

# Chain-of-Discussion: A Multi-Model Framework for Complex Evidence-Based Question Answering

Anonymous ACL submission

## Abstract

Open-ended question answering requires models to find appropriate evidence to form well-reasoned, comprehensive and helpful answers. In practical applications, models also need to engage in extended discussions on potential scenarios closely relevant to the question. With augmentation of retrieval module, open-source Large Language Models (LLMs) can produce coherent answers often with different focuses, but are still sub-optimal in terms of reliable evidence selection and in-depth question analysis. In this paper, we propose a novel Chain-of-Discussion framework to leverage the synergy among multiple open-source LLMs aiming to provide **more correct** and **more comprehensive** answers for open-ended QA, although they are not strong enough individually. Our experiments show that discussions among multiple LLMs play a vital role in enhancing the quality of answers. We will release our data and code for further research.

## 1 Introduction

Large Language Models (LLMs) have demonstrated remarkable language generation capabilities (Brown et al., 2020; Touvron et al., 2023; OpenAI, 2023), propelling advancements in various understanding/generation tasks, including open-domain question answering (QA) (Song et al., 2024). However, for complex open-ended question answering, which plays an important role in human-AI interaction, LLMs may still produce output with hallucination and often deliver inferior performance compared to short-form QA (Huang et al., 2023a). This task usually requires LLMs to analyze the questions first, retrieve evidence accordingly, then form a long-form answer which is expected to be correct and well-reasoned with details and proper evidence supported. It has a wide range of applications, from legal consultations and medical advice to education support and financial

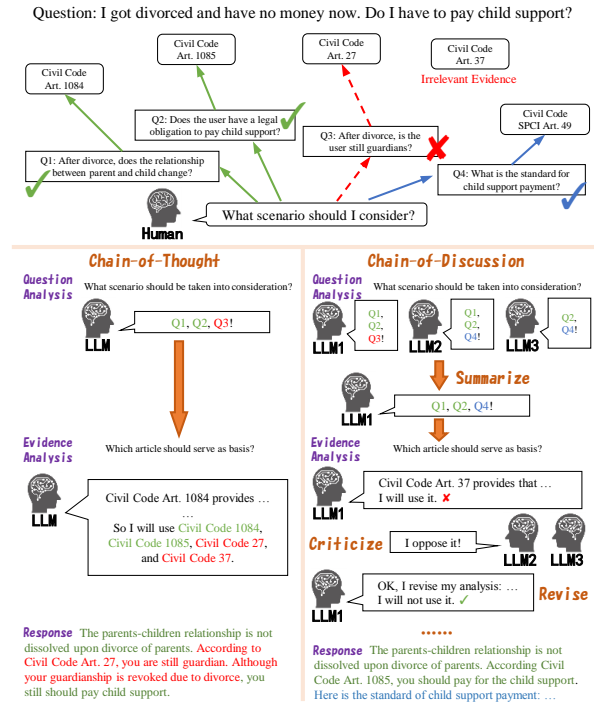


Figure 1: The process of Chain-of-Discussion, compared with chain-of-thought. The green parts are necessary to answer user’s question. Blue parts indicate closely related to the question, useful for detailed/extended discussions. Red parts are irrelevant content that should be avoided.

analysis, where users may pose various complex and knowledge-intensive questions.

Although current LLMs can produce long and coherent texts (Peng et al., 2024), the complex open-ended QA is still an admittedly challenging task, even with augmented retrieval modules. The challenges primarily arise from two aspects.

Firstly, retrieval models are not entirely perfect, inevitably with noise in the retrieval results. As a legal consultation example shown in Figure 1, the model is required to respond to a question regarding the necessity of child support payments. Due to the semantic similarity between obligations for supporting children (financially) and raising/protecting children (physically), the retrieval

model may wrongly return law articles pertaining to guardianship qualifications. LLMs usually cannot filter all these noisy evidence, which may propagate and lead to incomplete analysis, wrong reasoning paths, biased opinions and finally problematic or even misleading answers.

Secondly, we expect LLMs output correct responses and consistent explanations, providing useful suggestions about potential scenarios not directly mentioned in questions but indeed helpful for users' current/near-future situations. In Figure 1, when responding to a question about the obligation to pay child support for a user facing financial difficulties, the model should also remind her/him of the standards for child support payments and ways to negotiate for a reduction given her/his current situation. This is even hard for humans where one should have access to proper evidence, e.g., necessary or closely related law articles, and accordingly provide kind reminders with reasonable explanations. Let alone LLMs without specific training/fine-tuning, which usually focus on the specific facts literally appearing in the questions.

In this work, we focus on the complex evidence-based question answering (CEBQA) task, a typical example of the open-ended QA tasks. We collect a high quality CEBQA dataset consisting of 200 carefully annotated legal consultation questions in the field of marriage and family affairs. To address the challenges, we propose a novel chain-of-thought (CoT) framework, Chain-of-Discussion (CoD), encouraging multiple LLMs to summarize, criticize and revise each other's output to reach a well-supported and helpful response.

Our motivations are two-fold. First, different LLMs may have different intrinsic knowledge and reasoning capabilities due to different training data. Thus, multiple LLMs can be less possible to make errors concurrently than a single LLM. Recent works (Zhang et al., 2023) show checking the consistency across multiple LLMs helps reduce output hallucinations. Specifically, we propose a criticize-and-revise framework, which requires multiple LLMs to discuss and reach a consensus for a better response. For questions that need to involve helpful scenarios or possible extensions, we guess multiple LLMs may provide a diverse set of perspectives to address these possibilities. We thus propose a summarizing step to gather different but helpful perspectives from multiple LLMs, which will eventually form comprehensive and detailed

responses based on the summarized analyses.

Different from existing multi-model interaction works (Chan et al., 2024; Zhang et al., 2023) using strong closed-source LLMs, e.g., GPT-4 (OpenAI, 2023), we decide to take a challenge to study how to best exploit the small-scaled open-source LLMs, e.g., around 7B parameters, for a shared objective, while pushing the boundary of research regarding multi-model interaction.

Our main contributions are as follows: (1) We collect a high-quality CEBQA dataset consisting of 200 legal consultation questions in Chinese with carefully annotated evidence and answers. (2) We propose a novel chain-of-discussion framework, i.e., summarize-criticize-revise, which harnesses the synergy among multiple open-source LLMs to generate more accurate and helpful responses. (3) Both GPT-4-based and evidence-centric evaluations demonstrate our framework can help small-scaled LLMs benefit from each other and improve the overall quality in terms of correctness and comprehensiveness.

## 2 Related Works

**Retrieval-Augmented Generation** Lewis et al. (2020) initially propose the paradigm of retrieval-augmented generation (RAG), which can effectively reduce hallucinations within the texts generated by LLMs. RAG offers a vital solution to mitigate the problem of LLMs lacking domain-specific knowledge, thereby enhancing the credibility of LLMs (Gao et al., 2023). In the RAG paradigm, models typically undergo multiple generation steps to achieve the final results. For a user input, models first run a retriever to scan the store of evidence to select several documents as reference. Subsequently, models should determine when and whether to use each evidence document before generating (Izacard et al., 2022; Shi et al., 2023b; Yu et al., 2023; Trivedi et al., 2023).

