraction as Machine Reading Comprehension, in: *Empirical Methods in Natural Language Processing*, 2020.
- [86] S. Yang, D. Feng, L. Qiao, Z. Kan and D. Li, Exploring Pre-trained Language Models for Event Extraction and Generation, in: *Association for Computational Linguistics*, 2019.
- [87] S. Zheng, W. Cao, W. Xu and J. Bian, Doc2EDAG: An End-to-End Document-level Framework for Chinese Financial Event Extraction, in: *Empirical Methods in Natural Language Processing and the International Joint Conference on Natural Language Processing*, 2019.
- [88] D. Lou, Z. Liao, S. Deng, N. Zhang and H. Chen, MLBiNet: A Cross-Sentence Collective Event Detection Network, in: *Association for Computational Linguistics and the International Joint Conference on Natural Language Processing*, 2021.
- [89] X. Liu, H.-Y. Huang and Y. Zhang, Open Domain Event Extraction Using Neural Latent Variable Models, in: *Association for Computational Linguistics*, 2019.
- [90] D. Rusu, J. Hodson and A. Kimball, Unsupervised Techniques for Extracting and Clustering Complex Events in News, in: *EVENTS*, 2014.
- [91] H. Cao, J. Li, F. Su, F. Li, H. Fei, S. Wu, B. Li, L. Zhao and D. Ji, OneEE: A One-Stage Framework for Fast Overlapping and Nested Event Extraction, in: *Proceedings of the 29th International Conference on Computational Linguistics*, 2022.
- [92] Y. Ma, Z. Wang, Y. Cao, M. Li, M. Chen, K. Wang and J. Shao, Prompt for Extraction? PAIE: Prompting Argument Interaction for Event Argument Extraction, in: *Association for Computational Linguistics*, 2022.
- [93] J. Peng, W. Yang, F. Wei, L. He, L. Yao and H. Lv, Event co-occurrences for prompt-based generative event argument extraction, *Scientific Reports* (2024).
- [94] G. Zhang, H. Zhang, Y. Wang, R. Li, H. Tan and J. Liang, Hyperspherical multi-prototype with optimal transport for event argument extraction, in: *The Association for Computational Linguistics*, 2024.
- [95] J. Zhao and H. Yang, Trigger-free Event Detection via Derangement Reading Comprehension, *CoRR* (2022).
- [96] H. Huang, Y. Chen, C. Lin, R. Huang, Q. Zheng and Y. Qin, A multi-graph representation for event extraction, *Artificial Intelligence* (2024).
- [97] A. Balali, M. Asadpour and S.H. Jafari, COFEE: A Comprehensive Ontology for Event Extraction from text, *Computer Speech & Language* (2021).
- [98] M.E. Peters, M. Neumann, M. Iyyer, M. Gardner, C. Clark, K. Lee and L. Zettlemoyer, Deep Contextualized Word Representations, in: *The North American Chapter of the Association for Computational Linguistics*, 2018.
- [99] T. Brown, B. Mann, N. Ryder, M. Subbiah, J.D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell et al., Language Models are Few-shot Learners, *Advances in neural information processing systems* (2020).
- [100] N.P. Jouppi, C. Young, N. Patil, D. Patterson, G. Agrawal, R. Bajwa, S. Bates, S. Bhatia, N. Boden, A. Borchers et al., In-datacenter Performance Analysis of a Tensor Processing Unit, in: *international symposium on computer architecture*, 2017.
