

# Exploiting LLMs’ Reasoning Capability to Infer Implicit Concepts in Legal Information Retrieval

Hai-Long Nguyen<sup>1</sup>, Tan-Minh Nguyen<sup>1</sup>, Duc-Minh Nguyen<sup>2</sup>  
 Thi-Hai-Yen Vuong<sup>1</sup>, Ha-Thanh Nguyen<sup>3</sup>, Xuan-Hieu Phan<sup>1</sup>  
 Ken Satoh<sup>3</sup>

<sup>1</sup> VNU University of Engineering and Technology, Hanoi

<sup>2</sup> RMIT University Vietnam

<sup>3</sup> National Institute of Informatics, Japan  
 long.nh@vnu.edu.vn

## Abstract

Statutory law retrieval is a typical problem in legal language processing, that has various practical applications in law engineering. Modern deep learning-based retrieval methods have achieved significant results for this problem. However, retrieval systems relying on semantic and lexical correlations often exhibit limitations, particularly when handling queries that involve real-life scenarios, or use the vocabulary that is not specific to the legal domain. In this work, we focus on overcoming this weaknesses by utilizing the logical reasoning capabilities of large language models (LLMs) to identify relevant legal terms and facts related to the situation mentioned in the query. The proposed retrieval system integrates additional information from the term-based expansion and query reformulation to improve the retrieval accuracy. The experiments on COLIEE 2022 and COLIEE 2023 datasets show that extra knowledge from LLMs helps to improve the retrieval result of both lexical and semantic ranking models. The final ensemble retrieval system outperformed the highest results among all participating teams in the COLIEE 2022 and 2023 competitions.

## 1 Introduction

The statute law retrieval problem is a typical task in Legal NLP. The problem takes as input a natural language query, which could be a question, legal statement, or a specific scenario. The output of the problem is relevant legal articles or segments extracted from articles containing information that address the given query. In countries that follow the statute law system, this problem is the vital core of legal reference systems or legal search engines which could serve for many people, from legal experts to non-experts. These AI-based tools could help legal practitioners reduce the amount of time and resources on paperwork.

Previous studies have done well in using lexical features through statistical models combined with semantic features from Transformer-based models to address and improve retrieval results (Kim et al. 2022; Goebel et al. 2023; Louis and Spanakis 2022). Given the high number of candidates, the legal retrieval task is often split into two steps: lexical ranking and semantic re-ranking. However, a significant issue arises when queries containing content with no lexical overlap with the gold standard articles, which leads

to the gold standard articles being eliminated from the first retrieval phase. Moreover, different query types exist, including legal statements and specific scenarios as shown in Table 1. Legal statements are typically concise and may not require logical reasoning for comprehension. Specific scenarios describe conflicts of rights between parties, involving more complex underlying logical reasoning.

The emergence of Large Language Models (LLMs) has opened up a new era of development in the field of AI, especially NLP (Zhao et al. 2023; Yang et al. 2023; Bommasani et al. 2021). To date, LLMs have effectively addressed various general NLP tasks, including text summarization and machine translation (Huang et al. 2023; Yao et al. 2023; Laban et al. 2023). For legal domains, these models have shown their abilities to improve the performance of downstream tasks (Yu, Quartey, and Schilder 2022; Trautmann, Petrova, and Schilder 2022; Zhou, Huang, and Wu 2023). However, previous work relied on prompting techniques or the in-context learning capability of LLMs to tackle these tasks straightforwardly. In contrast, this work leverages the strength of LLMs to explore hidden logical reasoning of queries as additional information to retrieval models.

The proposed method is inspired by how legal experts search legal documents. When searching for legal documents that relate to a real-life situation, legal experts typically begin by finding which behaviours the subjects perform, which legal issues are described in the situation, and which legal concepts they relate to. Then, they search for legal documents related to those legal issues or concepts. Identifying relevant legal issues and concepts related to real-life situations is equivalent to the process of query expansion in terms of computer science. However, this process is not merely semantic matching but requires inference based on the interaction between entities in the event. Recent studies have shown that large language models such as GPT-4 (Achiam et al. 2023) and Gemini have basic inference capabilities, although they have not reached human-level inference (Achiam et al. 2023; Rane, Choudhary, and Rane 2024; Wang et al. 2024). This has motivated this research to utilize LLMs to explore hidden logical semantics of queries for query expansion.

The main contributions of this research include identify-

Table 1: An example of a specific scenario query in the COLIEE dataset.