In this work, we face more complex challenges than RAG. While the model filters out irrelevant evidence, it also needs to retain evidence relevant to potential scenarios. Sometimes, determining which evidence can be used for potential scenarios and which are irrelevant is also challenging for humans.

**Chain-of-Thought** Previous works demonstrate that LLMs have a promising capability to decompose a complex question into several intermediate steps (Wei et al., 2022; Kojima et al., 2022). By segmenting the original question, LLMs can focus

on handling each simple sub-question at each step, thus yield more accurate results (Zhou et al., 2023). The CoT framework is now widely employed in diverse practical NLP applications (Zelikman et al., 2022; Shi et al., 2023a; Wang et al., 2023). Previous works also employ CoT in the self-correction process of LLMs, which aims to re-generate better outputs. For instance, in Chain-of-Verification, the model generates several queries to verify its original answer, and then revise the answer based on the verification results (Dhuliawala et al., 2023). Most of these efforts perform self-checking based on a single model. However, we study a novel CoT framework for multi-model interactive checking and re-generating.

### 3 Preliminaries

**Task Definition** In CEBQA tasks, given a user’s question  $q$  and a store of evidence documents  $\mathcal{D}$ , a model should analyze  $q$  first, find necessary evidence  $\mathcal{D}_q = \{d_1, \dots, d_t\}$  from  $\mathcal{D}$  accordingly and generate a paragraph  $r$  as the final response. For instance, in the legal consultation task, users may ask what to do given her/his current situation. The model should find supportive evidence from a store of law articles or previous legal cases, and generate a helpful and detailed response.

Specifically, we expect the generated responses to meet the requirements in terms of correctness and comprehensiveness. (1) **Correctness**: The responses should be based on the evidence that can support to answer the questions, and refrain from employing irrelevant evidence or misinterpreting evidence out of context. (2) **Comprehensiveness**: The responses should engage in discussions about potential scenarios that would be relevant/helpful to users, even not explicitly mentioned in questions.

We note that it is hard to guarantee all the retrieved evidence pieces can be perfectly used to answer the question. Therefore, similar to RAG, models should filter out irrelevant evidence. However, it is more challenging for models to carefully retain the evidence that can be used for discussions about potential scenarios, even though the evidence may not directly support answering the question.

**Baseline Framework: CoT** Previous works have revealed that the CoT prompt can enhance the ability of LLMs to handle complex reasoning tasks (Wei et al., 2022; Kojima et al., 2022). Inspired by these works, we employ a multi-step prompt to stimulate LLMs to generate more correct

while comprehensive answers.

We initially prompt LLMs to analyze the question  $q$ , including identifying the possible role of users, understanding explicit and implicit demands of users, and determining what types of evidence is needed to answer the question. The generated analysis of question can be denoted as  $a_q^{\text{que}}$ .

The next step is to judge whether each evidence document can serve as a potential basis for responding to the question  $q$ . Here, we employ a prompt to feed the LLM with question  $q$ , analysis  $a_q^{\text{que}}$  of the question, and a specific evidence document  $d_i$ . The LLM then need to analyze whether  $a_{d_i}^{\text{evi}}$  can be used to address the issues raised in  $q$  and whether evidence  $d_i$  can probably be used to respond or not.

The LLM with parameters  $\theta$  should finally respond to the question  $q$  according to question analysis  $a_q^{\text{que}}$  and evidence analysis  $\{a_{d_i}^{\text{evi}}\}_i$ , based on the evidence document set  $\mathcal{D}_q$ :

$$r = f(q, \mathcal{D}_q, a_q^{\text{que}}, \{a_{d_1}^{\text{evi}}, \dots, a_{d_t}^{\text{evi}}\} | \theta).$$

As observed in our pilot study, one small-scaled LLM could generate fluent answers, but often with incomplete analysis or wrong reasoning paths.

### 4 CoD: Summarize, Criticize, and Revise

Our Chain-of-Discussion framework leverages interactive discussions among multiple LLMs, thereby addressing potential shortcomings in individual’s intrinsic knowledge.

Similar to the baseline, we employ a two-stage analyzing pipeline that instructs LLMs to analyze the question and evidence separately. To address the correctness and comprehensiveness of generated answers, during question analysis, we encourage models to read and summarize others’ analyses so as to consider more scenarios closely relevant to the question, in the purpose of augmenting the comprehensiveness. During the stage of evidence analysis, we require all other LLMs to **criticize** the evidence analysis of each LLM. Subsequently, the model will read others’ critique and determine whether to **revise** its own analysis or not. The model finally generate a correct and more helpful response based on the summarized question analysis and revised evidence analysis.

#### 4.1 Stage 1: Question Analysis

Formally, suppose there are  $n$  accessible LLMs, denoted as  $M_1, \dots, M_n$ . For a given question  $q$  and the retrieved evidence  $\mathcal{D}_q$ , we aim to employ

the target LLM  $M_k$  to generate a response, with the assistance of the remaining LLMs.

We first instruct the LLMs to analyze the question, including facts mentions in  $q$ , primary needs of the user, and potential scenarios associated with the question. We observe that LLMs may perform poorly in analyzing potential scenarios when solely relying on their intrinsic knowledge, especially those models that have not been pre-trained or supervised fine-tuned on domain-specific data. Thus, we argue that the evidence documents  $\mathcal{D}_q$  can serve as vital cues about the potential scenarios not mentioned in  $q$ .

Different LLMs can have varying preferences in analyzing the potential scenarios. Therefore, we believe that by integrating the outputs of multiple LLMs, we can take more helpful scenarios into account, thus improve the **comprehensiveness** of question analysis. We prompt each LLM  $M_i$  to analyze the question  $q$ , with retrieved evidence  $\mathcal{D}_q$  as a reference:  $a_{q, M_i}^{\text{que}} = f_{\text{que}}(q, \mathcal{D}_q | \theta_{M_i})$ .

We then employ the target LLM  $M_k$  to **summarize** the question analyses of all models, according to following instructions:

- **Consistency:** If the majority of LLMs provide similar analyses regarding a fact in the question or a potential scenario, then it is likely to be correct. You can include it in the summary.

- **Comprehensiveness:** If a minority of LLMs hold a particular viewpoint in their analyses with reasons, it does not imply its unreliability. You should scrutinize this content, assessing its logical coherence and relevance to the question.

The summarized question analysis can be  $a_q^{\text{que}} = f_{\text{sum}}(q, a_{q, M_1}^{\text{que}}, \dots, a_{q, M_n}^{\text{que}} | \theta_{M_k})$ .

## 4.2 Stage 2: Evidence Analysis

Incorporating many irrelevant evidence documents as input would inevitably introduce noise, which could deteriorate the model performance. Thus, we should discern which evidence document should be used to address the question. For an evidence document  $d_j \in \mathcal{D}_q$ , we prompt the target model  $M_k$  to analyze it based on the question and question analysis:  $\hat{a}_{d_j}^{\text{evi}} = f_{\text{evi}}(d_j, q, a_q^{\text{que}} | \theta_{M_k})$ .

However, a single LLM might generate hallucinated outputs (Li et al., 2023b; Huang et al., 2023a), and incorrectly assess the relevance between evidence documents and the given question. Inspired by previous work (Zhang et al., 2023), we propose a multi-party discussion framework to improve the

quality of evidence analysis.

