

LexRAG: Benchmarking Retrieval-Augmented Generation in Multi-Turn Legal Consultation Conversation

Haitao Li*

DCST, Tsinghua University &
Quan Cheng Laboratory
Beijing, China
liht22@mails.tsinghua.edu.cn

Yifan Chen*

DCST, Beijing University of Posts and
Telecommunications
Beijing, China
chenyifan@bupt.edu.cn

Yiran Hu*

Tsinghua University
Beijing, China
huyr21@mails.tsinghua.edu.cn

Qingyao Ai

Quan Cheng Laboratory &
DCST, Tsinghua University
Beijing, China
aiqy@tsinghua.edu.cn

Junjie Chen

Quan Cheng Laboratory &
DCST, Tsinghua University
Beijing, China
chenjj826@gmail.com

Xiaoyu Yang

Tsinghua University
Beijing, China
y15011462822@163.com

Jianhui Yang

Tsinghua University
Beijing, China
yangjh23@mails.tsinghua.edu.cn

Yueyue Wu[†]

Quan Cheng Laboratory &
DCST, Tsinghua University
Beijing, China
wuyueyue@mail.tsinghua.edu.cn

Zeyang Liu[‡]

Quan Cheng Laboratory &
Shandong University
Beijing, China
zeyangliu@sdu.edu.cn

Yiqun Liu

Quan Cheng Laboratory &
DCST, Tsinghua University
Beijing, China
yiqunliu@tsinghua.edu.cn

ABSTRACT

Retrieval-augmented generation (RAG) has proven highly effective in improving large language models (LLMs) across various domains. However, there is no benchmark specifically designed to assess the effectiveness of RAG in the legal domain, which restricts progress in this area. To fill this gap, we propose LexRAG, the first benchmark to evaluate RAG systems for multi-turn legal consultations. LexRAG consists of 1,013 multi-turn dialogue samples and 17,228 candidate legal articles. Each sample is annotated by legal experts and consists of five rounds of progressive questioning. LexRAG includes two key tasks: (1) Conversational knowledge retrieval, requiring accurate retrieval of relevant legal articles based on multi-turn context. (2) Response generation, focusing on producing legally sound answers. To ensure reliable reproducibility, we develop LexiT, a legal RAG

toolkit that provides a comprehensive implementation of RAG system components tailored for the legal domain. Additionally, we introduce an LLM-as-a-judge evaluation pipeline to enable detailed and effective assessment. Through experimental analysis of various LLMs and retrieval methods, we reveal the key limitations of existing RAG systems in handling legal consultation conversations. LexRAG establishes a new benchmark for the practical application of RAG systems in the legal domain, with its code and data available at <https://github.com/CSHaitao/LexRAG>.

ACM Reference Format:

Haitao Li, Yifan Chen*, Yiran Hu*, Qingyao Ai, Junjie Chen, Xiaoyu Yang, Jianhui Yang, Yueyue Wu, Zeyang Liu, and Yiqun Liu. 2025. LexRAG: Benchmarking Retrieval-Augmented Generation in Multi-Turn Legal Consultation Conversation. In *Proceedings of ACM Conference (Conference'17)*. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

*These authors contributed equally to this work.

[†]Corresponding author

[‡]Corresponding author

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

Conference'17, July 2017, Washington, DC, USA

© 2025 Association for Computing Machinery.

ACM ISBN 978-x-xxxx-xxxx-x/YY/MM... \$15.00

<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 INTRODUCTION

Recently, Retrieval-augmented generation (RAG) has gained significant attention as a powerful approach to improving the performance of large language models (LLMs). By integrating the strengths of information retrieval with generative models, RAG enables the generation of more accurate, relevant, and contextually appropriate responses based on documents retrieved from up-to-date, reliable sources. While RAG has demonstrated success in various domains, its application in legal domain remains underexplored.

Compared to general domains, RAG faces greater challenges in the legal domain. First, legal consultations are more complex,

often involving progressively unfolding issues. The users posing questions typically lack sufficient legal knowledge, requiring clarification, confirmation, and correction of details through multiple turns of dialogue. RAG systems must handle irrelevant information from previous interactions and effectively manage abrupt topic shifts. Moreover, in each turn, the relevance of a question to legal knowledge is not simply determined by lexical or semantic similarity [18, 21]. The model needs to consider the context for reasoning, identifying the legal logic and focus of the question to determine the relevant knowledge.

Although some benchmarks have been created to evaluate LLMs in the legal domain, they typically focus on simple tasks, such as legal case retrieval [25, 27] and judgment prediction [39], failing to capture the complexity that RAG faces in real-world legal scenarios. To fill this gap, we introduce LexRAG, a benchmark designed for RAG in multi-turn legal consultation conversations. It consists of 1,013 multi-turn consultation samples and includes 17,228 candidate articles. Each sample comprises five rounds of questions, with responses annotated by legal experts. In each conversation, the LLM must effectively incorporate previous turns and resolve pronoun references to understand the current query and ensure logical consistency. Additionally, LLMs need to handle abrupt topic shifts, which increase complexity and can degrade retrieval and generation quality as the dialogue history grows.

In LexRAG, we evaluate two key tasks of RAG systems: (1) Conversational Knowledge Retrieval, which assesses the system’s ability to retrieve relevant information from a large document corpus based on multi-turn context. (2) Response Generation, which tests its ability to generate accurate, contextually rich answers. To enable reproducible automated evaluation, we provide an easy-to-use toolkit LexiT, that includes the complete implementation of components for RAG systems in the legal domain. Moreover, we have carefully designed an LLM-as-a-judge evaluation pipeline within the toolkit to enable effective, fine-grained assessment. We conduct a comprehensive evaluation of various LLMs and retrieval methods, offering an in-depth analysis of the current limitations and shortcomings of RAG systems in the legal domain. Our findings highlight key challenges and suggest future directions for advancing RAG in the legal domain.

In summary, our contributions are three-fold:

- (1) **First Benchmark for RAG system in Legal Domain.** To the best of our knowledge, LexRAG is the first benchmark specifically designed to evaluate RAG in the legal domain. This benchmark provides a standardized platform for evaluating retrieval and generation capabilities in complex legal consultation conversations. It not only advances legal AI technologies but also lays the foundation for the future development of RAG across various domains.
- (2) **Open-Source Legal RAG Evaluation Toolkit.** In addition to the dataset, we provide LexiT, a dedicated toolkit for RAG in the legal domain. This toolkit includes various implementations of modules such as processors, retrievers, and generators for RAG systems. Additionally, we have carefully designed the LLM-as-a-judge evaluation framework to enable effective and fine-grained automated assessment.

This contributes to the advancement of research in the legal domain by enabling consistent and comparable evaluations.

- (3) **Systematic Evaluation and Analysis.** Through rigorous evaluation of several LLMs and retrieval methods, we analyze the strengths and limitations of current RAG systems in the legal domain. These observations offer valuable insights and highlight areas for further improvement, providing a roadmap for enhancing RAG-based legal consultation systems in the future.

