

# Learning Interpretable Legal Case Retrieval via Knowledge-Guided Case Reformulation

Anonymous ACL submission

## Abstract

Legal case retrieval for sourcing similar cases is critical in upholding judicial fairness. Different from general web search, legal case retrieval involves processing lengthy, complex, and highly specialized legal documents. Existing methods in this domain often overlook the incorporation of legal expert knowledge, which is crucial for accurately understanding and modeling legal cases, leading to unsatisfactory retrieval performance. This paper introduces KELLER, a legal knowledge-guided case reformulation approach based on large language models (LLMs) for effective and interpretable legal case retrieval. By incorporating professional legal knowledge about crimes and law articles, we enable large language models to accurately reformulate the original legal case into concise sub-facts of crimes, which contain the essential information of the case. Extensive experiments on two legal case retrieval benchmarks demonstrate superior retrieval performance and robustness on complex legal case queries of KELLER over existing methods.

## 1 Introduction

Legal case retrieval is vital for legal experts to make informed decisions by thoroughly analyzing relevant precedents, which upholds justice and fairness (Hamann, 2019). This practice is crucial in both common law and civil law systems globally (Lastres, 2015; Harris, 2002). In civil law, although following past cases (known as "stare decisis") is not mandatory, judges are still highly advised to consider previous cases to improve the accuracy and trustworthiness of their judgments.

In legal case retrieval, both the query and the document are structured legal cases, distinguishing the task from other information retrieval (IR) tasks. Specifically, as shown in Figure 1, a legal case document comprises several sections, such as procedure, facts, and the court’s decision, making it much longer than typical queries and passages in

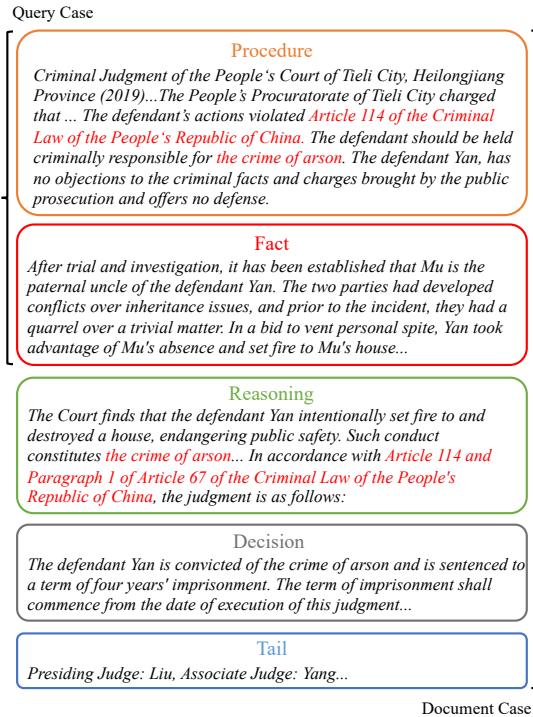


Figure 1: The query case and candidate document case examples. The query case typically contains only partial content since it has not been adjudicated. Extractable crimes and law articles are highlighted in red.

the standard ad-hoc search tasks. Its average text length often exceeds the maximum input limits of popular retrievers, such as 512 tokens (Devlin et al., 2019). Moreover, a legal case may encompass multiple, distinct criminal behaviors. Comprehensively considering all criminal behaviors of a legal case is important in determining its matching relevance with a query case. However, these key criminal descriptions are usually dispersed throughout the lengthy contents, which can significantly affect the effectiveness of traditional long document modeling strategies like FirstP and MaxP (Dai and Callan, 2019) in the legal domain.

To tackle the challenge of comprehending long and complex legal cases, previous works mainly

057 fall into two categories. The first approach focuses  
058 on expanding the context window size (Xiao et al.,  
059 2021) or splitting legal cases into passages (Shao  
060 et al., 2020). However, given the specialized and  
061 complex nature of legal texts, merely increasing the  
062 context window size still proves insufficient for sig-  
063 nificantly improving the retrieval performance. In  
064 contrast, the second approach performs direct text  
065 summarization (Askari and Verberne, 2021; Tang  
066 et al., 2023) or embedding-level summarization (Yu  
067 et al., 2022) on the legal case, aiming to only keep  
068 the most crucial information for assessing the rele-  
069 vance between legal cases. However, they typically  
070 only rely on heuristic rules or the models’ inher-  
071 ent knowledge for summarization. As the legal  
072 domain is highly specialized, existing approaches  
073 that overlook professional legal knowledge (e.g.,  
074 law articles) are likely to perform inaccurate sum-  
075 marization.

076 In this paper, we present a Knowledge-guidEd  
077 case reformuLation approach for LEgal case Re-  
078 trieval, named KELLER. Our main idea is to lever-  
079 age professional legal knowledge to guide large  
080 language models (LLMs) to summarize the corre-  
081 sponding key sub-facts for the crimes of the legal  
082 cases, and then directly learn to model case rele-  
083 vance based on these crucial and concise sub-facts.

084 Due to the specialization and complexity of the  
085 legal case, it is quite challenging to directly sum-  
086 marize the corresponding key sub-facts for all the  
087 crimes from the legal case, even using advanced  
088 LLMs (Tang et al., 2023). To address this problem,  
089 we propose a two-step legal knowledge-guided  
090 prompting method, as illustrated in the left side  
091 of Figure 2. In the initial step, we prompt LLM to  
092 extract all of the crimes and law articles contained  
093 in the legal case and then perform post-processing  
094 on them to establish correct mappings between  
095 the crimes and law articles by referring to the le-  
096 gal expert database. In the next step, we prompt  
097 LLM with the extracted “crime-article ” pairs to  
098 summarize the sub-fact of the crime from the le-  
099 gal case. The intermediate law articles, serving  
100 as high-level abstractions of the actual criminal  
101 events, can largely reduce the difficulty of identi-  
102 fying the corresponding sub-fact for the crime and  
103 improve accuracy. Figure 5 shows an example of  
104 three summarized sub-facts from a legal case.

105 Then, we directly model the case relevance  
106 based on these sub-facts because they are not only  
107 the most crucial information for relevance judg-  
108 ment in legal case retrieval but are also concise

109 enough to meet the text length limitations of popu-  
110 lar pre-trained retrieval models. For the comprehen-  
111 sive consideration of effectiveness, efficiency, and  
112 interoperability, we adopt the simple *MaxSim* and  
113 *Sum* operators to aggregate the relevance scores  
114 between query and document sub-facts to get the fi-  
115 nal case relevance score. The model is trained with  
116 dual-level contrastive learning to comprehensively  
117 capture the matching signals at the case level and  
118 the sub-fact level. On two widely-used datasets, we  
119 show that KELLER achieves new state-of-the-art  
120 results in both zero-shot and fine-tuning settings.  
121 Remarkably, KELLER demonstrates substantial  
122 improvements in handling complex queries.

123 Our main contributions can be summarized as:

124 (1) We propose to leverage professional legal  
125 knowledge about crimes and law articles to equip  
126 LLM with much-improved capabilities for summa-  
127 rizing essential sub-facts from complex cases.

128 (2) We suggest performing simple *MaxSim* and  
129 *Sum* aggregation directly on those refined sub-facts  
130 to achieve effective and interpretable legal retrieval.

131 (3) We introduce dual-level contrastive learning  
132 that enables the model to capture multi-granularity  
133 matching signals from both case-level and sub-fact-  
134 level for enhanced retrieval performance.

