

# TKGT: Redefinition and A New Way of Text-to-Table Tasks Based on Real World Demands and Knowledge Graphs Augmented LLMs

Anonymous ACL submission

## Abstract

The task of text-to-table receives widespread attention, but its importance and difficulty are underestimated. Existing works use simple datasets like those from table-to-text tasks and employ methods that ignore domain structures. As a bridge between raw text and statistical analysis, the text-to-table task faces challenges from more complex semi-structured texts that refer to certain domain topics in the real world with obvious entities and events, especially from those of social sciences. In this paper, we analyse the limitation of previous datasets with methods and redefine the text-to-table task, based on which we propose a new dataset called **CPL** (Chinese Private Lending) of case judgments from a real world legal academic project. We further propose TKGT (**T**ext-**KG**-**T**able), a two stages domain-aware pipeline, which firstly generates domain knowledge graphs (KGs) classes semi-automatically from raw text with the mixed information extraction (Mixed-IE) method, then adopts the hybrid retrieval augmented generation (Hybird-RAG) method to transform it to tables for downstream needs under the guidance of KGs classes. Experiment results show that TKGT achieves state-of-the-art (SOTA) performance on both traditional datasets and the CPL. Our code and data are available at <https://anonymous.4open.science/r/TKGT-4755>.

## 1 Introduction

Extracting structured information from unstructured or semi-structured text is significant for Natural Language Processing (NLP), as it means extracting valuable information through rule-based, statistical, or deep learning (DL) methods to compress texts and facilitate downstream application (Li et al., 2023a; Sui et al., 2024; Pan et al., 2024).

\*Equal contribution.

†Co-corresponding authors.

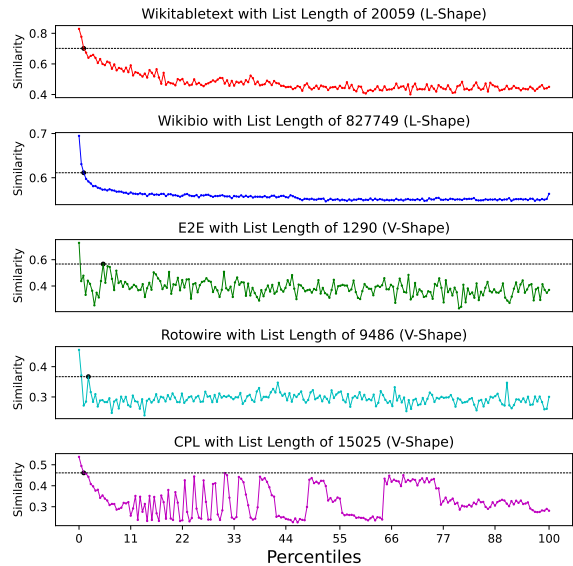


Figure 1: Statistical results of four text-to-table datasets and our **CPL**. The horizontal axis represents the percentile of the ordered word frequency lists, and the vertical axis represents the maximum similarity between each word and datasets’ field sets. The intersection point is the maximum value point after 1% of each list.

Recently, with the development of deep learning (DL) especially the LLMs, some works explore the potential for Transformer models to revolutionize traditional IE (Lu et al., 2022; Wang et al., 2023; Ni et al., 2023), while some directly focus on transforming raw text to structured forms such as KGs (Kommineni et al., 2024; Meyer et al., 2023), mind maps (Jain et al., 2024), and tables (Wu et al., 2021; Li et al., 2023b; Sundar et al., 2024; Deng et al., 2024), among which table is the most popular form.

However, the importance and difficulty of text-to-table tasks are underestimated. Datasets currently used are often structurally simple, fictional, and not from *real world demands*. As shown in Table 1, the first four datasets used in current text-to-table tasks share features that the numbers of average words per document and fields are small.

Datasets	DN	OT	TW	AW/D	TFW(%)	TF	TVTF
Wikitabletext	13318	Entity	185111	13.90	50.04%	2443	2262 / 791 / 1022
Wikibio	728221	Entity	70257683	96.48	45.22%	2996	2771 / 1400 / 1406
E2E	51426	Entity	1152364	22.41	49.04%	7	7 / 7 / 7
Rotowire	4853	Event	1637820	337.49	39.97%	33	33 / 33 / 33
<b>CPL</b>	850	Event	1149207	1105.94	65.58%	97	97 / 97

Table 1: Profiles of five datasets, first four ones in which are originally from table-to-text tasks (Wiseman et al., 2017; Novikova et al., 2017; Bao et al., 2018; Lebret et al., 2016) respectively and pre-processed by (Wu et al., 2021). Abbreviations are used for title, in which DN means document numbers, OT means object type, TW means total words, AW/D means average words per document, TFW means proportions after filtering, TF means total fields and are divided into three parts of train, validation, test respectively in TVTF. CPL has no validation set.

In addition, the two datasets from Wikipedia are essentially relationship extraction (RE) due to the lack of determined and refined fields. Recent work (Deng et al., 2024) proposes a new dataset that generates summary tables of sports competitions from commentary text. However, such a task is still distant from real-world applications.

In contrast, tabular data are important foundations for quantitative statistical analysis, holding tremendous value in various fields, including business intelligence (Vidal-García et al., 2019), natural sciences (Hey et al., 2009), and social sciences (King, 2014). For social scientists adopting the computational social science (CSS) paradigm (Lazer et al., 2009), there is an increasingly urgent need to efficiently extract meaningful information from unstructured or semi-structured texts and store it as tabular data (Gentzkow et al., 2019). This demand is expanding from CSS fields, such as economics (Ash and Hansen, 2023), political science (Grossman and Pedahzur, 2020), and law (Ashley, 2017), to digital humanities disciplines, including history and literature (Michel et al., 2011). Therefore, we redefine the requirements of text-to-table tasks and propose a new dataset called CPL (in Section 2) to fill the gap between existing datasets and real-world demands.

Besides, corresponding methods on previous data remain problems. Text-to-table is initially modeled as Seq2Seq tasks (Wu et al., 2021; Li et al., 2023b), embedding tokens to data-driven learn inner similarities and generate table rows end-to-end; Further researches include inferring table fields (Sundar et al., 2024) before traversing texts with RE and merging finally (Deng et al., 2024). Some works also utilize structures of text and hope to reduce difficulty through segmentation (Jain et al., 2024). After the emergence of LLMs, question and

answer (Q&A) is explored as an approach for IE (Wang et al., 2023; Ni et al., 2023). However, existing works ignore the importance and difficulty of building table fields and treat them as known or just extract triples by simply crawling, which are only applicable to simple formats since that identifying valuable information in complex texts and building fields themselves require professional efforts. Besides, it’s challenging to guarantee completeness, especially for long texts whose valuable points may scatter globally or in the disguise of multiple perspectives that are common in real world.

We propose TKG (Text-KG-Table), a two-stage text-to-table method with KGs as middle-ware. **In the first stage**, the Mixed-IE method based on regulations, statistics, and DL is used to obtain topic keywords and to construct domain KGs sketch, based on which users can better understand the datasets and easily form uninstantiated KGs adapting to downstream tabular needs using LLMs. **In the second stage**, based on dynamic prompts and Hybrid-RAG supported by descriptions of empty KGs classes, table content can be filled with LLMs Q&A. Through experiments, TKG achieves SOTA performance on both traditional datasets and CPL. Our contributions are summarized as follows:

- Redefine the characteristics and requirements of text-to-table tasks in a more standardized manner and introduce the CPL, a new and highly challenging manually completed dataset in the field of law.
- Propose the two-stage TKG, filling the gap in how to obtain table fields based on domain topic structures and use the Hybrid-RAG to fill the table with Q&A. We also demonstrate its SOTA performance through experiments.

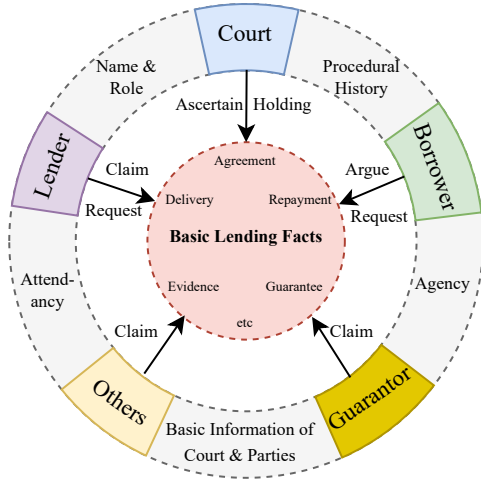


Figure 2: Overview of CPL dataset, which include five role types who have own view for the case facts in the outer layer, and a key contents set of the case in the inner layer.

