

# Zero-Shot On-the-Fly Event Schema Induction

Anonymous EMNLP submission

## Abstract

001 What are the events involved in a pandemic  
002 outbreak? What steps should be taken when  
003 planning a wedding? The answers to these  
004 questions can be found by collecting many  
005 documents on the complex event of interest,  
006 extracting relevant information, and analyzing  
007 it. We present a new approach<sup>1</sup> in which  
008 large language models are utilized to gener-  
009 ate source documents that allow predicting,  
010 given a high-level event definition, the specific  
011 events, arguments, and relations between them  
012 to construct a schema that describes the com-  
013 plex event in its entirety. Using our model,  
014 complete schemas on any topic can be gener-  
015 ated on-the-fly without any data collection  
016 needed, i.e., in a zero-shot manner. More-  
017 over, we develop efficient methods to extract  
018 pertinent information from texts and demon-  
019 strate, in a series of experiments, that these  
020 schemas are considered to be more complete  
021 than human-curated ones in the majority of ex-  
022 amined scenarios. Finally, we show that this  
023 framework is comparable in performance with  
024 previous supervised schema induction meth-  
025 ods that rely on collecting real texts, while be-  
026 ing more general and flexible by avoiding the  
027 need to use a predefined ontology.

## 028 1 Introduction

029 Event processing refers to tracking, analyzing, and  
030 drawing conclusions from streams of information  
031 about events. This event analysis aims at identi-  
032 fying meaningful events (such as opportunities or  
033 threats) in real-time situations and responding ap-  
034 propriately. Event processing can also be utilized  
035 to gain a deep understanding of the specific steps,  
036 arguments, and relations between them that are in-  
037 volved in a complex event. The information above  
038 can be consolidated into a graphical representation  
039 called an *event schema* (Li et al., 2021).

<sup>1</sup>Our code and data will be made publicly available upon acceptance.

040 Consider the example schema of kidnapping pre-  
041 sented in Fig. 1. This representation of events and  
042 participants assists in gaining an understanding of  
043 the complex event of kidnapping and could help  
044 composing a reaction plan if needed.

045 The NLP community has devoted much effort to  
046 understanding events that are described in a docu-  
047 ment or in a collection of documents for this pur-  
048 pose. These efforts include identifying event trig-  
049 gers (Lu and Roth, 2012; Huang et al., 2018; Wad-  
050 den et al., 2019; Han et al., 2019), extracting event  
051 arguments (Punyakanok et al., 2008; Peng et al.,  
052 2016; Lin et al., 2020; Zhang et al., 2021a), and pre-  
053 dicting the relations between events, e.g., temporal,  
054 coreference, causal or hierarchical relations (Do  
055 et al., 2012; Lee et al., 2012; Glavaš et al., 2014;  
056 Caselli and Vossen, 2017; Ning et al., 2018; Wang  
057 et al., 2020; Zhang et al., 2020a).

058 Previous works on event schema induction re-  
059 lied on the information extracted from collected  
060 documents to build the schema graph. For instance,  
061 Li et al. (2020) learn an auto-regressive language  
062 model (LM) over paths in the instance graphs de-  
063 picting events, arguments and relations of instances  
064 of the complex events, and later on construct a  
065 schema graph by merging the top  $k$  ranked paths.  
066 However, their approach requires access to many  
067 documents on each topic of interest, which can be  
068 extremely laborious and time consuming to obtain.

069 Our goal, on the other hand, is to allow creat-  
070 ing schemas on-the-fly by taking as input only  
071 the name of the complex event of interest (like a  
072 “pandemic outbreak” or an “armed robbery”). To  
073 avoid collecting many documents on the topic of  
074 the schema, we utilize pre-trained auto-regressive  
075 text generation models, specifically GPT-3 (Brown  
076 et al., 2020), to generate texts on the desired topic  
077 (examples presented in Fig. 2). These documents  
078 are then processed to extract pertinent informa-  
079 tion, from which a schema is constructed. The fact  
080 that we do not collect any data makes our learning

## A Kidnapping Schema

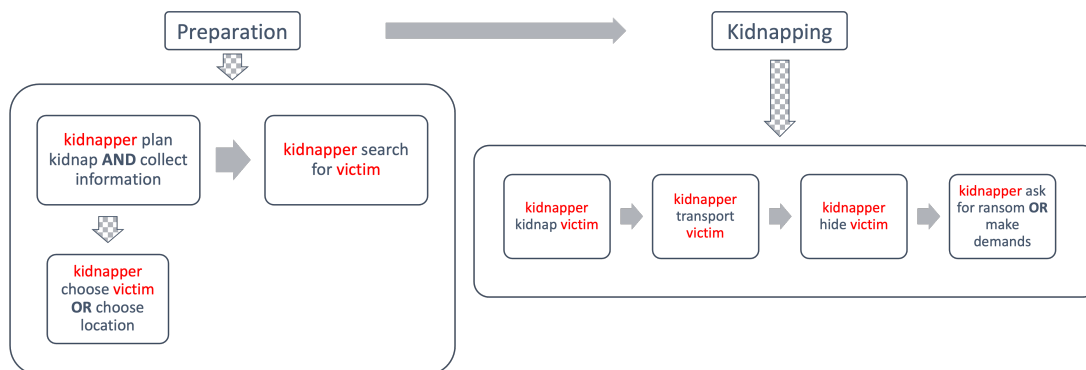


Figure 1: An example schema for the event of “Kidnapping”. The gray arrows present temporal relations and the checkered arrows present hierarchical relations (PARENT-CHILD).

framework zero-shot since we do not rely on any human-collected articles or example schemas.

In addition to making the induction faster by eliminating the need to collect data, we also made the information extraction process faster by implementing new and efficient methods for identifying temporal and hierarchical relations between events mentioned in the text. These two steps are the most time consuming in the process of schema induction and could take up to 2 hours each. Accepting the whole text as input instead of two sentences at each time, the proposed model shortens the inference time significantly to several minutes without enduring a major loss in performance.

The process of generating texts is explained in Section §3, and the process of extracting relevant and salient information is described in Section §4, then we introduce the construction of the schema graph in Section §5. To evaluate our zero-shot schema generator we conduct experiments on a benchmark dataset for schema induction (LDC2020E25) and provide a new dataset for further evaluation called Schema-11. Additionally, we design a subject-matter expert Turing test, a.k.a. Feigenbaum test (Feigenbaum, 2003), to determine whether our algorithm could mimic experts’ response towards several common complex event scenarios. The experiments and results are presented in Section §6.

The contributions of our work include:

1. Predicting an entire schema given the name of a complex event without collecting data.
2. Implementing a novel and efficient One-Pass approach for identifying temporal and hierar-

chical relations between events.

3. Presenting a method for automatically inducing logical relations between events based on temporal relations.
4. Offering a Feigenbaum test for evaluation on a new schema dataset, Schema-11.

## 2 Related Work

**Schema induction:** Early schema induction efforts focused on identifying the triggers and participants of atomic events without considering relations between atomic events that comprise complex schemas (Chambers, 2013; Cheung et al., 2013; Nguyen et al., 2015; Sha et al., 2016; Yuan et al., 2018). More recent work focuses on inducing schemas for pairs of events (Li et al., 2020) and multiple events (Zhang et al., 2021b; Li et al., 2021), but they require access to large corpora for the induction process. In this work, we induce schemas on-the-fly in a zero-shot manner. As is standard in state-of-the-art (SOTA) works (Li et al., 2020, 2021; Wen et al., 2021), we output all the information about relations between events and arguments extracted from the text, in addition to logical and hierarchical relations not studied previously.