2 RELATED WORK

2.1 Legal Applications of LLMs

Large language models (LLMs) have demonstrated strong potential for application in the legal domain [19, 24]. Several studies have thoroughly reviewed the current applications of LLMs in the legal domain, highlighting their vast potential in areas like legal consultation and trial assistance [17]. Additionally, researchers have developed LLMs specifically for the legal domain through continued pretraining and fine-tuning [42]. For example, ChatLaw [9] is built on the Anima-33B model and fine-tuned with a large dataset that includes legal news, statutes, judicial interpretations, legal consultations, exam questions, and court judgments. Meanwhile, LexiLaw [20] is further trained on ChatGLM to offer accurate and reliable legal consultation for legal professionals, students, and the general public. It excels in interpreting legal clauses, analyzing cases, and understanding regulations.

2.2 Retrieval-Augmented Generation

Naive LLMs can suffer from hallucinations or provide outdated answers when handling domain-specific tasks or recent information [33, 40]. RAG addresses this issue by first retrieving relevant information from external knowledge sources, improving the LLM’s accuracy and ensuring timely, up-to-date responses [2, 10]. The RAG workflow consists of three steps. First, the retriever fetches relevant information from an external knowledge base. Next, the retrieved information is combined with the original query to create an augmented prompt. Then, the generator produces the response based on the augmented prompt. In recent years, RAG has been widely applied across various fields. For example, in question-answering systems, LLMs enhance their ability to handle complex queries and generate more accurate responses by integrating a retrieval mechanism.

2.3 Multi-Turn Conversation

Dialogue systems are designed to facilitate continuous communication between humans and machines by understanding context and generating coherent responses [6, 30, 46]. These systems are typically classified into task-oriented [12] and open-domain systems [13]. Task-oriented systems help users complete specific tasks, such as booking hotels or checking the weather, while open-domain systems engage users in conversations on a wide variety of topics. The main challenge for these systems is generating coherent and diverse responses to maintain a natural and smooth conversation. With the development of deep learning and pre-trained models, LLM-based multi-turn dialogue systems have shown excellent performance [41]. These models, pre-trained on large corpora, acquire

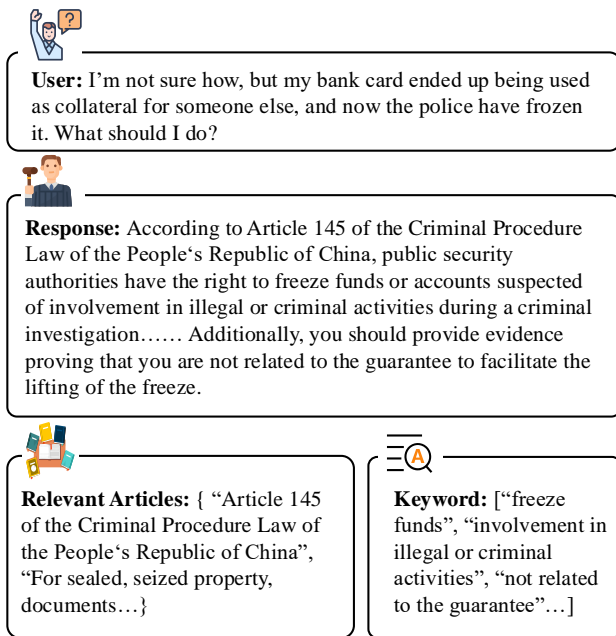


Figure 1: An example of a legal consultation in LexRAG.

rich linguistic and world knowledge, enabling them to generate more natural and contextually relevant responses. However, LLMs still face challenges in multi-turn dialogues, such as context understanding, dialogue state tracking, reasoning and planning, and response consistency.

2.4 Benchmarks in Legal Domain

In the legal domain, evaluation benchmarks are crucial for the development of LLMs. Evaluation benchmarks for LLMs are essential for assessing their performance on legal tasks [24, 44]. Researchers have developed benchmarks to evaluate LLMs' capabilities in tasks like legal reasoning, text classification, and question answering. For example, LexGLUE [5] is an English-language benchmark that standardizes the evaluation of models across various legal NLP tasks, such as text classification and case judgment prediction. LexEval introduces a legal cognition taxonomy and organizes 14,150 tasks to systematically evaluate LLMs' abilities in the legal domain. Moreover, Li et al. [23] introduced LegalAgentBench, which evaluates LLM agents specifically in the legal domain.

3 LEXRAG

In this section, we provide a detailed overview of LexRAG, including task definition, characteristics, data construction, and RAG toolkit.

3.1 Overview

LexRAG is the first benchmark designed to evaluate the performance of RAG in the legal domain, covering two key tasks: conversational knowledge retrieval and response generation. The dataset contains 1,013 multi-turn conversations, each with 5 rounds of questions and responses. Each conversation is carefully annotated

by legal experts to ensure accuracy and professionalism. Additionally, it includes 17,228 candidate legal articles across various legal domains, such as civil, criminal, contract, and intellectual property law. Figure 1 illustrates an example of legal consultations in LexRAG. In addition to the questions and responses, legal experts also identify and annotate relevant legal articles and keywords within the responses. Given that legal terminology has precise meanings, we also use the accuracy of keyword as an evaluation metric. LexRAG provides a standardized evaluation platform to advance RAG applications in the legal field and support the development of high-quality legal consultation systems.

3.2 Task Definition

LexRAG is designed to evaluate two fundamental tasks: (1) **Conversational Knowledge Retrieval** and (2) **Response Generation**. Compared to general domains, both tasks present unique challenges inherent to the legal domain.

For conversational knowledge retrieval task, the RAG system must identify legal articles relevant to the current query while considering the context. Formally, given a multi-turn legal dialogue history $H = \{q_1, r_1, \dots, q_t\}$, where q_t represents the user's question and r_t is the response at turn t . The objective is to retrieve a set of relevant legal articles $A_t = \{a_1, \dots, a_n\}$ from a predefined legal corpus \mathcal{D} . The retrieved articles should provide authoritative references for answering q_t . Unlike web search tasks that primarily rely on keyword matching or semantic similarity, this task in the legal domain introduces additional complexities. The retrieval models must not only understand the explicit query but also deduce the implicit legal intent behind it. For example, a user might ask a seemingly simple everyday question, but the final answer may involve referencing a complex series of interrelated statutes. Therefore, the retrieval system must go beyond simple keyword matching and instead focus on a nuanced understanding of legal concepts and relationships.

For Response Generation task, the LLM needs to generate contextually coherent and legally accurate response a_t based on the dialogue history H and retrieved legal articles A_t . In addition to the inherent challenges of multi-turn dialogues, such as anaphora resolution, context dependency, and topic shifts, this task requires LLMs to accurately interpret the legal requirements embedded in the query and apply the retrieved legal information precisely. In summary, LexRAG requires a deep integration of legal knowledge, multi-turn dialogue management, and advanced retrieval mechanisms.