First, we instruct each LLM, excluding  $M_k$ , to **criticize** the evidence analysis  $\hat{a}_{d_j}^{\text{evi}}$ . Each critic model  $M_i$  should explicitly output whether it holds opinions contrary to  $\hat{a}_{d_j}^{\text{evi}}$ , which are denoted as  $c_i^{d_j}$ . In this work, we employ a revising threshold  $\delta$ . If the proportion of opposite opinions in the critiques exceeds  $\delta$ , the target model needs to **revise** its evidence analysis:  $a_{d_j}^{\text{rev}} = f_{\text{rev}}(q, d_j, a_q^{\text{que}}, \hat{a}_{d_j}^{\text{evi}} | \{c_i^{d_j}\}_i, \theta_{M_k})$ .

We assume that the critique requiring to revise can be reliable only when a majority of critic models achieve a consensus. Otherwise, we retain the original evidence analysis. Formally, we collect the evidence analysis as following:

$$a_{d_j}^{\text{evi}} = \begin{cases} \hat{a}_{d_j}^{\text{evi}}, & \text{if } \frac{|\{c_i | c_i = \text{opposite}\}|}{|\{c_i\}|} \leq \delta; \\ a_{d_j}^{\text{rev}}, & \text{otherwise.} \end{cases}$$

## 4.3 Response Generation

For a fair comparison, we employ prompts similar to those of the baseline framework to generate responses. We denote the response as  $r = f_{\text{ans}}(q, \mathcal{D}_q, a_q^{\text{que}}, \{a_{d_1}^{\text{evi}}, \dots, a_{d_n}^{\text{evi}}\} | \theta_{M_k})$ .

## 5 Experiments

As there is no existing dataset for CEBQA tasks, here we delve into the legal consultation task and collect the first CEBQA dataset. In civil law systems like China, all legal activities, including legal consultation, should be based on *law articles* (or *judicial interpretations*), which can be naturally considered as the evidence store in our framework.

### 5.1 Data Collection

To reflect the diversity in practical scenarios, we focus on the fields of *marriage, family affairs, and inheritance*, covering various legal disputes, e.g., *divorce, custody, contracts, property*, etc. From a pool of 609 questions collected from real users with consultants' responses through Web Search Engines, we employ an legal expert to select 200 questions to ensure their semantic distinctiveness and coverage. See more details in Ethics Statement.

**Evidence Annotation** We construct the evidence store based on all 1,749 articles of the *Civil Code, Civil Procedure Law* and their judicial interpretations in China. For each question, we consider three types of articles: *necessary, optional*, and

338 *not required*. *Necessary* articles are highly rele- 386  
339 vant to the question, while *optional* ones can be 387  
340 basis for discussions of potential scenarios (more 388  
341 details in Appendix B). We ensure there are 5 arti- 389  
342 cles retrieved for each question, and on average, 390  
343 each contains 1.52 *necessary*, 1.23 *optional*, and 391  
344 2.25 *not required* articles. It means about 45% of 392  
345 the retrieved articles are not required at all. 393

346 **Data Quality** We employ 6 annotators with back- 394  
347 ground in civil law to manually check the questions, 395  
348 answers and articles. Annotators are instructed to 396  
349 correct all typos but retain the informal expressions 397  
350 in questions. Note that there are many omissions or 398  
351 slight word-order inversions in the questions, pos-  
352 ing a challenge to models’ reasoning capabilities.

353 Annotators also examine the correctness and log-  
354 ical coherence of the responses. For problematic  
355 ones, annotators are encouraged to discuss and  
356 reach a consensus for modifications, otherwise,  
357 leave them as they are. Averagely, it takes about 20  
358 hours per annotator to examine 100 instances.

## 359 5.2 Experimental Setup

360 As our main focus is to investigate whether  
361 small open-source LLMs can collaborate through  
362 summarize-criticize-revise, we study four open-  
363 source fine-tuned LLMs, Baichuan2-7B (Baichuan,  
364 2023), Deepseek-7B (DeepSeek-AI, 2024), Qwen-  
365 7B (Bai et al., 2023), and Xverse-7B<sup>1</sup>, which are  
366 four of the best 7B-parameter LLMs performing on  
367 CMMLU (Li et al., 2023a). When we use a specific  
368 LLM as the target model, the other three LLMs are  
369 expected to generate diverse question analyses and  
370 criticize the evidence analysis of target model.

371 To examine if close-source LLMs can still bene-  
372 fit from our CoD, we also test with *gpt-3.5-turbo-*  
373 *1106*, *gemini-1.0-pro-latest*, and *claude-3-haiku-*  
374 *20240307* similarly to the open-source group.

375 We note that the two stages in Chain-of-  
376 Discussion are independent of each other. There-  
377 fore, we can investigate how they contribute to the  
378 ultimate performance by the following settings:

379 **Single-model baselines (BS):** One LLM takes a  
380 query and several retrieved articles as input and  
381 performs question analysis, article analysis, and  
382 response generation in a vanilla CoT manner.

383 **Only Stage 1 (S1):** All LLMs produce question  
384 analysis. The target LLM summarizes these analy-  
385 ses, and proceeds to the rest by itself.

**Only Stage 2 (S2):** Three other LLMs criticize  
the article analysis generated by the target LLM.  
The question analysis and the final response are  
generated by the target LLM on its own.

**Chain-of-Discussion (S1S2):** All LLMs involve  
into both question analysis and article analysis.  
Eventually, the target LLM produces the response  
by itself.

We employ each LLM as the target model, repli-  
cating the experimental settings. We report the  
performance for each LLM as the target role. More  
details and prompt templates are in Appendix A  
and E.

**Evaluation Metrics** Different from conventional  
QA tasks, the responses in CEBQA tasks can  
consist of several hundred or even thousands of  
words, which are also knowledge intensive and  
complex in structures, containing facts and causal  
relations to be verified. Therefore, it is impossible  
to employ the popular metrics such as F1 or exact  
match (Joshi et al., 2017; Rajpurkar et al., 2018).

Following previous works (Liu et al., 2023; Chan  
et al., 2024), we employ GPT-4 to evaluate the  
quality of generated responses, with the expert-  
written responses, necessary and optional articles  
as reference. We prompt `gpt4-turbo-0125` to  
score a response in an integer between 1 and 10  
based on correctness and comprehensiveness. If  
there is no clear reason to indicate a response is  
significantly better or worse than human-written  
ones, a score of around 7 should be given (Scoring  
prompt is in Appendix F).

## 5.3 Main Results

Table 1 shows the evaluation results produced by  
GPT-4. Comparing the baseline CoT setting (BS)  
and Chain-of-Discussion (S1S2) in the open-source  
group, we can find **each LLM can obtain improve-  
ments from discussions with other LLMs**, with  
Baichuan2-7B increased by +0.340, Deepseek-7B  
by +0.115, Qwen-7B by +0.120, Xverse-7B by +  
0.110. We think other LLMs bring more related  
aspects according to their own strengths into dis-  
cussions while the baseline CoT setting has to rely  
on one LLM only. We also find employing multi-  
model discussion on both stages can bring more  
improvement than using it on one stage only.