- [101] M.P. Polak and D. Morgan, Extracting Accurate Materials Data from Research Papers with Conversational Language Models and Prompt Engineering, *Nature Communications* (2024).
- [102] B. Li, G. Fang, Y. Yang, Q. Wang, W. Ye, W. Zhao and S. Zhang, Evaluating ChatGPT's Information Extraction Capabilities: An Assessment of Performance, Explainability, Calibration, and Faithfulness, *CoRR* (2023).
- [103] Y. Liu, T. Han, S. Ma, J. Zhang, Y. Yang, J. Tian, H. He, A. Li, M. He, Z. Liu et al., Summary of Chatgpt-related Research and Perspective Towards the Future of Large Language Models, *Meta-Radiology* (2023).
- [104] C. Kirti, A. Chattopadhyay, A. Anand and P. Guha, Deciphering Storytelling Events: A Study of Neural and Prompt-Driven Event Detection in Short Stories, in: *2023 International Conference on Asian Language Processing (IALP)*, 2023.
- [105] O. Sharif, M. Basak, T. Parvin, A. Scharfstein, A. Bradham, J.T. Borodovsky, S.E. Lord and S.M. Preum, Characterizing Information Seeking Events in Health-Related Social Discourse (2024).
- [106] S. Li, Q. Zhan, K. Conger, M. Palmer, H. Ji and J. Han, GLEN: General-Purpose Event Detection for Thousands of Types, in: *The 2023 Conference on Empirical Methods in Natural Language Processing*, 2023.
- [107] X.F. Zhang, C. Blum, T. Choji, S. Shah and A. Vempala, ULTRA: Unleash LLMs' Potential for Event Argument Extraction through Hierarchical Modeling and Pair-wise Refinement (2024).
- [108] J. Peng, H. Sun, W. Yang, F. Wei, L. He and L. Wang, One Small and One Large for Document-level Event Argument Extraction, *arXiv preprint arXiv:2411.05895* (2024).
- [109] J. Gao, H. Zhao, W. Wang, C. Yu and R. Xu, EventRL: Enhancing Event Extraction with Outcome Supervision for Large Language Models, *arXiv preprint arXiv:2402.11430* (2024).
- [110] N. Mihindukulasooriya, S. Tiwari, C.F. Enguix and K. Lata, Text2kgbench: A Benchmark for Ontology-driven Knowledge Graph Generation from Text, in: *International Semantic Web Conference*, 2023.
- [111] L. Yao, J. Peng, C. Mao and Y. Luo, Exploring Large Language Models for Knowledge Graph Completion, *arXiv preprint arXiv:2308.13916* (2023).
- [112] Y. Zhu, X. Wang, J. Chen, S. Qiao, Y. Ou, Y. Yao, S. Deng, H. Chen and N. Zhang, LLMs for Knowledge Graph Construction and Reasoning: Recent Capabilities and Future Opportunities, *World Wide Web* (2024).
- [113] Y. Zhang, Z. Chen, W. Zhang and H. Chen, Making Large Language Models Perform Better in Knowledge Graph Completion (2023).
- [114] F.M. Suchanek, G. Kasneci and G. Weikum, Yago: A Core of Semantic Knowledge, in: *The Web Conference*, 2007.
- [115] Y. Luan, D. Wadden, L. He, A. Shah, M. Ostendorf and H. Hajishirzi, A General Framework for Information Extraction using Dynamic Span Graphs, in: *North American Chapter of the Association for Computational Linguistics*, 2019.