## 135 2 Related Work

136 **Legal case retrieval.** Existing legal case retrieval  
137 methods are categorized into statistical and neural  
138 models. Statistical models, notably the BM25 algo-  
139 rithm, can be enhanced by incorporating legal ex-  
140 pert knowledge such as legal summarization (Tran  
141 et al., 2020; Askari and Verberne, 2021), issue ele-  
142 ments (Zeng et al., 2005) and ontology (Saravanan  
143 et al., 2009). Neural models have been advanced  
144 through deep learning and the use of pre-trained  
145 language models (Devlin et al., 2019; Zhong et al.,  
146 2019; Chalkidis et al., 2020; Zhang et al., 2023).  
147 Recent advancements in this domain include the  
148 design of specialized pre-training tasks tailored for  
149 legal case retrieval, which yields remarkable im-  
150 provements in retrieval metrics (Li et al., 2023a;  
151 Ma et al., 2023b).

152 Due to the limitations of neural models in  
153 handling long texts, researchers mainly focus on  
154 processing lengthy legal documents by isolating  
155 the "fact description" section and truncating it  
156 to fit the model’s input constraints (Ma et al.,  
157 2021; Yao et al., 2022; Ma et al., 2023b; Li et al.,  
158 2023a). To overcome the long-text problem, some

159 other strategies include segmenting texts into  
160 paragraphs for interaction modeling (Shao et al.,  
161 2020), employing architectures like Longformer  
162 for extensive pre-training on legal texts (Xiao et al.,  
163 2021), and transforming token-level inputs into  
164 sentence-level encoding (Yu et al., 2022).

165  
166 **Query rewriting with LLMs.** Recently, re-  
167 searchers naturally employ LLMs to enhance  
168 the effectiveness of query rewriting (Zhu et al.,  
169 2023; Mao et al., 2023; Ma et al., 2023a; Wang  
170 et al., 2023; Jagerman et al., 2023). For instance,  
171 HyDE (Gao et al., 2023) creates pseudo passages  
172 for better query answers, integrating them into  
173 a vector for retrieval, while Query2Doc (Wang  
174 et al., 2023) employs few-shot methods to gener-  
175 ate precise responses. Furthermore, Jagerman  
176 et al. (2023) explores LLMs’ reasoning capacities  
177 to develop "Chain-of-Thoughts" responses for com-  
178 plex queries. However, the above methods struggle  
179 with legal case retrieval, where both queries and  
180 documents are lengthy cases. In the legal domain,  
181 PromptCase (Tang et al., 2023) attempts to address  
182 this by summarizing case facts within 50 words,  
183 but this approach often misses important details as  
184 many cases feature multiple independent facts.

### 185 3 Methodology

186 In this section, we first introduce some basic con-  
187 cepts in legal case retrieval. Then we delve into the  
188 three core parts of our KELLER, including legal  
189 knowledge-guided case reformulation, relevance  
190 modeling, and dual-level contrastive learning.

#### 191 3.1 Preliminaries

192 In legal case retrieval, both queries and candidate  
193 documents are real structured legal cases that can  
194 extend to thousands of tokens in length. Figure 1  
195 shows an illustration of the typical case structure.  
196 Specifically, a case usually contains several sec-  
197 tions, including *procedure*, *fact*, *reasoning*, *decisi-*  
198 *on*, and *tail*. Notably, the candidate documents  
199 are completed legal cases that have been through  
200 the adjudication process and therefore contain all  
201 sections. In contrast, the query cases are not yet  
202 adjudicated, so they usually only include the *proce-*  
203 *dure* and *fact* sections.

204 Formally, given a query case  $q$  and a set of docu-  
205 ment cases  $D$ , the objective of legal case retrieval is  
206 to calculate a relevance score  $s$  between the query  
207 case and each document case in  $D$ , and then rank

the document cases accordingly. 208

#### 209 3.2 Knowledge-Guided Case Reformulation

210 When assessing the relevance between two legal  
211 cases, the key facts of their crimes are the most  
212 crucial things for consideration. Therefore, given  
213 the complexity of the original legal cases which  
214 makes direct learning challenging, we try to first  
215 refine the legal cases into shorter but more essential  
216 “crime-fact” snippets. For example, we can get such  
217 a snippet from the case shown in Figure 1, whose  
218 crime is “*the crime of arson*” and the fact is “*Yan*  
219 *took advantage of Mu’s absence and set fire ...*”.

220 However, the description of a crime and its  
221 corresponding facts are often scattered throughout  
222 the lengthy case, and a single case may contain  
223 multiple crimes and facts, significantly com-  
224 plicating the extraction process. To tackle this  
225 problem, we propose a two-step prompting method  
226 leveraging professional legal knowledge to guide  
227 LLM to achieve accurate extraction. 228

229 **Crime and law article extraction.** First, we  
230 prompt LLM to extract all crimes and all law  
231 articles from the case. This step is relatively  
232 straightforward for LLM, as each crime and law  
233 article is a distinct, identifiable element within the  
234 text. For example, the extracted crime and law  
235 article for the case shown in Figure 1 are “*the*  
236 *crime of arson*” and “*Article 114 and Paragraph 1*  
237 *of Article 67 of the Criminal Law of the People’s*  
238 *Republic of China*”, respectively. Our extraction  
239 prompt is shown in Appendix B. 240

241 **Post-Processing.** The extracted law articles may  
242 just be the titles. We then expand these titles into  
243 full articles by gathering their detailed provision  
244 content from the Web based on the titles. Then,  
245 we establish a mapping between each crime and  
246 its relevant law articles by referring to a database  
247 built by our legal experts. Note that the correlation  
248 between specific crimes and their corresponding  
249 legal articles is objective, as it is clearly defined by  
250 law. After post-processing, we can obtain all the  
251 “crime-articles” pairs for a legal case. 252

253 **Fact summarization.** Next, we leverage the  
254 extracted crimes and their relevant law articles to  
255 guide LLM in summarizing the specific facts of  
256 each crime from the original legal case. The law  
257 articles, serving as high-level abstractions of the  
258 actual criminal events, can considerably simplify

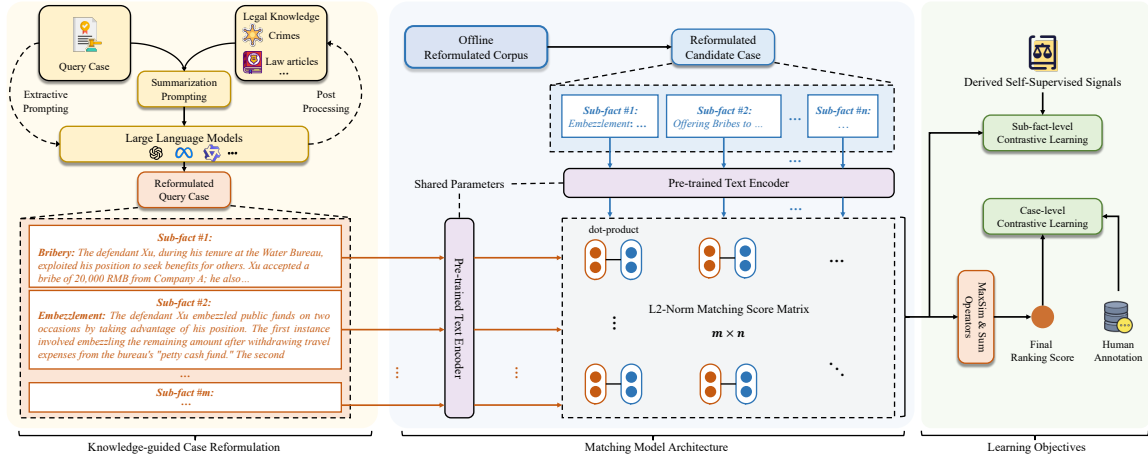


Figure 2: Overview of KELLER. We first perform legal knowledge-guided prompting to reformulate the legal cases into a series of crucial and concise sub-facts. Then, we directly model the case relevance based on the sub-facts. The model is trained at both the coarse-grained case level and the fine-grained sub-fact level via contrastive learning.

the task of identifying the corresponding specific facts. The prompt for fact summarization is shown in Appendix B.2.