	Japanese	English
<b>Query</b>	Aが自己所有の動産甲をBに賃貸し引き渡していた場合において、CがBのもとから甲を窃取したときは、Aは、Cに対して、占有回収の訴えによって甲の返還を求めることができる。	A has leased and delivered movable X owned by A to B. If C steals X from B, A may demand C to return X by an action for recovery of possession.
<b>LLMs generation</b>		
Legal terms	占有回収の訴え	Action for recovery of possession
Re-write query	このクエリに関連する法的概念は「占有回収訴訟」です。占有回収訴訟とは、動産を不法に占有する者に対して、その返還を求める訴訟です。本件では、CがBのもとから甲を窃取した場合、AはCに対して占有回収訴訟によって甲の返還を求めることができます。なぜなら、Aは甲の所有者であり、Cは甲を不法に占有しているからです。ただし、例外として、善意かつ無過失で動産を取得した者は、たとえその動産が盗品であったとしても、占有回収訴訟によって返還を求められることはありません（善意取得）。	The legal concept relevant to this query is “possession recovery litigation.” A possession recovery lawsuit is a lawsuit that seeks the return of movable property against a person who illegally possesses it. In this case, if C steals Party A from B, A can demand the return of Party A by filing a possession recovery lawsuit against C. This is because A is the owner of Party A, and C is in illegal possession of Party A. However, as an exception, a person who acquires movable property in good faith and without fault will not be required to return it through a possession recovery lawsuit, even if the movable property is stolen (good faith acquisition).
<b>Relevant articles</b>	第百八十一条 占有権は、代理人によって取得することができる。	Possessory rights may be acquired through an agent.

ing the current weaknesses of lexical and semantic feature-based matching models and proposing two query expansion methods using LLMs to extract the underlying topic of that query. The term generation and query-reformulation techniques are used and injected in both lexical and semantic ranking models. The proposed query expansion methods enable the retrieval system to effectively surpass the top-performing approaches of the COLIEE 2022 and 2023 datasets.

The paper will be structured as follows: the next section presents some related studies on the statute law retrieval problem and the application of LLMs to this problem, section 3 presents proposed method and experiments. Finally, the section 5 concludes the paper and highlights some future work.

## 2 Related Work

Legal NLP has attracted much interest from researchers and companies because of its potential and wide range of applications (Chalkidis and Kampas 2019; Zhong et al. 2020). Statutory article retrieval is one of the core problems in this field, playing an essential role in search engines and reference systems. A traditional approach for query-article matching is the term-based model (Kim, Rabelo, and Goebel 2019; Tran, Nguyen, and Satoh 2019). However, the trend rapidly changes to deep learning models involving word embedding (Landthaler et al. 2016; Kayalvizhi, Thenmozhi, and Aravindan 2019), document vector embedding (Sugathadasa et al. 2019), and contextual understanding (Nguyen et al. 2024; Nguyen et al. 2023).

Recently, a growing number of studies have focused on large language models (LLMs) in legal NLP (Sun 2023; Lai et al. 2023). These models can address downstream tasks such as question answering (Yu, Quartey, and Schilder 2022), legal judgment prediction (Trautmann, Petrova, and Schilder 2022), legal summarization (Pont et al. 2023; Gesnouin et al. 2024) based on fine-tuning and prompt engineering techniques. However, there is space for exploiting LLMs to address legal information retrieval. Zhou et al. (Zhou, Huang, and Wu 2023) utilized LLMs to extract significant content from legal cases and incorporate it into retrieval models. Unlike (Zhou, Huang, and Wu 2023), which involved a manual annotation process with the assistance of legal experts to collect salient content, this work does not rely on human-labeled data for LLMs. This distinction enhances the reproducibility and openness of the research.

Query expansion is one of the typical methods to shorten the gap between query and document, further improving retrieval performance. The expansion process mainly relies on external knowledge-based and pseudo-relevance feedback to enrich query information. Relevance documents can serve as additional knowledge in query expansion. However, this approach does not apply to our work as the relevant legal documents are not provided with the query. Therefore, we generate legal terms and a new query based on the original query to serve as additional knowledge for the retrieval system. Recently studies have leveraged generative models to rewrite the query via fine-tuning with labeled data (Imani et al. 2019; Zheng et al. 2020; Zheng et al. 2021). Large language models (LLMs) (Team et al. 2023; Touvron et al. 2023; Achiam et al. 2023) with billions of

parameters are trained on massive volumes of data. LLMs could address downstream tasks via prompt engineering or instructions. Wang et al. (Wang, Yang, and Wei 2023) utilized *gpt-3.5-turbo* to generate a pseudo document from a given query using few-shot prompting. Jagerman et al. (Jagerman et al. 2023) studied the performance of various prompting techniques in query expansion using Flan-family models. These works used common information retrieval datasets such as MS-MACRO (Nguyen et al. 2016) and BEIR (Thakur et al. 2021) as evaluation data, proving the effectiveness of LLMs-based query expansion. This paper focuses on inferring implicit concepts based on a given legal query. To the best of our knowledge, we are the first to exploit general LLMs as a query rewriter in the legal statute law retrieval problem.