Although CoD can enhance all LLMs, the CoD-  
augmented Baichuan2-7B, Qwen-7B, or Xverse-  
7B still can not outperform Deepseek-7B under its  
baseline setting, with around 0.5 left behind. This

<sup>1</sup><https://huggingface.co/xverse/XVERSE-7B-Chat>

Target LLM	Setting	Avg. Score	$\Delta$ Score
Open-source LLMs Group			
Baichuan2-7B	BS	5.750	–
	S1	6.030	+0.280
	S2	5.935	+0.185
	S1S2	<b>6.090</b>	+0.340
Deepseek-7B	BS	6.465	–
	S1	6.505	+0.040
	S2	6.480	+0.015
	S1S2	<b>6.580</b>	+0.115
Qwen-7B	BS	5.835	–
	S1	5.890	+0.055
	S2	5.815	-0.020
	S1S2	<b>5.955</b>	+0.120
Xverse-7B	BS	6.015	–
	S1	5.995	-0.020
	S2	6.030	+0.015
	S1S2	<b>6.125</b>	+0.110
Close-source LLMs Group			
GPT-3.5-turbo	BS	6.895	–
	S1S2	<b>6.955</b>	+0.060
Gemini-1.0-pro	BS	6.940	–
	S1S2	<b>6.975</b>	+0.035
Claude-3-haiku	BS	7.280	–
	S1S2	<b>7.300</b>	+0.020

Table 1: The average scores of each target LLM and each setting evaluated by GPT-4. The upper is evaluated within the open-source LLMs group, while the bottom is within the closed-source LLMs group.

may indicate that the quality of responses primarily relies on the inherent ability of LLMs to analyze context and then to generate.

As shown in the bottom of Table 1, even for powerful closed-source LLMs like Claude-3, our CoD framework can still bring improvement consistently, but the improvements are admittedly small compared to the open-source group.

Through a manual check, we also find that these powerful closed-source LLMs may also produce erroneous analyses of questions, resulting in incorrect responses. We believe that regardless the model scales, our CoD framework can leverage the varied knowledge and capabilities of multiple models to achieve better performance.

## 6 Discussions

### 6.1 Evidence-Centric Evaluation

Besides overall evaluation by GPT-4, we wonder if our Chain-of-Discussion framework can enhance the comprehensiveness and correctness of the model output. When discussing the details of questions or potential scenarios, LLMs should refer to *necessary* or *optional* evidence. Hence, we can

Target LLM	Setting	N-Acc%	O-Acc%
Baichuan2-7B	BS	58.26	50.14
	S1	60.03	<u>50.67</u>
	S2	<u>61.86</u>	50.25
	S1S2	<b>63.17</b>	<b>52.38</b>
Deepseek-7B	BS	75.93	59.27
	S1	<u>76.36</u>	<u>59.70</u>
	S2	76.12	59.23
	S1S2	<b>76.79</b>	<b>59.80</b>
Qwen-7B	BS	69.87	60.98
	S1	70.31	61.63
	S2	<u>70.64</u>	<u>63.65</u>
	S1S2	<b>71.29</b>	<b>64.20</b>
Xverse-7B	BS	74.00	63.95
	S1	74.24	<u>64.72</u>
	S2	<u>75.67</u>	64.44
	S1S2	<b>76.16</b>	<b>65.35</b>

Table 2: The Macro average N-Acc and O-Acc results of each target LLM and each setting. The highest scores are made **bold**, while the second underlined.

assess the correctness and comprehensiveness of responses by the accuracy of evidence documents.

We accordingly design two metrics of accuracy, N-Acc and O-Acc, to assess the correctness and comprehensiveness, respectively. We utilize the *not required* articles as negative samples. For N-Acc, we employ the *necessary* articles as positive samples, while the *optional* articles for O-Acc. We employ heuristic methods to examine if a response has used an article (details in Appendix C).

We compute the Macro average N-Acc and O-Acc across all examples, shown in Table 2. We can see these LLMs exhibit differently on their own, and when performing a simple majority vote among all LLMs, we could get 77.62% for N-Acc, confirming our hypothesis that **it is possible to harness the synergy among multiple open-source LLMs for a better performance**. Compared to the baselines (BS), our CoD (S1S2) can achieve around a 2% improvement on both N-Acc and O-Acc for Baichuan2-7B, Qwen-7B, and Xverse-7B. Even for Deepseek-7B, the best in GPT-4 based evaluation, our framework still brings improvements of 0.86% and 0.53% to N-Acc and O-Acc, respectively.

The results indicate that introducing multi-model discussions during both question analysis and evidence analysis **helps target LLMs better refer to correct evidence**. This also explains why CoD can improve the quality of model responses.

Comparing the results under BS, S1, and S2, we find that involving multiple LLMs in a single stage can actually enhance both correctness and comprehensiveness. Specifically, employing multi-model

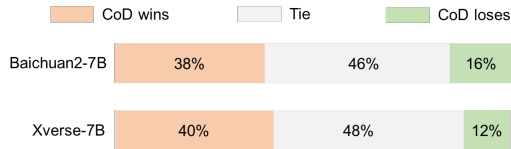


Figure 2: Human preference evaluation, comparing the CoD settings of Baichuan2-7B and Xverse-7B to their corresponding baseline settings across 50 randomly sampled examples.

discussions in question analysis contributes more to comprehensiveness, while introducing other models in evidence analysis (BS vs. S2) helps filtering irrelevant evidence thus brings more improvement in correctness.

## 6.2 Human Evaluation

We further examine the quality of generated responses through a win-rate analysis by human experts. We randomly sample 50 examples, and ask two human experts to examine whether the responses of CoD-enhanced models are better than those of baselines. We use the responses of Baichuan2-7B and Xverse-7B, which are models obtaining the most and the least improvement from CoD, respectively.

As shown in Figure 2, in 38% of the examples, CoD can bring improvements to the responses of Baichuan2-7B, while 40% responses generated by Xverse+CoD surpass Xverse’s baseline setting. For Baichuan2-7B, only 16% of responses deteriorated after introducing the CoD framework, while for Xverse, only 12% worsened.

**Qualitative Analysis.** We also manually examine 30 responses of both models. The main error types are *misunderstanding key legal terms* (60%), *discussing not-required articles* (15%), and *inconsistent explanations* (10%). The latter is often caused by term misunderstanding. This shows that there may not be sufficient legal knowledge in 7B models to distinguish very similar terms. Comparing our CoD and baseline setting, we can see CoD can correct improper articles for 6 out 30 cases, leading to better results, through the cooperation mechanism.

**Correlation between GPT-4 Evaluation and Human** Although GPT-4 has been deemed as a reliable evaluator in various LLM evaluations (Liu et al., 2023), we still want to examine how well GPT-4 preferences correlate to human experts in our challenging task. For each example from the randomly sampled 50 cases: (1) if GPT-4 assigns a higher score for the response generated by the CoD setting, we denote this evaluation result as +1; (2)

if GPT-4 assigns a higher score for the response by the baseline setting, we denote this evaluation result as -1; Otherwise, we denote it as 0. We use the same principle to annotate the evaluation results of human expert.

The Pearson correlation coefficient between GPT-4 scores and human expert evaluation is 0.5863 with  $p = 7.7 \times 10^{-6}$ . This indicates a relatively high correlation between the two evaluation methods, where we could consider GPT-4 evaluation as a reliable manner for our work.