Through our legal knowledge-guided reformulation, we can accurately distill a series of crimes and their corresponding specific facts from the originally lengthy legal cases. Finally, we form a *sub-fact* snippet, with the crime as the title and its facts as the main body. These refined sub-facts are not only the most crucial information for relevance judgment in legal case retrieval but are also concise enough to meet the text length limitations of popular pre-trained retrieval models. Please note that, since the required legal knowledge is present in criminal case documents from mainstream countries (e.g., China and the United States), our approach is actually internationally applicable. Our materials in Appendix D further prove this.

### 3.3 Relevance Modeling

We directly model the relevance of legal cases using the refined sub-facts, rather than relying on the full text of the original legal cases. Specifically, given a query case  $q = \{q_1, \dots, q_m\}$  and a candidate case  $d = \{d_1, \dots, d_n\}$ , where  $q_i$  represents the  $i$ -th sub-fact of  $q$  and  $d_j$  represents the  $j$ -th sub-fact of  $d$ . We utilize a pre-trained text encoder to encode them:

$$\begin{aligned} E_{q_i} &= \text{Pool}_{[\text{CLS}]}(\text{Encoder}(q_i)), \\ E_{d_j} &= \text{Pool}_{[\text{CLS}]}(\text{Encoder}(d_j)), \end{aligned} \quad (1)$$

where  $\text{Pool}_{[\text{CLS}]}$  means extracting the embedding output at the [CLS] token position. Then, we com-

pute the similarity matrix  $M_{m \times n}$  using the L2-norm dot product. Each element  $M_{i,j}$  of  $M$  is the similarity calculated between the normalized embeddings of the  $i$ -th sub-fact in the reformulated query case and  $j$ -th sub-fact in the reformulated document case:

$$M_{i,j} = \text{Sim}(E_{q_i}, E_{d_j}) = \text{Norm}(E_{q_i}) \cdot \text{Norm}(E_{d_j}^T). \quad (2)$$

Finally, we aggregate this similarity matrix to derive the matching score. There are various sophisticated choices for aggregation, such as using attention or kernel pooling (Xiong et al., 2017). In this paper, we opt to employ the *MaxSim* and *Sum* operators (Khattab and Zaharia, 2020):

$$s_{q,d} = \sum_{i=1}^m \text{Max}_{j=1}^n M_{i,j}, \quad (3)$$

where  $s_{q,d}$  is the final predicted relevance score. We choose these two operators because of their advantages in effectiveness, efficiency, and interpretability over the other aggregation approaches for our scenario:

**(1) Effectiveness:** Typically, each query’s sub-fact  $q_i$  matches one document sub-fact  $d_j$  at most in practice, which is well-suited for *MaxSim* of applying the Max operation across all document’s sub-facts for a given query’s sub-fact. For instance, considering a query sub-fact about “*drug trafficking*”, and the document sub-facts about “*drug trafficking*” and “*the discovery of privately stored guns and ammunition*”, only the “*drug trafficking*” sub-fact of the document is relevant for providing matching evidence. In contrast, using soft aggregation meth-



ods (e.g., kernel pooling (Xiong et al., 2017)) may introduce additional noise in this scenario.

(2) **Efficiency:** *Maxsim* and *Sum* operations on tensors are quite efficient for both re-ranking and large-scale top- $k$  retrieval supported by multi-vector-based Approximate Nearest Neighbor algorithms (Khattab and Zaharia, 2020). This high efficiency is important for meeting the low-latency requirements of the practical use.

(3) **Interpretability:** *MaxSim* provides clear interpretability by revealing the quantitative contribution of each query and document sub-fact towards the final relevance score, which can aid in understanding the ranking strategies and justifying the retrieval results. We further illustrate this advantage by studying a real case in Section 4.6.

### 3.4 Dual-Level Contrastive Learning

We incorporate matching signals from both the coarse-grained case level and the fine-grained sub-fact level to comprehensively enhance the model performance in legal case matching.

**Case-level contrastive learning.** At the case level, we consider directly optimizing toward the final matching score between the query case and the document cases. Specifically, we employ the classical ranking loss function to promote the relevance score between the query and the positive document while reducing it for negative documents:

$$\mathcal{L}_R = -\log \frac{\exp(s_{q,d^+}/\tau)}{\exp(s_{q,d^+}/\tau) + \sum_{d^-} \exp(s_{q,d^-}/\tau)}, \quad (4)$$

where  $d^+$  is the positive document of the query  $q$  and each  $d^-$  is from the in-batch negatives.  $\tau$  is a temperature parameter.

**Sub-fact-level contrastive learning.** At the sub-fact level, we incorporate intermediate relevance signals among sub-facts to fine-grainedly enhance the model’s effectiveness in understanding sub-facts’ content and their matching relationships. However, only the case-level relevance labels are available in the dataset. Naively considering all the sub-fact pairs between the query and the positive documents as positives and all the sub-fact pairs between the query and the negative documents as negatives will introduce substantial false positive and negative noise. To mitigate this issue, we propose a heuristic strategy to obtain high-quality relevance labels for the query’s sub-facts  $\{q_1, \dots, q_m\}$ . The

core idea of this strategy is to combine the case-level relevance and the charges of each sub-fact to accurately identify true positive and negative samples. We introduce the details of this strategy in Appendix C due to the space limitation.

After getting the sub-fact level relevance labels, we also adopt the ranking loss function for sub-fact level contrastive learning:

$$\mathcal{L}_S = -\log \frac{\exp(s_{M_{i,j^+}}/\tau)}{\exp(s_{M_{i,j^+}}/\tau) + \sum_{J^-} \exp(s_{M_{i,j^-}}/\tau)}, \quad (5)$$

where  $M_{i,j^+}$  are the similarity score between  $q_i$  and its positive document.  $M_{i,j^-}$  are the similarity score between  $q_i$  and its negative document sub-fact.  $J^-$  is the collection of all negative document sub-facts for  $q_i$ . The final learning objective is the combination of  $\mathcal{L}_R$  and  $\mathcal{L}_S$ :

$$\mathcal{L} = \mathcal{L}_R + \alpha \mathcal{L}_S, \quad (6)$$

where  $\alpha$  is a hyper-parameter to adjust the weights of two losses.