### 3 Query expansion for legal document retrieval

As mentioned in section 1, to address complex queries with multiple layers of implicit semantics, reasoning capabilities are required. This section will describe the prompting patterns on Large Language Models (LLMs) to leverage their basic reasoning abilities and extract key information. Additionally, the legal document retrieval system described in section 3.3 will integrate ranking results from multiple models and is designed to entirely leverage information obtained from the prompting process.

#### 3.1 LLMs-based legal term extraction

Due to their generation capability and knowledge across various domains, LLMs are utilized to generate the key terms which represents the general topics of the query. Legal terms related to the query are extracted using zero-shot prompting techniques on LLMs. Since the experimented dataset is the COLIEE dataset with the original version written in Japanese, the prompting sentence will be written in Japanese to avoid information loss. Furthermore, to facilitate the processing of LLMs' output, the prompting sentence includes additional instructions for LLMs output in JSON format. The prompting pattern and instruction in English version are described in listing 1, the one with Japanese version is presented at the listing 3 in the Appendix.

Listing 1: Prompting pattern for legal term extraction (English version)

```

1 Given a legal situation, find the relevant facts and legal concepts
2 relevant to that situation: { query }
3
4 The output must be formatted as a JSON instance conforming to
5 the JSON schema below. For example, the schema
6 {
7   properties: {
8     foo: {
9       title: Foo,
10      description: a list of strings,
11      type: array,
12      items: {type: string}
13    },
14    required: [foo]
  }

```

To simplify the experiments and verify the impact of the expanded information, the collection of legal terms will be directly concatenated with the query at the lexical layer. Since this concatenation disrupts the semantic integrity at the end of the query, lexical-based retrieval models will be most appropriate for this query-expansion type. The BM25 (Robertson, Zaragoza, and others 2009) model which is one of the most commonly lexical-based ranking models is utilized. The process of using BM25 to rank this type of expanded queries will be described in section 3.3.

#### 3.2 LLMs-based query augmentation

Directly concatenating legal terms into the query does not ensure semantic coherence for the query, which may cause confusion for semantic-based ranking models. To preserve the semantic integrity, the legal-style oriented query reformulation is performed by zero-shot prompting technique on LLMs. The suggested prompting pattern written in English is detailed in listing 2, the Japanese one is showed at the listing 4.

Listing 2: Prompting pattern for redescribed query (English version)

```

1 Given a legal situation, extract the relevant facts and legal
2 concepts relevant to that situation: { query }

```

The set of reformulated queries will be used to train the semantic matching model, further supporting information about the legal speciality of the content mentioned in the original query. The specific method to leverage the prompting results from LLMs will be described in section 3.3.

#### 3.3 Retrieval system utilizing expanded query

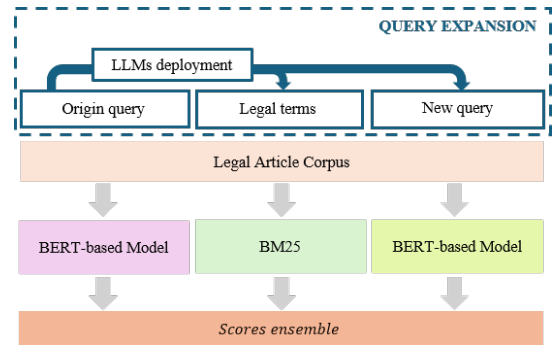


Figure 1: Query expansion supported retrieval system

To combine the information from both prompting methods described in section 3.1 and 3.2, both lexical-based ranking and semantic-based ranking models will be utilized and ensemble.

Assuming the input query is denoted as  $q$  and the legal corpus, denoted as set  $D$ , including  $m$  legal documents:  $D = \{d_1, d_2, \dots, d_m\}$ . After the legal-term extraction, the obtained terms form the set  $T_q = \{t_1, t_2, \dots, t_n\}$ . These terms are concatenated with the query at the lexical level to create a new query:

$$q_{\text{term-expand}} = \text{concat}(q, t_1, t_2, \dots, t_n)$$

The BM25 model is employed to calculate the lexical relevance between the expanded query  $q_{\text{term-expand}}$  and the legal document  $d_i$ . The relevance score calculated by the BM25 model is denoted as:

$$R_{\text{BM25}} = \text{BM25}(q_{\text{term-expand}}, d_i)$$

For the semantic-based ranking component, we utilized ranking models based on the BERT architecture, fine-tuned by sequence classification tasks. Experimenting with the simplest backbone model such as BERT will help assess the contribution of the additional information from the query expansion process. Two distinct BERT models are employed in this proposed retrieval system, one for ranking the original query and the other for the reformulated query. These two BERT models are described as follows:

- Let  $q_{\text{origin}}$  be the original query. The first semantic ranking model denoted as  $\text{BERT}_{\text{origin}}$ , evaluates the semantic relevance between the original query  $q_{\text{origin}}$  and the legal document  $d_i$ , denoted as:

$$R_{\text{ori}} = \text{BERT}_{\text{origin}}(q_{\text{origin}}, d_i)$$

- Let  $q_{\text{reformulated}}$  be the reformulated query, derived from LLMs. The second BERT model, denoted as  $\text{BERT}_{\text{reformulated}}$ , evaluates the semantic relevance between the reformulated query  $q_{\text{reformulated}}$  and the legal document  $d_i$ , denoted as:

$$R_{\text{reform}} = \text{BERT}_{\text{reformulated}}(q_{\text{reformulated}}, d_i)$$

The architecture of the retrieval system, which includes three ranking models, is illustrated in Figure 1. The final relevance scores from all three ranking models will be weighted ensembled using the equation 1.

$$R_{\text{final}} = \alpha * R_{\text{ori}} + \beta * R_{\text{bm25}} + \gamma * R_{\text{reform}} \quad (1)$$

$s.t : \alpha + \beta + \gamma = 1$

where:

- $R_{\text{ori}}$  represents the correlation score between the original query and the article as calculated by the BERT-based ranking model.
- $R_{\text{bm25}}$  is the correlation score between the query, which has been concatenated with legal term outputs from LLMs, and the article, as calculated by the BM25 ranking algorithm.
- $R_{\text{reform}}$  is the correlation score between the query that has been re-described using LLMs and the article, also calculated by the BERT-based ranking model.
- $\alpha$ ,  $\beta$ , and  $\gamma$  are the weights assigned to the BERT model with the original query, the BM25 model with the legal-term expanded query, and the BERT model with the re-described query, respectively. The optimal values of these three weights are selected through a grid-search process on a validation set.

### 3.4 Inference strategy

The legal document retrieval task is indeed about ranking the documents in the law corpus based on the relevancy

between each document and the input query. Meanwhile, the training process for the BERT model uses the objective function of the Sequence Classification downstream task, classifying the (*query-article*) pair into *positive* or *negative* labels. Therefore, directly using the relevance score from the BERT model often does not yield the highest effectiveness for the retrieval system. To address this drawback, after obtaining the  $R_{\text{final}}$  for each (*query-article*) pair, a post-processing phase will be conducted. In the post-processing phase, the relevance scores for all (*query-article*) pairs corresponding to the same query will be min-max normalized. Subsequently, a common optimal threshold will be determined through a grid-search process on the validation set. All articles with a relevance score exceeding this threshold will be selected as relevant to the query. Specifically, the post-processing procedure for producing the final set of relevant documents is described in Figure 2.

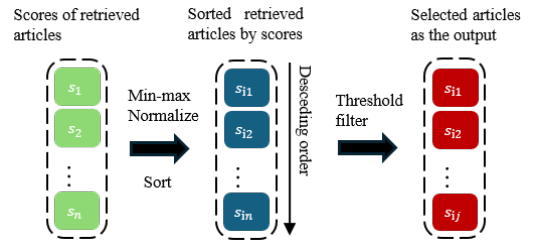


Figure 2: Retrieval task post-processing

## 4 Experiment and result

### 4.1 COLIEE dataset introduction

To evaluate the effectiveness of the proposed retrieval system, the COLIEE 2022<sup>1</sup> and COLIEE 2023<sup>2</sup> datasets for the Statute Law Retrieval task will be utilized. Both datasets employ a corpus derived from the Japanese statute law, containing 782 legal provisions. The labeled set of COLIEE

Table 2: Statute law corpus token statistics

	#word per articles	
	en	jp
<b>Min</b>	1	6
<b>Max</b>	655	755
<b>Average</b>	71.83	93.34

2022 and COLIEE 2023 respectively includes 887 and 996 instances. The content of queries in both datasets often describes a legal statement, legal situation or a real-life legal scenario. Table 2 summarizes the number of samples and the length of each query through the statistical indices *minimum*, *maximum*, and *average*.