3.3 Characteristics

LexRAG is designed as a comprehensive and reliable benchmark with the following key characteristics:

Legal Expertise. All responses in LexRAG are carefully annotated and reviewed by experienced legal experts to ensure accuracy and reliability. Additionally, the seed questions are sourced from legal consultation platforms, reflecting real-world legal practices.

Multi-Turn. In LexRAG, each conversation consists of five interactive turns. User queries often involve anaphora resolution, clarification, and topic shifts. This requires the system to effectively track conversation history and adapt to the evolving legal context.

Diversity. LexRAG covers a broad range of real-world legal issues, including 27 query types such as traffic accidents, personal injury, and debt disputes. The retrieval corpus includes 17,728 legal provisions from 222 statutes and regulations, ensuring comprehensive legal coverage for thorough evaluation.

Citation-Based Grounding A key feature of LexRAG is its focus on legal citation. Most responses explicitly reference legal articles, ensuring alignment with authoritative sources. This approach enhances transparency, verifiability, and highlights the importance of accurate knowledge retrieval in legal consultation.

3.4 Data Construction

In this section, we introduce the construction process of LexRAG, including data sources, preprocessing, human annotation, and data analysis.

3.4.1 Data Source and Preprocessing. To construct LexRAG, we collected 222 commonly used legal statutes in China, ensuring each was from its latest version. We standardized the formatting of legal provisions and created a structured retrieval corpus with 17,228 legal articles. Then, we collect seed questions to guide human annotators in structuring and annotating the conversation. These questions are sourced from real-world legal consultation platforms¹ to ensure relevance and authenticity. We thoroughly review and exclude queries containing personal information, sensitive content, or legally irrelevant inquiries.

3.4.2 Human Annotation. Our annotation team consists of 11 legal experts from China, all of whom have passed the Chinese Judicial Examination and possess extensive legal experience. The team includes six males and five females. Before starting the annotation process, we signed legally binding agreements with all members to ensure compliance with legal standards and protect their rights.

Training. To ensure dataset quality, we provided systematic training for all legal experts before annotation. We developed a comprehensive annotation guideline, clearly defining the annotation standards and procedures. Additionally, we provided 10 examples to facilitate a better understanding of the annotation requirements. Each annotator was required to complete 10 pilot tasks and receive feedback and guidance from senior legal experts, who are the creators of the annotation guidelines. Only those who achieved a pass rate above 90% were permitted to proceed to the formal annotation phase.

Annotation. The annotation process begins with an initial seed question. In the subsequent turns, annotators are encouraged to naturally expand the conversation, ensuring that new questions logically follow the existing conversation threads.

To support the annotation process, we provide a convenient annotation toolkit to annotators. This toolkit uses the BM25 [32] algorithm to retrieve 30 legal articles relevant to the current question from the legal corpus, providing annotators with valuable references. Additionally, annotators have direct access to the full legal corpus, allowing them to manually select the most relevant legal articles for each question. Then, annotators must provide detailed responses based on their legal expertise. They are also required to highlight keywords in their responses and annotate them with

the corresponding legal articles for review and analysis. To reduce the annotation workload, we use GPT-4o-mini to pre-generate 10 rounds of derivative questions from the initial seed question, covering different perspectives. These generated questions serve as examples, providing inspiration for annotators. To ensure the diversity and originality of the dataset, direct copying is strictly prohibited.

Review. We implemented a thorough review process to ensure the quality and reliability of the annotated data. Our gold annotators, who created the annotation guidelines, performed cross-validation of each annotation from multiple perspectives. Specifically, they evaluated whether the questions were logically coherent and legally valid, whether the responses were accurate and aligned with legal principles, whether the cited legal articles were relevant and correctly referenced, and whether key terms were appropriately annotated. Any annotations that did not meet the required standards were reviewed by a senior legal expert to ensure they followed legal standards and best practices. If any issues were found, the data point was sent back for revision and clarification. This process continued until both annotators agreed. Only high-quality annotations were included in the final dataset. To fairly reward annotators for their expertise, we paid \$0.42 per validated question-response pair. With 5,065 dialogues created, the total payment amounted to \$2,110.

3.4.3 Annotation Guideline. To ensure the quality, consistency, and reliability of LexRAG, we have implemented a rigorous validation and annotation process based on the following principles and standards. Specifically, the annotators follow the annotation pattern of “parsing the question–identifying relevant legal articles–generating answers–formulate new questions–simulate real-life scenarios”.

- **Parsing the Question.** The seed questions in LexRAG are sourced from real-world legal consultation platforms, meaning they often focus on real-life issues rather than legal facts. As a result, directly answering the questions may lead to inaccurate responses or a failure to capture the true intent behind the query. To address this, annotators first parse the real-life issues into key legal terms. For example, if a user asks, “My girlfriend was already pregnant when we were together! Can the child be registered in the household?” The annotator can derive legal terms such as “household registration”, “birth certificate” and “child out of wedlock” based on their legal knowledge. These terms are then used to guide the retrieval of relevant legal articles.
- **Identifying Relevant Legal Articles.** Based on the legal terms derived in the previous step, annotators can use our provided retrieval toolkit or keyword matching to identify relevant legal articles from the candidate database. For example, for the question above, the most relevant provision is Article 7 of the “Household Registration Regulations of the People’s Republic of China”.
- **Generating Answers.** Based on the legal logic of syllogism, annotators are encouraged to respond by referencing relevant legal articles. For example, for the question above, the generated response is: “According to Article 7 of the ‘Household Registration Regulations of the People’s Republic of

¹<https://www.12348.gov.cn/>

Table 1: Basic statistic of LexRAG.

Statistic	#Number
Total Conversations	1,013
Total Queries	5,065
Total Legal Articles	17,728
Avg. Query Length	19.43
Avg. Response Length	165.92
Avg. Relevant Articles per Query	1.09
Avg. Keywords per Query	3.57

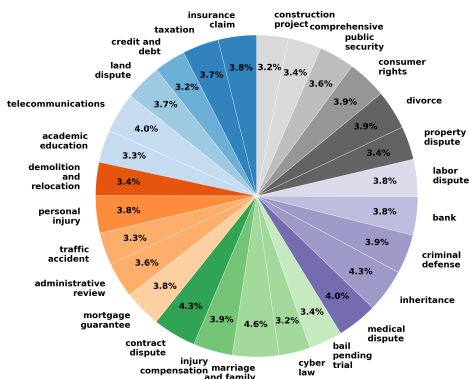


Figure 2: The distribution of query types.

China, children have the right to register for household registration regardless of whether they were born within or outside of marriage. As long as the child’s birth complies with national birth policies and relevant supporting documents (such as the birth certificate, parents’ ID cards, and household registration book) are provided, the child can legally be registered. The legitimacy of household registration is not affected, even for children born out of wedlock.”