## 4 Experiments

### 4.1 Experimental Setup

**Dataset and evaluation metrics.** We conduct extensive experiments on two widely-used datasets: LeCaRD (Ma et al., 2021) and LeCaRDv2 (Li et al., 2023b), whose statistics are listed in Appendix A.1. Considering the limited number of queries in LeCaRD, we directly evaluate all the queries of LeCaRD using the best model trained on LeCaRDv2, thereby avoiding the need for dataset split. Following the previous studies (Li et al., 2023a,b), we regard label=3 in LeCaRD and label $\geq$ 2 in LeCaRDv2 as positive. For the query whose candidate documents are all annotated as positive, we supplement the candidate pool by sampling 10 document cases from the top 100-150 BM25 results. To exclude the effect of unlabeled potential positives in the corpus, we rank the candidate pools and adopt MAP, P@k (k=3), and NDCG@k (k=3, 5, 10) as our evaluation metrics.

**Baselines.** We compare KELLER against the following baselines across three categories. The first is *traditional probabilistic models*, including TF-IDF and BM25. The second is *ranking methods based on pre-trained language models*, including BERT (Devlin et al., 2019), RoBERTa (Liu et al., 2019), BGE (Xiao et al., 2023) and SAILER (Li

Table 1: Main results of the fine-tuned setting on LeCaRD and LeCaRDv2. “†” indicates our approach outperforms all baselines significantly with paired t-test at  $p < 0.05$  level. The best results are in bold.

Model	LeCaRD					LeCaRDv2				
	MAP	P@3	NDCG@3	NDCG@5	NDCG@10	MAP	P@3	NDCG@3	NDCG@5	NDCG@10
<i>Traditional ranking baselines</i>										
BM25	47.30	40.00	64.45	65.59	69.15	55.20	48.75	72.11	72.51	79.85
TF-IDF	42.59	36.19	58.14	59.98	63.37	55.19	47.92	71.38	72.70	75.04
<i>PLM-based neural ranking baselines</i>										
BERT	53.83	50.79	73.19	73.43	75.54	60.66	53.12	77.78	78.73	80.85
RoBERTa	55.79	53.33	74.40	74.33	76.70	59.75	53.12	78.15	78.97	80.70
BGE	54.98	53.33	74.29	74.09	75.65	60.64	51.87	76.99	78.43	80.90
SAILER	57.98	56.51	77.55	77.04	79.41	60.62	54.58	78.67	78.99	81.41
<i>Neural ranking baselines designed for long text</i>										
BERT-PLI	48.16	43.80	65.74	68.14	71.32	55.34	46.67	71.62	73.68	76.63
Lawformer	54.58	50.79	73.19	73.43	75.54	60.17	54.17	78.23	78.99	81.40
<i>Case reformulation with LLMs</i>										
PromptCase	59.71	55.92	78.75	78.44	80.71	62.25	54.19	78.51	79.07	81.26
KELLER	<b>66.84†</b>	<b>57.14</b>	<b>81.24†</b>	<b>82.42†</b>	<b>84.67†</b>	<b>68.29†</b>	<b>63.13†</b>	<b>84.97†</b>	<b>85.63†</b>	<b>87.61†</b>

et al., 2023a). The third is *ranking methods designed for handling long (legal) text*, including BERT-PLI (Shao et al., 2020), Lawformer (Xiao et al., 2021), and PromptCase (Tang et al., 2023).

**Implementations.** We introduce the selected language models, hyperparameter settings and other details in Appendix A.2.

## 4.2 Main Results

The main results are as shown in Table 1 and we have the following observations:

(1) **KELLER outperforms all baseline methods across all metrics on both datasets.** Compared with previous methods tailored for the long-text problem, KELLER employs knowledge-guided case reformulation to address the challenge of long-text comprehension. This demonstrates the effectiveness of separating comprehension and matching tasks in the domain of legal case retrieval.

(2) **After fine-tuning on legal case retrieval datasets, the performance gap between general-purpose and retrieval-oriented PLMs becomes less distinct.** This observation may stem from two reasons. First, the scarcity of training data in the legal case retrieval task can induce overfitting to annotation signals, which hampers the model’s generalization capabilities. Second, Naive truncation of lengthy texts can make the model’s inputs lose sufficient matching signals, leading to inconsistencies between relevance annotations and matching evidence.

(3) **We observe that these long-text-oriented baseline methods do not show significant advantages.** Despite BERT-PLI and Lawformer pro-

cessing more text than other methods, their input capacity was still insufficient for the average length of legal cases. Handling both long-text processing and complex semantic understanding within one retriever presents a significant challenge. To address this issue, our approach offloads a portion of the long-text comprehension task via knowledge-guided case reformulation and improves the ranking performance.

## 4.3 Zero-shot Evaluation

Considering the inherent data scarcity problem in legal case retrieval, we evaluate the zero-shot performance (i.e., without fine-tuning on the training set of LeCaRDv2) of models on LeCaRDv2.

Results are shown in Table 2 and we find that KELLER consistently outperforms baselines in both zero-shot and fine-tuning settings. Upon comparing the performance of each method under zero-shot and fine-tuned settings, we observe that most methods benefit from fine-tuning except SAILER. Intuitively, models trained in a general domain or task could be enhanced through fine-tuning. In specific domains, continued fine-tuning of models generally does not lead to a significant decrease in performance. We posit that the unexpected outcomes in the SAILER model primarily arise from overfitting the limited data used for fine-tuning, which impairs the generalization capabilities established in the pre-training phase.

## 4.4 Ablation Study

We design the following six ablations: (1) *KGCR*→*NS*: We replace our Knowledge-Guided Case Reformulation (KGCR) with a Naive Sum-

Table 2: Zero-shot performance on LeCaRD and LeCaRDv2. “†” indicates our approach outperforms all baselines significantly with paired t-test at  $p < 0.05$  level. The best results are in bold.

Model	LeCaRD					LeCaRDv2				
	MAP	P@3	NDCG@3	NDCG@5	NDCG@10	MAP	P@3	NDCG@3	NDCG@5	NDCG@10
<i>General PLM-based baselines</i>										
BERT	42.92	37.78	60.11	61.37	64.10	56.46	52.08	75.82	77.05	79.39
RoBERTa	51.50	47.62	69.21	71.07	73.60	57.89	52.08	75.48	76.33	78.38
Lawformer	42.80	38.41	59.46	61.61	64.13	55.05	49.58	74.42	74.31	76.96
<i>Retrieval-oriented pre-training baselines</i>										
BGE	51.81	47.62	68.57	69.91	72.61	57.21	50.42	73.59	75.36	77.80
SAILER	60.62	56.19	79.93	78.99	81.41	62.80	55.00	79.38	81.17	83.83
KELLER	<b>64.17†</b>	<b>57.78</b>	<b>80.47</b>	<b>81.43†</b>	<b>84.36†</b>	<b>65.87†</b>	<b>61.67†</b>	<b>83.33†</b>	<b>83.75†</b>	<b>86.06†</b>

marization (NS), which produces case summaries without hierarchical structure. We subsequently optimize the dual encoders with this text as the input. (2)  $MS \rightarrow Mean$ : We replace *MaxSim* and *Sum* (MS) with *Mean* to capture the average relevance of each sub-fact in the candidate cases to the query. (3)  $MS \rightarrow NC$ : We Naively Concatenate (NC) all the reformulated sub-facts into a text sequence and subsequently optimize the dual-encoders. (4)  $MS \rightarrow KP$ : We employ kernel pooling (Xiong et al., 2017) on the score matrix to capture relevance signals. (5) *w/o sfCL*: Training without the sub-fact-level contrastive learning. (6) *w/o SfCL*: Training without the case-level contrastive learning.