## 2 Redefinition and CPL Dataset

The CPL dataset is from a real-world academic project, whose raw texts are from the China Judgments Online (CJO)<sup>1</sup>, collected manually by legal experts (Appendix A).

### 2.1 Redefinition of Text-to-Table

For the convenience of further research, we redefine the requirements of text-to-table as follows. **Firstly**, table fields must be limited and refined to serve a practical need rather than unrestricted key-value pairs. **Secondly**, topic information must be clear, and content can be modeled as multi-attribute entities or multi-entity events. **Thirdly**, the information is relatively complex, requiring certain writing formats with logic for clear organization.

As shown in Figure 1, the maximum similarity curves of Rotowire, E2E and CPL present a **V-shape** pattern that first decreases, then rebounds and oscillates after one percent position at lists, which indicates that there exists not only the field information at the front of lists, but also the shared structural information dissimilar with fields on the semantic meaning. In contrast, curves of the two datasets from Wikipedia consistently decrease as **L-shape**, indicating no obvious structural information and explaining why the field numbers of the two datasets are so large and inconsistent in Table 1.

<sup>1</sup>CJO, established by the Supreme People’s Court of the People’s Republic of China(SPC), allows the public to freely search, read, download, and analyze cases.<https://wenshu.court.gov.cn/>

## 2.2 Statistics of CPL

The CPL dataset contains 850 judgment documents and corresponding tables. **Firstly**, it is a typical event-type dataset, which including one lender, one court, at least one borrower, zero or several guarantors, and other roles like witnesses (Figure 2). As shown in Table 1, it has 1149209 words in total and 1105.94 in each document on average. Fields in this dataset are scalable to fit multiple lending in a case. Actual field numbers depends on specific text contents and exceeds 220 overall, which are ultimately abstracted into 97 core fields considering reusable concepts such as interest and penalty sharing attributes like start date and interest rate. **Secondly**, to reduce the complexity of subsequent works, we filter out stop words and stop position tag, leaving behind 753610 core words, accounting for 65.58% of the total, which is much higher than the other four datasets filtered based on the same strategy (Appendix C). **Thirdly**, as shown in Figure 1, this dataset shows a significant **V-shape**, which is similar to the other two datasets with table structure (Rotowire and E2E). In short, this dataset has longer text, more complex field structures, higher word quality, and distinct semi-structured features.

## 3 TKG Two-Stages Pipeline

### 3.1 Overview

As illustrated in Figure 3, TKG uses KGs classes as middleware to transform raw texts to tables through two stages. The first stage aims at semi-automatically assisting users to better understand datasets with the Mixed-IE methods, based on which LLMs can be used to mine the topic information and construct domain models in the form of KGs classes without instantiating. The second stage adopts the Hybrid-RAG method to extract values under the guidance of KGs classes and interpret them into tables with specialized fields according to downstream needs using dynamic prompts.

### 3.2 Mixed-IE Assisted KGs Generation

As illustrated in Figure 3 (a), ① represents regulations and seed knowledge from human and ② represents the relevant inner knowledge of LLMs from pre-training, based on which ③ and ④ pre-processes the dataset such as section segmentation, tokenization, position tagging, named entity recognition (NER), and feature distribution statistics as well as filtering, to obtain lists of high term frequency (TF) and document frequency (DF). ⑤ con-

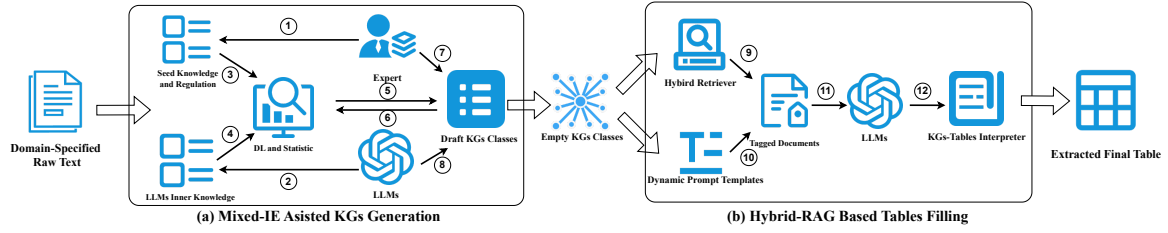


Figure 3: Overview of two-stages pipeline of TKGT.

209 structs domain models in the form of KGs classes  
 210 with the joint efforts of both human expert (7)  
 211 and LLMs (8), who also check the quality of KGs and  
 212 iterate it (6) to get final KGs classes. Here follows  
 213 the details of regulations, statistics, and DL, especially  
 214 LLMs methods, separately.

### 215 3.2.1 Regulations

216 Regulations refer to the structure, format, and logic,  
 217 which help to decompose complex texts into multiple  
 218 independent parts, reducing overall complexity. **Firstly**,  
 219 for general writing sense, writers produce texts logically,  
 220 such as the *What-Why-How Principle*, which is the inner  
 221 structure meaning different parts undertake different  
 222 functions with different information. **Secondly**, complex  
 223 texts usually adopt explicit structures like hierarchical  
 224 sections of academic papers to show inner logic clearly  
 225 to readers. **Finally**, shared elements are usually fixed  
 226 in the same positions, such as titles, author names,  
 227 and dates in certain lines. For instance, CPL judgment  
 228 documents contain the logic of legal trial and usually  
 229 adopt ordered positional words to present them more  
 230 clearly, as shown in Appendix B. By decomposing based  
 231 on regulations, the difficulties of subsequent work can  
 232 be greatly reduced. Thus, if users want to retrieve  
 233 identity information, the best choice is to perform  
 234 small-scale retrieval in the corresponding section.  
 235  
 236

### 237 3.2.2 Statistics

238 Purposes of statistics are ensuring the completeness  
 239 of IE to minimize losses of key words and exploring  
 240 topic and structure information. With mature NLP  
 241 toolkits and specified filtering, TF and DF reflect both  
 242 target information of a domain dataset. As shown in  
 243 Figure 1, after calculating the semantic similarity of  
 244 words and table fields, documents with the potential  
 245 for tabulation (Rotowire, E2E, and CPL) will exhibit  
 246 a **V-shape** pattern. By manually checking frequency  
 247 lists, it can be found that the first one percent of  
 248 the front parts of lists contain almost all keywords,  
 249 while the bottom part of

250 **V-shape** contains structure words dissimilar with  
 251 fields. Through statistics, users can quickly extract  
 252 keywords from large text sets and serve for LLMs  
 253 and human experts, greatly reducing the difficulty  
 254 of constructing KGs classes with completeness.

### 255 3.2.3 LLMs and KGs

256 An important trend of text-to-table is to break  
 257 down the original end-to-end paradigm into multiple  
 258 stages like (Deng et al., 2024) using triplets as  
 259 middleware. Compared to the topic-ignoring crawling  
 260 paradigm of triples, KGs can better model entities  
 261 and events, logically organize different roles and  
 262 adapt to downstream tabular needs. TKGT statistics  
 263 overall datasets to obtain relevant KGs classes, which  
 264 logically conducts retrieving values of certain objects'  
 265 fields in the second stage. This not only conforms to  
 266 more interpretable human methodology but is also  
 267 more accurate and complete. However, considering  
 268 that KGs generation itself is a difficult task and  
 269 existing research results only demonstrate the  
 270 possibility of using LLM to assist human experts in  
 271 generation (Meyer et al., 2023; Kommineni et al.,  
 272 2024), we simplify it as a *slack classes mining task*  
 273 with aims of reducing human expert participation.  
 274 That is, we do not instantiate KGs and only abstract  
 275 them as a set of classes with two types of *role entity*  
 276 and *relation/action* as shown in Appendix C.  
 277

### 278 3.3 Hybrid-RAG Based Table Filling

279 As illustrated in Figure 3 (b), (9) and (10) use  
 280 KGs classes from the first stage to dynamically rewrite  
 281 prompt templates and guide the hybrid retriever  
 282 respectively, combining with documents tagged in  
 283 the first stage to avoid unnecessary queries as LLMs  
 284 inputs (11). With inputs containing a set of retrieved  
 285 original texts as evidence and prompts, LLMs can  
 286 get certain values of the KGs classes (12) and  
 287 transform them to table form through the KGs-table  
 288 interpreter.