- Formulate New Questions.** Due to the difficulty of obtaining real-world multi-turn consultation dialogues, especially those involving personal privacy, annotators are encouraged to expand the questions as much as possible. We also use GPT-4o-mini to generate questions from different perspectives, serving as a reference. These questions are intended to inspire annotators and must not be used directly. For example, based on the previous question, the next query could be: “What documents are needed for household registration?”
- Simulate Real-life Scenarios.** Finally, annotators need to modify the questions to better align with real-life scenarios, including replacing nouns with pronouns and making the language more conversational. For example, the question “What documents are needed for household registration?” can be rephrased as “What documents are needed to register the his household?”

When annotators encounter uncertainties during the annotation process, they should refer to relevant authoritative legal documents, terminology glossaries, or consult legal experts directly to clarify

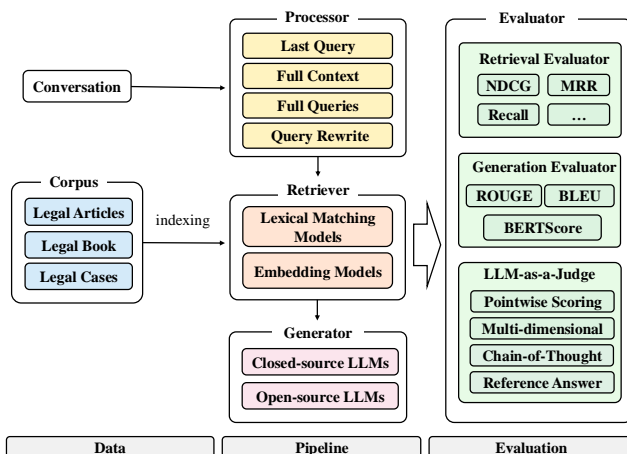


Figure 3: Overview of LexiT Components.

any ambiguities. All decisions made during the annotation process, particularly those following expert consultation, must be clearly documented. This ensures transparency and traceability of decisions, providing a basis for future reviews or revisions and maintaining consistency and standardization. We encourage annotators to actively provide feedback, propose suggestions for improving the annotation process, or elaborate on any challenges encountered during annotation. The annotation guidelines will be regularly reviewed and updated based on this feedback to meet the evolving needs of LexRAG and enhance annotation quality.

3.5 Data Analysis

After careful manual review, LexRAG ultimately contains 1,013 multi-turn conversations, each with 5 interaction rounds. Table 1 presents the basic statistics of LexRAG. The average response length is notably longer than the query length, suggesting that user queries are typically brief, while responses are designed to offer more comprehensive and detailed information.

As shown in Figure 2, LexRAG contains 27 distinct conversation types, which are determined by the seed questions. We observe that the questions are fairly evenly distributed across these types, indicating that LexRAG covers a wide range of legal domains and exhibits diversity. Overall, LexRAG provides a rich and representative sample for evaluating retrieval and generation capabilities in the legal domain.

4 LEXIT

To advance RAG system research in the legal domain, we’ve proposed LexiT, a modular and scalable RAG toolkit for legal researchers. Although there are some general-domain RAG toolkits available, they do not support multi-turn conversations and evaluations tailored to the legal domain [14]. As shown in Figure 3, LexiT consists of three components: Data, Pipeline, and Evaluation. It integrates all elements of the RAG process into a unified framework and supports standalone applications. This modular design enhances flexibility and allows for high customizability in evaluating different legal scenarios.

Table 2: The Prompt Template used in LLM-as-a-judge.

You are an experienced legal expert responsible for evaluating the quality of legal consultation responses. As an impartial and rigorous evaluator, please assess the AI assistant’s response objectively. You will evaluate based on the following five key dimensions:

Factuality: Whether the information provided in the response is accurate, based on reliable facts and legal texts.

User Satisfaction: Whether the response meets the user’s question and needs, and provides a comprehensive and appropriate answer to the question.

Clarity: Whether the response is clear and understandable, and whether it uses concise language and structure so that the user can easily understand it.

Logical Coherence: Whether the response maintains overall consistency and logical coherence between different sections, avoiding self-contradiction.

Completeness: Whether the response provides sufficient information and details to meet the user’s needs, and whether it avoids omitting important aspects.

Longer responses are not necessarily better. The ideal response is short while still meeting the above requirements.

You will be provided with the user’s multi-turn conversation, a reference answer, and the AI assistant’s response to the final question in the conversation. When starting your evaluation, please follow these steps:

1. Compare the AI assistant’s response with the reference answer, highlighting shortcomings and providing further explanations.
2. Evaluate each dimension strictly according to the scoring criteria outlined above. All dimensions must adhere to the high standard of the reference answer, avoiding inflated scores.
3. Combine the evaluations from all dimensions to assign an overall score between 1 and 10. The final score should reflect the overall performance across all dimensions and not be unduly influenced by a single strength.
4. Provide strict and consistent scoring, following the rules below. In general, the higher the quality of the model’s response, the higher the score.

Scoring Standards:

1-2 points: The model provides severe factual errors, incorrect or irrelevant legal texts and interpretations, or completely unrelated responses. The language is confusing, overly long, or incomprehensible, and the structure is extremely complex, causing user confusion. The answer lacks logical coherence, with incoherent reasoning and contradictions, and fails to provide any valid information. Key details are missing.

.....

As an example, the reference answer can score 8 points.

Please provide a detailed evaluation for each dimension, followed by the corresponding score. All scores should be integers. The final evaluation should be returned in the following format.

.....

Data. The data component consists of two key elements: input conversations and corpora. The conversation format can be either single-turn or multi-turn. Single-turn conversations are simple QA dialog, while multi-turn conversations provide previous dialogue history as context. For the corpora, we collect raw data from three different sources. In addition to Legal Articles, which serve as the candidate corpus in this paper, Legal Books and Legal Cases are also included in the toolkit for researchers’ convenience. Specifically, Legal Articles contains 17,228 provisions from various Chinese statutory laws. Legal Book refers to the National Unified Legal Professional Qualification Examination Counseling Book, which consists of 15 topics and 215 chapters, totaling 26,951 provisions. Legal Cases includes 2,370 officially published guiding cases in China. In the future, we plan to expand the corpus with more legal data.

Pipeline. The pipeline component consists of processor, retriever, and generator. The processor is responsible for converting the conversation into queries used by the retriever. There are several strategies for constructing the query, including using the

last question, the entire conversation context, or the entire query history. Moreover, we also predefined a query rewrite strategy, which employs an LLM to integrate all necessary context into a clear, standalone question. Users can easily customize the preprocessing strategy by inheriting and modifying the relevant classes. For the retriever, we integrate various popular retrieval methods. For lexical matching, we use the Pyserini [29] library to implement BM25 [32] and QLD [43]. For dense retrieval, we support advanced models such as BGE [7] and GTE. Users can encode vectors using locally loaded models or API calls. We employ the Faiss [15] for index construction, ensuring compatibility with mainstream indexing formats. In the generator module, we leverage vLLM [16] and Huggingface² to support mainstream LLMs. LexiT also supports flexible prompt customization by combining queries with retrieved content, enabling users to easily adjust generation strategies.