Results are shown in Table 3 and we can observe:

(1) Every ablation strategy results in a decline in the model’s performance, demonstrating the effectiveness of each module within KELLER. This outcome indicates that KELLER’s architecture is both comprehensive and synergistic, with each module contributing to the model’s overall performance.

(2) The replacement of the KGCR module exhibits the most significant impact on performance. This highlights the pivotal role of the KGCR module in KELLER. The KGCR module decomposes cases into structured sub-facts, which are crucial for the model’s learning process.

(3) Among different aggregation strategies,  $MS \rightarrow Mean$  demonstrates the least performance degradation. This is primarily because the dataset mainly consists of simple cases with single charges, where *Mean* and *MS* become essentially equivalent. Conversely,  $MS \rightarrow NC$  exhibits the most notable performance decline. This is mainly because the model no longer maintains a cross-matching architecture after the concatenation operation. Merging multiple facts into a single representation negatively impacts representation learning.

Table 3: Results of ablation study on LeCaRDv2.

Strategy	MAP	P@3	NDCG@3	NDCG@5	NDCG@10
<i>Effect of knowledge-guided case reformulation</i>					
<i>KGCR</i> →NS	61.91	55.13	79.50	79.11	81.47
<i>Effect of different aggregation strategy</i>					
<i>MS</i> → <i>Mean</i>	67.15	61.81	81.58	84.42	86.74
<i>MS</i> → <i>NC</i>	63.35	57.92	80.37	81.99	84.04
<i>MS</i> → <i>KP</i>	65.47	60.06	79.87	83.61	85.39
<i>Effect of contrastive learning</i>					
<i>w/o SfCL</i>	67.39	61.93	81.24	84.73	86.91
<i>w/o CaCL</i>	67.18	61.67	82.76	84.45	86.51
KELLER	<b>68.29</b>	<b>63.13</b>	<b>84.97</b>	<b>85.63</b>	<b>87.61</b>

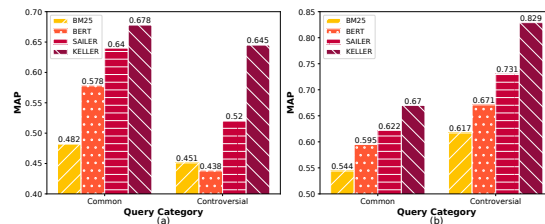


Figure 3: Evaluation on different query types. We evaluate four models on (a) LeCaRD and (b) LeCaRDv2.

#### 4.5 Evaluations on Different Query Types

We investigate the two query types presented in both LeCaRD and LeCaRDv2: *common* and *controversial*. Common queries are similar to initial trials, and controversial queries to retrials, which are typically more complex and require additional expert review. We evaluated multiple models on these query types. Notably, SAILER’s performance declined after fine-tuning, so we included its zero-shot results for comparison, alongside the fine-tuned outcomes of other models. Results as shown in Figure 3 and we find:

(1) KELLER outperformed other models on both query types, showing more substantial gains in controversial queries with improvements of 24.04% and 13.41% in the LeCaRD and LeCaRDv2

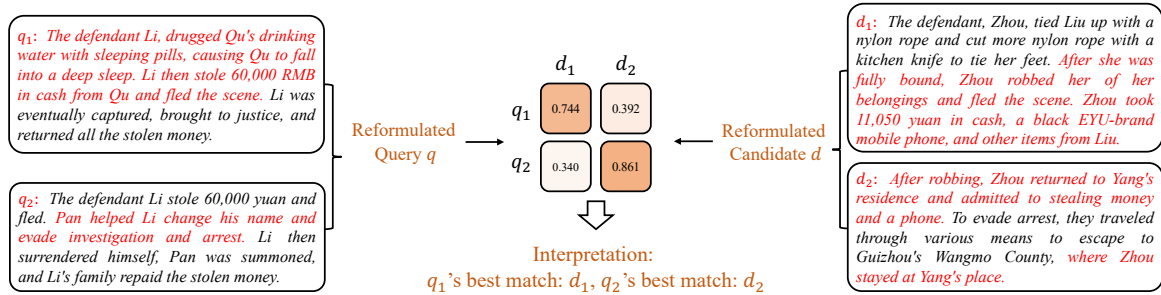


Figure 4: An example of the interpretability of KELLER. We can observe that each sub-fact of the query finds a correct match in the candidate document (in red).

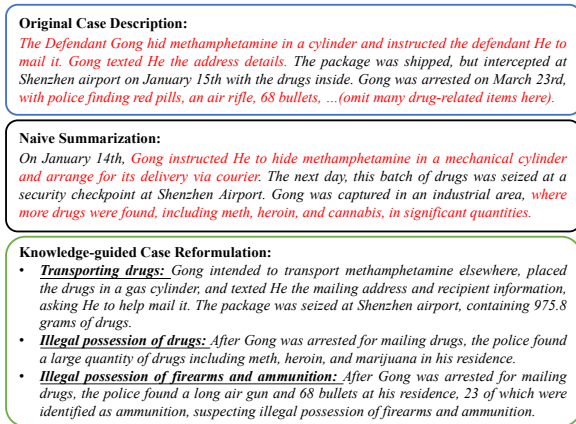


Figure 5: Comparison of the original text, naive summarization, and our proposed knowledge-guided case reformulation. The original text is manually abbreviated due to its length. Important sentences are marked in red.

535 datasets, respectively. This enhanced performance  
 536 is credited to KELLER’s novel case reformulation,  
 537 which simplifies complex scenarios into sub-facts,  
 538 aiding in better comprehension and matching.

539 (2) In the LeCaRD dataset, lexical-based models  
 540 showed consistent performance across different  
 541 queries, unlike representation-based models  
 542 which varied significantly. For example, BERT  
 543 outperformed BM25 on common queries but was  
 544 less effective on controversial ones, a difference  
 545 attributed to the models’ limited ability to handle  
 546 multifaceted cases. KELLER’s cross-matching  
 547 architecture successfully addresses this limitation.

#### 548 4.6 Case Studies

549 **Case reformulation.** We provide an illustrative  
 550 comparison between the original case description,  
 551 naive summarization, and our knowledge-guided  
 552 case reformulation in Figure 5. The case centers  
 553 on complex issues of drug transport and  
 554 firearm possession. Most details focus on drug  
 555 transportation, with brief mentions of firearms

556 found at the defendant’s residence towards the  
 557 end. Given the 512-token limit of most retrievers,  
 558 crucial information about the firearms is often  
 559 inaccessible. While naive summarization captures  
 560 the main points, it overlooks specifics about  
 561 the firearms in the context of drug offenses. In  
 562 contrast, our KGCR method segments the case  
 563 into three topics—drug transportation, illegal drug  
 564 possession, and illegal firearms possession—thus  
 565 detailing each criminal aspect comprehensively.

566 **Interpretability.** In KELLER, each sub-fact in  
 567 a query represents a specific intent of the query,  
 568 with the highest match score from a candidate case  
 569 indicating how well this intent is met. KELLER  
 570 allows users to see which sub-fact in a candidate  
 571 case matches their intent. For example, in a case  
 572 involving robbery and harboring crimes shown in  
 573 Figure 4, KELLER accurately matches sub-facts  
 574 in the query to those in the candidate case, demon-  
 575 strating the alignment of KELLER’s scoring with  
 576 the underlying legal facts of the case. The matching  
 577 is shown in a matrix, where the positions  $(q_1, d_1)$   
 578 and  $(q_2, d_2)$  highlight the defendant’s actions in the  
 579 query and the candidate case, respectively, estab-  
 580 lishing a direct correlation between the computed  
 581 scores and the case ranking.