## 6.3 Limitations of Open-Source LLMs

As shown in Table 1, using multi-model interaction at only one stage fails to enhance Xverse-7B (S1) or Qwen-7B (S2). We guess this may be due to the small size of the model scales, 7B here.

After a manual check, we find the main difficulties for those small scaled open-source LLMs: 1) those 7B LLMs may not have sufficient legal/background knowledge and reasoning capabilities, thus often fail to discern errors in the analysis of others and integrate them into the summary; 2) the small scaled LLMs do not always follow instructions, and often refuse to modify the errors in evidence analysis.

**Hallucination Propagation in Summary.** Table 3 illustrates an example where Baichuan-7B incorporates inaccurate analysis from DeepSeek-7B into the summary. We find that both LLMs can correctly identify the primary intention from the question, i.e., determining whether providing support for parents is a prerequisite for inheriting. Unfortunately, Deepseek-7B mistakenly engages in analysis regarding whether parents pay support for children, while Baichuan2-7B fails to identify this as irrelevant to the question.

Although both raising children and supporting parents are financial support and life care between two generations in a family, it is a commonsense, for people with legal background, that these are two different concepts and legally, there is no causal relationship between them. We think both LLMs actually rely on shallow semantic similarity to judge if a question analysis is relevant to user’s question.

**Capturing Nuance between Related Terms.** We guess that open-source LLMs with around 7B parameters may not have sufficient knowledge and enough reasoning capabilities to distinguish related but different terms. Table 4 shows an example where Qwen-7B fails to recognize that 复婚 (rec-

**Question:** Can a daughter/son still be eligible to inherit if she/he has not supported her/his parents?

*Original question analysis:*

**Baichuan2:** This question involves the relationship between inheritance rights and filial duties. The user may be concerned about whether a person who has not fulfilled filial duties still qualifies for inheritance. This may entail rules regarding the recognition of inheritance rights, the forfeiture and reinstatement of inheritance rights, and the role of filial duties in the determination of inheritance rights.

*Question analyses of other LLMs:*

**Qwen:** ... **Xverse:** ...

**Deepseek:** ... The user may want to clarify whether not supporting parents will affect her/his eligibility to inherit parents' estate. Additionally, user may also inquire about her/his entitlement to demand child support payments from her/his parents, if the parents refuse to raise children. ...

*Summarized question analysis:*

**Baichuan2:** This question involves the relationship between inheritance rights and filial duties. The user may wonder if her/his child retains eligibility for inheritance despite failing to fulfill filial duties. This may entail rules regarding the recognition of inheritance rights, the forfeiture and reinstatement of inheritance rights, and the role of filial duties in the determination of inheritance rights. The user may also want to know the regulations regarding whether parent who does not raise children is required to pay child support.

Table 3: An example of hallucinated summary of question analysis. Violet texts are correct analysis, while the texts with yellow background are hallucinated parts.

584 oncile and remarry) and 再婚 (remarry with other  
585 person) represent different marriage relationships,  
586 comply with different articles or provisions, and  
587 should not be used interchangeably. However, in  
588 our CoD setting, the discussions with other mod-  
589 els can, to some extent, have the chance to reduce  
590 the hallucination caused by the target model's poor  
591 reasoning capabilities, thus help to reduce logical  
592 errors in Qwen-7B's responses.

**Reasoning with Knowledge.** Similarly, we find  
593 the lack of good reasoning capabilities and knowl-  
594 edge may also limit the performance of LLMs in  
595 revising evidence analysis. LLMs tend to focus on  
596 the literal differences between questions and law  
597 articles. Users often omit commonsense informa-  
598 tion in their questions, whereas law articles exhibit  
599 a rigorous and detailed structure. The stylistic dif-  
600 ferences between questions and evidence may lead  
601 to the misjudgments of LLMs (more examples and  
602 analysis in Appendix D.2).  
603

### 6.3.1 Model Preference of Revising

604 We find different LLMs have varying preferences  
605 regarding whether to modify the evidence analysis.

606 To quantitatively study this issue, we initially  
607 employ Baichuan2-7B to generate the original evi-  
608

**Question:** What is the difference between *resumption of marital relationship* (复婚) and *remarriage with other person* (再婚)?

**Article:** Article 1046 A man and a woman shall enter into marriage freely and voluntarily. ...

Article 1083 Where, after divorce, both the man and the woman voluntarily intend to resume their marital relationship, they shall file for re-registration...

**Qwen-7B:** ... According to Article 1083, both parties should be voluntary for *resumption of marital relationship*, while there is no such limitation for *remarriage with other person*. ...

**Qwen-7B+CoD:** ... According to Article 1046, whether it is *Fuhun* or *Zaihun*, both parties need to do so voluntarily.

Table 4: A case of Qwen-7B obtaining improvement from CoD. Violet texts are correct analysis, while the texts with yellow background are hallucinated parts.

609 dence analyses. Deepseek-7B and Xverse-7B then  
610 play the role of critics to determine which analysis  
611 should be revised. Then, we use Baichuan2-7B and  
612 Qwen-7B to revise these analyses, respectively. We  
613 find that Baichuan2-7B successfully revises 96.5%  
614 of the analyses, while Qwen-7B can only revise  
615 56.1% of them.

616 To mitigate the influence of target model selec-  
617 tion, we also use Qwen-7B to provide original evi-  
618 dence analyses, still with Deepseek-7B and Xverse-  
619 7B as the critics. Similarly, Baichuan2-7B can  
620 revise 92.5% of the analyses, but Qwen-7B only  
621 revise 67.2% of them.

622 We argue that an LLM's preference for refus-  
623 ing to revise may lead to a failure to obtain bet-  
624 ter evidence analysis based on the critiques. Con-  
625 sequently, it might limit our Chain-of-Discussion  
626 framework to bringing more improvement as ex-  
627 pected. The preference of LLMs can be af-  
628 fected by supervised fine-tuning and reward mod-  
629 eling (Ouyang et al., 2022; Rafailov et al., 2023).  
630 We hope to study the effect of supervised training  
631 on Chain-of-Discussion in future.

## 7 Conclusions

632 In this work, we proposed a novel reasoning frame-  
633 work, Chain-of-Discussion, for complex evidence-  
634 based question answering tasks. The CoD frame-  
635 work involves multiple LLMs in discussions to  
636 achieve more correct and comprehensive responses  
637 with less hallucination and more supportive evi-  
638 dence. Experiments on a challenging legal consul-  
639 tation dataset show CoD can effectively improve  
640 the performance of open-source LLMs by encour-  
641 aging them to discuss and criticize.  
642



## 643 Limitations

644 Our proposed framework is designed to generate  
645 correct and comprehensive answers to respond  
646 complex questions. When used for providing le-  
647 gal advisory services, this technique can produce  
648 helpful responses to help people with needs, but it  
649 still cannot guarantee all responses are completely  
650 correct. Hence, this techniques should be used with  
651 cautions for further applications.

652 Our dataset is designed and annotated to reflect  
653 the nature of CEBQA tasks, which requires models  
654 to generate detailed analysis to each closely rel-  
655 evant scenarios of the user’s question. However,  
656 our annotated results may be inevitably not perfect  
657 from the professional perspectives of experts in  
658 civil law. Thus it should be used with caution and  
659 for research purpose only.