289 **3.3.1 Structure-Aware Hybrid-RAG**  
 290 We create an algorithm for scheduling the RAG  
 291 process with KGs, which is easy to understand and  
 292 adapt to other variants.

---

**Algorithm 1** KG Object Label Filling Algorithm

---

```

1: Initialize an empty KG object
2: while the KG object contains empty labels do
3:   if no entity in KG has filled labels then
4:     Select the entity with highest centrality
5:   else
6:     Calculate the ratio  $\frac{\text{Count}(\text{Label}|\text{Unfilled})}{\text{Count}(\text{Label})}$ 
       for each entity
7:     Select the entity with the highest ratio of
       unfilled labels
8:   end if
9:   if the selected entity's name label is not filled
       then
10:    Search and extract the entity name
11:   else
12:    Randomly select one unfilled label
13:    Search and extract information for the un-
       filled label
14:   end if
15:   if the information is found then
16:    Fill the searched information to the label
17:   else
18:    Fill 'Bad Information' to the label
19:   end if
20: end while

```

---

293 **3.3.2 Rewriting Prompt Dynamically**

294 We also utilize our KG design for query rewrit-  
 295 ing and summarizing relevant information before  
 296 passing them into the IE prompt. For query rewrit-  
 297 ing, we describe the relations between the "to-be-  
 298 extracted entity" and the label values of its adjacent  
 299 entities in the prompt. An example prompt is pro-  
 300 vided, asking the query rewriting model to generate  
 301 a search query for retrieving relevant information.  
 302 For information summarization, we describe the  
 303 same relations between the "to-be-extracted entity"  
 304 and the label values of its adjacent entities in the  
 305 prompt, asking the summarization model to retain  
 306 information that might be useful for answering the  
 307 user's question as shown in Appendix D.

308 **4 Experiments**

309 This section introduces the experimental setup and  
 310 results of TKGT's two stages respectively.

311 **4.1 Setup**

312 **Datasets.** As shown in Table 1, experiments use  
 313 datasets of Rotowire and E2E with table structure  
 314 processed by (Wu et al., 2021) and the CPL dataset  
 315 whose details are at Section 2 for more complex  
 316 challenges.

317 **Baselines and Models.** Considering the exten-  
 318 sive exploration of instruction following for various  
 319 LLMs (Ni et al., 2023; Deng et al., 2024), we pick  
 320 several popular LLMs as processors and focus on  
 321 the performance of TKGT on different datasets. Ta-  
 322 ble 3 shows baselines and models used. (1) *For*  
 323 *first stage*, we choose LLaMA3-70B<sup>2</sup> to test the  
 324 ability of KGs classes generation, comparing it  
 325 with two naive solutions: pure LLM with naive  
 326 prompt, and LLM with the same prompt template  
 327 of TKGT's using In-Context-Learning (ICL) and  
 328 Chain-of-Thought (CoT) but without statistical re-  
 329 sults. (2) *For the second stage*, for the demands of  
 330 deploying LLMs on consumer-grade GPUs in many  
 331 social science scenarios, we choose ChatGLM3-  
 332 6B<sup>3</sup> to test the ability of table extraction. We also  
 333 fine-tune it with LoRA (Hu et al., 2021) and com-  
 334 pare it with mainstream and SOTA commercial  
 335 LLM of GPT series<sup>4</sup>.

336 **Metrics.** (1) *For the first stage*, we develop an  
 337 evaluation method for the quality of KGs gener-  
 338 ation aiming at using LLMs to assist humans in  
 339 constructing domain KGs. We also recruit a group  
 340 of graduate students with knowledge in law and  
 341 computer science as referees. For the target dataset,  
 342 a set of fields is predefined by humans, and weights  
 343 are assigned to each field on average or based on  
 344 importance, which sum to 1. By checking the gen-  
 345 erated fields one by one with the target fields, we  
 346 can accumulate scores according to the rules in  
 347 Table 2, whose core principle is whether humans  
 348 can be inspired naturally by the fields generated by  
 349 LLMs. (2) *For the second stage*, metric follows the  
 350 F1 score at three levels defined in (Wu et al., 2021)

351 **4.2 Results of TKGT's First Stage**

352 Since TKGT's first stage is semi-automatic, results  
 353 can be iteratively improved by feedback from hu-  
 354 man and LLMs, making it difficult to reproduce.  
 355 Therefore, we only present results of first iteration,  
 356 in which TKGT provides predefined few-shot  
 357 templates and Mix-IE results, guiding LLMs to

<sup>2</sup><https://github.com/meta-llama/llama3>

<sup>3</sup><https://github.com/THUDM/ChatGLM3>

<sup>4</sup><https://openai.com/index/gpt-4-research/>

Matching Degree	Relationship of G&T Fields	Scoring Rules
Totally Match	Match in form or semantics	Obtain the total score of target field only once.
Including	Be a neighboring parent concept	Obtain 75% of the sum of all target fields.
Included	Be a neighboring sub concept	If parent concept is separable, obtain the field score divided by the number of categories each; If not, gain 25%.
Not Match	Completely different	No score.

Table 2: Metrics for the quality of KGs generated by TKGt’s first stage, in which *Relationship of G&T Fields* means the best-matching pair of one generated field and one target field. *Neighboring* refers to the ability to naturally infer parent/child concepts from subsequent textual information.

Stage	Method	Detail
First Stage	Zero-shot	LLaMA3-70B
	Few-shot	LLaMA3-70B & Prompt Template
	<b>TKGT-Stage-1</b>	LLaMA3-70B & Prompt Template & Statistics
Second Stage	Commercial LLM	GPT-3.5-turbo
	SOTA Commercial LLM	GPT-4-turbo
	Open-Source LLM	ChatGLM3-6B
	<b>TKGT-Stage-2</b>	ChatGLM3-6B & LoRA Tuning & RAG & KGs

Table 3: Experiment baselines of TKGt and details. LLaMA3-70B is one of the largest and most powerful open-source LLMs. ChatGLM3-6B is a popular medium-sized open-source LLM. GPT series contain the most popular commercial LLMs.

Subset	Model	The first column F1			Table header F1			Data cell F1			Error
		Exact	Chrf	BERT	Exact	Chrf	BERT	Exact	Chrf	BERT	
Team	Sent-level RE	85.28	87.12	93.65	85.54	87.99	87.53	77.17	79.10	87.48	0.00
	Doc-level RE	84.90	86.73	93.44	85.46	88.09	87.99	75.66	77.89	87.82	0.00
	Seq2Seq	94.71	94.93	97.35	86.07	89.18	88.90	82.97	84.43	90.62	0.49
	Seq2Seq-c	94.97	95.20	97.51	86.02	89.24	89.05	83.36	84.76	90.80	0.00
	Seq2Seq&set	<b>96.80</b>	<b>97.10</b>	<b>98.45</b>	86.00	89.48	93.11	84.33	85.68	<b>91.30</b>	0.00
	T-(No RAG)-T*	72.38	72.84	73.41	<b>100.0</b>	<b>100.0</b>	<b>100.0</b>	64.42	65.53	66.84	0.00
	T-KG-T*	91.44	91.83	93.26	<b>100.0</b>	<b>100.0</b>	<b>100.0</b>	<b>85.03</b>	<b>87.58</b>	91.21	0.00
Player	Sent-level RE	89.05	93.00	90.98	86.36	89.38	93.07	79.59	83.42	85.35	0.00
	Doc-level RE	89.26	93.28	91.19	87.35	90.22	97.30	80.76	84.64	86.50	0.00
	Seq2Seq	92.16	93.89	93.60	87.82	91.28	94.44	81.96	84.19	88.66	7.40
	Seq2Seq-c	92.31	94.00	93.71	87.78	91.26	94.41	82.53	84.74	88.97	0.00
	Seq2Seq&set	92.83	94.48	<b>96.43</b>	88.02	91.60	95.08	83.51	85.75	<b>90.93</b>	0.00
	T-(No RAG)-T*	67.51	69.29	69.22	<b>100.0</b>	<b>100.0</b>	<b>100.0</b>	64.27	66.25	66.94	0.00
	T-KG-T*	<b>93.05</b>	<b>94.59</b>	95.18	<b>100.0</b>	<b>100.0</b>	<b>100.0</b>	<b>88.26</b>	<b>90.18</b>	90.39	0.00

Table 4: Results of baselines, pure LLMs prompts, and our TKGt model on Rotowire. We show the F1 score based on exact match (Exact), chrf score (Chrf), and BERTScore (BERT) respectively. GLM3-6B refers to the pre-trained ChatGLM3-6B model without any finetuning. \* refers to the finetuned IE model tuned on the respective IE finetuning dataset we created based on the corresponding dataset.