<sup>1</sup><https://sites.ualberta.ca/~rabelo/COLIEE2022/>

<sup>2</sup><https://sites.ualberta.ca/~rabelo/COLIEE2023/>

Table 3: Token statistics

		COLIEE 2022		COLIEE 2023	
#sample	-	887	109	996	101
Min	EN	6	13	6	11
	JP	10	18	10	15
Max	EN	149	91	149	27
	JP	202	125	202	109
Average	EN	39	42	39	42
	JP	51	52	51	54

Table 4: Recall score of BM25 model in Japanese-version training dataset of COLIEE 2022 and COLIEE 2023 with original query and term-expanded query using two LLMs models: Gemini and GPT-4.

Datasets	Top-k	Origin	Gemini	GPT-4
COLIEE 2022	30	0.7590	0.8423	0.8394
	50	0.7878	0.8735	0.8663
	100	0.8394	0.9098	0.9061
	200	0.8836	0.9432	0.9366
	500	<b>0.9418</b>	<b>0.9782</b>	<b>0.9776</b>
COLIEE 2023	30	0.8320	0.8530	0.8516
	50	0.8592	0.8834	0.8767
	100	0.8955	0.9182	0.9139
	200	0.9316	0.9479	0.9421
	500	<b>0.9736</b>	<b>0.9790</b>	<b>0.9785</b>

## 4.2 Implementation detail

In this section, the implementation, training, inference, and ensemble of model results are discussed. Firstly, for the BM25 model, we utilize the rank-bm25 library<sup>3</sup>. Before being processed by the BM25 model, articles and queries written in Japanese are word-based tokenized using the Konoha library<sup>4</sup> and the MeCab tokenizer. The prompting process is performed on GPT-4 (Achiam et al. 2023) and Gemini (Team et al. 2023).

The training dataset provides sample which contains a query and a list of gold relevant articles. For training the BERT-based ranking model, the training samples are transformed into pairwise training samples. The transformed set contains a query and an article with two labels 0 or 1 representing not-relevant or relevant. The negative samples are chosen from top candidates according to the BM25’s relevance score. To achieve the optimal ratio between negative and positive samples, the top-30 most relevant candidates according to the BM25 model are utilized for creating negative samples for the queries. The Huggingface library<sup>5</sup> is employed for creating and training the BERT-based model with the downstream task sequence classification. The pre-trained Multilingual BERT<sup>6</sup> is being used for initializing the BERT-based retrieval model. The experimented mod-

<sup>3</sup>[https://github.com/dorianbrown/rank\\_bm25](https://github.com/dorianbrown/rank_bm25)

<sup>4</sup><https://github.com/himkt/konoha?tab=readme-ov-file>

<sup>5</sup><https://huggingface.co>

<sup>6</sup><https://huggingface.co/google-bert/bert-base-multilingual-cased>

els are fine-tuned with the training data in 3 epochs. There are two ranking BERT-based models in the proposed retrieval system which are the ranking model fine-tuned with the origin query and the ranking model fine-tuned with the re-described query. The results in section 4.3 will demonstrate that using two separate BERT-based ranking models for original and re-described queries is more effective than a single fine-tuned BERT model handling both.

Table 5: Experiment results on private test of COLIEE 2022 dataset.

Teams/ Models	Private Test		
	F2	P	R
HUKB	0.8200	<b>0.8180</b>	0.8405
OVGU	0.7790	0.7781	0.8054
JNLP	0.7699	0.6865	0.8378
UA	0.7638	0.8073	0.7641
<b>Proposed Retrieval System</b>	<b>0.8449</b>	0.7935	<b>0.9005</b>

Table 6: Experiment results on private test of COLIEE 2023 dataset.

Teams/ Models	Private Test		
	F2	P	R
CAPTAIN	0.757	0.726	0.792
JNLP	0.745	0.645	<b>0.867</b>
NOWJ	0.727	0.682	0.772
HUKB	0.6473	0.6248	0.6708
<b>Proposed Retrieval System</b>	<b>0.7601</b>	<b>0.7254</b>	0.8366

Before the training phase, a validation set is separated from 20% of the labelled set, and the remaining 80% of samples are formed into the training set. After the training the retrieval model, a grid-search process is conducted on the validation set for choosing the best  $\alpha$ ,  $\beta$ ,  $\gamma$  and *threshold*.

## 4.3 Experiment Result

The recall score of two term-expansion methods are presented in the table 4. According to that table, the recall scores when using the term-expanded queries instead of the original one are significantly higher at all top-k levels, with both Gemini and GPT-4 terms-expansion. This proves that adding legal terms will make the queries more lexically relevant to the gold standard articles, facilitating the BM25 model to rank more accurately. Additionally, the recall score of the Gemini-based term-expansion achieved a higher recall score than the GPT-4. Therefore, the BM25’s relevance score with Gemini’s expanded queries is used in the final retrieval system.