660 We also note that the proposed framework in-  
661 volves multiple LLMs to generate for several  
662 rounds. Straightly using commercial APIs may  
663 lead to more promising generated results and cost  
664 less time. However, our aim is to validate how  
665 to better and more efficiently exploit the synergy  
666 among small LLMs, without relying on larger  
667 LLMs. We pioneer to expand the border of in-  
668 vestigation about multi-model interaction to the  
669 small open-source LLMs.

## 670 Ethics Statement

671 In this work, we collect the legal consultation  
672 question-answer pairs through Web Search Engine.  
673 All these data are publicly accessible on the inter-  
674 net. Additionally, these websites all allow crawlers  
675 of search engines to automatically access their web  
676 pages.

677 To protect the privacy of the consultants in real  
678 world, we manually check all the instances, and re-  
679 move those which may indicate the name, address,  
680 phone number of a specific users. Note that infor-  
681 mation such as gender, age, and family composition  
682 serves as crucial background in judicial practice.  
683 Since this information does not identify a specific  
684 individual and does not result in privacy leaks, we  
685 do not discard data where the users mention their  
686 gender, age, or family composition.

## 687 References

688 Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang,  
689 Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei  
690 Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin,

Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu, 691  
Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren, 692  
Xuancheng Ren, Chuanqi Tan, Sinan Tan, Jianhong 693  
Tu, Peng Wang, Shijie Wang, Wei Wang, Sheng- 694  
guang Wu, Benfeng Xu, Jin Xu, An Yang, Hao Yang, 695  
Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu, 696  
Hongyi Yuan, Zheng Yuan, Jianwei Zhang, Xingx- 697  
uan Zhang, Yichang Zhang, Zhenru Zhang, Chang 698  
Zhou, Jingren Zhou, Xiaohuan Zhou, and Tianhang 699  
Zhu. 2023. [Qwen technical report](#). 700

Baichuan. 2023. [Baichuan 2: Open large-scale lan- 701  
guage models](#). *arXiv preprint arXiv:2309.10305*. 702

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie 703  
Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind 704  
Neelakantan, Pranav Shyam, Girish Sastry, Amanda 705  
Askell, Sandhini Agarwal, Ariel Herbert-Voss, 706  
Gretchen Krueger, Tom Henighan, Rewon Child, 707  
Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, 708  
Clemens Winter, Christopher Hesse, Mark Chen, Eric 709  
Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, 710  
Jack Clark, Christopher Berner, Sam McCandlish, 711  
Alec Radford, Ilya Sutskever, and Dario Amodei. 712  
2020. [Language models are few-shot learners](#). 713

Chi-Min Chan, Weize Chen, Yusheng Su, Jianxuan Yu, 714  
Wei Xue, Shanghang Zhang, Jie Fu, and Zhiyuan Liu. 715  
2024. [Chateval: Towards better LLM-based eval- 716  
uators through multi-agent debate](#). In *The Twelfth 717  
International Conference on Learning Representa- 718  
tions*. 719

Wanxiang Che, Yunlong Feng, Libo Qin, and Ting 720  
Liu. 2021. [N-LTP: An open-source neural language 721  
technology platform for Chinese](#). In *Proceedings of 722  
the 2021 Conference on Empirical Methods in Nat- 723  
ural Language Processing: System Demonstrations*, 724  
pages 42–49, Online and Punta Cana, Dominican Re- 725  
public. Association for Computational Linguistics. 726

DeepSeek-AI. 2024. [Deepseek llm: Scaling open- 727  
source language models with longtermism](#). *arXiv 728  
preprint arXiv:2401.02954*. 729

Shehzaad Dhuliawala, Mojtaba Komeili, Jing Xu, 730  
Roberta Raileanu, Xian Li, Asli Celikyilmaz, and 731  
Jason Weston. 2023. [Chain-of-verification reduces 732  
hallucination in large language models](#). 733

Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, 734  
Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, Qianyu Guo, 735  
Meng Wang, and Haofen Wang. 2023. [Retrieval- 736  
augmented generation for large language models: A 737  
survey](#). 738

Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, 739  
Zhangyin Feng, Haotian Wang, Qianglong Chen, 740  
Weihua Peng, Xiaocheng Feng, Bing Qin, and Ting 741  
Liu. 2023a. [A survey on hallucination in large lan- 742  
guage models: Principles, taxonomy, challenges, and 743  
open questions](#). 744

Quzhe Huang, Mingxu Tao, Chen Zhang, Zhenwei An, 745  
Cong Jiang, Zhibin Chen, Zirui Wu, and Yansong 746  
Feng. 2023b. [Lawyer llama technical report](#). 747



859	Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. <a href="#">Llama 2: Open foundation and fine-tuned chat models</a> .	At Stage 2, the revising threshold $\delta$ is set to 0.66. To enhance the quality of model outputs, we employ in-context learning (Brown et al., 2020). For each step, due to the limit of context length, we construct 2 human-annotated examples for in-context learning. We ask LLMs to regenerate for fragmentary outputs.	915
860			916
861			917
862			918
863			919
864			920
865			921
866			
867			
868			
869			
870			
871			
872	Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2023. <a href="#">Interleaving retrieval with chain-of-thought reasoning for knowledge-intensive multi-step questions</a> . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 10014–10037, Toronto, Canada. Association for Computational Linguistics.	<b>B Details of Evidence Annotation</b>	922
873		Following previous work (Huang et al., 2023b), we train a classification model to retrieve relevant articles. We fine-tune RoBERTa-large (Liu et al., 2019) with 80K examples. Each example consists of one question and 1~5 articles. For each consultation question, we keep the top 10 articles with the highest probability scores predicted by the classifier.	923
874			924
875			925
876			926
877			927
878			928
879			929
880	Jianing Wang, Qiushi Sun, Nuo Chen, Xiang Li, and Ming Gao. 2023. <a href="#">Boosting language models reasoning with chain-of-knowledge prompting</a> .	To avoid the model failing to retrieve articles that should serve as the basis, we employ 6 annotators with legal backgrounds to supplement missing articles. They are then asked to assess whether each article should/can serve as the basis to respond the given question based on the following principles:	930
881			931
882			932
883	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed H. Chi, Quoc V Le, and Denny Zhou. 2022. <a href="#">Chain of thought prompting elicits reasoning in large language models</a> . In <i>Advances in Neural Information Processing Systems</i> .		933
884			934
885			935
886			936
887			937
888	Wenhao Yu, Dan Iter, Shuohang Wang, Yichong Xu, Mingxuan Ju, Soumya Sanyal, Chenguang Zhu, Michael Zeng, and Meng Jiang. 2023. <a href="#">Generate rather than retrieve: Large language models are strong context generators</a> . In <i>The Eleventh International Conference on Learning Representations</i> .		938
889			939
890			940
891			941
892			942
893			943
894	Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah Goodman. 2022. <a href="#">STar: Bootstrapping reasoning with reasoning</a> . In <i>Advances in Neural Information Processing Systems</i> .		944
895			945
896			946
897			
898	Jiaxin Zhang, Zhuohang Li, Kamalika Das, Bradley Malin, and Sricharan Kumar. 2023. <a href="#">SAC<sup>3</sup>: Reliable hallucination detection in black-box language models via semantic-aware cross-check consistency</a> . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 15445–15458, Singapore. Association for Computational Linguistics.	We assign relevant scores of 2, 1, and 0 to the three categories of articles respectively. If the average score of an article exceeds 1.66, it will be regarded as a <i>necessary</i> one. And the articles with average scores less than 0.67 can be regarded as <i>not required</i> , while the remaining ones are <i>optional</i> .	947
899			948
900			949
901			950
902			951
903			952
904			953
905	Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans, Claire Cui, Olivier Bousquet, Quoc V Le, and Ed H. Chi. 2023. <a href="#">Least-to-most prompting enables complex reasoning in large language models</a> . In <i>The Eleventh International Conference on Learning Representations</i> .	To imitate retrieval-augmented generation, we provide five articles for each question, supposing them as retrieval results. We first keep all the <i>necessary</i> and <i>optional</i> articles. We then select <i>not required</i> articles in descending order of the probability scores predicted by the classifier.	954
906			955
907			956
908			957
909			958
910			
911			
912	<b>A Details of Generation</b>	<b>C Rules to Identify the Article Reference</b>	959
913	During generation, we set the temperature to 0.3, the repetition penalty to 1.05, and the top-p to 0.8.	We examine whether the responses use an article as the basis by following rules:	960
914			961