## 582 5 Conclusion

583 In this paper, we introduce KELLER, a ranking  
 584 model that effectively retrieves legal cases with  
 585 high interpretability. KELLER structures legal doc-  
 586 uments into hierarchical texts using LLMs and de-  
 587 termines relevance through a cross-matching mod-  
 588 ule. Our tests on two expert-annotated datasets  
 589 validate its effectiveness. In the future, we will  
 590 enhance KELLER by incorporating additional spe-  
 591 cialized knowledge and generative models to refine  
 592 performance and produce language explanations.  
 593



## 6 Limitations

**External Knowledge base Construction.** Our method requires constructing a legal knowledge base to assist in case reformulation, which introduces an extra step compared to the out-of-the-box dense retrievers. This issue is common in most domain-specific knowledge-enhanced methods.

**Computing Efficiency.** Our approach needs to call large language models when processing the query case, which may bring additional computational costs. In our experiments, we have employed techniques such as vLLM to achieve high-speed inference. Furthermore, we believe that with ongoing advancements in techniques in both hardware and algorithms, the computational of utilizing LLMs for processing individual query cases online will be acceptable. For example, Llama3-8B can achieve a speed exceeding 800 tokens per second on the Groq platform, while recent inference services provided by Qwen and DeepSeek require less than \$0.0001 per 1,000 tokens.

## 7 Ethical Discussion

The application of artificial intelligence in the legal domain is sensitive, requiring careful examination and clarification of the associated ethical implications. The two datasets utilized in our experimental analysis have undergone anonymization processes, particularly with regard to personally identifiable information such as names.

Although KELLER demonstrates superior performance on two human-annotated datasets, its recommendations for similar cases may sometimes be imprecise when dealing with intricate real-world queries. Additionally, the case databases in existing systems may not consistently include cases that fully satisfy user requirements. The choice to reference the retrieved cases should remain at the discretion of the experts.

## References

Arian Askari and Suzan Verberne. 2021. [Combining lexical and neural retrieval with longformer-based summarization for effective case law retrieval](#). In *Proceedings of the Second International Conference on Design of Experimental Search & Information REtrieval Systems, Padova, Italy, September 15-18, 2021*, volume 2950 of *CEUR Workshop Proceedings*, pages 162–170. CEUR-WS.org.

Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei

Huang, et al. 2023. Qwen technical report. *arXiv preprint arXiv:2309.16609*. 643  
644

Ilias Chalkidis, Manos Fergadiotis, Prodromos Malakasiotis, Nikolaos Aletras, and Ion Androutsopoulos. 2020. [LEGAL-BERT: the muppets straight out of law school](#). *CoRR*, abs/2010.02559. 645  
646  
647  
648

Zhuyun Dai and Jamie Callan. 2019. Deeper text understanding for ir with contextual neural language modeling. In *Proceedings of the 42nd international ACM SIGIR conference on research and development in information retrieval*, pages 985–988. 649  
650  
651  
652  
653

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, pages 4171–4186. Association for Computational Linguistics. 654  
655  
656  
657  
658  
659  
660  
661  
662  
663

Luyu Gao, Xueguang Ma, Jimmy Lin, and Jamie Callan. 2023. [Precise zero-shot dense retrieval without relevance labels](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023*, pages 1762–1777. Association for Computational Linguistics. 664  
665  
666  
667  
668  
669  
670

Hanjo Hamann. 2019. The german federal courts dataset 1950–2019: from paper archives to linked open data. *Journal of empirical legal studies*, 16(3):671–688. 671  
672  
673  
674

Bruce V Harris. 2002. Final appellate courts overruling their own "wrong" precedents: the ongoing search for principle. *Law Quarterly Review*, 118(July 2002):408–427. 675  
676  
677  
678

Rolf Jagerman, Honglei Zhuang, Zhen Qin, Xuanhui Wang, and Michael Bendersky. 2023. [Query expansion by prompting large language models](#). *CoRR*, abs/2305.03653. 679  
680  
681  
682

Omar Khattab and Matei Zaharia. 2020. Colbert: Efficient and effective passage search via contextualized late interaction over bert. In *Proceedings of the 43rd International ACM SIGIR conference on research and development in Information Retrieval*, pages 39–48. 683  
684  
685  
686  
687  
688

Steven A Lastres. 2015. Rebooting legal research in a digital age. 689  
690

Haitao Li, Qingyao Ai, Jia Chen, Qian Dong, Yueyue Wu, Yiqun Liu, Chong Chen, and Qi Tian. 2023a. [SAILER: structure-aware pre-trained language model for legal case retrieval](#). In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR 2023, Taipei, Taiwan, July 23-27, 2023*, pages 1035–1044. ACM. 691  
692  
693  
694  
695  
696  
697  
698