According to the implementation detail, the retrieval system is experimented with the COLIEE 2022 and 2023 datasets. Because the BERT’s training sample pair is generated by filtering from BM25’s candidates, the relevance score derived from BERT-based ranking model need to be ensembled with the relevance score from BM25 model to

Table 7: Evaluation of some retrieval system’s variations. The BM25<sub>legalterm</sub> is the BM25 model worked with the term-expanded query instead of original query.

Retrieval system’s variations	2022 Private Test			2023 Private Test		
	F2	P	R	F2	P	R
BM25 <sub>legalterm</sub> + BERT <sub>origin</sub>	0.8305	<b>0.8107</b>	0.8680	0.7356	0.7227	0.7871
BM25 <sub>legalterm</sub> + BERT <sub>reform</sub>	0.8201	0.7928	0.8515	0.7050	0.6787	0.7574
BM25 <sub>legalterm</sub> + BERT <sub>origin</sub> + BERT <sub>reform</sub>	<b>0.8449</b>	0.7953	<b>0.9005</b>	<b>0.7643</b>	<b>0.7478</b>	<b>0.8218</b>

achieve full potential. The table 7 shows the results of three variants of retrieval system including:

- (1) The system that ensemble BM25 model with term-expanded queries and the BERT model with original queries.
- (2) The system that ensemble BM25 model with term-expanded queries and the BERT model with reformulated queries.
- (3) The system that ensemble BM25 model with term-expanded queries, BERT model with original queries and another BERT model with reformulated queries.

The retrieval performance of the system (2) is lower than system (1). This decrease occurs because using LLMs to reformulated queries make the query’s distribution be changed significantly. However, the combination of three models achieved the highest F2-score which is 0.8449 for the COLIEE 2022 dataset and 0.7643 for the COLIEE 2023 dataset. This observation suggests that despite having different distribution with the original query, the reformulated query contributes different-perspective information for the retrieval system.

The table 5 (COLIEE 2022) and table 6 (COLIEE 2023) show the F2, precision and recall of the retrieval system (3) and results of all teams that participated the competition that years. As indicated on the table, the proposed retrieval system achieved an F2 score that was 2.49% higher than the best-performing team in COLIEE 2022, and 0.3% higher than CAPTAIN team (best performance team) in COLIEE 2023. This result demonstrates the effectiveness of retrieval system which utilizes legal-specific knowledge extracted from LLMs and the way for integrating the knowledge into the retrieval system.

## 5 Conclusion

In this study, two query expansion methods including legal-term extraction and query reformulation are used to exploiting the reasoning capability and general knowledge from LLMs. The information extracted from these two prompting method are injected into lexical-based and semantic-based ranking models. The experiments in COLIEE 2022 and 2023 dataset show that the generated legal information from LLMs significantly improves retrieval results and outperforms the results of the teams participating in the COLIEE competition in 2022 and 2023.

## Acknowledgments

Hai-Long Nguyen was funded by the Master, PhD Scholarship Programme of Vingroup Innovation Foundation

(VINIF), code VINIF.2023.ThS.075.

## Appendix

Due to the original data is written in Japanese, the prompting pattern used will be in Japanese. Listing 3 describes the prompting pattern of the legal term extraction process and listing 4 shows the prompting pattern for generating the re-described query.

Listing 3: Prompting pattern for legal term extraction (Japanese version)

```

1 法的な状況が与えられた場合、その状況に関連す
2 る適切な事実と法的概念を抽出します: {query}
3
4 出力は、以下のJSONスキーマに準拠するJSONイン
5 スタンスとしてフォーマットする必要があります
6 す。例として、スキーマ
7 {
8   properties: {
9     foo: {
10       title: Foo,
11       description: a list of strings,
12       type: array,
13       items: {type: string}
14     }},
15   required: [foo]
16 }
17 の場合、オブジェクト {foo: [bar, baz]}
```

Listing 4: Prompting pattern for redescription query (Japanese version)

```

1 法的な状況が与えられた場合、その状況に関連す
2 る適切な事実と法的概念を抽出します: {query}
```

## References

- Achiam, J.; Adler, S.; Agarwal, S.; Ahmad, L.; Akkaya, I.; Aleman, F. L.; Almeida, D.; Altenschmidt, J.; Altman, S.; Anadkat, S.; et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Bommasani, R.; Hudson, D. A.; Adeli, E.; Altman, R.; Arora, S.; von Arx, S.; Bernstein, M. S.; Bohg, J.; Bosselut, A.; Brunskill, E.; et al. 2021. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*.
- Chalkidis, I., and Kampas, D. 2019. Deep learning in law: early adaptation and legal word embeddings trained on large corpora. *Artificial Intelligence and Law* 27(2):171–198.