Model	Avg. Lens. of base	Avg. Lens. of CoD
Baichuan2-7B	201.20	200.78
Deepseek-7B	214.29	219.09
Qwen-7B	236.55	259.03
Xverse-7B	236.61	244.60

Table 5: The average lengths of responses generated by different LLMs in baseline setting and CoD setting.

- If the article number appears in the response, we believe the LLM has used this article as a reference.
- We segment the responses into sentence (Chen et al., 2021) and calculate the longest common subsequence (LCS) between each sentence and the article’s content. If the length of the longest LCS exceeds one-third of the article, we believe the LLM has referenced this article.
- Otherwise, the article is considered not to serve as a reference.

## D More Results and Analysis

### D.1 Lengths of Generated Responses

In Table 5 we provide the average lengths of generated responses from each LLM in different settings. We find that the lengths of responses under different experimental settings can be similar. After checking the correlation coefficient between the difference in response lengths and the difference in GPT-4 scores, we believe that the GPT-4 evaluator does not show preference for longer responses.

### D.2 Examples required Reasoning with Commonsense Knowledge

As shown in the *original analysis* of Table 6, LLMs tend to focus on the literal differences between questions and law articles. We humans have the background knowledge that only courts have the authority to revoke guardianship, thus often omit this information in our questions. But, LLMs often struggle with such differences between the questions and rigorous law articles, thus may not yield correct analysis.

### D.3 Significant Test

We apologize that we have not provided significance test. We have to admit that the GPT-4 evaluation is indeed too expensive. A single evaluation of all results under all settings can cost \$150. It

**Question:** Do you still need to pay child support after having your guardianship revoked?

**Article:** *Article 37* Parents, children, and spouses who support the wards in the form of child support, support for elderly parents, or spousal support in accordance with the law shall continue to perform such obligations after they are disqualified by the people’s courts as guardians.

**Original analysis:** *Article 37* explicitly stipulates [the content of Article 37]. This article pertains to guardianship and child support, but since the question does not mention revocation by the People’s Court, this article should not be used as a basis.

*Low-quality modification:*

**Revised analysis:** *Article 37* explicitly stipulates [the content of Article 37]. This article pertains to guardianship and child support. However, the user does not explicitly say who revokes her/his guardianship. Thus, this article should not be used as a basis.

*High-quality modification:*

**Revised analysis:** Article 37 stipulates that the revocation of guardianship does not affect existing obligations to pay child support. Thus, this article should be used as a basis.

Table 6: Failed and successful cases for revising evidence analysis. Red texts are the key basis of the question. Violet texts are correct analysis, while the texts with yellow background are hallucinated parts.

is difficult for us to repeat the evaluation multiple times, since it may cost thousands of dollars. 1000

In this work, we use a low temperature when generating. And we set the temperature to 0.0 for GPT-4 evaluation. 1001

We randomly sample 20 examples and employ Baichuan2-7B to repeat the experiments of the baseline setting for 5 times with different seeds. We then use GPT-4 to score the responses of every run. The standard error of the average scores is 0.0245, which is much smaller than the improvement brought by CoD. 1002

We also repeat the experiments on Baichuan2-7B for 5 times under the setting of CoD, using the same subset. The standard error of average scores yielded by GPT-4 is 0.0510, which is also smaller than the improvement brought by CoD. 1003

## E Prompts of Chain-of-Discussion 1004

### E.1 Prompt of Question Analysis 1005

To obtain the question analysis, we employ the prompts as below ([ ] indicates the English translation of the prompts in Chinese): 1006

你是一个民法领域的专家，你需要从法律专业的角度分析一名咨询者提出的问题涉及哪些关键点。在分析问题之后，你还要分析检索器提供的参考法条是否能作为分析该问题的依据。请你紧紧围绕咨询者的问题进行分析，不要过度设想潜在的、与问题不相关的场景。 1007

[ You are an expert in the field of civil law. You need to analyze the key points of a question posed by a consultant 1008