- [116] J. Liu, L. Min and X. Huang, An Overview of Event Extraction and its Applications, *arXiv preprint arXiv:2111.03212* (2021).
- [117] S.T.R. Rizvi, D. Mercier, S. Agne, S. Erkel, A. Dengel and S. Ahmed, Ontology-based Information Extraction from Technical Documents., in: *ICAART*, 2018.
- [118] K. Guo, D. Diefenbach, A. Gourru and C. Gravier, Wikidata as a Seed for Web Extraction, in: *The Web Conference*, 2023.
- [119] C. Walker, S. Strassel, J. Medero and K. Maeda, ACE 2005 Multilingual Training Corpus, *Linguistic Data Consortium, Philadelphia* (2006).
- [120] Z. Li, S. Zhao, X. Ding and T. Liu, EEG: Knowledge Base for Event Evolutionary Principles and Patterns, in: *Chinese National Conference on Social Media Processing*, 2017.
- [121] S. Abdollahi, S. Gottschalk and E. Demidova, EventKG+ Click: A Dataset of Language-specific Event-centric User Interaction Traces, in: *CLEOPATRA Workshop co-located with the 17th Extended Semantic Web Conference*, 2020.
- [122] Q. Ning, B. Zhou, Z. Feng, H. Peng and D. Roth, CogCompTime: A Tool for Understanding Time in Natural Language, in: *Empirical Methods in Natural Language Processing*, 2018.
- [123] M.D. Ma, J. Sun, M. Yang, K.-H. Huang, N. Wen, S. Singh, R. Han and N. Peng, EventPlus: A Temporal Event Understanding Pipeline, in: *North American Chapter of the Association for Computational Linguistics*, 2021.
- [124] H. Wang, H. Zhang, M. Chen and D. Roth, Learning Constraints and Descriptive Segmentation for Subevent Detection, in: *Empirical Methods in Natural Language Processing*, 2021.
- [125] P. Rajpurkar, R. Jia and P. Liang, Know What You Don't Know: Unanswerable Questions for SQuAD, in: *Association for Computational Linguistics*, 2018.
- [126] L. Qiu, H. Zhou, Y. Qu, W. Zhang, S. Li, S. Rong, D. Ru, L. Qian, K. Tu and Y. Yu, QA4IE: a Question Answering based Framework for Information Extraction, in: *International Semantic Web Conference*, 2018.
- [127] X. Li, F. Yin, Z. Sun, X. Li, A. Yuan, D. Chai, M. Zhou and J. Li, Entity-Relation Extraction as Multi-Turn Question Answering, in: *Association for Computational Linguistics*, 2019.
- [128] A. Delpeuch, OpenTapioca: Lightweight Entity Linking for Wikidata, in: *Wikidata Workshop co-located with International Semantic Web Conference*, 2020.
- [129] B. Xu, Y. Zhang, J. Liang, Y. Xiao, S.-w. Hwang and W. Wang, Cross-lingual Type Inference, in: *International Conference on Database Systems for Advanced Applications*, 2016.
- [130] M.D. Ward, N.W. Metternich, C. Carrington, C. Dorff, M. Gallop, F.M. Hollenbach, A. Schultz, S. Weschle et al., Geographical Models of Crises: Evidence from ICEWS, *Design for Cross-Cultural Activities* (2012).
- [131] J.L. Martinez-Rodriguez, A. Hogan and I. Lopez-Arevalo, Information Extraction Meets the Semantic Web: a Survey, *Semantic Web* (2020).
- [132] P.-L.H. Cabot and R. Navigli, REBEL: Relation Extraction by End-to-end Language Generation, in: *Empirical Methods in Natural Language Processing*, 2021.
- [133] X. Han, T. Gao, Y. Yao, D. Ye, Z. Liu and M. Sun, OpenNRE: An Open and Extensible Toolkit for Neural Relation Extraction, in: *Empirical Methods in Natural Language Processing and the International Joint Conference on Natural Language Processing*, 2019.
- [134] K. Verspoor, J. Cohn, S. Mniszewski and C. Joslyn, A Categorization Approach to Automated Ontological Function Annotation, *Protein Science* (2006).
- [135] G. Rossiello, F. Chowdhury, N. Mihindukulasooriya, O. Cornec and A. Gliozzo, KnowGL: Knowledge Generation and Linking from Text (2022).
- [136] T. Liu, Y.E. Jiang, N. Monath, R. Cotterell and M. Sachan, Autoregressive Structured Prediction with Language Models, in: *Empirical Methods in Natural Language Processing*, 2022.
- [137] Z. Zhong and D. Chen, A Frustratingly Easy Approach for Entity and Relation Extraction, in: *North American Chapter of the Association for Computational Linguistics*, 2021.
- [138] D. Ye, Y. Lin, P. Li and M. Sun, Packed Levitated Marker for Entity and Relation Extraction, in: *Association for Computational Linguistics*, 2022.
- [139] J. Chen, H. Lin, X. Han and L. Sun, Honey or Poison? Solving the Trigger Curse in Few-shot Event Detection via Causal Intervention, in: *Empirical Methods in Natural Language Processing*, 2021.
- [140] M. Tong, B. Xu, S. Wang, M. Han, Y. Cao, J. Zhu, S. Chen, L. Hou and J. Li, DocEE: A Large-scale and Fine-grained Benchmark for Document-level Event Extraction, *Association for Computational Linguistics*, 2022.
- [141] T. Ge, L. Cui, B. Chang, Z. Sui, F. Wei and M. Zhou, Eventwiki: A Knowledge Base of Major Events, in: *International Conference on Language Resources and Evaluation*, 2018.
- [142] J. Ellis, J. Getman, D. Fore, N. Kuster, Z. Song, A. Bies and S.M. Strassel, Overview of Linguistic Resources for the TAC KBP 2015 Evaluations: Methodologies and Results., in: *Tac*, 2015.
- [143] J. Zheng, F. Cai, W. Chen, W. Lei and H. Chen, Taxonomy-aware Learning for Few-Shot Event Detection, in: *The Web Conference*, 2021.
- [144] S. Zhou, B. Yu, A. Sun, C. Long, J. Li and J. Sun, A Survey on Neural Open Information Extraction: Current Status and Future Directions, in: *International Joint Conference*, 2022.
- [145] V. Perot, K. Kang, F. Luisier, G. Su, X. Sun, R.S. Boppana, Z. Wang, J. Mu, H. Zhang and N. Hua, LMDX: Language Model-based Document Information Extraction and Localization (2024).
- [146] N. Heist, S. Hertling, D. Ringler and H. Paulheim, Knowledge Graphs on the Web-An Overview., *Knowledge Graphs for eXplainable Artificial Intelligence* (2020).