**Script learning:** Early script learning works concentrated on chains of events with a single protagonist (Chambers and Jurafsky, 2008, 2009; Jans et al., 2012; Rudinger et al., 2015; Granroth-Wilding and Clark, 2016) and later extended to multiple protagonists (Pichotta and Mooney, 2014; Peng and Roth, 2016; Pichotta and Mooney, 2016; Modi, 2016; Weber et al., 2018, 2020; Zhang et al.,

2020b). All of these works assume there exists a single line of events that describes all occurrences within a complex event. This work does not limit itself to generating single-chained schemas. We also consider more complex graphs as schema outputs. In addition, none of these works deals with zero-shot scenarios that do not require training data.

**Pre-trained generation models:** Large-scale pre-trained text generation models such as GPT-3 (Brown et al., 2020), BART (Lewis et al., 2020), T5 (Raffel et al., 2020), i.a. have been used in NLP to solve many tasks. These models are often seen as few-shot learners (Brown et al., 2020) and therefore used as inference methods. However, these text generation models are not explicitly trained to perform inference, rather they are trained to produce the most likely sequence of words to proceed a certain prompt, similar to language models (which is also how researchers refer to them).

In a recently published paper, Wang et al. (2021b) used GPT-3 to generate training data in a few-shot inference paradigm by querying the model with a prompt and a few examples in order to create additional examples with a desired label. Later they used the generated data to fine-tune standard pre-trained models to perform inference. Our work, however, uses these large pre-trained LMs only as text generators. We generate documents on a particular topic and use it as a corpus for extracting the topic’s schema. We rely on the intuition that the generated text will include salient and stereotypical information that is expected to be mentioned in the context of the topic (e.g., for a topic of “planning a wedding,” we assume most documents will include the event “order catering”).

### 3 Data Generation

The schema induction process begins with generating texts using large LMs as text generation models. These texts are joined to form a knowledge base for the schema, including all of the potential information that the schema may present. One could, of course, create this knowledge base by crawling the web for real news articles or Wikipedia entries related to a certain topic.

We argue, however, that in addition to the obvious advantages of not having to rely on the availability of data online and not having to crawl the entire web for relevant documents on each topic, the generated data from these large generative models is more efficient in reporting salient events than

	Generated Text	Real Text
# events / # tokens	0.1252	0.0631
# arguments / # tokens	0.0545	0.0301

Table 1: The ratio of relevant events and relevant argument roles identified in generated text and real text for the scenario of IED attack.

random events described in the news, i.e., generated texts are more likely to mention important information than real documents do.

Our analysis shows that the generated stories contain a higher percentage of relevant tokens than existing real news articles that are used for schema induction. To demonstrate this phenomenon, we compare manually gathered documents with those that are automatically generated for the event of Improvised Explosive Device (IED) attack (Li et al., 2021). To identify salient events and arguments concerning IED attacks, we adopt the DARPA KAIROS Phase 1 (v3.0) ontology<sup>2</sup>, a fine-grained ontology for schema learning, with 24 entity types, 67 event types, and 85 argument roles.

We calculated the number of relevant event triggers and arguments identified in the text, where a relevant mention is one whose type appears in the ontology. The results shown in Table 1 demonstrate that the quality of the generated texts in terms of conciseness and appearance of important details is higher than that of real texts. For example, the ratio of relevant events per token is more than twice as high in generated texts as it is in real texts. Hence we are able to not only generate a schema for every given topic without putting any effort in searching the web, but the information we generate is also better suited for our end task of depicting all of the important aspects of a complex event.

Given a topic for which we want to create a schema, we generate multiple texts that discuss the topic event using the OpenAI GPT-3 API<sup>3</sup> with the Davinci-instruct-beta-v3 model. We use three prompting methods to generate documents of diverse genres as follows:

**News articles:** We begin by generating a headline using the prompt: “Write a news headline about *topic*.” The output from this prompt, which appears between quotation marks, is then used in the following prompt: “Write a news story titled *headline*.”

<sup>2</sup>The full ontology definition can be accessed at this link: <https://bit.ly/3mIWJoN>.

<sup>3</sup><https://openai.com/blog/openai-api/>.

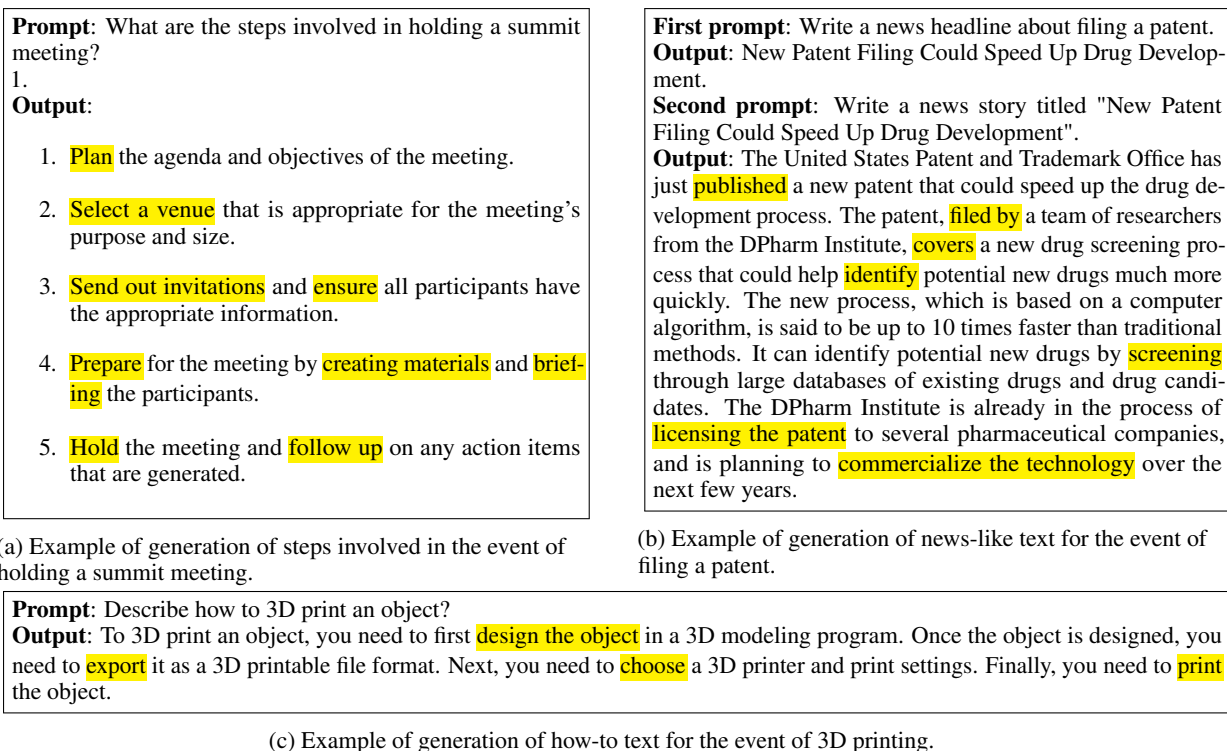


Figure 2: Examples of generated texts using different prompting methods. The highlighted text is relevant events that will be extracted in the information extraction step.