²<https://huggingface.co>

Table 3: Retrieval Performance of different methods on LexRAG using Recall(%) and nDCG(%) metrics. The best results are highlighted in bold.

Retriever	Processor	Recall@1	Recall@3	Recall@5	Recall@10	NDCG@1	NDCG@3	NDCG@5	NDCG@10
BM25	Last Query	5.64	10.60	13.80	18.75	6.13	9.11	10.52	12.21
BM25	Full Context	4.89	11.20	15.02	21.28	5.31	8.92	10.58	12.70
BM25	Full Queries	3.82	7.89	11.86	17.86	4.15	6.58	8.30	10.36
BM25	Query Rewrite	5.73	10.95	14.13	18.84	6.21	9.35	10.74	12.36
BGE-base	Last Query	9.86	19.26	24.40	31.41	10.70	16.22	18.46	20.84
BGE-base	Full Context	6.04	13.09	17.48	25.26	6.55	10.61	12.50	15.17
BGE-base	Full Queries	5.40	12.15	16.64	24.22	5.86	9.75	11.72	14.32
BGE-base	Query Rewrite	9.89	19.19	24.46	31.66	10.74	16.17	18.48	20.92
GTE-Qwen2-1.5B	Last Query	11.37	21.35	26.55	33.13	12.34	18.23	20.46	22.68
GTE-Qwen2-1.5B	Full Context	7.98	16.42	21.77	29.93	8.67	13.45	15.72	18.45
GTE-Qwen2-1.5B	Full Queries	7.11	15.26	20.19	27.71	7.72	12.47	14.58	17.10
GTE-Qwen2-1.5B	Query Rewrite	11.46	21.37	26.60	33.33	12.44	18.29	20.53	22.81
text-embedding-3	Last Query	10.07	18.02	21.91	27.84	10.94	15.56	17.25	19.26
text-embedding-3	Full Context	8.80	17.46	22.80	30.71	9.56	14.53	16.81	19.49
text-embedding-3	Full Queries	6.89	13.86	18.29	24.75	7.48	11.46	13.38	15.55
text-embedding-3	Query Rewrite	10.20	17.97	21.97	28.08	11.08	15.58	17.30	19.39

Evaluation. The evaluation module consists of three key components: the retrieval evaluator, the generation evaluator, and the LLM-as-a-judge. The retrieval evaluator assesses the relevance and accuracy of retrieved documents, supporting the calculation of mainstream automated metrics such as NDCG [36], Recall, MRR [38], Precision, and F1. The generation evaluator measures the consistency between generated responses and reference answers, supporting automated metrics like ROUGE [28], BLEU [31], METEOR [4], and BERTScore [45].

While current automated metrics are useful, they often fail to capture key aspects such as fluency, logical coherence, and factuality, making it difficult to meet the demands of multi-dimensional evaluation criteria. Human evaluation, often considered the gold standard, is time-consuming and labor-intensive, making large-scale assessments difficult. Therefore, we introduce LLM-as-a-judge to enable efficient multi-dimensional automated evaluation. As LLM capabilities continue to advance, they have been widely adopted as evaluators, demonstrating high consistency with human assessments [8, 22, 26]. However, evaluating legal texts remains particularly challenging due to the need for a deep understanding of legal nuances and complex reasoning. To overcome this, we carefully designed the LLM judge evaluation framework within our toolkit to ensure the professionalism and reliability of legal text assessments.

As shown in Figure 3, the LLM-as-a-judge has four key features:

- **Pointwise Scoring.** We use a pointwise scoring method due to its enhanced flexibility and scalability. Specifically, the LLM judge assigns a score from 1 to 10 to each response, considering the dialogue context, the current question, and the reference answer. This method enables a more detailed evaluation of each response while ensuring consistency across the same criteria.
- **Multi-dimensional Evaluation.** Inspired by Wang et al. [35], we develop five evaluation dimensions: Factuality, User Satisfaction, Clarity, Logical Coherence, and Completeness, each with detailed explanations and scoring standards. We also remind the LLM judges that longer responses are not always better, to mitigate potential biases.

- **Chain-of-Thought Reasoning.** To obtain more reliable evaluation results, the LLM-as-a-judge evaluation framework incorporates chain-of-thought reasoning [37]. Specifically, LLM judges first compare the generated response with the reference answer, identify shortcomings, and provide further explanations. Then, they evaluate each dimension based on the established scoring criteria. Finally, the LLM judges combine the evaluations from all dimensions to generate an overall score.
- **Reference-based Evaluation.** Due to the specialized knowledge required for legal evaluations, we provide the LLM judges with human expert-annotated responses as references. These reference answers serve as a baseline, with a score of 8 representing the standard for a well-constructed answer.

In Table 2, we provide the prompt template used in LLM-as-a-judge, which includes the evaluation criteria, chain-of-thought process, scoring standards, and output format requirements.

5 CONVERSATIONAL KNOWLEDGE RETRIEVAL

In this section, we evaluate the performance of different processing strategies and retrieval models in LexRAG.

5.1 Experimental Setting

We evaluate several popular retrieval models, including BM25 [32], BGE-base-zh [7], GTE-Qwen2-1.5B-instruct³, and text-embedding-3-small⁴. These models cover lexical matching and dense retrieval techniques, making them representative. We report commonly used evaluation metrics including Recall and nDCG, evaluated at positions @1, @3, @5 and @10.

For the processor, we test four different strategies.

- **Last Query.** Using the last query in the conversation as input to the retriever.

³<https://huggingface.co/Alibaba-NLP/gte-Qwen2-1.5B-instruct>

⁴<https://platform.openai.com/docs/guides/embeddings>

- **Full Context.** Using the entire conversation as input to the retriever.
- **Full Queries.** Using all queries in the conversation as input to the retriever.
- **Query Rewrite.** Using GPT-4o-mini to turn the relevant context into a clear, standalone question. Specific prompts and examples can be found on our GitHub.