## Appendix A. Prompts of L-SEE

This appendix shows the prompts used by L-SEE as described in Sections 4.1 and 4.2.

### A.1. Prompt for Event Classification

Fig. 16 shows the LLM prompt template we use for L-SEE's event classification as described in Section 4.1.

#### LLM Prompt for Event Classification

##### **Prompt:**

*Your task is to analyse the sentence and classify events that are in the sentence.*

*An event is identified by an action or a mention of an event.*

*You will only consider events that are likely to have their own Wikipedia page.*

*Note: The events that you should identify are links in Wikipedia, they may not be referred to directly by name in the sentence but a specific word or phrase in the sentence may link to the event. E.g., in "Senator McCain also got 10% higher approval rating compared to 2010", 2010 is a link to the event "United States Senate elections, 2010" even though it is not mentioned directly in the sentence.*

*For example, in the sentence "John married Mary in Paris on 12th December 2020 during the Parisian Unrests", the events are "marriage" and "unrests". However, "married" is not to be considered an event as it is unlikely to have its own Wikipedia page.*

*You are to select event types from the following list of event types and return it as a list of strings of event types:*

*{All Event Classes C}*

*This is your task:*

*Sentence:*

*{Text t}*

Fig. 16. Illustration of a structured prompt provided to the LLM for event classification.

## A.2. Prompt for Relation Extraction

Fig. 17 shows the LLM prompt template we use for L-SEE's relation extraction as described in Section 4.2.

### LLM Prompt for Relation Extraction

#### Prompt:

Your task is to extract the properties of the events that are in a given sentence and their values. You will only consider properties that are likely to be associated with the given event classes. Extract the properties of the events and return a JSON object with the event classes as the keys and the properties as the values.

The property values can be dates, entities, or quantities. If there is no specific value for a property, you must not include it in the JSON object.

The extracted property values must fit their respective property types. For example, if the property is "date", the value must be able to be formatted as a date (e.g. "12th December 2019" or "2019" in the case of a year).

Similarly, if the property is "location", the value must be a location. If there are multiple values for a property, you must include all the values in a list.

Consider the following example:

Sentence:

John married Mary on the first day of the start of the COVID-19 pandemic, on 12th December 2019. It was only a few days later that in the winter of 2019, the German-French War destroyed the cities of Paris and Berlin.

Event classes and their potential properties:

- Pandemic: city, startDate
- MilitaryConflict: city, date, participant

Output:

```
{
  "Pandemic": {
    "startDate": ["12th December 2019"]
  },
  "MilitaryConflict": {
    "city": ["Paris", "Berlin"],
    "date": ["2019"]
  }
}
```

This is your task:

Sentence:

{Text  $t$ }

Event classes and their potential properties:

{Event classes  $C_t$  and properties  $P_{C_t}$ }

Fig. 17. Illustration of a structured prompt provided to the LLM for property extraction.