The output from the second prompt is added to the pool of generated texts. The process is repeated 30 times. See example in Fig. 2b.

**How-to articles:** For this genre type, we use the prompt: “Describe how to *topic*.” The process is repeated 30 times and all the generated texts are added to the pool. See example in Fig. 2c.

**Direct step-by-step schema:** Here we use the prompt: “What are the steps involved in *topic*? 1.”<sup>4</sup> to allow GPT-3 to generate a schema directly. We run this process once. See example in Fig. 2a.

Generating documents in various genres enables our model to induce comprehensive schemas on any given topic. Considering that some events are more likely to be in the news (e.g., elections, pandemic outbreaks) while others are more technical in nature and are hence less newsworthy (such as earning a Ph.D. degree or planning a wedding), we generate diverse texts and then use a ranking model to choose the most relevant documents. The ranking process includes embedding the texts and the topic with the model of Reimers and Gurevych (2019), then cosine similarity is calculated between each text and the topic embeddings, and only the 30

<sup>4</sup>The “1.” in the prompt is for GPT-3 to automatically complete the steps.

texts closest to the topic are selected together with the output from the direct step-by-step schema.

The following section describes the next step in generating a schema: extracting all the relevant information from the selected texts.

## 4 Information Extraction

For each document, we extract event triggers, arguments and relations between the events that are important and relevant to the schema topic. We do not work with a predefined ontology that defines in advance what events and arguments are essential, so we extract all the information and later filter it down to include just the most frequent items. Here are the steps involved in extracting the information:

1. **Semantic Role Labeling (SRL):** We use the SOTA SRL system<sup>5</sup> trained on CoNLL12 (Pradhan et al., 2012) and Nombank dataset (Meyers et al., 2004) to extract both verb and nominal event triggers and arguments.
2. **Named Entity Recognition (NER):** We employ the SOTA NER model to extract and map entities (potential arguments of events)

<sup>5</sup>Details for the SRL and NER systems were removed for anonymity and will be published upon acceptance.



into entity types defined in the CoNLL 2002 dataset (Tjong Kim Sang, 2002) and the LORELEI project (Strassel and Tracey, 2016).

3. **Constituency Parsing:** Since the arguments extracted by SRL can be clauses and long phrasal nouns, we employ the constituency parsing model from AllenNLP<sup>6</sup> for argument head word extraction. For example, in this sentence “The first passengers rescued from a helicopter that ditched in the North Sea have arrived at hospital,” the ARGM-LOC for “ditched” is “in the North Sea.” However, the NER model can only extract “North Sea” instead of “in the North Sea,” and thus we use the parser to match the argument to its type.

4. **Coreference Resolution:** We use the SOTA model (Yu et al., 2020) for event and entity coreference resolution to identify within-document coreferential relations.

5. **Temporal Relation Extraction:** We first try to use SOTA models (Ning et al., 2019; Zhou et al., 2021) to predict the temporal relations<sup>7</sup> between all possible pairs of extracted events but since the SOTA models accept two sentences containing events as input, the inference time<sup>8</sup> for an  $n$ -event document is  $\mathcal{O}(n^2)$ , making the schema induction process several hours long. We develop a One-Pass model<sup>9</sup> that takes the document as input and uses the contextual representation of events to predict relations between them. As shown in Table 2, the inference time is shortened 63-186 times on average, while the performance of the One-Pass model is comparable to SOTA models.

6. **Hierarchical Relation Extraction:** The extremely long inference time of SOTA models for predicting hierarchical relations (PARENT-CHILD, CHILD-PARENT, COREF, NOREL) (Zhou et al., 2020; Wang et al., 2021a) also impairs the efficiency of our schema induction system. Thus we use the same One-Pass methodology to extract hierarchical relations.

<sup>6</sup><https://demo.allennlp.org/constituency-parsing>.

<sup>7</sup>The possible temporal relations (start-time comparison) are: BEFORE, AFTER, EQUAL and VAGUE.

<sup>8</sup>The inference time is mostly spent on obtaining the contextual representation of events using large fine-tuned LMs.

<sup>9</sup>We take advantage of the recently developed BigBird (Zaheer et al., 2020) that handles long sequences with sparse attention mechanism.

Corpus	Model	Metrics		
		$F_1$ score	Speed	GPU Memory
HiEve	Zhou et al. (2020)	0.489	-	-
	Wang et al. (2021a)	0.522	41.68s	4515MiB
	One-Pass model	0.472	0.65s	2941MiB
MATRES	Ning et al. (2019)	0.767	30.12s	4187MiB
	Zhou et al. (2021)	0.821	89.36s	9311MiB
	One-Pass model	0.768	0.48s	2419MiB

Table 2: Performance comparison between our One-Pass model and SOTA models for event temporal and hierarchical relation extraction. We report  $F_1$  scores on benchmark datasets (HiEve for hierarchical relations, MATRES for temporal relations), speed (average inference time for 100 relations), and required GPU memory during inference. The One-Pass models are 63-186 times faster than SOTA models and take up only 26%-65% of the GPU memory required by SOTA models, while being comparable in performance.

We observe that the inference time is greatly shortened, and the One-Pass model achieves comparable results to previous models, and it takes up less GPU memory (see Table 2).

After processing the data using the procedure described above, we get a list of events, their arguments, and relations between the events. We concentrate on events and relations that frequently appear in the generated texts since we assume those are the most important to add to the schema (without having any other source of information that could identify what is salient). The next section describes the process of building a schema.

## 5 Schema Induction

To consolidate the information extracted from the previous step, we build a schema as follows:

**Make a list of events and relations:** To compare similar event mentions in different texts, we compare the event trigger itself (whether they are the same verb or coreferential verbs<sup>10</sup>) and the NER types of its arguments. For example, the trigger “(take) precautions” appeared in 5 documents generated for the topic of Pandemic Outbreak. In two documents the subject of the verb phrase “take precautions” was “residents”, in another two it was “people” and in the last one, it was “public”. Nevertheless, the NER type is identical in all cases (PER), and thus we set the frequency of “(take) precautions” to 5. Similarly, we calculate the frequency of the temporal and hierarchical relations.

<sup>10</sup>We only consider coreferential relations if they appeared in more than 2 documents.

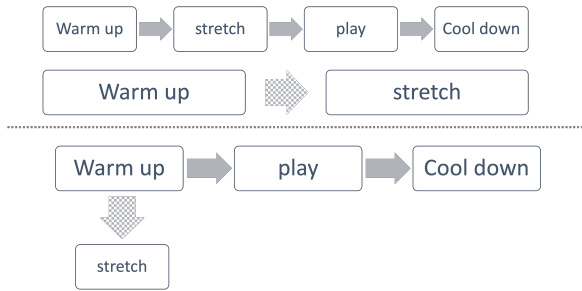


Figure 3: Example of amending a timeline in the schema of “Sports Games”. The timeline at the top includes events of two different levels (“warm up” is the parent of “stretch”), hence it is rectified to include only events of the same level like the timeline at the bottom. Gray arrows mark temporal relations, and checkered arrows denote PARENT-CHILD.

We consider only the top-30 most frequent events and relations for the schema and continue to the next step.

**Construct timelines:** We construct the longest timelines from the list of temporal relations. This list is a list of tuples  $(A, B)$ , indicating that event  $A$  happened before event  $B$ . To construct a timeline, we search recursively for the longest chains of the following form  $(A, B), (B, C), (A, C)$  and so on.