5.2 Retrieval Result

Table 3 presents the performance of different retrieval models and processing strategies on LexRAG. Based on the experimental results, we draw the following conclusions:

- **Comparing Different Retrieval Models.** Dense retrieval methods outperform traditional lexical matching methods like BM25. Overall, GTE-Qwen2-1.5B-instruct achieved the best results. This can be attributed to the challenge that queries in multi-turn consultations often involve pronouns, making basic lexical matching insufficient for identifying relevant legal articles.
- **Comparing Different Process Strategies.** For dense retrieval methods, the query rewrite strategy typically produces the best results. This is likely because it integrates relevant information while minimizing the influence of irrelevant data. Moreover, the last query strategy performs better than using all queries or all contexts. We speculate that this is due to the inclusion of previous conversation content without filtering, which may introduce noise and distort the query’s semantics, ultimately reducing performance. For lexical matching models, such as BM25, the full context strategy generally achieves the best recall results. This is likely because providing more context helps reduce the ambiguity caused by pronouns and other context-dependent terms, improving the retrieval of relevant legal articles. Given these findings, we recommend adjusting processing strategies to align with the strengths of each retrieval method, ensuring optimal performance in different scenarios.
- **Existing LLMs Still Struggle with Conversation Knowledge Retrieval in the Legal Domain.** Overall, current methods perform suboptimally in conversational knowledge retrieval task. Even with the best combination of model and processing strategy, the highest achieved Recall@10 is only 33.33%. This result highlights the challenging of LexRAG, demonstrating that existing retrieval models struggle to effectively handle the nuances of legal consultations. This gap presents an opportunity for the community to create more specialized models that can better address the unique challenges posed by legal contexts.

6 RESPONSE GENERATION

In this section, we report the performance of different LLMs in response generation task.

6.1 Experimental Setting

We evaluated several popular models: GLM-4-flash [11], GLM-4 [11], GPT-3.5-turbo (gpt-3.5-turbo-1106) [1], GPT-4o-mini (gpt-4o-mini-2024-07-18) [1], Qwen-2.5-72B-Instruct [3], LLaMA-3.3-70B-Instruct [34], and Claude-3.5-sonnet (claude-3-5-sonnet-20241022). To reduce the risk of sampling variability, we set the temperature for all LLMs to 0.

We evaluated the performance of LLMs under three settings, simulating ideal and noisy scenarios:

- **Zero Shot.** The LLM generates answers without referencing legal knowledge, relying solely on its internal knowledge and reasoning abilities.
- **Retriever.** The model generates answers using the top 5 documents retrieved by the retriever. In our experiments, we use the GTE-Qwen2-1.5B-instruct combined with query rewriting strategy, as this combination achieved the best recall rate.
- **Reference.** The model generates answers with relevant legal articles annotated by legal experts. This evaluates the LLM’s ability to solve the current issue under ideal knowledge conditions.

We use keyword accuracy and LLM judge scores as evaluation metrics. Since legal terms often have unique meanings, a higher keyword accuracy indicates that the response covers more key legal knowledge. In LLM-as-a-judge, we use the open-source LLM **Qwen-2.5-72B-Instruct** as the evaluator to ensure reproducibility.

6.2 Generation Result

Table 4 reports the performance of LexRAG under different LLMs and settings. Based on the experimental results, we have the following observations:

- In terms of keyword accuracy, the performance under the reference setting is the best, followed by the retriever setting, while the zero-shot setting performs the worst. This indicates that current LLMs lack sufficient legal knowledge to generate relevant response. When provided with relevant legal knowledge, LLMs can generate responses that include more keywords.
- Surprisingly, we observe that in the LLM judge score, the retriever setting does not consistently lead to performance improvements. In contrast, the reference setting consistently results in higher LLM judge scores. We believe this discrepancy occurs because when LLMs are provided with noisy or incomplete legal provisions, their limited legal knowledge prevents them from accurately referencing and analyzing the information, ultimately leading to lower scores. These results suggest that advanced legal consultation systems cannot solely rely on retrieval techniques. To achieve optimal performance, it is crucial to also enhance the foundational LLM’s understanding of legal concepts and reasoning.
- Overall, we observe that Qwen-2.5-72B-Instruct achieved the best performance, followed by GLM-4. This may be due to the fact that these LLMs were developed by the Chinese community, which may make them better suited to legal consultation conversations in the Chinese legal domain. However, even the best-performing LLMs still struggle to achieve a score of 8 in legal consultation scenarios. Given the complexity and precision required in the legal field, we recommend that the community focus on developing AI technologies specifically tailored to the unique needs and nuances of legal contexts.

7 LIMITATION

Although LexRAG advances the evaluation of RAG systems in the legal domain, there are still some limitations that need to be further addressed. First, LexRAG primarily focuses on Chinese legal

Table 4: The Accuracy and LLM judge score of different baselines on LexRAG. The best results are highlighted in bold.

Model	type	1-turn		2-turn		3-turn		4-turn		5-turn		ALL	
		Accuracy	LLM	Accuracy	LLM	Accuracy	LLM	Accuracy	LLM	Accuracy	LLM	Accuracy	LLM
GLM-4-Flash	Zero	0.3431	6.11	0.3534	6.86	0.3738	6.87	0.3737	6.88	0.3726	6.82	0.3633	6.71
GLM-4-Flash	Retriever	0.3403	5.92	0.3670	6.75	0.3783	6.78	0.3794	6.83	0.3820	6.77	0.3694	6.61
GLM-4-Flash	Reference	0.5843	6.52	0.4776	7.06	0.4610	6.99	0.4451	6.93	0.4382	6.89	0.4812	6.88
GLM-4	Zero	0.3468	6.40	0.3462	7.08	0.3782	7.13	0.3809	7.15	0.3836	7.16	0.3671	6.98
GLM-4	Retriever	0.3713	6.24	0.3726	6.87	0.3981	6.90	0.3934	6.92	0.3905	6.88	0.3851	6.76
GLM-4	Reference	0.6151	6.76	0.5423	7.27	0.5208	7.30	0.4906	7.27	0.4862	7.25	0.5310	7.17
GPT-3.5-turbo	Zero	0.3016	6.10	0.3032	6.63	0.3173	6.54	0.3218	6.47	0.3335	6.49	0.3154	6.45
GPT-3.5-turbo	Retriever	0.3217	5.88	0.3057	6.41	0.3220	6.38	0.3278	6.31	0.3231	6.30	0.3200	6.26
GPT-3.5-turbo	Reference	0.5063	6.53	0.4055	6.90	0.3970	6.74	0.3862	6.63	0.3946	6.65	0.4179	6.69
GPT-4o-mini	Zero	0.2982	5.95	0.2962	6.48	0.3195	6.39	0.3075	6.28	0.3219	6.27	0.3086	6.28
GPT-4o-mini	Retriever	0.3308	5.92	0.3395	6.51	0.3411	6.38	0.3445	6.32	0.3468	6.33	0.3405	6.29
GPT-4o-mini	Reference	0.5249	6.39	0.4265	6.83	0.4063	6.62	0.3948	6.47	0.3953	6.48	0.4295	6.56
Qwen-2.5-72B	Zero	0.3583	6.83	0.4037	7.37	0.4260	7.32	0.4271	7.33	0.4266	7.33	0.4083	7.24
Qwen-2.5-72B	Retriever	0.3723	6.46	0.4097	7.24	0.4296	7.23	0.4249	7.27	0.4359	7.28	0.4144	7.09
Qwen-2.5-72B	Reference	0.6045	7.14	0.5260	7.45	0.5186	7.49	0.5117	7.41	0.5015	7.37	0.5324	7.37
Llama-3.3-70B	Zero	0.2556	4.98	0.2695	5.63	0.2846	5.46	0.2800	5.21	0.2894	5.22	0.2758	5.30
Llama-3.3-70B	Retriever	0.2735	5.26	0.2755	5.69	0.2861	5.47	0.2850	5.32	0.2884	5.18	0.2817	5.38
Llama-3.3-70B	Reference	0.5468	5.83	0.4583	6.30	0.4459	6.02	0.4423	5.85	0.4454	5.83	0.4677	5.97
Claude-3.5-sonnet	Zero	0.2464	5.60	0.2856	6.03	0.2989	5.95	0.305	5.86	0.3064	5.91	0.2884	5.87
Claude-3.5-sonnet	Retriever	0.3667	6.42	0.3436	6.90	0.3554	6.88	0.3604	6.79	0.3597	6.77	0.3571	6.75
Claude-3.5-sonnet	Reference	0.5030	6.26	0.4304	6.60	0.4039	6.32	0.3786	6.12	0.3840	6.19	0.4199	6.30