- Gesnoui, J.; Tannier, Y.; Da Silva, C. G.; Tapory, H.; Brier, C.; Simon, H.; Rozenberg, R.; Woehrel, H.; Yakaabi, M. E.; Binder, T.; et al. 2024. Llamandement: Large language models for summarization of french legislative proposals. *arXiv preprint arXiv:2401.16182*.
- Goebel, R.; Kano, Y.; Kim, M.-Y.; Rabelo, J.; Satoh, K.; and Yoshioka, M. 2023. Summary of the competition on legal information, extraction/entailment (coliee) 2023. In *Proceedings of the Nineteenth International Conference on Artificial Intelligence and Law*, 472–480.
- Huang, H.; Wu, S.; Liang, X.; Wang, B.; Shi, Y.; Wu, P.; Yang, M.; and Zhao, T. 2023. Towards making the most of llm for translation quality estimation. In *CCF International Conference on Natural Language Processing and Chinese Computing*, 375–386. Springer.
- Imani, A.; Vakili, A.; Montazer, A.; and Shakery, A. 2019. Deep neural networks for query expansion using word embeddings. In *Advances in Information Retrieval: 41st European Conference on IR Research, ECIR 2019, Cologne, Germany, April 14–18, 2019, Proceedings, Part II 41*, 203–210. Springer.
- Jagerman, R.; Zhuang, H.; Qin, Z.; Wang, X.; and Bender-sky, M. 2023. Query expansion by prompting large language models. *arXiv preprint arXiv:2305.03653*.
- Kayalvizhi, S.; Thenmozhi, D.; and Aravindan, C. 2019. Legal assistance using word embeddings. In *FIRE (Working Notes)*, 36–39.
- Kim, M.-Y.; Rabelo, J.; Goebel, R.; Yoshioka, M.; Kano, Y.; and Satoh, K. 2022. Coliee 2022 summary: Methods for legal document retrieval and entailment. In *JSAI International Symposium on Artificial Intelligence*, 51–67. Springer.
- Kim, M.-Y.; Rabelo, J.; and Goebel, R. 2019. Statute law information retrieval and entailment. In *Proceedings of the Seventeenth International Conference on Artificial Intelligence and Law, ICAIL '19*, 283–289. New York, NY, USA: Association for Computing Machinery.
- Laban, P.; Kryściński, W.; Agarwal, D.; Fabbri, A. R.; Xiong, C.; Joty, S.; and Wu, C.-S. 2023. Summedits: Measuring llm ability at factual reasoning through the lens of summarization. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, 9662–9676.
- Lai, J.; Gan, W.; Wu, J.; Qi, Z.; and Yu, P. S. 2023. Large language models in law: A survey. *arXiv preprint arXiv:2312.03718*.
- Landthaler, J.; Walzl, B.; Holl, P.; and Matthes, F. 2016. Extending full text search for legal document collections using word embeddings. In *JURIX*, 73–82.
- Louis, A., and Spanakis, G. 2022. A statutory article retrieval dataset in french. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 6789–6803.
- Nguyen, T.; Rosenberg, M.; Song, X.; Gao, J.; Tiwary, S.; Majumder, R.; and Deng, L. 2016. MS MARCO: A human generated machine reading comprehension dataset. In Besold, T. R.; Bordes, A.; d’Avila Garcez, A. S.; and Wayne, G., eds., *Proceedings of the Workshop on Cognitive Computation: Integrating neural and symbolic approaches 2016 co-located with the 30th Annual Conference on Neural Information Processing Systems (NIPS 2016), Barcelona, Spain, December 9, 2016*, volume 1773 of *CEUR Workshop Proceedings*. CEUR-WS.org.
- Nguyen, T.-M.; Nguyen, X.-H.; Mai, N.-D.; Hoang, M.-Q.; Nguyen, V.-H.; Nguyen, H.-V.; Nguyen, H.-T.; and Vuong, T.-H.-Y. 2023. Nowj1@ alqac 2023: Enhancing legal task performance with classic statistical models and pre-trained language models. *arXiv preprint arXiv:2309.09070*.
- Nguyen, C.; Nguyen, P.; Tran, T.; Nguyen, D.; Trieu, A.; Pham, T.; Dang, A.; and Nguyen, L.-M. 2024. Captain at coliee 2023: Efficient methods for legal information retrieval and entailment tasks. *arXiv preprint arXiv:2401.03551*.
- Pont, T. D.; Galli, F.; Loreggia, A.; Pisano, G.; Rovatti, R.; and Sartor, G. 2023. Legal summarisation through llms: The prodigit project. *arXiv preprint arXiv:2308.04416*.
- Rane, N.; Choudhary, S.; and Rane, J. 2024. Gemini versus chatgpt: Applications, performance, architecture, capabilities, and implementation. *Performance, Architecture, Capabilities, and Implementation (February 13, 2024)*.
- Robertson, S.; Zaragoza, H.; et al. 2009. The probabilistic relevance framework: Bm25 and beyond. *Foundations and Trends® in Information Retrieval* 3(4):333–389.
- Sugathadasa, K.; Ayesha, B.; de Silva, N.; Perera, A. S.; Jayawardana, V.; Lakmal, D.; and Perera, M. 2019. Legal document retrieval using document vector embeddings and deep learning. In *Intelligent Computing: Proceedings of the 2018 Computing Conference, Volume 2*, 160–175. Springer.
- Sun, Z. 2023. A short survey of viewing large language models in legal aspect. *arXiv preprint arXiv:2303.09136*.
- Team, G.; Anil, R.; Borgeaud, S.; Wu, Y.; Alayrac, J.-B.; Yu, J.; Soricut, R.; Schalkwyk, J.; Dai, A. M.; Hauth, A.; et al. 2023. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*.
- Thakur, N.; Reimers, N.; Rücklé, A.; Srivastava, A.; and Gurevych, I. 2021. BEIR: A heterogeneous benchmark for zero-shot evaluation of information retrieval models. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*.
- Touvron, H.; Lavril, T.; Izacard, G.; Martinet, X.; Lachaux, M.-A.; Lacroix, T.; Rozière, B.; Goyal, N.; Hambro, E.; Azhar, F.; et al. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.
- Tran, V.; Nguyen, M. L.; and Satoh, K. 2019. Building legal case retrieval systems with lexical matching and summarization using a pre-trained phrase scoring model. In *Proceedings of the Seventeenth International Conference on Artificial Intelligence and Law, ICAIL '19*, 275–282. New York, NY, USA: Association for Computing Machinery.
- Trautmann, D.; Petrova, A.; and Schilder, F. 2022. Legal prompt engineering for multilingual legal judgement prediction. *arXiv preprint arXiv:2212.02199*.