**Fix timelines according to hierarchical relations:** We build a hierarchy of the events using the hierarchical relation list<sup>11</sup> and change the timelines so that they will only include events that appear in the same level of hierarchy (see example in Fig. 3).

**Add logical relations:** The final step is to combine the timelines and hierarchies into a single graph using logical relations (AND/OR). When observing two timelines with discrepancies between the order of events, we place a logical AND between them since we interpret this discrepancy as both events occurring at the same time or there is no significance to the order of those two events. For example, in Fig. 4, the events “call demands” and “clash officers” appear in different orders in different documents, hence we conclude that they occur simultaneously or interleaved. We use a logical OR to mark different outcomes or events that can happen simultaneously but not necessarily. For example, in Fig. 4, the events “police disperse crowd” and “government urge to exercise restraint” may both occur or either one of them occurs.

The final output is a schema graph that contains

<sup>11</sup>We only consider PARENT-CHILD and CHILD-PARENT relations that appear in more than 2 documents.

all the events, arguments, and the temporal, hierarchical and logical relations between the events. This schema generating model can also be used to extend the scope of existing schemas by further querying the model on more specific topics. For example, the schema in Fig. 1 does not cover the consequences of kidnapping, probably because the LM did not attend to this aspect. Hence an analyst can input another topic (e.g. consequences of kidnapping) to further develop the schema. Similarly, analysts can generate schemas for very specific events (e.g., kidnapping in a political setting).

Next, we provide an in-depth experimentation for the proposed schema induction framework.

## 6 Experiments

### 6.1 Data

We conduct experiments on a dataset for general schema learning released by LDC (LDC2020E25). The corpus includes 84 types of complex events, such as Cyber Attack, Farming and Recycling. This dataset includes ground-truth schemas created by LDC annotators.

In addition to the LDC dataset, we also collected human generated schemas for 11 complex events (denoted henceforth as the Schema-11 dataset)<sup>12</sup>. These schemas were generated by four human experts<sup>13</sup> that were instructed to write a schema on each topic based on their commonsense knowledge that includes a list of events, relations<sup>14</sup>, arguments and their NER types<sup>15</sup>.

### 6.2 Evaluation

We follow Li et al. (2021) to use instance coverage and last event prediction to evaluate our method on the LDC dataset; for the Schema-11 dataset, we ask human testers to assess the completeness and soundness of both human- and automatically-generated schemas.

**Coverage and Prediction** A common evaluation method in schema induction and script prediction is to calculate the recall of events and relations

<sup>12</sup>The topics are: Bombing Attack, Business Change, Civil Unrest, Disaster and Rescue, Elections, International Conflict, Kidnapping, Mass Shooting, Pandemic Outbreak, Sports Games, and Terrorism Attack.

<sup>13</sup>Graduate students who are familiar with the research topic of schema induction and are not the authors of this paper.

<sup>14</sup>No restrictions were placed for the annotators. For example, in one case, an annotator mentioned causal relations that are not covered in our framework.

<sup>15</sup>The annotators are familiar with SRL annotations (e.g., ARG0, ARG1, etc.) and NER types (e.g., PER, ORG, etc.).

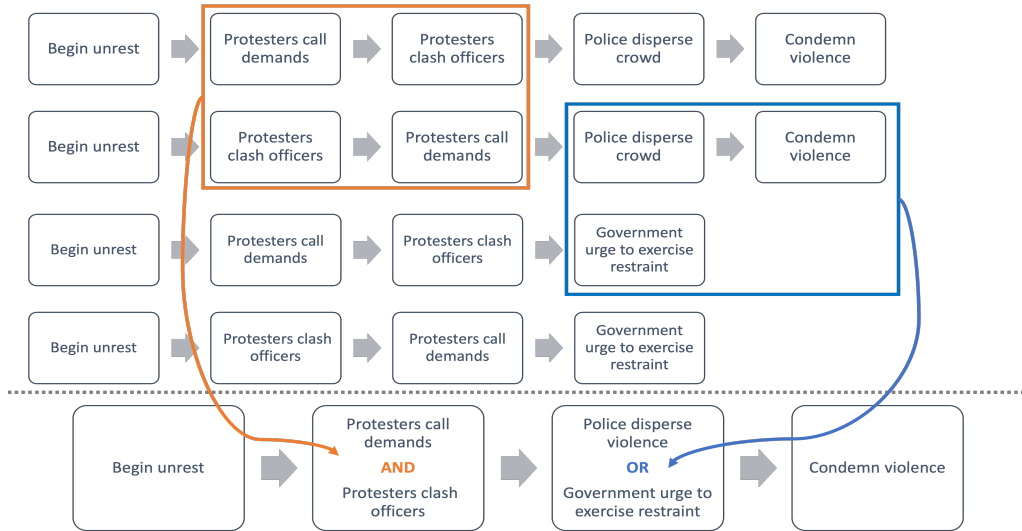


Figure 4: An example of integrating timelines and logical relations in the schema of Civil Unrest. The four upper timelines are the ones extracted from the generated texts and the lower one is their merger into a single timeline with logical relations.

427 predicted by the model, assuming the human an- 454  
 428 notators’ results are gold labels (coverage) and to 455  
 429 calculate the accuracy in predicting the final out- 456  
 430 come of a scenario (prediction). Li et al. (2021), 457  
 431 for example, calculated the accuracy of predict- 458  
 432 ing the last event type of the LDC schemas. Here 459  
 433 we present the results of predicting the last events 460  
 434 using event triggers, instead of event types. 461

435 **Feigenbaum Test** We show human testers two 462  
 436 schemas on each topic in the Schema-11 dataset 463  
 437 (see example in §A). One schema was automati- 464  
 438 cally generated by our model, and the other was 465  
 439 randomly sampled from the Schema-11 corpus<sup>16</sup>. 466

440 We ask the testers to determine which events 467  
 441 and relations are valid to appear in the schema 468  
 442 (soundness) and the following questions: which 469  
 443 schema is more complete in the sense of including 470  
 444 all the events needed to describe the topic, and 471  
 445 which schema, in their opinion, was generated by a 472  
 446 human expert (as opposed to a machine). 473

### 447 6.3 Results 474

448 **Coverage** We calculate the intersection between 475  
 449 events in the generated schemas and the gold 476  
 450 schemas in two ways: (a) match event triggers, 477  
 451 and (b) match event triggers and synonyms of the 478  
 452 events in the gold schemas (synonym coverage)<sup>17</sup>. 479  
 453 We believe that calculating synonym coverage is a 480

<sup>16</sup>In some cases we combine two randomly sampled schemas because the length of the human schemas tend to be shorter than the automatically generated ones.