scenarios, which limits its applicability in broader multilingual contexts. We plan to release an updated version supporting English in future iterations to expand its scope and enhance its cross-language evaluation capabilities. Second, due to the privacy and security constraints of real-world multi-turn consultation dialogues, the subsequent dialogue data in LexRAG is primarily annotated by legal experts. While this strategy ensures data quality and legality, it does not fully reflect the diversity and non-standardized interaction scenarios that may occur in real-world legal dialogues. To address this issue, future research will explore ways to leverage simulated data and artificial intelligence technologies while ensuring privacy protection, to better capture the complexity and demands of real-world multi-turn legal consultation conversations.

8 CONCLUSION

In this paper, we introduce LexRAG, a benchmark specifically designed to evaluate RAG systems in multi-turn legal consultation conversations. LexRAG comprises 1,013 consultations and 17,228 candidate legal articles, offering a comprehensive platform for assessing both conversational knowledge retrieval and response generation within the legal domain. In addition, we present LexIT, an open-source evaluation toolkit that provides a set of tools for automated, reproducible assessments of RAG systems in legal contexts. This toolkit enables detailed, fine-grained evaluations of various LLMs and retrieval methods, contributing to the advancement of AI applications in the legal field. In the future, we plan to develop RAG technologies more tailored to legal scenarios and expand LexRAG to support additional languages and legal systems, fostering the global advancement of intelligent judicial technologies.

REFERENCES

[1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal

Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774* (2023).

[2] Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. 2023. Self-rag: Learning to retrieve, generate, and critique through self-reflection. *arXiv preprint arXiv:2310.11511* (2023).

[3] 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* (2023).

[4] Satyanjeev Banerjee and Alon Lavie. 2005. METEOR: An automatic metric for MT evaluation with improved correlation with human judgments. In *Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*. 65–72.

[5] Ilias Chalkidis, Abhik Jana, Dirk Hartung, Michael Bommarito, Ion Androutsopoulos, Daniel Martin Katz, and Nikolaos Aletras. 2021. LexGLUE: A benchmark dataset for legal language understanding in English. *arXiv preprint arXiv:2110.00976* (2021).

[6] Hongshen Chen, Xiaorui Liu, Dawei Yin, and Jiliang Tang. 2017. A survey on dialogue systems: Recent advances and new frontiers. *Acm Sigkdd Explorations Newsletter* 19, 2 (2017), 25–35.

[7] Jianlv Chen, Shitao Xiao, Peitian Zhang, Kun Luo, Defu Lian, and Zheng Liu. 2024. Ege m3-embedding: Multi-lingual, multi-functionality, multi-granularity text embeddings through self-knowledge distillation. *arXiv preprint arXiv:2402.03216* (2024).

[8] Zhumin Chu, Qingyao Ai, Yiteng Tu, Haitao Li, and Yiqun Liu. 2024. PRE: A Peer Review Based Large Language Model Evaluator. *arXiv:2401.15641 [cs.LG]* <https://arxiv.org/abs/2401.15641>

[9] Jiayi Cui, Zongjian Li, Yang Yan, Bohua Chen, and Li Yuan. 2023. ChatLaw: Open-Source Legal Large Language Model with Integrated External Knowledge Bases. *arXiv e-prints* (2023), arXiv–2306.

[10] Wenqi Fan, Yujuan Ding, Liangbo Ning, Shijie Wang, Hengyun Li, Dawei Yin, Tat-Seng Chua, and Qing Li. 2024. A survey on rag meeting llms: Towards retrieval-augmented large language models. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. 6491–6501.

[11] Team GLM, Aohan Zeng, Bin Xu, Bowen Wang, Chenhui Zhang, Da Yin, Dan Zhang, Diego Rojas, Guanyu Feng, Hanlin Zhao, et al. 2024. Chatglm: A family of large language models from glm-130b to glm-4 all tools. *arXiv preprint arXiv:2406.12793* (2024).

[12] Ehsan Hosseini-Asl, Bryan McCann, Chien-Sheng Wu, Semih Yavuz, and Richard Socher. 2020. A simple language model for task-oriented dialogue. *Advances in Neural Information Processing Systems* 33 (2020), 20179–20191.

[13] Minlie Huang, Xiaoyan Zhu, and Jianfeng Gao. 2020. Challenges in building intelligent open-domain dialog systems. *ACM Transactions on Information Systems (TOIS)* 38, 3 (2020), 1–32.