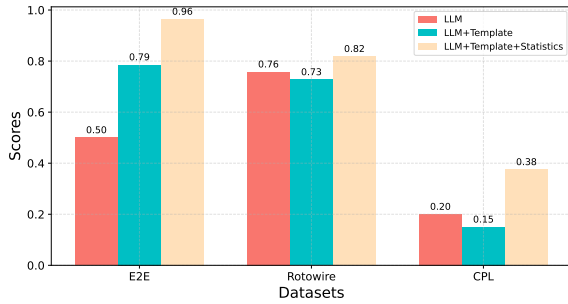


Figure 4: Results of TKGT’s first stage.

generate KGs classes for three datasets. Besides, we add ablation elements to it, removing Mix-IE results and few-shot templates. We run 10 times each and submit outputs to a group of human judge with metrics to obtain the best result.

As shown in Figure 4, TKGT achieves the best performance on all datasets, which proves that our method can extract more complete domain models. We observe that scores decrease as the complexity (numbers of fields and structures) of the dataset increases, and TKGT get 0.96 and 0.82 on E2E and Rotowire respectively, indicating that TKGT can generate almost complete structures for traditional datasets. Furthermore, as for Rotowire and CPL, the method with Few-shot templates but without results from Mix-IE gets even lower scores than pure LLM, which means templates without top keywords hinder LLM’s ability to exert its inner knowledge and proves the importance of Mixed-IE. Finally, TKGT performs poorly without iteration, proposing further research challenges.

### 4.3 Results of TKGT’s Second Stage

As shown in Table 4 and the first half of Table 5, our TKGT pipeline achieves near SOTA performance with minimal dataset-specific engineering for the Rotowire dataset. Our KG-based design avoids generating incorrect table headers and mismatched table shapes, achieving perfect scores in table header F1 and Error compared to previous methods. The relatively low F1 scores for the first column (Team name) extraction are due to the model’s difficulty in identifying ‘home team’ and ‘visiting team’ from their positions in the text. We achieve SOTA performance on nearly all metrics. We did not use any RAG technique in the ablation experiment because both the E2E and Rotowire data are short and lack a specific writing style, where RAG might cause more information loss than precision gain. Comparing ‘T-(No RAG)-T’ and ‘T-KG-T’ shows

the benefits of our KG-guided query, query-rewrite, and summarizing pipeline.

We compare TKGT with larger commercial LLMs on CPL dataset. Despite the base model’s limitations, T-KG-T performs comparably to more advanced models like GPT-4-Turbo using naive RAG, showcasing the effectiveness of our KG-guided methods. Fine-tuning the IE model is crucial for ‘Text-to-Table’ tasks, initially ensuring adherence to the output format, then distinguishing between valid and invalid information cases, and finally accurately extracting valid information. Our KG-guided query, query-rewrite, and summarizing pipeline enhance the model’s ability to deliver accurate information by reducing unnecessary context and adding relevant information, ultimately achieving state-of-the-art performance.

## 5 Related Work

### 5.1 Text-to-Table Works in Social Science

Text-to-table works in social science are more engineering-oriented, meeting needs of text-as-data (Ash and Hansen, 2023), which involves four core empirical tasks: ① measure document similarity (Cagé et al., 2020; Kelly et al., 2021); ② concept detection (Shapiro et al., 2022; Angelico et al., 2022); ③ how concepts are related (Thorsrud, 2020; Ash et al., 2024); ④ associate text to metadata (Ke et al., 2019)). Traditional methods of structuring is manual coding, such as Chang et al. (2021) spending years coding 170 dimensions of property law in 128 jurisdictions to draw the legal family. With the development of NLP, structuring tasks become semi-automated or even fully-automated (Grimmer et al., 2022). Luo et al. (2017) propose an Transformer-based method to simultaneously model charge prediction and relevant article extraction tasks. Mentzingen et al. (2024) first develop a two-stage cascade classifier model that predicts regulatory decisions, based on textual features extracted from the original documents by ML and proceedings’ metadata.

### 5.2 Text-to-Table Works in Computer Science

The research paradigm of text-to-table officially originated from Wu et al. (2021), which uses datasets from table-to-text and an end-to-end sequence generation mode based on the BART model. All rows are generated at once, and the results are controlled using table constraints and column embedding. Li et al. (2023b) improves it by point-

Dataset	Model	The first column F1			Data cell F1			Error
		Exact	Chrf	BERT	Exact	Chrf	BERT	
E2E	NER	85.28	87.12	93.65	85.54	87.99	87.53	0.00
	Seq2Seq	84.90	86.73	93.44	85.46	88.09	87.99	0.49
	Seq2Seq-c	94.71	94.93	97.35	86.07	89.18	88.90	0.00
	Seq2Seq&set	94.97	95.20	<b>97.51</b>	86.02	89.24	89.05	0.00
	T-(No RAG)-T (GLM3-6B*)	74.34	76.07	78.92	71.39	73.15	74.07	0.00
	T-KG-T (GLM3-6B*)	<b>95.14</b>	<b>95.87</b>	96.12	<b>92.17</b>	<b>93.79</b>	<b>92.83</b>	0.00
CPL	T-(Naive RAG)-T (GPT3.5)	84.26	82.67	66.28	79.43	67.73	55.01	0.00
	T-(Naive RAG)-T (GPT4)	<b>93.41</b>	<b>91.73</b>	80.52	90.27	<b>88.62</b>	78.70	0.00
	T-(No RAG)-T (GLM3-6B)	9.97	0.89	0.95	4.60	1.95	0.86	0.00
	T-(Naive RAG)-T (GLM3-6B)	12.25	11.31	11.98	8.87	2.19	1.98	0.00
	T-KG-T (GLM3-6B*)	91.33	88.79	<b>82.68</b>	<b>90.79</b>	87.58	<b>82.45</b>	0.00

Table 5: Results of baselines, pure LLMs prompts, and our TKGT model on CPL. F1 scores are same as Table 4. GLM3-6B refers to the pretrained ChatGLM3-6B model without any finetuning. GLM3-6B\* refers to the finetuned IE model tuned on the respective IE finetuning dataset we created based on the corresponding dataset.

ing out the order-insensitive property of rows and adopted a fast method of generating all rows in parallel after generating the header. Sundar et al. (2024) abandons the end-to-end paradigm and adopts a two-stage approach of generating table frameworks and content separately and switches to use conditional Q&A for IE. Deng et al. (2024) further innovates by proposing a new benchmark and uses LLMs prompt engineering to extract triples from the original text and merge them into tables.

### 5.3 LLMs Prompt and Knowledge Graphs

Prompt originated from the GPT-3 series (Brown et al., 2020), whose works focus on engineering experience and practice, such as the various prompt techniques listed in (Liu et al., 2023). In addition, Sahoo et al. (2024) combines prompt and fine-tuning to explain the essence of instruction following. Wang et al. (2023) further explores the potential of fine-tuned LLMs in IE. As for KGs, recent works explore how to use LLMs to empower the construction of KGs. Meyer et al. (2023) first explores the potential of LLMs to generate KGs in multiple engineering fields, Ni et al. (2023) elucidates the complementary relationship between LLMs and KGs, and Kommineni et al. (2024) proposes a semi-automatic pipeline method using LLMs to assist human experts in generating KGs as the latest research.