<sup>17</sup>Implemented using the NLTK WordNet Python package. 481

454 better evaluation methodology to avoid errors such 455  
 456 as considering different verbs describing the same 457  
 458 action as different (e.g., “buy” and “acquire”) than 459  
 459 using a predefined ontology of event types such 460  
 460 as the one used in Li et al. (2021). The reason is 461  
 461 twofold: firstly, any predefined ontology is limited 462  
 462 to certain scenarios and it may impair the variety 463  
 463 of events extracted; and secondly the typing mech- 464  
 464 anism may also inflict errors to the schema. 465

466 From the results in Table 3, we can observe that 467  
 467 despite the difficulty of exact matching, our model 468  
 468 can cover 23.73% events in the gold annotations, 469  
 469 showing that the generated text has a good cover- 470  
 470 age of events required in the schemas. And if we 471  
 471 use synonym coverage as our metric, we achieve 472  
 472 a promising coverage of 36.35% while the state- 473  
 473 of-the-art supervised event graph model (Li et al., 474  
 474 2021) covers 54.84% using limited event types. 475  
 475 Furthermore, with the high quality event represen- 476  
 476 tations obtained from the One-Pass model and the 477  
 477 proposed logical relation induction algorithm, our 478  
 478 method covers 14.09% of all the relations anno- 479  
 479 tated in the gold schemas, whereas the best per- 480  
 480 formance achieved by the event graph model is 481  
 481 44.44%. The high coverage of the SOTA method 482  
 482 can be attributed to the joint modeling of multiple 483  
 483 relations using graph neural networks, which is 484  
 484 impracticable in our zero-shot settings.

**Prediction** In the prediction task, our schemas 485  
 485 are able to predict the final outcome in 46.42% of 486  
 486 the cases for the LDC schemas (see Tab. 4). This 487

	Ours			(Li et al., 2021)
	Coverage	Coverage (Synonym)	Total	Coverage
Event Match	23.73	36.35	36.35	54.84
Temporal Relations	1.99	5.80		
Hierarchical Relations	0.14	0.91	14.09	44.44
Logical Relations	4.56	7.38		

Table 3: Coverage results for the LDC dataset. The first row presents the percentage of events that appeared in both the LDC schemas and the automatically generated schemas (out of events in LDC schemas), and the three bottom rows calculates the same metric for relations of different types. Total is the sum of all three types of relations.

Model	Accuracy
Event Language Model	49.7
Sequential Pattern Mining	47.8
Human Schema	20.5
Event Graph Model	52.0
Zero-Shot Schema	28.5
Zero-Shot Schema Synonym	46.4

Table 4: Experimental results for last event prediction in the LDC dataset. The top 4 results are from (Li et al., 2021), and the metric is HITS@1 where the events are typed based on a predefined ontology.

result is extremely impressive when it is compared with Li et al. (2021) since they predict event types instead of verbs, which is a much easier task due to the fact that the set of possible answers is limited.

**Feigenbaum test** In the soundness experiments, where the testers are asked to decide which events and relations are valid to appear in the schema, it turns out that human generated schemas contain 7.14% invalid events and 15.4% invalid relations on average. For the automatically generated schemas, 6.06% of the events and 22.9% of the relations are considered to be invalid on average. We observe that the average percentage of valid events is higher in the automated schemas, yet the soundness of induced relations is relative inferior.

For the completeness results, in 4 cases the testers agreed that the automatically generated schemas are more complete; in 3 cases they claimed that the human schemas are more complete; and the result is a tie in the remaining 4 cases. The distribution of votes for completeness is presented in Tab. 5. Hence our automatically generated schemas are of comparable quality to human generated ones in the sense of completeness.

Finally, in the Feigenbaum test, where testers are asked to decide whether a schema is generated by a human or a machine, 8 out of 11 times they correctly identify the human-generated schema, 1 incorrectly, and 2 ties. Some of the testers who

	S1	S2	S3	S4	S5	S6	S7	S8	S9	S10	S11
Human	4	0	1	1	1	2	1	1	0	3	1
Automatic	2	3	4	2	1	1	1	1	4	0	1

Table 5: Distribution of votes for which is the more complete schema for Schema-11 dataset.

succeeded in their guesses mentioned that it was easy to determine which schema was automatically generated since it tends to be longer and more complete. Although in the test, the machine-generated schemas fail to deceive the testers into misidentifying them as human generated ones, the experiments shed light on future directions, e.g., keeping the most salient events in the schema, improving the accuracy of temporal and hierarchical relation extraction, developing reliable approaches for causal relation extraction, and so forth. The full results from the Feigenbaum test appear in §B.

## 7 Conclusion

We propose a method to generate schemas given the sole input of a topic. We use GPT-3 to generate texts of diverse genres and a pipeline of information extraction tools to obtain relevant information before inducing logical relations and integrating the events and relations into a schema graph. To improve the efficiency of the pipeline, we implement One-Pass models for event temporal and hierarchical relations that achieve comparable performances with SOTA models but require far less inference time and GPU memory space. To evaluate our framework, we conduct experiments on the benchmark LDC dataset to show that our schemas cover a decent amount of pertinent information and display comparable ability for event prediction with supervised approaches. Although our proposed method fails the Feigenbaum test on Schema-11, we observe a very high percentage of valid events and relations and the testers endorsed the completeness of our machine-generated schemas.



## 8 Ethical Consideration

The proposed schema induction method does not present any direct societal implications. As is observed in [Abid et al. \(2021\)](#), the text generated by GPT-3 might include undesired social bias. Extracting events and relations from text with such social bias might potentially propagate the bias to the induced schemas. Besides, there are risks of malicious or unintended harmful uses of the generated schemas, for instance, the system might be used to inquire about making a bomb or contriving a terrorist attacks. Yet we believe that the proposed method can benefit various downstream NLP/NLU tasks like event prediction, task-oriented dialogue agents ([Andreas et al., 2020](#)) and risk detection ([Pohl et al., 2012](#)).

## References

Abubakar Abid, Maheen Farooqi, and James Zou. 2021. Large language models associate muslims with violence. *Nature Machine Intelligence*, 3(6):461–463.

Jacob Andreas, John Bufe, David Burkett, Charles Chen, Josh Clausman, Jean Crawford, Kate Crim, Jordan DeLoach, Leah Dorner, Jason Eisner, Hao Fang, Alan Guo, David Hall, Kristin Hayes, Kellie Hill, Diana Ho, Wendy Iwaszuk, Smriti Jha, Dan Klein, Jayant Krishnamurthy, Theo Lanman, Percy Liang, Christopher H. Lin, Ilya Lintsbakh, Andy McGovern, Aleksandr Nisnevich, Adam Pauls, Dmitriy Petters, Brent Read, Dan Roth, Subhro Roy, Jesse Rusak, Beth Short, Div Slomin, Ben Snyder, Stephon Striplin, Yu Su, Zachary Tellman, Sam Thomson, Andrei Vorobev, Izabela Witoszko, Jason Wolfe, Abby Wray, Yuchen Zhang, and Alexander Zotov. 2020. [Task-oriented dialogue as dataflow synthesis](#). *Transactions of the Association for Computational Linguistics*, 8:556–571.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.