699	Haitao Li, Yunqiu Shao, Yueyue Wu, Qingyao Ai,	<a href="#">models</a> . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023</i> , pages 9414–9423. Association for Computational Linguistics.	755
700	Yixiao Ma, and Yiqun Liu. 2023b. Lecardv2: A		756
701	large-scale chinese legal case retrieval dataset. <i>arXiv preprint arXiv:2310.17609</i> .		757
702			758
703	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Man-		759
704	dar Joshi, Danqi Chen, Omer Levy, Mike Lewis,	Chaojun Xiao, Xueyu Hu, Zhiyuan Liu, Cunchao Tu,	760
705	Luke Zettlemoyer, and Veselin Stoyanov. 2019.	and Maosong Sun. 2021. Lawformer: A pre-trained	761
706	<a href="#">Roberta: A robustly optimized BERT pretraining</a>	language model for chinese legal long documents. <i>AI</i>	762
707	<a href="#">approach</a> . <i>CoRR</i> , abs/1907.11692.	<i>Open</i> , 2:79–84.	763
708	Xinbei Ma, Yeyun Gong, Pengcheng He, Hai Zhao,	Shitao Xiao, Zheng Liu, Peitian Zhang, and Niklas	764
709	and Nan Duan. 2023a. <a href="#">Query rewriting for</a>	Muennighof. 2023. <a href="#">C-pack: Packaged resources</a>	765
710	<a href="#">retrieval-augmented large language models</a> . <i>CoRR</i> ,	<a href="#">to advance general chinese embedding</a> . <i>CoRR</i> ,	766
711	abs/2305.14283.	abs/2309.07597.	767
712	Yixiao Ma, Yunqiu Shao, Yueyue Wu, Yiqun Liu,	Chenyang Xiong, Zhuyun Dai, Jamie Callan, Zhiyuan	768
713	Ruizhe Zhang, Min Zhang, and Shaoping Ma. 2021.	Liu, and Russell Power. 2017. <a href="#">End-to-end neural</a>	769
714	Lecard: a legal case retrieval dataset for chinese law	<a href="#">ad-hoc ranking with kernel pooling</a> . In <i>Proceedings</i>	770
715	system. In <i>Proceedings of the 44th international</i>	<i>of the 40th International ACM SIGIR Conference on</i>	771
716	<i>ACM SIGIR conference on research and development</i>	<i>Research and Development in Information Retrieval,</i>	772
717	<i>in information retrieval</i> , pages 2342–2348.	<i>Shinjuku, Tokyo, Japan, August 7-11, 2017</i> , pages	773
718	Yixiao Ma, Yueyue Wu, Weihang Su, Qingyao Ai,	55–64. ACM.	774
719	and Yiqun Liu. 2023b. <a href="#">Caseencoder: A knowledge-</a>	Feng Yao, Chaojun Xiao, Xiaozhi Wang, Zhiyuan Liu,	775
720	<a href="#">enhanced pre-trained model for legal case encoding</a> .	Lei Hou, Cunchao Tu, Juanzi Li, Yun Liu, Weixing	776
721	In <i>Proceedings of the 2023 Conference on Empirical</i>	Shen, and Maosong Sun. 2022. <a href="#">LEVEN: A large-</a>	777
722	<i>Methods in Natural Language Processing, EMNLP</i>	<a href="#">scale chinese legal event detection dataset</a> . In <i>Find-</i>	778
723	<i>2023, Singapore, December 6-10, 2023</i> , pages 7134–	<i>ings of the Association for Computational Linguistics:</i>	779
724	7143. Association for Computational Linguistics.	<i>ACL 2022, Dublin, Ireland, May 22-27, 2022</i> , pages	780
725	Kelong Mao, Zhicheng Dou, Fengran Mo, Jiewen Hou,	183–201. Association for Computational Linguistics.	781
726	Haonan Chen, and Hongjin Qian. 2023. <a href="#">Large lan-</a>	Weijie Yu, Zhongxiang Sun, Jun Xu, Zhenhua Dong,	782
727	<a href="#">guage models know your contextual search intent: A</a>	Xu Chen, Hongteng Xu, and Ji-Rong Wen. 2022.	783
728	<a href="#">prompting framework for conversational search</a> . In	<a href="#">Explainable legal case matching via inverse optimal</a>	784
729	<i>Findings of the Association for Computational Lin-</i>	<a href="#">transport-based rationale extraction</a> . In <i>SIGIR '22:</i>	785
730	<i>guistics: EMNLP 2023, Singapore, December 6-10,</i>	<i>The 45th International ACM SIGIR Conference on</i>	786
731	<i>2023</i> , pages 1211–1225. Association for Computa-	<i>Research and Development in Information Retrieval,</i>	787
732	tional Linguistics.	<i>Madrid, Spain, July 11 - 15, 2022</i> , pages 657–668.	788
733	Manavalan Saravanan, Balaraman Ravindran, and Shiv-	ACM.	789
734	ani Raman. 2009. Improving legal information retrie-	Yiming Zeng, Ruili Wang, John Zeleznikow, and Eliz-	790
735	val using an ontological framework. <i>Artificial</i>	abeth A. Kemp. 2005. <a href="#">Knowledge representation</a>	791
736	<i>Intelligence and Law</i> , 17:101–124.	<a href="#">for the intelligent legal case retrieval</a> . In <i>Knowledge-</i>	792
737	Yunqiu Shao, Jiaxin Mao, Yiqun Liu, Weizhi Ma, Ken	<i>Based Intelligent Information and Engineering Sys-</i>	793
738	Satoh, Min Zhang, and Shaoping Ma. 2020. Bert-pli:	<i>tems, 9th International Conference, KES 2005, Mel-</i>	794
739	Modeling paragraph-level interactions for legal case	<i>bourne, Australia, September 14-16, 2005, Proceed-</i>	795
740	retrieval. In <i>IJCAI</i> , pages 3501–3507.	<i>ings, Part I</i> , volume 3681 of <i>Lecture Notes in Com-</i>	796
741	Yanran Tang, Ruihong Qiu, and Xue Li. 2023. <a href="#">Prompt-</a>	<i>puter Science</i> , pages 339–345. Springer.	797
742	<a href="#">based effective input reformulation for legal case</a>	Kun Zhang, Chong Chen, Yuanzhuo Wang, Qi Tian, and	798
743	<a href="#">retrieval</a> . In <i>Databases Theory and Applications -</i>	Long Bai. 2023. <a href="#">Cfgl-lcr: A counterfactual graph</a>	799
744	<i>34th Australasian Database Conference, ADC 2023,</i>	learning framework for legal case retrieval. In <i>Pro-</i>	800
745	<i>Melbourne, VIC, Australia, November 1-3, 2023, Pro-</i>	<i>ceedings of the 29th ACM SIGKDD Conference on</i>	801
746	<i>ceedings</i> , volume 14386 of <i>Lecture Notes in Com-</i>	<i>Knowledge Discovery and Data Mining</i> , pages 3332–	802
747	<i>puter Science</i> , pages 87–100. Springer.	3341.	803
748	Vu Tran, Minh Le Nguyen, Satoshi Tojo, and Ken Satoh.	Haoxi Zhong, Zhengyan Zhang, Zhiyuan Liu, and	804
749	2020. Encoded summarization: summarizing docu-	Maosong Sun. 2019. <a href="#">Open chinese language pre-</a>	805
750	ments into continuous vector space for legal case	<a href="#">trained model zoo</a> . Technical report.	806
751	retrieval. <i>Artificial Intelligence and Law</i> , 28:441–	Yutao Zhu, Huaying Yuan, Shuting Wang, Jiongnan Liu,	807
752	467.	Wenhan Liu, Chenlong Deng, Zhicheng Dou, and	808
753	Liang Wang, Nan Yang, and Furu Wei. 2023.	Ji-Rong Wen. 2023. <a href="#">Large language models for infor-</a>	809
754	<a href="#">Query2doc: Query expansion with large language</a>	<a href="#">mation retrieval: A survey</a> . <i>CoRR</i> , abs/2308.07107.	810

Table 4: Basic statistics of the datasets.

Dataset	LeCaRD	LeCaRDv2
# Train queries	-	640
# Test queries	107	160
# Documents	9,195	55,192
Average query length	445	4,499
Average doc length	7,446	4,768
Average golden docs / query	10.39	13.65

## A More Details for Experimental Setup

### A.1 Datasets

The statistics of both datasets are listed in Table 4. LeCaRD comprises 107 queries and 10,700 candidate cases. LeCaRDv2, a more extensive collection, includes 800 queries and 55,192 candidate cases.

### A.2 Implementation Details

For baseline models, we employ the default parameter settings of Okapi-BM25 in the implementation of BM25. For ranking methods based on PLMs, a uniform learning rate of  $1e-5$  and a batch size of 128 are consistently applied. In BERT-PLI, the numbers of queries and candidate case segments are set to 3 and 4, respectively, with a maximum segment length of 256. For Lawformer, the maximum text input length is set to 3,072, optimized using a learning rate of  $1e-5$  and a batch size of 64.

In KELLER, we employ the Qwen-72B-Chat (Bai et al., 2023), which is currently one of the best open-source Chinese LLMs, to perform case reformulation. We do not choose OpenAI API due to concerns about reproducibility and high cost. All prompts, except for the case description, are input as system prompts. In the ranking model, the maximum number of crimes per case is capped at 4, which meets the needs of most cases. We adopt the pre-trained retriever SAILER as the text encoder. The  $\tau$  in the contrastive learning is 0.01, and the  $\alpha$  in the final loss function is 0.9. We conduct model training with a learning rate of  $1e-5$  and a batch size of 128. All experiments are conducted on four Nvidia Tesla A100-40G GPUs. All the source code and data will be shared at <https://github.com/hide-for-blind-review> if the paper gets accepted.