### 5.4 IE and RAG

Retrieval-Augmented Generation (RAG) aims to enhance the factual accuracy of Large Language

Models (LLMs) by incorporating relevant textual information, thereby expanding the knowledge base of the training data and reducing hallucination problems (Gao et al., 2024). Khattab et al. (2023) was one of the pioneering works utilizing the in-context learning ability of LLMs to perform knowledge-intensive information retrieval tasks in the form of question-answering. Subsequent research has made various improvements to RAG, such as introducing new data structures for retrieval data (Luo et al., 2023; He et al., 2024) and developing more efficient retrieval pipelines. These advancements include hybrid retrieval methods (Gao et al., 2022), fine-tuning embeddings (Shi et al., 2023), reranking (Yu et al., 2023), and iterative retrieval processes (Cheng et al., 2023).

## 6 Conclusion

We first review the research field of text-to-table, point out the shortcomings of existing datasets with statistical methods, and redefine the core requirements of this task more comprehensively. Secondly, we propose a social science dataset CPL from real-world structuring requirements, which presents new challenges to the field due to its complexity and semi-structured nature. In addition, to address the shortcomings of existing text-to-table methods that overlook topic and structural information, we propose a two-stage pipeline called TKGT using KGs classes as middleware and demonstrate its SOTA performance through experiments.



## 507 Limitations

508 Although the TKG T pipeline we propose covers  
509 the entire process of text-to-table task, it cannot be  
510 fully automated in the first stage. On the one hand,  
511 this is limited by the current capabilities of LLMs;  
512 On the other hand, academic level complex text  
513 extraction tasks are extremely challenging even  
514 for untrained humans. One possible solution is  
515 to build the first stage as a more comprehensive  
516 and powerful agent, and explore a more powerful  
517 initialization framework that balances universality  
518 and practicality. This is also one of our future tasks.

## 519 Ethics Statement

520 This work does not adopt AI assistants. The  
521 four datasets we use are entirely from the MIT  
522 license open-source pre-processing results of previ-  
523 ous work (Wu et al., 2021), while the CPL dataset is  
524 sourced from the official judgment documents pub-  
525 licly available on the CJO, which complies with the  
526 requirement of transparency in court rulings. The  
527 CPL dataset involves real person names and other  
528 information. In order to further ensure privacy  
529 and ensure the accuracy of named entity recog-  
530 nition during data pre-processing, we randomly  
531 replaced the person names using existing named  
532 entity recognition techniques (He and Choi, 2021).  
533 In addition, all experiments in this work followed  
534 the expected purpose of their research. Therefore,  
535 to the best of the author’s knowledge, we believe  
536 that this work will not bring any additional risks.

## 537 References

538 Cristina Angelico, Juri Marcucci, Marcello Miccoli,  
539 and Filippo Quarta. 2022. [Can we measure inflation  
540 expectations using twitter?](#) *Journal of Econometrics*,  
541 228(2):259–277.

542 Elliott Ash, Germain Gauthier, and Philine Widmer.  
543 2024. [Relatio: Text semantics capture political  
544 and economic narratives.](#) *Political Analysis*,  
545 32(1):115–132.

546 Elliott Ash and Stephen Hansen. 2023. [Text algorithms  
547 in economics.](#) *Annual Review of Economics*, 15(Vol-  
548 ume 15, 2023):659–688.

549 Kevin D. Ashley. 2017. *Artificial Intelligence and Legal  
550 Analytics: New Tools for Law Practice in the Digital  
551 Age*. Cambridge University Press, Cambridge.

552 Junwei Bao, Duyu Tang, Nan Duan, Zhao Yan, Yuanhua  
553 Lv, Ming Zhou, and Tiejun Zhao. 2018. Table-to-  
554 text: Describing table region with natural language.

*In Proceedings of the AAAI conference on artificial  
intelligence*, volume 32. 555  
556

Tom Brown, Benjamin Mann, Nick Ryder, Melanie  
Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind  
Neelakantan, Pranav Shyam, Girish Sastry, Amanda  
Askell, et al. 2020. Language models are few-shot  
learners. *Advances in neural information processing  
systems*, 33:1877–1901. 557  
558  
559  
560  
561  
562

Julia Cagé, Nicolas Hervé, and Marie-Luce Viaud. 2020.  
[The production of information in an online world.](#)  
*The Review of Economic Studies*, 87(5):2126–2164. 563  
564  
565

Yun-chien Chang, Nuno Garoupa, and Martin T Wells.  
2021. [Drawing the legal family tree: An empirical  
comparative study of 170 dimensions of property  
law in 129 jurisdictions.](#) *Journal of Legal Analysis*,  
13(1):231–282. 566  
567  
568  
569  
570

Xin Cheng, Di Luo, Xiuying Chen, Lemao Liu,  
Dongyan Zhao, and Rui Yan. 2023. [Lift yourself  
up: Retrieval-augmented text generation with self  
memory.](#) *Preprint*, arXiv:2305.02437. 571  
572  
573  
574

Zheye Deng, Chunkit Chan, Weiqi Wang, Yuxi Sun,  
Wei Fan, Tianshi Zheng, Yauwai Yim, and Yangqiu  
Song. 2024. [Text-tuple-table: Towards information  
integration in text-to-table generation via global tuple  
extraction.](#) *arXiv preprint arXiv:2404.14215*. 575  
576  
577  
578  
579

Federal Judicial Center FJC. 2020. *Judicial Writing  
Manual: A Pocket Guide for Judges*. Lulu.com. 580  
581

Luyu Gao, Xueguang Ma, Jimmy Lin, and Jamie Callan.  
2022. [Precise zero-shot dense retrieval without rele-  
vance labels.](#) *Preprint*, arXiv:2212.10496. 582  
583  
584

Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia,  
Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, Meng Wang,  
and Haofen Wang. 2024. [Retrieval-augmented gener-  
ation for large language models: A survey.](#) *Preprint*,  
arXiv:2312.10997. 585  
586  
587  
588  
589

Matthew Gentzkow, Bryan Kelly, and Matt Taddy.  
2019. [Text as data.](#) *Journal of Economic Literature*,  
57(3):535–574. 590  
591  
592

Justin Grimmer, Margaret E. Roberts, and Brandon M.  
Stewart. 2022. *Text as Data: A New Framework for  
Machine Learning and the Social Sciences*. Princeton  
University Press, Princeton. 593  
594  
595  
596

Jonathan Grossman and Ami Pedahzur. 2020. [Political  
science and big data: Structured data, unstructured  
data, and how to use them.](#) *Political Science Quar-  
terly*, 135(2):225–257. 597  
598  
599  
600

Han He and Jinho D Choi. 2021. The stem cell hypoth-  
esis: Dilemma behind multi-task learning with trans-  
former encoders. *arXiv preprint arXiv:2109.06939*. 601  
602  
603

Xiaoxin He, Yijun Tian, Yifei Sun, Nitesh V. Chawla,  
Thomas Laurent, Yann LeCun, Xavier Bresson, and  
Bryan Hooi. 2024. [G-retriever: Retrieval-augmented  
generation for textual graph understanding and ques-  
tion answering.](#) *Preprint*, arXiv:2402.07630. 604  
605  
606  
607  
608