Tommaso Caselli and Piek Vossen. 2017. [The event StoryLine corpus: A new benchmark for causal and temporal relation extraction](#). In *Proceedings of the*

*Events and Stories in the News Workshop*, pages 77–86, Vancouver, Canada. Association for Computational Linguistics. 601–602–603

Nathanael Chambers. 2013. [Event schema induction with a probabilistic entity-driven model](#). In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 1797–1807, Seattle, Washington, USA. Association for Computational Linguistics. 604–605–606–607–608–609

Nathanael Chambers and Dan Jurafsky. 2008. [Unsupervised learning of narrative event chains](#). In *Proceedings of ACL-08: HLT*, pages 789–797, Columbus, Ohio. Association for Computational Linguistics. 610–611–612–613

Nathanael Chambers and Dan Jurafsky. 2009. [Unsupervised learning of narrative schemas and their participants](#). 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*, pages 602–610, Suntec, Singapore. Association for Computational Linguistics. 614–615–616–617–618–619–620–621

Jackie Chi Kit Cheung, Hoifung Poon, and Lucy Vanderwende. 2013. [Probabilistic frame induction](#). In *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 837–846, Atlanta, Georgia. Association for Computational Linguistics. 622–623–624–625–626–627–628

Quang Do, Wei Lu, and Dan Roth. 2012. [Joint inference for event timeline construction](#). In *Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning*, pages 677–687, Jeju Island, Korea. Association for Computational Linguistics. 629–630–631–632–633–634–635

Edward A Feigenbaum. 2003. Some challenges and grand challenges for computational intelligence. *Journal of the ACM (JACM)*, 50(1):32–40. 636–637–638

Goran Glavaš, Jan Šnajder, Marie-Francine Moens, and Parisa Kordjamshidi. 2014. [HiEve: A corpus for extracting event hierarchies from news stories](#). In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC’14)*, pages 3678–3683, Reykjavik, Iceland. European Language Resources Association (ELRA). 639–640–641–642–643–644–645

Mark Granroth-Wilding and Stephen Clark. 2016. [What happens next? event prediction using a compositional neural network model](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 30(1). 646–647–648–649

Rujun Han, Qiang Ning, and Nanyun Peng. 2019. [Joint event and temporal relation extraction with shared representations and structured prediction](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 434–444, Hong Kong, China. Association for Computational Linguistics. 650–651–652–653–654–655–656–657–658

659	Lifu Huang, Heng Ji, Kyunghyun Cho, Ido Dagan, Sebastian Riedel, and Clare Voss. 2018. <a href="#">Zero-shot transfer learning for event extraction</a> . In <i>Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 2160–2170, Melbourne, Australia. Association for Computational Linguistics.	717
660		718
661		
662		
663		719
664		720
665		721
666	Bram Jans, Steven Bethard, Ivan Vulić, and Marie Francine Moens. 2012. <a href="#">Skip n-grams and ranking functions for predicting script events</a> . In <i>Proceedings of the 13th Conference of the European Chapter of the Association for Computational Linguistics</i> , pages 336–344, Avignon, France. Association for Computational Linguistics.	722
667		723
668		724
669		725
670		726
671		
672		
673	Heeyoung Lee, Marta Recasens, Angel Chang, Mihai Surdeanu, and Dan Jurafsky. 2012. <a href="#">Joint entity and event coreference resolution across documents</a> . In <i>Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning</i> , pages 489–500, Jeju Island, Korea. Association for Computational Linguistics.	727
674		728
675		729
676		730
677		731
678		
679		
680		
681	Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. <a href="#">BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension</a> . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 7871–7880, Online. Association for Computational Linguistics.	732
682		733
683		734
684		735
685		736
686		737
687		738
688		739
689		740
690	Manling Li, Sha Li, Zhenhailong Wang, Lifu Huang, Kyunghyun Cho, Heng Ji, Jiawei Han, and Clare Voss. 2021. <a href="#">The future is not one-dimensional: Complex event schema induction by graph modeling for event prediction</a> . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 5203–5215, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	741
691		742
692		743
693		744
694		745
695		746
696		
697		
698		
699	Manling Li, Qi Zeng, Ying Lin, Kyunghyun Cho, Heng Ji, Jonathan May, Nathanael Chambers, and Clare Voss. 2020. <a href="#">Connecting the dots: Event graph schema induction with path language modeling</a> . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 684–695, Online. Association for Computational Linguistics.	747
700		748
701		749
702		750
703		751
704		752
705		753
706		754
707	Ying Lin, Heng Ji, Fei Huang, and Lingfei Wu. 2020. <a href="#">A joint neural model for information extraction with global features</a> . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 7999–8009, Online. Association for Computational Linguistics.	755
708		756
709		757
710		758
711		759
712		760
713	Wei Lu and Dan Roth. 2012. <a href="#">Automatic event extraction with structured preference modeling</a> . In <i>Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 835–844, Jeju Island, Korea. Association for Computational Linguistics.	717
714		718
715		
716		
	Adam Meyers, Ruth Reeves, Catherine Macleod, Rachel Szekely, Veronika Zielinska, Brian Young, and Ralph Grishman. 2004. <a href="#">Annotating noun argument structure for NomBank</a> . In <i>Proceedings of the Fourth International Conference on Language Resources and Evaluation (LREC’04)</i> , Lisbon, Portugal. European Language Resources Association (ELRA).	719
		720
		721
		722
		723
		724
		725
		726
	Ashutosh Modi. 2016. <a href="#">Event embeddings for semantic script modeling</a> . In <i>Proceedings of The 20th SIGNLL Conference on Computational Natural Language Learning</i> , pages 75–83, Berlin, Germany. Association for Computational Linguistics.	727
		728
		729
		730
		731
	Kiem-Hieu Nguyen, Xavier Tannier, Olivier Ferret, and Romaric Besançon. 2015. <a href="#">Generative event schema induction with entity disambiguation</a> . In <i>Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 188–197, Beijing, China. Association for Computational Linguistics.	732
		733
		734
		735
		736
		737
		738
		739
		740
	Qiang Ning, Zhili Feng, Hao Wu, and Dan Roth. 2018. <a href="#">Joint reasoning for temporal and causal relations</a> . In <i>Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 2278–2288, Melbourne, Australia. Association for Computational Linguistics.	741
		742
		743
		744
		745
		746
	Qiang Ning, Sanjay Subramanian, and Dan Roth. 2019. <a href="#">An improved neural baseline for temporal relation extraction</a> . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 6203–6209, Hong Kong, China. Association for Computational Linguistics.	747
		748
		749
		750
		751
		752
		753
		754
	Haoruo Peng and Dan Roth. 2016. <a href="#">Two discourse driven language models for semantics</a> . In <i>Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 290–300, Berlin, Germany. Association for Computational Linguistics.	755
		756
		757
		758
		759
		760
	Haoruo Peng, Yangqiu Song, and Dan Roth. 2016. <a href="#">Event detection and co-reference with minimal supervision</a> . In <i>Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing</i> , pages 392–402, Austin, Texas. Association for Computational Linguistics.	761
		762
		763
		764
		765
		766
	Karl Pichotta and Raymond Mooney. 2014. <a href="#">Statistical script learning with multi-argument events</a> . In <i>Proceedings of the 14th Conference of the European Chapter of the Association for Computational Linguistics</i> , pages 220–229, Gothenburg, Sweden. Association for Computational Linguistics.	767
		768
		769
		770
		771
		772