- Wang, S.; Wei, Z.; Choi, Y.; and Ren, X. 2024. Can llms reason with rules? logic scaffolding for stress-testing and improving llms. *arXiv preprint arXiv:2402.11442*.
- Wang, L.; Yang, N.; and Wei, F. 2023. Query2doc: Query expansion with large language models. In *The 2023 Conference on Empirical Methods in Natural Language Processing*.
- Yang, J.; Jin, H.; Tang, R.; Han, X.; Feng, Q.; Jiang, H.; Zhong, S.; Yin, B.; and Hu, X. 2023. Harnessing the power of llms in practice: A survey on chatgpt and beyond. *ACM Transactions on Knowledge Discovery from Data*.
- Yao, B.; Jiang, M.; Yang, D.; and Hu, J. 2023. Empowering llm-based machine translation with cultural awareness. *arXiv preprint arXiv:2305.14328*.
- Yu, F.; Quartey, L.; and Schilder, F. 2022. Legal prompting: Teaching a language model to think like a lawyer. *arXiv preprint arXiv:2212.01326*.
- Zhao, W. X.; Zhou, K.; Li, J.; Tang, T.; Wang, X.; Hou, Y.; Min, Y.; Zhang, B.; Zhang, J.; Dong, Z.; et al. 2023. A survey of large language models. *arXiv preprint arXiv:2303.18223*.
- Zheng, Z.; Hui, K.; He, B.; Han, X.; Sun, L.; and Yates, A. 2020. Bert-qe: contextualized query expansion for document re-ranking. *arXiv preprint arXiv:2009.07258*.
- Zheng, Z.; Hui, K.; He, B.; Han, X.; Sun, L.; and Yates, A. 2021. Contextualized query expansion via unsupervised chunk selection for text retrieval. *Inf. Process. Manage.* 58(5).
- Zhong, H.; Xiao, C.; Tu, C.; Zhang, T.; Liu, Z.; and Sun, M. 2020. How does nlp benefit legal system: A summary of legal artificial intelligence. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 5218–5230.
- Zhou, Y.; Huang, H.; and Wu, Z. 2023. Boosting legal case retrieval by query content selection with large language models. In *Proceedings of the Annual International ACM SIGIR Conference on Research and Development in Information Retrieval in the Asia Pacific Region*, 176–184.