609	Tony Hey, Stewart Tansley, and Kristin Tolle. 2009. <i>The Fourth Paradigm: Data-Intensive Scientific Discovery</i> . Microsoft Research, Redmond, Washington.	In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 5358–5370.	663
610			664
611			
612	Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. <i>arXiv preprint arXiv:2106.09685</i> .	Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2023. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. <i>ACM Computing Surveys</i> , 55(9):1–35.	665
613			666
614			667
615			668
616			669
617	Parag Jain, Andreea Marzoca, and Francesco Piccinno. 2024. Structsum generation for faster text comprehension. <i>arXiv preprint arXiv:2401.06837</i> .	Yaojie Lu, Qing Liu, Dai Dai, Xinyan Xiao, Hongyu Lin, Xianpei Han, Le Sun, and Hua Wu. 2022. Unified structure generation for universal information extraction. <i>arXiv preprint arXiv:2203.12277</i> .	670
618			671
619			672
620	Zheng Tracy Ke, Bryan T. Kelly, and Dacheng Xiu. 2019. Predicting returns with text data. (26186). DOI: 10.3386/w26186.	Bingfeng Luo, Yansong Feng, Jianbo Xu, Xiang Zhang, and Dongyan Zhao. 2017. Learning to predict charges for criminal cases with legal basis. In <i>Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing</i> , page 2727–2736, Copenhagen, Denmark. Association for Computational Linguistics.	674
621			675
622			676
623	Bryan Kelly, Dimitris Papanikolaou, Amit Seru, and Matt Taddy. 2021. Measuring technological innovation over the long run. <i>American Economic Review: Insights</i> , 3(3):303–320.	Ziyang Luo, Can Xu, Pu Zhao, Xiubo Geng, Chongyang Tao, Jing Ma, Qingwei Lin, and Daxin Jiang. 2023. Augmented large language models with parametric knowledge guiding. <i>Preprint</i> , arXiv:2305.04757.	677
624			678
625			679
626			680
627	Omar Khattab, Keshav Santhanam, Xiang Lisa Li, David Hall, Percy Liang, Christopher Potts, and Matei Zaharia. 2023. Demonstrate-search-predict: Composing retrieval and language models for knowledge-intensive nlp. <i>Preprint</i> , arXiv:2212.14024.	Hugo Mentzingen, Nuno Antonio, and Victor Lobo. 2024. Joining metadata and textual features to advise administrative courts decisions: a cascading classifier approach. <i>Artificial Intelligence and Law</i> , 32(1):201–230.	681
628			682
629			683
630			684
631			
632			
633	Gary King. 2014. Restructuring the social sciences: Reflections from harvard’s institute for quantitative social science. <i>PS, Political Science Politics</i> , 47(1):165–172.	LP Meyer, C Stadler, J Frey, N Radtke, K Junghanns, R Meissner, G Dziwis, K Bulert, and M Martin. 2023. Llm-assisted knowledge graph engineering: experiments with chatgpt (2023). In <i>conference proceedings of AI-Tomorrow-23</i> , volume 29, pages 6–2023.	685
634			686
635			687
636			688
637	Diederik P. Kingma and Jimmy Ba. 2017. Adam: A method for stochastic optimization. <i>Preprint</i> , arXiv:1412.6980.	Jean-Baptiste Michel, Yuan Kui Shen, Aviva Presser Aiden, Adrian Veres, Matthew K. Gray, THE GOOGLE BOOKS TEAM, Joseph P. Pickett, Dale Hoiberg, Dan Clancy, Peter Norvig, Jon Orwant, Steven Pinker, Martin A. Nowak, and Erez Lieberman Aiden. 2011. Quantitative analysis of culture using millions of digitized books. <i>Science</i> , 331(6014):176–182.	689
638			690
639			691
640	Vamsi Krishna Kommineni, Birgitta König-Ries, and Sheeba Samuel. 2024. From human experts to machines: An llm supported approach to ontology and knowledge graph construction. <i>arXiv preprint arXiv:2403.08345</i> .	Xuanfan Ni, Piji Li, and Huayang Li. 2023. Unified text structuralization with instruction-tuned language models. <i>arXiv preprint arXiv:2303.14956</i> .	692
641			693
642			694
643			
644			
645	David Lazer, Alex Pentland, Lada Adamic, Sinan Aral, Albert-László Barabási, Devon Brewer, Nicholas Christakis, Noshir Contractor, James Fowler, Myron Gutmann, Tony Jebara, Gary King, Michael Macy, Deb Roy, and Marshall Van Alstyne. 2009. Computational social science. <i>Science</i> , 323(5915):721–723.	Jekaterina Novikova, Ondřej Dušek, and Verena Rieser. 2017. The e2e dataset: New challenges for end-to-end generation. <i>arXiv preprint arXiv:1706.09254</i> .	695
646			696
647			697
648			698
649			699
650			700
651	Rémi Lebret, David Grangier, and Michael Auli. 2016. Neural text generation from structured data with application to the biography domain. <i>arXiv preprint arXiv:1603.07771</i> .	Shirui Pan, Linhao Luo, Yufei Wang, Chen Chen, Jipu Wang, and Xindong Wu. 2024. Unifying large language models and knowledge graphs: A roadmap. <i>IEEE Transactions on Knowledge and Data Engineering</i> .	701
652			702
653			703
654			704
655	Peng Li, Yeye He, Dror Yashar, Weiwei Cui, Song Ge, Haidong Zhang, Danielle Rifinski Fainman, Dongmei Zhang, and Surajit Chaudhuri. 2023a. Table-gpt: Table-tuned gpt for diverse table tasks. <i>arXiv preprint arXiv:2310.09263</i> .	Pranab Sahoo, Ayush Kumar Singh, Sriparna Saha, Vinija Jain, Samrat Mondal, and Aman Chadha. 2024. A systematic survey of prompt engineering in large language models: Techniques and applications. <i>arXiv preprint arXiv:2402.07927</i> .	705
656			706
657			707
658			708
659			709
660	Tong Li, Zhihao Wang, Liangying Shao, Xuling Zheng, Xiaoli Wang, and Jinsong Su. 2023b. A sequence-to-sequence&set model for text-to-table generation.		710
661			711
662			712
			713
			714
			715
			716
			717
			718

719 Adam Hale Shapiro, Moritz Sudhof, and Daniel J. Wil-  
720 son. 2022. *Measuring news sentiment*. *Journal of*  
721 *Econometrics*, 228(2):221–243.

722 Weijia Shi, Sewon Min, Michihiro Yasunaga, Min-  
723 joon Seo, Rich James, Mike Lewis, Luke Zettle-  
724 moyer, and Wen tau Yih. 2023. *Replug: Retrieval-*  
725 *augmented black-box language models*. *Preprint*,  
726 arXiv:2301.12652.

727 Yuan Sui, Mengyu Zhou, Mingjie Zhou, Shi Han, and  
728 Dongmei Zhang. 2024. Table meets llm: Can large  
729 language models understand structured table data?  
730 a benchmark and empirical study. In *Proceedings*  
731 *of the 17th ACM International Conference on Web*  
732 *Search and Data Mining*, pages 645–654.

733 Anirudh Sundar, Christopher Richardson, and Larry  
734 Heck. 2024. gtbls: Generating tables from text  
735 by conditional question answering. *arXiv preprint*  
736 *arXiv:2403.14457*.

737 Leif Anders Thorsrud. 2020. *Words are the new num-*  
738 *bers: A newsy coincident index of the business*  
739 *cycle*. *Journal of Business Economic Statistics*,  
740 38(2):393–409.

741 Javier Vidal-García, Marta Vidal, and Rafael Hernández  
742 Barros. 2019. *Computational Business Intelligence,*  
743 *Big Data, and Their Role in Business Decisions in*  
744 *the Age of the Internet of Things*, page 1048–1067.  
745 IGI Global.

746 Xiao Wang, Weikang Zhou, Can Zu, Han Xia, Tianze  
747 Chen, Yuansen Zhang, Rui Zheng, Junjie Ye,  
748 Qi Zhang, Tao Gui, et al. 2023. Instructuie: multi-  
749 task instruction tuning for unified information extrac-  
750 tion. *arXiv preprint arXiv:2304.08085*.

751 Sam Wiseman, Stuart M Shieber, and Alexander M  
752 Rush. 2017. Challenges in data-to-document genera-  
753 tion. *arXiv preprint arXiv:1707.08052*.

754 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien  
755 Chaumond, Clement Delangue, Anthony Moi, Pier-  
756 ric Cistac, Tim Rault, Rémi Louf, Morgan Funtow-  
757 icz, Joe Davison, Sam Shleifer, Patrick von Platen,  
758 Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu,  
759 Teven Le Scao, Sylvain Gugger, Mariama Drame,  
760 Quentin Lhoest, and Alexander M. Rush. 2020. *Hug-*  
761 *gingface’s transformers: State-of-the-art natural lan-*  
762 *guage processing*. *Preprint*, arXiv:1910.03771.

763 Xueqing Wu, Jiacheng Zhang, and Hang Li. 2021. Text-  
764 to-table: A new way of information extraction. *arXiv*  
765 *preprint arXiv:2109.02707*.