773	Karl Pichotta and Raymond Mooney. 2016. <a href="#">Learning statistical scripts with lstm recurrent neural networks</a> . <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , 30(1).	
774		
775		
776		
777	Daniela Pohl, Abdelhamid Bouchachia, and Hermann Hellwagner. 2012. Automatic sub-event detection in emergency management using social media. In <i>Proceedings of the 21st international conference on world wide web</i> , pages 683–686.	
778		
779		
780		
781		
782	Sameer Pradhan, Alessandro Moschitti, Nianwen Xue, Olga Uryupina, and Yuchen Zhang. 2012. <a href="#">CoNLL-2012 shared task: Modeling multilingual unrestricted coreference in OntoNotes</a> . In <i>Joint Conference on EMNLP and CoNLL - Shared Task</i> , pages 1–40, Jeju Island, Korea. Association for Computational Linguistics.	
783		
784		
785		
786		
787		
788		
789	Vasin Punyakanok, Dan Roth, and Wen-tau Yih. 2008. <a href="#">The importance of syntactic parsing and inference in semantic role labeling</a> . <i>Computational Linguistics</i> , 34(2):257–287.	
790		
791		
792		
793	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. <a href="#">Exploring the limits of transfer learning with a unified text-to-text transformer</a> . <i>Journal of Machine Learning Research</i> , 21(140):1–67.	
794		
795		
796		
797		
798		
799	Nils Reimers and Iryna Gurevych. 2019. <a href="#">Sentence-BERT: Sentence embeddings using Siamese BERT-networks</a> . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 3982–3992, Hong Kong, China. Association for Computational Linguistics.	
800		
801		
802		
803		
804		
805		
806		
807	Rachel Rudinger, Pushpendre Rastogi, Francis Ferraro, and Benjamin Van Durme. 2015. <a href="#">Script induction as language modeling</a> . In <i>Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing</i> , pages 1681–1686, Lisbon, Portugal. Association for Computational Linguistics.	
808		
809		
810		
811		
812		
813	Lei Sha, Sujian Li, Baobao Chang, and Zhifang Sui. 2016. <a href="#">Joint learning templates and slots for event schema induction</a> . In <i>Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 428–434, San Diego, California. Association for Computational Linguistics.	
814		
815		
816		
817		
818		
819		
820		
821	Stephanie Strassel and Jennifer Tracey. 2016. <a href="#">LORELEI language packs: Data, tools, and resources for technology development in low resource languages</a> . In <i>Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC’16)</i> , pages 3273–3280, Portorož, Slovenia. European Language Resources Association (ELRA).	
822		
823		
824		
825		
826		
827		
828		
	Erik F. Tjong Kim Sang. 2002. <a href="#">Introduction to the CoNLL-2002 shared task: Language-independent named entity recognition</a> . In <i>COLING-02: The 6th Conference on Natural Language Learning 2002 (CoNLL-2002)</i> .	829 830 831 832 833
	David Wadden, Ulme Wennberg, Yi Luan, and Hananeh Hajishirzi. 2019. <a href="#">Entity, relation, and event extraction with contextualized span representations</a> . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 5784–5789, Hong Kong, China. Association for Computational Linguistics.	834 835 836 837 838 839 840 841 842
	Haoyu Wang, Muhao Chen, Hongming Zhang, and Dan Roth. 2020. <a href="#">Joint constrained learning for event-event relation extraction</a> . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 696–706, Online. Association for Computational Linguistics.	843 844 845 846 847 848
	Haoyu Wang, Hongming Zhang, Muhao Chen, and Dan Roth. 2021a. <a href="#">Learning constraints and descriptive segmentation for subevent detection</a> . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 5216–5226, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	849 850 851 852 853 854 855
	Zirui Wang, Adams Wei Yu, Orhan Firat, and Yuan Cao. 2021b. <a href="#">Towards zero-label language learning</a> . <i>arXiv preprint arXiv:2109.09193</i> .	856 857 858
	Noah Weber, Niranjana Balasubramanian, and Nathanael Chambers. 2018. <a href="#">Event representations with tensor-based compositions</a> . <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , 32(1).	859 860 861 862 863
	Noah Weber, Rachel Rudinger, and Benjamin Van Durme. 2020. <a href="#">Causal inference of script knowledge</a> . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 7583–7596, Online. Association for Computational Linguistics.	864 865 866 867 868 869
	Haoyang Wen, Ying Lin, Tuan Lai, Xiaoman Pan, Sha Li, Xudong Lin, Ben Zhou, Manling Li, Haoyu Wang, Hongming Zhang, Xiaodong Yu, Alexander Dong, Zhenhailong Wang, Yi Fung, Piyush Mishra, Qing Lyu, Dídac Surís, Brian Chen, Susan Windisch Brown, Martha Palmer, Chris Callison-Burch, Carl Vondrick, Jiawei Han, Dan Roth, Shih-Fu Chang, and Heng Ji. 2021. <a href="#">RESIN: A dockerized schema-guided cross-document cross-lingual cross-media information extraction and event tracking system</a> . In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies: Demonstrations</i> , pages 133–143, Online. Association for Computational Linguistics.	870 871 872 873 874 875 876 877 878 879 880 881 882 883 884



885 Xiaodong Yu, Wenpeng Yin, and Dan Roth. 2020.  
886 Paired representation learning for event and entity  
887 coreference. *arXiv preprint arXiv:2010.12808*.

888 Quan Yuan, Xiang Ren, Wenqi He, Chao Zhang, Xinhe  
889 Geng, Lifu Huang, Heng Ji, Chin-Yew Lin, and Ji-  
890 awei Han. 2018. Open-schema event profiling for  
891 massive news corpora. In *Proceedings of the 27th  
892 ACM International Conference on Information and  
893 Knowledge Management*, pages 587–596.

894 Manzil Zaheer, Guru Guruganesh, Kumar Avinava  
895 Dubey, Joshua Ainslie, Chris Alberti, Santiago On-  
896 tanon, Philip Pham, Anirudh Ravula, Qifan Wang,  
897 Li Yang, and Amr Ahmed. 2020. **Big bird: Trans-  
898 formers for longer sequences**. In *Advances in  
899 Neural Information Processing Systems*, volume 33,  
900 pages 17283–17297. Curran Associates, Inc.

901 Hongming Zhang, Muhao Chen, Haoyu Wang,  
902 Yangqiu Song, and Dan Roth. 2020a. **Analogous  
903 process structure induction for sub-event sequence  
904 prediction**. In *Proceedings of the 2020 Conference  
905 on Empirical Methods in Natural Language Process-  
906 ing (EMNLP)*, pages 1541–1550, Online. Associa-  
907 tion for Computational Linguistics.

908 Hongming Zhang, Haoyu Wang, and Dan Roth. 2021a.  
909 **Zero-shot Label-aware Event Trigger and Argu-  
910 ment Classification**. In *Findings of the Association  
911 for Computational Linguistics: ACL-IJCNLP 2021*,  
912 pages 1331–1340, Online. Association for Computa-  
913 tional Linguistics.

914 Li Zhang, Qing Lyu, and Chris Callison-Burch. 2020b.  
915 **Reasoning about goals, steps, and temporal ordering  
916 with WikiHow**. In *Proceedings of the 2020 Con-  
917 ference on Empirical Methods in Natural Language  
918 Processing (EMNLP)*, pages 4630–4639, Online. As-  
919 sociation for Computational Linguistics.

920 Yi Zhang, Sujay Kumar Jauhar, Julia Kiseleva, Ryen  
921 White, and Dan Roth. 2021b. **Learning to decom-  
922 pose and organize complex tasks**. In *Proceedings of  
923 the 2021 Conference of the North American Chap-  
924 ter of the Association for Computational Linguistics:  
925 Human Language Technologies*, pages 2726–2735,  
926 Online. Association for Computational Linguistics.

927 Ben Zhou, Qiang Ning, Daniel Khashabi, and Dan  
928 Roth. 2020. **Temporal common sense acquisition  
929 with minimal supervision**. In *Proceedings of the  
930 58th Annual Meeting of the Association for Compu-  
931 tational Linguistics*, pages 7579–7589, Online. As-  
932 sociation for Computational Linguistics.