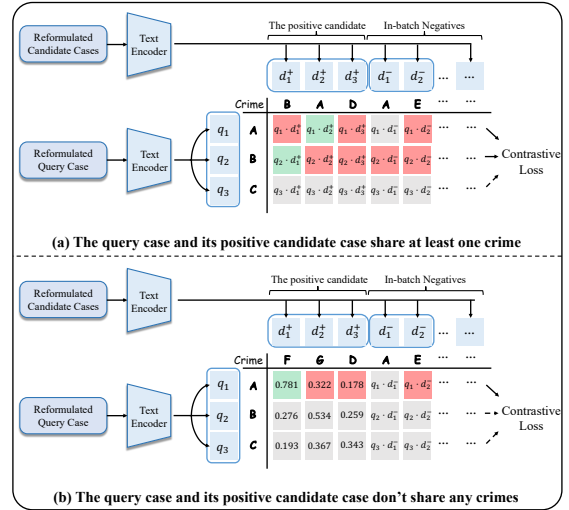


Figure 6: Illustration of our proposed sub-fact-level contrastive learning. The green and red squares represent the positive pairs and negative pairs, respectively. The gray squares are the discarded pairs that are not used for training. The blue rounded rectangles encompass blue squares belonging to the same query/document case. {A, ..., G} are crimes.

## B Prompts

### B.1 Extraction Prompt

**Extraction Prompt:** *You are now a legal expert, and your task is to find all the crimes and law articles in the procuratorate’s charges (or court judgments) from the provided case. The output format is one line each for crimes and law articles, two lines in total. Multiple crimes (law articles) are separated by semicolons.*

### B.2 Summarization Prompt

**Summarization Prompt:** *You are now a legal expert, and you are good at analyzing lengthy legal case texts containing multiple circumstances of crime. Your task is to concisely summarize the causes, procedures, and outcomes associated with a specified crime, ensuring each part does not exceed 100 words.*

[Crime]: the specific crime name

[Law Articles]: the specific provisions of law articles

## C Strategy to Obtain Sub-Fact-Level Relevance Labels

Specifically, for a positive document  $d^+$  of query  $q$ , we first check whether any of the document sub-facts share the same crimes as any of the query sub-facts:

- If it exists, as shown in Figure 6(a), for a query sub-fact  $q_i$ , we treat the document sub-facts that share the same crime as the positives (e.g., the green rectangles in columns  $d_1^+$ ,  $d_2^+$ , and  $d_3^+$ ), and all the other document sub-facts as negatives (e.g., the red rectangles in columns  $d_1^+$ ,  $d_2^+$ , and  $d_3^+$ ). If the crime of  $q_i$  is different from any of the document sub-facts, we will not include  $q_i$  for training (e.g., the gray rectangles in row  $q_3$ ).
- If not, as shown in Figure 6 (b), we select the  $(q_i, d_j^+)$  which has the highest similarity score as a positive training pair (e.g., the green rectangle), and retain any  $(q_i, d_k^+ (k \neq j))$  as negatives (e.g., the red rectangles in columns  $d_2^+$  and  $d_3^+$ ). All the other query and document sub-fact pairs are discarded (e.g., the gray rectangles in columns  $d_1^+$ ,  $d_2^+$ , and  $d_3^+$ ).

Then, for a negative document  $d^-$  of one query sub-fact  $q_i$ , we first check whether  $q_i$  has one positive sample.

- If not, we discard all the document sub-facts because there doesn't exist a positive sample for contrastive learning (e.g., the gray rectangles of row  $q_3$  in Figure 6 (a) and (b)).
- If it exists, we further check whether one of its document sub-facts  $d_j^-$  shares the same crime as a  $q_i$ .
  1. Both  $d_j^-$  and  $q_i$  are implicated to the same crime. we will include all  $(q_i, d_k^- (k \neq j))$  as negatives (e.g., the red rectangles of column  $d_1^-$  and  $d_2^-$  in Figure 6 (a) and (b)). All the other sub-facts are discarded to avoid introducing false negatives (e.g., the gray rectangles of  $(q_1, d_1^-)$  in Figure 6 (a) and (b)).
  2. None of  $d_j^-$  and  $q_i$  pertain to the same crime. We will include all  $(q_i, d_j^-)$  as negatives (e.g., the red rectangles of  $(q_2, d_1^-)$  and  $(q_2, d_2^-)$  in Figure 6 (a)).

## D Case Format of Other Regions

To demonstrate the international applicability of our method, we use U.S. legal documents as examples. Figure 7 and Figure 8 depict the formats of a U.S. indictment and a judgment document, respectively. It is evident that the legal knowledge required by our method (a combination of charges and law articles in this paper) is commonly

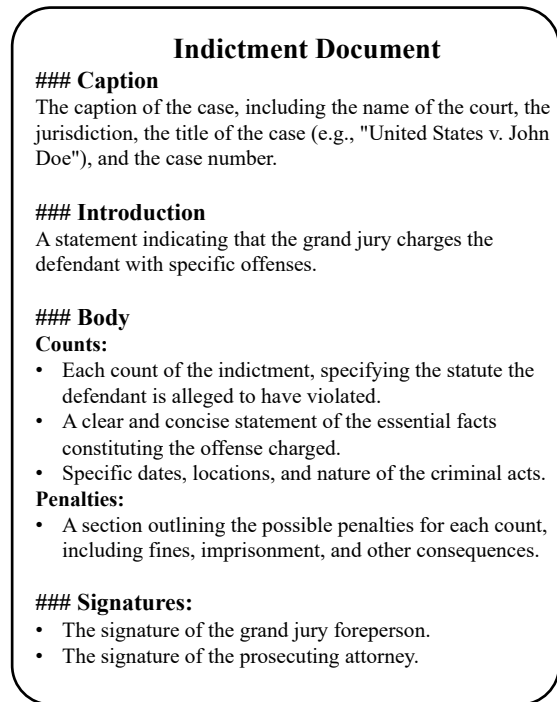


Figure 7: Illustration of the indictment document of US.

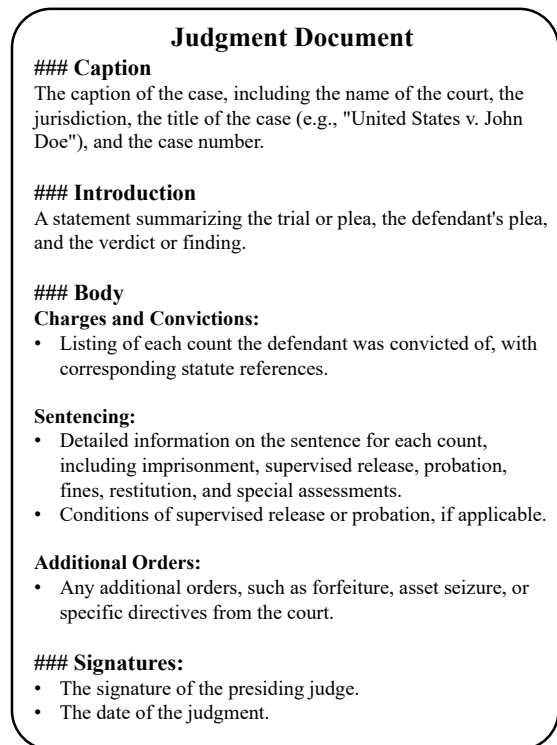


Figure 8: Illustration of the judgment document of US.

present in the body sections of these documents. our method can be applied to reformulate legal texts in documents from other jurisdictions similarly, thereby enhancing their performance of legal case retrieval.

904  
905  
906  
907  
908