- [14] Jiajie Jin, Yutao Zhu, Xinyu Yang, Chenghao Zhang, and Zhicheng Dou. 2024. FlashRAG: A Modular Toolkit for Efficient Retrieval-Augmented Generation Research. arXiv:2405.13576 [cs.CL] <https://arxiv.org/abs/2405.13576>
- [15] Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2019. Billion-scale similarity search with GPUs. *IEEE Transactions on Big Data* 7, 3 (2019), 535–547.
- [16] Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient Memory Management for Large Language Model Serving with PagedAttention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*.
- [17] Jinqi Lai, Wensheng Gan, Jiayang Wu, Zhenlian Qi, and S Yu Philip. 2024. Large language models in law: A survey. *AI Open* (2024).
- [18] Haitao Li, Qingyao Ai, Jia Chen, Qian Dong, Yueyue Wu, Yiqun Liu, Chong Chen, and Qi Tian. 2024. 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*. 1035–1044.
- [19] Haitao Li, Qingyao Ai, Jia Chen, Qian Dong, Zhijing Wu, Yiqun Liu, Chong Chen, and Qi Tian. 2024. BLADE: Enhancing Black-box Large Language Models with Small Domain-Specific Models. arXiv preprint arXiv:2403.18365 (2024).
- [20] Haitao Li, Qingyao Ai, Qian Dong, and Yiqun Liu. 2024. Lexilaw: A Scalable Legal Language Model for Comprehensive Legal Understanding. <https://github.com/CSHaitao/Lexilaw>
- [21] Haitao Li, Qingyao Ai, Xinyan Han, Jia Chen, Qian Dong, Yiqun Liu, Chong Chen, and Qi Tian. 2024. DELTA: Pre-train a Discriminative Encoder for Legal Case Retrieval via Structural Word Alignment. arXiv preprint arXiv:2403.18435 (2024).
- [22] Haitao Li, Junjie Chen, Qingyao Ai, Zhumin Chu, Yujia Zhou, Qian Dong, and Yiqun Liu. 2024. Calibraeval: Calibrating prediction distribution to mitigate selection bias in llms-as-judges. arXiv preprint arXiv:2410.15393 (2024).
- [23] Haitao Li, Junjie Chen, Jingli Yang, Qingyao Ai, Wei Jia, Youfeng Liu, Kai Lin, Yueyue Wu, Guozhi Yuan, Yiran Hu, et al. 2024. LegalAgentBench: Evaluating LLM Agents in Legal Domain. arXiv preprint arXiv:2412.17259 (2024).
- [24] Haitao Li, You Chen, Qingyao Ai, Yueyue Wu, Ruizhe Zhang, and Yiqun Liu. 2024. Lexeval: A comprehensive chinese legal benchmark for evaluating large language models. arXiv preprint arXiv:2409.20288 (2024).
- [25] Haitao Li, You Chen, Zhekai Ge, Qingyao Ai, Yiqun Liu, Quan Zhou, and Shuai Huo. 2024. Towards an In-Depth Comprehension of Case Relevance for Better Legal Retrieval. In *JSAI International Symposium on Artificial Intelligence*. Springer, 212–227.
- [26] Haitao Li, Qian Dong, Junjie Chen, Huixue Su, Yujia Zhou, Qingyao Ai, Ziyi Ye, and Yiqun Liu. 2024. Llm-as-judges: a comprehensive survey on llm-based evaluation methods. arXiv preprint arXiv:2412.05579 (2024).
- [27] Haitao Li, Yunqiu Shao, Yueyue Wu, Qingyao Ai, Yixiao Ma, and Yiqun Liu. 2024. LeCARDv2: A Large-Scale Chinese Legal Case Retrieval Dataset. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval* (Washington DC, USA) (SIGIR '24). Association for Computing Machinery, New York, NY, USA, 2251–2260. <https://doi.org/10.1145/3626772.3657887>
- [28] Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*. 74–81.
- [29] Jimmy Lin, Xueguang Ma, Sheng-Chieh Lin, Jheng-Hong Yang, Ronak Pradeep, and Rodrigo Nogueira. 2021. Pyserini: An Easy-to-Use Python Toolkit to Support Replicable IR Research with Sparse and Dense Representations. arXiv:2102.10073 [cs.IR] <https://arxiv.org/abs/2102.10073>
- [30] Jinjie Ni, Tom Young, Vlad Pandealea, Fuzhao Xue, and Erik Cambria. 2023. Recent advances in deep learning based dialogue systems: A systematic survey. *Artificial intelligence review* 56, 4 (2023), 3055–3155.
- [31] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*. 311–318.
- [32] Stephen Robertson, Hugo Zaragoza, et al. 2009. The probabilistic relevance framework: BM25 and beyond. *Foundations and Trends® in Information Retrieval* 3, 4 (2009), 333–389.
- [33] SM Tonmoy, SM Zaman, Vinija Jain, Anku Rani, Vipula Rawte, Aman Chadha, and Amitava Das. 2024. A comprehensive survey of hallucination mitigation techniques in large language models. arXiv preprint arXiv:2401.01313 (2024).
- [34] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. arXiv preprint arXiv:2302.13971 (2023).
- [35] Jiayin Wang, Fengran Mo, Weizhi Ma, Peijie Sun, Min Zhang, and Jian-Yun Nie. 2024. A User-Centric Multi-Intent Benchmark for Evaluating Large Language Models. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*. 3588–3612.
- [36] Yining Wang, Liwei Wang, Yuanzhi Li, Di He, and Tie-Yan Liu. 2013. A theoretical analysis of NDCG type ranking measures. In *Conference on learning theory*. PMLR, 25–54.
- [37] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems* 35 (2022), 24824–24837.
- [38] Andrew Worcester and Ted Haines. 2004. Advanced statistics: understanding medical record review (MRR) studies. *Academic emergency medicine* 11, 2 (2004), 187–192.
- [39] Chaojun Xiao, Haoxi Zhong, Zhipeng Guo, Cunchao Tu, Zhiyuan Liu, Maosong Sun, Yansong Feng, Xianpei Han, Zhen Hu, Heng Wang, et al. 2018. Cail2018: A large-scale legal dataset for judgment prediction. arXiv preprint arXiv:1807.02478 (2018).
- [40] Xiaohui Xie, Qian Dong, Bingning Wang, Feiyang Lv, Ting Yao, Weinan Gan, Zhijing Wu, Xiangsheng Li, Haitao Li, Yiqun Liu, et al. 2023. T2ranking: A large-scale chinese benchmark for passage ranking. In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 2681–2690.
- [41] Zihao Yi, Jiarui Ouyang, Yuwen Liu, Tianhao Liao, Zhe Xu, and Ying Shen. 2024. A Survey on Recent Advances in LLM-Based Multi-turn Dialogue Systems. arXiv preprint arXiv:2402.18013 (2024).
- [42] Shengbin Yue, Wei Chen, Siyuan Wang, Bingxuan Li, Chenchen Shen, Shujun Liu, Yuxuan Zhou, Yao Xiao, Song Yun, Xuanjing Huang, et al. 2023. Disc-lawllm: Fine-tuning large language models for intelligent legal services. arXiv preprint arXiv:2309.11325 (2023).
- [43] ChengXiang Zhai. 2008. Statistical language models for information retrieval. *Synthesis lectures on human language technologies* 1, 1 (2008), 1–141.
- [44] Ruizhe Zhang, Haitao Li, Yueyue Wu, Qingyao Ai, Yiqun Liu, Min Zhang, and Shaoping Ma. 2024. Evaluation ethics of llms in legal domain. arXiv preprint arXiv:2403.11152 (2024).
- [45] Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019. BERTscore: Evaluating text generation with bert. arXiv preprint arXiv:1904.09675 (2019).
- [46] Hao Zhou, Chujie Zheng, Kaili Huang, Minlie Huang, and Xiaoyan Zhu. 2020. Kd-Conv: A Chinese multi-domain dialogue dataset towards multi-turn knowledge-driven conversation. arXiv preprint arXiv:2004.04100 (2020).