766 Zichun Yu, Chenyan Xiong, Shi Yu, and Zhiyuan Liu.  
767 2023. *Augmentation-adapted retriever improves gen-*  
768 *eralization of language models as generic plug-in*.  
769 *Preprint*, arXiv:2305.17331.

770 Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan  
771 Ye, Zheyuan Luo, and Yongqiang Ma. 2024. *Llamafac-*  
772 *tory: Unified efficient fine-tuning of 100+ language*  
773 *models*. *Preprint*, arXiv:2403.13372.

## A Details of CPL Dataset

In order to study private lending in China, such as the changing patterns of lending behavior, the logic and efficiency of trial, and the policy effects of interest rate regulation, a real-world academic project obtains CPL judgements from the CJO and conducts manual structuring of these judgements. The main goal of this work is to extract the content of each judgment as comprehensively as possible into a structured format in a table.

The project carries out this work through the following steps. **Firstly**, design the format of the table. In different countries, the logic of trials and the writing of judgements are basically the same (FJC, 2020). The core logic of the court’s trial is to accurately grasp the claims and grounds of the litigants surrounding the same lending behavior facts, and the court makes its determination and judgment accordingly. And the CPL judgments have a consistent structure. Therefore, the project reassemble the content of the judgement into a  $(2 \times n) \times 5$  format, as shown in Figure 2. The 2 represents the two major dimensions: Basic Information of Court and Parties and Basic Lending Facts. The  $n$  represents the specific content under each dimension. The 5 represents the five main entities: court, borrower, lender, guarantor, and others. **Secondly**, set over 200 fields and corresponding value ranges by reading judgements and sorting out relevant legal norms. These fields basically cover the core elements of trial, such as the *Elemental Trial Guide*<sup>5</sup> and the *Model Texts of Written Civil Complaints and Statements of Defense*<sup>6</sup>, indicating that this work is thorough and scientific. The Excel table for manual data collection is constructed by professors and graduate students in law. **Thirdly**, complete text-to-table manually. The project recruit undergraduate students with a legal background and conduct a two-week training. The work is carried out in a one-by-one format, with one undergraduate student collecting and one graduate student reviewing.

This project recruited students and compensated them based on the work-study standards of their respective universities. It provided participants with the full text of instructions, including disclaimers

<sup>5</sup>Issued by The High People’s Court of Shandong Province, <http://ytzy.sdcourt.gov.cn/ytzy/yhfzyshj/zxht39/sfwj/6518994/index.html>

<sup>6</sup>Issued by the Supreme People’s Court, the Ministry of Justice, and the All China Lawyers Association, <https://pkulaw.com/ch1/1b4f90e3dcf35b36bdfb.html>



of any risks. The data collection protocol was approved by an ethics review board. The subjects included in CPL dataset are Chinese citizens, primarily from Shanghai, Zhejiang Province, and Anhui Province. We obtained authorization from the project leader to use the CPL dataset.

## B Structure of CPL Judgement

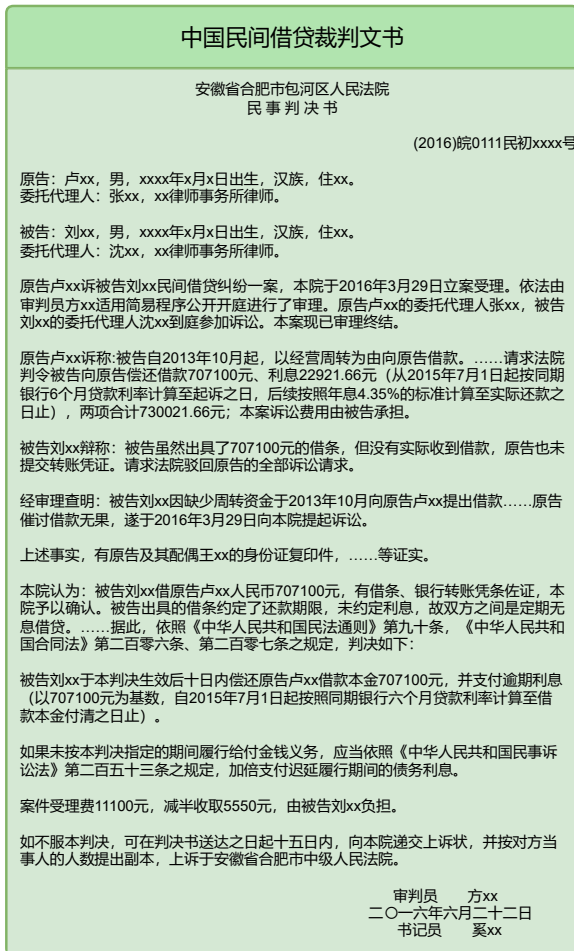


Figure 5: CPL Judgement Demo (Chinese Version).

Due to the issuance of *Specifications for Preparing Civil Judgments by the People's Courts*<sup>7</sup> and the *Style of Civil Litigation Documents*<sup>8</sup> by SPC, CPL judgments have a consistent structure (Figure 5 and Figure 6): ① Basic information of the court, such as the name of the court, the name of the judgment, and the case number; ② Parties and their basic information (e.g., name, address, role); ③ Procedural history; ④ Claims, facts, and grounds of the parties; ⑤ Evidence and facts identified by

<sup>7</sup><https://pkulaw.com/ch1/4c13be0c1802426abdfb.html?way=listView>

<sup>8</sup><https://www.court.gov.cn/susong.html>

the court; ⑥ Grounds, judicial basis, and main body of judgment; ⑦ Signatory information, such as the information of the trial personnel and the closed date.

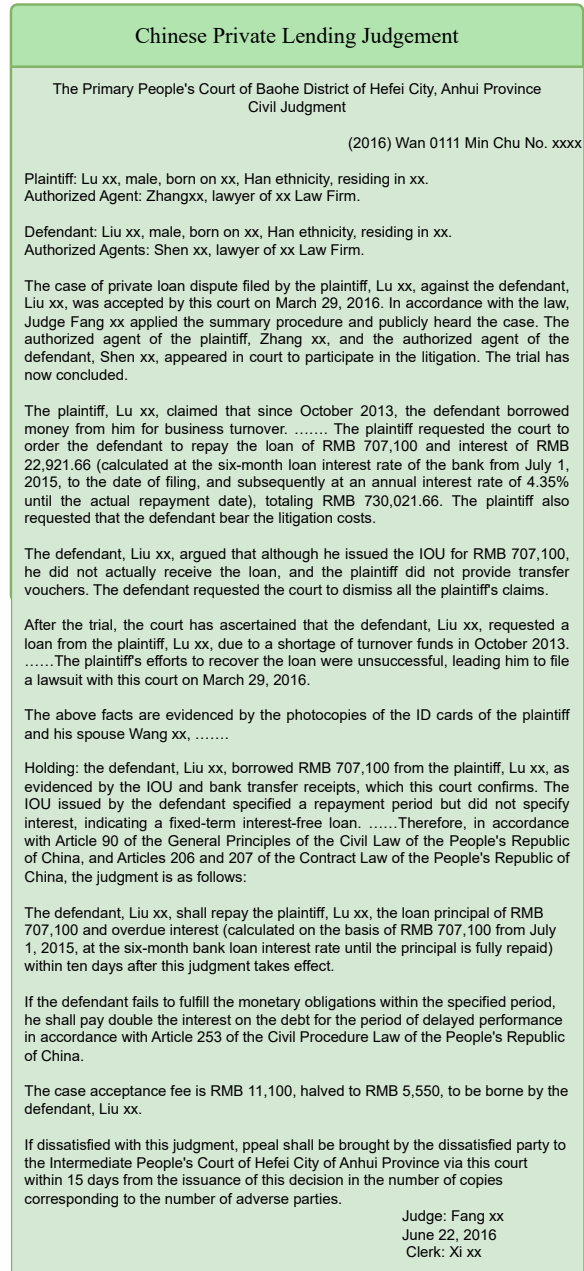


Figure 6: CPL Judgement Demo (English Version).

## C Details of TKGT's First Stage

**Slack classes.** To simplify KGs, we abstract it as two basic classes of *role entity classes* and *relation/action classes*. The former can represent any entity such as humans or objects, while the latter broadly represents relationships or behaviors that require multi-party participation.