1150	法条: [[article]]	are required to provide specific legal articles as the basis for	1208
1151	[ Articles: [[article]]]	your answer, informing the consultant of their rights, obliga-	1209
1152	律师对于法条的分析: [[art_ana]]	tions conferred by the law, or actions prohibited by it. Before	1210
1153	[ Lawyer's analysis for the articles: [[art_ana]]]	answering the question, you can refer to some reference ar-	1211
1154	接下来, 请先用简洁的语言点评律师对	ticles provided by the retriever. However, please note that	1212
1155	于[[cur_art_id]]的分析。之后, 请你判断他的分析是否	the articles provided by the retriever may not necessarily be	1213
1156	误解了法条的内容。	helpful in answering the consultant's question; they may also	1214
1157	[ Next, please provide a concise critique of the lawyer's	be irrelevant to the question. Therefore, you need to analyze	1215
1158	analysis of [[cur_art_id]]. Then, determine whether their	the factual background of the issue involved, then determine	1216
1159	analysis misconstrues the content of the legal article. ]	whether each article can serve as a basis for answering the	1217
1160	<b>E.5 Prompt of revising</b>	question. Please do not consider all the reference articles	1218
1161	When revising, we employ following prompts:	provided by the retriever as the basis, nor cite any articles	1219
1162	你是一名律师, 你对于某个法条是否有助于解答某个	outside the reference ones as evidence. During your response,	1220
1163	法律咨询问题进行了点评。一些法学专家认为你的点评	focus closely on the consultant's question, avoiding overly	1221
1164	中存在对法条内容的理解、法条与问题之间的关联性等	imagining potential scenarios unrelated to the issue.]	1222
1165	角度存在错误。你需要参考你对问题的分析, 修改你对	咨询者的问题是“[[question]]”	1223
1166	法条的分析。	[ The consultant's question is: "[[question]]"]	1224
1167	[ You are a lawyer who has provided an assessment of	下面是检索器提供的参考法条: [[articles]]	1225
1168	whether a certain legal article is helpful in addressing a spe-	[ Below are the reference articles provided by the retriever:	1226
1169	cific legal consultation question. Some legal experts believe	[[articles]]]	1227
1170	there are errors in your assessment regarding understanding	接下来, 请你回答咨询者提出的问题“[[question]]”你	1228
1171	the content of the legal article and its relevance to the question.	需要先对该问题的关键点进行分析, 然后判断各个参考	1229
1172	You need to revise your analysis of the legal article based on	法条是否有助于解答该问题。最后请你使用与该问题有	1230
1173	your analysis of the issue. ]	关的部分法条作为依据, 给出详细的回答。回答过程中	1231
1174	问题: [[question]]	禁止使用参考法条之外的内容。	1232
1175	[ Question: [[question]]]	[ Next, please answer the question posed by the consul-	1233
1176	法条: [[article]]	tant "[[question]]" You need to analyze the key points of the	1234
1177	[ Articles[[article]]]	question first, then determine whether each reference article	1235
1178	律师对于问题的分析: [[que_ana]]	is helpful in answering it. Finally, please provide a detailed	1236
1179	[ Lawyer's Analysis of the Question: [[que_ana]]]	answer using relevant portions of the articles as the basis. Use	1237
1180	律师对于法条的分析: [[art_ana]]	of content outside the reference articles is prohibited during	1238
1181	[ Lawyer's Analysis of the Legal Article: [[art_ana]]]	the response. ]	1239
1182	专家点评: [[critiques]]	问题分析: [[que_ana]]	1240
1183	[ Expert Critiques: [[critiques]]]	[ Question Analysis: [[que_ana]]]	1241
1184	接下来, 请你重写一份更为正确的法条分析。在重写	法条分析: [[art_ana]]	1242
1185	后的法条分析的结尾, 请你按照你的分析, 评估一下该	[ Article Analysis: [[art_ana]]]	1243
1186	法条是否可能有助于解答问题。	回答:	1244
1187	[ Next, please rewrite a more accurate analysis of the legal	[ Answer:]	1245
1188	article. At the end of the rewritten analysis of the legal article,	<b>F Scoring Prompt of GPT-4</b>	1246
1189	evaluate whether the legal article may indeed be helpful in	Following CritiqueLLM (Ke et al., 2023), we em-	1247
1190	addressing the question based on your analysis. ]	ploy a reference-based prompt to instruct GPT-4 to	1248
1191	<b>E.6 Prompt of response</b>	assess the overall quality of the responses generated	1249
1192	Finally, to response to the user's question, we use	by open-source LLMs. We use the human-written	1250
1193	following prompts:	response and the <i>necessary</i> and <i>option</i> articles as	1251
1194	你是一个法律专家, 你需要从法律专业的角度回答	reference. The prompt is shown as below:	1252
1195	咨询者提出的问题。你需要以具体的法条为依据回答问	[Instruction]	1253
1196	题, 并告诉咨询者法律赋予他哪些权利和义务, 或者禁	Please act as an impartial judge and evaluate the	1254
1197	止他实施哪些举措。在回答问题之前, 你可以参考检索	quality of the response provided by an AI assis-	1255
1198	器提供的一些参考法条。但请注意, 检索器提供的法条	tant to the user question displayed below. Your	1256
1199	并不一定都有助于回答咨询者提出的问题, 它也可能与	evaluation should consider factors such as the logi-	1257
1200	提问者的问题无关。因此, 你需要对问题涉及的事实背		
1201	景进行分析, 再判断各个法条是否能够作为回答问题的		
1202	依据。请你不要将检索器提供的全部参考法条都当作依		
1203	据, 也不要引用参考法条之外的其他法条作为依据。在		
1204	回答的过程中, 请你紧紧围绕提问者的问题进行讨论,		
1205	不要过度设想潜在的、与问题不相关的情形。		
1206	[ You are a legal expert, and you need to answer the ques-		
1207	tion posed by the consultant from a legal perspective. You		

1258 cality, helpfulness, relevance, accuracy, depth, and  
1259 whether using irrelevant articles beyond the refer-  
1260 ence articles as a basis. Begin your evaluation by  
1261 providing a short explanation. You will be given  
1262 several reference articles, a high-quality reference  
1263 answer and the assistant’s answer. Be as objective  
1264 as possible. You should first provide your explana-  
1265 tion IN CHINESE, then you must rate the response  
1266 on a scale of 1 to 10 by STRICTLY following  
1267 the below MAPPING for the relation between the  
1268 scores and response quality:

1269 1. The score 1~2 stands for very chaotic or ab-  
1270 sence of answer, and the AI assistant completely  
1271 failed to answer the user’s question. Serious logi-  
1272 cal and factual errors might also be included in this  
1273 term. The gap between the AI assistant’s answer  
1274 and the high-quality reference answer is huge and  
1275 insuperable.

1276 2. The score 3~4 indicates fragment-like re-  
1277 sponses from AI assistant’s answer. It did not  
1278 provide answers in proper grammar, fluency, or  
1279 accuracy. Citing irrelevant articles and resulting in  
1280 a redundant output also falls under this scenario.  
1281 There are obvious gaps between the high-quality  
1282 reference answer and the AI assistant’s response.

1283 3. The score 5~6 indicates for existence of  
1284 minute disadvantage from the AI assistant’s answer  
1285 compared to the high-quality reference answer. Yet  
1286 the AI assistant did provide an average answer. The  
1287 AI assistant either did not fully address the question,  
1288 or was somewhat short of logicity, helpfulness,  
1289 relevance, depth, or detailedness. The disadvan-  
1290 tages from the AI assistant’s answer overwhelm its  
1291 advantages.

1292 4. The score 7~8 indicates the AI assistant pro-  
1293 vided a good answer as well as the high-quality  
1294 reference answer, addressing the question, with  
1295 good helpfulness, relevance, accuracy, depth, cre-  
1296 ativity, and enough details. The response of AI  
1297 assistant does not include any irrelevant articles be-  
1298 yond the reference articles. The AI assistant might  
1299 have flaws compared to the reference answer, but  
1300 that does not overwhelm the above advantages.

1301 5. The score 9~10 indicates the AI assistant re-  
1302 sponded better than the provided reference answer  
1303 in most aspects, fully achieved the instruction goal,  
1304 provided more detailed analysis, and have unique  
1305 advantages to the reference answer.

1306 You should give scores around 7 if you do not  
1307 find obvious advantages or disadvantages. You  
1308 should seriously consider the above guide before

1309 give lowest and highest scores such as 1 or 10, and  
1310 avoid such situation if you do not have sound expla-  
1311 nations. Avoid any positional biases and ensure that  
1312 the order in which the responses were presented  
1313 does not influence your decision. Do not allow the  
1314 length of the responses to influence your evaluation.  
1315 Do not favor certain names of the assistants. AND  
1316 again, VERY IMPORTANTLY, after you provide  
1317 your explanation IN CHINESE, you must rate the  
1318 response strictly following this format: “Rating:  
1319 [[Number]]”, for example: Rating: [[5]].

[User’s Question]

{{QUESTION}}

[The Start of Reference Articles]

{{ARTICLES}}

[The End of Reference Articles]

[The Start of Reference Answer]

{{GOLDEN RESPONSE}}

[The End of Reference Answer]

[The Start of Assistant’s Answer]

{{LLM’S RESPONSE}}

[The End of Assistant’s Answer]