933 Ben Zhou, Kyle Richardson, Qiang Ning, Tushar Khot,  
934 Ashish Sabharwal, and Dan Roth. 2021. **Tempo-  
935 ral reasoning on implicit events from distant super-  
936 vision**. In *Proceedings of the 2021 Conference of  
937 the North American Chapter of the Association for  
938 Computational Linguistics: Human Language Tech-  
939 nologies*, pages 1361–1371, Online. Association for  
940 Computational Linguistics.

## A Feigenbaum Test Details 941

The experiment took place online through filling 942  
a Google Form and involved 11 annotators. Each 943  
annotator got 3-4 scenarios to annotate. The in- 944  
structions for the survey appear in Figure 5. An 945  
example scenario and the questions of the survey 946  
are presented in Figures 6, 7, 8, and 9. 947

## B Feigenbaum Test Results 948

In this section we present all the results from the ex- 949  
periments on the dataset Schema-11. Table 6 shows 950  
the distribution of answers for the question “which 951  
schema is more complete?” (same as depicted in 952  
Table 5), Table 7 presents the distribution of an- 953  
swers for the question “which schema was gener- 954  
ated by a human?” together with the correct answer 955  
written in the bottom row, and Table 8 presents 956  
the percentage of invalid events and relations deter- 957  
mined by the majority vote of the annotators in the 958  
automatic schema and the human schema. 959



# Feigenbaum Test - Scenario 11

This form mainly focuses on the evaluation of machine generated schema. Given a certain scenario, the schema includes stereotypical events and the relations between them, for instance, within scenario "acquiring a PhD degree", a schema would typically includes "publish papers," "attend conferences," "write PhD thesis" and "defend PhD thesis." And there is also a "before" relation between "write PhD thesis" and "defend PhD thesis." Besides, we also have "SuperSub" relation that means hierarchical relation between events, and "AND"/"OR" relation that means the two events must happen together/either of the events may happen.

We've asked a group of people to generate schemas from their commonsense knowledge. Given two schemas per scenario, your task is to determine whether you can distinguish the machine generated schema from the human generated one. And also provide your insights on the completeness and soundness of each schema.

For completeness, we would like you to tell us which schema is more complete.

For soundness, we would like you to tell us for each event and relation listed, whether it is valid for this scenario.

Most importantly, we would like to know which schema you think is generated by human.

Figure 5: Instructions for the Feigenbaum test.

	S1	S2	S3	S4	S5	S6	S7	S8	S9	S10	S11
Human	4	0	1	1	1	2	1	1	0	3	1
Automatic	2	3	4	2	1	1	1	1	4	0	1

Table 6: Completeness results. The table presents the number of votes that were recorded for which schema is more complete - the human generated schema or the automatically generated schema. The majority vote is highlighted in yellow.

	S1	S2	S3	S4	S5	S6	S7	S8	S9	S10	S11
A	1	1	3	0	0	2	0	2	2	1	1
B	5	2	2	3	2	1	2	0	2	2	1
Correct Answer	B	B	B	B	B	A	B	A	A	B	B

Table 7: Feigenbaum test results. The annotators guesses which schema (A or B) was generated by humans. The number of votes for each option appear along with the correct answer in the bottom row. The correct majority guesses are marked with green and incorrect with red.

	S1	S2	S3	S4	S5	S6	S7	S8	S9	S10	S11
Invalid Events (Auto.)	0	0	0	0	0	8.33	0	7.69	0	14.28	0
Invalid Relations (Auto.)	46.15	16.66	25	25	0	23.52	0.4	11.76	12.5	22.22	46.15
Invalid Events (Human)	0	0	14.28	14.28	0	0	0	0	0	0	0
Invalid Relations (Human)	7.69	50	15.38	15.38	0	6.25	0	11.11	0	10	7.69

Table 8: Invalidity results. The table presents the percentage of invalid events and relations determined by the human annotators for each schema and scenario.

### Scenario 11: Terrorism Attack (A)

#### Events:

1. event: kill, arg0: {PER, ORG, VEH, WEA}, arg1: PER
2. event: injure, arg0: {PER, ORG, VEH, WEA}, arg1: PER
3. event: detonate, arg0: PER, arg1: WEA
4. event: come, arg1: attack
5. event: open, arg0: {PER, ORG}, arg1: fire
6. event: wound, arg0: {PER, ORG, VEH, WEA}, arg1: PER
7. event: strike, arg0: {PER, ORG, WEA}
8. event: claim, arg0: ORG, arg1: responsibility
9. event: leave, arg0: {PER, VEH}
10. event: attack, arg0: {PER, ORG}
11. event: choose, arg0: {PER, ORG}, arg1: {PER, ORG, GPE}
12. event: select, arg0: {PER, ORG}, arg1: method
13. event: acquire, arg0: {PER, ORG}, arg1: WEA
14. event: carry out, arg0: PER

#### Relations:

1. before: 3->8
2. before: 3->5
3. before: 1->4
4. before: 1->9
5. before: 2->4
6. before: 2->9
7. before: 6->4
8. before: 6->9
9. before: 11->12->13->14->10
10. OR: 8,5
11. OR: 4,9
12. AND: 1,2,6
13. supersub: 10->7->1,2

Figure 6: An example schema in the topic of Terrorism Attack. This schema was generated automatically (information that was unknown to the annotators).

## Scenario 11: Terrorism Attack (B)

### Events:

1. event: find, arg0: PER, arg1: ORG, arg-loc:LOC
2. event: emerge, arg0: ORG, arg1: ORG
3. event: fade, arg0: ORG, arg-tmp: TMP
4. event: reemerge, arg0: ORG, arg-tmp: TMP, arg-loc: LOC
5. event: lead, arg0: ORG, arg1: losses
6. event: lost, arg0: ORG, arg1: LOC
7. event: declare, arg0: GPE, arg1: ORG
8. event: kill, arg0: GPE, arg1: PER
9. event: plan, arg0: PER
10. event: executes, arg0: PER
11. event: injures, arg0: the attack, arg1: PER
12. event: kills, arg0: the attack, arg1: PER
13. event: damages, arg0: the attack, arg1: infrastructure
14. event: calls, arg0: PER, arg1: PER
15. event: arrive, arg0: PER
16. event: treat, arg0: PER, arg1: PER
17. event: take, arg0: PER, arg1: PER
18. event: reports, arg0: PER
19. event: claims, arg0: the group, arg1: responsibility

### Relations:

1. before: 9->10
2. before: 10->11
3. before: 10->12
4. before: 10->13
5. before: 10->14
6. before: 14->15->16->17
7. before: 10->18
8. before: 10->19
9. before: 1->3
10. before: 3->4
11. AND: 1->2
12. cause: 5->6
13. cause: 8->7

Figure 7: An example schema in the topic of Terrorism Attack. This schema was generated by a human (information that was unknown to the annotators).

Which schema is more complete? \*

(A)

(B)

---

Which one do you think is generated by human? \*

(A)

(B)

Figure 8: Questions that were asked about the completeness of the schemas and the generator of the schema.

For each EVENT in schema (B), select if it is valid to appear in the scenario. \*

	valid	invalid
Row 1	<input type="radio"/>	<input type="radio"/>
Row 2	<input type="radio"/>	<input type="radio"/>

Figure 9: Questions about the validity of the events appearing in one of the schemas. This question was asked on both schemas and on the relations appearing in the schemas too.