**Toolkits.** We used existing NLP methods in TKGT.



849 For Chinese, we use Hanlp’s (He and Choi, 2021)  
 850 sentence splitter as well as its integrated tokenizer,  
 851 position tagger, and Chinese NER model. As  
 852 for English, we use nltk’s tokenizer and posi-  
 853 tion tagger. *As for stop Words*, we use Chinese  
 854 stop words from [https://blog.csdn.net/qq\\_33772192/article/details/91886847](https://blog.csdn.net/qq_33772192/article/details/91886847) and En-  
 855 glish stop words from spaCy<sup>9</sup>. *As for stop position*  
 856 *taggers*, due to the differences in the categories of  
 857 parts of speech between Chinese and English, we  
 858 choose positions to use based on the CTB tag set  
 859 for Chinese, while the positions to disable based  
 860 on the NLTK tag set for English as follows.  
 861

```
862 [
863     used_pos_zh = ["NR", "NN", "CD", "VV",
864                  "NT", "FW", "AD", "JJ" ],
865     stop_pos_en = ["CC", "DT", "EX", "IN",
866                  "MD", "PDT", "POS", "PRP",
867                  "RP", "SYM", "TO", "UH",
868                  "WDT" , "WP" ]
869 ]
```

## 870 D Prompt Example

### 871 D.1 Information Extraction Prompt

872 We design the prompt to contain 3 parts as the IE  
 873 task the model would complete would also follow  
 874 three key steps: First, the assistant checks if the pro-  
 875 vided paragraph contains the attribute values corre-  
 876 sponding to the role; if not, it responds with 'Bad  
 877 Information'. Second, if the paragraph contains  
 878 the relevant attribute values, the assistant extracts  
 879 and provides the value according to the specified  
 880 requirements. Third, the assistant responds to the  
 881 user’s question in the format of the provided in-  
 882 context examples. Each example outlines the role,  
 883 attribute, related context, value scope, question,  
 884 and answer, ensuring the assistant’s responses are  
 885 precise and consistent. The relatively low F1 scores  
 886 for the first column (Team name) extraction are due  
 887 to the model’s difficulty in identifying 'home team'  
 888 and 'visiting team' from their positions in the text.

<sup>9</sup><https://spacy.io/>

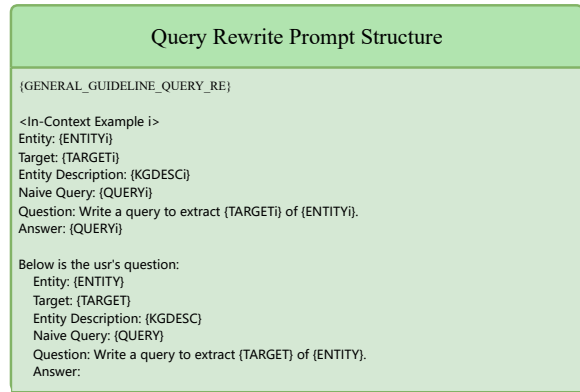


Figure 8: Structure of Query Rewrite Prompt.

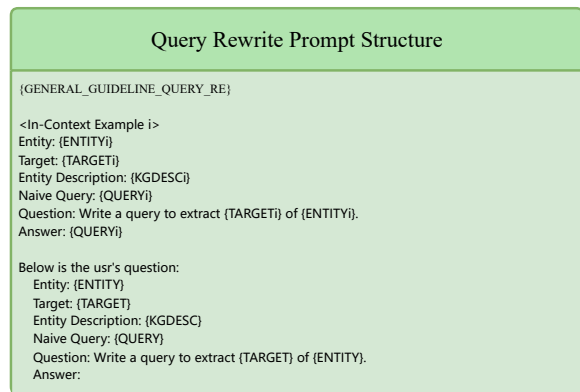


Figure 9: Structure of Information Summary Prompt.

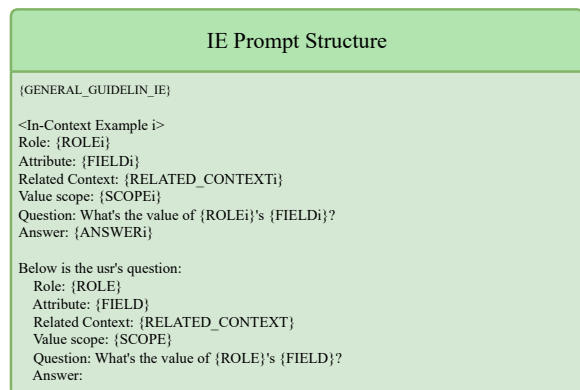


Figure 7: Structure of Information Retrieving Prompt.

## 889 E Fine-tuning Setting

### 890 E.1 Fine-tuning Parameter and Setting

891 We use the open-source library LLaMA-Factory  
 892 (Zheng et al., 2024) to fine-tune all models. LoRA  
 893 (Hu et al., 2021) is used as the fine-tuning. The  
 894 pre-trained weights are downloaded from the hug-  
 895 gingface library (Wolf et al., 2020). We load the  
 896 models with FP16 as the precision and optimize  
 897 them with an Adam optimizer (Kingma and Ba,

2017). The learning rate is set to  $1e-4$  with cosine decay and the batch size is 2 per device. The maximum length for the input and generated sentence concatenation is 2048. We warm up the model with 3,000 steps and evaluate the model every 500 steps. A linear scheduler is also used. The LoRA rank is set to 16, and the  $\alpha$  is set to 32.

## E.2 Fine-tuning Data Preparation

In this subsection, we detail the data collection process for fine-tuning the Information Extraction (IE) model. Our approach to constructing the fine-tuning dataset aligns with the structure of the 'TKGT' framework. The IE model is employed only at the IE stage following 'Query Generation', 'Query Rewrite', and 'Information Summarizing'. To ensure consistency between the fine-tuning data and inference stages, we utilize pre-trained models for query rewriting and information summarizing.

Fine-tuning the IE model is crucial for enhancing the performance of 'Text-To-Table' tasks. Initially, the model learns to adhere to the specified output format. Subsequently, it differentiates between cases containing valid information (Good Information Case) and those that do not (Bad Information Case). Finally, the model identifies and extracts valid information accurately.

The fine-tuning dataset is composed in the following format:

```
[
  {"instruction": <ie task id>,
   "input": <ie prompt>,
   "output": <ground truth>},
  ...
]
```

## F Computing Cost

### F.1 Cost of Stage 1 Inference

Although we can measure the coverage of zero-shot and few-shot performance of KG generation, constructing an accurate domain-specific KG for information extraction depends on human expert judgment, the complexity of the text data, and the granularity of the information designed to be extracted to form the outcome table. For the E2E and Rotowire datasets, we report that LLaMa3-70B is able to construct acceptable KG classes with a single prompt. However, for more complex datasets like CPL, it requires significantly more iterations and human expert involvement in constructing the KG.

### F.2 Cost of Stage 2 Inference

We can estimate the cost of stage 2 inference following the T-KG-T pipeline. For each document, suppose there are  $n$  variables in total and  $m$  variables are 'easy and obvious'<sup>10</sup> that can be easily extracted. For every variable that needs to go through the pipeline for extraction, it must undergo 'Query Rewrite', 'Information Summarization', 'Information Retrieving', and 'Information Extraction' processes, totaling 3 prompts and 1 retrieval. The algorithm ensures that each variable goes through the pipeline at most once.

Therefore, to extract the document, we would need a maximum of  $3 \times (n - m)$  model requests and  $n - m$  retrievals<sup>11</sup>. For a typical CPL document, we extract around 150 variables, which implies an upper bound of 450 prompts and 150 retrieval actions. This translates to approximately 8 minutes of runtime on a single RTX 3090 GPU.

<sup>10</sup>In the CPL case, variables like 'case ID', 'court name', and 'date' are always in the same place in the legislation document (typically, these values are placed at a fixed location in the title, before the first paragraph, or at the end).

<sup>11</sup>The number of model requests and retrievals depends on the document's content. For example, if the defendant has not appeared in court, logic is added to avoid extracting variables that would only have non-null values when the defendant is present in court.